

PRINCIPLES OF HELIOPHYSICS:
a textbook on the universal processes
behind planetary habitability

by

*Karel Schrijver, Fran Bagenal, Tim Bastian, Jürg Beer, Mario Bisi, Tom Bogdan, Steve Bougher,
David Boteler, Dave Brain, Guy Brasseur, Don Brownlee, Paul Charbonneau, Ofer Cohen,
Uli Christensen, Tom Crowley, Debrah Fischer, Terry Forbes, Tim Fuller-Rowell, Marina Galand,
Joe Giacalone, George Gloeckler, Jack Gosling, Janet Green, Steve Guetersloh, Viggo Hansteen,
Lee Hartmann, Mihaly Horanyi, Hugh Hudson, Norbert Jakowski, Randy Jokipii, Margaret Kivelson,
Dietmar Krauss-Varban, Norbert Krupp, Judith Lean, Jeff Linsky, Dana Longcope, Daniel Marsh,
Mark Miesch, Mark Moldwin, Luke Moore, Sten Odenwald, Merav Opher, Rachel Osten,
Matthias Rempel, Hauke Schmidt, George Siscoe, Dave Siskind, Chuck Smith, Stan Solomon,
Tom Stallard, Sabine Stanley, Jan Sojka, Kent Tobiska, Frank Toffoletto, Alan Tribble,
Vytenis Vasylunas, Richard Walterscheid, Ji Wang, Brian Wood, Tom Woods, and Neal Zapp*

Principles of heliophysics, V 1.2, October 30, 2019

Contents

Preface (do read it!) *page* vii

Introduction

1	Stars, planets, planetary systems, and the local cosmos	1
1.1	Preparing for the future	1
1.2	Considering planetary habitability	2
1.3	Heliophysics: unification, coupling, exploration	3
1.4	The language of heliophysics	7
1.5	A timeline of exploration of planetary systems	11

Foundations

2	Neutrals, ions, and photons	13
2.1	Conditions in the local cosmos	13
2.2	Gravitationally stratified atmospheres and stellar winds	16
2.3	Photons, collisions, ionization, and differentiation	27
2.4	On collisions and currents, and on neutrals and pickup ions	34
2.5	Sources of plasma	43
3	MHD, field lines, and reconnection	44
3.1	Introduction	44
3.2	(Magneto-)Hydrodynamics	49
	3.2.1 MHD equations, individual terms, and special cases	52
	3.2.2 The induction equation	55
3.3	Waves in magnetized plasmas	59
3.4	MHD, magnetic field lines and reconnection	61
3.5	A few notes about conditions	70
	3.5.1 Solar atmosphere <i>vs.</i> terrestrial magnetosphere	70

3.5.2	Heliosphere	71
4	Dynamos of Sun-like stars and Earth-like planets	72
4.1	Dynamo settings	72
4.1.1	Earth and other terrestrial planets	77
4.1.2	The Sun and other stars	78
4.2	Dynamo principles	80
4.3	Essentials of fluid motions in dynamos	82
4.4	Insights from approximate stellar dynamo models	86
4.5	Mean-field dynamo models	88
4.6	Dynamos in other stars	101
4.7	Dynamos in terrestrial planets	104
5	Flows, shocks, obstacles, and currents	109
5.1	Introductory overview	109
5.2	Low-velocity interactions versus shocks	112
5.3	Elementals of shocks and other discontinuities	114
5.4	The magnetized solar wind and the Parker spiral	120
5.5	Flow-based interactions in heliophysics	122
5.5.1	Solar-wind stream interactions	122
5.5.2	A non-conducting body without atmosphere	128
5.5.3	Flow around a conducting body	130
5.5.4	Plasma flow around a permanently magnetized body	132
5.5.5	A closed magnetosphere	133
5.5.6	The open magnetosphere	135
5.5.7	Solar wind-magnetosphere-ionosphere interaction	140
5.5.8	A large-scale flow impinging on a fast outflow	146
6	Magnetic (in-)stability and energy pathways	148
6.1	Introduction	148
6.1.1	Introducing solar flares and coronal mass ejections	150
6.1.2	Introducing geospace (sub-)storms	153
6.2	Terrestrial magnetospheric disturbances	155
6.2.1	Energy pathways and reservoirs	155
6.2.2	What leads to explosive energy releases?	157
6.2.3	Terrestrial magnetospheric substorms	160
6.2.4	Terrestrial magnetic storms	161
6.3	Solar impulsive events	163
6.3.1	The magnetic reservoir	163
6.3.2	Two-dimensional force-free models	165
6.3.3	Three-dimensional force-free models	167
6.3.4	Formation of the pre-eruption field	170
6.3.5	Observed signatures of flares and CMEs	172

6.4	Magnetic instabilities and reconnection	178
7	Torques and tides	182
7.1	Introduction	182
7.2	Magnetic torques	184
7.2.1	Stellar winds and magnetic braking	184
7.2.2	Planetary magnetospheric torque	185
7.2.3	Magneto-rotational coupling	188
7.2.4	Disk winds	190
7.3	Gravitational tides	192
7.3.1	Spin-orbit interactions	192
7.3.2	Orbital interaction	196
7.4	Planetary atmospheric tides	198
8	Particle orbits, transport, and acceleration	199
8.1	Single particle motion	200
8.2	Phase space density and Liouville's theorem	206
8.3	The collisionless Boltzmann equation	207
8.4	Particle scattering and transport	210
8.4.1	Solar energetic particles	214
8.4.2	Galactic cosmic rays	216
8.5	Particle acceleration in shocks	220
8.6	Relativistic particles in planetary radiation belts	229
8.6.1	Electron acceleration mechanisms	229
8.6.2	Proton acceleration in the radiation belt	233
8.6.3	Radiation belt losses at Earth	234
9	Convection, heating, conduction, and radiation	236
9.1	Convective and radiative energy transport	236
9.2	Heating and cooling of the solar outer atmosphere	241
9.3	Magnetic activity and atmospheric radiation	245
Comparative eco-astrophysics		
10	Evolution of stars, their activity, and their asterospheres	247
10.1	Evolution of stars	247
10.2	Stellar activity and its evolution	251
10.2.1	Overall activity level	251
10.2.2	Flares	253
10.2.3	Rotation rates	256
10.2.4	Stellar infancy: birth to the zero-age main sequence	258
10.2.5	Stellar teenage years: ZAMS - 1 Gyr	259

10.2.6	Stellar adulthood: 1-5 Gyr	260
10.3	Evolution of astrospheres	261
10.3.1	Effects of a variable ISM on heliospheric structure	261
10.3.2	Long-term evolution of stellar winds	265
10.3.3	Astrospheric field patterns in time	270
11	Formation of stars and planets	273
11.1	(Exo-)Planets and (exo-)planetary systems	275
11.1.1	Exoplanet formation	278
11.1.2	Exoplanet migration	280
11.1.3	Exoplanet geology	281
11.1.4	Exoplanets and binary star systems	283
11.2	Formation and early evolution of stars and disks	283
11.2.1	Observations of star-forming processes	283
11.2.2	Properties of young stars	287
11.2.3	The rotation rate of very young stars	291
11.2.4	Protoplanetary disks and gravity	292
11.2.5	Dust-disk evolution	295
11.2.6	Disk evaporation	298
12	Evolving irradiance, atmospheres, and habitability	300
12.1	Evolving planetary habitability	300
12.1.1	Earth's formative phase	300
12.1.2	The habitable zone	302
12.1.3	Oxygen, methane, and carbon dioxide over time	304
12.1.4	Water over time	306
12.2	Atmospheres and climates of Venus, Earth, and Mars	307
12.3	Irradiance, orbits, spin, and climate	310
12.3.1	Atmospheric effects and albedo	310
12.3.2	Orbital changes	318
12.4	Planetary atmospheres, geological activity, and stellar winds	320
12.4.1	On time scales beyond millions of years	320
12.4.2	On time scales of up to several millennia	328
13	Evolving upper atmospheres and iono-magnetospheres	330
13.1	Maintaining ionospheres	334
13.1.1	Ionization	334
13.1.2	Recombination	336
13.1.3	Venus and Mars	339
13.2	Setting geospace climate	341
13.2.1	Geospace climate response to solar photon irradiation	341
13.2.2	Geospace climate at earlier terrestrial ages	346
13.2.3	Geospace climate and Earth's magnetic field	349

13.2.4	Geospace climate dependence on the solar wind	352
14	Magnetic fields and cosmic rays over time	354
14.1	Long-term energetic-particle exposure of Earth	355
14.1.1	Generation of cosmogenic radionuclides	355
14.1.2	Transport and deposition of cosmogenic radionuclides	359
14.2	Radionuclides as proxies of magnetic variability	360
14.2.1	Geomagnetic field	362
14.2.2	Solar variability	363
14.2.3	Very-long time scale variability in cosmic-ray exposure	365
14.3	Exposure to supernovae	366
15	Applied heliophysics, <i>mutatis mutandis</i>, ...	367

Activities

16	Compilation of activities found throughout the text	372
-----------	--	------------

<i>Version history</i>	400
<i>List of Illustrations</i>	401
<i>List of Boxes and Tables</i>	404
<i>Bibliography</i>	405

Preface (do read it!)

The five volumes on Heliophysics (four published in printed form by Cambridge University Press, and one online at the Heliophysics Summer School website) contain in total 1919 pages of text and figures, in 56 topical chapters (Vol. I: Schrijver and Siscoe (2011); Vol. II: Schrijver and Siscoe (2012b); Vol. III: Schrijver and Siscoe (2012a); Vol. IV: Schrijver *et al.* (2016); Vol. V: Schrijver and Siscoe (2015)). The present volume presents a selection of these texts, while adding new text as connecting or summarizing material, with an overall text length that is less than one-fifth of the original textbooks. The topics in

Heliophysics

helio-, pref., on the Sun and environs, from the Greek helios.
physics, n., the science of matter and energy and their interactions.

Heliophysics is the

- *comprehensive new term for the science of the Sun - Solar System Connection.*
- *exploration, discovery, and understanding of our space environment.*
- *system science that unites all of the linked phenomena in the region of the cosmos influenced by a star like our Sun.*

Heliophysics concentrates on the Sun and its effects on Earth, the other planets of the solar system, and the changing conditions in space. Heliophysics studies the magnetosphere, ionosphere, thermosphere, mesosphere, and upper atmosphere of the Earth and other planets. Heliophysics combines the science of the Sun, corona, heliosphere and geospace. Heliophysics encompasses cosmic rays and particle acceleration, space weather and radiation, dust and magnetic reconnection, solar activity and stellar cycles, aeronomy and space plasmas, magnetic fields and global change, and the interactions of the solar system with our galaxy.

From NASA's "Heliophysics. The New Science of the Sun - Solar System Connection: Recommended Roadmap for Science and Technology 2005 - 2035."

Table 0.1.

this volume are organized to emphasize universal processes from a perspective that draws attention to what provides Earth (and similar (exo-)planets) with a relatively stable setting in which life as we know it can thrive. This text aims to serve as a textbook-style volume for which the original Heliophysics books are the extended ‘readers’ with much more detail, and domain-specific topical chapters. Note that references from the original texts were omitted here (see the original volumes for those); references for new texts can be found in the Bibliography, where also source references to figures are provided as needed.

This volume is intended for students in physical sciences in later years of their university training and for beginning graduate students in fields of solar, stellar, (exo-)planetary, and planetary-system sciences. This contrasts with the intended audiences for the Heliophysics volumes which included the community of mid-to-advanced graduate students, the cohort of early postdoctoral researchers, and those professional researchers looking for review-like introductions into fields of heliophysics adjacent to their own. In targeting the audience of advanced undergraduate and beginning graduate students, many of the deeply technical details discussed in the original volumes were omitted, introductions were broadened, and the emphasis was placed on processes rather than on details of equations, states, or numerical experiments.

Throughout this work the original text from the Heliophysics volumes is directly quoted, following a volume and chapter reference, where between double quotation marks, but with equations, units (here cgs-Gaussian throughout with a few exceptions ^[i]), and symbols modified where needed for homogeneity throughout this work, with edits (and some corrections) shown between brackets, with many parenthetical notes removed, and with citations of the professional literature left out (and those to other sections in the books modified as appropriate).

The source texts in the series of Heliophysics books are referenced as H-#[roman]:#[arabic].#[arabic]. For example, Vol. I, Section 9 in Chapter 2 would be referred to as H-I:2.9. The original sources of all of the figures can be found in the figure captions of the Heliophysics books, but for many here a reference to the original publication is included for figures not made by the Heliophysics authors but whose original authors have given permission to have their artwork used in this volume. A few figures were replaced by color versions or by alternative figures.

ⁱ A good resource for unit conversions (and many other things related to plasma physics) is the online NRL plasma formulary.

Table 0.2. *Chapters and their authors in the Heliophysics series sorted by theme (continued on the next page), not showing introductory chapters.*

Universal and fundamental processes, diagnostics, and methods	
I.2. Introduction to heliophysics	<i>T. Bogdan</i>
I.3. Creation and destruction of magnetic field	<i>M. Rempel</i>
I.4. Magnetic field topology	<i>D. Longcope</i>
I.5. Magnetic reconnection	<i>T. Forbes</i>
I.6. Structures of the magnetic field	<i>M. Moldwin et al.</i>
II.3 In-situ detection of energetic particles	<i>G. Gloeckler</i>
II.4 Radiative signatures of energetic particles	<i>T. Bastian</i>
II.7 Shocks in heliophysics	<i>M. Opher</i>
II.8 Particle acceleration in shocks	<i>D. Krauss-Varban</i>
II.9 Energetic particle transport	<i>J. Giacalone</i>
II.11 Energization of trapped particles	<i>J. Green</i>
IV.11 Dusty plasmas	<i>M. Horányi</i>
IV.12 Energetic-particle environments in the solar system	<i>N. Krupp</i>
IV.13 Heliophysics with radio scintillation and occultation	<i>M. Bisi</i>
Stars, their planetary systems, planetary habitability, and climates	
III.3 Formation and early evol. of stars and proto-planetary disks	<i>L. Hartmann</i>
III.4 Planetary habitability on astronomical time scales	<i>D. Brownlee</i>
III.11 Astrophysical influences on planetary climate systems	<i>J. Beer</i>
III.12 Assessing the Sun-climate relationship in paleoclimate records ...	<i>T. Crowley</i>
III.14 Long-term evolution of the geospace climate	<i>J. Sojka</i>
III.15 Waves and transport processes in atmosph. and oceans	<i>R. Walterscheid</i>
IV.5 Characteristics of planetary systems	<i>D. Fischer & J. Wang</i>
IV.7 Climates of terrestrial planets	<i>D. Brain</i>
The Sun, its dynamo, and its magnetic activity; past, present, and future	
I.8. The solar atmosphere	<i>V. Hansteen</i>
II.5 Observations of solar and stellar eruptions, flares, and jets	<i>H. Hudson</i>
II.6 Models of coronal mass ejections and flares	<i>T. Forbes</i>
III.2 Long-term evolution of magnetic activity of Sun-like stars	<i>C. Schrijver</i>
III.5 Solar internal flows and dynamo action	<i>M. Miesch</i>
III.6 Modeling solar and stellar dynamos	<i>P. Charbonneau</i>
III.10 Solar irradiance: measurements and models	<i>J. Lean & T. Woods</i>
IV.2 Solar explosive activity throughout the evol. of the solar system	<i>R. Osten</i>

Activities for the reader

New here is the inclusion of 200 'activities' in the form of problems, exercises, explorations, literature readings, and 'what if' challenges. Many contain additional information complementing the main text, so I suggest you read

Table 0.2. (Continued from the previous page) Chapters and their authors in the Heliophysics series sorted by theme, not showing introductory chapters.

Astro-/heliospheres, the interstellar environment, and galactic cosmic rays	
I.7. Turbulence in space plasmas	<i>C. Smith</i>
I.9. Stellar winds and magnetic fields	<i>V. Hansteen</i>
III.8 The structure and evolution of the 3D solar wind	<i>J. Gosling</i>
III.9 The heliosphere and cosmic rays	<i>J. Jokipii</i>
IV.3 Astrospheres, stellar winds, and the interst. medium	<i>B. Wood & J. Linsky</i>
IV.4 Effects of stellar eruptions throughout astrospheres	<i>O. Cohen</i>
Dynamics and environments of planets, moons, asteroids, and comets	
I.10. Fundamentals of planetary magnetospheres	<i>V. Vasyliūnas</i>
I.11. Solar-wind magnetosphere coupling	<i>F. Toffoletto & G. Siscoe</i>
I.13. Comparative planetary environments	<i>F. Bagenal</i>
II.10 Energy conversion in planetary magnetospheres	<i>V. Vasyliūnas</i>
III.7 Planetary fields and dynamos	<i>U. Christensen</i>
IV.6 Planetary dynamos: updates and new frontiers	<i>S. Stanley</i>
IV.10 Moons, asteroids, and comets interact. with their surround.	<i>M. Kivelson</i>
Planetary upper atmospheres	
I.12. On the ionosphere and chromosphere	<i>T. Fuller-Rowell & C. Schrijver</i>
II.12 Flares, CMEs, and atmospheric responses	<i>T. Fuller-Rowell & S. Solomon</i>
III.13 Ionospheres of the terrestrial planets	<i>S. Solomon</i>
III.16 Solar variability, climate, and atmosph. photochemistry	<i>G. Brasseur et al.</i>
IV.8 Upper atmospheres of the giant planets	<i>L. Moore et al.</i>
IV.9 Aeronomy of terrestrial upper atmospheres	<i>D. Siskind & S. Bougher</i>
Technological and societal impacts of space weather phenomena	
II.2 Introduction to space storms and radiation	<i>S. Odenwald</i>
II.13 Energetic particles and manned spaceflight	<i>S. Guetersloh & N. Zapp</i>
II.14 Energetic particles and technology	<i>A. Tribble</i>
V.2 Space weather: impacts, mitigation, forecasting	<i>S. Odenwald</i>
V.3 Commercial space weather in response to societal needs	<i>W. Tobiska</i>
V.4 The impact of space weather on the electric power grid	<i>D. Boteler</i>
V.5 Radio waves for communication and ionospheric probing	<i>N. Jakowski</i>

A:0 them as you go along, if not on first reading, then at least on review. {A:[0]}
 Some were developed by the teachers for the Heliophysics Summer School but most are newly created specifically for this volume. They are meant to let the reader look up a definition, to introduce a moment of reflection on an equation or figure, to see connections to similarities elsewhere, to get a feel for the magnitude of things or the relative importance of processes, or to consider

⁰ Activity: Exercises are flagged like this, with continuous numbering throughout the volume. These 'activities' are shown both at the bottom of the pages where included as well as extracted from the main text and repeated as a compilation Ch. 16.

Heliophysics and space weather:

'Space weather' is the term used to describe an ensemble of changing conditions in the vicinity of Earth and, by extension, any other body in a planetary system, typically occurring on time scales up to a few days. Often, the term is implicitly taken to refer also to the conditions from the solar dynamo outward to the furthest reaches of the heliosphere that are involved in space weather around Earth. Much of what is described in this volume therefore concerns space weather: heliophysics contains the science of space weather. However, where the science of space weather focuses on phenomena that can impact society through short-term variability, this text takes the long view by putting the spotlight on evolutionary changes in the states of star-planet systems. As such, this text does describe the foundational processes of space weather, but is not concerned with the impacts of space weather on technological infrastructure, does not address the challenges of forecasting space weather, and skips coupling mechanisms such as ground-induced currents (GICs) associated with geomagnetic disturbances and ground-level enhancements (GLEs) of energetic particles. This choice of focus is motivated by my desire to introduce the reader to the science of heliophysics from the perspective of habitability on time scales on which stellar and planetary atmospheres change, and indeed up to time scales on which stars and their planets evolve, and to do that in a relatively compact form. As you go through this text, you should realize that many of the processes described here have consequences for society, ranging from system design choices to potentially substantial failures in one or more of the infrastructures that we have come to rely on, including continuous and reliable electric power, positional information, and means of communication. Interruptions in quality or availability of any of these can have substantial consequences that may be costly or life-threatening on scales that may involve single individuals or populations of millions. Descriptions of the impacts of space weather can be found in the Heliophysics books in Chapters H-II:2, H-II:13, H-II:14, H-V:2, H-V:3, H-V:4, and H-V:5; another resource is a 'roadmap' document (Schrijver *et al.*, 2015) that reviews the state of our knowledge of space weather and its technological and societal impacts, and what is needed to advance our abilities to forecast space weather.

Table 0.3.

what would happen under conditions other than those encountered in our Solar System; they are not meant to particularly exercise mathematical skills. At the end of the book, in Activity 200, the reader is asked to reflect back on all the processes that are involved in the habitability of Earth and, by analogy, of exoplanets elsewhere in the Universe.

Terminology

As you go through this volume, you will encounter words that have somewhat different meanings in different communities. For example, 'convection' is often used in the magnetospheric community to describe movement that in astrophysics would be referred to as 'advection', while 'convection' in that discipline is reserved for overturning plasma motions involved in the transport

of thermal energy. Another example is that of the word 'dynamo' which in astrophysics and planetary sciences is used to describe the ensemble of processes maintaining a magnetic field against decay, often with an alternating temporal character. In ionospheric physics, it is often used for processes where differential motions of (neutral plus ionized) gas and magnetic field exchange energy through work.

You will also note that terms may describe locations or something in a location, or a property of what is in that location. For example, 'ionosphere' may sound like a location descriptor but actually refers to only the ionized medium in an atmosphere (with 'thermosphere' used for the overlapping neutral environment). The term 'chromosphere', which describes a stellar environment in some respects not dissimilar to an ionosphere-thermosphere, encompasses both the ionized and neutral components; it is often used as an indicator of a volume above a stellar surface in a certain thermal range, but is defined formally (as you will see later) by the properties of the radiative transfer of the medium.

Finally, there are words like 'late-type star' that have nothing to do with a temporal attribute, but which survived an older era where the nature of stars was not yet understood and where cooler was erroneously interpreted as older.

I hope that all terms are properly defined where first used. Here, I want to raise your awareness that as you talk to colleagues in other disciplines they may not only be puzzled by processes that you study, but that communication may be hampered by misinterpretation of the terms that you use: language can be a very precise tool, but only if the user is aware of how the listener/reader may interpret the words that are being used.

A few notes on other resources

The figures published in the Heliophysics book series are available on-line at the website of the Heliophysics Summer School, where you can also find labs (with instruction manuals) and many recorded lectures sorted by theme (in part hosted on youtube).

There is no subject index in this volume: this being an online and electronic book, the search tools of web browsers and pdf viewers provide a more effective and entirely comprehensive alternative.

This volume focuses on processes, not on their measurements. For an introduction to some of the aspects of remote and *in-situ* sensing within the Heliophysics series I refer to the following chapters that focus on that aspect in particular: Chs. H-II:3, H-II:4, H-IV:5, and H-IV:13.

Navigating this pdf

References to sections, figures, tables, and equations in this book are shown surrounded by a red box, pointers to the bibliography are surrounded by a green box, and references to web pages are shown surrounded by a blue box. Clicking on a box gets you to that reference. How you get back to reading where you left off depends on how you are viewing this file and on what type of device. For example, this web site shows a list of keyboard shortcuts to move around the pdf version of this book with Acrobat Reader. Using that on a Mac, you can return to the page you came from by pressing [command + left-arrow] after clicking on a link to a figure, section, or activity.

Corrections and updates

This version of the textbook is subject to corrections and updates. I welcome input from students, teachers, and colleagues: if you see a typo or an explanation that you think is in error, or if you believe a serious update is in order, please email me! Be as specific as you can about where the text is that you think should be changed, what to change it to, and why it needs such a change. Your input will help improve this text for all users.

Acknowledgements

I thank all the Heliophysics authors and other teachers in the Heliophysics Summer School for the skill with which they taught me as one of the participants in the School as well as for their patience with me as one of the editors of the book series. This volume was supported by the Johannes Geiss Fellowship of the International Space Science Institute. I am indebted to Jürg Beer, Paul Charbonneau, Terry Forbes, Marina Galand, Dana Longcope, Sten Odenwald, and Matthias Rempel for their insightful comments on an earlier version of the manuscript. A special thanks goes to Tom Bogdan who worked through the entire draft volume and made many insightful suggestions.

Karel Schrijver, November 1, 2019

1

Stars, planets, planetary systems, and the local cosmos

1.1 Preparing for the future

By the time you reach the end of this book, you will have the basic set of tools of scientific imagination involved in understanding what couples stars and planets. What you will learn is universal, literally: it does not matter which stars and planets we speak of: whether of those few nearby or of the many distant ones. Nor does it matter whether they are those few that we are long familiar with or the many that we know about, so far, only in a statistical sense. It does not matter either whether your particular interests lie within the Solar System or beyond it: the same principles apply in our local cosmos as in the most distant planetary system we shall ever have access to.

But looking forward to your science-based career, whether as a researcher or as a teacher, as a journalist or as a politician, you need to be familiar with what is known. That is particularly true in order to discover something new. And to appreciate the value of a discovery, you need to know how to apply what you know to what is not (yet) known. You will need to imagine things no one has ever seen, but not arbitrarily: science demands that you come up with what appears most probable, not merely with things that are possible. Richard Feynman (in *The meaning of it all*) said it this way: 'It is surprising that people do not believe that there is imagination in science. It is a very interesting kind of imagination, unlike that of the artist. The great difficulty is in trying to imagine something that you have never seen, that is consistent in every detail with what has already been seen, and that is different from what has been thought of.'

The pace at which exoplanets are being discovered is simply amazing. What we can learn from them, and from our Solar System, offers so many opportunities to learn yet more. I realize that going through the first nine chapters will be hard, because they have to build your foundation, because they cover so many different branches of science, and because they look at things so different

from everyday life. But these first nine chapters look for commonalities, for 'universal processes' that help create in your mind a virtual laboratory: in the astronomical sciences we cannot turn dials to explore things under different conditions, but we can compare different environments and look for what they have in common and for what sets them apart.

The final six chapters require prior digestion of the first nine. These final six invite you to imagine, scientifically, Earth in the distant past and future, Earth-like planets in a variety of orbits around Sun-like stars, and the space environments and climates of tropospheres of exo-worlds. Future discoveries have their beginnings in lessons from the past:

1.2 Considering planetary habitability

Planetary system are, statistically speaking, about as common as stars. We have learned a lot about stars over the century that followed the realization that they are huge nuclear fusion reactors and that most, like the Sun, function also as giant dynamos. In contrast, firm evidence that planetary systems are common companions to stars was only obtained within the past two decades. It is therefore no surprise that much still needs to be learned about how planets form, how planetary systems evolve, and what the conditions are near planetary surfaces (if indeed a solid or liquid surface exists). The combination of exoplanetary science and the study of the local cosmos is enlightening us as much about the history and future of our Solar System as about the growing number of planetary systems that have been observed in some detail. Whether life exists anywhere beyond Earth remains to be established, but scientists are making rapid headway in knowing about the conditions that life on Earth has been subjected to since its genesis and also the conditions that any life on any other planet would be subjected to depending on the properties of their central star and companion planets.

Heliophysics deals with all of the aspects of 'living with a star' on time scales from fractions of a second to billions of years. The series of Heliophysics books offers an introduction to a large cross section of that vast scientific field. In the present volume, we focus on the universal processes that tie together the branches of heliophysics with particular emphasis on those processes that are relevant to what one might describe as 'planetary habitability'. With life having been found on only a single astrophysical body we do not have a particularly well-considered concept of what 'planetary habitability' might mean, of course. But we have an intuitive feel for it: a long-lived planet orbiting a long-lived star, with a fairly substantial planetary atmosphere that is neither too hot nor too cool to allow chemistry to be complex (and, in many minds, restricting that to chemistry that involves liquid water), shielded well enough

(but not necessarily perfectly) from energetic radiation (both electromagnetic and particulate) by that atmosphere and by a planetary magnetic field. The star's irradiance onto a 'habitable planet' should not vary too much, comet and asteroid impacts should be limited, atmospheric erosion slow, . . . If that sounds like we are describing the Sun-Earth system then that is no surprise: we know it has made the Earth habitable to a diversity of life that is on one hand astoundingly diverse and on the other – at the molecular level – remarkably homogeneous.

As the number of known exoplanetary systems is bound to keep growing rapidly, and as our instrumentation and methods are bringing exoplanetary atmospheric science within our grasp, it is clear that our understanding of the Solar System and its central star provide crucial guidance to the study of 'planetary habitability' and – some day rather soon, one should anticipate – the study of extraterrestrial life. That expectation has guided the selection of topics covered in this volume.

[H-I:2.9] ^[ii] “If we gaze upon the uncountable array of stars strewn across the vault of the heavens, one may know that the remarkable things one will come to know about heliophysics in the pages that follow are presently unfolding around those very stars and planetary systems that give light to the night sky. Heliophysics is truly a universal science.”

1.3 Heliophysics: unification, coupling, exploration

[H-I:2.1] “Walk along an island beach on a clear, breezy, cloudless night, or stand on the spine of a barren mountain ridge after sunset, and behold the firmament of stars {A:^[1]} glittering against the coal-black sky above. They fill the sky with their timeless, brilliant flickering [(mostly caused by the terrestrial atmosphere)]. With binoculars or even a small telescope one finds that even the lacy dark matrix between the vast sea of stars is populated with still more stars that are simply too faint to be seen with the naked eye. Within the Milky Way galaxy, that stretches from horizon to horizon, the density of stars against the background sky is even greater. A:1

Each twinkling point of light is a star not too unlike our own Sun. The Sun is an ordinary star that features so prominently in our lives and on the pages of [the Heliophysics book series, as in this volume,] because of its proximity.

¹ Activity: Look up what type of astrophysical body is a true 'star'. Contrast that to 'white dwarf star', 'brown dwarf star', and 'neutron star': none of these are true stars in their present state and only two of which have ever been. 'Brown dwarfs' take up the mass interval between true stars and (exo)planets.

ⁱⁱ Throughout this work the original text from the Heliophysics volumes is directly quoted (with edits between brackets) like this: H-#[roman]:#[arabic].#[arabic]. So, for example, Vol. I, Section 9 in Chapter 2 would be referred to as H-I:2.9.

The next closest star, α Centauri (which is a triple system in which Proxima Centauri is currently the closest to Earth), is almost a million times farther away (at 4.22 light years), and the others are farther still. We may now say with some confidence that many of the stars are surrounded by planets of various sizes. {A:[2]} Some of these orbital companions are so immense that they are stars in their own right: double-star systems are quite common. {A:[3]}

A:2

A:3

With the same measure of confidence we may assert that most of these stars possess magnetic fields; that these magnetic fields create hot outer atmospheres, or coronae, that drive magnetized winds from their stars; and that these variable plasma winds blow past the orbiting planets, distorting their individual magnetospheres, and push outward against the surrounding interstellar medium. Where the ram pressure of the stellar wind becomes comparable to the surrounding pressure (gas, magnetic, and cosmic ray) of the interstellar medium, a bow shock forms. This serves to mark the farthest extent of the mechanical impact of the star on its surrounding environment: a sphere of influence, so to speak. [iii]

Our Sun has and does all of these things, and we refer to the sphere of influence carved out by the solar wind as our heliosphere. It is not really spherical and it varies in extent with solar activity. But in broad terms we may safely think of it as extending about 100 times further from the Sun than the Earth's orbit. We have yet to agree on the name for such spheres of influence around the other stars (for which astrospheres has been proposed), but there can be little doubt that such environments are as commonplace as the many points of light we see strewn across the sky on a dark and cloudless night."

[H-I:2.2] "Heliophysics encompasses the study of the various physical processes that take place within the sphere of influence of the Sun (*i.e.*, the heliosphere), and by analogy, those environments surrounding most other typical stars. But heliophysics also defines a specific method of study. This method embraces a holistic connected-system approach. It emphasizes a comparative context in which to understand a process by the many facets it presents in its various incarnations throughout the heliosphere. Taken together, each diverse

² Activity: Look up the definition of 'planet'. Note that, formally, the term 'planet' has only been defined by the International Astronomical Union for bodies within the Solar System; the term 'exoplanet' is reserved for bodies like planets in other planetary systems, although for these, and certainly for the joint collective, the term 'planets' is often used.

³ Activity: Many 'stars' we see in the night sky are binaries, including, for example, the brightest star in the night sky, Sirius (α CMa). More complex multiple-star systems may be less frequent, but are nonetheless common. Look up the example of Castor (α Gem) for an example of a sextenary, and then explore some more on multiple star systems in general.

ⁱⁱⁱ See Table 1.1 for definitions of the most common descriptors used for domains in, or phenomena related to, heliophysics. For a more extensive glossary of terms used in this volume, see, for example this NASA site; for a glossary of terms related to space weather, see this NOAA site.

- *active region*: a bipolar area of relatively strong magnetic flux, mostly consisting of magnetic *plage* and, by definition, containing one or more *sunspots* at some point in its evolution (*cf.* Fig. 4.4)
- *ast(e)rosphere*: equivalent of a heliosphere around another star (Sect. 10.3)
- *chromosphere*: domain above the Sun's visible 'surface', with temperatures around 10,000–20,000K (see Table 2.3)
- *corona*: the hottest domain of the Sun's atmosphere, at ≥ 1 MK (see Table 2.3)
- *coronal hole*: formally a coronal region that is dark in X-rays and EUV; generally identified with a region where the Sun's magnetic field is 'open', *i.e.*, reaches into the heliosphere (*e.g.*, Sect. 2.2)
- *coronal loop*: a high-temperature atmosphere within the Sun's corona, constrained to the volume of a magnetic 'flux tube' (*e.g.*, Sect. 3.4)
- *coronal mass ejection*: impulsive expulsion of magnetized material from a star into an astrosphere (*e.g.*, Fig. 5.1, Sect. 6.1)
- *current sheet*: defined in Table 3.1
- *exosphere*: outermost domain of an atmosphere in which collisions are rare and ballistic trajectories dominate for constituent particles (*e.g.*, Sect. 2.3)
- *facula* and *bright point*: a small flux tube in near-photospheric layers viewed towards the solar limb or disk center, respectively (*e.g.*, Sect. 9.1)
- *flare*: impulsive conversion of magnetic energy in a stellar atmosphere into thermal and non-thermal particles and bulk plasma motion, and appearing as a brightening over much of the stellar spectrum, although not significantly in total stellar brightness except for the most energetic events (*e.g.*, Sect. 6.1)
- *flux tube/rope*: defined in Table 3.1
- *heliosphere*: the extended region where the solar wind dominates over the interstellar medium (*e.g.*, Fig. 5.1)
- *ionosphere*: the ionized component of a planetary atmosphere, largely overlapping with the thermosphere (*e.g.*, Sect. 2.3)
- *magnetosphere*: a magnetic environment, generally of a planet, in which the intrinsic or induced magnetic field of the central body dominates over external fields or flows (*e.g.*, Fig. 5.1)
- *magnetospheric (sub-)storm*: a global disturbance in a planetary magnetic field driven by the solar wind, released through fast reconnection processes within the planetary magnetic field (*e.g.*, Sect. 6.2)
- *mesosphere*: at Earth, layer between stratosphere and thermosphere (Sect. 2.2)
- *photosphere*: 'surface' of a star, at the rapid transition from opaque to transparent
- *stratosphere*: at Earth, the domain between troposphere and mesosphere where temperature rises with height and convection is rare (*e.g.*, Sect. 2.2)
- *solar cycle*: quasi-cyclic variation in the number of sunspots seen on the solar surface when averaged over time scales of months (*e.g.*, Fig. 4.5)
- *(spectral, total) solar irradiance*: solar input into a planetary atmospheric system in the form of photons (*e.g.*, Sect. 2.2)
- *sunspot*: a 'flux tube' in the near-surface layers, with suppressed internal convection, and large enough to cool and appear dark (*e.g.*, Sect. 4.1.2)
- *thermosphere*: outer layers of a planetary atmosphere in which the temperature increases with height (*e.g.*, Sect. 2.2), specifically the neutral particles
- *(solar) transition region*: a domain between chromosphere and corona with a very strong temperature gradient dominated by conduction (see Sect. 9.2)
- *troposphere*: the lower layers of a planetary atmosphere (*e.g.*, Sect. 2.2)

Table 1.1. *Basic glossary for domains and phenomena in heliophysics.*

facet serves to fill out a complete and physically satisfying picture of a given process or phenomenon.

The physical processes and phenomena that we will encounter in [this volume] are themselves especially diverse. They include the rapid and efficient energization of thermal particles to suprathermal energies, the generation and annihilation of magnetic field, stellar variability and activity cycles, space weather, turbulent transport of energy and momentum, [the coupling between ionized and neutral atmospheres, and atmospheric chemistry,] to name just a few. Heliophysics fills a critical need to establish a unified science that connects these seemingly unrelated concepts in a manner that emphasizes complementarity over individuality, function over form, and generality over specificity.

Along with unification, coupling provides the second principal pillar upon which heliophysics rests. The heliosphere is a collection of coupled systems. It is fortunate that many of the linkages essentially operate only in one direction. That is to say, system A impacts system B, but B has little influence on A. Under these circumstances it is expedient to treat system A independent of the behavior of system B. This provides a certain economy of effort and scale, and it often reduces the (apparent!) complexity of a problem. For example, complex geomagnetic activity has no impact on solar flares, and the solar wind does not influence the Sun's cyclic variability.

Linkages, especially when several are present and working at cross purposes, can lead to confusion and spirited debate over what is a root cause and what is simply a resulting effect. The cause and effect relationship between solar flares and coronal mass ejections is a good case in point. Consider, for example, what the purported cause and effect relationship might be between a sore throat and a fever. Because a sore throat often starts before a fever develops one might be tempted to assign the effect to the fever and take the sore throat to be the cause. Fortunately, medical research informs us that both are effects and the root cause is the influenza virus. Heliophysics is needed to play this very same role in sorting out the appropriate relationships (or lack thereof) between any variety of physical effects that often occur contemporaneously throughout the heliosphere.

Solar variability does influence our climate here on Earth. This fact is certainly not negotiable in a purely scientific context and is arguably one of the most important linkages between the Sun and the Earth. Satellites have confirmed that the solar irradiance is variable on time scales from minutes to decades. The fluctuations are greatest on the shortest time scales. Day-to-day irradiance changes are on the order of a percent versus tenths of a percent over a solar cycle. The magnitude and sense of irradiance trends over centuries and

millennia are currently difficult to determine with any measure of certainty. Slow but steady progress on this question is being made through the studies of paleoclimate records. Over much longer time scales, stellar evolution theory provides assurances that significant changes in solar irradiance have taken, and will take, place with dramatic impacts on our climate and way of life.

What is debatable, however, is precisely what the direct relationship is between solar variability and climate change over any particular time scale, or epoch, of interest. For example, various opinions have been advanced that span the entire gamut from wholly inconsequential to complete solar responsibility for the gradual warming of the planet that has been observed since the middle of the 20th Century. Yet, it may not even make sense to speak of direct relationships between drivers and the behavior of systems which are as nonlocal, nonlinear and plagued by various hystereses as is our climate here on Earth.

The third and final pillar upon which heliophysics rests is the exploration of Earth's neighborhood in space. As a space-faring civilization we have visited all the planets, [several asteroids and] comets and numerous planetary satellites. We have ventured to the boundaries of the heliosphere and have flown through various parts of our magnetosphere. We have a spacecraft [that passed] the Pluto/Charon system [and after that flew by Kuiper-belt object 2014 MU69, colloquially known as Ultima Thule]. Heliophysics enables our exploration to be successful and at the same time gains in knowledge and understanding from our exploration initiatives.

In summary, heliophysics is the systems-mediated study of the physical processes that take place within the Sun's sphere of influence. It is based upon the three pillars of unification of physical processes and phenomena, coupling of distinct physical systems, and the exploration of our neighborhood in space. And it is broadly applicable to the environments around most ordinary stars."

1.4 The language of heliophysics

[H-I:2.3] "The language of heliophysics is mathematics. And the body of literature from which heliophysics draws its substance and in turn records its accomplishments is the physics of magnetized plasmas. With only a rudimentary knowledge of a language, a literature is incomprehensible, except, perhaps in translation. And even in translation so much of the original meaning and the nuance the author wished to convey are inevitably lost, or worse, misinterpreted by even the most conscientious translator.

The most precise, and intellectually demanding, literary prose of heliophysics assigns a phase-space distribution function to each individual species of particle. By a species one may simply mean free electrons, protons, or oxygen molecules,

or even photons. In some applications it might be necessary to distinguish between oxygen molecules in different excited (vibrational, rotational and electronic) states, or between iron atoms at different stages of ionization, or between different senses of photon polarization. In any case, the evolution of each distribution function is obtained by setting the total time derivative equal to the net production/loss of an individual species by various collisional or radiative processes. Such evolution equations are commonly referred to as Boltzmann, or Vlasov equations. When there is no net gain or loss, then Liouville's Theorem asserts that the vanishing of the total time derivative of the distribution function conserves the phase space density for each species [(see Sect. 8.3)].

In specifying the total time derivative it is necessary to determine the forces acting upon a given particle species. For uncharged particles, gravitational attraction is the only important consideration. Accordingly, to the system of equations for the individual distribution functions one must add Poisson's equation in order to specify the gravitational field based on the mass distribution provided by those particles with mass. Charged particles are also subject to electromagnetic interactions. Thus we must also include Maxwell's equations to deduce the electric and magnetic fields based on the distribution of charges and currents provided by the charged particle species. {A:[4]}

A:4

In principle, this suffices to provide a complete description of the grammar and syntax of heliophysics at a very elegant, learned and precise level. In practice the task of following through with this program (a) is prohibitively difficult with or without the assistance of the computer, (b) is subject to the problem that the initial conditions are not known with any degree of certainty, (c) is complicated by the fact that many of the collisional and radiative transition probabilities are not even approximately known, and (d) requires that certain conditions be fulfilled so that electromagnetic interactions can be separated into large-scale fields and small-scale collisions. Finally, this comprehensive description usually provides far more information than is usually necessary for comparing with observations or understanding the predictions of a theory over specified temporal and spatial scales.

At the opposite extreme from the scholarly literary prose is the common vernacular. For heliophysics, if high literary prose centers on Poisson, Maxwell, Boltzmann and Vlasov, then the vernacular is single-fluid, ideal, magnetohydro-

⁴ Activity: Remind yourself of Maxwell's equations that are mathematical renditions of these properties: (1) electric monopoles are linked with an electric field; (2) there are no magnetic monopoles; (3) variations in the magnetic field are associated with an electric field; and (4) a magnetic field implies either steady currents or time-dependent electric fields, or both. Good news: once we have reached magnetohydrodynamics in Ch. 3, Maxwell's equations are in principle superfluous as they are contained within the MHD equations; if you are interested in how that works, see here (Sections 1.1.1–1.1.9). By the way, a really useful resource for all things related to plasma physics (and how to convert between different unit systems) is the online NRL Plasma Formulary.

dynamics, or MHD for short (see Ch. 3). MHD is a continuum fluid description that does not distinguish between particle species, averages (in some sense) over particle collisions, ignores radiative effects altogether, and is based on velocity moments of the underlying distribution functions. It retains Poisson without modification, but takes certain liberties with Maxwell. Boltzmann and Vlasov drop out of the picture entirely.

MHD can be rigorously derived from Poisson/Maxwell/Boltzmann/Vlasov under various conditions that are not altogether unreasonable for very many heliophysical applications. Usually this involves following the behavior of a physical process or phenomenon over coarse-grained spatial and temporal scales. In other words, it is a useful, and indeed often very accurate, description of the 'big picture'. Because of its relative simplicity, ideal MHD provides a useful context in which to interpret and understand the behavior of magnetized plasmas at a basic and often extremely intuitive level. On the other hand, ideal MHD is often applied to processes or phenomena to which it does not actually apply. Generally speaking, if collisional and radiative relaxation times are short compared to the coarse-grained time scale of interest then ideal MHD is likely to be a reasonable option. But 'gotchas' are always present.

The successful derivation of the MHD equations requires a closure prescription, which may be regarded as a consequence of the familiar, 'no free lunch' maxim. Closure entails specifying a tractable procedure to determine the pressure tensor (second-order velocity moment) in terms of the fluid density (zeroth-order velocity moment), the bulk fluid velocity (first-order velocity moment), and the magnetic field. The so-called polytropic approximation—in which the pressure is a scalar proportional to the particle number density raised to a specified power—is the simplest option. A power law index of unity corresponds to an isothermal process (constant temperature). A power law index equal to the ratio of specific heats describes an isentropic (constant specific entropy) process that also manages to conserve energy. More complicated options are possible and are often tailored to accommodate specific situations. A successful and accurate closure scheme is inevitably based on some additional *a priori* knowledge of the behavior of the particle trajectories, or the general nature of the particle distribution functions.

In contrast to the Poisson-Maxwell-Boltzmann-Vlasov description, ideal MHD is a system of nine partial differential equations for nine dependent variables [(shown in Table 3.3 and discussed in Ch. 3)]: the gravitational potential, the fluid density and pressure, the fluid velocity (3 components) and the magnetic field (3 components). These equations are (a) the Poisson equation to describe gravity, (b) the continuity equation expressing the conservation of mass, (c) the closure relation to specify the pressure tensor, (d) the equation for

the conservation of momentum, or the force-balance equation (3 components), and (e) the magnetic induction equation (3 components).

Of course, between ideal MHD and the Poisson-Maxwell-Boltzmann-Vlasov description lies a vast real estate filled with a plethora of compromise or hybrid descriptions. The number of such schemes is limited only by the imagination and ingenuity of the investigators. Multi-fluid treatments allow for individual densities, velocities and pressures associated with different particle species or groupings of particle species, but retain a single gravitational potential and magnetic field applicable to every fluid. This formulation is useful when the time scales of interest are short compared to characteristic inter-species collisional relaxation times, but long compared to the analogous intra-species times.

Another intermediate scheme employs high-order moment closures. These schemes are necessary when the species distribution functions deviate significantly from the fully-relaxed Maxwellian. Often this situation occurs when significant spatial gradients are imposed on the system. Additional partial differential equations are then used to describe the time-evolution of the components of the pressure tensor. The closure is postponed to the next higher level of the heat flux tensor (third-order velocity moment), or in extreme circumstances to even higher-order moments.

Hybrid schemes treat some species as fluids and retain a Boltzmann-Vlasov - or kinetic - description for others. Indeed even a single species of particle may be partitioned in such a fashion that some of the particles are treated kinetically (generally the high energy suprathermal tail of the distribution function) while the remainder are described as a fluid (the thermal core of the distribution). Such schemes are particularly useful in describing the energization of charged particles; [we see this in action in Ch. 8].

In summary, there is a bewildering array of schemes that are presently invoked to describe the behavior of magnetized plasmas in the heliosphere. They encompass an extremely wide range of complexity. Each is specifically tailored to a given physical process and phenomenon. They are not simply interchangeable, but have their own individual strengths and weaknesses. One should always choose the simplest description that will suffice for understanding the problem in hand. Use all the information and knowledge you have at your disposal about the nature and behavior of a physical system in selecting a scheme. If the heliophysics concepts can be adequately framed in the common vernacular, then eschew the sophisticated flowery prose unless nothing less will do.”

1.5 A timeline of exploration of planetary systems

NASA's Heliophysics Division within the Science Mission Directorate was previously known as the Sun-Earth Connections Division. That earlier name reflected that much of its research focused on how solar activity impacts our home planet. As probes explored ever more of the solar system, researchers realized that learning about the science of terrestrial space weather and of the evolution of Earth's climate system was boosted by the incorporation of discoveries from around the solar system; the name change of the Division reflected the shift to a broader perspective that was already taking place in the research community. As exoplanets were found to be more common than stars, the application of the science of heliophysics to the exploration and understanding of processes in exoplanetary systems, and in particular to exoplanetary habitability, presents a natural development of the discipline. The multi-disciplinary science arena that looks into star-planet couplings has accelerated rapidly alongside astronomical exploration.

In 1969, half a century ago, astronauts first landed on Earth's sole moon. The first successful robotic landers touched down on the much more distant Venus and Mars in 1970 and 1976, respectively, and in the same decade spacecraft flybys provided the first, fleeting close-ups of Jupiter and Saturn. It was not until two decades later, however, that missions that explicitly targeted these giant planets revealed how fundamentally distinct these worlds are from our own.

The Galileo satellite started exploring the Jupiter system in late 1995, swinging by moon after moon. The Cassini-Huygens mission reached Saturn in 2004, exploring the giant planet, its rings and satellites, and even sending a lander onto Titan, the only moon in the solar system with a substantial atmosphere. These spacecraft uncovered a fascinating diversity of environments on dozens of moons: many are cold worlds enrobed in miles-thick ice; some with volcanoes spewing molten rock but others whose volcanoes somehow gush liquid water or nitrogen; and then there is Titan with its seas of liquid methane and ethane. Their pictures are as stunning and diverse as the scientific discoveries enabled by these spacecraft. The far reaches of the Solar System continue to offer surprises: dwarf planets Haumea and Makemake, objects in the distant Kuiper belt, were not discovered until 2004 and 2005, respectively.

As the close-up exploration of the largest planets in the solar system got underway, a revolution was about to befall astronomers looking much further out. It started in 1995 with the announcement of the first exoplanet, now known as 51 Pegasi b, orbiting a star like our own Sun. There are now over 4,000 exoplanets on the books ^[iv] (more than half found with NASA's *Kepler*

^{iv} See <https://exoplanetarchive.ipac.caltech.edu> and <http://exoplanet.eu>.

satellite), but the number expected to exist is vastly larger: by carefully quantifying what our available methods can and cannot observe, scientists estimate that there are over a hundred billion planetary systems in our Milky Way galaxy alone, with perhaps of order ten billion planets with some similarity to Earth.

Apart from its very existence, 51 Pegasi b had another surprise in store: at 150 Earth masses and orbiting its star almost 20 times closer than Earth does the Sun, this hot Jupiter should not have existed by theories of the time. These and many subsequent observations have changed our ideas on how planetary systems form and evolve: we now realize that orbits can change so that planets may be discovered well away from where they formed; planets can engage in gravitational fights that can cause losers to be ejected as lone 'nomads' into interstellar space; planets exist that have two stars to cast twin shadows on their surfaces; . . . Many planets orbit their stars at distances where water, if there is any, may exist in liquid form on their surface for billions of years, as on Earth where it enabled the development of life.

These discoveries have intensified the astronomers' hunt for extraterrestrial life in which also solar-system scientists participate. Organic molecules cause the haze in the icy-cold atmosphere of Saturn's Titan and are vented in cryovolcanic plumes rising from the ice-locked deep ocean of nearby Enceladus. There are many sizable moons and dwarf planets in the solar system that are rich in water, although much of it is frozen solid. The combination of liquids and organics in many places around our solar system fuels theories of life and plans for space missions designed to look for it near to home.

But exoplanet astronomers have the advantage of the vast number of systems. Their challenge is that even the largest telescopes can image exoplanets no better than as an unresolved blur the size of the instrumental point spread function, if indeed they can separate the reflected light from the exoplanets from the light of the stars that they orbit. In fact, most of what we learn about exoplanets comes from analyzing how their star's light is modified in brightness or color by the exoplanets, either by adding some reflected starlight or by taking away some light should they move in front of their star during their orbit. Careful study of these effects as observed with the most powerful telescopes can reveal which gases contribute to the changes. This is already happening, but it will receive a big boost from future telescopes being built, including NASA's James Webb Space Telescope planned for launch in 2021. So much was discovered in the most recent few decades; what will the next several decades bring?

2

Neutrals, ions, and photons

2.1 Conditions in the local cosmos

The local cosmos discussed in this book exhibits an enormous diversity of conditions. Figure 2.1 is one perspective of this in its comparison of number densities and temperatures: densities range over more than 28 orders of magnitude (more than the contrast between solid rock and the 'vacuum' of low-Earth orbit) and temperatures over 5 orders of magnitude. The magnetic field, another crucial parameter that is explored starting in Ch. 3, provides another dimension and adds its own physical processes. All together, these physical parameters cover a wide range of states that include solids, liquids, gases, and ionized and magnetized particle ensembles called plasmas.

Matter in most of the domain of heliophysics is at least electrically conducting, but generally at least partially or even fully ionized as will be abundantly clear from the chapters in this volume. Ionization can be a consequence of high-speed collisions between particles in a hot medium and/or of high energies in the thermal radiation associated with high temperatures. A hot medium can result from the transport and conversion of different forms of energy where a balance of thermal sources and sinks may only be reached at high temperatures. Examples of such settings are the interior and the atmospheric domains of the Sun. In these environments, internal collisional ionization and recombination, as well as excitation and de-excitation processes dominate in balancing ionization and recombination rates. Alternatively, ionization can be the result of impacts of externally-generated high-energy particles (such as solar energetic particles or particles accelerated in a planetary magnetosphere) or be caused by irradiation by solar photons of sufficiently high energy (typically X-ray and [extreme] ultraviolet) such as occurs in planetary ionospheres and cometary tails. {A:[5]}

A:5

⁵ Activity: Planetary lower atmospheres are dominated by molecular substances, transitioning to atomic elements with a relatively low admixture of ions and electrons as one moves up through the ionospheres and thermospheres, while magnetospheres and the solar outer atmosphere and wind are

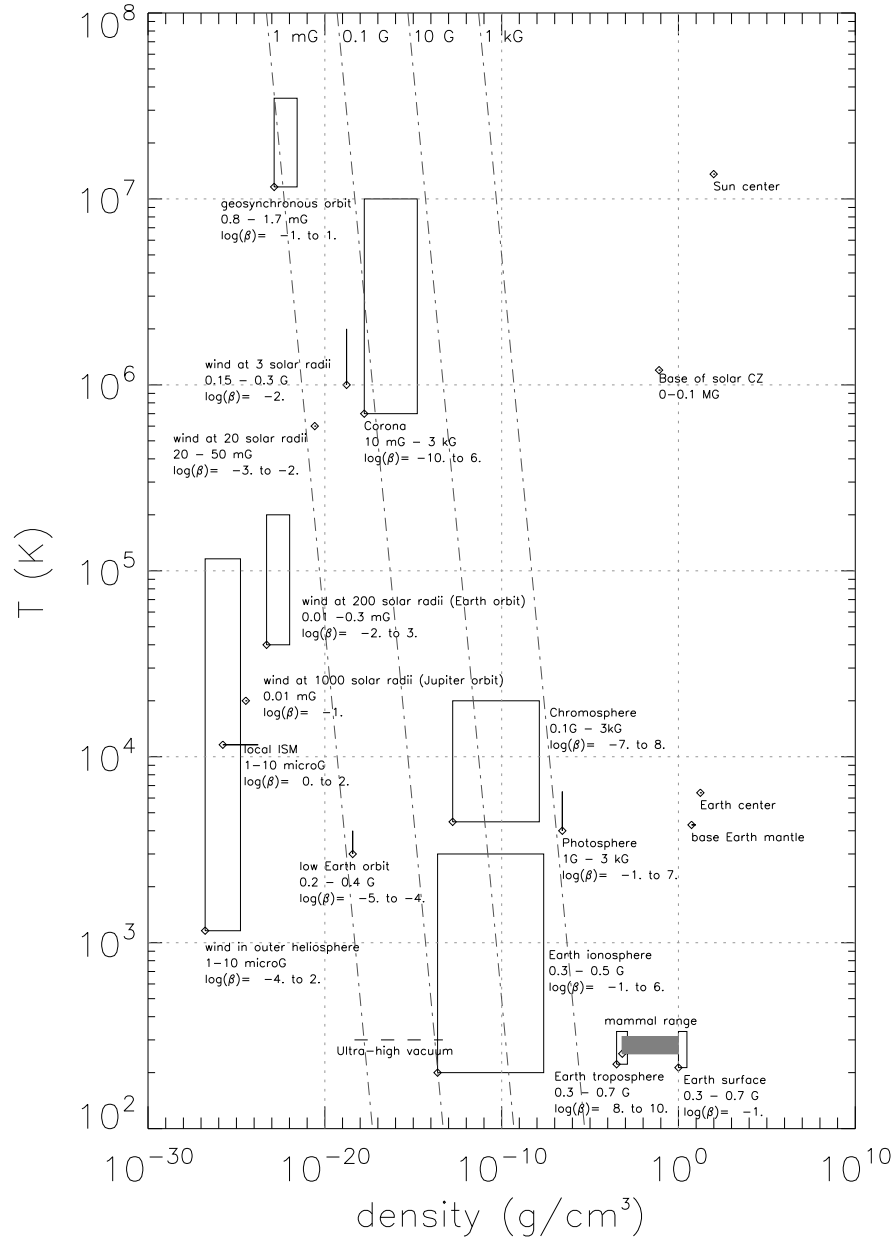


Fig. 2.1. Temperature versus mass density for a variety of conditions within the local cosmos. Some typical ranges are indicated, and labeled with magnetic field strengths (in Gauss) found in that domain, followed by estimated ranges of the plasma β , i.e., the ratio of energy density in plasma over that in the magnetic field (Eq. 3.24), in this scaling for a fully ionized hydrogen-dominated plasma. [Fig. H-I:1.1]

Table 2.1. *Present characteristics and climates of the terrestrial planets.*
 [Modified after Table H-III:7.1, with added surface gravity, escape velocity, and
 escape energies E_{esc} for protons and atomic oxygen.]

	Venus	Earth	Mars
Radius	6050 km	6400 km	3400 km
Orbital radius	0.72 AU*	1 AU	1.52 AU
Rotation period	243 days	24 hours	24.6 hours
Surface gravity	8.9 m/s ²	9.8 m/s ²	3.7 m/s ²
Escape velocity	10 km/s	11 km/s	5 km/s
E_{esc} for H ⁺ , O	0.5, 9 eV	0.6, 10 eV	0.1, 2 eV
Surface temp.	740 K	288 K	210 K
Surface pressure	92 bar	1 bar	7 mbar
Composition	96% CO ₂ 3.5% N ₂	78% N ₂ 21% O ₂	95% CO ₂ 2.7% N ₂
H ₂ O content	20 ppm	10,000 ppm	210 ppm
Precipitation	None at surface	Rain, frost, snow	Frost
Circulation	1 cell/hemisphere; quiet at surface but very active aloft	3 cells/hemisphere; local and regional storms	1 cell/hemisphere or patchy circulation; global dust storms
Maximum surface wind	~3 m/s	>100 m/s	~ 30 m/s
Seasons	None	Comparable northern and southern seasons	Southern summer more extreme

* An AU, or Astronomical Unit, is the average distance between Sun and Earth.

Much of what is described in this volume deals with the physics of magnetized plasmas, and much of that physics is approximated by a description known as magnetohydrodynamics, or MHD, as introduced in Ch. 3. In the present chapter, however, we first look at the more familiar situation of neutral gases, also because many of the phenomena discussed in this volume occur in the layers of planetary atmospheres for which the concept of hydrodynamics – in which magnetic field is ignored – gives us a good starting point. Later in this chapter, we focus on where ionization becomes important. For now disregarding the effects of magnetic fields, the limits of pure ('non-magneto-') hydrodynamics are reached in the high tenuous layers of planetary atmospheres where collisions are infrequent and other processes enter into our discussion, such as chemical differentiation subject to gravity or even outflows from the body in question.

A great variety of phenomena in the local cosmos have their foundation in

comprised predominantly of charged particles. Compare thermal kinetic energies in different settings with molecular binding energies of, say, water and carbon dioxide. Also compare the energy of X-ray and EUV photons with ionization energies of atomic hydrogen and oxygen. See Tables 2.1, 2.3, and 2.4 for conditions in different settings.

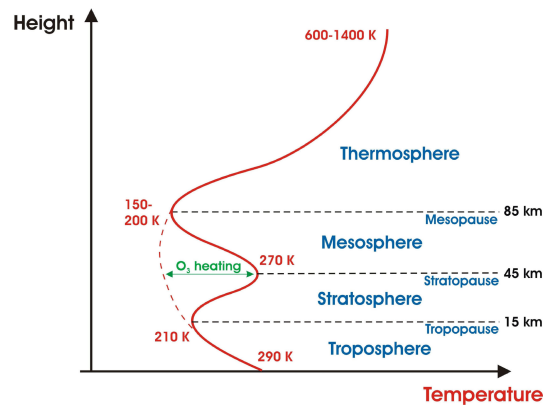


Fig. 2.2. Average vertical temperature profile through Earth's atmosphere. The general shape of the temperature profile is reasonably consistent to the point where it can be used to define the four main neutral atmosphere 'layers', from the troposphere to the thermosphere. The temperature of the uppermost layer, the thermosphere, increases steeply with altitude due to absorption of solar extreme ultraviolet (EUV) and far-ultraviolet (FUV) radiation. The thermosphere and upper mesosphere are partially ionized by the same EUV radiation, which varies by a factor of three over the solar cycle, and by auroral particle precipitation. [The effect of absorption by ozone is specifically highlighted. Fig. H-I:12.1]

the electrical conductivity of the media within which they occur. This may be in the generation and maintenance of magnetic field deep inside the Sun and in most of the planets, in the many phenomena driven by the interaction of the magnetized flow of the solar wind with Solar-System bodies, or even in the processes in the ionized domains of atmospheres of many of these bodies. In most situations discussed in this volume, that conductivity has its origin not in the metallic behavior of the medium as it does deep inside Earth but rather in the ionization of matter: whereas in metals ions are relatively immobile and share some of their electrons, in a plasma both the ions and the electrons are entirely unbound on microscopic scales. This chapter introduces electrical conductivity in a magnetized medium, here looking at plasma with a low degree of ionization; fully ionized plasmas are discussed in Ch. 3.

2.2 Gravitationally stratified atmospheres and stellar winds

Among the planetary atmospheres in the Solar System, those of Venus and Mars are most similar to those of Earth. The abundances of their primary constituents – mostly CO₂ and N₂, and, on Earth, N₂ and O₂ – are compared in Table 2.1. Note that the order of the most abundant components as well as the absolute base pressures differ markedly.

A sketch of the Earth's atmospheric vertical thermal structure is shown

Table 2.2. *Extent and important species for upper atmospheric regions of terrestrial planets [Table H-IV:7.3; added planetary radii R_p (km)].*

	Venus $R_{p,\text{♀}} = 6052$	Earth $R_{p,\text{♁}} = 6378$	Mars $R_{p,\text{♂}} = 3396$
Thermosphere	~120-250 km CO ₂ , CO, O, N ₂	~85-500 km O ₂ , He, N ₂	~80-200 km CO ₂ , N ₂ , CO
Ionosphere	~150-300 km O ₂ ⁺ , O ⁺ , H ⁺	~75-1,000 km NO ⁺ , O ⁺ , H ⁺	~80-450 km O ₂ ⁺ , O ⁺ , H ⁺
Exosphere	~250-8,000 km H	~500-10,000 km H, (He, CO ₂ , O)	~200-30,000 km H, (O)

in Fig. 2.2. The temperature gradually drops from the surface – where the bulk of the conversion of solar irradiance into heat occurs – through the troposphere due to adiabatic expansion. At greater altitudes the absorption of short-wavelength sunlight by tenuous gas that is less efficient in cooling through radiation leads to increased temperatures in the stratosphere (mainly by photons between about 2,000 Å and 3,000 Å) and in the thermosphere (for wavelengths mostly short-ward of 2,000 Å). Energy leaves the Earth’s atmospheric domains mainly by infrared radiation from the lower regions, which also leads to a decrease in temperature above the stratosphere by radiation from the mesosphere. The densities in the thermosphere are so low, and the dominant chemical constituents such inefficient radiators, that downward thermal conduction exceeds radiative losses above about 100 km (see Ch. H-IV:9). Table 2.2 compares the properties of the upper atmospheres of the three terrestrial planets (with significant atmospheres), *i.e.*, the thermospheres, the ionized constituents referred to as the ionospheres that largely overlap with the thermospheres, and the exospheres beyond that; the reasons for the apparent chemical mismatch between the neutral molecular and the ionized components are discussed in Ch. 13. ^[v]

For the Sun’s atmosphere, there is a comparable pattern of temperature with height: moving upward, the temperature drops throughout the lower atmosphere (the ‘photosphere’ from which the bulk of the solar irradiance is emitted; also referred to as the ‘solar surface’ by astronomers, despite the fact that the Sun is entirely gaseous throughout), but then increases again in the chromosphere (extending a few thousand km above the photosphere) and then shoots up to form an extended, hot corona. Some of the physical properties of these domains (along with a rough definition of the terms) are summarized in Table 2.3. The reasons behind this similarity in pattern are

^v This volume focuses on terrestrial planets; we refer to Ch. H-IV:8 for an introduction to the upper atmospheres of the giant planets.

partly the same, partly completely different. A similarity is that energy is most efficiently radiated from the low, dense atmospheric layers, and poorly from high, tenuous layers where conductive redistribution plays an important role. But the heat input differentiates the two: the solar chromosphere and corona are not heated by absorption of photons from the solar surface (which is thermodynamically impossible because the atmospheric temperature is higher than the surface temperature) but by dissipation of electrical currents and a variety of waves running through the plasma (both generated by the convective flows below the solar surface, and coupled into the outer atmosphere via the Sun's magnetic field; see Sect. 9.3). The amount of energy converted in the solar outer atmosphere from chromosphere to corona and solar wind is a function of the instantaneous magnetic activity. This activity exhibits an 11-year quasi-cyclic pattern that is often referred to as 'the sunspot cycle' because it was discovered from multi-decade records of sunspot counts.

The Sun's radiative input into the Earth's atmosphere (known as the spectral irradiance, $S(\lambda)$) exhibits a significant variability depending on solar magnetic activity (Figure 2.3). The overall emission from the solar photosphere varies little with magnetic activity, that from the warm chromosphere mildly, and that from the hot corona strongly. {A:[6]} As a result, the relative variability in $S(\lambda)$ through the solar cycle increases markedly short-ward of about 3,000 Å with $(S_{\max} - S_{\min})/S_{\min}$ climbing from about one part in 1,000 or less long-ward of that to near unity short-ward of 1,000 Å. The absorption of the most variable segment of the spectral irradiance in Earth's atmosphere occurs primarily above about 50 – 100 km (Fig. 2.4), causing the high atmosphere to evolve strongly in temperature and density in response to the solar sunspot cycle (see Fig. 2.5), further modulated as Earth goes through its weakly elliptical orbit around the Sun and its rotation about a tilted axis, and with variable contributions from geomagnetic activity (see Chs. 12 and 13).

A:6

Wavelengths short-ward of about 2400 Å and about 1250 Å can dissociate O₂ and N₂, respectively, and short-ward of about 900 Å can ionize, *e.g.*, O atoms. Consequently, the atomic and ionic components in the Earth's atmosphere do not show up significantly below around 100 km in altitude because all ionizing and dissociating wavelengths have been absorbed by that depth into the atmosphere (see Ch. 13); above that altitude, the abundances of the ionic and atomic components all reflect solar, orbital, and diurnal cycles.

The density stratification in much of the lower atmosphere of the terrestrial planets (defined as below about 100 km for Earth) can be understood to first

⁶ Activity: Look up and compare images of the Sun's magnetic field and atmosphere in different phases of the solar cycle, such as those obtained with the *HMI* and *AIA* instruments on *SDO* (NASA's *Solar Dynamics Observatory*). Note that such images are typically in false color, and with non-linear intensity scales to accommodate brightness contrasts.

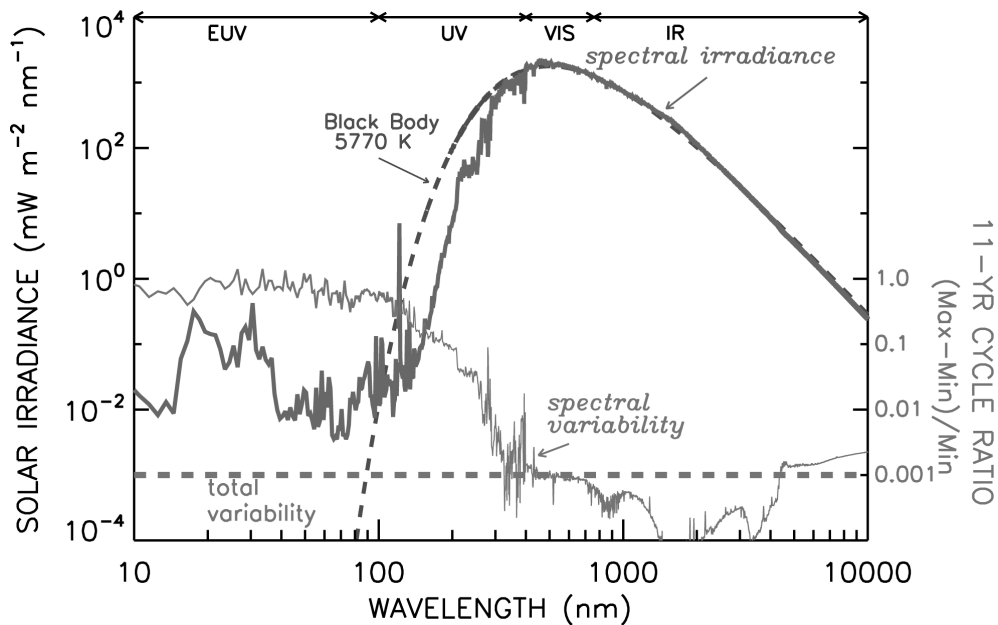


Fig. 2.3. Comparison of the solar spectrum and the black body spectrum for radiation at 5770 K (the approximate temperature of the Sun's visible surface). Also shown is an estimate of the variability of the solar spectrum during the 11-y solar cycle, inferred from measurements (at wavelength below 4000 Å) and models (at longer wavelengths) and, for reference (dashed line), the solar cycle 0.1% change in the total solar irradiance. [Fig. H-III:10.1]

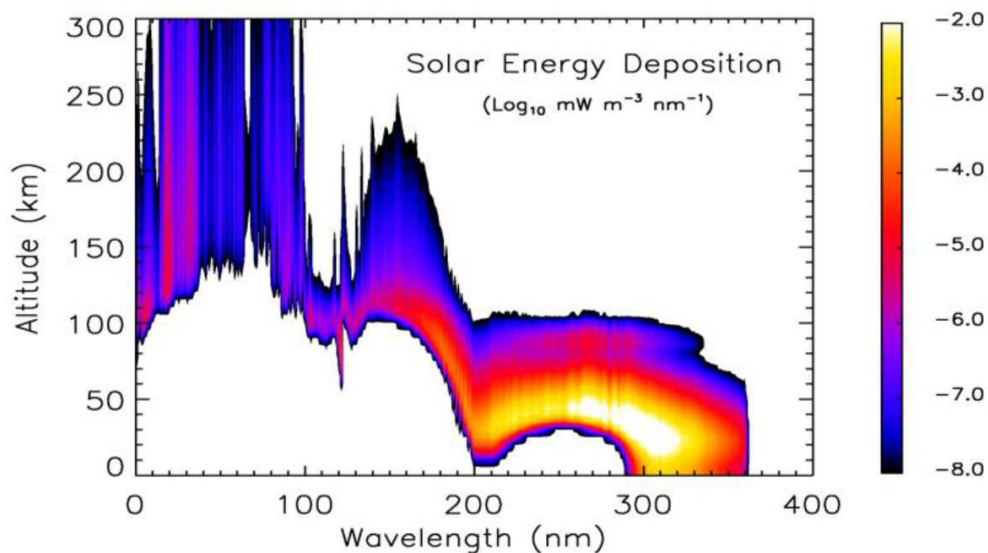


Fig. 2.4. [Altitude of penetration of the solar radiation as a function of wavelength [from X-rays through 3600 Å]. The color range shows the amount of energy deposited in the different layers of the atmosphere for the different parts of the solar spectrum (on a logarithmic scale, in units of $\text{mW}/\text{m}^3/\text{nm}$ [or $10^{-3} \text{ erg/s}/\text{cm}^3/\text{Å}$]). [Fig. H-III:13.3]

Table 2.3. *Basic parameters for, and definitions of, domains in the solar atmosphere. Note that all regions of the solar atmosphere are very inhomogeneous and that these values are only meant to give a rough idea of their magnitudes. [Table H-I:8.1, here converted to cgs-Gaussian units, and with solar properties added. n_e and n_H are the densities of electron and neutral hydrogen; the plasma β is defined in Eq. (3.24)]*

Region	n [cm ⁻³]	n_e/n_H	T [K]	B [Gauss]	β
Photosphere ¹	10 ¹⁷	10 ⁻⁴	6 10 ³	1 – 1500	> 10
Chromosphere ²	10 ¹³	10 ⁻³	2 10 ⁴ – 10 ⁴	10 – 100	10 – 0.1
Transition region ³	10 ⁹	1	10 ⁴ – 10 ⁶	1 – 10	10 ⁻²
Corona ⁴	10 ⁸	1	10 ⁶	1 – 10	10 ⁻² – 1

Sun: radius $R_\odot = 710^5$ km; surface gravity $g_\odot = 274$ m/s²; bolometric luminosity $L_{\text{bol}} = 410^{33}$ erg/s; effective temperature $T_{\text{eff}} = 5772$ K, defined such that $L_{\text{bol}} \equiv \sigma T_{\text{eff}}^4 4\pi R_\odot^2$

Definitions: ¹ the *photosphere* is the layer from which the bulk of the electromagnetic radiation leaves the Sun (this layer has an optical thickness $\tau_\nu \lesssim 1$ in the near-UV, visible, and near-IR spectral continua, but it is optically thick in all but the weakest spectral lines); ² the *chromosphere* is optically thin in the near-UV, visible, and near-IR continua, but optically thick in strong spectral lines - it is often associated with temperatures around 10,000 – 20,000 K; ³ the *transition region* is a thermal domain between chromosphere and corona in which thermal conduction leads to a steep temperature gradient; ⁴ the *corona* is optically very thin over the entire electromagnetic spectrum except at radio wavelengths and in a few spectral lines - the term is often used to describe the solar outer atmosphere out to a few solar radii with temperatures exceeding ~ 1 MK.

order by looking at the behavior of a stationary gas subject only to gravity. [H-I:12.2.1] “The frequent collisions of molecules in a gas close to thermal equilibrium enable the Maxwellian [velocity distribution (with corresponding exponential energy distribution)] of the individual particles to be characterized by the basic fluid properties of pressure, p , temperature, T , number density, n , and mass density, ρ , that are related by the perfect gas law:

$$p = nkT = (\rho/m)kT = \rho\mathcal{R}T/\mu, \quad (2.1)$$

where k [(1.4 $\times 10^{-16}$ erg/K)] and \mathcal{R} [(8.31 $\times 10^7$ erg/K/mol)] are the Boltzmann and universal gas constants, respectively, and m is the mean molecular mass [while μ the mean molecular mass in atomic units]. {A:[7]} {A:[8]} The

A:7

⁷ Activity: Consider why for a fully-ionized, hydrogen-dominated plasma we see $p = 2nkT$. For the answer, see below Eq. (2.7).

fluid concept of pressure in the atmosphere represents the weight of the column of gas above. A:8

The neutral gas under the influence of the planet's gravitational force gives rise to the concept of hydrostatic balance, which states that the change in pressure with height, dp , is closely balanced by the weight of the fluid, $nmgdh$ (where m is the mean molecular mass in [grams] and h is the height), under the action of the planet's gravitational acceleration, g . The concept is expressed mathematically as:

$$\frac{dp}{dh} = -\rho g = -p/H_p. \quad (2.2)$$

This basic equation describes the exponential decrease in gas density with altitude, and results in the concept of the pressure scale height,

$$H_p = kT/mg, \quad (2.3)$$

which represents the [height difference] through which the gas pressure [in an isothermal atmosphere] will decrease by a factor of $e = \exp(1)$. Earth's upper atmosphere extends for about a dozen scale heights above 100 km altitude, with scale heights changing from about 5 km to 50 km with increasing altitude, as the temperature increases from about 180 K to over 1000 K (see Fig. 2.2). A:9

[The] quasi-equilibrium implied by hydrostatic balance does not exclude the possibility of vertical winds. The assumption simply demands that the rate of [flow] is such that the atmosphere adjusts at a comparable rate. The term quasi-hydrostatic balance is the more correct expression in the case of accommodating vertical winds in the system. [...] Vertical winds in Earth's upper atmosphere of the order of 100 m/s can be accommodated within the quasi-hydrostatic assumption." A:10

The quasi-hydrostatic description applies not only to planetary atmospheres A:11

- ⁸ Activity: At the solar surface we see a mean 'molecular mass' of $m \approx 1.3m_p$ while in the fully-ionized corona $m \approx 0.6m_p$ (for proton mass m_p). Explain why. (A hint: see Fig. 2.10.)
- ⁹ Activity: Compute scale heights H_p in the Earth's atmosphere for molecular nitrogen (the dominant component) at a range of temperatures, and compare these with the value $H_{p\odot}$ for the atomic hydrogen-dominated gas in the solar photosphere, and for the CO₂-rich atmospheres of Venus and Mars. Use the data in Tables 2.1 and 2.3. Consider how the value of $H_{p\odot}/R_\odot$ contributes to the appearance of the Sun as having a well-defined surface. Also, consider why neutral, atomic hydrogen dominates in the solar photosphere (see Fig. 2.10 for the answer).
- ¹⁰ Activity: One way to quantify the 'strength' of storms in different planetary atmospheres is to compare the dynamic pressure ρv^2 for the maximum surface winds listed in Table 2.1. Compare those values with the dynamic pressure in the solar wind using Table 2.4. Note: 1 bar = 10⁶ dyne/cm².
- ¹¹ Activity: The fastest flows (in any direction, not only in the vertical, gravitationally stratified direction) that can be accommodated in a quasi-hydrostatic situation can be estimated from the fact that the gas pressure $p = \zeta nkT$ (with $\zeta = 1$ for a neutral gas and $\zeta = 2$ for a fully ionized hydrogen gas) should well exceed the flow's dynamic pressure ρv^2 . Look at Fig. 2.1 and add the horizontal lines where the two pressure terms are equal for a variety of flow velocities and corresponding temperatures; compare with the conditions discussed later in this chapter for the solar wind.

but is also used for the interior of the gas giants, for the interior and lower-atmosphere of the Sun, and – as we shall see later – even inside magnetic ‘containers’ in the solar atmosphere that are known as flux tubes (see Table 3.1 for a definition), one incarnation of which are ‘coronal loops’ – which is the general term describing the emitting structures seen in EUV and X-ray images of the Sun’s hot outer atmosphere. Table 2.3 summarizes characteristic physical parameters for the domains within the solar atmosphere from the solar ‘surface’ (photosphere) up into the corona. These numbers should be seen as characteristic values only: all these domains span a few orders of magnitude in density and all are very dynamic at any given location, while moreover the solar magnetic field plays a key role in them as it structures multitudes of adjacent distinct atmospheres along magnetic field bundles (Sect. 3.4). The solar corona is visible at X-ray and EUV wavelengths up to several hundred thousand kilometers. The coronal plasma is mostly contained in magnetic structures relatively low down, but increasingly with height the gas pressure forces the magnetic field to ‘open’ into the heliosphere. The plasma on ‘open field’ streams out to form the solar wind, resulting in a low-density and consequently dark lower coronal region known as a ‘coronal hole’.

The quasi-hydrostatic description even forms a useful, albeit very crude, approximation for that part of extended atmosphere of the Sun that is the inner-heliospheric domain of the solar wind: whereas there is in fact an outflow, this ‘vertical wind’ leaves the stratification nearly hydrostatic for many solar radii above the solar surface, as we shall see shortly.

Table 2.4 summarizes a few characteristics of the solar wind near Earth orbit. Outside of dynamic coronal mass ejections (Ch. 6), the solar wind is predominantly in one of two states, referred to as the ‘fast wind’ and the ‘slow wind’. These states originate from distinct environments on the Sun, and because the Sun rotates underneath the radially outflowing wind, slow and fast streams unavoidably interact – see Sect. 5.5.1.1. For what follows here, we focus on domains where only one of these types of wind prevail for several days, which is the time it takes to flow from Sun to Earth (the geometry of the magnetic field that it carries is discussed in Sect. 5.4). {A:^[12]}

The medium of the heliosphere is fundamentally distinct from that of the lower 100 km of the terrestrial atmosphere: the solar wind is primarily made up of hydrogen with a lesser amount of helium, is hot and therefore almost fully ionized, and is threaded by a magnetic field. The dynamics of the solar wind and the ways in which it interacts with planetary magnetospheres is

A:12

¹² Activity: With the values in Table 2.4, how long do the slow and fast solar-wind streams take to reach Earth? How many degrees does the Sun rotate between the moment these wind streams leave the Sun and the moment they arrive at Earth? How long for Neptune? Given that the wind flows out essentially radially, what is the apparent direction of the wind relative to the direction of the Sun as observed from the orbiting Earth (with an orbital velocity of about 30 km/s)?

Table 2.4. *Basic parameters of the fast and slow solar wind [near Earth; modified after Table H-I:9.1. Notes: (1) subscripts 'e', 'p', and 'i' are used to denote electrons, protons, and other ions, respectively; (2) v_A denotes the Alfvén velocity; (3) 'FIP' stands for 'first ionization potential'; 'low-FIP' is a group of elements with first-ionization potentials below 10 eV.]*

Property (1 AU)	Slow wind	Fast wind
Speed	430 ± 100 km/s	700 – 900 km/s
Density	$\simeq 10$ cm ⁻³	$\simeq 3$ cm ⁻³
Flux	$(3.5 \pm 2.5) \times 10^8$ cm ⁻² s ⁻¹	$(2 \pm 0.5) \times 10^8$ cm ⁻² s ⁻¹
Magnetic field	60 ± 30 μ G	60 ± 30 μ G
Temperatures	$T_p = (4 \pm 2) \times 10^4$ K $T_e = (1.3 \pm 0.5) \times 10^5$ K $> T_p$	$T_p = (2.4 \pm 0.6) \times 10^5$ K $T_e = (1 \pm 0.2) \times 10^5$ K $< T_p$
Anisotropies	T_p isotropic	$T_{p\perp} > T_{p\parallel}$
Structure	filamentary, highly variable	uniform, slow changes
Composition	He/H $\simeq 1 - 30\%$ low-FIP enhanced	He/H $\simeq 5\%$ near-photospheric
Minor species	n_i/n_p variable $T_i \simeq T_p$ $v_i \simeq v_p$	n_i/n_p constant $T_i \simeq (m_i/m_p)T_p$ $v_i \simeq v_p + v_A$
Associated with	streamers, transiently open field	coronal holes

modulated by that magnetic field, but the basic stratification and flow of the solar wind can be understood by looking at only the second of these characteristics: because it is hot and ionized, the electrons in the solar wind are very efficient at conducting heat, and that is all it takes to understand how it can lead unavoidably to a fast wind that can escape solar gravity. It is not simply an 'evaporation' off the Sun; after all, even at some millions of degrees, [H-I:9] "the sound speed c_s — essentially the mean ion speed — is much smaller than the [escape speed v_{esc} which can be derived by equating a particle's kinetic energy with its gravitational potential energy at the surface:

$$v_{\text{esc}} = \sqrt{2GM/r}. \quad (2.4)$$

For the solar corona, the sound and escape speeds are]

$$c_s \approx \sqrt{kT/m} \approx 100 \text{ km/s} \ll v_{\text{esc}} = \sqrt{2GM_\odot/R_\odot} = 618 \text{ km/s}, \quad (2.5)$$

where k is Boltzmann's constant, T the coronal temperature, m the mean particle mass, G the universal gravitational constant, and M_\odot and R_\odot the solar mass and radius, respectively.

Mass and momentum balance radially away from the Sun [in an assumed

uniform, strictly radial flow] at heliocentric distance r can be written

$$\frac{d}{dr}(\rho v 4\pi r^2) = 0 \quad (2.6)$$

$$\rho v \frac{dv}{dr} = -\frac{dp}{dr} - \rho \frac{GM_\odot}{r^2}, \quad (2.7)$$

with ρ the mass density, v the flow speed. Then $p = 2nkT$ is the gas pressure in an electron–proton plasma with n representing the electron or proton number density, and $\rho = mn$ where m is the mean particle mass which is given by $m \approx m_p/2$ for an electron–proton plasma. {A:[13]}

A:13

[The] consequence of the thermal conduction in a million degree corona is to extend the corona; *i.e.*, the temperature falls off slowly with distance from the Sun. Thus, in a hypothetical *static* atmosphere, we find a pressure at infinity given by

$$\frac{dp}{dr} = -nm \frac{GM_\odot}{r^2}, \quad (2.8)$$

$$p(r) = p_0 \exp \left[-\frac{mGM_\odot}{2k} \int_{R_\odot}^r \frac{dr}{r^2 T(r)} \right]. \quad (2.9)$$

Thus, if the temperature falls less rapidly than $1/r$, we find that $\lim_{r \rightarrow \infty} p(r) > 0$, we expect a non-vanishing pressure at infinity when the corona is extended. In particular, we find that for reasonable temperatures and densities n_0, T_0 at the 'coronal base' this pressure is much larger than any conceivable interstellar pressure.

[The observed slow decrease of temperature with distance from the Sun, caused by the efficient thermal conduction that is mostly carried by electrons, implies that the solar wind must expand supersonically into interstellar space. For a spherically symmetric, single-fluid, isothermal outflow,] the equations of mass and momentum conservation (Eqs. 2.6, 2.7) can be rewritten to give {A:[14]}

A:14

¹³ Activity: The momentum balance in Eq. (2.7) describes a radially-flowing wind over a non-rotating Sun. In reality, the Sun is rotating, and the magnetic field reaching into the heliosphere enforces the wind to co-rotate with the Sun, out to a distance where it becomes too weak to enforce such co-rotation. Show that for a sufficiently slowly rotating Sun, ignoring the centrifugal force is warranted. At what rotation period of a star like the Sun does the centrifugal force at, say, $2R_\odot$ counteract gravity by more than 10%? The centrifugal force in the wind would have been important for the very young Sun, see Sec. 10.2.1. Moreover, in the early phases of star-disk systems, centrifugal forces may be important in driving a cold wind; see Sect. 7.2.4.

¹⁴ Activity: What powers the solar wind in the basic model discussed here? To see the answer, rewrite Eq. (2.10) to an energy equation with the terms for the kinetic and potential energy in the Sun's gravitational field, plus a term that reflects the work done by the expanding gas both geometrically and by acceleration; the energy for that expansion in the isothermal approximation is provided by the thermal conduction by the electron population. The real-world solar wind is not isothermal, certainly not far from the Sun (compare the coronal temperatures in Table 2.3 with near-Earth wind properties in Table 2.4), and moreover is provided some additional power (in the form of heating and pressure) by waves and turbulence.

$$\frac{1}{v} \frac{dv}{dr} \left\{ v^2 - \frac{2kT}{m_p} \right\} = \left\{ \frac{4kT}{m_p r} - \frac{GM_\odot}{r^2} \right\} \quad (2.10)$$

[The solar wind starts slow, but is supersonic further out in the heliosphere; such a] transonic wind passes through a critical point at

$$r_c = \frac{m_p GM_\odot}{4kT} \quad \text{where} \quad v_c = \sqrt{\frac{2kT}{m_p}} \quad (2.11)$$

[(note the dependence on stellar mass). Formally, the equations allow such a flow] to match *any* pressure as $r \rightarrow \infty$ [although in reality the reach of the flow is limited by the existence of an interstellar medium (Sec. 5.5.8)]. {A:[15]} A:15

Let us examine this transonic wind solution in somewhat greater detail. If we integrate the force balance, Eq. (2.7), from the coronal base to the critical point r_c we find a density ρ_c at the critical point given by

$$\rho_c = \rho(r_c) = \rho_0 \exp \left\{ -\frac{m_p GM_\odot}{2kTR_\odot} + \frac{3}{2} \right\}. \quad (2.12)$$

Note that this density is almost exactly the same as if there had been *no* solar wind flow: *The subsonic corona in the solar wind is essentially stratified as a static atmosphere.*

We can also find the resultant mass flux for the wind by examining the density and the velocity at the critical point:

$$(nv)_r = n_c v_c \frac{r_c^2}{r^2} \propto \rho_0 T^{-3/2} \exp \left[-\frac{C}{T} \right] \quad (2.13)$$

where ρ_0 is the density at the coronal base [and C a constant]. *The mass flux is proportional to the density at the coronal base and depends exponentially on the coronal temperature.*” The actual solar wind is not only driven by thermal conduction from the coronal environment (which supplies energy for the work of driving the wind against gravity), but also by magnetic waves, known as Alfvén waves, whose fluctuations act as an additional pressure term, and whose dissipation aids in heating far above the solar surface, all of which is particularly important for the fast wind streams; more on that in Sect. H-I:9.5. Another note on more detail is found in Sect. H-I:9.6, which begins to explain why for a more realistic solar wind description that also allows for helium,

¹⁵ Activity: In principle, Eq. (2.10) allows for an inflow: where v is negative, dv/dr needs to be of opposite sign also. This inflow, accelerating from infinity towards the star, is known as Bondi accretion. However, such inflow is unlikely to occur as an isothermal flow from infinity because the interstellar medium is typically cold, with low ionization and thus low heat conductivity by electrons. Consequently, compression would raise the temperature of the inflow. Moreover, be aware that the quasi-hydrostatic approximation fails for the inner regions of such an infall, starting already well outside the critical point! Note that there is another class of solutions, namely a ‘solar breeze’: starting at low speed and never becoming transonic. Where does a ‘solar breeze’ reach its maximum velocity?

the exponential dependence of the solar mass loss on temperature is much weakened into a power-law dependence of temperature.

Note that it is not only the efficient thermal conduction per se that leads to a significant solar wind, but also the high temperature and low particle mass, and that that is the reason for the contrast with Earth's atmosphere. In Eq. (2.2) gravity is approximated by a constant, leading to a formal solution for the pressure stratification of the terrestrial atmosphere that tends to zero exponentially even for an isothermal atmosphere; this is not a bad approximation for an atmosphere in which the pressure scale height (at most some 50 km) is well below 1% of the planet's radius, so gravity changes little even over many scale heights above the surface. But in the hot corona, the pressure scale height for the hydrogen-dominated gas at ≈ 2 MK is about $0.15R_\odot$, so gravity diminishes noticeably in the first few pressure scale heights, hence its distance dependence needs to be reflected in Eq. (2.7). The relatively weaker gravity (and the correspondingly reduced escape energy) at large heights leads to a transonic wind at coronal temperatures.

On a side note (to which we return in Ch. 11), the same equation Eq. (2.7) also informs us about an accelerating inflow (for which $vdv/dr > 0$ as both $v < 0$ and $dv/dr < 0$) enabling the formation of stars and planetary systems: gravity can win out over a pressure difference on very large scales in the Galaxy on which stars form, because now gravity in fact is built up by the infalling matter itself so that M_\odot needs to be replaced to read

$$\rho v \frac{dv}{dr} = -\frac{dp}{dr} - \rho \frac{G}{r^2} \int_0^r \rho 4\pi r^2 dr. \quad (2.14)$$

[H-III:3] “To make a star of a given mass M from a gas with temperature T , gravity must overcome the pressure support. [One way to estimate the required properties of a cloud involved in the initiation of star formation is to look at Eq. (2.14) and see when conditions cannot remain in a stationary balance, *i.e.*, when $v = 0$ cannot be maintained. That occurs when] the radius R of the protostellar cloud exceeds]

$$R \gtrsim \frac{GM}{c_s^2} = \frac{GM\mu m_p}{kT}, \quad (2.15)$$

where c_s is the sound speed and m_p is the mass of the hydrogen atom. Taking a mean molecular weight $\mu = 2.3$, appropriate for molecular hydrogen plus helium, and a typical cold molecular cloud temperature of $T = 10$ K, Eq. (2.15) implies that a solar mass star must collapse from a cloud of radius $R \sim 2 \times 10^4$ astronomical units (*i.e.*, Sun-Earth distances; shorthand] AU).”

You will see the logic used in these examples applied throughout this book, and indeed astrophysics in general: approximations in functional forms, simplifi-

cations about geometries, and order of magnitude estimates are used throughout to aid in the basic understanding what is going on. With these tools, analytical and – far more commonly – numerical solutions become interpretable in terms of the basic, common processes. How much can be simplified to show the basics, however, depends on the environment: heliophysics, as is physics in general, is about simplifying as much as is allowed, but no more.

2.3 Photons, collisions, ionization, and differentiation

In our everyday lives we can get away with taking it for granted that the atmosphere around us is the same no matter where we are. Moreover, we may take it to be true that this atmosphere is a mixture of mostly N_2 and O_2 . And that this atmosphere is a very poor electrical conductor and that its winds are unaffected by the planetary magnetic field. As it turns out, none of these properties that we take for granted apply outside of the domain where we live: the chemical mixture depends on height in planetary atmospheres and is affected by the variable spectral irradiance from the Sun's outer atmosphere, ions and thus electrical conductivity are important in most of the local cosmos, and magnetic fields influence flows and vice versa almost everywhere in space. In this section, we focus on the processes that make the atmospheric composition dependent on location, primarily altitude. In the next section, Section 2.4, we start looking at the role of ions in electrical conductivity and flows, although the role of magnetic fields in that is the focus of Ch. 3.

The scale height for different atmospheric constituents depends on the molecular or atomic mass, and is thus in principle different for different chemicals. But as long as the mixing by winds and (turbulent) convection is fast enough compared to the time scale by which the chemical separation can occur by diffusive settling, the atmospheric composition will remain uniform, and all major species will share the same scale height. When collisions become relatively infrequent above the homopause (at about 100 km for Earth), and diffusive settling exceeds mixing by flows, separation of chemicals by molecular mass occurs; see Ch. 13. The rate of separation depends on the diffusion coefficients, which themselves depend on chemical species and density, and on the chemical reactions that couple species (and, in the ionosphere, also through ion-neutral interactions), relative to turbulent mixing efficiency; see the discussion in Ch. H-IV:9.

Still higher in the atmosphere, where collision frequencies become so low that the mean free path approaches or exceeds the formal pressure scale height, the description of the medium as an ideal gas fails. That environment, where particles essentially move ballistically over long distances subject only to gravity (still disregarding any effects of electric and magnetic fields), is known as the

exosphere. The exospheric base height can be estimated by looking at collision frequencies.

The characteristic frequency at which a particle in a non-magnetized plasma or a non-ionized gas of identical particles, all characterized by a temperature T and at particle density n , collides with other such particles is given by

$$\nu = \sigma_{cc} v_{\text{rel}} n = \sigma_{cc} \left(\frac{kT}{m} \right)^{1/2} n, \quad (2.16)$$

where σ_{cc} is the mutual collision cross section and v_{rel} is the velocity of one particle relative to another. In computing the mean free path, the velocity cancels out, leaving only the density as a variable:

$$\lambda_{\text{mfp}} = \frac{v_{\text{rel}}}{\nu} = \frac{1}{\sigma_{cc} n}. \quad (2.17)$$

By way of example, let us look at neutral atoms with a collisional cross section of order, say, $3 \times 10^{-16} \text{ cm}^2$ (as for hydrogen atoms). For these, a density of $3 \times 10^8 \text{ cm}^{-3}$ (reached at roughly 500 km in Earth's atmosphere, depending on solar activity) would correspond to $\lambda_{\text{mfp}} \approx 100 \text{ km}$. This order-of-magnitude estimate shows that this density in the Earth's atmosphere roughly forms the point at which a vertically moving atom could jump over a scale height, or essentially through the bulk of overlying matter, so where the assumption that we can work with the medium as a gas of electrically neutral particles fails; this is about the point where the Earth's atmosphere transitions into an exosphere where neutral atoms move essentially ballistically.

On the Sun, in contrast, the neutral hydrogen population could still be described by hydrodynamics at that density because of the much larger scales involved, if matter were largely neutral there; however, that density is reached only in the corona where high temperatures cause hydrogen and helium to be fully ionized (see Table 2.3), and collisions occur via long-range electromagnetic forces between charged particles (see Table 3.4 for mean-free path estimates in an ionized medium, which, with Eq. (2.17), shows the larger effective collision cross section for Coulomb collisions). Lower down in the solar atmosphere where neutrals do dominate, the mean free path lengths are significantly smaller: the plasma throughout the Sun up to the inner corona behaves like a gas in which (often turbulent) flows counter gravitational separation. There are fractionation effects deep inside the Sun where mixing by flows is negligible on solar evolutionary time scales. Chemical differentiation is also seen in the atmosphere in the minority species, specifically determined by the energy required for first ionization of the atom (see Fig. H-I:9.2); this differentiation, not by diffusive settling but likely related to MHD waves and by EUV and

X-ray irradiation of the chromosphere from the higher atmosphere, is still inadequately understood and not further discussed here.

Below the Earth's exosphere and above the mesosphere, in a domain ranging from roughly 110 km to around 500 km in altitude, *i.e.*, throughout much of the bulk of the thermosphere, lies a domain where collisions are frequent enough that the gas approximation is largely valid but not frequent enough to maintain uniform mixing of the chemicals that make up the terrestrial atmosphere up to that height: the atmosphere up to heights of about 110 km [H-I:12.3] "is known as the homosphere and is constantly being mixed by turbulent wave eddies. It is only at altitudes above about 110 km that turbulent mixing gives way to molecular mixing processes, where each species begins to be distributed vertically under its own pressure scale height or hydrostatic balance, see Eq. (2.2). A heavy species, such as carbon dioxide, will decrease in concentration with height more rapidly than a lighter species, such as atomic oxygen (see Fig. 2.5). Each species, i , will have its own characteristic scale height H_{pi} , where $H_{pi} = kT/m_i g$, which is the vertical distance a species will decrease in partial pressure and number density by a [factor of e]. The upper atmosphere differs from the lower atmosphere in this respect such that the mean mass of the fluid will change with altitude, as well as other gas parameters such as the specific heat, c_p . [...]

A:16

The vertical distribution of species also has a global seasonal/latitudinal structure from large scale [...] inter-hemispheric circulation from summer to winter. Closure of this circulation drives an upwelling of material across

¹⁶ Activity: **'What if' scenarios:** If you would like to think well 'outside the box' of things explicitly discussed in this book and in the Heliophysics volumes then consider this in the following chapters as you go along: what are things like when settings change? You could think of exoplanets with different host stars, orbits, and atmospheres, but there will be limited guidance by what we actually know from the literature. (1) For an example not too far from home, you could consider Titan, the only moon (natural satellite) in the solar system with a substantial atmosphere that is mostly N₂ (some 97%) and CH₄ (much of the remainder). *Activity:* Find Titan's equivalent values for the quantities listed in Table 2.1. *Further reading:* You can find publications on the (photo-)chemistry of its atmosphere leading to an ionosphere rich in HCNH⁺ and C₂H₅⁺. The chemical network in the high atmosphere leads to heavy organic molecules and aerosols that are deposited onto Titan's frozen surface and into its hydrocarbon lakes. Titan orbits within Saturn's magnetosphere, generally shielded from the direct impacts of the solar wind. However, the solar wind causes Saturn's magnetosphere to be highly asymmetric, and thus the environment through which Titan orbits is highly dependent on its orbital phase. Cosmic rays and energetic particles from Saturn's magnetosphere penetrate deep into Titan's atmosphere causing ionization and influencing chemical pathways. Titan has no intrinsic magnetic field (*i.e.*, no functioning magnetic dynamo) but an induced magnetosphere that changes as the moon orbits the rotating giant planet Saturn. There may be subsurface areas of liquid water, a water-ammonia mixture, or different mixtures in different locations and at different depths. Life might exist under these circumstances, and the traditional definition of 'habitability' as involving liquid surface water may need rethinking as we learn more. (2) For something far from home, consider the compact 7(?) -planet system of TRAPPIST-1 (see Fig. 15.1) on which much is being written: execute an ADS search for refereed papers with 'TRAPPIST-1' in the title. Task: let your imagination wander, read up on some of these things, and see how processes discussed in this volume apply to environments that are very different from those for Earth even though they are in some sense 'terrestrial'. Keep a running list of your thoughts as you read along for use later on!

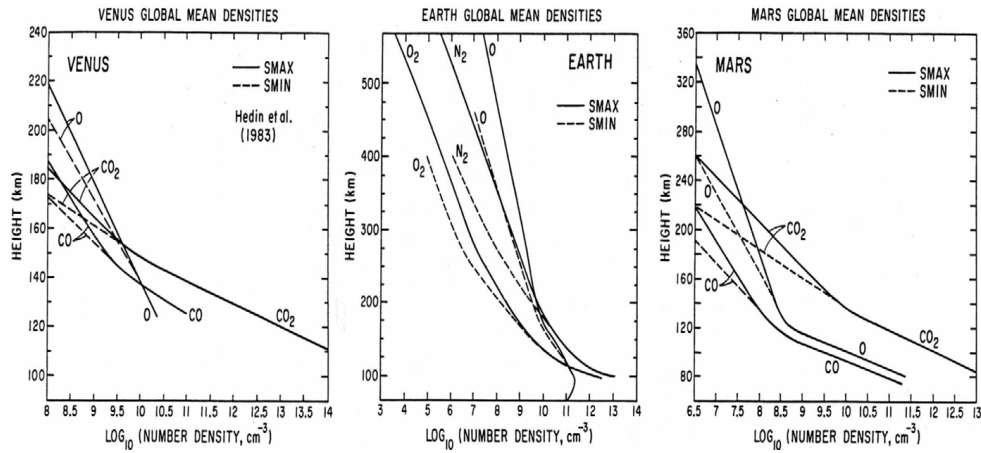


Fig. 2.5. Comparison of the global mean vertical profiles of the major species in the neutral upper atmospheres of a) Venus, b) Earth, and c) Mars for low and high solar activity. SMIN and SMAX indicate solar minimum and maximum conditions. Note that the turbopause heights (where turbulent mixing and diffusive separation are comparable) are 135, 110, and 125 km for Venus, Earth, and Mars, respectively. [Note: the International Space Station orbits at an altitude of ~ 400 km. Figs. H-I:12.2, H-IV:9.1; source: Bougher and Roble (1991).]

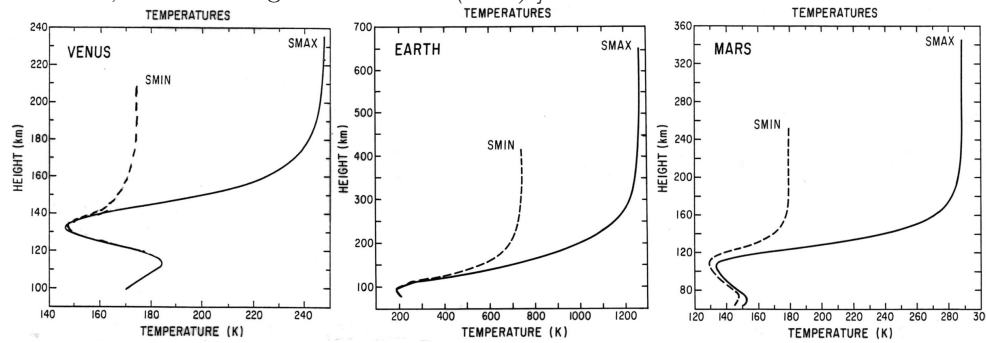


Fig. 2.6. Three planet global mean temperature profiles for solar minimum (SMIN) and maximum (SMAX) conditions. [Note the differences in horizontal and vertical scales. Fig. H-IV:9.3; source: Bougher and Roble (1991).]

surfaces of constant pressure in the summer hemisphere and a downwelling in the winter hemisphere. The upwelling causes the heavier molecular rich gas, which had diffusively separated at lower altitudes, to be transported upwards to increase the mean molecular mass in summer. In winter the downwelling reduces the mean mass.”

The seasonal changes in insolation and the resulting circulations subject to Coriolis forces on the rotating planet are modulated by the effects of space weather. These effects include the X-ray and (E)UV part of the solar

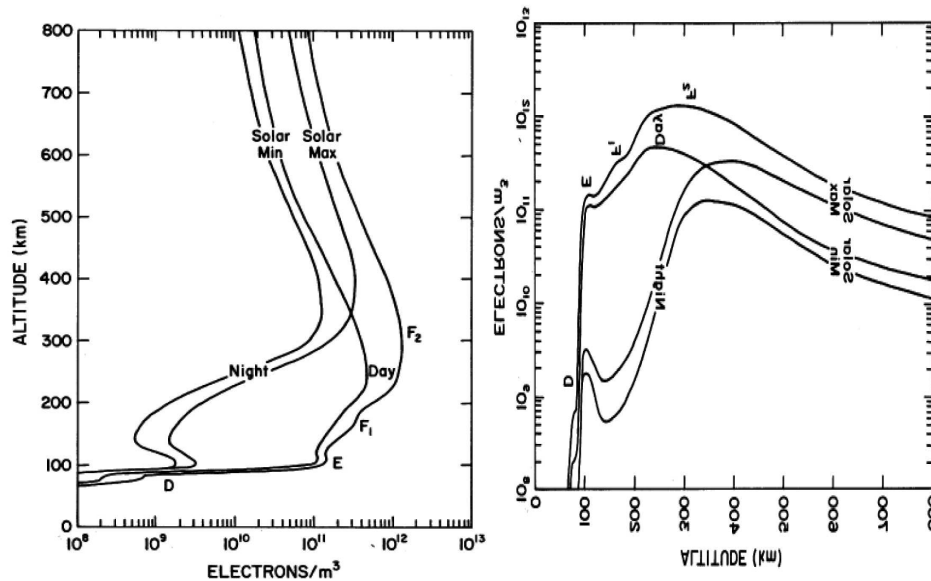


Fig. 2.7. [left:] Overview of the altitude distribution of Earth's ionosphere for daytime and nighttime conditions, at high and low solar activity. [Fig. H-III:13.1] [right: One of multiple different conventions between planetary scientists and astrophysicists is that the height coordinate is usually displayed vertically for planetary scientists and horizontally for stellar scientists. This flipped and rotated version of the figure conveys the difference in appearance.]

spectral irradiation, dissipation of electrical currents, and energetic particles precipitating from the magnetosphere. All of these (and others discussed in Chs. 5, 6 and 14) lead to heating, ionization, and dissociation of the high atmosphere. [H-III:13.1] “Early investigation of the terrestrial ionosphere through its effect on radio waves resulted in description by means of layers, principally the *D*, *E*, and *F* layers, the latter subdivided into *F*₁ and *F*₂. This terminology continues to influence our current concept of the nature of energy deposition in atmospheres, although the misleading term ‘layer’ has given way to ‘region’. The term ‘layer’ arose from the observation of systematic variation in the height at which the critical frequency of reflection occurs in ionospheric radio sounding; this method cannot detect ionization above the peak of a region, which explains the appearance of layers. Radar and spacecraft measurements now give a more complete picture of peaks and valleys and reveal the complex morphology of the ionosphere. [...] An overview of the altitude dependence and variability of Earth's ionosphere is given in Figure 2.7, showing the diurnal and solar-cycle changes and the locations of the named regions.”

[H-III:13.1] “An additional historical artifact in terminology is the word ionosphere itself. Because the atmospheric ionization was discovered before the

neutral thermosphere in which it is contained, anything above the stratosphere is often referred to as the ionosphere, resulting in a common misconception that this region of the atmosphere is mostly ionized. In fact, it is mostly neutral, ranging from less than a part in a million ionized during the day at 100 km altitude to about 1% ionized at the exobase (~ 600 km, depending on solar activity; compare Fig. 2.10). Even at 1000 km, there is only of the order of 10% ionization. At several thousand km, where ions (mostly protons) finally become dominant, the region is defined as the plasmasphere. [. . . In the bulk of the terrestrial ionosphere] O^+ is the most important ion, particularly in the extensive F_2 region above ~ 200 km. The F_1 region from ~ 150 to ~ 200 km appears as a mere plateau in the profile, but is distinguished by a transition to molecular ions, particularly NO^+ . The low levels of N_2^+ , given the dominance of N_2 at these altitudes, is noteworthy. {A:^[17]} The E region from ~ 100 to ~ 150 km exhibits a small peak, dominated by O_2^+ and NO^+ .”

A:17

[H-I:12.4.1] “Much of the external sources of heating, ionization, and dissociation of a planetary atmosphere comes from the absorption of photons or particles impinging on the neutral atmosphere. The physics defining the altitude profile of the three processes is the same. For example, the rate of ionization, q [($cm^{-3}s^{-1}$)], by solar radiation intensity, $I(h)$ [($erg\ cm^{-2}\ s^{-1}$)], at some height in the atmosphere of number density, $n(h)$, can be expressed as a product of four terms:

$$q = \sigma_a I(h) n(h) \eta_i, \quad (2.18)$$

where σ_a [(cm^2) is the atomic] absorption cross section [for a wavelength interval matching that of I ,] and η_i [(erg^{-1})] is the ionizing efficiency; η_i could equally be the heating or dissociation efficiency. The intensity of the radiation gradually decreases along the path through the atmosphere starting from an initial intensity of $I(h = \infty)$. The altitude deposition profile depends on the absorption coefficient and on the atmospheric number density, which varies exponentially with height. Clearly the product of the intensity of the radiation, I , that decreases as the source penetrates the atmosphere, and on the atmospheric number density, $n(h)$, that increases with increasing depth into the atmosphere, must reach a maximum at some altitude or, more correctly, at some pressure level [(except, of course, for visible wavelengths for which the atmosphere is largely transparent, in which case the surface absorption and reflection need to be taken into account)]. The level of penetration is referred

¹⁷ Activity: The fact that concentrations of atomic nitrogen are not shown in Fig. 2.5 should make you wonder given that molecular nitrogen is the most common species in the troposphere. Why is atomic nitrogen rare in the upper atmosphere? Hint: compare the molecular binding energies of nitrogen, oxygen, and water.

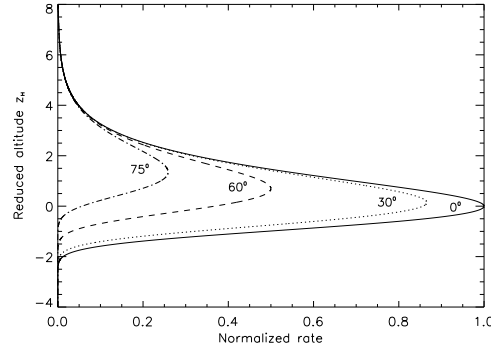


Fig. 2.8. The vertical profile of the classical Chapman profile appropriate for heating, ionization, or dissociation in a stratified hydrostatic atmosphere irradiated from above, relative to a reference height, shown for different slant angles χ . [Fig. H-I:12.4]

to as the optical depth, τ , which is expressed mathematically as

$$\tau = \sigma_a n(h) \frac{H_p(h)}{\cos(\chi)}, \quad (2.19)$$

where the product of the number density $n(h)$ at height h with the scale height $H_p(h)$ at that level represents the integrated content of a column of gas above that point, and χ is the angle from the zenith at which the radiation penetrates a planar atmosphere. [The above expression is valid as long as the curvature of the atmosphere can be neglected, so for angles $\chi \lesssim 75^\circ$.] {A:[18]} A:18

The profile of the rate of heating, ionization, or dissociation from these processes takes the form of the classical Chapman profile, as depicted in Fig. 2.8, and is given mathematically by

$$q(h) = I_\infty \exp \left[-\sigma_a n(h) \frac{H_p(h)}{\cos(\chi)} \right] \eta_i \sigma_a n(h). \quad (2.20)$$

{A:[19]} The peak of the profile is at unit optical depth, which depends on A:19 the mass of atmosphere above traversed by the energetic photon or particle. This corresponds to a fixed pressure level for a given angle of incidence. The

¹⁸ Activity: Optical depth is an integral over absorption along a line of sight, and thus as useful for incoming as for outgoing radiation. Explain why the layers contributing most to the light from the solar photosphere are geometrically higher as you look away from disk center. What can you infer about the stratification of the solar atmosphere from the fact that the Sun (emitting close to black-body radiation over much of the optical spectrum) is brightest near disk center, darkening towards the limb? What follows directly from the fact that, on average, the solar corona seen in X-rays or extreme ultraviolet (EUV) has essentially double the intensity just outside the solar limb compared to that just inside the limb when there are no active regions along these lines of sight?

¹⁹ Activity: You can think of the optical depth as the mean number of absorbers within the cross-section along a photon's path from infinity to height h . The probability of suffering zero absorptions, and thus making it to h , is $\exp(-\tau)$. The intensity at h is then an integral from infinity over the expected number of absorptions along the way. Combine that with Eq. (2.18) to derive Eqs. (2.19) and (2.20).

depth of penetration into the atmosphere of a photon or particle in pressure coordinates therefore does not change with the gas temperature or the degree of thermal expansion. Even with the changing heating over the solar cycle or during a [magnetospheric] storm that might cause a thermal expansion of the atmospheric gas, that same radiation will still penetrate and produce heating or ionization at the same pressure level. The altitude associated with that pressure and the local number density would, of course, be different since they depend explicitly on gas temperature.” {A:[20]}

A:20

2.4 On collisions and currents, and on neutrals and pickup ions

The terrestrial upper atmosphere is coupled to the Earth’s magnetic field through the ionized component of the atmosphere (referred to as the ionosphere) that is in turn collisionally coupled to the neutral molecular and atomic medium within which it is embedded. The dynamics of these couplings in the overall system of solar wind, magnetosphere, and ionosphere are discussed mostly in later Chs. 5, 6, and 13. Here, we look at the consequence of the ionized medium threaded by a dynamic magnetic field and embedded in moving neutral gas: electrical currents. In the terrestrial atmosphere, the effects depend sensitively on the magnetic latitude because of the orientation of the magnetic field: at high latitudes, where the field is predominantly vertical, the connection with the magnetosphere dominates and the dissipated power can lead to substantial heating. At mid and low latitudes, where the field is mostly horizontal, internal processes dominate that provide less dissipative power than at higher magnetic latitudes, but that do contribute to transport of plasma.

A moving electrical charge subject to a magnetic field experiences a Lorentz force perpendicular to its velocity and to the magnetic field, in a direction that depends on the sign of the charge. Also allowing for an electrical field to be present, the total force equals:

$$\mathbf{F}_L = m \frac{d\mathbf{v}}{dt} = q\mathbf{E} + \frac{q}{c} \mathbf{v} \times \mathbf{B}. \quad (2.21)$$

In case $\mathbf{E} = \mathbf{0}$ and in the absence of collisions, electrons and ions thus would spiral about the magnetic field line in opposite directions [(much more on that in Sect. 8.1)]. Their gyration radii and frequencies are very different because of their difference in mass and thermal velocity (see Table 2.5). Where the

²⁰ Activity: A similar expression to Eq. (2.20) derived for photons holds for energetic particles (from, say, 1 keV/nucleon to 1 GeV/nucleon) losing their energy when propagating into a relatively dense medium (from the Earth’s magnetosphere into its atmosphere, from the solar corona into its chromosphere or photosphere, or from interplanetary space into a spacecraft hull). Such energetic particles can penetrate a medium up to a column density of a few grams per cm². Very roughly, estimate how far down that is into Earth’s atmosphere, into the Martian atmosphere, into the solar lower atmosphere, and into an aluminum shell of a spacecraft. (See, *e.g.*, Sects. H-II:1.6, H-II:13.4, and H-II:14.4.)

Table 2.5. Selected plasma quantities, mostly for thermal motions in fully-ionized plasmas, generally from the NRL Plasma Formulary^a.

Name/Symbol	Value ^b	Description
Frequencies and rates		
f_{ge} (Hz)	$2.8 \cdot 10^6 B$	electron gyrofrequency
f_{gi} (Hz)	$1.5 \cdot 10^3 B \frac{Z}{\mu}$	ion gyrofrequency
Thermal collision frequencies for fully ionized plasmas:		
ν_{ee} (s ⁻¹)	$3.6 \frac{n_e}{T_e^{3/2}} \ln(\Lambda)$	electron-electron collision rate
ν_{ii} (s ⁻¹)	$0.06 \frac{n_i}{T_i^{3/2}} \frac{Z^4}{\mu^{1/2}} \ln(\Lambda)$	ion-ion collision rate
ν_{ei} (s ⁻¹)	$\approx 0.5 \nu_{ee}$	electron-ion collision rate
Thermal length scales		
r_{ge} (cm)	$0.022 \frac{T_e^{1/2}}{B}$	electron gyroradius
r_{gi} (cm)	$0.95 \frac{T_i^{1/2}}{B} \frac{\mu^{1/2}}{Z}$	ion gyroradius
λ_D (cm)	$6.9 \frac{T_e^{1/2}}{n^{1/2}}$	Debye length
Thermal velocities		
v_{Te} (km/s)	$3.9 T_e^{1/2}$	electron thermal velocity
v_{Ti} (km/s)	$0.091 T_i^{1/2} \frac{1}{\mu^{1/2}}$	ion thermal velocity
c_s (km/s)	$0.091 T_e^{1/2} \frac{\gamma^{1/2} Z^{1/2}}{\mu^{1/2}}$	ion sound velocity
v_A (km/s)	$2.2 \cdot 10^6 B \frac{1}{\mu^{1/2} n_i^{1/2}}$	Alfvén velocity
Dimensionless numbers		
electron Hall coeff.	$\frac{f_e}{\nu_{ee}} \approx 8 \cdot 10^5 \frac{B T_e^{3/2}}{n_e}$	electron gyro- to collision frequency ^c
ion Hall coefficient	$\frac{f_{ii}}{\nu_i} \approx 2.2 \cdot 10^4 \frac{B T_i^{3/2}}{n_i Z^3 \mu^{1/2}}$	ion gyro- to collision frequency ^c
plasma β	$3.5 \cdot 10^{-15} \frac{n T}{B^2}$	thermal to magnetic energy

^a <https://www.nrl.navy.mil/ppd/content/nrl-plasma-formulary>; ^b in cgs-Gaussian units. Symbols: B magnetic field strength (Gauss); T temperature in Kelvin, n density in cm⁻³; γ the adiabatic index; $\ln(\Lambda)$ the Coulomb logarithm (typically in the range of 10 to 20, *cf.* Table 3.4); μ ion mass in units of the proton mass; $n_{e,i}$ electron or ion density; Z ion charge state; ^c for $\ln(\Lambda) \approx 10$.

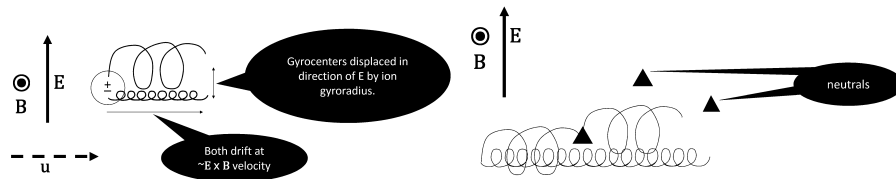


Fig. 2.9. Schematic of interactions of plasma with neutrals. Left: Initial motion of pickup ions and electrons. The gray circle represents a neutral composed of a positively charged ion and a negatively charged electron. The directions of plasma flow velocity, \mathbf{u} , of the magnetic field, \mathbf{B} , and of the electric field, \mathbf{E} , are indicated. In the image, following dissociation, the ion path starts upward and the electron path starts downward. Although initial motion is along \mathbf{E} for the ion, the Lorentz force causes the path to twist, resulting in motion around \mathbf{B} at the ion cyclotron period, leading to a net drift at a velocity of $\mathbf{E} \times \mathbf{B}/B^2$. The electron initially moves in the $-\mathbf{E}$ direction. Its motion also rotates around \mathbf{B} , but at the electron cyclotron frequency. The net effect is a transient current in the direction of \mathbf{E} . Right: Schematic of the effect of collisions with neutrals for a case with the collision frequency of order the ion cyclotron frequency. Triangles represent neutrals. The effect of collisions is to slow the motion in the $\mathbf{E} \times \mathbf{B}$ direction of the ions but not of the electrons [(which have other collision and gyrofrequencies)] and to displace the ions in the direction of \mathbf{E} . A net current arises, with one component along $-\mathbf{E} \times \mathbf{B}$ (a Hall current) and one component along \mathbf{E} (a Pedersen current). [Fig. H-IV:10.3]

gyration radii are well below the gradients in the magnetic field, these opposite circular motions do not lead to a net current in the absence of collisions. However, when field gradients are substantial within the gyration radii of the particles (most readily for the ions, in particular the more energetic ones) the particles drift perpendicular to the field in directions opposite for opposite charges, thus leading to a current; one important heliophysical setting in which this occurs is in the Earth's inner magnetosphere, where the gradient drift of primarily the energetic ions leads to the 'ring current' (see Sect. 8.1).

In the variety of settings in heliophysics, collisions may occur among the electron and ion populations (see Ch. 3 for that), or with neutral particles (the focus here). In ionospheres, the neutral particles are atoms and molecules of a body's atmosphere. In, say, the environments of comets, planetary rings, or in the outer heliospheric solar wind the neutral particles, in contrast, may be either dust particles, escaping atmospheric gas, or inflowing neutral interstellar atoms.

Let us start with a collision in which no charge-transfer occurs in a setting where the charged particle senses both a magnetic and electric field. In each such collision of an electron or ion with a neutral particle, the gyro-motion of

the electron or ion involved is modified. Because of the opposite charges of the electron and ion populations, they attempt to gyrate about the magnetic field in opposite directions as they are accelerated by the electric field; consequently, they exhibit a net drift perpendicular to the magnetic field with ions and electrons moving in the same direction and at the same rate. There is no net current (see the left-hand side of Fig. 2.9), but if there are collisions roughly at the same frequency as the gyrofrequency (different for particles of different masses), the situation changes fundamentally: collisions interrupt the gyromotion, and this results in a net separation of the charges. A graphic example, discussed in more detail below, is given in Fig. 2.9(right).

If collisions are very infrequent, or to be precise if the electrons or ions can gyrate about the magnetic field many times between collisions, the electrical conductivity across the magnetic field is very low. If collisions are very frequent for both ions and electrons, hardly any charge separation can occur between collisions, and the electrical conductivity perpendicular to the magnetic field is also very low. Peak perpendicular conductivity depends on the direction relative to the electric field and is reached depending on the ratios of collision and gyro-frequencies, as shown below.

Collisions between the populations of charged and neutral particles in the presence of a magnetic field while allowing for bulk flows is described through multiple equations. One of these captures the transfer of momentum that affects the force balance (touched upon towards the end of this section) almost entirely by looking at ions because they carry the bulk of the mass. Another accommodates electrical currents that arise from the differential behavior of the ions and electrons subject to the magnetic field. A third describes the energy transfer through the collisional effects formulated as Ohmic dissipation in the energy balance.

How collision frequencies influence currents in the ionosphere/thermosphere, where the neutral component is the most common, can be approximated as follows (collisions between charged particles are ignored here because collisions with the abundant neutrals are far more common in the bulk of the terrestrial ionosphere). [H-I:12.6] “If we take the magnetic field to be aligned with the z axis, then the generalized Ohm’s law [(the derivation of which is shown for a fully ionized plasma in Ch. 3)], $\mathbf{j} = \boldsymbol{\Sigma}_e \cdot \mathbf{E}_0$ (where \mathbf{E}_0 is the total electric field: $\mathbf{E}_0 = \mathbf{E} + \frac{1}{c} \mathbf{v} \times \mathbf{B}$), contains the conductivity tensor

$$\boldsymbol{\Sigma}_e = \begin{pmatrix} \sigma_P & \sigma_H & 0 \\ -\sigma_H & \sigma_P & 0 \\ 0 & 0 & \sigma_{\parallel} \end{pmatrix}, \quad (2.22)$$

where the Pedersen ($\perp \mathbf{B}, \parallel \mathbf{E}_{\perp}$), Hall ($\perp \mathbf{B}, \perp \mathbf{E}_{\perp}$), and parallel ($\parallel \mathbf{B}$)

conductivities are given [by:

$$\sigma_P = \frac{n_e e c}{B} \left(\left[\frac{M_e}{1 + M_e^2} \right] + \left(\frac{e}{q_i} \right) \left[\frac{M_i}{1 + M_i^2} \right] \right); \quad (2.23)$$

$$\sigma_H = \frac{n_e e c}{B} \left(\left[\frac{M_e^2}{1 + M_e^2} \right] - \left(\frac{e}{q_i} \right) \left[\frac{M_i^2}{1 + M_i^2} \right] \right); \quad (2.24)$$

$$\sigma_{\parallel} = \frac{n_e e c}{B} \left(M_e + \left(\frac{e}{q_i} \right) M_i \right) \quad (2.25)$$

(where the equations from Sect. H-I:12.6 were rewritten to the above by using the expression for $\omega_{e,i}$ below). For characteristic values of these conductivities in the terrestrial ionosphere, see Figure H-I:12.5. Here, $M_{e,i} = \omega_{e,i}/\nu_{e,i}$ are the electron and ion magnetizations, with $\omega_{e,i} = |q_{e,i}|B/m_{e,i}c$ the electron and ion (with charge q_i) gyro-frequencies around the field of strength B , $m_{e,i}$ are the electron and ion masses, ν_{en} and ν_{in} the electron-neutral and ion-neutral collision frequencies.] The effect of the collisions is to rotate the net current from the direction of \mathbf{E} at high altitudes towards the negative $\mathbf{E} \times \mathbf{B}$ direction at low altitudes. {A:^[21]} In the terrestrial ionosphere, the current and dissipation reach a peak at the altitude where the Pedersen and Hall conductivities are equal, around 125 km. For high-frequency currents, like those that may occur in the solar chromosphere, the dissipation may increase markedly (see Sect. H-I:12.8). Note that σ_P is generally dominated by the ion term.”

A:21

The collisional coupling between ions and neutrals causes momentum exchange (through the drag force that works to reduce the velocity difference between these two populations) and energy dissipation (in the form of Joule heating). [H-I:12.6] “The electrodynamic properties can be conveniently separated into a high [magnetic] latitude region, where the current flow in the ionosphere is connected to the magnetospheric current system, and a mid and low latitude region, where the majority of the current flow and polarization electric fields are controlled internally by the thermosphere-ionosphere conductivity and dynamics.” [H-I:12.6.1] “In the ionosphere, currents flowing perpendicular to the magnetic field are produced by electric fields and neutral winds. Although collisions between ions and the neutral gas are relatively infrequent [in Earth’s upper atmosphere] above ~ 160 km, they are sufficient to accelerate the neutral component, *i.e.*, the thermosphere, at high latitudes to

²¹ Activity: Work through Eqs. (2.22-2.25) to confirm that the effect of the collisions of charged particles in the ionosphere with the neutral thermospheric component is to rotate the net current from the direction of \mathbf{E} at high altitudes towards the negative $\mathbf{E} \times \mathbf{B}$ direction at low altitudes. As the expressions assume \mathbf{B} to be in the z direction, you could chose \mathbf{v} in the x direction to describe a horizontal velocity near the geomagnetic pole. In that same coordinate system, what is the direction of the current at about 125 km in the daytime terrestrial ionosphere where $\sigma_P \approx \sigma_H$ (see, *e.g.*, Fig. H-I:12.5).

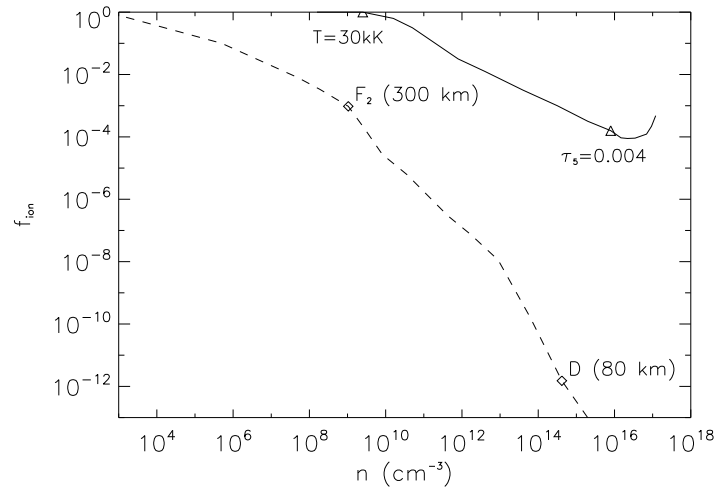


Fig. 2.10. Comparison of densities, n (cm^{-3}), and ionization fractions, f_{ion} , for a characteristic dayside ionosphere (dashed) and mean chromosphere (solid). The diamonds mark the mean values for the ionospheric D and F_2 regions, centered on about 80 km and 300 km, respectively. The triangles denote the base of the chromosphere (defined here as at a continuum optical depth of $\tau_5 = 0.004$) and the top of the chromosphere (where the temperature exceeds 30 000 K). [Fig. H-I:12.13]

many hundreds of m/s over periods of tens of minutes or more [to speeds well in excess of those associated with solar heating]. [...] At low altitude, ~ 100 km, the ions are forced to move with the neutral gas, whether stationary or moving. The large-scale wind system at this altitude is driven by the tidal and planetary waves propagating from the lower-atmospheric terrestrial weather system, and the mass of the atmosphere is such that ion drag has little or no impact on the neutral dynamics. The altitude range between 100 and 160 km altitude is the narrow altitude range that is responsible for most of the dissipation of electromagnetic energy from the magnetosphere. The neutral dynamics and conductivity in this boundary region between space and atmospheric plasma are critical.” [H-I:12.6.2] “At mid and low [magnetic] latitudes the electric fields [in Earth’s ionosphere] arise largely from internal dynamo processes driven by the conversion of neutral wind kinetic energy to electromagnetic energy, and are typically an order of magnitude smaller (a few mV/m) than high-latitude fields. The energy involved is also much smaller. The importance of the small electric fields at low latitudes is no longer the Joule heating and momentum dissipation, but rather their role in the redistribution of plasma.” Some of these effects are touched upon generically in Sect. 5.5.7, with a more comprehensive discussion for Earth’s ionosphere in Sect. H-I:12.6.

In much of the discussion of magnetized plasma in the Sun's interior and atmosphere in subsequent chapters, the Hall and Pedersen conductivities are often assumed to be negligible. A similar approximation is often seen in the study of the heliosphere and planetary magnetospheres. The Sun's chromosphere, however, is an environment with a strong neutral population and with collision frequencies not so high that Pedersen and Hall conductivities are effectively ignorably small. The chromosphere is located immediately above the photosphere (which itself has a thickness of roughly a single scale height of about 100 km), and extends over a height range of some 2,500 km, spanning roughly a dozen pressure scale heights in a highly dynamic setting that is strongly patterned by the magnetic field, before the transition region is reached in which the temperature rapidly rises to coronal values.

[H-I:12.8.3] “The Earth's ionosphere has a range of degrees of ionization, starting from the essentially neutral troposphere below, reaching an ionization fraction of about $10^{-4} - 10^{-3}$ around 200 km in height, and exceeding a few percent by 1 000 km. In the case of the chromosphere, the ionization fraction starts at about 10^{-4} around photospheric heights, drops through 10^{-5} through the classical 'temperature minimum' around 500 km in height, and then increases through a few percent around 1 500 km in height, continuing to near-complete ionization in the solar corona. Figure 2.10 compares the densities and ionization fractions for mean states characteristic of the ionosphere and chromosphere. Note that the neutral densities in the D–F₂ ionospheric region are comparable to those in the chromosphere, but the ion densities are at least 1 000 times lower at any given neutral density, resulting in a much weaker ion-neutral coupling in the ionosphere than in the chromosphere.

Let us look back at Eqs. (2.23)-(2.25) and assess their meaning for both chromosphere and ionosphere. In the limit of a weak magnetic field or a high collision frequency, the ion and electron magnetizations $M_{e,i} = \omega_{e,i}/\nu_{e,i} \rightarrow 0$, $\sigma_P \rightarrow \sigma_{\parallel}$, $\sigma_H \rightarrow 0$; hence, currents are more readily aligned with the *electric* field, as expected. As the collision frequencies with the neutral population decrease, the above expressions would have current and *magnetic* field aligned (as both $\sigma_{P,H} \rightarrow 0$) [...]

In the chromosphere of a solar [sunspot] region, $M_e(500 < h < 2000 \text{ km}) = \mathcal{O}(100)$, decreasing rapidly towards the photosphere to $M_e(h = 0) = \mathcal{O}(0.01)$ at the solar surface. Some studies find the proton magnetization to remain below unity throughout the chromosphere, up to the transition into the corona (these findings depend on the atmospheric domain, of course, and on the models used [...]). Consequently, the bulk of the active-region chromosphere has an anisotropic conductivity of at least a factor of 10 difference between the

field-aligned and transverse components. Conduction in the corona is almost exclusively field aligned (and thus essentially free of Lorentz forces), while photospheric conduction is nearly isotropic. [...]” {A:[22]} A:22

Now, let us look at different environments, and illustrate not only currents but also the effects of momentum transfer. [H-IV:10.3.2] “At comets and in the vicinity of moons, such as Io and Enceladus, that are significant sources of neutral gas, various processes that convert neutral atoms or molecules into ions are important to consider. Neutrals can be ionized by photons (photoionization) or by collisions with other particles, typically electrons (impact ionization). An additional process that affects the interaction region is charge exchange. In this process, a neutral gives up a charge to an ion. The original ion, now neutral, carries off its incident momentum while the original neutral becomes an ion at rest in the frame of the neutral gas.

The ions introduced into the plasma by ionization of neutrals modify the bulk properties of the plasma. Consider a situation in which the neutrals are at rest relative to [a location] towards which the plasma flows at (bulk) velocity \mathbf{u} . Photoionization and impact ionization add mass to the plasma whereas charge exchange between the ionized or neutral form of the same element does not change the mass density. All three processes slow the bulk flow because the new ions must be accelerated so that their average motion matches that of the bulk plasma and the process extracts momentum from the incident plasma. These processes also change the thermal energy of the plasma and may modify the plasma composition. The complex effects associated with pickup can significantly modify the interaction region surrounding a moon or a comet.

The relation between pickup and currents is shown schematically in Fig. 2.9a. The newly ionized ion senses the electric field of the flowing plasma and begins to move in the direction of this electric field. The electron that has separated from the ion is initially accelerated in the opposite direction. After one gyroperiod, the average separation of the gyrocenters of the two charges is close to one ion gyroradius

$$r_{\text{gi}} = m_{\text{ion}} v_{\text{ion}} c / qB \quad (2.26)$$

where m_{ion} is the ion mass, v_{ion} is its thermal velocity, and q is its charge. {A:[23]} The result of the separation of charges is to produce a transient A:23

²² Activity: Look up what defines a ‘sunspot’ and what an ‘active region’. A record of sunspot counts over many decades is shown in Fig. 4.5: what is the typical latitudinal range over which sunspots and sunspot groups occur? In Sect. 10.3.2 you will read about high-latitude and even polar starspots on rapidly-rotating, active stars, as the Sun would have been in its first few hundred million years.

²³ Activity: Use Table 2.5 to show that Eq. 2.26 yields r_{gi} for thermal motions. Then estimate energies of non-thermal particles so that their r_{gi} are comparable to the scale of the geomagnetic field (important for the terrestrial ring-current, which is a manifestation of particles drifting across the magnetic field because the heavy, energetic ones sense the gradient in the field strength; see Sect. 3.4) or perturbations in the heliospheric field (important for incoming cosmic rays, see Ch. 14). Compare

current density in the direction of the electric field. If the pickup is occurring at a rate \dot{n} , where \dot{n} is the number of ionizations per unit volume and time, then the pickup current [density] is

$$j_{\text{pickup}} = q\dot{n}r_{\text{gi}}. \quad (2.27)$$

Because pickup current flows across the background field, a cloud of pickup ions acts much like a solid conducting obstacle in the flow and imposes the same types of perturbations, *i.e.*, it slows and diverts the incident flow” in a way outlined in Ch. 5.

In this volume, we do not go into the behavior of dusty plasmas. The interested reader is referred to Ch. H-IV:11, which introduces the subject as follows: [H-IV:11.1] “The study of dusty plasmas bridges a number of traditionally separate subjects, for example, celestial mechanics, mechanics of granular materials, and plasma physics. Dust particles, typically micron and submicron sized solid objects, immersed in plasmas and UV radiation collect electrostatic charges and respond to electromagnetic forces in addition to all the other forces acting on uncharged grains. Simultaneously, dust can alter its plasma environment by acting as a possible sink and/or source of electrons and ions. Dust particles in plasmas are unusual charge carriers. They are many orders of magnitude heavier than any other plasma particles, and they can have many orders of magnitude larger (negative or positive) time-dependent charges. Dust particles can communicate non-electromagnetic effects, including gravity, neutral gas and plasma drag, and radiation pressure to the plasma electrons and ions. Their presence can influence the collective plasma behavior by altering the traditional plasma wave modes and by triggering new types of waves and instabilities. Dusty plasmas represent the most general form of space, laboratory, and industrial plasmas. Interplanetary space, comets, planetary rings, asteroids, the Moon, and aerosols in the atmosphere, are all examples where electrons, ions, and dust particles coexist.” {A:[24]}

A:24

with values in Table 3.4, compare these to mean-free path lengths there, and bear these results in mind going into Ch. 3.

²⁴ Activity: An intriguing property of dust is that, if the particles are small enough, radiation pressure is important in their momentum equation. Assuming neutral dust particles, estimate at what (density-dependent) size photon pressure from solar illumination exceeds solar gravity (note that this is independent of distance to the Sun for a completely transparent solar wind). There is a surprise here for dust of any size: the orbital motion of the dust causes photon absorption (and assumed isotropic re-radiation of that energy) to lead to a ‘brake’ on the orbital velocity, causing larger dust particles to spiral inward; look up ‘Poynting-Robertson drag’ to see how that works. From this, realize that dust needs to be continually replenished somehow in the Solar System, generally by impact collisions and by disintegrating comets.

2.5 Sources of plasma

There are many sources of plasma around the heliosphere: all it takes is some neutral medium subjected to sufficient energy to ionize particles. The bulk source medium can be the gas in the Sun's surface layers that is largely neutral, but dissipation of magnetic waves and the acceleration of particles in electric fields cause heating and ionization of the Sun's outer atmosphere. The larger planets have neutral atmospheres of which the top layers are ionized by solar radiation and by suprathermal particle precipitation (which can be of magnetospheric or solar origin). Moons may be large enough to have their own atmosphere (as is the case for Titan at Saturn), and those without significant atmospheres may still have some matter around their surfaces because these are subjected to sputtering by the solar wind or, for moons within planetary magnetospheres, by magnetospheric particles, or matter may be supplied by geysers (as on Enceladus at Saturn) or volcanoes (as on Io at Jupiter) that contribute molecules (including SO_2 , SO , S_2 , H_2S , ...) as well as atoms. Comets have a coma of gas that sublimates off the nucleus, along with dust. And dusty material is around in the rings of all the giant planets. Whereas the magnetized and ionized components of the interstellar medium cannot penetrate each other (as discussed in Chs. 3 and 5), neutral interstellar-medium particles can make it deep into the heliosphere, following free-fall trajectories in the collisionless environment until they are subjected to a charge-exchange collision with solar-wind ions.

3

MHD, field lines, and reconnection

3.1 Introduction

[H-II:1] “Absent the magnetic field, neither solar activity nor magnetic storms – the solar and terrestrial sources of [variable conditions referred to as space weather ^[vi] – would exist. . . .] Although in principle fossil magnetic fields could have remained from the creation of the Solar System, this appears not to be the case. Witness the 22-year magnetic cycle of the Sun and the reversals of the Earth’s magnetic field. On shorter time scales, the magnetic topography of the solar surface changes so rapidly that it must be monitored constantly as input for space weather forecasts. {A:[25]}

A:25

[The contrast between magnetic variability and gravitational persistence has its origin in the sources of the two fields: the magnetic field, \mathbf{B} , has its origin in a variable source, namely the relative motion of differently charged particles, while the gravitational field, \mathbf{g} , springs forth from a conserved (positive definite) source.] The conserved source of the gravitational field is mass, as can be seen in the [non-relativistic] field equations that apply to the gravitational field:

$$\nabla \cdot \mathbf{g} = -4\pi G\rho, \quad \nabla \times \mathbf{g} = \mathbf{0}, \quad (3.1)$$

where G is the gravitational constant and ρ is the mass density. Thus, gravity is determined by the amount of mass present and its distribution. Because mass is conserved and the gravitational force causes matter to collapse into systems in which the gravitational force is almost perfectly balanced by thermal pressure or inertial forces, gravitationally organized matter tends to be stable

²⁵ Activity: Look up magnetic maps of the solar surface (such as made with the HMI instrument on NASA’s Solar Dynamics Observatory) and make a movie at an image cadence of a few hours; one option to do so is to use HeliViewer. Compare one for 2013–2014 (near cycle maximum, with multiple sunspot groups dispersing into the surrounding network of small-scale flux patterns) to 2017–2018 (around cycle minimum, with only the small scales on the disk).

^{vi} For introductions to the impacts of space weather on society and its technological infrastructure we refer to Chs. H-II:2, H-II:12, and H-V:1-5.

over eons [...] In contrast, the pertinent field equations for the magnetic field are

$$\nabla \cdot \mathbf{B} = 0, \quad \nabla \times \mathbf{B} = \frac{4\pi}{c} \mathbf{j} \quad (3.2)$$

[(the second expression holds if all velocities involved are well below light speed).] The source term for the magnetic field in these equations is electrical current, \mathbf{j} , which, unlike mass, is not a conserved quantity [(although $\nabla \cdot \mathbf{j} = 0$) and which can point in any direction]. Thus we see that \mathbf{B} is a product of dynamo or other magnetohydrodynamic (MHD) processes that generate current in real time. The crucial distinction is that unlike the gravitational field, which is in effect a byproduct of a conserved, definite quantity of mass and so is inherently persistent, the magnetic field is generated by a variety of plasma motions in the Sun, in the solar wind, and in planetary magnetospheres on time scales shorter than what would be needed to reach an equilibrated state. Hence, the local cosmos is constantly adjusting and attempting to relax, but it never gets to such a quasi-stationary state. The consequence of this is what we call weather, including [...] space weather.”

[H-II:1] “There is an important difference regarding the types of volumes that the gravitational and magnetic tension forces organize. The gravitational field has no shielding currents ($\nabla \times \mathbf{g} = \mathbf{0}$) [because its source is the positive-definite mass density ($\nabla \cdot \mathbf{g} = -4\pi G\rho$); consequently, gravity] has no discontinuities because that would require an infinite mass density. Hence, the gravitational field is relatively homogeneous; it varies smoothly and continuously in space. On the other hand, [owing to the fact that electrical charges can be of either sign, a magnetic field can contain] shielding currents ($\nabla \times \mathbf{B} = \frac{4\pi}{c} \mathbf{j}$) which spontaneously form discontinuities,” that are commonly referred to as current sheets [(see Table 3.1 for a definition)] despite the fact that their geometry is generally quite complex in the local cosmos. The combination of the distinct behaviors of gravitational and magnetic forces yields a rich diversity of phenomena in the local cosmos and beyond that emerge from the ‘universal processes’ captured in the MHD description of magnetized plasma.

Among the universal concepts in heliophysics one pair stands out in particular, namely that of *magnetic lines of force* – or commonly *magnetic field lines* – and of their *reconnection*. {A:^[26]} Field lines are abstractions; they are A:26

²⁶ Activity: Throughout this volume we use ‘field line’ only for lines of force of the magnetic field. The concept can be applied to any field, however, including a flow field (such as in Fig. 4.10(C), then often referred to as streamlines) and the gravitational field. Field lines of \mathbf{B} and \mathbf{g} are fundamentally different in one key respect: a magnetic field line never ends (because there are no magnetic monopoles) while gravitational field lines start from a point of mass. What are starting points and/or endpoints (if any) of a system of electrical current (see Activity 28 for the answer)? And of electrical fields? As to magnetic field lines, note that there are drawings in this book, as in many other resources, where field lines are shown to start from one polarity and end on another. As magnetic field has no monopoles, such drawings should not be misread to mean that field lines end, but only that their

Structures in the magnetic field:

Current sheet: *Examples on large scales: heliospheric current sheet; magnetospheric current sheet. E.g., Fig. 5.4. [H-I:6.2]* “Our focus here is mainly on current sheets in the form of tangential discontinuities or rotational discontinuities that evolve into tangential-like discontinuities. Tangential discontinuities are non-propagating surfaces across which no magnetic flux passes as the magnetic field changes direction or strength or both, while total (magnetic plus thermal) pressure is continuous. [...] current sheets (tangential discontinuities) inevitably form in naturally occurring turbulent plasmas; [...] they] form in the corona through the expansion of magnetic flux tubes that poke out of the photosphere [and] expand until at some altitude they press against each other forming a beehive pattern of flux tubes separated by current sheets [unless and until the currents dissipate and the field becomes potential]. [...] Interplanetary space is a honeycomb of outwardly advecting current sheets. [...] In the magnetospheric case, the solar wind snags magnetic field lines from the planet’s two poles on the sunward side and stretches them anti-sunward to form the characteristic two-lobe magnetospheric tail across which the magnetosphere’s analog of the heliospheric current sheet separates the two lobes.”

Flux tube, flux rope: *Examples: compact sunspots, pores, ‘bright points’; as an entity, they are bounded by a current sheet. E.g., Fig. 9.3(top) and 9.1. [H-I:6.4]* “A flux tube is the volume enclosed by a set of field lines that intersect a simple closed curve. The frozen-in flux condition of ideal MHD describes a parcel of plasma threaded by magnetic field lines as a conserved entity whose motion can be followed.” In the solar photosphere, flux tubes may emerge as preformed entities, or may form from by ‘convective collapse’. A flux rope is a flux tube twisted about itself (and thus carrying an internal net current); many magnetic configurations emerge into the solar photosphere as flux ropes; many form in the corona by the dynamics of the reconnecting field; coronal mass ejections inject ropes into the heliosphere (there known as ‘magnetic clouds’) while others form by reconnection across current sheets; at magnetopause, flux ropes (‘flux transfer events’) form by reconnection; and flux ropes (‘plasmoids’) form by reconnection across the magnetospheric current sheet.

Cell: *Examples: planetary magnetosphere; heliosphere. E.g., Fig. 5.1. [H-I:6.7]* “Magnetic fields tied to gravitating bodies will expand to fill all space unless prevented from doing so. [...] the magnetic field’s expansionist ambition is checked by some other magnetic field-bearing plasma expanding from somewhere else. Each magnetic field is therefore encased within a definable volume, which we refer to as a cell. In the Sun’s case, the cell is the heliosphere. In the other cases mentioned, the cells are planetary magnetospheres. [...] The cellular structure] is like a Russian nesting doll in which one cell is encased within another. [...] within the heliosphere, the scale sizes of the objects already mentioned cover seven orders of magnitude”

Table 3.1. *Current sheets, flux tubes, and cells [H-I:6.1] “make up the common forms of heliophysical magnetic structures that exist on MHD time and distance scales (we are not concerned here with kinetic-scale structures that inhabit the dissipation range of turbulence).” Ch. H-I:6 is dedicated to a description of these magnetic structures.*

1-dimensional virtual devices that are used to outline the geometry of magnetic structures in the local cosmos, in a way in which the tangent of the field line anywhere along it has the same direction as the field there while the local field line density is equivalent to the magnetic field strength.

In a vacuum, magnetic field lines have no intrinsic temporal continuity. For example, consider the field between and surrounding magnets or electrical wires at time t_0 and again at time t_1 after having moved the magnets or wires into new positions. The lines of force used to visualize the field at times t_0 and t_1 are completely independent, the result only of the magnetic fields at the two instances in combination with two sets of points, one for t_0 and one for t_1 , selected by a researcher from which to compute the lines of force. In a plasma, however, field lines can be thought of as structures whose continuity in time derives from the ionized matter that is contained in the flux bundle or tube that is centered on the field line. In our thinking, we should map these ‘lines of force’ to their 3-dimensional equivalent, the ‘flux tube’: as long as the ions and electrons once contained never move out of the flux tube, the field line has some temporal continuity. Whenever matter does migrate out of the flux tube, the attribute of continuity for the field line fails. However, if the locations where this occurs are compact compared to the field line’s length, one can think of field lines – that can never end in the divergence-free magnetic field except close onto other field lines – as being cut and connected onto another field line. Where that happens, the concept of ‘magnetic reconnection’ is then introduced to salvage that of the ‘field line’ as something that has an identity over time, at least while matter remains constrained to within the flux tube. A:27

Field lines and their reconnection are but two of the concepts related to a variety of processes that occur in ionized gases (‘plasmas’) that are threaded by magnetic field. We come across such processes in the vastly different environments of the solar interior and of the far reaches of the heliosphere, and in the depths of planets as well as in the most tenuous parts of their outer atmospheres. Temperatures and densities (and, as we shall discuss later in this chapter, magnetic field strengths) differ by many orders of magnitude;

rendering in the diagram is incomplete, *i.e.*, merely terminated for simplicity, for lack of information, or to restrict the discussion to a particular region of interest.

²⁷ Activity: An exception of sorts to the fact that field lines cannot begin or end in a divergence-free field lies in field lines that carry a ‘null point’, which is a point where the magnetic field goes to zero. Draw field lines around a pair of aligned but opposing magnetic dipoles in 2d and identify the null point(s). Then make a 3d rendering from a perspective away from the line connecting the dipoles. Visualize only the set of field lines going through the null(s) (these surfaces are called ‘fans’). Such renderings with charges, nulls, fans (and their intersections, the ‘spines’) are useful tools in analyses of potential magnetic fields of a mixture of charges (such as bipolar solar regions on the solar surface). Consider how such ‘fans’ from nulls would not conflict with the field being divergence-free (the concept of ‘measure’ in set theory helps). For an introduction to the topology of the magnetic field, see Ch. H-I:4.

a summary of some of the conditions encountered in the local cosmos that surrounds us is visually represented in Figure 2.1.

In everyday life we tend to ignore the Earth's magnetic field, but we can do so only because of the low temperatures in which we live (which renders most material, except metals, non-conducting) combined with the high densities; together, as we shall see in more detail later in this chapter, these conditions make the forces exerted by the terrestrial magnetic field utterly negligible in our day-to-day affairs, except where we take special care to uncover them, such as in magnetic compasses. Conditions are markedly different, however, in the layers underneath the atmosphere of the Sun, throughout the extended solar atmosphere, and in the outermost reaches of atmospheres of all bodies in the Solar System: there, magnetism is effectively coupled to matter while the inertia of that matter is in much of the domain significantly lower in comparison to magnetic forces than in our daily settings. There, the magnetic field is an important player that adds a significant force to compete with pressure, gravity, and inertia. It provides a medium for a variety of waves (which this text merely touches upon), and changes the transport of thermal energy and energetic particles. Add to that the fact that the magnetic field is evolving on a range of spatio-temporal scales, and you have a source of continual change in conditions throughout the local cosmos.

The mathematical formulation of what happens in a magnetized plasma is often simplified through an ensemble approximation that is equivalent to the hydrodynamics used in the description of gases, but here including the magnetic field in what is called magnetohydrodynamics, or MHD for short. MHD is a description of the multitude of constituent particles in the local cosmos that relies on statistical averaging carried out by the medium itself, namely through interactions that lead to essentially Maxwellian velocity distributions, often assumed to be isotropic (but in some formulations distinct for directions along and perpendicular to the magnetic field) and for velocity equilibrium between electrons and ions. To this, a few other assumptions are made about local conditions: processes described by MHD assume that ion and electron interactions as well as their gyrations about the field occur on scales that are small compared to the gradients in the magnetic field while at the same time large compared to a distance (known as the Debye length) over which electrical charges can exist unshielded by other particles, with velocities well below relativistic, and only allowing for wave-like phenomena that are slow enough that electrical neutrality is achieved well within any time scale of interest and that are slow compared to the plasma frequency and electron/ion gyro-frequencies. However, interactions between particles should be infrequent enough that the medium should allow the electron and ion populations to move

differentially with relative ease, *i.e.*, conditions should allow the medium to conduct electrical currents rather effectively.

MHD treats the ionized medium as a fluid by working with ensemble properties. In hydrodynamics this is generally allowed because of a high frequency of molecular collisions relative to the time scales of the processes on macroscopic scales. In many environments in heliophysics, however, collisions can be so rare that distances between collisions can be comparable to the scale of the system under consideration, while the solar wind is entirely collisionless beyond a few dozen radii from the Sun . . . and yet MHD has been shown to be a useful approximation. The key factor in making MHD useful is that the medium should not be able to maintain a significant electric field in its own reference frame. Even if collisions are rare in such a medium, long-range flights of the particles are impeded: the gyration of particles about the magnetic field reduces the scale of flight perpendicular to the field, while wave-particle interactions have a similar effect along the field. Consequently, the movement of individual charged particles in a plasma is coupled to the collective of its environment, resulting in a fluid-like behavior even if collisions are rare.

However, where binary interactions are important in the MHD description, the anisotropy imposed by the magnetic field does affect what approximations can be made. Most importantly, these effects are seen on gas pressure and viscosity. In a collisional plasma, these terms are generally essentially isotropic and thus described by scalars. But in a collisionless plasma, pressure and viscosity are anisotropic, and thus are approximated by tensors. In this volume, we generally use a scalar for pressure, and capture anisotropy in conductivity in Hall and Pedersen terms (see below and Ch. 2).

3.2 (Magneto-)Hydrodynamics

The equations of magnetohydrodynamics, or MHD, are based on the assumption that the plasma can be described as a continuum; see Table 3.2 for a very concise description of what that entails. The approximations used here lead to six equations that describe magnetized plasma subject to gravity, as shown in Table 3.3 (note that processes involving radiative transfer are largely omitted from this volume). Five of these are essentially equations of hydrodynamics, namely continuity, momentum, energy, gravity, and the equation of state (EOS), with two important modifications: the magnetic, or Lorentz, force $(1/c)\mathbf{j} \times \mathbf{B}$ ⑥ is added in the momentum description, and there are additional terms ⑩ in the energy equation. We return to these terms and equations below, and discuss the additional equation, namely the induction equation Eq. (3.3), which

Philosophy of magnetohydrodynamics:

The fundamental assumption underlying the MHD equations as shown in Table 3.3, and the principal criterion to judge the applicability of that MHD approximation under given circumstances, is that the medium can be suitably described as a continuum. This presents us with a statistical criterion: MHD can be applied beyond a fiducial length, say L , such that there are sufficient particles in a volume L^3 such that statistical means – like density, mean velocity, pressure and so forth – have small variances or fluctuations about them. Within that volume, collisions (or wave-particle interactions) result in average properties of the medium that transform the need to describe each particle separately in its interaction with all others into an enormously truncated set of descriptions of statistical averages. This truncation is known as ‘closure’: the continuum description requires a closure relation at some level that relates an unknown high-order moment of the full particle distribution function, such as pressure, to lower-order moments (see Sect. 8.3 for more on that). An equation of state, as in Eq. (3.8), is predicated on there being ample collisions to isotropize the random motions and achieve a thermodynamic equilibrium, with its characteristic Maxwellian velocity distribution (or more than one if a multi-fluid description is used). The MHD equations as in Table. 3.3 describe a 5-moment continuum closure scheme using mass density, temperature, pressure, energy density, and velocity. As collisions become less frequent one is required to enforce closure at higher levels, examples of which lead to, *e.g.*, Eqs. (3.11) and (3.27). More generally dielectric and magnetization properties of the material enter in the definitions of D and H . Therefore, if by some other means (*e.g.*, by observation) you know how to close the moments (like in Sect. 2.2 for the solar wind by using the observationally motivated approximation that the temperature is constant throughout the heliosphere and the pressure is an isotropic scalar) then you can use the continuum fluid description to answer some questions even about a medium where collisions are a rare thing.

Table 3.2. *MHD approximation and the concept of ‘closure’.*

A:28

couples the magnetic field to macroscopic flows and microscopic collisions, in some detail in Section 3.2.2. {A:[28]}

In order to assess the validity of the assumption made to derive the MHD equations for the vastly different conditions with which heliophysics concerns itself we can look at a variety of dimensionless numbers. Table 2.5 lists frequently used length and time scales, as well as some commonly used ratios, some of which have been given a name. Some of these are pertinent to microscopic, particle-level conditions and some are pertinent to macroscopic, system-level conditions. We introduce them here only briefly – most will be looked at explicitly later on – in order to give you an impression of which types of processes or relative scales are important. For example, we can look at the length scale on which ions gyrate around the local magnetic field relative to the gradients in the field to assess whether the ions sense the magnetic field in an

²⁸ Activity: Note that $\nabla \cdot \mathbf{B} = 0$ is not needed to complement the MHD equations in Table 3.3 as long as the initial condition satisfies that equation. Take the divergence of Eq. (3.3) to prove that. Use the same operation on $\nabla \times \mathbf{B} = \frac{4\pi}{c} \mathbf{j}$ to show that currents in MHD have no sources or sinks.

Magnetohydrodynamics:		
Induction	$\frac{\partial \mathbf{B}}{\partial t} = \nabla \times \overset{\textcircled{1}}{(\mathbf{v} \times \mathbf{B})} - \nabla \times (\eta \overset{\textcircled{2}}{\nabla \times \mathbf{B}})$	(3.3)
Continuity	$\frac{\partial \rho}{\partial t} + (\mathbf{v} \cdot \nabla) \rho = -\rho \overset{\textcircled{3}}{\nabla \cdot \mathbf{v}} + (S - L) \overset{\textcircled{a}}{}$	(3.4)
Momentum	$\rho \frac{\partial \mathbf{v}}{\partial t} + \rho (\mathbf{v} \cdot \nabla) \mathbf{v} = +\overset{\textcircled{4}}{\rho \mathbf{g}} - \overset{\textcircled{5}}{\nabla p} + \frac{1}{4\pi} \overset{\textcircled{6}}{(\nabla \times \mathbf{B}) \times \mathbf{B}}$ $+ \overset{\textcircled{7}}{\nabla \cdot \boldsymbol{\tau}} - \mathbf{v} \overset{\textcircled{b}}{(S - L)} + (\mathbf{S}_p - \mathbf{L}_p) \overset{\textcircled{c}}{}$	(3.5)
Internal energy	$\rho \frac{\partial e}{\partial t} + \rho (\mathbf{v} \cdot \nabla) e = -p \overset{\textcircled{8}}{\nabla \cdot \mathbf{v}} + \nabla \cdot (\overset{\textcircled{9}}{\kappa \nabla T}) + (Q_\nu + Q_\eta) \overset{\textcircled{10}}{}$	(3.6)
Gravity	$\nabla \cdot \mathbf{g} = \nabla^2 \Phi = -4\pi G \rho$	(3.7)
EOS	$p = (\gamma - 1) \rho e$	(3.8)
Complemented by initial and boundary conditions		
Online resources:		
Plasma physics:	NRL Plasma Formulary	
Vector calculus:	Wikipedia	
Introduction to MHD	'Essential magnetohydrodynamics for astrophysics' (Spruit, 2013)	

Table 3.3. *Equations of magnetohydrodynamics for a fully-ionized plasma, ignoring radiative energy transport and radiation pressure, to be complemented by initial and boundary conditions to specify the solution. Symbols: \mathbf{B} magnetic field; \mathbf{v} fluid velocity; $e = C_V T$ specific internal energy; p gas pressure; ρ mass density; Φ the gravitational potential and G Newton's gravitational constant; \mathbf{g} gravity, $\tau_{ik} = 2\rho\nu (\Lambda_{ik} - \frac{1}{3}\delta_{ik}\nabla \cdot \mathbf{v})$ the viscous stress tensor with the deformation tensor $\Lambda_{ik} = \frac{1}{2} (\frac{\partial v_i}{\partial x_k} + \frac{\partial v_k}{\partial x_i})$; Q_ν viscous heating; and $Q_\eta = \eta \mathbf{j}^2 = \eta (c/4\pi)^2 (\nabla \times \mathbf{B})^2$ the resistive (Ohmic) dissipation; ν , η_e and κ represent the viscosity, magnetic diffusivity, and the thermal conductivity tensor (which is highly anisotropic, with heat most effectively conducted by electrons moving along the magnetic field); $\gamma = C_p/C_V$ is the adiabatic index, the ratio of specific heats for constant pressure and constant volume. In an ideal, mono-atomic gas with 3 degrees of freedom $\gamma = 5/3$. S , L , \mathbf{S}_p and \mathbf{L}_p are source and loss terms for mass and momentum by introduction or loss of ions from a non-ionized reservoir.*

ensemble sense such as required of a fluid or whether higher-order descriptions are needed. Or one can ask whether length scales involved are large enough that the plasma can be viewed as not having significant charge separation; the length scale on which electrostatic potential of any particle is effectively shielded by the surrounding plasma is known as the Debye length. Or one can look at the ratio of the average time between collisions and the time needed to complete one gyration around the magnetic field in order to assess whether the magnetic field can effectively be followed by the charged particles and whether the Hall current needs to be considered.

3.2.1 MHD equations, individual terms, and special cases

First, let us briefly review what the MHD equations express, the role of the individual terms, and some special cases:

- Eq. (3.3): • The induction equation (a combination of Faraday's law with Ohm's law, see Sect. 3.2.2) states that any local change in the magnetic field is associated with a 'curl', or 'circulation', in the component of plasma flows working perpendicular to the magnetic field and/or to the slippage of plasma relative to the magnetic field through finite diffusivity. Note that this form of the induction equation is linear in \mathbf{B} so that if $\mathbf{B}(t = 0) = \mathbf{0}$ then no field can arise at a later time. Sect. 3.2.2 touches on the fact that some terms were ignored to arrive at this form, some of which can act as a source term for magnetic field; this is not further discussed in this volume as the clouds out of which stars and planetary systems form initially are threaded by a galactic seed field from the outset (interested readers could look for 'battery effects', including the 'Biermann battery').
- Eq. (3.4): • Continuity requires that the local plasma density changes only because of flow through a volume and by compression or dilation in doing so.
- Eq. (3.5): • The momentum (or force) equation (Newton's second law in volumetric form) summarizes how the plasma velocity is affected, as in hydrodynamics, by gravity, pressure gradients, and viscosity, but here also by the Lorentz force associated with currents flowing across the magnetic field.
- Eq. (3.6): • The local energy density (here shown in a per-mass formulation of the first law of thermodynamics) is affected by flows, including compression or dilation, thermal conduction, and by viscous and resistive heating.
- Eq. (3.7): • As mass is a positive definite quantity, it can only strengthen gravity, which can be represented by the gradient of a potential.
- Eq. (3.8): • The equation of state couples pressure, density, and internal energy.
- Eqs. (3.3, 3.5): • The induction and momentum equations are derived from the (mass-weighted) difference (see Sect. 3.2.2) and the sum of the equations of motions

for the electrons and for the ions, each of which includes a term for their collisional coupling.

- ①: • In case the term ② is negligible, Eq. (3.3) describes what is known as ‘ideal MHD’. In this case (see Sect. 3.4) the plasma and magnetic field must move with each other for velocity components perpendicular to the magnetic field, whereas plasma movement along the field is not affected by that field. In this condition, the field is said to be ‘frozen in’ the plasma. In such a state, the lines of force (‘field lines’) are advected with the flow while unable to break their connectivity between any plasma elements along their length; in non-ideal, or resistive, MHD such connections can be broken through a process known as ‘reconnection’. The concepts of field lines and reconnection are described in Section 3.4.
- ②: • This term quantifies the effects of resistivity on the magnetic field by the dissipation and diffusion of the electrical current $\mathbf{j} = \frac{c}{4\pi} \nabla \times \mathbf{B}$. If the magnetic diffusivity η is constant throughout the medium, then term ② can be rewritten as $\eta \nabla \times (\nabla \times \mathbf{B}) = -\eta \nabla^2 \mathbf{B}$ (because $\nabla \cdot \mathbf{B} = \mathbf{0}$), which shows that it causes the magnetic field to decay diffusively; in the absence of ①, such as in a stationary plasma, this makes Eq. (3.3) a diffusion equation for decaying magnetic field.
- ③: • For an incompressible fluid, ρ is constant as material flows throughout the volume under study, which consequently means that $\nabla \cdot \mathbf{v} = \mathbf{0}$, *i.e.*, that the velocity field is divergence free, and – unless there are terms like ① to consider – Eq. (3.4) vanishes from the set. That also removes term ⑧ from Eq. (3.6), so that the energy density of the medium can only change by thermal conduction ⑨ and by viscous and resistive dissipation ⑩ (disregarding here, as we do throughout this chapter, the effects of radiation).
A:29
- ⑤: • As formulated here, the isotropic part of the pressure tensor is expressed as a scalar, while the other terms are captured in the stress tensor. If only this scalar term is carried, then particle microscopic velocity distributions are taken to be isotropic.

²⁹ Activity: The so-called ‘Boussinesq approximation’ is intermediate to fully compressible and incompressible, and in principle internally inconsistent: it assumes a fluid for which (and in numerical codes replaces Eq. (3.4) by) $\nabla \cdot \mathbf{v} = \mathbf{0}$ but allows density variations in the term in the force balance that includes gravity (and thus allows for buoyancy). This approximation works well if the flow can be characterized as ‘nearly incompressible’. For settings where the scale of the density stratification is large compared to processes of interest the incompressible approximation can be valid; in such settings, compressibility becomes only important in structures like shock waves, but is ignorable if the flows are much slower than the sound and Alfvén speeds. Advanced, for the curious: In planetary atmospheric envelopes and stellar interiors alike, zones of relatively low temperature under relatively strong gravity are highly stratified compared to the scales of flows within them. In such settings, numerical codes have been developed under the ‘anelastic’ approximation. This approximation provides a better description of the density in stratified settings than the pure Boussinesq one while filtering out sound waves that would require much higher spatio-temporal resolution of the code. This article by Durran and Arakawa (2007) introduces and compares several ‘anelastic’ approximations.

- ⑥: • Term ⑥ measures the interaction of the Lorentz force and the plasma flow. The vector product $(\nabla \times \mathbf{B}) \times \mathbf{B}$ can be reformulated (see Eq. 3.23) into the sum of a pressure-like term (that works to expand unless countered) and a term that is equivalent to a tension (which works to straighten unless countered), showing that the magnetic field in a plasma behaves as if it were both like a gas and like a flexing rod or taut string. There is a special class of magnetic fields in which currents run parallel to the magnetic field; in that case $(\nabla \times \mathbf{B}) \times \mathbf{B} = \mathbf{0}$, *i.e.*, there is no Lorentz force, and these are consequently referred to as 'force-free fields', of which the potential field is a special (and lowest-energy) state. As the field is parallel to the current, there is a scalar field α such that $\nabla \times \mathbf{B} = \alpha \mathbf{B}$. If α is a uniform constant, this field is called 'linear force free' (which is mathematically easier to work with, but does not develop in general astrophysical settings); if not, the corresponding field is a 'non-linear force-free' field (to which we return in Sect. 6.3).
- ⑨: • As for term ② with uniform magnetic diffusivity η , here a uniform thermal conductivity κ would allow rewriting of term ⑨ to be proportional to $\nabla^2 T$, quantifying diffusion of thermal energy.
- ①, ②, ③: • These terms reflect source and loss terms for mass and momentum density per unit volume through, *e.g.*, (de-)ionization of neutrals (including charge exchange) that are important, for example, where comets add gas and dust or around geysers on low-gravity moons (Sec. 2.4), or to the inflow of neutral matter into the solar wind from outside the heliosphere.
- $\mathbf{B} = \mathbf{0}$: • A field-free state (or a non-conducting, and thus current-free, gas in which the field does not apply force to the gas; see also under 'Potential' below) transforms Eqs. (3.3)–(3.8) into regular hydrodynamic equations.
- $\mathbf{v} = \mathbf{0}$: • A static plasma is described by Eqs. (3.3)–(3.8) without terms ①, ③, ⑦, ⑧, and Q_ν in ⑩. Moreover, with no flows, no change can occur that involves bulk flows, so that, for example, the lefthand side of the momentum Eq. (3.5) has to equal $\mathbf{0}$. This yields an equation for magnetohydrostatic balance in which gravity, pressure gradient, and Lorentz force sum to zero.
- $\frac{\partial}{\partial t} = \mathbf{0}$: • Stationary situation in which none of the variables can change. In particular, $\frac{\partial \mathbf{v}}{\partial t} = \mathbf{0}$ is a situation with stationary flows, which can be maintained only for limited times.
- Potential: • In the case of a potential field, there are no currents in the system, *i.e.*, $\nabla \times \mathbf{B} = \mathbf{0}$. Consequently, term ⑥ vanishes because there is no Lorentz force. Term ② also vanishes, leaving only term ① in the righthand side of induction Eq. (3.3) (equivalent to the infinitely conducting case of ideal MHD with frozen-in field, or in which the field is maintained from outside of a current-free volume). To see to full consequence of this state, however, we

need to realize that $\nabla \times \mathbf{B} = \mathbf{0}$ means that there is a magnetic potential Φ_m such that $\mathbf{B} = -\nabla\Phi_m$ from $\nabla^2\Phi_m = 0$. Such a Laplace equation, once the boundary condition is specified, has a unique solution. And for a current-free system with fixed boundary conditions that, in turn, means that \mathbf{B} cannot change in time, such that term ① then implies that there is a scalar field Ψ such that $\mathbf{v} \times \mathbf{B} = \nabla\Psi$, of which one particular case has $\mathbf{v} \parallel \mathbf{B}$.

Force free: • See at ⑥ above in this listing.

{A:^[30]}

A:30

3.2.2 The induction equation

[H-I:3.2] “The induction equation, Eq. (3.3), arises from Ohm’s law combined with the non-relativistic approximation of the Maxwell equations. In its most general form Ohm’s law is a relation between electric current, electric field, magnetic field, plasma motions and electron pressure gradients. Ohm’s law is derived from an equation of motion for electrons in which the interaction with ions (defining the bulk motion of the plasma with velocity \mathbf{v} [because the ions, here taken to be dominated by singly-ionized species, by far outweigh the electrons]) is described through a collisional drag term related to the differential motion:

$$n_e m_e \frac{d\mathbf{v}_e}{dt} = -n_e e (\mathbf{E} + \frac{1}{c} \mathbf{v} \times \mathbf{B}) - \nabla p_e + n_e m_e \frac{\mathbf{v} - \mathbf{v}_e}{\tau_{ei}}. \quad (3.9)$$

Here \mathbf{v}_e denotes the electron velocity, τ_{ei} the collision time between electrons and ions, e the electron charge, m_e the electron mass, n_e the electron density, and p_e the electron pressure” (omitting gravity). \mathbf{E} and \mathbf{B} are the electric and magnetic vector fields.

By noting that the differential velocity between ions and electrons is proportional to the current,

$$\mathbf{j} = n_e e (\mathbf{v} - \mathbf{v}_e) \quad (3.10)$$

we can reformulate Eq. (3.9), when combined with the analogous version for

³⁰ Activity: **What if radiative transfer were included?** The MHD equations in Table 3.3 do not incorporate electromagnetic radiation. In a sufficiently dense medium, in which the photon mean free path is small compared to plasma and field gradients, energy transport by electromagnetic radiation can be described by a diffusion equation. Where the mean-free path is long, however, energy can ‘jump’ between different locations without (or with weak) coupling to the intermediate medium, in a manner that depends on wavelength as well as on atomic properties. With that in mind, contrast the solar interior to its atmosphere; a cloud-free planetary atmosphere to a (partially) clouded one; and (maybe once you get to Ch. 11) initial to later phases of star formation and of protoplanetary disks. In the context of that question and other assumptions going into the MHD equations in Table 3.3: Why is the solar chromosphere the hardest part of the solar interior and atmosphere to describe? And what makes a terrestrial ionosphere hard to capture in equations? Some of the answers to these questions will come as you read along. For an introduction to radiative transfer in stellar atmospheres, see this freely available online text by Rutten (2003): URL.

the ions, to yield a formulation of Ohm's law (here ignoring electron inertia and assuming pressure to be a scalar, *i.e.*, isotropic; compare with Section 3.4):

$$\mathbf{j} = \frac{\tau_{ei}n_e e^2}{m_e} \left(\mathbf{E} + \frac{1}{c} \mathbf{v} \times \mathbf{B} \right) - \frac{\tau_{ei} e}{m_e c} \mathbf{j} \times \mathbf{B} + \frac{\tau_{ei} e}{m_e} \nabla p_e. \quad (3.11)$$

A:31 {A:[31]} Note that when the electric field expressed through the electron pressure gradient is ignored, this equation can be rewritten to an equivalent Ohm's law discussed for the ionosphere in Ch. 2 that has a conductivity tensor with components as in Eqs. (2.23)-(2.25), from which terms with M_i disappear for the fully-ionized plasma because there are only ion-electron collisions, and in which the electron magnetization subject to collisions with neutrals, M_e , is replaced by $M_{ei} = eB/(m_e \nu_{ei})$ for $\nu_{ei} = 1/\tau_{ei}$. In other words, the Hall term in Eq. (3.11) takes care of the anisotropic part of the conductivity in the fully ionized plasma. The Pedersen current, directed along \mathbf{E} is part of the first term on the right-hand side. {A:[32]}

A:32

Specifically, [H-I:3.2] “the second term on the right-hand side describes the Hall current, which becomes important if the collision time is longer than the electron gyration time, *i.e.*, when $\tau_{ei} \omega_L > 1$, where $\omega_L = eB/m_e$ denotes the Larmor (or [electron] gyro-)frequency. The Hall term leads to anisotropic plasma conductivity with respect to the magnetic field direction and is typically important in low-density plasmas in which τ_{ei} can be very large”. In many settings in heliophysics, the last two terms in Eq. (3.11) are ignored “(unless high-frequency plasma oscillations are considered), leading to the simplified Ohm's law

$$\mathbf{j} = \sigma_e \left(\mathbf{E} + \frac{1}{c} \mathbf{v} \times \mathbf{B} \right) \quad (3.12)$$

with the plasma conductivity

$$\sigma_e = \frac{\tau_{ei} n_e e^2}{m_e}. \quad (3.13)$$

Using Ampère's law, $\nabla \times \mathbf{B} = \frac{4\pi}{c} \mathbf{j}$, yields for the electric field in the laboratory frame

$$\mathbf{E} = -\frac{1}{c} \mathbf{v} \times \mathbf{B} + \frac{c}{4\pi\sigma_e} \nabla \times \mathbf{B} \quad (3.14)$$

³¹ Activity: Formulate the ion equivalent of Eq. (3.9) (remember Newton's third law, and with $m_e/m_i \rightarrow 0$) and derive Eq. (3.11) (using $m_e \mathbf{v}_i + m_i \mathbf{v}_e = m_i \mathbf{v}_i + m_e \mathbf{v}_e + m_i [\mathbf{v}_e - \mathbf{v}_i] + m_e [\mathbf{v}_i - \mathbf{v}_e]$) and also the corresponding momentum equation (absent gravity): $\rho d\mathbf{v}/dt = \mathbf{j} \times \mathbf{B} - \nabla p$. Then add gravity and compare to Eq. (3.5).

³² Activity: Demonstrate, for a fully-ionized single-species plasma, the equivalence of Eq. (3.11) and $\mathbf{j} = \Sigma_e \cdot (\mathbf{E} + \frac{1}{c} \mathbf{v} \times \mathbf{B})$ with Eqs. (2.22)-(2.25).

leading to the induction equation through one of the Maxwell equations:

$$\frac{\partial \mathbf{B}}{\partial t} = -c \nabla \times \mathbf{E} = \nabla \times (\mathbf{v} \times \mathbf{B} - \eta \nabla \times \mathbf{B}) \quad (3.15)$$

with the magnetic diffusivity

$$\eta = \frac{c^2}{4\pi\sigma_e}. \quad (3.16)$$

In MHD, the equations are typically expressed in terms of the magnetic field \mathbf{B} and flows \mathbf{v} , with electric fields and currents eliminated from the system. This is done primarily out of mathematical convenience, since formulating the problem in terms of currents leads to intractable equations involving integrals of the currents over the entire volume under study.”

Whether a formulation in terms of the magnetic field or electrical currents is more convenient also depends on the inhomogeneity and anisotropy of the conductivity. The most extreme example is electrical engineering where cables give full control over the current, and thus a current-based description is clearly the method of choice. A formulation in terms of currents can be easier to work with also when currents can only flow along the field or are restricted to relatively thin layers with high conductivity, such as is the case in the ionosphere. In most MHD problems with highly conducting fluids, however, there is no *a priori* control over where currents flow, so that dealing with the magnetic field is typically the better choice. Because of this, solar and heliospheric physicists generally use arguments primarily based on the magnetic field; in space physics, however, and in particular in ionospheric physics, currents are often discussed.

Of interest to the induction equation Eq. (3.3) is the relative importance of the advection and diffusion-like terms on the right-hand side. One way to assess that is to reformulate it into characteristic scales and a frequently occurring dimensionless number: [H-I:3.2.3] “Let L_t be a typical length-scale and v_t a characteristic velocity of the problem. Expressing the time in units of L_t/v_t and the spatial derivatives in the induction equation Eq. (3.3) in units of L_t leads to the dimensionless form of the induction equation

$$\frac{\partial \mathbf{B}}{\partial t} = \nabla \times \left(\mathbf{v} \times \mathbf{B} - \frac{1}{\mathcal{R}_m} \nabla \times \mathbf{B} \right) \quad (3.17)$$

with the magnetic Reynolds number

$$\mathcal{R}_m \equiv \frac{v_t L_t}{\eta}. \quad (3.18)$$

The limit $\mathcal{R}_m \ll 1$ is referred to as diffusion dominated regime, in which the

(dimensional) induction equation reduces to a diffusion equation of the form

$$\frac{\partial \mathbf{B}}{\partial t} = \eta \nabla^2 \mathbf{B}. \quad (3.19)$$

Here we made the additional simplifying assumption of a constant magnetic diffusivity η . Assuming that the magnetic field has a typical length scale L_t , we can estimate [its decay time scale:]

$$\tau_d \sim \frac{L_t^2}{\eta}. \quad (3.20)$$

The limit $\mathcal{R}_m \gg 1$ is referred to as the advection-dominated regime, in which the induction equation reduces to the equation of ideal MHD (except for possible boundary layers where diffusivity could be still important)

$$\frac{\partial \mathbf{B}}{\partial t} = \nabla \times (\mathbf{v} \times \mathbf{B}). \quad (3.21)$$

Expanding the expression of the right-hand side of this ideal induction equation leads to

$$\frac{\partial \mathbf{B}}{\partial t} = -(\mathbf{v} \cdot \nabla) \mathbf{B} + (\mathbf{B} \cdot \nabla) \mathbf{v} - \mathbf{B} (\nabla \cdot \mathbf{v}). \quad (3.22)$$

While the first term on the right-hand side describes the advection of magnetic field, the last two terms describe the amplification by shear (second term) and compression (third term)."

Of interest to the momentum equation Eq (3.5) is that a vector identity allows us to reformulate the Lorentz force [H-I:3.2] "to equal:

$$\frac{1}{4\pi} (\nabla \times \mathbf{B}) \times \mathbf{B} = -\nabla \frac{B^2}{8\pi} + \frac{1}{4\pi} (\mathbf{B} \cdot \nabla) \mathbf{B}, \quad (3.23)$$

which shows that the Lorentz force is a sum of an isotropic pressure-like force and a tension force related to the curvature of the field" (note that both of these are insensitive to a reversal of the direction of the magnetic field). Because the pressure and tension terms, as does therefore the full Lorentz force, scale as $\mathcal{O}(B_t^2/L_t)$ (where the subscript 't' denotes a typical value of the quantity) they can be compared in magnitude to the pressure gradient force $\mathcal{O}(p_t/L_t)$; the ratio of magnetic and gas pressure terms in Eq. (3.5) yields an often-used dimensionless number in heliophysics, the plasma β :

$$\beta \equiv \frac{8\pi p}{B^2}. \quad (3.24)$$

A:33

{A:[33]}

³³ Activity: Look back at Fig. 2.1 and review the ranges shown of the value of the plasma β from Eq. (3.24) to get a feel for where plasma pressure gradients might dominate magnetic pressure gradients or vice versa. Add lines for unit plasma β for field strengths of 1 μG (as found in the outer

3.3 Waves in magnetized plasmas

Before we proceed with a discussion of field lines and reconnection, we look into an important aspect of a magnetized plasma, namely how it carries waves. Waves are important, among other things, in communicating information about changes in the field's structure or in boundary conditions or the effects of obstacles embedded in flows, while moreover they transport energy. [H-IV:10.2.1] "The waves that carry information through a magnetized plasma differ from the sound waves of a neutral gas, partly because of the anisotropy imposed on the fluid by a magnetic field and partly because the waves must be capable of carrying currents that modify the properties of both matter and magnetic field. The properties of such waves can be derived from the MHD [equations] by analyzing the evolution of small perturbations.

Consider a uniform plasma with constant pressure and density (p and ρ) whose center of mass is at rest ($\mathbf{v} = 0$). Assume that a constant background field (\mathbf{B}) is present and that neither sources nor losses need be considered. Small departures from this background state are taken to vary with space (\mathbf{x}) and time (t) as $e^{i(\mathbf{k}\cdot\mathbf{x}-\omega t)}$. Here, \mathbf{k} is the wave vector and ω is the angular frequency of the wave. Perturbations occur in density $d\rho$, velocity $d\mathbf{v}$, pressure dp , current \mathbf{j} , and field \mathbf{b} . Terms linear in small quantities in Eqs. (3.4) and 3.5) satisfy

$$-\omega d\rho + \rho \mathbf{k} \cdot d\mathbf{v} = 0 \quad (3.25)$$

$$-\omega \rho d\mathbf{v} = -\mathbf{k} dp + \frac{1}{4\pi} \mathbf{b}(\mathbf{k} \cdot \mathbf{B}) - \frac{1}{4\pi} \mathbf{k}(\mathbf{b} \cdot \mathbf{B}). \quad (3.26)$$

[If we assume an isentropic (*i.e.*, adiabatic and reversible) process, then Eq. 3.8 becomes $p\rho^{-5/3} = \text{constant}$, so that] the pressure perturbation in terms of the density perturbation is

$$dp/p = \gamma d\rho/\rho \quad (3.27)$$

and [the ideal induction equation Eq. (3.3 (with $\eta \equiv 0$))]

$$\omega \mathbf{b} = d\mathbf{v}(\mathbf{k} \cdot \mathbf{B}) - \mathbf{B}(\mathbf{k} \cdot d\mathbf{v}). \quad (3.28)$$

The solutions to Eqs. (3.25) to (3.28) are the roots of the equation

$$(\omega^2 - v_A^2 k^2 \cos^2 \theta)[\omega^4 - \omega^2 k^2 (c_s^2 + v_A^2) + k^4 v_A^2 c_s^2 \cos^2 \theta] = 0, \quad (3.29)$$

where θ is the angle between \mathbf{k} and \mathbf{B} , and the Alfvén speed (v_A) and the

heliosphere and interstellar medium; see Chs. 5 and 10) and for 0.1 MG (considered characteristic of the field strength of flux bundles at the bottom of the solar convective envelope where the principal processes in the solar dynamo are considered to operate; see Ch. 4).

sound speed (c_s) have been introduced. These quantities characterize the speed of propagation of waves in a magnetized plasma and are defined by

$$v_A^2 = B^2/4\pi\rho, \quad (3.30)$$

$$c_s^2 = \gamma p/\rho. \quad (3.31)$$

The sound speed has the form familiar for a neutral gas. The Alfvén speed is a second natural wave speed characteristic of a magnetized plasma. Just as [one can work with the dimensionless] sonic Mach number as the ratio of the flow speed to the sound speed, it is useful to define a dimensionless Mach number, the Alfvénic Mach number

$$M_A = v/v_A, \quad (3.32)$$

related to the Alfvén speed.

As mentioned [at Eq. (3.23)], the quantity $B^2/8\pi$ is the pressure exerted by the magnetic field, so both of the basic wave speeds are proportional to the square root of a pressure divided by a density. [... When $\beta = 8\pi p/B^2 \ll 1$,] magnetic effects dominate the effects of the thermal plasma, but in a high- β plasma, the plasma effects dominate.

Equation (3.29) is of sixth order in ω/k with three pairs of roots. One pair results from setting the first factor in Eq. (3.29) to zero; the resulting dispersion relation is

$$(\omega^2 - v_A^2 k^2 \cos^2 \theta) = 0. \quad (3.33)$$

This solution describes waves referred to as Alfvén waves. For this dispersion relation to apply, the magnetic perturbation must be perpendicular to both \mathbf{B} and \mathbf{k} (see Fig. H-IV:10.1a). This orientation implies that to first order in small quantities, the Alfvén wave does not change the field magnitude [$(\mathbf{B} + \mathbf{b})^2 = B^2 + 2(\mathbf{B} \cdot \mathbf{b})^2 + b^2 \approx B^2$]. The wave phase speed is $v_{\text{ph}} = \omega/k$ and $v_{\text{ph}} = \pm v_A \cos \theta$. [Wave packets] carry information at the group velocity, $\mathbf{v}_g = \nabla_k \omega$, where the subscript on the gradient indicates that the derivatives are taken in \mathbf{k} space; the solution is $v_g = \pm \hat{\mathbf{B}} v_A$ where $\hat{\mathbf{B}}$ is a unit vector along the background field. The remarkable property of these waves is that they carry information only along the background field, and they bend the field without changing its magnitude. These properties are of considerable importance in interpreting the interaction of a flowing plasma with the solid bodies of the Solar System [(discussed in Ch. 5)].

Eq. (3.29) has two more pairs of roots, the zeroes of the fourth order polynomial in square brackets in Eq. (3.29), *i.e.*, the solutions

$$v_{\text{ph}}^2 = \omega^2/k^2 = \frac{1}{2} \left(c_s^2 + v_A^2 \pm [(c_s^2 + v_A^2)^2 - 4v_A^2 c_s^2 \cos^2 \theta]^{1/2} \right). \quad (3.34)$$

The solutions (two pairs, one positive and one negative, of roots) correspond to what are unimaginatively referred to as fast-mode (or magnetosonic) and slow-mode waves. The wave perturbations of both modes may have magnetic perturbations along and across \mathbf{B} (see Fig. H-IV:10.1b). Perturbations along \mathbf{B} change the field magnitude and the thermal pressure. The fast mode changes of thermal and magnetic pressure are in phase with each other; this implies that the total pressure fluctuates. The slow mode changes of thermal and magnetic pressure are in antiphase, and the total pressure fluctuations are very small. For waves propagating along the background field ($\cos\theta = \pm 1$), the solutions to Eq. (3.34) are c_s^2 and v_A^2 , with the larger of the two applying to the fast mode. For waves propagating at right angles to the background field ($\cos\theta = 0$), the [solutions] are $c_s^2 + v_A^2$ and 0, indicating that only fast-mode waves propagate across the field.” {A:[34]}

A:34

3.4 MHD, magnetic field lines and reconnection

[H-I:4.1] “One of the most idiosyncratic aspects of space physics is the central role assigned to magnetic field lines. Particularly in studies of the Sun, the heliosphere and the magnetosphere, magnetic field lines are treated as full-fledged physical objects with their own dynamics. The electrical current, when needed, is derived *from* the magnetic field lines. These practices appear at odds to the basic approach, followed in elementary electrodynamics, of deriving the magnetic field *from* a current distribution and treating magnetic field lines at best as fictitious curiosities. However, physical laws such as Ampère’s law (without displacement current [because velocities are assumed to be well below relativistic]),

$$\nabla \times \mathbf{B} = \frac{4\pi}{c} \mathbf{j}, \quad (3.35)$$

do not attribute a causative nature to either side of the equality; they simply state the equality of two quantities. So either approach to satisfying Eq. (3.35), beginning with either \mathbf{j} or \mathbf{B} , is a valid one.

[The central role of \mathbf{B} in space physics has been furthered tremendously by the introduction of the concept of the field line.] A magnetic field line, sometimes called a line of force, is a space-curve $\mathbf{r}(\ell)$ which is everywhere tangent to the local magnetic field vector, $\mathbf{B}(\mathbf{x})$. This description can be cast as the differential equation

$$\frac{d\mathbf{r}}{d\ell} = \frac{\mathbf{B}[\mathbf{r}(\ell)]}{|\mathbf{B}[\mathbf{r}(\ell)]|}, \quad (3.36)$$

whose solution, starting from some initial point $\mathbf{r}(0)$, is a magnetic field line.

³⁴ Activity: Compare values for c_s and v_A for the environments listed in Table 3.4.

[...] A field line is a curve, and therefore has zero volume. A *flux tube* may be constructed by bundling together a group of field lines. The net flux, Φ , of the tube is the integral $\int \mathbf{B} \cdot d\mathbf{a}$ over any surface pierced by the entire tube. Because $\nabla \cdot \mathbf{B} = 0$, the tube must have the same flux at every cross section.

The only way, in general, to find a field line is to integrate the differential Eq. (3.36). A solution to the field line equation, Eq. (3.36), can in principle be found for a magnetic field at any instant. What is not immediately evident is why such a curve should be physically significant, even if one concedes that the magnetic field itself is significant. There is, in fact, no single reason that field lines will be significant under general circumstances — this is why students are often warned not to attribute undue importance to them. There are, however, numerous circumstances arising in space physics whereby a magnetic field line can achieve a degree of [utility]. The following is a brief list of the most common, applicable to a wide variety of plasma regimes from general ((i)), to the fluid regime ((ii)), to MHD ((iii)), to ideal MHD ((iv)).

- (i) *General: single particle motion.* Subject to no other forces than a relatively stationary magnetic field, [the guiding center of a charged particle will remain tied to a single field line while the particle gyrates about that] according to its mass and charge [(discussed in detail in Sect. 8.1)]. Drifts will displace the particle's guiding center by several gyro-radii after it has traversed a length comparable to the field's curvature radius or gradient scale. Global scales of space plasmas are typically much, much greater than the gyro-radii of their electrons, and to a lesser extent of their heavier ions (the Earth's geomagnetic ring current is a counterexample to this[; see Activity 23]). Waves in the field may scatter particles (important in, *e.g.*, the Earth's radiation belts, [and for solar and galactic cosmic rays propagating through the heliosphere, see Ch. 8]), but this too is generally unimportant. Field lines therefore serve as excellent approximations of the electron orbits. [...]
- (ii) *Fluid regime: thermal conductivity and solar coronal loops:* In a diffuse, high-temperature plasma, thermal energy is conducted principally by electrons. When electrons are strongly magnetized (*i.e.*, the cyclotron frequency is much greater than the collision frequency) their orbits will follow field lines over long distances between collisions at which point they scatter a perpendicular distance no greater than a single gyro-radius. The huge disparity between parallel and perpendicular scattering distances makes thermal conductivity highly anisotropic. Consequently, heat is conducted parallel to the magnetic field far more readily than perpendicular to the field.

Due to this anisotropic conductivity, heat deposited somewhere in a plasma is rapidly and efficiently conducted to all points on the same field line, at least while collision frequencies remain relatively low. [In the coronal setting, for example, the] plasma β is also generally low, so plasma flows are mechanically confined by the field. This means that a bundle of field lines will behave as a one-dimensional autonomous atmosphere, at least as long as reconnection is relatively unimportant. [...]

- (iii) *MHD: Alfvén wave propagation:* Low-frequency waves in a magnetized plasma [(see Section 3.3)] comprise three branches: slow magnetosonic, fast magnetosonic and shear Alfvén waves. The group velocity of the shear Alfvén wave is exactly parallel to the local magnetic field. In the limit of very short wavelengths, any small localized disturbance will therefore propagate along a path following a magnetic field line. This means that a given field line will 'learn' of perturbations anywhere along itself at the Alfvén speed. In this sense the magnetic field line has a dynamical integrity similar to that of a piece of string. Indeed, it is common to derive the Alfvén speed intuitively using the analogy to a string under tension. [...]
- (iv) *Ideal MHD: frozen-in field lines:* [...] At its simplest, the frozen-in-field-line theorem states that if two fluid elements lie on a common field line at one time, then they lie on a common field line at all times past and future. This follows directly from the ideal induction equation, ([Eq. (3.3) with $\eta \equiv 0$]), and from the fact that fluid elements move at the same velocity \mathbf{v} that appears in it."

In Ch. H-I:4 you can see why the mathematics of ideal MHD is such [H-I:4.1] "that differentiation along a field line is interchangeable with differentiation along a flow trajectory. From this it follows that a field line linking two fluid elements can be traced either before or after following the flow of those elements. That is a restatement of the [frozen-in-field-line] theorem introduced above. One can thereafter imagine 'labeling' all the fluid elements along a given field line and then following those fluid elements as they move at their own velocities, \mathbf{v} . These material elements, which are manifestly real, will trace out a single field line at all times, so that [the field line is a useful concept in thinking about plasma motions. Wherever $\eta \neq 0$ in Eq. (3.3) field lines lose their nature as coherent entities; more on this below where we discuss reconnection.]"

Field lines and flux tubes have taken on a remarkable degree of utility in the thinking of many working in the various branches of heliophysics. [H-I:2.5]

“The motion of plasma along the magnetic field does not stress the field and incurs no dynamic back-reaction on the plasma through the action of the Lorentz force. Magnetic field lines therefore serve as conduits for moving energy, mass, momentum and energetic particles from point to point in the heliosphere. Heliophysics accordingly focuses on the magnetic connectivity of the Earth to the Sun, of the magnetotail to the polar caps, of the Io plasma torus to the Jovian magnetic field, and so forth. Magnetic field lines are truly the interstate highway system, the Autobahn network, the autostrada web of the heliosphere.” {A:[35]}

A:35

Field lines as true, persistent entities have their greatest utility in ideal (non-resistive) MHD. But ideal MHD, in which field lines always connect the same parcels of plasma, fails when magnetic diffusivity becomes important in the MHD approximation, or when the basic assumptions of MHD itself fail on the smallest time or length scales. Then field is no longer ‘frozen in’ wherever that happens, and the very concept of continuity of field lines in space and time loses validity. Failure of the field-line concept as it is discussed above is captured by the term ‘reconnection.’ This term, widely used, turns out to be very loosely defined. [H-II:1] “It can be used to refer to the changing connectivity in a vacuum potential field as much as to the decoupling of particle motions from the background magnetic field by any number of concepts, ranging from inertia to wave-particle interactions, or from resistivity to infinitesimal current sheets. It is thus as much a culturally accepted term for something that we really do not understand, as a descriptor of a well-understood consequence: we can say that reconnection occurs whenever the approximation of frozen-in flux fails.”

Non-ideal MHD sees reconnection as a consequence of resistivity. [H-I:5.2.2] “To determine a realistic resistivity for a collisionless plasma requires consideration of the generalized Ohm’s law. For a fully ionized plasma it can be written as

$$\mathbf{E} = -\frac{1}{c} \mathbf{v} \times \mathbf{B} + \frac{\mathbf{j}}{\sigma_e} + \frac{m_e}{ne^2} \left[\frac{\partial \mathbf{j}}{\partial t} + \nabla \cdot (\mathbf{v} \mathbf{j} + \mathbf{j} \mathbf{v}) \right] + \frac{\mathbf{j} \times \mathbf{B}}{nec} - \frac{\nabla \cdot \mathbf{p}_e}{ne}, \quad (3.37)$$

where $\mathbf{v} \mathbf{j}$ and $\mathbf{j} \mathbf{v}$ are dyadic tensors [(with components $v_n j_m$ and $v_m j_n$)] and \mathbf{p}_e is the electron stress tensor. Term (i) on the right-hand side of this equation is the convective electric field, while the term (ii) is the field associated with Ohmic dissipation caused by electron-ion collisions. The conductivity, σ_e , is the inverse of the electrical resistivity, η . The next group of terms (iii)

³⁵ Activity: In solar physics, flux tubes are commonly used as an approximation of the state of the magnetic field in near-photospheric layers: embedded in a field-free atmosphere is a bundle of field separated from its surroundings by a thin current envelope. Assuming an ideal plasma without flows, show that the atmosphere within the tube is in hydrostatic equilibrium regardless of the path of the flux tube through the atmosphere. Show how pressure balance (incorporating both gas and magnetic components) determines the cross section of the tube.

describes the effects of electron inertia [(which is ignored in Eq. (3.11) as another approximation of Ohm's law, while the latter describes also simplifies pressure by assuming it to be isotropic). The next term, (iv),] is the Hall effect. Ion inertia is negligible because the large mass of the ions means that they do not contribute significantly to a change in the current density. Finally, the term (v) includes the electron gyro-viscosity, which is considered by many to be important at [any point where the magnetic field vanishes, *i.e.*, at a] magnetic null. For a partially ionized plasma, collisions between charged particles and neutrals lead to additional terms associated with ambipolar diffusion.

Although all of the terms on the right-hand side of the generalized Ohm's law, other than the first, allow field lines to slip through the plasma, they do not all produce dissipation. For example, the inertial terms in (iii) do not cause the entropy of the plasma to increase. Thus, even though one may speak of inertial effects as creating an effective resistivity, this resistivity does not necessarily lead to dissipation.

Which terms are important in a particular situation depends not only upon the plasma parameters, but also upon the length and time scales for variations of these parameters. For magnetic reconnection, we normally want to know which non-ideal terms are likely to be significant within the current sheet where the frozen-flux condition is violated. Because each non-ideal (*i.e.*, diffusion) term in the generalized Ohm's law contains either a spatial or temporal gradient, we can estimate the significance of any particular term by computing the gradient scale-length, L_t , required to make the term as large as the value of the convective electric field, $\frac{1}{c}\mathbf{v} \times \mathbf{B}$, outside the diffusion region.

Consider, for example, the three inertial components of term (iii). If we assume that $\nabla \approx 1/L_t$, $|\mathbf{j}| \approx (c/4\pi)B_t/L_t$ and $\partial/\partial t \approx v_t/L_t$, say, where L_t is a typical length-scale and v_t a typical velocity, then these three components of (iii) will be of the same order as the convective electric field if

$$\frac{cm_e}{4\pi ne^2} \frac{v_t B_t}{L_t^2} \approx \frac{V_t B_t}{c}. \quad (3.38)$$

In other words, in order for the inertial terms to be important in a current sheet, its thickness ℓ_{inertia} should be

$$\ell_{\text{inertia}} \approx \left(\frac{c^2 m_e}{4\pi n e^2} \right)^{1/2} \approx \lambda_e, \quad (3.39)$$

where

$$\lambda_e = \frac{c}{\omega_{pe}} = \left(\frac{c^2 m_e}{4\pi n e^2} \right)^{1/2} = 5.3 \times 10^5 n^{-1/2} \quad (3.40)$$

is the electron-inertial length or skin-depth [(which characterizes the depth in

a plasma into which electromagnetic radiation can penetrate)], c is the speed of light and

$$\omega_{pe} = (4\pi n e^2 / m_e)^{1/2} = 5.6 \cdot 10^4 n_e^{1/2} \quad (3.41)$$

is the electron plasma frequency.

Similarly, for the Hall term (iv)

$$\frac{B_t^2}{4\pi n e L_t} \approx \frac{V_t B_t}{c} \quad (3.42)$$

or

$$\ell_{\text{Hall}} \approx \frac{c}{M_A} \left(\frac{\mu m_p}{4\pi n e^2} \right)^{1/2} \approx \frac{\lambda_i}{M_A}, \quad (3.43)$$

where

$$\lambda_i = \frac{c}{\omega_{pi}} = \left(\frac{\mu c^2 m_p}{4\pi n e^2} \right)^{1/2} = 2.3 \times 10^7 \left(\frac{\mu}{n} \right)^{1/2} \quad (3.44)$$

is the ion-inertial length or skin-depth [(below which ions decouple from electrons, and the magnetic field may no longer be frozen into the plasma overall but instead into the electron fluid), and $\mu = m_i/m_p$]. The Alfvén Mach number equals [$M_A = V_t/v_A$,] and $\omega_{pi} = (4\pi n e^2/m_i)^{1/2}$ the ion plasma frequency, and v_A the Alfvén speed [(see Eq. 3.30 and Table 2.5)].

For the electron-stress term (v) we can write

$$\frac{n k T_t}{n e L_t} \approx \frac{V_t B_t}{c} \quad (3.45)$$

if we assume $|\mathbf{p}_e| \approx n k T_e$ and $T_e \approx T_i \approx T$. Solving for L_t leads to

$$\ell_{\text{stress}} \approx \frac{\beta^{1/2}}{M_A} r_{gi}, \quad (3.46)$$

where [the plasma β is given by Eq. (3.24) and the ion-gyro radius for the average thermal velocity (v_{Ti}) equals $r_{gi} = (k T m_i)^{1/2} c / e B$; see also Table 2.5.]

Finally, for the collision term (ii), \mathbf{j}/σ_e ,

$$\frac{c B_t}{4\pi \sigma_e L_t} \approx \frac{M_A v_A B_t}{c}. \quad (3.47)$$

[where σ_e^{-1} is] also the magnetic diffusivity, η . Using Spitzer's formula for the collisional resistivity, η , of a plasma (see [H-I:3]) we obtain

$$\eta = \frac{(k m_e T_e)^{1/2}}{n e^2 \lambda_{\text{mfp}}}, \quad (3.48)$$

where

$$\lambda_{\text{mfp}} = \frac{3}{4\pi^{1/2}} \frac{(k T_e)^2}{n e^4 \ln \Lambda} = 1.1 \times 10^5 \frac{T_e^2}{n \ln \Lambda} \quad (3.49)$$

is the mean-free path for electron-ion collisions. Combining these expressions with those for the electron and ion inertial lengths we obtain [an estimate for the length scale below which effects of collisions become important to field diffusion:]

$$\ell_{\text{collisions}} \approx \frac{\beta^{1/2} \lambda_e \lambda_i}{M_A \lambda_{\text{mfp}}}. \quad (3.50)$$

Note that the length-scale, $\ell_{\text{collision}}$, of the spatial variations required to achieve significant field-line diffusion is inversely proportional to the mean-free path, λ_{mfp} . As λ_{mfp} increases, the diffusion caused by collisions becomes less effective, and increasingly sharper gradients are required to maintain the size of the dissipation term, \mathbf{j}/σ_e . {A:[36]}

A:36

[Table 3.4 lists] various plasma parameters along with the characteristic scale-lengths for four different regions where reconnection is thought to occur. The parameter L_s is the global (system-level) scale-size of the region, and the fundamental quantities from which all other parameters are derived are the density n , temperature T , and magnetic field B . For convenience, we assume that the Alfvén Mach number M_A is unity and that the electron and ion temperatures are roughly equal. The most extreme plasma environments listed in Table 3.4 occur in the magnetosphere, which is completely collisionless, and in the solar interior which is highly collisional.

In addition to the parameters discussed above, Table 3.4 also lists the value of the *Debye length* [whose expression is shown in Table 2.5.] The number of particles within a Debye sphere (*i.e.*, $4\pi n \lambda_D^3/3$) must be larger than unity in order for the generalized Ohm’s law to hold. Otherwise, the collective behavior which characterizes a plasma breaks down. The number of particles in a Debye sphere for the environments shown in Table 3.4 ranges from 10^{14} for the magnetosphere to only about four for the solar interior at the base of the convection zone. Also shown in the table is the Lundquist number, L_u , which is the same as the magnetic Reynolds number, \mathcal{R}_m [introduced in Eq. (3.18)], when the flow and Alfvén speeds are the same. For a collisional plasma the Lundquist number based on L_s can be expressed as

$$L_u = \frac{v_A}{v_d} = \frac{L_s v_A}{\eta} = 2. \times 10^8 \frac{L_s T_e^{3/2} B_t}{(\mu n)^{1/2} \ln(\Lambda)} \quad (3.51)$$

[...] In the expression on the right, η has been replaced by Spitzer’s formula for the electrical resistivity of collisional plasma.

The characteristic scale-lengths in Table 3.4 provide an indication of which

³⁶ Activity: Units: this text uses cgs-Gaussian units. In other texts (including many of the Heliophysics chapters) you will find SI units. Look into conversions from one system to another (for example with the online NRL Plasma Formulary).

Table 3.4. *Comparison of order-of-magnitude plasma parameters in different environments (cgs-Gaussian units – i.e., length scales in cm, n in cm^{-3} , T in K, B in Gauss, electric fields in statV cm^{-1}). [Modified after H-I:5, merging two tables in SI units]*

Parameter	Laboratory experiment ¹	Terrestrial magnetosphere ²	Solar corona ³	Solar interior ⁴
region scale L_s	10	10^9	10^{10}	10^9
density n_t	10^{14}	10^{-1}	10^9	10^{23}
temperature T_t	10^5	10^7	10^6	10^6
field strength B_t	10^3	10^{-4}	10^2	10^5
Debye length λ_D	10^{-4}	10^5	10^{-1}	10^{-8}
ion gyro radius r_{gi}	10^{-1}	10^7	10	10^{-2}
ion inertial length λ_i	1	10^8	10^3	10^{-4}
Coulomb logarithm $\ln(\Lambda)$	11	33	19	3
coll. mean-free path λ_{mfp} ⁵	1	10^{18}	10^6	10^{-7}
$\ell_{inertia}(\lambda_e)$	10^{-2}	10^6	10	10^{-6}
$\ell_{Hall}(\lambda_i)$	1	10^8	10^3	10^{-4}
ℓ_{stress}	10^{-1}	10^7	10^{-1}	1
$\ell_{collision}$	10^{-2}	10^{-5}	10^{-5}	10^{-1}
plasma β	10^{-2}	10^{-1}	10^{-4}	10^4
Lundquist no. $L_u(\approx \mathcal{R}_m)$	10^3	10^{14}	10^{14}	10^{10}
Dreicer field E_D	10^{-1}	10^{-17}	10^{-6}	10^7
$E_A(= v_A B/c)$	1	10^{-6}	10^1	1
$E_{SP}(= E_A/\sqrt{\mathcal{R}_m})$	10^{-2}	10^{-13}	10^{-7}	10^{-6}

¹ The Magnetic Reconnection eXperiment (MRX) at Princeton Plasma Physics Laboratory; ² plasma sheet; ³ above a solar active region; ⁴ at the base of the solar convection envelope [at a depth of about 200,000 km around which many consider primary dynamo mechanisms to operate; ⁵ note that this is a purely collisional mean-free path, ignoring other couplings that may occur through the magnetic field].

terms in the generalized Ohm's law of Eq. (3.37) are likely to be important for reconnecting current sheets. As with MHD shocks and turbulence, the large-scale dynamics of the flow cause the current sheet to thin until it reaches a length-scale where field-line diffusion is effective. Thus, in principle, the term with the largest characteristic length-scale in Table 3.4 is the one that will be most important. Because the Hall term has the largest length in every environment except the solar interior, one might conclude that it is generally the most important. However, this conclusion does not take into consideration the fact that the Hall term tends to zero in the region of a magnetic null point

or sheet. The Hall term on its own does not contribute directly to reconnection, since it freezes the magnetic field to the electron flow. [...] An excessively small scale does indicate that any process associated with that term is unlikely to be important. Therefore, on this basis, we can conclude that collisional diffusion is not important in the terrestrial magnetosphere or the solar corona, and that the electron-inertial terms and the Hall term are not important in the solar interior. [...] On the other hand, if a term is not associated with an obviously 'excessively small' scale, it is difficult to know whether a particular term is really as important as suggested by its relative length scale; evaluating such cases] requires a complete analysis of the kinetic dynamics, which is a rather formidable task.

Although the collision length-scale, $\ell_{\text{collision}}$, is equally small in both the magnetosphere and the corona, the general importance of collisions for these two regions is quite different. In the magnetosphere the collision-mean-free path, λ_{mfp} , is nine orders of magnitude larger than the global scale-size, L_s , but in the corona it is four orders of magnitude smaller than the global scale. Thus, we can be confident that collisional transport theory applies to large-scale structures in the corona even though it is not applicable within thin current sheets or dissipation layers. By contrast, in the magnetosphere, collisions are so few that collisional transport theory does not apply at any scale.

Another important issue concerning the applicability of collisional theory is the strength of the electric field in a frame moving with the plasma. If this field exceeds the *Dreicer electric field* defined by

$$E_D = \frac{e \ln(\Lambda)}{\lambda_D^2} = \frac{4\pi e^3}{k} \frac{\ln(\Lambda) n}{T_e} = 10^{-11} \frac{n \ln(\Lambda)}{T_e}, \quad (3.52)$$

runaway acceleration of electrons will occur. The most likely location for the production of runaway electrons in a reconnection process is in a thin current sheet that forms at the null point. This field could be as large as the convective electric field based on the Alfvén speed, that is

$$E_A = \frac{1}{c} v_A B_t = 7.2 \frac{B_t^2}{(\mu n)^{1/2}}, \quad (3.53)$$

or as low as the Sweet-Parker electric field

$$E_{\text{SP}} = \frac{E_A}{\mathcal{R}_m^{1/2}}, \quad (3.54)$$

where \mathcal{R}_m is the magnetic Reynolds number based on the inflow Alfvén speed (*i.e.*, the inflow Lundquist number). As shown in Table 3.4, the Dreicer field in the magnetosphere is much smaller than E_A or E_{SP} , so runaway electrons will always be generated by reconnection there. On the other hand, in the solar

interior the Dreicer field is so large that runaway electrons never occur. In the intermediate regimes of the laboratory and the solar corona, the Dreicer field lines between E_A and E_{SP} , so perhaps runaway electrons are only produced when very fast reconnection occurs.

Even in completely collisionless environments like the Earth's magnetosphere, it is still sometimes possible to express the relation between electric field and current density in terms of an anomalous resistivity. For example, [...] the electron inertial terms (iii) in the generalized Ohm's law of Eq. (3.37) lead to an anomalous resistivity

$$\frac{1}{\sigma_e^*} = \frac{\pi B_\perp}{2ne}, \quad (3.55)$$

where B_\perp is the field normal to the current sheet. This resistivity is derived solely from a consideration of the particle orbits, and in the magnetotail current sheet it may be larger than any anomalous resistivity due to wave-particle interactions. A typical example of the latter is the anomalous resistivity due to ion-acoustic waves".

For some discussion of reconnection in two and three dimensions in the Heliophysics books, see H-I:5.3 and H-I:5.4. More on the effects of reconnection follows in Chs. 6 and 8.

3.5 A few notes about conditions

3.5.1 Solar atmosphere vs. terrestrial magnetosphere

The scale lengths estimated for the importance of terms in Ohm's law in Table 3.4 are very much smaller than the scale of the corona itself and even compared to any active region, but importantly also very much smaller than the angular resolution achievable by imaging instruments (currently about 1 arcsec or ~ 700 km for space-based EUV imagers). Consequently, the scale on which reconnection occurs in the corona is not observed, while the effects of such reconnection become apparent in the magnetic geometry and plasma atmospheres on scales well above the reconnection itself.

In contrast, in the terrestrial magnetosphere all but the length scale, $\ell_{\text{collision}}$, on which collisional effects could contribute significantly as a term in the generalized Ohm's law are large enough that spacecraft can scan reconnection regions as they fly through, while constellations of spacecraft can probe reconnection in multiple dimensions.

Another significant distinction is that in the terrestrial magnetosphere the ion-gyro radius (particularly for relatively heavy and energetic particles) is not small compared to the scale on which these particles probe the magnetic field. This is an important cause behind what is known as the ring current

(see Sect. 8.1) that is largely carried precisely by such particles. For solar conditions, in contrast, such effects of particle gyration are not directly evident on any observable scale.

3.5.2 Heliosphere

[H-I:11.3] “Adopting typical solar wind values near Earth of $n_t = 5$ particles/cm³ for density, $v_t = 400$ km/s solar wind speed and $B_t = 50 \mu\text{G}$ magnetic field strength (values consistent with Table [2.4]) we can evaluate the expected energy density of the solar wind, which can be broken down into three components: flow, magnetic and thermal. [...] The flow energy density is estimated to be

$$e_{v,\odot} \equiv \left(\frac{1}{2} \rho v^2 \right)_{\text{sw}} \approx \frac{1}{2} m_p n v_{\text{sw}}^2 \approx 7 \times 10^{-9} \left(\frac{n(\text{cm}^{-3})}{5} \right) \left(\frac{v(\text{km/s})}{400} \right)^2 \text{ erg/cm}^3. \quad (3.56)$$

The energy density of the solar wind’s magnetic field is

$$e_{B,\odot} \equiv \left\langle \frac{B^2}{2} \right\rangle \approx 1.0 \times 10^{-10} \left(\frac{B(\mu\text{G})}{50} \right)^2 \approx 0.015 e_{v,\odot} \text{ erg/cm}^3, \quad (3.57)$$

while the thermal energy density using values from Table [2.4] is

$$\begin{aligned} e_{T,\odot} \equiv \left\langle \frac{3}{2} n k (T_e + T_i) \right\rangle &\approx 2.5 \times 10^{-10} \left(\frac{n(\text{cm}^{-3})}{6} \right) \left(\frac{T_e(\text{K})}{1.2 \times 10^5} + \frac{T_i(\text{K})}{1.4 \times 10^5} \right) \\ &\approx 0.03 e_{v,\odot} \text{ erg/cm}^3 \end{aligned} \quad (3.58)$$

where $T_{i,e}$ are the solar wind ion and electron temperatures. Taking the [values from Table 2.4], the above [These] estimates show that the bulk of the energy in the solar wind at Earth is in the flow:” $e_{v,\odot} \sim 30 e_{B,\odot} \sim 70 e_{T,\odot}$ {A:[37]} A:37

³⁷ Activity: Make comparisons of energy densities for the solar wind as in Sec. 3.5.2 at other bodies in the Solar System (using Table 5.2). Why comparisons of energy densities in planetary magnetic fields (Table 5.2) and in the surrounding solar wind are informative is discussed in Ch. 5. Why would you expect the flow energy density and the magnetic field energy density to be comparable at only a few solar radii from the Sun?

4

Dynamos of Sun-like stars and Earth-like planets

4.1 Dynamo settings

Stellar and planetary dynamos thrive wherever sufficiently vigorous flows of a conducting medium transport substantial thermal energy in an adequately spinning body. The energy transported has to come from a reservoir that may date back to the formation of the body (in planets or very young stars) or may have its origin in nuclear fusion (in stars) or in solidification – the latter often accompanied by chemical separation – or nuclear fission (in terrestrial planets). The flows that transport the energy may be dominated by Coriolis forces (in planets where flows are slow compared to spin rates) or by stratification (including chemical gradients in planets, while in stars pressure gradients of the compressible medium limit how far matter can efficiently rise before overturning). The amount of energy transported is regulated by the source in the deep interior as well as by the sink at the top of the dynamo region. In Sun-like stars that sink is the stellar surface, and the properties of radiative transfer through these surface layers are important in determining the internal structure of the entire star as it balances the energy produced by nuclear fusion with its luminosity. In a planet like Earth, the energy transport in the dynamo region of the core is determined to a large extent by the convective motions in the enveloping mantle that transport heat to where it is ultimately lost through the surface.

[H-I:3.3] “The formal difference between the type of dynamos that we are interested in here and the self-excited dynamos in power plants is the homogeneous distribution of conductivity (that would lead to a short-circuit situation) that does not put any constraints on electric currents (electric wires could be considered as special cases of a highly inhomogeneous conductivity distribution). For this reason these dynamos are also called homogeneous fluid dynamos.”

In stars, [H-III:5.1] “[t]hermonuclear fusion in their cores converts matter into

thermal energy and electromagnetic radiation which, in the Sun, is transported outward via the diffusion of photons. In the solar envelope, the plasma becomes more opaque as the temperature drops, which inhibits radiative diffusion and steepens the temperature gradient relative to the adiabatic temperature gradient. The stratification soon becomes superadiabatic and thermal convection [gradually] takes over as the primary mechanism for transporting energy to the solar photosphere where it is radiated into space. {A:[38]} [All stars A:38 with a mass of somewhat above that of the Sun or less than that have such a convective envelope during their 'main-sequence' (equilibrated hydrogen-fusing) phase (see Figs. 4.1 and 4.2); the least massive stars are fully convective. All of these stars power a dynamo during the longest-lived mature phase, and all stars do during their initial birth phases and in the last phases of their lives, both of which are short compared to the mature phase (Ch. 10). Stars cool enough to have a convective envelope reaching into their surface layers are known as 'cool stars'. [vii] {A:[39]} A:39

The solar convection zone occupies approximately the outer 30% of the Sun by radius. It is here where [a small fraction of the] internal energy of the plasma is converted to kinetic energy and then [a small fraction of that] to magnetic energy, aided by radiation and gravity. Radiative heating [of the bottom

³⁸ Activity: The transition from radiative diffusion to convective enthalpy transport at the bottom of the convective envelope is gradual: the fraction of total energy carried as a diffusive flux gradually drops while that of the enthalpy flux smoothly increases, making convection the dominant transport about 35,000 km above the bottom of the convective envelope, or roughly after a single pressure scale height (see Sect. 4.3). Can you think of other terms that would be involved in the energy transport equation in a stellar convective envelope? A fair idea of the answer, along with a quantitative comparison of the relative importance of the processes involved in carrying energy through the convective envelope, can be found, for example, in this analysis by Brun *et al.* (2004), in particular their Fig. 3 (note that they show transport by convection that is resolved by their model and by (parameterized) unresolved - 'subgrid-scale' - convection).

³⁹ Activity: Figure 4.2 is a brightness-color diagram (known as a Hertzsprung-Russell, or HR, diagram) using typical astronomical units: absolute visual magnitude M_V , which is a logarithmic measure of stellar brightness, and spectral color $B - V$, which is the logarithm of the ratio of two brightnesses measured in different color bands (often using logarithmic brightness B and V , or less commonly R for blue, visual, and red). The table in that figure maps spectral type (see footnote vii), $B - V$, effective temperature T_{eff} and a correction factor BC that relates visual and bolometric brightnesses (see equations below that table). Using this information, estimate stellar radii R_* for Sirius A, ϵ Eri, 61 Cyg A, and AD Leo, realizing that $L_* = (\sigma T_{\text{eff}}^4)(4\pi R_*^2)$, with the Stefan-Boltzmann constant $\sigma = 5.7 \times 10^{-5} \text{ erg/cm}^2/\text{sec}/\text{deg}^4$. Sketch a double-logarithmic $L - T_{\text{eff}}$ version of the HR diagram and draw lines of constant radius in it. Then compare that to Fig. 10.1.

vii Astronomers characterize the properties of stars based on their spectrum. The overall shape gives an indication of the surface temperature, while details of spectral lines (generally in absorption, but some in emission) provide finer detail used in classification schemes. One such scheme frequently used is that of 'spectral type' in the Morgan-Keenan (MK) scheme: only after the classes were introduced was a monotonic mapping to temperature established, going from hot to cooler: O, B, A, F, G, K, M, L , and T (with the last two fairly recent additions for very cool, very faint stars, with T reaching the domain of 'brown dwarfs'). The letter is followed by a subclass from 0 to 9, and commonly an indicator of 'luminosity class': a roman numeral indicative of the size of the star: I, II, III, IV, and V for supergiants, bright giants, giants, subgiants, and main-sequence or dwarf stars. The term 'main sequence' refers to a band in brightness-color diagrams, such as Fig. 4.2, within which stars spend most of their lives, as long as they are steadily fusing hydrogen into helium.

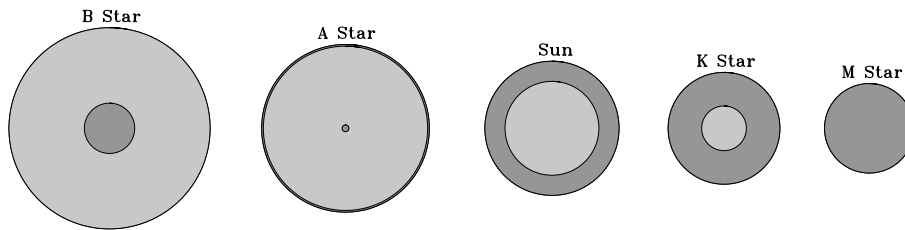


Fig. 4.1. Schematic representation of the radiative (light grey) and convective (dark grey) internal structure of main-sequence stars. The thickness of the outer convection zone for the A-star is here greatly exaggerated; drawn to scale it would be thinner than the black circle delineating the stellar surface on this drawing. Relative stellar sizes are also not to scale: a B0 V star has a radius of $\sim 7.5 R_{\odot}$, and an M0 V star has a radius of $\sim 0.6 R_{\odot}$, i.e., 12 times smaller. [Fig. H-III:2.10]

of the] convection zone and radiative cooling in the photosphere maintain a superadiabatic temperature gradient that sustains convective motions by means of buoyancy. In a rotating star, convection transports momentum as well as energy, establishing shearing flows and global circulations. These mean flows work together with turbulent convection to amplify and organize magnetic fields through hydromagnetic dynamo action, giving rise to the rich display of magnetic activity so striking in modern solar observations.”

The Sun’s large scale magnetic field exhibits a quasi-periodic modulation on a roughly 11-year basis during which the level of magnetic activity waxes and wanes as a pattern of activity migrates from mid to low latitudes, then to pick up again at higher latitudes, with some temporal overlap in the early and late phases of these cycles. For stars like the Sun, the mean level of activity as expressed by the surface-averaged absolute magnetic flux density ranges over more than three orders of magnitude, depending on the stellar rotation rate, age, and internal structure (more on that in Sect. 9.3; see also H-III:2). [H-III:6.1] “[T]he *existence* of solar and stellar magnetic fields is in itself not really surprising; any large-scale fossil field present at the time of stellar formation would still be there today at almost its initial strength, because the Ohmic dissipation timescale is extremely large for most astrophysical objects [(Eq. 3.20)]. The challenge is instead to reproduce the various observed spatiotemporal patterns [...], most notably the cyclic polarity reversals on decadal timescales.”

As to planetary dynamos, [H-III:7.1] “[s]pace missions revealed that most planets in the Solar System have internal magnetic fields (see Ch. H-I:13), but there are exceptions (Venus, Mars). Some planets seem to have had a field that is now extinguished (*e.g.*, Mars). In many cases with an active dynamo

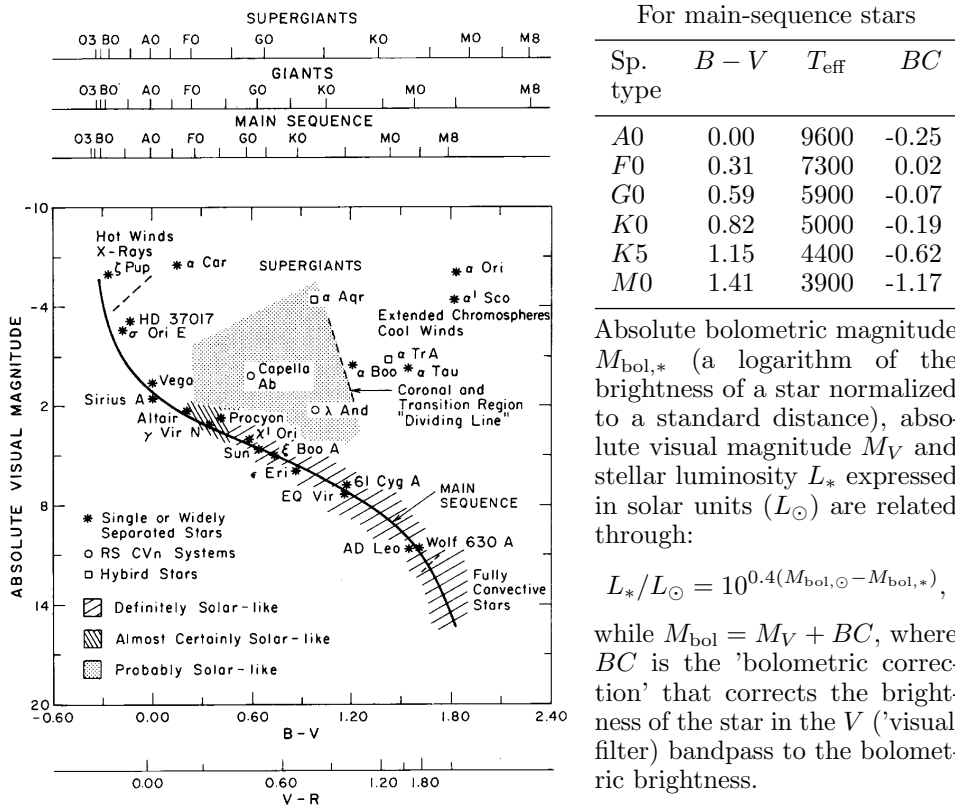


Fig. 4.2. A Hertzsprung-Russell diagram showing stars with substantial magnetic activity in shaded or hatched domains, which are distinguished in groups of solar-likeness as indicated in the legend. The main sequence where stars spend most of their lifetime fusing hydrogen into helium in their cores is indicated by a solid curve; well above that lies the domain of the supergiant stars, with the giant star domain in between. Also indicated is the region where massive winds occur and where hot coronal plasma appears to be absent. Some frequently studied stars (both magnetically active and nonactive) are identified by name. The axes above the main panel show the spectral types (see footnote vii) for supergiant, giant, and main-sequence stars for the corresponding spectral color index $B - V$ or corresponding $V - R$ index. [Fig. H-III:2.8, with an added information panel on the right; figure source: Linsky (1985).]

the axial dipole dominates the field at the planetary surface, but Uranus and Neptune are exceptions. Saturn is special because its field is extremely symmetric with respect to the planet's rotation axis. The field strengths at the planetary surfaces differ by a factor of 1000 between Mercury and Jupiter [(cf. Table 5.3)]. A full understanding of this diversity in the morphology and strength of planetary magnetic fields is still lacking, but a number of promising ideas have been suggested and backed up by dynamo simulations. Some of the differences can be explained by a systematic dependence of the dynamo

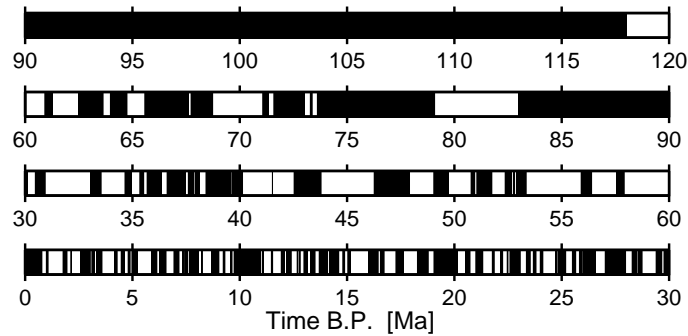


Fig. 4.3. Polarity of the geomagnetic field for the past 120 million years, with time running backward from left to right in each row (before present - B.P., i.e., 1950 - in units of millions of years). Dark regions indicate times when the dipole polarity was the same as today, in white regions it has been opposite. [Fig. H-III:7.4]

behavior on parameters such as rotation rate or energy flux, whereas others seem to require qualitative differences in the structure and dynamics of the planetary dynamo.”

[H-III:7.4.1] “Earth serves as the prototype for the terrestrial planets. [...] There is a core with radius $R_{\text{core}} \approx 0.55R_{\text{planet}}$, [the outer part of which is liquid]. The small inner core, with a radius $0.35R_{\text{core}}$, is [solid]. The core appears to consist predominantly of iron. [...]

The total internal heat flow at the Earth’s surface is $4.6 \cdot 10^{20}$ erg/s (although a large number, it is only 0.03% of the total power coming into the Earth’s atmosphere by insolation). Roughly one half of it is balanced by the heat generated by the decay of uranium, thorium and the potassium isotope ^{40}K inside the Earth. The remainder of the heat flow is due to the cooling of the Earth. The loss of gravitational potential energy associated with the contraction of the Earth contributes a modest amount, but is much less important than it is in young stars or in gas planets. How much of the Earth’s heat flow comes from the core is rather uncertain. Recent estimates that are based on different lines of evidence mostly fall into the range $(0.5 - 1.5) \cdot 10^{20}$ erg/s, although values as low as $(0.3 - 0.4) \cdot 10^{20}$ erg/s have also been discussed. Most of the radioactive elements reside in the silicate crust and mantle. Some amount of potassium may be present in the core, but the majority of the core heat must be due to cooling. It is important to note that the heat loss from the core is regulated by the slow solid-state convection in the mantle. The core, which convects vigorously in comparison to the mantle and which is thermally well-mixed, delivers as much heat as the mantle is able to carry away.”

4.1.1 Earth and other terrestrial planets

Solar System bodies that have a present-day active dynamo include Mercury and Earth among the terrestrial planets, the jovian moon Ganymede, and all the giant planets; see Table 5.3 for their global properties.

[H-I:3.1.1] “The surface magnetic field of the Earth has a strength of about 0.5 G with mainly dipolar character. [The dipole axis is tilted by a variable amount over time with respect to the axis of rotation, such that the magnetic north pole has wandered from as far south as about 70 degrees in geographic latitude to within a few degrees from the geographic north pole over the past two centuries]. From studies of rock magnetism (when rocks cool below the Curie point they preserve the magnetic field that was present in them at that time) it is known that the Earth had a magnetic field over the past 3.5×10^9 years and that the strength and orientation of the field varied significantly on time scales of 10^3 to 10^4 years. A given polarity typically dominates for about 200 000 years with quick reversals on a time scale of a few thousand years in between [(see Fig. 4.3)]. While the orientation of the axis of the dipole changes significantly with time, the dipole moment is aligned with the axis of rotation when averaged over $\sim 10^4$ years.”

In contrast to the case of cool stars, [H-III:7.4.1] “[r]adiative heat transfer is not an issue in planetary cores, but liquid metal is a good thermal conductor. The heat flux that can be transported by conduction along an adiabatic temperature gradient, $(dT/dr)_{ad} = T/H_T$, is sometimes called the ‘adiabatic heat flow’ (T is absolute temperature, $H_T = c_p/(\zeta g)$ is the temperature scale height with c_p the heat capacity, ζ the thermal expansivity and g the gravitational acceleration). In terrestrial planets, the adiabatic heat flow can be a large fraction of the actual heat flow, or it may exceed the actual heat flow, in which case at least the top layers of the core would be thermally stable. Near the top of Earth’s core approximately $(0.3 - 0.4) 10^{20}$ erg/s can be conducted along the adiabat, *i.e.*, close to the minimum estimates for the entire core heat flow. But even if all the heat flux near the core-mantle boundary were carried by conduction, a convective dynamo can exist thanks to the inner core. At the inner core boundary, the adiabatic temperature profile of the convecting outer core crosses the melting point of iron. The latter increases with pressure more steeply than the adiabatic gradient, which is the reason why the Earth’s core freezes from the center rather than from above. As the core cools, the inner core grows with time by freezing iron onto its outer boundary. This has two important implications for driving the dynamo. The latent heat that is released upon solidification is an effective heat source, which contributes to the heat budget approximately the same amount as the bulk cooling of the core. [...] A second, perhaps more important effect is that the light elements

in the outer core are preferentially rejected when iron freezes onto the inner core. Hence, they become concentrated in the residual fluid near the inner core boundary. This layering is gravitationally unstable because of the reduced density, which leads to compositional convection that homogenizes the light elements in the bulk of the fluid core. Compositional convection contributes as much as, or more than, thermal convection to the driving of the geodynamo in recent geological times.

Most predictions for the inner core growth rate imply that the inner core did not exist for most of the history of the Earth. Rather, it would have nucleated between 0.5 and 2 billion years ago. In the absence of an inner core, only thermal convection by secular cooling of the fluid core (and perhaps radioactive heating) can drive a dynamo, which is less efficient than the present-day setting. A change in the geomagnetic field properties might be expected upon the nucleation of the inner core, but no clear indication for such an event has been found in the paleomagnetic record.”

[H-III:7.4.2] “No direct evidence on the existence or non-existence of a solid inner core is available for any planet other than Earth. But the possible absence of an inner core could explain why Venus and Mars do not have an active dynamo. On Earth, mantle convection reaches the surface in the form of plate tectonics, which is a fairly efficient mode of removing heat from the interior. None of the other terrestrial planets have plate tectonics. In their cases, mantle convection is confined to the region below the lithosphere, a rigid lid of some 100 – 300 km thickness through which heat must be transported by conduction. Without plate tectonics, the heat flow is expected to be significantly lower not only at the surface, but also at the top of the core, where it is very probably subadiabatic. If no inner core exists to provide latent heat, it is then subadiabatic throughout the core. Furthermore, compositional convection is also unavailable to drive a dynamo. The slower cooling of the planetary interior in the absence of plate tectonics concurs with the idea that an inner core has not (yet) nucleated in the cases of Mars and Venus. Early in the planets’ history the cooling rate was probably much higher and the associated core heat flow large enough for thermal convection. The demise of the dynamo must have occurred when the declining heat flow dropped below the conductive threshold.”

For discussion of dynamos in non-terrestrial planets, see Ch. H-III:7.

4.1.2 The Sun and other stars

[H-I:3.1.2] “The Sun shows magnetic field on all observable scales [(Fig. 4.4)] with a significant range in field strength, from individual sunspots with magnetic

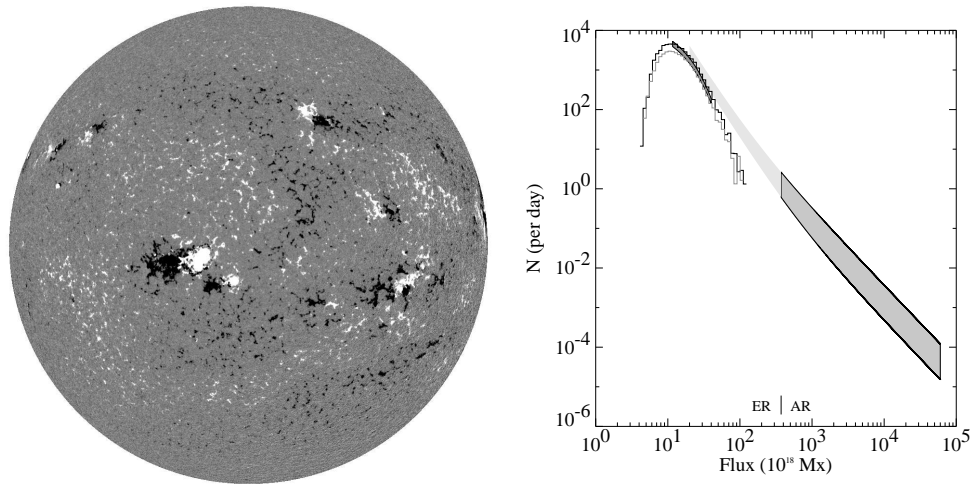


Fig. 4.4. Left: First solar magnetic map (magnetogram) of the current millennium, taken by SOHO's MDI on 2001/01/01 00:03 UT. The magnetogram (with white/black for negative/positive line of sight polarity) shows a variety of active regions, embedded in patches of largely unipolar enhanced supergranular network, mixed-polarity quiet Sun regions, and low-flux polar caps (weak at this near-maximum phase of the cycle, and weakened further in the line-of-sight flux map because of projection effects on the near-vertical magnetic field). Right: Distribution function of emerging magnetic bipolar regions on the Sun, showing the emergence frequency per day per flux interval of 10^{18} Mx, estimated for the entire solar surface. The shaded region on the right envelopes the range of the active-region spectrum for solar cycle 22 (for half-year intervals around sunspot minimum and maximum). The histograms on the left are for the ephemeral regions; the shaded band shows where observations are least affected by spatial (lower cutoff) and temporal (upper cutoff) biases. The spectrum for regions below $\sim 10^{19}$ Mx has yet to be determined; the cutoff here is caused by the limited resolution of the SOHO/MDI magnetograph. [Fig. H-III:2.1]

field strengths of 2 500 to 3 000 G to the average field strength of the global field of only a few Gauss. {A:[40]} {A:[41]}

A:40

The most prominent feature of solar magnetism is the 11-year sunspot cycle (if one considers the field reversals the full period is 22 years), which is reflected in the changing number of sunspots appearing on the surface of the Sun. In the

A:41

⁴⁰ Activity: How is the Sun's magnetic field observed? Look up the effects on photons propagating through a plasma threaded by a magnetic field. This results in the 'Zeeman effect' of line splitting and of circular and linear polarization. For relatively weak field or relatively low wavelengths, the Zeeman splitting of 'magnetically sensitive' spectral lines is generally less than the thermal line width (and less than the Doppler width for rapidly rotating stars), so that what is in principle line splitting for individual atoms becomes line broadening when averaging over populations of atoms and over entire stellar disks.

⁴¹ Activity: Compare a series of solar magnetograms over the past ~ 22 years (using, e.g., SOHO/MDI and SDO/HMI observations). How do the magnetic patterns change over time in terms of overall activity, latitudinal distribution, polarity patterns on the northern versus southern hemisphere, ...?

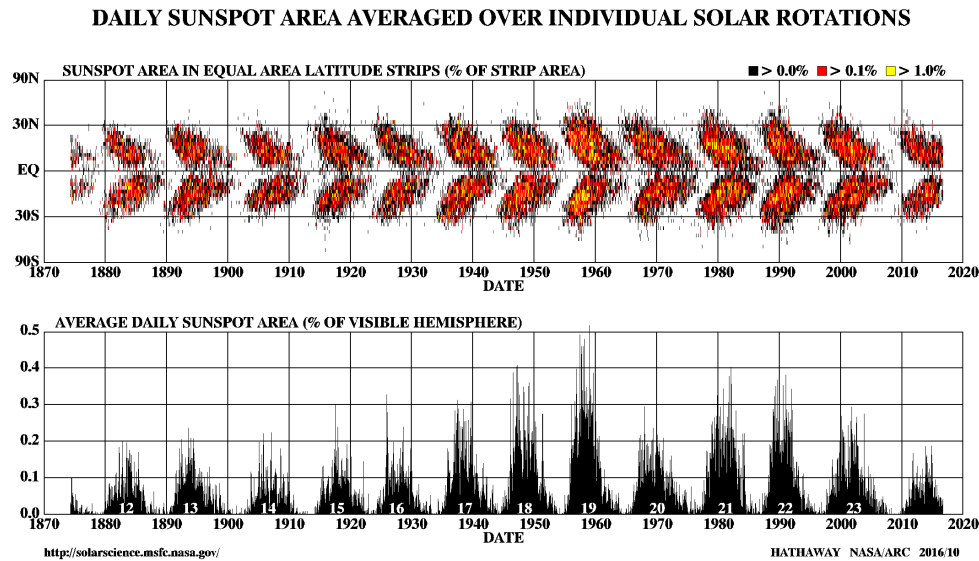


Fig. 4.5. ‘Butterfly diagram’ showing sunspot latitudes (top) and total fractional area coverage (bottom) as a function time. [updated with data through 2018]

beginning of a cycle spots appear at latitudes of about 35° , while close to the end they appear almost at the equator. This property is commonly summarized in the so-called solar butterfly diagram [(Fig. 4.5)]. During the epoch of minimum, the large-scale field of the Sun is most dipolar; the reversal of the poles takes place during solar maximum. On a longer time scale the magnetic activity changes significantly in amplitude and is interrupted by epochs of 100 – 200 years in duration where sunspots are [infrequent or] completely absent [(such as during the Maunder Minimum, about 1645–1715). . . .] Observations of the stellar luminosity or of chromospheric (UV/optical) and coronal (X-ray) emission show that a majority of solar-like stars are magnetically active and around a third to a half show cyclic activity with periods in the range from 3 to 30 years.”

4.2 Dynamo principles

[H-IV:6.1] “Dynamo action refers to the conversion of mechanical energy into electromagnetic energy through induction. In [stars and in planets alike], the mechanical energy is supplied by fluid motions in electrically conducting regions inside [these bodies] and the electromagnetic energy produces the observed [...] magnetic fields. A dynamo is referred to as *self-sustaining* if it does

not require any external magnetic field contributions for regeneration (except initially for a starting seed field).

The fundamental equation governing this induction process is known as the *Magnetic Induction Equation* [Eq. (3.3) in Table 3.3; its derivation and its limitations are described in Sect. 3.2.2. That equation is complemented by the requirement that the currents and the driving flows that are associated with the magnetic field are entirely contained within the body, and that the transition to outside the body for the field is smooth (compare H-I:3.3). . . .]

By inspecting the two terms on the right hand side of Eq. (3.3) we see that magnetic field can grow or decay in time through two processes. The first term ① involves interactions of the velocity and magnetic fields through electromagnetic induction and acts as a source/sink term for field generation. The second term ② represents diffusion due to Ohmic dissipation. To ensure magnetic field does not decay away in time, field must be generated as fast as or faster than its diffusion. A necessary condition for self-sustained dynamo action is therefore that the induction term ① be larger than the diffusion term ② in Eq. (3.3). By using characteristic scales for the variables in the Magnetic Induction Equation (*i.e.*, B_t for the magnetic field scale, v_t for the velocity scale and L_t for a length scale) we derive a common measure of the ratio of field generation to field diffusion known as the *magnetic Reynolds Number*: $\mathcal{R}_m \equiv \frac{v_t L_t}{\eta}$, see Eq. (3.18).]

Upon first glance, it seems reasonable that the magnetic Reynolds number must be larger than unity for dynamo action to be possible. However, more rigorous theoretical analyses suggest that the lower bound for \mathcal{R}_m is instead closer to π^2 and planetary numerical dynamo simulations typically find \mathcal{R}_m must be larger than $\sim 20 - 50$ for self-sustained dynamo action to occur. These higher values are due to the complexities in the velocity field morphologies that cannot be captured in the simple estimate given in Eq. (3.18):” after all, it is a big leap from small-scale field generated on the scale of the flow (such as sketched in Fig. 4.6) to a large-scale field. In cool stars, \mathcal{R}_m typically far exceeds critical values for dynamo action because of the large scales and relatively fast flows involved (see Sect. H-III:5.3.2).

A perspective of what actually supplies the energy to power the dynamo is provided by integrating the induction equation Eq. (3.3) over the object’s volume to establish the total energy in the system:

$$\frac{d}{dt} \int_V \frac{B^2}{4\pi} dV = - \oint_{\partial V} \mathbf{S} \cdot \hat{\mathbf{n}} dS - \eta \int_V j^2 dV - \int_V \mathbf{v} \cdot (\mathbf{j} \times \mathbf{B}) dV. \quad (4.1)$$

The first term on the right is the Poynting flux $\mathbf{S} = (1/4\pi)\mathbf{B} \times (\mathbf{v} \times \mathbf{B})$, which is the energy via the electromagnetic field through a surface into or out of the system across the closed boundary surface ∂V (ignorable if the stellar

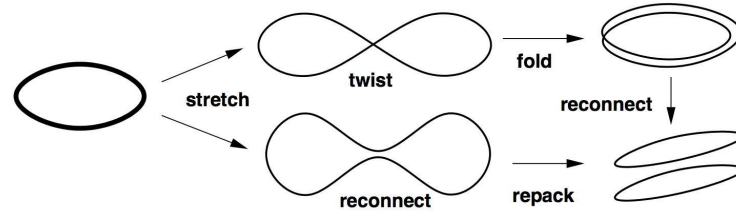


Fig. 4.6. Illustration of two possible flux-rope dynamos. In both cases the field amplification takes place during the stretch operation. The twist-fold (top) and reconnect-repack (bottom) steps are required to remap the amplified flux-rope into the original volume element so that the process can be repeated. Magnetic diffusivity is essential to allow for the topology change required to close the cycle. Each cycle increases the field strength by a factor of 2. [Fig. H-1:3.3]

wind does not take too much power away compared to the total). The second term is the dissipative loss (assuming here that η is uniform). The final term shows that the magnetic energy in the system can be maintained against the dissipative losses only if there are sufficient flows working against – *i.e.*, have an antiparallel component relative to – the Lorentz force $\mathbf{F} = (1/c)\mathbf{j} \times \mathbf{B}$. {A:[42]}

A:42

4.3 Essentials of fluid motions in dynamos

In essence, to drive a large-scale stellar or planetary dynamo, the magnetic field must be subjected to a combination of flow components of a different nature that have their origin in convection and rotation. [H-1:3.3.4] “Fig. 4.6 illustrates the basic ingredients required to amplify a closed magnetic field loop. After a full cycle, the magnetic field strength and the flux have doubled (two loops, each with the original magnetic flux) and the process can be repeated. This very simple illustration points out already a few fundamental properties of a dynamo process. To be able to remap the magnetic field configuration into the original volume element, three-dimensional motions are required. Amplification through stretching is possible in a strictly two-dimensional domain, but there is no way to move the resulting field to return to the right-hand side of the image. The two examples also point out the crucial role of diffusivity in changing the topology of the field. The ‘stretch-twist-fold’ mechanism (excluding diffusive steps) leads to loops of increased complexity, while the ‘stretch-reconnect-repack’ process explicitly involves magnetic diffusivity and ends up with two flux ropes [(see Table 3.1 for a definition)] of similar topology. A reconnection

⁴² Activity: Work through how Eq. (4.1) is obtained by taking the dot product of Eq. (3.3) with \mathbf{B} , integrating over the total volume of the system, and assuming no Poynting flux or currents (or at most only a force-free field) leave the volume. Use vector identities ($\mathbf{a} \cdot (\nabla \times \mathbf{b}) = (\nabla \times \mathbf{a}) \cdot \mathbf{b} - \nabla \cdot (\mathbf{a} \times \mathbf{b})$), $\mathbf{a} \cdot (\mathbf{b} \times \mathbf{c}) = \mathbf{b} \cdot (\mathbf{c} \times \mathbf{a})$, Eq. (3.2), and Gauss’s theorem. For other vector calculus identities, see here.

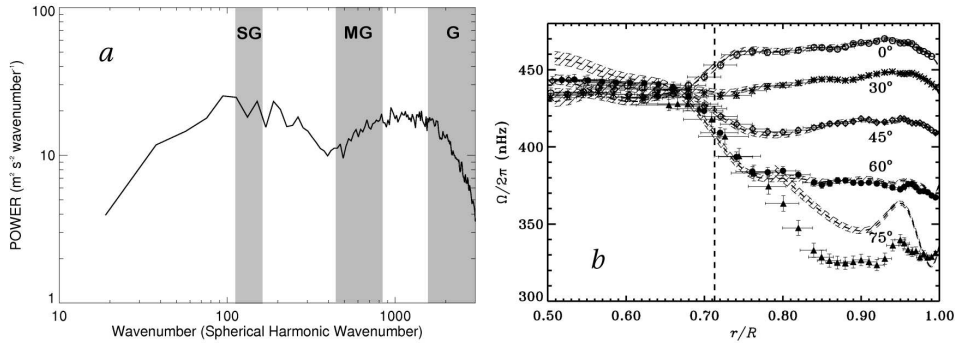


Fig. 4.7. Panel a: Power spectrum of the convective velocity field in the solar photosphere obtained from Doppler measurements, plotted as a function of spherical harmonic degree ℓ . Mean flows and p -modes are filtered out. The falloff beyond $\ell \sim 1500$ reflects the resolution limit of the Michelson Doppler Imager (MDI) instrument onboard the SOHO spacecraft from which these data were obtained and is therefore artificial. Shaded areas indicate the approximate size ranges of supergranulation (SG), mesogranulation (MG) and granulation (G). Note that the expected granulation spectral peak at $\ell \sim 4400$, corresponding to $L \sim 1$ Mm, is not resolved. Panel b: The solar internal rotation profile inferred from helioseismic inversions. Angular velocity $\Omega/2\pi$ is shown as a function of fractional radius r/R_\odot for several latitudes as indicated. Symbols and dashed lines denote different inversion methods, known as subtractive optimally localized averages (SOLA) and regularized least squares (RLS) respectively. Vertical $1-\sigma$ error bars (SOLA) and bands (RLS) are indicated and horizontal bars reflect the resolution of the inversion kernels. The vertical dashed line indicates the base of the convection zone. [Fig. H-III:5.1; panel a is based on data from this source: Hathaway et al. (2000); source panel b: Thompson et al. (2003).]

step at the end of the 'stretch-twist-fold' process leads to a similar result. In the case of the 'stretch-twist-fold' dynamo the sign of the twist does not matter."

The driving flow of dynamos in stars and planets is energy-transporting convection. [H-III:5.2] "Thermal convection is familiar to most of us from our daily experience; warm air rises and cooler air sinks. When a fluid is heated from below it overturns, provided the temperature gradient is large enough, which here means that it must not only be greater than the adiabatic temperature gradient (the Schwarzschild criterion) but it must also overcome stabilizing influences such as thermal and viscous diffusion, rotation, compositional gradients (the Ledoux criterion), and magnetic flux. An intuitive way to think about convection (and to derive the Schwarzschild and Ledoux criteria) is to consider a small isolated volume, or *parcel*, of fluid that will buoyantly rise like a hot air balloon if its density is less than that of its surroundings or sink like a stone if its density is greater (the parcel is assumed to be in pressure equilibrium with its surroundings so density and temperature are anticorrelated). [For a compressible medium, t]his is the conceptual framework

behind mixing length theory which goes on to say that the parcel will lose its identity, dispersing into the background, after traveling a vertical distance of order a pressure scale height H_p . [...]

With this intuitive picture in mind, we may expect that the vertical scale of solar convection should vary tremendously from the deep convection zone where the stratification is relatively gentle ($H_p \sim 35$ Mm) to the solar surface layers where the density and pressure drop precipitously ($H_p \sim 36$ km) as [radiation escapes freely into space]. The associated drop in temperature near the surface triggers the recombination of hydrogen and other ions, which modifies the opacity, decreases the particle number density, and releases latent heat, altering the thermodynamics (in particular the equation of state and the specific heats) and contributing to the convective enthalpy transport. Add in radiative energy transfer and the result is what we call solar granulation; the continually shifting pattern (lifetime ~ 5 min) of small-scale convection cells (with a horizontal extent ~ 1 Mm) that blankets the solar surface and accounts for the dappled appearance of the solar photosphere (Fig. H-I:8.3).” {A:[43]}

A:43

Also the global-scale flows are important in the solar dynamo. The solar surface exhibits a differential rotation: the equator rotates faster than the poles, with a smooth latitudinal gradient between these. {A:[44]} [H-III:5.2.3] “Helioseismology now reveals that this monotonic decrease in angular velocity with increasing latitude persists throughout the convection zone, with an abrupt transition to nearly uniform rotation in the radiative interior (Fig. 4.7b). The transition region near the base of the convection zone is known as the solar tachocline [...]. There is also a less dramatic but no less significant *near-surface shear layer* in which the rotation rate systematically decreases by about 10-20 nHz from $r = 0.96R_\odot$ to the photosphere. This is most apparent at low latitudes but may also occur at higher latitudes. [...]

A:44

A:45

⁴³ Activity: Look up sample images of solar granulation, the most easily detectable pattern of convection reaching into the solar surface layers. What are the characteristic length and time scales of granulation? Also look up the larger-scale flow patterns of mesogranulation and supergranulation.

⁴⁴ Activity: Estimate the time it takes for the solar equator to execute one more full rotation than the poles in the same time.

⁴⁵ Activity: Helioseismology uses resonant waves that run through the solar interior. These pressure-mode (or p mode) waves (generated by the turbulent convective motions) probe a range of depths depending on the wavelength and resonance conditions. At depth, downward traveling waves refract upward as the sound speed increases with temperature. If their frequency is below the ‘acoustic cutoff period’ around the photosphere upward traveling waves are reflected back into the interior, even as they are detectable around their upper turning point both in brightness (by compression and dilation) and velocity (through the Doppler effect on spectral lines). The combination of refraction and reflection leads to a cavity in which resonances occur. Intuitively, the cutoff frequency comes about because if the wavelength of a pressure wave exceeds a few pressure scale heights, there is essentially no restoring pressure force as the bulk of the atmospheric mass is simply lifted and lowered in response to the wave. Based on that argument, make a rough estimate of the acoustic cutoff period for the solar photosphere at around 5800 K (a later Activity will let you develop the relevant equations for an isothermal atmosphere). Waves with shorter periods continue to travel upward, while those with longer periods mostly reflect but partly tunnel through into the hotter chromosphere.

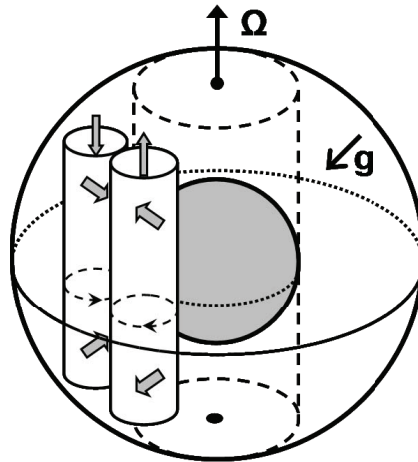


Fig. 4.8. Columnar convection in a rotating spherical shell near onset. The inner core tangent cylinder is shown by broken lines. Under Earth's core conditions the columns would be much thinner and very numerous. [Fig. H-III:7.6]

A:46

The striking difference in the rotation profile of the convective envelope and that of the radiative interior implicates convection as the primary source of the differential rotation. Furthermore, it tells us that giant cells are large enough and slow enough to be influenced by the rotation of the star. The magnitude of nonlinear advection relative to the Coriolis force $[(\boldsymbol{\Omega} \times \mathbf{v})]$ is quantified by the [Rossby number:

$$N_R = \frac{v_t}{\Omega L_t}, \quad (4.2)$$

where v_t and L_t are characteristic velocity and length scale, respectively.] In the deep solar convection zone it is of order unity or less whereas it is much greater than unity in the solar surface layers. Coriolis-induced velocity correlations in the convection redistribute angular momentum via the Reynolds stress, generating a substantial rotational shear: $\Delta\Omega/\Omega \sim 30\%$ where $\Omega(r, \theta)$ is the angular velocity and $\Delta\Omega$ is the angular velocity difference between equator and pole. {A:[47]} Furthermore, the nature of the redistribution is such that

A:47

⁴⁶ Activity: To hear how helioseismology can measure rotation rates of stars (and, with enough different modes, of layers within stars) you can do the following experiment: Hold up a bell dangling from a string, strike it, and listen. Then twist up the string, let the bell spin freely, hit it and listen once more. The modulation in intensity that you hear for the spinning bell results from the beat of the Doppler effect working differentially on waves running with and against the spin direction. This is the essence of how helioseismology measures the Sun's internal rotation.

⁴⁷ Activity: If we take the Sun's polar field – averaging at cycle minimum at about 5 Gauss – how long would it take to wind that field into a strength of some 10^5 G – which is the estimated minimum field strength for flux bundles to survive their rise through the convective motions in the Sun's envelope –

the angular velocity increases away from the rotation axis, $\partial\Omega/\partial d_\theta > 0$ where $d_\theta = r \sin(\theta)$ is [the distance to the axis of rotation]. This is in stark contrast to the behavior one would expect from isotropic turbulent diffusion (if $\Delta\Omega/\Omega \ll 1$) or from fluid parcels that tend to locally conserve their angular momentum as they move ($\partial\Omega/\partial d_\theta < 0$), [which would behave as sketched in Fig. 4.8]. Giant cells must be a global phenomenon distinct from supergranulation.”

These solar flow patterns are in striking contrast to what fluid motions in the planets are thought to look like: [H-III:5.5.4] “the latter often tend to be quasi-two-dimensional. This is largely a consequence of rapid rotation. Planets are smaller than stars and generally spin faster (with the exception of compact remnants such as pulsars). In the fluid cores and mantles of terrestrial planets and the extended atmospheres of many gas giant planets, the convective time scales are much longer than the rotation period, implying very low Rossby numbers [... T]his gives rise to elongated, quasi-2D convective structures such that the flow is relatively invariant in the direction parallel to the rotation axis (Fig. 4.8). In the atmospheres and oceans of terrestrial planets, on the other hand, quasi-2D dynamics arises simply by virtue of the geometry; global-scale horizontal motions are confined to thin spherical shells.”

4.4 Insights from approximate stellar dynamo models

Astrophysical dynamos have been studied for many decades, and whereas the fundamental ingredients may be known, there is no proper theory of dynamo action in stars and planets: there is no validated dynamo model that matches all stellar observations or that has been demonstrated to successfully forecast the Sun’s magnetism over multiple sunspot cycles, nor do planetary dynamo models successfully reproduce, for example, the quasi-irregular reversals in the terrestrial magnetic field. Nonetheless, dynamo concepts do guide our thinking as to the important ingredient processes as well as the possible internal structure and dynamics of both the magnetic field and the plasma/magma flows involved. The remainder of this chapter is an exploration of some of these to create a sense of how dynamos in stars and planets are thought to function.

[H-III:6.1] “All solar and stellar dynamo models to be considered in this chapter operate within a sphere of electrically conducting fluid embedded in vacuum. We restrict ourselves here to axisymmetric mean-field-like models, in the sense that we will be setting and solving evolutionary equations for the large-scale magnetic field, and subsume the effects of small-scale fluid motions and magnetic fields into coefficients of these partial differential equations.

if the rotational shear were maximally used and if no back-reaction on that flow occurred? Hint: remember the field line stretch-and-fold from Fig. 4.6, look at the illustration in Fig. 4.9, and consider ‘compound interest’.

Working in spherical polar coordinates (r, θ, ϕ) , we begin by writing:

$$\mathbf{v}(r, \theta) = \mathbf{v}_p(r, \theta) + d_\theta \Omega(r, \theta) \hat{\mathbf{e}}_\phi, \quad (4.3)$$

$$\mathbf{B}(r, \theta, t) = \nabla \times (A(r, \theta, t) \hat{\mathbf{e}}_\phi) + B(r, \theta, t) \hat{\mathbf{e}}_\phi, \quad (4.4)$$

where $d_\theta = r \sin(\theta)$, \mathbf{v}_p is a notational shortcut for the component of the large scale flow in meridional planes, and Ω is the angular velocity of rotation, which in the solar interior varies with both depth and latitude, and is now well-constrained by helioseismology. Note that in this prescription neither of these large-scale flow components is time dependent. This **kinematic approximation** is an assumption that is tolerably well-supported observationally. Substituting these expressions in the MHD induction equation in Eq. (3.3) allows separation into two coupled 2D partial differential equations for the scalar functions A and B defining respectively the poloidal and toroidal components of the magnetic field:

$$\frac{\partial A}{\partial t} = \eta \left(\nabla^2 - \frac{1}{d_\theta^2} \right) A - \frac{\mathbf{v}_p}{d_\theta} \cdot \nabla (d_\theta A), \quad (4.5)$$

$$\begin{aligned} \frac{\partial B}{\partial t} = & \eta \left(\nabla^2 - \frac{1}{d_\theta^2} \right) B + \frac{1}{d_\theta} \frac{\partial (d_\theta B)}{\partial r} \frac{\partial \eta}{\partial r} - \\ & d_\theta \nabla \cdot \left(\frac{B}{d_\theta} \mathbf{v}_p \right) + d_\theta (\nabla \times (A \hat{\mathbf{e}}_\phi)) \cdot \nabla \Omega, \end{aligned} \quad (4.6)$$

where we retain the possibility that η varies with depth. The shearing term ($\propto \nabla \Omega$) on the right-hand side of Eq. (4.6) acts as a source of toroidal field. However, no such source term appears in Eq. (4.5). This is the essence of Cowling's theorem which in fact guarantees that an axisymmetric flow of the general form given by Eq. (4.3) *cannot* act as a dynamo for an axisymmetric magnetic field as described by Eq. (4.4). The construction of solar and stellar dynamo models, therefore, hinges critically on the addition of an extraneous source term in Eq. (4.5). The physical origin of this source term is what fundamentally distinguishes the various classes of solar and stellar dynamo models described [below].

Shearing of the poloidal magnetic field into a strong toroidal component by differential rotation [(as illustrated in Fig. 4.9)] is an essential ingredient of all solar cycle models discussed below. The growing magnetic energy of the toroidal field is supplied by the kinetic energy of the rotational shearing motion, which makes for an attractive field amplification mechanism, because in the Sun and stars the available supply of rotational kinetic energy is immense (unless the dynamo were entirely confined to a very thin layer, for example the tachocline, [the shear layer just below the Sun's convective envelope into which convection overshoots]). Moreover, a strong, axisymmetric and temporally

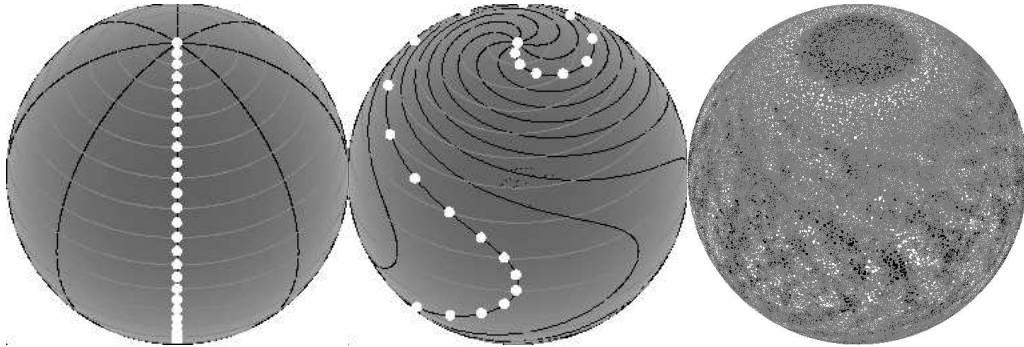


Fig. 4.9. Left and center: Visualization of the effects of differential rotation and equator-to-pole meridional flow for Sun-like conditions: lines of equal longitude (with markers) are distorted into a spiral pattern. The center panel shows the distorted lines after 3 months. Right: Simulated magnetogram for a star like the Sun, [simulated with a flux-transport model with parameters as observed for the Sun,] but with an active-region emergence rate 30 times larger. The simulated star is shown from a latitude of 40° to better show the polar-cap field structure. [Fig. H-III:2.3]

quasi-steady internal differential rotation is likely responsible for the observed high degree of axisymmetry observed in the Sun’s magnetic field on spatial scales comparable to its radius. This situation is very different from that encountered in planetary core dynamos, where differential rotation is believed to be much weaker, and energetics pose a much stronger constraint on dynamo action. Lacking the large-scale organization provided by differential rotation, planetary core dynamos also tend to produce non-axisymmetric large-scale fields. The one outstanding exception appears to be Saturn, and indeed in this case the high axisymmetry of the observed surface field may well reflect the symmetrizing action of differential rotation in the envelope overlying the metallic-hydrogen core. The important point remains that in the solar dynamo context, the assumption of an axisymmetric large-scale magnetic field is consistent with the observed and helioseismically-inferred axisymmetry and quasi-steadiness of internal differential rotation.”

4.5 Mean-field dynamo models

[H-III:6.2.1] “Turbulence at a high magnetic Reynolds number \mathcal{R}_m [(Eq. 3.18)] is known to be quite effective at producing a lot of small-scale magnetic fields, where ‘small-scale’ is roughly $\mathcal{R}_m^{-1/2}$ times the length scale of the flow. In addition, under certain conditions, solar/stellar convective turbulence can also produce magnetic fields with a mean component building up on large spatial scales. These **mean-field dynamo models** remain arguably the most

'popular' descriptive models for dynamo action in the Sun and stars, but also in planetary metallic cores, stellar accretion disks, and even galactic disks.

Under the assumption that a good separation of scales exists between the large-scale 'laminar' magnetic field $\overline{\mathbf{B}}$ and the flow $\overline{\mathbf{v}}$, and the small-scale turbulent field \mathbf{B}' and flow \mathbf{v}' , it becomes possible to express the inductive and diffusive action of the turbulence on $\overline{\mathbf{B}}$ in terms of the statistical properties of the small-scale flow and field. ^{A:48} $\{\text{A:}^{[48]}\}$ The corresponding theory of mean-field electrodynamics is discussed in detail in Ch. H-I:3. The turbulent flow introduces on the right-hand side of the induction equation Eq. (3.3) a term of the form $\nabla \times \overline{\mathcal{E}}$, where $\overline{\mathcal{E}}$ is a mean-electromotive force."

For a quick introduction to the origin of that term we can see what happens when [H-I:3.4] "we decompose the magnetic field into a large-scale 'mean' field and the small-scale components through an averaging procedure. We assume in the following that the averaging procedure obeys the Reynolds rules: For any function f and g decomposed as $f = \overline{f} + f'$ and $g = \overline{g} + g'$, where the bar indicates the averaged and the prime the fluctuating quantity, we require that

$$\overline{\overline{f}} = \overline{f} \longrightarrow \overline{f'} = 0 \quad (4.7)$$

$$\overline{f + g} = \overline{f} + \overline{g} \quad (4.8)$$

$$\overline{f'g'} = \overline{f'g} \longrightarrow \overline{f'g'} = 0 \quad (4.9)$$

$$\overline{\partial f / \partial x_i} = \partial \overline{f} / \partial x_i \quad (4.10)$$

$$\overline{\partial f / \partial t} = \partial \overline{f} / \partial t. \quad (4.11)$$

The averaging procedures that are of interest in the context of mean-field theory are the ensemble average (meaning a chaotic system is averaged over several representations of the chaotic system) and the longitudinal average, in which $\overline{\mathbf{B}}$ reflects the axisymmetric component of the large-scale magnetic field (multipole series with $m = 0$)." [H-I:3.4.1] "In order to derive an equation for the time evolution of the mean field we apply the averaging procedure to the induction equation Eq. (3.3) which leads to

$$\frac{\partial \overline{\mathbf{B}}}{\partial t} = \nabla \times (\overline{\mathbf{v}' \times \mathbf{B}'} + \overline{\mathbf{v}} \times \overline{\mathbf{B}} - \eta \nabla \times \overline{\mathbf{B}}). \quad (4.12)$$

The new term which enters this equation compared to the original induction

⁴⁸ Activity: Consider that the assumption of a separation of scales as for the 'mean field dynamo theory' is also made in hydrodynamics when 'internal energy' (which includes the kinetic energy of the random motions of the gas particles) is 'separated from' the 'kinetic energy' of bulk motion. This assumption is commonly made with little consideration of why it works: there must be a scale that is small compared to flows of interest but large enough that low-order moments of the velocity distribution (like temperature and pressure) are defined by so many particles that there is negligible random noise when determined for a 'small' volume. Consider that in the context of the words in Table 3.2.

equation is the second order correlation electromotive force (EMF)

$$\overline{\mathcal{E}} \equiv \overline{\mathbf{v}' \times \mathbf{B}'} . \quad (4.13)$$

While the fluctuating velocity component \mathbf{v}' is assumed to be known (kinematic approach), \mathbf{B}' has to be computed from the induction equation. An equation for \mathbf{B}' can be derived by subtracting the mean-field induction equation Eq. (4.12) from the microscopic induction equation Eq. (3.3), which leads to

$$\frac{\partial \mathbf{B}'}{\partial t} = \nabla \times (\mathbf{v}' \times \overline{\mathbf{B}} + \overline{\mathbf{v}} \times \mathbf{B}' - \eta \nabla \times \mathbf{B}' + \mathbf{v}' \times \mathbf{B}' - \overline{\mathbf{v}' \times \mathbf{B}'}) . \quad (4.14)$$

It is in general only possible to solve this equation by making strong assumptions, primarily because of the terms that are quadratic in the fluctuating quantities (closure problem).” For a more detailed description, see Sect. H-I:3.4.3. Here, we proceed with one particular such assumption that leads to the conclusion that for [H-III:6.2.1] “mildly inhomogeneous and near-isotropic turbulence, $\overline{\mathcal{E}}$ can be expressed in terms of the large-scale field $\overline{\mathbf{B}}$ as:

$$\overline{\mathcal{E}} = \alpha \overline{\mathbf{B}} + \beta \nabla \times \overline{\mathbf{B}} , \quad (4.15)$$

with

$$\alpha = -\frac{1}{3} \tau_{\text{corr}} \overline{\mathbf{v}' \cdot (\nabla \times \mathbf{v}')} \quad [\text{cm s}^{-1}] , \quad \beta = \frac{1}{3} \tau_{\text{corr}} \overline{\mathbf{v}'^2} \quad [\text{cm}^2 \text{s}^{-1}] , \quad (4.16)$$

where τ_{corr} is the correlation time for the turbulent flow. Note that the α -term is proportional to the (negative) kinetic helicity $[(\mathbf{v}' \cdot (\nabla \times \mathbf{v}'))]$ of the turbulence, which requires a break of reflectional symmetry. In stellar interiors and planetary metallic cores alike, this anisotropy is provided by the Coriolis force. Small-scale turbulence thus impacts the induction equation for the mean-field in two ways: it introduces a field-aligned electromotive force (the α -term), which acts as a source term and is called the ‘ α -effect’, and an enhanced ‘turbulent diffusion’ (the β -term), associated with the folding action of the turbulent flow. In principle, the α and β coefficients can be calculated from the lowest-order statistics of the turbulent flow. In practice, more often than not they are chosen *a priori*, although with care taken to embody in these choices what can be learned from mean-field theory. {A:[49]}

A:49

⁴⁹ Activity: Take the mean-field induction equation Eq. (4.12) and the expression for Eq. (4.13) as in Eq. (4.15) to find a mean-field form of the general induction equation Eq. (3.3). Group the ‘diffusive’ terms together. Estimate the order of magnitude of the advection, α , and diffusive terms. For these order of magnitude comparisons, approximate for the mean field $\nabla \approx 1/R_{\odot}$; in solar near-surface layers ‘small-scale’ random walk leads to $\beta \approx 300 \text{ km}^2/\text{s}$; the large-scale advective term of the surface meridional flow has an average value of order 5 m/s (peaking at about 15 m/s); estimate τ_{corr} from this value of β with Eq. (4.16), which corresponds to the characteristic evolutionary time scale of the dispersing supergranular convection; with that, estimate α using the characteristic supergranulation length scale of 30,000 km; then compare the order-of-magnitude values of the three terms expressed as time scales for the magnetic field. Note that the ‘turbulent diffusivity’ β far exceeds the ‘resistive diffusivity’ η in stellar dynamos (and see Activity 51 how the above helps in understanding how surface flux dispersal can be described quite well by a random-walk diffusive description).

Under mean-field dynamo theory, Eqs. (4.5)–(4.6) are now taken to apply to an axisymmetric large-scale mean magnetic field. With the inclusion of the mean-field α -effect and turbulent diffusivity, scaling all lengths in terms of the radius R of star or planet, and time in terms of the diffusion time

$$\tau_d = R^2/\eta \quad (4.17)$$

based on the (turbulent) diffusivity in the convective envelope, these expressions become

$$\frac{\partial A}{\partial t} = \eta \left(\nabla^2 - \frac{1}{d_\theta^2} \right) A - \frac{\mathcal{R}_m}{d_\theta} \mathbf{v}_p \cdot \nabla (d_\theta A) + C_\alpha \alpha B, \quad (4.18)$$

$$\begin{aligned} \frac{\partial B}{\partial t} = & \eta \left(\nabla^2 - \frac{1}{d_\theta^2} \right) B + \frac{1}{d_\theta} \frac{\partial (d_\theta B)}{\partial r} \frac{\partial \eta}{\partial r} - \mathcal{R}_m d_\theta \nabla \cdot \left(\frac{B}{d_\theta} \mathbf{v}_p \right) + \\ & C_\Omega d_\theta (\nabla \times (A \hat{\mathbf{e}}_\phi)) \cdot (\nabla \Omega) + C_\alpha \hat{\mathbf{e}}_\phi \cdot \nabla \times [\alpha \nabla \times (A \hat{\mathbf{e}}_\phi)]. \end{aligned} \quad (4.19)$$

We continue to use the symbol η for the total diffusivity, with the understanding that within the convective envelope this now includes the (dominant) contribution from the β -term of mean-field theory. Three non-dimensional numbers have materialized:

$$C_\alpha = \frac{\alpha_t R}{\eta}, \quad C_\Omega = \frac{\Omega_t R^2}{\eta}, \quad \mathcal{R}_m = \frac{u_t R}{\eta}, \quad (4.20)$$

with α_t , u_t , and Ω_t as reference values for the α -effect, meridional flow and envelope rotation, respectively. The quantities C_α and C_Ω are *dynamo numbers*, measuring the importance of inductive versus diffusive effects on the right-hand side of Eqs. (4.18)–(4.19). The magnetic Reynolds number \mathcal{R}_m here measures the relative importance of advection versus diffusion in the transport of A and B in meridional planes. ^{A:50} Structurally, Eqs. (4.18)–(4.19) only differ from Eqs. (4.5)–(4.6) by the presence of two new source terms on the right-hand side, both associated with the α -effect. The appearance of this term in Eq. (4.18) is crucial for evading Cowling’s theorem.”

In what follows in this section, we first look at a simplified, linear mean-field dynamo model to illustrate the geometry and temporal evolution. Later, we look at non-linearities that lead to amplification and saturation of the field, and to the modulation of the magnetic cycles. First, the linear model: [H-III:6.2.1.1] “In constructing mean-field dynamos for the Sun, it has been a common procedure to neglect meridional circulation, because it is a very weak flow. It is also customary to drop the α -effect term on the right-hand side of Eq. (4.19) on the grounds that with $R \simeq 7 \times 10^{10}$ cm, $\Omega_t \sim 10^{-6}$ rad s $^{-1}$, and $\alpha_t \sim 10^2$ cm s $^{-1}$, one finds $C_\alpha/C_\Omega \sim 10^{-3}$, independently of the assumed (and

⁵⁰ Activity: Relate the Rossby number in Eq. (4.2) to the dynamo number C_Ω and the magnetic Reynolds number \mathcal{R}_m in Eq. 4.20: $N_R = \mathcal{R}_m^2/C_\Omega$.

poorly constrained) value for η . Equations (4.18)—(4.19) then reduce to the so-called **$\alpha\Omega$ dynamo equations**. In the spirit of producing a model that is solar-like we use a fixed value $C_\Omega = 2.5 \times 10^4$, obtained by assuming [an equatorial angular velocity of] $\Omega_{\text{Eq}} \simeq 10^{-6} \text{ rad s}^{-1}$ and $\eta = 50 \text{ km}^2\text{s}^{-1}$, which leads to a diffusion time $\tau_d = R^2/\eta \simeq 300 \text{ yr}$.

For the total magnetic diffusivity, we use a steep but smooth variation of η from a high value (η_{CZ}) in the convection zone to a low value (η_{core}) in the underlying core [...]. A typical profile is shown in Fig. 4.10A (dash-dotted line). In practice, the core-to-envelope diffusivity ratio $\Delta\eta \equiv \eta_{\text{core}}/\eta_{\text{CZ}}$ is treated as a model parameter, with of course $\Delta\eta \ll 1$, because we associate η_{core} with the microscopic magnetic diffusivity, and η_{CZ} with the presumably much larger mean-field turbulent diffusivity. Taking at face values estimates from mean-field theory, one should have $\Delta\eta \sim 10^{-9}$ to 10^{-6} . The solutions discussed below have $\Delta\eta = 10^{-3}$ to 10^{-1} , which is still small enough to illustrate important effects of radial gradients in total magnetic diffusivity.

All solar dynamo models discussed in this chapter utilize the helioseismically-calibrated solar-like parametrization of solar differential rotation [...]. The corresponding angular velocity contour levels are plotted in Fig. 4.10B. Such a solar-like differential rotation profile is quite complex from the point of view of dynamo modelling, in that it is characterized by multiple partially overlapping shear regions: a rotational shear layer, straddling the core-envelope interface, known as the *tachocline*, with a strong positive radial shear in its equatorial regions and an even stronger negative radial shear in its polar regions, as well as a significant latitudinal shear throughout the convective envelope and extending partway into the tachocline; for a tachocline of half-thickness $w/R_\odot = 0.05$, the mid-latitude latitudinal shear at $r/R_\odot = 0.7$ is comparable in magnitude to the equatorial radial shear, and its potential contribution to toroidal field production cannot be casually dismissed.

For the dimensionless function $\alpha(r, \theta)$ we use an expression [... that] concentrates the α -effect in the bottom half of the envelope, and lets it vanish smoothly below, just as the net magnetic diffusivity does (see Fig. 4.10A). Various lines of argument point to an α -effect peaking in the bottom half of the convective envelope, because there the convective turnover time is commensurate with the solar rotation period, a most favorable setup for the type of toroidal field twisting at the root of the α -effect (see Fig. H-I:3.5). [The choice made here for $\alpha(r, \theta)$ scales with latitude as $\cos\theta$, which] reflects the hemispheric dependence of the Coriolis force, which also suggests that the α -effect should be positive in the Northern hemisphere. The dimensionless number C_α , which measures the strength of the α -effect, is treated as a free parameter of the model. [...]

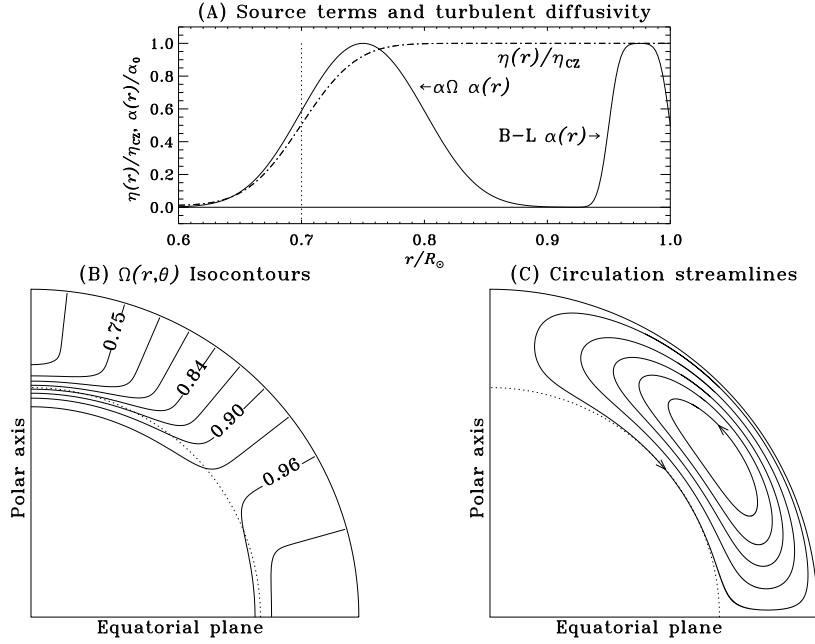


Fig. 4.10. Various 'ingredients' for the dynamo models constructed in this chapter. Part (A) shows radial profiles of the total magnetic diffusivity η and poloidal source [terms: $\alpha(r)$ for the $\alpha\Omega$ dynamo and for the Babcock-Leighton (B-L) dynamo]. Part (B) shows contour levels of the rotation rate $\Omega(r, \theta)$ normalized to its surface equatorial value. The dotted line is the core-envelope interface at $r/R_{\odot} = 0.7$. Part (C) shows streamlines of meridional circulation, included in some of the dynamo models discussed below. [Helioseismic studies suggest that the meridional flow in the Sun is more complex than a single 'roll' of the flow, but that there may be (at least) two stacked on top of each other. A key point for a flux-transport dynamo is that the meridional flow at the base of the convective envelope is equatorward. Fig. H-III:6.1]

In such linear $\alpha\Omega$ models the onset of dynamo activity turns out to be controlled by the product of C_{α} and C_{Ω} :

$$D \equiv C_{\alpha} \times C_{\Omega} = \frac{\alpha_t \Omega_t R^3}{\eta_{CZ}^2}. \quad (4.21)$$

with positive growth rates materializing above a threshold value known as the *critical dynamo number*. [...]

Figure 4.11 shows half a cycle of the dynamo solution, in the form of snapshots of the toroidal (gray scale) and poloidal eigenfunctions (field lines) in a meridional plane, with the symmetry axis defined by the stellar rotation oriented vertically. The four frames are separated by a phase interval $\varphi = \pi/3$, so that panel (D) is identical to panel (A) except for reversed magnetic

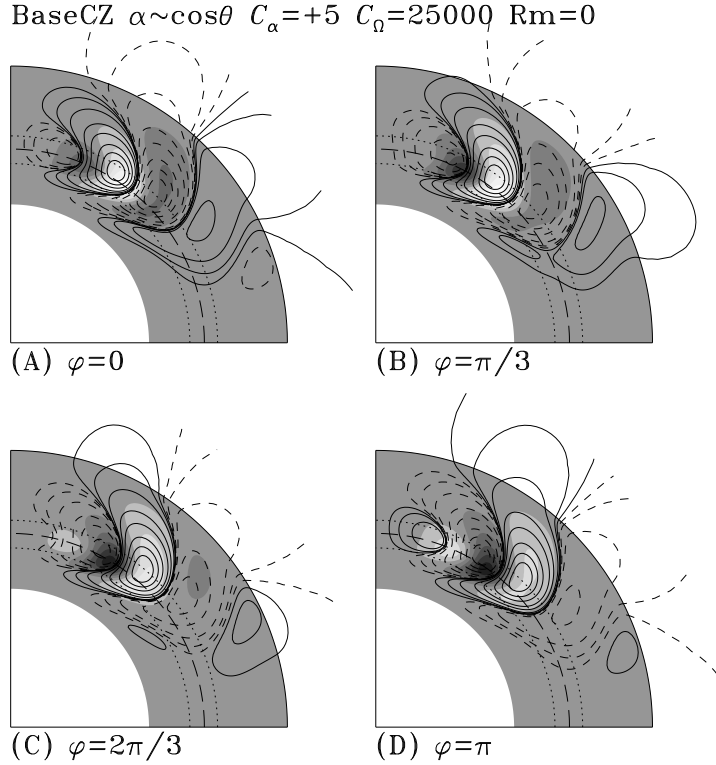


Fig. 4.11. Four snapshots in meridional planes of our minimal linear $\alpha\Omega$ dynamo solution with defining parameters $C_\Omega = 25000$, $\Delta\eta = 0.1$, and $\eta_{CZ} = 50 \text{ km}^2/\text{s}$. With $C_\alpha = +5$, this is a mildly supercritical solution, with oscillation frequency $\omega \simeq 300 \tau_d^{-1}$ (see Eq. 4.17). The toroidal field is plotted as filled contours (gray to black for negative B , gray to white for positive B , normalized to the peak strength and with increments $\Delta B = 0.2$), on which poloidal field lines are superimposed (solid for clockwise-oriented field lines, dashed for counter-clockwise orientation). The long-dashed line is the core-envelope interface at $r/R_\odot = 0.7$. [Fig. H-III:6.2]

polarities in both magnetic components [halfway through the cycle with period $P_{\text{cycle}} = 2\pi/\omega$]. Such linear eigensolutions leave the absolute magnitude of the magnetic field undetermined, but the relative magnitude of the poloidal to toroidal components is found to scale approximately as $|C_\alpha/C_\Omega|$.

The [model's magnetic field] is concentrated in the vicinity of the core-envelope interface, and has very little amplitude in the underlying, low-diffusivity radiative core. This is due to the oscillatory nature of the solution, which restricts penetration into the core to a distance of the order of the electromagnetic skin depth $\ell_{\text{skin}} = \sqrt{2\eta_{\text{core}}/\omega}$. Having assumed $\eta_{CZ} = 50 \text{ km}^2\text{s}^{-1}$, with $\Delta\eta = 0.1$, a dimensionless dynamo frequency $\omega \simeq 300$ corresponds to $3 \times 10^{-8} \text{ s}^{-1}$, so that $\ell_{\text{skin}}/R \simeq 0.026$, quite small indeed.

Careful examination of Fig. 4.11(A)→(D) also reveals that the toroidal-poloidal flux systems present in the shear layer first show up at high latitudes, and then *migrate equatorward* to finally disappear at mid-latitudes in the course of the half-cycle. These *dynamo waves* travel in a direction given by $\alpha\nabla\Omega \times \hat{\mathbf{e}}_\phi$, *i.e.*, along contours of equal angular velocity, a result known as the *Parker-Yoshimura sign rule*. Here with a negative $\partial\Omega/\partial r$ in the high-latitude region of the tachocline, a positive α -effect results in an equatorward propagation of the dynamo wave, in qualitative agreement with the observed equatorward drift of the latitudes of sunspot emergences as the solar cycle unfolds (see Fig. 4.5).”

[H-III:6.2.1.2] “Obviously, the exponential growth characterizing supercritical linear solutions must stop once the Lorentz force associated with the growing magnetic field becomes dynamically significant for the inductive flow. Because the solar surface and internal differential rotation show little variation with the phase of the solar cycle, it is usually assumed that magnetic back-reaction occurs at the level of the α -effect. In the mean-field spirit of *not* solving dynamical equations for the small-scales, it has become common practice to introduce an *ad hoc* algebraic nonlinear quenching of α directly on the mean-toroidal field B by writing:

$$\alpha \rightarrow \alpha(B) = \frac{\alpha_t}{1 + (B/B_{\text{eq}})^2}. \quad (4.22)$$

where $B_{\text{eq}} = (4\pi\rho u_t^2)^{1/2}$ is the equipartition field strength, of order 10^4 G at the base of the solar convective envelope. Needless to say, this simple **α -quenching** formula is an *extreme* oversimplification of the complex interaction between flow and field that is known to characterize MHD turbulence, but its wide usage in solar dynamo modeling makes it the nonlinearity of choice for the illustrative purpose of this [chapter: with this description, the only MHD equation that needs solving to experiment with dynamo action – as we do here – is the induction equation Eq. (3.3) that is now subjected to a parameterized coupling between the small-scale flow and field that may or may not be an appropriate approximation of reality. Note that α can, and in many models now is, time dependent, leading to what is called ‘dynamical α -quenching’.]

Introducing α -quenching in our model renders the $\alpha\Omega$ dynamo equations nonlinear, so that solutions are now obtained as initial-value problems starting from an arbitrary seed field of very low amplitude, in the sense that $B \ll B_{\text{eq}}$ everywhere in the domain. [...] At early times, $B \ll B_{\text{eq}}$ and the equations are effectively linear, leading to exponential growth [...]. Eventually, however, B becomes comparable to B_{eq} in the region where the α -effect operates, leading to a break in exponential growth, and eventual saturation.

The saturation energy level increases with increasing C_α , an intuitively satisfying behavior because solutions with larger C_α have a more vigorous

poloidal source term. The cycle frequency for these solutions is very nearly independent of the dynamo number, and is slightly *smaller* than the frequency of the linear critical mode (here by some 10 – 15%), a behavior that is typical of kinematic α -quenched mean-field dynamo models. Yet the overall form of the dynamo solutions very closely resembles that of the linear eigenfunctions plotted in Fig. 4.11.”

[H-III:6.2.1.3] “The α -quenching expression in Eq. (4.22) implies that dynamo action saturates once the mean, dynamo-generated large-scale magnetic field reaches an energy density comparable to that of the driving small-scale turbulent fluid motions. However, various calculations and numerical simulations have indicated that long before the mean toroidal field B reaches this strength, the helical turbulence reaches equipartition with the *small-scale* turbulent component of the magnetic field. Such calculations also suggest that the ratio between the small-scale and mean magnetic components should itself scale as $\mathcal{R}_m^{1/2}$, where $\mathcal{R}_m = v_t L_t / \eta$ is a magnetic Reynolds number based on the turbulent speed but *microscopic* magnetic diffusivity. This then leads to the alternative quenching expression

$$\alpha \rightarrow \alpha(B) = \frac{\alpha_t}{1 + \mathcal{R}_m (B/B_{\text{eq}})^2} , \quad (4.23)$$

known in the literature as **strong α -quenching** or *catastrophic quenching* (see Ch. H-I:3 in Vol. I). Because $\mathcal{R}_m \sim 10^8$ in the solar convection zone, this leads to quenching of the α -effect for very low amplitudes of the mean magnetic field, of order 0.1 G. Even though significant field amplification is likely in the formation of a toroidal flux rope from the dynamo-generated magnetic field, we are now a very long way from the $10^4 - 10^5$ G demanded by simulations [needed for buoyantly rising flux ropes to survive emergence and to eventually lead to] sunspot formation.

[One] way out of this difficulty exists in the form of **interface dynamos**. The idea is beautifully simple: to produce and store the toroidal field away from where the α -effect is operating. [...] in a situation where a radial shear and α -effect are segregated on either side of a discontinuity in magnetic diffusivity taken to coincide with the core-envelope interface, the constant coefficient, cartesian form of the $\alpha\Omega$ dynamo equations support solutions in the form of traveling surface waves localized on the discontinuity in diffusivity. For supercritical dynamo waves, the ratio of peak toroidal field strength on either side of the discontinuity surface is found to scale as $(\eta_{\text{CZ}}/\eta_{\text{core}})^{-1/2}$. With the core diffusivity η_{core} equal to the microscopic value, and if the envelope diffusivity is of turbulent origin so that $\eta_{\text{CZ}} \sim L_t v_t$, then the toroidal field strength ratio scales as $\sim (v_t L_t / \eta_{\text{core}})^{1/2} \equiv \mathcal{R}_m^{1/2}$. This is precisely the factor needed to bypass strong α -quenching, at least as embodied in Eq. (4.23).”

So far, this discussion has ignored the large-scale flow system known as meridional circulation. Such a flow [H-III:6.2.1.4] “is unavoidable in turbulent, compressible rotating convective shells. The $\sim 15 \text{ m s}^{-1}$ poleward flow observed at the surface has been detected helioseismically, down to $r/R_\odot \simeq 0.85$ without significant departure from the poleward direction, except locally and very close to the surface, in the vicinity of active region belts. Mass conservation requires an equatorward flow deeper down [(helioseismic measurements suggest that there may be two meridional overturning cells stacked within the convective envelope, but confirmation is still pending of what is a challenging measurement close to the noise levels of helioseismology)].

Meridional circulation can bodily transport the dynamo-generated magnetic field (terms $\propto \mathbf{v}_p \cdot \nabla$ in Eqs. (4.5)–(4.6)). At low circulation speeds, the primary effect is a Doppler shift of the dynamo wave, leading to a small change in the cycle period and equatorward concentration of the activity belts. However, for a (presumably) solar-like equatorward return flow that is vigorous enough, it can overpower the Parker-Yoshimura propagation rule and produce equatorward propagation no matter what the sign of the α -effect is. The behavioral turnover from dynamo wave-like solutions sets in when the circulation speed in the dynamo region becomes comparable to the propagation speed of the dynamo wave. In this advection-dominated regime, the cycle period loses sensitivity to the assumed turbulent diffusivity value, and becomes determined primarily by the circulation’s turnover time. Solar cycle models achieving equatorward migration of activity belts in this manner are often called **flux transport dynamos**. [...]

One interesting consequence [of meridional circulation] is that induction of the toroidal field is now effected primarily by the *latitudinal* shear within the tachocline, with the radial shear, although larger in magnitude, playing a lesser role because $B_r/B_\theta \ll 1$. The meridional flow also has a profound impact on the magnetic field evolution at $r = R$, as it concentrates the poloidal field in the polar regions. This leads to a large amplification factor through magnetic flux conservation, so [these dynamo models] are typically characterized by very large polar field strengths, here some 20% of the toroidal field magnitude in the tachocline, even though we have here $C_\alpha/C_\Omega = 10^{-6}$. This concentrated poloidal field, when advected downwards to the polar regions of the tachocline, is responsible for the strong polar branch often seen in the time-latitude diagram of dynamo solutions including a rapid meridional flow. This difficulty can be alleviated, at least in part, by a number of relatively minor modifications to the model, such as the addition of a high- η subsurface layer, or displacement of the meridional flow cell towards lower latitudes, thus reducing the degree of polar convergence. [...]

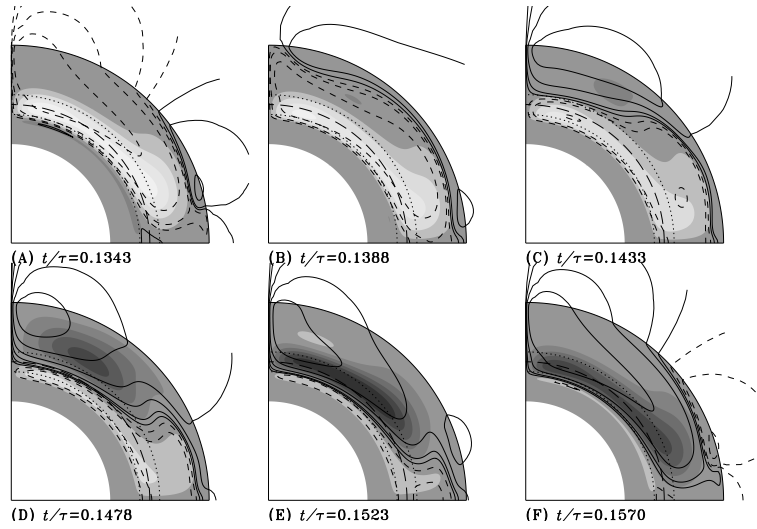


Fig. 4.12. [Snapshots covering half a cycle of a Babcock-Leighton dynamo solution. The grey-scale coding of the toroidal field and poloidal field lines is as in Fig. 4.11.] This solution uses the same differential rotation, magnetic diffusivity and meridional circulation profile as for the advection-dominated $\alpha\Omega$ solution but now with the non-local surface source term as shown in Fig. 4.10(A) in the curve labeled 'B-L', with parameter values $C_\alpha = 5$, $C_\Omega = 5 \times 10^4$, $\Delta\eta = 0.003$, $\mathcal{R}_m = 840$. Note the strong amplification of the surface polar fields, and the latitudinal stretching of poloidal field lines by the meridional flow at the core-envelope interface. [Fig. H-III:6.6]

Yet another incarnation of solar cycle models is based on active region decay and dispersal. These go back to 1961 when Babcock [H-III:6.2.2] “suggested that the polarity reversals of the high-latitude surface magnetic field are driven by the accumulation of magnetic fields released at low latitudes by the decay of bipolar magnetic regions. Figure 4.9 shows a numerical simulation illustrating this process, which leads to the buildup of a net poloidal hemispheric flux because the trailing member of the pair tends to be located at higher latitudes than the leading component, a pattern known as Joy’s rule, and therefore are subjected to less transequatorial dissipative flux cancellation than the leading members of the bipolar pair. Babcock went on to argue that in conjunction with shearing by differential rotation, this could explain the observed patterns of solar cycle polarity reversals. In subsequent years [he] turned this idea into a *bona fide* solar cycle model, known since as the **Babcock-Leighton model**. [...] The key point, from the dynamo perspective, is that the Babcock-Leighton mechanism taps into the (formerly) toroidal flux in the bipolar magnetic region to produce a poloidal magnetic component, and so can act as a source term on the right-hand side of Eq. (4.5). [...] {A:[51]} {A:[52]}

A:51

⁵¹ Activity: One of the basic concepts behind Babcock’s idea is that magnetic field at the solar surface

To the degree that a positive dipole moment is being produced from a toroidal field that is positive in the N-hemisphere, this is operationally equivalent to a positive α -effect in mean-field theory. In both cases the Coriolis force is the agent imparting a twist on a magnetic field; with the α -effect this process occurs on the small spatial scales and operates on individual magnetic field lines. In contrast, the Babcock-Leighton mechanism operates on the large scales, the twist being imparted via the Coriolis force acting on the flow generated along the axis of a buoyantly rising magnetic flux tube that, upon emergence, gives rise to sunspot pairs. {A:[53]}

A:52

A:53

Numerous dynamo models based on this mechanism of poloidal field regeneration have been constructed, based on the axisymmetric mean-field dynamo equations but with the α -effect replaced by a suitably designed source term on the right-hand side of Eq. (4.18). One important difference with the mean-field $\alpha\Omega$ models considered earlier is that the two source regions are now spatially segregated: production of the toroidal field takes place in or near the tachocline, as before, but now production of the poloidal field is restricted to the surface layers. A transport mechanism is then required to link the two source regions for a dynamo loop to operate. [... Most Babcock-Leighton models use the meridional circulation for this, which acts] as a form of conveyor belt,

is largely advected like a scalar quantity. Consequently, the field disperses in the random motions of the surface convection (with an equivalent diffusion coefficient of $D \approx 250 \text{ km}^2/\text{s}$) subject to the large-scale advection of the differential rotation and meridional flow. To see how this can be, use the ideal version of Eq. (3.3) and assume that the field is always vertical to the surface (a good approximation to the observed photospheric field, except during emergence and cancellation; a result of the buoyancy of flux bundles – see Activity (35) – and show it is equivalent to Eq. (3.4) for the advection of a scalar (without the source and loss terms in that version). Note that this formulation is linear, so that you can think about N and S polarities as diffusing separately, then to sum to obtain the net result; this helps visualize why the active-region tilt angle is important in reversing the polar fields from cycle to cycle. Question: with this value of D , what is the characteristic time scale for flux to disperse over the solar surface (hint: Eq. 3.20)? With that in mind, how important is the meridional advection from equator to pole (with a characteristic velocity of 10 m/s) in transporting the field within the duration of a solar cycle? Remember that Activity 49 shows how the diffusion coefficient β associated with (super-)granular random walk adds to the molecular/resistive diffusion coefficient η .

⁵² Activity: Joy's rule, that the leading polarities (in the direction of rotation) of active regions emerge statistically closer to the equator than the trailing polarities (see Fig. 4.4, for example), is the reason why eventually flux of the trailing polarity builds up a polar cap that reaches its maximum strength at cycle minimum. For some interval around that time, the bulk of the heliospheric field originates from the polar caps. Estimate the total flux in the solar wind (assuming an isotropic flux density at Earth orbit; use Table 2.4). This is the equivalent of only a few large active regions (Fig. 4.4) although it is in fact composed of a fraction of the flux from the ensemble of all bipolar regions emerging over a cycle.

⁵³ Activity: Why are solar photospheric flux tubes buoyant (hint: look back at Activity (35)? What is the maximum density contrast between interior and exterior? For an essentially evacuated flux tube at the solar surface, show that the buoyancy force per unit length (causing the tube to buoy towards vertical) dominates the dynamic pressure force exerted by a convective flow (which could bend the tube away from vertical) of $v = 1 \text{ km/s}$ for any tube with diameter $2a$ exceeding just a few km. Indeed, observations show flux tubes to be essentially vertical to the photosphere (except around emergence and collisional cancellation when magnetic curvature forces of the field arching from one polarity to the opposite one are strong).

concentrating to high latitudes the surface magnetic fields released by the decay of active regions, and dragging it down to the tachocline where shearing by differential rotation leads to the buildup of a new toroidal flux system, and thus to the onset of a new sunspot cycle. [...]

Figure 4.12 shows a series of meridional-plane snapshots of one such Babcock-Leighton dynamo solution, covering one sunspot cycle and starting approximately at sunspot maximum (based on magnetic energy as a proxy for sunspot number). Surface poloidal flux from the current cycle has begun to build up at low latitudes, and is rapidly swept to the pole, with polarity reversal of the polar field taking place shortly thereafter (panel B). As with the advection-dominated $\alpha\Omega$ solution discussed above, this solution is characterized by strong surface polar fields resulting from the poleward transport by the meridional flow of the poloidal component produced at lower latitudes, and the equatorward propagation of the toroidal field in the tachocline is also driven by the meridional flow. The turnover time of the meridional flow is here again the primary determinant of the cycle period. With $\eta = 30 \text{ km}^2 \text{ s}^{-1}$, this solution has a nicely solar-like half-period of 12.4 yr. All in all, this is once again a reasonable representation of the cyclic spatiotemporal evolution of the solar large-scale magnetic field.” {A:[54]}

A:54

There are yet other non-linearities that can be considered in solar/stellar dynamo models. For example, in [H-III:6.2.3] “the presence of stratification and rotation, a number of hydrodynamical (HD) and magnetohydrodynamical (MHD) instabilities associated with the presence of a strong toroidal field in the stably stratified, radiative portion of the tachocline can lead to the growth of disturbances with a net helicity, which under suitable circumstances can produce a toroidal electromotive force, and therefore act as a source of poloidal field. Different types of solar cycle models have been constructed in this manner. In nearly all cases the resulting dynamo models end up being described by something closely resembling the axisymmetric mean-field dynamo equations, the novel poloidal field regeneration mechanisms being once again subsumed in an α -effect-like source term appearing on the right-hand side of Eq. (4.18).” More on this in Sect. H-III:6.2.3. {A:[55]}

A:55

The models discussed thus far all lead to a steadily repeating magnetic cycle,

⁵⁴ Activity: For the curious: This work by Lemerle and Charbonneau (2017) describes an interesting dynamo experiment in which the Babcock-Leighton concept is combined with surface flux transport modeling (see Activity 51) to create a quasi-regular dynamo in which convection-induced fluctuations on the tilt angle of emerging active regions (perturbations on Joy’s rule, see Activity 52) provide the stochastic noise that can lead to cycle-to-cycle differences and even extended periods of weak cycling (as in the Maunder Minimum period for the Sun), something also reported on by Karak and Miesch (2017) in this paper.

⁵⁵ Activity: Make a summary of the essential distinctions between the dynamo concepts discussed up to this point: $\alpha\Omega$ (with or without α quenching, which itself can be strong/catastrophic or not); interface; flux-transport; and Babcock-Leighton.

where appropriate after an initial growth phase. The Sun, however, displays a rather erratic modulation of its activity from one cycle to the next, [H-III:6.3] “and certain aspects of the observed fluctuations may actually hold important clues as to the physical nature of the dynamo process.” Section H-III:6.3 discusses some of these processes, including those that could be responsible for long-term modulations of the solar cycle pattern (such as the Maunder Minimum): stochastic effects (with strong evidence from both models and observations for the importance of the scatter on tilt angles of active regions that reflect the influence of random convective flows during the rise of the flux to the surface), back reaction of the field on the flow patterns, and time delays through transport processes. Section H-III:6.3.5 discusses issues related to the forecasting of the solar cycle based on precursor signatures.

4.6 Dynamos in other stars

[H-III:6.4] “Figure 4.1 illustrates, in schematic form, the internal structure of main-sequence stars, more specifically the presence or absence of convection zones. A *G*-star like the Sun has a thick outer convection zone, spanning the outer 30% in radius. As one moves to lower masses, the relative thickness of the convective envelope increases until, somewhere around spectral type *M5*, stars become fully convective [(see Fig. 4.2 for an HR diagram and indications of spectral types)]. Moving from the Sun to higher masses, the convective envelope becomes ever thinner, until somewhere around spectral type *A0* it essentially vanishes. However, at around the same spectral type hydrogen [fusion] switches from the proton-proton (or *p-p*) chain to the CNO cycle, for which nuclear reaction rates are much more sensitively dependent on temperature. Core energy release becomes strongly depth-dependent, leading to convectively unstable temperature gradients. The resulting small convective core grows in size as one moves up to larger masses. In an early *B*-star of solar metallicity, the convective core spans the inner 25% or so in radius of the star.

Main-sequence stars of the *O* and *B* [spectral type] combine vigorous core-convection and high rotation rates, which makes dynamo action more than likely. This expectation has been amply confirmed by 3D MHD numerical simulations of dynamo action in the convective cores of massive stars. [...] All these core dynamo models have one thing in common: the large [diffusivity contrast $\eta_{\text{core}}/\eta_{\text{envelope}}$ between the convective core and the stably stratified envelope] leads to a ‘trapping’ of the magnetic field in the lower part of the radiative envelope, a direct consequence of the difficulty experienced by an externally-imposed magnetic field to diffusively penetrate a good electrical conductor [(analogous to, but here the inverse of, the ‘skin depth’ issue that was discussed in Section 4.5 for a cooler star) ...] This long-recognized property

of stellar core dynamos represents a rather formidable obstacle to be bypassed if the magnetic fields generated by dynamo action in convective cores are to become observable at the stellar surface [... In fact, in] a time-dependent situation where the core dynamo 'turns on' at or shortly before [a young star settles into a stable equilibrium represented by] the arrival on the zero-age main sequence, the time needed for the magnetic field to resistively diffuse to the surface can become larger than the star's main-sequence lifetime, for masses in excess of about $5 M_{\odot}$." (More on formation and evolution of stars in Ch. 10.)

[H-III:6.4.2] "Stars with spectral types ranging from late-*B* to early-*F* ^[viii] stand out as the least likely to support dynamo action, because they lack a convective region of substantial size. This squares well with various lines of observations; in particular, main-sequence *A*-stars are among the most 'magnetically quiet' stars in the HR diagram. A subset of late-*B* and *A* stars, namely the slowly-rotating, chemically peculiar *Ap/Bp* stars, do show strong magnetic fields, but even those show no sign of anything even mildly analogous to solar activity. The single pattern of temporal evolution noted is a decrease, by factors of 2 to 3, in the overall strength of the surface field, most prominent in the early stages of main-sequence evolution. This seems compatible with the idea of diffusive decay of residual higher-degree eigenmodes, and slow decreases associated with flux conservation as the stars slowly expand in the course of their main-sequence evolution. For these reasons, the fossil field hypothesis remains the favored explanatory model for the magnetic field of *Ap* stars. It is also quite striking that the high field strength observed in *Ap* stars (a few times 10^4 G), in magnetized white dwarfs ($\sim 10^9$ G), and in the most intensely magnetized neutron stars ($\sim 10^{15}$ G) all amount to [a] total surface magnetic flux $\sim 10^{27}$ Mx, lending support to the idea that these high fields can be understood from simple flux-freezing arguments [along an evolutionary timeline for these objects] (see Ch. H-I:3). [...]"

[H-III:6.4.3] "Until strong evidence to the contrary is brought to the fore, we are allowed to assume that late-type stars ^{viii} with a thick convective envelope overlying a radiative core host a solar-type dynamo. Observationally, a lot of what we know regarding dynamo activity in solar-type stars comes from the Mt. Wilson CaH+K survey[, a survey that focuses on a pair of strong resonance lines, which are known as the H and K lines, so named by Fraunhofer during early spectroscopic studies, and which were later found to be associated with singly ionized calcium; their signal reflects the chromospheric activity

^{viii} In stellar parlance, 'late' means 'cooler' and 'early' hotter. 'late-*B*' thus refers to *Bn*-type stars on the cooler side of the HR diagram, with digits *n* closer to 9 than to 0. 'Late type stars' is often used synonymously with 'cool stars', which refers to stars with convective envelopes immediately below their surface; see Fig. 4.2.

of a star.] Two important pieces of information can be extracted from these data, as constraints on dynamo models. The first is [that the overall level of Ca H+K emission, which is taken as a measure of overall magnetic flux in the photosphere,] is found to increase with rotation up to 5 – 10 times the solar rotation rate, after which saturation sets in (see Ch. 10). The second is of course the cycle period, for [the minority of stars that exhibit a regular cycle.]

The preponderance of strong magnetic field concentrated at high latitude in rapidly rotating solar-type stars (see Ch. 10) is also a potentially interesting discriminant. This can arise through channelling of buoyantly rising toroidal flux ropes along the polar axis [prior to surfacing], or efficient poleward transport of surface magnetic flux [after surfacing. . . .]” ^[ix]

[H-III:6.4.4] “With fully convective stars we encounter potential deviations from a solar-type dynamo mechanism; without a stably stratified tachocline and radiative core to store and amplify toroidal flux ropes, the Babcock-Leighton mechanism, the tachocline α -effect and the flux-tube α -effect all become problematic. Mean-field models based on the turbulent α -effect remain viable, but the dynamo behavior becomes dependent on the presence and strength of internal differential rotation, about which we really don’t know very much in stars other than the Sun. The full-sphere MHD simulations of an ‘M-star in a box’ are particularly interesting in this respect, as they indicate that fully convective stars do produce significant internal differential rotation and well-defined patterns of hemispheric kinetic helicity, both supporting the growth of a spatially well-organized large-scale magnetic component.

Moving to even cooler stars, as the luminosity drops and surface temperature falls below a few thousand K, the magnetic Reynolds number in the surface layers is expected to eventually fall back towards values approaching unity [because of the low degree of ionization at such temperatures]. Small-scale turbulent dynamo action may shut down, with magnetic activity then reflecting only the operation of a deep-seated, large-scale dynamo. Whether this transition is sharp or gradual, and whether it leads to well-defined observational signatures, remain open questions. There is certainly no *a priori* reason to presume that dynamo action should cease. Indeed, in some ways rapidly rotating very low-mass stars are getting closer to the physical parameter regime characterizing the geodynamo.”

^{ix} As computers continue to grow more powerful, 3D MHD dynamo simulations are advancing towards generating cycling large-scale fields in modeled stellar convection zones. An entry point for that literature is provided, for example, in the by Charbonneau (2014). His contribution to Living Reviews in Solar Physics (Charbonneau, 2010) provides a description of advanced Babcock-Leighton type models that can now take observed magnetograms to provide forecasts of long-term trends of solar activity.

4.7 Dynamos in terrestrial planets

[H-III:7.5] “Planetary dynamos share with stellar dynamos that the basic physical concept for their description is that of convection-driven magnetohydrodynamic flow in a rotating spherical shell combined with the associated magnetic induction effects. [...]

Inside a shell of depth d with an electrical conductivity σ_e the fluid must move with a sufficiently large characteristic velocity v_t , so that the magnetic Reynolds number [in Eq. (3.18)] exceeds a critical value $\mathcal{R}_{m,crit}$ in order to have a self-sustained dynamo. The flow pattern must also be favorable for dynamo action, which requires a certain complexity. In particular helical (corkscrew-type) motion with a large-scale order in the distribution of right-handed and left-handed helices is suitable. The Coriolis force plays a significant part in the force balance of the fluid motion and influences the pattern of convection. With this the requirement for ‘flow complexity’ seems to be satisfied and self-sustained dynamo action is possible above $\mathcal{R}_{m,crit} \approx 40 - 50$.

At greater depth in the solar convection zone, the magnetic Reynolds number reaches values of order 10^9 for molecular values of the magnetic diffusivity. In the geodynamo \mathcal{R}_m is approximately 1000. This fairly moderate value allows for the direct numerical simulation of the magnetic field evolution without the need to use an ‘effective diffusivity’ or a parameterization of the induction process through a turbulent α -effect. [...]

The density in the Sun varies by many orders of magnitude and the convection region spans many density scale heights. The density changes associated with radial motion are thought to be important. Flow helicity [$(\mathbf{v} \cdot (\nabla \times \mathbf{v}))$] arises in the Sun because of the action of the Coriolis force on rising expanding and sinking contracting parcels of plasma. Strong magnetic flux tubes have their own dynamics, because the reduction of fluid pressure that compensates magnetic pressure reduced their density and makes them buoyant. In contrast, the dynamo region in Jupiter covers approximately one density scale height and much less in terrestrial planets. The two compressibility effects mentioned before probably do not play a significant role in planetary dynamos. Present geodynamo models usually neglect the small density variation and assume incompressible flow in the Boussinesq approximation (where density differences are only taken into account for the calculation of buoyancy forces; see Activity 29).

Many models of the solar dynamo assume that most of magnetic field generation occurs at the tachocline, the shear layer between the radiative deep interior and the convection zone of the Sun. For planetary dynamos the process of magnetic field generation is thought to occur in the bulk of the convecting layer. [...] The relevant equation of motion for an incompressible

fluid [in a corotating frame of reference] is

$$\rho \frac{\partial \mathbf{v}}{\partial t} + \rho(\mathbf{v} \cdot \nabla) \mathbf{v} + 2\rho\Omega \hat{\mathbf{e}}_z \times \mathbf{v} = \rho\zeta T g \hat{\mathbf{e}}_r - \nabla p' + \frac{1}{4\pi} (\nabla \times \mathbf{B}) \times \mathbf{B} + \rho\nu \nabla^2 \mathbf{v}, \quad (4.24)$$

where \mathbf{v} is velocity, Ω rotation rate, ρ density, p' non-hydrostatic pressure, ν kinematic viscosity, ζ thermal expansivity, g gravity, T temperature, \mathbf{B} magnetic field, r radius and z the direction parallel to the rotation axis. The terms in Eq. (4.24) describe, in order, the linear and non-linear parts of inertial forces, Coriolis force, buoyancy force, pressure gradient force, Lorentz force, and viscous force [(compare with Eq. (3.5))].

In the non-magnetic and rapidly rotating case, the primary force balance is between the pressure gradient force and the Coriolis force (geostrophic balance), similar as for large-scale weather systems in the Earth's atmosphere. [Assuming a stationary flow, and ignoring all other terms on the right of] Eq. (4.24) and taking the curl, we arrive at the Taylor-Proudman theorem, which predicts the flow to be two-dimensional with $\partial \mathbf{v} / \partial z = 0$. The only type of perfectly geostrophic flow in a sphere, *i.e.*, a flow that satisfies this condition, is the differential rotation of cylinders that are co-aligned with the rotation axis (geostrophic cylinders). Such a flow can neither transport heat in the radial direction, nor can it act as a dynamo. Convection requires motion away from and towards the rotation axis. This must violate the Taylor-Proudman theorem, because a column of fluid that is aligned with the z -direction will then stretch or shrink because it is bounded by the outer surface of the sphere. Hence the velocity cannot be independent from z . The necessity to violate the Taylor-Proudman theorem inhibits convection and requires that some other force, such as viscous friction, must enter the force balance. In order for viscosity to do so, the length scale of the flow must become small, at least in one direction. But the flow maintains a nearly geostrophic structure as far as possible. At the onset of convection it takes the form of columns aligned with the rotation axis (Fig. 4.8). They surround the inner core tangent cylinder like pins in a roller bearing. The tangent cylinder is parallel to the z -axis and touches the inner core at the equator. It separates the fluid core into dynamically distinct regions.

The primary circulation is around the axes of these columns. However, in addition there is a net flow along the column axes which diverges from the equatorial plane in anticyclonic vortices and converges towards the equatorial plane in columns with a cyclonic sense of rotation. The combination implies a coherently negative flow helicity in the northern hemisphere and positive helicity in the southern hemisphere, [which] can serve as an efficient dynamo of the α^2 -type.

When the motion becomes more vigorous at highly supercritical convection

and when a strong magnetic field is generated, other forces such as inertia (advection of momentum) and the Lorentz force can affect the flow. However, one difference between the solar dynamo and planetary dynamos is the different role of inertial forces versus the Coriolis force. Their ratio is measured by the Rossby number (Eq. 4.2). Deep in the solar convection zone $N_R \approx 1$ when the pressure scale height is taken for L_t . With typical estimates for the flow velocity in the Earth's core (1 mm s^{-1}), the Rossby number is of order 10^{-6} when a global scale such as the core radius or shell thickness is used for L_t . Therefore, fluid motion in the geodynamo is often considered to be largely unaffected by inertial forces. The general force balance is believed to be that between Coriolis force, pressure gradient force, Lorentz forces and buoyancy forces. However, at small scales inertial forces may become important also in planetary dynamos and can potentially feed back on the large scale flow.

Like rotation, the presence of an imposed uniform magnetic field inhibits convection in an electrically conducting fluid. However, the combination of a magnetic field and rotation reduces the impeding influence that either effect has separately. This constructive interference is most efficient when the Coriolis force and the Lorentz force are in balance. [... Applied to dynamos, it is argued that when the Coriolis force exceeds the Lorentz force the field will strengthen, and when the Lorentz force exceeds the Coriolis force the convection will weaken. Hence, it is assumed that the field equilibrates when the forces match (referred to as a magnetostrophic balance). The field strength inside the geodynamo or in Jupiter's dynamo seems to agree with that argument.] However, numerical dynamo simulations put some doubt on its validity." More on that in Ch. H-III:6.

[H-III:7.6.3] "The stretching of magnetic field lines by differential rotation in the case of the solar dynamo, particularly at the tachocline, is thought to be of major importance for the generation of a toroidal magnetic field that is much stronger than the poloidal field. In most geodynamo models, in contrast, differential rotation does not contribute much to the total kinetic energy and the toroidal and poloidal magnetic field components have similar strength. As mentioned before, the flow is strongly organized by rotational forces and the vortices are elongated in the z -direction. Even at a highly supercritical Rayleigh number [(which measures the time scale of conductive relative to convective transport)] and in the presence of a strong magnetic field, the flow outside the inner core tangent cylinder is reminiscent of the helical convection columns found at onset. Inside the tangent cylinder, the flow pattern is different and often exhibits a rising plume near the polar axis (Fig. 4.13b). {A:[56]} The

A:56

⁵⁶ Activity: The dynamo model in Fig. 4.13 is characterized by five dimensionless numbers. Two, the magnetic Reynolds number and the Rossby number, are defined in Eqs. (3.18) and (4.2), respectively. Look up the meaning of the other three: Ekman, magnetic Prandtl, and Rayleigh. These three

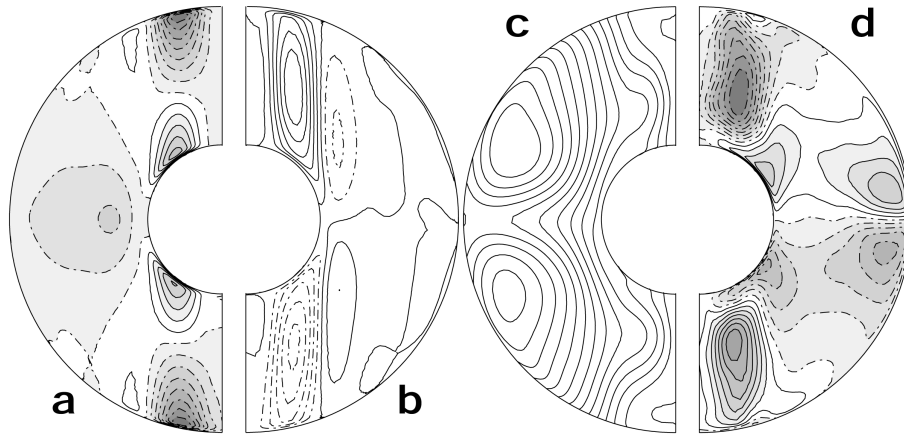


Fig. 4.13. Time-averaged axisymmetric components of velocity and magnetic-field components for a planetary dynamo model with [Rayleigh number] $R_a^* = 0.225$, [Ekman number] $E = 3 \times 10^{-4}$, $P_r = 1$, [magnetic Prandtl number] $P_m = 3$, [Reynolds number] $R_m \approx 250$ and [Rossby number] $N_R \approx 0.1$. The grey-scale indicates absolute intensity. (a) Azimuthal velocity, broken lines are for retrograde flow, (b) streamlines of meridional velocity, full lines for clockwise circulation, (c) poloidal magnetic field lines, (d) azimuthal (toroidal) magnetic field, broken lines westward directed field. [Fig. H-III:7.8]

plume is accompanied with a strong vortex motion (called a ‘thermal wind’) with a retrograde sense of rotation near the outer surface changing to prograde rotation at depth (Fig. 4.13a), because the Coriolis force acts on the associated converging flow near the inner core boundary and diverging flow near the outer boundary.

[...] There is general agreement that the axial dipole field is generated from the axisymmetric toroidal field by an α -effect associated with the helical flow in the convection columns outside the tangent cylinder. In mean-field theory as it is used in astrophysics, the α -effect is associated with unresolved turbulent eddies. In the geodynamo models a ‘macroscopic’ α -effect is observed.

The mechanism for generating the axisymmetric toroidal field is less clear and both an α -effect and differential rotation seem to play a role. Often two flux bundles in the azimuthal direction are found outside the tangent cylinder, with opposite polarity north and south of the equatorial plane (Fig. 4.13d). [T]hey are generated from the axisymmetric poloidal field by a similar macroscopic α -effect associated with the helical convection columns (α^2 -dynamo). Other

are important numbers when computing the flows and their coupling, but not encountered until this Figure in this text because the solar dynamo models that were discussed and shown in the figures are kinematic, relying on a given, not consistently computed, flow pattern. The dynamo model behind Fig. 4.13, in contrast, computes flow, field, and their interaction, and thus also needs these remaining three dimensionless numbers specified.

authors show that the Ω -effect (the shearing of poloidal field by differential rotation) contributes strongly to the generation of axisymmetric toroidal field, even though the kinetic energy in the differential rotation is rather limited. While in weakly driven numerical dynamo models the regions inside the tangent cylinder, north and south of the inner core, are nearly quiescent, vigorous flow is found here in more strongly driven models. In these cases a strong axisymmetric toroidal field is found inside the tangent cylinder region, produced by the shearing of poloidal field lines in the polar vortex (Fig. 4.13a,c,d).”

A:57

{A:^[57]}

⁵⁷ Activity: Summarize the contrast between the dynamo of a terrestrial planet with that of stars as discussed in this Ch. 4: consider, among others, flow speed, rotation period, stratification, differential rotation, and meridional advection.

5

Flows, shocks, obstacles, and currents

5.1 Introductory overview

Much of what happens in the heliosphere originates in the interaction of an object or a plasma-filled volume with flows of magnetized plasma directed at it. The scales of the phenomena discussed here range from comets and asteroids up to the entire heliosphere where the solar wind couples to the interstellar medium. The interactions often involve shocks, such as in cases when a fast solar wind stream catches up with a significantly slower one, or where the solar wind envelops a magnetosphere. In other settings, they may involve smooth sub-Alfvénic adjustments in the flow, such as happens around many of the moons orbiting within the plasma-filled magnetospheres of the giant planets. The flow of magnetized plasma around a body may be affected by a magnetic field that it induces in that body's conducting deep interior or near-surface shell, or that induced field may add to an already present dynamo-generated field. The flow may pick up matter from an object's outer atmosphere through reconnection processes, or through ionization of neutral matter that enters it from outside, such as happens when the solar wind engulfs a comet or because of interstellar-medium neutrals entering the heliosphere.

Despite this great variety of conditions, common patterns emerge. These are the focus of this chapter which reviews the effects of a plasma flow around objects from two different perspectives. One is to look into what happens to the external flow, the other is concerned with what happens to the atmosphere or magnetic field of the body that is cocooned by that flow. The first can be summarized by looking at what response is induced in the body by the magnetized external flow. Electrodynamics teaches us that the moving external magnetized plasma induces a current system in a conductor. In the extreme of a perfect conductor, that induced current corresponds to a magnetic field that counters the continuation of the external field inside the conductor so that there is no net field there. In the opposite extreme of a perfect insulator no

field is induced and the external field permeates the body as in its surroundings. Intermediate conductivity induces an intermediate response.

Should the conductivity be limited to part of the body, the resulting field external to that body provides clues as to the location of the conducting medium and the magnitude of the conductivity. For example, a planetary system body with a relatively small conducting core that is enveloped by a non-conducting shell will have an induced field near that core that is comparable, but opposite, to the external field. The strength of that field (generally largely bipolar when a large-scale field flows by the object) decreases through the envelope towards the body's surface so that the external signature may be weak. If there is a conducting layer in a body that lies at or just above or below the surface (such as an electrolyte-laden ocean or a substantial ionosphere) the induced field can be strong if the conductance is high, leading to a net field outside the body that is markedly distinct from that of the incoming plasma. If the body has a substantial, sustained intrinsic magnetic field (*i.e.*, a dynamo in its interior, and thus a conducting volume somewhere in its interior) an induced field distorts the intrinsic field. For the incoming flow, the 'object' that it encounters is bounded by the permanent and/or induced field, and thus lies anywhere between the body's surface (or high in the atmosphere, if it has one) or where the magnetic field of that body is strong enough to withstand the momentum of the incoming flow of ionized, magnetized matter.

The consequence of the dynamo and induced field for the incoming flow is communicated by waves. If the incoming flow is relatively slow, specifically if it is sub-magnetosonic (see Sect. 3.3), the flow can be deflected well before approaching the 'object' and largely flow around it; how much flows around it depends on the magnetic field. If, in contrast, the flow comes in super-magnetosonically, the incoming flow is unaware of the object until forced to realize its presence, either (essentially) at the surface of the body or where its magnetic field is strong enough; there will be a shock that decelerates and deflects the flow. The interplay between the incoming field and the body's magnetic field through compression and reconnection drives magnetospheric and ionospheric processes, resulting in much of what we know as space weather, further modified by planetary rotation.

But 'obstacles' to flows in the local cosmos are not limited to planets, moons, and other bodies. One example of another type of obstacle is the outflowing wind from a Sun-like star as it is encountered by the interstellar medium. Another is that of (relatively) fast wind streams and coronal mass ejections that catch up with slower wind plasma ahead of them; such pileups often include shocks, here with the potential of plasma becoming compressed in the collision zone rather than flowing around the 'obstacle' because the scales

are such that there generally is no way 'around' the 'obstacles' within the characteristic time scale of the passage of such flows through a large part of the heliosphere.

Generally, when a flow interacts with another volume of magnetized plasma, the magnetic field is distorted in both volumes (if a field exists, which in heliophysics is commonly the case) and magnetic (Lorentz) forces play into the balance of forces, as expressed in the momentum equation Eq. (3.5). One can view this as a result of the pressure and tension forces ascribed to the magnetic field or, equivalently, to induced currents - the equations do not care about our perspective in this matter (see Ch. 3).

One key differentiating factor in how the flow and the enveloped volume interact is whether the magnetic fields in these two domains can connect or not. In the ideal-MHD approximation, in which effects of resistivity are ignored, the induction equation Eq. (3.3), through the frozen-in flux paradigm, leads to the conclusion that the two plasmas involved cannot interpenetrate: the flow moves around the impacted plasma as a wind that flows past a solid object. This can still lead to very complex dynamics, as diverse as for a wind flowing past a flag or around a supersonic jet-plane. If the magnetic field can reconnect, however, the plasmas can interact in entirely different ways, that include, for example, a variety of magnetospheric phenomena. The differentiator here is not solely the plasma resistivity but the effect of such resistivity within the interaction time scale of the flow passing by the enveloped volume as expressed by the magnetic Reynolds number (Eq. 3.18). The geometry of the interaction is also set by the Alfvén Mach number (Eq. 3.32).

In Sect. 11.2.2 we encounter another type of flow into a magnetosphere. In very young stars that are still surrounded by a gaseous disk, matter spirals gradually towards a still growing star. Close to the star, this accreting matter will diffuse into, and then be locked onto the magnetic field of the star. This allows the material to 'fall' through the stellar magnetosphere, while being channeled by the magnetic field, to end up near the surface of the star in what are known as 'accretion columns' (sketched in Fig. 11.7). But that is for later.

This chapter takes you through the following situations throughout the heliosphere: {A:^[58]}

A:58

- Sects. 5.2, 5.3 and 5.4 are *introductions*: they discuss, respectively, **low-**

⁵⁸ Activity: To build a comparison of the different conditions, keep pen(cil) and paper at hand to sketch the various configurations as you read about how flows interact with bodies in the planetary system. Working in a reference frame in which the body is at rest, assume a spherical object, and let a flow move past it from left to right. Then prepare to make drawings in two orthogonal planes: the first plane is defined by the flow vector and the magnetic field carried in the flow (you may assume the field to be normal to the flow), while the second plane is normal to the first. Draw streamlines of the flow and subsequently add magnetic field lines. If you are good at 3-D renderings, also try a visualization such as in Fig. 5.11.

velocity interactions versus shocks, the **elementals of shocks and discontinuities**, and the **magnetized solar wind and the Parker spiral** that forms as the wind flows out from the rotating Sun. Staying with the *solar wind*, Sect. 5.5.1 reviews **solar-wind stream interactions**.

- Next come discussions of *flows around bodies in the heliosphere*, beginning in Sect. 5.5.2 with a **non-conducting body without atmosphere**, then in Sect. 5.5.3 a **flow around a conducting body** without an intrinsic magnetic field.
- By Sect. 5.5.4 we reach *bodies with dynamos* and look at **plasma flow around a permanently magnetized body**, after which we can discuss *magnetospheres*: first, in Sect. 5.5.5 a **closed magnetosphere** which exists only in the world of ideal MHD, but then in Sect. 5.5.6 we introduce the **open magnetosphere** such as happens in the real world.
- Finally, we move on to what happens *within the magnetosphere of a planet or moon* as a consequence of the variable solar wind coupling to the body's magnetic field: in Sect. 5.5.7 we talk about the overall system of **solar wind-magnetosphere-ionosphere interaction**, including the effects of rotation and advection.
- Finally, we return to the solar wind, looking at the outermost regions of the heliosphere, where the outflow meets the interstellar medium: Sect. 5.5.8 explores what happens when we have a **flow impinging on a fast outflow** but now on scales such that the flow can find a way around the outflowing plasma, in contrast to what happens in the case of wind streams interacting with comparable scales discussed in Sect. 5.5.1.

5.2 Low-velocity interactions versus shocks

In view of the great diversity in conditions encountered throughout the local cosmos [H-IV:10.2] “it may seem unlikely that general rules can describe the interaction regions. We are rescued from the need to treat each case as totally distinct by recognizing that physical theories often incorporate a small set of dimensionless parameters that control important aspects of a system, even if such properties as spatial scale, temperature, and flow velocity vary by many orders of magnitude. For a flowing plasma incident on an obstacle, the form of the interaction depends critically on how the flow speed is related to the speed of waves that transmit information about changes of plasma properties from one part of the system to another. An analogy to waves in neutral gases helps to clarify the concept. In the frame of an airplane in flight, the atmosphere flows onto the plane at some velocity, call it v . As the gas encounters the plane, pressure perturbations develop. Pressure perturbations launch sound waves that travel at the sound speed, c_s . If such waves can move away [in the

Table 5.1. *Properties of the plasmas upstream of various planets and other bodies of the Solar System. [Listed are the Alfvén and magnetosonic Mach numbers and the plasma β (Eq. 3.24). Table H-IV:10.1].*

Obstacle	Ambient plasma	M_A	M_{ms}	β
Io, Europa, Ganymede	jovian magnetosph.	< 1	< 1	> 1
Asteroids	solar wind	> 1	> 1	~ 1
Comets	solar wind	> 1	> 1	~ 1
Moon	Earth's magnetosphere or solar wind	either > 1 or < 1	either > 1 or < 1	~ 1 or < 1

forward direction] from the plane, they can divert the atmosphere upstream of the plane. But the waves are swept back toward the plane at the flow speed of the plasma. Only if $v < c_s$ is it possible for the waves to begin to divert the atmosphere well upstream of the plane. If $v > c_s$, as for a supersonic jet, the waves pile up in front of the plane, causing a shock to develop upstream. Only downstream of the shock is the flow diverted. Assuming that the plane is large compared with distances characteristic of atmospheric properties, the parameter that determines whether or not a shock will form is the (dimensionless) sonic Mach number of the surrounding atmosphere, v/c_s . [Shocks are described in Sect. 5.3.]

In a plasma, much as in a neutral gas, compressional perturbations develop when there is an obstacle in the flow. [...] Having identified [in Section 3.3] some of the waves that carry information through a magnetized plasma, we are now able to introduce the dimensionless parameters that help us understand aspects of flow and field perturbations. The magnetosonic Mach number (M_{ms}) is the ratio of the flow speed to the fast mode speed, taken as $(c_s^2 + v_A^2)^{1/2}$. M_{ms} reveals whether or not a shock is likely to form upstream in the flow. When $M_{ms} < 1$, compressional waves can travel upstream from the obstacle faster than the flowing plasma can sweep them back. These waves, moving upstream, can divert the incident flow around the obstacle, much as the bow wave of a ship diverts water to the sides, and no shock develops. However, as in the situation discussed in the context of supersonic flight, if $M_{ms} > 1$, compressional waves are unable to propagate upstream faster than they are swept back by the flow. They pile up to form a shock. Most bodies in the super-magnetosonic solar wind [...] create shocks standing somewhat upstream on their sunward sides. Downstream of the shock, plasma is heated, compressed, and diverted around the obstacle.

The Alfvén Mach number (M_A , [Eq. 3.32]) is the ratio of the speed with which the ambient plasma flows towards an obstacle divided by the Alfvén

speed. We will see that this quantity controls the shape of the interaction region in planes containing the unperturbed plasma flow and the background magnetic field. The plasma beta (β , [see Eq. 3.24]) is the ratio of the thermal pressure to the magnetic pressure. This quantity enables us to understand how significantly the magnetic field structure can be modified by changes of the plasma pressure.

The plasma environment differs greatly among the small bodies of the Solar System. Some of the bodies are embedded in the solar wind, others in the plasma of a planetary magnetosphere, and some [...] move from one environment to another [(such as Earth's Moon, which spends part of each lunar orbit in Earth's magnetotail and the rest of the month in the solar wind)]. Table 5.1 lists some plasma properties relevant to the environment of selected bodies."

5.3 Elementals of shocks and other discontinuities

Shocks that we discussed up to here develop when a flow speed exceeds the speed of waves that can serve as a warning to the flow that an obstacle lies ahead. Shocks can also form if non-linear effects in the propagation of a wave become important, such as when a wave runs into a medium in which strong gradients in density or temperature cause the wave amplitude to grow more rapidly than dissipation can limit that growth; examples of such shock waves are found in upward traveling pressure waves in atmospheres, including the Earth's and the Sun's, and also in very large and long-lived wind streams in the heliosphere.

[H-II:7.2] "In the small-amplitude limit, the profile of a magnetohydrodynamic (MHD) wave does not change as it propagates, but even a small-amplitude wave will eventually distort due to *wave steepening*. The wave steepening happens when gradients of pressure, density and temperature become so large that dissipative processes (*e.g.*, viscosity, thermal conduction) are no longer negligible. In the steady state, a steady wave-shape – a *shock wave* – is formed in which the steepening effect of nonlinear advective terms balance the broadening effects of dissipation. The shock waves move at speeds larger than the ambient intrinsic speed, which for magnetized ionized matter in the heliosphere, is the magnetosonic speed. If the shock moves much faster than the magnetosonic wave, it is called a strong shock; if it moves just slightly faster, it is called a weak shock. The dissipation inside the shock front leads to a gradual conversion of the energy being carried by the wave into heat. In the heliospheric plasma, we have collisionless shocks in which the thermalization happens through wave-particle interactions. [...]

A propagating wave described by the ideal fluid equations leads to infinite

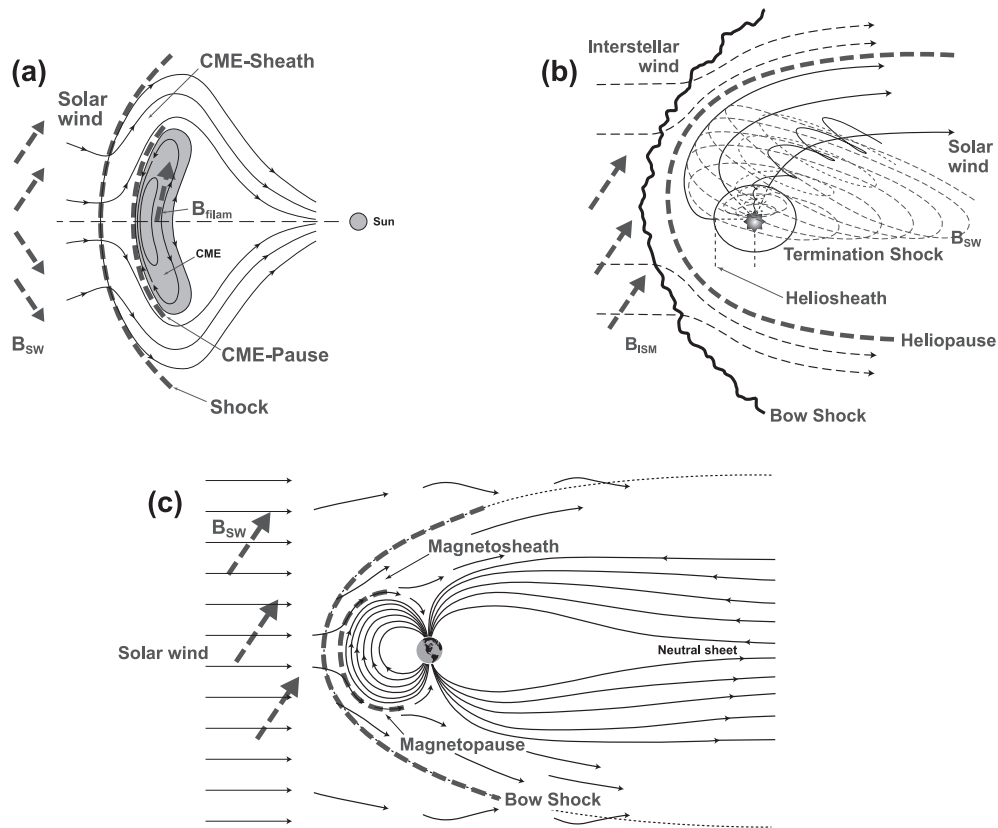


Fig. 5.1. Schematic comparison of shocks around [coronal mass ejections (CMEs, magnetically-driven explosions from the solar atmosphere into the heliosphere)], the heliosphere, and the terrestrial magnetosphere. The figure shows some of the types of shocks and sheaths that exist in the heliosphere and their universal basic structures: (a) a CME; (b) the outer heliosphere, and (c) Earth's magnetosphere. The same basic structures appear: shocks where the solar wind becomes subsonic; the sheaths that separate the subsonic solar wind from the obstacle ahead; and the 'pause' where there is a pressure equilibrium between the subsonic solar wind and the obstacle's environment. In the case of a CME these three structures are the shock, CME-sheath, CME-pause and the obstacle is the magnetic filament that drives the CME. In the case of the outer heliosphere the structures are the termination shock, heliosheath, and heliopause. The obstacle is the interstellar wind and the magnetic field it is carrying. If the interstellar wind is supersonic there is an additional shock, the bow shock. In the case of the Earth's magnetosphere the structures are the shock, the magnetosheath, the magnetopause and the obstacle is the Earth's dipolar magnetic field. [Fig. H-II:7.1]

gradients in a finite time. There is no solution for the ideal MHD equations. This is not surprising: ideal equations are valid when scales of variations are larger than the mean free path. The breakdown in ideal equations occurs in a very thin region, while the fluid equations are valid everywhere else. In this

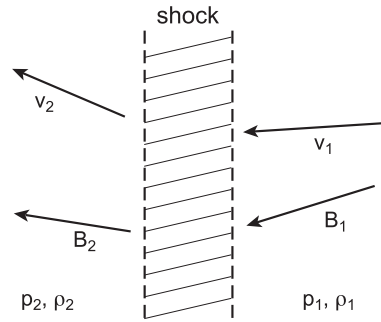


Fig. 5.2. Diagram showing the region upstream (left) and downstream of a shock. [Fig. H-II:7.3]

very thin region, it is difficult to describe the plasma in detail. The simple picture is a discontinuity dividing two roughly uniform fluids. An important aspect is that the simple picture of a discontinuity dividing two roughly uniform fluids is not usually applicable in a plasma. Shocks can involve turbulence for example. For this initial discussion, we make the simplifying assumption that there is a planar discontinuity of zero thickness that separates two uniform fluids, as depicted in Figure 5.2. We also assume that the shock is stationary [or, in other words, that we are in the co-moving frame of reference . . .] The transition must be such as to conserve mass, magnetic flux, and energy. The MHD jump conditions are independent of the physics of the shock itself and are known as the *Rankine-Hugoniot jump conditions*.”

[H-II:7.3] “It is straightforward to obtain the Rankine-Hugoniot jump conditions from [Maxwell’s equations and] the MHD equations. Assuming steady state in the frame of reference of the shock, the equation for the conservation of mass [in Eq. (3.4) in the absence of sources and sinks,] gives

$$\rho_1 \mathbf{v}_1 \cdot \hat{\mathbf{e}}_{\perp} = \rho_2 \mathbf{v}_2 \cdot \hat{\mathbf{e}}_{\perp}, \quad (5.1)$$

[(where $\hat{\mathbf{e}}_{\perp}$ is a unit-length vector pointing in the direction normal to the shock)] or in a different notation

$$\{\rho \mathbf{v} \cdot \hat{\mathbf{e}}_{\perp}\} = 0, \quad (5.2)$$

where the symbol $\{\dots\}$ represents differences between the two sides of the discontinuity.

Conservation of momentum, [with Eq. (3.5) without sources, sinks, or viscosity,] yields

$$\left\{ \rho \mathbf{v} (\mathbf{v} \cdot \hat{\mathbf{e}}_{\perp}) + \left(p \hat{\mathbf{e}}_{\perp} + \frac{B^2}{8\pi} \hat{\mathbf{e}}_{\perp} - \frac{(\mathbf{B} \cdot \hat{\mathbf{e}}_{\perp})}{4\pi} \mathbf{B} \right) \right\} = 0. \quad (5.3)$$

Conservation of energy, [...] results in

$$\left\{ \left(\frac{1}{2} \rho v^2 + \frac{\gamma p}{\gamma - 1} \right) (\mathbf{v} \cdot \hat{\mathbf{e}}_{\perp}) + \frac{c}{4\pi} (\mathbf{E} \times \mathbf{B}) \cdot \hat{\mathbf{e}}_{\perp} \right\} = 0. \quad (5.4)$$

[Note that $\mathbf{S} = (c/4\pi)\mathbf{E} \times \mathbf{B}$ is the Poynting flux, which measures the directional energy transfer in an electromagnetic field; compare with Eq. (4.1) where that is expressed for a plasma with infinite conductivity, as it is below in Eq. (5.10).]

Conservation of magnetic flux, [...] gives

$$\{\mathbf{B} \cdot \hat{\mathbf{e}}_{\perp}\} = 0. \quad (5.5)$$

The equation

$$\nabla \times \mathbf{E} = -\frac{1}{c} \frac{\partial \mathbf{B}}{\partial t} \quad (5.6)$$

[for a steady state] can be written as

$$\{\mathbf{E} \times \hat{\mathbf{e}}_{\perp}\} = \mathbf{0}. \quad (5.7)$$

Let us consider, now, the normal \perp and the tangential \parallel components relative to the shock's surface so that the jump conditions can be written as:

$$\left\{ \rho v_{\perp}^2 + p + \frac{B_{\parallel}^2}{8\pi} \right\} = 0 \quad (5.8)$$

$$\left\{ \rho \mathbf{v}_{\parallel} v_{\perp} - \frac{\mathbf{B}_{\parallel} B_{\perp}}{4\pi} \right\} = 0 \quad (5.9)$$

$$\left\{ \left(\frac{1}{2} \rho v^2 + \frac{\gamma p}{\gamma - 1} + \frac{B^2}{4\pi} \right) v_{\perp} - (\mathbf{v} \cdot \mathbf{B}) \frac{B_{\perp}}{4\pi} \right\} = 0 \quad (5.10)$$

$$\{B_{\perp}\} = 0 \quad (5.11)$$

$$\{\mathbf{v}_{\perp} \times \mathbf{B}_{\parallel} + \mathbf{v}_{\parallel} \times \mathbf{B}_{\perp}\} = \mathbf{0}. \quad (5.12)$$

Equations [(5.2) and] (5.8)–(5.12) are the Rankine-Hugoniot jump conditions that describe all types of shocks” and also allow for three types of discontinuities that are not shocks. An example of the heating associated with shocks is given in Fig. 5.3. [x] {A:[59]}

A:59

⁵⁹ Activity: Write the Eqs. (5.2) and (5.8)–(5.12) for the hydrodynamic limit, and derive the temperature ratio between the post- and pre-shock media. You should find that the density contrast $r_{\rho} = ((\gamma + 1)M_s^2)/(2 + (\gamma - 1)M_s^2)$ and the pressure ratio $r_p = (2\gamma M_s^2 - (\gamma - 1))/(\gamma + 1)$ where $M_s = v_1/c_{s1}$ for sound speed c_s . Note there is a maximum value for r_{ρ} but not for r_p as function of M_s . What are the values for r_{ρ} and r_p for $\gamma = 5/3$ for $M_s \downarrow 1$ and $M_s \gg 1$?

^x A note on terminology: a “parallel shock” propagates along the magnetic field, *i.e.*, has the vector $\hat{\mathbf{e}}_{\perp}$ normal to the shock front aligned along the magnetic field, or $\hat{\mathbf{e}}_{\perp} \parallel \mathbf{B}$. A “perpendicular shock” has $\hat{\mathbf{e}}_{\perp} \perp \mathbf{B}$.

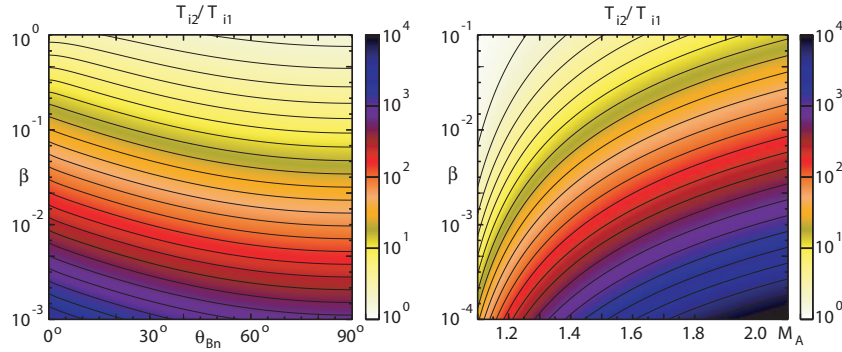


Fig. 5.3. Iso-contours of shock heating, expressed as the ratio between downstream to upstream ion temperature T_{i2}/T_{i1} , as a function of shock-normal angle $\theta_{B\perp}$ [(the angle between the shock normal and the upstream magnetic field)] (fixed $M_A = 2$) and Alfvén Mach number M_A (fixed $\theta_{B\perp} = 45^\circ$) for low β plasmas. Derived from standard Rankine-Hugoniot conditions for fast shocks, assuming a specific heat ratio $\gamma = 5/3$. The graphs show that for a wide range of angles, there can be very substantial downstream heating at sufficiently low plasma β , as present in much of the solar corona. Such extreme heating may help form a seed population for further acceleration [that will be discussed in Ch. 8]. [Fig. H-II:8.1]

[H-II:7.4] “Discontinuities can be classified as either contact or rotational discontinuities. Contact discontinuities happen when there is no flow across the discontinuity, *i.e.*, $v_\perp = 0$, but $\{\rho\} \neq 0$. A classic example is the contact discontinuity of a mix of vinegar and olive oil. If $\{B_\perp\} \neq 0$ at a contact discontinuity then only the density changes across the discontinuity, which is rarely observed in plasmas. A tangential discontinuity occurs when $\{B_\perp\} = 0$, then $\{v_\parallel\} \neq 0$ and $\{B_\parallel\} \neq 0$ and $\{p + B^2/8\pi\} = 0$. This means that the fluid velocity and magnetic field in this case are parallel to the surface of the discontinuity but change in magnitude and direction, and that the sum of thermal and magnetic pressures is constant. [... Large heliophysical examples of tangential discontinuities with $\{B_\perp\} = 0$ are the heliospheric current sheet and the magnetospheric current sheet (illustrated in Fig. 5.4).]

A rotational discontinuity occurs when $\{v_\perp\} \neq 0$ and $\{\rho\} = 0$. From the jump conditions this implies that $\{v_\perp\} = 0$ and $\{p + B_\parallel^2/8\pi\} = 0$ so $\mathbf{v}_1 \cdot \hat{\mathbf{e}}_\perp = \mathbf{v}_2 \cdot \hat{\mathbf{e}}_\perp = v_\perp$ and $\rho_1 = \rho_2$. After some math, we find that $v_\perp^2 = B_\perp^2/4\pi\rho$, and that B_\parallel remains constant in magnitude but rotates in the plane of the discontinuity. [...]

The Rankine-Hugoniot jump conditions have 12 variables. Four upstream parameters are specified (ρ , v , B_\parallel , and B_\perp), so we have 7 equations for

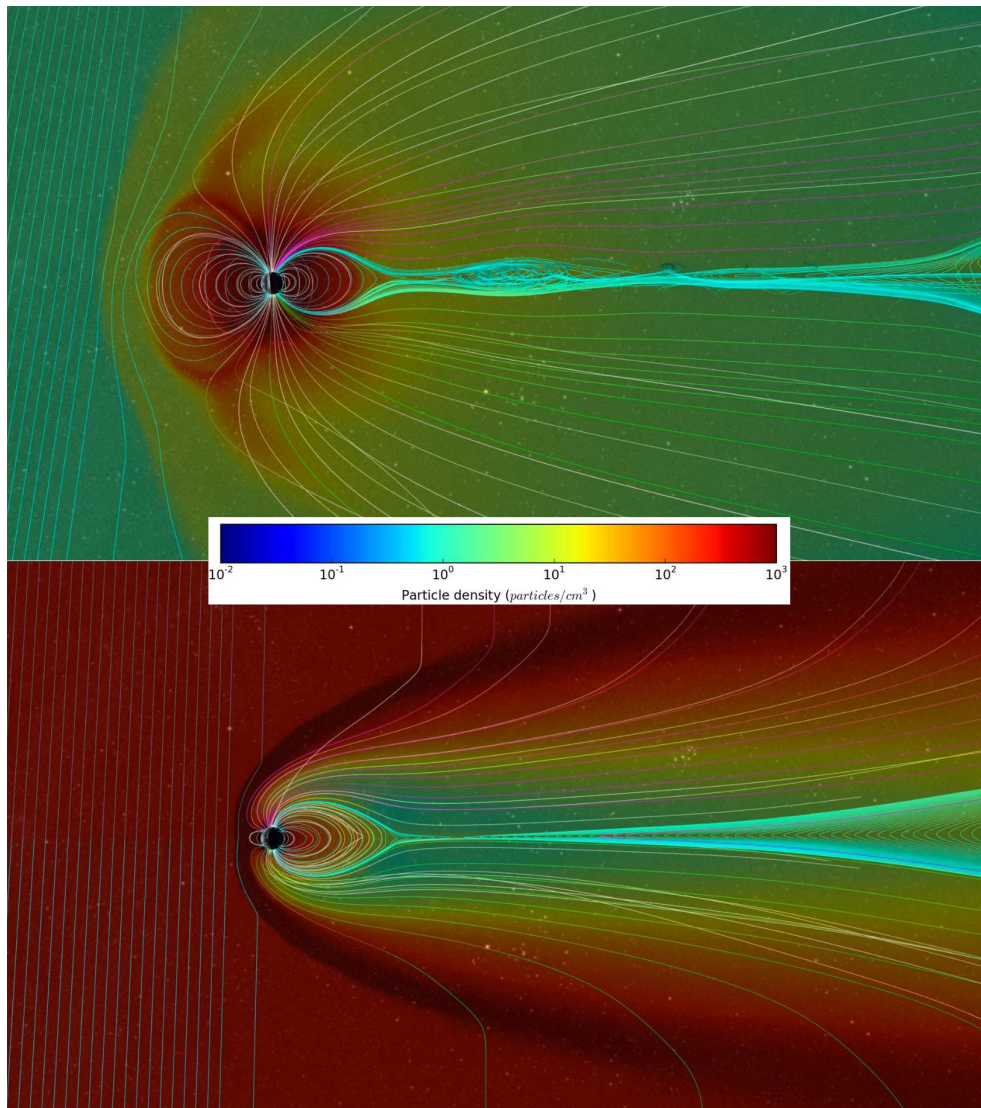


Fig. 5.4. An MHD simulation of the Earth’s magnetic field subject to (top) a typical background solar wind and (bottom) about 3.5 h after subjecting that field to a coronal mass ejection similar to that thought to have hit the Earth after the Carrington-Hodgson solar flare of 1859. The current sheet, extending from the cusp in the field in the downwind direction, clearest in the bottom panel, is the interface between oppositely-directed magnetic fields. The color bar in the center shows how the particle density is encoded in color. Source: first and last frames of a NASA animation, based on this study by (Ngwira et al., 2014). See Fig. H-I:6.3 for an comparison of the heliospheric and magnetospheric field configurations.

8 unknowns. Therefore we need to specify one more quantity, namely the strength of the shock ρ_2/ρ_1 .”

Examples of tangential discontinuities are the heliopause and planetary magnetospheres when there is little reconnection (see Sect. 5.5.5). A rotational discontinuity occurs for example when reconnection at the magnetopause is relatively efficient (Sect. 5.5.6). Shocks occur upstream of where the solar wind meets planetary system objects (Sects. 5.5.2–5.5.6) as well as where it encounters the interstellar medium (Sect. 5.5.8), and also where fast wind streams plow into slower ones ahead as well as at the leading edge of relatively fast explosions called coronal mass ejections (Sect. 5.5.1).

5.4 The magnetized solar wind and the Parker spiral

Before we can discuss the interplay of the solar wind with objects throughout the heliosphere, we need to introduce two properties of the solar wind itself: the geometry of the magnetic field that it carries, and the consequence of wind gusts running at different velocities. First, the magnetic field:

[H-1:9.2] "Let us briefly consider [a steady-state] outflow of ionized, magnetized gas from a *rotating* star with [a magnetic field that scales with distance from the star like a monopole, which is to say with a radially-flowing wind that stretches the field out from its effective base that lies, say, a few stellar radii above its actual surface ^[xi]]. The salient aspects of such a flow are found even when only considering the equatorial plane and restricting attention to solutions where all variables are functions of r only. We make use of spherical coordinates r, ϕ, θ .

We find that the equation of mass conservation (for plasma mass density ρ moving at velocity \mathbf{v} and carrying a field \mathbf{B}) then can be written

$$\frac{1}{r^2} \frac{\partial}{\partial r} (\rho v_r r^2) = 0, \quad (5.13)$$

while the ϕ component of the momentum equation is given by

$$\rho \left(v_r \frac{\partial v_\phi}{\partial r} + v_\phi \frac{v_r}{r} \right) = \frac{1}{4\pi} \left(B_r \frac{\partial B_\phi}{\partial r} + B_\phi \frac{B_r}{r} \right) \quad (5.14)$$

or

$$\rho v_r \frac{1}{r} \frac{d}{dr} (r v_\phi) = \frac{1}{4\pi} B_r \frac{1}{r} \frac{d}{dr} (r B_\phi) \quad (5.15)$$

Mass conservation implies that $\rho v_r r^2$ is constant, while the divergence-free magnetic field requires that $B_r r^2$ is constant. Multiplying Eq. (5.15) with r^3 we see that

$$r v_\phi - \frac{B_r r^2}{\rho v_r r^2} \frac{1}{4\pi} r B_\phi = \text{constant} = L. \quad (5.16)$$

^{xi} This base is what is meant by the term 'source surface' introduced in Activity 62; within that surface, the field is approximated as corotating rigidly with the star.

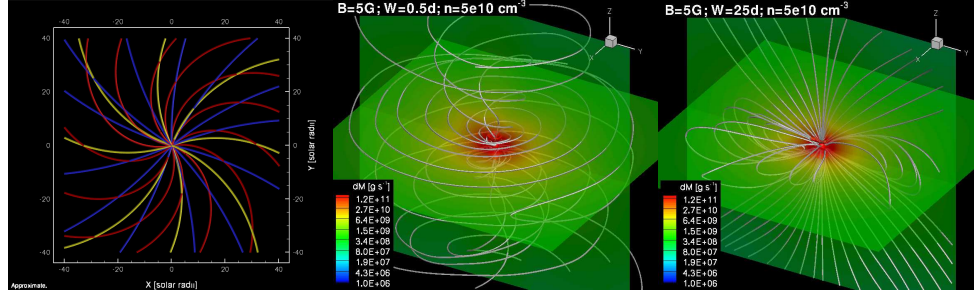


Fig. 5.5. Left: Conceptual display of different stellar-wind magnetic field spirals for a Sun with a 4.6 day rotation period (red), a 10 day period (yellow), and a 26 day period (blue), as a function of distance in solar radii. Center/right: Results from numerical simulations for the stellar coronae of solar analogs with rotation period of 0.5 day (middle) and 25 days (right), [but assuming the same base field strengths and densities. Fig. H-IV:4.1; sources: left from Cohen et al. (2012) and center and right from Cohen and Drake (2014).]

Under the assumptions above, the [ideal] induction equation is

$$\frac{1}{r} \frac{d}{dr} (r[v_r B_\phi - v_\phi B_r]) = 0. \quad (5.17)$$

For a star rotating with angular velocity Ω , radius R_s and with a [simplified monopolar] field so that $B_{\phi s} \approx 0$ [(with the index 's' meaning at the solar wind base or source surface)] we find that the induction equation implies

$$r(v_r B_\phi - v_\phi B_r) = \text{constant} \approx -R_s(R_s \Omega) B_{rs} = -\Omega r^2 B_r. \quad (5.18)$$

We can now solve Eq. (5.16) and Eq. (5.18) for v_ϕ and B_ϕ and using $M_A^2 = (v_r^2/v_A^2)$ with $v_A^2 = (B_r^2/4\pi\rho)$ we find

$$v_\phi = \Omega r \frac{M_A^2(L/r^2\Omega) - 1}{M_A^2 - 1}; \quad B_\phi = -\frac{B_r \Omega r}{v_r} \left(\frac{1 - (L/r^2\Omega)}{M_A^2 - 1} \right) M_A^2. \quad (5.19)$$

Both expressions show that we must have $1 - (L/r^2\Omega) = 0$ when $M_A^2 - 1 = 0$ [(formally, going to zero in such a way that their ratio remains finite)]. We define $r \equiv r_A$ where $M_A^2 = 1$. Thus, we must have $L = r_A^2 \Omega$. Notice [that v_r tends to a constant for large r , and thus $M_A^2 \propto \Omega r^2$ and as a result

$$\text{for large } r : v_\phi \approx \frac{\Omega r_A^2}{r} \rightarrow 0; \quad B_\phi \approx -\frac{B_r \Omega r}{v_r}, \quad (5.20)$$

while

$$\text{close to the star} : v_\phi \approx \Omega r; \quad B_\phi \approx -\frac{B_r \Omega r}{v_A}. \quad (5.21)$$

In] other words, the magnetic field and stellar wind rotate like a solid body out

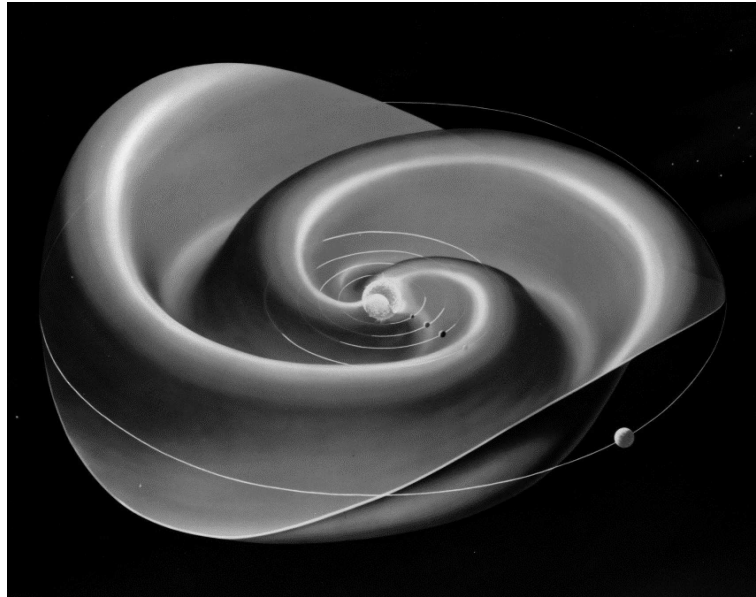


Fig. 5.6. The heliospheric current sheet forms a warped, undulating structure extending from the top ridge of the helmet streamer belt [...] that sweeps by the planets as the Sun rotates once per 27 days (synodic period). The magnetic field changes direction across the current sheet. [image source; see also Fig. H-I:9.3]

to the critical point r_A where the radial flow speed is equal to the 'radial' Alfvén speed. Beyond this point the field is pulled along the wind into a spiral, the *Parker spiral*, as the flow becomes nearly radial far from the star". Figure 5.5 shows this spiral and also two MHD simulations of stellar winds (discussed in Sect. 10.3.3). Note that whereas B_r decreases as $\sim 1/r^2$, B_ϕ decreases as $\sim 1/r$. {A:[60]} {A:[61]}

A:60

A:61

5.5 Flow-based interactions in heliophysics

5.5.1 Solar-wind stream interactions

5.5.1.1 A 1D model of high-speed stream evolution

The solar wind stretches the high-coronal field nearly radially into the heliosphere, there to deform subject to solar rotation into the Parker spiral. The

⁶⁰ Activity: Compare the radial dependence of the magnetic fields in this solar wind model with the values listed in Table 5.2. Also: use these dependences to demonstrate that the plasma β tends to a constant value far from the Sun.

⁶¹ Activity: At what distance from the Sun does the above solar-wind model have $|B_r| = |B_\phi|$ for typical values of the slow and fast solar wind? What are typical values for B_ϕ/B_r at 1 AU, 5 AU (cf. Fig. 5.9), and at the ice giants?

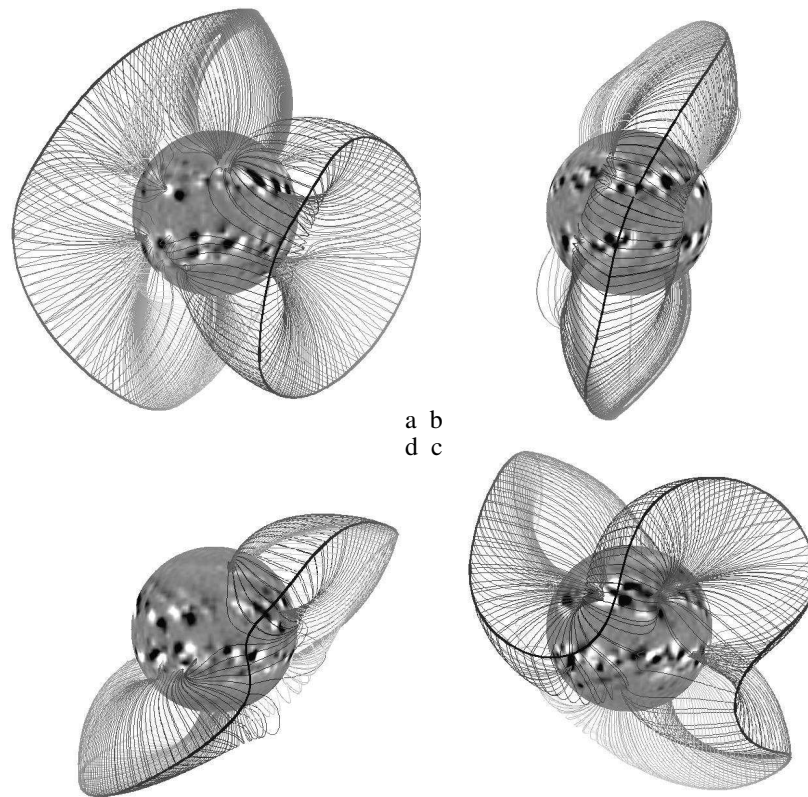


Fig. 5.7. The Sun's surface magnetic field is comprised of a multitude of dipolar regions of widely different fluxes, whose numbers wax and wane with the solar cycle. The large-scale coronal magnetic field, the foundation of the heliospheric field, expands from regions of partly open magnetic field that enclose the closed-field corona. This diagram shows the global topology of the Sun's field in a so-called potential-field source-surface approximation. In particular, it shows four realizations of the 'streamer belt' for a solar magnetic model. Shown are four phases of the simulated magnetic cycle: clockwise from the top left, $t = 3.1, 3.6, 4.5, 6.0$ years into a sunspot cycle of 11. years. Each panel shows a magnetogram of the solar surface, the neutral line(s) at the source surface, and the highest closed field lines that reach up to the neutral line(s); the lines are colored so that the darkest colors are nearest to the 'observer.' The panels show, clockwise, an example of a near-quadrupolar situation; a strongly tilted dipolar case; a strongly warped current sheet; and another nearly dipolar case with less tilt relative to the solar equator. [Fig. H-I:8.1]

magnetic field in the solar wind has its roots in the two magnetic polarities on the solar surface. Although the solar surface field has myriad adjacent regions of either polarity, further from the Sun the low orders dominate, so that the heliospheric field often resembles that of a distorted bipolar pattern stretched out far beyond the planets. In this field the opposite polarities straddle a

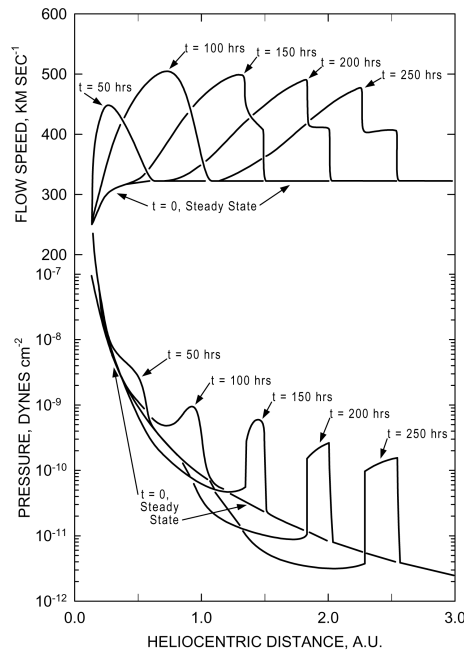


Fig. 5.8. Snapshots of solar wind flow speed and pressure as functions of heliocentric distance at different times during the outward evolution of a high-speed stream as calculated using a simple 1D gas-dynamic code. After obtaining a steady-state solar wind expansion that produced a flow speed of 325 km/s far from the Sun, a high-speed stream was introduced into the calculation by linearly increasing and then decreasing the temperature (and thus also the pressure) by a factor of four at the inner boundary from 0.14 AU over an interval of 100 hrs. [Fig. H-III:8.4]

transition in field direction in a warped skirt around the Sun, known as the heliospheric current sheet (sketched in its fundamental properties in Fig. 5.6; see Fig. 5.7 for a representation of an approximating (potential) field model at the foundation of that current sheet based on solar observations). {A:^[62]}

⁶² Activity: The solar wind stretches the high-coronal magnetic field into the heliosphere into a roughly radial field below the Alfvén radius. This enables an analogy with electrostatics: the field of electric charges placed above a flat perfect conductor can be computed by placing mirror charges opposite to the conducting surface, which then naturally has the electric field perfectly normal to the conducting surface. Analogously, in a magneto-static consideration above the spherical Sun, the magnetic field can be approximated by placing mirror 'charges' on a sphere at distance d_{SS}^2 which then has the field perfectly radial at d_{SS} . This is called the 'source surface model' with empirically $d_{SS} \approx 2.5R_{\odot}$ (where that 'source surface' is taken as the foundation of the heliospheric field; the virtual surface with mirror charges used to compute the potential field below d_{SS} is then at d_{SS}^2). This model (introduced by Schatten *et al.* (1969)) works remarkably well below d_{SS} on large scales. The heliospheric field is approximated by a radial continuation from that source surface, then subject to the Parker spiral. For illustration, simplify the source-surface model by a 2-d sketch involving a line of charges and another of mirror charges. Sketch the equivalent of the foundation of the heliospheric current sheet and examples of 'closed' field lines (the equivalent of coronal loops closing back onto the solar surface) and 'open' field lines (the equivalent of field stretched out into the heliosphere), at the base of which we find dark 'coronal holes' in X-ray images of the Sun.

Before going into this geometry, let us look into the simplified case of a 1-D ^{A:62} radial outflow.

Different magnetic regions on the Sun lead to different speeds in the wind. [Outside of eruptive phases,] this typically manifests itself in so-called fast streams and slow streams (see Table 2.4 for their properties). Because the Sun rotates, the radially flowing fast and slow winds cannot avoid but to run into each other. [H-III:8.5] “Because radially aligned parcels of plasma within a stream originate from different locations on the Sun, they are threaded by different magnetic field lines and thus cannot interpenetrate one another [without reconnection, and such reconnection proceeds relatively slowly in the solar wind compared to the characteristic time the wind takes to traverse much of the heliosphere]. Figure 5.8, which shows the result of a simple 1D gas-dynamic simulation, illustrates the basic reasons why high-speed streams evolve with increasing heliocentric distance. The rising portion of the high-speed stream steepens kinematically with increasing heliocentric distance because gas (plasma) at the peak of the stream is traveling faster than the slower plasma ahead. As the speed profile steepens, material within the stream is rearranged; parcels of plasma on the rising-speed portion of the stream are compressed, causing an increase in pressure there, while parcels of plasma on the falling-speed portion of the stream are increasingly separated, producing a rarefaction.

It is common to refer to the compression on the leading edge of a high-speed stream as an interaction region. Being a region of high pressure, the interaction region expands into the plasma both ahead and behind at the fast mode speed (actually at the sound speed in the calculation shown in Figure 5.8). The leading edge of the interaction region is called a forward wave because it propagates in the direction of the solar wind flow; the trailing edge is called a reverse wave because it propagates sunward in the solar wind rest frame but is carried away from the Sun by the highly supersonic flow of the wind. Pressure gradients associated with these waves produce an acceleration of the slow wind ahead and a deceleration of the high-speed wind within the stream. The net result of the interaction is to limit the steepening of the stream and to transfer momentum and energy from the fast wind to the slow wind. [...]

As long as the amplitude of a high-speed stream is sufficiently small, it gradually dampens with increasing heliocentric distance in the manner just described. However, when the difference in speed between the slow wind ahead and the peak of the stream is more than about twice the fast mode (sound) speed the stream initially steepens faster than the forward and reverse pressure waves can expand into the surrounding plasma; thus in such cases the interaction region at first narrows with increasing heliocentric distance. The nonlinear

rise in pressure associated with this squeezing eventually causes the forward and reverse waves bounding the interaction region to steepen into shocks. Because shocks (Sect. 5.3) propagate faster than the fast mode (sound) speed, the interaction region can expand once shock formation occurs. Observations reveal that relatively few stream interaction regions are bounded by shocks at 1 AU, but that most are near the equatorial plane at heliocentric distances beyond about 3 AU because the fast mode (sound) speed generally decreases with increasing distance from the Sun. At heliocentric distances beyond about 5 – 10 AU a large fraction of the mass and magnetic field flux in the solar wind at low heliographic latitudes is found within expanding compression regions bounded by shock waves on the rising portions of strongly damped high-speed streams. The basic structure of the solar wind near the solar equatorial plane in the distant heliosphere thus differs considerably from that observed near Earth. Stream amplitudes are severely reduced, and short-wavelength structure is damped out. The dominant structures at low latitudes (*i.e.*, within the band of variable wind ^[xii]) in the outer heliosphere are expanding compression regions that interact and merge with one another to form what are commonly called global merged interaction regions, GMIRs.”

5.5.1.2 Stream evolution in two and three dimensions

[H-III:8.5] “Should the coronal expansion be time-independent but inhomogeneous in heliocentric latitude and longitude, stream evolution proceeds similarly at all longitudes, but the state of a stream’s evolution varies with longitude. Because of solar rotation, the interaction region on the leading edge of a high-speed stream is wound into a spiral that at any particular heliocentric distance is inclined to the radial direction at an angle intermediate to that of the magnetic field threading the slow and fast wind flows respectively, as illustrated in Figure 5.9. The entire pattern of interaction co-rotates with the Sun and the compression region is known as a corotating interaction region, CIR. It is important to note, however, that it is only the pattern that co-rotates with the Sun because each parcel of solar wind plasma moves radially outward in this simple picture, except within the interaction region itself where both radial and transverse deflections of the flow occur. Because a CIR is inclined relative to the radial direction the pressure gradients associated with the interaction region have both radial and azimuthal components. With increasing heliocentric distance the forward wave propagates both anti-sunward and westward (in the direction of planetary motion about the Sun), whereas the reverse wave propagates both sunward (in the rest frame of the average solar wind) and

^{xii} The solar wind originating from high latitudes is typically fast as long as there are polar cap fields, *i.e.*, in phases around solar minimum. The solar wind from mid-to-low latitudes is a mixture of fast and slow streams, particularly around solar maximum.

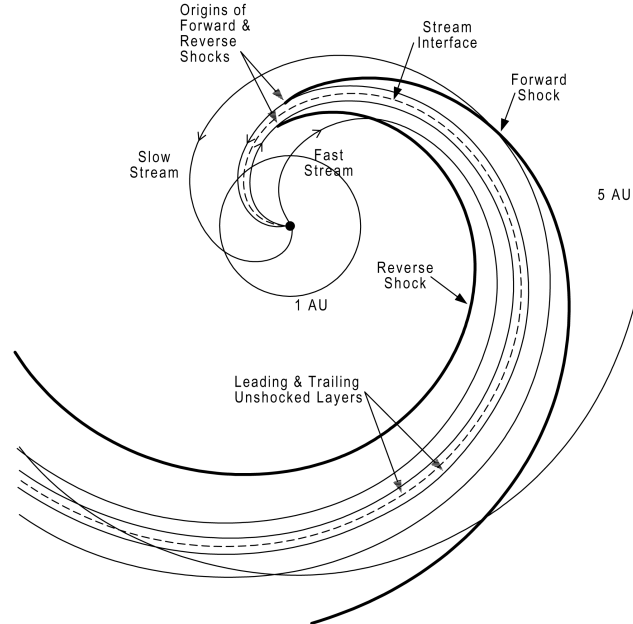


Fig. 5.9. Idealized schematic illustrating the basic structure of a corotating interaction region in the solar equatorial plane. The dashed line threading the middle of the corotating interacting region (CIR) denotes the stream interface and the solid heavy lines indicate the forward and reverse shocks. Plasma immediately surrounding the stream interface is compressed, but not shocked. [Note that in this stationary model moving a radial cut through the Sun in the clock-wise direction is equivalent to going forward in time at any given angle. Fig. H-III:8.5; source: Crooker et al. (1999).]

eastward. As a result, the slow wind is accelerated outward and deflected westward within the interaction region and the fast wind is decelerated and deflected eastward there, thus accounting for the characteristic westward and then eastward flow deflections commonly associated with interaction regions on the leading edges of high-speed streams (see Figure H-III:8.3 and related discussion). One consequence of the transverse deflections is that they partially relieve the pressure build-up induced by stream steepening by allowing the plasma to slip aside. Thus solar wind streams steepen less rapidly than is predicted by the simple 1D simulation shown in Figure 5.8.”

[H-III:8.5] “There is, of course, a three-dimensional aspect to stream evolution that becomes most apparent at heliocentric distances beyond about 3 – 4 AU and at latitudes away from the solar equatorial plane. [O]bservations have revealed (1) that the reverse shocks on the trailing edges of CIRs are observed both within the low-latitude band of solar wind variability and at latitudes $10^\circ - 20^\circ$ above that band, whereas the forward shocks on the leading edges of

corotating interaction regions are generally confined to the low-latitude band itself; and (2) that in addition to the flow deflections already discussed, the slow wind is usually deflected in both solar hemispheres toward the opposite hemisphere at the forward shocks, whereas the fast wind is usually deflected poleward at the reverse shocks.” For more details, see Sect. H-III:8.5.

5.5.2 *A non-conducting body without atmosphere*

Next, we look at one of the simplest setups for a flow encountering a body: a non-conducting sphere moving relative to a low-density magnetized plasma. The non-conducting body has no intrinsic magnetic field, and no currents can be induced in it by the magnetized plasma through which it moves. For a low-density plasma, this means that no signal is sent upstream from that body that could modify the flow heading towards the body: the gas pressure is insignificant and the magnetic field is not affected by the non-conducting body, moving through it without generating reflected waves that might move upstream from the body. Consequently, the upstream plasma that is on a collision course with the body will, in essence, simply crash onto the body, while the plasma to the sides of that body continues to flow without noticing the object at all [(Fig. 5.10a)]. This is true regardless of whether the body is moving sub-Alfvénically or super-Alfvénically relative to the incoming plasma. Examples of rather non-conducting bodies in the Solar System are Earth’s Moon, Jupiter’s moon Callisto and Saturn’s moon Rhea (which are subjected to the sub-Alfvénical flow of the giant planets’ magnetospheric plasma throughout their orbits), and also many asteroids, particularly the S-type, or silicate-rich ‘rocky’ ones (note that MHD does not apply for asteroids small compared to gyro-radii of solar wind ions; among other things, this means there is no upstream shock, as that is a collective phenomenon).

In such situations, there will be a wake behind the body that is void of plasma immediately downstream of the body. This void has two primary effects. One is that the plasma pressure has dropped away, and consequently plasma will propagate into the void to refill it (at about the slow-mode speed), taking matter from an outward propagating domain behind a rarefaction front that moves out at essentially the fast-mode wave speed (somewhat anisotropically); this leads to a wake in density behind the body, forming a somewhat asymmetric conical V-shape, albeit with the wings ending on a terminator-like ring defined by where the incoming plasma is just tangent to the body’s surface. The other effect is that because the contribution of the plasma to the total pressure falls away immediately behind the object, the field is somewhat strengthened (mainly by a motion perpendicular to the plane spanned by field and flow

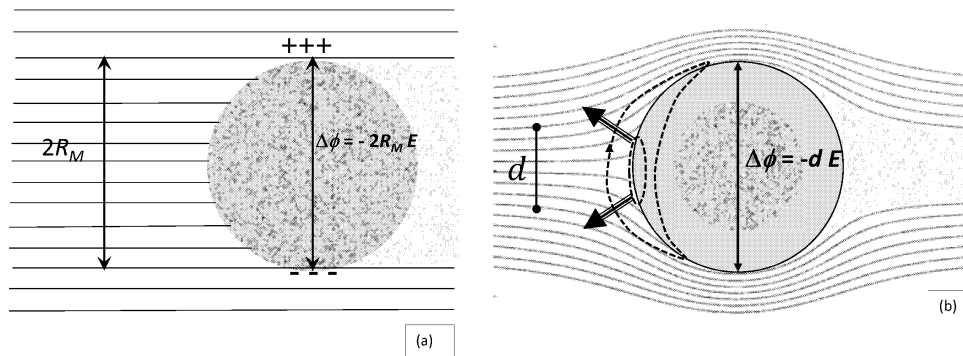


Fig. 5.10. Schematics of plasma flow (shown by lines of flow) at velocity \mathbf{v} from the left onto (a) a non-conducting body and (b) a conducting body [with radius R_M]. In the plasma, \mathbf{B} is into the paper, \mathbf{E} is $-\mathbf{v}/c \times \mathbf{B}$ in both cases. Diagram (a) shows that a non-conducting body builds up surface charge that imposes a potential drop $\Delta\phi = -2R_M E$ across the diameter, producing an electric field that opposes the solar wind electric field. Diagram (b) shows the response of a conducting body that does not build up surface charge. Conducting paths allow current (shown schematically as a dashed line) to flow through the body and close in the incident flow. Heavy banded arrows identify the orientation of the resultant $\mathbf{j} \times \mathbf{B}$ force that diverts part of the incident flow. Because much of the incident flow has been diverted, the potential drop across the body is only $\Delta\phi = d E$, where $d < R_M$ is the distance in the incident flow between the flow lines that just graze the body. The electric field that penetrates the body is a fraction of the upstream field determined by the fraction of the upstream flow that impacts the surface. In the wake region, gray in both diagrams, the plasma pressure is reduced and the magnetic pressure is increased relative to the upstream values. [Fig. H-IV:10.2]

vectors) to regain total pressure balance, to adjust again further downstream as plasma refills the void. In the case of Earth's Moon in the solar wind, the void persists up to about a dozen lunar radii downstream. {A:[63]} A:63

Note that the Moon is not a perfect insulator, and in fact has something akin to a weak ionosphere because the incoming solar wind ionizes some of the surface dust, and the ongoing process of ionization of such 'pick-up particles' is associated with a current that can send a magnetic signal upstream; see Chs. H-IV:10 and H-IV:11.

⁶³ Activity: For the solar wind flowing onto a non-conducting sphere, use estimates of wave speeds to sketch the density wake, the slow-mode refilling, and the fast-mode rarefaction front in a plane defined by the flow vector and the field vector, and in a plane defined by the flow vector and perpendicular to the field. You may compare the result with measurements for the case of the Moon (in Fig. H-IV:10.7).

5.5.3 Flow around a conducting body

When a conducting body moves through a magnetized plasma, the $-\mathbf{v}/c \times \mathbf{B}$ electric field associated with the relative motion induces a potential drop across the body. Because that body is conducting, a current flows to attempt to neutralize the charge buildup that would occur in the absence of such a current. That current closes through the incoming plasma in such a way that the associated Lorentz forces act to bend the plasma around the object (Fig. 5.10b), which, equivalently, sets up an induced magnetic field that, at infinite conductivity, would keep the external field entirely outside the conducting body. The conducting medium can be a metallic core (which in the case of the moons of the giant planets is generally too small to detect with significance) or a mantle ocean of water with dissolved electrolytes, *e.g.*, salts (which is seen on multiple moons of the giant planets, such as the Galilean moon Europa at Jupiter which is discussed in H-IV:10.5.2), a magma layer (as is inferred for another of the Galilean moons, Io) or an ionosphere in the upper layers of the body's atmosphere (such as in the case of Venus discussed in Sect. 13.1.3).

Europa moves sub-Alfvénically within the jovian magnetosphere. Because its orbit is inclined by about 10° relative to Jupiter's magnetic dipole moment, Europa senses a changing magnetic field throughout its orbit, so that not only a current system is induced by its motion, but that current system (and thus its associated perturbation magnetic field) evolves through the orbit. These changes (slightly modified by pickup ions playing their part) have revealed where the current flows, and thereby the existence of a conducting liquid underneath the non-conducting ice mantle.

The induced current system, or equivalently the induced magnetic field, sends out information about the obstacle into the plasma ahead. These waves, led by the magnetosonic fast-mode type, modify the upstream flow so that it can begin to deflect well ahead of the body. Part of the incoming flow may impact upon the surface if the conduction is not infinite. The rest of the flow is diverted around the body, leading to a narrower wake behind the body than in the case of an insulating body.

If the flow is coming in super-magnetosonically, however, no significant 'warning signal' can move upstream so that much of the flow will impact the body or flow very close to it, as it would in the case of large iron-rich asteroids (as has been argued, for example, for asteroid Ida, despite it being characterized as an S-type). But for many asteroids the scale is too small for MHD to apply, so the analogy with larger bodies fails in multiple respects.

Venus and Mars do not have active dynamos, but they do have conducting ionospheres. [H-I:13.6] "The magnetic structure surrounding Mars and Venus

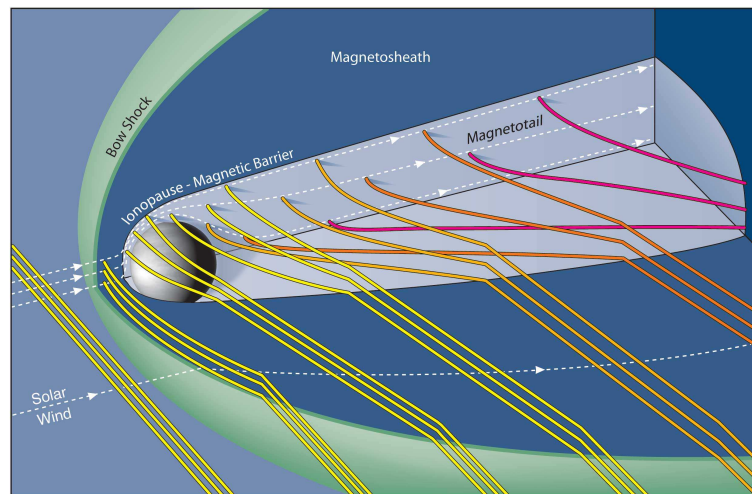


Fig. 5.11. Sketch of the draping of tubes of solar magnetic flux around a conducting ionosphere such as that of Venus. The flux tubes are slowed down and sink into the wake to form a tail. [color version of Fig. H-I:13.12]

is similar to that around magnetized objects because the interaction causes the magnetic field of the solar wind to drape around the planet. The draped field stretches out downstream (away from the Sun), forming a magnetotail (Fig. 5.11). The symmetry of the magnetic configuration within such a tail is governed by the orientation of the magnetic field in the incident solar wind, and that orientation changes with time. For example, if the interplanetary magnetic field (IMF) is oriented northward, the symmetry plane of the tail is in the east-west direction and the northern lobe field points away from the Sun while the southern lobe field points towards the Sun. A southward oriented IMF would reverse these polarities, and other orientations would produce rotations of the tail's plane of symmetry. {A:^[64]}

A:64

The solar wind brings in magnetic flux tubes that pile up at high altitudes at the dayside ionopause where, depending on the solar wind dynamic pressure $[(\rho v^2)]$, they may either remain for extended times, thus producing a magnetic barrier that diverts the incident solar wind, or may penetrate to low altitudes in localized bundles. Such localized bundles of magnetic flux are often highly twisted structures stretched out along the direction of the magnetic field. [These bundles] may be dragged deep into the atmosphere, possibly carrying away significant amounts of atmosphere.

⁶⁴ Activity: On the largest scales, there may be a long magnetotail to the entire heliosphere, that may even be oblate because of the tension force of the interstellar magnetic field. Although alternative views propose a much shorter tail, making the heliosphere more like a bubble, it is illustrative to see how such a moderate flattening by the interstellar magnetic field might work. Have a look at, *e.g.*, McComas *et al.* (2013), in particular their Figure 9.

While Mars' remarkably strong remanent magnetism ^[xiii] extends its influence > 1000 km from the surface, the overall interaction of the solar wind with Mars is more atmospheric than magnetospheric. Mars interacts with the solar wind principally through currents that link to the ionosphere, but there are portions of the surface over which local magnetic fields block the access of the solar wind to low altitudes. It has been suggested that 'mini-magnetospheres' extending up to 1000 km form above the regions of intense crustal magnetization in the southern hemisphere; these mini-magnetospheres protect portions of the atmosphere from direct interaction with the solar wind. [...]"

5.5.4 Plasma flow around a permanently magnetized body

Ganymede, orbiting sub-magnetosonically in Jupiter's magnetosphere, is the only moon with a substantial, large-scale internally maintained magnetic field. Of the planets, Earth and the giant planets all have magnetic fields sustained by dynamos, but in contrast to Ganymede and its surrounding plasma, they all move super-magnetosonically relative to the solar wind. In all of these cases, the bodies' magnetic fields are the primary 'obstacle' to the plasma flowing around it. All deflect the plasma stream around them. In the ideal-MHD approximation, the field-carrying plasma should flow around the magnetic obstacle, with the distance out to which the body's field can withstand the inflowing field dependent on the relative strength of the forces exerted (balancing magnetic fields and plasma inertial forces). In a realistic, non-ideal case, reconnection between the fields is important, which depends on the plasma parameters and on the relative directions of the two fields involved. In case the relative motion corresponds to a super-magnetosonic flow, a shock front develops; upstream of that, the inflowing plasma (generally the incoming solar wind) is, so to speak, unaware of the existence of the obstacle ahead, while the flow is deflected only after going through the shock, then moving around the obstacle at a reduced speed. This can still be faster than the Alfvén speed; see H-IV:10.4, which leads to a strong bending back of the wind flow around Earth into a bullet shape, in contrast to a V-shaped pattern for a sub-magnetosonic flow (*cf.* Fig. H-IV:10.4).

For those planets with a magnetic field of their own, *i.e.*, those with a dynamo, the solar wind leads to a shock-enveloped, asymmetrically-stretched magnetosphere. [H-I:10.2] "In the most general context, we consider a *central object*: a distinct well-defined body held together (in most cases) by its gravity. It is immersed in a tenuous *external medium*, assumed to be sufficiently ionized so it behaves like a plasma. The *magnetosphere* is then the region of

^{xiii} 'Remanent magnetism' is defined as the magnetic field that remains after the magnetizing field is removed.

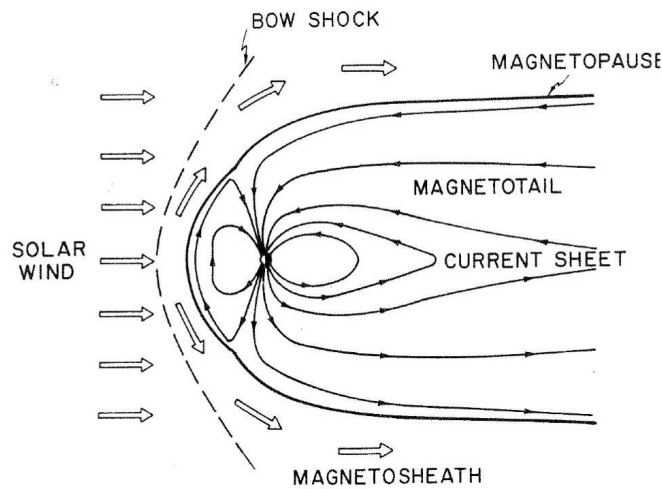


Fig. 5.12. Schematic view of a magnetically closed magnetosphere, cut in the noon-midnight meridian plane. Open arrows: solar wind bulk flow. Solid lines within magnetosphere: magnetic field lines (direction appropriate for Earth). [Fig. H-I:10.1]

space around the central object within which the object’s magnetic field has a dominant influence on the dynamics of the local medium. An alternative and in some ways more precise view is to regard the magnetosphere as the region enclosed by its bounding surface, the *magnetopause*, the latter being defined as the discontinuity of the magnetic field where its direction changes: inside it is controlled by the magnetic field of the central object, while outside it is determined primarily by the magnetic field of the distant external medium. This definition is particularly useful for the magnetospheres of planets in the solar wind: the continual variability of the interplanetary magnetic field direction in contrast to the relative constancy of the planetary magnetic dipole allows in most cases an easy observational identification of the magnetopause.”

5.5.5 A closed magnetosphere

[H-I:10.3] “The basic configuration of a prototypical planetary magnetosphere is sketched in Figure 5.12. Many of its characteristic structures can be understood on the basis of a simple model that takes into account only the two ingredients indispensable for the formation of a magnetosphere: the solar wind (mass density ρ_{sw} , bulk velocity \mathbf{v}_{sw}) and the planetary magnetic field (dipole moment $\mu_{\text{p}} = B_{\text{p}}R_{\text{p}}^3$, with B_{p} the surface magnetic field strength at the equator and R_{p} the radius of the planet). As a consequence of constraints imposed by the magnetohydrodynamic (MHD) approximation [...], the boundary surface

between the solar wind and the planetary magnetic field — the magnetopause — is nearly impermeable both to plasma and to magnetic field, resulting in a clear separation between the two distinct regions of space: the magnetosphere itself, within which the magnetic field lines from the planet are confined and from which the solar wind plasma is excluded, and the exterior region beyond the magnetopause, to which the plasma that comes from the solar wind is confined. This simple *closed magnetosphere* is only a first-order approximation (in reality the magnetopause is not completely impermeable but allows, under certain conditions, some penetration of plasma and of magnetic field to produce the *open magnetosphere* described in Sect. 5.5.6); it does, however, describe fairly accurately the size and shape of the main structures.

The solar wind flow, initially directed away from the Sun, must be diverted around the magnetosphere, as indicated in Figure 5.12. Because the initial flow speed is supersonic and super-Alfvénic (faster than both the speed of sound and the Alfvén speed v_A), the solar wind is first slowed down, deflected, and heated at a detached *bow shock* standing upstream of the magnetopause (analogous to the sonic boom in supersonic aerodynamic flow past an obstacle). The region between the bow shock and the magnetopause, within which the plasma from the solar wind is flowing around the magnetosphere, gradually speeding up and cooling, is called the *magnetosheath* [...]. {A:[65]}

A:65

The location of the magnetopause is determined primarily by the requirement of pressure balance: the total pressure (plasma plus magnetic) must have the same value on both sides of the discontinuity. In the simple closed magnetosphere considered here, the plasma pressure inside the magnetopause and the magnetic pressure outside are both neglected. The exterior pressure then scales as the linear momentum flux density in the undisturbed solar wind, $\rho_{sw}v_{sw}^2$ (often called the *dynamic pressure* of the solar wind), and is maximum in the sub-solar region, where the plasma near the magnetopause is almost stagnant. The interior pressure scales as the magnetic pressure of the dipole field, $(1/8\pi)(\mu_p/r^3)^2$ with μ_p the magnetic dipole moment of the planet, and thus varies strongly with distance from the planet. Equating the two gives an estimate for the distance R_{mp} of the sub-solar magnetopause:

$$R_{mp} = \frac{(\xi\mu_p)^{1/3}}{(8\pi\rho_{sw}v_{sw}^2)^{1/6}} \quad (5.22)$$

⁶⁵ Activity: Use Eqs. (5.2) and (5.8-5.12) to show that in the case of a strong shock (in which the thermal energy of the solar wind upstream of the bow shock can be ignored) the temperature just downwind of the bow shock is given by $(3m_p/32k)v_{sw}^2$ for a wind speed of v_{sw} , and that the density contrast across the shock is a factor of 4 (show that is true anywhere along the shock). Use this to estimate the angle from the upwind direction out to which the flow remains supersonic just inside the shock front (remembering that the transverse component of the velocity is unaffected by the shock).

where ξ is a numerical factor to correct for the added field from magnetopause currents ($\xi \simeq 2$ to first approximation). {A:[66]} A:66

The distance given in Eq. (5.22) (with various choices of ξ) is often called the Chapman-Ferraro distance. [Here, we] consistently use the symbol R_{CF} for the distance *defined* by Eq. (5.22) with $\xi = 2$, *i.e.*, for the nominal distance of the sub-solar magnetopause predicted by pressure balance; [the symbol R_{mp} is reserved] for the *actual* distance of the sub-solar magnetopause in any particular context. Thus, $R_{mp} \simeq R_{CF}$ in the present case of a simple closed magnetosphere but not necessarily in the case of more general models. {A:[67]} {A:[68]} A:67

The pressure balance condition, combined with assumptions about the sources of the magnetic field within the magnetosphere, may be used to calculate not only the distance to the sub-solar point, but also the complete shape of the magnetopause surface (for discussion of such models at Earth, see Ch. H-I:11). Typically the magnetopause is roughly spherical on the dayside of the planet, facing into the solar wind flow (the effective center of the sphere being located behind the planet, very roughly at a distance $\sim 0.5 R_{mp}$), and is elongated in the anti-sunward direction. A:68

The magnetopause distance R_{mp} may be regarded as the characteristic scale for the size of a magnetosphere. Equal to R_{CF} in the case of negligible plasma pressure and no magnetic field sources other than the planetary dipole inside the magnetosphere, R_{mp} can be readily calculated from Eq. (5.22) given only a few basic parameters of the system. In the case that the plasma pressure or a non-dipolar field in the outer regions of the magnetosphere are not negligible, the qualitative effects on R_{mp} can still be estimated from pressure balance, as illustrated in Figure 5.13: (a) the actual distance R_{mp} is larger than the nominal distance R_{CF} (the value $\xi = 2$ instead of $\xi = 1$ is in fact a consequence of the non-dipolar field from the magnetopause currents), (b) a change of solar wind dynamic pressure produces a larger change of magnetopause distance — the magnetosphere is less 'stiff' if plasma pressure in the interior is significant.”

5.5.6 The open magnetosphere

[H-I:10.3.3] “At the location of the planets, the interplanetary magnetic field is weak in the sense that the energy density of the magnetic field is very small in comparison to the kinetic energy density of solar wind bulk flow, or

⁶⁶ Activity: What is the expression for the temperature of the gas at the stagnation point on the magnetopause assuming that the flow continues adiabatically after the shock (*i.e.*, that it conserves the sum of bulk kinetic and thermal energies)? What is the value for $v_{sw} = 800$ km/s.

⁶⁷ Activity: Use Eq. (5.22) to show the scaling of R_{CF} with orbital radius, planetary magnetic field, and planetary radius.

⁶⁸ Activity: With the fastest recorded solar-wind gusts at $v_{sw} \approx 2500$ km/s, what is the required plasma density to push the magnetopause to within geosynchronous orbit according to Eq. (5.22)?

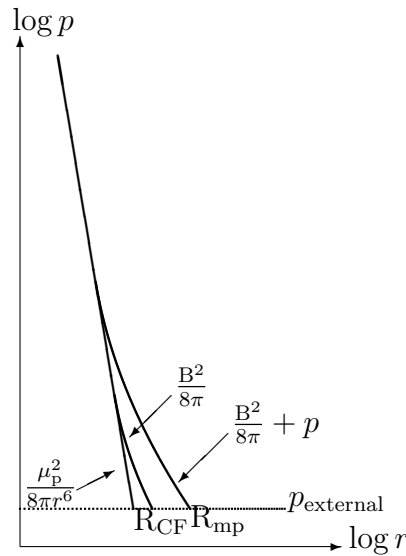


Fig. 5.13. Variation of total pressure (magnetic plus plasma) with distance from the planet and its relation to the radial distance of the sub-solar magnetopause. Compared are the relationship in Eq. (5.22) to a schematic representation of a more realistic plasma-filled, non-dipolar planetary magnetic field. [Fig. H-I:10.2]

equivalently $v_A^2 \ll v_{sw}^2$ [(see Sect. 3.5.2)]. The flow of solar wind plasma past the magnetospheric obstacle deforms the magnetic field lines within the magnetosheath and drapes them around the magnetopause, a process well modeled at Earth (*cf.* Sect. H-I:11.4). The magnetic field is amplified and may become dynamically no longer negligible as the magnetopause is approached, but the total pressure is in general not greatly modified, an increase of magnetic pressure often being offset by a decrease of plasma pressure. One might therefore anticipate that the effect of the interplanetary magnetic field on planetary magnetospheres should be minimal.

What is overlooked in the above discussion is the possibility that, through the process of *magnetic reconnection* [...], the magnetic field lines from the planet may become connected with those of the interplanetary magnetic field, to produce the magnetically *open magnetosphere*, sketched in Fig. 5.14 for the simplest case of the interplanetary magnetic field parallel to the planetary dipole moment. The magnetopause is now no longer impermeable to the magnetic field, and as a consequence it no longer need be impermeable to plasma, either. {A:[69]}

A:69

⁶⁹ Activity: Illustrative diagrams like Fig. 5.14 typically show the interplanetary magnetic field (IMF) as lying within the $x - z$ plane of such diagrams. In reality, the three IMF components $B_{x,y,z}$ are typically of comparable magnitude. Moreover, the orientation of the Earth's magnetic axis relative to the incoming wind changes in the course of the year. Consider how the diagrams should look when drawn in three dimensions for a few different combinations of $B_{x,y,z}$. Look up the 'Russell-McPherron

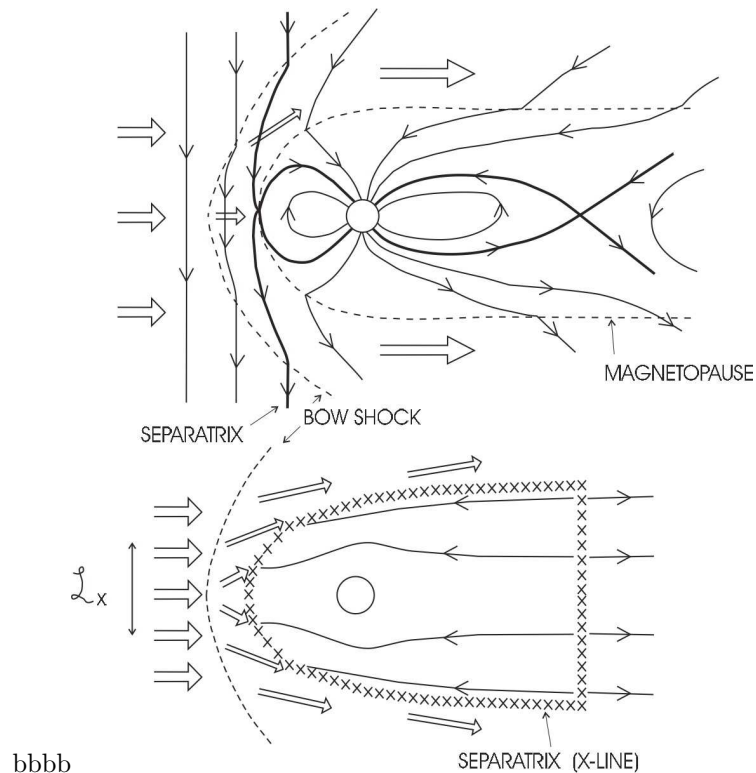


Fig. 5.14. Schematic representation of a magnetically open magnetosphere. Top: cut in noon-midnight meridian plane; thick lines are magnetic field lines within the 'separatrix surfaces' that separate open from closed or open from interplanetary field lines; other conventions same as in Fig. 5.12. Bottom: cut in equatorial plane; a line of \times symbols represents intersection with the two branches of the separatrix; solid lines are streamlines of magnetospheric plasma flow, and \mathcal{L}_x represents the projection of the dayside magnetic reconnection region along streamlines into the solar wind. [Fig. H-I:10.3]

The modifications of the magnetospheric system implied by the open character of the magnetosphere are in some ways minor, in other ways very far-reaching. The location and shape of the dayside magnetopause is for the most part not greatly modified (in agreement with the expectations above). The component B_{\perp} of the magnetic field normal to the magnetopause is in general small compared to the magnitude of the field, $|B_{\perp}| \ll |B|$ (so much so that it is often difficult to establish by direct observation that $B_{\perp} \neq 0$, and much of the evidence for an open magnetosphere has been indirect). On the other

effect' which attributes the semi-annual variations in geomagnetic activity largely to the relative orientation of the Earth's bipole axis: maximum geomagnetic activity around the equinoxes, minimum around solstices.

Table 5.2. *Properties of the solar wind near the planets [after Table H-I:13.2]. Plasma β values assume a solar-wind temperature of 1.5 MK.*

Planet	Distance d_p (AU) ^a	Solar wind density (cm^{-3})	B_{IMF} (μG) ^c	$\approx \beta$	$\approx v_A$ (km/s)
Mercury	0.39	53	410	2	120
Venus	0.72	14	140	4	80
Earth	1	7 ^b	80	6	70
Mars	1.52	3	50	6	60
Jupiter	5.2	0.2	10	10	50
Saturn	9.5	0.07	6	10	50
Uranus	19	0.02	3	10	50
Neptune	30	0.006	2	10	50

^a 1 AU = $1.5 \cdot 10^8$ km; ^b The density of the solar wind fluctuates by about a factor of 5 about typical values of $\rho_{\text{sw}} \sim (7 \text{ cm}^{-3})/d_p^2$; ^c mean values. [...]

Table 5.3. *[Intrinsic magnetic fields of Solar System bodies. After Table H-I:13.3, with planetary rotation periods P_p and planetary radii R_p].*

	Gany- mede	Mer- cury	Earth	Jupiter	Saturn	Uranus	Nep- tune
$B_{\text{dip,eq}}$ ^a	7.2 mG	3 mG	0.31 G	4.3 G	0.21 G	0.23 G	0.14 G
$B_{\text{max}}/B_{\text{min}}$ ^b	2	2	2.8	4.5	4.6	12	9
dipole tilt ^c	-4°	$\sim 10^\circ$	11.2°	-9.4°	-0.0°	-59°	-47°
dipole offset ^d	-	-	0.076	0.119	0.038	0.352	0.485
obliquity ^e	0°	0°	23.5°	3.1°	26.7°	97.9°	29.6°
$\delta\phi_{\text{sw}}$ ^f	90°	90°	$67-114^\circ$	$87-93^\circ$	$64-117^\circ$	$8-172^\circ$	$60-120^\circ$
P_p (h)	171	4223.	24	9.9	10.7	17.2	16.1
$R_p/R_{p,\oplus}$	0.41	0.38	1	11.2	9.4	4.0	3.9

^a Surface field at dipole equator. Values derived from modeling the magnetic field as an offset dipole; ^b ratio of maximum surface field to minimum, which equals to 2 for a centered dipole field (this ratio tends to increase with the planet's oblateness); ^c angle between the magnetic and rotation axes (positive values correspond to magnetic field directed north at the equator; the magnetic dip poles of the Earth's field are currently located at 86°N and 65°S latitudes and moving about 10° per century); ^d values (in planetary radii, R_p); ^e the inclination of a planet's spin equator to the ecliptic plane; ^f range of angle between the radial direction from the Sun and the planet's rotation axis over an orbital period (in Ganymede's case, the angle is between the corotational flow and the moon's spin axis).

hand, the total amount of open magnetic flux Φ_M of one polarity can (at least at Earth) become comparable to the maximum amount that could reasonably be expected to be open (estimated as $\sim \mu_p/R_{\text{mp}}$, the dipole flux beyond the

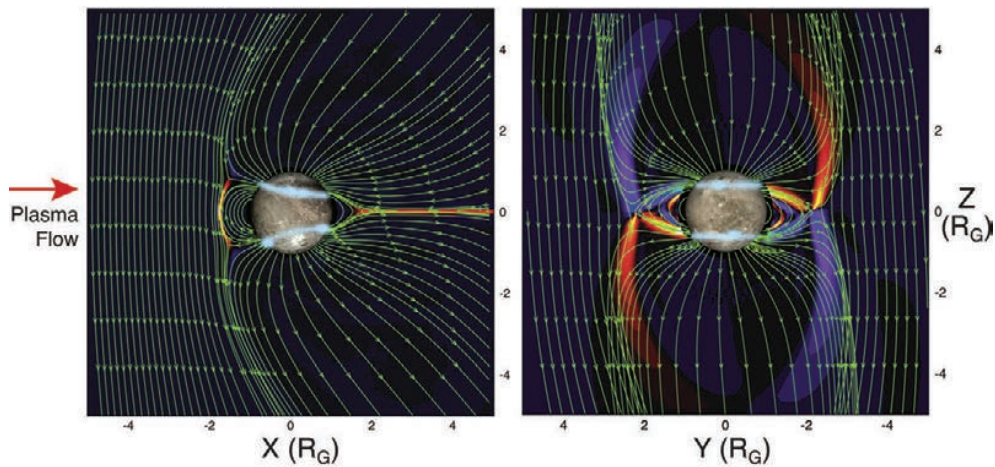


Fig. 5.15. Numerical model of the magnetosphere of Ganymede with the satellite and the location of auroral emissions superimposed. [The shaded areas show the current density perpendicular to the plane; yellow-red out of the plane, purple-blue into the plane.] Left: view looking at the anti-Jupiter side of Ganymede. Right: View looking in the direction of the plasma flow at the upstream side (orbital trailing side) of Ganymede with Jupiter to the left. The shaded areas show the regions of currents parallel to the magnetic field; yellow-red anti-parallel, purple-blue parallel]. [Fig. H-I:13.11; source: Jia *et al.* (2008).]

distance of the sub-solar magnetopause); despite $|B_{\perp}| \ll |B|$, this is possible if the effective length of the magnetotail is much larger than R_{mp} . [... That length is in large part determined by the efficiency of the reconnection process. This] depends greatly on the relative orientation of magnetic fields on the two sides of the magnetopause, one result of which is that the open character of the magnetosphere is most pronounced when the interplanetary magnetic field is parallel to the planetary dipole moment (*i.e.*, anti-parallel to the dipole magnetic field in the equatorial plane), $\mathbf{B}_{\text{sw}} \cdot \boldsymbol{\mu}_{\text{p}} > 0$. Because the direction of the interplanetary magnetic field is highly variable on all time scales, this can lead to pronounced time-varying changes of magnetospheric configuration as well as energy input and dissipation [...].”

Although diagrams of the terrestrial environment such as Fig. 5.14 generally include the bow shock, the processes discussed here are generic and apply equally when a magnetized plasma flows sub-Alfvénically around a magnetized body. One example is a numerical simulation of Ganymede orbiting within Jupiter’s magnetosphere as shown in Fig. 5.15. {A:^[70]}

A:70

⁷⁰ Activity: Use Fig. 5.15 to estimate the Alfvén velocity in the jovian magnetosphere near Ganymede. First, estimate the flow speed of the incoming plasma relative to the moon realizing that the plasma is sub-corotating by about 80% of the speed of corotation with Jupiter at Ganymede’s orbit. Then use the geometry of the field shown in the figure to estimate the Alfvén velocity. The coordinate

5.5.7 Solar wind-magnetosphere-ionosphere interaction

5.5.7.1 Fundamental principles

[H-I:10.4.1] “A well-known consequence of the MHD approximation is a constraining relation between the plasma bulk flow and the magnetic field: plasma elements that are initially on a common field line remain on a common field line as they are carried by the bulk flow. Because magnetic field lines in the magnetosphere of a planet connect to the ionosphere of the planet, any discussion of plasma flow in the magnetosphere immediately involves questions of magnetosphere-ionosphere interaction. The field lines extend in fact into the interior of the planet, which in many cases is highly conducting electrically; hence it might seem that the magnetospheric flow should be constrained by the planet itself. This does not happen, however, because [...] most planets possess an electrically neutral (and effectively non-conducting) atmosphere, sandwiched between the ionosphere and the planetary interior. Although very thin in comparison to the radius of the planet, this layer suffices to break the MHD constraints and thus allows the plasma in the ionosphere and the magnetosphere to move without being necessarily attached to the planet; without such an insulating layer, much of the magnetospheric dynamics as we know it would not be possible.

While the plasma in the ionosphere can thus move relative to the planet, it remains constrained to move more or less together with the plasma in the magnetosphere. The conventional formulation, however, describes the plasma flow rather differently in the two regions. The magnetosphere is treated, to first approximation at least, as an MHD medium, with the electric field \mathbf{E} related to the plasma bulk flow \mathbf{v} by the MHD approximation and with the electric current \mathbf{j} related to plasma pressure by stress balance. The ionosphere is treated, on the other hand, as a moving conductor (the conductivity results primarily from collisions between the ions and the neutral particles, planetary ionospheres being for the most part weakly ionized), with \mathbf{j} related by a conductivity tensor to $\mathbf{E} + \mathbf{v}_\perp \times \mathbf{B}/c$, where \mathbf{v}_\perp is the bulk velocity of the neutral medium. [...] A few comments on this coupling:]

(1) As long as $v_A^2/c^2 \ll 1$ (*i.e.*, the inertia of the plasma is dominated by the rest mass of the plasma particles, not by the relativistic energy-equivalent mass of the magnetic field), \mathbf{v} produces \mathbf{E} but \mathbf{E} does not produce \mathbf{v} . The primary quantity physically is thus the plasma bulk flow, established by appropriate stresses. The electric field is the result of the flow, not the cause; its widespread use in calculations is primarily for mathematical convenience [...].

(2) The electric current in the ionosphere is not an Ohmic current in the

system for the simulation has the y -axis pointing towards Jupiter and the z -axis aligned with the jovian spin axis, and the units are expressed in Ganymede radii.

physical sense, and its conventional expression by the 'ionospheric Ohm's law' [as discussed around Eq. (2.22)] has only a mathematical significance [...]. Physically, the current is determined by the requirement that the Lorentz force balance the collisional drag between the plasma and the neutral atmosphere when their bulk flow velocities differ. [...] The current in the ionosphere is thus governed by stress balance in the same way as the current in the magnetosphere [... while neglecting the time derivative term in the momentum equation.]

(3) Underlying the neglect of the time derivative (acceleration) terms in the momentum equations is the implicit assumption that any imbalance between the mechanical and the magnetic stresses (which, fundamentally, is what determines the acceleration of the plasma) produces a bulk flow that acts to reduce the imbalance, which then becomes negligible over a characteristic time scale (easily shown to be of the order of the Alfvén wave travel time across a typical spatial scale \mathcal{L} , *e.g.*, along a field line). The theory is thus applicable only to systems that are stable and evolve on time scales much longer than \mathcal{L}/v_A .

(4) For [phenomena well above the scale above plasma oscillations, so slow compared to $1/\omega_p$ and large compared to $\lambda_e \equiv c/\omega_p$,] (where ω_p is the electron plasma frequency and λ_e the electron inertial length, also known as collisionless skin depth), \mathbf{j} adjusts itself to become equal to $(c/4\pi)\nabla \times \mathbf{B}$ and not the other way around; although \mathbf{B} is in principle determined from a given \mathbf{j} by Maxwell's equations (on a time scale of light travel time, $\sim \mathcal{L}/c$), in a large-scale plasma any $\mathbf{j} \neq (c/4\pi)\nabla \times \mathbf{B}$ is immediately (on a time scale $\sim 1/\omega_p$) changed by the action of the displacement-current electric field on the free electrons in the plasma. The current continuity condition $\nabla \cdot \mathbf{j} = 0$ is thus satisfied automatically; there is no physics in current closure — what is often discussed under that rubric is in reality the coupling of the Maxwell stresses along different portions of a field line.”

5.5.7.2 Corotation

[H-I:10.4.2] “Corotation with the planet is the simplest pattern of plasma flow in a planetary magnetosphere and one that plays a major role particularly in the magnetospheres of the giant planets. [...] If the planet possesses an insulating atmosphere, the rotation of the planet itself has no *direct* effect on plasma flow in the magnetosphere, as discussed in Sect. 5.5.7.1. What does affect plasma flow is the motion of the neutral upper atmosphere (thermosphere) at altitudes of the ionosphere (where the neutral and the ionized components coexist and interact). 'Corotation with the planet' is therefore not quite an accurate description. What really is meant is co-motion with the upper atmosphere, which in turn is then assumed to corotate with the planet, for reasons unrelated to the magnetic field: vertical transport of horizontal linear momentum from

the planet to the neutral atmosphere (*e.g.*, by collisional or eddy viscosity and similar processes), together with an assumed small relative amplitude of neutral winds.

Any difference between the bulk flow of the neutral medium and the ionized component of the plasma in the ionosphere results in a collisional drag that must be balanced by the Lorentz force; without it, the drag force would soon bring the plasma to flow with the (much more massive) neutral medium. The Lorentz force in the ionosphere is coupled to a corresponding Lorentz force in the magnetosphere, which in turn must be balanced locally by an appropriate mechanical stress. The net result is that departure from corotation requires a mechanical stress in the magnetosphere to balance the plasma-neutral drag in the ionosphere; conversely, plasma will corotate if the stress in question is negligibly small. (It is fairly obvious that the direction of the stress must be more or less azimuthal, opposed to the direction of rotation.) Quantitatively, the requirements for corotation of magnetospheric plasma may be expressed by four conditions:

(1) Planet-atmosphere coupling: This is simply the assumption, discussed above, that the upper atmosphere effectively corotates with the planet.

(2) Plasma-neutral coupling in the ionosphere: the collisional drag of the neutral medium on the plasma must be sufficiently strong to ensure $\mathbf{v} \simeq \mathbf{v}_n$. The quantitative condition is derived in principle [from the momentum equation of the ionospheric plasma (horizontal components only)]

$$\frac{\partial \rho \mathbf{v}}{\partial t} + \dots = \mathbf{j} \times \mathbf{B}/c - \nu_{in} \rho (\mathbf{v} - \mathbf{v}_n) \quad (5.23)$$

with the left-hand side set to zero (where ν_{in} is the ion-neutral collision frequency),] but with one complication: what is relevant for the interaction with the magnetosphere is not the local current density \mathbf{j} but the current per unit length integrated over the extent of the ionosphere in altitude z , *i.e.*, the height-integrated current $\mathbf{I} \equiv \int \mathbf{j} dz$. A direct integration of Eq. (5.23) over height, however, is not simple because \mathbf{v} varies strongly with z (even when \mathbf{v}_n is independent of z , as usually assumed). The horizontal electric field, on the other hand, is essentially constant over the entire (relatively thin) height range of the ionosphere, from continuity of tangential components implied by Faraday's law. It is thus convenient to first express \mathbf{j} by [the ionospheric Ohm's law as discussed around Eq. (2.22)] and then integrate over height to obtain

$$\mathbf{I}_\perp = (B/c) \left[\Sigma_P \hat{\mathbf{B}} \times (\mathbf{v}_0 - \mathbf{v}_n) - \Sigma_H (\mathbf{v}_0 - \mathbf{v}_n)_\perp \right] \quad (5.24)$$

where \mathbf{v}_0 is the plasma flow at the top side of the ionosphere, related to the

electric field by

$$\mathbf{E}^* = -(\mathbf{v}_0 - \mathbf{v}_n) \times \mathbf{B}/c \quad \text{or equivalently} \quad \mathbf{E} = -\mathbf{v}_0 \times \mathbf{B}/c. \quad (5.25)$$

Σ_P and Σ_H are the height-integrated Pedersen and the Hall conductances, Σ_P being the more important one for magnetosphere-ionosphere interactions (Hall currents close within the ionosphere, to first approximation).

Obviously, to ensure $\mathbf{v}_0 \simeq \mathbf{v}_n$, the ionospheric conductance Σ_P must be sufficiently large in relation to the height-integrated current \mathbf{I} , which scales as the current per unit length in the magnetosphere and hence ultimately as the mechanical stresses in the magnetosphere. For a more precise criterion, one must consider a specific process. [...]

(3) MHD coupling from ionosphere to magnetosphere along magnetic field lines: conditions (1) and (2) ensure merely that the plasma corotates at the top side of the ionosphere, at the foot of a magnetic flux tube within the magnetosphere. For corotation to extend into the magnetosphere itself, the MHD constraining relation between the flow and the magnetic field must hold. [...]

(4) Stress balance to maintain centripetal acceleration in the magnetosphere: If conditions (1), (2), and (3) are satisfied, the plasma will be corotating at least as far as the components of \mathbf{v} perpendicular to \mathbf{B} are concerned, but the flow parallel to \mathbf{B} remains unconstrained. For the entire flow to be corotational, one further condition must be satisfied: there must exist a radial stress to balance the centripetal acceleration of the corotating plasma. In most cases, this stress is produced by the corotation itself, as the magnetic field lines are pulled out until their tension force becomes sufficiently strong to balance the centripetal acceleration.”

[H-I:10.4.4] “The four conditions for corotation are all, in essence, local conditions at a given magnetic flux tube. Deviations from corotational flow when one or another of these conditions is no longer satisfied need not, therefore, be global but can be confined to limited regions. Typically, plasma flow in any particular magnetosphere may follow corotation in the inner regions, out to a critical radial distance in the equatorial plane, and then deviate significantly from corotation at larger distances. The critical distance depends on which of the four conditions is violated and by which process.” Section H-I:10.4.4 provides more discussion. {A:[71]}

A:71

⁷¹ Activity: Consider differences and similarities between ‘corotation’ in a planetary magnetosphere and in the solar wind, including (a) the absence of a sufficiently neutral atmosphere in the Sun to decouple the motions between internal and heliospheric fields (associated with a concept called ‘line tying’, which we touch upon in Sec. 6.3.1.1), and (b) the very term ‘corotation’ which to a heliospheric physicist does *not* include the component of $\mathbf{v} \parallel \mathbf{B}$ but is limited to the pattern of the field, not the plasma itself.

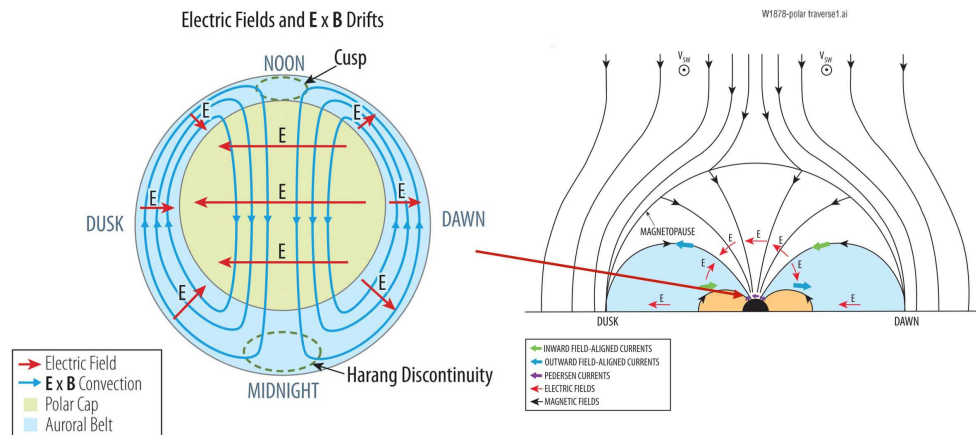


Fig. 5.16. Schematic diagram of magnetospheric convection over the Earth's north polar region (note the diagram extends in latitude only to the lower edge of the auroral belt). Left: streamlines of the plasma bulk flow and associated electric field (the Sun is towards the top). Right: Magnetic field and current systems for the northern hemisphere, for a southward interplanetary magnetic field, viewed from behind the Earth looking towards the Sun. See Section H-I:10.4.3 for a detailed description and Sect. H-I:11.6 for corresponding MHD model results for the electric potential. For a schematic representation of magnetospheric convection throughout the magnetosphere, see Fig. H-I:13.4. [Source: NASA ;compare: Fig. H-I:10.5]

5.5.7.3 Magnetospheric convection

[H-I:10.4.3] “Magnetospheric convection may be considered the other canonical pattern (besides corotation) of plasma flow in a planetary magnetosphere, one that plays an overwhelmingly important role in the magnetosphere of Earth. The basic concept is that the flow of solar wind plasma past the magnetosphere imparts some of its motion to plasma in the outermost regions of the magnetosphere, either directly by MHD coupling along open field lines or through an unspecified tangential drag near the magnetopause.

By continuity of mass and magnetic flux transport, the flow then extends into the region of closed field lines or the interior of the magnetosphere, setting up a large-scale circulation pattern (which has some superficial resemblance to, but no real physical commonality with, what is called convection in ordinary fluid dynamics). Figure 5.16 illustrates the pattern, projected on the top-side ionosphere of the planet: shown on the left-hand side are the streamlines of the plasma bulk flow \mathbf{v} , which are also the equipotentials of the electric field according to Eq. (5.25). The lines of the electric field and the associated Pedersen currents are shown on the right-hand side, along with a sketch of the implied Birkeland (*i.e.*, magnetic field-aligned) currents. The yellow region is the polar cap, identified with the region of open field lines in the

open magnetosphere; otherwise it represents the mapping (along field lines) of the boundary region where the solar wind motion is being imparted to magnetospheric plasma. The equatorial-plane counterpart of the flow outside the polar cap has been sketched in Fig. 5.14 (bottom).

A quantitative global measure of the strength of magnetospheric convection is the EMF (maximum line integral of the electric field) across the polar cap, \mathcal{E}_{PC} . Its physical meaning is that of the rate of magnetic flux transport (advection) through the polar cap. In an open magnetosphere, $c\mathcal{E}_{\text{PC}}$ equals the rate of reconnection of magnetic flux between the interplanetary and the planetary magnetic fields. Numerous empirical studies at Earth have shown that for a southward interplanetary magnetic field (*i.e.*, $(\mathbf{B}_{\text{sw}} \cdot \hat{\boldsymbol{\mu}}_{\text{p}}) > 0$), \mathcal{E}_{PC} can be related to solar wind parameters approximately as

$$c\mathcal{E}_{\text{PC}} \simeq v_{\text{sw}} (\mathbf{B}_{\text{sw}} \cdot \hat{\boldsymbol{\mu}}_{\text{p}}) \mathcal{L}_{\times} \quad (5.26)$$

where \mathcal{L}_{\times} is a length that typically is a fraction (~ 0.2 to ~ 0.5) of the magnetopause radius R_{mp} . When comparing different magnetospheres, one often supposes that the ratio $\mathcal{L}_{\times}/R_{\text{mp}}$ is a more or less universal constant. Physically, \mathcal{L}_{\times} may be looked at as the length of the reconnection X-line on the magnetopause at which the magnetic field lines from the solar wind and from the planet first become interconnected, projected along the streamlines of the magnetosheath flow back into the undisturbed solar wind, as illustrated in Fig. 5.14.” {A:[72]} {A:[73]} {A:[74]} {A:[75]} A:72

⁷² Activity: The magnetospheric magnetic field cycle starts for the field with (1) day-side reconnection to the field in the wind, is then (2) followed by being dragged towards night-side, from there (3) moving into the magnetotail, and after (4) reconnection in the current sheet the field (5) moves back towards the day-side to replenish (at least on average over longer periods) the flux lost from there in the reconnection process. That loop, called the Dungey cycle, can be visualized from Fig. 5.14 if the succession of drawn field lines is interpreted as a sequence of events for a single field line (and realizing that step (5) has to occur over lower magnetic latitudes to avoid the field that is at the same time involved in step (2)). But during the cycle, the planet rotates underneath, dragging the ionospheric plasma onto which the magnetospheric field connects with it. To see for which planets this process is important, estimate for each of the planets: (a) the model-based magnetopause distance R_{CF} , (b) the time it takes to move from step (1) to a phase somewhere around steps (3) and (4), and (c) how many turns the planet has made in the meantime. Assume the following: that the solar wind speed averages to roughly the same value at all the planets (say $v_{\text{sw}} = 400$ km/s), and that the flow of the plasma in the deep magnetosheath carries the field from front to back over, say, $3R_{\text{CF}}$ at $0.1v_{\text{sw}}$. Info in Table 5.3. If you also want to do Ganymede, realize that it has spin-orbit synchronization within the jovian magnetosphere. A:73

⁷³ Activity: Consider the possible equivalent of a Dungey cycle for the heliospheric field subject to reconnection with an interstellar magnetic field. What would happen in case there were no coronal heating? A:74

⁷⁴ Activity: **What if?:** Many so-called ‘hot Jupiters’ have been found among the exoplanet population: giant planets that orbit very close to their parent stars. What would the estimated magnetopause distance $R_{\text{CF,hJ}}$ be if Jupiter were orbiting the present-day Sun at 0.05 AU? For a younger Sun (see Ch. 12) the solar wind would have been stronger, pushing $R_{\text{CF,hJ}}$ to below the orbital radius of Ganymede; describe what that would mean for this ‘hot-Ganymede’ moon? A:75

⁷⁵ Activity: Advanced, for the curious: Things get more complicated when objects are smaller than the gyration radii of particles in the flow, or when ionization processes occur when neutral particles from an ‘atmosphere’ move into an approaching flow, or both. If you are interested in seeing how

5.5.8 A large-scale flow impinging on a fast outflow

An example of colliding plasmas on the largest scale in heliophysics involves the interplanetary medium into which the heliosphere is moving. Figure 5.1 [H-IV:3.3] “shows the prevailing picture of the global heliosphere. The structure is characterized by three flow discontinuities: the termination shock (TS), the heliopause (HP), and the bow shock (BS). The solar wind density and therefore its ram pressure falls off as r^{-2} , where r is the distance from the Sun. The wind speed is supersonic and super-Alfvénic, so when the ram pressure falls to the pressure of the ambient [interstellar medium (ISM)], the result is a shock, specifically the ellipsoidal TS in Fig. 5.1, where the flow is decelerated.

If the ISM flow is super-Alfvénic, it also encounters a shock as it approaches the Sun, specifically the roughly hyperboloid shaped bow shock in Fig. 5.1, where the ISM flow is decelerated to subsonic speeds. However, the $v_{\text{ISM}} = 23 - 27 \text{ km s}^{-1}$ interstellar flow happens to yield an Alfvénic Mach number of $M_A \approx 1$, making the existence or nonexistence of a bow shock very much an open question. Much depends on the strength and orientation of the ISM magnetic field, B_{ISM} . The higher B_{ISM} is (and the more perpendicular to the ISM flow), the lower M_A should be, and less likely that there is a bow shock.

Even the seemingly small uncertainty in v_{ISM} is enough to make a difference. For many years the best assessments were believed to be the $v_{\text{ISM}} = 26.3 \pm 0.4 \text{ km s}^{-1}$ measurement of the ISM neutral He flowing through the Solar System by *Ulysses* and the $v_{\text{ISM}} = 25.7 \pm 0.5 \text{ km s}^{-1}$ measurement from ISM absorption lines. With these relatively high values, heliospheric modelers favored $M_A > 1$, implying the existence of a bow shock. However, later He flow measurements and a new analysis of ISM absorption line data have yielded lower velocities. Specifically, measurements of neutral He flow from *IBEX* suggest $v_{\text{ISM}} = 23.2 \pm 0.3 \text{ km s}^{-1}$, and $v_{\text{ISM}} = 23.84 \pm 0.90 \text{ km s}^{-1}$ from ISM absorption lines.

This has been enough for many to argue that $M_A < 1$ should be preferred, though [it has been argued] that including He^+ density in the calculation of sound and Alfvén speeds instead of just assuming a pure proton plasma would still suggest $M_A > 1$ even if $v_{\text{ISM}} \approx 23 \text{ km s}^{-1}$. With M_A so close to 1, it is possible that the issue will not be fully resolved until an interstellar probe mission of some sort is sent out to this region. However, with M_A so close to 1 it is also possible that secondary physical processes (*e.g.*, charge exchange interactions with neutral particles) make it fundamentally ambiguous whether any boundary that may exist out there should be called a true bow shock, or whether we should instead refer to it as a ‘bow wave’.

Regardless of whether or not a bow shock exists, strong plasma interactions

these complications play out, have a look at this study of comet 67P by Behar *et al.* (2017) using observations by the Rosetta spacecraft.

prevent the ISM plasma from mixing with the solar wind plasma. The roughly paraboloid heliopause in Fig. 5.1, lying between the termination shock and the bow shock (or wave) is the contact discontinuity separating the two plasma flows. Representing the boundary between solar wind and ISM plasma, the heliopause is generally considered the true boundary of the heliosphere.”

[H-IV:3.4] “The basic structure in Fig. 5.1 is mostly defined by plasma interactions. The local ISM is partly neutral, but collisional mean free paths for neutrals are large compared to the size of the heliosphere, so their effects on heliospheric structure were long ignored. In essence, the assumption was that neutrals pass through the heliosphere unimpeded, feeling only the Sun’s gravity and photo-ionizing flux. However, in reality neutrals do participate in heliospheric interactions through charge exchange (CX). The CX interactions end up providing ways to remotely explore the heliosphere that would be impossible if the local ISM were fully ionized.

A CX interaction is a rather simple process by which an electron hops from a neutral atom to a neighboring ion (*e.g.*, $\text{H}^0 + \text{H}^+ \rightarrow \text{H}^+ + \text{H}^0$). Mean free paths for CX for most neutral ISM atoms are short enough that they do experience significant CX losses on their way through the heliosphere. The exceptions are the noble gases, which have low CX cross sections, explaining why neutral He flowing through the solar system is considered the best local probe of the undisturbed ISM flow.

Modeling neutrals in the heliosphere is very difficult because CX sends the neutrals wildly out of thermal and ionization equilibrium with the ambient plasma. Including neutrals in hydrodynamic models of the global heliosphere therefore requires either a fully kinetic treatment of the neutrals, or at least a sophisticated multi-fluid approach. The earliest models that could treat neutrals properly were from the 1990s. These models demonstrated that through CX, neutrals could have significant effects on heliospheric structure. The ISM protons are heated, compressed, deflected, and decelerated as they approach the heliopause, and thanks to CX [with the neutral component of that incoming ISM] the proton properties are at least partially imprinted on the neutral hydrogen as well, creating what has been called a ‘hydrogen wall’ of higher density [neutral hydrogen] around the heliosphere, in between the heliopause and the bow shock.” Section 10.3 (and Ch. H-IV:3) discusses stellar observations and inferences about astrospheres.

For us living deep inside the heliosphere, the consequences of the solar wind sculpting out a cavity in the interstellar medium are limited as no perturbations in the solar wind or its magnetic field can propagate against the super-Alfvénic wind. Nonetheless, the heliosphere that is shaped by this interplay does affect our exposure to cosmic rays that traverse it, see Ch. 14.

6

Magnetic (in-)stability and energy pathways

6.1 Introduction

Instabilities occur when mild perturbations to some energy reservoir provide access to an energy conversion pathway into a significantly reduced state of that reservoir. This can happen because the pathway itself develops (such as when a condition for fast reconnection is met), because the energy reservoir changes in content (as external sources insert energy), because the surrounding conditions change (for example, the direction of the solar wind magnetic field or the makeup of the solar magnetic landscape), or because a sufficiently large perturbation occurs (such as by variations in solar wind speed or the passage of a (shock) wave associated with another impulsive event). One analogy of a purely mechanical nature is the fall of a ball that is somehow nudged over the edge of a bowl; in that process, gravitational energy contained in the reservoir (the elevated ball in the bowl) is converted into the kinetic energy of the ball's fall, and ultimately into heat and waves (that themselves eventually dissipate into the microscopic kinetic energy of heat) as the ball hits the floor. The magnetic field of volumes within the Sun's atmosphere and of planetary magnetospheres can similarly destabilize: when deformed from a potential state, the added energy may be gradually dissipated thereby avoiding a (large-scale) instability or it may be stored for some time in a growing reservoir, then to be converted impulsively through a variety of pathways, eventually ending in kinetic or electromagnetic energy that is extracted from the magnetic field and the plasma that it holds. Tracking energy reservoirs and flows is often helpful in understanding processes.

The primary storage reservoir for what eventually develops into a solar impulsive event or a magnetospheric (sub-)storm is the distortion of the magnetic field away from a potential state. This elevated energy is often attributed in our thinking to electrical currents, but is stored throughout the

distorted magnetic field. {A:[76]} In quantitative terms, the maximum energy available for an impulsive event is the volume-integrated field energy in excess of the minimum level. The latter is often taken to be the potential field B_{pot} matching the observed surface field on the Sun or a reference dipole field B_{dip} for a planet, *i.e.*, A:76

$$E_{B,\text{res}} = \frac{1}{8\pi} \int \left[B^2 - (B_{\text{pot,dip}})^2 \right] dV, \quad (6.1)$$

although that potential-field energy level may not practically be achievable (see discussions of helicity in solar conditions, or consider continued stressing of the geomagnetic field by a sustained solar wind). Note that this energy is an integral quantity: the local difference $B^2 - (B_{\text{pot,dip}})^2$ quantifies the change in local energy density, but a location with a high value, for example, should not be taken as a location where strong non-potentiality originates.

The observable signatures of magnetically-dominated instabilities in the solar atmosphere and in planetary magnetospheres have led to the development of a colorful and often unclear and ambiguous array of terms, generally introduced well before the processes themselves were understood and incorporated into an overall view. Present-day understanding ascribes the impulsive and decay phases of flares, coronal mass ejections (CMEs), and terrestrial magnetospheric (sub-)storms and their counterparts in other planetary magnetospheres to a loss of a quasi-equilibrium in, or a departure from, a quasi-steady evolution of the magnetic field that leads to a rapid increase in the rate of reconnection. The latter is associated, among other things, with the acceleration of populations of ions and electrons that lead to observable emissions in a wide range of wavelengths in the electromagnetic spectrum, among them the terrestrial auroral emissions and their solar counterpart, the flare ribbons. {A:[77]} A:77

Not only do these impulsive phenomena share many physical processes, they are also links in the chain of Sun-Earth connections: many of the more energetic solar field destabilizations are associated with both flares and CMEs (see Table 6.1), while CMEs that envelop a planet that has an intrinsic magnetic

⁷⁶ Activity: Consider how a non-potential state can arise or be strengthened in the solar atmosphere and in a magnetosphere, including the roles of plasma motions and induction. Eq. (4.1) is illustrative for the overall energy budget.

⁷⁷ Activity: The processes of electromagnetic radiation from a plasma involve three fundamentally distinct processes: bound-bound, free-bound (radiative recombination), and free-free (Bremsstrahlung) emission. Aurorae and flare ribbons are caused by collisions of downward-propagating, energetic charged particles with the atmosphere below. Aurorae observed from the ground include both free-bound and bound-bound emission from ions and molecules, respectively (there is X-ray emission, too, but that does not penetrate to ground level). Look up which ions and molecules dominate in the terrestrial aurora, and which emission processes are involved with these. See also Activity 121 for the contrast with solar coronal emission.

Table 6.1. *Solar flare classifications. [Listed are the GOES flare class, the corresponding flare-peak irradiance at the top of Earth's atmosphere, the class and surface footprint based on chromospheric emission patterns, the fraction of such events associated with coronal mass ejections (CMEs), and the characteristic frequency of such events during the maximum and minimum of a typical solar cycle. Table H-II:5.1]*

GOES class	1-8Å peak (W/m^2) ($kerf/cm^2/s$)	H α class (percent)	H α Area (<i>Millionths</i> of hemisphere)	CME fraction ^a	Events/year (<i>cycle max./min.</i>)
A	$>10^{-8}$	-	-	-	-
B	$>10^{-7}$	S	<200	-	-
C	$>10^{-6}$	1	>200	0.2	>2000/300
M	$>10^{-5}$	2	>500	0.5	300/20
X	$>10^{-4}$	3	>1200	0.9	10/one?
-	$>10^{-3}$	4	>1200	1.0	few?/none?

^a (approximate values)

field often trigger (immediate or delayed) magnetospheric activity. {A:[78]}
{A:[79]}

A:78

A:79

6.1.1 Introducing solar flares and coronal mass ejections

[H-II:5.1] “A solar flare is narrowly defined as a sudden atmospheric brightening, traditionally in chromospheric H α emission [(at 656 nm, associated with a 3 \rightarrow 2 level transition of hydrogen atoms, and thus the lowest-energy transition in the Balmer series)] but more practically now as a coronal soft X-ray source [(Fig. 6.1 summarizes the common names used for wavelength bands from radio to gamma rays)]. The physical processes resulting in a flare include restructurings of the magnetic field, non-thermal particle acceleration, and plasma flows. Flares have intimate relationships with other observable phenomena such as filament eruptions, jets, and coronal mass ejections [...]

The energy released in a solar flare is dominated by particle acceleration,

⁷⁸ Activity: The average speed of a CME between Sun and Earth is close to 500 km/s while the fastest have speeds exceeding 3000 km/s. How long are the transit times from Sun to Earth? Compare the average and peak CME speeds to typical wind speeds (Table 2.4). Describe qualitatively what happens in the interaction with slow and fast wind streams for average CMEs and for the fastest CMEs.

⁷⁹ Activity: The phenomena discussed in Ch. 6 are all part of what is referred to as space weather. To explore how aspects of space weather are quantified review this NOAA site that lists the types of ‘storms,’ their potential effects, and their approximate frequency within a solar cycle. For current space weather conditions, forecasts, and more see this site of the Space Weather Prediction Center.

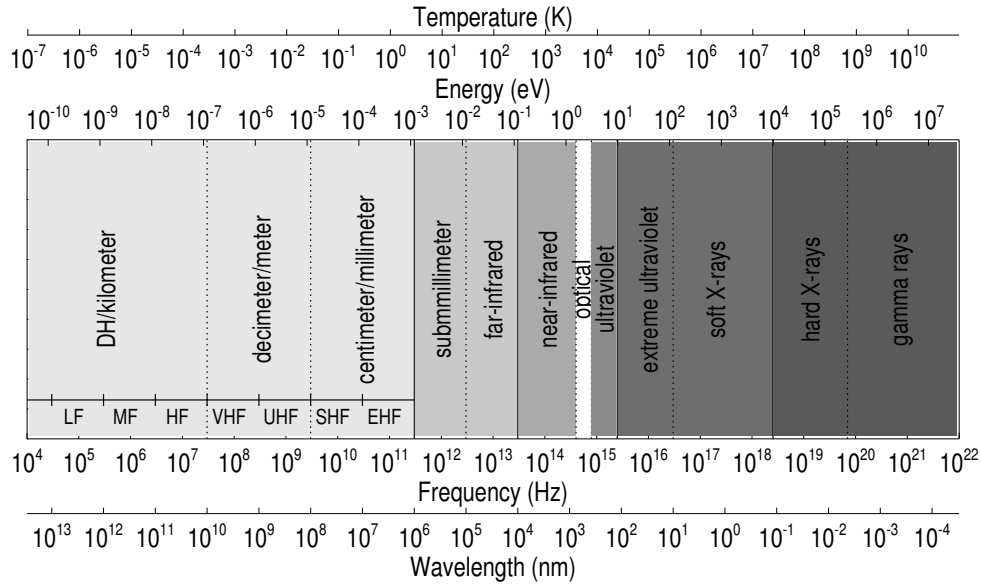


Fig. 6.1. Overview of the electromagnetic spectrum with energy in electron volts and the equivalent temperature in Kelvin (top axes), frequency in Hertz, and wavelength in nanometers (bottom axes). Note that the AM band lies in the low- and medium-frequency (LF-MF) range and the FM band in the very-high frequency (VHF) range. [Fig. H-II:4.1]

both of electrons and ions. This means that the most direct observations are in the X-ray and γ -ray domains; note that non-thermal processes also usually dominate the emission signatures in the radio range (10^7 - 10^{12} Hz; meter-submillimeter wavelengths). Please refer to Ch. H-II:4 for a fuller discussion of the remote-sensing signatures. We will simply comment here that in general the hard X-ray spectrum ($h\nu \gtrsim 10$ keV [or wavelengths shortward of about 1 Å]) is dominated by electrons of this energy or greater, while the soft X-ray spectrum ($h\nu \lesssim 10$ keV) is dominated by the free-bound and bound-bound transitions of a thermal plasma with assumed Maxwellian distribution functions, and also usually assuming the electron and ion temperatures to be equal, *i.e.*, $T_e = T_i$. The free-bound process (radiative recombination) may also contribute to the hard X-ray spectrum under certain conditions.”

A common observed pattern, most frequently in eruptive flares associated with coronal mass ejections into the heliosphere, is that a volume of the corona over a magnetic polarity inversion line expands explosively (often involving a large-scale shock front) as the hard X-ray and γ -ray emission brightens impulsively, with two (or more) ribbon-like brightenings at chromospheric and

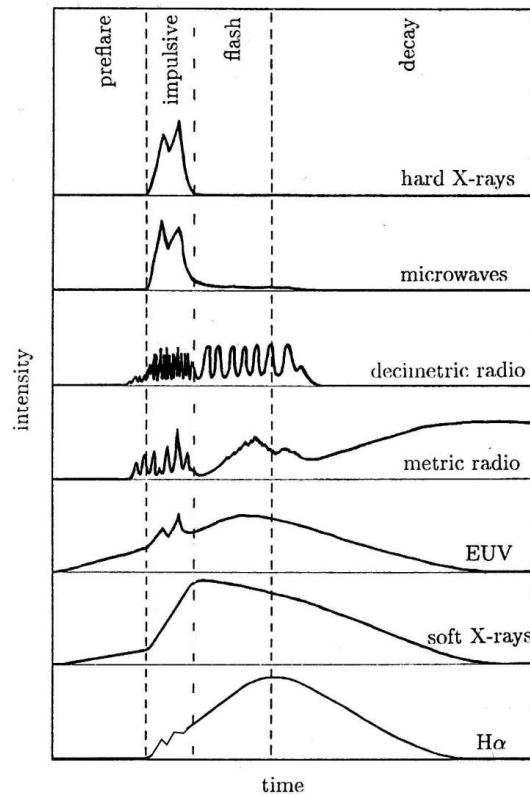


Fig. 6.2. Schematic view of the evolution of solar flare emissions in different wavelengths, showing the intermingling of impulsive-phase and gradual-phase signatures across the spectrum. Note the wide variety of radio signatures. [In wide wavelength bands in the visible, the emission peaks in the impulsive phase of the flare. Fig. H-II:5.1; source: Benz (2002).]

photospheric levels propagating away from the polarity inversion line, with a coronal mass ejection moving away while behind it the corona fills with heated plasma from the lower atmosphere, which then cools by radiation in the soft X-ray and EUV bands and by conduction into the lower, cooler atmosphere. See Fig. 6.2 for the characteristic evolution of a flare in wavelength space, and Fig. 6.3 for a sketch of the various emissions throughout the EM spectrum.

[H-II:5.2.1] “The modern view of [solar flares] is via the soft X-ray monitoring by the GOES and other ‘operational’ spacecraft. We now routinely classify solar flares by their GOES classes: A, B, C, M, and X in decades, with the X class signifying 1-8Å energy fluxes greater than 10^{-4} W/m², on the order of 0.01% of the solar luminosity. Table 6.1 summarizes these and other properties, with very approximate correspondences between the [chromospheric] H α and

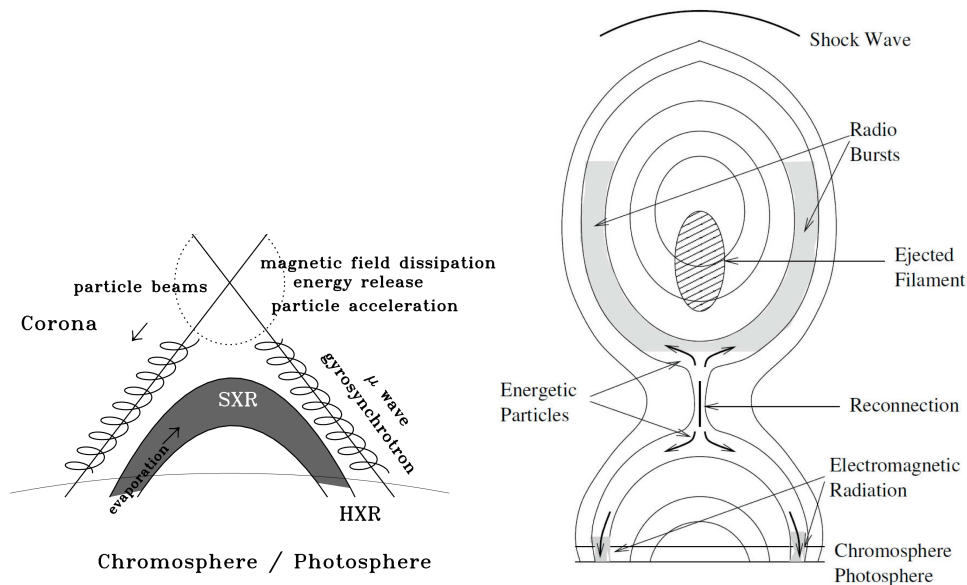


Fig. 6.3. (left) Schematic arrangement in the outer atmosphere of the Sun or a comparable cool star indicating the flow of energy during a flare: a flare involves magnetic reconnection high in the atmosphere which accelerates particles, leading to motion along field lines upward away from or downward towards the visible surface. Resulting emissions include hard X-rays (HXR), soft X-rays (SXR), and microwave emission. [Fig. H-IV:2.1] (right) Concepts of particle acceleration and emissions in a solar [eruptive] event. [Fig. H-IV:12.6; source: (Kallenrode, 2003).]

GOES X-ray systems, and very approximate ranges for the number of flares that occur per year at maximum and minimum of the solar cycle.”

6.1.2 Introducing geospace (sub-)storms

In contrast to a flare or CME observed by imaging the electromagnetic radiation, a terrestrial magnetic storm is typically observed by sampling magnetic field changes at the Earth’s surface or tracked by monitoring energetic particles and their effects (such as in aurorae). A terrestrial [H-II:10.2] “magnetic storm is defined nowadays by the time variation of the geomagnetic Dst index, illustrated schematically in Fig. 6.4. The Dst index is a measure of a quasi-uniform magnetic disturbance field near the Earth, aligned with the dipole axis (northward for $Dst > 0$), such as would be produced by a ring of electric current (westward if $Dst < 0$) near the equatorial plane. A prolonged (hours to days) interval of negative Dst values constitutes a magnetic storm. The peak negative excursion is often taken as a measure of storm intensity: Dst -30 nT to

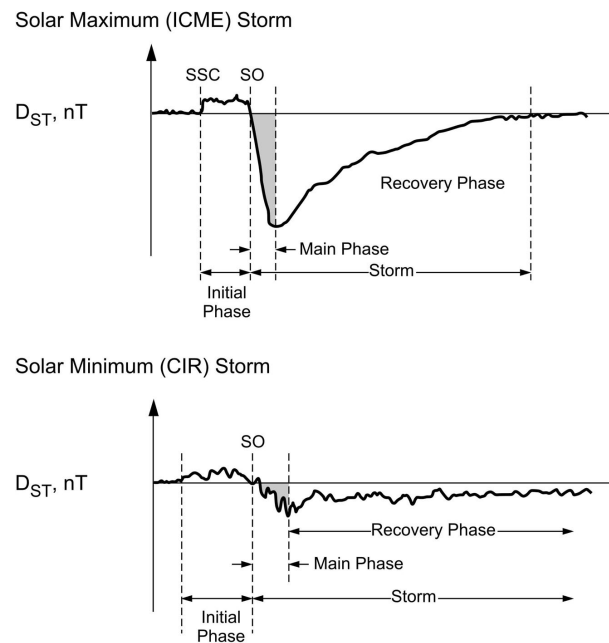


Fig. 6.4. Schematic time history of geomagnetic field variation for two characteristic magnetic storms. Time range: several days. Vertical variation range: $\sim 100 - 200$ nT ($\sim 1 - 2$ mG). SSC: storm sudden commencement. SO: storm onset. The top panel shows the storm development in response to a characteristic interplanetary coronal mass ejection (ICME), and the bottom panel that for the passage of a corotating interaction region (CIR). [Fig. H-II:10.1; source: Tsurutani et al. (2006).]

-50 nT are weak storms, -50 nT to -100 nT moderate, and over -100 nT intense; storms over -300 nT occur at most a few times during a solar cycle (Earth's dipole field at the equator is $\sim 31,000$ nT, [or 0.31 Gauss,] for comparison) [...]

[T]he field depression quantified by Dst is the result of plasma pressure that inflates the dipole field. The essential phenomenon of the magnetic storm is thus the addition of a large amount of plasma energy to the dipolar field region of the magnetosphere. Furthermore, it is now well established that this energy addition results from a particular condition in the solar wind: 'a sufficiently intense and long-lasting interplanetary convection electric field', meaning $-\mathbf{v} \times \mathbf{B}/c$, for the [interplanetary magnetic field's (IMF's)] southward component.

In contrast to the magnetic storm, there is much less unanimity on what defines a magnetospheric substorm. Probably the most spectacular phenomenon, and the one most widely used as a unifying concept, is the auroral substorm, [...], which has a characteristic temporal development. Early to show up are] the auroral forms (light-emitting regions) during what is called the *expansion*

phase of the substorm: beginning with an initial brightening at the lowest latitudes near midnight (*onset*), the aurora intensifies greatly, becomes very complex in spatial structure (*auroral breakup*) and expands, predominantly westward and poleward but also eastward, eventually subsiding in a *recovery phase*. This auroral development is accompanied by strong geomagnetic disturbances (commonly reaching [a Dst index of] ~ 1500 nT [or 0.015 G] and more), with a spatial distribution almost as complex as that of the aurora but describable roughly as equivalent to a current above the Earth (*auroral electrojet*) that is westward near and before midnight and eastward after midnight. [Essentially the same sequence occurs simultaneously in the two hemispheres,] at the (more or less) magnetically conjugate locations [...]"

6.2 Terrestrial magnetospheric disturbances

6.2.1 Energy pathways and reservoirs

[H-II:10.3.2] “For the magnetosphere and the upper regions of the ionosphere [...] the solar wind is the only significant external source of energy available [because these regions are completely transparent to solar radiation ...] An interior source of energy available for a planetary magnetosphere is planetary rotation [...]"

When considering the solar wind as the energy source, only the kinetic energy of plasma bulk flow is of importance; the thermal and magnetic energies of the solar wind can be neglected [...: they are relatively small to begin with (Sect. 3.5.2) while moreover] at the bow shock they are overwhelmed by additional thermal and magnetic energies extracted from the flow. Furthermore, to transfer magnetic energy across the magnetopause requires [...] a tangential component of the electric field which interacts with the magnetopause current to extract more mechanical energy from the plasma [...] The interplanetary magnetic field does exert a dominant influence on energy conversion processes in a planetary magnetosphere, but primarily by control of magnetic reconnection processes and open field lines [...]"

[H-II:10.3.3] “The following are among the principal loss and dissipation processes in planetary magnetospheres, energy being lost primarily to the atmosphere in (1) and (2) and being removed outside the system (to ‘infinity’) in (3) and (4):

(1) **Collisional and Joule heating in the ionosphere.** If the bulk flow of plasma differs from the bulk flow of the neutral atmosphere (usually as a consequence of magnetospheric dynamics), there is energy dissipation given by $\mathbf{E}^* \cdot \mathbf{J}$, where \mathbf{E}^* is the electric field in the frame of reference of the neutral atmosphere. This is commonly referred to as ‘ionospheric Joule heating’; [...] it

is primarily frictional heating by collisions between plasma and neutral particles, Joule heating in the true physical sense ($\mathbf{E}' \cdot \mathbf{J}$, where \mathbf{E}' is the electric field in the frame of reference of the plasma) contributing only a small fraction of the total. The energy is removed from the magnetic field and converted (via kinetic energy of relative bulk flow as an intermediary) to heat (thermal energy), with the heating rate per unit volume partitioned approximately equally between plasma and neutrals.

(2) **Charged-particle precipitation.** Energetic charged particles that enter the atmosphere from above are usually said to be *precipitating*. They penetrate the atmosphere to a depth that increases with increasing energy, until their energy is lost, going partly into heating the atmosphere and partly into ionization or other interactions.

A:80

One source of precipitating particles is simple loss from the radiation belts or from the ring current and plasma sheet regions; $\{A:^{[80]}\}$ the energy deposited in the atmosphere is taken from the mechanical (thermal) energy of the respective magnetospheric particle populations. In addition to these particles that precipitate merely because their velocity vectors are oriented in the appropriate direction, there are other sources of precipitating charged particles, in which the energy and the intensity of the particles have been enhanced by an acceleration process. In particular, the auroral phenomena that occur in nearly all of the planetary magnetospheres observed to date are generally interpreted as resulting from some special acceleration process that supplies the required intensities of precipitating charged particles. A widely accepted model, developed from extensive studies at Earth and applied to aurora at Jupiter and at Saturn, ascribes auroral acceleration to Birkeland (magnetic-field-aligned) electric currents accompanied by electric fields parallel to the magnetic field; the rate of energy supply to the precipitating particles is $E_{\parallel} J_{\parallel}$, hence the added energy is taken out of the magnetic field (in this model, an aurora occurs only when the Birkeland current is directed upward, corresponding to electron motion downward). Auroral acceleration has also been associated with intense Alfvénic turbulence (which contains fluctuating Birkeland currents) [...]

(3) **Emission of electromagnetic radiation.** A variety of processes in planetary magnetospheres produce electromagnetic radiations of various types: atomic and molecular line emissions (from the aurora and from magnetospheric interactions with plasma and neutral tori), radio waves (wideband and narrowband), a veritable zoo of plasma waves, and even X-rays (bremsstrahlung from precipitating electrons and, possibly, nuclear line emissions excited by very energetic precipitating particles). [...] As far as the energetics of planetary

⁸⁰ Activity: Look up locations and properties of the Earth's (a) electron and proton radiation belts, (b) ring current, and (c) plasma sheet.

magnetospheres are concerned, however, the amount of energy involved is negligibly small for most emissions, with only a few exceptions (UV radiation from the Io torus at Jupiter).

(4) **Energetic neutral particle escape.** Neutral particles that remain within a magnetosphere must be gravitationally bound to the planet; plasma particles within the magnetosphere, on the other hand, typically have speeds that exceed (often by a large factor) the gravitational escape speed — plasma is held within the magnetosphere by the magnetic field, not by gravity [...] Charge-exchange collisions between ions and neutrals, in which the outgoing neutral has the velocity of the incoming ion and vice versa, thus produce fast neutrals that escape from the system immediately, with their kinetic energy. This process represents a loss (generally by quite significant amounts) both of neutral particles and of energy from the magnetosphere.

(5) **Dissipation processes in the magnetosphere.** In regions of the magnetosphere with major departures from the MHD approximation (particularly where magnetic reconnection is occurring) dissipative processes such as Joule heating associated with effective resistivity may be significant. The primary effect is not energy loss but enhancement of conversion from magnetic to thermal energy.”

[H-II:10.3.4] “The field approach to energy implies that energy may be regarded as *stored* in space [...] The primary reservoir of stored mechanical energy in a planetary magnetosphere is the thermal energy of its various plasma structures, especially the *plasma sheet* of the magnetotail or magnetodisk, the *ring current*, and the plasma and neutral *tori* associated with the planet’s moons; the kinetic energy of bulk flow of magnetospheric plasma also plays a role, particularly for plasma tori and in the case of rapid changes [...]

The primary reservoir of stored electromagnetic energy of importance for a planetary magnetosphere is the energy of the magnetic field [...] Because the energy of the planetary dipole field itself does not change (except on time scales of the secular variation, $\sim 10^2 - 10^3$ years for Earth) and thus has no effect on the energetics of the magnetosphere, a convenient measure of stored electromagnetic energy is the energy of the total magnetic field minus the (unchanging) energy of the dipole field, [reflected in Eq. 6.1].

The stored gravitational energy can be changed only by a net radial displacement of matter; any such effects in the magnetosphere are for the most part negligible in comparison to changes of mechanical or magnetic energy.”

6.2.2 What leads to explosive energy releases?

[H-II:10.5] “The discussion so far has ignored time variations and has proceeded on the tacit assumptions that all the energy supply, conversion, and dissipations

processes are more or less in balance. There is no general requirement for this to be the case, and in fact often it is not the case [...] The prototypical example is kinetic energy from the solar wind being converted into magnetic energy of the magnetotail at an increased rate due to enhanced dayside reconnection (in response to changed solar-wind conditions), but the rate of removal by conversion of magnetic energy into mechanical energy of magnetospheric plasma plus escape down the magnetotail not being equally enhanced (for reasons that need to be identified); in this case, the magnetic energy reservoir increases with time and reaches a point at which (again, for reasons that need to be identified) the magnetic energy content can no longer be maintained but must be converted to other forms.”

First, let us look at topological changes involved in magnetospheric processes. [H-II:10.5.1] “[M]agnetic flux transport and the increase of magnetic energy by stretching the field play an important role in supplying energy to the magnetosphere. Non-equilibrium configurations of the magnetotail that change the magnetic topology and allow different paths of flux transport are therefore of particular interest.

A simple sketch of a model widely invoked to interpret magnetospheric substorms at Earth is shown in Fig. 6.5, which displays a time sequence of magnetospheric configurations. Each panel shows the magnetic field line configuration in the noon-midnight meridian plane (left) as well as the configuration of magnetic singular X [lines (where field vectors of opposite directions cross; see, for example, Fig. 5.14)] and O lines [(around which field lines loop)] in the equatorial plane (right) and projected to the ionosphere (top); the equatorial projection, [...] is essential for describing the three-dimensional structure of the magnetic field. Panel 1 is the simplest topology of the open magnetosphere [(compare with Fig. 5.14)]. In panel 2, a small volume usually called a *plasmoid* appears deep within the closed-field-line region, bounded on the earthward side by a newly formed *near-Earth X-line* (NEXL) and threaded by magnetic field lines that encircle the attached O-line; ideally, the field lines are confined within the plasmoid and connect neither to the Earth nor to the solar wind (what the real topology is, however, is still uncertain). For the ideal topology, the plasmoid can be visualized in three dimensions as shaped roughly like a banana, oriented approximately dawn to dusk and tapering to zero thickness at both ends, with the X-line on its surface and the O-line running through the middle of its volume. The plasmoid grows (panel 3) by magnetic reconnection until it touches the separatrix of the open field lines (panel 4, *onset of lobe reconnection*); afterward (panel 5), the plasmoid is on interplanetary field lines and is carried away (presumably) by the solar wind.

A model of topological changes for a rotation-dominated magnetosphere

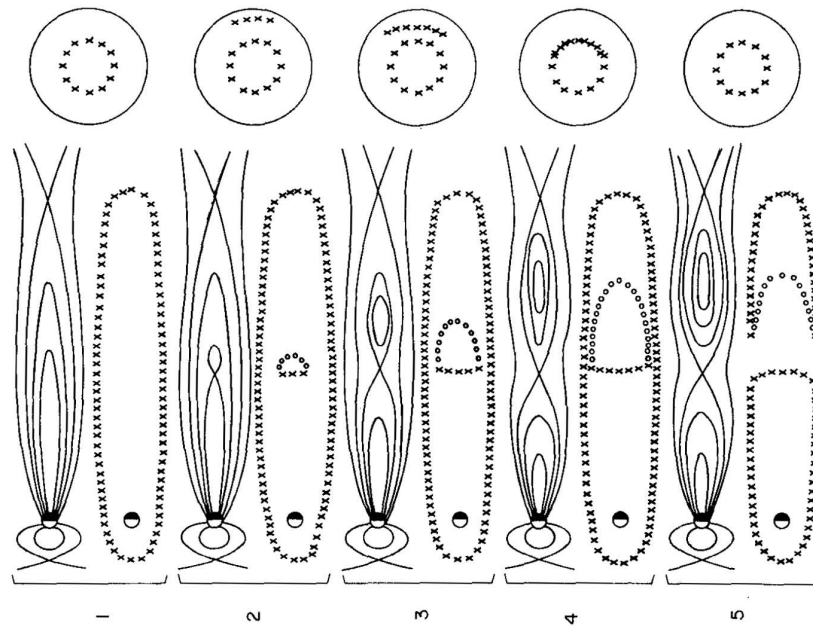


Fig. 6.5. Possible changes of the magnetic field topology in the magnetotail of a solar-wind-dominated magnetosphere. The diagram is shown [with the Sun's direction at the bottom] to facilitate comparisons with diagrams of filament eruptions [...]. Each panel in the sequence shows a side view of the magnetic field (left), the outline of the X lines [where field of opposite directions meet] seen from above the north pole (right), and a top-down view of the mapping of the reconnection region onto the Earth (top). [Compare this to the sketch of a solar eruption in Fig. 6.6; Fig. H-II:10.5; source: Vasyliunas (1976).]

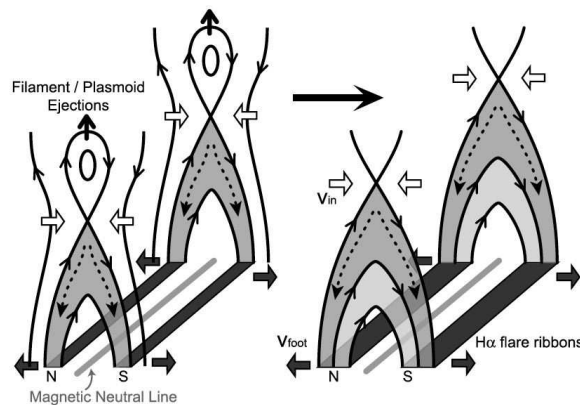


Fig. 6.6. How the ribbon motion sweeps out magnetic field during the reconnection process in the standard model. [Compare this solar eruption to the sketch of the magnetospheric substorm (righthand side of each panel) in Fig. 6.5; Fig. H-II:5.16; source: Asai et al. (2004).]

[differs] only in three respects: (1) the time sequence has been translated into an azimuthal-angle sequence, (2) field lines are stretched by the outflow of plasma from an internal magnetospheric source (planetary/magnetospheric wind) [...], (3) there are no counterparts to panels 4 and 5, since field lines connected to the solar wind are not considered. [...]"

Next, let us look at the role of instabilities in causing rapid changes in topology. [H-II:10..5.2] "Instabilities have attracted much attention as a possible way of inducing rapid change from equilibrium to non-equilibrium configurations – an alternative to straightforward evolution to non-equilibrium as the result of changing boundary conditions. [...]"

(1) **Tearing-mode instabilities.** 'Tearing mode' is a generic term for instabilities that result in the reconnection of initially oppositely directed magnetic fields. They are obvious candidates for initiating topological changes of the magnetotail (in particular, those envisaged in Fig. 6.5).

(2) **Current-driven instabilities.** The concept that a sufficiently intense electric current may bring about its own breakdown, by creating conditions that impede current flow, was first suggested [...] as a model for solar flares. Under the name 'current disruption' it has been widely discussed as a model for substorm onset and expansion. Various instabilities that develop when the current density exceeds some threshold value have been proposed.

(3) **Interchange and ballooning instabilities.** Interchange instabilities which do not appreciably change the magnetic field are thought to be essential for plasma transport in rotation-dominated magnetospheres. Ballooning instabilities can be viewed roughly as interchange that does change the magnetic field. As a model for substorms, they have been invoked particularly at the transition between the dipole field and the magnetotail, in several variants."

6.2.3 Terrestrial magnetospheric substorms

A substorm can be summarized as a two-stage process. [H-II:10.6.1] "Stage 1 (growth phase): as a consequence of a southward interplanetary magnetic field, the configuration of the magnetosphere changes, its magnetic field becoming highly stretched (increased magnetic flux in the magnetotail, reduced flux in the nightside equatorial region). Stage 2 (expansion phase, initiated by the onset): the magnetic field changes to more nearly dipolar (increased flux on the nightside), and there is enhanced energy input and dissipation to the inner magnetosphere and the ionosphere/atmosphere; the process occurs on dynamical time scales (comparable to or shorter than wave travel times) and is accompanied (most probably) by changes of magnetic topology.

[In terms of energy flow paths: during stage 1, power input from the bulk flow kinetic energy of the solar wind is enhanced and is appreciably larger than

the sum of power outputs due to heating, bulk motion, and plasma escape from the geomagnetic system. During stage 2, energy flow into mechanical energy of plasma and particularly plasma heating are enhanced; power flow through plasma and field escape and field reconfiguration presumably are enhanced in connection with topological changes]

The substorm growth phase is in essence the increase of open magnetic flux in the magnetosphere, which occurs for a two-fold reason. First, the flux addition rate at the dayside reconnection region increases as the solar wind transports more magnetic flux, of the sense opposite to the terrestrial dipole flux, toward the magnetosphere; the reasons for this are assumed to lie in the physics of magnetic reconnection. Second, the flux return rate at the nightside reconnection region does *not* increase to match the addition rate; the reasons for this are not at all well understood. [...] Within the magnetosphere, the net effect of the substorm growth phase is to remove magnetic flux from the nightside magnetosphere by flow toward the dayside reconnection region and to add magnetic flux to the magnetotail (enhanced stretching of magnetotail field lines).

The substorm expansion phase does return the magnetic flux, rapidly and spectacularly, from the magnetotail to the nightside magnetosphere (dipolarization of a previously stretched tail-like field); given that plasma in the magnetotail beyond a distance typically $\sim 15 - 20$ Earth radii is observed to flow away from Earth, the process must almost unavoidably proceed by topological changes of the type sketched in Fig. 6.5. The energy input into plasma, energetic charged particles, and the aurora can be largely accounted for by adiabatic compression and Birkeland current effects. What remains highly controversial is how the process starts and why it is so sudden and catastrophic [...]

A further complication is the question of external versus internal influences. That the growth phase is initiated by changing solar wind conditions is the consensus view. The onset and expansion phase, on the other hand, are regarded by the majority as basically the result of internal dynamical processes, although subject to solar wind influences (*e.g.*, if the system is evolving toward instability, it may be pushed over the threshold by a change in the solar wind). A substantial minority, however, considers the substorm onset intrinsically as triggered by a solar wind change (typically toward a more northward interplanetary magnetic field)."

6.2.4 Terrestrial magnetic storms

[H-II:10.6.2] "Our understanding of magnetic storms has been decisively influenced by a remarkable theoretical result, the Dessler-Parker-Sckopke theorem,

which relates the external magnetic field at the location of a dipole to properties of the plasma trapped in the field of the dipole. [T]he theorem states that $\mathbf{b}(0)$, the magnetic disturbance field of external origin at the location of a dipole of moment $\boldsymbol{\mu}_B$ [in an undisturbed state], satisfies

$$\boldsymbol{\mu}_B \cdot \mathbf{b}(0) = 2U_K \quad (6.2)$$

where U_K is the total kinetic energy content of plasma in the magnetosphere. What is remarkable is that the right-hand side does not depend on the spatial distribution, the partition between bulk-flow and thermal energy, or any properties of the energy spectrum. [...]

Although $\mathbf{b}(0)$ nominally is evaluated at the center of the Earth, it is also equal to the (vector) average of $\mathbf{b}(\mathbf{r})$ over the surface of the globe (by a theorem for solutions of Laplace's equation, satisfied within the globe by each Cartesian component). The Dst index is the average, over a low-latitude strip of the globe, of the disturbance field component aligned with the dipole; after some corrections (chiefly removing the contribution from induced earth currents), -Dst may be considered a reasonable proxy for the left-hand side of Eq. (6.2), as long as $\text{Dst} < 0$. The Dessler-Parker-Sckopke theorem then provides a method of inferring the plasma energy content simply from the value of the Dst index. [...] Direct *in situ* observations have established that the greater part of the energy resides in what is called the ring current region.

Geomagnetic storms, particularly the intense ones, are characterized by unusually large amounts of energy stored as mechanical energy of plasma in the ring current region, in comparison to other storage regions. This implies that during the development of an intense storm the power [going from the magnetic reservoir to the ring-current plasma kinetic reservoir] is unusually large, on the average. Whether this enhanced conversion rate from magnetic energy into mechanical energy of ring current plasma [...] results from a different interaction process or simply from a different time sequence of solar wind parameters is an unresolved question. More specifically, can the energy for storms be supplied by a sequence of substorms (perhaps unusually frequent and/or unusually intense), or is some other process required? A related question is that of *geoeffectiveness*: when interplanetary structures such as CME's impinge on the Earth, under what conditions do they produce intense magnetic storms? (prolonged southward B_{sw} is one that is well established).”
 {A:[81]}

A:81

⁸¹ Activity: The energy processed by the magnetosphere during a magnetic storm is of order $E_{\text{storm}} = 5 \times 10^{23} - 5 \times 10^{24}$ erg from moderate storm to superstorm. Compare that to an order of magnitude estimate of the energy $E_{\text{mag}, \oplus}$ contained in the geomagnetic field (by, say, using a scale of $3R_{\oplus}$ and a characteristic field strength of 0.1 G) and with the incoming total energy E_{sw} of the solar wind during the storm period (with typical conditions for the fast solar wind and an active cross section of

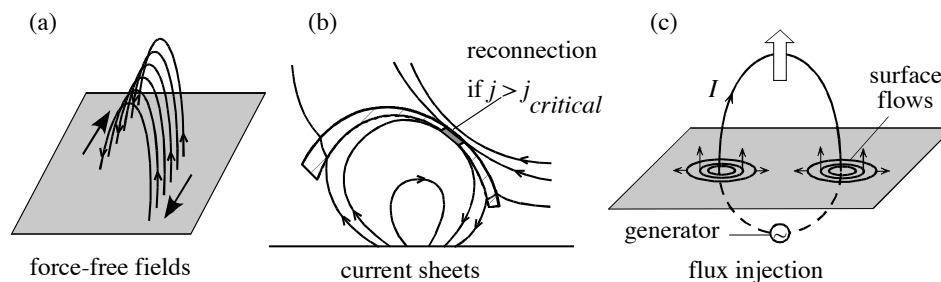


Fig. 6.7. Schematic illustration of three different types of models that use magnetic energy to power a flare or CME. Panel (a): Magnetic energy is stored in the corona in the form of field-aligned currents that eventually become unstable. Panel (b): Magnetic energy is stored in the corona in the form of a thin current sheet that is suddenly dissipated when a micro-instability is triggered within the sheet. Panel (c): An example of a directly driven flare model. Here magnetic flux is suddenly injected from the convection zone into the corona at the onset of the flare or CME. Such a model produces a well-organized flow pattern during the impulsive phase (small arrows at surface in panel c). [Fig. H-II:6.8]

6.3 Solar impulsive events

6.3.1 The magnetic reservoir

6.3.1.1 Storage models

[H-II:6.2.1] “Although it is generally agreed that flares and CMEs derive their energy from the Sun’s magnetic field, exactly how the magnetic energy is extracted remains uncertain. One possibility is that a flare or CME occurs when a slowly evolving coronal magnetic field reaches a point where a stable equilibrium is no longer possible. The slow evolution of the corona is driven by the changes continually occurring in the photospheric field as a result of solar convection; [in phases before solar impulsive events, these processes build up the stored magnetic energy]. The equilibrium may disappear altogether or, alternatively, a stable equilibrium may simply become unstable. The continual emergence of new flux from the convection zone and the shuffling of the footpoints of closed coronal field lines increase the free magnetic energy in the corona. Eventually, these stresses may exceed a threshold beyond which a stable equilibrium cannot be maintained, and the field erupts. Models based on this principle are often referred to as *storage models*.

[...] Because the plasma in the photosphere is almost 10^9 times denser than the plasma in the corona, it is difficult for disturbances in the tenuous corona to have much effect on the photosphere and the deeper layers below it. Field

πR_{CF}^2 , and a storm duration of 1-10 h). What are the values of $E_{storm}/E_{mag,\oplus}$ and E_{storm}/E_{sw} ? Compare these values to solar equivalents when you reach Activity 85.

lines mapping from the corona to the photosphere are thus said to be 'inertially line-tied' which means that the footpoints of coronal field lines are essentially stationary over the time scale of the eruption [...]

Unlike models of confined flares, models of CMEs must be able to explain not only the release of magnetic energy, but also how mass is ejected into interplanetary space. During a CME, magnetic field lines mapping from the ejected plasma to the photosphere are stretched outwards to form an extended, open field structure [...]" that resembles the sketches on the left of each of the panels in Fig. 6.5: plasmoids leaving the magnetosphere have been compared to filaments erupting as part of a CME.

6.3.1.2 Directly driven models

[H-II:6.2.2] Some researchers "have proposed models that produce a sudden energy release in the corona by means of a surface or sub-surface current generator. In contrast to storage models, there is no build-up of magnetic energy in the corona prior to onset. Instead, there is a sudden injection of current or magnetic flux into the corona from below. As a rule, the models do not address the mechanism that leads to the sudden injection of current or flux. They simply posit that such an injection occurs, and then model the consequences of such an injection for the corona." Section H-II:6.2.2 describes some of these models and their problems when compared to observations and physical conditions in the Sun; these are not further discussed here.

6.3.1.3 Pre-eruption current sheet models

[H-II:6.2.3] "Because the magnetic energy in the corona is much larger than the thermal and gravitational energies, the magnetic force ($\mathbf{j} \times \mathbf{B}$) cannot, in general, be balanced by gravity or by a gas pressure gradient. Thus, as a rule, the coronal field will tend to be force-free, meaning that the current will flow along the direction of the magnetic field (*cf.* Fig. 6.7a). An exception to this rule occurs when a current sheet is present. In this case gas pressure within the sheet balances the strong magnetic field outside. If the current sheet is sufficiently thin, then the high temperature or density within the sheet may not be detectable. Thus the corona could still have the appearance of a plasma with a low gas to magnetic pressure ratio (*i.e.*, plasma $\beta \ll 1$). Figure 6.7b shows a flare model with such a current sheet, where a micro-instability within the sheet triggers an eruption.

Prior to onset, the current sheet grows as a consequence of the emergence of new magnetic flux into a pre-existing magnetic loop as shown in Fig. 6.7b. As the current sheet grows, it eventually reaches a point where a micro-instability is triggered because the current density exceeds some critical value. Once the

micro-instability occurs, the electrical resistivity of the plasma in the sheet dramatically increases, and rapid reconnection ensues.”

6.3.2 Two-dimensional force-free models

[H-II:6.2.4] “[M]any storage models use configurations that have currents flowing parallel to the magnetic field in the pre-eruption state. Thus, there is no [net] magnetic force anywhere in the configuration prior to eruption. To explain an eruption, such models need to show how a strong magnetic force can rapidly appear as a result of the slow evolution of the photospheric boundary conditions.

To illustrate the basic principles, we first consider a relatively simple flux-rope model [. . . for which the external field is] prescribed by

$$B_y + iB_x = \frac{2iA_0\lambda(h^2 + \lambda^2)\sqrt{(\zeta^2 + p^2)(\zeta^2 + q^2)}}{\pi(\zeta^2 - \lambda^2)(\zeta^2 + h^2)\sqrt{(\lambda^2 + p^2)(\lambda^2 + q^2)}}, \quad (6.3)$$

where $\zeta = x + iy$ and A_0 is the photospheric magnetic flux, or, equivalently, the magnetic vector potential at the origin. In this expression h is the height of the flux rope above the surface and p and q are the lower and upper tips of a vertical current sheet below the flux rope as shown in Fig. 6.8. The parameter λ is the half-distance between two photospheric field sources located at $\zeta = \pm\lambda$ on the surface [. . .]

Application of the frozen-flux condition at the surface of the flux rope determines the current in the rope. This condition keeps the magnetic flux between the flux rope and the surface constant in time. It also ensures that during an eruption there is no flow of energy into the corona if the normal component of the field at the base remains invariant. Consequently, the current in the flux rope is prescribed by

$$I = \frac{c\lambda A_0}{2\pi h} \frac{\sqrt{(h^2 - p^2)(h^2 - q^2)}}{\sqrt{(\lambda^2 + p^2)(\lambda^2 + q^2)}}. \quad (6.4)$$

This current decreases with time during an eruption as magnetic energy is converted into kinetic energy. This decrease becomes apparent only when the formula giving the dependence of q upon h and p is incorporated into the above expression [(for references, see Sect. H-II:6.2.4)].

The magnetic field configuration is shown in Fig. 6.8 for three different sets of parameters. The surface at $y = 0$ corresponds to the photosphere, and the boundary condition at this surface is

$$A(x, 0) = A_0 \mathcal{H}(\lambda - |x|), \quad (6.5)$$

where \mathcal{H} is the Heavyside step-function and A_0 is the value of A at the origin.

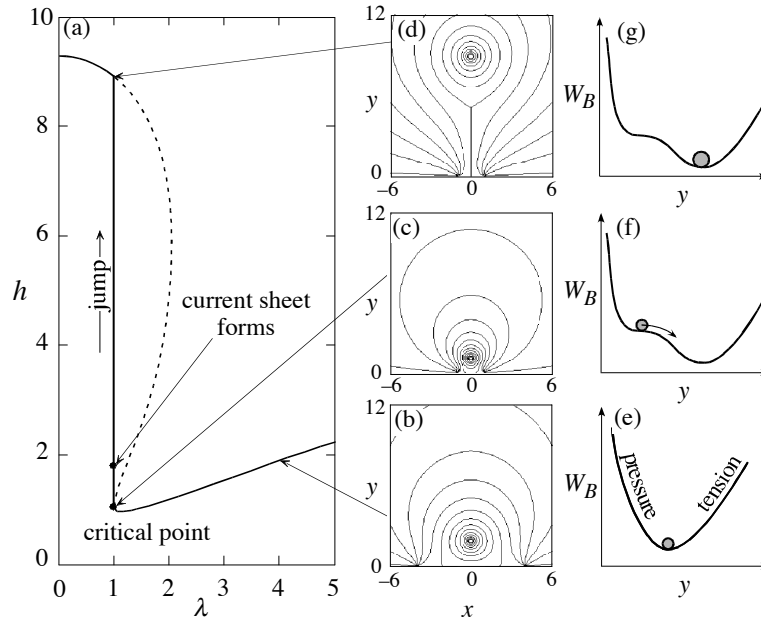


Fig. 6.8. Ideal-MHD evolution of a two-dimensional arcade containing a magnetic flux rope. Panel (a) shows the equilibrium curve for the flux rope height, h , in normalized units, as function of the source separation half-distance λ . Panels (b), (c), and (d) show the magnetic field configuration at three different locations on the equilibrium curve, and panels (e), (f) and (g) show the corresponding energy schematic for each configuration. The case shown is for a flux rope radius of 0.1 in normalized units. [Fig. H-II:6.9; source: Forbes and Priest (1995).]

This boundary condition corresponds to two sources of opposite polarity located at $x = \pm\lambda$.

[...] Depending on the choice of model parameters, there may be three equilibria, one equilibrium, or no equilibrium for a given set of parameters. In situations with three equilibria the magnetic energy of each equilibrium is different. For the isolated equilibrium shown in Fig. 6.8b the flux rope sits in an energy well as shown in Fig. 6.8e. If the flux rope is pushed downward toward the surface, compression of the magnetic field between the flux rope and the surface creates an upward force. If the flux rope is pulled upward away from the surface, magnetic tension from the overlying arcade creates a downward force. Line-tying plays a key role in creating the equilibrium because it prevents field lines from being pushed into, or pulled out of, the surface when the flux rope is perturbed.

An evolutionary sequence is created by assuming that the distance between the two sources at $\pm\lambda$ decreases at a rate that is much slower than the Alfvén

time-scale in the corona. A flux rope located on the lower portion of the equilibrium curves shown in Fig. 6.8a will erupt when the distance between the line sources becomes less than the height of the flux rope. When this location is reached, the unstable and stable equilibria coincide as shown in Fig. 6.8g. Once equilibrium is lost, the flux rope rapidly moves upwards. In the absence of reconnection ($p = 0$) the flux rope does not escape, but, instead, reaches a new equilibrium position with a vertical current sheet, as shown in Fig. 6.8d.

In the absence of any reconnection the amount of energy released by the loss of equilibrium is quite small, less than 5% [...] Thus, while the loss of equilibrium can account for the rapid onset of an eruption, it cannot, by itself, account for the large amount of energy released. For this, magnetic reconnection is needed. [...] For typical coronal conditions a very modest rate of reconnection is sufficient to allow escape. For reconnection rates corresponding to an inflow Alfvén Mach number, M_A , > 0.05 (at the midpoint of the current sheet sides) the flux rope can escape without any deceleration [...]"

MHD simulations are needed to analyze such an eruption with more realism, and also to understand the role of waves, including shocks that develop when the eruption speed exceeds the propagation speeds of any of the possible MHD waves. More on this in the Heliophysics books.

6.3.3 Three-dimensional force-free models

[H-II:6.2.6] “It will probably come as no surprise that three-dimensional models are considerably more complex than two-dimensional ones. Three-dimensional field configurations are subject to a much greater number of instabilities. The helical ideal-MHD kink mode is an example of an inherently three-dimensional instability that does not exist in two dimensions. The dynamical evolution that occurs in three-dimensions is also more complicated. Fully nonlinear three-dimensional MHD turbulence can occur and magnetic reconnection exhibits new features that have no counterpart in two dimensions. Nevertheless, despite these additional complications, the underlying principles of the three-dimensional storage models remain the same.

[...] In order to show the relation of the relatively simple two-dimensional model of the previous section with these three-dimensional models, we take a reductionist approach. That is, we start with a very simple three-dimensional configuration and then sequentially add new features that increase its complexity. We start with the simple toroidal flux rope shown in Fig. 6.9. The antiparallel orientation of the current flowing on the opposite sides of the ring produces a repulsive force similar to the force between two parallel wires with antiparallel currents. For a small minor radius, a , this force, sometimes referred

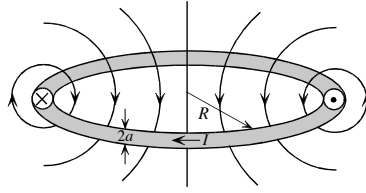


Fig. 6.9. An isolated toroidal flux rope. The flux rope has a major radius, R , a minor radius, a , and carries a net toroidal current I . The antiparallel orientation of the current flowing on the opposite sides of the torus creates an outward force in the radial direction. [Fig. H-II:6.16]

to as the hoop force, is approximately

$$F \propto \frac{I^2}{R} \ln(R/a), \quad (6.6)$$

where I is the flux-rope current, R is the major radius, and a is the minor radius of torus. The right-hand side of the above expression is the lowest order term of an expansion in the parameter a/R , so the expression is only valid for $a \ll R$ [...]

Just as for two-dimensional storage models, the three-dimensional models assume that the time scale of the eruption is so fast that any additional input of magnetic energy after the eruption starts, is completely negligible. Therefore, the flux associated with the flux rope current is conserved. In the limit that a/R tends to zero, the flux-rope current is roughly

$$I \approx \frac{I_0 R_0}{R \ln(R/a)}, \quad (6.7)$$

where I_0 and R_0 are initial values. If one considers the torus configuration as an initial state that subsequently evolves in response to the force, then R will increase to infinity, but as it does, so I will decrease to zero. In the process the magnetic energy associated with the flux rope's initial current is converted into the kinetic energy of the expanding plasma ring.

To create an equilibrium one must add an additional magnetic field of the proper orientation and strength. In tokamak terminology such a field is called a strapping field. [... Whereas it is possible to create a stable equilibrium by an appropriately shaped] strapping field, an alternative possibility that is more appropriate for a storage model is to introduce a line-tying surface as shown in Fig. 6.10. The effect of line-tying can be modeled by introducing a fictitious image current below the surface. With the introduction of this additional current, a new equilibrium appears which, unlike the previous one,

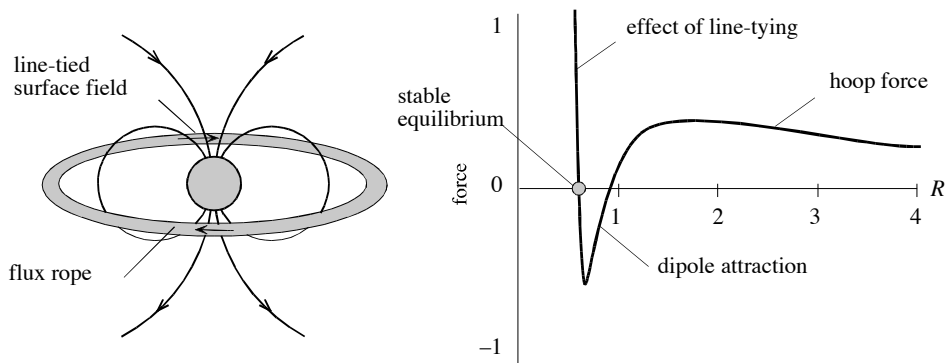


Fig. 6.10. A stable toroidal equilibrium. (left) The addition of a line-tying surface representing the surface of the Sun creates the possibility of a stable equilibrium. Surface currents (which can be modeled using an image current) create an additional magnetic field component that gives rise to a second equilibrium position as shown on the right. The new equilibrium is stable because displacements away from it produce a restoring force. [Fig. H-II:6.18]

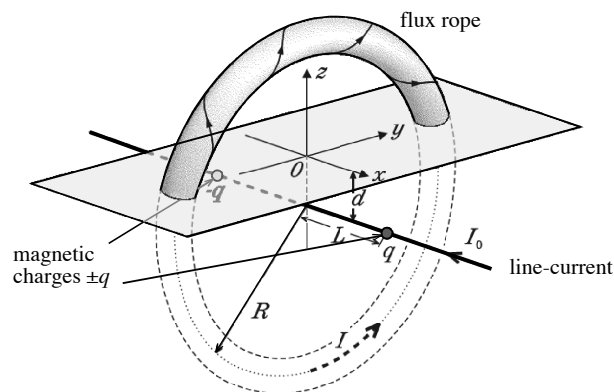


Fig. 6.11. The three-dimensional flux-rope model of Titov and Démoulin. The coronal magnetic field is produced by three different sources consisting of a flux rope current, a pair of magnetic charges, and a line current. The source regions located below the surface are fictitious constructs used to create the coronal field. The model does not prescribe the form of the subsurface field. [Fig. H-II:6.19; source: Titov and Démoulin (1999).]

is stable. Stabilization is achieved because line-tying prevents field lines from being pushed into, or pulled out of, the surface [...]

Although we now have an eruptive model with some degree of three-dimensionality, it still has the drawback that the flux rope is not itself anchored

to the solar surface. An analytical configuration that does have this property is shown in Fig. 6.11; [...] it consists of a toroidal flux rope that intersects the photospheric surface. The flux rope, with current I , is held in equilibrium by an overlying arcade (not shown in the figure) which is produced by subsurface magnetic charges $\pm q$ located along the centerline at a depth d below the photospheric surface at $z = 0$. Finally, there is a subsurface line current lying along the centerline. The strength of the current, I_0 , flowing in this subsurface line controls the pitch of the coronal magnetic field. When I_0 is varied from small to large values, the configuration changes gradually from a highly twisted flux rope resembling a slinky to one that resembles a sheared arcade without a flux rope.

Although the magnetic field of [what is known as a Titov-Démoulin] configuration is still azimuthally symmetric about the centerline of the torus, the solar surface no longer shares this symmetry. Instead the surface is a flat plane that intersects the flux rope torus at some arbitrary position without influencing the field structure. Thus, any line-tied evolution of this configuration away from the initial state necessarily creates a highly asymmetrical configuration. An example of what such a configuration looks like is shown in Fig. 6.12. This figure shows two different views of an iso-current surface of the current density obtained from a simulation. This simulation starts with an unstable Titov and Démoulin configuration that is given a small perturbation. Within a few Alfvén scale times the configuration evolves into the kinked, omega-shaped flux rope shown in the figure. For this particular case, the initial instability is actually a helical kink instability rather than the torus instability discussed previously. However, it is possible to construct unstable Titov and Démoulin configurations that are unstable to the torus instability rather than the helical kink [...]"

6.3.4 Formation of the pre-eruption field

[H-II:6.2.7] “An important issue that the above flux rope models do not address is the creation and growth of the magnetic stress that causes the field to erupt. It could be that most of the stress build-up occurs in the convection zone before the field emerges into the corona. Alternatively, it may be that the field emerges in a nearly unstressed, current-free state, and that the stress subsequently develops in response to the observed surface flows. In practice both possibilities are likely to occur at least at some level.

[Among the three-dimensional simulations that address this issue is one] called the *breakout model*. The evolution of this model is shown in Fig. 6.13. The initial state consists of a quadrupolar magnetic field that carries no current, so it contains no free-magnetic energy. Slowly shearing the central arcade

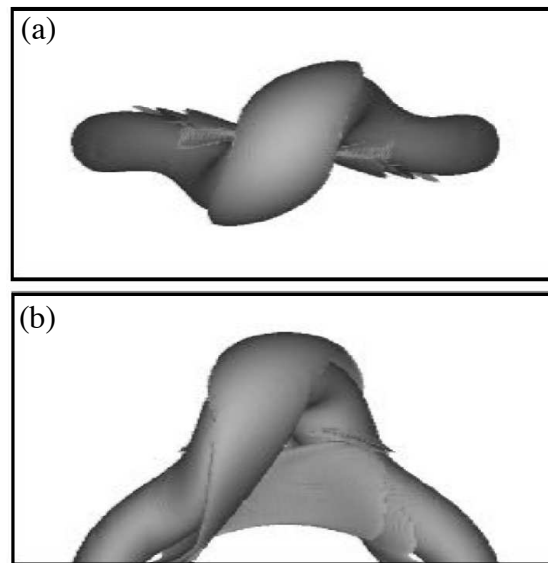


Fig. 6.12. Top view (a) and side view (b) of constant current density surfaces from a simulation for an unstable Titov and Démoulin equilibrium. [Fig. H-II:6.20; source: Török *et al.* (2004).]

around the equator gives rise to a set of stressed loops that push outward against the overlying arcade. As this happens, a curved, horizontal current sheet forms at high altitude at the pre-existing x-line. Eventually, the stresses build up to a level that causes an eruption. The nature of the mechanism that triggers the eruption has not yet been fully resolved, but it is likely that it consists of some kind of combination of both ideal and non-ideal processes.

[An alternative to this model is where a flux rope emerges into a pre-stressed field from below the photosphere.] Generally, the flux rope will tend to erupt once there are one or two turns in the portion of it that has emerged into the corona. However, if the flux rope emerges into a pre-existing arcade, the strength and orientation of this arcade also has a strong effect on whether an eruption occurs or not.

[...] One of the important issues that [various] studies address is the effect of mass loading on the emergence of a flux-rope into the low-density corona. Most of the CME models discussed in the previous section are based on the supposition that a flux rope exists in the corona prior to onset, but it is not obvious how such a structure could be formed. Formation of the flux rope within the convection zone followed by its buoyant rise into the corona immediately encounters the problem that mass cannot easily drain out of concave-upward portions of the magnetic field. Unless there is a way for the

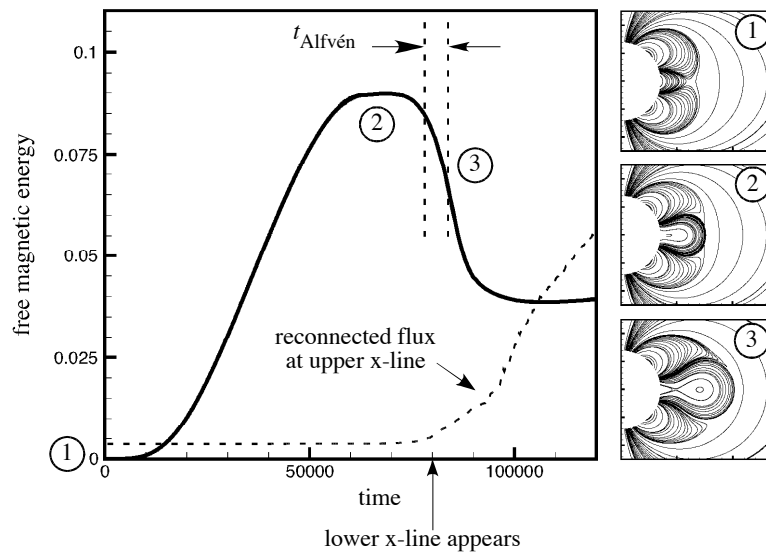


Fig. 6.13. Numerical simulation of a storage model. The panel at left shows the free magnetic energy as a function time, while the three panels at right show contours of the magnetic flux surfaces at three different times. [Fig. H-II:6.21; source for inset panels: MacNeice *et al.* (2004).]

mass to drain out of the flux rope, the rope will remain half buried in the solar surface. One way around this difficulty is to suppose that the flux rope does not exist prior to the emergence of magnetic flux, but instead forms in the corona by a combination of converging flows and slow reconnection. Most dynamo models, however, predict that large-scale flux ropes will form near the base of the convection zone and then rise buoyantly to the solar surface to form an active region. Thus, this solution to the mass-loading problem involves both the destruction and reformation of the flux rope below and above the surface.”

6.3.5 Observed signatures of flares and CMEs

[H-II:5.2.2] “The release of energy can either be ‘impulsive’, with time scales sometimes faster than 1 s, or ‘gradual.’ The impulsive and gradual signatures of a flare extend across the entire electromagnetic spectrum in a complicated way, as illustrated in Fig. 6.2. The terminology may not seem appropriate when one considers a slowly developing flare-like event, such as a quiet-Sun filament eruption; in such a case the ‘impulsive phase’ may take tens of minutes to evolve, and the hard X-ray emission may be below the detection level. Thus we don’t know how ‘impulsive’ the energy release really is in such an event, but in other respects it has the morphology of an ordinary active-region flare.”

[H-IV:2.1.2] “An individual flare can be divided into two main phases: impulsive and gradual. This generally refers to the timing of emissions relative to the processes thought to be occurring in the flare. In the standard picture, the initial energy conversion caused by magnetic reconnection powers particle acceleration and possibly – depending on the energetics and on the magnetic configurations – a mass ejection. The downward-directed particles become trapped in loops and emit non-thermal incoherent radio emission ([compare Fig. 6.3]). Coherently emitting particles can be traveling either upwards out of the atmosphere or downwards into the atmosphere. Once the trapped particles precipitate from the magnetic trap, they deposit their energy in or just above the photosphere, producing thick-target non-thermal bremsstrahlung emission. This energy deposition results in the heating of the photospheric material to temperatures near 10^4 K, and emissions from FUV lines. All of this is associated with the impulsive phase of the flare. The flow of energy at this point proceeds back into the upper atmosphere, with line emission from the lower chromosphere. Thermal X-ray emission occurs as well. As the energy input into the system decreases, emissions of all flare components return to the pre-flare level.”

[H-II:5.2.2] “We understand the impulsive and gradual phases to show the main energy [conversion out of the magnetic reservoir] and its aftermath (secondary effects), with the proviso that it is really not just that simple. The most prominent ‘aftermath’ is the action of coronal magnetic loops as an energy reservoir, with cooling time scales that can approach hours. This reservoir function is often described as the ‘Neupert effect’: the coronal manifestations of a flare tend to lag behind its chromospheric ones. This results from the finite time scale associated with the coronal density increase during the impulsive phase, via the process of ‘chromospheric evaporation.’ {A:[82]} The decay time scale reflects its slower cooling and return to the lower atmosphere. The new material in the corona could be seen in the coronal emission lines, via free-free emission at radio wavelengths, or via free-free emission at soft X-ray wavelengths [...]” {A:[83]} A:82
A:83

[H-IV:2.1.2.1] “The Neupert Effect relationship [...] was formulated originally to describe the integral relationship between markers in a solar flare

⁸² Activity: ‘Chromospheric evaporation’ is a misnomer because there is no phase transition involved: the heating of chromospheric material from $\approx 10^4$ K to of order $\approx 5 \times 10^6$ K causes the pressure and the associated pressure scale height to increase. What are the pre-heating and post-expansion scale heights for the above temperatures? How do these compare to the solar radius?

⁸³ Activity: For a given temperature, coronal soft X-ray brightness scales essentially with the square of the particle density. Why? Let a given coronal loop have an initial loop-top density n_0 at temperature T_0 and let an impulsive heating event change these to n_1 and T_1 . With $T_{0,1}$ within the range of about 0.4–30 MK the radiative losses scale as $P(T) \propto T^{-2/3}$. If the temperature changes from 1 MK to 5 MK and the density increases by a factor of 15, show that the ratio of radiative cooling time scales is close to unity. Conductive losses into the lower atmosphere, however, are larger at higher temperatures; why?

corresponding to the action of non-thermal particles, and the response from the atmosphere to the deposition of energy from these particles as it appears in coronal radiation. Written more generally,

$$L_{\text{gradual}}(t) \sim C_{\lambda\lambda'} \int_{t_0}^t L_{\text{impulsive}}(t') dt', \quad (6.8)$$

where $L_{\text{impulsive}}(t)$ is the time variation of an impulsive phase process which diagnoses the presence and action of particles accelerated in the explosive event (for stellar studies usually radio gyrosynchrotron, transition region FUV emission lines, or photospheric UV-optical continuum emissions), and $L_{\text{gradual}}(t)$ is the intensity corresponding to the gradual phase (usually coronal emission, but some chromospheric emission lines display the Neupert effect as well). The interpretation is that the gradual phase emission is responding to the buildup of energy that occurs as a result of the energy deposition being diagnosed by the impulsive phase emission. [... Note that] not all solar flares follow the standard flare scenario”: it appears to hold for some 80% of large flares, but overall for about half of all flares. The value of $C_{\lambda\lambda'}$ depends on the wavelength bands [λ and λ'] used for both the impulsive and gradual phases.”

[H-II:5.2.2] “The different atmospheric layers have a high degree of interconnectedness. Because a flare marks a transition between one quasi-stable configuration and another, the ordinary law of hydrostatic equilibrium dictates the run of pressure up through the atmosphere. A flare increases the gas pressure in the corona, at the expense of magnetic energy, and this can readily be detected at all levels). The hydrostatic scale height for pressure is given by $2k_{\text{B}}T_e/mg_{\odot}$, where k_{B} is the Boltzmann constant, T_e the temperature, m the mean molecular weight, and g_{\odot} the surface gravitational acceleration. For a flare temperature of 10^7 K, this scale height is a large fraction of the solar radius, much larger than the flare loop structures. Thus the vertical structure is isobaric in the upper chromospheric and coronal regions, and the chromosphere acts as a reservoir of mass to maintain this isobaric state as the flare loops cool, [lose pressure, and drain into the chromosphere] quasi-statically.”

[H-II:5.3] “In the photospheric spectrum we see solar flares as brief flashes of white light and UV continuum. At present these sources are often not resolved either in space (Mm scales) or time (few sec scales). The bright emission regions are embedded in the ‘ribbon’ regions that become more prominent in the chromospheric and EUV coronal lines. In the coronal emissions one sees bright coronal loops developing slowly, with those from the highest temperatures appearing first and then cooling down through generally longer wavelengths, while at the same time shrinking in length. [...]

Solar flares are not luminous on the scale of the total solar irradiance (‘solar

constant'), although they may produce a localized brightening seen against the bright photosphere. The powerful flare of November 4, 2003 was the first that could actually be detected in the total solar irradiance, by the radiometer on board the SORCE spacecraft. The signal, at roughly 5σ significance, amounted to about 300 ppm of the total signal, or 0.3 millimagnitudes in astronomical terms. There is a solar background noise level for such a measurement due to convection and oscillations; this amounts to some 50-100 ppm spread out over a bandwidth of a few mHz.

The localized brightening of a flare is much easier to see, of course, via an image even in white light. Carrington [was the first to see a solar flare. He described] his 1859 discovery as resembling the brilliance of Vega (α Lyrae), for example. [The photospheric brightening is a major fraction of a flare's energy budget.] Soft X-ray emission, for example, contains only 5-10% as much luminosity. This gradual component [...] results from a thermal distribution (hot gas) for which the X-ray emission itself is a dominant cooling term. The non-thermal tail of the X-ray spectrum ($h\nu > 10$ keV), on the other hand, is due to bremsstrahlung from stopping particles. The bremsstrahlung mechanism is very inefficient, providing a fraction of order 10^{-5} of the energy losses. The rest of the energy winds up in longer-wavelength radiation, notably the visible/UV continuum.

We must also consider the bulk kinetic energy [involved in major solar impulsive events: CME kinetic energies can rival [total photon losses] in such cases. In rare cases a CME can occur in the absence of a major perturbation of the lower atmosphere. [...] The partition of energy in a flare/CME event remains unclear physically and hard to determine observationally.

The impulsive phase of a flare marks the period of intense energy release and strong non-thermal effects, including the launching of the CME. The traditional observational tools for the impulsive phase are hard X-ray emission and gyrosynchrotron emission at cm to mm radio wavelengths. The hard X-rays normally show two dominant footpoints embedded in ribbon regions of opposite magnetic polarity, but we do not presently understand why there are normally just two. The sources are compact and rapidly variable, and we associate them with the UV and white-light continuum emissions that also come from the footpoint regions. Other wavelengths show impulsive emission components as well as gradual ones. A clear impulsive-phase signature also appears even in the total irradiance, but rarely exceeds the background variability [...]

The hard X-ray spectrum above about 10 keV plays a central role in our understanding of the impulsive phase because the collisional energy losses of the bremsstrahlung-emitting electrons rival the total flare energy itself. This relationship can be established directly by inverting the hard X-ray spectrum,

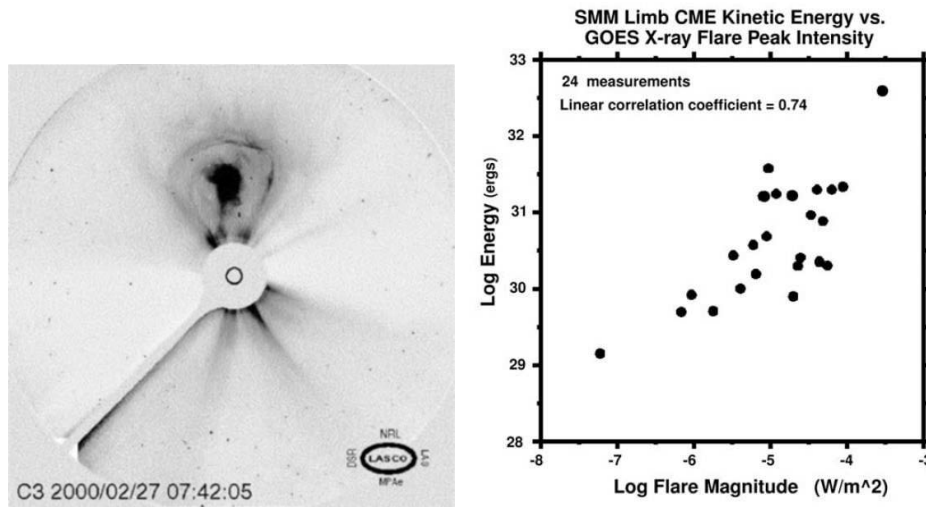


Fig. 6.14. Left: Coronagraph observation of a CME that nicely shows the three-part structure: front, cavity, and (the bright core) filament (this is a file image taken from the LASCO database, presented in a reverse greyscale). Right: Correlation between inferred CME kinetic energy and peak GOES soft X-ray flux. [Fig. H-II:5.5; source: Burkepile et al. (2004).]

under model assumptions. The 'collisional thick target model' envisions a black-box accelerator of 10-100 keV electrons in the corona, with a directed beam penetrating to the chromosphere or even photosphere to excite UV and visible-light emission. This simple model has become less tenable as spatial resolution improves, since the WL/UV brightenings [observed with newer instruments] imply beams with extreme intensity [...]

'Gradual phase' refers to the thermal emission from the hot coronal material evaporated during the impulsive phase, plus the strong transition-region and chromospheric emissions driven by the cooling of these coronal loops. The loops connecting the roughly parallel ribbons form a semi-cylindrical *arcade* structure, divided into many unresolved loops. [...] The hot regions eventually cool to form the H α loop prominence system, whence thermal instability leads to the phenomenon of 'coronal rain'. The cooling also corresponds to shrinkage, as the gas pressure diminishes; shrinkage may also relate to the gradual release of energy as the coronal equilibrium returns to a stable configuration. This is the process termed 'dipolarization' in the geomagnetic community [...]"

[H-II:5.6.4] "The expanding motions of flare ribbons provided one of the first clues to what we think of as the standard reconnection model of a flare [T]hese motions can be interpreted as an electric field. This is a motional or 'convective' electric field given by $\mathbf{E} = -\mathbf{v} \times \mathbf{B}/c$, and it is often taken as a

measure of the reconnection rate. [T]he rate the ribbons sweep out the field should correspond in some sense to the rate at which energy is released during reconnection, and that at the same time the field guides the particle or heat flux responsible for the ribbon excitation.”

[H-II:10.2] “[T]he resemblance of terrestrial substorms to a two-ribbon solar flare, with ribbons of opposite magnetic polarity, has been repeatedly remarked upon” (compare Figs. 6.5 and 6.6). Note, however, that the process leading up to the event is entirely different: the unstable field configuration is built up by stressing and/or flux injection from below in the case of solar eruptions but driven via the wind-magnetosphere interaction from the outside in the case of terrestrial substorms.

[H-II:5.3] “Major flare events almost invariably involve the ‘opening’ of the magnetic field as a CME; see Table 6.1 for the statistics [...] Observationally, [...] we often see a characteristic three-part structure: front, cavity, and filament (Fig. 6.14). This pattern makes it clear that the CME originated in a filament cavity near the surface of the Sun. A filament cavity consists of long, basically horizontal field, presumably more intense than its overlying ‘tie-down’ field that is more potential [...]

Modern images in coronal emissions such as soft X-rays allow a comparison of the coronal state before and after a CME event. Such comparisons revealed ‘dimmings,’ readily interpreted as the evacuation of the mass of the corona by the CME eruption. The soft X-ray dimmings presumably correspond to the coronal depletions found via similar before/after comparisons of the visible corona.” {A:[84]} {A:[85]} {A:[86]} {A:[87]} {A:[88]} A:84

⁸⁴ Activity: Describe what is seen in Fig. 6.14: how can a CME be imaged, and why is that best done from space, or from a very high mountain top? Argue why the CME in this image is not likely to envelop Earth. What would an Earth-bound CME look like? Can you differentiate that from one moving in the opposite direction? A:85 A:86 A:87

⁸⁵ Activity: The energy processed during a strong to intense solar flare and CME is of order $E_{\text{flare}} = 10^{30} - 10^{33}$ erg. Compare that to an order of magnitude estimate of the energy $E_{\text{AR},\odot}$ contained in the field of an active-region core (by, say, using a scale of 30,000 km and a characteristic average magnetic flux density of 300 G). What is the value of $E_{\text{flare}}/E_{\text{AR},\odot}$? How does this compare to $E_{\text{storm}}/E_{\text{mag},\oplus}$ and $E_{\text{storm}}/E_{\text{sw}}$ in Activity 81? A:88

⁸⁶ Activity: One phenomenon associated with many CMEs is a so-called ‘coronal dimming’, in which a large fraction of the quiet-Sun solar corona fades for some time. Think about the potential causes: temperature change (so the signal moves from one bandpass to another), quasi-adiabatic expansion of the coronal field, and entrainment of coronal plasma in the erupting CME. Estimate the volume of quiet-Sun corona (at a density of some 10^7 cm^{-3}) that would need to be involved if it were to move out with an erupting field configuration if that made up, say, 25% of the erupting mass of, for example, 10^{15} g.

⁸⁷ Activity: For a sense of scale: how many nuclear bombs are needed to match the energy released in a large solar flare of 10^{32} erg?

⁸⁸ Activity: Advances in numerical capabilities are making a big difference in understanding magnetic instabilities, how and where associated plasma heating occurs, and how combinations of plasma flows and a variety of temperatures in plasmas along a line of sight through the optically thin corona lead to observables. Such work shows how apparently non-thermal signatures in spectra can emerge from line-of-sight integration through thermal plasmas. If you would like to learn more about how observables based on numerical work help guide the interpretation of real-world observables, a paper (with illuminating graphics) by Cheung *et al.* (2019) provides a good example.

6.4 Magnetic instabilities and reconnection

One of the mechanisms thought to be involved in the destabilization of magnetic configurations is reconnection. Fast reconnection is often accompanied by shocks, and both the motions in the reconnecting field and the shocks themselves contribute to energy conversion into a mixture of thermal and non-thermal populations. There is a vast literature on the topic, including how such processes contribute to instabilities in the solar corona and the magnetosphere. Here we touch only on the fundamentals, specifically the concepts involved in steady 2-dimensional reconnection; more comprehensive material (and references to further reading) moving towards the time-dependent and 3-dimensional real world is provided in H-I:5 and H-II:6.

Let us start with a highly simplified configuration, generally referred to as 'Sweet-Parker reconnection' of steady 2-dimensional reconnection in an incompressible plasma in a current sheet with [system-level] length scale L_e [H-I:5.3.1] "as shown in [the left panel of] Fig. 6.15. Under these conditions [...] the speed of the plasma flowing into the current sheet is [approximately

$$v_e = \left(\frac{v_{Ae}\eta}{L_e} \right)^{1/2} \quad (6.9)$$

where $v_{Ae} = B_e/\sqrt{4\pi\rho_e}$ is the Alfvén speed in the inflow region. The outflow speed of the plasma from the current sheet is the local Alfvén speed V_{Ae} .] The reconnection rate in two dimensions is measured by the electric field at the reconnection site. This electric field is perpendicular to the plane of Fig. 6.15, and it prescribes the rate at which magnetic flux is transported from one topological domain to another. In two-dimensional steady-state models this electric field is uniform in space. Therefore, the Alfvén Mach number, $M_{Ae} = v_e/v_{Ae}$, provides a quantitative measure of the reconnection rate, normalized by the characteristic electric field $v_{Ae}B_e$. [...]

In astrophysical and space plasmas [...] Sweet-Parker reconnection is usually too slow to account for phenomena such as geomagnetic substorms or solar flares. [A later model, known as 'Petschek reconnection', was developed to ensure much faster reconnection by encasing the] current sheet in an exterior field with global scale length L_e , [and by introducing] two pairs of standing slow-mode shocks radiating outwards from the tip of the current sheet as shown in [the righthand panel of] Fig. 6.15. In Petschek's solution most of the energy conversion comes from these shocks which accelerate and heat the plasma to form two hot outflow jets.

Petschek also assumed that the magnetic field in the inflow region was current free and that there were no sources of field at large distances. These assumptions, together with the trapezoidal shape of the inflow region created

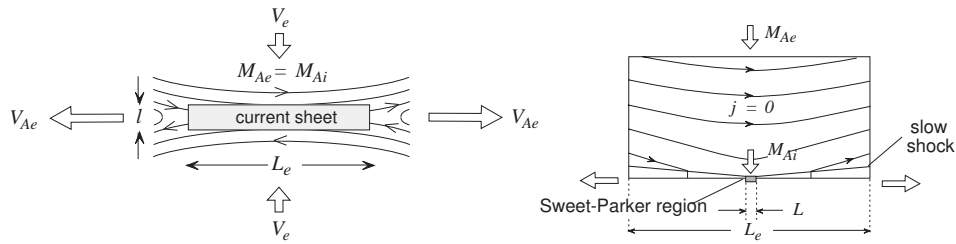


Fig. 6.15. [(left)] The Sweet-Parker field configuration. Plasma flows into the upper and lower sides of a current sheet of length L_e , but must exit through the narrow tips of the sheet of width l . Because the field is assumed to be uniform in the inflow region, the external Alfvén Mach number, $M_{Ae} = v_e/v_{Ae}$, at large distance is the same as the internal Alfvén Mach number, M_{Ai} , at the midpoint edge of the current sheet. [Fig. H-I:5.2] [(right)] Petschek's field configuration. Here the length, L , of the Sweet-Parker current sheet is much shorter than the global scale length, L_e , and the magnetic field in the inflow is nonuniform. Two pairs of standing slow-mode shocks extend outwards from the central current sheet. Petschek's model assumes that the current density in the inflow region is zero and that there are no external sources of field at large distance. [Fig. H-I:5.3]

by the slow shocks, lead to a logarithmic decrease of the magnetic field as the inflowing plasma approaches the Sweet-Parker current sheet. This variation of the field leads in turn to Petschek's formula for the maximum reconnection rate, namely

$$M_{Ae[\text{Max}]} = \pi / (8 \ln(L_e v_{Ae} / \eta)) \quad (6.10)$$

where [...] M_{Ae} is the Alfvén Mach number in the region far upstream of the current sheet as shown in Fig. 6.15. [...] The] Petschek reconnection rate is many orders of magnitude greater than the Sweet-Parker rate, and for most space and laboratory applications Petschek's formula predicts that $M_{Ae} \approx 10^{-1}$ to 10^{-2} . [...]

It is not always appreciated that Petschek's reconnection model is a particular solution of the MHD equations which applies only when [...] the flows into the reconnection region be set up spontaneously without external forcing [and] that there be no external source of field in the inflow region. In other words, the field must be just the field produced by the currents in the diffusion region and the slow shocks. In many applications of interest neither of these conditions is met." [H-I:5.3.1] "Even in circumstances where Petschek's model would be expected to apply it apparently does not. [Numerical simulations suggest that it only does in case of a nonuniform, localized resistivity. This] does not contradict Petschek's model because the model makes no explicit assumption about whether the resistivity is uniform or not. It is equally valid for both cases because it assumes only that the region where resistivity is important

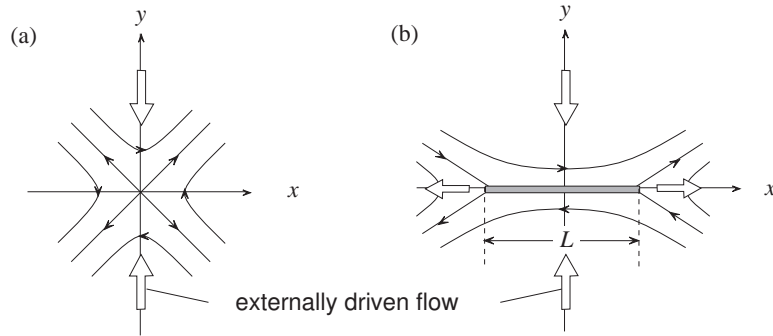


Fig. 6.16. Syrovatskii's field configuration. Unlike Petschek's configuration even when local sources of current are absent. The application of external driving creates a current sheet whose length, L , depends on the temporal history of the driving and the rate at which reconnection operates. The fastest reconnection rate occurs when L is equal to the external scale length, L_e . [Fig. H-I:5.4]

is localized. The numerical experiments carried out to date imply that the diffusion region can only be localized by enhancing the resistivity near the \times -line. Whether there might be other ways to localize the diffusion region (e.g., a non-uniform viscosity) remains unknown.”

[H-I:5.3.1] “An alternative approach to reconnection in current sheets was [developed by considering] what happens when a weak flow impinges on an \times -line in a strongly magnetized plasma as indicated in Fig. 6.16. The imposed flow creates a current sheet which achieves a steady-state when the rate of field line diffusion through the sheet matches the speed of the flow. [...]

For a steady-state MHD model the spatial variation of the field in the inflow region is the key quantity which [sets the] reconnection rate. [...] For any such model, the electric field is uniform and perpendicular to the plane of the field. Thus, outside the diffusion region $E_o = -v_y B_x / c$ where E_o is a constant, v_y is the inflow along the axis of symmetry (y axis in Fig. 6.16), and B_x is the corresponding field. Thus the inflow Alfvén Mach number, M_{Ae} , at large distances can be expressed as

$$M_{Ae} = M_{Ai} \left(\frac{B_i}{B_e} \right)^2 \quad (6.11)$$

where M_{Ai} is the Alfvén Mach number at the current sheet, B_i is the magnetic field at the edge of the current sheet, and B_e is the magnetic field at large distance.

In Syrovatskii's model the field along the inflow axis of symmetry is

$$B_x = B_i (1 + y^2 / L^2)^{1/2} \quad (6.12)$$

where B_i is the field at the current sheet, y is the coordinate along the inflow axis, and L is the length of the current sheet. Combining (6.12) with (6.11) yields

$$M_{Ae} = M_{Ai}/(1 + L_e^2/L^2) \quad (6.13)$$

which has its maximum value when $L = L_e$. Thus the maximum reconnection rate in Syrovatskii's model scales [the same as] the Sweet-Parker model.

By comparison, the field in Petschek's model along this axis varies approximately as

$$B_x = B_i \frac{1 - (4/\pi)M_{Ae} \ln(L_e/y)}{1 - (4/\pi)M_{Ae} \ln(L_e/l)} \quad (6.14)$$

where l is the current sheet thickness. (This expression for the field is only a rough estimate since the actual variation in the region $y < L$ is more complex.) Evaluating this at $y = L_e$ and substituting the result into Eq. (6.11) gives

$$M_{Ai} = M_{Ae}/[1 - (4/\pi)M_{Ae} \ln(L_e/l)]^2 \quad (6.15)$$

The Sweet-Parker theory can be used to eliminate L_e/l , so as to obtain an expression for M_{Ae} [which] has a maximum value as given by equation (6.10). [...]"

Even a much longer summary than the above could not be conclusive: [H-I:5.5] "There are many aspects of magnetic reconnection that have yet to be explored. Even well long studied topics such as steady-state two-dimensional reconnection are not fully understood. Many questions remain about how time-dependent reconnection works in impulsively driven phenomena such as solar flares and geomagnetic substorms. For example, during the impulsive phase of eruptive solar flares the current sheet where reconnection occurs can grow at a rate that exceeds the Alfvén time scale of the system. This rapid growth means that no steady-state reconnection theory applies during the impulsive phase, and there are virtually no theories that predict how the reconnection rate scales with plasma resistivity in such a situation. [Another large challenge to our understanding of reconnection is that, despite] growing evidence that the reconnection process in both solar flares and the terrestrial magnetosphere may be turbulent, there are only few studies that address the issue of turbulent reconnection. The occurrence of plasma turbulence in a highly-structured environment poses a severe challenge to large scale numerical simulations, so progress in this area may be slow for sometime to come."

7

Torques and tides

7.1 Introduction

Rotation and revolution are key properties for stars, planets, and indeed the entirety of planetary systems. For example, rotation is one of the essential ingredients of dynamo action in stars and planets, while climates are determined to a large extent by the Coriolis forces associated with planetary spin in combination with the overall duration of insolation on a planet's dayside, by the orbital eccentricity, and by the planetary obliquity, *i.e.*, the tilt of the planetary spin axis relative to the orbital plane. But spin rates evolve: stars slow their rotation because of their magnetic activity, while planetary rotation can change subject to tidal coupling with their moons, with their own atmospheres, and with their central stars. The latter, tidal synchronization of planetary spin and orbital motion, is likely common among exoplanets in the habitable zones of intrinsically faint M-type dwarf star because these exoplanets would have to have very tight orbits.

Also revolution is subject to change: planetary orbits need to be stable over long times to offer long-term habitability, but – as we discuss in Sects. 7.3.2 and 11.1 – orbits need to evolve for juvenile planets to grow efficiently and also, for instance, to transport water through a planetary system across the ice line from the cold outer reaches to the habitable inner domain. ^[xiv] And if orbital dynamics (quantified in angular momentum) could not be efficiently transported through gas and dust, planetary systems and their central stars could not form as they are observed to do (see Sects 7.2.3, 7.2.4, and 11.2).

[H-III:1.3] “Transport of angular momentum through the coupling of distant concentrations of mass occurs either through gravitational tides, by magnetic stresses, or by flows. Gravitational coupling has obviously played an important part in the spin-orbit synchronization of the Earth's single moon. This coupling

^{xiv} The temperature at which H₂O freezes into a solid in a proto-planetary disk is dependent on the partial pressure of the water vapor in the overall gas mixture, and is typically expected to lie in the range of 145 K to 170 K.

continues to be important as a stabilizer for the direction of the Earth's spin axis, even as it causes the precession of that axis with associated climatic effects (Chs. H-III:11 and H-III:12). [This chapter reviews these processes, and more:] Tidal forces also act significantly on Jupiter's moon Europa and Saturn's moon Enceladus, in which it appears to result in liquid water in their interiors, which makes these moons interesting objects to study from an exo-biological perspective. Tidal spin-orbit coupling also leads to the formation of short-period, highly active binary stars (like the so-called RS CVn type systems). [...]

[This chapter also highlights a]ngular momentum transport via the magnetic field [which] is important in the coupling of proto-stars and young T Tauri stars to their surrounding disks and magnetized stellar winds. {A:^[89]}

After the early formation phases of a planetary system, the loss of stellar angular momentum continues through a stellar wind, leading to magnetic braking of the stellar rotation and the concomitant gradual decrease in stellar activity with age. In tidally interacting binaries with one or more magnetically active components, the loss of spin angular momentum by a stellar wind drains the orbital angular momentum reservoir, eventually leading to the merger of the component stars, leaving an old but rapidly-spinning single star (like FK Comae). [The consequences of these couplings are discussed in Chs. 10 and 11.] A:89

Angular momentum transport by flows inside astrophysical bodies is the cause of the near-rigid rotation with latitude and depth of the solar interior. But the models of the full convective envelopes of stars and giant planets need to advance significantly before we can use their results in, *e.g.*, magnetohydrodynamic dynamo models in which the non-rigid rotation and other large-scale circulations appear to be crucial."

Even photons are involved in a form of tidal action: [H-III:15.2] "Atmospheric tides are the response to periodic astronomical forcing. Atmospheric tides [on Earth] are forced primarily by the thermal heating due to the absorption of solar radiation by ozone and water vapor. These tides have periods which are the length of a mean solar day and its harmonics. [...]"

This chapter briefly introduces each of these processes and the settings in which they are important, but the consequences of evolving orbital motions and spin rates are left for later chapters.

⁸⁹ Activity: Look up what T Tauri stars are, and what differentiates the 'classical' T Tauri star from the 'weak-line' variant.

7.2 Magnetic torques

7.2.1 Stellar winds and magnetic braking

The solar wind discussed in Ch. 5 not only carries mass away from the Sun, but also angular momentum. [H-IV:4.2] “The conventional mechanism for stellar spin down is that stars lose angular momentum to the magnetized stellar wind in the concept called ‘magnetic braking’. In this process, the mass flux carried by the accelerating stellar wind conserves angular momentum as long as the wind speed is below the Alfvén speed, $v_A = B/\sqrt{4\pi\rho}$ (in cgs units of cm s^{-1}), where B is the local magnetic field strength, and ρ is the local mass density. Once the wind speed equals the Alfvén speed [(at the ‘Alfvén radius’; see Sect. 5.4) the wind is effectively decoupled] from the star. Another way to look at this process is to think of the magnetic field lines as rods [up to the Alfvén radius, beyond which the wind flows out essentially radially as if flung free from the star only at that radius (which in reality is a gradual process), dragging the magnetic field into a Parker spiral.] As a result, each field line applies a torque on the star and spins it down. This torque is proportional to the momentum of the wind at the Alfvén point, to the stellar rotation rate, and to the distance of the Alfvén point (the lever arm that applies the torque). The imaginary surface that contains all the Alfvén points is called the ‘Alfvén surface’ and the integral of the mass flux through this surface is the mass loss rate, \dot{M} , of the star to the stellar wind. For a spherically symmetric wind, and [a magnetic field that is close to uniformly distributed across the Alfvén surface,] we can calculate the total torque on the star and the total angular momentum loss rate, \dot{J} :

$$\dot{J} = -\Omega_* \dot{I}_{\text{shell}} = -\frac{2}{3} \Omega_* \dot{M} r_A^2, \quad (7.1)$$

where Ω_* is the stellar rotation rate, [\dot{I}_{shell} is the moment of inertia of a uniform shell of mass \dot{M} and radius r_A , and] r_A is the average distance to the Alfvén surface, and we assume constant moment of inertia [for the star (*i.e.*, we assume the time scale for angular momentum loss in this expression is short relative to the evolution of the internal structure of the star, which is appropriate for the long-lived ‘mature’ phase of the star, see Ch. 10). Note that Eq. (7.1) shows that the near co-rotation out to r_A^2 causes the solar wind to carry a factor of r_A^2/r_\odot^2 more angular momentum away from the star than is contained in the mass that is actually leaving the stellar surface.]

From Eq. (7.1) we see that the mass-loss rate is necessary to estimate the spin-down rate of a star. However, stellar winds of cool, Sun-like stars are very weak and cannot be directly observed (see Ch. 10), which makes it challenging to estimate \dot{J} as a necessary input for stellar evolution models [...] Based on

[measurements supported by modeling (described in Sect. 10.3.2) mass-loss rates in Sun-like stars] fall in the range between $10^{-15} - 10^{-11} M_{\odot} \text{ yr}^{-1}$ (the present-day solar mass-loss rate is $(2 - 3) \times 10^{-14} M_{\odot} \text{ yr}^{-1}$). However, stars can also lose mass via CMEs. In the case of the Sun, each CME carries some $10^{13} - 10^{17} \text{ g}$ into space, with an annual integrated mass-loss via CMEs of several percents of the ambient mass-loss. Therefore, CMEs on the Sun play very little role in the solar mass-loss. This role can become significant if the CME rate is higher by a factor of 10 or more. In this case, CMEs can even dominate the stellar mass-loss.” But we know very little of CMEs of stars other than the Sun, or even of the Sun in its distant past, so this area is left for future exploration.

Let us make a few comparisons of energy budgets and time scales, using rough approximations only: The above-mentioned solar mass-loss rate can be combined with numbers in Table 2.4 to estimate the power needed to drive the flow of the solar wind (bulk kinetic energy), the Alfvén radius, and with that the rate at which rotational energy is drained from the Sun. Assuming for this estimate a constant mean wind velocity, a temperature of 1.5 MK for the high corona, an isotropic heliospheric magnetic field strength (approximating the field as radial), a total characteristic power associated with all forms of coronal radiative losses driven by solar magnetic activity of order $\approx 10^5 \text{ erg/cm}^2/\text{s}$ (averaged over a solar cycle), and a moment of inertia of $I_{\odot} \approx 7 \times 10^{53} \text{ g cm}^2$ (see Ch. 10), one can conclude that (1) the solar wind kinetic energy flux amounts to of order 1/10th of the coronal radiative losses, (2) the characteristic Alfvén radius is roughly 20 solar radii, and (3) the time scale for magnetic braking, *i.e.*, the ratio of angular momentum to loss rate of angular momentum for the present-day Sun is of order 10 Gyr. {A:[90]}

A:90

In very rapidly spinning stars, the centrifugal forces that we have ignored for the solar wind, also need to be taken into account. An example where these dominate the process in the case of a cold wind is discussed in Sect. 7.2.4.

7.2.2 Planetary magnetospheric torque

We can make a similar comparison of energy budgets and time scales for the solar wind that delivers power to Earth’s magnetosphere and induces a torque on Earth’s rotation. First, let us look at the energy, then at the angular momentum. [H-II:10.4.1] “The net rate of energy extraction (power) \mathcal{P}_{sw} from the solar wind flow is equal to the difference of the solar wind kinetic energy flux across two surfaces A perpendicular to the Sun-planet line, surface 1 ahead

⁹⁰ Activity: Verify the numbers in the conclusions about stellar magnetic braking for the present-day Sun at the end of Sect. 7.2.1.

of the bow shock and surface 2 far downstream of the entire interaction,

$$\begin{aligned}\mathcal{P}_{\text{sw}} &= \frac{1}{2} \int_1 \rho v^3 dA - \frac{1}{2} \int_2 \rho v^3 dA \\ &= \frac{1}{2} \int \rho v (\bar{v}_1^2 - \bar{v}_2^2) dA \\ &= \dot{M}_{\text{ft}} \bar{v} \Delta v\end{aligned}\quad (7.2)$$

(subscripts 'sw' on ρ and v have been omitted, for simplicity), and the total force F is similarly equal to the difference of the linear momentum flux,

$$F = \int_1 \rho v^2 dA - \int_2 \rho v^2 dA = \dot{M}_{\text{ft}} \Delta v, \quad (7.3)$$

where $\Delta v \equiv \bar{v}_1 - \bar{v}_2$ and $\bar{v} \equiv (\bar{v}_1 + \bar{v}_2)/2$ (bars indicate suitable averages) and

$$\dot{M}_{\text{ft}} = \int_1 \rho v dA \simeq \int_2 \rho v dA \quad (7.4)$$

is the amount of mass per unit time flowing through the region of interaction between the solar wind and the magnetosphere, to be distinguished from \dot{M}_{sw} , the mass input rate from the solar wind into the magnetosphere. (Note: magnetic and thermal contributions to solar wind energy and momentum flux have been neglected as small in comparison to those of the bulk flow; see Sect. 3.5.2.) Combining Eqs. (7.2) and (7.3) yields a relation between the power and the force (in the direction of solar wind flow), $\{A:^{[91]}\}$

A:91

$$\mathcal{P}_{\text{sw}} = F \bar{v}, \quad (7.5)$$

which [has been used] to estimate the energy input into the terrestrial magnetosphere, under the assumption that the relevant force F is the tangential (magnetotail) force acting primarily on the nightside, F_{MT} (see Sect. H-I:10.3.2)."

[H-I:10.3.2] "While pressure from the external medium thus accounts for the formation and shape of the magnetosphere on the dayside of the planet, it cannot by itself explain the formation of the *magnetotail* on the night side. This structure, shown also in Figure 5.12, is a region of magnetic field lines pulled out into an elongated tail in the anti-sunward direction, with the magnetic field reversing direction between the two sides of a *current sheet* or *plasma sheet* in the equatorial region. To form this structure one needs an appropriate stress: a tension force pulling away from the planet. If we choose a closed volume bounded by a surface just outside the magnetopause plus a cross-section of the

⁹¹ Activity: For comparison: what is the approximate ratio of forces exerted on the Earth of the total solar irradiance onto the Earth's surface (ignoring albedo, and assuming isotropic radiation from the atmosphere) to the solar-wind pressure on the magnetopause? That ratio shows why solar sails are designed for photon pressure rather than solar-wind dynamic pressure (note that some are designed to couple to induced electromagnetic effects, not dynamic pressure).

magnetotail (vertical cut at the right edge of Figure 5.12) and evaluate the force [...], the total tension force F_{MT} is given by the integral over the cross-section and the total pressure force F_{MP} by the integral over the magnetopause:

$$F_{\text{MT}} \simeq (B_t^2/8\pi) A_t \quad F_{\text{MP}} \simeq \rho_{\text{sw}} v_{\text{sw}}^2 A_t \quad (7.6)$$

where B_t is the mean magnetic field strength and A_t the cross-sectional area of the magnetotail (typically, A_t exceeds πR_{MP}^2 by a factor 3 to 4). Both F_{MP} and F_{MT} are directed away from the Sun and are exerted ultimately on the planet." [H-II:10.4.1] "Note: if F [in Eq. (7.5)] is equated to the pressure force F_{MP} on the entire magnetopause, it can be shown that the associated \mathcal{P} does not go into the magnetosphere but represents the power expended in irreversible heating at the bow shock.

Calculating the power extracted from planetary rotation is somewhat simpler. The angular momentum of the rotating planet is $I_p \Omega_p$ and the kinetic energy of rotation is $\frac{1}{2} I_p \Omega_p^2$, where I_p is the moment of inertia and Ω_p the angular frequency of rotation [of the planet]. With \mathcal{T} the torque on the planet (component along the rotation axis),

$$\mathcal{P}_{\text{rot}} = \frac{d}{dt} \left(\frac{1}{2} I_p \Omega_p^2 \right) = \Omega_p \frac{d}{dt} (I_p \Omega_p) = \mathcal{T} \Omega_p, \quad (7.7)$$

a relation between the power and the torque, completely analogous to Eq. (7.5). [...]

What happens to the linear momentum extracted from the solar wind flow is well understood: it is transferred to and exerts an added force on the massive planet. The angular momentum extracted from the rotation of the planet, on the other hand, can only be removed to 'infinity,' and identifying the mechanism by which it is transported away is indispensable for understanding the interaction. There are several possibilities:

(a) In magnetospheres with a significant interior source \dot{M} of plasma (from moons or planetary rings), angular momentum can be advected by the outward transport of mass [as long as the planet's rotation period is below the orbital period of the plasma source]. For the simple example of plasma corotating rigidly out to a distance R_c and coasting freely beyond R_c , angular momentum is transported outward at the rate $\dot{M} R_c^2 \Omega_p$, hence from Eq. (7.7) the extracted power is

$$\mathcal{P}_{\text{rot}} \simeq \dot{M} \Omega_p^2 R_c^2, \quad (7.8)$$

one half of which goes into the kinetic energy of bulk flow of the outflowing plasma (in this model), and the remainder is available for powering other magnetospheric processes (proposed for the magnetosphere of Jupiter).

(b) If the solar wind exerts a tangential force on the magnetosphere, it will

also exert a torque whenever the distribution of the force is not symmetric about the plane containing the solar wind velocity and the planetary rotation axis. The torque may be estimated as $\mathcal{T} \sim R_{\text{MP}}\Delta F$, where R_{MP} is the distance to the dayside magnetopause and ΔF is the difference between the force on the dawn and on the dusk side; this gives the ratio of power from rotation to power from solar wind flow as

$$\mathcal{P}_{\text{rot}}/\mathcal{P}_{\text{sw}} \sim (\Delta F/F) (\Omega_{\text{p}}R_{\text{MP}}/v_{\text{sw}}) . \quad (7.9)$$

In a slowly rotating magnetosphere such as Earth, $\Omega_{\text{p}}R_{\text{MP}}/v_{\text{sw}} \equiv \epsilon \ll 1$ and one also expects $\Delta F/F$ to scale as $\sim \epsilon$; hence the power extracted from rotation by the solar-wind torque is negligible.” [H-II:10.4.1] “In principle, Ω_{p} decreases with time as the result of the torque, but in practice the rate of decrease is completely negligible. The time for the present magnetospheric torque to reduce appreciably the planet’s rate of rotation is several orders of magnitude longer than the [age of the Universe], both at Jupiter and at Earth; for the latter, this implies that the magnetospheric torque is much smaller than the lunar tidal torque.”

[H-II:10.4.1] “(c) In a rapidly rotating open magnetosphere, on the other hand, magnetic field lines that extend from the planet into the solar wind may become twisted (by a process analogous to the formation of the Parker spiral in the solar wind), creating a Maxwell stress that transports angular momentum outward into the solar wind. This mechanism of extracting energy from planetary rotation was proposed for Jupiter (where it is now considered not important in comparison to mass outflow) and for Uranus.

(d) If the magnetic moment of the planet is tilted relative to the rotation axis, electromagnetic waves that carry away angular momentum may be generated by the rotation. This is generally believed to be the primary mechanism for energy loss from pulsars but is negligible for systems that are very small in comparison to c/Ω , the radius of the speed-of-light cylinder (which is the case for all planets in our Solar System and their magnetospheres).”

7.2.3 Magneto-rotational coupling

As we saw in the case of the stellar wind, magnetic fields can support tension (Sect. 3.2.2) and thereby can essentially enforce co-rotation of gases at different distances from a star, at least out to where the field is strong compared to the inertial forces associated with the plasma. This not only holds for outflows such as stellar winds, but also in systems where matter is ‘descending’ onto the star, such as in very young proto-planetary systems where material has shaped itself into a disk spinning around an accreting star. More on that process in Ch. 11, but let us look at what a magnetic field that threads such a disk can do:

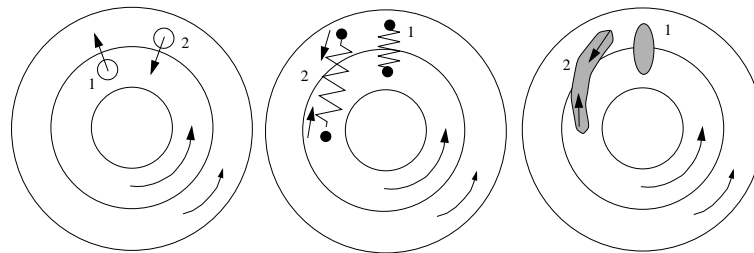


Fig. 7.1. Schematic treatment of angular momentum transfer in a shearing disk with angular velocity decreasing outwards [as the orbiting material approximates Keplerian orbits]. An initially radial field is perturbed azimuthally (left panel); these azimuthal perturbations grow due to the shear in the disk (middle panel [going from time label '1' to time label '2']). In the case of gravitational instability (right panel), an excess of material gets sheared out by the differential rotation; the gravitational attraction on the sheared excess (spiral arm pattern) exerts a restoring force in the same sense as the magnetic case, again transferring angular momentum outward. [Fig. H-III:3.4; source: Hartmann (2009).]

a field can be an effective agent in transporting angular momentum outwards, thereby enabling the gas in the disk to spiral inwards and thus help form the star.

[H-III:3.2] “As shown in the middle panel of Figure 7.1, if the magnetic field lines are thought of schematically as springs tying adjacent disk annuli together, then as differential rotation continually separates the regions ‘tied’ to the field (*e.g.*, evolution from [‘1’ to ‘2’]), the ‘springs’ or field lines become stretched [and bent], and the resultant [tension] forces will work in the direction of spinning up the outer annulus while spinning down the inner annulus.

The magnetic fields shown in the top-down view of the middle panel of Figure 7.1 cannot be stretched indefinitely; at some point there will be reconnection and diffusion as the flow becomes turbulent. [In that case, an] initially vertical field is perturbed radially; these radial perturbations grow due to the shear in the disk; and eventually the field lines become so stretched that they pinch off and develop into full turbulence.

Although there is currently some controversy over the efficiency of this ‘magneto-rotational instability’, or MRI, it seems very likely that it provides a sufficiently effective means of promoting accretion in astrophysical disks – provided, of course, that the magnetic field can couple effectively to the gas; there must be a sufficient population of ions and electrons to collide rapidly enough with neutral gas to make the MRI work. Protostellar disks are problematic in this regard: with much or most of their mass heavily shielded

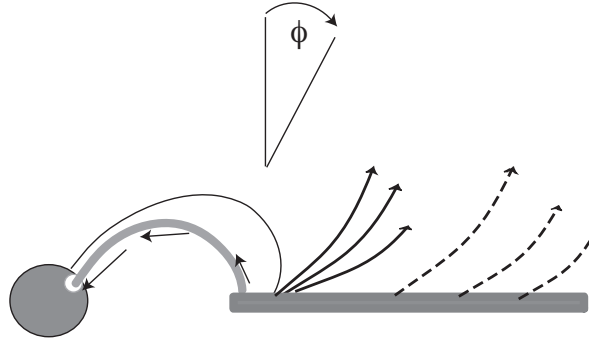


Fig. 7.2. Schematic structure for a connected system of accretion disk, stellar wind, and stellar magnetosphere. Magnetic fields which, owing to a finite magnetic diffusivity, penetrate the disk inside the co-rotation radius (where the angular velocity of the rotating disk matches the angular velocity of the star) allow material to accrete (gray curve); fields penetrating the disk outside of corotation help provide a spindown torque (solid dark curve). In the \times -wind model, the wind arises from the disk just at corotation (curved solid arrows), while disk wind models involve mass loss from a wider range of disk radii (dashed arrows). Magnetic field lines pitched at angles $\Phi_e > 30^\circ$ allow for rapid, cold mass loss (see text). [Fig. H-III:3.7]

from ionizing radiation, and possessing temperatures far too low to effectively ionize even low-ionization potential metals like Na and K, it seems highly unlikely that the MRI can account for (at least low-mass) star formation on its own.” More on this in Ch. 11.

7.2.4 Disk winds

An alternative to transporting angular momentum outward through an accretion disk going against, and thereby enabling, matter to spiral inward is to remove angular momentum by a variant of a stellar wind, namely one that is cold and propelled by centrifugal forces. [H-III:3.3] “The basic version of the cold, magnetically-driven wind takes advantage of the rapid disk rotation to fling material outward (and later collimate it). Near the disk it is assumed that the magnetic pressure is much larger than the gas pressure. In this limit, the magnetic fields are stiff at the launching region, *i.e.*, corotation of the inner wind is assured. In this case the energy (Bernoulli) constant of the motion [for a unit of mass] becomes

$$E = \frac{v_\phi^2}{2} + c_s^2 \ln \rho - \frac{1}{2} \Omega_o^2 r^2 - \frac{GM_*}{(r^2 + z^2)^{1/2}} = \frac{v_\phi^2}{2} + c_s^2 \ln \rho - \Phi_e, \quad (7.10)$$

where v_ϕ is the poloidal velocity, Ω_o is the (Keplerian) angular velocity of the disk in which the magnetic field is rooted, c_s is the (assumed isothermal) sound

speed, and Φ_e is an effective potential term including the effects of rotation and magnetic fields; [the terms in the central expression measure kinetic energy (first and third), change in internal energy in an isothermal process (second), and gravitational potential energy (fourth) at a distance r from the rotation axis of the disk, and height z above that disk]. The behavior of the flow depends upon the form of Φ_e , which in turn depends upon the geometry of the flow.

In the case of a perfectly vertical field, perpendicular to the disk, any material which flows outward must be propelled initially by gas pressure; the Keplerian rotation is of course insufficient by itself to drive outflow. The atmospheric structure is nearly hydrostatic until one reaches a radial distance such that

$$c_s^2 \sim \frac{GM_*}{(r^2 + z^2)^{1/2}}, \quad (7.11)$$

in analogy with a Parker thermal wind [(see Sect. 2.2 around Eq. 2.11)]. When the gas is cold, the flow 'starts' only at large radii; the flow interior to this must pass through many scale heights of density, resulting in negligible outflow.

In contrast, a field line tipped away from the rotation axis can effectively drive a cold flow, taking advantage of the $\frac{1}{2}\Omega_o^2 r^2$ term in Eq. (7.10). Neglecting thermal pressure,

$$E = \frac{1}{2}v_\phi^2 - \Phi_e, \quad (7.12)$$

where the 'effective' potential is

$$\Phi_e = -\frac{GM_*}{r_o} \left[\frac{1}{2} \frac{r^2}{r_o^2} + \frac{r_o}{(r^2 + z^2)^{1/2}} \right]. \quad (7.13)$$

Consider now a small displacement along the field line, with a coordinate given by s , and

$$ds^2 = dr^2 + dz^2. \quad (7.14)$$

At the base of the flow, the disk material is rotating at the local Keplerian velocity. This is an equilibrium state, because $d\Phi_e/ds = 0$ at $z = 0$. However, if $\partial^2\Phi_e/\partial s^2 < 0$, this equilibrium is *unstable*; any small perturbation along the field line will result in an increased (outward) poloidal velocity from Eq. (7.12). If θ is the angle between the field line and the disk plane, the critical stability criterion

$$\frac{\partial^2\Phi_e}{\partial s^2} = 0 \quad (r = r_o, z = 0) \quad (7.15)$$

requires $\tan^2\theta_c = 3$, or $\theta_c = 60^\circ$. Disk magnetic field lines which are tipped away from the rotation axis by an angle greater than 30° result in an unstable

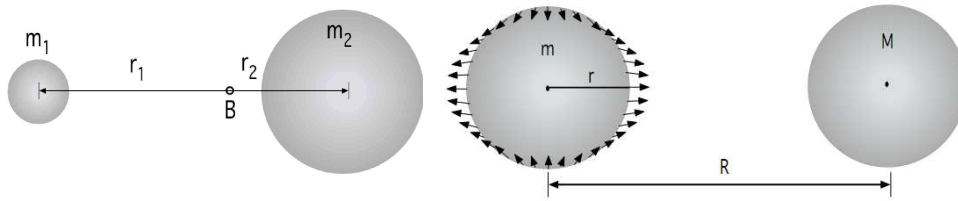


Fig. 7.3. Left: Two bodies orbiting around the barycenter B. Right: Tidal acceleration induced by the body with mass M on the body with mass m with the distance R between their centers. [Fig. H-III:11.3]

equilibrium, and rapid outflow will commence at the disk.” This flow carries angular momentum away from the disk.

7.3 Gravitational tides

7.3.1 Spin-orbit interactions

[H-III:11.2.1] “A well known gravitational influence is the tidal force of Moon and Sun on Earth. [Similar tides occur in other planet-moon systems throughout the Solar System – and has led to spin-orbit synchronization for most of the major moons – and also in binary stars and in star-planet pairs with relatively tight orbits – more on that below.] To calculate the tidal acceleration, let us consider two masses M and m with the distance R between their centers as shown in Figure 7.3. According to Newton’s law of gravitation, the mass m feels the gravitational acceleration a :

$$g = -G \frac{M}{R^2} \quad (7.16)$$

However, each point of a body with mass m and radius r feels a different gravitational acceleration depending on the effective distance to mass M which ranges from $R - r$ to $R + r$. For the two extreme cases we find:

$$g = -G \frac{M}{(R \pm r)^2} = -G \frac{M}{R^2(1 \pm r/R)^2}. \quad (7.17)$$

[In cases for which] r is much smaller than R this equation can be expanded into a Taylor series:

$$\frac{1}{(1+x)^2} = 1 - 2x + 3x^2 - \dots, \quad (7.18)$$

$$g = -G \frac{M}{R^2} \pm G \frac{2M}{R^2} \frac{r}{R} \mp \dots \quad (7.19)$$

The tidal acceleration a_t is the difference between the effective and the gravitational acceleration $\{A:[92]\}$:

A:92

$$a_t \approx \pm G \frac{2M}{R^2} \frac{r}{R}. \quad (7.20)$$

Note that a_t decreases with the third power of R . As a result of this the tidal accelerations are relatively small. On Earth the tidal acceleration is about $1.1 \times 10^{-6} \text{ m s}^{-2}$ due to the Moon and $0.5 \times 10^{-6} \text{ m s}^{-2}$ due to the Sun compared to the gravitational acceleration of about 10 m s^{-2} . This corresponds to an expected lunar tidal effect of about 70 cm. In reality, the average tide is about 30 cm because of a slight deformation of the Earth. In the case of the Sun, the tidal effects caused by the planets are very small; [t]he largest effects are due to Venus and Jupiter with a theoretical tide in the order of 1 mm.

As a result of the friction between the tide and the planet, the rotation [and revolution tend towards synchronizing]. In the case of Earth this [results in a slowing down of the spin rate by] about one second per year. Some 2.5 billion years ago the length of a day was only about 6 hours. Because the angular momentum must be conserved this leads to a corresponding increase in the distance between Moon and Earth (4 cm per year) as measured by laser technique. $\{A:[93]\}$ The tidal friction generates a power of $3 \times 10^{19} \text{ erg/s}$ which is mostly dissipated in the ocean. There are indications that this tidal power affects the global ocean circulation which plays a crucial role in the climate system by transporting energy from low to high latitudes. The tides act also in the atmosphere causing changes in pressure, temperature, and wave propagation.

A:93

There are climatic effects on Earth related to the lunar tides. The plane in which the moon moves is inclined to the ecliptic by about 5° . The points where the lunar orbit crosses the ecliptic are called nodes. As a result of the gravitational force of the Sun on the Moon the orbital spin axis of the Moon precesses, which leads to a continuous slight shift of the nodes. After 18.6 years the nodes are back to their original position. $\{A:[94]\}$ The inclination of

A:94

⁹² Activity: Looking only at gravitational forces, how close to a solar-mass object would the Earth need to be to be pulled apart by tides? Whereas this is impossible with the Sun, an Earth-sized planet could be pulled apart if it approached a white dwarf or neutron star (and something like that is involved in 'contaminating' some white-dwarf atmospheres with heavy elements). An object of lesser density can be pulled apart, however, during a sufficiently close approach to the Sun: estimate at what distance (ignoring tensile strengths, spin, and orbital forces) comet 67P (with a mass of about 10^{16} g and characteristic dimension of 3 km) would have to come to the Sun to be broken up. Some Sun-grazing comets (such as the Kreutz family) have been observed to go through this breakup process.

⁹³ Activity: Consider what it means for solar eclipses that the Moon is moving away from the Earth: at some future time, the Moon will be so far away that no more total solar eclipses can occur anywhere on Earth. Assuming the Moon continues to move away at 4 cm/yr, roughly when will the last total solar eclipse occur? Confirm that the answer is somewhat more than 600 million years.

⁹⁴ Activity: If you are interested in solar eclipses, and wonder why the saros cycle has a slightly different length from the lunar nodal period, have a look here.

the Moon's rotation axis has an effect on the amplitude of the tides. The amplitude of the lunar nodal tide is only about 5% of the daily diurnal tide but integrated in space and time it becomes significant. The 18.6 yr cycle and sometimes also its second subharmonic of 74 yr have been found in the arctic ocean temperature and sea ice extent and in drought records.

The dynamics of a multibody system such as the Solar System is largely determined by gravitation. The bodies orbit around the barycenter. In the case of a two-body system with a large body (Sun) and a small body (planet) the orbit is an ellipse with the large body in one of the focal points. {A:^[95]}

A:95

In a multibody system (Solar System) the gravitational interaction between the bodies disturbs slightly their orbital parameters. For example the planets (mainly Jupiter and Saturn) change the eccentricity of the Earth's orbit with periodicities of about 100,000 and 400,000 years which has an effect on the amount of solar radiation received from the Sun" determined by the orbital eccentricity; more on that in Sect. 12.3.2 around Eq. (12.10).

The effects of the solar tides on the Earth's orbit are negligible, but that will not stay that way. Late in the life of the Sun, as it runs out of fuel (see Ch. 10), the Sun will swell up into what is known as a red giant. In fact, its [H-III:4.11] "diameter increases by approximately two orders of magnitude. The physical expansion of giant stars results in the assimilation of many of the planets that may have formed in their formerly habitable zones [(defined as the distance from the star where liquid water can be present on a planet's surface)]. In the case of the Sun, current predictions indicate that the Sun will expand (Fig. 7.4) to nearly 1 AU, engulfing both Venus and Mercury. Because the Sun loses over 40% of its mass during [its phases as a red giant], Earth's orbit will actually expand to conserve angular momentum. This seemingly places Earth just beyond the presently modeled maximum diameter of the Sun but detailed modeling indicates that Earth will be assimilated into the Sun because of tidal effects. Tidal forces raise a bulge in the Sun's upper layers that follows Earth and provide a retarding force that causes Earth's orbit to decay. {A:^[96]}

A:96

Earth is totally vaporized by this process due to the power generated by its ~ 25 km/s entry into the Sun's upper atmospheric layers. If Earth had formed 15% further from the Sun it would have escaped assimilation. Mars and all other planets are well beyond the effects of gas drag and tidal effects and are

⁹⁵ Activity: One of the ways in which exoplanets are detected is to look spectroscopically at the displacements of the star about the barycenter of the exoplanetary system. How large is the velocity amplitude, and how large the associated Doppler shift at visible wavelengths, for the Sun-Jupiter system?

⁹⁶ Activity: What is the upper limit to the Sun's rotation rate in this phase? Formulate your arguments. You may ignore solar mass loss in this estimate. Use Fig. 10.5. This upper limit shows that the Sun's outer layers are rotating (much?) more slowly than the Earth is orbiting it, so that the tidal bulge on the Sun will be traveling through, and dissipating energy within, the solar outer envelope.

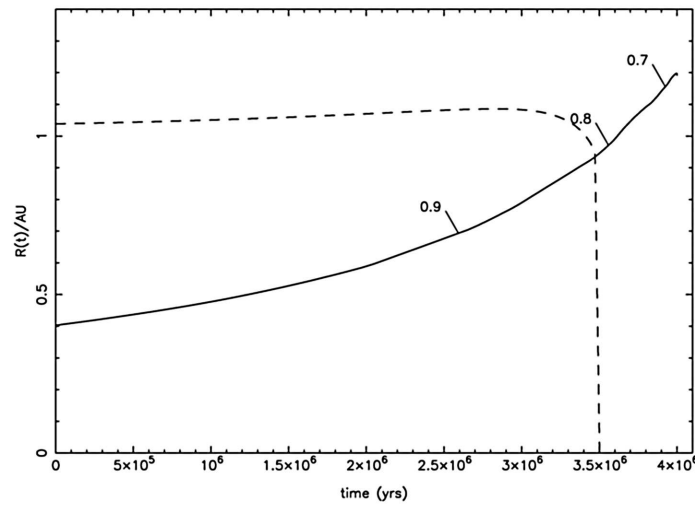


Fig. 7.4. The diameter of the future red-giant Sun (solid; in AU; [the labels along the curve show the Sun's mass at the time expressed in present-day solar masses]) and the size of Earth's orbit (dashed) during the 4 million years leading up to the phase when the Sun reaches maximum brightness. Earth's orbit expands slightly as the Sun loses mass but the Sun expands to the point where tidal drag causes Earth's orbit to decay and intersect the Sun's upper layers. These calculations predict that Earth will be destroyed in the Sun's atmosphere 7.59 billion years from present. [Fig. H-III:4.6; source: Schröder and Connon Smith (2008).]

safe from total destruction although they are severely heated and rendered lifeless during the Sun's red giant phase.”

This effect of orbital synchronization by gravitational tides occurs for all close-in planets, and has a particular consequence if we look at the evolution of the billions of planetary systems throughout the Galaxy. [H-III:4.11] “In the far future, the Universe will look quite different than it does at present. All massive bright stars will have evolved and become invisible. Only the slowly evolving and faint M stars will persevere. After several tens of billions of years of Galactic evolution, questions about habitability will only concern the bodies that remain, these faint low mass red stars and planets that orbit them in thin habitable zones close to their surfaces[; . . .] they are the most numerous stars now in the Universe and [. . .] in the long term they will be the only stars in the Universe. Compared to the Sun these low mass stars offer new challenges to understanding habitability. Although faint, they have pronounced flare activity which generates both UV and energetic particle fluxes capable of harassing life. Due to their faintness, their habitable zones are so close to the stars that planets can be tidally locked with one side always facing out to space. This can cause thin atmospheres to freeze out on the dark side of planets although

sufficiently thick atmospheres may be able to adequately distribute heat and prevent this calamity.”

Tides also have their consequences on the multitude of double stars: roughly one in two of the stars seen in the sky are pairs of stars orbiting their joint center of gravity. [H-III:2.5] “The gravitational tides in binaries with periods of order a week or less (depending on stellar masses and radii) are so strong that the orbital and rotational periods of these stars are synchronized on time scales much less than the main-sequence life time. Because any cool-star components of such binaries lose angular momentum through their wind, they will tend to spin down, but the tidal coupling replenishes the lost rotational angular momentum from the reservoir of orbital angular momentum. This causes the orbital separation to shrink, the locked orbital and rotational periods to decrease, and – counterintuitively – the activity to increase with age until eventually the stars merge into a single, rapidly-rotating but old star (forming the class of FK Comae stars).” {A:[97]}

A:97

Another consequence of gravitational tides in the case of a tilted rotation axis relative to the orbital plane is precession. [H-III:11.3] “The precession is a wobbling of the Earth’s axis of rotation which is caused by the tidal forces associated with the Moon and the Sun. Because the Earth is spinning, its shape deviates slightly from a sphere leading to an equatorial bulge. Tidal forces act on the bulge and force the axis to precess. The periods of [Earth’s] precession range from 19,000 to 24,000 years.” {A:[98]}

A:98

7.3.2 Orbital interaction

Differential gravitational forces are also thought to be of major importance in the formative phases of planetary systems, specifically acting between clumps of matter once these have condensed within the spinning accretion disk, and in even earlier phases when gravity may have led to unstable situations in which relatively dense areas may form by contraction and compression. [H-III:3.2] “As shown in the right-most part of Figure 7.1, [such early gaseous concentrations will be] sheared due to the differential rotation. The gravitational attraction of one ‘end’ of the spiral arm pulls on the other; this has the effect of

⁹⁷ Activity: Between the phases of tidally-locked binaries and merged binaries are (semi-)contact binaries in which mass transfer can occur as one of the binary components becomes larger than its ‘Roche lobe’, either because the star swells up in late evolutionary phases (see Ch. 10) or because the orbit shrinks by ‘magnetic braking’. Now or after reading Ch. 10, look up the definition of ‘Roche lobe’ and the properties of RS CVn, Algol, W UMa, and FK Com objects as characteristic phases in the evolution of close binary stars towards single stars (with the details, and the class names, dependent on the masses of the two components).

⁹⁸ Activity: The Earth’s equatorial bulge is nowadays used to keep satellites in a Sun-synchronous orbit, which is useful for satellites that need to scan the entire surface of the Earth, and also to enable Earth-orbiting satellites to have an uninterrupted view of the Sun throughout the year. Look up how this works.

accelerating the outer material at the expense of decelerating the inner material – *i.e.*, transferring angular momentum outward. [One among several distinct mechanisms (see Ch. 11)], it appears that this mechanism [of gravitational instability (GI)] will prevent most of the mass from remaining in the disk, but instead will allow accretion toward the central object.”

These same gravitational forces are likely to play a major role in enabling forming planets to grow into giants: growing planets set up wave-like density disturbances in large spirals, and the interaction of the growing planet with the matter in these spirals can cause all of the constituent parts to change their orbits, as long as the mass in the disk is not too small compared to that in the growing planet. [H-IV:1.2] “For example, the combination of observations and numerical experiments suggests that gas giants accumulate up to a few hundred Earth masses of material – first the solids and then increasingly rapidly gases – within a matter of a few million years. This process is aided in its efficiency by the migration of growing planets within the young planetary system: planets are not bound to their initial orbits, but can migrate either inward or outward, subject to gravitational interactions, thus having access to a large volume of the primordial disk from which to collect material. Interestingly, it appears that it is the very collection process of matter onto the growing planet that causes mass redistributions within the disk so that their tidal effects can make planets migrate, particularly if other planets are forming elsewhere in the system, while the gravitational coupling between multiple young planets in eccentric orbits can scatter bodies around (both in distance from their central star and in orbital inclination).”

[H-IV:5.7.2] “The realization that exoplanets are mobile during the early stages of formation has led to many studies of dynamical interactions. The details of migration and the parking mechanisms that [can lead to] gas giant planets just a few stellar radii away from their host stars are an active area of research. In the younger primordial disk with significant gas and dust density, the planet embryos will clear gaps in the disk. In this case, material can pile up at both the inner and outer edges of the gap. When the disk mass at the edges of one of these gaps is comparable to the mass of the planet embryo the disk will exert a torque that causes the planet [to either tighten or widen its orbit around the parent star, *i.e.*, causes the planet] to migrate. The outer edge of the disk causes inward migration while the inner edge of the disk can produce outward migration. When multiple planet embryos exist in the disk it is possible for the outer embryo to become locked into a resonant orbit with the inner planet, a process called convergent migration. As the disk clears, convergent migration can leave planets in resonant orbits that persist stably over the lifetime of the star. This effect is especially powerful for resonances

where the ratio of the orbital periods ($P_{\text{outer}}/P_{\text{inner}}$) is close to an integer number, N . Planets with small N are said to be in mean-motion resonance (MMR) and the exchange of angular momentum between MMR planets is flagged by oscillations in eccentricity and orbital periods.”

7.4 Planetary atmospheric tides

Apart from magnetic torques and gravitational tides, there is also a class of tides associated with irradiation. [H-III:15.7.1.2] “In general terms, tides are the periodic response to periodic astronomical forcing. In the [Earth’s] atmosphere, by far the dominant forcing agent is thermal excitation by solar radiation, although forcing by latent heat release [(*e.g.*, cloud formation)] can also be important. The dominant atmospheric tides are the diurnal tide and the semi-diurnal tide at double the frequency. In the lower and middle atmosphere, tides are excited primarily by the absorption of solar UV radiation by stratospheric ozone and solar near-IR radiation by tropospheric water vapor. The diurnal tide is forced about one-third by water vapor absorption and about two-thirds by ozone absorption. The semidiurnal tide is predominately forced by ozone absorption. Although the diurnal component of the diurnal variation of solar heating is stronger than the semidiurnal component, there is a rough parity between the two because the semidiurnal tide responds more efficiently to ozone forcing than does the diurnal tide. This is because the region of ozone forcing is fairly deep and main semidiurnal modes with their comparatively long vertical wavelengths respond in phase over the forcing regions, while the diurnal tide with its fairly short wavelengths experience a degree of phase cancellation.”

In [H-III:15.12.2] “[p]lanets with thick atmosphere [...], atmospheric tides can affect rotation. It is speculated that all planets [in the Solar System] formed with similar rotation rates and spun in the prograde sense (aligned with the total angular momentum of the Solar System). Gravitational torques can de-spin rotation toward synchronous rotation, but cannot produce retrograde rotation. The torques acting on the solar tidal bulge and coupling with the solid planet, however, can cause retrograde rotation and this is what may have produced the retrograde rotation of Venus. The present state of Venus is thought to be an equilibrium between gravitational and thermal atmospheric tidal torques. Clearly the resonances supported by planetary atmospheres can affect where equilibrium states might be found and thus the speed of retrograde rotation.”

More on tides and other large-scale wave phenomena in both oceans and atmospheres can be found in Ch. H-III:15.

8

Particle orbits, transport, and acceleration

Deep inside stars and planets energy is exchanged between particles (including photons) so frequently that the distribution of velocities of the ions and electrons in stars, and of the atoms and molecules in planets, are essentially pure Gaussians (and thus the distributions of the magnitudes of the velocities pure Maxwellians) around the mean bulk velocity. With sufficient collisional interactions in a neutral medium, or in an ionized medium in the absence of magnetism, the mean bulk flows of different species in a mixture tend to be equal. The presence of a magnetic field, in contrast, is associated with a difference in bulk motions between negatively-charged electrons and positively-charged ions. This, in turn, leads to collisional interactions that convert the kinetic energy of bulk population motions into random kinetic energy, *i.e.*, the dissipation of electrical current equates to heating.

Where collisional time scales grow to time scales approaching those of physical processes, or even exceed these, velocity distributions can deviate from Maxwellians. The populations of non-thermal particles of most interest in the context of heliophysics are those of the highest energies. Among these are radiation-belt particles, but also those that originate from outside the Earth's environment, and referred to as 'cosmic rays', which encompass solar energetic particles (SEPs), galactic cosmic rays (GCRs) and 'anomalous cosmic rays' (ACRs). {A:[99]} [xv]

A:99

[H-II:8.1] "To understand the ubiquitous presence of energetic particles it is important to realize that except for planetary ionospheres and the lowest layers

⁹⁹ Activity: The so-called 'anomalous cosmic rays' have a complex history: originally neutral particles in the interstellar medium, ionized by charge-exchange or photo-ionization in the solar wind, advected to the heliospheric extremes there to be accelerated. Important though they are as diagnostics of the outer heliosphere and the enveloping sheath-shock structure, they are not discussed in this volume. You can look them up for an interesting read . . . after finishing this chapter. See Fig. 8.5 for where they appear in the energy spectrum.

^{xv} For an introduction to how energetic particles are detected and their properties determined, see Ch. H-II:3.

of the Sun's corona and below, most plasmas in the heliosphere are basically collisionless. That is, the mean free path of charged particles is larger than most scales of interest. For example, in the undisturbed solar wind, the mean free path for ions is of the order of 1 AU [(see also Table 3.4)]. The lack of such collisions means that there exists no primary mechanism that forces the particles to assume thermalized Maxwellian distributions. In fact, observed distributions, often on top of thermal (colder) approximate 'core' Maxwellians, almost universally contain energetic tails, which usually can be described by power laws. In real-world plasmas, there is a multitude of processes responsible for generating such supra-thermal and high-energy tails; usually, so-called wave-particle interactions are involved."

This chapter touches on various aspects of how energy can be converted from large-scale dynamics of magnetized plasma into an increased energy content in the thermal reservoir, the energetic-particle reservoir, or both, as well as on the transport and loss of such energy once in these reservoirs. This chapter covers topics as diverse as GCR transport inward through the heliosphere to SEP transport outward from the corona; all of these topics have to do with conversion or transport of energy. The chapter starts with motions of individual particles and their transport within magnetic environments, then moves to mechanisms by which their energies can change to become so-called 'energetic particles'. A description of how energy from non-thermal particles is deposited into the thermal energy reservoir with particular focus on the solar corona is partitioned off into Ch. 9.

Flares, CMEs as well as magnetospheric (sub-)storms extract their energy from what has been somehow stored in the magnetic field. This extraction is typically enabled by the phenomenon of reconnection, and both the total flux involved and the rate at which reconnection proceeds help set the magnitude of energetic-particle events. The chapter touches on reconnection and shocks, which are essential ingredients in both heating and impulsive phenomena, but only introduces the basics of these complicated processes, which remain far from understood.

8.1 Single particle motion

[H-II:11.2.1] "The motion of every individual charged particle in the heliosphere can be described by the Lorentz force equation, Eq. (2.21). [...]" [H-II:9.2.2] "For the simplest case of no electric field and a constant magnetic field in the z

direction, the solution to Eq. (2.21) is straightforward. It is given by ^[xvi]:

$$v_x = +v \sin \alpha \cos(\omega_g t - \phi) ; v_y = -v \sin \alpha \sin(\omega_g t - \phi) ; v_z = +v \cos \alpha, \quad (8.1)$$

where $\omega_g = qB/(mc)$ is the cyclotron (gyro-)frequency, α is called the pitch angle (note that our definition is such that $\alpha = 0$ implies the particle is moving directly along the magnetic field), ϕ is the phase angle, and v is the magnitude of the particle velocity.”

[H-II:11.2.1] “A very important aspect of the Lorentz equation when discussing particle acceleration is that the electric field may change the energy of the particle but the magnetic field does not. This relation is shown by taking the dot product of the Lorentz equation with \mathbf{v} giving:

$$\mathbf{F} \cdot \mathbf{v} = q(\mathbf{v} \cdot \mathbf{E}) + \mathbf{v} \cdot (\mathbf{v} \times \mathbf{B}), \quad \text{or} \quad \frac{dW}{dt} = q(\mathbf{v} \cdot \mathbf{E}), \quad (8.2)$$

where W is the kinetic energy.

[In realistic situations magnetic and electric fields rarely occur in separate and uniform configurations. Even in the simple case of a dipole potential field,] the motion separates into three oscillatory types occurring at increasingly slower timescales, [visualized together in Fig. 8.1(right)]. On the fastest timescale, a particle gyrates around the field line as described above.

The second oscillatory type motion in the dipole relates to the particle’s velocity parallel to the magnetic field. As the particle follows the field line towards the poles, it moves through a gradient because the magnetic dipole field increases [when the particle approaches the planetary or solar] surface. The effect of this gradient is to convert the parallel motion of the particle into perpendicular motion as shown schematically in Fig. 8.1(left). As the particle moves toward the pole, the gradient effectively creates a Lorentz force opposite to the parallel motion. Eventually, the parallel velocity will go to zero and then reverse direction [ultimately] causing the particle to bounce between the southern and northern poles. The point at which the parallel velocity goes to zero is called the mirror point and the oscillation between the two poles is referred to as the bounce motion.

[In the case of a planetary magnetosphere dominated by a dipole, the particle will circle the planet] in an oscillatory manner known as drift motion. The azimuthal drift is caused by the radial gradient of the dipole field. Intuitively, this drift can be attributed to the changing gyroradius in different magnetic field strengths. In the stronger magnetic field the gyroradius will decrease and in the weaker field the gyroradius will increase creating the orbit shown in the

^{xvi} Note that the gyrating charged particle emits gyro-synchrotron radiation, thereby losing energy, so that this orbital motion approximated by Eq. (8.1) – and thus also in Eqs. (8.4) and (8.5) – is not sustained indefinitely.



Fig. 8.1. (left) Schematic diagram showing the Lorentz force as a particle moves into the magnetic field gradient at Earth's poles. [Fig. H-II:11.2] (right) Schematic diagram of particle motion in a dipole magnetic field. [Fig. H-II:11.4]

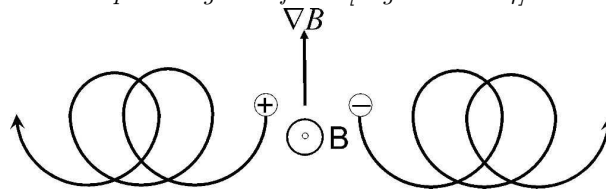


Fig. 8.2. Schematic diagram for the gradient-B drift. [Fig. H-II:11.3]

schematic of Fig. 8.2. As protons and electrons gyrate in opposite directions, they also drift in opposite directions.”

More generally speaking, for particles moving through a magnetic field with a mixture of waves and turbulence, the latter two processes transition from drifts to scattering – to which we turn later in this chapter.

Guiding center motion

[H-II:11.2.2] “Often, particle motion can be described by separating it into a drift velocity with gyromotion superimposed as in the examples provided here.

$\mathbf{E} \times \mathbf{B}$ drift: The $\mathbf{E} \times \mathbf{B}$ drift can be defined by including a uniform electric field in the Lorentz equation and separating the equation into components parallel and perpendicular to the magnetic field. [Let $\mathbf{B} = B\hat{\mathbf{z}}$.] In the parallel direction the Lorentz equation becomes

$$m\dot{v}_z = qE_z, \quad (8.3)$$

where E_z is the component of the electric field parallel to the magnetic field. This equation simply describes a particle accelerating along the magnetic field. In the perpendicular direction, assuming $\mathbf{E} = E_x\hat{\mathbf{x}} + E_z\hat{\mathbf{z}}$ [(so that $E_y = 0$)],

the Lorentz equation becomes

$$\dot{v}_x = +\omega_g v_y + \frac{q}{m} E_x ; \dot{v}_y = -\omega_g v_x. \quad (8.4)$$

Taking the second derivative of the velocity gives

$$\ddot{v}_x = -\omega_g^2 v_x ; \ddot{v}_y = -\omega_g^2 (v_y + \frac{E_x}{B}). \quad (8.5)$$

These equations describe gyration superimposed on [a] drift in the $\mathbf{E} \times \mathbf{B}$ direction. {A:[100]} A:100

General force drift: [If for the electrical force $q\mathbf{E}$ we substituted another force \mathbf{F} (such as gravity)] into the $\mathbf{E} \times \mathbf{B}$ drift equation creates a general force equation,

$$\mathbf{v}_F = \frac{1}{\omega_g} \left(\frac{\mathbf{F}}{m} \times \frac{\mathbf{B}}{B} \right). \quad (8.6)$$

This equation can be used to define the drift velocity caused by any general force. Other types of drift include curvature drift caused by a centrifugal force related to the curvature of the dipole field lines, polarization drift that results from a slowly varying electric field, and a gravitational drift. [.]”

Let us quantify the gradient and curvature drifts: [H-II:9.2.3] “For the special case in which $\nabla \times \mathbf{B} = 0$, these [are] given by:

$$\mathbf{v}_G = \frac{cW_\perp}{qB^3} \mathbf{B} \times \nabla |\mathbf{B}| \quad (8.7)$$

$$\mathbf{v}_C = \frac{2cW_\parallel}{qB^3} \mathbf{B} \times \nabla |\mathbf{B}| = \frac{2cW_\parallel}{qR_c^2 B^2} \mathbf{R}_c \times \mathbf{B} \quad (8.8)$$

where $W_\perp = (1/2)mv_\perp^2$ and $W_\parallel = (1/2)mv_\parallel^2$. Note that these expressions are for the case of non-relativistic particles. [The final expression in Eq. (8.8) is added to explicitly show the dependence on the curvature radius \mathbf{R}_c of the field; this latter expression holds also in a non-potential field.] {A:[101]} {A:[102]} A:101
A:102

However, in most applications of interest $\nabla \times \mathbf{B} \neq 0$. A more-general expression for the particle drift can be derived by expanding the magnetic field about the smallness parameter r_g/ℓ_t where r_g is the particle gyroradius and ℓ_t

¹⁰⁰ Activity: Use Eq. (8.5) to formulate (in a general vector expression) the magnitude of the $\mathbf{E} \times \mathbf{B}$ drift (in case you need a hint: assume the velocity can be described by an oscillatory component plus a constant drift).

¹⁰¹ Activity: Rewrite Eqs. (8.7) and (8.8) to show that the drift velocity scales as the product of the particle’s velocity and the gyroradius relative to the typical length scale in the gradient of the field, *i.e.*, as $v(r_g/\ell_t)$.

¹⁰² Activity: Why do you think that bounce and drift motions are commonly ignored for the solar corona but are of dominant importance in the terrestrial magnetosphere? Hint: look at Table 3.4.

is the characteristic scale of the variation of the magnetic field. The resulting guiding center drift velocity, in the non-relativistic limit, is given by:

$$\mathbf{v}_{gc} = \left[v_{\parallel} + \frac{cW_{\perp}}{qB} \hat{\mathbf{e}}_B \cdot (\nabla \times \hat{\mathbf{e}}_B) \right] \hat{\mathbf{e}}_B + \frac{cW_{\perp}}{qB^2} \hat{\mathbf{e}}_B \times \nabla |\mathbf{B}| + \frac{2cW_{\parallel}}{qB} \hat{\mathbf{e}}_B \times (\hat{\mathbf{e}}_B \cdot \nabla) \hat{\mathbf{e}}_B \quad (8.9)$$

A:103 where $\hat{\mathbf{e}}_B = \mathbf{B}/B$. ^[xvii] The gradient and curvature drifts are associated with the last two terms in this equation, which are in the direction normal to the magnetic field; however, it is important to note that there exists a component of the drift *along* the magnetic field in addition to these. {A:[103]}

When Equation (8.9) is averaged over an isotropic distribution of particles, one obtains the drift velocity $\mathbf{v}_d = (cmv^2/q)\nabla \times (\mathbf{B}/B^2)$, which is commonly used in models of cosmic-ray transport.”

The gradient-curvature drift in the terrestrial magnetosphere causes one of the primary mechanisms often discussed as an agent in space weather: the ring current. Seen from above the geographic north pole, positive particles drift clockwise and negative particles drift counter-clockwise. This differential motion leads to a westward ring current between about 2 and 9 Earth radii. This ring current is associated with a largely dipolar magnetic field with direction opposite to the Earth’s field. The variability of this current is caused by the injection of particles into, and leakage from, the magnetosphere associated with solar-wind variability. The Dst (disturbance storm time) index used in space weather characterizations quantifies the strength of the ring current. The variation in the surface magnetic field at Earth owing to the ring current is of order 0.1 – 0.23 mG (see Table H-I:13.5). The phenomenon of a ring current is captured in an MHD description in principle, but because of the interest in how particles of different energies and anisotropic pitch-angle distributions behave, the inner-magnetospheric ring current is generally studied with a custom ring-current model that then is coupled to MHD magnetospheric and solar-wind models. {A:[104]}

A:104

[H-II:11.2.3] “The Lorentz equation and drift velocity derivations provide a feel for how single particles [behave . . . but the analysis of] satellite measurements requires a more generalized view of particle motion because detectors

¹⁰³ Activity: Use a vector identity to show that the final term in Eq. (8.9) transforms into the central expression in Eq. (8.8) for a potential field.

¹⁰⁴ Activity: Estimate the orbital period associated with the drift velocity as in Eq. (8.7 for a purely equatorial motion for a proton with kinetic energy of 0.3 MeV orbiting, respectively, at 2 and 10 planetary radii r_p for, for example, Mercury, Earth (where the ring current is contained roughly within these distances), and Jupiter. Use the equatorial field strengths B_e as in Table 5.3 and $B(r) = B_e(r_p/r)^3$ for the equatorial dipole field. Is the non-relativistic approximation warranted for this proton? And for an electron of the same energy? Compare the relative size of the terrestrial ring current with the Chapman-Ferraro distance. How does this comparison work out for Mercury and what does that imply?

^{xvii} Note that Eq. (8.9) is a corrected version of Eq. (H-II:9.8).

do not measure the position and velocity of every particle in space to be propagated forward in time using the Lorentz equation. To this end, it is instructive to describe particle motion using aspects of the motion that are conserved when time variations of the magnetic field are slow. For charged particles in the magnetosphere, there are three such invariants associated with the gyro, bounce, and drift motion. Assuming that the invariants are conserved confines the particle location to within a shell [in a dipolar field such as that in the inner magnetosphere] about Earth.

First invariant: The first invariant is associated with the gyromotion of the particle about the field line and is given by:

$$\mu_m = \frac{p_{\perp}^2}{2mB}. \quad (8.10)$$

Here p_{\perp} is the relativistic momentum in the direction perpendicular to the magnetic field, m is the rest mass [...], and B is the field strength.

Second invariant: The second invariant corresponds to the bounce motion of a particle along a field line and is given by:

$$J = \oint p_{\parallel} ds, \quad (8.11)$$

where p_{\parallel} is the particle momentum parallel to the magnetic field and ds is the distance a particle travels along the field line. It is convenient to rewrite the second invariant in terms of only the magnetic field geometry by the following manipulation. If no parallel forces act on a particle then momentum is conserved along a bounce path and $J = 2pI$ where p is momentum and

$$I = \int_{s_m}^{s'_m} \left(1 - \frac{B(s)}{B_m}\right)^{1/2} ds. \quad (8.12)$$

Here s_m is the distance of the particle mirror point, $B(s)$ is the field strength at point s , and B_m is the mirror point magnetic field strength. If the first invariant is conserved then K , as defined below, is also conserved.

$$K = \frac{J}{2\sqrt{2m\mu_m}} = I\sqrt{B_m} = \int_{s_m}^{s'_m} (B_m - B(s))^{1/2} ds \quad (8.13)$$

[...]

Third invariant: The third and final invariant corresponds to the drift motion of a particle [and is given by:

$$\Phi = \oint A_{\Phi} dl = \int \mathbf{B} dS. \quad (8.14)$$

In this equation A_{Φ} is the magnetic vector potential, dl is the curve along which lies the guiding center drift shell of the electron, \mathbf{B} is the magnetic field

and dS is area.] Therefore, conservation of this invariant requires that an electron gyration always encloses the same amount of magnetic flux as it drifts [...] In a dipole field this is equivalent to saying that the electron remains at fixed radial distance. The Roederer L parameter, commonly written as L^* , is another useful form of the third invariant [often used for the terrestrial magnetosphere]:

$$L^* = \frac{2\pi\mu_p}{\Phi R_E}, \quad (8.15)$$

where μ_p is the magnetic moment of the Earth's dipole field. The L^* parameter is the radial distance to the equatorial location where an electron would be found if all external magnetic fields were slowly turned off leaving only the internal dipole field."

8.2 Phase space density and Liouville's theorem

[H-II:11.2.4] "Two more concepts are needed to finally interpret particle measurements from satellites: phase space density and Liouville's Theorem." [H-II:9.3.1] "The number of particles per phase-space volume is known as the phase-space distribution function, f , which is a function of the 6-dimensions of phase space and time $(\mathbf{p}, \mathbf{r}, t)$, where \mathbf{p} is the particle momentum vector ($\mathbf{p} = m\mathbf{v}$). The number density of particles at a given location at a given time, $n(\mathbf{r}, t)$ is related to the phase space distribution function by:

$$n(\mathbf{r}, t) = \int f(\mathbf{p}, \mathbf{r}, t) d^3\mathbf{p}, \quad (8.16)$$

where $d^3\mathbf{p}$ is the volume element of phase space. For example, for a Cartesian geometry $d^3\mathbf{p} = dp_x dp_y dp_z$ and for a spherical geometry it is $d^3\mathbf{p} = d\phi \sin\alpha d\alpha p^2 dp$ [...] (α is the pitch angle).

The differential intensity [...] is related to the phase-space distribution function by

$$J = p^2 f. \quad (8.17)$$

Sometimes this is written as dJ/dE . This has units of particles per area, per time, per energy, per solid angle. If one integrates J over energy and solid angle (*i.e.*, a spacecraft detector with a given acceptance cone that sums over all energy channels), the result is the *flux density* of particles, or the number of particles crossing per area per time."

[H-II:11.2.4] "Our interest in working with phase space density is that it can be used to understand how collections of particles move rather than individual particles. More specifically, Liouville's Theorem states that as the system evolves or moves along a trajectory in phase space the density

must remain constant. The proof of this theorem is illustrated intuitively by considering a volume of phase space. As the particles in the volume are subjected to forces their position and momentum will change but the trajectories of particles in phase space can never cross. Trajectories crossing would imply the physical impossibility that two particles with the same position and momentum subjected to the same forces go in different directions. Thus, the particles act as an incompressible fluid [in phase space]. As they move, the volume can change shape but the density remains the same.

At first glance, Liouville's Theorem seems to be an esoteric statement but in fact its application is quite powerful. The particle flux (number of particles per $\text{cm}^2 \text{s str keV}$) measured by a particle detector on a satellite, $J(E, \alpha, \varphi, \mathbf{x})$ where E is the energy, α is the pitch angle, φ is the gyro-phase, and \mathbf{x} is the position, can be directly related to the phase space density through the relation $J(E, \alpha, \varphi, \mathbf{x}) = f(\mathbf{x}, \mathbf{p})/p^2$. Liouville's theorem states that the phase space density does not change as the particles move along a trajectory. We also know that if time variations of the magnetic field are slow, a particle's trajectory must move along a contour of constant adiabatic invariants. Putting these two concepts together means that $f(\mu_m, J, L^*, \varphi_1, \varphi_2, \varphi_3)$ wherever it is measured must remain constant. (Here $\varphi_{1,2,3}$ are phase angles associated with each invariant. For simplicity, it is generally assumed that the phase space density does not vary with the phase angles.) Any change of phase space density implies that one of the invariants is broken. In fact, acceleration mechanisms always violate an invariant. Thus, an increase in phase space density expressed as a function of the adiabatic invariants is a sign that acceleration has occurred. Flux measurements, in contrast, can change simply because the magnetic field topology has changed making these data very difficult to interpret."

8.3 The collisionless Boltzmann equation

Let us start with a general view of what we can do with the phase space density, looking specifically at what it takes to change it (or at what it takes to maintain it so that the adiabatic invariants can be applied to particle trajectories). You can review this section quickly on first pass, then revisit this when you reach the end of Sect. 8.4.2.

Throughout the heliosphere, we can generally ignore collisions between charged particles, particularly for the particles residing in supra-thermal tails of velocity distributions. Consequently, the distribution function for heliospheric charged particles generally satisfies the collisionless Boltzmann equation, which is a continuity equation in the 6D space of momentum and location coordinates $\mathbf{w} = [\mathbf{p}, \mathbf{r}]$ and time: $\partial f / \partial t + \nabla \cdot (f \dot{\mathbf{w}}) = 0$ (a 6D mathematical equivalent of Eq. 3.4, absent sources and sinks), or in another formulation ([using the

index notation so that, for example, the vector for space coordinates is written $x_i = \mathbf{x} = x_x \hat{x} + x_y \hat{y} + x_z \hat{z}$, as for] momentum p_i (or velocity v_i); with acceleration a_i , and with implied summation over repeated indices):

$$\frac{\partial f}{\partial t} = -\frac{p_i}{m} \frac{\partial f}{\partial x_i} - F_i \frac{\partial f}{\partial p_i} + S - L, \quad (8.18)$$

where the components of the force F_i are given by

$$F_i \equiv ma_i = q \left(E_i + \frac{1}{mc} \epsilon_{ijk} p_j B_k \right) + \mathcal{S}_i. \quad (8.19)$$

Here, \mathcal{S}_i is a placeholder for any other force or sum of forces that may apply, including gravity; even radiative energy losses or gains could be incorporated (although we ignore these here). Sources S and losses L could represent couplings to other reservoirs, such as neutral atoms or dust, which could happen through charge exchange or photoionization. We ignore these terms further in this chapter. {A:[105]}

A:105

Note that low-order velocity moments of the Boltzmann equation for combinations of interacting particle populations yield the equations of fluid dynamics. Take, for example, the case of a fully ionized hydrogen plasma with phase-space densities f_e and f_i for electrons and ions. The suitable combinations of the Boltzmann equations Eq. (8.18) for these phase space densities after multiplication by mv^α and integration over velocity space for $\alpha = 0, 1, 2$ yield, respectively, the continuity equation Eq. (3.4), the momentum equation Eq. (3.5), and the energy equation Eq. (3.6). A complete set of fluid dynamics equations would continue with ever higher moments until the entire phase space density has been described, but that is not practical. Instead, the series is commonly truncated by some approximation, known as 'closure'; see also Table 3.2.

A persistent electric field (such as in reconnection processes, see Sect. 6.4), for example, can change a particle's energy when that is accelerated along the field. Forces that can change the energy of particle populations need to be retained explicitly in whatever we do with Boltzmann's equation. Fluctuations in the magnetic fields in space and time (such as in Alfvén waves), in contrast, do not change a particle's energy (more on that in Sect. 8.4): they do scatter a particle in pitch angle. Repeated scattering in a perturbation field that is symmetric in the probability of scattering a particle in either direction can be described as diffusion. With that realization, Eq. (8.18) can be reformulated in a quasi-linear approximation by separating large-scale trends from small-scale fluctuations, denoting the large-scale average flow \mathbf{u} , and capturing the net effects of the small scale fluctuations in diffusion terms {A:[106]} :

A:106

¹⁰⁵ Activity: Verify that without sources and losses, Eq. (8.18) – also known as Vlasov's equation – is a reformulation of Liouville's theorem, *i.e.*, $df/dt = 0$.

¹⁰⁶ Activity: **Advanced:** If you are interested in the origin of the terms in Eq. (8.20) you could review

$$\frac{\textcircled{a}}{\partial t} = \frac{\partial}{\partial x_i} \left[\textcircled{b} \frac{\partial f}{\partial x_j} - \textcircled{c} f \right] + \frac{\partial}{\partial p_i} \left[\textcircled{d} \frac{\partial f}{\partial p_j} - \textcircled{e} f \right] + [\textcircled{f} - \text{L}], \quad (8.20)$$

$$\text{with } F_i = -\frac{1}{3} p_i \frac{\partial u_j}{\partial x_j} + \dots \quad (8.21)$$

Here, u is the mean velocity of the scatterers, which equals the flow speed of the bulk thermal plasma provided that comparable power resides in waves traveling in opposite directions. The explicitly listed term in F_i above represents the adiabatic momentum change.

The expressions \textcircled{b} & \textcircled{c} and \textcircled{d} & \textcircled{e} in Eq. (8.20) reflect fluxes in physical space and in momentum space, respectively. Terms \textcircled{b} and \textcircled{d} reflect diffusive processes; in geometric space with diffusion parameter κ_{ij} (with diagonal elements describing diffusion parallel and perpendicular to the magnetic field, and off-diagonal elements quantifying particle drifts) and in momentum or velocity space with diffusion parameter D_{ij} (which includes, among other things, pitch-angle scattering that does not affect the particles' energy); whereas we can more readily appreciate the symmetry between these two spaces, the physics of the scattering processes now lies hidden in the two diffusion tensors (see, for example, Sect. 8.4). Terms \textcircled{c} and \textcircled{e} reflect advection, in geometric space in \textcircled{c} and in momentum space (by the forces acting on the medium) in \textcircled{e} (but note the mixed partial derivatives in \textcircled{e} which means different groupings are possible).

Eq. (8.20) informs us on how particles move in the coupled 6-dimensional realm of geometric space and velocity space. When we talk about the transport of either solar energetic particles or galactic cosmic ray particles through the heliosphere, we look primarily at transport in geometric space which involves terms \textcircled{b} and \textcircled{c} : transport is affected by scattering and advection (in addition to, *e.g.*, geometric expansion in a spherical geometry, which involves term \textcircled{e}).

When we look into acceleration (and thus also heating) mechanisms, such as for shocks in Sect. 8.5, we need to figure out how particles move about in momentum space, *i.e.*, using expressions \textcircled{d} & \textcircled{e} , while crossing the shock in geometric space with expressions \textcircled{b} & \textcircled{c} . It often helps to focus on parts of the overall function f . For example, for the bulk of the plasma well within the thermal range of Maxwellian distributions, with relatively short mean-free paths compared to system scales and with frequent interactions with the collective,

classic papers with fairly 'intuitive' introductions to the equation by, *e.g.*, Harm Moraal (1976) or Luke Drury (1983), the latter also including how term \textcircled{d} arises. Or you could look at the paper by John Quenby (1984) which also describes the so-called 'force-field solution' that you will find in Sec. 14.1.1 on cosmogenic radionuclides.

we have the MHD description (Ch. 3) and, for shocks, the Rankine-Hugoniot jump conditions (Sect. 5.3). In contrast, for a 'contaminant' population of solar energetic particles and galactic cosmic rays moving through, but to first order not interacting with, the background plasma flow of the solar wind, but being scattered by perturbations in its magnetic field, Eq. (8.20) provides a powerful tool, as we shall see next. Other descriptions below focus on narrow parts of the overall phase-space distribution, such as the supra-thermal particles that interact at shallow angles with a shock and scatter in the collective of particles around it, and for a sub-population of quite energetic particles with long mean-free paths that bounce back and forth across a shock in a ping-pong fashion as they are scattered by waves. With this perspective, let us look at how this all describes the propagation of energetic particles through the heliosphere and the creation of energetic particles at shocks.

8.4 Particle scattering and transport

[H-II:9.2.4] “To this point we have considered only smoothly varying electric and magnetic fields as compared to the radius of gyration of the particles, $r_g = v/\omega_g$. For such cases, the particle speed and pitch angle change very slowly compared to the cyclotron period. However, when the typical scale of the variation in the fields, L_t , is of the order of r_g , the speed, phase, and pitch angle can undergo more rapid changes. This leads to a form of scattering that is loosely analogous to classical scattering, although it differs in important ways. For instance, the particles do not collide off of one another, as in the lower portions of Earth’s atmosphere, nor do they collide off of large targets, like photons moving through a dense gas, but rather, they scatter off of irregularities in the magnetic field. Formally one can solve the equations of motion under the approximation that the amplitude of the magnetic fluctuations are small and show that there exists a resonance condition, $v_{\parallel} \sim L_t \omega_g$, for which the equations become undetermined. At such instances, the particle is said to ‘scatter’ and it reverses its pitch angle and its phase angle becomes randomized. [...]

Because particle scattering is a stochastic process, it is most useful to perform a statistical analysis on a large number, or *ensemble*, of charged particles. The relationship between the average particle motion and the magnetic field can be determined from the quasi-linear theory. It is found that the dynamical behavior of the distribution function obeys the standard diffusion equation in classical statistical physics. [...] [H-II:9.3.2] “It is important to keep in mind that this equation is strictly valid only for time scales that are long compared to the time in between scatterings (the scattering time) and spatial scales that are large to the distance traveled between scatterings (the mean-free path).”

[H-II:9.3.3] Because “the magnetic field in space exists in a highly electrically conductive plasma, the field moves with the flow of the plasma (it is said to be ‘frozen in’. In the limit of ideal magnetohydrodynamics (MHD), which is the limit we are concerned with for energetic-particle transport, there is no electric field in the frame moving with the plasma. Thus, as a charged particle scatters off of a magnetic irregularity, its energy in the frame of reference moving with the plasma remains unchanged. [Strictly speaking, this assumes that the magnetic field is stationary in this frame of reference which is factually incorrect because of the presence of waves with a variety of phase and group velocities, but a good approximation in the case of the transport of energetic particles that move much faster than the waves (*i.e.*, $v \gg v_A$, where v_A is the Alfvén speed).] From the perspective of such fast particles, the magnetic fluctuations, which provide the scattering centers, move with the bulk plasma. [...] In an inertial frame relative to which the plasma moves with a velocity u , the evolution of f satisfies the advection-diffusion] equation, which in one spatial dimension is given by”

$$\frac{\partial f}{\partial t} = \frac{\partial}{\partial x} \left(\kappa \frac{\partial f}{\partial x} \right) - u \frac{\partial f}{\partial x}, \quad (8.22)$$

[H-II:9.3.2] “where κ is the diffusion coefficient. For the case of charged particles moving in an irregular magnetic field, κ is related to the statistical properties of the magnetic field, in particular, its power spectrum.” Here, we interpret f as integrated over all velocity space, so looking only at total numbers as a function of space. Thus, Eq. (8.22) is the 1D version of Eq. (8.20) for a case with constant u , in which the velocity integral for expression (d) disappears because scattering under these conditions does not change the overall energy of the population and (e) because there are no other forces assumed to act on the plasma.

[H-II:9.3.2] “We note that for Eq. [8.22] we have assumed that the distribution function varies only in one spatial direction. This should not be confused with [...] the restriction on particle motion arising from fields that vary with only one spatial coordinate. By using Eq. (8.22), we have already assumed that the process is diffusive. If, for example, x is taken to be the direction normal to a mean magnetic field, then the use of this equation implies that the field must be fully three dimensional for cross-field diffusion to take place. The key is that the field is fully three dimensional but it is also statistically homogeneous in space.” [H-II:9.3.4] “In two dimensions, there are two diffusion coefficients, one for each direction (plus cross terms which we can ignore for now). Consider the motion of particles in a turbulent magnetic field ^[xviii] whose

^{xviii} This volume does not go into the generation and properties of turbulence; for an introduction within the context of heliophysics, see H-I:7.

average points along the z direction. Then, for example, in the x - z plane, the diffusion equation (neglecting the advection term discussed above and cross terms) is given by:

$$\frac{\partial f}{\partial t} = \frac{\partial}{\partial x} \left(\kappa_{\perp} \frac{\partial f}{\partial x} \right) + \frac{\partial}{\partial z} \left(\kappa_{\parallel} \frac{\partial f}{\partial z} \right), \quad (8.23)$$

where κ_{\perp} and κ_{\parallel} are the diffusion coefficients across the magnetic field and along it, respectively.

Because the time τ_s it takes for a charged particle in the heliosphere [or magnetosphere] to scatter is generally much longer than the time it takes to gyrate about a magnetic field (*i.e.*, $\omega_g \tau_s \gg 1$), particles tend to move much more closely along the magnetic field than across it. As such, κ_{\perp} is usually assumed to be much smaller than κ_{\parallel} . For this reason, many analyses simply neglect perpendicular transport. However, it is important to note that in many astrophysical plasmas of interest, perpendicular transport is the most important [...]

The motion of a particle across a magnetic field occurs in two ways: (1) the actual transfer of particles from one magnetic field line to the next resulting from scattering, or across the field arising from drifts, and (2) the motion of particles along magnetic lines of force that themselves meander in space in the direction(s) normal to the mean magnetic field. [...]"

[H-II:9.3.6] "In addition to scattering and advection with the flow, the particle speed itself can change. Principally, this can happen in two ways: (1) by scattering within a spatially varying flow [(*i.e.*, by term ③ of Eq. (8.20))], or (2) by diffusing in energy space because of collisions with randomly moving scattering centers. The latter of these two [(related to D_{ij} in term ④ of Eq. 8.20, and related to dispersion in scatterer speeds)] is called second-order Fermi acceleration, or stochastic acceleration. This is an interesting topic, but is not considered in our discussion here. We examine further the first case.

Consider a particle moving in a given direction in an inertial frame which then scatters. Energy is conserved in the local plasma frame, but in the inertial frame the particle either gains or loses energy depending on whether it is moving initially against or with the flow \mathbf{u} . Suppose that at one scattering, it initially moves against the flow, and gains energy in the inertial frame (this is a head-on collision). When it next scatters, it will be moving initially with the flow and will lose energy. If the flow is everywhere uniform, then the particle loses the energy it gained in the previous scattering and there is no net energy gain. But, if the second scatter occurs at a different flow speed, there is a net change in the particle's energy. The term that accounts for this behavior is

given by

$$\frac{p}{3} \nabla \cdot \mathbf{u} \frac{\partial f}{\partial p} \quad (8.24)$$

[(as in term ③ in Eq. 8.20 and the expression for F following it)]. Particles gain energy if this term is negative and lose energy if it is positive.

A particularly good example of this is particle acceleration at a shock. Consider the energy of a particle in a frame of reference moving with the shock. As a particle scatters in the flow behind the shock, it loses energy because the particle was initially moving with the flow. The particle then returns upstream where it scatters off of the incoming upstream flow leading to a gain in energy. The energy lost by the downstream scattering event is smaller than the energy gained by the upstream scattering event because the upstream flow speed is larger than that downstream. Thus, there is a net energy gain, which leads to an acceleration of particles [(more on that in Sect. 8.5)]. Note that at a shock, the flow goes from large to small (in the shock frame) so that the divergence is negative and Eq. (8.24) is negative, giving rise to acceleration.

It is also noteworthy that the energy change term is positive for the case of a constant radial solar wind speed. So, all charged particles *lose energy* in the adiabatically expanding solar wind!”

[H-II:9.3.7] “The resulting superposition of the terms that we have discussed above, lead to the cosmic-ray transport equation. It is given by

$$\frac{\partial f}{\partial t} = \frac{\partial}{\partial x_i} \left[\kappa_{ij} \frac{\partial f}{\partial x_j} \right] - u_i \frac{\partial f}{\partial x_i} + \frac{p}{3} \frac{\partial u_i}{\partial x_i} \left[\frac{\partial f}{\partial p} \right] + S - L \quad (8.25)$$

[(which is very nearly the diffusion version of the collisionless Boltzmann equation of Eq. 8.20, but with $D_{ij} = 0$).] Note that we have written the diffusion coefficient κ_{ij} in its full tensor form [...]

The cosmic-ray equation is remarkably general. It has been used widely in most discussions of cosmic-ray transport and acceleration over more than three decades. It is a good approximation provided there is sufficient scattering to keep the pitch-angle distribution nearly isotropic ^[xix], and if the particles move substantially faster than the speed of both the background fluid and the characteristic speed of the MHD waves contained in the plasma.”

[H-II:9.4] “All of the quantities in the transport equation, except for the diffusion tensor, are directly observed by spacecraft or can be accurately determined by using the hydromagnetic approximation. Consequently, determining transport coefficients poses a fundamental challenge in the modeling of cosmic rays.

In general, the diffusion tensor κ_{ij} is related to the magnetic field vector B_i ,

^{xix} This should not to be confused with anisotropic diffusion resulting when $\kappa_{\perp} \neq \kappa_{\parallel}$.

the diffusion coefficients parallel and perpendicular to the mean field, κ_{\perp} and κ_{\parallel} , and the antisymmetric diffusion coefficient, κ_A , as

$$\kappa_{ij} = \kappa_{\perp} \delta_{ij} - \frac{(\kappa_{\perp} - \kappa_{\parallel}) B_i B_j}{B^2} + \epsilon_{ijk} \kappa_A \frac{B_k}{B}, \quad (8.26)$$

where δ_{ij} is the Kronecker delta function ($\delta_{ij} = 1$ if $i = j$ and $\delta_{ij} = 0$ if $i \neq j$), and ϵ_{ijk} is the Levi-Civita symbol: $\epsilon_{ijk} = 1$, or -1 if (i, j, k) is an even or odd permutation of $(1, 2, 3)$, respectively, and $\epsilon_{ijk} = 0$ if any index is repeated. We have also introduced the antisymmetric diffusion coefficient κ_A . Note that the symmetric terms reflect the diffusion due to small-scale turbulent fluctuations; in contrast, the antisymmetric term contains the particle drifts caused by the spatial variations of the large-scale magnetic field.”

8.4.1 Solar energetic particles

[H-II:9.5.1] “A particularly simple, yet illustrative example of the use of the cosmic-ray transport equation is the evolution of impulsively released particles from a point source. This is presumably a reasonable representation of the physics of solar-energetic particle transport after their release onto open magnetic field lines following their rapid acceleration in the vicinity of a solar flare. Of course, we must recognize that the earliest arriving particles suffer very little pitch-angle scattering, and therefore, the transport equation is not useful for describing these particles, but is adequate to describe the long-time behavior.

A proper treatment of the impulsive SEP problem should necessarily include, as a minimum, the effects of diffusion, advection with the solar wind, and adiabatic cooling. Spherical coordinates with the origin at the Sun would be a good choice. The resulting equation, even when simplified by making various assumptions about the choice of parameters can be impossible to solve analytically. For our purposes here, which is simply for illustration and by no means is meant to be directly comparable to SEP observations, it suffices to consider a Cartesian geometry, a constant diffusion coefficient, and to neglect both advection with the flow and energy change. The result is simply [Eq. (8.22) with $v = 0$], which is the 1D diffusion equation. The solution for an impulsive injection of particles at $x = 0$ at time, $t = 0$ is given by

$$f(x, t) = \frac{N_0}{\sqrt{4\pi\kappa t}} \exp\left(-\frac{x^2}{4\kappa t}\right), \quad (8.27)$$

where N_0 is the number of particles released.

Figure 8.3 shows a plot of the distribution of particles, given by Eq. (8.27), at the location $x = 1$ AU, as a function of time (in days). The diffusion coefficient

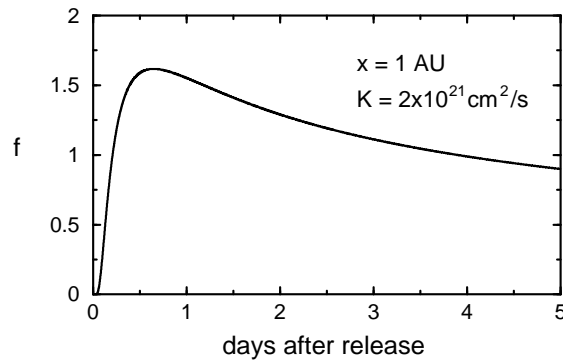


Fig. 8.3. Solution to the one-dimensional diffusion equation for a point-source release at a position 1 AU away from an observer: $f(1, t)$ from Eq. (8.27). [Fig. H-II:9.10]

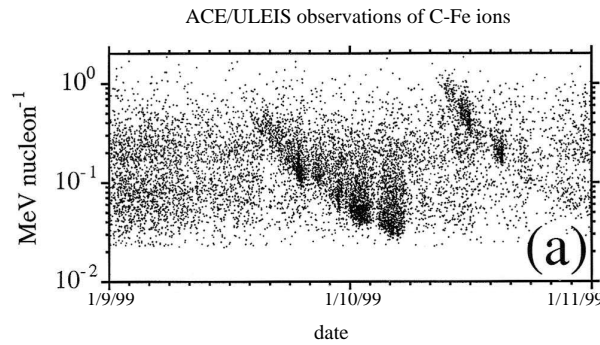


Fig. 8.4. A solar energetic particle (SEP) event, associated with an impulsive solar flare, seen by ACE/ULEIS. Each dot represents the detection of a particle by the detector. Two distinct events are shown. [Fig. H-II:9.11; source: Mazur et al. (2000).]

was taken to be $\kappa = 2 \times 10^{21} \text{ cm}^2/\text{s}$, and $N_0 = 10^{14}$. If, for example, these are 10-MeV protons, then the corresponding mean-free path would be about 0.1 AU. This profile has similarities to those seen at 1 AU following a flare or CME on the Sun [...]. An example of an impulsive-like solar-energetic particle event observed at 1 AU by the ACE spacecraft (ULEIS instrument) is shown in Fig. 8.4. Each dot represents a detection by the instrument of an individual particle. Plotted is the particle kinetic energy versus time. The earliest arriving particles are the ones with the highest energy since they move with the highest speed. The slower ones arrive later. This velocity dispersion leads to the characteristic profile shown in the figure.

It is clear from Fig. 8.4 that particles released at the Sun and observed near Earth undergo pitch-angle scattering in the inner heliosphere, because at any

A:107

given time there is a range of particle energies detected. That is, high energy particles can arrive later because they have scattered in the medium between the source and the observer. Thus, the 'thickness' of the comma-shaped particle event seen in the middle of this figure is related to the scattering frequency of the particles. {A:¹⁰⁷} Aside from this, however, there are many features in this event that are difficult to explain with a diffusive-advection-energy change approach [...]

It is noteworthy to point out another feature of the event shown in Fig. 8.4. There are intermittent dropouts in intensity during each of the two distinct events shown. These dropouts have been interpreted as resulting from the passage of alternatively filled and empty 'tubes' of particle flux by the spacecraft. The connection to the source, *i.e.*, the flare site, determines which field lines are populated with particles and which are not. [...]

These observations indicate that solar-energetic particles associated with impulsive solar flares undergo little cross-field transport, otherwise, these intermittent dropouts would not exist. This, of course, leads to the interesting puzzle of why galactic cosmic rays, or other types of energetic particles, do not exhibit such behavior. The answer is simply that the energetic particles in impulsive SEP events were relatively recently injected into the system and therefore have not had time to scatter sufficiently to become more spatially uniform. GCRs, however, have spent much more time in the Solar System (see Section 8.4.2). Thus, impulsive SEP events reveal the early time behavior of a collection of energetic charged particles moving in the heliospheric magnetic field.”

8.4.2 Galactic cosmic rays

[H-II:9.5.2] “GCRs are cosmic rays that pervade interstellar space and enter the heliosphere from the outside. The vast majority of them are swept out of the heliosphere before ever reaching Earth’s orbit. [... For] the purpose of a simple illustration of modulation, consider the steady-state Parker transport equation in one-dimensional spherical coordinates given by

$$\frac{\partial f}{\partial t} = \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \kappa \frac{\partial f}{\partial r} \right) - v \frac{\partial f}{\partial r} + \frac{2vp}{3r} \frac{\partial f}{\partial p} = 0, \quad (8.28)$$

[Note that this derives from Eq. (8.20) with (a) set to zero and for $D_{ij} = 0$ and uses that in spherical symmetry with an assumed constant solar wind speed

¹⁰⁷ Activity: For isotropic diffusion from a point source into 3-d space, the equivalent to the 1-d version of Eq. (8.27) is $f(r, t) = (N_0/(4\pi\kappa t)^{3/2}) \exp[-r^2/(4\kappa t)]$. Assuming that the particles of different energies 'scatter' off the same irregularities and that the diffusion coefficient is independent of position, use this approximation to estimate the release time at the Sun for the first event in Fig. 8.4, as well as the equivalent mean free path λ_{mfp} for a diffusion coefficient of $\kappa = \lambda_{\text{mfp}}^2/2\tau_s$ for a typical time between scatterings τ_s for a population of protons. Hint: remember Fig. 8.3.

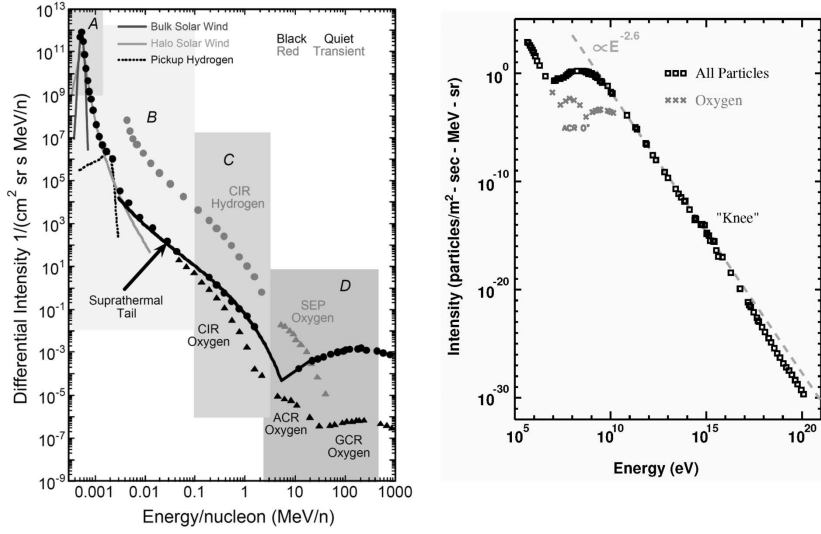


Fig. 8.5. Energy spectra of energetic particles in the heliosphere (left) and for cosmic rays (right). The curves illustrate the energy spectra during quiet time and disturbed solar wind conditions. The dots and triangles represent the supra-thermal part of the spectrum and the particles accelerated at corotating interaction regions (CIRs), galactic cosmic rays (GCRs), and the anomalous cosmic rays (ACRs) together with Solar Energetic Particles (SEPs). The right figure shows the high-energy part of the galactic cosmic ray energy spectrum to TeV energies. Note the characteristic peak at about 10 MeV and the $E^{-2.6}$ power-law dependence for energies above the peak. [Fig. H-IV:12.2]

u (and neglecting gravity) one has $\partial v / \partial x_i = 2v/r$] (this simple illustration neglects the effect of the heliosheath and termination shock). Here we have taken the diffusion tensor to be symmetric and $\kappa_{rr} = \kappa$.

It is convenient to rewrite Eq. (8.28) in the following form:

$$\frac{1}{r^2} \frac{\partial}{\partial r} r^2 \left(\kappa \frac{\partial f}{\partial r} - v f \right) + \frac{2v}{3rp^2} \frac{\partial}{\partial p} (p^3 f) = 0. \quad (8.29)$$

Generally this equation is not easy to solve, but if we assume that the last term on the left (describing the energy change of diffusing particles) is negligible, the resulting equation is readily solved to yield

$$f(r, p) = f(R, p) \exp \left(- \int_r^R \frac{v}{\kappa(r', p)} dr' \right). \quad (8.30)$$

Equation (8.30) gives an exponential decay of particles from the source ($r = R$) inward, into the the Solar System (where $r < R$). Moreover, it is reasonable to expect the diffusion coefficient to increase with momentum p so that

higher-energy particles have a larger diffusion coefficient than lower-energy particles. Thus, higher-energy particles have a longer exponential-decay length, or diffusive skin depth, than do lower energy ones. Thus, they more easily reach the inner heliosphere than lower-energy cosmic rays. This leads to a turnover in the spectrum that is due to modulation. This is in qualitative agreement with the observed cosmic-ray spectrum at Earth as shown in Fig. 8.5.”

The GCR intensity at a given orbital distance from the Sun is not a constant but varies with the solar cycle. [H-II:9.5.4] “Shown in Fig. 8.6 is the daily count of neutrons produced by the impact of cosmic rays on the upper atmosphere, from ground-based neutron monitors. This is an indirect measure of the cosmic-ray flux in near-Earth orbit. The time-intensity profile shows a clear 11-year cycle that is coincident with the sunspot-number cycle. During periods of high solar activity, sunspot maximum, the cosmic-ray flux is low, and during periods of low solar activity, or solar minimum, the cosmic-ray flux is high. In addition to this, there is also 22-year cycle present (the alternating ‘leveled’ *vs.* ‘rounded’ cosmic-ray flux), which, as we discuss below, is related to the drift motions of cosmic rays.

The increased modulation during periods of solar maximum is related to a combination of effects related to the shedding of magnetic flux by the Sun at solar maximum. On the one hand, increased solar activity leads to more magnetic turbulence which decreases the diffusion coefficient in the outer heliosphere leading to more modulation. On the other hand, and in addition to this, the merging of more numerous transient shocks and coronal mass ejections in the distant heliosphere creates magnetic barriers (so-called global merged interaction regions, or GMIRs) which also reduce the transport of cosmic rays into the inner heliosphere. There is a lower level of magnetic turbulence and fewer magnetic barriers for cosmic rays to propagate through during solar minimum. This is a qualitative explanation for the 11-year cosmic-ray cycle and its relation to the sunspot number cycle.

The 22-year cosmic-ray cycle seen in Fig. 8.6 is related to the 22-year solar magnetic polarity cycle [because the] polarity of the Sun’s magnetic field is important for the cosmic-ray drift that arises from the antisymmetric part of the diffusion tensor in Parker’s transport equation. [...]

Including the drifts of cosmic rays has led to the widely accepted paradigm for cosmic ray transport shown in Fig. 8.7. Drift motions for protons during two different solar polarity cycles are shown. During the period in which the solar magnetic field spirals outward in the north and inward in the south ($A > 0$, left panel) the GCR protons drift into the heliosphere from the polar regions of the heliosphere and outward along the heliospheric current sheet (which separates the two hemispheres and where the field reverses direction, hence

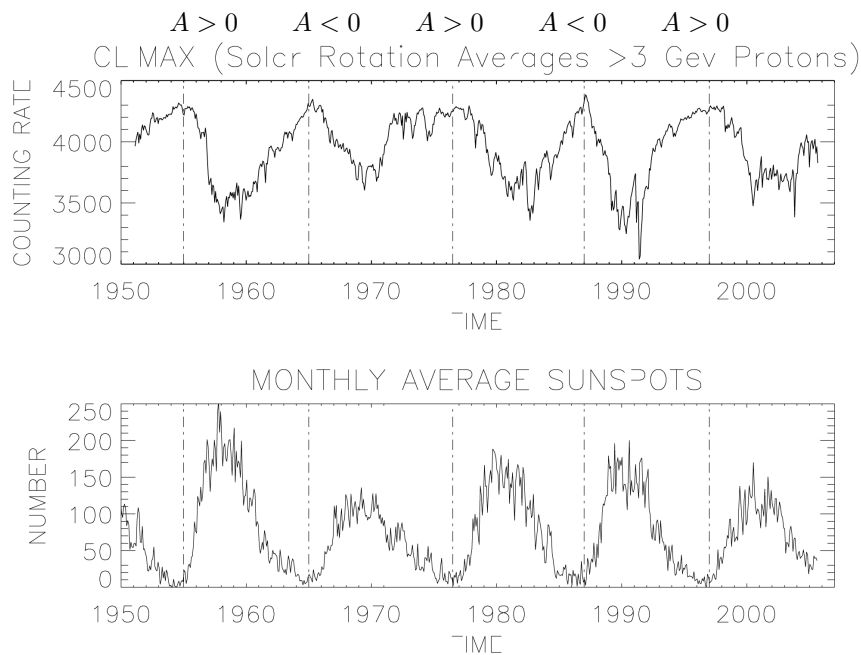


Fig. 8.6. [top] Climax neutron monitor daily count rate of neutrons produced by the interaction of a primary cosmic ray with Earth's atmosphere, which quantifies the modulation of galactic cosmic rays [near Earth orbit] during five sunspot cycles (shown in the bottom panel). Note the alternation in the cosmic-ray maxima between sharply peaked and more-rounded shapes. [The meaning of A is defined in Fig. 8.7. Fig. H-III:9.4]

the term 'current sheet'). During the opposite polarity, in which the solar field is inward in the north and outward in the south ($A < 0$, right panel), galactic cosmic-ray protons drift into the heliosphere along the current sheet. Note that in addition to the drift along the current sheet, there is also a gradient- B drift along the termination shock resulting from the jump in the magnetic field strength across the shock.

The explanation for the alternating leveled and rounded cosmic ray intensity involves both the drift motions of the cosmic rays shown in Fig. 8.7, and the 'waviness' of the heliospheric current sheet due to the offset of the solar magnetic axis and its rotation axis. When the 'tilt' is large, the current sheet is very warped, whereas, when it is small, the current sheet is much flatter (imagine the current sheet [forming above the rotating Sun with a tilted axial dipole, as in Fig. 5.7]). The current sheet is generally known to be relatively flat during the center of the solar cycle minimum. So, during the cycle in which the cosmic rays come into the heliosphere along the current sheet, only when it is very flat will the full cosmic-ray flux be reached at Earth's orbit.

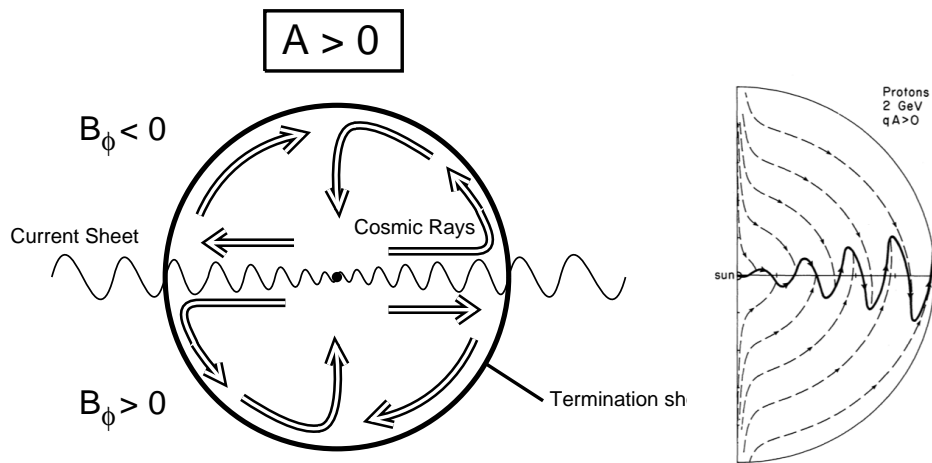


Fig. 8.7. [(left) Simplified side view of the heliosphere, with the current sheet depicted by wavy lines, to illustrate the drift] motion of cosmic rays in the heliosphere for [one solar magnetic polarity. The two polarities of the heliospheric] field are separated by the heliospheric current sheet. The value of $A > 0$ is for the period during which the solar magnetic field is outward in the north and inward in the south. During the next sunspot cycle, $A < 0$, the heliospheric field polarities are reversed, along with the direction of cosmic-ray advection.] The termination of the solar wind is also shown. [This is a cropped version of Fig. H-II:9.14.] [(right)] Cosmic-ray drift motions in a Parker spiral magnetic field with a current sheet. The arrows shown correspond to the time when the northern-hemisphere heliospheric magnetic field is outward from the Sun ($[A > 0$, as in] 1975, 1996) for positively charged particles. The arrows [in both panels] reverse for the alternate sign of the magnetic field ($[A < 0$, as in] 1986, 2007) and for the opposite sign of the particle's electric charge. [Fig. H-III:9.8; source right panel: (Jokipii and Thomas, 1981).]

Thus, during this phase, the cosmic-ray intensity will exhibit a rounded or 'peaked' time-intensity profile. When the cosmic rays come in along the poles of the heliosphere, the full intensity is reached much sooner and remains at a high level throughout solar minimum, and hence, during this phase, the time-intensity profile is more level, or flat."

8.5 Particle acceleration in shocks

Shocks provide an effective means to increase the kinetic energy of individual particles. For an ensemble of particles, this may shift their Maxwellian velocity distribution to a higher temperature, may distort that Maxwellian outside its core range, or lead to pronounced high-energy tails. Some form of shock heating and shock acceleration may play a role in processes as diverse as coronal heating (see Ch. 9) and the formation of solar energetic particles. [H-II:8.3.2] "Some of the processes that heat particles at thermal energies will also elevate

the energy at the upper part of the range. However, such enhancements are often only by a more or less constant, relatively minor factor. An example of this is the adiabatic heating of ions due to the magnetic field compression at the relatively narrow oblique or nearly-perpendicular shocks associated with the compression of the plasma” (for a quantitative example for low- β conditions as they occur, for example, in the solar corona, see Fig. 5.3). [H-II:8.2] “[T]he shocks of most interest to particle acceleration are MHD fast mode shocks, which compress both the density and the magnetic field. [Here, acceleration may span several orders of magnitude, and may significantly alter the shape of the energy distribution. However,] slow mode shocks may also play a role under certain circumstances.”

[H-III:8.3.1] “All mechanisms that contribute to the acceleration of charged particles at shocks rely on the particle orbits in the spatial and temporal features of electric and magnetic field environment of the shock. Roughly speaking, such processes are called kinetic when they go beyond the fluid (MHD) properties of the shock, when they are related to the scales associated with the charged particle motion, and when they require some self-consistent back-reaction between the charged particles and the plasma, *e.g.*, in the form of wave generation. For the highest particle energies, gyro-radii are so large that the size of the shock transition and even that of many local waves no longer matter. Conversely, for the thermal and so-called supra-thermal particles (just above the thermal energy to several thermal energies), the intrinsic shock scales and locally-generated waves do matter. As a consequence, the intrinsic shock scales and associated mechanisms play an important role not only for the general dissipation at the shock (the conversion to thermal energy), but also in providing a first, background level of energetic particles from ‘seed particles’ in the thermal and supra-thermal energy range. [...]

The two most important scales in collisionless shocks are the proton inertial length $\lambda_{\text{pi}} = c/\omega_{\text{pi}}$ (see Eq. 3.44) and the proton gyro-radius $r_{\text{gp}} = m_{\text{p}}v_{\text{c}}/eB$, which are related via the proton beta by $r_{\text{gp}}/\lambda_{\text{pi}} = \sqrt{\beta_{\text{p}}}$. [... The width of the transition for many shocks is the larger of λ_{pi} and the distance v_1/ω_{gp} which the upstream flow, moving at speed v_1 in the normal incidence frame (NIF, see Fig. 8.8), travels during the time $1/\omega_{\text{gp}}$ for a single gyration of a proton.] Exceptions are the almost perpendicular shock [see footnote x on terminology], which can be cyclically reforming and steepen to electron scales, and quasi-parallel shocks, which are not only reforming, but at sufficient Mach number have extended regions of steepening upstream waves, and highly non-linear turbulence downstream.”

[H-II:8.3.2] “In most shocks in the heliosphere, the thermalization of the upstream flow is primarily achieved via the ion dynamics, whereas the electrons

mostly 'just go along for the ride,' *i.e.*, they move almost adiabatically, with some subsequent scattering that fills otherwise inaccessible regions in the downstream velocity space. Any heating of the electrons (which can be quite small) is important in regulating the so-called cross-shock potential, because much of the electron phase space needs to be confined to the downstream by a potential, to prevent escape of the highly mobile electrons and to preserve overall charge neutrality. [...]

In typical shocks of the interplanetary medium, and in planetary bow shocks, it has been established that the reflection and gyration of the incoming ions plays a dominant role. At oblique shocks, part of the incoming ion phase space is reflected, but then convected back into the downstream. That is, after reflection, at sufficient Mach number, any upstream-directed parallel velocity of most thermal and even of many supra-thermal particles is not sufficient to overcome the general plasma drift into the shock. Much of the converted flow energy is initially stored in these gyrating ions, which during this process have attained elevated perpendicular temperatures from the magnetic field jump. Depending on parameters, it may take a while before these protons are thermalized downstream, typically in Alfvén wave turbulence driven by the temperature anisotropy $T_{\perp} > T_{\parallel}$. Generally speaking, the closer to perpendicular the shock, the more difficult it is for both particles and waves to escape upstream.

In contrast, in quasi-parallel shocks reflected (and partially gyrating) ions also play a role, but they can much more easily escape upstream against the flow, because the magnetic field direction is close to the shock normal. There, they generate both obliquely-propagating, compressional fast-mode waves, and parallel-propagating Alfvén waves. These waves can grow to large, non-linear amplitudes while convected back towards the shock, where the beam density and growth rate are largest. However, below Alfvén mach numbers of about $M_A < 2.8$, the majority of resonantly generated waves are no longer convected back and therefore do not steepen as easily and do not impact the shock any longer, thus resulting in fewer ions making it upstream to generate waves in the first place. [The] resulting lower level of turbulence also has a negative impact on ion acceleration to higher energies.”

[H-II:8.3.3] “For most heliospheric shocks, proton acceleration is of prime interest. Protons can easily reach energies of tens, if not hundreds of MeV, and as such have a large range of societal consequences such as malfunction or destruction of equipment in space, and posing danger to astronauts or crew and passengers of high-flying aircraft. Electrons, on the other hand, are rarely accelerated to comparable fluxes at these energies, except perhaps at processes well inside the magnetosphere that periodically lead to huge enhancements of trapped populations (see Section 8.6).”

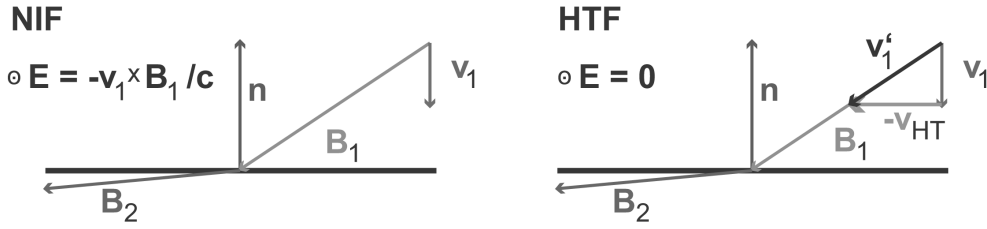


Fig. 8.8. Comparison of the normal-incidence frame (NIF) and de Hoffman-Teller frame (HTF) at fast-mode shocks. The NIF is the shock frame in which the upstream flow is aligned with the shock normal. [...] Transformation to the HTF is along the plane shock surface until the upstream flow vector coincides with the magnetic field. [...] Fig. H-II:8.2]

[H-II:8.4] “For ions, and for the energy range typically observed in the heliosphere, it is well accepted that two distinct acceleration mechanisms are at play:” (a) shock-drift acceleration and (b) diffusive shock acceleration.

(a) Shock-drift acceleration

[H-II:8.4.2] “[R]eflection of a portion of the incoming proton phase space [by the shock], and subsequent convection downstream, is the prime mechanism that eventually provides the [acceleration and heating] at quasi-perpendicular shocks. Even at highly oblique shocks, a small fraction of these ions will have sufficient parallel speed to make it upstream instead of being convected downstream, but the flux of such ions is strongly diminishing [with increasing angle between the shock normal and the upstream magnetic field,] making upstream wave generation increasingly difficult. Although the thermal proton gyroradius is typically comparable to the shock width, and that of supra-thermal ions clearly larger than the shock transition, surprisingly, many ions approximately behave adiabatically in simple shock transitions with sufficiently homogeneous upstream and downstream fields. A portion of the ion phase space then gains energy through their gyromotion under consideration of the shock electric fields. The family of such processes is called shock-drift acceleration (SDA).”

SDA [H-II:8.4.1] “is a ‘kinematic’ process in the sense that the particles simply perform their usual, mostly adiabatic orbits in the given, static or average electric and magnetic fields of the shock transition, neglecting any scattering. [...]” [H-II:8.4.2] “Consider a steady-state, one-dimensional shock. In this case, in the normal-incidence frame (NIF, see Fig. 8.8), there will be an out-of-plane electric field given by the cross-product of the upstream flow and magnetic field [...] $\mathbf{E}_p = -\mathbf{v}_1 \times \mathbf{B}_1/c$. [The MHD Rankine-Hugoniot jump conditions are such that the strength of this electric field is the same upstream

A:108

as downstream of the shock. {A:[108]} This motional electric field is aligned with the direction of both the curvature and gradient drifts associated with the jump in \mathbf{B} across the shock (see Sect. 8.1, and around Eqs. (8.7) and (8.9)), in such a way that ions gain energy by the gradient drift and lose energy through curvature drift.] It turns out that at quasi-perpendicular shocks, gradient drift wins out for most ions, which then gain energy proportional to the distance they drift along \mathbf{E}_p .”

A perspective change to another inertial system provides an alternative interpretation: in the so-called de Hoffmann-Teller frame (HTF, see Fig. 8.8) a translational velocity of $v_{HT} = v_1 \tan \theta_{B\perp}$ is introduced so that the upstream flow and magnetic field are aligned. As a result, there is no motional electric field in this reference frame, and only energy conservation and magnetic moment conservation come into play. In the HTF [H-II:8.4.2] “energy is conserved in the absence of other processes, and the only allowed change absent scattering is between the perpendicular and parallel velocity components. For close to perpendicular shocks, the field-aligned velocity component becomes increasingly larger due to the transformation into the HTF. [...] Because the perpendicular energy gain under magnetic moment conservation is simply a factor based on B_2/B_1 , only ions with sufficient initial perpendicular energy may exchange large fractions of their velocity components, while slowing down significantly or reflecting in the magnetic field gradient and in the cross-shock potential. Subsequent back-transformation shows that they have gained energy proportional to the squared transformation velocity. While this energy gain can be huge close to $\theta_{B\perp} \sim 90^\circ$, an increasingly smaller subset of phase space has sufficient perpendicular energy to effectively participate.”

Other mechanisms have been proposed, for electrons and ions alike, some, like ‘shock surfing’ acceleration (SSA) rely on the differences in gyro-radii and on the cross-shock potential; see Sect. H-II:8.4 for some more information. The challenge with all is that without additional scattering mechanisms, all these processes are too limited in the portion of phase space that is affected, and in the amount of energy gain, to explain the large, highly-energized populations often observed. It has been argued, however, that mechanisms like SDA and SSA can add energy for particles already energized by another mechanism, or that they provide the seed particles for such other mechanism to continue the energization. Turbulence or particle-induced waves may provide the scattering required to access more of the phase space.

¹⁰⁸ Activity: Review the Rankine-Hugoniot jump conditions (Eqs. 5.2 and 5.8–5.12) and show that the motional electric field $\mathbf{E}_p = -\mathbf{v}_1 \times \mathbf{B}_1/c$ is constant across a steady-state, one-dimensional shock.

(b) Diffusive shock acceleration

Diffusive shock acceleration (also known as first-order Fermi acceleration) [H-II:8.4.1] “is of ‘kinetic’ nature, in the sense that wave-particle interactions play the decisive role. As explained above, reflected or otherwise energized ions can easily escape into the upstream at quasi-parallel shocks, where they self-consistently” generate waves. Once grown sufficiently, these waves, and existing turbulence, diffusively scatter particles into a population that ranges from the far upstream to the far downstream. As the scattering centers converge owing to the compression associated with the shock, repeated scattering results in energization until they escape from the shock zone.

[H-II:8.4.3] “First-order Fermi acceleration produces a power-law distribution and intensities that depend on the shock strength (compression ratio ρ_2/ρ_1). Power-law distributions are as ubiquitous for SEPs as they are in cosmic plasmas, in general. [The] restricted temporal and spatial dimensions available lead to an upper cut-off of the spectra at high energies – typically between 10 MeV and 100 MeV for SEPs escaping interplanetary shocks. [...] For particles that are already significantly faster than the flow speed, the associated momentum gain of a returning particle is: $\delta p/p = (v_1 - v_2)/v(\cos\theta - \cos\theta')$, where the prime denotes the new pitch angle. [Note that the particles involved are quite energetic and therefore have a mean free path length exceeding the shock width; thus they sense the shock as a delta function, with the value of $(v_1 - v_2)$ in the above expression reflecting the step associated with term \textcircled{c} in Eq. (8.20).]

If one now assumes an almost isotropic distribution of particles, one can average over all pitch angles, and the cos terms simply convert into a constant factor. One then proceeds to calculate the probability of escape downstream (which is simply given by the ratio of the downstream to upstream flux) versus the probability of an acceleration cycle. From the calculation it follows that the particle distribution assumes a power law with index q , which depends on the shock compression ratio: $q = 3r/(r - 1)$, where from mass continuity in the assumed one-dimensional shock: $r = v_1/v_2 = n_2/n_1$, *i.e.*, the compression ratio between the downstream and upstream densities. [...]

Because waves that make up efficient scattering centers should be generated self-consistently by the energetic ions, must exist for extended regions upstream and downstream of the shock, and should not be convected towards or away from the shock too quickly, diffusive shock acceleration is most efficient and easiest understood for fairly high Mach number, almost parallel shocks. Conversely, it is much less understood how this process can be so efficient at the low-to-medium Mach number, oblique shocks that make up most interplanetary

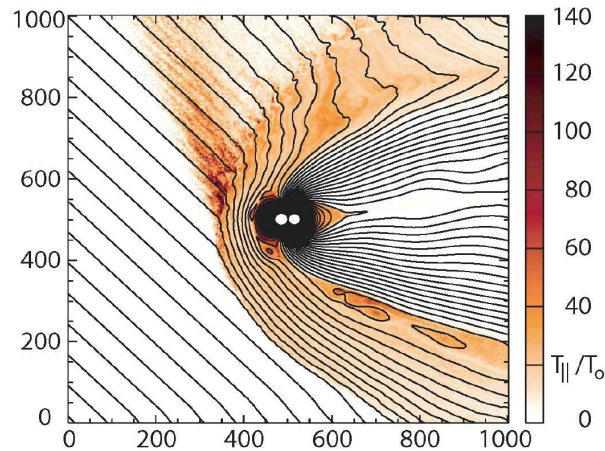


Fig. 8.9. Example of a two-dimensional (2-D) hybrid simulation of the solar wind – magnetosphere interaction. Shown are magnetic field lines (upstream IMF angle $\theta = 45^\circ$) and, the normalized parallel ion temperature T_{\parallel} , as a proxy of ion acceleration. As well-documented in many observations of the Earth’s bow shock, the ion foreshock starts close to $\theta_{B\perp} = 45^\circ$ with energized and back-streaming ions, and simultaneous excitation of waves (visible in the field line undulations). Conversely, at this scale, and with the number of pseudo-particles used in the simulation, there are virtually no upstream ions at larger shock-normal angles. [Fig. H-II:8.3; source: Krauss-Varban et al. (2008).]

shocks. In particular, at nearly perpendicular shocks, diffusive acceleration may require effective scattering across the magnetic field.”

There are multiple challenges to overcome in the study of diffusive shock acceleration, including the large number of particles that need to be tracked in numerical models, the relative roles of the various processes and their dependence on specific geometries, the generation of adequate turbulence to scatter particles, the complexities of the self-generated upstream wave field with multiple possible modes, the escape of the particles that have been energized from the upstream wave field as well as their further propagation through the turbulence in the solar-wind field, and the role of second-order Fermi acceleration in which particles scatter off counter-propagating waves, and, of course, the vast range of scales that needs to be treated. [H-II:8.4.4] “It is also known that multiple shocks generate a much more efficient acceleration environment. Not only does the first shock leave a much more turbulent and seed-particle rich upstream for the following shock, but particles may scatter multiple times in both shocks. Of course, the upstream seed particle spectrum and background turbulence are highly variable in the solar wind in general, and will have an impact on achieved fluxes.”

The Earth's bow shock

[H-II:8.5.1] “[P]lanetary bow shocks are of finite size, and as such, any production of energetic particles is both localized and highly non-local: some regions (*i.e.*, the quasi-parallel portion) are much more able to easily generate energetic ions, while any ions propagating upstream, or waves excited upstream of the oblique portion are quickly convected to a different portion of the finite-size bow shock, or around the obstacle, altogether. The general scale size of the Earth’s bow shock is of the order of $20R_E$ (Earth radii); the stand-off distance is [typically some $15R_E$ T]here is an ion foreshock that starts somewhere below $\theta_{B\perp} \sim 45^\circ$ and permeates the quasi-parallel domain, while the faster electrons form a foreshock boundary close to the perpendicular shock.

Figure 8.9 shows a snapshot of a 2-D bow shock simulation to further demonstrate this point. [. . .] The turbulence upstream and downstream of the quasi-parallel portion is clearly visible, as is the large enhancement of upstream-propagating, energetic protons. Conversely, there is virtually no upstream activity at or beyond 45° . [This approximate description conforms with the general state of the terrestrial environment and] also illustrates why there is so little activity upstream of the oblique portion beyond $\theta_{B\perp} \sim 45^\circ$: any ions that manage to make it upstream of the oblique portion, and any waves generated there, are either convected into the quasi-parallel portion of the bow shock, or instead move past the finite-sized obstacle altogether.”

Interplanetary shocks

[H-II:8.5] “Both co-rotating interaction regions and CME-driven shocks are capable of accelerating charged particles; however, not surprisingly, the largest events are associated with the fastest CMEs and can reach Alfvén Mach numbers of 5 to 6, and occasionally even higher. These Mach numbers are comparable to the Earth’s bow shock; yet, energetic particle energies and fluxes observed at the bow shock are almost dismal compared to those at the largest CME-driven events. Yet, while the Earth’s bow shock virtually always generates upstream energetic ions, the same cannot be said for IP shocks. [. . .] Finally, the heliospheric termination shock is also generally viewed as capable of producing highly-energized ions.”

[H-II:8.5.2] “Interplanetary (IP) shocks have a great variety of strength, and most of them are actually not particularly active when it comes to energetic particles. At the other extreme are IP shocks that are associated with strong solar energetic particle events (SEPs). Today, it is thought that SEPs are

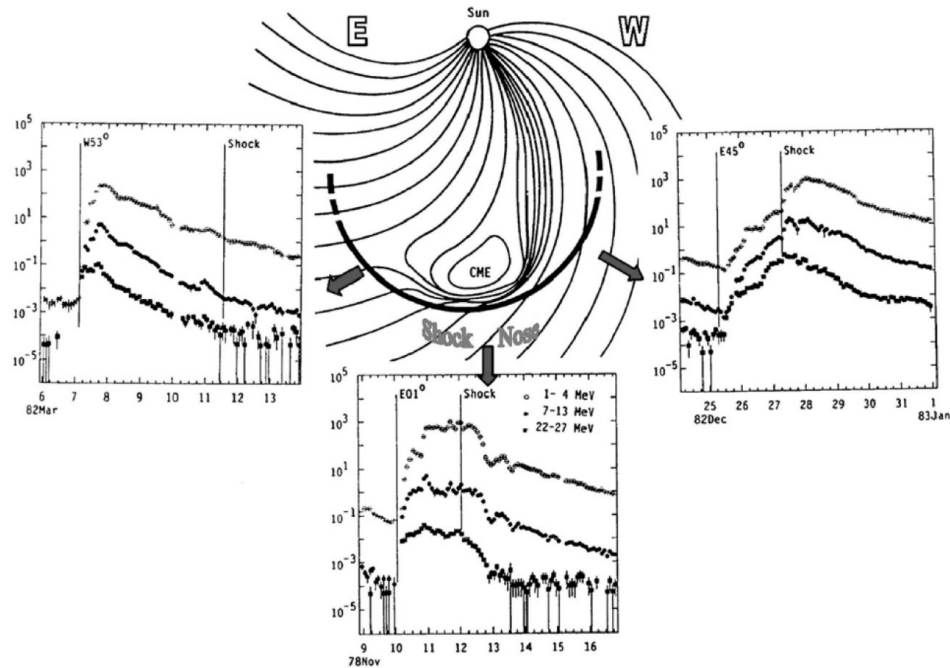


Fig. 8.10. Intensity-time plots of particle fluxes ejected from three different solar longitudes with respect to the nose of the shock front. [H-IV:12.7; source: Reames (2013).]

generated both by flare processes deep in the solar corona, and by shocks driven by coronal mass ejections (CMEs) [...]. In fact, it is estimated that almost all of the magnetic energy released in flares goes into energetic particles, with perhaps approximately equal share between the ions and electrons. These particles show up as 'prompt' events when observed at Earth: extremely energetic ions can traverse the distance from the Sun in minutes, with little delay compared to observed X-ray flare signatures at the Sun.

Conversely, so-called 'gradual' solar energetic particle events are generally accepted to be associated with coronal and interplanetary (IP) shock acceleration, driven by coronal mass ejections (CMEs). Even in this case, the most energetic particles are produced when the shock is in the corona, with resulting hard spectra that are observed at Earth within tens of minutes. However, production of energetic ions continues to 1 AU and beyond, and peak fluxes, with a softer spectrum, often arrive at Earth with the shock itself – historically called energetic storm particle (ESP) events [(Fig. 8.10)]. {A:^[109]} {A:^[110]}

A:109

¹⁰⁹ Activity: Interpret the flux profiles as a function of time shown in Fig. 8.10 for the three different perspectives (with Earth in the direction of the arrows for three different events). Argue for the

[F]orward propagating interplanetary (fast mode) shocks near 1 AU typically have an Alfvén Mach number $M_A < 3$, [only rarely reaching $M_A > 4$. D]espite their low Mach number, about one half of [the sample of] observed shocks had identifiable (albeit relatively low energy) upstream, energized ions and associated waves. [The] distribution of these ions [are] fairly isotropic, whereas in only a few cases, upstream beams were observed. This behavior also extends to higher energies and may be interpreted as a consequence of the large spatio-temporal scales of IP shocks, which rarely allows one to see the initial evolution of wave-particle interactions. While the large scales provide an important clue, and energetic seed particles may play an additional role, currently no scenario self-consistently accounts for the observed energetic ion environment of the weaker and oblique shocks.”

A:110

8.6 Relativistic particles in planetary radiation belts

8.6.1 Electron acceleration mechanisms

Earth’s electron radiation belts are located at $L \approx 3 - 10$, typically peaking around $L \approx 4 - 5$ (where L is the radial distance in Earth radii where magnetic field lines of an unperturbed dipole would cross the magnetic equator). [H-II:11.4.1] “Many acceleration mechanisms have been proposed to explain electron radiation belt flux increases at Earth but their exact contributions are still debated. Proposed acceleration mechanisms are often separated into two categories: internal (or local) source acceleration and external source acceleration. External source acceleration mechanisms are so named because they move electrons from outside geosynchronous orbit (at $6.6 R_E$) to the inner magnetosphere accelerating electrons through the transport process. They operate over large spatial and temporal scales that violate the particles’ third adiabatic invariant. Internal source acceleration mechanisms, on the other hand, locally accelerate electrons in the inner magnetosphere inside of $6.6 R_E$. They operate on fast timescales and small spatial scales and violate all three adiabatic invariants. The most prominent of the proposed mechanisms in each category are listed below. [...]”

differences in timing of the solar event (first vertical line in each panel) and the passage of the shock (second line) relative to the timing of the peak fluxes.

¹¹⁰ Activity: Energetic particles gyrate around the field lines in the solar wind. Roughly from what longitude region should we expect energetic particles to reach Earth that were created in flares or in shocks close to the Sun? From what longitudinal region at the Sun should we expect energetic particles in ‘gradual events’ to originate. What is roughly the delay between a flare/CME and ‘prompt’ particle storms? Explain the wide range of delays that can occur for flare/CME and ‘gradual’ particle storms?

8.6.1.1 External acceleration mechanisms

[H-II:11.4.1.1] “The manner in which external mechanisms accelerate particles can be illustrated starting with the assumption that the first adiabatic invariant [(Eq. (8.10))] is conserved. These mechanisms move electrons radially inward where the magnetic field is stronger. Because μ_m is conserved during the transport process, the increase in field strength requires that the particles’ perpendicular energy also increase. The total energy gain is directly related to the amount of radial transport. The relationship between transport and acceleration is easy to describe using the conservation of the first adiabatic invariant but the explanation hides the complex physics of the acceleration. Ultimately, it is an electric field that transports and accelerates the electrons because the magnetic field cannot change the particle energy. What separates the acceleration mechanisms is the exact form and timescale of that electric field. The electric field in both shock-induced acceleration and substorm induced acceleration is a large-scale inductive electric field that sweeps through the magnetosphere as the global magnetic field changes.

The shock-induced electric field is caused by the compression of the magnetosphere as shocked solar wind passes Earth. [...] However, such large sudden events are rare [while] smaller more pervasive compressions do not contribute significantly to electron radiation belt flux increases. Thus, shock acceleration is usually only discussed for specific events and not the very common flux increases that occur with most geomagnetic storms.

The substorm electric field is produced when the stretched magnetotail is pinched off near $10 R_E$ and the remaining plasma is hurled Earthward resulting in a more dipolar magnetic field configuration. [With this mechanism, numerical models have difficulty in transporting electrons inside of $10 R_E$; it may be that substorms] contribute to a seed population of electrons at large radial distance but some other mechanism, such as radial diffusion, is necessary to bring the electrons into the inner magnetosphere. Hence, much of the acceleration debate focused on radial diffusion.

In the case of radial diffusion, the electric field is that of ultra-low frequency (ULF; 300 Hz to 3 kHz) waves that continuously agitate the magnetosphere. [...] The basic premise of the mechanism is that electric fluctuations induce small random perturbations of the electrons’ position causing them to diffuse radially throughout the magnetosphere. The process is similar to diffusion in a gas only in this case the random walk motion of the particles is caused by electric fields instead of collisions. [...]

[Time] varying fields fluctuating specifically at the same frequency of an electron drifting about Earth [cause] rapid acceleration through a ‘drift resonance’. Figure 8.11 gives a pictorial explanation of an electron drift resonance. [...]

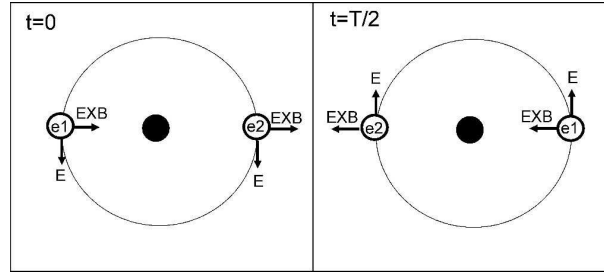


Fig. 8.11. Schematic diagram of an electron in drift resonance with a ULF wave. The left panel shows two electrons labeled e_1 and e_2 , the direction of the wave electric field, and the direction of the particle's $\mathbf{E} \times \mathbf{B}$ drift at time $t = 0$. The right panel shows the same properties half a wave period and electron drift period later. [Fig. H-II:11.7]

The electron drifting around Earth labeled e_1 [...] experiences an azimuthal electric field that continuously moves it inward causing the electron to gain energy. However, the electron, e_2 , that began at time $t = 0$ on the opposite side would have seen an electric field that pushed it radially outward in the same manner. Thus, the drift resonance causes electrons to diffuse radially inward and outward and decelerates as well as accelerates electrons.

[...] If electrons are acted on by [a randomly varying electric field, the net energy gain of the distribution of electrons] is determined by the phase-space density as a function of L . If electrons are uniformly distributed in L then the same number of electrons moves inward and gain energy as those that move outward and lose energy with no net energy gain. If the slope of f versus L is positive more particles move inward and gain energy than particles move outward and lose energy and the distribution of electrons gains energy. If the slope of f versus L is negative then the opposite occurs. [...] The radial diffusion has been described by an approximation to the Fokker-Plank equation]

$$\frac{\partial f(L, \mu_m, K, t)}{\partial t} = L^2 \frac{\partial}{\partial L} \left(\frac{D}{L^2} \frac{\partial}{\partial L} [f(L, \mu_m, K, t)] \right). \quad (8.31)$$

Here $f(L, \mu_m, K, t)$ [(with K defined in Eq. 8.13)] is the phase-space density of electrons and D is the diffusion coefficient which is calculated separately for electric and magnetic field perturbations. [Later,] the theory was revisited and elaborated to include higher-order resonances caused by electron drift motion in more realistic non-dipolar fields that increase diffusion. However, doubt about the ability of radial diffusion to fully explain observations led to the development of new competing ideas regarding electron acceleration including the internal source acceleration mechanisms.”

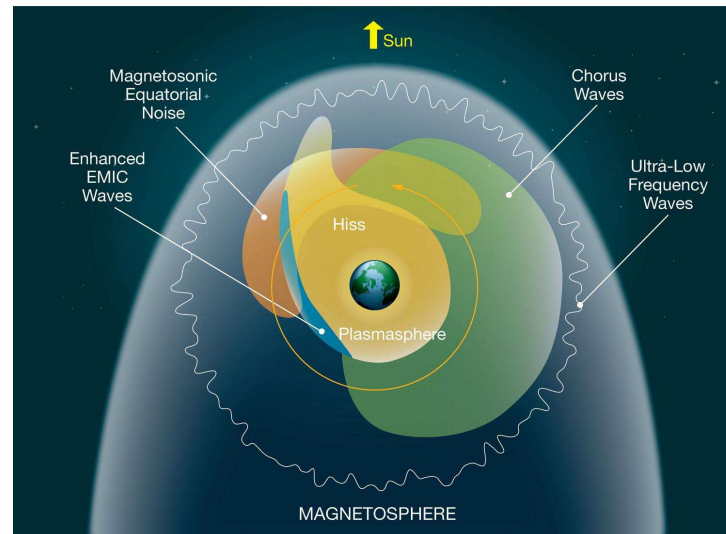


Fig. 8.12. Characteristic wave types within a magnetosphere, here visualized for the terrestrial case, viewed from above the arctic. Credit: NASA's Goddard Space Flight Center/Mary Pat Hrybyk-Keith.

8.6.1.2 Internal source acceleration mechanisms

[H-II:11.4.1.2] “The internal source acceleration mechanisms discussed here accelerate electrons through interaction with the electric field of a VLF (3 kHz to 30 kHz) wave. The interaction is similar to the ULF wave resonance, but in this case the resonance occurs between the wave electric field and the gyration of the particle about the magnetic field instead of the drift about Earth. The EMIC-Chorus wave mechanism assumes the interaction with the wave can be described as a random walk diffusive process very similar to radial diffusion [(Fig. 8.12 illustrates characteristic domains of various waves in the terrestrial magnetosphere)]. This assumption is only valid when wave amplitudes are small. The non-linear whistler wave acceleration mechanisms describe how electrons interact with a monochromatic set of large amplitude waves when diffusion is no longer valid. {A:^{[111]}}}

A:111

The resonance between an electron and a VLF wave can be illustrated by considering a VLF wave propagating at an angle θ from the direction of the magnetic field with magnetic and electric field perturbations perpendicular to the direction of propagation. The electron gyrating about the magnetic field will experience a constant electric field from the wave when the gyrofrequency of the electron equals the Doppler-shifted frequency of the wave [...]. In contrast

¹¹¹ Activity: Look up the properties of (whistler) chorus waves, (ELF/VLF) hiss, and EMIC waves.

to the ULF wave resonance, the VLF wave resonance will affect both the electron's energy and pitch angle. [...]

The *Chorus-EMIC wave mechanism* proposes that electrons interact with both whistler Chorus and EMIC waves as the electron drifts about Earth in such a fortuitous way that the distribution is steadily pushed to higher energy. In this model, EMIC waves at dusk interact with electrons to produce an isotropic pitch angle distribution. The electrons continue their drift to the dawn side of the magnetosphere where Chorus waves are predominantly found. The diffusion curves for Chorus waves are such that an isotropic distribution will diffuse towards higher energy and larger pitch angles. The energized electrons now peaked near 90 degrees continue around to the dusk side of the magnetosphere where the EMIC waves are found. The EMIC waves interact with the electrons to again produce an isotropic pitch angle distribution but with no energy loss. This isotropic distribution is now primed to interact with the Chorus waves once again and gain energy. Because the electrons traverse the magnetosphere in less than 10 minutes, the mechanism can effectively increase the energy over periods of days.

The *non-linear whistler wave mechanisms* assume that electrons are energized through a resonant interaction with whistler waves. However, the previous diffusion model requires that wave amplitudes are small in order for diffusion to be an adequate approximation. If this is not the case, the interaction must be described in a more detailed manner. [... Under] the right conditions a 100 keV electron could be accelerated to MeV energies within minutes. These mechanisms have yet to be compared in detail with observations or included in any kind of global model of electron flux. However, new measurements of whistler waves suggest that the small amplitude assumption is very often invalid making non-linear modeling an active area of interest.”

8.6.2 Proton acceleration in the radiation belt

[H-II:11.4.3] “The structure and temporal variability of the proton radiation belt is strikingly different from its electron counterpart. Yet, some of the same mechanisms are proposed to explain the acceleration of these particles. The protons normally form only one belt [around $L = 1 - 3$] with fluxes that peak near $L = 1.5$ and they tend to be more stable. However, during highly geomagnetically active periods, such as brought about by the passage of a large shock and sometimes an accompanying solar energetic particle (SEP) event, fast and dramatic changes occur. Often these changes mean a complete reconfiguration where entirely new, sometimes transient proton belts are formed that may last days to years.

Simulations of proton motion in both analytical and MHD magnetic field

models suggest that the new proton belts are formed when protons are transported radially inward by large induced electric fields that arise as a large shock passes the Earth. The mechanism is almost the same as proposed for some electron radiation belt acceleration events at Earth, except that forming a new proton belt requires an additional source of protons from the solar wind. Often large shocks are accompanied by very high fluxes of protons that are released from the Sun and further accelerated by the shock. Normally, Earth's magnetic field acts as a protective bubble that only allows these solar protons to enter over the polar caps where they are absorbed into the atmosphere. However, as the shock passes Earth, the magnetic field is distorted such that the accompanying protons can gain access to the inner magnetosphere and become trapped in the field. Once trapped, they are swept up by the induced field and pushed to small radial distances and higher energies to form a new belt."

8.6.3 Radiation belt losses at Earth

A [H-II:11.5.1] "survey of electron radiation belt changes [...] found that only 53% of storms cause radiation belt flux levels to increase even though these storms signify increased energy input to the magnetosphere. In 19% of storms the flux actually decreased and in 28% the flux did not change [by more than a factor of two]. The variable response to energy input suggests that loss and acceleration rates are often comparable and ultimately compete to determine final flux levels. [...] The mechanisms that have been proposed to explain the loss of relativistic electrons are: drift out the magnetopause boundary, outward radial diffusion, and scattering into the atmosphere. Scattering can be caused by interactions with a thin current sheet, EMIC waves, or Chorus waves.

Loss of electrons through the magnetopause boundary occurs when the drift paths of electrons are altered as the magnetic field changes from a quiet time configuration to more disturbed conditions. During quiet times, the drift motion of an electron starting in the magnetotail is dominated by an electric potential field directed from dawn to dusk [(see Fig. 5.16)] that moves electrons Earthward. As the electrons get closer to Earth, the magnetic radial gradient causes a westward drift. Some of these drift trajectories will cross Earth's magnetopause and the electron will be swept away by the solar wind. Closer to Earth, the trajectories of the electrons will be dominated by the gradient drift. Undisturbed, electrons in this near-Earth region will simply drift about continuously on closed almost circular paths. [Model] results suggest that during geomagnetic storms, most of the outer electron radiation belts are emptied into the solar wind and replaced by an entirely new belt of accelerated electrons. [This plausible suggestion has not been observationally. It appears]

that loss to the magnetopause was not an adequate explanation for electron flux depletions observed during more quiet conditions because the flux of energetic protons on similar drift paths did not decrease.

Radial diffusion [...] has also been proposed as a loss process. Radial diffusion acts to reduce gradients by pushing particles from high phase-space density to low phase-space density. The outermost closed drift orbit of the radiation belts represents a very steep gradient where the phase-space density goes to zero. If ULF waves are present, then radial diffusion will push particles outward to the magnetopause. [...]

Losses into the atmosphere occur when some mechanism scatters electrons to smaller pitch angles causing them to travel farther down the field line and collide with the neutral atmosphere. The current sheet, that forms in the magnetotail as the lobes are stretched and forced together by solar wind dynamics, is an effective scattering region. Scattering occurs when the magnetotail becomes stretched to the point that an electron bouncing along a field line can no longer make it around the kinked field without violating its first invariant. Traversing the kink changes the particle's pitch angle. Under certain conditions the pitch angle changes can be described as a diffusive process. [The] significance of this loss contribution has yet to be verified.

Chorus and EMIC waves [...] may also cause rapid loss into the atmosphere. Whether or not the waves produce net acceleration or loss depends on the initial gradients of the electron distribution as a function of pitch angle. [If] the appropriate distribution exists, EMIC waves are expected to cause losses on the timescales of several hours to a day. Whistler Chorus may cause losses on timescales of one day, but these estimates are sensitive to parameters such as the cold plasma density. Loss rates may increase to timescales less than a day during storm main phase when the plasma density is expected to vary [...]"

[H-II:11/5/2] "The proton losses from the radiation belts have not been analyzed in the same details as the dramatic formation of new belts. New belts last from days to years. Mechanisms proposed to explain the disappearance of these belts include scattering caused by the kinked field, and interaction with EMIC waves. [There is no firm understanding of the proton loss mechanisms, and] it may be that more than one mechanism plays a role in each event."

9

Convection, heating, conduction, and radiation

One topic at the foundations of heliophysics remains to be introduced before we move on to comparative stellar and planetary astrophysics: which processes lead to a relatively steady background heating of the solar atmosphere in 'quiescence', *i.e.*, outside of obvious impulsiveness? These processes have their origin in the convection that occurs below the solar surface and in the diversity of waves that are generated by these convective motions.

9.1 Convective and radiative energy transport

The solar convection zone persists from a depth of about 200,000 km all the way to the surface. [H-I:8.1] "Looking at the solar photosphere, we see the top of the convection zone in the form of granulation: Hot gas rising from the solar interior as part of the energy transport process reaches a position where the opacity is no longer sufficient to prevent the escape of radiation. The gas radiates and cools, and in doing so loses its buoyancy and descends. A:112 {A:[112]} At the surface the gas density is of order 10^{-7} g/cm³ while the pressure is of order 10^5 dyne/cm², but decreasing exponentially with height with a scale height of some 100 km. (This small scale height is the reason that the solar limb appears sharp as viewed with the naked eye.) Granular cells have dimensions of order 1 Mm, but numerical simulations indicate that convective length scales rapidly become larger as one proceeds below the solar surface. A:113 {A:[113]} These motions, ultimately driven by the requirement

¹¹² Activity: What drives tropospheric convection? Why is there no significant convection in the stratosphere (consider the role of ozone)? Is there an equivalent of a stratosphere on Venus? On Mars? Is there lower atmospheric convection? Formulate your arguments. The Web can help. The answers are 'yes', but not with a role for ozone except on Earth.

¹¹³ Activity: The scale of the granulation in the photosphere of the Sun (and analogously of other cool stars) follows from a comparison of energy loss by radiation (effective once the plasma can radiate into space, with a time scale of order 20s) and supply by upflows. Work through this estimate: just below the photosphere, the largest contribution to energy being carried upward resides in latent heat of recombination of ionized hydrogen (with an ionization fraction of order 0.1);

that the energy generated by nuclear fusion in the Sun's core be transported in the most efficient manner, represent a vast reservoir of 'mechanical' energy.

Looking closer, we see that granulation is not the only phenomenon visible at the solar surface. The quiet and semi-quiet photosphere is also threaded by magnetic fields that appear as bright points, as well as darker micro-pores and pores. These small-scale magnetic structures are, while able to modify photospheric emission, subject to granular flows and seem to be passively carried by the convective motions. Bright points are organized in a honeycomb-like pattern on the solar surface with a size scale larger than granulation, roughly 20 Mm; this pattern defines the so-called supergranular network and is suggestive of convective cells larger than granulation extending deeper into the solar interior.

Convective flows are also known to generate the perturbations that drive solar oscillations. Oscillations, sound waves, with frequencies mainly in the band centered roughly at 3 mHz or 5 minutes are omnipresent in the solar photosphere and are collectively known as p-modes ('p' for pressure). These p-modes are a subject in their own right and studies of their properties have given solar physicists a unique tool in gathering information on solar structure — the variation of the speed of sound c_s , the rotation rate, and other important quantities — at depths far below those accessible through direct observations.

[...] {A:^[114]}

A:114

On average the photospheric gas pressure of $p_g = 10^5$ dyne/cm² is much greater than the pressure represented by an average unsigned magnetic field strength of 1 – 10 Gauss ($p_B = B^2/8\pi < 10$ dyne/cm²) that is observed. However, in the largely isothermal chromosphere, the gas pressure falls exponentially with a scale height of some 200 km, while the magnetic field strength falls off much less rapidly, even as the field expands [from the compact flux tubes that characterize it in the photosphere] and fills all space. Thus, depending on the

balance that with photospheric black-body radiation; use this to derive a minimum upward flow v_z needed to balance radiative losses. Then match timescales, and use that $v_h \leq c_s = (kT/m_p)^{1/2}$: for overturning convective flows, the horizontal time scale of ℓ_h/v_h should equal the vertical one ℓ_z/v_z , for a typical horizontal granular scale ℓ_h and overturning depth $\ell_z \approx H_p \approx 400$ km somewhat below the photosphere. (This argument is developed in H-III:5.2.1)

¹¹⁴ Activity: Sound waves in an isothermal, hydrostatically-stratified atmosphere are 'evanescent', *i.e.*, non-propagating, at frequencies below the acoustic cutoff frequency $\omega_a = c_s/2H_p$. Can you argue intuitively why (think of the need for a restoring force roughly within a wavelength)? Estimate the value of ω_a for the solar atmosphere. Why are p-modes only observed at frequencies below about ω_a ? Now derive an approximate dispersion relation in simplified geometry: Start from Eqs. (3.4), (3.5) and (2.3) for a hydrostatic 1-d plasma (mind the sign of g) and combine them retaining terms to first order for perturbations $\rho = \rho_0 + \rho_a$ and $p = p_0 + p_a$, where ρ_0 and p_0 describe the background stratified atmosphere at rest. Then factor out the exponential growth of the amplitude with height by substituting $v = u \exp(z/2H_p)$ and use $u \propto \exp(i[kz - \omega t])$ to obtain a dispersion relation that has propagating waves (real values of k) only for frequencies above ω_a (a somewhat different approach can be found in H-I:8.3). Lower frequency waves reflect and can form standing p-modes if they meet the criteria for global resonance, while higher frequency waves can propagate, but will steepen (readily into shock waves) as they propagate into the lower-density atmosphere.

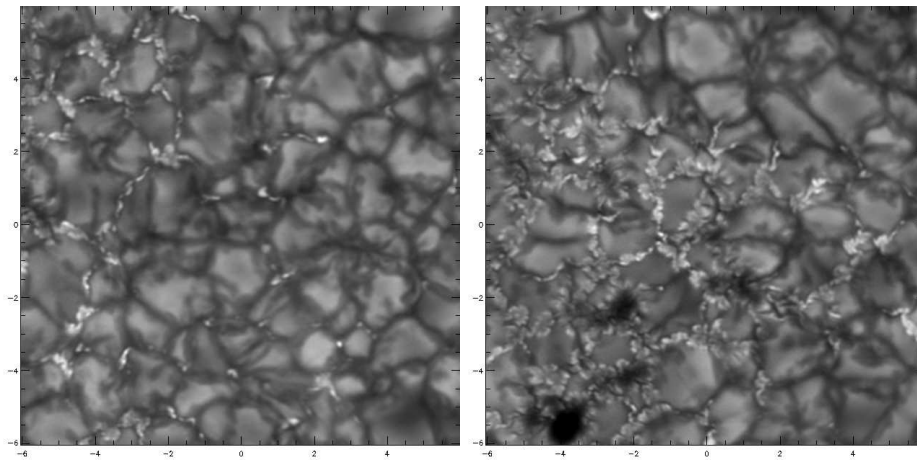


Fig. 9.1. The image on the left shows a typical quiet photospheric region observed in the G-band with the Swedish 1-meter Solar Telescope. The image on the right shows a plage region where the total amount of magnetic flux penetrating the photosphere is larger. The axes of both panels are numbered in arcseconds measured on the Sun; 1 arcsec is approximately 725 km. The G-band near 4300 Å contains several spectral lines, notably the lines of the CH-molecule, and is formed near the solar surface (the height where the optical depth $\tau_{5000\text{Å}} = 1$); the granulation and intergranular lanes some 100 km above this height, bright points some 200 km below — as explained in the text. Bright points are regions of enhanced magnetic field embedded between the granular motions. Notice also that bright points are pulled into ribbons and may fill the entire intergranular lane. The image on the right shows a photospheric plage region. Notice the large number of phenomena showing complex structure; ribbons, 'flowers', micropores, as well as isolated and seemingly simple bright points. The magnetic field in this image is in places strong enough to perturb granulation dynamics and the granules appear 'abnormal' while displaying a slower evolution than in the quieter photosphere. [Fig. H-I:8.3]

actual magnetic field geometry, the magnetic pressure and energy density will surpass the gas pressure some 1500 km or so above the photosphere in the mid chromosphere. Another 1000 km, or 5 scale heights above the level where $\beta = 1$ (see Eq. 3.24), the plasma's ability to radiate becomes progressively worse, while the dominance of the magnetic field becomes steadily greater. As we will explain below, any given heat input in this region cannot be radiated away, and will invariably raise the gas temperature to 1 MK or greater; a corona is formed. A corona that is bound to follow the evolution of magnetic field as the field in turn is bound to photospheric driving.”

[H-I:8.2] “In Fig. 9.1 we show typical images of the quiet to semi-quiet photosphere as well as a plage region.^[xx] These images are made in the so

^{xx} 'Plage' is formally the bright chromospheric area over regions of enhanced magnetism in the photosphere, but is commonly also used to describe the interior of active regions, *i.e.*, regions of strongly

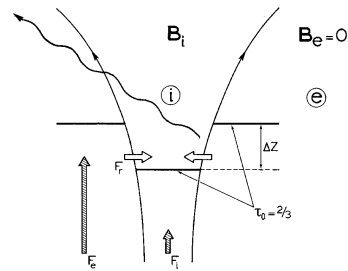


Fig. 9.2. Concept of the magnetohydrostatic flux-tube model. One level of constant optical depth in the continuum, $\tau_0 = 2/3$, is shown, with the Wilson depression Δz . The hatched arrows F_i and F_e stand for the flux densities in the (non-radiative) energy flows inside and outside the flux tubes, respectively. The horizontal arrows indicate the influx of radiation into the transparent top part of the tube. The resulting bright walls are best seen in observations toward the solar limb (as seen along the oblique wavy arrow. [From Schrijver and Zwaan (2000)]

called G-band centered around 4300 Å which is formed some 100 km above the nominal photosphere. [...] {A:^[115]}

A:115

In [the left panel of] Fig. 9.1 we show examples of simple bright points in a fairly quiet region of the photosphere, *i.e.*, a region in which the magnetic field is too weak to significantly modify granular motions. Isolated bright points are seen to be constrained to intergranular lanes and do not seem to have any internal structure on the scales that are visible on this resolution. Isolated bright points appear to be passively advected towards the periphery of supergranular cells, where they gather and form the collection of bright points that define the supergranular network mentioned above. [...]

The right panel [of Fig. 9.1] shows a region of stronger average field strength, a *plage* region, than that found in the left-hand panel. Flux concentrations with larger spatial extent are embedded in (micro-)pores with distinctly dark centers. Such small dark micropores may be the smallest manifestation of the phenomena that produce pores and sunspots. The bright points in the plage region are not simple points but are seen to have structure and appear to modify the granular flow itself: The image shows that granules near network bright points and in plage regions appear smaller, have lower contrast and in addition display slower temporal behavior than granules in weak field regions. Coalescing bright points in plage and network regions can form dark centers

¹¹⁵ Activity: The 'G band' is a narrow spectral interval centered on electronic transitions of the CH molecule, mixed in with spectral lines from multiple metals. For the interested: look up the 'Fraunhofer lines' and their designations. This old nomenclature from the days in which the solar spectrum was first studied is still used for some of these 'lines', most frequently for the Ca II H and K lines, the Na D lines, and the G band.

enhanced mean magnetic field in the photosphere, outside of, but commonly in the immediate vicinity of, sunspots.

and thus become micropores if their density is large enough, or indeed brighten again if the granular flow breaks them apart into smaller flux elements. [...]”

Substantial bundles of strong-field magnetic flux form dark pores and even larger ones form the dark sunspots. Relatively smaller bundles on the other hand are bright compared to the surrounding photosphere, particularly when seen somewhere between the center of the solar disk and the edge (or ‘solar limb’). Qualitatively, this can be understood by the combination of radiative transport and MHD. Let us start with an essentially vertical flux bundle with an intrinsic field of ≈ 1 kG as commonly observed for such solar field structures. Such a tube has settled into pressure equilibrium between inside and outside, but as the field adds its own pressure, the gas pressure inside a tube at a given height is lower than outside. {A:[116]} {A:[117]} Radiative losses lead to cooling inside and outside the tube, but the field inside lowers the ability for convection to resupply heat relative to outside, so that the interior of the tube is cooler. Being cooler and less dense, the level at which the optical depth reaches unity inside the tube is lower than outside, leading to a ‘depression’ in the solar ‘surface’ (called the Wilson depression). The line of sight into the tube from a somewhat slanted perspective looks down on the flux-tube walls, which shows layers effectively under the surface, and thus brighter than the surrounding photosphere (appearing as what are known as ‘faculae’; compare Fig. 9.2), [provided the tube is relatively narrow compared to the characteristic photon mean free path. For such a narrow tube, looking down into it shows a ‘bright point’ as we view the tube’s photosphere that lies below, and is somewhat hotter, than the surrounding photosphere.] If the tube is wide, however, such as is the case for a pore or sunspot, the sideways heat transport cannot keep the atmosphere near the wall as hot as deep into the wall, causing a cooler and thus darker wall around a dark interior, which is the manifestation of the penumbra and umbra of a sunspot. {A:[118]} {A:[119]} (A

A:116

A:117

A:118

A:119

¹¹⁶ Activity: Estimate the gas pressure contrast between inside and outside for narrow 1 kG tube in thermal equilibrium at the solar effective temperature.

¹¹⁷ Activity: Stars have a range of surface gravities, typically increasing monotonically along the main sequence towards lower effective temperatures, and substantially lower in evolved (‘giant’ and ‘supergiant’) stars than in main-sequence stars. Qualitative insight is provided by the following exercise: using the concepts of optical depth (and the fact that the stellar photosphere is around unit optical depth for continuum emission) and hydrostatic equilibrium (Ch. 2), show that the photospheric pressure would scale proportional with gravity in the idealized case of an isothermal atmosphere. In reality, radiative transport and convective motions modify that scaling for a real non-isothermal atmosphere, but the trend is in the correct direction. With this insight, argue for the trend of intrinsic field strength of photospheric magnetic concentrations with gravity: from ~ 1.4 kG in mid F-type dwarf stars to ~ 3.2 kG in late K-type dwarf stars, and well below 1 kG for cool giants.

¹¹⁸ Activity: The transition from bright to dark magnetic structures occurs at a scale of roughly 200 – 300 km. What does that say about the typical photon-mean free path ℓ_{ph} in the photosphere? Compare that value to the corresponding pressure scale height, and argue why $\ell_{\text{ph}} \gtrsim H_p$ just at the photosphere.

¹¹⁹ Activity: Explain why observed field strengths inside flux tubes exceed the equipartition field strength (field strength in an imaginary completely evacuated tube) at the level of the external photosphere.

description of numerical models of radiative magneto-convection that reveals the quantitative details is given in Sect. H-I:8.2.1). Flux tubes typically rise into the photosphere with lower intrinsic field strengths, but the process of radiative cooling and hampered internal energy transport from below then lead to a 'convective collapse' by which a more or less isolated flux tube forms (from small faculae to large sunspots) with final field strengths of order 1 – 3 kG (larger at the center of larger structures). {A:^[120]}

A:120

9.2 Heating and cooling of the solar outer atmosphere

The motions of the Sun's near-surface convection are a good fraction of the local sound speed, and thereby they generate a lot of acoustic power. Waves at frequencies below ω_a are trapped inside the Sun and can, in resonance patterns, set up one of millions of p-modes. But because the solar atmosphere above the surface has a temperature reversal into the chromosphere and corona, some degree of tunneling occurs, while higher-frequency waves can simply propagate upward through the atmosphere. All such waves steepen as they move into lower-density regions. The enhanced radiative losses during compression phases as well as dissipation through heating in the developing shocks together cause energy conversion from wave motion into thermal energy. Some of the heating of the solar chromosphere is due to such acoustic phenomena. However, the most pronounced non-radiative heating occurs at locations where magnetic field penetrates the solar surface.

Owing to the insulating atmosphere, the magnetic field of the Earth's dynamo couples to the terrestrial magnetosphere only through induction. The Sun's magnetic field, in contrast, is directly coupled from interior to atmosphere through the conduction of the plasma throughout. Consequently, the movements of the plasma in the convective envelope, from meridional circulation all the way down to sub-granular motions, drive processes in the solar atmosphere from the photosphere out to the distant heliosphere on time scales commensurate to the length scale involved, *i.e.*, from a decade down to minutes (Fig. 9.3).

The convective motions of the solar plasma, and the waves that these drive, couple into the magnetic field that threads the photosphere. The wave-like plasma motions transform into various magnetohydrodynamic wave types, while the convective motions 'braid' the higher field by the random walk of

¹²⁰ Activity: When the total solar irradiance (TSI) is smoothed over time scales of, say, a week, the Sun is brighter at sunspot maximum than at sunspot minimum, but when looking at TSI curves with a resolution of a day or so, the presence of large sunspots leads to dips when these are near the central meridian. Explain this qualitatively by the mix of faculae, pores, and spots in and around active regions. Look up TSI curves in different phases of the solar cycle.

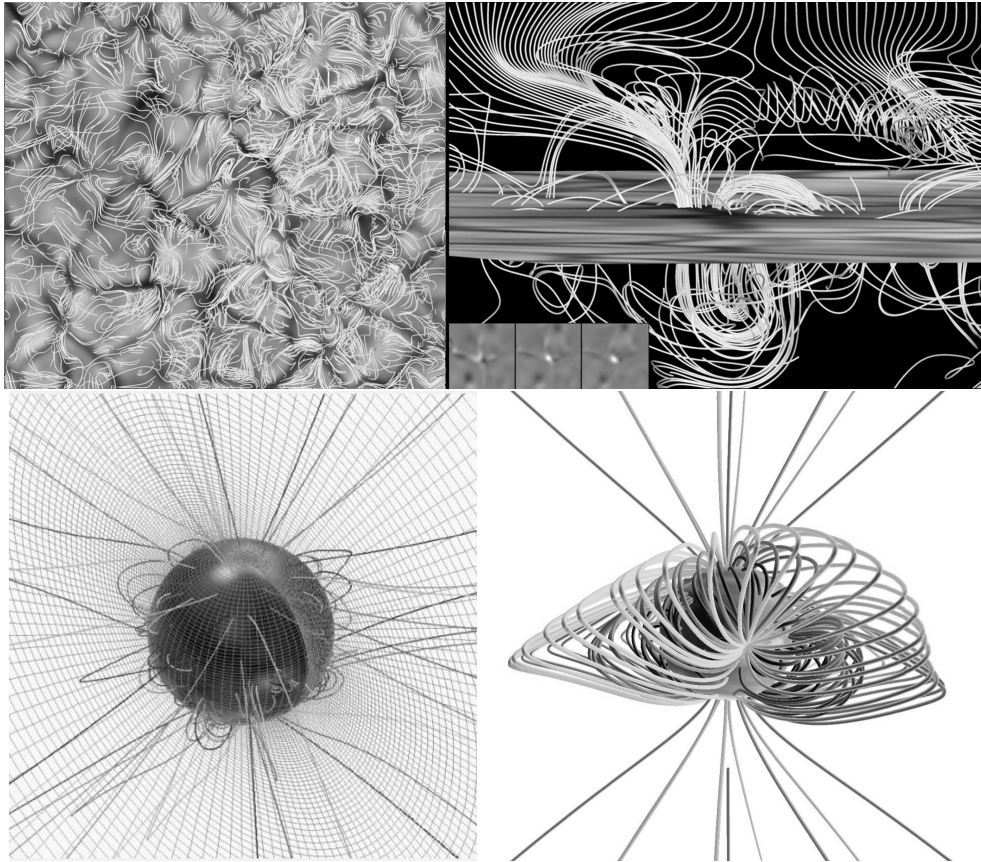


Fig. 9.3. Illustration of the multitude of scales in the solar magnetic field. The top panel shows a model computation of the magnetic field on the scale of the dominant convective motions at the solar surface, the 1000-km scale of the granulation (see also Fig. 9.1); the left panel is a top view of the solar surface with sample magnetic field lines overplotted, while the right panel shows a vertical cut through one of the convection cells to illustrate how the field in this model threads the surface sometimes multiple times, evolving on a time scale of a few minutes. The lower two panels show models of the global solar field, tracing field lines up to the cusps of the streamers that outline the topologically distinct regions of closed field and the field that is open to the heliosphere; this global scale field evolves on time scales of months to a decade. [Fig. H-I:1.2; sources for: top (from Abbett, 2007) and bottom panels.]

the footpoints 'line-tied' to the convective motions. In the higher atmosphere, these lead to wave interference and mode coupling, resonances, and field discontinuities (in a cascade of 'current sheets'). Add to that the insertion of new magnetic field emerging into the atmosphere as well as the removal from surface layers through reconnection between opposite polarities. All of these processes lead to high gradients in field and dynamics, and those to dissipation

into heat (in the literature generally differentiated into low-frequency, braiding-dominated 'DC' mechanisms in contrast to higher-frequency, wave-dominated 'AC' mechanisms). It remains an open question as to which of the proposed heating mechanisms dominates (or more comprehensively the question of which mechanism dominates under what specific conditions), but observations make it clear that the magnetic field is both the conduit and the agent involved in heating the solar atmosphere.

These processes lead to the $\sim 20,000$ K chromosphere and at sufficient height to the multi-million degree corona. In order to understand why such high temperatures arise, we need to [H-I:8.4] “realize that the temperature of a plasma is set not only by the heat dissipated but also by the plasma’s ability to lose energy. {A:[121]}

A:121

The coronal plasma has essentially three possible ways to shed energy:

- (i) Through optically thin radiation, mostly from carbon, oxygen, nitrogen, iron, and neon (and at lower temperatures from effectively thin hydrogen Lyman- α), given by

$$\Lambda(T_e) = n_e n_H f(T_e) \quad (9.1)$$

where n_e and n_H are the electron and total hydrogen densities and $f_{\text{rad}}(T_e)$ is a function of temperature dependent mainly on line emission and, at higher temperatures, on thermal bremsstrahlung. {A:[122]}

A:122

- (ii) Through thermal conduction along the magnetic field, with a conduction coefficient

$$\kappa(T_e) = -\kappa_0 T_e^{5/2} \nabla_{\parallel} T_e. \quad (9.2)$$

- (iii) The magnetically open corona can also lose energy through the acceleration of a solar wind. [...]

In short, when the plasma is dense, $n_e n_H$ is large and variations in the heat input can be dealt with by small changes in the plasma temperature which will remain on order 10^4 K or less (similar to the photospheric radiation temperature). Conduction on the other hand is very inefficient at these temperatures. However, the density drops exponentially with height, with a

¹²¹ Activity: The processes of electromagnetic radiation from a plasma involve three fundamentally distinct processes: bound-bound, free-bound (radiative recombination), and free-free (Bremsstrahlung) emission. The Sun’s coronal emission, caused by collisions of ions with thermal electrons, is dominated the first, except for flares when the last is also important; why? Which ions are typically strong contributors to the coronal X-ray and EUV emission from an active region at ~ 3 MK? Hint: combine elemental abundances with ionization energies (such as given here). For this rough estimate, ignore oscillator strengths for the transitions involved. For the solar corona under most conditions, the dominant radiative losses are from C, N, O (below about 0.5 MK), and Fe (above about 0.5 MK).

¹²² Activity: Eq. (9.1) contains a product of electron and hydrogen densities, but hydrogen is fully ionized at coronal temperatures and thus has no spectral lines that can be excited through collisions with electrons. Why is it acceptable to express it this way?

scale height of only some hundreds of kilometers for a 10^4 K plasma. The efficiency of radiative losses therefore drops very rapidly with height and *any* mechanical energy input will raise the temperature of the plasma. The temperature will continue to rise until thermal conduction can balance the energy input. Because thermal conduction varies with a high power of the temperature this does not happen until the plasma has reached 1 MK or so. Thus, we expect any and every heating mechanism to give coronal temperatures of this order [...]”

The energy equation balances the thermal energy content of the plasma with energy input by heating (ϵ_{heat}), loss by radiation, and transport through thermal conduction. Looking back at Eq. (3.6), most terms vanish in this approximation, while radiative losses, there not yet introduced, need to be added. Realizing that energy conduction occurs along the magnetic field, and with the cross section of a coronal loop is inversely proportional to the field strength ($A(s) \propto 1/B(s)$, with s the coordinate along the coronal field loop) we have:

$$\rho \frac{de}{dt} + p \nabla \cdot \mathbf{v} = \epsilon_{\text{heat}} - \frac{1}{A(s)} \frac{d}{ds} \left(A(s) \kappa \frac{dT}{ds} \right) - n_e n_H f_{\text{rad}}(T_e). \quad (9.3)$$

We can use the following rough approximations: $\kappa \approx 8 \times 10^{-7} T_e^{5/2}$ erg/cm/s/K and $f_{\text{rad}}(T_e) \approx 1.5 \times 10^{-18} T_e^{-2/3}$ erg/cm³/s for T_e in the range of 0.4 MK to 30 MK. The terms on the left of this equation are small in case the flows along the loop are sufficiently slow while any short-term variability in the heating should be relatively small compared to the internal energy content of the local plasma [for a loop in quasi-hydrostatic equilibrium]. Using that approximation along with Eq. (3.5) describes the general appearance of any slowly evolving coronal loop (but note that the radiative losses should not be approximated as above once looking into the ‘transition region’, *i.e.*, at the base of coronal loops where the temperature transitions from coronal to chromospheric values). Eq. (9.3) can be used to estimate, for example, the conductive time scale τ_{cond} (the ratio of thermal energy to the rate at which conduction over a thermal gradient transports energy) or the radiative time scale τ_{rad} (the ratio of internal energy to the time scale of radiative energy loss). The ratio of these two shows the conditions in which conduction is more important than radiation, or vice versa. This is most useful by combining these time scales with a relation that emerges from the modeling of quasi-static loops as a whole, made (by Rosner *et al.* (1978)) in the approximation of a constant cross section:

$$T_{a,6} \approx 2.8(n_{a,10} L_9)^{1/2} \quad ; \quad T_{a,6} \approx 7.3(\epsilon_{\text{heat}} L_9^2)^{2/7}, \quad (9.4)$$

for a loop apex temperature of $T_{a,6}$ MK, half length L_9 (in units of 10^9 cm) and

apex density $n_{a,10}$ (in units of 10^{10} cm^{-3}). Using the lefthand equation yields:

$$\frac{\tau_{\text{cond}}}{\tau_{\text{rad}}} = \frac{0.3}{T_{a,6}^{1/6}}. \quad (9.5)$$

This shows that for coronal loops as a whole, conduction tends to be somewhat more important than radiation in their response to energy input fluctuations, but that both need to be involved in modeling. {A:[123]}

A:123

9.3 Magnetic activity and atmospheric radiation

[H-III:2.2.3.2] “The magnetic field in the solar atmosphere is associated with the transport and dissipation of non-thermal energy; about one part in 10^4 of the Sun’s luminosity is radiated from the quiet chromosphere, and an order of magnitude less than that from the corona. For the most active stars [– or, largely equivalently, for the youngest stars (see Ch. 12) –] in contrast, a total of about 1% of the luminosity can be converted into outer-atmospheric heating. [...]

When measured for relatively large areas – *i.e.*, when averaging over an ensemble of similar atmospheric components – the radiative losses from the outer atmosphere increase with the [unsigned] magnetic flux density at the base. A variety of heating mechanisms has been proposed for the chromosphere, the corona, or – for many scenarios – both. Non-thermal energy is likely deposited into the corona in the form of electrical currents that are the result of the motion of the field’s photospheric footpoints that are moved about by convective flows. The cascade of such currents to smaller scales, and the details of the eventual dissipation continue to be debated, as is the relative importance of wave dissipation. For the chromosphere, the situation is even less clear: waves of both predominantly magnetic and predominantly acoustic nature have been proposed to play a dominant role, but numerical simulations suggest that electrical currents and reconnection phenomena contribute if not dominate. [...] With the high degree of structure in the magnetic field within the chromosphere, different mechanisms may dominate in different environments. [...]

The chromospheric and coronal emissions [scale essentially] as power laws with each other and with the average magnetic flux density of the underlying field: $F_i \propto |\mathbf{B}|^{b_i}$. The power-law index b_i between radiative and magnetic

¹²³ Activity: Use Eq. (9.4) to estimate typical volumetric heating rates for a coronal region over ‘quiet Sun’ (*i.e.*, outside of active regions; with coronal field strengths of order 20 G, loop-top temperatures of ≈ 1 MK, and loop half lengths $L \sim 4 \cdot 10^9$ cm) and for an active (sunspot-bearing) region (with coronal field strengths of order 200 G, loop-top temperatures of ≈ 3 MK, and loop half lengths $L \sim 15 \cdot 10^9$ cm). Compare these to the thermal energies also estimated from Eq. (9.4) and also compare plasma to field pressures (*i.e.*, , compute values of plasma- β).

flux densities appears to be an essentially monotonic function of the formation temperature of the radiation observed, increasing from about 0.5 for chromospheric emission from $\sim 15,000$ K plasma to just over unity for X-ray emission from ~ 3 MK plasma; these power laws hold over a contrast in X-ray surface flux densities from $100\times$ below the quiet Sun to $100\times$ above the active Sun, spanning a total of nearly five orders of magnitude (much of which will be covered by the Sun over its lifetime [... (see Ch. 12)]).

The chromospheric and coronal heating of the Sun and of stars like the Sun are a function only of the magnetic flux density [...]. In other words, once the magnetic field is in the stellar atmosphere, the dissipation of that energy and the distribution of the energy over the outer-atmospheric domains are independent of stellar properties: stars with masses from about $0.09 M_{\odot}$ (equivalent to ≈ 90 Jupiter masses) to a few solar masses, with radii of $< 0.5 R_{\odot}$ to $> 50 R_{\odot}$, and with coronal X-ray flux densities ranging over a factor of 10^5 all adhere to the same scaling relationship within the measurement uncertainties and the intrinsic stellar variability.”

10

Evolution of stars, their activity, and their asterospheres

10.1 Evolution of stars

This section summarizes stellar evolution in the context of heliophysics, *i.e.*, for stars like the Sun (defined below), stopping short of the final evolutionary phases (white dwarfs and neutron stars, as well as black holes and the supernovae on the path to their creation). [H-III:2.3.1] “Here, we introduce only some principles, terminology, and properties needed within the present context:

In the strict definition, a star is a self-gravitating body in which gravity is countered by gas pressure that is maintained by nuclear fusion balancing the loss of thermal energy through the stellar surface. Before a star forms, a contracting cloud forms opaque but still nebulous Herbig-Haro objects associated with collapsing clouds, and then pre-main sequence T Tauri stars (the subject of Ch. 11). Once a balance between contraction and internal pressure has been found, stars are on the ‘main sequence’, where they spend by far the largest fraction of their lifetime. The term main sequence refers to the well-defined clustering of stars in any one of a variety of Hertzsprung-Russell diagrams, in which the stellar luminosity or a logarithmic equivalent (the [absolute] ‘magnitude’) is plotted against the surface temperature or some filter ratio that measures the relative brightness in differently-colored filters (often the $B - V$ value is used, referring to the [absolute] Blue and Visible magnitudes, respectively). Examples of such brightness-color diagrams (often referred to as H-R diagrams) are shown in Figs. 4.2 and 10.1(left). Stars are generally characterized by their color, or an equivalent descriptor of their spectral properties called ‘spectral type’ (see the top of Fig. 4.2). [Stars on the high-temperature, left side of the HR diagram are called ‘early’ and those toward the right side ‘late’; surface temperature and mass decrease monotonically along the main sequence through the spectral type series: O, B, A, F, G, K, M.]

When stars run out of hydrogen fuel in their cores, they evolve off the main

sequence in the H-R diagram (see Fig. 10.1(left)) to become giant or supergiant stars. Their eventual fate depends on their mass: low-mass stars fade into ever-cooling white dwarfs, heavier stars eject some of their outer layers, while very heavy stars become supernovae and leave neutron stars or black holes behind. Objects that are too light to sustain hydrogen fusion during any stage of their evolution (although they may have phases with deuterium fusion) are called 'brown dwarfs', which have masses of [$\lesssim 0.07M_{\text{Sun}}$ or] $\lesssim 75M_{\text{Jupiter}}$. These cool very slowly, taking billions of years to lose their thermal energy. Even cooler objects merge into the realm of the (heavy) jovian planets.

Before stars reach the main sequence, they migrate through the H-R diagram from the top right (as red giants), initially moving down (to become red subgiants), then curving towards the main sequence (increasing their temperature to become orange, yellow, white, or even blue stars) with a much weaker change in their luminosity than during their initial contraction phase [(see Figs. 10.5 and 11.6)]. All stars cooler than a surface temperature of about 10,000 K [(roughly from spectral type late-A)] have a 'convective envelope,' or mantle, immediately below their surfaces, and the coolest stars, be they young or old, are fully convective. All of these stars make up the ensemble of cool stars, [all of which display some degree of magnetic activity.] Beneath the convective mantles, if any, lie the 'radiative interiors' in which energy is transported diffusively by photons; fusion occurs within this interior in the deep 'core' of main-sequence stars (see Fig. 4.1 for a graphic comparison – not to scale – of internal structure along the main sequence).

The evolutionary time scales are a sensitive function of mass. A star with a mass of, say, three solar masses evolves towards the main sequence in a few million years[, exhibiting magnetic activity except in the final birth phase near the main sequence.] On the main sequence, where they stay for 'only' ~ 0.4 Gyr, these stars have no magnetic activity, and they only resume magnetic activity after they evolve off the main sequence when they develop convective envelopes again for another 100 million years or so, until they rapidly evolve into what eventually [becomes a white dwarf after ejecting much of the outer layers; a star heavier than about nine solar masses ultimately] explodes as a supernova. A star of solar mass [(M_{\odot})] remains magnetically active to some degree throughout the ~ 10 Gyr that it spends on the main sequence, {A:^{[124]}}} and during the subsequent ~ 0.8 Gyr giant phase (its maximum radius may

A:124

¹²⁴ Activity: What fraction of the Sun's hydrogen would need to be converted to helium to keep it at (roughly) its current brightness throughout the time it spends on the main sequence? Once core hydrogen is consumed, the stellar internal structure changes considerably, enough to ignite fusion in higher layers as the star moves into its giant phase. Use $E = mc^2$.

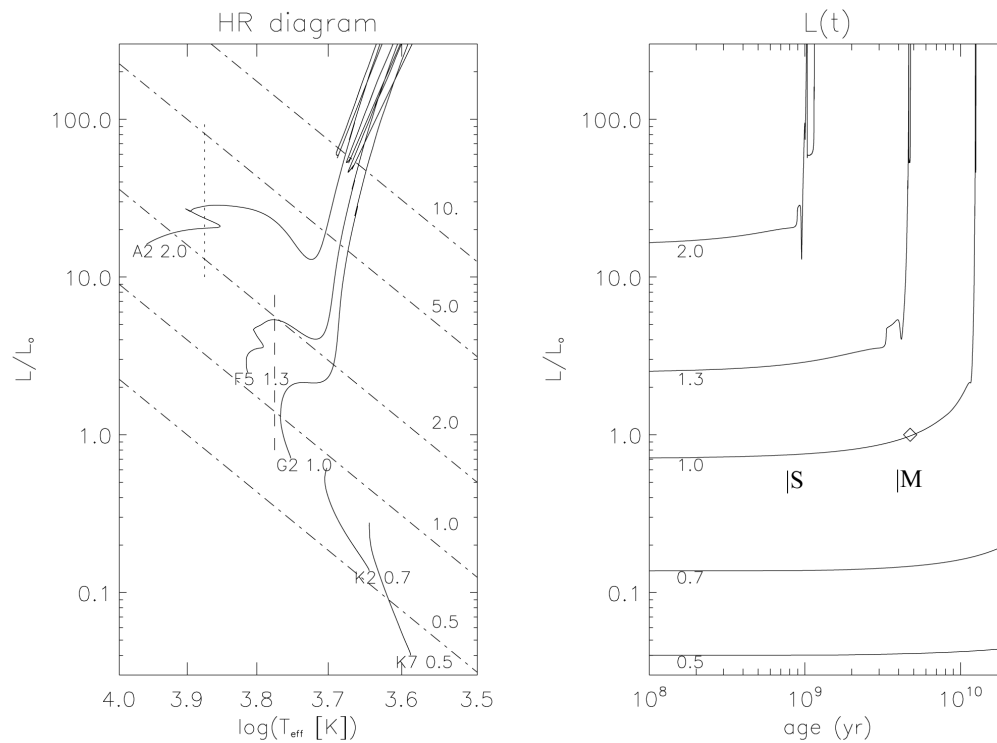


Fig. 10.1. Evolutionary diagrams for luminosity, surface temperature, and age from the mature, main-sequence phase onward. The diagram on the left relates the stellar luminosity L (in present-day solar units) with the surface effective temperature (T_{eff} ; K) in a Hertzsprung-Russel diagram (for an initial helium abundance of $Y = 0.2734$ and 'metal' abundance (everything heavier than helium) of $Z = 0.0198$). Evolutionary tracks start on the 'zero-age main sequence' (ZAMS), and are labeled with the spectral type on the main sequence and the stellar mass (in solar units). The dashed line segment indicates where dynamo action reaches its full strength; for shallower convective envelopes (warmer stars), the activity level weakens until it has dropped by a factor of 100 at the dotted line segment relative to a Sun-like star at the same angular velocity. The slanted dashed-dotted lines indicate stellar radii, with labels in solar units. The diagram on the right shows the evolution of the stellar luminosity with stellar age (yr since ZAMS). The diamond shows the present-day Sun (see Fig. 10.2 for details on the Sun's red-giant phase). The approximate ages for which the oldest fossils of single-cell microbial life (S) and multi-cellular plants and animals (M) have been found on Earth are indicated (see Ch. H-III:4). [Fig. H-III:2.9]

reach ~ 0.99 AU, and the maximum luminosity is likely to be around $5,200 L_{\odot}$) until it evolves into an ever-cooling white dwarf [after phases as giant star, ending in a series of pulses as the star gasps for fuel, during which time an appreciable fraction of its outer layers is ejected (compare Fig. 10.2)]. An

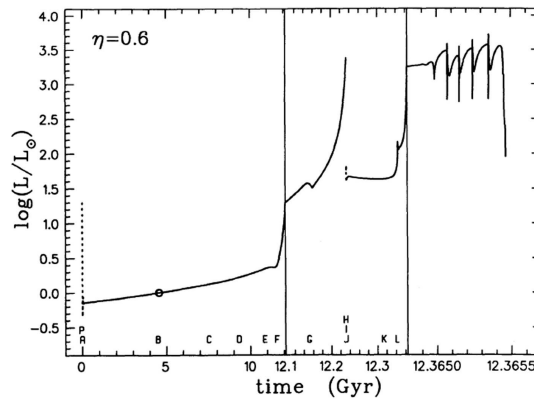


Fig. 10.2. Evolution of the luminosity of the Sun over its full life span. The first 12 billion years show the gradual brightness as hydrogen is depleted in the core of the Sun's main-sequence lifetime as a hydrogen-fusing star (cf. Fig 10.1). The large luminosity increases and pulses that follow this period include both the red-giant and asymptotic giant branch (AGB) phases when the Sun swells in size and loses appreciable mass to space. [Fig. H-III:4.5; source: Sackmann et al. (1993)].

$0.2 M_{\odot}$ M9 brown dwarf takes ~ 1 Gyr merely to contract to the main sequence, changing little in effective temperature as it descends in the H-R diagram.

During this evolution, the stellar luminosity and its associated color ([characterized by the 'effective temperature', and with it the spectral irradiance]) continually change. Examples of evolutionary changes are given in Fig. 10.1 for stars from $0.5 M_{\odot}$ to $2.0 M_{\odot}$.

The young Sun should have been some 25% fainter at the start of the Archean Eon (at $\sim 3.8 \times 10^9$ yr ago, when life is presumed to have originated) than the current mature Sun according to stellar-evolution models. This should have resulted in a much cooler Earth, covered in ice. Yet, geological records show that there was liquid water on Earth even in the first billion years of its existence. How this could happen continues to be studied. The greenhouse effect as a result of the high concentration of carbon dioxide may have compensated the lower energy input from the young Sun. Alternatively, the Sun may have been significantly more massive, and therefore brighter, early in its life; if there has been substantial mass loss in a strong wind in the first billion years, this paradox would be resolved. [Analyses suggest that it is possible] that a more massive young Sun (brighter, and with somewhat tighter planetary orbits) is compatible with the present internal structure for a Sun with a mass up to about 1.07 present solar masses. [One such model starts] with a mass of 1.07 solar masses, has an initial irradiance at Earth that is 5% higher than at present (compared to 50% lower for the Standard Solar Model),

would subsequently decrease to about 10% lower, and then increase again to the present value.”

10.2 Stellar activity and its evolution

10.2.1 Overall activity level

[H-III:2.3.2] “The defining properties of stellar magnetic activity are the existence of variable coronal (X-ray) and chromospheric (UV-optical) emissions. These characteristics are observed for a wide variety of stars [(indicated in Fig. 4.2). . . .] In single stars or in wide binaries, the activity level measured by emission from the chromospheres or coronae of these stars, or by the coverage by starspots, increases monotonically with increasing angular velocity to rotation periods as short as a few days. Rather than using the rotation period per se, however, studies of the rotation-activity relationships frequently use the Rossby number:

$$N_R = \frac{v_t}{2\Omega \sin(\theta)L_t} \sim \frac{P_{\text{rot}}}{\tau_{\text{conv}}}, \quad (10.1)$$

which is defined such that it measures the relative importance of the inertial to Coriolis forces ($\mathbf{v} \cdot \nabla \mathbf{v}$ and $\boldsymbol{\Omega} \times \mathbf{v}$, respectively) acting on a parcel of plasma of scale L_t moving with velocity \mathbf{v}_t in a rotating system with angular velocity Ω ; the Rossby number is an important metric in the theory of astrophysical dynamos, see around Eq. 4.2 and also Activity 50]. The central expression is a definition that includes the latitude, which is often neglected when estimating the global effect of rotation. When, moreover, the convective turnover time scale $\tau_{\text{conv}} = \pi L_t / v_t$ for characteristic length scales and velocities of the deepest (largest and slowest) convective motions in a stellar convective envelope is introduced, the commonly used final expression results. When using the Rossby number [estimated for the deepest layers of the convective envelope], the activity is seen to increase with rotation up to a value of $N_R \sim 0.1$ (see Fig. H-III:2.11). {A:^[125]}

A:125

¹²⁵ Activity: A cautionary intermezzo: Sect. 9.3 gives power-law scalings between radiative losses from chromospheres and coronae over stellar surface areas with mean magnetic flux densities over these areas (which hold approximately without changes for areas up to entire hemispheres). The values of the power-law indices in these relationships depend on the formation temperatures of the diagnostics used (thus, for example, steepening towards higher-energy X-ray channels), while published values also depend on the correction for a reference level (there is a minimum or ‘basal’ level of chromospheric emission that needs to be subtracted first but different authors use different corrections). This dependence on the details of the diagnostics used are one cause behind the somewhat different power-law scaling between coronal and chromospheric radiative losses you find in the literature. There are other reasons why you may find other approximate parameterizations. For one thing, although the scaling in rotation-activity diagrams between a relative brightness in terms of luminosity or surface flux density ($L_i/L_{\text{bol}} \equiv F_i/F_{\text{bol}}$) versus Rossby number works fairly well, it does not work perfectly, and other authors, using other stellar samples, might prefer using F_i versus P_{rot} . As long as the stellar sample contains stars of rather comparable internal properties, the choice of metric does not matter, but for more diverse samples, scalings with these properties matter – no simple

For even more rapidly rotating stars, the activity reaches a saturation level, and for stars with rotation periods of only a fraction of a day, supersaturation sets in, with activity decreasing with increasing angular velocity. It appears that when proceeding towards shorter rotation periods, the coronal activity saturates first, followed by chromospheric activity, and finally by starspot coverage. This has led to the suggestion that different processes set in at successively shorter rotation periods: centrifugal stripping (see Activity 13) of the high corona, saturation of the level to which non-thermal heating can be extracted from the near-surface convection or deposited into the chromosphere, and finally saturation of the dynamo process itself possibly by the coupling of the magnetic field and the plasma flows (see Section 4.5) or because the Coriolis force changes the large-scale circulation patterns that are involved in efficient dynamo action.

Main-sequence stars warmer than the Sun have shallower convective envelopes. Their magnetic activity is markedly suppressed compared to cooler stars with the same rotation period. This has been argued to be either because of the shallowness of their envelope or because of the short average turnover time of convection resulting in little influence of the Coriolis force that otherwise would introduce a preferential direction into the system. By spectral type F2 significant magnetic activity is observed, which rapidly increases in efficiency towards G0 as the convective envelope becomes deeper and the time scales of deep convective motions approach or exceed the [characteristic mean] rotation period.

Magnetic fields are observed along the main sequence as far down as we have been able to identify and apply Zeeman sensitive spectral lines, *i.e.*, down to at least M9.5. At that point we have already reached the brown dwarfs, *i.e.*, astrophysical objects that are too small to have sustained hydrogen fusion in their cores.

For stars above the main sequence, activity is seen both in stars that have recently formed and are still contracting to the main sequence (pre-main-sequence stars, which include fully convective T Tauri stars) and stars that have exhausted their core hydrogen supply and are moving away from the main sequence, once again en route to a fully-convective giant phase, now sustained by nuclear fusion of helium and heavier elements in either their core or in shells surrounding a burned out core.

multiplicative scaling seems to lead to a single tight rotation-activity relationship for all cool stars. Other reasons for differing results from different studies include the fact that the relationships are not simple power laws and fits thus depend on the parameter range covered in stellar samples, and, of course, uncertainties in models for, *e.g.*, stellar ages, and intrinsic stellar variability combined with relatively small samples. You could review, for example, this study by Booth *et al.* (2017) for more discussion and for references.

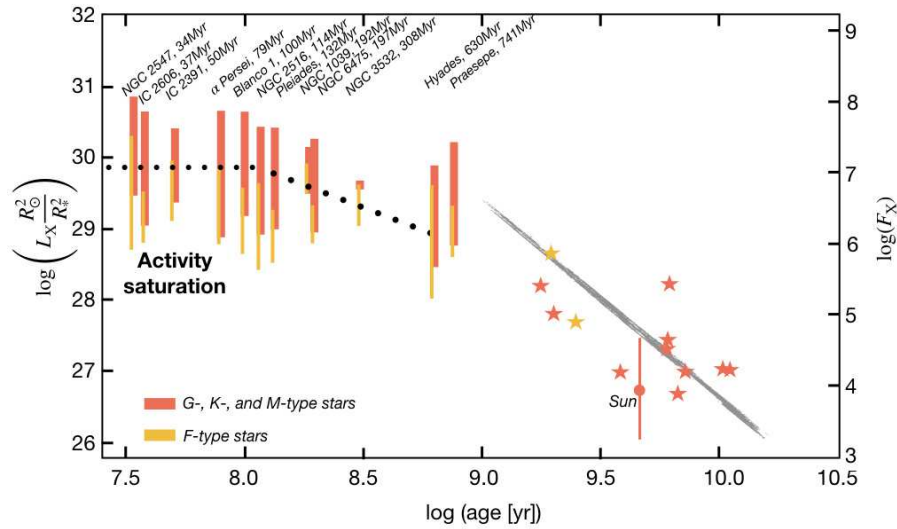


Fig. 10.3. Relationship between stellar coronal X-ray brightness (in a passband from 6 \AA to 60 \AA , expressed as a luminosity scaled to the equivalent solar surface area on the left, and as a surface flux density on the right) versus age. Bars show data for open clusters: in orange for the range of G-, K-, and M-type stars observed, and in yellow for the generally less active F-type stars. Individual stars with ages above 1 Gyr are shown by stars, and the Sun as a filled circle (with its range over the solar cycle). The relationship shown by the gray line segment is a fit by Booth et al. (2017). The dotted line is a relationship for stars with $(B - V) \in [0.56, 0.79]$. [After Booth et al. (2017) and references therein.]

During their main-sequence phases, cool stars exhibit a variety of activity patterns. A clear activity cycle, as exhibited by the Sun ever since the Maunder minimum [(a period from about 1645 CE to 1715 CE during which sunspots were very infrequent and sometimes absent for multiple years)], is relatively rare, even for solar analogs: only roughly 60% of solar-like stars show a clear activity cycle, and the reasons for this and for those that set the cycle duration are still being researched (see Fig. 10.4).

A few Sun-like stars in the solar neighborhood are so-called flat-activity stars, showing no clear cycle at all, yet they rotate with a period similar to that of the Sun. Such stars have been argued to be in a state similar to the solar Maunder minimum [...]"

10.2.2 Flares

[H-III:2.3.3] “Solar flares define power laws in spectra of frequency, N_f , versus peak brightness or overall energy, E_f . The spectrum of $N_f(E_f)$ can be approx-

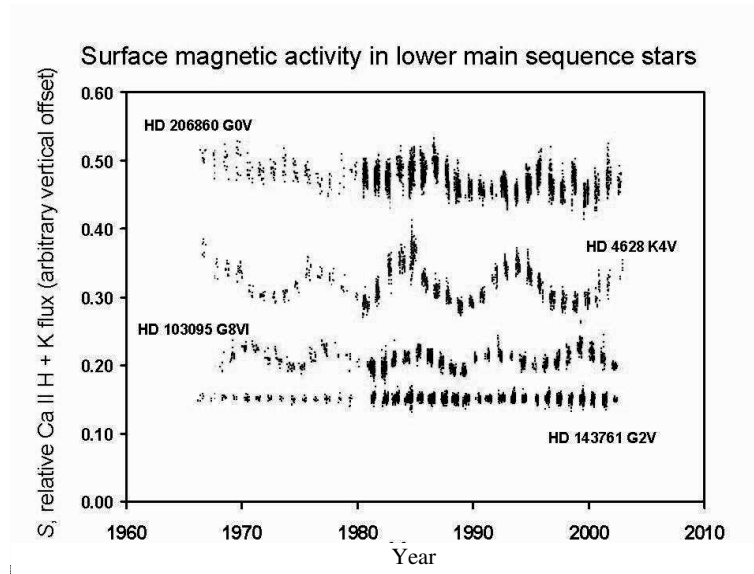


Fig. 10.4. Examples of chromospheric activity cycles (as observed in the H and K resonance lines of singly ionized Ca, or Ca II H+K). Surface magnetic activity records of four stars on or near the lower main sequence from a survey begun by O.C. Wilson in 1966 at Mount Wilson Observatory. [The ratio of the flux in the emission cores of singly-ionized calcium lines in the violet (the Fraunhofer H and K lines at 393.3 and 396.8 nm) and photospheric flux in nearby regions of the spectrum, necessarily integrated over the unresolved stellar disks, is used as a proxy for stellar magnetic activity.] The strength of the H and K fluxes increases as the coverage by and intensity of magnetic surface features increases; on the Sun the H and K fluxes vary nearly in phase with the sunspot cycle. The four records show the counterpart of the Sun approximately 2 billion years ago (upper curve, HD 206860; $P_{\text{rot}} = 4.7$ d), and then three Sun-like stars, which show records similar to the present-day Sun, HD 4628 ($P_{\text{rot}} = 38$ d), HD 103095 ($P_{\text{rot}} = 31$ d, or $P_{\text{rot}} = 60$ d) and HD 143761 ($P_{\text{rot}} = 21$ d). Both HD 4628 and HD 103095 display decadal periodicities similar to the sunspot cycle. The star HD 143761 may be in a state like the Sun's Maunder Minimum. The star HD 103095 is an extremely old (approximately 10 billion years) metal deficient subdwarf, and is shown as an example of the persistence of decadal magnetic activity cycles in a star of extreme age compared to the Sun. The spectral types are listed next to each record's star name. Arbitrary vertical shifts in the average value of the H and K relative fluxes have been applied to show the records without overlap; the offsets are 0.0 (HD 143761), 0.02 (HD 103095), 0.09 (HD 4628) and 0.15 (HD 206860). [Fig. H-III:2.12]

imated by a power law $N_f(E_f)dE_f \propto E_f^{\alpha_f}dE_f$ with $\alpha_f \approx -2$; the value of α_f reported in the literature depends on the instrument and wavelengths used, and on the sample used (active region flares, EUV quiet-Sun brightenings, etc.), and ranges from about -2.4 to -1.5 . The flare energies studied range from $\sim 10^{24}$ ergs to $\sim 10^{32}$ ergs [for the present-day Sun].

The relatively small solar flares drown into a quasi-steady background emission if the Sun is observed as a star. It is not surprising, therefore, that stellar flare spectra are limited to large flares that stand out above the surface-integrated X-ray fluxes. [O]bservations of F through M type main-sequence stars [reveal] ubiquitous power laws with power-law indices near $\alpha_f = -2$ (with a possible mild steepening from cool to warm stars). Flare X-ray [energies for some cool stars] range up to 10^{35} ergs, *i.e.*, up to ~ 1000 times brighter than the largest solar flares, with no evidence for a cutoff energy. [For the more active stars in the population,] flare frequencies for energies exceeding 10^{32} ergs scale proportionally to the time-averaged X-ray emission, saturating as the X-ray activity saturates, and contribute some 10% of the total X-ray luminosity. [... A]dopting $\alpha_f \equiv -2$, and using a characteristic solar X-ray luminosity [around cycle maximum] of $L_{X,\odot} = 3 \cdot 10^{27} \text{ erg s}^{-1}$, supports a scaling for the frequency of large flares with energy E_f exceeding a threshold value of $E_{f,32}^*$ (in units of 10^{32} ergs, characteristic of a large solar flare) of

$$N_f^*(E_f > E_{f,32}^*) \approx 0.26 \left(\frac{L_X}{L_{X,\odot}} \right)^{0.95 \pm 0.1} \left(\frac{1}{E_{f,32}^*} \right) / \text{day}. \quad (10.2)$$

[Note that this relationship derived from observations of very active stars lies some two orders of magnitude above the cycle-averaged frequency distribution observed for the time-average present-day Sun. This mismatch remains a mystery, as does the problem of establishing whether flares of $> 10^{33}$ ergs can still occur on the current Sun, or whether that was only possible in the distant past. For young, active stars we can use the above expression to find that when] the Sun was only 0.1 Gyr old [...] flares with energies exceeding 10^{35} ergs would likely have occurred once per week, and those with energies exceeding 10^{38} ergs may have occurred about once per decade. {A:[126]} A:126

It appears that quiescent activity and flaring activity on stars scale with each other, as also seen in the rise and fall of quiescent and impulsive heating through the solar cycle. One result of this is that more active stars have a stronger high-temperature coronal component, so that the effective X-ray 'color temperature' or spectral hardness increases with activity. It also appears that larger flares are associated with higher characteristic temperatures, going from solar micro-flares to large flares on very active cool stars [...]"

¹²⁶ Activity: Note that integration over the power in flares as parameterized in Eq. (10.2) diverges when the lower and upper limits extend to $[0, \infty]$. Consider what processes could be at play in introducing cutoffs to the integral on either side. The answer remains under study: it is not clear over what range Eq. (10.2) holds its slope, or what determines the energy of the 'largest flare', or how and how much relatively tiny 'nano-flares' contribute to coronal heating. But considering the possibilities should prove educational.

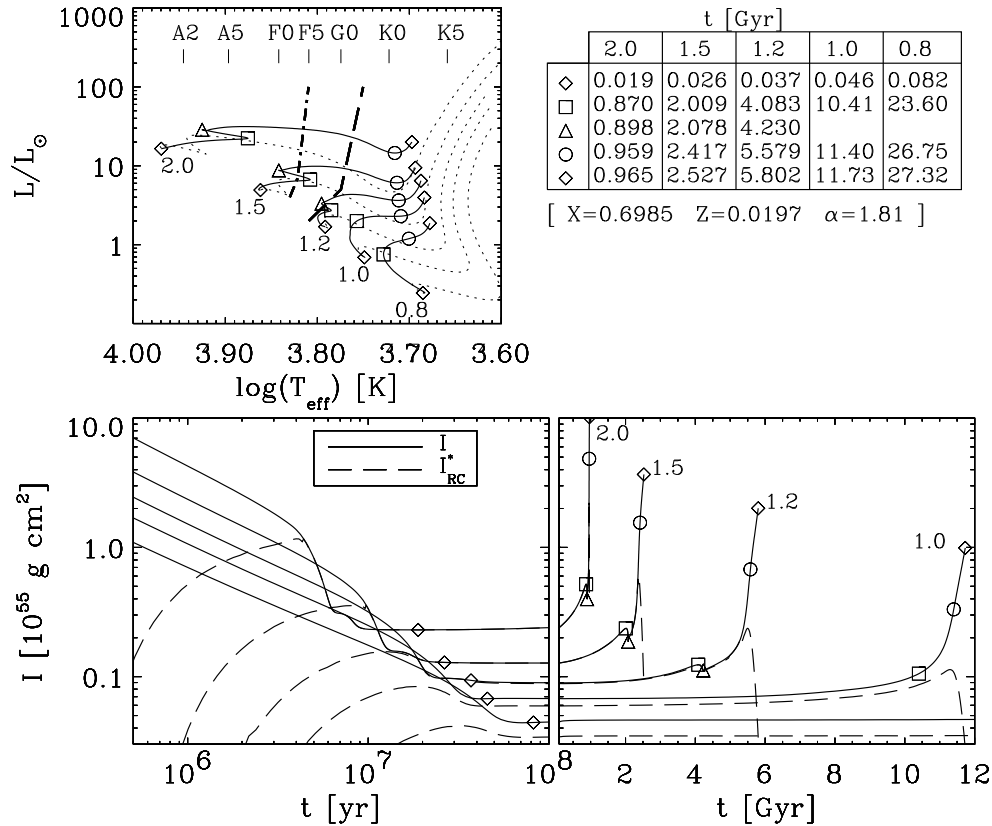


Fig. 10.5. Evolutionary tracks (top panel) for late-type stars of various masses, from the pre-main sequence to main sequence (dotted curves in the top panel), and from there to the base of the giant branch (solid curves in the top panel). The diamonds indicate the zero-age main sequence (at the lower-left end of the solid curves) and the end of model computations. Stellar masses are given in units of the solar mass. The dashed-dotted curve marks the onset of envelope convection. Ages at selected points along the tracks are listed in the table in the top right of this figure for stellar masses indicated in the top row. The evolutionary variations of moments of inertia of the entire star (solid curves) and of the radiative interior below any convective envelope (dashed curves) are shown in the lower panels. [Fig. H-III:2.15]

10.2.3 Rotation rates

[H-III:2.5] “The primary stellar property that determines the level of magnetic activity is the rate of rotation. The rate at which a star spins is influenced by the evolutionary changes in (1) the moment of inertia, (2) the angular momentum loss through a stellar wind, and (3) the angular momentum exchange in tidally interacting binaries.

- (1) The evolutionary changes in the global moment of inertia are readily

computed from stellar evolutionary models (see the example in Fig. 10.5). These changes amount to several orders of magnitude during the first tens of millions of years of a star (when magnetic coupling with surrounding accretion disks are also important, see Ch. 11) and the final fraction of a Gyr, but during the main-sequence phase, they are generally negligibly small compared to the loss of angular momentum through the outflowing wind.

(2) The outflowing stellar wind is coupled to the stellar magnetic field, which introduces a relatively long arm over which the stellar wind can extract angular momentum, so that it eventually carries far more than its own specific angular momentum [(see Sect. 7.2.1)]. The torque on the star is applied by the magnetic field into the stellar interior, and the rapid convective motions cause the angular momentum to be extracted from the entire convective envelope. How much radial and differential rotation this sets up within the convective envelope remains under active study, but the argument is generally made that the convective envelope spins down as a whole. The coupling to the radiative interior underneath it occurs somehow by coupling to a primordial field, wave exchange, or slow flows. In rapidly evolving stars with shallow convective envelopes this may lead to a (temporary) strong differential rotation between envelope and interior. For the Sun, however, helioseismic measurements have shown that interior and envelope rotate at very nearly the same rate, with the interior matching the angular velocity of the differentially rotating envelope at a latitude of about 30° (see Fig. 4.7b).

The angular momentum loss leads to a spin down relative to the evolution in which only the total moment of inertia I is evolved as the star ages. During the main-sequence phase, I changes little (Fig. 10.5) so that most of the change of P_{rot} with age [is associated with magnetic braking. $\{A:[127]\}$ $\{A:[128]\}$ A:127

After the first Gyr when dynamo saturation is important, the dependence of rotation rate on age t for Sun-like stars can be approximated by what is often referred to as the Skumanich law:] A:128

$$\Omega_{\text{rot}} \propto t^{-1/2}. \quad (10.3)$$

¹²⁷ Activity: Use Fig. 10.5 to estimate the mass of the least-massive post-main-sequence star (say, past the phase marked by a turn towards cooler surfaces marked by the squares) that could exist in the present-day Universe.

¹²⁸ Activity: Fig. 10.5 can be used to illustrate how astronomers determine the ages of 'open clusters' of stars (other than by the modern means of asteroseismology): assuming that all stars in a cluster are formed at about the same time, the shape of the HR diagram for stars in a cluster reveals the age when compared to theoretical evolutionary tracks as in the panel on the upper left. Try this: assume the stars are all 955 Myr old as in the open cluster called NGC 2355, then mark the approximate positions of the stars at that age in the upper-left panel of Fig. 10.5 estimating also where stars of intermediate masses might show up, and realize how the turnoff from the main sequence in such a cluster HR diagram reveals the age of the cluster. Open clusters all have the low-mass end of this HR diagram in common, so even if the distance to a cluster is not known, the distance can be determined by shifting that low-mass tail to overlap with that of a cluster of known distance. Also: look up the definition of 'open cluster' in contrast to a 'globular cluster'.

A:129

For the present-day Sun, the time scale of angular momentum loss is ~ 1 Gyr. [...] {A:[129]}

(3) [...] Even though the Sun is a single star, there are interesting lessons to be learned from close binaries. [...] It is the population of tidally-interacting binaries [...] that unambiguously showed us that activity is related causally to rotation, and only indirectly to stellar age [...]” the combination of angular momentum loss by magnetic activity of one or both of the binary components and their tidal interaction drains the system of angular momentum, which results in tightening of the orbits and that, in turn, to a spin-up of the stars’ rotation rates and an increase of activity level with advancing age, contrary to what happens in single stars (see Ch. 7).

At present, there are no instruments capable of detecting stellar CMEs, leaving us only – for the time being, at least – with the Sun as the guiding star, and with assumptions for scalings of CME properties for other stars with different properties and ages.

10.2.4 Stellar infancy: birth to the zero-age main sequence

[H-IV:2.2.1] “Although the age range in this first age category is only a percent or so of the total main-sequence lifetime of the star, there are several important steps to the life of the star which occur during this time. [...] For this section we concentrate on ages ranging from stellar birth to the time it takes the star to reach the zero-age main sequence, at which point the star is in stable hydrostatic equilibrium, and there is negligible contribution to the stellar luminosity from any accretion-related processes. This time scale is a function of stellar mass, being approximately 50 Myr for a solar-mass star, and longer than 160 Myr for a star of $0.5 M_{\odot}$ or less. For the purposes of discussion, and because stellar ages can be uncertain by factors of two or more, we include stars of ages up to ~ 100 Myr. [...]

Magnetic activity in general is at a high level in these young stars because of their rapid rotation, but the interpretation can be confused by other processes occurring in the system which have similar observational characteristics to magnetic reconnection processes [as discussed in Ch. 11. ...] Flares some 100-1000 times more energetic than the biggest solar flares occur roughly once a week on these young, rapidly spinning stars. [...] Over slightly longer evolutionary time scales there is a decrease in flare rate.”

¹²⁹ Activity: Start with Eq. (7.1) and note that the mass loss \dot{M} can be expressed in terms of the Alfvén speed v_A and the radial field strength B_A at the Alfvén radius. Then make the approximation that for a thermally-driven wind (in which centrifugal forces can be neglected) the Alfvén speed $v_A \approx c_s$, and take cool-star winds to all have comparable temperatures. Show that Eqs. (10.3) and (7.1) together imply that $B_A \propto \Omega$. The combination of the latter with a relationship between the photospheric field and rotation rate implicitly constrains stellar field geometries as it connects B_A and B_{surf} .

10.2.5 Stellar teenage years: ZAMS - 1 Gyr

[H-IV:2.2.2] “At this phase in a star’s evolution, rapid rotation is still an important factor, although it has declined since the star’s youth. According to [Eq. (10.3), a ‘teenage’] solar mass star would have a rotation rate that is only a factor of 2–7 above the Sun’s present-day rotation; activity that accompanies the faster rotation should be enhanced, but below the extremes represented by the youngest stars. {A:^[130]} [G- and K-type main-sequence stars spin down faster than their M-type counterparts,] so by these ages M dwarfs dominate the samples of active stars. The general decrease in activity levels compared to the extremes seen at young ages means that capturing flaring activity on stars of this age range (with the exception of M dwarfs) is more difficult to do systematically, and consequently there is a heavy bias towards the lower mass end in observations of flares on stars of this age range. The fact that M dwarfs are the most common type of star based on mass functions also contributes to this bias. There are open clusters (notably the Hyades at an age of ~ 800 Myr) which are nearby enough for sensitive studies of explosive events, although they are spatially dispersed compared to star forming regions and this makes it difficult to capture more than one or two objects in the field of view of typical astronomical telescopes. A:130

The possible dependence of stellar flare rate on evolutionary age can be explored by combining scaling relations between flare frequency and underlying coronal emission with those relating coronal and chromospheric emission, and others describing the decline of chromospheric emission with time. [The empirical scaling in Eq. (10.2)] between coronal flare rate and underlying stellar X-ray luminosity [appears to hold also for stars] with ages in this age range[. . . There are] scalings between coronal emission and different chromospheric emission indicators for cool main-sequence dwarfs, $L_X \propto L_{\text{chrom}}^y$ where $y \sim 1.5$ for C IV emission [from triply-ionized carbon typical of the transition region], $y \sim 2$ for Ca II HK emission and $y \sim 3$ for Mg II h emission[, the latter two being characteristic of singly-ionized Ca and Mg which are strong emitters from the chromosphere. . . . Chromospheric emission declines with rotation rate, which can be transformed into a relationship with stellar age using Eq. (10.3) to be roughly

$$L_{\text{chrom}} \propto t^{-1/2}. \quad (10.4)$$

Simplifying] these relations to

$$N_f(> E_{f,c}) \propto L_X, \quad (10.5)$$

$$L_X \propto L_{\text{chrom}}^y, \quad (10.6)$$

¹³⁰ Activity: Estimate the coronal soft X-ray brightness for a Sun-like star in its ‘teenage years’ (Sect. 10.2.5) relative to that of the present-day Sun.

where y takes on different values depending on the chromospheric emission being considered, and [with Eq. (10.4) ^[xxi]], suggests that the flare rate may decline with age anywhere from $N_f(> E_{f,c}) \propto t^{-0.75}$ to $N_f(> E_{f,c}) \propto t^{-1.5}$. [The] above scaling between flare rate and coronal luminosity cannot be used to 'correct' the flare rate of [young, active stars] to the solar flare rate via their coronal luminosity. This suggests a breakdown in the validity of a scaling relation approach [at ages, and commensurate rotation rates, between the 'teenage years' and the Sun's present age].

Single G stars in this age range exhibit flares at least as powerful as the largest solar flares, but occurring several times per day. [One example is κ Cet, a G5V star with an age of 300–400 Myr that exhibits 6.7 flares per day with energies of at least 10^{32} erg. The fraction of time that stars are clearly flaring in their coronal X-ray emission tends to decrease with age, from approximately 10% around 1 Myr to about 3% approaching 1 Gyr (see Fig. H-IV:2.8).]

10.2.6 Stellar adulthood: 1-5 Gyr

[H-IV:2.2.3] “The Solar System, and thus the Sun's, age measurement of 4568 Myr fits squarely within the 'stellar adulthood' phase of its life. Detections of flares on stars in this age range are much fewer. The decline of flaring with age is generally assumed to follow the trends of other activity indicators, but whether this is in fact the case is an open question. Evidence that magnetic activity may not decline monotonically at Gyr ages comes from a few sources: [...] chromospheric activity in M dwarfs [may] not decline in the 1-10 Gyr range as fast as predicted based on extrapolating from objects with ages < 1 Gyr. [For stars older than a few Gyr, it appears that there is no evidence for further decay in quiescent chromospheric activity after the major decline in activity seen] in objects at ages of the Hyades and earlier (0.6 Gyr), [while] for clusters of about 2 Gyr and older (up to 4.5 Gyr) the same activity level was seen. [...]

Because the flare rate is expected to be low on older stars, a systematic search for flares in an older stellar population needs a large number of stars, and involves a relatively long stare coupled with fast cadence to detect and resolve the flaring emission from any other variability. The *Keplers* spacecraft's

^{xxi} Note that the power laws shown in this chapter relating stellar activity, wind, and rotation are not all consistent with each other. This is not all attributable to the sensitivity to the diagnostics used (see Activity 125), which tells us something is missing in how these various parameters really scale with each other, but observations and/or theory have yet to reveal what it is that is missing. Part of the discrepancy is likely the use of different stellar samples in different studies; compare, for example, the slopes of the power-law fits in Fig. 10.3: the slope for $L_X(t > 800 \text{ Myr})$ differs depending on the age range of the stars that is included in a study. Another reason may be a change in the dependence of the loss of mass and angular momentum somewhere around the age of 1 Gyr – come back to this after reading Sect. 10.3.2.

exquisite photometry can be re-purposed from finding evidence of transiting extrasolar planets around stars to looking for rare short-timescale flaring events on the stars themselves. [Energetic flares have been found in G-type main-sequence stars, even on apparently single solar-type stars] with rotation periods of 21.8 and 25.3 days, near the solar value, and thus approximately solar age. The energetics of these flares is large, with minimum flare energies in the range 10^{33} erg, and extending up to 10^{36} erg. [...]” {A:[131]}

A:131

10.3 Evolution of astrospheres

10.3.1 Effects of a variable ISM on heliospheric structure

[H-IV:3.1] “The solar wind does not expand indefinitely. Eventually it runs into the interstellar medium (ISM), the extremely low particle density environment that exists in between the stars [(*cf.* Fig. 10.6)]. Our Sun is moving relative to the ISM that surrounds the planetary system, so we see a flow of interstellar matter in the heliocentric rest frame, coming from the direction of the constellation Ophiuchus. The interaction between the solar wind and the ISM flow determines the large scale structure of our heliosphere, which basically defines the extent of the solar wind’s reach into our Galactic environment. Other stars are naturally surrounded by their own ‘astrospheres’ (alternatively ‘asterospheres’) defined by the strength of their stellar winds, the nature of the ISM in their Galactic neighborhoods, and their relative motion.”

[H-IV:3.5] “The Sun is now traveling through the ISM at a rate of 16–20 [parsecs (or pc, a unit of 3.26 light years)] per million years (Myr) compared to the average motion of nearby stars about the Galactic Center. {A:[132]}

The ISM has densities ranging from 10^4 cm^{-3} or higher in dense molecular clouds down to about 0.005 cm^{-3} in very low-density hot gas regions. Because the heliosphere will contract or expand by large factors when the Sun enters such high- or low-density regions, it is important to investigate when such environmentally driven changes could have occurred and will possibly occur by considering the Sun’s historic and future path through the ISM.

A:132

At present, the heliosphere resides inside of the partially ionized [local interstellar cloud (LIC)], with properties likely similar to other warm partially ionized clouds within 15 pc of the Sun. The Sun likely entered this cluster of local warm clouds about 1 Myr ago. However, on a larger scale, the Sun

¹³¹ Activity: The minimum flare energy given here is instrumental, not intrinsic. Argue why the empirical lower limit of flares detectable by an instrument like *Kepler* is limited to of order 10^{33} ergs. Note that this lower limit exceeds the energies observed (to date, at least) for solar flares.

¹³² Activity: Just to get an impression of relative velocities: compare the average speed of the Solar System relative to the local ISM to the speed of 828,000 km/h with which the Solar System orbits the Galactic center.

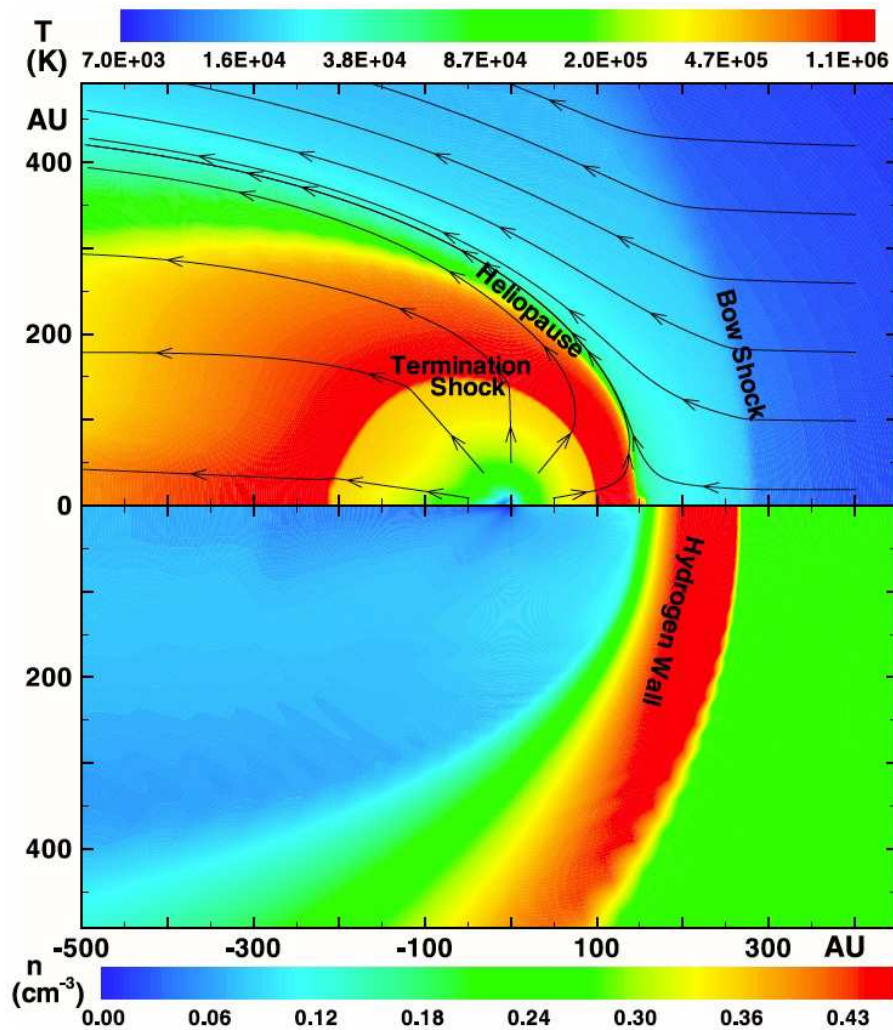


Fig. 10.6. A 2.5D axisymmetric, hydrodynamic model of the heliosphere and the surrounding ISM. The upper panel shows plasma temperature and ISM streamlines, and the bottom panel shows neutral hydrogen density. [Note: the effective solar wind plasma temperature, increasing with heliocentric distance, is dominated by the energy density of pickup ions from several AU outward to the termination shock. Fig. H-IV:3.6; source: Müller and Zank (2004).]

actually lies in a region called the Local Cavity, or Local Bubble, which is ~ 200 pc across and is filled mostly with fully ionized, low-density ISM. [No evidence has been found that the Sun has traveled through significantly denser regions over the last 30 Myr (500 pc) until about 7.5 Myr ago (120 pc) when the Sun was at the edge of the Local Cavity. It appears that the Sun will]

leave the LIC in less than 3000 yr. What will be the properties of this new environment? [...]

Our ideas concerning the properties of the gas located between the warm local ISM clouds have undergone a radical change in the last 20 years. The gas between the clouds, extending out to roughly 100 pc from the Sun in what is now called the Local Cavity was originally assumed to be hot (roughly 10^6 K), fully ionized, and low density (roughly 0.005 cm^{-3}). This conclusion was based upon the predictions of the classical models and observations of diffuse soft X-ray emission consistent with the properties of the hot gas. This picture has since been complicated by the realization that X-ray emission from charge-exchange (CX) reactions between the solar wind ions and inflowing interstellar neutral hydrogen can explain much of the observed diffuse X-ray emission, except for the Galactic pole regions. [It has instead been] proposed that the Local Cavity is an old supernova remnant with photo-ionized gas at a temperature of about 20,000 K. The likely photo-ionizing sources are the hot stars ϵ CMa and β CMa and nearby hot white dwarfs like Sirius B. {A:[133]}
{A:[134]}

A:133

A:134

How will the heliosphere change as the Sun passes through very different regions of the interstellar medium? [...] Figure 10.7 compares today's heliosphere properties with the Sun located inside of the partially ionized warm LIC to a model computed for the Sun surrounded by 10^6 K fully ionized interstellar plasma. The main difference between these models is that the hydrogen wall does not exist when the inflowing interstellar gas contains no neutral hydrogen atoms. The locations of the termination shock (TS), heliopause (HS), and bow shock (BS) are determined by pressure balance between the solar wind ram pressure and the thermal and ram pressure of the surrounding interstellar gas. In this comparison, the locations of the TS, HP, and BS are about the same in the two models because the high temperature and low density of the interstellar gas produce a pressure that is about the same as in the LIC.

When the Sun enters a region of much higher density or speed, and therefore higher ram pressure, the effect is to compress the heliosphere. For example, a model for $n_{\text{HI}} = 15 \text{ cm}^{-3}$, roughly 100 times that of the LIC, has a TS at 9.8 AU such that Uranus would move in and out of the TS and Neptune would be surrounded by hot, shocked plasma beyond the HP (upwind) or heliotail

¹³³ Activity: With average values for solar wind density and velocity (assuming a radial outflow at constant velocity and with a density as specified in Table 2.4), at what distance from the Sun does the solar wind dynamic pressure equal the interstellar total pressure for estimated values of $B_{\text{LISM}} \approx 3 \mu\text{G}$, $T_{\text{LISM}} \approx 6500 \text{ K}$, and $n_{\text{p,LISM}} \approx 0.06 \text{ cm}^{-3}$ and $n_{\text{H,LISM}} \approx 0.18 \text{ cm}^{-3}$ (see, e.g., Sect. H-IV:3.2)?

¹³⁴ Activity: Given present-day parameters for the ISM as in Activity 133, where would the heliopause be, very approximately, for the range of ISM densities given in Sect. 10.3.1, assuming a present-day spherically-symmetric, constant-velocity solar wind, and the same temperature for the ISM? Compare your result with Fig. 10.7. Look back to Sect. 5.5.8 for some of the physics involved.

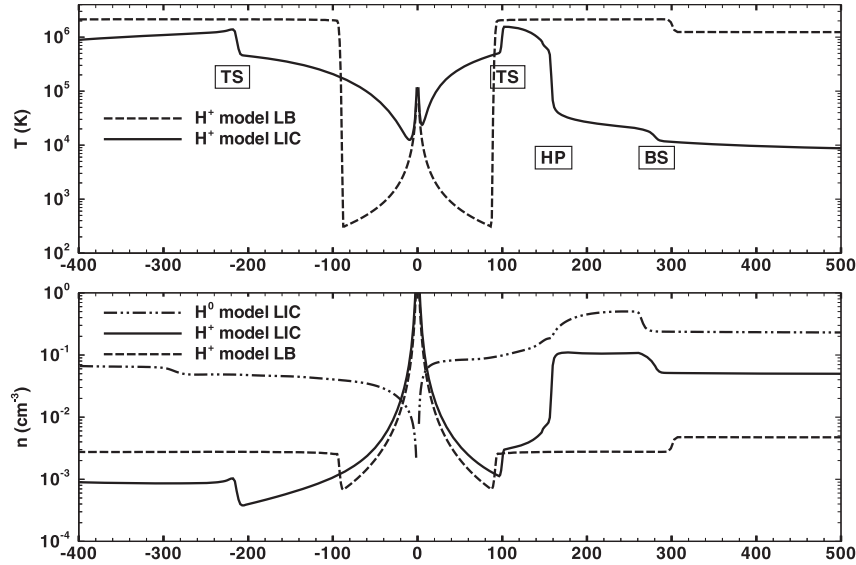


Fig. 10.7. Top: Plots of the temperature vs. distance in Sun-Earth distances (astronomical units, or AU) relative to the Sun (interstellar flow upwind direction to the right and downwind to the left) for a heliosphere model with the Sun located inside of the local interstellar cloud (LIC; solid line) or inside a 10^6 K hot interstellar medium (dashed line, LB). The heliosphere in the LIC model has a termination shock (TS), heliopause (HP) and bow shock (BS) structure. Bottom: Density structures for the LIC neutral hydrogen (solid line), LIC protons (dot-dash line), and hot [Local Bubble (LB)] interstellar model protons (dashed line). Note that the hydrogen wall at 150–280 AU exists when the heliosphere is located inside partially neutral interstellar gas but not when it is inside fully ionized interstellar gas. [Fig. H-IV:3.8; source: Müller et al. (2009).]

(downwind). Models of the heliosphere inside of a high-speed interstellar wind with corresponding high ram pressure would compress the heliosphere in a similar way. [A potential] cloud encounter that results in a stellar astrosphere being compressed to less than the size of the star's habitable zone [... has been described as a 'de-screening event'. This] should happen when a star encounters an interstellar cloud with a number density of $600(M_{\odot}/M)^2 \text{ cm}^{-3}$, where M is the mass of the star. Only the densest ISM clouds are capable of this de-screening, with such clouds being relatively rare. The densest clouds are cold ($T \sim 100$ K) molecular clouds, with many of the refractory elements depleted onto dust grains. In addition to increased GCR exposure

(see Ch. 14), a de-screening event caused by a molecular cloud encounter would also expose planetary atmospheres to high fluxes of interstellar dust, with potentially dramatic consequences [that include potential 'snowball Earth' states for climate. Given what we know about mass-loss and wind properties of stars (as discussed below), it appears] that habitable zone planets orbiting stars significantly less massive than the Sun (with spectral types of late K to M) are virtually never exposed to de-screening events, but de-screening may happen occasionally for stars with the Sun's mass or larger. However, these calculations assumed that the relative velocity of these encounters is only 10 km s^{-1} . Assuming a faster encounter speed would increase the estimated frequency of de-screening events."

10.3.2 Long-term evolution of stellar winds

In Sect. 5.5.8 we described how neutrals moving toward to the heliosphere leads to a 'hydrogen wall' outside of the heliopause through charge-exchange collisions in that region (see Figs. 5.1 and 10.6). {A:^[135]} [H-IV:3.4] "The importance of this hydrogen wall is that it is actually detectable in UV spectra from [the *Hubble Space Telescope (HST)*], not only around the Sun but around other stars as well. {A:^[136]}"

The effect of heliospheric and astrospheric absorption on stellar H I Lyman- α spectra [(emitted in the stellar atmosphere by de-excitation from the first excited to the ground state in neutral hydrogen atoms and absorbed en route to Earth by the inverse process)] is described by Fig. 10.8, showing the journey of Lyman- α photons from the star to the observer. Most of the absorption

¹³⁵ Activity: For charge exchange only, and assuming (very approximately, as done initially (Holzer, 1972) decades ago) a velocity-independent cross section for resonant-charge exchange of solar wind protons with ISM neutral hydrogen of $\sigma_{CX} \approx 2 \cdot 10^{-15} \text{ cm}^2$, what fraction of H^0 , looking at the population after passing through the 'hydrogen wall' and moving in a straight line towards the Sun, would reach Earth orbit for present-day slow wind conditions? In reality, other processes are major players: radiation pressure (for neutral hydrogen primarily by repeated Lyman α absorption followed by isotropic re-emission) pushes outward on the atoms, and photo-ionization in the Sun's EUV and X-rays presents a significant loss term. It appears that Lyman α radiation pressure on ISM H^0 just balances solar gravity, see Schwadron *et al.* (2013); for a significantly younger Sun, ISM H^0 would never reach Earth orbit. The combined effects of these processes would render an IBEX-like mission to learn about the ISM H^0 around a young Sun pointless except during times of passage through dense interstellar clouds.

¹³⁶ Activity: Argue why the heliospheric hydrogen wall has a thickness L_{HW} that is, within a factor of a few, comparable to, but less than, the distance d_{HP} from the Sun to the heliopause. For a simple estimate, use a circular 'cookie tin' geometry to approximate conditions at the heliosphere's 'nose', with the incoming flow through the top being decelerated by the gas pressure (ignore magnetic effects here), and accelerated sideways out by the same pressure, combining scale estimates based on the continuity and momentum equations; focus only on the flow into the heliopause and assume no bow shock (see Sect. 5.5.8). Check that this gives it just enough of a total column depth with the charge-exchange cross section from Activity 135 so that a useful fraction of ISM neutral hydrogen can indeed be made part of the flow in heliosheath. See this study by Wood *et al.* (2002) for simulated astrospheres and their hydrogen walls, and some images for different stars and their speeds through the ISM. How would the thickness of the hydrogen wall change for a much higher speed of a planetary system through its LISM?

is by interstellar gas in the line of sight from the star to the Sun, but the astrosphere and heliosphere provide additional absorption on the left and right sides of the interstellar absorption, respectively. The effect of the hydrogen wall around the Sun is to provide additional red-shifted absorption on the right side of the interstellar absorption feature because the neutral hydrogen gas in the solar hydrogen wall is slowed down and deflected relative to the inflowing interstellar gas. Conversely, the absorption by the hydrogen wall gas around the star is seen as blue-shifted relative to the interstellar flow from our perspective outside the astrosphere, and is therefore seen on the left side of the absorption line. [... By way of an example observation, the] bottom panel of Fig. 10.8 shows the HI (and [equivalent deuterium] DI) Lyman- α spectrum of the lower-brightness component of [the star] α Cen B. Most of the intervening HI and DI between us and the star is interstellar, but the ISM cannot account for all of the HI absorption. As mentioned above, the red-shifted excess on the right side is heliospheric and the blue-shifted excess on the left is astrospheric.”

[H-IV:3.7.2] “Currently, the only way coronal winds can be detected around other stars is through astrospheric Lyman- α absorption, but the number of astrospheric Lyman- α detections is still very limited. [... These measurements need to be interpreted in terms of models that] are extrapolated from a heliospheric model that successfully reproduces heliospheric absorption, specifically a multi-fluid model. These models assume the same ISM characteristics as the heliospheric model, with the exception of the ISM flow speed in the stellar rest frame, v_{ISM} , which can be computed using our knowledge of the local ISM flow vector and each star’s unique space motion vector. [...]

The astrospheric models are computed assuming different stellar wind densities, corresponding to different mass-loss rates, and the Lyman- α absorption predicted by these models is compared with the data to see which best matches the observed astrospheric absorption. [...] In order to look for some correlation between coronal activity and wind strength, Fig. 10.9 shows mass-loss rates (per unit surface area) plotted versus F_X [(the ratio of X-ray luminosity to surface area)], focusing only on the main-sequence stars. For the low-activity stars, mass loss increases with activity in a manner consistent with the $\dot{M} \propto F_X^{1.34 \pm 0.18}$ power-law relation shown in the figure. For the ξ Boo binary, in which (like α Cen) the two members of the binary share the same astrosphere, Fig. 10.9 indicates how the binary’s combined wind strength of $\dot{M} = 5\dot{M}_{\odot}$ is most consistent with the other measurements if 90% of the wind is ascribed to ξ Boo B, and only 10% to ξ Boo A.

For $F_X < 10^6$ erg cm $^{-2}$ s $^{-1}$, mass loss appears to increase with activity.

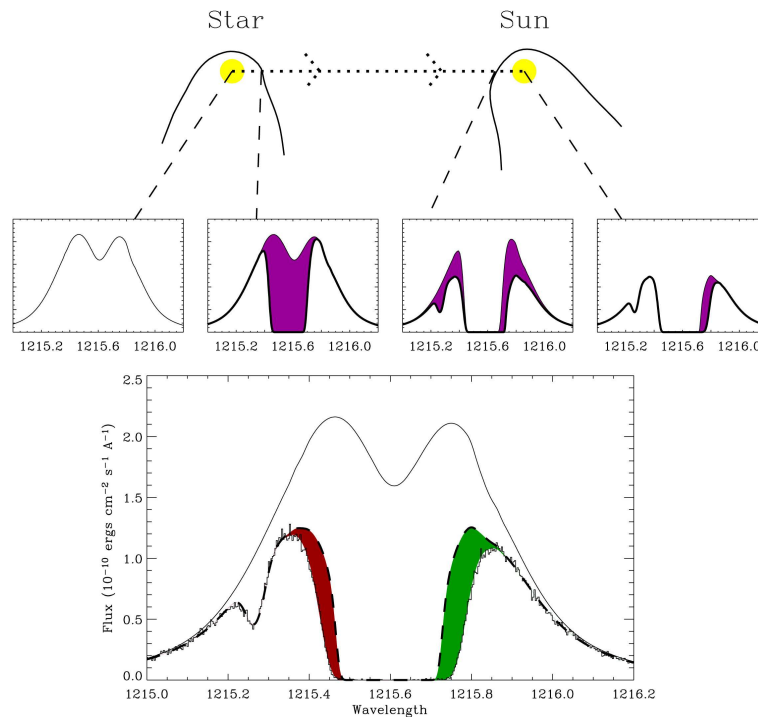


Fig. 10.8. Top panel: The journey of a Lyman- α photon from a star through its astrosphere, the interstellar medium, and the heliosphere. Middle panels from left to right: The Lyman- α emission line emitted by the star, absorption due to the stellar astrosphere, additional absorption due to the interstellar medium, and additional absorption due to the heliosphere. Bottom panel: HST Lyman- α spectrum of α Cen B, showing broad HI absorption at 1215.6 \AA and DI absorption at 1215.25 \AA . The upper solid line is the assumed stellar emission profile and the dashed line is the ISM absorption alone. The excess absorption is due to heliospheric HI (green shading, vertical lines) and astrospheric HI (red shading, horizontal lines). [Fig. H-IV:3.7; source: Wood (2004).]

{A:¹³⁷} However, above $F_X = 10^6 \text{ erg cm}^{-2} \text{ s}^{-1}$ (i.e., for more active, and thus generally younger stars) this relation seems to fail, a boundary identified as the 'Wind Dividing Line' in Fig. 10.9. Highly active stars above this limit appear to have surprisingly weak winds. This is suggested not only by the two solar-like G stars above the limit, ξ Boo A and π^1 UMa, but also by the two active M dwarfs above the limit, which have very modest mass-loss rates. (For Proxima Cen we only have an upper limit of $\dot{M} < 0.2\dot{M}_\odot$, while for EV Lac $\dot{M} = 1\dot{M}_\odot$.) The apparent failure of the wind/corona correlation to the right

¹³⁷ Activity: Show that the total power lost in X-rays from the present-day solar corona (estimated from Fig. 10.3 or 10.9) is roughly twice the total power lost in the solar wind (using the expressions in Sec. 3.5.2), and that these numbers would have been comparable for the young Sun at the 'wind dividing line' if the characteristic wind speed would have been the same.

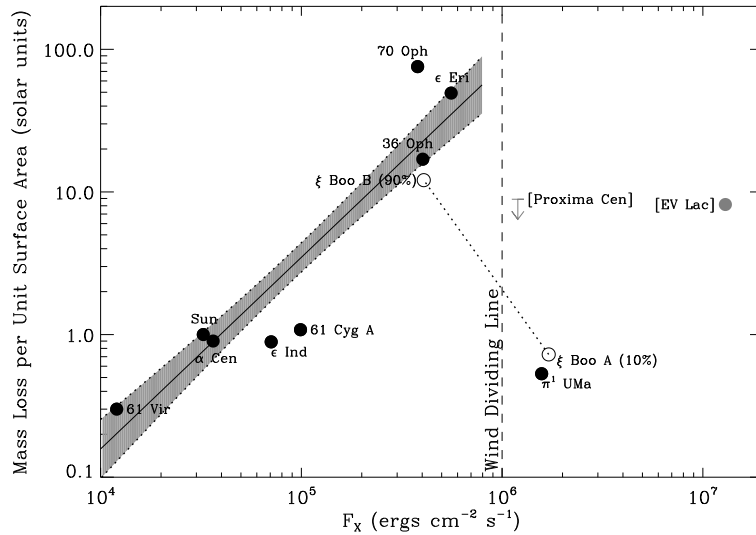


Fig. 10.9. A plot of mass-loss rate (per unit surface area) versus X-ray surface flux density for all main-sequence stars with measured winds. Most of these have spectral types of G (like the Sun) or (cooler) K, but the two with square-bracketed labels are (much cooler) tiny M dwarf stars. Separate points are plotted for the two members of the ξ Boo binary, assuming ξ Boo B accounts for 90% of the binary's wind, and ξ Boo A only accounts for 10%. A power law, $\dot{M} \propto F_X^{1.34 \pm 0.18}$, is fitted to the less active stars where a wind/corona relation seems to exist, but this relation seems to fail for stars to the right of the 'Wind Dividing Line' in the figure. [Fig. H-IV:3.12; source: Wood et al. (2014).]

of the 'Wind Dividing Line' may indicate a fundamental change in magnetic field topology at that stellar activity level.

[...] Sophisticated spectroscopic and polarimetric techniques are also available for studying stellar [surface] magnetic fields. One interesting discovery is that very active stars usually have stable, long-lived polar starspots, in contrast to the solar example where sunspots are only observed at low latitudes. Perhaps the polar spots are indicative of a particularly strong dipolar magnetic field that envelopes the entire star and inhibits stellar wind flow, thereby explaining why very active stars have surprisingly weak winds. Strong toroidal fields are also often observed for active stars. {A:^{[138]}}}

A:138

Given that young stars are more active than old stars, the correlation between mass loss and activity indicated in Fig. 10.9 implies an anti-correlation of mass loss with age. [One parameterization of this is given by] $F_X \propto t^{-1.7 \pm 0.3}$ [$\Omega^{-3.4 \pm 0.6}$]. Combining this with the power-law relation from Fig. 10.9 yields

¹³⁸ Activity: Place the coronal activity level corresponding to the 'wind dividing line' in the rotation-age diagram in Fig. 10.3, and consider possible consequences for that diagram.

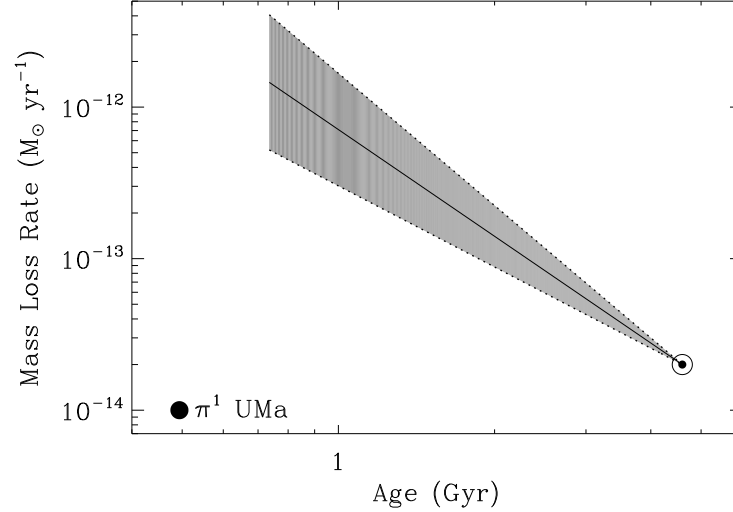


Fig. 10.10. The mass-loss history of the Sun inferred from the power-law relation in Fig. 10.9. The truncation of the relation in Fig. 10.9 means that the mass-loss/age relation is truncated as well. The low mass-loss measurement for π^1 UMa suggests that the wind weakens at $t \approx 0.7$ Gyr as one goes back in time. [Fig. H-IV:3.13; source: Wood *et al.* (2005).]

the following relation between mass-loss rate and age:

$$\dot{M} \propto t^{-2.3 \pm 0.6} [\propto \Omega^{-4.6 \pm 1.2}] \quad (10.7)$$

[where the final expression above between brackets links to the intrinsic dependence on rotation rate via Eq. (10.3).] Fig. 10.10 shows what this relation suggests for the history of the solar wind, and for the history of winds from any solar-like star for that matter. The truncation of the power-law relation in Fig. 10.10 near $F_X = 10^6$ erg cm $^{-2}$ s $^{-1}$ leads to the mass-loss/age relation in Fig. 10.10 being truncated as well at about $t = 0.7$ Gyr. The plotted location of π^1 UMa in Fig. 10.10 indicates what the solar wind may have been like at times earlier than $t = 0.7$ Gyr.

[Fig. 10.10 indicates] that solar-like coronal winds can be up to two orders of magnitude stronger than the current solar wind at $t \approx 1$ Gyr. This makes it more likely that the erosive effects of stellar winds play an important role in planetary atmosphere evolution at these later ages” (see Ch. 12). {A:^[139]} A:139

¹³⁹ Activity: Estimate the size of the heliosphere and the terrestrial magnetopause distance for a young Sun at an age of 700 Myr, assuming unchanged LISM conditions and geomagnetic properties.

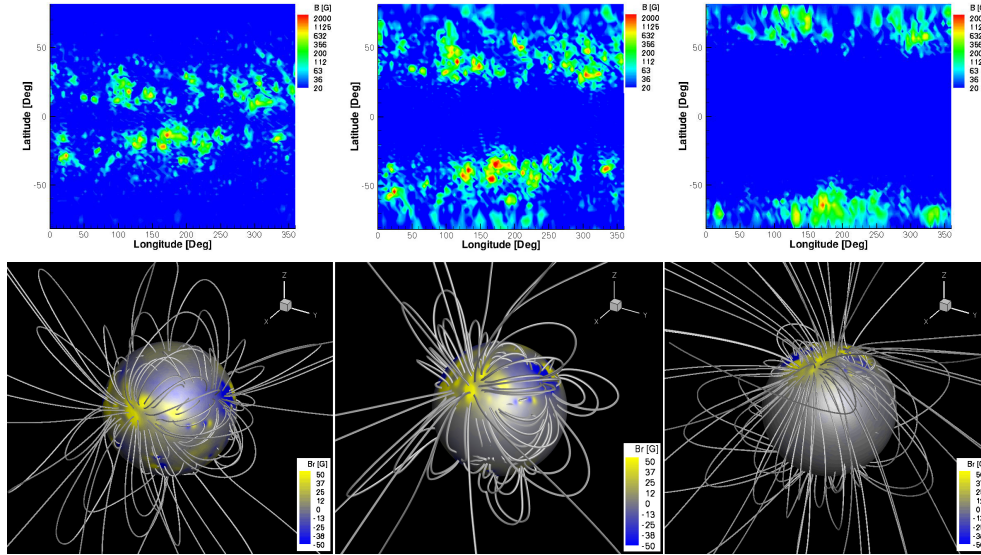


Fig. 10.11. Top row: A map of the solar photospheric radial magnetic field (magnetogram) during Carrington Rotation 1958 (January 2000, solar maximum period) shown on the left. The middle and right panels show manipulation of the original map, where the active regions have been shifted by 30 and 60 degrees toward the poles, respectively. [Fig. H-IV:4.4; source: Cohen et al. (2012)] Bottom row: The three-dimensional magnetic field corresponding to the surface distribution of the photospheric radial magnetic field (shown on a sphere of $r = R_{\odot}$) during solar maximum (left), and for manipulated photospheric field with the active regions shifted by 30 (middle) and 60 (right) degrees toward the poles, as shown in the row above. [Fig. H-IV:4.5]

10.3.3 Astrospheric field patterns in time

[H-IV:4.1.1] “The extent and structure of astrospheres is determined by the radially-expanding super-Alfvénic stellar wind that drags the stellar magnetic field from the stellar corona through interplanetary space, until the wind is stopped by the Interstellar Medium (ISM). It is also determined by the rotation of the star. As a result, each astrospheric magnetic field (AMF) line has one end (or ‘footpoint’) attached to the stellar surface, while its location at each point in the astrosphere, $\mathbf{r}(r, \theta, \phi)$ (for co-latitude θ), is given by the following formula. It describes a spiral shape and is known as the ‘Parker Spiral’ [(Sect. 5.4, compare with Eq. 5.20)]:

$$\mathbf{B}(r) = B_0 \left(\frac{r_0}{r} \right)^2 \left[\hat{\mathbf{e}}_r + \frac{(r - r_0)\Omega \sin \theta}{v_w} \hat{\mathbf{e}}_{\phi} \right]. \quad (10.8)$$

Here Ω is the stellar rotation rate (angular velocity), v_w is stellar wind speed (which is here assumed to be radial and fixed in time and space); r_0 is the

actual base point of the AMF, and is at a reference distance from the stellar surface at which we assume the stellar wind is fully developed and has achieved its asymptotic speed and radial direction; B_0 is the magnetic field magnitude at that point. We can see that the radial component of the AMF has an r^{-2} dependence, while the azimuthal component has only a r^{-1} dependence. As a result, throughout most of the astrospheres, the AMF is dominated by the azimuthal field, which is a function of Ω , except for high latitudes (small θ) where the AMF lines are nearly radial.

Over time, stellar rotation periods evolve from less than one day for very active, young stars to about 20-100 days for older, main-sequence stars like the Sun. For very fast-rotating stars, the AMF spiral is completely dominated by the azimuthal component: the field is highly compressed, and its azimuthal component dominates even at relatively small distances from the star and inside the stellar corona, which typically extends to 10-20 stellar radii. In this case, even extended closed magnetic loops can be bent as a result of the fast rotation. This effect can have implications for the triggering of very strong stellar flares, and for the mass-loss rate of the star to the stellar wind. The right panel in Figure 5.5 shows how the compression of the AMF spiral changes for different stellar rotation periods. The other two panels show the AMF lines close to the star (up to 24 stellar radii). It can be clearly seen that the field lines are nearly radial for the slow, solar-like rotation period of 25 days, while the field lines are strongly bent in the azimuthal direction for a fast rotation period of half a day. {A:[140]}

A:140

Equation (10.8) describes how a given magnetic field line changes with distance for a given value of B_0 at its base (r_0), and a given asymptotic stellar wind speed v_w . However, the AMF is formed by a collection of field lines that are defined by some spherical distribution of B_0 at the base of the stellar corona. This distribution depends on the topology of the stellar magnetic field at a given time. In addition, the value of v also varies as it empirically depends on the expansion of the magnetic flux tubes and on the non-uniform distribution of B_0 . [...]

Over time, stellar activity appears at different latitudes, while changing in magnitude as the behavior of surface magnetic activity is highly tied to the rotation rate. Young active stars seem to have very strong large-scale magnetic fields with magnitude of several kilo-Gauss. For reference, the Sun's dipole field strength is of the order of 5-10 G, and while the magnetic flux density within active regions can be high (ranging up to well over a kiloGauss in sunspots), solar active regions are rather small in size. In addition, magnetic activity in active stars tends to appear at high-latitude, polar regions. This behavior is

¹⁴⁰ Activity: Section 10.3.3 mentions a solar rotation period of 25 d while the caption to Fig. 5.6 mentions 27 d. What is the reason for using these two different values in the different contexts?

most likely related to the role of the fast stellar rotation in the stellar dynamo and meridional magnetic flux circulation. [...]

The appearance of stellar activity described above reflects a change in the distribution of B_0 . Therefore, it affects the shape of the AMF and the astrospheric volume. It is not clear how v_w changes for young stars as we cannot directly measure stellar winds of 'cool stars', *i.e.*, stars with a convective envelope beneath their surfaces such as in the case of the Sun. Some techniques to estimate mass-loss rates from cool stars are [outlined above]. However, these estimates do not separate the stellar wind speed from the density, so it cannot be obtained independently. Another cause for the lack of estimates for stellar wind speeds of cool stars is the incomplete theory about the solar wind acceleration. In order to demonstrate how the change in the photospheric field affects the three-dimensional structure, Figure 10.11 shows the distribution of the photospheric magnetic field and the shape of the three-dimensional magnetic field close to the Sun. The top-left panel is obtained using actual data of the photospheric field during a period of high solar activity. In the other two panels in the top row, the original data was manipulated, so that the active regions have been shifted by 30 and 60 degrees, respectively, towards higher latitudes to mimic the activity distribution of young active stars. It can be seen in the bottom panels that the large-scale field topology changes dramatically even if only the positions of the active regions are changed."

We can presently observe CMEs only in the heliosphere. For some discussion on CMEs in different astrospheres and their potential (but currently speculative) role in stellar angular momentum loss and stellar spin down, see Sect. H-IV:4.2.

11

Formation of stars and planets

A star like the Sun begins its life within a relatively dense concentration of molecular gas, called a cloud 'core', somewhere in the interstellar medium. The density in such molecular clouds is of order $10^2 - 10^6 \text{ cm}^{-3}$, to be compared with, *e.g.*, the density of the local interstellar medium of roughly 10^{-1} cm^{-3} . [H-III:3.1] "The mechanisms by which molecular clouds of many solar masses break up into stellar mass pieces are a matter of debate; probably turbulence generated in the process of forming the cloud produces the denser fragments which accrete to form stars [... G]iven the large sizes of protostellar clouds, they almost certainly contain enough angular momentum to form disks of substantial size and mass; thus, a major part of the story of star formation involves moving matter from a disk into a small, spherical protostar. {A:[141]}

A:141

To make a star of a given mass M_* from a gas with temperature T_c , gravity must overcome the pressure support; this means that [the protostellar cloud must have a] radius $R_c \gtrsim 2 \times 10^4$ astronomical units (AU; [see the argumentation around Eq. 2.15]). {A:[142]} {A:[143]} We see pre-stellar dense

A:142

A:143

¹⁴¹ Activity: How many Earth masses of elements heavier than carbon are contained in a solar mass cloud of solar composition? Most of that material in the original cloud ended up inside the Sun, of course. What fraction, roughly, of the original cloud would need to remain in the disk to ultimately form the planets? Why are the answers to these two questions largely independent of each other (think about what mostly makes Jupiter and Saturn).

¹⁴² Activity: Compare a size of $R_c \gtrsim 2 \times 10^4$ astronomical units to distances between stars in star-forming regions. Express that distance in light years and in parsecs, and compare those to the distances to the nearest stars for the present-day Sun.

¹⁴³ Activity: Another way of formulating Eq. (2.15) is to say that the mass of the cloud must exceed a certain value. Reformulate Eq. (2.15) as function of cloud temperature T_c , cloud density n_c , and stellar mass M_* (Note: this is similar to what is known as the Jeans Mass, which is commonly derived from energy imbalance or by a comparison of sound and free-fall time scales in a perturbation analysis). This shows that $M_* \sim f M_\odot T_c^{3/2} / n_c^{1/2}$. Derive the value of the constant $f \approx 2$ assuming, for simplicity, that the gas consists predominantly of molecular hydrogen. For n_c of order 100 cm^{-3} estimate M_* for $T_c \approx 10\text{K}$, characteristic of present-day molecular clouds (realizing this is a rough order-of-magnitude estimate). Early in the life of the Universe, with only H and He in the mix, the interstellar gas lacked many of the strong emission lines of heavier elements, could therefore not cool as efficiently, leaving interstellar clouds significantly warmer, roughly of order 100 K. Use the derived

concentrations of this size with properties such that they are likely to be on the verge of gravitational collapse. As these cloud cores have sizes $\sim 10^6$ times larger than the final radius of any resulting star, it is clear that virtually all of the angular momentum of the initial cloud must be transferred somewhere else; in general, it must be to a circumstellar disk. In this way, the formation of stars necessarily leaves behind material which can in principle form planets.” The initial phase of star formation, and the clearing of the dust-rich environment of the protoplanetary disk happens on a time scale of just a few million years, as we shall see in Sect. 11.2.5, and this means that much of the growth phase of planets, or at least the sizable planetesimals that later coalesce to form fully-grown planets, must be completed by then.

[H-IV:5.1] “Confirmed and candidate exoplanets number in the thousands and search techniques include Doppler measurements, transit photometry, microlensing, direct [(and since 2019 also interferometric)] imaging and astrometry. Each detection technique has some type of observational incompleteness that imposes a biased view of the underlying population of exoplanets. In some cases, statistical corrections can be applied. For example, transiting planets can only be observed if the orbital inclination is smaller than a few degrees from an edge-on configuration. However, with the reasonable assumption of randomly oriented orbits, a geometrical correction can be applied to determine the occurrence rate for all orbital inclinations. In other cases, there is simply no information about the underlying population and it is not possible to apply a meaningful correction. For example, the number of planets with a similar mass (or radius) and a similar intensity of intercepted stellar flux as our Earth is not secure at this time because the number of confirmed detections for this type of planet [and orbit is too] small. {A:[144]} {A:[145]} ”

A:144

A:145

As a result of the sample biases and observational incompleteness for each discovery technique, our view of exoplanet architectures is fuzzy at best. There are no cases beyond the Solar System where the entire parameter

expression to show that this favors the formation of much heavier stars, even when starting from a higher density of order 10^4 cm^{-3} . This review by Johnson (2019) discusses how this contributed to the evolution of elements heavier than H and He (known as ‘metals’ to astronomers) over the history of the Universe.

¹⁴⁴ Activity: Estimate the orbital Doppler swings and the fractional dimming during transits observed from afar of Mercury, Earth, and Jupiter around the Sun. Also estimate how close a Jupiter-like exoplanet (with an albedo of 0.5) should orbit for the fractional bolometric dimming during a secondary eclipse (when the planet moves behind the star) to be about 1 millimagnitude (which is the noise level for the telescope of the *Kepler* spacecraft for a 13th magnitude star at 1-minute exposure times; consider at what wavelength range the contrast is optimal). Use, *e.g.*, this fact sheet. Compare the Doppler signals with the thermal widths of spectral lines, and consider what to use as reference wavelengths. How large is the Doppler swing added to the stellar signals owing to Earth’s orbit around the Sun?

¹⁴⁵ Activity: Look up and summarize the principles of the five detection methods of exoplanets, and consider what the strengths, weaknesses, and technological challenges are for each method. Note: activities 95 and 144 ask about Doppler signals and transit photometry.

space for orbiting planets has been observed. Instead, we piece together an understanding of exoplanet architectures by counting planets in the regimes where techniques are robust and then we estimate correction factors when possible. When drawing conclusions about the statistics of exoplanets, it is helpful to understand the incompleteness in this underlying patchwork of orbital parameter space.” Sections H-IV:5.2-5.6 provide brief descriptions of the methods and their limitations.

[H-IV:5.7.4] “Our view of exoplanets is still skewed by the observational sensitivities of the techniques that we use. However, the discoveries that have been made have helped us to revise our understanding of planet formation and the formation of the Solar System. We see that planet formation is a chaotic process and that disks are sculpted by gravitational interactions to a greater extent than we appreciated by considering our Solar System. We now know that almost every star has planets and that planet formation is far more robust than astronomers expected.”^[xxii] Although we do not touch on the process of forming binary star systems (or higher multiplets), the outcome of the evolution of a molecular cloud to a planetary system often involves fragmentation of the cloud into two or more stars: roughly one in every two ‘stars’ visible in the sky is, in fact, a double or higher-multiple star.

It may be counterintuitive, but our knowledge of the evolution of the formative phases of stars by accretion from spinning disks of gas and dust that contracted out of huge molecular clouds has been helped greatly by the hunt for exoplanets and their story of formation. It is for that reason that this chapter begins with a very concise summary of what has been learned about (exo-)planetary systems, thereafter to go ‘back in time’ to the gaseous phases of the protoplanetary disks and how the gases in these formed the central stars, and how some planets ended up being ejected from the forming planetary system. {A:^[146]}

A:146

11.1 (Exo-)Planets and (exo-)planetary systems

[H-III:3.9] “The ultimate stage of disk evolution, in addition to accretion and photo-evaporation, involves the growth of planetary bodies. We now know [over 3,000 (by late 2019) exoplanetary systems], with the number continually increasing. Of course, the first major surprise was the discovery of Jupiter-mass

¹⁴⁶ Activity: Star-forming regions and disks around young stars are best observed in the near-infrared region of the spectrum. Look into what wavelengths are often used for such observations, and consider why (‘Why is the sky blue?’), given that dust sizes in the interstellar medium peak around a few tenths of a micron.

^{xxii} Recent reviews on the making of planets in general and on giant planet formation and migration, see *Space Science Reviews* (2018) volume 214, pages 38 (by Paardekooper and Johansen) and 60 (by Lammer and Blanc, referred to as ^a below), as well as ‘One of ten billion Earths’ by Schrijver (2018).

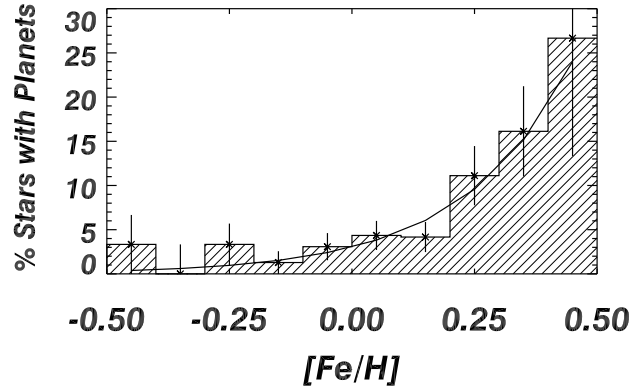


Fig. 11.1. High metallicity stars are more likely to host gas-giant planets than sub-solar metallicity stars. ['Fe/H' denotes the ratio of the abundance of Fe relative to hydrogen, while the brackets mean that the logarithm has been taken of that ratio normalized to the solar value, so that the Sun would have a value of 0, and the scale reaches a factor of about three down and up on the left and right side, respectively; Fig. H-IV:5.12]

bodies at very small orbital radii. This emphasized the almost certain necessity of inward *migration*, as it appears unlikely that disks can be sufficiently massive at 0.02 – 0.1 [Sun-Earth distances (or astronomical units: AU)] to form such objects (unless disks are gravitationally unstable all the way to the central star). The other major surprise was how eccentric [many of the orbits of close-in, large exoplanets] are. These two features are probably related, especially if planet-planet scattering is responsible for much of the inward migration. Before discussing migration further, it is useful to consider how the planets would form in the first place.

The two major scenarios of planet formation are those of core accretion [(starting with solids, and – once heavy enough, should that indeed occur – also accumulating gases)] and gaseous gravitational instability. [...] One strong piece of evidence for core accretion versus gaseous gravitational fragmentation” comes from one [H-IV:5.7.4] “of the first observed statistical correlations established[:] gas-giant planets form more frequently around metal-rich stars.^[xxiii] This planet-metallicity correlation [(see Fig. 11.1)] was used as evidence for core accretion as the formation mechanism for gas-giant exoplanets that orbit closer than a few AU around their host main sequence stars. [...] Interestingly, a similar correlation with host star metallicity has not been identified for

^{xxiii} Note that astronomers have a habit of referring to all elements heavier than helium as 'metals'.

smaller Neptune-like [...] planets. However, it remains unclear whether such correlation exists for rocky planets. [...]"

[H-III:3.9] "In the core-accretion model for giant planet formation, solid bodies accumulate via collisions until the resulting core is sufficiently large that its gravity can pull in surrounding gas. There is some concern that core accretion might proceed too slowly to explain the observed disk clearing on timescales as short as 1 – 2 Myr in significant numbers of stars. There are two potential bottlenecks in the process. One is the formation of km-sized planetesimals from cm-sized objects. Such bodies are thought to be held together lightly – too large for effective sticking and too small for gravity to become important – and, as bodies of different sizes have different velocities due to gas drag, collisions between these objects might shatter them rather than build them up. In addition, [there is growing evidence that young planets in material-rich disks are subject to rapid inward migration, which would] require fast agglomeration, especially for Earth-sized objects, though studies suggest that this inference of rapid migration may not always be correct. Various schemes of dust concentration might help avoid shattering by reducing relative motions and increasing densities, perhaps through vortices or eddies or other turbulent structures.

Once km-sized planetesimals are made, collisions among them can lead to the building of terrestrial planets and giant-planet cores. The remaining bottleneck[, at least significant for giant planets,] is that of accumulating gas. The energy released by accretion of planetesimals and gravitational contraction of the envelope must be radiated by the outer envelope. If the opacity of the envelope is large, it must extend to large radii; in turn, this can limit the gas available for accretion, which must lie close enough that the tidal forces of the central star do not overcome the protoplanet's gravity. [It appears] that, with sufficiently massive cores, giant planets can form within 1 Myr for an opacity $\sim 2\%$ of interstellar values, [because the opacity (dominated by dust) may be reduced due] to rainout of solid materials in the planetary envelope; as grain growth almost certainly precedes core formation, reduced dust opacity is an extremely plausible assumption. [...]

There is general agreement that terrestrial planets generally (fully) form later than the giant planets; gas drag is important in early stages but the final growth may well occur after gas removal from the disk. [O]nce growth to km-sized planetesimals has occurred, gravitational effects become important. At first the planetesimals grow by gravitational focusing; as they grow, eventually they excite or stir up other bodies, making their relative velocities larger and limiting accretion. The result is thought to be a set of 'oligarchic' protoplanets with relatively similar masses (at least locally). After the oligarchs have swept

up most of the available material, interactions between them dominate the subsequent evolution, with large impacts a major feature. This indicates that the final state of terrestrial planet systems is difficult to predict, as it is the result of chaotic growth.

Even after the terrestrial planets are essentially fully formed, significant system evolution can occur, simply because multi-body gravitating systems are generally not stable. A particularly interesting possibility is long-term evolution and migration due to interactions of an outer system of gas/ice giants with the planetesimals left in the outer disk, objects formed in regions with such low densities that growth to large bodies was not possible. [T]here probably has been outward migration of at least Neptune in our outer Solar System, based on the analysis of resonant structure in our own planetesimal system – the Kuiper Belt. One possible mechanism for explaining this migration is giant planet-planetesimal interactions. Such gravitational perturbations can result in the system becoming dynamically unstable, resulting in ejection and scattering of many planetesimals into high-eccentricity orbits; this has been suggested, in the so-called 'Nice' model, as an explanation for the late heavy bombardment seen in the impact history of the Moon (cf., [Sec. 12.1.1] ...)]. It has also been proposed as an explanation of the stunted growth of Mars, and moreover of the chemical gradient in the asteroid belt: silicate-rich and carbon/water-rich populations should have been differentiated by distance to the Sun (across the 'ice line' where the temperature would have been low enough to create water-rich asteroids only further out; see footnote xiv), but actually show much overlap of the populations, albeit with a clear trend for the average chemical makeup as function of orbital radius. It is argued that this smoothed trend of what should have formed as a clear chemical segregation was introduced by gravitational interaction with migrating gas/ice giants (in what is referred to as the 'Grand Tack' model, in which the Jupiter-Saturn pair first migrated inward and subsequently outward). {A:[147]} {A:[148]}

A:147

A:148

11.1.1 Exoplanet formation

The solar nebula theory holds that the Sun and its attending planetary system formed out of a cloud of gas and dust (with dust making up, on average, about 1% of the total mass^a) that contracted into a spinning disk, with most matter migrating towards the center to form a star even as much of the angular momentum ended up in the orbiting planets that formed out of the cool disk material before the remainder of the gases were somehow cleared out (more

¹⁴⁷ Activity: Figure 11.2 shows a curved 'snow line' (or 'ice line'). What is the reason behind that?

¹⁴⁸ Activity: Look up the 'Grand Tack' model and review the likely consequences for the growing Mars, for the asteroid belt, and for water distribution by scattered asteroids into the inner solar system.

on that below). [H-IV:5.7.1] “The solar nebula theory provides a theoretical description for the formation of the Solar System. Indeed, it has been said that this model is so elegant, that it is hard to imagine that it could be wrong. The solar nebula theory neatly explains most observations: the planets closest to the Sun form in a hot environment and as a consequence these planets are small and comprised of refractory elements (*i.e.*, elements [whose solids] withstand high temperatures); the more massive gas giants form beyond the ice line (a distance where it is cold enough for dust grains to be coated with icy mantles) where the feeding ground is more voluminous; jovian planets have moons that were either captured or that form as mini-solar-systems; the planets all orbit in the same direction in the disk because they inherit the same angular momentum vector; the Solar System is littered with leftover debris such as asteroids and comets. The theory supports the idea first suggested by Kant and Laplace that the proto-Sun was surrounded by a primordial spinning disk of dust and gas. All of the material that makes up the Sun drained through this disk. [...]

The mass of the protoplanetary disk is a fraction of the stellar mass and evolves with the central star. Our understanding of the physics and chemistry of protoplanetary disks is distilled in Fig. 11.2. The temperature is about 1500 K near the inner part of the disk and along the flared outer layers. These high temperatures are too hot for grain growth, but a few AU from the protostar, the disk mid-plane is cool enough for icy grains to stick and grow. The opacity of the disk is set by the dust, which gradually decouples from the gas and settles toward the mid-plane, increasing transparency of the disk over time. {A:^[149]}

A:149

Protoplanetary disks provide the initial conditions for planet formation. [...] In the first phase of planet formation, the planet grows by runaway accretion of solid material. The second phase of growth is very slow; both solid and gas accretion are nearly time-independent and this phase sets the planet formation timescale. Once the planet core reaches a mass of about $10M_{\oplus}$, [if indeed it succeeds in that,] the third phase of runaway gas accretion begins, growing the planet mass from 10 to a few hundred M_{\oplus} . [It has been] estimated that gas-giant planet formation should take roughly 10 Myr. However, observations of protoplanetary disks in the 1990s presented a conundrum: the primordial disks appear to be nearly ubiquitous around stars that are 1 Myr; at 2 Myr only

¹⁴⁹ Activity: Figure 11.2 shows a clearing near the central star. This is associated with the magnetic field of the rotating star. Consider what processes are at play there and the role of the following: accretion rate, ionization fraction, diffusion of field into the ionized gaseous disk, orbital and angular velocities, the corotation radius, winding up of magnetic field that connects the star to the disk, centrifugal force, etc. There is no easy concept for this: you can look at the literature of MHD models of T Tauri accretion disks to see how complex the coupling is. Store your thoughts: the star-disk interaction leading to the clearing is discussed in Sect. 11.2.2.

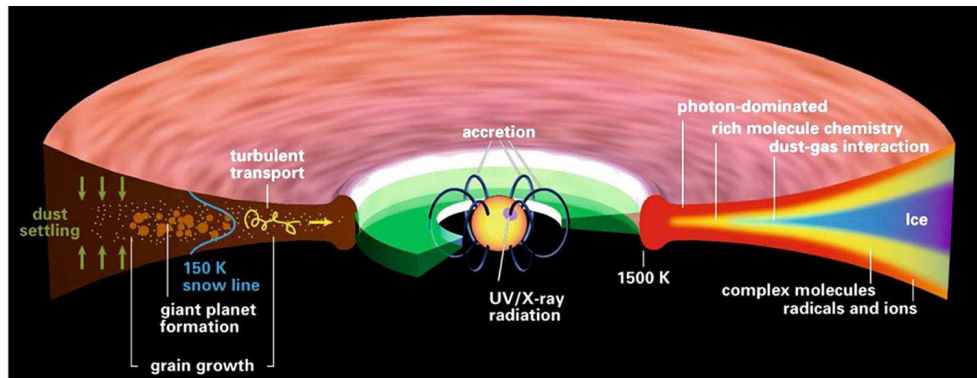


Fig. 11.2. A sketch of the structure and processes of protoplanetary disks [with ages in the range of about 1–5 Myr. Fig. H-IV:5.9; source: Henning and Semenov (2013).]

about half of young stars have disks and by 10 Myr, the disks are essentially gone. Figure 11.9 shows the fraction of protoplanetary disks found in young cluster stars.

One triumph that emerged from the discovery of exoplanets was a solution to the disagreement between theory and observations for the formation timescale of gas-giant planets. The first detected gas-giant planets orbited close to their host stars providing evidence that exoplanets could undergo orbital migration. Thus, planets were not restricted to a planetesimal feeding ground at a fixed orbital radius; instead, the planet embryos are pushed around in the disk by planet-planet interactions and tidal torques. The access to a wider part of the disk suggests a wider feeding zone for more rapid accretion of planetesimals that would shorten the second phase of gas-giant planet formation [...]

11.1.2 Exoplanet migration

In Sect. 7.3.2 we already described the possibility that planets can change their orbits in the formation phase of a planetary system by tidal interaction with the surrounding disk. This is one scenario by which, for example, giant planets may ultimately find themselves orbiting their parent star at distances much closer than where they could readily form. [H-IV:5.7.2] “Another way to push exoplanets inward is through gravitational encounters. There are several proposed mechanisms that excite orbital eccentricity including secular migration, planet-planet scattering, and Kozai perturbation in which gravitational interactions result in coupled variations in orbital inclination and eccentricity. High eccentricity planets with a small enough periastron passage eventually experience tidal circularization and can end up in short-period orbits.

Different migration mechanisms predict distinct observables. A particularly interesting observable is stellar obliquity, the relative angle between the stellar rotation vector and the vector normal to the planet orbital plane. The stellar obliquity can be measured by observing the Rossiter-McLaughlin effect. This effect is caused by a transiting object blocking some of the light from a rotating star. [In the case of a prograde, low-obliquity planet, the transiting planet first] crosses the approaching limb of the rotating star, decreasing the contribution of blue-shifted light in the spectral line and a few hours later the planet crosses the receding limb of the rotating star, decreasing the contribution of red-shifted light. The systematic decrement of Doppler-shifted light in the composite spectral lines results in a distortion of line profile, which is (mis)interpreted as a change in the radial velocity of the star. The shape of the Rossiter-McLaughlin curve during transit is entirely dependent on the stellar obliquity. Consequently, the stellar obliquity is determined by modeling the anomalous radial velocity signals during a transiting event. {A:[150]}

A:150

Disk-driven migration is expected to produce a small stellar obliquity whereas gravitational encounters that temporarily pump up the orbital eccentricity of gas-giant planets should result in a wide range of stellar obliquities including retrograde orbits. The latter has been observed for many transiting planets suggesting that high eccentricity mechanisms drive gas-giant planets inward. However, it has also been suggested that the observed stellar obliquity range may reflect a primordial stellar obliquity due to interactions between proto-planetary disk and a companion star. Interestingly, the small stellar obliquity of low-mass multi-planet systems suggests well-aligned vectors for the stellar spin and planetary orbits. It is certainly possible that gas-giant and low-mass planets migrate by different mechanisms.

In summary, the most important revisions to the solar nebula model and our understanding of planet formation can be attributed to one source: the addition of dynamical interactions between planets and the primordial disk. These dynamical interactions speed up the accretion timescales, produce mean-motion resonances, scatter planets out of the disk into non-coplanar orbits that can be detected by the Rossiter-McLaughlin effect and even eject some planets.”

11.1.3 Exoplanet geology

Studies suggest that there may be [H-IV:5.7.2] “two characteristic planet radii ($1.7R_{\oplus}$ and $3.9R_{\oplus}$) that divide planets into three populations: terrestrial

¹⁵⁰ Activity: Sketch and describe the observable spectral signatures of transiting planets for orbits of different obliquity (including effectively retrograde planets). Also: estimate transit times for planets around of solar-mass star at distances such as Mercury, Earth, Jupiter, and Neptune. Use, *e.g.*, this fact sheet.

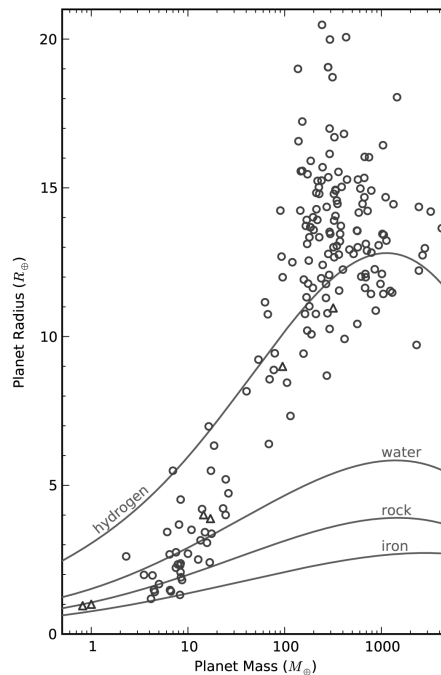


Fig. 11.3. Masses and radii of well-characterized exoplanets (circles) and Solar-System planets (triangles). Curves show models for idealized planets consisting of pure hydrogen, water, rock (Mg_2SiO_4) or iron. [Fig. H-IV:5.11; source: Howard et al. (2013).]

planets, gas-dwarf planets and gas-giant planets. [...] Both the mass-radius relationship and the transition radius from rocky to non-rocky planets help us to better understand the formation history of small planets. Planets that form *in situ* in the inner part of the disk would consist primarily of rocky materials and possibly a primordial H/He atmosphere. In comparison, planets that have undergone significant migration should contain more volatile materials such as astrophysical ice (H_2O , CO , and NH_3). The debate of whether close-in planets form *in situ* or migrate should eventually gain evidence from studies of exoplanet atmospheres that add constraints on their chemical composition.”

[H-IV:5.7.3] “Thousands of planet candidates were discovered by the *Kepler* mission, allowing for precise measurements of exoplanet radii. The combination of the radius and mass measurements (either from the Doppler technique or from transit timing variations) provides a mean density for hundreds of exoplanets and allow us to begin considering the bulk composition of unseen planets that orbit stars hundreds of light years away from us. The varying bulk composition

of exoplanets results in different curves that cut through the mass-radius parameter space shown in Fig. 11.3.

Planets with radii smaller than 4 times that of the Earth can exhibit a remarkable diversity of compositions. [...] Planets smaller than 1.5 Earth radii increase in density with increasing radius and seem to have a composition that is consistent with rock. Planets with radii between 1.5 and 4 times the radius of the Earth showed decreasing density with increasing radius, suggesting that the larger planet radius is a product of gaseous envelopes. [...] The significant amount of scatter in the mass-radius parameter space suggests a large diversity in planet composition at a given radius.”

With the growing number of exoplanet detections, one more thing has become abundantly clear: whereas the Solar System suggests a marked division between the four terrestrial planets (with masses of one Earth mass or less) and the four giant planets (with masses of 14.5 to 318 Earth masses), the exoplanet population overall has no such division, showing a continuum of masses from low to high^a.

11.1.4 Exoplanets and binary star systems

[H-IV:5.7.4] “Many stars in the solar neighborhood are components of multiple-star systems, [and exoplanets have been found orbiting one of the two components while others have distant circumbinary orbits. ...] The occurrence rate of circumbinary planets is estimated to be $\sim 10\%$ assuming the orbital plane of circumbinary planets roughly align with the binary orbital plane. The occurrence rate could be much higher if the orientation of planet orbits is more isotropic.

It is expected that planet formation may be impeded in systems where the binary stars have small separations (*e.g.*, $\sim 10 - 200$ AU). This is supported both by simulations and observations that find a smaller fraction of exoplanets in binary star systems. It is not surprising that the dynamics of binary star systems stir things up and challenges planet formation. What is surprising is that the planets exist there at all.”

11.2 Formation and early evolution of stars and disks

11.2.1 Observations of star-forming processes

Before the discovery that exoplanetary systems were about as common as stars (reached in the first decade of the 21st century) astronomers struggled to understand how angular momentum from the contracting pre-stellar cloud could be removed so that a star could form at all. There were studies on

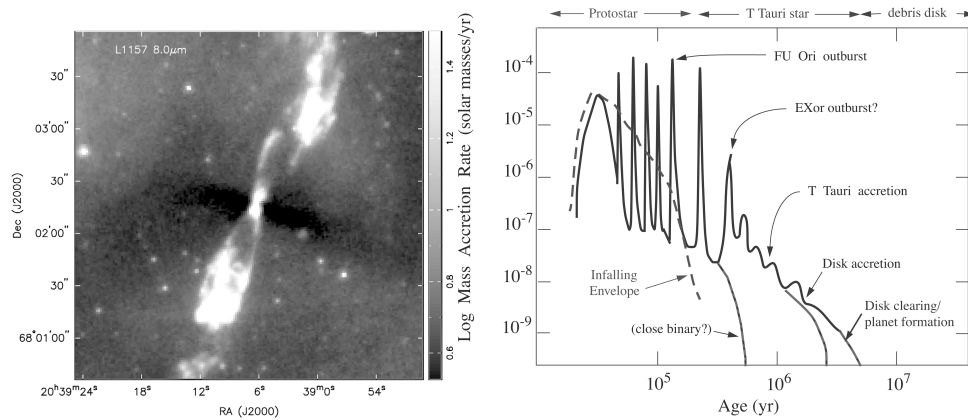


Fig. 11.4. Left: An $8\ \mu\text{m}$ image of an accreting low-mass protostar. The darker, filamentary region running east-west (horizontally in the image) represents dust extinguishing the background radiation; this indicates that the densest, most massive region of the material falling in to make disk and star is far from spherically-symmetric. The bright regions running north-south (top to bottom) are due to protostellar continuum emission reflected from dust and molecular emission lines excited by a high-velocity, bipolar outflow thought to be driven from the innermost regions of the protostellar accretion disk. Right: Schematic diagram of a likely accretional history of a typical low-mass star. The dashed curve indicates the expected rate of infall of matter from the protostellar envelope (e.g., dense region indicated in the left-hand panel). The solid curve suggests a possible variation of accretion through the protostellar disk onto the central star, which may be steady at the earliest times but is subject to strong variations in accretion (so-called FU Ori outbursts). In this picture, material piles up in the disk due to the infall rate being higher than the disk can smoothly pass on to the central star; this leads to episodic bursts of accretion which drain the excess disk mass. Finally, after infall ceases, slower, more steady accretion occurs during the T Tauri phase, which may cease because either a binary companion or planets accrete the remaining mass. This results in 'clearing' the disk, i.e., removing most of the small dust and apparently most of the gas. Finally, secondary production of small amounts of dust can occur during the debris disk stage, when solid bodies collide and shatter. [Fig. H-III:3.2; source: Hartmann (2009).]

how Alfvén waves could carry angular momentum away, how fragmentation into multiple star systems could deal with the problem, or how winds from magnetized disks could extract angular momentum. Realizing that much of the angular momentum is left behind in the planetary system reduced the magnitude of the problem tremendously, and many of the earlier ideas about where the angular momentum would end up have been left behind or now form a lesser challenge to the formation scenario of stars. [H-III:11.2.2] “If we assumed that there were no planets and all the angular momentum resides in the Sun, this leads to an increase in the angular velocity by about a factor of 35. As a result the Sun would spin around its axis in about 18 hours instead of

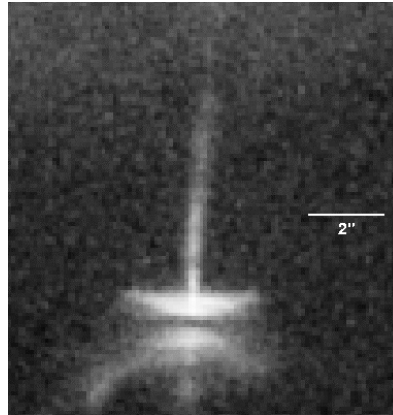


Fig. 11.5. Optical image of the accreting young star HH 30, showing the upper surfaces of its dusty disk in scattered light (the dark lane is due to dust extinction of the central star by the disk), along with an optical, high-velocity, bipolar jet. For scale, 2 arcsec = 280 Astronomical Units. [Fig. H-III:3.3]

27 days. Because the assumption of a homogeneous sphere underestimates the effective moment of inertia, the Sun would complete a rotation within about 12 hours which corresponds to a velocity of the photosphere in the order of 100 km s^{-1} instead of the observed 2 km s^{-1} .

Nowadays, the view is that [H-III:3.1] “[t]he accretion disk is basically an engine in which angular momentum is transferred outward to ever-decreasing amounts of material while the majority of the mass moves inward to the center. {A:^[151]} In the case of at least moderately-ionized disks, it seems increasingly certain that magnetic turbulence provides the necessary angular momentum transport for accretion. The low ionization of protostellar disks is likely to render this mechanism ineffective over significant radial regions; gravitational torques can come to the rescue, moving most of the cloud mass into the central regions in any event. However, gravitational torques alone will leave a sizable amount of mass in the disk, of order 10%–30% of the central star mass. As this is much larger than estimated by many techniques, and substantially more than assumed in many models of planet formation, it may be necessary for additional angular momentum transport to occur via magnetic turbulence. On longer timescales, the remaining disk gas is probably removed by some mechanism of ejection due to stellar X-rays and extreme-ultraviolet (EUV) heating.

A:151

This picture of star formation has considerable observational support. Cold

¹⁵¹ Activity: Estimate the total values and ratios of mass and angular momentum in the planetary system and in the Sun (use Fig. 10.5).

clouds of the mass and size indicated in Eq. (2.15) are seen in star-forming regions, some with already growing protostars (Figure 11.4, left). We also observe extended circumstellar disks around many young stars (Figure 11.5). The masses of these disks are at least $\sim 1\%$ that of the central star; with radii of hundreds of AU, they clearly must contain most of the system angular momentum. {A:[152]}

A:152

The implication of this picture is that most of the mass of a star must pass through its disk; that is, stars are most directly formed by disk *accretion*. As shown schematically in the right side of Figure 11.4, disk accretion may not be steady if it cannot keep up with the infall to the disk; instead, early stellar evolution may be punctuated by outbursts of very rapid accretion followed by extended periods of slow mass addition. There is observational evidence for such accretional outbursts in the FU Orionis objects; their properties suggest that disks are likely to be quite massive, at least in early stages.” {A:[153]}

A:153

Magnetic fields are good candidates for the transport of angular momentum in a disk, as the differentially rotating disk (trying to have matter orbit the forming star in Keplerian orbits) would stretch embedded field, causing a back-reaction that works to reduce the differential rotation (see Sect. 7.2.3). That can work, provided that the disk material has a sufficient degree of ionization so that the field and the gases can effectively couple. Another process that can contribute to the transport of angular momentum is gravitational coupling; see Sect. 11.2.4 for a description of this process.

[H-III:3.3] “The other major potential mechanism of disk angular momentum transport is that of winds. It is now thought that most of the angular momentum of disks results in expansion of the outer disk rather than simply being lost in a wind; however, because [Sun-like,] low-mass stars become slowly-rotating early in their existence (Sect. 11.2.3), it is quite possible that winds from the innermost disk regions play a central role in regulating the rotation of protostars.

Young stars with disks often eject powerful, collimated, bipolar winds or jets. These outflows are clearly the result of disk accretion. We can say this confidently because a) young stars without disks do not show this phenomenon, and b) mass ejection rates, as best we can determine, clearly scale with the accretion rate. Indeed, in the case of the most powerful low-mass outflows

¹⁵² Activity: Iron, oxygen, and silicon make up three quarters of the Earth’s mass. Iron is some 30% of the total. In the interstellar medium, iron makes up about 1 part in 1,000 of total mass. How many Earth-equivalents of iron does a circumstellar disk with a mass of 1% of the Sun contain?

¹⁵³ Activity: Think about similarities and differences with Solar-System magnetic instabilities as discussed in Ch. 6 when reading about things like FU-Orionis outbursts and ‘ballooning out’ of magnetic field in ejections of mass from corona and disk, likely driven by necessarily failing attempts of the forces at play to impose corotation.

– those of the FU Ori objects – accretion is the only energy source large enough to account for the necessary driving.

The high degree of collimation seen in many jets (*e.g.*, Figure 11.5) favors magnetic fields, as well-developed theory shows that rotating fields can provide the necessary collimation. Moreover, the observed outflows or jets are relatively cold; that is, the sound speeds of the gas are well below escape velocity, making thermal acceleration unimportant; and thus magnetic acceleration is not only attractive but probably necessary. What is not clear is whether *outer* disk regions exhibit outflows, at least at a sufficiently significant level to affect disk evolution. [...]

Using the basic theory of magneto-centrifugal acceleration, spatially-resolved kinematics – expansion, rotation – of jets can be used to infer the origin of the outflow, below currently resolvable scales. Observations of jets using the *Hubble Space Telescope* have suggested that the source region for the observed optical jets is ~ 0.2 to 2 AU. These estimates must be regarded as uncertain, as it is very difficult to detect the jet rotation; the analysis must assume no asymmetries in the flow, which may be questionable, given the probable presence of complex internal shocks needed to heat the radiating jet gas.

While outflows clearly emerge from the inner disk, there is little evidence for significant mass loss from outer disks, which could take away significant amounts of angular momentum. In addition, there are difficulties with assuming that the disk wind dominates angular momentum transport even in the inner disk. Removing all the angular momentum by the wind involves removing all the accretion energy in the wind as well, leaving no remaining energy to radiate; but this is problematic, because some rapidly-accreting pre-main sequence disks are self-luminous. It seems more plausible that other mechanisms – the gravitational and magneto-rotational instabilities – dominate the angular momentum transport of disks, with the winds being a byproduct of accretion. However, the slow rotation of low-mass protostars may require a powerful wind from the innermost regions to remove the final amount of angular momentum (Sect. 11.2.3).”

11.2.2 Properties of young stars

[H-III:3.4] “Solar-type stars begin their lives with only modestly-larger radii than [in the state into which they settle as ‘mature’ stars (referred to as the ‘main sequence’ phase; [*e.g.*, Fig. 4.2)]. This is a consequence of (a) the need to have a significant gas opacity to trap thermal energy, and thus produce enough pressure to halt collapse, and (b) the fact that most of the energy of accretion is radiated outward rather than being trapped. Item (b) is ensured in general

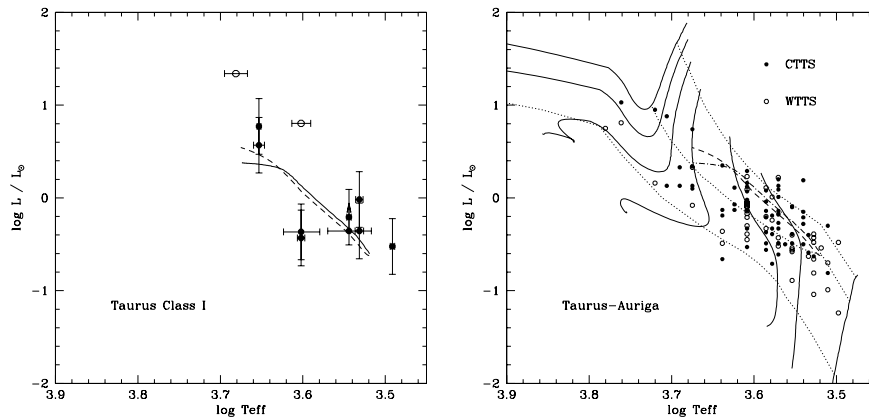


Fig. 11.6. Hertzsprung-Russell diagram positions of Taurus protostars (left) and young ([pre-main sequence] T Tauri) stars (right). These plots of two observed quantities – the stellar luminosity L (in solar units) and the effective temperature T_{eff} – can be used directly to infer the stellar radius R via the equation $L = 4\pi R^2 \sigma T_{\text{eff}}^4$, and indirectly the stellar mass via evolutionary tracks. Left: Solid and dashed curves correspond to theoretical estimates of initial protostellar radii ('birth lines') as a function of effective temperature (which corresponds roughly to mass). The open circles denote objects in which most of the luminosity derives from accretion, not stellar photospheric radiation. The agreement between theory and observation is reasonably satisfactory given the uncertainties, showing that low-mass protostars do indeed begin their existence with radii only a few times larger than that of the Sun's. Right: Standard stellar evolutionary tracks compared with observed HR diagram positions of T Tauri stars in the Taurus–Auriga star-forming region. The dashed lines show approximate isochrones for 1 Myr and 10 Myr, assuming contraction from very large radii, along with the birth lines of the left-hand panel. Ages of young solar-type stars are thus determined by the amount they have descended in the HR diagram from the birth line, due to gravitational contraction. [Fig. H-III:3.8; source: Hartmann (2009).]

by the very high opacity of the protostar compared with the infalling material, and in particular by the angular momentum of the protostellar core, which makes much (most) of the material land first on the disk rather than onto the central star.

In the absence of energy input, [so prior to the initiation of nuclear fusion, the star-to-be] contracts on the Kelvin–Helmholtz time scale

$$\tau_{\text{KH}} = \frac{3 GM_*^2}{7 R_* L_*} \quad (11.1)$$

where R_* is the protostellar radius and L_* its luminosity. This is basically the ratio of the internal energy divided by the rate at which energy is being lost, with the numerical coefficient set in this case by the assumption that the

star is completely convective. {A:[154]} More detailed calculations indicate A:154 that during protostellar accretion, the protostellar luminosity and radius have roughly those values which would yield a Kelvin-Helmholtz contraction time of the same order as the timescale for infall. In the case of the protostellar cloud described above, this timescale is $\sim R/c_s$, or a few times 10^5 yr.

For low-mass protostars, fusion of deuterium can play an important role in stopping protostellar contraction at early times. Deuterium fusion occurs at a significant rate when the central temperature reaches $\sim 10^6$ K; this results when $R_*/M_* \sim 5R_\odot/M_\odot$ for a completely convective star. However, as D has a very low abundance, its fusion represents a significant energy source for only a modest time at low masses and very short times for higher-mass, higher-luminosity objects. The result is that stars of masses $\lesssim 0.5M_\odot$ may be detected initially near the D main sequence in the Hertzsprung-Russell (HR) diagram (see Sect. 10.1 and Figs. 4.2 and 11.6), but the youngest higher-mass objects will be found below this 'birth line'). After D is exhausted, the solar-type star will then undergo Kelvin-Helmholtz contraction until it reaches the main sequence, as shown in Figure 11.6. {A:[155]} A:155

Stellar ages for very young stars are estimated from Kelvin-Helmholtz contraction timescales. The accuracy of these estimates depends mainly on uncertainties in two quantities: the stellar mass and the 'starting' radius for KH contraction (left-hand panel of Figure 11.6). Masses are mostly estimated from theoretical evolutionary tracks, though progress is being made in calibrating these from binary orbits and disk rotation; currently there are significant uncertainties for the lowest-mass stars. For higher masses, calibrations are better but the starting radius or birth-line position is uncertain, as it depends upon the precise thermal content of accreted matter rather than on the occurrence of D fusion (see Figure 11.6). For solar mass stars, the upshot is that ages are uncertain by a factor of two or more for Kelvin-Helmholtz estimates at ~ 1 Myr, and perhaps 30% at 10 Myr.

Stellar magnetic fields and activity are important for understanding the angular momentum 'problems' [...]. In brief, large areas of the photospheres

¹⁵⁴ Activity: The internal energy of the star in Eq. (11.1) is derived from the so-called 'virial theorem' which states that the total gravitational energy E_{grav} is related to the total thermal energy E_{thermal} as $E_{\text{grav}} = -2E_{\text{thermal}}$ if $\gamma = 5/3$ as for a monoatomic ideal gas. Derive this from Eq. (3.5) assuming a field-free stationary state for a spherically symmetric ball of gas: $dp/dr = -GM(r)\rho/r^2$. One way to do so is to multiply both sides by $4\pi r^3$, integrate (in part 'by parts') from center to surface (where $p(R)$ essentially vanishes, and realizing that the internal energy per unit volume of the gas is given by $u = p/(\gamma - 1)$ for an adiabatic exponent γ). The result is equivalent to the virial theorem. Eq. (11.1) can be used for the present-day Sun to show that continued gravitational contraction cannot support the solar energy budget over the age estimated for the Earth based on radio-nuclide dating (note a factor of two difference between thermal and gravitational time scales). What is the present-day value of τ_{KH} in Eq. (11.1) for the Sun?

¹⁵⁵ Activity: Draw lines of equal radius (as multiples of the solar value) in Fig. 11.6, using $\log(T_{\text{eff},\odot}) = 3.762$.

of very young stars are covered with strong magnetic fields, with $B \sim 2$ kG and covering (or filling) factors of tens of percent. Polar dark spots seem to be typical, though there are significant spots at other latitudes, and the spot areas/fields are not axisymmetric – explaining why there is often substantial rotational modulation of the optical/near-IR stellar photospheric emission. [...] The variability of the rotationally-modulated starspot-produced light curves – on timescales of days, weeks, months, years – indicates that the fields are not fossil in origin but are produced by some sort of stellar dynamo. [...] The large-scale (dipolar) magnetic field strengths of these stars are important in understanding the interface between the accretion disk and the stellar photosphere. [...] While Zeeman *broadening* clearly demonstrates the existence of 2 kG photospheric fields over substantial areas of the star, the low measurements or upper limits of *polarization* suggest that there must be substantial [polarity] reversals to cancel out the net polarization; this would seem to indicate that the fields are of higher order than dipole, and thus that the large-scale (dipolar) component may be relatively weak, [although it appears that there are] non-negligible large scale fields nonetheless. {A:40}

An important consequence of the large magnetic fields of pre-main sequence stars is that the stellar magnetospheric pressure and torques truncate the disk accretion disks well above the stellar photosphere [(as sketched in Figs. 11.2 and 11.7)]. Magnetospheres are certainly present, given the strong fields found empirically. Moreover, it is clear from observations that [young, still fully convective, pre-main sequence] T Tauri stars accrete through their magnetospheres. The high $H\alpha$ emission and the strongly Doppler-broadened $H\alpha$ emission line profiles of accreting T Tauri stars are convincingly explained by some type of quasi-radial infall; this implies that the rapid rotation and slow radial drift of accreting material in the disk must be disrupted, most plausibly by the stellar magnetosphere (Figure 11.7). The magnitude of the observed velocity line widths can be explained only if the stellar magnetic field is strong enough to truncate the disk at least a few stellar radii above the photosphere, allowing the essentially freely-infalling gas to develop a large gravitationally-produced velocity.

In addition to broad emission lines, accreting T Tauri stars exhibit significant amounts of excess continuum emission at wavelengths running from the far-ultraviolet through the optical region. This ultraviolet-optical continuum emission is most plausibly explained as radiation produced in the accretion shock at the base of the magnetosphere, where the material in near-freefall comes to rest at the stellar photosphere. As described in the previous paragraph, it appears that the disk must be magnetospherically truncated at a few stellar radii above the photosphere; this implies that most of the energy generated by

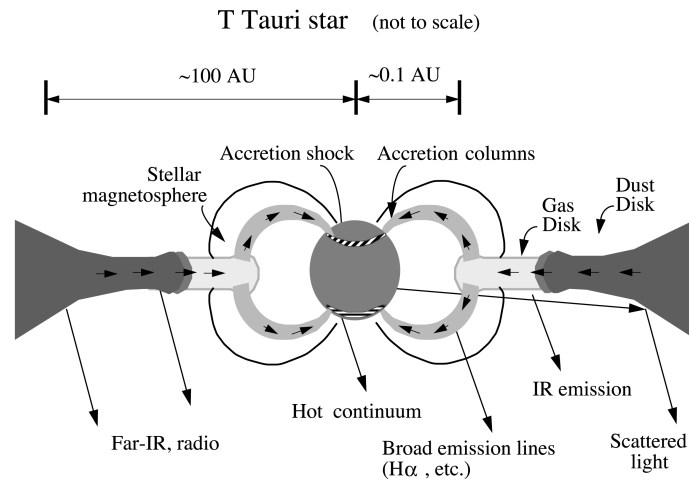


Fig. 11.7. Schematic representation of magnetosphere-disk interaction in low-mass, pre-main sequence (T Tauri) stars, with diagnostics of specific regions labeled. [Compare this to the discussion of the 'open magnetosphere' of planets in Sect. 5.5.6 and the effects of corotation and its failure in Sect. 5.5.7.2. Fig. H-III:3.10; source: Hartmann (2009).]

accretion will be radiated in this accretion shock. Estimates of mass accretion rates \dot{M} for T Tauri stars are thus generally based on setting the UV-optical emission excess luminosity $L_{\text{acc}} \sim GM_*\dot{M}/R_*$.

11.2.3 The rotation rate of very young stars

[H-III:3.5] “One of the most striking problems of angular momentum transport is that very slowly-rotating low-mass stars are produced by accretion from rapidly-rotating disks. In general, T Tauri stars of masses $\lesssim 1M_{\odot}$ rotate at rates from a few tens of percent to less than ten percent of their breakup values. {A:^{[156]}} The problem of producing slowly rotating stars somewhat older is made much more difficult by the apparent requirement of spinning down the star at the same time it is accreting high-angular-momentum material. Of course, if magnetic stellar winds were intrinsically powerful enough to spin down stars rapidly, there is no problem; but spindown does not seem to be extremely rapid in non-accreting stars, at least not on timescales needed to explain the slow rotation in stars of ages ~ 1 Myr.

One possible option is that the magnetospheric coupling between the star and its disk transfers the angular momentum outwards at the necessary rate. However, there are difficulties with applying this model. In the first place

¹⁵⁶ Activity: Derive the expression for the breakup rotation rate of stars as function of mass and radius. What is the value for the Sun? Ignore distortion from spherical symmetry for this estimate.

accretion, which is observed in essentially all T Tauri stars with detectable inner disks, basically requires magnetic field lines tied to disk material inside of corotation; this spins down the gas so that it can accrete, spinning up the star. Spindown of the accreting star requires magnetic fields connected to the disk outside of corotation; thus, to explain T Tauri stars one would like one set of stellar magnetic field lines to be connected inside of corotation, and another outside of corotation, and somehow balance the angular momentum addition due to the accretion with coupling to the outer disk. Numerical simulations indicate that a quasi-steady state [with both types of connections] may be possible with a large enough turbulent diffusivity, but whether such diffusivities are realistic is unknown. [However, some] estimates of inner gas-disk radii are significantly inside of corotation, raising the question as to whether there is a strong enough large-scale magnetic field to effectively couple to the outer disk for spindown.

T Tauri magnetospheres are probably best thought of as a series of individual magnetic loops, not all of which are filled with accreting gas; this makes it easier to explain the very small covering factors of the hot (shocked) continuum regions on the stellar photosphere of order $\lesssim 1\%$. As at least some of the loops (if not most) must connect to the disk interior to corotation, it is almost certainly the case that magnetic field lines must tend to become twisted. Such twists rapidly lead to a 'ballooning out' of closed field lines, with eventual opening up of field lines and possible ejection of mass, with reconnection following. [...]” Such processes would appear to make the angular momentum transfer from star to disk even less efficient, although processes related to winds and waves complicate the modeling and our understanding of how things work (as discussed in Ch. H-III:3.5). It is also possible that heating of part of gas coming into the stellar magnetosphere enables a hot stellar wind from within the star's magnetosphere, which could lead to efficient magnetic braking. Clearly, for now, the loss of angular momentum from the material accreting onto the protostar remains an area of study.

11.2.4 *Protoplanetary disks and gravity*

[H-III:3.6] “The mechanisms of angular momentum transport determine the mass distribution within the protoplanetary disk. It is important to understand whether gravitational instabilities dominate this transport, in which case accretion onto the central star is likely to decay away with time, leaving a relatively massive disk behind; or whether another mechanism not tied to gravity can reduce disk mass distributions leading to the epoch of planet formation.

The one non-gravitational mechanism of angular momentum transport that

we currently understand (at some level) is the magneto-rotational instability (MRI; Sect. 7.2.3). It is possible that the upper layers of the otherwise cold disk can be non-thermally ionized by stellar X-rays [(as suggested in Fig. 11.2)] and cosmic rays [entering from outside the system], to the extent that a significant amount of mass and angular momentum transport can occur. If large amounts of the disk can be activated magnetically in this way, then the disk can behave essentially as a standard viscous disk, with most of the mass at large radii. However, X-ray and cosmic ray ionization are insufficient if small dust grains, which can absorb ions and electrons very efficiently, are not heavily depleted. [... Whereas *Spitzer IRS (Infra-Red Spectrograph)* spectra suggest levels of depletion of 10^{-2} to 10^{-3} from interstellar medium values of small dust, it appears that] depletions of order 10^{-4} are needed for the MRI to operate robustly in upper disk layers.

As discussed earlier, it is plausible if not likely that protostellar disks are initially gravitationally unstable, given the need to accrete most of the mass of the central star through the disk and likely limited MRI transport in cold disks. If the MRI is inefficient, the disk could settle into a state of marginal gravitational instability, with the Toomre parameter

$$Q = \frac{c_s \Omega_e}{\pi G \Sigma} \sim 1.4 \quad (11.2)$$

where Ω is the Keplerian (presumed to be the epicyclic) angular frequency and Σ is the disk surface mass density [(the epicyclic frequency is the frequency at which a radially displaced parcel oscillates within the disk)]. The Q parameter basically results from satisfying two conditions: one, that gravity can overcome resisting gas pressure forces; and two, that gravity is stronger than the effects of angular momentum in opposing collapse. Larger values of Q mean that the disk is gravitationally-stable, while smaller values of Q indicate strong instability. In many instances disks tend to self-regulate; strong instabilities tend to produce heating via shocks which raise c_s and thus increase Q , until the sound speed rises sufficiently that the instabilities heating the gas begin to decay.

Even if the MRI is reasonably well activated by non-thermal ionization, it may easily be insufficient over the 1 – 10 AU region to transport all the mass viscously; this could result in the general picture in which a 'magnetically dead' zone of the disk is sandwiched radially by MRI-active regions at small and large radii.

To develop this further, consider estimates of the mass distribution of the solar nebula. [Figure 11.8 compares two different estimates of the so-called 'minimum mass solar nebula' (MMSN), one using the so-called 'Nice'

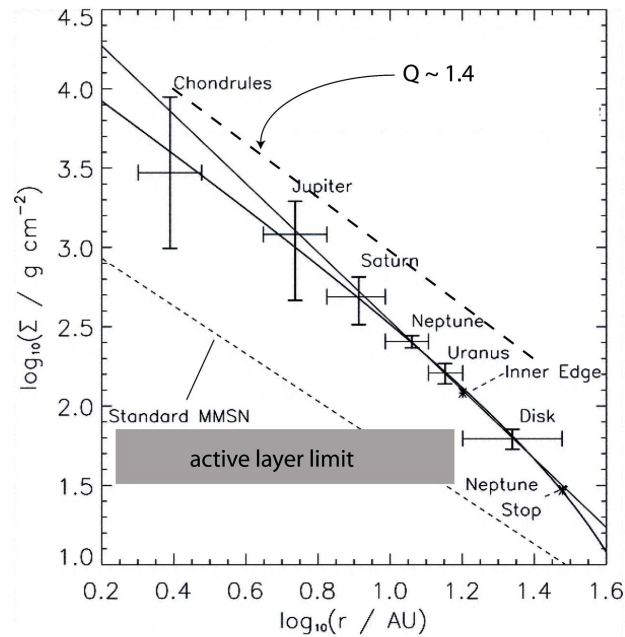


Fig. 11.8. Two estimates of the minimum mass solar nebula. The lower, light-dashed curve indicates the usual estimate, derived from the current position of the giant planets and accounting for the missing light elements; the solid curves show a higher estimate based on the initial positions of the giant planets assumed in a model which has substantial outward migration of the giant planets. Limits on the expected MRI-active surface density due to non-thermal ionization and on the surface density expected for a marginally gravitationally-unstable disk (the dashed line showing the condition for the critical value of the 'Toomre Q ' parameter – see Eq. (11.2) – are also shown. [Fig. H-III:3.12; source: Desch (2007).]

model (falling not far below the gravitational instability result with Eq. 11.2), which posits substantial inward migration of Jupiter and Saturn and outward migration of the Uranus and Neptune from their original positions, the other an older version based on the current positions of the giant planets. Both estimates lie] above the maximum $\Sigma \sim 100 \text{ g cm}^{-2}$ estimated for non-thermal ionization by cosmic rays in the most optimistic scenario. While either version of the MMSN must be considered uncertain, the possibility that the solar nebula had a '[magnetically] dead zone' must clearly be considered. [...]

The consequence of a disk structure with a 'dead zone', as described in the previous paragraph, may be highly time-variable accretion during the protostellar phase. [The gravitational instability, GI] can be relatively efficient in transferring mass inward at large disk radii but tends to become inefficient at small radii; conversely, the MRI becomes increasingly important at small radii, especially at high mass accretion rates. If matter moving inward under

GI dissipates enough energy locally in the inner disk, it can 'turn on' the MRI thermally, resulting in an onrush of mass onto the central star. This picture has been invoked to explain the FU Orionis outbursts [during which] of order $10^{-2}M_{\odot}$ gets dumped onto the central low-mass star over timescales $\sim 10^2$ yr. It is difficult to explain the FU Ori outbursts without having a large amount of disk mass at a few AU, well above that of the standard MMSN.

The possibility of gravitational instability [makes one reconsider] the possibility of forming giant planets directly through gravitational fragmentation [rather than by the core-accretion scenario described near the top of Sect. 11.1]. This suggestion runs into difficulty, however, because a low Q is not enough; the disk must be able to cool on something like an orbital period P_{orb} to continue fragmenting; otherwise perturbations shear out and transport angular momentum instead. This poses a problem for protostellar disks because they are so cold, and thus do not cool rapidly. The cooling timescale t_c for an optically-thick disk [...] is basically the energy content divided by the blackbody radiation loss. Numerically, for temperatures below 170 K, one finds

$$\frac{t_c}{P_{\text{orb}}} \sim 10^4 \left(\frac{M}{M_{\odot}} \right)^{3/2} R_{10}^{-9/4} \quad (11.3)$$

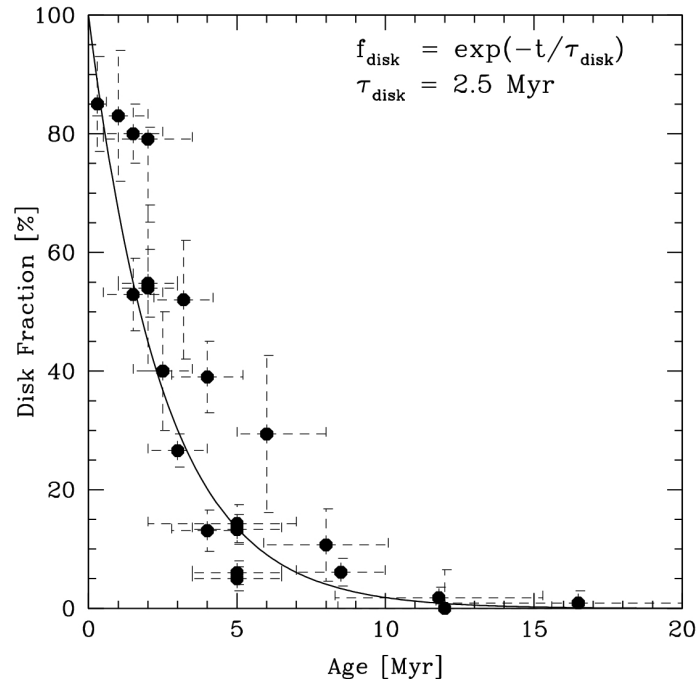
(with R_{10} is a characteristic scale in units of 10 AU), which poses an obvious difficulty for fragmentation in that the cooling time far exceeds the Keplerian period. (Things change on distance scales ~ 100 AU or larger, because the disk typically becomes optically thin, and thus cools much more rapidly than indicated by the above equation.)

Even if fragmentation could occur after infall ceases, one would still expect it to be more important early on, when the disk is more massive. It is not obvious how initial gravitational instability would explain the observed clearing of disks over millions of years."

11.2.5 Dust-disk evolution

[H-III:3.7] "In the core accretion model for the formation of giant planets, and in all models of terrestrial planet formation, dust grains grow from sub-micron sizes to thousands of km. A starting point for thinking about how planets grow from disks is then considering observations of the evolution of disk dust, detected through its emission.

Figure 11.9 shows the estimated fractions of young stars in various groups with large dust-disk excesses as a function of age. [...] The overall result is that optically-thick dust disks (with the opacity probably dominated by particles of μm size or a bit less) disappear on timescales of a few Myr. While less is known about the presence of gas in the inner 10–20 AU, clearing of small



*Fig. 11.9. Primordial disk fractions of stars in young clusters. These observations show that the dust disks only last for a few million years. [Fig. H-IV:5.10; source: Mamajek (2009). Note that since this figure was made the ages for very young stars systematically increased by ~ 50 -100% for many of these young samples (see, e.g., work by Bell *et al.* (2013), by Bell *et al.* (2015); by Pecaut *et al.* (2012), and by Pecaut and Mamajek (2016).), so that the e-folding timescale for primordial disk loss is something like 4-5 Myr rather than the ~ 2 Myr shown in the figure (Eric Mamajek, private communication).]*

dust particles seems generally accompanied by removal or disappearance of gas as well. It is important to emphasize that there is no single timescale for disk clearing. Some (inner) disks disappear immediately, perhaps because of disk disruption by a binary companion; others take a few Myr; a small percentage last for 10 Myr.

The disk can 'disappear' in one of three ways; mass can be accreted, ejected, or condensed into large bodies. It is difficult to accrete all the mass of the disk, as some must be left behind to take up the angular momentum; the outer disk is likely to expand over time and evolve on continually slower timescales. Evaporation of the disk may be important, though it is thought to take place over longer timescales than this (Sect. 11.2.6). Perhaps the strongest evidence for coagulation into larger bodies is the detection, either through spectral energy distribution fitting and/or imaging, of systems with substantive outer

disks but inner disk holes or gaps. This is consistent with the idea that settling, grain growth, coagulation, and formation of large solid bodies occurs fastest in the inner disk, where the surface densities are largest.

Dust grains in the disk generally are thought to evolve to larger sizes, with a decreasing population of small grains with increasing age. During this overall growth, dust is expected to settle vertically and drift radially. {A:^[157]} This evolution of dust in size and position in the disk can reduce and ultimately eliminate infrared excess emission, consistent with the observed disappearance of dusty disk emission over millions of years (Figure 11.9). In principle, dust growth can be extremely rapid: [the disk interior to about 10 AU may become] optically thin on timescales of 0.1 to 1 Myr, as dust particles settle and coalesce into larger bodies. The evolution of the mm fluxes is slower because of longer timescales of accumulation in the outer disk, with substantial reductions in mm-wave emission on timescales of 10 Myr. One might expect that turbulence would lengthen settling and growth timescales, but [it may actually stimulate] growth due to turbulent mixing. [...]

A:157

As particles grow in size, many effects converge to make evolution uncertain. For example, the difference in velocities between objects within an order of magnitude of meter size can result in their complete shattering or disruption. Turbulent eddies or whirlpools might help collect these objects at low velocities so that they can accrete, or alternatively disperse them more widely.

Can core accretion proceed fast enough to explain the observed disk clearing on timescales as short as [several] Myr? One problem is the formation of km-sized planetesimals from cm-sized objects. Such bodies are thought to be held together lightly – too large for effective sticking and too small for gravity to become important – and, as bodies of differing sizes have differing velocities due to gas drag, collisions between these objects might shatter them rather than build them up. Another problem is that the so-called Type I inward migration due to torques between the disk and the body is very rapid, making it important to grow quickly at ~ 1 Earth mass to avoid falling into the central star on a timescale < 1 Myr. These estimates have usually been made in the 'minimum' MMSN (Figure 11.8); the timescale for inward migration is inversely proportional to the surface density, so gravitationally-unstable disks may pose even bigger problems in this regard.

Once km-sized planetesimals are made, collisions among them can lead to the building of terrestrial planets and giant planet cores. The remaining bottleneck is that of accumulating gas which depends upon the opacity; larger opacities make it difficult for the growing planet to lose energy and allow additional

¹⁵⁷ Activity: The key mechanism by which dust is expected to settle into the center of an accretion disk is hydrodynamic drag. Explain how this works. Consider orbital inclination and effects of gas pressure, gravity, and stratification.

material to be accreted. A reduction in opacity due to grain growth and depletion would help considerably in this regard.”

11.2.6 Disk evaporation

[H-III:3.8] “As the planets are overabundant in heavy elements relative to the Sun, it is clear that most of the original gas in the solar nebula has been lost. Of course some of it accreted into the Sun, but it is unlikely that all of this material was removed in this way. For some time it was thought that a powerful solar-type wind was responsible for gas removal from the nebula. However, we now realize that the strong mass loss we see is not a solar wind but a disk wind; more importantly, the wind material is ejected perpendicular to, not into, the disk (Figure 11.5).

The high-energy radiation emitted by T Tauri stars provides a mechanism by which the gas of the disk can be evaporated rather than accreted. In this case, rather than generating stellar mass loss from the star via a coronal wind, one can generate disk mass loss from a much lower temperature wind because the material is ejected from much farther out in the gravitational potential field, where the escape velocity is very much smaller than at the stellar surface. Using the usual Parker wind formula (*e.g.*, Eq. 7.11), and assuming photoionization and thus heating to a typical temperature of $\sim 10^4\text{K}$, the sonic point occurs for

$$R_s \sim \frac{GM_*}{2c_s^2} \sim 3.6 \frac{M_*}{M_\odot} \text{ AU}, \quad (11.4)$$

where the mean molecular weight is 0.67, appropriate for a gas of cosmic abundance with ionized hydrogen and neutral helium. Thus, ionizing photons have the potential for removing disk gas at radii of a few to ten AU from the central star.

To see the essential physics of the problem with a minimum of geometrical complication, assume that a volume of $4\pi R^3$ must be ionized, where R is a characteristic radius of escape. This estimate is justified because the gas must maintain its ionization over the disk to a distance comparable to its escape radius to flow out of the gravitational potential well. The balance between photoionization and recombination leads to

$$\Phi_i = 4\pi R^3 n_e n_p \alpha_B, \quad (11.5)$$

where Φ_i is the flux of ionizing photons from the central source, n_e and n_p are the electron and proton densities, respectively, and α_B is the Case B recombination rate for hydrogen. ^[xxiv] Assuming complete ionization of hydrogen, the mass

^{xxiv} ‘Case B’ recombination considers only recombinations in which the recombined electron transitions

loss rate is

$$\dot{M} \sim 10^{-9} \Phi_{i,41}^{1/2} R_{10}^{1/2} M_{\odot} \text{ yr}^{-1}, \quad (11.6)$$

where $\Phi_{i,41}$ is the Lyman continuum photon flux in units of 10^{41} s^{-1} and R_{10} is a characteristic scale of the flow in units of 10 AU. This estimate illustrates the potential of photo-evaporation to remove disk gas over evolutionarily interesting timescales. Much more sophisticated treatments of the outflow have been considered, but this illustrates the basic result. [...]

Unfortunately, the true ionizing fluxes of young stars are not really known because interstellar absorption prevents direct detection [..., but there are observational results that suggest] evaporation of disks due to stellar magnetic activity occurs on timescales of order 10 Myr or more. Whether photo-evaporation plays a major role in the strong disk evolution from 1–10 Myr remains unclear.

Disks close to a hot luminous star can be photo-evaporated rapidly due not only to EUV (Lyman continuum) radiation but also by far-UV ($\sim 1000 \text{ \AA}$) radiation, which can heat the gas to temperatures $\sim 1000 \text{ K}$ as electrons are driven off grains. The FUV radiation thus can drive a wind off the outer disk, and may be more important in many systems if most of the disk mass resides at large distances. [...] Although the solar nebula appears to have been 'polluted' by ejection from a supernova, it is not clear that it was close enough to the massive star such that FUV radiation was important in evaporating Solar System gas." {A:[158]}

A:158

¹⁵⁸ Activity: For further study/reading: Most stars are born in groups of substantial numbers (often in what are called 'open clusters'). In such clusters, stars of a range of masses are formed (statistically yielding the 'initial mass function'). The heaviest among these evolve fastest, and if heavy enough can end their lives in a 'supernova'. The open cluster is eventually pulled apart by the 'galactic tides', which limits the exposure of planetary systems to nearby supernovae and to gravitational perturbation of the orbits of the planets. Look up the terms between quotation marks. The occurrence of a nearby supernova appears consistent with several properties of the solar system, including one of several possible means for the early melting of small bodies (as reflected in what are known as 'chondrules'). Look at this study by Portegies Zwart *et al.* (2018) for more on this.

to the ground state via intermediate transitions; a direct transition to the ground state would emit a photon that could be absorbed and lead to ionization in which case no net recombination would have occurred.

12

Evolving irradiance, atmospheres, and habitability

12.1 Evolving planetary habitability

Historical records are too short for us to see first-hand accounts of Earth in a significantly different climatic state than the present one. There are, of course, reports on the relatively recent moderate (but nonetheless impactful) excursions from the mean climatic state, such as the Medieval Warm Period, the Little Ice Age, and the modern-day onset of global warming, but there have been much larger changes over the life of the planet. Substantial modifications of climate in the past have been attributed to the formative processes of Earth and asteroid impacts, to the evolving spectral output of the Sun and the stripping effects of the solar wind, to orbital changes in response to the gravitational pull by the giant planets, to the torque applied by the Moon, to geological and geochemical activity over eons (including the geodynamo), and – last but by no means least – to the emergence and evolution of life. This chapter provides brief introductions to each of these drivers of the terrestrial atmosphere and its climate system. This provides insight into the diversity of conditions on Earth over time, while also setting the stage for appreciating the challenge of establishing the 'habitability' of planets elsewhere in the universe.

12.1.1 *Earth's formative phase*

[H-III:4.5] “Earth’s formation, like that of the other solid planets, occurred by accretion of solid materials. The processes began with particles of dust, but collision and sticking processes rapidly led to the formation of larger and larger bodies. An important aspect of the growth of rocky planets is the amount of a planet’s mass that is accreted in the form of large chunks. The accretional growth process yields a number of Moon to Mars-sized 'embryos' in a given radial region of the nebula. The final assembly of a rocky planet involves both the accretion of numerous large embryos as well as gravitational ejection of

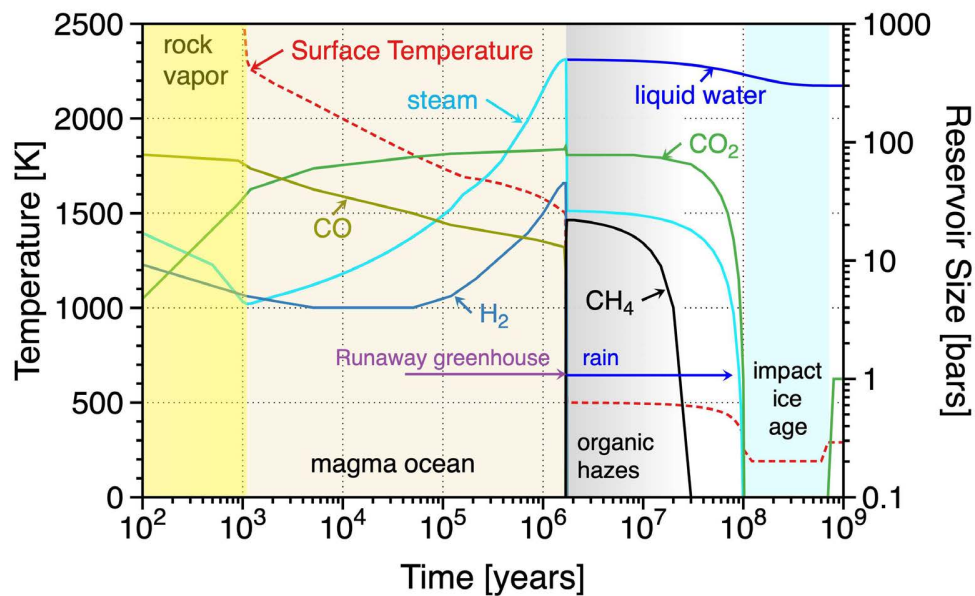


Fig. 12.1. The Earth's surface temperature and above-surface reservoirs of water and carbon dioxide after the Moon-forming collision. The surface temperature drops below 1000 K after a few million years when Earth's steam atmosphere condenses, and it drops below 500 K to habitable conditions after about 100 million years when most of the atmospheric CO_2 is incorporated into the mantle. [This is an updated version, courtesy of Kevin Zahnle (July 2019), of the originally used Fig. H-III:4.3; the latter was from this source: Zahnle et al. (2007).]

some of them to other locales [(including, as we now know through gravitational microlensing, out of a planetary system altogether and into interstellar space)].

This formation mode that includes impacts of very large bodies is indicated both by the numerical simulation of accretion processes and by evidence that our Moon formed as the result of the impact of a Mars-sized body with the growing Earth [...] {A:¹⁵⁹} Following lunar formation, Earth's post-impact atmosphere of vaporized silicates may have condensed in ~ 1000 yr (Fig. 12.1). The heat of the impact would have melted and partly vaporized Earth's mantle, but the resulting silicate magma ocean may have solidified in only a few million years. Once the magma ocean crystallized, cooling conditions would have allowed the large amount of water vapor injected into the atmosphere to condense and thus reduce the extreme greenhouse warming of the early Earth

¹⁵⁹ Activity: For an impression of order-of-magnitude numbers, estimate the energy involved in a collision between an Earth-mass body and an Mars-mass body at an impact velocity of, say, 14 km/s. Ignoring the energy going into the formation of the Moon in such a process, but rather assuming all mass and energy remain within the newly formed body, estimate the average temperature increase if all energy were distributed throughout half of the volume of the mantle, and that that material has a specific heat of approximately 1.5×10^7 erg/g/K.

A:160

and allow surface temperatures to drop below 1000 K. {A:[160]} Even with most of the water condensed, the atmosphere would still retain ~ 100 bars of CO_2 whose greenhouse warming would keep the Earth's surface temperature at ~ 500 K, even though the early Sun was $\sim 30\%$ fainter than its present brightness. The final lowering of the Earth's surface temperature to habitable conditions requires transfer of most of the atmospheric CO_2 to the mantle and crust, a process that can happen over a timescale of 10 – 100 Myr [... by the process of weathering (more on that below).]

[...] Earth's oldest known rocks, whose properties could provide information about the early Earth, are just less than 3.9 billion years. This is a curious age: the Earth's oldest surviving rocks formed just after a rock-destroying time period known as the Late Heavy Bombardment or LHB. [...] The origin of the LHB has long remained a mystery. Solar system formation models as well as the observed crater record suggests that the LHB was not just the tail end of the planetary accretion process. The presence of heavily cratered regions on other bodies, including Mars, suggest that the LHB may have been a Solar-System wide process. [...] The 'Nice Hypothesis' [(named after the city in France at whose university this hypothesis was first formulated)] suggests that a dramatic rearrangement of the outer planets gravitationally perturbed a large number of cometary bodies into orbits that penetrated the inner Solar System and cratered the surfaces of all Solar-System bodies [(see Sect. 11.1) ...]"

12.1.2 The habitable zone

[H-IV:4.3] "One of the most important requirements for life as we know it is water. The ability to retain surface water is the general basis of the concept of the Habitable Zone (HZ). As most commonly used, the habitable zone is an estimate of the range of distances from a star where an Earth-like planet can maintain surface water for extended periods of time. {A:[161]}

A:161

While a number of factors, including greenhouse gases, tilt of spin axis, planet composition, surface gravity and cloud properties can be important

¹⁶⁰ Activity: Make an order of magnitude estimate of the cooling time of Earth's atmosphere after impact of a Mars-mass body: assume an impact velocity of 14 km/s, that all kinetic energy remains within the near-surface layers and atmosphere; an optically thick atmosphere of vaporized silicate; and a characteristic temperature of the radiating vapor of, say, 2000 K.

¹⁶¹ Activity: Although the definition of 'habitability' commonly involves the requirement of liquid surface water, some definitions are more relaxed. Perhaps other surface liquids can serve as agents in support of life (such as ethane and methane lakes and seas on that cover 1.6 million square kilometers, or 2% of the surface, of Saturn's moon Titan) or perhaps subsurface water (as encapsulated seas or even globe-spanning layers) can support life. With that in mind, explore the moons of the giant planets that are thought to meet at least the condition of large reservoirs of some liquid somewhere, in particular: Europa, Callisto, Ganymede, and Io at Jupiter, Enceladus and Titan at Saturn, and Triton at Neptune. Which three power sources are thought to be most important in maintaining liquid states on giant-planet moons?

for habitability, the primary factor considered for the habitable zone is the most fundamental, just the distance from the star (see below around Eq. 12.3). For the present-day Sun, the habitable zone is generally considered to be the range from just inside Earth's orbit to a region near or just beyond Mars' orbit. The inner boundary is where surface water is lost to space by either a runaway greenhouse effect or the 'moist greenhouse' effect. In a full runaway, the surface temperature can exceed the critical point of water (374°C), *i.e.*, the temperature where liquid water and steam have the same density and are not distinguishable from each other. Due to the extreme greenhouse warming caused by an ocean mass of water vapor, the surface temperatures on an Earth-like planet can reach the melting points of rocks. In comparison, the moist greenhouse is gentle and occurs when the partial pressure of water vapor at high altitudes becomes sufficiently elevated so that a substantial flux of water can be transported into the stratosphere and beyond. At high altitudes, H_2O is decomposed by UV photolysis and the liberated hydrogen ultimately escapes to space.

The outer edge of the habitable zone occurs when surface water freezes. A commonly quoted limit is 1.37 AU based on the onset of formation of CO_2 ice clouds. A more extended limit of 1.67 AU is based on the maximum greenhouse warming that could occur in a cloud-free CO_2 - H_2O atmosphere. The highest estimate and perhaps an upper limit is 2.4 AU based on a combination of cloud altitudes and particle sizes that could optimize radiative warming by CO_2 clouds [...]

For planets, the conventional habitable zone moves outward with time as their central stars brighten. Typical stars brighten by a factor of ~ 2.5 during their main-sequence lifetimes, the periods of their lives when they are stable stars fusing hydrogen to form helium. Main sequence stars of all mass brighten by a similar fraction as the ratio of He/H in their cores increases with time. At present, the Sun is nearing half its main-sequence lifetime and it is brightening at a rate of about 10% per billion years, and is currently about 30% more luminous than it was 4 billion years ago. More massive and less massive stars brighten at higher and lower rates proportionate to their total main-sequence lifetimes (cf., Fig. 10.1) [...]

The habitable zone concept becomes more complex when the ability to have photosynthesis is considered. A more restrictive consideration of surface habitability by organisms similar to plants and animals is the photosynthetic Habitable Zone or pHZ. Photosynthesis requires atmospheric levels of CO_2 above some critical limit, approximately 10 ppm for known plants. The pHZ of a given star (see Fig. 12.2 for the case of the Sun) narrows over time as the star gets brighter. The inner edge moves outwards and the outer edge

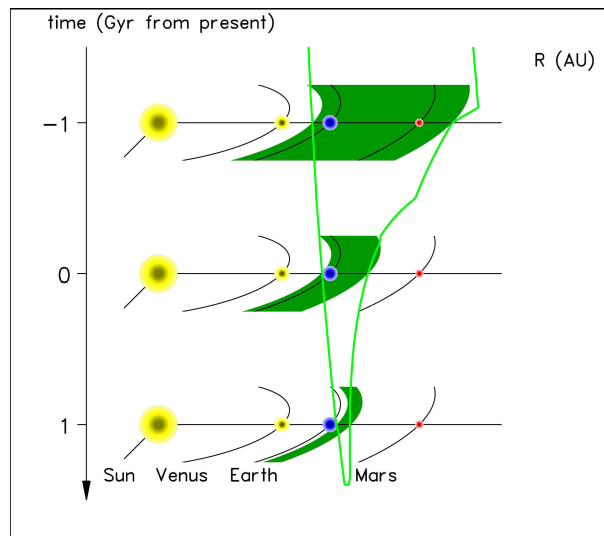


Fig. 12.2. The photosynthetic habitable zone (pHZ) over time, from 1 Gyr in the past to 1 Gyr in the future. The inner edge of the pHZ moves outwards as the Sun becomes brighter with age and the outer edge moves inwards as surface warming leads to decline of CO_2 in the atmosphere to the point where photosynthesis is not possible. [Fig. H-III:4.1; modified from source: Franck *et al.* (2001); also here on the web.]

moves inwards. For a planet with land and surface water, weathering processes remove CO_2 from the atmosphere. The process involves sequestering CO_2 in carbonates and this becomes increasingly more effective as [brightening] stars produce warmer planetary surfaces. This process can cause the pHZ to shrink to zero. Estimates indicate that the Sun's pHZ will shrink to zero width when the Sun reaches an age of 6.5 billion years. The Earth, even now close to the inner edge of the habitable zone, will be left behind the moving pHZ, and lose most of its surface water, long before this time [...]” {A:[162]}

A:162

12.1.3 Oxygen, methane, and carbon dioxide over time

[H-III:4.7] “Prior to 2.4 billion years ago, the Earth’s atmosphere was essentially devoid of free oxygen. Although it was being produced by photosynthetic organisms such as cyanobacteria as well as the photolysis of water vapor, it was efficiently removed from the atmosphere as it oxidized compounds on the surface and in the atmosphere. Before this time, the atmosphere was dominated

¹⁶² Activity: **For the curious:** Photosynthesis depends on the chemicals involved and as such is sensitive to the spectral energy distribution of the star. You could search the literature on developments in this area, but for stars substantially different from our Sun that work remains hypothetical. Here is a possible entry point. Look up where the main absorption bands of chlorophyll and β -carotene lie relative to the solar spectrum at sea level. How does the solar spectrum change under water for, for example, flora in the oceans?

by N₂ but it contained appreciable amounts of CO₂, water vapor and probably moderate amounts of CH₄, possibly up to the percent level. There is abundant evidence for low oxygen abundance on the early Earth including the oxidation state of various minerals, including iron oxide [...]

The crossover, *i.e.*, the appearance of oxygen and simultaneous loss of methane, occurred 2.4 billion years ago. At this time the Earth entered a severe ice age, also called a 'Snowball Earth' episode, during which the planet surface cooled to the point where ice formed at equatorial latitudes. It seems likely that this unusual cooling event was related to the rapid loss of significant greenhouse warming previously associated with the presence of methane.”
 {A:^[163]}

A:163

[H-III:4.9] “Photosynthesis is the primary means by which life on Earth derives energy from the Sun. The complex chemical processes involved with photosynthesis depend on the availability CO₂ in the atmosphere and CO₂ can be considered an essential 'food of life' on our planet. CO₂ on the present Earth is controlled by biogeochemical processes but in the future, as the Sun becomes brighter, the atmospheric CO₂ abundance will decline below the minimum (~10 ppm) amount needed to support plants. The end of CO₂ will mark the end of plants and animals that depend on direct contact with Earth's natural atmosphere [...]

We currently have major concerns with the CO₂ increase from burning fossil fuels and its global warming effects. However, this is a short-term problem. Ultimately, all of the atmospheric CO₂ will become locked up in carbonates and removed from the atmosphere. Even now, most of the CO₂ that has ever been in the atmosphere is already in carbonates. CO₂ is the dominant gas in the atmospheres of Venus and Mars and it must have been a major gas in the Earth's atmosphere before it declined due to carbonate formation. If Earth's total carbonate content were decomposed, it would yield over 20 atmospheres of CO₂, over 4×10^4 times the present CO₂ content of the atmosphere. As the Sun gets brighter, as all stars do as the hydrogen content of their cores is consumed, the Earth's surface temperature will increase and the CO₂ will decline as more and more is sequestered into carbonates. The removal process is related to weathering of rocks, a process whose rate increases with increasing temperature. The presence of silicates, water and atmospheric CO₂, leads to the formation of carbonates. Presently, this process is dominated by biological processes such as the formation of shells, corals, and microscopic organisms such as foraminifera. {A:^[164]}

A:164

When atmospheric CO₂ is sufficiently depleted, Earth will have lost an important factor that has promoted the long-term stability of its surface

¹⁶³ Activity: Sect. 10.3.1 describes the possibility of the Solar System moving through dense, cold interstellar clouds, which could greatly enhance the dust environment of Earth. Review the study

temperature. Over Earth's history, the abundance of carbon dioxide and its greenhouse warming effects have varied in ways that have counteracted changes in atmospheric temperature. When Earth cools over long time periods, the CO₂ abundance can rise and promote greenhouse warming. When Earth warms, the CO₂ abundance can decline and promote cooling. This effect is called the carbonate-silicate cycle and it is a case of negative feedback where change is resisted leading to stability. [...] Carbon is removed by weathering but is involved in a cycle because it is ultimately reintroduced back into the atmosphere. Carbonate deposits in the ocean floor are subducted beneath continents on ~100 Myr time scale where they are thermally decomposed and release CO₂ back into the atmosphere via volcanism. The CO₂ sink depends on weathering and carbonate deposition and the CO₂ source depends on subduction, an ongoing process associated with plate tectonics.”

12.1.4 Water over time

It appears that much of the water on the terrestrial planets may have been transported to the inner parts of the Solar System frozen within asteroids that were scattered from further out during phases of orbital changes of the giant planets, and then to Earth in collisions. Venus and Mars have lost their oceans a long time ago, as discussed in Sect. 12.4.1. Earth, too, will eventually lose the bulk of its water [H-III:4.10] “when a critical threshold brightness is reached [in the Sun's evolution]. Ocean loss is a drastic change for a planet, and for Earth it will mean a change to a seemingly 'unearth-like' state, a planet more like Mars than the blue planet of its past. [...] Even without oceans, Earth will probably always have regional ponds or lakes fed by water derived from the mantle. The mantle is a reservoir that may contain several ocean-masses of water.

The most likely fate of Earth's oceans is loss by the 'moist greenhouse' effect, a process that occurs at present but at a very low rate. In this process, water is transported through the troposphere and stratosphere to heights where its hydrogen can be liberated by photolysis with solar UV photons

by Pavlov *et al.* (2005) for the potential effects on terrestrial climate, including periods of strong glaciation and potentially the triggering of a 'Snowball Earth' state.

¹⁶⁴ Activity: Consider the evolving CO₂ content of the atmosphere of a lifeless terrestrial planet. Which of the following parameters would influence the atmospheric CO₂ content over time: (1) atmospheric mass and composition, (2) chemical composition of seas and oceans, (3) continent sizes and placement, (4) fractional coverage by liquids in seas and oceans, (5) motion through, and density of, local interstellar medium, (6) orbital obliquity, (7) orbital period (length of the planetary 'year'), (8) planetary mass, (9) planetary radius, (10) planetary spin obliquity, (11) planetary spin rate (length of the planetary 'day'), (12) planets elsewhere in the planetary system, (13) plate tectonics, (14) properties of moons, (15) spectral type of the central star, (16) stellar spin rate. Formulate your arguments for each. You may want to read on in Ch. 12 and return here later to complete the activity.

[(Sect. 12.4.1)]. Near the exosphere the liberated hydrogen escapes to space, and forms Earth's L geocorona. {A:[165]} This process currently occurs at a rate of only a meter of ocean in a billion years due to the very low abundance of water vapor in the stratosphere. As the Sun warms, the partial pressure of water in the upper atmosphere rises and the timescale for water loss shortens. Modeling of this process indicates that the moist greenhouse effect will begin severely depleting the Earth's oceans in about a billion years or less. If surface water is not largely depleted by the rather gentle moist greenhouse process in roughly 3 billion years, a much more severe process will take over when the Sun is about 35% brighter than it is at present (Fig. 10.2, also Fig. 10.1). In a runaway, increasing temperatures introduce more greenhouse gas thereby providing positive feedback. This full runaway greenhouse advances to the critical point of water where density of water vapor equals the density of liquid water. In a runaway, the enormous amount of water vapor in the atmosphere produces greenhouse warming sufficient to melt surface rocks. Either the moist-greenhouse or the runaway-greenhouse process will result in the Earth's loss of its oceans to space and our planet will spend over half of its total life as an ocean-free planet, at least initially covered with salt and very oxidized rocks. [...]

The loss of oceans is likely to also lead to the end of plate tectonics. Hydrated minerals have lower melting points and in several ways the presence of water promotes the sinking of oceanic crust to subduct beneath continents. Without oceans it is expected that plate movement will stop and Earth – like all other planets in the Solar System – will cease to have subduction and the drift of continents. Without subduction, the Earth's major mechanism for cycling CO₂ back into the atmosphere will be lost.” When plate tectonics stops, this may also have major consequences for the planetary dynamo. {A:[166]}

12.2 Atmospheres and climates of Venus, Earth, and Mars

In this volume, we focus on the terrestrial planets Venus, Earth, and Mars [H-IV:7.1] “because they are thought to have been habitable at their surfaces at some point during Solar System history. They formed under similar conditions, with early atmospheres that were more similar than they are today. The present day climates of Venus and Mars provide a useful contrast to that of Earth, and exploration of the root causes for differences in the present climates of all three planets allows us to better understand the processes that control climate on terrestrial exoplanets. Their current climates are summarized in Table 2.1 [...]

¹⁶⁵ Activity: Look up what constitutes the geocorona.

¹⁶⁶ Activity: What role does plate tectonics likely play in dynamos in terrestrial planets? Reminder: Sect. 4.1.1.

A:167

Despite their large differences in mass, the atmospheres of Venus and Mars have similar bulk compositions, with carbon dioxide (CO₂) comprising ~95% by volume, followed by molecular nitrogen (~3%) and argon (~1%). Earth's atmosphere, by contrast, is composed mainly of nitrogen and oxygen, followed by argon. {A:[167]} Earth's atmospheric composition likely mirrored that of Venus and Mars early on, but much of Earth's atmospheric CO₂ now resides in carbonates on the ocean floors, leaving nitrogen as the most common constituent. Earth's abundant atmospheric oxygen is believed to have been contributed by photosynthetic bacteria.

The surface temperatures of the three planets also differ widely, in part due to the distance of each planet from the Sun and in part due to the quantity of greenhouse gases in each atmosphere. Earth is the only of the three planets with a surface temperature (and pressure) appropriate for liquid water to be stable for long periods of time, thanks to ~30 K of greenhouse warming. The Cytherean atmosphere is too hot for water to exist as liquid at the surface, while the Martian atmosphere has too low a surface pressure (liquid water would sublime, except at the lowest elevations). The atmosphere of Venus is very dry, indicating that any surface water driven into the atmosphere by the high temperatures no longer resides there. The atmospheric water content at Mars is an order of magnitude larger than at Venus and, given the low atmospheric pressure, is often nearly saturated. Despite the near 100% Martian relative humidity, Earth still has roughly 50 times more water molecules (per number of particles of atmosphere) than Mars. The composition, temperature, and water content lead to different forms of precipitation on the three planets. Earth has a variety of forms of water precipitation, while Mars has carbon dioxide and water frost. Venus has no precipitation at the surface due to its high temperatures; any precipitation that forms higher in the atmosphere would turn to vapor before reaching the ground.

Circulation patterns on the three planets also differ. Earth possesses three circulation cells in each hemisphere, leading to prevailing winds organized by latitude. The circulation results, in a simplified sense, from an equator-to-pole temperature gradient that causes warm air to rise at the equator and fall at the poles. Earth's rotation provides a Coriolis influence that breaks the circulation cells into three regions, keeping the warmest air relatively confined at low latitudes. Venus, by contrast, rotates very slowly. Thus, heat is transferred efficiently from the equator to polar regions, leading to uniform surface temperatures as a function of latitude and local time. Mars rotates at nearly the same rate as Earth but has only one circulation cell per hemisphere, though there are some arguments to suggest that while there is a net circulation,

¹⁶⁷ Activity: The Earth's argon is predominantly Argon-40, whereas that in the universe at large, as in the Sun, is Argon-36. What is the source of Argon-40 in Earth's atmosphere?

air tends to move in localized regional cells. Air at the surface of Mars moves sufficiently quickly to drive dust devil activity, while the surface of Venus is very still. At higher altitudes on Venus, however, the atmosphere super-rotates on timescales of days.

While Earth's seasonal variations, caused by a 23.5° tilt relative to its orbital plane, will be well known to the reader, seasonal variations on Venus and Mars are substantially different. Venus has nearly no seasonal variation due to a very small ($\sim 3^\circ$) axis tilt. Mars has a tilt of 25° , similar to that of Earth, but the planet's greater orbital eccentricity (a 21% difference between the perihelion and aphelion distances compared to 1.4% and 3.3% for Venus and Earth, respectively) leads to shorter and more intense summers in the southern hemisphere compared to the north. Strong heating during southern summer drives enhanced dust devil activity, which can couple across circulation cell boundaries and grow into planet-encompassing dust storms that last several weeks."

[H-IV:7.2] "[A]bundant evidence points to changes in the climate of all three terrestrial planets on a variety of timescales. Here, we focus on evidence for climate change over tens of thousands of years or longer. [...] The most compelling evidence for climate change on Venus comes from measurements of the isotopes deuterium and hydrogen in the atmosphere today. Deuterium is far scarcer than hydrogen in the atmospheres of all planets. However, the ratio of deuterium to hydrogen (D/H) in the Venus atmosphere – about 2 deuterium atoms for every 100 hydrogen atoms – is more than 100 times the same ratio calculated for Earth and most other Solar System objects. There is little reason to expect that Venus formed with a D/H ratio significantly different from that of Earth, so we infer that the D/H ratio on Venus increased after the planet formed. Specifically, it is thought that hydrogen atoms (possibly from a primordial ocean ^[xxv]) preferentially escaped the planet's gravity compared to deuterium and were lost to space [...] – water was dissociated in the atmosphere and the hydrogen removed to space. [...]

[E]vidence for climate change on Earth is abundant and comes in many different forms. [...] The terrestrial climate record [derived from a diversity of sources (including growth rates of tree rings and corals, isotope ratios, gases trapped in air pockets in ice, geochemistry, fossils and sediments)] suggests that Earth's climate varies on many timescales, with departures in temperature of as much as $10\text{--}15^\circ\text{C}$ over Earth's history. There are many inferred cold (glaciation) and warm periods that have been tied with changes in atmospheric conditions and diversity of life. Similarly, there are a few major changes in

^{xxv} Alternatives to a primordial ocean for Venus include a more recent reservoir, or perhaps one that continues to be replenished, from volcanic outgassing – possibly clustered in major events – or cometary supplies – see work cited in this review by Marcq *et al.* (2018)

atmospheric composition, the most notable of which is the oxygenation of the terrestrial atmosphere more than two billion years ago, likely caused by the rise of oxygen-producing bacteria and the subsequent depletion of sinks for oxygen at Earth's surface. Analysis of the size and depth of fossilized raindrop imprints in sedimentary rock even suggests that Earth's surface pressure has varied by as much as a factor of two over 2.7 billion years. Taken together, the evidence provides a caution against interpreting the present day climates of other terrestrial planets too finely, and assuming only monotonic changes in planetary climates over billion year timescales. At the same time, one of the most notable aspects of the terrestrial record is the fact that water has existed as liquid at the surface for most of the planet's history, suggesting that despite short term deviations Earth's climate has been relatively stable over its history, in likely contrast to Venus and Mars.

Mars also provides several lines of evidence suggesting past climate that differs from today. [...] These include dry dendritic (branching) river valley networks, river delta deposits, possible regions of sedimentary rock, smoothed and rounded rocks imaged by Mars rovers, and possible ancient ocean shorelines. These features all suggest an ancient Mars where liquid water was abundant and active in shaping the surface of the planet. Further, highly eroded crater rims and a paucity of small craters relative to what might be expected from the abundance of large craters suggest that the ancient atmosphere was much more efficient at eroding surface features (*i.e.*, thicker) than today – perhaps as thick as 0.5-3 bars, or even more. [...] A] number of Martian atmospheric isotope ratios (D/H, $^{38}\text{Ar}/^{36}\text{Ar}$, $^{13}\text{C}/^{12}\text{C}$, $^{15}\text{N}/^{14}\text{N}$, $^{18}\text{O}/^{16}\text{O}$) point to the stripping of atmospheric particles to space over billions of years, similar to the inference drawn from D/H measurements at Venus. [...] T]he isotope ratios suggest that 50-90% of the total atmospheric content has been removed to space from stripping processes alone.”

12.3 Irradiance, orbits, spin, and climate

12.3.1 Atmospheric effects and albedo

In this section we first look into equilibrium temperatures in the absence of a planetary atmosphere, and then proceed to see how an atmosphere modifies such an equilibrium. We will focus on planets orbiting our Sun, but the same arguments hold for exoplanets orbiting other stars, of course. [H-III:11.2.3] “The fraction of the solar luminosity L_{\odot} that is absorbed by a planet is given by the ratio of the planet's cross section πR_p^2 to the area $4\pi d_p^2$ of a sphere containing the planet at distance d_p from the Sun, corrected for the albedo a

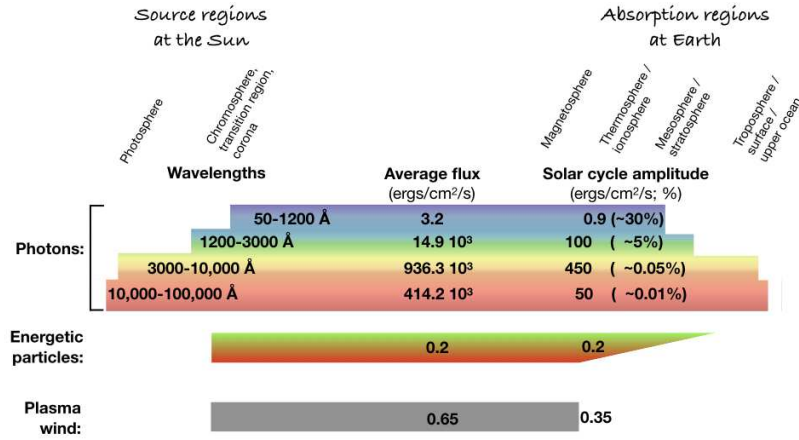


Fig. 12.3. The flow of energy from the Sun to the Earth is compared for photons in four different wavelength bands, energetic particles, and the plasma wind. The numbers are approximate energies with their variations during an 11-yr solar activity cycle, in $\text{erg}/\text{cm}^2/\text{s}$. Visible radiation connects the surfaces of the Sun and Earth while ultraviolet radiation connects their atmospheres. Particles and plasma connect the outer solar atmosphere primarily with Earth's magnetosphere and high-latitude upper atmosphere. [After Fig. H-III:10.3]

([the fraction of total incoming power that is reflected]):

$$\mathcal{P}_{\text{p|a}} = L_{\odot} \frac{\pi R_{\text{p}}^2}{4\pi d_{\text{p}}^2} (1 - a) \equiv \mathcal{P}_{\odot} (1 - a) \quad (12.1)$$

[where \mathcal{P}_{\odot} is defined as the total energy per unit time intercepted by the planetary disk. ...] [H-III:11.3.4] “The albedo is defined as the ratio of diffusely reflected to incident electromagnetic radiation and, therefore, lies in the interval 0 – 1. It is difficult to determine the total albedo of a planet because it is highly variable, ranging from less than 0.1 for water and forests to more than 0.8 for fresh snow. On Earth, the largest contribution comes from the clouds which cover about 50% of its surface. For the Earth an average albedo of 0.3 is usually assumed [...].”

[H-III:11.2.3] “If we assume as a first approximation that a planet is an atmosphere-free black body and that the climate machine distributes the incoming solar radiation uniformly [(i.e., that effects of very low or very high spin rates can be ignored)], the emitted power is given by the law of Stefan-Boltzmann:

$$\mathcal{P}_{\text{emi}} = 4\pi R_{\text{p}}^2 \sigma T_{\text{e|0}}^4. \quad (12.2)$$

Under steady state conditions absorption and emission are equal and the temperature $T_{\text{e|0}}$ [(with index 0 to indicate absence of an atmosphere as we

Table 12.1. Comparison of the calculated temperatures of the planets for different combinations of planetary albedo a and stellar luminosity L in the absence of atmospheres, compared with the observed temperatures.

[Table H-IV:11.3]

Planet	Distance (AU)	Effective temperature absent an atmosphere (°C)			Observed temperature (°C)
		$a = 0.5$ $L = 0.8$	$a = 0.3$ $L = L_{\odot}$	$a = 0.1$ $L = 1.3L_{\odot}$	
Mercury	0.38	77	130	175	180 to 420
Venus	0.72	-10	30	66	460
Earth	1	-50	-18	11	15
Mars	1.52	-95	-65	-40	-87 to 5
Jupiter	5.2	-175	-160	-150	-130
Saturn	9.54	-200	-190	-180	-180
Uranus	19.18	-220	-215	-210	-210
Neptune	30.06	-230	-225	-220	-210

have here)] can be calculated:

$$T_{e|0} = \left(\frac{L_{\odot}(1-a)}{16\pi\sigma d_p^2} \right)^{1/4}. \quad (12.3)$$

Note that the temperature of a planet does not depend on its size [...]

A:168

A:169

{A:[168]} {A:[169]}
 In Table 12.1 the calculated equilibrium temperatures for the eight planets in the absence of atmospheres are compared to the measured ones. [...] Overall there is a reasonable agreement between the estimated and the observed temperatures. The largest discrepancy is observed for Venus. The reason is that Venus has a very dense atmosphere which consists for 96% of CO₂ [(see Table 2.1)] with clouds of SO₂ generating the strongest greenhouse effect in the Solar System. In the case of Earth, the difference between calculated (using the present values $a = 0.3$ and solar luminosity) and measured mean global temperature is 33°C. This difference is also due to the natural greenhouse effect. It is important to note that the Earth needs the natural greenhouse effect to be habitable, but not necessarily an additional anthropogenic increase.

¹⁶⁸ Activity: At what distance would an Earth-equivalent exoplanet need to orbit an $0.6 M_{\odot}$ M0 V star to reach the same global 'equilibrium temperature', all other things being equal? You may disregard effects associated with the difference in the stellar spectral energy distribution on the exoplanet, but you should not ignore the bolometric correction in estimating the total stellar irradiance. How long would a year last on such a planet compared to Earth's? Use Fig. 4.2. Note: such close-in planets are subject to very strong tidal forces that will synchronize spin and orbital periods, causing these exoplanets to lose their day-night cycles. That, in turn, invalidates your estimate – why?

¹⁶⁹ Activity: Beyond the furthest planet: The New Horizons spacecraft flew by Kuiper Belt Object 2014 MU₆₉ on 2019/01/01, the most distant body visited by a spacecraft to date, at an orbital distance of ~ 44 AU. Estimate the surface temperature of 2014 MU₆₉, which has an albedo of ~ 0.1 . Compare your estimate to the observed temperature in this paper by Stern *et al.* (2019).

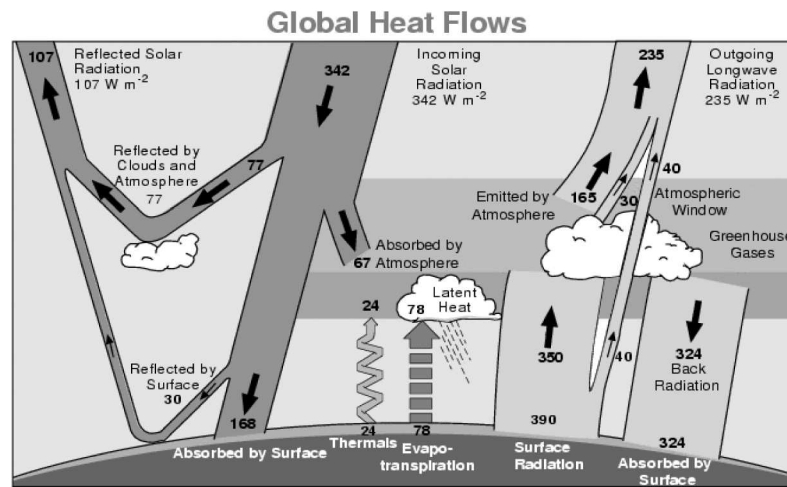


Fig. 12.4. Exchanges of solar (shortwave) and terrestrial (longwave) energy in the Earth's atmosphere. The flow of energy is expressed in W m^{-2} ($\text{kiloerg/cm}^2/\text{s}$), averaged over the entire Earth surface (i.e., over the day-night cycle). [Fig. H-III:16.2; source: Kiehl and Trenberth (1997).]

The range of observed temperatures on Mars is very large because Mars has only a very thin atmosphere (0.3 millibar compared to 1 bar of Earth) and no liquid water to transport and distribute energy. Jupiter is considerably warmer than calculated (-110°C instead of -160°C). Most likely, this difference is due to gravitational contraction which provides an additional power at least as large as the solar insolation."

The value of the planetary albedo is determined by the properties of the planetary surface and, if present, the planetary atmosphere. For the Earth, not surprisingly, the impact of the atmosphere in setting the overall albedo has been studied in great detail. [H-III:16.2] "Perhaps, the best way to represent the exchanges of radiative energy in the atmosphere is to refer to Figure 12.4. This figure shows that [... a large fraction of the infrared radiation from the planetary surface] is absorbed by greenhouse gases in the atmosphere. These gases, whose temperature is lower than the surface temperature, re-emit radiation both towards space and towards the Earth's surface [...]"

For an illustrative first-order approximation of the sensitivity of the ground-level climate to the greenhouse effect of an atmosphere, we can look at a highly simplified version of the energy flow in which we disregard the mechanical and evapo-transpiration energies and assume that the atmosphere radiates equally in the upward and downward directions (which is a significant oversimplification as you can infer from Fig. 12.4), as sketched in Fig. 12.5. [H-IV:7.3] "[The

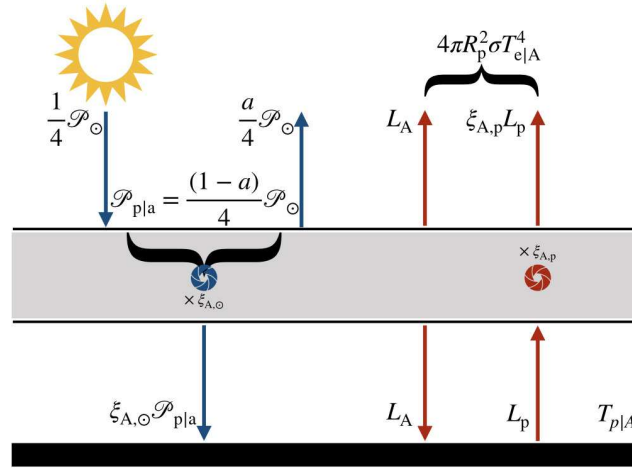


Fig. 12.5. Simplified model of radiative exchange in the [atmosphere of a planet with a partially transparent atmosphere that balances incoming solar power \mathcal{P} and outgoing planetary luminosity L]. The optical transmission in the atmosphere is represented by $\xi_{A,\odot}$ in the shortwave ([solar, shown in blue]) spectral region and by $\xi_{A,p}$ in the longwave ([planetary, shown in red]) part of the spectrum. [This figure is a modified version of Fig. H-III:16.3 for consistency with symbols used in this text.]

e]ffective temperature $T_{e|A}$ can be related to surface temperature of a planet with an atmosphere under a few assumptions. Here, [one may approximate the planetary surface temperature in the presence of an atmosphere by]

$$T_{p|A} = (1 + \tau_A)^{1/4} T_{e|A}, \quad (12.4)$$

where $T_{p|A}$ is the [ground-level or] surface temperature and τ_A is the optical depth of the atmosphere.” Although proper radiative transfer is essential to quantify τ_A and thereby how atmospheric properties combine to set the ultimate planetary temperature, let us here look at a very [H-III:16.2] “simple model of the radiative transfer processes described above [but allowing for a simple wavelength-dependent effect in radiative transport]; we represent the atmosphere by a single layer of radiatively active gases whose optical transmission is noted $\xi_{A,\odot}$ and $\xi_{A,p}$ for shortwave (incoming solar) and longwave ([emitted planetary]) radiation, respectively. The radiative shortwave solar [power] and the longwave surface [luminosity] are noted by symbols $\mathcal{P}_{p|a}$ and L_p , respectively. L_A represents the radiative [power] emitted by the atmospheric layer and $T_{p|A}$, an indicator of the [planet’s] climate, represents the surface (ground) temperature. From Figure 12.5, we derive the energy balance at the top of the atmosphere [and at the planets’s surface, respectively:]

$$\mathcal{P}_{p|a} = (1 - a)\mathcal{P}_{\odot}/4 = \xi_{A,p}L_p + L_A ; \quad \xi_{A,\odot}\mathcal{P}_{p|a} = L_p - L_A. \quad (12.5)$$

We deduce that the surface temperature [in the presence of an atmosphere] is given by

$$T_{\text{p|A}} = T_{\text{e|0}} \left(\frac{1 + \xi_{\text{A},\odot}}{1 + \xi_{\text{A},\text{p}}} \right)^{1/4}, \quad (12.6)$$

where the planetary equilibrium temperature $T_{\text{e|0}}$ [is given by Eq. (12.3). For Earth, this] is equal to 255 K (-18°C) for $\mathcal{P}_\odot/4 = 342 \text{ Wm}^{-2}$ and for an albedo $a = 0.31$. Assuming that the atmosphere is approximately transparent to solar radiation, so that the shortwave transmission $\xi_{\text{A},\odot}$ is close to 1.0, and adopting a longwave transmission $\xi_{\text{A},\text{p}}$ of 0.2, the surface temperature [in the presence of its atmosphere] becomes

$$T_{\text{p|A}} = T_{\text{e|0}} \left(\frac{2.0}{1.2} \right)^{1/4} = 289 \text{ K [or } 16^\circ \text{ C]}. \quad (12.7)$$

The value calculated by this simple model, tuned by approximating choices for $\xi_{\text{A},\odot}$ and $\xi_{\text{A},\text{p}}$ [(and suggesting an effective atmospheric optical depth in Eq. 12.4 of $\tau_{\text{A}} \approx 0.67$)], is in agreement with the observed temperature $T_{\text{p|A,obs}}$ (288 K). More refined models account in greater detail for wavelength-dependent radiative transfer, vertical and horizontal heat transport in the atmosphere, energy and water exchanges at the Earth's surface. Absorption coefficients for different molecules in different spectral regions are measured in the laboratory. [...] The simple conceptual model presented here can, however, be used to estimate to a first approximation the change in the surface temperature that would result, for example from a relative change in the solar input \mathcal{P}_\odot of 0.1%. We derive easily that, for constant $\xi_{\text{A},\odot}$ and $\xi_{\text{A},\text{p}}$,

$$\frac{\Delta T_{\text{p|A}}}{T_{\text{p|A}}} = \frac{\Delta \mathcal{P}_\odot}{4 \mathcal{P}_\odot}. \quad (12.8)$$

For $T_{\text{p|A}} = 288 \text{ K}$, we obtain a surface temperature change $\Delta T_{\text{p|A}}$ of 0.07 K for a solar-cycle TSI variation of $1500 \text{ erg/cm}^{-2}/\text{s}$. The amplitude of the solar variation is therefore a factor of 10 smaller than the surface temperature trend observed since the beginning of the industrial era. However, over a period of a decade or so, the solar signal should be significant compared to human-driven temperature trends, and should therefore be taken into consideration in the analysis of temperature records. [Studies] have shown that, even if the global temperature variation associated with solar forcing is small, changes in temperature patterns become significant at the regional scale.

A more accurate treatment requires that the transmission functions and the atmospheric emissivity change with the chemical composition of the atmosphere in response to Sun-induced climatic changes, that dynamical feedbacks be taken into account and that the influence of the ocean be considered. [...]" Many

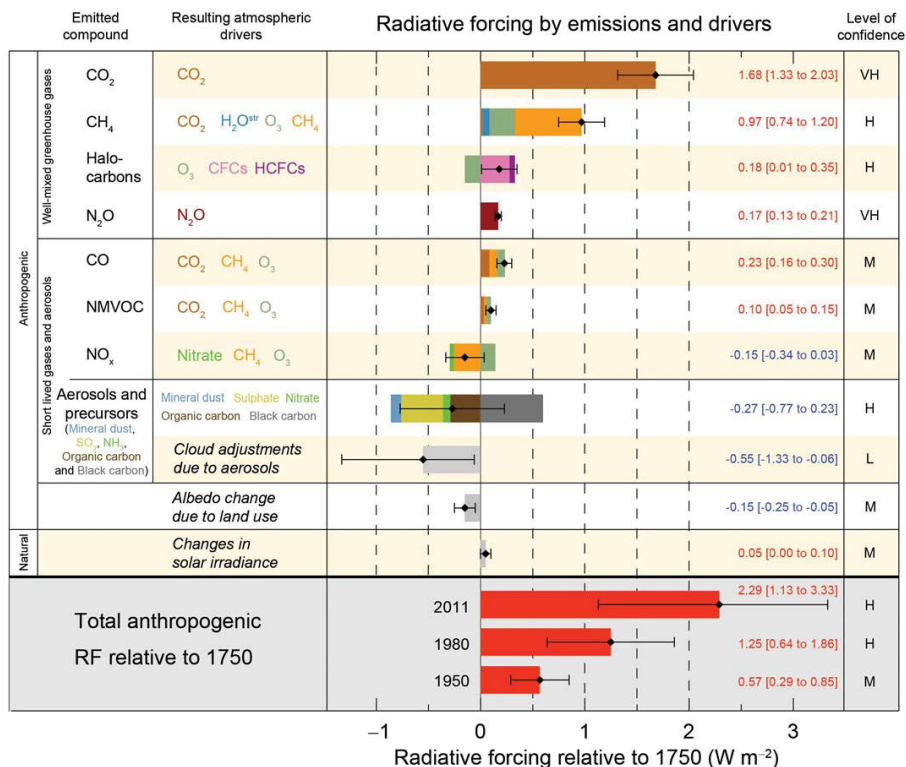


Fig. 12.6. Radiative forcing estimates in 2011 relative to 1750 and aggregated uncertainties for the main drivers of climate change. Values are global average radiative forcing (RF), partitioned according to the emitted compounds or processes that result in a combination of drivers. The best estimates of the net RF are shown as black diamonds with corresponding uncertainty intervals; the numerical values are provided on the right, together with the confidence level in the net forcing (VH - very high, H - high, M - medium, L - low, VL - very low). Albedo forcing due to black carbon on snow and ice is included in the black carbon aerosol bar. Small forcings due to contrails (0.05 W m^{-2} , including contrail induced cirrus), and HFCs, PFCs and SF₆ (total 0.03 W m^{-2}) are not shown. Concentration-based RFs for gases can be obtained by summing the like-colored bars. Volcanic forcing is not included. Total anthropogenic radiative forcing is provided for three different years relative to 1750. [From IPCC (2013).]

details go into establishing the 'radiative forcing' of an atmosphere, see, for example, the publications of the Intergovernmental Panel on Climate Change (IPCC) from which Figure 12.6 was taken to illustrate the often counteracting effects of atmospheric constituents on radiative forcing. {A:^[170]}

A:170

¹⁷⁰ Activity: Show that the simple model in Fig. 12.5 yields an estimate consistent with Earth's global temperature rise of about one degree (observed between 1850 and 2010) based on the increase in anthropogenic radiative forcing as shown in Figure 12.6 within the uncertainty indicated in that figure.

Equation (12.3) illustrates [H-IV:7.3] “four main ways in which planetary climate can be altered. First, the amount of radiation from the star (L_{\odot}) can change. The solar constant at Earth varies by only $\sim 0.1\%$ over the course of a solar cycle. [Evolutionary stellar-structure models] suggest that the Sun is $\sim 30\%$ brighter today than it was when the terrestrial planet[s formed (Fig. 10.1) ...]

Second, changes in the albedo (a) of a planet will change the amount of incident energy absorbed by the surface (and atmosphere). Variation in cloud cover, the extent of polar ices, vegetation, or wind blown dust, for example, can all change the albedos of the terrestrial planets, and will have an influence on the atmospheric energy budget. Venus has an albedo of ~ 0.9 , while the albedos of Earth (~ 0.3) and Mars (~ 0.25) are considerably lower. [...] {A:^{[171]}} A:171

Third, characteristics of a planet’s orbit and rotation influence its energy budget. The amount of solar radiation encountering a planet varies with average orbital distance (d_p), with the result that Venus encounters roughly double the energy that Earth does, while Mars encounters $\sim 45\%$. Ellipticity of the orbit [(see Eq. 12.10) causes incident energy to vary] between 36% and 52% over a Martian year due to Mars’ relatively high orbital ellipticity. This explains why the southern summer at Mars (near perihelion) is more extreme than the northern summer. [Realizing that incident solar energy is not uniformly distributed by the thin atmospheres of the terrestrial planets, it will be clear that t]ilt also influences the amount of sunlight that reaches each part of a planet’s surface, making some portions of the planet cold and other portions warm. This effect influences where ices form at the surface, removing some gases from the atmosphere and changing albedo in some locations. Chaotic changes in the eccentricity, obliquity, and spin precession of Mars and Earth over periods of tens to hundreds of thousands of years are thought to contribute to climate variations (Sect. 12.3.2), though the range of variation in both orbital properties (especially tilt) and climate is estimated to be larger at Mars due to the lack of a large Moon.

Fourth, the amount of radiation-absorbing atmosphere (*i.e.*, greenhouse gases) influences surface temperatures. [...] The thick CO_2 atmosphere of Venus provides more than 500 K of greenhouse warming compared to the theoretical surface temperature in the absence of an atmosphere. Earth’s atmosphere provides approximately 30 K of greenhouse warming. This warming, while much smaller than at Venus, is crucial to keeping our average surface temperature above the freezing point of water, making life and many aspects of our climate possible. The atmosphere of Mars, while dominated by CO_2 , is

¹⁷¹ Activity: Compare the values of \mathcal{P}_{abs} from Eq. (12.1) for Venus and Earth. Explain qualitatively why Venus’ surface temperature exceeds Earth’s, then read on for the answer.

too thin to provide substantial greenhouse warming today. The temperature is warmed only ~ 5 K due to greenhouse gases [...].”

12.3.2 Orbital changes

The physical basis of orbital changes and of tidal effects on planetary rotation were discussed in Ch. 7. Here, we look in some detail at the orbital effects on climate. [H-III:11.3.2] “[T]he distance d_p is a prime parameter for the temperature of a planet. [S]olar power decreases with the square of the distance or in other words that the relative change of the temperature is 1/2 of the relative change of the distance:

$$\frac{\Delta T_{e|0}}{T_{e|0}} = -\frac{1}{2} \frac{\Delta d_p}{d_p}. \quad (12.9)$$

[...] Because all the planets have elliptical orbits the distance is continuously changing. The eccentricity ranges from 0.0068 for Venus to 0.2056 for Mercury. The eccentricity of the Earth’s orbit is 0.017. That means the distance between Earth and Sun is 1.017 AU at the aphelion compared to 0.983 AU at the perihelion. This difference results in a change of insolation by about 10^5 erg/cm²/s” and would result in $\Delta T_{e|0} \approx 5$ K throughout the year, but that is strongly dampened by the thermal inertia of Earth’s oceans and land masses.

But not only does ellipticity of orbits lead to seasonal changes, the orbits actually evolve over time. [H-III:11.3.3] “[A]ll the bodies in the Solar System are gravitationally coupled. This was known already since Newton’s time. [...] It] was Milutin Milankovic who, for the first time, worked out the mathematical details of these disturbances [...] There are three orbital parameters of the Earth which are affected by the other planets, the Sun, and the Moon.

(1) Orbital eccentricity: [...] Integration over a full year] reveals the following relationship between the relative change in the annual amount of solar radiation S received by Earth and the relative change in the eccentricity e :

$$\frac{dS}{S} = \frac{e^2}{(1-e^2)^{3/2}} \frac{de}{e}. \quad (12.10)$$

The largest change in e (0.06) which the Earth experienced over the past million years (Fig. 12.7) therefore leads to a very small change of 0.36 % in the annual mean insolation which corresponds to a mean global forcing of less than 10^3 erg/cm²/s [(compare Fig. 12.6 for present-day forcings)]. The changes in the eccentricity occur on time scales of 100,000 and 400,000 years. It is interesting to note that it is exactly this small change in the eccentricity which seems responsible for the 100,000-year cycle in the sequence of glacial

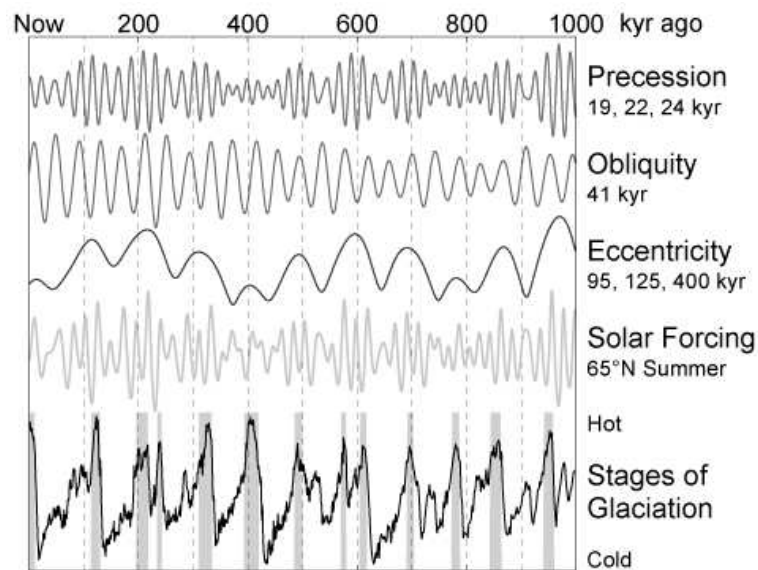


Fig. 12.7. Earth's orbital parameters for the past million years. The first three panels show the three orbital parameters influenced by the other planets (mainly Jupiter and Saturn) and the moon (precession). The fourth panel exhibits the calculated solar forcing at 65° N. The lowermost panel shows sea level changes derived from stable isotope measurements on benthic foraminifera indicating glacial (cold) and interglacial (warm, grey bands) periods. [Fig. H-III:11.6]

and interglacial periods during the past 1,000,000 years (Fig. 12.7). This is a nice example that climate is a non-linear system and that even a small forcing can cause a large effect if feedback mechanisms are involved. Such a feedback mechanism could be that although a larger eccentricity does not change the mean annual insolation much, but with it the seasonality changes: colder summers on the northern hemisphere may result in a reduced melting of the winter snow enlarging the ice sheets and the albedo which further reduces the effective insolation.

(2) Obliquity: The tilt angle of the Earth's spin axis relative to the ecliptic plane varies between 22.1° and 24.5° with a periodicity of about 41,000 years. Contrary to the eccentricity changes the obliquity does not change the total amount of received solar radiation but only its latitudinal distribution. The larger the obliquity the stronger is the seasonality. A smaller obliquity reduces both the mean insolation and the summer insolation at high latitudes, thereby providing favorable conditions for ice ages.

(3) Precession: [...] Because the Earth is spinning, its shape deviates slightly from a sphere leading to an equatorial bulge. Tidal forces act on the

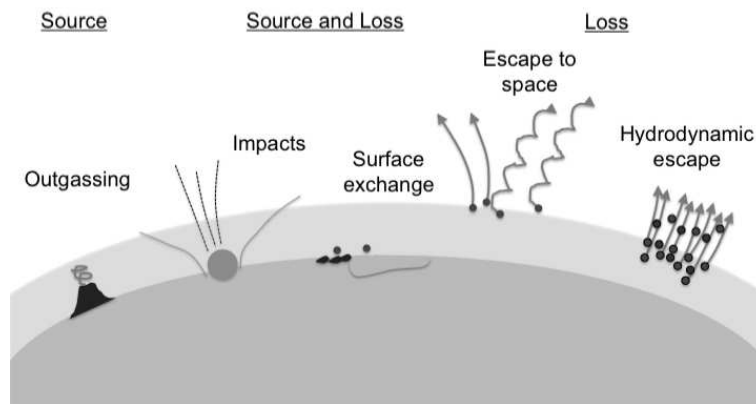


Fig. 12.8. Source and loss mechanisms for planetary atmospheres. [Fig. H-IV:7.3]

bulge and force the [rotation] axis to precess. The periods of precession range from 19,000 to 24,000 years.

The calculated values of the three orbital parameters are plotted in Fig. 12.7, together with the corresponding summer insolation at 65° N, a latitude which is considered as critical for the formation of ice sheets as a result of cold summers. The bottom panel shows a compilation of $\delta^{18}\text{O}$ records from deep-sea sediments. Benthic foraminifera live in the deep sea and form CaCO_3 shells. After death, the shells are buried in the sediment layer by layer for millions of years. Measuring the $^{18}\text{O}/^{16}\text{O}$ isotope ratio with a mass spectrometer relative to a standard, expressed as $\delta^{18}\text{O}$, reflects the sea level. Water evaporating from the sea preferentially contains the lighter molecules H_2^{16}O . If the evaporated water stays on the continents forming glacial ice sheets the ocean becomes depleted in ^{16}O . Warm interglacial periods are indicated by grey bands. They normally last 10,000 to 20,000 years and occur with a typical periodicity of 100,000 years when the eccentricity is large.”

12.4 Planetary atmospheres, geological activity, and stellar winds

12.4.1 On time scales beyond millions of years

[H-IV:7.4] “[Planetary s]urface temperature and climate are strongly affected by the amount of greenhouse gases in an atmosphere, which can be viewed as a combination of the total number of particles in an atmosphere (surface pressure) and its composition. Several mechanisms are capable of changing atmospheric abundance and composition (Fig. 12.8) [...]

Volcanic outgassing from planetary interiors is thought to be the primary

source for the terrestrial planet atmospheres we observe today. Water vapor is the most common gas released in terrestrial eruptions, followed by CO₂. Other commonly released gases include sulfur dioxide, nitrogen, argon, methane, and hydrogen. Outgassing should be a declining source of atmospheric particles over Solar System history, as the interior heat required to generate volcanic activity declines. [Earth is evidently volcanically active today, Venus likely is (although without signatures of active plate tectonics), while there is no direct evidence for ongoing activity for Mars.]

Atoms and molecules can be exchanged between a planet's surface layers and its atmosphere via a variety of processes and over many timescales. For example, changes in temperature can increase condensation rates to the surface, forming surface liquids or ices (evident on Earth and Mars). Chemical reactions (weathering) can also remove particles from the atmosphere, and is typically most effective in warm or wet environments (evident on Venus, Earth, and Mars). Adsorption removes atmospheric particles that stick to surface materials. Most or all of these processes can be considered to be reversible. Release of particles back to the atmosphere can involve changes in temperature, chemical reactions (including reactions with sunlight), and geologic events that allow subsurface reservoirs access to the atmosphere.

All planetary atmospheres are subject to impact from asteroids, comets, dust, and even atoms and molecules. Impactors of all sizes can deliver volatile species to an atmosphere (*e.g.*, impact delivery is responsible for at least part of Earth's water inventory as well as meteoritic layers observed in terrestrial planet ionospheres). Impacts can also remove atmospheric particles via collisions, and sufficiently large impactors can additionally accelerate atmospheric particles via impact vapor plumes and lofted surface material [...]. Monte Carlo simulations suggest impacts have resulted in a net gain of atmospheric gases for Earth and Mars over Solar System history, and a net loss for Venus. {A:^[172]}

A:172

Hydrodynamic escape occurs when a light species escapes (thermally) in sufficient abundance that it becomes equivalent to a net upward wind, and drags heavier species with it through collisions. This process is usually enabled by high solar EUV flux or another form of heating. It should have been significant for all of the terrestrial planets during the first few hundred million years after formation, stripping away most of their primordial atmospheres. [...]

The removal of atmospheric particles to space from the upper layers of the atmosphere is commonly referred to as escape to space. This term typically excludes impacts by asteroids, meteoroids, and comets, and hydrodynamic escape

¹⁷² Activity: To get an idea of scales: estimate the size of a comet that would double the CO₂ content of Earth's atmosphere. How does that compare to, *e.g.* comet 1P, the target of the *Giotto* mission, and 67C, the target of the *Rosetta* mission?

is also often listed as a distinct process. Here, escape to space encompasses a set of approximately six processes, all of which provide escape energy to atmospheric particles. The energy is ultimately provided (sometimes directly, and sometimes indirectly) through interaction with the parent star and stellar wind. [...] It is currently thought that atmospheric escape has played an important role in the evolution of the climates of both Venus and Mars by altering atmospheric pressure and trace gas abundance.”

[H-IV:7.5] “All particles escaping from a planetary atmosphere share three characteristics. The first is that they have sufficient energy to escape the gravity of the planet[, which means that their velocity should exceed the escape speed (Eq. 2.4 with r set to the radial distance from which the escape occurs, typically the exobase, discussed below). The values listed in Table 2.1 show that Mars has a much lower escape speed than Earth or Venus.]

A second characteristic of an escaping particle is that it is unlikely to collide with other particles after acquiring sufficient escape energy. In planetary atmospheres, the region above which collisions are unlikely is termed the exobase, and is loosely defined as the location where the mean free path of a particle [(Eq. 2.17)] is equal to an atmospheric scale height [(Eq. 2.3) ...]

Finally, any escaping particles must not be confined to the planet by planetary magnetic fields. This requires either that an escaping particle be neutral, that the planet lack a magnetic field, or that any magnetic fields are weak enough that energized charged particles are able to easily traverse magnetic field lines. Venus lacks a measurable global magnetic field like that of Earth. Mars also lacks a global magnetic field but possesses localized regions of strongly magnetized crust that may locally trap energized atmospheric ions.

Due to the highly collisional nature of planetary lower atmospheres, escape is generally limited to three regions of the upper atmosphere: the thermosphere, the exosphere, and the ionosphere. The altitude and composition of these regions are summarized for each planet in Table 2.2 [...]”

[H-IV:7.6] “A number of mechanisms are capable of giving atmospheric particles sufficient energy to escape from a planet [(see Fig.12.9.) Neutral particles can escape an atmosphere in one of three ways: (1) *Jeans escape*, (2) *photochemical escape*, and (3) *atmospheric sputtering*. [Ion loss processes can be grouped into three additional] categories: (4) *ion outflow*, (5) *ion pickup*, and (6) *bulk plasma escape*.]

(1) *Jeans (or thermal) escape* occurs because some fraction of neutral particles near the exobase will have sufficient energy to escape simply because the particles have a thermal distribution. Neutral temperatures near the exobase of all three planets are sufficiently low ($\sim 250\text{--}1000$ K) that only species with small

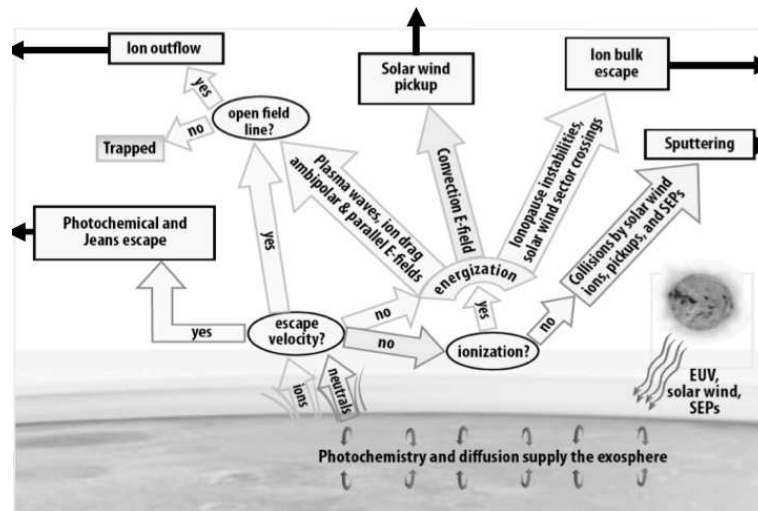


Fig. 12.9. Flowchart showing pathways to energization and escape of particles from a planetary atmosphere. [... Fig. H-IV:7.4]

mass (H, D, and He) can escape via this mechanism in significant quantity. The process should be more efficient for Mars (due to its low gravity) and for Earth (due to its higher exobase temperature) than for Venus.

(2) *Photochemical escape* refers to the escape of fast neutral particles energized by sunlight-driven chemical reactions. These reactions typically involve dissociative recombination of an ionized molecule with a nearby electron, resulting in two fast neutral atoms. Photochemical escape fluxes depend upon ionospheric molecular densities near the exobase, as well as electron density and temperature. Photochemistry is thought to be the dominant loss process for neutral species more massive than hydrogen and helium at Mars. Fast atoms produced photochemically at Venus and Earth are typically not energetic enough to escape the larger gravity.

(3) *Atmospheric sputtering* occurs when atmospheric particles near the exobase receive sufficient energy from collisions to escape. Collisions occur when energetic incident particles (often ionospheric particles accelerated by electric fields near the planet) encounter the exobase. There are no unambiguous observations that sputtering is actively occurring at any of the terrestrial planets, [but] it may have been important earlier in Solar System history, especially for unmagnetized planets. [...]

(4) *Ion outflow* refers to the acceleration of low energy particles out of the ionosphere via plasma heating and outward directed charge separation (ambipolar) electric fields. In this case the ion acceleration can occur below the exobase, where collisions maintain a more fluid-like behavior. Ion outflow is the

only significant ion loss process for the terrestrial atmosphere, and encompasses a number of processes referred to in the terrestrial literature, including wave heating, polar wind, and auroral outflow. [Analogues of these processes should be active for Venus and Mars.]

(5) *Ion pickup* refers to the situation where a neutral particle is ionized (via photons, electron impact, or charge exchange) and accelerated away from the planet by a motional electric field ($\mathbf{E} = -\frac{1}{c}\mathbf{v} \times \mathbf{B}$). Ion pickup occurs primarily for ionized exospheric neutrals (though some ionized thermospheric neutrals near the exobase region may escape via pickup as well). The motional electric field is usually supplied by the solar wind, so that the process is most relevant for compact magnetospheres unshielded by strong planetary magnetic fields (Venus and Mars) [...]

(6) *Bulk plasma escape* refers to any process which removes spatially localized regions of the ionosphere *en masse*. Bulk escape is relevant for unmagnetized planets, where the external plasma flow can create magnetic and/or velocity shear with the ionosphere. A popular example involves the Kelvin-Helmholtz instability, which may form at the ionopause of Venus or Mars and steepen into waves which eventually detach from the ionosphere. Other bulk escape processes are possible as well, such as transport via plasmoid-style flux ropes that may remove ionospheric plasma from Martian crustal magnetic field regions. {A:^[173]}

A:173

[Based on models and observations,] the present-day global escape rate for Venus is estimated to be 10^{24} - 10^{26} s⁻¹. The escape rate for Earth is 10^{25} - 10^{27} s⁻¹, and for Mars is 10^{24} - 10^{26} s⁻¹. [Normalized per unit area, these rates] are on the order of 10^6 - 10^9 cm⁻² s⁻¹. [These escape rates are a very small fraction of the column densities in the present-day atmospheres that range from 10^{23} - 10^{27} cm⁻², but they may be substantial when accumulated over ~ 4 billion years ($\sim 10^{17}$ s)]. For this latter point the two orders of magnitude uncertainty in escape rates are crucial; they are the difference between heliophysical drivers being the main loss mechanism for planetary atmospheres and merely an afterthought in determining present-day [atmospheres. ...]

Finally, it is important to keep in mind that escape to space not only influences atmospheric abundance but also atmospheric composition, which can be important in planetary evolution. One example is the aridity of the Cytherean atmosphere. The loss of atmospheric water is attributed to dissociation of the water in the atmosphere by sunlight, and the subsequent escape to space of oxygen. Water is only a trace gas in planetary atmospheres,

¹⁷³ Activity: What is the basis of the Kelvin-Helmholtz instability? This instability also occurs between the terrestrial magnetosphere and magnetopause flow because the magnetic tension is not strong enough to stabilize the developing waves. Why is this geospace phenomenon not listed as a process for 'bulk outflow'?

Table 12.2. *Properties of the solar wind and interplanetary magnetic field (IMF) at terrestrial planets. [Table H-IV:7.4]*

	Venus	Earth	Mars
IMF strength	~0.10–0.12 mG	~0.06 mG	~0.03 mG
Solar wind speed	~400 km/s	~400 km/s	~400 km/s
Solar wind density	10–15 cm ⁻³	~6 cm ⁻³	~1–3 cm ⁻³
Alfvén speed	~70 km/s	~55 km/s	~45 km/s
Mach number	5–7	6–8	8–10
H ⁺ gyroradius	~1500 km	~2500 km	~5000 km
H ⁺ gyroradius / R_p	0.5	0.4	~3

but is an important greenhouse gas and is extremely important for habitability. So even if escape to space does not appreciably change atmospheric thickness, it may contribute in important ways to climate. Interestingly, the escape rates listed above, when converted to precipitable microns of water, amount to global layers of water only centimeters thick. More than this is assumed to have been lost from Venus, suggesting either that escape rates have changed over time (and are low today) or that other processes (such as impacts) have been important for removing water.”

[H-IV:7.7] “Observations, simulations, and common sense all tell us that atmospheric escape rates are not constant, and are influenced by a number of heliophysical drivers that vary on both short and long timescales. [...] The three main drivers are photons, charged particles, and electromagnetic fields. Photons deposit energy in atmospheres when they are absorbed by atmospheric particles. Extreme ultraviolet (EUV) and soft X-ray photons (generated in the solar corona and chromosphere, and not to be confused with solar luminosity) provide the dominant energy source in upper atmospheric regions. Charged particles in the solar wind also supply energy to planetary upper atmospheres and plasma environments. Table 12.2 summarizes some of the relevant quantities of the solar wind at each terrestrial planet. While density and velocity can each vary independently, studies of solar wind influences on atmospheric escape (especially the induced magnetospheres of Venus and Mars) typically use solar wind pressure (ρv^2) as the organizing quantity. Finally, the solar wind carries a magnetic field, which creates a convection electric field (\mathbf{E}_{sw}) in the frame of the planet that depends upon solar wind velocity and interplanetary magnetic field (IMF) strength and orientation ([see around Eq. 3.11]). Magnetic and electric fields organize charged particle motion, and electric fields accelerate charged particles; both effects influence the ability of charged particles to escape a planet’s atmosphere.

The external drivers of atmospheric escape vary on four main timescales. Billion-year timescales are associated with the age of the Sun, and both theoretical calculations and observations of Sun-type stars suggest that all three drivers should have declined in intensity with age (see Figs. 10.3, 10.10, and 13.5). EUV flux varies by factors of several over a solar cycle (from solar minimum to solar maximum), and solar wind pressure varies by factors of 2 – 10. The IMF, in particular, is a function of the solar rotation period, and all three drivers also vary on more rapid timescales of minutes to hours.

Variability in the heliophysical drivers should influence atmospheric escape rates. In general, an increase in solar EUV fluxes (*e.g.*, a transition from solar minimum to solar maximum) is expected to result in an increase in loss rates of neutral particles. Energy from solar photons heats the upper atmospheric neutrals, so that Jeans escape rates should increase with solar EUV. This is likely to be true at Mars, but not at Earth where hydrogen escape from the exobase is limited not by the available energy, but by the supply (via diffusion) of particles from lower altitudes. Jeans escape should be negligible at Venus today, but may have been significant in the past if either exobase temperatures or solar EUV fluxes were much higher. Energy from solar photons is also used to drive the chemical reactions necessary for photochemical escape, so that contemporary Martian photochemical escape should vary with EUV flux. Neutral escape rates should be largely insensitive to changes in both the solar wind and the IMF, except for sputtering rates from Venus and Mars, which are thought to be dominated by re-impacting atmospheric pickup ions and will therefore increase as the pickup ion population increases in response to changes in solar EUV.

Ion escape rates should also vary with the three drivers. An increase in solar wind pressure will cause a corresponding decrease in the size of the magnetospheric cavity at all terrestrial planets, effectively lowering the pressure balance altitude between the solar wind and planetary obstacle to the flow. For Mars, with an extended neutral corona, an increase in solar wind pressure exposes significant additional high-altitude neutrals to ionization and stripping by the solar wind (via electron impact and charge exchange). The IMF, by contrast, chiefly organizes the trajectories of escaping particles at Venus and Mars; large-gyroradius pickup ions are preferentially accelerated away from the planet in regions where \mathbf{E}_{sw} points away from the planet. At Earth, the orientation of the IMF affects the location and extent of cusp regions, from which outflowing ions escape. EUV fluxes have a more indirect effect. In total, one might expect the ion escape rate to increase at solar maximum due to the additional energy input from EUV. At unmagnetized Venus and Mars, however, the increased ionospheric content deflects the solar wind around the

planet at higher altitudes and can prevent the interplanetary magnetic field from entering the ionosphere. The escape of heavy ion species (which are concentrated at lower ionospheric altitudes) via pickup and bulk escape may therefore remain roughly constant, or even decrease during solar maximum periods, even as lighter ion species escape more efficiently.” {A:[174]}

A:174

[H-IV:7.8] “A number of characteristics of a terrestrial planet itself influence the properties and energetics of upper atmospheric reservoirs for escape, including transient events such as dust storms (*e.g.*, for Mars), or longer-lived phenomena such as gravity waves that couple the lower and upper atmospheres. In the context of heliophysics, the nature of a planet’s intrinsic magnetic field is of the greatest relevance. [... Present-day Earth has an intrinsic magnetosphere] that deflects the solar wind at large distances from the planet ($\sim 10 R_{\oplus}$). There is an induced magnetosphere at Venus that deflects the solar wind at much closer distances ($\sim 1.3 R_{\oplus}$), and a similarly-sized (with respect to the planet) induced magnetosphere at Mars punctuated by ‘mini-magnetospheres’ tied to specific regions of the crust and that rotate with the planet. [...]

When considering the total atmospheric loss from a planet, it has often been assumed that the presence of a magnetic field results in lower escape rates. [However, we mentioned] that the measured atmospheric escape rates for Venus, Earth, and Mars are comparable within the current uncertainties. It has recently been proposed that magnetic fields, rather than shielding a planetary atmosphere from stripping by the solar wind, actually collect solar wind energy and transfer it to the ionosphere along field lines. Global magnetic field lines converge near the cusps, so that the energy is more spatially concentrated than for unmagnetized planets. The escape rate for a given planet may be comparable when it is magnetized, or even greater because planetary magnetic fields extend much further than the planet’s atmosphere, giving it a larger energy-collecting cross-section in the solar wind. One key difference with magnetized planets is that the concentrated energy in cusp regions is likely to lead to more efficient removal of heavy species.

There are a few caveats: [the estimated planetary atmospheric escape rates are quite uncertain, not all solar wind energy collected by a planet need go into removing atmospheric particles, and accelerated ions in Earth’s cusps may not escape the planet. Clearly multiple issues need to be understood] before we can determine whether magnetic fields protect an atmosphere from being lost.”

¹⁷⁴ Activity: Make a table summarizing which atmospheric loss processes work on each of the terrestrial planets. Which two processes are most effective for the present-day Earth based on the description in Sect. 12.4.1?

12.4.2 On time scales of up to several millennia

The Sun's variability has affected Earth's climate and atmospheric composition on astronomical time scales, but a multitude of studies looking for both causes and effects on shorter time scales suggest that the [H-III:12.1] "conclusion at the time of this writing with respect to the importance of low-frequency solar variability in the most recent decades, and perhaps up to centuries, might be 'Perhaps, but probably small'. The main reasons why uncertainties persist regarding this issue include these:

- (i) The ~ 150 -year instrumental record is too short to draw definitive statistical conclusions about the connection of any relation existing on the multi-decadal time scale.
- (ii) Forcing from anthropogenic greenhouse gases represent a significant overprint on trends since about 1850 CE. Because to first order the trends in proxies for solar activity indices and in greenhouse gas concentrations are similar, there is a statistical degeneracy which leads to ambiguous, and thus potentially misleading, conclusions unless great care is taken.
- (iii) A similar problem of statistical degeneracy applies to the Little Ice Age interval of cool conditions during the last millennium (main phase about 1450 – 1850 CE), when mountain glaciers advanced in many regions and planetary temperatures were about 0.5°C lower. During the Little Ice Age, solar activity, as inferred from changes in radiogenic isotopes such as ^{14}C and ^{10}Be , appears to have varied similar to pulses in volcanism and slightly lower carbon dioxide levels. Ignoring this similarity in patterns of variability in internal and external (in planetary terms) climate drivers can lead to erroneous conclusions. {A:^[175]}

A:175

[...] " There are, however, fingerprints of solar variability that locally stand out. For example, [H-III:16.7] "[v]ariations in solar radiation over the 11-yr cycle as well as over the 27-d solar rotation period have substantial effects in the upper atmosphere where energetic photons penetrate and directly initiate photochemical effects. In the stratosphere and the troposphere, above which shortwave radiation is absorbed, the direct impact of solar variability becomes less pronounced. Solar signal in ozone and temperature, however, can be derived from observations above approximately 25 km altitude. Below this height, the situation becomes more complex because other dynamical signals such as those produced by climatic modes of variability (*e.g.*, El Niño) interfere with possible variations resulting from solar variability.

¹⁷⁵ Activity: Human impacts on climate appear not to be limited to the Industrial Revolution! Have a look at a study by Koch *et al.* (2019): they argue that the large population reduction in the Americas following the arrival of European conquerors and settlers, and the resulting reforestation of abandoned agricultural lands, was a significant part of the change in atmospheric CO_2 in the late 16th Century and in the 17th Century.

Several mechanisms have been proposed to explain a plausible relation between solar variability and the observed 11-yr dynamical variability in the lower atmosphere. One of them is associated with disturbances produced in the upper atmosphere and resulting from ozone variations generated by changes in shortwave solar radiation. A second mechanism is linked to the ocean-surface response to 11-yr changes in the total solar irradiance. Observed weather patterns correlated with solar forcing could result from both downward-propagating disturbances produced in the stratosphere and upward-propagating perturbations generated at the surface of the ocean. To capture the amplifying mechanisms producing a dynamical response of the troposphere to solar variability, atmospheric models must therefore account for photochemical processes in the upper atmosphere and, at the same time, must be coupled to an ocean module. Despite many remaining uncertainties, much progress has been made in the last years to better understand how solar variability could potentially affect the climate system, particularly on decadal timescales.”

{A:^[176]}

A:176

¹⁷⁶ Activity: Compile a list of all the processes involved in setting a planetary climate system that reflects at least all those mentioned in Chs. 11 and 12. You can assimilate relevant processes from Activity 164 here as a start.

13

Evolving upper atmospheres and iono-magnetospheres

[H-IV:9] “As one moves up in altitude in a planetary atmosphere, several important changes in composition and structure are apparent. Most notably, as a consequence of hydrostatic equilibrium, the gas density decreases, *i.e.*, the air becomes ‘thinner’. [...] With decreasing density, the frequency of collisions between atmospheric molecules decreases to the point where bulk motions such as turbulence are no longer able to mix the atmosphere. Instead, molecular diffusion becomes the more rapid process and this also leads to a composition change whereby the lighter constituents, typically atomic species such as atomic oxygen, diffuse upward more rapidly than their heavier counterparts such as O₂, N₂ or CO₂. The region where the atmosphere is well mixed is known as the homosphere; the region where diffusive separation dominates is known as the heterosphere. [...]” [H-IV:9.1.1] Because “molecular diffusion coefficients (D) vary inversely as [the square root of the] molecular mass, the molecular diffusion velocities are greater for the lighter constituents and smaller for heavier constituents. Furthermore, they vary inversely as the total density (*i.e.*, diffusion of a gas is more rapid if collisions are less frequent), thus D increases with altitude.” {A:^[177]}

A:177

High in planetary atmospheres is also the region where the temperatures rise (Fig. 2.6) as a result of inefficient cooling while absorbing solar UV to X-rays. This absorption also acts to break chemical bonds and to liberate electrons from their orbits, thus creating the ionospheres. Earth’s upper neutral atmosphere is dominated by N₂ up to about 200 km (see Fig. 2.5), and the overlapping ionosphere is dominated by NO⁺ and O₂⁺ (see Fig. 13.2). Up to roughly 150 km and 200 km, respectively, Venus’ and Mars’ neutral atmospheres are dominated by CO₂ while O₂⁺ dominates in the corresponding ionic components in the

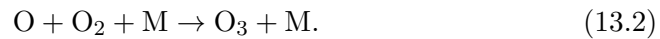
¹⁷⁷ Activity: The approximate scaling of the molecular diffusion coefficient D with molecular mass m and particle density n follows from energy equilibrium of the constituent particles. Formulate D as function of the collisional cross section σ and of temperature and density in the case of self-diffusion, *i.e.*, for molecules diffusing among themselves. For a mixture of components, mutual diffusion needs to be considered.

lower layers of their ionospheres. In the next 150 km above that, atomic oxygen is the dominant species in all three neutral atmospheres, while O^+ dominates for Earth, O_2^+ yields dominance to O^+ after the first 50 km or so for Venus, and O_2^+ dominates for Mars [(see Tables 2.1 and 2.2)]. High in this domain, ions are lifted higher than simple estimates of pressure scale heights might suggest because of the ambipolar effects associated with the free electrons. The compositional differences of the ionospheres are a consequence of the different pathways for photo-dissociation of molecules by solar radiation, which in turn feed a number of chemical reactions in the atmosphere.

An example of photo-dissociation is provided by the photo-dissociation of molecular oxygen (O_2)



which leads to the formation of two oxygen atoms. These atoms may react with molecular oxygen to produce ozone molecules (O_3)



Here, M represents a 'third body' (*e.g.*, N_2 , O_2 , Ar), which removes the thermal energy released by this exothermic reaction." [H-III:16.4] "This photochemical process constitutes the only significant ozone production mechanism above 20 km altitude" in Earth's atmosphere. [H-III:16.3] "In this example, the rate of ozone production [per unit volume] is directly proportional to the rate at which oxygen molecules are photo-dissociated:

$$\Pi(O_3) = 2J_{O_2}[O_2], \quad (13.3)$$

where J_{O_2} represents the photo-dissociation coefficient of O_2 and $[O_2]$ the number density of this molecule. The photo-dissociation frequency depends on the [local intensity of the solar radiation after having traversed the higher absorbing layers ($I(\lambda, z, \chi)$ for wavelength λ , height z , and slant or zenith angle χ)] and on the ability of the molecule to absorb solar photons at particular wavelengths. This last parameter is generally expressed as a wavelength-dependent absorption cross-section $\sigma_X(\lambda)$, which can also vary with temperature. In more general terms, the photo-dissociation frequency of a molecule X is expressed as an integral over all wavelengths that contribute to the decomposition of the molecule. The upper bound of this integral corresponds to the minimum energy required to break the molecular bond. The probability that the absorption of a photon leads to the dissociation of molecule X is expressed by the quantum efficiency η_X , which also varies with wavelength and in some cases with temperature. Thus,

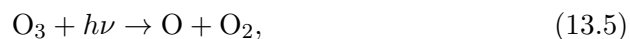
$$J_X = \int \sigma_X(\lambda, T(z)) I(\lambda, z, \chi) \eta_X(\lambda, T(z)) d\lambda. \quad (13.4)$$

A:178

{A:[178]} The solar actinic flux I must be calculated by a radiative transfer model that accounts for (1) absorption processes, (2) multiple scattering by air molecules and atmospheric particles, (3) cloud radiative transfer and (4) surface reflection. When considering upper and middle atmosphere processes, the most important contribution to photo-dissociation is the direct solar flux, so that the value of the actinic flux can be approximated by considering only absorption processes. In the lower atmosphere, multiple scattering and specifically cloud effects cannot be ignored. [...] The depth of penetration of solar radiation varies substantially with wavelength (Figure 2.4). [...]

The relative amplitude of the changes in the solar flux over the 11-yr solar cycle or the 27-d mean synodic solar rotation period decreases with increasing wavelengths (see Figures 2.3 and 12.3) and, as a result, the influence of solar variability is considerably more pronounced in the upper atmosphere [(where EUV and FUV are absorbed)] than in the lower layers [(where longer wavelength radiation is absorbed)]. Strong solar signals associated with the solar cycle are visible in the thermospheric temperature and air density, with impacts, for example, on satellite drag [(and, of course, on ionospheric densities)]. Substantial changes have also been reported in the concentration of nitric oxide (NO); these changes, however, are also related to the modulation of energetic particle precipitation associated with geomagnetic activity. Solar-related changes in the temperature, water vapor and polar mesospheric clouds have also been reported in the mesosphere. In the stratosphere, solar-driven changes in temperature and ozone concentrations have been observed. The influence of solar variability in the troposphere is [touched upon in Ch. 12]. A major forcing function for many of these changes is the variation of photo-dissociation rates. Together with the solar-induced changes in atmospheric heating resulting from the absorption of solar radiation by ozone and molecular oxygen, atmospheric models designed to simulate the response of the atmosphere account for the changes in the photo-dissociation coefficients of the different chemical compounds.”

Because ozone is an efficient absorber of solar UV radiation in the stratosphere it has received much attention. Having been generated by reaction (13.2), ozone is in principle [H-III:16.4] “photo-dissociated



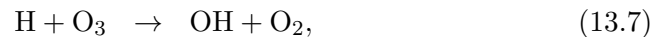
but, in most cases, this reaction does not constitute a net loss for stratospheric ozone because the oxygen atoms that result from this photo-decomposition usually recombine with molecular oxygen (reaction 13.2) to reproduce ozone. The net loss of ozone results from the reaction between oxygen atoms and

¹⁷⁸ Activity: Work through the units of Eqs. (13.3) and (13.4) to show that η_X is an efficiency per unit energy per unit wavelength.

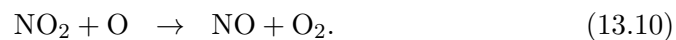
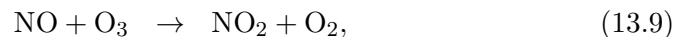
ozone molecules that produce two oxygen molecules



The simple scheme presented here provides a first-order description of the ozone chemistry in the stratosphere and mesosphere. Photochemical models that account only for [the above] reactions tend to substantially overestimate the concentration of ozone in the middle atmosphere, as shown by numerous atmospheric observations. The discrepancy can be eliminated by considering several additional reactions that catalyze (*i.e.*, accelerate) the net loss mechanism represented by reaction (13.6). [T]he presence of the hydrogen atoms and hydroxyl radicals, produced in the upper atmosphere from the photo-dissociation of water vapor (H_2O), could generate an efficient catalytic cycle such as

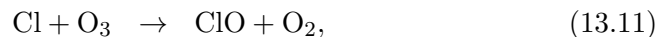


[T]he most effective ozone destruction in the stratosphere results from a catalytic cycle involving nitrogen oxides



NO is produced in the stratosphere by the oxidation of nitrous oxide (N_2O), a long-lived compound released from soils by bacterial activity. It can also be produced in the upper layers of the atmosphere by the dissociation and ionization of molecular nitrogen (N_2) by energetic particles.

Additional destruction mechanisms must be considered including catalytic processes involving halogen compounds including chlorine (Cl) and bromine (Br). For example



Before the 1960s, the contribution of this cycle was relatively small. However, its importance has grown in the last decades as the atmospheric abundance of Cl has increased steadily due to the production of industrially manufactured chlorofluorocarbons (CFCs). The atmospheric lifetime of CFCs varies typically from 50 to 100 years, so that anthropogenic chlorine will remain for several decades in the stratosphere.

In the cold polar regions, and specifically in Antarctica, ozone can be efficiently destroyed in a layer between 12 km and 25 km where polar stratospheric clouds are formed. The solid or liquid tiny particles inside these thin and

often invisible clouds that are present during winter provide surfaces for heterogeneous chemical reactions to operate. Chemical chlorine reservoirs such as HCl and ClONO₂, which are very slow to react in the gas phase, are rapidly converted on the surface of these cloud particles to form less stable molecules such as Cl₂ or HOCl. Large quantities of reactive chlorine atoms (Cl) are liberated via photo-dissociation as soon as the Sun returns in early spring. This chlorine activation leads to rapid ozone destruction with the formation of the springtime Antarctic ozone hole in September and October. These mechanisms are less efficient in the Arctic, where the winter temperature is usually 10° – 15°C higher than at the opposite pole, and the presence of polar stratospheric clouds is therefore less frequent.

A full description of the ozone behavior requires that large-scale transport processes are taken into consideration, specifically in the lower stratosphere, where the photochemical lifetime of this molecule becomes much longer than the time constant associated with transport. Below approximately 25 km altitude, ozone can be regarded as a quasi-inert tracer that is more sensitive to advection and mixing processes than to photochemical transformations. This highlights why the global ozone distribution in the atmosphere is strongly affected by the meridional circulation, and specifically why the ozone column abundance reaches a maximum value at high latitudes at the end of the winter. The poleward meridional circulation transports ozone towards the Arctic where it accumulates from December to April before it is slowly destroyed by photochemical processes after the Sun returns in early spring. [...] The same dynamical process occurs in the southern hemisphere with a lag of 6 months. However, ozone does not easily penetrate poleward of 60°S due to the existence of a strong dynamical barrier provided by the intense southern polar vortex. The ozone maximum is therefore located in a latitude band located at about 60°S. Large-scale planetary waves that characterize the northern hemisphere winter dynamics do not allow the northern hemisphere polar region to be isolated from lower latitudes as is the case in the less dynamically disturbed Southern hemisphere stratosphere. The ozone maximum in the Northern hemisphere is thus located near the Pole.”

13.1 Maintaining ionospheres

13.1.1 Ionization

[H-III:13.2] “The ionosphere is created by ionizing radiation, including extreme ultraviolet (EUV) and X-ray photons from the Sun, and [– in magnetized planets –] corpuscular radiation that is mostly energetic electrons, [which for Earth occurs mostly] at high magnetic latitude as auroral ‘precipitation.’ The

solar photon output at these wavelengths, from $\sim 10 \text{ \AA}$ to the H I Lyman- α line ^[xxvi] at 1216 \AA , varies by factors ranging from ~ 2 to > 100 over the 11-yr solar activity cycle (*cf.*, Fig. 2.3), and is additionally variable on shorter time scales, including especially the 27-d solar rotation period. This causes dramatic variations in the temperature and density of the thermosphere and ionosphere. Changes in the solar wind and in the interplanetary magnetic field also affect the thermosphere/ionosphere through geomagnetic perturbations that result in transfer of energy from the magnetosphere, both in the form of auroral particle ionization and in the form of heat from the resulting currents imposed in the polar regions. An additional form of energy transfer is the generation of energetic electrons released in the ionization process. These electrons, referred to as photo-electrons in the case of photo-ionization and secondary electrons in the case of particle-impact ionization, have enough energy to excite, dissociate, and further ionize the neutral atmosphere as well as heat the ambient plasma. Solar ionization and its byproducts provide most of the ionization and heating of the thermosphere, and account for most of its 11-yr cyclic and 27-d rotational variation, but geomagnetic activity accounts for much of the shorter term variation on time scales from hours to days.

The details of ionospheric formation can be explained through examination of the photo-ionization and photo-absorption cross sections of thermospheric constituents. The ionization continua of N_2 , O, and O_2 all peak in the vicinity of 600 \AA at tens of Megabarns ($1 \text{ Mb} = 10^{-18} \text{ cm}^2$). This causes their energy to be deposited largely in the F_1 region (compare [Figs. 2.4 and 2.7, and see Eq. 2.20 and the text preceding it].) Short-ward of 600 \AA , cross sections decrease and the radiation penetrates to lower altitude. At Earth, the intense solar He II emission at 304 \AA deposits most of its energy near 150 km and $10 - 100 \text{ \AA}$ soft X-rays can penetrate to 100 km. Most of the E region is produced by longer wavelength radiation, particularly the C III line at 977 \AA and the H I Lyman- β line at 1027 \AA . These do not have enough energy to ionize N_2 and O but penetrate through gaps in the N_2 absorption spectrum to ionize O_2 to O_2^+ . Longward of 1030 \AA , only the important minor species NO has a low enough ionization potential to be ionized by solar radiation. H I Lyman- α happens to fall at a low point in the O_2 absorption spectrum and so penetrates below 90 km, where ionization of NO to NO^+ and subsequent products create the D region. Thus, while the Chapman production function [(Eq. 2.20)] is

^{xxvi} Ions are denoted by their electrical charge, such as doubly-ionized C: C^{2+} . The line spectrum of such an ion is identified by a roman numeral that is one higher than the ionization charge, so the spectrum of C^{2+} is written in shorthand as C III; the numeral I is reserved for the spectrum of the neutral species, *e.g.*, CI for neutral carbon. Some spectral sequences have a proper name associated with them: for example, the H I Lyman sequence is a series of spectral lines absorbed or emitted when excited electrons transition from or to the ground state, respectively.

approximately correct for any species at each wavelength, ionized regions are created by the superposition of many such functions [...]"

13.1.2 Recombination

[H-III:13.3] "Positive ions have generally fast collision rates with electrons, so one would suppose that ionospheric production would be balanced by recombination and that the ions would be short-lived after sunset. However, atomic ions colliding with electrons have the problem common to all two-body reactions that a single atom is unlikely to result, because there is nothing to carry away surplus kinetic energy. Photon emission following collision of an atomic ion with an electron can stabilize the resulting atom; this radiative recombination is quite slow, with rate coefficients of the order of $10^{-12} \text{ cm}^3 \text{ s}^{-1}$. Although radiative recombination occurs and is important in the highest reaches of the ionosphere, it is insufficient as a loss mechanism for ions and electrons given their observed *F* region densities. Because the solar ionization frequency is $\sim 10^{-6} \text{ s}^{-1}$ at 1 AU, ion densities would be several orders of magnitude larger than observed if radiative recombination were the only loss mechanism. [What commonly happens is that atomic ions yield their charge to molecular ions in order to undergo rapid dissociative recombination, while in addition there is loss through diffusive transport.] *Dissociative recombination*, schematically $XY^+ + e^- \rightarrow X + Y$, has rate coefficients of the order of $10^{-7} \text{ cm}^3 \text{ s}^{-1}$ and is the fundamental loss mechanism for ions in dense planetary ionospheres. In order for an atomic ion to become a molecular ion, *atom-ion interchange*, schematically $X^+ + YZ \rightarrow XY + Z^+$, or *charge exchange*, schematically $X^+ + YZ \rightarrow X + YZ^+$, must occur. Charge exchange reactions are typically fast if energetically possible, but atom-ion interchange rates depend on the nature of the reacting molecule, because a bond must be broken. {A:^[179]}

A:179

In regions of the atmosphere where molecules dominate, recombination chemistry is simplified because it is essentially a balance between ionization and dissociative recombination. [A] common approximation is the use of an effective recombination rate coefficient α_{eff} , the ion density-weighted average of the ion recombination rates. In photochemical equilibrium, the production rate [per unit volume] $\Pi(e^-) = \alpha_{\text{eff}}[M^+][e^-]$, where M^+ is the sum of the ions, and where square brackets denote number densities. Assuming charge neutrality,

¹⁷⁹ Activity: Consider the similarities and differences between the charge-exchange reactions described here and two- and three-body gravitational interactions, specifically what is needed for the capture of interplanetary spacecraft into closed orbits, or the capture of planetary bodies as moons of planets. For the latter, look up the concepts proposed for the capture of Triton, the largest moon of Neptune, orbiting that planet in a retrograde orbit (which implies it has to involve a capture well after the formation of the planet).

this yields

$$[e^-] = (\Pi(e^-)/\alpha_{\text{eff}})^{1/2}. \quad (13.13)$$

{A:[180]} Applying the Chapman production function for solar radiation A:180 [(Eq. 2.20)] to obtain Π results in a Chapman 'layer', considering as above the caveats associated with use of that term. Thus, in molecular ionospheres, electron density varies approximately as the square root of the ionization rate profile. Eq. (13.13) is a particularly useful form for auroral ionization, where electrons (and sometimes protons or heavier ions) penetrate to ~ 100 km or deeper into Earth's atmosphere. {A:[181]} A:181

Although the F_2 region has some of the morphological appearance of this type of layer, it is at the wrong altitude, and in the atom-dominated region. It is not a Chapman layer at all, but a result of diffusive processes. O^+ has an increasingly long lifetime as altitude increases and the molecular fraction of the thermosphere decreases. Above 200 km it becomes subject to diffusion, but is still chemically controlled up to the peak of the F_2 layer near 300 km. Above this altitude, ambipolar diffusion takes over, where 'ambipolar' refers to the effect of electrical attraction between the ions and nearly massless electrons, resulting in a scale height for O^+ about twice that of O. {A:[182]} The A:182 F_2 region varies in response to thermospheric winds and electric fields, so the mid-latitude and equatorial ionosphere can be greatly influenced by auroral processes at high latitudes through their effect on thermospheric dynamics.

Figure 13.1(left) is provided as a guide to understanding the ion-neutral

¹⁸⁰ Activity: Note the equivalence between Eq. (13.13) and Eq. (11.5) for a volumetric ionization rate of $\Pi(e^-) \propto \Phi_i/(4\pi R^3)$. This means that α_{eff} is, in effect, for a 'case B' recombination, *i.e.*, excluding the possibility that emitted photons in recombination are absorbed to lead to another ionization event. Consider what could happen to avoid that. Also see a parallel with the formulation of what can be viewed as the inverse in Eq. (9.3): for a stationary, isothermal case, the 'incoming' volumetric heating ϵ_{heat} balances the outgoing radiation $n_e n_H f_{\text{rad}}$ in which the product of ion and electron densities is a measure for the number of collisions leading to excitation, to compare with the ionizing radiation in the ionosphere which balances the recombination in which the product of ion and electron densities is a measure for the number of collisions leading to recombination.

¹⁸¹ Activity: One might think that collisions between particles that can 'bond' and thereby be taken out of a population under study, such as electrons and positively-charged ions that combine into a neutral particle, might have a good analogy in how flux concentrations in the solar photosphere behave: the concentrations perform a random walk and in collisions opposite magnetic polarities 'cancel', *i.e.*, disappear from the population of magnetic charges. Yet the scaling behavior between the strength of the source (the total of emerging bipoles per unit time) and sinks (the total of canceling flux per unit time) is different: the square root dependence reflected in Eq. (13.13) does not show up, but instead a near-linear dependence appears (as shown here by Schrijver (2001)). Consider the reasons: when the Sun's activity increases, flux concentrations grow larger by collision thereby countering the increase in collision frequency expected; larger concentrations are less mobile within the evolving convective motions; fragmentation and coagulation are seeking a balance; while in general the large-scale meridional flow aids in separating polarities (a process that is countered in an ionosphere by the tendency towards charge neutrality).

¹⁸² Activity: Show for the simplified case of a fully-ionized static gas that the scale height for ions is twice that for the corresponding atoms in a neutral atmosphere by combining the momentum equations for ions and electrons in comparison to that equation for a neutral species. And remind yourself how this is consistently incorporated in the MHD equations.

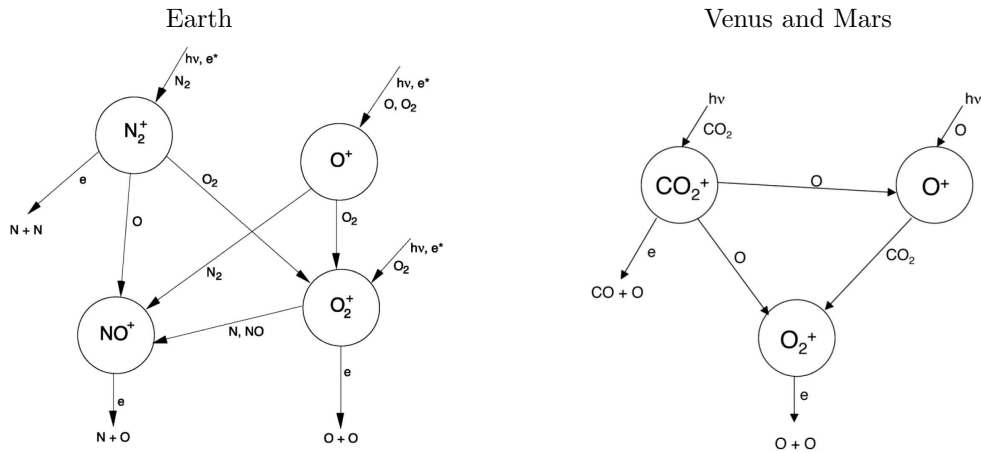


Fig. 13.1. Simplified diagram of ionospheric chemistry in the upper atmosphere of Earth (left) and of Venus and Mars (right). [The relative ionization potential of these species is roughly indicated by the position in this diagram: higher position indicates higher ionization potential. Fig. H-III:13.4]

chemical processes described above. It is a greatly simplified schematic, but contains the essential species and reactions necessary to describe ionospheric photochemistry from 100 – 600 km. Ionization occurs primarily on the three major species N_2 , O , and O_2 by photon, photo-electron, and auroral particle impact. N_2^+ quickly loses its charge through dissociative recombination, atom-ion interchange with O to make NO^+ , and, at lower altitude, charge exchange with O_2 to make O_2^+ . Thus it is always low in density and negligible in the absence of production. O^+ loses its charge by molecule-ion interchange with N_2 and O_2 . Reaction with N_2 is slow, $\sim 10^{-12} \text{ cm}^3 \text{ s}^{-1}$, because of the high strength of the triple $N \equiv N$ bond. Reaction with O_2 is faster, $\sim 10^{-11} \text{ cm}^3 \text{ s}^{-1}$, but there is far less O_2 available. This is why O^+ is long-lived in the Earth's ionosphere, and why the F_2 region exists. O_2^+ loses its charge through dissociative recombination or through reaction with the odd-nitrogen species NO and N , which control the balance between O_2^+ and NO^+ in the E region. NO^+ , daughter of all the above and the 'terminal ion', is subject only to dissociative recombination. Figure 13.1(left) thus describes a mechanism for dissociating molecular gases. Ionization goes in the top [involving high energies], and dissociation comes out the bottom, because dissociative recombination is the only significant way out [of the cascade that is accessible given the energies involved].

Figure 13.1(left) [neglects the effects of hydrogen at high altitudes, where that] reacts by charge exchange with O^+ to make H^+ . N^+ is also a significant minor ion, created by photo-dissociative ionization of N_2 , that is neglected

here. Doubly-ionized species are also ignored in this simplification. Although ground state $O^+(4s)$ does not have enough energy to make N_2^+ , metastable $O^+(2d)$ and $O^+(2p)$ are created by photon and electron impact ionization, and these can charge exchange with N_2 to form N_2^+ . It is possible that vibrational excitation of thermospheric N_2 can also accelerate the reaction of O^+ with N_2 .

In Earth's E region, there is a complex interplay between O_2^+ and NO^+ , due to the involvement of odd-nitrogen species with the ion chemistry, because O_2^+ is converted to NO^+ by reaction with NO and N . NO in particular is highly variable with solar activity, geomagnetic activity, and location, so this is a complicated problem. Older empirical and theoretical models which assumed that O_2^+ is the dominant E region ion, due to its production by solar H I Lyman- β radiation, have been superseded by evidence that NO^+ is generally observed to be the dominant E region ion, and considerable recent observational and modeling advances in understanding the high levels of NO and its importance to radiative cooling as well as ion chemistry have occurred.

In the D region, ion chemistry is entirely different due to the higher neutral density which allows three-body attachment, particularly $2O_2 + e^- \rightarrow O_2 + O_2^-$. This sets in motion a complicated negative-ion chain involving carbon, nitrogen, and hydrogen compounds, including water, that finally results in mutual neutralization of negative and positive ions, schematically $M^+ + M^- \rightarrow M + M$. NO^+ created by H I Lyman- α ionization also initiates an involved positive-ion sequence, again involving hydration processes. [...]"

13.1.3 Venus and Mars

[H-III:13.4] "The terrestrial planets Venus, Earth, and Mars, are so named because of their fundamental similarity, and are presumed to have had common elemental origins. However, their subsequent evolution differed, due to their differing distance from the Sun, the smaller size of Mars, and the lack of rotation of Venus ([see Ch. 12]). Thus, their atmospheres are entirely different, and so are their upper atmospheres and ionospheres. Early exploration of Venus and Mars found that instead of persistent, high-altitude, F_2 -type ionospheres, these planets had less dense, lower-altitude ionospheres ([Fig. 2.5]) that more resembled Chapman 'layers', that were greatly attenuated at night, and consisted mostly of O_2^+ and other molecular ions. The presence of O_2^+ seems especially perplexing, because Earth is the planet we generally associate with the unusual and quite reactive oxygen molecule. At higher altitude, O^+ becomes an important species in the ionospheres of Venus and Mars, as on Earth, but at significantly lower density and without the same degree of persistence throughout the night. CO_2^+ is a minor ion on both planets. There

is a basic similarity in their ionospheres, despite the vastly different density of their lower atmospheres. [...]

The reason that the ionospheres of Venus and Mars are different from that of Earth is that the molecular compositions of their atmospheres are different, and therefore the compositions of their thermospheres are different ([Fig. 2.5]). Table [2.1] gives a simple overview of the abundance of the primary atmospheric gases in the three terrestrial planets.

Aside from the large differences in surface pressure, the atmospheres of Venus and Mars are similar in composition, and N_2 is an important species on all three planets. N_2 requires more energy to dissociate than the oxygen compounds, however, so at thermospheric altitudes, atomic oxygen becomes important on all three planets. The Venus and Mars thermospheres are distinguished by high levels of CO_2 (and also CO) due to the underlying atmospheric composition, as shown in [Fig. 2.5].

On Earth, O^+ is a long-lived species in the high ionosphere because the $O^+ + N_2$ reaction is so slow and there are few other molecules to react with to make a short-lived molecular ion. On Venus and Mars, the reaction $O^+ + CO_2$ is quite fast, $\sim 10^{-9} \text{ cm}^3 \text{ s}^{-1}$, because CO_2 is much less strongly bound than N_2 . (The triple bond in $N \equiv N$ is [among the strongest known chemical bonds in nature, along with the triple bond in CO].) The reaction of $O^+ + CO_2$ produces O_2^+ , which is also produced by the reaction $CO_2^+ + O$. Ionization of the major thermospheric gases on Venus and Mars, O and CO_2 , is thus quickly converted into O_2^+ , which dissociatively recombines, resulting in the observed ionospheric morphology, lacking a significant F_2 region. A simplified schematic of these processes is shown in Figure 13.1(right). Thus, curiously, although life-supporting Earth is the planet associated with O_2 , Venus and Mars are the planets with O_2^+ ionospheres.

The F_2 ionosphere is unique to Earth among the known planets. This is due to its peculiar atmosphere, lacking in CO_2 , dominated by N_2 , and carrying its oxygen in unusual and reactive states. Venus and Mars have nitrogen as well, but carbon and oxygen dominate their upper atmospheres, so it has little effect. Earth has a significant carbon budget, and once had much higher levels of CO_2 in its atmosphere, but most of its carbon is currently locked up in the crust in the form of carbonate rocks. Thus, the F_2 ionosphere may be a recent event in the history of Earth, an artifact of geology and biology.” {A:^[183]}

A:183

¹⁸³ Activity: Review Figs. 2.5, 13.1, and 13.2 and think through the dominant reactions described in Sects. 13.1.2 and 13.1.3.

13.2 Setting geospace climate

13.2.1 Geospace climate response to solar photon irradiation

[H-III:14.3] “The solar spectrum provides a [relatively] stable irradiance of $\sim 1360 \text{ W/m}^2$ ($\sim 1.36 \cdot 10^6 \text{ erg/s/cm}^2$) to the Earth’s upper atmosphere. Geospace [(the region of space near Earth down to, and including, the ionosphere/thermosphere))] responds to [roughly 2–6 ppm of that] this fraction of the solar spectrum lies between 30 and 3600 Å which extends from the X-ray through the ultraviolet part of the spectrum. The photons in this spectral range may ionize atoms and molecules or may deposit their energy directly into the thermal reservoir of the upper atmosphere. These processes are responsible for the ‘climate’ of the geospace-atmosphere interface whose regions are labeled the ionosphere and the thermosphere (IT). The IT in this sense is only weakly dependent on the Earth’s magnetic field or the solar wind. In this ‘climate’ scenario, the role of the terrestrial dipole field can be viewed as defining the boundary for the plasmasphere [(the inner magnetosphere filled with low-energy, cool plasma)] and then with the solar wind the magnetosphere. In the case of the Earth’s sister planets, Venus and Mars, the absence of a significant intrinsic magnetic field confirms that the IT development has been based on these three photochemical processes.

Now let us explore the question how geospace climate responds to extremes of the solar photon radiation. Figures 13.2 and 13.3 provide a comparison of the [1-D Global Averaged Ionosphere and Thermosphere (GAIT)] solar cycle climate of the thermosphere’s neutral densities, of the ionosphere’s plasma densities, and of the neutral and plasma temperatures respectively. Each panel is shown as a function of pressure level defined by $Z = \log p_0/p$, where p_0 is the reference pressure of $0.5 \cdot 10^{-3} \text{ dyn/cm}^2$. The corresponding altitude scale is also provided. For reference to observations, the dashed lines where present in Figure 13.2 and 13.3 correspond to profiles obtained from the MSIS-90 empirical model of the thermosphere.

A simple interpretation of this solar minimum to solar maximum climate change in the IT is that the effective IT energy deposition has almost quadrupled; hence, the neutral atmosphere, which at these heights is in hydrostatic equilibrium, leads to a hotter thermosphere; compare T_n in the two panels of Figure 13.3. In turn, the hotter thermosphere has redistributed neutrals now with relatively higher densities at higher altitudes; compare neutral densities in the two left-hand panels of Figure 13.2 using the right side altitude scale. [...] A secondary but also important additional effect is that the composition is also being modified because of the different neutral masses. For the ionosphere, the consequences can readily be seen by comparing the two right-hand panels of Figure 13.2. [...] A comparison of the *E* and *F* layer peak density provides

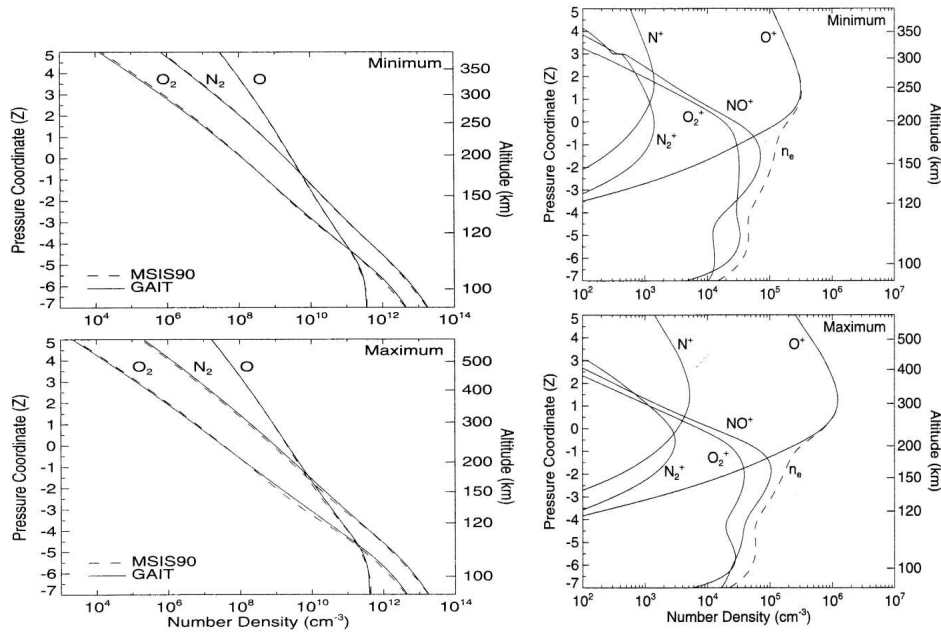


Fig. 13.2. Left: Global mean number density profiles under the same conditions for the three major neutral species in Earth's upper atmosphere (N_2 , O_2 , and O) calculated using the GAIT model (solid lines) and the MSIS-90 empirical model (dashed lines). Right: Global mean number density profiles for five ion species in Earth's ionosphere (O^+ , NO^+ , O_2^+ , N^+ , and N_2^+) and the total electron density (n_e) calculated using the [1-dimensional Global Averaged Ionosphere and Thermosphere (GAIT)] model. Solar minimum (top) and solar maximum (bottom; assuming quiet geomagnetic conditions with $A_p = 4$). The discontinuity observed in the NO^+ and N_2^+ profiles at $Z = 3$ corresponds to where the photo-electron calculation stops. [Fig. H-III:14.1; source: Smithtro and Sojka (2005).]

another useful scaling law, rule-of-thumb, in that the F -layer density scales linearly with the appropriate photon wavelength energy flux while that of the E layer is more like a square root dependence on energy flux [(compare with Eq. 13.13)].

That the different ionospheric layers respond differently to the solar spectrum creates the problem of deciding what the most suitable solar spectral representation is. In fact, even over the limited solar cycle energy flux range of a factor of about 4, the spectrum itself is variable and the E and F layers respond to different parts of the spectrum. {A:^[184]} The thermosphere is a somewhat better integrator as seen [from a] study in which four distinctly

A:184

¹⁸⁴ Activity: Trace which part of the solar spectrum provides the predominant power to the E and F layers of the terrestrial ionosphere and overlapping thermospheric regions, and note that the power going into the F layer exhibits a larger variation over the solar cycle than that going into the E region. See Sect. 13.1.1 and Figs. 2.3, 2.4 and 2.7.

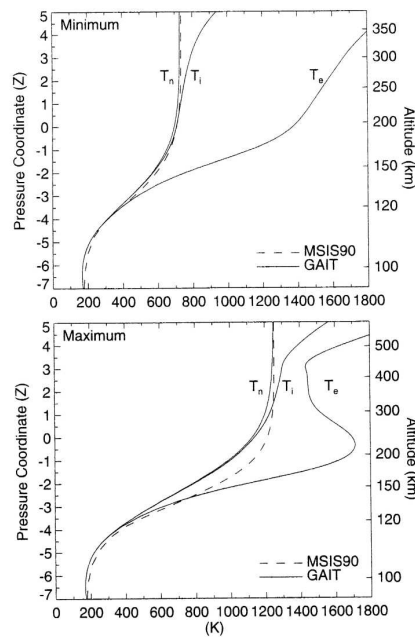


Fig. 13.3. Global mean temperature profiles of Earth's ionosphere/thermosphere calculated using the GAIT model [(solid lines). The three profiles correspond to neutral (T_n), ion (T_i), and electron (T_e) gases. The dashed lines show T_n for the MSIS-90 empirical model.] Solar minimum (top) and solar maximum (bottom) assuming quiet geomagnetic conditions ($A_p = 4$). [Fig. H-III:14.2; source: Smithtro and Sojka (2005).]

different representations of the solar spectrum were used as drivers for the GAIT model. As each spectral model was run over the solar cycle range of 2–8 erg/s/cm² [going into the ionospheric/thermospheric height range,] the GAIT exospheric temperature was determined. The results are that the GAIT-model thermosphere responded linearly to each spectral model, and the same linear dependence is found for each. Note that the exospheric neutral temperature refers to the asymptotic, altitude-independent, temperature found at higher altitudes, see Figure 13.3 for specific solar minimum and maximum examples. The exosphere refers to the ionospheric plasma whose composition is light ions of hydrogen and helium that is located in altitude above the F layer.

[By combining various irradiance models computed for solar activity levels throughout the sunspot cycle] it is possible to extrapolate how the thermospheric exospheric temperature would trend for lower and higher levels of the solar EUV flux. [... The] procedure assumed that a linear dependence existed in the relevant EUV energy flux between solar minimum and solar maximum. An index S_{EUV} is defined to be 0 at solar minimum (energy flux

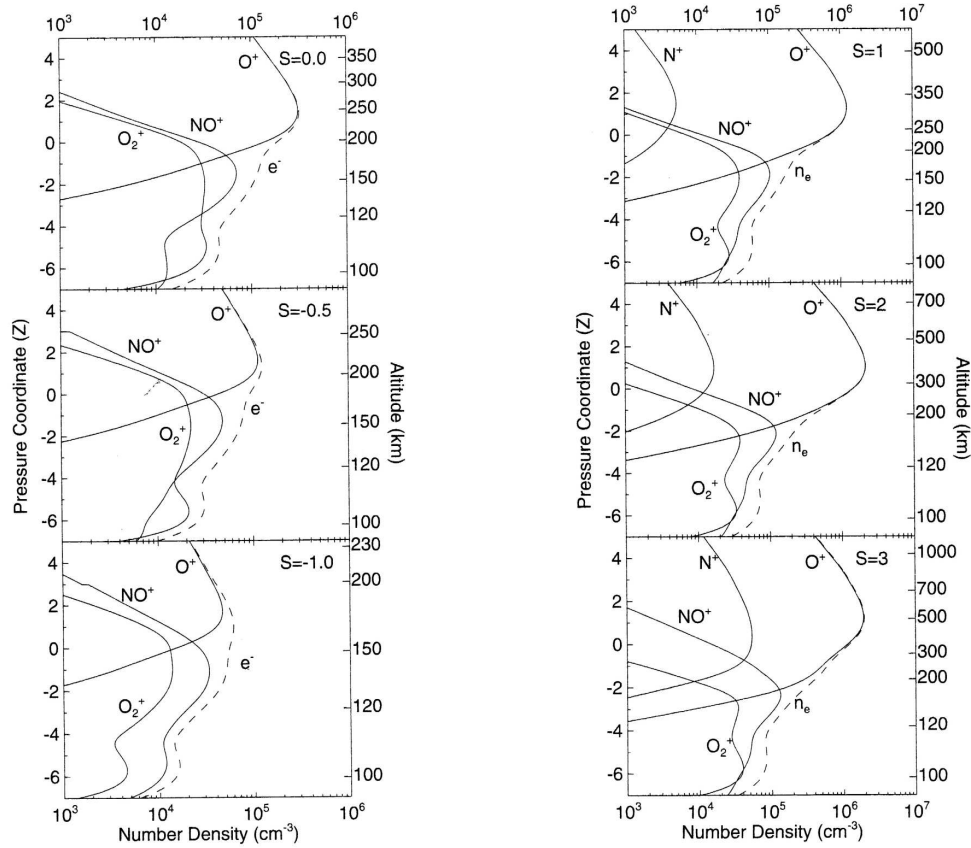


Fig. 13.4. Global mean concentration of the ion and electron (dashed line) gases in Earth's upper atmosphere, calculated using the GAIT model for six different levels of the solar activity increasing from $S_{\text{EUV}} = -1.0$ to $S_{\text{EUV}} = 3$, clockwise from the bottom left. The profiles are plotted as a function of the pressure coordinate, $Z = \log p_0/p$, with the corresponding altitudes provided on the right-hand axis. [Fig. H-III:14.4; source: Smithtro and Sojka (2005).]

of 3 erg/s/cm^2) and $S_{\text{EUV}} = 1$ at solar maximum (energy flux of 7 erg/s/cm^2). Then Maunder-Minimum type conditions correspond to $S_{\text{EUV}} < 0$ and grand maximum values correspond to $S_{\text{EUV}} > 1$. Note that the specific response to the solar cycle of each wavelength is different, hence S_{EUV} is applied to each wavelength separately to generate extreme solar spectra." [H-III:14.3.1] "[T]he Maunder Minimum S_{EUV} value would be between $S_{\text{EUV}} = -0.5$ and -1.0 . Figure 13.4 shows the GAIT ionospheric plasma composition for solar minimum ($S_{\text{EUV}} = 0$), $S_{\text{EUV}} = -0.5$, and $S_{\text{EUV}} = -1.0$. The earlier trends concerning the E and F layer are continued as the S_{EUV} value [is lowered. . .] The most significant ionospheric modification during the Maunder-Minimum

period is that the molecular ion NO^+ peak below 200 km becomes significant; compare the top-left panel for $S_{\text{EUV}} = 0$ and the bottom-left panel $S_{\text{EUV}} = -1$ in Figure 13.4. This additional structure in the electron density profile is referred to as the F_1 layer. Indeed, in Figure 13.4 the $S_{\text{EUV}} = -1.0$ case almost has this F_1 electron density equal to that of the higher altitude F_2 peak. [...] These Maunder-Minimum scenarios provide significant problems for modern-day technology. For example: (1) using the ionosphere to propagate radio waves over the horizon is restricted to much lower frequencies because the maximum ionospheric density has decreased; and (2) because the F_1 layer is located significantly lower than the F_2 layer propagation, paths for radio waves are also modified significantly. [On the other hand, (3)] with less ionospheric density in the path of GPS radio waves, the adverse role of the ionosphere in geolocation analysis is reduced.” {A:^[185]}

A:185

Now for a much more active Sun: [H-III:14.3.2] “[a value of $S_{\text{EUV}} = 3$ can be used to characterize] the upper range of enhanced solar EUV flux to simulate grand-maximum type conditions. The grand maximum existed between 1100 and 1250 CE. An S_{EUV} value of 3 corresponds to doubling the solar maximum solar energy flux from 7 erg/s/cm^2 to just over 14 erg/s/cm^2 . The right-hand column in Figure 13.4 shows the GAIT-model ionospheric plasma distributions at solar maximum ($S_{\text{EUV}} = 1$), $S_{\text{EUV}} = 2$, and $S_{\text{EUV}} = 3$ from top panel to bottom panel. In all cases, the F_2 layer is the dominant layer with O^+ the dominant ion. As predicted from the normal solar cycle trend, this layer will rise, in this case from 300 km ($S_{\text{EUV}} = 1$) to about 500 km ($S_{\text{EUV}} = 3$). The F_2 peak density does not increase linearly with S_{EUV} ! Between $S_{\text{EUV}} = 2$ and $S_{\text{EUV}} = 3$ the F_2 peak density has remained at $2 \times 10^6 \text{ cm}^{-3}$.

This maximum in the F_2 peak density is by far the most significant change in the geospace climate in response to solar photon radiation. The processes responsible for this effect are: (i) the production of neutral O and, hence, its concentration has non-linearly decreased at altitudes at which the F_2 peak is created now as the thermosphere heats up as S_{EUV} increasing from 2 to 3; (ii) the O^+ production rate does increase linearly as S_{EUV} increases from 2 to 3; and (iii) the competition between these two processes leads to a maximum peak F_2 density at $S_{\text{EUV}} = 2$, and then as S_{EUV} increases further even a slight decrease in the peak density. The consequences for modern-day technologies under enhanced solar maximum, grand maximum conditions are: (1) the changing altitude of the F_2 layer leads to modified radio wave propagation paths; (2) that

¹⁸⁵ Activity: Upward traveling radio waves with frequencies below the plasma frequency are evanescent within the ionosphere and are reflected downward, thus enabling ‘over the horizon’ or ‘skywave’ communication. Look up the plasma frequency (Eq. 3.41), typical ionospheric electron densities within the ionosphere, and resulting values for the radio frequencies useful for such communication. See, *e.g.*, Fig. 6.1 for an overview of the EM spectrum with an indication of various radio bands (including what differentiates propagation of AM and FM bands).

the peak F_2 density saturates only slightly above solar maximum values implies that the 'radio' reflection characteristics of the ionosphere are consistent with today's 'radio climate'; (3) the impact on trans-ionospheric radio applications such as GPS geolocation is somewhat adverse since the total electron content (the electron density integrated over a line of sight, *i.e.*, a column density) continues to increase even though the F_2 peak density becomes constant; and (4) because the ionosphere is significantly more dense, the absolute magnitude of plasma density irregularities would increase which would lead to greater scintillation problems with radio propagation."

[H-III:14.3.3] "In modeling the ionosphere and thermosphere as the solar EUV energy flux is changed, there are at least two impacts of significance for the outer reaches of geospace. First, assuming that the magnetosphere is somewhat similar to the state that we are familiar with, then the IT contributes plasma to the magnetosphere/plasmasphere and second, the IT electrical conductivity is a component of the magnetosphere-ionosphere (M-I) electrical coupling. Under the Maunder-Minimum type conditions, ionospheric outflow of plasma into the magnetosphere/plasmasphere will decrease because the ionospheric topside is colder and less dense. Under extreme conditions such as $S_{\text{EUV}} = -1$, the composition may also begin to change from atomic to molecular. In contrast, under $S_{\text{EUV}} = 2$ and upward, during grand-maximum conditions with hotter topside, the outflow would increase and be very much O^+ dominated. Note that in these GAIT-type modeling studies the light ion, H^+ , has not been included, and therefore the remarks pertain to O^+ and heavier molecular ions. In contrast, the ionospheric conductivity changes are smaller because the major contribution comes from the E layer whose composition remains molecular. However, the decreasing dayside conductivity during Maunder-Minimum conditions would raise issues about how this impacts the M-I electric circuit response, *i.e.*, would this modify present-day concepts of voltage versus current generator descriptions of the M-I system? Under the grand maximum with enhanced conductivities and also the assumption of increased solar wind energy, would M-I coupling be characterized by significantly enhanced currents and electric field? Both scenarios would probably impact the morphology of auroral displays! This may lead to the most significant human experience of the geospace climate."

13.2.2 Geospace climate at earlier terrestrial ages

Now let us look at the far larger range of solar activity as that evolved over the 4.6 Gyr history of Venus, Earth and Mars. [H-III:14.4] "In earlier times, the solar EUV was more intense and, consequently, the thermosphere was much hotter, leading to the dominance of significantly different processes. [In

Sun-like stars of different ages			
Name	Spectral type	P_{rot} (days)	Age (Gyr)
EK Dra	G1.5 V	2.7	0.1*
π^1 UMa	G1.5 V	4.9	0.3
χ^1 Ori	G1 V	5.2	0.3
κ^1 Cet	G5 V	9.2	0.7
β Com	G0 V	12.	1.6
Sun	G2 V	25.4	4.6
β Hyi	G2 IV	~ 28	6.7

* [Another study reports] an age of 0.03 – 0.05 Gyr.

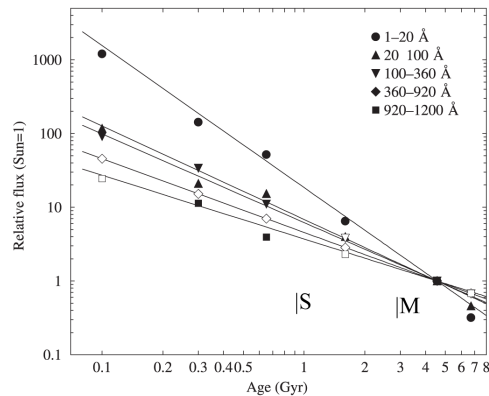


Fig. 13.5. Spectral radiance versus age of solar-type stars (identified in the table on the left, with spectral type, rotation period, and estimated age), in solar units. Measurements are shown by filled symbols; missing data (open symbols) are derived from power-law fits (solid lines) for passbands from 1 to 1200Å. The approximate ages for which the oldest fossils of single-cell microbial life (S) and multi-cellular plants and animals (M) have been found on Earth are indicated. [Fig. H-III:2.14; source: Ribas et al. (2005).]

particular, we look at] an early period when atomic hydrogen was in a blow-off phase as well as periods when high escape rates for the fastest particles in the energy distribution ('Jeans escape') of heavier species like H_2 , He, C, N, O existed [(compare Ch. 12). S]tudies also show that IR-radiating molecules like CO_2 , NO, OH, etc., control the exospheric temperature that, in turn, controls the Jeans escape rates for the neutral constituents. Hence, the results depend not only on a knowledge of solar EUV but also on the contribution of molecules such as CO_2 and H_2O in the earlier terrestrial atmosphere [that can be addressed, for example, with] a diffusive-gravitational equilibrium and thermal balance model to study heating of the earlier thermosphere.

In an initial simulation, this model was used to evaluate the terrestrial exospheric temperature over the past 4.6 Gyr. Significant assumptions were made that the present-day composition as well as that of the lower atmosphere up to 90 km were the same then as they are today. The increased solar flux values at earlier ages [were estimated as] summarized in [Fig. 13.5].

Figure 13.6 shows [the simulated] history of the Earth's exospheric temperature. Assuming that the blow-off temperature for atomic hydrogen is about 5000 K, the Earth's first Gyr would exhibit a markedly different upper atmosphere where even the atomic atoms and molecular hydrogen would be

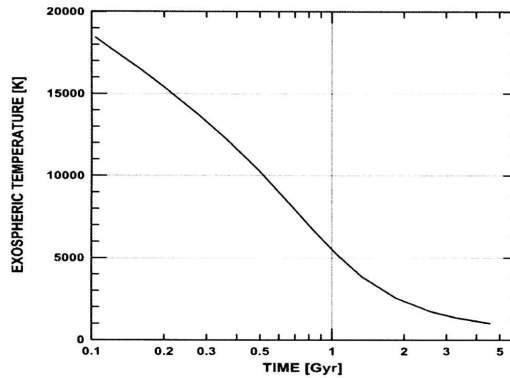


Fig. 13.6. Evolution of the exospheric temperature, assuming Earth's present atmospheric composition, over the planet's history as a function of the solar XUV flux for a strongly limited hydrogen blow-off rate. [Fig. H-III:14.5; source: Kulikov et al. (2007).]

Historical values of the solar EUV fluxes relative to the present-day value.
[Table H-III:14.1]

Time	Solar flux multiplier
3.5 Gyr ago	factor ~ 6
3.8 Gyr ago	factor ~ 10
4.33 Gyr ago	factor ~ 50
4.5 Gyr ago	factor ~ 100

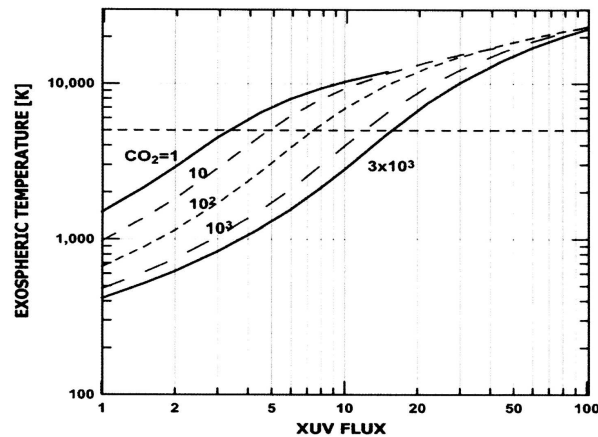


Fig. 13.7. Earth's exospheric temperatures for different levels of CO_2 abundance in units of PAL (Present Atmospheric Level: 1 PAL for $\text{CO}_2 = 330$ ppm) in the thermosphere as a function of solar XUV flux. The numbers by the curves correspond to CO_2 volume mixing ratios expressed in PAL. The horizontal dashed line shows the blow-off temperature of atomic hydrogen. [The Table lists estimated XUV flux levels relative to the present day. Fig. H-III:14.6; source: Kulikov et al. (2007).]

approaching their thermal escape speeds. [T]his simple model becomes a rough estimate when in the life of the Sun (smoothing over time scales long compared to solar cycles) the exospheric temperatures exceed 10,000 K.

The major assumption that would be questioned for these earlier Earth ages would be the density of the IR radiating molecules such as CO_2 . In significantly earlier times, these would be expected to be larger. If this is the case, then

their role in 'cooling' the thermosphere would increase. Defining the CO₂ mixing ratio relative to present atmospheric level (PAL) as 1, [...] Figure 13.7 shows how, indeed, significant increases in CO₂ will cool the upper atmosphere. In this figure, the 'XUV flux' is the scaling ratio of earlier age solar EUV compared to today. The current situation is shown at unit XUV flux. This shows that increasing CO₂ by a factor of 10 (10 PAL) leads to a drop of almost 600 K in the exospheric temperature from 1600 K to 1000 K. [...] Further work has shown that] by solar fluxes that are about 5 times the present average EUV energy flux of 5.1 erg/s/cm², the composition of the upper thermosphere will be dramatically different from today as the Jeans escape mechanism becomes effective for hydrogen as well as other atomic species. [Model studies (see H-III:12.3.3 for some details and references)] indicate that at earlier ages of the Earth's upper atmosphere-ionosphere, the response to increased solar EUV flux was a heating of this part of geospace with the following impacts on the geospace climate: (1) for exospheric temperatures of 5000 K and above, the upper atmospheric composition would be dramatically different due to Jeans escape fluxes of hydrogen and other atomic species; (2) the altitude of the F₂ layer would increase to heights above 2000 km; (3) the F₂ layer peak density would become constant; (4) the total electron content of the ionosphere would increase linearly with increasing solar EUV flux.

The geospace climate would change from an ionospheric standpoint when the solar energy flux slightly exceeds levels of the present-day solar maximum. From the thermospheric point of view, the geospace climate would change when the solar flux EUV reaches about 20 erg/s/cm² (three times the present-day solar maximum values). [...]” {A:^[186]}

A:186

13.2.3 Geospace climate and Earth's magnetic field

[H-III:14.5] “[T]he outer boundary of geospace is defined to be the magnetosphere, and specifically, the magnetopause. It is created by the solar wind that interacts with an intrinsic property of the Earth, the magnetic field. Consequently, in this section questions concerning how long-term trends of the solar wind and Earth's magnetic field are considered in discussing the long-term geospace climate. Of specific interest are the conditions under which the geospace would be dramatically changed. [...]”

[H-III:14.5.1] “Perhaps the geospace response to flips in the Earth's dominant dipolar field is the most frequently discussed geospace 'what if' scenario. Geological evidence obtained in the last century has clearly proven that the

¹⁸⁶ Activity: Estimate when the solar EUV flux dropped to a level that the thermospheric climate became comparable to the present-day state; and summarize the ionospheric changes over geological time scales as far as the models discussed here is concerned.

Earth's magnetic field, especially its dominant dipole component, has reversed many times during geological times. The most recent reversal occurred 0.78 Myr ago. Prior to this reversal the 6 most recent occurred at 0.99, 1.07, 1.19, 1.2, 1.77, and 1.95 Myr ago. [Reversals] occur at quite irregular intervals with the shortest time between reversals being at the Cobb Mountain reversal pair separated by only about 10,000 years. [... R]eversals occur when the dipole field strength is relatively weak.

[...] The specific 'N-S' or 'S-N' dipole orientation itself would not introduce significant geospace climate changes. Perhaps, the most obvious would be that the solar wind northward versus southward reconnection morphology would be reversed. What is significantly more important would be the magnitude of the Earth's field and the orientation of its dipole component."

[H-III:14.5.2] "Over the past 100 years, the Earth's dipole moment has decreased by about 5% from 8.3×10^{25} to 7.8×10^{25} erg/G, while three-thousand years ago, it was at almost 12×10^{25} erg/G, at its highest value during the Holocene era. From [Ch. 5 it is clear that] it is the balance between the Earth's magnetic field and the solar wind pressure that determines the outer boundaries of the magnetosphere/geospace. Hence, a larger (or smaller) dipole moment with otherwise the same solar wind conditions would increase (or reduce) the size of geospace. In turn, this would reduce (or increase) the size of the polar cap, and auroral regions would move poleward (or equatorward). However, a 5% change in the [virtual axial dipole moment (VADM)] would probably not have a dominant impact on geospace because the solar wind pressure varies by more than this over its normal solar cycle. Considering earlier times when the VADM did decrease to values as low as, if not lower than, [a quarter of the present-day value,] the geospace climate may well have been dramatically different, especially during solar maximum type conditions. The magnetosphere would have been severely reduced, and in volume regions such as the plasmasphere it would have almost been reduced to ionospheric altitudes and in the 'open' polar regions would extend to mid-latitudes. The effectiveness of plasma sheet energization processes would also have been changed, causing impacts on ring currents, electrojets, as well as the visible aurora. Perhaps, the energy transfer to geospace would simply decrease as the magnetosphere's cross section to the solar wind decreases, and consequently, all internal energy processes would be similarly scaled down.

The extreme scenario of the dipole reversal is the idea that the VADM for a time period is extremely small, approximately zero. If the higher-order multipole terms are also negligible, then the Earth's atmosphere is unprotected. But this is the Venus and Mars type scenario and extensive analysis has been done on these planetary atmospheres. At present, the scientific techniques that

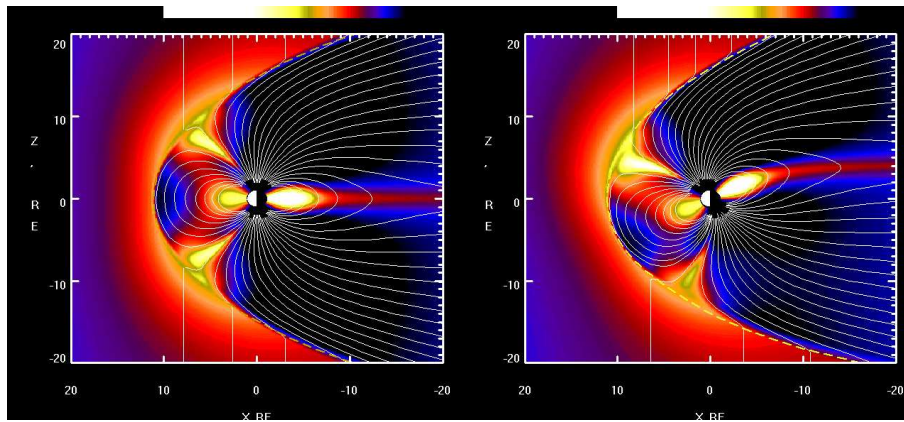


Fig. 13.8. Field lines plotted in the noon-midnight meridian plane for an untitled planetary dipole embedded in the solar wind (left) and a dipole tilted in the noon-midnight meridian by 30° (right). The colors show the difference in the magnitudes of the field in the $x - z$ plane of a model field including the effects of a solar wind compared to that of a dipole field; positive differences are shown in blue to black, negative differences in red through yellow to white. [source; see also Fig. H-III:14.13]

provide information on the reversals are unable to be specific on this question, but a near-zero magnetic field appears to last no longer than a few thousand years, if that.”

[H-III:14.5.3] “The scenarios for the geospace climate dependence on tilt angle between the Earth’s rotational axis and dipole axis provide vivid geometries of geospace regions such as the plasmasphere, plasma sheet, cusps, auroral zones, and open/closed field line regions. For extreme tilt angles, a significant question would be how rapidly these region can evolve and replenish themselves. [...] Figure 13.8 provides a pair of noon-midnight cross sections through the Earth’s magnetosphere for a 0° and 30° tilt. In the left panel, all the conventional magnetosphere regions can be identified and their evolution over a day would be shown at each time, as seen in this panel for a constant solar wind. Our present-day tilt scenario is somewhat different; in the northern hemisphere it is approximately 10° while in the southern hemisphere it is almost 15° .

However, even with this tilt, the fundamental magnetospheric regions found in Figure 13.8(left) are present all day with relatively small wobbles in the geocentric-solar-ecliptic coordinate system (GSE; x , Earth-Sun line; z , ecliptic north pole) of this figure. Both cusps are dayside and wobble in latitude. The plasmaspheric equatorial plane is that of the ‘average’ dipole and would wobble in 24 hours about the GSE- x axis. Even today, the concept of the plasmasphere’s ‘average’ dipole orientation is not fully explored since it is well

known that the Earth's equatorial fields are not well represented by a pure dipole component.

Over time scales of decades and more, the tilt angle as well as its geographic longitude wander. Indeed, this has been identified as a major factor in complicating the historic auroral observation data base. For example, when an aurora was observed at lower mid-latitudes as described in Sect. 13.2, was this due to an especially strong or geoeffective solar storm (CME) or did the Earth's dipole tilt have a particularly large value at that time, making this terrestrial location a much higher geomagnetic latitude?

The right-hand panel in Figure 13.8 shows the magnetospheric geometry for a specific 'UT' during northern-hemisphere summer solstice when the tilt can reach some 30° . At other times of the day, as the Earth rotates, this geometry changes significantly. Six hours earlier or later, the $x - z$ GSE cross section might look similar to the symmetric geometry in the left panel. However, the cusps would be displaced in the y GSE direction and the plasma sheet would have a large tilt in $y - z$ GSE cross section. As the tilt angle increases beyond 35° , would the normal diurnal independences of the magnetospheric morphologies remain? For example, would auroral zones still be referred to as a north and a south auroral oval? In the extreme case of a tilt approaching 90° , does the plasma sheet in the $x - z$ GSE cross section have two plasma sheets at certain UTs? Under these conditions with the same VADM and solar wind, dramatically different geospace climate would be observed in the form of auroral sightings as well as terrestrial magnetic field records of the electrojets and ring currents."

13.2.4 Geospace climate dependence on the solar wind

What are the effects of the long-term evolution of the solar wind on the geospace climate? [H-III:14.4] "Most of today's knowledge of the early Sun's history, normally referred to times that the Sun reached its zero-age main sequence (ZAMS) has been obtained from studies of Sun-like stars, *i.e.*, main-sequence G and K stars. [The] time dependences for the solar wind velocity (v_{sw}) and density (n_{sw}) at 1 AU [have been coarsely approximated by]:

$$v_{\text{sw}} = v_* \left[1 + \frac{t}{\tau_{\text{sw}}} \right]^{-0.4} ; \quad n_{\text{sw}} = n_* \left[1 + \frac{t}{\tau_{\text{sw}}} \right]^{-1.5}, \quad (13.14)$$

where $v_* = 3200 \text{ km/s}$, $n_* = 2.4 \times 10^4 \text{ cm}^{-3}$ and $\tau_{\text{sw}} = 2.56 \times 10^7 \text{ yr}$. [...]"

A:187

¹⁸⁷ Activity: Assuming a similar geomagnetic field, use the expressions in Eq. (13.14) to derive an estimate of the magnetopause distance over time. Show that for a young Sun this comes down to $\sim 1.25R_\oplus$ (with Eq. 5.22).

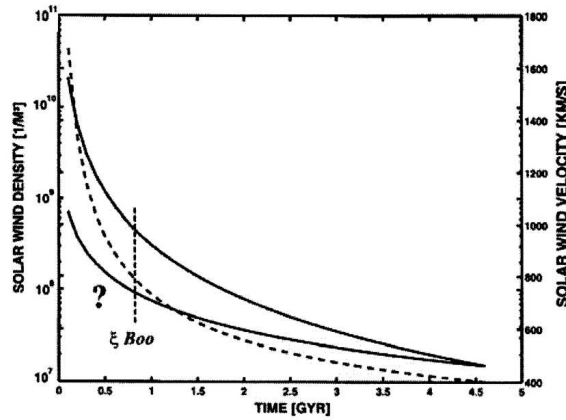


Fig. 13.9. Evolution of the observation-based minimum and maximum stellar wind densities (in units of m^{-3}) scaled to 1 AU (left scale; solid lines) obtained from several nearby solar-like stars. On the right scale is the solar wind speed for the stellar wind evolution (dashed line). [The star ξ Boo sits on the 'wind dividing line'; mass loss appears strongly reduced for stars younger than that, contrary to the simple extrapolation shown here, see Fig. 10.9. Fig. H-III:14.14; source: Lundin et al. (2007).]

Figure 13.9 shows the large spread in the range of the n_{sw} dependence since the Sun reached ZAMS about 4.6 Gyr ago. Note that the [above approximations and the curves] shown in Figure 13.9 have as a reference a present-day n_{sw} of 20 cm^{-3} ($2 \times 10^7 \text{ m}^{-3}$) and a $v_{sw} = 400 \text{ km/s}$. Today's solar cycle and solar storms have periods when the density can almost be a factor of 10 higher and the velocity reaches 1000 km/s. These enhanced conditions are associated with storms and superstorms in geospace that can persist for days while the solar wind remains perturbed. If the Earth's intrinsic magnetic field were then as it is today, would geospace at 2 to 3 Gyr ago be in a continuous superstorm state? Figure 13.9 shows that at these times the solar wind's pressure would permanently be at, or exceed, superstorm solar wind conditions. Would the auroral phenomena be permanent displays and exist to very low latitudes, or would perhaps M-I coupling require an unsustainable flow of ionospheric plasma into the magnetosphere? The past geospace climate over the Holocene, *i.e.*, the human time period, was not significantly affected by long-term changes in the solar wind while at very early ages it could well have been a very illuminating dynamic M-I coupling environment."

14

Magnetic fields and cosmic rays over time

Energetic particles can affect electronics components and presents a health hazard for astronauts, particularly when outside the magnetosphere, such as en route to the Moon or to Mars (as described in Chs. H-II:13 and H-II:14). The energetic particles discussed in this chapter originate at the Sun, in the solar wind, and in the Galaxy beyond the heliosphere. Observation of their variability tells us about their sources and about conditions they encountered between their origin and their detection.

Collisions of the most energetic among these particles with bodies in the Solar System lead to the formation of radionuclides that subsequently decay with half-lives of various durations. If such radionuclides are stored in suitably 'stratified' natural archives – such as in long-lived snow deposits or growth rings of trees – their concentrations measured through such archives shine a light on intensities and variability in times from before instrumental records. Even unstratified 'archives', such as lunar and meteoric rocks, provide information as dosimeters in which depth profiles of cosmogenic radionuclide contain information on particle energies and fluxes, and the different decay time scales some information on the integrated exposure history.

In the present chapter, we focus on changes in exposure over time scales up to billions of years, and on what these changes tell us about solar activity, galactic cosmic rays, the state of the heliosphere, and the terrestrial magnetic field.

H-IV:12] “Energetic particles in the energy range between 1 eV to 10^{20} eV can be found everywhere in our Solar System as sketched in Fig. 14.1. Their sources can be either outside our Solar System from galactic [and extra-galactic interstellar space or inside our Solar System, including the Sun, interplanetary space, and planetary magnetospheres]. The types of energetic particles range from electrons to charged atoms and molecules to neutral atoms and molecules

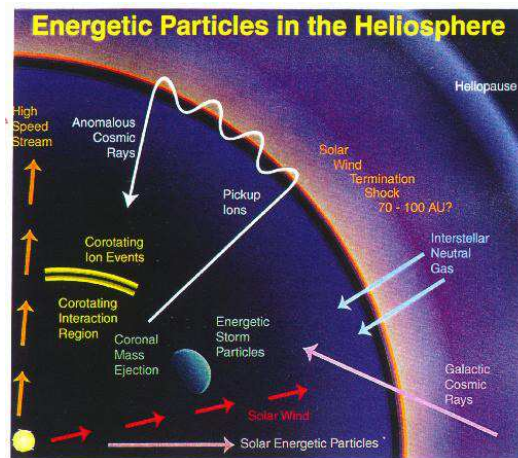


Fig. 14.1. Sources of energetic particles in the heliosphere. [H-IV:12.1; source.]

as well as dust particles. Fig. 8.5 shows the particle intensity versus energy of various types of energetic particles (left) and for cosmic rays (right).”

14.1 Long-term energetic-particle exposure of Earth

The exposure of Earth to energetic particles over many millions of years can be derived from the study of various cosmogenic radionuclides found in rocks (including those brought back from the Moon by the Apollo astronauts) and ice deposits (in Greenland and Antarctica). On time scales up to thousands of years we find signals in the biosphere (primarily in tree rings).

14.1.1 Generation of cosmogenic radionuclides

[H-III:11.4.2] “Cosmogenic radionuclides are produced by nuclear interactions of the galactic cosmic rays (GCR) with atoms (N, O, Ar) in the atmosphere; the contribution of solar cosmic rays is negligible because of their low energies.”

[H-III:13.2.1] “Galactic cosmic radiation originates outside our Solar System but generally within our Milky Way galaxy and is treated as [isotropic. This radiation consists of atoms that] have been ionized and accelerated to very high energies, probably by [shock fronts of] supernova remnants. The GCR population consists of about 87% protons and 12% α particles, with the remaining 1-2% heavier nuclei with charges ranging from 3 (lithium) to about 28 (nickel). Ions heavier than nickel are also present, but they are rare. Electrons and positrons constitute about 1% of the overall GCR.”

[H-III:11.4.2] “To reach the atmosphere, the GCR have to propagate through

the heliosphere which forms a bubble with a radius of about 100 AU around the Sun that is filled with solar plasma carrying magnetic field [as discussed in Ch. 8. It is appropriate] to use the transport equation [Eq. (8.25)] to parameterize the intensity of the GCR, however the so-called force-field approximation has proven to be a [reasonable simplification] near Earth. This approximation describes the modulation effect of the Sun on the energy spectrum of the GCR in terms of a parameter Φ called the solar modulation function” and comes about as follows:

[H-III:9.5] “If one assumes spherical symmetry, and then considers only high energy cosmic rays, for which the dimensionless modulation quantity rv/κ , which measures the strength of the modulation, is small, a very simple analytic solution can be obtained. The form of this solution corresponds exactly to that obtained for charged particles influenced by [an electric] field with a potential given as a function of heliocentric radius r by

$$\Phi(r) \propto \int_r^D [v_{\text{sw}}/\kappa] dr \quad (14.1)$$

[with the solar wind velocity v_{sw} assumed constant and κ an equivalent radial diffusion coefficient, also assumed to be a simple constant, for the full expression in Eq. (8.20).] {A:¹⁸⁸} Note that this is not a real electrostatic potential because it affects positively and negatively charged particles in the same way. [Moreover, observations show a strong dependence of the cosmic-ray intensity on heliographic latitude which cannot develop in the force-field approximation.] Attempts to fit the data yield values of $\Phi \approx 300$ MeV near 1 AU. Because of the use of an effective potential energy, this approximation is called the ‘force-field’ solution. [...]” [H-III:11.4.2] “The solar modulation function Φ basically corresponds to the average energy lost by a cosmic ray proton on its way to the Earth.

A:188

Figure 14.2 shows the differential energy spectrum of the GCR proton flux for different levels of solar activity. A value of $\Phi = 0$ MeV corresponds to the local interstellar spectrum outside the heliosphere. [The spectrum shown here is an estimate (dependent on the model approximation and on the properties of the interstellar medium and the solar wind) made before Voyagers 1 and 2 had reached the interstellar medium, which appears to have happened in late 2018.] Figure 14.2 shows that the shielding effects of the open solar magnetic field and the advecting solar wind are most pronounced at the low energy end of the spectrum. As a consequence, GCR particles above about 20 GeV are hardly affected by the varying heliospheric magnetic field.

Before reaching [the terrestrial atmosphere], the cosmic ray particles have to

¹⁸⁸ Activity: Advanced: If you are interested in how Eq. (8.20) can be approximated by something like Eq. (14.1) you can find the origin of this transformation in a study by Gleeson and Axford (1968).

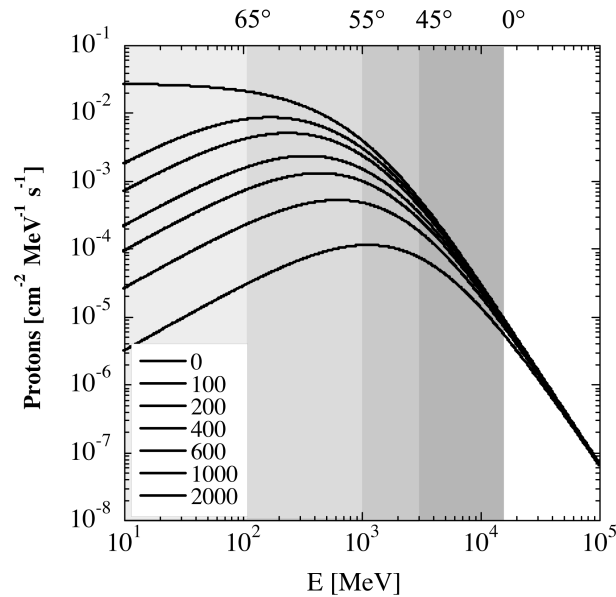


Fig. 14.2. Differential galactic cosmic ray proton fluxes for different levels of solar activity ranging from a value of the solar modulation function $\Phi = 0$ MeV (Eq. 14.1), corresponding to the local interstellar spectrum arriving at Earth without any solar influence, to $\Phi = 2000$ MeV which corresponds to a very active Sun. There are similar curves for cosmic ray alpha particles and heavier nuclides. The vertical bands illustrate the effect of the geomagnetic field which cuts off all protons approaching vertically with an energy below about 100 MeV for a geomagnetic latitude of 65° ; below 1 GeV for 55° , and below 3 GeV for 45° . At 0° the cut-off energy is 13.9 GeV for the present geomagnetic field. [Fig. H-III:11.11]

overcome a second barrier, the geomagnetic field. This field prevents particles with too low rigidity (momentum per unit charge) from reaching the top of the atmosphere. In a first approximation, the geomagnetic field is considered as a dipole and in this case the cut-off rigidity depends only on the angle of incidence and the geomagnetic latitude. At low latitudes the cut-off rigidity for vertical incidence is presently ~ 14.9 GV. This means that a cosmic ray proton needs a kinetic energy of at least 14 GeV ($14.9 \cdot m_p c^2$) to reach the top of the atmosphere (see shaded bands in Fig. 14.2) [...]

If a primary cosmic ray particle makes its way through the heliosphere and the geomagnetic field and enters the atmosphere it will interact quickly with an atomic nucleus of oxygen, nitrogen, or argon. Because the energies of incoming particles are generally very high, only part of their kinetic energy is transferred to the first atom they hit. They continue their travel and hit a few more atoms until their energy is dissipated. Each collision results in the generation of secondary particles covering the full spectrum of hadrons and leptons, which

Table 14.1. *Main properties for some cosmogenic radionuclides including nuclear production reactions and globally-averaged production rates for the present geomagnetic field strength and a solar modulation of $\Phi = 550$ MeV. (EC: electron capture). All nuclear reactions are induced by high-energy secondary particles generated by the primary cosmic-ray particles (so-called spallation reactions). The only exception is ^{14}C which is almost totally produced by thermal neutrons interacting with nitrogen. [Table H-III:11.4]*

Isotope	half life (yr)	decay	target	nuclear reaction	production rate ($\text{cm}^{-2}\text{s}^{-1}$)
^{14}C	5730	β^-	N,O	$^{14}\text{N}(\text{n,p})^{14}\text{C}$ $^{16}\text{O}(\text{p,3p})^{14}\text{C}$ $^{16}\text{O}(\text{n,2p1n})^{14}\text{C}$	2.02
^{10}Be	1.5×10^6	β^-	N,O	$^{14}\text{N}(\text{n,3p2n})^{10}\text{Be}$ $^{14}\text{N}(\text{p,4p1n})^{10}\text{Be}$ $^{18}\text{O}(\text{n,4p3n})^{10}\text{Be}$ $^{18}\text{O}(\text{p,5p2n})^{10}\text{Be}$	0.018
^{36}Cl	0.30×10^6	β^- , EC	Ar	$^{40}\text{Ar}(\text{n,1p4n})^{36}\text{Cl}$ $^{40}\text{Ar}(\text{p,2p3n})^{36}\text{Cl}$ $^{36}\text{Ar}(\text{n,p})^{36}\text{Cl}$	0.0019

either decay or interact with other atoms of the atmosphere. In this way a cascade of secondary particles develops which can be simulated using Monte Carlo techniques. Table 14.1 shows the different production reactions for the radionuclides ^{14}C , ^{10}Be , and ^{36}Cl , and the resulting mean global production rates for the present geomagnetic field intensity and a solar modulation function equal to $\Phi = 550$ MeV. {A:^[189]}

A:189

The simulations show that the majority of the secondaries are neutrons followed by protons. Both, in turn, collide with atmospheric atoms initiating spallation reactions that generate the cosmogenic nuclides that are archived for us in ice (^{10}Be , ^{36}Cl) or tree rings (^{14}C). In addition, the cosmic-ray produced neutrons have been monitored continuously since 1951 by so-called neutron monitors. [These measurements show that whenever] the magnetic activity is high (at high sunspot count) the shielding is strong and the neutron flux is low. [...] Many studies have shown that the 11-yr and longer-term variations are faithfully reproduced in the cosmogenic data, and they and the neutron

¹⁸⁹ Activity: To appreciate how little radionuclide material there is to work with, compute the global annual production in kg for ^{14}C and ^{10}Be . That production rate puts roughly one ^{14}C atom per 10^{12} atoms of ^{12}C in living tissue through uptake of atmospheric CO_2 by plants and their subsequent consumption by animals.

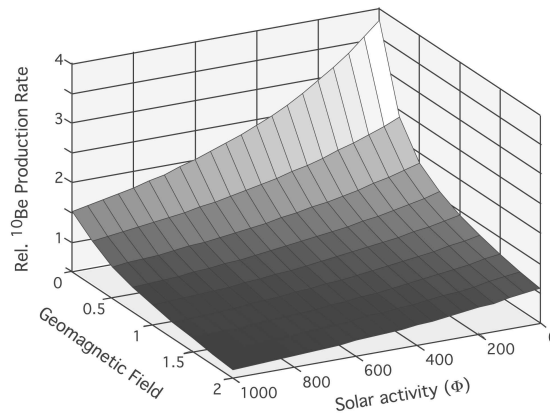


Fig. 14.3. Dependence of the ^{10}Be production rate on the geomagnetic field intensity (in units relative to today's field) and the solar activity (expressed by the solar modulation function Φ , Eq. 14.1). The production rate is normalized to the present strength of the geomagnetic field and solar activity corresponding to a solar modulation function of 550 MeV (matching the long-term average). [Fig. H-III:11.12]

monitor data have been inter-calibrated to yield a continuous cosmic-ray record for the past 10,000 years.

As an example, the combined effect of solar activity and geomagnetic field on the relative production rate of cosmogenic nuclides is shown for ^{10}Be in Fig. 14.3. The relative dipole component of the geomagnetic field μ_r varies between 0 and 2, 1 being the present field. For $\mu_r = 1$ and $\Phi = 550$ MeV (long-term average), the ^{10}Be production rate is normalized to 1. It should be noted that the dependence of the production rate on μ_r and Φ is nonlinear.”

{A:[190]}

A:190

14.1.2 Transport and deposition of cosmogenic radionuclides

[H-III:11.4.2] “The fate of a cosmogenic nuclide after its production in the atmosphere depends strongly on its geochemical properties. Within a short time, ^{10}Be becomes attached to aerosols and follows their pathways. ^{14}C on the other hand oxidizes to $^{14}\text{CO}_2$ and is exchanged between atmosphere, biosphere and ocean. After a mean residence time of 1 to 2 years, ^{10}Be is removed from the atmosphere mainly by wet precipitation. The flux F of cosmogenic nuclides from the atmosphere into, for example, a polar ice sheet is proportional to the atmospheric production rate Π : $F = \psi \Pi$. Locally and temporally, ψ can

¹⁹⁰ Activity: The ^{10}Be production rate for Mars would be about 2.5 times the terrestrial rate if it had a terrestrial atmosphere. Show why based on data in this text.

vary due to changes in the atmospheric transport and deposition processes. The degree of variability depends very much on how well the atmosphere is mixed. In the case of ^{14}C , the large atmospheric $^{14}\text{CO}_2$ reservoir leads to an atmospheric residence time of 6 to 7 years and therefore to a complete mixing. In the case of the aerosol-bound nuclides the residence times are shorter, roughly 1 – 2 years; mixing in the troposphere is not complete. After deposition, some of the nuclides become incorporated into natural archives such as ice sheets, glaciers, sediments, and tree rings.

For our purpose, a useful archive stores the complete flux of nuclides from the atmosphere in a stratigraphically undisturbed way and records the time accurately. Excellent archives in this respect are ice sheets which directly collect the atmospheric precipitation containing ^{10}Be . Typically, they cover the last several 10^5 years with a time resolution per sample ranging from 1 year at the top to decades or centuries near the bottom. However, due to the flow characteristics of ice, dating is difficult, especially in the deeper part of ice cores.

Tree rings represent an ideal archive for the atmospheric $^{14}\text{C}/^{12}\text{C}$ ratio. So far, by chronologically matching trees of different ages, the atmospheric $^{14}\text{C}/^{12}\text{C}$ ratio has been reconstructed back to approximately [14 kyr BP (before present, relative to 1950 CE)]. Potentially, the full range covered with today's measuring techniques (40 to 50 kyr) will be traceable in tree rings in the future.”

A:191

{A:[191]}

14.2 Radionuclides as proxies of magnetic variability

[H-III:11.4.2.1] “What can be learned by measuring cosmogenic nuclides in ice? [...] In an archive, changes in the concentration can result from changes either in the production rate Π or in the Earth-system processes ψ (transport and deposition). Changes due to radioactive decay can be corrected for, if a reliable time scale is available. Changes in the production rate can be caused by heliomagnetic and geomagnetic modulation of the cosmic-ray flux. Episodic solar proton events can cause short but intense cosmic radiation, but do not contribute much to the total production rate due to the relatively low proton energies. Changes in the system on the other hand are related to the atmospheric transport and mixing processes as well as to the local precipitation rate.

The question arises how the different causes of concentration changes can be separated. A straightforward answer to this question is to combine several nu-

¹⁹¹ Activity: Over the past century the concentration of ^{14}C in the biosphere has been dropping considerably because of fossil-fuel burning (why?). Express in functional form how this leads to an ambiguity in ^{14}C dating if no other information on the age of an object is known.

clide records from different sites. Comparing ^{10}Be with ^{14}C permits separating production from system effects. Changes in the production rate due to helio- and geomagnetic modulation of the cosmic-ray flux are reflected both in ^{10}Be and in ^{14}C in a very similar way. Changes within the Earth system, however, are expected to affect ^{10}Be and ^{14}C in a completely different way because the geochemical behavior of these nuclides is fundamentally different. [... Of course, this argument is useful only over] the range of the radiocarbon dating (last 50 kyr) and requires a high precision [record of $\Delta^{14}\text{C}$ that] is not yet available for the period 13 – 50 kyr BP. The next step is to separate heliomagnetic and geomagnetic signals. In principle, these two signals could be separated by looking at two radionuclide records, one from the equator and one from the regions of the magnetic poles. Without latitudinal atmospheric mixing, the record from the magnetic pole would only reflect solar modulation because geomagnetic shielding disappears at high latitudes, whereas the signal in the equatorial record would be dominated by geomagnetic modulation. However, as a result of atmospheric mixing, this is not the case.

Solar modulation effects have been found in cores from Greenland and Antarctica. The same is true for geomagnetic modulation effects like for the Laschamp event at about 40 kyr BP, when the magnetic dipole field was close to zero. This event is present in the high latitude ice-cores from the Arctic and from Antarctica (GRIP, Vostok, Byrd, Dome C, Taylor Dome). Radionuclide records from low-latitude ice cores are still rare [and have a smaller potential due to dating and other problems.]

Another approach is to assume that solar modulation effects generally occur on shorter time scales than geomagnetically induced production changes. Applying low-pass filters with cut-off frequencies in the range of $1/2000$ and $1/3000\text{ yr}^{-1}$ on cosmogenic nuclide fluxes provides production signals in good agreement with paleomagnetic intensity records based on remanence measurements. {A:^[192]}

A:192

The task of separating the different causes of variability observed in radionuclide records is complicated by the fact that some of the causes are coupled. For example, changes in solar activity affect atmospheric processes and possibly also induce, to a smaller extent, climatic changes. Therefore, additional information from other measured parameters should be included to obtain a complete and consistent picture of what happened during the period of investigation. In the following, we discuss how the intensity of the geomagnetic dipole field and the solar variability can be derived from cosmogenic nuclides.”

¹⁹² Activity: Look up 'paleomagnetic dating' in relation to the 'remanence measurements' mentioned in Sect. 14.2.

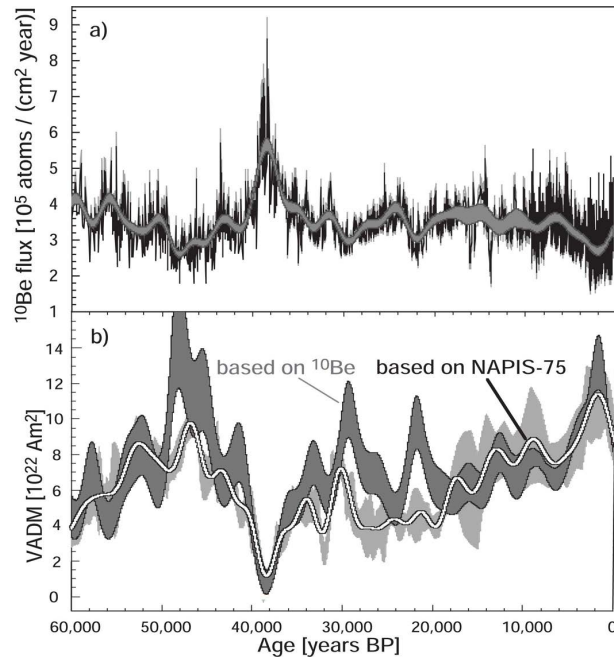


Fig. 14.4. Comparison of (a) ^{10}Be data with (b) the geomagnetic dipole field for the past 60,000 years. Panel (a) shows a compilation of ^{10}Be data from the GRIP and GISP ice cores in Greenland. Panel (b) compares the dipole field derived from ^{10}Be (panel a) to that from remanence data (NAPIS-75) measured in ocean sediment cores. [Fig. H-III:11.14]

14.2.1 Geomagnetic field

Fig. 14.4a shows [H-III:11.4.2.2] “a compilation of ^{10}Be data from the GRIP and the GISP ice cores drilled in central Greenland [...] covering the past 60,000 years. To correct for the lower precipitation rate during glacial times (10–60 kyr BP) the ^{10}Be flux has been calculated and smoothed (gray band). The plot shows a significant peak at about 40 kyr BP. To check whether the smoothed curve does reflect the geomagnetic dipole field as expected from Fig. 14.4 the corresponding changes in the dipole field intensity have been calculated based on its relationship with the ^{10}Be production shown in Fig. 14.3. The result is compared in Fig. 14.4b with the completely independent reconstruction NAPIS-75 which was derived from remanence measurements in Atlantic sediment cores. Overall the agreement is good and confirms that the ^{10}Be peak at 40 kyr BP corresponds to the Laschamp event when the dipole field intensity was almost zero but did not reverse.”

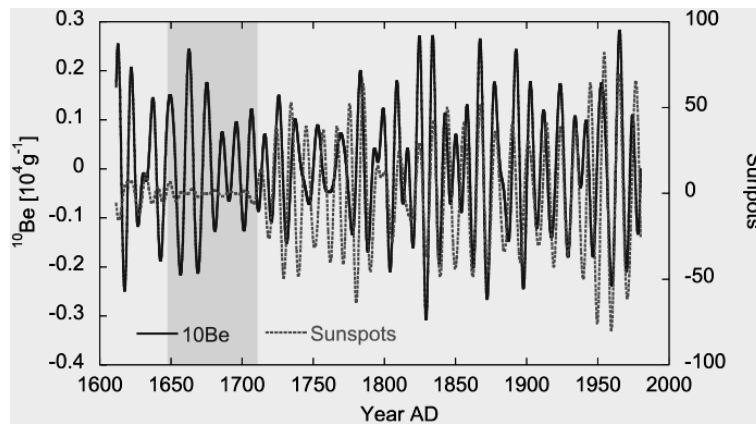


Fig. 14.5. Comparison of the ^{10}Be concentration measured in the Dye 3 ice core from Greenland with the sunspot number after applying a band-pass filter (8 – 16 years). Note that during the Maunder Minimum 1645 – 1715 (shaded area) when almost no sunspots were observed ^{10}Be shows a clear 11-yr sunspot cycle. [Fig. H-III:11.15; source: Beer et al. (1994).]

14.2.2 Solar variability

[H-III:11.4.3] “We return now to the discussion of solar variability and discuss to what extent cosmogenic radionuclides can expand our knowledge about long-term solar variability. In a first step, we compare annual ^{10}Be data with the sunspot record which represents the longest observational data of solar variability. A resolution of one year is about the limit because it corresponds to the mean travel time for a ^{10}Be atom produced in the atmosphere to reach the Earth surface where it is stored in, for example, an ice sheet. Fig. 14.5 shows a comparison of the ^{10}Be concentration from Dye 3, Greenland, with the sunspot number. Both records have been band-pass filtered (8 – 16 yr). While during the Maunder Minimum (shaded area between 1645 and 1715) hardly any sunspots were observed, the solar dynamo clearly continued to produce open magnetic field modulating the cosmic rays and the ^{10}Be production.

The overall good agreement between ^{10}Be and sunspot numbers gives us confidence to extend the time interval over the Holocene, *i.e.*, about the last 10,000 years. During this period the climate was relatively stable compared to glacial times and therefore we can assume that transport and deposition effects did not disturb the production signal in the ^{10}Be record. This assumption is confirmed by global circulation model (GCM) runs which show that the transport effects were relatively stable during the climatic conditions prevailing during the Holocene. So, indeed, to a first approximation they can be neglected. This is not the case for the geomagnetic field which exhibits significant long-term changes.

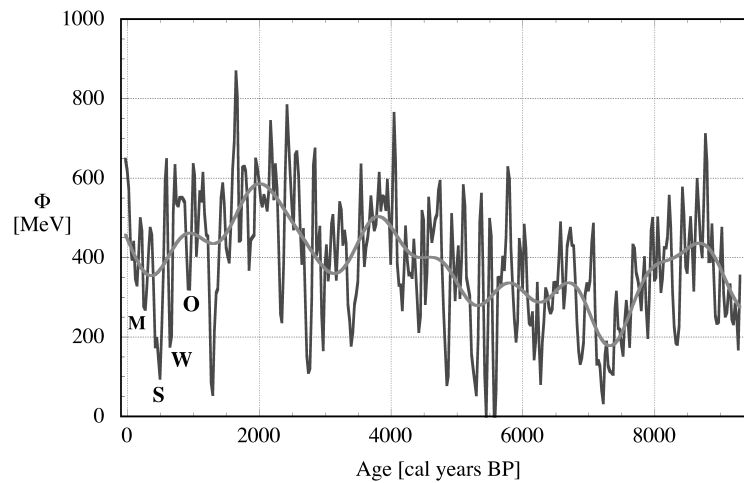


Fig. 14.6. Solar modulation function Φ from the present (0 BP corresponds to 1950) back to 9350 BP. The black curve shows data that have been low-pass filtered with a cut-off of 150 years; the smooth grey curve with 1000 years. The most recent solar minima are indicated: M: Maunder; S: Spörer; W: Wolf, and O: Oort. [Fig. H-III:11.16]

Using Monte Carlo simulations, the effect of the geomagnetic dipole field has been removed and we are left with the solar modulation function Φ [...] The data of Fig. 14.6 have been low-pass filtered with a 150 yr cutoff. The most striking features of the Φ record are the many distinct minima which correspond to grand solar minima such as the Maunder (M), Spörer (S), Wolf (W), and Oort (O) minima. The fact that Φ never reaches zero means that there is always some residual open magnetic flux; in other words the solar dynamo seems to weaken from time to time, but, as a close inspection of the unfiltered data shows, it never stops. The two exceptions in Fig. 14.6 are due to uncertainties in the data.

The maxima are less pronounced. It is interesting to note that the present level of solar activity is comparatively high, although there were earlier periods with similar or possibly even higher activity around 2000, 4000, and 9000 BP. There is also a clear long-term trend indicated by the thick line that is low-pass filtered with a cut-off of 1000 years.”

For periods covering 10^5 yr to over 10^6 yr the obvious radionuclide archives are ocean sediments that go back many millions of years, such as ^{10}Be and ^{26}Al . The price one pays for the long records is the reduced temporal resolution owing to the very small sedimentation rates and additional processes related to the transport of the radionuclide into the sediment. As is the case for tree ring records, these radionuclides are sequentially stored as they are taken from the atmosphere, in contrast to radionuclides measured in rocks, which we discuss

next, that are continually produced and that therefore provide only an integral measure of the production rate. The only time information available comes from the different half-lives. It is therefore important to measure as many distinct radionuclides as feasible.

14.2.3 Very-long time scale variability in cosmic-ray exposure

[H-III:9.2.4.2] “The record of the galactic cosmic-ray flux on a million-year time scale can be inferred from induced nuclear reactions in extraterrestrial matter of known exposure geometry, such as lunar rocks or meteorites. Nuclear reactions produce a variety of radioactive and stable nuclei that can be measured and related to the incident cosmic ray flux. The radionuclides ^{81}Kr (2.1×10^2 yr half-life), ^{36}Cl (3.0×10^5 yr), ^{26}Al (7.2×10^5 yr), ^{10}Be (1.6×10^6 yr) and ^{53}Mn (3.7×10^6 yr) represent a good set of monitors for cosmic ray flux variations on this time scale. Among the chondritic meteorites which were studied extensively, the production rates of the above radionuclides can vary because of differences in size and shielding conditions. These, when analyzed, reflect a constant ($\pm 10 - 15\%$) galactic flux over the $10^5 - 10^7$ yr time scale, which matches the average present-day flux.”

[H-III:9.2.4.3] “There are few radioisotopes with appropriate half-lives that can be used for [times scale of $10^7 - 10^9$ yr] and only ^{129}I (1.6×10^7 yr) and ^{40}K (1.3×10^9 yr) have been studied so far.

Chondritic meteorites cannot be used to study variations in the cosmic ray flux on longer time scales, because their exposure ages (time between being formed and striking the Earth) are typically less than a few tens of million years. Fortunately, there are numerous recovered iron meteorites which were exposed in space as small bodies for up to two billion years since being formed, and which are well suited for this purpose. The measurement of all three isotopes of potassium permits the detection of the cosmic-ray-produced component which is superimposed to potassium initially present in the meteorite. For the period of 0.2 – 1.0 Gyr ago, essentially constant ^{38}Ar production rates are observed, and agreement between ages determined from ^{38}Ar and from ^{40}K and ^{41}K .” {A:^[193]} It appears (Wieler *et al.*, 2013) that not much has changed in terms of long-term variability in solar activity or long-term trends in GCR fluxes coming into the heliosphere over the past billion years, to within a factor of ~ 1.5 , based on a variety of radionuclides in meteoritic samples

A:193

¹⁹³ Activity: How are the decrease of stellar rotation speed, magnetic activity, and mass-loss rate on long time scales compatible with the ‘essentially constant’ GCR exposure over the past ≈ 1 Gyr? The answer has to do with the fact that the Sun is already an aged star, and can be traced to its relatively weak magnetic braking over the past 1 Gyr, and thus relatively little decrease in coronal activity and mass-loss rate. The limited impact on GCRs at Earth orbit over time also suggests that the heliospheric variability (leading to diffusive GCR scattering) has not changed too much. Estimate the changes over time using Eqs. (10.3) and (10.7), and Fig. 10.3.

and terrestrial sediments combined . . . but realize that variability below time scales of hundreds of thousands to millions of years cannot be detected within these records.

14.3 Exposure to supernovae

There is evidence of at least one nearby supernova in the very early formation phases of the Solar System when the terrestrial planets had yet to fully take shape. Then, the Sun was still embedded within its birth cluster, and it appears that one of its heavy siblings exploded prior to the cluster falling apart. {A:[194]} There are some stars relatively nearby that are candidates to go supernova in the distant future, but none so close that the explosion would directly affect the solar system. Indirect effects, however, are possible: in the case of a blast wave from a nearby supernova, for present-day conditions of the solar wind, pressure [H-IV:3.5] “balance between the supernova shock wave and the solar wind produces extreme heliosphere models that have the same physical structures as the models with the heliopause at 1.4 times the distance of the termination shock in the upwind direction but with both located very close to the Sun. [A] supernova located at ≈ 9 pc from the Sun would create a heliopause that penetrates to within 1 AU, subjecting the Earth to an infusion of supernova debris including iron and other heavy atoms. {A:[195]} The discovery of the radioisotope ^{60}Fe with a half-life of 1.5 Myr in a deep-sea ferromanganese crust and dated to 2.8 ± 0.4 Myr ago indicates that a nearby supernova explosion likely occurred [around that time (and perhaps another (Wallner *et al.*, 2016) some 6 – 9 Myr ago).]

The effect on the Earth of a nearby supernova and the effect on more distant planets from supernovae at distances up to 30 pc will include an increase in the amount of neutral hydrogen atoms, dust, supernova metals, and Galactic cosmic rays reaching the planet’s atmosphere. The latter would influence the planet’s magnetosphere and change the planet’s atmospheric chemistry, including the important molecule ozone.”

¹⁹⁴ Activity: Heavy stars evolve much faster than low-mass stars, and can, if heavy enough, explode in a supernova even as lower-mass stars and their planetary system forming within the same molecular cloud are still in their formative phases. Look up lifetimes and evolutionary pathways for stars of different masses. Also look up properties of clusters of stars in star-forming regions.

¹⁹⁵ Activity: Estimate how much stronger the dynamic pressure of the incoming supernova wave front needs to be than the present-day IMF, assuming comparable solar-wind properties, to push the heliospheric boundary to within 1 AU. See Activity 133.

Applied heliophysics, *mutatis mutandis*,

...

Our advancing understanding of the processes that are part of the network of Sun-planet connections in heliophysics is being applied and tested in innovative studies of star-exoplanet couplings. A small sampling from the already extensive and rapidly growing literature illustrates the fertile and diverse fields that heliophysics finds in astrophysics at large:

- For ultracool stars near the boundary of the stellar and substellar regimes, where coronal heating fails as the photosphere decouples from the magnetic field because of weak ionization, an analogy has been identified with Jupiter's magnetosphere in which failure of co-rotation introduces stresses in the field that ultimately lead to an auroral radio signature that can have a counterpart in ultracool stars (e.g., this paper by Schrijver (2009), and this work by Pineda *et al.* (2017)).
- For cool stars in general, signatures of stellar winds are evasive other than indirectly through magnetic braking over many millions of years, but where such a wind collides with the interstellar medium a neutral-hydrogen wall forms whose optical depth for, e.g., Lyman α radiation provides a field-independent measurement of mass-loss rates (see Wood *et al.* (2005)).
- Another signature of stellar winds may be found in observations of the bowshock around the magnetosphere/exosphere of a transiting exoplanet during pre-transit phases (see Cauley *et al.* (2015)).
- The stellar winds and their coupling with exoplanetary atmospheres can be modeled using codes developed for our Solar System, providing insight into star-planet couplings and exoplanetary magnetospheric activity for systems well outside the parameter domain of our own home in the local cosmos (such as for TRAPPIST-1 – see Fig. 15.1 – see Garraffo *et al.* (2017)).
- The exposure of exoplanets to the stellar equivalent of solar energetic particles is being estimated by applying astrospheric models that include turbulence by which energetic particles are scattered, which enables, for example, quan-

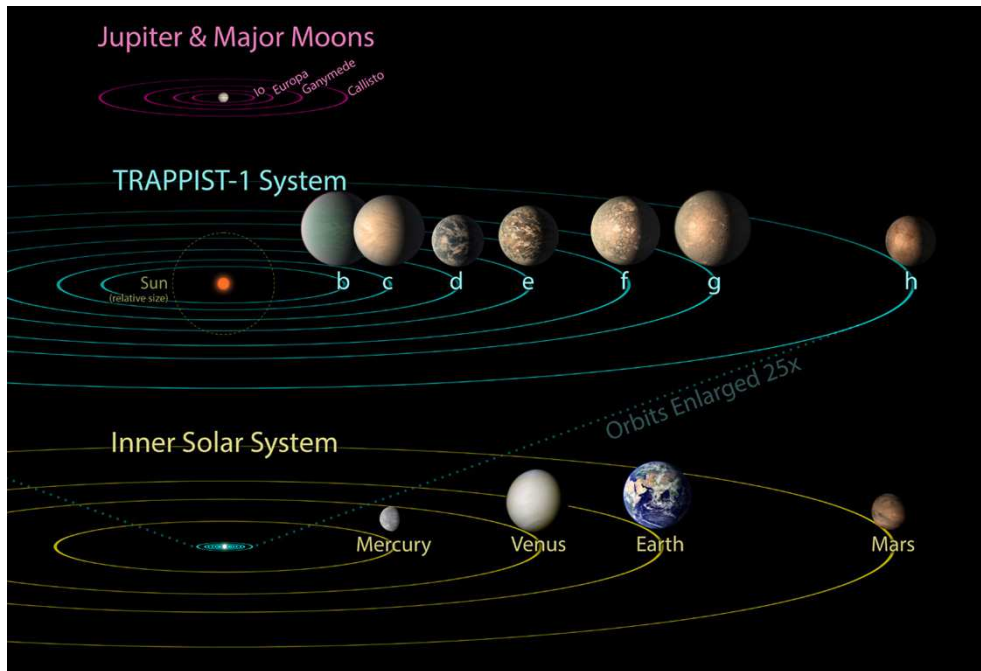


Fig. 15.1. Comparison of the orbits of the major moons of Jupiter, of the known TRAPPIST-1 planets, and the inner Solar System. The planets and moons are shown on the same scale, but vastly enlarged relative to the orbital scales. Courtesy: NASA/JPL-Caltech

titative estimates of particle radiation for exoplanets in stellar habitable zones (see Frascchetti *et al.* (2019)).

- The galactic cosmic ray exposure for exoplanets can be quantified based on modified heliospheric models, such as for the TRAPPIST-1 planets that orbit deep inside a strong-field astrosphere (see Struminsky *et al.* (2018)).
- The effects of stellar radiation and of stellar and galactic cosmic rays on the chemistry of exoplanetary atmospheres is being analyzed guided by, and calibrated against, solar-terrestrial spectral models and terrestrial tropospheric and ionospheric models (see Scheucher *et al.* (2018)).
- Evolution of the chemical makeup of planetary atmospheres subject to differences in insolation, rotation, atmospheric and oceanic circulations, and chemical weathering provide insight of the impacts of each of these on planetary habitability (see Jansen *et al.* (2019)) and help guide target selection for the search for biosignatures. The same is true on evolutionary timescales for the role of plate tectonics and volcanism (see Foley and Smye (2018)).

It should be no surprise that stellar and exoplanetary-system sciences conversely continue to provide crucial information for heliophysics:

- The very formation of stars and their planetary systems is being observed, models being refined, and results compared to empirical evidence from within the Solar System (e.g., see Lammer and Blanc (2018), and other papers in that volume).
- Over the past few decades, observations of stars of various spectral types and evolutionary phases have revealed the fundamental ingredients for stellar dynamos – rotation and convection –and quantitative information on how these set stellar atmospheric activity (e.g., see Schrijver and Zwaan (2000)) and how the Sun’s activity and wind have changed over its lifetime (see Güdel (2007), and see Ó Fionnagáin *et al.* (2019)), and from that how the Earth’s magnetopause distance would have changed over time (see Pognan *et al.* (2018)).
- The multitude of exoplanetary systems continues to clarify the formative processes of planetary systems in general, and of the Solar System in particular, with strong evidence for migration of planetary orbits, the impact of that on the formation of Mars, and the possible cause of the Late Heavy Bombardment and the transport of water to the terrestrial planets by gravitational scattering (e.g., see O’Brien *et al.* (2014)).

The exchange of concepts, knowledge, and models between the domains of heliophysics, stellar astrophysics, and exoplanetary sciences is only in its beginning phases: expanding observational and computational capabilities will propel these fields forward, catalyzed by a joint approach. One such area that requires a joint approach is already developing under the name of ‘transit light source effect’ (see Rackham *et al.* (2018), and Rackham *et al.* (2019b), and Rackham *et al.* (2019a)): exoplanetary transit spectroscopy (e.g., see Deming *et al.* (2019)) to study the chemical makeup and dynamics of an exoplanetary atmosphere is unavoidably linked with the analysis of the non-magnetic atmosphere (see Dravins *et al.* (2017a), and Dravins *et al.* (2017b)) and of the magnetic structures on the stellar atmosphere (e.g., see Pinhas *et al.* (2018); and Zhang *et al.* (2018)). Thus as we learn about exoplanetary atmospheres we will at the same time learn about starspots and stellar active regions at a resolution that has so far been unobtainable. That will provide information on how stellar dynamos structure their magnetic field from the scale of starspots upward that may prove critical to the development and validation of a predictive model for solar and stellar magnetic activity. That, in turn, will diminish the uncertainties on properties of stellar winds and their impacts over time on planetary atmospheres.

The purpose this book is to give you fundamental insights into the couplings between the Sun, its planets, and the interstellar medium, with particular focus on Venus, Earth, and Mars. The text could only give you an introduction to the multitude of aspects of the field of heliophysics and the links to other disciplines, although many of the related fields – such as nuclear physics, geophysics, biochemistry, meteorology, and radiative transfer – were only mentioned or implied in passing. With this book, you have the basic tools in hand to explore the Solar System over time, from its formation to its ultimate demise (although it stops short of the very end when the Sun transforms into a white dwarf, while the Earth may end up either evaporated and blown into interstellar space or pulled apart and spiraling into the white dwarf's atmosphere, or a combination of these). This book also gives you the tools to look outward and into the future: you have the basic concepts available now to explore the multitude of new worlds that are being discovered, analyzed, and inspected for potential signs of life (such as the TRAPPIST-1 system sketched in Fig. 15.1). It is an exciting time for all of that. Why not start exploring with these final five activities? {A:[196]} {A:[197]} {A:[198]} {A:[199]}

A:196

A:197

A:198

¹⁹⁶ Activity: **Observing exoplanetary atmospheres:** (1) Approximate the contrast $C_{*,p}(\lambda)$ between exoplanetary atmospheric radiation and stellar surface radiation, assuming both star and planet radiate as black bodies (using Planck's law $B(\lambda, T)$; ignoring center-to-limb effects), as function of wavelength, of the temperatures of star (T_*) and planet (T_p), and of the effective radii of star ($R_*(\lambda)$) and planet ($R_p(\lambda)$). Show that $C_{*,p}(\lambda = 10 \mu\text{m}) \approx (20 * R_*(\lambda)/R_p(\lambda))^2$ for $T_* = T_\odot$ and $T_p = T_\oplus$. (2) What is $C_{\odot,\oplus}(\lambda = 10 \mu\text{m})$? This value shows how hard it is to separate stellar and planetary signals (it is easier for closer-in, warmer, planets and for larger planets, such as hot Jupiters). (3) Why does the IR domain of the wavelength spectrum provide optimal access to the exoplanetary spectrum using a secondary eclipse (when the planet moves behind the star)? (4) At what wavelength does $B(\lambda, T)d\lambda$ peak for a planet at $T_p = T_\oplus$ (use Wien's displacement law)? (5) For transit spectroscopy, in contrast, optical wavelengths are most suitable for G- and K-type stars; why? See this tutorial by Deming *et al.* (2019) on exoplanet transit spectroscopy for answers and for much more on this topic.

¹⁹⁷ Activity: **Exoplanetary atmospheric spectroscopy:** How does wavelength-dependent transparency of an exoplanetary atmosphere lead to wavelength-dependent transit depths $\mathcal{T}(\lambda)$, and thereby yield spectral signatures of atmospheric chemicals? Basically, the apparent radius of exoplanet-plus-atmosphere depends on wavelength because the atmospheric opacity does. But the transit depth also depends on whether there are features on the stellar disk within the transit path. This provides information on, *e.g.*, starspot properties. Sketch how these two signals combine into the observed transit depth signal $\mathcal{T}(\lambda, t)$ over a transit. Consider how one might go about disentangling these two signals. See the reference in Activity 196 for more information.

¹⁹⁸ Activity: **Comparative heliophysics:** Look up the properties of the stars α CMa A, the Sun, and TRAPPIST-1. Consider how the following properties differ for a planet orbiting each of these stars within the continuously habitable zone (unconfirmed to exist in the case of α CMa A): (a) the size and color of the star, (b) the maximum possible age of the planet, (c) the duration of the orbital year and constraints on the length of the planetary day, (d) possible constraints on the planetary dynamo (subject to what we know about these at present), (e) the Alfvén Mach number of the stellar wind, (f) the magnetopause distance (assuming comparable planetary dynamos), (g) constraints on loss of planetary water, (h) the potential of measuring interstellar neutral hydrogen from an orbit near that planet, (i) the spectrum of the stellar and galactic cosmic rays (assuming the same spectrum external to the planetary system).

 {A:[200]}

A:199

A:200

¹⁹⁹ Activity: **A study on energetic particles in TRAPPIST-1:** TRAPPIST-1 is a very different world from our Solar System. The central star - itself only first observed in 1999 - is merely 1/8th the size of the Sun, only slightly larger than Jupiter. Its brightness is almost 2,000 times less than that of the Sun. The star is orbited by seven known exoplanets (first published on in 2016), much like Earth in size and mass, but all very close to their star. At least three of these seven planets are estimated to orbit within the liquid-water habitable zone. You can start reading up on TRAPPIST-1 using an ADS search, but for this Activity review this study by Garraffo *et al.* (2017) on the astrosphere, and this study by Frascetti *et al.* (2019) of the possibly very intense radiation environment of the planets. The role of heliophysics in this is evident throughout these studies: (1) identify the processes you have read about in this book that are elements of these studies. (2) Use what you learned in this text to explain why the wind is mostly sub-Alfvénic around the seven planets. (3) With dynamic pressures 3 to 6 orders of magnitude higher than for Earth, what does that do for the planetary magnetopause distances? (4) Although this system is diminutive, its astrosphere is potentially huge: estimate the distance to the astropause assuming the system is subject to ISM conditions similar to those for the Solar System.

²⁰⁰ Activity: **Arriving at Earth's climate from scratch:** In Activity 176 you compiled a list of all the processes involved in setting a planetary climate system that reflected at least all those mentioned in Chs. 11 and 12. Now, complement that list with the additional topics discussed in Ch. 14. Do not forget to add relevant thoughts from your notes for Activity 16! Then review that list and flag those processes that are beneficial to life as we know it on Earth and those that are detrimental to it. The duality of many, perhaps most, of the entries on your list should make you think about how our Earth, as it is in its present state, is a consequence of a remarkable interplay of often simultaneously beneficial and detrimental processes, including, perhaps, a series of fortuitous developments. Consider the extraordinary challenge of thinking about 'habitability' of any of the other thousands of exoplanets found to date, including what the phrase 'habitability' itself adds to that challenge given how little we know about life itself. Better yet, write an essay on this to share with fellow students, with teachers, and perhaps a much larger readership. After all, science is about communicating your thoughts and discoveries, as much to your peers as to society at large.

16

Compilation of activities found throughout the text

Chapter 1

- 1-p. 3: Look up what type of astrophysical body is a true 'star'. Contrast that to 'white dwarf star', 'brown dwarf star', and 'neutron star': none of these are true stars in their present state and only two of which have ever been. 'Brown dwarfs' take up the mass interval between true stars and (exo)planets.
- 2-p. 3: Look up the definition of 'planet'. Note that, formally, the term 'planet' has only been defined by the International Astronomical Union for bodies within the Solar System; the term 'exoplanet' is reserved for bodies like planets in other planetary systems, although for these, and certainly for the joint collective, the term 'planets' is often used.
- 3-p. 3: Many 'stars' we see in the night sky are binaries, including, for example, the brightest star in the night sky, Sirius (α CMa). More complex multiple-star systems may be less frequent, but are nonetheless common. Look up the example of Castor (α Gem) for an example of a sextenary, and then explore some more on multiple star systems in general.
- 4-p. 8: Remind yourself of Maxwell's equations that are mathematical renditions of these properties: (1) electric monopoles are linked with an electric field; (2) there are no magnetic monopoles; (3) variations in the magnetic field are associated with an electric field; and (4) a magnetic field implies either steady currents or time-dependent electric fields, or both. Good news: once we have reached magnetohydrodynamics in Ch. 3, Maxwell's equations are in principle superfluous as they are contained within the MHD equations; if you are interested in how that works, see here (Sections 1.1.1–1.1.9). By the way, a really useful resource for all things related to plasma physics (and how to convert between different unit systems) is the online NRL Plasma Formulary.

Chapter 2

- 5-p. 13: Planetary lower atmospheres are dominated by molecular substances, transitioning to atomic elements with a relatively low admixture of ions and electrons as one moves up through the ionospheres and thermospheres, while magnetospheres and the solar outer atmosphere and wind are comprised predominantly of charged particles. Compare thermal kinetic energies in different settings with molecular binding energies of, say, water and carbon dioxide. Also compare the energy of X-ray and EUV photons with ionization energies of atomic hydrogen and oxygen. See Tables 2.1, 2.3, and 2.4 for conditions in different settings.
- 6-p. 18: Look up and compare images of the Sun's magnetic field and atmosphere

- in different phases of the solar cycle, such as those obtained with the *HMI* and *AIA* instruments on *SDO* (NASA's *Solar Dynamics Observatory*). Note that such images are typically in false color, and with non-linear intensity scales to accommodate brightness contrasts.
- 7-p. 20: Consider why for a fully-ionized, hydrogen-dominated plasma we see $p = 2nkT$. For the answer, see below Eq. (2.7).
- 8-p. 20: At the solar surface we see a mean 'molecular mass' of $m \approx 1.3m_p$ while in the fully-ionized corona $m \approx 0.6m_p$ (for proton mass m_p). Explain why. (A hint: see Fig. 2.10.)
- 9-p. 21: Compute scale heights H_p in the Earth's atmosphere for molecular nitrogen (the dominant component) at a range of temperatures, and compare these with the value $H_{p\odot}$ for the atomic hydrogen-dominated gas in the solar photosphere, and for the CO₂-rich atmospheres of Venus and Mars. Use the data in Tables 2.1 and 2.3. Consider how the value of $H_{p\odot}/R_\odot$ contributes to the appearance of the Sun as having a well-defined surface. Also, consider why neutral, atomic hydrogen dominates in the solar photosphere (see Fig. 2.10 for the answer).
- 10-p. 21: One way to quantify the 'strength' of storms in different planetary atmospheres is to compare the dynamic pressure ρv^2 for the maximum surface winds listed in Table 2.1. Compare those values with the dynamic pressure in the solar wind using Table 2.4. Note: 1 bar = 10^6 dyne/cm².
- 11-p. 21: The fastest flows (in any direction, not only in the vertical, gravitationally stratified direction) that can be accommodated in a quasi-hydrostatic situation can be estimated from the fact that the gas pressure $p = \zeta nkT$ (with $\zeta = 1$ for a neutral gas and $\zeta = 2$ for a fully ionized hydrogen gas) should well exceed the flow's dynamic pressure ρv^2 . Look at Fig. 2.1 and add the horizontal lines where the two pressure terms are equal for a variety of flow velocities and corresponding temperatures; compare with the conditions discussed later in this chapter for the solar wind.
- 12-p. 22: With the values in Table 2.4, how long do the slow and fast solar-wind streams take to reach Earth? How many degrees does the Sun rotate between the moment these wind streams leave the Sun and the moment they arrive at Earth? How long for Neptune? Given that the wind flows out essentially radially, what is the apparent direction of the wind relative to the direction of the Sun as observed from the orbiting Earth (with an orbital velocity of about 30 km/s)?
- 13-p. 24: The momentum balance in Eq. (2.7) describes a radially-flowing wind over a non-rotating Sun. In reality, the Sun is rotating, and the magnetic field reaching into the heliosphere enforces the wind to co-rotate with the Sun, out to a distance where it becomes too weak to enforce such co-rotation. Show that for a sufficiently slowly rotating Sun, ignoring the centrifugal force is warranted. At what rotation period of a star like the Sun does the centrifugal force at, say, $2R_\odot$ counteract gravity by more than 10%? The centrifugal force in the wind would have been important for the very young Sun, see Sec. 10.2.1. Moreover, in the early phases of star-disk systems, centrifugal forces may be important in driving a cold wind; see Sect. 7.2.4.
- 14-p. 24: What powers the solar wind in the basic model discussed here? To see the answer, rewrite Eq. (2.10) to an energy equation with the terms for the kinetic and potential energy in the Sun's gravitational field, plus a term that reflects the work done by the expanding gas both geometrically and by acceleration; the energy for that expansion in the isothermal approximation is provided by the thermal conduction by the electron population. The real-world solar wind is not isothermal, certainly not far from the Sun (compare the coronal temperatures

in Table 2.3 with near-Earth wind properties in Table 2.4), and moreover is provided some additional power (in the form of heating and pressure) by waves and turbulence.

- 15-p. 25: In principle, Eq. (2.10) allows for an inflow: where v is negative, dv/dr needs to be of opposite sign also. This inflow, accelerating from infinity towards the star, is known as Bondi accretion. However, such inflow is unlikely to occur as an isothermal flow from infinity because the interstellar medium is typically cold, with low ionization and thus low heat conductivity by electrons. Consequently, compression would raise the temperature of the inflow. Moreover, be aware that the quasi-hydrostatic approximation fails for the inner regions of such an inflow, starting already well outside the critical point! Note that there is another class of solutions, namely a 'solar breeze': starting at low speed and never becoming transonic. Where does a 'solar breeze' reach its maximum velocity?
- 16-p. 29: **'What if' scenarios:** If you would like to think well 'outside the box' of things explicitly discussed in this book and in the Heliophysics volumes then consider this in the following chapters as you go along: what are things like when settings change? You could think of exoplanets with different host stars, orbits, and atmospheres, but there will be limited guidance by what we actually know from the literature. (1) For an example not too far from home, you could consider Titan, the only moon (natural satellite) in the solar system with a substantial atmosphere that is mostly N_2 (some 97%) and CH_4 (much of the remainder). *Activity:* Find Titan's equivalent values for the quantities listed in Table 2.1. *Further reading:* You can find publications on the (photo-)chemistry of its atmosphere leading to an ionosphere rich in $HCNH^+$ and $C_2H_5^+$. The chemical network in the high atmosphere leads to heavy organic molecules and aerosols that are deposited onto Titan's frozen surface and into its hydrocarbon lakes. Titan orbits within Saturn's magnetosphere, generally shielded from the direct impacts of the solar wind. However, the solar wind causes Saturn's magnetosphere to be highly asymmetric, and thus the environment through which Titan orbits is highly dependent on its orbital phase. Cosmic rays and energetic particles from Saturn's magnetosphere penetrate deep into Titan's atmosphere causing ionization and influencing chemical pathways. Titan has no intrinsic magnetic field (*i.e.*, no functioning magnetic dynamo) but an induced magnetosphere that changes as the moon orbits the rotating giant planet Saturn. There may be subsurface areas of liquid water, a water-ammonia mixture, or different mixtures in different locations and at different depths. Life might exist under these circumstances, and the traditional definition of 'habitability' as involving liquid surface water may need rethinking as we learn more. (2) For something far from home, consider the compact 7(?) -planet system of TRAPPIST-1 (see Fig. 15.1) on which much is being written: execute an ADS search for refereed papers with 'TRAPPIST-1' in the title. Task: let your imagination wander, read up on some of these things, and see how processes discussed in this volume apply to environments that are very different from those for Earth even though they are in some sense 'terrestrial'. Keep a running list of your thoughts as you read along for use later on!
- 17-p. 31: The fact that concentrations of atomic nitrogen are not shown in Fig. 2.5 should make you wonder given that molecular nitrogen is the most common species in the troposphere. Why is atomic nitrogen rare in the upper atmosphere? Hint: compare the molecular binding energies of nitrogen, oxygen, and water.
- 18-p. 33: Optical depth is an integral over absorption along a line of sight, and thus as useful for incoming as for outgoing radiation. Explain why the layers contributing most to the light from the solar photosphere are geometrically higher

- as you look away from disk center. What can you infer about the stratification of the solar atmosphere from the fact that the Sun (emitting close to black-body radiation over much of the optical spectrum) is brightest near disk center, darkening towards the limb? What follows directly from the fact that, on average, the solar corona seen in X-rays or extreme ultraviolet (EUV) has essentially double the intensity just outside the solar limb compared to that just inside the limb when there are no active regions along these lines of sight?
- 19-p. 33: You can think of the optical depth as the mean number of absorbers within the cross-section along a photon's path from infinity to height h . The probability of suffering zero absorptions, and thus making it to h , is $\exp(-\tau)$. The intensity at h is then an integral from infinity over the expected number of absorptions along the way. Combine that with Eq. (2.18) to derive Eqs. (2.19) and (2.20).
- 20-p. 33: A similar expression to Eq. (2.20) derived for photons holds for energetic particles (from, say, 1 keV/nucleon to 1 GeV/nucleon) losing their energy when propagating into a relatively dense medium (from the Earth's magnetosphere into its atmosphere, from the solar corona into its chromosphere or photosphere, or from interplanetary space into a spacecraft hull). Such energetic particles can penetrate a medium up to a column density of a few grams per cm^2 . Very roughly, estimate how far down that is into Earth's atmosphere, into the Martian atmosphere, into the solar lower atmosphere, and into an aluminum shell of a spacecraft. (See, *e.g.*, Sects. H-II:1.6, H-II:13.4, and H-II:14.4.)
- 21-p. 38: Work through Eqs. (2.22-2.25) to confirm that the effect of the collisions of charged particles in the ionosphere with the neutral thermospheric component is to rotate the net current from the direction of \mathbf{E} at high altitudes towards the negative $\mathbf{E} \times \mathbf{B}$ direction at low altitudes. As the expressions assume \mathbf{B} to be in the z direction, you could choose \mathbf{v} in the x direction to describe a horizontal velocity near the geomagnetic pole. In that same coordinate system, what is the direction of the current at about 125 km in the daytime terrestrial ionosphere where $\sigma_P \approx \sigma_H$ (see, *e.g.*, Fig. H-I:12.5).
- 22-p. 40: Look up what defines a 'sunspot' and what an 'active region'. A record of sunspot counts over many decades is shown in Fig. 4.5: what is the typical latitudinal range over which sunspots and sunspot groups occur? In Sect. 10.3.2 you will read about high-latitude and even polar starspots on rapidly-rotating, active stars, as the Sun would have been in its first few hundred million years.
- 23-p. 41: Use Table 2.5 to show that Eq. 2.26 yields r_{gi} for thermal motions. Then estimate energies of non-thermal particles so that their r_{gi} are comparable to the scale of the geomagnetic field (important for the terrestrial ring-current, which is a manifestation of particles drifting across the magnetic field because the heavy, energetic ones sense the gradient in the field strength; see Sect. 3.4) or perturbations in the heliospheric field (important for incoming cosmic rays, see Ch. 14). Compare with values in Table 3.4, compare these to mean-free path lengths there, and bear these results in mind going into Ch. 3.
- 24-p. 42: An intriguing property of dust is that, if the particles are small enough, radiation pressure is important in their momentum equation. Assuming neutral dust particles, estimate at what (density-dependent) size photon pressure from solar illumination exceeds solar gravity (note that this is independent of distance to the Sun for a completely transparent solar wind). There is a surprise here for dust of any size: the orbital motion of the dust causes photon absorption (and assumed isotropic re-radiation of that energy) to lead to a 'brake' on the orbital velocity, causing larger dust particles to spiral inward; look up 'Poynting-Robertson drag' to see how that works. From this, realize that dust needs to

be continually replenished somehow in the Solar System, generally by impact collisions and by disintegrating comets.

Chapter 3

- 25-p. 44: Look up magnetic maps of the solar surface (such as made with the HMI instrument on NASA's Solar Dynamics Observatory) and make a movie at an image cadence of a few hours; one option to do so is to use HelioViewer. Compare one for 2013–2014 (near cycle maximum, with multiple sunspot groups dispersing into the surrounding network of small-scale flux patterns) to 2017–2018 (around cycle minimum, with only the small scales on the disk).
- 26-p. 45: Throughout this volume we use 'field line' only for lines of force of the magnetic field. The concept can be applied to any field, however, including a flow field (such as in Fig. 4.10(C), then often referred to as streamlines) and the gravitational field. Field lines of \mathbf{B} and \mathbf{g} are fundamentally different in one key respect: a magnetic field line never ends (because there are no magnetic monopoles) while gravitational field lines start from a point of mass. What are starting points and/or endpoints (if any) of a system of electrical current (see Activity 28 for the answer)? And of electrical fields? As to magnetic field lines, note that there are drawings in this book, as in many other resources, where field lines are shown to start from one polarity and end on another. As magnetic field has no monopoles, such drawings should not be misread to mean that field lines end, but only that their rendering in the diagram is incomplete, *i.e.*, merely terminated for simplicity, for lack of information, or to restrict the discussion to a particular region of interest.
- 27-p. 47: An exception of sorts to the fact that field lines cannot begin or end in a divergence-free field lies in field lines that carry a 'null point', which is a point where the magnetic field goes to zero. Draw field lines around a pair of aligned but opposing magnetic dipoles in 2d and identify the null point(s). Then make a 3d rendering from a perspective away from the line connecting the dipoles. Visualize only the set of field lines going through the null(s) (these surfaces are called 'fans'). Such renderings with charges, nulls, fans (and their intersections, the 'spines') are useful tools in analyses of potential magnetic fields of a mixture of charges (such as bipolar solar regions on the solar surface). Consider how such 'fans' from nulls would not conflict with the field being divergence-free (the concept of 'measure' in set theory helps). For an introduction to the topology of the magnetic field, see Ch. H-I:4.
- 28-p. 49: Note that $\nabla \cdot \mathbf{B} = 0$ is not needed to complement the MHD equations in Table 3.3 as long as the initial condition satisfies that equation. Take the divergence of Eq. (3.3) to prove that. Use the same operation on $\nabla \times \mathbf{B} = \frac{4\pi}{c} \mathbf{j}$ to show that currents in MHD have no sources or sinks.
- 29-p. 53: The so-called 'Boussinesq approximation' is intermediate to fully compressible and incompressible, and in principle internally inconsistent: it assumes a fluid for which (and in numerical codes replaces Eq. (3.4) by) $\nabla \cdot \mathbf{v} = \mathbf{0}$ but allows density variations in the term in the force balance that includes gravity (and thus allows for buoyancy). This approximation works well if the flow can be characterized as 'nearly incompressible'. For settings where the scale of the density stratification is large compared to processes of interest the incompressible approximation can be valid; in such settings, compressibility becomes only important in structures like shock waves, but is ignorable if the flows are much slower than the sound and Alfvén speeds. Advanced, for the curious: In planetary atmospheric envelopes and stellar interiors alike, zones of relatively low temperature under relatively strong gravity are highly stratified compared to the scales of flows within them.

In such settings, numerical codes have been developed under the 'anelastic' approximation. This approximation provides a better description of the density in stratified settings than the pure Boussinesq one while filtering out sound waves that would require much higher spatio-temporal resolution of the code. This article by Durran and Arakawa (2007) introduces and compares several 'anelastic' approximations.

- 30-p. 55: **What if radiative transfer were included?** The MHD equations in Table 3.3 do not incorporate electromagnetic radiation. In a sufficiently dense medium, in which the photon mean free path is small compared to plasma and field gradients, energy transport by electromagnetic radiation can be described by a diffusion equation. Where the mean-free path is long, however, energy can 'jump' between different locations without (or with weak) coupling to the intermediate medium, in a manner that depends on wavelength as well as on atomic properties. With that in mind, contrast the solar interior to its atmosphere; a cloud-free planetary atmosphere to a (partially) clouded one; and (maybe once you get to Ch. 11) initial to later phases of star formation and of protoplanetary disks. In the context of that question and other assumptions going into the MHD equations in Table 3.3: Why is the solar chromosphere the hardest part of the solar interior and atmosphere to describe? And what makes a terrestrial ionosphere hard to capture in equations? Some of the answers to these questions will come as you read along. For an introduction to radiative transfer in stellar atmospheres, see this freely available online text by Rutten (2003): URL.
- 31-p. 56: Formulate the ion equivalent of Eq. (3.9) (remember Newton's third law, and with $m_e/m_i \rightarrow 0$) and derive Eq. (3.11) (using $m_e \mathbf{v}_i + m_i \mathbf{v}_e = m_i \mathbf{v}_i + m_e \mathbf{v}_e + m_i[\mathbf{v}_e - \mathbf{v}_i] + m_e[\mathbf{v}_i - \mathbf{v}_e]$) and also the corresponding momentum equation (absent gravity): $\rho d\mathbf{v}/dt = \mathbf{j} \times \mathbf{B} - \nabla p$. Then add gravity and compare to Eq. (3.5).
- 32-p. 56: Demonstrate, for a fully-ionized single-species plasma, the equivalence of Eq. (3.11) and $\mathbf{j} = \Sigma_e \cdot (\mathbf{E} + \frac{1}{c} \mathbf{v} \times \mathbf{B})$ with Eqs. (2.22)-(2.25).
- 33-p. 58: Look back at Fig. 2.1 and review the ranges shown of the value of the plasma β from Eq. (3.24) to get a feel for where plasma pressure gradients might dominate magnetic pressure gradients or vice versa. Add lines for unit plasma β for field strengths of $1 \mu\text{G}$ (as found in the outer heliosphere and interstellar medium; see Chs. 5 and 10) and for 0.1 MG (considered characteristic of the field strength of flux bundles at the bottom of the solar convective envelope where the principal processes in the solar dynamo are considered to operate; see Ch. 4).
- 34-p. 60: Compare values for c_s and v_A for the environments listed in Table 3.4.
- 35-p. 63: In solar physics, flux tubes are commonly used as an approximation of the state of the magnetic field in near-photospheric layers: embedded in a field-free atmosphere is a bundle of field separated from its surroundings by a thin current envelope. Assuming an ideal plasma without flows, show that the atmosphere within the tube is in hydrostatic equilibrium regardless of the path of the flux tube through the atmosphere. Show how pressure balance (incorporating both gas and magnetic components) determines the cross section of the tube.
- 36-p. 67: Units: this text uses cgs-Gaussian units. In other texts (including many of the Heliophysics chapters) you will find SI units. Look into conversions from one system to another (for example with the online NRL Plasma Formulary).
- 37-p. 71: Make comparisons of energy densities for the solar wind as in Sec. 3.5.2 at other bodies in the Solar System (using Table 5.2). Why comparisons of energy densities in planetary magnetic fields (Table 5.2) and in the surrounding solar wind are informative is discussed in Ch. 5. Why would you expect the flow energy

density and the magnetic field energy density to be comparable at only a few solar radii from the Sun?

Chapter 4

- 38-p. 72: The transition from radiative diffusion to convective enthalpy transport at the bottom of the convective envelope is gradual: the fraction of total energy carried as a diffusive flux gradually drops while that of the enthalpy flux smoothly increases, making convection the dominant transport about 35,000 km above the bottom of the convective envelope, or roughly after a single pressure scale height (see Sect. 4.3). Can you think of other terms that would be involved in the energy transport equation in a stellar convective envelope? A fair idea of the answer, along with a quantitative comparison of the relative importance of the processes involved in carrying energy through the convective envelope, can be found, for example, in this analysis by Brun *et al.* (2004), in particular their Fig. 3 (note that they show transport by convection that is resolved by their model and by (parameterized) unresolved – ‘subgrid-scale’ – convection).
- 39-p. 72: Figure 4.2 is a brightness-color diagram (known as a Hertzsprung-Russell, or HR, diagram) using typical astronomical units: absolute visual magnitude M_V , which is a logarithmic measure of stellar brightness, and spectral color $B - V$, which is the logarithm of the ratio of two brightnesses measured in different color bands (often using logarithmic brightness B and V , or less commonly R for blue, visual, and red). The table in that figure maps spectral type (see footnote vii), $B - V$, effective temperature T_{eff} and a correction factor BC that relates visual and bolometric brightnesses (see equations below that table). Using this information, estimate stellar radii R_* for Sirius A, ϵ Eri, 61 Cyg A, and AD Leo, realizing that $L_* = (\sigma T_{\text{eff}}^4)(4\pi R_*^2)$, with the Stefan-Boltzmann constant $\sigma = 5.7 \times 10^{-5} \text{ erg/cm}^2/\text{sec/deg}^4$. Sketch a double-logarithmic $L-T_{\text{eff}}$ version of the HR diagram and draw lines of constant radius in it. Then compare that to Fig. 10.1.
- 40-p. 78: How is the Sun’s magnetic field observed? Look up the effects on photons propagating through a plasma threaded by a magnetic field. This results in the ‘Zeeman effect’ of line splitting and of circular and linear polarization. For relatively weak field or relatively low wavelengths, the Zeeman splitting of ‘magnetically sensitive’ spectral lines is generally less than the thermal line width (and less than the Doppler width for rapidly rotating stars), so that what is in principle line splitting for individual atoms becomes line broadening when averaging over populations of atoms and over entire stellar disks.
- 41-p. 78: Compare a series of solar magnetograms over the past ~ 22 years (using, *e.g.*, *SOHO/MDI* and *SDO/HMI* observations). How do the magnetic patterns change over time in terms of overall activity, latitudinal distribution, polarity patterns on the northern versus southern hemisphere, ...?
- 42-p. 81: Work through how Eq. (4.1) is obtained by taking the dot product of Eq. (3.3) with \mathbf{B} , integrating over the total volume of the system, and assuming no Poynting flux or currents (or at most only a force-free field) leave the volume. Use vector identities ($\mathbf{a} \cdot (\nabla \times \mathbf{b}) = (\nabla \times \mathbf{a}) \cdot \mathbf{b} - \nabla \cdot (\mathbf{a} \times \mathbf{b})$, $\mathbf{a} \cdot (\mathbf{b} \times \mathbf{c}) = \mathbf{b} \cdot (\mathbf{c} \times \mathbf{a})$), Eq. (3.2), and Gauss’s theorem. For other vector calculus identities, see here.
- 43-p. 84: Look up sample images of solar granulation, the most easily detectable pattern of convection reaching into the solar surface layers. What are the characteristic length and time scales of granulation? Also look up the larger-scale flow patterns of mesogranulation and supergranulation.
- 44-p. 84: Estimate the time it takes for the solar equator to execute one more full rotation than the poles in the same time.

- 45-p. 84: Helioseismology uses resonant waves that run through the solar interior. These pressure-mode (or p mode) waves (generated by the turbulent convective motions) probe a range of depths depending on the wavelength and resonance conditions. At depth, downward traveling waves refract upward as the sound speed increases with temperature. If their frequency is below the 'acoustic cutoff period' around the photosphere upward traveling waves are reflected back into the interior, even as they are detectable around their upper turning point both in brightness (by compression and dilation) and velocity (through the Doppler effect on spectral lines). The combination of refraction and reflection leads to a cavity in which resonances occur. Intuitively, the cutoff frequency comes about because if the wavelength of a pressure wave exceeds a few pressure scale heights, there is essentially no restoring pressure force as the bulk of the atmospheric mass is simply lifted and lowered in response to the wave. Based on that argument, make a rough estimate of the acoustic cutoff period for the solar photosphere at around 5800 K (a later Activity will let you develop the relevant equations for an isothermal atmosphere). Waves with shorter periods continue to travel upward, while those with longer periods mostly reflect but partly tunnel through into the hotter chromosphere.
- 46-p. 84: To hear how helioseismology can measure rotation rates of stars (and, with enough different modes, of layers within stars) you can do the following experiment: Hold up a bell dangling from a string, strike it, and listen. Then twist up the string, let the bell spin freely, hit it and listen once more. The modulation in intensity that you hear for the spinning bell results from the beat of the Doppler effect working differentially on waves running with and against the spin direction. This is the essence of how helioseismology measures the Sun's internal rotation.
- 47-p. 85: If we take the Sun's polar field – averaging at cycle minimum at about 5 Gauss – how long would it take to wind that field into a strength of some 10^5 G – which is the estimated minimum field strength for flux bundles to survive their rise through the convective motions in the Sun's envelope – if the rotational shear were maximally used and if no back-reaction on that flow occurred? Hint: remember the field line stretch-and-fold from Fig. 4.6, look at the illustration in Fig. 4.9, and consider 'compound interest'.
- 48-p. 89: Consider that the assumption of a separation of scales as for the 'mean field dynamo theory' is also made in hydrodynamics when 'internal energy' (which includes the kinetic energy of the random motions of the gas particles) is 'separated from' the 'kinetic energy' of bulk motion. This assumption is commonly made with little consideration of why it works: there must be a scale that is small compared to flows of interest but large enough that low-order moments of the velocity distribution (like temperature and pressure) are defined by so many particles that there is negligible random noise when determined for a 'small' volume. Consider that in the context of the words in Table 3.2.
- 49-p. 90: Take the mean-field induction equation Eq. (4.12) and the expression for Eq. (4.13) as in Eq. (4.15) to find a mean-field form of the general induction equation Eq. (3.3). Group the 'diffusive' terms together. Estimate the order of magnitude of the advection, α , and diffusive terms. For these order of magnitude comparisons, approximate for the mean field $\nabla \approx 1/R_{\odot}$; in solar near-surface layers 'small-scale' random walk leads to $\beta \approx 300 \text{ km}^2/\text{s}$; the large-scale advective term of the surface meridional flow has an average value of order 5 m/s (peaking at about 15 m/s); estimate τ_{corr} from this value of β with Eq. (4.16), which corresponds to the characteristic evolutionary time scale of the dispersing supergranular

- convection; with that, estimate α using the characteristic supergranulation length scale of 30,000 km; then compare the order-of-magnitude values of the three terms expressed as time scales for the magnetic field. Note that the 'turbulent diffusivity' β far exceeds the 'resistive diffusivity' η in stellar dynamos (and see Activity 51 how the above helps in understanding how surface flux dispersal can be described quite well by a random-walk diffusive description).
- 50-p. 91: Relate the Rossby number in Eq. (4.2) to the dynamo number C_Ω and the magnetic Reynolds number \mathcal{R}_m in Eq. 4.20: $N_R = \mathcal{R}_m^2 / C_\Omega$.
- 51-p. 97: One of the basic concepts behind Babcock's idea is that magnetic field at the solar surface is largely advected like a scalar quantity. Consequently, the field disperses in the random motions of the surface convection (with an equivalent diffusion coefficient of $D \approx 250 \text{ km}^2/\text{s}$) subject to the large-scale advection of the differential rotation and meridional flow. To see how this can be, use the ideal version of Eq. (3.3) and assume that the field is always vertical to the surface (a good approximation to the observed photospheric field, except during emergence and cancellation; a result of the buoyancy of flux bundles – see Activity (35) – and show it is equivalent to Eq. (3.4) for the advection of a scalar (without the source and loss terms in that version). Note that this formulation is linear, so that you can think about N and S polarities as diffusing separately, then to sum to obtain the net result; this helps visualize why the active-region tilt angle is important in reversing the polar fields from cycle to cycle. Question: with this value of D , what is the characteristic time scale for flux to disperse over the solar surface (hint: Eq. 3.20)? With that in mind, how important is the meridional advection from equator to pole (with a characteristic velocity of 10 m/s) in transporting the field within the duration of a solar cycle? Remember that Activity 49 shows how the diffusion coefficient β associated with (super-)granular random walk adds to the molecular/resistive diffusion coefficient η .
- 52-p. 97: Joy's rule, that the leading polarities (in the direction of rotation) of active regions emerge statistically closer to the equator than the trailing polarities (see Fig. 4.4, for example), is the reason why eventually flux of the trailing polarity builds up a polar cap that reaches its maximum strength at cycle minimum. For some interval around that time, the bulk of the heliospheric field originates from the polar caps. Estimate the total flux in the solar wind (assuming an isotropic flux density at Earth orbit; use Table 2.4). This is the equivalent of only a few large active regions (Fig. 4.4) although it is in fact composed of a fraction of the flux from the ensemble of all bipolar regions emerging over a cycle.
- 53-p. 99: Why are solar photospheric flux tubes buoyant (hint: look back at Activity (35)? What is the maximum density contrast between interior and exterior? For an essentially evacuated flux tube at the solar surface, show that the buoyancy force per unit length (causing the tube to buoy towards vertical) dominates the dynamic pressure force exerted by a convective flow (which could bend the tube away from vertical) of $v = 1 \text{ km/s}$ for any tube with diameter $2a$ exceeding just a few km. Indeed, observations show flux tubes to be essentially vertical to the photosphere (except around emergence and collisional cancellation when magnetic curvature forces of the field arching from one polarity to the opposite one are strong).
- 54-p. 100: For the curious: This work by Lemerle and Charbonneau (2017) describes an interesting dynamo experiment in which the Babcock-Leighton concept is combined with surface flux transport modeling (see Activity 51) to create a quasi-regular dynamo in which convection-induced fluctuations on the tilt angle of emerging active regions (perturbations on Joy's rule, see Activity 52) provide the

stochastic noise that can lead to cycle-to-cycle differences and even extended periods of weak cycling (as in the Maunder Minimum period for the Sun), something also reported on by Karak and Miesch (2017) in this paper.

- 55-p. 100: Make a summary of the essential distinctions between the dynamo concepts discussed up to this point: $\alpha\Omega$ (with or without α quenching, which itself can be strong/catastrophic or not); interface; flux-transport; and Babcock-Leighton.
- 56-p. 106: The dynamo model in Fig. 4.13 is characterized by five dimensionless numbers. Two, the magnetic Reynolds number and the Rossby number, are defined in Eqs. (3.18) and (4.2), respectively. Look up the meaning of the other three: Ekman, magnetic Prandtl, and Rayleigh. These three are important numbers when computing the flows and their coupling, but not encountered until this Figure in this text because the solar dynamo models that were discussed and shown in the figures are kinematic, relying on a given, not consistently computed, flow pattern. The dynamo model behind Fig. 4.13, in contrast, computes flow, field, and their interaction, and thus also needs these remaining three dimensionless numbers specified.
- 57-p. 107: Summarize the contrast between the dynamo of a terrestrial planet with that of stars as discussed in this Ch. 4: consider, among others, flow speed, rotation period, stratification, differential rotation, and meridional advection.

Chapter 5

- 58-p. 111: To build a comparison of the different conditions, keep pen(cil) and paper at hand to sketch the various configurations as you read about how flows interact with bodies in the planetary system. Working in a reference frame in which the body is at rest, assume a spherical object, and let a flow move past it from left to right. Then prepare to make drawings in two orthogonal planes: the first plane is defined by the flow vector and the magnetic field carried in the flow (you may assume the field to be normal to the flow), while the second plane is normal to the first. Draw streamlines of the flow and subsequently add magnetic field lines. If you are good at 3-D renderings, also try a visualization such as in Fig. 5.11.
- 59-p. 117: Write the Eqs. (5.2) and (5.8)–(5.12) for the hydrodynamic limit, and derive the temperature ratio between the post- and pre-shock media. You should find that the density contrast $r_\rho = ((\gamma + 1)M_s^2)/(2 + (\gamma - 1)M_s^2)$ and the pressure ratio $r_p = (2\gamma M_s^2 - (\gamma - 1))/(\gamma + 1)$ where $M_s = v_1/c_{s1}$ for sound speed c_s . Note there is a maximum value for r_ρ but not for r_p as function of M_s . What are the values for r_ρ and r_p for $\gamma = 5/3$ for $M_s \downarrow 1$ and $M_s \gg 1$?
- 60-p. 121: Compare the radial dependence of the magnetic fields in this solar wind model with the values listed in Table 5.2. Also: use these dependences to demonstrate that the plasma β tends to a constant value far from the Sun.
- 61-p. 121: At what distance from the Sun does the above solar-wind model have $|B_r| = |B_\phi|$ for typical values of the slow and fast solar wind? What are typical values for B_ϕ/B_r at 1 AU, 5 AU (*cf.* Fig. 5.9), and at the ice giants?
- 62-p. 122: The solar wind stretches the high-coronal magnetic field into the heliosphere into a roughly radial field below the Alfvén radius. This enables an analogy with electrostatics: the field of electric charges placed above a flat perfect conductor can be computed by placing mirror charges opposite to the conducting surface, which then naturally has the electric field perfectly normal to the conducting surface. Analogously, in a magneto-static consideration above the spherical Sun, the magnetic field can be approximated by placing mirror 'charges' on a sphere at distance d_{SS}^2 which then has the field perfectly radial at d_{SS} . This is called the 'source surface model' with empirically $d_{SS} \approx 2.5R_\odot$ (where that

'source surface' is taken as the foundation of the heliospheric field; the virtual surface with mirror charges used to compute the potential field below d_{SS} is then at d_{SS}^2). This model (introduced by Schatten *et al.* (1969)) works remarkably well below d_{SS} on large scales. The heliospheric field is approximated by a radial continuation from that source surface, then subject to the Parker spiral. For illustration, simplify the source-surface model by a 2-d sketch involving a line of charges and another of mirror charges. Sketch the equivalent of the foundation of the heliospheric current sheet and examples of 'closed' field lines (the equivalent of coronal loops closing back onto the solar surface) and 'open' field lines (the equivalent of field stretched out into the heliosphere), at the base of which we find dark 'coronal holes' in X-ray images of the Sun.

- 63-p. 128: For the solar wind flowing onto a non-conducting sphere, use estimates of wave speeds to sketch the density wake, the slow-mode refilling, and the fast-mode rarefaction front in a plane defined by the flow vector and the field vector, and in a plane defined by the flow vector and perpendicular to the field. You may compare the result with measurements for the case of the Moon (in Fig. H-IV:10.7).
- 64-p. 130: On the largest scales, there may be a long magnetotail to the entire heliosphere, that may even be oblate because of the tension force of the interstellar magnetic field. Although alternative views propose a much shorter tail, making the heliosphere more like a bubble, it is illustrative to see how such a moderate flattening by the interstellar magnetic field might work. Have a look at, *e.g.*, McComas *et al.* (2013), in particular their Figure 9.
- 65-p. 134: Use Eqs. (5.2) and (5.8-5.12) to show that in the case of a strong shock (in which the thermal energy of the solar wind upstream of the bow shock can be ignored) the temperature just downwind of the bow shock is given by $(3m_p/32k)v_{sw}^2$ for a wind speed of v_{sw} , and that the density contrast across the shock is a factor of 4 (show that is true anywhere along the shock). Use this to estimate the angle from the upwind direction out to which the flow remains supersonic just inside the shock front (remembering that the transverse component of the velocity is unaffected by the shock).
- 66-p. 134: What is the expression for the temperature of the gas at the stagnation point on the magnetopause assuming that the flow continues adiabatically after the shock (*i.e.*, that it conserves the sum of bulk kinetic and thermal energies)? What is the value for $v_{sw} = 800$ km/s.
- 67-p. 135: Use Eq. (5.22) to show the scaling of R_{CF} with orbital radius, planetary magnetic field, and planetary radius.
- 68-p. 135: With the fastest recorded solar-wind gusts at $v_{sw} \approx 2500$ km/s, what is the required plasma density to push the magnetopause to within geosynchronous orbit according to Eq. (5.22)?
- 69-p. 136: Illustrative diagrams like Fig. 5.14 typically show the interplanetary magnetic field (IMF) as lying within the $x - z$ plane of such diagrams. In reality, the three IMF components $B_{x,y,z}$ are typically of comparable magnitude. Moreover, the orientation of the Earth's magnetic axis relative to the incoming wind changes in the course of the year. Consider how the diagrams should look when drawn in three dimensions for a few different combinations of $B_{x,y,z}$. Look up the 'Russell-McPherron effect' which attributes the semi-annual variations in geomagnetic activity largely to the relative orientation of the Earth's bipole axis: maximum geomagnetic activity around the equinoxes, minimum around solstices.
- 70-p. 139: Use Fig. 5.15 to estimate the Alfvén velocity in the jovian magnetosphere near Ganymede. First, estimate the flow speed of the incoming plasma relative to the moon realizing that the plasma is sub-corotating by about 80% of the speed

of corotation with Jupiter at Ganymede's orbit. Then use the geometry of the field shown in the figure to estimate the Alfvén velocity. The coordinate system for the simulation has the y -axis pointing towards Jupiter and the z -axis aligned with the jovian spin axis, and the units are expressed in Ganymede radii.

71-p. 143: Consider differences and similarities between 'corotation' in a planetary magnetosphere and in the solar wind, including (a) the absence of a sufficiently neutral atmosphere in the Sun to decouple the motions between internal and heliospheric fields (associated with a concept called 'line tying', which we touch upon in Sec. 6.3.1.1), and (b) the very term 'corotation' which to a heliospheric physicist does *not* include the component of $\mathbf{v} \parallel \mathbf{B}$ but is limited to the pattern of the field, not the plasma itself.

72-p. 145: The magnetospheric magnetic field cycle starts for the field with (1) day-side reconnection to the field in the wind, is then (2) followed by being dragged towards night-side, from there (3) moving into the magnetotail, and after (4) reconnection in the current sheet the field (5) moves back towards the day-side to replenish (at least on average over longer periods) the flux lost from there in the reconnection process. That loop, called the Dungey cycle, can be visualized from Fig. 5.14 if the succession of drawn field lines is interpreted as a sequence of events for a single field line (and realizing that step (5) has to occur over lower magnetic latitudes to avoid the field that is at the same time involved in step (2)). But during the cycle, the planet rotates underneath, dragging the ionospheric plasma onto which the magnetospheric field connects with it. To see for which planets this process is important, estimate for each of the planets: (a) the model-based magnetopause distance R_{CF} , (b) the time it takes to move from step (1) to a phase somewhere around steps (3) and (4), and (c) how many turns the planet has made in the meantime. Assume the following: that the solar wind speed averages to roughly the same value at all the planets (say $v_{sw} = 400$ km/s), and that the flow of the plasma in the deep magnetosheath carries the field from front to back over, say, $3R_{CF}$ at $0.1v_{sw}$. Info in Table 5.3. If you also want to do Ganymede, realize that it has spin-orbit synchronization within the jovian magnetosphere.

73-p. 145: Consider the possible equivalent of a Dungey cycle for the heliospheric field subject to reconnection with an interstellar magnetic field. What would happen in case there were no coronal heating?

74-p. 145: **What if?:** Many so-called 'hot Jupiters' have been found among the exoplanet population: giant planets that orbit very close to their parent stars. What would the estimated magnetopause distance $R_{CF,hJ}$ be if Jupiter were orbiting the present-day Sun at 0.05 AU? For a younger Sun (see Ch. 12) the solar wind would have been stronger, pushing $R_{CF,hJ}$ to below the orbital radius of Ganymede; describe what that would mean for this 'hot-Ganymede' moon?

75-p. 145: Advanced, for the curious: Things get more complicated when objects are smaller than the gyration radii of particles in the flow, or when ionization processes occur when neutral particles from an 'atmosphere' move into an approaching flow, or both. If you are interested in seeing how these complications play out, have a look at this study of comet 67P by Behar *et al.* (2017) using observations by the Rosetta spacecraft.

Chapter 6

76-p. 148: Consider how a non-potential state can arise or be strengthened in the solar atmosphere and in a magnetosphere, including the roles of plasma motions and induction. Eq. (4.1) is illustrative for the overall energy budget.

77-p. 149: The processes of electromagnetic radiation from a plasma involve three

fundamentally distinct processes: bound-bound, free-bound (radiative recombination), and free-free (Bremsstrahlung) emission. Aurorae and flare ribbons are caused by collisions of downward-propagating, energetic charged particles with the atmosphere below. Aurorae observed from the ground include both free-bound and bound-bound emission from ions and molecules, respectively (there is X-ray emission, too, but that does not penetrate to ground level). Look up which ions and molecules dominate in the terrestrial aurora, and which emission processes are involved with these. See also Activity 121 for the contrast with solar coronal emission.

- 78-p. 149: The average speed of a CME between Sun and Earth is close to 500 km/s while the fastest have speeds exceeding 3000 km/s. How long are the transit times from Sun to Earth? Compare the average and peak CME speeds to typical wind speeds (Table 2.4). Describe qualitatively what happens in the interaction with slow and fast wind streams for average CMEs and for the fastest CMEs.
- 79-p. 149: The phenomena discussed in Ch. 6 are all part of what is referred to as space weather. To explore how aspects of space weather are quantified review this NOAA site that lists the types of 'storms,' their potential effects, and their approximate frequency within a solar cycle. For current space weather conditions, forecasts, and more see this site of the Space Weather Prediction Center.
- 80-p. 156: Look up locations and properties of the Earth's (a) electron and proton radiation belts, (b) ring current, and (c) plasma sheet.
- 81-p. 162: The energy processed by the magnetosphere during a magnetic storm is of order $E_{\text{storm}} = 5 \times 10^{23}$ – 5×10^{24} erg from moderate storm to superstorm. Compare that to an order of magnitude estimate of the energy $E_{\text{mag},\oplus}$ contained in the geomagnetic field (by, say, using a scale of $3R_{\oplus}$ and a characteristic field strength of 0.1 G) and with the incoming total energy E_{sw} of the solar wind during the storm period (with typical conditions for the fast solar wind and an active cross section of πR_{CF}^2 , and a storm duration of 1-10 h). What are the values of $E_{\text{storm}}/E_{\text{mag},\oplus}$ and $E_{\text{storm}}/E_{\text{sw}}$? Compare these values to solar equivalents when you reach Activity 85.
- 82-p. 173: 'Chromospheric evaporation' is a misnomer because there is no phase transition involved: the heating of chromospheric material from $\approx 10^4$ K to of order $\approx 5 \times 10^6$ K causes the pressure and the associated pressure scale height to increase. What are the pre-heating and post-expansion scale heights for the above temperatures? How do these compare to the solar radius?
- 83-p. 173: For a given temperature, coronal soft X-ray brightness scales essentially with the square of the particle density. Why? Let a given coronal loop have an initial loop-top density n_0 at temperature T_0 and let an impulsive heating event change these to n_1 and T_1 . With $T_{0,1}$ within the range of about 0.4-30 MK the radiative losses scale as $P(T) \propto T^{-2/3}$. If the temperature changes from 1 MK to 5 MK and the density increases by a factor of 15, show that the ratio of radiative cooling time scales is close to unity. Conductive losses into the lower atmosphere, however, are larger at higher temperatures; why?
- 84-p. 177: Describe what is seen in Fig. 6.14: how can a CME be imaged, and why is that best done from space, or from a very high mountain top? Argue why the CME in this image is not likely to envelop Earth. What would an Earth-bound CME look like? Can you differentiate that from one moving in the opposite direction?
- 85-p. 177: The energy processed during a strong to intense solar flare and CME is of order $E_{\text{flare}} = 10^{30}$ – 10^{33} erg. Compare that to an order of magnitude estimate of the energy $E_{\text{AR},\odot}$ contained in the field of an active-region core (by, say, using

a scale of 30,000 km and a characteristic average magnetic flux density of 300 G). What is the value of $E_{\text{flare}}/E_{\text{AR},\odot}$? How does this compare to $E_{\text{storm}}/E_{\text{mag},\oplus}$ and $E_{\text{storm}}/E_{\text{sw}}$ in Activity 81?

86-p. 177: One phenomenon associated with many CMEs is a so-called 'coronal dimming', in which a large fraction of the quiet-Sun solar corona fades for some time. Think about the potential causes: temperature change (so the signal moves from one bandpass to another), quasi-adiabatic expansion of the coronal field, and entrainment of coronal plasma in the erupting CME. Estimate the volume of quiet-Sun corona (at a density of some 10^7 cm^{-3}) that would need to be involved if it were to move out with an erupting field configuration if that made up, say, 25% of the erupting mass of, for example, 10^{15} g .

87-p. 177: For a sense of scale: how many nuclear bombs are needed to match the energy released in a large solar flare of 10^{32} erg ?

88-p. 177: Advances in numerical capabilities are making a big difference in understanding magnetic instabilities, how and where associated plasma heating occurs, and how combinations of plasma flows and a variety of temperatures in plasmas along a line of sight through the optically thin corona lead to observables. Such work shows how apparently non-thermal signatures in spectra can emerge from line-of-sight integration through thermal plasmas. If you would like to learn more about how observables based on numerical work help guide the interpretation of real-world observables, a paper (with illuminating graphics) by Cheung *et al.* (2019) provides a good example.

Chapter 7

89-p. 183: Look up what T Tauri stars are, and what differentiates the 'classical' T Tauri star from the 'weak-line' variant.

90-p. 185: Verify the numbers in the conclusions about stellar magnetic braking for the present-day Sun at the end of Sect. 7.2.1.

91-p. 186: For comparison: what is the approximate ratio of forces exerted on the Earth of the total solar irradiance onto the Earth's surface (ignoring albedo, and assuming isotropic radiation from the atmosphere) to the solar-wind pressure on the magnetopause? That ratio shows why solar sails are designed for photon pressure rather than solar-wind dynamic pressure (note that some are designed to couple to induced electromagnetic effects, not dynamic pressure).

92-p. 192: Looking only at gravitational forces, how close to a solar-mass object would the Earth need to be to be pulled apart by tides? Whereas this is impossible with the Sun, an Earth-sized planet could be pulled apart if it approached a white dwarf or neutron star (and something like that is involved in 'contaminating' some white-dwarf atmospheres with heavy elements). An object of lesser density can be pulled apart, however, during a sufficiently close approach to the Sun: estimate at what distance (ignoring tensile strengths, spin, and orbital forces) comet 67P (with a mass of about 10^{16} g and characteristic dimension of 3 km) would have to come to the Sun to be broken up. Some Sun-grazing comets (such as the Kreutz family) have been observed to go through this breakup process.

93-p. 193: Consider what it means for solar eclipses that the Moon is moving away from the Earth: at some future time, the Moon will be so far away that no more total solar eclipses can occur anywhere on Earth. Assuming the Moon continues to move away at 4 cm/yr, roughly when will the last total solar eclipse occur? Confirm that the answer is somewhat more than 600 million years.

94-p. 193: If you are interested in solar eclipses, and wonder why the saros cycle has a slightly different length from the lunar nodal period, have a look here.

95-p. 194: One of the ways in which exoplanets are detected is to look spectroscopi-

cally at the displacements of the star about the barycenter of the exoplanetary system. How large is the velocity amplitude, and how large the associated Doppler shift at visible wavelengths, for the Sun-Jupiter system?

96-p. 194: What is the upper limit to the Sun's rotation rate in this phase? Formulate your arguments. You may ignore solar mass loss in this estimate. Use Fig. 10.5. This upper limit shows that the Sun's outer layers are rotating (much?) more slowly than the Earth is orbiting it, so that the tidal bulge on the Sun will be traveling through, and dissipating energy within, the solar outer envelope.

97-p. 196: Between the phases of tidally-locked binaries and merged binaries are (semi-)contact binaries in which mass transfer can occur as one of the binary components becomes larger than its 'Roche lobe', either because the star swells up in late evolutionary phases (see Ch. 10) or because the orbit shrinks by 'magnetic braking'. Now or after reading Ch. 10, look up the definition of 'Roche lobe' and the properties of RS CVn, Algol, W UMa, and FK Com objects as characteristic phases in the evolution of close binary stars towards single stars (with the details, and the class names, dependent on the masses of the two components).

98-p. 196: The Earth's equatorial bulge is nowadays used to keep satellites in a Sun-synchronous orbit, which is useful for satellites that need to scan the entire surface of the Earth, and also to enable Earth-orbiting satellites to have an uninterrupted view of the Sun throughout the year. Look up how this works.

Chapter 8

99-p. 199: The so-called 'anomalous cosmic rays' have a complex history: originally neutral particles in the interstellar medium, ionized by charge-exchange or photoionization in the solar wind, advected to the heliospheric extremes there to be accelerated. Important though they are as diagnostics of the outer heliosphere and the enveloping sheath-shock structure, they are not discussed in this volume. You can look them up for an interesting read ... after finishing this chapter. See Fig. 8.5 for where they appear in the energy spectrum.

100-p. 203: Use Eq. (8.5) to formulate (in a general vector expression) the magnitude of the $\mathbf{E} \times \mathbf{B}$ drift (in case you need a hint: assume the velocity can be described by an oscillatory component plus a constant drift).

101-p. 203: Rewrite Eqs. (8.7) and (8.8) to show that the drift velocity scales as the product of the particle's velocity and the gyroradius relative to the typical length scale in the gradient of the field, *i.e.*, as $v(r_g/\ell_t)$.

102-p. 203: Why do you think that bounce and drift motions are commonly ignored for the solar corona but are of dominant importance in the terrestrial magnetosphere? Hint: look at Table 3.4.

103-p. 204: Use a vector identity to show that the final term in Eq. (8.9) transforms into the central expression in Eq. (8.8) for a potential field.

104-p. 204: Estimate the orbital period associated with the drift velocity as in Eq. (8.7) for a purely equatorial motion for a proton with kinetic energy of 0.3 MeV orbiting, respectively, at 2 and 10 planetary radii r_p for, for example, Mercury, Earth (where the ring current is contained roughly within these distances), and Jupiter. Use the equatorial field strengths B_e as in Table 5.3 and $B(r) = B_e(r_p/r)^3$ for the equatorial dipole field. Is the non-relativistic approximation warranted for this proton? And for an electron of the same energy? Compare the relative size of the terrestrial ring current with the Chapman-Ferraro distance. How does this comparison work out for Mercury and what does that imply?

105-p. 208: Verify that without sources and losses, Eq. (8.18) – also known as Vlasov's equation – is a reformulation of Liouville's theorem, *i.e.*, $df/dt = 0$.

106-p. 208: **Advanced:** If you are interested in the origin of the terms in Eq. (8.20)

- you could review classic papers with fairly 'intuitive' introductions to the equation by, *e.g.*, Harm Moraal (1976) or Luke Drury (1983), the latter also including how term (d) arises. Or you could look at the paper by John Quenby (1984) which also describes the so-called 'force-field solution' that you will find in Sec. 14.1.1 on cosmogenic radionuclides.
- 107-p. 215: For isotropic diffusion from a point source into 3-d space, the equivalent to the 1-d version of Eq. (8.27) is $f(r, t) = (N_0/(4\pi\kappa t)^{3/2}) \exp[-r^2/(4\kappa t)]$. Assuming that the particles of different energies 'scatter' off the same irregularities and that the diffusion coefficient is independent of position, use this approximation to estimate the release time at the Sun for the first event in Fig. 8.4, as well as the equivalent mean free path λ_{mfp} for a diffusion coefficient of $\kappa = \lambda_{\text{mfp}}^2/2\tau_s$ for a typical time between scatterings τ_s for a population of protons. Hint: remember Fig. 8.3.
- 108-p. 223: Review the Rankine-Hugoniot jump conditions (Eqs. 5.2 and 5.8–5.12) and show that the motional electric field $\mathbf{E}_p = -\mathbf{v}_1 \times \mathbf{B}_1/c$ is constant across a steady-state, one-dimensional shock.
- 109-p. 228: Interpret the flux profiles as a function of time shown in Fig. 8.10 for the three different perspectives (with Earth in the direction of the arrows for three different events). Argue for the differences in timing of the solar event (first vertical line in each panel) and the passage of the shock (second line) relative to the timing of the peak fluxes.
- 110-p. 228: Energetic particles gyrate around the field lines in the solar wind. Roughly from what longitude region should we expect energetic particles to reach Earth that were created in flares or in shocks close to the Sun? From what longitudinal region at the Sun should we expect energetic particles in 'gradual events' to originate. What is roughly the delay between a flare/CME and 'prompt' particle storms? Explain the wide range of delays that can occur for flare/CME and 'gradual' particle storms?
- 111-p. 231: Look up the properties of (whistler) chorus waves, (ELF/VLF) hiss, and EMIC waves.

Chapter 9

- 112-p. 236: What drives tropospheric convection? Why is there no significant convection in the stratosphere (consider the role of ozone)? Is there an equivalent of a stratosphere on Venus? On Mars? Is there lower atmospheric convection? Formulate your arguments. The Web can help. The answers are 'yes', but not with a role for ozone except on Earth.
- 113-p. 236: The scale of the granulation in the photosphere of the Sun (and analogously of other cool stars) follows from a comparison of energy loss by radiation (effective once the plasma can radiate into space, with a time scale of order 20s) and supply by upflows. Work through this estimate: just below the photosphere, the largest contribution to energy being carried upward resides in latent heat of recombination of ionized hydrogen (with an ionization fraction of order 0.1); balance that with photospheric black-body radiation; use this to derive a minimum upward flow v_z needed to balance radiative losses. Then match timescales, and use that $v_h \leq c_s = (kT/m_p)^{1/2}$: for overturning convective flows, the horizontal time scale of ℓ_h/v_h should equal the vertical one ℓ_z/v_z , for a typical horizontal granular scale ℓ_h and overturning depth $\ell_z \approx H_p \approx 400$ km somewhat below the photosphere. (This argument is developed in H-III:5.2.1)
- 114-p. 237: Sound waves in an isothermal, hydrostatically-stratified atmosphere are 'evanescent', *i.e.*, non-propagating, at frequencies below the acoustic cutoff frequency $\omega_a = c_s/2H_p$. Can you argue intuitively why (think of the need for a

restoring force roughly within a wavelength)? Estimate the value of ω_a for the solar atmosphere. Why are p-modes only observed at frequencies below about ω_a ? Now derive an approximate dispersion relation in simplified geometry: Start from Eqs. (3.4), (3.5) and (2.3) for a hydrostatic 1-d plasma (mind the sign of g) and combine them retaining terms to first order for perturbations $\rho = \rho_0 + \rho_a$ and $p = p_0 + p_a$, where ρ_0 and p_0 describe the background stratified atmosphere at rest. Then factor out the exponential growth of the amplitude with height by substituting $v = u \exp(z/2H_p)$ and use $u \propto \exp(i[kz - \omega t])$ to obtain a dispersion relation that has propagating waves (real values of k) only for frequencies above ω_a (a somewhat different approach can be found in H-I:8.3). Lower frequency waves reflect and can form standing p-modes if they meet the criteria for global resonance, while higher frequency waves can propagate, but will steepen (readily into shock waves) as they propagate into the lower-density atmosphere.

- 115-p. 238: The 'G band' is a narrow spectral interval centered on electronic transitions of the CH molecule, mixed in with spectral lines from multiple metals. For the interested: look up the 'Fraunhofer lines' and their designations. This old nomenclature from the days in which the solar spectrum was first studied is still used for some of these 'lines', most frequently for the Ca II H and K lines, the Na D lines, and the G band.
- 116-p. 240: Estimate the gas pressure contrast between inside and outside for narrow 1 kG tube in thermal equilibrium at the solar effective temperature.
- 117-p. 240: Stars have a range of surface gravities, typically increasing monotonically along the main sequence towards lower effective temperatures, and substantially lower in evolved ('giant' and 'supergiant') stars than in main-sequence stars. Qualitative insight is provided by the following exercise: using the concepts of optical depth (and the fact that the stellar photosphere is around unit optical depth for continuum emission) and hydrostatic equilibrium (Ch. 2), show that the photospheric pressure would scale proportional with gravity in the idealized case of an isothermal atmosphere. In reality, radiative transport and convective motions modify that scaling for a real non-isothermal atmosphere, but the trend is in the correct direction. With this insight, argue for the trend of intrinsic field strength of photospheric magnetic concentrations with gravity: from ~ 1.4 kG in mid F-type dwarf stars to ~ 3.2 kG in late K-type dwarf stars, and well below 1 kG for cool giants.
- 118-p. 240: The transition from bright to dark magnetic structures occurs at a scale of roughly 200 – 300 km. What does that say about the typical photon-mean free path ℓ_{ph} in the photosphere? Compare that value to the corresponding pressure scale height, and argue why $\ell_{ph} \gtrsim H_p$ just at the photosphere.
- 119-p. 240: Explain why observed field strengths inside flux tubes exceed the equipartition field strength (field strength in an imaginary completely evacuated tube) at the level of the external photosphere.
- 120-p. 240: When the total solar irradiance (TSI) is smoothed over time scales of, say, a week, the Sun is brighter at sunspot maximum than at sunspot minimum, but when looking at TSI curves with a resolution of a day or so, the presence of large sunspots leads to dips when these are near the central meridian. Explain this qualitatively by the mix of faculae, pores, and spots in and around active regions. Look up TSI curves in different phases of the solar cycle.
- 121-p. 243: The processes of electromagnetic radiation from a plasma involve three fundamentally distinct processes: bound-bound, free-bound (radiative recombination), and free-free (Bremsstrahlung) emission. The Sun's coronal emission, caused by collisions of ions with thermal electrons, is dominated the first, except

- for flares when the last is also important; why? Which ions are typically strong contributors to the coronal X-ray and EUV emission from an active region at ~ 3 MK? Hint: combine elemental abundances with ionization energies (such as given here). For this rough estimate, ignore oscillator strengths for the transitions involved. For the solar corona under most conditions, the dominant radiative losses are from C, N, O (below about 0.5 MK), and Fe (above about 0.5 MK).
- 122-p. 243: Eq. (9.1) contains a product of electron and hydrogen densities, but hydrogen is fully ionized at coronal temperatures and thus has no spectral lines that can be excited through collisions with electrons. Why is it acceptable to express it this way?
- 123-p. 245: Use Eq. (9.4) to estimate typical volumetric heating rates for a coronal region over 'quiet Sun' (*i.e.*, outside of active regions; with coronal field strengths of order 20 G, loop-top temperatures of ≈ 1 MK, and loop half lengths $L \sim 4 \cdot 10^9$ cm) and for an active (sunspot-bearing) region (with coronal field strengths of order 200 G, loop-top temperatures of ≈ 3 MK, and loop half lengths $L \sim 15 \cdot 10^9$ cm). Compare these to the thermal energies also estimated from Eq. (9.4) and also compare plasma to field pressures (*i.e.*, compute values of β).

Chapter 10

- 124-p. 248: What fraction of the Sun's hydrogen would need to be converted to helium to keep it at (roughly) its current brightness throughout the time it spends on the main sequence? Once core hydrogen is consumed, the stellar internal structure changes considerably, enough to ignite fusion in higher layers as the star moves into its giant phase. Use $E = mc^2$.
- 125-p. 251: A cautionary intermezzo: Sect. 9.3 gives power-law scalings between radiative losses from chromospheres and coronae over stellar surface areas with mean magnetic flux densities over these areas (which hold approximately without changes for areas up to entire hemispheres). The values of the power-law indices in these relationships depend on the formation temperatures of the diagnostics used (thus, for example, steepening towards higher-energy X-ray channels), while published values also depend on the correction for a reference level (there is a minimum or 'basal' level of chromospheric emission that needs to be subtracted first but different authors use different corrections). This dependence on the details of the diagnostics used are one cause behind the somewhat different power-law scaling between coronal and chromospheric radiative losses you find in the literature. There are other reasons why you may find other approximate parameterizations. For one thing, although the scaling in rotation-activity diagrams between a relative brightness in terms of luminosity or surface flux density ($L_i/L_{\text{bol}} \equiv F_i/F_{\text{bol}}$) versus Rossby number works fairly well, it does not work perfectly, and other authors, using other stellar samples, might prefer using F_i versus P_{rot} . As long as the stellar sample contains stars of rather comparable internal properties, the choice of metric does not matter, but for more diverse samples, scalings with these properties matter – no simple multiplicative scaling seems to lead to a single tight rotation-activity relationship for all cool stars. Other reasons for differing results from different studies include the fact that the relationships are not simple power laws and fits thus depend on the parameter range covered in stellar samples, and, of course, uncertainties in models for, *e.g.*, stellar ages, and intrinsic stellar variability combined with relatively small samples. You could review, for example, this study by Booth *et al.* (2017) for more discussion and for references.
- 126-p. 255: Note that integration over the power in flares as parameterized in Eq. (10.2) diverges when the lower and upper limits extend to $[0, \infty]$. Consider what processes could be at play in introducing cutoffs to the integral on either side.

The answer remains under study: it is not clear over what range Eq. (10.2) holds its slope, or what determines the energy of the 'largest flare', or how and how much relatively tiny 'nano-flares' contribute to coronal heating. But considering the possibilities should prove educational.

- 127-p. 257: Use Fig. 10.5 to estimate the mass of the least-massive post-main-sequence star (say, past the phase marked by a turn towards cooler surfaces marked by the squares) that could exist in the present-day Universe.
- 128-p. 257: Fig. 10.5 can be used to illustrate how astronomers determine the ages of 'open clusters' of stars (other than by the modern means of asteroseismology): assuming that all stars in a cluster are formed at about the same time, the shape of the HR diagram for stars in a cluster reveals the age when compared to theoretical evolutionary tracks as in the panel on the upper left. Try this: assume the stars are all 955 Myr old as in the open cluster called NGC 2355, then mark the approximate positions of the stars at that age in the upper-left panel of Fig. 10.5 estimating also where stars of intermediate masses might show up, and realize how the turnoff from the main sequence in such a cluster HR diagram reveals the age of the cluster. Open clusters all have the low-mass end of this HR diagram in common, so even if the distance to a cluster is not known, the distance can be determined by shifting that low-mass tail to overlap with that of a cluster of known distance. Also: look up the definition of 'open cluster' in contrast to a 'globular cluster'.
- 129-p. 257: Start with Eq. (7.1) and note that the mass loss \dot{M} can be expressed in terms of the Alfvén speed v_A and the radial field strength B_A at the Alfvén radius. Then make the approximation that for a thermally-driven wind (in which centrifugal forces can be neglected) the Alfvén speed $v_A \approx c_s$, and take cool-star winds to all have comparable temperatures. Show that Eqs. (10.3) and (7.1) together imply that $B_A \propto \Omega$. The combination of the latter with a relationship between the photospheric field and rotation rate implicitly constrains stellar field geometries as it connects B_A and B_{surf} .
- 130-p. 258: Estimate the coronal soft X-ray brightness for a Sun-like star in its 'teenage years' (Sect. 10.2.5) relative to that of the present-day Sun.
- 131-p. 260: The minimum flare energy given here is instrumental, not intrinsic. Argue why the empirical lower limit of flares detectable by an instrument like *Kepler* is limited to of order 10^{33} ergs. Note that this lower limit exceeds the energies observed (to date, at least) for solar flares.
- 132-p. 261: Just to get an impression of relative velocities: compare the average speed of the Solar System relative to the local ISM to the speed of 828,000 km/h with which the Solar System orbits the Galactic center.
- 133-p. 263: With average values for solar wind density and velocity (assuming a radial outflow at constant velocity and with a density as specified in Table 2.4), at what distance from the Sun does the solar wind dynamic pressure equal the interstellar total pressure for estimated values of $B_{\text{LISM}} \approx 3 \mu\text{G}$, $T_{\text{LISM}} \approx 6500 \text{ K}$, and $n_{\text{p,LISM}} \approx 0.06 \text{ cm}^{-3}$ and $n_{\text{H,LISM}} \approx 0.18 \text{ cm}^{-3}$ (see, *e.g.*, Sect. H-IV:3.2)?
- 134-p. 263: Given present-day parameters for the ISM as in Activity 133, where would the heliopause be, very approximately, for the range of ISM densities given in Sect. 10.3.1, assuming a present-day spherically-symmetric, constant-velocity solar wind, and the same temperature for the ISM? Compare your result with Fig. 10.7. Look back to Sect. 5.5.8 for some of the physics involved.
- 135-p. 265: For charge exchange only, and assuming (very approximately, as done initially (Holzer, 1972) decades ago) a velocity-independent cross section for resonant-charge exchange of solar wind protons with ISM neutral hydrogen of

$\sigma_{\text{CX}} \approx 2 \cdot 10^{-15} \text{ cm}^2$, what fraction of H^0 , looking at the population after passing through the 'hydrogen wall' and moving in a straight line towards the Sun, would reach Earth orbit for present-day slow wind conditions? In reality, other processes are major players: radiation pressure (for neutral hydrogen primarily by repeated Lyman α absorption followed by isotropic re-emission) pushes outward on the atoms, and photo-ionization in the Sun's EUV and X-rays presents a significant loss term. It appears that Lyman α radiation pressure on ISM H^0 just balances solar gravity, see Schwadron *et al.* (2013); for a significantly younger Sun, ISM H^0 would never reach Earth orbit. The combined effects of these processes would render an IBEX-like mission to learn about the ISM H^0 around a young Sun pointless except during times of passage through dense interstellar clouds.

136-p. 265: Argue why the heliospheric hydrogen wall has a thickness L_{HW} that is, within a factor of a few, comparable to, but less than, the distance d_{HP} from the Sun to the heliopause. For a simple estimate, use a circular 'cookie tin' geometry to approximate conditions at the heliosphere's 'nose', with the incoming flow through the top being decelerated by the gas pressure (ignore magnetic effects here), and accelerated sideways out by the same pressure, combining scale estimates based on the continuity and momentum equations; focus only on the flow into the heliopause and assume no bow shock (see Sect. 5.5.8). Check that this gives it just enough of a total column depth with the charge-exchange cross section from Activity 135 so that a useful fraction of ISM neutral hydrogen can indeed be made part of the flow in heliosheath. See this study by Wood *et al.* (2002) for simulated astrospheres and their hydrogen walls, and some images for different stars and their speeds through the ISM. How would the thickness of the hydrogen wall change for a much higher speed of a planetary system through its LISM?

137-p. 266: Show that the total power lost in X-rays from the present-day solar corona (estimated from Fig. 10.3 or 10.9) is roughly twice the total power lost in the solar wind (using the expressions in Sec. 3.5.2), and that these numbers would have been comparable for the young Sun at the 'wind dividing line' if the characteristic wind speed would have been the same.

138-p. 268: Place the coronal activity level corresponding to the 'wind dividing line' in the rotation-age diagram in Fig. 10.3, and consider possible consequences for that diagram.

139-p. 269: Estimate the size of the heliosphere and the terrestrial magnetopause distance for a young Sun at an age of 700 Myr, assuming unchanged LISM conditions and geomagnetic properties.

140-p. 271: Section 10.3.3 mentions a solar rotation period of 25 d while the caption to Fig. 5.6 mentions 27 d. What is the reason for using these two different values in the different contexts?

Chapter 11

141-p. 273: How many Earth masses of elements heavier than carbon are contained in a solar mass cloud of solar composition? Most of that material in the original cloud ended up inside the Sun, of course. What fraction, roughly, of the original cloud would need to remain in the disk to ultimately form the planets? Why are the answers to these two questions largely independent of each other (think about what mostly makes Jupiter and Saturn).

142-p. 273: Compare a size of $R_c \gtrsim 2 \times 10^4$ astronomical units to distances between stars in star-forming regions. Express that distance in light years and in parsecs, and compare those to the distances to the nearest stars for the present-day Sun.

143-p. 273: Another way of formulating Eq. (2.15) is to say that the mass of the

cloud must exceed a certain value. Reformulate Eq. (2.15) as function of cloud temperature T_c , cloud density n_c , and stellar mass M_* (Note: this is similar to what is known as the Jeans Mass, which is commonly derived from energy imbalance or by a comparison of sound and free-fall time scales in a perturbation analysis). This shows that $M_* \sim f M_\odot T_c^{3/2} / n_c^{1/2}$. Derive the value of the constant $f \approx 2$ assuming, for simplicity, that the gas consists predominantly of molecular hydrogen. For n_c of order 100 cm^{-3} estimate M_* for $T_c \approx 10\text{K}$, characteristic of present-day molecular clouds (realizing this is a rough order-of-magnitude estimate). Early in the life of the Universe, with only H and He in the mix, the interstellar gas lacked many of the strong emission lines of heavier elements, could therefore not cool as efficiently, leaving interstellar clouds significantly warmer, roughly of order 100 K. Use the derived expression to show that this favors the formation of much heavier stars, even when starting from a higher density of order 10^4 cm^{-3} . This review by Johnson (2019) discusses how this contributed to the evolution of elements heavier than H and He (known as 'metals' to astronomers) over the history of the Universe.

144-p. 274: Estimate the orbital Doppler swings and the fractional dimming during transits observed from afar of Mercury, Earth, and Jupiter around the Sun. Also estimate how close a Jupiter-like exoplanet (with an albedo of 0.5) should orbit for the fractional bolometric dimming during a secondary eclipse (when the planet moves behind the star) to be about 1 millimagnitude (which is the noise level for the telescope of the *Kepler* spacecraft for a 13th magnitude star at 1-minute exposure times; consider at what wavelength range the contrast is optimal). Use, *e.g.*, this fact sheet. Compare the Doppler signals with the thermal widths of spectral lines, and consider what to use as reference wavelengths. How large is the Doppler swing added to the stellar signals owing to Earth's orbit around the Sun?

145-p. 274: Look up and summarize the principles of the five detection methods of exoplanets, and consider what the strengths, weaknesses, and technological challenges are for each method. Note: activities 95 and 144 ask about Doppler signals and transit photometry.

146-p. 275: Star-forming regions and disks around young stars are best observed in the near-infrared region of the spectrum. Look into what wavelengths are often used for such observations, and consider why ('Why is the sky blue?'), given that dust sizes in the interstellar medium peak around a few tenths of a micron.

147-p. 278: Figure 11.2 shows a curved 'snow line' (or 'ice line'). What is the reason behind that?

148-p. 278: Look up the 'Grand Tack' model and review the likely consequences for the growing Mars, for the asteroid belt, and for water distribution by scattered asteroids into the inner solar system.

149-p. 279: Figure 11.2 shows a clearing near the central star. This is associated with the magnetic field of the rotating star. Consider what processes are at play there and the role of the following: accretion rate, ionization fraction, diffusion of field into the ionized gaseous disk, orbital and angular velocities, the corotation radius, winding up of magnetic field that connects the star to the disk, centrifugal force, etc. There is no easy concept for this: you can look at the literature of MHD models of T Tauri accretion disks to see how complex the coupling is. Store your thoughts: the star-disk interaction leading to the clearing is discussed in Sect. 11.2.2.

150-p. 280: Sketch and describe the observable spectral signatures of transiting planets for orbits of different obliquity (including effectively retrograde planets).

- Also: estimate transit times for planets around of solar-mass star at distances such as Mercury, Earth, Jupiter, and Neptune. Use, *e.g.*, this fact sheet.
- 151-p. 285: Estimate the total values and ratios of mass and angular momentum in the planetary system and in the Sun (use Fig. 10.5).
- 152-p. 285: Iron, oxygen, and silicon make up three quarters of the Earth's mass. Iron is some 30% of the total. In the interstellar medium, iron makes up about 1 part in 1,000 of total mass. How many Earth-equivalents of iron does a circumstellar disk with a mass of 1% of the Sun contain?
- 153-p. 286: Think about similarities and differences with Solar-System magnetic instabilities as discussed in Ch. 6 when reading about things like FU-Orionis outbursts and 'ballooning out' of magnetic field in ejections of mass from corona and disk, likely driven by necessarily failing attempts of the forces at play to impose corotation.
- 154-p. 288: The internal energy of the star in Eq. (11.1) is derived from the so-called 'virial theorem' which states that the total gravitational energy E_{grav} is related to the total thermal energy E_{thermal} as $E_{\text{grav}} = -2E_{\text{thermal}}$ if $\gamma = 5/3$ as for a monoatomic ideal gas. Derive this from Eq. (3.5) assuming a field-free stationary state for a spherically symmetric ball of gas: $dp/dr = -GM(r)\rho/r^2$. One way to do so is to multiply both sides by $4\pi r^3$, integrate (in part 'by parts') from center to surface (where $p(R)$ essentially vanishes, and realizing that the internal energy per unit volume of the gas is given by $u = p/(\gamma - 1)$ for an adiabatic exponent γ). The result is equivalent to the virial theorem. Eq. (11.1) can be used for the present-day Sun to show that continued gravitational contraction cannot support the solar energy budget over the age estimated for the Earth based on radio-nuclide dating (note a factor of two difference between thermal and gravitational time scales). What is the present-day value of τ_{KH} in Eq. (11.1) for the Sun?
- 155-p. 289: Draw lines of equal radius (as multiples of the solar value) in Fig. 11.6, using $\log(T_{\text{eff},\odot}) = 3.762$.
- 156-p. 291: Derive the expression for the breakup rotation rate of stars as function of mass and radius. What is the value for the Sun? Ignore distortion from spherical symmetry for this estimate.
- 157-p. 297: The key mechanism by which dust is expected to settle into the center of an accretion disk is hydrodynamic drag. Explain how this works. Consider orbital inclination and effects of gas pressure, gravity, and stratification.
- 158-p. 299: For further study/reading: Most stars are born in groups of substantial numbers (often in what are called 'open clusters'). In such clusters, stars of a range of masses are formed (statistically yielding the 'initial mass function'). The heaviest among these evolve fastest, and if heavy enough can end their lives in a 'supernova'. The open cluster is eventually pulled apart by the 'galactic tides', which limits the exposure of planetary systems to nearby supernovae and to gravitational perturbation of the orbits of the planets. Look up the terms between quotation marks. The occurrence of a nearby supernova appears consistent with several properties of the solar system, including one of several possible means for the early melting of small bodies (as reflected in what are known as 'chondrules'). Look at this study by Portegies Zwart *et al.* (2018) for more on this.

Chapter 12

- 159-p. 301: For an impression of order-of-magnitude numbers, estimate the energy involved in a collision between an Earth-mass body and an Mars-mass body at an impact velocity of, say, 14 km/s. Ignoring the energy going into the formation of the Moon in such a process, but rather assuming all mass and energy remain

- within the newly formed body, estimate the average temperature increase if all energy were distributed throughout half of the volume of the mantle, and that that material has a specific heat of approximately 1.5×10^7 erg/g/K.
- 160-p. 301: Make an order of magnitude estimate of the cooling time of Earth's atmosphere after impact of a Mars-mass body: assume an impact velocity of 14 km/s, that all kinetic energy remains within the near-surface layers and atmosphere; an optically thick atmosphere of vaporized silicate; and a characteristic temperature of the radiating vapor of, say, 2000 K.
- 161-p. 302: Although the definition of 'habitability' commonly involves the requirement of liquid surface water, some definitions are more relaxed. Perhaps other surface liquids can serve as agents in support of life (such as ethane and methane lakes and seas on that cover 1.6 million square kilometers, or 2% of the surface, of Saturn's moon Titan) or perhaps subsurface water (as encapsulated seas or even globe-spanning layers) can support life. With that in mind, explore the moons of the giant planets that are thought to meet at least the condition of large reservoirs of some liquid somewhere, in particular: Europa, Callisto, Ganymede, and Io at Jupiter, Enceladus and Titan at Saturn, and Triton at Neptune. Which three power sources are thought to be most important in maintaining liquid states on giant-planet moons?
- 162-p. 303: **For the curious:** Photosynthesis depends on the chemicals involved and as such is sensitive to the spectral energy distribution of the star. You could search the literature on developments in this area, but for stars substantially different from our Sun that work remains hypothetical. Here is a possible entry point. Look up where the main absorption bands of chlorophyll and β -carotene lie relative to the solar spectrum at sea level. How does the solar spectrum change under water for, for example, flora in the oceans?
- 163-p. 305: Sect. 10.3.1 describes the possibility of the Solar System moving through dense, cold interstellar clouds, which could greatly enhance the dust environment of Earth. Review the study by Pavlov *et al.* (2005) for the potential effects on terrestrial climate, including periods of strong glaciation and potentially the triggering of a 'Snowball Earth' state.
- 164-p. 305: Consider the evolving CO₂ content of the atmosphere of a lifeless terrestrial planet. Which of the following parameters would influence the atmospheric CO₂ content over time: (1) atmospheric mass and composition, (2) chemical composition of seas and oceans, (3) continent sizes and placement, (4) fractional coverage by liquids in seas and oceans, (5) motion through, and density of, local interstellar medium, (6) orbital obliquity, (7) orbital period (length of the planetary 'year'), (8) planetary mass, (9) planetary radius, (10) planetary spin obliquity, (11) planetary spin rate (length of the planetary 'day'), (12) planets elsewhere in the planetary system, (13) plate tectonics, (14) properties of moons, (15) spectral type of the central star, (16) stellar spin rate. Formulate your arguments for each. You may want to read on in Ch. 12 and return here later to complete the activity.
- 165-p. 306: Look up what constitutes the geocorona.
- 166-p. 307: What role does plate tectonics likely play in dynamos in terrestrial planets? Reminder: Sect. 4.1.1.
- 167-p. 307: The Earth's argon is predominantly Argon-40, whereas that in the universe at large, as in the Sun, is Argon-36. What is the source of Argon-40 in Earth's atmosphere?
- 168-p. 312: At what distance would an Earth-equivalent exoplanet need to orbit an $0.6 M_{\odot}$ M0 V star to reach the same global 'equilibrium temperature', all other

- things being equal? You may disregard effects associated with the difference in the stellar spectral energy distribution on the exoplanet, but you should not ignore the bolometric correction in estimating the total stellar irradiance. How long would a year last on such a planet compared to Earth's? Use Fig. 4.2. Note: such close-in planets are subject to very strong tidal forces that will synchronize spin and orbital periods, causing these exoplanets to lose their day-night cycles. That, in turn, invalidates your estimate – why?
- 169-p. 312: Beyond the furthest planet: The New Horizons spacecraft flew by Kuiper Belt Object 2014 MU₆₉ on 2019/01/01, the most distant body visited by a spacecraft to date, at an orbital distance of ~ 44 AU. Estimate the surface temperature of 2014 MU₆₉, which has an albedo of ~ 0.1 . Compare your estimate to the observed temperature in this paper by Stern *et al.* (2019).
- 170-p. 315: Show that the simple model in Fig. 12.5 yields an estimate consistent with Earth's global temperature rise of about one degree (observed between 1850 and 2010) based on the increase in anthropogenic radiative forcing as shown in Figure 12.6 within the uncertainty indicated in that figure.
- 171-p. 317: Compare the values of \mathcal{P}_{abs} from Eq. (12.1) for Venus and Earth. Explain qualitatively why Venus' surface temperature exceeds Earth's, then read on for the answer.
- 172-p. 321: To get an idea of scales: estimate the size of a comet that would double the CO₂ content of Earth's atmosphere. How does that compare to, e.g. comet 1P, the target of the *Giotto* mission, and 67C, the target of the *Rosetta* mission?
- 173-p. 324: What is the basis of the Kelvin-Helmholtz instability? This instability also occurs between the terrestrial magnetosphere and magnetopause flow because the magnetic tension is not strong enough to stabilize the developing waves. Why is this geospace phenomenon not listed as a process for 'bulk outflow'?
- 174-p. 326: Make a table summarizing which atmospheric loss processes work on each of the terrestrial planets. Which two processes are most effective for the present-day Earth based on the description in Sect. 12.4.1?
- 175-p. 328: Human impacts on climate appear not to be limited to the Industrial Revolution! Have a look at a study by Koch *et al.* (2019): they argue that the large population reduction in the Americas following the arrival of European conquerors and settlers, and the resulting reforestation of abandoned agricultural lands, was a significant part of the change in atmospheric CO₂ in the late 16th Century and in the 17th Century.
- 176-p. 328: Compile a list of all the processes involved in setting a planetary climate system that reflects at least all those mentioned in Chs. 11 and 12. You can assimilate relevant processes from Activity 164 here as a start.

Chapter 13

- 177-p. 330: The approximate scaling of the molecular diffusion coefficient D with molecular mass m and particle density n follows from energy equilibrium of the constituent particles. Formulate D as function of the collisional cross section σ and of temperature and density in the case of self-diffusion, *i.e.*, for molecules diffusing among themselves. For a mixture of components, mutual diffusion needs to be considered.
- 178-p. 331: Work through the units of Eqs. (13.3) and (13.4) to show that η_X is an efficiency per unit energy per unit wavelength.
- 179-p. 336: Consider the similarities and differences between the charge-exchange reactions described here and two- and three-body gravitational interactions, specifically what is needed for the capture of interplanetary spacecraft into closed orbits, or the capture of planetary bodies as moons of planets. For the latter,

- look up the concepts proposed for the capture of Triton, the largest moon of Neptune, orbiting that planet in a retrograde orbit (which implies it has to involve a capture well after the formation of the planet).
- 180-p. 337: Note the equivalence between Eq. (13.13) and Eq. (11.5) for a volumetric ionization rate of $\Pi(e^-) \propto \Phi_i/(4\pi R^3)$. This means that α_{eff} is, in effect, for a 'case B' recombination, *i.e.*, excluding the possibility that emitted photons in recombination are absorbed to lead to another ionization event. Consider what could happen to avoid that. Also see a parallel with the formulation of what can be viewed as the inverse in Eq. (9.3): for a stationary, isothermal case, the 'incoming' volumetric heating ϵ_{heat} balances the outgoing radiation $n_e n_H f_{\text{rad}}$ in which the product of ion and electron densities is a measure for the number of collisions leading to excitation, to compare with the ionizing radiation in the ionosphere which balances the recombination in which the product of ion and electron densities is a measure for the number of collisions leading to recombination.
- 181-p. 337: One might think that collisions between particles that can 'bond' and thereby be taken out of a population under study, such as electrons and positively-charged ions that combine into a neutral particle, might have a good analogy in how flux concentrations in the solar photosphere behave: the concentrations perform a random walk and in collisions opposite magnetic polarities 'cancel', *i.e.*, disappear from the population of magnetic charges. Yet the scaling behavior between the strength of the source (the total of emerging bipoles per unit time) and sinks (the total of canceling flux per unit time) is different: the square root dependence reflected in Eq. (13.13) does not show up, but instead a near-linear dependence appears (as shown here by Schrijver (2001)). Consider the reasons: when the Sun's activity increases, flux concentrations grow larger by collision thereby countering the increase in collision frequency expected; larger concentrations are less mobile within the evolving convective motions; fragmentation and coagulation are seeking a balance; while in general the large-scale meridional flow aids in separating polarities (a process that is countered in an ionosphere by the tendency towards charge neutrality).
- 182-p. 337: Show for the simplified case of a fully-ionized static gas that the scale height for ions is twice that for the corresponding atoms in a neutral atmosphere by combining the momentum equations for ions and electrons in comparison to that equation for a neutral species. And remind yourself how this is consistently incorporated in the MHD equations.
- 183-p. 340: Review Figs. 2.5, 13.1, and 13.2 and think through the dominant reactions described in Sects. 13.1.2 and 13.1.3.
- 184-p. 342: Trace which part of the solar spectrum provides the predominant power to the *E* and *F* layers of the terrestrial ionosphere and overlapping thermospheric regions, and note that the power going into the *F* layer exhibits a larger variation over the solar cycle than that going into the *E* region. See Sect. 13.1.1 and Figs. 2.3, 2.4 and 2.7.
- 185-p. 343: Upward traveling radio waves with frequencies below the plasma frequency are evanescent within the ionosphere and are reflected downward, thus enabling 'over the horizon' or 'skywave' communication. Look up the plasma frequency (Eq. 3.41), typical ionospheric electron densities within the ionosphere, and resulting values for the radio frequencies useful for such communication. See, *e.g.*, Fig. 6.1 for an overview of the EM spectrum with an indication of various radio bands (including what differentiates propagation of AM and FM bands).
- 186-p. 349: Estimate when the solar EUV flux dropped to a level that the thermospheric climate became comparable to the present-day state; and summarize the

ionospheric changes over geological time scales as far as the models discussed here is concerned.

187-p. 352: Assuming a similar geomagnetic field, use the expressions in Eq. (13.14) to derive an estimate of the magnetopause distance over time. Show that for a young Sun this comes down to $\sim 1.25R_{\oplus}$ (with Eq. 5.22).

Chapter 14

188-p. 356: Advanced: If you are interested in how Eq. (8.20) can be approximated by something like Eq. (14.1) you can find the origin of this transformation in a study by Gleeson and Axford (1968).

189-p. 357: To appreciate how little radionuclide material there is to work with, compute the global annual production in kg for ^{14}C and ^{10}Be . That production rate puts roughly one ^{14}C atom per 10^{12} atoms of ^{12}C in living tissue through uptake of atmospheric CO_2 by plants and their subsequent consumption by animals.

190-p. 359: The ^{10}Be production rate for Mars would be about 2.5 times the terrestrial rate if it had a terrestrial atmosphere. Show why based on data in this text.

191-p. 360: Over the past century the concentration of ^{14}C in the biosphere has been dropping considerably because of fossil-fuel burning (why?). Express in functional form how this leads to an ambiguity in ^{14}C dating if no other information on the age of an object is known.

192-p. 361: Look up 'paleomagnetic dating' in relation to the 'remanence measurements' mentioned in Sect. 14.2.

193-p. 365: How are the decrease of stellar rotation speed, magnetic activity, and mass-loss rate on long time scales compatible with the 'essentially constant' GCR exposure over the past ≈ 1 Gyr? The answer has to do with the fact that the Sun is already an aged star, and can be traced to its relatively weak magnetic braking over the past 1 Gyr, and thus relatively little decrease in coronal activity and mass-loss rate. The limited impact on GCRs at Earth orbit over time also suggests that the heliospheric variability (leading to diffusive GCR scattering) has not changed too much. Estimate the changes over time using Eqs. (10.3) and (10.7), and Fig. 10.3.

194-p. 366: Heavy stars evolve much faster than low-mass stars, and can, if heavy enough, explode in a supernova even as lower-mass stars and their planetary system forming within the same molecular cloud are still in their formative phases. Look up lifetimes and evolutionary pathways for stars of different masses. Also look up properties of clusters of stars in star-forming regions.

195-p. 366: Estimate how much stronger the dynamic pressure of the incoming supernova wave front needs to be than the present-day IMF, assuming comparable solar-wind properties, to push the heliospheric boundary to within 1 AU. See Activity 133.

Chapter 15

196-p. 369: **Observing exoplanetary atmospheres:** (1) Approximate the contrast $C_{*,p}(\lambda)$ between exoplanetary atmospheric radiation and stellar surface radiation, assuming both star and planet radiate as black bodies (using Planck's law $B(\lambda, T)$; ignoring center-to-limb effects), as function of wavelength, of the temperatures of star (T_*) and planet (T_p), and of the effective radii of star ($R_*(\lambda)$) and planet ($R_p(\lambda)$). Show that $C_{*,p}(\lambda = 10 \mu\text{m}) \approx (20 * R_*(\lambda) / R_p(\lambda))^2$ for $T_* = T_{\odot}$ and $T_p = T_{\oplus}$. (2) What is $C_{\odot, \oplus}(\lambda = 10 \mu\text{m})$? This value shows how hard it is to separate stellar and planetary signals (it is easier for closer-in, warmer, planets and for larger planets, such as hot Jupiters). (3) Why does the

IR domain of the wavelength spectrum provide optimal access to the exoplanetary spectrum using a secondary eclipse (when the planet moves behind the star)? (4) At what wavelength does $B(\lambda, T)d\lambda$ peak for a planet at $T_p = T_\oplus$ (use Wien's displacement law)? (5) For transit spectroscopy, in contrast, optical wavelengths are most suitable for G- and K-type stars; why? See this tutorial by Deming *et al.* (2019) on exoplanet transit spectroscopy for answers and for much more on this topic.

197-p. 369: **Exoplanetary atmospheric spectroscopy:** How does wavelength-dependent transparency of an exoplanetary atmosphere lead to wavelength-dependent transit depths $\mathcal{T}(\lambda)$, and thereby yield spectral signatures of atmospheric chemicals? Basically, the apparent radius of exoplanet-plus-atmosphere depends on wavelength because the atmospheric opacity does. But the transit depth also depends on whether there are features on the stellar disk within the transit path. This provides information on, *e.g.*, starspot properties. Sketch how these two signals combine into the observed transit depth signal $\mathcal{T}(\lambda, t)$ over a transit. Consider how one might go about disentangling these two signals. See the reference in Activity 196 for more information.

198-p. 369: **Comparative heliophysics:** Look up the properties of the stars α CMa A, the Sun, and TRAPPIST-1. Consider how the following properties differ for a planet orbiting each of these stars within the continuously habitable zone (unconfirmed to exist in the case of α CMa A): (a) the size and color of the star, (b) the maximum possible age of the planet, (c) the duration of the orbital year and constraints on the length of the planetary day, (d) possible constraints on the planetary dynamo (subject to what we know about these at present), (e) the Alfvén Mach number of the stellar wind, (f) the magnetopause distance (assuming comparable planetary dynamos), (g) constraints on loss of planetary water, (h) the potential of measuring interstellar neutral hydrogen from an orbit near that planet, (i) the spectrum of the stellar and galactic cosmic rays (assuming the same spectrum external to the planetary system).

199-p. 369: **A study on energetic particles in TRAPPIST-1:** TRAPPIST-1 is a very different world from our Solar System. The central star - itself only first observed in 1999 - is merely 1/8th the size of the Sun, only slightly larger than Jupiter. Its brightness is almost 2,000 times less than that of the Sun. The star is orbited by seven known exoplanets (first published on in 2016), much like Earth in size and mass, but all very close to their star. At least three of these seven planets are estimated to orbit within the liquid-water habitable zone. You can start reading up on TRAPPIST-1 using an ADS search, but for this Activity review this study by Garraffo *et al.* (2017) on the astrosphere, and this study by Fraschetti *et al.* (2019) of the possibly very intense radiation environment of the planets. The role of heliophysics in this is evident throughout these studies: (1) identify the processes you have read about in this book that are elements of these studies. (2) Use what you learned in this text to explain why the wind is mostly sub-Alfvénic around the seven planets. (3) With dynamic pressures 3 to 6 orders of magnitude higher than for Earth, what does that do for the planetary magnetopause distances? (4) Although this system is diminutive, its astrosphere is potentially huge: estimate the distance to the astropause assuming the system is subject to ISM conditions similar to those for the Solar System.

200-p. 369: **Arriving at Earth's climate from scratch:** In Activity 176 you compiled a list of all the processes involved in setting a planetary climate system that reflected at least all those mentioned in Chs. 11 and 12. Now, complement that list with the additional topics discussed in Ch. 14. Do not forget to add

relevant thoughts from your notes for Activity 16! Then review that list and flag those processes that are beneficial to life as we know it on Earth and those that are detrimental to it. The duality of many, perhaps most, of the entries on your list should make you think about how our Earth, as it is in its present state, is a consequence of a remarkable interplay of often simultaneously beneficial and detrimental processes, including, perhaps, a series of fortuitous developments. Consider the extraordinary challenge of thinking about 'habitability' of any of the other thousands of exoplanets found to date, including what the phrase 'habitability' itself adds to that challenge given how little we know about life itself. Better yet, write an essay on this to share with fellow students, with teachers, and perhaps a much larger readership. After all, science is about communicating your thoughts and discoveries, as much to your peers as to society at large.

Version history

- 1.1 2019/10/24: Original version.
- 1.2 2019/10/30: Deleted cover figure and compressed large-volume figures to fit into arXiv's 15MB limit. Added explicit references and associated bibliography, corrected a few URLs to figure sources, and made the Astrophysics Data System (ADS) the primary path to the literature in URLs throughout. Added version history. Added explanation of Heliophysics volume numbers to Preface.

List of Illustrations

2.1	Temperature-density diagram for heliophysics.	14
2.2	Average vertical temperature profile through Earth's atmosphere.	16
2.3	Solar <i>vs.</i> 5770 K black body spectrum, and solar-cycle variability.	19
2.4	Altitude of penetration of the solar radiation as a function of wavelength.	19
2.5	Profiles of major species in upper atmospheres of Venus, Earth, Mars.	30
2.6	Temperature profiles for solar min. and max. for Venus, Earth, Mars.	30
2.7	Earth's ionosphere for day- and nighttime, at high and low solar activity.	31
2.8	Classical Chapman profile.	33
2.9	Interactions of plasma with neutrals.	36
2.10	Characteristic densities and ionization in ionosphere and chromosphere.	39
4.1	Radiative and convective internal structure of main-sequence stars.	74
4.2	Activity across the Hertzsprung-Russell diagram.	75
4.3	Polarity of the geomagnetic field for the past 120 million years.	76
4.4	Solar magnetogram, and frequency of emerging magnetic bipoles.	79
4.5	Butterfly diagram: sunspot latitudes versus time.	80
4.6	Two possible flux-rope dynamos.	82
4.7	Photospheric convection spectrum and solar internal rotation.	83
4.8	Columnar convection in a rotating spherical shell.	85
4.9	Effects of large-scale flows, and a simulated magnetogram of a young Sun.	88
4.10	Ingredients of dynamo models.	93
4.11	A minimal linear $\alpha\Omega$ dynamo solution.	94
4.12	A Babcock-Leighton dynamo solution.	98
4.13	Velocity and magnetic field for a planetary dynamo model.	107
5.1	Shocks around CMEs, the heliosphere, and Earth's magnetosphere.	115
5.2	Diagram: upstream and downstream of a shock.	116
5.3	Iso-contours of shock heating as a function of $\theta_{B\perp}$ and M_A .	118
5.4	Magnetospheric field subject to a Carrington-level storm.	119
5.5	Spiraling asterospheric magnetic field for different rotation periods.	121
5.6	Sketch of the heliospheric current sheet.	122
5.7	The Sun's magnetic field and its extension into the heliosphere.	123
5.8	Solar wind flow speed and pressure (1D model) for a high-speed stream.	124
5.9	A corotating interaction region in the solar equatorial plane.	127
5.10	Plasma flow onto (a) a non-conducting body and (b) a conducting body.	129

5.11	Draping solar magnetic flux tubes around a conducting ionosphere.	131
5.12	A magnetically closed magnetosphere.	133
5.13	Total pressure (magnetic plus plasma) with distance from a planet.	136
5.14	A magnetically open magnetosphere.	137
5.15	Numerical model of the magnetosphere of Ganymede.	139
5.16	Magnetospheric convection on the Earth's surface.	144
6.1	Overview of the electromagnetic spectrum.	151
6.2	Evolution of flare emissions, and concepts of particle acceleration.	152
6.3	Sketch of the flow of energy during a flare.	153
6.4	Geomagnetic field variation for two characteristic magnetic storms.	154
6.5	Topology of the magnetotail in a wind-dominated magnetosphere.	159
6.6	Ribbon motion in solar flare reconnection.	159
6.7	Three different ways to use magnetic energy to power a flare or CME.	163
6.8	Ideal-MHD evolution of a 2D arcade with a flux rope.	166
6.9	An isolated toroidal flux rope.	168
6.10	A stable toroidal equilibrium.	169
6.11	The three-dimensional flux-rope model).	169
6.12	Current density surfaces for an unstable Titov-Démoulin equilibrium.	171
6.13	Numerical simulation of a storage model.	172
6.14	Observation of a CME, and CME kinetic energy <i>vs.</i> flare class.	176
6.15	The Sweet-Parker and Petschek field configurations.	179
6.16	Syrovatskii's field configuration.	180
7.1	Angular momentum transfer in a shearing disk.	189
7.2	Connected accretion disk, stellar wind, and stellar magnetosphere.	190
7.3	Two bodies orbiting their barycenter, and their tidal acceleration.	192
7.4	Distant-future diameter of the Sun and the size of Earth's orbit.	195
8.1	Particle motions in magnetic field.	202
8.2	Schematic diagram for the gradient-B drift.	202
8.3	1-dimensional particle diffusion from a point-source.	215
8.4	SEP event associated with an impulsive solar flare.	215
8.5	Energy spectra of energetic particles in the heliosphere and for GCRs.	217
8.6	Modulation of galactic cosmic rays during five sunspot cycles.	219
8.7	Drift motion of cosmic rays in a Parker-spiral heliospheric field.	220
8.8	Normal-incidence frame (NIF) <i>vs.</i> de Hoffman-Teller frame (HTF).	223
8.9	2-D hybrid simulation of the solar wind – magnetosphere interaction.	226
8.10	Particle flux profiles from three different solar longitudes.	228
8.11	An electron in drift resonance with a ULF wave.	231
8.12	Characteristic wave types within a magnetosphere.	232
9.1	Granulation in quiet-Sun and plage observed in the G-band.	238
9.2	Solar flux-tube model.	239
9.3	The multitude of scales in the solar magnetic field.	242
10.1	Luminosity, surface temperature, and age of evolving Sun-like stars.	249
10.2	Evolution of the luminosity of the Sun over its full life span.	250
10.3	Age-activity relationship for main-sequence stars.	253
10.4	Examples of cycles in stellar chromospheric activity.	254
10.5	Moments of inertia of evolving stars.	256
10.6	A 2.5D axisymmetric, hydrodynamic model of the heliosphere.	262

10.7	Heliosphere-ISM models for different ISM parameters.	264
10.8	Lyman- α : the journey of photons; the emission line; and α Cen B	267
10.9	Mass-loss rate versus X-ray surface flux density for main sequence stars.	268
10.10	The inferred mass-loss history of the Sun.	269
10.11	Surface field maps and associated heliospheric fields.	270
11.1	Planet-metallicity correlation for gas-giant planets.	276
11.2	Structure and processes of protoplanetary disks.	280
11.3	Masses and radii of selected exoplanets and Solar-System planets.	282
11.4	Accreting protostar and likely accretional history of a low-mass star.	284
11.5	Optical image of the accreting young star HH 30.	285
11.6	HR diagram: Taurus protostars and young (T Tauri) stars.	288
11.7	Magnetosphere-disk interaction in low-mass, pre-main sequence stars.	291
11.8	Estimates of the minimum mass solar nebula.	294
11.9	Primordial disk fractions of stars in young clusters.	296
12.1	Earth's temperature and H ₂ O and CO ₂ after the Moon-forming collision.	301
12.2	The photosynthetic habitable zone (pHZ) over time.	304
12.3	Sun-Earth energy flow: photons, energetic particles, solar wind.	311
12.4	Exchanges of solar and terrestrial energy in the Earth's atmosphere.	313
12.5	Simplified model of radiative exchange in the Earth's atmosphere.	314
12.6	Radiative forcing (RF) bar chart for Earth's climate.	316
12.7	Earth's orbital parameters for the past million years.	319
12.8	Source and loss mechanisms for planetary atmospheres.	320
12.9	Flowchart for energization and escape of atmospheric particles.	323
13.1	Ionospheric chemistry in the upper atmospheres of Earth, Venus, Mars.	338
13.2	Model of the Earth's ionospheric density profiles.	342
13.3	Model for the mean Earth-ionospheric/thermospheric temperature profiles.	343
13.4	Global mean ion and electron profiles for Earth <i>vs.</i> solar activity.	344
13.5	Spectral radiance versus age of solar-type stars.	347
13.6	Evolution of the exospheric temperature.	348
13.7	Earth's exospheric temperatures as function of CO ₂ and solar XUV flux.	348
13.8	Field lines for untilted and tilted planetary dipoles in the solar wind.	351
13.9	Evolution of the minimum and maximum stellar wind densities at 1 AU.	353
14.1	Sources of energetic particles in the heliosphere.	355
14.2	Differential GCR proton fluxes for different levels of solar activity.	357
14.3	¹⁰ Be production rate <i>vs.</i> geomagnetic field strength and solar activity.	359
14.4	¹⁰ Be data and geomagnetic dipole field for the past 60,000 years.	362
14.5	Bandpass-filtered ¹⁰ Be concentration and sunspot number.	363
14.6	Solar modulation function Φ from the present to 9350 BP	364
15.1	Systems to scale: Jupiter, TRAPPIST-1, and the inner Solar System.	368

List of Boxes and Tables

0.1	Heliophysics: definition.	vii
0.2	Heliophysics chapters sorted by theme (1).	ix
0.2	Heliophysics chapters sorted by theme (2).	x
0.3	Heliophysics and space weather	xi
1.1	Basic glossary.	5
2.1	Climates of terrestrial planets.	15
2.2	Chemical species in upper atmospheres.	17
2.3	Domains in the solar atmosphere, with fundamental properties.	20
2.4	Fast and slow solar wind; basic parameters.	23
2.5	Selected plasma quantities; summary.	35
3.1	Basic magnetic structures: sheets, tubes, cells.	46
3.2	MHD approximation and the concept of 'closure'.	50
3.3	MHD equations.	51
3.4	Plasma parameters in different environments.	68
5.1	Plasma properties upstream of solar-system bodies.	113
5.2	Properties of the solar wind near the planets.	138
5.3	Intrinsic magnetic fields of Solar System bodies.	138
6.1	Solar flare classifications.	150
12.1	Temperatures of the planets for different albedo and stellar luminosity.	312
12.2	Solar wind and interplanetary magnetic field at the terrestrial planets.	325
14.1	Production reactions and rates for cosmogenic radionuclides.	358

Bibliography

- Abbett, W. P.: 2007, *ApJ* 665(2), 1469, doi:10.1086/519788
- Asai, Ayumi, Yokoyama, Takaaki, Shimojo, Masumi, Masuda, Satoshi, Kurokawa, Hiroki, & Shibata, Kazunari: 2004, *ApJ* 611(1), 557, doi:10.1086/422159
- Beer, J., Baumgartner, S. T., Dittrich-Hannen, B., Hauenstein, J., Kubik, P., Lukasczyk, C., Mende, W., Stellmacher, B., & Suter, M.: 1994, in J. M. Pap, C. Frohlich, H. S. Hudson, and S. K. Solanki (Eds.), *Invited Papers from IAU Colloquium 143: The Sun as a Variable Star: Solar and Stellar Irradiance Variations*, 291
- Behar, E., Nilsson, H., Alho, M., Goetz, C., & Tsurutani, B.: 2017, *Mon. Not. R. Astron. Soc.* 469, S396, doi:10.1093/mnras/stx1871
- Bell, Cameron P. M., Mamajek, Eric E., & Naylor, Tim: 2015, *Mon. Not. R. Astron. Soc.* 454(1), 593, doi:10.1093/mnras/stv1981
- Bell, Cameron P. M., Naylor, Tim, Mayne, N. J., Jeffries, R. D., & Littlefair, S. P.: 2013, *Mon. Not. R. Astron. Soc.* 434(1), 806, doi:10.1093/mnras/stt1075
- Benz, A: 2002, *Plasma Astrophysics, second edition*, Kluwer, Springer, Astrophysics and Space Science Library, Vol. 279
- Booth, R. S., Poppenhaeger, K., Watson, C. A., Silva Aguirre, V., & Wolk, S. J.: 2017, *Mon. Not. R. Astron. Soc.* 471(1), 1012, doi:10.1093/mnras/stx1630
- Bougher, S. W. & Roble, R. G.: 1991, *J. Geophys. Research* 96(A7), 11045, doi:10.1029/91JA01162
- Brun, Allan Sacha, Miesch, Mark S., & Toomre, Juri: 2004, *ApJ* 614(2), 1073, doi:10.1086/423835
- Burkepile, J. T., Hundhausen, A. J., Stanger, A. L., St. Cyr, O. C., & Seiden, J. A.: 2004, *Journal of Geophysical Research (Space Physics)* 109(A3), A03103, doi:10.1029/2003JA010149
- Cauley, P. Wilson, Redfield, Seth, Jensen, Adam G., Barman, Travis, Endl, Michael, & Cochran, William D.: 2015, *ApJ* 810(1), 13, doi:10.1088/0004-637X/810/1/13
- Charbonneau, Paul: 2010, *Living Reviews in Solar Physics* 7(1), 3, doi:10.12942/lrsp-2010-3
- Charbonneau, Paul: 2014, *Ann. Rev. Astron. Astrophys.* 52, 251, doi:10.1146/annurev-astro-081913-040012
- Cheung, M. C. M., Rempel, M., Chintzoglou, G., Chen, F., Testa, P., Martínez-Sykora, J., Sainz Dalda, A., DeRosa, M. L., Malanushenko, A., Hansteen, V., De Pontieu, B., Carlsson, M., Gudiksen, B., & McIntosh, S. W.: 2019, *Nature Astronomy* 3, 160, doi:10.1038/s41550-018-0629-3
- Cohen, O. & Drake, J. J.: 2014, *ApJ* 783(1), 55, doi:10.1088/0004-637X/783/1/55

- Cohen, O., Drake, J. J., & Kóta, J.: 2012, *ApJ* 760(1), 85, doi:10.1088/0004-637X/760/1/85
- Crooker, N. U., Gosling, J. T., Bothmer, V., Forsyth, R. J., Gazis, P. R., Hewish, A., Horbury, T. S., Intriligator, D. S., Jokipii, J. R., Kóta, J., Lazarus, A. J., Lee, M. A., Lucek, E., Marsch, E., Posner, A., Richardson, I. G., Roelof, E. C., Schmidt, J. M., Siscoe, G. L., Tsurutani, B. T., & Wimmer-Schweingruber, R. F.: 1999, *Space Sci. Rev.* 89, 179, doi:10.1023/A:1005253526438
- Deming, Drake, Louie, Dana, & Sheets, Holly: 2019, *PASPs* 131(995), 013001, doi:10.1088/1538-3873/aae5c5
- Desch, S. J.: 2007, *ApJ* 671(1), 878, doi:10.1086/522825
- Dravins, Dainis, Ludwig, Hans-Günter, Dahmén, Erik, & Pazira, Hiva: 2017a, *A&A* 605, A90, doi:10.1051/0004-6361/201730900
- Dravins, Dainis, Ludwig, Hans-Günter, Dahmén, Erik, & Pazira, Hiva: 2017b, *A&A* 605, A91, doi:10.1051/0004-6361/201730901
- Drury, L. Oc.: 1983, *Reports on Progress in Physics* 46(8), 973, doi:10.1088/0034-4885/46/8/002
- Durrant, Dale R. & Arakawa, Akio: 2007, *Comptes Rendus Mecanique* 335(9), 655, doi:10.1016/j.crme.2007.08.010
- Foley, Bradford J. & Smye, Andrew J.: 2018, *Astrobiology* 18(7), 873, doi:10.1089/ast.2017.1695
- Forbes, T. G. & Priest, E. R.: 1995, *ApJ* 446, 377, doi:10.1086/175797
- Franck, S., Block, A., Bloh, W., Bounama, C., Garrido, I., & Schellnhuber, H. J.: 2001, *Naturwissenschaften* 88(10), 416, doi:10.1007/s001140100257
- Fraschetti, F., Drake, J. J., Alvarado-Gómez, J. D., Moschou, S. P., Garraffo, C., & Cohen, O.: 2019, *ApJ* 874(1), 21, doi:10.3847/1538-4357/ab05e4
- Garraffo, Cecilia, Drake, Jeremy J., Cohen, Ofer, Alvarado-Gómez, Julian D., & Moschou, Sofia P.: 2017, *ApJL* 843(2), L33, doi:10.3847/2041-8213/aa79ed
- Gleeson, L. J. & Axford, W. I.: 1968, *ApJ* 154, 1011, doi:10.1086/149822
- Güdel, Manuel: 2007, *Living Reviews in Solar Physics* 4(1), 3, doi:10.12942/lrsp-2007-3
- Hartmann, Lee: 2009, *Accretion Processes in Star Formation: Second Edition*, Cambridge University Press, Cambridge, UK
- Hathaway, D. H., Beck, J. G., Bogart, R. S., Bachmann, K. T., Khatri, G., Petitto, J. M., Han, S., & Raymond, J.: 2000, *Solar Phys.* 193, 299, doi:10.1023/A:1005200809766
- Henning, Thomas & Semenov, Dmitry: 2013, *Chemical Reviews* 113(12), 9016, doi:10.1021/cr400128p
- Holzer, Thomas E.: 1972, *J. Geophys. Research* 77(28), 5407, doi:10.1029/JA077i028p05407
- Howard, Andrew W., Sanchis-Ojeda, Roberto, Marcy, Geoffrey W., Johnson, John Asher, Winn, Joshua N., Isaacson, Howard, Fischer, Debra A., Fulton, Benjamin J., Simukoff, Evan, & Fortney, Jonathan J.: 2013, *Nature* 503(7476), 381, doi:10.1038/nature12767
- IPCC: 2013, in T.F. Stocker, D. Qin, G.-K. Plattner, M. Tignor, S.K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex, and P.M. Midgley (Eds.), *Summary for Policymakers*, Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA
- Jansen, Tiffany, Scharf, Caleb, Way, Michael, & Del Genio, Anthony: 2019, *ApJ* 875(2), 79, doi:10.3847/1538-4357/ab113d
- Jia, Xianzhe, Walker, Raymond J., Kivelson, Margaret G., Khurana, Krishan K., & Linker, Jon A.: 2008, *Journal of Geophysical Research (Space Physics)* 113(A6), A06212, doi:10.1029/2007JA012748

- Johnson, Jennifer A.: 2019, *Science* 363(6426), 474, doi:10.1126/science.aau9540
- Jokipii, J. R. & Thomas, B.: 1981, *ApJ* 243, 1115, doi:10.1086/158675
- Kallenrode, M. B.: 2003, *Journal of Physics G Nuclear Physics* 29(5), 965
- Karak, Bidya Binay & Miesch, Mark: 2017, *ApJ* 847(1), 69, doi:10.3847/1538-4357/aa8636
- Kiehl, J. T. & Trenberth, Kevin E.: 1997, *Bulletin of the American Meteorological Society* 78(2), 197, doi:10.1175/1520-0477(1997)078<0197:EAGMEB>2.0.CO;2
- Koch, Alexander, Brierley, Chris, Maslin, Mark M., & Lewis, Simon L.: 2019, *Quaternary Science Reviews* 207, 13, doi:10.1016/j.quascirev.2018.12.004
- Krauss-Varban, Dietmar, Li, Yan, & Luhmann, Janet G.: 2008, in Gang Li, Qiang Hu, Olga Verkhoglyadova, Gary P. Zank, R. P. Lin, and J. Luhmann (Eds.), *American Institute of Physics Conference Series*, Vol. 1039, p. 307
- Kulikov, Yuri N., Lammer, Helmut, Lichtenegger, Herbert I. M., Penz, Thomas, Breuer, Doris, Spohn, Tilman, Lundin, Rickard, & Biernat, Helfried K.: 2007, *Space Sci. Rev.* 129(1-3), 207, doi:10.1007/s11214-007-9192-4
- Lammer, Helmut & Blanc, Michel: 2018, *Space Sci. Rev.* 214(2), 60, doi:10.1007/s11214-017-0433-x
- Lemerle, Alexandre & Charbonneau, Paul: 2017, *ApJ* 834(2), 133, doi:10.3847/1538-4357/834/2/133
- Linsky, J. L.: 1985, *Solar Phys.* 100, 333, doi:10.1007/BF00158435
- Lundin, Rickard, Lammer, Helmut, & Ribas, Ignasi: 2007, *Space Sci. Rev.* 129(1-3), 245, doi:10.1007/s11214-007-9176-4
- MacNeice, P., Antiochos, S. K., Phillips, A., Spicer, D. S., DeVore, C. R., & Olson, K.: 2004, *ApJ* 614(2), 1028, doi:10.1086/423887
- Mamajek, Eric E.: 2009, in Tomonori Usuda, Motohide Tamura, and Miki Ishii (Eds.), *American Institute of Physics Conference Series*, Vol. 1158 of *American Institute of Physics Conference Series*, p. 3
- Marcq, Emmanuel, Mills, Franklin P., Parkinson, Christopher D., & Vandaele, Ann Carine: 2018, *Space Sci. Rev.* 214(1), 10, doi:10.1007/s11214-017-0438-5
- Mazur, J. E., Mason, G. M., Dwyer, J. R., Giacalone, J., Jokipii, J. R., & Stone, E. C.: 2000, *ApJL* 532(1), L79, doi:10.1086/312561
- McComas, D. J., Dayeh, M. A., Funsten, H. O., Livadiotis, G., & Schwadron, N. A.: 2013, *ApJ* 771(2), 77, doi:10.1088/0004-637X/771/2/77
- Moraal, H.: 1976, *Space Sci. Rev.* 19(6), 845, doi:10.1007/BF00173707
- Müller, H. R., Frisch, P. C., Fields, B. D., & Zank, G. P.: 2009, *Space Sci. Rev.* 143(1-4), 415, doi:10.1007/s11214-008-9448-7
- Müller, H. R. & Zank, G. P.: 2004, *Journal of Geophysical Research (Space Physics)* 109(A7), A07104, doi:10.1029/2003JA010269
- Ngwira, Chigomezyo M., Pulkkinen, Antti, Kuznetsova, Maria M., & Glocer, Alex: 2014, *Journal of Geophysical Research (Space Physics)* 119(6), 4456, doi:10.1002/2013JA019661
- Ó Fionnagáin, D., Vidotto, A. A., Petit, P., Folsom, C. P., Jeffers, S. V., Marsden, S. C., Morin, J., do Nascimento, J. D., & BCooll Collaboration: 2019, *Mon. Not. R. Astron. Soc.* 483(1), 873, doi:10.1093/mnras/sty3132
- O'Brien, David P., Walsh, Kevin J., Morbidelli, Alessandro, Raymond, Sean N., & Mandell, Avi M.: 2014, *Icarus* 239, 74, doi:10.1016/j.icarus.2014.05.009
- Pavlov, Alexander A., Toon, Owen B., Pavlov, Anatoli K., Bally, John, & Pollard, David: 2005, *Geophys. Res. Lett.* 32(3), L03705, doi:10.1029/2004GL021890
- Pecaut, Mark J. & Mamajek, Eric E.: 2016, *Mon. Not. R. Astron. Soc.* 461(1), 794, doi:10.1093/mnras/stw1300
- Pecaut, Mark J., Mamajek, Eric E., & Bubar, Eric J.: 2012, *ApJ* 746(2), 154,

- doi:10.1088/0004-637X/746/2/154
- Pineda, J. Sebastian, Hallinan, Gregg, & Kao, Melodie M.: 2017, *ApJ* 846(1), 75, doi:10.3847/1538-4357/aa8596
- Pinhas, Arazi, Rackham, Benjamin V., Madhusudhan, Nikku, & Apai, Dániel: 2018, *Mon. Not. R. Astron. Soc.* 480(4), 5314, doi:10.1093/mnras/sty2209
- Pognan, Quentin, Garraffo, Cecilia, Cohen, Ofer, & Drake, Jeremy J.: 2018, *ApJ* 856(1), 53, doi:10.3847/1538-4357/aaaebb
- Portegies Zwart, S., Pelupessy, I., van Elteren, A., Wijnen, T. P. G., & Lugaro, M.: 2018, *A&A* 616, A85, doi:10.1051/0004-6361/201732060
- Quenby, J. J.: 1984, *Space Sci. Rev.* 37(3-4), 201, doi:10.1007/BF00226364
- Rackham, Benjamin, Pinhas, Arazi, Apai, Dániel, Haywood, Raphaëlle, Cegla, Heather, Espinoza, Néstor, Teske, Johanna, Gully-Santiago, Michael, Rau, Gioia, Morris, Brett M., Angerhausen, Daniel, Barclay, Thomas, Carone, Ludmila, Cauley, P. Wilson, de Wit, Julien, Domagal-Goldman, Shawn, Dong, Chuanfei, Dragomir, Diana, Giampapa, Mark S., Hasegawa, Yasuhiro, Hinkel, Natalie R., Hu, Renyu, Jordán, Andrés, Kitiashvili, Irina, Kreidberg, Laura, Lisse, Carey, Llama, Joe, López-Morales, Mercedes, Mennesson, Bertrand, Molaverdikhani, Karan, Osip, David J., & Quintana, Elisa V.: 2019a, *Bull. Am. Astron. Soc.* 51(3), 328
- Rackham, Benjamin V., Apai, Dániel, & Giampapa, Mark S.: 2018, *ApJ* 853(2), 122, doi:10.3847/1538-4357/aaa08c
- Rackham, Benjamin V., Apai, Dániel, & Giampapa, Mark S.: 2019b, *Astron. J.* 157(3), 96, doi:10.3847/1538-3881/aaf892
- Reames, Donald V.: 2013, *Space Sci. Rev.* 175(1-4), 53, doi:10.1007/s11214-013-9958-9
- Ribas, Ignasi, Guinan, Edward F., Güdel, Manuel, & Audard, Marc: 2005, *ApJ* 622(1), 680, doi:10.1086/427977
- Rosner, R., Tucker, W. H., & Vaiana, G. S.: 1978, *ApJ* 220, 643, doi:10.1086/155949
- Rutten, Robert J.: 2003, *Radiative Transfer in Stellar Atmospheres*, Univ. Utrecht, The Netherlands
- Sackmann, I. Juliana, Boothroyd, Arnold I., & Kraemer, Kathleen E.: 1993, *ApJ* 418, 457, doi:10.1086/173407
- Schatten, Kenneth H., Wilcox, John M., & Ness, Norman F.: 1969, *Solar Phys.* 6(3), 442, doi:10.1007/BF00146478
- Scheucher, Markus, Grenfell, J. L., Wunderlich, F., Godolt, M., Schreier, F., & Rauer, H.: 2018, *ApJ* 863(1), 6, doi:10.3847/1538-4357/aacf03
- Schrijver, Carolus J.: 2001, *ApJ* 547(1), 475, doi:10.1086/318333
- Schrijver, Carolus J.: 2009, *ApJL* 699(2), L148, doi:10.1088/0004-637X/699/2/L148
- Schrijver, C. J., Bagenal, F., & Sojka, J. J. (Eds.): 2016, *Heliophysics: Active Stars, their Astrospheres, and Impacts on Planetary Environments*, Cambridge University Press, Cambridge, UK, (Volume IV)
- Schrijver, Carolus J., Kauristie, Kirsti, Aylward, Alan D., Denardini, Clezio M., Gibson, Sarah E., Glover, Alexi, Gopalswamy, Nat, Grande, Manuel, Hapgood, Mike, Heynderickx, Daniel, Jakowski, Norbert, Kalegaev, Vladimir V., Lapenta, Giovanni, Linker, Jon A., Liu, Siqing, Mandrini, Cristina H., Mann, Ian R., Nagatsuma, Tsutomu, Nandy, Dibyendu, Obara, Takahiro, Paul O'Brien, T., Onsager, Terrance, Opgenoorth, Hermann J., Terkildsen, Michael, Valladares, Cesar E., & Vilmer, Nicole: 2015, *Advances in Space Research* 55(12), 2745, doi:10.1016/j.asr.2015.03.023
- Schrijver, Carolus J. & Siscoe, George L. (Eds.): 2011, *Heliophysics: Plasma Physics of the Local Cosmos*, Cambridge University Press, Cambridge, UK, (Volume I)
- Schrijver, Carolus J. & Siscoe, George L. (Eds.): 2012a, *Heliophysics: Evolving Solar Activity and the Climates of Space and Earth*, Cambridge University Press,

- Cambridge, UK, (Volume III)
- Schrijver, Carolus J. & Siscoe, George L. (Eds.): 2012b, *Heliophysics: Space Storms and Radiation: Causes and Effects*, Cambridge University Press, Cambridge, UK, (Volume II)
- Schrijver, Carolus J. & Siscoe, George L. (Eds.): 2015, *Heliophysics: Space Weather and Society*, published online at NASA's Heliophysics Summer School site, <https://cpaess.ucar.edu/sites/default/files/heliophysics/documents/HSS5.pdf> (Volume V)
- Schrijver, Carolus J. & Zwaan, Cornelis: 2000, *Solar and Stellar Magnetic Activity*, Cambridge University Press, Cambridge, UK
- Schrijver, Karel: 2018, *One of ten billion Earths: How we Learn about our Planet's Past and Future from Distant Exoplanets*, Oxford University Press, Oxford, UK
- Schröder, K. P. & Connon Smith, Robert: 2008, *Mon. Not. R. Astron. Soc.* 386(1), 155, doi:10.1111/j.1365-2966.2008.13022.x
- Schwadron, N. A., Moebius, E., Kucharek, H., Lee, M. A., French, J., Saul, L., Wurz, P., Bzowski, M., Fuselier, S. A., Livadiotis, G., McComas, D. J., Frisch, P., Gruntman, M., & Mueller, H. R.: 2013, *ApJ* 775(2), 86, doi:10.1088/0004-637X/775/2/86
- Smithtro, C. G. & Sojka, J. J.: 2005, *Journal of Geophysical Research (Space Physics)* 110(A8), A08305, doi:10.1029/2004JA010781
- Spruit, H. C.: 2013, arXiv e-prints arXiv:1301.5572
- Stern, S. A., Weaver, H. A., Spencer, J. R., Olkin, C. B., Gladstone, G. R., Grundy, W. M., Moore, J. M., Cruikshank, D. P., Elliott, H. A., McKinnon, W. B., Parker, J. Wm., Verbiscer, A. J., Young, L. A., Aguilar, D. A., Albers, J. M., Andert, T., Andrews, J. P., Bagenal, F., Banks, M. E., Bauer, B. A., Bauman, J. A., Bechtold, K. E., Beddingfield, C. B., Behrooz, N., Beisser, K. B., Benecchi, S. D., Bernardoni, E., Beyer, R. A., Bhaskaran, S., Bierson, C. J., Binzel, R. P., Birath, E. M., Bird, M. K., Boone, D. R., Bowman, A. F., Bray, V. J., Britt, D. T., Brown, L. E., Buckley, M. R., Buie, M. W., Buratti, B. J., Burke, L. M., Bushman, S. S., Carcich, B., Chaikin, A. L., Chavez, C. L., Cheng, A. F., Colwell, E. J., Conard, S. J., Conner, M. P., Conrad, C. A., Cook, J. C., Cooper, S. B., Custodio, O. S., Dalle Ore, C. M., Deboy, C. C., Dharmavaram, P., Dhingra, R. D., Dunn, G. F., Earle, A. M., Egan, A. F., Eisig, J., El-Maarry, M. R., Engelbrecht, C., Enke, B. L., Ercol, C. J., Fattig, E. D., Ferrell, C. L., Finley, T. J., Firer, J., Fischetti, J., Folkner, W. M., Fosbury, M. N., Fountain, G. H., Freeze, J. M., Gabasova, L., Glaze, L. S., Green, J. L., Griffith, G. A., Guo, Y., Hahn, M., Hals, D. W., Hamilton, D. P., Hamilton, S. A., Hanley, J. J., Harch, A., Harmon, K. A., Hart, H. M., Hayes, J., Hersman, C. B., Hill, M. E., Hill, T. A., Hofgartner, J. D., Holdridge, M. E., Horányi, M., Hosadurga, A., Howard, A. D., Howett, C. J. A., Jaskulek, S. E., Jennings, D. E., Jensen, J. R., Jones, M. R., Kang, H. K., Katz, D. J., Kaufmann, D. E., Kavelaars, J. J., Keane, J. T., Keleher, G. P., Kinczyk, M., Kochte, M. C., Kollmann, P., Krimigis, S. M., Kruizinga, G. L., Kusnierkiewicz, D. Y., Lahr, M. S., Lauer, T. R., Lawrence, G. B., Lee, J. E., Lessac-Chenen, E. J., Linscott, I. R., Lisse, C. M., Lunsford, A. W., Mages, D. M., Mallder, V. A., Martin, N. P., May, B. H., McComas, D. J., McNutt, R. L., Mehoke, D. S., Mehoke, T. S., Nelson, D. S., Nguyen, H. D., Núñez, J. I., Ocampo, A. C., Owen, W. M., Oxtton, G. K., Parker, A. H., Pätzold, M., Pelgrift, J. Y., Pelletier, F. J., Pineau, J. P., Piquette, M. R., Porter, S. B., Protopapa, S., Quirico, E., Redfern, J. A., Regiec, A. L., Reitsema, H. J., Reuter, D. C., Richardson, D. C., Riedel, J. E., Ritterbush, M. A., Robbins, S. J., Rodgers, D. J., Rogers, G. D., Rose, D. M., Rosendall, P. E., Runyon, K. D., Ryschkewitsch, M. G., Saina, M. M., Salinas, M. J., Schenk, P. M., Scherrer, J. R., Schlei, W. R.,

- Schmitt, B., Schultz, D. J., Schurr, D. C., Scipioni, F., Sepan, R. L., Shelton, R. G., Showalter, M. R., Simon, M., Singer, K. N., Stahlheber, E. W., Stanbridge, D. R., Stansberry, J. A., Steffl, A. J., Strobel, D. F., Stothoff, M. M., Stryk, T., Stuart, J. R., Summers, M. E., Tapley, M. B., Taylor, A., Taylor, H. W., Tedford, R. M., Throop, H. B., Turner, L. S., Umurhan, O. M., Van Eck, J., Velez, D., Versteeg, M. H., Vincent, M. A., Webbert, R. W., Weidner, S. E., Weigle, G. E., Wendel, J. R., White, O. L., Whittenburg, K. E., Williams, B. G., Williams, K. E., Williams, S. P., Winters, H. L., Zangari, A. M., & Zurbuchen, T. H.: 2019, *Science* 364(6441), aaw9771, doi:10.1126/science.aaw9771
- Struminsky, A. B., Sadvovskii, A. M., & Zharikova, M. S.: 2018, *Geomagnetism and Aeronomy* 58(8), 1108, doi:10.1134/S0016793218080169
- Thompson, Michael J., Christensen-Dalsgaard, Jørgen, Miesch, Mark S., & Toomre, Juri: 2003, *Ann. Rev. Astron. Astrophys.* 41, 599, doi:10.1146/annurev.astro.41.011802.094848
- Titov, V. S. & Démoulin, P.: 1999, *A&A* 351, 707
- Török, T., Kliem, B., & Titov, V. S.: 2004, *A&A* 413, L27, doi:10.1051/0004-6361:20031691
- Tsurutani, Bruce T., Gonzalez, Walter D., Gonzalez, Alicia L. C., Guarnieri, Fernando L., Gopalswamy, Nat, Grande, Manuel, Kamide, Yohsuke, Kasahara, Yoshiya, Lu, Gang, Mann, Ian, McPherron, Robert, Soraas, Finn, & Vasyliunas, Vytenis: 2006, *Journal of Geophysical Research (Space Physics)* 111(A7), A07S01, doi:10.1029/2005JA011273
- Vasyliunas, V. M.: 1976, in *Magnetospheric Particles and Fields*, p. 99
- Wallner, A., Feige, J., Kinoshita, N., Paul, M., Fifield, L. K., Golser, R., Honda, M., Linnemann, U., Matsuzaki, H., Merchel, S., Rugel, G., Tims, S. G., Steier, P., Yamagata, T., & Winkler, S. R.: 2016, *Nature* 532(7597), 69, doi:10.1038/nature17196
- Wieler, Rainer, Beer, Jürg, & Leya, Ingo: 2013, *Space Sci. Rev.* 176(1-4), 351, doi:10.1007/s11214-011-9769-9
- Wood, Brian E.: 2004, *Living Reviews in Solar Physics* 1, 2, doi:10.12942/lrsp-2004-2
- Wood, Brian E., Müller, Hans-Reinhard, Redfield, Seth, & Edelman, Eric: 2014, *ApJL* 781(2), L33, doi:10.1088/2041-8205/781/2/L33
- Wood, Brian E., Müller, Hans-Reinhard, Zank, Gary P., & Linsky, Jeffrey L.: 2002, *ApJ* 574(1), 412, doi:10.1086/340797
- Wood, B. E., Müller, H. R., Zank, G. P., Linsky, J. L., & Redfield, S.: 2005, *ApJL* 628(2), L143, doi:10.1086/432716
- Zahnle, Kevin, Arndt, Nick, Cockell, Charles, Halliday, Alex, Nisbet, Euan, Selsis, Franck, & Sleep, Norman H.: 2007, *Space Sci. Rev.* 129(1-3), 35, doi:10.1007/s11214-007-9225-z
- Zhang, Zhanbo, Zhou, Yifan, Rackham, Benjamin V., & Apai, Dániel: 2018, *Astron. J.* 156(4), 178, doi:10.3847/1538-3881/aade4f