

Contents

Preface	xiii
1 Introduction	1
1.1 An overview of the observations	5
• Stars 5 • The Galaxy 11 • Other galaxies 19 ▷ Elliptical galaxies 20 ▷ Spiral galaxies 25 ▷ Lenticular galaxies 28 ▷ Irregular galaxies 28 • Open and globular clusters 29 • Groups and clusters of galaxies 30 • Black holes 32	
1.2 Collisionless systems and the relaxation time	33
• The relaxation time 34	
1.3 The cosmological context	37
• Kinematics 38 • Geometry 39 • Dynamics 40 • The Big Bang and inflation 45 • The cosmic microwave background 48	
Problems	52
2 Potential Theory	55
2.1 General results	56
▷ The potential-energy tensor 59	
2.2 Spherical systems	60
• Newton's theorems 60 ▷ Potential energy of spherical systems 63 • Potentials of some simple systems 63 ▷ Point mass 63 ▷ Homogeneous sphere 63 ▷ Plummer model 65 ▷ Isochrone potential 65 ▷ Modified Hubble model 66 ▷ Power-law density model 68 ▷ Two-power density models 70	
2.3 Potential-density pairs for flattened systems	72
• Kuzmin models and generalizations 72 • Logarithmic potentials 74 • Poisson's equation in very flattened systems 77	
2.4 Multipole expansion	78
2.5 The potentials of spheroidal and ellipsoidal systems	83
• Potentials of spheroidal shells 84 • Potentials of spheroidal systems 87 • Potentials of ellipsoidal systems 94 ▷ Ferrers potentials 95 ▷ Potential-energy tensors of ellipsoidal systems 95	

2.6	The potentials of disks	96
	• Disk potentials from homoeoids 96 ▷ The Mestel disk 99 ▷ The exponential disk 100 ▷ Thick disks 102 • Disk potentials from Bessel functions 103 ▷ Application to axisymmetric disks 106 • Disk potentials from logarithmic spirals 107 • Disk potentials from oblate spheroidal coordinates 109	
2.7	The potential of our Galaxy	110
	▷ The bulge 111 ▷ The dark halo 112 ▷ The stellar disk 112 ▷ The interstellar medium 112 ▷ The bulge as a bar 117	
2.8	Potentials from functional expansions	118
	▷ Bi-orthonormal basis functions 120 ▷ Designer basis functions 120	
2.9	Poisson solvers for N-body codes	122
	• Direct summation 123 ▷ Softening 123 • Tree codes 125 ▷ Cartesian multipole expansion 127 • Particle-mesh codes 129 ▷ Periodic boundary conditions 131 ▷ Vacuum boundary conditions 132 ▷ Mesh refinement 135 ▷ P ³ M codes 135 • Spherical-harmonic codes 136 • Simulations of planar systems 137	
	Problems	137
3	The Orbits of Stars	142
3.1	Orbits in static spherical potentials	143
	▷ Spherical harmonic oscillator 147 ▷ Kepler potential 147 ▷ Isochrone potential 149 ▷ Hyperbolic encounters 153 • Constants and integrals of the motion 155	
3.2	Orbits in axisymmetric potentials	159
	• Motion in the meridional plane 159 • Surfaces of section 162 • Nearly circular orbits: epicycles and the velocity ellipsoid 164	
3.3	Orbits in planar non-axisymmetric potentials	171
	• Two-dimensional non-rotating potential 171 • Two-dimensional rotating potential 178 • Weak bars 188 ▷ Lindblad resonances 188 ▷ Orbits trapped at resonance 193	
3.4	Numerical orbit integration	196
	• Symplectic integrators 197 ▷ Modified Euler integrator 197 ▷ Leapfrog integrator 200 • Runge–Kutta and Bulirsch–Stoer integrators 201 • Multistep predictor-corrector integrators 202 • Multivalued integrators 203 • Adaptive timesteps 205 • Individual timesteps 206 • Regularization 208 ▷ Burdet–Heggie regularization 208 ▷ Kustaanheimo–Stiefel (KS) regularization 210	
3.5	Angle-action variables	211
	• Orbital tori 212 ▷ Time averages theorem 215 ▷ Action space 216 ▷ Hamilton–Jacobi equation 217 • Angle-action variables for spherical potentials 220 • Angle-action variables for flattened axisymmetric potentials 226 ▷ Stäckel potentials 226	

▷ Epicycle approximation 231	• Angle-action variables for a non-rotating bar 234	• Summary 236	
3.6 Slowly varying potentials			237
• Adiabatic invariance of actions 237	• Applications 238		
▷ Harmonic oscillator 238	▷ Eccentric orbits in a disk 240		
▷ Transient perturbations 240	▷ Slow growth of a central black hole 241		
3.7 Perturbations and chaos			243
• Hamiltonian perturbation theory 243	• Trapping by resonances 246	▷ Levitation 250	• From order to chaos 253
▷ Irregular orbits 256	▷ Frequency analysis 258	▷ Liapunov exponents 260	
3.8 Orbits in elliptical galaxies			262
• The perfect ellipsoid 263	• Dynamical effects of cusps 263		
• Dynamical effects of black holes 266			
Problems			268
4 Equilibria of Collisionless Systems			274
4.1 The collisionless Boltzmann equation			275
• Limitations of the collisionless Boltzmann equation 278	▷ Finite stellar lifetimes 278	▷ Correlations between stars 279	• Relation between the DF and observables 280
▷ An example 282			
4.2 Jeans theorems			283
• Choice of f and relations between moments 285	▷ DF depending only on H 285	▷ DF depending on H and L 286	▷ DF depending on H and L_z 286
4.3 DFs for spherical systems			287
• Ergodic DFs for systems 288	▷ Ergodic Hernquist, Jaffe and isochrone models 290	▷ Differential energy distribution 292	
• DFs for anisotropic spherical systems 293	▷ Models with constant anisotropy 294	▷ Osipkov–Merritt models 297	
▷ Other anisotropic models 298	▷ Differential-energy distribution for anisotropic systems 299	• Spherical systems defined by the DF 299	▷ Polytropes and the Plummer model 300
▷ The isothermal sphere 302	▷ Lowered isothermal models 307	▷ Double-power models 311	▷ Michie models 312
4.4 DFs for axisymmetric density distributions			312
• DF for a given axisymmetric system 312	• Axisymmetric systems specified by $f(H, L_z)$ 314	▷ Fully analytic models 314	▷ Rowley models 318
▷ Rotation and flattening in spheroids 320	• The Schwarzschild DF 321		
4.5 DFs for razor-thin disks			329
• Mestel disk 329	• Kalnajs disks 330		
4.6 Using actions as arguments of the DF			333
• Adiabatic compression 335	▷ Cusp around a black hole 336	▷ Adiabatic deformation of dark matter 337	

4.7 Particle-based and orbit-based models	338
• N-body modeling 339 ▷ Softening 341 ▷ Instability and chaos 341 • Schwarzschild models 344	
4.8 The Jeans and virial equations	347
• Jeans equations for spherical systems 349 ▷ Effect of a central black hole on the observed velocity dispersion 350 • Jeans equations for axisymmetric systems 353 ▷ Asymmetric drift 354 ▷ Spheroidal components with isotropic velocity dispersion 356 • Virial equations 358 ▷ Scalar virial theorem 360 ▷ Spherical systems 361 ▷ The tensor virial theorem and observational data 362	
4.9 Stellar kinematics as a mass detector	365
• Detecting black holes 366 • Extended mass distributions of elliptical galaxies 370 • Dynamics of the solar neighborhood 372	
4.10 The choice of equilibrium	376
• The principle of maximum entropy 377 • Phase mixing and violent relaxation 379 ▷ Phase mixing 379 ▷ Violent relaxation 380 • Numerical simulation of the relaxation process 382	
Problems	387
5 Stability of Collisionless Systems	394
5.1 Introduction	394
• Linear response theory 396 • Linearized equations for stellar and fluid systems 398	
5.2 The response of homogeneous systems	401
• Physical basis of the Jeans instability 401 • Homogeneous systems and the Jeans swindle 401 • The response of a homogeneous fluid system 403 • The response of a homogeneous stellar system 406 ▷ Unstable solutions 410 ▷ Neutrally stable solutions 411 ▷ Damped solutions 412 • Discussion 416	
5.3 General theory of the response of stellar systems	417
• The polarization function in angle-action variables 418 • The Kalnajs matrix method 419 • The response matrix 421	
5.4 The energy principle and secular stability	423
• The energy principle for fluid systems 423 • The energy principle for stellar systems 427 • The relation between the stability of fluid and stellar systems 431	
5.5 The response of spherical systems	432
• The stability of spherical systems with ergodic DFs 432 • The stability of anisotropic spherical systems 433 ▷ Physical basis of the radial-orbit instability 434 • Landau damping and resonances in spherical systems 437	
5.6 The stability of uniformly rotating systems	439
• The uniformly rotating sheet 439 • Kalnajs disks 444 • Maclaurin spheroids and disks 449	

Problems	450
6 Disk Dynamics and Spiral Structure	456
6.1 Fundamentals of spiral structure	458
• Images of spiral galaxies 460	
• Spiral arms at other wave-lengths 462	
▷ Dust 464	
▷ Relativistic electrons 465	
▷ Molecular gas 465	
▷ Neutral atomic gas 465	
▷ HII regions 467	
• The geometry of spiral arms 468	
▷ The strength and number of arms 468	
▷ Leading and trailing arms 469	
▷ The pitch angle and the winding problem 471	
▷ The pattern speed 474	
• The anti-spiral theorem 477	
• Angular-momentum transport by spiral-arm torques 478	
6.2 Wave mechanics of differentially rotating disks	481
• Preliminaries 481	
▷ Kinematic density waves 481	
▷ Resonances 484	
• The dispersion relation for tightly wound spiral arms 485	
▷ The tight-winding approximation 485	
▷ Potential of a tightly wound spiral pattern 486	
▷ The dispersion relation for fluid disks 488	
▷ The dispersion relation for stellar disks 492	
• Local stability of differentially rotating disks 494	
• Long and short waves 497	
• Group velocity 499	
• Energy and angular momentum in spiral waves 503	
6.3 Global stability of differentially rotating disks	505
• Numerical work on disk stability 505	
• Swing amplifier and feedback loops 508	
▷ The swing amplifier 508	
▷ Feedback loops 512	
▷ Physical interpretation of the bar instability 513	
• The maximum-disk hypothesis 515	
• Summary 517	
6.4 Damping and excitation of spiral structure	518
• Response of the interstellar gas to a density wave 518	
• Response of a density wave to the interstellar gas 522	
• Excitation of spiral structure 524	
▷ Excitation by companion galaxies 524	
▷ Excitation by bars 525	
▷ Stationary spiral structure 525	
▷ Excitation of intermediate-scale structure 526	
6.5 Bars	528
• Observations 528	
▷ The pattern speed 531	
• Dynamics of bars 533	
▷ Weak bars 534	
▷ Strong bars 535	
▷ The vertical structure of bars 536	
▷ Gas flow in bars 536	
▷ Slow evolution of bars 539	
6.6 Warping and buckling of disks	539
• Warps 539	
▷ Kinematics of warps 540	
▷ Bending waves with self-gravity 542	
▷ The origin of warps 544	
• Buckling instability 548	
Problems	552
7 Kinetic Theory	554
7.1 Relaxation processes	555
▷ Relaxation 555	
▷ Equipartition 556	
▷ Escape 556	
▷ Inelastic	

encounters 557	▷ Binary formation by triple encounters 557	
	▷ Interactions with primordial binaries 558	
7.2	General results	559
	• Virial theorem 559	
	• Liouville's theorem 561	
	• Reduced distribution functions 563	
	• Relation of Liouville's equation to the collisionless Boltzmann equation 565	
7.3	The thermodynamics of self-gravitating systems	567
	• Negative heat capacity 567	
	• The gravothermal catastrophe 568	
7.4	The Fokker–Planck approximation	573
	• The master equation 573	
	• Fokker–Planck equation 574	
	▷ Weak encounters 574	
	▷ Local encounters 576	
	▷ Orbit-averaging 577	
	• Fluctuation-dissipation theorems 578	
	• Diffusion coefficients 580	
	▷ Heating of the Galactic disk by MACHOs 583	
	• Relaxation time 586	
	• Numerical methods 588	
	▷ Fluid models 588	
	▷ Monte Carlo methods 592	
	▷ Numerical solution of the Fokker–Planck equation 593	
	▷ N-body integrations 594	
	▷ Checks and comparisons 595	
7.5	The evolution of spherical stellar systems	596
	• Mass loss from stellar evolution 600	
	• Evaporation and ejection 602	
	▷ The maximum lifetime of a stellar system 605	
	• Core collapse 606	
	• After core collapse 609	
	• Equipartition 612	
	• Tidal shocks and the survival of globular clusters 615	
	• Binary stars 616	
	▷ Soft binaries 618	
	▷ Hard binaries 620	
	▷ Reaction rates 621	
	• Inelastic encounters 625	
	• Stellar systems with a central black hole 629	
	▷ Consumption of stars by the black hole 629	
	▷ The effect of a central black hole on the surrounding stellar system 631	
7.6	Summary	633
	Problems	634
8	Collisions and Encounters of Stellar Systems	639
8.1	Dynamical friction	643
	▷ The validity of Chandrasekhar's formula 646	
	• Applications of dynamical friction 647	
	▷ Decay of black-hole orbits 647	
	▷ Galactic cannibalism 649	
	▷ Orbital decay of the Magellanic Clouds 650	
	▷ Dynamical friction on bars 651	
	▷ Formation and evolution of binary black holes 652	
	▷ Globular clusters 654	
8.2	High-speed encounters	655
	▷ Mass loss 657	
	▷ Return to equilibrium 657	
	▷ Adiabatic invariance 658	
	• The distant-tide approximation 658	
	• Disruption of stellar systems by high-speed encounters 661	
	▷ The catastrophic regime 662	
	▷ The diffusive regime 663	
	▷ Disruption of open clusters 664	
	▷ Disruption of binary stars 665	
	▷ Dynamical constraints on MACHOs 668	
	▷ Disk and bulge shocks 669	
	▷ High-speed interactions in clusters of galaxies 672	

8.3 Tides	674
• The restricted three-body problem 675 • The sheared-sheet or Hill's approximation 678 ▷ The epicycle approximation and Hill's approximation 679 ▷ The Jacobi radius in Hill's approximation 680 • Tidal tails and streamers 681	
8.4 Encounters in stellar disks	685
• Scattering of disk stars by molecular clouds 687 • Scattering of disk stars by spiral arms 691 • Summary 695	
8.5 Mergers	695
• Peculiar galaxies 696 • Grand-design spirals 698	
• Ring galaxies 699 • Shells and other fine structure 701	
• Starbursts 705 • The merger rate 708	
Problems	710
9 Galaxy Formation	716
9.1 Linear structure formation	717
• Gaussian random fields 719 ▷ Filtering 720 ▷ The Harrison–Zeldovich power spectrum 721 • Gravitational instability in the expanding universe 722 ▷ Non-relativistic fluid 722 ▷ Relativistic fluid 726	
9.2 Nonlinear structure formation	733
• Spherical collapse 733 • The cosmic web 735 • Press–Schechter theory 739 ▷ The mass function 744 ▷ The merger rate 746	
• Collapse and virialization in the cosmic web 748	
9.3 N-body simulations of clustering	751
• The mass function of halos 752 • Radial density profiles 753	
• Internal dynamics of halos 756 ▷ The shapes of halos 756	
▷ Rotation of halos 757 ▷ Dynamics of halo substructure 759	
9.4 Star formation and feedback	760
▷ Reionization 760 ▷ Feedback 761 ▷ Mergers, starbursts and quiescent accretion 762 ▷ The role of central black holes 764	
▷ Origin of the galaxy luminosity function 765	
9.5 Conclusions	765
Problems	766
Appendices	
A Useful numbers	770
B Mathematical background	771
• Vectors 771 • Curvilinear coordinate systems 773	
• Vector calculus 775 • Fourier series and transforms 778	
• Abel integral equation 780 • Schwarz's inequality 780	
• Calculus of variations 781 • Poisson distribution 781	
• Conditional probability and Bayes's theorem 782 • Central limit theorem 783	

C	Special functions	785
	• Delta function and step function 785 • Factorial or gamma function 786 • Error function, Dawson’s integral, and plasma dispersion function 786 • Elliptic integrals 787 • Legendre functions 788 • Spherical harmonics 789 • Bessel functions 790	
D	Mechanics	792
	• Single particles 792 • Systems of particles 794 • Lagrangian dynamics 797 • Hamiltonian dynamics 797 ▷ Hamilton’s equations 797 ▷ Poincaré invariants 799 ▷ Poisson brackets 800 ▷ Canonical coordinates and transformations 800 ▷ Extended phase space 803 ▷ Generating functions 803	
E	Delaunay variables for Kepler orbits	805
F	Fluid mechanics	807
	• Basic equations 807 ▷ Continuity equation 807 ▷ Euler’s equation 808 ▷ Energy equation 810 ▷ Equation of state 811 • The ideal gas 812 • Sound waves 813 ▷ Energy and momentum in sound waves 814 • Group velocity 817	
G	Discrete Fourier transforms	818
H	The Antonov–Lebovitz theorem	822
I	The Doremus–Feix–Baumann theorem	823
J	Angular-momentum transport in disks	825
	• Transport in fluid and stellar systems 825 • Transport in a disk with stationary spiral structure 826 • Transport in perturbed axisymmetric disks 828 • Transport in the WKB approximation 829	
K	Derivation of the reduction factor	830
L	The diffusion coefficients	833
M	The distribution of binary energies	838
	• The evolution of the energy distribution of binaries 838 • The two-body distribution function in thermal equilibrium 839 • The distribution of binary energies in thermal equilibrium 839 • The principle of detailed balance 841	
	References	842
	Index	857

Preface

Always majestic, often spectacularly beautiful, galaxies are the fundamental building blocks of the Universe. The inquiring mind cannot help asking how they formed, how they function, and what will become of them in the distant future. The principal tool used in answering these questions is stellar dynamics, the study of the motion of a large number of point masses orbiting under the influence of their mutual self-gravity. The main aim of this book is to provide the reader with an understanding of stellar dynamics at the level required to carry out research into the formation and dynamics of galaxies.

Galaxies are not only important in their own right, but also provide powerful tools for investigating some of the most important and fundamental problems in physics: our current expectation that the great majority of the matter in the universe is made up of weakly interacting elementary particles of unknown nature arose from studies of the outer reaches of galaxies; the standard theory of the origin of structure in the universe, which rests on exotic hypotheses such as inflation and vacuum energy, is tested and challenged by observations of the structure of galaxies; and galaxies are frequently used as enormous laboratories to study the laws of physics in extreme conditions.

The study of galactic dynamics carries the student to the frontiers of knowledge faster than almost any other branch of theoretical physics, in part because the fundamental issues in the subject are easy to understand for anyone with an undergraduate training in physics, and in part because theorists are scrambling to keep pace with a flood of new observations.

The tools required to reach the research frontier in galactic dynamics are for the most part ones developed in other fields: classical, celestial, and Hamiltonian mechanics, fluid mechanics, statistical mechanics, and plasma physics provide the most relevant backgrounds, and, although there is little need for quantum mechanics, the mathematical techniques developed in an introductory quantum mechanics course are in constant use. Brief summaries of the required background material are given in Appendices B (mathematics), C (special functions), D (mechanics), and F (fluid mechanics).

This book has been designed for readers with a standard undergraduate preparation in physics. By contrast, we have assumed no background in astronomy, although the context of many discussions will be clearer to the reader who has a broad grasp of basic astronomy and astrophysics, at the level of Shu (1982). A brief summary of the relevant observations is provided in §1.1. Introductions to galactic astronomy are given in Marochnik & Suchkov (1996), Elmegreen (1998) and Sparke & Gallagher (2000). For a comprehensive description of the properties of galaxies and other stellar systems, see *Galactic Astronomy* (Binney & Merrifield 1998), which is a companion to the present book and will be referred to simply as BM.

A one-semester graduate course on galaxies might be based on the following sections from the two books:

Galactic Astronomy: 1.1, 1.2, 2.1, 2.2, 2.3, 3.6, 3.7, 4.1, 4.2, 4.3, 4.6, 9.1, 10.3, 11.1

Galactic Dynamics: 1.1, 1.2, 2.1, 2.2, 2.3, 2.9, 3.1, 3.2, 3.3, 3.4, 4.1, 4.3, 4.7, 4.8, 4.9, 5.1, 5.2, 6.1, 7.1, 7.2, 7.3, 8.1, 8.3, 8.5

This selection assumes that the students will be exposed in other courses to the material on cosmology (GD §§1.3 and 9.1–9.3), stellar structure (GA §§3.1–3.5 and 5.1–5.4), and the interstellar medium (GA §§8.1–8.4 and 9.2–9.6).

The first edition of this book appeared in 1987, and after two decades major revisions were in order, both to accommodate the many important advances in the field and to reflect changes in the perspectives of the authors. The present edition makes more extensive use of Lagrangian and Hamiltonian mechanics, and includes new theoretical topics such as basis-function expansions for potential theory, angle-action variables, orbit-based methods of constructing stellar systems, linear response theory, stability analysis using the Kalnajs matrix method and energy principles, energy and angular-momentum transport in disks, the fluctuation-dissipation theorem, and the sheared sheet. N-body simulations of stellar systems have grown enormously in importance and sophistication, and we have added extensive descriptions of modern numerical methods for evaluating the gravitational field and following orbits.

The last two chapters of the first edition have been eliminated; much of the material of the old Chapter 9 is now covered in *Galactic Astronomy*, while the topics from the old Chapter 10 have been integrated into earlier chapters.

The most dramatic change in this subject since the publication of the first edition has been the development of a theory of structure formation in the universe of remarkable elegance and predictive power, based on concepts such as inflation, cold dark matter, and vacuum energy. The study of galactic structure and the study of cosmology are now inseparable, and one of our struggles was to incorporate the rich new insights that cosmology has brought to galactic dynamics without writing a superfluous textbook on cosmology. We have chosen to devote the final chapter to a short but self-contained outline of the contemporary theory of large-scale structure and galaxy formation. We introduce the theory of random fields and the linear theory of the growth of fluctuations in the universe, and give a simple treatment of nonlinear structure formation through the spherical-collapse model and extended Press–Schechter theory. The final, speculative section summarizes our still incomplete understanding of how the complex physics of baryons gives rise to the galaxies we see about us.

There are problems at the end of each chapter, many of which are intended to elucidate topics that are not fully covered in the main text. Their

degree of difficulty is indicated by a number in square brackets at the start of each problem, ranging from [1] (easy) to [3] (difficult).

One vexing issue in astrophysical notation is how to indicate approximate equality. We use “=” to denote equality to several significant digits, “ \approx ” for equality to order of magnitude, and “ \simeq ” for everything in between. The ends of proofs are indicated by a sideways triangle, “ \triangleleft ”.

Although we have tried to keep jargon to a minimum, we were unable to resist the economy of a few abbreviations: “distribution function” is written throughout the book as “DF,” “cosmic microwave background” as “CMB,” “initial mass function” as “IMF,” “interstellar medium” as “ISM,” “line-of-sight velocity distribution” as “LOSVD,” and “root mean square” as “RMS.”

We are deeply indebted to many colleagues, both for thoughtful comments on the first edition and for their patient and enthusiastic support during the preparation of the second.

Errors and inaccuracies in the first edition were pointed out to us by Ed Bertschinger, Carlo Del Noce, David Earn, Stefan Engström, Andreas Ernst, Gerry Gilmore, Chris Hunter, Doug Johnstone, Konrad Kuijken, Ari Laor, Blane Little, Thomas Lydon, Phil Mahoney, Kap-Soo Oh, Maria Petrou, Sterl Phinney, Gerald Quinlan, James Rhoads, George Rybicki, Jerry Sellwood, Yue Shen, David Sher, Min-Su Shin, Noam Soker, Guo-xuan Song, S. Sridhar, Björn Sundelius, Maria Sundin, Peter Teuben, Alexey Yurchenko, Rosemary Wyse, and Harold Zepolsky. Marc Kamionkowski and the students and postdocs in theoretical astrophysics at Caltech solved most of the problems and worked with us to improve them. Advice, data, and specially adapted or constructed figures were provided by Ron Allen, Lia Athanasoula, Jeremy Bailin, Rainer Beck, Kirk Borne, Daniela Calzetti, Michele Cappellari, Roc Cutri, John Dubinski, Doug Finkbeiner, Chris Flynn, Marc Freitag, Ortwin Gerhard, Mirek Giersz, Oleg Gnedin, Bill Harris, Clovis Hopman, Adrian Jenkins, Avi Loeb, Robert Lupton, John Magorrian, David Malin, Tom Oosterloo, Michael Perryman, Fred Rasio, Michael Regan, Jerry Sellwood, Tom Statler, Max Tegmark, Jihad Touma, Rien van de Weygaert, and Donghai Zhao. Major portions of this edition were read in draft form and critiqued by Luca Ciotti, Ken Freeman, Douglas Hoggie, Chris Hunter, Marc Kamionkowski, Jerry Sellwood, and Alar Toomre. We have benefited greatly from a powerful suite of plain \TeX macros written by Nigel Dowrick.

We also thank the many institutions that have provided support and hospitality to us during the writing of both editions of this book. These include the Institute of Astronomy, Cambridge; the Massachusetts Institute of Technology; the Max Planck Institute for Astrophysics, Garching; the Kapteyn Astronomical Institute, Groningen; Merton College and the Department of Physics, Oxford; the Institute for Advanced Study, Princeton; Princeton University Observatory; the Weizmann Institute of Science, Rehovot; and the Canadian Institute for Theoretical Astrophysics, Toronto.

We are greatly indebted to the Smithsonian/NASA Astrophysics Data System (see adsabs.harvard.edu) and the arXiv e-print service (arXiv.org),

which have revolutionized access to the astronomy literature.

Finally we thank our families for their support and understanding over the years in which *Galactic Dynamics* has encroached on times such as weekends and vacations that are properly reserved for family life.

July 2007

James Binney
Scott Tremaine

1

Introduction

A **stellar system** is a gravitationally bound assembly of stars or other point masses. Stellar systems vary over more than fourteen orders of magnitude in size and mass, from binary stars, to star clusters containing 10^2 to 10^6 stars, through galaxies containing 10^5 to 10^{12} stars, to vast clusters containing thousands of galaxies.

The behavior of these systems is determined by Newton's laws of motion and Newton's law of gravity,¹ and the study of this behavior is the branch of theoretical physics called **stellar dynamics**. Stellar dynamics is directly related to at least three other areas of theoretical physics. Superficially, it is closest to celestial mechanics, the theory of planetary motions—both involve the study of orbits in a gravitational field—however, much of the formalism of celestial mechanics is of little use in stellar dynamics, since it is based on perturbation expansions that do not converge when applied to most stellar systems. The most fundamental connections of stellar dynamics are with classical statistical mechanics, since the number of stars in a star cluster or galaxy is often so large that a statistical treatment of the dynamics is necessary. Finally, many of the mathematical tools that have been developed to study stellar systems are borrowed from plasma physics, which also involves the study of large numbers of particles interacting via long-range forces.

¹As yet, there is no direct evidence for stellar systems in which relativistic effects are important, although such systems are likely to be present at the centers of galaxies.

For an initial orientation, it is useful to summarize a few orders of magnitude for a typical stellar system, the one to which we belong. Our Sun is located in a stellar system called the **Milky Way** or simply **the Galaxy**. The Galaxy contains four principal constituents:

- (i) There are about 10^{11} stars, having a total mass $\simeq 5 \times 10^{10}$ solar masses (written $5 \times 10^{10} \mathcal{M}_{\odot}$; $1 \mathcal{M}_{\odot} = 1.99 \times 10^{30}$ kg).² Most of the stars in the Galaxy travel on nearly circular orbits in a thin disk whose radius is roughly 10^4 parsecs ($1 \text{ parsec} \equiv 1 \text{ pc} \equiv 3.086 \times 10^{16}$ m), or 10 kiloparsecs (kpc). The thickness of the disk is roughly 0.5 kpc and the Sun is located near its midplane, about 8 kpc from the center.
- (ii) The disk also contains gas, mostly atomic and molecular hydrogen, concentrated into clouds with a wide range of masses and sizes, as well as small solid particles (“dust”), which render interstellar gas opaque at visible wavelengths over distances of several kpc. Most of the atomic hydrogen is neutral rather than ionized, and so is denoted **HI**. Together, the gas and dust are called the **interstellar medium** (ISM). The total ISM mass is only about 10% of the mass in stars, so the ISM has little direct influence on the dynamics of the Galaxy. However, it plays a central role in the chemistry of galaxies, since dense gas clouds are the sites of star formation, while dying stars eject chemically enriched material back into the interstellar gas. The nuclei of the atoms in our bodies were assembled in stars that were widely distributed through the Galaxy.
- (iii) At the center of the disk is a black hole, of mass $\simeq 4 \times 10^6 \mathcal{M}_{\odot}$. The black hole is sometimes called Sagittarius A* or Sgr A*, after the radio source that is believed to mark its position, which in turn is named after the constellation in which it is found.
- (iv) By far the largest component, both in size and mass, is the **dark halo**, which has a radius of about 200 kpc and a mass of about $10^{12} \mathcal{M}_{\odot}$ (both these values are quite uncertain). The dark halo is probably composed of some weakly interacting elementary particle that has yet to be detected in the laboratory. For most purposes, the halo interacts with the other components of the Galaxy only through the gravitational force that it exerts, and hence stellar dynamics is one of the few tools we have to study this mysterious yet crucial constituent of the universe.

The typical speed of a star on a circular orbit in the disk is about 200 km s^{-1} . It is worth remembering that 1 km s^{-1} is almost exactly 1 pc (actually 1.023) in 1 megayear ($1 \text{ megayear} \equiv 1 \text{ Myr} = 10^6 \text{ years}$). Thus the time required to complete one orbit at the solar radius of 8 kpc is 250 Myr. Since the age of the Galaxy is about 10 gigayears ($1 \text{ gigayear} \equiv 1 \text{ Gyr} = 10^9 \text{ yr}$), most disk stars have completed over forty revolutions, and it is reasonable to assume that the Galaxy is now in an approximately steady state. The steady-state

² See Appendix A for a tabulation of physical and astronomical constants, and Tables 1.1, 1.2 and 2.3 for more precise descriptions of the properties of the Galaxy.

approximation allows us to decouple the questions of the present-day *equilibrium* and *structure* of the Galaxy, to which most of this book is devoted, from the thornier issue of the *formation* of the Galaxy, which we discuss only in the last chapter of this book.

Since the orbital period of stars near the Sun is several million times longer than the history of accurate astronomical observations, we are forced to base our investigation of Galactic structure on what amounts to an instantaneous snapshot of the system. To a limited extent, the snapshot can be supplemented by measurements of the angular velocities (or **proper motions**) of stars that are so close that their position on the sky has changed noticeably over the last few years; and by **line-of-sight velocities** of stars, measured from Doppler shifts in their spectra. Thus the *positions* and *velocities* of some stars can be determined, but their *accelerations* are almost always undetectable with current observational techniques.

Using the rough values for the dimensions of the Galaxy given above, we can estimate the mean free path of a star between collisions with another star. For an assembly of particles moving on straight-line orbits, the mean free path is $\lambda = 1/(n\sigma)$, where n is the number density and σ is the cross-section. Let us make the crude assumption that all stars are like the Sun so the cross-section for collision is $\sigma = \pi(2R_{\odot})^2$, where $R_{\odot} = 6.96 \times 10^8 \text{ m} = 2.26 \times 10^{-8} \text{ pc}$ is the solar radius.³ If we spread 10^{11} stars uniformly over a disk of radius 10 kpc and thickness 0.5 kpc, then the number density of stars in the disk is 0.6 pc^{-3} and the mean free path is $\lambda \simeq 2 \times 10^{14} \text{ pc}$. The interval between collisions is approximately λ/v , where v is the random velocity of stars at a given location. Near the Sun, the random velocities of stars are typically about 50 km s^{-1} . With this velocity, the collision interval is about $5 \times 10^{18} \text{ yr}$, over 10^8 times longer than the age of the Galaxy. Evidently, near the Sun collisions between stars are so rare that they are irrelevant—which is fortunate, since the passage of a star within even 10^3 solar radii would have disastrous consequences for life on Earth. For similar reasons, hydrodynamic interactions between the stars and the interstellar gas have a negligible effect on stellar orbits.

Thus, each star's motion is determined solely by the gravitational attraction of the mass in the galaxy—other stars, gas, and dark matter. Since the motions of weakly interacting dark-matter particles are also determined by gravitational forces alone, the tools that we develop in this book are equally applicable to both stars and dark matter, despite the difference of 70 or more orders of magnitude in mass.

We show in §1.2 that a useful first approximation for the gravitational field in a galaxy is obtained by imagining that the mass is continuously distributed, rather than concentrated into discrete mass points (the stars

³This calculation neglects the enhancement in collision cross-section due to the mutual gravitational attraction of the passing stars, but this increases the collision rate by a factor of less than 100, and hence does not affect our conclusion. See equation (7.195).

and dark-matter particles) and clouds (the gas). Thus we begin Chapter 2 with a description of Newtonian potential theory, developing methods to describe the smoothed gravitational fields of stellar systems having a variety of shapes. In Chapter 3 we develop both quantitative and qualitative tools to describe the behavior of particle orbits in gravitational fields. In Chapter 4 we study the statistical mechanics of large numbers of orbiting particles to find equilibrium distributions of stars in phase space that match the observed properties of galaxies, and learn how to use observations of galaxies to infer the properties of the underlying gravitational field.

The models constructed in Chapter 4 are **stationary**, that is, the density at each point is constant in time because the rates of arrival and departure of stars in every volume element balance exactly. Stationary models are appropriate to describe a galaxy that is many revolutions old and hence presumably in a steady state. However, some stationary systems are unstable, in that the smallest perturbation causes the system to evolve to some quite different configuration. Such systems cannot be found in nature. Chapter 5 studies the stability of stellar systems.

In Chapter 6 we describe some of the complex phenomena that are peculiar to galactic disks. These include the beautiful spiral patterns that are usually seen in disk galaxies; the prominent bar-like structures seen at the centers of about half of all disks; and the warps that are present in many spiral galaxies, including our own.

Even though stellar collisions are extremely rare, the gravitational fields of passing stars exert a series of small tugs that slowly randomize the orbits of stars. Gravitational encounters of this kind in a stellar system are analogous to collisions of molecules in a gas or Brownian motion of small particles in a fluid—all these processes drive the system towards energy equipartition and a thermally relaxed state. Relaxation by gravitational encounters operates so slowly that it can generally be neglected in galaxies, except very close to their centers (see §1.2); however, this process plays a central role in determining the evolution and present form of many star clusters. Chapter 7 describes the kinetic theory of stellar systems, that is, the study of the evolution of stellar systems towards thermodynamic equilibrium as a result of gravitational encounters. The results can be directly applied to observations of star clusters in our Galaxy, and also have implications for the evolution of clusters of galaxies and the centers of galaxies.

Chapter 8 is devoted to the interplay between stellar systems. We describe the physics of collisions and mergers of galaxies, and the influence of the surrounding galaxy on the evolution of smaller stellar systems orbiting within it, through such processes as dynamical friction, tidal stripping, and shock heating. We also study the effect of irregularities in the galactic gravitational field—generated, for example, by gas clouds or spiral arms—on the orbits of disk stars.

Throughout much of the twentieth century, galaxies were regarded as “island universes”—distinct stellar systems occupying secluded positions in

space. Explicitly or implicitly, they were seen as isolated, permanent structures, each a dynamical and chemical *unit* that was formed in the distant past and did not interact with its neighbors. A major conceptual revolution in **extragalactic astronomy**—the study of the universe beyond the edges of our own Galaxy—was the recognition in the 1970s that this view is incorrect. We now believe in a model of **hierarchical galaxy formation**, the main features of which are that: (i) encounters and mergers of galaxies play a central role in their evolution, and in fact galaxies are formed by the mergers of smaller galaxies; (ii) even apparently isolated galaxies are surrounded by much larger dark halos whose outermost tendrils are linked to the halos of neighboring galaxies; (iii) gas, stars, and dark matter are being accreted onto galaxies up to the present time. A summary of the modern view of galaxy formation and its cosmological context is in Chapter 9.

1.1 An overview of the observations

1.1.1 Stars

The luminosity of the Sun is $L_{\odot} = 3.84 \times 10^{26}$ W. More precisely, this is the **bolometric luminosity**, the total rate of energy output integrated over all wavelengths. The bolometric luminosity is difficult to determine accurately, in part because the Earth’s atmosphere is opaque at most wavelengths. Hence astronomical luminosities are usually measured in one or more specified wavelength bands, such as the **blue** or *B* band centered on $\lambda = 450$ nm; the **visual** or *V* band at $\lambda = 550$ nm; the *R* band at $\lambda = 660$ nm; the near-infrared *I* band at $\lambda = 810$ nm; and the infrared *K* band centered on the relatively transparent atmospheric window at $\lambda = 2200$ nm = $2.2 \mu\text{m}$, all with width $\Delta\lambda/\lambda \simeq 0.2$ (see Binney & Merrifield 1998, §2.3; hereafter this book is abbreviated as BM). For example, the brightest star in the sky, Sirius, has luminosities

$$L_V = 22 L_{\odot V} \quad ; \quad L_R = 15 L_{\odot R}, \quad (1.1)$$

while the nearest star, Proxima Centauri, has luminosities

$$L_V = 5.2 \times 10^{-5} L_{\odot V} \quad ; \quad L_R = 1.7 \times 10^{-4} L_{\odot R}. \quad (1.2)$$

This notation is usually simplified by dropping the subscript from L_{\odot} , when the band to which it refers is clear from the context.

Luminosities are often expressed in a logarithmic scale, by defining the **absolute magnitude**

$$M \equiv -2.5 \log_{10} L + \text{constant}. \quad (1.3)$$

The constant is chosen separately and arbitrarily for each wavelength band. The solar absolute magnitude is

$$M_{\odot B} = 5.48 \quad ; \quad M_{\odot V} = 4.83 \quad ; \quad M_{\odot R} = 4.42. \quad (1.4)$$

Sirius has absolute magnitude $M_V = 1.46$, $M_R = 1.47$, and Proxima Centauri has $M_V = 15.5$, $M_R = 13.9$. The **flux** from a star of luminosity L at distance d is $f = L/(4\pi d^2)$, and a logarithmic measure of the flux is provided by the **apparent magnitude**

$$m \equiv M + 5 \log_{10}(d/10 \text{ pc}) = -2.5 \log_{10} [L(10 \text{ pc}/d)^2] + \text{constant}; \quad (1.5)$$

thus, the absolute magnitude is the apparent magnitude that the star would have if it were at a distance of 10 parsecs. Note that *faint* stars have *large* magnitudes. Sirius is at a distance⁴ of (2.64 ± 0.01) pc and has apparent magnitude $m_V = -1.43$, while Proxima Centauri is at (1.295 ± 0.004) pc and has apparent magnitude $m_V = 11.1$. The faintest stars visible to the naked eye have $m_V \simeq 6$, and the limiting magnitude of the deepest astronomical images at this time is $m_V \simeq 29$. The apparent magnitudes m_V and m_R are often abbreviated simply as V and R .

The **distance modulus** $m - M = 5 \log_{10}(d/10 \text{ pc})$ is often used as a measure of distance.

The **color** of a star is measured by the ratio of the luminosity in two wavelength bands, for example by L_R/L_V or equivalently by $M_V - M_R = m_V - m_R = V - R$. Sirius has color $V - R = -0.01$ and Proxima Centauri has $V - R = 1.67$. Stellar spectra are approximately black-body and hence the color is a measure of the temperature at the surface of the star.

A more precise measure of the surface temperature is the **effective temperature** T_{eff} , defined as the temperature of the black body with the same radius and bolometric luminosity as the star in question. If the stellar radius is R , then the Stefan–Boltzmann law implies that the bolometric luminosity is

$$L = 4\pi R^2 \sigma T_{\text{eff}}^4, \quad (1.6)$$

where $\sigma = 5.670 \times 10^{-8} \text{ W m}^{-2} \text{ K}^{-4}$. The relation between color and effective temperature is tabulated in BM §3.4 and shown in Figure 1.1.

A third measure of the surface temperature of a star is its **spectral class**, which is assigned on the basis of the prominence of various absorption lines in the stellar spectrum. In order of decreasing temperature, the spectral classes are labeled O, B, A, F, G, K, M, L, and T, and each class is divided into ten subclasses by the numbers 0, 1, ..., 9. Thus a B0 star is slightly

⁴Throughout this book, quoted errors are all 1 standard deviation, i.e., there is a probability of 0.68 that the actual error is less than the quoted error. The reader is warned that astronomical errors are often dominated by unknown systematic effects and thus tend to be underestimated.

cooler than an O9 star. Using this scheme, experienced observers can determine the effective temperature of a star to within about 10% from a quick examination of its spectrum. For example, Sirius has spectral class A1 and effective temperature 9500 K, while Proxima Centauri has spectral class M5 and $T_{\text{eff}} = 3000$ K. The Sun is a quite ordinary G2 star, with $T_{\text{eff}} = 5780$ K.

The ultraviolet emission from the hottest stars ionizes nearby interstellar gas, forming a sphere of ionized hydrogen called an **HII region** (BM §8.1.3). The brightest star-like objects in other galaxies are often HII regions shining in emission lines, rather than normal stars shining from thermal emission (see, for example, Plate 1).

The **color-magnitude** diagram is a plot of absolute magnitude against color; since color is related to effective temperature, each point on this diagram corresponds to a unique luminosity, effective temperature, and stellar radius (through eq. 1.6). In older work spectral type sometimes replaces color, since the two quantities are closely related, and in this case the plot is called a **Hertzsprung–Russell** or **HR** diagram (BM §3.5). This simple diagram has proved to be the primary point of contact between observations and the theory of stellar structure and evolution.

The distribution of stars in the color-magnitude diagram depends on the age and chemical composition of the sample of stars plotted. Astronomers refer to all elements beyond helium in the periodic table as “metals”; with the exception of lithium, such elements are believed to be formed in stars rather than at the birth of the universe (see §1.3.5). To a first approximation the abundances of groups of elements vary in lockstep since they are formed in the same reaction chain and injected into interstellar space by the same type of star. At an even cruder level the chemical composition of a star can be approximately specified by a single number Z , the **metallicity**, which is the fraction by mass of all elements heavier than helium. Similarly, the fractions by mass of hydrogen and helium are denoted X and Y ($X + Y + Z = 1$). The Sun’s initial composition was $X_{\odot} = 0.71$, $Y_{\odot} = 0.27$, $Z_{\odot} = 0.019$.

Figure 1.1 shows the color-magnitude diagram for about 10^4 nearby stars. The most prominent feature is the well-defined band stretching from $(B - V, M_V) \simeq (0, 0)$ to $(B - V, M_V) \simeq (1.5, 11)$. This band, known as the **main sequence**, contains stars that are burning hydrogen in their cores. In this stage of a star’s life, the mass—and to a lesser extent, chemical composition—uniquely determine both the effective temperature and the luminosity, so stars remain in a fixed position on the color-magnitude diagram. Main-sequence stars are sometimes called **dwarf** stars, to distinguish them from the larger giant stars that we discuss below. The main sequence is a mass sequence, with more massive stars at the upper left (high luminosity, high temperature, blue color) and less massive stars at the lower right (low luminosity, low temperature, red color). The most and least luminous main-sequence stars in this figure, with absolute magnitudes $M \simeq -2$ and $+12$, have masses of about $10 \mathcal{M}_{\odot}$ and $0.2 \mathcal{M}_{\odot}$, respectively (BM Table 3.13). Objects smaller than about $0.08 \mathcal{M}_{\odot}$ never ignite hydrogen in their cores, and

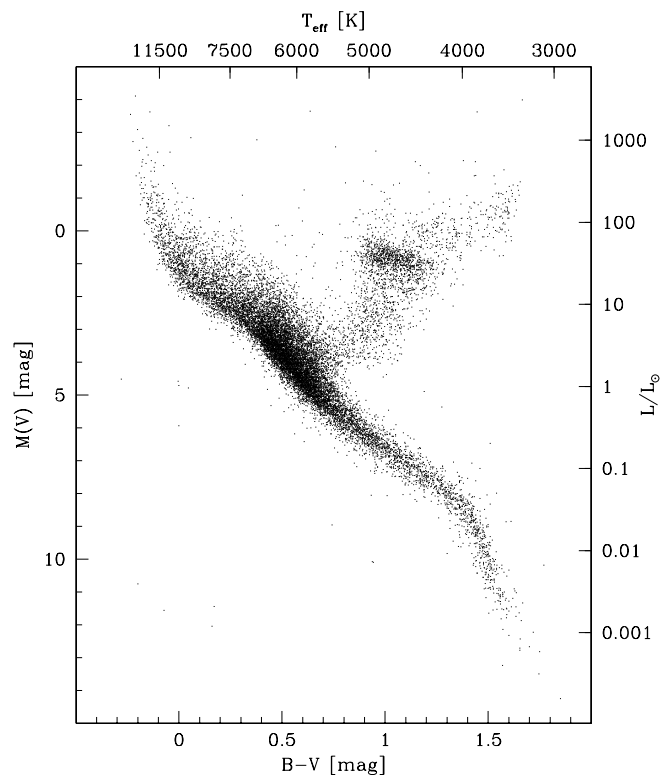


Figure 1.1 The color-magnitude diagram for over 10^4 nearby stars. Close binary stars have been excluded, since the companion contaminates the color and magnitude measurements. Most of the stars fall on the main sequence, which runs from upper left to lower right. The red-giant branch runs upward and to the right from the main sequence. The red-clump stars form a prominent concentration in the middle of the red-giant branch. The giants are chosen from a larger volume than the other stars to enhance their numbers and thereby show the structure of the giant branch more clearly. The absolute magnitudes are in the V band and the colors are based on $B - V$, which gives the ratio of fluxes between $\lambda = 450$ nm and 550 nm (BM Table 2.1). The right axis shows the V -band luminosity in solar units, and the top axis shows approximate values of the effective temperature, which follows from the color and luminosity because the stars are approximate black bodies. From Perryman et al. (1995).

are visible mainly from the radiation they emit as they contract and cool. The luminosity of these objects, known as **brown dwarfs**, therefore depends on both mass and age and they do not form a one-parameter sequence like their more massive siblings.

The massive, high-luminosity stars at the upper end of the main sequence exhaust their fuel rapidly and hence are short-lived ($\lesssim 100$ Myr for stars with $M_V = -2$), while stars at the lower end of the main sequence

burn steadily for much longer than the current age of the universe. The Sun has a lifetime of 10 Gyr on the main sequence.

From equation (1.6), stars that are luminous and cool (upper right of the color-magnitude diagram) must be large, while stars that are dim and hot (lower left of the diagram) must be small. Since the main sequence crosses from upper left to lower right, this argument suggests that radius is not a strong function of luminosity along the main sequence. The mass-radius relation in Table 3.13 of BM bears this out: between $M \simeq 0$ and $M \simeq 10$, a factor of 10^4 in luminosity, the radius varies by only a factor of six, from $3 R_\odot$ to $0.5 R_\odot$.

The color-magnitude diagram contains a handful of dim blue stars around $(B - V, M) \simeq (0, 12)$. These are **white dwarfs**, stars that have exhausted their nuclear fuel and are gradually cooling to invisibility. As their location in the diagram suggests, white dwarfs are very small, with radii of order $10^{-2} R_\odot$. White dwarfs are so dense that the electron gas in the interior of the star is degenerate; in other words, gravitational contraction is resisted, not by thermal pressure as in main-sequence stars, but rather by the Fermi energy of the star's cold, degenerate electron gas.

Figure 1.1 also contains a prominent branch slanting up and to the right from the main sequence, from $(B - V, M) \simeq (0.3, 4)$ to $(B - V, M) \simeq (1.5, -1)$. These are **red giants**, stars that have exhausted hydrogen in their cores and are now burning hydrogen in a shell surrounding an inert helium core. As their location in the color-magnitude diagram suggests—red therefore cool, yet very luminous—red giants are much larger than main-sequence stars; the stars at the tip of the giant branch have radii $\gtrsim 100 R_\odot$. In contrast to the main sequence, on which stars remain in a fixed position determined by their mass, the red-giant branch is an evolutionary sequence: stars climb the giant branch from the main sequence to its tip, over an interval of about 1 Gyr for stars like the Sun.

The prominent concentration in the middle of the red-giant branch, near $(B - V, M) \simeq (1, 1)$, is called the **red clump**. This feature arises from a later evolutionary stage, which happens to coincide with the red-giant branch for stars of solar metallicity. Red-clump stars have already ascended the giant branch to its tip and returned, settling at the red clump when they begin burning helium in their cores (see below).

Red giants are rare compared to main-sequence stars because the red-giant phase in a star's life is much shorter than its main-sequence phase. Nevertheless, red giants are so luminous that they dominate the total luminosity of many stellar systems. Another consequence of their high luminosity is that a far larger fraction of red giants is found in flux-limited samples than in volume-limited samples. For example, over half of the 100 brightest stars are giants, but none of the 100 nearest stars is a giant.

Figure 1.2 illustrates the color-magnitude diagram of a typical globular star cluster (§1.1.4). The advantage of studying a star cluster is that all of its members lie at almost the same distance, and have the same age and

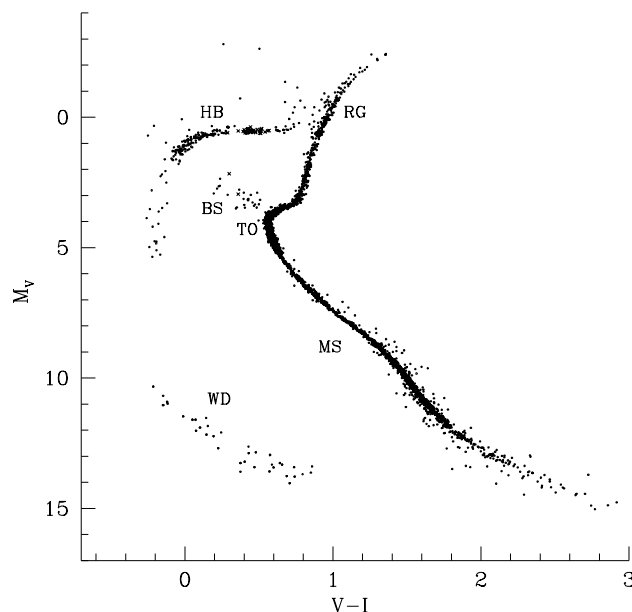


Figure 1.2 The color-magnitude diagram for metal-poor globular clusters. The horizontal axis is the color $V - I$. Labels denote the main sequence (MS), the main-sequence turnoff (TO), red giants (RG), horizontal branch (HB), blue stragglers (BS), and white dwarfs (WD). This is a composite diagram in which data from five globular clusters have been combined selectively to emphasize the principal sequences; thus the relative numbers of different types of stars are not realistic. From data supplied by W. E. Harris; see also Harris (2003).

chemical composition. Thus age and composition differences and distance errors do not blur the diagram, so the sequences are much sharper than in Figure 1.1 (BM §6.1.2). In this diagram the main sequence stretches from $(V - I, M_V) \simeq (0.6, 4)$ to $(V - I, M_V) \simeq (2.4, 14)$. In contrast to Figure 1.1, the main sequence terminates sharply at $M_V \simeq 4$ (the **turnoff**); the more luminous, bluer part of the main sequence that is seen in the sample of nearby stars is absent in the cluster, because such stars have lifetimes shorter than the cluster age. Figure 1.2 shows a few stars situated along the extrapolation of the main sequence past the turnoff point; these “blue stragglers” may arise from collisions and mergers of stars in the dense core of the cluster or mass transfer between the components of a binary star (page 628). The white dwarfs are visible near $(V - I, M_V) \simeq (0.4, 13)$, and the tip of the red-giant branch lies at $(V - I, M_V) \simeq (1.3, -2)$.

As a star evolves, it climbs the giant branch until, at the tip of the giant branch, helium starts to burn in its core. The stars then evolve rapidly

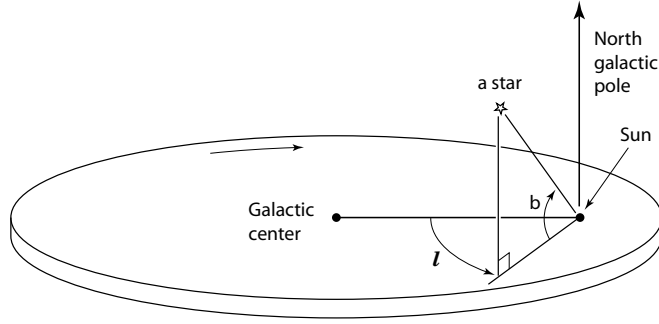


Figure 1.3 A schematic picture of the Sun’s location in the Galaxy, illustrating the Galactic coordinate system. An arrow points in the direction of Galactic rotation, which is clockwise as viewed from the north Galactic pole.

to the **horizontal branch** (the sequence of stars near $M_V \simeq 0$, stretching from $V - I \simeq 0$ to 0.8), where they remain until the helium in the core is exhausted. The form of the horizontal branch depends on the metallicity: as the metallicity increases from very low values ($Z \simeq 0.01Z_\odot$) the horizontal branch shortens and moves to the right, until at near-solar metallicity it is truncated to the red clump seen in Figure 1.1.

Stars in metal-poor globular clusters are among the oldest objects in the Galaxy. Fits of theoretical models of stellar evolution to the color-magnitude diagrams of metal-poor globular clusters yield ages of (12.5 ± 1.5) Gyr (Krauss & Chaboyer 2003). This result is consistent with the age of the universe determined from measurements of the cosmic background radiation, $t_0 = (13.7 \pm 0.2)$ Gyr (eq. 1.77), if the globular clusters formed when the universe was about 1 Gyr old.

1.1.2 The Galaxy

Most of the stars in the Galaxy lie in a flattened, roughly axisymmetric structure known as the **Galactic disk**. On clear, dark nights the cumulative light from the myriad of faint disk stars is visible as a luminous band stretching across the sky, which is the source of the name “Milky Way” for our Galaxy. The midplane of this disk is called the **Galactic plane** and serves as the equator of **Galactic coordinates** (ℓ, b) , where ℓ is the **Galactic longitude** and b is the **Galactic latitude** (BM §2.1.2). The Galactic coordinate system is a heliocentric system in which $\ell = 0, b = 0$ points to the Galactic center and $b = \pm 90^\circ$ points to the **Galactic poles**, normal to the disk plane (see Figure 1.3).

The Sun is located at a distance R_0 from the center of the Galaxy; the best current estimate $R_0 = (8.0 \pm 0.5)$ kpc comes from the orbits of stars near the black hole that is believed to mark the center (see §1.1.6 and Eisenhauer

et al. 2003). One measure of the distribution of stars in the Galactic disk is the surface brightness, the total stellar luminosity emitted per unit area of the disk (see Box 2.1 for a more precise definition). Observations of other disk galaxies suggest that the surface brightness is approximately an exponential function of radius,

$$I(R) = I_d \exp(-R/R_d). \quad (1.7)$$

The **disk scale length** R_d is difficult to measure in our Galaxy because of our position within the disk. Current estimates place R_d between about 2 and 3 kpc. Thus the Sun lies farther from the Galactic center than about 75–90% of the disk stars. The resulting concentration of luminosity towards the Galactic center is not apparent to the naked eye, since interstellar dust absorbs the light from distant disk stars (the optical depth in the V band along a line of sight in the Galactic midplane is unity at a distance of only about 0.7 kpc); however, in the infrared, where dust extinction is unimportant, our position near the edge of the disk is immediately apparent from the strong concentration of light in the direction of the constellation Sagittarius (at the center of the image in Plate 2). By contrast, the Galaxy is nearly transparent in the direction of the Galactic poles, which greatly facilitates studying the extragalactic universe.

The stars of the disk travel in nearly circular orbits around the Galactic center. The speed of a star in a circular orbit of radius R in the Galactic equator is denoted $v_c(R)$, and a plot of $v_c(R)$ versus R is called the **circular-speed curve**. The circular speed at the solar radius R_0 is

$$v_0 \equiv v_c(R_0) = (220 \pm 20) \text{ km s}^{-1}. \quad (1.8)$$

A strong additional constraint on v_0 and R_0 comes from the angular motion of the radio source Sgr A* relative to extragalactic sources: if Sgr A* coincides with the black hole at the Galactic center, and if this black hole is at rest in the Galaxy—both very plausible assumptions, but not certainties—then the angular speed of the Sun is $v_0/R_0 = (236 \pm 1) \text{ km s}^{-1}/(8 \text{ kpc})$ (Reid & Brunthaler 2004).

The **Local Standard of Rest** (LSR) is an inertial reference frame centered on the Sun and traveling at speed v_0 in the direction of Galactic rotation. Since most nearby disk stars are on nearly circular orbits, their velocities relative to the Local Standard of Rest are much smaller than v_0 . For example, the Sun's velocity relative to the LSR (the **solar motion**) is (BM §10.3.1)

$$13.4 \text{ km s}^{-1} \text{ in the direction } \ell = 28^\circ, b = 32^\circ. \quad (1.9)$$

The root-mean-square (RMS) velocity of old disk stars relative to the Local Standard of Rest is 50 km s^{-1} , larger than the Sun's velocity but still small compared to the circular speed v_0 .

In the direction perpendicular to the Galactic plane (usually called the “vertical” direction), the density of stars falls off exponentially,

$$\rho(R, z) = \rho(R, 0)e^{-|z|/z_d(R)}, \quad (1.10)$$

where z is the distance from the midplane and $z_d(R)$ is the **scale height** at radius R .⁵ The thickness z_d of the Galactic disk depends on the age of the stars that are being examined. Older stellar populations have larger scale heights, probably because stochastic gravitational fields due to spiral arms and molecular clouds gradually pump up the random velocities of stars (see §8.4). In the solar neighborhood, the scale height ranges from $\lesssim 100$ pc for the young O and B stars to $\simeq 300$ pc for the stars with ages of order 10 Gyr that constitute the bulk of the disk mass.

A more accurate representation of the vertical structure of the disk is obtained by superimposing two populations with densities described by equation (1.10): the **thin disk** with $z_d \simeq 300$ pc, and the **thick disk** with $z_d \simeq 1$ kpc (BM Figure 10.25). The stars of the thick disk are older and have a different chemical composition from those of the thin disk—thick-disk stars have lower metallicities, and at a given metallicity they have higher abundances of the α nuclides (^{16}O , ^{20}Ne , ^{24}Mg , ^{28}Si , etc.; see BM §5.2.1 and Figure 10.17) relative to ^{56}Fe . The surface density of the thick disk is about 7% of that of the thin disk, so in the midplane, thin-disk stars outnumber thick-disk stars by about 50:1. The thick disk was probably created when the infant thin disk was shaken and thickened by an encounter with a smaller galaxy early in its history.

The enhanced α nuclides found in the thick disk are the signature of stars formed early in the history of the disk, for the following reason. The interstellar gas is polluted with heavy elements by two main processes: (i) “core-collapse” supernovae, arising from the catastrophic gravitational collapse of massive stars, which lag star formation by no more than ~ 40 Myr, and produce ejecta that are rich in α nuclides; (ii) “thermonuclear” or Type Ia supernovae, which are caused by runaway nuclear burning on the surface of white-dwarf stars in binary systems, lag star formation by of order 0.5–10 Gyr, and produce mostly nuclei near ^{56}Fe . Thus the chemical composition of thick-disk stars suggests that the thick disk formed in less than about 1 Gyr. In contrast, it appears that stars in the thin disk have formed at a steady rate throughout the lifetime of the Galaxy.

Throughout this book, we shall distinguish the **kinematics** of a stellar system—the observational description of the positions and motions of the stars in the system—from its dynamics—the interpretation of these motions in terms of physical laws (forces, masses, etc.). Thus, the description of the

⁵ This formula has a discontinuous slope at $z = 0$, which reflects the gravitational attraction of the much thinner gas layer on the stars. The vertical distribution of stars in a thin disk is explored theoretically in Problem 4.22.

Galaxy in this subsection has so far been kinematic. The simplest approximate dynamical description of the Galaxy is obtained by assuming that its mass distribution is spherical. Let the mass interior to radius r be $M(r)$. From Newton’s theorems (§2.2.1) the gravitational acceleration at radius r is equal to that of a point whose mass is the same as the total mass interior to r ; thus the inward acceleration is $GM(r)/r^2$, where the **gravitational constant** $G = 6.674 \times 10^{-11} \text{ m}^3 \text{ kg}^{-1} \text{ s}^{-2}$. The central or **centripetal** acceleration required to hold a body in a circular orbit with speed v_0 is v_0^2/r . Thus the mass interior to the solar radius R_0 in this crude model is

$$M(R_0) = \frac{v_0^2 R_0}{G} = 9.0 \times 10^{10} \mathcal{M}_\odot \left(\frac{v_0}{220 \text{ km s}^{-1}} \right)^2 \left(\frac{R_0}{8 \text{ kpc}} \right). \quad (1.11)$$

The approximation that the mass distribution is spherical is reasonable for the dark halo, but not for the flat stellar disk. Better models suggest that this estimate is probably high by about 30%, since a disk requires less mass to produce a given centripetal acceleration (see Figure 2.17).

Most of our understanding of stellar astrophysics comes from observations of stars within a few hundred parsecs of the Sun. This distance is much smaller than the disk scale length, and hence it is reasonable to assume that the distribution of properties of these stars (chemical compositions, ages, masses, kinematics, fraction of binary stars, etc.) is constant within this region, even though there may be large-scale gradients in these properties across the Galactic disk. To formalize this assumption, we define the **solar neighborhood** to be a volume centered on the Sun that is much smaller than the overall size of the Galaxy but large enough to contain a statistically useful sample of stars. The concept is somewhat imprecise but nevertheless extremely useful. The appropriate size of the volume depends on which stars we wish to investigate: for white dwarfs, which are both common and dim, the “solar neighborhood” may consist of a sphere of radius only 30 pc centered on the Sun, while for the luminous but rare O and B stars, the solar neighborhood may be considered to extend as far as 1–2 kpc from the Sun.

Our best estimate of the inventory of the solar neighborhood is summarized in Table 1.1. The category “visible stars” includes all main-sequence and giant stars. The category “stellar remnants” includes white dwarfs and neutron stars, while “ISM” (interstellar medium) includes atomic and molecular hydrogen, ionized gas, and a small contribution from interstellar dust. The volume density and luminosity density are quoted in the Galactic midplane and the surface density and surface brightness are integrated over a column perpendicular to the Galactic plane, extending to ± 1.1 kpc from the midplane. “Dynamical” denotes determinations of the total volume or surface density from the dynamics of disk stars (see §4.9.3). The dynamically determined volume density in the midplane is consistent with the observed density in stars and gas to within about 10%, so there is no evidence for a significant component of dark matter in the disk—in other words, the inventory in Table 1.1 appears to be complete. The dynamically determined

Table 1.1 Inventory of the solar neighborhood

component	volume density ($\mathcal{M}_\odot \text{pc}^{-3}$)	surface density ($\mathcal{M}_\odot \text{pc}^{-2}$)	luminosity density ($L_\odot \text{pc}^{-3}$)	surface brightness ($L_\odot \text{pc}^{-2}$)
visible stars	0.033	29	0.05	29
stellar remnants	0.006	5	0	0
brown dwarfs	0.002	2	0	0
ISM	0.050	13	0	0
total	0.09 ± 0.01	49 ± 6	0.05	29
dynamical	0.10 ± 0.01	74 ± 6	–	–

NOTES: Volume and luminosity densities are measured in the Galactic midplane and surface density is the total within ± 1.1 kpc of the plane. Luminosity density and surface brightness are given in the R band. Dynamical estimates are from §4.9.3. Most other entries are taken from Flynn et al. (2006).

surface density appears to be higher than the surface density in stars and gas, by $(25 \pm 9) \mathcal{M}_\odot \text{pc}^{-2}$; if significant, this excess probably represents the contribution of the dark halo. The dark halo also contributes to the volume density in the midplane, but this contribution is undetectably small.

A stellar system is often characterized by its **mass-to-light ratio**, which we denote by Υ and write in units of the solar ratio, $\Upsilon_\odot = \mathcal{M}_\odot / L_\odot$. According to Table 1.1, the mass-to-light ratio of the solar neighborhood in the R band is $\Upsilon_R \simeq 2 \Upsilon_\odot$ in the midplane and $\simeq 2.5 \Upsilon_\odot$ after integrating to ± 1.1 kpc from the plane. The second value is higher because the scale height z_d of luminous young stars is smaller than that of older, dimmer stars.

In addition to the disk, the Galaxy contains a **bulge**, a small, amorphous, centrally located stellar system that is thicker than the disk and comprises $\sim 15\%$ of the total luminosity. The Galactic bulge is clearly visible at the center of the disk in infrared images of the Galaxy (Plate 2). The evolutionary history, kinematics, and chemical composition of bulge stars are quite different from those of disk stars near the Sun. The bulge stars are believed to date from near the time of formation of the Galaxy, whereas the disk stars have a wide range of ages, since star formation in the disk is an ongoing process. While disk stars in the solar neighborhood are found in nearly circular orbits with speeds $v_c(R) \simeq 220 \text{ km s}^{-1}$ and RMS velocity relative to this speed of only 50 km s^{-1} , the velocity vectors of bulge stars are randomly oriented, with RMS velocity $\simeq 150 \text{ km s}^{-1}$. The bulge stars exhibit a wide range of metallicities, spread around a median metallicity of about $0.4Z_\odot$ (Zoccali et al. 2003), substantially smaller than the metallicity of young stars in the solar neighborhood—presumably because the interstellar gas from which the local disk stars form has steadily become more and more polluted by the metal-rich debris of exploding supernovae.

By analogy to the statistical-mechanical concept that temperature is proportional to mean-square velocity, a stellar population like the disk in

which the random velocities are much smaller than the ordered or mean velocity is said to be **cool**, while the bulge population, in which the random velocities are larger than the mean velocity, is said to be **hot**. A hypothetical disk in which the stars move on precisely circular orbits would be **cold**.

Although the distribution of bulge stars is symmetric about the Galactic midplane, the bulge is somewhat brighter and thicker on one side of the Galactic center (longitude $\ell > 0$) than on the other. This asymmetry arises because the bulge is triaxial: the lengths of the two principal axes that lie in the Galactic plane are in the ratio 3:1, and the triaxial structure extends to about 3 kpc from the center. The long axis is oriented about 20° from the line between the Galactic center and the Sun (§2.7e). Thus the bulge is brighter and thicker at positive longitudes simply because that side is closer to the Sun. Because the bulge is triaxial it is also sometimes called a “bar” and the Milky Way is said to be a barred galaxy (see §1.1.3).

About 1% of the stellar mass in the Galaxy is contained in the **stellar halo**, which contains old stars of low metallicity (median about $0.02Z_\odot$). The stellar halo has little or no mean rotation, and a density distribution that is approximately spherical and a power-law function of radius, $\rho \propto r^{-3}$, out to at least 50 kpc. The metal-poor globular clusters that we describe below (§1.1.4) are members of the stellar halo. The low metallicity of this population suggests that it was among the first components of the Galaxy to form. Much of the halo comprises the debris of disrupted stellar systems, such as globular clusters and small satellite galaxies.

The dark halo is the least well understood of the Galaxy’s components. We have only weak constraints on its composition, shape, size, mass, and local density. A wide variety of candidates for the dark matter have been suggested, most falling into one of two broad classes: (i) some unknown elementary particle—the preferred candidates are WIMPs, an acronym for weakly interacting massive particles, but there are also more exotic possibilities such as axions; (ii) non-luminous macroscopic objects, such as neutron stars or black holes, which are usually called MACHOs, for massive compact halo objects. Measurements of the optical depth to gravitational lensing through the halo exclude MACHOs in the range 10^{-7} – $30 M_\odot$ as the dominant component of the dark halo (Alcock et al. 2001; Tisserand et al. 2007), and indirect dynamical arguments (§8.2.2e) suggest that more massive compact objects are also excluded. On the other hand, hypothetical massive, neutral, weakly interacting particles could be formed naturally in the early universe in approximately the numbers required to make a substantial contribution to the overall density. Thus most physicists and astronomers believe that the dark halo is probably composed of WIMPs. Ordinary matter—stars, dust, interstellar gas, MACHOs, etc., whether luminous or dark—derives almost all of its mass from baryons and hence is usually referred to as **baryonic matter** to distinguish it from **non-baryonic matter** such as WIMPs.⁶

⁶ A baryon is a strongly interacting fermion. The word derives from *barus*, the Greek

The formation of flat astrophysical systems such as the solar system or a galaxy disk requires dissipation, which removes energy but conserves angular momentum and therefore leads naturally to a rapidly rotating thin disk. Since WIMPs cannot dissipate energy, the dark halo is expected to be approximately spherical. Numerical simulations of the formation of dark halos suggest that they are triaxial rather than precisely spherical, with minor-to-major axis ratios of 0.4–0.6, but there is little direct observational evidence on halo shapes (§9.3.3).

The total size and mass of the Galaxy’s halo can be constrained by the kinematics of distant globular clusters and nearby galaxies. Using this method Wilkinson & Evans (1999) find a best-fit mass of $2 \times 10^{12} \mathcal{M}_\odot$, with a **median** or **half-mass radius** (the radius containing half the total mass) of 100 kpc; however, these values are very uncertain and masses as small as $2 \times 10^{11} \mathcal{M}_\odot$ or as large as $5 \times 10^{12} \mathcal{M}_\odot$ are allowed. A reasonable guess of the total mass of the Galaxy inside 100 kpc radius is

$$M(r < 100 \text{ kpc}) = 5\text{--}10 \times 10^{11} \mathcal{M}_\odot. \quad (1.12)$$

The mass distribution in the dark halo is equally uncertain at smaller radii: the halo contribution to the radial gravitational force at the solar radius, which determines the circular speed, could lie anywhere from less than 10% to almost 50% of the total force without violating the observational constraints (§6.3.3). The uncertain halo mass distribution inside R_0 implies an uncertain halo density at R_0 , which is a significant concern to experimentalists hoping to detect the dark matter in laboratory experiments (Gaitskell 2004).

It is useful to parametrize the relative amounts of dark and luminous matter in a stellar system by the mass-to-light ratio. Stellar systems composed entirely of stars usually have mass-to-light ratios Υ_R in the range $1\text{--}10\Upsilon_\odot$, depending on the age and chemical composition of the stars, while systems composed entirely of dark matter would have $\Upsilon \rightarrow \infty$. The largest mass-to-light ratios known, $\Upsilon_R \sim 500\Upsilon_\odot$, occur in dwarf spheroidal galaxies (page 24). In the R band, the luminosity of the Galaxy is $3 \times 10^{10} L_\odot$, so its mass-to-light ratio is $\sim 60\Upsilon_\odot$, with large uncertainties ($7\text{--}170\Upsilon_\odot$) because of the uncertain mass of the dark halo.

A summary of properties of the Galaxy is provided in Table 1.2; for more details see §2.7.

for heavy. The use of the term “baryonic matter” for ordinary matter is conventional, but less than ideal for several reasons: (i) ordinary matter includes electrons, which are leptons, not baryons; (ii) the unknown dark-matter particle is likely to be even heavier than any baryonic particle; (iii) it is not clear whether to count neutrinos as “ordinary” matter, because they have been known for decades, or as dark matter, because they have mass and interact only weakly so they are WIMPs. Usually the latter choice is made.

Table 1.2 Properties of the Galaxy

Global properties:	
disk scale length R_d	(2.5 ± 0.5) kpc
disk luminosity	$(2.5 \pm 1) \times 10^{10} L_\odot$
bulge luminosity	$(5 \pm 2) \times 10^9 L_\odot$
total luminosity	$(3.0 \pm 1) \times 10^{10} L_\odot$
disk mass	$(4.5 \pm 0.5) \times 10^{10} \mathcal{M}_\odot$
bulge mass	$(4.5 \pm 1.5) \times 10^9 \mathcal{M}_\odot$
dark halo mass	$(2^{+3}_{-1.8}) \times 10^{12} \mathcal{M}_\odot$
dark halo half-mass radius	(100^{+100}_{-80}) kpc
disk mass-to-light ratio Υ_R	$(1.8 \pm 0.7) \Upsilon_\odot$
total mass-to-light ratio Υ_R	$(70^{+100}_{-63}) \Upsilon_\odot$
black-hole mass	$(3.9 \pm 0.3) \times 10^6 \mathcal{M}_\odot$
Hubble type	Sbc
Solar neighborhood properties:	
solar radius R_0	(8.0 ± 0.5) kpc
circular speed v_0	(220 ± 20) km s ⁻¹
angular speed: from v_0/R_0	(27.5 ± 3) km s ⁻¹ kpc ⁻¹
from Sgr A*	(29.5 ± 0.2) km s ⁻¹ kpc ⁻¹
disk density ρ_0	$(0.09 \pm 0.01) \mathcal{M}_\odot \text{pc}^{-3}$
disk surface density Σ_0	$(49 \pm 6) \mathcal{M}_\odot \text{pc}^{-2}$
disk thickness Σ_0/ρ_0	500 pc
scale height z_d (old stars)	300 pc
rotation period $2\pi/\Omega_0$	(220 ± 30) Myr
vertical frequency $\nu_0 = \sqrt{4\pi G\rho_0}$	$(2.3 \pm 0.1) \times 10^{-15}$ Hz = (70 ± 4) km s ⁻¹ kpc ⁻¹
vertical period $2\pi/\nu_0$	87 Myr
Oort's A constant	(14.8 ± 0.8) km s ⁻¹ kpc ⁻¹
Oort's B constant	$-(12.4 \pm 0.6)$ km s ⁻¹ kpc ⁻¹
epicycle frequency $\kappa_0 = \sqrt{-4B(A-B)}$	(37 ± 3) km s ⁻¹ kpc ⁻¹
radial dispersion of old stars	(38 ± 2) km s ⁻¹
vertical dispersion of old stars	(19 ± 2) km s ⁻¹
RMS velocity of old stars	(50 ± 3) km s ⁻¹
escape speed $v_e(R_0)$	(550 ± 50) km s ⁻¹

NOTES: See §2.7 and §§10.1 and 10.3 of BM for more detail. Luminosities are in the R band at $\lambda = 660$ nm. The halo mass and half-mass radius are taken from Wilkinson & Evans (1999). The angular speed of the central black hole (Sgr A*) relative to an extragalactic frame is from Reid & Brunthaler (2004). The density in the midplane of the disk, ρ_0 , and the surface density Σ_0 are taken from Table 1.1. The scale height z_d is defined by equation (1.10). The RMS velocity of old stars is the square root of the sum of the squared dispersions along the three principal axes of the velocity-dispersion tensor (BM Table 10.2). Escape speed is from Smith et al. (2007).

1.1.3 Other galaxies

The nearest known galaxy to our own is the **Sagittarius dwarf galaxy**, (see Table 4.3 of BM or van den Bergh 2000 for a list of nearby galaxies). The Sagittarius galaxy has a total luminosity $L \simeq 2 \times 10^7 L_\odot$ and is located on the opposite side of the Galaxy from us, about 24 kpc from the Sun and 16 kpc from the Galactic center. The line of sight to the Sagittarius dwarf passes only 15° from the Galactic center, so the galaxy is masked by the dense star fields of the Galactic bulge, and thus was discovered only in 1994 (from anomalies in the kinematics of what were thought to be bulge stars). The orbit of Sagittarius carries it so close to the center of the Galaxy that it is being disrupted by the Galactic tidal field, and a tidal tail or streamer—a trail of stars torn away from the main body of the galaxy—can be traced across most of the sky (Figure 8.10).

Our next nearest neighbor is the **Large Magellanic Cloud** or LMC. Although some 50 times as luminous than Sagittarius, with luminosity $L_R \simeq 1 \times 10^9 L_\odot$, the LMC is still a relatively modest galaxy. The LMC is 45–50 kpc from the Sun and is visible to the naked eye in the southern hemisphere as a faint patch of light (see Plates 2 and 11). Because of its proximity and its location at relatively high Galactic latitude, where foreground contamination and dust obscuration are small ($b = -30^\circ$), the LMC provides a unique laboratory for studies of interstellar gas and dust, stellar properties, and the cosmological distance scale. Also visible to the naked eye is the **Small Magellanic Cloud**, located 20° from the LMC on the sky, 20% further away, and with 20% of its luminosity. It is likely that the two Clouds are a former binary system that has been disrupted by tidal forces from the Galaxy.

The nearest large disk galaxy similar to our own is called the **Andromeda galaxy**, **M31**, or **NGC 224** (see Plate 3 and Hodge 1992). M31 is more than ten times as far away as the LMC ($d \simeq 740$ kpc) and more than ten times as luminous ($L \simeq 4 \times 10^{10} L_\odot$). Only the central parts of M31 are visible to the naked eye, but deep telescopic images show that its stellar disk extends across more than six degrees on the sky.

Our Galaxy is just one member of a vast sea of some 10^9 galaxies stretching to a distance of several thousand megaparsecs (1 megaparsec \equiv 1 Mpc = 10^6 pc = 3.086×10^{22} m). The determination of distances to these galaxies is one of the most important tasks in extragalactic astronomy, since many of the properties derived for a galaxy depend on the assumed distance. Methods for measuring galaxy distances are described in detail in Chapter 7 of BM. For our purposes it is sufficient to note that in a universe that is homogeneous and isotropic, the relative velocity v between two galaxies that are separated by a large distance r is given by the **Hubble law**,

$$v = H_0 r, \quad (1.13)$$

where H_0 is the **Hubble constant** (see §1.3.1). Our universe is nearly homogeneous and isotropic on large scales, so the velocity field or **Hubble flow**

implied by the Hubble law is approximately correct: the only significant error comes from random velocities of a few hundred km s^{-1} that are generated by the gravitational acceleration from small-scale irregularities in the cosmic matter density. Thus, for example, the distance of a galaxy with a velocity of 7000 km s^{-1} is known to within about 5% once the Hubble constant is known.

By comparing the flux from Cepheid variable stars in the Large Magellanic Cloud and more distant galaxies, Freedman et al. (2001) deduce that

$$H_0 = (72 \pm 8) \text{ km s}^{-1} \text{ Mpc}^{-1}; \quad (1.14)$$

while measurements of small fluctuations in the cosmic background radiation (§1.3.5 and Spergel et al. 2007) give

$$H_0 = (73.5 \pm 3.2) \text{ km s}^{-1} \text{ Mpc}^{-1}. \quad (1.15)$$

When precision is required, we shall write the Hubble constant as

$$\begin{aligned} H_0 &\equiv 70 h_7 \text{ km s}^{-1} \text{ Mpc}^{-1} \\ &= 2.268 h_7 \times 10^{-18} \text{ Hz}, \\ H_0^{-1} &= 13.97 h_7^{-1} \text{ Gyr}, \end{aligned} \quad (1.16)$$

where the dimensionless parameter h_7 is probably within 10% of unity. Any uncertainty in the Hubble constant affects the whole distance scale of the universe and hence is reflected in many of the average properties of galaxies; for example, the mean density of galaxies scales as h_7^3 and the mean luminosity of a galaxy of a given type scales as h_7^{-2} .

If galaxies suffered no acceleration due to external gravitational forces, the distance between any two galaxies would be a linear function of time. Combined with the Hubble law (1.13), this assumption implies that the distance between any two galaxies was zero at a time H_0^{-1} (the **Hubble time**) before the present. The Hubble time provides a rough estimate of the age of the universe. The actual age is somewhat different because the relative velocities of galaxies are decelerated by the gravitational attraction of baryonic and dark matter and accelerated by vacuum energy (see §1.3.3); these effects nearly cancel at present, so the best estimate of the age, $t_0 = 13.7 \text{ Gyr}$ (eq. 1.77), is accidentally very close to the Hubble time.

Galaxies can usefully be divided into four main types according to the **Hubble classification system**—see BM §4.1.1 or Sandage & Bedke (1994) for a more complete description.

(a) Elliptical galaxies These are smooth, featureless stellar systems containing little or no cool interstellar gas or dust and little or no stellar disk. The galaxy M87 shown in Plate 4 is a classic example of this type. The stars

in most elliptical galaxies are old, having ages comparable to the age of the universe, consistent with the absence of gas from which new stars can form.

The fraction of luminous galaxies that are elliptical depends on the local density of galaxies, ranging from about 10% in low-density regions to over 40% in the centers of dense clusters of galaxies (BM §4.1.2).

As the name suggests, the contours of constant surface brightness, or **isophotes**, of elliptical galaxies are approximately concentric ellipses, with axis ratio b/a ranging from 1 to about 0.3. The **ellipticity** is $\epsilon \equiv 1 - b/a$. In the Hubble classification system, elliptical galaxies are denoted by the symbols E0, E1, etc., where a galaxy of type En has axis ratio $b/a = 1 - n/10$. The most elongated elliptical galaxies are type E7. Since we see only the projected brightness distribution, it is impossible to determine directly whether elliptical galaxies are axisymmetric or triaxial; however, indirect evidence strongly suggests that both shapes are present (BM §§4.2 and 4.3).

The surface brightness of an elliptical galaxy falls off smoothly with radius, until the outermost parts are undetectable against the background sky brightness. Because galaxies do not have sharp outer edges, their sizes must be defined with care. One useful measure of size is the **effective radius** R_e , the radius of the isophote containing half of the total luminosity⁷ (or the geometric mean of the major and minor axes of this isophote, if the galaxy is elliptical). The effective radius is correlated with the luminosity of the elliptical galaxy, ranging from 20 kpc for a giant galaxy such as M87 (Plate 4) to 0.2 kpc for a dwarf such as M32 (Plate 3).

The Hubble classification is based on the ellipticity of the isophotes near the effective radius. In many galaxies the isophotes become more elliptical at large radii; thus, for example, M87 is classified as E0 but the isophotal axis ratio is only 0.5 in its outermost parts.

Several empirical formulae have been used to fit the surface-brightness profiles of ellipticals. One of the most successful is the **Sérsic law**

$$I_m(R) = I(0) \exp(-kR^{1/m}) = I_e \exp\{-b_m[(R/R_e)^{1/m} - 1]\}; \quad (1.17)$$

here $I(R)$ is the surface brightness at radius R and I_e is the surface brightness at the effective radius R_e . The parameter m is the **Sérsic index**, which is correlated with the luminosity of the elliptical galaxy, luminous ellipticals having $m \simeq 6$ and dim ones having $m \simeq 2$. The middle of this range is $m = 4$, which defines the **de Vaucouleurs** or $R^{1/4}$ law (de Vaucouleurs 1948). The function b_m must be determined numerically from the condition $\int_0^{R_e} dR R I_m(R) = \frac{1}{2} \int_0^\infty dR R I_m(R)$; but the fitting formula $b_m = 2m - 0.324$ has fractional error $\lesssim 0.001$ over the range $1 < m < 10$ (see Ciotti & Bertin 1999 for properties of Sérsic laws). For $m = 1$ the Sérsic law reduces to the

⁷ The effective radius is measured on the plane of the sky, and is not to be confused with the half-light or median radius (page 17), the radius of a sphere containing half the luminosity.

exponential profile (1.7) that describes the surface-brightness distribution of disk galaxies.

The total luminosity of a galaxy is difficult to define precisely because the outer parts are too faint to measure. One approach is to define a **model luminosity** by fitting the surface-brightness profile to a Sérsic or de Vaucouleurs profile and then estimating the luminosity as $L = \int d^2\mathbf{R} I_m(R)$.

The luminosities of elliptical galaxies range over a factor of 10^8 , from almost $10^{12} L_\odot$ for the very luminous galaxies found at the centers of massive clusters of galaxies, to $\lesssim 10^4 L_\odot$ for the dimmest dwarf galaxies. The **luminosity function** $\phi(L)$ describes the relative numbers of galaxies of different luminosities, and is defined so that $\phi(L) dL$ is the number of galaxies in the luminosity interval $L \rightarrow L + dL$ in a representative unit volume of the universe. A convenient analytic approximation to $\phi(L)$ is the **Schechter law** (BM §4.1.3),

$$\phi(L) dL = \phi_* \left(\frac{L}{L_*} \right)^\alpha \exp(-L/L_*) \frac{dL}{L_*}, \quad (1.18)$$

where $\phi_* \simeq 4.9 \times 10^{-3} h_7^3 \text{Mpc}^{-3}$, $\alpha = -1.1$, and $L_* \simeq 2.9 \times 10^{10} h_7^{-2} L_\odot$ in the R band (Brown et al. 2001). The concept of a “universal” luminosity function embodied in the Schechter law is no more than a good first approximation: in fact the luminosity function is known to depend on both galaxy type and environment (BM §4.1.3).

The average R -band luminosity density derived from equation (1.18) is

$$j_R = \int dL L \phi(L) = \phi_* L_* \int_0^\infty dx x^{\alpha+1} e^{-x} = (\alpha + 1)! \phi_* L_*, \quad (1.19)$$

where the factorial function is defined for non-integer arguments in Appendix C.2. For the parameters given above, $j_R = 1.5 \times 10^8 h_7 L_\odot \text{Mpc}^{-3}$, with an uncertainty of about 30%.

Most luminous elliptical galaxies exhibit little or no rotation, even those with large ellipticity; this is in contrast to stars or other gravitating gas masses, which must be spherical if they do not rotate and flattened when rotating. Among dimmer elliptical galaxies, however, rotation and flattening do appear to be correlated (see §4.4.2c and Faber et al. 1997). This distinction between luminous and dim ellipticals may arise because the most recent mergers of luminous galaxies have been “dry,” that is, between progenitors containing little or no gas, while the recent mergers of low-luminosity ellipticals have involved gas-rich systems. Whether or not this interpretation is correct, the different rotational properties of high-luminosity and low-luminosity elliptical galaxies illustrate that stellar systems can exhibit a much greater variety of equilibria than gaseous systems such as stars.

Each star in an elliptical galaxy orbits in the gravitational field of all the other stars and dark matter in the galaxy. The velocities of individual stars

can be measured in only a few nearby galaxies, but in more distant galaxies the overall distribution of stellar velocities along the line of sight can be determined from the Doppler broadening of lines in the integrated spectrum of the galaxy. The most important parameter describing this distribution is the RMS line-of-sight velocity σ_{\parallel} , sometimes called simply the velocity dispersion (eq. 4.25).

The luminosity, velocity dispersion, and size of elliptical galaxies are correlated. Astronomers usually plot this correlation using not the luminosity but the average surface brightness within the effective radius, which is simply $\bar{I}_e \equiv \frac{1}{2}L/(\pi R_e^2)$. Then if we plot the positions of a sample of elliptical galaxies in the three-dimensional space with coordinates $\log_{10} \bar{I}_e$, $\log_{10} R_e$, and $\log_{10} \sigma_{\parallel}$, they are found to lie on a two-dimensional surface, the **fundamental plane** (see BM §4.3.4 and §4.9.2), given by

$$\log_{10} R_e = 1.24 \log_{10} \sigma_{\parallel} - 0.82 \log_{10} \bar{I}_e + \text{constant}, \quad (1.20)$$

with an RMS scatter of 0.08 in $\log_{10} R_e$ or 0.07 in $\log_{10} \sigma_{\parallel}$ (Jørgensen et al. 1996).

The properties of galaxies are determined both by the fundamental plane and by their distribution within that plane. Let us think of the space with coordinates $(\log_{10} \bar{I}_e, \log_{10} R_e, \log_{10} \sigma_{\parallel})$ as a fictitious three-dimensional space, and imagine observing the distribution of galaxies from a distance. If the line of sight to the observer in this fictitious space lies close to the fundamental plane, the observer will find that galaxies lie close to a line in the two-dimensional space normal to the line of sight. This distribution of galaxies can be thought of as a projection of the distribution in the fundamental plane. The most important of these projections are:

- (i) The **Faber–Jackson law** (BM §4.3.4),

$$\log_{10} \left(\frac{\sigma_{\parallel}}{150 \text{ km s}^{-1}} \right) \simeq 0.25 \log_{10} \left(\frac{L_R}{10^{10} h_7^{-2} L_{\odot}} \right). \quad (1.21)$$

Thus the velocity dispersion of an L_{\star} galaxy is $\sigma_{\parallel} \simeq 200 \text{ km s}^{-1}$. The RMS scatter in the Faber–Jackson law is about 0.1 in $\log_{10} \sigma_{\parallel}$ (Davies et al. 1983).

- (ii) The **Kormendy relation**

$$\log_{10} \left(\frac{\bar{I}_{e,R}}{1.2 \times 10^3 L_{\odot} \text{ pc}^{-2}} \right) = -0.8 \log_{10} \left(\frac{R_e}{h_7^{-1} \text{ kpc}} \right). \quad (1.22a)$$

Here $\bar{I}_{e,R}$ denotes the mean R -band surface brightness interior to R_e . The RMS scatter is less than 0.25 in $\log_{10} \bar{I}_e$. The Kormendy relation implies that

$$\log_{10} \left(\frac{L_R}{7.7 \times 10^9 h_7^{-2} L_{\odot}} \right) = 1.2 \log_{10} \left(\frac{R_e}{h_7^{-1} \text{ kpc}} \right). \quad (1.22b)$$

Thus more luminous galaxies are larger, but have lower surface brightness.

Careful dynamical modeling (§4.9.2) allows us to determine the mass-to-light ratio Υ in elliptical galaxies. These studies show that at radii less than $\sim R_e$ the mass-to-light ratio is not strongly dependent on radius, and consistent with the mass-to-light ratio that we would expect from the observed stellar population (Cappellari et al. 2006). Thus the contribution of dark matter to the mass inside R_e is $\lesssim 30\%$. The mass-to-light ratio is also tightly correlated with σ_e , the luminosity-weighted velocity dispersion within R_e :

$$\Upsilon_I = (3.80 \pm 0.2) \Upsilon_\odot \times \left(\frac{\sigma_e}{200 \text{ km s}^{-1}} \right)^{0.84 \pm 0.07} \quad (1.23)$$

with an intrinsic scatter of only 13%.

Just as stars are found in gravitationally bound systems such as galaxies, many galaxies are found in bound systems called **groups or clusters of galaxies** (see §1.1.5). The largest clusters of galaxies are several Mpc in radius and contain thousands of galaxies. The most luminous galaxy in a large cluster—more often called a **rich cluster**—is often exceptional, in that it is (i) several times more luminous than any other cluster galaxy, and much more luminous than one would expect from the Schechter law (1.18) ($L/L_\star \sim 3\text{--}10$); (ii) at rest in the center of the cluster; (iii) surrounded by a dim stellar halo that extends out to ~ 1 Mpc. Galaxies with these unique characteristics are called **brightest cluster galaxies**; the nearest example is M87 in the Virgo cluster (Plate 4).⁸ The existence of an extended dim halo is also the defining property of **cD galaxies** (BM §4.3.1); in practice, the terms “brightest cluster galaxy” and “cD galaxy” are often used interchangeably. The halo probably arises from stars that have been stripped from individual cluster galaxies by tidal forces and now orbit independently in the cluster’s gravitational field. Brightest cluster galaxies are believed to form during the hierarchical assembly of the cluster from smaller subunits (Dubinski 1998).

The dimmest elliptical galaxies are also unusual. In general, dim ellipticals have higher surface brightness than luminous ellipticals, a manifestation of the Kormendy relation (1.22). However, at luminosities $\lesssim 10^9 L_\odot$ a distinct family of **diffuse dwarf elliptical** or **dwarf spheroidal** galaxies appears, with much larger effective radii and lower surface brightnesses than “normal” ellipticals of the same luminosity (Mateo 1998). Dwarf spheroidal galaxies are difficult to detect because their surface brightness is much less than that of the night sky; nearby dwarf spheroidals are discovered because their brightest stars produce a slight enhancement in star counts that are otherwise dominated by foreground stars belonging to our own Galaxy.

⁸ The most luminous galaxy in the Virgo cluster is actually the E2 galaxy M49=NGC 4472, rather than M87. The Virgo cluster has a complex structure, consisting of two main concentrations, a dominant one near M87 and a smaller one near M49. These probably represent two merging sub-clusters, each with its own brightest cluster galaxy.

There are at least 20 dwarf spheroidal galaxies within 200 kpc, and given the limited sky coverage of existing surveys the actual number may be 50–100 (Belokurov et al. 2007). All appear to be satellites orbiting the Galaxy. Their luminosities range from $2 \times 10^7 L_{\odot}$ to $\lesssim 10^4 L_{\odot}$. The dwarf spheroidals offer a unique probe of dark matter in galaxies, for the following reason. In more luminous galaxies, both baryonic matter (stars and gas) and dark matter contribute comparable amounts to the total mass within the visible stellar system; thus, disentangling their effects to isolate the properties of the dark-matter distribution at small radii is difficult. In some dwarf spheroidal galaxies, however, dark matter contributes 90% or more of the total mass, even at the center of the galaxy, so the dynamics is determined entirely by the gravitational field of the dark matter.

The distribution of mass in the dark halos of ellipticals can be constrained by several methods, including: (i) The kinematics of tracer particles such as globular clusters or planetary nebulae, which typically sample radii from 10–30 kpc (Côté et al. 2003; Romanowsky et al. 2003). This approach relies on the statistical analysis of the positions and velocities of hundreds or thousands of objects, assuming they are found at random orbital phases. (ii) Diffuse X-ray emission from hot gas around the galaxy (Mathews & Brighenti 2003). Measurements of the emissivity and temperature distribution, combined with the plausible assumption that the gas is in hydrostatic equilibrium, can be used to constrain the distribution of the dark matter out to ~ 30 kpc in isolated galaxies. The same technique can be applied to brightest cluster galaxies out to much larger radii, but in this case we are measuring the combined dark-matter distribution of the galaxy and the cluster. (iii) Kinematics of satellite galaxies. This technique is similar in principle to the use of globular clusters or planetary nebulae; the satellite galaxies have the advantage that they sample much larger radii, from 100–400 kpc, but the disadvantage that generally no more than one satellite is detected around a given galaxy, so the method yields only an average of the dark-matter distribution over many galaxies (Prada et al. 2003). (iv) Weak gravitational lensing, in which the gravitational field of a nearby galaxy distorts the images of distant background galaxies (Schneider 2006); once again, this method requires averaging over a large sample of lensing galaxies.

The preliminary conclusion from these studies is that luminous, isolated elliptical galaxies contain dark halos that are much larger—both in size and in mass—than the stellar systems they surround. Within uncertainties of at least a factor of two, the halos extend to ~ 300 kpc and contain ~ 10 times the mass in stars.

(b) Spiral galaxies These are galaxies, like the Milky Way and M31, that contain a prominent disk composed of stars, gas, and dust. The disk contains **spiral arms**, filaments in which stars are continuously being formed. The same spiral arms are seen in the old stars that dominate the mass of the

disk (see Figure 6.1 and §6.1.2). The spiral arms vary greatly in their shape, length and prominence from one galaxy to another but are always present.

In low-density regions of the universe, about 60% of all luminous galaxies are spirals, but the fraction drops to $\lesssim 10\%$ in dense regions such as the cores of galaxy clusters (BM §4.1.2).

The surface brightness in spiral galaxy disks, which traces the radial distribution of stars, obeys the exponential law (1.7) (de Jong 1996). A typical disk scale length is $R_d \simeq 2h_7^{-1}$ kpc, but scale lengths range from $1h_7^{-1}$ kpc to more than $10h_7^{-1}$ kpc. The typical central surface brightness is $I_d \sim 100 L_\odot \text{pc}^{-2}$ (BM Figure 4.52). The interstellar gas in spiral galaxy disks often extends to much larger radii than the stars, probably because star formation is suppressed when the gas surface density falls below a critical value (see Plates 5, 6 and BM §8.2.8).

Using the 21-cm line of interstellar neutral hydrogen, the circular-speed curves $v_c(R)$ of spiral galaxies can be followed out to radii well beyond the outer edge of the stellar distribution. The circular-speed curves of luminous spirals are nearly flat out to the largest radii at which they can be measured, often a factor of two or more larger than the edge of the stellar disk (BM §8.2.4). If most of the mass of the galaxy were in stars, we would expect the circular-speed curve at these large radii to fall as $v_c(R) = (GM/R)^{1/2}$ where M is the total stellar mass (see Figure 2.17). The inescapable conclusion is that the mass of the galaxy at these large radii is dominated by the dark halo rather than the stars.

Typical circular speeds of spirals are between 100 and 300 km s^{-1} . Just as the velocity dispersion of elliptical galaxies is related to their luminosity by the Faber–Jackson law (1.21), the rotation rate of spirals in the flat part of the circular-speed curve is related to their luminosity by the **Tully–Fisher law** (BM §7.3.4; Sakai et al. 2000),

$$\log_{10} \left(\frac{L_R h_7^2}{10^{10} L_\odot} \right) = 3.5 \log_{10} \left(\frac{v_c}{200 \text{ km s}^{-1}} \right) + 0.5; \quad (1.24)$$

the RMS scatter in this relation is about 0.14 in $\log_{10} L_R$. The slope of the Tully–Fisher law is a function of the wavelength band in which the luminosity is measured, ranging from $\simeq 3$ in the B band centered at $0.45 \mu\text{m}$ to $\simeq 4$ in the K band centered at $2.2 \mu\text{m}$. Applied to our Galaxy, using $v_c = (220 \pm 20) \text{ km s}^{-1}$ from equation (1.8), we find $L_R = (4.4 \pm 1.5) \times 10^{10} L_\odot$, consistent with Table 1.2.

Like the Milky Way, most spiral galaxies contain a bulge, a centrally concentrated stellar system that has a smooth or amorphous appearance—quite unlike that of the disk, which exhibits spiral arms, dust lanes, concentrations of young stars, and other structure. The origin of bulges is not well understood: some resemble small elliptical galaxies and presumably formed in the same way, while others resemble thickened disks, and may have formed from

the disk through dynamical processes (see §6.6.2 and Kormendy & Kennicutt 2004). Bulges and elliptical galaxies are sometimes called **spheroidal stellar systems** or **spheroids**, even though their shapes are not necessarily close to mathematical spheroids (page 76)—in particular, many of them are probably not axisymmetric.

The luminosity of the bulge relative to that of the disk is correlated with many other properties of the galaxy, such as the fraction of the disk mass in gas, the color of the disk, and how tightly the spiral arms are wound. This correlation is the basis of a sub-division in the Hubble classification system, which breaks up spiral galaxies into four classes or types, called Sa, Sb, Sc, Sd (Sandage & Bedke 1994). Along the sequence Sa→Sd, (i) the relative luminosity of the bulge decreases; (ii) the spiral arms become more loosely wound; (iii) the relative mass of gas increases; and (iv) the spiral arms become more clumpy, so individual patches of young stars and HII regions become more prominent. This sequence is illustrated by comparing the images of M104 (Plate 7), which is classified Sa; M81 (Plate 8), which is classified Sab (i.e., between Sa and Sb); the Sb galaxy M31 (Plate 3); the Sbc galaxies M51, M63, and M100 (Plates 1, 9, and 17); the Sc galaxy M101 (Plate 18), and the Scd galaxy M33 (Plate 19). The Milky Way is type Sbc.

The Hubble classification also divides spiral galaxies into “normal” and “barred” categories. The bar is an elongated, smooth stellar system that is reminiscent of a rigid paddle or stirrer rotating at the center of the galactic disk. The bar can be thought of as a triaxial bulge, and in practice there is no clear distinction between these two categories of stellar system: for example, a “bar” in a face-on galaxy might well be classified as a “bulge” if the galaxy were viewed edge-on. Further properties of bars are described in §6.5.

A classic barred galaxy is NGC 1300 (Plate 10) although most bars are less prominent than the one in this galaxy. Other barred galaxies are shown in Figures 6.27 and 6.28. Our own Galaxy and its neighbor, the Large Magellanic Cloud (Plate 11), are both barred. About half of all spirals are barred, and bars appear in all of the Hubble classes Sa, Sb, Sc, Sd, where their presence is indicated by inserting the letter “B” into the notation (SBa, SBb, etc.). Elliptical galaxies do not have bars.

The first evidence for dark halos in spiral galaxies came from circular-speed curves; as we have discussed, neutral-hydrogen rotation curves in some spirals remain flat out to as much as 10 times the scale length of the stellar disk, which implies that the mass at large radii is dominated by dark matter rather than stars. At much larger radii, $\gtrsim 100$ kpc, the distribution of dark matter can be measured by the same techniques that are used for ellipticals, in particular satellite galaxy kinematics and weak gravitational lensing. Within the large uncertainties, the data are consistent with the hypothesis that the size and mass of the dark halos that surround luminous spiral galaxies are the same as those surrounding isolated ellipticals—about 300 kpc in radius and containing ~ 10 times the stellar mass of the galaxy.

For most spiral galaxies, the relative contributions of dark and luminous matter within the visible stellar system are difficult to disentangle (§6.3.3). However, in some low-luminosity and low surface-brightness spirals, dark matter appears to dominate the mass at all radii (Swaters et al. 2003). Like the dwarf spheroidals, these galaxies provide valuable probes of the properties of dark halos on small scales.

(c) Lenticular galaxies These are transition objects between elliptical and spiral galaxies: like spirals, they contain a rapidly rotating disk, a bulge, and sometimes a bar, and the disk obeys the exponential surface-brightness law (1.7) characteristic of spirals. Like ellipticals, they have little or no cool gas or recent star formation, are smooth and featureless in appearance, and exhibit no spiral structure. The absence of young stars is a consequence of the absence of gas, since this is the raw material from which stars are formed.

Lenticulars are rare in low-density regions, but comprise almost half of the galaxies in the high-density centers of galaxy clusters (BM Figure 4.10). This correlation suggests that lenticulars may be spirals that have been depleted of interstellar gas by interactions with the hot gas in the cluster (van Gorkom 2004).

Lenticulars are labeled in the Hubble classification by the notation S0, or SB0 if barred. The transition from ellipticals to lenticulars to spirals is smooth and continuous, so there are S0 galaxies that might well be classified as E7 and others that could be Sa (Sandage & Bedke 1994).

(d) Irregular galaxies Along the sequence from Sc to Sd, galaxies become progressively less luminous and their spiral structure becomes less well defined. These trends continue beyond Sd: we find low-luminosity (“dwarf”) disk galaxies in which the young stars are arranged chaotically rather than in spirals. These are called “irregular” galaxies and are denoted in the Hubble classification by Sm or Im, the prototypes of these two classes being the Large and Small Magellanic Clouds.

Irregular galaxies are extremely common—more than a third of our neighbors are of this type—but they do not feature prominently in most galaxy catalogs, because any flux-limited catalog is biased against intrinsically dim systems.

In irregulars the circular speed is a linear function of radius (corresponding to a constant angular speed) over most of the stellar disk, reaching a maximum of $\sim 50\text{--}70\text{ km s}^{-1}$ near the edge of the disk. These properties are in sharp contrast to luminous spiral galaxies, in which the circular speed is much higher and the circular-speed curve is nearly flat.

Much of the luminosity of irregular galaxies is emitted by massive young stars and large HII regions. These systems are extremely gas-rich: the interstellar gas in their disks often contains more than 30% of the mass in stars. Their irregular appearance arises partly because the optical emission is dominated by a relatively small number of luminous young stars and HII regions, and partly because the circular speed in the disk is not that much larger than the turbulent velocities in the interstellar gas ($\sim 10\text{ km s}^{-1}$).

A minority of galaxies are assigned to the “irregular” bin simply because they fit nowhere else: these include spiral or elliptical galaxies that have been violently distorted by a recent encounter with a neighbor (see Plate 12), galaxies in the last stages of merging, and galaxies that are undergoing an intense burst of star formation that overwhelms the stellar population that usually determines the classification.

It is convenient to think of the Hubble classification as a sequence $E \rightarrow S0 \rightarrow Sa \rightarrow Sb \rightarrow Sc \rightarrow Sd \rightarrow Sm \rightarrow Im$. Galaxies near the beginning of this sequence are called **early**, while those near the end are **late**. Thus the term “early-type galaxies” refers to ellipticals and lenticulars; an Sa galaxy is an “early-type spiral,” while an Sc or Sd galaxy is a “late-type spiral,” etc. This terminology is a fossil of the initial incorrect belief that the Hubble sequence was an evolutionary or time sequence.

1.1.4 Open and globular clusters

A typical galaxy contains many small stellar systems of between 10^2 and 10^6 stars. These systems are called **star clusters** and can be divided into two main types.

Open clusters are irregular stellar systems that contain $\sim 10^2$ to 10^4 stars (see Table 1.3, Plate 13, and BM §6.2). New open clusters are formed continuously in the Galactic disk, and most of the ones we see are younger than 1 Gyr (Figure 8.5). Older clusters are rare because most have been disrupted, probably by gravitational shocks from passing interstellar gas clouds (§8.2.2c). There are over 1000 cataloged open clusters out of an estimated 10^5 throughout the Galaxy. It is likely that most of the stars in the Galactic disk formed in open clusters that have since dissolved.

Globular clusters are much more massive stellar systems, containing 10^4 – 10^6 stars in a nearly spherical distribution (BM §§4.5 and 6.1, and Plate 14; see Ashman & Zepf 1998 and Carney & Harris 2001 for reviews). Globular clusters do not contain gas, dust, or young stars. Our Galaxy contains about 150 globular clusters, but large elliptical galaxies such as M87 can contain as many as 10 000 (Plate 4). Unlike open clusters, the Galaxy’s globular clusters are old, and are believed to be relics of the formation of the Galaxy itself.⁹ The metallicity appears to be the same for all the stars in a given cluster—presumably because the cluster formed from a well-mixed gas cloud—but different clusters have a wide range of metallicity, from only $0.005Z_{\odot}$ to nearly solar. The spatial distribution and the kinematics of a group of clusters are correlated with the metallicity, and for many purposes the clusters in our Galaxy can be divided into two groups (Zinn 1985): a roughly spherical population that contains 80% of the clusters, shows little

⁹ It is a mystery why young globular clusters are absent in our Galaxy but common in many others, such as M31, the Large Magellanic Cloud, and galaxies that have undergone recent mergers.

or no rotation and has metallicity $Z < 0.1Z_{\odot}$, and is associated with the stellar halo; and a flattened population that contains the remaining 20%, has $Z > 0.1Z_{\odot}$, exhibits rapid rotation, and is associated with the disk and bulge. This bimodal distribution of metallicities is present in the globular-cluster systems of other galaxies as well (Gebhardt & Kissler-Patig 1999).

The stellar density in the center of a globular cluster is extremely high: a typical value is $10^4 \mathcal{M}_{\odot} \text{pc}^{-3}$, compared with $0.05 \mathcal{M}_{\odot} \text{pc}^{-3}$ in the solar neighborhood. Because globular clusters have strong or high **central concentration** (the central density is much larger than the mean density) three different measures of the radius are usually quoted for globular clusters: the **core radius**, where the surface brightness has fallen to half its central value; the median or half-light radius, the radius of a sphere that contains half of all the luminosity; and the **limiting** or **tidal radius**, the outer limit of the cluster where the density drops to zero. Typical values of these and other cluster parameters are given in Table 1.3.

Luminous globular clusters emit as much light as dwarf spheroidal galaxies. However, a dwarf spheroidal galaxy is a very low surface-brightness object with a half-light radius of $\sim 300 \text{pc}$, while a luminous globular cluster has a much smaller radius ($\sim 3 \text{pc}$) and a correspondingly higher surface brightness. A handful of exceptionally luminous globular clusters, such as ω Centauri in our Galaxy and G1 in M31, may be the dense centers of tidally disrupted galaxies (Freeman 1993).

Globular clusters are among the simplest stellar systems: they are spherical, they have no dust or young stars to obscure or confuse the observations, they appear to have no dark matter other than low-luminosity stars, and they are **dynamically old**: a typical star in a globular cluster has completed many orbits ($\sim 10^4$) since the cluster was formed. Thus globular clusters provide the best physical realization we have of the **gravitational N-body problem**, which is to understand the evolution of a system of N point masses interacting only by gravitational forces (Chapter 7).

1.1.5 Groups and clusters of galaxies

Galaxies are not distributed uniformly in the universe. They belong to a rich hierarchy of structure that includes binary galaxies, small groups of a few galaxies in close proximity, enormous voids in which the number density of galaxies is greatly depleted, filaments and walls stretching for tens of Mpc, and rare giant clusters containing thousands of galaxies (see §9.2.2 and Mulchaey, Dressler, & Oemler 2004). Only on scales $\gtrsim 100 \text{Mpc}$ is the distribution of galaxies statistically homogeneous.

Associations that contain only a handful of galaxies are called groups while bigger associations are called clusters of galaxies (see Plates 15 and 16). The dividing line between groups and clusters is arbitrary, since the distribution of properties is continuous from one class to the other.

Table 1.3 Parameters of globular and open clusters

	globular	open
central density ρ_0	$1 \times 10^4 \mathcal{M}_\odot \text{pc}^{-3}$	$10 \mathcal{M}_\odot \text{pc}^{-3}$
core radius r_c	1 pc	1 pc
half-mass radius r_h	3 pc	2 pc
tidal radius r_t	35 pc	10 pc
central velocity dispersion σ_0	6 km s^{-1}	0.3 km s^{-1}
crossing time r_h/σ_0 (line-of-sight)	0.5 Myr	7 Myr
mass-to-light ratio Υ_R	$2\Upsilon_\odot$	$1\Upsilon_\odot$
mass M	$2 \times 10^5 \mathcal{M}_\odot$	$300 \mathcal{M}_\odot$
lifetime	10 Gyr	300 Myr
number in the Galaxy	150	10^5

NOTES: Values for globular clusters are medians from the compilation of Harris (1996). Values for open clusters are from Figure 8.5, Piskunov et al. (2007), and other sources.

The galaxies within ~ 1 Mpc are members of the **Local Group**. The two dominant members of this group are the Galaxy and M31. Dozens of smaller galaxies, mostly satellites of the two dominant galaxies, are also members (see BM §4.1.4 and van den Bergh 2000). The Local Group is believed to be a physical system rather than a chance superposition because the density of galaxies in this region is substantially higher than average, and because the Galaxy and M31 are approaching one another rather than receding with the Hubble flow. It is believed that the gravitational attraction between these two galaxies slowed and then reversed their recession, and that they will eventually merge into a single giant stellar system (see Box 3.1 and Figure 8.1).

Like star clusters, groups and clusters of galaxies may be regarded for many purposes as assemblies of masses orbiting under their mutual gravitational attraction, except that now the masses are galaxies rather than stars. However, there are two important differences between the dynamics of star clusters and galaxy groups or clusters. First, groups and clusters of galaxies are **dynamically young**: a typical galaxy in even the largest and most populous clusters has completed only a few orbits since the cluster formed, and in many smaller groups, including the Local Group, galaxies are still falling towards the group center for the first time. Second, the fractional volume of a group or cluster that is occupied by galaxies ($\gtrsim 10^{-3}$) is much larger than the fractional volume of a star cluster that is occupied by stars ($\approx 10^{-19}$). Thus collisions between galaxies in a cluster are much more frequent than collisions between stars in a star cluster (see Chapter 8).

Clusters of galaxies are the largest equilibrium structures in the universe. They arose from the gravitational collapse of rare high peaks in the fluctuating density field of dark matter in the early universe (§9.2). Conse-

quently their properties provide a sensitive probe of cosmological parameters. Clusters also offer a unique probe of the distribution of dark matter on large scales. The mass distribution in clusters of galaxies can be measured by many complementary methods, including (i) statistical analysis of the velocities and positions of large numbers of galaxies in the cluster; (ii) measurements of the X-ray emissivity and temperature of hot gas in the cluster; (iii) distortion of the images of background galaxies by weak gravitational lensing; (iv) strong gravitational lensing, which can produce multiple images of background galaxies near the center of the cluster; (v) the **Sunyaev–Zeldovich effect**, which is a slight depression in the measured temperature of the cosmic microwave background at the locations of clusters, caused by Compton scattering of photons by electrons in the hot cluster gas.

The biggest clusters have masses $\sim 10^{15} \mathcal{M}_\odot$ within 2 Mpc of their centers and velocity dispersions of $\sim 1000 \text{ km s}^{-1}$. The mass-to-light ratios are

$$\Upsilon_R \simeq (200 \pm 50) h_7 \Upsilon_\odot, \quad (1.25)$$

(Fukugita, Hogan, & Peebles 1998), with no detectable dependence on cluster properties such as velocity dispersion or total population.

Most of the baryons in clusters of galaxies are in the hot gas. The mass in gas is a fraction $0.11 h_7^{-3/2}$ of the total mass (Allen, Schmidt & Fabian 2002), while the mass in stars is only about $0.02 h_7^{-1}$ of the total. Thus the fraction of the total mass that resides in baryons is

$$f_b = 0.13 \pm 0.02, \quad (1.26)$$

with the remaining 87% comprised of WIMPs or other non-baryonic dark matter. In structures as large as clusters it is difficult to imagine how baryons and non-baryonic dark matter could be segregated (in contrast to individual galaxies, where the baryons have concentrated at the center of the dark halo to form the visible stars). Thus the baryon-to-total mass ratio f_b that is found in clusters should be a fair sample of the universe as a whole.

1.1.6 Black holes

Dynamical studies of the centers of galaxies reveal that they often contain “massive dark objects”—concentrations of 10^6 – $10^9 \mathcal{M}_\odot$ contained within a few pc of the center (BM §11.2.2). The best-studied of these objects, the one at the center of the Galaxy, has a mass of $(3.9 \pm 0.3) \times 10^6 \mathcal{M}_\odot$ contained within a radius less than 0.001 pc. Astronomers believe that these objects must be black holes, for two main reasons. First, dynamical arguments show that no long-lived astrophysical system other than a black hole could be so massive and so small (§7.5.2). Second, many galaxies contain strong sources of non-stellar radiation at their centers, called **active galactic nuclei** or

AGN; the most luminous and rare of these, the quasars, can achieve luminosities of $10^{13} L_{\odot}$ and outshine their host galaxies by two orders of magnitude (see BM §4.6.2 and Krolik 1999). By far the most plausible power source for AGN is accretion onto a massive black hole, and the demography of massive dark objects in galaxy centers is roughly consistent with the hypothesis that these are dormant AGN.

It appears that most galaxies—or at least most early-type galaxies—contain a central black hole. The mass of the black hole typically amounts to ≈ 0.001 – 0.002 of the total mass of the stars in the host galaxy (Häring & Rix 2004). Another correlation that is more directly observable is between the black-hole mass M_{\bullet} and the velocity dispersion σ_{\parallel} near the center of the host galaxy,

$$\log_{10} \left(\frac{M_{\bullet} h_7}{10^8 M_{\odot}} \right) = (4 \pm 0.3) \log_{10} \left(\frac{\sigma_{\parallel}}{200 \text{ km s}^{-1}} \right) + (0.2 \pm 0.1). \quad (1.27)$$

The RMS scatter in this relation is $\lesssim 0.3$ in $\log_{10} M_{\bullet}$ (Tremaine et al. 2002).

Massive black holes are probably formed at the centers of galaxies. When galaxies merge, their black holes are dragged to the center of the merged galaxy by dynamical friction (§8.1.1a). If the resulting binary black hole is so tightly bound that it continues to decay by gravitational radiation, the two black holes will eventually merge. The final stages of this merger could provide a powerful source of gravitational radiation (§8.1.1e; Begelman, Blandford & Rees 1980).

1.2 Collisionless systems and the relaxation time

There is a fundamental difference between galaxies and the systems that are normally dealt with in statistical mechanics, such as molecules in a box. This difference lies in the nature of the forces that act between the constituent particles. The interaction between two molecules is short-range: the force is small unless the molecules are very close to each other, when it becomes strongly repulsive. Consequently, molecules in a diffuse gas are subject to violent and short-lived accelerations as they collide with one another, interspersed with much longer periods when they move at nearly constant velocity. In contrast, the gravitational force that acts between the stars of a galaxy is long-range.

Consider the force from the stars in the cone shown in Figure 1.4 on a star at the apex of the cone. The force from any one star falls off with distance r as r^{-2} , but if the density of stars is uniform, the number of attracting stars per unit length of the cone increases as r^2 . Let us call a factor of two interval in radius an **octave**, by analogy with the musical octave. Then each octave in radius, from r to $2r$, has a length proportional to r , so each octave attracts

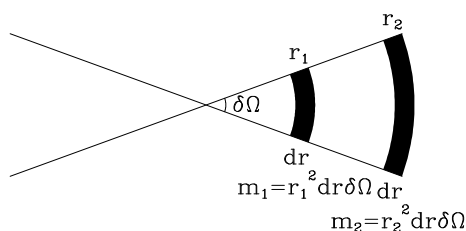


Figure 1.4 If the density of stars were everywhere the same, the stars in each of the shaded segments of a cone would contribute equally to the force on a star at the cone's apex. Thus the acceleration of a star at the apex is determined mainly by the large-scale distribution of stars in the galaxy, not by the star's nearest neighbors.

the star at the apex with a force proportional to $r^{-2} \times r^2 \times r = r$. This simple argument shows that the force on the star at the apex is dominated by the most distant stars in the system, rather than by its closest neighbors. Of course, if the density of attracting stars were exactly spherical, the star at the apex would experience no net force because it would be pulled equally in all directions. But in general the density of attracting stars falls off in one direction more slowly than in the opposing direction, so the star at the apex is subject to a net force, and this force is determined by the structure of the galaxy on the largest scale. Consequently—in contrast to the situation for molecules—the force on a star does not vary rapidly, and each star may be supposed to accelerate smoothly through the force field that is generated by the galaxy as a whole. In other words, for most purposes we can treat the gravitational force on a star as arising from a smooth density distribution rather than a collection of mass points.

1.2.1 The relaxation time

We now investigate this conclusion more quantitatively, by asking how accurately we can approximate a galaxy composed of N identical stars of mass m as a smooth density distribution and gravitational field. To answer this question, we follow the motion of an individual star, called the **subject star**, as its orbit carries it once across the galaxy, and seek an order-of-magnitude estimate of the difference between the actual velocity of this star after this interval and the velocity that it would have had if the mass of the other stars were smoothly distributed. Suppose the subject star passes within distance b of another star, called the **field star** (Figure 1.5). We want to estimate the amount $\delta\mathbf{v}$ by which the encounter deflects the velocity \mathbf{v} of the subject star. In §3.1d we calculate $\delta\mathbf{v}$ exactly, but for our present purposes an approximate estimate is sufficient. To make this estimate we shall assume that $|\delta\mathbf{v}|/v \ll 1$, and that the field star is stationary during the encounter. In this case $\delta\mathbf{v}$ is perpendicular to \mathbf{v} , since the accelerations parallel to \mathbf{v} average to zero. We may calculate the magnitude of the velocity change, $\delta v \equiv |\delta\mathbf{v}|$, by assuming that the subject star passes the field star on a straight-line trajectory, and integrating the perpendicular force F_{\perp} along this trajectory. We place the origin of time at the instant of closest approach of the two stars,

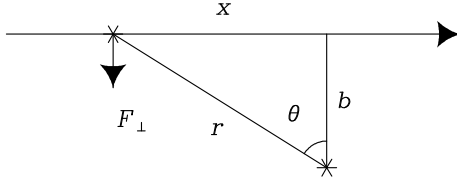


Figure 1.5 A field star approaches the subject star at speed v and impact parameter b . We estimate the resulting impulse to the subject star by approximating the field star's trajectory as a straight line.

and find in the notation of Figure 1.5,

$$F_{\perp} = \frac{Gm^2}{b^2 + x^2} \cos \theta = \frac{Gm^2 b}{(b^2 + x^2)^{3/2}} = \frac{Gm^2}{b^2} \left[1 + \left(\frac{vt}{b} \right)^2 \right]^{-3/2}. \quad (1.28)$$

But by Newton's laws

$$m\dot{\mathbf{v}} = \mathbf{F} \quad \text{so} \quad \delta v = \frac{1}{m} \int_{-\infty}^{\infty} dt F_{\perp}, \quad (1.29)$$

and we have

$$\delta v = \frac{Gm}{b^2} \int_{-\infty}^{\infty} \frac{dt}{[1 + (vt/b)^2]^{3/2}} = \frac{Gm}{bv} \int_{-\infty}^{\infty} \frac{ds}{(1 + s^2)^{3/2}} = \frac{2Gm}{bv}. \quad (1.30)$$

In words, δv is roughly equal to the acceleration at closest approach, Gm/b^2 , times the duration of this acceleration $2b/v$. Notice that our assumption of a straight-line trajectory breaks down, and equation (1.30) becomes invalid, when $\delta v \simeq v$; from equation (1.30), this occurs if the impact parameter $b \lesssim b_{90} \equiv 2Gm/v^2$. The subscript 90 stands for a 90-degree deflection—see equation (3.51) for a more precise definition.

Now the surface density of field stars in the host galaxy is of order $N/\pi R^2$, where N is the number of stars and R is the galaxy's radius, so in crossing the galaxy once the subject star suffers

$$\delta n = \frac{N}{\pi R^2} 2\pi b db = \frac{2N}{R^2} b db \quad (1.31)$$

encounters with impact parameters in the range b to $b + db$. Each such encounter produces a perturbation $\delta \mathbf{v}$ to the subject star's velocity, but because these small perturbations are randomly oriented in the plane perpendicular to \mathbf{v} , their mean is zero.¹⁰ Although the mean velocity change is zero, the mean-square change is not: after one crossing this amounts to

$$\sum \delta v^2 \simeq \delta v^2 \delta n = \left(\frac{2Gm}{bv} \right)^2 \frac{2N}{R^2} b db. \quad (1.32)$$

¹⁰ Strictly, the mean change in velocity is zero only if the distribution of perturbing stars is the same in all directions. A more precise statement is that the mean change in velocity is due to the smoothed-out mass distribution, and we ignore this because the goal of our calculation is to determine the *difference* between the acceleration due to the smoothed mass distribution and the actual stars.

Integrating equation (1.32) over all impact parameters from b_{\min} to b_{\max} , we find the mean-square velocity change per crossing,

$$\Delta v^2 \equiv \int_{b_{\min}}^{b_{\max}} \sum \delta v^2 \simeq 8N \left(\frac{Gm}{Rv} \right)^2 \ln \Lambda, \quad (1.33a)$$

where the factor

$$\ln \Lambda \equiv \ln \left(\frac{b_{\max}}{b_{\min}} \right) \quad (1.33b)$$

is called the **Coulomb logarithm**. Our assumption of a straight-line trajectory breaks down for impact parameters smaller than b_{90} , so we set $b_{\min} = f_1 b_{90}$, where f_1 is a factor of order unity. Our assumption of a homogeneous distribution of field stars breaks down for impact parameters of order R , so we set $b_{\max} = f_2 R$. Then

$$\ln \Lambda = \ln \left(\frac{R}{b_{90}} \right) + \ln(f_2/f_1). \quad (1.34)$$

In most systems of interest $R \gg b_{90}$ (for example, in a typical elliptical galaxy $R/b_{90} \gtrsim 10^{10}$), so the fractional uncertainty in $\ln \Lambda$ arising from the uncertain values of f_1 and f_2 is quite small, and we lose little accuracy by setting $f_2/f_1 = 1$.

Thus encounters between the subject star and field stars cause a kind of diffusion of the subject star's velocity that is distinct from the steady acceleration caused by the overall mass distribution in the stellar system. This diffusive process is sometimes called **two-body relaxation** since it arises from the cumulative effect of myriad two-body encounters between the subject star and passing field stars.

The typical speed v of a field star is roughly that of a particle in a circular orbit at the edge of the galaxy,

$$v^2 \approx \frac{GNm}{R}. \quad (1.35)$$

If we eliminate R from equation (1.33a) using equation (1.35), we have

$$\frac{\Delta v^2}{v^2} \approx \frac{8 \ln \Lambda}{N}. \quad (1.36)$$

If the subject star makes many crossings of the galaxy, the velocity \mathbf{v} will change by roughly Δv^2 at each crossing, so the number of crossings n_{relax} that is required for its velocity to change by of order itself is given by

$$n_{\text{relax}} \simeq \frac{N}{8 \ln \Lambda}. \quad (1.37)$$

The **relaxation time** may be defined as $t_{\text{relax}} = n_{\text{relax}} t_{\text{cross}}$, where $t_{\text{cross}} = R/v$ is the **crossing time**, the time needed for a typical star to cross the galaxy once. Moreover $\Lambda = R/b_{90} \approx Rv^2/(Gm)$, which is $\approx N$ by equation (1.35). Thus our final result is

$$t_{\text{relax}} \simeq \frac{0.1N}{\ln N} t_{\text{cross}}. \quad (1.38)$$

After one relaxation time, the cumulative small kicks from many encounters with passing stars have changed the subject star’s orbit significantly from the one it would have had if the gravitational field had been smooth. In effect, after a relaxation time a star has lost its memory of its initial conditions. Galaxies typically have $N \approx 10^{11}$ stars and are a few hundred crossing times old, so for these systems stellar encounters are unimportant, except very near their centers. In a globular cluster, on the other hand, $N \approx 10^5$ and the crossing time $t_{\text{cross}} \approx 1$ Myr (Table 1.3), so relaxation strongly influences the cluster structure over its lifetime of 10 Gyr.

In all of these systems the dynamics over timescales $\lesssim t_{\text{relax}}$ is that of a **collisionless system** in which the constituent particles move under the influence of the gravitational field generated by a smooth mass distribution, rather than a collection of mass points. Non-baryonic dark matter is also collisionless, since both weak interactions and gravitational interactions between individual WIMPs are negligible in any galactic context.

In most of this book we focus on collisionless stellar dynamics, confining discussion of the longer-term evolution that is driven by gravitational encounters among the particles to Chapter 7.

1.3 The cosmological context

This section provides a summary of the aspects of cosmology that we use in this book. For more information the reader can consult texts such as Weinberg (1972), Peebles (1993), and Peacock (1999).

To a very good approximation, the universe is observed to be homogeneous and isotropic on large scales—here “large” means $\gtrsim 100$ Mpc, which is still much smaller than the characteristic “size” of the universe, the **Hubble length** $c/H_0 = 4.3h_7^{-1}$ Gpc where 1 Gpc = 10^9 pc = 10^3 Mpc and c is the speed of light. Therefore a useful first approximation is to average over the small-scale structure and treat the universe as *exactly* homogeneous and isotropic. Of course, the universe does not appear isotropic to all observers: an observer traveling rapidly with respect to the local matter will see galaxies approaching in one direction and receding in another. Therefore we define a set of **fundamental observers**, who are at rest with respect to the matter around them.¹¹ The universe is expanding, so we may synchronize the

¹¹ A more precise definition is that a fundamental observer sees no dipole component in the cosmic microwave background radiation (§1.3.5).

clocks of the fundamental observers by setting them to the same time at the moment when the universal homogeneous density has some particular value. This procedure enables us to define a unique **cosmic time** (Gunn 1978).

The random velocities of galaxies, and the velocities of stars within galaxies (a few hundred km s^{-1}), are small compared to the relative velocities of galaxies that are separated by 100 Mpc (many thousand km s^{-1}). Thus, on scales large enough that the assumption of homogeneity is accurate, any observer who moves with a typical star in a galaxy is a fundamental observer.

1.3.1 Kinematics

Consider the triangle defined by three nearby fundamental observers. As the universe evolves, the triangle may change in size, but cannot change in shape or orientation—in the contrary case, it would define a preferred direction, thereby violating the isotropy assumption. Thus, if $r_{ij}(t)$ is the length of the side joining observers i and j at cosmic time t , we must have $r_{ij}(t) = r_{ij}(t_0)a(t)$, where $a(t)$ is independent of i and j . Since this argument holds for all fundamental observers, the distance between any two of them must have the form

$$r(t) = r(t_0)a(t), \quad (1.39)$$

where the **scale factor** $a(t)$ is a universal function, which we may normalize so that $a(t_0) = 1$ at the present cosmic time t_0 . The relative velocity of the two observers is

$$v(t) = \frac{dr}{dt} = r(t_0)\dot{a}(t) = r(t)\frac{\dot{a}(t)}{a(t)} \equiv r(t)H(t), \quad (1.40)$$

where $H(t)$ is the **Hubble parameter**. At the present time, $H(t_0) \equiv H_0$ is the Hubble constant, and equation (1.40) is a statement of the Hubble law (1.13). Thus we see that (i) the Hubble law is a consequence of homogeneity and isotropy; (ii) in a homogeneous, isotropic universe the Hubble law remains true at all times but the Hubble “constant” varies with cosmic time.

Next consider a photon that at cosmic time t passes a fundamental observer, who observes it to have frequency ν . After an infinitesimal time interval dt , the photon has traveled a distance $dr = c dt$ and hence is overtaking a second fundamental observer who is moving away from the first at speed $dv = H(t)dr = H(t)c dt$. This observer will measure a different frequency for the photon because of the Doppler shift. The measured frequency will be $\nu(1 - dv/c) = \nu[1 - H(t)dt]$; the use of this first-order formula is justified because dv is infinitesimal. Thus the frequency of a propagating photon as measured by a local fundamental observer decreases at a rate

$$\frac{d\nu}{dt} = -H(t)\nu \quad \text{or} \quad \frac{\dot{\nu}}{\nu} = -\frac{\dot{a}}{a}, \quad \text{thus} \quad \nu(t) \propto \frac{1}{a(t)}. \quad (1.41)$$

In words, a photon emitted by a fundamental observer at frequency ν_e and wavelength $\lambda_e = c/\nu_e$, and received by a second fundamental observer at the present time t_0 , will be observed to have frequency ν_0 and wavelength λ_0 given by

$$\frac{\nu_e}{\nu_0} = \frac{\lambda_0}{\lambda_e} = \frac{a(t_0)}{a(t_e)} = \frac{1}{a(t_e)} \equiv 1 + z. \quad (1.42)$$

Here z is the **redshift**. Redshift is often used instead of time t to describe the cosmic time of an event, since redshift is directly observable from the wavelengths of known spectral lines, whereas the relation between time and redshift depends on the cosmological model (Figure 1.7).

These derivations assume only that spacetime is locally Euclidean, and thus they are correct even in a curved spacetime, so long as it is homogeneous and isotropic.

1.3.2 Geometry

Let the position of any fundamental observer be labeled by time-independent coordinates (q_1, q_2, q_3) . At a given cosmic time t , the distance dl between the observers at (q_1, q_2, q_3) and $(q_1 + dq_1, q_2 + dq_2, q_3 + dq_3)$ can be written in the form

$$dr^2 = a^2(t)h_{ij}dq_i dq_j, \quad (1.43)$$

where we have used the summation convention (page 772), and the metric tensor h_{ij} (cf. eq. B.13) must be independent of time in a homogeneous, isotropic universe. It can be shown that homogeneity and isotropy also imply that the q_i can be chosen such that equation (1.43) takes the form of the **Robertson–Walker metric**,

$$dr^2 = a^2(t) \left[\frac{dx^2}{1 - kx^2/x_u^2} + x^2(d\theta^2 + \sin^2\theta d\phi^2) \right]. \quad (1.44)$$

Here θ and ϕ are the usual angles in spherical coordinates (Appendix B.2), x is a radial coordinate, x_u is a constant called the **radius of curvature**, and k is $+1$, 0 or -1 . Since x remains fixed as the fundamental observers recede from one another, it is called a **comoving coordinate**.

In the case $k = 0$, the metric (1.44) corresponds to ordinary spherical polar coordinates (cf. eq. B.32), so at a given cosmic time the geometry of the universe is that of ordinary Euclidean or **flat** space. In the case $k = +1$, (1.44) is the three-dimensional generalization of the metric on the surface of a sphere of radius x_u , $dr^2 = dx^2/(1 - x^2/x_u^2) + x^2d\phi^2$, where x is the perpendicular distance from the polar axis to the point in question. This case is said to represent a **closed universe**, since the volume of space is finite (Problem 1.7). The case $k = -1$ has no analogous 2-surface embedded in Euclidean 3-space (e.g., Weinberg 1972). It represents an **open universe** with infinite volume.

1.3.3 Dynamics

The evolution of the scale factor $a(t)$ is determined by the equations of general relativity and the equation of state of the material in the universe. We shall assume that all of the major components of this material can be described as (possibly relativistic) fluids. To derive the equations governing $a(t)$ we then need only one result from relativity: that a fluid with inertial mass density ρ and pressure p has a gravitational mass density (Problem 9.5)

$$\rho' = \rho + \frac{3p}{c^2}. \quad (1.45)$$

By isotropy, the universe is spherically symmetric as viewed by any fundamental observer. Now draw a sphere of radius r around such an observer, where r is large enough that the approximation of homogeneity and isotropy is valid, but small enough that Newtonian physics applies within it. As shown at the beginning of this section, in practice this means $100 \text{ Mpc} \ll r \ll 4000 \text{ Mpc}$. By analogy with Newton's famous theorem that a body experiences no gravitational force from a spherical shell of matter outside it (§2.2.1), we ignore the effects of material outside the sphere. Then Newton's law of gravity tells us that a fundamental observer on the surface of the sphere is accelerated towards its center at a rate

$$\frac{d^2 r}{dt^2} = -\frac{GM}{r^2}, \quad (1.46)$$

where M is the gravitational mass inside the sphere. Note that there are no pressure forces, since $\nabla p = 0$ by homogeneity. Since $M = \rho' V$, where $V = \frac{4}{3}\pi r^3$, and $r = r_0 a(t)$, we may rewrite equation (1.46) as

$$\frac{\ddot{a}}{a} = -\frac{4\pi G \rho'}{3} = -\frac{4\pi G}{3} \left(\rho + \frac{3p}{c^2} \right). \quad (1.47)$$

To integrate this equation, we need to know how p and ρ vary with the scale factor $a(t)$. The internal energy of the sphere, including its rest-mass energy, is $U = \rho c^2 V$. The material satisfies $dU + p dV = 0$ (eq. F.22 with $dS = 0$, since there is no heat flow in a homogeneous, isotropic universe), so

$$c^2 d(\rho V) + p dV = 0 \quad \text{or} \quad d\rho + \left(\rho + \frac{p}{c^2} \right) \frac{dV}{V} = 0. \quad (1.48)$$

Since $V \propto a^3(t)$, we have $dV/V = 3 da/a$, and equations (1.47) and (1.48) can be combined to eliminate p :

$$\frac{\ddot{a}}{a} = \frac{4\pi G}{3} \left(2\rho + a \frac{d\rho}{da} \right). \quad (1.49)$$

After multiplying by $a\dot{a}$ this equation can be integrated to yield

$$\dot{a}^2 - \frac{8\pi G\rho}{3}a^2 = 2E, \quad (1.50)$$

where E is a constant of integration, analogous to the Newtonian energy.

These equations can also be derived directly from general relativity. The relativistic derivation also connects the geometry to the energy density, by relating the parameters of the Robertson–Walker metric (1.44) to the integration constant E :

$$2E = -\frac{kc^2}{x_{\text{u}}^2}. \quad (1.51)$$

Equations (1.44), (1.50), and (1.51) specify the **Friedmann–Robertson–Walker** or **FRW** model of the universe.

When $k = 0$, space is flat, $E = 0$, and the density equals the **critical density**

$$\rho_{\text{c}}(t) \equiv \frac{3\dot{a}^2}{8\pi G a^2} = \frac{3H^2(t)}{8\pi G}. \quad (1.52)$$

If we define the **density parameter**

$$\Omega(t) \equiv \frac{\rho(t)}{\rho_{\text{c}}(t)} = \frac{8\pi G\rho(t)}{3H^2(t)}, \quad (1.53)$$

then equation (1.50) can be written

$$\Omega^{-1} - 1 = \frac{3E}{4\pi G\rho a^2} = -\frac{3kc^2}{8\pi G\rho a^2 x_{\text{u}}^2}. \quad (1.54)$$

This result implies that if $\Omega < 1$ at any time, it always remains so; this case corresponds to a universe that is open ($k = -1$) and infinite. In contrast, if Ω exceeds unity it always remains so, and we have a universe that is closed ($k = +1$) and finite. Finally, if $\Omega = 1$ at any instant it is unity for all time, and the universe is always flat. Thus the geometry of the universe is determined by its mass content, and an open universe cannot turn into a closed one or vice versa.

The present value of the critical density is

$$\rho_{\text{c}0} = \rho_{\text{c}}(t_0) = 9.204 \times 10^{-27} h_7^2 \text{ kg m}^{-3} = 1.3599 \times 10^{11} h_7^2 \mathcal{M}_{\odot} \text{ Mpc}^{-3}. \quad (1.55)$$

We have parametrized galaxies and other stellar systems by their mass-to-light ratio. Since the universe is homogeneous, the mass-to-light ratio measured on sufficiently large scales must be the same everywhere at a given cosmic time. The present R -band luminosity density j_R is given by equation (1.19), so the R -band mass-to-light ratio Υ_R is related to the density parameter by

$$\Upsilon_R = \frac{\rho_{\text{c}0}\Omega_0}{j_R} = (900 \pm 300) h_7 \Omega_0 \Upsilon_{\odot}; \quad (1.56)$$

the subscript “0” on Ω indicates the density parameter at the present epoch.

An obvious next step is to compare this result to observed mass-to-light ratios and thereby estimate Ω_0 . Unfortunately, the total mass-to-light ratios of individual galaxies are quite uncertain, because the total mass contained in their dark halos is difficult to determine. However, as we argued after equation (1.26), it is likely that the mixture of baryonic and non-baryonic dark matter in rich clusters of galaxies is representative of the universe as a whole, so we might hope that the mass-to-light ratios of rich clusters are a reasonable approximation to the total mass-to-light ratios of galaxies. Taking $\Upsilon_R \simeq (200 \pm 50) h_7 \Upsilon_\odot$ from equation (1.25), we conclude that $\Omega_0 = 0.22 \pm 0.09$. This argument is subject to at least two possible biases: first, the galaxy population in rich clusters has a higher fraction of ellipticals than average, and hence fewer luminous young stars; second, most of the baryons in rich clusters are in the form of hot gas, which does not contribute to the R -band luminosity, but in isolated galaxies this gas might cool to form additional stars. Both of these effects should increase the mass-to-light ratio in clusters relative to isolated galaxies, and hence lead us to overestimate Ω_0 . Thus we can conclude from this argument only that

$$\Omega_{m0} \lesssim 0.3. \quad (1.57)$$

The subscript “m” on Ω_0 is a reminder that this value refers only to the density in non-relativistic matter that can collapse along with the baryons into clusters of galaxies. Any uniformly distributed component of the density, such as a population of relativistic particles or vacuum energy, is excluded.

The inequality (1.57) encapsulates one of the fundamental conclusions of modern cosmology: the most “natural” model, a matter-dominated flat universe in which $\Omega = 1$ at all times, is excluded by the observations.

To solve the differential equations that describe FRW models, we need to know the equation of state relating the pressure p to the density ρ for each component of the universe. For our purposes a sufficiently general parametrization is

$$p = w\rho c^2, \quad (1.58)$$

where w is a constant. If the equation of state has this form, equation (1.48) can be integrated to yield

$$\rho \propto V^{-1-w} \propto a^{-3(1+w)}. \quad (1.59)$$

Three major components contribute to the dynamics of the universe:

- (i) Non-relativistic matter. This has $p \ll \rho c^2$ so $w = 0$. We label the corresponding density $\rho_m(t)$, and for brevity we simply call this component “matter.” In this case there is no distinction between the inertial mass density ρ_m and the gravitational density $\rho'_m = \rho_m$ (eq. 1.45). From equation (1.59) $\rho_m \propto a^{-3}$, as expected from conservation of mass.

- (ii) Radiation and other massless or highly relativistic particles. We label this density $\rho_\gamma(t)$, and call this component “radiation.” In this case $p = \frac{1}{3}\rho c^2$ so $w = \frac{1}{3}$, and from equation (1.45) the gravitational attraction is twice as strong as non-relativistic matter with the same density: $\rho'_\gamma = 2\rho_\gamma$. As the universe expands, the radiation density declines as $\rho_\gamma \propto a^{-4}$ from equation (1.59). Physically, this dependence arises because the number of photons is conserved so the number density declines as a^{-3} , and their frequency and thus the energy per photon decay as a^{-1} (eq. 1.42).
- (iii) A hypothetical energy density ρ_Λ associated with the vacuum (Carroll, Press, & Turner 1992). This must be accompanied by a *negative* pressure $p = -\rho_\Lambda c^2$ (i.e., a tension) because the energy-momentum tensor of the vacuum must be proportional to the Minkowski metric if the vacuum is to appear the same to all observers, regardless of their relative motion. In the parametrization of equation (1.58), vacuum energy therefore has $w = -1$. Equation (1.59) shows that as the universe expands, ρ_Λ is independent of the scale factor, as expected since it is a universal constant.

A remarkable feature of vacuum energy is that it exerts repulsive gravitational forces—equations (1.45) and (1.58) show that gravity is repulsive for any medium with $w < -\frac{1}{3}$. Consequently, the gravity from such a medium tends to accelerate rather than decelerate the expansion of the universe.

Vacuum energy plays a significant role in cosmology only if the vacuum-energy density ρ_Λ is comparable to the current critical density ρ_c (eq. 1.55). There is no motivation from fundamental physics for a vacuum-energy density of this magnitude: the theoretical prejudice is that either ρ_Λ has a very large value, or else is exactly zero on account of some unidentified symmetry. There is no known mechanism that would favor a value of ρ_Λ comparable to ρ_c . Moreover, because the critical density evolves with time while the vacuum-energy density does not, any approximate coincidence between ρ_Λ and ρ_c must be a special feature of the present epoch. These difficulties have led physicists to explore quantum fields with more general behavior than vacuum energy, under the general heading of **dark energy**.

By analogy with equation (1.53), we define

$$\Omega_{m0} \equiv \frac{\rho_{m0}}{\rho_{c0}} \quad ; \quad \Omega_{\gamma0} \equiv \frac{\rho_{\gamma0}}{\rho_{c0}} \quad ; \quad \Omega_{\Lambda0} \equiv \frac{\rho_{\Lambda0}}{\rho_{c0}} \quad (1.60)$$

to be the present densities of matter, radiation, and vacuum energy in units of the critical density. With this notation $\Omega_{m0} + \Omega_{\gamma0} + \Omega_{\Lambda0} = \Omega_0$. Then equation (1.50) can be rewritten as

$$\dot{a}^2 = H_0^2 [1 + \Omega_{m0}(a^{-1} - 1) + \Omega_{\gamma0}(a^{-2} - 1) + \Omega_{\Lambda0}(a^2 - 1)], \quad (1.61)$$

which can be integrated to yield a formula for the time dependence of the scale factor $a(t)$:

$$H_0 t = \int_0^{a(t)} \frac{a \, da}{\sqrt{\Omega_{\gamma 0} + \Omega_{m 0} a + (1 - \Omega_{m 0} - \Omega_{\gamma 0} - \Omega_{\Lambda 0}) a^2 + \Omega_{\Lambda 0} a^4}}. \quad (1.62)$$

This integral can be evaluated analytically or numerically for arbitrary values of $\Omega_{m 0}$, $\Omega_{\gamma 0}$, and $\Omega_{\Lambda 0}$, but it is more illuminating to examine special cases:

- (i) A flat, matter-dominated universe (the **Einstein-de Sitter universe**) has $\Omega_{\gamma 0} = \Omega_{\Lambda 0} = 0$, $\Omega_{m 0} = 1$, so

$$a(t) \propto t^{2/3} \quad ; \quad \rho_m(t) = \frac{1}{6\pi G t^2}; \quad (1.63)$$

the second equation follows when the first is substituted into equation (1.50) with $E = 0$.

- (ii) A flat, radiation-dominated universe has

$$a(t) \propto t^{1/2} \quad ; \quad \rho_\gamma(t) = \frac{3}{32\pi G t^2}. \quad (1.64)$$

- (iii) A flat universe dominated by vacuum energy has

$$a(t) \propto \exp(H_0 t) = \exp\left[\left(\frac{8}{3}\pi G \rho_\Lambda\right)^{1/2} t\right] \quad ; \quad \rho_\Lambda = \frac{3H_0^2}{8\pi G} = \text{constant}. \quad (1.65)$$

At the present time, $\Omega_{\gamma 0} \simeq 10^{-4}$ (eq. 1.72), so the evolution of the universe is determined by $\Omega_{m 0}$ and $\Omega_{\Lambda 0}$ except at very early times. Thus the properties of the universe can be parametrized on a diagram such as Figure 1.6. Lines on this figure mark the boundary between models that have open or closed geometries ($k = -1$ or $k = +1$), and between models that expand forever and those that collapse at some future time. We have also marked off models that have no initial singularity: as we follow these ‘‘bounce’’ models back in time from the present, the repulsion from the vacuum energy becomes so strong that the expansion rate \dot{a} reaches zero at some time t_b and is negative for $t < t_b$. Physically, this means that the universe was contracting for $t < t_b$, coasted to a halt because of increasing repulsion by the vacuum energy, and then began the expansion that continues at the present time. Such models are excluded by observations because they predict a maximum redshift, $z \simeq 2$, much smaller than the largest observed redshifts, $z \gtrsim 6$.

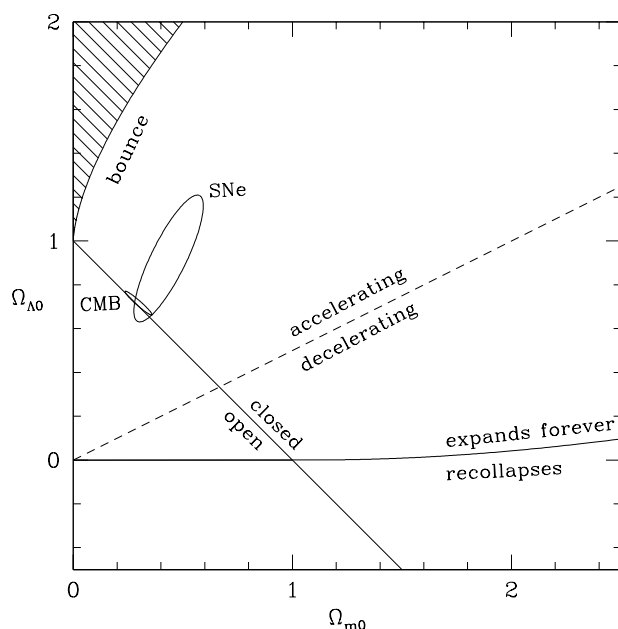


Figure 1.6 Characteristics of FRW models of the universe in which the current radiation density $\Omega_{\gamma 0}$ is negligible. The solutions are parametrized by the current matter density Ω_{m0} and vacuum energy $\Omega_{\Lambda 0}$, both relative to the critical density (1.52). The lines divide models in which the geometry is open from those with closed geometry; models that will expand forever from those that will eventually collapse; and models in which the expansion is accelerating ($\ddot{a} > 0$) from those in which it is decelerating. The shaded region denotes models with no initial singularity or Big Bang; these bounced from a collapsing to an expanding state at some non-zero value $a(t) < 1$ of the scale factor. The large oval marked SNe is the $1\text{-}\sigma$ error ellipse from measurements of distant supernovae (Riess et al. 2004), and the small oval labeled CMB is the $1\text{-}\sigma$ error ellipse from measurements of fluctuations in the cosmic microwave background combined with local measurements of the Hubble constant (Spergel et al. 2007). For further discussion see page 50.

1.3.4 The Big Bang and inflation

There is strong evidence that the universe was much hotter and denser in the past—this evidence includes the existence of the cosmic microwave background and the primordial abundances of the light elements (see §1.3.5). This is consistent with the discussion of the preceding paragraph, which shows that all FRW models that are consistent with observations begin from an initial singularity or **Big Bang**. Immediately after the Big Bang the universe satisfied several striking constraints:

- (i) A matter- or radiation-dominated FRW universe always evolves away from $\Omega = 1$: equation (1.54) shows that $|\Omega^{-1} - 1|$ grows in proportion to $1/\rho a^2$, which grows as $a(t)$ or $a^2(t)$, respectively. At present Ω is not

far from unity, so soon after the Big Bang, when $a(t)$ was much smaller than now, Ω must have been extremely close to unity. This fine-tuning of the initial conditions is called the “flatness problem.”

- (ii) The universe is homogeneous on large scales. It is natural to ask whether this property could be the result of physical processes occurring shortly after the Big Bang. Since information cannot propagate faster than the speed of light, the size of the largest causally connected region is given by the distance a photon can propagate since the Big Bang. Light travels at speed $c = dr/dt$, where dr is the distance. In the Robertson–Walker metric (1.44), the comoving coordinate of a photon that is moving towards the origin therefore satisfies

$$\frac{dx}{dt} = -\frac{c}{a(t)} \left(1 - \frac{kx^2}{x_u^2}\right)^{1/2}. \quad (1.66)$$

We have seen that in the early universe, $|\Omega - 1|$ must have been very small, so the geometry is nearly flat and we may set $k = 0$. Thus a photon that is emitted at t_i and arrives at the origin at t has come from a comoving coordinate

$$x = c \int_{t_i}^t \frac{dt}{a(t)}. \quad (1.67)$$

The comoving radius of the region that has been in causal contact since the Big Bang is called the **particle horizon** $x_h(t)$, and is obtained from equation (1.67) by letting t_i shrink to zero. At early times the universe is expected to be radiation-dominated, since $\rho_\gamma(t) \propto a^{-4}(t)$ while $\rho_m \propto a^{-3}$ and ρ_Λ is constant. In a flat radiation-dominated universe, the scale factor is $a(t) = bt^{1/2}$, where b is a constant (eq. 1.64). Thus we obtain

$$x_h(t) = \frac{2ct^{1/2}}{b}, \quad (1.68)$$

which shows that the comoving horizon shrinks to zero as $t \rightarrow 0$; this in turn implies that right after the Big Bang different parts of the universe are not causally connected, so the large-scale homogeneity of the universe must be imposed as an initial condition. This uncomfortable situation is called the “horizon problem.”

- (iii) The rich structure in the present universe—galaxies, clusters of galaxies, etc.—has grown by gravitational instability. The spectrum of perturbations that seeded this growth must be inserted as an initial condition (the “structure problem”).

Remarkably, all three of these problems can be resolved at one stroke by a single powerful assumption: that the early universe underwent a phase of accelerated expansion or **inflation**. Inflation arises when the dominant contributor to the mass density has an equation of state with $w < -\frac{1}{3}$ (eq. 1.58).

For example, suppose that there is an early inflationary phase in which the universe is dominated by a large vacuum energy, which has $w = -1$. Then the scale factor grows exponentially with time, as described by equation (1.65), but with much larger constants H_0 and ρ_Λ . As the universe inflates, the density of non-relativistic matter falls as $a^{-3}(t)$, and the density of photons and other relativistic matter falls as $a^{-4}(t)$, while the vacuum-energy density remains constant. Hence the dynamics rapidly becomes completely dominated by the vacuum energy, equation (1.65) applies with ever greater precision, and the density parameter Ω tends to unity. More precisely, if the inflationary phase lasts for n e-foldings of the scale factor, then at the end of this phase $|\Omega - 1|$ will be smaller by a factor $\sim \exp(-2n)$ than it was at the beginning. Thus inflation naturally produces Ω very close to unity, thereby solving the flatness problem. Moreover, since $|\Omega - 1|$ is zero at the end of inflation with exponential precision, it is plausible to assume that it is exactly zero for all practical purposes, even today. Thus inflation strongly suggests that $\Omega_0 = 1$; in other words, $\Omega_{\gamma 0} + \Omega_{m0} + \Omega_{\Lambda 0} = 1$. Since $\Omega_{\gamma 0}$ is negligible, the universe must lie on the line $\Omega_{m0} + \Omega_{\Lambda 0} = 1$ that separates open from closed universes in Figure 1.6, just as the observations seem to indicate.

Next, inflation solves the horizon problem: equation (1.65) tells us that $\dot{a} \propto a$, so the integral in equation (1.67) is $\int dt/a = \int da/(a\dot{a}) \propto \int da/a^2$, which diverges as $a \rightarrow 0$. Thus the region that is in causal contact becomes arbitrarily large if the inflationary phase begins early enough, so all regions in the observable universe could have been in causal contact in the interval between the Big Bang and the onset of inflation.

Inflation also predicts that when quantum fluctuations in the matter density are inflated past the event horizon, they are frozen into classical density fluctuations that provide the initial conditions for the growth of structure in the universe. These fluctuations enjoy a number of properties that greatly simplify the study of structure formation in the later universe: (i) they are nearly scale-invariant, in the sense that the RMS fluctuations in the gravitational potential are independent of scale; (ii) they form a Gaussian random field; (iii) they are adiabatic, that is, the entropy per particle is constant (see §9.1 for further discussion of these concepts).

Finally, inflation solves the **monopole problem**: point topological defects known as monopoles arise naturally in grand unified theories of particle physics, and the predicted density of these objects is far larger than allowed by observational constraints. Inflation solves this problem by diluting the density of monopoles—like non-relativistic matter, this density falls as $a^{-3}(t)$ —to an undetectably small value.

Particle physics has no difficulty embracing particle fields that could have caused inflation. Inflation is thought to end when a phase transition converts the inflating matter to ordinary matter and radiation (in this context, “ordinary” includes WIMPs or other non-baryonic dark matter). Any

matter or radiation present before inflation was diluted to negligible densities by this time. The newly formed matter and radiation will be in thermal equilibrium, with density parameter $\Omega = 1$ and density fluctuations that are Gaussian, adiabatic, and nearly scale-invariant—precisely the conditions we need to explain many of the properties of the observed universe.

Thus the inflationary hypothesis not only solves the horizon, flatness, and monopole problems but also provides simple, well-defined initial conditions that enable quantitative predictions about the growth of structure in the universe. It should be kept in mind, however, that despite its central role in modern cosmology, the inflationary hypothesis is still unsupported by any evidence from other arenas in theoretical or experimental physics.

1.3.5 The cosmic microwave background

Following inflation, the universe was radiation-dominated, so the scale factor grew as $a(t) \propto t^{1/2}$ and the density fell as $\rho \propto a^{-4}$ (eq. 1.64). All particle species were in thermal equilibrium at a temperature $T(t)$, which declined approximately as a^{-1} . Once the temperature dropped below 100 MeV, 10^{-4} s after the Big Bang, the constituents of this hot plasma consisted of relativistic electrons, positrons, neutrinos, and photons, and non-relativistic protons, neutrons, and perhaps WIMPs.

As the universe continued to expand and cool, the collision time between particles grew faster than the expansion time, so particles began to drift out of thermal equilibrium. In particular, weakly interacting particles such as neutrinos dropped out of thermal equilibrium at $T \simeq 10^{10}$ K = 0.86 MeV, about 1 s after the Big Bang. At this point the neutron/proton ratio, which had been kept in equilibrium by weak interactions, was frozen at about 0.2. The free neutrons then began to decay, with e-folding time 886 s. However, long before this decay process was completed, nucleosynthesis began: below 10^9 K, at $t \simeq 100$ s, $k_B T$ was much smaller than the deuteron binding energy, so deuterium began to accumulate. Eventually its abundance grew large enough for deuterium to burn to tritium and then helium. Nucleosynthesis was essentially complete 200 s after the Big Bang, leaving most of the nucleons as hydrogen (75% by mass) or ^4He (25% by mass), with traces of deuterium, ^3He and ^7Li . The abundance of deuterium, in particular, depends sensitively on the density of baryons at a given temperature, which is determined by the baryon-to-photon ratio η . Thus measurements of the abundance of deuterium in primordial astrophysical systems such as intergalactic clouds can be used to determine η ; measurements of the cosmic microwave background radiation (see below) determine the current photon density, and from these two quantities we can determine the baryon density. This method yields a current density parameter for baryons (Yao et al. 2006)

$$\Omega_{\text{b}0} = (0.042 \pm 0.004) h_7^{-2}. \quad (1.69)$$

As the universe expanded further, the density in radiation and relativistic matter (photons and neutrinos) continued to decline as a^{-4} , while the density in non-relativistic matter (mostly protons, electrons, and helium nuclei) declined as a^{-3} . Eventually, at redshift and time

$$z_{\gamma m} \simeq 3100, \quad t_{\gamma m} \simeq 6 \times 10^4 \text{ yr}, \quad (1.70)$$

the density of matter equaled that of radiation (see Figure 1.7). At this point the matter was still fully ionized. However, as the universe continued to expand, at

$$z_d \simeq 1100, \quad t_d \simeq 4 \times 10^5 \text{ yr}, \quad (1.71)$$

the electrons and protons combined to form neutral atomic hydrogen. This **decoupling** or **recombination epoch** was a milestone in the history of the universe for two reasons: (i) before recombination, the ionized baryonic plasma was locked to the photons by Thomson scattering, while after decoupling, the baryonic matter could move relative to the photons so the assembly of bound baryonic structures such as galaxies could begin; (ii) the universe became transparent.

Recombination occurred rather quickly: the fractional RMS dispersion in the scale factor at which photons suffer their last scattering is less than 10%. Thus we can imagine ourselves to be surrounded by an opaque **last-scattering surface** that hides the Big Bang from us.¹²

At recombination, the photons had a black-body spectrum, and this spectrum was preserved even after the universe became transparent: the photon frequencies all decline as $a^{-1}(t)$ (eq. 1.42) so the spectrum remained black-body with a temperature that declined in the same way.

The relic black-body radiation from the last-scattering surface, the **cosmic microwave background** or **CMB**, was discovered in 1965. It dominates the night sky at wavelengths in the range millimeters to centimeters. The spectrum is accurately black-body, with temperature $T = (2.725 \pm 0.001) \text{ K}$, and this finding provides compelling evidence that the universe arose from a hot, dense initial state—it is only in such a state that the photons can be thermalized in less than a Hubble time. The energy density of the CMB photons corresponds to a density parameter $5.04 \times 10^{-5} h_7^{-2}$; to compute the total density in radiation we must add to this the energy density contributed by relic neutrinos (this cosmic neutrino background has not yet been detected, but is a firm prediction of Big Bang cosmology). Thus the current energy density in radiation is

$$\Omega_{\gamma 0} = 8.48 \times 10^{-5} h_7^{-2}. \quad (1.72)$$

¹² The last-scattering surface is analogous to a stellar photosphere, except we are in a cavity in the middle of optically thick material rather than outside a sphere of optically thick material.

The CMB is remarkably close to isotropic: apart from a dipole term that arises from the velocity of the solar system with respect to a fundamental observer ($368 \pm 2 \text{ km s}^{-1}$), the largest fractional RMS anisotropies are $\lesssim 10^{-4}$. These are believed to arise from the primordial fluctuations introduced by inflation, and the power spectrum of these anisotropies provides an exquisitely sensitive probe of many of the fundamental parameters of the universe.

In particular, assuming a FRW universe currently dominated by non-relativistic matter and vacuum energy, and $\Omega_0 = 1$ as predicted by inflation, the power spectrum of CMB fluctuations strongly constrains the Hubble constant, the density in baryons, and the matter density (Spergel et al. 2007):

$$h_7 = 1.05 \pm 0.05, \quad \Omega_{\text{b}0} = (0.0455 \pm 0.0015) h_7^{-2}, \quad \Omega_{\text{m}0} = 0.237 \pm 0.034. \quad (1.73)$$

The vacuum-energy density is then $\Omega_{\Lambda 0} = 1 - \Omega_{\text{m}0} = 0.763 \pm 0.034$. This result implies that the dynamics of the universe is currently dominated by vacuum energy—in other words, the universe appears to have entered a second period of inflation.

With these parameters, the density in non-relativistic matter is much larger than the density in baryons. Thus there must be non-baryonic dark matter that contributes roughly 19% of the critical density. Note also that the present density in vacuum energy exceeds the density in matter; since the former is independent of scale factor and the latter scales as $a^{-3} = (1+z)^3$, equality occurred quite recently (Figure 1.7), at

$$z_{\text{m}\Lambda} = \left(\frac{\Omega_{\Lambda 0}}{\Omega_{\text{m}0}} \right)^{1/3} - 1 = 0.48 \pm 0.09. \quad (1.74)$$

The fluctuation spectrum deduced from the CMB measurements is also approximately scale-invariant, again as predicted by inflation.

The parameters in equation (1.73) are consistent with a wide variety of astronomical measurements, including the following: (i) The Hubble constant is consistent with the best direct estimate of the distance scale, using Cepheid distances, which yields $h_7 = 1.03 \pm 0.11$ (eq. 1.14). Moreover, if measurements of the CMB power spectrum are combined with this measurement of the Hubble constant, then the assumption that $\Omega_0 = 1$ can be eliminated, and the data yield $\Omega_0 = 1.014 \pm 0.017$, consistent with the prediction of inflation that $\Omega = 1$. (ii) The fraction of the total mass composed of baryons is

$$f_{\text{b}} = \frac{\Omega_{\text{b}0}}{\Omega_{\text{m}0}} = 0.17 \pm 0.01, \quad (1.75)$$

in reasonable agreement with the estimate from clusters of galaxies, 0.13 ± 0.02 (eq. 1.26). (ii) The baryon density $\Omega_{\text{b}0}$ is consistent with the result from nucleosynthesis, equation (1.69). (iii) The matter density $\Omega_{\text{m}0}$ is consistent with measurements of the geometry of the universe from supernovae, which

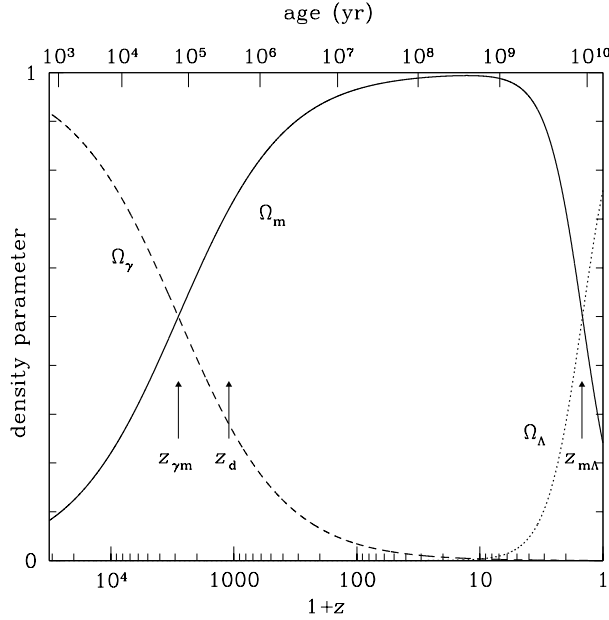


Figure 1.7 The fraction of the critical density provided by radiation (Ω_γ), matter (Ω_m), and vacuum energy (Ω_Λ) as a function of redshift z or scale factor $a(t) = (1+z)^{-1}$ (bottom axis) or age (top axis). The values shown are for a flat FRW universe with present matter density $\Omega_{m0} = 0.24$, vacuum-energy density $\Omega_{\Lambda0} = 0.76$ (eq. 1.73), and radiation density $\Omega_{\gamma0}$ determined from the temperature of the CMB (eq. 1.72). The redshifts of equal densities of matter and radiation (eq. 1.70), decoupling (eq. 1.71), and equal densities of matter and vacuum energy (eq. 1.74) are marked.

yield $\Omega_{m0} = 0.29 \pm 0.04$ for a flat universe (see Figure 1.6 and Riess et al. 2004). (iv) The matter density Ω_{m0} , together with equation (1.56), predicts that the average mass-to-light ratio in the universe is

$$\Upsilon_R = (220 \pm 80)\Upsilon_\odot, \quad (1.76)$$

consistent with the observed mass-to-light ratio of large clusters of galaxies, $\Upsilon_R = (210 \pm 50)\Upsilon_\odot$. (v) Given Ω_{m0} , $\Omega_{\Lambda0}$, and $\Omega_{\gamma0}$ (eq. 1.72), the temporal evolution of the scale factor is entirely determined by equation (1.61) (see Figure 1.7). The age of the universe is thus found to be

$$t_0 = (13.73 \pm 0.16) \text{ Gyr}, \quad (1.77)$$

consistent with the ages of globular-cluster stars (§1.1.1) and radioactive dating of the oldest stars (Cayrel et al. 2001).

This approximate but encouraging agreement among quite different methods of measuring the same cosmological parameters strongly suggests that we live in a universe with a flat geometry containing three main components: vacuum energy or some other field with a similar equation of state ($\sim 76\%$); non-baryonic dark matter ($\sim 20\%$), and baryons ($\sim 4\%$)—for a detailed inventory see Fukugita, Hogan, & Peebles (1998) and Fukugita & Peebles (2004).

There are unsatisfying aspects to this picture. In particular, there is no independent experimental evidence or strong theoretical justification for three of the central ingredients of this cosmological model: vacuum energy, the particle(s) that comprise the non-baryonic dark matter, and the field that drives inflation. There is also no explanation of why we happen to live at the special epoch when the densities in matter and vacuum energy are similar. Much work remains to be done.

Problems

1.1 [2] In principle, the density of matter in the solar neighborhood can be measured from its effects on planetary orbits. Assume that the solar system is permeated by a uniform medium of density $0.1 M_{\odot} \text{pc}^{-3}$ (Table 1.1). Estimate the rate of precession of the perihelion of Neptune (orbital radius 4.5×10^{12} m) due to the perturbing force from this medium, and compare your result to the minimum measurable precession, which is ≈ 0.01 arcsec yr^{-1} . An answer to within an order of magnitude is sufficient.

1.2 [2] (a) The **luminosity density** $j(\mathbf{r})$ of a stellar system is the luminosity per unit volume at position \mathbf{r} . For a transparent spherical galaxy, show that the surface brightness $I(R)$ (Box 2.1) and luminosity density $j(r)$ are related by the formula

$$I(R) = 2 \int_R^{\infty} dr \frac{rj(r)}{\sqrt{r^2 - R^2}}. \quad (1.78)$$

(b) What is the surface brightness of a transparent spherical galaxy with luminosity density $j(r) = j_0(1 + r^2/b^2)^{-5/2}$ (this is the Plummer model of §2.2.2c)?

(c) Invert equation (1.78) using Abel's formula (eq. B.72) to obtain

$$j(r) = -\frac{1}{\pi} \int_r^{\infty} \frac{dR}{\sqrt{R^2 - r^2}} \frac{dI}{dR}. \quad (1.79)$$

(d) Determine numerically the luminosity density in a spherical galaxy that follows the $R^{1/4}$ surface-brightness law. Plot $\log_{10} j(r)$ versus $\log_{10} r/R_e$, where R_e is the effective radius.

1.3 [2] The **strip brightness** $S(x)$ is defined so that $S(x) dx$ is the total luminosity in a strip of width dx that passes a distance x from the projected center of the system.

(a) Show that in a transparent, spherical system

$$S(x) = 2 \int_x^{\infty} dR \frac{RI(R)}{\sqrt{R^2 - x^2}}, \quad (1.80)$$

where $I(R)$ is the surface brightness at radius R .

(b) Show that the luminosity density and the total luminosity interior to r are related to the strip brightness by (Plummer 1911)

$$j(x) = -\frac{1}{2\pi x} \frac{dS}{dx} \quad ; \quad L(r) = -2 \int_0^r dx x \frac{dS}{dx}. \quad (1.81)$$

1.4 [2] An axisymmetric transparent galaxy has luminosity density that is constant on spheroids $R^2 + z^2/q^2$ having axis ratio q . A distant observer located on the symmetry axis of the galaxy sees an image with circular isophotes and central surface brightness I_n . A second distant observer, observing the galaxy from a line of sight that is inclined by an angle i to the symmetry axis, sees an image with elliptical isophotes with axis ratio $Q < 1$ and central surface brightness I_0 .

(a) What is the relation between I_0 , I_n , and Q ? Hint: the answers are different for oblate ($q < 1$) and prolate ($q > 1$) galaxies.

(b) What is the relation between q , Q , and i ?

(c) Assuming that galaxies are oriented randomly, what fraction are seen from a line of sight that lies within 10° of the symmetry axis? From within 10° of the equatorial plane?

1.5 [1] (a) Why is the estimated mass-to-light ratio of clusters of galaxies, equation (1.25), proportional to the assumed value of the Hubble constant h_7 ?

(b) Dark matter was discovered by Zwicky (1933), who compared the mass-to-light ratio of the Coma cluster of galaxies (as measured by the virial theorem, §4.8.3) with the mass-to-light ratios of the luminous parts of spiral galaxies as measured by circular-speed curves, and concluded that there was 400 times as much dark matter as luminous matter in the Coma cluster. However, Zwicky's conclusion was based on a Hubble constant $H_0 = 558 \text{ km s}^{-1} \text{ Mpc}^{-1}$. How would his conclusion about the ratio of dark to luminous matter have been affected had he used the correct value of the Hubble constant, which is smaller by a factor of eight?

1.6 [1] (a) Associated with the vacuum-energy density ρ_Λ is the characteristic timescale $(G\rho_\Lambda)^{-1/2}$. What is its value for the cosmological parameters in equation (1.73), and what is its physical significance?

(b) Einstein's original formulation of general relativity included a contribution from vacuum energy, which he called the **cosmological constant** and parametrized by

$$\Lambda \equiv \frac{8\pi G\rho_\Lambda}{c^2}. \quad (1.82)$$

$\Lambda^{-1/2}$ is a characteristic length. What is its value for the cosmological parameters in equation (1.73)?

1.7 [2] Prove that the volume of a closed FRW universe is $2\pi^2 a^3(t) x_u^3$.

1.8 [1] Einstein proposed a static FRW universe, that is, one in which the scale factor $a(t) = a_0 = \text{constant}$ and the Hubble constant $H_0 = 0$.

(a) If the radiation density is negligible, prove that the matter density in this universe equals twice the vacuum-energy density.

(b) Suppose that the scale factor is perturbed from a_0 by a small amount, $a(t) = a_0 + \epsilon a_1(t)$ with $\epsilon \ll 1$. Show that $a_1(t)$ grows exponentially, so the static universe is unstable, and derive the growth rate.

1.9 [1] Assuming a flat FRW universe with parameters given by equation (1.73) at the present time, what is the value of the Hubble parameter in the distant future? If this is less than its present value H_0 , why is the universe said to be accelerating?

1.10 [1] Suppose that some of the dark matter is composed of iron asteroids of density $\rho = 8 \text{ g cm}^{-3}$ and radius r that are uniformly distributed throughout intergalactic space. If the density in this form is $\Omega_a = 0.01$, find an approximate lower limit on r from the condition that the universe is not opaque, i.e., that we can see distant quasars. Your answer need be correct only to within a factor of two or so.

1.11 [2] Reproduce Figure 1.6.

1.12 [2] Write a function to do the following task: given values of the Hubble constant H_0 , the density parameters Ω_{m0} , $\Omega_{\gamma0}$, $\Omega_{\Lambda0}$, and a specified redshift z , find the age of the universe at that redshift. For a flat universe with parameters given by equations (1.72) and (1.73), what is the age at redshifts $z = 1000$, $z = 1$, and $z = 0$?

1.13 [2] The redshift at which the densities in matter and radiation are equal is $z_{\gamma m}$. For a flat FRW universe containing matter, radiation, and vacuum energy, with parameters not too far from our own, prove that (a)

$$1 + z_{\gamma m} = \frac{\Omega_{m0}}{\Omega_{\gamma0}} = 1.18 \times 10^4 \Omega_{m0} h_7^2 \quad (1.83)$$

(hint: use eq. 1.72); (b) the age of the universe at $z_{\gamma m}$ is

$$t_{\gamma m} = \frac{2(2 - \sqrt{2})}{3H_0} \frac{\Omega_{\gamma0}^{3/2}}{\Omega_{m0}^2}; \quad (1.84)$$

(c) the comoving horizon at $z_{\gamma m}$ (eq. 1.67) is

$$x_{\gamma m} = 2(\sqrt{2} - 1) \frac{c}{H_0} \frac{\Omega_{\gamma0}^{1/2}}{\Omega_{m0}} = \frac{32.7 \text{ Mpc}}{\Omega_{m0} h_7^2}. \quad (1.85)$$

Evaluate $z_{\gamma m}$, $t_{\gamma m}$, and $x_{\gamma m}$ for the parameters of equation (1.73).

1.14 [1] The universe was opaque before decoupling at $z > z_d \simeq 1100$ (eq. 1.71) because the ionized baryonic plasma had a high optical depth to Thomson scattering. For $z < z_d$ the electrons and protons recombined to form neutral atoms and the universe became transparent. Somewhere between $z \sim 20$ and $z \sim 6$ high-energy photons from newly formed quasars reionized most of the intergalactic medium. Why is the universe not opaque for $z \lesssim 6$?

2

Potential Theory

Much of the mass of a galaxy resides in stars. To compute the gravitational potential of a large collection of stars, we should in principle simply add the point-mass potentials of all the stars together. Of course, this is not practicable for the $\approx 10^{11}$ stars in a typical galaxy, and for most purposes it is sufficient to model the potential by smoothing the mass density in stars on a scale that is small compared to the size of the galaxy, but large compared to the mean distance between stars. In particular, in §1.2.1 we saw that we obtain an excellent approximation to the orbit of a single star in a galaxy by treating the star as a test particle that moves in a smooth potential of this kind. In this chapter we show how the force field of such an idealized galaxy can be calculated.

The chapter is divided into eight sections. We start by reviewing some general results on potential theory, and then in §2.2 we specialize to discuss the simplest potentials, those of spherical bodies. In §2.3 we describe flattened density distributions that also have simple potentials. These special systems give us insight into the potentials of real galaxies and provide useful prototypes, but they are not adequate for accurate modeling. Therefore in §2.4 and §2.5 we describe a variety of general techniques for computing the potentials of aspherical bodies. The potentials of razor-thin disks form an important limiting case, which we discuss in §2.6. In §2.7 we use these results to examine the potential of the Milky Way. In §2.8 and §2.9 we describe how gravitational potentials are found in computer simulations of stellar systems.

The discussion in §§2.4 to 2.6 is rather mathematical, and readers who are willing to take a few results on trust may prefer to move straight from §2.3 to §2.7.

2.1 General results

Our goal is to calculate the force $\mathbf{F}(\mathbf{x})$ on a particle of mass m_s at position \mathbf{x} that is generated by the gravitational attraction of a distribution of mass $\rho(\mathbf{x}')$. According to Newton's inverse-square law of gravitation, the force $\mathbf{F}(\mathbf{x})$ may be obtained by summing the small contributions

$$\delta\mathbf{F}(\mathbf{x}) = Gm_s \frac{\mathbf{x}' - \mathbf{x}}{|\mathbf{x}' - \mathbf{x}|^3} \delta m(\mathbf{x}') = Gm_s \frac{\mathbf{x}' - \mathbf{x}}{|\mathbf{x}' - \mathbf{x}|^3} \rho(\mathbf{x}') d^3\mathbf{x}' \quad (2.1)$$

to the overall force from each small element of volume $d^3\mathbf{x}'$ located at \mathbf{x}' . Thus

$$\mathbf{F}(\mathbf{x}) = m_s \mathbf{g}(\mathbf{x}) \quad \text{where} \quad \mathbf{g}(\mathbf{x}) \equiv G \int d^3\mathbf{x}' \frac{\mathbf{x}' - \mathbf{x}}{|\mathbf{x}' - \mathbf{x}|^3} \rho(\mathbf{x}') \quad (2.2)$$

is the **gravitational field**, the force per unit mass.

If we define the **gravitational potential** $\Phi(\mathbf{x})$ by

$$\Phi(\mathbf{x}) \equiv -G \int d^3\mathbf{x}' \frac{\rho(\mathbf{x}')}{|\mathbf{x}' - \mathbf{x}|}, \quad (2.3)$$

and notice that

$$\nabla_{\mathbf{x}} \left(\frac{1}{|\mathbf{x}' - \mathbf{x}|} \right) = \frac{\mathbf{x}' - \mathbf{x}}{|\mathbf{x}' - \mathbf{x}|^3}, \quad (2.4)$$

we find that we may write \mathbf{g} as

$$\begin{aligned} \mathbf{g}(\mathbf{x}) &= \nabla_{\mathbf{x}} \int d^3\mathbf{x}' \frac{G\rho(\mathbf{x}')}{|\mathbf{x}' - \mathbf{x}|} \\ &= -\nabla\Phi, \end{aligned} \quad (2.5)$$

where for brevity we have dropped the subscript \mathbf{x} on the gradient operator ∇ .

The potential is useful because it is a scalar field that is easier to visualize than the vector gravitational field but contains the same information. Also, in many situations the easiest way to obtain \mathbf{g} is first to calculate the potential and then to take its gradient.

If we take the divergence of equation (2.2), we find

$$\nabla \cdot \mathbf{g}(\mathbf{x}) = G \int d^3\mathbf{x}' \nabla_{\mathbf{x}} \cdot \left(\frac{\mathbf{x}' - \mathbf{x}}{|\mathbf{x}' - \mathbf{x}|^3} \right) \rho(\mathbf{x}'). \quad (2.6)$$

Now

$$\nabla_{\mathbf{x}} \cdot \left(\frac{\mathbf{x}' - \mathbf{x}}{|\mathbf{x}' - \mathbf{x}|^3} \right) = -\frac{3}{|\mathbf{x}' - \mathbf{x}|^3} + \frac{3(\mathbf{x}' - \mathbf{x}) \cdot (\mathbf{x}' - \mathbf{x})}{|\mathbf{x}' - \mathbf{x}|^5}. \quad (2.7)$$

When $\mathbf{x}' - \mathbf{x} \neq 0$ we may cancel the factor $|\mathbf{x}' - \mathbf{x}|^2$ from top and bottom of the last term in this equation to conclude that

$$\nabla_{\mathbf{x}} \cdot \left(\frac{\mathbf{x}' - \mathbf{x}}{|\mathbf{x}' - \mathbf{x}|^3} \right) = 0 \quad (\mathbf{x}' \neq \mathbf{x}). \quad (2.8)$$

Therefore, any contribution to the integral of equation (2.6) must come from the point $\mathbf{x}' = \mathbf{x}$, and we may restrict the volume of integration to a small sphere of radius h centered on this point. Since, for sufficiently small h , the density will be almost constant through this volume, we can take $\rho(\mathbf{x}')$ out of the integral. The remaining terms of the integrand may then be arranged as follows:

$$\begin{aligned} \nabla \cdot \mathbf{g}(\mathbf{x}) &= G\rho(\mathbf{x}) \int_{|\mathbf{x}' - \mathbf{x}| \leq h} d^3\mathbf{x}' \nabla_{\mathbf{x}} \cdot \left(\frac{\mathbf{x}' - \mathbf{x}}{|\mathbf{x}' - \mathbf{x}|^3} \right) \\ &= -G\rho(\mathbf{x}) \int_{|\mathbf{x}' - \mathbf{x}| \leq h} d^3\mathbf{x}' \nabla_{\mathbf{x}'} \cdot \left(\frac{\mathbf{x}' - \mathbf{x}}{|\mathbf{x}' - \mathbf{x}|^3} \right) \\ &= -G\rho(\mathbf{x}) \int_{|\mathbf{x}' - \mathbf{x}| = h} d^2\mathbf{S}' \cdot \frac{(\mathbf{x}' - \mathbf{x})}{|\mathbf{x}' - \mathbf{x}|^3}. \end{aligned} \quad (2.9a)$$

The last step in this sequence uses the divergence theorem to convert the volume integral into a surface integral (eq. B.43). Now on the sphere $|\mathbf{x}' - \mathbf{x}| = h$ we have $d^2\mathbf{S}' = (\mathbf{x}' - \mathbf{x})h d^2\Omega$, where $d^2\Omega$ is a small element of solid angle. Hence equation (2.9a) becomes

$$\nabla \cdot \mathbf{g}(\mathbf{x}) = -G\rho(\mathbf{x}) \int d^2\Omega = -4\pi G\rho(\mathbf{x}). \quad (2.9b)$$

If we substitute from equation (2.5) for $\nabla \cdot \mathbf{g}$, we obtain **Poisson's equation** relating the potential Φ to the density ρ ;

$$\nabla^2\Phi = 4\pi G\rho. \quad (2.10)$$

This is a differential equation that can be solved for $\Phi(\mathbf{x})$ given $\rho(\mathbf{x})$ and an appropriate boundary condition.¹ For an isolated system the boundary condition is $\Phi \rightarrow 0$ as $|\mathbf{x}| \rightarrow \infty$. The potential given by equation (2.3) automatically satisfies this boundary condition. Poisson's equation provides

¹ Using the physically correct boundary condition is essential: for example, if $\Phi(\mathbf{x})$ is a solution, then so is $\Phi(\mathbf{x}) + \mathbf{k} \cdot \mathbf{x}$, with \mathbf{k} an arbitrary constant vector, but the corresponding gravitational fields differ by \mathbf{k} .

a route to Φ and then to \mathbf{g} that is often more convenient than equations (2.2) or (2.3). In the special case $\rho = 0$ Poisson's equation becomes **Laplace's equation**,

$$\nabla^2\Phi = 0. \quad (2.11)$$

If we integrate both sides of equation (2.10) over an arbitrary volume containing total mass M , and then apply the divergence theorem (eq. B.43), we obtain

$$4\pi GM = 4\pi G \int d^3\mathbf{x} \rho = \int d^3\mathbf{x} \nabla^2\Phi = \int d^2\mathbf{S} \cdot \nabla\Phi. \quad (2.12)$$

This result is **Gauss's theorem**, which states that *the integral of the normal component of $\nabla\Phi$ over any closed surface equals $4\pi G$ times the mass contained within that surface.*

Since \mathbf{g} is determined by the gradient of a potential, the gravitational field is conservative, that is, the work done against gravitational forces in moving two stars from infinity to a given configuration is independent of the path along which they are moved, and is defined to be the potential energy of the configuration (Appendix D.1). Similarly, the work done against gravitational forces in assembling an arbitrary continuous distribution of mass $\rho(\mathbf{x})$ is independent of the details of how the mass distribution was assembled, and is defined to be equal to the **potential energy** of the mass distribution. An expression for the potential energy can be obtained by the following argument.

Suppose that some of the mass is already in place so that the density and potential are $\rho(\mathbf{x})$ and $\Phi(\mathbf{x})$. If we now bring in a additional small mass δm from infinity to position \mathbf{x} , the work done is $\delta m \Phi(\mathbf{x})$. Thus, if we add a small increment of density $\delta\rho(\mathbf{x})$, the change in potential energy is

$$\delta W = \int d^3\mathbf{x} \delta\rho(\mathbf{x})\Phi(\mathbf{x}). \quad (2.13)$$

According to Poisson's equation the resulting change in potential $\delta\Phi(\mathbf{x})$ satisfies $\nabla^2(\delta\Phi) = 4\pi G(\delta\rho)$, so

$$\delta W = \frac{1}{4\pi G} \int d^3\mathbf{x} \Phi \nabla^2(\delta\Phi). \quad (2.14)$$

Using the divergence theorem in the form (B.45), we may write this as

$$\delta W = \frac{1}{4\pi G} \int \Phi \nabla(\delta\Phi) \cdot d^2\mathbf{S} - \frac{1}{4\pi G} \int d^3\mathbf{x} \nabla\Phi \cdot \nabla(\delta\Phi), \quad (2.15)$$

where the surface integral vanishes because $\Phi \propto r^{-1}$, $|\nabla\delta\Phi| \propto r^{-2}$ as $r \rightarrow \infty$, so the integrand $\propto r^{-3}$ while the total surface area $\propto r^2$. But $\nabla\Phi \cdot \nabla(\delta\Phi) = \frac{1}{2}\delta(\nabla\Phi \cdot \nabla\Phi) = \frac{1}{2}\delta(|\nabla\Phi|^2)$. Hence

$$\delta W = -\frac{1}{8\pi G} \delta \left(\int d^3\mathbf{x} |\nabla\Phi|^2 \right). \quad (2.16)$$

If we now sum up all of the contributions δW , we have a simple expression for the potential energy,

$$W = -\frac{1}{8\pi G} \int d^3\mathbf{x} |\nabla\Phi|^2. \quad (2.17)$$

To obtain an alternative expression for W , we again apply the divergence theorem and replace $\nabla^2\Phi$ by $4\pi G\rho$ to obtain

$$W = \frac{1}{2} \int d^3\mathbf{x} \rho(\mathbf{x})\Phi(\mathbf{x}). \quad (2.18)$$

The potential-energy tensor In §4.8.3 we shall encounter the tensor \mathbf{W} that is defined by

$$W_{jk} \equiv - \int d^3\mathbf{x} \rho(\mathbf{x}) x_j \frac{\partial\Phi}{\partial x_k}, \quad (2.19)$$

where ρ and Φ are the density and potential of some body, and the integral is to be taken over all space. We now deduce some useful properties of \mathbf{W} , which is known as the **Chandrasekhar potential-energy tensor**.²

If we substitute for Φ from equation (2.3), \mathbf{W} becomes

$$W_{jk} = G \int d^3\mathbf{x} \rho(\mathbf{x}) x_j \frac{\partial}{\partial x_k} \int d^3\mathbf{x}' \frac{\rho(\mathbf{x}')}{|\mathbf{x}' - \mathbf{x}|}. \quad (2.20)$$

Since the range of the integration over \mathbf{x}' does not depend on \mathbf{x} , we may carry the differentiation inside the integral to find

$$W_{jk} = G \int d^3\mathbf{x} \int d^3\mathbf{x}' \rho(\mathbf{x}) \rho(\mathbf{x}') \frac{x_j(x'_k - x_k)}{|\mathbf{x}' - \mathbf{x}|^3}. \quad (2.21a)$$

Furthermore, since \mathbf{x} and \mathbf{x}' are dummy variables of integration, we may relabel them and write

$$W_{jk} = G \int d^3\mathbf{x}' \int d^3\mathbf{x} \rho(\mathbf{x}') \rho(\mathbf{x}) \frac{x'_j(x_k - x'_k)}{|\mathbf{x} - \mathbf{x}'|^3}. \quad (2.21b)$$

Finally, on interchanging the order of integration in equation (2.21b) and adding the result to equation (2.21a), we obtain

$$W_{jk} = -\frac{1}{2}G \int d^3\mathbf{x} \int d^3\mathbf{x}' \rho(\mathbf{x}) \rho(\mathbf{x}') \frac{(x'_j - x_j)(x'_k - x_k)}{|\mathbf{x}' - \mathbf{x}|^3}. \quad (2.22)$$

²Subramanyan Chandrasekhar (1910–1995) was educated in India and England and spent most of his career at the University of Chicago. He discovered the Chandrasekhar limit, the maximum mass of a white dwarf star, and elucidated the concept of dynamical friction in astrophysics (§8.1). He shared the 1983 Nobel Prize in Physics with W. A. Fowler.

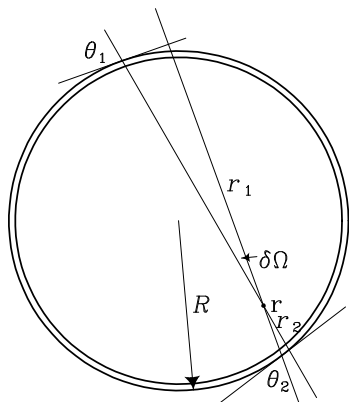


Figure 2.1 Proof of Newton's first theorem.

From this expression we draw the important inference that the tensor \mathbf{W} is **symmetric**, that is, that $W_{jk} = W_{kj}$. If the body is flattened along some axis, say the x_3 axis, W_{33} will be smaller than the other components because for most pairs of matter elements, $|x_3 - x'_3| < |x_1 - x'_1|$ or $|x_2 - x'_2|$.

When we take the **trace** of both sides of equation (2.22), we find

$$\begin{aligned} \text{trace}(\mathbf{W}) &\equiv \sum_{j=1}^3 W_{jj} = -\frac{1}{2}G \int d^3\mathbf{x} \rho(\mathbf{x}) \int d^3\mathbf{x}' \frac{\rho(\mathbf{x}')}{|\mathbf{x}' - \mathbf{x}|} \\ &= \frac{1}{2} \int d^3\mathbf{x} \rho(\mathbf{x}) \Phi(\mathbf{x}). \end{aligned} \quad (2.23)$$

Comparing this with equation (2.18) we see that $\text{trace}(\mathbf{W})$ is simply the total gravitational potential energy W . Taking the trace of (2.19) we have

$$W = - \int d^3\mathbf{x} \rho \mathbf{x} \cdot \nabla \Phi, \quad (2.24)$$

which provides another useful expression for the potential energy of a body.

2.2 Spherical systems

2.2.1 Newton's theorems

Newton proved two results that enable us to calculate the gravitational potential of any spherically symmetric distribution of matter easily:

Newton's first theorem *A body that is inside a spherical shell of matter experiences no net gravitational force from that shell.*

Newton's second theorem *The gravitational force on a body that lies outside a spherical shell of matter is the same as it would be if all the shell's matter were concentrated into a point at its center.*

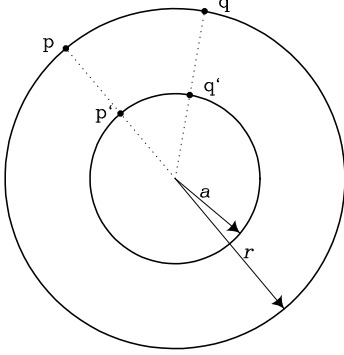


Figure 2.2 Proof of Newton's second theorem.

Figure 2.1 illustrates the proof of Newton's first theorem. Consider the cone associated with a small solid angle $\delta\Omega$ centered on the point \mathbf{r} . This cone intersects the spherical shell of matter at two points, at distances r_1 and r_2 from \mathbf{r} . Elementary geometrical considerations assure us that the angles θ_1 and θ_2 are equal, and therefore that the masses δm_1 and δm_2 contained within $\delta\Omega$ where it intersects the shell are in the ratio $\delta m_1/\delta m_2 = (r_1/r_2)^2$. Hence $\delta m_2/r_2^2 = \delta m_1/r_1^2$ and a particle placed at \mathbf{r} is attracted equally in opposite directions. Summing over all cones centered on \mathbf{r} , one concludes that the body at \mathbf{r} experiences no net force from the shell.<

An important corollary of Newton's first theorem is that the gravitational potential inside an empty spherical shell is constant because $\nabla\Phi = -\mathbf{g} = 0$. Thus we may evaluate the potential $\Phi(\mathbf{r})$ inside the shell by calculating the integral expression (2.3) for \mathbf{r} located at any interior point. The most convenient place for \mathbf{r} is the center of the shell, for then all points on the shell are at the same distance R , and one immediately has

$$\Phi = -\frac{GM}{R}. \quad (2.25)$$

The proof of his second theorem eluded Newton for more than ten years. Yet with hindsight it is easy. The trick (Figure 2.2) is to compare the potential Φ at a point \mathbf{p} located a distance r from the center of a spherical inner shell of mass M and radius a ($r > a$), with the potential Φ' at a point \mathbf{p}' located a distance a from the center of an outer shell of mass M and radius r . Consider the contribution $\delta\Phi$ to the potential at \mathbf{p} from the portion of the inner shell with solid angle $\delta\Omega$ located at \mathbf{q}' . Evidently

$$\delta\Phi = -\frac{GM}{|\mathbf{p} - \mathbf{q}'|} \frac{\delta\Omega}{4\pi}. \quad (2.26a)$$

But the contribution $\delta\Phi'$ of the matter in the outer shell near \mathbf{q} to the potential at \mathbf{p}' is

$$\delta\Phi' = -\frac{GM}{|\mathbf{p}' - \mathbf{q}|} \frac{\delta\Omega}{4\pi}. \quad (2.26b)$$

Finally, as $|\mathbf{p} - \mathbf{q}'| = |\mathbf{p}' - \mathbf{q}|$ by symmetry, it follows that $\delta\Phi = \delta\Phi'$, and then by summation over all points \mathbf{q} and \mathbf{q}' that $\Phi = \Phi'$. But we already know that $\Phi' = -GM/r$, therefore $\Phi = -GM/r$, which is exactly the potential that would be generated by concentrating the entire mass of the inner shell at its center.◁

Alternative proofs of Newton's theorems, using the machinery of spherical harmonics, are given in §2.4.

From Newton's first and second theorems, it follows that the gravitational attraction of a spherical density distribution $\rho(r')$ on a unit mass at radius r is entirely determined by the mass interior to r :

$$\mathbf{F}(r) = -\frac{GM(r)}{r^2}\hat{\mathbf{e}}_r, \quad (2.27a)$$

where

$$M(r) = 4\pi \int_0^r dr' r'^2 \rho(r'). \quad (2.27b)$$

The total gravitational potential may be considered to be the sum of the potentials of spherical shells of mass $dM(r) = 4\pi\rho(r)r^2dr$, so we may calculate the gravitational potential at \mathbf{r} generated by an arbitrary spherically symmetric density distribution $\rho(r')$ by adding the contributions to the potential produced by shells (i) with $r' < r$, and (ii) with $r' > r$. In this way we obtain

$$\begin{aligned} \Phi(r) &= -\frac{G}{r} \int_0^r dM(r') - G \int_r^\infty \frac{dM(r')}{r'} \\ &= -4\pi G \left[\frac{1}{r} \int_0^r dr' r'^2 \rho(r') + \int_r^\infty dr' r' \rho(r') \right]. \end{aligned} \quad (2.28)$$

It is worthwhile to check that the force $\mathbf{F} = -\nabla\Phi$ obtained from (2.28) recovers the simple expression (2.27a).

An important property of a spherical matter distribution is its **circular speed** $v_c(r)$, defined to be the speed of a particle of negligible mass (a **test particle**) in a circular orbit at radius r . We may readily evaluate v_c by equating the gravitational attraction $|\mathbf{F}|$ from equation (2.27a) to the centripetal acceleration v_c^2/r :

$$v_c^2 = r|\mathbf{F}| = r \frac{d\Phi}{dr} = \frac{GM(r)}{r}. \quad (2.29)$$

The associated angular frequency is called the **circular frequency**

$$\Omega \equiv \frac{v_c}{r} = \sqrt{\frac{GM(r)}{r^3}}. \quad (2.30)$$

The circular speed and frequency measure the mass interior to r . Another important quantity is the **escape speed** v_e defined by³

$$v_e(r) \equiv \sqrt{2|\Phi(r)|}. \quad (2.31)$$

A star at r can escape from the gravitational field represented by Φ only if it has a speed at least as great as $v_e(r)$, for only then does its (positive) kinetic energy $\frac{1}{2}v^2$ exceed the absolute value of its (negative) potential energy Φ . The escape speed at r depends on the mass both inside and outside r .

Potential energy of spherical systems The simplest expression for the potential energy of a spherical body is obtained from equation (2.24). Substituting equation (2.27a) and integrating over all directions of \mathbf{r} , we obtain

$$W = -4\pi G \int_0^\infty dr r \rho(r) M(r). \quad (2.32)$$

It is straightforward to show (see Problem 2.2) that the potential-energy tensor of a spherical body is **diagonal**, that is, $W_{jk} = 0$ for $j \neq k$, and has the form

$$W_{jk} = \frac{1}{3}W\delta_{jk}, \quad (2.33)$$

where δ_{ij} is unity if $i = j$ and zero otherwise. Such tensors are said to be **isotropic**.

2.2.2 Potentials of some simple systems

It is instructive to discuss the potentials generated by several simple density distributions:

(a) **Point mass** In this case

$$\Phi(r) = -\frac{GM}{r} \quad ; \quad v_c(r) = \sqrt{\frac{GM}{r}} \quad ; \quad v_e(r) = \sqrt{\frac{2GM}{r}}. \quad (2.34)$$

Potentials of this form, and orbits within them, are frequently referred to as **Keplerian** because Kepler first understood that $v_c \propto r^{-1/2}$ in the solar system.

(b) **Homogeneous sphere** If the density is some constant ρ , we have $M(r) = \frac{4}{3}\pi r^3 \rho$ and

$$v_c = \sqrt{\frac{4\pi G\rho}{3}}r. \quad (2.35)$$

³This result is correct only if the potential $\Phi(r) \rightarrow 0$ as $r \rightarrow \infty$ —we have assumed this so far but for systems with very extended mass distributions other zero points may be necessary (cf. eq. 2.62).

Thus in this case the circular speed rises linearly with radius, and the orbital period of a mass on a circular orbit is

$$T = \frac{2\pi r}{v_c} = \sqrt{\frac{3\pi}{G\rho}}, \quad (2.36)$$

independent of the radius of its orbit. The inverse of the angular frequency of a circular orbit is

$$\frac{r}{v_c} = \sqrt{\frac{3}{4\pi G\rho}} = 0.4886(G\rho)^{-1/2}. \quad (2.37)$$

If a small mass is released from rest at radius r in the gravitational field of a homogeneous sphere, its equation of motion is

$$\frac{d^2 r}{dt^2} = -\frac{GM(r)}{r^2} = -\frac{4\pi G\rho}{3}r, \quad (2.38)$$

which is the equation of motion of a harmonic oscillator of angular frequency $2\pi/T$. Therefore no matter what is the initial value of r , the test mass will reach $r = 0$ in a quarter of a period, or in a time

$$\sqrt{\frac{3\pi}{16G\rho}} = 0.767(G\rho)^{-1/2}. \quad (2.39)$$

The times in equations (2.37) and (2.39) are rather similar, and this suggests that the time taken for a particle to complete a significant fraction of its orbit is $\sim (G\rho)^{-1/2}$, independent of the size and shape of the orbit. This result also holds for inhomogeneous systems, so long as ρ is replaced by the mean density $\bar{\rho}$ interior to the particle's current radius. Thus we estimate the crossing time (sometimes also called the **dynamical time**) to be

$$t_{\text{cross}} \simeq t_{\text{dyn}} \simeq (G\bar{\rho})^{-1/2} \quad (2.40)$$

and shall use this as a measure of the characteristic time associated with the orbital motion of a star. Note that a complete orbital period is larger than t_{cross} by a factor $\simeq 2\pi \simeq 6$.

The potential energy of a homogeneous sphere of radius a and density ρ is conveniently obtained from equation (2.32). We have $M(r) = \frac{4}{3}\pi\rho r^3$, and therefore

$$W = -\frac{16\pi^2}{3}G\rho^2 \int_0^a dr r^4 = -\frac{16}{15}\pi^2 G\rho^2 a^5 = -\frac{3}{5}\frac{GM^2}{a}. \quad (2.41)$$

Sometimes it is useful to characterize the size of a system that lacks a sharp boundary by quoting the **gravitational radius**, which is defined as

$$r_g \equiv \frac{GM^2}{|W|}. \quad (2.42)$$

For a homogeneous sphere of radius a , $r_g = \frac{5}{3}a$.

From equation (2.28) the gravitational potential of a homogeneous sphere of radius a is

$$\Phi(r) = \begin{cases} -2\pi G\rho(a^2 - \frac{1}{3}r^2) & (r < a), \\ -\frac{4\pi G\rho a^3}{3r} & (r > a). \end{cases} \quad (2.43)$$

(c) Plummer model We might expect that in many spherical systems the density is roughly constant near the center, and falls to zero at large radii. The potential of a system of this type would be proportional to $r^2 + \text{constant}$ at small radii and to r^{-1} at large radii. A simple potential with these properties is the **Plummer model**

$$\Phi = -\frac{GM}{\sqrt{r^2 + b^2}}. \quad (2.44a)$$

The linear scale of the system that generates this potential is set by the **Plummer scale length** b , while M is the system's total mass.

From equation (B.53) for ∇^2 in spherical polar coordinates we have

$$\nabla^2\Phi = \frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Phi}{dr} \right) = \frac{3GMb^2}{(r^2 + b^2)^{5/2}}. \quad (2.45)$$

Thus from Poisson's equation (2.10) we have that the density corresponding to the potential (2.44a) is

$$\rho(r) = \frac{3M}{4\pi b^3} \left(1 + \frac{r^2}{b^2} \right)^{-5/2}. \quad (2.44b)$$

The potential energy of a Plummer model is

$$W = -\frac{3\pi GM^2}{32b}. \quad (2.46)$$

Plummer (1911) used the potential-density pair that is described by equations (2.44) to fit observations of globular clusters. We shall encounter it again in §4.3.3a as a member of the family of stellar systems known as polytropes.

(d) Isochrone potential The position of a star orbiting in a Plummer potential cannot be given in terms of elementary functions. However, in Chapter 3 we shall see that all orbits are analytic in the **isochrone potential**

$$\Phi(r) = -\frac{GM}{b + \sqrt{b^2 + r^2}}, \quad (2.47)$$

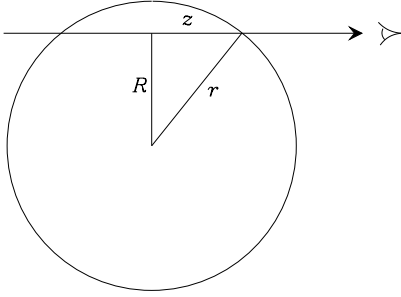


Figure 2.3 Projection of a spherical body along the line of sight.

which owes its name to a property of its orbits that we shall derive in §3.1c. By equation (2.29) we have for the circular speed at radius r

$$v_c^2(r) = \frac{GMr^2}{(b+a)^2a}, \quad (2.48a)$$

where

$$a \equiv \sqrt{b^2 + r^2}. \quad (2.48b)$$

When r is large $v_c \simeq \sqrt{GM/r}$, as required for a system of finite mass and extent. By Poisson's equation the density associated with the isochrone potential is

$$\rho(r) = \frac{1}{4\pi G} \frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Phi}{dr} \right) = M \left[\frac{3(b+a)a^2 - r^2(b+3a)}{4\pi(b+a)^3a^3} \right]. \quad (2.49)$$

Thus the central density is

$$\rho(0) = \frac{3M}{16\pi b^3}, \quad (2.50)$$

and at large radii the density tends to

$$\rho(r) \simeq \frac{bM}{2\pi r^4} \quad (r \gg b). \quad (2.51)$$

(e) Modified Hubble model The surface brightnesses of many elliptical galaxies may be approximated over a large range of radii by the **Hubble-Reynolds law**

$$I_H(R) = \frac{I_0}{(1 + R/R_H)^2}, \quad (2.52)$$

or the $R^{1/4}$ law (eq. 1.17 with $m = 4$). It is possible to solve for the spherical luminosity density $j(r)$ that generates a given axisymmetric brightness

Box 2.1: Definitions of surface brightness

Two definitions of surface brightness are used in the astronomical literature. Consider the radiative power that enters a telescope of aperture dA at normal incidence from directions in the element of solid angle $d^2\Omega$ due to the distribution of luminosity density $j(\mathbf{r})$. We take the coordinate origin to lie at the center of the telescope aperture. Then in the distance range $(r, r + dr)$ the volume that lies within the given solid-angle element is $r^2 dr d^2\Omega$, and the associated luminosity is $dL = j(\mathbf{r}) r^2 dr d^2\Omega$. A fraction $dA/(4\pi r^2)$ of this luminosity enters the telescope, so the power received in the given solid-angle element from the given distance range is $(4\pi)^{-1} j(\mathbf{r}) dr d^2\Omega dA$. Summing over all distances and dividing by $d^2\Omega dA$ we find that the flux per unit solid angle is

$$\hat{I} = \frac{1}{4\pi} \int dr j(\mathbf{r}).$$

The surface brightness \hat{I} might be reported in units of $\text{W m}^{-2} \text{sr}^{-1}$.

The second definition of surface brightness is the integral

$$I = \int dr j(\mathbf{r}).$$

I is the galaxy's luminosity per unit area when viewed from the given direction, and is conveniently given in units of solar luminosities per square parsec. Surface brightnesses are often reported by observers as so many magnitudes per square arcsecond, for example $20 \text{ mag arcsec}^{-2}$, meaning that in one square arcsecond of the image as much radiation is received as from a 20th magnitude star. In this book we follow the second convention and call I the "surface brightness."

distribution $I(R)$ (see Problem 1.2). However, the resulting formulae for the luminosity distribution of a galaxy that obeys either of these empirical laws are cumbersome (Hubble 1930). Fortunately, the simple luminosity density

$$j_{\text{h}}(r) = j_0 \left(1 + \frac{r^2}{a^2}\right)^{-3/2}, \quad (2.53)$$

where a is a constant, gives rise to a surface-brightness distribution that is similar to I_{H} (Rood et al. 1972). In fact, in the notation of Figure 2.3 we have that

$$I_{\text{h}}(R) = 2 \int_0^{\infty} dz j_{\text{h}}(r) = 2j_0 \int_0^{\infty} dz \left(1 + \frac{R^2}{a^2} + \frac{z^2}{a^2}\right)^{-3/2}. \quad (2.54)$$

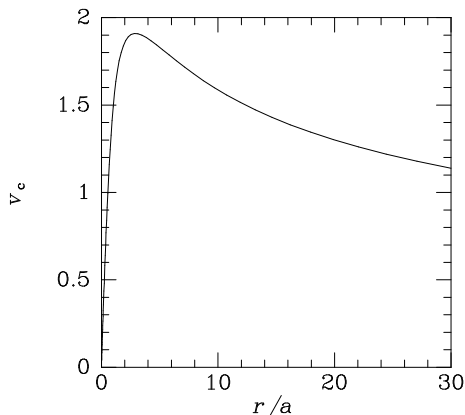


Figure 2.4 Circular speed versus radius for a body whose projected density follows the modified Hubble model (2.55). The circular speed v_c is plotted in units of $\sqrt{Gj_0\Upsilon a^2}$.

Using the substitution $y \equiv z/\sqrt{a^2 + R^2}$, we obtain the **modified Hubble model**

$$I_h(R) = \frac{2j_0a}{1 + R^2/a^2} \int_0^\infty \frac{dy}{(1 + y^2)^{3/2}} = \frac{2j_0a}{1 + R^2/a^2}. \quad (2.55)$$

Thus $I_h(a) = \frac{1}{2}I_h(0)$, so a is the core radius as defined on page 30. At large R , $I_h(R) \propto R^{-2}$, just as in the Hubble–Reynolds law (2.52). However, the Hubble–Reynolds law and the modified Hubble model behave quite differently near the center: while the luminosity density j_h is well behaved at the origin, in the Hubble–Reynolds law, $dI_H/dR \neq 0$ at $R = 0$, which implies that the luminosity density $j_H(r)$ diverges as $r \rightarrow 0$ (Hubble 1930).

Using equation (2.28) we can calculate the potential that would be generated by the modified Hubble model if its mass were distributed in the same way as its light. If $\rho(r) = \Upsilon j(r)$, where Υ is the constant mass-to-light ratio in the galaxy, one has

$$M_h(r) = 4\pi a^3 \Upsilon j_0 \left[\ln \left(\frac{r}{a} + \sqrt{1 + \frac{r^2}{a^2}} \right) - \frac{r}{a} \left(1 + \frac{r^2}{a^2} \right)^{-1/2} \right], \quad (2.56)$$

$$\Phi_h = -\frac{GM_h(r)}{r} - \frac{4\pi G \Upsilon j_0 a^2}{\sqrt{1 + (r/a)^2}}. \quad (2.57)$$

One feature of both the Hubble–Reynolds law and the modified Hubble model is that the mass diverges logarithmically at large r ; from equation (2.56) $M_h \simeq 4\pi a^3 \Upsilon j_0 [\ln(2r/a) - 1]$ for $r \gg a$. In practice, a galaxy must have a finite mass, so $j(r)$ must fall below $j_h(r)$ at sufficiently large r . Nevertheless, the potential Φ_h is finite, and in fact rather nearly equal to $-GM_h(r)/r$ whenever $r \gg a$. This behavior indicates that from the gravitational point of view the density distribution of equation (2.53) behaves much like a point mass at large radii. The circular speed is shown in Figure 2.4. It peaks at $r = 2.9a$ and then falls nearly as steeply as in the Keplerian case.

(f) Power-law density model Many galaxies have luminosity profiles that approximate a power law over a large range in radius. Consider the

structure of a system in which the mass density drops off as some power of the radius:

$$\rho(r) = \rho_0 \left(\frac{r_0}{r} \right)^\alpha. \quad (2.58)$$

The surface density of this system is

$$\Sigma(R) = \frac{\rho_0 r_0^\alpha}{R^{\alpha-1}} \frac{(-\frac{1}{2})! (\frac{\alpha-3}{2})!}{(\frac{\alpha-2}{2})!}. \quad (2.59)$$

We assume that $\alpha < 3$, since only in this case is the mass interior to r finite, namely

$$M(r) = \frac{4\pi\rho_0 r_0^\alpha}{3-\alpha} r^{3-\alpha}. \quad (2.60)$$

From equations (2.60) and (2.29) the circular speed is given by

$$v_c^2(r) = \frac{4\pi G\rho_0 r_0^\alpha}{3-\alpha} r^{2-\alpha}. \quad (2.61)$$

In Chapter 1 we saw that the circular-speed curves of many galaxies are remarkably flat. Equation (2.61) suggests that the mass density in these galaxies is proportional to r^{-2} , corresponding to $\alpha = 2$. In §4.3.3(b) we shall find that this is the density profile characteristic of a stellar-dynamical model called the “singular isothermal sphere.”

The potential difference between radius r and the reference radius r_0 is

$$\begin{aligned} \Phi(r) - \Phi(r_0) &= G \int_{r_0}^r dr' \frac{M(r')}{r'^2} = \frac{4\pi G\rho_0 r_0^\alpha}{3-\alpha} \int_{r_0}^r dr' r'^{(1-\alpha)} \\ &= \begin{cases} \frac{v_c^2(r_0) - v_c^2(r)}{\alpha - 2} & \text{for } \alpha \neq 2 \\ v_c^2 \ln(r/r_0) & \text{for } \alpha = 2 \end{cases} \end{aligned} \quad (2.62)$$

Equation (2.60) shows that $M(r)$ diverges at large r for all $\alpha < 3$. However, such models are still useful because by Newton’s first theorem, the mass exterior to any radius r does not affect the dynamics interior to r . For $3 > \alpha > 2$, $v_c^2(r)$ decreases with increasing r (eq. 2.61), so by (2.62) there is only a finite potential difference between radius r and infinity. In fact, the escape speed $v_e(r)$ from radius r is given by

$$v_e^2(r) = 2[\Phi(\infty) - \Phi(r)] = 2\frac{v_c^2(r)}{\alpha - 2} \quad (\alpha > 2). \quad (2.63)$$

For $\alpha < 2$, $\Phi(r)$ grows without limit as $r \rightarrow \infty$, so the interpretation of Φ as the energy per unit mass required to remove a particle to infinity is no longer valid. The potential is nevertheless a useful concept because the

gravitational field is given by $\mathbf{g} = -\nabla\Phi$. Since the light distributions of elliptical galaxies suggest $\alpha \simeq 3$ at large r (see eq. 2.53), while the flatness of the circular-speed curves in spiral galaxies suggest $\alpha \simeq 2$, it is clear that the escape speeds of galaxies are very uncertain.

(g) Two-power density models The luminosity density of many elliptical galaxies can be approximated as a power law in radius at both the largest and smallest observable radii, with a smooth transition between these power laws at intermediate radii (BM §4.3.1). Numerical simulations of the clustering of dark-matter particles suggest that the mass density within a dark halo has a similar structure (§9.3). For these reasons much attention has been devoted to models in which the density is given by

$$\rho(r) = \frac{\rho_0}{(r/a)^\alpha (1+r/a)^{\beta-\alpha}}. \quad (2.64)$$

With $\beta = 4$ these models have particularly simple analytic properties, and are known as **Dehnen models** (BM eq. 4.15; Dehnen 1993; Tremaine et al. 1994). BM Table 4.5 gives formulae for the projected surface density of Dehnen models for the cases $\alpha = 0, 1, \frac{3}{2}$, and 2. The model with $\alpha = 1$ and $\beta = 4$ is called a **Hernquist model** (Hernquist 1990), while that with $\alpha = 2$ and $\beta = 4$ is called a **Jaffe model** (Jaffe 1983). Dehnen models with α in the range (0.6, 2) provide reasonable models of the centers of elliptical galaxies (BM §4.3.1).

Dark halos can be modeled by equation (2.64) with $\beta \simeq 3$ and α in the range (1, 1.5). The model with $(\alpha, \beta) = (1, 3)$ is called the **NFW model** after Navarro, Frenk, & White (1995). The NFW formula contains two free parameters: ρ_0 and a . Navarro, Frenk, & White (1996) showed that the values taken by these parameters for the halos that formed in their simulations were strongly correlated, so the halos were essentially members of a one-parameter family. The conventional choice of parameter is the distance r_{200} from the center of the halo at which the mean density is 200 times the cosmological critical density, ρ_c (eq. 1.52).⁴ A more physical choice of parameter is the mass interior to r_{200} , which is $M = 200\rho_c \frac{4}{3}\pi r_{200}^3$. The **concentration** of the halo is $c \equiv r_{200}/a$. The central result of Navarro, Frenk, & White (1996) is that at a given value of M (and therefore r_{200}), halos show a relatively small scatter in c . The mean value of c falls from $\simeq 16$ to $\simeq 6$ as the halo mass increases from $\sim 3 \times 10^{11} \mathcal{M}_\odot$ to $\sim 3 \times 10^{15} \mathcal{M}_\odot$.

According to (2.64) the mass interior to radius r is

$$M(r) = 4\pi\rho_0 a^3 \int_0^{r/a} ds \frac{s^{2-\alpha}}{(1+s)^{\beta-\alpha}}. \quad (2.65)$$

⁴ In §9.2.1 we shall show that r_{200} approximately divides the interior region in which material has crossed the halo at least once from the exterior one in which matter is still falling in to the halo for the first time.

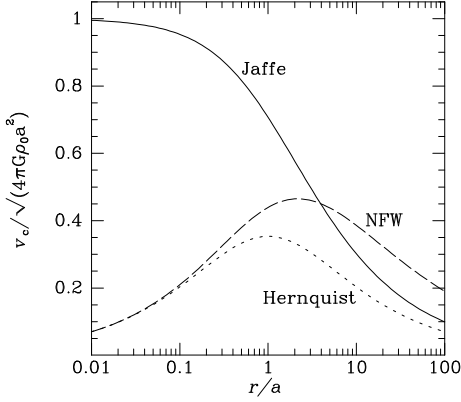


Figure 2.5 Circular speed versus radius for the Jaffe, Hernquist, and NFW models.

This integral is readily evaluated for integer values of α, β . For the important cases of the Jaffe, Hernquist, and NFW models, we have

$$M(r) = 4\pi\rho_0 a^3 \times \begin{cases} \frac{r/a}{1+r/a} & \text{for a Jaffe model} \\ \frac{(r/a)^2}{2(1+r/a)^2} & \text{for a Hernquist model} \\ \ln(1+r/a) - \frac{r/a}{1+r/a} & \text{for a NFW model.} \end{cases} \quad (2.66)$$

In the Jaffe and Hernquist models the mass asymptotes to a finite value as $r \rightarrow \infty$, while in the NFW model the mass diverges logarithmically with r . The circular speed in each model, which follows immediately from (2.66) and (2.29), is plotted in Figure 2.5. From equation (2.29) the potentials of the three models are

$$\begin{aligned} \Phi &= -G \int_r^\infty dr \frac{M(r)}{r^2} \\ &= -4\pi G \rho_0 a^2 \times \begin{cases} \ln(1+a/r) & \text{for a Jaffe model} \\ \frac{1}{2(1+r/a)} & \text{for a Hernquist model} \\ \frac{\ln(1+r/a)}{r/a} & \text{for a NFW model.} \end{cases} \end{aligned} \quad (2.67)$$

From these formulae and equations (2.32) and (2.42) it is straightforward to show that the Jaffe and Hernquist models have gravitational radii $r_g = 2a$ and $6a$, respectively, while for the NFW model r_g is undefined.

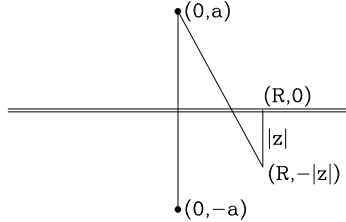


Figure 2.6 At the point $(R, -|z|)$ below Kuzmin's disk, the potential is identical with that of a point mass located distance a above the disk's center.

2.3 Potential-density pairs for flattened systems

Later in this chapter we will show how to obtain the gravitational potentials of systems of arbitrary shape. However, we shall find that the calculation of the gravitational potential and field generated by a given distribution of matter is often an arduous task that leads to cumbersome formulae involving special functions, or numerical calculations. Fortunately, for many purposes it suffices to represent a galaxy by a simple model that has the same gross structure as the galaxy. In this section we describe families of potentials that are generated by fairly simple and realistic axisymmetric density distributions. These potentials help us to understand how the gravitational potential of an initially spherical body is affected by flattening, and in later chapters we shall use several of these potentials to illustrate features of dynamics in axisymmetric galaxies.

2.3.1 Kuzmin models and generalizations

Consider the axisymmetric potential

$$\Phi_K(R, z) = -\frac{GM}{\sqrt{R^2 + (a + |z|)^2}} \quad (a \geq 0). \quad (2.68a)$$

As Figure 2.6 indicates, at points with $z < 0$, Φ_K is identical with the potential of a point mass M located at the point $(R, z) = (0, a)$; and when $z > 0$, Φ_K coincides with the potential generated by a point mass at $(0, -a)$. Hence $\nabla^2 \Phi_K$ must vanish everywhere except on the plane $z = 0$. By applying Gauss's theorem (2.12) to a flat volume that contains a small portion of the plane $z = 0$, we conclude that Φ_K is generated by the surface density

$$\Sigma_K(R) = \frac{aM}{2\pi(R^2 + a^2)^{3/2}}. \quad (2.68b)$$

The potential-density pair of equations (2.68) was introduced by Kuzmin (1956), but it is often referred to as "Toomre's model 1" because it became widely known in the West only after Toomre (1963) unknowingly re-derived it.

Consider next the potential

$$\Phi_{\text{M}}(R, z) = -\frac{GM}{\sqrt{R^2 + (a + \sqrt{z^2 + b^2})^2}}. \quad (2.69a)$$

When $a = 0$, Φ_{M} reduces to Plummer's spherical potential (2.44a), and when $b = 0$, Φ_{M} reduces to Kuzmin's potential of a razor-thin disk (2.68a). Thus, depending on the choice of the two parameters a and b , Φ_{M} can represent the potential of anything from an infinitesimally thin disk to a spherical system. If we calculate $\nabla^2\Phi_{\text{M}}$, we find that the mass distribution with which it is associated is (Miyamoto & Nagai 1975)

$$\rho_{\text{M}}(R, z) = \left(\frac{b^2 M}{4\pi}\right) \frac{aR^2 + (a + 3\sqrt{z^2 + b^2})(a + \sqrt{z^2 + b^2})^2}{[R^2 + (a + \sqrt{z^2 + b^2})^2]^{5/2}(z^2 + b^2)^{3/2}}. \quad (2.69b)$$

In Figure 2.7 we show contour plots of $\rho_{\text{M}}(R, z)$ for various values of b/a . When $b/a \simeq 0.2$, these are qualitatively similar to the light distributions of disk galaxies, although there are quantitative differences. For example, from equation (2.69b) we have that $\rho(R, 0) \propto R^{-3}$ when R is large, whereas the brightness profiles of disks fall off at least as fast as $\exp(-R/R_{\text{d}})$ (eq. 1.7).

Since Poisson's equation is linear in Φ and ρ , the difference between any two potential-density pairs is itself a potential-density pair. Therefore, if we differentiate a potential-density pair with respect to one of its parameters, we obtain a new potential-density pair. For example, Toomre (1963) derived a family of potential-density pairs by differentiating $\Phi_{\text{K}}(R, z)/a$ n times with respect to a^2 . Similarly, Satoh (1980) obtained a series of spherical potential-density pairs by differentiating b^{-2} times the Plummer potential and density (eq. 2.44) n times with respect to b^2 . He flattened these potentials by replacing $r^2 + b^2$ with $R^2 + (a + \sqrt{z^2 + b^2})^2$, and in the limit $n \rightarrow \infty$ obtained

$$\Phi_{\text{S}}(R, z) = -\frac{GM}{S}, \quad (2.70a)$$

where

$$S^2 \equiv R^2 + z^2 + a \left(a + 2\sqrt{z^2 + b^2} \right). \quad (2.70c)$$

The corresponding density distribution follows from Poisson's equation:

$$\rho_{\text{S}}(R, z) = \frac{ab^2 M}{4\pi S^3(z^2 + b^2)} \left[\frac{1}{\sqrt{z^2 + b^2}} + \frac{3}{a} \left(1 - \frac{R^2 + z^2}{S^2} \right) \right]. \quad (2.70b)$$

Figure 2.8 shows that at large b/a the isodensity surfaces Φ_{S} are more nearly elliptical than those of Φ_{M} .

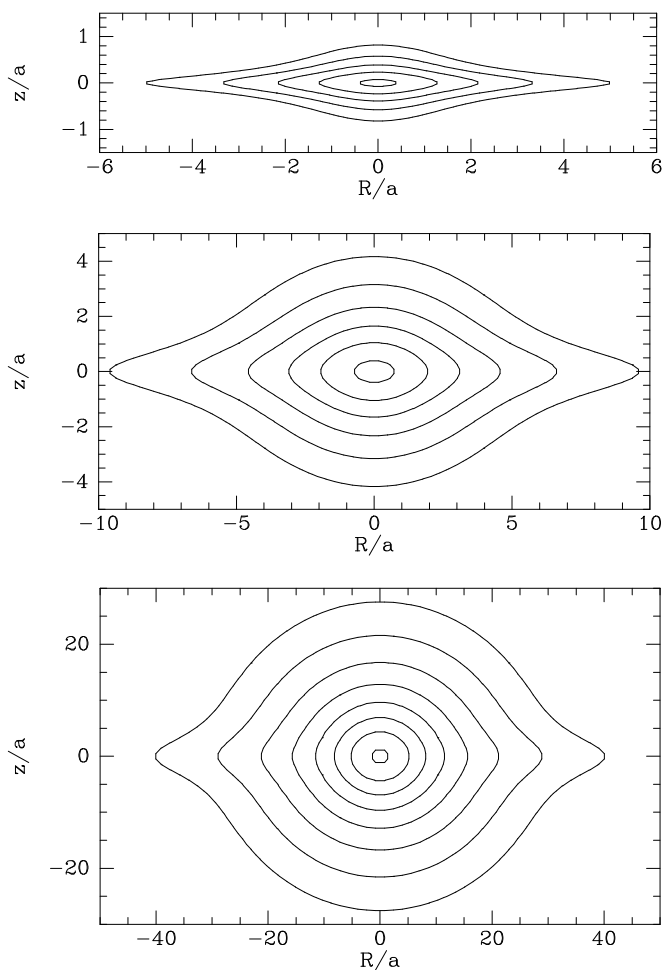


Figure 2.7 Contours of equal density in the (R, z) plane for the Miyamoto–Nagai density distribution (2.69b) when: $b/a = 0.2$ (top); $b/a = 1$ (middle); $b/a = 5$ (bottom). There are two contours per decade, and the highest contour levels are $0.3M/a^3$ (top), $0.03M/a^3$ (middle), and $0.001M/a^3$ (bottom).

2.3.2 Logarithmic potentials

Since the Kuzmin and other models in the previous subsection all have finite mass, the circular speed associated with these potentials falls off in Keplerian fashion, $v_c \propto R^{-1/2}$, at large R . However, in Chapter 1 it was shown that the circular-speed curves of spiral galaxies tend to be flat at large radii. If at large R , $v_c = v_0$, a constant, then $d\Phi/dR = v_0^2/R$, and hence $\Phi \propto$

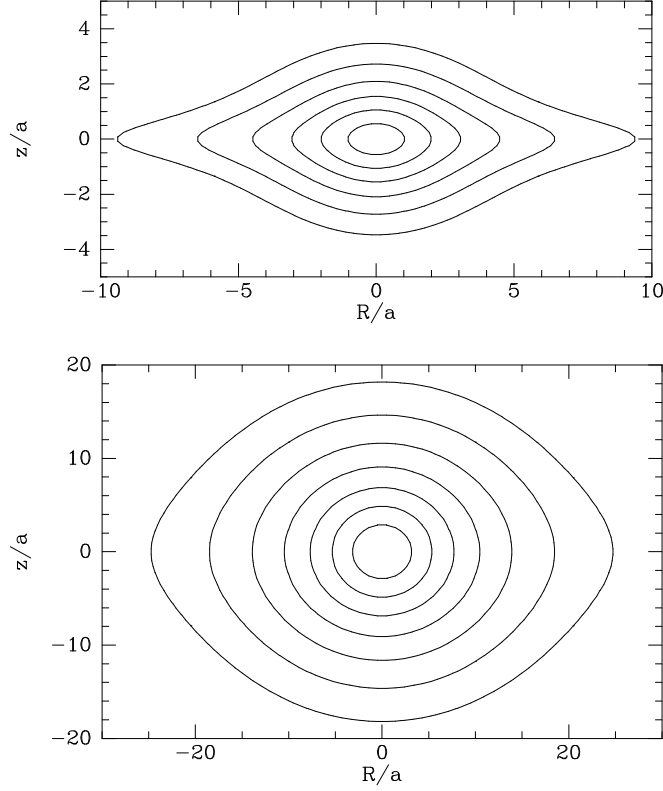


Figure 2.8 Contours of equal density in the (R, z) plane for Satoh's density distribution (2.70c) when: $b/a = 1$ (top); $b/a = 10$ (bottom). There are two contours per decade, and the highest contour levels are $0.1M/a^3$ (top), and $0.001M/a^3$ (bottom).

$v_0^2 \ln R + \text{constant}$ in this region. Therefore, consider the potential

$$\Phi_L = \frac{1}{2}v_0^2 \ln \left(R_c^2 + R^2 + \frac{z^2}{q_\Phi^2} \right) + \text{constant}, \quad (2.71a)$$

where R_c and v_0 are constants, and q_Φ is the axis ratio of the equipotential surfaces. The circular speed at radius R in the equatorial plane of Φ_L is

$$v_c = \frac{v_0 R}{\sqrt{R_c^2 + R^2}}. \quad (2.71b)$$

The density distribution to which Φ_L corresponds is

$$\rho_L(R, z) = \frac{v_0^2}{4\pi G q_\Phi^2} \frac{(2q_\Phi^2 + 1)R_c^2 + R^2 + (2 - q_\Phi^{-2})z^2}{(R_c^2 + R^2 + z^2 q_\Phi^{-2})^2}. \quad (2.71c)$$

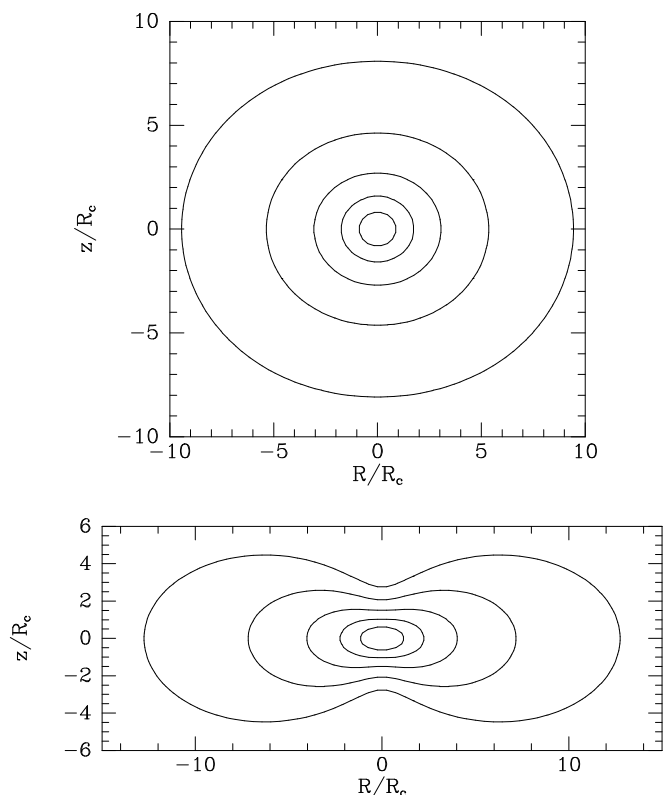


Figure 2.9 Contours of equal density in the (R, z) plane for ρ_L (eq. 2.71c) when $q_\Phi = 0.95$ (top), $q_\Phi = 0.7$ (bottom). There are two contours per decade and the highest contour level is $0.1v_0^2/(GR_c^2)$. When $q_\Phi = 0.7$ the models are unphysical because the density is negative near the z axis for $|z| \gtrsim 7R_c$.

At small R and z , ρ_L tends to a constant value, and when R or $|z|$ is large, ρ_L falls off as R^{-2} or z^{-2} .

The equipotential surfaces of Φ_L are spheroids⁵ of axial ratio q_Φ , but Figure 2.9 shows that the equidensity surfaces are rather flatter and can deviate strongly from spheroids. In fact, if we define the axial ratio q_ρ of the isodensity surfaces by the ratio z_m/R_m of the distances down the z and R axes at which a given isodensity surface cuts the z axis and the x or y axis, we find

$$q_\rho^2 \simeq \frac{1 + 4q_\Phi^2}{2 + 3/q_\Phi^2} \quad (r \ll R_c) \quad (2.72a)$$

⁵ A **spheroid** is the surface generated by rotating an ellipse about one of its principal axes. An **oblate**, or flattened, spheroid is generated if the axis of rotation is the minor axis, and a **prolate**, or elongated, spheroid is generated if this axis is the major axis.

or

$$q_\rho^2 \simeq q_\Phi^4 \left(2 - \frac{1}{q_\Phi^2} \right) \quad (r \gg R_c). \quad (2.72b)$$

Outside the core, the flattening $1 - q_\Phi$ of the potential is only about a third that of the density distribution: $1 - q_\Phi \simeq \frac{1}{3}(1 - q_\rho)$. The density ρ_L becomes negative on the z axis when $q_\Phi < 1/\sqrt{2} = 0.707$.

2.3.3 Poisson's equation in very flattened systems

In any axisymmetric system with density $\rho(R, z)$, Poisson's equation can be written (eq. B.52)

$$\frac{\partial^2 \Phi}{\partial z^2} = 4\pi G \rho(R, z) + \frac{1}{R} \frac{\partial}{\partial R} (R F_R), \quad (2.73)$$

where $F_R = -\partial\Phi/\partial R$ is the radial force. Now consider, for example, the Miyamoto–Nagai potential-density pair given by equations (2.69). As the parameter $b \rightarrow 0$, the density distribution becomes more and more flattened, and at fixed R the density in the plane $z = 0$ becomes larger and larger as b^{-1} . However, the radial force F_R remains well behaved as $b \rightarrow 0$; indeed, in the limit $b = 0$, $F_R = -\partial\Phi_K/\partial R$, where $\Phi_K(R, z)$ is simply the Kuzmin potential (2.68a). Thus, near $z = 0$ the first term on the right side of equation (2.73) becomes very large compared to the second, and Poisson's equation simplifies to the form

$$\frac{\partial^2 \Phi(R, z)}{\partial z^2} = 4\pi G \rho(R, z). \quad (2.74)$$

This result applies to almost any thin disk system. It implies that the vertical variation of the potential at a given radius R depends only on the density distribution at that radius. Effectively, this means that the solution of Poisson's equation in a thin disk can be decomposed into two steps: (i) Approximate the thin disk as a surface density layer of zero thickness and determine the potential in the plane of the disk $\Phi(R, 0)$ using the models of this section or the more general techniques of §2.6. (ii) At each radius R solve equation (2.74) to find the vertical variation of $\Phi(R, z)$.

Thus we have

$$\Phi(R, z) = \Phi(R, 0) + \Phi_z(R, z) \quad (2.75a)$$

where

$$\Phi_z(R, z) \equiv 4\pi G \int_0^z dz' \int_0^{z'} dz'' \rho(R, z'') + a(R)z. \quad (2.75b)$$

The constant of integration, a , is zero if the disk is symmetric around the equatorial plane.

2.4 Multipole expansion

In the last section we encountered a number of axisymmetric density distributions that give rise to potentials of known form. By adding a few of these distributions together one can obtain quite a wide range of model galaxies that have readily available potentials. However, for many purposes one requires a systematic procedure for calculating the potential of an arbitrary density distribution to whatever accuracy one pleases. The next few sections are devoted to this task.

The first such technique, based on spherical harmonics, works best for systems that are neither very flattened nor very elongated. Hence it is a good method for calculating the potentials of bulges and dark-matter halos (§§1.1.2 and 1.1.3).

Our first step is to obtain the potential of a thin spherical shell of variable surface density. Since the shell has negligible thickness, the task of solving Poisson's equation $\nabla^2\Phi = 4\pi G\rho$ reduces to that of solving Laplace's equation $\nabla^2\Phi = 0$ inside and outside the shell, subject to suitable boundary conditions at infinity, at the origin, and on the shell. Now in spherical coordinates Laplace's equation is (eq. B.53)

$$\frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial \Phi}{\partial r} \right) + \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial \Phi}{\partial \theta} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial^2 \Phi}{\partial \phi^2} = 0. \quad (2.76)$$

This may be solved by the method of **separation of variables**. We seek special solutions that are the product of functions of one variable only:

$$\Phi(r, \theta, \phi) = R(r)P(\theta)Q(\phi). \quad (2.77a)$$

Substituting equation (2.77a) into (2.76) and rearranging, we obtain

$$\frac{\sin^2 \theta}{R} \frac{d}{dr} \left(r^2 \frac{dR}{dr} \right) + \frac{\sin \theta}{P} \frac{d}{d\theta} \left(\sin \theta \frac{dP}{d\theta} \right) = -\frac{1}{Q} \frac{d^2 Q}{d\phi^2}. \quad (2.77b)$$

The left side of this equation does not depend on ϕ , and the right side does not depend on r or θ . It follows that both sides are equal to some constant, say m^2 . Hence

$$-\frac{1}{Q} \frac{d^2 Q}{d\phi^2} = m^2, \quad (2.78a)$$

$$\frac{\sin^2 \theta}{R} \frac{d}{dr} \left(r^2 \frac{dR}{dr} \right) + \frac{\sin \theta}{P} \frac{d}{d\theta} \left(\sin \theta \frac{dP}{d\theta} \right) = m^2. \quad (2.78b)$$

Equation (2.78a) may be immediately integrated to

$$Q(\phi) = Q_m^+ e^{im\phi} + Q_m^- e^{-im\phi}. \quad (2.79a)$$

We require Φ to be a periodic function of ϕ with period 2π , so m can take only integer values. Since equations (2.78) depend only on m^2 , we could restrict our attention to non-negative values of m without loss of generality. However, a simpler convention is to allow m to take both positive and negative values, so the second exponential in equation (2.79a) becomes redundant, and we may write simply

$$Q = Q_m e^{im\phi} \quad (m = \dots, -1, 0, 1, \dots). \quad (2.79b)$$

Equation (2.78b) can be written

$$\frac{1}{R} \frac{d}{dr} \left(r^2 \frac{dR}{dr} \right) = \frac{m^2}{\sin^2 \theta} - \frac{1}{P \sin \theta} \frac{d}{d\theta} \left(\sin \theta \frac{dP}{d\theta} \right). \quad (2.80)$$

Since the left side of this equation does not depend on θ and the right side does not depend on r , both sides must equal some constant, which we write as $l(l+1)$. Thus equation (2.80) implies

$$\frac{d}{dr} \left(r^2 \frac{dR}{dr} \right) - l(l+1)R = 0, \quad (2.81a)$$

and in terms of $x \equiv \cos \theta$,

$$\frac{d}{dx} \left[(1-x^2) \frac{dP}{dx} \right] - \frac{m^2}{1-x^2} P + l(l+1)P = 0. \quad (2.81b)$$

Two linearly independent solutions of equation (2.81a) are

$$R(r) = Ar^l \quad \text{and} \quad R(r) = Br^{-(l+1)}. \quad (2.82)$$

The solutions of equation (2.81b) are associated Legendre functions $P_l^m(x)$ (see Appendix C.5). Physically acceptable solutions exist only when l is an integer. Without loss of generality we can take l to be non-negative, and then physically acceptable solutions exist only for $|m| \leq l$. When $m = 0$ the solutions are simply polynomials in x , called Legendre polynomials $P_l(x)$.

Rather than write out the product $P_l^m(\cos \theta)e^{im\phi}$ again and again, it is helpful to define the spherical harmonic $Y_l^m(\theta, \phi)$, which is equal to $P_l^m(\cos \theta)e^{im\phi}$ times a constant chosen so the Y_l^m satisfy the orthogonality relation (see eq. C.44)

$$\begin{aligned} \int d^2\Omega Y_l^{m*}(\mathbf{\Omega}) Y_{l'}^{m'}(\mathbf{\Omega}) &\equiv \int_0^\pi d\theta \sin \theta \int_0^{2\pi} d\phi Y_l^{m*}(\theta, \phi) Y_{l'}^{m'}(\theta, \phi) \\ &= \delta_{ll'} \delta_{mm'}, \end{aligned} \quad (2.83)$$

where we have used $\mathbf{\Omega}$ as a shorthand for (θ, ϕ) and $d^2\Omega$ for $\sin \theta d\theta d\phi$. The spherical harmonics with $l \leq 2$ are listed in equation (C.50).

Putting all these results together, we have from equations (2.77a), (2.79b), and (2.82) that

$$\Phi_{lm}(r, \mathbf{\Omega}) = \left(A_{lm}r^l + B_{lm}r^{-(l+1)} \right) Y_l^m(\mathbf{\Omega}) \quad (2.84)$$

is a solution of $\nabla^2\Phi = 0$ for all non-negative integers l and integer m in the range $-l \leq m \leq l$.

Now let us apply these results to the problem of determining the potential of a thin shell of radius a and surface density $\sigma(\mathbf{\Omega})$. We write the potential internal and external to the shell as

$$\Phi_{\text{int}}(r, \mathbf{\Omega}) = \sum_{l=0}^{\infty} \sum_{m=-l}^l \left(A_{lm}r^l + B_{lm}r^{-(l+1)} \right) Y_l^m(\mathbf{\Omega}) \quad (r \leq a), \quad (2.85a)$$

and

$$\Phi_{\text{ext}}(r, \mathbf{\Omega}) = \sum_{l=0}^{\infty} \sum_{m=-l}^l \left(C_{lm}r^l + D_{lm}r^{-(l+1)} \right) Y_l^m(\mathbf{\Omega}) \quad (r \geq a). \quad (2.85b)$$

The potential at the center must be non-singular, so $B_{lm} = 0$ for all l, m . Similarly, the potential at infinity must be zero, so $C_{lm} = 0$ for all l, m . Furthermore, $\Phi_{\text{ext}}(a, \mathbf{\Omega})$ must equal $\Phi_{\text{int}}(a, \mathbf{\Omega})$ because no work can be done in passing through an infinitesimally thin shell. Hence from equations (2.85) we have

$$\sum_{l=0}^{\infty} \sum_{m=-l}^l A_{lm}a^l Y_l^m(\mathbf{\Omega}) = \sum_{l=0}^{\infty} \sum_{m=-l}^l D_{lm}a^{-(l+1)} Y_l^m(\mathbf{\Omega}). \quad (2.86)$$

The coefficients $A_{lm}a^l$ etc. of each spherical harmonic Y_l^m on the two sides of equation (2.86) must be equal, as can be shown by multiplying both sides of the equation by $Y_{l'}^{m'*}(\mathbf{\Omega})$, integrating over $\mathbf{\Omega}$, and using the orthogonality relation (2.83). Therefore, from equation (2.86) we have

$$D_{lm} = A_{lm}a^{2l+1}. \quad (2.87)$$

Now let us expand the surface density of the thin shell as

$$\sigma(\mathbf{\Omega}) = \sum_{l=0}^{\infty} \sum_{m=-l}^l \sigma_{lm} Y_l^m(\mathbf{\Omega}), \quad (2.88)$$

where the σ_{lm} are numbers yet to be determined. To obtain the coefficient $\sigma_{l'm'}$, we multiply both sides of equation (2.88) by $Y_{l'}^{m'*}(\mathbf{\Omega})$ and integrate over $\mathbf{\Omega}$. With equation (2.83) we find

$$\int d^2\Omega Y_{l'}^{m'*}(\mathbf{\Omega}) \sigma(\mathbf{\Omega}) = \sigma_{l'm'}. \quad (2.89)$$

Since $Y_0^0 = 1/\sqrt{4\pi}$, $\sigma_{00} = M/(2a^2\sqrt{\pi})$, where M is the mass of the shell.

Gauss's theorem (2.12) applied to a small piece of the shell tells us that

$$\left(\frac{\partial\Phi_{\text{ext}}}{\partial r}\right)_{r=a} - \left(\frac{\partial\Phi_{\text{int}}}{\partial r}\right)_{r=a} = 4\pi G\sigma(\boldsymbol{\Omega}), \quad (2.90)$$

so inserting equations (2.85) and (2.88) into equation (2.90), we obtain

$$\begin{aligned} -\sum_{l=0}^{\infty} \sum_{m=-l}^l \left((l+1)D_{lm}a^{-(l+2)} + lA_{lm}a^{l-1} \right) Y_l^m(\boldsymbol{\Omega}) = \\ 4\pi G \sum_{l=0}^{\infty} \sum_{m=-l}^l \sigma_{lm} Y_l^m(\boldsymbol{\Omega}). \end{aligned} \quad (2.91)$$

Once again the coefficients of Y_l^m on each side of the equation must be identical, so with (2.87) we have

$$A_{lm} = -4\pi G a^{-(l-1)} \frac{\sigma_{lm}}{2l+1} \quad ; \quad D_{lm} = -4\pi G a^{l+2} \frac{\sigma_{lm}}{2l+1}. \quad (2.92)$$

Collecting these results together, we have from equations (2.85) that

$$\begin{aligned} \Phi_{\text{int}}(r, \boldsymbol{\Omega}) &= -4\pi G a \sum_{l=0}^{\infty} \left(\frac{r}{a}\right)^l \sum_{m=-l}^l \frac{\sigma_{lm}}{2l+1} Y_l^m(\boldsymbol{\Omega}), \\ \Phi_{\text{ext}}(r, \boldsymbol{\Omega}) &= -4\pi G a \sum_{l=0}^{\infty} \left(\frac{a}{r}\right)^{l+1} \sum_{m=-l}^l \frac{\sigma_{lm}}{2l+1} Y_l^m(\boldsymbol{\Omega}), \end{aligned} \quad (2.93)$$

where the σ_{lm} are given by equation (2.89).

Finally we evaluate the potential of a solid body by breaking it down into a series of spherical shells. We let $\delta\sigma_{lm}(a)$ be the σ -coefficient of the shell lying between a and $a + \delta a$, and $\delta\Phi(r, \boldsymbol{\Omega}; a)$ be the corresponding potential at r . Then we have by equation (2.89)

$$\delta\sigma_{lm}(a) = \int_0^\pi d\theta \sin\theta \int_0^{2\pi} d\phi Y_l^{m*}(\boldsymbol{\Omega}) \rho(a, \boldsymbol{\Omega}) \delta a \equiv \rho_{lm}(a) \delta a. \quad (2.94)$$

Substituting these values of σ_{lm} into equations (2.93) and integrating over all a , we obtain the potential at r generated by the entire collection of shells:

$$\begin{aligned} \Phi(r, \boldsymbol{\Omega}) &= \sum_{a=0}^r \delta\Phi_{\text{ext}} + \sum_{a=r}^{\infty} \delta\Phi_{\text{int}} \\ &= -4\pi G \sum_{l,m} \frac{Y_l^m(\boldsymbol{\Omega})}{2l+1} \left(\frac{1}{r^{l+1}} \int_0^r da a^{l+2} \rho_{lm}(a) + r^l \int_r^\infty \frac{da}{a^{l-1}} \rho_{lm}(a) \right). \end{aligned} \quad (2.95)$$

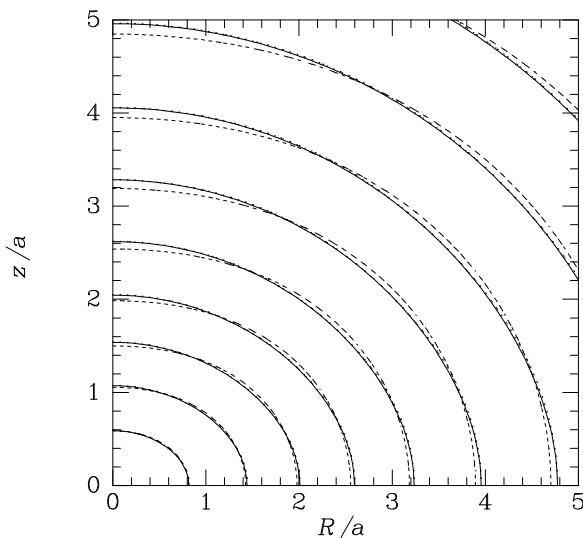


Figure 2.10 Equipotentials of Satoh's density distribution (2.70c) with $b/a = 1$. Full curves show the exact equipotentials computed from equation (2.70a), and dashed curves show the estimate provided by equation (2.95) with the sum over l extending to $l = 2$. Contours based on the sum to $l = 8$ are also plotted (dotted contours) but almost overlie the full curves.

This equation gives the potential generated by the body as an expansion in **multipoles**: the terms associated with $l = m = 0$ are the **monopole** terms, those associated with $l = 1$ are **dipole** terms, those with $l = 2$ are **quadrupole** terms, and those with larger l are 2^l -poles. Similar expansions occur in electrostatics (e.g., Jackson 1999). The monopole terms are the same as in equation (2.28) for the potential of a spherical system.⁶ Since there is no gravitational analog of negative charge, pure dipole or quadrupole gravitational potentials cannot arise, in contrast to the electrostatic case. In fact, if one places the origin of coordinates at the center of mass of the system, the dipole term vanishes identically outside any matter distribution. While the monopole terms generate a circular-speed curve $v_c(r) = \sqrt{GM(r)/r}$ that never declines with increasing r more steeply than in the Keplerian case ($v_c \propto r^{-1/2}$), over a limited range in r the higher-order multipoles may cause the circular speed to fall more steeply with increasing radius.

As an illustration of the effectiveness of the multipole expansion, we show in Figure 2.10 the contours of Satoh's potential $\Phi_S(R, z)$ (eq. 2.70a), together with the approximations to this potential that one obtains from equation (2.95) if one includes only terms with $l \leq 2$ or 8. The flexibility

⁶ Thus the spherical-harmonic expansion provides an alternative proof of Newton's first and second theorems.

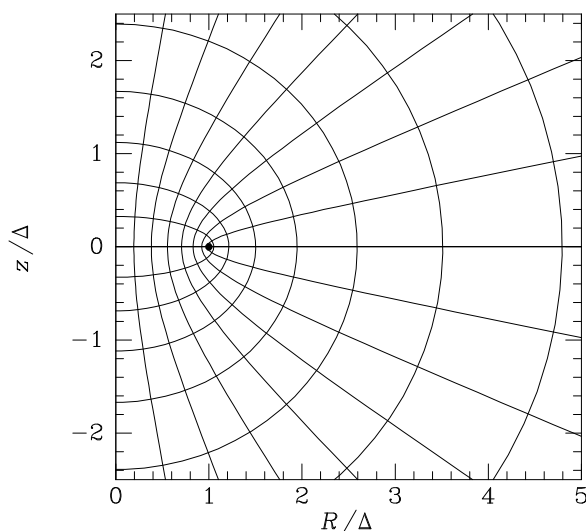


Figure 2.11 Curves of constant u and v in the (R, z) plane. Semi-ellipses are curves of constant u , and hyperbolae are curves of constant v . The common focus of all curves is marked by a dot. In order to ensure that each point has a unique v -coordinate, we exclude the disk ($z = 0$, $R \leq \Delta$) from the space to be considered.

of the multipole expansion makes it a powerful tool for numerical work, and it plays a central role in some of the Poisson solvers for N-body codes that will be described in §2.9. However, multipole expansions are poorly suited for modeling the potentials of disks.

2.5 The potentials of spheroidal and ellipsoidal systems

Many galaxies have nearly spheroidal or ellipsoidal equidensity surfaces (BM §4.3.3). Moreover, Newton's theorems for spherical bodies can be generalized to include spheroidal and ellipsoidal bodies, so models with isodensity surfaces of this shape are relatively easy to construct. Finally, as the axis ratio shrinks to zero a spheroid becomes a disk, and thus we can obtain the potentials of razor-thin disks as a limiting case of spheroids.

In this section we develop efficient techniques for calculating the potentials of such objects. In §§2.5.1 and 2.5.2 we derive formulae for oblate (i.e., flattened) spheroidal systems, and in §2.5.3 we briefly discuss ellipsoidal systems. Results for prolate spheroidal systems can be obtained either by adapting our derivations for oblate systems, or by specializing the results for ellipsoidal systems. The principal formulae for all three geometries are given in Tables 2.1 and 2.2.

2.5.1 Potentials of spheroidal shells

Consider a system in which the isodensity surfaces are similar concentric spheroids. We align our coordinate system such that the z axis is the minor axis of the spheroids, and designate the principal-axis lengths of the spheroids by a and c , with $a \geq c$. Spheroidal bodies call for spheroidal coordinates, so we consider the form of Laplace's equation, $\nabla^2\Phi = 0$, in **oblate spheroidal coordinates**. These coordinates employ the usual azimuthal angle ϕ of cylindrical coordinates, but replace the coordinates (R, z) with new coordinates (u, v) that are defined by

$$R = \Delta \cosh u \sin v \quad ; \quad z = \Delta \sinh u \cos v \quad (u \geq 0, 0 \leq v \leq \pi), \quad (2.96)$$

where Δ is a constant. Figure 2.11 shows the curves of constant u and v in the (R, z) plane. The curves $u = \text{constant}$ are confocal half-ellipses with foci at $(R, z) = (\Delta, 0)$, namely

$$\frac{R^2}{\cosh^2 u} + \frac{z^2}{\sinh^2 u} = \Delta^2. \quad (2.97)$$

Similarly, the curves $v = \text{constant}$ coincide with the hyperbolae

$$\frac{R^2}{\sin^2 v} - \frac{z^2}{\cos^2 v} = \Delta^2 \quad (2.98)$$

formed by the normals to these ellipses.

From equations (2.96) we see that if we increase one of u , v , and ϕ by a small amount while holding the other two coordinates constant, the point (u, v, ϕ) moves parallel to the three orthogonal unit vectors $\hat{\mathbf{e}}_u$, $\hat{\mathbf{e}}_v$, $\hat{\mathbf{e}}_\phi$ by the distances $h_u \delta u$, $h_v \delta v$ and $h_\phi \delta \phi$, where the scale factors are

$$h_u = h_v = \Delta \sqrt{\sinh^2 u + \cos^2 v} \quad ; \quad h_\phi = \Delta \cosh u \sin v. \quad (2.99)$$

Hence the gradient of a potential Φ may be expressed in these coordinates as (eq. B.39)

$$\nabla\Phi = \frac{1}{\Delta \sqrt{\sinh^2 u + \cos^2 v}} \left[\frac{\partial\Phi}{\partial u} \hat{\mathbf{e}}_u + \frac{\partial\Phi}{\partial v} \hat{\mathbf{e}}_v \right] + \frac{1}{\Delta \cosh u \sin v} \frac{\partial\Phi}{\partial \phi} \hat{\mathbf{e}}_\phi. \quad (2.100)$$

We further have (eq. B.54)

$$\begin{aligned} \nabla^2\Phi = & \frac{1}{\Delta^2(\sinh^2 u + \cos^2 v)} \left[\frac{1}{\cosh u} \frac{\partial}{\partial u} \left(\cosh u \frac{\partial\Phi}{\partial u} \right) \right. \\ & \left. + \frac{1}{\sin v} \frac{\partial}{\partial v} \left(\sin v \frac{\partial\Phi}{\partial v} \right) \right] + \frac{1}{\Delta^2 \cosh^2 u \sin^2 v} \frac{\partial^2\Phi}{\partial \phi^2}. \end{aligned} \quad (2.101)$$

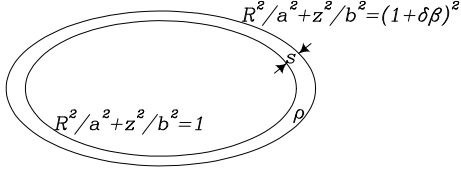


Figure 2.12 A homoeoid of density ρ is bounded by the surfaces $R^2/a^2 + z^2/b^2 = 1$ and $R^2/a^2 + z^2/b^2 = (1 + \delta\beta)^2$. The perpendicular distance s between the bounding surfaces varies with position around the homoeoid.

We concentrate on potentials that are functions $\Phi(u)$ of only the “radial” coordinate u . For potentials of this class, $\nabla^2\Phi = 0$ reduces to

$$\frac{d}{du} \left(\cosh u \frac{d\Phi}{du} \right) = 0. \quad (2.102)$$

Hence either

$$\Phi = \Phi_0 \quad (\text{a constant}), \quad (2.103a)$$

or $d\Phi/du = A \operatorname{sech} u$, where A is a constant. Integrating this last equation we find

$$\Phi = -A \sin^{-1}(\operatorname{sech} u) + B, \quad (2.103b)$$

where B is a constant.

For u large, $\operatorname{sech} u \rightarrow \Delta/r \rightarrow 0$, where r is the usual spherical radius. So a potential of the form (2.103b) varies as

$$\Phi \simeq -A \operatorname{sech} u + B \rightarrow -\frac{A\Delta}{r} + B, \quad (2.104)$$

Hence, if we set $B = 0$ and $A = G\delta M/\Delta$, the potential given by equation (2.103b) tends to zero at infinity like the gravitational potential of a shell of mass δM . Thus we are led to consider the potential defined by

$$\Phi = -\frac{G\delta M}{\Delta} \times \begin{cases} \sin^{-1}(\operatorname{sech} u_0) & (u < u_0), \\ \sin^{-1}(\operatorname{sech} u) & (u \geq u_0). \end{cases} \quad (2.105)$$

This potential is everywhere continuous and solves $\nabla^2\Phi = 0$ everywhere except on the spheroid $u = u_0$ (see eqs. 2.103). Hence it is the gravitational potential of a shell of material on the surface $u = u_0$. This shell has principal semi-axes of lengths $a \equiv \Delta \cosh u_0$ and $c \equiv \Delta \sinh u_0$. Hence the shell’s **eccentricity**

$$e \equiv \sqrt{1 - \frac{c^2}{a^2}} = \operatorname{sech} u_0, \quad (2.106)$$

and we may rewrite equations (2.105) as

$$\Phi = -\frac{G\delta M}{ae} \times \begin{cases} \sin^{-1}(e) & (u < u_0), \\ \sin^{-1}(\operatorname{sech} u) & (u \geq u_0). \end{cases} \quad (2.107)$$

We can find the surface density of the shell $u = u_0$ by applying Gauss's theorem (2.12) to the potential (2.107). Since $\nabla\Phi = 0$ inside the shell, by equation (2.100) the surface density of the shell is

$$\begin{aligned}\Sigma(v) &= \frac{\hat{\mathbf{e}}_u \cdot \nabla\Phi}{4\pi G} = \frac{1}{4\pi G \Delta \sqrt{\sinh^2 u_0 + \cos^2 v}} \left(\frac{d\Phi}{du} \right)_{u=u_0+} \\ &= \frac{\delta M}{4\pi a^2 \sqrt{1 - e^2 \sin^2 v}},\end{aligned}\quad (2.108)$$

where evaluation at u_0+ denotes the limiting value as $u \rightarrow u_0$ from above.

Equation (2.108) has a simple physical interpretation: $\Sigma(v)$ is the surface density of the thin shell of uniform density ρ that is bounded by the two surfaces β and $\beta + \delta\beta$ of the set of similar spheroids,

$$\frac{R^2}{a^2} + \frac{z^2}{c^2} = \beta^2, \quad (2.109)$$

where $\delta\beta$ and δM are related by (2.113) below. *Proof:* The small perpendicular \mathbf{s} in Figure 2.12 runs between the shell's inner and outer skins. Thus $\mathbf{s} = s\nabla\beta/|\nabla\beta|$, and at any point on the surface we have $\delta\beta = (\mathbf{s} \cdot \nabla\beta) = s|\nabla\beta|$. Hence $s = \delta\beta/|\nabla\beta|$ and the surface density of the shell is

$$\tilde{\Sigma} = \rho s = \frac{\rho\delta\beta}{|\nabla\beta|} = \left(\frac{R^2}{a^4} + \frac{z^2}{c^4} \right)^{-1/2} \rho\delta\beta. \quad (2.110)$$

Finally, writing $R = \beta a \sin v$, $z = \beta c \cos v$, and $e = \sqrt{1 - c^2/a^2}$, we find

$$\tilde{\Sigma} = \frac{a\sqrt{1 - e^2} \rho \delta\beta}{\sqrt{1 - e^2 \sin^2 v}}. \quad (2.111)$$

The volume inside an oblate spheroidal shell with semi-axis lengths $a\beta$ and $c\beta$ is

$$V = \frac{4}{3}\pi a^2 c \beta^3 = \frac{4}{3}\pi a^3 \beta^3 \sqrt{1 - e^2}, \quad (2.112)$$

so the mass that is bounded by the surfaces β and $\beta + \delta\beta$ is

$$\delta M = 4\pi\rho a^3 \sqrt{1 - e^2} \beta^2 \delta\beta. \quad (2.113)$$

We can now set $\beta = 1$ and substitute (2.113) into equation (2.111) to find that $\tilde{\Sigma} = \Sigma$, where Σ is given by equation (2.108).◀

We have shown that the potential (2.107) is generated by a thin shell of uniform density that is bounded by similar spheroids of eccentricity e . We call such a shell a **thin homoeoid** and have:

Homoeoid theorem *The exterior isopotential surfaces of a thin homoeoid are the spheroids that are confocal with the shell itself. Inside the shell the potential is constant.*

The homoeoid theorem applies only to a *thin* homoeoid. But it immediately yields a remarkable property of any homoeoid—that is, of any shell of constant density, no matter how thick, whose inner and outer surfaces are similar (*not* confocal) spheroids:

Newton’s third theorem *A mass that is inside a homoeoid experiences no net gravitational force from the homoeoid.*

Proof: Break the given homoeoid into a series of thin homoeoids. The interior of the thick homoeoid lies in the interior of each of its component thin homoeoids, and the interior potential of each thin homoeoid is constant. Hence the aggregate interior potential is constant and generates no gravitational force.◁

Newton’s first theorem, for spherical systems, thus emerges as a special case of Newton’s third theorem for spheroidal systems. Newton’s second theorem, for spherical systems, has no analog for spheroidal systems because the potential outside a spheroidal body *does* depend on the distribution of matter within it.

These theorems help us to understand qualitatively the potential of an inhomogeneous spheroidal body. Each shell of the body makes a contribution to the potential that is constant interior to the shell and on the shell, and gradually becomes rounder as one moves outward from the shell. This tendency of isopotentials to become spherical at large radii manifested itself already in §2.4 in the rapid radial decay of the higher multipole components of the gravitational potential. The shape of the isopotential surface at a distance r from the center of an inhomogeneous spheroidal body represents a compromise between the rather round contributions of the central shells, and the more aspherical contributions of the shells just interior to r . Thus, if the body is very centrally concentrated, the isopotentials near its edge will be nearly round, while a more homogeneous spheroidal body will have more flattened isopotentials.

2.5.2 Potentials of spheroidal systems

We now use equation (2.107) to calculate the gravitational potential of a body whose isodensity surfaces are the similar spheroids

$$\text{constant} = m^2 \equiv R^2 + \frac{z^2}{1 - e^2}, \quad (2.114)$$

i.e., a body in which $\rho = \rho(m^2)$. Comparing equations (2.109) and (2.114), we see that $m = \beta a$, so the mass of the shell between m and $m + \delta m$ is given by equation (2.113) as

$$\delta M = 4\pi\rho(m^2)\sqrt{1 - e^2} m^2 \delta m. \quad (2.115)$$

For a spherical system ($e = 0$) this reduces to the familiar formula $\delta M = 4\pi\rho r^2\delta r$.

There is a unique family of confocal spheroids such that one member of the family coincides with the homoeoid labeled by m . Let $u_m(R_0, z_0)$ be the label of the member of this family that passes through the point (R_0, z_0) at which the potential is required (see eq. 2.122 below for an explicit formula). Then if (R_0, z_0) lies inside the homoeoid m , we have on setting $a = m$ in equation (2.107) and substituting for δM from equation (2.115) that the contribution of m to the potential at (R_0, z_0) is

$$\delta\Phi_{\text{int}} \equiv \delta\Phi(R_0, z_0) = -4\pi G\rho(m^2)m\delta m \frac{\sqrt{1-e^2}}{e} \sin^{-1}(e). \quad (2.116a)$$

Similarly, if (R_0, z_0) lies outside the homoeoid,

$$\delta\Phi_{\text{ext}} \equiv \delta\Phi(R_0, z_0) = -4\pi G\rho(m^2)m\delta m \frac{\sqrt{1-e^2}}{e} \sin^{-1}(\text{sech } u_m). \quad (2.116b)$$

The potential of the entire body is the sum of contributions (2.116) from all the homoeoids that make up the body. If we define

$$\psi(m) \equiv \int_0^{m^2} dm^2 \rho(m^2), \quad (2.117)$$

the sum of the $\delta\Phi_{\text{int}}$ is

$$\sum_{m>m_0} \delta\Phi_{\text{int}} = -2\pi G \frac{\sqrt{1-e^2}}{e} \sin^{-1}(e) [\psi(\infty) - \psi(m_0)], \quad (2.118)$$

where m_0 is the label of the homoeoid that passes through (R_0, z_0) :

$$m_0^2 \equiv R_0^2 + \frac{z_0^2}{1-e^2}. \quad (2.119)$$

Similarly,

$$\sum_{m<m_0} \delta\Phi_{\text{ext}} = -2\pi G \frac{\sqrt{1-e^2}}{e} \int_0^{m_0^2} dm^2 \rho(m^2) \sin^{-1}(\text{sech } u_m). \quad (2.120)$$

Integrating equation (2.120) by parts,

$$\begin{aligned} \sum_{m<m_0} \delta\Phi_{\text{ext}} &= -2\pi G \frac{\sqrt{1-e^2}}{e} \\ &\times \left\{ [\psi(m) \sin^{-1}(\text{sech } u_m)]_{m=0}^{m_0} - \int_{m=0}^{m_0} \frac{\psi(m) d \text{sech } u_m}{\sqrt{1-\text{sech}^2 u_m}} \right\}. \end{aligned} \quad (2.121)$$

The quantity u_m appearing in equation (2.121) is a function $u_m(R_0, z_0)$ by virtue of the condition that u_m label the spheroid through (R_0, z_0) that is confocal with the homoeoid $m = \text{constant}$. Let the Δ parameter of the confocal family of spheroids containing m be Δ_m , and let u_* be the label of the homoeoid m within this family. Then $m = \Delta_m \cosh u_*$ and $\sqrt{1 - e^2} m = \Delta_m \sinh u_*$, so $\Delta_m = me$, and we have from equation (2.97) that

$$\frac{R_0^2}{\Delta_m^2 \cosh^2 u_m} + \frac{z_0^2}{\Delta_m^2 \sinh^2 u_m} = 1, \quad (2.122)$$

which implies
$$\frac{R_0^2}{1 + \sinh^2 u_m} + \frac{z_0^2}{\sinh^2 u_m} = m^2 e^2.$$

This is the required equation for u_m . Thus, in particular, $m = 0$ implies $\sinh u_m = \infty$, and $m = m_0$ implies $\sinh u_m = \sqrt{1 - e^2}/e$. Inserting these limits into equation (2.121), and adding the result to equation (2.116), we find

$$\begin{aligned} \Phi(R_0, z_0) = & -2\pi G \frac{\sqrt{1 - e^2}}{e} \\ & \times \left(\psi(\infty) \sin^{-1} e - \int_{\sinh u_m = \sqrt{1 - e^2}/e}^{\infty} \psi(m) \frac{d \sinh u_m}{1 + \sinh^2 u_m} \right). \end{aligned} \quad (2.123)$$

We can simplify this equation by defining a new variable of integration

$$\tau \equiv a_0^2 e^2 \left[\sinh^2 u_m - \left(\frac{1}{e^2} - 1 \right) \right], \quad (2.124)$$

where a_0 is any constant. Then equation (2.122) becomes

$$\frac{R_0^2}{\tau + a_0^2} + \frac{z_0^2}{\tau + c_0^2} = \frac{m^2}{a_0^2} \quad (c_0 \equiv \sqrt{1 - e^2} a_0), \quad (2.125a)$$

and equation (2.123) becomes

$$\begin{aligned} \Phi(R_0, z_0) = & -2\pi G \frac{\sqrt{1 - e^2}}{e} \\ & \times \left(\psi(\infty) \sin^{-1} e - \frac{a_0 e}{2} \int_0^{\infty} d\tau \frac{\psi(m)}{(\tau + a_0^2) \sqrt{\tau + c_0^2}} \right). \end{aligned} \quad (2.125b)$$

The integral in this equation gives the contributions to Φ from homoeoids for which (R_0, z_0) is an exterior point, with $\tau = 0$ corresponding to the homoeoid that touches (R_0, z_0) , and large τ corresponding to small homoeoids.

It is instructive to apply equations (2.125) to the determination of the interior potential of a homogeneous spheroid of density ρ_0 and eccentricity e

Table 2.1 Formulae for the dimensionless quantities $I \equiv a_2 a_3 a_1^{-1} \int_0^\infty d\tau \Delta^{-1}$ and $A_i \equiv a_1 a_2 a_3 \int_0^\infty d\tau \Delta^{-1} (a_i^2 + \tau)^{-1}$ that occur in equations (2.128) and Table 2.2. $\Delta^2(\tau) \equiv \prod_{i=1}^3 (a_i^2 + \tau)$. The functions $F(\theta, k)$ and $E(\theta, k)$ are elliptic integrals (Appendix C.4).

	$a_1 = a_2 > a_3$ (oblate)	$a_1 = a_2 < a_3$ (prolate)	$a_1 > a_2 > a_3$ (triaxial)
	$e \equiv \sqrt{1 - a_3^2/a_1^2}$	$e \equiv \sqrt{1 - a_1^2/a_3^2}$	$k \equiv \sqrt{\frac{a_1^2 - a_2^2}{a_1^2 - a_3^2}}; k'^2 \equiv 1 - k^2; \theta \equiv \cos^{-1}\left(\frac{a_3}{a_1}\right)$
I	$2 \frac{\sqrt{1-e^2}}{e} \sin^{-1} e$	$\frac{1}{e} \ln\left(\frac{1+e}{1-e}\right)$	$2 \frac{a_2 a_3}{a_1^2} \frac{F(\theta, k)}{\sin \theta}$
A_1	$\frac{\sqrt{1-e^2}}{e^2} \left[\frac{\sin^{-1} e}{e} - \sqrt{1-e^2} \right]$	$\frac{1-e^2}{e^2} \left[\frac{1}{1-e^2} - \frac{1}{2e} \ln\left(\frac{1+e}{1-e}\right) \right]$	$2 \frac{a_2 a_3}{a_1^2} \frac{F(\theta, k) - E(\theta, k)}{k^2 \sin^3 \theta}$
A_2	$= A_1$	$= A_1$	$2 \frac{a_2 a_3}{a_1^2} \frac{E(\theta, k) - k'^2 F(\theta, k) - (a_3/a_2) k^2 \sin \theta}{k^2 k'^2 \sin^3 \theta}$
A_3	$2 \frac{\sqrt{1-e^2}}{e^2} \left[\frac{1}{\sqrt{1-e^2}} - \frac{\sin^{-1} e}{e} \right]$	$2 \frac{1-e^2}{e^2} \left[\frac{1}{2e} \ln\left(\frac{1+e}{1-e}\right) - 1 \right]$	$2 \frac{a_2 a_3}{a_1^2} \frac{(a_2/a_3) \sin \theta - E(\theta, k)}{k'^2 \sin^3 \theta}$

Table 2.2 Potentials and potential-energy tensors of ellipsoidal bodies

Thin shell	$\Phi(\mathbf{x}_{\text{int}}) = -\frac{Ga_1}{2a_2 a_3} I(\mathbf{a}) M_{\text{shell}}$	$\Phi(\mathbf{x}_{\text{ext}}) = -\frac{Ga_1'}{2a_2' a_3'} I(\mathbf{a}') M_{\text{shell}}$
Homogeneous	$\Phi(\mathbf{x}_{\text{int}}) = -\pi G \rho [I(\mathbf{a}) a_1^2 - \sum_{i=1}^3 A_i(\mathbf{a}) x_i^2]$	$W_{ij} = -\frac{8}{15} \pi^2 G \rho^2 a_1 a_2 a_3 A_i a_i^2 \delta_{ij}$
	$\Phi(\mathbf{x}_{\text{ext}}) = -\pi G \rho \frac{a_1 a_2 a_3}{a_1' a_2' a_3'} [I(\mathbf{a}') a_1'^2 - \sum_{i=1}^3 A_i(\mathbf{a}') x_i^2]$	$W = -\frac{8}{15} \pi^2 G \rho^2 a_1^3 a_2 a_3 I$
Inhomogeneous	$\Phi(\mathbf{x}) = -\pi G \frac{a_2 a_3}{a_1} \int_0^\infty \frac{d\tau}{\Delta} \{ \psi(\infty) - \psi[m(\tau, \mathbf{x})] \}$	$W_{ij} = -2\pi^2 G \frac{a_2 a_3}{a_1^4} \mathcal{S} A_i a_i^2 \delta_{ij} \quad W = -2\pi^2 G \frac{a_2 a_3}{a_1^2} \mathcal{S} I \delta_{ij}$

NOTES: I and A_i as in Table 2.1. \mathbf{x}_{int} and \mathbf{x}_{ext} denote points on the interior or exterior of the ellipsoidal shell or body. If $\sum_{i=1}^3 x_i^2 / [a_i^2 + \lambda(\mathbf{x})] = 1$, then $a_i'^2 \equiv a_i^2 + \lambda(\mathbf{x})$; $\Delta^2(\tau) \equiv \prod_{i=1}^3 (a_i^2 + \tau)$; $m^2(\tau, \mathbf{x}) \equiv a_1^2 \sum_{i=1}^3 x_i^2 / (a_i^2 + \tau)$; $\psi(m) \equiv \int_0^{m^2} \rho(\mathbf{x}) dm^2(0, \mathbf{x})$; $\mathcal{S} \equiv \int_0^\infty dm^2 \rho(m^2) \int_0^{m^2} dm'^2 m' \rho(m'^2) = \frac{1}{2} \int_0^\infty dm [\psi(\infty) - \psi(m)]^2$.

that has semi-axes of lengths a_1 and $a_3 = \sqrt{1 - e^2} a_1$. For this case, equation (2.117) yields

$$\psi(m) = \rho_0 \times \begin{cases} m^2 & (m^2 < a_1^2), \\ a_1^2 & (m^2 \geq a_1^2). \end{cases} \quad (2.126)$$

Equation (2.125b) takes on a particularly simple form if we set the arbitrary constant a_0 equal to a_1 . If (R_0, z_0) lies inside the spheroid, $m(\tau)$ is always smaller than a_0 . Hence, we may substitute from equation (2.125a) and (2.126) into (2.125b), to obtain

$$\begin{aligned} \Phi(R_0, z_0) = & -2\pi G \rho_0 a_1^2 \frac{\sqrt{1 - e^2}}{e} \left[\sin^{-1} e \right. \\ & \left. - \frac{a_1 e}{2} \int_0^\infty \frac{d\tau}{(\tau + a_1^2) \sqrt{\tau + a_3^2}} \left(\frac{R_0^2}{\tau + a_1^2} + \frac{z_0^2}{\tau + a_3^2} \right) \right]. \end{aligned} \quad (2.127)$$

This potential is quadratic in the coordinates and may be written

$$\Phi(\mathbf{x}) = -\pi G \rho_0 (I a_1^2 - A_1 R^2 - A_3 z^2), \quad (2.128)$$

where the dimensionless coefficients I and A_i are given in Table 2.1. An expression for the exterior potential of the homogeneous spheroid is given in Table 2.2.

With the help of equations (2.117) and (2.125) we obtain the gravitational field generated by a spheroidal system as

$$\mathbf{g} = -\nabla\Phi = -\pi G \sqrt{1 - e^2} a_0 \int_0^\infty d\tau \frac{\rho(m^2) \nabla m^2}{(\tau + a_0^2) \sqrt{\tau + c_0^2}}, \quad (2.129a)$$

where

$$\nabla m^2 = 2a_0^2 \left(\frac{R}{\tau + a_0^2} \hat{\mathbf{e}}_R + \frac{z}{\tau + c_0^2} \hat{\mathbf{e}}_z \right). \quad (2.129b)$$

We may use equation (2.129) to find the circular speed $v_c(R)$ in the equatorial plane of an oblate spheroidal galaxy. The radial component of the field (2.129) is

$$g_R(R, z) = -2\pi G \sqrt{1 - e^2} a_0^3 R \int_0^\infty d\tau \frac{\rho(m^2)}{(\tau + a_0^2)^2 \sqrt{\tau + c_0^2}}. \quad (2.130)$$

In the equatorial plane $z = 0$, equation (2.125a) yields

$$m = \frac{a_0 R}{\sqrt{\tau + a_0^2}}. \quad (2.131a)$$

Hence

$$\frac{d\tau}{(\tau + a_0^2)^2} = -\frac{2m}{R^2 a_0^2} dm, \quad (2.131b)$$

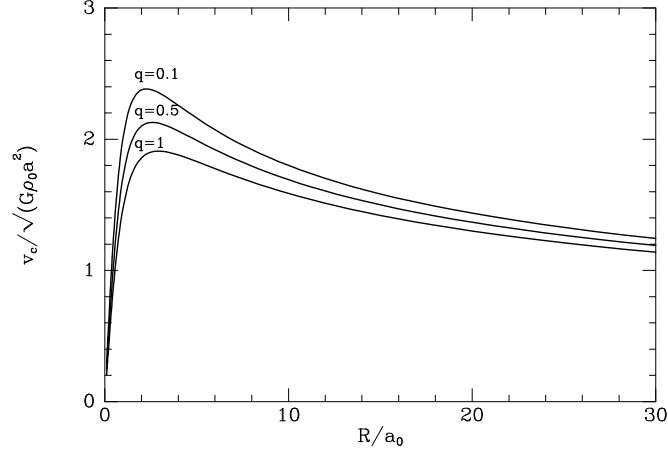


Figure 2.13 Circular speed versus radius for three bodies with the same face-on projected density profile (the modified Hubble model, eq. 2.133) but different axis ratios $q = c/a$. Though all three bodies have the same mass inside a spheroid of given semi-major axis, v_c increases with flattening $1 - q$.

and equation (2.130) yields

$$v_c^2(R) = -Rg_R(R, 0) = 4\pi G\sqrt{1 - e^2} \int_0^R dm \frac{m^2 \rho(m^2)}{\sqrt{R^2 - m^2 e^2}}. \quad (2.132)$$

Let us see how these formulae work out in a specific case. Consider the oblate spheroidal density distribution

$$\rho(m^2) = \rho_0 \left[1 + \left(\frac{m}{a_0} \right)^2 \right]^{-3/2}, \quad (2.133)$$

where a_0 is the core radius and the parameter e that appears in the definition (2.114) of m is the eccentricity of the system. In the limit $e \rightarrow 0$ this reduces to the modified Hubble model (2.53). We substitute for ρ in equation (2.132) to obtain

$$v_c^2(R) = 4\pi G \rho_0 a_0^3 \frac{\sqrt{1 - e^2}}{e} \int_0^R \frac{m^2 dm}{(a_0^2 + m^2)^{3/2} \sqrt{R^2/e^2 - m^2}}. \quad (2.134)$$

By making the substitution

$$m = \frac{R \sin \theta}{e \sqrt{1 + (R/ea_0)^2 \cos^2 \theta}} \quad (2.135)$$

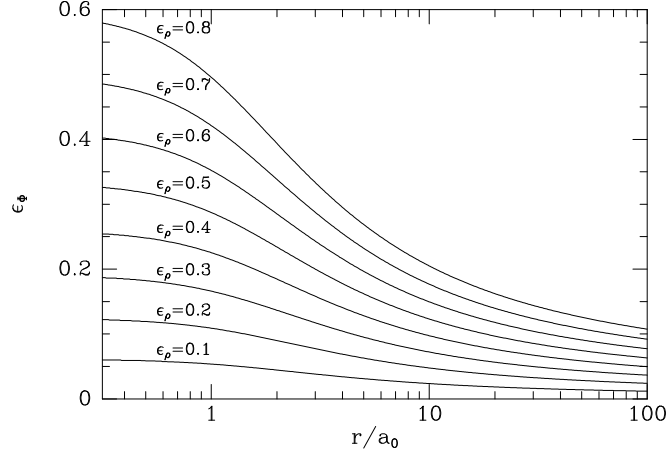


Figure 2.14 The ellipticity ϵ_ϕ of an equipotential surface versus the surface's semi-major axis length r . Each curve is labeled by the ellipticity $\epsilon_\rho = 1 - q$ of the body with density (2.133) that generates the corresponding potential. Notice the rapidity with which the equipotential surfaces become spherical at large r/a_0 .

one may show that the integral of equation (2.134) equals

$$\frac{ek}{R} [F(\theta_m, k) - E(\theta_m, k)], \quad (2.136)$$

where F and E are incomplete elliptic integrals (see Appendix C.4),

$$k \equiv \frac{R}{\sqrt{e^2 a_0^2 + R^2}}, \quad \text{and} \quad \theta_m \equiv \sin^{-1} \sqrt{\frac{e^2 a_0^2 + R^2}{a_0^2 + R^2}}. \quad (2.137a)$$

Hence

$$v_c^2(R) = 4\pi G \rho_0 a_0^3 \frac{\sqrt{1-e^2}}{R} k [F(\theta_m, k) - E(\theta_m, k)]. \quad (2.137b)$$

We may use this result to investigate how strongly a galaxy's circular speed is affected by its shape. In Figure 2.13 we plot the circular-speed curves of three galaxies whose density profiles are given by equation (2.133) for axis ratio $q = \sqrt{1-e^2} = 1$ (spherical system), $q = 0.5$ (E5 galaxy), and $q = 0.1$ (the flatness characteristic of disk galaxies). The central density has in each case been adjusted so as to hold constant the mass $M(a)$ interior to the spheroid of semi-major axis a . The peak circular speed of the $q = 0.1$ model is about 20% higher than that of the spherical system because flattening the system increases the radial component of the force between a given mass element and a test mass in the equatorial plane.

For the density distribution defined by equation (2.133), Figure 2.14 shows the ellipticity ϵ_Φ of the isopotential surfaces for several values of the ellipticity $\epsilon_\rho \equiv 1 - q$ of the density distribution. One sees that in the core $r < a_0$, $\epsilon_\Phi \gtrsim \frac{1}{2}\epsilon_\rho$, while at a few core radii $\epsilon_\Phi \approx \frac{1}{3}\epsilon_\rho$, and at $r \gg a$, ϵ_Φ rapidly approaches zero. Thus, in the region containing the bulk of the mass, the potential is generally flattened only about a third as much as the density, just as we found for logarithmic potentials in §2.3.2.

2.5.3 Potentials of ellipsoidal systems

The problem of calculating the gravitational potential of a body whose isodensity surfaces are similar, coaxial ellipsoids challenged some of the best minds of the eighteenth and nineteenth centuries—see Chandrasekhar (1969) for details. The general problem was solved by George Green (1793–1841), a Nottingham millwright, in 1835 using rather specialized geometrical methods that we shall not present here. Green’s results are natural generalizations of the ones we deduced above by a more accessible technique for axisymmetric systems. We now summarize the results for triaxial systems and refer readers to Kellogg (1953) or Chandrasekhar (1969) for proofs.

On surfaces of constant density, the variable

$$m^2 \equiv a_1^2 \sum_{i=1}^3 \frac{x_i^2}{a_i^2} \quad (2.138)$$

is constant, where (x_1, x_2, x_3) are Cartesian coordinates and a_1, a_2, a_3 are the semi-axes of the ellipsoid. A thin shell of uniform density, whose inner and outer skins are the surfaces m and $m + \delta m$, generates an exterior potential that is constant on the ellipsoidal surfaces

$$m^2 = a_1^2 \sum_{i=1}^3 \frac{x_i^2}{a_i^2 + \tau}, \quad (2.139)$$

where $\tau \geq 0$ labels the surfaces. (This is a straightforward extension of the homoeoid theorem proved in §2.5.1.) There is no gravitational field inside such a shell.

We may find the gravitational potential of any body in which $\rho = \rho(m^2)$ by breaking the body down into thin triaxial homoeoids: the triaxial analog of equation (2.125b) is

$$\Phi(\mathbf{x}) = -\pi G \frac{a_2 a_3}{a_1} \int_0^\infty d\tau \frac{\psi(\infty) - \psi(m)}{\sqrt{(\tau + a_1^2)(\tau + a_2^2)(\tau + a_3^2)}}, \quad (2.140)$$

where $\psi(m)$ is again defined by (2.117) and $m = m(\mathbf{x}, \tau)$ through equation (2.139). Merritt & Fridman (1996) give expressions derived from (2.140)

for the gravitational potentials and fields of triaxial generalizations of the Dehnen models of §2.2.2g.

(a) Ferrers potentials A particularly simple application of equation (2.140) is to the case in which

$$\rho(m^2) = \begin{cases} \rho_0 \left(1 - m^2/a_1^2\right)^n & \text{for } m \leq a_1 \\ 0 & \text{for } m > a_1, \end{cases} \quad (2.141)$$

where $m = m(\mathbf{x})$ through equation (2.138). By (2.117) we now have

$$\psi(\infty) - \psi(m) = \frac{\rho_0 a_1^2}{n+1} \left(1 - \frac{m^2}{a_1^2}\right)^{n+1} \quad (m \leq a_1). \quad (2.142)$$

Hence the internal potential of a body whose density is of the form (2.141) is

$$\begin{aligned} \Phi(\mathbf{x}) = & -\frac{\pi G \rho_0 a_1 a_2 a_3}{n+1} \int_0^\infty \frac{d\tau}{\sqrt{(\tau + a_1^2)(\tau + a_2^2)(\tau + a_3^2)}} \\ & \times \left(1 - \sum_{i=1}^3 \frac{x_i^2}{\tau + a_i^2}\right)^{n+1}. \end{aligned} \quad (2.143)$$

If n is an integer, the bracket involving \mathbf{x} in equation (2.143) can be multiplied out, and the potential at any point obtained as a sum of terms of the form $A_{pqr} x_1^p x_2^q x_3^r$, where the coefficients A_{pqr} are independent of \mathbf{x} . Potentials of this simple form are ideally suited to numerical studies of orbits in triaxial galaxies, such as we shall describe in §3.3. We shall refer to these as **Ferrers potentials**.

The $n = 0$ Ferrers potential arises from a homogeneous ellipsoid with semi-axes a_1, a_2, a_3 . Expressions for the interior and exterior potentials of such bodies can be derived from equations (2.140) and (2.143) and are given in Tables 2.1 and 2.2.

(b) Potential-energy tensors of ellipsoidal systems Roberts (1962) showed that for ellipsoidal bodies equation (2.22) has a remarkably simple form:

$$W_{jk} = -\pi^2 G \frac{a_2 a_3}{a_1^2} \left(\frac{a_j}{a_1}\right)^2 A_j \delta_{jk} \int_0^\infty dm [\psi(\infty) - \psi(m)]^2, \quad (2.144)$$

where the A_j are given in Table 2.1. Notice that the right side of equation (2.144) comprises a constant times the product of two factors: (i) a factor $(a_j/a_1)^2 A_j \delta_{jk}$ that depends only on the axial ratios (a_2/a_1) etc.; and (ii) a factor $\int dm [\psi(\infty) - \psi(m)]^2$ that is independent of the body's ellipticity and the same for all components of the tensor; this integral can be evaluated from a knowledge of the radial density structure alone. In particular, ratios

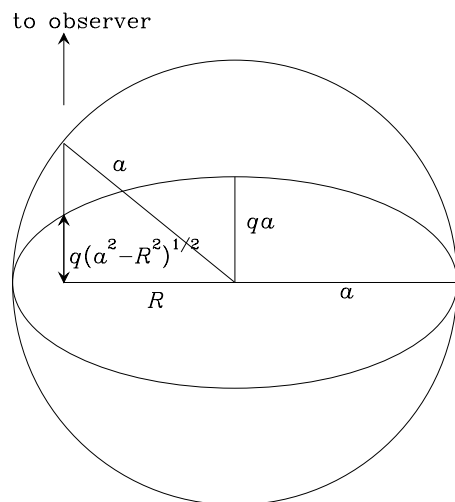


Figure 2.15 A spheroid of axis ratio q and semi-major axis a is viewed along a line of sight that cuts the spheroid's equatorial plane perpendicularly at radius R . This line of sight cuts through the spheroid for a distance $2q\sqrt{a^2 - R^2}$.

of potential-energy terms, for example W_{11}/W_{33} , depend only on the body's ellipticity, and are *entirely independent of the radial density structure* so long as the density is stratified on similar ellipsoids. We shall exploit this useful result in §4.8.3. In Table 2.2 we give expressions for the potential-energy tensors of homogeneous ellipsoids.

2.6 The potentials of disks

Most of the light emitted by a typical spiral galaxy comes from a thin disk. Thus we anticipate that a substantial fraction of the galaxy's mass is concentrated in the disk, and it is therefore important to be able to calculate efficiently the gravitational field of a thin disk. We begin by investigating the potential of an idealized axisymmetric disk of zero thickness.

2.6.1 Disk potentials from homoeoids

We may consider any axisymmetric disk to be a very flat spheroid and use the formulae of §2.5.1 to obtain the potential. A homogeneous spheroid of density ρ , semi-axes a and c , and axial ratio $q = c/a$ has mass $M(a) = \frac{4}{3}\pi\rho qa^3$ (eq. 2.112) and surface density (see Figure 2.15)

$$\Sigma(a, R) = 2\rho q\sqrt{a^2 - R^2}, \quad (2.145)$$

where R is the usual cylindrical radius. Differentiating these expressions with respect to a , we obtain the mass $\delta M(a)$ and the surface density $\delta\Sigma(a, R)$ of

the thin homoeoid of density ρ , semi-major axis a , thickness δa , and axial ratio q :

$$\delta M(a) = 4\pi\rho qa^2\delta a \quad ; \quad \delta\Sigma(a, R) = \frac{2\rho qa}{\sqrt{a^2 - R^2}}\delta a. \quad (2.146)$$

If we now let q tend to zero while holding $2\rho qa \equiv \Sigma_0$ constant, we obtain the mass and surface density of an infinitely flattened homoeoid:

$$\delta M(a) = 2\pi\Sigma_0 a \delta a \quad ; \quad \delta\Sigma(a, R) = \frac{\Sigma_0 \delta a}{\sqrt{a^2 - R^2}}. \quad (2.147)$$

We may construct a razor-thin disk of known surface density $\Sigma(R)$ by finding the family of homoeoids whose combined surface density equals $\Sigma(R)$ at all R . In mathematical language, we have to find the function $\Sigma_0(a)$ that satisfies the integral equation

$$\Sigma(R) = \sum_{a \geq R} \delta\Sigma(a, R) = \int_R^\infty da \frac{\Sigma_0(a)}{\sqrt{a^2 - R^2}}. \quad (2.148a)$$

This is an Abel integral equation. Its solution is (eq. B.72)

$$\Sigma_0(a) = -\frac{2}{\pi} \frac{d}{da} \int_a^\infty dR \frac{R\Sigma(R)}{\sqrt{R^2 - a^2}}. \quad (2.148b)$$

Note that $\Sigma_0(a)$ is *not* the same function as $\Sigma(R)$. In particular, some of the mass that lies interior to radius R comes from homoeoids having $a > R$. By Newton's third theorem, this portion of matter does not contribute to the gravitational force at R , because a point in the equatorial plane at radius R is an interior point of all homoeoids with $a > R$. Thus two disks can have identical surface-density distributions for $R' < R$ and yet have very different force fields at R . In this respect disks differ from spherical distributions of mass, for which the force at r_0 depends only on the density at $r < r_0$. In fact, the surface density of a disk at $R' > R$ affects the attraction at R because the annulus of material exterior to R actually pulls a star placed at radius R outward, thus partially compensating the inward attraction of the interior matter. For example on the perimeter of a sharp-edged disk, the circular speed can be much higher than at the edge of a spherical body, or a more extended disk, with the same mass interior to this point.

We now calculate the potential of a thin disk by adding the potentials of the thin homoeoids into which we have decomposed it. By Gauss's theorem, the gravitational field is discontinuous across a sheet of finite surface density, but the potential is continuous. Consequently, the potential in the equatorial plane differs infinitesimally from the potential just above or below the disk.

Therefore we need only calculate the potential at points that are external to all homoeoids, and take the limit $z \rightarrow 0$ to find the potential in the plane.

Equation (2.107) gives the potential outside a homoeoid of mass M . For a completely flattened homoeoid, we have $e = 1$ and the mass is given in terms of Σ_0 by the first of equations (2.147). Hence, at location (R, z) the potential of this homoeoid is

$$\delta\Phi(R, z) = -2\pi G\Sigma_0\delta a \sin^{-1}(\operatorname{sech} u), \quad (2.149)$$

where u is determined by the equations $R = a \cosh u \sin v$, $z = a \sinh u \cos v$ (eq. 2.96). Eliminating v from these equations by using $1 = \cos^2 v + \sin^2 v$, we obtain a quadratic equation for $\cosh^2 u$:

$$a^2 \cosh^4 u - (R^2 + z^2 + a^2) \cosh^2 u + R^2 = 0. \quad (2.150)$$

The root that we require is the one with $\cosh^2 u \geq 1$, which is

$$\begin{aligned} \cosh^2 u &= \frac{1}{2a^2} \left[R^2 + z^2 + a^2 + \sqrt{(R^2 + z^2 + a^2)^2 - 4a^2 R^2} \right] \\ &= \frac{1}{4a^2} \left[\sqrt{z^2 + (a+R)^2} + \sqrt{z^2 + (a-R)^2} \right]^2. \end{aligned} \quad (2.151)$$

Taking the square root of both sides and substituting the result into equation (2.149), we obtain

$$\delta\Phi = -2\pi G\Sigma_0 \delta a \sin^{-1} \left(\frac{2a}{\sqrt{z^2 + (a+R)^2} + \sqrt{z^2 + (a-R)^2}} \right). \quad (2.152)$$

Finally, the potential of an axisymmetric disk of arbitrary surface-density profile is obtained by combining this result with equations (2.148), which decompose a disk into homoeoids. We have (Cuddeford 1993)

$$\Phi(R, z) = 4G \int_0^\infty da \sin^{-1} \left(\frac{2a}{\sqrt{+} + \sqrt{-}} \right) \frac{d}{da} \int_a^\infty dR' \frac{R'\Sigma(R')}{\sqrt{R'^2 - a^2}}. \quad (2.153a)$$

where

$$\sqrt{\pm} \equiv \sqrt{z^2 + (a \pm R)^2}. \quad (2.153b)$$

An alternative form is obtained by integrating by parts:

$$\begin{aligned} \Phi(R, z) &= -2\sqrt{2}G \int_0^\infty da \frac{[(a+R)/\sqrt{+}] - [(a-R)/\sqrt{-}]}{\sqrt{R^2 - z^2 - a^2 + \sqrt{+}\sqrt{-}}} \\ &\quad \times \int_a^\infty dR' \frac{R'\Sigma(R')}{\sqrt{R'^2 - a^2}}. \end{aligned} \quad (2.154)$$

This form does not require differentiation of the term that depends on the surface density. Equations (2.153a) and (2.154) are numerically convenient because the inner integrals depend only on a , and thus can be tabulated on a grid of values of a at the outset. Then only a single integral is required to evaluate Φ at each fresh point (R, z) . Moreover, the integrands do not oscillate in sign, so numerical integrations converge rapidly.

We are particularly interested in the value of the potential in the equatorial plane. Consider first the case $a > R$. For small z , it is easy to see that $\sqrt{\pm z} \rightarrow (a \pm R) + O(z^2)$, with the consequence that the numerator of the fraction in the first line in equation (2.154) vanishes like z^2 as $z \rightarrow 0$. In the same limit the denominator vanishes too, but more slowly, like z . Hence, for $a > R$ the integrand tends to zero with z and we readily find that

$$\Phi(R, 0) = -4G \int_0^R \frac{da}{\sqrt{R^2 - a^2}} \int_a^\infty dR' \frac{R'\Sigma(R')}{\sqrt{R'^2 - a^2}}. \quad (2.155)$$

To obtain the circular speed at radius R in the equatorial plane, the natural procedure is to differentiate this expression with respect to R , which is the upper limit of the outer integral. The usual formula for differentiating such an integral requires us to evaluate the integrand at $a = R$, when it diverges. Instead of resolving this problem we return to equation (2.153a), in which we may put z to zero without any awkwardness, to obtain

$$\Phi(R, 0) = 4G \int_0^\infty da \sin^{-1} \left(\frac{2a}{(a+R) + |a-R|} \right) \frac{d}{da} \int_a^\infty dR' \frac{R'\Sigma(R')}{\sqrt{R'^2 - a^2}}. \quad (2.156)$$

The argument of \sin^{-1} is unity for $R < a$ and a/R otherwise. Consequently, when we differentiate with respect to R we obtain

$$v_c^2(R) = R \frac{\partial \Phi}{\partial R} = -4G \int_0^R da \frac{a}{\sqrt{R^2 - a^2}} \frac{d}{da} \int_a^\infty dR' \frac{R'\Sigma(R')}{\sqrt{R'^2 - a^2}}. \quad (2.157)$$

(a) The Mestel disk As a simple application of equation (2.157) consider a disk in which the surface density is given by

$$\Sigma(R) = \begin{cases} \frac{v_0^2}{2\pi GR} & \text{for } R < R_{\max} \\ 0 & \text{otherwise,} \end{cases} \quad (2.158)$$

where v_0 and R_{\max} are constants with dimensions of velocity and length, respectively. For this surface density the inner integral in equation (2.157) is proportional to $\cosh^{-1}(R_{\max}/a)$, so its derivative with respect to a follows from

$$\frac{d}{da} \cosh^{-1}(R_{\max}/a) = -\frac{R_{\max}}{a\sqrt{R_{\max}^2 - a^2}}. \quad (2.159)$$

We let the outer radius of the disk, R_{\max} , tend to infinity. In this limit the disk becomes the **Mestel disk** (Mestel 1963) and the derivative (2.159) becomes $-1/a$. Substituting this value into equation (2.157), we find that the circular speed in the Mestel disk is

$$v_c^2 = \frac{2v_0^2}{\pi} \int_0^R \frac{da}{\sqrt{R^2 - a^2}} = v_0^2. \quad (2.160)$$

Hence, the circular speed of a disk in which the surface density is inversely proportional to radius is independent of radius. Moreover, for this surface-density law, $v_c(R)$ is given by the simple formula

$$v_c^2(R) = \frac{GM(R)}{R}, \quad (2.161a)$$

where

$$M(R) = 2\pi \int_0^R dR' R' \Sigma(R') = \frac{v_0^2 R}{G} \quad (2.161b)$$

is the mass interior to R . This is precisely analogous to equation (2.29) for a spherical system. Although we have argued that for general disks the circular speed is affected by the mass exterior to R , for the particular case of the Mestel disk the simple formula (2.161a) happens to give the correct answer—we are aware of no other disks with this property.

(b) The exponential disk The surface-brightness profiles of many galactic disks are approximately exponential in form (eq. 1.7). Let us use the results just derived to calculate the potential that such a disk would generate if its surface mass density were also exponential. Setting

$$\Sigma(R) = \Sigma_0 e^{-R/R_d}, \quad (2.162)$$

we use equation (C.69) to show that the inner integral in equations (2.153a) to (2.157) is

$$\int_a^\infty dR' \frac{R' \Sigma_0 e^{-R'/R_d}}{\sqrt{R'^2 - a^2}} = \Sigma_0 a K_1(a/R_d), \quad (2.163)$$

where K_1 is a modified Bessel function (Appendix C.7). Figure 2.16 shows contours of the potential that one obtains by substituting this formula into equation (2.154).

The potential in the equatorial plane is given by substituting (2.163) into equation (2.155):

$$\begin{aligned} \Phi(R, 0) &= -4G\Sigma_0 \int_0^R da \frac{aK_1(a/R_d)}{\sqrt{R^2 - a^2}} \\ &= -\pi G\Sigma_0 R [I_0(y)K_1(y) - I_1(y)K_0(y)], \end{aligned} \quad (2.164a)$$

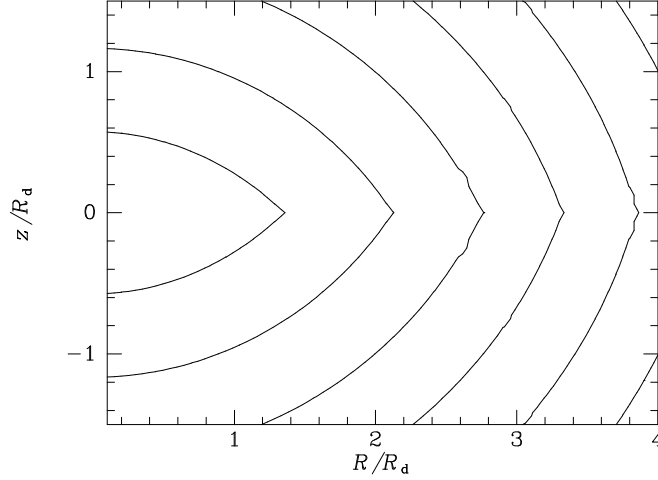


Figure 2.16 Contours in the (R, z) plane of constant potential Φ for a razor-thin exponential disk. The contour levels are GM_d/R_d divided by 1.5, 2, 2.5, \dots , where $M_d = 2\pi\Sigma_0 R_d^2$ is the mass of the disk.

where

$$y \equiv \frac{R}{2R_d}, \quad (2.164b)$$

we have used equation (C.70), and the I_n are modified Bessel functions (Appendix C.7).

If we differentiate equation (2.164a) with respect to R , we obtain the circular speed of the exponential disk (Freeman 1970):

$$v_c^2(R) = R \frac{\partial \Phi}{\partial R} = 4\pi G \Sigma_0 R_d y^2 [I_0(y)K_0(y) - I_1(y)K_1(y)]. \quad (2.165)$$

In Figure 2.17 we show this circular speed together with the circular speed of the spherical body that has as much mass $M_s(r)$ interior to $r = R$ as the exponential disk, that is,

$$\begin{aligned} M_s(R) = M_d(R) &= 2\pi \int_0^R dR' R' \Sigma_0 e^{-R'/R_d} \\ &= 2\pi \Sigma_0 R_d^2 \left[1 - e^{-R/R_d} \left(1 + \frac{R}{R_d} \right) \right]. \end{aligned} \quad (2.166)$$

The exponential disk achieves a peak circular speed that is about 15% higher than that of the equivalent spherical distribution. The dotted line in Figure 2.17 gives the Keplerian circular speed for a system in which the entire mass of the disk is concentrated at the center. Notice that the disk's circular

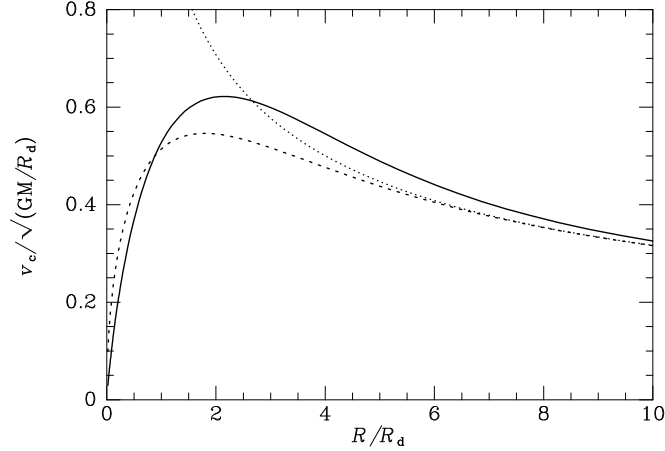


Figure 2.17 The circular-speed curves of: an exponential disk (full curve); a point with the same total mass (dotted curve); the spherical body for which $M(r)$ is given by equation (2.166) (dashed curve).

speed approaches the Keplerian speed only slowly and from above, whereas the circular speed of the equivalent spherical distribution tends rapidly to the Keplerian speed from below.

(c) Thick disks Although galactic disks are thin in the sense that the density falls off much faster perpendicular to the equatorial plane than in the radial direction within the plane, it is frequently essential to take into account the non-zero thickness of the disk in the perpendicular direction. For example, the dynamics of the solar neighborhood, and of the Sun itself, would be very different if the Galactic disk were razor-thin rather than some hundreds of parsecs thick. The techniques for deriving the potential of a spheroidal system that are described in §2.5.1 are not directly applicable to this problem because the density of a galactic disk is generally not constant on spheroids. The multipole expansion described in §2.4 is also not suitable, because it converges slowly for flat systems. Here we describe how to calculate the potential of a disk when the density is of the form

$$\rho(R, z) = \Sigma(R)\zeta(z). \quad (2.167)$$

Physically, this formula implies that a cross-section through the disk always has the same shape, no matter at what radius it is taken. In particular, the characteristic scale height of the disk is independent of R , an assumption that is in reasonable agreement with observations of edge-on disk galaxies (BM §4.4.3).

It is convenient to normalize the function $\zeta(z)$ in equation (2.167) such that $\int dz \zeta(z) = 1$. With this normalization, $\Sigma(R)$ is the total surface density

and $\Sigma(R)\zeta(z) dz$ is the surface density of the layer of material of thickness dz that lies a distance z above the equatorial plane. Let $\Phi_0(R, z)$ be the potential that would be generated by a razor-thin disk with surface density $\Sigma(R)$ that lay in the plane $z = 0$. Then the potential generated at (R, z) by the layer at distance z' from the plane is

$$d\Phi(R, z) = dz' \Phi_0(R, z - z')\zeta(z'). \quad (2.168)$$

Adding the contributions to the disk's potential from every layer, we obtain for the overall potential

$$\Phi(R, z) = \int_{-\infty}^{\infty} dz' \zeta(z')\Phi_0(R, z - z'), \quad (2.169)$$

where Φ_0 can be obtained from either (2.153a) or (2.154).

Consider, for example, the exponential disk (2.162). Then inserting into equation (2.153a) the derivative $-\Sigma_0(a/R_d)K_0(a/R_d)$ with respect to a of equation (2.163), and then substituting the resulting value of Φ_0 into equation (2.169), we find

$$\Phi(R, z) = -\frac{4G\Sigma_0}{R_d} \int_{-\infty}^{\infty} dz' \zeta(z') \int_0^{\infty} da \sin^{-1} \left(\frac{2a}{\sqrt{+} + \sqrt{-}} \right) a K_0(a/R_d), \quad (2.170)$$

where $\sqrt{\pm}$ is defined by equation (2.153b) with z replaced by $z - z'$. We shall use this formula in §2.7 below.

2.6.2 Disk potentials from Bessel functions

Disk galaxies generally contain non-axisymmetric features such as a bar or spiral arms. Hence it is essential to know how to calculate the potential of a flattened, non-axisymmetric system. Several methods may be employed and none is ideal for every problem. Evans & de Zeeuw (1992) show how the results of §2.5.3 may be used to obtain the potential of razor-thin, elliptical disks in analogy with the work of §2.6.1. However, the technique they present is not easy to apply, so in the rest of this section we present three alternatives, each with its own strengths when applied to particular problems. All these methods can also be applied to axisymmetric disks, and all yield alternatives to the formulae of §2.5.

Above and below an isolated razor-thin disk the gravitational potential satisfies Laplace's equation, $\nabla^2\Phi = 0$, with appropriate boundary conditions on the disk and at infinity. In cylindrical coordinates Laplace's equation is (eq. B.52)

$$\frac{1}{R} \frac{\partial}{\partial R} \left(R \frac{\partial \Phi}{\partial R} \right) + \frac{1}{R^2} \frac{\partial^2 \Phi}{\partial \phi^2} + \frac{\partial^2 \Phi}{\partial z^2} = 0. \quad (2.171)$$

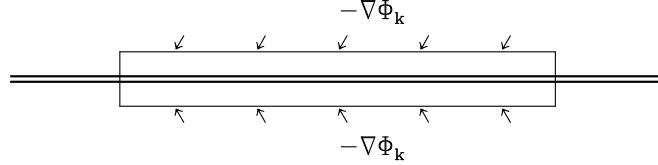


Figure 2.18 The disk mass within the box shown in cross-section equals $-(4\pi G)^{-1}$ times the integral of the normal component of $\nabla\Phi_{km}$ over the surface of the box. The horizontal component of $\nabla\Phi_{km}$ is due to the gravitational attraction from the rest of the galaxy.

Writing

$$\Phi(R, z) = J(R)F(\phi)Z(z), \quad (2.172)$$

we obtain by the method of separation of variables (see §2.4)

$$\frac{1}{J(R)R} \frac{d}{dR} \left(R \frac{dJ}{dR} \right) + \frac{1}{F(\phi)R^2} \frac{d^2 F}{d\phi^2} = -\frac{1}{Z(z)} \frac{d^2 Z}{dz^2} = -k^2, \quad (2.173)$$

where k is an arbitrary real or complex number. Thus

$$\begin{aligned} 0 &= \frac{d^2 Z}{dz^2} - k^2 Z, \\ 0 &= \frac{1}{J(R)R} \frac{d}{dR} \left(R \frac{dJ}{dR} \right) + \frac{1}{F(\phi)R^2} \frac{d^2 F}{d\phi^2} + k^2. \end{aligned} \quad (2.174)$$

The first of equations (2.174) may be immediately integrated to

$$Z(z) = e^{\pm kz}. \quad (2.175)$$

When multiplied by R^2 , the second of equations (2.174) separates into equations for J and F :

$$\begin{aligned} 0 &= \frac{d^2 F}{d\phi^2} + m^2 F \\ 0 &= R \frac{d}{dR} \left(R \frac{dJ}{dR} \right) + k^2 R^2 J(R) - m^2 J. \end{aligned} \quad (2.176)$$

The first of these equations trivially yields $F \propto e^{im\phi}$, where m is an integer. If we make the substitution $u \equiv kR$, the second equation simplifies to

$$u \frac{d}{du} \left(u \frac{dJ}{du} \right) + (u^2 - m^2)J(u) = 0. \quad (2.177)$$

We require the solution of this equation to be finite at $u = 0$ ($R = 0$). This solution is conventionally written $J_m(u) = J_m(kR)$ and is called the Bessel function of order m (Appendix C.7). Summarizing these results, we have that the functions

$$\Phi_{\pm}(R, z) = e^{i(m\phi \pm kz)} J_m(kR) \quad (2.178)$$

are solutions of $\nabla^2 \Phi = 0$.

Now consider the function

$$\Phi_{km}(R, z) = e^{im\phi - k|z|} J_m(kR), \quad (2.179a)$$

where k is real and positive. Φ_{km} satisfies all the conditions required for it to be the potential generated by an isolated density distribution: it is non-singular at $R = 0$, periodic in ϕ and vanishes at large distances from the origin. Furthermore, for $z > 0$, Φ_{km} coincides with Φ_- , and for $z < 0$, Φ_{km} coincides with Φ_+ . Therefore, Φ_{km} solves $\nabla^2 \Phi = 0$ everywhere except in the plane $z = 0$. At $z = 0$, Φ_{km} does not satisfy Laplace's equation because its gradient suffers a discontinuity. Figure 2.18 illustrates how we may use Gauss's theorem (eq. 2.12) to evaluate the surface density $\Sigma_{km}(R, \phi)$ of the sheet that generates this discontinuity. We have that

$$\lim_{z \rightarrow 0^+} \frac{\partial \Phi_{km}}{\partial z} = -k e^{im\phi} J_m(kR) \quad \text{and} \quad \lim_{z \rightarrow 0^-} \frac{\partial \Phi_{km}}{\partial z} = +k e^{im\phi} J_m(kR). \quad (2.180)$$

The integral of $\nabla \Phi_{km}$ over the closed unit surface that is shown in the figure must equal $4\pi G \Sigma_{km}$ from which it follows that

$$\Sigma_{km}(R, \phi) = -\frac{k}{2\pi G} e^{im\phi} J_m(kR). \quad (2.179b)$$

We now use equations (2.179) to find the potential generated by a disk of arbitrary surface density $\Sigma(R, \phi)$. If we can find functions $S_m(k)$ such that

$$\begin{aligned} \Sigma(R, \phi) &= \sum_{m=-\infty}^{\infty} \int_0^{\infty} dk S_m(k) \Sigma_{km}(R, \phi) \\ &= -\frac{1}{2\pi G} \sum_{m=-\infty}^{\infty} \int_0^{\infty} dk k S_m(k) e^{im\phi} J_m(kR), \end{aligned} \quad (2.181)$$

then we will have

$$\begin{aligned} \Phi(R, \phi, z) &= \sum_{m=-\infty}^{\infty} \int_0^{\infty} dk S_m(k) \Phi_{km}(R, \phi, z) \\ &= \sum_{m=-\infty}^{\infty} \int_0^{\infty} dk S_m(k) J_m(kR) e^{im\phi - k|z|}. \end{aligned} \quad (2.182)$$

Multiplying (2.181) through by $e^{-im'\phi}$ and averaging over ϕ , we obtain

$$\Sigma_{m'}(R) \equiv \frac{1}{2\pi} \int_0^{2\pi} d\phi e^{-im'\phi} \Sigma(R, \phi) = -\frac{1}{2\pi G} \int_0^{\infty} dk k S_{m'}(k) J_{m'}(kR). \quad (2.183)$$

Equation (2.183) states that $S_m(k)$ is the m th-order Hankel transform of $-2\pi G\Sigma_m$ (eq. C.60b). Hankel transforms have properties that are similar to those of the familiar Fourier transforms (Appendix B.4). In particular, they may be inverted by use of equation (C.60a). We find

$$S_m(k) = -2\pi G \int_0^\infty dR R J_m(kR) \Sigma_m(R). \quad (2.184)$$

When we eliminate $S_m(k)$ between this equation and (2.182), we obtain finally

$$\Phi(R, \phi, z) = -2\pi G \sum_{m=-\infty}^{\infty} \int_0^\infty dk e^{im\phi - k|z|} J_m(kR) \int_0^\infty dR' R' J_m(kR') \Sigma_m(R'). \quad (2.185)$$

Application to axisymmetric disks Potential-density pairs for axisymmetric disks are obtained by setting $z = 0$ in equation (2.182) and restricting the sum to the case $m = 0$. We have (Toomre 1963)

$$\Phi(R, 0) = \int_0^\infty dk S_0(k) J_0(kR). \quad (2.186)$$

Differentiating with respect to R and using the identity $dJ_0(x)/dx = -J_1(x)$ (eq. C.58), we obtain

$$v_c^2(R) = R \frac{\partial \Phi}{\partial R} = -R \int_0^\infty dk k S_0(k) J_1(kR). \quad (2.187)$$

Substituting for $S_0(k)$ from equation (2.184) this can be rewritten

$$v_c^2(R) = 2\pi GR \int_0^\infty dk k J_1(kR) \int_0^\infty dR' R' \Sigma(R') J_0(kR'). \quad (2.188)$$

Applying to equation (2.187) the inversion formula for Hankel transforms (eqs. C.60) we find

$$S_0(k) = - \int_0^\infty dR' v_c^2(R') J_1(kR'). \quad (2.189)$$

Substituting this expression for S_0 into equation (2.181), we have

$$\Sigma(R) = \frac{1}{2\pi G} \int_0^\infty dk k J_0(kR) \int_0^\infty dR' v_c^2(R') J_1(kR'). \quad (2.190)$$

Comparison of this formula with our expression (2.188) for the reverse relation reveals complete symmetry between the quantities $v_c(R)$ and $2\pi GR\Sigma(R)$

(Kalnajs 1999). Hence our mathematics seems to be saying that it is as easy to determine a disk's surface density from measurements of its circular speed, as to obtain the circular speed from the surface density. Unfortunately, observational constraints destroy this symmetry. The key point is that the left side of either equation (2.188) or (2.190) can be determined at any given value of R only if the variable on the right side can be measured out to radii at which its value becomes negligible. The surface density declines rapidly with radius, so equation (2.188) can be used to derive accurate values of v_c . Circular speeds, by contrast, decline little if at all out to the largest observable radii. Consequently, in practice we cannot obtain the data needed to determine Σ accurately from equation (2.190).

2.6.3 Disk potentials from logarithmic spirals

An alternative technique for finding non-axisymmetric potential-density pairs was introduced by Kalnajs (1971). The potential $\Phi(R, \phi)$ at any point in the plane of a disk is

$$\begin{aligned}\Phi(R, \phi) &= -G \int dR' R' \int d\phi' \frac{\Sigma(R', \phi')}{|\mathbf{x} - \mathbf{x}'|} \\ &= -G \int_0^\infty dR' R' \int_0^{2\pi} d\phi' \frac{\Sigma(R', \phi')}{\sqrt{R'^2 + R^2 - 2RR' \cos(\phi' - \phi)}}.\end{aligned}\quad (2.191)$$

The integral in this expression can be simplified if we define a new radial coordinate,

$$u \equiv \ln R, \quad (2.192)$$

and introduce the **reduced potential** V and the **reduced surface density** S by

$$\begin{aligned}R^{1/2}\Phi &\equiv V(u, \phi) = e^{u/2}\Phi[R(u), \phi] \\ R^{3/2}\Sigma &\equiv S(u, \phi) = e^{3u/2}\Sigma[R(u), \phi].\end{aligned}\quad (2.193)$$

With these substitutions (2.191) becomes

$$V(u, \phi) = -G \int_{-\infty}^\infty du' \int_0^{2\pi} d\phi' K(u - u', \phi - \phi') S(u', \phi'), \quad (2.194a)$$

where

$$K(u - u', \phi - \phi') \equiv \frac{1}{\sqrt{2}\sqrt{\cosh(u - u') - \cos(\phi - \phi')}}. \quad (2.194b)$$

Now consider the reduced potential $V_{\alpha m}(u, \phi)$ that is generated by the particular reduced surface density

$$S_{\alpha m}(u, \phi) = e^{i(\alpha u + m\phi)}, \quad (2.195)$$

where α is a real number and m is an integer. We have

$$\begin{aligned} V_{\alpha m}(u, \phi) &= -G \int_{-\infty}^{\infty} du' \int_0^{2\pi} d\phi' K(u - u', \phi - \phi') e^{i(\alpha u' + m\phi')} \\ &= -G e^{i(\alpha u + m\phi)} \int_{-\infty}^{\infty} du' \int_0^{2\pi} d\phi' K(u - u', \phi - \phi') e^{i[\alpha(u' - u) + m(\phi' - \phi)]}. \end{aligned} \quad (2.196)$$

If we change to new variables of integration $u'' \equiv u - u'$ and $\phi'' \equiv \phi - \phi'$, equation (2.196) becomes

$$V_{\alpha m} = -GN(\alpha, m)e^{i(\alpha u + m\phi)}, \quad (2.197)$$

where

$$\begin{aligned} N(\alpha, m) &\equiv \int_{-\infty}^{\infty} du'' \int_0^{2\pi} d\phi'' K(u'', \phi'') e^{-i(\alpha u'' + m\phi'')} \\ &= \pi \frac{(\frac{1}{2}m - \frac{3}{4} + \frac{1}{2}i\alpha)! (\frac{1}{2}m - \frac{3}{4} - \frac{1}{2}i\alpha)!}{(\frac{1}{2}m - \frac{1}{4} + \frac{1}{2}i\alpha)! (\frac{1}{2}m - \frac{1}{4} - \frac{1}{2}i\alpha)!}. \end{aligned} \quad (2.198)$$

The kernel $N(\alpha, m)$ is real and even in both α and m .⁷ The reduced potential generated by an arbitrary linear combination

$$S(u, \phi) \equiv \sum_{m=-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{d\alpha}{2\pi} A_m(\alpha) e^{i(\alpha u + m\phi)} \quad (2.199a)$$

of surface densities of the form (2.195) is

$$V(u, \phi) = -G \sum_{m=-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{d\alpha}{2\pi} N(\alpha, m) A_m(\alpha) e^{i(\alpha u + m\phi)}. \quad (2.199b)$$

Furthermore, equation (2.199a) states that $A_m(\alpha)$ is nothing but the Fourier transform of the reduced surface density $S(u, \phi)$ (eqs. B.64 and B.67). Consequently,

$$A_m(\alpha) = \frac{1}{2\pi} \int_{-\infty}^{\infty} du \int_0^{2\pi} d\phi S(u, \phi) e^{-i(\alpha u + m\phi)}. \quad (2.199c)$$

So we may use equations (2.199) to obtain the potential in the plane $z = 0$ that is generated by any distribution of surface density. Since $\alpha u + m\phi = \alpha \ln R + m\phi$ is constant on **logarithmic spirals**, equations (2.199) determine the potential of a disk by decomposing the density into spirals. Note that this derivation does not produce an expression for the value of the potential away from the plane $z = 0$ (Problem 2.19 remedies this defect).

⁷ These statements can be proved using equations (C.12) and (C.14).

2.6.4 Disk potentials from oblate spheroidal coordinates

In some situations we require the potential of a disk with a sharp outer edge. If the outer edge is circular, this problem can be efficiently solved using the oblate spheroidal coordinates of §2.5.1.⁸ We substitute $\Phi = U(u)V(v)e^{im\phi}$ into Laplace's equation using (2.101), and separate variables in the usual way. Then we find that U and V satisfy

$$\frac{m^2}{\sin^2 v} - \frac{1}{V \sin v} \frac{d}{dv} \left(\sin v \frac{dV}{dv} \right) = l(l+1), \quad (2.200a)$$

$$\frac{m^2}{\cosh^2 u} + \frac{1}{U \cosh u} \frac{d}{du} \left(\cosh u \frac{dU}{du} \right) = l(l+1), \quad (2.200b)$$

where $l(l+1)$ is the separation constant. The left side of equation (2.200a) is the same as the right side of equation (2.80) with $V(v)$ substituted for $P(\theta)$. Furthermore, the boundary conditions at $v = 0$ or π (i.e., along the z axis) that must be satisfied by a physically acceptable function $V(v)$ are the same as the conditions we imposed on $P(\theta)$. Hence

$$V(v, \phi) = V_{lm} Y_l^m(v, \phi), \quad (2.201)$$

where V_{lm} is a constant and Y_l^m is the spherical harmonic defined by equation (C.42). Since the potential must be symmetrical about the plane of the disk $v = \pi/2$, we must restrict ourselves to values of l and m for which $Y_l^m(v, \phi)$ is an even function of $\cos v$. Hence we require $l - m$ to be even.

If we change the independent variable in equation (2.200b) to $x = i \sinh u$, the equation becomes the associated Legendre equation (2.81b) with x now pure imaginary. If the potential is to vanish at infinity (large u), it must be proportional to the solution of equation (2.81b) that vanishes at large x . This is written $Q_l^m(x)$ (see Appendix C.5). Hence the functions

$$\Phi_{lm}(u, v, \phi) \equiv \frac{V_{lm}}{Q_l^m(0)} Q_l^m(i \sinh u) Y_l^m(v, \phi) \quad (l - m \text{ even}) \quad (2.202)$$

satisfy Laplace's equation everywhere outside the excluded disk $u = 0$ ($z = 0$, $R \leq \Delta$) and vanish at infinity. However, there is a discontinuity in the gradient of Φ_{lm} on the excluded disk. By Gauss's theorem, this discontinuity is generated by a surface density $\Sigma_{lm}(v, \phi)$; thus from equation (2.100) we have

$$\begin{aligned} 4\pi G \Sigma_{lm}(v, \phi) &= 2(\hat{\mathbf{e}}_u \cdot \nabla \Phi_{lm})_{u \rightarrow 0^+} = \frac{2}{\Delta |\cos v|} \lim_{u \rightarrow 0^+} \left(\frac{\partial \Phi_{lm}}{\partial u} \right) \\ &= 2V_{lm} i \lim_{x \rightarrow 0} \left[\frac{d \ln Q_l^m(x)}{dx} \right] \frac{Y_l^m(v, \phi)}{\Delta |\cos v|}, \end{aligned} \quad (2.203)$$

⁸In the more general case in which the outer edge is an ellipse, one uses ellipsoidal coordinates (Morse & Feshbach 1953; Tremaine 1976b; Evans & de Zeeuw 1992).

where x approaches zero along the positive imaginary axis. Using equations (C.31) and (C.13) to evaluate the limit in equation (2.203), we then find (Hunter 1963)

$$\Sigma_{lm} = -\frac{2V_{lm}}{\pi^2 G \Delta g_{lm}} \frac{Y_l^m(v, \phi)}{|\cos v|}, \quad (2.204a)$$

where

$$g_{lm} \equiv \frac{(l+m)!(l-m)!}{2^{2l-1} \left[\left(\frac{l+m}{2}\right)! \left(\frac{l-m}{2}\right)! \right]^2}. \quad (2.204b)$$

A general disk potential, which is a sum over l and m of potentials of the form (2.202), is generated by the surface density $\Sigma(v, \phi)$ that is the sum of surface densities $\Sigma_{lm}(v, \phi)$. According to equation (2.204a), $-2V_{lm}/(\pi^2 G \Delta g_{lm})$ is the coefficient of $Y_l^m(v, \phi)$ when $|\cos v|\Sigma(v, \phi)$ is expanded in spherical harmonics. Thus with the orthogonality relation (C.44) we have

$$\frac{2V_{lm}}{\pi^2 G \Delta g_{lm}} = -\int_0^{2\pi} d\phi \int_0^\pi dv \sin v |\cos v| \Sigma(v, \phi) Y_l^{m*}(v, \phi). \quad (2.205)$$

The integrand in equation (2.205) is symmetrical about $v = \pi/2$ when $l - m$ is even, so we may restrict the v integration to the range $(0, \pi/2)$ and double the result. Hence

$$V_{lm} = -\frac{\pi^2 G g_{lm}}{\Delta} \int_0^{2\pi} d\phi \int_0^\Delta dR R \Sigma(R, \phi) Y_l^{m*}(\sin^{-1}(R/\Delta), \phi). \quad (2.206)$$

All the techniques in this section are special cases of a general method for finding disk potential-density pairs that is described by Qian (1992).

2.7 The potential of our Galaxy

In this section we investigate the gravitational field of our own galaxy, the Milky Way. The Galaxy is made of several components, the disk, the bulge, the stellar halo, and the dark halo. The mix of stars, gas and dark matter that makes up a galaxy such as our own varies from component to component and is even likely to depend on location within each component.

Ideally, we should rely solely on dynamical tracers, such as the velocity fields of gas and stars and observations of gravitational lensing, to map out the distribution of mass in the Galaxy. Sadly, at the present time such a project is unfeasible.

Since we are not yet in a position to model the Galactic density and gravitational field in a purely dynamical way, we flesh out the available dynamical constraints with photometric information. In particular, we simply

assume that each component has a mass-to-light ratio Υ that is independent of position. For the reason given above, this procedure is arbitrary and unsatisfactory, but it yields concrete Galactic potentials, which make testable predictions regarding the kinematics of stars and gas. Proceeding in this spirit, we now investigate models of our Galaxy, following Dehnen & Binney (1998a) and BM §10.6.

The brightness distribution of each component is assumed to be similar to those of external galaxies (BM §§4.3 and 4.4), while the size and total luminosity of each component is determined from photometry and star counts, or by fitting to the available dynamical constraints. We do not model the stellar halo here since its contribution to the potential is negligible.

The models are constrained by fitting to the following data (cf. Table 1.2):

- (i) The circular-speed curve $v_c(R)$ for an assumed value of the solar circular speed, $v_0 \equiv v_c(R_0)$. Since this curve is determined from the line-of-sight velocities of tracers such as HI clouds and Cepheid stars, the circular-speed curve depends on v_0 , which must be determined by other methods.
- (ii) The values of the Oort constants (Table 1.2, eq. 3.83, and BM §10.3.3).
- (iii) The total surface density within 1.1 kpc of the Galactic plane near the Sun, $\Sigma_{1.1}(R_0)$, and the contribution of the disk to this density (Table 1.1).
- (iv) The velocity dispersion of bulge stars in Baade's window, a line of sight that passes ~ 500 pc from the Galactic center in which absorption by intervening dust is unusually low. We take this dispersion to be $117 \pm 15 \text{ km s}^{-1}$.
- (v) The total mass within 100 kpc of the Galactic center (eq. 1.12).
- (vi) The solar radius $R_0 = 8$ kpc.

The functional forms assumed for each of the Galaxy's components are as follows.

(a) The bulge The density of this component is assumed to be

$$\rho_b(R, z) = \rho_{b0} \left(\frac{m}{a_b} \right)^{-\alpha_b} e^{-m^2/r_b^2}, \quad (2.207a)$$

where

$$m = \sqrt{R^2 + z^2/q_b^2}. \quad (2.207b)$$

For $q_b < 1$ this is an oblate, spheroidal power-law model that is truncated at an outer radius r_b . Its potential is conveniently calculated from equations (2.125) with

$$e = \sqrt{1 - q_b^2} \quad ; \quad \psi(m) = \rho_{b0} \int_0^m dm^2 \left(\frac{m}{a_b} \right)^{-\alpha_b} e^{-m^2/r_b^2}. \quad (2.208)$$

Near-infrared photometry (BM §10.2.1) suggests values for three of the parameters, $\alpha_b = 1.8$, $q_b = 0.6$, $r_b = 1.9$ kpc, and without loss of generality,

we can set $a_b = 1$ kpc. The parameter ρ_{b0} , and hence the mass of the bulge, are determined by fitting the dynamical constraints.

(b) The dark halo By extending the spherical two-power-law models of §2.2.2f to oblate models, the density of this component is taken to have the form

$$\rho_h(R, z) = \rho_{h0} \left(\frac{m}{a_h} \right)^{-\alpha_h} \left(1 + \frac{m}{a_h} \right)^{\alpha_h - \beta_h}, \quad (2.209)$$

where m is again given by equation (2.207b) with q_b replaced by q_h . The potential of this component, in which the density varies as $r^{-\alpha_h}$ for $r \ll a_h$ and $r^{-\beta_h}$ at large r , can also be obtained from equation (2.125). Clearly, photometry provides no guidance as to the values of any of the parameters in equation (2.209); all five parameters ρ_{h0} , a_h , α_h , β_h , and q_h can only be determined by fitting the dynamical constraints. The data we use have little sensitivity to q_h , and we arbitrarily set it to 0.8.

(c) The stellar disk The density of the stellar disk is assumed to fall off exponentially with radius R , as in equation (1.7), and to depend on distance from the midplane z through the sum of two exponentials, representing the thin and thick disks described on page 13—this dependence on z is motivated by observations such as those of Gilmore & Reid (1983), shown in BM Figure 10.25. Mathematically,

$$\rho_d(R, z) = \Sigma_d e^{-R/R_d} \left(\frac{\alpha_0}{2z_0} e^{-|z|/z_0} + \frac{\alpha_1}{2z_1} e^{-|z|/z_1} \right), \quad (2.210)$$

where $\alpha_0 + \alpha_1 = 1$, Σ_d is the central surface density, R_d is the disk scale length, and $z_0 = 0.3$ kpc and $z_1 = 1$ kpc are scale heights for the thin and thick components. The potential generated by this density distribution is given by equation (2.170) with $\zeta(z)$ replaced by the expression in large brackets in equation (2.210).

(d) The interstellar medium The disk formed by a galaxy's interstellar medium (ISM) is thinner and more extended radially than the galaxy's stellar disk (see, for example, BM Figures 8.25 and 9.19). In the case of the Milky Way there is a hole of radius ~ 4 kpc at the center of the disk of the ISM (BM Figure 9.19). These observations are crudely represented by taking the density of the ISM to be

$$\rho_g(R, z) = \frac{\Sigma_g}{2z_g} \exp \left(-\frac{R}{R_g} - \frac{R_m}{R} - \frac{|z|}{z_g} \right), \quad (2.211)$$

with $R_m = 4$ kpc and $z_g = 80$ pc. The parameters Σ_g and R_g are related to the parameters Σ_d and R_d of equation (2.210) by the assumption that $R_g = 2R_d$ and that the ISM contributes 25% of the total disk surface density

Table 2.3 Parameters of Galaxy models

Parameter	Model I	Model II
R_d/kpc	2	3.2
$(\Sigma_d + \Sigma_g)/\mathcal{M}_\odot \text{pc}^{-2}$	1905	536
$\rho_{b0}/\mathcal{M}_\odot \text{pc}^{-3}$	0.427	0.3
$\rho_{h0}/\mathcal{M}_\odot \text{pc}^{-3}$	0.711	0.266
α_h	-2	1.63
β_h	2.96	2.17
a_h/kpc	3.83	1.90
$M_d/10^{10} \mathcal{M}_\odot$	5.13	4.16
$M_b/10^{10} \mathcal{M}_\odot$	0.518	0.364
$M_{h,<10 \text{ kpc}}/10^{10} \mathcal{M}_\odot$	2.81	5.23
$M_{h,<100 \text{ kpc}}/10^{10} \mathcal{M}_\odot$	60.0	55.9
$v_e(R_0)/\text{km s}^{-1}$	520	494
f_b	0.05	0.04
f_d	0.60	0.33
f_h	0.35	0.63

NOTES: In both models $0.75\Sigma(R_0)$ is contributed by stars, of which $0.05\Sigma(R_0)$ is in the thick disk. Interstellar gas accounts for the remaining $0.25\Sigma(R_0)$. The thin and thick disks have the same scale length R_d , while the gas disk has scale length $2R_d$ and a central hole of radius $R_m = 4 \text{ kpc}$. The thicknesses of the disks are $z_0 = 300 \text{ pc}$, $z_1 = 1 \text{ kpc}$, $z_g = 80 \text{ pc}$. In both models the bulge parameters are $a_b = 1 \text{ kpc}$, $\alpha_b = 1.8$, $r_b = 1.9 \text{ kpc}$, $q_b = 0.6$, while the halo axis ratio $q_h = 0.8$. The quantity $v_e(R_0)$ is the escape speed from the solar neighborhood; f_b , f_d and f_h are the fractions of the radial force supplied by bulge, disk and halo at $R_0 = 8 \text{ kpc}$. These are slightly modified forms of Models 1 and 4 of Dehnen & Binney (1998a).

at the solar radius, R_0 . The potential implied by equation (2.211) is best found from equation (2.154).

Dehnen & Binney (1998a) found that fits to the constraints described above could be obtained for a wide range of models made up of the components (a) to (d). The most important single parameter for determining the properties of a model is the scale length of the stellar disk, R_d . In §1.1.2 we estimated that R_d lies between 2 and 3 kpc. When R_d is at the lower end of this range, the disk dominates the gravitational field out to beyond the solar radius, whereas when $R_d = 3 \text{ kpc}$, the halo dominates at all radii. It is useful to examine the properties of two extreme models, namely the most and the least halo-dominated models; we designate them Models I and II and list their parameters in Table 2.3.

Model I has a small scale length, $R_d = 2 \text{ kpc}$, and gives rise to the isopotential surfaces and circular-speed curves shown in Figures 2.19 and 2.20. At small radii the halo density is $\rho_h \propto r^{-\alpha_h} = r^2$, which is the smallest value of α_h allowed by the fitting program—with this disk scale length, the best fit has the smallest possible halo contribution near the center. Figure 2.20 illus-

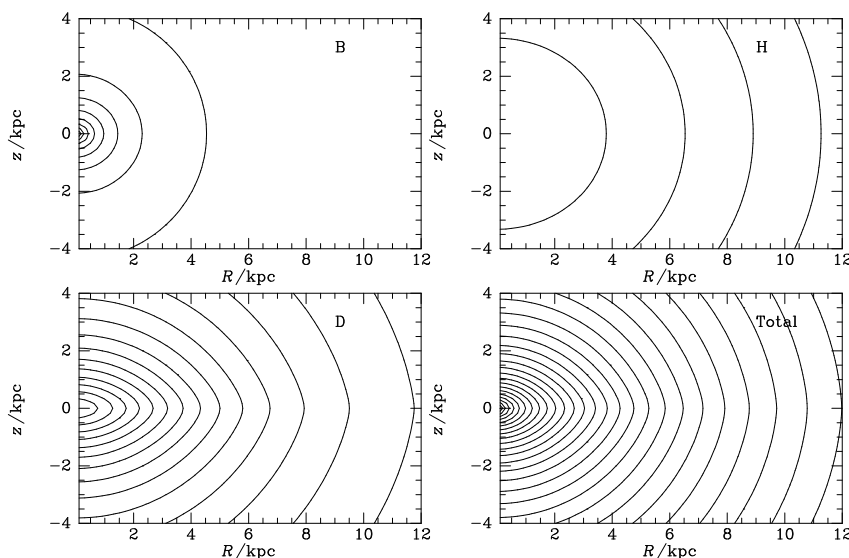


Figure 2.19 The lower right panel shows equipotential contours of a model of the Galaxy with $R_d = 2$ kpc (Model I). Contour levels are $(-0.5, -1, -1.5 \dots) \times (100 \text{ km s}^{-1})^2$. The top left panel shows the potential of the bulge, while the potentials of halo and disk are shown at top right and lower left, respectively. From top left to lower right the potentials at $(R, z) = (8 \text{ kpc}, 0)$ are $-0.28, -10.2, -2.98, -13.46 \times (100 \text{ km s}^{-1})^2$.

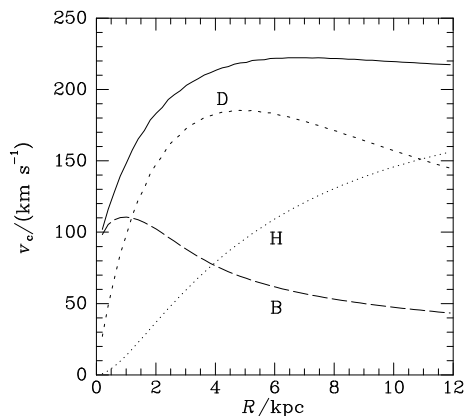


Figure 2.20 The full curve shows the circular-speed curve of Model I, whose potential is contoured in Figure 2.19. The contributions from the bulge, halo and disk are shown by the long-dashed, dotted and short-dashed curves, respectively. Notice that the total circular speed is given by the sum in quadrature of the circular speeds of the components.

trates the dynamical importance of the disk and bulge interior to the solar radius, showing that at such radii the halo makes only a small contribution to the overall circular speed—since $v_c^2 \propto g$, the contribution to the gravitational force is even smaller. This dominance is reflected in the contour plots of Figure 2.19 by the much closer packing of the equipotential contours of the bulge (top left panel) and disk (lower left panel) than those of the halo (top right panel). The equipotential surfaces of the disk are naturally more

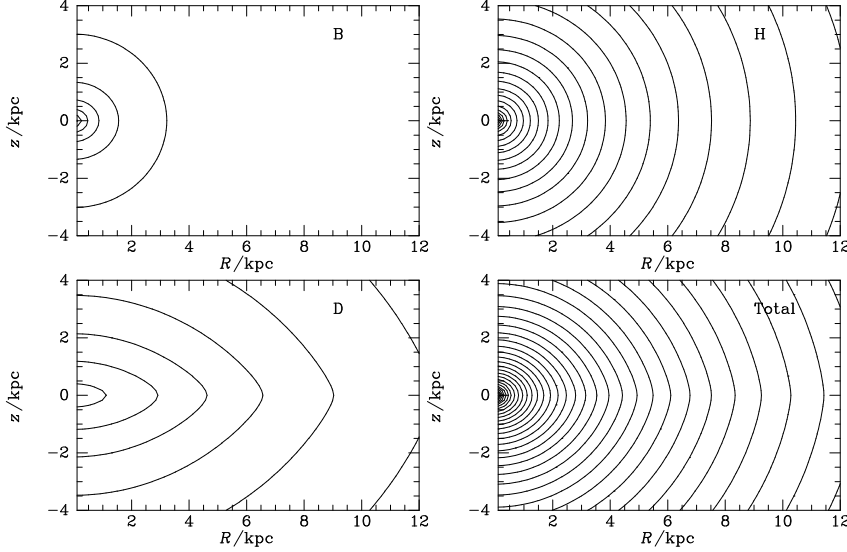


Figure 2.21 Equipotential contours of the halo-dominated Galaxy model, Model II, which has $R_d = 3.2$ kpc. The top left panel shows the potential of the bulge, while the potentials of halo and disk are shown at top right and lower left, respectively. The lower right panel shows the overall potential with contour levels $(-0.5, -1, -1.5 \dots) \times (100 \text{ km s}^{-1})^2$. From top left to lower right the potentials at $(R, z) = (8 \text{ kpc}, 0)$ are $-0.20, -9.83, -2.19, -12.21 \times (100 \text{ km s}^{-1})^2$.

highly flattened than those of either the bulge or the halo, so the equipotential surfaces of the total potential are most flattened at radii $r \sim 5$ kpc, where the disk's potential dominates.

Although in Model I the disk dominates the gravitational field (potential gradient) at R_0 , the halo makes by far the largest contribution to the total potential at all radii. For example, at the Sun's location the halo contributes $-10.2 \times (100 \text{ km s}^{-1})^2$ to the overall potential, while the disk and bulge together contribute only $-3.26 \times (100 \text{ km s}^{-1})^2$. The large contribution from the halo reflects the its enormous mass, most of it beyond R_0 . Just how much mass the halo contains is ill-determined because the Galaxy's circular speed $v_c(R)$ is uncertain beyond $\approx 2R_0$.

Figures 2.21 and 2.22 analyze the potential of Model II, a model that has a larger disk scale length, $R_d = 3.2$ kpc. As Figure 2.22 shows, in this model the halo dominates the circular speed at all radii. It does so because it is much more centrally concentrated than the halo of Model I: at small r its density rises towards the center as $r^{-1.63}$ rather than falling as in Model I. At the solar position the escape speed in this model is $v_e(R_0) = 494 \text{ km s}^{-1}$, which is observationally indistinguishable from $v_e(R_0) = 520 \text{ km s}^{-1}$ in Model I; both are consistent with the observational estimate $v_e(R_0) = (550 \pm 50) \text{ km s}^{-1}$ in Table 1.2.

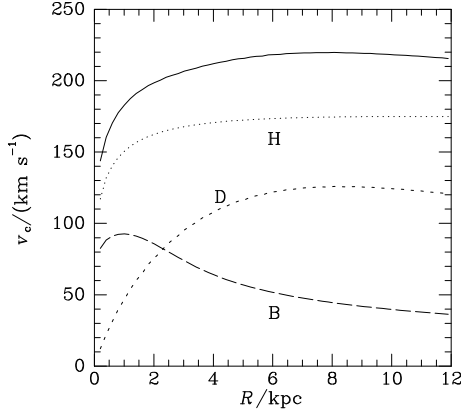


Figure 2.22 The full curve shows the circular-speed curve of Model II. The contributions from the bulge, halo and disk are shown by the long-dashed, dotted and short-dashed curves, respectively.

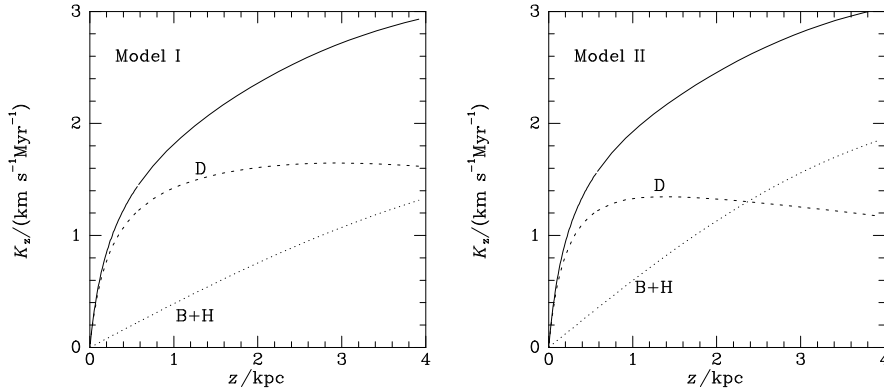


Figure 2.23 In each panel the full curve shows as a function of z in a Galactic model the force towards the galactic plane, $K_z = \partial\Phi/\partial z$ at $R_0 = 8$ kpc. The contributions from the bulge plus halo and disk are shown by the dotted and dashed curves, respectively. The left panel is for disk-dominated model, Model I (Figures 2.19 and 2.20), while the right panel is for Model II (Figures 2.21 and 2.22).

One of the striking conclusions from these models is that the relative contributions of the disk and the halo to the interior mass and the circular speed at R_0 are very uncertain. As R_d varies from 3.2 kpc to 2 kpc, the mass of the dark halo inside 10 kpc decreases by nearly a factor 2 and the fraction of the gravitational force at R_0 contributed by the halo falls⁹ from 0.63 to 0.35. Similar uncertainties are encountered in models of external disk galaxies (van Albada et al. 1985; Sellwood 1999). This degeneracy between the disk and halo parameters has to be resolved by bringing other observational constraints or dynamical arguments to bear, such as those obtained from measurements of the dynamics of galactic bars (§6.5.2e), and

⁹Models in which most of the force at $\sim 2R_d$ comes from the disk are called “maximum-disk models”—see §6.3.3.

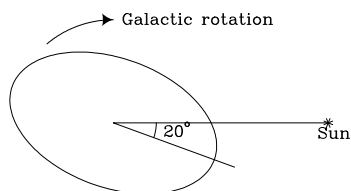


Figure 2.24 Schematic diagram of the Galactic bar.

the optical depth to gravitational microlensing towards the Galactic center (Bissantz & Gerhard 2002; Famaey & Binney 2005).

Even though the halo of Model II dominates the circular speed at R_0 , Figure 2.23 shows that the vertical force towards the disk, $K_z = \partial\Phi/\partial z$, is dominated by the disk within ~ 2 kpc of the plane. Even a relatively low-mass disk can dominate K_z in this way because a disk's contribution to K_z rises extremely quickly near the plane, where the density of disk material is high. Above one scale height, ~ 200 pc, the disk's contribution to K_z flattens off to the nearly constant value $2\pi G\Sigma(R)$ (cf. Problem 2.3). By contrast, in both panels of Figure 2.23 the halo's contribution to K_z (dotted curves) rises nearly linearly with z out to several kiloparsecs above the Sun. Notice how similar the full curves in Figure 2.23 are: despite the very different contributions to v_c from disk and halo in the two models, the shape of the observationally measurable quantity K_z (§4.9.3) is almost the same in the two cases.

(e) The bulge as a bar There is both kinematic and photometric evidence that the Milky Way's bulge is in fact a bar, that is, a highly elongated, rapidly rotating stellar system (§6.5 and BM §§9.4 and 10.1). From the vantage point of the Sun, it is hard to determine the precise shape of the bar, but, as sketched in Figure 2.24, the bar is believed to extend to a Galactocentric radius ~ 3 kpc, with its longest axis inclined by about 20° to the line from the Sun to the Galactic center (Bissantz & Gerhard 2002). The lengths of the bar's semi-axes are roughly in the ratios 1 : 0.3 : 0.3.

Both photometric studies of the bar itself and comparisons with bars in other galaxies suggest that the isodensity surfaces deviate significantly from ellipsoids (López-Corredoira et al. 2005). Nonetheless, when considering the impact that the bar has on the Galaxy's gravitational potential, it is useful to start by approximating the isodensity surfaces by ellipsoids for in this case we can obtain the potential from equation (2.140)—if a more exact result were required, one could expand the difference $\rho(\mathbf{x}) - \rho_e(\mathbf{x})$ between the actual density distribution ρ and the elliptical model ρ_e in spherical harmonics, and obtain the small correction to the potential from equation (2.95). In this spirit, we estimate the effect of the bar on the Galaxy's potential by fashioning a bar out of the axisymmetric bulge of Model I as follows.

In equations (2.207) we increase the scale radius r_b from 1.9 kpc to 3 kpc, and redefine m by $m^2 = x^2 + (y^2 + z^2)/q^2$, where x runs along the bar's long axis. We adopt $q = 0.35$ and increase ρ_0 so that the bar has roughly the same mass as the original bulge. In Figure 2.25 the full curves show the

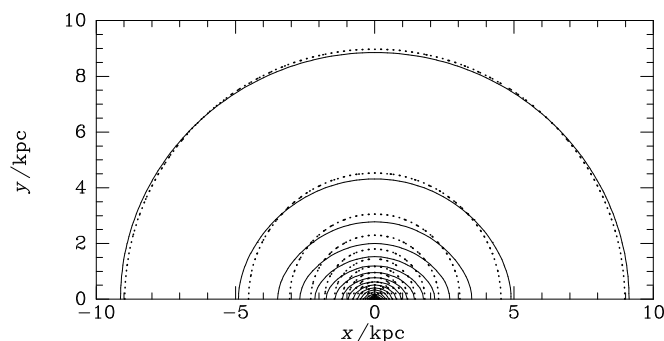


Figure 2.25 Full curves show the intersection with the Galactic plane of the isopotential surfaces of a model of the galactic bar. For comparison the dotted curves show the same curves for the axisymmetric bulge of Model I.

intersection with the Galactic plane of the bar's isopotential surfaces, while the dotted curves show the corresponding curves for the bulge of Model I. As expected, the bar's isopotential surfaces are elongated. The effect is very small near the solar circle but appreciable at $R \lesssim 5$ kpc. On account of this elongation, the potential now generates tangential forces: Along a radius that makes an angle of 45° with the bar's long axis, the ratio F_ϕ/F_r of the tangential and radial forces from the bar falls from 0.4 at the center to 0.27 at 2 kpc and 0.125 at 4 kpc. From Figure 2.20 we see that the bulge dominates F_r at $R \lesssim 1$ kpc, so in this crude model tangential forces are very important at small radii. Conversely, at $R \simeq 4$ kpc the bulge contributes only 11% of F_r , so at that radius F_ϕ is only $\sim 1\%$ of F_r . Nonetheless, the tangential forces that the bar induces can be dynamically significant for resonant orbits as far out as the solar circle because along such orbits the effect of F_ϕ can accumulate over several periods (§3.7.2 and Dehnen 2000a).

We conclude that accurate models of the triaxial bar are needed to understand the dynamics of the Milky Way at $R \lesssim 2\text{--}3$ kpc, and possibly beyond.

2.8 Potentials from functional expansions

A common theme of many of the methods we have described is the expansion of the gravitational potential and density in a set of functions that are potential-density pairs. We shall encounter such methods again in §2.9.4 as efficient tools for N-body simulation, and in §5.3.2, where we study linear response theory for stellar systems. In this section we re-examine these techniques from a general standpoint.

The basic idea of §2.3 was to approximate a real galactic density distribution by a density for which the potential is known analytically. Only a

small number of different functional forms for the density were presented, so with them the potential of a given galaxy could not be calculated to arbitrary accuracy. In §§2.6.2 to 2.6.4 we described a second approach: representing the system's density as an infinite sum of density distributions of known potential—for example the Bessel-function distributions $\Sigma_{km}(R, \phi)$ that are defined by equations (2.179). This second approach has the advantage that the potential and density can be approximated to arbitrarily high accuracy. However, a good approximation will require that a large number of terms are taken in the sum for the potential, because none of the density distributions of individual terms resembles real galaxies.

Here we show how to combine the best aspects of these two approaches. In mathematical language, we find pairs of (possibly complex) **basis functions** $\Phi_\beta(\mathbf{x})$ and $\rho_\beta(\mathbf{x})$ for $\beta = 1, 2, \dots$, that satisfy

$$\nabla^2 \Phi_\beta = 4\pi G \rho_\beta, \quad (2.212)$$

and determine coefficients a_β such that the density of the system under study can be written as the sum

$$\rho(\mathbf{x}) = \sum_{\beta} a_{\beta} \rho_{\beta}(\mathbf{x}); \quad (2.213a)$$

then the system's potential $\Phi(\mathbf{x})$ is given by

$$\Phi(\mathbf{x}) = \sum_{\beta} a_{\beta} \Phi_{\beta}(\mathbf{x}). \quad (2.213b)$$

We determine the coefficients a_β as follows. We multiply equation (2.213a) by $-\Phi_\alpha^*$ and integrate over all space to obtain

$$s_\alpha = \sum_{\beta} M_{\alpha\beta} a_\beta, \quad (2.214a)$$

where

$$s_\alpha = - \int d^3\mathbf{x} \Phi_\alpha^*(\mathbf{x}) \rho(\mathbf{x}) \quad ; \quad M_{\alpha\beta} = - \int d^3\mathbf{x} \Phi_\alpha^*(\mathbf{x}) \rho_\beta(\mathbf{x}). \quad (2.214b)$$

The elements of the matrix \mathbf{M} have a simple physical interpretation: $M_{\alpha\beta}$ is minus the potential energy of the density distribution ρ_β in the gravitational potential Φ_α^* . Using Poisson's equation (2.212) and applying the divergence theorem (B.45), we can show that \mathbf{M} is a Hermitian matrix:¹⁰

$$\begin{aligned} M_{\alpha\beta} &= -\frac{1}{4\pi G} \int d^3\mathbf{x} \Phi_\alpha^* \nabla^2 \Phi_\beta \\ &= -\frac{1}{4\pi G} \oint \Phi_\alpha^* \nabla \Phi_\beta \cdot d^2\mathbf{S} + \frac{1}{4\pi G} \int d^3\mathbf{x} \nabla \Phi_\alpha^* \cdot \nabla \Phi_\beta. \end{aligned} \quad (2.215)$$

¹⁰ In the language of linear operators, the natural inner product on the space spanned by $\{\Phi_\alpha\}$ is $(f, g) = -(4\pi G)^{-1} \int d^3\mathbf{x} f^* \nabla^2 g$. Then $(f, g) = (g, f)^*$ is minus the interaction energy of the density distributions associated with the potentials f^* and g .

The surface term vanishes when the integral is taken over all space, because Φ_α falls off at least as fast as r^{-1} . This result shows that $M_{\alpha\beta} = M_{\beta\alpha}^*$ so \mathbf{M} is Hermitian. Note that \mathbf{M} does not depend on the galactic mass distribution, so it can be computed once and for all after the basis potentials Φ_α have been chosen, and subsequently used time and again to follow the evolution of the potentials of many different galaxies, or of a single dynamically evolving galaxy.

The coefficients a_β can now be found by solving the linear equation (2.214a). We can choose the basis functions to facilitate this process. There are two strategies for exploiting this freedom.

(a) Bi-orthonormal basis functions We choose the basis functions such that \mathbf{M} is the unit matrix, so equation (2.214a) has the trivial solution $a_\alpha = s_\alpha$. From equation (2.214b) we see that this requires that the basis functions are **bi-orthonormal**:

$$-\int d^3\mathbf{x} \Phi_\alpha^*(\mathbf{x}) \rho_\beta(\mathbf{x}) = \delta_{\alpha\beta}. \quad (2.216)$$

A straightforward way of ensuring bi-orthonormality is to require that the Φ_α are eigenfunctions of the Hermitian operator ∇^2 —in this case $\rho_\alpha \propto \Phi_\alpha$, and the orthogonality of Φ_α and ρ_β is assured by the usual theorem that the eigenfunctions of a Hermitian operator are mutually orthogonal. Examples include the three functional expansions for disk potentials described in §§2.6.2 to 2.6.4. In §5.3.2 we shall use bi-orthonormal functional expansions to investigate the stability of stellar systems.

(b) Designer basis functions In principle \mathbf{a} and \mathbf{s} are vectors of infinite dimension. In practice it is necessary to truncate them and work with finite-dimensional vectors and matrices. The philosophy of **designer basis functions** is to choose the basis functions so that the galaxy can be accurately represented by the smallest possible number of them; the computational savings of having smaller vectors and matrices can more than offset the disadvantage that \mathbf{M} no longer has a simple form. The approach is most easily described by a concrete example. We examine the important special case in which each basis function is the product of a function of radius r and a spherical harmonic:

$$\Phi_{\mathbf{n}}(\mathbf{x}) = F_{\mathbf{n}}(r) Y_l^m(\boldsymbol{\Omega}), \quad (2.217)$$

where the potentials are enumerated by the vector $\mathbf{n} = (n, l, m)$, which has integer components, rather than the index α used above; here n indexes an infinite number of radial functions for each spherical harmonic. Since the spherical harmonics (eq. C.44) are orthogonal, \mathbf{M} is now block diagonal, that is, $M_{\mathbf{n}'\mathbf{n}} = 0$ for $l' \neq l$ or $m' \neq m$. Thus the equations (2.214a) can be solved separately for each spherical-harmonic index pair (l, m) .

Applying the Laplacian operator (C.49) to equation (2.217), we find with Poisson's equation (2.212) that $\rho_{\mathbf{n}}$ is given by

$$4\pi G \rho_{\mathbf{n}}(\mathbf{x}) = \nabla^2 \Phi_{\mathbf{n}} = Y_l^m(\boldsymbol{\Omega}) \left[\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{dF_{\mathbf{n}}}{dr} \right) - \frac{l(l+1)}{r^2} F_{\mathbf{n}} \right]. \quad (2.218)$$

For *any* choice of the radial function $F_{\mathbf{n}}$, equations (2.217) and (2.218) yield a potential-density pair. We now exploit this freedom to arrange for the first few functions $\Phi_{\mathbf{n}}$ in the set to provide a good approximation to the potentials of real galaxies. One way of doing this is to set $F_{000}(r)$ equal to the potential of one of the two-power-law models described in §2.2.2g. For guidance in choosing F_{0lm} for $l > 0$ we turn to the solution of Poisson's equation that we obtained in terms of spherical harmonics in §2.4. Considering the coefficient of Y_l^m in equation (2.95), we see that the terms proportional to r^l will be dominant at small r and the terms proportional to $r^{-(l+1)}$ will be dominant at large r . These observations imply that a promising choice for F_{0lm} is

$$F_{0lm}(r) = F_l(r) \equiv \frac{r^l}{(1 + r/r_0)^{2l+1}} \quad (l > 0), \quad (2.219)$$

where r_0 is a suitable scale radius. If we also apply this equation for $l = 0$, then $F_0(r)$ is proportional to the potential of a Hernquist model (eq. 2.67), consistent with our earlier argument that F_{000} should be the potential of a two-power-law galaxy model.

We must still expand our basis to include functions $F_{\mathbf{n}}(r)$ for each given (l, m) and $n > 0$ to describe accurately the radial dependence. To this end we write

$$F_{nl}(r) \equiv F_{\mathbf{n}}(r) = U_n(r)F_l(r), \quad (2.220)$$

where $U_0(r) = 1$ and the set of U_n is complete. A possible choice would be

$$U_n(r) = u^n \quad \text{where} \quad u \equiv \frac{r/r_0 - 1}{r/r_0 + 1} \quad (-1 \leq u < 1). \quad (2.221)$$

When our choice of basis functions is inserted in equation (2.214b), we find with equation (2.218) that the matrix \mathbf{M} and the vector \mathbf{s} take the forms

$$\begin{aligned} M_{\mathbf{n}'\mathbf{n}} &= -(4\pi G)^{-1} \delta_{l'l} \delta_{m'm} \int dr F_{n'l}^*(r) \left[\frac{d}{dr} \left(r^2 \frac{d}{dr} \right) - l(l+1) \right] F_{nl}(r) \\ s_{\mathbf{n}'} &= - \int d^3\mathbf{x} Y_l^{m*}(\boldsymbol{\Omega}) F_{n'l}^*(r) \rho(\mathbf{x}). \end{aligned} \quad (2.222)$$

It is possible to combine the best features of bi-orthonormal and designer basis functions. For example, if we choose $F_{nl}(r) = U_n(r)F_l(r)$, where $F_l(r)$ is given by (2.219) and $U_n(r)$ is an appropriate polynomial in the variable u defined by equation (2.221), the designer basis functions are bi-orthonormal (Hernquist & Ostriker 1992).

2.9 Poisson solvers for N-body codes

Our understanding of stellar dynamics has been profoundly advanced by **N-body codes**, computer programs that follow the motion of a large number of masses under their mutual gravitational attraction. Our discussion of these codes is in three parts: here we discuss the algorithms they use to find forces, in §3.4 we discuss the algorithms they use to move particles, and in §4.7.1 we discuss their general principles.

Realistic systems often contain many more particles than it is feasible to follow in a computer—for example, the Milky Way contains in excess of 10^{11} stars and 10^{69} GeV/ (mc^2) dark-matter particles of mass m . Currently even the largest computers cannot work efficiently with more than $\sim 10^{10}$ particles. Therefore there are two distinct types of N-body calculation with very different methodologies and problems. **Collisional N-body codes** simulate the evolution of a system with N_* stars by numerically integrating the equations of motion of exactly $N = N_*$ particles. **Collisionless N-body codes** simulate the evolution of N_* stars by following the motion of $N \ll N_*$ particles.

Collisional N-body codes are used to model systems in which relaxation is important (§1.2 and Chapter 7), in the sense that the relaxation time (1.38) is less than the duration of the numerical integration. They must accurately follow close encounters between stars, the formation and evolution of binary and triple stars, etc. These are challenging tasks because of the wide range of time and length scales involved. For example, a globular cluster may have binary stars with periods as short as hours, and yet evolve on a timescale of 10 Gyr.

Collisionless codes are easier to write but harder to understand and justify. They are used to model systems over times much shorter than the relaxation time. Since the number of particles $N \ll N_*$, the relaxation time in the model system is much smaller than that in the system being modeled; the philosophy is that the model is nevertheless accurate because the duration of the integration is much less than the relaxation time of either the real system or the model.

In essence, a collisionless code attempts to mimic the evolution of a system that contains *infinitely* many particles. Consequently, the system's density distribution is to be thought of as a continuum $\rho(\mathbf{x}, t)$, and the actual locations of the particles that model the system in the computer should be regarded as Monte-Carlo samplings of the probability-density distribution in position and velocity (§4.1).

N-body simulations—whether collisional or collisionless—use a simple principle: from the current positions of the particles, we derive the gravitational force on each particle. Then we use this force to advance the position and momentum of each particle for a short time, and find new forces. A major challenge is to produce code that will efficiently calculate the gravitational forces on a large number of bodies. We call such code a **Poisson solver**.

Poisson solvers for collisional problems can be designed to be essentially perfect: in §7.4.6d we shall see that the conceptually difficult part of a collisional simulation is integrating the orbits of particles. By contrast, Poisson solvers for collisionless simulations are fundamentally limited by particle noise: since we do not really know $\rho(\mathbf{r})$ but only the locations of a finite number of particles that sample $\rho(\mathbf{r})$, we cannot recover $\Phi(\mathbf{r})$ accurately no matter how much potential theory we know. On account of this fundamental limitation, Poisson solvers for collisionless simulations inevitably involve a compromise between inadequate resolution and excessive statistical noise. Some solvers are undoubtedly better than others, but the solver that is best suited to a given simulation depends on the scientific problem that is to be addressed. In this section we describe the most important types of Poisson solver, and discuss their strengths and weaknesses.

2.9.1 Direct summation

Consider evaluating the force on particle α by simply summing the contributions from all the other particles in the simulation,

$$\mathbf{F}_\alpha = \sum_{\beta \neq \alpha} Gm_\beta \frac{\mathbf{r}_\beta - \mathbf{r}_\alpha}{|\mathbf{r}_\beta - \mathbf{r}_\alpha|^3}. \quad (2.223)$$

Each such force evaluation involves the calculation of $N-1$ distances $|\mathbf{r}_\beta - \mathbf{r}_\alpha|$. Each distance can be used twice, once for the contribution of particle β to the force on particle α and once for the force from particle α on particle β . So, if forces are evaluated by direct summation a minimum of $\frac{1}{2}N(N-1)$ distances have to be evaluated, where $N \gg 1$ is the number of particles in the simulation. Thus, the work per timestep increases with N as N^2 . Below we shall see that there are vastly more efficient ways of evaluating forces, which scale as $N \ln N$.

The first N -body calculations were restricted to small values of N , for which the difference between N^2 and $N \ln N$ is not large, so direct summation was a viable option. For the values of $N > 10\,000$ now current, $N^2/(N \ln N) > 10^3$, so direct summation is in principle not an attractive strategy. Notwithstanding this fact it is still in use because its simple formulae can be encoded in hardware, so special-purpose processors can calculate forces between tens of thousands of particles at acceptable speeds (Makino et al. 2003).

Softening If two particles, say α and β , approach each other closely, the term in equation (2.223) that describes the force $\mathbf{F}_{\alpha\beta}$ between them becomes large. This phenomenon is problematic for both collisional and collisionless calculations, for different reasons.

In the collisional case the divergence of $\mathbf{F}_{\alpha\beta}$ as $\mathbf{r}_\alpha \rightarrow \mathbf{r}_\beta$ is a real physical effect, but one that is computationally awkward because it implies that the

equations of motion of particles α and β have to be integrated with very small timesteps. Unless sophisticated workarounds such as “regularization” (§3.4.7) are employed to handle this situation, close encounters can bring the integration virtually to a halt. “Softening” of the force-law as described below is an expedient that keeps the integration moving along at an acceptable pace at the price of a loss of realism.

In the collisionless case the divergence of $\mathbf{F}_{\alpha\beta}$ predicted by equation (2.223) is entirely unphysical, for the mass distribution we are trying to model is inherently smooth. The divergence is an artifact of the Monte-Carlo sampling of the density distribution. In this case softening can enhance rather than detract from the realism of the simulation.

Softening involves replacing equation (2.223) with

$$\mathbf{F}_{\alpha} = \sum_{\beta \neq \alpha} Gm_{\beta} S_{\mathbf{F}}(|\mathbf{r}_{\beta} - \mathbf{r}_{\alpha}|) \frac{\mathbf{r}_{\beta} - \mathbf{r}_{\alpha}}{|\mathbf{r}_{\beta} - \mathbf{r}_{\alpha}|}. \quad (2.224)$$

This equation differs from (2.223) only in the replacement of r^{-2} by $S_{\mathbf{F}}(r)$, where $\mathbf{r} = \mathbf{r}_{\beta} - \mathbf{r}_{\alpha}$. $S_{\mathbf{F}}(r)$, the **force softening kernel**, is some function that tends to r^{-2} for values of its argument bigger than the **softening length** ϵ , and tends smoothly to zero for small values of its argument. Equation (2.224) has the desirable features that: (i) the force exerted by particle β on particle α , $Gm_{\alpha}m_{\beta}S_{\mathbf{F}}(r)\mathbf{r}/r$, is equal and opposite to the force exerted by α on β , so Newton’s third law is satisfied; (ii) the force between any two particles acts along the line that joins them; (iii) the force approaches the usual gravitational force at large separations; (iv) the force between two particles at the same location is zero.

$S_{\mathbf{F}}(r)$ is the derivative of another function $S(r)$, the **softening kernel**, which appears in the equation for the potential at the location of particle α :

$$\Phi_{\alpha} \equiv \Phi(\mathbf{r}_{\alpha}) = \sum_{\beta \neq \alpha} Gm_{\beta} S(|\mathbf{r}_{\beta} - \mathbf{r}_{\alpha}|). \quad (2.225)$$

A widely used form of S is

$$S(\mathbf{r}) = -\frac{1}{\sqrt{r^2 + \epsilon^2}}. \quad (2.226)$$

In this case, the gravitational potential of each particle is just that of a Plummer sphere of scale length ϵ (eq. 2.44a). In a collisionless simulation this makes physical sense: the mass that is represented by a particle is in reality distributed through a non-zero volume.¹¹ The density law of a Plummer sphere (eq. 2.44b) places most of the sphere’s mass inside $r = \epsilon$, so the natural choice is to make ϵ of order the inter-particle separation.

¹¹ Note that in this model the softened force between two particles is the force between a point mass and a Plummer sphere, not the force between two Plummer spheres.

In collisionless simulations the choice of softening kernel and softening length is a compromise between maximizing the smoothness of the force-field and minimizing the degradation of the spatial resolution caused by the softening. Despite its popularity, the choice (2.226) of S is not optimal in this sense, because the density of a Plummer sphere falls off too slowly with radius. A better choice is

$$S(\mathbf{r}) = -\frac{r^2 + \frac{3}{2}\epsilon^2}{(r^2 + \epsilon^2)^{3/2}}, \quad (2.227)$$

which amounts to replacing the potential of each particle with that of a sphere of radius ϵ in which the density is proportional to $(r^2 + \epsilon^2)^{-7/2}$. See Dehnen (2001) for a discussion of the merits of various softening kernels.

An ideal Poisson solver would allow ϵ to vary as a particle moves from a high-density to a low-density region. In practice most contemporary Poisson solvers use fixed softening due to technical difficulties associated with energy conservation.

We now consider algorithms that enable one to evaluate the sum (2.225) to high accuracy in a number of arithmetic operations that scales as $N \ln N$, rather as N^2 .

2.9.2 Tree codes

Our discussion is based on Dehnen (2000b). The particles are first organized into a structure that was introduced by Barnes & Hut (1986): we place an imaginary cube around the simulation and divide it into eight equal sub-cubes. If any sub-cube contains more than one particle, we divide it in turn into eight equal cubes, and so on until every cube contains at most one particle. This hierarchy of cubes forms an **oct-tree**, and the original cube is called the **root** of the tree. Each cube after the root has a parent cube, and seven sibling cubes. Any cube that contains more than one particle has eight child cubes. There is a clear analogy with a real tree, whose trunk divides into great boughs, which divide into branches, which divide into twigs, which ultimately carry leaves. Particles are the leaves of our oct-tree.

Suppose now that we locate the center of mass of the particles in each cube and evaluate the sums

$$M_0 \equiv \sum_{\alpha} m_{\alpha} \quad ; \quad M_{ij} = \sum_{\alpha} m_{\alpha} x_i^{\alpha} x_j^{\alpha} \quad ; \quad M_{ijk} = \sum_{\alpha} m_{\alpha} x_i^{\alpha} x_j^{\alpha} x_k^{\alpha} \quad (2.228)$$

and so forth, where the sums run over the particles in the cube and \mathbf{x}^{α} is the position vector of particle α relative to the cube's center of mass. This hierarchy of sums is closely related to the multipole expansion of §2.4, and these quantities are called the **Cartesian multipole moments**. The three in equation (2.228) are the monopole, quadrupole, and octopole moments.

Box 2.2: Scaling of tree codes

Consider the work involved in determining the potential at the location of one of N simulation particles. For simplicity we assume that the system is not very inhomogeneous, although the argument we give can be generalized. The sum over cubes for the potential includes some cubes (“leaves”) that contain only one particle, the rest being “branches” that contain more than one particle. All branch cubes have opening angle $\leq \theta_o$, while most of the leaf cubes have opening angle $> \theta_o$. If the simulation is representing a system of linear size L , the leaf cubes have a characteristic linear size $l = L/N^{1/3}$ and lie within distance $\sim D = l/\theta_o$ of the given particle; hence there are $\sim (D/l)^3 = \theta_o^{-3}$ leaf cubes.

Suppose we now increase the number of simulation particles by a factor 8. Then branch cubes will continue to be branch cubes and will not require further subdivision because the angles they subtend are $< \theta_o$. Hence the extra work that comes with increased N arises mainly from the θ_o^{-3} leaf cubes that previously had only one particle; most will now need to be subdivided. Thus every time we increase N by a factor 8 we have to do a fixed amount of additional work. It now follows that the work required to determine the potential at the location of a given particle is proportional to $\ln N$, and the work involved in determining the forces on all particles scales as $N \ln N$.

The object $M_i = \sum_{\alpha} m_{\alpha} x_i^{\alpha}$ that might have been listed in equation (2.228) is the cube’s dipole moment, and it vanishes identically by virtue of the fact that \mathbf{x}^{α} is the position of α relative to the center of mass. If the cube is one of the tree’s leaves, $M_0 = m_{\alpha}$ and all the other multipoles vanish. There are of order $N \log_2 N$ cubes in the tree, so the labor involved in setting it up scales as $N \ln N$.

Once the tree has been constructed, we can evaluate the sum (2.225) that yields the potential at any location \mathbf{x} to good accuracy with only $\sim \ln N$ operations as follows. We first consider evaluating the force as the sum of the forces due to the monopole, quadrupole, octopole, \dots , 2^k -pole etc., of the root cube, up to some predetermined order $k = K$. This sum would converge rapidly and could be truncated after a few terms without unacceptable inaccuracy if \mathbf{x} were remote from the root cube in the sense that the root cube subtended a small **opening angle** θ_o , say $\theta_o \lesssim 1^\circ$, at \mathbf{x} . In general, \mathbf{x} will lie within the root cube, so the desired condition will be not be satisfied. However, we may then instead consider evaluating the potential by summing the multipole contributions from each of the root’s eight child cubes. If any of these subtends an angle at \mathbf{x} that exceeds θ_o , we consider the sum of the multipoles of its children, and so on recursively until the whole potential is obtained as a sum of the multipoles from cubes that

subtend angles at \mathbf{x} smaller than θ_o . Some of these final cubes will be leaves, which are deemed to subtend zero angle, because they contain only a single particle. The key point is that as the number of particles in the simulation is increased, the number of cubes over which we sum to get the force at \mathbf{x} increases much more slowly—see Box 2.2.

Cartesian multipole expansion To derive explicit formulae for the potential in terms of the multipoles of cubes, we now develop some of the theory of Cartesian multipoles. We consider the contribution to the gravitational potential Φ at position \mathbf{X} in cube A that is generated by the matter in cube B. We have

$$\Phi(\mathbf{X}) = G \int_{\text{cube B}} d^3\mathbf{Y} S(\mathbf{X} - \mathbf{Y}) \rho(\mathbf{Y}), \quad (2.229)$$

where ρ is the matter density and S is the softening kernel (eq. 2.225). We specialize to the case of well-separated cubes and write $\mathbf{X} = \hat{\mathbf{X}} + \mathbf{x}$ and $\mathbf{Y} = \hat{\mathbf{Y}} + \mathbf{y}$, where $\hat{\mathbf{X}}$ and $\hat{\mathbf{Y}}$ are the centers of mass of cubes A and B, respectively. Then we Taylor expand $S(\mathbf{X} - \mathbf{Y})$ in powers of $\mathbf{x} - \mathbf{y}$ around $\hat{\mathbf{X}} - \hat{\mathbf{Y}}$:

$$\begin{aligned} S(\mathbf{X} - \mathbf{Y}) &= S(\hat{\mathbf{X}} - \hat{\mathbf{Y}}) + \sum_{i=1}^3 (\mathbf{x} - \mathbf{y})_i \left. \frac{\partial S(\mathbf{R})}{\partial R_i} \right|_{\mathbf{R}=\hat{\mathbf{X}}-\hat{\mathbf{Y}}} \\ &\quad + \frac{1}{2!} \sum_{i,j=1}^3 (\mathbf{x} - \mathbf{y})_i (\mathbf{x} - \mathbf{y})_j \left. \frac{\partial^2 S(\mathbf{R})}{\partial R_i \partial R_j} \right|_{\mathbf{R}=\hat{\mathbf{X}}-\hat{\mathbf{Y}}} + \dots \quad (2.230) \\ &= \sum_n \frac{1}{n!} (\mathbf{x} - \mathbf{y})^{(n)} \cdot [\partial_{(n)} S(\mathbf{R})]_{\mathbf{R}=\hat{\mathbf{X}}-\hat{\mathbf{Y}}}. \end{aligned}$$

Here the last line employs a notation in which $\mathbf{r}^{(n)} \equiv r_{i_1} r_{i_2} \dots r_{i_n}$ is a product of n components of the vector \mathbf{r} and $\partial_{(n)}$ is the corresponding product of partial derivatives with respect to these components, while the centered dot implies summation over all the n indices i_k (a total of 3^n terms). Substituting (2.230) into (2.229), we obtain

$$\Phi(\hat{\mathbf{X}} + \mathbf{x}) = G \sum_n \frac{1}{n!} \left[\int d^3\mathbf{y} \rho(\hat{\mathbf{Y}} + \mathbf{y}) (\mathbf{x} - \mathbf{y})^{(n)} \right] \cdot [\partial_{(n)} S(\mathbf{R})]_{\mathbf{R}=\hat{\mathbf{X}}-\hat{\mathbf{Y}}}. \quad (2.231)$$

Using the expansion

$$(\mathbf{x} - \mathbf{y})^{(n)} = \sum_{k=0}^n \frac{n!}{k!(n-k)!} \mathbf{x}^{(k)} \mathbf{y}^{(n-k)}, \quad (2.232)$$

we have

$$\Phi(\hat{\mathbf{X}} + \mathbf{x}) = G \sum_n \sum_{k=0}^n \frac{1}{k!(n-k)!} \mathbf{x}^{(k)} \mathbf{M}^{(n-k)} \cdot [\partial_{(n)} S(\mathbf{R})]_{\mathbf{R}=\hat{\mathbf{X}}-\hat{\mathbf{Y}}}, \quad (2.233)$$

Box 2.3: Derivatives of a spherically symmetric softening kernel

Consider the case of a spherically symmetric softening kernel $S(\mathbf{R})$, i.e., one of the form $S(|\mathbf{R}|)$. In this case the evaluation of the first few tensors $\partial_{(n)}S$ is straightforward:

$$\begin{aligned}\partial_{(1)}S(R) &= S' \frac{R_i}{R} \\ \partial_{(2)}S(R) &= \left(S'' - \frac{S'}{R}\right) \frac{R_i R_j}{R^2} + \frac{S'}{R} \delta_{ij} \\ \partial_{(3)}S(R) &= \left(S''' - 3\frac{S''}{R} + 3\frac{S'}{R^2}\right) \frac{R_i R_j R_k}{R^3} \\ &\quad + \left(S'' - \frac{S'}{R}\right) \frac{\delta_{ij} R_k + \delta_{jk} R_i + \delta_{ki} R_j}{R^2} \\ \partial_{(4)}S(R) &= \left(S^{iv} - 6\frac{S'''}{R} + 9\frac{S''}{R^2} - 9\frac{S'}{R^3}\right) \frac{R_i R_j R_k R_l}{R^4} \\ &\quad + \left(S''' - 3\frac{S''}{R} + 3\frac{S'}{R^2}\right) (\delta_{ij} R_k R_l + \delta_{jk} R_i R_l + \delta_{ki} R_j R_l \\ &\quad + \delta_{il} R_j R_k + \delta_{lj} R_i R_k + \delta_{lk} R_i R_j) / R^3 \\ &\quad + \left(S'' - \frac{S'}{R}\right) \frac{\delta_{ij} \delta_{kl} + \delta_{ik} \delta_{jl} + \delta_{il} \delta_{jk}}{R^2}\end{aligned}$$

where

$$\mathbf{M}^{(k)} \equiv \int_{\text{cube B}} d^3\mathbf{y} \rho(\hat{\mathbf{Y}} + \mathbf{y}) \mathbf{y}^{(k)}. \quad (2.234)$$

When the density distribution is simulated with particles, $\rho(\mathbf{r}) = \sum_{\alpha} m_{\alpha} \delta(\mathbf{r} - \mathbf{r}^{\alpha})$ is a sum of delta functions centered on the locations of the particles. In this case the integral in (2.234) can be immediately evaluated and yields the objects defined by equations (2.228).

Equation (2.233) expresses the potential generated by cube B at an arbitrary point in cube A as a power series in the components of \mathbf{x} , the separation between the point of observation and A's center of mass. The coefficients in this power series are made up of (a) the multipole moments of B, which were evaluated during tree construction, and (b) the derivatives $\partial_{(n)}S$ of the softening kernel evaluated on the vector \mathbf{R} that joins the centers of mass of the two cubes, which can be evaluated once and for all at the start of the integration. Once these quantities are all to hand, we can quickly evaluate the force on every particle in cube A by summing the series with the appropriate value of \mathbf{x} . The derivatives $\partial_{(n)}S$ are explicitly evaluated for a spherically symmetric kernel in Box 2.3.

An important simplification arises because the same derivatives occur in the analogous series for the potential that cube A generates at a point within cube B. Dehnen (2000b) explains how one can exploit this fact, both to accelerate the computation and to achieve exact momentum conservation.

Unlike some of the codes that we are about to present, a tree code employs no grid. This feature enables tree codes to handle problems such as galaxy mergers, in which dense stellar systems move through a large volume of nearly empty space—most grid-based codes cannot handle such problems efficiently because they would waste most of their computing resources covering low-density regions with the same high-resolution grid that they require to handle the dense stellar systems. Another merit of tree codes is that they are readily parallelized by, for example, delegating to a different processor construction of each of the trees that emerge from the root's eight child cubes.

2.9.3 Particle-mesh codes

A wide variety of N -body codes solve Poisson's equation using estimates of the density at a set of regularly spaced points, the nodes of a mesh or grid. The simplest mesh is a Cartesian grid and we concentrate on this case, although most of the principles we describe carry over to other meshes, such as ones based on spherical polar coordinates. More detail can be found in Hockney & Eastwood (1988).

The process of estimating the density on a grid from the positions of a large number of particles that trace the density is called **mass assignment** because it involves assigning the mass of each particle to one or more nearby nodes. The simplest algorithm is to assign all of the mass of each particle to the nearest node. This procedure, which is known as the **nearest grid point (NGP)** scheme, is rarely used for two reasons. First, the NGP scheme only samples the density distribution crudely. Second, as explained on page 134 below, for technical reasons associated with momentum conservation the NGP scheme leads to solutions of Poisson's equation in which forces change discontinuously midway between nodes.

Better mass-assignment schemes spread the mass of each particle over several nearby nodes. Box 2.4 describes a hierarchy of widely used schemes that spread the mass over $1, 2^D, 3^D, \dots$ nodes of a D -dimensional grid. If we decide on one node, we have the NGP scheme. If we want to split a particle's mass over 2^D nodes, the resulting scheme is called the **cloud in cell (CIC)** mass-assignment scheme because it may be visualized by expanding every mass point into a homogeneous cube with side length equal to the grid spacing, and assigning every mass element in this cube to the nearest node (Figure 2.26). If we split a particle's mass over 3^D nodes, we have the **triangular shaped cloud (TSC)** mass-assignment scheme, so called because the scheme may be derived by regarding the particle to be a cloud

Box 2.4: Mass assignment schemes

A mass-assignment scheme is a function $W_{\mathbf{p}}(\mathbf{r})$ that gives the fraction of the mass of a particle at \mathbf{r} that is assigned to node \mathbf{p} . Here \mathbf{p} is a triple of integers and the node's location is $\mathbf{x}_{\mathbf{p}} = (p_1, p_2, p_3)\Delta$, with Δ the grid spacing. The simplest schemes for a three-dimensional Cartesian grid are the product of three functions, each evaluated on one component of \mathbf{r}

$$W_{\mathbf{p}}(\mathbf{r}) = w(x - p_1\Delta)w(y - p_2\Delta)w(z - p_3\Delta).$$

The condition that all the mass is assigned somewhere is $\sum_{\mathbf{p}} W_{\mathbf{p}}(\mathbf{r}) = 1$, which is assured if $M_0 \equiv \sum_i w(x - i\Delta) = 1$. Additional constraints arise from demanding that the grid-scale fluctuations are as small as possible at large distances. Suppose that we are modeling a one-dimensional potential with softening kernel $S(x)$. Then the potential due to a mass m at x is $\Phi(x') = mS(x' - x)$. The potential computed from a grid is

$$\Phi_{\mathbf{g}}(x') = m \sum_i w(x - i\Delta)S(x' - i\Delta).$$

If $|x' - x| \gg \Delta$, we can Taylor expand $S(x' - i\Delta)$ in powers of $x - i\Delta$:

$$\Phi_{\mathbf{g}}(x') = m \sum_i \sum_{n=0}^{\infty} \frac{1}{n!} w(x - i\Delta) S^{(n)}(x' - x) (x - i\Delta)^n,$$

where $S^{(n)}$ is the n th derivative of S . If condition $M_0 = 1$ is satisfied, the $n = 0$ term equals the exact potential, so grid-scale fluctuations are minimized if

$$M_n \equiv \sum_i (x - i\Delta)^n w(x - i\Delta) = \text{constant} \quad \text{for } n = 1, 2, \dots$$

If the mass is distributed between two nodes,

$$w(x) = \begin{cases} 1 - |x|/\Delta & \text{for } |x| < \Delta, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

satisfies $M_0 = 1$, $M_1 = 0$. This is the cloud in cell (CIC) mass-assignment scheme. If the mass is distributed between three nodes,

$$w(x) = \begin{cases} \frac{3}{4} - |x|^2/\Delta^2 & \text{for } |x| < \frac{1}{2}\Delta \\ \frac{1}{2}(\frac{3}{2} - |x|/\Delta)^2 & \text{for } \frac{1}{2}\Delta < |x| < \frac{3}{2}\Delta \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

satisfies $M_0 = 1$, $M_1 = 0$, $M_2 = \frac{1}{4}\Delta^2$. This is the triangular shaped cloud (TSC) mass-assignment scheme. There is an alternative three-node scheme in which $M_2 = 0$, but in this scheme $w(x)$ is discontinuous at $|x|/\Delta = \frac{1}{2}, \frac{3}{2}$, whereas in the TSC scheme $w(x)$ is continuous and differentiable.

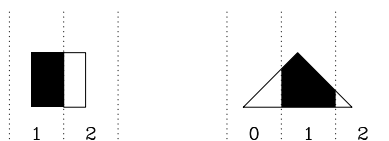


Figure 2.26 The CIC (left) and TSC (right) mass-assignment schemes in one dimension. The dotted lines mark the boundaries between cells of the grid, and the rectangle (CIC) and triangle (TSC) represent the mass of the particle. Mass is assigned to cells in proportion to the fraction of the area of the rectangle or triangle that lies in each cell—for example, the shaded portions are assigned to cell 1.

in which the mass density falls linearly with coordinate differences from the location of the particle (Figure 2.26).

In the NGP scheme, the mass assigned to a node changes discontinuously as a particle moves through the grid. In the CIC scheme, the mass changes continuously but with discontinuous first derivative. In the TSC scheme, both the mass and its derivative change continuously. In view of the role a mass-assignment scheme plays in force interpolation—see page 134 below—the progressive increase in smoothness as one proceeds up this hierarchy of mass-assignment schemes is valuable. This advantage must be set against two disadvantages of proceeding up the hierarchy: (i) in three dimensions the computational cost increases rapidly as the number of nodes to which mass is assigned increases from 1 for NGP, to 8 for CIC, 27 for TSC, etc.; (ii) spreading mass over many nodes involves a loss of spatial resolution. In practice the most widely used schemes are CIC and TSC.

Given the density on the mesh, we are ready to solve Poisson’s equation for the potential. To do this we must specify boundary conditions, which in practice are one of two types: **periodic** and **vacuum**. Periodic boundary conditions are used when the grid is imagined to be one cell of an infinite lattice; these are natural for cosmological simulations. With vacuum boundary conditions we require that at large distances $\Phi \rightarrow -GM/r$, where M is the total mass on the grid. These are appropriate for simulations of isolated stellar systems.

(a) Periodic boundary conditions We solve Poisson’s equation by representing ∇^2 as a finite-difference operator. For example, we can make the approximation

$$\begin{aligned}
 (\nabla^2\Phi)_{lmn} \simeq & (\Phi_{l+1,m,n} + \Phi_{l-1,m,n} + \Phi_{l,m+1,n} + \Phi_{l,m-1,n} \\
 & + \Phi_{l,m,n+1} + \Phi_{l,m,n-1} - 6\Phi_{l,m,n})/\Delta^2,
 \end{aligned}
 \tag{2.235}$$

where Δ is the mesh spacing. This substitution converts the Poisson equation into a large system of linear equations for the Φ_{lmn} in terms of the known ρ_{lmn} . The system obtained by equating (2.235) to $4\pi G\rho_{lmn}$ is most

efficiently solved using discrete Fourier transforms (DFTs; see Appendix G and Problem 2.23).

(b) Vacuum boundary conditions Consider now how to solve Poisson's equation with vacuum boundary conditions. We use the integral form of Poisson's equation (2.3), or its smoothed form (2.229), in which the softening kernel S already contains our choice of boundary conditions. We write

$$\Phi(\mathbf{r}) = G \sum_{\mathbf{p}} S(\mathbf{r} - \mathbf{r}_{\mathbf{p}}) m_{\mathbf{p}}, \quad (2.236)$$

where $m_{\mathbf{p}}$ is the mass assigned to node \mathbf{p} . The labor required to evaluate Φ at the location of every particle is proportional to the number of particles times the number of nodes. A much more efficient approach is to evaluate Φ only at the nodes, and then to interpolate to the locations of particles. At the nodes, equation (2.236) reads

$$\Phi_{\mathbf{q}} = G \sum_{\mathbf{p}} S_{\mathbf{q}-\mathbf{p}} m_{\mathbf{p}}, \quad (2.237)$$

where \mathbf{p} and \mathbf{q} are three-vectors with integer components in the range $(0, K - 1)$, with K the number of nodes along a side of the computational box. In equation (2.237), each component of S 's subscript ranges from $-(K - 1)$ (when $q_1 = 0$ and $p_1 = K - 1$) to $K - 1$ (when $q_1 = K - 1$ and $p_1 = 0$). Thus S 's subscript has components that take on $2K - 1$ values, while the subscript of m has components that range over only K possible values. Discrete Fourier transforms provide an extremely efficient way of solving equations of this form, but before we can apply them, we have to make the problem symmetrical in the subscripts, and arrange for them to take 2^n values, for some integer n . We do this by taking $K = 2^{n-1}$ and padding m out to a matrix of size $(2K)^3$, whose subscripts cover the range $(-K, K - 1)$ with $m_{\mathbf{p}}$ set to zero when any component of the subscript lies outside the physically significant range, $(0, K - 1)$. Finally we eliminate negative indices by making S and m periodic functions of each index with period $2K$. Thus extended, equation (2.237) takes the form of a discrete convolution. So we take the DFT of both sides of this equation, and, with the aid of the discrete Fourier convolution theorem (eq. G.5), we obtain

$$\hat{\Phi}_{\mathbf{k}} = G \hat{S}_{\mathbf{k}} \hat{m}_{\mathbf{k}}. \quad (2.238)$$

According to this equation, the DFT of Φ can be obtained from that of m simply by multiplying by the DFT of S . With $\hat{\Phi}_{\mathbf{k}}$ in hand, Φ is easily obtained by doing an inverse DFT. Moreover, the DFT of S can be calculated once and for all at the start of the simulation and stored for use at each timestep.

Box 2.5: James’s Fourier Poisson solver

We describe the method of James (1977) for imposing vacuum boundary conditions on the system (2.235) without doubling the range of all the indices. This method is both fast and economical with computer memory.

We start from the density values on a cubical grid and imagine the values to represent electric charge density rather than mass density. Next, we find the potential that would be generated by the charges attached to the interior of the grid if the boundary of the grid were grounded—that is, if the potential Φ vanished on the boundary. This potential is the solution of Poisson’s equation for the given charge density at interior points of the grid, subject to the boundary condition $\Phi = 0$ on the grid boundary. We obtain the required solution by writing Φ in terms of its sine transform (eq. G.7a)

$$\Phi_{lmn} = \sum_{\alpha\beta\gamma=1}^{K-1} \Phi^{\alpha\beta\gamma}(\text{SSS}) \sin\left(\frac{\pi l\alpha}{K}\right) \sin\left(\frac{\pi m\beta}{K}\right) \sin\left(\frac{\pi n\gamma}{K}\right), \quad (1)$$

where $0 \leq l, m, n \leq K$. The sine transform automatically ensures that Φ_{lmn} vanishes on the boundary. Substituting this expression for Φ_{lmn} into the numerical approximation (2.235) for the Laplacian, and then equating this to $4\pi G\rho_{lmn}$ with ρ_{lmn} expressed in terms of its sine transform, we discover that each amplitude $\Phi^{\alpha\beta\gamma}(\text{SSS})$ is simply a multiple of the corresponding amplitude $\rho^{\alpha\beta\gamma}(\text{SSS})$. Since the sine transform is its own inverse (eq. G.7b), it is a simple matter to recover the Φ_{lmn} from the $\Phi^{\alpha\beta\gamma}(\text{SSS})$.

The potential we obtain in this way differs from the one we require because the walls bear charges that differ from those specified in the original problem. (When the walls of the grid are grounded, charges flow along the grounding cables until Φ vanishes on the boundary.) We use Gauss’s theorem to determine the actual surface density of wall charges Σ —since the potential vanishes outside the box, $\Sigma = -g_{\perp}/(4\pi G)$, where $g_{\perp} = |\mathbf{n} \cdot \nabla\Phi|$ is the magnitude of the force field just inside the walls of the grounded box. Finally, we convolve the resulting charge distribution on the walls with a softening kernel to obtain the potential that it generates when vacuum boundary conditions are applied, and subtract this potential from the potential generated by the mass in the grounded box, to get the potential that the originally specified density distribution generates with vacuum boundary conditions. The softening kernel should be the inverse of the discrete Laplacian used in the calculation of $\Phi^{\alpha\beta\gamma}(\text{SSS})$ —see James (1977) for details.

The key to the success of this procedure is that the convolution of the wall charge density with the softening kernel is simple because the charge density is hollow—see Appendix G or Magorrian (2007) for details.

This approach to imposing vacuum boundary conditions is very wasteful of computer memory. Box 2.5 describes a clever and much more economical alternative.

Once the potential has been obtained at every node, we use numerical differentiation to obtain the gravitational field $\mathbf{g}_\mathbf{n}$ at the nodes: for example, the component parallel to the first axis might be

$$\hat{\mathbf{e}}_1 \cdot \mathbf{g}_\mathbf{n} = -\frac{1}{2\Delta}(\Phi_{(n_1+1, n_2, n_3)} - \Phi_{(n_1-1, n_2, n_3)}). \quad (2.239)$$

More generally

$$\mathbf{g}_\mathbf{j} = \sum_{\mathbf{n}} \mathbf{A}_\mathbf{n} \Phi_{\mathbf{j}+\mathbf{n}}, \quad (2.240)$$

where $\mathbf{A}_\mathbf{n} = -\mathbf{A}_{-\mathbf{n}}$ defines some numerical differentiation scheme. Values of the forces on individual particles are then obtained by interpolation from these values.

In any such scheme, it is desirable that the sum of the forces on all the particles is zero to machine precision, because this is a prerequisite for conservation of the simulation's total momentum; if momentum is not conserved the system is liable to rocket off the grid. Consider therefore the sum of all forces. We let $W_\mathbf{p}(\mathbf{r})$ denote the fraction of the mass of a particle that lies at \mathbf{r} that is assigned to node \mathbf{p} , so

$$m_\mathbf{p} = \sum_{\text{particles } \beta} m_\beta W_\mathbf{p}(\mathbf{r}_\beta). \quad (2.241)$$

We introduce similar functions $Q_\mathbf{j}(\mathbf{r})$ to describe the interpolation scheme used to calculate forces on particles given the field on the mesh: let the gravitational field at \mathbf{r} be

$$\mathbf{g}(\mathbf{r}) = \sum_{\text{nodes } \mathbf{j}} Q_\mathbf{j}(\mathbf{r}) \mathbf{g}_\mathbf{j}. \quad (2.242)$$

Then with equation (2.237), (2.240) and (2.241), the sum of all particle forces is

$$\begin{aligned} \mathbf{F}_{\text{tot}} &= \sum_{\text{particles } \alpha} m_\alpha \mathbf{g}(\mathbf{r}_\alpha) = \sum_{\alpha} m_\alpha \sum_{\text{nodes } \mathbf{j}} Q_\mathbf{j}(\mathbf{r}_\alpha) \sum_{\mathbf{n}} \mathbf{A}_\mathbf{n} \Phi_{\mathbf{j}+\mathbf{n}} \\ &= G \sum_{\mathbf{j}\mathbf{p}} \sum_{\alpha\beta} m_\alpha m_\beta Q_\mathbf{j}(\mathbf{r}_\alpha) W_\mathbf{p}(\mathbf{r}_\beta) \sum_{\mathbf{n}} \mathbf{A}_\mathbf{n} S_{\mathbf{j}+\mathbf{n}-\mathbf{p}}. \end{aligned} \quad (2.243)$$

The sum over \mathbf{n} , which is the numerical derivative of the softening kernel, is antisymmetric in \mathbf{j} and \mathbf{p} , as can be seen by replacing the dummy variable \mathbf{n} by $-\mathbf{n}$ and exploiting the antisymmetry of $\mathbf{A}_\mathbf{n}$ and the symmetry of $S_\mathbf{q}$. In view of this antisymmetry, \mathbf{F}_{tot} vanishes if $Q_\mathbf{p} = W_\mathbf{p}$, and in fact only this

choice of $Q_{\mathbf{p}}$ guarantees symmetry of the sum over $\alpha\beta$ for all possible particle locations. Therefore, we conclude that net momentum can be conserved if and only if we use the same scheme to interpolate forces from the grid as we use to assign mass to the grid. In particular, if we use the NGP scheme, forces change discontinuously midway between nodes.

(c) Mesh refinement Stellar systems, whether galaxies, galactic nuclei, star clusters or clusters of galaxies, tend to develop small, tightly bound condensations. These objects are often of considerable astronomical interest. For example, the visible stars of the Milky Way are a dense condensation in a much larger dark-matter halo. The ideal Poisson solver would determine the gravitational field of a structure no matter how small the structure becomes.

In the CIC and TSC mass-assignment schemes, the force between two particles that are separated by no more than the mesh spacing is significantly smaller than the value predicted by Newton's inverse-square law, regardless of the value of ϵ that is used in the softening kernel. Hence the dynamics of structures smaller than a few mesh spacings cannot be faithfully followed.

Stellar systems usually have large volumes in which the density is low, and small volumes in which it is high. Hence achieving higher spatial resolution by subdividing the grid (**mesh refinement**) throughout the volume occupied by a system is undesirable: in low-density regions a fine mesh is both computationally expensive and pointless. If higher resolution is to be attained by mesh refinement, we must refine locally and adaptively—at each timestep, high-density regions must be located and equipped with a finer mesh. Handling an adaptive mesh requires intricate code, and if the mesh has a complex geometry, it will not be possible to use DFTs to evaluate the sum (2.237) efficiently. For a discussion of adaptive-mesh codes see Knebe, Green, & Binney (2001).

(d) P³M codes The **Particle-particle-particle-mesh** or **P³M** technique provides a simpler way of enhancing the spatial resolution of a PM code. The idea behind a P³M code is to use the standard PM algorithm to calculate the contribution from distant particles to the force \mathbf{F} on a particle, while using the direct sum (2.224) to get the contribution to \mathbf{F} from particles that lie in the same or adjacent cells. The softening kernel $S_{\mathbf{F}}$ used for the direct sum must be carefully chosen to provide only the difference between the Newtonian force and the force already obtained from the grid.

The efficiency of a P³M code depends sensitively on the number of particles, N_{\max} , in the most densely populated cell, since the cost of evaluating the direct sum scales as N_{\max}^2 rather than as $N_{\max} \ln(N_{\max})$. P³M codes have been extensively used for cosmological simulations (Efstathiou et al. 1985). In these simulations, galaxy formation leads gradually to higher and higher densities, so N_{\max} eventually becomes very large and the simulations grind to a halt. This problem can be resolved by using either a tree algorithm (Bode & Ostriker 2003) or a separate P³M implementation (Couchman 1991) to handle sums within populous cells.

2.9.4 Spherical-harmonic codes

We saw in §2.4 and Figure 2.10 that good approximations to the potentials of stellar systems that are not too far from spherical can be obtained from the first few terms of the spherical-harmonic expansion (2.95).

Two rather different approaches to the numerical implementation of equation (2.95) have been widely used. In the first (McGlynn 1984; Bontekoe 1988), we cover the computational volume with a spherical mesh centered on the estimated location of the center of the system. We identify the particles that lie within each spherical shell around this center, say between radii $a_i - \Delta/2$ and $a_i + \Delta/2$, and then evaluate the integrals (2.94) up to $l = l_{\max}$, using the formula

$$m_{lm}(a_i) \simeq \sum_{a_i - \Delta/2 < r_\alpha \leq a_i + \Delta/2} m_\alpha Y_l^{m*}(\theta_\alpha, \phi_\alpha). \quad (2.244)$$

Once the $m_{lm}(a_i)$ have been evaluated for each shell i , we use an interpolation algorithm to construct a continuous function $m_{lm}(a)$. Then we can numerically estimate at r the value of the big brackets in equation (2.95):

$$Q_{lm}(r) \equiv \frac{1}{r^{l+1}} \int_0^r da a^{(l+2)} m_{lm}(a) + r^l \int_r^\infty \frac{da}{a^{l-1}} m_{lm}(a). \quad (2.245)$$

The potential at the location of any particle can be calculated from the sum

$$\Phi(\mathbf{r}_\alpha) = -4\pi G \sum_{l=0}^{l_{\max}} \sum_{m=-l}^l \frac{Y_l^m(\theta_\alpha, \phi_\alpha)}{2l+1} Q_{lm}(r_\alpha). \quad (2.246)$$

The second approach uses the functional-expansion technique introduced in §2.8 (Saha 1993). When setting up the simulation we choose a set of functions $F_{nl}(r)$ in which to do the expansion, and for the first few values of l evaluate the matrix \mathbf{M} that is defined by the first of equations (2.222). Finally we store the inverses of these matrices. Then at each timestep we approximate the quantities that are defined by the second of equations (2.222) as sums over particles

$$s_{\mathbf{n}} \simeq \sum_{\alpha} m_\alpha Y_l^{m*}(\theta_\alpha, \phi_\alpha) F_{nl}(r_\alpha), \quad (2.247)$$

where $\mathbf{n} \equiv (n, l, m)$. Now we solve the linear equations (2.214) for the $a_{\mathbf{n}}$ and have that the potential is

$$\Phi(\mathbf{r}) = \sum_{\mathbf{n}} a_{\mathbf{n}} Y_l^m(\boldsymbol{\Omega}) F_{nl}(r). \quad (2.248)$$

When properly constructed—in particular when the code includes all allowed values of m for each l —spherical-harmonic codes conserve angular momentum (relative to the chosen center) to machine accuracy, but they do not conserve linear momentum.

2.9.5 Simulations of planar systems

Studies of the dynamics of thin disks often confine the particles to a plane. In this case the softening length ϵ can be interpreted as the characteristic thickness of the disk: if two particles pass at a fixed vertical separation Δz from one another, their mutual interaction potential is given by equation (2.225) with \mathbf{r}_α and \mathbf{r}_β two-dimensional vectors, and ϵ set equal to Δz in equation (2.226).

A two-dimensional mesh that is K cells on a side consumes less computer memory and CPU time than the corresponding three-dimensional mesh by a factor K , so PM codes are much more attractive in two dimensions than in three. Both rectangular (Hohl & Hockney 1969) and polar (Sellwood 1983) meshes have been widely used. On going from three dimensions to two, the functional expansion technique yields a similar efficiency gain because the number of subscripts n, l etc., that need be summed over decreases from three to two (Earn & Sellwood 1995).

Problems

2.1 [1] Show that the gravitational potential energy of a spherical system of finite mass in which the density satisfies $\lim_{r \rightarrow 0} \rho r^{5/2} = 0$ can be written

$$W = -\frac{G}{2} \int_0^\infty dr \frac{M^2(r)}{r^2}, \quad (2.249)$$

where $M(r)$ is the mass interior to radius r .

2.2 [1] Prove that the Chandrasekhar potential-energy tensor for any spherical body has the form $W_{jk} = \frac{1}{3}W\delta_{jk}$, where W is the potential energy. Hint: start from equation (2.19).

2.3 [1] Show that the potential of an infinite razor-thin sheet of surface density Σ in the plane $z = 0$ is $\Phi = 2\pi G\Sigma|z| + \text{constant}$, (a) using Gauss's theorem, and (b) from Poisson's equation.

2.4 [1] (Suggested by A. Toomre) Show that $\Phi = \ln[r(1 + |\cos(\theta)|)]$ solves Laplace's equation everywhere except when $r = 0$ or $\theta = \pi/2$. By applying Gauss's theorem near $\theta = \pi/2$, find the potential of the Mestel disk (2.158) in the limit $R_{\max} \rightarrow \infty$.

2.5 [2] The **finite Mestel disk** is a razor-thin disk with surface density $\Sigma(R)$ such that (i) $\Sigma(R) = 0$ for all $R > R_0$; (ii) the circular speed is $v_c(R) = v_0 = \text{constant}$ for all $R < R_0$. The surface density of the finite Mestel disk was first derived by Mestel (1963) but here we describe a short, elegant derivation due to Brada & Milgrom (1995).

(a) Consider a spherical mass distribution with density $\rho(r) = \frac{1}{2}Ar^{-2}$ for $r < R_0$ and zero for $r > R_0$, where A is a constant. Argue that the circular speed is independent of radius and independent of R_0 so long as $r < R_0$.

(b) Now squash the sphere along one direction, so that its isodensity surfaces are spheroids with axis ratio q . Argue that the circular speed in the equatorial plane of the squashed system is still independent of radius and independent of R_0 so long as the equatorial radius $R < R_0$. Hint: use Newton's third theorem.

(c) By considering the limit $q \rightarrow 0$, show that a disk with surface density

$$\Sigma(R) = \begin{cases} (A/R) \cos^{-1}(R/R_0) & (R < R_0) \\ 0 & (R > R_0) \end{cases} \quad (2.250)$$

has a flat circular-speed curve for $R < R_0$.

(d) Show that $v_0^2 = \pi^2 GA$. Hint: let $R_0 \rightarrow \infty$ and compare to the infinite Mestel disk.

2.6 [1] Defining **prolate spheroidal coordinates** (u, v) by $R = a \sinh u \sin v$, $z = a \cosh u \cos v$, where $a > 0$ is a constant, show that $R^2 + (a + |z|)^2 = a^2(\cosh u + |\cos v|)^2$. Hence show that the potential (2.68a) of the Kuzmin disk can be written

$$\Phi_K(u, v) = -\frac{GM}{a} \frac{\cosh u - |\cos v|}{\sinh^2 u + \sin^2 v}. \quad (2.251)$$

In §3.5.3 we show that this potential is an example of a Stäckel potential, in which orbits admit an extra isolating integral.

2.7 [2] Astronauts orbiting an unexplored planet find that (i) the surface of the planet is precisely spherical and centered on $r = 0$; and (ii) the potential exterior to the planetary surface is $\Phi = -GM/r$ exactly, that is, there are no non-zero multipole moments other than the monopole. Can they conclude from these observations that the mass distribution in the interior of the planet is spherically symmetric? If not, give a simple example of a non-spherical mass distribution that would reproduce the observations.

2.8 [1] (Suggested by L. Ciotti) If a transparent, spherical stellar system has constant mass-to-light ratio Υ , prove that the potential at radius r is

$$\Phi(r) = -\frac{2G\Upsilon}{r} \int_0^r dx S(x), \quad (2.252)$$

where $S(x)$ is the strip brightness defined in Problem 1.3.

2.9 [1] If a transparent, spherical stellar system has constant mass-to-light ratio Υ , prove that the central potential is (Ciotti 1991)

$$\Phi(0) = -4G\Upsilon \int_0^\infty dR I(R), \quad (2.253)$$

where $I(R)$ is the surface brightness at projected radius R .

2.10 [2] Consider an axisymmetric body whose density distribution is $\rho(R, z)$ and total mass is $M = \int d^3\mathbf{r} \rho(R, z)$. Assume that the body has finite extent, $\rho(R, z) = 0$ for $r^2 = R^2 + z^2 > r_{\max}^2$, and is symmetric about its equator, that is, $\rho(R, -z) = \rho(R, z)$.

(a) Show that at distances large compared to r_{\max} , the potential arising from this body can be written in the form

$$\Phi(R, z) \simeq -\frac{GM}{r} - \frac{G}{4} \frac{(R^2 - 2z^2)}{r^5} \int d^3\mathbf{r}' \rho(R', z')(R'^2 - 2z'^2), \quad (2.254)$$

where the error is of order $(r_{\max}/r)^2$ smaller than the second term.

(b) Show that at large distances from an exponential disk with surface density $\Sigma(R) = \Sigma_0 \exp(-R/R_d)$, the potential has the form

$$\Phi(R, z) \simeq -\frac{GM}{r} \left[1 + \frac{3R_d^2(R^2 - 2z^2)}{2r^4} + O(R_d^4/r^4) \right], \quad (2.255)$$

where M is the mass of the disk.

2.11 [2] Show that the potential energy of an exponential disk is $W \simeq -11.627 G\Sigma_0^2 R_d^3$. Show further that if all stars move on circular orbits, the disk's angular momentum is $J \simeq 17.462 G^{1/2} \Sigma_0^{3/2} R_d^{7/2}$ and its kinetic energy is $K \simeq 5.813 G\Sigma_0^2 R_d^3$. Hence show that for this disk the dimensionless spin parameter $\lambda \equiv J|E|^{1/2}/GM^{5/2} \simeq 0.4255$, where $E = K + W$ is the total energy.

2.12 [1] The r^{-1} dependence of the gravitational potential on distance arises because the graviton, which carries the gravitational field, is massless. If the graviton had a mass m_g , the gravitational potential due to a body of mass M would be $\Phi(r) = -GM e^{-\alpha r}/r$, where $\alpha = m_g c/\hbar$ (the **Yukawa potential**), which reduces to the Newtonian potential in the limit $\alpha \rightarrow 0$. What is the analog of Poisson's equation (2.10) for the Yukawa potential?

2.13 [2] Prove that the external potentials and gravitational fields of any two confocal spheroids of uniform density and equal mass are everywhere the same.

2.14 [2] Use equation (2.140) to show that a prolate body with density $\rho = \rho_0(1 + R^2/a_1^2 + z^2/a_3^2)^{-2}$, where $a_3 > a_1$, generates the potential

$$\Phi(u, v) = -\pi G \rho_0 a_1^2 a_3 \int_0^\infty d\tau \frac{\sqrt{a_3^2 + \tau}}{(\tau + a_3^2 + \lambda)(\tau + a_3^2 + \mu)}, \quad (2.256)$$

where (u, v) are oblate spheroidal coordinates defined by equation (2.96) with $\Delta^2 = a_3^2 - a_1^2$, and we have written $\lambda \equiv \Delta^2 \sinh^2 u$, $\mu \equiv -\Delta^2 \cos^2 v$. Decompose the integral in (2.256) into partial fractions to show (without evaluating the integrals) that Φ is a Stäckel potential

$$\Phi(\lambda, \mu) = \frac{H(\lambda) - H(\mu)}{\lambda - \mu}, \quad (2.257)$$

where H is a continuous function (de Zeeuw 1985 and §3.5.3). Finally, show that

$$\Phi(u, v) = -\frac{2\pi G a_1^2 a_3 \rho_0}{\Delta^2} \frac{f(\Delta \sinh u) - f(i\Delta \cos v)}{\sinh^2 u + \cos^2 v}, \quad (2.258a)$$

where

$$f(z) \equiv z \tan^{-1}(z/a_3). \quad (2.258b)$$

Hint: to ensure convergence of the integrals, you may wish to add $(\tau + a_3^2)^{-\frac{1}{2}}$ to one of the integrands and subtract it from the other. The body with this potential is called the **perfect prolate spheroid**, because it is the only prolate axisymmetric density distribution of constant ellipticity that has a Stäckel potential.

2.15 [1] Show that the central potential of a thin axisymmetric disk is

$$\Phi(0, 0) = -2\pi G \int_0^\infty dR \Sigma(R). \quad (2.259)$$

Hint: use equation (C.68).

2.16 [1] Prove that the potential $\Phi(r)$ is a non-decreasing function of r in any spherical system. Does the same conclusion hold in an axisymmetric razor-thin disk? If so, prove it; if not, find a counter-example.

2.17 [2] (Suggested by M. Merrifield) An axisymmetric disk is seen edge-on and has projected mass per unit length $\mu(X)$. Show that its surface density is

$$\Sigma(R) = -\frac{1}{\pi} \int_R^\infty \frac{dX}{\sqrt{X^2 - R^2}} \frac{d\mu}{dX}, \quad (2.260)$$

and that its potential is

$$\Phi(R, z) = 2G \int_0^\infty dX \frac{d\mu}{dX} \sin^{-1} \left(\frac{2X}{\sqrt{z^2 + (R+X)^2} + \sqrt{z^2 + (R-X)^2}} \right). \quad (2.261)$$

Hint: use (2.153a).

2.18 [1] Use equation (2.190) to show that the razor-thin disk for which the circular speed is given by

$$v_c^2 = \frac{v_0^2}{\sqrt{1 + R^2/a^2}} \quad (2.262a)$$

has surface density (Toomre 1963)

$$\Sigma(R) = \frac{v_0^2}{2\pi G R} \left(1 - \frac{1}{\sqrt{1 + a^2/R^2}} \right). \quad (2.262b)$$

Show that these formulae correspond to the Mestel disk in an appropriate limit. Show that the surface density (2.68b) of the Kuzmin disk is obtained when $\Sigma(R)$ is differentiated with respect to a^2 , and hence recover the circular speed of this model.

2.19 [3] We have derived relations between the potential and surface density of non-axisymmetric disks by solving Laplace's equation in cylindrical coordinates (§2.6.2) and oblate spheroidal coordinates (§2.6.4). Derive a relation of this kind by solving Laplace's equation in spherical coordinates, and show that the result is identical with the formula derived using logarithmic spirals (eq. 2.199). Hint: you may need associated Legendre functions $P_\lambda^m(x)$, where λ is a real number. See also equations (C.12) and (C.31).

2.20 [2] Show that the circular speed $v_c(R)$ in a thin axisymmetric disk of surface density $\Sigma(R)$ may be written in the form (Mestel 1963)

$$v_c^2(R) = \frac{GM(R)}{R} + 2G \sum_{k=1}^{\infty} \alpha_k \left[\frac{2k+1}{R^{2k+1}} \int_0^R dR' \Sigma(R') R'^{2k+1} - 2kR^{2k} \int_R^{\infty} dR' \frac{\Sigma(R')}{R'^{2k}} \right], \quad (2.263)$$

where

$$\alpha_k \equiv \int_0^\pi d\theta P_{2k}(\cos \theta) = \pi \left[\frac{(2k)!}{2^{2k}(k!)^2} \right]^2. \quad (2.264)$$

Hint: start with equation (2.3) and expand $|\mathbf{x} - \mathbf{x}'|^{-1}$ in Legendre polynomials using equation (C.35).

2.21 [2] Show that the potential of an axisymmetric disk with surface density $\Sigma(R)$ is

$$\Phi(R, z) = -\frac{2G}{\sqrt{R}} \int_0^\infty dR' \Sigma(R') k K(k) \sqrt{R'}, \quad (2.265)$$

where $K(k)$ is a complete elliptic integral (Appendix C.4) and

$$k^2 \equiv \frac{4RR'}{(R+R')^2 + z^2}. \quad (2.266)$$

Hint: start with equation (2.3) and show that $|\mathbf{x} - \mathbf{x}'|^2 = 4RR'k^{-2}\{1 - k^2 \cos^2[\frac{1}{2}(\phi - \phi')]\}$. Note that the integral (2.265) has a logarithmic singularity when $z = 0$ and $R' \rightarrow R$, which requires some care when the integral is being evaluated numerically.

2.22 [3] (Suggested by H. Dejonghe) Prove that the surface density $\Sigma(x, y)$ and potential $\Phi(x, y)$ in a disk occupying the plane $z = 0$ are related by

$$\Sigma(x', y') = \frac{1}{4\pi^2 G} \iint \frac{dx dy}{|\mathbf{x} - \mathbf{x}'|} \left(\frac{\partial^2 \Phi}{\partial x^2} + \frac{\partial^2 \Phi}{\partial y^2} \right). \quad (2.267)$$

2.23 [1] Consider the discrete form of Poisson's equation with periodic boundary conditions that is obtained by using the approximation (2.235) for the value of the Laplacian on a grid with K mesh cells on a side. Write Φ_{lmn} and ρ_{lmn} in terms of their DFTs,

$$\Phi_{\mathbf{r}} = \sum_{\mathbf{k}} \hat{\Phi}_{\mathbf{k}} e^{2\pi i \mathbf{k} \cdot \mathbf{r} / K} \quad ; \quad \rho_{\mathbf{r}} = \sum_{\mathbf{k}} \hat{\rho}_{\mathbf{k}} e^{2\pi i \mathbf{k} \cdot \mathbf{r} / K}, \quad (2.268)$$

where \mathbf{r} and \mathbf{k} are vectors with integer components in the range $(0, K - 1)$. Show that

$$\Phi_{\mathbf{k}} = \frac{2\pi G \rho_{\mathbf{k}} \Delta^2}{\cos(2\pi k_x / K) + \cos(2\pi k_y / K) + \cos(2\pi k_z / K) - 3}. \quad (2.269)$$

2.24 [2] In some numerical simulations of spherical stellar systems, spherical symmetry is enforced by treating each of the N stars as a "superstar," that is, a spherical shell containing a large number of stars with randomly oriented orbits but the same radius, radial velocity, and scalar angular momentum. Let m_n and r_n be the mass and radius of the n^{th} superstar. Assume that the superstars are sorted in order of increasing radius, that is, $r_1 < r_2 < \dots < r_N$.

(a) Show that the force on superstar n is

$$\mathbf{F}_n = -\hat{\mathbf{e}}_r \frac{GM_n}{r_n^2} \quad \text{where} \quad M_n \equiv \sum_{j=1}^{n-1} m_j. \quad (2.270)$$

Thus show that the force on every star can be computed in $O(N)$ steps once the stars are sorted in radius, in contrast to direct summation (§2.9.1), which requires $O(N^2)$ steps.

(b) A set of N superstars can be sorted into increasing order in $O(N \ln N)$ steps (e.g., Press et al. 1986). Thus, the initial calculation of the forces for an unsorted system requires $O(N \ln N)$ steps. Show that force calculations at subsequent timesteps require only $O(N)$ steps.

3

The Orbits of Stars

In this chapter we examine the orbits of individual stars in gravitational fields such as those found in stellar systems. Thus we ask the questions, “What kinds of orbits are possible in a spherically symmetric, or an axially symmetric potential? How are these orbits modified if we distort the potential into a bar-like form?” We shall obtain analytic results for the simpler potentials, and use these results to develop an intuitive understanding of how stars move in more general potentials.

In §§3.1 to 3.3 we examine orbits of growing complexity in force fields of decreasing symmetry. The less symmetrical a potential is the less likely it is that we can obtain analytic results, so in §3.4 we review techniques for integrating orbits in both a given gravitational field, and the gravitational field of a system of orbiting masses. Even numerically integrated orbits in gravitational fields of low symmetry often display a high degree of regularity in their phase-space structures. In §3.5 we study this structure using analytic models, and develop analytic tools of considerable power, including the idea of adiabatic invariance, which we apply to some astronomical problems in §3.6. In §3.7 we develop Hamiltonian perturbation theory, and use it to study the phenomenon of orbital resonance and the role it plays in generating orbital chaos. In §3.8 we draw on techniques developed throughout the chapter to understand how elliptical galaxies are affected by the existence of central stellar cusps and massive black holes at their centers.

All of the work in this chapter is based on a fundamental approximation:

although galaxies are composed of stars, we shall neglect the forces from individual stars and consider only the large-scale forces from the overall mass distribution, which is made up of thousands of millions of stars. In other words, we assume that the gravitational fields of galaxies are *smooth*, neglecting small-scale irregularities due to individual stars or larger objects like globular clusters or molecular clouds. As we saw in §1.2, the gravitational fields of galaxies *are* sufficiently smooth that these irregularities can affect the orbits of stars only after many crossing times.

Since we are dealing only with gravitational forces, the trajectory of a star in a given field does not depend on its mass. Hence, we examine the dynamics of a particle of unit mass, and quantities such as momentum, angular momentum, and energy, and functions such as the Lagrangian and Hamiltonian, are normally written per unit mass.

3.1 Orbits in static spherical potentials

We first consider orbits in a static, spherically symmetric gravitational field. Such fields are appropriate for globular clusters, which are usually nearly spherical, but, more important, the results we obtain provide an indispensable guide to the behavior of orbits in more general fields.

The motion of a star in a centrally directed gravitational field is greatly simplified by the familiar law of conservation of angular momentum (see Appendix D.1). Thus if

$$\mathbf{r} = r\hat{\mathbf{e}}_r \quad (3.1)$$

denotes the position vector of the star with respect to the center, and the radial acceleration is

$$\mathbf{g} = g(r)\hat{\mathbf{e}}_r, \quad (3.2)$$

the equation of motion of the star is

$$\frac{d^2\mathbf{r}}{dt^2} = g(r)\hat{\mathbf{e}}_r. \quad (3.3)$$

If we remember that the cross product of any vector with itself is zero, we have

$$\frac{d}{dt} \left(\mathbf{r} \times \frac{d\mathbf{r}}{dt} \right) = \frac{d\mathbf{r}}{dt} \times \frac{d\mathbf{r}}{dt} + \mathbf{r} \times \frac{d^2\mathbf{r}}{dt^2} = g(r)\mathbf{r} \times \hat{\mathbf{e}}_r = 0. \quad (3.4)$$

Equation (3.4) says that $\mathbf{r} \times \dot{\mathbf{r}}$ is some constant vector, say \mathbf{L} :

$$\mathbf{r} \times \frac{d\mathbf{r}}{dt} = \mathbf{L}. \quad (3.5)$$

Of course, \mathbf{L} is simply the angular momentum per unit mass, a vector perpendicular to the plane defined by the star's instantaneous position and

velocity vectors. Since this vector is constant, we conclude that the star moves in a plane, the **orbital plane**. This finding greatly simplifies the determination of the star's orbit, for now that we have established that the star moves in a plane, we may simply use plane polar coordinates (r, ψ) in which the center of attraction is at $r = 0$ and ψ is the azimuthal angle in the orbital plane. In terms of these coordinates, the Lagrangian per unit mass (Appendix D.3) is

$$\mathcal{L} = \frac{1}{2}[\dot{r}^2 + (r\dot{\psi})^2] - \Phi(r), \quad (3.6)$$

where Φ is the gravitational potential and $g(r) = -d\Phi/dr$. The equations of motion are

$$0 = \frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{r}} - \frac{\partial \mathcal{L}}{\partial r} = \ddot{r} - r\dot{\psi}^2 + \frac{d\Phi}{dr}, \quad (3.7a)$$

$$0 = \frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{\psi}} - \frac{\partial \mathcal{L}}{\partial \psi} = \frac{d}{dt}(r^2\dot{\psi}). \quad (3.7b)$$

The second of these equations implies that

$$r^2\dot{\psi} = \text{constant} \equiv L. \quad (3.8)$$

It is not hard to show that L is actually the length of the vector $\mathbf{r} \times \dot{\mathbf{r}}$, and hence that (3.8) is just a restatement of the conservation of angular momentum. Geometrically, L is equal to twice the rate at which the radius vector sweeps out area.

To proceed further we use equation (3.8) to replace time t by angle ψ as the independent variable in equation (3.7a). Since (3.8) implies

$$\frac{d}{dt} = \frac{L}{r^2} \frac{d}{d\psi}, \quad (3.9)$$

equation (3.7a) becomes

$$\frac{L^2}{r^2} \frac{d}{d\psi} \left(\frac{1}{r^2} \frac{dr}{d\psi} \right) - \frac{L^2}{r^3} = -\frac{d\Phi}{dr}. \quad (3.10)$$

This equation can be simplified by the substitution

$$u \equiv \frac{1}{r}, \quad (3.11a)$$

which puts (3.10) into the form

$$\frac{d^2u}{d\psi^2} + u = \frac{1}{L^2u^2} \frac{d\Phi}{dr} (1/u). \quad (3.11b)$$

The solutions of this equation are of two types: along **unbound** orbits $r \rightarrow \infty$ and hence $u \rightarrow 0$, while on **bound** orbits r and u oscillate between finite limits. Thus each bound orbit is associated with a periodic solution of this equation. We give several analytic examples later in this section, but in general the solutions of equation (3.11b) must be obtained numerically.

Some additional insight is gained by deriving a “radial energy” equation from equation (3.11b) in much the same way as we derive the conservation of kinetic plus potential energy in Appendix D; we multiply (3.11b) by $du/d\psi$ and integrate over ψ to obtain

$$\left(\frac{du}{d\psi}\right)^2 + \frac{2\Phi}{L^2} + u^2 = \text{constant} \equiv \frac{2E}{L^2}, \quad (3.12)$$

where we have used the relation $d\Phi/dr = -u^2(d\Phi/du)$.

This result can also be derived using Hamiltonians (Appendix D.4). From (3.6) we have that the momenta are $p_r = \partial\mathcal{L}/\partial\dot{r} = \dot{r}$ and $p_\psi = \partial\mathcal{L}/\partial\dot{\psi} = r^2\dot{\psi}$, so with equation (D.50) we find that the Hamiltonian per unit mass is

$$\begin{aligned} H(r, p_r, p_\psi) &= p_r\dot{r} + p_\psi\dot{\psi} - \mathcal{L} \\ &= \frac{1}{2}\left(p_r^2 + \frac{p_\psi^2}{r^2}\right) + \Phi(r) \\ &= \frac{1}{2}\left(\frac{dr}{dt}\right)^2 + \frac{1}{2}\left(r\frac{d\psi}{dt}\right)^2 + \Phi(r). \end{aligned} \quad (3.13)$$

When we multiply (3.12) by $L^2/2$ and exploit (3.9), we find that the constant E in equation (3.12) is simply the numerical value of the Hamiltonian, which we refer to as the energy of that orbit.

For bound orbits the equation $du/d\psi = 0$ or, from equation (3.12)

$$u^2 + \frac{2[\Phi(1/u) - E]}{L^2} = 0 \quad (3.14)$$

will normally have two roots u_1 and u_2 between which the star oscillates radially as it revolves in ψ (see Problem 3.7). Thus the orbit is confined between an inner radius $r_1 = u_1^{-1}$, known as the **pericenter** distance, and an outer radius $r_2 = u_2^{-1}$, called the **apocenter** distance. The pericenter and apocenter are equal for a circular orbit. When the apocenter is nearly equal to the pericenter, we say that the orbit has small **eccentricity**, while if the apocenter is much larger than the pericenter, the eccentricity is said to be near unity. The term “eccentricity” also has a mathematical definition, but only for Kepler orbits—see equation (3.25a).

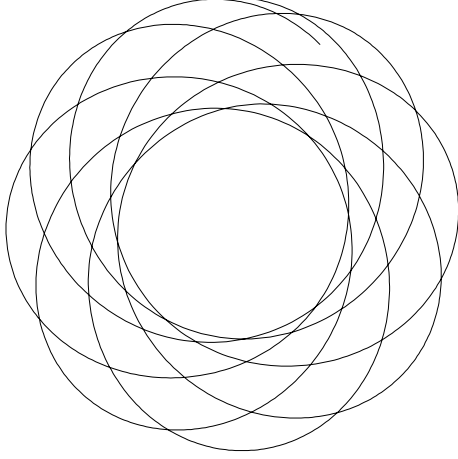


Figure 3.1 A typical orbit in a spherical potential (the isochrone, eq. 2.47) forms a rosette.

The **radial period** T_r is the time required for the star to travel from apocenter to pericenter and back. To determine T_r we use equation (3.8) to eliminate $\dot{\psi}$ from equation (3.13). We find

$$\left(\frac{dr}{dt}\right)^2 = 2(E - \Phi) - \frac{L^2}{r^2}, \quad (3.15)$$

which may be rewritten

$$\frac{dr}{dt} = \pm \sqrt{2[E - \Phi(r)] - \frac{L^2}{r^2}}. \quad (3.16)$$

The two possible signs arise because the star moves alternately in and out. Comparing (3.16) with (3.14) we see that $\dot{r} = 0$ at the pericenter and apocenter distances r_1 and r_2 , as of course it must. From equation (3.16) it follows that the radial period is

$$T_r = 2 \int_{r_1}^{r_2} \frac{dr}{\sqrt{2[E - \Phi(r)] - L^2/r^2}}. \quad (3.17)$$

In traveling from pericenter to apocenter and back, the azimuthal angle ψ increases by an amount

$$\Delta\psi = 2 \int_{r_1}^{r_2} \frac{d\psi}{dr} dr = 2 \int_{r_1}^{r_2} \frac{L}{r^2} \frac{dt}{dr} dr. \quad (3.18a)$$

Substituting for dt/dr from (3.16) this becomes

$$\Delta\psi = 2L \int_{r_1}^{r_2} \frac{dr}{r^2 \sqrt{2[E - \Phi(r)] - L^2/r^2}}. \quad (3.18b)$$

The **azimuthal period** is

$$T_\psi = \frac{2\pi}{|\Delta\psi|} T_r; \quad (3.19)$$

in other words, the mean angular speed of the particle is $2\pi/T_\psi$. In general $\Delta\psi/2\pi$ will not be a rational number. Hence the orbit will not be closed: a typical orbit resembles a rosette and eventually passes close to every point in the annulus between the circles of radii r_1 and r_2 (see Figure 3.1 and Problem 3.13). There are, however, two and only two potentials in which all bound orbits are closed.

(a) Spherical harmonic oscillator We call a potential of the form

$$\Phi(r) = \frac{1}{2}\Omega^2 r^2 + \text{constant} \quad (3.20)$$

a spherical harmonic oscillator potential. As we saw in §2.2.2b, this potential is generated by a homogeneous sphere of matter. Equation (3.11b) could be solved analytically in this case, but it is simpler to use Cartesian coordinates (x, y) defined by $x = r \cos \psi$, $y = r \sin \psi$. In these coordinates, the equations of motion are simply

$$\ddot{x} = -\Omega^2 x \quad ; \quad \ddot{y} = -\Omega^2 y, \quad (3.21a)$$

with solutions

$$x = X \cos(\Omega t + \epsilon_x) \quad ; \quad y = Y \cos(\Omega t + \epsilon_y), \quad (3.21b)$$

where X, Y, ϵ_x , and ϵ_y are arbitrary constants. Every orbit is closed since the periods of the oscillations in x and y are identical. The orbits form ellipses centered on the center of attraction. The azimuthal period is $T_\psi = 2\pi/\Omega$ because this is the time required for the star to return to its original azimuth. During this time, the particle completes two in-and-out cycles, so the radial period is

$$T_r = \frac{1}{2}T_\psi = \frac{\pi}{\Omega}. \quad (3.22)$$

(b) Kepler potential When the star is acted on by an inverse-square field $g(r) = -GM/r^2$ due to a point mass M , the corresponding potential is $\Phi = -GM/r = -GMu$. Motion in this potential is often called **Kepler motion**. Equation (3.11b) becomes

$$\frac{d^2u}{d\psi^2} + u = \frac{GM}{L^2}, \quad (3.23)$$

the general solution of which is

$$u(\psi) = C \cos(\psi - \psi_0) + \frac{GM}{L^2}, \quad (3.24)$$

where $C > 0$ and ψ_0 are arbitrary constants. Defining the orbit's **eccentricity** by

$$e \equiv \frac{CL^2}{GM} \quad (3.25a)$$

and its **semi-major axis** by

$$a \equiv \frac{L^2}{GM(1 - e^2)}, \quad (3.25b)$$

equation (3.24) may be rewritten

$$r(\psi) = \frac{a(1 - e^2)}{1 + e \cos(\psi - \psi_0)}. \quad (3.26)$$

An orbit for which $e \geq 1$ is unbound, since $r \rightarrow \infty$ as $(\psi - \psi_0) \rightarrow \pm \cos^{-1}(-1/e)$. We discuss unbound orbits in §3.1d below. Bound orbits have $e < 1$ and along them r is a periodic function of ψ with period 2π , so the star returns to its original radial coordinate after exactly one revolution in ψ . Thus bound Kepler orbits are closed, and one may show that they form ellipses with the attracting center at one focus. The pericenter and apocenter distances are

$$r_1 = a(1 - e) \quad \text{and} \quad r_2 = a(1 + e). \quad (3.27)$$

In many applications, equation (3.26) for r along a bound Kepler orbit is less convenient than the parameterization

$$r = a(1 - e \cos \eta), \quad (3.28a)$$

where the parameter η is called the **eccentric anomaly** to distinguish it from the **true anomaly**, $\psi - \psi_0$. By equating the right sides of equations (3.26) and (3.28a) and using the identity $\cos \theta = (1 - \tan^2 \frac{1}{2}\theta)/(1 + \tan^2 \frac{1}{2}\theta)$, it is straightforward to show that the true and eccentric anomalies are related by

$$\sqrt{1 - e} \tan \frac{1}{2}(\psi - \psi_0) = \sqrt{1 + e} \tan \frac{1}{2}\eta. \quad (3.29)$$

Equation (3.26) gives alternative relations between the two anomalies.

Taking $t = 0$ to occur at pericenter passage, from $L = r^2\dot{\psi}$ we have

$$t = \int_{\psi_0}^{\psi} \frac{d\psi}{\dot{\psi}} = \int d\psi \frac{r^2}{L} = \frac{a^2}{L} \int_0^{\eta} d\eta \frac{d\psi}{d\eta} (1 - e \cos \eta)^2. \quad (3.30)$$

Evaluating $d\psi/d\eta$ from (3.29), integrating, and using trigonometrical identities to simplify the result, we obtain finally

$$t = \frac{a^2}{L} \sqrt{1 - e^2} (\eta - e \sin \eta) = \frac{T_r}{2\pi} (\eta - e \sin \eta), \quad (3.28b)$$

where the second equality follows because the bracket on the right increases by 2π over an orbital period. This is called **Kepler's equation**, and the quantity $2\pi t/T_r$ is sometimes called the **mean anomaly**. Hence

$$T_r = T_\psi = \frac{a^2}{L} \sqrt{1 - e^2} = 2\pi \sqrt{\frac{a^3}{GM}}, \quad (3.31)$$

where the second equality uses (3.25b).

From (3.12) the energy per unit mass of a particle on a Kepler orbit is

$$E = -\frac{GM}{2a}. \quad (3.32)$$

To unbind the particle, we must add the **binding energy** $-E$.

The study of motion in nearly Kepler potentials is central to the dynamics of planetary systems (Murray & Dermott 1999).

We have shown that a star on a Kepler orbit completes a radial oscillation in the time required for ψ to increase by $\Delta\psi = 2\pi$, whereas a star that orbits in a harmonic-oscillator potential has already completed a radial oscillation by the time ψ has increased by $\Delta\psi = \pi$. Since galaxies are more extended than point masses, and less extended than homogeneous spheres, a typical star in a spherical galaxy completes a radial oscillation after its angular coordinate has increased by an amount that lies somewhere in between these two extremes; $\pi < \Delta\psi < 2\pi$ (cf. Problem 3.17). Thus, we expect a star to oscillate from its apocenter through its pericenter and back in a shorter time than is required for one complete azimuthal cycle about the galactic center.

It is sometimes useful to consider that an orbit in a non-Kepler force field forms an approximate ellipse, though one that **precesses** by $\psi_p = \Delta\psi - 2\pi$ in the time needed for one radial oscillation. For the orbit shown in Figure 3.1, and most galactic orbits, this precession is in the sense opposite to the rotation of the star itself. The angular velocity Ω_p of the rotating frame in which the ellipse appears closed is

$$\Omega_p = \frac{\psi_p}{T_r} = \frac{\Delta\psi - 2\pi}{T_r}. \quad (3.33)$$

Hence we say that Ω_p is the **precession rate** of the ellipse. The concept of closed orbits in a rotating frame of reference is crucial to the theory of spiral structure—see §6.2.1, particularly Figure 6.12.

(c) Isochrone potential The harmonic oscillator and Kepler potentials are both generated by mass distributions that are qualitatively different from the mass distributions of galaxies. The only known potential that could be generated by a realistic stellar system for which all orbits are analytic is the isochrone potential of equation (2.47) (Hénon 1959).

Box 3.1: Timing the local group

The nearest giant spiral galaxy is the Sb galaxy M31, at a distance of about (740 ± 40) kpc (BM §7.4.1). Our galaxy and M31 are by far the two largest members of the Local Group of galaxies. Beyond these, the next nearest prominent galaxies are in the Sculptor and M81 groups, at a distance of 3 Mpc. Thus the Local Group is an isolated system.

The line-of-sight velocity of the center of M31 relative to the center of the Galaxy is -125 km s^{-1} (for a solar circular speed $v_0 = 220 \text{ km s}^{-1}$, eq. 1.8); it is negative because the two galaxies are approaching one another. It seems that gravity has halted and reversed the original motion of M31 away from the Galaxy. Since M31 and the Galaxy are by far the most luminous members of the Local Group, we can treat them as an isolated system of two point masses, and estimate their total mass (Kahn & Woltjer 1959; Wilkinson & Evans 1999). Moreover, the original Hubble recession corresponded to an orbit of zero angular momentum, so we expect the angular momentum of the current orbit to be negligible. Thus we assume that the eccentricity $e = 1$.

We may now apply equations (3.28) for a Kepler orbit. Taking the log of both equations, differentiating with respect to η , and taking the ratio, we obtain

$$\frac{d \ln r}{d \ln t} = \frac{t}{r} \frac{dr}{dt} = \frac{e \sin \eta (\eta - e \sin \eta)}{(1 - e \cos \eta)^2}. \quad (1)$$

We set $e = 1$, and require that $r = 740$ kpc, $dr/dt = -125 \text{ km s}^{-1}$, and $t = 13.7$ Gyr, the current age of the universe (eq. 1.77). Inserting these constraints in (1) gives a nonlinear equation for η , which is easily solved numerically to yield $\eta = 4.29$. Then equations (3.28) yield $a = 524$ kpc and $T_r = 16.6$ Gyr, and equation (3.31) finally yields $M = 4.6 \times 10^{12} \mathcal{M}_\odot$ for the total mass of M31 and the Galaxy. The uncertainty in this result, assuming that our model is correct, is probably about a factor of 1.5.

This calculation assumes that the vacuum-energy density ρ_Λ is zero. Inclusion of non-zero ρ_Λ is simple (Problem 3.5); with parameters from equations (1.52) and (1.73), the required mass M increases by 15%.

The luminosity of the Galaxy in the R band is $3 \times 10^{10} L_\odot$ (Table 1.2) and M31 is about 1.5 times as luminous (BM Table 4.3); thus, if our mass estimate is correct, the mass-to-light ratio for the Local Group is $\Upsilon_V \simeq 60 \Upsilon_\odot$. This is far larger than expected for any normal stellar population, and the total mass is far larger than the masses within the outer edges of the disks of these galaxies, as measured by circular-speed curves. Thus the Kahn–Woltjer timing argument provided the first direct evidence that most of the mass of the Local Group is composed of dark matter. For a review see Peebles (1996).

Box 3.2: The eccentricity vector for Kepler orbits

The orbit of a test particle in the Kepler potential can also be found using vector methods. Since the angular momentum per unit mass $\mathbf{L} = \mathbf{r} \times \mathbf{v}$ is constant in any central field $g(r)$, with the equation of motion (3.3) and the vector identity (B.9) we have

$$\frac{d}{dt}(\mathbf{v} \times \mathbf{L}) = \frac{d\mathbf{v}}{dt} \times \mathbf{L} = g(r)\hat{\mathbf{e}}_r \times (\mathbf{r} \times \mathbf{v}) = g(r)[(\hat{\mathbf{e}}_r \cdot \mathbf{v})\mathbf{r} - r\mathbf{v}]. \quad (1)$$

The time derivative of the unit radial vector is

$$\frac{d\hat{\mathbf{e}}_r}{dt} = \frac{d}{dt}\left(\frac{\mathbf{r}}{r}\right) = \frac{\mathbf{v}}{r} - \frac{\mathbf{r} \cdot \mathbf{v}}{r^3}\mathbf{r} = \frac{1}{r^2}[r\mathbf{v} - (\hat{\mathbf{e}}_r \cdot \mathbf{v})\mathbf{r}]. \quad (2)$$

Comparing equations (1) and (2) we have

$$\frac{d}{dt}(\mathbf{v} \times \mathbf{L}) = -g(r)r^2\frac{d\hat{\mathbf{e}}_r}{dt}. \quad (3)$$

If and only if the field is Kepler, $g(r) = -GM/r^2$, this equation can be integrated to yield

$$\mathbf{v} \times \mathbf{L} = GM(\hat{\mathbf{e}}_r + \mathbf{e}), \quad (4)$$

where \mathbf{e} is a vector constant, or integral of motion (see §3.1.1). Taking the dot product of \mathbf{L} with equation (4), we find that $\mathbf{e} \cdot \mathbf{L} = 0$, so \mathbf{e} lies in the orbital plane. Taking the dot product of \mathbf{r} with equation (4) and using the vector identity (B.8), we have

$$L^2 = GM(r + \mathbf{e} \cdot \mathbf{r}). \quad (5)$$

If we now define ψ to be an azimuthal angle in the orbital plane, with \mathbf{e} at azimuth ψ_0 , then $\mathbf{e} \cdot \mathbf{r} = er \cos(\psi - \psi_0)$, where $e = |\mathbf{e}|$, and equation (5) can be rewritten

$$r = \frac{L^2}{GM} \frac{1}{1 + e \cos(\psi - \psi_0)}, \quad (6)$$

which is the same as equations (3.25b) and (3.26) for a Kepler orbit if we identify e with the eccentricity. It is therefore natural to call the vector constant \mathbf{e} the **eccentricity vector**, also sometimes called the Laplace or Runge–Lenz vector. The eccentricity vector has length equal to the eccentricity and points from the central mass towards the pericenter. The direction of the eccentricity vector is called the **line of apsides**.

Orbits in other central fields have integrals of motion analogous to the scalar eccentricity, but they do not have vector integrals analogous to the eccentricity vector, because orbits in non-Kepler potentials are not closed.

It is convenient to define an auxiliary variable s by

$$s \equiv -\frac{GM}{b\Phi} = 1 + \sqrt{1 + \frac{r^2}{b^2}}. \quad (3.34)$$

Solving this equation for r , we find that

$$\frac{r^2}{b^2} = s^2 \left(1 - \frac{2}{s}\right) \quad (s \geq 2). \quad (3.35)$$

Given this one-to-one relationship between s and r , we may employ s as a radial coordinate in place of r . The integrals (3.17) and (3.18b) for T_r and $\Delta\psi$ both involve the infinitesimal quantity

$$dI \equiv \frac{dr}{\sqrt{2(E - \Phi) - L^2/r^2}}. \quad (3.36)$$

When we use equation (3.35) to eliminate r from this expression, we find

$$dI = \frac{b(s-1)ds}{\sqrt{2Es^2 - 2(2E - GM/b)s - 4GM/b - L^2/b^2}}. \quad (3.37)$$

As the star moves from pericenter r_1 to apocenter r_2 , s varies from the smaller root s_1 of the quadratic expression in the denominator of equation (3.37) to the larger root s_2 . Thus, combining equations (3.17) and (3.37), the radial period is

$$T_r = \frac{2b}{\sqrt{-2E}} \int_{s_1}^{s_2} ds \frac{(s-1)}{\sqrt{(s_2-s)(s-s_1)}} = \frac{2\pi b}{\sqrt{-2E}} \left[\frac{1}{2}(s_1 + s_2) - 1\right], \quad (3.38)$$

where we have assumed $E < 0$ since we are dealing with bound orbits. But from the denominator of equation (3.37) it follows that the roots s_1 and s_2 obey

$$s_1 + s_2 = 2 - \frac{GM}{Eb}, \quad (3.39a)$$

and so the radial period

$$T_r = \frac{2\pi GM}{(-2E)^{3/2}}, \quad (3.39b)$$

exactly as in the Kepler case (the limit of the isochrone as $b \rightarrow 0$). Note that T_r depends on the energy E but not on the angular momentum L —it is this unique property that gives the isochrone its name.

Equation (3.18b), for the increment $\Delta\psi$ in azimuthal angle per cycle in the radial direction, yields

$$\begin{aligned} \Delta\psi &= 2L \int_{s_1}^{s_2} \frac{dI}{r^2} = \frac{2L}{b\sqrt{-2E}} \int_{s_1}^{s_2} ds \frac{(s-1)}{s(s-2)\sqrt{(s_2-s)(s-s_1)}} \\ &= \pi \operatorname{sgn}(L) \left(1 + \frac{|L|}{\sqrt{L^2 + 4GMb}}\right), \end{aligned} \quad (3.40)$$

where $\text{sgn}(L) = \pm 1$ depending on the sign of L . From this expression we see that

$$\pi < |\Delta\psi| < 2\pi. \quad (3.41)$$

The only orbits for which $|\Delta\psi|$ approaches the value 2π characteristic of Kepler motion are those with $L^2 \gg 4GMb$. Such orbits never approach the core $r \lesssim b$ of the potential, and hence always move in a near-Kepler field. In the opposite limit, $L^2 \ll 4GMb$, $|\Delta\psi| \rightarrow \pi$; physically this implies that low angular-momentum orbits fly straight through the core of the potential. In fact, the behavior $|\Delta\psi| \rightarrow \pi$ as $L \rightarrow 0$ is characteristic of any spherical potential that is not strongly singular at $r = 0$ —see Problem 3.19.

Inserting equations (3.39b) and (3.40) into equation (3.19), we have that the azimuthal period of an isochrone orbit is

$$T_\psi = \frac{4\pi GM}{(-2E)^{3/2}} \frac{\sqrt{L^2 + 4GMb}}{|L| + \sqrt{L^2 + 4GMb}}. \quad (3.42)$$

(d) Hyperbolic encounters In Chapter 7 we shall find that the dynamical evolution of globular clusters is largely driven by gravitational encounters between stars. These encounters are described by unbound Kepler orbits.

Let $(\mathbf{x}_M, \mathbf{v}_M)$ and $(\mathbf{x}_m, \mathbf{v}_m)$ be the positions and velocities of two point masses M and m , respectively; let $\mathbf{r} = \mathbf{x}_M - \mathbf{x}_m$ and $\mathbf{V} = \dot{\mathbf{r}}$. Then the separation vector \mathbf{r} obeys equation (D.33),

$$\left(\frac{mM}{M+m}\right)\ddot{\mathbf{r}} = -\frac{GMm}{r^2}\hat{\mathbf{e}}_r \quad \text{or} \quad \mu\ddot{\mathbf{r}} = -\frac{G(M+m)\mu}{r^2}\hat{\mathbf{e}}_r. \quad (3.43)$$

This is the equation of motion of a fictitious particle, called the reduced particle, which has mass $\mu = Mm/(M+m)$ and travels in the Kepler potential of a fixed body of mass $M+m$ (see Appendix D.1). If $\Delta\mathbf{v}_m$ and $\Delta\mathbf{v}_M$ are the changes in the velocities of m and M during the encounter, we have

$$\Delta\mathbf{v}_M - \Delta\mathbf{v}_m = \Delta\mathbf{V}. \quad (3.44a)$$

Furthermore, since the velocity of the center of mass of the two bodies is unaffected by the encounter (eq. D.19), we also have

$$M\Delta\mathbf{v}_M + m\Delta\mathbf{v}_m = 0. \quad (3.44b)$$

Eliminating $\Delta\mathbf{v}_m$ between equations (3.44) we obtain $\Delta\mathbf{v}_M$ as

$$\Delta\mathbf{v}_M = \frac{m}{M+m}\Delta\mathbf{V}. \quad (3.45)$$

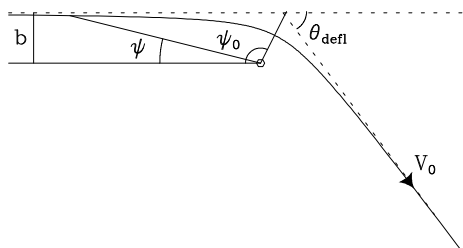


Figure 3.2 The motion of the reduced particle during a hyperbolic encounter.

We now evaluate $\Delta \mathbf{V}$.

Let the component of the initial separation vector that is perpendicular to the initial velocity vector $\mathbf{V}_0 = \mathbf{V}(t = -\infty)$ have length b (see Figure 3.2), the **impact parameter** of the encounter. Then the conserved angular momentum per unit mass associated with the motion of the reduced particle is

$$L = bV_0. \quad (3.46)$$

Equation (3.24), which relates the radius and azimuthal angle of a particle in a Kepler orbit, reads in this case,

$$\frac{1}{r} = C \cos(\psi - \psi_0) + \frac{G(M + m)}{b^2 V_0^2}, \quad (3.47)$$

where the angle ψ is shown in Figure 3.2. The constants C and ψ_0 are determined by the initial conditions. Differentiating equation (3.47) with respect to time, we obtain

$$\begin{aligned} \frac{dr}{dt} &= Cr^2 \dot{\psi} \sin(\psi - \psi_0) \\ &= CbV_0 \sin(\psi - \psi_0), \end{aligned} \quad (3.48)$$

where the second line follows because $r^2 \dot{\psi} = L$. If we define the direction $\psi = 0$ to point towards the particle as $t \rightarrow -\infty$, we find on evaluating equation (3.48) at $t = -\infty$,

$$-V_0 = CbV_0 \sin(-\psi_0). \quad (3.49a)$$

On the other hand, evaluating equation (3.47) at this time we have

$$0 = C \cos \psi_0 + \frac{G(M + m)}{b^2 V_0^2}. \quad (3.49b)$$

Eliminating C between these equations, we obtain

$$\tan \psi_0 = -\frac{bV_0^2}{G(M + m)}. \quad (3.50)$$

But from either (3.47) or (3.48) we see that the point of closest approach is reached when $\psi = \psi_0$. Since the orbit is symmetrical about this point, the angle through which the reduced particle's velocity is deflected is $\theta_{\text{defl}} = 2\psi_0 - \pi$ (see Figure 3.2). It proves useful to define the **90° deflection radius** as the impact parameter at which $\theta_{\text{defl}} = 90^\circ$:

$$b_{90} \equiv \frac{G(M+m)}{V_0^2}. \quad (3.51)$$

Thus

$$\theta_{\text{defl}} = 2 \tan^{-1} \left(\frac{G(M+m)}{bV_0^2} \right) = 2 \tan^{-1}(b_{90}/b). \quad (3.52)$$

By conservation of energy, the relative speed after the encounter equals the initial speed V_0 . Hence the components $\Delta \mathbf{V}_{\parallel}$ and $\Delta \mathbf{V}_{\perp}$ of $\Delta \mathbf{V}$ parallel and perpendicular to the original relative velocity vector \mathbf{V}_0 are given by

$$\begin{aligned} |\Delta \mathbf{V}_{\perp}| &= V_0 \sin \theta_{\text{defl}} = V_0 |\sin 2\psi_0| = \frac{2V_0 |\tan \psi_0|}{1 + \tan^2 \psi_0} \\ &= \frac{2V_0(b/b_{90})}{1 + b^2/b_{90}^2}, \end{aligned} \quad (3.53a)$$

$$\begin{aligned} |\Delta \mathbf{V}_{\parallel}| &= V_0(1 - \cos \theta_{\text{defl}}) = V_0(1 + \cos 2\psi_0) = \frac{2V_0}{1 + \tan^2 \psi_0} \\ &= \frac{2V_0}{1 + b^2/b_{90}^2}. \end{aligned} \quad (3.53b)$$

$\Delta \mathbf{V}_{\parallel}$ always points in the direction opposite to \mathbf{V}_0 . By equation (3.45) we obtain the components of $\Delta \mathbf{v}_M$ as

$$|\Delta \mathbf{v}_{M\perp}| = \frac{2mV_0}{M+m} \frac{b/b_{90}}{1 + b^2/b_{90}^2}, \quad (3.54a)$$

$$|\Delta \mathbf{v}_{M\parallel}| = \frac{2mV_0}{M+m} \frac{1}{1 + b^2/b_{90}^2}. \quad (3.54b)$$

$\Delta \mathbf{v}_{M\parallel}$ always points in the direction opposite to \mathbf{V}_0 . Notice that in the limit of large impact parameter b , $|\Delta \mathbf{v}_{M\perp}| = 2Gm/(bV_0)$, which agrees with the determination of the same quantity in equation (1.30).

3.1.1 Constants and integrals of the motion

Any stellar orbit traces a path in the six-dimensional space for which the coordinates are the position and velocity \mathbf{x}, \mathbf{v} . This space is called **phase space**.¹ A **constant of motion** in a given force field is any function

¹In statistical mechanics phase space usually refers to position-momentum space rather than position-velocity space. Since all bodies have the same acceleration in a given gravitational field, mass is irrelevant, and position-velocity space is more convenient.

$C(\mathbf{x}, \mathbf{v}; t)$ of the phase-space coordinates and time that is constant along stellar orbits; that is, if the position and velocity along an orbit are given by $\mathbf{x}(t)$ and $\mathbf{v}(t) = d\mathbf{x}/dt$,

$$C[\mathbf{x}(t_1), \mathbf{v}(t_1); t_1] = C[\mathbf{x}(t_2), \mathbf{v}(t_2); t_2] \quad (3.55)$$

for any t_1 and t_2 .

An **integral of motion** $I(\mathbf{x}, \mathbf{v})$ is any function of the phase-space coordinates alone that is constant along an orbit:

$$I[\mathbf{x}(t_1), \mathbf{v}(t_1)] = I[\mathbf{x}(t_2), \mathbf{v}(t_2)]. \quad (3.56)$$

While every integral is a constant of the motion, the converse is not true. For example, on a circular orbit in a spherical potential the azimuthal coordinate ψ satisfies $\psi = \Omega t + \psi_0$, where Ω is the star's constant angular speed and ψ_0 is its azimuth at $t = 0$. Hence $C(\psi, t) \equiv t - \psi/\Omega$ is a constant of the motion, but it is not an integral because it depends on time as well as the phase-space coordinates.

Any orbit in any force field always has six independent constants of motion. Indeed, since the initial phase-space coordinates $(\mathbf{x}_0, \mathbf{v}_0) \equiv [\mathbf{x}(0), \mathbf{v}(0)]$ can always be determined from $[\mathbf{x}(t), \mathbf{v}(t)]$ by integrating the equations of motion backward, $(\mathbf{x}_0, \mathbf{v}_0)$ can be regarded as six constants of motion.

By contrast, orbits can have from zero to five integrals of motion. In certain important cases, a few of these integrals can be written down easily: in any static potential $\Phi(\mathbf{x})$, the Hamiltonian $H(\mathbf{x}, \mathbf{v}) = \frac{1}{2}v^2 + \Phi$ is an integral of motion. If a potential $\Phi(R, z, t)$ is axisymmetric about the z axis, the z -component of the angular momentum is an integral, and in a spherical potential $\Phi(r, t)$ the three components of the angular-momentum vector $\mathbf{L} = \mathbf{x} \times \mathbf{v}$ constitute three integrals of motion. However, we shall find in §3.2 that even when integrals exist, analytic expressions for them are often not available.

These concepts and their significance for the geometry of orbits in phase space are nicely illustrated by the example of motion in a spherically symmetric potential. In this case the Hamiltonian H and the three components of the angular momentum per unit mass $\mathbf{L} = \mathbf{x} \times \mathbf{v}$ constitute four integrals. However, we shall find it more convenient to use $|\mathbf{L}|$ and the two independent components of the unit vector $\hat{\mathbf{n}} = \mathbf{L}/|\mathbf{L}|$ as integrals in place of \mathbf{L} . We have seen that $\hat{\mathbf{n}}$ defines the orbital plane within which the position vector \mathbf{r} and the velocity vector \mathbf{v} must lie. Hence we conclude that the two independent components of $\hat{\mathbf{n}}$ restrict the star's phase point to a four-dimensional region of phase space. Furthermore, the equations $H(\mathbf{x}, \mathbf{v}) = E$ and $|\mathbf{L}(\mathbf{x}, \mathbf{v})| = L$, where L is a constant, restrict the phase point to that two-dimensional surface in this four-dimensional region on which $v_r = \pm\sqrt{2[E - \Phi(r)] - L^2/r^2}$ and $v_\psi = L/r$. In §3.5.1 we shall see that this surface is a torus and that the sign ambiguity in v_r is analogous to the sign ambiguity in the z -coordinate

of a point on the sphere $r^2 = 1$ when one specifies the point through its x and y coordinates. Thus, given E , L , and $\hat{\mathbf{n}}$, the star's position and velocity (up to its sign) can be specified by two quantities, for example r and ψ .

Is there a fifth integral of motion in a spherical potential? To study this question, we examine motion in the potential

$$\Phi(r) = -GM \left(\frac{1}{r} + \frac{a}{r^2} \right). \quad (3.57)$$

For this potential, equation (3.11b) becomes

$$\frac{d^2u}{d\psi^2} + \left(1 - \frac{2GMa}{L^2} \right) u = \frac{GM}{L^2}, \quad (3.58)$$

the general solution of which is

$$u = C \cos \left(\frac{\psi - \psi_0}{K} \right) + \frac{GMK^2}{L^2}, \quad (3.59a)$$

where

$$K \equiv \left(1 - \frac{2GMa}{L^2} \right)^{-1/2}. \quad (3.59b)$$

Hence

$$\psi_0 = \psi - K \operatorname{Arccos} \left[\frac{1}{C} \left(\frac{1}{r} - \frac{GMK^2}{L^2} \right) \right], \quad (3.60)$$

where $t = \operatorname{Arccos} x$ is the multiple-valued solution of $x = \cos t$, and C can be expressed in terms of E and L by

$$E = \frac{1}{2} \frac{C^2 L^2}{K^2} - \frac{1}{2} \left(\frac{GMK}{L} \right)^2. \quad (3.61)$$

If in equations (3.59b), (3.60) and (3.61) we replace E by $H(\mathbf{x}, \mathbf{v})$ and L by $|\mathbf{L}(\mathbf{x}, \mathbf{v})| = |\mathbf{x} \times \mathbf{v}|$, the quantity ψ_0 becomes a function of the phase-space coordinates which is constant as the particle moves along its orbit. Hence ψ_0 is a fifth integral of motion. (Since the function $\operatorname{Arccos} x$ is multiple-valued, a judicious choice of solution is necessary to avoid discontinuous jumps in ψ_0 .) Now suppose that we know the numerical values of E , L , ψ_0 , and the radial coordinate r . Since we have four numbers—three integrals and one coordinate—it is natural to ask how we might use these numbers to determine the azimuthal coordinate ψ . We rewrite equation (3.60) in the form

$$\psi = \psi_0 \pm K \cos^{-1} \left[\frac{1}{C} \left(\frac{1}{r} - \frac{GMK^2}{L^2} \right) \right] + 2nK\pi, \quad (3.62)$$

where $\cos^{-1}(x)$ is defined to be the value of $\text{Arccos}(x)$ that lies between 0 and π , and n is an arbitrary integer. If K is irrational—as nearly all real numbers are—then by a suitable choice of the integer n , we can make ψ modulo 2π approximate any given number as closely as we please. Thus for any values of E and L , and any value of r between the pericenter and apocenter for the given E and L , an orbit that is known to have a given value of the integral ψ_0 can have an azimuthal angle as close as we please to any number between 0 and 2π .

On the other hand, if K is rational these problems do not arise. The simplest and most important case is that of the Kepler potential, when $a = 0$ and $K = 1$. Equation (3.62) now becomes

$$\psi = \psi_0 \pm \cos^{-1} \left[\frac{1}{C} \left(\frac{1}{r} - \frac{GM}{L^2} \right) \right] + 2n\pi, \quad (3.63)$$

which yields only two values of ψ modulo 2π for given E , L and r .

These arguments can be restated geometrically. The phase space has six dimensions. The equation $H(\mathbf{x}, \mathbf{v}) = E$ confines the orbit to a five-dimensional subspace. The vector equation $\mathbf{L}(\mathbf{x}, \mathbf{v}) = \text{constant}$ adds three further constraints, thereby restricting the orbit to a two-dimensional surface. Through the equation $\psi_0(\mathbf{x}, \mathbf{v}) = \text{constant}$ the fifth integral confines the orbit to a one-dimensional curve on this surface. Figure 3.1 can be regarded as a projection of this curve. In the Kepler case $K = 1$, the curve closes on itself, and hence does not cover the two-dimensional surface $H = \text{constant}$, $\mathbf{L} = \text{constant}$. But when K is irrational, the curve is endless and densely covers the surface of constant H and \mathbf{L} .

We can make an even stronger statement. Consider any volume of phase space, of any shape or size. Then if K is irrational, the fraction of the time that an orbit with given values of H and \mathbf{L} spends in that volume does not depend on the value that ψ_0 takes on this orbit.

Integrals like ψ_0 for irrational K that do not affect the phase-space distribution of an orbit, are called **non-isolating integrals**. All other integrals are called **isolating integrals**. The examples of isolating integrals that we have encountered so far, namely, H , \mathbf{L} , and the function ψ_0 when $K = 1$, all confine stars to a five-dimensional region in phase space. However, there can also be isolating integrals that restrict the orbit to a six-dimensional subspace of phase space—see §3.7.3. Isolating integrals are of great practical and theoretical importance, whereas non-isolating integrals are of essentially no value for galactic dynamics.

3.2 Orbits in axisymmetric potentials

Few galaxies are even approximately spherical, but many approximate figures of revolution. Thus in this section we begin to explore the types of orbits that are possible in many real galaxies. As in Chapter 2, we shall usually employ a cylindrical coordinate system (R, ϕ, z) with origin at the galactic center, and shall align the z axis with the galaxy's symmetry axis.

Stars whose motions are confined to the equatorial plane of an axisymmetric galaxy have no way of perceiving that the potential in which they move is not spherically symmetric. Therefore their orbits will be identical with those we discussed in the last section; the radial coordinate R of a star on such an orbit oscillates between fixed extrema as the star revolves around the center, and the orbit again forms a rosette figure.

3.2.1 Motion in the meridional plane

The situation is much more complex and interesting for stars whose motions carry them out of the equatorial plane of the system. The study of such general orbits in axisymmetric galaxies can be reduced to a two-dimensional problem by exploiting the conservation of the z -component of angular momentum of any star. Let the potential, which we assume to be symmetric about the plane $z = 0$, be $\Phi(R, z)$. Then the motion is governed by the Lagrangian

$$\mathcal{L} = \frac{1}{2} [\dot{R}^2 + (R\dot{\phi})^2 + \dot{z}^2] - \Phi(R, z). \quad (3.64)$$

The momenta are

$$p_R = \dot{R} \quad ; \quad p_\phi = R^2\dot{\phi} \quad ; \quad p_z = \dot{z}, \quad (3.65)$$

so the Hamiltonian is

$$H = \frac{1}{2} \left(p_R^2 + \frac{p_\phi^2}{R^2} + p_z^2 \right) + \Phi(R, z). \quad (3.66)$$

From Hamilton's equations (D.54) we find that the equations of motion are

$$\dot{p}_R = \ddot{R} = \frac{p_\phi^2}{R^3} - \frac{\partial\Phi}{\partial R}, \quad (3.67a)$$

$$\dot{p}_\phi = \frac{d}{dt}(R^2\dot{\phi}) = 0, \quad (3.67b)$$

$$\dot{p}_z = \ddot{z} = -\frac{\partial\Phi}{\partial z}. \quad (3.67c)$$

Equation (3.67b) expresses conservation of the component of angular momentum about the z axis, $p_\phi = L_z$ (a constant), while equations (3.67a) and (3.67c) describe the coupled oscillations of the star in the R and z -directions.

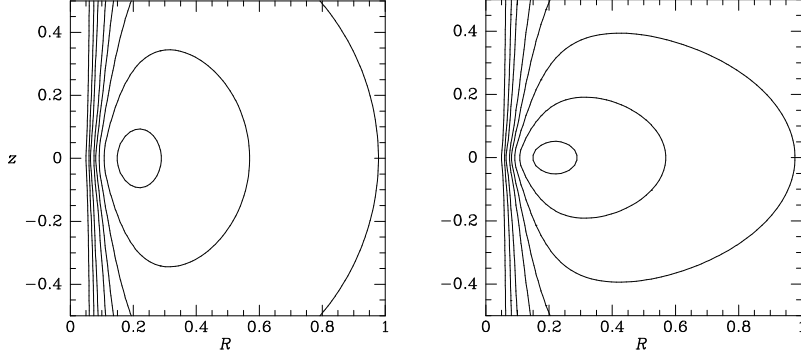


Figure 3.3 Level contours of the effective potential of equation (3.70) when $v_0 = 1$, $L_z = 0.2$. Contours are shown for $\Phi_{\text{eff}} = -1, -0.5, 0, 0.5, 1, 1.5, 2, 3, 5$. The axis ratio is $q = 0.9$ in the left panel and $q = 0.5$ in the right.

After replacing p_ϕ in (3.67a) by its numerical value L_z , the first and last of equations (3.67) can be written

$$\ddot{R} = -\frac{\partial\Phi_{\text{eff}}}{\partial R} \quad ; \quad \ddot{z} = -\frac{\partial\Phi_{\text{eff}}}{\partial z}, \quad (3.68a)$$

where

$$\Phi_{\text{eff}} \equiv \Phi(R, z) + \frac{L_z^2}{2R^2} \quad (3.68b)$$

is called the **effective potential**. Thus the three-dimensional motion of a star in an axisymmetric potential $\Phi(R, z)$ can be reduced to the two-dimensional motion of the star in the (R, z) plane (the **meridional plane**) under the Hamiltonian

$$H_{\text{eff}} = \frac{1}{2}(p_R^2 + p_z^2) + \Phi_{\text{eff}}(R, z). \quad (3.69)$$

Notice that H_{eff} differs from the full Hamiltonian (3.66) only in the substitution of the constant L_z for the azimuthal momentum p_ϕ . Consequently, the numerical value of H_{eff} is simply the orbit's total energy E . The difference $E - \Phi_{\text{eff}}$ is the kinetic energy of motion in the (R, z) plane, equal to $\frac{1}{2}(p_R^2 + p_z^2)$. Since kinetic energy is non-negative, the orbit is restricted to the area in the meridional plane satisfying the inequality $E \geq \Phi_{\text{eff}}$. The curve bounding this area is called the **zero-velocity curve**, since the orbit can only reach this curve if its velocity in the (R, z) plane is instantaneously zero.

Figure 3.3 shows contour plots of the effective potential

$$\Phi_{\text{eff}} = \frac{1}{2}v_0^2 \ln\left(R^2 + \frac{z^2}{q^2}\right) + \frac{L_z^2}{2R^2}, \quad (3.70)$$

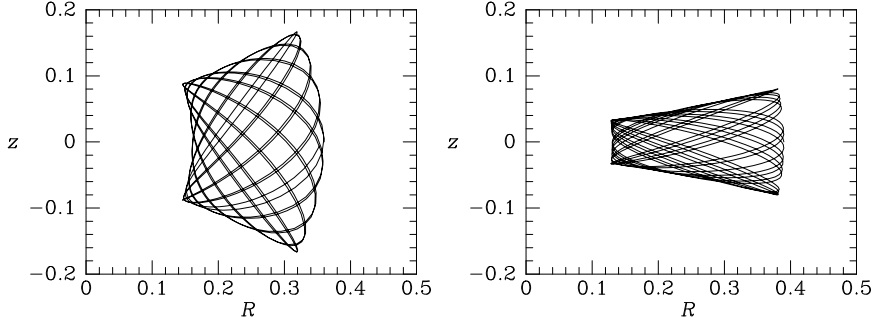


Figure 3.4 Two orbits in the potential of equation (3.70) with $q = 0.9$. Both orbits are at energy $E = -0.8$ and angular momentum $L_z = 0.2$, and we assume $v_0 = 1$.

for $v_0 = 1$, $L_z = 0.2$ and axial ratios $q = 0.9$ and 0.5 . This resembles the effective potential experienced by a star in an oblate spheroidal galaxy that has a constant circular speed v_0 (§2.3.2). Notice that Φ_{eff} rises very steeply near the z axis, as if the axis of symmetry were protected by a **centrifugal barrier**.

The minimum in Φ_{eff} has a simple physical significance. The minimum occurs where

$$0 = \frac{\partial \Phi_{\text{eff}}}{\partial R} = \frac{\partial \Phi}{\partial R} - \frac{L_z^2}{R^3} \quad ; \quad 0 = \frac{\partial \Phi_{\text{eff}}}{\partial z}. \quad (3.71)$$

The second of these conditions is satisfied anywhere in the equatorial plane $z = 0$ on account of the assumed symmetry of Φ about this plane, and the first is satisfied at the **guiding-center radius** R_g where

$$\left(\frac{\partial \Phi}{\partial R} \right)_{(R_g, 0)} = \frac{L_z^2}{R_g^3} = R_g \dot{\phi}^2. \quad (3.72)$$

This is simply the condition for a circular orbit with angular speed $\dot{\phi}$. Thus the minimum of Φ_{eff} occurs at the radius at which a circular orbit has angular momentum L_z , and the value of Φ_{eff} at the minimum is the energy of this circular orbit.

Unless the gravitational potential Φ is of some special form, equations (3.68a) cannot be solved analytically. However, we may follow the evolution of $R(t)$ and $z(t)$ by integrating the equations of motion numerically, starting from a variety of initial conditions. Figure 3.4 shows the result of two such integrations for the potential (3.69) with $q = 0.9$ (see Richstone 1982). The orbits shown are of stars of the same energy and angular momentum, yet they look quite different in real space, and hence the stars on these orbits must move through different regions of phase space. Is this because the equations of motion admit a third isolating integral $I(R, z, p_R, p_z)$ in addition to E and L_z ?

3.2.2 Surfaces of section

The phase space associated with the motion we are considering has four dimensions, R , z , p_R , and p_z , and the four-dimensional motion of the phase-space point of an individual star is too complicated to visualize. Nonetheless, we can determine whether orbits in the (R, z) plane admit an additional isolating integral by use of a simple graphical device. Since the Hamiltonian $H_{\text{eff}}(R, z, p_R, p_z)$ is constant, we could plot the motion of the representative point in a three-dimensional reduced phase space, say (R, z, p_R) , and then p_z would be determined (to within a sign) by the known value E of H_{eff} . However, even three-dimensional spaces are difficult to draw, so we simply show the points where the star crosses some plane in the reduced phase space, say the plane $z = 0$; these points are called **consequents**. To remove the sign ambiguity in p_z , we plot the (R, p_R) coordinates only when $p_z > 0$. In other words, we plot the values of R and p_R every time the star crosses the equator going upward. Such plots were first used by Poincaré and are called **surfaces of section**.² The key feature of the surface of section is that, even though it is only two-dimensional, no two distinct orbits at the same energy can occupy the same point. Also, any orbit is restricted to an area in the surface of section defined by the constraint $H_{\text{eff}} \geq \frac{1}{2}\dot{R}^2 + \Phi_{\text{eff}}$; the curve bounding this area is often called the zero-velocity curve of the surface of section, since it can only be reached by an orbit with $p_z = 0$.

Figure 3.5 shows the (R, p_R) surface of section at the energy of the orbits of Figure 3.4: the full curve is the zero-velocity curve, while the dots show the consequents generated by the orbit in the left panel of Figure 3.4. The cross near the center of the surface of section, at $(R = 0.26, p_R = 0)$, is the single consequent of the **shell orbit**, in which the trajectory of the star is restricted to a two-dimensional surface. The shell orbit is the limit of orbits such as those shown in Figure 3.4 in which the distance between the inner and outer boundaries of the orbit shrinks to zero.

In Figure 3.5 the consequents of the orbit of the left panel of Figure 3.4 appear to lie on a smooth curve, called the **invariant curve** of the orbit. The existence of the invariant curve implies that some isolating integral I is respected by this orbit. The curve arises because the equation $I = \text{constant}$ restricts motion in the two-dimensional surface of section to a one-dimensional curve (or perhaps to a finite number of discrete points in exceptional cases). It is often found that for realistic galactic potentials, orbits do admit an integral of this type. Since I is in addition to the two classical integrals H and p_ϕ , it is called the **third integral**. In general there is no analytic expression for I as a function of the phase-space variables, so it is called a **non-classical integral**.

² A surface of section is defined by some arbitrarily chosen condition, here $z = 0, p_z > 0$. Good judgment must be used in the choice of this condition lest some important orbits never satisfy it, and hence do not appear on the surface of section.

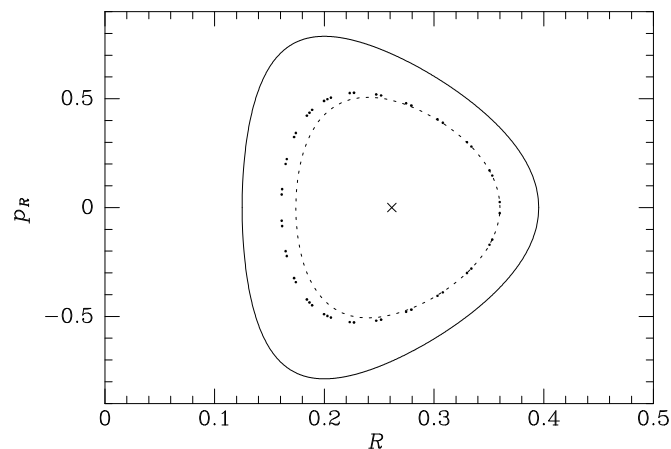


Figure 3.5 Points generated by the orbit of the left panel of Figure 3.4 in the (R, p_R) surface of section. If the total angular momentum L of the orbit were conserved, the points would fall on the dashed curve. The full curve is the zero-velocity curve at the energy of this orbit. The \times marks the consequent of the shell orbit.

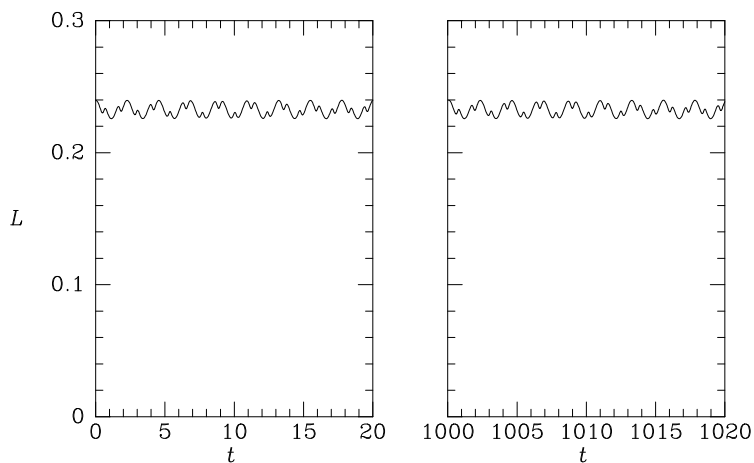


Figure 3.6 The total angular momentum is almost constant along the orbit shown in the left panel of Figure 3.5. For clarity $L(t)$ is plotted only at the beginning and end of a long integration.

We may form an intuitive picture of the nature of the third integral by considering two special cases. If the potential Φ is spherical, we know that the total angular momentum $|\mathbf{L}|$ is an integral. This suggests that for a nearly spherical potential—this one has axis ratio $q = 0.9$ —the third integral may be approximated by $|\mathbf{L}|$. The dashed curve in Figure 3.5 shows the curve

on which the points generated by the orbit of the left panel of Figure 3.4 would lie if the third integral were $|\mathbf{L}|$, and Figure 3.6 shows the actual time evolution of $|\mathbf{L}|$ along that orbit—notice that although $|\mathbf{L}|$ oscillates rapidly, its mean value does not change even over hundreds of orbital times. From these two figures we see that $|\mathbf{L}|$ is an approximately conserved quantity, even for orbits in potentials that are significantly flattened. We may think of these orbits as approximately planar and with more or less fixed peri- and apocenter radii. The approximate orbital planes have a fixed inclination to the z axis but precess about this axis, at a rate that gradually tends to zero as the potential becomes more and more nearly spherical.

The second special case is when the potential is separable in R and z :

$$\Phi(R, z) = \Phi_R(R) + \Phi_z(z). \quad (3.73)$$

Then the third integral can be taken to be the energy of vertical motion

$$H_z = \frac{1}{2}p_z^2 + \Phi_z(z). \quad (3.74)$$

Along nearly circular orbits in a thin disk, the potential is approximately separable, so equation (3.74) provides a useful expression for the third integral. In §3.6.2b we discuss a more sophisticated approximation to the third integral for orbits in thin disks.

3.2.3 Nearly circular orbits: epicycles and the velocity ellipsoid

In disk galaxies many stars are on nearly circular orbits, so it is useful to derive approximate solutions to equations (3.68a) that are valid for such orbits. We define

$$x \equiv R - R_g, \quad (3.75)$$

where $R_g(L_z)$ is the guiding-center radius for an orbit of angular momentum L_z (eq. 3.72). Thus $(x, z) = (0, 0)$ are the coordinates in the meridional plane of the minimum in Φ_{eff} . When we expand Φ_{eff} in a Taylor series about this point, we obtain

$$\Phi_{\text{eff}} = \Phi_{\text{eff}}(R_g, 0) + \frac{1}{2} \left(\frac{\partial^2 \Phi_{\text{eff}}}{\partial R^2} \right)_{(R_g, 0)} x^2 + \frac{1}{2} \left(\frac{\partial^2 \Phi_{\text{eff}}}{\partial z^2} \right)_{(R_g, 0)} z^2 + O(xz^2). \quad (3.76)$$

Note that the term that is proportional to xz vanishes because Φ_{eff} is assumed to be symmetric about $z = 0$. The equations of motion (3.68a) become very simple in the **epicycle approximation** in which we neglect all terms in Φ_{eff} of order xz^2 or higher powers of x and z . We define two new quantities by

$$\kappa^2(R_g) \equiv \left(\frac{\partial^2 \Phi_{\text{eff}}}{\partial R^2} \right)_{(R_g, 0)} \quad ; \quad \nu^2(R_g) \equiv \left(\frac{\partial^2 \Phi_{\text{eff}}}{\partial z^2} \right)_{(R_g, 0)}, \quad (3.77)$$

for then equations (3.68a) become

$$\ddot{x} = -\kappa^2 x, \quad (3.78a)$$

$$\ddot{z} = -\nu^2 z. \quad (3.78b)$$

According to these equations, x and z evolve like the displacements of two harmonic oscillators, with frequencies κ and ν , respectively. The two frequencies κ and ν are called the **epicycle** or **radial frequency** and the **vertical frequency**. If we substitute from equation (3.68b) for Φ_{eff} we obtain³

$$\kappa^2(R_g) = \left(\frac{\partial^2 \Phi}{\partial R^2} \right)_{(R_g, 0)} + \frac{3L_z^2}{R_g^4} = \left(\frac{\partial^2 \Phi}{\partial R^2} \right)_{(R_g, 0)} + \frac{3}{R_g} \left(\frac{\partial \Phi}{\partial R} \right)_{(R_g, 0)}, \quad (3.79a)$$

$$\nu^2(R_g) = \left(\frac{\partial^2 \Phi}{\partial z^2} \right)_{(R_g, 0)}. \quad (3.79b)$$

Since the circular frequency is given by

$$\Omega^2(R) = \frac{1}{R} \left(\frac{\partial \Phi}{\partial R} \right)_{(R, 0)} = \frac{L_z^2}{R^4}, \quad (3.79c)$$

equation (3.79a) may be written

$$\kappa^2(R_g) = \left(R \frac{d\Omega^2}{dR} + 4\Omega^2 \right)_{R_g}. \quad (3.80)$$

Note that the radial and azimuthal periods (eqs. 3.17 and 3.19) are simply

$$T_r = \frac{2\pi}{\kappa} \quad ; \quad T_\psi = \frac{2\pi}{\Omega}. \quad (3.81)$$

Very near the center of a galaxy, where the circular speed rises approximately linearly with radius, Ω is nearly constant and $\kappa \simeq 2\Omega$. Elsewhere Ω declines with radius, though rarely faster than the Kepler falloff, $\Omega \propto R^{-3/2}$, which yields $\kappa = \Omega$. Thus, in general,

$$\Omega \lesssim \kappa \lesssim 2\Omega. \quad (3.82)$$

Using equations (3.19) and (3.81), it is easy to show that this range is consistent with the range of $\Delta\psi$ given by equation (3.41) for the isochrone potential.

³ The formula for the ratio κ^2/Ω^2 from equations (3.79) was already known to Newton; see Proposition 45 of his *Principia*.

It is useful to define two functions

$$\begin{aligned} A(R) &\equiv \frac{1}{2} \left(\frac{v_c}{R} - \frac{dv_c}{dR} \right) = -\frac{1}{2} R \frac{d\Omega}{dR}, \\ B(R) &\equiv -\frac{1}{2} \left(\frac{v_c}{R} + \frac{dv_c}{dR} \right) = - \left(\Omega + \frac{1}{2} R \frac{d\Omega}{dR} \right), \end{aligned} \quad (3.83)$$

where $v_c(R) = R\Omega(R)$ is the circular speed at radius R . These functions are related to the circular and epicycle frequencies by

$$\Omega = A - B \quad ; \quad \kappa^2 = -4B(A - B) = -4B\Omega. \quad (3.84)$$

The values taken by A and B at the solar radius can be measured directly from the kinematics of stars in the solar neighborhood (BM §10.3.3) and are called the **Oort constants**.⁴ Taking values for these constants from Table 1.2, we find that the epicycle frequency at the Sun is $\kappa_0 = (37 \pm 3) \text{ km s}^{-1} \text{ kpc}^{-1}$, and that the ratio κ_0/Ω_0 at the Sun is

$$\frac{\kappa_0}{\Omega_0} = 2 \sqrt{\frac{-B}{A - B}} = 1.35 \pm 0.05. \quad (3.85)$$

Consequently the Sun makes about 1.3 oscillations in the radial direction in the time it takes to complete an orbit around the galactic center. Hence its orbit does not close on itself in an inertial frame, but forms a rosette figure like those discussed above for stars in spherically symmetric potentials.

The equations of motion (3.78) lead to two integrals, namely, the one-dimensional Hamiltonians

$$H_R \equiv \frac{1}{2}(\dot{x}^2 + \kappa^2 x^2) \quad ; \quad H_z \equiv \frac{1}{2}(\dot{z}^2 + \nu^2 z^2) \quad (3.86)$$

of the two oscillators. Thus if the star's orbit is sufficiently nearly circular that our truncation of the series for Φ_{eff} (eq. 3.76) is justified, then the orbit admits three integrals of motion: H_R , H_z , and p_ϕ . These are all isolating integrals.

From equations (3.75), (3.77), (3.78), and (3.86) we see that the Hamiltonian of such a star is made up of three parts:

$$H = H_R(R, p_R) + H_z(z, p_z) + \Phi_{\text{eff}}(R_g, 0). \quad (3.87)$$

⁴ Jan Hendrik Oort (1900–1992) was Director of Leiden Observatory in the Netherlands from 1945 to 1970. In 1927 Oort confirmed Bertil Lindblad's hypothesis of galactic rotation with an analysis of the motions of nearby stars that established the mathematical framework for studying Galactic rotation. With his student H. van de Hulst, he predicted the 21-cm line of neutral hydrogen. Oort also established the Netherlands as a world leader in radio astronomy, and showed that many comets originate in a cloud surrounding the Sun at a distance $\sim 0.1 \text{ pc}$, now called the Oort cloud.

Thus the three integrals of motion can equally be chosen as (H_R, H_z, p_ϕ) or (H, H_z, p_ϕ) , and in the latter case H_z , which is a classical integral, is playing the role of the third integral.

We now investigate what the ratios of the frequencies κ , Ω and ν tell us about the properties of the Galaxy. At most points in a typical galactic disk (including the solar neighborhood) $v_c \simeq \text{constant}$, and from (3.80) it is easy to show that in this case $\kappa^2 = 2\Omega^2$. In cylindrical coordinates Poisson's equation for an axisymmetric galaxy reads

$$\begin{aligned} 4\pi G\rho &= \frac{1}{R} \frac{\partial}{\partial R} \left(R \frac{\partial \Phi}{\partial R} \right) + \frac{\partial^2 \Phi}{\partial z^2} \\ &\simeq \frac{1}{R} \frac{dv_c^2}{dR} + \nu^2, \end{aligned} \quad (3.88)$$

where in the second line we have approximated the right side by its value in the equatorial plane and used equation (3.79b). If the mass distribution were spherical, we would have $\Omega^2 \simeq GM/R^3 = \frac{4}{3}\pi G\bar{\rho}$, where M is the mass and $\bar{\rho}$ is the mean density within the sphere of radius R about the galactic center. From the plot of the circular speed of an exponential disk shown in Figure 2.17, we know that this relation is not far from correct even for a flat disk. Hence, at a typical point in a galaxy such as the Milky Way

$$\frac{\nu^2}{\kappa^2} \simeq \frac{3}{2}\rho/\bar{\rho}. \quad (3.89)$$

That is, the ratio ν^2/κ^2 is a measure of the degree to which the galactic material is concentrated towards the plane, and will be significantly greater than unity for a disk galaxy. From Table 1.1 we see that at the Sun $\rho \simeq 0.1 \mathcal{M}_\odot \text{pc}^{-3}$, so the Sun's vertical period of small oscillations is $2\pi/\nu \simeq 87 \text{Myr}$. For $v_c = 220 \text{km s}^{-1}$ and $R_0 = 8 \text{kpc}$ (Table 1.2) we find $\bar{\rho} = 0.039 \mathcal{M}_\odot \text{pc}^{-3}$. Equation (3.89) then yields $\nu/\kappa \simeq 2.0$ for the Sun.

From equation (3.88) it is clear that we expect $\Phi_{\text{eff}} \propto z^2$ only for values of z small enough that $\rho_{\text{disk}}(z) \simeq \text{constant}$, i.e., for $z \ll 300 \text{pc}$ at R_0 . For stars that do not rise above this height, equation (3.78b) yields

$$z = Z \cos(\nu t + \zeta), \quad (3.90)$$

where Z and ζ are arbitrary constants. However, the orbits of the majority of disk stars carry these stars further above the plane than 300 pc (Problem 4.23). Therefore the epicycle approximation does not provide a reliable guide to the motion of the majority of disk stars in the direction perpendicular to the disk. The great value of this approximation lies rather in its ability to describe the motions of stars *in* the disk plane. So far we have described only the radial component of this motion, so we now turn to the azimuthal motion.

Equation (3.78a), which governs the radial motion, has the general solution

$$x(t) = X \cos(\kappa t + \alpha), \quad (3.91)$$

where $X \geq 0$ and α are arbitrary constants. Now let $\Omega_g = L_z/R_g^2$ be the angular speed of the circular orbit with angular momentum L_z . Since $p_\phi = L_z$ is constant, we have

$$\begin{aligned} \dot{\phi} &= \frac{p_\phi}{R^2} = \frac{L_z}{R_g^2} \left(1 + \frac{x}{R_g}\right)^{-2} \\ &\simeq \Omega_g \left(1 - \frac{2x}{R_g}\right). \end{aligned} \quad (3.92)$$

Substituting for x from (3.91) and integrating, we obtain

$$\phi = \Omega_g t + \phi_0 - \gamma \frac{X}{R_g} \sin(\kappa t + \alpha), \quad (3.93a)$$

where

$$\gamma \equiv \frac{2\Omega_g}{\kappa} = -\frac{\kappa}{2B}, \quad (3.93b)$$

where the second equality is derived using (3.84). The nature of the motion described by these equations can be clarified by erecting Cartesian axes (x, y, z) with origin at the **guiding center**, $(R, \phi) = (R_g, \Omega_g t + \phi_0)$. The x and z coordinates have already been defined, and the y coordinate is perpendicular to both and points in the direction of rotation.⁵ To first order in the small parameter X/R_g we have

$$\begin{aligned} y &= -\gamma X \sin(\kappa t + \alpha) \\ &\equiv -Y \sin(\kappa t + \alpha). \end{aligned} \quad (3.94)$$

Equations (3.91) and (3.94) are the complete solution for an equatorial orbit in the epicycle approximation. The motion in the z -direction is independent of the motion in x and y . In the (x, y) plane the star moves on an ellipse called the **epicycle** around the guiding center (see Figure 3.7). The lengths of the semi-axes of the epicycle are in the ratio

$$\frac{X}{Y} = \gamma^{-1}. \quad (3.95)$$

For a harmonic oscillator potential $X/Y = 1$ and for a Kepler potential $X/Y = \frac{1}{2}$; the inequality (3.82) shows that in most galactic potentials

⁵In applications to the Milky Way, which rotates clockwise when viewed from the north Galactic pole, either \hat{e}_z is directed towards the south Galactic pole, or (x, y, z) is a left-handed coordinate system; we make the second choice in this book.

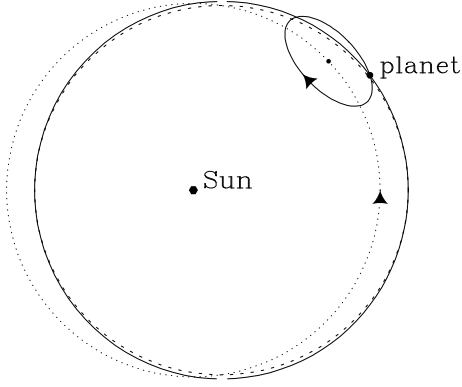


Figure 3.7 An elliptical Kepler orbit (dashed curve) is well approximated by the superposition of motion at angular frequency κ around a small ellipse with axis ratio $\frac{1}{2}$, and motion of the ellipse's center in the opposite sense at angular frequency Ω around a circle (dotted curve).

$Y > X$, so the epicycle is elongated in the tangential direction.⁶ From equation (3.85), $X/Y \simeq 0.7$ in the solar neighborhood. The motion around the epicycle is in the opposite sense to the rotation of the guiding center around the galactic center, and the period of the epicycle motion is $2\pi/\kappa$, while the period of the guiding-center motion is $2\pi/\Omega_g$.

Consider the motion of a star on an epicyclic orbit, as viewed by an astronomer who sits at the guiding center of the star's orbit. At different times in the orbit the astronomer's distance measurements range from a maximum value Y down to X . Since by equation (3.95), $X/Y = \kappa/(2\Omega_g)$, these measurements yield important information about the galactic potential. Of course, the epicycle period is much longer than an astronomer's lifetime, so we cannot in practice measure the distance to a given star as it moves around its epicycle. Moreover, in general we do not know the location of the guiding center of any given star. But we can measure v_R and $v_\phi(R_0) - v_c(R_0)$ for a group of stars, each of which has its own guiding-center radius R_g , as they pass near the Sun at radius R_0 . We now show that from these measurements we can determine the ratio $2\Omega/\kappa$. We have

$$\begin{aligned} v_\phi(R_0) - v_c(R_0) &= R_0(\dot{\phi} - \Omega_0) = R_0(\dot{\phi} - \Omega_g + \Omega_g - \Omega_0) \\ &\simeq R_0 \left[(\dot{\phi} - \Omega_g) - \left(\frac{d\Omega}{dR} \right)_{R_g} x \right]. \end{aligned} \quad (3.96a)$$

With equation (3.92) this becomes

$$v_\phi(R_0) - v_c(R_0) \simeq -R_0 x \left(\frac{2\Omega}{R} + \frac{d\Omega}{dR} \right)_{R_g}. \quad (3.96b)$$

⁶ Epicycles were invented by the Greek astronomer Hipparchus (190–120 BC) to describe the motion of the planets about the Sun. Hipparchus also measured the distance to the Moon and discovered the precession of the Earth's spin axis. Epicycles—the first known perturbation expansion—were not very successful, largely because Hipparchus used circular epicycles with $X/Y = 1$. If only he had used epicycles with the proper axis ratio $X/Y = \frac{1}{2}$!

If we evaluate the coefficient of the small quantity x at R_0 rather than R_g , we introduce an additional error in $v_\phi(R_0)$ which is of order x^2 and therefore negligible. Making this approximation we find

$$v_\phi(R_0) - v_c(R_0) \simeq -x \left(2\Omega + R \frac{d\Omega}{dR} \right)_{R_0}. \quad (3.96c)$$

Finally using equations (3.83) to introduce Oort's constants, we obtain

$$v_\phi(R_0) - v_c(R_0) \simeq 2Bx = \frac{\kappa}{\gamma} x = \frac{\kappa}{\gamma} X \cos(\kappa t + \alpha). \quad (3.97)$$

Averaging over the phases α of stars near the Sun, we find

$$\overline{[v_\phi - v_c(R_0)]^2} = \frac{\kappa^2 X^2}{2\gamma^2} = 2B^2 X^2. \quad (3.98)$$

Similarly, we may neglect the dependence of κ on R_g to obtain with equation (3.84)

$$\overline{v_R^2} = \frac{1}{2} \kappa^2 X^2 = -2B(A - B)X^2. \quad (3.99)$$

Taking the ratio of the last two equations we have

$$\frac{\overline{[v_\phi - v_c(R_0)]^2}}{\overline{v_R^2}} \simeq \frac{-B}{A - B} = -\frac{B}{\Omega_0} = \frac{\kappa_0^2}{4\Omega_0^2} = \gamma^{-2} \simeq 0.46. \quad (3.100)$$

In §4.4.3 we shall re-derive this equation from a rather different point of view and compare its predictions with observational data.

Note that the ratio in equation (3.100) is the *inverse* of the ratio of the mean-square azimuthal and radial velocities relative to the guiding center: by (3.95)

$$\frac{\overline{y^2}}{\overline{x^2}} = \frac{\frac{1}{2}(\kappa Y)^2}{\frac{1}{2}(\kappa X)^2} = \gamma^2. \quad (3.101)$$

This counter-intuitive result arises because one measure of the RMS tangential velocity (eq. 3.101) is taken with respect to the guiding center of a single star, while the other (eq. 3.100) is taken with respect to the circular speed at the star's instantaneous radius.

This analysis also leads to an alternative expression for the integral of motion H_R defined in equation (3.86). Eliminating x using equation (3.97), we have

$$H_R = \frac{1}{2} \dot{x}^2 + \frac{1}{2} \gamma^2 [v_\phi(R_0) - v_c(R_0)]^2. \quad (3.102)$$

3.3 Orbits in planar non-axisymmetric potentials

Many, possibly most, galaxies have non-axisymmetric structures. These are evident near the centers of many disk galaxies, where one finds a luminous stellar bar—the Milky Way possesses just such a bar (BM §10.3). Non-axisymmetry is harder to detect in an elliptical galaxy, but we believe that many elliptical galaxies, especially the more luminous ones, are triaxial rather than axisymmetric (BM §4.2). Evidently we need to understand how stars orbit in a non-axisymmetric potential if we are to model galaxies successfully.

We start with the simplest possible problem, namely, planar motion in a non-rotating potential.⁷ Towards the end of this section we generalize the discussion to two-dimensional motion in potentials whose figures rotate steadily, and in the next section we show how an understanding of two-dimensional motion can be exploited in problems involving three-dimensional potentials.

3.3.1 Two-dimensional non-rotating potential

Consider the logarithmic potential (cf. §2.3.2)

$$\Phi_L(x, y) = \frac{1}{2}v_0^2 \ln \left(R_c^2 + x^2 + \frac{y^2}{q^2} \right) \quad (0 < q \leq 1). \quad (3.103)$$

This potential has the following useful properties:

- (i) The equipotentials have constant axial ratio q , so the influence of the non-axisymmetry is similar at all radii. Since $q \leq 1$, the y axis is the minor axis.
- (ii) For $R = \sqrt{x^2 + y^2} \ll R_c$, we may expand Φ_L in powers of R/R_c and find

$$\Phi_L(x, y) \simeq \frac{v_0^2}{2R_c^2} \left(x^2 + \frac{y^2}{q^2} \right) + \text{constant} \quad (R \ll R_c), \quad (3.104)$$

which is just the potential of the two-dimensional harmonic oscillator. In §2.5 we saw that gravitational potentials of this form are generated by homogeneous ellipsoids. Thus for $R \lesssim R_c$, Φ_L approximates the potential of a homogeneous density distribution.

- (iii) For $R \gg R_c$ and $q = 1$, $\Phi_L \simeq v_0^2 \ln R$, which yields a circular speed $v_c \simeq v_0$ that is nearly constant. Thus the radial component of the force generated by Φ_L with $q \simeq 1$ is consistent with the flat circular-speed curves of many disk galaxies.

The simplest orbits in Φ_L are those that are confined to $R \ll R_c$; when Φ_L is of the form (3.104), the orbit is the sum of independent harmonic motions

⁷This problem is equivalent to that of motion in the meridional plane of an axisymmetric potential when $L_z = 0$.

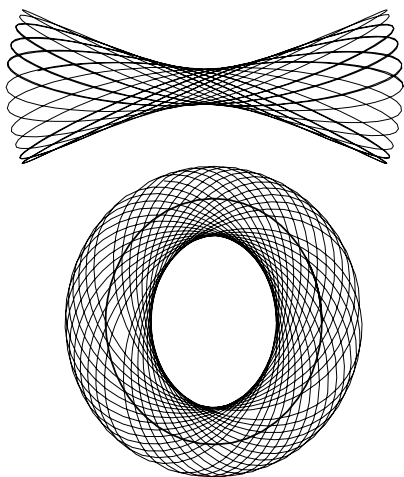


Figure 3.8 Two orbits of a common energy in the potential Φ_L of equation (3.103) when $v_0 = 1$, $q = 0.9$ and $R_c = 0.14$: top, a box orbit; bottom, a loop orbit. The closed parent of the loop orbit is also shown. The energy, $E = -0.337$, is that of the isopotential surface that cuts the long axis at $x = 5R_c$.

parallel to the x and y axes. The frequencies of these motions are $\omega_x = v_0/R_c$ and $\omega_y = v_0/qR_c$, and unless these frequencies are **commensurable** (i.e., unless $\omega_x/\omega_y = n/m$ for some integers n and m), the star eventually passes close to every point inside a rectangular box. These orbits are therefore known as **box orbits**.⁸ Such orbits have no particular sense of circulation about the center and thus their time-averaged angular momentum is zero. They respect two integrals of the motion, which we may take to be the Hamiltonians of the independent oscillations parallel to the coordinate axes,

$$H_x = \frac{1}{2}v_x^2 + \frac{1}{2}v_0^2 \frac{x^2}{R_c^2} \quad ; \quad H_y = \frac{1}{2}v_y^2 + \frac{1}{2}v_0^2 \frac{y^2}{q^2 R_c^2}. \quad (3.105)$$

To investigate orbits at larger radii $R \gtrsim R_c$, we must use numerical integrations. Two examples are shown in Figure 3.8. Neither orbit fills the elliptical zero-velocity curve $\Phi_L = E$, so both orbits must respect a second integral in addition to the energy. The upper orbit is still called a box orbit because it can be thought of as a distorted form of a box orbit in the two-dimensional harmonic oscillator. Within the core the orbit's envelope runs approximately parallel to the long axis of the potential, while for $R \gg R_c$ the envelope approximately follows curves of constant azimuth or radius.

In the lower orbit of Figure 3.8, the star circulates in a fixed sense about the center of the potential, while oscillating in radius. Orbits of this type are called **loop orbits**. Any star launched from $R \gg R_c$ in the tangential direction with a speed of order v_0 will follow a loop orbit. If the star is launched at speed $\sim v_0$ at a large angle to the tangential direction, the annulus occupied by the orbit will be wide, while if the launch angle is small, the annulus is narrow. This dependence is analogous to the way in which

⁸ The curve traced by a box orbit is sometimes called a **Lissajous figure** and is easily displayed on an oscilloscope.

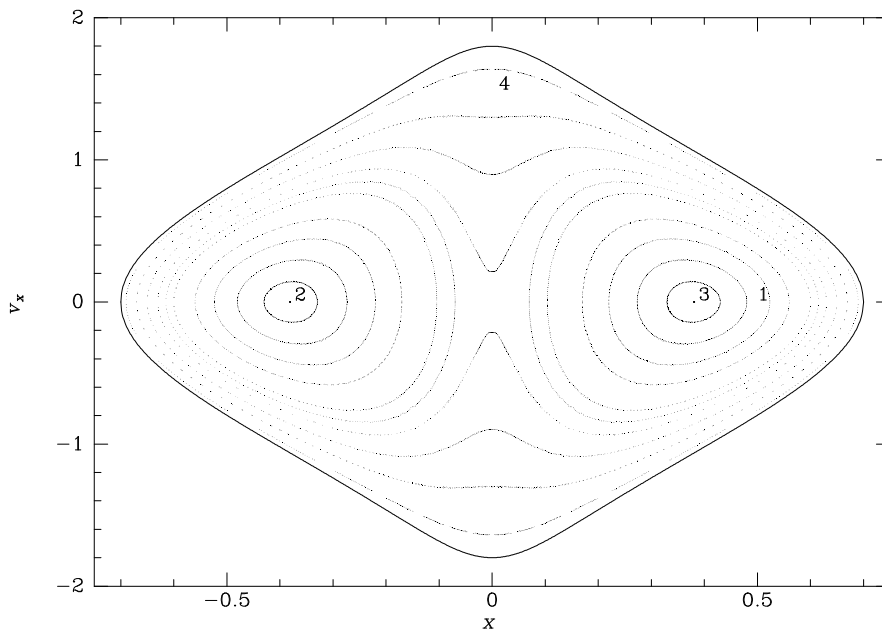


Figure 3.9 The (x, \dot{x}) surface of section formed by orbits in Φ_L of the same energy as the orbits depicted in Figure 3.8. The isopotential surface of this energy cuts the long axis at $x = 0.7$. The curves marked 4 and 1 correspond to the box and loop orbits shown in the top and bottom panels of Figure 3.8.

the thickness of the rosette formed by an orbit of given energy in a planar axisymmetric potential depends on its angular momentum. This analogy suggests that stars on loop orbits in Φ_L may respect an integral that is some sort of generalization of the angular momentum p_ϕ .

We may investigate these orbits further by generating a surface of section. Figure 3.9 is the surface of section $y = 0, \dot{y} > 0$ generated by orbits in Φ_L of the same energy as the orbits shown in Figure 3.8. The boundary curve in this figure arises from the energy constraint

$$\frac{1}{2}\dot{x}^2 + \Phi_L(x, 0) \leq \frac{1}{2}(\dot{x}^2 + \dot{y}^2) + \Phi_L(x, 0) = H_{y=0}. \quad (3.106)$$

Each closed curve in this figure corresponds to a different orbit. All these orbits respect an integral I_2 in addition to the energy because each orbit is confined to a curve.

There are two types of closed curve in Figure 3.9, corresponding to the two basic types of orbit that we have identified. The lower panel of Figure 3.8 shows the spatial form of the loop orbit that generates the curve marked 1 in Figure 3.9. At a given energy there is a whole family of such orbits that differ in the width of the elliptical annuli within which they are confined—see Figure 3.10. The unique orbit of this family that circulates in

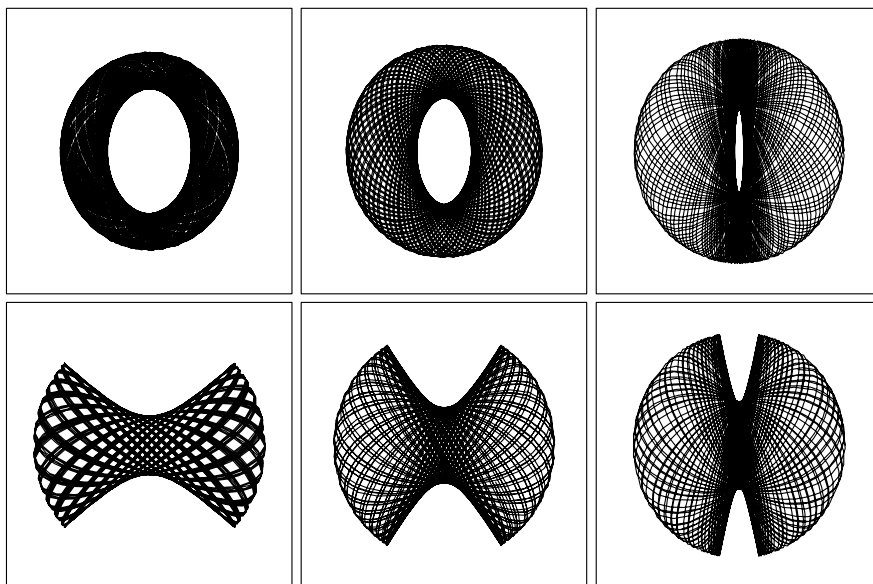


Figure 3.10 A selection of loop (top row) and box (bottom row) orbits in the potential $\Phi_L(q = 0.9, R_c = 0.14)$ at the energy of Figures 3.8 and 3.9.

an anti-clockwise sense and closes on itself after one revolution is the **closed loop orbit**, which is also shown at the bottom of Figure 3.8. In the surface of section this orbit generates the single point 3. Orbits with non-zero annular widths generate the curves that loop around the point 3. Naturally, there are loop orbits that circulate in a clockwise sense in addition to the anti-clockwise orbits; in the surface of section their representative curves loop around the point 2.

The second type of closed curve in Figure 3.9 corresponds to box orbits. The box orbit shown at the top of Figure 3.8 generates the curve marked 4. All the curves in the surface of section that are symmetric about the origin, rather than centered on one of the points 2 or 3, correspond to box orbits. These orbits differ from loop orbits in two major ways: (i) in the course of time a star on any of them passes arbitrarily close to the center of the potential (in the surface of section their curves cross $x = 0$), and (ii) stars on these orbits have no unique sense of rotation about the center (in the surface of section their curves are symmetric about $x = 0$). The outermost curve in Figure 3.9 (the zero-velocity curve) corresponds to the orbit on which $y = \dot{y} = 0$; on this orbit the star simply oscillates back and forth along the x axis. We call this the **closed long-axis orbit**. The curves interior to this bounding curve that also center on the origin correspond to less and less elongated box orbits. The bottom row of Figure 3.10 shows this progression from left to right. Notice the strong resemblance of the most eccentric loop

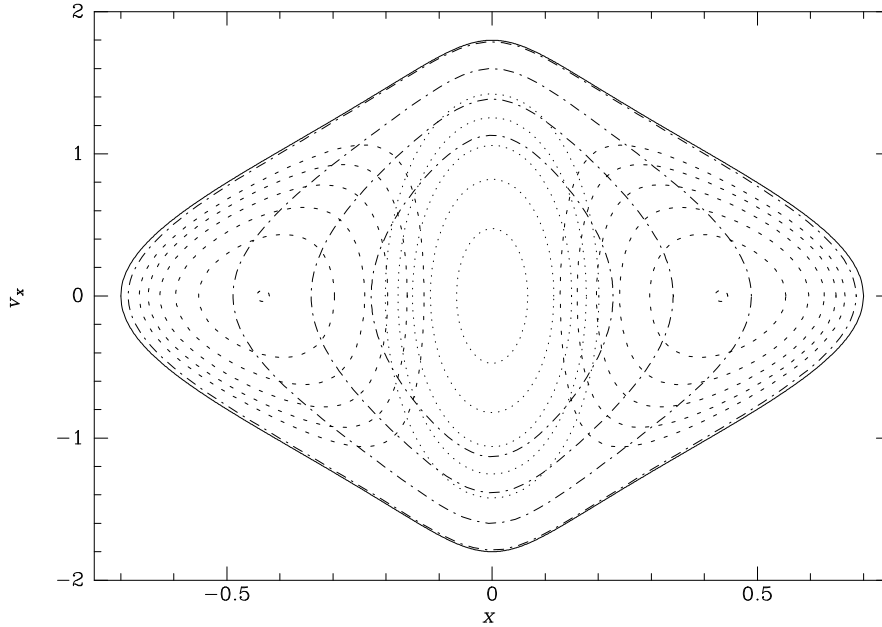


Figure 3.11 The appearance of the surface of section Figure 3.9 if orbits conserved (a) angular momentum (eq. 3.107; dashed curves), or (b) H_x (eq. 3.105; inner dotted curves), or (c) H'_x (eq. 3.108; outer dot-dashed curves).

orbit in the top right panel to the least elongated box orbit shown below it. The big difference between these orbits is that the loop orbit has a fixed sense of circulation about the center, while the box orbit does not.

It is instructive to compare the curves of Figure 3.9 with the curves generated by the integrals that we encountered earlier in this chapter. For example, if the angular momentum p_ϕ were an integral, the curves on the surface of section $y = 0$, $\dot{y} > 0$ would be given by the relation

$$(p_\phi)_{y=0} = x\dot{y} = x\sqrt{2[E - \Phi_L(x, 0)] - \dot{x}^2}. \quad (3.107)$$

These curves are shown as dashed curves in Figure 3.11. They resemble the curves in Figure 3.9 near the closed loop orbits 2 and 3, thus supporting our suspicion that the integral respected by loop orbits is some generalization of angular momentum. However, the dashed curves do not reproduce the curves generated by box orbits. If the extra integral were the Hamiltonian H_x of the x -component of motion in the harmonic potential (3.105), the curves in Figure 3.9 would be the dotted ellipses near the center of Figure 3.11. They resemble the curves in Figure 3.9 that are generated by the box orbits only in that they are symmetrical about the x axis. Figure 3.11 shows that a better approximation to the invariant curves of box orbits is provided by contours

of constant

$$H'_x \equiv \frac{1}{2}\dot{x}^2 + \Phi(x, 0). \quad (3.108)$$

H'_x may be thought of as the Hamiltonian associated with motion parallel to the potential's long axis. In a sense the integrals respected by box and loop orbits are analogous to H'_x and p_ϕ , respectively.

Figures 3.8 and 3.9 suggest an intimate connection between **closed orbits** and **families of non-closed orbits**. We say that the clockwise closed loop orbit is the **parent** of the family of clockwise loop orbits. Similarly, the closed long-axis orbit $y = 0$ is the parent of the box orbits.

The closed orbits that are the parents of orbit families are all **stable**, since members of their families that are initially close to them remain close at all times. In fact, we may think of any member of the family as engaged in stable oscillations about the parent closed orbit. A simple example of this state of affairs is provided by orbits in an axisymmetric potential. In a two-dimensional axisymmetric potential there are only two stable closed orbits at each energy—the clockwise and the anti-clockwise circular orbits.⁹ All other orbits, having non-zero eccentricity, belong to families whose parents are these two orbits. The epicycle frequency (3.80) is simply the frequency of small oscillations around the parent closed orbit.

The relationship between stable closed orbits and families of non-closed orbits enables us to trace the evolution of the orbital structure of a potential as the energy of the orbits or the shape of the potential is altered, simply by tracing the evolution of the stable closed orbits. For example, consider how the orbital structure supported by Φ_L (eq. 3.103) evolves as we pass from the axisymmetric potential that is obtained when $q = 1$ to the barred potentials that are obtained when $q < 1$. When $q = 1$, p_ϕ is an integral, so the surface of section is qualitatively similar to the dashed curves in Figure 3.11. The only stable closed orbits are circular, and all orbits are loop orbits. When we make q slightly smaller than unity, the long-axis orbit becomes stable and parents a family of elongated box orbits that oscillate about the axial orbit. As q is diminished more and more below unity, a larger and larger portion of phase space comes to be occupied by box rather than loop orbits. Comparison of Figures 3.9 and 3.12 shows that this evolution manifests itself in the surface of section by the growth of the band of box orbits that runs around the outside of Figure 3.12 at the expense of the two bull's-eyes in that figure that are associated with the loop orbits. In real space the closed loop orbits become more and more elongated, with the result that less and less epicyclic motion needs to be added to one of these closed orbits to fill in the hole at its center and thus terminate the sequence of loop orbits. The erosion of the bull's-eyes in the surface of section is associated with this process.

The appearance of the surface of section also depends on the energy of its orbits. Figure 3.13 shows a surface of section for motion in $\Phi_L(q =$

⁹ Special potentials such as the Kepler potential, in which all orbits are closed, must be excepted from this statement.

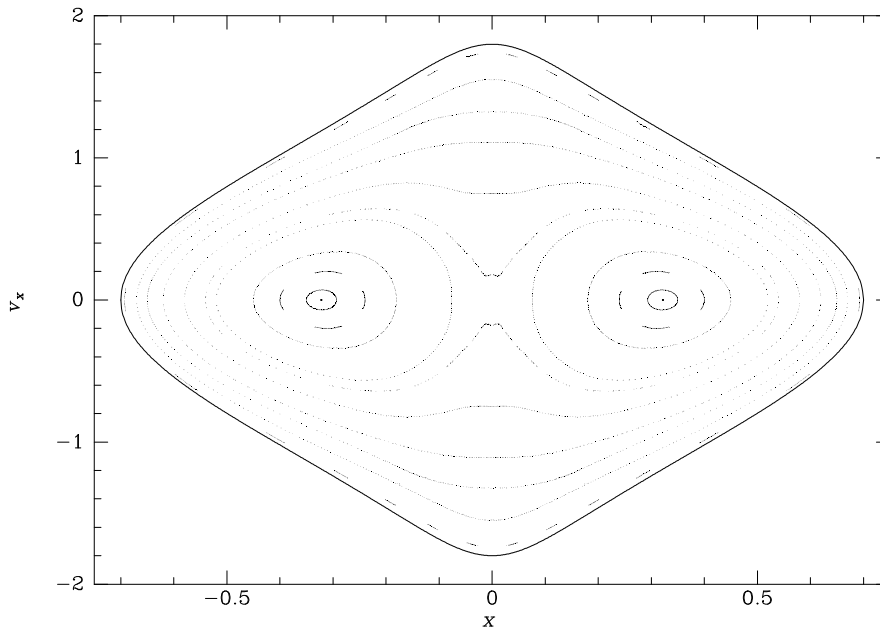


Figure 3.12 When the potential Φ_L is made more strongly barred by diminishing q , the proportion of orbits that are boxes grows at the expense of the loops: the figure shows the same surface of section as Figure 3.9 but for $q = 0.8$ rather than $q = 0.9$.

0.9, $R_c = 0.14$) at a lower energy than that of Figure 3.9. The changes in the surface of section are closely related to changes in the size and shape of the box and loop orbits. Box orbits that reach radii much greater than the core radius R_c have rather narrow waists (see Figure 3.10), and closed loop orbits of the same energy are nearly circular. If we consider box orbits and closed loop orbits of progressively smaller dimensions, the waists of the box orbits become steadily less narrow, and the closed orbits become progressively more eccentric as the dimensions of the orbits approach R_c . Eventually, at an energy E_c , the closed loop orbit degenerates into a line parallel to the short axis of the potential. Loop orbits do not exist at energies less than E_c . At $E < E_c$, all orbits are box orbits. The absence of loop orbits at $E < E_c$ is not unexpected since we saw above (eq. 3.105) that when $x^2 + y^2 \ll R_c^2$, the potential is essentially that of the two-dimensional harmonic oscillator, none of whose orbits are loops. At these energies the only closed orbits are the short- and the long-axis closed orbits, and we expect both of these orbits to be stable. In fact, the short-axis orbit becomes unstable at the energy E_c at which the loop orbits first appear. One says that the stable short-axis orbit of the low-energy regime **bifurcates** into the stable clockwise and anti-clockwise loop orbits at E_c . Stable closed orbits often appear in pairs like this.

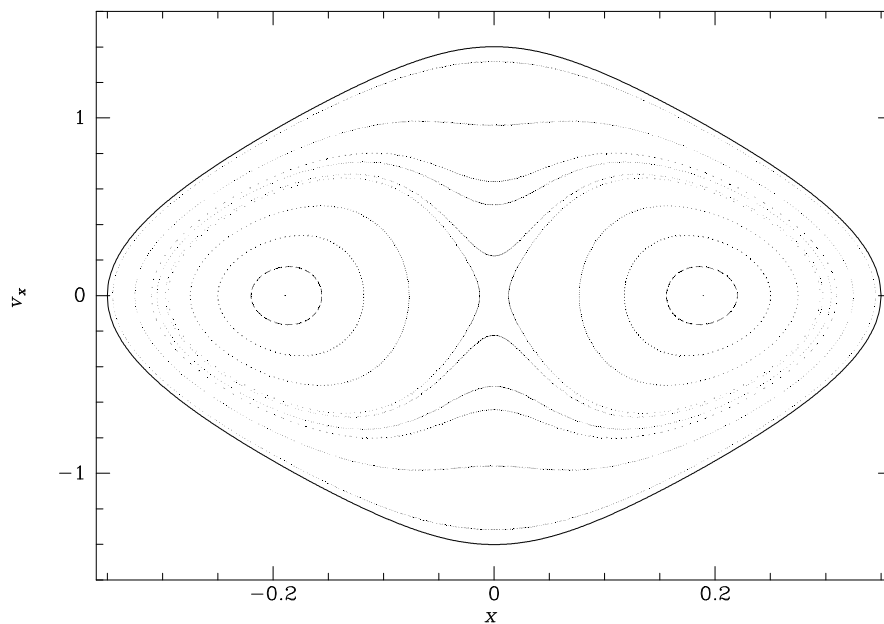


Figure 3.13 At low energies in a barred potential a large fraction of all orbits are boxes: the figure shows the same surface of section as Figure 3.9 but for the energy whose isopotential surface cuts the x axis at $x = 0.35$ rather than at $x = 0.7$ as in Figure 3.9.

Many two-dimensional barred potentials have orbital structures that resemble that of Φ_L . In particular:

- (i) Most orbits in these potentials respect a second integral in addition to energy.
- (ii) The majority of orbits in these potentials can be classified as either loop orbits or box orbits. The loop orbits have a fixed sense of rotation and never carry the star near the center, while the box orbits have no fixed sense of rotation and allow the star to pass arbitrarily close to the center.
- (iii) When the axial ratio of the isopotential curves is close to unity, most of the phase space is filled with loop orbits, but as the axial ratio changes away from unity, box orbits fill a bigger fraction of phase space.

Although these properties are fairly general, in §3.7.3 we shall see that certain barred potentials have considerably more complex orbital structures.

3.3.2 Two-dimensional rotating potential

The figures of many non-axisymmetric galaxies rotate with respect to inertial space, so we now study orbits in rotating potentials. Let the frame of reference in which the potential Φ is static rotate steadily at angular velocity Ω_b , often called the **pattern speed**. In this frame the velocity is $\dot{\mathbf{x}}$

and the corresponding velocity in an inertial frame is $\dot{\mathbf{x}} + \boldsymbol{\Omega}_b \times \mathbf{x}$. Thus the Lagrangian is

$$\mathcal{L} = \frac{1}{2} |\dot{\mathbf{x}} + \boldsymbol{\Omega}_b \times \mathbf{x}|^2 - \Phi(\mathbf{x}). \quad (3.109)$$

Consequently, the momentum is

$$\mathbf{p} = \frac{\partial \mathcal{L}}{\partial \dot{\mathbf{x}}} = \dot{\mathbf{x}} + \boldsymbol{\Omega}_b \times \mathbf{x}, \quad (3.110)$$

which is just the momentum in the underlying inertial frame. The Hamiltonian is

$$\begin{aligned} H_J &= \mathbf{p} \cdot \dot{\mathbf{x}} - \mathcal{L} \\ &= \mathbf{p} \cdot (\mathbf{p} - \boldsymbol{\Omega}_b \times \mathbf{x}) - \frac{1}{2} p^2 + \Phi \\ &= \frac{1}{2} p^2 + \Phi - \boldsymbol{\Omega}_b \cdot (\mathbf{x} \times \mathbf{p}), \end{aligned} \quad (3.111)$$

where we have used the vector identity (B.8). Since \mathbf{p} coincides with the momentum in an inertial frame, $\mathbf{x} \times \mathbf{p} = \mathbf{L}$ is the angular momentum and $\frac{1}{2} p^2 + \Phi$ is the Hamiltonian H that governs the motion in the inertial frame. Hence, (3.111) can be written

$$H_J = H - \boldsymbol{\Omega}_b \cdot \mathbf{L}. \quad (3.112)$$

Since $\Phi(\mathbf{x})$ is constant in the rotating frame, H_J has no explicit time dependence, and its derivative along any orbit $dH_J/dt = \partial H_J/\partial t$ vanishes (eq. D.56). Thus H_J is an integral, called the **Jacobi integral**: in a rotating non-axisymmetric potential, neither H nor \mathbf{L} is conserved, but the combination $H - \boldsymbol{\Omega}_b \cdot \mathbf{L}$ is conserved. From (3.111) it is easy to show that the constant value of H_J may be written as

$$\begin{aligned} E_J &= \frac{1}{2} |\dot{\mathbf{x}}|^2 + \Phi - \frac{1}{2} |\boldsymbol{\Omega}_b \times \mathbf{x}|^2 \\ &= \frac{1}{2} |\dot{\mathbf{x}}|^2 + \Phi_{\text{eff}}, \end{aligned} \quad (3.113)$$

where the effective potential

$$\begin{aligned} \Phi_{\text{eff}}(\mathbf{x}) &\equiv \Phi(\mathbf{x}) - \frac{1}{2} |\boldsymbol{\Omega}_b \times \mathbf{x}|^2 \\ &= \Phi(\mathbf{x}) - \frac{1}{2} [|\boldsymbol{\Omega}_b|^2 |\mathbf{x}|^2 - (\boldsymbol{\Omega}_b \cdot \mathbf{x})^2]. \end{aligned} \quad (3.114)$$

In deriving the second line we have used the identity (B.10). The effective potential is the sum of the gravitational potential and a repulsive **centrifugal potential**. For $\boldsymbol{\Omega}_b = \Omega_b \hat{\mathbf{e}}_z$, this additional term is simply $-\frac{1}{2} \Omega_b^2 R^2$ in cylindrical coordinates.

With equation (3.111) Hamilton's equations become

$$\begin{aligned} \dot{\mathbf{p}} &= -\frac{\partial H_J}{\partial \mathbf{x}} = -\nabla \Phi - \boldsymbol{\Omega}_b \times \mathbf{p} \\ \dot{\mathbf{x}} &= \frac{\partial H_J}{\partial \mathbf{p}} = \mathbf{p} - \boldsymbol{\Omega}_b \times \mathbf{x}, \end{aligned} \quad (3.115)$$

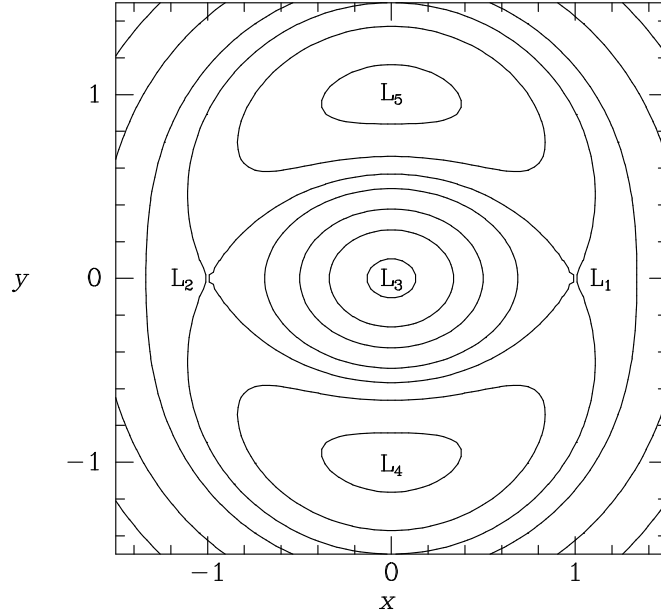


Figure 3.14 Contours of constant effective potential Φ_{eff} when the potential is given by equation (3.103) with $v_0 = 1$, $q = 0.8$, $R_c = 0.1$, and $\Omega_b = 1$. The point marked L_3 is a minimum of Φ_{eff} , while those marked L_4 and L_5 are maxima. Φ_{eff} has saddle points at L_1 and L_2 .

where we have used the identity (B.40). Eliminating \mathbf{p} between these equations we find

$$\begin{aligned}\ddot{\mathbf{x}} &= -\nabla\Phi - 2\boldsymbol{\Omega}_b \times \dot{\mathbf{x}} - \boldsymbol{\Omega}_b \times (\boldsymbol{\Omega}_b \times \mathbf{x}) \\ &= -\nabla\Phi - 2\boldsymbol{\Omega}_b \times \dot{\mathbf{x}} + |\boldsymbol{\Omega}_b|^2 \mathbf{x} - \boldsymbol{\Omega}_b(\boldsymbol{\Omega}_b \cdot \mathbf{x}).\end{aligned}\quad (3.116)$$

Here $-2\boldsymbol{\Omega}_b \times \dot{\mathbf{x}}$ is known as the **Coriolis force** and $-\boldsymbol{\Omega}_b \times (\boldsymbol{\Omega}_b \times \mathbf{x})$ is the **centrifugal force**. Taking the gradient of the last line of equation (3.114), we see that (3.116) can be written in the simpler form

$$\ddot{\mathbf{x}} = -\nabla\Phi_{\text{eff}} - 2\boldsymbol{\Omega}_b \times \dot{\mathbf{x}}. \quad (3.117)$$

The surface $\Phi_{\text{eff}} = E_J$ is often called the **zero-velocity surface**. All regions in which $\Phi_{\text{eff}} > E_J$ are forbidden to the star. Thus, although the solution of the differential equations for the orbit in a rotating potential may be difficult, we can at least define forbidden regions into which the star cannot penetrate.

Figure 3.14 shows contours of Φ_{eff} for the potential Φ_L of equation (3.103). Φ_{eff} is characterized by five stationary points, marked L_1 to L_5 ,

at which $\nabla\Phi_{\text{eff}} = 0$. These points are sometimes called **Lagrange points** after similar points in the restricted three-body problem (Figure 8.6). The central stationary point L_3 in Figure 3.14 is a minimum of the potential and is surrounded by a region in which the centrifugal potential $-\frac{1}{2}\Omega_b^2 R^2$ makes only a small contribution to Φ_{eff} . At each of the four points L_1 , L_2 , L_4 , and L_5 , it is possible for a star to travel on a circular orbit while appearing to be stationary in the rotating frame, because the gravitational and centrifugal forces precisely balance. Such orbits are said to **corotate** with the potential. The stationary points L_1 and L_2 on the x axis (the long axis of the potential) are saddle points, while the stationary points L_4 and L_5 along the y axis are maxima of the effective potential. Stars with values of E_J smaller than the value Φ_c taken by Φ_{eff} at L_1 and L_2 cannot move from the center of the potential to infinity, or indeed anywhere outside the inner equipotential contour that runs through L_1 and L_2 . By contrast, a star for which E_J exceeds Φ_c , or any star that is initially outside the contour through L_1 and L_2 , can *in principle* escape to infinity. However, it cannot be assumed that a star of the latter class will *necessarily* escape, because the Coriolis force prevents stars from accelerating steadily in the direction of $-\nabla\Phi_{\text{eff}}$.

We now consider motion near each of the Lagrange points L_1 to L_5 . These are stationary points of Φ_{eff} , so when we expand Φ_{eff} around one of these points $\mathbf{x}_L = (x_L, y_L)$ in powers of $(x - x_L)$ and $(y - y_L)$, we have

$$\begin{aligned} \Phi_{\text{eff}}(x, y) &= \Phi_{\text{eff}}(x_L, y_L) + \frac{1}{2} \left(\frac{\partial^2 \Phi_{\text{eff}}}{\partial x^2} \right)_{\mathbf{x}_L} (x - x_L)^2 \\ &+ \left(\frac{\partial^2 \Phi_{\text{eff}}}{\partial x \partial y} \right)_{\mathbf{x}_L} (x - x_L)(y - y_L) + \frac{1}{2} \left(\frac{\partial^2 \Phi_{\text{eff}}}{\partial y^2} \right)_{\mathbf{x}_L} (y - y_L)^2 + \dots \end{aligned} \quad (3.118)$$

Furthermore, for any bar-like potential whose principal axes lie along the coordinate axes, $\partial^2 \Phi_{\text{eff}} / \partial x \partial y = 0$ at \mathbf{x}_L by symmetry. Hence, if we retain only quadratic terms in equation (3.118) and define

$$\xi \equiv x - x_L \quad ; \quad \eta \equiv y - y_L, \quad (3.119)$$

and

$$\Phi_{xx} \equiv \left(\frac{\partial^2 \Phi_{\text{eff}}}{\partial x^2} \right)_{\mathbf{x}_L} \quad ; \quad \Phi_{yy} \equiv \left(\frac{\partial^2 \Phi_{\text{eff}}}{\partial y^2} \right)_{\mathbf{x}_L}, \quad (3.120)$$

the equations of motion (3.117) become for a star near \mathbf{x}_L ,

$$\ddot{\xi} = 2\Omega_b \dot{\eta} - \Phi_{xx} \xi \quad ; \quad \ddot{\eta} = -2\Omega_b \dot{\xi} - \Phi_{yy} \eta. \quad (3.121)$$

This is a pair of linear differential equations with constant coefficients. The general solution can be found by substituting $\xi = X \exp(\lambda t)$, $\eta = Y \exp(\lambda t)$,

where X , Y , and λ are complex constants. With these substitutions, equations (3.121) become

$$(\lambda^2 + \Phi_{xx})X - 2\lambda\Omega_b Y = 0 \quad ; \quad 2\lambda\Omega_b X + (\lambda^2 + \Phi_{yy})Y = 0. \quad (3.122)$$

These simultaneous equations have a non-trivial solution for X and Y only if the determinant

$$\begin{vmatrix} \lambda^2 + \Phi_{xx} & -2\lambda\Omega_b \\ 2\lambda\Omega_b & \lambda^2 + \Phi_{yy} \end{vmatrix} = 0. \quad (3.123)$$

Thus we require

$$\lambda^4 + \lambda^2 (\Phi_{xx} + \Phi_{yy} + 4\Omega_b^2) + \Phi_{xx}\Phi_{yy} = 0. \quad (3.124)$$

This is the **characteristic equation** for λ . It has four roots, which may be either real or complex. If λ is a root, $-\lambda$ is also a root, so if there is any root that has non-zero real part $\text{Re}(\lambda) = \gamma$, the general solution to equations (3.121) will contain terms that cause $|\xi|$ and $|\eta|$ to grow exponentially in time; $|\xi| \propto \exp(|\gamma|t)$ and $|\eta| \propto \exp(|\gamma|t)$. Under these circumstances essentially all orbits rapidly flee from the Lagrange point, and the approximation on which equations (3.121) rest breaks down. In this case the Lagrange point is said to be **unstable**.

When all the roots of equation (3.124) are pure imaginary, say $\lambda = \pm i\alpha$ or $\pm i\beta$, with $0 \leq \alpha \leq \beta$ real, the general solution to equations (3.121) is

$$\begin{aligned} \xi &= X_1 \cos(\alpha t + \phi_1) + X_2 \cos(\beta t + \phi_2), \\ \eta &= Y_1 \sin(\alpha t + \phi_1) + Y_2 \sin(\beta t + \phi_2), \end{aligned} \quad (3.125)$$

and the Lagrange point is stable, since the perturbations ξ and η oscillate rather than growing. Substituting these equations into the differential equations (3.121), we find that X_1 and Y_1 and X_2 and Y_2 are related by

$$Y_1 = \frac{\Phi_{xx} - \alpha^2}{2\Omega_b\alpha} X_1 = \frac{2\Omega_b\alpha}{\Phi_{yy} - \alpha^2} X_1, \quad (3.126a)$$

and

$$Y_2 = \frac{\Phi_{xx} - \beta^2}{2\Omega_b\beta} X_2 = \frac{2\Omega_b\beta}{\Phi_{yy} - \beta^2} X_2. \quad (3.126b)$$

The following three conditions are necessary and sufficient for both roots λ^2 of the quadratic equation (3.124) in λ^2 to be real and negative, and hence for the Lagrange point to be stable:

$$\begin{aligned} \text{(i)} \quad & \lambda_1^2 \lambda_2^2 = \Phi_{xx}\Phi_{yy} > 0, \\ \text{(ii)} \quad & \lambda_1^2 + \lambda_2^2 = -(\Phi_{xx} + \Phi_{yy} + 4\Omega_b^2) < 0, \\ \text{(iii)} \quad & \lambda^2 \text{ real} \Rightarrow (\Phi_{xx} + \Phi_{yy} + 4\Omega_b^2)^2 > 4\Phi_{xx}\Phi_{yy}. \end{aligned} \quad (3.127)$$

At saddle points of Φ_{eff} such as L_1 and L_2 , Φ_{xx} and Φ_{yy} have opposite signs, so these Lagrange points violate condition (i) and are always unstable. At a minimum of Φ_{eff} , such as L_3 , Φ_{xx} and Φ_{yy} are both positive, so conditions (i) and (ii) are satisfied. Condition (iii) is also satisfied because it can be rewritten in the form

$$(\Phi_{xx} - \Phi_{yy})^2 + 8\Omega_b^2(\Phi_{xx} + \Phi_{yy}) + 16\Omega_b^4 > 0, \quad (3.128)$$

which is satisfied whenever both Φ_{xx} and Φ_{yy} are positive. Hence L_3 is stable.

For future use we note that when Φ_{xx} and Φ_{yy} are positive, we may assume $\Phi_{xx} < \Phi_{yy}$ (since the x axis is the major axis of the potential) and we have already assumed that $\alpha < \beta$, so we can show from (3.124) that

$$\alpha^2 < \Phi_{xx} < \Phi_{yy} < \beta^2. \quad (3.129)$$

Also, when $\Omega_b^2 \rightarrow 0$, α^2 tends to Φ_{xx} , and β^2 tends to Φ_{yy} .

The stability of the Lagrange points at maxima of Φ_{eff} , such as L_4 and L_5 , depends on the details of the potential. For the potential Φ_L of equation (3.103) we have

$$\Phi_{\text{eff}} = \frac{1}{2}v_0^2 \ln \left(R_c^2 + x^2 + \frac{y^2}{q^2} \right) - \frac{1}{2}\Omega_b^2(x^2 + y^2), \quad (3.130)$$

so L_4 and L_5 occur at $(0, \pm y_L)$, where

$$y_L \equiv \sqrt{\frac{v_0^2}{\Omega_b^2} - q^2 R_c^2}, \quad (3.131)$$

and we see that L_4, L_5 are present only if $\Omega_b < v_0/(qR_c)$. Differentiating the effective potential again we find

$$\begin{aligned} \Phi_{xx}(0, y_L) &= -\Omega_b^2(1 - q^2) \\ \Phi_{yy}(0, y_L) &= -2\Omega_b^2 \left[1 - q^2 \left(\frac{\Omega_b R_c}{v_0} \right)^2 \right]. \end{aligned} \quad (3.132)$$

Hence $\Phi_{xx}\Phi_{yy}$ is positive if the Lagrange points exist, and stability condition (i) of (3.127) is satisfied. Deciding whether the other stability conditions hold is tedious in the general case, but straightforward in the limit of negligible core radius, $\Omega_b R_c/v_0 \ll 1$ (which applies, for example, to Figure 3.14). Then $\Phi_{xx} + \Phi_{yy} + 4\Omega_b^2 = \Omega_b^2(1 + q^2)$, so condition (ii) is satisfied. A straightforward calculation shows that condition (iii) holds—and thus that L_4 and L_5 are stable—providing $q^2 > \sqrt{32} - 5 \simeq (0.810)^2$. For future use we note that for small R_c , and to leading order in the ellipticity $\epsilon = 1 - q$, we have

$$\alpha^2 = 2\epsilon\Omega_b^2 = -\Phi_{xx} \quad ; \quad \beta^2 = 2(1 - 2\epsilon)\Omega_b^2 = 2\Omega_b^2 + O(\epsilon). \quad (3.133)$$

Equations (3.125) describing the motion about a stable Lagrange point show that each orbit is a superposition of motion at frequencies α and β around two ellipses. The shapes of these ellipses and the sense of the star's motion on them are determined by equations (3.126). For example, in the case of small R_c and ϵ , the α -ellipse around the point L_4 is highly elongated in the x - or ξ -direction (the tangential direction), while the β -ellipse has $Y_2 = -X_2/\sqrt{2}$. The star therefore moves around the β -ellipse in the sense opposite to that of the rotation of the potential. The β -ellipse is simply the familiar epicycle from §3.2.3, while the α -ellipse represents a slow tangential wallowing in the weak non-axisymmetric component of Φ_L .

Now consider motion about the central Lagrange point L_3 . From equations (3.126) and the inequality (3.129), it follows that $Y_1/X_1 > 0$. Thus the star's motion around the α -ellipse has the same sense as the rotation of the potential; such an orbit is said to be **prograde** or **direct**. When $\Omega_b^2 \ll |\Phi_{xx}|$, it is straightforward to show from equations (3.124) and (3.126) that $X_1 \gg Y_1$ and hence that this prograde motion runs almost parallel to the long axis of the potential—this is the long-axis orbit familiar to us from our study of non-rotating bars. Conversely the star moves around the β -ellipse in the sense opposite to that of the rotation of the potential (the motion is **retrograde**), and $|X_2| < |Y_2|$. When $\Omega_b^2/|\Phi_{xx}|$ is small, the β -ellipse goes over into the short-axis orbit of a non-rotating potential. A general prograde orbit around L_3 is made up of motion on the β -ellipse around a guiding center that moves around the α -ellipse, and conversely for retrograde orbits.

We now turn to a numerical study of orbits in rotating potentials that are not confined to the vicinity of a Lagrange point. We adopt the logarithmic potential (3.103) with $q = 0.8$, $R_c = 0.03$, $v_0 = 1$, and $\Omega_b = 1$. This choice places the corotation annulus near $R_{CR} = 30R_c$. The Jacobi integral (eq. 3.112) now plays the role that energy played in our similar investigation of orbits in non-rotating potentials, and by a slight abuse of language we shall refer to its value E_J as the “energy.” At radii $R \lesssim R_c$ the two important sequences of stable closed orbits in the non-rotating case are the long- and the short-axis orbits. Figure 3.15 confirms the prediction of our analytic treatment that in the presence of rotation these become oval in shape. Orbits of both sequences are stable and therefore parent families of non-closed orbits.

Consider now the evolution of the orbital structure as we leave the core region. At an energy E_1 , similar to that at which loop orbits first appeared in the non-rotating case, pairs of prograde orbits like those shown in Figure 3.16 appear. Only one member of the pair is stable. When it first appears, the stable orbit is highly elongated parallel to the short axis, but as the energy is increased it becomes more round. Eventually the decrease in the elongation of this orbit with increasing energy is reversed, the orbit again becomes highly elongated parallel to the short axis and finally disappears

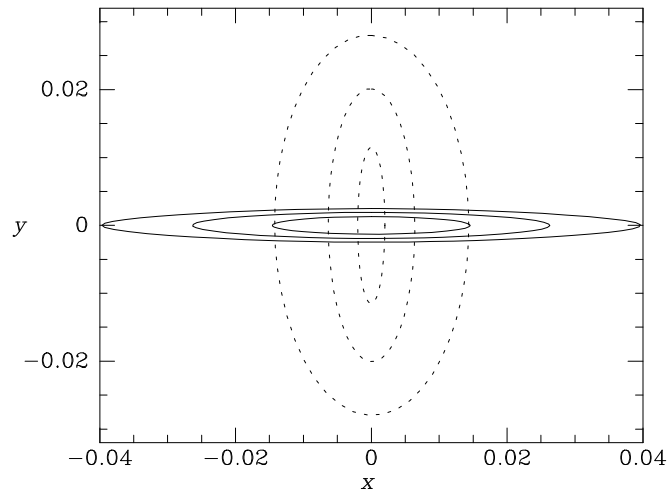


Figure 3.15 In the near-harmonic core of a rotating potential, the closed orbits are elongated ellipses. Stars on the orbits shown as full curves circulate about the center in the same sense as the potential's figure rotates. On the dashed orbits, stars circulate in the opposite sense. The x axis is the long axis of the potential.

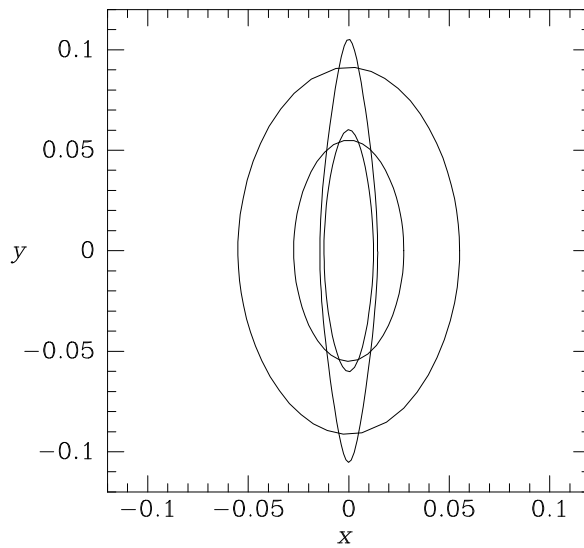


Figure 3.16 Closed orbits at two energies higher than those shown in Figure 3.15. Just outside the potential's near-harmonic core there are at each energy two prograde closed orbits aligned parallel to the potential's short axis. One of these orbits (the less elongated) is stable, while the other is unstable.

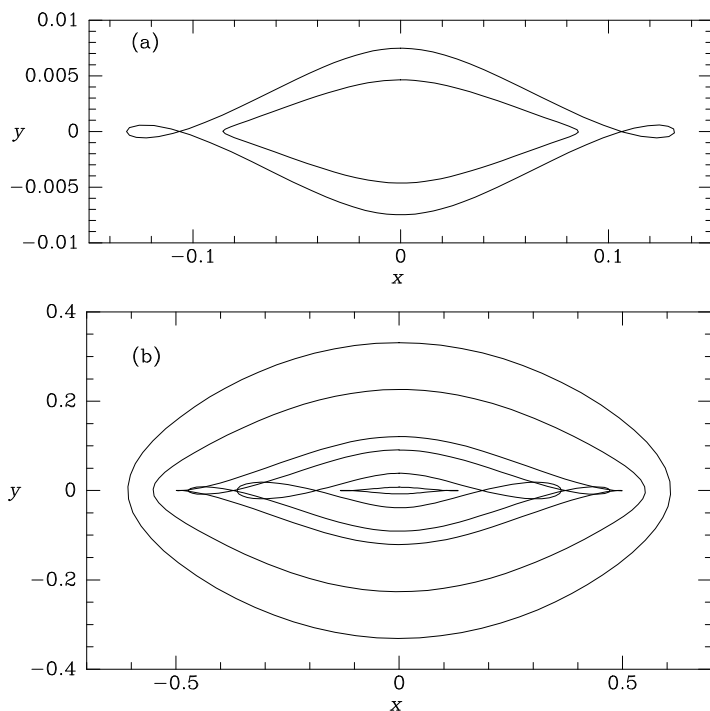


Figure 3.17 Near the energy at which the orbit pairs shown in Figure 3.16 appear, the closed long-axis orbits develop ears. Panel (a) shows orbits at energies just below and above this transition. Panel (b) shows the evolution of the closed long-axis orbits at higher energies. Notice that in panel (a) the x - and y -scales are different. The smallest orbit in panel (b) is the larger of the two orbits in panel (a).

along with its unstable companion orbit at an energy E_2 .¹⁰ In the notation of Contopoulos & Papayannopoulos (1980) these stable orbits are said to belong to the **sequence \mathbf{x}_2** , while their unstable companions are of the **sequence \mathbf{x}_3** .

The sequence of long-axis orbits (often called the **sequence \mathbf{x}_1**) suffers a significant transition near E_2 . On the low-energy side of the transition the long-axis orbits are extremely elongated and lens shaped (smaller orbit in Figure 3.17a). On the high-energy side the orbits are self-intersecting (larger orbit in Figure 3.17a). As the energy continues to increase, the orbit's ears become first more prominent and then less prominent, vanishing to form a cusped orbit (Figure 3.17b). At still higher energies the orbits become approximately elliptical (largest orbit in Figure 3.17b), first growing rounder and then adopt progressively more complex shapes as they approach

¹⁰ In the theory of weak bars, the energies E_1 and E_2 at which these prograde orbits appear and disappear are associated with the first and second inner Lindblad radii, respectively (eq. 3.150).

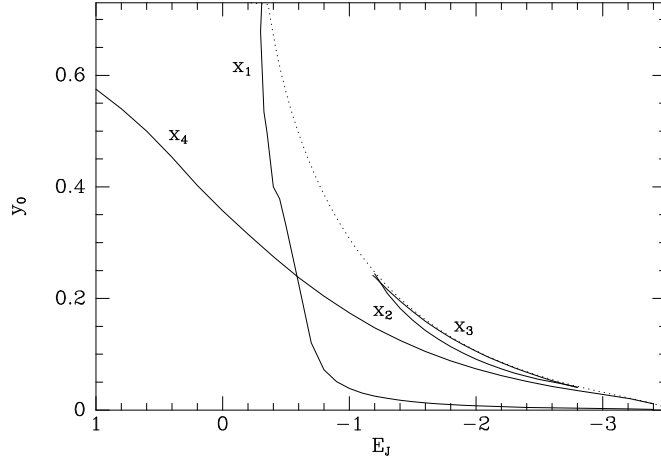


Figure 3.18 A plot of the Jacobi constant E_J of closed orbits in $\Phi_L(q = 0.8, R_c = 0.03, \Omega_b = 1)$ against the value of y at which the orbit cuts the potential's short axis. The dotted curve shows the relation $\Phi_{\text{eff}}(0, y) = E_J$. The families of orbits x_1 – x_4 are marked.

the corotation region in which the Lagrange points L_1 , L_2 , L_4 , and L_5 are located.

In the vicinity of the corotation annulus, there are important sequences of closed orbits on which stars move around one of the Lagrange points L_4 or L_5 , rather than about the center.

Essentially all closed orbits that carry stars well outside the corotation region are nearly circular. In fact, the potential's figure spins much more rapidly than these stars circulate on their orbits, so the non-axisymmetric forces on such stars tend to be averaged out. One finds that at large radii prograde orbits tend to align with the bar, while retrograde orbits align perpendicular to the bar.

These results are summarized in Figure 3.18. In this figure we plot against the value of E_J for each closed orbit the distance y at which it crosses the short axis of the potential. Each sequence of closed orbits generates a continuous curve in this diagram known as the **characteristic curve** of that sequence.

The stable closed orbits we have described are all associated with substantial families of non-closed orbits. Figure 3.19 shows two of these. As in the non-rotating case, a star on one of these non-closed orbits may be considered to be executing stable oscillations about one of the fundamental closed orbits. In potentials of the form (3.103) essentially all orbits belong to one of these families. This is not always true, however, as we explain in §3.7.

It is important to distinguish between orbits that enhance the elongation of the potential and those that oppose it. The overall mass distribution of a

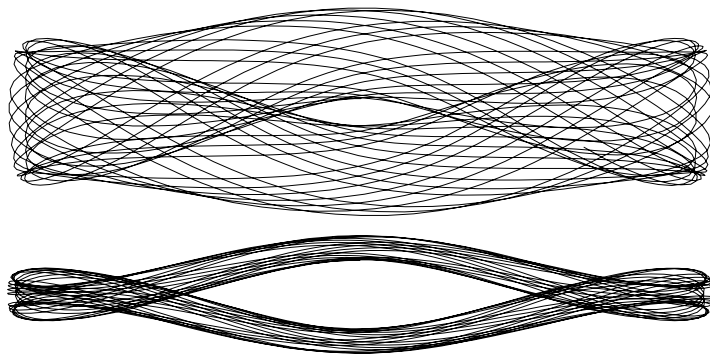


Figure 3.19 Two non-closed orbits of a common energy in the rotating potential Φ_L .

galaxy must be elongated in the same sense as the potential, which suggests that most stars are on orbits on which they spend the majority of their time nearer to the potential's long axis than to its short axis. Interior to the corotation radius, the only orbits that satisfy this criterion are orbits of the family parented by the long-axis orbits, which therefore must be the most heavily populated orbits in any bar that is confined by its own gravity. The shapes of these orbits range from butterfly-like at radii comparable to the core radius R_c , to nearly rectangular between R_c and the **inner Lindblad radius** (see below), to oval between this radius and corotation.

To an observer in an inertial frame of reference, stars on orbits belonging to the long-axis family circulate about the center of the potential in the same sense as the potential rotates. One part of the circulation seen by such an observer is due to the rotation of the frame of reference in which the potential is static. A second component of circulation is due to the mean streaming motion of such stars when referred to the rotating frame of the potential. Both components of circulation diminish towards zero if the angular velocity of the potential is reduced to zero. Near corotation the dominant component arises from the rotation of the frame of reference of the potential, while at small radii the more important component is the mean streaming motion of the stars through the rotating frame of reference.

3.3.3 Weak bars

Before we leave the subject of orbits in planar non-axisymmetric potentials, we derive an analytic description of loop orbits in weak bars.

(a) Lindblad resonances We assume that the figure of the potential rotates at some steady pattern speed Ω_b , and we seek to represent a general loop orbit as a superposition of the circular motion of a guiding center and small oscillations around this guiding center. Hence our treatment of orbits

in weak non-axisymmetric potentials will be closely related to the epicycle theory of nearly circular orbits in an axisymmetric potential (§3.2.3).

Let (R, φ) be polar coordinates in the frame that rotates with the potential, such that the line $\varphi = 0$ coincides with the long axis of the potential. Then the Lagrangian is

$$\mathcal{L} = \frac{1}{2}\dot{R}^2 + \frac{1}{2}[R(\dot{\varphi} + \Omega_b)]^2 - \Phi(R, \varphi), \quad (3.134)$$

so the equations of motion are

$$\ddot{R} = R(\dot{\varphi} + \Omega_b)^2 - \frac{\partial\Phi}{\partial R}, \quad (3.135a)$$

$$\frac{d}{dt}[R^2(\dot{\varphi} + \Omega_b)] = -\frac{\partial\Phi}{\partial\varphi}. \quad (3.135b)$$

Since we assume that the bar is weak, we may write

$$\Phi(R, \varphi) = \Phi_0(R) + \Phi_1(R, \varphi), \quad (3.136)$$

where $|\Phi_1/\Phi_0| \ll 1$. We divide R and φ into zeroth- and first-order parts

$$R(t) = R_0 + R_1(t) \quad ; \quad \varphi(t) = \varphi_0(t) + \varphi_1(t) \quad (3.137)$$

by substituting these expressions into equation (3.135) and requiring that the zeroth-order terms should sum to zero. Thus

$$R_0(\dot{\varphi}_0 + \Omega_b)^2 = \left(\frac{d\Phi_0}{dR}\right)_{R_0} \quad \text{and} \quad \dot{\varphi}_0 = \text{constant}. \quad (3.138)$$

This is the usual equation for centrifugal equilibrium at R_0 . If we define $\Omega_0 \equiv \Omega(R_0)$, where

$$\Omega(R) \equiv \pm \sqrt{\frac{1}{R} \frac{d\Phi_0}{dR}} \quad (3.139)$$

is the circular frequency at R in the potential Φ_0 , equation (3.138) for the angular speed of the guiding center (R_0, φ_0) becomes

$$\dot{\varphi}_0 = \Omega_0 - \Omega_b, \quad (3.140)$$

where $\Omega_0 > 0$ for prograde orbits and $\Omega_0 < 0$ for retrograde ones. We choose the origin of time such that

$$\varphi_0(t) = (\Omega_0 - \Omega_b)t. \quad (3.141)$$

The first-order terms in the equations of motion (3.135) now yield

$$\ddot{R}_1 + \left(\frac{d^2\Phi_0}{dR^2} - \Omega^2\right)_{R_0} R_1 - 2R_0\Omega_0\dot{\varphi}_1 = -\left(\frac{\partial\Phi_1}{\partial R}\right)_{R_0}, \quad (3.142a)$$

$$\ddot{\varphi}_1 + 2\Omega_0 \frac{\dot{R}_1}{R_0} = -\frac{1}{R_0^2} \left(\frac{\partial \Phi_1}{\partial \varphi} \right)_{R_0}. \quad (3.142b)$$

To proceed further we must choose a specific form of Φ_1 ; we set

$$\Phi_1(R, \varphi) = \Phi_b(R) \cos(m\varphi), \quad (3.143)$$

where m is a positive integer, since any potential that is an even function of φ can be expanded as a sum of terms of this form. In practice we are mostly concerned with the case $m = 2$ since the potential is then barred. If $\varphi = 0$ is to coincide with the long axis of the potential, we must have $\Phi_b < 0$.

So far we have assumed only that the angular velocity $\dot{\varphi}_1$ is small, not that φ_1 is itself small. Allowing for large excursions in φ_1 will be important when we consider what happens at resonances in part (b) of this section, but for the moment we assume that $\varphi_1 \ll 1$ and hence that $\varphi(t)$ always remains close to $(\Omega_0 - \Omega_b)t$. With this assumption we may replace φ by φ_0 in the expressions for $\partial\Phi_1/\partial R$ and $\partial\Phi_1/\partial\varphi$ to yield

$$\ddot{R}_1 + \left(\frac{d^2\Phi_0}{dR^2} - \Omega^2 \right)_{R_0} R_1 - 2R_0\Omega_0\dot{\varphi}_1 = - \left(\frac{d\Phi_b}{dR} \right)_{R_0} \cos [m(\Omega_0 - \Omega_b)t], \quad (3.144a)$$

$$\ddot{\varphi}_1 + 2\Omega_0 \frac{\dot{R}_1}{R_0} = \frac{m\Phi_b(R_0)}{R_0^2} \sin [m(\Omega_0 - \Omega_b)t]. \quad (3.144b)$$

Integrating the second of these equations, we obtain

$$\dot{\varphi}_1 = -2\Omega_0 \frac{R_1}{R_0} - \frac{\Phi_b(R_0)}{R_0^2(\Omega_0 - \Omega_b)} \cos [m(\Omega_0 - \Omega_b)t] + \text{constant}. \quad (3.145)$$

We now eliminate $\dot{\varphi}_1$ from equation (3.144a) to find

$$\ddot{R}_1 + \kappa_0^2 R_1 = - \left[\frac{d\Phi_b}{dR} + \frac{2\Omega\Phi_b}{R(\Omega - \Omega_b)} \right]_{R_0} \cos [m(\Omega_0 - \Omega_b)t] + \text{constant}, \quad (3.146a)$$

where

$$\kappa_0^2 \equiv \left(\frac{d^2\Phi_0}{dR^2} + 3\Omega^2 \right)_{R_0} = \left(R \frac{d\Omega^2}{dR} + 4\Omega^2 \right)_{R_0} \quad (3.146b)$$

is the usual epicycle frequency (eq. 3.80). The constant in equation (3.146a) is unimportant since it can be absorbed by a shift $R_1 \rightarrow R_1 + \text{constant}$.

Equation (3.146a) is the equation of motion of a harmonic oscillator of natural frequency κ_0 that is driven at frequency $m(\Omega_0 - \Omega_b)$. The general solution to this equation is

$$R_1(t) = C_1 \cos(\kappa_0 t + \alpha) - \left[\frac{d\Phi_b}{dR} + \frac{2\Omega\Phi_b}{R(\Omega - \Omega_b)} \right]_{R_0} \frac{\cos [m(\Omega_0 - \Omega_b)t]}{\Delta}, \quad (3.147a)$$

where C_1 and α are arbitrary constants, and

$$\Delta \equiv \kappa_0^2 - m^2(\Omega_0 - \Omega_b)^2. \quad (3.147b)$$

If we use equation (3.141) to eliminate t from equation (3.147a), we find

$$R_1(\varphi_0) = C_1 \cos\left(\frac{\kappa_0\varphi_0}{\Omega_0 - \Omega_b} + \alpha\right) + C_2 \cos(m\varphi_0), \quad (3.148a)$$

where

$$C_2 \equiv -\frac{1}{\Delta} \left[\frac{d\Phi_b}{dR} + \frac{2\Omega\Phi_b}{R(\Omega - \Omega_b)} \right]_{R_0}. \quad (3.148b)$$

If $C_1 = 0$, $R_1(\varphi_0)$ becomes periodic in φ_0 with period $2\pi/m$, and thus the orbit that corresponds to $C_1 = 0$ is a closed loop orbit. The orbits with $C_1 \neq 0$ are the non-closed loop orbits that are parented by this closed loop orbit. In the following we set $C_1 = 0$ so that we may study the closed loop orbits.

The right side of equation (3.148a) for R_1 becomes singular at a number of values of R_0 :

(i) **Corotation resonance.** When

$$\Omega_0 = \Omega_b, \quad (3.149)$$

$\dot{\varphi}_0 = 0$, and the guiding center corotates with the potential.

(ii) **Lindblad resonances.** When

$$m(\Omega_0 - \Omega_b) = \pm\kappa_0, \quad (3.150)$$

the star encounters successive crests of the potential at a frequency that coincides with the frequency of its natural radial oscillations. Radii at which such resonances occur are called **Lindblad radii** after the Swedish astronomer Bertil Lindblad (1895–1965). The plus sign in equation (3.150) corresponds to the case in which the star overtakes the potential, encountering its crests at the resonant frequency κ_0 ; this is called an **inner Lindblad resonance**. In the case of a minus sign, the crests of the potential sweep by the more slowly rotating star, and R_0 is said to be the radius of the **outer Lindblad resonance**.

There is a simple connection between these two types of resonance. A circular orbit has two natural frequencies. If the star is displaced radially, it oscillates at the epicycle frequency κ_0 . On the other hand, if the star is displaced azimuthally in such a way that it is still on a circular orbit, then it will continue on a circular orbit displaced from the original one. Thus the star is neutrally stable to displacements of this form; in other words, its natural azimuthal frequency is zero. The two types of resonance arise when the

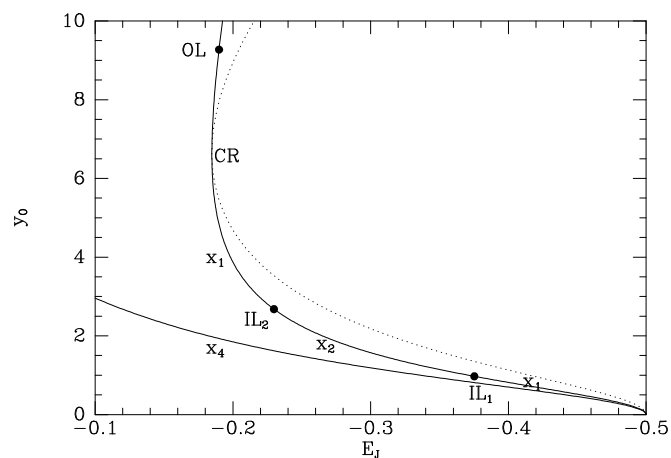


Figure 3.20 The full curves are the characteristic curves of the prograde (upper) and retrograde (lower) circular orbits in the isochrone potential (2.47) when a rotating frame of reference is employed. The dashed curve shows the relation $\Phi_{\text{eff}}(0, y) = E_J$, and the dots mark the positions of the Lindblad resonances when a small non-axisymmetric component is added to the potential.

forcing frequency seen by the star, $m(\Omega_0 - \Omega_b)$, equals one of the natural frequencies $\pm\kappa_0$ and 0.

Figure 6.11 shows plots of Ω , $\Omega + \frac{1}{2}\kappa$ and $\Omega - \frac{1}{2}\kappa$ for two circular-speed curves typical of galaxies. A galaxy may have 0, 1, 2, or more Lindblad resonances. The Lindblad and corotation resonances play a central role in the study of bars and spiral structure, and we shall encounter them again Chapter 6.

From equation (3.148a) it follows that for $m = 2$ the closed loop orbit is aligned with the bar whenever $C_2 > 0$, and is aligned perpendicular to the bar when $C_2 < 0$. When R_0 passes through a Lindblad or corotation resonance, the sign of C_2 , and therefore the orientation of the closed loop orbits, changes.

It is interesting to relate the results of this analytic treatment to the orbital structure of a strong bar that we obtained numerically in the last subsection. In this connection it is helpful to compare Figure 3.18, which shows data for a barred potential, with Figure 3.20, which describes orbits in an axisymmetric potential viewed from a rotating frame. The full curves in Figure 3.20 show the relationship between the Jacobi constant E_J and the radii of prograde and retrograde circular orbits in the isochrone potential (2.47). As in Figure 3.18, the dotted curve marks the relation $\Phi_{\text{eff}}(0, y) = E_J$. There are no orbits in the region to the right of this curve, which touches the curve of the prograde circular orbits at the corotation resonance, marked CR in the figure. If in the given frame we were to add a small

non-axisymmetric component to the potential, the orbits marked by large dots would lie at the Lindblad resonances (from right to left, the first and second inner Lindblad resonances and the outer Lindblad resonance marked OL). We call the radius of the first inner Lindblad resonance¹¹ R_{IL1} , and similarly R_{IL2} , R_{OL} , and R_{CR} for the radii of the other Lindblad resonances and of corotation. Equations (3.148) with $C_1 = 0$ describe nearly circular orbits in a weakly barred potential. Comparing Figure 3.20 with Figure 3.18, we see that nearly circular retrograde orbits belong to the family x_4 . Nearly circular prograde orbits belong to different families depending on their radius. Orbits that lie within R_{IL1} belong to the family x_1 . In the radius range $R_{\text{IL1}} < R < R_{\text{IL2}}$ the families x_2 and x_3 exist and contain orbits that are more circular than those of x_1 . We identify the orbits described by (3.148) with orbits of the family x_2 as indicated in Figure 3.20, since the family x_3 is unstable. In the radius range $R > R_{\text{IL2}}$, equations (3.148) with $C_1 = 0$ describe orbits of the family x_1 . Thus equations (3.148) describe only the families of orbits in a barred potential that are parented by a nearly circular orbit. However, when the non-axisymmetric component of the potential is very weak, most of phase space is occupied by such orbits. As the non-axisymmetry of the potential becomes stronger, families of orbits that are not described by equations (3.148) become more important.

(b) Orbits trapped at resonance When R_0 approaches the radius of either a Lindblad resonance or the corotation resonance, the value of R_1 that is predicted by equations (3.148) becomes large, and our linearized treatment of the equations of motion breaks down. However, one can modify the analysis to cope with these resonances. We now discuss the necessary modifications for the case of the corotation resonance. The case of the Lindblad resonances is described in Goldreich & Tremaine (1981).

The appropriate modification is suggested by our investigation of orbits near the Lagrange points L_4 and L_5 in the potential Φ_L (eq. 3.103), when the radius is large compared to the core radius R_c and the ellipticity $\epsilon = 1 - q$ approaches zero. In this limit the non-axisymmetric part of the potential is proportional to ϵ , so we have an example of a weak bar when $\epsilon \rightarrow 0$. We found in §3.3.2 that a star's orbit was a superposition of motion at frequencies α and β around two ellipses. In the limit $\epsilon \rightarrow 0$, the β -ellipse represents the familiar epicyclic motion and will not be considered further. The α -ellipse is highly elongated in the azimuthal direction, with axis ratio $|Y_1/X_1| = \sqrt{2\epsilon}$, and its frequency is small, $\alpha = \sqrt{2\epsilon}\Omega_b$.

These results suggest we consider the approximation in which R_1 , \dot{R}_1 , and $\dot{\varphi}_1$ are small but φ_1 is not. Specifically, if the bar strength Φ_1 is proportional to some small parameter that we may call ϵ , we assume that φ_1 is of order unity, R_1 is of order $\epsilon^{1/2}$, and the time derivative of any quantity is smaller than that quantity by of order $\epsilon^{1/2}$. Let us place the guiding

¹¹ Also called the **inner inner Lindblad resonance**.

Box 3.3: The donkey effect

An orbiting particle that is subject to weak tangential forces can exhibit unusual behavior. To illustrate this, suppose that the particle has mass m and is in a circular orbit of radius r , with angular speed $\dot{\phi} = \Omega(r)$ given by $r\Omega^2(r) = d\Phi/dr$ (eq. 3.7a). Now let us imagine that the particle experiences a small force, F , directed parallel to its velocity vector. Since the force is small, the particle remains on a circular orbit, which slowly changes in radius in response to the force. To determine the rate of change of radius, we note that the angular momentum is $L(r) = mr^2\Omega$ and the torque is $N = rF = \dot{L}$. Thus

$$\dot{r} = \frac{dr}{dL} \dot{L} = \frac{F/m}{2\Omega + r d\Omega/dr} = -\frac{F}{2mB}; \quad (1)$$

where $B(r) = -\Omega - \frac{1}{2}r d\Omega/dr$ is the function defined by equation (3.83). The azimuthal angle accelerates at a rate

$$\ddot{\phi} = \frac{d\Omega}{dr} \dot{r} = -\frac{2A\dot{r}}{r}, \quad (2)$$

where $A(r) = -\frac{1}{2}rd\Omega/dr$ (eq. 3.83). Combining these results,

$$r\ddot{\phi} = \frac{A}{mB}F. \quad (3)$$

This acceleration in azimuthal angle can be contrasted to the acceleration of a free particle under the same force, $\ddot{x} = F/m$. Thus the particle behaves as if it had an inertial mass mB/A , which is negative whenever

$$-2 < \frac{d \ln \Omega}{d \ln R} < 0. \quad (4)$$

Almost all galactic potentials satisfy this inequality. Thus the orbiting particle behaves as if it had negative inertial mass, accelerating in the opposite direction to the applied force.

There are many examples of this phenomenon in galactic dynamics, which has come to be called the **donkey effect**: to quote Lynden-Bell & Kalnajs (1972), who introduced the term, “in azimuth stars behave like donkeys, slowing down when pulled forwards and speeding up when held back.”

The simplest example of the donkey effect is an Earth satellite subjected to atmospheric drag: the satellite sinks gradually into a lower orbit with a larger circular speed and shorter orbital period, so the drag force speeds up the angular passage of the satellite across the sky.

center at L_5 [$\Omega(R_0) = \Omega_b$; $\varphi_0 = \pi/2$] and use equation (3.146b) to write the equations of motion (3.142) as

$$\ddot{R}_1 + (\kappa_0^2 - 4\Omega_0^2) R_1 - 2R_0\Omega_0\dot{\varphi}_1 = -\frac{\partial\Phi_1}{\partial R}, \quad (3.151a)$$

$$\ddot{\varphi}_1 + 2\Omega_0\frac{\dot{R}_1}{R_0} = -\frac{1}{R_0^2}\frac{\partial\Phi_1}{\partial\varphi}. \quad (3.151b)$$

According to our ordering, the terms on the left side of the first line are of order $\epsilon^{3/2}$, $\epsilon^{1/2}$, and $\epsilon^{1/2}$, respectively, while the term on the right side is of order ϵ . All the terms on the second line are of order ϵ . Hence we may simplify the first line by keeping only the terms of order $\epsilon^{1/2}$:

$$(\kappa_0^2 - 4\Omega_0^2) R_1 - 2R_0\Omega_0\dot{\varphi}_1 = 0. \quad (3.152)$$

Substituting equation (3.152) into equation (3.151b) to eliminate R_1 , we find

$$\ddot{\varphi}_1 \left(\frac{\kappa_0^2}{\kappa_0^2 - 4\Omega_0^2} \right) = -\frac{1}{R_0^2} \frac{\partial\Phi_1}{\partial\varphi} \Big|_{(R_0, \varphi_0 + \varphi_1)}. \quad (3.153)$$

Substituting from equation (3.143) for Φ_1 we obtain with $m = 2$

$$\ddot{\varphi}_1 = -\frac{2\Phi_b}{R_0^2} \left(\frac{4\Omega_0^2 - \kappa_0^2}{\kappa_0^2} \right) \sin[2(\varphi_0 + \varphi_1)]. \quad (3.154)$$

By inequality (3.82) we have that $4\Omega_0^2 > \kappa_0^2$. Also we have $\Phi_b < 0$ and $\varphi_0 = \pi/2$, and so equation (3.154) becomes

$$\frac{d^2\psi}{dt^2} = -p^2 \sin\psi, \quad (3.155a)$$

where

$$\psi \equiv 2\varphi_1 \quad \text{and} \quad p^2 \equiv \frac{4}{R_0^2} |\Phi_b(R_0)| \frac{4\Omega_0^2 - \kappa_0^2}{\kappa_0^2}. \quad (3.155b)$$

Equation (3.155a) is simply the equation of a pendulum. Notice that the singularity in R_1 that appeared at corotation in equations (3.148) has disappeared in this more careful analysis. Notice also the interesting fact that the stable equilibrium point of the pendulum, $\varphi_1 = 0$, is at the *maximum*, not the minimum, of the potential Φ_1 (Box 3.3). If the integral of motion

$$E_p = \frac{1}{2}\dot{\psi}^2 - p^2 \cos\psi \quad (3.156)$$

is less than p^2 , the star oscillates slowly or **librates** about the Lagrange point, whereas if $E_p > p^2$, the star is not trapped by the bar but **circulates**

about the center of the galaxy. For small-amplitude librations, the libration frequency is p , consistent with our assumption that the oscillation frequency is of order $\epsilon^{1/2}$ when Φ_b is of order ϵ . Large-amplitude librations of this kind may account for the rings of material often seen in barred galaxies (page 538).

We may obtain the shape of the orbit from equation (3.152) by using equation (3.156) to eliminate $\dot{\varphi}_1 = \frac{1}{2}\dot{\psi}$:

$$R_1 = -\frac{2R_0\Omega_0\dot{\varphi}_1}{4\Omega_0^2 - \kappa_0^2} = \pm \frac{2^{1/2}R_0\Omega_0}{4\Omega_0^2 - \kappa_0^2} \sqrt{E_p + p^2 \cos(2\varphi_1)}. \quad (3.157)$$

We leave as an exercise the demonstration that when $E_p \gg p^2$, equation (3.157) describes the same orbits as are obtained from (3.148a) with $C_1 = 0$ and $\Omega \neq \Omega_b$.

The analysis of this subsection complements the analysis of motion near the Lagrange points in §3.3.2. The earlier analysis is valid for small oscillations around a Lagrange point of an arbitrary two-dimensional rotating potential, while the present analysis is valid for excursions of any amplitude in azimuth around the Lagrange points L_4 and L_5 , but only if the potential is nearly axisymmetric.

3.4 Numerical orbit integration

In most stellar systems, orbits cannot be computed analytically, so effective algorithms for numerical orbit integration are among the most important tools for stellar dynamics. The orbit-integration problems we have to address vary in complexity from following a single particle in a given, smooth galactic potential, to tens of thousands of interacting stars in a globular cluster, to billions of dark-matter particles in a simulation of cosmological clustering. In each of these cases, the dynamics is that of a Hamiltonian system: with N particles there are $3N$ coordinates that form the components of a vector $\mathbf{q}(t)$, and $3N$ components of the corresponding momentum $\mathbf{p}(t)$. These vectors satisfy Hamilton's equations,

$$\dot{\mathbf{q}} = \frac{\partial H}{\partial \mathbf{p}} \quad ; \quad \dot{\mathbf{p}} = -\frac{\partial H}{\partial \mathbf{q}}, \quad (3.158)$$

which can be written as

$$\frac{d\mathbf{w}}{dt} = \mathbf{f}(\mathbf{w}, t), \quad (3.159)$$

where $\mathbf{w} \equiv (\mathbf{q}, \mathbf{p})$ and $\mathbf{f} \equiv (\partial H/\partial \mathbf{p}, -\partial H/\partial \mathbf{q})$. For simplicity we shall assume in this section that the Hamiltonian has the form $H(\mathbf{q}, \mathbf{p}) = \frac{1}{2}p^2 + \Phi(\mathbf{q})$, although many of our results can be applied to more general Hamiltonians. Given a phase-space position \mathbf{w} at time t , and a **timestep** h , we require an

algorithm—an **integrator**—that generates a new position \mathbf{w}' that approximates the true position at time $t' = t + h$. Formally, the problem to be solved is the same whether we are following the motion of a single star in a given potential, or the motion of 10^{10} particles under their mutual gravitational attraction.

The best integrator to use for a given problem is determined by several factors:

- How smooth is the potential? The exploration of orbits in an analytic model of a galaxy potential places fewer demands on the integrator than following orbits in an open cluster, where the stars are buffeted by close encounters with their neighbors.
- How cheaply can we evaluate the gravitational field? At one extreme, evaluating the field by direct summation in simulations of globular cluster with $\gtrsim 10^5$ particles requires $O(N^2)$ operations, and thus is quite expensive compared to the $O(N)$ cost of orbit integrations. At the other extreme, tree codes, spherical-harmonic expansions, or particle-mesh codes require $O(N \ln N)$ operations and thus are comparable in cost to the integration. So the integrator used in an N-body simulation of a star cluster should make the best possible use of each expensive but accurate force evaluation, while in a cosmological simulation it is better to use a simple integrator and evaluate the field more frequently.
- How much memory is available? The most accurate integrators use the position and velocity of a particle at several previous timesteps to help predict its future position. When simulating a star cluster, the number of particles is small enough ($N \lesssim 10^5$) that plenty of memory should be available to store this information. In a simulation of galaxy dynamics or a cosmological simulation, however, it is important to use as many particles as possible, so memory is an important constraint. Thus for such simulations the optimal integrator predicts the future phase-space position using only the current position and gravitational field.
- How long will the integration run? The answer can range from a few crossing times for the simulation of a galaxy merger to 10^5 crossing times in the core of a globular cluster. Long integrations require that the integrator does not introduce any systematic drift in the energy or other integrals of motion.

Useful references include Press et al. (1986), Hairer, Lubich, & Wanner (2002), and Aarseth (2003).

3.4.1 Symplectic integrators

(a) Modified Euler integrator Let us replace the original Hamiltonian $H(\mathbf{q}, \mathbf{p}) = \frac{1}{2}p^2 + \Phi(\mathbf{q})$ by the time-dependent Hamiltonian

$$H_h(\mathbf{q}, \mathbf{p}, t) = \frac{1}{2}p^2 + \Phi(\mathbf{q})\delta_h(t), \quad \text{where} \quad \delta_h(t) \equiv h \sum_{j=-\infty}^{\infty} \delta(t - jh) \quad (3.160)$$

is an infinite series of delta functions (Appendix C.1). Averaged over a time interval that is long compared to h , $\langle H_h \rangle \simeq H$, so the trajectories determined by H_h should approach those determined by H as $h \rightarrow 0$.

Hamilton's equations for H_h read

$$\dot{\mathbf{q}} = \frac{\partial H_h}{\partial \mathbf{p}} = \mathbf{p} \quad ; \quad \dot{\mathbf{p}} = -\frac{\partial H_h}{\partial \mathbf{q}} = -\nabla\Phi(\mathbf{q})\delta_h(t). \quad (3.161)$$

We now integrate these equations from $t = -\epsilon$ to $t = h - \epsilon$, where $0 < \epsilon \ll h$. Let the system have coordinates (\mathbf{q}, \mathbf{p}) at time $t = -\epsilon$, and first ask for its coordinates $(\bar{\mathbf{q}}, \bar{\mathbf{p}})$ at $t = +\epsilon$. During this short interval \mathbf{q} changes by a negligible amount, and \mathbf{p} suffers a kick governed by the second of equations (3.161). Integrating this equation from $t = -\epsilon$ to ϵ is trivial since \mathbf{q} is fixed, and we find

$$\bar{\mathbf{q}} = \mathbf{q} \quad ; \quad \bar{\mathbf{p}} = \mathbf{p} - h\nabla\Phi(\mathbf{q}); \quad (3.162a)$$

this is called a **kick step** because the momentum changes but the position does not. Next, between $t = +\epsilon$ and $t = h - \epsilon$, the value of the delta function is zero, so the system has constant momentum, and Hamilton's equations yield for the coordinates at $t = h - \epsilon$

$$\mathbf{q}' = \bar{\mathbf{q}} + h\bar{\mathbf{p}} \quad ; \quad \mathbf{p}' = \bar{\mathbf{p}}; \quad (3.162b)$$

this is called a **drift step** because the position changes but the momentum does not. Combining these results, we find that over a timestep h starting at $t = -\epsilon$ the Hamiltonian H_h generates a map $(\mathbf{q}, \mathbf{p}) \rightarrow (\mathbf{q}', \mathbf{p}')$ given by

$$\mathbf{p}' = \mathbf{p} - h\nabla\Phi(\mathbf{q}) \quad ; \quad \mathbf{q}' = \mathbf{q} + h\mathbf{p}'. \quad (3.163a)$$

Similarly, starting at $t = +\epsilon$ yields the map

$$\mathbf{q}' = \mathbf{q} + h\mathbf{p} \quad ; \quad \mathbf{p}' = \mathbf{p} - h\nabla\Phi(\mathbf{q}'). \quad (3.163b)$$

These maps define the “kick-drift” or “drift-kick” **modified Euler integrator**. The performance of this integrator in a simple galactic potential is shown in Figure 3.21.

The map induced by any Hamiltonian is a canonical or symplectic map (page 803), so it can be derived from a generating function. It is simple to confirm using equations (D.93) that the generating function $S(\mathbf{q}, \mathbf{p}') = \mathbf{q} \cdot \mathbf{p}' + \frac{1}{2}h\mathbf{p}'^2 + h\Phi(\mathbf{q})$ yields the kick-drift modified Euler integrator (3.163a).

According to the modified Euler integrator, the position after timestep h is

$$\mathbf{q}' = \mathbf{q} + h\mathbf{p}' = \mathbf{q} + h\mathbf{p} - h^2\nabla\Phi(\mathbf{q}), \quad (3.164)$$

while the exact result may be written as a Taylor series,

$$\mathbf{q}' = \mathbf{q} + h\dot{\mathbf{q}}(t=0) + \frac{1}{2}h^2\ddot{\mathbf{q}}(t=0) + O(h^3) = \mathbf{q} + h\mathbf{p} - \frac{1}{2}h^2\nabla\Phi(\mathbf{q}) + O(h^3). \quad (3.165)$$

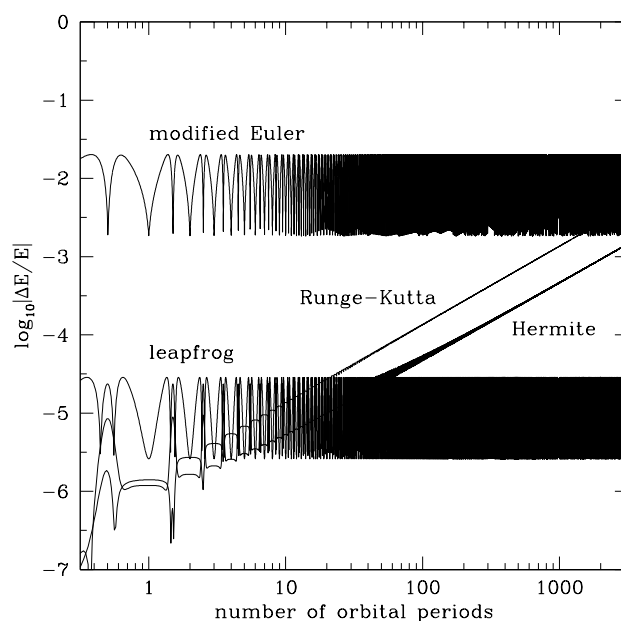


Figure 3.21 Fractional energy error as a function of time for several integrators, following a particle orbiting in the logarithmic potential $\Phi(r) = \ln r$. The orbit is moderately eccentric (apocenter twice as big as pericenter). The timesteps are fixed, and chosen so that there are 300 evaluations of the force or its derivatives per period for all of the integrators. The integrators shown are kick-drift modified Euler (3.163a), leapfrog (3.166a), Runge-Kutta (3.168), and Hermite (3.172a–d). Note that (i) over moderate time intervals, the errors are smallest for the fourth-order integrators (Runge-Kutta and Hermite), intermediate for the second-order integrator (leapfrog), and largest for the first-order integrator (modified Euler); (ii) the energy error of the symplectic integrators does not grow with time.

The error after a single step of the modified Euler integrator is seen to be $O(h^2)$, so it is said to be a **first-order** integrator.

Since the mappings (3.163) are derived from the Hamiltonian (3.160), they are symplectic, so either flavor of the modified Euler integrator is a **symplectic integrator**. Symplectic integrators conserve phase-space volume and Poincaré invariants (Appendix D.4.2). Consequently, if the integrator is used to advance a series of particles that initially lie on a closed curve in the (q_i, p_i) phase plane, the curve onto which it moves the particles has the same line integral $\oint p_i dq_i$ around it as the original curve. This conservation property turns out to constrain the allowed motions in phase space so strongly that the usual tendency of numerical orbit integrations to drift in energy (sometimes called **numerical dissipation**, even through the energy can either decay or grow) is absent in symplectic integrators (Hairer,

Lubich, & Wanner 2002).

Leapfrog integrator By alternating kick and drift steps in more elaborate sequences, we can construct higher-order integrators (Yoshida 1993); these are automatically symplectic since they are the composition of maps (the kick and drift steps) that are symplectic. The simplest and most widely used of these is the **leapfrog** or **Verlet** integrator in which we drift for $\frac{1}{2}h$, kick for h and then drift for $\frac{1}{2}h$:

$$\mathbf{q}_{1/2} = \mathbf{q} + \frac{1}{2}h\mathbf{p} ; \mathbf{p}' = \mathbf{p} - h\nabla\Phi(\mathbf{q}_{1/2}) ; \mathbf{q}' = \mathbf{q} + \frac{1}{2}h\mathbf{p}' . \quad (3.166a)$$

This algorithm is sometimes called “drift-kick-drift” leapfrog; an equally good form is “kick-drift-kick” leapfrog:

$$\mathbf{p}_{1/2} = \mathbf{p} - \frac{1}{2}h\nabla\Phi(\mathbf{q}) ; \mathbf{q}' = \mathbf{q} + h\mathbf{p}_{1/2} ; \mathbf{p}' = \mathbf{p} - \frac{1}{2}h\nabla\Phi(\mathbf{q}') . \quad (3.166b)$$

Drift-kick-drift leapfrog can also be derived by considering motion in the Hamiltonian (3.160) from $t = -\frac{1}{2}h$ to $t = \frac{1}{2}h$.

The leapfrog integrator has many appealing features: (i) In contrast to the modified Euler integrator, it is second- rather than first-order accurate, in that the error in phase-space position after a single timestep is $O(h^3)$ (Problem 3.26). (ii) Leapfrog is **time reversible** in the sense that if leapfrog advances the system from (\mathbf{q}, \mathbf{p}) to $(\mathbf{q}', \mathbf{p}')$ in a given time, it will also advance it from $(\mathbf{q}', -\mathbf{p}')$ to $(\mathbf{q}, -\mathbf{p})$ in the same time. Time-reversibility is a constraint on the phase-space flow that, like symplecticity, suppresses numerical dissipation, since dissipation is not a time-reversible phenomenon (Roberts & Quispel 1992; Hairer, Lubich, & Wanner 2002). (iii) A sequence of n leapfrog steps can be regarded as a drift step for $\frac{1}{2}h$, then n kick-drift steps of the modified Euler integrator, then a drift step for $-\frac{1}{2}h$; thus if $n \gg 1$ the leapfrog integrator requires negligibly more work than the same number of steps of the modified Euler integrator. (iv) Leapfrog also needs no storage of previous timesteps, so is economical of memory.

Because of all these advantages, most codes for simulating collisionless stellar systems use the leapfrog integrator. Time-reversible, symplectic integrators of fourth and higher orders, derived by combining multiple kick and drift steps, are described in Problem 3.27 and Yoshida (1993).

One serious limitation of symplectic integrators is that they work well only with fixed timesteps, for the following reason. Consider an integrator with fixed timestep h that maps phase-space coordinates \mathbf{w} to $\mathbf{w}' = \mathbf{W}(\mathbf{w}, h)$. The integrator is symplectic if the function \mathbf{W} satisfies the symplectic condition (D.78), which involves the Jacobian matrix $g_{\alpha\beta} = \partial W_{\alpha} / \partial w_{\beta}$. Now suppose that the timestep is varied, by choosing it to be some function $h(\mathbf{w})$ of location in phase space, so $\mathbf{w}' = \mathbf{W}[\mathbf{w}, h(\mathbf{w})] \equiv \widetilde{\mathbf{W}}(\mathbf{w})$. The Jacobian matrix of $\widetilde{\mathbf{W}}$ is not equal to the Jacobian matrix of \mathbf{W} , and in general will not satisfy the symplectic condition; in words, a symplectic integrator with fixed timestep is generally no longer symplectic once the timestep is varied.

Fortunately, the geometric constraints on phase-space flow imposed by time-reversibility are also strong, so the leapfrog integrator retains its good behavior if the timestep is adjusted in a time-reversible manner, even though the resulting integrator is no longer symplectic. Here is one way to do this: suppose that the appropriate timestep h is given by some function $\tau(\mathbf{w})$ of the phase-space coordinates. Then we modify equations (3.166a) to

$$\begin{aligned} \mathbf{q}_{1/2} &= \mathbf{q} + \frac{1}{2}h\mathbf{p} \quad ; \quad \mathbf{p}_{1/2} = \mathbf{p} - \frac{1}{2}h\nabla\Phi(\mathbf{q}_{1/2}), \\ t' &= t + \frac{1}{2}(h + h'), \\ \mathbf{p}' &= \mathbf{p}_{1/2} - \frac{1}{2}h'\nabla\Phi(\mathbf{q}_{1/2}) \quad ; \quad \mathbf{q}' = \mathbf{q}_{1/2} + \frac{1}{2}h'\mathbf{p}'. \end{aligned} \quad (3.167)$$

Here h' is determined from h by solving the equation $u(h, h') = \tau(\mathbf{q}_{1/2}, \mathbf{p}_{1/2})$, where $\tau(\mathbf{q}, \mathbf{p})$ is the desired timestep at (\mathbf{q}, \mathbf{p}) and $u(h, h')$ is any symmetric function of h and h' such that $u(h, h) = h$; for example, $u(h, h') = \frac{1}{2}(h + h')$ or $u(h, h') = 2hh'/(h + h')$.

3.4.2 Runge–Kutta and Bulirsch–Stoer integrators

To follow the motion of particles in a given smooth gravitational potential $\Phi(\mathbf{q})$ for up to a few hundred crossing times, the fourth-order Runge–Kutta integrator provides reliable transportation. The algorithm is

$$\begin{aligned} \mathbf{k}_1 &= h\mathbf{f}(\mathbf{w}, t) \quad ; \quad \mathbf{k}_2 = h\mathbf{f}(\mathbf{w} + \frac{1}{2}\mathbf{k}_1, t + \frac{1}{2}h), \\ \mathbf{k}_3 &= h\mathbf{f}(\mathbf{w} + \frac{1}{2}\mathbf{k}_2, t + \frac{1}{2}h) \quad ; \quad \mathbf{k}_4 = h\mathbf{f}(\mathbf{w} + \mathbf{k}_3, t + h), \\ \mathbf{w}' &= \mathbf{w} + \frac{1}{6}(\mathbf{k}_1 + 2\mathbf{k}_2 + 2\mathbf{k}_3 + \mathbf{k}_4) \quad ; \quad t' = t + h. \end{aligned} \quad (3.168)$$

The Runge–Kutta integrator is neither symplectic nor reversible, and it requires considerably more memory than the leapfrog integrator because memory has to be allocated to $\mathbf{k}_1, \dots, \mathbf{k}_4$. However, it is easy to use and provides fourth-order accuracy.

The **Bulirsch–Stoer** integrator is used for the same purposes as the Runge–Kutta integrator; although more complicated to code, it often surpasses the Runge–Kutta integrator in performance. The idea behind this integrator is to estimate $\mathbf{w}(t + h)$ from $\mathbf{w}(t)$ using first one step of length h , then two steps of length $h/2$, then four steps of length $h/4$, etc., up to 2^K steps of length $h/2^K$ for some predetermined number K . Then one extrapolates this sequence of results to the coordinates that would be obtained in the limit $K \rightarrow \infty$. Like the Runge–Kutta integrator, this integrator achieves speed and accuracy at the cost of the memory required to hold intermediate results. Like all high-order integrators, the Runge–Kutta and Bulirsch–Stoer integrators work best when following motion in smooth gravitational fields.

3.4.3 Multistep predictor-corrector integrators

We now discuss more complex integrators that are widely used in simulations of star clusters. We have a trajectory that has arrived at some phase-space position \mathbf{w}_0 at time t_0 , and we wish to predict its position \mathbf{w}_1 at t_1 . The general idea is to assume that the trajectory $\mathbf{w}(t)$ is a polynomial function of time $\mathbf{w}^{\text{poly}}(t)$, called the **interpolating polynomial**. The interpolating polynomial is determined by fitting to some combination of the present position \mathbf{w}_0 , the past positions, $\mathbf{w}_{-1}, \mathbf{w}_{-2}, \dots$ at times t_{-1}, t_{-2}, \dots , and the present and past phase-space velocities, which are known through $\dot{\mathbf{w}}_j = \mathbf{f}(\mathbf{w}_j, t_j)$. There is no requirement that \mathbf{f} is derived from Hamilton's equations, so these methods can be applied to any first-order differential equations; on the other hand they are not symplectic.

If the interpolating polynomial has order k , then the error after a small time interval h is given by the first term in the Taylor series for $\mathbf{w}(t)$ not represented in the polynomial, which is $O(h^{k+1})$. Thus the order of the integrator is k .¹²

The **Adams–Bashforth** multistep integrator takes \mathbf{w}^{poly} to be the unique k th-order polynomial that passes through \mathbf{w}_0 at t_0 and through the k points $(t_{-k+1}, \dot{\mathbf{w}}_{-k+1}), \dots, (t_0, \dot{\mathbf{w}}_0)$.

Explicit formulae for the Adams–Bashforth integrators are easy to find by computer algebra; however, the formulae are too cumbersome to write here except in the special case of equal timesteps, $t_{j+1} - t_j = h$ for all j . Then the first few Adams–Bashforth integrators are

$$\mathbf{w}_1 = \mathbf{w}_0 + h \begin{cases} \dot{\mathbf{w}}_0 & (k = 1) \\ \frac{3}{2}\dot{\mathbf{w}}_0 - \frac{1}{2}\dot{\mathbf{w}}_{-1} & (k = 2) \\ \frac{23}{12}\dot{\mathbf{w}}_0 - \frac{4}{3}\dot{\mathbf{w}}_{-1} + \frac{5}{12}\dot{\mathbf{w}}_{-2} & (k = 3) \\ \frac{55}{24}\dot{\mathbf{w}}_0 - \frac{59}{24}\dot{\mathbf{w}}_{-1} + \frac{37}{24}\dot{\mathbf{w}}_{-2} - \frac{3}{8}\dot{\mathbf{w}}_{-3} & (k = 4). \end{cases} \quad (3.169)$$

The case $k = 1$ is called **Euler's integrator**, and usually works rather badly.

The **Adams–Moulton** integrator differs from Adams–Bashforth only in that it computes the interpolating polynomial from the position \mathbf{w}_0 and the phase-space velocities $\dot{\mathbf{w}}_{-k+2}, \dots, \dot{\mathbf{w}}_1$. For equal timesteps, the first few Adams–Moulton integrators are

$$\mathbf{w}_1 = \mathbf{w}_0 + h \begin{cases} \dot{\mathbf{w}}_1 & (k = 1) \\ \frac{1}{2}\dot{\mathbf{w}}_1 + \frac{1}{2}\dot{\mathbf{w}}_0 & (k = 2) \\ \frac{5}{12}\dot{\mathbf{w}}_1 + \frac{2}{3}\dot{\mathbf{w}}_0 - \frac{1}{12}\dot{\mathbf{w}}_{-1} & (k = 3) \\ \frac{3}{8}\dot{\mathbf{w}}_1 + \frac{19}{24}\dot{\mathbf{w}}_0 - \frac{5}{24}\dot{\mathbf{w}}_{-1} + \frac{1}{24}\dot{\mathbf{w}}_{-2} & (k = 4). \end{cases} \quad (3.170)$$

¹² Unfortunately, the term “order” is used both for the highest power retained in the Taylor series for $\mathbf{w}(t)$, t^k , and the dependence of the one-step error on the timestep, h^{k+1} ; fortunately, both orders are the same.

Since $\dot{\mathbf{w}}_1$ is determined by the unknown phase-space position \mathbf{w}_1 through $\dot{\mathbf{w}}_1 = \mathbf{f}(\mathbf{w}_1, t_1)$, equations (3.170) are nonlinear equations for \mathbf{w}_1 that must be solved iteratively. The Adams–Moulton integrator is therefore said to be **implicit**, in contrast to Adams–Bashforth, which is **explicit**.

The strength of the Adams–Moulton integrator is that it determines \mathbf{w}_1 by *interpolating* the phase-space velocities, rather than by extrapolating them, as with Adams–Bashforth. This feature makes it a more reliable and stable integrator; the cost is that a nonlinear equation must be solved at every timestep.

In practice the Adams–Bashforth and Adams–Moulton integrators are used together as a **predictor-corrector** integrator. Adams–Bashforth is used to generate a preliminary value \mathbf{w}_1 (the prediction or P step), which is then used to generate $\dot{\mathbf{w}}_1 = \mathbf{f}(\mathbf{w}_1, t_1)$ (the evaluation or E step), which is used in the Adams–Moulton integrator (the corrector or C step). This three-step sequence is abbreviated as PEC. In principle one can then iterate the Adams–Moulton integrator to convergence through the sequence PECEC \cdots ; however, this is not cost-effective, since the Adams–Moulton formula, even if solved exactly, is only an approximate representation of the differential equation we are trying to solve. Thus one usually stops with PEC (stop the iteration after evaluating \mathbf{w}_1 twice) or PECE (stop the iteration after evaluating $\dot{\mathbf{w}}_1$ twice).

When these methods are used in orbit integrations, the equations of motion usually have the form $\dot{\mathbf{x}} = \mathbf{v}$, $\dot{\mathbf{v}} = -\nabla\Phi(\mathbf{x}, t)$. In this case it is best to apply the integrator only to the second equation, and to generate the new position \mathbf{x}_1 by analytically integrating the interpolating polynomial for $\mathbf{v}(t)$ —this gives a formula for \mathbf{x}_1 that is more accurate by one power of h .

Analytic estimates (Makino 1991) suggest that the one-step error in the Adams–Bashforth–Moulton predictor-corrector integrator is smaller than the error in the Adams–Bashforth integrator by a factor of 5 for $k = 2$, 9 for $k = 3$, 13 for $k = 4$, etc. These analytic results, or the difference between the predicted and corrected values of \mathbf{w}_1 , can be used to determine the longest timestep that is compatible with a prescribed target accuracy—see §3.4.5.

Because multistep integrators require information from the present time and $k - 1$ past times, a separate startup integrator, such as Runge–Kutta, must be used to generate the first $k - 1$ timesteps. Multistep integrators are not economical of memory because they store the coefficients of the entire interpolating polynomial rather than just the present phase-space position.

3.4.4 Multivalued integrators

By differentiating the equations of motion $\dot{\mathbf{w}} = \mathbf{f}(\mathbf{w})$ with respect to time, we obtain an expression for $\ddot{\mathbf{w}}$, which involves second derivatives of the potential, $\partial^2\Phi/\partial q_i\partial q_j$. If our Poisson solver delivers reliable values for these second derivatives, it can be advantageous to use $\ddot{\mathbf{w}}$ or even higher time derivatives of \mathbf{w} to determine the interpolating polynomial $\mathbf{w}^{\text{poly}}(t)$. Algorithms that

employ the second and higher derivatives of \mathbf{w} are called **multivalued integrators**.

In the simplest case we set $\mathbf{w}^{\text{poly}}(t)$ to the k th-order polynomial that matches \mathbf{w} and its first k time derivatives at t_0 ; this provides $k+1$ constraints for the $k+1$ polynomial coefficients and corresponds to predicting $\mathbf{w}(t)$ by its Taylor series expansion around t_0 . A more satisfactory approach is to determine $\mathbf{w}^{\text{poly}}(t)$ from the values taken by \mathbf{w} , $\dot{\mathbf{w}}$, $\ddot{\mathbf{w}}$, etc., at both t_0 and t_1 . Specifically, for even k only, we make $\mathbf{w}^{\text{poly}}(t)$ the k th-order polynomial that matches \mathbf{w} at t_0 and its first $\frac{1}{2}k$ time derivatives at both t_0 and t_1 —once again this provides $1 + 2 \times \frac{1}{2}k = k + 1$ constraints and hence determines the $k + 1$ coefficients of the interpolating polynomial. The first few integrators of this type are

$$\mathbf{w}_1 = \mathbf{w}_0 + \begin{cases} \frac{1}{2}h(\dot{\mathbf{w}}_0 + \dot{\mathbf{w}}_1) & (k = 2) \\ \frac{1}{2}h(\dot{\mathbf{w}}_0 + \dot{\mathbf{w}}_1) + \frac{1}{12}h^2(\ddot{\mathbf{w}}_0 - \ddot{\mathbf{w}}_1) & (k = 4) \\ \frac{1}{2}h(\dot{\mathbf{w}}_0 + \dot{\mathbf{w}}_1) + \frac{1}{10}h^2(\ddot{\mathbf{w}}_0 - \ddot{\mathbf{w}}_1) \\ \quad + \frac{1}{120}h^3(\dddot{\mathbf{w}}_0 + \dddot{\mathbf{w}}_1) & (k = 6). \end{cases} \quad (3.171)$$

Like the Adams–Moulton integrator, all of these integrators are implicit, and in fact the first of these formulae is the same as the second-order Adams–Moulton integrator in equation (3.170). Because these integrators employ information from only t_0 and t_1 , there are two significant simplifications compared to multistep integrators: no separate startup procedure is needed, and the formulae look the same even if the timestep is variable.

Multivalued integrators are sometimes called **Obreshkov** (or Obrechhoff) or **Hermite** integrators, the latter name arising because they are based on Hermite interpolation, which finds a polynomial that fits specified values of a function and its derivatives (Butcher 1987).

Makino & Aarseth (1992) and Makino (2001) recommend a fourth-order multivalued predictor-corrector integrator for star-cluster simulations. Their predictor is a single-step, second-order multivalued integrator, that is, a Taylor series including terms of order h^2 . Writing $d\mathbf{v}/dt = \mathbf{g}$, where \mathbf{g} is the gravitational field, their predicted velocity is

$$\mathbf{v}_{p,1} = \mathbf{v}_0 + h\mathbf{g}_0 + \frac{1}{2}h^2\dot{\mathbf{g}}_0. \quad (3.172a)$$

The predicted position is obtained by analytically integrating the interpolating polynomial for \mathbf{v} ,

$$\mathbf{x}_{p,1} = \mathbf{x}_0 + h\mathbf{v}_0 + \frac{1}{2}h^2\mathbf{g}_0 + \frac{1}{6}h^3\dot{\mathbf{g}}_0. \quad (3.172b)$$

The predicted position and velocity are used to compute the gravitational field and its time derivative at time t_1 , \mathbf{g}_1 and $\dot{\mathbf{g}}_1$. These are used to correct the velocity using the fourth-order formula (3.171):

$$\mathbf{v}_1 = \mathbf{v}_0 + \frac{1}{2}h(\mathbf{g}_0 + \mathbf{g}_1) + \frac{1}{12}h^2(\dot{\mathbf{g}}_0 - \dot{\mathbf{g}}_1); \quad (3.172c)$$

in words, \mathbf{v}_1 is determined by the fourth-order interpolating polynomial $\mathbf{v}^{\text{poly}}(t)$ that satisfies the five constraints $\mathbf{v}^{\text{poly}}(t_0) = \mathbf{v}_0$, $\dot{\mathbf{v}}^{\text{poly}}(t_i) = \mathbf{g}_i$, $\ddot{\mathbf{v}}^{\text{poly}}(t_i) = \dot{\mathbf{g}}_i$ for $i = 0, 1$.

To compute the corrected position, the most accurate procedure is to integrate analytically the interpolating polynomial for \mathbf{v} , which yields:

$$\mathbf{x}_1 = \mathbf{x}_0 + h\mathbf{v}_0 + \frac{1}{20}h^2(7\mathbf{g}_0 + 3\mathbf{g}_1) + \frac{1}{60}h^3(3\dot{\mathbf{g}}_0 - 2\dot{\mathbf{g}}_1). \quad (3.172d)$$

The performance of this integrator, often simply called the Hermite integrator, is illustrated in Figure 3.21.

3.4.5 Adaptive timesteps

Except for the simplest problems, any integrator should have an **adaptive timestep**, that is, an automatic procedure that continually adjusts the timestep to achieve some target level of accuracy. Choosing the right timestep is one of the most challenging tasks in designing a numerical integration scheme. Many sophisticated procedures are described in publicly available integration packages and numerical analysis textbooks. Here we outline a simple approach.

Let us assume that our goal is that the error in \mathbf{w} after some short time τ should be less than $\epsilon|\mathbf{w}_0|$, where $\epsilon \ll 1$ and \mathbf{w}_0 is some reference phase-space position. We first move from \mathbf{w} to \mathbf{w}_2 by taking two timesteps of length $h \ll \tau$. Then we return to \mathbf{w} and take one step of length $2h$ to reach \mathbf{w}_1 . Suppose that the correct position after an interval $2h$ is \mathbf{w}' , and that our integrator has order k . Then the errors in \mathbf{w}_1 and \mathbf{w}_2 may be written

$$\mathbf{w}_1 - \mathbf{w}' \simeq (2h)^{k+1}\mathbf{E} \quad ; \quad \mathbf{w}_2 - \mathbf{w}' \simeq 2h^{k+1}\mathbf{E}, \quad (3.173)$$

where \mathbf{E} is an unknown error vector. Subtracting these equations to eliminate \mathbf{w}' , we find $\mathbf{E} \simeq (\mathbf{w}_1 - \mathbf{w}_2)/[2(2^k - 1)h^{k+1}]$. Now if we advance for a time τ , using $n \equiv \tau/h'$ timesteps of length h' , the error will be

$$\Delta = nh'^{k+1}\mathbf{E} = (\mathbf{w}_1 - \mathbf{w}_2) \frac{\tau h'^k}{2(2^k - 1)h^{k+1}}. \quad (3.174)$$

Our goal that $|\Delta| \lesssim \epsilon|\mathbf{w}_0|$ will be satisfied if

$$h' < h_{\text{max}} \equiv \left(2(2^k - 1) \frac{h}{\tau} \frac{\epsilon|\mathbf{w}_0|}{|\mathbf{w}_1 - \mathbf{w}_2|} \right)^{1/k} h. \quad (3.175)$$

If we are using a predictor-corrector scheme, a similar analysis can be used to deduce h_{max} from the difference of the phase-space positions returned by the predictor and the corrector, without repeating the entire predictor-corrector sequence.

3.4.6 Individual timesteps

The density in many stellar systems varies by several orders of magnitude between the center and the outer parts, and as a result the crossing time of orbits near the center is much smaller than the crossing time in the outer envelope. For example, in a typical globular cluster the crossing time at the center is $\lesssim 1$ Myr, while the crossing time near the tidal radius is ~ 100 Myr. Consequently, the timestep that can be safely used to integrate the orbits of stars is much smaller at the center than the edge. It is extremely inefficient to integrate *all* of the cluster stars with the shortest timestep needed for *any* star, so integrators must allow individual timesteps for each star.

If the integrator employs an interpolating polynomial, the introduction of individual timesteps is in principle fairly straightforward. To advance a given particle, one uses the most recent interpolating polynomials of all the other particles to predict their locations at whatever times the integrator requires, and then evaluates the forces between the given particle and the other particles.

This procedure makes sense if the Poisson solver uses direct summation (§2.9.1). However, with other Poisson solvers there is a much more efficient approach. Suppose, for example, that we are using a tree code (§2.9.2). Then before a single force can be evaluated, *all* particles have to be sorted into a tree. Once that has been done, it is comparatively inexpensive to evaluate large numbers of forces; hence to minimize the computational work done by the Poisson solver, it is important to evaluate the forces on many particles simultaneously. A **block timestep** scheme makes this possible whilst allowing different timesteps for different particles, by quantizing the timesteps. We now describe how one version of this scheme works with the leapfrog integrator.

We assign each particle to one of $K + 1$ classes, such that particles in class k are to be advanced with timestep $h_k \equiv 2^k h$ for $k = 0, 1, 2, \dots, K$. Thus h is the shortest timestep (class 0) and $2^K h$ is the longest (class K). The Poisson solver is used to evaluate the gravitational field at the initial time t_0 , and each particle is kicked by the impulse $-\frac{1}{2}h_k \nabla\Phi$, corresponding to the first part of the kick-drift-kick leapfrog step (3.166b). In Figure 3.22 the filled semicircles on the left edge of the diagram symbolize these kicks; they are larger at the top of the diagram to indicate that the strength of the kicks increases as 2^k . Then every particle is drifted through time h , and the Poisson solver is used only to find the forces on the particles in class 0, so these particles can be kicked by $-h \nabla\Phi$, which is the sum of the kicks at the end of their first leapfrog step and the start of their second.

Next we drift all particles through h a second time, and use the Poisson solver to find the forces on the particles in both class 0 and class 1. The particles of class 0 are kicked by $-h \nabla\Phi$, and the particles of class 1 are kicked by $-h_1 \nabla\Phi = -2h \nabla\Phi$. After an interval $3h$ the particles in class 0 are kicked, after $4h$ the particles in classes 0, 1 and 2 are kicked, etc. This

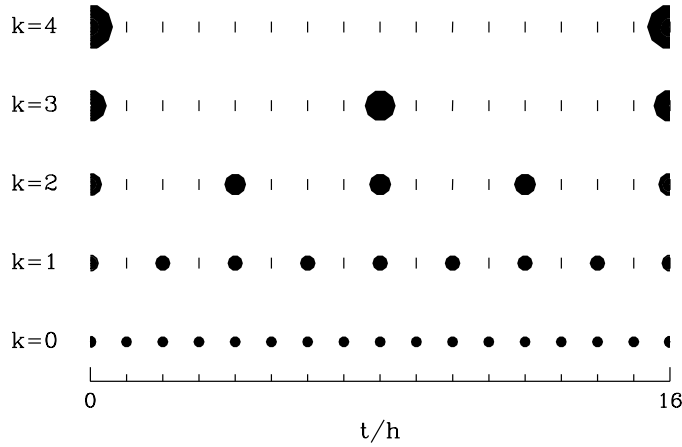


Figure 3.22 Schematic of the block timestep scheme, for a system with 5 classes of particles, having timestep h (class $k = 0$), $2h, \dots, 16h$ (class $k = K = 4$). The particles are integrated for a total time of $16h$. Each filled circle or half-circle marks the time at which particles in a given class are kicked. Each vertical bar marks a time at which particles in a class are paused in their drift step, without being kicked, in order to calculate their contribution to the kick given to particles in lower classes. The kicks at the start and end of the integration, $t = 0$ and $t = 16h$, are half as strong as the other kicks, and so are denoted by half-circles.

process continues until all particles are due for a kick, after a time $h_K = 2^K h$. The final kick for particles in class k is $-\frac{1}{2}h_k \nabla \Phi$, which completes 2^{K-k} leapfrog steps for each particle. At this point it is prudent to reconsider how the particles are assigned to classes in case some need smaller or larger timesteps.

A slightly different block timestep scheme works well with a particle-mesh Poisson solver (§2.9.3) when parts of the computational domain are covered by finer meshes than others, with each level of refinement being by a factor of two in the number of mesh points per unit length (Knebe, Green, & Binney 2001). Then particles are assigned timesteps according to the fineness of the mesh they are in: particles in the finest mesh have timestep $\Delta t = h$, while particles in the next coarser mesh have $\Delta t = 2h$, and so on. Particles on the finest mesh are drifted through time $\frac{1}{2}h$ before the density is determined on this mesh, and the Poisson solver is invoked to determine the forces on this mesh. Then the particles on this mesh are kicked through time h and drifted through time $\frac{1}{2}h$. Then the same drift-kick-drift sequence is used to advance particles on the next coarser mesh through time $2h$. Now these particles are ahead in time of the particles on the finest mesh. This situation is remedied by again advancing the particles on the finest mesh by h with the drift-kick-drift sequence. Once the particles on the two finest

meshes have been advanced through time $2h$, we are ready to advance by $\Delta t = 4h$ the particles that are the next coarser mesh, followed by a repeat of the operations that were used to advance the particles on the two finest meshes by $2h$. The key point about this algorithm is that at each level k , particles are first advanced ahead of particles on the next coarser mesh, and then the latter particles jump ahead of the particles on level k so the next time the particles on level k are advanced, they are catching up with the particles of the coarser mesh. Errors arising from moving particles in a gravitational field from the surroundings that is out-of-date are substantially canceled by errors arising from moving particles in an ambient field that has run ahead of itself.

3.4.7 Regularization

In any simulation of a star cluster, sooner or later two particles will suffer an encounter having a very small impact parameter. In the limiting case in which the impact parameter is exactly zero (a **collision orbit**), the equation of motion for the distance r between the two particles is (eq. D.33)

$$\ddot{r} = -GM/r^2, \quad (3.176)$$

where M is the sum of the masses of the two particles. This equation is singular at $r = 0$, and a conscientious integrator will attempt to deal with the singularity by taking smaller and smaller timesteps as r diminishes, thereby bringing the entire N-body integration grinding to a halt. Even in a near-collision orbit, the integration through pericenter will be painfully slow. This problem is circumvented by transforming to a coordinate system in which the two-body problem has no singularity—this procedure is called **regularization** (Stiefel & Schiefel 1971; Mikkola 1997; Heggie & Hut 2003; Aarseth 2003). Standard integrators can then be used to solve the equations of motion in the regularized coordinates.

(a) Burdet–Heggie regularization The simplest approach to regularization is time transformation. We write the equations of motion for the two-body problem as

$$\ddot{\mathbf{r}} = -GM \frac{\mathbf{r}}{r^3} + \mathbf{g}, \quad (3.177)$$

where \mathbf{g} is the gravitational field from the other $N - 2$ bodies in the simulation, and change to a fictitious time τ that is defined by

$$dt = r d\tau. \quad (3.178)$$

Denoting derivatives with respect to τ by a prime we find

$$\dot{\mathbf{r}} = \frac{d\tau}{dt} \frac{d\mathbf{r}}{d\tau} = \frac{1}{r} \mathbf{r}' \quad ; \quad \ddot{\mathbf{r}} = \frac{d\tau}{dt} \frac{d}{d\tau} \frac{1}{r} \mathbf{r}' = \frac{1}{r^2} \mathbf{r}'' - \frac{r'}{r^3} \mathbf{r}'. \quad (3.179)$$

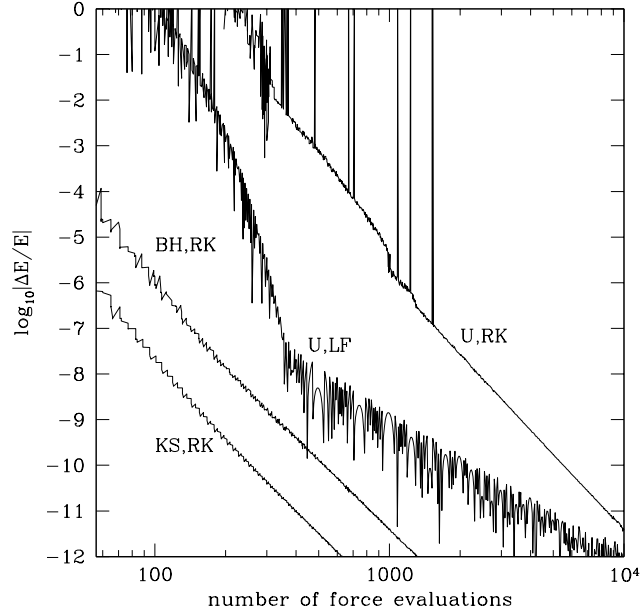


Figure 3.23 Fractional energy error from integrating one pericenter passage of a highly eccentric orbit in a Keplerian potential, as a function of the number of force evaluations. The orbit has semi-major axis $a = 1$ and eccentricity $e = 0.99$, and is integrated from $r = 1, \dot{r} < 0$ to $r = 1, \dot{r} > 0$. Curves labeled by “RK” are followed using a fourth-order Runge–Kutta integrator (3.168) with adaptive timestep control as described by Press et al. (1986). The curve labeled “U” for “unregularized” is integrated in Cartesian coordinates, the curve “BH” uses Burdet–Heggie regularization, and the curve “KS” uses Kustaanheimo–Stiefel regularization. The curve labeled “U,LF” is followed in Cartesian coordinates using a leapfrog integrator with timestep proportional to radius (eq. 3.167). The horizontal axis is the number of force evaluations used in the integration.

Substituting these results into the equation of motion, we obtain

$$\mathbf{r}'' = \frac{r'}{r} \mathbf{r}' - GM \frac{\mathbf{r}}{r} + r^2 \mathbf{g}. \quad (3.180)$$

The eccentricity vector \mathbf{e} (eq. 4 of Box 3.2) helps us to simplify this equation. We have

$$\begin{aligned} \mathbf{e} &= \mathbf{v} \times (\mathbf{r} \times \mathbf{v}) - GM \hat{\mathbf{e}}_r \\ &= |\mathbf{r}'|^2 \frac{\mathbf{r}}{r^2} - \frac{r'}{r} \mathbf{r}' - GM \frac{\mathbf{r}}{r}, \end{aligned} \quad (3.181)$$

where we have used $\mathbf{v} = \dot{\mathbf{r}} = \mathbf{r}'/r$ and the vector identity (B.9). Thus equation (3.180) can be written

$$\mathbf{r}'' = |\mathbf{r}'|^2 \frac{\mathbf{r}}{r^2} - 2GM \frac{\mathbf{r}}{r} - \mathbf{e} + r^2 \mathbf{g}. \quad (3.182)$$

The energy of the two-body orbit is

$$E_2 = \frac{1}{2}v^2 - \frac{GM}{r} = \frac{|\mathbf{r}'|^2}{2r^2} - \frac{GM}{r}, \quad (3.183)$$

so we arrive at the regularized equation of motion

$$\mathbf{r}'' - 2E_2\mathbf{r} = -\mathbf{e} + r^2\mathbf{g}, \quad (3.184)$$

in which the singularity at the origin has disappeared. This must be supplemented by equations for the rates of change of E_2 , \mathbf{e} , and t with fictitious time τ ,

$$E_2' = \mathbf{g} \cdot \mathbf{r}' \quad ; \quad \mathbf{e}' = 2\mathbf{r}(\mathbf{r}' \cdot \mathbf{g}) - \mathbf{r}'(\mathbf{r} \cdot \mathbf{g}) - \mathbf{g}(\mathbf{r} \cdot \mathbf{r}') \quad ; \quad t' = r. \quad (3.185)$$

When the external field \mathbf{g} vanishes, the energy E_2 and eccentricity vector \mathbf{e} are constants, the equation of motion (3.184) is that of a harmonic oscillator that is subject to a constant force $-\mathbf{e}$, and the fictitious time τ is proportional to the eccentric anomaly (Problem 3.29).

Figure 3.23 shows the fractional energy error that arises in the integration of one pericenter passage of an orbit in a Kepler potential with eccentricity $e = 0.99$. The error is plotted as a function of the number of force evaluations; this is the correct economic model if force evaluations dominate the computational cost, as is true for N-body integrations with $N \gg 1$. Note that even with $\gtrsim 1000$ force evaluations per orbit, a fourth-order Runge–Kutta integrator with adaptive timestep is sometimes unable to follow the orbit. Using the same integrator, Burdet–Heggie regularization reduces the energy error by almost five orders of magnitude.

This figure also shows the energy error that arises when integrating the same orbit using leapfrog with adaptive timestep (eq. 3.167) in unregularized coordinates. Even though leapfrog is only second-order, it achieves an accuracy that substantially exceeds that of the fourth-order Runge–Kutta integrator in unregularized coordinates, and approaches the accuracy of Burdet–Heggie regularization. Thus a time-symmetric leapfrog integrator provides much of the advantage of regularization without coordinate or time transformations.

(b) Kustaanheimo–Stiefel (KS) regularization An alternative regularization procedure, which involves the transformation of the coordinates in addition to time, can be derived using the symmetry group of the Kepler problem, the theory of quaternions and spinors, or several other methods (Stiefel & Scheifele 1971; Yoshida 1982; Heggie & Hut 2003). Once again we use the fictitious time τ defined by equation (3.178). We also define a four-vector $\mathbf{u} = (u_1, u_2, u_3, u_4)$ that is related to the position $\mathbf{r} = (x, y, z)$ by

$$\begin{aligned} u_1^2 &= \frac{1}{2}(x+r)\cos^2\psi & u_2 &= \frac{yu_1 + zu_4}{x+r} \\ u_4^2 &= \frac{1}{2}(x+r)\sin^2\psi & u_3 &= \frac{zu_1 - yu_4}{x+r}, \end{aligned} \quad (3.186)$$

where ψ is an arbitrary parameter. The inverse relations are

$$x = u_1^2 - u_2^2 - u_3^2 + u_4^2; \quad y = 2(u_1u_2 - u_3u_4); \quad z = 2(u_1u_3 + u_2u_4). \quad (3.187)$$

Note that $r = u_1^2 + u_2^2 + u_3^2 + u_4^2$. Let Φ_e be the potential that generates the external field $\mathbf{g} = -\nabla\Phi_e$. Then in terms of the new variables the equation of motion (3.177) reads

$$\begin{aligned} \mathbf{u}'' - \frac{1}{2}E\mathbf{u} &= -\frac{1}{4}\frac{\partial}{\partial\mathbf{u}}(|\mathbf{u}|^2\Phi_e), \\ E &= \frac{1}{2}v^2 - \frac{GM}{r} + \Phi_e = 2\frac{|\mathbf{u}'|^2}{|\mathbf{u}|^2} - \frac{GM}{|\mathbf{u}|^2} + \Phi_e, \\ E' &= |\mathbf{u}|^2\frac{\partial\Phi_e}{\partial t} \quad ; \quad t' = |\mathbf{u}|^2, \end{aligned} \quad (3.188)$$

When the external force vanishes, the first of equations (3.188) is the equation of motion for a four-dimensional harmonic oscillator.

Figure 3.23 shows the fractional energy error that arises in the integration of an orbit with eccentricity $e = 0.99$ using KS regularization. Using the same integrator, the energy error is more than an order of magnitude smaller than the error using Burdet–Heggie regularization.

3.5 Angle-action variables

In §3.1 we introduced the concept of an integral of motion and we saw that every spherical potential admitted at least four integrals I_i , namely, the Hamiltonian and the three components of angular momentum. Later we found that orbits in flattened axisymmetric potentials frequently admit three integrals, the classical integrals H and p_ϕ , and the non-classical third integral. Finally in §3.3 we found that many orbits in planar non-axisymmetric potentials admitted a non-classical integral in addition to the Hamiltonian.

In this section we explore the advantages of using integrals as coordinates for phase space. Since elementary Newtonian or Lagrangian mechanics restricts our choice of coordinates to ones that are rarely integrals, we work in the more general framework of Hamiltonian mechanics (Appendix D). For definiteness, we shall assume that there are three independent coordinates (so phase space is six-dimensional) and that we have three analytic isolating integrals $I_i(\mathbf{x}, \mathbf{v})$. We shall focus on a particular set of canonical coordinates, called **angle-action** variables; the three momenta are integrals, called “actions,” and the conjugate coordinates are called “angles.” An orbit fortunate enough to possess angle-action variables is called a **regular orbit**.

We start with a number of general results that apply to any system of angle-action variables. Then in a series of subsections we obtain explicit

expressions for these variables in terms of ordinary phase-space coordinates for spherical potentials, flattened axisymmetric potentials and planar, non-axisymmetric potentials. The section ends with a description of how actions enable us to solve problems in which the gravitational potential evolves slowly.

Angle-action variables cannot be defined for many potentials of practical importance for galactic dynamics. Nonetheless, the conceptual framework of angle-action variables proves extremely useful for understanding the complex phenomena that arise in potentials that do not admit them.

The discussion below is heuristic and non-rigorous; for a precise and elegant account see Arnold (1989).

3.5.1 Orbital tori

Let us denote the angle-action variables by $(\boldsymbol{\theta}, \mathbf{J})$. We assume that the momenta $\mathbf{J} = (J_1, J_2, J_3)$ are integrals of motion. Then Hamilton's equations (D.54) for the motion of the J_i read

$$0 = \dot{J}_i = -\frac{\partial H}{\partial \theta_i}. \quad (3.189)$$

Therefore, the Hamiltonian must be independent of the coordinates $\boldsymbol{\theta}$, that is $H = H(\mathbf{J})$. Consequently, we can trivially solve Hamilton's equations for the θ_i as functions of time:

$$\dot{\theta}_i = \frac{\partial H}{\partial J_i} \equiv \Omega_i(\mathbf{J}), \quad \text{a constant} \quad \Rightarrow \quad \theta_i(t) = \theta_i(0) + \Omega_i t. \quad (3.190)$$

So everything lies at our feet if we can install three integrals of motion as the momenta of a system of canonical coordinates.¹³

We restrict our attention to bound orbits. In this case, the Cartesian coordinates x_i cannot increase without limit as the θ_i do (eq. 3.190). From this we infer that the x_i are periodic functions of the θ_i . We can scale θ_i so that \mathbf{x} returns to its original value after θ_i has increased by 2π . Then we can expand \mathbf{x} in a Fourier series (Appendix B.4)

$$\mathbf{x}(\boldsymbol{\theta}, \mathbf{J}) = \sum_{\mathbf{n}} \mathbf{X}_{\mathbf{n}}(\mathbf{J}) e^{i\mathbf{n} \cdot \boldsymbol{\theta}}, \quad (3.191)$$

where the sum is over all vectors \mathbf{n} with integer components. When we eliminate the θ_i using equation (3.190), we find that the spatial coordinates are

¹³ To be able to use the J_i as a set of momenta, they must satisfy the canonical commutation relations (D.71), so we require $[J_i, J_j] = 0$; functions satisfying this condition are said to be **in involution**. For example, the components of angular momenta are not in involution: $[L_x, L_y] = L_z$, etc.

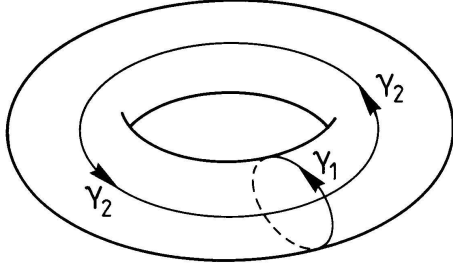


Figure 3.24 Two closed paths on a torus that cannot be deformed into one another, nor contracted to single points.

Fourier series in time, in which every frequency is a sum of integer multiples of the three **fundamental frequencies** $\Omega_i(\mathbf{J})$ that are defined by equation (3.190). Such a time series is said to be **conditionally periodic** or **quasiperiodic**.¹⁴ For example, in spherical potentials (§3.1) the periods T_r and T_ψ are inverses of such fundamental frequencies: $T_i = 2\pi/\Omega_i$. The third fundamental frequency is zero because the orbital plane is fixed in space—see §3.5.2.

An orbit is said to be **resonant** when its fundamental frequencies satisfy a relation of the form $\mathbf{n} \cdot \boldsymbol{\Omega} = 0$ for some integer triple $\mathbf{n} \neq \mathbf{0}$. Usually this implies that two of the frequencies are commensurable, that is the ratio Ω_i/Ω_j is a rational number ($-n_j/n_i$).

Consider the three-surface (i.e., volume) of fixed \mathbf{J} and varying $\boldsymbol{\theta}$. This is a cube of side-length 2π , and points on opposite sides must be identified since we have seen that incrementing, say, θ_1 by 2π while leaving θ_2, θ_3 fixed brings one back to the same point in phase space. A cube with faces identified in this way is called a three-torus by analogy with the connection between a rectangle and a two-torus: if we sew together opposite edges of a rectangular sheet of rubber, we generate the doughnut-shaped inner tube of a bicycle tire.

We shall find that these three-tori are in many respects identical with orbits, so it is important to have a good scheme for labeling them. The best set of labels proves to be the Poincaré invariants (Appendix D.4.2)

$$J'_i \equiv \frac{1}{2\pi} \iint d\mathbf{q} \cdot d\mathbf{p} = \frac{1}{2\pi} \iint \sum_{j=1,3} dq_j dp_j, \quad (3.192)$$

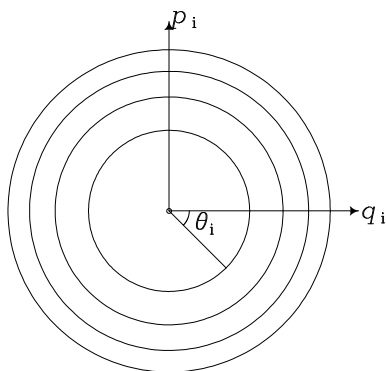
where the integral is over any surface that is bounded by the path γ_i on which θ_i increases from 0 to 2π while everything else is held constant (Figure 3.24). Since angle-action variables are canonical, $d\mathbf{q} \cdot d\mathbf{p} = d\boldsymbol{\theta} \cdot d\mathbf{J}$ (eq. D.84), so

$$J'_i = \frac{1}{2\pi} \iint_{\text{interior of } \gamma_i} d\boldsymbol{\theta} \cdot d\mathbf{J} = \frac{1}{2\pi} \iint_{\text{interior of } \gamma_i} d\theta_i dJ_i. \quad (3.193)$$

¹⁴ Observers of binary stars use the term quasiperiod more loosely. Our usage is equivalent to what is meant by a quasicrystal: a structure whose Fourier transform is discrete, but in which there are more fundamental frequencies than independent variables (in our case one, t , in a quasicrystal three x, y, z).

Box 3.4: Angle-action variables as polar coordinates

The figure shows the intersection with a coordinate plane of some of the nested orbital tori of a two-dimensional harmonic oscillator. The coordinates q_i, p_i have been scaled such that the tori appear as circles. The values of the action J_i on successive tori are chosen to be $0, 1, 2, \dots$ (in some suitable units), so, by equation (3.192), the areas inside successive tori are $0, 2\pi, 4\pi, \dots$. Hence, the radii $r = (q_i^2 + p_i^2)^{1/2}$ of successive circles are $\sqrt{2} \times (0, 1, \sqrt{2}, \sqrt{3}, \dots)$. In general the radius of



the circle associated with the torus on which J_i takes the value J' is $r = \sqrt{2J'}$. In this plane, the angle variable θ_i is closely analogous to the usual azimuthal angle. Hence, angle-action variables are closely analogous to plane polar coordinates, the major difference being that coordinate circles are labeled not by radius but by $\sqrt{2}$ times the area they enclose. The generating function for the transformation from (θ_i, J_i) to (q_i, p_i) is given in Problem 3.31.

As Box 3.4 explains, angle-action variables are a kind of polar coordinates for phase space, and have a coordinate singularity within the domain of integration. We must exclude this from the domain of integration before we use Green's theorem to convert the surface integral in (3.193) into a line integral. The value of our surface integral is unchanged by excluding this point, but when we use Green's theorem (eq. B.61) on the original domain less the excluded point, we obtain two line integrals, one along the curve γ_i and one along the boundary that surrounds the excluded point—along this second boundary, J_i takes some definite value, J_i^c , say, and θ_i takes all values in the range $(0, 2\pi)$. Thus we have

$$J'_i = \frac{1}{2\pi} \left(\oint_{\gamma_i} J_i d\theta_i - \oint_{J_i=J_i^c} J_i d\theta_i \right) = J_i - J_i^c. \quad (3.194)$$

This equation shows that the label J'_i defined by equation (3.192) will be identical with our original action coordinate J_i providing we set $J_i = 0$ at the coordinate singularity that marks the center of the angle-action coordinate system. We shall henceforth assume that this choice has been made.

In practical applications we often evaluate the integral of equation (3.192) using phase-space coordinates that have no singularity within the

domain of integration. Then we can replace the surface integral with a line integral that is easier to evaluate:

$$J_i = \frac{1}{2\pi} \oint_{\gamma_i} \mathbf{p} \cdot d\mathbf{q}. \quad (3.195)$$

(a) Time averages theorem In Chapter 4 we shall make extensive use of a result that we can now prove.

Time averages theorem *When a regular orbit is non-resonant, the average time that the phase point of a star on that orbit spends in any region D of its torus is proportional to the integral $V(D) = \int_D d^3\boldsymbol{\theta}$.*

Proof: Let f_D be the function such that $f_D(\boldsymbol{\theta}) = 1$ when the point $\boldsymbol{\theta}$ lies in D , and is zero otherwise. We may expand f_D in a Fourier series (cf. eq. 3.191)

$$f_D(\boldsymbol{\theta}) = \sum_{\mathbf{n}=-\infty}^{\infty} F_{\mathbf{n}} \exp(i\mathbf{n} \cdot \boldsymbol{\theta}). \quad (3.196)$$

Now

$$\int_{\text{torus}} d^3\boldsymbol{\theta} f_D(\boldsymbol{\theta}) = \int_D d^3\boldsymbol{\theta} = V(D). \quad (3.197a)$$

With equation (3.196) we therefore have

$$V(D) = \int_{\text{torus}} d^3\boldsymbol{\theta} f_D(\boldsymbol{\theta}) = \sum_{\mathbf{n}=-\infty}^{\infty} F_{\mathbf{n}} \prod_{k=1}^3 \int_0^{2\pi} d\psi \exp(in_k\psi) = (2\pi)^3 F_{\mathbf{0}}. \quad (3.197b)$$

On the other hand, the fraction of the interval $(0, T)$ during which the star's phase point lies in D is

$$\tau_T(D) = \frac{1}{T} \int_0^T dt f_D[\boldsymbol{\theta}(t)], \quad (3.198)$$

where $\boldsymbol{\theta}(t)$ is the position of the star's phase point at time t . With equations (3.190) and (3.196), equation (3.198) becomes

$$\begin{aligned} \tau_T(D) &= \frac{1}{T} \sum_{\mathbf{n}} e^{i\mathbf{n} \cdot \boldsymbol{\theta}(0)} \int_0^T dt F_{\mathbf{n}} e^{i(\mathbf{n} \cdot \boldsymbol{\Omega})t} \\ &= F_{\mathbf{0}} + \frac{1}{T} \sum_{\mathbf{n} \neq \mathbf{0}} e^{i\mathbf{n} \cdot \boldsymbol{\theta}(0)} F_{\mathbf{n}} \frac{e^{i(\mathbf{n} \cdot \boldsymbol{\Omega})T} - 1}{i\mathbf{n} \cdot \boldsymbol{\Omega}}. \end{aligned} \quad (3.199)$$

Thus

$$\lim_{T \rightarrow \infty} \tau_T(D) = F_{\mathbf{0}} = \frac{V(D)}{(2\pi)^3}, \quad (3.200)$$

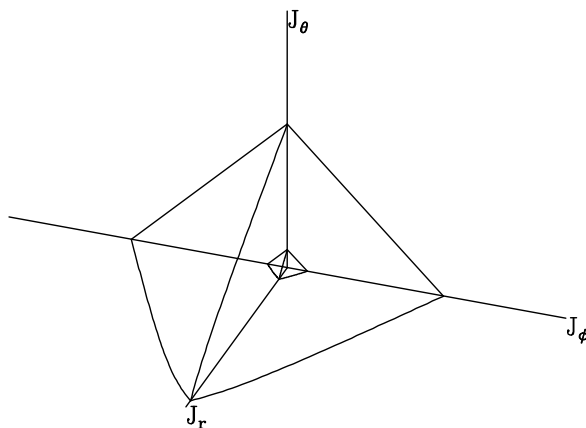


Figure 3.25 The action space of an axisymmetric potential. Two constant-energy surfaces are shown for the spherical isochrone potential (2.47). The surfaces $H = -0.5(GM/2b)$ and $H = -0.03(GM/2b)$ are shown (eq. 3.226) with the axes all scaled to length $5\sqrt{GMb}$.

which completes the proof.◁

Note that if the orbit is resonant, $\mathbf{n} \cdot \boldsymbol{\Omega}$ vanishes for some $\mathbf{n} \neq 0$ and the second equality in (3.199) becomes invalid, so the theorem cannot be proved. In fact, if $\Omega_i : \Omega_j = m : n$, say, then by equations (3.190) $I_4 \equiv n\theta_i - m\theta_j$ becomes an isolating integral that confines the star to a spiral on the torus. We shall see below that motion in a spherical potential provides an important example of this phenomenon.

(b) Action space In Chapter 4 we shall develop the idea that galaxies are made up of orbits, and we shall find it helpful to think of whole orbits as single points in an abstract space. Any isolating integrals can serve as coordinates for such a representation, but the most advantageous coordinates are the actions. We define **action space** to be the imaginary space whose Cartesian coordinates are the actions. Figure 3.25 shows the action space of an axisymmetric potential, when the actions can be taken to be generalizations of the actions for spherical potentials listed in Table 3.1 below. Points on the axes represent orbits for which only one of the integrals (3.192) is non-zero. These are the closed orbits. The origin represents the orbit of a star that just sits at the center of the potential. In each octant, surfaces of constant energy are approximate planes; by equation (3.190) the local normal to this surface is parallel to the vector $\boldsymbol{\Omega}$. Every point in the positive quadrant $J_r, J_\theta \geq 0$, all the way to infinity, represents a bound orbit.

A region R_3 in action space represents a group of orbits. Let the volume of R_3 be V_3 . The volume of six-dimensional phase space occupied by the orbits is

$$V_6 = \int_{R_6} d^3\mathbf{x} d^3\mathbf{v}, \quad (3.201)$$

where R_6 is the region of phase space visited by stars on the orbits of R_3 . Since the coordinate set $(\mathbf{J}, \boldsymbol{\theta})$ is canonical, $d^3\mathbf{x}d^3\mathbf{v} = d^3\mathbf{J}d^3\boldsymbol{\theta}$ (see eq. D.81) and thus

$$V_6 = \int_{R_6} d^3\mathbf{J}d^3\boldsymbol{\theta}. \quad (3.202)$$

But for any orbit the angle variables cover the range $(0, 2\pi)$, so we may immediately integrate over the angles to find

$$V_6 = (2\pi)^3 \int_{R_3} d^3\mathbf{J} = (2\pi)^3 V_3. \quad (3.203)$$

Thus the volume of a region of action space is directly proportional to the volume of phase space occupied by its orbits.

(c) Hamilton–Jacobi equation The transformation between any two sets of canonical coordinates can be effected with a generating function (Appendix D.4.6). Let $S(\mathbf{q}, \mathbf{J})$ be the (unknown) generating function of the transformation between angle-action variables and ordinary phase space coordinates (\mathbf{q}, \mathbf{p}) such as $\mathbf{q} = \mathbf{x}$, $\mathbf{p} = \mathbf{v}$. Then (eq. D.93)

$$\boldsymbol{\theta} = \frac{\partial S}{\partial \mathbf{J}} \quad ; \quad \mathbf{p} = \frac{\partial S}{\partial \mathbf{q}}, \quad (3.204)$$

where \mathbf{p} and $\boldsymbol{\theta}$ are now to be considered functions of \mathbf{q} and \mathbf{J} . We can use $S(\mathbf{q}, \mathbf{J})$ to eliminate \mathbf{p} from the usual Hamiltonian function $H(\mathbf{q}, \mathbf{p})$ and thus express H as a function

$$H\left(\mathbf{q}, \frac{\partial S}{\partial \mathbf{q}}(\mathbf{q}, \mathbf{J})\right).$$

of (\mathbf{q}, \mathbf{J}) . By moving along an orbit, we can vary the q_i while holding constant the J_i . As we vary the q_i in this way, H must remain constant at the energy E of the orbit in question. This suggests that we investigate the partial differential equation

$$H\left(\mathbf{q}, \frac{\partial S}{\partial \mathbf{q}}\right) = E \quad \text{at fixed } \mathbf{J}. \quad (3.205)$$

If we can solve this **Hamilton–Jacobi equation**, our solution should contain some arbitrary constants K_i —we shall see below that we usually solve the equation by the method of separation of variables (e.g., §2.4) and the constants are separation constants. We identify the K_i with functions of the actions as follows. Eliminating \mathbf{p} from equation (3.195) we have

$$J_i = \frac{1}{2\pi} \oint_{\gamma_i} \frac{\partial S}{\partial \mathbf{q}} \cdot d\mathbf{q} = \frac{\Delta S(\mathbf{K})}{2\pi}. \quad (3.206)$$

This equation states that J_i is proportional to the increment in the generating function when one passes once around the torus on the i th path— S , like the magnetic scalar potential around a current-carrying wire, is a multivalued function. The increment in S , and therefore J_i , depends on the integration constants that appear in S , so these are functions of the actions.

Once the Hamilton–Jacobi equation has been solved and the integrals in (3.206) have been evaluated, S becomes a known function $S(\mathbf{q}, \mathbf{J})$ and we can henceforth use equations (3.204) to transform between angle-action variables and ordinary phase-space coordinates. In particular, we can integrate orbits trivially by transforming the initial conditions into angle-action variables, incrementing the angles, and transforming back to ordinary phase-space coordinates.

Let us see how this process works in a simple example. The Hamiltonian of a two-dimensional harmonic oscillator is

$$H(\mathbf{x}, \mathbf{p}) = \frac{1}{2}(p_x^2 + p_y^2 + \omega_x^2 x^2 + \omega_y^2 y^2). \quad (3.207)$$

Substituting in $p_x = \partial S / \partial x$, $p_y = \partial S / \partial y$ (eq. 3.204), the Hamilton–Jacobi equation reads

$$\left(\frac{\partial S}{\partial x}\right)^2 + \left(\frac{\partial S}{\partial y}\right)^2 + \omega_x^2 x^2 + \omega_y^2 y^2 = 2E, \quad (3.208)$$

where S is a function of x, y and \mathbf{J} . We solve this partial differential equation by the method of separation of variables.¹⁵ We write $S(x, y, \mathbf{J}) = S_x(x, \mathbf{J}) + S_y(y, \mathbf{J})$ and rearrange the equation to

$$\left(\frac{\partial S_x}{\partial x}\right)^2 + \omega_x^2 x^2 = 2E - \left(\frac{\partial S_y}{\partial y}\right)^2 - \omega_y^2 y^2. \quad (3.209)$$

The left side does not depend on y and the right side does not depend on x . Consequently, each side can only be a function of \mathbf{J} , which we call $K^2(\mathbf{J})$ because it is evidently non-negative:

$$K^2 \equiv \left(\frac{\partial S_x}{\partial x}\right)^2 + \omega_x^2 x^2. \quad (3.210)$$

It follows that

$$S_x(x, \mathbf{J}) = K \int^x dx' \epsilon \sqrt{1 - \frac{\omega_x^2 x'^2}{K^2}},$$

¹⁵ When this method is applied in quantum mechanics and in potential theory (e.g., §2.4) one usually assumes that the dependent variable is a *product* of functions of one variable, rather than a sum of such functions as here.

where ϵ is chosen to be ± 1 so that $S_x(x, \mathbf{J})$ increases continuously along a path over the orbital torus. Changing the variable of integration, we have

$$\begin{aligned} S_x(x, \mathbf{J}) &= \frac{K^2}{\omega_x} \int d\psi' \sin^2 \psi' \quad \text{where} \quad x = -\frac{K}{\omega_x} \cos \psi \\ &= \frac{K^2}{2\omega_x} (\psi - \frac{1}{2} \sin 2\psi). \end{aligned} \quad (3.211)$$

Moreover, $p_x = \partial S / \partial x = \epsilon K \sqrt{1 - \omega_x^2 x^2 / K^2} = K \sin \psi$, so both x and p_x return to their old values when ψ is incremented by 2π . We infer that incrementing ψ by 2π carries us around the path γ_x that is associated with J_x through equation (3.192). Thus equation (3.206) now yields

$$J_x = \frac{\Delta S}{2\pi} = \frac{\Delta S_x}{2\pi}. \quad (3.212)$$

Equation (3.211) tells us that when ψ is incremented by 2π , S_x increases by $K^2\pi/\omega_x$. Hence,

$$J_x(x, p_x) = \frac{K^2}{2\omega_x} = \frac{p_x^2 + \omega_x^2 x^2}{2\omega_x}, \quad (3.213)$$

where the last equality follows from (3.210) with $\partial S_x / \partial x$ replaced by p_x . The solution for $J_y(y, p_y)$ proceeds analogously and yields

$$J_y(y, p_y) = \frac{2E - K^2}{2\omega_y} = \frac{p_y^2 + \omega_y^2 y^2}{2\omega_y}. \quad (3.214)$$

Comparing with equation (3.207), we find that

$$H(\mathbf{J}) = \omega_x J_x + \omega_y J_y. \quad (3.215)$$

Notice from (3.215) that $\Omega_x \equiv \partial H / \partial J_x = \omega_x$ and similarly for Ω_y .

Finally we determine the angle variables from the second of equations (3.204). The obvious procedure is to eliminate both K and ψ from equation (3.211) in favor of J_x and x . In practice it is expedient to leave ψ in and treat it as a function of J_x and x :

$$S_x(x, \mathbf{J}) = J_x(\psi - \frac{1}{2} \sin 2\psi) \quad \text{where} \quad \cos \psi = -\sqrt{\frac{\omega_x}{2J_x}} x, \quad (3.216)$$

so

$$\begin{aligned} \theta_x &= \frac{\partial S}{\partial J_x} = \psi - \frac{1}{2} \sin 2\psi + J_x(1 - \cos 2\psi) \frac{\partial \psi}{\partial J_x} \\ &= \psi - \frac{1}{2} \sin 2\psi + \sin^2 \psi \cot \psi \\ &= \psi. \end{aligned} \quad (3.217)$$

Thus the variable ψ that we introduced for convenience in doing an integral is, in fact, the angle variable conjugate to J_x . Problem 3.33 explains an alternative, and sometimes simpler, route to the angle variables.

3.5.2 Angle-action variables for spherical potentials

We now derive angle-action variables for any spherical potential. These are useful not only for strictly spherical systems, but also for axisymmetric disks, and serve as the starting point for perturbative analyses of mildly aspherical potentials. To minimize confusion between ordinary spherical polar coordinates and angle variables, in this section we reserve ϑ for the usual polar angle, and continue to use θ_i for the variable conjugate to J_i .

The Hamilton–Jacobi equation (3.205) for the potential $\Phi(r)$ is

$$\begin{aligned} E &= \frac{1}{2} |\nabla S|^2 + \Phi(r) \\ &= \frac{1}{2} \left[\left(\frac{\partial S}{\partial r} \right)^2 + \left(\frac{1}{r} \frac{\partial S}{\partial \vartheta} \right)^2 + \left(\frac{1}{r \sin \vartheta} \frac{\partial S}{\partial \phi} \right)^2 \right] + \Phi(r), \end{aligned} \quad (3.218)$$

where we have used equation (B.38) for the gradient operator in spherical polar coordinates. We write the generating function as $S(\mathbf{x}, \mathbf{J}) = S_r(r, \mathbf{J}) + S_\vartheta(\vartheta, \mathbf{J}) + S_\phi(\phi, \mathbf{J})$ and solve (3.218) by separation of variables. With the help of equation (3.204) we find

$$L_z^2 = \left(\frac{\partial S_\phi}{\partial \phi} \right)^2 = p_\phi^2, \quad (3.219a)$$

$$L^2 - \frac{L_z^2}{\sin^2 \vartheta} = \left(\frac{\partial S_\vartheta}{\partial \vartheta} \right)^2 = p_\vartheta^2, \quad (3.219b)$$

$$2E - 2\Phi(r) - \frac{L^2}{r^2} = \left(\frac{\partial S_r}{\partial r} \right)^2 = p_r^2. \quad (3.219c)$$

Here we have introduced two separation constants, L and L_z . We assume that $L > 0$ and choose the sign of L_z so that $L_z = p_\phi$; with these conventions L and L_z prove to be the magnitude and z -component of the angular-momentum vector (Problem 3.20). Taking the square root of each equation and integrating, we obtain a formula for S :

$$\begin{aligned} S(\mathbf{x}, \mathbf{J}) &= \int_0^\phi d\phi L_z + \int_{\pi/2}^\vartheta d\vartheta \epsilon_\vartheta \sqrt{L^2 - \frac{L_z^2}{\sin^2 \vartheta}} \\ &\quad + \int_{r_{\min}}^r dr \epsilon_r \sqrt{2E - 2\Phi(r) - \frac{L^2}{r^2}}, \end{aligned} \quad (3.220)$$

where ϵ_ϑ and ϵ_r are chosen to be ± 1 such that the integrals in which they appear increase monotonically along a path over the orbital torus. The lower limits of these integrals specify some point on the orbital torus, and are arbitrary. It is convenient to take r_{\min} to be the orbit's pericenter radius.

To obtain the actions from equation (3.206) we have to evaluate the change in S as we go round the orbital torus along curves on which only one

of the coordinates is incremented. The case of changing ϕ is easy: on the relevant curve, ϕ increases by 2π , so (3.220) states that $\Delta S = 2\pi L_z$ and

$$J_\phi = L_z. \quad (3.221)$$

We call J_ϕ the **azimuthal action**. Consider next the case of changing ϑ . Let ϑ_{\min} be the smallest value that ϑ attains on the orbit, given by

$$\sin \vartheta_{\min} = \frac{|L_z|}{L}, \quad (3.222)$$

$\vartheta_{\min} \leq \pi/2$. Then we start at $\pi/2$, where the integrand peaks, and integrate to $\pi - \vartheta_{\min}$, where it vanishes. We have now integrated over a quarter period of the integrand, so the whole integral is four times the value from this leg,

$$J_\vartheta = \frac{2}{\pi} \int_{\pi/2}^{\pi - \vartheta_{\min}} d\vartheta \sqrt{L^2 - \frac{L_z^2}{\sin^2 \vartheta}} = L - |L_z|. \quad (3.223)$$

We call J_ϑ the **latitudinal action**. The evaluation of J_r from equations (3.206) and (3.220) proceeds similarly and yields

$$J_r = \frac{1}{\pi} \int_{r_{\min}}^{r_{\max}} dr \sqrt{2E - 2\Phi(r) - \frac{L^2}{r^2}}, \quad (3.224)$$

where r_{\max} is the radius of the apocenter— r_{\min} and r_{\max} are the two roots of the radical—and J_r is the **radial action**.

An important example is that of the isochrone potential (2.47), which encompasses both the Kepler and spherical harmonic potentials as limiting cases. One finds that (Problem 3.41)

$$J_r = \frac{GM}{\sqrt{-2E}} - \frac{1}{2} \left(L + \frac{1}{2} \sqrt{L^2 - 4GMb} \right). \quad (3.225)$$

If we rewrite this expression as an equation for the Hamiltonian $H_I = E$ as a function of the actions, we obtain

$$H_I(\mathbf{J}) = - \frac{(GM)^2}{2 \left[J_r + \frac{1}{2} (L + \sqrt{L^2 + 4GMb}) \right]^2} \quad (L = J_\theta + |J_\phi|). \quad (3.226a)$$

Differentiating this expression with respect to the actions, we find the frequencies (eq. 3.190):

$$\begin{aligned} \Omega_r &= \frac{(GM)^2}{\left[J_r + \frac{1}{2} (L + \sqrt{L^2 + 4GMb}) \right]^3} \\ \Omega_\vartheta &= \frac{1}{2} \left(1 + \frac{L}{\sqrt{L^2 + 4GMb}} \right) \Omega_r \quad ; \quad \Omega_\phi = \text{sgn}(J_\phi) \Omega_\vartheta. \end{aligned} \quad (3.226b)$$

It is straightforward to check that these results are consistent with the radial and azimuthal periods determined in §3.1c.¹⁶ In the limit $b \rightarrow 0$, the isochrone potential becomes the Kepler potential and all three frequencies become equal. The corresponding results for the spherical harmonic oscillator are obtained by examining the limit $b \rightarrow \infty$ (Problem 3.36).

J_θ and J_ϕ occur in equations (3.226) only in the combination $L = J_\theta + |J_\phi|$, and in fact, the Hamiltonian for any spherical potential is a function $H(J_r, L)$. Therefore we elevate L to the status of an action by making the canonical transformation that is defined by the generating function (eq. D.93)

$$S' = \theta_\phi J_1 + \theta_\vartheta (J_2 - |J_1|) + \theta_r J_3, \quad (3.227)$$

where (J_1, J_2, J_3) are new angle-action coordinates. Differentiating with respect to the old angles, we discover the connection between the new and old actions:

$$\begin{aligned} J_\phi &= \frac{\partial S'}{\partial \theta_\phi} = J_1 \quad \Rightarrow \quad J_1 = L_z, \\ J_\vartheta &= \frac{\partial S'}{\partial \theta_\vartheta} = J_2 - |J_1| \quad \Rightarrow \quad J_2 = J_\vartheta + |J_\phi| = L, \\ J_r &= \frac{\partial S'}{\partial \theta_r} = J_3. \end{aligned} \quad (3.228)$$

Thus the new action J_2 is L as desired. Differentiating S' with respect to the new actions we find that the new angles are

$$\theta_1 = \theta_\phi - \text{sgn}(J_1)\theta_\vartheta \quad ; \quad \theta_2 = \theta_\vartheta \quad ; \quad \theta_3 = \theta_r. \quad (3.229)$$

Equation (3.224) can be regarded as an implicit equation for the Hamiltonian $H(\mathbf{J}) = E$ in terms of $J_3 = J_r$ and $J_2 = L$. Since J_1 does not appear in this equation, the Hamiltonian of *any* spherical potential must be of the form $H(J_2, J_3)$. Thus $\Omega_1 = \partial H / \partial J_1 = 0$ for all spherical potentials, and the corresponding angle θ_1 is an integral of motion. In §3.1 we saw that any spherical potential admits four isolating integrals. Here we have recovered this result from a different point of view: three of the integrals are the actions (J_1, J_2, J_3) , and the fourth is the angle θ_1 .

From Figure 3.26 we see that for orbits with $L_z > 0$ the inclination of the orbital plane $i = \frac{1}{2}\pi - \vartheta_{\min}$, while when $L_z < 0$, $i = \frac{1}{2}\pi + \vartheta_{\min}$. Combining these equations with (3.222) we find that

$$i = \cos^{-1}(L_z/L) = \cos^{-1}(J_1/J_2). \quad (3.230)$$

¹⁶ A minor difference is that in the analysis of §3.1c, the angular momentum L could have either sign. Here $L = |\mathbf{L}|$ is always non-negative, while L_z can have either sign.

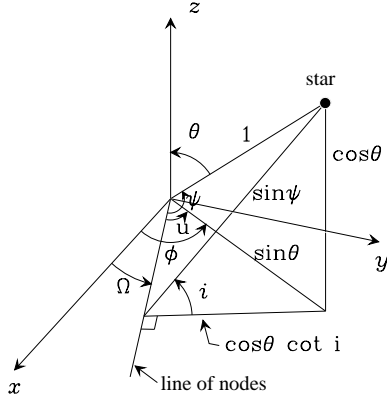


Figure 3.26 Angles defined by an orbit. The orbit is confined to a plane whose normal makes an angle i , the **inclination**, with the z axis. The orbital plane intersects the xy plane along the **line of nodes**. The ascending node is the node at which $z > 0$, and the angle Ω is the longitude of the ascending node. Elementary trigonometry shows that $u = \sin^{-1}(\cot i \cot \vartheta) = \phi - \Omega$ and that $\cos \vartheta = \sin i \sin \psi$, where ψ is the angle between the line of nodes and the radius vector to the star.

We now obtain explicit expressions for the angle variables of any spherical potential by evaluating $\partial S / \partial J_i = \theta_i$, where S is derived from equation (3.220) by replacing E with $H(J_2, J_3)$, L with J_2 , and L_z with J_1 . We have

$$S = \phi J_1 + \int_{\pi/2}^{\vartheta} d\vartheta \epsilon_{\vartheta} \sqrt{J_2^2 - \frac{J_1^2}{\sin^2 \vartheta}} + \int_{r_{\min}}^r dr \epsilon_r \sqrt{2H(J_2, J_3) - 2\Phi(r) - \frac{J_2^2}{r^2}}. \quad (3.231)$$

Figure 3.26 helps us to interpret our final result. It depicts the star after it has passed the line of nodes, moving upward. At this instant, $\dot{\vartheta} < 0$, and we must choose $\epsilon_{\vartheta} = -1$ to make the first integral of equation (3.231) increasing. We therefore specialize to this case, and using (3.230) find

$$\begin{aligned} \theta_1 &= \frac{\partial S}{\partial J_1} = \phi + \operatorname{sgn}(J_1) \int_{\pi/2}^{\vartheta} \frac{d\vartheta}{\sin \vartheta \sqrt{\sin^2 \vartheta \sec^2 i - 1}} \\ &= \phi - u, \end{aligned} \quad (3.232a)$$

where¹⁷

$$\sin u \equiv \cot i \cot \vartheta. \quad (3.232b)$$

Figure 3.26 demonstrates that the new variable u is actually $\phi - \Omega$ and thus that $\theta_1 = \Omega$, the **longitude of the ascending node**.¹⁸ Thus θ_1 is constant because the line of nodes is fixed. If the potential were not spherical, but

¹⁷ This follows because

$$\begin{aligned} d[\sin^{-1}(\cot i \cot \vartheta)] &= -\csc^2 \vartheta \cot i d\vartheta / \sqrt{1 - \cot^2 i \cot^2 \vartheta} \\ &= \operatorname{sgn}(\cos i) / (\sin \vartheta \sqrt{\sin^2 \vartheta \sec^2 i - 1}). \end{aligned} \quad (3.233)$$

¹⁸ Equation (3.232b) has two solutions in $(0, 2\pi)$ and care must be taken to choose the correct solution.

Table 3.1 Angle-action variables in a spherical potential

actions	$J_\phi = L_z$; $J_\vartheta = L - L_z $; J_r
angles	$\theta_\phi = \Omega + \text{sgn}(L_z)\theta_\vartheta$; θ_ϑ ; θ_r
Hamiltonian	$H(J_\vartheta + J_\phi , J_r)$
frequencies	$\Omega_\phi = \text{sgn}(L_z)\Omega_\vartheta$; Ω_ϑ ; Ω_r
actions	$J_1 = L_z$; $J_2 = L$; $J_3 = J_r$
angles	$\theta_1 = \Omega$; $\theta_2 = \theta_\vartheta$; $\theta_3 = \theta_r$
Hamiltonian	$H(J_2, J_3)$
frequencies	$\Omega_1 = 0$; $\Omega_2 = \Omega_\vartheta$; $\Omega_3 = \Omega_r$
actions	$J_a = L_z$; $J_b = L$; $J_c = J_r + L$
angles	$\theta_a = \Omega$; $\theta_b = \theta_\vartheta - \theta_r$; $\theta_c = \theta_r$
Hamiltonian	$H(J_b, J_c - J_b)$
frequencies	$\Omega_a = 0$; $\Omega_b = \Omega_\vartheta - \Omega_r$; $\Omega_c = \Omega_r$

NOTES: The Delaunay variables (J_a, J_b, J_c) are defined in Appendix E. When possible, actions and angles are expressed in terms of the total angular momentum L , the z -component of angular momentum L_z , the radial action J_r , and the longitude of the ascending node Ω (Figure 3.26). Unfortunately, Ω is also used for the frequency corresponding to a given action, but in this case it is always accompanied by a subscript. The Hamiltonian is $H(L, J_r)$.

merely axisymmetric, θ_1 would not be constant and the orbital plane would precess.

Next we differentiate equation (3.231) to obtain θ_3 . Only the third term, which is equal to S_r , depends on J_3 . Thus we have

$$\theta_3 = \left(\frac{\partial S_r}{\partial J_3} \right)_{J_2} = \left(\frac{\partial S_r}{\partial H} \right)_{J_2} \left(\frac{\partial H}{\partial J_3} \right)_{J_2} = \left(\frac{\partial S_r}{\partial H} \right)_{J_2} \Omega_3, \quad (3.234)$$

where the last step follows from equation (3.190). Similarly,

$$\theta_2 = \left(\frac{\partial S}{\partial J_2} \right)_{J_3} = \left(\frac{\partial S_\vartheta}{\partial J_2} \right)_{J_3} + \left(\frac{\partial S_r}{\partial H} \right)_{J_2} \left(\frac{\partial H}{\partial J_2} \right)_{J_3} + \left(\frac{\partial S_r}{\partial J_2} \right)_H. \quad (3.235)$$

We eliminate $\partial S_r / \partial H$ using equation (3.234),

$$\theta_2 = \left(\frac{\partial S_\vartheta}{\partial J_2} \right)_{J_3} + \frac{\Omega_2}{\Omega_3} \theta_3 + \left(\frac{\partial S_r}{\partial J_2} \right)_H. \quad (3.236)$$

From equation (3.231) with $\epsilon_\vartheta = -1$, it is straightforward to show that

$$\left(\frac{\partial S_\vartheta}{\partial J_2} \right)_{J_1} = \sin^{-1} \left(\frac{\cos \vartheta}{\sin i} \right). \quad (3.237)$$

Now let ψ be the angle measured in the orbital plane from the line of nodes to the current position of the star. From Figure 3.26 it is easy to see that $\cos \vartheta = \sin i \sin \psi$; thus

$$\left(\frac{\partial S_\vartheta}{\partial J_2} \right)_{J_1} = \psi. \quad (3.238)$$

The other two partial derivatives in equations (3.234) and (3.235) can only be evaluated once $\Phi(r)$ has been chosen. In the case of the isochrone potential (2.47), we have

$$\left(\frac{\partial S_r}{\partial H}\right)_{J_2} = \int_{r_{\min}}^r dI \quad ; \quad \left(\frac{\partial S_r}{\partial J_2}\right)_H = -J_2 \int_{r_{\min}}^r \frac{dI}{r^2} \quad (3.239)$$

where dI is defined by (3.36). Hence the integrals to be performed are just indefinite versions of the definite integrals that yielded T_r and $\Delta\psi$ in §3.1c. The final answers are most conveniently expressed in terms of an auxiliary variable η that is defined by (cf. eqs. 3.28, 3.32 and 3.34)

$$s = 2 + \frac{c}{b}(1 - e \cos \eta) \quad \text{where} \quad \begin{cases} c \equiv \frac{GM}{-2H} - b, \\ e^2 \equiv 1 - \frac{J_2^2}{GMc} \left(1 + \frac{b}{c}\right), \\ s \equiv 1 + \sqrt{1 + r^2/b^2}. \end{cases} \quad (3.240)$$

Then one has¹⁹

$$\begin{aligned} \theta_3 &= \eta - \frac{ec}{c+b} \sin \eta \\ \theta_2 &= \psi + \frac{\Omega_2}{\Omega_3} \theta_3 - \tan^{-1} \left(\sqrt{\frac{1+e}{1-e}} \tan\left(\frac{1}{2}\eta\right) \right) \\ &\quad - \frac{1}{\sqrt{1+4GMb/J_2^2}} \tan^{-1} \left(\sqrt{\frac{1+e+2b/c}{1-e+2b/c}} \tan\left(\frac{1}{2}\eta\right) \right). \end{aligned} \quad (3.241)$$

Thus in the case of the isochrone potential we can analytically evaluate all three angle variables from ordinary phase-space coordinates (\mathbf{x}, \mathbf{v}) .

To summarize, in an arbitrary spherical potential two of the actions can be taken to be the total angular momentum L and its z -component L_z , and one angle can be taken to be the longitude of the ascending node Ω . The remaining action and angles can easily be determined by numerical evaluation of the integral (3.224) and integrals analogous to those of equation (3.239). In the isochrone potential, all angle-action variables can be obtained analytically from the ordinary phase-space coordinates (\mathbf{x}, \mathbf{v}) . The analytic relations among angle-action variables in spherical potentials are summarized in Table 3.1.

¹⁹ In numerical work, care must be taken to ensure that the branch of the inverse trigonometric functions is chosen so that the angle variables increase continuously.

3.5.3 Angle-action variables for flattened axisymmetric potentials

In §3.2 we used numerical integrations to show that most orbits in flattened axisymmetric potentials admit three isolating integrals, only two of which were identified analytically. Now we take up the challenge of identifying the missing “third integral” analytically, and deriving angle-action variables from it and the classical integrals. It proves possible to do this only for special potentials, and we start by examining the potentials for which we *have* obtained action integrals for clues as to what a promising potential might be.

(a) Stäckel potentials In §3.3 we remarked that box orbits in a planar non-rotating bar potential resemble Lissajous figures generated by two-dimensional harmonic motion, while loop orbits have many features in common with orbits in axisymmetric potentials. Let us examine these parallels more closely. The orbits of a two-dimensional harmonic oscillator admit two independent isolating integrals, $H_x = \frac{1}{2}(p_x^2 + \omega_x^2 x^2)$ and $H_y = \frac{1}{2}(p_y^2 + \omega_y^2 y^2)$ (eq. 3.207). At each point in the portion of the (x, y) plane visited by the orbit, the particle can have one of four momentum vectors. These momenta arise from the ambiguity in the signs of p_x and p_y when we are given only E_x and E_y , the values of H_x and H_y : $p_x(x) = \pm\sqrt{2E_x - \omega_x^2 x^2}$; $p_y(y) = \pm\sqrt{2E_y - \omega_y^2 y^2}$. The boundaries of the orbit are the lines on which $p_x = 0$ or $p_y = 0$.

Consider now planar orbits in a axisymmetric potential $\Phi(r)$. These orbits fill annuli. At each point in the allowed annulus two momenta are possible: $p_r(r) = \pm\sqrt{2(E - \Phi) - L_z^2/r^2}$, $p_\phi = L_z$. The boundaries of the orbit are the curves on which $p_r = 0$.

These examples have a number of important points in common:

- (i) The boundaries of orbits are found by equating to zero one canonical momentum in a coordinate system that reflects the symmetry of the potential.
- (ii) The momenta in this privileged coordinate system can be written as functions of only one variable: $p_x(x)$ and $p_y(y)$ in the case of the harmonic oscillator; and $p_r(r)$ and $p_\phi = L_z$ (which depends on neither coordinate) in the case of motion in a axisymmetric potential.
- (iii) These expressions for the momenta are found by splitting the Hamilton–Jacobi equation $H - E = 0$ (eq. 3.205) into two parts, each of which is a function of only one coordinate and its conjugate momentum. In the case of the harmonic oscillator, $0 = H - E = \frac{1}{2}|\mathbf{p}|^2 + \frac{1}{2}(\omega_x^2 x^2 + \omega_y^2 y^2) - E = H_x(p_x, x) + H_y(p_y, y) - E$. In the case of motion in an axisymmetric potential, $0 = r^2(H - E) = r^2[\frac{1}{2}p_r^2 + \Phi(r)] - r^2 E + \frac{1}{2}p_\phi^2$.

The first of these observations suggests that we look for a curvilinear coordinate system whose coordinate curves run parallel to the edges of numerically integrated orbits, such as those plotted in Figure 3.4. Figure 3.27 shows that

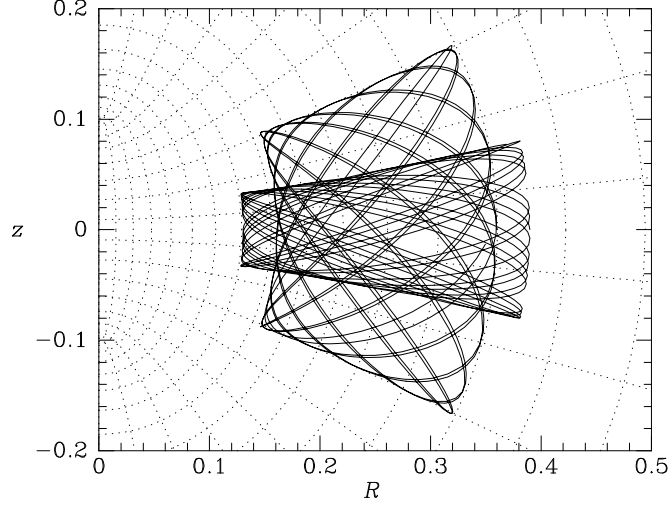


Figure 3.27 The boundaries of orbits in the meridional plane approximately coincide with the coordinate curves of a system of spheroidal coordinates. The dotted lines are the coordinate curves of the system defined by (3.242) and the full curves show the same orbits as Figure 3.4.

the (u, v) coordinate system defined by

$$R = \Delta \sinh u \sin v \quad ; \quad z = \Delta \cosh u \cos v \quad (3.242)$$

achieves this goal to high accuracy: the orbits of Figure 3.4 can be approximately bounded top and bottom by curves of constant v and right and left by curves of constant u .²⁰

Now that we have chosen a coordinate system, item (iii) above suggests that we next write the Hamiltonian function in terms of u , v , and their conjugate momenta. The first step is to write the Lagrangian as a function of the new coordinates and their time derivatives. By an analysis that closely parallels the derivation of equations (2.99) we may show that

$$|\dot{\mathbf{x}}|^2 = \Delta^2 (\sinh^2 u + \sin^2 v) (\dot{u}^2 + \dot{v}^2) + \Delta^2 \sinh^2 u \sin^2 v \dot{\phi}^2, \quad (3.243)$$

and the Lagrangian is

$$\mathcal{L} = \frac{1}{2} \Delta^2 \left[(\sinh^2 u + \sin^2 v) (\dot{u}^2 + \dot{v}^2) + \sinh^2 u \sin^2 v \dot{\phi}^2 \right] - \Phi(u, v). \quad (3.244)$$

The momenta are (eq. D.49)

$$\begin{aligned} p_u &= \Delta^2 (\sinh^2 u + \sin^2 v) \dot{u} \quad ; \quad p_v = \Delta^2 (\sinh^2 u + \sin^2 v) \dot{v} \\ p_\phi &= \Delta^2 \sinh^2 u \sin^2 v \dot{\phi}, \end{aligned} \quad (3.245)$$

²⁰ Note that *prolate* spheroidal coordinates are used to fit the boundaries of orbits in *oblate* potentials.

so the Hamiltonian is

$$H(u, v, p_u, p_v, p_\phi) = \frac{p_u^2 + p_v^2}{2\Delta^2(\sinh^2 u + \sin^2 v)} + \frac{p_\phi^2}{2\Delta^2 \sinh^2 u \sin^2 v} + \Phi(u, v). \quad (3.246)$$

Since H has no explicit dependence on time, it is equal to some constant E . Likewise, since H is independent of ϕ , the azimuthal momentum p_ϕ is constant at some value L_z .

The examples of motion in harmonic and circular potentials suggest that we seek a form of $\Phi(u, v)$ that will enable us to split a multiple of the equation $H(u, v, p_u, p_v, L_z) = E$ into a part involving only u and p_u and a part that involves only v and p_v . Evidently we require that $(\sinh^2 u + \sin^2 v)\Phi$ be of the form $U(u) - V(v)$, i.e., that²¹

$$\Phi(u, v) = \frac{U(u) - V(v)}{\sinh^2 u + \sin^2 v}, \quad (3.247)$$

for then we may rewrite $H = E$ as

$$2\Delta^2 [E \sinh^2 u - U(u)] - p_u^2 - \frac{L_z^2}{\sinh^2 u} = \frac{L_z^2}{\sin^2 v} + p_v^2 - 2\Delta^2 [E \sin^2 v + V(v)]. \quad (3.248)$$

It can be shown that potentials of the form (3.247) are generated by bodies resembling real galaxies (see Problems 2.6 and 2.14), so there are interesting physical systems for which (3.248) is approximately valid. Potentials of this form are called **Stäckel potentials** after the German mathematician P. Stäckel.²² Our treatment of these potentials will be restricted; much more detail, including the generalization to triaxial potentials, can be found in de Zeeuw (1985).

If the analogy with the harmonic oscillator holds, p_u will be a function only of u , and similarly for p_v . Under these circumstances, the left side of equation (3.248) does not depend on v , and the right side does not depend on u , so both sides must equal some constant, say $2\Delta^2 I_3$. Hence we would then have

$$p_u = \pm \sqrt{2\Delta^2 [E \sinh^2 u - I_3 - U(u)] - \frac{L_z^2}{\sinh^2 u}}, \quad (3.249a)$$

$$p_v = \pm \sqrt{2\Delta^2 [E \sin^2 v + I_3 + V(v)] - \frac{L_z^2}{\sin^2 v}}. \quad (3.249b)$$

²¹ The denominator of equation (3.247) vanishes when $u = 0$, $v = 0$. However, we may avoid an unphysical singularity in Φ at this point by choosing U and V such that $U(0) = V(0)$.

²² Stäckel showed that the *only* coordinate system in which the Hamilton–Jacobi equation for $H = \frac{1}{2}p^2 + \Phi(\mathbf{x})$ separates is confocal ellipsoidal coordinates. The usual Cartesian, spherical and cylindrical coordinate systems are limiting cases of these coordinates, as is the (u, v, ϕ) system.

It is a straightforward exercise to show that the analogy with the harmonic oscillator *does* hold, by direct time differentiation of both sides of equations (3.249), followed by elimination of \dot{u} and \dot{p}_u with Hamilton's equations (Problem 3.37). Thus the quantity I_3 defined by equations (3.249) is an integral. Moreover, we can display I_3 as an explicit function of the phase-space coordinates by eliminating E between equations (3.249) (Problem 3.39).

Equations (3.249) enable us to obtain expressions for the actions J_u and J_v in terms of the integrals E , I_3 and L_z , the last of which is equal to J_ϕ as in the spherical case. Specifically

$$\begin{aligned} J_u &= \frac{1}{\pi} \int_{u_{\min}}^{u_{\max}} du \sqrt{2\Delta^2 [E \sinh^2 u - I_3 - U(u)] - \frac{L_z^2}{\sinh^2 u}}, \\ J_v &= \frac{1}{\pi} \int_{v_{\min}}^{v_{\max}} dv \sqrt{2\Delta^2 [E \sin^2 v + I_3 + V(v)] - \frac{L_z^2}{\sin^2 v}}, \\ J_\phi &= L_z, \end{aligned} \quad (3.250)$$

where u_{\min} and u_{\max} are the smallest and largest values of u at which the integrand vanishes, and similarly for v_{\min} and v_{\max} .

As in the spherical case, we obtain expressions for the angle variables by differentiating the generating function $S(u, v, \phi, J_u, J_v, J_\phi)$ of the canonical transformation between angle-action variables and the (u, v, ϕ) system. We take S to be the sum of three parts S_u , S_v and S_ϕ , each of which depends on only one of the three coordinate variables. The gradient of S_u with respect to u is just p_u , so S_u is just the indefinite integral with respect to u of (3.249a). After evaluating S_v and S_ϕ analogously, we use the chain rule to differentiate $S = \sum_i S_i$ with respect to the actions (cf. the derivation of eq. 3.234):

$$\begin{aligned} \theta_u &= \frac{\partial S}{\partial J_u} = \sum_{i=u,v} \left(\frac{\partial S_i}{\partial H} \Omega_u + \frac{\partial S_i}{\partial I_3} \frac{\partial I_3}{\partial J_u} \right), \\ \theta_v &= \sum_{i=u,v} \left(\frac{\partial S_i}{\partial H} \Omega_v + \frac{\partial S_i}{\partial I_3} \frac{\partial I_3}{\partial J_v} \right), \\ \theta_\phi &= \sum_{i=u,v} \left(\frac{\partial S_i}{\partial H} \Omega_\phi + \frac{\partial S_i}{\partial I_3} \frac{\partial I_3}{\partial L_z} \right) + \phi. \end{aligned} \quad (3.251)$$

The partial derivatives in these expressions are all one-dimensional integrals that must in general be done numerically.

The condition (3.247) that must be satisfied by an axisymmetric Stäckel potential is very restrictive because it requires that a function of two variables can be written in terms of two functions of one variable. Most potentials that admit a third integral do not satisfy this condition. In particular the logarithmic potential Φ_L (2.71a) that motivated our discussion is not of Stäckel form: we can find a system of spheroidal coordinates that approximately bounds

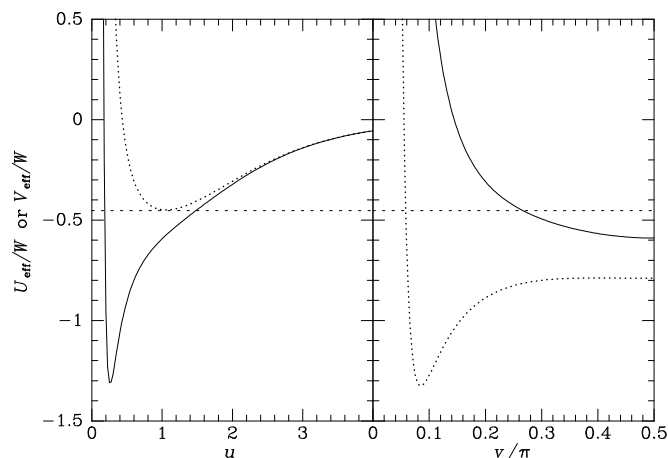


Figure 3.28 Plots of the effective potentials U_{eff} (left) and V_{eff} (right) that are defined by equations (3.252) and (3.253) for $\Delta = 0.6a_3$ and $L_z = 0.05a_3\sqrt{W}$. Curves are shown for $I_3 = -0.1W$ (full) and $I_3 = 0.1W$ (dotted).

any given orbit, but in general different orbits require different coordinate systems.

As an example of the use of equations (3.249) we investigate the shapes they predict for orbits in the potential obtained by choosing in (3.247)

$$U(u) = -W \sinh u \tan^{-1} \left(\frac{\Delta \sinh u}{a_3} \right) \quad (3.252)$$

$$V(v) = W \sin v \tanh^{-1} \left(\frac{\Delta \sin v}{a_3} \right),$$

where W , Δ , and a_3 are constants.²³ An orbit of specified E and I_3 can explore all values of u and v for which equations (3.249) predict positive p_u^2 and p_v^2 . This they will do providing E is larger than the largest of the “effective potentials”

$$U_{\text{eff}}(u) \equiv \frac{L_z^2}{2\Delta^2 \sinh^4 u} + \frac{I_3 + U(u)}{\sinh^2 u}, \quad (3.253a)$$

$$V_{\text{eff}}(v) \equiv \frac{L_z^2}{2\Delta^2 \sin^4 v} - \frac{I_3 + V(v)}{\sin^2 v}. \quad (3.253b)$$

Figure 3.28 shows these potentials for two values of I_3 and all other parameters fixed. Consider the case in which the energy takes the value $-0.453W$

²³ With these choices for U and V , the potential (3.247) becomes the potential of the perfect prolate spheroid introduced in Problem 2.14.

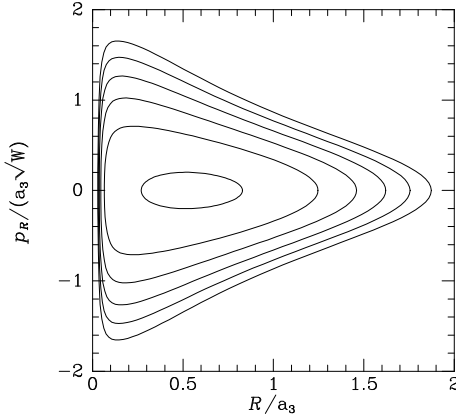


Figure 3.29 Surface of section at $E = -0.5W$ and $L_z = 0.05a_3\sqrt{W}$ constructed from equations (3.249) and (3.252) with $\Delta = 0.6a_3$.

(dashed horizontal line). Then for $I_3 = 0.1W$ (dotted curves), only a single value of u ($u = 1$) is permitted, so the orbit is confined to a segment of an ellipse in the meridional plane—this is a shell orbit. By contrast all values of $|v|$ larger than the intersection of the dashed and dotted curves in the right panel are permitted: these start at $|v| = 0.059\pi$. Consequently, the orbit covers much of the ellipse $u = 1$ (which in three dimensions is a spheroid).

Consider now the case in which $I_3 = -0.1W$ (full curves in Figure 3.28). Now a wide range is permitted in u ($0.17 < u < 1.48$) and a smaller range in v ($|v| > 0.27\pi$). Physically, lowering I_3 transfers some of the available energy from motion perpendicular to the potential's equatorial plane into the star's radial oscillation.

In §3.2.2 we detected the existence of non-classical integrals by plotting surfaces of section. It is interesting to see how I_3 structures surfaces of section. If we were to plot the (u, p_u) surface of section, the consequents of a given orbit (definite values of E, L_z, I_3) would lie on the curve in the (u, p_u) plane whose equation is (3.249a). This equation is manifestly independent of v , so the surface of section would look the same regardless of whether it was for $v = 0$, $v = 0.1$, or whatever. To get the structure of the (R, p_R) surfaces of section that we plotted in §3.2.2, for each allowed value of u we get $p(u)$ from (3.249a) and $p(v)$ from (3.249b) with $v = \pi/2$, and then obtain (R, p_R) from the (u, v, p_u, p_v) coordinates by inverting the transformations (2.96) and (3.245). Figure 3.29 shows a surface of section generated in this way.

In §3.2.1 we saw that motion in the meridional plane is governed by a Hamiltonian $H(R, z, p_R, p_z)$ in which L_z occurs as a parameter and the phase space is four-dimensional. In this space the orbital tori are ordinary two-dimensional doughnuts, and a surface of section is simply a cross-section through a nested sequence of such tori: each invariant curve marks the intersection of a two-dimensional doughnut with the two-dimensional surface of section.

(b) Epicycle approximation In §3.2.3 we used the epicycle approxima-

tion to obtain solutions to the equations of motion that are approximately valid for nearly circular orbits in an axisymmetric potential. Here we obtain the corresponding approximate angle-action variables. In cylindrical coordinates the Hamilton–Jacobi equation (3.205) is

$$\frac{1}{2}\left(\frac{\partial S}{\partial R}\right)^2 + \frac{1}{2R^2}\left(\frac{\partial S}{\partial \phi}\right)^2 + \frac{1}{2}\left(\frac{\partial S}{\partial z}\right)^2 + \Phi(R, z) = E. \quad (3.254)$$

As in equation (2.75a) we assume that Φ is of the form $\Phi_R(R) + \Phi_z(z)$; the radial dependence of $\Phi_z(z)$ is suppressed because the radial motion is small in the epicycle approximation. We further assume that S is of the form $S(\mathbf{J}, R, \phi, z) = S_R(\mathbf{J}, R) + S_\phi(\mathbf{J}, \phi) + S_z(\mathbf{J}, z)$. Now we use the method of separation of variables to split equation (3.254) up into three parts:

$$\begin{aligned} E_z &= \frac{1}{2}\left(\frac{\partial S_z}{\partial z}\right)^2 + \Phi_z(z) \quad ; \quad L_z^2 = \left(\frac{\partial S_\phi}{\partial \phi}\right)^2 \\ E - E_z &= \frac{1}{2}\left(\frac{\partial S_R}{\partial R}\right)^2 + \Phi_R(R) + \frac{L_z^2}{2R^2}, \end{aligned} \quad (3.255)$$

where E_z and L_z are the two constants of separation. The first equation of this set leads immediately to an integral for $S_z(z)$

$$S_z(z) = \int_0^z dz' \epsilon_z \sqrt{2[E_z - \Phi_z(z')]}, \quad (3.256)$$

where ϵ_z is chosen to be ± 1 such that the integral increases monotonically along the path. If, as in §3.2.3, we assume that $\Phi_z = \frac{1}{2}\nu^2 z^2$, where ν is a constant, then our equation for S_z becomes essentially the same as the first of equations (3.211), and by analogy with equations (3.213) and (3.216), we have

$$J_z = \frac{E_z}{\nu} \quad ; \quad z = -\sqrt{\frac{2J_z}{\nu}} \cos \theta_z. \quad (3.257)$$

The second of equations (3.255) trivially yields

$$S_\phi(\mathbf{J}, \phi) = L_z \phi, \quad (3.258)$$

and it immediately follows that $J_\phi = L_z$. The last of equations (3.255) yields

$$2(E - E_z) = \left(\frac{\partial S_R}{\partial R}\right)^2 + 2\Phi_{\text{eff}}(R), \quad (3.259a)$$

where (cf. eq. 3.68b)

$$\Phi_{\text{eff}}(R) \equiv \Phi_R(R) + \frac{J_\phi^2}{2R^2}. \quad (3.259b)$$

The epicycle approximation involves expanding Φ_{eff} about its minimum, which occurs at the radius $R_g(J_\phi)$ of the circular orbit of angular momentum J_ϕ ; with x defined by $R = R_g + x$, the expansion is $\Phi(R) = E_c(J_\phi) + \frac{1}{2}\kappa^2 x^2$, where $E_c(J_\phi)$ is the energy of the circular orbit of angular momentum J_ϕ and κ is the epicycle frequency defined in equation (3.77). Inserting this expansion into (3.259a) and defining $E_R \equiv E - E_z - E_c$, we have

$$2E_R = \left(\frac{\partial S_R}{\partial R}\right)^2 + \kappa^2 x^2, \quad (3.260)$$

which is the same as equation (3.210) with K^2 replaced by $2E_R$, x by R , and ω_x by κ . It follows from equations (3.213), (3.216) and (3.217) that

$$J_R = \frac{E_R}{\kappa} \quad ; \quad S_R(\mathbf{J}, R) = J_R(\theta_R - \frac{1}{2} \sin 2\theta_R) \quad ; \quad R = R_g - \sqrt{\frac{2J_R}{\kappa}} \cos \theta_R. \quad (3.261)$$

The last of these equations is equivalent to equation (3.91) if we set $\theta_R = \kappa t + \alpha$ and $X = -(2J_R/\kappa)^{1/2}$.

Finally, we find an expression for θ_ϕ . With equations (3.258) and (3.261) we have

$$\begin{aligned} \theta_\phi &= \frac{\partial S}{\partial J_\phi} = \frac{\partial S_\phi}{\partial J_\phi} + \frac{\partial S_R}{\partial J_\phi} = \phi + J_R(1 - \cos 2\theta_R) \frac{\partial \theta_R}{\partial J_\phi} \\ &= \phi + 2J_R \sin^2 \theta_R \frac{\partial \theta_R}{\partial J_\phi}. \end{aligned} \quad (3.262)$$

The derivative of θ_R has to be taken at constant J_R, J_z, R, ϕ , and z . We differentiate the last of equations (3.261) bearing in mind that both R_g and κ are functions of J_ϕ :

$$0 = \frac{dR_g}{dJ_\phi} + \frac{1}{2\kappa} \frac{d\kappa}{dJ_\phi} \sqrt{\frac{2J_R}{\kappa}} \cos \theta_R + \sqrt{\frac{2J_R}{\kappa}} \sin \theta_R \frac{\partial \theta_R}{\partial J_\phi}. \quad (3.263)$$

By differentiating $R_g^2 \Omega_g = J_\phi$ with respect to R_g we may show with equation (3.80) that

$$\frac{dR_g}{dJ_\phi} = \frac{\gamma}{\kappa R_g}, \quad (3.264)$$

where $\gamma = 2\Omega_g/\kappa$ is defined by equation (3.93b). Inserting this relation into (3.263) and using the result to eliminate $\partial \theta_R / \partial J_\phi$ from (3.262), we have finally

$$\theta_\phi = \phi - \frac{\gamma}{R_g} \sqrt{\frac{2J_R}{\kappa}} \sin \theta_R - \frac{J_R}{2} \frac{d \ln \kappa}{dJ_\phi} \sin 2\theta_R. \quad (3.265)$$

This expression should be compared with equation (3.93a). If we set $\theta_\phi = \Omega_g t + \phi_0$, $\theta_R = \kappa t + \alpha + \pi$, and $X = (2J_R/\kappa)^{1/2}$ as before, the only difference

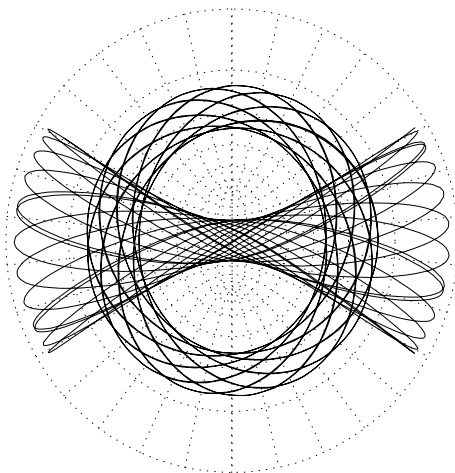


Figure 3.30 The boundaries of loop and box orbits in barred potentials approximately coincide with the curves of a system of spheroidal coordinates. The figure shows two orbits in the potential Φ_L of equation (3.103), and a number of curves on which the coordinates u and v defined by equations (3.267) are constant.

between the two equations is the presence of a term proportional to $\sin 2\theta_R$ in equation (3.265). For nearly circular orbits, this term is smaller than the term proportional to $\sin \theta_R$ by $\sqrt{J_R/J_\phi}$ and represents a correction to equation (3.92) that makes the $(\theta_R, \theta_\phi, J_R, J_\phi)$ coordinates canonical (Dehnen 1999a).

It is worth noting that when $J_R \neq 0$, the frequency associated with ϕ is not the circular frequency, Ω_g . To see this, recall that the Hamiltonian $H = E_R + E_c + E_z$, and $E_R = \kappa J_R$, while $dE_c/dJ_\phi = \Omega_g$, so

$$\Omega_\phi = \frac{\partial H}{\partial J_\phi} = \frac{d\kappa}{dJ_\phi} J_R + \Omega_g. \quad (3.266)$$

3.5.4 Angle-action variables for a non-rotating bar

The (u, v) coordinate system that allowed us to recover angle-action variables for flattened axisymmetric potentials enables us to do the same for a planar, non-rotating bar. This fact is remarkable, because we saw in §3.3 that these systems support two completely different types of orbit, loops and boxes. Figure 3.30 makes it plausible that the (u, v) system can provide analytic solutions for both loops and boxes, by showing that the orbits plotted in Figure 3.8 have boundaries that may be approximated by curves of constant u and v (cf. the discussion on page 226). We can explore this idea quantitatively by defining

$$x = \Delta \sinh u \sin v \quad ; \quad y = \Delta \cosh u \cos v \quad (3.267)$$

and then replacing R by x and z by y in the formulae of the previous subsection. Further setting $\phi = L_z = 0$ we find by analogy with equations (3.249) that

$$p_u = \pm \Delta \sinh u \sqrt{2[E - U_{\text{eff}}(u)]} \quad ; \quad p_v = \pm \Delta \sin v \sqrt{2[E - V_{\text{eff}}(v)]} \quad (3.268a)$$

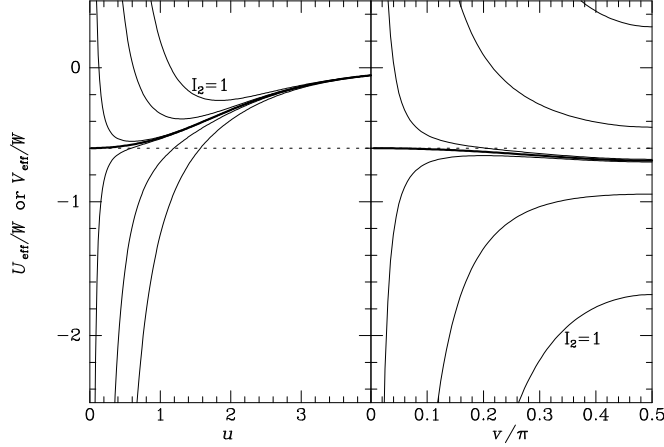


Figure 3.31 The effective potentials defined by equations (3.268b) when U and V are given by equations (3.252). The curves are for $I_2 = 1, 0.25, 0.01, 0, -0.01, -0.25$ and -1 , with the largest values coming on top in the left panel and on the bottom in the right panel. The thick curves are for $I_2 = 0$.

where

$$U_{\text{eff}}(u) = \frac{I_2 + U(u)}{\sinh^2 u} \quad ; \quad V_{\text{eff}}(v) = -\frac{I_2 + V(v)}{\sin^2 v}. \quad (3.268b)$$

Here U and V are connected to the gravitational potential by equation (3.247) as before and I_2 is the constant of separation analogous to I_3 .

An orbit of specified E and I_2 is confined to values of u and v at which both $E \geq U_{\text{eff}}$ and $E \geq V_{\text{eff}}$. Figure 3.31 shows the effective potentials as functions of their coordinates for several values of I_2 when U and V are chosen to be the functions specified by equations (3.252). In each panel the thick curve is for $I_2 = 0$, with curves for $I_2 > 0$ lying above this in the left panel, and below it on the right. Since the curves of U_{eff} have minima only when $I_2 > 0$, there is a lower limit on the star's u coordinate only in this case. Consequently, stars with $I_2 \leq 0$ can reach the center, while stars with $I_2 > 0$ cannot reach the center. This suggests that when $I_2 \leq 0$ the orbit is a box orbit, while when $I_2 > 0$ it is a loop orbit. Comparison of the right and left panels confirms this conjecture by showing that when $I_2 > 0$ (upper curves on left and lower curves on right), the minimum value of U_{eff} is greater than the maximum of V_{eff} . Hence when $I_2 > 0$ the condition $E > V_{\text{eff}}(v)$ imposes no constraint on v and the boundaries of the orbit are the ellipses $u = u_{\text{min}}$ and $u = u_{\text{max}}$ on which $E = U_{\text{eff}}$. When $I_2 \leq 0$, by contrast, the curves on the right tend to ∞ as $v \rightarrow 0$, so sufficiently small values of v are excluded and the boundaries of the orbit are the ellipse $u = u_{\text{max}}$ on which $E = U_{\text{eff}}(u)$ and the hyperbola $|v| = v_{\text{min}}$ on which $E = V_{\text{eff}}(v)$.

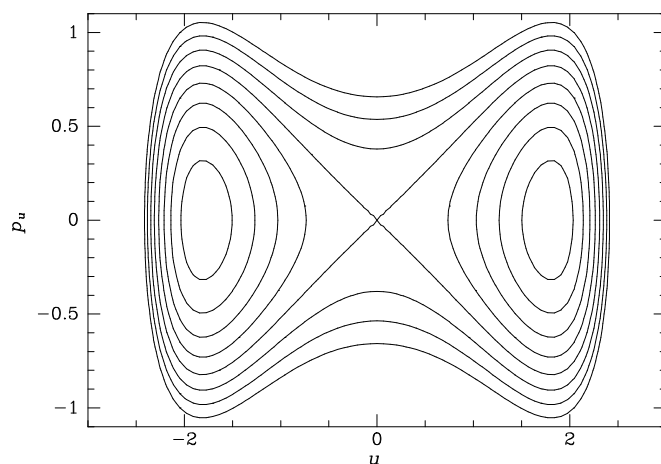


Figure 3.32 The (u, p_u) , $v = 0$ surface of section for motion at $E = -0.25$ in the Stäckel potential defined by equations (3.247) and (3.252) with $\Delta = 0.6$ and $a_3 = 1$. Each curve is a contour of constant I_2 (eqs. 3.268). The invariant curves of box orbits ($I_2 = -0.6, -0.4, \dots$) run round the outside of the figure, while the bull's-eyes at right are the invariant curves of anti-clockwise loop orbits. Temporarily suspending the convention that loops always have $u > 0$, we show the invariant curves of clockwise loops as the bull's-eyes at left.

Figure 3.32 shows the (u, p_u) surface of section, which is in practice nothing more than a contour plot of the integral $I_2(E, u, p_u)$ with E fixed (eq. 3.268a). Each contour shows the curve in which an orbital torus is sliced by the surface of section. As in Figure 3.9, for example, there are two different types of contour, namely those generated by the tori of loop orbits (which come in pairs, because there are both clockwise and anti-clockwise circulating loops), and those generated by the tori of box orbits, which envelop all the tori of the loop orbits.

3.5.5 Summary

We have made a considerable investment in the theory of angle-action variables, which is repaid by the power of these variables in investigations of a wide variety of dynamical problems. This power arises from the following features:

- (i) Angle-action variables are canonical. In particular, the phase-space volume $d^3\theta d^3\mathbf{J}$ is the same as the phase-space volume $d^3\mathbf{q}d^3\mathbf{p}$ for any other set of canonical variables (\mathbf{q}, \mathbf{p}) , including the usual Cartesian coordinates (\mathbf{x}, \mathbf{v}) .
- (ii) Every set of angle-action variables $(\boldsymbol{\theta}, \mathbf{J})$ is associated with a Hamiltonian²⁴ $H(\mathbf{J})$, and orbits in this Hamiltonian have the simple form

²⁴If a given set of angle-action variables is associated with $H(\mathbf{J})$, then it is also as-

$\mathbf{J} = \text{constant}$, $\boldsymbol{\theta} = \boldsymbol{\Omega}t + \text{constant}$. A Hamiltonian that admits angle-action variables is said to be **integrable**. The simplicity of angle-action variables makes them indispensable for investigating motion in non-integrable Hamiltonians by using perturbation theory. This technique will be used to explore chaotic orbits in §3.7, and the stability of stellar systems in Chapter 5.

- (iii) In the next section we shall see that actions are usually invariant during slow changes in the Hamiltonian.

3.6 Slowly varying potentials

So far we have been concerned with motion in potentials that are time-independent in either an inertial or a rotating frame. It is sometimes necessary to consider how stars move in potentials that are time-dependent. The nature of the problem posed by a time-varying potential depends on the speed with which the potential evolves. In this section we shall confine ourselves to potentials that evolve slowly, in which case angle-action variables enable us to predict how a stellar system will respond to changes in the gravitational field that confines it. Such changes occur when:

- (i) Encounters between the individual stars at the core of a dense stellar system (such as a globular cluster or galaxy center) cause the core to evolve on a timescale of order the relaxation time (1.38), which is much longer than the orbital times of individual stars (§7.5).
- (ii) Stars of galaxies and globular clusters lose substantial quantities of mass as they gradually evolve and shed their envelopes into interstellar or intergalactic space (Box 7.2).
- (iii) Gas settles into the equatorial plane of a pre-existing dark halo to form a spiral galaxy. In this case the orbits of the halo's dark-matter particles will undergo a slow evolution as the gravitational potential of the disk gains in strength.

Potential variations that are slow compared to a typical orbital frequency are called **adiabatic**. We now show that the actions of stars are constant during such adiabatic changes of the potential. For this reason actions are often called **adiabatic invariants**.

3.6.1 Adiabatic invariance of actions

Suppose we have a sequence of potentials $\Phi_\lambda(\mathbf{x})$ that depend continuously on the parameter λ . For each fixed λ we assume that angle-action variables could be constructed for Φ_λ . That is, we assume that at all times phase space is filled by arrays of nested tori on which the phase points of individual stars

sociated with $\tilde{H}(\mathbf{J}) \equiv f[H(\mathbf{J})]$, where f is any differentiable function. Thus, a set of angle-action variables is associated with *infinitely many* Hamiltonians.

travel. We consider what happens when λ is changed from its initial value, say $\lambda = \lambda_0$, to a new value λ_1 . After this change has occurred, each star's phase point will start to move on a torus of the set that belongs to Φ_{λ_1} . In general, two stellar phase points that started out on the same torus of Φ_{λ_0} will end up on two different tori of Φ_{λ_1} . But if λ is changed very slowly compared to all the characteristic times $2\pi/\Omega_k$ associated with motion on each torus, all phase points that are initially on a given torus of Φ_{λ_0} will be equally affected by the variation of λ . This statement follows from the time averages theorem of §3.5.1a, which shows that all stars spend the same fraction of their time in each portion of the torus; hence, all stars are affected by slow changes in Φ_λ in the same way. Thus all phase points that start on the same torus of Φ_{λ_0} will end on a single torus of Φ_{λ_1} . Said in other language, any two stars that are initially on a common orbit (but at different phases) will still be on a common orbit after the slow variation of λ is complete.

Suppose the variation of λ starts at time $t = 0$ and is complete by time τ , and let \mathbf{H}_t be the time-evolution operator defined in equation (D.55). Then we have just seen that \mathbf{H}_τ , which is a canonical map (see Appendix D.4.4), maps tori of Φ_{λ_0} onto tori of Φ_{λ_1} . These facts guarantee that actions are adiabatically invariant, for the following reason. Choose three closed curves γ_i , on any torus M of Φ_{λ_0} that through the integrals (3.195) generate the actions J_i of this torus. Then, since \mathbf{H}_τ is the endpoint of a continuous deformation of phase space into itself, the images $\mathbf{H}_\tau(\gamma_i)$ of these curves are suitable curves along which to evaluate the actions J'_i of $\mathbf{H}_\tau(M)$, the torus to which M is mapped by \mathbf{H}_τ . But by a corollary to the Poincaré invariant theorem (Appendix D.4.2), we have that if γ is any closed curve and $\mathbf{H}_\tau(\gamma)$ is its image under the canonical map \mathbf{H}_τ , then

$$\oint_{\mathbf{H}_\tau(\gamma)} \mathbf{p} \cdot d\mathbf{q} = \oint_\gamma \mathbf{p} \cdot d\mathbf{q}. \quad (3.269)$$

Hence $J'_i = J_i$, and the actions of stars do not change if the potential evolves sufficiently slowly.

It should be stressed that any action J_i with fundamental frequency $\Omega_i = 0$ is not an adiabatic invariant. For example, in a spherical potential, J_2 and J_3 are normally adiabatic invariants, but J_1 is not (Table 3.1).

3.6.2 Applications

We illustrate these ideas with a number of simple examples. Other applications of adiabatic invariants will be found in Binney & May (1986), Lichtenberg & Lieberman (1992), and §4.6.1.

(a) Harmonic oscillator We first consider the one-dimensional harmonic oscillator whose potential is

$$\Phi = \frac{1}{2}\omega^2 x^2. \quad (3.270)$$

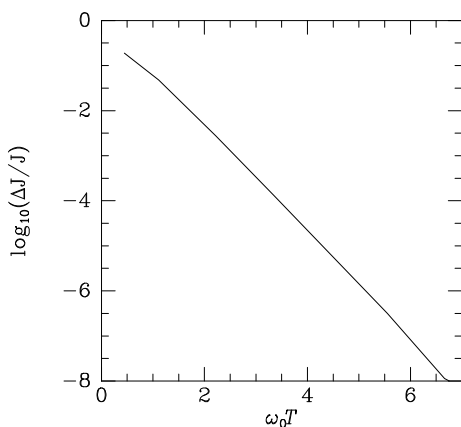


Figure 3.33 Checking the invariance of the action (3.271) when the natural frequency of a harmonic oscillator is varied according to equation (3.277). ΔJ is the RMS change in the action on integrating the oscillator's equation of motion from $t = -20T$ to $t = 20T$, using eight equally spaced phases. The RMS change in J declines approximately as $\Delta J \propto \exp(-2.8\omega_0 T)$.

By equation (3.213) the action is

$$J = \frac{1}{2\omega} [p^2 + (\omega x)^2] = \frac{H}{\omega}, \quad (3.271)$$

where $H(x, p) = \frac{1}{2}p^2 + \frac{1}{2}\omega^2 x^2$. The general solution of the equations of motion is $x(t) = X \cos(\omega t + \phi)$. In terms of the amplitude of oscillation X we have

$$J = \frac{1}{2}\omega X^2. \quad (3.272)$$

Now suppose that the oscillator's spring is slowly stiffened by a factor $s^2 > 1$, so the natural frequency increases to

$$\omega' = s\omega. \quad (3.273)$$

By the adiabatic invariance of J , the new amplitude X' satisfies

$$\frac{1}{2}\omega' X'^2 = J = \frac{1}{2}\omega X^2. \quad (3.274)$$

Thus the amplitude is diminished to

$$X' = \frac{X}{\sqrt{s}}, \quad (3.275)$$

while the energy, $E = \omega J$, has increased to²⁵

$$E' = \omega' J = s\omega J = sE. \quad (3.276)$$

²⁵ The simplest proof of this result uses quantum mechanics. The energy of a harmonic oscillator is $E = (n + \frac{1}{2})\hbar\omega$ where n is an integer. When ω is slowly varied, n cannot change discontinuously and hence must remain constant. Therefore $E/\omega = E'/\omega'$. Of course, for galaxies n is rather large.

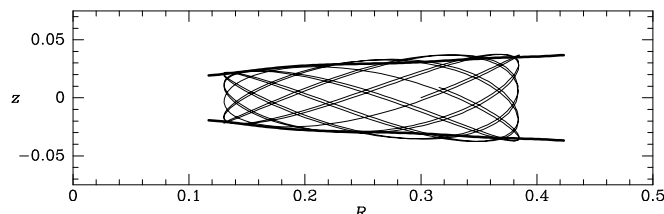


Figure 3.34 The envelope of an orbit in the effective potential (3.70) with $q = 0.5$ (light curve) is well modeled by equation (3.279) (heavy curves).

We now ask how rapidly we can change the frequency ω without destroying the invariance of J . Let ω vary with time according to

$$\omega(t) = \pi\sqrt{3 + \operatorname{erf}(t/T)}. \quad (3.277)$$

Thus the frequency changes from $\omega = \omega_0 \equiv \sqrt{2}\pi$ at $t \ll -T$ to $\omega = 2\pi = \sqrt{2}\omega_0$ at $t \gg T$. In Figure 3.33 we show the results of numerically integrating the oscillator's equation of motion with $\omega(t)$ given by equation (3.277). We plot the RMS difference ΔJ between the initial and final values of J for eight different phases of the oscillator at $t = -20T$. For $\omega_0 T \gtrsim 2$, J changes by less than half a percent, and for $\omega_0 T \gtrsim 4$, J changes by less than 3×10^{-5} . We conclude that the potential does not have to change very slowly for J to be well conserved. In fact, one can show that the fractional change in J is in general less than $\exp(-\omega T)$ for $\omega T \gg 1$ (Lichtenberg & Lieberman 1992).

(b) Eccentric orbits in a disk Consider the shapes shown in Figure 3.4 of the orbits in the meridional plane of an axisymmetric galaxy. On page 167 we remarked that disk stars in the solar neighborhood oscillate perpendicular to the galactic plane considerably more rapidly than they oscillate in the radial direction. Therefore, if we take the radial coordinate $R(t)$ of a disk star to be a known function of time, we may consider the equation of motion (3.67c) of the z -coordinate to describe motion in a slowly varying potential. If the amplitude of the z -oscillations is small, we may expand $\partial\Phi/\partial z$ about $z = 0$ to find

$$\ddot{z} \simeq -\omega^2 z \quad \text{where} \quad \omega(t) \equiv \left(\frac{\partial^2 \Phi}{\partial z^2} \right)_{[R(t), 0]}^{1/2} \equiv \sqrt{\Phi_{zz}[R(t), 0]}. \quad (3.278)$$

If the action integral of this harmonic oscillator is conserved, we expect the amplitude $Z(R)$ to satisfy (see eqs. 3.273 and 3.275)

$$Z(R) = Z(R_0) \left(\frac{\Phi_{zz}(R_0, 0)}{\Phi_{zz}(R, 0)} \right)^{1/4}. \quad (3.279)$$

Figure 3.34 compares the prediction of (3.279) with the true shape of an orbit in the effective potential (3.70). Evidently the behavior of such orbits can be accurately understood in terms of adiabatic invariants.

(c) Transient perturbations Consider the motion of a star on a loop

orbit in a slowly varying planar potential $\Phi(R, \phi)$. The relevant action is

$$J_\phi = \frac{1}{2\pi} \int_0^{2\pi} d\phi p_\phi. \quad (3.280)$$

We now conduct the following thought experiment. Initially the potential Φ is axisymmetric. Then $p_\phi = L_z$ is an integral, and we can trivially evaluate the integral in (3.280) to obtain $J_\phi = L_z$. We now slowly distort the potential in some arbitrary fashion into a new axisymmetric configuration. At the end of this operation, the azimuthal action, being adiabatically invariant, still has value J_ϕ and is again equal to the angular momentum L_z . Thus the star finishes the experiment with the same angular momentum with which it started,²⁶ even though its instantaneous angular momentum, p_ϕ , was changing during most of the experiment. Of course, if the potential remains axisymmetric throughout, p_ϕ remains an integral at all times and is exactly conserved no matter how rapidly the potential is varied.

A closely related example is a slowly varying external perturbation of a stellar system, perhaps from the gravitational field of an object passing at a low angular velocity. If the passage is slow enough, the actions are adiabatically invariant, so the distribution of actions in the perturbed system will be unchanged by the encounter. In other words, adiabatic encounters, even strong ones, have no lasting effect on a stellar system (§8.2c).

(d) Slow growth of a central black hole As our final application of the adiabatic invariance of actions, we consider the evolution of the orbit of a star near the center of a spherical galaxy, as a massive black hole grows by slowly accreting matter (Goodman & Binney 1984). A more complete treatment of the problem is given in §4.6.1d. We assume that prior to the formation of the hole, the density of material interior to the orbit can be taken to be a constant, so the potential is that of the spherical harmonic oscillator. It is then easy to show that the star's Hamiltonian can be written (Problem 3.36)

$$H = \Omega_r J_r + \Omega_\phi J_\phi = 2\Omega J_r + \Omega J_\phi, \quad (3.281)$$

where $\Omega = \Omega_\phi = \frac{1}{2}\Omega_r$ is the circular frequency, and $J_\phi = L$ is the magnitude of the angular-momentum vector. The radii r_{\min} and r_{\max} of peri- and apocenter are the roots of

$$0 = \frac{J_\phi^2}{2r^2} + \frac{1}{2}\Omega^2 r^2 - H \quad \Rightarrow \quad 0 = r^4 - \frac{2H}{\Omega^2} r^2 + \frac{J_\phi^2}{\Omega^2}. \quad (3.282)$$

²⁶ This statement does not apply for stars that switch from loop to box orbits and back again as the potential is varied (Binney & Spergel 1983; Evans & Collett 1994). These stars will generally be on highly eccentric orbits initially.

Hence, the axis ratio of the orbit is

$$q_H = \frac{r_{\min}}{r_{\max}} = \left(\frac{H/\Omega^2 - [(H/\Omega^2)^2 - (J_\phi/\Omega)^2]^{1/2}}{H/\Omega^2 + [(H/\Omega^2)^2 - (J_\phi/\Omega)^2]^{1/2}} \right)^{1/2} \quad (3.283)$$

$$= \left(\frac{2J_r + J_\phi - 2[J_r(J_r + J_\phi)]^{1/2}}{2J_r + J_\phi + 2[J_r(J_r + J_\phi)]^{1/2}} \right)^{1/2}.$$

Multiplying top and bottom of the fraction by the top, this last expression reduces to

$$q_H = \frac{1}{J_\phi} [2J_r + J_\phi - 2\sqrt{J_r(J_r + J_\phi)}]. \quad (3.284)$$

When the hole has become sufficiently massive, the Hamiltonian may be taken to be that for Kepler motion (eq. E.6) and the orbit becomes an ellipse with the black hole at the focus rather than the center of the ellipse. A similar calculation yields for the axis ratio of this ellipse

$$q_K = \left[1 - \left(\frac{r_{\max} - r_{\min}}{r_{\max} + r_{\min}} \right)^2 \right]^{1/2} = \frac{J_\phi}{J_r + J_\phi}. \quad (3.285)$$

When J_r/J_ϕ is eliminated between equations (3.284) and (3.285), we find

$$q_K = \frac{4q_H}{(1 + q_H)^2}. \quad (3.286)$$

For example, if $q_H = 0.5$ is the original axis ratio, the final one is $q_K = 0.889$, and if initially $q_H = 0.75$, then finally $q_K = 0.980$. Physically, an elongated ellipse that is centered on the black hole distorts into a much rounder orbit with the black hole at one focus.

For any orbit in a spherical potential the mean-square radial speed is

$$\overline{v_r^2} = \frac{\Omega_r}{\pi} \int_0^{\pi/\Omega_r} dt v_r^2 = \frac{\Omega_r}{\pi} \int_{r_{\min}}^{r_{\max}} dr v_r = \Omega_r J_r. \quad (3.287a)$$

Similarly, the mean-square tangential speed is

$$\overline{v_t^2} = \frac{\Omega_\phi}{2\pi} \int_0^{2\pi/\Omega_\phi} dt (R\dot{\phi})^2 = \frac{\Omega_\phi}{2\pi} \int_0^{2\pi} d\phi p_\phi = \Omega_\phi J_\phi. \quad (3.287b)$$

Since the actions do not change as the hole grows, the change in the ratio of the mean-square speeds is given by

$$\frac{(\overline{v_r^2}/\overline{v_t^2})_K}{(\overline{v_r^2}/\overline{v_t^2})_H} = \frac{(\Omega_r/\Omega_\phi)_K}{(\Omega_r/\Omega_\phi)_H} = \frac{1}{2}. \quad (3.288)$$

Consequently, the growth of the black hole increases the star's tangential velocity much more than it does the radial velocity, irrespective of the original eccentricity of the orbit. In §4.6.1a we shall investigate the implications of this result for measurements of the stellar velocity dispersion near a black hole, and show how the growth of the black hole enhances the density of stars in its vicinity.

3.7 Perturbations and chaos

Analytic solutions to a star's equations of motion exist for only a few simple potentials $\Phi(\mathbf{x})$. If we want to know how stars will move in a more complex potential, for example one estimated from observational data, two strategies are open to us: either solve the equations of motion numerically, or obtain an approximate analytic solution by invoking perturbation theory, which involves expressing the given potential as a sum of a potential for which we can solve the equations of motion analytically and a (one hopes) small additional term.

Even in the age of fast, cheap and convenient numerical computation, perturbative solutions to the equations of motion are useful in two ways. First, they can be used to investigate the stability of stellar systems (§5.3). Second, they give physical insight into the dynamics of orbits. We start this section by developing perturbation theory and sketching some of its astronomical applications; then we describe the phenomenon of orbital chaos, and show that Hamiltonian perturbation theory helps us to understand the physics of this phenomenon.

3.7.1 Hamiltonian perturbation theory

In §3.3.3 we derived approximate orbits in the potential of a weak bar, by treating the potential as a superposition of a small non-axisymmetric potential and a much larger axisymmetric one. Our approach involved writing the orbit $\mathbf{x}(t)$ as a sum of two parts, one of which described the circular orbit of a guiding center, while the other described epicyclic motion. We worked directly with the equations of motion. Angle-action variables enable us to develop a more powerful perturbative scheme, in which we work with scalar functions rather than coordinates, and think of the orbit as a torus in phase space rather than a time-ordered series of points along a trajectory. For more detail see Lichtenberg & Leiberman (1992).

Let H^0 be an integrable Hamiltonian, and consider the one-parameter family of Hamiltonians

$$H^\beta \equiv H^0 + \beta h, \quad (3.289)$$

where $\beta \ll 1$ and h is a Hamiltonian with gradients that are comparable in magnitude to those of H^0 . Let $(\boldsymbol{\theta}^\beta, \mathbf{J}^\beta)$ be angle-action variables for H^β . These coordinates are related to the angle-action variables of H^0 by a canonical transformation. As $\beta \rightarrow 0$ the generating function S (Appendix D.4.6) of this transformation will tend to the generating function of the identity transformation, so we may write

$$S(\boldsymbol{\theta}^\beta, \mathbf{J}^0) = \boldsymbol{\theta}^\beta \cdot \mathbf{J}^0 + s^\beta(\boldsymbol{\theta}^\beta, \mathbf{J}^0), \quad (3.290)$$

where s^β is $O(\beta)$, and (eq. D.94)

$$\mathbf{J}^\beta = \frac{\partial S}{\partial \boldsymbol{\theta}^\beta} = \mathbf{J}^0 + \frac{\partial s^\beta}{\partial \boldsymbol{\theta}^\beta} \quad ; \quad \boldsymbol{\theta}^0 = \boldsymbol{\theta}^\beta + \frac{\partial s^\beta}{\partial \mathbf{J}^0}. \quad (3.291)$$

Substituting these equations into (3.289), we have

$$\begin{aligned}
H^\beta(\mathbf{J}^\beta) &= H^0(\mathbf{J}^0) + \beta h(\boldsymbol{\theta}^0, \mathbf{J}^0) \\
&= H^0\left(\mathbf{J}^\beta - \frac{\partial s^\beta}{\partial \boldsymbol{\theta}^\beta}\right) + \beta h\left(\boldsymbol{\theta}^\beta + \frac{\partial s^\beta}{\partial \mathbf{J}^0}, \mathbf{J}^\beta - \frac{\partial s^\beta}{\partial \boldsymbol{\theta}^\beta}\right) \\
&= H^0(\mathbf{J}^\beta) - \boldsymbol{\Omega}^0(\mathbf{J}^\beta) \cdot \frac{\partial s^\beta}{\partial \boldsymbol{\theta}^\beta} + \beta h(\boldsymbol{\theta}^\beta, \mathbf{J}^\beta) + O(\beta^2),
\end{aligned} \tag{3.292}$$

where $\boldsymbol{\Omega}^0$ is the derivative of H^0 with respect to its argument. We next expand h and s^β as Fourier series in the periodic angle variables (Appendix B.4):

$$h(\boldsymbol{\theta}^\beta, \mathbf{J}^\beta) = \sum_{\mathbf{n}} h_{\mathbf{n}}(\mathbf{J}^\beta) e^{i\mathbf{n} \cdot \boldsymbol{\theta}^\beta}; \quad s^\beta(\boldsymbol{\theta}^\beta, \mathbf{J}^0) = i \sum_{\mathbf{n}} s_{\mathbf{n}}^\beta(\mathbf{J}^0) e^{i\mathbf{n} \cdot \boldsymbol{\theta}^\beta}, \tag{3.293}$$

where $\mathbf{n} = (n_1, n_2, n_3)$ is a triple of integers. Substituting these expressions into (3.292) we find

$$H^\beta(\mathbf{J}^\beta) = H^0(\mathbf{J}^\beta) + \beta h_0 + \sum_{\mathbf{n} \neq \mathbf{0}} \left(\beta h_{\mathbf{n}} + \mathbf{n} \cdot \boldsymbol{\Omega}^0 s_{\mathbf{n}}^\beta \right) e^{i\mathbf{n} \cdot \boldsymbol{\theta}^\beta} + O(\beta^2). \tag{3.294}$$

In this equation $\boldsymbol{\Omega}^0$ and $h_{\mathbf{n}}$ are functions of \mathbf{J}^β , while $s_{\mathbf{n}}$ is a function of \mathbf{J}^0 , but to the required order in β , \mathbf{J}^0 can be replaced by \mathbf{J}^β .

Since the left side of equation (3.294) does not depend on $\boldsymbol{\theta}^\beta$, on the right the coefficient of $\exp(i\mathbf{n} \cdot \boldsymbol{\theta}^\beta)$ must vanish for all $\mathbf{n} \neq 0$. Hence the Fourier coefficients of S are given by

$$s_{\mathbf{n}}^\beta(\mathbf{J}) = -\frac{\beta h_{\mathbf{n}}(\mathbf{J})}{\mathbf{n} \cdot \boldsymbol{\Omega}^0(\mathbf{J})} + O(\beta^2) \quad (\mathbf{n} \neq 0). \tag{3.295}$$

The $O(\beta)$ part of equation (3.295) defines the generating function of a canonical transformation. Let $(\boldsymbol{\theta}', \mathbf{J}')$ be the images of $(\boldsymbol{\theta}^0, \mathbf{J}^0)$ under this transformation. Then we have shown that

$$H^\beta(\mathbf{J}^\beta) = H'(\mathbf{J}') + \beta^2 h'(\boldsymbol{\theta}', \mathbf{J}'), \tag{3.296a}$$

where

$$H'(\mathbf{J}') \equiv H^0(\mathbf{J}') + \beta h_0(\mathbf{J}') \tag{3.296b}$$

and h' is a function involving second derivatives of H^0 and first derivatives of h .

The analysis we have developed can be used to approximate orbits in a given potential. As we saw in §3.2.2, if we know an integral other than the Hamiltonian of a system with two degrees of freedom, we can calculate the curve in a surface of section on which the consequents of a numerically

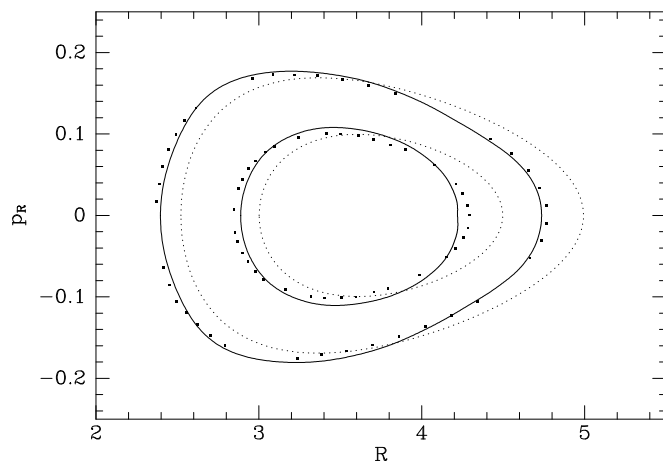


Figure 3.35 A surface of section for orbits in a flattened isochrone potential. The density distribution generating the potential has axis ratio $q = 0.7$. The points are the consequents of numerically calculated orbits. The dotted curves show the orbital tori for the spherical isochrone potential that have the same actions as the numerically integrated orbits. The full curves show the result of using first-order Hamiltonian perturbation theory to deform these tori.

integrated orbit should lie. Since \mathbf{J}' differs from the true action by only $O(\beta^2)$ it should provide an approximate integral of motion, and it is interesting to compare the invariant curve that it yields with the consequents of a numerically integrated orbit. Figure 3.35 is a surface of section for orbits in a flattened isochrone potential. The density distribution that generates this potential is obtained by replacing r by $\sqrt{R^2 + z^2/q^2}$ and M by M/q in equations (2.48b) and (2.49). The axis ratio q has been set equal to 0.7. The dots show the consequents of numerically integrated orbits. The dotted curves show the corresponding invariant curves for the spherical isochrone. The full curves show the results of applying first-order perturbation theory to the spherical isochrone to obtain better approximations to invariant curves.

The full curves in Figure 3.35 fit the numerical consequents much better than the dotted curves, but the fit is not perfect. An obvious strategy for systematically improving our approximation to the true angle-action variables is to use our existing machinery to derive from (3.296a) a second canonical transformation that would enable us to write H as a sum of a Hamiltonian $H''(\mathbf{J}'')$ that is a function of new actions \mathbf{J}'' and a yet smaller perturbation $\beta^4 h''$. After we have performed k transformations, the angle-dependent part of H^β will be of order β^{2k} . In practice this procedure is unlikely to work because after each application the “unperturbed” frequencies of the orbit change from $\boldsymbol{\Omega}' = \partial H'/\partial \mathbf{J}'$ to $\boldsymbol{\Omega}'' = \partial H''/\partial \mathbf{J}''$, and sooner or later we will find that $\mathbf{n} \cdot \boldsymbol{\Omega}''$ is very close to zero for some \mathbf{n} , with the consequence that the corresponding term in the generating function (3.295) becomes large.

This is the problem of **small divisors**. Fortunately, in many applications the coefficients $h_{\mathbf{n}}$ in the numerators of (3.295) decline sufficiently quickly as $|\mathbf{n}|$ increases that for most orbits $|\beta^{2k} h_{\mathbf{n}}/\mathbf{n} \cdot \boldsymbol{\Omega}|$ is small for all \mathbf{n} .

Box 3.5 outlines how the so-called **KAM theory** enables one to overcome the problem of small divisors for most tori, and for them construct a convergent series of canonical transformations that yield the angle-action variables of H^β to arbitrarily high accuracy for sufficiently small β .

3.7.2 Trapping by resonances

Figure 3.36, like Figure 3.35, is a surface of section for motion in a flattened isochrone potential, but the axis ratio of the mass distribution that generates the potential is now $q = 0.4$ rather than $q = 0.7$. The consequents of two orbits are shown together with the approximations to the invariant curves of these orbits that one obtains from the angle-action variables of the spherical isochrone potential with (full curves) and without (dotted curves) first-order perturbation theory. The inner full invariant curve is not very far removed from the inner loop of orbital consequents, but the outer full invariant curve does not even have the same shape as the crescent of consequents that is generated by the second orbit. The deviation between the outer full invariant curve and the consequents is an example of **resonant trapping**, a phenomenon intimately connected with the problem of small divisors that was described above.

To understand this connection, consider how the frequencies of orbits in the flattened isochrone potential are changed by first-order perturbation theory. We obtain the new frequencies by differentiating equation (3.296b) with respect to the actions. Figure 3.37 shows the resulting ratio $\Omega_r/\Omega_\vartheta$ as a function of J_ϑ at the energy of Figure 3.36. Whereas $\Omega_r > \Omega_\vartheta$ for all unperturbed orbits, for some perturbed orbit the resonant condition $\Omega_r - \Omega_\vartheta = 0$ is satisfied. Consequently, if we attempt to use equation (3.295) to refine the tori that generate the full curves in Figure 3.36, small divisors will lead to large distortions in the neighborhood of the resonant torus. These distortions will be unphysical, but they are symptomatic of a real physical effect, namely a complete change in the way in which orbital tori are embedded in phase space. The numerical consequents in Figure 3.36, which mark cross-sections through two tori, one before and one after the change in the embedding, make the change apparent: one torus encloses the shell orbit whose single consequent lies along $p_R = 0$, while the other torus encloses the resonant orbit whose single consequent lies near $(R, p_R) = (2.2, 0.38)$.

Small divisors are important physically because they indicate that a perturbation is acting with one sign for a long time. If the effects of a perturbation can accumulate for long enough, they can become important, even if the perturbation is weak. So if $\mathbf{N} \cdot \boldsymbol{\Omega}$ is small for some \mathbf{N} , then the term $h_{\mathbf{N}}$ in the Hamiltonian can have big effects even if it is very small.

Box 3.5: KAM theory

Over the period 1954–1967 Kolmogorov, Arnold and Moser demonstrated that, notwithstanding the problem of small divisors, convergent perturbation series *can* be constructed for Hamiltonians of the form (3.289). The key ideas are (i) to focus on a single invariant torus rather than a complete foliation of phase space by invariant tori, and (ii) to determine at the outset the frequencies $\boldsymbol{\Omega}$ of the torus to be constructed (Lichtenberg & Lieberman 1992). In particular, we specify that the frequency ratios are far from resonances in the sense that $|\mathbf{n} \cdot \boldsymbol{\Omega}| > \alpha |\mathbf{n}|^{-\gamma}$ for all \mathbf{n} and some fixed, non-negative numbers α and γ . We map an invariant torus of H^0 with frequencies $\boldsymbol{\Omega}$ into an invariant torus of H^β by means of the generating function

$$S(\boldsymbol{\theta}^\beta, \mathbf{J}^0) = \boldsymbol{\theta}^\beta \cdot (\mathbf{J}^0 + \mathbf{j}) + s^\beta(\boldsymbol{\theta}^\beta, \mathbf{J}^0). \quad (1)$$

This differs from (3.290) by the addition of a term $\boldsymbol{\theta}^\beta \cdot \mathbf{j}$, where \mathbf{j} is a constant of order β . Proceeding in strict analogy with the derivation of equations (3.295) and (3.296b), we find that if the Fourier coefficients of s^β are chosen to be

$$s_{\mathbf{n}}^\beta = -\frac{\beta h_{\mathbf{n}}}{\mathbf{n} \cdot \boldsymbol{\Omega}} \quad (\mathbf{n} \neq 0), \quad (2)$$

then we obtain a canonical transformation to new coordinates $(\boldsymbol{\theta}', \mathbf{J}')$ in terms of which H^β takes the form (3.296a) with

$$H'(\mathbf{J}') = H^0(\mathbf{J}') + \beta h_0(\mathbf{J}') - \mathbf{j} \cdot \boldsymbol{\Omega}. \quad (3)$$

We now choose the parameter \mathbf{j} such that the frequencies of H' are still the old frequencies $\boldsymbol{\Omega}$, which were far from any resonance. That is, we choose \mathbf{j} to be the solution of

$$\beta \frac{\partial h_0}{\partial J_j} = \mathbf{j} \cdot \frac{\partial \boldsymbol{\Omega}}{\partial J_j} = \sum_i j_i \cdot \frac{\partial^2 H^0}{\partial J_i \partial J_j}. \quad (4)$$

This linear algebraic equation will be soluble provided the matrix of second derivatives of H^0 is non-degenerate. With \mathbf{j} the solution to this equation, the problem posed by H^β in the $(\boldsymbol{\theta}', \mathbf{J}')$ coordinates differs from our original problem only in that the perturbation is now $O(\beta^2)$. Consequently, a further canonical transformation will reduce the perturbation to $O(\beta^4)$ and so on indefinitely. From the condition $|\mathbf{n} \cdot \boldsymbol{\Omega}| > \alpha |\mathbf{n}|^{-\gamma}$ one may show that the series of transformations converges.

We now use this idea to obtain an analytic model of orbits near resonances. Our working will be a generalization of the discussion of orbital trapping at Lindblad resonances in §3.3.3b. For definiteness we shall assume

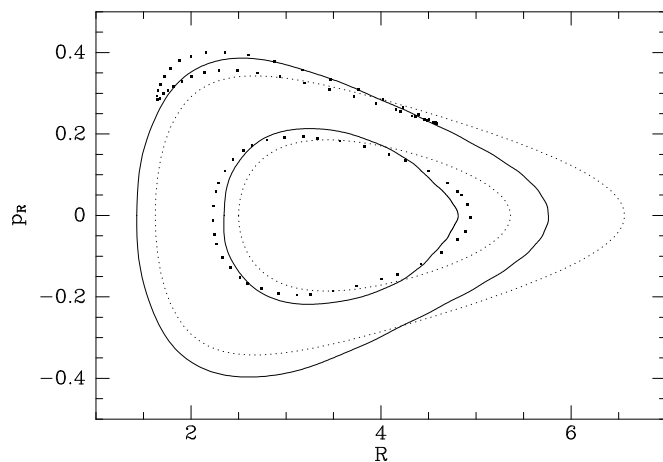


Figure 3.36 The same as Figure 3.35 except that the density distribution generating the potential now has axis ratio $q = 0.4$.

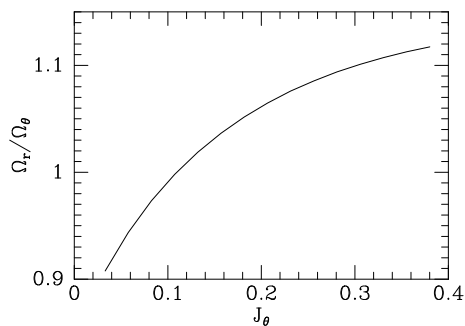


Figure 3.37 The ratio of the frequencies in first-order perturbation theory for a star that moves in a flattened isochrone potential.

that there are three actions and three angles. The resonance of H^0 is characterized by the equation $\mathbf{N} \cdot \boldsymbol{\Omega} = 0$, and $(\boldsymbol{\theta}, \mathbf{J})$ are angle-action variables for the unperturbed Hamiltonian. Then in the neighborhood of the resonant orbit the linear combination of angle variables $\phi_s \equiv \mathbf{N} \cdot \boldsymbol{\theta}$ will evolve slowly, and we start by transforming to a set of angle-action variables that includes the **slow angle** ϕ_s . To do so, we introduce new action variables I_s , I_{f1} , and I_{f2} through the generating function

$$S = (\mathbf{N} \cdot \boldsymbol{\theta})I_s + \theta_1 I_{f1} + \theta_2 I_{f2}. \quad (3.297)$$

Then (eq. D.93)

$$\begin{aligned} \phi_s &= \frac{\partial S}{\partial I_s} = \mathbf{N} \cdot \boldsymbol{\theta} & J_1 &= \frac{\partial S}{\partial \theta_1} = N_1 I_s + I_{f1} \\ \phi_{f1} &= \theta_1 & J_2 &= N_2 I_s + I_{f2} \\ \phi_{f2} &= \theta_2 & J_3 &= N_3 I_s. \end{aligned} \quad (3.298)$$

Since the old actions are functions only of the new ones, H^0 does not acquire any angle dependence when we make the canonical transformation, and the Hamiltonian is of the form

$$H(\phi, \mathbf{I}) = H^0(\mathbf{I}) + \beta \sum_{\mathbf{n}} h_{\mathbf{n}}(\mathbf{I}) e^{i\mathbf{n} \cdot \phi}, \quad (3.299)$$

where it is to be understood that H^0 is a different function of \mathbf{I} than it was of \mathbf{J} and similarly for the dependence on \mathbf{I} of $h_{\mathbf{n}}$. We now argue that any term in the sum that contains either of the **fast angles** ϕ_{f1} and ϕ_{f2} has a negligible effect on the dynamics—these terms give rise to forces that rapidly average to zero. We therefore drop all terms except those with indices that are multiples of $\mathbf{n} = \pm(1, 0, 0)$, including $\mathbf{n} = \mathbf{0}$. Then our approximate Hamiltonian reduces to

$$H(\phi, \mathbf{I}) = H^0(\mathbf{I}) + \beta \sum_k h_k(\mathbf{I}) e^{ik\phi_s}. \quad (3.300)$$

Hamilton's equations now read

$$\begin{aligned} \dot{I}_s &= -i\beta \sum_k k h_k(\mathbf{I}) e^{ik\phi_s} & ; & & \dot{\phi}_s &= \Omega_s + \beta \sum_k \frac{\partial h_k}{\partial I_s} e^{ik\phi_s} \\ \dot{I}_{f1} &= 0 & ; & & \dot{I}_{f2} &= 0, \end{aligned} \quad (3.301)$$

where $\Omega_s \equiv \partial H^0 / \partial I_s$. So I_{f1} and I_{f2} are two constants of motion and we have reduced our problem to one of motion in the (ϕ_s, I_s) plane. Eliminating \mathbf{I} between equations (3.298) and (3.301) we find that although all the old actions vary, two linear combinations of them are constant:

$$N_2 J_1 - N_1 J_2 = \text{constant} \quad ; \quad N_3 J_2 - N_2 J_3 = \text{constant}. \quad (3.302)$$

We next take the time derivative of the equation of motion (3.301) for ϕ_s . We note that Ω_s , but not its derivative with respect to I_s , is small because it vanishes on the resonant torus. Dropping all terms smaller than $O(\beta)$,

$$\ddot{\phi}_s \simeq \frac{\partial \Omega_s}{\partial I_s} \dot{I}_s = -i\beta \frac{\partial \Omega_s}{\partial I_s} \sum_k k h_k e^{ik\phi_s}. \quad (3.303)$$

If we define

$$V(\phi_s) \equiv \beta \frac{\partial \Omega_s}{\partial I_s} \sum_k h_k(\mathbf{I}) e^{ik\phi_s}, \quad (3.304)$$

where \mathbf{I} is evaluated on the resonant torus, then we can rewrite (3.303) as

$$\ddot{\phi}_s = -\frac{dV}{d\phi_s}. \quad (3.305)$$

This is the equation of motion of an oscillator. If V were proportional to ϕ_s^2 , the oscillator would be harmonic. In general it is an anharmonic oscillator, such as a pendulum, for which $V \propto \cos \phi_s$. The oscillator's energy invariant is

$$E_p \equiv \frac{1}{2} \dot{\phi}_s^2 + V(\phi_s). \quad (3.306)$$

V is a periodic function of ϕ_s , so it will have some maximum value V_{\max} , and if $E_p > V_{\max}$, ϕ_s circulates because equation (3.306) does not permit $\dot{\phi}_s$ to vanish. In this case the orbit is not resonantly trapped and the torus is like the ones shown in Figure 3.36 from first-order perturbation theory. If $E_p < V_{\max}$, the angle variable is confined to the range in which $V \leq E_p$; the orbit has been trapped by the resonance. On trapped orbits ϕ_s librates with an amplitude that can be of order unity, and at a frequency of order $\sqrt{\beta}$, while I_s oscillates with an amplitude that cannot be bigger than order $\sqrt{\beta}$. Such orbits generate the kind of torus that is delineated by the crescent of numerical consequents in Figure 3.36. We obtain an explicit expression for the resonantly induced change ΔI_s by integrating the equation of motion (3.301) for I_s :

$$\begin{aligned} \Delta I_s &= - \left(\frac{\partial \Omega_s}{\partial I_s} \right)^{-1} \int d\phi_s \frac{\partial V / \partial \phi_s}{\dot{\phi}_s} \\ &= \pm \left(\frac{\partial \Omega_s}{\partial I_s} \right)^{-1} \sqrt{2[E_p - V(\phi_s)]}, \end{aligned} \quad (3.307)$$

where (3.306) has been used to eliminate $\dot{\phi}_s$.

The full curve in Figure 3.38 shows the result of applying this model of a resonantly trapped orbit to the data depicted in Figure 3.36. Since the model successfully reproduces the gross form of the invariant curve on which the consequents of the trapped orbit lie, we infer that the model has captured the essential physics of resonant trapping. The discrepancies between the full curve and the numerical consequents are attributable to the approximations inherent in the model.

Levitation We now describe one example of an astronomical phenomena that may be caused by resonant trapping of stellar orbits. Other examples are discussed by Tremaine & Yu (2000). In our discussion we shall employ J_r , J_ϑ and J_ϕ to denote the actions of a mildly non-spherical potential that are the natural extensions of the corresponding actions for spherical systems that were introduced in §3.5.2.

The disk of the Milky Way seems to be a composite of two chemically distinct disks, namely the thin disk, to which the Sun belongs, and a thicker, more metal-poor disk (page 13). Sridhar & Touma (1996) have suggested that resonant trapping of the orbits of disk stars may have converted the Galaxy's original thin disk into the thick disk. The theory of hierarchical galaxy formation described in Chapter 9 predicts that the Galaxy was originally dominated by collisionless dark matter, which is not highly concentrated towards the plane. Consequently, the frequency Ω_ϑ at which a

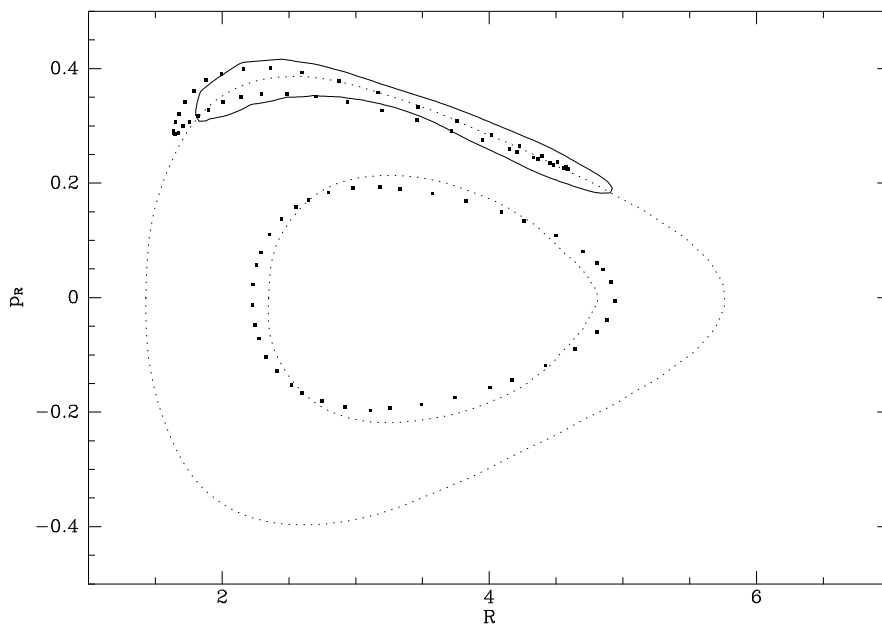


Figure 3.38 Perturbation theory applied to resonant trapping in the flattened isochrone potential. The points are the consequents shown in Figure 3.36, while the full curves in that figure are shown dotted here. The full curve shows the result of using (3.307) to model the resonantly trapped orbit.

star oscillated perpendicular to the plane was originally smaller than the frequency Ω_r of radial oscillations—see equation (3.82). As more and more baryonic material accumulated near the Galaxy’s equatorial plane, the ratio $\Omega_\vartheta/\Omega_r$ rose slowly from a value less than unity to its present value. For stars such as the Sun that are on nearly circular orbits within the plane, Ω_r and Ω_ϑ are equal to the current epicycle and vertical frequencies κ and ν , respectively, so now $\Omega_\vartheta/\Omega_r \simeq 2$ (page 167). It follows that the resonant condition $\Omega_r = \Omega_\vartheta$ has at some stage been satisfied for many stars that formed when the inner Galaxy was dark-matter dominated.

Let us ask what happens to a star in the disk as the disk slowly grows and $\Omega_\vartheta/\Omega_r$ slowly increases. At any energy, the first stars to satisfy the resonant condition $\Omega_r = \Omega_\vartheta$ will have been those with the largest values of Ω_ϑ , that is, stars that orbit close to the plane, and have $J_\vartheta \simeq 0$. In an (R, p_R) surface of section, such orbits lie near the zero-velocity curve that bounds the figure (§3.2.2) because J_ϑ increases as one moves in towards the central fixed point on $p_R = 0$. Hence, the resonant condition will first have been satisfied on the zero-velocity curve, and it is here that the resonant island seen in Figure 3.38 first emerged as the potential flattened. As mass accumulated in the disk, the island moved inwards, and, depending on the values of E and L_z , finally disappeared near the central fixed point.

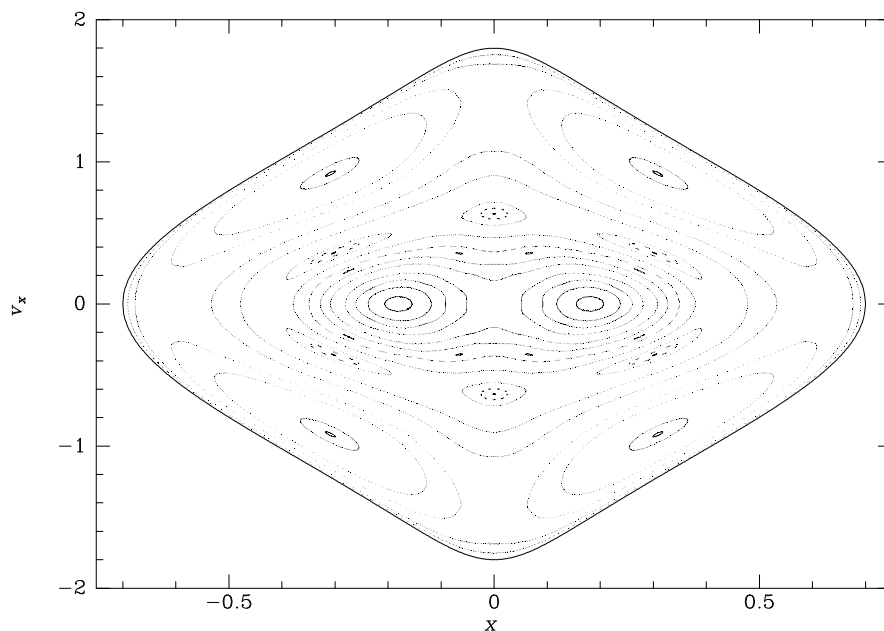


Figure 3.39 A surface of section for motion in Φ_L (eq. 3.103) with $q = 0.6$.

When the advancing edge of a resonant island reaches the star's phase-space location, there are two possibilities: either (a) the star is trapped by the resonance and its phase-space point subsequently moves within the island, or (b) its phase-space point suddenly jumps to the other side of the island. Which of (a) or (b) occurs in a particular case depends on the precise phase of the star's orbit at which the edge reaches it. In practice it is most useful to discard phase information and to consider that either (a) or (b) occurs with appropriate probabilities P_a and $P_b = 1 - P_a$. The value of P_a depends on the speed with which the island is growing relative to the speed with which its center is moving (Problem 3.43); it is zero if the island is shrinking.

We have seen that the resonant island associated with $\Omega_r = \Omega_\theta$ first emerged on the zero-velocity curve, which in a thin disk is highly populated by stars. Most of these stars were trapped as the island grew. They then moved with the island as the latter moved in towards the central fixed point. The stars were finally released as the island shrank somewhere near that point. The net effect of the island's transitory existence is to convert radial action to latitudinal action, thereby shifting stars from eccentric, planar orbits to rather circular but highly inclined ones. Hence, a hot thin disk could have been transformed into a thick disk.

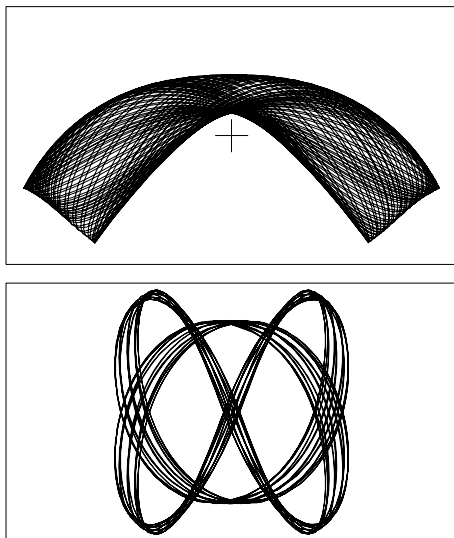


Figure 3.40 The appearance in real space of a banana orbit (top) and a fish orbit (bottom). In the upper panel the cross marks the center of the potential. Resonant box orbits of these types are responsible for the chains of islands in Figure 3.39. The banana orbits generate the outer chain of four islands, and the fish orbits the chain of six islands further in.

3.7.3 From order to chaos

Figure 3.39 is a surface of section for motion in the planar barred potential Φ_L that is defined by equation (3.103) with $q = 0.6$ and $R_c = 0.14$. It should be compared with Figures 3.9 and 3.12, which are surfaces of section for motion in Φ_L for more nearly spherical cases, with $q = 0.9$ and 0.8 . In Figure 3.39 one sees not only the invariant curves of loop and box orbits that fill the other two figures, but also a number of “islands”: a set of four large islands occupies much of the outer region, while a set of six islands of varying sizes is seen further in. In the light of our discussion of resonant trapping, it is natural to refer to the orbits that generate these islands as resonantly trapped box orbits. Figure 3.40 shows what these orbits look like in real space. We see that the outer islands are generated by “banana” orbits in which the x - and y -oscillations are trapped in a $\Omega_x:\Omega_y = 1:2$ resonance (the star oscillates through one cycle left to right while oscillating through two cycles up and down). Similarly, the inner chain of six islands is associated with a “fish” orbit that satisfies the resonance condition $\Omega_x:\Omega_y = 2:3$.

The islands in Figure 3.39 can be thought of as orbits in some underlying integrable Hamiltonian H^0 that are trapped by a resonance arising from a perturbation. This concept lacks precision because we do not know what H^0 actually is. In particular, Hamiltonians of the form $H_q(\mathbf{x}, \mathbf{v}) = \frac{1}{2}v^2 + \Phi_L(\mathbf{x})$ are probably not integrable for any value of the axis ratio q other than unity. Therefore, we cannot simply assume that $H^0 = H_{0.8}$, say. On the other hand, Figure 3.12, which shows the surface of section for $q = 0.8$, contains no resonant islands—all orbits are either boxes or loops—which we know from our study of Stäckel potentials in §3.5.4 is compatible with an integrable potential. So we can *define* an integrable Hamiltonian H^0 that differs very little from $H_{0.8}$ as follows. On each of the invariant tori that

appears in Figure 3.12 we set $H^0 = H_{0.8}$, and at a general phase-space point we obtain the value of H^0 by a suitable interpolation scheme from nearby points at which $H^0 = H_{0.8}$.

The procedure we have just described for defining H^0 (and thus the perturbation $h = H - H^0$) suffers from the defect that it is arbitrary: why start from the invariant tori of $H_{0.8}$ rather than $H_{0.81}$ or some other Hamiltonian? A numerical procedure that might be considered less arbitrary has been described by Kaasalainen & Binney (1994). In any event, it is worth bearing in mind in the discussion that follows that H^0 and h are not uniquely defined, and one really ought to demonstrate that for a given H the islands that are predicted by perturbation theory are reasonably independent of H^0 . As far as we know, no such demonstration is available.

If we accept that the island chains in Figure 3.39 arise from box orbits that are resonantly trapped by some perturbation h on a Stäckel-like Hamiltonian H^0 , two questions arise. First, “are box orbits trapped around resonances other than the 1:2 and 2:3 resonances that generate the banana and fish orbits of Figure 3.40?” Certainly infinitely many resonances are available to trap orbits because as one moves along the sequence of box orbits from thin ones to fat ones, the period of the y -oscillations is steadily growing in parallel with their amplitude, while the period of the x -oscillations is diminishing for the same reason.²⁷ In fact, the transition to loop orbits can be associated with resonant trapping by the 1:1 resonance, so between the banana orbits and the loop orbits there is not only the 2:3 resonance that generates the fishes, but also the 4:5, 5:6, . . . , resonances. In the potential Φ_L on which our example is based, the width of the region in phase space in which orbits are trapped by the $m:n$ resonance diminishes rapidly with $|m + n|$ and the higher-order resonances are hard to trace in the surface of section—but the 4:5 resonance can be seen in Figure 3.39.

The second question is “do resonances occur within resonant islands?” Consider the case of the banana orbits shown in Figure 3.40 as an example. Motion along this orbit is quasiperiodic with two independent frequencies. One independent frequency Ω_b is associated with motion along the bow-shaped closed orbit that runs through the heart of the banana, while the other is the frequency of libration Ω_l about this closed orbit. The libration frequency decreases as one proceeds along the sequence of banana orbits from thin ones to fat ones, so infinitely many resonant conditions $\Omega_b:\Omega_l = m:n$ will be satisfied within an island of banana orbits. In the case of Φ_L there is no evidence that any of these resonances traps orbits, but in another case we might expect trapping to occur also within families of resonantly trapped orbits.

This discussion is rather disquieting because it implies that the degree to which resonant trapping causes the regular structure of phase space inherited from the underlying integral potential H^0 to break up into islands depends

²⁷ The period of a nonlinear oscillator almost always increases with amplitude.

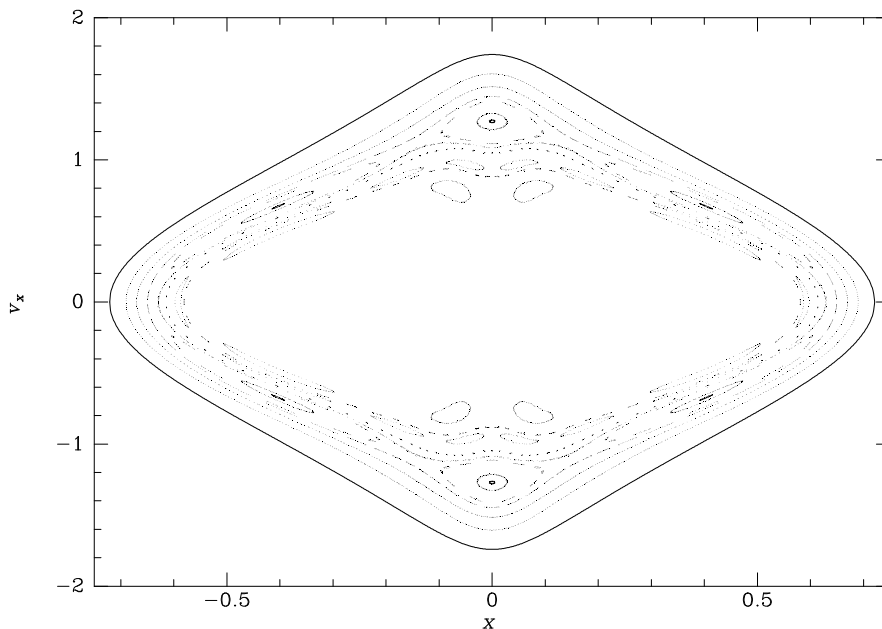


Figure 3.41 Surface of section for motion in the potential Φ_N of equation (3.309) with $R_e = 3$. The inner region has been blanked out and is shown in expanded form in Figure 3.42.

on the detailed structure of the perturbation h . Since we have no unique way of defining h we cannot compute its Fourier coefficients and cannot predict how important islands will be.

We illustrate this point by examining motion in a potential that is closely related to Φ_L in which resonant trapping is *much* more important (Binney 1982). In polar coordinates equation (3.103) for Φ_L reads

$$\Phi_L(R, \phi) = \frac{1}{2}v_0^2 \ln \left[R_c^2 + \frac{1}{2}R^2(q^{-2} + 1) - \frac{1}{2}R^2(q^{-2} - 1) \cos 2\phi \right]. \quad (3.308)$$

The potential

$$\Phi_N(R, \phi) = \frac{1}{2}v_0^2 \ln \left[R_c^2 + \frac{1}{2}R^2(q^{-2} + 1) - \frac{1}{2}R^2(q^{-2} - 1) \cos 2\phi - \frac{R^3}{R_e} \cos 2\phi \right], \quad (3.309)$$

where R_e is a constant, differs from Φ_L only by the addition of $(R^3/R_e) \cos 2\phi$ to the logarithm's argument. For $R \ll R_e$ this term is unimportant, but as R grows it makes the isopotential curves more elongated. Let us set $R_e = 3$, $R_c = 0.14$, and $q = 0.9$, and study the surface of section generated by orbits

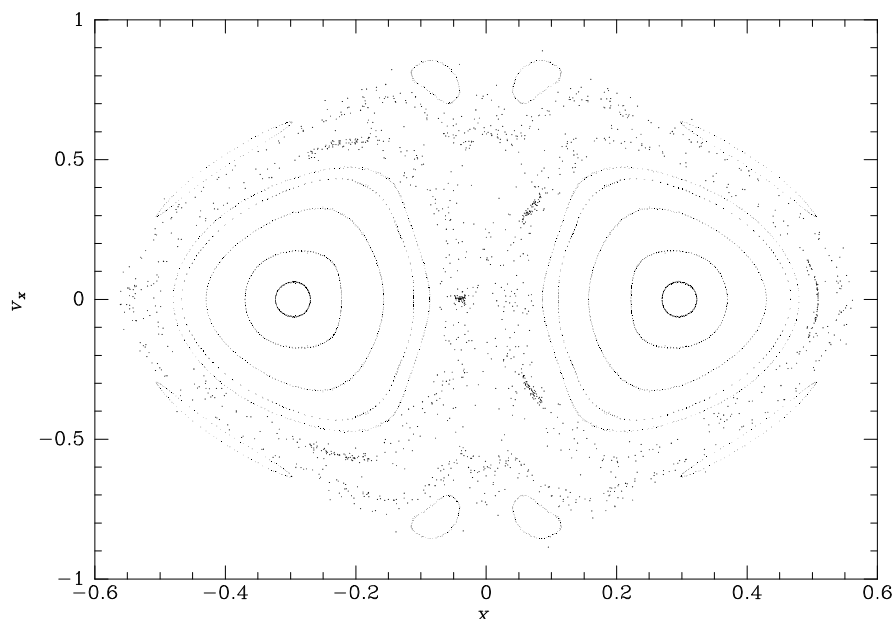


Figure 3.42 The inner part of the surface of section shown in Figure 3.41—the chain of eight islands around the edge is the innermost chain in Figure 3.41. In the gap between this chain and the bull’s-eyes are the consequents of two irregular orbits.

in Φ_N that is most nearly equivalent to the surface of section for Φ_L with the same values of R_c and q that is shown in Figure 3.9. Figure 3.41 shows the outer part of this surface of section. Unlike Figure 3.9 it shows several chains of islands generated by resonantly trapped box orbits. The individual islands are smaller than those in Figure 3.39, and the regions of untrapped orbits between chains of islands are very thin. Figure 3.42 shows the inner part of the same surface of section. In the gap between the region of regular box orbits that is shown in Figure 3.41 and the two bull’s-eyes associated with loop orbits, there is an irregular fuzz of consequents. These consequents belong to just two orbits but they do not lie on smooth curves; they appear to be randomly scattered over a two-dimensional region. Since the gap within which these consequents fall lies just on the boundary of the loop-dominated region, we know that it contains infinitely many resonant box orbits. Hence, it is natural to conclude that the breakdown of orbital regularity, which the random scattering of consequents betrays, is somehow caused by more than one resonance simultaneously trying to trap an individual orbit. One says that the orbits have been made irregular by **resonance overlap**.

(a) Irregular orbits We now consider in more detail orbits whose consequents in a surface of section do not lie on a smooth curve, but appear to be irregularly sprinkled through a two-dimensional region. If we take the

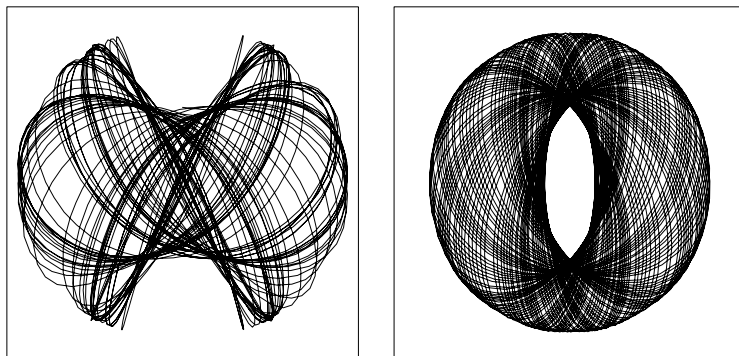


Figure 3.43 Two orbits from the surface of section of Figure 3.42. The left orbit is not quasiperiodic, while the right one is.

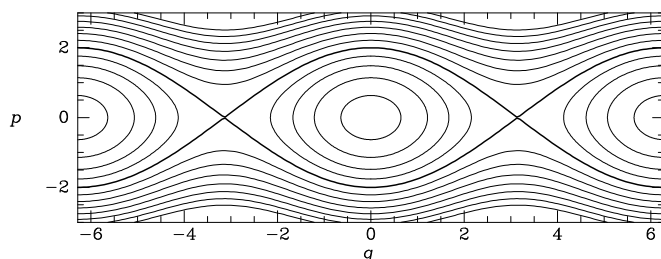


Figure 3.44 Trapped and circulating orbits in a phase plane. The homoclinic orbit, shown by the heavy curve, divides the trapped orbits, which form a chain of islands, from the circulating orbits, whose consequents lie on the wavy lines at top and bottom.

Fourier transform of the time dependence of some coordinate, for example $x(t)$, along such an orbit, we will find that the orbit is not quasiperiodic; the Fourier transform $X(\omega)$ (eq. B.69) has contributions from frequencies that are not integer linear combinations of two or three fundamental frequencies. Figure 3.43 shows the appearance in real space of an orbit that is not quasiperiodic (left) and one that is (right). The lack of quasiperiodicity gives the orbit a scruffy, irregular appearance, so orbits that are not quasiperiodic are called **chaotic** or **irregular orbits**.

There are generally some irregular orbits at the edge of a family of resonantly trapped orbits. Figure 3.44 is a sketch of a surface of section through such a region of phase-space when all orbits are quasiperiodic. The islands formed by the trapped orbits touch at their pointed ends and there are invariant curves of orbits that circulate rather than librate coming right up to these points. The points at which the islands touch are called **hyperbolic fixed points** and the invariant curves that pass through these points are generated by **homoclinic orbits**. In the presence of irregular orbits, the

islands of trapped orbits do not quite touch and the invariant curves of the circulating orbits do not reach right into the hyperbolic fixed point. Consequently there is space between the resonant islands and the region of the circulating orbits. Irregular orbits fill this space.

A typical irregular orbit alternates periods when it is resonantly trapped with periods of circulation. Consequently, if one Fourier transforms $x(t)$ over an appropriate time interval, the orbit may appear quasiperiodic, but the fundamental frequencies that would be obtained from the transform by the method of Box 3.6 would depend on the time interval chosen for Fourier transformation.

If the islands in a chain are individually small, it can be very hard to decide whether an orbit is librating or circulating, or doing both on an irregular pattern.

When it is available, a surface of section is the most effective way of diagnosing the presence of resonantly trapped and irregular orbits. Unfortunately, surfaces of section can be used to study three-dimensional orbits only when an analytic integral other than the Hamiltonian is known, as in the case of orbits in an axisymmetric potential (§3.2). Two other methods are available to detect irregular orbits when a surface of section cannot be used.

(b) Frequency analysis By numerically integrating the equations of motion from some initial conditions, we obtain time series $x(t)$, $y(t)$, etc., for each of the phase-space coordinates. If the orbit is regular, these time series are equivalent to those obtained by substituting $\boldsymbol{\theta} = \boldsymbol{\theta}_0 + \boldsymbol{\Omega}t$ in the Fourier expansions (3.191) of the coordinates. Hence, the frequencies Ω_i may be obtained by Fourier transforming the time series and identifying the various linear combinations $\mathbf{n} \cdot \boldsymbol{\Omega}$ of the fundamental frequencies that occur in the Fourier transform (Box 3.6; Binney & Spergel 1982). If a single system of angle-action variables covers the entire phase space (as in the case of Stäckel potentials), the actions J_i of the orbit that one obtains from a given initial condition \mathbf{w} are continuous functions $J(\mathbf{w})$ of \mathbf{w} , so the frequencies $\Omega_i = \partial H / \partial J_i$ are also continuous functions of \mathbf{w} . Consequently, if we choose initial conditions \mathbf{w}_α at the nodes of some regular two-dimensional grid in phase space, the frequencies will vary smoothly from point to point on the grid. If, by contrast, resonant trapping is important, the actions of orbits will sometimes change discontinuously between adjacent grid points, because one orbit will be trapped, while the next is not. Discontinuities in \mathbf{J} give rise to discontinuities in $\boldsymbol{\Omega}$. Moreover, the resonance that is entrapping orbits will be apparent from the ratios $r_a \equiv \Omega_2 / \Omega_1$ and $r_b \equiv \Omega_3 / \Omega_1$. Hence a valuable way of probing the structure of phase space is to plot a dot at (r_a, r_b) for each orbit obtained by integrating from a regular grid of initial conditions \mathbf{w}_α (Laskar 1990; Dumas & Laskar 1993).

Figure 3.45 shows an example of such a plot of frequency ratios. The orbits plotted were integrated in the potential

$$\Phi(\mathbf{x}) = \frac{1}{2} \ln[x^2 + (y/0.9)^2 + (z/0.7)^2 + 0.1]. \quad (3.310)$$

Box 3.6: Numerical determination of orbital frequencies

The determination of orbital frequencies Ω_i from a numerically integrated orbit is not entirely straightforward because (i) the orbit is integrated for only a finite time interval $(0, T)$, and (ii) the function $x(t)$ is sampled only at discrete times $t_0 = 0, \dots, t_{K-1} = T$, which we shall assume to be equally spaced. Let $\Delta = t_{i+1} - t_i$. Then a “line” $X e^{i\omega t}$ in $x(t)$ contributes to the discrete Fourier transform (Appendix G) an amount

$$\begin{aligned} \hat{x}_p &= X \sum_{k=0}^{K-1} e^{ik\Delta(\omega - \omega_p)} \\ &= X e^{i\alpha u} \frac{\sin \pi u}{\sin(\pi u/K)}, \end{aligned} \quad \text{where} \quad \begin{cases} \omega_p \equiv \frac{2\pi p}{K\Delta}, \\ u \equiv K\Delta(\omega - \omega_p)/(2\pi), \\ \alpha \equiv \pi(K-1)/K. \end{cases} \quad (1)$$

$|\hat{x}_p|$ is large whenever the sine in the denominator vanishes, which occurs when $\omega_p \simeq \omega + 2\pi m/\Delta$, where m is any integer. Thus peaks can arise at frequencies far from ω ; a peak in $|\hat{x}_p|$ that is due to a spectral line far removed from ω is called an **alias** of the line. Near to a peak we can make the approximation $\sin(\pi u/K) \simeq \pi u/K$, so $|\hat{x}_p|$ declines with distance u from the peak only as u^{-1} .

Orbital frequencies can be estimated by fitting equation (1) to the data and thus determining ω . The main difficulty with this procedure is confusion between spectral lines—this confusion can arise either because two lines are nearby, or because a line has a nearby alias. One way to reduce this confusion is to ensure a steeper falloff than u^{-1} by multiplying the original time sequence by a “window” function $w(t)$ that goes smoothly to zero at the beginning and end of the integration period (Press et al. 1986; Laskar 1990). Alternatively, one can identify peaks in the second difference of the spectrum, defined by $\hat{x}_p'' = \hat{x}_{p+1} + \hat{x}_{p-1} - 2\hat{x}_p$. One can show that for $u/K \ll 1$ the contribution to \hat{x}_p'' of a line is

$$\hat{x}_p'' = \frac{2XK}{\pi} \frac{e^{i\alpha u} \sin \pi u}{u(u^2 - 1)}, \quad (2)$$

which falls off as u^{-3} . The frequency, etc., of the line can be estimated from the ratio of the \hat{x}_p'' on either side of the line’s frequency.

Ω_i was defined to be the non-zero frequency with the largest amplitude in the spectrum of the i th coordinate, and 10 000 orbits were obtained by dropping particles from a grid of points on the surface $\Phi(\mathbf{x}) = 0.5$. Above and to the right of the center of the figure, the points are organized into regular ranks that reproduce the grid of initial conditions in slightly distorted form.

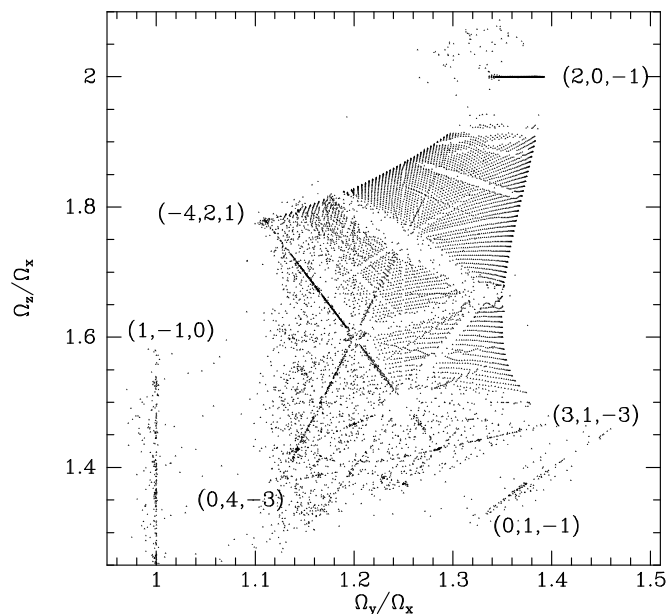


Figure 3.45 The ratios of orbital frequencies for orbits integrated in a three-dimensional non-rotating bar potential.

We infer that resonant trapping is unimportant in the phase-space region sampled by these initial conditions. Running through the ranks we see several depopulated lines, while both within the ranks and beyond other lines are conspicuously heavily populated: orbits that have been resonantly trapped produce points that lie along these lines. The integers n_i in the relevant resonant condition $\mathbf{n} \cdot \boldsymbol{\Omega} = 0$ are indicated for some of the lines.

In some parts of Figure 3.45, for example the lower left region, the grid of initial conditions has become essentially untraceable. The disappearance of the grid indicates that irregular motion is important. In fact, the frequencies Ω_i are not well defined for an irregular orbit, because its time series, $x(t)$, $y(t)$, etc., are not quasiperiodic. When software designed to extract the frequencies of regular orbits is used on a time series that is not quasiperiodic, the frequencies returned vary erratically from one initial condition to the next and the resulting points in the plane of frequency ratios scatter irregularly.

(c) Liapunov exponents If we integrate Hamilton's equations for some time t , we obtain a mapping \mathbf{H}_t of phase space onto itself. Let \mathbf{H}_t map the phase space point \mathbf{w}_0 into the point \mathbf{w}_t . Points near \mathbf{w}_0 will be mapped to points that lie near \mathbf{w}_t , and if we confine our attention to a sufficiently small region around \mathbf{w}_0 , we may approximate \mathbf{H}_t by a linear map of the neighborhood of \mathbf{w}_0 into a neighborhood of \mathbf{w}_t . We now determine this map. Let \mathbf{w}'_0 be a point near \mathbf{w}_0 , and $\delta\mathbf{w}(t) = \mathbf{H}_t\mathbf{w}'_0 - \mathbf{H}_t\mathbf{w}_0$ be the difference

between the phase-space coordinates of the points reached by integrating Hamilton's equations for time t from the initial conditions \mathbf{w}'_0 and \mathbf{w}_0 . Then the equations of motion of the components of $\delta\mathbf{w}$ are

$$\begin{aligned}\dot{\delta\mathbf{x}} &= \left(\frac{\partial H}{\partial \mathbf{v}}\right)_{\mathbf{w}'_t} - \left(\frac{\partial H}{\partial \mathbf{v}}\right)_{\mathbf{w}_t} \simeq \left(\frac{\partial^2 H}{\partial \mathbf{w} \partial \mathbf{v}}\right)_{\mathbf{w}_t} \cdot \delta\mathbf{w} \\ \dot{\delta\mathbf{v}} &= -\left(\frac{\partial H}{\partial \mathbf{x}}\right)_{\mathbf{w}'_t} + \left(\frac{\partial H}{\partial \mathbf{x}}\right)_{\mathbf{w}_t} \simeq -\left(\frac{\partial^2 H}{\partial \mathbf{w} \partial \mathbf{x}}\right)_{\mathbf{w}_t} \cdot \delta\mathbf{w},\end{aligned}\quad (3.311)$$

where the approximate equality in each line involves approximating the first derivatives of H by the leading terms in their Taylor series expansions. Equations (3.311) are of the form

$$\frac{d\delta\mathbf{w}}{dt} = \mathbf{M}_t \cdot \delta\mathbf{w} \quad \text{where} \quad \mathbf{M}_t \equiv \begin{pmatrix} \frac{\partial^2 H}{\partial \mathbf{x} \partial \mathbf{v}} & \frac{\partial^2 H}{\partial \mathbf{v} \partial \mathbf{v}} \\ -\frac{\partial^2 H}{\partial \mathbf{x} \partial \mathbf{x}} & -\frac{\partial^2 H}{\partial \mathbf{v} \partial \mathbf{x}} \end{pmatrix}. \quad (3.312)$$

For any initial vector $\delta\mathbf{w}_0$ these equations are solved by $\delta\mathbf{w}_t = \mathbf{U}_t \cdot \delta\mathbf{w}_0$, where \mathbf{U}_t is the matrix that solves

$$\frac{d\mathbf{U}_t}{dt} = \mathbf{M}_t \cdot \mathbf{U}_t. \quad (3.313)$$

We integrate this set of ordinary coupled linear differential equations from $\mathbf{U}_0 = \mathbf{I}$ in parallel with Hamilton's equations of motion for the orbit. Then we are in possession of the matrix \mathbf{U}_t that describes the desired linear map of a neighborhood of \mathbf{w}_0 into a neighborhood of \mathbf{w}_t . We perform a "singular-value decomposition" of \mathbf{U}_t (Press et al. 1986), that is we write it as a product $\mathbf{U}_t = \mathbf{R}_2 \cdot \mathbf{S} \cdot \mathbf{R}_1$ of two orthogonal matrices \mathbf{R}_i and a diagonal matrix \mathbf{S} .²⁸ \mathbf{U}_t conserves phase-space volume (page 803), so it never maps any vector to zero and the diagonal elements of \mathbf{S} are all non-zero. In fact they are all positive because \mathbf{U}_t evolves continuously from the identity, and their product is unity. A useful measure of the amount by which \mathbf{U}_t shears phase space is the magnitude s of the largest element of \mathbf{S} . The **Liapunov exponent** of the orbit along which (3.313) has been integrated is defined to be

$$\lambda = \lim_{t \rightarrow \infty} \frac{\ln s}{t}. \quad (3.314)$$

²⁸ Any linear transformation of an N -dimensional vector space can be decomposed into a rotation, a rescaling in N perpendicular directions, and another rotation. \mathbf{R}_1 rotates axes to the frame in which the coordinate directions coincide with the scaling directions. \mathbf{S} effects the rescaling. \mathbf{R}_2 first rotates the coordinate directions back to their old values and then effects whatever overall rotation is required.

Since the scaling s is dimensionless, the Liapunov exponent λ has dimensions of a frequency. In practice one avoids integrating (3.313) for long times because numerical difficulties would be encountered once the ratio of the largest and smallest numbers on the diagonal of \mathbf{S} became large. Instead one integrates along the orbit for some time t_1 to obtain a value s_1 , and then sets \mathbf{U}_t back to the identity and continues integrating for a further time t_2 to obtain s_2 , after which \mathbf{U}_t is again set to the identity before the integration is continued. After N such steps one estimates λ from

$$\lambda \simeq \frac{\sum_i^N \ln s_i}{\sum_i^N t_i}. \quad (3.315)$$

Using this procedure one finds that along a regular orbit $\lambda \rightarrow 0$, while along an irregular orbit λ is non-zero.

Angle-action variables enable us to understand why λ is zero for a regular orbit. A point near \mathbf{w}_0 will have angles and actions that differ from those of \mathbf{w}_0 by small amounts $\delta\theta_i, \delta J_i$. The action differences are invariant as we move along the orbit, while the angle differences increase linearly in time due to differences in the frequencies Ω_i of the orbits on which our initial point and \mathbf{w}_0 lie. Consequently, the scalings s_i associated with angle differences increase linearly in time, and, by (3.314), the Liapunov exponent is $\lambda = \lim_{t \rightarrow \infty} t^{-1} \ln t = 0$.

If the Liapunov exponent of an orbit is non-zero, the largest scaling factor s must increase exponentially in time. Thus in this case initially neighboring orbits diverge exponentially in time. It should be noted, however, that this exponential divergence holds only so long as the orbits remain close in phase space: the definition of the Liapunov exponent is in terms of the linearized equations for orbital perturbations. The approximations involved in deriving these equations will soon be violated if the solutions to the equations are exponentially growing. Hence, we cannot conclude from the fact that an orbit's Liapunov exponent is non-zero that an initially neighboring orbit will necessarily stray far from the original orbit.

3.8 Orbits in elliptical galaxies

Elliptical galaxies nearly always have cusps in their central density profiles in which $\rho \sim r^{-\alpha}$ with $0.3 \lesssim \alpha \lesssim 2$ (BM §4.3.1). Black holes with masses $\sim 0.2\%$ of the mass of the visible galaxy are believed to reside at the centers of these cusps (§1.1.6 and BM §11.2.2). Further out the mass distributions of many elliptical galaxies are thought to be triaxial (BM §4.3.3). These features make the orbital dynamics of elliptical dynamics especially rich, and illustrate aspects of galaxy dynamics that we have already discussed in this chapter (Merritt & Fridman 1996; Merritt & Valluri 1999).

3.8.1 The perfect ellipsoid

A useful basic model of the orbital dynamics of a triaxial elliptical galaxy is provided by extensions to three dimensions of the two-dimensional Stäckel potentials of §3.5.4 (de Zeeuw 1985). The simplest three-dimensional system that generates a Stäckel potential through Poisson's equation is the **perfect ellipsoid**, in which the density is given by

$$\rho(\mathbf{x}) = \frac{\rho_0}{(1+m^2)^2} \quad \text{where} \quad m^2 \equiv \frac{x^2 + (y/q_1)^2 + (z/q_2)^2}{a_0^2}. \quad (3.316)$$

In this formula q_1 and q_2 are the axis ratios of the ellipsoidal surfaces of constant density, and a_0 is a scale length. At radii significantly smaller than a_0 , the density is approximately constant, while at $r \gg a_0$ the density falls off $\propto r^{-4}$. Since these asymptotic forms differ from those characteristic of elliptical galaxies, we have to expect the orbital structures of real galaxies to differ in detail from that of the perfect ellipsoid, but nevertheless the model exhibits much of the orbital structure seen in real elliptical galaxies.

By an analysis similar to that used in §3.5.4 to explore the potential of a planar bar, one can show that the perfect ellipsoid supports four types of orbit. Figure 3.46 depicts an orbit of each type. At top left we have a box orbit. The key feature of a box orbit is that it touches the isopotential surface for its energy at its eight corners. Consequently, the star comes to rest for an instant at these points; a box orbit is conveniently generated numerically by releasing a star from rest on the equipotential surface. The potential's longest axis emerges from the orbit's convex face. The other three orbits are all **tube orbits**: stars on these orbits circulate in a fixed sense around the hole through the orbit's center, and are never at rest. The most important tube orbits are the short-axis loops shown at top right, which circulate around the potential's shortest axis. These orbits are mildly distorted versions of the orbits that dominate the phase space of a flattened axisymmetric potential. The tube orbits at the bottom of Figure 3.46 are called outer (left) and inner long-axis tube orbits, and circulate around the longest axis of the potential. Tube orbits around the intermediate axis are unstable. All these orbits can be quantified by a single system of angle-action coordinates $(J_\lambda, J_\mu, J_\nu)$ that are generalizations of the angle-action coordinates for spherical potentials $(J_r, J_\vartheta, J_\phi)$ of Table 3.1 (de Zeeuw 1985).

3.8.2 Dynamical effects of cusps

The most important differences between a real galactic potential and the best-fitting Stäckel potential are at small radii. Box orbits, which alone penetrate to arbitrarily small radii, are the most affected by these differences. The box orbits of a given energy form a two-parameter family: the parameters can be taken to be an orbit's axis ratios. Resonant relations $\mathbf{n} \cdot \boldsymbol{\Omega} = 0$ between the fundamental frequencies of an orbit are satisfied at various points

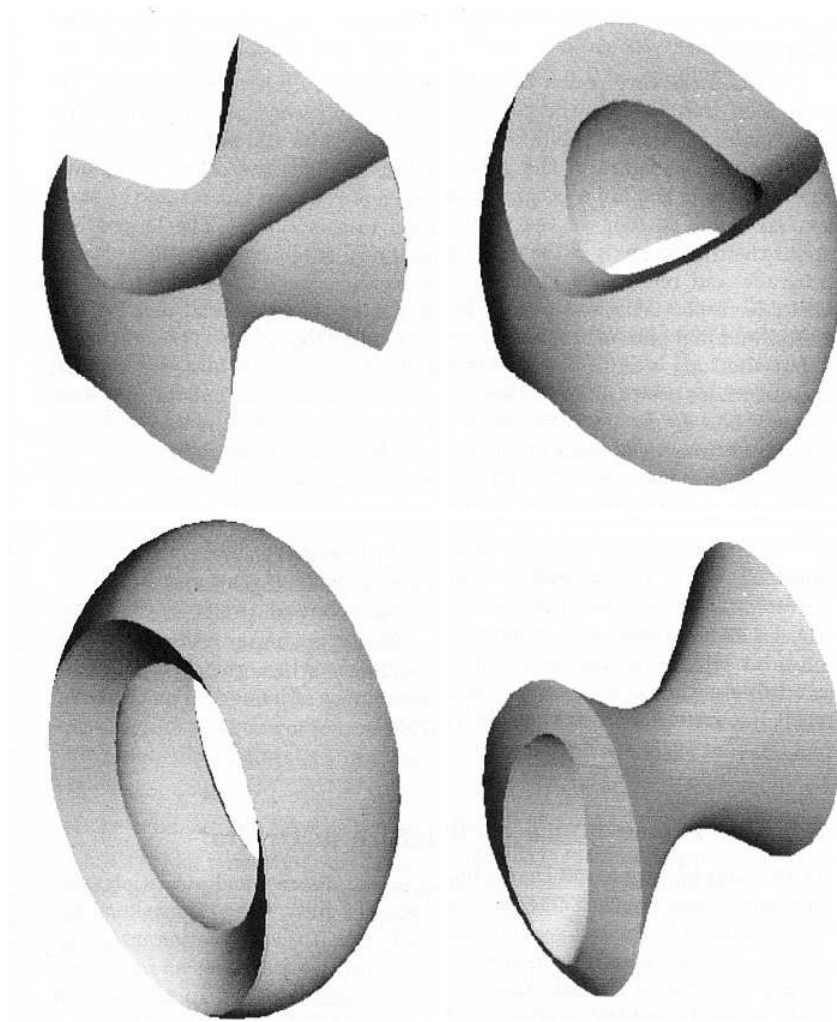


Figure 3.46 Orbits in a non-rotating triaxial potential. Clockwise from top left: (a) box orbit; (b) short-axis tube orbit; (c) inner long-axis tube orbit; (d) outer long-axis tube orbit. From Statler (1987), by permission of the AAS.

in parameter space, but in a Stäckel potential none of these resonances traps other orbits. We expect perturbations to cause some resonances to become trapping. Hence it is no surprise to find that in potentials generated by slightly cusped mass distributions, significant numbers of orbits are trapped by resonances. (In Figure 3.45 we have already encountered extensive resonant trapping of box orbits in a triaxial potential that differs from a Stäckel potential.)

A regular orbit on which the three angle variables satisfy the condition $\mathbf{n} \cdot \boldsymbol{\Omega} = 0$ is a two-dimensional object since its three actions are fixed, and one of its angles is determined by the other two. Consequently, the orbit occupies a surface in real space. A generic resonantly trapped orbit is a three-dimensional structure because it has a finite libration amplitude around the resonant orbit. In practice the amplitude of the libration is usually small, with the result that the orbit forms a sheet of small but finite thickness around the resonant orbit. It is found that stable resonant box orbits are **centrophobic**, that is, they avoid the galactic center (Merritt & Valluri 1999).

Steepening the cusp in the galaxy's central density profile enhances the difference between the galactic potential and the best-fitting Stäckel model and thus the importance of resonances. More and more resonances overlap (§3.7.3) and the fraction of irregular orbits increases.

The existence of large numbers of irregular orbits in elliptical galaxies is likely to have important but imperfectly understood astronomical implications because irregular orbits display a kind of creep or diffusion. To understand this phenomenon, imagine that there is a clean distinction between regular and irregular regions of $2N$ -dimensional phase space. The regular region is occupied by regular orbits and is strictly off-limits to any irregular orbit, while the irregular region is off-limits to regular orbits. However, while each regular orbit is strictly confined to its N -dimensional torus and never trespasses on the territory of a different regular orbit, over time an irregular orbit explores at least some of the irregular region of phase space. In fact, the principal barrier to an irregular orbit's ability to wander is walls formed by regular orbits. In the case $N = 2$ of two-dimensional motion, the energetically accessible part of phase space is three-dimensional, while the walls formed by regular orbits are two-dimensional. Hence such a wall can completely bound some portion of irregular phase space, and forever exclude an irregular orbit from part of irregular phase space. In the case $N = 3$ that is relevant for elliptical galaxies, the energetically accessible region of phase space is five-dimensional while the wall formed by a regular orbit is three-dimensional. Since the boundary of a five-dimensional volume is a four-dimensional region, it is clear that no regular orbit can divide the irregular region of phase space into two. Hence, it is believed that given enough time an irregular orbit with $N \geq 3$ degrees of freedom will eventually visit every part of the irregular region of phase space.

The process by which irregular orbits wander through phase space is called **Arnold diffusion** and is inadequately understood. Physically, it probably involves repeated trapping by a multitude of high-order resonances. In elliptical galaxies and the bars of barred disk galaxies, the rate of Arnold diffusion may be comparable to the Hubble time and could be a major factor in determining the rate of galactic evolution.

If the timescale associated with Arnold diffusion were short enough, galaxy models would need to include only one irregular orbit. The phase-

space density f_{irr} contributed by this orbit would be the same at all points on the energy hypersurface $H(\mathbf{x}, \mathbf{v}) = E$ except in the regular region of phase space, where f_{irr} would vanish.²⁹ It is not yet clear how galaxy modeling is best done when the timescale for Arnold diffusion is comparable to the Hubble time.

3.8.3 Dynamical effects of black holes

Introducing even a small black hole at the center of a triaxial galaxy that has a largely regular phase space destroys much of that regularity. There is a simple physical explanation of this phenomenon (Gerhard & Binney 1985; Merritt & Quinlan 1998).

Consider a star on the box orbit shown at top left in Figure 3.46. Each crossing time the star passes through the orbit's waist on an approximately rectilinear trajectory, and is deflected through some angle θ_{defl} by the black hole's gravitational field. If M is the mass of the hole, and v and b are, respectively, the speed and the distance from the galactic center at which the star would have passed the waist had the hole not deflected it, then from equation (3.52) we have that

$$\theta_{\text{defl}} = 2 \tan^{-1} \left(\frac{GM}{bv^2} \right). \quad (3.317)$$

The speed v will be similar for all passages, but the impact parameter b will span a wide range of values over a series of passages. For any value of M , no matter how small, there is a chance that b will be small enough for the star to be scattered onto a significantly different box orbit.

The tensor virial theorem (§4.8.3) requires that the velocity dispersion be larger parallel to the longest axis of a triaxial system than in the perpendicular directions. Repeated scattering of stars by a nuclear black hole will tend to make the velocity dispersion isotropic, and thus undermine the orbital support for the triaxiality of the potential. If the potential loses its triaxiality, angular momentum will become a conserved quantity, and every star will have a non-zero pericentric distance. Hence stars will no longer be exposed to the risk of coming arbitrarily close to the black hole, and stars will disappear from the black hole's menu.

Let us assume that the distribution of a star's crossing points is uniform within the waist and calculate the expectation value of the smallest value taken by r in N passages. Let the area of the waist be πR^2 . Then the probability of there being n crossing points in a circle of radius r is given by the Poisson distribution (Appendix B.8) as

$$P(n|r) = \frac{(Nr^2/R^2)^n}{n!} e^{-Nr^2/R^2}. \quad (3.318)$$

²⁹ See Häfner et al. (2000) for a method of exploiting the uniformity of f_{irr} in galaxy modeling.

The probability that the closest passage lies in $(r, r + dr)$ is the probability that there are zero passages inside r and a non-zero number of passages in the surrounding annulus, has area $2\pi r dr$. Thus this probability is

$$dP = (1 - e^{-2Nrdr/R^2})e^{-Nr^2/R^2} \simeq \frac{2Nrdr}{R^2}e^{-Nr^2/R^2}. \quad (3.319)$$

The required expectation value of r_1 is now easily calculated:

$$\langle r_1 \rangle = \int dr \frac{2Nr^2}{R^2}e^{-Nr^2/R^2} = \sqrt{\frac{\pi}{N}} \frac{R}{2}. \quad (3.320)$$

From equation (3.317) the deflection that corresponds to $\langle r_1 \rangle$ is

$$\theta_{\text{def,max}} = 2 \tan^{-1} \left(\frac{2\sqrt{N}GM}{\sqrt{\pi}v^2R} \right). \quad (3.321)$$

Two empirical correlations between galactic parameters enable us to estimate $\theta_{\text{def,max}}$ for a star that reaches maximum radius R_{max} in an elliptical galaxy with measured line-of-sight velocity dispersion σ_{\parallel} . First we take the black hole's mass M from the empirical relation (1.27). In the galaxy's lifetime τ we have $N \simeq \sigma_{\parallel}\tau/2R_{\text{max}}$, and we relate R_{max} to D_n , the diameter within which the mean surface brightness of an elliptical galaxy is 20.75 mag arcsec⁻² in the B band: D_n is correlated with σ_{\parallel} such that (BM eq. 4.43)

$$D_n = 5.2 \left(\frac{\sigma_{\parallel}}{200 \text{ km s}^{-1}} \right)^{1.33} \text{ kpc}. \quad (3.322)$$

With these relations, (3.321) becomes

$$\theta_{\text{def,max}} \simeq 2 \tan^{-1} \left[0.08 \frac{D_n^{3/2}}{R_{\text{max}}^{3/2}} \frac{R_{\text{max}}}{R} \frac{\sigma_{\parallel}^2}{v^2} \left(\frac{\sigma_{\parallel}}{200 \text{ km s}^{-1}} \right)^{0.5} \left(\frac{\tau}{10 \text{ Gyr}} \right)^{1/2} \right]. \quad (3.323)$$

For the moderately luminous elliptical galaxies that are of interest here, D_n is comparable to, or slightly larger than, the effective radius (Dressler et al. 1987), and thus similar to the half-mass radius $r_h = 1.3R_e$ for the $R^{1/4}$ profile. Thus for the majority of stars $D_n/R_{\text{max}} \simeq 1$. From Figure 3.46 we estimate $R_{\text{max}}/R \simeq 10$. To estimate the ratio σ_{\parallel}/v we deduce from equations (2.66) and (2.67) that for a Hernquist model with scale radius a the potential drop $\Delta\Phi = \Phi(a) - \Phi(0)$ between $r_h = 2.41a$ and the center is $0.71GM_{\text{gal}}/a$, so $v^2 = 2\Delta\Phi = 1.4GM_{\text{gal}}/a$. From Figure 4.4 we see that $\sigma_{\parallel} \simeq 0.2\sqrt{GM_{\text{gal}}/a}$, so $(\sigma_{\parallel}/v)^2 \simeq 35$. Inserting these values into equation (3.323) we find $\theta_{\text{def,max}} \simeq 2.6^\circ$. Scattering by such a small angle will probably not undermine a galaxy's triaxiality, but stars with smaller apocenter distances R_{max} will be deflected through significant angles, so it is likely that the black hole will erode triaxiality in the galaxy's inner parts (Norman, May, & van Albada 1985; Merritt & Quinlan 1998).

Problems

3.1 [1] Show that the radial velocity along a Kepler orbit is

$$\dot{r} = \frac{GM_e}{L} \sin(\psi - \psi_0), \quad (3.324)$$

where L is the angular momentum. By considering this expression in the limit $r \rightarrow \infty$ show that the eccentricity e of an unbound Kepler orbit is related to its speed at infinity by

$$e^2 = 1 + \left(\frac{Lv_\infty}{GM} \right)^2. \quad (3.325)$$

3.2 [1] Show that for a Kepler orbit the eccentric anomaly η and the true anomaly $\psi - \psi_0$ are related by

$$\cos(\psi - \psi_0) = \frac{\cos \eta - e}{1 - e \cos \eta} \quad ; \quad \sin(\psi - \psi_0) = \sqrt{1 - e^2} \frac{\sin \eta}{1 - e \cos \eta}. \quad (3.326)$$

3.3 [1] Show that the energy of a circular orbit in the isochrone potential (2.47) is $E = -GM/(2a)$, where $a = \sqrt{b^2 + r^2}$. Let the angular momentum of this orbit be $L_c(E)$. Show that

$$L_c = \sqrt{GMb} \left(x^{-1/2} - x^{1/2} \right), \quad \text{where} \quad x \equiv -\frac{2Eb}{GM}. \quad (3.327)$$

3.4 [1] Prove that if a homogeneous sphere of a pressureless fluid with density ρ is released from rest, it will collapse to a point in time $t_{\text{ff}} = \frac{1}{4} \sqrt{3\pi/(2G\rho)}$. The time t_{ff} is called the **free-fall time** of a system of density ρ .

3.5 [3] Generalize the timing argument in Box 3.1 to a universe with non-zero vacuum-energy density. Evaluate the required mass of the Local Group for a universe of age $t_0 = 13.7$ Gyr with (a) $\Omega_{\Lambda 0} = 0$; (b) $\Omega_{\Lambda 0} = 0.76$, $h_7 = 1.05$. Hints: the energy density in radiation can be neglected. The solution requires evaluation of an integral similar to (1.62).

3.6 [1] A star orbiting in a spherical potential suffers an arbitrary instantaneous velocity change while it is at pericenter. Show that the pericenter distance of the ensuing orbit cannot be larger than the initial pericenter distance.

3.7 [2] In a spherically symmetric system, the apocenter and pericenter distances are given by the roots of equation (3.14). Show that if $E < 0$ and the potential $\Phi(r)$ is generated by a non-negative density distribution, this equation has either no root, a repeated root, or two roots (Contopoulos 1954). Thus there is at most one apocenter and pericenter for a given energy and angular momentum. Hint: take the second derivative of $E - \Phi$ with respect to $u = 1/r$ and use Poisson's equation.

3.8 [1] Prove that circular orbits in a given potential are unstable if the angular momentum per unit mass on a circular orbit decreases outward. Hint: evaluate the epicycle frequency.

3.9 [2] Compute the time-averaged moments of the radius, $\langle r^n \rangle$, in a Kepler orbit of semi-major axis a and eccentricity e , for $n = 1, 2$ and $n = -1, -2, -3$.

3.10 [2] $\Delta\psi$ denotes the increment in azimuthal angle during one complete radial cycle of an orbit.

(a) Show that in the potential (3.57)

$$\Delta\psi = \frac{2\pi L}{\sqrt{-2Er_a r_p}}, \quad (3.328)$$

where r_a and r_p are the apo- and pericentric radii of an orbit of energy E and angular momentum L . Hint: by contour integration one can show that for $A > 1$, $\int_{-\pi/2}^{\pi/2} d\theta / (A + \sin\theta) = \pi/\sqrt{A^2 - 1}$.

(b) Prove in the epicycle approximation that along orbits in a potential with circular frequency $\Omega(R)$,

$$\Delta\psi = 2\pi \left(4 + \frac{d \ln \Omega^2}{d \ln R} \right)^{-1/2}. \quad (3.329)$$

(c) Show that the exact expression (3.328) reduces for orbits of small eccentricity to (3.329).

3.11 [1] For what spherically symmetric potential is a possible trajectory $r = ae^{b\psi}$?

3.12 [2] Prove that the mean-square velocity is on a bound orbit in a spherical potential $\Phi(r)$ is

$$\langle v^2 \rangle = \left\langle r \frac{d\Phi}{dr} \right\rangle, \quad (3.330)$$

where $\langle \cdot \rangle$ denotes a time average.

3.13 [2] Let $\mathbf{r}(s)$ be a plane curve depending on the parameter s . Then the **curvature** is

$$K = \frac{|\mathbf{r}' \times \mathbf{r}''|}{|\mathbf{r}'|^3}, \quad (3.331)$$

where $\mathbf{r}' \equiv d\mathbf{r}/ds$. The local radius of curvature is K^{-1} . Prove that the curvature of an orbit with energy E and angular momentum L in the spherical potential $\Phi(r)$ is

$$K = \frac{L d\Phi/dr}{2^{3/2} r [E - \Phi(r)]^{3/2}}. \quad (3.332)$$

Hence prove that no orbit in any spherical mass distribution can have an inflection point (in contrast to the cover illustration of Goldstein, Safko, & Poole 2002).

3.14 [1] Show that in a spherical potential the vertical and circular frequencies ν and Ω (eqs. 3.79) are equal.

3.15 [1] Prove that at any point in an axisymmetric system at which the local density is negligible, the epicycle, vertical, and circular frequencies κ , ν , and Ω (eqs. 3.79) are related by $\kappa^2 + \nu^2 = 2\Omega^2$.

3.16 [1] Using the epicycle approximation, prove that the azimuthal angle $\Delta\psi$ between successive pericenters lies in the range $\pi \leq \Delta\psi \leq 2\pi$ in the gravitational field arising from any spherical mass distribution in which the density decreases outwards.

3.17 [3] The goal of this problem is to prove the results of Problem 3.16 without using the epicycle approximation (Contopoulos 1954).

(a) Using the notation of §3.1, show that

$$E - \Phi - \frac{L^2}{2r^2} = (u_1 - u)(u - u_2) \left\{ \frac{1}{2}L^2 + \Phi[u, u_1, u_2] \right\}, \quad (3.333)$$

where $u_1 = 1/r_1$ and $u_2 = 1/r_2$ are the reciprocals of the pericenter and apocenter distances of the orbit respectively, $u = 1/r$, and

$$\Phi[u, u_1, u_2] = \frac{1}{u_1 - u_2} \left[\frac{\Phi(u_1) - \Phi(u)}{u_1 - u} - \frac{\Phi(u) - \Phi(u_2)}{u - u_2} \right]. \quad (3.334)$$

This expression is a second-order divided difference of the potential Φ regarded as a function of u , and a variant of the mean-value theorem of calculus shows that $\Phi[u, u_1, u_2] = \frac{1}{2}\Phi''(\bar{u})$ where \bar{u} is some value of u in the interval (u_1, u_2) . Then use the hint in Problem 3.7 and equation (3.18b) to deduce that $\Delta\psi \leq 2\pi$ when the potential Φ is generated by a non-negative, spherically symmetric density distribution.

(b) A lower bound on $\Delta\psi$ can be obtained from working in a similar manner with the function

$$\chi(\omega) = \frac{2\omega\Phi}{L}, \quad \text{where } \omega \equiv \frac{L}{r^2}. \quad (3.335)$$

Show that

$$\frac{2\omega E}{L} - \chi(\omega) - \omega^2 = (\omega_1 - \omega)(\omega - \omega_2) \{1 + \chi[\omega, \omega_1, \omega_2]\}, \quad (3.336)$$

where $\omega_1 = L/r_1^2$, $\omega_2 = L/r_2^2$ and $\chi[\omega, \omega_1, \omega_2]$ is a second-order divided difference of $\chi(\omega)$. Now deduce that $\Delta\psi \geq \pi$ for any potential in which the circular frequency $\Omega(r)$ decreases outwards.

3.18 [1] Let $\Phi(R, z)$ be the Galactic potential. At the solar location, $(R, z) = (R_0, 0)$, prove that

$$\frac{\partial^2 \Phi}{\partial z^2} = 4\pi G \rho_0 + 2(A^2 - B^2), \quad (3.337)$$

where ρ_0 is the density in the solar neighborhood and A and B are the Oort constants. Hint: use equation (2.73).

3.19 [3] Consider an attractive power-law potential, $\Phi(r) = Cr^\alpha$, where $-1 \leq \alpha \leq 2$ and $C > 0$ for $\alpha > 0$, $C < 0$ for $\alpha < 0$. Prove that the ratio of radial and azimuthal periods is

$$\frac{T_r}{T_\psi} = \begin{cases} 1/\sqrt{2+\alpha} & \text{for } \alpha > 0 & \text{for nearly circular orbits} \\ 1/2, & \text{for } \alpha > 0 & \\ 1/(2+\alpha), & \text{for } \alpha < 0 & \text{for nearly radial orbits.} \end{cases} \quad (3.338)$$

What do these results imply for harmonic and Kepler potentials?

Hint: depending on the sign of α use a different approximation in the radical for v_r . For $b > 0$, $\int_1^\infty dx/(x\sqrt{x^b-1}) = \pi/b$ (see Touma & Tremaine 1997).

3.20 [1] Show that in spherical polar coordinates the Lagrangian for motion in the potential $\Phi(\mathbf{x})$ is

$$\mathcal{L} = \frac{1}{2}[\dot{r}^2 + (r\dot{\theta})^2 + (r \sin \theta \dot{\phi})^2] - \Phi(\mathbf{x}). \quad (3.339)$$

Hence show that the momenta p_θ and p_ϕ are related to the the magnitude and z -component of the angular-momentum vector \mathbf{L} by

$$p_\phi = L_z \quad ; \quad p_\theta^2 = L^2 - \frac{L_z^2}{\sin^2 \theta}. \quad (3.340)$$

3.21 [3] Plot a (y, \dot{y}) , $(x = 0, \dot{x} > 0)$ surface of section for motion in the potential Φ_L of equation (3.103) when $q = 0.9$ and $E = -0.337$. Qualitatively relate the structure of this surface of section to the structure of the (x, \dot{x}) surface of section shown in Figure 3.9.

3.22 [3] Sketch the structure of the (x, \dot{x}) , $(y = 0, \dot{y} > 0)$ surface of section for motion at energy E in a Kepler potential when (a) the (x, y) coordinates are inertial, and (b) the coordinates rotate at 0.75 times the circular frequency Ω at the energy E . Hint: see Binney, Gerhard, & Hut (1985).

3.23 [3] The Earth is flattened at the poles by its spin. Consequently orbits in its potential do not conserve total angular momentum. Many satellites are launched in inclined, nearly circular orbits only a few hundred kilometers above the Earth's surface, and their orbits must remain nearly circular, or they will enter the atmosphere and be destroyed. Why do the orbits remain nearly circular?

3.24 [2] Let $\hat{\mathbf{e}}_1$ and $\hat{\mathbf{e}}_2$ be unit vectors in an inertial coordinate system centered on the Sun, with $\hat{\mathbf{e}}_1$ pointing away from the Galactic center (towards $\ell = 180^\circ$, $b = 0$) and $\hat{\mathbf{e}}_2$ pointing towards $\ell = 270^\circ$, $b = 90^\circ$. The mean velocity field $\mathbf{v}(\mathbf{x})$ relative to the Local Standard of Rest can be expanded in a Taylor series,

$$v_i = \sum_{j=1}^2 H_{ij} x_j + O(x^2). \quad (3.341)$$

(a) Assuming that the Galaxy is stationary and axisymmetric, evaluate the matrix \mathbf{H} in terms of the Oort constants A and B .

(b) What is the matrix \mathbf{H} in a rotating frame, that is, if $\hat{\mathbf{e}}_1$ continues to point to the center of the Galaxy as the Sun orbits around it?

(c) In a homogeneous, isotropic universe, there is an analogous 3×3 matrix \mathbf{H} that describes the relative velocity \mathbf{v} between two fundamental observers separated by \mathbf{x} . Evaluate this matrix in terms of the Hubble constant.

3.25 [3] Consider two point masses m_1 and $m_2 > m_1$ that travel in a circular orbit about their center of mass under their mutual attraction. (a) Show that the Lagrange point L_4 of this system forms an equilateral triangle with the two masses. (b) Show that motion near L_4 is stable if $m_1/(m_1 + m_2) < 0.03852$. (c) Are the Lagrange points L_1, L_2, L_3 stable? See Valtonen & Karttunen (2006).

3.26 [2] Show that the leapfrog integrator (3.166a) is second-order accurate, in the sense that the errors in \mathbf{q} and \mathbf{p} after a timestep h are $O(h^3)$.

3.27 [2] Forest & Ruth (1990) have devised a symplectic, time-reversible, fourth-order integrator of timestep h by taking three successive drift-kick-drift leapfrog steps of length $ah, bh,$ and ah where $2a + b = 1$. Find a and b . Hint: a and b need not both be positive.

3.28 [2] Confirm the formulae for the Adams–Bashforth, Adams–Moulton, and Hermite integrators in equations (3.169), (3.170), and (3.171), and derive the next higher order integrator of each type. You may find it helpful to use computer algebra.

3.29 [1] Prove that the fictitious time τ in Burdet–Heggie regularization is related to the eccentric anomaly η by $\tau = (T_r/2\pi a)\eta + \text{constant}$, if the motion is bound ($E_2 < 0$) and the external field $\mathbf{g} = 0$.

3.30 [1] We wish to integrate numerically the motions of N particles with positions \mathbf{x}_i , velocities \mathbf{v}_i , and masses m_i . The particles interact only by gravitational forces (the gravitational N-body problem). We are considering using several possible integrators: modified Euler, leapfrog, or fourth-order Runge–Kutta. Which of these will conserve the total momentum $\sum_{i=1}^N m_i \mathbf{v}_i$? Which will conserve the total angular momentum $\sum_{i=1}^N m_i \mathbf{x}_i \times \mathbf{v}_i$? Assume that all particles are advanced with the same timestep, and that forces are calculated exactly. You may solve the problem either analytically or numerically.

3.31 [2] Show that the generating function of the canonical transformation from angle-action variables (θ_i, J_i) to the variables (q_i, p_i) discussed in Box 3.4 is

$$S(q, J) = \mp \frac{1}{2} q \sqrt{2J - q^2} \pm J \cos^{-1} \left(\frac{q}{\sqrt{2J}} \right). \quad (3.342)$$

3.32 [1] Let $\epsilon(R)$ and $\ell(R)$ be the specific energy and angular momentum of a circular orbit of radius R in the equatorial plane of an axisymmetric potential.

(a) Prove that

$$\frac{d\ell}{dR} = \frac{R\kappa^2}{2\Omega} \quad ; \quad \frac{d\epsilon}{dR} = \frac{1}{2} R\kappa^2, \quad (3.343)$$

where Ω and κ are the circular and epicycle frequencies.

(b) The energy of a circular orbit as a function of angular momentum is $\epsilon(\ell)$. Show that $d\epsilon/d\ell = \Omega$ in two ways, first from the results of part (a) and then using angle-action variables.

3.33 [2] The angle variables θ_i conjugate to the actions J_i can be implicitly defined by the coupled differential equations $dw_\alpha/d\theta_i = [w_\alpha, J_i]$, where w_α is any ordinary phase-space

coordinate. Using this result, show that the angle variable for the harmonic oscillator, $H = \frac{1}{2}(p^2 + \omega^2 q^2)$, may be written

$$\theta(x, p) = -\tan^{-1} \left(\frac{p}{\omega q} \right). \quad (3.344)$$

Hint: the action is $J = H/\omega$.

3.34 [2] Consider motion for $L_z = 0$ in the Stäckel potential (3.247).

(a) Express I_3 as a function of u , v , p_u , and p_v .

(b) Show that $H \cos^2 v + I_3 = \frac{1}{2}(p_v^2/\Delta^2) - V$.

(c) Show that $[H, I_3] = 0$.

(d) Hence show that J_u and J_v are in involution, that is $[J_u, J_v] = 0$. Hint: if $f(a, b)$ is any differentiable function of two variables, and A is any differentiable function of the phase-space variables, then $[A, f] = [A, a](\partial f/\partial a) + [A, b](\partial f/\partial b)$.

3.35 [2] A particle moves in a one-dimensional potential well $\Phi(x)$. In angle-action variables, the Hamiltonian has the form $H(J) = cJ^{4/3}$ where c is a constant. Find $\Phi(x)$.

3.36 [2] Obtain the Hamiltonian and fundamental frequencies as functions of the actions for the three-dimensional harmonic oscillator by examining the limit $b \rightarrow \infty$ of equations (3.226).

3.37 [2] For motion in a potential of the form (3.247), obtain

$$\dot{p}_u = \frac{2E \sinh u \cosh u - dU/du}{\sinh^2 u + \sin^2 v} + \frac{L_z^2 \cosh u}{\Delta^2 \sinh^3 u (\sinh^2 u + \sin^2 v)}, \quad (3.345)$$

where (u, v) are the prolate spheroidal coordinates defined by equations (3.242), by (a) differentiating equation (3.249a) with respect to t and then using $\dot{u} = \partial H/\partial p_u$, and (b) from $\dot{p}_u = -\partial H/\partial u$.

3.38 [2] For the coordinates defined by equation (3.267), show that the integral defined by equations (3.268) can be written

$$I_2 = \frac{\sinh^2 u [\frac{1}{2}(p_v^2/\Delta^2) - V] - \sin^2 v [\frac{1}{2}(p_u^2/\Delta^2) + U]}{\sinh^2 u + \sin^2 v}. \quad (3.346)$$

Show that in the limit $\Delta \rightarrow 0$, $u \rightarrow \infty$ we have $\Delta \sinh u \rightarrow \Delta \cosh u \rightarrow R$ and $v \rightarrow \pi/2 - \phi$, where R and ϕ are the usual polar coordinates. Hence show that in this limit $2\Delta^2 I_2 \rightarrow L_z^2$.

3.39 [2] Show that the third integral of an axisymmetric Stäckel potential can be taken to be

$$I_3(u, v, p_u, p_v, p_\phi) = \frac{1}{\sinh^2 u + \sin^2 v} \times \left[\sinh^2 u \left(\frac{p_v^2}{2\Delta^2} - V \right) - \sin^2 v \left(\frac{p_u^2}{2\Delta^2} + U \right) \right] + \frac{p_\phi^2}{2\Delta^2} \left(\frac{1}{\sin^2 v} - \frac{1}{\sinh^2 u} \right). \quad (3.347)$$

Hint: generalize the work of Problem 3.38.

3.40 [1] Show that when orbital frequencies are incommensurable, adiabatic invariance of actions implies that closed orbits remain closed when the potential is adiabatically deformed. An initially circular orbit in a spherical potential Φ does not remain closed when Φ is squashed along any line that is not parallel to the orbit's original angular-momentum vector. Why does this statement remain true no matter how slowly Φ is squashed?

3.41 [2] From equations (3.39b) and (3.190), show that the radial action J_r of an orbit in the isochrone potential (2.47) is related to the energy E and angular momentum L of this orbit by

$$J_r = \sqrt{GMb} \left[x^{-\frac{1}{2}} - f(L) \right], \quad (3.348)$$

where $x \equiv -2Eb/(GM)$ and f is some function. Use equation (3.327) to show that $f(L) = (\sqrt{l^2 + 1} - l)^{-1} = \sqrt{l^2 + 1} + l$, where $l \equiv |L|/(2\sqrt{GMb})$, and hence show that the isochrone Hamiltonian can be written in the form (3.226a).

3.42 [2] Angle-action variables are also useful in general relativity. For example, the relativistic analog to the Hamilton–Jacobi equation (3.218) for motion in the point-mass potential $\Phi(r) = -GM/r$ is

$$E^2 \left(\frac{1 + \frac{1}{4}r_S/r}{1 - \frac{1}{4}r_S/r} \right)^2 = c^4 + \frac{c^2}{(1 + \frac{1}{4}r_S/r)^4} \left[\left(\frac{\partial S}{\partial r} \right)^2 + \left(\frac{1}{r} \frac{\partial S}{\partial \vartheta} \right)^2 + \left(\frac{1}{r \sin \vartheta} \frac{\partial S}{\partial \phi} \right)^2 \right], \quad (3.349)$$

where $r_S \equiv 2GM/c^2$ is the **Schwarzschild radius**, the energy per unit mass E includes the rest-mass energy c^2 , and the equations are written in the isotropic metric, i.e., ds^2 at any point is proportional to its Euclidean form (Landau & Lifshitz 1999).

(a) Show that the Hamiltonian can be written in the form

$$H(p_r, p_\vartheta, p_\phi) = \frac{1 - \frac{1}{4}r_S/r}{1 + \frac{1}{4}r_S/r} \sqrt{c^4 + \frac{c^2 p^2}{(1 + \frac{1}{4}r_S/r)^4}}, \quad (3.350)$$

where $p^2 = p_r^2 + p_\vartheta^2/r^2 + p_\phi^2/(r \sin \vartheta)^2$.

(b) For systems in which relativistic effects are weak, show that the Hamiltonian can be written in the form

$$H = c^2 + H_{\text{Kep}} + H_{\text{gr}} + \mathcal{O}(c^{-4}), \quad (3.351)$$

where $H_{\text{Kep}} = \frac{1}{2}p^2 - GM/r$ is the usual Kepler Hamiltonian and

$$H_{\text{gr}} = \frac{1}{c^2} \left(\frac{G^2 M^2}{2r^2} - \frac{p^4}{8} - \frac{3GMp^2}{2r} \right). \quad (3.352)$$

(c) To investigate the long-term effects of relativistic corrections on a Kepler orbit, we may average H_{gr} over an unperturbed Kepler orbit. Show that this average may be written

$$\langle H_{\text{gr}} \rangle = \frac{G^2 M^2}{c^2 a^2} \left(\frac{15}{8} - \frac{3}{\sqrt{1 - e^2}} \right), \quad (3.353)$$

where a and e are the semi-major axis and eccentricity. Hint: use the results of Problem 3.9.

(d) Show that relativistic corrections cause the argument of pericenter ω to precess by an amount

$$\Delta\omega = \frac{6\pi GM}{c^2 a(1 - e^2)} \quad (3.354)$$

per orbit. Hint: convert $\langle H_{\text{gr}} \rangle$ to angle-action variables using Table E.1 and use Hamilton's equations.

3.43 [2] The Hamiltonian $H(\mathbf{x}, \mathbf{p}; \lambda)$, where λ is a parameter, supports a family of resonant orbits. In the (x_1, p_1) surface of section, the family's chain of islands is bounded by orbits with actions $J_1 \equiv (2\pi)^{-1} \oint dx_1 p_1 = J_\pm(\lambda)$, where $J_+ > J_-$. Let λ increase sufficiently slowly for the actions of non-resonant orbits to be conserved, and assume that $J'_+ > J'_- > 0$, where a prime denotes differentiation with respect to λ . Show that, as λ grows, an orbit of unknown phase and action slightly larger than J_+ will be captured by the resonance with probability $P_c = 1 - J'_-/J'_+$. Hint: exploit conservation of phase-space volume as expressed by equation (4.10).

4

Equilibria of Collisionless Systems

In §1.2 we introduced the idea that stellar systems may be considered to be collisionless: we obtain a good approximation to the orbit of any star by calculating the orbit that it would have if the system's mass were smoothly distributed in space rather than concentrated into nearly point-like stars. Eventually, the true orbit deviates significantly from this model orbit, but in systems with more than a few thousand stars, the deviation is small for a time $\lesssim t_{\text{relax}}$ that is much larger than the crossing time t_{cross} . In fact, for a galaxy t_{relax} is usually much larger even than the age of the universe, so the approximation that the potential is smooth provides a complete description of the dynamics.

In this chapter we consider model stellar systems that would be perfect equilibria if t_{relax} were arbitrarily large. Such models are the primary tool for comparisons of observations and theory of galaxy dynamics. In Chapter 7 we shall see that they are also applicable to globular clusters, even though t_{relax} is significantly smaller than the cluster's age, so long as it is recognized that the equilibrium evolves slowly, on a timescale of order t_{relax} .

We assume throughout that the stellar systems we examine consist of N identical point masses, which might be stars or dark-matter particles. Although unrealistic, this assumption greatly facilitates our work and has no impact on the validity of our results.

In §4.1 we derive the equation that allows us to find equilibria, and discuss its connection to observational data. In §4.2 we show that solutions of

the equation can be readily found if integrals of motion in the galactic potential are known, and in §§4.3 to 4.5 we use such solutions to study models with a variety of symmetries. In §4.6 we show that it is advantageous to express solutions in terms of action integrals. Unfortunately, in many practical cases insufficient integrals are known to obtain relevant solutions, so in §§4.7 and 4.8 we discuss alternative strategies, starting with heavily numerical approaches and moving on to approximate techniques that are based on moments of the fundamental equation. In §4.9 we draw on techniques developed throughout the chapter to hunt for massive black holes and dark halos in galaxies using observations of the kinematics of their stars. In §4.10 we address the question “what determines the distribution of stars in a galaxy?” This is a difficult question to which we shall have to return in Chapter 9.

4.1 The collisionless Boltzmann equation

When modeling a collisionless system such as an elliptical galaxy, it is neither practical nor worthwhile to follow the orbits of each of the galaxy’s billions of stars. Most testable predictions depend on the probability of finding a star in the six-dimensional phase-space volume $d^3\mathbf{x}d^3\mathbf{v}$ around the position \mathbf{x} and velocity \mathbf{v} . Therefore we define the **distribution function** (or DF for short) f such that $f(\mathbf{x}, \mathbf{v}, t)d^3\mathbf{x}d^3\mathbf{v}$ is the probability that at time t a randomly chosen star, say star 1, has phase-space coordinates in the given range. Since by assumption all stars are identical, this probability is the same for stars 2, 3, \dots , N . By virtue of its definition f is normalized such that

$$\int d^3\mathbf{x}d^3\mathbf{v} f(\mathbf{x}, \mathbf{v}, t) = 1, \quad (4.1)$$

where the integral is over all phase space.

Let $\mathbf{w} = (\mathbf{x}, \mathbf{v})$ be the usual Cartesian coordinates, and consider an arbitrary region \mathcal{V} of phase space. The probability of finding star 1 in \mathcal{V} is $P = \int_{\mathcal{V}} d^6\mathbf{w} f(\mathbf{w})$. Let \mathbf{W} represent some arbitrary set of phase-space coordinates, and let $F(\mathbf{W})$ be the corresponding DF; that is the probability of finding star 1 in \mathcal{V} is $P = \int_{\mathcal{V}} d^6\mathbf{W} F(\mathbf{W})$. If \mathcal{V} is small enough, f and F will be approximately constant throughout it, and we can take them outside the integrals for P . Thus

$$P = f(\mathbf{w}) \int_{\mathcal{V}} d^6\mathbf{w} = F(\mathbf{W}) \int_{\mathcal{V}} d^6\mathbf{W}. \quad (4.2)$$

If the coordinates \mathbf{W} are canonical, equation (D.81) implies that $\int_{\mathcal{V}} d^6\mathbf{w} = \int_{\mathcal{V}} d^6\mathbf{W}$. Substituting this relation into (4.2), we conclude that $F(\mathbf{W}) = f(\mathbf{w})$. Therefore, the DF has the same numerical value at a given phase-space point in *any* canonical coordinate system. This invariance enables us

henceforth to treat $\mathbf{w} = (\mathbf{q}, \mathbf{p})$ as an arbitrary system of canonical coordinates.

Any given star moves through phase space, so the probability of finding it at any given phase-space location evolves with time. We now derive the differential equation that is satisfied by f as a consequence of this evolution. As f evolves, probability must be conserved, in the same way that mass is conserved in a fluid flow. The conservation of fluid mass is described by the continuity equation (F.3)

$$\frac{\partial \rho}{\partial t} + \frac{\partial}{\partial \mathbf{x}} \cdot (\rho \dot{\mathbf{x}}) = 0, \quad (4.3)$$

where ρ and $\dot{\mathbf{x}} = \mathbf{v}$ are the density and velocity of the fluid. The analogous equation for the conservation of probability in phase space is

$$\frac{\partial f}{\partial t} + \frac{\partial}{\partial \mathbf{w}} \cdot (f \dot{\mathbf{w}}) = 0. \quad (4.4)$$

We now use Hamilton's equations (D.54) to eliminate $\dot{\mathbf{w}} = (\dot{\mathbf{q}}, \dot{\mathbf{p}})$. The second term in equation (4.4) becomes

$$\begin{aligned} \frac{\partial}{\partial \mathbf{q}} \cdot (f \dot{\mathbf{q}}) + \frac{\partial}{\partial \mathbf{p}} \cdot (f \dot{\mathbf{p}}) &= \frac{\partial}{\partial \mathbf{q}} \cdot \left(f \frac{\partial H}{\partial \mathbf{p}} \right) - \frac{\partial}{\partial \mathbf{p}} \cdot \left(f \frac{\partial H}{\partial \mathbf{q}} \right) \\ &= \frac{\partial f}{\partial \mathbf{q}} \cdot \frac{\partial H}{\partial \mathbf{p}} - \frac{\partial f}{\partial \mathbf{p}} \cdot \frac{\partial H}{\partial \mathbf{q}} \\ &= \dot{\mathbf{q}} \cdot \frac{\partial f}{\partial \mathbf{q}} + \dot{\mathbf{p}} \cdot \frac{\partial f}{\partial \mathbf{p}}, \end{aligned} \quad (4.5)$$

where we have used the fact that $\partial^2 H / \partial \mathbf{q} \partial \mathbf{p} = \partial^2 H / \partial \mathbf{p} \partial \mathbf{q}$. Substituting this result into equation (4.4) we obtain the **collisionless Boltzmann equation**¹

$$\frac{\partial f}{\partial t} + \dot{\mathbf{q}} \cdot \frac{\partial f}{\partial \mathbf{q}} + \dot{\mathbf{p}} \cdot \frac{\partial f}{\partial \mathbf{p}} = 0, \quad (4.6)$$

which is a partial differential equation for f as a function of six phase-space coordinates and time.

Equation (4.6) can be rewritten in a number of forms, each of which is useful in different contexts. Equation (4.5) enables us to write

$$\begin{aligned} 0 &= \frac{\partial f}{\partial t} + \frac{\partial f}{\partial \mathbf{q}} \cdot \frac{\partial H}{\partial \mathbf{p}} - \frac{\partial f}{\partial \mathbf{p}} \cdot \frac{\partial H}{\partial \mathbf{q}} \\ &= \frac{\partial f}{\partial t} + [f, H], \end{aligned} \quad (4.7)$$

¹ Often also called the Vlasov equation, although it is a simplified version of an equation derived by L. Boltzmann in 1872. See Hénon (1982).

where the square bracket is a Poisson bracket (eq. D.65).

An alternative form of the collisionless Boltzmann equation can be derived by extending to six dimensions the concept of the convective or Lagrangian derivative (see eq. F.8). We define

$$\frac{df}{dt} \equiv \frac{\partial f}{\partial t} + \dot{\mathbf{w}} \cdot \frac{\partial f}{\partial \mathbf{w}}; \quad (4.8)$$

df/dt represents the rate of change of the local probability density as seen by an observer who moves through phase space with a star. Comparison of equations (4.6) and (4.7) shows that $\dot{\mathbf{w}} \cdot (\partial f / \partial \mathbf{w}) = [f, H]$, so the convective derivative can also be written

$$\frac{df}{dt} = \frac{\partial f}{\partial t} + [f, H], \quad (4.9)$$

and the collisionless Boltzmann equation (4.6) is simply

$$\frac{df}{dt} = 0. \quad (4.10)$$

In words, the flow through phase space of the probability fluid is incompressible; the phase-space density f of the fluid around a given star always remains the same.² In contrast to flows of incompressible fluids such as water, the density will generally vary greatly from point to point in phase space; the density is constant as one follows the flow around a particular star but the density around different stars can be quite different.

In terms of inertial Cartesian coordinates, in which $H = \frac{1}{2}v^2 + \Phi(\mathbf{x}, t)$ with Φ the gravitational potential, the collisionless Boltzmann equation reads

$$\frac{\partial f}{\partial t} + \mathbf{v} \cdot \frac{\partial f}{\partial \mathbf{x}} - \frac{\partial \Phi}{\partial \mathbf{x}} \cdot \frac{\partial f}{\partial \mathbf{v}} = 0. \quad (4.11)$$

In cylindrical coordinates we have (eq. 3.66) $H = \frac{1}{2}(p_R^2 + p_\phi^2/R^2 + p_z^2) + \Phi$ so with (4.7) the collisionless Boltzmann equation becomes³

$$\begin{aligned} \frac{\partial f}{\partial t} + p_R \frac{\partial f}{\partial R} + \frac{p_\phi}{R^2} \frac{\partial f}{\partial \phi} + p_z \frac{\partial f}{\partial z} - \left(\frac{\partial \Phi}{\partial R} - \frac{p_\phi^2}{R^3} \right) \frac{\partial f}{\partial p_R} \\ - \frac{\partial \Phi}{\partial \phi} \frac{\partial f}{\partial p_\phi} - \frac{\partial \Phi}{\partial z} \frac{\partial f}{\partial p_z} = 0. \end{aligned} \quad (4.12)$$

² A simple example of an incompressible flow in phase space is provided by an idealized marathon race in which all runners travel at constant speeds: at the start of the course, the spatial density of runners is large but they travel at a wide variety of speeds; at the finish, the density is low, but at any given time all runners passing the finish line have nearly the same speed.

³ A reader unconvinced of the usefulness of the Hamiltonian formalism should try deriving either (4.12) or (4.14) directly from (4.11).

To obtain the Hamiltonian for motion in spherical polar coordinates we replace in (3.218) $\partial S/\partial r$ by p_r , $\partial S/\partial\theta$ by p_θ and $\partial S/\partial\phi$ by p_ϕ and find

$$H = \frac{1}{2} \left(p_r^2 + \frac{p_\theta^2}{r^2} + \frac{p_\phi^2}{r^2 \sin^2 \theta} \right) + \Phi. \quad (4.13)$$

Using this expression in (4.7) we find

$$\begin{aligned} \frac{\partial f}{\partial t} + p_r \frac{\partial f}{\partial r} + \frac{p_\theta}{r^2} \frac{\partial f}{\partial \theta} + \frac{p_\phi}{r^2 \sin^2 \theta} \frac{\partial f}{\partial \phi} - \left(\frac{\partial \Phi}{\partial r} - \frac{p_\theta^2}{r^3} - \frac{p_\phi^2}{r^3 \sin^2 \theta} \right) \frac{\partial f}{\partial p_r} \\ - \left(\frac{\partial \Phi}{\partial \theta} - \frac{p_\phi^2 \cos \theta}{r^2 \sin^3 \theta} \right) \frac{\partial f}{\partial p_\theta} - \frac{\partial \Phi}{\partial \phi} \frac{\partial f}{\partial p_\phi} = 0. \end{aligned} \quad (4.14)$$

Conversion to rotating coordinates is discussed in Problem 4.1.

4.1.1 Limitations of the collisionless Boltzmann equation

(a) Finite stellar lifetimes The physical basis of the collisionless Boltzmann equation is conservation of the objects that are described by the DF. Stars are not really conserved because they are born and die, so their flow through phase space would be more accurately described by an equation of the type

$$\frac{df}{dt} = \frac{\partial f}{\partial t} + \mathbf{v} \cdot \frac{\partial f}{\partial \mathbf{x}} - \frac{\partial \Phi}{\partial \mathbf{x}} \cdot \frac{\partial f}{\partial \mathbf{v}} = B - D, \quad (4.15)$$

where $B(\mathbf{x}, \mathbf{v}, t)$ and $D(\mathbf{x}, \mathbf{v}, t)$ are the rates per unit phase-space volume at which stars are born and die. In the collisionless Boltzmann equation, $B - D$ is set to zero. This is a useful approximation to the truth if and only if $B - D$ is smaller in magnitude than terms on the left of equation (4.15). The term $\mathbf{v} \cdot \partial f/\partial \mathbf{x}$ is of order $v f/R$, where v and R are the characteristic speed and radius in the galaxy. The ratio R/v is simply the crossing time t_{cross} (§1.2). Similarly, $\partial \Phi/\partial \mathbf{x}$ is of order the characteristic acceleration a , so the term $(\partial \Phi/\partial \mathbf{x}) \cdot (\partial f/\partial \mathbf{v})$ is of order $a f/v$. Since $a \approx v/t_{\text{cross}}$, the two last terms in the middle section of equation (4.15) are of order f/t_{cross} . Thus consider the ratio

$$\gamma = \left| \frac{B - D}{f/t_{\text{cross}}} \right|. \quad (4.16)$$

The collisionless Boltzmann equation is valid if $\gamma \ll 1$, which requires that the fractional change in the number of stars per crossing time is small.

The significance of this criterion can be clarified by some concrete examples. We consider two contrasting stellar types: M dwarfs, which have masses $\lesssim 0.5 \mathcal{M}_\odot$ and live longer than the age of the universe; and O stars, which have masses $\gtrsim 20 \mathcal{M}_\odot$ and lifetimes $\lesssim 10$ Myr (BM Tables 3.13 and 5.3). In an elliptical galaxy the rate of formation of M dwarfs is negligible, and

even the oldest M dwarfs have not had time to evolve significantly. Hence, the collisionless Boltzmann equation will apply accurately to the DF of M dwarfs ($\gamma \leq 0.01$). Now consider the contrasting case of O stars in the Milky Way. These stars have lifetimes significantly shorter than a crossing time ~ 100 Myr. In fact, an O star will scarcely move from its birthplace before it dies, and the phase-space distribution of such stars will depend entirely on the processes that govern star formation, and not at all on the collisionless Boltzmann equation ($\gamma \simeq 10$). In between these extremes, the collisionless Boltzmann equation will apply quite accurately to main-sequence populations in the Milky Way less massive than $\sim 1.5 \mathcal{M}_\odot$, since these stars live for $\gtrsim 1$ Gyr, which will generally be some tens of crossing times. In certain circumstances the collisionless Boltzmann equation may even be applied to a population of short-lived objects, for example planetary nebulae in an elliptical galaxy, because the phase-space distributions of the objects' births and deaths are to a good approximation identical, so $B - D \simeq 0$.

(b) Correlations between stars The average number density of stars in an infinitesimal volume of phase space is Nf . However, in practice all we can hope to measure is the number density in some volume of phase space large enough to contain many stars. The natural assumption to make is that the density in such a volume is simply $N\bar{f}$, where \bar{f} is the average of f within this volume.⁴ However, this assumption will only be correct if the positions of stars in phase space are uncorrelated: that is, knowing that star 1 is at \mathbf{w} makes it neither more nor less likely that another star, say star 2, is at an adjacent phase-space location \mathbf{w}' . Mathematically, we assume that the probability of finding star 1 in the volume $d^6\mathbf{w}$ at \mathbf{w} and star 2 in $d^6\mathbf{w}'$ at \mathbf{w}' is simply the product $f(\mathbf{w})d^6\mathbf{w} f(\mathbf{w}')d^6\mathbf{w}'$ of the probabilities of finding star 1 at \mathbf{w} and star 2 at \mathbf{w}' —in §7.2.4 we shall call such distributions “separable”. When the assumption of separability holds, the probability $P_{\mathcal{V}}(k)$ that we will find k stars in a given volume \mathcal{V} of phase space is given by the Poisson distribution (Appendix B.8)

$$P_{\mathcal{V}}(k) = \frac{\mu^k}{k!} e^{-\mu} \quad \text{where} \quad \mu \equiv N\bar{f}\mathcal{V}. \quad (4.17)$$

It is easy to show that the mean number of stars predicted by this probability distribution is $\langle k \rangle = N\bar{f}\mathcal{V}$. Thus $N\bar{f}$ is indeed the expectation value of the stellar number density, if the DF is separable. Two obvious corollaries are that the mean mass within \mathcal{V} is

$$\langle m \rangle = M\bar{f}(\mathbf{w})\mathcal{V}, \quad (4.18)$$

where M is the total mass of the stellar system, and the mean luminosity emitted within \mathcal{V} is

$$\langle l \rangle = L\bar{f}(\mathbf{w})\mathcal{V}, \quad (4.19)$$

⁴This function is sometimes called the **coarse-grained** DF. The standard DF is then called the **fine-grained** DF to eliminate any danger of confusion with \bar{f} .

where L is the system's luminosity.

In reality, the presence of star 1 at \mathbf{x} always increases the probability that star 2 will be found at some nearby position \mathbf{x}' because stars attract one another. Hence, the assumption that the probability distribution of individual stars is separable is never strictly valid. In Chapter 7 we shall explore the effect of such correlations on the evolution of stellar systems. However, in this chapter we assume that separability holds, as it very nearly does for many stellar systems, because the force on a star from its neighbors is very much smaller than the force from the rest of the system.

4.1.2 Relation between the DF and observables

At any fixed position \mathbf{x} , the integral

$$\nu(\mathbf{x}) \equiv \int d^3\mathbf{v} f(\mathbf{x}, \mathbf{v}) \quad (4.20)$$

gives the probability per unit volume of finding a particular star at \mathbf{x} , regardless of its velocity. Multiplying by the total number N of stars in the population, we obtain the real-space number density of stars

$$n(\mathbf{x}) \equiv N\nu(\mathbf{x}). \quad (4.21)$$

In the Galaxy $n(\mathbf{x})$ can in principle be determined from star counts, and thus $\nu(\mathbf{x})$ can be derived from $n(\mathbf{x})$. In other galaxies it is not usually possible to count stars, but we can derive $\nu(\mathbf{x})$ from the luminosity density $j(\mathbf{x}) = L\nu(\mathbf{x})$, where L is the luminosity of the stellar population (BM §4.2.3).

It is often convenient to modify the definition of the DF so that $f d^6\mathbf{w}$ represents not the probability of finding a given star in the phase-space volume $d^6\mathbf{w}$, but rather the expected number, total mass, or total luminosity of the stars in $d^6\mathbf{w}$. These modifications correspond to multiplying f by N , M , or L , respectively. Ideally these different definitions would be reflected in different notations for the DF. In practice the definition is usually clear from the context, and f is conventionally used to denote all of these quantities.

Dividing f by ν we obtain the probability distribution of stellar velocities at \mathbf{x}

$$P_{\mathbf{x}}(\mathbf{v}) = \frac{f(\mathbf{x}, \mathbf{v})}{\nu(\mathbf{x})}, \quad (4.22)$$

which can be directly measured near the Sun (BM §10.3). In external galaxies $P_{\mathbf{x}}$ can be probed through the **line-of-sight velocity distribution** (LOSVD; BM §11.1), which gives for a particular line of sight through the galaxy the fraction $F(v_{\parallel})dv_{\parallel}$ of the stars that have line-of-sight velocity within dv_{\parallel} of v_{\parallel} . Almost all galaxies are sufficiently far away that all vectors from the observer to a point \mathbf{x} in the galaxy are very nearly parallel to the fixed unit

vector $\hat{\mathbf{s}}$ from the observer to the center of the galaxy. Then $x_{\parallel} \equiv \hat{\mathbf{s}} \cdot \mathbf{x}$ and $v_{\parallel} \equiv \hat{\mathbf{s}} \cdot \mathbf{v}$ are the components of \mathbf{x} and \mathbf{v} parallel to the line of sight. We also define $\mathbf{x}_{\perp} \equiv \mathbf{x} - x_{\parallel}\hat{\mathbf{s}}$ and $\mathbf{v}_{\perp} \equiv \mathbf{v} - v_{\parallel}\hat{\mathbf{s}}$ to be the components of \mathbf{x} and \mathbf{v} in the plane of the sky. The relation between $P_{\mathbf{x}}(\mathbf{v})$ and $F(\mathbf{x}_{\perp}, v_{\parallel})$ is

$$\begin{aligned} F(\mathbf{x}_{\perp}, v_{\parallel}) &= \frac{\int dx_{\parallel} \nu(\mathbf{x}) \int d^2\mathbf{v}_{\perp} P_{\mathbf{x}}(v_{\parallel}\hat{\mathbf{s}} + \mathbf{v}_{\perp})}{\int dx_{\parallel} \nu(\mathbf{x})} \\ &= \frac{\int dx_{\parallel} d^2\mathbf{v}_{\perp} f(\mathbf{x}, \mathbf{v})}{\int dx_{\parallel} d^3\mathbf{v} f(\mathbf{x}, \mathbf{v})}. \end{aligned} \quad (4.23)$$

The LOSVD is frequently quantified by two numbers, the mean line-of-sight velocity \bar{v}_{\parallel} and the dispersion σ_{\parallel} about this mean. We have

$$\begin{aligned} \bar{v}_{\parallel}(\mathbf{x}_{\perp}) &\equiv \int dv_{\parallel} v_{\parallel} F(\mathbf{x}_{\perp}, v_{\parallel}) = \frac{\int dx_{\parallel} d^3\mathbf{v} v_{\parallel} f(\mathbf{x}, \mathbf{v})}{\int dx_{\parallel} d^3\mathbf{v} f(\mathbf{x}, \mathbf{v})} \\ &= \frac{\int dx_{\parallel} \nu(\mathbf{x}) \hat{\mathbf{s}} \cdot \bar{\mathbf{v}}}{\int dx_{\parallel} \nu(\mathbf{x})}, \end{aligned} \quad (4.24a)$$

where we have defined the **mean velocity** at location \mathbf{x}

$$\bar{\mathbf{v}}(\mathbf{x}) \equiv \int d^3\mathbf{v} \mathbf{v} P_{\mathbf{x}}(\mathbf{v}) = \frac{1}{\nu(\mathbf{x})} \int d^3\mathbf{v} \mathbf{v} f(\mathbf{x}, \mathbf{v}). \quad (4.24b)$$

The **line-of-sight velocity dispersion** is defined to be

$$\begin{aligned} \sigma_{\parallel}^2(\mathbf{x}_{\perp}) &\equiv \int dv_{\parallel} (v_{\parallel} - \bar{v}_{\parallel})^2 F(\mathbf{x}_{\perp}, v_{\parallel}) \\ &= \frac{\int dx_{\parallel} d^3\mathbf{v} (\hat{\mathbf{s}} \cdot \mathbf{v} - \bar{v}_{\parallel})^2 f(\mathbf{x}, \mathbf{v})}{\int dx_{\parallel} d^3\mathbf{v} f(\mathbf{x}, \mathbf{v})}. \end{aligned} \quad (4.25)$$

The line-of-sight velocity dispersion is determined both by the variation in the mean velocity $\bar{v}_{\parallel}(\mathbf{x})$ along the line of sight, and the spread in stellar velocities at each point in the galaxy around $\bar{\mathbf{v}}(\mathbf{x})$. This spread is characterized by the **velocity-dispersion tensor**

$$\begin{aligned} \sigma_{ij}^2(\mathbf{x}) &\equiv \frac{1}{\nu(\mathbf{x})} \int d^3\mathbf{v} (v_i - \bar{v}_i)(v_j - \bar{v}_j) f(\mathbf{x}, \mathbf{v}) \\ &= \overline{v_i v_j} - \bar{v}_i \bar{v}_j. \end{aligned} \quad (4.26)$$

The velocity-dispersion tensor is manifestly symmetric, so we know from matrix algebra that at any point \mathbf{x} we may choose a set of orthogonal axes $\hat{\mathbf{e}}_i(\mathbf{x})$ in which σ^2 is diagonal, that is, $\sigma_{ij}^2 = \sigma_{ii}^2 \delta_{ij}$ (no summation over i , and $\delta_{ij} = 1$ for $i = j$ and zero otherwise). The ellipsoid that has the

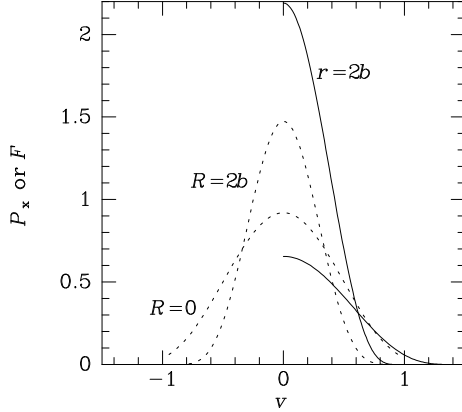


Figure 4.1 The full curves show the velocity distributions $P_{\mathbf{x}}(v)$ at the center of a Plummer model (lower curve) and at $r = 2b$ (upper curve). The dashed curves show the LOSVD $F(v_{\parallel})$ along two lines of sight, $R = 0, 2b$. The DF of the Plummer model is given by equation (4.91b).

diagonalizing coordinate axes $\hat{\mathbf{e}}_i(\mathbf{x})$ for its principal axes and σ_{11} , σ_{22} and σ_{33} for its semi-axis lengths is called the **velocity ellipsoid** at \mathbf{x} .

To determine the relation between the velocity-dispersion tensor and the line-of-sight velocity dispersion, we let $u(\mathbf{x}) \equiv \hat{\mathbf{s}} \cdot \bar{\mathbf{v}}(\mathbf{x}) - \bar{v}_{\parallel}$ be the difference between the mean velocity parallel to the line of sight at \mathbf{x} and the mean velocity for the entire line of sight. Then we can rewrite (4.25) as follows:

$$\begin{aligned} \sigma_{\parallel}^2(\mathbf{x}_{\perp}) &= \frac{\int dx_{\parallel} d^3\mathbf{v} [\hat{\mathbf{s}} \cdot (\mathbf{v} - \bar{\mathbf{v}}) + u]^2 f(\mathbf{x}, \mathbf{v})}{\int dx_{\parallel} d^3\mathbf{v} f(\mathbf{x}, \mathbf{v})} \\ &= \frac{\int dx_{\parallel} \nu(\mathbf{x}) (\hat{\mathbf{s}} \cdot \boldsymbol{\sigma}^2 \cdot \hat{\mathbf{s}} + u^2)}{\int dx_{\parallel} \nu(\mathbf{x})}, \end{aligned} \quad (4.27)$$

where we introduce the notation $\hat{\mathbf{s}} \cdot \boldsymbol{\sigma}^2 \cdot \hat{\mathbf{s}} \equiv \sum_{ij} \hat{s}_i \sigma_{ij}^2 \hat{s}_j$.

These results show that once ν , $\bar{\mathbf{v}}$ and $\boldsymbol{\sigma}^2$ are known at each point in a model, the observable quantities v_{\parallel} and σ_{\parallel}^2 can be determined for that model. This fact makes ν , $\bar{\mathbf{v}}$ and σ_{ij}^2 , all functions of \mathbf{x} , vital links between observations and theoretical models. Moreover, we shall see in §4.8 that in equilibrium stellar systems there are simple relations between these quantities and the gravitational field (the Jeans equations).

Notice that while v_{\parallel} depends only on the mean velocity field $\bar{\mathbf{v}}(\mathbf{x})$, there are contributions to σ_{\parallel}^2 from both $\boldsymbol{\sigma}^2$ and $\bar{\mathbf{v}}$. Moreover, both contributions are inherently positive, so σ_{\parallel}^2 is in general larger than the average of the intrinsic squared velocity dispersion $\hat{\mathbf{s}} \cdot \boldsymbol{\sigma}^2 \cdot \hat{\mathbf{s}}$ along the line of sight.

One shortcoming of \bar{v}_{\parallel} and σ_{\parallel}^2 as probes of the dynamics of a galaxy is that they are hard to measure accurately because they are sensitive to the contributions of the small number of high-velocity stars.

An example In §4.3.3a we shall encounter an exceptionally simple model system called the Plummer model. This is a non-rotating spherical system

in which the velocity distribution $P_{\mathbf{x}}$ depends only on $v \equiv |\mathbf{v}|$, and the gravitational potential is given by equation (2.44a). The full curves in Figure 4.1 show $P_{\mathbf{x}}(v)$ at the center of the system and at $r = 2b$, where b is the Plummer scale length. Notice that $P_{\mathbf{x}}$ vanishes for speeds larger than the escape speed $\sqrt{2|\Phi(\mathbf{x})|}$ (eq. 2.31). At small radii, where $|\Phi|$ is relatively large, a graph of $P_{\mathbf{x}}$ versus v is wide and gently peaked, while at large radii, where $|\Phi|$ is much smaller, a plot of $P_{\mathbf{x}}(v)$ shows a high, narrow peak.

The LOSVD $F(v_{\parallel})$ along a line of sight through a Plummer model depends on the projected distance $R = |\mathbf{x}_{\perp}|$ between the line of sight and the center of the model because it is a weighted mean of the velocity distributions $P_{\mathbf{x}}(v)$ for different points along the line of sight (eq. 4.23). The dashed curves in Figure 4.1 show the LOSVD for the line of sight through the center, and one further out. Notice that the LOSVD at each radius is more centrally peaked than the velocity distribution at that radius. There are two reasons for this phenomenon. First, $F(v_{\parallel})$ is depressed at large values of v_{\parallel} by the integral over \mathbf{v}_{\perp} in (4.23) because the range of allowed values of \mathbf{v}_{\perp} diminishes rapidly as v_{\parallel} approaches the escape speed. A subsidiary effect is that the line of sight through the center samples points that are physically far from the cluster center, where the velocity distribution $P_{\mathbf{x}}$ is narrowly peaked around $v = 0$.

4.2 Jeans theorems

In §3.1.1 we introduced the concept of an integral of motion in a given stationary potential $\Phi(\mathbf{x})$. According to equation (3.56), a function of the phase-space coordinates $I(\mathbf{x}, \mathbf{v})$ is an integral if and only if

$$\frac{d}{dt}I[\mathbf{x}(t), \mathbf{v}(t)] = 0 \quad (4.28)$$

along any orbit. With the equations of motion this becomes

$$\frac{dI}{dt} = \frac{\partial I}{\partial \mathbf{x}} \cdot \frac{d\mathbf{x}}{dt} + \frac{\partial I}{\partial \mathbf{v}} \cdot \frac{d\mathbf{v}}{dt} = 0, \quad \text{or} \quad \mathbf{v} \cdot \frac{\partial I}{\partial \mathbf{x}} - \frac{\partial \Phi}{\partial \mathbf{x}} \cdot \frac{\partial I}{\partial \mathbf{v}} = 0. \quad (4.29)$$

Comparing this with equation (4.11), we see that the condition for I to be an integral is identical with the condition for I to be a steady-state solution of the collisionless Boltzmann equation. This leads to the following theorem, first stated by Jeans (1915).

Jeans theorem *Any steady-state solution of the collisionless Boltzmann equation depends on the phase-space coordinates only through integrals of motion in the given potential, and any function of the integrals yields a steady-state solution of the collisionless Boltzmann equation.*

Proof: Suppose f is a steady-state solution of the collisionless Boltzmann equation. Then, as we have just seen, f is an integral, and the first part of the theorem is proved. Conversely, if I_1 to I_n are n integrals, and if f is any function of n variables, then

$$\frac{d}{dt} f [I_1(\mathbf{x}, \mathbf{v}), \dots, I_n(\mathbf{x}, \mathbf{v})] = \sum_{m=1}^n \frac{\partial f}{\partial I_m} \frac{dI_m}{dt} = 0 \quad (4.30)$$

and f is seen to satisfy the collisionless Boltzmann equation. ◁

Many of the results of this chapter will be based on the second proposition stated by the Jeans theorem, namely, that any function of integrals solves the collisionless Boltzmann equation. However, the first of Jeans's propositions, the assurance that the DF of any steady-state galaxy must be a function of integrals, is only of limited use since the examples discussed in §§3.2 and 3.3 show that orbits often respect integrals for which we lack analytic expressions. In such cases the first part of the Jeans theorem simply tells us, unhelpfully, that the DF is a function of integrals whose form is unknown.

Fortunately the time averages theorem (page 215) enables us to show that if all orbits in a galaxy are regular, then we can forget about any non-isolating integrals:

Strong Jeans theorem *The DF of a steady-state stellar system in which almost all orbits are regular with non-resonant frequencies may be presumed to be a function only of three independent isolating integrals, which may be taken to be the actions.*

Proof: Any observable property involves averaging the DF over some non-zero region of phase space; we may formalize this by stating that all observations are based on moments of the form $\langle Q \rangle = \int d^3\mathbf{x} d^3\mathbf{v} Q f$, where $Q(\mathbf{x}, \mathbf{v})$ is some smooth function on phase space and f is the DF. Since almost all orbits are regular, we can assume that phase space is covered by angle-action coordinates $(\boldsymbol{\theta}, \mathbf{J})$ and we can also write

$$\langle Q \rangle = \int d^3\boldsymbol{\theta} d^3\mathbf{J} Q(\mathbf{x}, \mathbf{v}) f(\mathbf{x}, \mathbf{v}, t), \quad (4.31)$$

where (\mathbf{x}, \mathbf{v}) are to be interpreted as functions of $(\boldsymbol{\theta}, \mathbf{J})$. Since the stellar system is in a steady state, $\langle Q \rangle$ is time-independent, so

$$\langle Q \rangle = \overline{\langle Q \rangle} = \int d^3\boldsymbol{\theta} d^3\mathbf{J} Q(\boldsymbol{\theta}, \mathbf{J}) \bar{f}(\boldsymbol{\theta}, \mathbf{J}), \quad (4.32a)$$

where

$$\bar{f}(\boldsymbol{\theta}, \mathbf{J}) \equiv \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T dt f(\mathbf{x}, \mathbf{v}, t) \quad (4.32b)$$

is the time average of the DF at the given point in phase space. If $\delta^3\mathbf{J}$ and $\delta^3\boldsymbol{\theta}$ are small coordinate ranges, $\overline{f} \delta^3\boldsymbol{\theta} \delta^3\mathbf{J}$ is the probability that at a given instant a randomly chosen star has actions in $\delta^3\mathbf{J}$ and angles in $\delta^3\boldsymbol{\theta}$. By the time averages theorem this is $\delta^3\boldsymbol{\theta}/(2\pi)^3$ times the probability that the star has actions in $\delta^3\mathbf{J}$, which we denote $(2\pi)^3 f_{\mathbf{J}}(\mathbf{J}) \delta^3\mathbf{J}$. Eliminating \overline{f} from equation (4.32a) we may therefore write

$$\langle Q \rangle = \int d^3\boldsymbol{\theta} d^3\mathbf{J} Q(\boldsymbol{\theta}, \mathbf{J}) f_{\mathbf{J}}(\mathbf{J}), \quad (4.33)$$

which shows that any observable can be evaluated using only the time-independent DF $f_{\mathbf{J}}(\mathbf{J})$.⁴

In summary, the Jeans theorem tells us that if I_1, \dots, I_n are n independent integrals in a given potential, then any DF of the form $f(I_1), f(I_1, I_2), \dots, f(I_1, \dots, I_n)$ is a solution of the collisionless Boltzmann equation. The strong Jeans theorem tells us that if the potential of a steady-state galaxy is such that almost all orbits are regular, then for all practical purposes the galaxy may be represented by a DF of the form $f(I_1, I_2, I_3)$, where I_1, I_2, I_3 are three independent isolating integrals.

4.2.1 Choice of f and relations between moments

(a) DF depending only on H In any steady-state potential $\Phi(\mathbf{x})$, the Hamiltonian H is an integral of motion. Consequently, an equilibrium stellar system is obtained by taking f to be any non-negative function of H —DFs of this type are called **ergodic**.⁵ If the potential is constant in an inertial frame, H will be of the form $H = \frac{1}{2}v^2 + \Phi(\mathbf{x})$ and it follows that the mean velocity vanishes everywhere:

$$\overline{\mathbf{v}}(\mathbf{x}) = \frac{1}{\nu(\mathbf{x})} \int d^3\mathbf{v} \mathbf{v} f\left(\frac{1}{2}v^2 + \Phi\right) = 0, \quad (4.34)$$

where the second equality follows because the integrand is an odd function of \mathbf{v} and the integral is over all velocity space. A similar line of reasoning shows that the velocity-dispersion tensor is isotropic:

$$\sigma_{ij}^2 = \overline{v_i v_j} = \sigma^2 \delta_{ij}, \quad (4.35a)$$

where

$$\begin{aligned} \sigma^2(\mathbf{x}) &= \frac{1}{\nu(\mathbf{x})} \int dv_z v_z^2 \int dv_x dv_y f\left[\frac{1}{2}(v_x^2 + v_y^2 + v_z^2) + \Phi(\mathbf{x})\right] \\ &= \frac{4\pi}{3\nu(\mathbf{x})} \int_0^\infty dv v^4 f\left(\frac{1}{2}v^2 + \Phi\right). \end{aligned} \quad (4.35b)$$

⁵ In statistical mechanics and chaos theory the term “ergodic” denotes a system that uniformly explores its energy surface in phase space, which implies that the DF is uniform on the energy surface. In our usage the DF is ergodic but the motion of individual stars generally is not.

An equivalent statement is that the velocity ellipsoid is a sphere of radius σ . An **isotropic system** has a velocity-dispersion tensor that is everywhere isotropic. Thus every system with an ergodic DF is isotropic.

(b) DF depending on H and L If the potential $\Phi(\mathbf{x})$ is spherically symmetric, the three components of the angular-momentum vector \mathbf{L} are isolating integrals that can be included in the arguments of f in addition to the Hamiltonian. Since the potential is spherical, we shall confine ourselves to DFs that produce systems that have complete spherical symmetry; that is, in which the three components of \mathbf{L} can occur only through their contributions to $L = |\mathbf{L}|$.⁶ Thus we consider the case in which the DF is some non-negative function $f(H, L)$. Let v_r and \mathbf{v}_t (for tangential) be the components of \mathbf{v} parallel and perpendicular to the radial direction, so $v_t^2 = v_\theta^2 + v_\phi^2$ in spherical coordinates (r, θ, ϕ) . Then $L = rv_t$ and $H = \frac{1}{2}(v_r^2 + v_t^2) + \Phi(r)$, and the mean velocity is

$$\begin{aligned}\bar{v}_r &= \frac{1}{\nu} \int dv_r v_r \int d^2\mathbf{v}_t f[\tfrac{1}{2}(v_r^2 + v_t^2) + \Phi(r), rv_t] = 0, \\ \bar{\mathbf{v}}_t &= \frac{1}{\nu} \int d^2\mathbf{v}_t \mathbf{v}_t \int dv_r f[\tfrac{1}{2}(v_r^2 + v_t^2) + \Phi(r), rv_t] = 0.\end{aligned}\tag{4.36}$$

In both cases the integrals vanish because the integrand is an odd function of either v_r or \mathbf{v}_t . Similar considerations show that the velocity-dispersion tensor is diagonal in the (v_r, v_θ, v_ϕ) frame, with diagonal components

$$\begin{aligned}\sigma_r^2 \equiv \overline{v_r^2} &= \frac{1}{\nu} \int dv_r v_r^2 \int d^2\mathbf{v}_t f[\tfrac{1}{2}(v_r^2 + v_\theta^2 + v_\phi^2) + \Phi, rv_t], \\ &= \frac{2\pi}{\nu} \int_{-\infty}^{\infty} dv_r v_r^2 \int_0^{\infty} dv_t v_t f[\tfrac{1}{2}(v_r^2 + v_t^2) + \Phi, rv_t], \\ \sigma_\theta^2 \equiv \overline{v_\theta^2} &= \frac{1}{\nu} \int dv_\theta v_\theta^2 \int dv_\phi \int dv_r f[\tfrac{1}{2}(v_r^2 + v_\theta^2 + v_\phi^2) + \Phi, rv_t], \\ &= \frac{\pi}{\nu} \int_0^{\infty} dv_t v_t^3 \int_{-\infty}^{\infty} dv_r f[\tfrac{1}{2}(v_r^2 + v_t^2) + \Phi, rv_t], \\ \sigma_\phi^2 &= \sigma_\theta^2.\end{aligned}\tag{4.37}$$

In general $\sigma_\theta^2(r) \neq \sigma_r^2(r)$ because the dependence of f on v_t differs from its dependence on v_r .

We shall see in §4.8 that in a stellar system $\nu\sigma^2$ plays a role analogous to that of pressure in a fluid. The inequality of σ_r^2 and σ_θ^2 implies that in general different pressures act radially and tangentially in a spherical stellar system; in other words pressure is a tensor rather than a scalar.

(c) DF depending on H and L_z If the potential $\Phi(\mathbf{x})$ is axisymmetric, $L_z = Rv_\phi$ is an isolating integral that can be included in f alongside H .

⁶ Lynden-Bell (1960) discusses spherical systems that rotate.

Then, in cylindrical coordinates the mean velocity has components

$$\begin{aligned}\bar{v}_R &= \frac{1}{\nu} \int dv_R v_R \int dv_z \int dv_\phi f[\tfrac{1}{2}(v_R^2 + v_z^2 + v_\phi^2) + \Phi, Rv_\phi] = 0, \\ \bar{v}_z &= \frac{1}{\nu} \int dv_z v_z \int dv_R \int dv_\phi f[\tfrac{1}{2}(v_R^2 + v_z^2 + v_\phi^2) + \Phi, Rv_\phi] = 0, \\ \bar{v}_\phi &= \frac{1}{\nu} \int dv_\phi v_\phi \int dv_R \int dv_z f[\tfrac{1}{2}(v_R^2 + v_z^2 + v_\phi^2) + \Phi, Rv_\phi].\end{aligned}\quad (4.38)$$

The integrals for \bar{v}_R and \bar{v}_z vanish because the integrands are odd functions of v_R and v_z , respectively. The integral for \bar{v}_ϕ will also vanish if f is an even function of L_z . In general $f(H, L_z)$ can be decomposed into a part that is even in L_z and a part that is odd:

$$f(H, L_z) = f_+(H, L_z) + f_-(H, L_z) \quad (4.39a)$$

where

$$f_\pm(H, L_z) \equiv \tfrac{1}{2}[f(H, L_z) \pm f(H, -L_z)]. \quad (4.39b)$$

The even part of f will not contribute to $\nu\bar{v}_\phi$, while f_- will not contribute to the corresponding integral (4.20) for ν . The velocity-dispersion tensor is diagonal in the (v_R, v_z, v_ϕ) frame, with non-zero components

$$\begin{aligned}\overline{v_R^2} &= \sigma_R^2 = \frac{1}{\nu} \int dv_R v_R^2 \int dv_z \int dv_\phi f[\tfrac{1}{2}(v_R^2 + v_z^2 + v_\phi^2) + \Phi, Rv_\phi], \\ \sigma_z^2 &= \sigma_R^2, \\ \sigma_\phi^2 &= \frac{1}{\nu} \int dv_\phi (v_\phi - \bar{v}_\phi)^2 \int dv_R \int dv_z f[\tfrac{1}{2}(v_R^2 + v_z^2 + v_\phi^2) + \Phi, Rv_\phi].\end{aligned}\quad (4.40)$$

The pressure that acts in the azimuthal direction is in general different from that which acts in any direction within the meridional plane.

4.3 DFs for spherical systems

The simplest stellar systems are spherical. A study of spherical models not only provides a good introduction to the structure of more general systems, but is also of considerable practical interest because some elliptical galaxies and clusters of galaxies, and most globular clusters, are very nearly spherical. For simplicity we shall consider the case in which the system has only one stellar population, so all stars are identical and there is a single DF f . We shall also generally assume that the mass density that generates the system's gravitational potential is proportional to $\int d^3\mathbf{v} f$ —such systems are

called **self-consistent** because the density distribution determines the potential through Poisson's equation, and the potential must also determine the density consistently through the collisionless Boltzmann equation. Many of the models we study are readily generalized to the more realistic case of multiple stellar populations, each of which will have its own DF and its own contribution to the total mass density.

We shall find it convenient to define a new gravitational potential and a new energy. If Φ_0 is some constant, then let the **relative potential** Ψ and the **relative energy** \mathcal{E} of a star be defined by

$$\Psi \equiv -\Phi + \Phi_0 \quad \text{and} \quad \mathcal{E} \equiv -H + \Phi_0 = \Psi - \frac{1}{2}v^2. \quad (4.41)$$

In practice, we generally choose Φ_0 to be such that $f > 0$ for $\mathcal{E} > 0$ and $f = 0$ for $\mathcal{E} \leq 0$. If an isolated system extends to infinity, $\Phi_0 = 0$ and the relative energy is equal to the binding energy. The relative potential of an isolated system satisfies Poisson's equation in the form

$$\nabla^2 \Psi = -4\pi G\rho, \quad (4.42)$$

subject to the boundary condition $\Psi \rightarrow \Phi_0$ as $|\mathbf{x}| \rightarrow \infty$.

4.3.1 Ergodic DFs for systems

Suppose we observe a spherical stellar system that is confined by a known spherical potential $\Phi(r)$. Then it is possible to derive for the system a unique ergodic DF that depends on the phase-space coordinates only through the Hamiltonian $H(\mathbf{x}, \mathbf{v})$. In this section we shall express this DF as a function of the relative energy $f(\mathcal{E})$. To derive this DF we note that the probability density $\nu(r)$ can be written as the integral of f over all velocities. Since f depends on the magnitude v of \mathbf{v} and not its direction, we can immediately integrate over angular coordinates in velocity space to produce 4π . We then have

$$\nu(r) = 4\pi \int dv v^2 f(\Psi - \frac{1}{2}v^2) = 4\pi \int_0^\Psi d\mathcal{E} f(\mathcal{E}) \sqrt{2(\Psi - \mathcal{E})}, \quad (4.43)$$

where we have used equation (4.41) and assumed that the constant Φ_0 in the definition of \mathcal{E} has been chosen such that $f = 0$ for $\mathcal{E} \leq 0$. Since Ψ is a monotonic function of r in any spherical system (Problem 2.16), we can regard ν as a function of Ψ instead of r . Thus

$$\frac{1}{\sqrt{8\pi}} \nu(\Psi) = 2 \int_0^\Psi d\mathcal{E} f(\mathcal{E}) \sqrt{\Psi - \mathcal{E}}. \quad (4.44)$$

**Box 4.1: An isolated system
with an ergodic DF is spherical**

An obvious question is whether there is any *self-consistent* non-spherical stellar system with an ergodic DF—in mathematical terms, are there finite, non-spherical solutions of the equations $\rho(\mathbf{x}) = M \int d^3\mathbf{v} f[\frac{1}{2}v^2 + \Phi(\mathbf{x})]$ and $\nabla^2\Phi = 4\pi G\rho$, in which $\Phi \rightarrow 0$ as $|\mathbf{x}| \rightarrow \infty$?

The first of these equations shows that the density depends on position only through the potential Φ ; thus the surfaces of constant density and potential coincide, so within the system the potential is a function $\Phi(\rho)$ of the density. Now define

$$p(\rho) = - \int_0^\rho d\rho' \rho' \frac{d\Phi}{d\rho}(\rho').$$

Using Problem 4.4, it is easy to see that $d\Phi/d\rho < 0$, so $p(\rho) > 0$. Taking the gradient of this equation yields $\nabla p = -\rho\nabla\Phi$, which is the equation of hydrostatic equilibrium (F.12) for a barotropic fluid with equation of state $p(\rho)$ (cf. eq. F.27).

We may now employ **Lichtenstein's theorem** on the symmetries of self-gravitating fluids, which states (Lindblom 1992): *Consider an isolated, self-gravitating, barotropic fluid of finite extent that is in a steady state, so the velocity \mathbf{v} and density ρ at every point are independent of time. If there is an axis $\hat{\mathbf{e}}_z$ such that $\mathbf{v} \cdot \hat{\mathbf{e}}_z = 0$, then the density distribution has a symmetry plane perpendicular to $\hat{\mathbf{e}}_z$.*

For a static fluid ($\mathbf{v} = 0$), there must be a symmetry plane perpendicular to every axis, so the fluid must be spherically symmetric. Thus *all isolated, finite, static, self-gravitating, barotropic fluids must be spherical*. Since we have shown that a stellar system with an ergodic DF satisfies the same equations—Poisson's equation, the equation of hydrostatic equilibrium, and the equation of state $p(\rho)$ —we conclude that *any isolated, finite, stellar system with an ergodic DF must be spherical*.

Differentiating both sides with respect to Ψ , we obtain

$$\frac{1}{\sqrt{8\pi}} \frac{d\nu}{d\Psi} = \int_0^\Psi d\mathcal{E} \frac{f(\mathcal{E})}{\sqrt{\Psi - \mathcal{E}}}. \quad (4.45)$$

Equation (4.45) is an Abel integral equation having solution (B.74)

$$f(\mathcal{E}) = \frac{1}{\sqrt{8\pi^2}} \frac{d}{d\mathcal{E}} \int_0^\mathcal{E} \frac{d\Psi}{\sqrt{\mathcal{E} - \Psi}} \frac{d\nu}{d\Psi}. \quad (4.46a)$$

An equivalent formula is

$$f(\mathcal{E}) = \frac{1}{\sqrt{8\pi^2}} \left[\int_0^\mathcal{E} \frac{d\Psi}{\sqrt{\mathcal{E} - \Psi}} \frac{d^2\nu}{d\Psi^2} + \frac{1}{\sqrt{\mathcal{E}}} \left(\frac{d\nu}{d\Psi} \right)_{\Psi=0} \right]. \quad (4.46b)$$

This result is due to Eddington (1916b), and we shall call it **Eddington's formula**. It implies that, given a spherical density distribution, we can recover an ergodic DF that generates a model with the given density. However, we have no guarantee that the solution $f(\mathcal{E})$ to equations (4.46) will satisfy the physical requirement that it be nowhere negative. Indeed, we may conclude from equation (4.46a) that *a spherical density distribution $\nu(r)$ in the potential $\Phi(r)$ can arise from an ergodic DF if and only if*

$$\int_0^{\mathcal{E}} \frac{d\Psi}{\sqrt{\mathcal{E} - \Psi}} \frac{d\nu}{d\Psi}$$

is an increasing function of \mathcal{E} . Note that this result holds regardless of whether the potential is self-consistently generated by the DF.

(a) Ergodic Hernquist, Jaffe and isochrone models In §2.2.2g we introduced the Jaffe and Hernquist models, which are members of the family of two-power density models. We now use Eddington's formula to recover the ergodic DFs of these models, as well as the DF of the isochrone model that was introduced in §2.2.2d.

From the second of equations (2.66) we have that the mass of a Hernquist model is related to the scale density ρ_0 and radius a by $M = 2\pi\rho_0 a^3$. The density of the Hernquist model is non-zero at all finite radii, so we choose $\Phi_0 = \Phi(r \rightarrow \infty) = 0$. Then from equation (2.67) radius is related to potential by

$$\frac{r}{a} = \frac{1}{\tilde{\Psi}} - 1 \quad \text{where} \quad \tilde{\Psi} \equiv \frac{\Psi a}{GM} = -\frac{\Phi a}{GM}. \quad (4.47)$$

Using this result to eliminate r/a from equation (2.64) with $\alpha = 1$ and $\beta = 4$ we obtain

$$\nu(\Psi) = \frac{\rho}{M} = \frac{1}{2\pi a^3} \frac{\tilde{\Psi}^4}{1 - \tilde{\Psi}}. \quad (4.48)$$

Differentiating with respect to Ψ we have

$$\frac{d\nu}{d\Psi} = \frac{1}{2\pi a^2 GM} \frac{\tilde{\Psi}^3(4 - 3\tilde{\Psi})}{(1 - \tilde{\Psi})^2}. \quad (4.49)$$

Equation (4.46b) gives the DF to be

$$f_{\text{H}}(\mathcal{E}) = \frac{\sqrt{2}}{(2\pi)^3 (GM)^2 a} \int_0^{\mathcal{E}} \frac{d\Psi}{\sqrt{\mathcal{E} - \Psi}} \frac{2\tilde{\Psi}^2(6 - 8\tilde{\Psi} + 3\tilde{\Psi}^2)}{(1 - \tilde{\Psi})^3}. \quad (4.50)$$

The evaluation of this integral is straightforward but tedious. The final result is (Hernquist 1990)

$$f_{\text{H}}(\mathcal{E}) = \frac{1}{\sqrt{2}(2\pi)^3 (GMa)^{3/2}} \frac{\sqrt{\tilde{\mathcal{E}}}}{(1 - \tilde{\mathcal{E}})^2} \times \left[(1 - 2\tilde{\mathcal{E}})(8\tilde{\mathcal{E}}^2 - 8\tilde{\mathcal{E}} - 3) + \frac{3 \sin^{-1} \sqrt{\tilde{\mathcal{E}}}}{\sqrt{\tilde{\mathcal{E}}(1 - \tilde{\mathcal{E}})}} \right], \quad (4.51)$$

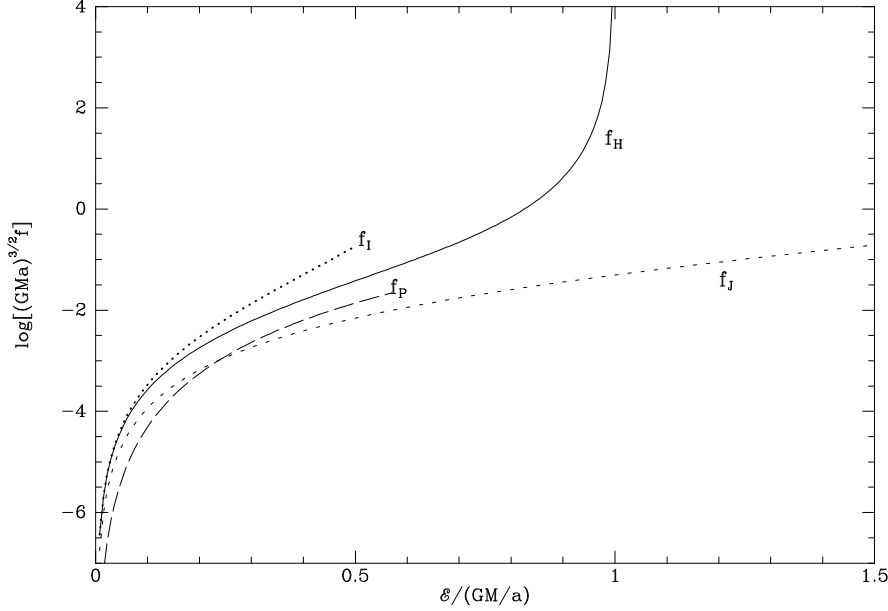


Figure 4.2 The ergodic DFs that generate stellar systems with the Hernquist (full curve) and Jaffe (dashed curve) density profiles. M is the model's mass and a is its scale length. The dotted curve shows the ergodic DF for the isochrone model, with a now denoting the scale length labeled b in equation (2.47). The long-dashed curve shows the DF of the Plummer model (eq. 4.83) with a the structure radius labeled b in equation (4.91.)

where $\tilde{\mathcal{E}} \equiv -Ea/GM$. The full curve in Figure 4.2 shows this DF.

An analogous calculation yields the DF of the Jaffe model. We now have

$$\frac{r}{a} = \frac{1}{e^{\tilde{\Psi}} - 1} \quad \text{so} \quad \nu = \frac{1}{4\pi a^3} e^{-2\tilde{\Psi}} (e^{\tilde{\Psi}} - 1)^4. \quad (4.52)$$

Differentiating this expression and then evaluating the integral of equation (4.46b) we find that the ergodic DF of the Jaffe model is (Jaffe 1983)

$$f_J(\mathcal{E}) = \frac{1}{2\pi^3(GMa)^{3/2}} \left[F_- \left(\sqrt{2\tilde{\mathcal{E}}} \right) - \sqrt{2} F_- \left(\sqrt{\tilde{\mathcal{E}}} \right) - \sqrt{2} F_+ \left(\sqrt{\tilde{\mathcal{E}}} \right) + F_+ \left(\sqrt{2\tilde{\mathcal{E}}} \right) \right], \quad (4.53)$$

where $F_{\pm}(z)$ is Dawson's integral (Appendix C.3).

The Jaffe DF is shown by the dashed curve in Figure 4.2. At the smallest values of \mathcal{E} , $f_J \rightarrow \frac{1}{2} f_H$ because at large radii each model potential tends to $-GM/r$ and for given model mass M , the density in the Jaffe model is half that in the Hernquist model. The DFs differ profoundly at large values

of \mathcal{E} because the Jaffe model has a divergent central potential, while the central potential of the Hernquist model is $-GM/a$. On account of the deeper potential well of the Jaffe model, stars are distributed through a larger volume of phase space than they are in the Hernquist model, so the DF is smaller at all energies.

From Eddington's formula one can determine the ergodic DF of the isochrone model (§2.2.2d and Hénon 1960b)

$$f_{\text{I}}(\mathcal{E}) = \frac{1}{\sqrt{2}(2\pi)^3(GMb)^{3/2}} \frac{\sqrt{\tilde{\mathcal{E}}}}{[2(1-\tilde{\mathcal{E}})]^4} \left[27 - 66\tilde{\mathcal{E}} + 320\tilde{\mathcal{E}}^2 - 240\tilde{\mathcal{E}}^3 \right. \\ \left. + 64\tilde{\mathcal{E}}^4 + 3(16\tilde{\mathcal{E}}^2 + 28\tilde{\mathcal{E}} - 9) \frac{\sin^{-1} \sqrt{\tilde{\mathcal{E}}}}{\sqrt{\tilde{\mathcal{E}}(1-\tilde{\mathcal{E}})}} \right]. \quad (4.54)$$

The dotted curve in Figure 4.2 shows this DF. At small binding energies $f_{\text{I}} \rightarrow f_{\text{H}}$ because both models have $\nu \propto r^{-4}$ at large r (cf. eqs. 2.51 and 2.64). At the largest binding energy attainable in the isochrone model, $\frac{1}{2}GM/b$, the phase-space density is finite, while to generate the model's central cusp the Hernquist DF has to diverge as \mathcal{E} tends to the central potential.

(b) Differential energy distribution In statistical physics an important role is played by the **density of states** g . In the limit of classical physics $g(E)$ is the volume of phase space per unit energy. We now evaluate g under the assumption that the system is spherical, so the potential is a function $\Phi(r)$ of r only. We have (eq. C.7)

$$g(E) = \int d^3\mathbf{x} d^3\mathbf{v} \delta(H - E) \\ = (4\pi)^2 \int_0^{r_{\text{m}}(E)} dr r^2 \int dv v^2 \delta(\frac{1}{2}v^2 + \Phi - E), \quad (4.55)$$

where $r_{\text{m}}(E)$ is the radius at which $\Phi = E$. To evaluate the inner integral, we change the integration variable to $\xi \equiv \frac{1}{2}v^2$. Then

$$g(E) = (4\pi)^2 \int_0^{r_{\text{m}}(E)} dr r^2 \int d\xi \sqrt{2\xi} \delta(\xi + \Phi - E) \\ = (4\pi)^2 \int_0^{r_{\text{m}}(E)} dr r^2 \sqrt{2(E - \Phi)}. \quad (4.56)$$

Applying this formula to the Hernquist model, we find that

$$g_{\text{H}}(E) = (4\pi)^2 a^3 \sqrt{2|E|} \int_1^A dX (X-1)^2 \left(\frac{A}{X} - 1 \right)^{1/2} \\ = (4\pi)^2 a^3 \sqrt{2|E|} \left[\sqrt{A-1} \left(\frac{1}{8}A^2 - \frac{5}{12}A - \frac{1}{3} \right) \right. \\ \left. + \frac{1}{8}A(A^2 - 4A + 8) \cos^{-1}(A^{-1/2}) \right], \quad (4.57a)$$

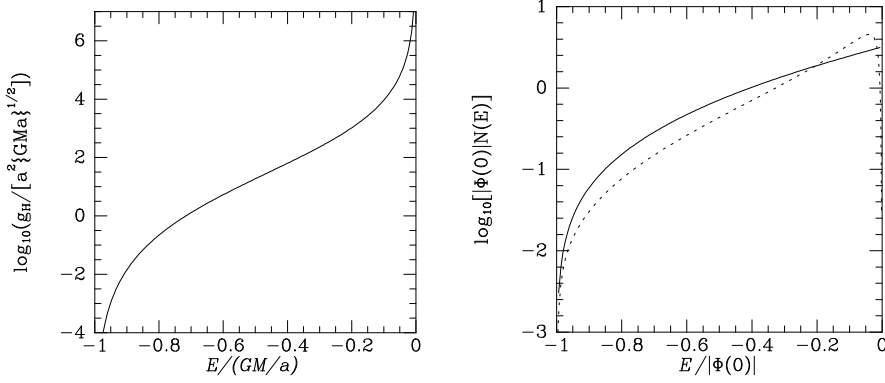


Figure 4.3 Left panel: the density of states $g_H(E)$ of the Hernquist model. Right panel: the differential energy distributions of the isotropic Hernquist model (full curve) and the $R^{1/4}$ model (dashed curve).

where

$$A \equiv \frac{GM}{a|E|}. \quad (4.57b)$$

The left panel in Figure 4.3 shows that $g_H(E)$ increases rapidly with E .

The **differential energy distribution**

$$N(E) \equiv g(E)f(E) \quad (4.58)$$

is such that the fraction of the system's stars that have energies in the range $(E + dE, E)$ is $N(E) dE$. For most realistic stellar systems the increase with E in g overwhelms the decrease of f , with the result that $N(E)$ is an increasing function of E . The right panel of Figure 4.3 illustrates this by showing the differential energy distributions of the Hernquist and $R^{1/4}$ models (eq. 1.17 with $m = 4$). In the case of the $R^{1/4}$ model⁷ $N(E) \propto e^{\beta E}$ is approximately valid for $E/|\Phi(0)| \gtrsim -0.8$. The rise in $N(E)$ with E reflects the fact that there are a large number of stars in the low-density envelope of a galaxy, and these stars are crowded into a relatively small range in binding energy.

4.3.2 DFs for anisotropic spherical systems

From Eddington's formula we can find an ergodic DF $f(H)$ that generates any given spherical density distribution $\nu(r)$ in a given potential $\Phi(r)$. However, there is no guarantee that this DF will satisfy the physical requirement $f \geq 0$. We now show that it *is* always possible find a non-negative DF if we consider DFs of the form $f(H, L)$, by building the system using only circular

⁷ In fact $N(E) \propto e^{\beta E}$ for all Sérsic models (Ciotti 1991).

orbits. By combining circular orbits of a given relative energy \mathcal{E}' with their angular-momentum vectors uniformly distributed over a sphere, we generate a spherical shell with the radius of circular orbits of energy \mathcal{E}' . Any density profile $\nu(r)$ can then be formed by adding such shells with a suitable radial weighting. We can express this idea mathematically by noting that the DF of the spherical shell is proportional to the product of two delta functions, $f_s(\mathcal{E}, L) = \delta(\mathcal{E} - \mathcal{E}')\delta[L - L_c(\mathcal{E}')]$, where $L_c(\mathcal{E}')$ is the angular momentum of a circular orbit of relative energy \mathcal{E}' . For a suitably chosen non-negative function $F(\mathcal{E}')$ the integral

$$f_c(\mathcal{E}, L) \equiv \int_0^{\mathcal{E}_{\max}} d\mathcal{E}' f_s(\mathcal{E}', L) F(\mathcal{E}') = F(\mathcal{E})\delta[L - L_c(\mathcal{E})] \quad (4.59)$$

is a DF that generates the required density distribution $\nu(r)$.

The circular-orbit DF f_c is associated with vanishing radial dispersion σ_r . If a non-negative ergodic DF $f_i(\mathcal{E})$ also exists, then these two DFs will be joined by a continuum of DFs of the form

$$f_\alpha \equiv \alpha f_i + (1 - \alpha) f_c \quad (0 \leq \alpha \leq 1). \quad (4.60)$$

As α increases the orbits become steadily more eccentric and σ_r increases to equality with σ_θ . This sequence may even continue to DFs with $\alpha > 1$ for which $\sigma_r > \sigma_\theta$, but such a continuation is not guaranteed: the more heavily we weight highly eccentric orbits, the more strongly $\nu(r)$ is constrained by the requirement that f_α be non-negative. Conversely, the circular-orbit DF will still be non-negative even when $\nu(r)$ is such that $f_i(\mathcal{E})$ is somewhere negative, and it is likely that we can construct some non-negative DFs that have $\sigma_r \neq 0$, even though an ergodic DF is not allowed.

We define the **anisotropy parameter** to be

$$\beta \equiv 1 - \frac{\sigma_\theta^2 + \sigma_\phi^2}{2\sigma_r^2} = 1 - \frac{\overline{v_\theta^2} + \overline{v_\phi^2}}{2\overline{v_r^2}}. \quad (4.61)$$

This parameter quantifies the system's degree of radial anisotropy: if all orbits are circular, $\sigma_r = 0$ and $\beta = -\infty$; if the DF is ergodic, $\beta = 0$; if all orbits are perfectly radial, $\sigma_\theta = \sigma_\phi = 0$ and $\beta = 1$. DFs with $\beta > 0$ are said to be **radially biased**, while those with $\beta < 0$ are **tangentially biased**. The value of β is determined by the way in which f depends on total angular momentum L .

To explore the effect of anisotropy on the structure of a stellar system, we shall find it useful to have available explicit expressions for DFs that all generate the same radial profile, but have different degrees of radial anisotropy.

(a) Models with constant anisotropy Models in which the anisotropy parameter takes some fixed non-zero value β at all radii can be generated by taking the DF to be of the form (Problem 4.6)

$$f(\mathcal{E}, L) = L^{-2\beta} f_1(\mathcal{E}), \quad (4.62)$$

where f_1 is an arbitrary non-negative function. In terms of the polar coordinates for velocity space,

$$v_r = v \cos \eta \quad ; \quad v_\theta = v \sin \eta \cos \psi \quad ; \quad v_\phi = v \sin \eta \sin \psi, \quad (4.63)$$

we may write

$$\nu(r) = \int d^3\mathbf{v} f(\mathcal{E}, L) = 2\pi \int_0^\pi d\eta \sin \eta \int_0^{\sqrt{2\Psi}} dv v^2 f\left(\Psi - \frac{1}{2}v^2, rv \sin \eta\right). \quad (4.64)$$

For a DF of the form (4.62) this expression may be rewritten

$$\nu(r) = \frac{2\pi I_\beta}{r^{2\beta}} \int_0^\infty dv v^{2-2\beta} f_1\left[\Psi(r) - \frac{1}{2}v^2\right], \quad (4.65a)$$

where

$$I_\beta \equiv \int_0^\pi d\eta \sin^{1-2\beta} \eta = \sqrt{\pi} \frac{(-\beta)!}{\left(\frac{1}{2} - \beta\right)!} \quad (\beta < 1). \quad (4.65b)$$

The integral in (4.65a) is similar to that occurring in equation (4.43), so after multiplying through by $r^{2\beta}$ we obtain an equation that is closely analogous to equation (4.44):

$$\frac{2^{\beta-1/2}}{2\pi I_\beta} r^{2\beta} \nu = \int_0^\Psi d\mathcal{E} \frac{f_1(\mathcal{E})}{(\Psi - \mathcal{E})^{\beta-1/2}} \quad (\beta < 1). \quad (4.66)$$

When the left side is considered to be a function of Ψ rather than r , this becomes an Abel integral equation (B.74a) so long as $\frac{1}{2} < \beta < \frac{3}{2}$. Moreover, when $\beta \leq \frac{1}{2}$, the equation can be reduced to an Abel equation by differentiating both sides one or more times with respect to Ψ . For example,

$$\frac{2^{\beta-1/2}}{2\pi I_\beta} \frac{d}{d\Psi} (r^{2\beta} \nu) = \left(\frac{1}{2} - \beta\right) \int_0^\Psi d\mathcal{E} \frac{f_1(\mathcal{E})}{(\Psi - \mathcal{E})^{\beta+1/2}}, \quad (4.67)$$

which is an Abel equation for $-\frac{1}{2} < \beta < \frac{1}{2}$. Hence we can analytically solve for $f_1(\mathcal{E})$ whenever we can express $r^{2\beta} \nu$ as a function of Ψ .

Equation (4.66) is exceptionally simple when $\beta = \frac{1}{2}$, corresponding to $\sigma_\theta^2 = \sigma_\phi^2 = \frac{1}{2}\sigma_r^2$. In this case the denominator of the integrand becomes a constant. Differentiating both sides with respect to Ψ we find

$$f_1(\Psi) = \frac{1}{2\pi^2} \frac{d}{d\Psi} (r\nu) \quad (\beta = \frac{1}{2}). \quad (4.68)$$

In the case of the Hernquist model, from equations (4.47) and (4.48) we find that

$$f_1(\mathcal{E}) = \frac{3\tilde{\mathcal{E}}^2}{4\pi^3 GM a}, \quad \tilde{\mathcal{E}} = \frac{\mathcal{E} a}{GM}. \quad (4.69)$$

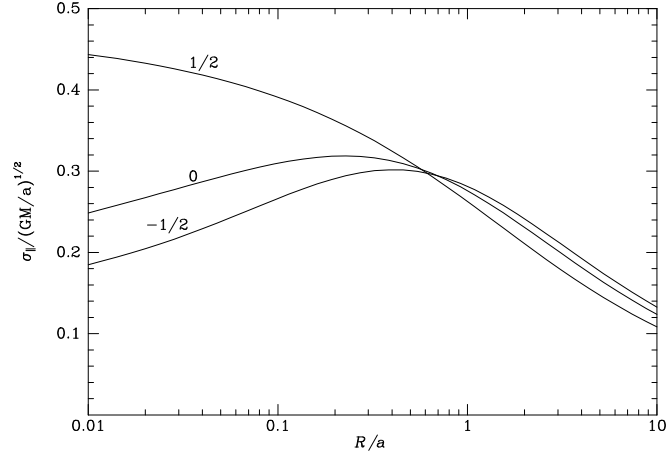


Figure 4.4 Line-of-sight velocity dispersion as a function of projected radius, from spatially identical systems that have different DFs. In each system the density and potential are those of the Hernquist model and the anisotropy parameter β of equation (4.61) is independent of radius. The curves are labeled by the relevant value of β . In the isotropic system, the velocity dispersion falls as one approaches the center (cf. Problem 4.14).

A contrasting case of almost equal simplicity is $\beta = -\frac{1}{2}$, corresponding to $\sigma_{\theta}^2 = \sigma_{\phi}^2 = \frac{3}{2}\sigma_r^2$. Then equation (4.66) becomes

$$\frac{1}{2\pi^2} \frac{\nu}{r} = \int_0^{\Psi} d\mathcal{E} f_1(\mathcal{E})(\Psi - \mathcal{E}). \quad (4.70)$$

Differentiating through twice with respect to Ψ we have

$$f_1(\Psi) = \frac{1}{2\pi^2} \frac{d^2(\nu/r)}{d\Psi^2} \quad (\beta = -\frac{1}{2}). \quad (4.71)$$

In the case of the Hernquist model, this yields

$$f_1(\mathcal{E}) = \frac{1}{4\pi^3(GMa)^2} \frac{d^2}{d\tilde{\mathcal{E}}^2} \left(\frac{\tilde{\mathcal{E}}^5}{(1 - \tilde{\mathcal{E}})^2} \right), \quad (4.72)$$

which one may easily show is non-negative for $\tilde{\mathcal{E}} \leq 1$.

Figure 4.4 shows the line-of-sight velocity dispersion σ_{\parallel} of a Hernquist model as a function of projected radius when the DF is (i) ergodic (eq. 4.50) labeled “0”; (ii) radially biased (eqs. 4.62 and 4.69) labeled $\frac{1}{2}$, and (iii) tangentially biased (eqs. 4.62 and 4.72) labeled $-\frac{1}{2}$. In the radially biased system, the central value of σ_{\parallel} is nearly twice that in the isotropic system, and more than twice that in the tangentially biased system. Conversely, at

large radii the tangentially biased system has the largest value of σ_{\parallel} , and the radially biased system the smallest.

It is easy to understand why radial bias increases line-of-sight velocity dispersion at small radii while tangential bias increases it at large: along the line of sight through the center, only v_r contributes to σ_{\parallel} ; conversely, in the outer envelope, where the galaxy's density profile is steeply falling, the dominant contribution to the dispersion comes from where the line of sight makes its closest approach to the center of the galaxy. At this point of closest approach, only the tangential velocity contributes to σ_{\parallel} , so a bias towards tangential velocities increases σ_{\parallel} . The tendency of radial bias to make σ_{\parallel} a steeper function of radius is general, and constitutes one of the most vexing sources of uncertainty in attempts to measure the mass distributions of stellar systems (§4.9).

(b) Osipkov–Merritt models Models of galaxy formation generally imply that β increases with radius, corresponding to a nearly ergodic DF near the center and a radially biased DF in the outer envelope (§4.10.3). Simple models in which β increases with radius can be constructed as follows (Osipkov 1979; Merritt 1985). We assume that f depends on \mathcal{E} and L only through the variable

$$Q \equiv \mathcal{E} - \frac{L^2}{2r_a^2}, \quad (4.73)$$

where r_a is a constant, called the **anisotropy radius**. In terms of the polar coordinates defined by equations (4.63), the definition of Q becomes

$$Q = \Psi - \frac{1}{2}v^2 \left(1 + \frac{r^2}{r_a^2} \sin^2 \eta \right). \quad (4.74)$$

We substitute the DF $f(\mathcal{E}, L) = f(Q)$ into equation (4.64) and replace the integration variable v with Q . At constant r and η , $dQ = -[1 + (r/r_a)^2 \sin^2 \eta] v dv$, and thus

$$\nu(r) = 2\pi \int_0^\pi d\eta \sin \eta \int_0^\Psi dQ f(Q) \frac{\sqrt{2(\Psi - Q)}}{[1 + (r/r_a)^2 \sin^2 \eta]^{3/2}}, \quad (4.75)$$

where we have imposed the condition $f(Q) = 0$ for $Q \leq 0$. We interchange the order of integrations in equation (4.75), and the inner integral becomes

$$\int_0^\pi d\eta \frac{\sin \eta}{[1 + (r/r_a)^2 \sin^2 \eta]^{3/2}} = \frac{2}{1 + (r/r_a)^2}. \quad (4.76)$$

Hence

$$\left(1 + \frac{r^2}{r_a^2} \right) \nu(r) = 4\pi \int_0^\Psi dQ f(Q) \sqrt{2(\Psi - Q)}. \quad (4.77)$$

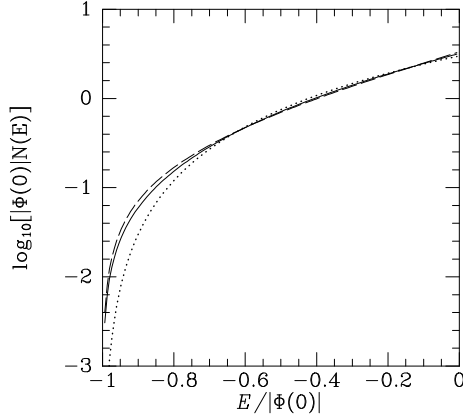


Figure 4.5 The differential energy distributions of the Hernquist models with constant anisotropy $\beta = -\frac{1}{2}$ (dashed curve), $\beta = 0$ (full curve) and $\beta = \frac{1}{2}$ (dotted curve).

But the right side of this equation is identical with the right side of equation (4.43), with \mathcal{E} replaced by Q . Hence by equation (4.46b) we have

$$f(Q) = \frac{1}{\sqrt{8\pi^2}} \left[\int_0^Q \frac{d\Psi}{\sqrt{Q-\Psi}} \frac{d^2\nu_Q}{d\Psi^2} + \frac{1}{\sqrt{Q}} \left(\frac{d\nu_Q}{d\Psi} \right)_{\Psi=0} \right], \quad (4.78a)$$

where

$$\nu_Q(r) \equiv \left(1 + \frac{r^2}{r_a^2} \right) \nu(r). \quad (4.78b)$$

A model for which the DF has the form $f(Q)$, where Q is given by equation (4.73), is called an **Osipkov–Merritt model**.

The anisotropy parameter of any model with a DF of the form $f(Q)$ is (Problem 4.13)

$$\beta(r) = \frac{1}{1 + r_a^2/r^2}. \quad (4.79)$$

This function rises from zero at $r \ll r_a$ (ergodic DF) to unity at $r \gg r_a$ (radial DF), and already exceeds 0.9 for $r > 3r_a$.

(c) Other anisotropic models The degree of anisotropy in the models described above is controlled by a single parameter, either β or r_a . We can obtain more flexible models by considering a DF of the form

$$f(\mathcal{E}, L) = G(\mathcal{E})h(x), \quad \text{where} \quad x \equiv \frac{L}{L_0 + L_c(\mathcal{E})} \quad (4.80)$$

with L_0 a free parameter and $L_c(\mathcal{E})$ is the angular momentum of the circular orbit with energy \mathcal{E} . The variable x increases from zero for radial orbits to $L_c/(L_0 + L_c)$ for circular orbits, and the model will be tangentially biased if h is an increasing function of x , and radially biased if h is a decreasing function of x . For any chosen function $h(x)$, one can numerically determine

$G(\mathcal{E})$ such that the model has a given radial density profile (Gerhard 1991). Hence this approach makes it possible to construct models that have the same density profile but a wide variety of functional forms of $\beta(r)$.

(d) Differential-energy distribution for anisotropic systems It proves useful to calculate the differential energy distribution of a model with DF $f(H, L)$. With $v_t = L/r$ the tangential speed, we have (cf. eq. 4.55)

$$\begin{aligned}
 N(E) &= \int d^3\mathbf{x} d^3\mathbf{v} \delta(E - H) f(H, L) \\
 &= 8\pi^2 \int dr r^2 \int dv_r dv_t v_t \delta(E - H) f(H, L) \\
 &= 8\pi^2 \int dr \int dv_r dL L \delta(E - H) f(H, L) \\
 &= (4\pi)^2 \int dr \int dL \frac{L f(E, L)}{\sqrt{2(E - \Phi) - L^2/r^2}},
 \end{aligned} \tag{4.81}$$

where the last equality uses $v_r^2 = 2(E - \Phi) - L^2/r^2$ to eliminate v_r in favor of E ; a factor of two arises because v_r can be positive or negative. The integral is taken over the region in (r, L) space in which the argument of the square root is non-negative. Figure 4.5 shows $N(E)$ for the Hernquist models with constant anisotropy parameters $\beta = -\frac{1}{2}$ (dashed curve), $\beta = 0$ (full curve) and $\beta = \frac{1}{2}$ (dotted curve). Even though the DFs of these models are very different, they yield extremely similar differential energy distributions. The only significant difference is at large binding energy, and in the sense that the radially biased model ($\beta = \frac{1}{2}$) has fewer very tightly bound stars than does the tangentially biased model. This result arises because stars that are on eccentric orbits contribute to the density at radii that are much smaller than the radius of the circular orbit that has the same energy. This difference may significantly affect the vulnerability of the central regions to disruption by external perturbers. Thus we expect that galaxies with radially biased DFs are more fragile than galaxies with ergodic DFs.

4.3.3 Spherical systems defined by the DF

In the previous subsection we found DFs that generated a given density profile. Not surprisingly, the resulting system had a simple functional form for the density in real space, and usually a much more complex expression for the density in phase space (cf. eq. 4.54). In this section we proceed in the reverse order: we choose the functional form of the density in phase space and then investigate what the system looks like in real space.

The Jeans theorem and the system's spherical symmetry allow us to assume that f is a function of the relative energy $\mathcal{E} = \Psi - \frac{1}{2}v^2$ and the

magnitude of the angular momentum, and to write Poisson's equation in the form

$$\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Psi}{dr} \right) = -4\pi G \rho = -4\pi G M \int d^3\mathbf{v} f \left(\Psi - \frac{1}{2}v^2, |\mathbf{r} \times \mathbf{v}| \right), \quad (4.82)$$

where M is the system's total mass. For any given function $f(\mathcal{E}, L)$, this is an integro-differential equation for $\Psi(r)$. Since we have to solve for the potential after we have chosen f , we cannot normalize f to have unit integral over all phase space as we have done hitherto, until after we have solved equation (4.82) for $\Psi(r)$. For this reason it is convenient to redefine our normalization so that the integral of f over phase space is the total mass; thus, in this section we shall assume that the mass density ρ is given by $\rho = \int d^3\mathbf{v} f$.

(a) Polytropes and the Plummer model A simple form for the DF is

$$f(\mathcal{E}) = \begin{cases} F \mathcal{E}^{n-3/2} & (\mathcal{E} > 0) \\ 0 & (\mathcal{E} \leq 0). \end{cases} \quad (4.83)$$

With this form of f we have for the density ρ at radii where $\Psi > 0$

$$\rho = 4\pi \int_0^\infty dv v^2 f(\Psi - \frac{1}{2}v^2) = 4\pi F \int_0^{\sqrt{2\Psi}} dv v^2 (\Psi - \frac{1}{2}v^2)^{n-3/2}. \quad (4.84)$$

If we make the substitution $v^2 = 2\Psi \cos^2 \theta$, this becomes

$$\rho = c_n \Psi^n \quad (\Psi > 0), \quad (4.85a)$$

where

$$c_n \equiv 2^{7/2} \pi F \left[\int_0^{\pi/2} d\theta \sin^{2n-2} \theta - \int_0^{\pi/2} d\theta \sin^{2n} \theta \right] = \frac{(2\pi)^{3/2} (n - \frac{3}{2})! F}{n!}. \quad (4.85b)$$

If c_n is to be finite, we must have $n > \frac{1}{2}$.

In these models the density rises as the n th power of the relative potential when $\Psi > 0$ and is, of course, zero when $\Psi \leq 0$. No finite ergodic stellar system is homogeneous, for this would correspond to $\rho \propto \Psi^0$ or $n = 0$, which would violate the constraint $n > \frac{1}{2}$ —in other words there is no stellar-dynamical analog of a self-gravitating sphere of incompressible liquid.

When we use equation (4.85a) to eliminate ρ from Poisson's equation, we find

$$\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Psi}{dr} \right) + 4\pi G c_n \Psi^n = 0. \quad (4.86)$$

Polytropic gases have an equation of state $p = K \rho^\gamma$, where K is a constant and γ is the ratio of principal specific heats (cf. eq. F.46). Thus the equation of hydrostatic equilibrium for a self-gravitating sphere of polytropic gas (eq. F.12),

$$\frac{dp}{dr} = -\rho \frac{d\Phi}{dr}, \quad (4.87a)$$

becomes

$$K\gamma\rho^{\gamma-2}\frac{d\rho}{dr} = \frac{d\Psi}{dr}. \quad (4.87b)$$

If we set the constant involved in the definition of Ψ such that $\Psi = 0$ on the edge of the system, equation (4.87b) yields on integration

$$\rho^{\gamma-1} = \frac{\gamma-1}{K\gamma}\Psi. \quad (4.88)$$

Equation (4.88) is the same as equation (4.85a) with

$$c_n = \left(\frac{\gamma-1}{K\gamma}\right)^{1/(\gamma-1)} \quad \text{and} \quad \gamma = 1 + \frac{1}{n}. \quad (4.89)$$

Hence *the density distribution of an ergodic stellar system with DF (4.83) is the same as that of a polytropic gas sphere with $\gamma = 1 + 1/n$* . For this reason stellar systems with DFs given by (4.83) are known as polytropes. Note that polytropic gas spheres with $\gamma > 3$ require $n < \frac{1}{2}$, so these spheres have no stellar-dynamical analogs. A full account of gaseous polytropes can be found in Chandrasekhar (1939) and Horedt (2004).

The simplest solutions of equation (4.86) are obtained by assuming that the density varies as a power of radius, $\rho \propto r^{-\alpha}$. Since $\rho \propto \Psi^n$ in a polytrope, we have $\Psi \propto r^{-\alpha/n}$. Inserting this ansatz into both sides of equation (4.86), and requiring the same power of r to occur on each side, we find that $\alpha = 2n/(n-1)$. Since the potential cannot decrease with radius faster than in the Keplerian case $\Psi \propto r^{-1}$, we have $\alpha/n \leq 1$, so solutions of this type are feasible only for $n \geq 3$. The mass contained within radius r is

$$M(r) = -\frac{r^2}{G}\frac{d\Psi}{dr} \propto r^{1-\alpha/n} = r^{(n-3)/(n-1)}, \quad (4.90)$$

which is independent of radius in the case $n = 3$ but otherwise tends to zero as $r \rightarrow 0$. Hence in the case $n = 3$ our power-law solution represents a massless halo orbiting in the potential of a point mass, but for $n > 3$ our solution is a self-gravitating system. In the limit $n \rightarrow \infty$ the potential becomes proportional to $\ln r$; we shall encounter this model below as the “singular isothermal sphere.” Problem 4.14 gives the velocity dispersion in power-law models.

To obtain models in which the central potential and density are finite, we eliminate r and Ψ from equation (4.86) in favor of the rescaled radial variables,

$$s \equiv \frac{r}{b} \quad \text{and} \quad \psi \equiv \frac{\Psi}{\Psi_0}, \quad \text{where} \quad b \equiv \left(\frac{4}{3}\pi G\Psi_0^{n-1}c_n\right)^{-1/2} \quad (4.91a)$$

is the scale radius and $\Psi_0 = \Psi(0)$. Now equation (4.86) takes the simple form

$$\frac{1}{s^2} \frac{d}{ds} \left(s^2 \frac{d\psi}{ds} \right) = \begin{cases} -3\psi^n & \psi > 0; \\ 0 & \psi \leq 0. \end{cases} \quad (4.91b)$$

This equation is known as the **Lane–Emden equation** after H. Lane and R. Emden who studied it in connection with polytropic gas spheres. The natural boundary conditions to impose on it are: (i) $\psi(0) = 1$ by definition; (ii) $d\psi/ds|_0 = 0$, since in the absence of a central singularity in the density, the gravitational force must vanish at the center.

For general n , (4.91b) cannot be solved in terms of elementary functions. However, there are two special cases for which simple analytical solutions are available. (i) When $n = 1$, equation (4.91b) becomes the linear Helmholtz equation familiar from the theory of spherical waves (see Problems 4.15 and 4.16); and (ii) when $n = 5$, we obtain a model discovered by Schuster (1883) that is worth describing in some detail because it provides the simplest plausible model of a self-consistent stellar system.

Consider the function

$$\psi = \frac{1}{\sqrt{1+s^2}}. \quad (4.92)$$

Differentiating with respect to s we find that

$$\frac{1}{s^2} \frac{d}{ds} \left(s^2 \frac{d\psi}{ds} \right) = -\frac{1}{s^2} \frac{d}{ds} \left(\frac{s^3}{(1+s^2)^{3/2}} \right) = -\frac{3}{(1+s^2)^{5/2}} = -3\psi^5. \quad (4.93)$$

Therefore ψ is a solution of equation (4.91b) with $n = 5$. Since ψ also satisfies the central boundary conditions, it represents a physically acceptable potential. In fact, it is a dimensionless form of equation (2.44a), the potential of the Plummer model introduced in §2.2.2c. That is, *the $n = 5$ polytrope is a Plummer model*. The density of this model, given by equation (2.44b), is everywhere non-zero, declining as r^{-5} for $r \gg b$. The total mass is finite, however, with value

$$M = \frac{1}{G} \left(r^2 \frac{d\Phi}{dr} \right)_{r \rightarrow \infty} = -\frac{b}{G} \left(s^2 \frac{d\Psi}{ds} \right)_{s \rightarrow \infty} = \frac{b\Psi_0}{G}. \quad (4.94)$$

In general, the extent of the outer parts of a polytropic model increases with n ; for $n < 5$ the density goes to zero at a finite radius, for $n = 5$ the density is non-zero everywhere but the total mass is finite, and for $n > 5$ the density falls off so slowly at large r that the mass is infinite.

(b) The isothermal sphere We have just seen that to every polytropic gas sphere with $\gamma < 3$, there corresponds a stellar-dynamical polytrope with index $n = 1/(\gamma - 1) > \frac{1}{2}$. Thus stellar polytropes with large n correspond to gaseous polytropes for which $\gamma \simeq 1$. Hence in the limit $n \rightarrow \infty$, the corresponding gaseous system has $\gamma = 1$, which implies that $p = K\rho$. This

is the equation of state of an isothermal gas. The equation governing the structure of a self-gravitating isothermal sphere of ideal gas can be derived by taking a suitable limit of the Lane–Emden equation as $n \rightarrow \infty$ (Hunter 2001), but a more illuminating derivation starts with the equation of hydrostatic equilibrium, which reads

$$\frac{dp}{dr} = \frac{k_B T}{m} \frac{d\rho}{dr} = -\rho \frac{d\Phi}{dr} = -\rho \frac{GM(r)}{r^2}, \quad (4.95a)$$

where k_B is Boltzmann's constant, p and T are the pressure and temperature of the gas, m is the mass per particle, and $M(r)$ is the total mass interior to radius r (eqs. F.12 and F.31). Multiplying equation (4.95a) through by $(r^2 m / \rho k_B T)$ and then differentiating with respect to r , we obtain

$$\frac{d}{dr} \left(r^2 \frac{d \ln \rho}{dr} \right) = -\frac{4\pi G m}{k_B T} r^2 \rho, \quad (4.95b)$$

where we have used the relationship $dM/dr = 4\pi r^2 \rho$.

Now suppose we have a stellar-dynamical system whose DF is

$$f(\mathcal{E}) = \frac{\rho_1}{(2\pi\sigma^2)^{3/2}} e^{\mathcal{E}/\sigma^2} = \frac{\rho_1}{(2\pi\sigma^2)^{3/2}} \exp\left(\frac{\Psi - \frac{1}{2}v^2}{\sigma^2}\right). \quad (4.96)$$

Then, integrating over all velocities, we find

$$\rho = \rho_1 e^{\Psi/\sigma^2}. \quad (4.97)$$

Poisson's equation for this system reads

$$\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{d\Psi}{dr} \right) = -4\pi G \rho, \quad (4.98)$$

or, with equation (4.97)

$$\frac{d}{dr} \left(r^2 \frac{d \ln \rho}{dr} \right) = -\frac{4\pi G}{\sigma^2} r^2 \rho. \quad (4.99a)$$

For future reference, note that if we eliminate ρ rather than Ψ between equations (4.97) and (4.98), we obtain

$$\frac{d}{dr} \left(r^2 \frac{d\Psi}{dr} \right) + 4\pi G \rho_1 r^2 e^{\Psi/\sigma^2} = 0, \quad (4.99b)$$

which is the analog of the governing equation of polytropes, (4.86).

Equations (4.95b) and (4.99a) are identical if we set

$$\sigma^2 = \frac{k_{\text{B}}T}{m}. \quad (4.100)$$

Therefore *the structure of an isothermal self-gravitating sphere of gas is identical with the structure of a collisionless system of stars whose DF is given by equation (4.96).*

A little thought shows why there is this correspondence between the gaseous and stellar-dynamical isothermal spheres. The distribution of velocities at each point in the stellar isothermal sphere is the **Maxwellian** or **Maxwell–Boltzmann distribution**

$$dn \propto \exp\left(-\frac{|\mathbf{v}|^2}{2\sigma^2}\right) d^3\mathbf{v}. \quad (4.101)$$

However, kinetic theory (e.g., Pathria 1972) tells us that this is also the equilibrium Maxwell–Boltzmann distribution which would obtain if the stars were allowed to bounce elastically off each other like the molecules of a gas. Therefore, if the DF of a system is given by equation (4.96), it is a matter of indifference whether the particles of the system collide with one another or not.

Notice that the correspondence between a gaseous polytrope with $\gamma > 1$ and the corresponding stellar-dynamical model is not as close as that between the two isothermal systems; the gas molecules always have a Maxwellian distribution (with temperature depending on radius), while the stellar velocity distribution is given by (4.83). Thus a stellar polytrope would be drastically altered if elastic collisions were allowed to occur between its stars.

The mean-square speed of the stars at a point in the isothermal sphere is

$$\overline{v^2} = \frac{\int_0^\infty dv v^4 \exp\left(\frac{\Psi - \frac{1}{2}v^2}{\sigma^2}\right)}{\int_0^\infty dv v^2 \exp\left(\frac{\Psi - \frac{1}{2}v^2}{\sigma^2}\right)} = 2\sigma^2 \frac{\int_0^\infty dx x^4 e^{-x^2}}{\int_0^\infty dx x^2 e^{-x^2}} = 3\sigma^2. \quad (4.102)$$

Thus $\overline{v^2}$ is independent of position. The dispersion in any one component of velocity, for example $(\overline{v_r^2})^{1/2}$, is equal to σ .

It is easy to find one solution of equation (4.99a). If we set $\rho = Cr^{-b}$, the left side of the equation is found to equal $-b$, while the right side equals $-(4\pi G/\sigma^2)Cr^{2-b}$. Therefore, we must set $b = 2$ and $C = \sigma^2/(2\pi G)$, which yields

$$\rho(r) = \frac{\sigma^2}{2\pi Gr^2}. \quad (4.103)$$

This solution describes the **singular isothermal sphere**. The mass interior to radius r , the circular speed and the gravitational potential are (eqs. 2.60, 2.61, and 2.62)

$$M(r) = \frac{2\sigma^2 r}{G} \quad ; \quad v_c(r) = \sqrt{2}\sigma \quad ; \quad \Phi(r) = 2\sigma^2 \ln(r) + \text{constant}, \quad (4.104)$$

and the surface density is (eq. 2.59)

$$\Sigma(R) = \frac{\sigma^2}{2GR}. \quad (4.105)$$

The singular isothermal sphere has infinite density at $r = 0$. To obtain a solution of equations (4.99) that is well behaved at the origin, it is convenient to define new dimensionless variables $\tilde{\rho}$ and \tilde{r} to replace ρ and r ; we define these in terms of the central density ρ_0 and the **King radius** r_0 by

$$\tilde{\rho} \equiv \frac{\rho}{\rho_0} \quad \text{and} \quad \tilde{r} \equiv \frac{r}{r_0}, \quad \text{where} \quad r_0 \equiv \sqrt{\frac{9\sigma^2}{4\pi G\rho_0}}. \quad (4.106)$$

We shall find that r_0 is the radius at which the projected density of the isothermal sphere falls to roughly half (in fact, 0.5013) of its central value, and because of this some authors call r_0 the core radius in analogy with the usual observational definition (page 30 and BM p. 366). In terms of our new variables, equations (4.99) become

$$\frac{d}{d\tilde{r}} \left(\tilde{r}^2 \frac{d \ln \tilde{\rho}}{d\tilde{r}} \right) = -9\tilde{r}^2 \tilde{\rho} \quad (4.107a)$$

or

$$\frac{d}{d\tilde{r}} \left[\tilde{r}^2 \frac{d(\Psi/\sigma^2)}{d\tilde{r}} \right] + 9\tilde{r}^2 \exp \left[\frac{\Psi(r) - \Psi(0)}{\sigma^2} \right] = 0. \quad (4.107b)$$

In Figure 4.6 we show the function $\tilde{\rho}(\tilde{r})$ obtained by numerically integrating equation (4.107a) from $\tilde{r} = 0$ outward, starting from the boundary conditions $\tilde{\rho}(0) = 1$ and $d\tilde{\rho}/d\tilde{r} = 0$. Notice that by about $\tilde{r} = 15$, $\tilde{\rho}(\tilde{r})$ is declining as a straight line in the log-log plot of Figure 4.6; in fact, the solution is approaching⁸ the singular isothermal sphere of equation (4.104), which in these variables has the form $\tilde{\rho} = \frac{2}{9}\tilde{r}^{-2}$. This is shown as a dotted line in the figure.

In Figure 4.6 we also plot the surface density $\Sigma(R)$ of the isothermal sphere in units of $\rho_0 r_0$. For $R \gg r_0$ the surface density asymptotes to that of the singular isothermal sphere (eq. 4.105).

⁸The asymptotic behavior of the isothermal sphere as $r \rightarrow \infty$ is described more accurately in Problem 7.6.

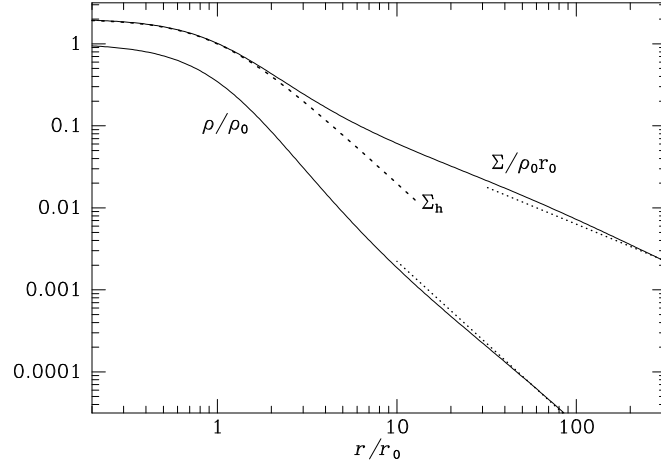


Figure 4.6 Volume (ρ/ρ_0) and projected ($\Sigma/\rho_0 r_0$) mass densities of the isothermal sphere. The dotted lines show the volume- and surface-density profiles of the singular isothermal sphere. The dashed curve shows the surface density of the modified Hubble model (4.109a).

If $M(r)$ is the mass interior to r , the circular speed at r is given by

$$v_c^2(r) = \frac{GM(r)}{r}. \quad (4.108a)$$

On integrating equation (4.99a), we find that

$$v_c^2 = -\sigma^2 \frac{d \ln \rho}{d \ln r}. \quad (4.108b)$$

In Figure 4.7 we plot $d \ln \rho / d \ln r$ for the isothermal sphere, which for $r \gg r_0$ tends to -2 . Thus the circular speed at large r is constant at $v_c = \sqrt{2}\sigma$, the value for the singular isothermal sphere (eq. 4.104).

At $\tilde{r} \lesssim 2$ ($r \lesssim 2r_0$) a useful approximation to $\tilde{\rho}(\tilde{r})$ is the modified Hubble model introduced in §2.2.2e,

$$\tilde{\rho}(\tilde{r}) \approx \tilde{\rho}_h(\tilde{r}) \equiv \frac{1}{(1 + \tilde{r}^2)^{3/2}}. \quad (4.109a)$$

The error in using this equation as an approximation to the isothermal sphere is less than 7% for $\tilde{r} < 4$. The surface density to which $\tilde{\rho}_h$ gives rise is (eq. 2.55)

$$\Sigma_h(\tilde{R}) = \frac{2}{1 + \tilde{R}^2}, \quad (4.109b)$$

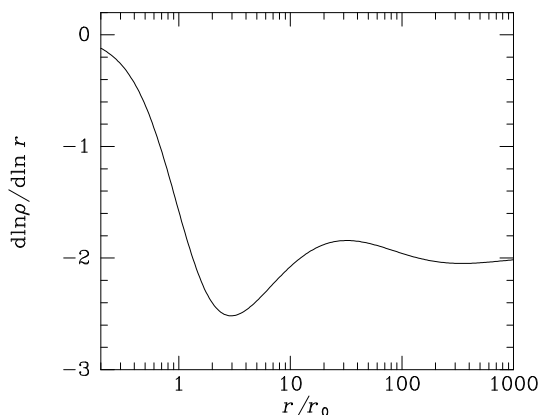


Figure 4.7 The logarithmic density gradient of the isothermal sphere. Note that the density gradient oscillates at large radii, a phenomenon that is explored in Problem 7.6.

where $\tilde{R} \equiv R/r_0$.

The density distribution $\tilde{\rho}_h$ does not fit the isothermal profile well at $\tilde{r} \gtrsim 3$ because it settles asymptotically to a logarithmic slope of -3 rather than -2 as is required by the isothermal profile. On the other hand, at large radii, $\tilde{\rho}_h$ has another use: when $\tilde{r} \gg 1$ the projected density (4.109b) to which it gives rise is very similar to the Hubble–Reynolds law (eq. 2.52), which fits the surface-brightness profiles of many elliptical galaxies rather well. Thus $\tilde{\rho}_h$ provides a simple analytical approximation to the inner parts of an isothermal sphere, or to the outer parts of a galaxy that obeys the Hubble–Reynolds law. It does *not* fit the outer parts of an isothermal sphere or the inner parts of the Hubble–Reynolds law. This dual application of $\tilde{\rho}_h$ has produced a certain amount of confusion in the literature.

From the astrophysical point of view, the isothermal sphere has a serious defect: its mass is infinite. Thus from equation (4.104), we have that $M \propto r$ at large r . No real astrophysical system can be modeled over more than a limited range of radii with a divergent mass distribution. On the other hand, the circular-speed curves of spiral galaxies (BM §8.2.4) are often remarkably flat out to great radii, and the divergent mass of the isothermal sphere is a useful reminder that we have little direct knowledge of the mass distribution in the outer parts of galaxies.

(c) Lowered isothermal models We seek a model that resembles the isothermal sphere at small radii, where the majority of stars have large values of the relative energy \mathcal{E} , but is less dense than the isothermal sphere at large radii, so its total mass is finite. We may obtain the DF f_K of such a model by simply diminishing the DF of the isothermal sphere at small values of \mathcal{E} . Thus we modify the DF (4.96) of the isothermal sphere in such a way that $f_K = 0$ for $\mathcal{E} \leq \mathcal{E}_0$. We may exploit the arbitrary constant Φ_0 in the definition (4.41) of \mathcal{E} to set the critical relative energy $\mathcal{E}_0 = 0$. Therefore $f_K(\mathcal{E})$ should be of the same form as equation (4.96) for $\mathcal{E} \gg 0$ and zero for

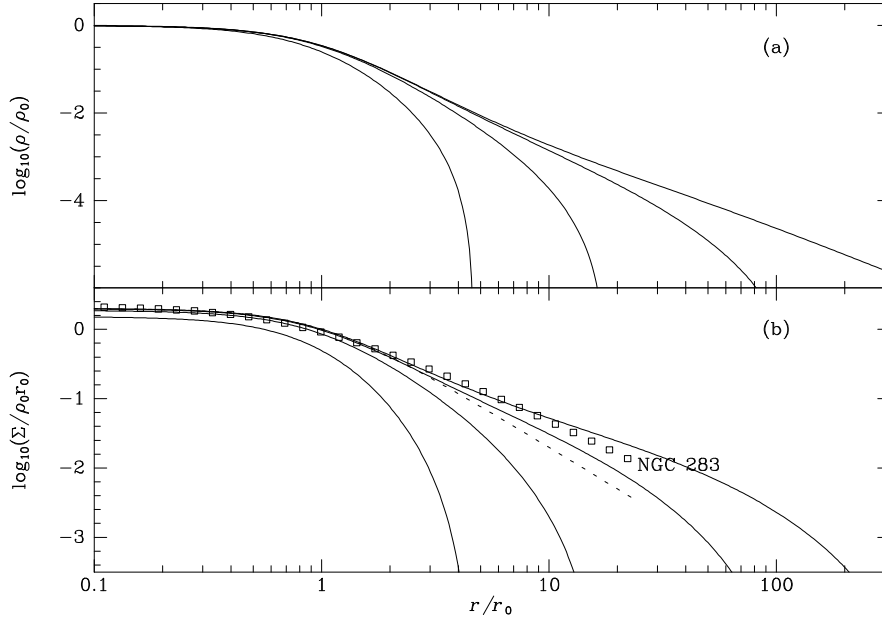


Figure 4.8 (a) Density profiles of four King models: from top to bottom the central potentials of these models satisfy $\Psi(0)/\sigma^2 = 12, 9, 6, 3$. (b) The projected mass densities of these models (full curves), and the projected modified Hubble model of equation (4.109b) (dashed curve). The squares show the surface brightness of the elliptical galaxy NGC 283 (Lauer et al. 1995).

$\mathcal{E} < 0$. A suitable function is

$$f_K(\mathcal{E}) = \begin{cases} \rho_1 (2\pi\sigma^2)^{-3/2} (e^{\mathcal{E}/\sigma^2} - 1) & \mathcal{E} > 0; \\ 0 & \mathcal{E} \leq 0. \end{cases} \quad (4.110)$$

This DF defines the family of **King models**.⁹ We now derive the density profiles and other properties of these models.

We proceed much as in the case of the isothermal sphere. Substituting into equation (4.110) for \mathcal{E} from equation (4.41) and integrating over all velocities, we obtain the density at any radius as

$$\begin{aligned} \rho_K(\Psi) &= \frac{4\pi\rho_1}{(2\pi\sigma^2)^{3/2}} \int_0^{\sqrt{2\Psi}} dv v^2 \left[\exp\left(\frac{\Psi - \frac{1}{2}v^2}{\sigma^2}\right) - 1 \right] \\ &= \rho_1 \left[e^{\Psi/\sigma^2} \operatorname{erf}\left(\frac{\sqrt{\Psi}}{\sigma}\right) - \sqrt{\frac{4\Psi}{\pi\sigma^2}} \left(1 + \frac{2\Psi}{3\sigma^2}\right) \right], \end{aligned} \quad (4.111)$$

⁹ DFs of the form (4.110) were actually introduced by Michie (1963) and studied in detail by Michie & Bodenheimer (1963), but King (1966) made them well known—see King (1981) for a discussion of their history.

where $\text{erf}(x)$ is the error function (Appendix C.3). Poisson's equation for Ψ may therefore be written

$$\frac{d}{dr} \left(r^2 \frac{d\Psi}{dr} \right) = -4\pi G \rho_1 r^2 \left[e^{\Psi/\sigma^2} \text{erf} \left(\frac{\sqrt{\Psi}}{\sigma} \right) - \sqrt{\frac{4\Psi}{\pi\sigma^2}} \left(1 + \frac{2\Psi}{3\sigma^2} \right) \right]. \quad (4.112)$$

We integrate this ordinary differential equation for $\Psi(r)$ outwards from $r = 0$, where we set $d\Psi/dr = 0$ and choose a value for Ψ . This value determines the central potential $\Phi(0)$ and total mass in the following implicit way. As we integrate equation (4.112) outward, $d\Psi/dr$ decreases, because initially $d\Psi/dr = 0$ and $d^2\Psi/dr^2 < 0$. As Ψ decreases towards zero, the range $(0, \sqrt{2\Psi})$ of speeds of stars at a given radius narrows, and the density of stars drops. Eventually at some radius r_t , when Ψ becomes equal to zero, the density vanishes. We call r_t the “tidal radius,” following the term observers use to denote the outermost limit of a cluster (page 30). The mass $M(r_t) = 4\pi \int_0^{r_t} dr r^2 \rho_K$ is the system's total mass, and the potential at the tidal radius is

$$\Phi(r_t) = -\frac{GM(r_t)}{r_t}. \quad (4.113)$$

The central potential is then $\Phi(0) = \Phi(r_t) - \Psi(0)$. The bigger the value $\Psi(0)$ from which we start our integration of equation (4.112), the greater will be the tidal radius, the total mass, and $|\Phi(0)|$.

Figure 4.8a shows the density profiles of King models obtained by integrating equation (4.112) from several values of $\Psi(0)$. The radial coordinate is marked in units of the King radius r_0 that is defined by equations (4.106). Figure 4.8b shows the projected density profiles $\Sigma_K(R)$ of the King models of Figure 4.8a. Notice that for some of these models r_0 is appreciably larger than the half-brightness, or core radius r_c , which is defined by the condition $\Sigma_K(r_c)/\Sigma_K(0) = \frac{1}{2}$ —see the discussion following equation (4.106). The dashed curve in Figure 4.8b shows the modified Hubble model (4.109b), which provides a moderately good fit to the projected surface density of the King model with central potential $\Psi(0) \simeq 8\sigma^2$. The squares in Figure 4.8b show the surface-brightness profile of the giant elliptical galaxy NGC 283. The King model with $\Psi(0) \simeq 10\sigma^2$ fits this profile fairly well.

The ratio of the tidal radius r_t to the King radius r_0 defines the **concentration** c through

$$c \equiv \log_{10}(r_t/r_0). \quad (4.114)$$

King models form a sequence that may be parameterized in terms of either c or $\Psi(0)/\sigma^2$. Figure 4.9 gives the relationship between c and $\Psi(0)/\sigma^2$. In the limit $c \rightarrow \infty$, $\Psi(0)/\sigma^2 \rightarrow \infty$, the sequence of King models goes over into the isothermal sphere.

Figure 4.10 shows the dependence on the concentration of the half-mass radius r_h (left panel) and the ratio r_h/r_g of the half-mass and gravitational radii (right panel). We see that although r_h/r_0 ranges over two orders of

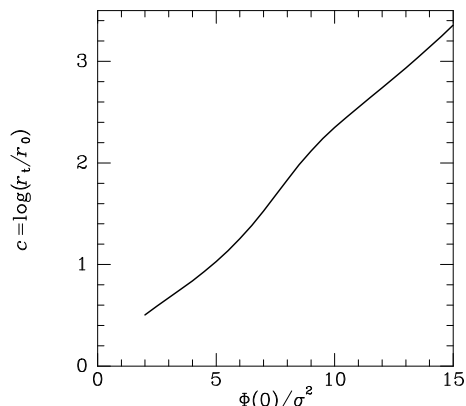


Figure 4.9 The relationship between the concentration c of a King model (eq. 4.114) and the central potential $\Psi(0)$ from which equation (4.112) is integrated.

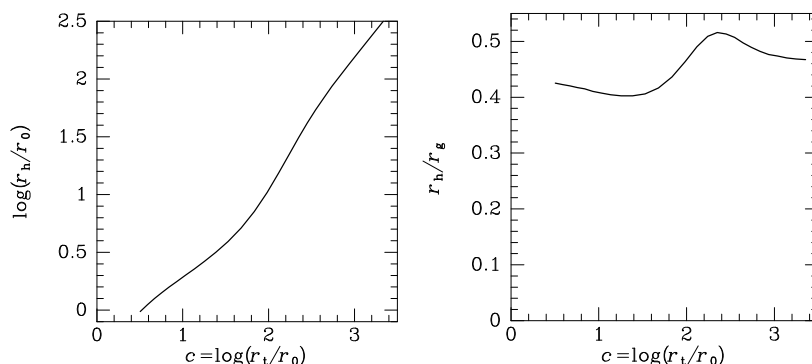


Figure 4.10 The half-mass radius r_h (left) and the ratio r_h/r_g of the half-mass radius to the gravitational radius (2.42) as a function of the concentration of a King model.

magnitude along the sequence of King models, r_h/r_g is confined to the interval (0.4, 0.51).

At each point on the sequence of King models, there is a two-parameter family of systems that are related to each other by changes of scale. Thus for any value of c there are models having any given values of two of the parameters r_0 , ρ_0 and σ , with the third then being determined by equation (4.106).

The parameter σ that occurs in the relations we have given for King models must not be confused with the actual velocity dispersion $\sigma_r = \sigma_\theta = \sigma_\phi$ of the stars of the system, or with the line-of-sight dispersion σ_\parallel . Figure 4.11 is a plot of σ_r/σ and σ_\parallel/σ for several King models. One sees that in all these models the velocity dispersion falls monotonically from the center outward, reaching zero at r_t . The velocity dispersion of the stars at r_t is zero because the potential energy of these stars is already equal to the largest energy allowed to any star.

The King DF (4.110) is only one of several possible modified isothermal

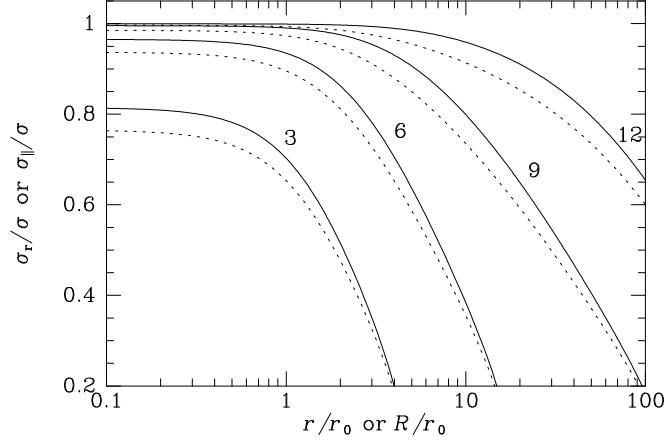


Figure 4.11 The one-dimensional velocity dispersion $\sigma_r = \sigma_\theta = \sigma_\phi$ at a given spatial radius r (full curves) and the RMS line-of-sight velocity σ_{\parallel} at projected radius R (dashed curves) for the King models shown in Figure 4.8. The curves are labeled by $\Psi(0)/\sigma^2$.

DFs. Woolley & Dickens (1961) discussed models for which f is given by (4.96) for $\mathcal{E} > 0$ and is zero otherwise, while Wilson (1975) (see also Hunter 1977) considered models generated by DFs of the form

$$f_W = \text{constant} \times \begin{cases} [e^{\mathcal{E}/\sigma^2} - 1 - (\mathcal{E}/\sigma^2)] & \text{for } \mathcal{E} > 0; \\ 0 & \text{otherwise.} \end{cases} \quad (4.115)$$

We now examine some spherical systems generated by DFs that depend on both \mathcal{E} and L , with the result that their velocity-dispersion tensors are anisotropic.

(d) Double-power models In Problem 4.6 it is shown that a DF of the form $L^\gamma f_1(\mathcal{E})$ generates a model in which the anisotropy parameter β (eq. 4.61) is at all radii equal to $-\frac{1}{2}\gamma$. By adding two DFs of this type, say $L^{\gamma_1} f_1(\mathcal{E})$ and $L^{\gamma_2} f_2(\mathcal{E})$, we can generate a model in which β is a function of radius. Generalizing this idea, we are led to consider DFs of the form $f(\mathcal{E}, L) = \sum_\gamma L^\gamma f_\gamma(\mathcal{E})$. We can take the idea of power-series expansion one step further by supposing that the functions f_γ can also be expanded as power series. The DF is then of the form

$$f(\mathcal{E}, L) = \sum_{\gamma\delta} \alpha_{\gamma\delta} L^\gamma \mathcal{E}^\delta, \quad (4.116)$$

where the $\alpha_{\gamma\delta}$ are arbitrary constants. The indices γ and δ are not necessarily integers. The density will be real and finite only if $\gamma > -2$, and a lower limit on δ will be required to ensure that the density diminishes sufficiently fast

as $r \rightarrow \infty$. Kent & Gunn (1982) investigated self-gravitating models with DFs that consist of a single term in the series of (4.116).

(e) Michie models A natural extension of King models to include velocity anisotropy is the family of **Michie models** defined by the DF

$$f_M(\mathcal{E}, L) = \begin{cases} \rho_1 (2\pi\sigma^2)^{-3/2} e^{-L^2/(2r_a^2\sigma^2)} (e^{\mathcal{E}/\sigma^2} - 1) & (\mathcal{E} > 0) \\ 0 & (\mathcal{E} \leq 0). \end{cases} \quad (4.117)$$

In the limit $r_a \rightarrow \infty$ this DF reduces to the DF (4.110) of the King models, and for $\mathcal{E} \gg \sigma^2$ the variables \mathcal{E} and L occur only through the variable Q that we defined in connection with Osipkov–Merritt models (eq. 4.73). Consequently, in a Michie model the velocity distribution is isotropic at the center, nearly radial in the outer parts, and the transition occurs near the anisotropy radius r_a . These models are fully described in Michie & Bodenheimer (1963).

4.4 DFs for axisymmetric density distributions

In §3.2 we saw that most orbits in an axisymmetric potential admit three isolating integrals, H , L_z and some third integral I_3 , for which we have an analytic expression only if the potential has the Stäckel form (§3.5.3). Since I_3 , when known, has a complicated functional form, models in which the DF is a function $f(H, L_z)$ of the two “classical” integrals are the only ones susceptible to analytic treatment.

4.4.1 DF for a given axisymmetric system

Equations (4.39) express a DF of the form $f(\mathcal{E}, L_z)$ as the sum of parts $f_{\pm}(\mathcal{E}, L_z)$ that are even and odd in L_z . We have seen that the probability density $\nu(R, z)$ is independent of f_- , while the azimuthal flux $\nu\bar{v}_\phi$ is independent of f_+ . Lynden–Bell (1962a) extended Eddington’s work on spherical systems to show that $f_+(\mathcal{E}, L_z)$ can be deduced if we are given the density $\nu(R, z)$ and its confining potential $\Phi(R, z)$, and that $f_-(\mathcal{E}, L_z)$ can be recovered if $\bar{v}_\phi(R, z)$ is also known. To prove these results, we use cylindrical coordinates (v_m, ψ, v_ϕ) for velocity space, with the polar axis in the azimuthal direction, so $v_R = v_m \cos \psi$, $v_z = v_m \sin \psi$. We have $d^3\mathbf{v} = v_m dv_m d\psi dv_\phi$ and the Jacobian determinant $\partial(\mathcal{E}, L_z)/\partial(v_\phi, v_m) = Rv_m$, so $d^3\mathbf{v} = R^{-1}d\mathcal{E}dL_zd\psi$. Hence

$$\nu(R, z) = \frac{4\pi}{R} \int_0^\Psi d\mathcal{E} \int_0^{R\sqrt{2(\Psi-\mathcal{E})}} dL_z f_+(\mathcal{E}, L_z^2), \quad (4.118)$$

where we have taken f_+ to be a function of L_z^2 rather than of L_z . We assume that ν is symmetrical about the equatorial plane, and is therefore a function

of $|z|$. At fixed R , $|z|$ is a monotone function of Ψ , so we can consider ν to be a function of (R, Ψ) rather than (R, z) . On this understanding we now differentiate (4.118) with respect to Ψ :

$$\frac{\partial}{\partial \Psi} \nu(R, \Psi) = 4\pi \int_0^\Psi d\mathcal{E} \frac{f_+[\mathcal{E}, 2(\Psi - \mathcal{E})R^2]}{\sqrt{2(\Psi - \mathcal{E})}}. \quad (4.119)$$

This integral equation for f_+ can be solved by taking Laplace transforms, so we multiply each side by $\exp(-s\Psi)$ and integrate over Ψ . After interchanging the order of the integrals over Ψ and \mathcal{E} we have

$$\begin{aligned} \int_0^\infty d\Psi e^{-s\Psi} \frac{\partial \nu}{\partial \Psi} &= 4\pi \int_0^\infty d\mathcal{E} e^{-s\mathcal{E}} \int_{\mathcal{E}}^\infty d\Psi e^{-s(\Psi - \mathcal{E})} \frac{f_+[\mathcal{E}, 2(\Psi - \mathcal{E})R^2]}{\sqrt{2(\Psi - \mathcal{E})}} \\ &= \frac{2\pi}{R} \int_0^\infty d\mathcal{E} e^{-s\mathcal{E}} \int_0^\infty du e^{-su/(2R^2)} g(\mathcal{E}, u), \end{aligned} \quad (4.120a)$$

where

$$g(\mathcal{E}, u) \equiv \frac{f_+(\mathcal{E}, u)}{\sqrt{u}}. \quad (4.120b)$$

The integrals on the right side of (4.120a) effect Laplace transforms of g with respect to each of its arguments. Let a hat denote Laplace transformation, so

$$\begin{aligned} \hat{g}(s, t) &\equiv \int_0^\infty d\mathcal{E} e^{-s\mathcal{E}} \int_0^\infty du e^{-tu} g(\mathcal{E}, u); \\ \hat{\nu}(R, s) &\equiv \int_0^\infty d\Psi e^{-s\Psi} \nu(R, \Psi). \end{aligned} \quad (4.121)$$

Then comparing the first of equations (4.121) with (4.120a) we see that $t = s/(2R^2)$. Integrating the left side of (4.120a) by parts and setting $R = \sqrt{s/2t}$, we have finally

$$\hat{g}(s, t) = \frac{s^{3/2}}{\pi(8t)^{1/2}} \hat{\nu} \left(\sqrt{\frac{s}{2t}}, s \right). \quad (4.122)$$

Since a function is uniquely defined by its Laplace transform, this result establishes that f_+ is determined by $\nu(R, z)$. The demonstration that f_- is determined by $\nu \bar{v}_\phi$ is similar. As in the case of Eddington's inversion, there is no guarantee that the recovered DF will be non-negative.

Two practical difficulties arise when attempting to use Lynden-Bell's result (4.122). First, only a very few models permit one analytically to replace z by Ψ in $\nu(R, z)$. Second, inverting the Laplace transforms in (4.122) usually requires analytic continuation of the function $\nu(R, \Psi)$ to the complex plane, which is problematic if $\nu(R, \Psi)$ is derived from observational data. On account of these difficulties, equation (4.122) has yielded disappointingly few

DFs in practice, although Lynden–Bell (1962a) was able to recover the DF of a rotating Plummer model. This situation improved when Hunter & Qian (1993) and Qian et al. (1995) showed that the DF could be obtained from a contour integral, without explicitly eliminating z between ν and Ψ , and without extensive analytic continuation.

4.4.2 Axisymmetric systems specified by $f(H, L_z)$

In §4.3.3 we investigated some spherical systems that were defined by their DFs. We now extend this approach to the axisymmetric case: we start with a simple functional form for f , and then use Poisson’s equation to recover the observable properties of the model. As in the spherical case, for most choices of the DF Poisson’s equation has to be solved numerically. We start by describing some exceptions to this rule. We adopt the convention that f is the phase-space density of mass.

(a) Fully analytic models Fricke (1951) expanded the DF of an axisymmetric galaxy in powers of \mathcal{E} and L_z . Since L_z can be positive or negative, it can occur only in integral powers, but we are free to include non-integral powers of \mathcal{E} . We use the term **Fricke component** to denote a DF of the form $\mathcal{E}^\gamma L_z^{2n}$, where the 2 in the exponent of L_z ensures that the component is even in L_z and thus has non-vanishing density. With equation (4.118) the density is

$$\rho(R, z) = \frac{4\pi}{R} \int_0^\Psi d\mathcal{E} \mathcal{E}^\gamma \int_0^{R\sqrt{2(\Psi-\mathcal{E})}} dL_z L_z^{2n}. \quad (4.123)$$

When we do the integral over L_z and change the second variable of integration to $x \equiv \mathcal{E}/\Psi$, we obtain

$$\begin{aligned} \rho(R, z) &= \frac{2^{n+5/2}}{2n+1} \pi R^{2n} \Psi^{\gamma+n+3/2} \int_0^1 dx x^\gamma (1-x)^{n+1/2} \\ &= 2^{n+3/2} \frac{\gamma!(n-\frac{1}{2})!}{(\gamma+n+\frac{3}{2})!} \pi R^{2n} \Psi^{\gamma+n+3/2}. \end{aligned} \quad (4.124)$$

Thus the density of a Fricke component is a power of R multiplied by a power of Ψ . Below we shall use (4.124) in the form

$$\rho = R^{2n} \Psi^\gamma \Leftrightarrow f = \frac{2^{-(n+3/2)} \gamma!}{\pi(n-\frac{1}{2})!(\gamma-n-\frac{3}{2})!} L_z^{2n} \mathcal{E}^{\gamma-n-3/2}. \quad (4.125)$$

Toomre (1982) discovered a closely related one-parameter sequence of models with remarkably simple analytic properties. The DFs of Toomre’s models are proportional to $L_z^{2n} e^{\mathcal{E}/\sigma^2}$ and their density profiles have the scale-free form $\rho(r, \theta) \propto r^{-2} S(\theta)$, where the function S depends on the parameter

n in the DF. In limiting cases these systems go over into (i) the isothermal sphere, (ii) Spitzer's isothermal sheet (Problem 4.21), and (iii) cold Mestel disks (see §4.5.1 below). Like Fricke components, for $n > 0$ Toomre's models have zero density on the z axis.

Fully analytical models that have finite central densities were discovered by Evans (1993, 1994) by assuming that the relative potential is a function $\Psi(m^2)$ of the spheroidal variable

$$m^2 = R_c^2 + R^2 + \frac{z^2}{q_\Phi^2}. \quad (4.126)$$

Inserting this form of Ψ into Poisson's equation, we find after some algebra that the density is given by

$$-4\pi G\rho = \left(4 + \frac{2}{q_\Phi^2}\right)\Psi' + \frac{4(m^2 - R_c^2)}{q_\Phi^2}\Psi'' + 4\left(1 - \frac{1}{q_\Phi^2}\right)\Psi''R^2, \quad (4.127)$$

where primes denote differentiation with respect to the argument of the function. Since Ψ' , Ψ'' , and Ψ are all functions only of m^2 , we may straightforwardly express Ψ' and Ψ'' as functions of Ψ . Once this has been done, equation (4.127) gives ρ as a function $\rho(R, \Psi)$, and the model's DF can be recovered from the formulae of Lynden-Bell (1962a) or Hunter & Qian (1993) as described in §4.4.1.

If Φ is either a power or the logarithm of m^2 , we can recover the DF more economically using Fricke's formula (4.125): if $\Psi = \Psi_a(R_c^2/m^2)^y$, where $y > 0$ and Ψ_a is a constant, then

$$\begin{aligned} m^2 &= R_c^2(\Psi/\Psi_a)^{-1/y}, \\ \Psi' &= -y\Psi_a R_c^{2y}(m^2)^{-(y+1)} = -y\frac{\Psi_a}{R_c^2}(\Psi/\Psi_a)^{1+1/y}, \\ \Psi'' &= y(y+1)\Psi_a R_c^{2y}(m^2)^{-(y+2)} = y(y+1)\frac{\Psi_a}{R_c^4}(\Psi/\Psi_a)^{1+2/y}, \end{aligned} \quad (4.128)$$

and our expression (4.127) for $\rho(R, \Psi)$ becomes

$$\begin{aligned} \rho(R, \Psi) &= \frac{y[2 - (2y+1)q_\Phi^{-2}]\Psi_a}{2\pi G R_c^2} \left(\frac{\Psi}{\Psi_a}\right)^{1+1/y} \\ &\quad + \frac{y(y+1)\Psi_a}{\pi q_\Phi^2 G R_c^2} \left(\frac{\Psi}{\Psi_a}\right)^{1+2/y} \left(1 + (1 - q_\Phi^2)\frac{R^2}{R_c^2}\right). \end{aligned} \quad (4.129)$$

This is just a sum of three Fricke components, so equation (4.125) enables us to say that the DF is

$$f_{\text{pow}}(\mathcal{E}, L_z) = A\mathcal{E}^{1/y-1/2} + B(1 + CL_z^2/\mathcal{E})\mathcal{E}^{2/y-1/2}, \quad (4.130a)$$

where

$$\begin{aligned} A &\equiv \frac{y[2 - (2y + 1)q_\Phi^{-2}](1 + 1/y)!}{(2\pi)^{5/2}(1/y - 1/2)!GR_c^2\Psi_a^{1/y}} \\ B &\equiv \frac{2y(y + 1)(1 + 2/y)!}{(2\pi)^{5/2}(2/y - 1/2)!q_\Phi^2GR_c^2\Psi_a^{2/y}} \end{aligned} \quad ; \quad C \equiv \frac{1 - q_\Phi^2}{R_c^2} \left(\frac{2}{y} - \frac{1}{2} \right). \quad (4.130b)$$

We shall call a model with this DF an **Evans model**. The model is unphysical if the phase-space density is ever negative. Negative densities do not arise so long as $1 \geq q_\Phi^2 \geq y + \frac{1}{2}$. For $y = \frac{1}{2}$ the only physical model is spherical—in fact it is the Plummer model. As y decreases, ever flatter models become possible.

In the limit $y \rightarrow 0$, the potential of an Evans model tends¹⁰ to the logarithmic potential $\Phi_L(R, z)$ of equation (2.71a), and the self-consistent DF becomes (Evans 1993)

$$f_{\log}(\mathcal{E}, L_z) = Ae^{\mathcal{E}/\sigma^2} + B(1 + CL_z^2)e^{2\mathcal{E}/\sigma^2}, \quad (4.131a)$$

where $\sigma^2 \equiv \frac{1}{2}v_0^2$ and

$$A \equiv \frac{2q_\Phi^2 - 1}{(2\pi)^{5/2}Gq_\Phi^2\sigma} \quad ; \quad B \equiv \frac{R_c^2}{\pi^{5/2}Gq_\Phi^2\sigma} \quad ; \quad C \equiv \frac{2(1 - q_\Phi^2)}{R_c^2\sigma^2}. \quad (4.131b)$$

Evans's DF (4.131a) has an obvious application as the DF of a population of dark-matter particles, since circular-speed curves tend to become flat at large radii, which implies that the potential is proportional to $\ln R$ (§1.1.3). The density of any luminous population falls off faster at large r . Evans observed that if the DF of luminous matter is

$$f_{\text{lum}} = \rho_0 R_c^p \left(\frac{p}{2\pi\sigma^2} \right)^{3/2} e^{p\mathcal{E}/\sigma^2}, \quad (4.132a)$$

where $p \geq 1$ is a constant, then the density of luminous matter confined by Φ_L is

$$\rho_{\text{lum}}(R, z) = \frac{\rho_0 R_c^p}{(R_c^2 + R^2 + z^2/q_\Phi^2)^{p/2}}, \quad (4.132b)$$

which is proportional to R^{-p} at large R . The combined density of luminous and dark matter will generate Φ_L if the DF of the dark matter is $f_{\text{dark}} = f_{\log} - f_{\text{lum}}$. Evans shows that f_{dark} is non-negative provided the fraction of the central mass density that is luminous is smaller than a specified function of q_Φ and p .

¹⁰ This is something of an oversimplification. We have $\ln x = \lim_{\alpha \rightarrow 0} \int_1^x dy/y^{1+\alpha} = \lim_{\alpha \rightarrow 0} \alpha^{-1}(1 - 1/x^\alpha)$, so $x^{-\alpha}$ tends to one minus a (small) multiple of $\ln x$.

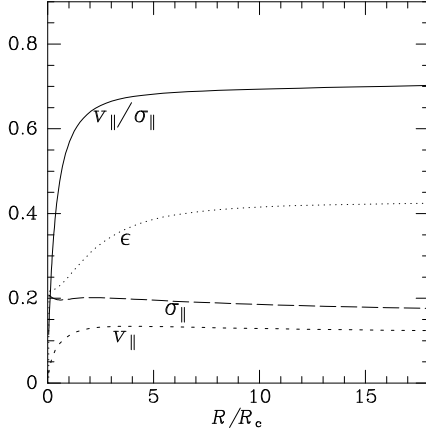


Figure 4.12 Projected quantities along the major axis of an Evans model that is an approximate isotropic rotator, seen edge-on (eq. 4.134). The shape of the isophotes is shown by $\epsilon = 1 - b/a$, where a and b are the isophote's intersections with the principal axes. The parameters of the model are $y = 0.09$, $q_\Phi = 0.85$ and $\alpha = 0.813$, and the velocities are in units of $\sqrt{\Psi_a}$.

The DF (4.132a) of luminous matter depends on \mathcal{E} only, with the consequence that the surfaces of constant luminosity density coincide with equipotentials. It is straightforward to write down DFs for luminous matter that depend on L_z as well as \mathcal{E} in such a way that the luminosity distribution is more flattened than the potential—see Evans (1993) for examples.

The DFs given by both (4.130a) and (4.131a) are even functions of L_z , so they produce non-rotating models. Rotating models can be produced by adding any DF f_- that is odd in L_z and small enough in magnitude to ensure that the composite DF is non-negative. For example, if $\alpha(x)$ is an odd function with absolute value less than unity, then we can choose f_- to be

$$f_-(\mathcal{E}, L_z) = \alpha(L_z)f_+(\mathcal{E}, L_z). \quad (4.133)$$

With this choice of f_- the DF is equal to $(1 + \alpha)f_+$.

A model in which $\sigma_\phi = \sigma_R = \sigma_z$ everywhere is called an **isotropic rotator** and can be said to be flattened by rotation alone. For any DF of the form $f(\mathcal{E}, L_z)$ we have $\sigma_R = \sigma_z$, and by choosing α such that $\sigma_\phi = \sigma_R$ at some particular point, we can generate a good approximation to an isotropic rotator.

If we apply to an Evans model the prescription of equation (4.133) with $\alpha(L_z) = \text{sgn}(L_z)\alpha_0$, where α_0 is a constant, the mean velocity becomes

$$\begin{aligned} \bar{v}_\phi &= \frac{4\pi\alpha_0}{\rho R^2} \int_0^\Psi d\mathcal{E} \int_0^{R\sqrt{2(\Psi-\mathcal{E})}} dL_z L_z f_{\text{pow}}(\mathcal{E}, L_z) \\ &= \frac{16\pi\alpha_0 y^2}{\rho} \Psi^{1/y+3/2} \left\{ \frac{A}{(2+y)(2+3y)} + \frac{B\Psi^{1/y}(4-y+4yCR^2)}{(4+3y)(16-y^2)} \right\}. \end{aligned} \quad (4.134)$$

Figure 4.12 shows v_\parallel , σ_\parallel and the isophote ellipticity ϵ along the major axis of an edge-on Evans model that is an approximate isotropic rotator. Both ϵ and $v_\parallel/\sigma_\parallel$ increase with distance from the center. This model is somewhat

unrealistic because our choice of $\alpha(x)$ has introduced a discontinuity in f at $L_z = 0$, which in turn causes a vortex along the symmetry axis of the model. Nevertheless, the projected velocities shown in Figure 4.12 reflect those of more realistic Evans models.

(b) Rowley models In general, when a DF is chosen for an axisymmetric system, we have to solve for the associated density and potential distributions numerically. We now describe an effective technique for this job, which differs significantly from the method we used to solve the corresponding problem in the spherical case. Prendergast & Tomer (1970), Jarvis & Freeman (1985) and Rowley (1988) discuss this problem.

As in the spherical case, the DF cannot be correctly normalized at the outset. So we write

$$f \propto \mathcal{F}(\mathcal{E}, L_z), \quad (4.135)$$

where \mathcal{F} is our chosen functional form. The model is specified by this form and the values Ψ_c and Ψ_t taken by the relative potential at the center and at an outer point, for example the point $(R_t, 0)$ in the equatorial plane at which the density first vanishes. We proceed iteratively. We guess what potential $\Phi_0(R, z)$ our DF will generate at a grid of points. We take the relative potential that appears in \mathcal{E} to be

$$\Psi_0(R, z) = \psi - \lambda \Phi_0(R, z), \quad (4.136)$$

where ψ and λ are constants that we choose such that $\Psi_0(0, 0) = \Psi_c$ and $\Psi_0(R_t, 0) = \Psi_t$. Then at each grid point we evaluate the density

$$\rho_1(R, z) \equiv \int d^3\mathbf{v} \mathcal{F}[\Psi_0(R, z) - \frac{1}{2}v^2, Rv_\phi]. \quad (4.137)$$

and from these values solve for the corresponding potential $\Phi_1(R, z)$. At the next iteration we use the relative potential $\Psi_1 = \psi - \lambda \Phi_1$, where ψ and λ are chosen afresh to ensure that Ψ_1 satisfies the same conditions at the origin and $(R_t, 0)$ as Ψ_0 did. We repeat this process until Ψ_n and ρ_n change very little between iterations. The final model has DF $\lambda \mathcal{F}$ and density $\lambda \rho_n$.

Proceeding in this way, Rowley (1988) constructed models using the functional form

$$\mathcal{F}(\mathcal{E}, L_z) = \begin{cases} e^{\chi/\sigma^2} & \text{for } \chi > 0, \\ 0 & \text{otherwise,} \end{cases} \quad (4.138a)$$

where

$$\chi \equiv \mathcal{E} + \omega L_z - \frac{1}{2}L_z^2/r_a^2. \quad (4.138b)$$

Here σ , ω and r_a are parameters. Let $v_m \equiv \sqrt{v_R^2 + v_z^2}$. Then χ can be written

$$\begin{aligned} \chi &= \Psi - \frac{1}{2}v_m^2 - \frac{1}{2}\left(1 + \frac{R^2}{r_a^2}\right)v_\phi^2 + \omega Rv_\phi \\ &= \Psi - \frac{1}{2}v_m^2 - \frac{1}{2}\left(1 + \frac{R^2}{r_a^2}\right)\left(v_\phi - \frac{\omega R}{1 + R^2/r_a^2}\right)^2 + \frac{\omega^2 R^2}{2(1 + R^2/r_a^2)}. \end{aligned} \quad (4.139)$$

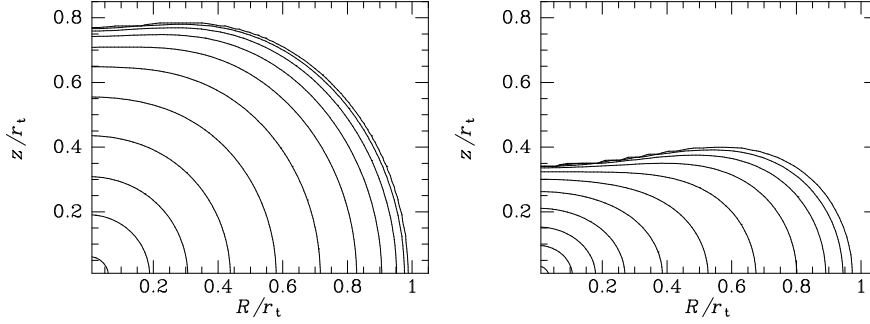


Figure 4.13 Density in the meridional plane of Rowley models that have $\omega R_t/\sigma = 2$ (left) and 4 (right). Both models have $\Psi(0,0) - \Psi(R_t,0) = 4\sigma^2$ and $r_a = 0.45R_t$.

Thus the DF (4.138) is a truncated Gaussian in v_R , v_z , and $v'_\phi \equiv v_\phi - \bar{v}_\phi$, where the mean speed is

$$\bar{v}_\phi(R) = \frac{\omega R}{1 + R^2/r_a^2}. \quad (4.140)$$

Since \bar{v}_ϕ depends only on R , the streaming velocity is independent of distance $|z|$ from the equatorial plane: the system is said to rotate “on cylinders”. Equation (4.140) shows that we have solid-body rotation $\bar{v}_\phi \simeq \omega R$ for $R \ll r_a$. At $R \gg r_a$ the mean streaming speed decays to zero as $1/R$.

Integrating over the truncated Gaussians we find that

$$\begin{aligned} \rho(R, z) &= \frac{4\pi\sigma^3}{\sqrt{1 + R^2/r_a^2}} \left(\sqrt{\pi/2} e^{\tilde{\Psi}} \operatorname{erf}(\sqrt{\tilde{\Psi}}) - \sqrt{2\tilde{\Psi}} \right) \\ \frac{\overline{v_R^2}}{\sigma^2} = \frac{\overline{v_z^2}}{\sigma^2} &= \left(1 + \frac{R^2}{r_a^2} \right) \frac{\overline{\tilde{v}_\phi^2}}{\sigma^2} = 1 - \frac{(2\tilde{\Psi})^{3/2}}{3 \left(\sqrt{\pi/2} e^{\tilde{\Psi}} \operatorname{erf}(\sqrt{\tilde{\Psi}}) - \sqrt{2\tilde{\Psi}} \right)}, \end{aligned} \quad (4.141a)$$

where $\tilde{v}_\phi \equiv v_\phi - \bar{v}_\phi$ and

$$\tilde{\Psi} \equiv \frac{1}{\sigma^2} \left(\Psi + \frac{\omega^2 R^2}{2(1 + R^2/r_a^2)} \right). \quad (4.141b)$$

For large values of $\tilde{\Psi}$ the density grows exponentially with Ψ as in the isothermal sphere, and the velocity dispersion in the meridional plane tends to the parameter σ . At $r \ll r_a$, the model is isotropic, but beyond the anisotropy radius r_a , the dispersion in the azimuthal direction falls below that in the meridional plane in close analogy with the behavior of an Osipkov–Merritt model. The mean-square velocity in the azimuthal direction is

$$\overline{v_\phi^2} = \overline{v_\phi^2} + \overline{\tilde{v}_\phi^2} = \frac{\overline{v_R^2} + \omega^2 R^2}{1 + R^2/r_a^2} \simeq \overline{v_R^2} \frac{1 + R^2\omega^2/\sigma^2}{1 + R^2/r_a^2}. \quad (4.142)$$

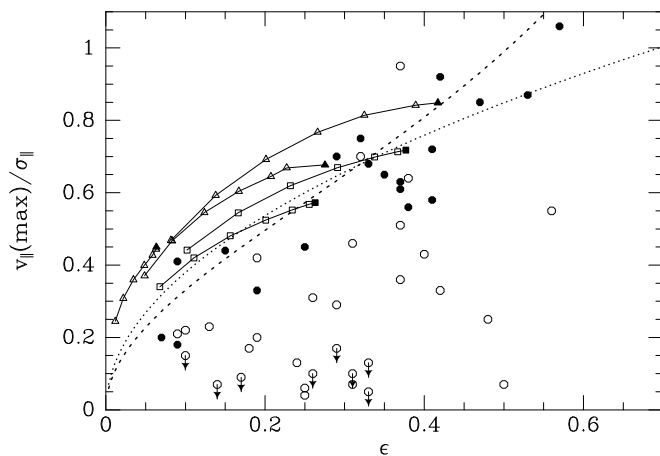


Figure 4.14 Circles show measured ratios of peak rotation speed to average velocity dispersion within half the effective radius R_e (page 21) for spheroidal systems (elliptical galaxies and the bulges of disk galaxies) from Davies et al. (1983). Filled circles are for systems less luminous than $M_R = -21.34 + 5 \log_{10} h_7$. The squares show this ratio for Evans models with $y = 0.09$ and $q_\Phi = 0.85$ and 0.9 ; the value of ϵ is for the isophote with semi-major axis length $15R_c$ and $\sigma_{||}$ is the mean value of the velocity dispersion on the major axis out to $15R_c$. The full squares are for edge-on models, and the inclination decreases by 10° between successive squares. The triangles show similar data for five Rowley models. The dotted curve shows the relation $v/\sigma = 1.2\sqrt{\epsilon}$ that is suggested by the Jeans equations in §4.8.2b, while the dashed curve shows $\pi/2^{3/2}$ times the relation (4.266a) (with $\alpha = \delta = 0$) that we derive in §4.8.3 from the tensor virial theorem. The observational values of v/σ for the less luminous galaxies are approximately consistent with the model predictions, while most luminous galaxies rotate much less rapidly for a given flattening.

In §4.8.3 we shall see that galaxies are flattened by an excess of kinetic energy in the equatorial plane relative to the meridional plane. When $\omega r_a/\sigma = 1$, the fraction on the right of (4.142) is unity, and there would be no excess, but when this dimensionless parameter exceeds unity, there is an excess. Hence we expect the flattening to be an increasing function of $\omega r_a/\sigma$. The other important dimensionless parameter of Rowley models is $\tilde{\Psi}(0, 0)$, which is analogous to the quantity that determines the concentration of a King model.

Figure 4.13 shows the density in the meridional plane of Rowley models that differ in their values of $\omega r_a/\sigma$. As this parameter increases the models rotate more rapidly and become more strongly flattened; the isodensity surfaces are not nearly ellipses and in projection these models have peanut shapes.

(c) Rotation and flattening in spheroids We have seen that mathematically the connection between the flattening of a spheroidal stellar system and its rotation rate is weak because the rotation rate is determined by f_- ,

while the density distribution is determined by f_+ . However, one might still anticipate a physical connection between flattening and rotation that arises from some aspect of the formation process. To establish a connection between the theoretical models and observational data, we must compute the properties of these models when viewed in projection with a given inclination angle i between the line of sight and the system's symmetry axis ($i = 90^\circ$ is edge-on, $i = 0$ is face-on). We assume that the mass-to-light ratio is independent of position, as is approximately true in the inner parts of elliptical galaxies (§4.9.2). The comparison is usually made in terms of two quantities: (i) the ellipticity $\epsilon = 1 - b/a$ of the isophotes at a given semi-major axis a , and (ii) the ratio v/σ of the peak value of the mean line-of-sight velocity \bar{v}_\parallel to the average of the line-of-sight dispersion σ_\parallel within some specified radius. The circles in Figure 4.14 show this measure for a number of elliptical galaxies and bulges of early-type disk galaxies. The circles are open if the spheroid is more luminous than $M_R = -21.34 + 5 \log_{10} h_7$ and otherwise filled. Although there is little or no correlation between v/σ and ϵ for the luminous systems, these variables are strongly correlated for the less luminous ones.

In Figure 4.14, each set of squares joined by lines shows v/σ for an Evans model (eq. 4.134) that is an approximate isotropic rotator, while the triangles show the analogous quantity for Rowley models. The filled symbols are for inclination 90° , the open symbols are for $i = 80^\circ, 70^\circ, \dots$. We see that when these models are viewed edge-on, the representative points of Evans models lie within the band that is populated by the less luminous spheroids, while the Rowley models lie at its upper edge. At smaller inclinations the Evans models lie along the edge of the band, while the Rowley models lie above it. We conclude that, for reasons that are not yet fully understood, on average low-luminosity spheroids are rotating slightly less rapidly than an isotropic rotator, whilst rotation plays almost no role in determining the shapes of luminous elliptical galaxies. In other words low-luminosity ellipticals are flattened by rotation but luminous ellipticals are not.

The dashed curve in Figure 4.14 shows¹¹ $\pi/2^{3/2} \simeq 1.1$ times a ratio of RMS rotation rate and velocity dispersion that we shall derive in §4.8.3c below. It provides a good fit to the data for low-luminosity spheroids and has historically been used to define the normalized rotation rate $(v/\sigma)^*$, which is the measured value of v/σ divided by the value predicted for the galaxy's ellipticity by this curve (BM §11.2.1).

4.4.3 The Schwarzschild DF

The stars in a galactic disk such as that of the Milky Way travel on nearly circular and coplanar orbits. Consequently, DFs that generate cool disks in which random velocities are much smaller than the circular speed are central to understanding disk galaxies. The mean radius of a star that is on a nearly

¹¹ The numerical factor is a relic of model-dependent assumptions in Binney (1978).

circular orbit near the equatorial plane is largely determined by its angular momentum L_z , or equivalently by its guiding-center radius $R_g(L_z)$ (eq. 3.72). In fact, in the epicycle approximation of §3.2.3 the mean radius is equal to the guiding-center radius. Thus the radial density profile of a cool disk is largely determined by the dependence of the DF upon L_z .

The difference

$$\Delta \equiv H - E_c(L_z) \quad (4.143)$$

between a star's energy and the energy $E_c(L_z)$ of the circular orbit with the same angular momentum is the energy associated with the star's gyrations around the guiding center. In a disk composed of many stars these oscillations—all with random phases—lead to a dispersion in velocities at each location in the disk, superimposed on the overall rotation. The distribution of stars with respect to Δ is what governs the velocity dispersion or temperature of the disk—in a cool disk all stars have small values of $\Delta/E_c(L_z)$. These considerations suggest that we examine DFs of the form

$$f(H, L_z) = S(L_z)T[\Delta/\sigma^2(L_z)], \quad (4.144)$$

where the function S is predominantly determined by the surface density $\Sigma(R)$, and the function T is chosen to fit the shape of the velocity distribution; $\sigma(L_z)$ determines the radial dependence of the velocity dispersion.

Unfortunately the DF (4.144) cannot reproduce all the properties of the solar neighborhood because observations show that $\overline{v_z^2}/\overline{v_R^2} \simeq 0.3 \neq 1$ (Table 1.2), while equation (4.40) shows that any DF that depends only on H and L_z requires these two dispersions to be equal. This situation motivates us to consider a more complex DF that also depends on the third integral I_3 that we studied in §3.2. We do not have an analytic expression for I_3 in a general axisymmetric potential, so we use the approximation (3.74)

$$I_3 \simeq H_z(z, \dot{z}, L_z) \equiv \frac{1}{2}\dot{z}^2 + \Phi_z(z, L_z), \quad (4.145)$$

where we have made the dependence of Φ_z on L_z explicit. We take advantage of equation (4.145) to generalize (4.144) to the form

$$f(H, L_z, I_3) \simeq S(L_z)T\left(\frac{\Delta}{\sigma_R^2}, \frac{H_z}{\sigma_3^2}\right), \quad (4.146)$$

where σ_R^2 and σ_3^2 are functions of L_z . An exponential function is an obvious choice for T , so we arrive at DFs of the form (Shu 1969)

$$f(H, L_z, H_z) = S(L_z) \exp\left(-\frac{\Delta}{\sigma_R^2} - \frac{H_z}{\sigma_3^2}\right). \quad (4.147)$$

If we introduce the epicycle approximation for motions parallel to the plane, then $\Delta \simeq H_R + H_z$, where H_R is given by equation (3.86) as

$$H_R \equiv \frac{1}{2}(\dot{x}^2 + \kappa^2 x^2) \quad (4.148)$$

with $x \equiv R - R_g$ and κ the epicycle frequency. Thus so long as σ_R and σ_3 are much smaller than the circular speed, we can write

$$f(H, L_z, I_3) \simeq S(L_z) \exp\left(-\frac{H_R}{\sigma_R^2} - \frac{H_z}{\sigma_z^2}\right), \quad (4.149a)$$

where

$$\sigma_z^2 \equiv \frac{\sigma_R^2 \sigma_3^2}{\sigma_R^2 + \sigma_3^2}. \quad (4.149b)$$

When we apply (4.149a) to the solar neighborhood, we encounter the problem that we do not directly measure the value of R_g for stars, so it is not straightforward to determine their x values. However, an observationally accessible quantity is the difference

$$\tilde{\mathbf{v}} \equiv \mathbf{v} - v_c(R)\hat{\mathbf{e}}_\phi \quad (4.150)$$

between the velocity of a star and the velocity of the circular orbit at the star's current location—in the solar neighborhood this is the Local Standard of Rest (LSR) (§1.1.2). We have from equation (3.97)

$$\tilde{v}_\phi = \frac{\kappa}{\gamma} x, \quad (4.151)$$

where $\gamma \equiv 2\Omega_g/\kappa$ (eq. 3.93b) and Ω_g is the circular frequency at R_g . Substituting equation (4.151) into (4.148), we have

$$H_R \simeq \frac{1}{2}(v_R^2 + \gamma^2 \tilde{v}_\phi^2). \quad (4.152)$$

Substituting this expression and equation (4.145) into equation (4.149a) we have finally

$$f \simeq f_{\text{Sch}} \equiv S(L_z) \exp\left(-\frac{v_R^2 + \gamma^2 \tilde{v}_\phi^2}{2\sigma_R^2(L_z)} - \frac{v_z^2 + 2\Phi_z(z, L_z)}{2\sigma_z^2(L_z)}\right). \quad (4.153)$$

The distribution f_{Sch} is called the **Schwarzschild DF**.¹²

We must now choose plausible forms for the three free functions $S(L_z)$, $\sigma_R(L_z)$ and $\sigma_z(L_z)$ that define the Schwarzschild DF. Equation (4.153) predicts that the distributions of v_R and v_z are superpositions of Gaussians, one for each value of L_z . If S , σ_R and σ_z vary sufficiently slowly with $L_z = R(v_c + \tilde{v}_\phi)$ that they change negligibly so long as $|\tilde{v}_\phi|$ is less than a

¹² Karl Schwarzschild (1873–1916) pioneered photographic photometry. In 1900 he published a lower limit on the radius of curvature of space. While serving in the German army from August 1914, he gave a quantum-mechanical explanation of the anomalous Stark effect in hydrogen and obtained the first and most important exact solution of Einstein's field equations.

few times σ_R , we may treat them as constants at a given position. Then $\sigma_R(Rv_c)$ and $\sigma_z(Rv_c)$ are simply the radial and vertical velocity dispersions at R . The oldest disk stars of the solar neighborhood have the largest random velocities, $\sigma_R \simeq 40 \text{ km s}^{-1}$, and even this value is much smaller than $v_c \simeq 220 \text{ km s}^{-1}$ (Table 1.2). Hence the approximation of constant σ_R and σ_z is a reasonable one, at least for cooler stellar populations.

We still have to choose $S(L_z)$, which we do by computing the surface density $\Sigma(R)$. Integrating over velocities in the approximation of constant σ_i , we find that the density of the disk is

$$\rho(R, z) = \int d^3\mathbf{v} f \simeq (2\pi)^{3/2} S(Rv_c) \left(\frac{\sigma_R^2 \sigma_z}{\gamma} \right)_{L_z=Rv_c} \exp\left(-\frac{\Phi_z(z, Rv_c)}{\sigma_z^2(Rv_c)}\right). \quad (4.154)$$

At any given radius, we therefore have $\rho(z) = \rho_0 \exp(-\Phi_z/\sigma^2)$, where $\rho_0(R)$ is the density in the midplane. In Problem 4.21 it is shown that when $\rho(z)$ takes this form and $\Phi_z(z)$ is determined self-consistently by the density in the disk, the surface density is $\Sigma = 4z_0\rho_0$, where $z_0 = \sigma_z(8\pi G\rho_0)^{-1/2}$ is the scale height of the disk. Thus

$$\Sigma(L_z/v_c) \simeq \Sigma(R) \simeq 4(2\pi)^{3/2} S(L_z) \left(\frac{\sigma_R^2 \sigma_z z_0}{\gamma} \right)_{L_z}. \quad (4.155)$$

This approximate result relating $S(L_z)$ and $\Sigma(R)$ holds only when the dispersions $\sigma_R, \sigma_z \ll v_0$. Nevertheless it is useful to use this result to eliminate $S(L_z)$ from equation (4.153), recognizing that the surface density derived from the resulting DF will not be exactly equal to $\Sigma(R)$. Thus we write

$$f_{\text{Schw}}(\mathbf{w}) \simeq \frac{\gamma \Sigma(L_z/v_c)}{4(2\pi)^{3/2} \sigma_R^2 \sigma_z z_0} \exp\left(-\frac{v_R^2 + \gamma^2 \tilde{v}_\phi^2}{2\sigma_R^2(L_z)} - \frac{v_z^2 + 2\Phi_z(z, L_z)}{2\sigma_z^2(L_z)}\right). \quad (4.156)$$

An important application of equation (4.156) is to the case of an exponential disk, $\Sigma(R) = \Sigma_0 \exp(-R/R_d)$, that rotates in a potential with a constant circular speed, so we can replace $v_c(R)$ with v_0 and set $\gamma = \sqrt{2}$. Observations of edge-on disks show that z_0 is approximately independent of radius, so from equation (4.302c) of Problem 4.21 it follows that $\sigma_z \propto \exp(-R/2R_d) \simeq \exp(-L_z/2R_d v_0)$. Let us assume that the ratio of dispersions σ_z/σ_R is independent of radius. Then the right side of equation (4.156) is completely determined and we can examine the predictions it makes. Figure 4.15 shows for the solar neighborhood the predicted distributions of azimuthal velocities \tilde{v}_ϕ for three populations of stars. The sharply peaked distribution is for an extremely cold population, which has $v_R^2/2 = 5 \text{ km s}^{-1}$, while the broadest distribution is for a population that has $v_R^2/2 = 30 \text{ km s}^{-1}$. The narrow distribution is almost Gaussian, and in

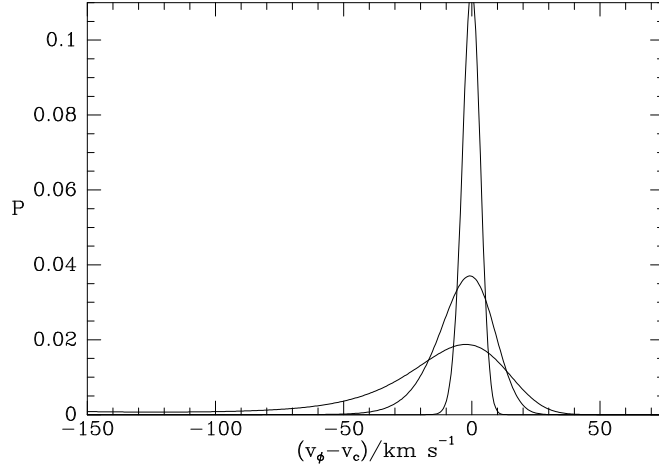


Figure 4.15 Three distributions of azimuthal velocities \tilde{v}_ϕ predicted for stellar populations in the solar neighborhood by the DF (4.156). The circular speed has been assumed to be $v_0 = 220 \text{ km s}^{-1}$ at all radii, $\sigma_R(L_z)$ and $\sigma_z(L_z)$ are taken to be proportional to $\exp[-L_z/(2v_0 R_d)]$, while $\Sigma = \Sigma_0 \exp(-R/R_d)$, with $R_0/R_d = 3.2$ (Table 1.2). The values of $\overline{v_R^2}^{-1/2}$ for the three populations are 5, 15 and 30 km s^{-1} , the largest value producing the widest spread in \tilde{v}_ϕ .

fact the three-dimensional velocity distribution of such a population would conform to the triaxial Gaussian model that Schwarzschild (1907) derived from observations of solar-neighborhood stars.

$$dn \propto \exp\left(-\frac{v_R^2 + \gamma^2 \tilde{v}_\phi^2}{2\sigma_R^2} - \frac{v_z^2}{2\sigma_z^2}\right) d^3\mathbf{v}, \quad (4.157)$$

where now σ_R and σ_z are constants.

The broadest distribution of azimuthal velocities in Figure 4.15 is extremely skew, with a long tail to highly negative values of \tilde{v}_ϕ and a sharp cutoff for $\tilde{v}_\phi > 0$. This asymmetry arises from two effects, both related to the exponential density profiles of stellar disks. Stars near the Sun that have $\tilde{v}_\phi > 0$ have more angular momentum than the LSR and thus have guiding centers at $R_g > R_0$, while stars with $\tilde{v}_\phi < 0$ have guiding centers at $R_g < R_0$. Since the surface density of stars declines exponentially with R , there are more stars in the latter class than in the former, and the distribution in velocity space therefore extends further towards negative \tilde{v}_ϕ than in the opposite direction. The second effect is that the velocity dispersion σ_R declines with R , so the fraction of the stars that are based at $R_g = R_0 - \delta R$ on eccentric orbits that bring them to the Sun with $\tilde{v}_\phi < 0$ is larger than the fraction of the stars based at $R_0 + \delta R$ that can reach the Sun with $\tilde{v}_\phi > 0$. Similarly, there are more Japanese than Nepalese in Oxford in the summer,

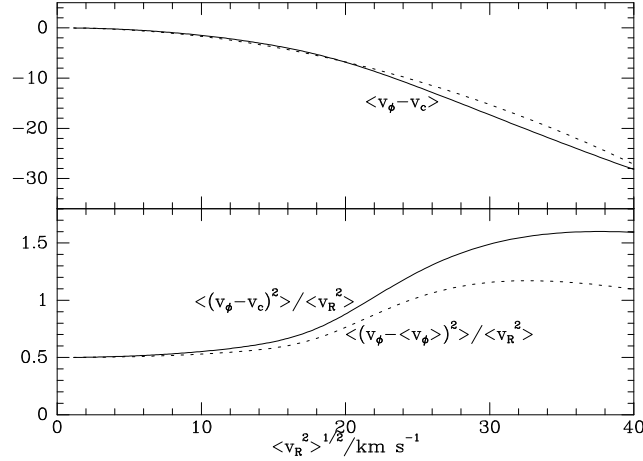


Figure 4.16 Upper panel: the mean value of \tilde{v}_ϕ as a function of $\overline{v_R^2}^{1/2}$ for azimuthal velocity distributions like those plotted in Figure 4.15. The dashed line shows a parabolic fit to the curve. Lower panel: for the same distributions the ratios $\overline{v_\phi^2}/\overline{v_R^2}$ (full) and $\overline{(v_\phi - \overline{v_\phi})^2}/\overline{v_R^2}$ (dashed).

both because the population of the Japan exceeds that of Nepal, and because Japanese have larger travel budgets than Nepalese.

Figure 4.16 shows the prediction of equation (4.156) for the dependence of various averages on the temperature $\overline{v_R^2}$ of the stellar population. The full curve in the upper panel shows that as $\overline{v_R^2}$ increases, the mean rotation rate of the population falls more and more below the circular speed. This phenomenon is called **asymmetric drift** (BM §10.3.1). In §4.8.2a we shall show that for any cool-disk DF we have to a good approximation that $\overline{v_\phi - v_c} \propto \overline{v_R^2}$. The dashed curve shows such a parabolic fit to the full curve. In the lower panel the full curve shows the ratio $\overline{(v_\phi - v_c)^2}/\overline{v_R^2}$, which is predicted by epicycle theory to be $\gamma^{-2} = 0.5$ (eq. 3.100). This ratio does start at 0.5, but by a radial dispersion of 10 km s^{-1} has risen by 10% and by a dispersion of 20 km s^{-1} , less than 10% of the circular speed, it has risen by nearly 80%. Thus the range of validity of equation (3.100) is surprisingly narrow. Problems 4.43 and 4.44 help to explain why.

Figure 4.17 compares the distribution of observed azimuthal velocities for F and G stars near the Sun from Nordström et al. (2004) with the prediction of the Schwarzschild DF for a population with the same value of $\overline{v_R^2}^{1/2} = 34 \text{ km s}^{-1}$. The model distribution exaggerates the skewness of the observed distribution and seriously overestimates the number of stars with $\tilde{v}_\phi \lesssim -60 \text{ km s}^{-1}$. These shortcomings reflect the breakdown of the epicycle approximation for stars on highly eccentric orbits. We can estimate the char-

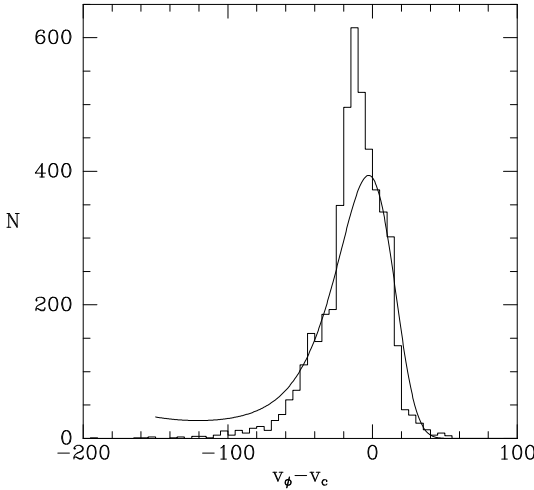


Figure 4.17 The distribution of v_ϕ components of 4787 F and G stars that have space velocities in Nordström et al. (2004). Stars with a high probability of having variable radial velocities are excluded. The smooth curve shows the distribution predicted by the Schwarzschild DF for a population with the same value of $\frac{\sigma_R^2}{v_R} = 34 \text{ km s}^{-1}$.

acteristic epicycle amplitude x of the stars in the Nordström et al. sample by equating the sample's radial velocity dispersion to the RMS radial velocity, averaged over time, of an individual star with epicycle amplitude x . Thus x follows from $\kappa_0 x / \sqrt{2} = 34 \text{ km s}^{-1}$, where $\kappa_0 = 37 \text{ km s}^{-1} \text{ kpc}^{-1}$ is the local epicycle frequency (Table 1.2), and we deduce $x \simeq 1.3 \text{ kpc}$. Thus sample stars typically cover a radial range of 2.6 kpc, which is as large as the disk's scale length R_d , and stars observed at $\tilde{v}_\phi < -60 \text{ km s}^{-1}$ have even larger epicycle diameters. For such large excursions the effective potential is not accurately harmonic, as the epicycle approximation requires. Much more satisfactory fits to the data can be obtained if one uses Shu's DF (4.147) upon which the Schwarzschild DF is based (Binney 1987; Kuijken & Tremaine 1991; Dehnen 1999b).

Figure 4.18 shows the density of F and G stars on four slices through velocity space, corresponding to $v_z = -30, -15, -5$ and 15 km s^{-1} . At $v_z = -30 \text{ km s}^{-1}$ the equi-density contours are reasonably elliptical, and qualitatively in agreement with the predictions of the Schwarzschild DF. At $v_z = -15 \text{ km s}^{-1}$ the third highest contour has a pronounced bulge around a local maximum at $(17, -43) \text{ km s}^{-1}$. This concentration, which cannot be due to noise because it involves a significant number of stars and is also visible in adjacent slices through velocity space, is qualitatively in conflict with the Schwarzschild DF. It is called the **Hercules star stream** (Famaey et al. 2005). A second local maximum is visible at $(30, -12) \text{ km s}^{-1}$; this is due to the **Hyades star stream** or **moving group**. The panel for $v_z = -5 \text{ km s}^{-1}$ shows a neighboring local maximum at $(6, -20) \text{ km s}^{-1}$ caused by the Pleiades star stream.

The differences between the structure of Figure 4.18 and that predicted by the Schwarzschild DF could have two explanations. The first possibility is that the stars are still dispersing from the associations in which they formed, with the result that the DF is still evolving towards a steady state, and does

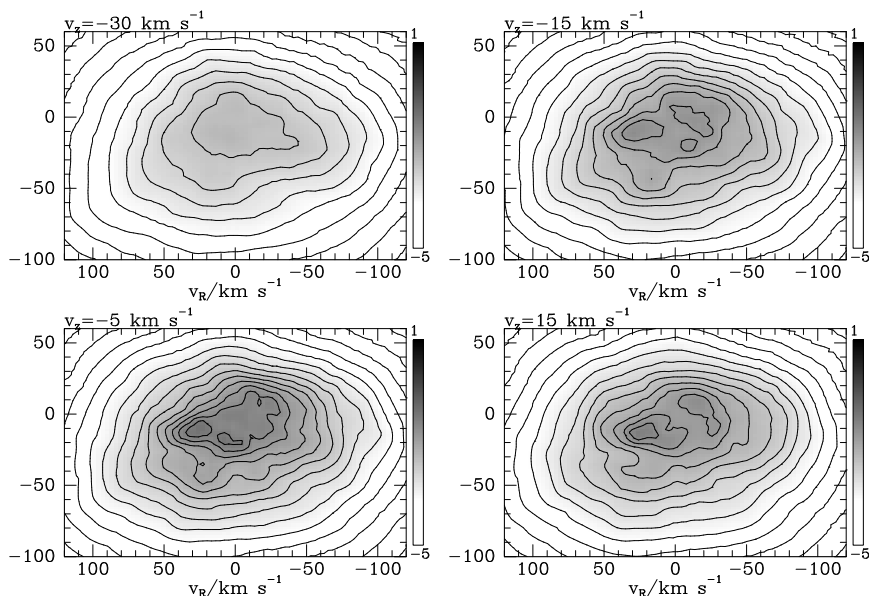


Figure 4.18 The density of solar-neighborhood stars in velocity space. Each panel is a slice through velocity space at the value of v_z given in the top left corner. The component \tilde{v}_ϕ is plotted vertically. The stellar density has been determined from the velocities of Nordström et al. (2004) for 4787 F and G stars using the FIEstAS algorithm of Ascasibar & Binney (2005). The velocities are relative to the LSR (page 12), so the Sun lies at $(v_R, \tilde{v}_\phi, v_z) = (-10, 5.2, 7.2)$.

not yet satisfy the Jeans theorem. However, from comparison of the spectra of stars to stellar models one can show that individual concentrations in Figure 4.18 contain stars of very different ages (Famaey et al. 2005), so the clumping cannot be due irregular star formation. The second possibility is that the structure of Figure 4.18 is due to stellar-dynamical processes. That is, the structure of Figure 4.18 implies that either (i) the Galactic potential is not axisymmetric, or (ii) it is not time-independent. In fact it is neither axisymmetric nor time-independent, because it has contributions from both the Galactic bar and spiral structure. Dehnen (2000a) and Fux (2001) have shown that the stars near $(v_R, \tilde{v}_\phi) \simeq (17, -43) \text{ km s}^{-1}$ are probably in resonance with the Galactic bar (BM §10.2). De Simone, Wu, & Tremaine (2004) show that transient spiral structure should generate horizontal striations in the panels of Figure 4.18, and at least in the bottom panels it is possible to imagine that such striations are present.

4.5 DFs for razor-thin disks

The majority of stars in a spiral galaxy lie in a thin disk. Thus models of disks of negligible thickness are both conceptually simple and directly applicable to real galaxies. These systems are in a sense the simplest ones in which the DF depends on the third integral I_3 in addition to H and L . Specifically, in a thin disk all stars have $I_3 = 0$ with the consequence that motion perpendicular to the galactic plane is prohibited and the system is perfectly planar. By reducing the model to a two-dimensional one, we can forget about the dependence of the DF on I_3 and write simply $f = f(H, L_z)$.

The problem of finding a DF that generates a disk with a given surface-density profile has much in common with the analogous problem for spherical systems because the surface-density distribution $\Sigma(R)$, like the density $\nu(r)$, is a function of only one variable, whereas the generic DF is a function of two integrals, $f(H, L_z)$ for a disk and $f(H, L)$ for a sphere. In particular, (i) there is a unique DF of the form $f(H)$ that generates a given surface-density profile $\Sigma(R)$; however, the associated model is of little physical interest because it does not rotate, whereas disk galaxies rotate rapidly. (ii) By analogy with the constant-anisotropy models described in §4.3.2b, we can posit that the DF is of the form $f = L_z^\alpha f_1(H)$ and obtain an integral equation for f_1 in terms of the given surface-density profile. (iii) By analogy with the approach of §4.3.2c we can assume that the DF depends on H and L_z only in some particular combination, such as $H - \Omega L_z$, where Ω is a parameter. Kalnajs (1976) discusses these approaches to the choice of DF. Here we simply present two of the most useful DFs that can be obtained in this way.

4.5.1 Mestel disk

In §2.6.1a we found that a disk with surface density

$$\Sigma(R) = \Sigma_0 \frac{R_0}{R} \quad (4.158)$$

has a circular speed v_c that is independent of radius and given by

$$v_c^2 = -R \frac{\partial \Psi}{\partial R} = 2\pi G \Sigma_0 R_0. \quad (4.159)$$

We set the arbitrary constant involved in the definition of the relative potential such that $\Psi(R_0) = 0$, and integrate equation (4.159) with respect to R , to find

$$\Psi(R) = -v_c^2 \ln(R/R_0). \quad (4.160)$$

Now following Toomre (1977a) consider the DF

$$f(\mathcal{E}, L_z) = \begin{cases} F(L_z/R_0 v_c)^q e^{\mathcal{E}/\sigma^2} & (L_z > 0) \\ 0 & (L_z \leq 0), \end{cases} \quad (4.161)$$

where F , q , and σ are all constants. Inserting equation (4.160) into equation (4.161) and integrating over all velocities in the plane, we find that the surface density produced by this DF in the potential (4.160) disk is

$$\begin{aligned}\Sigma'(R) &= F \left(\frac{R}{R_0 v_c} \right)^q \int_0^\infty dv_\phi v_\phi^q \int_{-\infty}^\infty dv_R \exp \left[-\frac{v_c^2}{\sigma^2} \ln \left(\frac{R}{R_0} \right) - \frac{v_R^2 + v_\phi^2}{2\sigma^2} \right] \\ &= F \left(\frac{R}{R_0 v_c} \right)^q \left(\frac{R}{R_0} \right)^{-v_c^2/\sigma^2} \int_0^\infty dv_\phi v_\phi^q e^{-v_\phi^2/2\sigma^2} \int_{-\infty}^\infty dv_R e^{-v_R^2/2\sigma^2} \\ &= 2^{q/2} \sqrt{\pi} \left(\frac{1}{2}q - \frac{1}{2} \right)! \left(\frac{R\sigma}{R_0 v_c} \right)^q \left(\frac{R}{R_0} \right)^{-v_c^2/\sigma^2} F \sigma^2.\end{aligned}\tag{4.162}$$

Comparing equations (4.158) and (4.162), we see that the DF of equation (4.161) will self-consistently generate the Mestel disk if we set

$$q = \frac{v_c^2}{\sigma^2} - 1 \quad \text{and} \quad F = \frac{\Sigma_0 v_c^q}{2^{q/2} \sqrt{\pi} \left(\frac{1}{2}q - \frac{1}{2} \right)! \sigma^{q+2}}.\tag{4.163}$$

The parameter q that appears in the DF (4.161) of the Mestel disk is a measure of the degree to which the disk is centrifugally supported: from equation (4.161) one may show that σ is the velocity dispersion \bar{v}_R^2 in the radial direction. The mean azimuthal velocity is

$$\bar{v}_\phi = \frac{\int d^2\mathbf{v} v_\phi f(\mathcal{E}, L_z)}{\int d^2\mathbf{v} f(\mathcal{E}, L_z)} = \frac{\int dv_\phi v_\phi^{q+1} e^{-v_\phi^2/2\sigma^2}}{\int dv_\phi v_\phi^q e^{-v_\phi^2/2\sigma^2}} = \frac{\sqrt{2} \left(\frac{1}{2}q \right)!}{\left(\frac{1}{2}q - \frac{1}{2} \right)!} \sigma.\tag{4.164}$$

For large q , $\bar{v}_\phi/\sigma = \sqrt{q}[1 + O(q^{-1})]$, all stars are on circular orbits, and $\bar{v}_\phi = v_c$.

4.5.2 Kalnajs disks

From equation (2.128) and Table 2.1 we have that the potential $\Phi(R)$ at radius R in the equatorial plane of a homogeneous oblate spheroid with eccentricity e , density ρ , and semi-axes of length a and $a_3 = a\sqrt{1-e^2}$ is

$$\Phi(R) = \frac{\pi G \rho a_3}{ae^2} \left(\frac{\sin^{-1} e}{e} - \sqrt{1-e^2} \right) R^2 + \text{constant}.\tag{4.165}$$

If we now flatten this spheroid down to a disk by letting $e \rightarrow 1$, while holding the central surface density $\Sigma_c \equiv 2\rho a_3$ constant, we obtain that

$$\Phi(R) = \frac{\pi^2 G \Sigma_c}{4a} R^2 + \text{constant} = \frac{1}{2} \Omega_0^2 R^2 + \text{constant},\tag{4.166}$$

where $\Omega_0 \equiv \sqrt{\frac{1}{2}\pi^2 G \Sigma_c / a}$ is the angular speed of a circular orbit. By equation (2.145) the surface density of our disk is

$$\Sigma(R) = \Sigma_c \sqrt{1 - \frac{R^2}{a^2}}. \quad (4.167)$$

Now consider the density distribution that arises from the DF¹³

$$f(\mathcal{E}, L_z) = \begin{cases} F [(\Omega_0^2 - \Omega^2)a^2 + 2(\mathcal{E} + \Omega L_z)]^{-1/2} & \text{for } [\dots] > 0, \\ 0 & \text{for } [\dots] \leq 0. \end{cases} \quad (4.168)$$

Since $\mathcal{E} = \Psi - \frac{1}{2}(v_\phi^2 + v_R^2)$ and $L_z = Rv_\phi$, we can also write the argument of the radical in equation (4.168) as

$$(\Omega_0^2 - \Omega^2)a^2 - (v_\phi - \Omega R)^2 - v_R^2 + 2\Psi + \Omega^2 R^2.$$

Hence at any radius R the DF (4.168) depends on the velocities only in the combination $v_R^2 + (v_\phi - \Omega R)^2$. Consequently, the distribution of azimuthal velocities in a model generated by this DF is symmetrical about $v_\phi = \Omega R$, which is therefore the mean azimuthal velocity at R . We choose the arbitrary constant involved in the definition of the relative potential such that

$$\Psi(R) = -\Phi(R) + \text{constant} = -\frac{1}{2}\Omega_0^2 R^2. \quad (4.169)$$

Substituting this form of Ψ into equation (4.168) and integrating over all velocities, we find the surface density $\Sigma'(R)$ generated by this DF in the potential of our disk to be

$$\Sigma'(R) = F \int_{v_{\phi 1}}^{v_{\phi 2}} dv_\phi \int_{v_{R1}}^{v_{R2}} \frac{dv_R}{\sqrt{(\Omega_0^2 - \Omega^2)(a^2 - R^2) - (v_\phi - \Omega R)^2 - v_R^2}}. \quad (4.170)$$

The limits v_{R1} , v_{R2} of the inner integral in equation (4.170) are just the values of v_R for which the integrand's denominator vanishes. Hence equation (4.170) is of the form

$$\Sigma'(R) = F \int_{v_{\phi 1}}^{v_{\phi 2}} dv_\phi \int_{-b}^b \frac{dv_R}{\sqrt{b^2 - v_R^2}} = \pi F \int_{v_{\phi 1}}^{v_{\phi 2}} dv_\phi = \pi F (v_{\phi 1} - v_{\phi 2}). \quad (4.171)$$

But $v_{\phi 1}$ and $v_{\phi 2}$ are just the roots of the quadratic equation

$$b^2 = (\Omega_0^2 - \Omega^2)(a^2 - R^2) - (v_\phi - \Omega R)^2 = 0, \quad (4.172)$$

¹³ Note that $-(\mathcal{E} + \Omega L_z)$ is the Hamiltonian in the frame that rotates at frequency Ω (eq. 3.112).

Box 4.2: Freeman's analytic bars

The **Freeman bars** are razor-thin elliptical disks that are stationary in a frame that rotates at angular speed Ω_b (Freeman 1966). In this frame, the outer boundary of the disk is elliptical, $x^2/a^2 + y^2/b^2 = 1$, and the surface density and potential are given by

$$\Phi(x, y) = \frac{1}{2}(\Omega_x^2 x^2 + \Omega_y^2 y^2) \quad ; \quad \Sigma(x, y) = \Sigma_c \sqrt{1 - x^2/a^2 - y^2/b^2}.$$

The family of Freeman bars is very rich, including bars with all possible axis ratios b/a and all pattern speeds such that $|\Omega_b| < \min(\Omega_x, \Omega_y)$. In the limit $b/a \rightarrow 1$ the Freeman bars reduce to the Kalnajs disks.

The Freeman bars provide the only known analytic models of bars. Regrettably, they have several unrealistic features. In particular, because the potential is quadratic in the coordinates, the equations of motion are linear, and every trajectory can be regarded as the superposition of motion in two ellipses. In fact these ellipses are simply the prograde and retrograde ellipses that surround the central Lagrange point, as described in equation (3.125) and Figure 3.15. These orbits can have quite different properties from the orbits in more realistic rotating potentials.

so

$$\Sigma'(R) = 2\pi F a \sqrt{\Omega_0^2 - \Omega^2} \sqrt{1 - \frac{R^2}{a^2}}. \quad (4.173)$$

Comparing equations (4.167) and (4.173) we see that if we set

$$F = \frac{\Sigma_c}{2\pi a \sqrt{\Omega_0^2 - \Omega^2}}, \quad (4.174)$$

then $\Sigma'(R) = \Sigma(R)$, so we have derived a self-consistent—though rather artificial—stellar-dynamical model of a flat disk galaxy, called a **Kalnajs disk**.

It is straightforward to verify that the mean angular speed Ω of the stars in a Kalnajs disk is independent of position, and relative to this mean speed the stars have isotropic velocity dispersion in the disk plane,

$$\overline{v_x^2} = \overline{v_y^2} = \frac{1}{3}a^2(\Omega_0^2 - \Omega^2)(1 - R^2/a^2). \quad (4.175)$$

Thus Kalnajs disks range from hot systems with $\Omega \ll \Omega_0$, in which the support against self-gravity comes from random motions, to cold systems with $\Omega \approx \Omega_0$, in which all stars move on nearly circular orbits and the random velocities are small.

4.6 Using actions as arguments of the DF

Hitherto we have focused on the use of the Hamiltonian H and various components of the angular momentum \mathbf{L} as the arguments of the DF. In §3.5 we saw that actions are constants of motion that describe orbits in integrable potentials with remarkable simplicity, and in this section we consider the advantages of using actions as arguments of the DF.

Several features make the actions the most convenient arguments for the DF:

- (i) By the Jeans theorem, the arguments of the DF for a steady-state galaxy should be isolating integrals. The space spanned by the integrals is called **integral space**. In a spherical potential H , L and L_z are isolating integrals and can serve as coordinates of integral space. However, in these coordinates the boundary of integral space does not have an analytic form. For example, the allowed values of L lie in the interval $[0, L_c(H)]$, where $L_c(H)$ is the angular momentum of a circular orbit of energy H . In contrast, the boundaries of action space are simple. For example if the actions are chosen to be (J_ϕ, J_θ, J_r) (Table 3.1), then the allowed region is the quadrant $J_\theta, J_r \geq 0$.
- (ii) The volume of space that is occupied by orbits with actions in $d^3\mathbf{J}$ is $(2\pi)^3 d^3\mathbf{J}$, whereas the phase-space volume associated with orbits in some range of H and L depends on the potential (eq. 4.288). Consequently, when we use actions as Cartesian coordinates for integral space, the density of stars in integral space is simply $(2\pi)^3 N$ times the DF, where N is the number of stars in the system.
- (iii) When the coordinates are actions, the DF is invariant under slow changes in the potential—see §4.6.1.

Let us examine more closely an action space for a spherical potential using the coordinates (J_ϕ, J_θ, J_r) (Figure 3.25). In galactic potentials the surfaces of constant H are each made up of two approximately planar triangles—in Kepler and harmonic potentials they are exactly so. At the point \mathbf{J} , the normal to the local surface of constant H is the vector $\boldsymbol{\Omega}$ whose components are the three characteristic frequencies of the orbit \mathbf{J} (eq. 3.190). In the simplest equilibrium models, which have ergodic DFs of the form $f(H)$, the density of stars in integral space is constant on these triangular surfaces.

These relations suggest a procedure for constructing the action-space DF for a galaxy with any desired triaxial density distribution (Binney 1987). We start by imagining a spherical galaxy that has the same “average” radial density profile as the galaxy in question, for example by computing the density along a line that makes equal angles with its three principal axes. We find the potential $\Phi(r)$ of this spherical galaxy, and then use Eddington’s formulae (4.46) to find the DF $f_0(H)$ of this system.

Now consider the DF

$$f(\mathbf{J}) = s(\mathbf{J})f_0(H), \quad (4.176)$$

where \mathbf{J} is the action vector in the potential $\Phi(r)$, and the **shift function** $s(\mathbf{J})$ shifts stars in action space but leaves their energy E unchanged. We now choose the shift function, following the precepts outlined above, to generate a DF that reproduces the axis ratios of the desired galaxy when confined by the spherical potential $\Phi(r)$.

Shifts at constant energy have relatively little effect on the radial density profile, so this process will also produce a galaxy with a radial distribution of stars that resembles the desired galaxy. Finally, we slowly change the gravitational potential from $\Phi(r)$ to the triaxial potential that is required for self-consistency, as determined by solving Poisson's equation. Since actions are invariant under slow changes in the potential, the DF will not be affected by this process. It is straightforward in principle to iterate this process by adjusting the shift function s and the ergodic DF $f_0(H)$ to bring the product of this process closer and closer to the desired set of properties.

The property required for a shift function—that it leave the energy distribution invariant—is simple to describe. The number of stars per unit energy is (cf. eq. 4.81)

$$N(E) = (2\pi)^3 \int d^3\mathbf{J} \delta[E - H(\mathbf{J})] f(\mathbf{J}). \quad (4.177)$$

The shift function in equation (4.176) will leave this distribution unchanged if

$$\int d^3\mathbf{J} \delta[E - H(\mathbf{J})] s(\mathbf{J}) = \int d^3\mathbf{J} \delta[E - H(\mathbf{J})]. \quad (4.178)$$

This can be rewritten in a more transparent form by assuming that the shift function depends on two of the actions—say, J_ϑ and J_ϕ —and the Hamiltonian H , thus eliminating the radial action. Since $\partial J_r / \partial H = \Omega_r^{-1}$, we have

$$\int \frac{dJ_\vartheta dJ_\phi}{\Omega_r} s(J_\vartheta, J_\phi, H) = \int \frac{dJ_\vartheta dJ_\phi}{\Omega_r}. \quad (4.179)$$

If we shift stars over each energy surface towards the J_r axis, we obtain a radially biased model, while pushing them away from this axis generates a tangentially biased model. If we push stars towards the J_ϕ axis, we flatten the system; in the extreme case in which all stars have been pushed onto the J_φ axis, the DF is $f(\mathbf{J}) = f_0(J_\varphi) \delta(J_r) \delta(J_\vartheta)$ and the system has become a razor-thin disk whose surface density $\Sigma(R)$ is determined by the function $f_0(J_\varphi)$. If f_0 is an even function, at each radius there will be equal numbers of stars orbiting in each sense around the z axis. If $f_0(J_\varphi) = 0$ for $J_\varphi < 0$, all stars will orbit in the same sense and the disk will be cold. We can heat this disk up while leaving it razor-thin by replacing $\delta(J_r)$ with a function such as a steep exponential. Similarly, we can give the disk finite thickness by replacing $\delta(J_\vartheta)$.

These examples show that in a spherical potential there is a transparent connection between the distribution of stars in action space and the shape and kinematics of the stellar distribution. Even in the case of a triaxial integrable potential, actions can be defined that are closely related to $(J_\phi, J_\vartheta, J_r)$ —see §3.8.1—and the relations we have described between the shape of the stellar system and the distribution of stars in action space continue to hold.

4.6.1 Adiabatic compression

In §3.6 we saw that actions are constant during slow changes in the confining potential Φ . At that stage we were only equipped to study the evolution of individual orbits as Φ evolved. Now we can discuss the corresponding evolution of entire stellar systems.

The key idea is that constancy of individual actions implies constancy of the system's DF $f(\mathbf{J})$. This invariance by no means implies that the system's density and velocity distributions are invariant, but it enables us to calculate their evolution relatively simply. We describe the procedure in the case that the system is at all times spherically symmetric, and assume that the slow evolution is driven by an external potential $\Phi_{\text{ext}}(\mathbf{r}, t)$ (say, due to a growing black hole or infalling gas). However, the scheme applies with minor modifications to any system that always has an integrable potential.

We assume that at time $t = 0$, $\Phi_{\text{ext}}(\mathbf{r}, 0) = 0$, and that we know the initial DF $f_0(H, L)$ and the corresponding self-consistent potential $\Phi_0(r)$. Given the final external potential $\Phi_{\text{ext}}(r, t)$, we determine the density and velocity distributions of the system when the potential is $\Phi_{\text{tot}} = \Phi_{\text{ext}} + \Phi_{\text{f}}$, where $\Phi_{\text{f}}(r)$ is the final contribution of our stellar system to the potential.

We make a first guess Φ_1 at the form of Φ_{tot} . To improve this guess we need to evaluate the integral

$$\rho_1(r) \equiv M \int d^3\mathbf{v} f(\mathbf{J}), \quad (4.180)$$

where M is the system's stellar mass. So we need to determine f at a grid of values (r, \mathbf{v}) . To do so we first determine the actions of the orbit in Φ_1 that has these initial conditions. Since the potential is spherical, we can take one action to be the total angular momentum L , which is trivially determined from r and \mathbf{v} , and then the radial action is given by equation (3.224) with Φ_1 substituted for Φ . The third action is L_z although this is not needed since the DF of a spherical system is independent of L_z . Next we use equation (3.224) again to find the energy E' of an orbit in the original potential Φ_0 that has these same actions. Then by the invariance of $f(\mathbf{J})$ we must have that $f(\mathbf{J}) = f_0(E', L)$. Having determined $f(\mathbf{J})$ we can carry out the integral in (4.180).

Once we have determined $\rho_1(r)$, we solve Poisson's equation for the corresponding potential $\Phi_1'(r)$, and thus obtain a revised estimate $\Phi_2 = \Phi_{\text{ext}} + \Phi_1'$ of the final overall potential. Then we repeat the procedure just described to determine the density distribution $\rho_2(r)$ that $f(\mathbf{J})$ generates in Φ_2 , and we iterate until the difference between Φ_n and Φ_{n+1} becomes negligible.

(a) Cusp around a black hole In §1.1.6 we saw that most luminous galaxies have a massive black hole at their centers, and in §3.6.2d we investigated how individual stellar orbits evolve as the mass M of the black hole slowly increases. We can now investigate how the growth of the black hole modifies the density and velocity dispersions within the system. We present only approximate formulae in order to obtain analytic results. More detailed treatments that use the technique just described can be found in Young (1980), Quinlan, Hernquist, & Sigurdsson (1995) and van der Marel (1999).

As in §3.6.2d we assume that the black hole forms in the core of an approximately isothermal system, so for the tightly bound stars of interest the initial Hamiltonian can be approximated by equation (3.281), and the DF is initially

$$f = \frac{\rho_0}{(2\pi\sigma^2)^{3/2}} e^{-H/\sigma^2} = \frac{\rho_0}{(2\pi\sigma^2)^{3/2}} e^{-\Omega(2J_r+L)/\sigma^2}, \quad (4.181)$$

where Ω and σ are the circular frequency and velocity dispersion within the core and we have set $\Psi(r) = -\frac{1}{2}\Omega^2 r^2$. Sufficiently close to the final black hole, the potential will be Keplerian, so from equation (E.6) we have

$$\frac{1}{2}v^2 - \frac{GM}{r} = H_K = -\frac{1}{2} \left(\frac{GM}{J_r + L} \right)^2. \quad (4.182)$$

Using this equation and $L = rv_t$ to eliminate J_r and L from equation (4.181) in favor of the total and tangential speeds v and v_t , we find that f can be written

$$f = \frac{\rho_0}{(2\pi\sigma^2)^{3/2}} \exp \left(-\frac{x_m^2}{\sqrt{x_m^2 - x^2}} \right) e^{x \sin \psi}, \quad (4.183a)$$

where

$$x \equiv \frac{\Omega r}{\sigma^2} v \quad ; \quad x_m \equiv \frac{\Omega r}{\sigma} \sqrt{\frac{2GM}{r\sigma^2}} \quad ; \quad v_t = v \sin \psi, \quad (4.183b)$$

and ψ is the angle between the radius and velocity vectors. Multiplying through by $d^3\mathbf{v} = 2\pi v^2 dv \sin \psi d\psi$ and integrating we find that

$$\rho(r) = \frac{2\rho_0}{\sqrt{\pi}} \left(\frac{GM}{\sigma^2 r} \right)^{3/2} \int_0^1 dy y^2 \exp \left(-\frac{x_m}{\sqrt{1-y^2}} \right) \int_0^\pi d\psi \sin \psi e^{x_m y \sin \psi}, \quad (4.184)$$

where $y \equiv x/x_m = v\sqrt{r/2GM}$. As we approach the center, $x_m \rightarrow 0$ and the inner and outer integrals in equation (4.184) tend to 2 and $\frac{1}{3}$, respectively, so the black hole distorts the original homogeneous stellar density into a cusp in which $\rho \propto r^{-3/2}$.

Multiplying f by v_r^2 or v_t^2 before integrating over all velocities, we obtain the ratio of velocity dispersions

$$\frac{\overline{v_r^2}}{\overline{v_t^2}} = \frac{\int_0^1 dy y^4 \exp(-x_m/\sqrt{1-y^2}) \int_0^\pi d\psi \sin \psi \cos^2 \psi e^{-x_m y \sin \psi}}{\int_0^1 dy y^4 \exp(-x_m/\sqrt{1-y^2}) \int_0^\pi d\psi \sin^3 \psi e^{-x_m y \sin \psi}}. \quad (4.185)$$

In the limit $x_m \rightarrow 0$, the integral over ψ on top tends to $\frac{2}{3}$, while that on the bottom tends to $\frac{4}{3}$, so $\overline{v_r^2}/\overline{v_t^2}$ tends to $\frac{1}{2}$, which implies an isotropic velocity distribution at small radii. This result is surprising, given that in §3.6.2d we showed that the mean value of v_t^2 grows more than does the mean value of v_r^2 for individual stars. The resolution of this apparent inconsistency is that eccentric orbits are pulled in towards the black hole more than circular orbits, thereby boosting the radial velocity dispersion at a given radius. Correspondingly, at larger radii the velocity distribution becomes strongly tangentially biased (Goodman & Binney 1984). The distribution must return to isotropy far from the black hole, but at such radii we cannot approximate the Hamiltonian with the Kepler Hamiltonian, and equation (4.185) is no longer valid.

(b) Adiabatic deformation of dark matter Simulations of the cosmological clustering of collisionless matter (§9.3) indicate that the radial density profiles of the structures that form in this way can be approximated by the NFW model (§2.2.2g). It is plausible that after these structures have formed, dissipation of energy caused many of their baryons to move inwards on a timescale that was long compared to the system's crossing time. The procedure we have described can be used to find how this infall modified the initial NFW halo.

Combining the photometry of the galaxy with a population-synthesis model, we can determine the mass density $\rho_{\text{ext}}(r)$ of the baryons in the present galaxy, and thus evaluate the potential $\Phi_{\text{ext}}(r)$ that drives the adiabatic deformation of the dark-matter halo. Hence, we can predict the current distribution of dark matter in luminous galaxies given the properties of halos in simulations that do not contain baryons (Sellwood & McGaugh 2005).

We take f_0 to be a DF that self-consistently generates the NFW profile of scale length 20 kpc, which is the estimated scale length of the Milky Way's dark-matter halo (Klypin, Zhao, & Somerville 2002). We consider three models: isotropic ($\beta = 0$), radially biased ($\beta = 0.5$) and tangentially biased ($\beta = -2$). Within these halos we slowly grow a spherical object that has the same radial density profile as the bulge and disk of the Galaxy—we approximate the disk as spherical for the sake of simplicity, recognizing that the growth of a flat disk will flatten the halo. We take the final mass of

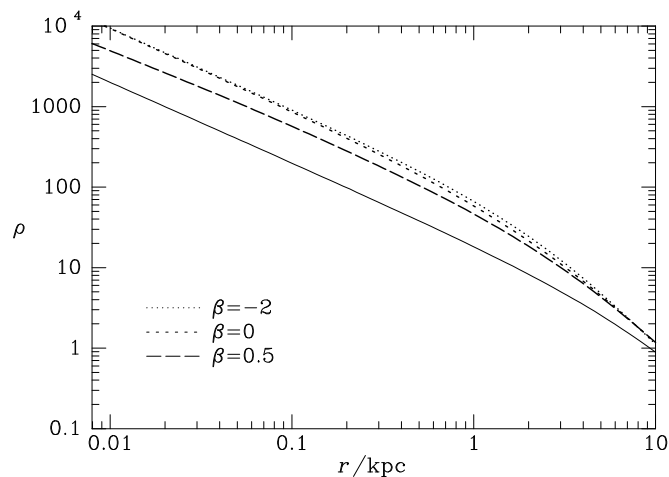


Figure 4.19 The full curve shows the original density profile of an NFW model with scale radius $a = 20$ kpc. The broken curves show the density profile of this system after a spherical representation of the Milky Way has grown slowly within it. The mass of this object is equal to the mass of the original NFW halo within 10 kpc. From data supplied by J. Magorrian.

the bulge and disk interior to 10 kpc to be equal to the original halo mass within this radius. Figure 4.19 shows the effect of this added mass on the dark halo. The full curve shows the halo's original density profile, while the short-dashed and dotted curves show the final density profile in the case of initially isotropic or tangentially biased halos. At 100 pc from the center, the introduction of the galaxy has increased the densities of these halos by a factor 4.25. The long-dashed curve, which is for the radially biased halo, shows that the density of this halo increases less; by a factor 2.6. In §9.4c we will discuss the relevance of these results for the theory of galaxy formation.

4.7 Particle-based and orbit-based models

So far we have built models of stellar systems by the explicit construction of a DF that depends on the integrals of motion. Such methods are mostly applicable when the system has a high degree of symmetry. Systems of lower symmetry—triaxial or time-dependent systems, or axisymmetric systems that depend on the third integral—are harder to model in this way because their DFs have no analytic expression. In this section we discuss techniques that can be used to obtain models of such systems. These techniques are powerful in that the class of systems to which they can be applied is very wide. Unfortunately, they rely on discrete samples of a probability

distribution with the result that dynamically interesting quantities can be obscured by discreteness noise.

4.7.1 N-body modeling

N-body modeling is one of the most flexible and widely used techniques for exploring the behavior of collisionless stellar systems. In §2.9 we discussed methods used to determine the forces on particles in such simulations, while in §3.4 we discussed algorithms for following the particle orbits under the influence of these forces. Here we explore the foundations of the N-body approach, which are not as simple as they may seem at first sight. We shall again denote by \mathbf{w} the location (\mathbf{x}, \mathbf{v}) of a general phase-space point.

The collisionless Boltzmann equation (4.10) is the governing equation of a collisionless system, and an N-body calculation is a device for numerically solving this partial differential equation in seven independent variables for $f(\mathbf{w}, t)$ given the initial DF $f(\mathbf{w}, 0)$. The equation states that the DF is constant along the single-particle trajectories of the Hamiltonian $H = \frac{1}{2}v^2 + \Phi(\mathbf{x}, t)$, so

$$f(\mathbf{w}, t) = f(\mathbf{w}_0, 0), \quad (4.186)$$

where $\mathbf{w}_0 \equiv \mathbf{w}(t = 0)$. Hence the value of f at any point \mathbf{w} can be determined from the initial values of f once trajectories are known.

To determine the trajectories we need the potential

$$\Phi(\mathbf{x}, t) = -GM \int d^6\mathbf{w}' \frac{f(\mathbf{w}', t)}{|\mathbf{x} - \mathbf{x}'|}, \quad (4.187)$$

where M is the mass of the system and we have used equation (4.18).

The central feature of N-body modeling of collisionless systems is that we evaluate the six-dimensional integral in (4.187) by Monte-Carlo sampling (e.g., Press et al. 1986). For any reasonable function $g(\mathbf{w})$ we have

$$\int d^6\mathbf{w} g(\mathbf{w}) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N g(\mathbf{w}_i) / f_s(\mathbf{w}_i), \quad (4.188)$$

where the points \mathbf{w}_i are randomly chosen by sampling the probability density $f_s(\mathbf{w})$, which can be any function that satisfies

$$f_s(\mathbf{w}) \geq 0 \quad ; \quad \int d^6\mathbf{w} f_s(\mathbf{w}) = 1. \quad (4.189)$$

Applying (4.188) to (4.187) we have for sufficiently large N

$$\Phi(\mathbf{x}, t) \simeq -\frac{GM}{N} \sum_{i=1}^N \frac{f(\mathbf{w}_i, t) / f_s(\mathbf{w}_i, t)}{|\mathbf{x} - \mathbf{x}_i|}, \quad (4.190)$$

where now we have allowed the sampling density to be time-dependent. The expression on the right of this equation is the gravitational potential generated by particles that have masses

$$m_i(t) = \frac{M}{N} \frac{f(\mathbf{w}_i, t)}{f_s(\mathbf{w}_i, t)} \quad (i = 1, \dots, N). \quad (4.191)$$

Thus the potential becomes

$$\Phi(\mathbf{x}, t) \simeq -G \sum_{i=1}^N \frac{m_i(t)}{|\mathbf{x} - \mathbf{x}_i(t)|}. \quad (4.192)$$

We can choose $f_s(\mathbf{w}, 0)$ to be any convenient function and draw a set of sampling points \mathbf{w}_i from this distribution. The corresponding initial masses $m_i(0)$ can be evaluated from equation (4.191) because we are given the functional form of the DF $f(\mathbf{w}, 0)$ as an initial condition.

We have attached a mass $m_i(0)$ to each sampling point. Now suppose we treat these masses as real particles, and advance them along the trajectories that are determined by H . Then by the collisionless Boltzmann equation, the DF at the locations of these masses will be time-independent (eq. 4.10) and we will have solved for the evolved DF $f(\mathbf{w}, t)$ at the location of each mass.

To continue advancing the particles, we need repeatedly to solve Poisson's equation, which we do by Monte-Carlo sampling. We can use equation (4.191) to evaluate the relevant masses only where we know the value of $f(\mathbf{w}, t)$, that is at the locations of the masses. Thus the scheme will be feasible only if the sampling points follow the same trajectories as the masses. We need to know what density $f_s(\mathbf{w}, t)$ these points are sampling. One such distribution is simple to obtain: we evolve the initial sampling density $f_s(\mathbf{w}, 0)$ with the collisionless Boltzmann equation for the same Hamiltonian that we used for the mass particles. If we do this, the sampling density associated with each sampling point will be independent of time because the Hamiltonian flow preserves phase-space volume. Since both $f[\mathbf{w}_i(t), t]$ and $f_s[\mathbf{w}_i(t), t]$ are time-independent, by equation (4.191) the masses m_i are also time-independent. The simplest procedure is to choose $f_s(\mathbf{w}, 0) = f(\mathbf{w}, 0)$ so that all particles have equal masses, but this is neither necessary nor always desirable.

Since with this procedure all particles have constant masses, and follow the same trajectories that real stars would, it is tempting to imagine that the sampling points are real stars, or at least groups of stars. But a more general interpretation is that one is integrating the partial differential equation (4.6) by the method of characteristics (e.g., Whitham 1974) and evaluating the integral in equation (2.3) by Monte-Carlo sampling. Variants of this basic

N-body techniques exist in which $f_s \neq f$ and the masses of particles vary in time (Leeuwin, Combes, & Binney 1993; Syer & Tremaine 1996).

(a) Softening Actually Monte-Carlo sampling is not well adapted to the integrand of equation (4.187), because the singularity in the integrand at $\mathbf{x} = \mathbf{x}'$ causes estimates of the force to have an inconveniently large scatter—occasionally a sampling point \mathbf{x}' will fall close to \mathbf{x} . For this reason N-body Poisson solvers generally eliminate the singularity by replacing $|\mathbf{x} - \mathbf{x}'|$ by a softening kernel $S(|\mathbf{x} - \mathbf{x}'|)$ (§2.9.1). Fundamentally, softening is a stratagem designed to increase the statistical accuracy of our numerical estimate of the potential $\Phi(\mathbf{x})$ at the cost of some systematic error. It has a convenient side-effect, however: it reduces the magnitude of the largest accelerations experienced by particles, which makes it possible to use longer timesteps when integrating the particles' equations of motion, and thus to reduce the computational cost.

(b) Instability and chaos In §3.7.3 we showed that orbits in some potentials are regular, while others are chaotic, in the sense that any small initial change $\delta\mathbf{w}$ in the phase-space coordinates eventually grows exponentially fast, $|\delta\mathbf{w}(t)| \approx |\delta\mathbf{w}(0)| \exp(t/t_L)$, where t_L is the Liapunov time. With a given particle number N , an N-body simulation becomes a Hamiltonian system with $3N$ degrees of freedom, and it is natural to ask whether solutions to the equations of motion of this complex system are also chaotic.

Following Goodman, Heggie & Hut (1993) we make a crude estimate of the Liapunov time in the absence of softening, under the assumption that all the simulation's particles have the same mass m . Suppose that soon after the simulation commences, particles 1 and 2 encounter one another at relative velocity v and impact parameter b . The resulting velocity impulse is given by equation (1.30) as $\Delta v \approx Gm/(bv)$. We now make an infinitesimal change in the initial position of star 1 that causes the impact parameter to change by δb . The resulting change in the velocity impulse is $\delta(\Delta v) \approx Gm \delta b/(b^2 v)$. Star 1 has a second encounter, this time with star 3, a time τ after its first encounter. If τ is short enough, the change in impact parameter b' in the second encounter will still be of order δb , but for large τ the change will be dominated by the drift in position caused by the velocity change from the first encounter. This drift is roughly $\delta(\Delta v)\tau$, so we may write

$$\delta b' \approx \left(1 + \frac{Gm\tau}{b^2 v}\right) \delta b; \quad (4.193)$$

the exact relation between $\delta b'$ and δb depends on the three-dimensional geometry of the encounters in a complex way, but this schematic expression is good enough for our purposes.

The magnification $\delta b'/\delta b$ is related to the Liapunov time by

$$\frac{\delta b'}{\delta b} \approx e^{\tau/t_L}; \quad (4.194)$$

thus

$$\frac{\tau}{t_L} \approx \ln \left(1 + \frac{Gm\tau}{b^2v} \right). \quad (4.195a)$$

For a given star, the time between encounters with impact parameter b is $\tau \approx 1/(nb^2v)$, where n is the number density of stars. Using this result to eliminate b^2v , we find that

$$t_L \approx \frac{\tau}{\ln(1 + Gmn\tau^2)} \approx \frac{\tau}{\ln(1 + \tau^2/t_{\text{cross}}^2)}, \quad (4.195b)$$

where we have used the definition (2.40) of the crossing time and the relation $\rho = nm$.

Equation (4.195) gives the Liapunov time for encounters with impact parameter $\sim b$. If all encounters had this impact parameter, this would be the Liapunov time of the orbit. In practice, the exponential divergence is dominated by those impact parameters that cause the most rapid exponential divergence, so the Liapunov time of the orbit is given by the minimum of equation (4.195) as a function of b , and therefore of τ . We have shown that $t_L/t_{\text{cross}} \simeq x/\ln(1+x^2)$, where $x = \tau/t_{\text{cross}}$. The function $x/\ln(1+x^2)$ has a broad minimum centered on $x = 2.0$, where it equals 1.2. Thus our crude calculation predicts that

$$t_L \approx t_{\text{cross}}. \quad (4.196)$$

The dominant encounters are those for which $R/v \approx t_{\text{cross}} \approx \tau \approx (nb^2v)^{-1}$. Hence the impact parameters of these dominant encounters satisfy $b \approx (nR)^{-1/2} \approx R/N^{1/2}$, which is much smaller than the typical interparticle separation $R/N^{1/3}$ for $N \gg 1$.

This calculation suggests that in the absence of softening *the Liapunov time in a stellar system is of the order of the crossing time, regardless of the total particle number N* . This phenomenon was discovered by Miller (1964) and is known as **Miller's instability**. Miller's instability is surprising, because in the limit $N \rightarrow \infty$, an N -body system should become collisionless, and its particles should orbit in a smooth potential: thus, if for example the smooth potential were spherical, all orbits would be regular and therefore have infinite Liapunov time, in conflict with our finding. More accurate analytic results (Goodman, Heggie & Hut 1993), verified by N -body simulations (Hemsendorf & Merritt 2002), yield an even more surprising result: t_L/t_{cross} actually *declines* slowly as N increases, so we must somehow reconcile the short Liapunov times of N -body simulations with our understanding of the nature of a collisionless system.

The resolution of this apparent paradox is that the Liapunov time describes the growth of infinitesimal perturbations to an orbit; it applies only so long as the perturbation is much smaller than the distance between stars (Valluri & Merritt 2000; Hut & Heggie 2002). Thus if we follow the motion of two particles separated by some small amount $\Delta\mathbf{x}_0$, their separation will

initially grow exponentially, $|\Delta\mathbf{x}| \approx |\Delta\mathbf{x}_0| \exp(t/t_L)$. However, this conclusion is valid only until $|\Delta\mathbf{x}|$ becomes comparable to the impact parameters of the dominant encounters, $b \sim R/N^{1/2}$. Beyond this point $|\Delta\mathbf{x}|$ will continue to grow, but only at a slower rate. In particular, it is reasonable to expect that only encounters with impact parameter $b \gtrsim |\Delta\mathbf{x}|$ will continue to contribute to the exponential growth of $\Delta\mathbf{x}$, while encounters with $b \lesssim |\Delta\mathbf{x}|$ will perturb the two stars independently. Encounters with $b \gtrsim R/N^{1/2}$ have $\tau/t_{\text{cross}} \approx (v/R)/(nb^2v) \approx (R/b)^2/N \lesssim 1$, so equation (4.195b) becomes $t_L \approx t_{\text{cross}}^2/\tau$ and the rate of divergence becomes

$$\frac{d|\Delta\mathbf{x}|}{dt} \approx \frac{|\Delta\mathbf{x}|}{t_L} \approx \frac{|\Delta\mathbf{x}|\tau}{t_{\text{cross}}^2} \approx \frac{|\Delta\mathbf{x}|}{t_{\text{cross}}} \frac{R^2}{b^2 N}. \quad (4.197)$$

The divergence is dominated by the encounters with the smallest impact parameters larger than the separation, so we set $b \approx |\Delta\mathbf{x}|$ to obtain

$$\frac{d|\Delta\mathbf{x}|}{dt} \approx \frac{R^2}{N t_{\text{cross}} |\Delta\mathbf{x}|}. \quad (4.198)$$

Integrating, we find

$$|\Delta\mathbf{x}|^2 \approx \frac{R^2 t}{N t_{\text{cross}}} \approx \frac{R^2 t}{t_{\text{relax}}}, \quad (4.199)$$

where we have written the relaxation time as $t_{\text{relax}} \approx N t_{\text{cross}}$ by dropping the Coulomb logarithm from equation (1.38). Thus the exponential growth of $|\Delta\mathbf{x}|$ that occurs when $|\Delta\mathbf{x}|$ is infinitesimal has been replaced by a much slower growth $\propto t^{1/2}$, such that $\Delta\mathbf{x}$ grows to of order the system size R in a relaxation time—a result we could have anticipated from the definition of the relaxation time.

Miller's instability is significantly weakened by softening, so long as the softening length is larger than the impact parameter of the dominant encounters, $b \simeq R/N^{1/2}$. For large N this is much smaller than the mean interparticle separation $R/N^{1/3}$. There is little to be gained by using a softening length ϵ that is much smaller than the interparticle separation, so for large N we can weaken the instability without compromising the resolution of the simulation by setting $R/N^{1/3} > \epsilon > R/N^{1/2}$.

Miller's instability raises a fundamental question about N-body simulations. The short Liapunov time implies that small errors in an N-body simulation are rapidly amplified. For example, consider a simulation of a stellar system with $N = 10^3$, for which the dominant encounters have impact parameters of order three percent of the system size R . If roundoff leads to a positional error of one bit in a numerical calculation with 16 decimal digits, this error will grow exponentially, until it is of order $0.03R$ in a time $\ln(0.03 \times 10^{16})t_L \approx \ln(0.03 \times 10^{16})t_{\text{cross}} \simeq 33t_{\text{cross}}$. In practice, the error growth is even more rapid, since the errors arising from the non-zero

timestep of the integration algorithm (§3.4) are usually much larger than roundoff error. Given this, why should we believe in the results of N-body simulations at all? We usually test a numerical algorithm by repeating it with a smaller timestep or other accuracy parameter until its results converge. On account of Miller’s instability, the positions and velocities at the end of an integration will not converge for any practical timestep. So what *does* it mean for an N-body integration to be “accurate”?

While no significance can be attached to the locations of individual particles in an N-body simulation, the statistical properties of the particle distribution *are* reproducible. This fact emerges most clearly from cosmological simulations of the clustering of dark matter (Chapter 9). In these simulations an initially nearly smooth particle distribution develops a complex hierarchical structure with a wide range of densities. The densest knots are equilibrium gravitationally bound structures, which we identify with the dark halos of galaxies. Experiments have shown that when completely different computer codes are used to evolve the particles from the same initial conditions, the masses, locations, and other properties of the halos that form in different simulations are very nearly the same (Frenk et al. 1999). This coincidence of results holds even when fundamentally different Poisson solvers are employed (for example a tree code and a particle-mesh code; §§2.9.2 and 2.9.3), when different integration algorithms are used, or when different sets of particle coordinates \mathbf{w}_i are used to sample the initial cosmic density field. Thus we are confident that the *statistical properties* of the endpoint of a well-designed N-body simulation are meaningful, even though the locations of individual particles have no physical significance.

4.7.2 Schwarzschild models

We now describe a powerful technique introduced by Martin Schwarzschild¹⁴ for constructing an equilibrium model of a stellar system. **Schwarzschild’s method** is intermediate between N-body models, which follow individual particles (**particle-based** methods) and analytic techniques based on the DF. Schwarzschild’s method combines orbits to create a stellar system and hence is called an **orbit-based** method.

We describe how to construct a steady-state galaxy that has a given three-dimensional density distribution $\rho(\mathbf{x})$. We divide the space occupied by the galaxy into K cells, such that the mass in the j th cell of volume V_j is $m_j = \rho(\mathbf{x}_j)V_j$. Next, we calculate the galaxy’s gravitational potential, and integrate a large number N of orbits in this potential for a time t that is much longer than the crossing time. These orbits should have a wide variety

¹⁴ Martin Schwarzschild (1912–1997) was the elder son of Karl Schwarzschild, the foremost German astronomer of his age. Educated in Göttingen, in 1936 he fled the Nazis, going first to Oslo and then Harvard and finally to Princeton, where he used electronic computers to establish the theory of stellar structure. He used high-altitude balloons to obtain high-resolution images and infrared spectra of the Sun, planets and M31. Most of his work in galactic dynamics was done after he had retired.

of initial conditions so that they sample all the phase space that is likely to be occupied in the galaxy. This set of orbits is called the **orbit library**. We note the fraction p_{ij} of the time t that the i th orbit spends in the j th cell. Suppose there were a large number of stars on each orbit, uniformly distributed in orbital phase, and that the total mass of stars on orbit i was $w_i M$, where w_i is a weight to be determined and M is the total galactic mass. Then the mass in the j th cell would be $M \sum_i w_i p_{ij}$. Consequently, this arrangement would constitute a valid steady-state dynamical model of the given galaxy providing we chose the weights such that

$$0 = \Delta_j \equiv m_j - M \sum_{i=1}^N w_i p_{ij}. \quad (4.200)$$

This is a set of K linear equations for the N unknown weights w_i . The condition $\sum_j m_j = M$ implies that $\sum_i w_i = 1$.

We have to insist on a solution in which all the weights are non-negative. In view of this restriction, it is not profitable to put N equal to K and solve these equations by the standard methods of linear algebra, since the resulting solution vector \mathbf{w} will almost certainly contain negative components. The way forward is to take $N \gg K$ —many more orbits than spatial cells—in which case points that satisfy the equations form a $N - K$ dimensional subspace of the N dimensional space of weight vectors \mathbf{w} . A non-negative solution $w_i \geq 0$ will exist if this sub-space passes through the region in which all coordinates are positive. If the subspace does reach this region, infinitely many non-negative solution vectors exist, and every one corresponds to a physically acceptable galaxy model. Thus either we find no solution, or we have an embarrassment of riches, and we must find a rationale for choosing one of the infinite set of possible solutions. This is normally done by choosing the solution that maximizes some **objective function** $\phi(\mathbf{w})$.

The simplest possibility is that the objective function is linear, $\phi(\mathbf{w}) = \sum_i \phi_i w_i$. Many commercial problems can be reduced to the problem of maximizing a linear objective function $\phi(\mathbf{w})$ of N variables $w_i \geq 0$, subject to K constraints of the form (4.200). A problem of this type is said to be an exercise in **linear programming** and sophisticated software exists to solve such problems with large numbers of variables. Schwarzschild (1979) constructed stellar models of triaxial galaxies by choosing a linear objective function, pretty much at random, and then using standard software to solve the resulting linear programming problem.¹⁵ Any linear objective function selects models that lie on the boundary of the allowed subspace, and therefore tends to concentrate most of the mass in a small fraction of the library's orbits, thereby producing a galaxy model with a very irregular DF.

¹⁵ The choice of objective function was unimportant to Schwarzschild because his goal was only to prove that *some* self-consistent triaxial galaxy models were possible.

The density distributions of individual orbits have square-root singularities at their edges. Consequently, the contribution p_{ij} of an orbit to a given cell depends strongly on whether the cell lies just inside or just outside the orbit. As a result the models are noisy, and sensitive to the choice of grid and orbit library. It is important that the orbit library combines a sufficiently wide variety of orbits with a reasonably dense sampling of phase space, for two reasons. First, if the orbit library is inadequate, there will be no solution to the constraint equations (4.200) that has non-negative weights. Second, we often wish to explore the whole range of galaxy models that are consistent with the observations, since this enables us to assign confidence intervals to derived quantities such as the mass-to-light ratio. In practice intuition and experience must be used to put together a high-quality library.

Many extensions of Schwarzschild's method are possible. To obtain a smoother distribution of orbit weights, we may use a nonlinear objective function. One such function is the entropy $S = -\sum_i w_i \ln w_i$. A simpler alternative is a quadratic objective function of the form $\phi(\mathbf{w}) = -\sum_i w_i^2/W_i$, where $W_i > 0$. The physical meaning of the W_i is seen by maximizing ϕ subject only to the constraint $\sum_i w_i = 1$; this yields $w_i \propto W_i$. Hence ϕ finds the solution that is in some sense closest to the weights $\{W_i\}$, which we choose to reflect our prejudices about the structure of the galaxy. One advantage of a quadratic objective function is that maximizing ϕ subject to the constraint equations is then an exercise in **quadratic programming**, and standard packages exist for such work.

Another useful extension of Schwarzschild's method is to model kinematic data. These will consist of measurements of the LOSVD at various points on the sky. The LOSVD at a given point is a linear function of the orbit weights w_i , so the χ^2 that describes the difference between the observed and model LOSVDs is a quadratic function of the w_i . Hence we can minimize χ^2 by maximizing the objective function $\phi(\mathbf{w}) = -\chi^2$ using quadratic programming. A major application is the search for black holes at the centers of galaxies (§4.9.1).

Although Schwarzschild's method was devised to model galaxies in which some of the isolating integrals are not analytic, it can also be useful when modeling simpler systems. Consider, for example, a spherical system. By defining a sufficiently dense grid (E_i, L_i) in energy-angular-momentum space, it should be possible to get a good fit to any given body of observational data under the assumption that the galaxy's DF is a sum of delta functions that are centered on the grid points:

$$f(E, L) = \sum_{i=1}^N w_i f_i \quad \text{where} \quad f_i(E, L) \equiv \delta(E - E_i)\delta(L - L_i), \quad (4.201)$$

and the w_i are weights to be determined. We then approximate the physical

density $\rho(r)$ as

$$\rho(r_j) = \sum_{i=1}^N w_i \rho_i(r_j) \quad \text{where} \quad \rho_i(r) \equiv \int d^3\mathbf{v} f_i(E, L) \quad (4.202)$$

is the density produced by the family of orbits that have the given energy and total angular momentum, but all possible orientations of the orbital plane. For given values of $\rho(r_j)$ these equations define a linear programming problem for the weights in the same way that equation (4.200) does. In fact, the only difference between these equations is that in one case the galaxy is decomposed into individual orbits, and in the other a symmetry principle is used to group orbits into families within which all orbits must have the same weight, and then the galactic density is written as a sum of the densities contributed by each family.

Schwarzschild modeling has been extensively used to search for massive black holes at the centers of luminous spheroids (e.g. Richstone & Tremaine 1985; van der Marel et al. 1998; Gebhardt et al. 2003), the results of which are summarized by the correlation (1.27) between black-hole mass and spheroid velocity dispersion. In §4.9.1 we shall see why reliable black-hole masses can be obtained only with sophisticated dynamical modeling of the observational data. Schwarzschild modeling has also been used to model the large-scale dynamics of early-type galaxies, thus constraining the mass densities and orbital distributions in these systems (Cappellari et al. 2006 and §4.9.2). Unfortunately, the inference of confidence intervals on the values of model parameters, such as black-hole masses and mass-to-light ratios, when Schwarzschild's method is used to fit a model to observational data, proves to be a subtle matter (Magorrian 2006), and published values should be treated with some caution.

4.8 The Jeans and virial equations

In §4.1.2 we saw that comparisons between theoretical models and observational data often center on velocity moments of the DF, such as $\bar{\mathbf{v}}$ and $\overline{v_i v_j}$. Calculating moments is easy if one knows the DF, but finding a DF that is compatible with a given probability density distribution $\nu(\mathbf{x})$ is not straightforward, and even if a DF can be found, it is often not unique. Therefore in this section we discuss techniques for inferring moments from stellar densities without actually recovering the DF. Dejonghe (1986) gives an extensive discussion of this problem.

Integrating equation (4.11) over all velocities, we obtain

$$\int d^3\mathbf{v} \frac{\partial f}{\partial t} + \int d^3\mathbf{v} v_i \frac{\partial f}{\partial x_i} - \frac{\partial \Phi}{\partial x_i} \int d^3\mathbf{v} \frac{\partial f}{\partial v_i} = 0, \quad (4.203)$$

where we have employed the summation convention (page 772). The range of velocities over which we are integrating does not depend on time, so the partial derivative $\partial/\partial t$ in the first term of this equation may be taken outside the integral. Similarly, since v_i does not depend on x_i , the partial derivative $\partial/\partial x_i$ in the second term of the equation may be taken outside the integral sign. Furthermore, the last term on the left side of the equation vanishes on application of the divergence theorem (eq. B.46), given that $f(\mathbf{x}, \mathbf{v}, t) = 0$ for sufficiently large $|\mathbf{v}|$, i.e., there are no stars that move infinitely fast. Recalling the definitions of the density ν (eq. 4.20) and the mean velocity $\bar{\mathbf{v}}$ (eq. 4.24b), we have that

$$\frac{\partial \nu}{\partial t} + \frac{\partial(\nu \bar{v}_i)}{\partial x_i} = 0. \quad (4.204)$$

Equation (4.204) differs from the continuity equation (F.3) only in that it describes conservation of probability rather than that of mass, and replaces the fluid velocity by the mean stellar velocity.

We now multiply equation (4.11) by v_j and integrate over all velocities, and obtain

$$\frac{\partial}{\partial t} \int d^3\mathbf{v} f v_j + \int d^3\mathbf{v} v_i v_j \frac{\partial f}{\partial x_i} - \frac{\partial \Phi}{\partial x_i} \int d^3\mathbf{v} v_j \frac{\partial f}{\partial v_i} = 0. \quad (4.205)$$

The last term on the left side can be transformed by applying the divergence theorem, using the fact that f vanishes for large $|\mathbf{v}|$:

$$\int d^3\mathbf{v} v_j \frac{\partial f}{\partial v_i} = - \int d^3\mathbf{v} \frac{\partial v_j}{\partial v_i} f = - \int d^3\mathbf{v} \delta_{ij} f = -\delta_{ij} \nu. \quad (4.206)$$

Thus equation (4.205) may be rewritten

$$\frac{\partial(\nu \bar{v}_j)}{\partial t} + \frac{\partial(\nu \bar{v}_i \bar{v}_j)}{\partial x_i} + \nu \frac{\partial \Phi}{\partial x_j} = 0. \quad (4.207)$$

This can be put into a more familiar form by subtracting from it \bar{v}_j times the equation of continuity (4.204) to yield

$$\nu \frac{\partial \bar{v}_j}{\partial t} - \bar{v}_j \frac{\partial(\nu \bar{v}_i)}{\partial x_i} + \frac{\partial(\nu \bar{v}_i \bar{v}_j)}{\partial x_i} = -\nu \frac{\partial \Phi}{\partial x_j}, \quad (4.208)$$

and then using the definition (4.26) of the velocity-dispersion tensor to eliminate $\bar{v}_i \bar{v}_j$. The result is an analog of Euler's equation (F.7) of fluid flow;

$$\nu \frac{\partial \bar{v}_j}{\partial t} + \nu \bar{v}_i \frac{\partial \bar{v}_j}{\partial x_i} = -\nu \frac{\partial \Phi}{\partial x_j} - \frac{\partial(\nu \sigma_{ij}^2)}{\partial x_i}. \quad (4.209)$$

The left side and the first term on the right side of equation (4.209) differ from terms in the ordinary Euler equation only by the replacement of the mass density by the probability density, and by the replacement of the fluid velocity by the mean stellar velocity. The last term on the right side of equation (4.209) represents something akin to the pressure force $-\nabla p$. More exactly, $-\nu\sigma_{ij}^2$ is a **stress tensor** that describes an anisotropic pressure. Since equations (4.204) and (4.209) were first applied to stellar dynamics by Jeans (1919), we call them the **Jeans equations**.¹⁶

Equations (4.204) and (4.209) are valuable because they relate observationally accessible quantities, such as the streaming velocity, velocity dispersion, and so forth. However, this is an incomplete set of equations in the following sense. If we know the potential $\Phi(\mathbf{x}, t)$ and the density $\nu(\mathbf{x}, t)$, we have nine unknown functions—the three components of $\bar{\mathbf{v}}$ and the six independent components of the symmetric tensor σ^2 —but only four equations—the scalar continuity equation and the three components of Euler's equation. Thus we cannot solve for $\bar{\mathbf{v}}$ and σ^2 without additional information. The reader may argue that if we multiply the collisionless Boltzmann equation (4.11) through by $v_i v_k$ and integrate over all velocities, we obtain a new set of differential equations for σ^2 which might supply the missing information. Unfortunately, these equations involve quantities like $\overline{v_i v_j v_k}$ for which we would require still further equations. Thus these additional equations are of no use unless we can in some way truncate or *close* this regression to ever higher moments of the velocity distribution. We shall find that closure is possible only in special circumstances, for example when the system is spherical and we know that its DF is ergodic, $f(H)$ (Box 4.3), or when the system is axisymmetric and its DF is of the form $f(H, L_z)$. The equations can also be closed for any Stäckel potential (van de Ven et al. 2003).

4.8.1 Jeans equations for spherical systems

To obtain the Jeans equations in spherical coordinates, we start from the collisionless Boltzmann equation in the form (4.14), which involves the canonical momenta

$$p_r = \dot{r} = v_r \quad ; \quad p_\theta = r^2 \dot{\theta} = r v_\theta \quad ; \quad p_\phi = r^2 \sin^2 \theta \dot{\phi} = r \sin \theta v_\phi. \quad (4.210)$$

We have

$$\int dp_r dp_\theta dp_\phi f = r^2 \sin \theta \int dv_r dv_\theta dv_\phi f = r^2 \sin \theta \nu. \quad (4.211)$$

We assume that the system is spherical and time-independent, so we can drop $\partial\Phi/\partial\theta$, $\partial\Phi/\partial\phi$, $\partial f/\partial t$ and $\partial f/\partial\phi$ from (4.14); we retain $\partial f/\partial\theta$ because any

¹⁶ They were originally derived by Maxwell, but he already has a set of equations named after him.

dependence of f on v_ϕ is likely to introduce θ -dependence through the last of equations (4.210) when v_ϕ is expressed in terms of p_ϕ . After simplification, equation (4.14) becomes

$$p_r \frac{\partial f}{\partial r} + \frac{p_\theta}{r^2} \frac{\partial f}{\partial \theta} - \left(\frac{d\Phi}{dr} - \frac{p_\theta^2}{r^3} - \frac{p_\phi^2}{r^3 \sin^2 \theta} \right) \frac{\partial f}{\partial p_r} + \frac{p_\phi^2 \cos \theta}{r^2 \sin^3 \theta} \frac{\partial f}{\partial p_\theta} = 0. \quad (4.212)$$

We now multiply by $p_r dp_r dp_\theta dp_\phi$ and integrate over all momenta. With equation (4.211) and similar results, and using the divergence theorem to eliminate derivatives with respect to the momenta, we find

$$\frac{\partial}{\partial r} (r^2 \sin \theta \overline{\nu p_r^2}) + \frac{\partial}{\partial \theta} (\sin \theta \overline{\nu p_r p_\theta}) + r^2 \sin \theta \nu \left(\frac{d\Phi}{dr} - \frac{\overline{p_\theta^2}}{r^3} - \frac{\overline{p_\phi^2}}{r^3 \sin^2 \theta} \right) = 0. \quad (4.213)$$

In any static spherical system, $\overline{p_r p_\theta} = r \overline{v_r v_\theta}$ must vanish because the DF is of the form $f(H, \mathbf{L})$, and is therefore an even function of v_r . Finally, dividing through by $r^2 \sin \theta$ and using equations (4.210) we obtain

$$\frac{d(\overline{\nu v_r^2})}{dr} + \nu \left(\frac{d\Phi}{dr} + \frac{2\overline{v_r^2} - \overline{v_\theta^2} - \overline{v_\phi^2}}{r} \right) = 0. \quad (4.214)$$

In terms of the anisotropy parameter of equation (4.61), equation (4.214) reads

$$\frac{d(\overline{\nu v_r^2})}{dr} + 2\frac{\beta}{r} \overline{\nu v_r^2} = -\nu \frac{d\Phi}{dr}. \quad (4.215)$$

Additional Jeans equations can be obtained by multiplying (4.212) by p_θ or p_ϕ , but these are not useful.

If the line-of-sight velocity dispersion has been measured as a function of radius, equation (4.215) can be used to constrain the radial dependence of β . The most direct approach is to assume a functional form for $\beta(r)$ and treat (4.215) as a first-order linear differential equation for $\overline{\nu v_r^2}$. The integrating factor is $\exp(2 \int dr \beta/r)$, so the solution can be written in closed form. Different choices of $\beta(r)$ yield different predictions for the line-of-sight velocity dispersion as a function of radius (see Problem 4.28), so β can be constrained by optimizing the fit between predictions obtained from (4.215) and the observed velocity-dispersion profile.

The case of constant non-zero β is particularly simple. Then the solution of (4.215) that satisfies the boundary condition $\lim_{r \rightarrow \infty} \overline{v_r^2} = 0$ is

$$\overline{v_r^2}(r) = \frac{1}{r^{2\beta} \nu(r)} \int_r^\infty dr' r'^{2\beta} \nu(r') \frac{d\Phi}{dr'}. \quad (4.216)$$

Effect of a central black hole on the observed velocity dispersion

We can use this equation to assess the impact of a central massive black hole

Box 4.3: Closure of the Jeans equations when the DF is ergodic

We have shown that the Jeans equations are not closed, in the sense that $\overline{v_r^2}$ and β cannot both be determined from ν and Φ . However, if the DF is known to be ergodic, $f(H)$, then $\beta = 0$ and $\overline{v_r^2}$ is determined by equation (4.216). Moreover, *all* of the n th-order velocity moments can be determined from $\overline{v_r^n}$ (Problem 4.29). A differential equation for $\overline{v_r^n}$ is obtained when we multiply (4.212) by $p_r^{n-1} dp_r dp_\theta dp_\phi$ and integrate over all momenta. For example, with $n = 4$ we find

$$\frac{d(\overline{\nu v_r^4})}{dr} = -3\nu \left(\frac{d\Phi}{v_r^2 dr} + \frac{\frac{2}{3}\overline{v_r^4} - \overline{v_r^2 v_\theta^2} - \overline{v_r^2 v_\phi^2}}{r} \right) \stackrel{\beta=0}{=} -3\nu \overline{v_r^2} \frac{d\Phi}{dr}, \quad (1)$$

where the first equality is valid for any spherical system and the second is obtained by assuming that $f = f(H)$ and using the relation $\overline{v_r^2 v_\theta^2} = \frac{1}{3}\overline{v_r^4}$ from (4.308). Once $\overline{v_r^2}(r)$ is known, we can solve (1) for $\overline{v_r^4}(r)$ and from that derive the other fourth-order moments. Then we can solve a similar equation for $\overline{v_r^6}(r)$ and so on up to whatever moment we desire. When moments up to order $n \sim 10$ have been determined, accurate predictions of LOSVDs can be made (Magorrian & Binney 1994). These predictions will be identical to those one could have obtained from Eddington's formula (4.46b) for f but will not enable us to check that f is non-negative.

on the host galaxy's velocity-dispersion profile. We assume that the galaxy has a constant mass-to-light ratio and is a Hernquist model of scale-length a —from equations (2.64) and (2.67), the density and potential are

$$\nu(r) = \frac{1}{2\pi a^2} \frac{1}{r(1+r/a)^3} \quad ; \quad \Phi(r) = -\frac{GM_g}{r+a} - \frac{G\mu M_g}{r}, \quad (4.217)$$

where $\mu = M_\bullet/M_g$ is the ratio of the black-hole mass M_\bullet to the galaxy mass M_g . Hence

$$\overline{v_r^2}(ax) = \frac{GM_g}{a} \frac{(1+x)^3}{x^{2\beta-1}} \int_x^\infty dx' \left(\frac{x'^{2\beta-1}}{(1+x')^5} + \frac{\mu x'^{2\beta-3}}{(1+x')^3} \right). \quad (4.218)$$

For integer values of 4β the integrals are elementary. For example with $y \equiv 1+x$ we have

$$\frac{a\overline{v_r^2}(ax)}{GM_g} = \begin{cases} 5(1+2\mu)x^2y^3 \ln(x/y) + \mu y^3(\frac{1}{3} - \frac{3}{2}x + 6x^2)/x \\ \quad + x^2[\frac{1}{4} + \frac{2}{3}y + \frac{3+\mu}{2}y^2 + 4(1+\mu)y^3]/y + xy^3 & (\beta = -\frac{1}{2}), \\ (1+6\mu)xy^3 \ln(y/x) - \mu y^3(3x - \frac{1}{2})/x \\ \quad - x[\frac{1}{4} + \frac{1}{3}y + \frac{1+\mu}{2}y^2 + (1+3\mu)y^3]/y & (\beta = 0), \\ 3\mu y^3 \ln(x/y) + 1/4y + \mu y(\frac{1}{2} + 2y + y^2/x) & (\beta = \frac{1}{2}). \end{cases} \quad (4.219)$$

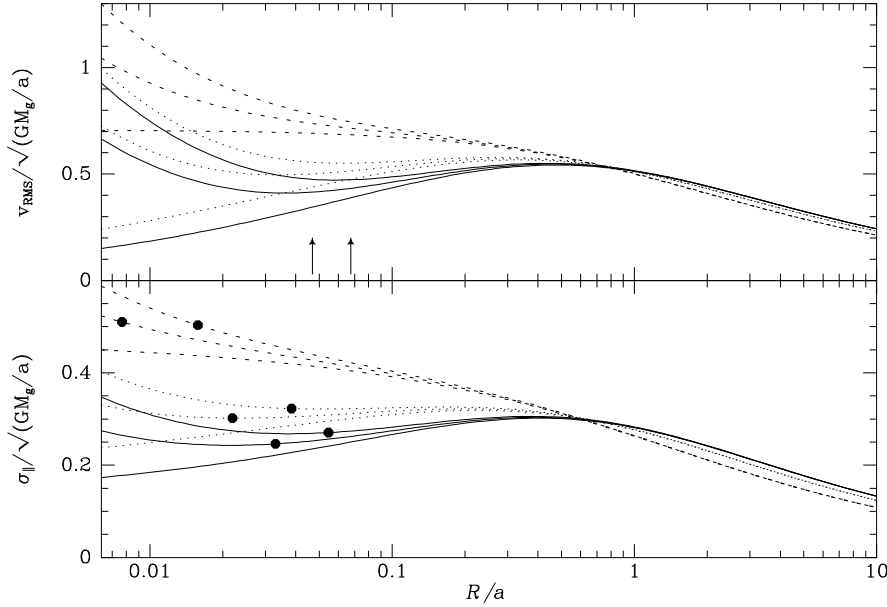


Figure 4.20 Velocity dispersion as a function of radius for three Hernquist models with a central black hole of mass 0, $0.002M_g$, or $0.004M_g$. The bottom panel shows line-of-sight dispersions, the top panel shows the RMS speed as a function of radius. The full curves are for tangential bias ($\beta = -0.5$), the dotted curves are for the isotropic model and the dashed curves are for radial bias ($\beta = 0.5$). The beads mark the radius of influence (eq. 4.220) of the black hole in each model, while the arrows mark the dynamical radius of the black hole, at which the interior mass of the galaxy equals the mass of the black hole.

The top panel of Figure 4.20 shows the RMS speed $v_{\text{RMS}} = (\overline{v_r^2} + \overline{v_\theta^2} + \overline{v_\phi^2})^{1/2}$ that follows from these formulae for $\mu = 0$, $\mu = 0.002$, and $\mu = 0.004$ (bottom to top). The full curves are for tangentially biased models, the dotted curves for isotropic models, and the dashed curves are for radially biased models. In each case the black hole causes the RMS speed to rise at small radii where its deep potential well speeds up the stars. The lower panel shows the associated line-of-sight dispersions. At small radii the upturn in σ_{\parallel} is much less sharp than that in v_{RMS} , because the signal from stars near the black hole is diluted by the light from foreground and background stars. Note also that the rise in dispersion associated with the black hole is difficult to distinguish from the rise in dispersion associated with radial anisotropy.

The black hole's **radius of influence** R_{infl} is defined to be the radius at which the Kepler speed due to the hole is equal to σ_{\parallel} . Quantitatively,

$$R_{\text{infl}} = \frac{GM_{\bullet}}{\sigma_{\parallel}^2(R_{\text{infl}})} = 11 \frac{M_{\bullet}}{10^8 M_{\odot}} \left(\frac{\sigma_{\parallel}}{200 \text{ km s}^{-1}} \right)^{-2} \text{ pc}. \quad (4.220)$$

In Figure 4.20 R_{infl} is marked by a black dot on each relevant curve of $\sigma_{\parallel}(R)$.

It can be seen that at R_{infl} the black hole has increased σ_{\parallel} by a few percent, and that the contribution to σ_{\parallel} from the black hole increases fairly rapidly interior to R_{infl} .

Another measure of the radial extent of the black hole's influence is the **dynamical radius** r_g of the black hole, at which the gravitational forces from the black hole and the galaxy are equal, or, equivalently, the radius within which the galactic mass is equal to the black-hole mass. The dynamical radius r_g , unlike the radius of influence R_{infl} , depends only on the galaxy's mass distribution and not its kinematics. Orbits with apocenters inside r_g will be nearly Keplerian. The vertical arrows in Figure 4.20 mark r_g for the two black-hole masses considered. For the tangential models, r_g is only slightly larger than R_{infl} , while in the radially biased models, $r_g \sim 7R_{\text{infl}}$.

This discussion demonstrates that a major obstacle to detecting a central black hole using stellar kinematics is the degeneracy between the mass of the black hole and velocity anisotropy. This degeneracy can be lifted by obtaining data with higher spatial resolution than assumed in Figure 4.20. Alternatively, we can exploit the information contained in the entire LOSVD rather than just its second moment (§4.9.1).

4.8.2 Jeans equations for axisymmetric systems

For simplicity we assume that the system under study is in a steady state and axisymmetric so all derivatives with respect to t and ϕ vanish. With these assumptions (4.12) becomes

$$p_R \frac{\partial f}{\partial R} + p_z \frac{\partial f}{\partial z} - \left(\frac{\partial \Phi}{\partial R} - \frac{p_\phi^2}{R^3} \right) \frac{\partial f}{\partial p_R} - \frac{\partial \Phi}{\partial z} \frac{\partial f}{\partial p_z} = 0. \quad (4.221)$$

We multiply this equation by p_R , integrate over the momenta $p_R = v_R$, $p_\phi = Rv_\phi$, $p_z = v_z$, and then express the momenta in terms of velocities. In close analogy with our derivation of equation (4.215) we obtain

$$\frac{\partial(\overline{\nu v_R^2})}{\partial R} + \frac{\partial(\overline{\nu v_R v_z})}{\partial z} + \nu \left(\frac{\overline{v_R^2} - \overline{v_\phi^2}}{R} + \frac{\partial \Phi}{\partial R} \right) = 0. \quad (4.222a)$$

When we multiply (4.221) by p_z or p_ϕ rather than p_R , we obtain

$$\frac{1}{R} \frac{\partial(R \overline{\nu v_R v_z})}{\partial R} + \frac{\partial(\overline{\nu v_z^2})}{\partial z} + \nu \frac{\partial \Phi}{\partial z} = 0, \quad (4.222b)$$

$$\frac{1}{R^2} \frac{\partial(R^2 \overline{\nu v_R v_\phi})}{\partial R} + \frac{\partial(\overline{\nu v_z v_\phi})}{\partial z} = 0. \quad (4.222c)$$

If we assume that the density $\nu(R, z)$ and the confining potential $\Phi(R, z)$ are known, equations (4.222) constitute three constraints on the six second-order

Box 4.4: Two useful formulae

If we obtain ν by integrating $f(H, L_z)$ over all velocities, the resulting expression will depend on z only through $\Phi(R, z)$. In these circumstances it is advantageous to consider ν to be a function of (R, Φ) and equation (4.223) yields

$$\overline{\nu v_R^2}(R, z) = \int_{\Phi(R, z)}^0 d\Phi' \nu(R, \Phi') \quad (1)$$

Multiplying equation (4.224) by ν and using (1) we obtain

$$\overline{\nu v_\phi^2} = \frac{\partial}{\partial R} \left(R \int_{\Phi}^0 d\Phi' \nu(R, \Phi') \right) + \nu R \frac{\partial \Phi}{\partial R}.$$

In the first term on the right we carry the factor R inside the integral and then use the standard formula

$$\frac{d}{dx} \int_{f(x)}^0 dy g(x, y) = \int_{f(x)}^0 dy \frac{\partial g}{\partial x} - g(x, f) \frac{df}{dx}$$

to establish that

$$\overline{\nu v_\phi^2} = \int_{\Phi}^0 d\Phi' \frac{\partial}{\partial R} [R\nu(R, \Phi')]. \quad (2)$$

velocity moments. Thus, just as in the spherical case, the Jeans equations are not closed. However, if the DF is known to be of the form $f(H, L_z)$, the mixed moments in these equations will vanish, $\overline{v_R^2} = \overline{v_z^2}$, and the third equation becomes trivial. So we have two equations for two unknowns, and the system is closed. Specifically, (4.222b) can be integrated to yield (Nagai & Miyamoto 1976)

$$\overline{v_R^2}(R, z) = \overline{v_z^2}(R, z) = \frac{1}{\nu(R, z)} \int_z^\infty dz' \nu(R, z') \frac{\partial \Phi}{\partial z'}. \quad (4.223)$$

Now that $\overline{v_R^2}$ is known, we can obtain $\overline{v_\phi^2}$ from (4.222a):

$$\overline{v_\phi^2}(R, z) = \overline{v_R^2} + \frac{R}{\nu} \frac{\partial(\nu \overline{v_R^2})}{\partial R} + R \frac{\partial \Phi}{\partial R}. \quad (4.224)$$

Proceeding similarly with higher-order Jeans equations obtained by multiplying (4.221) by $p_z^{k+1} p_\phi^{n-k-2}$ for $k = 0, 1, \dots, n-2$ and by $p_R p_\phi^{n-2}$, we can relate all the n -order moments to either $\nu(R, z)$ or $\nu \overline{v_\phi^2}(R, z)$, depending on whether n is even or odd (Magorrian & Binney 1994). These moments will be identical to those one would have obtained by using the Hunter–Qian algorithm to calculate $f(\mathcal{E}, L_z)$ from the same data (§4.4.1).

(a) Asymmetric drift Figure 4.17 shows that in the solar neighborhood

the distribution of high-velocity stars is strongly asymmetric, in the sense that there are more stars lagging the LSR than leading it. We saw on page 325 that this phenomenon is nicely explained by the surface-density and velocity-dispersion gradients in the disk, and more quantitatively by Figure 4.16, but from equation (4.222a) we can easily recover its most important aspect, which is the asymmetric drift (page 326)

$$v_a \equiv v_c - \bar{v}_\phi, \quad (4.225)$$

where v_c is the circular speed in the solar neighborhood. We consider the values of v_a of a sequence of stellar populations, each with its own value of $\overline{v_R^2}$.

We assume that the disk is in a steady state and is symmetric about its equator. Then, since the Sun lies close to the galactic equator, we may evaluate equation (4.222a) at $z = 0$. Since $\partial\nu/\partial z = 0$ by symmetry, we find

$$\frac{R}{\nu} \frac{\partial(\overline{\nu v_R^2})}{\partial R} + R \frac{\partial(\overline{v_R v_z})}{\partial z} + \overline{v_R^2} - \overline{v_\phi^2} + R \frac{\partial\Phi}{\partial R} = 0 \quad (z = 0). \quad (4.226)$$

Using equation (4.26) to replace $\overline{v_\phi^2}$ by the azimuthal velocity dispersion σ_ϕ^2 and using $R(\partial\Phi/\partial R) = v_c^2$, we obtain

$$\begin{aligned} \sigma_\phi^2 - \overline{v_R^2} - \frac{R}{\nu} \frac{\partial(\overline{\nu v_R^2})}{\partial R} - R \frac{\partial(\overline{v_R v_z})}{\partial z} &= v_c^2 - \overline{v_\phi^2} \\ &= (v_c - \bar{v}_\phi)(v_c + \bar{v}_\phi) = v_a(2v_c - v_a). \end{aligned} \quad (4.227)$$

If we neglect v_a compared to $2v_c$, we obtain Stromberg's **asymmetric drift equation**

$$v_a \simeq \frac{\overline{v_R^2}}{2v_c} \left[\frac{\sigma_\phi^2}{\overline{v_R^2}} - 1 - \frac{\partial \ln(\overline{\nu v_R^2})}{\partial \ln R} - \frac{R}{\overline{v_R^2}} \frac{\partial(\overline{v_R v_z})}{\partial z} \right]. \quad (4.228)$$

The value of the square bracket does not depend on the scale of the velocity-dispersion tensor $\overline{v_i v_j}$, but only on the ratios of its components. So if two populations have similar density distributions $\nu(R, z)$ and velocity ellipsoids of the same shape and orientation, the square bracket will take the same value for both populations. Hence in this case $v_a \propto \overline{v_R^2}$. Figure 4.21 shows that a relationship of this type holds for main-sequence stars near the Sun. The horizontal axis shows the dispersions in the velocities normal to the line of sight for stars in each population. The vertical axis shows the average amount by which the stars lag the azimuthal motion of the Sun. Each data point is for one bin in stellar color $B - V$. The redder bins contain older stars, which have larger dispersions S because stars are gradually accelerated by fluctuations in the gravitational potential (§8.4). The intersection of the

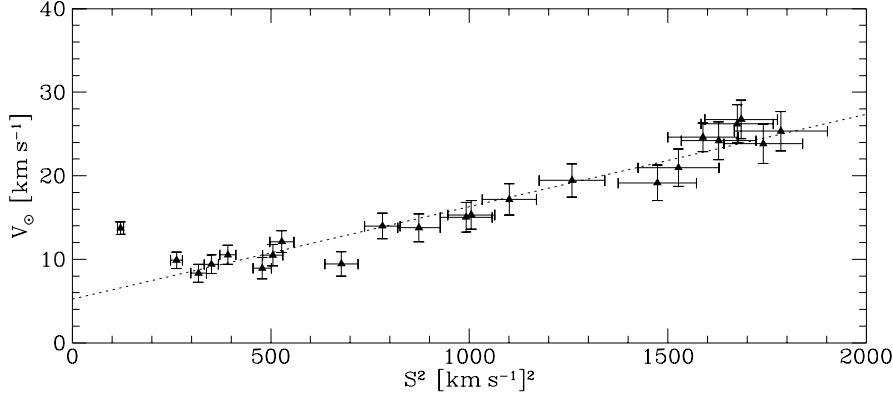


Figure 4.21 The asymmetric drift v_a for different stellar types is a linear function of the random velocity S^2 of each type. The vertical coordinate is actually $v_a + \tilde{v}_{\phi, \odot}$ where $\tilde{v}_{\phi, \odot}$ is the azimuthal velocity of the Sun relative to the LSR (after Dehnen & Binney 1998b).

best-fit line with $S = 0$, at $v_{\odot} = 5 \text{ km s}^{-1}$, represents the velocity of the Sun relative to the LSR.

It is interesting to compare the numerical value of the square bracket in equation (4.228) with the slope of the straight-line fit to the data in Figure 4.21. From BM Table 10.2 we adopt $\sigma_{\phi}^2/\overline{v_R^2} = 0.35$ and we assume that ν and $\overline{v_R^2}$ are both proportional to e^{-R/R_d} with $R_0/R_d = 3.2$ (Table 1.2)—this assumption regarding the radial dependence of the velocity dispersion is justified following equation (4.156). Then the bracket's first three terms sum to 5.8. The last term is problematic because its value depends on the orientation of the velocity ellipsoid at points just above the plane of our Galaxy, which is difficult to measure. Two extreme possibilities are that (i) the ellipsoid's principal axes are aligned with the coordinate directions of the (R, ϕ, z) system, and (ii) the principal axes are aligned with the coordinate directions of the (r, θ, ϕ) system centered on the galactic center. Orbit integrations (Binney & Spergel 1983) suggest that the truth lies nearly midway between these two possibilities. In the first case $\overline{v_R v_z}$ is independent of z and the term vanishes, and in the second $\overline{v_R v_z} \simeq (\overline{v_R^2} - \overline{v_z^2})(z/R)$ (see Problem 4.34) and the term contributes $-(1 - \overline{v_z^2}/\overline{v_R^2}) \simeq -0.8$. Averaging these values we estimate the value of the square bracket at 5.4 ± 0.4 , so $v_a \simeq \overline{v_R^2}/(82 \pm 6 \text{ km s}^{-1})$. From the data shown in Figure 4.21 one infers $v_a = \overline{v_R^2}/(80 \pm 5 \text{ km s}^{-1})$ in beautiful agreement with theory.

(b) Spheroidal components with isotropic velocity dispersion We know that if an axisymmetric system has a DF of the form $f(H, L_z)$ then two eigenvalues of the velocity-dispersion tensor σ^2 are equal (eq. 4.40). We now use the Jeans equations to predict the rotation rate of a spheroidal system in which all three eigenvalues of σ^2 are equal, that is, an isotropic rotator.

From the definition (4.26) of σ^2 and this assumption we have

$$\overline{v_\phi^2} = \overline{v_\phi^2} + \sigma_\phi^2 = \overline{v_\phi^2} + \overline{v_R^2}, \quad (4.229)$$

so equation (4.224) yields

$$\overline{v_\phi^2}(R, z) = R \frac{\partial \Phi}{\partial R} + \frac{R}{\nu} \frac{\partial(\nu \overline{v_R^2})}{\partial R}. \quad (4.230)$$

When we use equation (4.223) to eliminate $\nu \overline{v_R^2}$ we have

$$\overline{v_\phi^2}(R, z) = R \frac{\partial \Phi}{\partial R} + \frac{R}{\nu} \frac{\partial}{\partial R} \int_z^\infty dz' \nu(R, z') \frac{\partial \Phi}{\partial z'}. \quad (4.231)$$

Suppose both $\nu(R, z)$ and $\Phi(R, z)$ are constant on spheroids, which will be nearly true in many realistic cases. Then we can write $\nu(q_\nu^2 R^2 + \zeta)$ and $\Phi(q_\Phi^2 R^2 + \zeta)$ where $q_\nu < 1$ is the axis ratio of the isodensity surfaces, q_Φ is the axis ratio of the equipotentials, and $\zeta \equiv z^2$. Consequently, $\partial \nu / \partial R^2 = q_\nu^2 (\partial \nu / \partial \zeta)$ and $\partial \Phi / \partial R^2 = q_\Phi^2 (\partial \Phi / \partial \zeta)$. We convert the derivative in equation (4.231) into one with respect to R^2 , carry it under the integral sign, and use these relations to obtain

$$\begin{aligned} \overline{v_\phi^2}(R, z) &= R \frac{\partial \Phi}{\partial R} + \frac{2R^2}{\nu} \int_{z^2}^\infty d\zeta \left(q_\nu^2 \frac{\partial \nu}{\partial \zeta} \frac{\partial \Phi}{\partial \zeta} + q_\Phi^2 \nu \frac{\partial^2 \Phi}{\partial \zeta^2} \right) \\ &= R \frac{\partial \Phi}{\partial R} + (q_\nu^2 - q_\Phi^2) \frac{2R^2}{\nu} \int_{z^2}^\infty d\zeta \frac{\partial \nu}{\partial \zeta} \frac{\partial \Phi}{\partial \zeta} - 2R^2 q_\Phi^2 \frac{\partial \Phi}{\partial \zeta} \Big|_{z^2}, \end{aligned} \quad (4.232)$$

where the second equality is obtained by integrating by parts the term with the second derivative of Φ . We now observe that

$$2R^2 q_\Phi^2 \frac{\partial \Phi}{\partial \zeta} = 2R^2 \frac{\partial \Phi}{\partial R^2} = R \frac{\partial \Phi}{\partial R}. \quad (4.233)$$

Hence the last term on the right of (4.232) cancels the first term, and we have finally

$$\overline{v_\phi^2}(R, z) = (q_\nu^2 - q_\Phi^2) \frac{2R^2}{\nu} \int_{z^2}^\infty d\zeta \frac{\partial \nu}{\partial \zeta} \frac{\partial \Phi}{\partial \zeta}. \quad (4.234)$$

Since $\partial \Phi / \partial z > 0$ and $\partial \nu / \partial z < 0$, the integral is negative and $\overline{v_\phi} \propto \sqrt{q_\Phi^2 - q_\nu^2}$. If the isodensity surfaces coincide with the equipotentials, $q_\Phi = q_\nu$ and $\overline{v_\phi} = 0$, but normally the equipotentials are less flattened than the equidensity surfaces, and even a small excess in the flattening of the density distribution gives rise to appreciable rotation.

As an illustration of the use of equation (4.234), suppose Φ and ν are given by

$$\Phi = \frac{1}{2} v_0^2 \ln(R_c^2 + q_\Phi^2 R^2 + z^2) \quad ; \quad \nu = K(R_c^2 + q_\nu^2 R^2 + z^2)^{-3/2}, \quad (4.235)$$

where v_0 , R_c and K are constants. These functional forms are appropriate to the case of a galaxy that has a distribution of luminous matter consistent with a modified Hubble model (eq. 2.53) and an asymptotically flat circular-speed curve. Then the integral in equation (4.234) evaluates to

$$\begin{aligned} \int_{z^2}^{\infty} d\zeta \frac{\partial \nu}{\partial \zeta} \frac{\partial \Phi}{\partial \zeta} &= -\frac{3}{4} K v_0^2 \int_{z^2}^{\infty} \frac{d\zeta}{(R_c^2 + q_{\Phi}^2 R^2 + \zeta)(R_c^2 + q_{\nu}^2 R^2 + \zeta)^{5/2}} \\ &= -\frac{\frac{3}{2} K v_0^2}{(R_c^2 + q_{\Phi}^2 R^2 + z^2)^{5/2}} \frac{1}{\delta^4} \left(\frac{\sin^{-1} \delta}{\delta} + \frac{\frac{4}{3} \delta^2 - 1}{(1 - \delta^2)^{3/2}} \right). \end{aligned} \quad (4.236a)$$

where

$$\delta^2 \equiv \frac{(q_{\Phi}^2 - q_{\nu}^2) R^2}{R_c^2 + q_{\Phi}^2 R^2 + z^2}. \quad (4.236b)$$

To understand how this rather cumbersome formula works, we expand in powers of δ before substituting into equation (4.234) to obtain

$$\bar{v}_{\phi}^2(R, z) = \frac{3}{5} v_0^2 \left(\frac{R_c^2 + q_{\nu}^2 R^2 + z^2}{R_c^2 + q_{\Phi}^2 R^2 + z^2} \right)^{3/2} [\delta^2 + \frac{25}{14} \delta^4 + O(\delta^6)]. \quad (4.237)$$

At $R, z \ll R_c$, $\delta \propto R$ so $\bar{v}_{\phi} \propto R$ and there is solid-body rotation. Beyond R_c in the equatorial plane δ becomes independent of R and

$$\bar{v}_{\phi}/v_0 \rightarrow \sqrt{\frac{3}{5}(1 - q_{\nu}^2/q_{\Phi}^2)} (q_{\nu}/q_{\Phi})^{3/2}. \quad (4.238)$$

Observations indicate that within the effective radius of a typical luminous galaxy, the mass distribution is dominated by stars (Gerhard et al. 2001; Cappellari et al. 2006). From §2.3.2 we know that the ellipticity $\epsilon_{\nu} \equiv 1 - q_{\nu}$ of the density distribution that generates the logarithmic potential (4.235) is $\simeq 3\epsilon_{\Phi} = 3(1 - q_{\Phi})$. When we use this relation in equation (4.238), we find that $\bar{v}_{\phi}/v_0 = \sqrt{\frac{4}{5}\epsilon_{\nu}} + O(\epsilon_{\nu}^2)$. In Figure 4.14 the dotted curve shows the relationship $v/\sigma \propto \sqrt{\epsilon}$, and one sees that this proportionality provides a reasonable fit to the data for Evans models. It lies below the data for Rowley models, because these are not isotropic rotators. The filled circles, which show the data for low-luminosity spheroids, scatter around the dotted curve, although there is a tendency for the points to lie below the curve for $\epsilon \lesssim 0.5$ and above the curve at higher ellipticities. Thus these data suggest that low-luminosity spheroids are nearly isotropic rotators.

4.8.3 Virial equations

We obtained the Jeans equation (4.207) by multiplying the collisionless Boltzmann equation by v_j and integrating over all velocities. In this process an

equation in the six phase-space coordinates for a single scalar quantity f was reduced to partial differential equations for ν and the velocity moments in the three spatial coordinates. We now multiply equation (4.207) by x_k and integrate over all positions, thus converting these differential equations into a simple tensor equation relating global properties of the galaxy, such as total kinetic energy and mean-square streaming velocity.

We multiply equation (4.207) by Mx_k , where M is the total mass of the system. Then since the mass density is $\rho(\mathbf{x}) = M\nu(\mathbf{x})$, integrating over the spatial variables we find

$$\int d^3\mathbf{x} x_k \frac{\partial(\rho\bar{v}_j)}{\partial t} = - \int d^3\mathbf{x} x_k \frac{\partial(\rho\bar{v}_i\bar{v}_j)}{\partial x_i} - \int d^3\mathbf{x} \rho x_k \frac{\partial\Phi}{\partial x_j}. \quad (4.239)$$

The second term on the right side is the potential-energy tensor \mathbf{W} (eq. 2.19). The first term on the right side of equation (4.239) can be rewritten with the aid of the divergence theorem (B.45):

$$\int d^3\mathbf{x} x_k \frac{\partial(\rho\bar{v}_i\bar{v}_j)}{\partial x_i} = - \int d^3\mathbf{x} \delta_{ki} \rho\bar{v}_i\bar{v}_j = -2K_{kj}, \quad (4.240a)$$

where we have assumed that ρ vanishes at large radii and have defined the **kinetic-energy tensor**

$$K_{jk} \equiv \frac{1}{2} \int d^3\mathbf{x} \rho\bar{v}_j\bar{v}_k. \quad (4.240b)$$

With the help of equation (4.26) we split \mathbf{K} up into the contributions from ordered and random motion:

$$K_{jk} = T_{jk} + \frac{1}{2}\Pi_{jk}, \quad (4.241a)$$

where

$$T_{jk} \equiv \frac{1}{2} \int d^3\mathbf{x} \rho\bar{v}_j\bar{v}_k \quad ; \quad \Pi_{jk} \equiv \int d^3\mathbf{x} \rho\sigma_{jk}^2. \quad (4.241b)$$

The derivative with respect to time in equation (4.239) may be taken outside the integral sign because x_k does not depend on time. Finally, averaging the (k, j) and the (j, k) components of equation (4.239), we obtain

$$\frac{1}{2} \frac{d}{dt} \int d^3\mathbf{x} \rho (x_k\bar{v}_j + x_j\bar{v}_k) = 2T_{jk} + \Pi_{jk} + W_{jk}. \quad (4.242)$$

Here we have exploited the symmetry under exchange of indices of \mathbf{T} , $\mathbf{\Pi}$ (see eq. 4.241b) and \mathbf{W} (see eq. 2.22).

The left side of equation (4.242) may be brought to a more intuitive form if we define the tensor¹⁷ \mathbf{I} by

$$I_{jk} \equiv \int d^3\mathbf{x} \rho x_j x_k. \quad (4.243)$$

Differentiating \mathbf{I} with respect to time, we have

$$\frac{dI_{jk}}{dt} = \int d^3\mathbf{x} \frac{\partial \rho}{\partial t} x_j x_k. \quad (4.244)$$

With the continuity equation (4.204), the right side of this equation becomes

$$- \int d^3\mathbf{x} \frac{\partial(\rho \bar{v}_i)}{\partial x_i} x_j x_k = \int d^3\mathbf{x} \rho \bar{v}_i (x_k \delta_{ji} + x_j \delta_{ki}), \quad (4.245)$$

where the equality follows by an application of the divergence theorem. Substituting this expression back into equation (4.244) yields

$$\frac{dI_{jk}}{dt} = \int d^3\mathbf{x} \rho (x_k \bar{v}_j + x_j \bar{v}_k). \quad (4.246)$$

We now combine equations (4.242) and (4.246) to obtain the **tensor virial theorem**:

$$\frac{1}{2} \frac{d^2 I_{jk}}{dt^2} = 2T_{jk} + \Pi_{jk} + W_{jk}. \quad (4.247)$$

Equation (4.247) enables us to relate the gross kinematic and morphological properties of galaxies.¹⁸ In many applications the left side is simply zero since the galaxy is time-independent.

(a) Scalar virial theorem The trace of the potential-energy tensor is the system's total potential energy W (eq. 2.23). Equations (4.240b) show that $K \equiv \text{trace}(\mathbf{T}) + \frac{1}{2}\text{trace}(\mathbf{\Pi})$ is the total kinetic energy of the system. Thus, if the system is in a steady state, $\ddot{\mathbf{I}} = 0$, and the trace of equation (4.247) becomes

$$2K + W = 0. \quad (4.248)$$

Equation (4.248) is a statement of the **scalar virial theorem**.¹⁹ The kinetic energy of a stellar system with mass M is just $K = \frac{1}{2}M\langle v^2 \rangle$, where $\langle v^2 \rangle$ is

¹⁷ The tensor defined by equation (4.243) is sometimes called the “moment of inertia tensor” but we reserve this name for the related tensor that is defined by equation (D.41).

¹⁸ Equation (4.247) has here been derived from the collisionless Boltzmann equation, which is only valid for a collisionless system, but we shall find in §7.2.1 that an analogous result is valid for any system of N mutually gravitating particles. However, it should be noted that (4.247) applies only to *self-gravitating* systems. Similar results may be derived for systems embedded in an externally generated gravitational field; see Problems 3.12 and 4.38.

¹⁹ First proved by R. Clausius in 1870; Clausius also defined the **virial** of a system of N particles as $\sum_i^N m_i \mathbf{r}_i \cdot \mathbf{v}_i$. The theorem was first applied to stellar systems by Eddington (1916a). Einstein (1921) used it to estimate the mass of globular clusters.

the mean-square speed of the system's stars. Hence the virial theorem states that

$$\langle v^2 \rangle = \frac{|W|}{M} = \frac{GM}{r_g}, \quad (4.249a)$$

where r_g is the gravitational radius defined by equation (2.42). One often wishes to estimate $\langle v^2 \rangle$ without going to the trouble of calculating r_g . Spitzer (1969) noted that in simple stellar systems the half-mass radius r_h , which is easily measured, is tightly correlated with r_g . For example, the Jaffe and Hernquist models (§2.2.2g) have $r_h/r_g = \frac{1}{2}$ and 0.402, respectively, while for spherical galaxies that have radius-independent mass-to-light ratios and satisfy the Sérsic law (1.17) in projection, r_h/r_g ranges from 0.414 for $m = 2$ to 0.526 for $m = 6$ (Ciotti 1991). Moreover, we saw in §4.3.3c that along the sequence of King models r_h/r_g is confined to the interval (0.4, 0.51) (Figure 4.10). Hence, a useful approximation is

$$\langle v^2 \rangle = \frac{|W|}{M} \simeq 0.45 \frac{GM}{r_h}. \quad (4.249b)$$

If E is the energy of the system, we have from equation (4.248) that

$$E = K + W = -K = \frac{1}{2}W. \quad (4.250)$$

Thus if a system forms by collecting material together from a state of rest at infinity (in which state, $K = W = E = 0$), and then settles by any process into an equilibrium condition, it invests half of the gravitational energy that is released by the collapse in kinetic form, and in some way disposes of the other half in order to achieve a binding energy $E_b = -E$ equal to its kinetic energy. For example, suppose that our Galaxy formed by aggregating from an initial radius that was much larger than its present size. Then, since most of the galactic material is now moving at about $v_c \simeq 200 \text{ km s}^{-1}$, whether on circular orbits in the disk or on eccentric and highly inclined halo orbits, we have that $E_b = K \approx \frac{1}{2}M_g v_c^2$ of energy must have been released when the Galaxy formed, where M_g is the mass of the Galaxy. This argument suggests that as they form, galaxies radiate a fraction $\frac{1}{2}(v_c/c)^2 \simeq 3 \times 10^{-7}$ of their rest-mass energy.

(b) Spherical systems We may use the scalar virial theorem (4.248) to evaluate the mass-to-light ratio Υ of a non-rotating spherical galaxy under the assumption that Υ is independent of radius. We choose a coordinate system in which the line of sight to the galaxy center coincides with the x axis. Then the kinetic energy associated with motion in the x -direction is

$$K_{xx} = \frac{1}{2} \int d^3\mathbf{x} \rho v_x^2. \quad (4.251)$$

Rewriting this in terms of the luminosity density $j = \rho/\Upsilon$, we have

$$K_{xx} = \frac{1}{2}\Upsilon \int dy \int dz \int dx j \overline{v_x^2}. \quad (4.252)$$

The innermost integral in this expression yields the luminosity-weighted dispersion of the line-of-sight velocities at position (y, z) . Hence K_{xx} may be expressed in terms of the surface brightness $I(y, z)$ and the line-of-sight velocity dispersion $\sigma_{\parallel}(y, z)$ as

$$K_{xx} = \frac{1}{2}\Upsilon \int dy \int dz I(y, z) \sigma_{\parallel}^2(y, z) = \frac{1}{2}\Upsilon J, \quad (4.253a)$$

where J is the luminosity-weighted squared line-of-sight velocity dispersion

$$J \equiv 2\pi \int_0^{\infty} dR R I(R) \sigma_{\parallel}^2(R). \quad (4.253b)$$

Since the galaxy is assumed to be spherical and non-rotating, the total kinetic energy is

$$K = 3K_{xx} = \frac{3}{2}\Upsilon J. \quad (4.254)$$

On the other hand, from equation (1.79) we have

$$\rho(r) = -\frac{\Upsilon}{\pi} \int_r^{\infty} \frac{dR}{\sqrt{R^2 - r^2}} \frac{dI}{dR}. \quad (4.255)$$

When we use this relation in equation (2.32), we obtain W as

$$W = \Upsilon^2 \tilde{J}, \quad (4.256)$$

where \tilde{J} is an integral that depends only on $I(R)$. Using these results in (4.248), we obtain finally

$$\Upsilon = -3J/\tilde{J}. \quad (4.257)$$

Thus Υ may be obtained from measurements of $I(R)$ and $\sigma_{\parallel}(R)$. The nice feature of equation (4.257) is that it is valid no matter what velocity anisotropy there may be in the system. The main issues that must be addressed in any practical application are: (i) do the kinematic data extend far enough out for J to be calculated reliably?; (ii) is Υ really a constant throughout the system? (iii) is the system really spherical, or is it an axisymmetric system seen pole-on?

(c) The tensor virial theorem and observational data The analysis of observational data is facilitated if we reformulate the tensor virial theorem so that it involves the mean line-of-sight velocity \overline{v}_{\parallel} and the line-of-sight velocity dispersion σ_{\parallel} that are defined by equations (4.24) and (4.27). By

analogy with (4.24a) we have that the mean-square line-of-sight velocity at a location on the sky is

$$\overline{v_{\parallel}^2} = \frac{1}{\Sigma} \int dx_{\parallel} \int d^3\mathbf{v} v_{\parallel}^2 f, \quad (4.258a)$$

where

$$\Sigma = \int dx_{\parallel} \nu = \int dx_{\parallel} \int d^3\mathbf{v} f \quad (4.258b)$$

is the probability per unit area of finding a particular star along the line of sight. Hence, the line-of-sight velocity dispersion that is defined by equation (4.25) is related to the DF by

$$\begin{aligned} \sigma_{\parallel}^2 &= \overline{(v_{\parallel} - \bar{v}_{\parallel})^2} = \overline{v_{\parallel}^2} - \bar{v}_{\parallel}^2 \\ &= \frac{1}{\Sigma} \int dx_{\parallel} \int d^3\mathbf{v} v_{\parallel}^2 f - \frac{1}{\Sigma^2} \left(\int dx_{\parallel} \int d^3\mathbf{v} v_{\parallel} f \right)^2. \end{aligned} \quad (4.259)$$

We now integrate $\Sigma \overline{v_{\parallel}^2}$ over the sky and find that

$$\int d^2\mathbf{x} \Sigma \overline{v_{\parallel}^2} = \int d^2\mathbf{x} \Sigma (\sigma_{\parallel}^2 + \bar{v}_{\parallel}^2) = \int d^3\mathbf{x} \int d^3\mathbf{v} f v_{\parallel}^2 = \frac{2}{M} \sum_{ij} \hat{\mathbf{s}}_i K_{ij} \hat{\mathbf{s}}_j, \quad (4.260)$$

where M is the galaxy's mass, K_{ij} is defined by (4.240b), and $\hat{\mathbf{s}}$ is the unit vector in the direction of the line of sight. We shall find it convenient to introduce the notation $\langle q \rangle$ for the column-density weighted average of a quantity q over the sky. In terms of this notation (4.260) can be written

$$\langle \sigma_{\parallel}^2 \rangle + \langle \bar{v}_{\parallel}^2 \rangle = \frac{2}{M} \hat{\mathbf{s}} \cdot \mathbf{K} \cdot \hat{\mathbf{s}} \quad \text{where} \quad \langle q \rangle \equiv \int d^2\mathbf{x} \Sigma q. \quad (4.261)$$

Consider now the case of an axisymmetric galaxy that rotates around its symmetry axis (the z axis) and is seen edge-on. The x axis is taken to be parallel to the line of sight, so the yz plane is the sky plane. The tensors \mathbf{K} , \mathbf{T} , etc., are diagonal in these coordinates. Moreover, T_{xx} and T_{yy} are equal, T_{zz} vanishes because there is only azimuthal streaming motion, and $\hat{\mathbf{s}} \cdot \mathbf{K} \cdot \hat{\mathbf{s}} = K_{xx}$. The tensor virial theorem provides just two non-trivial equations

$$\begin{aligned} 2K_{xx} + W_{xx} = 0 \\ 2K_{zz} + W_{zz} = 0 \end{aligned} \quad \Rightarrow \quad \frac{M(\langle \sigma_{\parallel}^2 \rangle + \langle \bar{v}_{\parallel}^2 \rangle)}{\Pi_{zz}} = \frac{W_{xx}}{W_{zz}}. \quad (4.262)$$

As described in §4.1.2, the line-of-sight dispersion σ_{\parallel} has two components, one arising from the velocity-dispersion tensor $\sigma(\mathbf{x})$, and one from the variation in $\bar{\mathbf{v}}(\mathbf{x})$ along the line of sight. In the present notation equation (4.27) can be written

$$\Sigma \sigma_{\parallel}^2 = \int dx_{\parallel} \nu(\mathbf{x}) (\sigma_{xx}^2 + u^2), \quad (4.263)$$

where $u(\mathbf{x})$ is the difference between the component of $\bar{\mathbf{v}}(\mathbf{x})$ parallel to the line of sight and \bar{v}_{\parallel} . When this expression is integrated over the sky, we obtain

$$\frac{\Pi_{xx}}{M} = \langle \sigma_{\parallel}^2 \rangle - \int d^3\mathbf{x} \nu(\mathbf{x}) u^2. \quad (4.264)$$

We introduce the **global anisotropy parameter** δ to quantify the degree of deviation from isotropy by

$$\Pi_{zz} = (1 - \delta)\Pi_{xx}. \quad (4.265)$$

On using this equation to eliminate Π_{zz} from (4.262), we obtain

$$\frac{\langle \bar{v}_{\parallel}^2 \rangle}{\langle \sigma_{\parallel}^2 \rangle} = \frac{(1 - \delta)W_{xx}/W_{zz} - 1}{\alpha(1 - \delta)W_{xx}/W_{zz} + 1}, \quad (4.266a)$$

where

$$\alpha \equiv \frac{1}{\langle \bar{v}_{\parallel}^2 \rangle} \int d^3\mathbf{x} \nu(\mathbf{x}) u^2. \quad (4.266b)$$

α is a dimensionless number that does not depend on how rapidly the galaxy rotates but only on how the stellar density ν and streaming velocity $\bar{\mathbf{v}}$ vary within the (R, z) plane. For typical galaxy models, α lies in the range (0.05, 0.2) (Binney 2005).

The left side of equation (4.266a) can be determined from spectroscopic observations of an edge-on galaxy, while the ratio of components of \mathbf{W} that appears on the right side can be determined from photometry if we assume that the galaxy's light traces its mass. Hence we can use the equation to determine the global anisotropy parameter δ .

Furthermore, when a system's equidensity surfaces are similar spheroids, the components W_{ij} of its potential-energy tensor are given by (2.144), which implies that the ratio W_{xx}/W_{zz} depends only on the ellipticity ϵ of the spheroids, and is independent of the system's radial density profile. Hence, for such a system the right side of equation (4.266a) is a function of ϵ , α , and δ only. Figure 4.22 shows the resulting relation between the RMS rotation velocity and RMS velocity dispersion for two values of α (0.15, 0.1), and the marked values of δ . Also shown are the locations of 48 elliptical and lenticular galaxies from the SAURON survey (Cappellari et al. 2007), which has wide coverage of shapes and luminosities rather than being a photometrically complete sample. The great majority of the galaxies lie in the region $0 \leq \delta \leq 0.3$. In §4.4.2c we saw that as an oblate, rotating galaxy is moved from edge-on towards face-on position orientation, its representative point in a v/σ diagram moves to the left (Figure 4.14), and from Figure 4.22 it is clear that this movement will typically decrease the value of δ that we will infer for the galaxy. Hence, from Figure 4.22 we conclude that most elliptical galaxies have significantly anisotropic velocity-dispersion tensors, and that

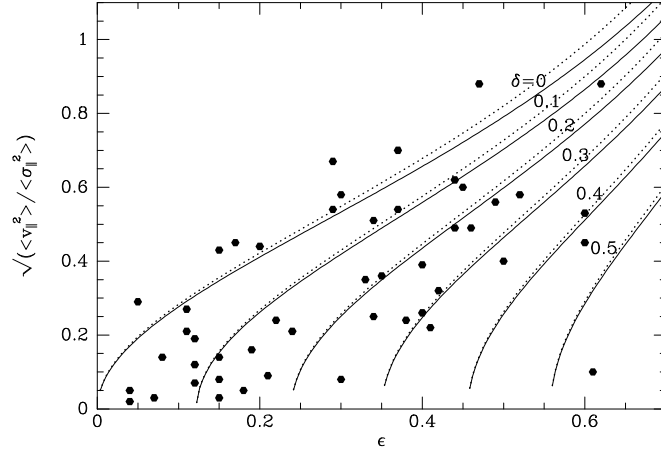


Figure 4.22 The ratio of the RMS line-of-sight streaming velocity to the RMS line-of-sight velocity dispersion in an edge-on spheroidal galaxy as a function of the ellipticity of the isodensity surfaces, for several values of the global anisotropy parameter δ . Full curves are for $\alpha = 0.15$ in equation (4.266a), while dotted curves are for $\alpha = 0.1$. The black dots show the locations of 48 luminous elliptical and lenticular galaxies from Cappellari et al. (2007).

this anisotropy is at least as important as rotation in determining the shapes of these objects. Naab, Jesseit, & Burkert (2006) discuss the implications of the measured values of δ for models of the formation of elliptical galaxies in mergers (§8.5).

4.9 Stellar kinematics as a mass detector

Knowing how mass is distributed in an astronomical system is fundamental to understanding what the system is, how it works, and how it formed. Since 1970 it has gradually emerged that “ordinary” baryonic matter, such as we are made of, contributes $\lesssim 20\%$ of the total matter of the universe (§1.3.5), with the rest consisting of some mysterious form of dark matter. It has also emerged that black holes with masses 10^6 to $10^9 \mathcal{M}_\odot$ reside at the centers of many galaxies (§1.1.6). The most convincing way to detect dark matter and black holes is through their gravitational fields. We have only two ways of detecting a gravitational field: through its action on photons (gravitational lensing; Schneider, Ehlers, & Falco 1999), and through its action on masses such as stars and gas clouds. Consequently, considerable effort has been expended on methods for inferring the gravitational field from the kinematics of a population of objects such as stars, globular clusters, or galaxies. These studies draw heavily on the theoretical tools that we have

assembled in this chapter. We discuss in turn, the detection of central black holes from stellar kinematics, probes for dark halos around elliptical galaxies, and measurement of the mass density in the solar neighborhood. The same physics underlies searches for both black holes and dark halos—we explain it in the context of black-hole searches.

4.9.1 Detecting black holes

Suppose that we wish to determine the mass of the black hole at the center of a spherical galaxy. From measurements of the surface brightness as a function of radius, we can determine the luminosity density $j(r)$. The pioneering work on this problem (Sargent et al. 1978) estimated the mass within radius r by setting $\beta = 0$ in equation (4.215) and writing

$$M(r) = -\frac{r}{G} \left(\sigma^2 \frac{d \ln j}{d \ln r} + \frac{d \sigma^2}{d \ln r} \right), \quad (4.267)$$

where $\sigma^2 = \overline{v_i^2}$ is the mean-square of any component of velocity, which can be determined from the line-of-sight velocity dispersion $\sigma_{\parallel}(R)$. If we were sure that β did indeed vanish, the equation would determine both the black-hole mass through $M_{\bullet} = \lim_{r \rightarrow 0} M(r)$, and the stellar mass $M_*(r) = M(r) - M_{\bullet}$. The problem with this methodology is that we have no guarantee that β vanishes, and Binney & Mamon (1982) showed that once β is allowed to be a free function of radius, a wide variety of mass profiles are consistent with given dispersion and surface-brightness profiles. In view of Figure 4.4 this result is not surprising: by varying β within modest limits, the central velocity-dispersion profile can be made to rise or fall at will, whether or not a black hole is present. Specifically, either inserting a black hole into an isotropic model, or making β turn sharply positive at small radii will cause $\sigma_{\parallel}(R)$ to turn upwards at small R (Figure 4.20); conversely, by causing β to turn negative at small radii we can mask the effect of a black hole.

To make progress it is vital to characterize the LOSVD with more than its dispersion, σ_{\parallel} . The left panels of Figure 4.23 show three LOSVDs for each of the models that contribute velocity-dispersion profiles to Figure 4.4. The right panels show the differences between each LOSVD and the Gaussian that has the same dispersion; the LOSVD is called “cuspy” if it lies above the equivalent Gaussian at zero velocity, and otherwise is said to be “flat-topped.” From top to bottom the panels show LOSVDs for the radially biased model ($\beta = \frac{1}{2}$), the isotropic model, and the tangentially biased model ($\beta = -\frac{1}{2}$). In the top left panel the LOSVD at $R = 4a$ for the radially biased model has a cusp at $v_{\parallel} = 0$ and lies above the equivalent Gaussian at small v_{\parallel} , as the top-right panel attests. The profile is cuspy because the LOSVD is dominated by the region in which the line of sight passes closest to the center and the density of stars is highest, and in this region stars on nearly

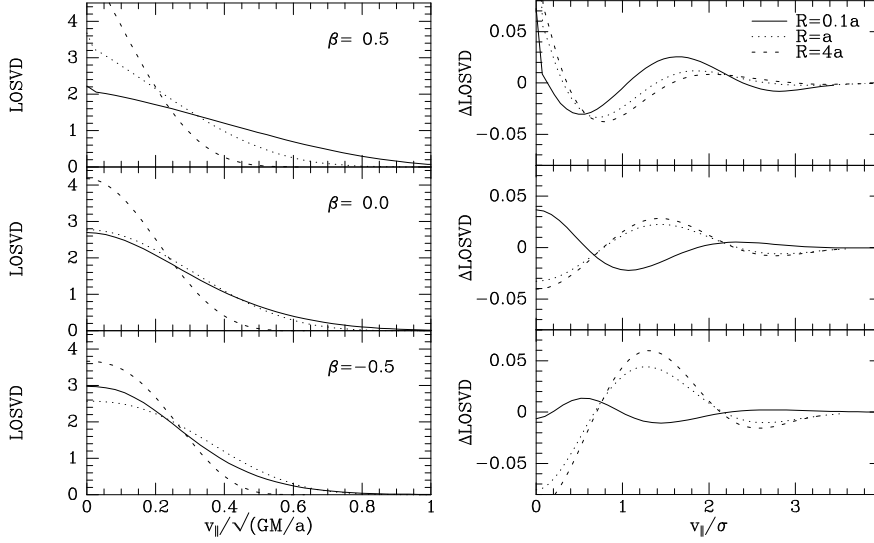


Figure 4.23 Left panels: LOSVDs for the Hernquist models plotted in Figure 4.4. From top to bottom the models have anisotropy parameter $\beta = \frac{1}{2}$, 0 and $-\frac{1}{2}$. In each panel profiles are shown for $R = 0.1a$, a and $4a$. The right panels show the deviations of each LOSVD from the Gaussian that has the same dispersion. From top to bottom the full curves have Gauss–Hermite parameters h_4 (BM §11.1.2) 0.001, 0.024 and 0.002; the dashed curves have $h_4 = 0.038$, -0.022 and -0.057 .

radial orbits have $v_{\parallel} \simeq 0$. In contrast, the bottom left panel shows that in the tangentially biased model the LOSVD at $R = 4a$ has a flat top, and the bottom-right panel shows that the LOSVD lies below the corresponding Gaussian at small v_{\parallel} . The model has a flat-topped profile because at the point of closest approach of the line of sight to the center, the multitude of stars on nearly circular orbits are seen at a variety of values of v_{\parallel} depending on the inclinations of their orbital planes to the line of sight.

The full curves, which show the LOSVDs at $R = 0.1a$, behave differently because these lines of sight are dominated by stars where the gradient in the stellar density is relatively shallow. Consequently the LOSVD is not so heavily dominated by the tangent point, and the radial velocities of stars make significant contributions to v_{\parallel} .

The shape of the LOSVD is frequently parameterized by the Gauss–Hermite coefficient h_4 (van der Marel & Franx 1993; Gerhard 1993; BM §11.1.2) and from the numbers given in the figure caption one sees that the radially biased model can be distinguished from the tangentially biased one if h_4 can be determined with an accuracy greater than ~ 0.03 . Such accuracy can be achieved only when the data have a high signal-to-noise ratio, and when the template star used in the analysis of the galaxy’s spectrum is well matched to the galaxy’s dominant stars (Houghton et al. 2006).

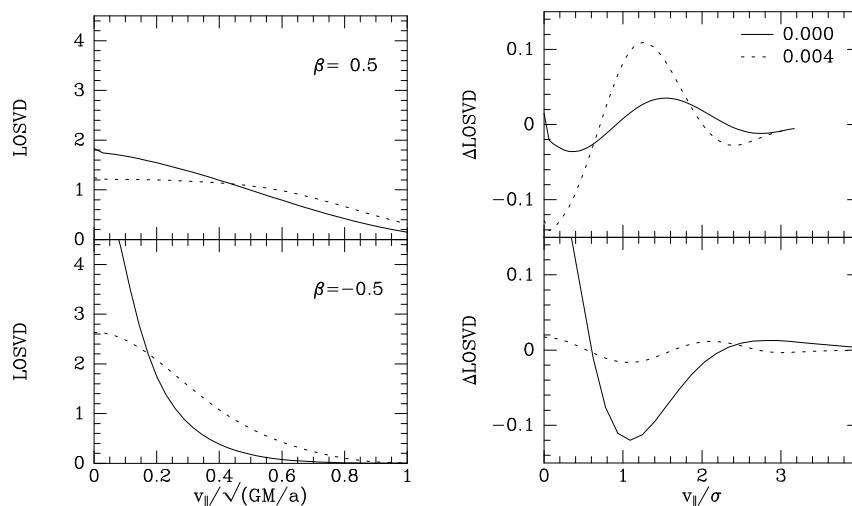


Figure 4.24 LOSVDs along $R = 0.01a$ through a Hernquist model that either does not (full line) or does (dotted curve) contain a black hole. The black-hole mass is a fraction 0.004 of the galaxy mass.

Figure 4.24 shows LOSVDs at $R = 0.01a$ for the radially (top panels) and tangentially biased models, both with and without a central black hole whose mass is 0.004 times the mass of the galaxy. The line-of-sight velocity dispersion profiles of these models are shown in Figure 4.20, which shows that the line of sight at $R = 0.01a$ passes well inside the black hole's radius of influence. The full curves in Figure 4.24 show the case of no black hole. The one in the top right panel shows that, in the absence of a black hole, the LOSVD of the radially biased model at this small projected radius is rather nearly Gaussian. The dashed line in the top left panel shows that when the black hole is added, the LOSVD becomes distinctly flat-topped, falling well below the Gaussian at small v_{\parallel} . This shape reflects the large contributions to v_{\parallel} from radial velocity at points that lie well in front of or behind the point of closest approach of the line of sight to the center. The full curve in the lower right panel shows that, in the absence of a black hole, the LOSVD of the tangentially biased model is sharply peaked around $v_{\parallel} = 0$; this shape reflects the tendency of nearly circular orbits with radii $r \gg R$ to cut the line of sight almost at right angles. The dashed curve in the bottom right panel shows that adding a black hole makes the LOSVD rather nearly Gaussian. The black hole reduces the number of stars seen at $v_{\parallel} \simeq 0$ by speeding up stars with apocenters at $r \lesssim R_{\text{infl}}$ (eq. 4.220).

These examples illustrate the following general effects:

- At $R > a$ the radial stellar gradient is large all along the line of sight, so radial bias makes the LOSVD cuspy while tangential bias makes it flat-topped.

- Along a line of sight that penetrates deeply into the central region in which the stellar density gradient is shallow, radial bias creates a Gaussian or flat-topped LOSVD while tangential bias makes the LOSVD cuspy.
- Along a line of sight that passes inside the black hole's radius of influence, the black hole makes the LOSVD more flat-topped than it otherwise would be, regardless of whether the model is radially or tangentially biased.

In a study of nearby elliptical galaxies (e.g., Gebhardt et al. 2003), one might have $R_{\text{infl}} \simeq 10$ pc and the galaxy might be at a distance $d \simeq 20$ Mpc. Then the angular size of R_{infl} will be 0.1 arcsec. This is only about a factor 2 larger than the best angular resolution achievable in photometry with the Hubble Space Telescope, and comparable to the best angular resolution at which spectroscopy can be done. Consequently, the deleterious effects of the point spread function on the data (BM §4.2.2) are important. In particular, the spectrum associated with the pixel that contains the galactic center will contain light emitted by stars that are very close to the black hole and moving at thousands of kilometers a second, as well as light from larger numbers of more distant stars that are moving much more slowly. Hence, the black hole will signal its presence by adding faint wings of large extent to the LOSVD. Data of the highest quality are required if these wings are to be detectable.

Elliptical galaxies exhibit a wide variety of surface brightness profiles near their centers (BM Figure 4.32). If at small radii $I(R) \propto R^{-\alpha}$, then the ease with which a central black hole can be detected from stellar kinematics increases with α because (i) the signal from stars near the black hole is less diluted by foreground and background stars when α is large, and (ii) it is easier to get high signal-to-noise data when the surface brightness is high. Galaxies have asymptotic slopes $0 \lesssim \alpha \lesssim 1$, and the Hernquist model lies at the lower end of this range, so in some respects we have focused on a particularly challenging example.

Several groups have hunted for central black holes in galaxies by fitting models to both the galaxy's photometry and the LOSVDs measured at different radii (e.g., Richstone & Tremaine 1985; van der Marel et al. 1998; Gebhardt et al. 2003). These studies have detected mass concentrations that are presumably black holes in many nearby galaxies, and have revealed the correlation (1.27) between black-hole mass and the velocity dispersion near the center of the galaxy. In most of this work the galaxies have been assumed to be axisymmetric, and Schwarzschild's modeling technique has been used to predict LOSVDs as a function of an assumed mass-to-light ratio Υ of the galaxy and an assumed black-hole mass M_{\bullet} . The values of Υ and M_{\bullet} are adjusted until a minimum is achieved in the χ^2 of the fit of the data to the observations. In general galaxies seem to exhibit mild radial biases, but there is a tendency for tangential bias to develop in the immediate vicinity of the black hole.

As we have seen, the detections of central black holes from stellar dynamics are challenging measurements: they require the highest possible spatial

resolution because the angular size of R_{infl} (eq. 4.220) is so small, and they require high-precision measurements of the shape of the LOSVD. For this reason, the two most secure black-hole masses are that of the Milky Way, which can be obtained by following the orbits of individual stars at distances of order 10^{-3} pc from the black hole (Eisenhauer et al. 2003), and that of NGC 4258, which is determined from observations of H_2O masers in the surrounding accretion disk (BM §7.2.4; Herrnstein et al. 2005). It is reassuring that (i) both objects lie on the standard mass-dispersion correlation (1.27), and (ii) black-hole mass measurements from gas kinematics define the same mass-dispersion correlation as measurements from stellar kinematics.

4.9.2 Extended mass distributions of elliptical galaxies

Measurements of the luminosities L , effective radii R_e , and central line-of-sight velocity dispersions σ_0 of spheroidal systems show that these variables are tightly correlated, such that spheroids lie on the “fundamental plane” in three-dimensional $(\log_{10} L, \log_{10} R_e, \log_{10} \sigma_0)$ space (§1.1.3a and BM §4.3.4). Photometric and spectroscopic measurements such as those we have discussed in connection with the determination of black-hole masses enable us to estimate mass as a function of radius, $M(r)$, within a given spheroidal system. In particular, recalling that half the galaxy’s light lies within R_e , we can determine the mass-to-light ratio $\Upsilon_e \equiv 2M(R_e)/L$ and ask how this varies within the fundamental plane. Crudely we expect from the scalar virial theorem (4.249b) that

$$\sigma_0^2 \approx \frac{GM(R_e)}{R_e} \simeq \frac{G\Upsilon_e L}{2R_e} = \pi G\Upsilon_e \bar{I} R_e, \quad (4.268)$$

where we have used the relation $L = 2\pi R_e^2 \bar{I}$ between the luminosity and \bar{I} , the mean surface brightness inside R_e . Taking logarithms we obtain

$$\log_{10} R_e = 2 \log_{10} \sigma_0 - \log_{10} \bar{I} - \log_{10} \Upsilon_e + \text{constant}. \quad (4.269)$$

Can this theoretical relation be reconciled with the empirical relation (1.20) that defines the fundamental plane?

Deviation of the data from the prediction of equation (4.269) is expected because high-mass spheroids are not simply scaled versions of low-mass spheroids, so σ_0^2 will not be precisely proportional to $GM(R_e)/R_e$ as equation (4.268) assumes. However, the violations of (4.269) expected from this cause are much smaller than those implied by (1.20), so we consider other ways of reconciling equation (4.269) with the fundamental plane.

The differences in the coefficients of $\log_{10} \sigma_0$ and $\log_{10} \bar{I}$ in (4.269) and (1.20) imply that Υ_e varies systematically with σ_0 and \bar{I} . If the masses of spheroids were mostly contributed by dark matter, these variations in

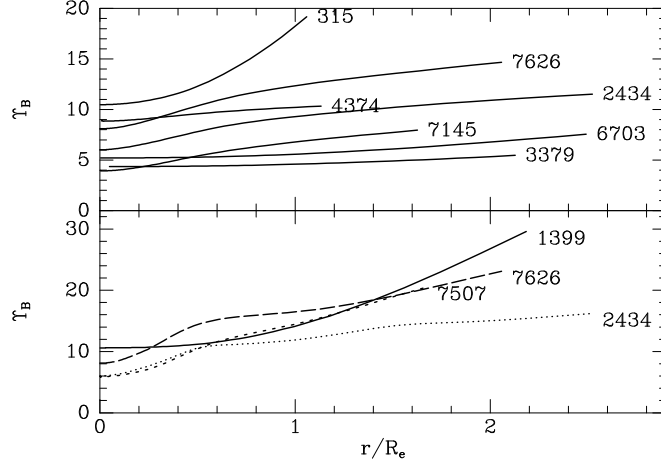


Figure 4.25 Upper panel: ratios $\Upsilon_B(r) = M(r)/L_B(r)$ of mass to light within radius r for seven elliptical galaxies with high-quality data—each curve is labeled by the galaxy’s NGC number. Lower panel: the local mass-to-light ratio $\rho(r)/j_B(r)$ in four such galaxies (after Gerhard et al. 2001).

Υ_e might arise from changes in the ratio of stars to dark matter. If, by contrast, the masses are mainly contributed by stars, the variations in Υ_e should reflect changes in characteristics, such as the age and metallicity, of the stellar population of each system.

Kronawitter et al. (2000) modeled photometry and spectroscopy of a sample of 21 round elliptical galaxies. The luminosity density of each galaxy was first determined from the photometry. Then the corresponding gravitational potential was calculated assuming some mass-to-light ratio Υ . To this was added a simple model $\Phi = \frac{1}{2}v_0^2 \ln(r^2 + r_c^2)$ for the additional contribution of dark matter to the potential. Then DFS $f_j(H, L)$ of the form (4.80) were found that produce the given luminosity density in the potential. Assuming that the actual DF was of the form $f = \sum_j a_j f_j$, σ_{\parallel} and h_4 were calculated at radii for which data were available. The coefficients a_j were used to optimize the fit to these data, and then the resulting value of $\chi^2(\Upsilon, v_0, r_c)$ was noted. The procedure was then repeated for different values of the stellar mass-to-light ratio Υ and the dark-halo parameters v_0 and r_c , to find the values that minimized χ^2 . The following conclusions emerged from this study (Gerhard et al. 2001). (i) Most galaxies show slight radial anisotropy, with $\beta \sim 0.2$. (ii) For $0.2 \lesssim r/R_e \lesssim 2$ the circular speeds vary by no more than $\sim 10\%$. (iii) The circular speeds define a Tully–Fisher relation $L \propto v_c^\alpha$ that has a similar, or slightly shallower, slope α to that in the relation (1.24) for spirals and a constant of proportionality that makes the ellipticals 0.6 mag fainter in the R band at a given value of v_c . (iv) Models in which all mass at $r \lesssim R_e$ is stellar are consistent with the data, although the spatial resolution

was not high enough to detect black holes at the centers of the galaxies. (v) The mass-to-light ratios Υ implied by these models are consistent with those predicted by stellar population models. In particular, the observed increase in Υ with L is consistent with the redder colors of more luminous galaxies. While a tendency for age to increase with L may contribute to the increase in Υ , the dominant effect is almost certainly increasing metal content. (vi) In some systems there is no evidence for dark matter out to $R = 2R_e$. In other galaxies dM/dL increases outwards by as much as a factor 3–4 between the center and $2R_e$ (Figure 4.25). A study by Cappellari et al. (2006) of 25 E and S0 galaxies additionally showed that (vii) the mass-to-light ratio inside R_e is tightly correlated with the velocity dispersion (eq. 1.23).

Beyond $R \sim R_e$ the low surface brightness of an elliptical galaxy makes it hard to obtain spectra that permit the reliable extraction of the LOSVD. As was described in §1.1.3a, we have three possible ways to probe the mass distributions of galaxies beyond R_e : (i) dynamics of test particles such as planetary nebulae and globular clusters; (ii) the hydrostatics of trapped X-ray emitting gas; (iii) weak gravitational lensing. Both X-ray observations and weak lensing show that many elliptical galaxies are surrounded by huge amounts of dark matter (Humphrey et al. 2006; Heymans et al. 2006). On the other hand, the kinematics of planetary nebulae around a few ellipticals are most straightforwardly interpreted as indicating that they do not have dark halos (Romanowsky et al. 2003). Ellipticals tend to exist in the most crowded areas of the universe (BM §4.1.2), so they are more likely than disk galaxies to have been affected by interactions with other galaxies (§8.2.2g). It is very possible that ellipticals that are not stationary at the centers of groups or clusters of galaxies have been largely stripped of their dark halos, while galaxies that sit at the center of a group or cluster have extensive halos that might be considered the property of the group or cluster rather than of the galaxy itself. Hence it is still unclear whether elliptical galaxies generally possess dark halos.

4.9.3 Dynamics of the solar neighborhood

A fundamental problem in Galactic structure, first studied by Kapteyn (1922), is to determine the density of matter in the disk near the Sun. The analysis of Oort (1932) and several subsequent workers was based on the Jeans equation (4.222b). The first term in this equation involves the mixed moment $\overline{v_R v_z}$. In our discussion of asymmetric drift (§4.8.2a) we saw that $\overline{v_R v_z}$ is probably smaller than, but on the order of, $(\overline{v_R^2} - \overline{v_z^2})(z/R)$. Thus if we assume, as in the discussion following equation (4.155), that $\overline{v_R^2}$ and $\overline{v_z^2}$ both decline with R as $\exp(-R/R_d)$, then we conclude that the first term in equation (4.222b) is constrained by

$$\left| \frac{1}{R} \frac{\partial(R\nu\overline{v_R v_z})}{\partial R} \right| \simeq \frac{2\nu}{R_d} \overline{v_R v_z} \lesssim \frac{2\nu z}{R_d} \frac{\overline{v_R^2} - \overline{v_z^2}}{R_0}. \quad (4.270)$$

The second term in equation (4.222b) is of order $\nu \overline{v_z^2}/z_0$, where $z_0 \ll R_0$ is the scale height of the disk. Hence the first term is smaller than the second by at least a factor $2zz_0/(R_d R_0) \lesssim 0.01$ and we may neglect it as Oort did. With this approximation we have

$$\frac{\partial(\nu \overline{v_z^2})}{\partial z} + \nu \frac{\partial \Phi}{\partial z} = 0. \quad (4.271)$$

This is the Jeans equation for a one-dimensional slab.

To relate the potential to the disk density we use Poisson's equation (3.88) in cylindrical coordinates. The first term vanishes when the circular-speed curve is flat, and it is very much smaller than the second term for any reasonable potential. Therefore we neglect it and have

$$4\pi G\rho = \frac{\partial^2 \Phi}{\partial z^2}; \quad (4.272)$$

once again the approximation for a one-dimensional slab. Combining equations (4.271) with (4.272), we now have

$$4\pi G\rho = -\frac{\partial}{\partial z} \left(\frac{1}{\nu} \frac{\partial(\nu \overline{v_z^2})}{\partial z} \right). \quad (4.273)$$

Here the mass density $\rho(z)$ is not necessarily proportional to the stellar density ν , which is that of any population of stars that is in a steady state.

If we could measure the run of density $\nu(z)$ and mean-square random velocity $\overline{v_z^2}(z)$ for any stellar population, from equation (4.273) we could read off the mass distribution $\rho(z)$. In practice statistical uncertainty in ν and $\overline{v_z^2}$ make it hard to estimate ρ reliably from (4.273).

An approach that relies on DFs rather than the Jeans equation probably provides a better way of determining the local mass density given the nature of the currently available observational data. At the end of §3.2.2 we saw that when a star's epicycle amplitude (the quantity X in eq. 3.91) is much smaller than the disk scale length R_d , the quantity H_z defined by equation (3.74) is an approximate isolating integral. By the Jeans theorem, we can assume that the DF that describes the distribution of stars in the (z, v_z) plane has the form $f(H_z)$. The density of stars is

$$\nu(z) = \int_{-\infty}^{\infty} dv_z f = 2 \int_0^{\infty} dv_z f\left(\frac{1}{2}v_z^2 + \Phi\right). \quad (4.274)$$

The function $\nu(z)$ we obtain by counting stars as a function of distance from the plane. The function $f(H_z)$ can be determined from the velocity distribution at $z = 0$. The fraction of stars in an interval dv_z is

$$P_0(v_z) dv_z = \frac{f\left(\frac{1}{2}v_z^2 + \Phi_0\right)}{\nu(0)} dv_z \quad \text{where} \quad \Phi_0 \equiv \Phi(0). \quad (4.275)$$

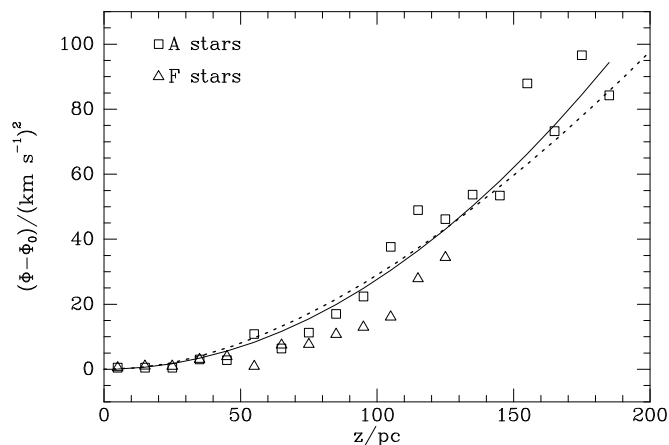


Figure 4.26 The change in gravitational potential near the Sun between the midplane and height z above the plane. Squares show values obtained with equations (4.276) from the A-star sample of Holmberg & Flynn (2000), while the triangles are from their F-star sample. The full curve is a least-squares fit of a parabola: it is the internal gravitational potential of a homogeneous slab with mass density $0.101 \mathcal{M}_{\odot} \text{pc}^{-3}$. The dashed curve shows the potential of Model I of §2.7.

Eliminating f between equations (4.274) and (4.275), we obtain

$$\frac{\nu(z)}{\nu(0)} = \frac{2}{\nu(0)} \int_0^{\infty} dv_z f\left[\left(\frac{1}{2}v_z^2 + \Phi - \Phi_0\right) + \Phi_0\right] = 2 \int_0^{\infty} dv_z P_0(u), \quad (4.276a)$$

where

$$u \equiv \sqrt{v_z^2 + 2(\Phi - \Phi_0)}. \quad (4.276b)$$

The right side of equation (4.276a) can be evaluated as a function of $\Phi - \Phi_0$, and then $\Phi(z) - \Phi_0$ is the value of $\Phi - \Phi_0$ at which the right side is equal to $\nu(z)/\nu(0)$.

The Hipparcos catalog,²⁰ which gives accurate three-dimensional locations and proper motions for stars in a sphere around the Sun of radius $\approx 120 \text{ pc}$, yields data to which we can apply equations (4.276). The catalog is complete only to a relatively bright apparent magnitude $V \sim 8$, and it is important that the volume around the Sun within which the selected stars lie is large enough to sample the vertical density profile $\nu(z)$. So one selects luminous, main-sequence stars ($M_V < 2.5$, $B - V < 0.6$), both because these stars are in the complete part of the catalog throughout the volume within which Hipparcos distances are accurate, and because they are young and therefore have a small scale height (Figure 8.11). Since the Hipparcos

²⁰ Hipparcos is the acronym of the HIgh Precision PARallax COLlecting Satellite. The acronym is inspired by the Greek astronomer Hipparchus.

catalog includes radial velocities only for a kinematically biased sample, P_0 has to be determined from proper motions alone.

The squares and triangles in Figure 4.26 show the values of $\Phi(z) - \Phi(0)$ that are obtained from equations (4.276) from two samples of tracer objects, namely A stars and F stars, that were extracted from the Hipparcos catalog by Holmberg & Flynn (2000). The density $\rho(z)$, which is our goal, is proportional to the double derivative of $\Phi(z)$, and it is clear from the figure that simple numerical differentiation will not yield credible results.

One way to proceed is to fit the data with a simple functional form. The simplest plan is to assume that $\Phi - \Phi_0 = \frac{1}{2}kz^2$, with k a parameter to be fitted to the data. The full curve in Figure 4.26 shows a least-squares fit of this formula to the data. The corresponding density is $\rho(0) = k/(4\pi G) = 0.101 \mathcal{M}_\odot \text{pc}^{-3}$. Using the same trial potential but a different approach to the observational data, Crézé et al. (1998) concluded that $\rho(0) = k/(4\pi G) = (0.076 \pm 0.015) \mathcal{M}_\odot \text{pc}^{-3}$. Holmberg & Flynn (2000) used a more complex trial potential but still one that depends only on an overall density parameter. They concluded that $\rho(0) = (0.10 \pm 0.01) \mathcal{M}_\odot \text{pc}^{-2}$. These values of $\rho(0)$ are in good agreement with each other and with the density expected to be contributed by stars and interstellar gas (Table 1.1). The dashed curve in Figure 4.26 shows the potential of Model I from §2.7, which is more physically plausible than the quadratic model shown by the full curve, but nonetheless provides a slightly less good fit to the data.

We have seen that the local density can be determined using only data in the Hipparcos catalog. Unfortunately, to determine the overall surface density $\Sigma = \int dz \rho(z)$ of the disk one has to use stars that are too distant to have accurate distances from Hipparcos. In practice one uses photometric distances to stars that lie towards the Galactic poles. In this direction, the line-of-sight velocity, which can be readily measured, is v_z , and the observations yield the distribution of stars in the (z, v_z) plane.

Kuijken & Gilmore (1989) obtained distances and line-of-sight velocities for 512 K dwarfs seen towards the south Galactic pole (BM §10.4.4). For various trial potentials $\Phi(z)$ they converted the measured density profile $\nu(z)$ into the functional dependence $\nu(\Phi)$. With ν expressed as a function of Φ , equation (4.274) becomes an Abel integral equation for $f(H_z)$ (eq. B.74):

$$\nu(z) = \int_{\Phi(z)}^{\infty} dH_z \frac{f(H_z)}{\sqrt{2(H_z - \Phi)}}. \quad (4.277)$$

The true dependence of Φ on z was estimated by varying the parameters in a simple functional form for $\Phi(z)$ until the likelihood of the measured distribution in the (z, v_z) plane was at a maximum. From this analysis Kuijken & Gilmore (1991) concluded that within 1.1 kpc of the plane the surface density of material is

$$\Sigma_{1.1}(R_0) = (71 \pm 6) \mathcal{M}_\odot \text{pc}^{-2}. \quad (4.278)$$

Holmberg & Flynn (2004) used a sample of K giants and an updated luminosity function for these stars and found that $\Sigma_{1.1}(R_0) = (74 \pm 6) \mathcal{M}_\odot \text{pc}^{-2}$, in agreement with (4.278).

Some part of $\Sigma_{1.1}$ must be contributed by the halo rather than the local disk. We can get an idea of the scale of this contribution by supposing that (i) the halo is spherical, (ii) the circular speed $v_c = v_0 = \text{constant}$, and (iii) without the halo, v_c would be falling in Keplerian fashion, $v_c = (GM_d/r)^{1/2}$ at R_0 . With these hypotheses, the halo mass $M(r)$ satisfies $G[M(r) + M_d] = rv_0^2$, so the halo density is

$$\begin{aligned} \rho_h &= \frac{1}{4\pi r^2} \frac{dM}{dr} = \frac{v_0^2}{4\pi G r^2} \\ &= 0.014 \mathcal{M}_\odot \text{pc}^{-3} \left(\frac{v_0}{200 \text{ km s}^{-1}} \right)^2 \left(\frac{R_0}{8 \text{ kpc}} \right)^{-2}, \end{aligned} \quad (4.279)$$

and the halo contributes only $\Sigma_{1.1}^h = (2.2 \text{ kpc}) \times \rho_h = 30.6 \mathcal{M}_\odot \text{pc}^{-2}$ to $\Sigma_{1.1}$. Thus $\Sigma_{1.1}$ must in fact be dominated by the disk not the halo, and the local disk must contribute more than $40 \mathcal{M}_\odot \text{pc}^{-2}$ to the local mass budget. The Sun lies in the transition region in which both disk and halo contribute significant mass.

4.10 The choice of equilibrium

One of the most important lessons of this chapter is that the range of equilibrium configurations accessible to a collisionless stellar system is large. The question arises “what determines the particular configuration to which a given stellar system settles?” Two classes of explanation are in principle possible. (i) The configuration actually adopted is favored or demanded by some fundamental physical principle, in the same way that the velocity distribution of an ideal gas always relaxes towards the Maxwell–Boltzmann distribution. (ii) The present configuration is a reflection of the particular initial conditions that gave rise to the system’s formation, in the same way that the shape of a particular stone in a field is due to particular circumstances rather than to any general physical principle. These two classes of explanation are not mutually exclusive. For example, the meandering course of the Mississippi River has been determined more by chance than by any fundamental principle, yet many characteristics of the river channel—typical sizes of elbows and the formation of oxbow lakes—can be understood in terms of simple physical arguments. Nevertheless, it is profitable to analyze the properties of galaxies in terms of first one and then the other of these points of view. We start by asking whether the present states of galaxies are simply more probable than any other configurations.

4.10.1 The principle of maximum entropy

In the 1890s J. W. Gibbs discovered that the standard relations between the thermodynamic variables of simple systems could be derived by hypothesizing that the probability that the system would be found to be in any small volume $d\tau$ of its phase space is proportional to $e^{-\beta H}d\tau$, where β is a parameter he identified as the inverse of the system's temperature, and H is the system's Hamiltonian. Since Gibbs's day there have been almost as many attempts to explain why this hypothesis works as there have been books written on statistical mechanics. After all these decades of debate, scientists are still far from agreement on this question. However, it is generally agreed that Gibbs's hypothesis can be derived from an alternative principle, that of **maximum entropy**: the thermodynamic relations for any physical system may be derived by seeking the probability density in phase space, p , that maximizes the **entropy**

$$S \equiv - \int_{\text{phase space}} d\tau p \ln p + \text{constant}, \quad (4.280)$$

subject to all relevant constraints. Can we derive the structures of galaxies from this principle?

The phase space of a galaxy of N stars is $6N$ -dimensional, and the infinitesimal $d\tau$ in equation (4.280) refers to an element of this phase space rather than an element of the phase space of a single star (§7.2.2). However, as was discussed in §4.1.1, we may neglect correlations between the particles of a collisionless system, so the probability $p d\tau$ associated with a range of configurations in the $6N$ -dimensional phase space of the whole galaxy is just the product of factors $f d^3\mathbf{x} d^3\mathbf{v}$ associated with individual stars. In these circumstances it is straightforward to show that

$$S = -N \int d^3\mathbf{x} d^3\mathbf{v} f \ln f + \text{constant}. \quad (4.281)$$

The obvious next step is to seek the form of f that maximizes S subject to given values of the galaxy's mass M and energy E . This calculation leads to the conclusion that S is extremized if and only if f is the DF (4.96) of the isothermal sphere (§7.3.2). However, we have seen that the isothermal sphere has infinite mass and energy, so the maximization of S subject to fixed M and E leads to a DF that is incompatible with finite M and E . The reason for this contradiction is that *no* DF with finite M and E maximizes S : if we constrain only M and E , configurations of arbitrarily large entropy can be constructed by suitable rearrangements of the galaxy's stars.²¹ The reason why this is so is easily explained.

²¹ It has sometimes been argued that the principle of increase of entropy dooms the universe to a "heat death," in which all matter is in a uniform, isothermal, maximum-entropy state. This argument is invalid once gravitational forces are included, since then there is no maximum-entropy state.

Suppose we have a spherical galaxy of total mass M and binding energy $|E|$. We mentally divide the system into a main body of mass M_1 and gravitational radius r_1 , and an outer envelope of radius r_2 that divides mass $M_2 \ll M$ among N_2 stars. By the virial theorem (4.250) the binding energy $|E_1|$ of the main body is $GM_1^2/2r_1$, and the binding energy of the envelope is $|E_2| \approx GM_1M_2/r_2$ since its stars orbit in a potential dominated by M_1 . Now imagine that we shrink the main body by a small fraction ϵ , so its radius changes from r_1 to $(1-\epsilon)r_1$. The shrinkage releases an energy $\Delta E \approx \epsilon GM_1^2/r_1$. We deposit this energy in the outer envelope, which swells in response to a radius r'_2 given by $|E'_2| = |E_2 + \Delta E| \approx GM_1M_2/r'_2$.

The velocity dispersion in the swollen envelope is $\sigma'_2 \approx \sqrt{GM_1/r'_2}$. Thus the volume $\mathcal{V} \approx \sigma_2^3 r_2^3$ of the region of phase space over which the representative points of the envelope's stars are distributed is $\mathcal{V} \approx (GM_1 r'_2)^{3/2}$ and the envelope's DF is $f \approx \mathcal{V}^{-1}$. Finally, by equation (4.281) the entropy of the envelope is

$$\begin{aligned} S &= -N_2 \int d^3\mathbf{x} d^3\mathbf{v} f \ln f + \text{constant} \approx N_2 \ln(\mathcal{V}) + \text{constant} \\ &\approx \frac{3}{2} N_2 \ln(r'_2) + \text{constant} \approx -\frac{3}{2} N_2 \ln(|E_2 + \Delta E|) + \text{constant}. \end{aligned} \quad (4.282)$$

It follows that the entropy of the envelope tends to infinity as ΔE tends to $|E_2|$. It is easy to see that the entropy of the main body changes by only a small amount as a result of the energy transfer. Hence the entropy of the combined system increases without limit. In other words, we can always increase the entropy of a self-gravitating system of point masses at fixed total mass and energy by increasing the system's degree of central concentration (page 30), because the entropy grows arbitrarily when a small fraction of the mass is put into a very large, diffuse envelope.

From this discussion we conclude that galaxies, unlike cold white-dwarf stars or an ideal gas, are not in thermodynamic equilibrium even though they are in dynamical equilibrium. Processes that disturb the dynamical equilibrium offer them the opportunity of moving to states of higher entropy, which will generally be characterized by a denser core and a more extensive envelope. These processes include spiral structure (Chapter 6), two-body relaxation (Chapter 7), and mergers (Chapter 8).

Because galaxies are not in thermodynamic equilibrium, if we wish to understand their present configurations we must investigate the initial conditions from which they started in life, and the rates at which various dynamical processes occur within them. The insights into initial conditions that cosmology provides are described in Chapter 9. Now we identify two physical mechanisms that play formative roles in the early lives of galaxies.

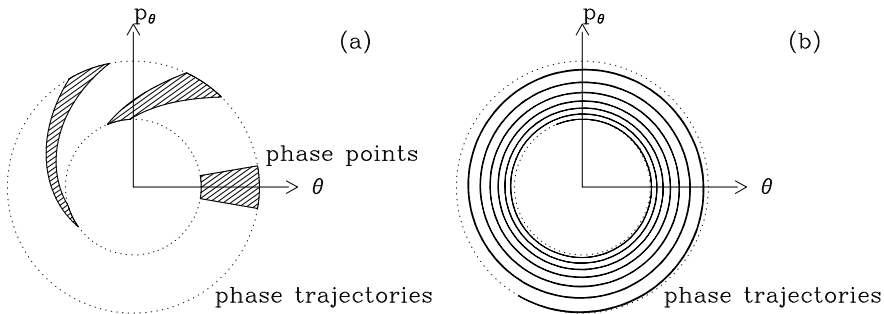


Figure 4.27 Schematic representation of how an initially compact group of phase points winds up into a larger region of lower coarse-grained phase-space density. (a) Initially the phase points fill the wedge on the axis $p_\theta = 0$. As time passes, and the phase-points move on circles like those shown, this wedge is drawn out into a band of ever-decreasing width. (b) After several crossing times the coarse-grained phase-space density is approaching uniformity in the annulus shown.

4.10.2 Phase mixing and violent relaxation

There are periods in the life of a galaxy when its material is far from a steady state, for example when it first forms by the collapse of a primordial density fluctuation (§9.2), or when it merges with another galaxy of comparable mass (§8.5). At these times of disequilibrium two related mechanisms come into play.

(a) Phase mixing Consider a collection of N independent pendulums that are all of the same length l , and therefore all have the same dynamical properties. Initially all the pendulums are swung back so they make angles θ with the vertical that are uniformly distributed in the interval $\theta_0 \pm \frac{1}{2}\Delta\theta$, where $\Delta\theta \ll \theta_0$. Each pendulum is given a small random velocity so the momenta $p = l^2\dot{\theta}$ are uniformly distributed in the range $\pm\Delta p$. When the pendulums are released, they start to oscillate. This situation is shown schematically in Figure 4.27—in reality the phase trajectories will be closed curves of constant energy, but not circles. The period of each pendulum depends on its amplitude, and therefore on its energy. The more energetic pendulums, which in Figure 4.27 move on circles of larger radius, oscillate slowest. Thus the patch formed in the figure by the phase-space points of the pendulums is gradually sheared out into a spiral of ever-diminishing pitch angle, which remains confined to the area between the curves of minimum and maximum energy.

The evolution of the whole system of pendulums may be described by the collisionless Boltzmann equation. According to this equation, the density of phase points in an infinitesimal volume around the phase point of any particular pendulum is constant. Consequently, the density of phase points in the spiral into which the occupied patch in Figure 4.27a is sheared is the same as the density in the original patch.

A macroscopic observer estimates the coarse-grained DF \bar{f} introduced in §4.1.1b by counting how many pendulums have phase-space coordinates in each of a number of cells of finite size. Initially $\bar{f} = f$, but at late times the spiral has been wound so tightly that any of the macroscopic observer's phase-space cells that intersects the spiral will contain both strips of the spiral and strips that are empty of phase points. When the observer works out the mean density within each of his cells, he finds that \bar{f} is constant throughout the annulus in phase space between the limiting energy curves. Since the area of the annulus is larger than the area of the small patch within which the phase points originally lay, \bar{f} is smaller than f .

The process that causes \bar{f} to decrease in this way as the pendulums get out of phase with one another is called **phase mixing**. An example is given on page 414.

The role of the collisionless Boltzmann equation in the relaxation of galaxies to equilibrium configurations is therefore rather subtle: it does not ensure that the empirically measurable DF \bar{f} is constant along stellar orbits, but rather it implies that \bar{f} can increase along an orbit only when the orbit mixes with other orbits that started from a higher phase-space density, and are themselves experiencing a decrease in \bar{f} . In particular the combination of the collisionless Boltzmann equation and phase mixing forces the maximum value of \bar{f} to decrease monotonically. The overall effect of regions of initially higher phase-space density mixing with regions of lower phase-space density is to reduce the fraction of the mass in the system that resides at values of \bar{f} larger than any given value.

The entropy defined by equation (4.281) is time-independent since f is constant along every orbit. However, if in equation (4.281) we replace f by \bar{f} , we obtain an entropy \bar{S} that in general increases in time; it can be shown that any decrease in the value of \bar{f} along orbits, such as occurs during phase mixing, causes \bar{S} to increase, just as the entropies of familiar thermodynamic systems increase when their different parts come into thermal equilibrium. In particular, we may state that no isolated collisionless stellar system A can evolve into a system B for which the entropy $\bar{S}(B) < \bar{S}(A)$. For developments of this idea, see Tremaine, Hénon, & Lynden-Bell (1986) and Dehnen (2005), and for its bearing on what dark matter might be, see Sellwood (2000a).

(b) Violent relaxation Phase mixing plays an important role during the relaxation of a galaxy of stars towards a steady state. But another relaxation process, known as **violent relaxation**, is also at work (Lynden-Bell 1967). While phase mixing changes the coarse-grained phase-space density near the phase point of each star, violent relaxation changes the energies of the stars themselves. When a star moves in a fixed potential Φ , its energy $E = \frac{1}{2}v^2 + \Phi$ is constant. But if Φ is a function $\Phi(\mathbf{x}, t)$ of both space and time, E is not constant. In fact (eq. D.10),

$$\frac{dE}{dt} = \left. \frac{\partial \Phi}{\partial t} \right|_{\mathbf{x}(t)}. \quad (4.283)$$

As a simple example, consider a star that is at rest at the center of a collapsing spherical protogalaxy. As the protogalaxy collapses, the potential well at its center becomes deeper. On the other hand, the velocity of the central star remains zero. Therefore the energy of this star decreases.

Other stars in the collapsing protogalaxy will gain energy. For example, consider a star that is initially located outside the half-mass radius of the system and is moving slowly radially outward. This star will be slow to respond to the overall collapse of the system, and by the time it is falling rapidly to the center, the system as a whole will be approaching its most compact configuration. Hence the star will acquire a lot of kinetic energy as it falls into the deep potential well at the center of the system. Later, by the time the star has passed close to the center and is on its way out again, the system will have re-expanded significantly and the potential well out of which it has to climb will be less deep than that into which it fell. Consequently, it will reach the potential at which it originally started with more kinetic energy than it had originally.

The change in the energy of a particular star during a collapse depends in a complex way on the initial position and velocity of that star—even in a spherical collapse, but more so in the generic collapse of a lumpy matter distribution—and the overall effect is to widen the range of energies of the stars. In this respect, a time-varying potential provides a relaxation mechanism analogous to collisions in a gas. However, there is an important distinction between relaxation in a gas and violent relaxation. Since the mass of the star whose energy is being followed does not appear in equation (4.283) or in the equation of motion that determines $\mathbf{x}(t)$, violent relaxation changes the energy per unit mass of a star that has a given initial position and velocity in a way that is independent of the star's mass. In contrast, we know from statistical mechanics that collisional relaxation tends to pump energy from the most massive particles to the least massive particles, thus establishing equipartition of energy (cf. page 583). It is sometimes desirable to be able to check that collisional relaxation through gravitational encounters is not playing an important role in a numerical simulation of a stellar system, and the best way to do this is to check that the distribution of energies per unit mass of stars is independent of mass.

The collisionless relaxation processes that we have described—phase mixing and violent relaxation—are distinct effects. For example, the coarse-grained phase-space density of an ensemble of non-interacting pendulums can be radically reduced by phase mixing although the energies of the individual pendulums are exactly constant. By contrast, while the phase-space density in the neighborhood of the phase point of a star that is at rest at the center of a collapsing protogalaxy does not change during a spherical collapse (the motion of stars that remain very close to the center of the protogalaxy is unaffected by the collapse), we have seen that violent relaxation does change the energy of these stars. Although these processes are conceptually distinct,

in §9.2.4 we shall see that in a collapsing system they work cooperatively, and indeed drive one another.

4.10.3 Numerical simulation of the relaxation process

With an N-body program (§3.4) we can investigate the relaxation of stellar systems experimentally. An N-body simulation is determined by the initial assignment of positions and velocities to its stars. Most simulations start from initial conditions that may be grouped into four broad categories.

- (i) At $t = 0$ the particles are distributed through a sharply bounded region with total energy $E = T + W < 0$, where T and W are the kinetic and potential energies. The outcome of such a **collapse simulation** is heavily influenced by the **virial ratio** $2T/|W|$. By the virial theorem (4.250) $2T/|W|$ will tend to unity, and the smaller it is initially the more violently the simulation will relax. The goal is to determine how the final equilibrium state depends on the initial conditions.
- (ii) At $t = 0$ the particles form an equilibrium axisymmetric disk. Each particle moves in the tangential direction with a velocity that is close to that required for centrifugal support of this disk, and has a small random radial motion. The goal of these **disk simulations** is to study the stability and long-term evolution of galactic disks.
- (iii) At $t = 0$ the particles are grouped into two galaxies. The positions and velocities of the particles are chosen so that each galaxy, treated as an isolated system, is in a steady state, and is approaching the other galaxy. Thus the galaxies are set to collide with one another. The goal of these **merger simulations** is to understand how galaxies merge and what merger remnants look like.
- (iv) At $t = 0$ the particles are nearly homogeneously distributed throughout a spherical volume, and are receding from the sphere's center with velocities that are approximately proportional to radius. These **cosmological simulations** model galaxy formation by gravitational clustering in an expanding universe.

We shall discuss disk simulations in §6.3.1, merger simulations in §8.5, and cosmological simulations in §9.3. Here we discuss collapse simulations, which mimic aspects of the formation of galaxies.

Figures 4.28 to 4.32 illustrate general features of collapse simulations following van Albada (1982), who studied the collapse of cold, clumpy distributions of stars. The initial conditions are obtained by first distributing 50 000 stars of unit total mass uniformly in a sphere of unit radius. Each star is then displaced from \mathbf{x} to $\mathbf{x} + \nabla\psi$ along the gradient of a Gaussian random field $\psi(\mathbf{x})$ —see §9.1.1 for an account of such fields. The initial velocity of the star is $\nabla\psi/t_{\text{ff}}$, where t_{ff} is the free-fall time (Problem 3.4) of the original sphere; the idea here is that both the displacement $\nabla\psi$ and the initial velocity result from the gravitational forces generated by small initial

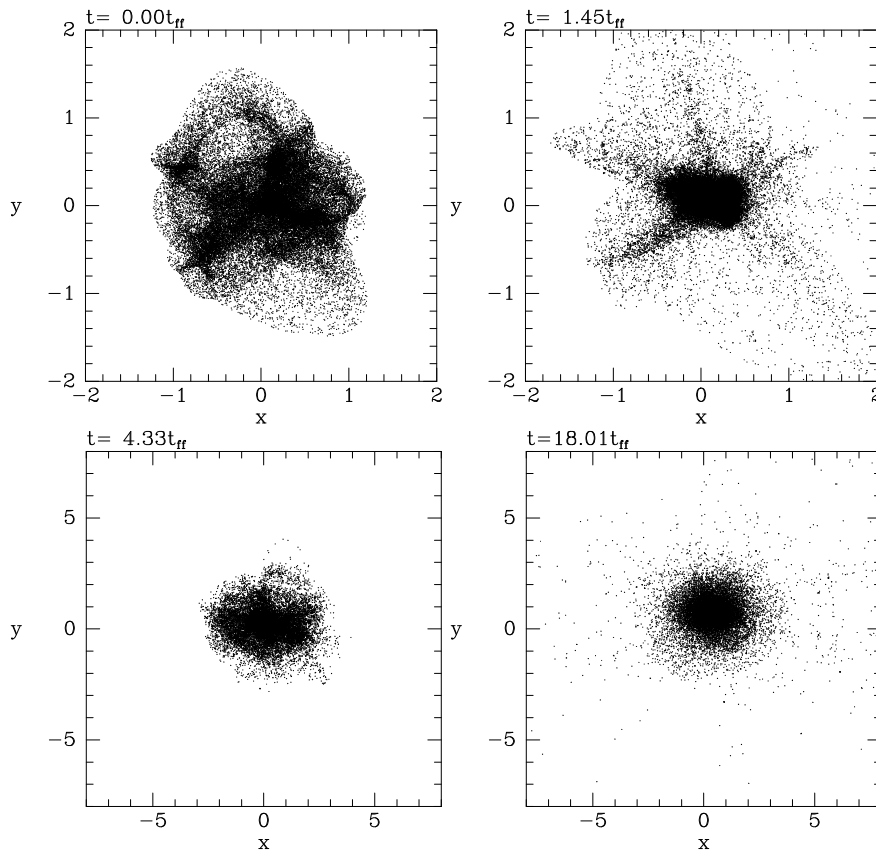


Figure 4.28 Four stages of a collapse simulation. Time is shown at the top of each panel in units such that the initial free-fall time is $\pi/2^{3/2} \simeq 1.1$ (Problem 3.4 with $GM = r = 1$). Top left: initially the 50 000 particles have positions and velocities obtained by shifting them from a homogeneous distribution within the unit sphere. Top right: gravity causes the system to fragment into lumps that fall together to form a tight minimum configuration. Bottom left: after several pulsations of ever-decreasing intensity, the core has settled to a quasi-steady state. Bottom right: after a much longer time a low-density halo of violently ejected stars is in place. Notice that the linear scale of the lower panels is four times that of the upper ones.

inhomogeneities. These forces have acted throughout the time t_{ff} required for the sphere to expand with the universe from infinite initial density to unit radius, and the simulation starts as the sphere is beginning to collapse. This prescription for choosing initial conditions approximately reproduces the results of the more sophisticated theory of structure formation that we shall outline in Chapter 9. The amplitude of ψ controls the initial virial ratio $2T/|W|$.

Figure 4.28 shows the spatial structure at four evolutionary stages. The

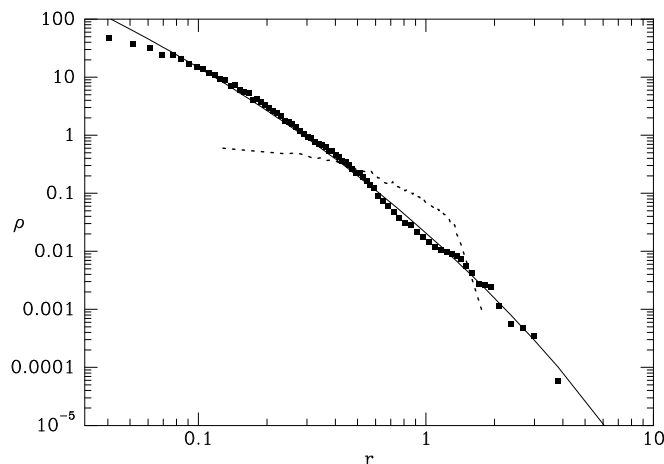


Figure 4.29 The radial density profile of the final configuration of the simulation shown in Figure 4.28 (squares), together with the simulation’s initial density profile (dashed line) and the density profile of the $R^{1/4}$ model that has the same half-mass radius (full curve).

system first contracts to a very compact configuration on a timescale t_{ff} , and then partially re-expands. After the elapse of two or three times t_{ff} , the center of the system has settled to a nearly steady state, and over time significant oscillations become confined to more and more peripheral regions. The smaller $2T/|W|$ is, the more pronounced is the initial collapse, and the higher the final central density is; in fact the final central density will be comparable to the density of the most compact configuration. Smaller values of $2T/|W|$ also lead to more mass being flung out into an extensive, low-density halo that takes significant time to come into equilibrium—this halo emerges between the bottom-left panel of Figure 4.28 at $t = 4.3t_{\text{ff}}$ and the bottom-right panel at $t = 18t_{\text{ff}}$.

Figure 4.29 shows the initial (dashed curve) and final (squares) density profiles of the simulation of Figure 4.28 together with the density profile of the $R^{1/4}$ model (full curve and eq. 1.17 with $m = 4$) that has the same half-mass radius $r_{\text{h}} = 1.35R_{\text{e}}$ as the simulation. This particular simulation, which started from $2T/|W| = 0.2$, generates a system that fits the $R^{1/4}$ model extremely well (van Albada 1982). Simulations with smaller values of $2T/|W|$ can be better fitted by the NFW profile (§2.2.2f).

Figure 4.30 shows the evolution of the differential energy distribution $N(E)$ that was introduced in §4.3.1b. Initially the energies of the particles lie in a narrow range, but this range is rapidly extended as violent relaxation causes particles to gain and lose energy. In the final configuration the most densely populated energies lie near the escape energy $E = 0$. This is just what we found to be the case in equilibrium models that are based on the Jeans theorem (cf. Figure 4.3).

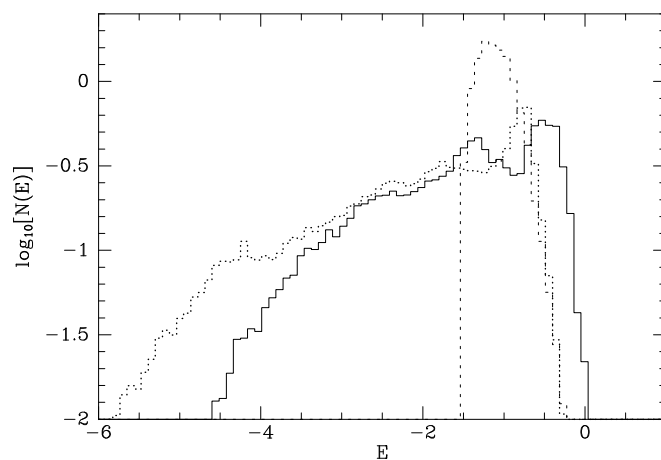


Figure 4.30 The evolution of the differential energy distribution of the model shown in Figure 4.28 at $t = 0$ (dashed curve), at $t = 1.45t_{\text{ff}}$ (dotted curve) and $t = 18t_{\text{ff}}$ (full curve). Energy is measured in units of GM/R_0 , where M is the system's mass and R_0 is the radius of the initial particle distribution.

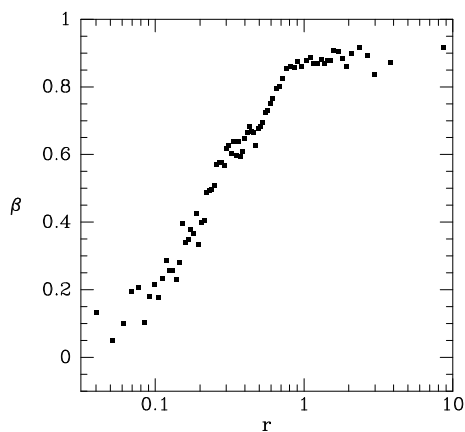


Figure 4.31 Anisotropy parameter $\beta \equiv 1 - \sigma_{\theta}^2/\sigma_r^2$ in the final configuration of the simulation shown in Figure 4.28.

As Figure 4.31 illustrates, the outer parts of the equilibrium system are strongly radially biased. This situation arises because the particles that populate the outer regions of the final system were accelerated onto their present orbits by fluctuations in the strong gravitational field that prevails near the galactic center. Consequently, most of these particles are on highly elongated orbits that pass close to the center of the system.

If the initial configuration is triaxial rather than spherical, the final system is also triaxial (Aarseth & Binney 1978). Figure 4.32 illustrates this phenomenon by showing the endpoint of a simulation that started from the triaxial configuration that is obtained by stretching the initial configuration

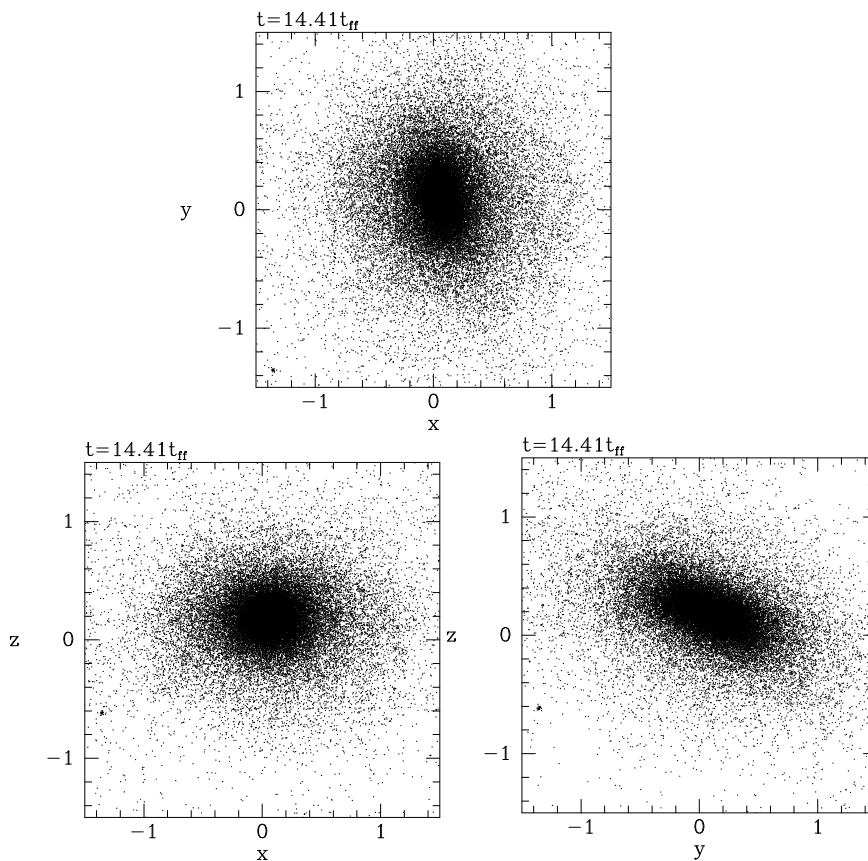


Figure 4.32 Three orthogonal projections of the final configuration to which 50 000 particles settled after they had been released from a cold, elliptical initial configuration. The initial conditions were generated by changing the initial conditions used to make Figure 4.28 from (\mathbf{x}, \mathbf{v}) to $(\mathbf{x} + \nabla\chi, \mathbf{v} + \nabla\chi/t_{ff})$, where $\chi \equiv 0.075(3x^2 - y^2 - 2z^2)$.

of Figure 4.28 parallel to the x axis and compressing it in the perpendicular directions. The velocity-dispersion tensors of ellipsoidal models formed in this way are everywhere anisotropic, in contrast to the case of spherical systems, in which the dispersion tensor is isotropic at small radii—at the center of the system shown in Figure 4.32 the principal velocity dispersions are in the ratios 1 : 0.80 : 0.68. When the initial configuration is slowly rotating, the figure of the final configuration also rotates (Wilkinson & James 1982). If the initial configuration is spherical but slowly rotating, the final system will be an oblate figure of rotation. When the initial rotation is rapid, the final state will be a prolate triaxial bar (Hohl & Zang 1979; Miller & Smith 1979).

Problems

4.1 [1] Show that in a frame that rotates with constant angular velocity $\boldsymbol{\Omega}$, with $\Phi_{\text{eff}} \equiv \Phi - \frac{1}{2}|\boldsymbol{\Omega} \times \mathbf{r}|^2$, the collisionless Boltzmann equation can be written

$$\frac{\partial f}{\partial t} + \mathbf{v} \cdot \nabla f - [2(\boldsymbol{\Omega} \times \mathbf{v}) + \nabla \Phi_{\text{eff}}] \cdot \frac{\partial f}{\partial \mathbf{v}} = 0. \quad (4.284)$$

Hint: see §3.3.2.

4.2 [2] Consider an infinite homogeneous system of collisionless zero-mass test particles in D -dimensional space. The particles have an isotropic velocity distribution $f(v)$. Initially the particles are subject to no forces. At $t = 0$ a gravitational potential well suddenly appears in a finite region of the space. Show that as $t \rightarrow \infty$, the density of unbound particles traveling through the well is smaller than the asymptotic density if $D = 1$, larger if $D = 3$, and unchanged if $D = 2$.

4.3 [1] A spherical mass distribution is immersed in a sea of collisionless test particles, which arrive with velocities $\mathbf{v} = (v, 0, 0)$ from the negative x -direction and are scattered by the gravitational field from the mass. Does the DF of the test particles satisfy the Jeans theorem? If so, write down the DF as a function of the integrals of motion; if not, explain why the Jeans theorem fails.

4.4 [1] Prove that the density of a spherical, ergodic, self-consistent stellar system must decrease outward. Hint: in the integral for ρ make Φ the integration variable.

4.5 [3] A spherical galaxy with a constant mass-to-light ratio and an ergodic DF has an $R^{1/4}$ surface-brightness profile (eq. 1.17 with $m = 4$). The luminosity-weighted line-of-sight velocity dispersion within the effective radius R_e is σ_e , that is

$$\sigma_e^2 = \frac{\int_0^{R_e} dR RI(R) \sigma_{\parallel}^2(R)}{\int_0^{R_e} dR RI(R)}. \quad (4.285)$$

Show that the total mass of the galaxy can be written in the form

$$M = k \frac{\sigma_e^2 R_e}{G}, \quad (4.286)$$

and evaluate k numerically.

4.6 [2] The DF of a spherical system is proportional to $L^\gamma f(\mathcal{E})$. Show that at all radii the anisotropy parameter is $\beta = -\frac{1}{2}\gamma$.

4.7 [1] Show that a Hernquist model with constant anisotropy $\beta = \frac{1}{2}$ has

$$N(E) = \frac{3a}{GM} \tilde{\mathcal{E}}^2 \left(\frac{1}{\tilde{\mathcal{E}}} - 1 \right)^2, \quad (4.287)$$

where $\tilde{\mathcal{E}} = \mathcal{E}a/(GM)$ and M and a are the mass and scale radius of the Hernquist model.

4.8 [2] Consider a spherical system with DF $f(\mathcal{E}, L)$. Let $N(\mathcal{E}, L)d\mathcal{E}dL$ be the fraction of stars with \mathcal{E} and L in the ranges $(\mathcal{E}, \mathcal{E} + d\mathcal{E})$ and $(L, L + dL)$.

(a) Show that

$$N(\mathcal{E}, L) = 8\pi^2 L f(\mathcal{E}, L) T_r(\mathcal{E}, L), \quad (4.288)$$

where T_r is the radial period defined by equation (3.17).

(b) A spherical system of test particles with ergodic DF surrounds a point mass. Show that the fraction of particles with eccentricities in the range $(e, e + de)$ is $2e de$.

4.9 [1] Consider the DF

$$f(\mathcal{E}, L) = \begin{cases} F\delta_+(L^2)(\mathcal{E} - \mathcal{E}_0)^{-1/2} & (\mathcal{E} > \mathcal{E}_0) \\ 0 & (\mathcal{E} \leq \mathcal{E}_0), \end{cases} \quad (4.289a)$$

where F and \mathcal{E}_0 are constants. Here $\delta_+(x) \equiv \delta(x - \epsilon)$, where ϵ is a small positive number; thus $\delta_+(x) = 0$ for $x \neq 0$ and $\int_0^\infty dx \delta_+(x) = 1$ (cf. Appendix C.1). Show that this DF self-consistently generates a model with density

$$\rho(r) = \begin{cases} Cr^{-2} & (r < r_0) \\ 0 & (r \geq r_0), \end{cases} \quad (4.289b)$$

where C is a constant and $\Psi(r_0) = \mathcal{E}_0$ (Fridman & Polyachenko 1984).

4.10 [2] Consider a spherical system in which at every radius the star density in velocity space is constant on ellipsoidal figures of rotation. Show that the DF has the form $f = f(Q)$, where $Q(\mathcal{E}, L)$ is defined by equation (4.73).

4.11 [2] Prove that the following DF generates a stellar system in which the density distribution is that of a homogeneous sphere of density ρ and radius a :

$$f(E, L) = \frac{9}{16\pi^4 G\rho a^5} \frac{1}{\sqrt{L^2/a^2 + \frac{4}{3}\pi G\rho a^2 - 2E}} \quad (L^2 < \frac{4}{3}\pi G\rho a^4). \quad (4.290)$$

Here it is understood that $f = 0$ when the argument of the square root is not positive, the DF is normalized so that $\int d^3\mathbf{x}d^3\mathbf{v} f = 1$, the potential $\Phi = 0$ at $r = 0$, and the system is isolated (Polyachenko & Shukhman 1973).

4.12 [2] Show that the Osipkov–Merritt model that self-consistently generates the Jaffe model has DF

$$f(Q) = \frac{1}{2\pi^3(GMa)^{3/2}} \left[F_- \left(\sqrt{2\tilde{Q}} \right) - \sqrt{2}F_- \left(\sqrt{\tilde{Q}} \right) - \sqrt{2} \left(1 + \frac{a^2}{2r_a^2} \right) F_+ \left(\sqrt{\tilde{Q}} \right) + \left(1 + \frac{a^2}{r_a^2} \right) F_+ \left(\sqrt{2\tilde{Q}} \right) \right], \quad (4.291)$$

where $\tilde{Q} = aQ/(GM)$.

4.13 [2] Show that when the DF of a spherical system depends only on the function $Q(\mathbf{x}, \mathbf{v})$ defined by equation (4.73), the ratio of the mean-square tangential and radial speeds is

$$\frac{\overline{v_t^2}}{\overline{v_r^2}} = \frac{2}{1 + (r/r_a)^2}. \quad (4.292)$$

4.14 [1] A self-consistent stellar system has an ergodic DF and a power-law density profile $\rho = \rho_0(r_0/r)^\alpha$ with $1 < \alpha < 3$. Show that the velocity dispersion is given by

$$\frac{\overline{v_r^2}(r)}{\overline{v_t^2}(r)} = \frac{2\pi G\rho_0 r_0^\alpha r^{2-\alpha}}{(3-\alpha)(\alpha-1)} \quad (\alpha \neq 2). \quad (4.293)$$

What does this formula become in the case $\alpha = 2$ of the singular isothermal sphere?

4.15 [1] Solve the Lane–Emden equation (4.91b) for the case $n = 1$, to show that

$$\psi = \begin{cases} \frac{\sin \sqrt{3}s}{\sqrt{3}s} & (s < \pi/\sqrt{3}) \\ \frac{\pi}{\sqrt{3}s} - 1 & (s \geq \pi/\sqrt{3}). \end{cases} \quad (4.294)$$

Show that the model's total mass is $M = \frac{1}{2}\Psi_0 G^{-3/2} \sqrt{\pi/c_1}$, where c_1 is defined by equation (4.85b).

4.16 [1] Consider a stellar system having a DF of the form $f \propto \mathcal{E}^{-1/2}$ for $\mathcal{E} > 0$ and zero otherwise. The density and potential are related by $\rho = c_1 \Psi$ for $\Psi > 0$, and zero otherwise (eq. 4.85). Prove that Poisson's equation is satisfied if the density of the system has the form

$$\rho(x, y, z) = A \cos\left(\frac{\pi x}{2L}\right) \cos\left(\frac{\pi y}{2L}\right) \cos\left(\frac{\pi z}{2L}\right), \quad (4.295)$$

for $|x|, |y|, |z| \leq L$ and zero otherwise, where $L^2 = 3\pi/(16Gc_1)$. Does this mean we can construct a cubical galaxy?

4.17 [2] An extension to the polytropes described in §4.3.3a is obtained by considering spherical stellar systems with the DF

$$f(E) = FE^{-n-3/2} \quad (E \geq 0); \quad (4.296)$$

here f is defined so that $\int d^3\mathbf{v} f = \rho$ and the potential is defined so that $\Phi = 0$ at the center of the system.

(a) Show that the density satisfies

$$\rho = d_n \Phi^{-n} \quad (\Phi > 0), \quad (4.297)$$

and evaluate d_n . What values of n are allowed?

(b) Show that the dimensionless radius $s \equiv r/b$ and potential $\phi \equiv \Phi/\Phi_0$ satisfy the equation

$$\frac{1}{s^2} \frac{d}{ds} \left(s^2 \frac{d\phi}{ds} \right) = 3\phi^{-n}, \quad (4.298)$$

where Φ_0 is arbitrary and $b \equiv (\frac{4}{3}\pi G \Phi_0^{-n-1} d_n)^{-1/2}$. What is the stellar system described by these equations in the limit $n \rightarrow \infty$?

(c) Show that these equations admit power-law solutions of the form $\rho \propto r^{-\alpha}$ for $n > 0$, with $\alpha = 2n/(1+n)$ so $0 < \alpha \leq 2$.

(d) In Problem 4.14 we found self-consistent power-law stellar systems with $1 < \alpha < 3$. Why does the current approach not find the systems with $2 < \alpha < 3$? Why does the approach in Problem 4.14 not find the solutions with $0 < \alpha \leq 1$?

4.18 [1] For a Maxwellian distribution of velocities with one-dimensional dispersion σ , show that: (a) the mean speed is $\bar{v} = (8/\pi)^{1/2}\sigma$; (b) the mean-square speed is $\overline{v^2} = 3\sigma^2$; (c) the mean-square of one component of velocity is $\overline{v_x^2} = \sigma^2$; (d) the mean-square relative speed of any two particles is $\overline{v_{\text{rel}}^2} = 6\sigma^2$; (e) the fraction of particles with $v^2 > 4\overline{v^2}$ is 0.00738.

4.19 [1] At large radii, the density in a Michie model (eq. 4.117) is dominated by stars with $\mathcal{E}/\sigma^2 \ll 1$. In this case, show that the density can be written in the form

$$\rho \propto \int_0^{\sqrt{2\Psi}} dv_t v_t \exp\left(-\frac{r^2 v_t^2}{2r_a^2 \sigma^2}\right) (2\Psi - v_t^2)^{3/2}. \quad (4.299)$$

Hence show that as the tidal radius r_t is approached, ρ tends to zero as

$$\rho \propto \Psi^{5/2} \propto (r_t - r)^{5/2}. \quad (4.300)$$

4.20 [1] Show that in a Kepler potential, the Schwarzschild DF (4.153) is equivalent to the **Rayleigh distribution** of eccentricities and inclinations,

$$dn \propto ei \exp\left(-\frac{e^2}{e_0^2} - \frac{i^2}{i_0^2}\right) de di. \quad (4.301)$$

What is the relation of e_0^2 to σ_R^2 and i_0^2 to σ_z^2 ?

4.21 [2] We may study the vertical structure of a thin axisymmetric disk by neglecting all radial derivatives and assuming that all quantities vary only in the coordinate z normal to the disk. Thus we adopt the form $f = f(E_z)$ for the DF, where $E_z \equiv \frac{1}{2}v_z^2 + \Phi(z)$. Show that for an isothermal disk in which $f = \rho_0(2\pi\sigma_z^2)^{-1/2} \exp(-E_z/\sigma_z^2)$, the approximate form (2.74) of Poisson's equation may be written

$$2 \frac{d^2\phi}{d\zeta^2} = e^{-\phi}, \quad \text{where } \phi \equiv \frac{\Phi}{\sigma_z^2}, \quad \zeta \equiv \frac{z}{z_0}, \quad \text{and } z_0 \equiv \frac{\sigma_z}{\sqrt{8\pi G\rho_0}}. \quad (4.302a)$$

By solving this equation subject to the boundary conditions $\phi(0) = \phi'(0) = 0$, show that the density in the disk is given by (Spitzer 1942)

$$\rho(z) = \rho_0 \operatorname{sech}^2(\frac{1}{2}z/z_0). \quad (4.302b)$$

Show further that the surface density of the disk is

$$\Sigma = \frac{\sigma_z^2}{2\pi G z_0} = 4\rho_0 z_0. \quad (4.302c)$$

4.22 [3] Determine the density $\rho(z)$ of an isothermal distribution of stars with dispersion σ_z in a disk that also contains a razor-thin layer of gas in the midplane, with surface density Σ_g . Hint: generalize the results of Problem 4.21.

4.23 [2] Using the one-dimensional approximation of Problem 4.21, write a numerical procedure that finds the fraction $F(z)$ of stars that reach a maximum height above the midplane that exceeds z . Show that $F(z_0) = 0.808$. Find $F(2z_0)$.

4.24 [3] Every star in a spherical system loses mass slowly and isotropically. If the initial DF is $f_0(\mathcal{E}, L)$, show that after every star has been reduced to a fraction p of its original mass, the DF will be

$$f_p(\mathcal{E}, L) = f_0(p^{-2}\mathcal{E}, L). \quad (4.303)$$

How has the density profile of the system changed? Hint: see Richstone & Potter (1982).

4.25 [2] A spherical stellar system has surface brightness $I(R)$, line-of-sight velocity dispersion $\sigma_{\parallel}(R)$, and constant mass-to-light ratio. From these functions, together with the Jeans equations, can we deduce uniquely the luminosity density $j(r)$, the radial dispersion profile, $\sigma_r^2(r)$, and the anisotropy parameter $\beta(r)$? Hints: (i) consider a change in the dispersions of the form $\Delta\sigma_r^2 = 2\epsilon/(jr^3)$, $\Delta\sigma_{\theta}^2 = -\frac{1}{2}\Delta\sigma_r^2$; (ii) the answer may depend on whether $I(R)$ and $\sigma_{\parallel}(R)$ are known over a finite range of radii or at *all* radii (see Dejonghe & Merritt 1992).

4.26 [3] A finite, spherical stellar system of test particles is confined by the potential $\Phi(r) = v_c^2 \ln r$.

(a) If the DF is ergodic, prove that the number of stars with line-of-sight velocity in the interval $(v_{\parallel}, v_{\parallel} + dv_{\parallel})$ is $n(v_{\parallel}) dv_{\parallel}$, where

$$n(v_{\parallel}) \propto \exp\left(-3v_{\parallel}^2/2v_c^2\right). \quad (4.304)$$

(b) If the DF is radial, that is, if all of the test particles have zero angular momentum, prove that the distribution of line-of-sight velocities is given by

$$n(v_{\parallel}) \propto E_1\left(v_{\parallel}^2/2v_c^2\right), \quad (4.305)$$

where E_1 is the exponential integral.

4.27 [1] Show that a self-gravitating isothermal stellar system with velocity dispersion σ , cylindrical symmetry, and non-singular, non-zero density ρ_0 at $R = 0$ has the density distribution

$$\rho(R, \phi, z) = \rho_0 \left(1 + \frac{\pi G \rho_0 R^2}{2\sigma^2}\right)^{-2}. \quad (4.306)$$

4.28 [1] In a spherical stellar system with mass profile $M(r)$, a stellar population with number density $n(r)$ has anisotropy parameter (4.61) of the form $\beta(r) = r^2/(r_a^2 + r^2)$, where r_a is a constant. Show that

$$\overline{v_r^2}(r) = \frac{G \int_r^\infty dr' [(r_a/r')^2 + 1] n(r') M(r')}{(r_a^2 + r^2) n(r)}. \quad (4.307)$$

4.29 [2] Let us write the general n th-order velocity moment in spherical coordinates in the form $\overline{v_\theta^{n-j} v_\phi^{j-k} v_r^k}$. Prove that when $f = f(H)$ (i) the moment vanishes if any of j , k , or n is odd; (ii) if j , k , and n are all even, then

$$\overline{v_\theta^{n-j} v_\phi^{j-k} v_r^k} = \frac{\overline{v_r^n} \left(\frac{k-1}{2}\right)! \left(\frac{n-j-1}{2}\right)! \left(\frac{j-k-1}{2}\right)!}{\pi \left(\frac{n-1}{2}\right)!}. \quad (4.308)$$

4.30 [2] When the DF of an axisymmetric system has the form $f(H, L_z)$, show that the n th-order velocity moments are related by

$$\overline{v_R^{n-j+2} v_z^{j-k-2} v_\phi^k} = \frac{n-j+1}{j-k-1} \overline{v_R^{n-j} v_z^{j-k} v_\phi^k}, \quad (4.309)$$

where j , k , and n are all even. Moments that contain odd powers of either v_R or v_z vanish. Hence the independent non-vanishing n th order moments may be taken to be $\overline{v_\phi^n}$, $\overline{v_\phi^{n-2} v_R^2}$, $\overline{v_\phi^{n-4} v_R^4}$, \dots

4.31 [1] Show that in a stellar-dynamical polytrope $\overline{v_r^2} \propto \Psi$. Show that for a Plummer model the coefficient of proportionality is $\frac{1}{6}$.

4.32 [2] The velocity dispersion in some axisymmetric stellar system is isotropic and a function $\sigma(\rho)$ of the density alone. Show that the mean azimuthal velocity must be a function $\overline{v_\phi}(R)$ of the cylindrical radius R only. Is this configuration physically plausible?

4.33 [2] A static, spherically symmetric stellar system with ergodic DF is confined by a spherical vessel of radius r_b . Show that $2K + W = 4\pi r_b^3 p$, where K and W are the system's kinetic and potential energies, and $p = \rho \overline{v_r^2}$ is the pressure exerted by the system on the vessel's walls.

4.34 [1] Suppose the principal axes of the velocity ellipsoid near the Sun are always parallel to the unit vectors of spherical coordinates. Then show that for $|z|/R$ small, $\overline{v_R v_z} \simeq (\overline{v_r^2} - \overline{v_\theta^2})(z/R)$.

4.35 [1] A stationary stellar system of negligible mass and finite extent is confined by the potential $\Phi(r) = v_c^2 \ln r + \text{constant}$.

(a) Prove that the mean-square velocity is $\langle v^2 \rangle = v_c^2$, independent of the shape, radial profile, or other properties of the stellar system. Hint: as in the derivation of the virial theorem, consider the behavior of $d^2 I / dt^2$, where in this case $I = r^2$.

(b) The singular isothermal sphere has the same potential (eq. 4.104), but in this system the mean-square velocity is $\langle v^2 \rangle = 3\sigma^2 = \frac{3}{2}v_c^2$. How is this consistent with the result of part (a)?

4.36 [2] The energy per unit mass of a star in a stationary stellar system can be written $\epsilon = \frac{1}{2}v^2 + \Phi$, where v is the speed of the star and $\Phi(\mathbf{x})$ is the gravitational potential. Prove that the total energy of the stellar system is

$$\tilde{E} = \frac{1}{3} M \langle \epsilon \rangle, \quad (4.310)$$

where M is the total mass of the system and $\langle \cdot \rangle$ denotes a mass-weighted average over the stars.

4.37 [1] Show that the part of equation (4.239) that is antisymmetric in j and k is equivalent to the law of conservation of angular momentum.

4.38 [1] Show that in the presence of an externally generated gravitational potential Φ_{ext} , the right side of equation (4.247) acquires an extra term:

$$V_{jk} \equiv -\frac{1}{2} \int d^3\mathbf{x} \left(x_k \frac{\partial \Phi_{\text{ext}}}{\partial x_j} + x_j \frac{\partial \Phi_{\text{ext}}}{\partial x_k} \right) \rho. \quad (4.311)$$

4.39 [2] In this problem we use the tensor virial theorem to connect the shape of a bar, its pattern speed, and the extent to which there is less motion parallel to the axis of figure rotation than in the perpendicular directions. Let the z axis coincide with a principal axis of the tensor \mathbf{I} (eq. 4.243), and suppose that the density distribution is stationary in a frame that rotates about this axis with angular frequency Ω . Show that at an instant when $I_{xy} = 0$, the left side of equation (4.247) is Ω^2 times the diagonal tensor with components $(I_{yy} - I_{xx})$, $(I_{xx} - I_{yy})$, and 0 along the diagonal. Hence show that

$$\Omega^2 = -\frac{(W_{xx} - W_{yy}) + 2(T_{xx} - T_{yy}) + (\Pi_{xx} - \Pi_{yy})}{2(I_{xx} - I_{yy})}, \quad (4.312a)$$

and if $T_{zz} = 0$,

$$\frac{v_0^2}{\sigma_0^2} = (1 - \delta) \frac{W_{xx} + W_{yy}}{W_{zz}} - 2, \quad (4.312b)$$

where $v_0^2 \equiv 2(T_{xx} + T_{yy})/M$, $\sigma_0^2 \equiv (\Pi_{xx} + \Pi_{yy})/2M$ and $(1 - \delta)(\Pi_{xx} + \Pi_{yy}) \equiv 2\Pi_{zz}$.

4.40 [1] Suppose that the Oort limit has been determined as described in §4.9.3 from observations of stars whose distances have been systematically overestimated by a factor λ . By what factor is the derived local mass density $\rho(0)$ in error, if the kinematics are derived from (a) radial velocities; (b) proper motions?

4.41 [2] Consider a hypothetical disk galaxy in which all the mass is contained in a central point mass. The disk density is negligible; more precisely, the disk consists of a population of stars of zero mass with RMS z -velocity σ_z that is independent of z . At radius R , the number density of these stars as a function of z is $\nu(z) = \nu(0) \exp(-z^2/2z_0^2)$, where $z_0 \ll R$ is a constant. (a) What is the relation between σ_z and z_0 ? (b) What does equation (4.273) predict for the local mass density if these stars are used as tracers? Why is the wrong answer obtained?

4.42 [2] Consider a time-independent, self-gravitating collisionless stellar system with slab symmetry, that is, a system in which the density ρ depends only on a single coordinate z . Prove that the system must be symmetric, that is, with a suitable choice of the coordinate origin $\rho(z) = \rho(-z)$.

4.43 [3] A natural model DF for a razor-thin axisymmetric disk is given by equation (4.147),

$$f(H, L_z) = S(L_z) \exp[-\Delta/\sigma_R^2(L_z)], \quad (4.313)$$

where $\Delta = H - E_c(L_z)$ and $E_c(L_z)$ is the energy of a circular orbit with angular momentum L_z . For a disk with surface density $\Sigma(R) \propto \exp(-R/R_d)$, in a potential with a constant circular speed v_0 , we may take $\sigma_R^2(L_z) \propto \exp[-L_z/(R_d v_0)]$ and $S(L_z) \propto \Sigma(L_z/v_0)/\sigma_R^2(L_z) \propto \text{constant}$ (see eq. 4.156 and the following discussion). With these assumptions, the DF in the solar neighborhood, at radius $R = R_0$, depends on the dimensionless parameters $\xi \equiv R_0/R_d$ and $b \equiv \sigma_R(R_0 v_0)/v_0$, which is $\ll 1$ in a cool disk. Show that in the solar neighborhood

$$\begin{aligned} \overline{v_\phi} &= v_0 + \left(\frac{1}{4} - \xi\right) \frac{w}{v_0} + \left(\frac{1}{32} + \frac{5}{12}\xi + \frac{3}{2}\xi^2 - \frac{9}{8}\xi^3\right) \frac{w^2}{v_0^3} + O(w^3), \\ \overline{(v_\phi - \overline{v_\phi})^2} &= \frac{1}{2}w^2 - \left(\frac{1}{8} + \xi - \frac{5}{4}\xi^2\right) \frac{w^2}{v_0^2} + O(w^3), \end{aligned} \quad (4.314)$$

where $w \equiv \overline{v_R^2}$. Relate the $O(w)$ terms to epicycle theory (§3.2.3) and Stromberg's asymmetric drift equation (§4.8.2a). The $O(w^2)$ terms provide convenient analytic estimates for the errors incurred in using the epicycle and Stromberg approximations. Hint: use computer algebra.

4.44 [2] (a) By taking a suitable moment of the collisionless Boltzmann equation, show that in a steady-state axisymmetric galaxy

$$\frac{\partial(\overline{\nu v_R^2 v_\phi})}{\partial R} + \frac{\partial(\overline{\nu v_R v_z v_\phi})}{\partial z} - \frac{\nu}{R} \left(\overline{v_\phi^3} - \overline{v_\phi} R \frac{\partial \Phi}{\partial R} \right) + \frac{2\nu \overline{v_R^2 v_\phi}}{R} = 0. \quad (4.315)$$

(b) Given that the system is symmetric in z , and that all odd moments of $v_\phi - \overline{v_\phi}$ vanish, so $0 = \overline{v_R^2 (v_\phi - \overline{v_\phi})}$ and $0 = \overline{(v_\phi - \overline{v_\phi})^3}$, etc., show that at $z = 0$

$$\overline{v_R^2} \left(\frac{\partial \overline{v_\phi}}{\partial R} + \frac{\overline{v_\phi}}{R} \right) - \frac{2}{R} \overline{v_\phi (v_\phi - \overline{v_\phi})^2} = 0. \quad (4.316)$$

Hence using equation (4.222a), show that (cf. eq. 3.100)

$$\frac{\sigma_\phi^2}{\sigma_R^2} \equiv \frac{\overline{(v_\phi - \overline{v_\phi})^2}}{\overline{v_R^2}} \simeq \frac{-B}{A - B}, \quad (4.317)$$

where A and B are the Oort constants (eq. 3.83). What is the most questionable assumption made in this derivation? Explain why violations of equation (4.317) increase with σ_R , and compare this result to the results of Problem 4.43.

4.45 [3] A rotating axisymmetric stellar system has a star density in velocity space that is constant on ellipsoids, that is, the DF at a given position depends on velocity \mathbf{v} only through the combination $Q = \sum_{i,j} s_{ij} (v_i - \overline{v}_i)(v_j - \overline{v}_j)$, where $\overline{\mathbf{v}} = \overline{v}_\phi \hat{\mathbf{e}}_\phi$ is the mean azimuthal velocity.²²

(a) If the DF depends only on \mathcal{E} and L_z , prove that the rotation curve must have the form

$$\overline{v}_\phi(R, z) = \frac{R}{a + bR^2}, \quad (4.318)$$

where a and b are constants.

(b) If the velocity distribution is isotropic (constant on spheres in velocity space), so $s_{ij} = s\delta_{ij}$, prove that the system rotates at constant angular velocity, that is, the constant b in equation (4.318) is zero.

(c) Prove that result (b) holds for any stationary DF, even if it depends on a third integral.

²² Systems of this kind were a major early focus of research in stellar dynamics, because Schwarzschild's observation that the velocity distribution was ellipsoidal in the solar neighborhood (§4.4.3) led theorists to explore the **ellipsoidal hypothesis** that the distribution was ellipsoidal at all points in the Galaxy. See Chandrasekhar (1942).

5

Stability of Collisionless Systems

5.1 Introduction

In this chapter we examine how an equilibrium stellar system responds to external forces. Does it ring like a bell? Does the disturbance grow larger and larger, like the tilt of a pencil that is initially balanced on its point? Or does the disturbance die away? In fact, we shall find that all three types of response can be present in stellar systems.

The stability of a dynamical system is a particular aspect of its response to external forces; loosely speaking, the system is unstable if a small perturbation causes a large response. Obviously, we would not expect to find in nature a system so delicately balanced that the slightest change causes it to evolve rapidly away from its initial state. Hence we should test any proposed stellar-dynamical configuration for stability before using it to model a real stellar system. It turns out that many equilibrium stellar systems *are* unstable. For example, the simplest model of a disk galaxy consists of a self-gravitating, razor-thin disk in which the stars move on precisely circular orbits, but we shall find that any cold disk, with no random motions, is violently unstable. It appears that a minimum level of random motion is needed to stabilize a self-gravitating disk (eq. 6.71).

Stability analyses can also aid us to interpret observations of astronomical systems. A classic example is furnished by Saturn's rings: Laplace (1829) showed that these could not be rigid bodies, as most astronomers then believed, because a solid ring would be unstable (Problem 5.1). A more recent example concerns the structure of disk galaxies. Dynamical models of low-luminosity disk galaxies often exhibit a fierce non-axisymmetric instability that results in the growth of a large bar-like structure in the central regions. Since these galaxies do not seem to be unstable, we infer that dark matter probably dominates the gravitational forces throughout the disk (see §6.3.4).

In Chapter 3 discussed the stability of orbits in a variety of fixed gravitational potentials. The instabilities that we shall encounter here are of a different type. They are caused by cooperative or collective effects, in which a density perturbation gives rise to extra gravitational forces, which deflect the stellar orbits in such a way that the original density perturbation is enhanced.

The response of a galactic disk is strongly influenced by differential rotation—the variation of angular speed with distance from the galactic center. Differential rotation shears out disk disturbances since the material at each radius is carried around the galaxy at a different rate. We shall defer this complication to Chapter 6, which is devoted to the dynamics of differentially rotating disks and the origin of spiral structure in galaxies. In this chapter we focus on static and uniformly rotating stellar systems.

Our task is made easier by the similarities between stellar systems and two other kinds of systems:

- (i) *Self-gravitating fluids.* As we saw in Chapter 4, many equilibrium stellar systems have close fluid analogs. The analogies arise because a fluid system is supported against gravity by gradients in the scalar pressure p , while a stellar system is supported by gradients in the stress tensor $-\nu\sigma_{ij}^2$ (see the discussion after eq. 4.209). Similarly, in this chapter we shall find it useful to draw analogies between the responses of self-gravitating fluid and stellar systems. A fluid system resists the gravitational collapse of a local density enhancement through pressure gradients, while a stellar system resists collapse because the spread in stellar velocities at every point tends to disperse any density enhancement before it has time to grow. Since the fluid systems are simpler, and important in their own right, we shall analyze fluid systems alongside the analogous stellar-dynamical ones.
- (ii) *Electrostatic plasmas.* Rarefied plasmas share with collisionless stellar systems the property that the mean field of the system is more important than forces from nearby particles (§1.2). Hence many of the techniques of plasma physics can be used in stellar dynamics. However, there is a fundamental difference: plasmas have both positive and negative charges, so they are neutral on large scales and can form static homogeneous equilibria. By contrast, gravity is always an attractive force,

so equilibrium gravitating systems *must* be inhomogeneous. This essential inhomogeneity of self-gravitating systems complicates the study of their stability.

In the remainder of this section we develop the machinery of linear response theory and linear perturbation theory that we shall use to analyze the stability and response of fluid and stellar systems. In §5.2 we examine the response of infinite, homogeneous, self-gravitating systems—this case is artificial but illuminates much of the behavior of more realistic fluid and stellar systems. In §5.3 we outline the methods used to evaluate the response of general stellar systems; in §5.4 we describe the energy principles that constrain the behavior of stellar systems, and in §5.5 we apply these tools to spherical systems. Finally, in §5.6 we examine the novel behavior seen in uniformly rotating stellar systems, as an introduction to the analysis of differentially rotating systems in Chapter 6.

5.1.1 Linear response theory

Some insight into the linear response of stellar or fluid systems can be achieved before actually solving the equations of motion. Let us examine an equilibrium system with density $\rho_s(\mathbf{x})$ that is forced by an external gravitational field $-\epsilon \nabla \Phi_e$, where $|\nabla \Phi_e|$ is of the same order as the gravitational field in the equilibrium system, and $\epsilon \ll 1$. The density distribution that would generate this field is $\epsilon \rho_e(\mathbf{x}, t)$ where

$$\nabla^2 \Phi_e = 4\pi G \rho_e. \quad (5.1)$$

Because the perturbation is weak, the response is linear and therefore also proportional to ϵ ; thus we may write the induced density perturbation in the system as $\epsilon \rho_{s1}(\mathbf{x}, t)$. The **response function** $R(\mathbf{x}, \mathbf{x}', \tau)$ defined by

$$\rho_{s1}(\mathbf{x}, t) = \int d^3 \mathbf{x}' dt' R(\mathbf{x}, \mathbf{x}', t - t') \rho_e(\mathbf{x}', t') \quad (5.2)$$

relates the response density $\rho_{s1}(\mathbf{x}, t)$ to the forcing density $\rho_e(\mathbf{x}', t')$. The response depends only on the difference between t and t' because the equilibrium system is time-independent, so the response at time t from an instantaneous impulse at t' can depend only on the lag $t - t'$. Causality requires that $R(\mathbf{x}, \mathbf{x}', \tau) = 0$ for $\tau < 0$ (the effect cannot precede the cause).

Both the forcing density ρ_e and the response density ρ_{s1} contribute to the gravitational potential, and it is the *total* perturbing potential $\Phi_1 \equiv \Phi_e + \Phi_{s1}$ that determines the dynamics of the system. The corresponding density is $\rho_1 = \rho_e + \rho_{s1}$. The **polarization function** $P(\mathbf{x}, \mathbf{x}', \tau)$ relates the response density ρ_{s1} to the total density ρ_1 :

$$\rho_{s1}(\mathbf{x}, t) = \int d^3 \mathbf{x}' dt' P(\mathbf{x}, \mathbf{x}', t - t') \rho_1(\mathbf{x}', t'). \quad (5.3)$$

Once again, $P(\mathbf{x}, \mathbf{x}', \tau) = 0$ for $\tau < 0$.

It is important to understand the physical difference between the response and polarization functions. The response function describes the density response to an external perturbing force, while the polarization function describes the density response to a total perturbing force, which includes any contributions from the self-gravity of the density response. If the self-gravity of the density response is negligible, then the polarization and response functions are identical.¹

The temporal Fourier transform of a function $y(\tau)$ that vanishes for $\tau < 0$ is $\tilde{y}(\omega)$, where (eq. B.71)

$$\tilde{y}(\omega) = \int_0^{\infty} d\tau y(\tau) e^{i\omega\tau} \quad ; \quad y(\tau) = \int_{ic-\infty}^{ic+\infty} \frac{d\omega}{2\pi} \tilde{y}(\omega) e^{-i\omega\tau}, \quad (5.6)$$

and the real number $c > 0$ is large enough so that $\int d\tau \exp(-c\tau)y(\tau)$ converges.

We shall mostly work with the Fourier transform of R and P , since then the convolution over time in equation (5.2) or (5.3) is simplified to a multiplication:

$$\tilde{\rho}_{s1}(\mathbf{x}, \omega) = \int d^3\mathbf{x}' \tilde{R}(\mathbf{x}, \mathbf{x}', \omega) \tilde{\rho}_e(\mathbf{x}', \omega) = \int d^3\mathbf{x}' \tilde{P}(\mathbf{x}, \mathbf{x}', \omega) \tilde{\rho}_1(\mathbf{x}', \omega). \quad (5.7)$$

The functions $\tilde{R}(\mathbf{x}, \mathbf{x}', \omega)$ and $\tilde{P}(\mathbf{x}, \mathbf{x}', \omega)$ can usually be analytically continued over the entire complex ω -plane, except for isolated poles. Problems 5.2 and 5.3 describe some of the general properties of these functions.

In practice the response function can be determined analytically only for simple systems, such as those in §5.2.4; more powerful numerical techniques are described in §5.3.

¹The term “polarization function” suggests an analogy with the electrostatics of macroscopic media (e.g., Jackson 1999). In a homogeneous medium, the electric field \mathbf{E} , macroscopic charge density ρ , polarization \mathbf{P} , and displacement \mathbf{D} are related by

$$\nabla \cdot \mathbf{D} = \rho \quad ; \quad \mathbf{P} = \epsilon_0 \chi \mathbf{E} \quad ; \quad \mathbf{D} = \epsilon_0 \mathbf{E} + \mathbf{P} = \epsilon \mathbf{E} \quad ; \quad \frac{\epsilon}{\epsilon_0} = 1 + \chi, \quad (5.4)$$

where ϵ_0 is the electric constant (permittivity of the vacuum), and ϵ/ϵ_0 and χ are the dielectric constant and susceptibility of the medium. Here \mathbf{D} arises from unbound charges and hence is analogous to the external or forcing density ρ_e ; $\epsilon_0 \mathbf{E}$ represents the total field and hence is analogous to ρ_1 , and $\epsilon_0 \mathbf{E} - \mathbf{D} = -\mathbf{P}$ represents the field arising from the response of the medium and is analogous to ρ_{s1} . For a uniform, time-independent gravitational field in a homogeneous medium, the polarization and response functions are scalars, so $\rho_{s1} = P\rho_1 = R\rho_e$. Thus the analogy implies that

$$P = \frac{\rho_{s1}}{\rho_1} \Leftrightarrow \frac{-\mathbf{P}}{\epsilon_0 \mathbf{E}} = -\chi \quad ; \quad R = \frac{\rho_{s1}}{\rho_e} \Leftrightarrow \frac{-\mathbf{P}}{\mathbf{D}} = -\frac{\epsilon_0 \chi}{\epsilon} = -\frac{\chi}{1 + \chi}. \quad (5.5)$$

The minus sign in the analogy $P \Leftrightarrow -\chi$ is significant: in a dielectric the movement of bound charges cancels the imposed field \mathbf{D} , whereas in a gravitational system polarization tends to enhance the applied field.

5.1.2 Linearized equations for stellar and fluid systems

The dynamics of self-gravitating collisionless stellar systems is described by the collisionless Boltzmann equation (4.7),

$$\frac{\partial f}{\partial t} + [f, H] = \frac{\partial f}{\partial t} + \mathbf{v} \cdot \frac{\partial f}{\partial \mathbf{x}} - \frac{\partial \Phi}{\partial \mathbf{x}} \cdot \frac{\partial f}{\partial \mathbf{v}} = 0, \quad (5.8)$$

and Poisson's equation (2.10),

$$\nabla^2 \Phi_s(\mathbf{x}, t) = 4\pi G \int d^3\mathbf{v} f(\mathbf{x}, \mathbf{v}, t). \quad (5.9)$$

In these equations $\Phi(\mathbf{x}, t)$ is the total potential, $H = \frac{1}{2}v^2 + \Phi(\mathbf{x}, t)$ is the Hamiltonian, $[\cdot, \cdot]$ is the Poisson bracket, and $\Phi_s(\mathbf{x}, t)$ is the gravitational potential of the stellar system itself, which may differ from the total potential $\Phi(\mathbf{x}, t)$ if there are external forces. Here and throughout this chapter, the DF $f(\mathbf{x}, \mathbf{v}, t)$ is defined to be the mass density of stars in phase space, as opposed to the definition in terms of probability density used in §4.1.

An isolated equilibrium stellar system is described by a time-independent DF $f_0(\mathbf{x}, \mathbf{v})$ and potential $\Phi_0(\mathbf{x})$ that are solutions of (5.8) and (5.9),

$$[f_0, H_0] = 0 \quad ; \quad \nabla^2 \Phi_0 = 4\pi G \int d^3\mathbf{v} f_0, \quad (5.10)$$

where $H_0 = \frac{1}{2}v^2 + \Phi_0(\mathbf{x})$.

Now imagine that this equilibrium system is subjected to a weak gravitational force arising from some external potential $\epsilon \Phi_e(\mathbf{x}, t)$, where $|\nabla \Phi_e|$ is of order $|\nabla \Phi_0|$ and $\epsilon \ll 1$. In response to this disturbance, the DF of the stellar system and the gravitational potential arising from its stars are modified to

$$f(\mathbf{x}, \mathbf{v}, t) = f_0(\mathbf{x}, \mathbf{v}) + \epsilon f_1(\mathbf{x}, \mathbf{v}, t) \quad ; \quad \Phi_s(\mathbf{x}, t) = \Phi_0(\mathbf{x}) + \epsilon \Phi_{s1}(\mathbf{x}, t), \quad (5.11)$$

and the total gravitational potential becomes

$$\Phi(\mathbf{x}, t) = \Phi_0(\mathbf{x}, t) + \epsilon \Phi_1(\mathbf{x}, t), \quad \text{where} \quad \Phi_1(\mathbf{x}, t) = \Phi_{s1}(\mathbf{x}, t) + \Phi_e(\mathbf{x}, t). \quad (5.12)$$

Hence the Hamiltonian in equation (5.8) becomes $H = H_0 + \epsilon \Phi_1$. Substituting these results into equations (5.8) and (5.9), we find that the terms that are independent of ϵ vanish by virtue of (5.10). Dropping the terms proportional to ϵ^2 since $\epsilon \ll 1$, we have

$$\frac{\partial f_1}{\partial t} + [f_1, H_0] + [f_0, \Phi_1] = 0 \quad ; \quad \nabla^2 \Phi_{s1} = 4\pi G \int d^3\mathbf{v} f_1. \quad (5.13)$$

The first of these equations is called the **linearized collisionless Boltzmann equation**; the second is simply Poisson's equation (5.9), except that it relates Φ_{s1} and f_1 , instead of Φ_s and f . Much of this chapter is devoted to analyzing solutions of these equations.

Equation (4.9) states that

$$\frac{df_1}{dt} \equiv \frac{\partial f_1}{\partial t} + [f_1, H_0] \quad (5.14)$$

is the rate of change of f_1 as seen by an observer moving through phase space along the unperturbed orbit. Hence the first two terms in the first of equations (5.13) can be replaced by df/dt . Integrating the equation we obtain

$$f_1(\mathbf{x}, \mathbf{v}, t) = - \int_{-\infty}^t dt' [f_0, \Phi_1]_{\mathbf{x}(t'), \mathbf{v}(t'), t'} , \quad (5.15)$$

where the Poisson bracket is evaluated along the unperturbed orbit $\mathbf{x}(t')$, $\mathbf{v}(t')$ that reaches \mathbf{x}, \mathbf{v} at time t .

We shall examine the response of fluid systems in parallel with stellar systems. Therefore we now consider the fluid-dynamical analogs of equations (5.13). A fluid system with density $\rho_s(\mathbf{x}, t)$, pressure $p(\mathbf{x}, t)$, and velocity $\mathbf{v}(\mathbf{x}, t)$ in a potential $\Phi(\mathbf{x}, t)$ obeys the continuity equation (F.3),

$$\frac{\partial \rho_s}{\partial t} + \nabla \cdot (\rho_s \mathbf{v}) = 0, \quad (5.16)$$

Euler's equation (F.10),

$$\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} = - \frac{1}{\rho_s} \nabla p - \nabla \Phi, \quad (5.17)$$

and Poisson's equation,

$$\nabla^2 \Phi_s = 4\pi G \rho_s, \quad \text{where} \quad \Phi = \Phi_s + \epsilon \Phi_e \quad (5.18)$$

is the sum of the potentials arising from the fluid and the forcing mass distribution.

We must also specify the equation of state relating p and ρ . For our purposes it is sufficient to use a simple equation of state. Thus, we will consider barotropic fluids (eq. F.27), in which the pressure is a function of the density only:

$$p(\mathbf{x}, t) = p[\rho_s(\mathbf{x}, t)]. \quad (5.19)$$

When the equation of state is barotropic, Euler's equation (5.17) can be replaced by

$$\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} = - \nabla (h + \Phi), \quad (5.20)$$

where the specific enthalpy is defined by (eq. F.29)

$$h(\rho_s) \equiv \int_0^{\rho_s} \frac{dp(\rho)}{\rho}. \quad (5.21)$$

In the absence of self-gravity, weak short-wavelength disturbances propagate through a barotropic fluid at the sound speed $v_s(\mathbf{x})$ (eq. F.50), defined by

$$v_s^2(\mathbf{x}) \equiv \left[\frac{dp(\rho)}{d\rho} \right]_{\rho_0(\mathbf{x})}. \quad (5.22)$$

An isolated equilibrium fluid is described by the time-independent density, enthalpy, velocity, and potential distributions $\rho_0(\mathbf{x})$, $h_0(\mathbf{x})$, $\mathbf{v}_0(\mathbf{x})$, and $\Phi_0(\mathbf{x})$, which are solutions of equations (5.16)–(5.20). Then we consider the response of the fluid to a weak external potential $\epsilon\Phi_e(\mathbf{x}, t)$:

$$\begin{aligned} \rho_s(\mathbf{x}, t) &= \rho_0(\mathbf{x}) + \epsilon\rho_{s1}(\mathbf{x}, t) & ; & \quad h(\mathbf{x}, t) = h_0(\mathbf{x}) + \epsilon h_1(\mathbf{x}, t), \\ \mathbf{v}(\mathbf{x}, t) &= \mathbf{v}_0(\mathbf{x}) + \epsilon\mathbf{v}_1(\mathbf{x}, t) & ; & \quad \Phi(\mathbf{x}, t) = \Phi_0(\mathbf{x}) + \epsilon\Phi_1(\mathbf{x}, t). \end{aligned} \quad (5.23)$$

Here, as before, $\Phi_1 = \Phi_{s1} + \Phi_e$ is the total perturbation in the gravitational potential, the sum of the external potential Φ_e and the potential Φ_{s1} arising from the density perturbation ρ_{s1} . Substituting equations (5.23) into the fluid equations, we find that the terms that are independent of ϵ sum to zero, and discarding terms proportional to ϵ^2 , we obtain

$$\frac{\partial\rho_{s1}}{\partial t} + \nabla \cdot (\rho_0\mathbf{v}_1) + \nabla \cdot (\rho_{s1}\mathbf{v}_0) = 0, \quad (5.24a)$$

$$\frac{\partial\mathbf{v}_1}{\partial t} + (\mathbf{v}_0 \cdot \nabla)\mathbf{v}_1 + (\mathbf{v}_1 \cdot \nabla)\mathbf{v}_0 = -\nabla(h_1 + \Phi_{s1} + \Phi_e), \quad (5.24b)$$

$$\nabla^2\Phi_{s1} = 4\pi G\rho_{s1}, \quad (5.24c)$$

$$h_1 = \frac{p_1}{\rho_0} = \left(\frac{dp}{d\rho} \right)_{\rho_0} \frac{\rho_{s1}}{\rho_0} = v_s^2 \frac{\rho_{s1}}{\rho_0}. \quad (5.24d)$$

Equations (5.24a) to (5.24d) constitute a complete set of linear equations that govern the response of a barotropic fluid to small perturbations.

5.2 The response of homogeneous systems

5.2.1 Physical basis of the Jeans instability

Consider a fluid of density ρ_0 and pressure p_0 , with no internal motions so $\mathbf{v}_0 = 0$. Now draw a sphere of radius r around any point and suppose that we compress this spherical region by reducing its radius to $(1-\alpha)r$, where $\alpha \ll 1$. We will deal only with order-of-magnitude arguments in this subsection, so the details of how the fluid is compressed and the exact shape of the perturbed density distribution are not important. To order of magnitude, the density perturbation is $\rho_1 \approx \alpha\rho_0$, and the pressure perturbation is $p_1 \approx (dp/d\rho)_0\alpha\rho_0 = \alpha v_s^2\rho_0$. The pressure force per unit mass is $\mathbf{F}_p = -\nabla p/\rho$, and our compressive perturbation therefore leads to an additional outward pressure force \mathbf{F}_{p1} , where $|\mathbf{F}_{p1}| \approx p_1/(\rho_0 r) \approx \alpha v_s^2/r$, where we have replaced the gradient ∇ by $1/r$. Similarly, the enhanced density of the perturbation gives rise to an extra inward gravitational force $|\mathbf{F}_{g1}| \approx GM\alpha/r^2$, where $M = \frac{4}{3}\pi\rho_0 r^3$ is the mass originally within r ; in other words, $|\mathbf{F}_{g1}| \approx G\rho_0 r\alpha$. If the net force is outward, the compressed fluid re-expands and the perturbation is stable; if the net force is inward, the fluid continues to contract and the perturbation is unstable. Thus, there is instability if $|\mathbf{F}_{g1}| > |\mathbf{F}_{p1}|$, or if

$$G\rho_0 r\alpha \gtrsim \alpha v_s^2/r, \quad \text{that is, if } r^2 \gtrsim \frac{v_s^2}{G\rho_0}. \quad (5.25)$$

The same approximate criterion can be derived by comparing the gravitational potential energy and internal energy in a volume V of size r or mass $M \approx \rho_0 r^3$. The potential energy is $W \approx -GM^2/r \approx -G\rho_0^2 r^5$, and the internal energy is $U \approx Mv_s^2 \approx \rho_0 v_s^2 r^3$. If the perturbation is localized within V , and there is no energy flow into V , then the sum of these two energies plus the bulk kinetic energy associated with the perturbation must be zero. Since the kinetic energy is positive, the perturbation can grow only if $W + U < 0$. This requires $r^2 \gtrsim v_s^2/(G\rho_0)$.

We conclude that *perturbations with a scale longer than $\approx v_s/(G\rho_0)^{1/2}$ are unstable*. This behavior is known as the **Jeans instability** (Jeans 1902). We now give a more careful treatment for homogeneous fluids and stellar systems; our goal is both to analyze the Jeans instability and to sharpen our tools for investigating the response of more realistic systems.

5.2.2 Homogeneous systems and the Jeans swindle

As we mentioned in §5.1, an infinite homogeneous gravitating system cannot be in static equilibrium. Nevertheless, it is useful to set up artificial homogeneous systems because their linear stability properties are relatively easy to analyze. Consider, for example, the response function $R(\mathbf{x}, \mathbf{x}', \tau)$ defined

in equation (5.2). In a homogeneous system this can depend on \mathbf{x} , \mathbf{x}' only through their difference, so the response function takes the form $R(\mathbf{x} - \mathbf{x}', \tau)$, and equation (5.2) simplifies to a convolution in both space and time.

We showed in §5.1.1 that equations involving convolutions in time can be analyzed with temporal Fourier transforms. Similarly, equations involving convolution in space are best analyzed using spatial Fourier transforms. We define the spatial Fourier transform of the function $g(\mathbf{x})$ by (eq. B.68)

$$\bar{g}(\mathbf{k}) = \int d^3\mathbf{x} g(\mathbf{x}) e^{-i\mathbf{k}\cdot\mathbf{x}} \quad ; \quad g(\mathbf{x}) = \int \frac{d^3\mathbf{k}}{(2\pi)^3} \bar{g}(\mathbf{k}) e^{i\mathbf{k}\cdot\mathbf{x}}. \quad (5.26)$$

We shall assume that all relevant functions vanish at spatial infinity so their Fourier transforms exist.

We replace $R(\mathbf{x}, \mathbf{x}', t)$ in equation (5.2) by $R(\mathbf{x} - \mathbf{x}', t)$, with a similar replacement in equation (5.3), multiply both equations by $\exp(-i\mathbf{k}\cdot\mathbf{x})$ and integrate over $d^3\mathbf{x}$, to find

$$\bar{\rho}_{s1}(\mathbf{k}, t) = \int dt' \bar{R}(\mathbf{k}, t - t') \bar{\rho}_e(\mathbf{k}, t') = \int dt' \bar{P}(\mathbf{k}, t - t') \bar{\rho}_1(\mathbf{k}, t'). \quad (5.27)$$

Taking a temporal Fourier transform, we obtain:

$$\tilde{\bar{\rho}}_{s1}(\mathbf{k}, \omega) = \tilde{\bar{R}}(\mathbf{k}, \omega) \tilde{\bar{\rho}}_e(\mathbf{k}, \omega) = \tilde{\bar{P}}(\mathbf{k}, \omega) \tilde{\bar{\rho}}_1(\mathbf{k}, \omega). \quad (5.28)$$

There is a simple relation between the polarization and response functions for homogeneous systems. Using the relation $\rho_1 = \rho_{s1} + \rho_e$ (cf. eq. 5.12), equation (5.28) can be rewritten as

$$\tilde{\bar{\rho}}_{s1}(\mathbf{k}, \omega) = \tilde{\bar{R}}(\mathbf{k}, \omega) \tilde{\bar{\rho}}_e(\mathbf{k}, \omega) = \tilde{\bar{P}}(\mathbf{k}, \omega) [\tilde{\bar{\rho}}_{s1}(\mathbf{k}, \omega) + \tilde{\bar{\rho}}_e(\mathbf{k}, \omega)]. \quad (5.29)$$

Combining these equations gives the desired relations:

$$\tilde{\bar{R}}(\mathbf{k}, \omega) = \frac{\tilde{\bar{P}}(\mathbf{k}, \omega)}{1 - \tilde{\bar{P}}(\mathbf{k}, \omega)} \quad ; \quad \tilde{\bar{P}}(\mathbf{k}, \omega) = \frac{\tilde{\bar{R}}(\mathbf{k}, \omega)}{1 + \tilde{\bar{R}}(\mathbf{k}, \omega)}. \quad (5.30)$$

We construct our artificial equilibrium for a homogeneous system by perpetrating what we shall call the **Jeans swindle**. Mathematically, the difficulty we must overcome is that if the density and pressure of the medium ρ_0 , p_0 are constant, and the mean velocity \mathbf{v}_0 is zero, then Euler's equation (5.17) implies that $\nabla\Phi_0 = 0$. The same conclusion follows from symmetry considerations—since the system is homogeneous, there is no preferred direction for the vector $\nabla\Phi_0$ to point. On the other hand, Poisson's equation (5.18) requires that $\nabla^2\Phi_0 = 4\pi G\rho_0$. These two requirements are inconsistent unless $\rho_0 = 0$. Physically, there are no pressure gradients in a

homogeneous medium to balance gravitational attraction. A similar inconsistency arises in an infinite homogeneous stellar system. We remove the inconsistency by the *ad hoc* assumption that Poisson's equation describes only the relation between the *perturbed* density and the *perturbed* potential, while the unperturbed potential is zero. An equivalent statement is that some fixed gravitational field from an unspecified source cancels $\nabla\Phi_0$. This assumption constitutes the Jeans swindle; it is a swindle, of course, because there is no justification for discarding the unperturbed gravitational field. However, the swindle is vindicated by the insight it provides, so long as its limitations are kept in mind—for further discussion see §5.2.5.²

5.2.3 The response of a homogeneous fluid system

We use the linearized fluid equations (5.24a) to (5.24d). The equilibrium state is $\rho_0 = \text{constant}$, $\mathbf{v}_0 = 0$, and the Jeans swindle lets us set $\Phi_0 = 0$. We then have

$$\frac{\partial\rho_{s1}}{\partial t} + \rho_0\nabla\cdot\mathbf{v}_1 = 0 \quad ; \quad \frac{\partial\mathbf{v}_1}{\partial t} = -\nabla(h_1 + \Phi_{s1} + \Phi_e), \quad (5.31a)$$

$$\nabla^2\Phi_{s1} = 4\pi G\rho_{s1} \quad ; \quad h_1 = v_s^2\rho_{s1}/\rho_0, \quad (5.31b)$$

and the sound speed v_s is a constant.

By taking the time derivative of the first of equations (5.31) and the divergence of the second, and eliminating \mathbf{v}_1 , Φ_{s1} , and h_1 with the aid of the other equations, we can combine equations (5.31) into the single equation

$$\frac{\partial^2\rho_{s1}}{\partial t^2} - v_s^2\nabla^2\rho_{s1} - 4\pi G\rho_0\rho_{s1} = 4\pi G\rho_0\rho_e. \quad (5.32)$$

We now take the spatial Fourier transform by multiplying by $\exp(-i\mathbf{k}\cdot\mathbf{x})$ and integrating over $d^3\mathbf{x}$. We then apply the divergence theorem (B.45) twice to eliminate ∇^2 , using the fact that the boundary terms are zero since all perturbed quantities are assumed to vanish at spatial infinity. We obtain

$$\frac{\partial^2\bar{\rho}_{s1}(\mathbf{k}, t)}{\partial t^2} + v_s^2k^2\bar{\rho}_{s1}(\mathbf{k}, t) - 4\pi G\rho_0\bar{\rho}_{s1}(\mathbf{k}, t) = 4\pi G\rho_0\bar{\rho}_e(\mathbf{k}, t). \quad (5.33)$$

A **mode** is a perturbation that can be sustained without external forces. Setting $\bar{\rho}_e(\mathbf{k}, t) = 0$ and substituting $\bar{\rho}_{s1}(\mathbf{k}, t) \propto \exp(-i\omega t)$ into equation (5.33), we find that modes must satisfy the dispersion relation

$$\omega^2 = \omega_0^2(k) \equiv v_s^2k^2 - 4\pi G\rho_0 = v_s^2(k^2 - k_J^2), \quad (5.34)$$

²The Jeans swindle is not needed in one (unrealistic) cosmological model. This is the static model proposed, and later repudiated, by Einstein, in which the gravitational attraction of the unperturbed density ρ_0 is canceled by vacuum energy ρ_Λ . See Problem 1.8.

where the **Jeans wavenumber** k_J is defined by

$$k_J^2 \equiv \frac{4\pi G\rho_0}{v_s^2}. \quad (5.35)$$

For $k > k_J$ (short wavelengths) $\omega^2 > 0$ and the solutions are oscillatory; for $k < k_J$ (long wavelengths) the solutions are exponentially growing or decaying.

Now suppose that the external density is given by $\bar{\rho}_e(\mathbf{k}, t) = \delta(t - t_0)$, where as usual $\delta(\tau)$ is the delta function (Appendix C.1). Then equation (5.27) implies that $\bar{\rho}_{s1}(\mathbf{k}, t)$ is equal to the response function $\bar{R}(\mathbf{k}, t - t_0)$; substituting this result into equation (5.33) yields

$$\frac{\partial^2 \bar{R}(\mathbf{k}, \tau)}{\partial \tau^2} + v_s^2 k^2 \bar{R}(\mathbf{k}, \tau) - 4\pi G\rho_0 \bar{R}(\mathbf{k}, \tau) = 4\pi G\rho_0 \delta(\tau), \quad (5.36)$$

subject to the causality condition $\bar{R}(\mathbf{k}, \tau) = 0$ for $\tau < 0$. We now solve this to determine the response function.

First consider the case of an oscillatory perturbation ($k > k_J$). For $\tau > 0$ the external potential is zero, so $\bar{R}(\mathbf{k}, \tau)$ must satisfy the dispersion relation (5.34),

$$\bar{R}(\mathbf{k}, \tau) = A \sin[\omega_0(k)\tau] + B \cos[\omega_0(k)\tau], \quad (5.37)$$

where A and B are constants chosen to satisfy the boundary conditions at $\tau = 0$. The response function must be a continuous function of τ , since the response density cannot change discontinuously even if subjected to an impulse from an external source. Since the response function is zero for $\tau < 0$, we must have $B = 0$. Now integrate equation (5.36) from $\tau = -\epsilon$ to $\tau = \epsilon$, where ϵ is a small positive number. Since the response function is continuous, the second and third terms on the left side vanish as $\epsilon \rightarrow 0$. Since $\int_{-\epsilon}^{\epsilon} d\tau \delta(\tau) = 1$ we have

$$\left. \frac{\partial \bar{R}(\mathbf{k}, \tau)}{\partial \tau} \right|_{-\epsilon}^{\epsilon} = 4\pi G\rho_0, \quad (5.38)$$

which requires $A\omega_0 = 4\pi G\rho_0$. Thus the response function for $k > k_J$ is

$$\bar{R}(\mathbf{k}, \tau) = \frac{4\pi G\rho_0}{\omega_0(k)} H(\tau) \sin[\omega_0(k)\tau] \quad (5.39)$$

where $H(\tau)$ is the step function (Appendix C.1). For $k < k_J$, $\omega_0(k)$ is purely imaginary, so it proves more convenient to use

$$\gamma_0^2(k) \equiv -\omega_0^2(k) = v_s^2(k_J^2 - k^2), \quad (5.40)$$

and a similar derivation shows that the response function is

$$\overline{R}(\mathbf{k}, \tau) = \frac{4\pi G \rho_0}{\gamma_0(k)} H(\tau) \sinh[\gamma_0(k)\tau]. \quad (5.41)$$

Taking the temporal Fourier transform (5.6) of our results we obtain

$$\widetilde{R}(\mathbf{k}, \omega) = \begin{cases} -\frac{4\pi G \rho_0}{\omega^2 - \omega_0^2(k)}, & k > k_J, \text{ Im}(\omega) > 0, \\ -\frac{4\pi G \rho_0}{\omega^2 + \gamma_0^2(k)}, & k < k_J, \text{ Im}(\omega) > \gamma_0(k). \end{cases} \quad (5.42)$$

The polarization function for a homogeneous system can be obtained by a closely analogous derivation (Problem 5.6). However, there is a simpler route: the general relation (5.30) between the polarization and response functions in homogeneous systems gives

$$\widetilde{P}(\mathbf{k}, \omega) = -\frac{4\pi G \rho_0}{\omega^2 - v_s^2 k^2}, \quad \text{Im}(\omega) > 0, \quad (5.43)$$

and the inverse Fourier transform gives

$$\overline{P}(\mathbf{k}, \tau) = \frac{4\pi G \rho_0}{v_s k} H(\tau) \sin(v_s k \tau). \quad (5.44)$$

This result can also be derived by considering equation (5.39) in the limit that self-gravity is negligible, which occurs when $k_J/k \rightarrow 0$. As we argued in the paragraph following equation (5.3), if self-gravity is negligible the polarization and response functions are identical. As $k_J \rightarrow 0$, $\omega_0(k) \rightarrow v_s k$, so (5.44) follows from (5.39).

The dispersion relation for modes, equation (5.34), can be obtained by setting the external density ρ_e in equation (5.29) to zero; thus

$$\widetilde{P}(\mathbf{k}, \omega) = 1. \quad (5.45)$$

The solutions of the dispersion relation are the singularities of the response function (5.30), since a mode has non-zero density response even though there is no external forcing.

We now describe these results in more physical terms. Equation (5.39) shows that an external impulse with wavenumber $k > k_J$ sets up a sinusoidal oscillation in the fluid, with frequency $\omega_0(k)$. When the density ρ_0 is small, the oscillation consists of sound waves, since the dispersion relation (5.34) reduces to that of a sound wave, $\omega^2(k) = v_s^2 k^2$ (eq. F.54). As the density is increased, the frequency ω decreases and the oscillation becomes more and more sluggish, until eventually ω reaches zero at $k = k_J$. When $k < k_J$,

the external impulse excites a growing density disturbance with temporal dependence $\propto \exp \gamma_0(k)t$ at large time (eq. 5.40). The presence of the growing solution implies that the system is unstable whenever $k^2 < k_J^2$.

In terms of the wavelength, the perturbation is unstable if λ exceeds the **Jeans length** $\lambda_J = 2\pi/k_J$, that is, if³

$$\lambda^2 > \lambda_J^2 = \frac{\pi v_s^2}{G\rho_0}. \quad (5.46)$$

We define the **Jeans mass** M_J as the mass originally contained within a sphere of diameter λ_J :

$$M_J = \frac{4\pi}{3}\rho_0\left(\frac{1}{2}\lambda_J\right)^3 = 2.92\frac{v_s^3}{G^{3/2}\rho_0^{1/2}}. \quad (5.47)$$

The Jeans instability in a fluid has a simple interpretation in terms of energy. The energy density of an ordinary sound wave is positive. However, the gravitational energy density of a sound wave is negative, because the enhanced potential energy in the compressed regions outweighs the reduced potential energy in the dilated regions. The Jeans instability sets in at the wavelength λ_J at which the net energy density becomes negative, so kinetic energy is available to feed the growing wave (see Problem 5.7).

5.2.4 The response of a homogeneous stellar system

The equilibrium state of an infinite, homogeneous stellar system is described by the DF $f_0(\mathbf{v})$, which is taken here to describe the mass density of stars in (\mathbf{x}, \mathbf{v}) phase space. We shall usually assume that this DF is Maxwellian,

$$f_0(\mathbf{v}) = \frac{\rho_0}{(2\pi\sigma^2)^{3/2}} e^{-v^2/(2\sigma^2)}. \quad (5.48)$$

Since the sound speed v_s in an ideal gas is closely related to its velocity dispersion (see eq. F.55), it is natural to define the Jeans wavenumber k_J of a homogeneous, Maxwellian stellar system by analogy with equation (5.35),

$$k_J^2 \equiv \frac{4\pi G\rho_0}{\sigma^2}. \quad (5.49)$$

We show below that this analogy is exact, in that perturbations to the stellar system with wavelength $\lambda < 2\pi/k_J$ are stable, and perturbations with

³ The dispersion relation (5.34) is similar to the dispersion relation for electrostatic plasma waves, $\omega_0^2(k) = v_s^2 k^2 + \omega_p^2$, where $\omega_p^2 = ne^2/(\epsilon_0 m)$ is the square of the plasma frequency, ϵ_0 is the electric constant, n is the electron number density and e and m are the electron charge and mass. In a plasma like charges repel, so the dispersion relation involves $+\omega_p^2$ instead of $-k_J^2 v_s^2$ and all wavelengths are stable.

$\lambda > 2\pi/k_J$ are unstable. However, the details of the behavior of fluid and stellar systems are otherwise quite different. In particular, the response of most homogeneous stellar systems to perturbations with $\lambda < 2\pi/k_J$ is strongly damped, in contrast to the oscillatory response of fluid systems. This interesting phenomenon of damping in a time-reversible system—apparent dissipation even though the collisionless Boltzmann equation contains no friction, viscosity, or other dissipative effects—is a fundamental but paradoxical property of most stellar systems, which we shall examine with some care.

To analyze the response of the homogeneous stellar system, we use the linearized collisionless Boltzmann and Poisson equations (5.13), and invoke the Jeans swindle to set $\Phi_0 = 0$. We have

$$\frac{\partial f_1}{\partial t} + \mathbf{v} \cdot \frac{\partial f_1}{\partial \mathbf{x}} - \frac{\partial}{\partial \mathbf{x}}(\Phi_{s1} + \Phi_e) \cdot \frac{\partial f_0}{\partial \mathbf{v}} = 0 \quad ; \quad \nabla^2 \Phi_{s1} = 4\pi G \int d^3 \mathbf{v} f_1, \quad (5.50)$$

where as usual Φ_{s1} is the perturbed potential due to the self-gravity of the system, and Φ_e is the perturbing external potential. We multiply these equations by $\exp(-i\mathbf{k} \cdot \mathbf{x})$, integrate over $d^3 \mathbf{x}$, and then use the divergence theorem (B.45) to eliminate spatial derivatives, assuming that f_1 , Φ_e , and Φ_{s1} vanish at infinity. We find

$$\begin{aligned} \frac{\partial \bar{f}_1}{\partial t} + i\mathbf{k} \cdot \mathbf{v} \bar{f}_1 &= i(\bar{\Phi}_{s1} + \bar{\Phi}_e) \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}}; \\ -k^2 \bar{\Phi}_{s1} &= 4\pi G \int d^3 \mathbf{v} \bar{f}_1 = 4\pi G \bar{\rho}_{s1}, \end{aligned} \quad (5.51)$$

where $\bar{f}_1(\mathbf{k}, \mathbf{v}, t)$ is the spatial Fourier transform of $f_1(\mathbf{x}, \mathbf{v}, t)$, defined as in equation (5.26), and ρ_{s1} is the perturbed density of the stellar system. The first of equations (5.51) has the integrating factor $\exp(i\mathbf{k} \cdot \mathbf{v} t)$; assuming that the perturbation vanishes as $t \rightarrow -\infty$, its solution is

$$\bar{f}_1(\mathbf{k}, \mathbf{v}, t) = i\mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}} \int_{-\infty}^t dt' e^{i\mathbf{k} \cdot \mathbf{v}(t'-t)} [\bar{\Phi}_{s1}(\mathbf{k}, t') + \bar{\Phi}_e(\mathbf{k}, t')]. \quad (5.52)$$

We now integrate over velocity, replace $\bar{\Phi}_{s1} + \bar{\Phi}_e$ by the total perturbing potential $\bar{\Phi}_1$, and replace this in turn by the perturbing density $\bar{\rho}_1$ using Poisson's equation. We then have

$$\bar{\rho}_{s1}(\mathbf{k}, t) = \int d^3 \mathbf{v} \bar{f}_1 = -\frac{4\pi G i}{k^2} \int d^3 \mathbf{v} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}} \int_{-\infty}^t dt' e^{i\mathbf{k} \cdot \mathbf{v}(t'-t)} \bar{\rho}_1(\mathbf{k}, t'). \quad (5.53)$$

Comparing with equation (5.27), we find that the polarization function for a homogeneous stellar system is

$$\bar{P}(\mathbf{k}, \tau) = -\frac{4\pi G i}{k^2} H(\tau) \int d^3 \mathbf{v} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}} e^{-i\mathbf{k} \cdot \mathbf{v} \tau}. \quad (5.54)$$

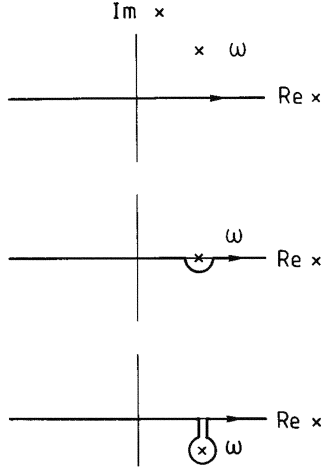


Figure 5.1 Integration contours in the complex x -plane.

In particular, for the Maxwellian DF (5.48),

$$\bar{P}(\mathbf{k}, \tau) = 4\pi G \rho_0 \tau H(\tau) e^{-(k\sigma\tau)^2/2}. \quad (5.55)$$

The polarization function decays to zero at $\tau \rightarrow \infty$, in contrast to the polarization of a fluid, which oscillates (eq. 5.44). This decay is a manifestation of phase mixing (§4.10.2).

We now turn to the response function. We take the temporal Fourier transform of the polarization function (5.54),

$$\tilde{P}(\mathbf{k}, \omega) = \int d\tau \bar{P}(\mathbf{k}, \tau) e^{i\omega\tau} = -\frac{4\pi G}{k^2} \int \frac{d^3\mathbf{v}}{\mathbf{k} \cdot \mathbf{v} - \omega} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}}, \quad \text{Im}(\omega) > 0. \quad (5.56)$$

The response function follows from equations (5.6) and (5.30),

$$\begin{aligned} \bar{R}(\mathbf{k}, \tau) &= \frac{1}{2\pi} \int_{ic-\infty}^{ic+\infty} d\omega \tilde{R}(\mathbf{k}, \omega) e^{-i\omega\tau} \\ &= \frac{1}{2\pi} \int_{ic-\infty}^{ic+\infty} d\omega \frac{\tilde{P}(\mathbf{k}, \omega) e^{-i\omega\tau}}{1 - \tilde{P}(\mathbf{k}, \omega)}, \end{aligned} \quad (5.57)$$

where c is a sufficiently large real number. To evaluate this integral, we wish to extend the integration path to a closed contour. For $\tau > 0$, this can be done using a large semi-circle in the lower half-plane. Consequently, the value of the integral is given by the sum over the residues of any poles that the integrand has below the original integration path. Specifically,

$$\bar{R}(\mathbf{k}, \tau) = -i \sum_{\mathbf{p}} R_{\mathbf{p}} e^{-i\omega_{\mathbf{p}}\tau}, \quad (5.58)$$

where R_p is the residue of $\tilde{R}(\mathbf{k}, \omega)$ at the pole $\omega = \omega_p$. The location of these poles determines the stability of the stellar system: if the maximum value of $\text{Im}(\omega_p)$ is positive, then $\exp(-i\omega_p\tau)$ grows exponentially, and the system is unstable; if zero, the system is stable but oscillates forever after a disturbance; and if negative all disturbances eventually decay exponentially, with the largest value of $\text{Im}(\omega_p)$ dominating the long-term response.

To find the poles, we first consider the upper half-plane $\text{Im}(\omega) > 0$. The numerator $\tilde{P}(\mathbf{k}, \omega)$ has no poles in the upper half-plane, since the integral (5.56) is non-singular in this region. Hence, any poles in the upper half-plane must arise from zeros of the denominator, i.e., from solutions of the dispersion relation

$$\tilde{P}(\mathbf{k}, \omega) = 1, \quad \text{Im}(\omega) > 0; \quad (5.59)$$

solutions of this equation yield unstable modes.

Equation (5.56) does not define $\tilde{P}(\mathbf{k}, \omega)$ for $\text{Im}(\omega) \leq 0$: in this region, the integrand of (5.57) is defined by the analytic continuation of $\tilde{P}(\mathbf{k}, \omega)$ to the lower half-plane. We next examine how to do this.

We set up Cartesian coordinates (v_1, v_2, v_3) in velocity space, with the 1-axis parallel to \mathbf{k} , and let $F_0(v_1) = \int dv_2 dv_3 f_0(\mathbf{v})$ and $x = kv_1$. Then equation (5.56) simplifies to

$$\tilde{P}(\mathbf{k}, \omega) = -\frac{4\pi G}{k^2} \int_{-\infty}^{\infty} dx \frac{F_0'(x/k)}{x - \omega}, \quad \text{Im}(\omega) > 0, \quad (5.60)$$

where $F_0'(v_1) = dF_0/dv_1$. Next, we consider x to be a complex variable, and recall that because we are integrating along the real x axis and $\text{Im}(\omega) > 0$, the integral in (5.60) is non-singular. This situation corresponds to the top panel in Figure 5.1. Now let us analytically continue the function in equation (5.60) to $\text{Im}(\omega) = 0$. Before doing so, we make the small semicircular deformation of the integration contour in the complex x -plane that is shown in the middle panel of Figure 5.1. This deformation does not change the value of the integral, but allows us to reduce $\text{Im}(\omega)$ to zero without crossing the integration contour. Similarly, to continue analytically equation (5.60) to $\text{Im}(\omega) < 0$ we make the keyhole-shaped deformation shown in the bottom panel of Figure 5.1, which once again does not change the value of the integral. Therefore the analytic continuation of the polarization function can be written as

$$\tilde{P}(\mathbf{k}, \omega) = -\frac{4\pi G}{k^2} \int_{\mathcal{L}} dx \frac{F_0'(x/k)}{x - \omega}, \quad (5.61)$$

where \mathcal{L} is the **Landau contour** shown in Figure 5.1 (Landau 1946). The Landau contour allows us to continue analytically $\tilde{P}(\mathbf{k}, \omega)$ over the entire complex ω -plane.

By the residue theorem

$$\int_{\mathcal{L}} dx \frac{F_0'(x/k)}{x - \omega} = \begin{cases} \int_{-\infty}^{\infty} dx \frac{F_0'(x/k)}{x - \omega} & (\text{Im}(\omega) > 0) \\ \wp \int_{-\infty}^{\infty} dx \frac{F_0'(x/k)}{x - \omega} + \pi i F_0'(\omega/k) & (\text{Im}(\omega) = 0) \\ \int_{-\infty}^{\infty} dx \frac{F_0'(x/k)}{x - \omega} + 2\pi i F_0'(\omega/k) & (\text{Im}(\omega) < 0), \end{cases} \quad (5.62)$$

where \wp denotes the Cauchy principal value (eq. C.6).

For the Maxwellian DF (5.48), $F_0(v_1) = \rho_0 \exp(-\frac{1}{2}v_1^2/\sigma^2)/\sqrt{2\pi\sigma^2}$, so the integrals on the right of equation (5.62) can be written

$$-\frac{\rho_0}{\sqrt{2\pi\sigma^3}k} \int_{-\infty}^{\infty} dx \frac{x e^{-x^2/(2k^2\sigma^2)}}{x - \omega}, \quad (5.63)$$

with the understanding that the principal value is to be taken if ω is real. Then equations (5.61) and (5.62) can be combined to give

$$\tilde{P}(\mathbf{k}, \sqrt{2}k\sigma w) = \frac{4\pi G\rho_0}{k^2\sigma^2} [1 + wZ(w)], \quad (5.64)$$

where the “plasma dispersion function” is

$$Z(w) = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} ds \frac{e^{-s^2}}{s - w} \quad (\text{Im}(w) > 0), \quad (5.65)$$

and its analytic continuation for $\text{Im}(w) \leq 0$. The analytic properties of the plasma dispersion function are described in Appendix C.3 and by Fried & Conte (1961).⁴

The dispersion relation follows from equation (5.59),

$$\frac{k^2}{k_J^2} = 1 + wZ(w), \quad \omega = \sqrt{2}k\sigma w, \quad (5.66)$$

where the Jeans wavenumber k_J is defined by equation (5.49).

In examining the dispersion relation (5.66) there are three cases to consider:

(a) Unstable solutions In this case $\text{Im}(\omega) > 0$. In order for the dispersion relation (5.66) to be satisfied, the imaginary part of $wZ(w)$ must vanish. Since the integration variable s in (5.65) is real,

$$wZ(w) = \frac{w}{\sqrt{\pi}} \int_{-\infty}^{\infty} ds \frac{(s - w^*) e^{-s^2}}{|s - w|^2} \quad (5.67)$$

⁴ Analogs for non-Maxwellian DFs—which can be simpler than the Maxwellian case—are discussed in Problem 5.8.

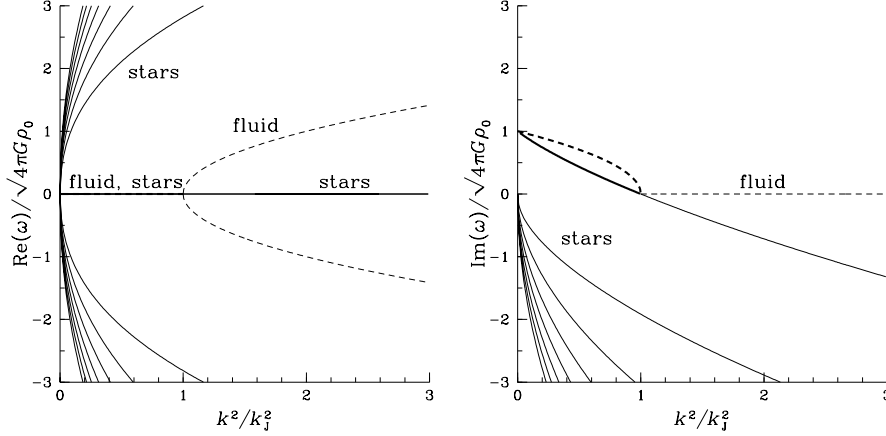


Figure 5.2 The dispersion relation for infinite homogeneous fluid and stellar systems, equations (5.34) and (5.59). The real and imaginary parts of the frequency are plotted separately, in units of $(4\pi G\rho_0)^{1/2}$. The Jeans wavenumber k_J is defined by equation (5.35) for fluids and equation (5.49) for stellar systems. Unstable branches are marked by heavy lines. As discussed on page 413, the curves with $\text{Im}(\omega) < 0$ describe Landau-damped waves rather than true modes.

so

$$\text{Im}[wZ(w)] = \frac{\text{Im}(w)}{\sqrt{\pi}} \int_{-\infty}^{\infty} ds \frac{se^{-s^2}}{|s-w|^2}. \quad (5.68)$$

This integral can vanish only if $\text{Re}(w) = 0$: if, for example, $\text{Re}(w) > 0$, the absolute value of the integrand at $s = s_0 > 0$ will always be larger than its value at $s = -s_0 < 0$, and the integral will be positive, while if $\text{Re}(w) < 0$ the integral will be negative. Hence unstable solutions must have w and ω imaginary—in other words, there are no **overstable** modes with $\text{Im}(\omega) > 0$ and $\text{Re}(\omega) \neq 0$. Setting $\omega = i\gamma$ in equation (5.66) and using equation (C.22), we obtain the dispersion relation for unstable modes,

$$\frac{k^2}{k_J^2} = 1 - \frac{\sqrt{\pi}\gamma}{\sqrt{2}k\sigma} \exp\left(\frac{\gamma^2}{2k^2\sigma^2}\right) \left[1 - \text{erf}\left(\frac{\gamma}{\sqrt{2}k\sigma}\right)\right]. \quad (5.69)$$

The growth rate γ is zero at $k = k_J$, which confirms that the definition (5.49) for the Jeans wavenumber in a Maxwellian stellar system does indeed separate stable and unstable wavenumbers. Unstable waves—solutions of (5.69) with $\gamma > 0$ —exist if and only if $k^2 < k_J^2$. The dispersion relation (5.69) is plotted in Figure 5.2, along with the analogous relation (5.34) for a fluid.

(b) Neutrally stable solutions Here $\text{Im}(\omega) = 0$, $\text{Re}(\omega) \neq 0$. From equation (5.62), for real argument the plasma dispersion function is

$$Z(w) = \frac{1}{\sqrt{\pi}} \wp \int_{-\infty}^{\infty} ds \frac{e^{-s^2}}{s-w} + i\sqrt{\pi}e^{-w^2}. \quad (5.70)$$

The principal value is real, so $\text{Im}[wZ(w)] = 0$ only when $\omega = 0$. Consequently, there are no undamped traveling waves. In this respect the stellar system is quite different from a fluid system, whose dispersion relation (5.34) allows undamped gravity-modified sound waves at all wavelengths less than the Jeans length.

(c) Damped solutions Since only waves having $k < k_J$ are unstable, and there are no undamped waves, all waves with $k > k_J$ must be damped. Numerical evaluation of $Z(w)$ (Figure 5.2) shows that there is an infinite number of solutions of (5.66) for a given wavenumber, all of them strongly damped, in the sense that $|\text{Im}(\omega)/\text{Re}(\omega)|$ is of order unity or greater (the minimum value of this ratio for any $k < k_J$ is 0.481). This damping is called **Landau damping**.⁵

Landau damping arises because of the singularity in the integrand of the polarization function (5.56) when

$$\mathbf{k} \cdot \mathbf{v} - \omega = 0. \quad (5.71)$$

When a star's velocity \mathbf{v} satisfies this equation, its position \mathbf{x} at time t , namely $\mathbf{x}(t) = \mathbf{x}_0 + \mathbf{v}t$, is such that its phase with respect to the wave (ω, \mathbf{k}) is constant: $\phi \equiv \mathbf{k} \cdot \mathbf{x} - \omega t = \mathbf{k} \cdot \mathbf{x}_0 + (\mathbf{k} \cdot \mathbf{v} - \omega)t$. That is, the star follows the wave in the same way that a surfer rides an ocean wave: the surfer may have a substantial velocity parallel to the crest of the wave, but perpendicular to the crest he tries to move at exactly the same speed as the crest; that is, he tries keep his phase with respect to the wave constant.

The surfer sets his phase such that the downhill direction is the direction in which the wave is running, for then the wave is doing work on him: the horizontal component of the buoyancy force \mathbf{F} that the water imparts to the surfboard accelerates him in the direction of the wave motion, so the rate of doing work on the surfer is $\mathbf{F} \cdot \mathbf{v} > 0$ (cf. eq. D.4). The surfer adjusts his speed parallel to the wave crests such that this work balances the dissipation of energy that arises from the board's motion through the water.

If the surfer crosses to the back side of the wave crest, he will do work *on* the wave, and will slow down. But one can imagine a speedboat holding its position on this side of the crest by driving endlessly up the retreating slope of the wave. The boat would then be steadily transferring energy to the wave. Thus a particle that is in resonance with a wave can either draw energy from the wave, or give energy to the wave, depending on its phase with respect to the wave.

In a hot stellar system stars have a continuous distribution of velocities, and few stars will be close to resonance with a given wave. Non-resonant stars will in quick succession do work on the wave when they are on a retreating slope, and gain energy from the wave when they are on an advancing slope,

⁵ Waves in electrostatic plasmas can exhibit much slower Landau damping, in the sense that $|\text{Im}(\omega)| \ll |\text{Re}(\omega)|$.

so on the average they will exchange negligible energy with the wave. By contrast, stars that are nearly resonant can suffer net gains or losses. Stars that are initially moving slightly faster than the wave may lose enough energy as they move up a retreating slope to fall back and become trapped near a trough of the wave. Conversely, stars that are initially moving slower than the wave may gain enough energy on an advancing slope of the wave to become trapped. If the energy given up by the fast stars exceeds that gained by the slow stars, the amplitude of the wave must increase to conserve energy. When the distribution of stellar velocities is Maxwellian, as assumed in our calculation above, there are more slow stars than fast stars, and overall the near-resonant stars gain energy from the wave, so the wave decays. Trapping of stars at wave troughs is a nonlinear process, but nevertheless the rate of energy transfer through this process is captured correctly by the linear mathematical analysis that we have used (Stix 1992).

Both Landau damping and phase mixing (§4.10.2) can damp waves in a stellar system, but the processes differ in several ways: (i) Landau damping is due to collective effects arising because of self-gravity, while phase mixing is a kinematic process that occurs even in systems with no self-gravity. (ii) In phase mixing the amplitude of the fluctuations in the DF does not decay—the decay in the amplitude of the fluctuations in spatial density arises because the fluctuations in the DF become more and more tightly wound in phase space. In contrast, Landau damping washes out the fluctuations in the DF. (iii) Phase mixing almost always leads to decaying density fluctuations, while Landau damping can lead to either growth or decay, depending on the equilibrium DF. (iv) Phase mixing is determined by the behavior of the numerator of the integrand in equation (5.57), while Landau damping is determined by the zeros of the denominator. (v) The rate of decay of the density fluctuations due to phase mixing depends on the initial perturbation and the equilibrium DF, while Landau damping always leads to an exponential decay. For example, if the DF is Maxwellian, the polarization function (5.55) decays *faster* than an exponential, so Landau “damping” actually *extends* the long-term survival of a perturbation.

A subtle but important point is that Landau-damped waves are *not* modes. Our analysis shows only that as $t \rightarrow \infty$, the density response is $\propto \exp(-i\omega_p t)$, where ω_p is determined from the dispersion relation (5.66) using the Landau contour. In contrast, a mode is a solution of the linearized collisionless Boltzmann and Poisson equations (5.50) that behaves like $\exp(-i\omega t)$ at *all* times. Solutions of the dispersion relation (5.59) are modes only in the part of the complex ω -plane where the polarization function exists, which for homogeneous Maxwellian stellar systems means $\text{Im}(\omega) > 0$ (unstable modes). The properties of Landau-damped waves are determined by the analytic continuation of the polarization function to $\text{Im}(\omega) < 0$, where modes do not exist. There *are* modes of a homogeneous stellar system with $k > k_J$, but these **van Kampen modes** have quite different properties from Landau-damped waves: they exist for all real ω and

$k > k_J$ and have singular DFs (Box 5.1).

The distinction between phase mixing and Landau damping can be illustrated by a numerical calculation of the response function for an infinite homogeneous stellar system.⁶ We start with equation (5.53) and replace the density $\bar{\rho}_1$ by $\bar{\rho}_{s1} + \bar{\rho}_e$ to obtain

$$\bar{\rho}_{s1}(\mathbf{k}, t) = -\frac{4\pi G i}{k^2} \int d^3\mathbf{v} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}} \int_{-\infty}^t dt' e^{i\mathbf{k}\cdot\mathbf{v}(t'-t)} [\bar{\rho}_{s1}(\mathbf{k}, t') + \bar{\rho}_e(\mathbf{k}, t')]. \quad (5.72)$$

Now set $\bar{\rho}_e(\mathbf{k}, t) = \delta(t)$ and assume that $\bar{\rho}_{s1}(\mathbf{k}, t) = 0$ for $t < 0$. In this case equation (5.27) shows that $\bar{\rho}_{s1}(\mathbf{k}, t)$ is just the response function $\bar{R}(\mathbf{k}, t)$. Thus

$$\begin{aligned} \bar{R}(\mathbf{k}, t) = & -\frac{4\pi G i}{k^2} \int d^3\mathbf{v} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}} e^{-i\mathbf{k}\cdot\mathbf{v}t} \\ & - \frac{4\pi G i}{k^2} \int d^3\mathbf{v} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}} \int_0^t dt' e^{i\mathbf{k}\cdot\mathbf{v}(t'-t)} \bar{R}(\mathbf{k}, t'), \end{aligned} \quad (5.73)$$

for $t > 0$, and zero for $t < 0$. For a Maxwellian DF,

$$\begin{aligned} \bar{R}(\mathbf{k}, t) = & \frac{4\pi G \rho_0 i}{(2\pi)^{3/2} \sigma^5 k^2} \int d^3\mathbf{v} \mathbf{k} \cdot \mathbf{v} e^{-v^2/(2\sigma^2) - i\mathbf{k}\cdot\mathbf{v}t} \\ & + \frac{4\pi G \rho_0 i}{(2\pi)^{3/2} \sigma^5 k^2} \int d^3\mathbf{v} \mathbf{k} \cdot \mathbf{v} e^{-v^2/(2\sigma^2)} \int_0^t dt' e^{i\mathbf{k}\cdot\mathbf{v}(t'-t)} \bar{R}(\mathbf{k}, t'). \end{aligned} \quad (5.74)$$

Now set up Cartesian coordinates (v_1, v_2, v_3) in velocity space, with the 1-axis parallel to \mathbf{k} . Then the integrations over v_2 and v_3 are immediate, the integral over v_1 is straightforward, and we have

$$\bar{R}(\mathbf{k}, t) = 4\pi G \rho_0 t e^{-(k\sigma t)^2/2} + 4\pi G \rho_0 \int_0^t dt' (t-t') e^{-k^2 \sigma^2 (t-t')^2/2} \bar{R}(\mathbf{k}, t'). \quad (5.75)$$

This is a Volterra integral equation, which is easily solved numerically (Press et al. 1986).⁷ The results are shown in Figure 5.3. The perturbation grows when $k < k_J$ and damps for $k > k_J$, as expected. In the limit $k/k_J \rightarrow \infty$ the damping is entirely due to phase mixing (cf. §4.10.2), the response function is equal to the polarization function of equation (5.55), and the density response decays as a Gaussian in time. As the effects of self-gravity increase, the qualitative nature of the damping changes: the damping becomes less rapid and the density response decays at the exponential rate characteristic of Landau damping—as mentioned earlier, Landau damping actually *slows* the decay relative to phase mixing.

⁶ Suggested by A. Toomre.

⁷ A simple DF for which the response function can be determined analytically is described in Problem 5.8.

Box 5.1: Van Kampen modes

If Landau-damped waves are not modes, then what *are* the modes of a homogeneous stellar system with $k > k_J$? We look for solutions of the linearized collisionless Boltzmann and Poisson equations (5.51) in which ω is real, the external potential $\bar{\Phi}_e = 0$, the perturbed DF $\bar{f}_1(\mathbf{k}, \mathbf{v}, t) = \bar{f}_1(\mathbf{k}, \mathbf{v}) \exp(-i\omega t)$, and the potential due to self-gravity $\bar{\Phi}_{s1}(\mathbf{k}, t) = \bar{\Phi}_{s1}(\mathbf{k}) \exp(-i\omega t)$. This requires

$$(\mathbf{k} \cdot \mathbf{v} - \omega)\bar{f}_1 - \bar{\Phi}_{s1} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}} = 0 \quad ; \quad -k^2 \bar{\Phi}_{s1} = 4\pi G \int d^3\mathbf{v} \bar{f}_1 = 4\pi G \bar{\rho}_{s1}. \quad (1)$$

An obvious solution to the first equation is

$$\bar{f}_a(\mathbf{k}, \mathbf{v}) = \frac{\bar{\Phi}_{s1}}{\mathbf{k} \cdot \mathbf{v} - \omega} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}}; \quad (2)$$

however, van Kampen (1955) pointed out that when $\mathbf{k} \cdot \mathbf{v} - \omega = 0$, the first of equations (1) does not constrain \bar{f}_1 , so the most general solution is $\bar{f}_a + \bar{f}_b$, where

$$\bar{f}_b(\mathbf{k}, \mathbf{v}) = cg(\mathbf{v}_\perp)\delta(\mathbf{k} \cdot \mathbf{v} - \omega). \quad (3)$$

Here c is a constant, \mathbf{v}_\perp is the component of \mathbf{v} that is perpendicular to \mathbf{k} , $g(\mathbf{v}_\perp)$ is an arbitrary function normalized such that $\int d^2\mathbf{v}_\perp g(\mathbf{v}_\perp) = 1$, and $\delta(x)$ is the delta function.

Integrating $\bar{f}_a + \bar{f}_b$ over velocity and using the second of equations (1), it is straightforward to show that the DF of the mode is

$$\begin{aligned} \bar{f}_1(\mathbf{k}, \mathbf{v}) = & \frac{a}{\mathbf{k} \cdot \mathbf{v} - \omega} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}} \\ & - \frac{k^3 a}{4\pi G} \left[1 + \frac{4\pi G}{k^2} \wp \int \frac{d^3\mathbf{v}'}{\mathbf{k} \cdot \mathbf{v}' - \omega} \mathbf{k} \cdot \frac{\partial f_0}{\partial \mathbf{v}'} \right] g(\mathbf{v}_\perp) \delta(\mathbf{k} \cdot \mathbf{v} - \omega), \end{aligned} \quad (4)$$

where a is an arbitrary constant and \wp denotes the principal value (eq. C.6). The corresponding potential is simply $\bar{\Phi}_{s1} = a$.

Singular modes of this form are known as van Kampen modes. It can be shown that the van Kampen modes are complete and that Landau-damped waves can be regarded as a superposition of van Kampen modes (van Kampen 1955; Case 1959; Stix 1992; Vandervoort 2003). The modes do not satisfy a dispersion relation since they exist for all real ω and $k > k_J$. Since these solutions are neither growing nor decaying, they resolve the apparent paradox of how all perturbations in a time-reversible system can decay, as appears to happen with Landau-damped waves. The van Kampen modes also demonstrate that a stellar system has a much richer—and more pathological—set of modes than a fluid.

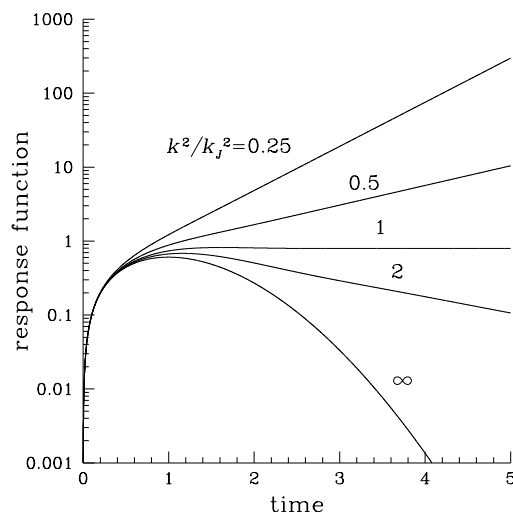


Figure 5.3 Response function of an infinite homogeneous Maxwellian stellar system, in units of $4\pi G\rho_0/(k\sigma)$ (eq. 5.75). Time is measured in units of $(k\sigma)^{-1}$. Curves are labeled by the value of k^2/k_J^2 .

5.2.5 Discussion

The stability of an infinite homogeneous stellar system is closely related to the stability of the analogous fluid system: in both cases there is instability if and only if the wavenumber of the disturbance is less than the Jeans wavenumber k_J , defined by equation (5.35) for fluids and (5.49) for a Maxwellian stellar system. However, the responses to perturbations with $k > k_J$ are quite different: the fluid supports undamped gravity-modified sound waves while the disturbance in the stellar system is strongly phase-mixed and Landau-damped.

These results rest on two linked foundations: the assumption that the system is infinite and homogeneous, and the Jeans swindle of neglecting the gravitational potential of the unperturbed system. It is time to pause and investigate the consequences of these assumptions.

The virial theorem (4.248) tells us that a stellar system of mass M and mean-square velocity v^2 has a characteristic size $\lambda_0 \approx GM/v^2$. In terms of the mean density $\rho \approx M/\lambda_0^3$, we have $\lambda_0^2 \approx v^2/G\rho$. However, from equation (5.46) the Jeans length is given by $\lambda_J^2 \approx v^2/G\rho$, which is the same as λ_0^2 to order of magnitude. Thus the Jeans length is comparable to the size of the system, and the assumption of homogeneity generally is not valid. Accordingly, the Jeans analysis does not establish that there is a real instability in an isolated stellar system. Nevertheless, the analysis is a cornerstone of stability theory for self-gravitating systems, for several reasons:

- The homogeneity assumption is valid and the Jeans swindle is legitimate on scales $\lambda \ll \lambda_0$, because the effects of the inhomogeneity and self-gravity of the equilibrium system are small on small scales. Therefore, we can conclude from our analysis that stationary stellar systems with DFs that are approximately Maxwellian are stable on small scales.

- The Jeans analysis can be used to investigate whether non-Maxwellian DFs introduce new instabilities on small scales. For example, we can ask whether there is an analog for stellar systems of the two-stream instability, which arises in a plasma containing two interpenetrating beams of electrons or ions with different mean velocities (Stix 1992). Because the stellar system is already unstable for $k < k_J$, the two-stream instability is distinct from the Jeans instability only if the DF $f_2(\mathbf{v}) \equiv \frac{1}{2}[f_0(\mathbf{v} + \mathbf{V}) + f_0(\mathbf{v} - \mathbf{V})]$ is unstable at wavenumbers where the DF $f_0(\mathbf{v})$ is stable. It is found that there is no two-stream instability in an infinite homogeneous medium when $f_0(\mathbf{v})$ is Maxwellian (Sweet 1963; Araki 1987), although related instabilities can occur in disks containing equal populations of stars rotating in opposite directions (Sellwood & Merritt 1994; Lovelace, Jore, & Haynes 1997).
- Generalizations of Jeans's analysis to relativistic fluids in a homogeneous expanding universe are central to cosmology, where structures such as galaxies arise from gravitational instabilities (§9.1.2b). In this case the stability analysis can be done self-consistently, without invoking the Jeans swindle.
- The physics of the Jeans instability enables us to interpret small-scale instabilities in rotating disk systems (see §5.6.1 and §6.2.3).
- The Jeans analysis provides insight into aspects of the behavior of more realistic stellar systems, such as the phenomenon of Landau damping, which is discussed in the context of spherical systems in §5.5.3.

5.3 General theory of the response of stellar systems

The response of realistic, finite, inhomogeneous stellar systems is substantially more complicated than the response of infinite homogeneous systems that we described in the last section. Fortunately, many of the tools that we have developed for homogeneous systems remain relevant. Most fundamentally, the linear response of any stellar system is still determined by the linearized collisionless Boltzmann and Poisson equations (5.13). A key to analyzing these systems is the insight that the natural variables for solving these two equations are different: the Poisson equation is simplest in coordinates that reflect the spatial symmetry of the equilibrium system, while the linearized collisionless Boltzmann equation is simplest in the angle-action variables of the equilibrium potential. Most of the work consists of transforming between these two sets of variables, which is numerically tedious but conceptually straightforward. Palmer (1994) provides a thorough review of this approach.

5.3.1 The polarization function in angle-action variables

Let us suppose that the potential $\Phi_0(\mathbf{x})$ of the equilibrium stellar system is integrable, in other words that there exist angle-action variables $(\boldsymbol{\theta}, \mathbf{J})$ in which the equilibrium Hamiltonian is independent of angles, $H = H_0(\mathbf{J})$. All spherical and two-dimensional axisymmetric disk potentials are integrable, as are the triaxial Stäckel potentials discussed in §3.5.3. The position and velocity of any star are functions of the actions and angles, which we may write as $\mathbf{x}(\boldsymbol{\theta}, \mathbf{J})$ and $\mathbf{v}(\boldsymbol{\theta}, \mathbf{J})$.

According to the Jeans theorem, the equilibrium DF may be assumed to depend only on the actions; thus $f_0 = f_0(\mathbf{J})$. To compute the polarization function we imagine that the stellar system is subjected to a weak gravitational force arising from some potential $\epsilon\Phi_1(\mathbf{x}, t)$, where $\epsilon \ll 1$ and Φ_1 includes the self-gravity of the response. The DF is then modified to

$$f(\boldsymbol{\theta}, \mathbf{J}, t) = f_0(\mathbf{J}) + \epsilon f_1(\boldsymbol{\theta}, \mathbf{J}, t). \quad (5.76)$$

Since f_0 and H_0 are independent of the angles $\boldsymbol{\theta}$, the linearized collisionless Boltzmann equation (5.13) becomes

$$\begin{aligned} 0 &= \frac{\partial f_1}{\partial t} + [f_1, H_0] + [f_0, \Phi_1] \\ &= \frac{\partial f_1}{\partial t} + \frac{\partial f_1}{\partial \boldsymbol{\theta}} \cdot \boldsymbol{\Omega} - \frac{\partial f_0(\mathbf{J})}{\partial \mathbf{J}} \cdot \frac{\partial \Phi_1}{\partial \boldsymbol{\theta}}, \end{aligned} \quad (5.77)$$

where

$$\boldsymbol{\Omega}(\mathbf{J}) = \frac{\partial H_0}{\partial \mathbf{J}} \quad (5.78)$$

gives the rate of change $\dot{\boldsymbol{\theta}}$ of the angles in the unperturbed system (eq. 3.190).

Any function of the phase-space coordinates must be a periodic function of the angles $\boldsymbol{\theta}$. Thus we can expand f_1 and Φ_1 in Fourier series (cf. eqs. B.62 and B.64),

$$f_1(\boldsymbol{\theta}, \mathbf{J}, t) = \sum_{\mathbf{m}} f_{\mathbf{m}}(\mathbf{J}, t) e^{i\mathbf{m} \cdot \boldsymbol{\theta}} \quad ; \quad \Phi_1[\mathbf{x}(\boldsymbol{\theta}, \mathbf{J}), t] = \sum_{\mathbf{m}} \Phi_{\mathbf{m}}(\mathbf{J}, t) e^{i\mathbf{m} \cdot \boldsymbol{\theta}}, \quad (5.79)$$

where \mathbf{m} denotes a triple of integers (m_1, m_2, m_3) and

$$\begin{aligned} f_{\mathbf{m}}(\mathbf{J}, t) &= \frac{1}{(2\pi)^3} \int d^3\boldsymbol{\theta}' e^{-i\mathbf{m} \cdot \boldsymbol{\theta}'} f_1(\boldsymbol{\theta}', \mathbf{J}, t), \\ \Phi_{\mathbf{m}}(\mathbf{J}, t) &= \frac{1}{(2\pi)^3} \int d^3\boldsymbol{\theta}' e^{-i\mathbf{m} \cdot \boldsymbol{\theta}'} \Phi_1[\mathbf{x}(\boldsymbol{\theta}', \mathbf{J}), t]. \end{aligned} \quad (5.80)$$

Equation (5.77) must be satisfied separately by each Fourier component. Thus

$$\frac{\partial f_{\mathbf{m}}}{\partial t} + i\mathbf{m} \cdot \boldsymbol{\Omega} f_{\mathbf{m}} = i\mathbf{m} \cdot \frac{\partial f_0(\mathbf{J})}{\partial \mathbf{J}} \Phi_{\mathbf{m}}. \quad (5.81)$$

We multiply this result by $\exp(i\omega t)$ and integrate from $t = 0$ to infinity. The first term can be integrated by parts, yielding

$$[f_{\mathbf{m}}e^{i\omega t}]_0^\infty + i(\mathbf{m} \cdot \boldsymbol{\Omega} - \omega) \int_0^\infty dt f_{\mathbf{m}}e^{i\omega t} = i\mathbf{m} \cdot \frac{\partial f_0(\mathbf{J})}{\partial \mathbf{J}} \int_0^\infty dt \Phi_{\mathbf{m}}e^{i\omega t}. \quad (5.82)$$

We assume that the system was initially in its equilibrium state, so $f_{\mathbf{m}} = \Phi_{\mathbf{m}} = 0$ for $t \leq 0$. We also assume that $\text{Im}(\omega) > c$, where c is a positive constant chosen large enough to ensure that $f_{\mathbf{m}}e^{i\omega t} \rightarrow 0$ as $t \rightarrow \infty$. Then the first term in equation (5.82) vanishes. The integrals in the other terms are simply the temporal Fourier transforms $\tilde{f}_{\mathbf{m}}(\mathbf{J}, \omega)$ and $\tilde{\Phi}_{\mathbf{m}}(\mathbf{J}, \omega)$ (cf. eq. 5.6). Thus

$$\tilde{f}_{\mathbf{m}}(\mathbf{J}, \omega) = \mathbf{m} \cdot \frac{\partial f_0(\mathbf{J})}{\partial \mathbf{J}} \frac{\tilde{\Phi}_{\mathbf{m}}(\mathbf{J}, \omega)}{\mathbf{m} \cdot \boldsymbol{\Omega} - \omega}, \quad \text{Im}(\omega) > c. \quad (5.83)$$

Multiplying equation (5.83) by $\exp(i\mathbf{m} \cdot \boldsymbol{\theta})$ and summing over \mathbf{m} gives the linearized DF $\tilde{f}_1(\boldsymbol{\theta}, \mathbf{J}, \omega)$. However, what we need for the polarization function is the density response $\tilde{\rho}_1(\mathbf{x}, \omega)$. The direct conversion from one to the other is straightforward—just integrate over velocity—but cumbersome, since we must first convert back from angle-action variables to Cartesian coordinates. In the next section, we obtain formulae that are both simpler and more useful, by recasting the problem in a matrix form.

5.3.2 The Kalnajs matrix method

In §2.8 we introduced the concept of a bi-orthonormal potential-density basis: two complete sets of basis functions $\Phi_\alpha(\mathbf{x})$ and $\rho_\alpha(\mathbf{x})$ such that

$$\begin{aligned} \nabla^2 \Phi_\alpha &= 4\pi G \rho_\alpha, \\ - \int d^3\mathbf{x} \Phi_\alpha^*(\mathbf{x}) \rho_\beta(\mathbf{x}) &= \delta_{\alpha\beta}, \end{aligned} \quad (5.84)$$

where α and β are labels for the members of the basis, and $\delta_{\alpha\beta}$ is 1 if $\alpha = \beta$ and zero otherwise. We expand the densities in equation (5.7) using the bi-orthonormal basis:

$$\tilde{\rho}_{s1}(\mathbf{x}, \omega) = \sum_\beta \tilde{d}_\beta(\omega) \rho_\beta(\mathbf{x}) \quad ; \quad \tilde{\rho}_1(\mathbf{x}, \omega) = \sum_\beta \tilde{g}_\beta(\omega) \rho_\beta(\mathbf{x}). \quad (5.85)$$

Because $\rho_\beta(\mathbf{x})$ and $\Phi_\beta(\mathbf{x})$ satisfy Poisson's equation (5.84), the potentials corresponding to $\tilde{\rho}_{s1}$ and $\tilde{\rho}_1$ are described by the same coefficients; thus

$$\tilde{\Phi}_{s1}(\mathbf{x}, \omega) = \sum_\beta \tilde{d}_\beta(\omega) \Phi_\beta(\mathbf{x}) \quad ; \quad \tilde{\Phi}_1(\mathbf{x}, \omega) = \sum_\beta \tilde{g}_\beta(\omega) \Phi_\beta(\mathbf{x}). \quad (5.86)$$

We substitute the expansions (5.85) into equation (5.7), multiply by $\Phi_\alpha^*(\mathbf{x})$, integrate with respect to $d^3\mathbf{x}$, and use the orthogonality relation (5.84) to obtain

$$\tilde{d}_\alpha(\omega) = \sum_\beta \tilde{P}_{\alpha\beta}(\omega) \tilde{g}_\beta(\omega), \quad (5.87a)$$

where

$$\tilde{P}_{\alpha\beta}(\omega) = - \int d^3\mathbf{x} d^3\mathbf{x}' \Phi_\alpha^*(\mathbf{x}) \tilde{P}(\mathbf{x}, \mathbf{x}', \omega) \rho_\beta(\mathbf{x}'). \quad (5.87b)$$

In more compact notation, we may write equation (5.87a) as

$$\tilde{\mathbf{d}}(\omega) = \tilde{\mathbf{P}}(\omega) \tilde{\mathbf{g}}(\omega), \quad (5.88)$$

where $\tilde{\mathbf{d}}$, $\tilde{\mathbf{g}}$ are vectors and $\tilde{\mathbf{P}}$ is a matrix. Thus we have replaced the integral equation (5.7) with a matrix equation, and the polarization function with a polarization matrix. In principle, the matrix has infinite dimension, since there are infinitely many basis functions. However, if the bi-orthonormal basis is chosen sensibly we can obtain an accurate representation of the response with a manageable subset of the basis (Kalnajs 1971, 1977).

Let $\tilde{f}_\beta(\mathbf{x}, \mathbf{v}, \omega) \exp(i\omega t)$ be the perturbation to the DF induced by the potential from a total perturbing density $\rho_\beta(\mathbf{x}) \exp(i\omega t)$. Then

$$\int d^3\mathbf{v} \tilde{f}_\beta(\mathbf{x}, \mathbf{v}, \omega) = \int d^3\mathbf{x}' \tilde{P}(\mathbf{x}, \mathbf{x}', \omega) \rho_\beta(\mathbf{x}'), \quad (5.89)$$

so the polarization matrix (5.87b) can be written

$$\tilde{P}_{\alpha\beta}(\omega) = - \int d^3\mathbf{x} d^3\mathbf{v} \Phi_\alpha^*(\mathbf{x}) \tilde{f}_\beta(\mathbf{x}, \mathbf{v}, \omega). \quad (5.90)$$

Now the transformation from Cartesian phase-space coordinates (\mathbf{x}, \mathbf{v}) to angle-action variables $(\boldsymbol{\theta}, \mathbf{J})$ is canonical, so $d^3\mathbf{x} d^3\mathbf{v} = d^3\boldsymbol{\theta} d^3\mathbf{J}$ (eq. D.81). We can therefore convert the integration variables in equation (5.90) to angle-action variables,

$$\tilde{P}_{\alpha\beta}(\omega) = - \int d^3\boldsymbol{\theta} d^3\mathbf{J} \Phi_\alpha^*[\mathbf{x}(\boldsymbol{\theta}, \mathbf{J})] \tilde{f}_\beta(\boldsymbol{\theta}, \mathbf{J}, \omega). \quad (5.91)$$

We next expand the potential basis functions in a Fourier series (cf. eq. 5.79)

$$\Phi_\alpha[\mathbf{x}(\boldsymbol{\theta}, \mathbf{J})] = \sum_{\mathbf{m}} \Phi_{\alpha, \mathbf{m}}(\mathbf{J}) e^{i\mathbf{m} \cdot \boldsymbol{\theta}}, \quad (5.92)$$

and employ equation (5.83) to evaluate \tilde{f}_β :

$$\tilde{P}_{\alpha\beta}(\omega) = - \sum_{\mathbf{m}, \mathbf{m}'} \int d^3\boldsymbol{\theta} d^3\mathbf{J} \mathbf{m} \cdot \frac{\partial f_0(\mathbf{J})}{\partial \mathbf{J}} \frac{\tilde{\Phi}_{\alpha, \mathbf{m}'}^*(\mathbf{J}) \tilde{\Phi}_{\beta, \mathbf{m}}(\mathbf{J})}{\mathbf{m} \cdot \boldsymbol{\Omega} - \omega} e^{i(\mathbf{m} - \mathbf{m}') \cdot \boldsymbol{\theta}}. \quad (5.93)$$

The integration over $d^3\boldsymbol{\theta}$ vanishes unless the integer triples \mathbf{m} and \mathbf{m}' are equal, so the expression simplifies to

$$\tilde{P}_{\alpha\beta}(\omega) = -(2\pi)^3 \sum_{\mathbf{m}} \int d^3\mathbf{J} \mathbf{m} \cdot \frac{\partial f_0(\mathbf{J})}{\partial \mathbf{J}} \frac{\tilde{\Phi}_{\alpha,\mathbf{m}}^*(\mathbf{J}) \tilde{\Phi}_{\beta,\mathbf{m}}(\mathbf{J})}{\mathbf{m} \cdot \boldsymbol{\Omega} - \omega}. \quad (5.94)$$

An equally valid form is obtained by applying the divergence theorem (B.45):⁸

$$\tilde{P}_{\alpha\beta}(\omega) = (2\pi)^3 \sum_{\mathbf{m}} \int d^3\mathbf{J} f_0(\mathbf{J}) \mathbf{m} \cdot \frac{\partial}{\partial \mathbf{J}} \left[\frac{\tilde{\Phi}_{\alpha,\mathbf{m}}^*(\mathbf{J}) \tilde{\Phi}_{\beta,\mathbf{m}}(\mathbf{J})}{\mathbf{m} \cdot \boldsymbol{\Omega} - \omega} \right]. \quad (5.95)$$

All these expressions should be restricted to $\text{Im}(\omega) > c$.

A mode requires that $\rho_{s1} = \rho_1$, so by (5.85) $\mathbf{d} = \mathbf{g}$, and by (5.88)

$$\tilde{\mathbf{P}}(\omega) \tilde{\mathbf{d}} = \tilde{\mathbf{d}}, \quad \text{Im}(\omega) > c; \quad (5.96)$$

in words, a mode with frequency ω_p exists if $\tilde{\mathbf{d}}$ is an eigenvector of $\tilde{\mathbf{P}}(\omega_p)$ with eigenvalue 1. This matrix equation is the direct analog of the dispersion relation (5.45) that we derived for homogeneous stellar systems.

5.3.3 The response matrix

The response matrix is defined by analogy to the polarization matrix (5.87b),

$$\tilde{R}_{\alpha\beta}(\omega) = - \int d^3\mathbf{x} d^3\mathbf{x}' \Phi_{\alpha}^*(\mathbf{x}) \tilde{R}(\mathbf{x}, \mathbf{x}', \omega) \rho_{\beta}(\mathbf{x}'). \quad (5.97)$$

If we expand the external density in the bi-orthonormal basis (cf. eq. 5.85)

$$\tilde{\rho}_e(\mathbf{x}, \omega) = \sum_{\beta} \tilde{h}_{\beta}(\omega) \rho_{\beta}(\mathbf{x}), \quad (5.98)$$

then as in equation (5.88),

$$\tilde{\mathbf{d}}(\omega) = \tilde{\mathbf{R}}(\omega) \tilde{\mathbf{h}}(\omega). \quad (5.99)$$

The response matrix can be determined from the polarization matrix by generalizing the arguments of §§5.2.2 and 5.2.4 from scalar functions of ω to matrices (Kalnajs 1971). The total perturbing density is the sum of the external density and the perturbed density of the stellar system; thus $\rho_1 = \rho_{s1} + \rho_e$ and $\mathbf{g} = \mathbf{d} + \mathbf{h}$. Using equations (5.88) and (5.99) we then have

$$\tilde{\mathbf{d}} = \tilde{\mathbf{P}}(\omega) \tilde{\mathbf{g}} = \tilde{\mathbf{P}}(\omega) \tilde{\mathbf{h}} + \tilde{\mathbf{P}}(\omega) \tilde{\mathbf{d}} = \tilde{\mathbf{P}}(\omega) \tilde{\mathbf{h}} + \tilde{\mathbf{P}}(\omega) \tilde{\mathbf{R}}(\omega) \tilde{\mathbf{h}}. \quad (5.100)$$

⁸The surface term vanishes; see Problem 5.11.

Since this relation holds for arbitrary $\tilde{\mathbf{h}}$, comparison with equation (5.99) implies that

$$\tilde{\mathbf{R}}(\omega) = \tilde{\mathbf{P}}(\omega) + \tilde{\mathbf{P}}(\omega)\tilde{\mathbf{R}}(\omega) \quad \text{or} \quad [\mathbf{I} - \tilde{\mathbf{P}}(\omega)]\tilde{\mathbf{R}}(\omega) = \tilde{\mathbf{P}}(\omega), \quad (5.101)$$

where \mathbf{I} is the identity matrix.

Equation (5.94) shows that the matrix elements of $\mathbf{P}(\omega)$ shrink to zero as $\text{Im}(\omega) \rightarrow \infty$. Thus there exists some positive constant c such that the series $\mathbf{I} + \tilde{\mathbf{P}} + \tilde{\mathbf{P}}^2 + \dots$ converges for $\text{Im}(\omega) > c$; in this case the series is equal to $[\mathbf{I} - \tilde{\mathbf{P}}(\omega)]^{-1}$, and we can invert equation (5.101) to obtain a relation between the polarization and response matrices,

$$\tilde{\mathbf{R}}(\omega) = [\mathbf{I} - \tilde{\mathbf{P}}(\omega)]^{-1}\tilde{\mathbf{P}}(\omega), \quad \text{Im}(\omega) > c. \quad (5.102)$$

This matrix equation is the analog of the scalar equation (5.30) for homogeneous systems.

From equations (5.99) and (5.102) the basis-function expansion of the response density is

$$\tilde{\mathbf{d}}(\omega) = [\mathbf{I} - \tilde{\mathbf{P}}(\omega)]^{-1}\tilde{\mathbf{P}}(\omega)\tilde{\mathbf{h}}(\omega), \quad \text{Im}(\omega) > c. \quad (5.103)$$

Taking the inverse Fourier transform (5.6) yields

$$\mathbf{d}(t) = \frac{1}{2\pi} \int_{ic-\infty}^{ic+\infty} d\omega \tilde{\mathbf{d}}(\omega) e^{-i\omega t} = \frac{1}{2\pi} \int_{ic-\infty}^{ic+\infty} d\omega [\mathbf{I} - \tilde{\mathbf{P}}(\omega)]^{-1} \tilde{\mathbf{P}}(\omega) \tilde{\mathbf{h}}(\omega) e^{-i\omega t}, \quad (5.104)$$

which gives the response density \mathbf{d} as a function of the external density \mathbf{h} .

As in §5.2.4, we evaluate the integral by closing the contour in the lower half-plane. Then by the residue theorem (cf. eq. 5.58)

$$\mathbf{d}(\tau) = -i \sum_{\mathbf{p}} \mathbf{d}_{\mathbf{p}} e^{-i\omega_{\mathbf{p}}\tau}, \quad (5.105)$$

where $\mathbf{d}_{\mathbf{p}}$ is the residue of the integrand at the pole $\omega = \omega_{\mathbf{p}}$. These contributions can arise from poles in either $\tilde{\mathbf{P}}(\omega)$ or $[\mathbf{I} - \tilde{\mathbf{P}}(\omega)]^{-1}$. The matrix $\tilde{\mathbf{P}}(\omega)$ is non-singular for $\text{Im}(\omega) > 0$ (see eq. 5.94). The poles of $[\mathbf{I} - \tilde{\mathbf{P}}(\omega)]^{-1}$ are determined by the equation

$$\tilde{\mathbf{P}}(\omega)\tilde{\mathbf{d}} = \tilde{\mathbf{d}}, \quad (5.106)$$

where $\tilde{\mathbf{P}}(\omega)$ is now the analytic continuation of the polarization matrix. For $\text{Im}(\omega) > 0$, equation (5.95) defines $\tilde{\mathbf{P}}(\omega)$ and values of ω that satisfy (5.106) correspond to unstable modes. Analytic continuation of $\tilde{\mathbf{P}}(\omega)$ to $\text{Im}(\omega) \leq 0$

may give rise to solutions of (5.106) that are not associated with modes, but rather Landau-damped disturbances.

These matrix equations can be simplified when the stellar system has certain symmetries. For example, when the equilibrium system is spherical the density of a mode can always be assumed to have the form $\rho_1(\mathbf{x}, t) = d(r)Y_l^m(\vartheta, \phi)\exp(-i\omega t)$, where $Y_l^m(\vartheta, \phi)$ is a spherical harmonic (Appendix C.6) and $d(r)$ and ω are independent of m (see Box 5.2). Thus it is sufficient to restrict the basis-function expansion to axisymmetric ($m = 0$) basis functions, and to a single value of the angular quantum number l at a time, repeating the calculation for $l = 0, 1, 2, \dots$ as required.

The Kalnajs matrix method has been used to determine the stability of a variety of stellar systems, including spherical systems, plane-parallel systems, and differentially rotating disks (see §5.5, page 432, and §6.3.1 respectively).

5.4 The energy principle and secular stability

Energy considerations often provide a powerful technique for determining the stability of dynamical systems. For example, a ball resting at the bottom of a hemispherical bowl must be stable, because all nearby positions of the ball have higher energy. In this case it is not necessary to solve the linearized equations of motion to determine that the configuration is stable. A similar criterion can be established for continuous systems: they are certain to be stable if all neighboring configurations with the same total mass have higher energy. A system that satisfies this criterion is said to enjoy **secular stability** (Lyttleton 1953; Hunter 1977). Secular stability is sufficient but not necessary for dynamical stability, because the lower energy states may not be dynamically accessible. An example of a system that is dynamically stable but secularly unstable is a ball that rolls around a bowl at constant height, maintaining a balance between centrifugal and gravitational forces. The ball's trajectory is dynamically stable (in the absence of friction) despite the existence of states of lower energy, because the ball cannot reach these states while conserving its angular momentum (see Problem 5.12).

To find energy-based stability criteria we must first find the energy change resulting from a small perturbation to the system. Once again, we shall examine self-gravitating fluid systems in parallel with stellar systems.

5.4.1 The energy principle for fluid systems

The rate at which work is done against an external force $\mathbf{F} = -\epsilon m \nabla \Phi_e$ in moving a particle of mass m is (cf. eq. D.4)

$$-\mathbf{F} \cdot \mathbf{v} = \epsilon m \mathbf{v} \cdot \nabla \Phi_e. \quad (5.107)$$

Box 5.2: Response of spherical systems

The polarization function of a spherical fluid or stellar system, $P(\mathbf{x}, \mathbf{x}', \tau)$, inherits certain symmetry properties from the spherical symmetry of the equilibrium system. In particular, the polarization function can depend on \mathbf{x} and \mathbf{x}' only through their magnitudes $r = |\mathbf{x}|$, $r' = |\mathbf{x}'|$, and the angle γ between the two vectors, defined by $\cos \gamma = \mathbf{x} \cdot \mathbf{x}' / (rr')$. Thus it can be expanded in a complete set of angular functions,

$$P(\mathbf{x}, \mathbf{x}', \tau) = \sum_{k=0}^{\infty} K_k(r, r', \tau) P_k(\cos \gamma), \quad (1)$$

where P_k is a Legendre polynomial (Appendix C.5). Using the expansion (C.47), this can be rewritten in spherical harmonics as

$$P(\mathbf{x}, \mathbf{x}', \tau) = \frac{4\pi}{2k+1} \sum_{k=0}^{\infty} \sum_{n=-k}^k K_k(r, r', \tau) Y_k^{n*}(\vartheta', \phi') Y_k^n(\vartheta, \phi), \quad (2)$$

where $\mathbf{x} = (r, \vartheta, \phi)$, $\mathbf{x}' = (r', \vartheta', \phi')$ in spherical coordinates. (To minimize confusion between spherical coordinates and angle variables, in this chapter we reserve ϑ for the usual polar angle, and continue to use θ_i for the variable conjugate to J_i .)

Now consider a density perturbation of the form

$$\rho_1(\mathbf{x}, t) = g(r) Y_l^m(\vartheta, \phi) \exp(-i\omega t). \quad (3)$$

Substituting this result and equation (2) into the definition (5.3) of the polarization function, and using the orthonormality of the spherical harmonics (eq. C.44), we find that the response density is

$$\rho_{s1}(\mathbf{x}, t) = \frac{4\pi}{2l+1} Y_l^m(\vartheta, \phi) e^{-i\omega t} \int dr' r'^2 f(r') \int_0^{\infty} d\tau K_l(r, r', \tau) e^{i\omega\tau}. \quad (4)$$

The condition for a mode is that $\rho_1 = \rho_{s1}$. Since both are proportional to $Y_l^m(\vartheta, \phi)$ we may always choose the modes of spherical fluid or stellar systems to have this angular dependence. Moreover, when written in this form the equation determining the modes is independent of the **azimuthal wavenumber** m ; thus the frequency ω and the radial density profile $g(r)$ of the mode are also independent of m . In other words there are $2m+1$ degenerate modes for each l . Because of this symmetry, in practice we need to calculate the response of a spherical system only to axisymmetric $m=0$ disturbances.

This argument does not yield the behavior of the DF in velocity space; for this see Barnes, Goodman, & Hut (1986).

The rate of doing work on the particle is the negative of this expression, so the total rate of doing work on a fluid system is

$$\frac{dE}{dt} = -\epsilon \int d^3\mathbf{x} \rho \mathbf{v} \cdot \nabla \Phi_e, \quad (5.108)$$

where ρ is the fluid density. We shall assume that the integration volume is large enough that the density is zero at its boundary, so matter and energy cannot flow into or out of the volume.

The energy principle is most useful for static fluids, in which the unperturbed velocity $\mathbf{v}_0 = 0$, so we shall restrict ourselves to systems of this kind. Since $\mathbf{v}_0 = 0$, the rate of doing work is zero to first order in the perturbation parameter ϵ : thus, the energy change must be a quadratic, not linear, function of the perturbation strength. To go to quadratic order, we write $\mathbf{v}(\mathbf{x}, t) = \epsilon \mathbf{v}_1(\mathbf{x}, t)$, so

$$\frac{dE}{dt} = -\epsilon^2 \int d^3\mathbf{x} \rho_0 \mathbf{v}_1 \cdot \nabla \Phi_e; \quad (5.109)$$

note that we have replaced the density $\rho(\mathbf{x}, t)$ by its unperturbed value $\rho_0(\mathbf{x})$ since we are dropping terms of order ϵ^3 .

The linearized Euler equation (5.24b) reads

$$\frac{\partial \mathbf{v}_1}{\partial t} = -\nabla(h_1 + \Phi_{s1} + \Phi_e). \quad (5.110)$$

We eliminate Φ_e from equations (5.109) and (5.110) to obtain

$$\frac{dE}{dt} = \epsilon^2 \int d^3\mathbf{x} \rho_0 \mathbf{v}_1 \cdot \left[\frac{\partial \mathbf{v}_1}{\partial t} + \nabla(\Phi_{s1} + h_1) \right]. \quad (5.111)$$

The contribution of the first term in square brackets can be rewritten as

$$\epsilon^2 \int d^3\mathbf{x} \rho_0 \mathbf{v}_1 \cdot \frac{\partial \mathbf{v}_1}{\partial t} = \frac{1}{2} \epsilon^2 \frac{d}{dt} \int d^3\mathbf{x} \rho_0 \mathbf{v}_1^2. \quad (5.112)$$

To evaluate the second and third terms, we apply the divergence theorem (B.45)—the boundary terms vanish, since $\rho_0 = 0$ at large distances—and use the linearized continuity equation (5.24a):

$$\begin{aligned} \epsilon^2 \int d^3\mathbf{x} \rho_0 \mathbf{v}_1 \cdot \nabla(\Phi_{s1} + h_1) &= -\epsilon^2 \int d^3\mathbf{x} (\Phi_{s1} + h_1) \nabla \cdot (\rho_0 \mathbf{v}_1) \\ &= \epsilon^2 \int d^3\mathbf{x} (\Phi_{s1} + h_1) \frac{\partial \rho_{s1}}{\partial t}. \end{aligned} \quad (5.113)$$

Using the linearized equation of state (5.24d) we can rewrite the term involving the enthalpy as

$$\begin{aligned} \epsilon^2 \int d^3\mathbf{x} h_1 \frac{\partial \rho_{s1}}{\partial t} &= \epsilon^2 \int d^3\mathbf{x} \left(\frac{dp}{d\rho} \right)_{\rho_0} \frac{\rho_{s1}}{\rho_0} \frac{\partial \rho_{s1}}{\partial t} \\ &= \frac{1}{2} \epsilon^2 \frac{d}{dt} \int \frac{d^3\mathbf{x}}{\rho_0} \left(\frac{dp}{d\rho} \right)_{\rho_0} \rho_{s1}^2. \end{aligned} \quad (5.114)$$

The gravitational potential is related to the density by Poisson's equation,

$$\Phi_{s1}(\mathbf{x}, t) = -G \int d^3\mathbf{x}' \frac{\rho_{s1}(\mathbf{x}', t)}{|\mathbf{x} - \mathbf{x}'|}; \quad (5.115)$$

thus the term in equation (5.113) involving the potential can be written

$$\begin{aligned} \epsilon^2 \int d^3\mathbf{x} \frac{\partial \rho_{s1}}{\partial t} \Phi_{s1} &= -\epsilon^2 G \int d^3\mathbf{x} d^3\mathbf{x}' \frac{\partial \rho_{s1}(\mathbf{x}, t)}{\partial t} \frac{\rho_{s1}(\mathbf{x}', t)}{|\mathbf{x} - \mathbf{x}'|} \\ &= -\epsilon^2 G \int d^3\mathbf{x} d^3\mathbf{x}' \frac{\partial \rho_{s1}(\mathbf{x}', t)}{\partial t} \frac{\rho_{s1}(\mathbf{x}, t)}{|\mathbf{x} - \mathbf{x}'|}, \end{aligned} \quad (5.116)$$

where in the last equation we have simply exchanged the integration variables \mathbf{x} and \mathbf{x}' . Averaging the two right sides of the last equation gives

$$-\frac{1}{2} \epsilon^2 G \frac{d}{dt} \int d^3\mathbf{x} d^3\mathbf{x}' \frac{\rho_{s1}(\mathbf{x}, t) \rho_{s1}(\mathbf{x}', t)}{|\mathbf{x} - \mathbf{x}'|}. \quad (5.117)$$

Combining equations (5.112), (5.114), and (5.117) we have

$$\frac{dE}{dt} = \frac{1}{2} \epsilon^2 \frac{d}{dt} \left(\int d^3\mathbf{x} \rho_0 v_1^2 + \int \frac{d^3\mathbf{x}}{\rho_0} \left| \frac{dp}{d\rho} \right|_{\rho_0} \rho_{s1}^2 - G \int d^3\mathbf{x} d^3\mathbf{x}' \frac{\rho_{s1} \rho'_{s1}}{|\mathbf{x} - \mathbf{x}'|} \right), \quad (5.118)$$

where we have replaced $dp/d\rho$ by $|dp/d\rho|$ since $dp/d\rho$ is equal to the square of the sound speed (eq. 5.22), which must be positive for any realistic equation of state. Hydrostatic equilibrium requires $\nabla p_0 = -\rho_0 \nabla \Phi_0$ (eq. 5.17), so $dp_0 = -\rho_0 d\Phi_0$ and $|dp/d\rho|_{\rho_0}$ can be replaced by $\rho_0 |d\Phi/d\rho|_0$.

Finally, assuming that the perturbations E , \mathbf{v}_1 and ρ_{s1} all vanish as $t \rightarrow -\infty$, we can integrate with respect to time to obtain our final expression for the energy,

$$\begin{aligned} E &= \frac{1}{2} \epsilon^2 \left[\int d^3\mathbf{x} \rho_0(\mathbf{x}) \mathbf{v}_1^2(\mathbf{x}, t) + \int d^3\mathbf{x} \left| \frac{d\Phi}{d\rho} \right|_0 \rho_{s1}^2(\mathbf{x}, t) \right. \\ &\quad \left. - G \int d^3\mathbf{x} d^3\mathbf{x}' \frac{\rho_{s1}(\mathbf{x}, t) \rho_{s1}(\mathbf{x}', t)}{|\mathbf{x} - \mathbf{x}'|} \right]. \end{aligned} \quad (5.119)$$

Notice that we have successfully derived an expression for the energy that is correct to order ϵ^2 using the linearized equations (5.24), which are correct only to order ϵ . Calculating the second-order perturbations in the density, velocity, etc., proved to be unnecessary for this task.

Now in any unstable mode, all of the perturbed quantities such as ρ_{s1} and \mathbf{v}_1 grow as $\exp(\gamma t)$, $\gamma > 0$. Therefore all of the terms on the right side of equation (5.119) grow as $\exp(2\gamma t)$. However, the energy E is constant and cannot grow; hence *the energy of any unstable mode must be zero*. Since the first term on the right side is positive, the sum of the second and third terms must be negative. Thus a sufficient condition for stability is that the sum of these two terms is non-negative for all functions $\rho_{s1}(\mathbf{x})$. This condition can be sharpened by restricting the range of possible functions for $\rho_{s1}(\mathbf{x})$ to those for which $\int d^3\mathbf{x} \rho_{s1}(\mathbf{x}) = 0$, since mass must be conserved in any physical perturbation. Thus we arrive at

Chandrasekhar's variational principle: *A barotropic fluid in static equilibrium with $d\rho(\rho)/d\rho > 0$ is stable if the quantity*

$$W[\rho_1] \equiv \int d^3\mathbf{x} \left| \frac{d\Phi}{d\rho} \right|_0 \rho_1^2(\mathbf{x}) - G \int d^3\mathbf{x} d^3\mathbf{x}' \frac{\rho_1(\mathbf{x})\rho_1(\mathbf{x}')}{|\mathbf{x} - \mathbf{x}'|} \quad (5.120)$$

is non-negative for all real functions $\rho_1(\mathbf{x})$ that conserve the total mass, $\int d^3\mathbf{x} \rho_1(\mathbf{x}) = 0$.

The constraint that the trial functions $\rho_1(\mathbf{x})$ must conserve the mass can be incorporated automatically by writing $\rho_1 = -\nabla \cdot (\rho_0 \boldsymbol{\xi})$ (eq. F.5) and examining arbitrary displacement vector fields $\boldsymbol{\xi}(\mathbf{x})$. The expression (5.120) is zero for $\boldsymbol{\xi} = \text{constant}$, as it must be, since this corresponds to a uniform displacement of the entire system (Problem 5.14).

This variational principle was first stated by Chandrasekhar (1963, 1964). The condition $W[\rho_1] \geq 0$ is necessary as well as sufficient for stability (Laval, Mercier, & Pellat 1965; Kulsrud & Mark 1970).

5.4.2 The energy principle for stellar systems

In stellar systems, the energy principle is most useful if the equilibrium DF is ergodic, that is, if it can be written in the form $f_0(H_0)$, where $H_0(\mathbf{x}, \mathbf{v}) \equiv \frac{1}{2}v^2 + \Phi_0(\mathbf{x})$ is the Hamiltonian for motion in the unperturbed potential Φ_0 (§4.2.1a). We restrict ourselves to systems of this kind; in practice these normally have spherical or plane-parallel symmetry (Box 4.1). We shall also assume that $f'_0(H_0) = df_0/dH_0 < 0$ everywhere, a condition that is satisfied by most realistic models of stellar systems.

By analogy with equation (5.108), the rate at which an external potential $\epsilon\Phi_e$ does work on a stellar system is

$$\frac{dE}{dt} = -\epsilon \int d^3\mathbf{x} d^3\mathbf{v} f(\mathbf{x}, \mathbf{v}, t) \mathbf{v} \cdot \nabla \Phi_e ; \quad (5.121)$$

where, as usual in this chapter, the DF $f(\mathbf{x}, \mathbf{v}, t)$ is defined to be the phase-space mass density.

Since $\epsilon \ll 1$ we can expand the DF in powers of ϵ , $f(\mathbf{x}, \mathbf{v}, t) = f_0(H_0) + \epsilon f_1(\mathbf{x}, \mathbf{v}, t) + \dots$. The contribution of f_0 to dE/dt is zero, since H_0 is an even function of \mathbf{v} while $\mathbf{v} \cdot \nabla \Phi_e$ is odd. Thus the dominant energy change is

$$\frac{dE}{dt} = -\epsilon^2 \int d^3\mathbf{x} d^3\mathbf{v} f_1(\mathbf{x}, \mathbf{v}, t) \mathbf{v} \cdot \nabla \Phi_e ; \quad (5.122)$$

once again the energy change depends quadratically on the perturbation strength.

To evaluate this formula we use the linearized collisionless Boltzmann equation (5.13),

$$\begin{aligned} 0 &= \frac{\partial f_1}{\partial t} + [f_1, H_0] + [f_0, \Phi_{s1} + \Phi_e] \\ &= \frac{\partial f_1}{\partial t} + [f_1, H_0] - f'_0(H_0) \mathbf{v} \cdot \nabla (\Phi_{s1} + \Phi_e). \end{aligned} \quad (5.123)$$

We may now eliminate Φ_e from equation (5.122):

$$\frac{dE}{dt} = -\epsilon^2 \int \frac{d^3\mathbf{x} d^3\mathbf{v}}{f'_0(H_0)} f_1 \left(\frac{\partial f_1}{\partial t} + [f_1, H_0] \right) + \epsilon^2 \int d^3\mathbf{x} d^3\mathbf{v} f_1 \mathbf{v} \cdot \nabla \Phi_{s1}. \quad (5.124)$$

The middle term can be written

$$\int \frac{d^3\mathbf{x} d^3\mathbf{v}}{f'_0(H_0)} f_1 [f_1, H_0] = \frac{1}{2} \int \frac{d^3\mathbf{x} d^3\mathbf{v}}{f'_0(H_0)} [f_1^2, H_0]. \quad (5.125)$$

Furthermore,

$$\frac{1}{f'_0(H_0)} [f_1^2, H_0] = [f_1^2, h(H_0)] \quad (5.126)$$

where $h(H_0) \equiv \int^{H_0} dx/f'_0(x)$ (eq. 5.210c). Now since f_1 vanishes at infinity, we know that $\int d^3\mathbf{x} d^3\mathbf{v} [f_1^2, h(H_0)] = 0$ (eq. 5.210a). Thus equation (5.125) is zero, and equation (5.124) simplifies to

$$\frac{dE}{dt} = -\frac{1}{2}\epsilon^2 \int \frac{d^3\mathbf{x} d^3\mathbf{v}}{f'_0(H_0)} \frac{\partial f_1^2}{\partial t} - \epsilon^2 \int d^3\mathbf{x} \Phi_{s1} \nabla \cdot \int d^3\mathbf{v} f_1 \mathbf{v}, \quad (5.127)$$

where the divergence theorem (B.45) has been used on the second term (the boundary terms vanish since the DF vanishes at large distances).

We next integrate the linearized collisionless Boltzmann equation (5.123) over velocity, to obtain an equation analogous to the Jeans continuity equation (4.204),

$$\frac{\partial \rho_{s1}}{\partial t} + \nabla \cdot \int d^3\mathbf{v} f_1 \mathbf{v} = 0, \quad (5.128)$$

where $\rho_{s1} = \int d^3\mathbf{v} f_1$. Therefore equation (5.127) can be rewritten as

$$\frac{dE}{dt} = -\frac{1}{2}\epsilon^2 \frac{d}{dt} \int \frac{d^3\mathbf{x} d^3\mathbf{v}}{f'_0(H_0)} f_1^2 + \epsilon^2 \int d^3\mathbf{x} \frac{\partial \rho_{s1}}{\partial t} \Phi_{s1}. \quad (5.129)$$

Since $f'(H_0) < 0$ by assumption, we may replace it with $-|f'(H_0)|$. Following the logic that took us from equation (5.116) to (5.117), and integrating with respect to time, we finally obtain

$$E = \frac{1}{2}\epsilon^2 \left[\int \frac{d^3\mathbf{x} d^3\mathbf{v}}{|f'_0(H_0)|} f_1^2(\mathbf{x}, \mathbf{v}, t) - G \int d^3\mathbf{x} d^3\mathbf{x}' \frac{\rho_{s1}(\mathbf{x}, t) \rho_{s1}(\mathbf{x}', t)}{|\mathbf{x} - \mathbf{x}'|} \right]. \quad (5.130)$$

As for a fluid system, the energy of any unstable mode must be zero. Thus a sufficient condition for stability is that E is always positive (or always negative; but this cannot occur since E is positive for any f_1 that is an odd function of \mathbf{v} —the first term in 5.130 is positive and the second is zero). Thus we have the

Variational principle for stellar systems: *A stellar system having an ergodic DF $f_0(H_0)$ with $f'_0(H_0) < 0$ is stable if the quantity*

$$W[f_1] \equiv \int \frac{d^3\mathbf{x} d^3\mathbf{v}}{|f'_0(H_0)|} f_1^2(\mathbf{x}, \mathbf{v}) - G \int \frac{d^3\mathbf{x} d^3\mathbf{v} d^3\mathbf{x}' d^3\mathbf{v}'}{|\mathbf{x} - \mathbf{x}'|} f_1(\mathbf{x}, \mathbf{v}) f_1(\mathbf{x}', \mathbf{v}') \quad (5.131)$$

is non-negative for all functions $f_1(\mathbf{x}, \mathbf{v})$ that conserve the total mass, that is, $\int d^3\mathbf{x} d^3\mathbf{v} f_1(\mathbf{x}, \mathbf{v}) = 0$. As in the case of fluid systems, the trivial perturbation $f_1 = -\boldsymbol{\xi} \cdot \nabla f_0$ corresponds to a uniform displacement of the entire system and yields $W[f_1] = 0$.

We can sharpen this result, because not all trial functions f_1 are physically plausible. In particular, f_1 must arise from the displacement of stars from their original phase-space positions and by analogy with equation (F.5) we may write

$$\epsilon f_1 = -\frac{\partial}{\partial \mathbf{x}} \cdot (f_0 \Delta \mathbf{x}) - \frac{\partial}{\partial \mathbf{v}} \cdot (f_0 \Delta \mathbf{v}). \quad (5.132)$$

Furthermore the displacements $(\Delta \mathbf{x}, \Delta \mathbf{v})$ must be due to the motion of the stars under the influence of some Hamiltonian. The time-evolution operator that describes Hamiltonian motion is canonical (see discussion on page 803), so the transformation from the original phase-space position (\mathbf{x}, \mathbf{v}) to the new position $(\mathbf{x}', \mathbf{v}') \equiv (\mathbf{x} + \Delta \mathbf{x}, \mathbf{v} + \Delta \mathbf{v})$ is a canonical transformation. The generating function for this transformation can be written $S(\mathbf{x}, \mathbf{v}') = \mathbf{x} \cdot \mathbf{v}' + \epsilon g(\mathbf{x}, \mathbf{v}')$. Using equations (D.93) we have

$$\mathbf{x}' = \mathbf{x} + \epsilon \frac{\partial g}{\partial \mathbf{v}'} \quad ; \quad \mathbf{v} = \mathbf{v}' + \epsilon \frac{\partial g}{\partial \mathbf{x}}, \quad (5.133)$$

or

$$\Delta \mathbf{x} = \epsilon \frac{\partial g(\mathbf{x}, \mathbf{v})}{\partial \mathbf{v}} \quad ; \quad \Delta \mathbf{v} = -\epsilon \frac{\partial g(\mathbf{x}, \mathbf{v})}{\partial \mathbf{x}}, \quad (5.134)$$

where we have replaced \mathbf{v}' by \mathbf{v} on the right side of these equations since the difference between the two velocities is small.

Equation (5.132) now becomes

$$f_1 = -\frac{\partial}{\partial \mathbf{x}} f_0 \frac{\partial g}{\partial \mathbf{v}} + \frac{\partial}{\partial \mathbf{v}} f_0 \frac{\partial g}{\partial \mathbf{x}} = [g, f_0]. \quad (5.135)$$

Moreover, all functions of this form conserve total mass, $\int d^3 \mathbf{x} d^3 \mathbf{v} f_1 = 0$ (by eq. 5.210a). The condition (5.135) severely restricts the allowable trial functions f_1 , and with this restriction Antonov (1960) was able to show that the variational principle (5.131) provided both necessary and sufficient conditions for stability. Thus we have

Antonov's variational principle: *A stellar system having an ergodic DF $f_0(H_0)$ with $f'_0(H_0) < 0$ is stable if and only if the quantity*

$$W_A[g] \equiv \int \frac{d^3 \mathbf{x} d^3 \mathbf{v}}{|f'_0(H_0)|} [g, f_0]^2 - G \int \frac{d^3 \mathbf{x} d^3 \mathbf{v} d^3 \mathbf{x}' d^3 \mathbf{v}'}{|\mathbf{x} - \mathbf{x}'|} [g, f_0]_{\mathbf{x}, \mathbf{v}} [g, f_0]_{\mathbf{x}', \mathbf{v}'} \quad (5.136)$$

is non-negative for all functions $g(\mathbf{x}, \mathbf{v})$.

We have already argued that $W_A[g] \geq 0$ is *sufficient* for stability. Antonov's original proof of this result takes a different but instructive route, which is traced in Problem 5.15. The proof that $W_A[g] \geq 0$ is *necessary* for stability is more delicate. Antonov's proof assumed that the stellar system has a complete set of modes, i.e., that every perturbation can be written as a sum of modes, but this assumption is difficult to justify for a stellar system (see, for example, Box 5.1). An alternative and more satisfactory proof is given by Kulsrud & Mark (1970); for a general discussion of energy-based criteria for stability of stellar systems see Bartholomew (1971) and references therein.

A simple corollary of Antonov's variational principle, proved in Problem 5.15, is that *the modes of a stellar system having an ergodic DF $f_0(H_0)$ with $f'_0(H_0) < 0$ have frequency ω such that ω^2 is real*. In other words the modes either have real frequencies, corresponding to undamped oscillations, or imaginary frequencies, corresponding to growing modes. There are no growing oscillations (overstabilities).

Finally, we state without proof

Goodman's variational principle: *An equilibrium stellar system in which the DF is invariant under velocity reversal, $f(\mathbf{x}, \mathbf{v}) = f(\mathbf{x}, -\mathbf{v})$, and the potential $\Phi_0(\mathbf{x})$ is integrable is unstable if the quantity*

$$\frac{-\int d^3 \mathbf{x} d^3 \mathbf{x}' \Phi^*(\mathbf{x}) \tilde{P}(\mathbf{x}, \mathbf{x}', is) \rho(\mathbf{x}')}{-\int d^3 \mathbf{x} \Phi^*(\mathbf{x}) \rho(\mathbf{x})} > 1 \quad (5.137)$$

for any $s > 0$. The restrictions on the trial functions $\rho(\mathbf{x})$, $\Phi(\mathbf{x})$ are that $\int d^3\mathbf{x} \rho(\mathbf{x}) = 0$, and $\Phi(\mathbf{x})$ is the gravitational potential corresponding to an isolated mass distribution with density $\rho(\mathbf{x})$.

The proof is given by Goodman (1988). Goodman's variational principle is mainly important for the conceptual insight it provides rather than as a practical tool for determining stability—it requires the polarization function, and once this is known it is not too hard to calculate the modes of the system directly rather than through a variational principle. Loosely speaking, Goodman's principle formalizes the intuitive feeling that a system is unstable if the response ρ_{s1} is greater than the stimulus ρ_1 .

5.4.3 The relation between the stability of fluid and stellar systems

The analogy between Antonov's variational principle (5.136) for stellar systems and Chandrasekhar's variational principle (5.120) for fluids can be made precise by

Antonov's first law: *A stellar system having an ergodic DF $f_0(H_0)$ with $f'_0(H_0) < 0$ is stable if the barotropic star with the same equilibrium density distribution is stable.*

To prove this, consider the function

$$\rho_1(\mathbf{x}) = \int d^3\mathbf{v} [g, f_0]. \quad (5.138)$$

We now use Schwarz's inequality (B.75), with $A = |f'_0(H_0)|^{1/2}$ and $B = [g, f_0]/|f'_0(H_0)|^{1/2}$, which yields

$$\int d^3\mathbf{v} \frac{[g, f_0]^2}{|f'_0(H_0)|} \geq \frac{(\int d^3\mathbf{v} [g, f_0])^2}{\int d^3\mathbf{v} |f'_0(H_0)|} = \frac{\rho_1^2}{\int d^3\mathbf{v} |f'_0(H_0)|}. \quad (5.139)$$

The denominator on the right side of (5.139) can be related to the unperturbed density $\rho_0(\Phi_0) = \int d^3\mathbf{v} f_0(\frac{1}{2}v^2 + \Phi_0)$ by differentiating both sides of the latter expression with respect to Φ_0 :

$$\left(\frac{d\rho}{d\Phi}\right)_0 = \int d^3\mathbf{v} f'_0(H_0) = - \int d^3\mathbf{v} |f'_0(H_0)|. \quad (5.140)$$

Combining equations (5.136)–(5.140), we obtain

$$W_A[g] \geq \int d^3\mathbf{x} \left|\frac{d\Phi}{d\rho}\right|_0 \rho_1^2 - G \int \frac{d^3\mathbf{x} d^3\mathbf{x}'}{|\mathbf{x} - \mathbf{x}'|} \rho_1(\mathbf{x}) \rho_1(\mathbf{x}'). \quad (5.141)$$

If the barotropic fluid with density $\rho_0(r)$ is stable, then the right side is non-negative by Chandrasekhar's variational principle. Hence $W_A[g] \geq 0$ and the stellar system is stable by Antonov's variational principle.

5.5 The response of spherical systems

The tools developed in §5.3 can be used to investigate the linear response and stability of a wide variety of stellar systems, so long as the Hamiltonian that describes motion in the equilibrium potential is integrable. The simplest of these are plane-parallel and spherical stellar systems. Plane-parallel systems are mostly used to model the properties of galactic disks in the direction normal to the galactic plane; an example is the isothermal sheet described in Problem 4.21. For stability analyses of plane-parallel systems see Kulsrud & Mark (1970), Mark (1971), Antonov (1971), Kalnajs (1973a), Fridman & Polyachenko (1984), Araki (1985), Mathur (1990), and Weinberg (1991b). For the sake of brevity, in this section we focus exclusively on spherical systems.

As shown in Box 5.2, the modes of spherical fluid or stellar systems can always be chosen to have density distributions with angular dependence proportional to a spherical harmonic $Y_l^m(\vartheta, \phi)$, and the frequency of the mode associated with (l, m) is independent of m . Modes with $l = m = 0$ are spherically symmetric and hence are called **radial** modes, while modes with $l \geq 1$ are called **non-radial**.

5.5.1 The stability of spherical systems with ergodic DFs

We begin by examining the stability of spherical systems in which the equilibrium DF is ergodic, $f_0 = f_0(H_0) = f_0[\frac{1}{2}v^2 + \Phi_0(r)]$. Once again we shall exploit the analogy with barotropic fluids.

Chandrasekhar's variational principle (eq. 5.120) can be used to derive a remarkably general result that was proved independently by Antonov (1962b) in the context of stellar systems and by Lebovitz (1965) for stars. This is the

Antonov–Lebovitz theorem: *All non-radial modes of a barotropic star with $dp(\rho)/d\rho > 0$ are stable.*

The proof is given in Appendix H.

This theorem shows that only radial modes are dangerous for the stability of a spherical star.

Now any spherical stellar system with an ergodic DF that satisfies $f_0'(H_0) < 0$ must have $d\rho_0/dr < 0$ —this follows from equation (5.140) and the observation that $d\Phi_0/dr > 0$ since the gravitational force in a spherical system is always radially inward. Thus an immediate consequence of the Antonov–Lebovitz theorem and Antonov's first law is

Antonov's second law: *All non-radial modes of a stellar system having an ergodic equilibrium DF $f_0(H_0)$ with $f_0'(H_0) < 0$ are stable.*

Unfortunately, Antonov's first law is not very helpful when we consider the stability of stellar systems to *radial* modes. The barotropic analogs of

realistic stellar systems are usually unstable to radial modes, so the theorem, which provides only sufficient but not necessary conditions for stability, does not constrain the stability of the analogous stellar systems. A much more powerful tool is provided by the

Doremus–Feix–Baumann theorem: *All radial modes of a stellar system with an ergodic equilibrium DF $f_0(H_0)$ and $f'_0(H_0) < 0$ are stable.*

The proof is given in Appendix I.⁹

The Doremus–Feix–Baumann theorem, combined with Antonov’s second law, shows that almost *all* realistic spherical stellar systems with an ergodic DF are stable. Thus, for example, all of the polytropic stellar systems (§4.3.3a) with $n > \frac{3}{2}$ are stable.¹⁰ In particular, the Plummer model ($n = 5$) is stable; this is important because the density distribution in the Plummer model is similar to the density distribution in many real stellar systems. The isothermal sphere (§4.3.3b) is also stable.

The ergodic DF for the isochrone model (eq. 2.47) is given by equation (4.54); as shown in Figure 4.2 this model has $f'_0(H_0) < 0$ and hence is stable, as are the Jaffe and Hernquist models shown in the same figure. The Jaffe and Hernquist models are special cases of the Dehnen models (eq. 2.64 with $\beta = 4$); it is straightforward to show that all Dehnen models with ergodic DFs have $f'_0(H_0) < 0$ and hence are stable (Tremaine et al. 1994). Finally, the DFs of King models (eq. 4.110) decrease with increasing energy, so all King models are stable.

5.5.2 The stability of anisotropic spherical systems

In most spherical stellar systems the DF depends both on energy and angular momentum, $f_0 = f_0(H_0, L)$. In such systems, the variational principles that we have proved cannot be used to establish stability to non-radial modes. However, a modified version of these principles can still be applied to radial modes. To show this, we note that in this case $\nabla\Phi_1 = \hat{\mathbf{e}}_r d\Phi_1/dr$, so

$$-\nabla\Phi_1 \cdot \frac{\partial f_0}{\partial \mathbf{v}} = -\frac{\partial\Phi_1}{\partial r} \frac{\partial f_0}{\partial v_r} = -\frac{\partial\Phi_1}{\partial r} v_r \frac{\partial f_0}{\partial H_0}. \quad (5.142)$$

When we substitute this expression into the left side of the linearized collisionless Boltzmann equation (5.13), no partial derivatives of the form $\partial f_0/\partial L$ appear. Hence all of the steps leading to the variational principles (5.131) and (5.136) can be carried out by simply replacing df_0/dH_0 by $\partial f_0/\partial H_0$ wherever it appears, and the Doremus–Feix–Baumann theorem can be extended to:

⁹ Our discussion here is restricted to linear stability. Nonlinear stability is discussed by Holm et al. (1985) and Rein (2007) among others.

¹⁰ Polytropes with $\frac{1}{2} < n < \frac{3}{2}$ are somewhat unrealistic because they have an integrable singularity in f_0 at the boundary $\mathcal{E} = 0$. However, numerical experiments (Hénon 1973a; Barnes, Goodman, & Hut 1986) indicate that these systems are also stable. Thus $f'_0(H_0) < 0$ is *sufficient* but not *necessary* for stability.

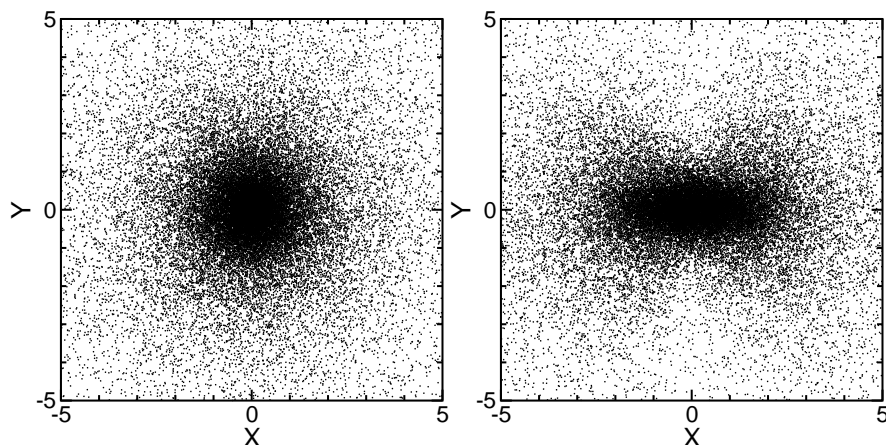


Figure 5.4 Instability in a spherical stellar system with the Hernquist density distribution (eq. 2.64). The DF has the Osipkov–Merritt form (§4.3.2c) with anisotropy radius $r_a = 0.3a$. The radial-orbit instability converts the initial spherical system on the left to the triaxial one on the right in a few crossing times. From Meza & Zamorano (1997).

All radial modes of a spherical stellar system with $\partial f_0/\partial H_0 < 0$ are stable.

For example, all Michie models (eq. 4.117) are stable to radial modes, as are the constant-anisotropy Hernquist models with $\beta = \pm \frac{1}{2}$ (eqs. 4.69 and 4.72).

The stability of anisotropic spherical systems to non-radial modes is more complicated. In contrast to ergodic systems, which are generally stable to non-radial perturbations because of Antonov’s second law, anisotropic systems with predominantly radial orbits are susceptible to the **radial-orbit instability**. As an example, Figure 5.4 shows the evolution of an initially spherical Hernquist model with anisotropy radius $r_a/a = 0.3$. There is a strong instability, which leads to a triaxial or bar-like final state, and persists in models with $r_a/a \simeq 1$.

Physical basis of the radial-orbit instability In a smooth, spherical galactic potential, stars can oscillate back and forth through the center on radial orbits. A time exposure of a star on such an orbit would show a bright, straight wire or rod, symmetric about the origin and reaching out to the star’s apocenter. The image would be faintest at the center of the galaxy, since the star travels fastest there, and brightest at apocenter, since the stellar velocity is temporarily zero at this point. Because of this geometry, rather than thinking of the orbit as having a single apocenter that rotates by π in one radial period, it is convenient to regard the orbit as having two apocenters on opposite sides of the galaxy, each of them stationary. The orientation can be specified by the azimuthal angle ψ_a of either apocenter.

Next consider a star of energy E on a nearly radial orbit, with small but

non-zero angular momentum L . Each apocenter precesses slowly, at a rate

$$\frac{d\psi_a}{dt} = \frac{\Delta\psi - \pi}{T_r} = \Omega_\phi - \frac{1}{2}\Omega_r, \quad (5.143)$$

where $\Delta\psi$ is the angle between successive apocenters, defined by equation (3.18b), T_r is the radial period (eq. 3.39b), and $\Omega_r = 2\pi/T_r$ and Ω_ϕ are the radial and azimuthal frequencies. A sequence of images of the trajectory, with exposure time long compared to T_r but short compared to the precession time $(\dot{\psi}_a)^{-1}$, would show a nearly straight wire that rotated slowly around its center at the rate $\dot{\psi}_a$. For nearly radial orbits in non-singular potentials¹¹

$$\Delta\psi = \pi + p(E)L + \text{higher-order terms in } L. \quad (5.144)$$

Thus the rate of precession of nearly radial orbits is

$$\frac{d\psi_a}{dt} = \frac{p(E)}{T_r(E)}L + \mathcal{O}(L^2), \quad (5.145)$$

where $T_r(E)$ is the radial period at zero angular momentum.

For example, in the isochrone potential (2.47), the angle between successive apocenters for an orbit with angular momentum L is¹²

$$\Delta\psi = \pi + \frac{\pi L}{\sqrt{L^2 + 4GMb}}, \quad (5.146)$$

so $p(E) = \frac{1}{2}\pi(GMb)^{-1/2}$.

Now imagine that the star is subjected to a weak non-radial gravitational field. Its angular momentum will change at an average rate $\dot{L} = N$, where N is the torque per unit mass exerted by this field, averaged over one radial period. Changes in the angular momentum affect the precession rate through equation (5.145); we have

$$\frac{d^2\psi_a}{dt^2} \simeq \frac{d}{dt} \frac{p(E)}{T_r(E)}L \simeq \frac{p(E)}{T_r(E)}N. \quad (5.147)$$

We have neglected the effect of the torque on the energy; this has a much smaller influence on the precession rate because $\dot{\psi}_a$ is proportional to L and L is small.

¹¹ This property does not hold for singular potentials; for example, in the Kepler potential the angle between successive apocenters is 2π , independent of L . The behavior of $\Delta\psi$ in singular power-law potentials is discussed in Problem 3.19.

¹² This result follows from equation (3.40), after replacing $\pi \operatorname{sgn}(L)$ by π since π and $-\pi$ are equivalent azimuths.

For comparison, the orientation of a rigid body subjected to the same torque would change at a rate given by

$$\frac{d^2\psi}{dt^2} = \frac{N}{I}, \quad (5.148)$$

where I is the moment of inertia per unit mass (eq. D.43). Thus we deduce that a near-radial orbit has an effective moment of inertia $I = T_r(E)/p(E)$; for example, in the isochrone potential $I(E) = 2\pi^{-1}T_r(E)\sqrt{GMb}$. In contrast to rigid bodies, the effective moment of inertia of nearly radial orbits can be negative, although it is positive for most galactic potentials (i.e., nearly radial orbits precess in the same direction that they revolve).¹³

The mass distribution in a spherical galaxy composed of stars on nearly radial orbits can be regarded as a collection of wires, resembling a hedgehog or porcupine. The wires have various lengths and orientations, and precess at different rates, but on average the mass distribution is spherical so there are no significant non-radial gravitational forces. Now suppose that for some reason the wires in a small solid angle are squeezed together, forming a clump with slightly higher density than the rest of the sphere. The enhanced gravitational force from the clump tends to attract more wires, whose additional mass promotes the growth of this clump. This tendency is counterbalanced by the precession of the wires in random directions, which tends to disperse the clump before it has time to grow. This competition between self-gravity and dispersion is the same one that appears in the Jeans instability.

We argued after equation (5.25) that the Jeans instability set in when the sum of the gravitational potential energy and the kinetic energy in a volume of homogeneous fluid became negative. Similarly, an approximate criterion for the radial-orbit instability in a spherical stellar system can be derived by comparing the potential and kinetic energies associated with precessional motion in a cone of opening angle ϑ . For simplicity we shall assume that the effective moment of inertia $I > 0$. The mass of the cone is $M \approx \rho r^3 \vartheta^2$ where ρ and r are the mean density and radius of the system, and the gravitational potential energy associated with the cone is $W \approx -GM^2/r \approx -G\rho^2 r^5 \vartheta^4$. The kinetic energy per unit mass associated with precession of the orbits in the cone is $\frac{1}{2}L^2/I$ (cf. eq. D.43), so the total kinetic energy is $K \approx ML^2/I \approx \rho r^3 \vartheta^2 L^2/I$. Instability sets in when $K + W \lesssim 0$, or when $\vartheta^2 \gtrsim L^2/(G\rho r^2 I)$.

An important difference between the Jeans instability and the radial-orbit instability is that infinite homogeneous systems are *always* unstable if the scale λ is sufficiently large. In contrast, the angle ϑ cannot exceed a value of order unity, so we expect that the system is stable if the typical angular momentum $L^2 \gtrsim G\rho r^2 I$, that is, tangential velocity dispersion suppresses the instability.

¹³ It is precisely because $I > 0$ for the logarithmic potential $\Phi_L(x, y)$ (eq. 3.103) that long-axis orbits are stable while short-axis orbits are unstable.

The existence of the radial-orbit instability was pointed out by Antonov (1973), who argued that any spherical system composed entirely of stars on radial orbits was unstable, and similar arguments were made by Lynden-Bell (1979) in the context of bar formation in disk galaxies. Polyachenko (1981) demonstrated the presence of the instability in N-body simulations. However, its importance was widely recognized only after it was rediscovered by Barnes (1985). Merritt & Aguilar (1985) and Meza & Zamorano (1997) have investigated the onset of the radial-orbit instability in simulations of Osipkov-Merritt models (§4.3.2b) with the Dehnen density distribution (2.64); they find that these models are unstable if the anisotropy radius r_a is less than about 40% of the half-mass radius r_h . The stability boundary can be determined more precisely by evaluating the complex frequencies of the modes of anisotropic spherical systems, using the Kalnajs matrix method (Polyachenko & Shukhman 1981; Palmer & Papaloizou 1987; Saha 1991; Weinberg 1991a; Saha 1992; see Merritt 1999 for a review).

The radial-orbit instability is sensitive to the details of the DF near $L = 0$ and the potential near the center: *any* spherical system in which the DF is unbounded as $L \rightarrow 0$ is unstable if the potential is smooth (Palmer & Papaloizou 1987); conversely, a small central point mass can suppress the instability in otherwise unstable systems (Palmer & Papaloizou 1988).

The radial-orbit instability is important because the dark halos of galaxies are believed to form hierarchically from the dissipationless collapse and merging of smaller sub-units (§9.2), and this process tends to produce systems with radially biased orbits. Thus it is likely that the radial-orbit instability operates during the collapse process. As this argument would suggest, N-body simulations of dark-halo formation in a cosmological context usually yield triaxial final states (Dubinski & Carlberg 1991; Bullock 2002; Bailin & Steinmetz 2005).

5.5.3 Landau damping and resonances in spherical systems

Stability theory addresses only one aspect of how a stellar system responds to external forces. A natural next question is whether a stable stellar system can sustain undamped oscillations, that is, can it ring like a bell? The answer is not obvious, since two simple systems we have relied on for guidance give quite different answers: the response of the infinite homogeneous stellar system that we analyzed in §5.2.4 is strongly damped at all wavenumbers greater than the Jeans wavenumber k_J , while barotropic self-gravitating systems (e.g., stars) exhibit a rich spectrum of undamped oscillations.

To examine this question, let us focus initially on radial oscillations of a spherical stellar system. The stars in such a system have radial frequencies Ω_r that lie in the range $\Omega_{r,\min}$ to $\Omega_{r,\max}$. The minimum radial frequency is non-zero so long as the system has finite extent, and the maximum radial frequency will be finite so long as the system has a smooth core.

Assume that the perturbed gravitational potential oscillates radially at some real frequency ω . The orbital radius of a star with actions \mathbf{J} oscillates with frequency $\Omega_r(\mathbf{J})$ and its harmonics $k\Omega_r(\mathbf{J})$, where k is an integer. Consequently the oscillation of the stellar system resonates with stars that satisfy the condition

$$k\Omega_r(\mathbf{J}) = \omega. \quad (5.149)$$

Resonant stars damp the oscillations of their host stellar system, because energy is transferred from the oscillation to the resonant stars, just as oscillations in an infinite homogeneous stellar system decay through Landau damping when the resonance condition (5.71) is satisfied. The collection of resonant stars can be thought of as a singular van Kampen mode, in which the perturbation is restricted to the surface of action space on which the resonant condition (5.149) is satisfied (Vandervoort 2003).

We define the **resonant spectrum** to be the set of all frequencies ω at which resonant stars are present. All oscillations with frequencies in the resonant spectrum are damped. Normally, the resonant spectrum is continuous. However, there may be frequencies that are not in the resonant spectrum. In particular, there is a gap in the spectrum for $-\Omega_{r,\min} < \omega < \Omega_{r,\min}$, and other gaps may exist if $|k|\Omega_{r,\max} < (|k| + 1)\Omega_{r,\min}$ for integer k . Oscillations with frequencies in these gaps will not be damped, since they are not in resonance with any stars in the system.

A more careful version of this argument is provided by Mathur (1990), who also proves that at least some spherical systems have radial modes with frequencies in the gaps of the continuous resonant spectrum. Thus spherical stellar systems *can* ring like a bell, although in most cases they do not.

These arguments can be extended to non-radial oscillations. In this case the general resonance condition is

$$\mathbf{m} \cdot \boldsymbol{\Omega}(\mathbf{J}) = \omega, \quad (5.150)$$

where \mathbf{m} is an integer triple and $\boldsymbol{\Omega}(\mathbf{J}) = \partial H_0 / \partial \mathbf{J}$ (eq. 3.190). In spherical systems we may choose the actions $\mathbf{J} = (J_1, J_2, J_3)$ to be respectively the z -component of angular momentum, the total angular momentum, and the radial action (Table 3.1). In this case $\Omega_1 = 0$, $\Omega_2 = \Omega_\theta$ is the azimuthal frequency, and $\Omega_3 = \Omega_r$ is the radial frequency.

The gravitational potential arising from a mode of a spherical system can be written in the form $\Phi_1(\mathbf{x}, t) = f(r)Y_l^m(\vartheta, \phi) \exp(-i\omega t)$ (see Box 5.2). When this potential is written as a Fourier series in angle-action variables, $\Phi_1(\mathbf{x}, t) = \sum_{\mathbf{m}} \Phi_{\mathbf{m}}(\mathbf{J}) \exp[i(\mathbf{m} \cdot \boldsymbol{\theta} - \omega t)]$, the only non-zero terms have $m_1 = m$ and $|m_2| \leq l$.¹⁴ Thus the resonant condition (5.150) for perturbations of order l can be written

$$m_2\Omega_\theta + m_3\Omega_r = \omega \quad (|m_2| \leq l). \quad (5.151)$$

¹⁴ The relation $m_1 = m$ follows because an azimuthal rotation by $\Delta\phi$ changes θ_1 by $\Delta\phi$ but leaves θ_2 and θ_3 unaffected; thus the two forms of Φ_1 change by $\exp(im\Delta\phi)$ and $\exp(im_1\Delta\phi)$. Since these two changes must be the same, $m = m_1$. The relation $|m_2| \leq l$

For $l = 0$ the oscillations are radial, $m_1 = m_2 = 0$, and we recover the resonant condition (5.149) with $m_3 = k$. Oscillations with $l = 1$ shift the core of the galaxy relative to its outer parts, and are sometimes called “sloshing” or “seiche” oscillations (see Problem 5.17).

Using the Kalnajs matrix method (§5.3.2), Weinberg (1994b) has shown that many spherical stellar systems exhibit low-frequency seiche oscillations that decay very slowly, having damping times of tens to hundreds of crossing times. Thus seiche oscillations may dominate the appearance of a galaxy long after the other effects of a disturbance have died away. Moreover the simple geometry and low frequency of seiche oscillations imply that they are easily excited by the strong tidal forces that are present during galaxy mergers (§8.5).

5.6 The stability of uniformly rotating systems

Stability analyses are more difficult for rotating systems than for spherical systems, for two main reasons. First, since rotating systems are usually flattened rather than spherically symmetric, both the equilibrium gravitational field corresponding to a given density distribution and the behavior of orbits in that field are more difficult to determine than in the spherical case. Second, and equally fundamental, the system now has a reservoir of rotational kinetic energy to feed any possible unstable modes.

Because of these complications, there are few general stability theorems for rotating stellar systems. Instead, we must develop our insight using simple models, N-body simulations, and numerical mode calculations. In this chapter we restrict ourselves to uniformly rotating systems, deferring the discussion of differentially rotating systems to the next chapter.

There is a rich literature, dating back to Newton, on uniformly rotating, self-gravitating fluid systems. As in the case of spherical systems, these provide useful analogs to rotating stellar systems. A brief summary of this classic topic in applied mathematics is given in §5.6.3.

5.6.1 The uniformly rotating sheet

We begin by investigating a fluid model that exhibits the effects of rotation and a flattened geometry in the simplest possible way. Our model consists of an infinite disk or sheet of zero thickness and constant surface density Σ_0 .

arises because a rotation to a new coordinate system (ϑ', ϕ') in which the equatorial plane of the coordinate system and the plane of a given orbit coincide changes $Y_l^m(\vartheta, \phi)$ into a sum of spherical harmonics $Y_l^{m'}(\vartheta', \phi')$ with the same l but different m' . Then after one radial period, $m'\phi'$ has increased by $2\pi m'\Omega_\vartheta/\Omega_r$, while $m_2\theta_2$ has increased by $2\pi m_2\Omega_\vartheta/\Omega_r$. Since these must be the same, $m' = m_2$ and thus $|m_2| \leq l$. See Tremaine & Weinberg (1984b).

The sheet occupies the plane $z = 0$, is uniform in the x - and y -directions, and rotates with constant angular velocity $\boldsymbol{\Omega} = \Omega \hat{\mathbf{e}}_z$.

We consider the response of the sheet to disturbances in its own plane. We do not examine bending or corrugation modes, deferring these to §6.6. The analysis is simplified if we work in a frame that rotates with the unperturbed sheet at $\boldsymbol{\Omega}$; thus the continuity equation (5.16), Euler's equation (5.17), and Poisson's equation (5.18) read

$$\frac{\partial \Sigma_d}{\partial t} + \nabla \cdot (\Sigma_d \mathbf{v}) = 0, \quad (5.152a)$$

$$\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} = -\frac{\nabla p}{\Sigma_d} - \nabla \Phi - 2\boldsymbol{\Omega} \times \mathbf{v} + \Omega^2(x\hat{\mathbf{e}}_x + y\hat{\mathbf{e}}_y), \quad (5.152b)$$

$$\nabla^2 \Phi = 4\pi G \Sigma \delta(z), \quad (5.152c)$$

where $\delta(z)$ is the delta function, $\mathbf{v}(x, y, t) = v_x(x, y, t)\hat{\mathbf{e}}_x + v_y(x, y, t)\hat{\mathbf{e}}_y$ is the velocity in the rotating frame, and the last two terms on the right side of equation (5.152b) are the Coriolis and centrifugal forces (eq. 3.116). The function $\Sigma(x, y, t)$ is the total surface density, composed of the surface density $\Sigma_d(x, y, t)$ of the disk and the surface density $\Sigma_e(x, y, t)$ of an external perturber that interacts with the disk only through gravitational forces. Note that the first two of equations (5.152) are defined only in the (x, y) plane but the third must hold throughout three-dimensional space. In a two-dimensional system of this kind, the pressure p is assumed to act only in the plane of the sheet and has dimensions of force per unit length. We assume that the equation of state is barotropic and write (cf. eq. 5.19)

$$p(x, y, t) = p[\Sigma_d(x, y, t)]. \quad (5.153)$$

In the unperturbed state $\Sigma = \Sigma_d = \Sigma_0$, $\mathbf{v} = 0$, and $p = p_0 = p(\Sigma_0)$. Equation (5.152a) is satisfied trivially, and equations (5.152b) and (5.152c) read

$$\nabla \Phi_0 = \Omega^2(x\hat{\mathbf{e}}_x + y\hat{\mathbf{e}}_y), \quad (5.154)$$

$$\nabla^2 \Phi_0 = 4\pi G \Sigma_0 \delta(z). \quad (5.155)$$

Since the sheet is uniform there is no preferred direction in the (x, y) plane. Hence the gravitational field $-\nabla \Phi_0$ must point in the z -direction, and it is easy to show that $\Phi_0 = 2\pi G \Sigma_0 |z|$ (Problem 2.3). Thus equation (5.154) cannot be satisfied as it stands: there are no pressure gradients or gravitational forces to balance the centrifugal force. To proceed further we must perpetrate a version of the Jeans swindle: we assume that the centrifugal force is balanced by a gravitational force that is produced by some unspecified mass distribution. The nature of this mass distribution does not concern us, since its only function is to ensure centrifugal balance in the equilibrium state.

We now consider a small perturbation of the form $\Sigma(x, y, t) = \Sigma_0 + \epsilon \Sigma_1(x, y, t)$, $\mathbf{v}(x, y, t) = \epsilon \mathbf{v}_1(x, y, t)$, etc., where ϵ is sufficiently small. Keeping only terms linear in ϵ in equations (5.152) and (5.153), we have

$$\frac{\partial \Sigma_{d1}}{\partial t} + \Sigma_0 \nabla \cdot \mathbf{v}_1 = 0, \quad (5.156a)$$

$$\frac{\partial \mathbf{v}_1}{\partial t} = -\frac{v_s^2}{\Sigma_0} \nabla \Sigma_{d1} - \nabla \Phi_1 - 2\boldsymbol{\Omega} \times \mathbf{v}_1, \quad (5.156b)$$

$$\nabla^2 \Phi_1 = 4\pi G \Sigma_1 \delta(z), \quad (5.156c)$$

where we have introduced the sound speed v_s defined by (cf. eq. 5.22)

$$v_s^2 = \left[\frac{dp(\Sigma)}{d\Sigma} \right]_{\Sigma_0}. \quad (5.157)$$

These equations are similar to equations (5.31) except that there is a Coriolis term in Euler's equation arising from the rotation, and Poisson's equation is modified to apply to a two-dimensional system.

To solve equations (5.156) we write $\Sigma_1(x, y, t) = \Sigma_a \exp[i(\mathbf{k} \cdot \mathbf{x} - \omega t)]$, $\Sigma_{d1}(x, y, t) = \Sigma_{da} \exp[i(\mathbf{k} \cdot \mathbf{x} - \omega t)]$, $\mathbf{v}_1(x, y, t) = (v_{ax} \hat{\mathbf{e}}_x + v_{ay} \hat{\mathbf{e}}_y) \exp[i(\mathbf{k} \cdot \mathbf{x} - \omega t)]$, and $\Phi_1(x, y, z = 0, t) = \Phi_a \exp[i(\mathbf{k} \cdot \mathbf{x} - \omega t)]$. With no loss of generality, we can choose the x axis to be parallel to \mathbf{k} , so $\mathbf{k} = k \hat{\mathbf{e}}_x$. First consider Poisson's equation (5.156c). For $z \neq 0$ we have $\nabla^2 \Phi_1 = 0$ but when $z = 0$, $\Phi_1 = \Phi_a \exp[i(kx - \omega t)]$. The only continuous function that satisfies both constraints and that approaches zero far from the sheet has the form

$$\Phi_1(x, y, z, t) = \Phi_a e^{i(kx - \omega t) - |kz|}. \quad (5.158)$$

To relate Φ_a to Σ_a we integrate equation (5.156c) from $z = -\zeta$ to $z = \zeta$, where ζ is a positive constant, and then let $\zeta \rightarrow 0$. Since $\partial^2 \Phi_1 / \partial x^2$ and $\partial^2 \Phi_1 / \partial y^2$ are continuous at $z = 0$, but $\partial^2 \Phi_1 / \partial z^2$ is not, the left side gives

$$\lim_{\zeta \rightarrow 0} \int_{-\zeta}^{\zeta} dz \nabla^2 \Phi_1 = \lim_{\zeta \rightarrow 0} \int_{-\zeta}^{\zeta} dz \frac{\partial^2 \Phi_1}{\partial z^2} = \lim_{\zeta \rightarrow 0} \left. \frac{\partial \Phi_1}{\partial z} \right|_{-\zeta}^{\zeta} = -2|k| \Phi_a e^{i(kx - \omega t)}. \quad (5.159)$$

The right side of (5.156c) gives

$$4\pi G \lim_{\zeta \rightarrow 0} \int_{-\zeta}^{\zeta} dz \Sigma_1 \delta(z) = 4\pi G \Sigma_a e^{i(kx - \omega t)}. \quad (5.160)$$

Hence $-2|k| \Phi_a = 4\pi G \Sigma_a$ or

$$\Phi_1(x, y, z, t) = -\frac{2\pi G \Sigma_a}{|k|} e^{i(kx - \omega t) - |kz|}. \quad (5.161)$$

Substituting for Σ_{d1} , Σ_1 , \mathbf{v}_1 , and Φ_1 in equations (5.156) we obtain

$$\begin{aligned}\omega\Sigma_{da} &= k\Sigma_0v_{ax}, \\ \omega v_{ax} &= \frac{v_s^2 k \Sigma_{da}}{\Sigma_0} - \frac{2\pi G \Sigma_a k}{|k|} + 2i\Omega v_{ay}, \\ \omega v_{ay} &= -2i\Omega v_{ax}.\end{aligned}\quad (5.162)$$

This set of equations can be solved for Σ_{da} to yield

$$\Sigma_{da} = \tilde{P}(\mathbf{k}, \omega)\Sigma_a, \quad \text{where} \quad \tilde{P}(\mathbf{k}, \omega) = \frac{2\pi G \Sigma_0 |k|}{4\Omega^2 - \omega^2 + v_s^2 k^2}, \quad (5.163)$$

where $\tilde{P}(\mathbf{k}, \omega)$ is the Fourier transform of the polarization function (eq. 5.7) that relates the response density Σ_{da} to the total density Σ_a .

In the absence of an external perturber, we require $\tilde{P}(\mathbf{k}, \omega) = 1$, which yields the dispersion relation for the uniformly rotating sheet,

$$\omega^2 = 4\Omega^2 - 2\pi G \Sigma_0 |k| + v_s^2 k^2. \quad (5.164)$$

The sheet is stable if $\omega^2 \geq 0$ and unstable if $\omega^2 < 0$.

Equation (5.164) is central to understanding the stability of disk systems. First consider the case in which the sheet is not rotating. If $\Omega = 0$, the sheet is unstable when $v_s^2 k^2 - 2\pi G \Sigma_0 |k| < 0$, in other words if

$$|k| < k_J \equiv \frac{2\pi G \Sigma_0}{v_s^2}, \quad (5.165)$$

where k_J may be thought of as the Jeans wavenumber for the sheet. There is evidently a direct analog to the classical Jeans instability (eq. 5.34) for a three-dimensional homogeneous medium: in both cases long wavelengths are subject to a gravitational instability.

Next consider the case in which the sheet rotates, but has zero sound speed. Now the sheet is unstable if $|k| > 2\Omega^2/(\pi G \Sigma_0)$. As $k \rightarrow \infty$ the perturbation grows as $\exp(\gamma t)$, where $\gamma^2 = -\omega^2 = 2\pi G \Sigma_0 |k|$. Thus the growth rate $\gamma \rightarrow \infty$ as $\lambda \rightarrow 0$: a cold disk is violently unstable on small scales. This result is easy to understand from the three-dimensional case: the growth rate of the Jeans instability in a three-dimensional fluid with zero sound speed is $\gamma = (4\pi G \rho_0)^{1/2}$ (eq. 5.34); the average density of the sheet within a sphere of radius λ is $\bar{\rho}_0 \sim \Sigma_0/\lambda$, so the growth rate of the Jeans instability in the sheet is $\gamma \sim (G\bar{\rho}_0)^{1/2} \sim (G\Sigma_0/\lambda)^{1/2} \sim (G\Sigma_0|k|)^{1/2}$.

Neither rotation nor pressure is able by itself to stabilize the sheet: a rotating sheet with zero sound speed is unstable at small wavelengths, and a non-rotating sheet with non-zero sound speed is unstable at large wavelengths. However, rotation and pressure working together *can* stabilize the

sheet: if both effects are present, the right side of the dispersion relation (5.164) is quadratic in k , with a minimum at the “most unstable wavenumber” $|k| = \pi G \Sigma_0 / v_s^2 = \frac{1}{2} k_J$. The sheet is stable at all wavelengths if this minimum is positive, which requires that

$$\frac{v_s \Omega}{G \Sigma_0} \geq \frac{1}{2} \pi = 1.5708. \quad (5.166)$$

Toomre (1964) has given a physical interpretation of this stability criterion, using arguments analogous to those used to interpret the Jeans instability in §5.2.1. Consider a small circular patch of the sheet. The radius of the patch is h and its mass $M = \pi \Sigma_0 h^2$. Now suppose that the patch radius is reduced to a fraction $(1 - \alpha)$ of its original value, where $\alpha \ll 1$. The resulting pressure perturbation will be $p_1 \approx \alpha p_0 \approx \alpha v_s^2 \Sigma_0$. The pressure force per unit mass is $\mathbf{F}_p = -\nabla p / \Sigma$, so the extra outward pressure force has magnitude $|\mathbf{F}_{p1}| \approx p_1 / (\Sigma_0 h) \approx \alpha v_s^2 / h$. Similarly, the compression leads to an extra inward gravitational force per unit mass \mathbf{F}_{g1} , where

$$|\mathbf{F}_{g1}| \approx GM\alpha/h^2 \approx G\Sigma_0\alpha. \quad (5.167)$$

In the absence of other effects, the sheet is expected to be stable if $|\mathbf{F}_{p1}|$ exceeds $|\mathbf{F}_{g1}|$, that is, if

$$h \lesssim \frac{v_s^2}{G\Sigma_0} \equiv h_1. \quad (5.168)$$

There are also internal motions in the patch, which arise from the rotation of the sheet. If we neglect external influences, the compressed region will tend to conserve spin angular momentum around its own center. The typical spin angular momentum per unit mass is $S \approx \Omega h^2$, where Ω is the angular speed of the sheet. The outward centrifugal force per unit mass is given by $|\mathbf{F}_c| \approx \Omega^2 h \approx S^2/h^3$. If S is conserved, the centrifugal force felt by each element is increased by the compression; the amount of the increase is $|\mathbf{F}_{c1}| \approx \alpha S^2/h^3 \approx \alpha \Omega^2 h$. Stability requires that $|\mathbf{F}_{c1}|$ exceeds $|\mathbf{F}_{g1}|$, or

$$h \gtrsim \frac{G\Sigma_0}{\Omega^2} \equiv h_u. \quad (5.169)$$

Equations (5.168) and (5.169) show that both small and large regions are stable, one through pressure and the other through centrifugal force. The instability is suppressed at intermediate radii if $h_u \lesssim h_1$, which requires

$$\frac{v_s \Omega}{G \Sigma_0} \gtrsim 1, \quad (5.170)$$

an order of magnitude statement of the stability criterion (5.166).

Although equation (5.166) was derived for a fluid sheet, a very similar stability criterion applies to the analogous stellar system. A razor-thin sheet of stars with a Maxwellian velocity distribution is stable if (Toomre 1964)

$$\frac{\sigma\Omega}{G\Sigma_0} \geq 1.68, \quad (5.171)$$

where σ is the one-dimensional velocity dispersion of the stars. We defer the derivation of this important result to §6.2.3, where it appears as a special case of Toomre's stability criterion for differentially rotating disks. The coefficient on the right of equation (5.171) differs by less than 7% from the coefficient in the fluid stability criterion (5.166), illustrating once again the close analogies between stellar systems and fluids.

The approximation that the sheet is razor-thin has greatly simplified the stability analysis. Nevertheless, it is possible to investigate the stability of more realistic sheets with three-dimensional structure. These sheets or disks are still uniform in the x and y directions, but have an equilibrium vertical structure $\rho_0(z)$ that is determined by the equation of state and the equation of hydrostatic equilibrium. The stability criteria derived by such analyses are generally very similar to equation (5.166). For example, Goldreich & Lynden-Bell (1965b) have analytically determined the stability of a uniformly rotating isothermal disk (equation of state $p = v_s^2\rho$, where v_s is a constant). They find that the disk is stable if

$$\frac{v_s\Omega}{G\Sigma_0} \geq 1.06, \quad (5.172)$$

a result that differs by only about 30% from equation (5.166). The reason that the idealized two-dimensional sheet works so well is that the most unstable wavelength is several times the characteristic disk thickness, and the behavior of perturbations with such relatively long wavelengths is insensitive to the details of the vertical structure.

In conclusion, the uniformly rotating sheet exhibits three important features: (i) a cold sheet is violently unstable; (ii) the sheet can be stabilized by a sound speed v_s or velocity dispersion σ that satisfies the stability criteria (5.166) or (5.171) for fluid or stellar systems, respectively; (iii) the stability properties of fluid and stellar sheets are very similar. We shall encounter all of these features again in more realistic models of uniformly and differentially rotating disks.

5.6.2 Kalnajs disks

This family, described in §4.5.2, comprises the simplest self-gravitating stellar disks and one of the few self-consistent stellar systems whose modes can

be studied analytically. Kalnajs disks have potential and surface density (eqs. 4.166 and 4.167)

$$\Sigma_0(R) = \Sigma_c \sqrt{1 - \frac{R^2}{a^2}} \quad ; \quad \Phi_0(R) = \frac{1}{2} \Omega_0^2 R^2, \quad (5.173)$$

where

$$\Omega_0^2 = \frac{\pi^2 G \Sigma_c}{2a} \quad (5.174)$$

and a is the radius of the disk edge. The mean angular speed Ω of the stars in a Kalnajs disk is independent of position, and relative to this mean speed the stars have isotropic velocity dispersion in the disk plane, given by equation (4.175).

The modes of razor-thin disks such as these can be separated into horizontal oscillations or **density waves**, in which the disk's surface density varies but it remains in its original plane, and vertical oscillations or **bending waves**, in which the surface density remains unchanged and the disk oscillates in the direction normal to its original plane. The horizontal and vertical modes of the Kalnajs disks were analyzed by Kalnajs (1972a) and Polyachenko (1977) respectively. The horizontal modes, on which we focus here, have potentials that can be written as

$$\Phi_l^m(R, \phi, t) = P_l^{|m|}(\eta) e^{i(m\phi - \omega t)} \quad \left(\begin{array}{l} l > 0, \quad |m| \leq l \\ l - m \text{ even} \end{array} \right). \quad (5.175)$$

Here $\eta = \sqrt{1 - R^2/a^2}$ and P_l^m is an associated Legendre function (Appendix C.5). The corresponding surface density can be determined from equations (2.202) and (2.204a); after adjusting to the notation of the present section, we have

$$\Sigma_l^m(R, \phi, t) = -\frac{2}{\pi^2 G a g_{lm} \eta} P_l^{|m|}(\eta) e^{i(m\phi - \omega t)} \quad \left(\begin{array}{l} l > 0, \quad |m| \leq l \\ l - m \text{ even} \end{array} \right), \quad (5.176)$$

where g_{lm} is defined by equation (2.204b).

There are three trivial zero-frequency modes of the Kalnajs disks. The two modes with $l = 1$, $m = \pm 1$ correspond to a uniform translation of the disk.¹⁵ The third zero-frequency mode has $l = 2$, $m = 0$ and corresponds to a rescaling of the outer radius a of the disk. There is also a non-trivial

¹⁵ The proof is simple. If we displace the origin to $\boldsymbol{\xi}$, the equilibrium potential in equation (5.173) becomes $\frac{1}{2} \Omega_0^2 (\mathbf{R} - \boldsymbol{\xi})^2$. Expanding to first order in the small quantity $\boldsymbol{\xi}$, the corresponding potential perturbation is just $\Phi_1(\mathbf{R}) = -\Omega_0^2 \mathbf{R} \cdot \boldsymbol{\xi}$. The potential perturbation arising from a displacement of the origin by $\boldsymbol{\xi} = \epsilon(\hat{\mathbf{e}}_x \pm i\hat{\mathbf{e}}_y)$ is therefore

$$\Phi_1(\mathbf{R}) = -\epsilon \Omega_0^2 (x \pm iy) \propto R e^{\pm i\phi}, \quad (5.177)$$

which is the same as equation (5.175) for $l = 1$, $m = \pm 1$ and $\omega = 0$.

mode with $l = 2$, $m = 0$ and non-zero frequency, corresponding to a stable pulsation in which the outer radius oscillates but the surface-density profile of the disk retains the form (5.173).

More interesting behavior is exhibited by the modes with quantum numbers $l = 2$, $m = \pm 2$, corresponding to a bar-like distortion of the disk. These modes can be analyzed relatively simply, without even using the collisionless Boltzmann equation (Kalnajs & Athanassoula–Georgala 1974).¹⁶

The surface density and potential of the mode with $l = m = 2$ are given by equations (5.175) and (5.176)¹⁷

$$\begin{aligned}\Phi_1(R, \phi, t) &= R^2 e^{i(2\phi - \omega t)} = (x + iy)^2 e^{-i\omega t} \\ \Sigma_1(R, \phi, t) &= -\frac{8}{3\pi^2 Ga} \frac{(x + iy)^2 e^{-i\omega t}}{\sqrt{1 - R^2/a^2}}.\end{aligned}\quad (5.178)$$

The unperturbed equations of motion are

$$\ddot{\mathbf{R}}_0 = -\nabla\Phi_0 = -\Omega_0^2 \mathbf{R}_0. \quad (5.179)$$

If we now add a weak potential perturbation equal to ϵ times (5.178), the perturbed orbit $\mathbf{R}(t)$ is governed by the equations

$$\ddot{\mathbf{R}} = -\nabla(\Phi_0 + \epsilon\Phi_1), \quad (5.180)$$

where the right side must be evaluated at $\mathbf{R}(t)$. We now assume that ϵ is sufficiently small that we may work to first order in ϵ and set $\mathbf{R}(t) = \mathbf{R}_0(t) + \epsilon\mathbf{R}_1(t)$. To first order in ϵ the right side of (5.180) is given by

$$\begin{aligned}[\nabla(\Phi_0 + \epsilon\Phi_1)]_{\mathbf{R}_0(t) + \epsilon\mathbf{R}_1(t)} &= (\nabla\Phi_0)_{\mathbf{R}_0(t) + \epsilon\mathbf{R}_1(t)} + \epsilon(\nabla\Phi_1)_{\mathbf{R}_0(t)} + O(\epsilon^2) \\ &= (\nabla\Phi_0)_{\mathbf{R}_0} + \epsilon[(\mathbf{R}_1 \cdot \nabla)\nabla\Phi_0]_{\mathbf{R}_0} + \epsilon(\nabla\Phi_1)_{\mathbf{R}_0} + O(\epsilon^2).\end{aligned}\quad (5.181)$$

Since $\nabla\Phi_0 = \Omega_0^2 \mathbf{R}$, we have $(\mathbf{R}_1 \cdot \nabla)\nabla\Phi_0 = \Omega_0^2 \mathbf{R}_1$. When we subtract the unperturbed equations of motion (5.179), the perturbed equations become

$$\ddot{\mathbf{R}}_1 + \Omega_0^2 \mathbf{R}_1 = -(\nabla\Phi_1)_{\mathbf{R}_0(t)}, \quad (5.182)$$

which is the equation of motion for a driven harmonic oscillator. We now write $\mathbf{R}_1(t) \equiv [x_1(t), y_1(t)]$ and seek the solution of this equation for the

¹⁶ Even the finite-amplitude oscillations corresponding to these quantum numbers can be described analytically: for real ω these are simply the Freeman bars discussed in Box 4.2.

¹⁷ We do not have to consider modes with $m < 0$ since $\text{Re}\{\Phi_a(R) \exp[i(m\phi - \omega t)]\} = \text{Re}\{\Phi_a^*(R) \exp[i(-m\phi + \omega^*t)]\}$. Thus every mode with $m < 0$ and frequency ω is also a mode with $m > 0$ and frequency $-\omega^*$.

unperturbed trajectory that passes through (x_i, y_i) with velocity (v_{xi}, v_{yi}) at time t_i . We find

$$x_1(t) = \frac{x_i + iy_i - i(v_{xi} + iv_{yi})/\Omega_0}{\omega^2 - 2\Omega_0\omega} e^{i(\Omega_0 t - \Omega_0 t_i - \omega t)} + \frac{x_i + iy_i + i(v_{xi} + iv_{yi})/\Omega_0}{\omega^2 + 2\Omega_0\omega} e^{i(-\Omega_0 t + \Omega_0 t_i - \omega t)}, \quad (5.183)$$

in the analogous equation for $y_1(t)$ the right side is multiplied by i . If we now set $t = t_i$ we obtain a relation between the displacement (x_1, y_1) and the unperturbed phase-space position of a star,

$$x_1(t) = -iy_1(t) = 2 \frac{\omega(x + iy) - 2i(v_x + iv_y)}{\omega(\omega^2 - 4\Omega_0^2)} e^{-i\omega t}. \quad (5.184)$$

To first order in the perturbation, the equation of continuity may be written in the form (eq. F.5)

$$\Sigma_1 + \nabla \cdot (\Sigma_0 \boldsymbol{\xi}) = 0, \quad (5.185)$$

where $\boldsymbol{\xi}$ is the mean displacement of the stars at a given position. We may write $\boldsymbol{\xi} = \bar{x}_1 \hat{\mathbf{e}}_x + \bar{y}_1 \hat{\mathbf{e}}_y$, where the bar, as usual, denotes the average over velocities of the stars at a given position.

The unperturbed mean velocity at (x, y) is $\bar{x}_0 = -\Omega y$, $\bar{y}_0 = \Omega x$. The mean displacement is obtained by averaging equation (5.184) over all velocities at a given point,

$$\bar{x}_1(t) = -i\bar{y}_1(t) = 2 \frac{(\omega + 2\Omega)(x + iy)}{\omega(\omega^2 - 4\Omega_0^2)} e^{-i\omega t}. \quad (5.186)$$

We now substitute into the continuity equation (5.185) using equation (5.173) for Σ_0 . It is straightforward to show that

$$\Sigma_1 = \frac{2\Sigma_c}{a^2} \frac{\omega + 2\Omega}{\omega(\omega^2 - 4\Omega_0^2)} \frac{(x + iy)^2 e^{-i\omega t}}{\sqrt{1 - R^2/a^2}}. \quad (5.187)$$

The frequency of the mode is determined by the requirement that equation (5.178) be consistent with (5.187). Eliminating Σ_1 from these equations and using equation (5.174), we find the equation for the frequencies of the modes with quantum numbers $(l, m) = (2, 2)$:

$$\omega^3 - \frac{5}{2}\Omega_0^2\omega + 3\Omega\Omega_0^2 = 0. \quad (5.188)$$

To analyze stability we rewrite this as $\omega^3 - \frac{5}{2}\Omega_0^2\omega = -3\Omega\Omega_0^2$. The cubic polynomial on the left has a local minimum of $-2(\frac{5}{6})^{3/2}\Omega_0^3$ at $\omega = (\frac{5}{6})^{1/2}\Omega_0$.

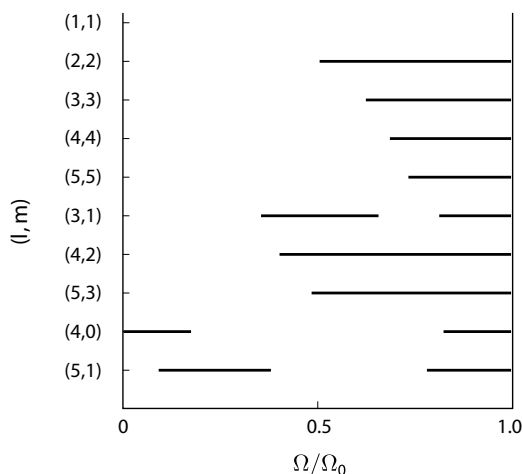


Figure 5.5 Stability of Kalnajs disks as a function of the degree of centrifugal support Ω/Ω_0 . The solid lines represent zones of instability.

If the right side is less than this local minimum, then there is only one real root, and one of the pair of complex roots corresponds to an unstable mode. Hence the disk is stable if and only if $-2(\frac{5}{6})^{3/2}\Omega_0^3 < -3\Omega\Omega_0^2$, or

$$\frac{\Omega}{\Omega_0} < \sqrt{\frac{125}{486}} = 0.507. \quad (5.189)$$

A convenient measure of the importance of rotation for the equilibrium structure of any self-gravitating body is the ratio t of the rotational kinetic energy T to the body's self-gravitational potential energy W :

$$t \equiv \frac{T}{|W|}. \quad (5.190)$$

The virial theorem (4.248) states that $K/|W| = \frac{1}{2}$, where $K = T + \frac{1}{2}\Pi$ is the sum of the rotational kinetic energy T and the kinetic energy in random motions $\frac{1}{2}\Pi$. Thus we have $0 \leq t \leq \frac{1}{2}$. For a Kalnajs disk, $t = \frac{1}{2}(\Omega/\Omega_0)^2$. Thus the Kalnajs disks are unstable to a bar-like mode if

$$t > 0.1286. \quad (5.191)$$

It is natural to speculate that this instability is related to the bars found in many disk galaxies, but uniformly rotating stellar systems such as Kalnajs disks can behave quite differently from the differentially rotating disks found in galaxies, and we shall find in Chapter 6 that bar-like instabilities in differentially rotating disks arise from a quite different mechanism.

Figure 5.5 shows the regions of instability for modes with quantum number $l \leq 5$, as determined by Kalnajs (1972a) using the linearized collisionless

Boltzmann equation. As l increases, the behavior becomes increasingly complicated: for a given pair of indices (l, m) , a Kalnajs disk has up to $l + 1$ different modes with different frequencies ω . The figure shows that all Kalnajs disks are unstable to one or more modes; however, there are stable composite systems consisting of superimposed Kalnajs disks.

5.6.3 Maclaurin spheroids and disks

Throughout this chapter we have stressed the analogies between the dynamics of stellar and fluid systems. This analogy extends to the bar-like instability of rapidly rotating Kalnajs disks, which is also present in uniformly rotating, self-gravitating bodies of incompressible fluid.

The study of rotating fluid bodies was initiated by Newton, who treated the Earth as a homogeneous, incompressible fluid body in order to estimate the expected flattening of the earth due to its rotation. His treatment was valid only for small rotation speeds, but in 1742 C. Maclaurin found an exact solution for the equilibrium of an incompressible, uniformly rotating fluid body (Problem 5.19). These **Maclaurin spheroids** form a one-parameter family of oblate (flattened) spheroids (§2.5.2), which can be parametrized by their angular momentum or angular speed, by the ratio of polar to equatorial axes, or by the parameter t defined in equation (5.190).

A surprise came in 1834 when C. G. Jacobi showed that the Maclaurin spheroids were not the only possible equilibrium shape for a uniformly rotating, incompressible fluid body: there were also equilibria in which the fluid surface was ellipsoidal, with all three principal axes different. In other words, while a non-rotating fluid body supported by its own self-gravity *must* be spherical by Lichtenstein's theorem (Box 4.1), a rotating fluid body supported by its own self-gravity need *not* be axisymmetric. It is a measure of the unexpectedness of this result that Jacobi's discovery came almost a century after Maclaurin's. The relationship of the **Jacobi ellipsoids**, as these configurations are now called, to the simpler Maclaurin spheroids remained obscure until the end of the nineteenth century, when Poincaré and others showed that the Jacobi ellipsoids are actually the preferred configurations of rapidly rotating fluid bodies because they have lower energy for fixed angular momentum and mass. In the language of §5.4, rapidly rotating Maclaurin spheroids with $0.1375 < t < 0.2738$ are dynamically stable but secularly unstable: if the fluid has even the tiniest viscosity, they will gradually evolve into Jacobi ellipsoids.

There is a simple physical basis for this instability. The kinetic energy of a uniformly rotating body is $\frac{1}{2}L^2/I$, where L is the angular momentum and I is the moment of inertia (eq. D.43). The deformation of the spheroid into an ellipsoid increases the moment of inertia and thus lowers the rotational kinetic energy. The deformation also raises the gravitational potential energy; however, for sufficiently large angular momentum the loss of rotational

kinetic energy is greater than the gain of potential energy, so the spheroid is secularly unstable.

For the most rapidly rotating Maclaurin spheroids, $0.2738 < t < 0.5$, the secular instability is overwhelmed by a more rapid dynamical instability to a bar-like mode, reminiscent of the dynamical instability that plagues the Kalnajs disks with $0.1286 < t < 0.5$ (eq. 5.191).

The complete spectrum of modes of the Maclaurin spheroids was first analyzed by Bryan (1888). More recent discussions appear in Lyttleton (1953), Chandrasekhar (1969), and Tassoul (1978). An even closer fluid analog to the Kalnajs disk is the **Maclaurin disk**, a razor-thin disk of two-dimensional fluid, in which the pressure acts only in the disk plane and the surface-density distribution is the same as in the Kalnajs disks. The Maclaurin disks become secularly unstable to a bar-like mode at $t = \frac{1}{8} = 0.125$, and dynamically unstable at $t = \frac{1}{4} = 0.25$ (Problem 5.20; Takahara 1976; Weinberg 1983).

Finally, in 1860 Riemann showed that even the Jacobi ellipsoids are only special members of a much larger family of triaxial equilibrium configurations of self-gravitating fluid; the **Riemann ellipsoids** not only rotate but have internal streaming motions (Chandrasekhar 1969). The Riemann ellipsoids draw attention to the distinction between the **material speed** at which the *matter* of a rotating triaxial body streams and the pattern speed, the angular speed at which the *figure* of the body rotates. This distinction is important for bars in disk galaxies because the pattern speed governs the gravitational interaction of the bar with the rest of the galaxy, while the material speed determines the observed velocities.

Problems

5.1 [1] For over 150 years, most astronomers believed that Saturn's rings were rigid bodies, until Laplace showed that a solid ring would be unstable. The same instability plagues Larry Niven's famous fictional planet *Ringworld*. Following Laplace, consider a rigid, circular wire of radius R and mass m that lies in the x - y plane, centered on a planet of mass $M \gg m$ at the origin. The wire rotates around the planet in the x - y plane at the Keplerian angular speed $\Omega_K = (GM/R^3)^{1/3}$. Show that this configuration is linearly unstable and find the growth rate of the instability.

5.2 [2] (a) Prove that the temporal Fourier transform of the polarization function satisfies

$$\tilde{P}^*(\mathbf{x}, \mathbf{x}', \omega) = \tilde{P}(\mathbf{x}, \mathbf{x}', -\omega^*), \quad (5.192)$$

where the star denotes complex conjugation.

(b) If $P(\mathbf{x}, \mathbf{x}', \tau) \exp(-c\tau) \rightarrow 0$ as $\tau \rightarrow \infty$ for some real constant $c > 0$, prove that $\tilde{P}(\mathbf{x}, \mathbf{x}', \omega)$ is analytic throughout the half-plane $\text{Im}(\omega) > c$. Hint: use causality and the residue theorem.

(c) Do analogous results also hold for the response function $\tilde{R}(\mathbf{x}, \mathbf{x}', \omega)$?

5.3 [1] (a) Let $\tilde{P}(\mathbf{x}, \mathbf{x}', \omega)$ denote the temporal Fourier transform of the polarization function of a stable system. We assume that $\tilde{P}(\mathbf{x}, \mathbf{x}', \omega) \rightarrow 0$ as $|\omega| \rightarrow \infty$ in the upper half-plane. Prove that

$$\int_{-\infty}^{\infty} d\omega' \frac{\tilde{P}(\mathbf{x}, \mathbf{x}', \omega')}{\omega' - \omega + i\eta} = 0, \quad (5.193)$$

where ω and ω' are real and η is a positive number. Hint: use Problem 5.2.

(b) Prove that

$$\tilde{P}(\mathbf{x}, \mathbf{x}', \omega) = -\frac{i}{\pi} \wp \int_{-\infty}^{\infty} d\omega' \frac{\tilde{P}(\mathbf{x}, \mathbf{x}', \omega')}{\omega' - \omega}, \quad (5.194)$$

where \wp denotes the Cauchy principal value, defined in equation (C.6). Hint: use the Plemelj identity, equation (C.5).

(c) Prove that for real ω , the polarization function can always be written in the form

$$\tilde{P}(\mathbf{x}, \mathbf{x}', \omega) \equiv \tilde{P}_+(\mathbf{x}, \mathbf{x}', \omega) + \tilde{P}_-(\mathbf{x}, \mathbf{x}', \omega), \quad (5.195)$$

where

$$\tilde{P}_{\pm}^*(\mathbf{x}', \mathbf{x}, \omega) = \pm \tilde{P}_{\pm}(\mathbf{x}, \mathbf{x}', \omega). \quad (5.196)$$

(d) Prove the **Kramers–Kronig relations**,

$$\tilde{P}_{\pm}(\mathbf{x}, \mathbf{x}', \omega) = -\frac{i}{\pi} \wp \int_{-\infty}^{\infty} d\omega' \frac{\tilde{P}_{\mp}(\mathbf{x}, \mathbf{x}', \omega')}{\omega' - \omega}. \quad (5.197)$$

The same results apply to $\tilde{R}(\mathbf{x}, \mathbf{x}', \omega)$.

5.4 [1] At typical sea-level conditions ($p = 1.01 \times 10^5 \text{ N m}^{-2}$ and $T = 15^\circ \text{ C}$), the density of air is 1.22 kg m^{-3} and the speed of sound is 340 m s^{-1} . Find (i) the fractional change in frequency due to the self-gravity of the air, for a sound wave with wavelength 1 meter; (ii) the Jeans length.

5.5 [2] Suppose that the gravitational potential due to a body of mass m is modified from the Newtonian form to the Yukawa potential, $\Phi(r) = -Gm \exp(-\alpha r)/r$. How does this modification affect the Jeans wavenumber for a homogeneous fluid (eq. 5.35) or stellar system (eq. 5.49)? Is the Jeans swindle needed in this analysis (KieSSLing 2003)? Hint: see Problem 2.12.

5.6 [1] Modify the derivation in equations (5.32)–(5.39) to find the polarization function $\bar{P}(\mathbf{k}, \tau)$ rather than the response function.

5.7 [3] Appendix F.3.1 derives the energy density of sound waves. Here we generalize these results to find the energy density of sound waves with self-gravity.

(a) Consider a homogeneous fluid with constant density ρ_0 , and implement the Jeans swindle by assuming that the equilibrium gravitational field $-\nabla\Phi_{s0}$ associated with this density is canceled by a fixed external field $-\nabla\Phi_e$, i.e., $\Phi_{s0} = -\Phi_e$. Starting from equation (F.26), show that the self-gravity of the fluid contributes an energy density $E_{g0} + \Delta E_g$, where E_{g0} is the gravitational energy density in the homogeneous fluid,

$$\Delta E_g = \frac{1}{2} \rho_0 \Delta\Phi_s - \frac{1}{2} \Delta\rho \Phi_{s0} + \frac{1}{2} \Delta\rho \Delta\Phi_s, \quad (5.198)$$

and $\Delta\rho \equiv \rho - \rho_0$, $\Delta\Phi_s \equiv \Phi_s - \Phi_{s0}$.

(b) Prove that

$$E'_g \equiv \frac{1}{2} \rho_0 \Delta\Phi_s - \frac{1}{2} \Phi_{s0} \Delta\rho \quad ; \quad F'_g \equiv \frac{1}{8\pi G} \left(\Phi_{s0} \frac{\partial^2 \Delta\Phi_s}{\partial x \partial t} - \frac{\partial \Phi_{s0}}{\partial x} \frac{\partial \Delta\Phi_s}{\partial t} \right) \quad (5.199)$$

are the non-wave contributions to the energy density ΔE_g , by showing that they satisfy the transport equation $\partial E'_g / \partial t + \partial F'_g / \partial x = 0$.

(c) The gravitational wave energy is $\Delta E_g - E'_g$, and the total wave energy density is $E_{w,g} \equiv E_w + \Delta E_g - E'_g$, where E_w is given by equation (F.62). For wavetrains of the form (F.64), prove that the average wave energy density is

$$\langle E_{w,g} \rangle = \frac{A_\rho^2}{4k^2 \rho_0} (\omega^2 + v_s^2 k^2 - 4\pi G \rho_0) + O(A^3) = \frac{A_\rho^2 \omega^2}{2k^2 \rho_0} + O(A^3), \quad (5.200)$$

where the last equality makes use of the dispersion relation (5.34). This result holds only for wavenumbers $k > k_J$, since disturbances with smaller wavenumber are unstable rather than oscillatory. All oscillatory wavetrains have positive energy density, but the energy density decreases to zero at the stability boundary $k^2 = 4\pi G \rho_0 / v_s^2 = k_J^2$.

(d) Waves with $k < k_J$ have the form $\rho_1 = A_\rho \exp(\lambda t) \cos(kx)$, $v_{x1} = A_v \exp(\lambda t) \sin(kx)$. Prove that

$$A_v = -\frac{\lambda}{\rho_0 k} A_\rho. \quad (5.201)$$

(e) Prove that all unstable disturbances have $\langle E_w \rangle = 0$, as required by energy conservation since the amplitude of the wave is growing as $\exp(\lambda t)$.

5.8 [2] An infinite homogeneous stellar system has density ρ_0 and DF

$$f_0(\mathbf{v}) = \frac{\rho_0 \theta}{\pi^2} \frac{1}{(v^2 + \theta^2)^2}, \quad (5.202)$$

where θ is a measure of the characteristic velocity (note that θ^2 is *not* the mean-square velocity, which diverges for this DF). Using the Jeans swindle, show that the polarization function is (Summers & Thorne 1991)

$$\bar{P}(\mathbf{k}, \tau) = 4\pi G \rho_0 H(\tau) \tau e^{-\theta k \tau}, \quad (5.203)$$

where $H(\tau)$ is the step function (Appendix C.1); that the Fourier transform of the polarization function is

$$\tilde{\bar{P}}(\mathbf{k}, \omega) = -\frac{4\pi G \rho_0}{(\omega + ik\theta)^2}, \quad \text{Im}(\omega) > 0; \quad (5.204)$$

and that the response function is

$$\bar{R}(\mathbf{k}, \tau) = \omega_0 H(\tau) e^{-\theta k \tau} \sinh(\omega_0 \tau), \quad (5.205)$$

where $\omega_0 = (4\pi G \rho_0)^{1/2}$. Hence show that the system is unstable if

$$k < k_J \equiv \left(\frac{4\pi G \rho_0}{\theta^2} \right)^{1/2}. \quad (5.206)$$

5.9 [2] The Jeans instability can be analyzed exactly in rotating systems, without invoking the Jeans swindle. Consider a homogeneous, self-gravitating, barotropic fluid of density ρ_0 , contained in an infinite cylinder of radius R_0 whose symmetry axis is the z axis. The cylinder walls and the fluid rotate at angular speed $\boldsymbol{\Omega} = \Omega \hat{\mathbf{e}}_z$.

(a) Show that the gravitational field inside the cylinder is

$$-\nabla \Phi_0 = -2\pi G \rho_0 (x \hat{\mathbf{e}}_x + y \hat{\mathbf{e}}_y). \quad (5.207)$$

(b) Using Euler's equation in a rotating frame (cf. eq. 5.152b), find the condition on Ω so that the fluid is in equilibrium.

(c) Now let $R_0 \rightarrow \infty$, or, what is equivalent, consider wavelengths $\lambda \ll R_0$, so the boundary condition due to the wall can be neglected. Working in the rotating frame, find the dispersion relation analogous to equation (5.34) for (i) waves propagating perpendicular to $\boldsymbol{\Omega}$; (ii) waves propagating parallel to $\boldsymbol{\Omega}$. Show that waves propagating perpendicular to $\boldsymbol{\Omega}$ are always stable, while waves propagating parallel to $\boldsymbol{\Omega}$ are stable if and only if the usual Jeans criterion $k < k_J$ is satisfied (Chandrasekhar 1961). See Lynden-Bell (1962b) for the analogous stellar system.

5.10 [3] In this problem we examine the linear modes of a self-gravitating sphere of incompressible liquid of density ρ and equilibrium radius R . The radius of the surface of the perturbed sphere can be written as $r(\theta, \phi, t) = R + \epsilon Y_l^m(\theta, \phi) \exp(-i\omega_l t)$, where $\epsilon \ll R$, $l = 1, 2, \dots$, and $|m| \leq l$. Show that the frequencies of these **Kelvin modes** satisfy (Chandrasekhar 1961)

$$\omega_l^2 = \frac{8\pi G \rho}{3} \frac{l(l-1)}{2l+1}. \quad (5.208)$$

Hint: the gravitational potential arising from the small distortion ϵ of the surface of the sphere is the same as the potential arising from a surface density $\epsilon \rho$.

5.11 [2] Let \mathcal{V} be the action space for a given Hamiltonian $H(\mathbf{J})$. For example, the actions (J_1, J_2, J_3) defined in Table 3.1 have an action space $-J_2 \leq J_1 \leq J_2, 0 \leq J_2 \leq \infty, 0 \leq J_3 \leq \infty$. Let \mathcal{S} be the surface bounding the action space, and let $\hat{\mathbf{n}}(\mathbf{J})$ be the normal to this surface at a given point. We now perturb the Hamiltonian by adding a potential $\Phi(\mathbf{x})$, which can be expanded in a Fourier series in the angles, $\Phi(\mathbf{x}) = \sum_{\mathbf{m}} \Phi_{\mathbf{m}}(\mathbf{J}) \exp(i\mathbf{m} \cdot \boldsymbol{\theta})$ (eq. 5.79). Prove that

$$\mathbf{m} \cdot \hat{\mathbf{n}}(\mathbf{J}) \Phi_{\mathbf{m}}(\mathbf{J}) = 0 \quad \text{everywhere on } \mathcal{S}. \quad (5.209)$$

This result is used in deriving equation (5.95). Hint: argue that particles cannot cross the boundaries of the action space.

5.12 [1] This problem analyzes a simple example of secular stability. We consider a hemispherical bowl of radius R , which rotates about a vertical axis with angular speed Ω . A particle slides inside the bowl. The particle is subject to a frictional force $-k(\mathbf{v} - \mathbf{v}_b)$, where \mathbf{v} is the velocity of the particle and \mathbf{v}_b is the velocity of the bowl. The coefficient of friction k may be assumed to be very small. The particle is initially at rest, at the bottom of the bowl, and then is given a small displacement.

(a) Prove that the particle returns to rest at the bottom of the bowl (secular stability) if and only if $\Omega < (g/R)^{1/2}$, where g is the acceleration due to gravity.

(b) If the motion is secularly unstable, what is the final fate of the particle?

5.13 [1] Let $A(\mathbf{x}, \mathbf{v}), B(\mathbf{x}, \mathbf{v})$, and $C(\mathbf{x}, \mathbf{v})$ be arbitrary functions that vanish as $|\mathbf{x}|, |\mathbf{v}| \rightarrow \infty$. Prove the following identities involving Poisson brackets:

$$\int d^3\mathbf{x} d^3\mathbf{v} [A, B] = 0; \quad (5.210a)$$

$$\int d^3\mathbf{x} d^3\mathbf{v} A[B, C] = \int d^3\mathbf{x} d^3\mathbf{v} C[A, B] = \int d^3\mathbf{x} d^3\mathbf{v} B[C, A]; \quad (5.210b)$$

$$[A, f(B)] = f'(B)[A, B] = [f'(B)A, B] \quad \text{for any function } f(B). \quad (5.210c)$$

5.14 [2] Prove that $W[\rho_1]$ in Chandrasekhar's variational principle (5.120) is zero for a uniform displacement, $\rho_1(\mathbf{x}) = -\nabla \cdot (\rho_0 \boldsymbol{\xi})$ with $\boldsymbol{\xi}$ constant.

5.15 [2] In this problem we describe Antonov's original proof of his variational principle for the stability of stellar systems with an ergodic DF having $f'_0(H_0) < 0$.

(a) Show that any DF $f_1(\mathbf{x}, \mathbf{v}, t)$ can be written in the form

$$f_1(\mathbf{x}, \mathbf{v}, t) = f_+(\mathbf{x}, \mathbf{v}, t) + f_-(\mathbf{x}, \mathbf{v}, t), \quad (5.211)$$

where f_+ is an even function of \mathbf{v} and f_- is an odd function of \mathbf{v} .

(b) If the external potential $\Phi_e = 0$, show that the linearized collisionless Boltzmann equation (5.123) can be written in the form

$$\frac{\partial f_+}{\partial t} + [f_-, H_0] = 0 \quad ; \quad \frac{\partial f_-}{\partial t} + [f_+, H_0] + f'_0(H_0)[H_0, \Phi_1] = 0, \quad (5.212)$$

where as usual $H_0(\mathbf{x}, \mathbf{v}) = \frac{1}{2}\mathbf{v}^2 + \Phi_0(\mathbf{x})$.

(c) Show that equations (5.212) can be combined into a single equation for f_- (Antonov 1960; Laval, Mercier, & Pellat 1965; Kulsrud & Mark 1970)

$$\frac{\partial^2 f_-}{\partial t^2} - [[f_-, H_0], H_0] + G f'_0(H_0) \left[H_0, \int \frac{d^3\mathbf{x}' d^3\mathbf{v}'}{|\mathbf{x} - \mathbf{x}'|} [f_-, H_0]_{\mathbf{x}', \mathbf{v}'} \right] = 0. \quad (5.213)$$

(d) Take the temporal Fourier transform (5.6) of equation (5.213), multiply by $\tilde{f}_-^*/f'_0(H_0)$, integrate over $d^3\mathbf{x}d^3\mathbf{v}$, and use equations (5.210) to obtain

$$\begin{aligned} \omega^2 \int \frac{d^3\mathbf{x}d^3\mathbf{v}}{|f'_0(H_0)|} |\tilde{f}_-|^2 &= \int \frac{d^3\mathbf{x}d^3\mathbf{v}}{|f'_0(H_0)|} |\tilde{f}_-, H_0|^2 \\ &- G \int \frac{d^3\mathbf{x}d^3\mathbf{v}d^3\mathbf{x}'d^3\mathbf{v}'}{|\mathbf{x} - \mathbf{x}'|} [\tilde{f}_-, H_0]_{\mathbf{x},\mathbf{v}}^* [\tilde{f}_-, H_0]_{\mathbf{x}',\mathbf{v}'}. \end{aligned} \quad (5.214)$$

(e) Prove that all modes of stellar systems with an ergodic DF have real ω^2 ; stable modes have $\omega^2 \geq 0$ and unstable modes have $\omega^2 < 0$.

(f) Using the substitution $\tilde{f}_- = f'_0(H_0)g$, prove that Antonov's variational principle $W_A[g] \geq 0$ (eq. 5.136) is sufficient for stability.

5.16 [3] Let $f_1(\mathbf{x}, \mathbf{v}, t)$ be the linear response of a spherical stellar system with an ergodic DF $f_0(H_0)$ to a weak gravitational potential $\Phi_1(\mathbf{x}, t)$, which includes both an external potential and the potential arising from the self-gravity of the response. Assume that $\Phi_1 \rightarrow 0$ as $t \rightarrow -\infty$.

(a) If Φ_1 varies slowly (i.e., on a timescale much greater than the orbital time) prove that

$$f_1(\mathbf{x}, \mathbf{v}, t) = f'_0(H_0) [\Phi_1(\mathbf{x}, t) - \langle \Phi_1 \rangle_{\mathbf{x},\mathbf{v}}], \quad (5.215)$$

where $\langle X \rangle_{\mathbf{x},\mathbf{v}}$ is the time average of X over the unperturbed orbit that passes through the phase-space point (\mathbf{x}, \mathbf{v}) (Lynden-Bell 1969). Hint: use equation (5.83).

(b) Show that the density response is

$$\rho_1(\mathbf{x}, t) = \int d^3\mathbf{v} f_1(\mathbf{x}, \mathbf{v}, t) = \left(\frac{d\rho}{d\Phi} \right)_0 \Phi_1(\mathbf{x}, t) - \int d^3\mathbf{v} f'_0(H_0) \langle \Phi_1 \rangle_{\mathbf{x},\mathbf{v}}. \quad (5.216)$$

(c) Show that the linear response of a static, barotropic fluid to a similar perturbation is given by the first of the two terms in the final expression in equation (5.216).

5.17 [1] A semicircular trough of radius a contains water of density ρ . The maximum depth of the water is $h \ll a$. The trough is gently moved back and forth until the water is excited into a sloshing or seiche mode, in which the surface of the water stays flat but its angle ϑ from the horizontal oscillates. What is the oscillation period for $|\vartheta| \ll 1$?

5.18 [2] A simple way to stabilize the Kalnajs disks is to imagine that they are embedded in a fixed axisymmetric gravitational field (say, due to the halo of the galaxy). Suppose that the fixed potential is $\Phi_h(R) = \frac{1}{2}h\Omega_0^2 R^2$, where the constant h is chosen so that the disk provides a fraction f_d of the total radial force in the equilibrium system. Show that a Kalnajs disk embedded in such a halo is stable to $l = m = 2$ modes if and only if

$$\frac{\Omega^2}{\Omega_h^2} < \frac{(8 - 3f_d)^3}{486f_d^2}, \quad (5.217)$$

where Ω_h is the angular speed of a cold disk. For what values of f_d are *all* of the Kalnajs disks stable?

5.19 [2] The gravitational potential in a Maclaurin spheroid is given by equation (2.128) and Table 2.1. Prove that $t = T/|W|$ (eq. 5.190) is related to the eccentricity e of the surface of the spheroid by

$$t = \frac{3}{2e^2} - 1 - \frac{3\sqrt{1-e^2}}{2e \sin^{-1} e}. \quad (5.218)$$

5.20 [3] The Maclaurin disks (§5.6.3) are razor-thin fluid disks. The fluid pressure p acts only in the disk plane and has a polytropic equation of state, $p = K\Sigma^3$, where Σ is the surface density. Maclaurin disks have the same surface-density distribution as Kalnajs disks, $\Sigma(R) = \Sigma_c(1 - R^2/a^2)^{1/2}$ (eq. 4.167), and rotate at uniform angular speed Ω .

(a) Prove that

$$\Omega^2 = \Omega_0^2 - \frac{3K\Sigma_c^2}{a^2}, \quad (5.219)$$

where Ω_0 is given by equation (4.166).

(b) Prove that the Maclaurin disks have modes with the same surface-density distribution as the modes of the Kalnajs disks (eq. 5.176), with the frequency ω of the mode having quantum numbers (l, m) given by (Takahara 1976)

$$\omega_r^3 - \omega_r \{4\Omega^2 + (l^2 + l - m^2)[\Omega_0^2(1 - g_{lm}) - \Omega^2]\} + 2m\Omega[\Omega_0^2(1 - g_{lm}) - \Omega^2] = 0, \quad (5.220)$$

where $\omega_r \equiv \omega - m\Omega$ and g_{lm} is defined by equation (2.204b).

(c) Prove that (i) the modes with $l = 1$ correspond to a uniform translation of the disk; (ii) the modes with $l = 2, m = \pm 2$ (the bar mode) are unstable if and only if $\Omega/\Omega_0 > 2^{-1/2}$.
Hint: for $l = 2, m = \pm 2$ one root is $\omega = 2m\Omega$.

6

Disk Dynamics and Spiral Structure

Galaxies contain disks for the same basic reason that planetary systems, planetary rings, accretion disks, and many other astrophysical systems are flat: gas can radiate energy but not angular momentum, and for a given distribution of angular momentum along an axis, the state of lowest energy is a flat disk perpendicular to that axis.

Galactic disks exhibit more interesting and complex behavior than hot, slowly rotating stellar systems such as elliptical galaxies as a result of the interplay of energy and angular momentum. Loosely speaking, systems such as galaxies tend to evolve to states of lower energy. The virial theorem tells us that a self-gravitating system can lose energy at constant mass M by contracting to a smaller radius. Uniform contraction of a rotating disk is not permissible because it would not conserve the angular momentum. However, the disk can release energy by a more subtle strategy. Suppose that a small element of mass m initially at radius R_i is moved to a circular orbit at a radius R_f that is much larger than the typical radius of material in the disk. The gravitational field of the disk is approximately Keplerian at large radii, so the angular momentum of the mass element in its new orbit is $m(GMR_f)^{1/2}$, which diverges as $R_f \rightarrow \infty$. On the other hand, the energy of the mass element in its new orbit is $-\frac{1}{2}GMm/R_f$, which approaches a constant value of zero as $R_f \rightarrow \infty$. Thus a small mass element can swallow most or all of the angular momentum of the disk with a negligible energetic penalty. Having thus been relieved of its angular momentum, the rest of the

disk can release free energy by shrinking.

A more precise statement is that it is energetically favorable for differentially rotating disks to transfer angular momentum outward and mass inward (Lynden-Bell & Pringle 1974). In fluid disks this transfer can be effected by viscous or magnetic stresses and, indeed, angular-momentum transfer through these stresses is believed to drive the accretion disks that power active galactic nuclei, cataclysmic variable stars, and many stellar X-ray sources. Disks that are collisionless stellar systems have neither viscosity nor magnetic fields, but can transfer angular momentum through gravitational torques that are generated by non-axisymmetric features such as spiral arms. It is the ability of such features to liberate energy stored in the disk that is the root cause of their formation in disk galaxies.

The disk of a spiral galaxy exhibits far richer behavior than a purely stellar disk because it contains both stars and gas. These two components interact in complex ways, not only through their mutual gravity but also because gas can be converted to stars—in particular, we shall find that star formation is strongly enhanced by the non-axisymmetric features that we are trying to study. A complete exposition of the dynamics of galactic disks would therefore require that we first understand stellar dynamics, gas dynamics, stellar evolution, and star formation; despite much progress, astronomers are still groping towards this goal.

The first part of this chapter is devoted to spiral structure, the most dramatic and beautiful aspect of disk dynamics in galaxies. After describing the phenomenology of spiral structure in §6.1, we investigate the propagation and evolution of density waves in §6.2. The analysis in this section is based on the fundamental approximation that the waves are tightly wound spirals. In §6.3 we relax this assumption, and show that numerical calculations of waves in stellar disks can be understood by thinking of the disk as a cavity within which resonant wave patterns are established. The growth and decay of waves in this cavity is determined by the competition between a built-in amplifier and absorber, in the same way that audio feedback is established between a microphone, amplifier, and loudspeaker. In §6.4 we apply this understanding of disk dynamics to spiral structure, where we must deal with the additional complications caused by interstellar gas and the birth of new stars from it. In §6.5 we discuss the structure and dynamics of the prominent bar-like structures that are found in the centers of many disk galaxies, including our own. Finally, in §6.6 we discuss the collective motions of stars in the direction perpendicular to the disk plane, which give rise to warps and other structures that dominate the appearance of the outermost parts of many disk galaxies.

6.1 Fundamentals of spiral structure

The majestic sweep of spiral arms across the face of a galaxy like M51 or M100 (Plates 1 and 17) is one of the most inspiring sights in the sky. Features such as these are more than just decorative frosting on a galactic disk. As the primary sites of star formation, spiral arms mold the properties of the stellar disk, as well as the chemical composition, dynamics, and thermal balance of the interstellar gas. Their strength and shape offer probes of the dynamics of the gas and stars in the disk. Spiral arms drive the long-term dynamical evolution of the galactic disk, through such processes as gravitational scattering of stars (see §8.4.2) and angular-momentum transport (see §§6.1.5, 6.2.6, and Appendix J). Spiral or bar-like structure in the inner disk may also feed massive black holes in galaxy centers with gas, thereby contributing to the luminosities of some active galactic nuclei.

Understanding the origin and evolution of spiral structure has proved to be one of the harder problems in astrophysics. The first major attack was made by the Swedish astronomer Bertil Lindblad, who struggled with this problem until his death in 1965. Lindblad correctly recognized that spiral structure arises from the interaction between the orbits and gravitational forces of the stars of the disk. In this view he stood almost alone; at the time of his death most astronomers believed that spiral structure is caused by the interstellar magnetic field, which we now know is not responsible for large-scale spiral structure (see Box 6.1). However, Lindblad's methods have been superseded by more powerful analytical and numerical tools.

At about the time of Lindblad's death, two major insights ignited the interest of the astronomical community in spiral structure, and established much of the needed theoretical framework for the study of both spiral structure and the stability of galactic disks (see Pasha 2002, 2004 for a historical account).

The first of these insights came from C. C. Lin and Frank Shu at MIT. They recognized that spiral structure could be viewed as a density wave, a periodic compression and rarefaction of the disk surface density that propagates through the disk in much the same way that waves propagate over the ocean surface. Lin and Shu also showed that many of the tools of wave mechanics could be used to study the properties of density waves in differentially rotating stellar disks. They combined these insights with a bold hypothesis due to Lindblad: that the spiral patterns in galaxies—or at least some galaxies—are long-lasting, in other words, that the appearance of the pattern remains stationary (unchanged except for an overall rotation) over many orbital periods. We shall call this the **stationary spiral structure hypothesis**. These concepts led to the **Lin–Shu hypothesis**, that *spiral structure is a stationary density wave*. In more mathematical terms, the Lin–Shu hypothesis is that spiral structure is a neutrally stable mode of the galactic disk, analogous to the modes of spherical stellar systems that we examined in Chapter 5.

Box 6.1: Magnetism and spiral structure

In the 1950s, most astronomers suspected that spiral structure was the result of interactions between the interstellar gas and magnetic field. This argument is no longer viable for the largest and most prominent spirals, because spiral structure is found to be present in the old disk stars, which are unaffected by interstellar magnetic fields (§6.1.2). However, the magnetic field could play some role in shaping small-scale spiral structure in gas and young stars. The following simple energy argument allows us to assess this possibility. The energy density due to a magnetic field B is $\frac{1}{2}B^2/\mu_0$ (Jackson 1999). The kinetic-energy density in a patch of the disk with an internal spread in velocity $\Delta\mathbf{v}$ and gas density ρ is $\frac{1}{2}\rho(\Delta\mathbf{v})^2$. If the patch has size ΔR then the rotation of the galaxy imparts a velocity spread of order $\Omega\Delta R$, where Ω is the circular angular speed; since the interstellar clouds have random motions as well, it is safe to assume that $\Delta v \gtrsim \Omega\Delta R$. Setting $\rho = 0.05 \mathcal{M}_\odot \text{pc}^{-3}$ (the local gas density according to Table 1.1) and equating the magnetic- and kinetic-energy density yields $\Delta R \lesssim 0.5 \text{ kpc} (B/\text{nT})$. Typical magnetic fields in galaxies are $\sim 0.5 \text{ nT}$ (Beck et al. 1996). Thus the interstellar magnetic field is not strong enough to play a role in large-scale spiral structure, although it may facilitate the formation of smaller features such as spurs (Kim & Ostriker 2002), cloud complexes (Parker 1966; Basu, Mouschovias, & Paleologou 1997) and nuclear rings (Beck et al. 1999).

The Lin–Shu hypothesis enabled theorists, for the first time, to make a wide variety of quantitative predictions for comparison with observations of spiral galaxies, and offered them the heady vision of computing the shapes and other properties of spiral galaxies from first principles. Unfortunately, we shall see that the Lin–Shu hypothesis is not correct, at least for most galaxies: large-scale spiral structure *is* a density wave—the most convincing evidence comes from near-infrared images of nearby galaxies, as we discuss in §6.1.2—but the spiral pattern is far from stationary.

The second major discovery was that galactic disks are remarkably responsive to small disturbances. This insight came from two quite different research projects. At Cambridge University, Peter Goldreich and Donald Lynden–Bell investigated the evolution of small-scale disturbances in a simplified model of a differentially rotating, self-gravitating fluid disk. They found the surprising result that generic initial disturbances in the disk were amplified—grew for a limited time while being sheared by the differential rotation—by a factor of ten or more, even when the disk itself was safely stable. Simultaneously, William Julian and Alar Toomre at MIT asked how a differentially rotating stellar disk would respond to the presence of a point mass traveling on a circular orbit within the disk. They found that the grav-

itational field of the mass induced a unexpectedly strong spiral-shaped wake or wave in the stars of the disk. Remarkably, the total mass enhancement associated with the wake could exceed the mass of the perturber by an order of magnitude or more. The results emerging from Cambridge and MIT reinforced one another, by showing that differentially rotating, self-gravitating disks responded vigorously both to temporary disturbances and to steady forcing, and by showing that this phenomenon was present in both fluid and stellar disks. The obvious inference was that strong spiral patterns are likely to result from a wide variety of causes—clumps of matter in the interstellar gas, tidal forces from companion or satellite galaxies, or (in hindsight, since dark halos were not then known) substructure in the dark-matter halo, etc.—and hence that spiral structure in galaxies may mostly be transitory rather than stationary.

The study of the dynamics of differentially rotating disks has been developed into an extensive formalism called **density-wave theory**, which is central to understanding all kinds of astrophysical disks. Density-wave theory is also the main tool for studying the gravitational stability of disk galaxies and other astrophysical disks.

The number of reviews of spiral-structure theory is disappointingly small for such an important subject (but see Marochnik & Suchkov 1996). Remarkably, the article by Toomre (1977a) is still worth careful reading, even after several decades.

6.1.1 Images of spiral galaxies

Spiral structure is closely related to the large-scale properties of galaxies. This relation is reflected in the Hubble classification of spirals, which was briefly described in §1.1.3 (see Sandage & Bedke 1994 or BM §4.1.1 for more details). The Hubble classification tells us that the properties of the spiral arms (how tightly they are wound, how smooth they are, etc.) are correlated with properties such as the luminosity of the bulge relative to the disk, and the relative masses in interstellar gas and stars.

Despite these general correlations, there is a great deal of variety in the spiral structure exhibited by galaxies, even those of the same Hubble type. Hence the best introduction to spiral structure is to study images of nearby spiral galaxies (Malin 1993; Sandage & Bedke 1994; see also the web-based galleries maintained by observatories such as the Anglo–Australian Observatory, the National Optical Astronomy Observatory, and the Hubble Space Telescope). As an introduction, consider some of the galaxies shown in the color insert:

(i) M100=NGC 4321 (Plate 17): This galaxy contains two major spiral arms, each of which can be traced for almost a full revolution around the galaxy's center. Since the arms are fairly open and there is only a small central bulge, the galaxy is classified as Sbc. Spirals exhibiting this high

degree of symmetry and large-scale coherence are rare, although they are over-represented in textbooks and image galleries because of their striking appearance. Galaxies such as these, having long, continuous, symmetric arms, are called **grand-design** spirals; presumably, they have been formed by some large-scale *global* process that involves the whole galaxy. Grand-design spirals almost always have two main arms; thus, the appearance of the galaxy remains approximately the same if the image is rotated 180° . The thin dark stripes are **dust lanes** caused by absorption of the galaxy's starlight in dense clouds of gas and dust (see BM §§4.1.1, 4.4.7, 8.3.2); notice that these generally follow the spiral arms. Other examples of grand-design spirals are M51 (Plate 1) and M81 (Plate 8). In looking at grand-design spiral structure, it is worth remembering that the human brain sometimes finds long-range order in images even when none exists; this is presumably a consequence of many millenia spent looking for tigers in the jungle.

(ii) M101=NGC 5457 (Plate 18): Here the spiral arms, though prominent and well-defined, are less regular than in M100. Individual arms can be followed only for about half a rotation, and the two arms seen in M100 have been replaced by multiple arms which give a visual impression of branching or bifurcating. Note the prominent dark dust lanes and the bright blue knots of star formation that follow the spiral arms. This **intermediate-scale** spiral structure is coherent over scales that are a significant fraction of the galaxy size, but not over the whole galaxy. In contrast to grand-design spirals, intermediate-scale spirals do not give the impression of being long-lived features. To quote Oort (1970), "Looking at the irregularities in the actual spiral galaxies one wonders whether the present spirals could *continue* to exist for such a large number of revolutions. The problem seems particularly acute for the outer parts of Sc spirals, like M101..." The spiral structure of the Milky Way is probably of this type.

(iii) M33=NGC 598 (Plate 19): This is one of the nearest spiral galaxies. It has little or no central bulge, and its arms are broken up into stars and HII regions (see §6.1.2e); hence it is classified Scd. The arms are less regular than in M101, though it is still classified as an intermediate-scale spiral.

(iv) M63=NGC 5055 (Plate 9): The appearance of this galaxy is quite different from M100 or M101. Each spiral arm can be followed over only a small angle and the overall spiral pattern is composed of many patchy arm segments—a "swirling hotch-potch of pieces of spiral arms," in the words of Goldreich & Lynden-Bell (1965a). Such galaxies are often called **flocculent** spirals, in contrast to grand-design or intermediate-scale spirals (Elmegreen 1981). In galaxies like this one, there is probably little or no causal connection between the arms on opposite sides of the galaxy, and a *local*, rather than global, origin seems likely. The distinction between flocculent and grand-design structure is independent of the Hubble type: like M51, M63 is classified as Sbc. Elmegreen & Elmegreen (1987) describe the arm classification of a large sample of galaxies.

(v) NGC 1300 (Plate 10): Here is one of the most dramatic barred spirals. It is classified SBb (the letter B stands for “bar”; thus NGC 1300 is a barred Sb galaxy). The two spiral arms are very symmetrical: they can easily be followed through 180° , and on deep images through almost a full circle; thus the galaxy is classified as a grand-design spiral. Note that (a) there are sharp, straight dust lanes that extend from the sides of the central nucleus to the end of the bar; (b) the spiral arms start at the tips of the bar; (c) at the start of each spiral arm there is a cluster of HII regions (which look like bright stars on this image), indicating rapid star formation. These are all common features of barred spiral galaxies.

(vi) NGC 6745 (Plate 12): This galaxy has a peculiar lopsided, warped appearance and an extended trail of young, blue stars at the lower right. These features are probably the result of a recent encounter with the smaller galaxy just visible at the bottom right corner of the image. The gravitational field from the small galaxy has compressed and shocked the interstellar gas in NGC 6745, leading to a burst of star formation, and has dragged out the trail of stars pointing back towards it. This image is included as a reminder that not all galaxies have the relatively symmetrical structure seen in the other examples here.

6.1.2 Spiral arms at other wavelengths

Images like those in Plates 17 and 10 are dominated by the light from luminous, young stars and HII regions, which delineate narrow, sharply defined spiral arms. In the apt phrase of Baade (1963), the HII regions are “strung out like pearls along the arms”. Since these stars live less than ~ 10 Myr, compared to a typical orbital period of 100 Myr in the disk, they cannot travel far from their birthplace. We conclude that *the star-formation rate in spiral arms is much higher than in the rest of the disk*.

Since the appearance of these images is so strongly affected by the young stellar population, it is important to examine spiral galaxies at other wavelengths. Images in near-infrared light are particularly instructive, for two reasons. First, absorption by dust is much less severe at longer (redder) wavelengths, so the distribution of stars is unobscured—the optical depth due to dust is a factor of 10 smaller at $2.2 \mu\text{m}$ (K band) than at $550 \text{ nm} = 0.55 \mu\text{m}$ (V band) (see Tables 2.1 and 3.21 of BM). Second, near-infrared emission ($\sim 1\text{--}5 \mu\text{m}$) is mostly due to giant stars that are continuously produced from the old main-sequence dwarf stars that dominate the mass. The contribution of young stars to the near-infrared flux is small, except in localized regions of rapid star formation (Rix & Rieke 1993; Rhoads 1998). Thus, *near-infrared light traces mass, while blue light traces star formation*; both tracers are valuable but the mass is more fundamental for the dynamics of spiral structure.

Figure 6.1 shows a near-infrared image of M51, for comparison with Plate 1. The spiral pattern in the infrared image is similar to that traced

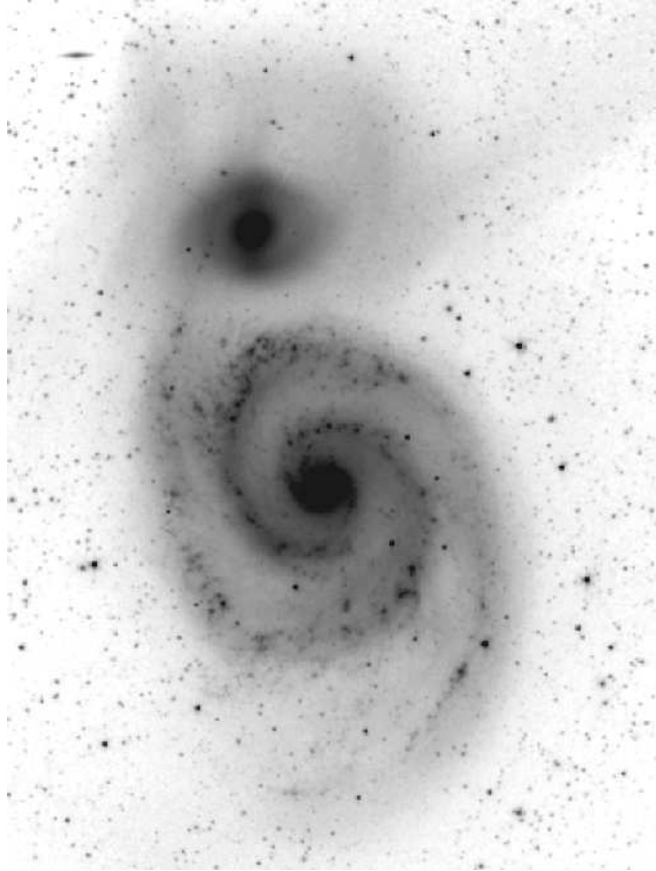


Figure 6.1 An image of M51 in the near-infrared ($3.6\ \mu\text{m}$), taken by the Spitzer Space Telescope (Calzetti et al. 2005). At this wavelength the dominant emission is from old stars, except for isolated point-like sources which are foreground stars or knots of rapid star formation. Note the faint, diffuse tidal streamers extending from the companion galaxy NGC 5195 (Toomre 1978). Credit: R. C. Kennicutt and the Spitzer Nearby Galaxies Survey.

by the young, blue stars, but the arms are smoother and broader—probably because the irregular spatial distribution of star formation sites has been phase-mixed away in the old stars. This difference between blue and red arms was discovered in a prescient paper by Zwicky (1955), and was confirmed and extended by Schweizer (1976). The presence of spiral structure in the old stars that dominate the mass is common to most grand-design spiral galaxies (Eskridge et al. 2002), and implies that *the entire stellar disk participates in the spiral pattern*; that is, there is a spiral pattern in *both* the surface density of the stellar disk (seen in near-infrared light) and in the star-formation

rate (seen in visible light). Typically, the arm traced by the young stars is displaced slightly inside the old-star arm.

Similar broad near-infrared spiral arms are seen in other grand-design spirals such as M100 and M81, and in galaxies with intermediate-scale arms such as M101. The presence of spirals in the old stars and hence in the disk mass is the strongest single piece of evidence that spiral structure in M51 and many other galaxies is a density wave, as envisaged by Lindblad, Lin, and Shu decades ago.

What is the physical connection between the spiral density wave and the spiral pattern in the young stars? Most likely, the gravitational field from the density wave deflects and squeezes the streamlines of the interstellar gas, thereby enhancing the gas density in the arms (see §6.4.1). In some cases the gas may even pass through a shock.¹ The star-formation rate is enhanced in these high-density regions, partly because there is simply more gas mass per unit volume, and partly because dense gas turns into stars faster (i.e., the star-formation rate per unit mass is higher). The spiral pattern in the young stars traces these regions of high star-formation rate.

In flocculent spirals, on the other hand, there is little or no spiral structure in the old, red disk stars (Elmegreen & Elmegreen 1984; Thornley 1996). This result suggests that the short spiral segments in flocculent galaxies may simply be local patches of star formation that have been sheared out into a spiral form by differential rotation, without significant spiral structure in the mass distribution.

It is helpful to have names for the spiral arms outlined by different tracers. The center of the **mass arm** is delineated by the maximum surface density at a given radius, which approximately coincides with the maximum near-infrared surface brightness (although gas can also contribute to the density). The **potential arm** is marked by the minimum of the gravitational potential, which is determined from the shape of the mass arm by solving Poisson's equation (see Problem 6.1). The center of the **gas arm** is the maximum of the gas density (see §6.4.1); and the **bright-star arm** is marked by the maximum of the luminosity density due to young stars.

Spiral structure can be seen in almost every component of the disk:

(a) Dust Several of the galaxies in the plates show thin, dark dust lanes, which curve along the spiral arms and partially mask the light from the luminous young stars that trace the arms. Barred galaxies also show straight dust lanes along the bar (e.g., NGC 1300 in Plate 10). Many of the properties of the dust lanes can be explained simply by assuming that the dust-to-gas ratio is constant, so the dust lane coincides with the gas arm. However, this model is almost certainly oversimplified, since the mechanisms that create

¹The concept of a shock from fluid mechanics (Landau & Lifshitz 2000) is only approximately valid in this context since the interstellar gas is concentrated in clouds that occupy only a small fraction of the disk volume; “shock” should be taken to mean “region of greatly enhanced cloud-cloud collision rate”.

and destroy interstellar dust—formation in stellar outflows, by coagulation, or accretion from the gas phase; destruction in shock waves from supernovae or young stars, by grain-grain collisions, or by sputtering—are all enhanced in spiral arms. See Draine (2004) for a review of interstellar dust.

In most grand-design spirals the dust lanes are displaced inside the bright-star arm. There is a natural explanation for this feature. Let us assume that the sense of rotation of the spiral density wave relative to the material in the disk is such that the gas and stars approach the density wave from the inside (in the language we shall develop later in this chapter, this occurs if spiral structure is trailing and inside corotation). If the dust lane coincides with the gas arm, it marks the location of maximum density in the interstellar gas. This high density promotes the formation of massive stars, but only after the time lag of a few Myr needed for the protostellar clouds to collapse. The bright-star arm is therefore displaced downstream (outside) from the gas arm. For a typical relative velocity between the gas and the spiral pattern of $\sim 100 \text{ km s}^{-1}$, the displacement will be a few hundred pc.

(b) Relativistic electrons The interstellar gas in galaxies contains both relativistic electrons and magnetic fields. The electrons emit synchrotron radiation as they spiral in the magnetic field, which is detectable as polarized non-thermal radio emission (BM §8.1.4). The rate of energy loss per unit volume is proportional to $n_e B^2$, where n_e is the electron density and B is the magnetic field strength. Since the magnetic field is frozen into the interstellar gas (Kulsrud 2005), compression of the gas in a spiral arm increases both n_e and B , and therefore can dramatically enhance the synchrotron emission. As this simple model would predict, well-defined radio spiral arms are seen in M51, lying inside the bright-star arms, and the magnetic field is oriented along the arms (Figure 6.2). In the nearby grand-design spiral M81 (Plate 8), however, the radio arms are broader and centered on the bright-star arms (Kaufman et al. 1989), possibly because they are strongest where supernova remnants enhance n_e . Moreover, the large-scale magnetic fields are strongest in the interarm region (Beck et al. 1996). Thus the relation of the radio arms to the other arm tracers remains somewhat unclear.

(c) Molecular gas Figure 6.3 shows the distribution of emission from the carbon monoxide (CO) molecule in M51, as traced by the rotational transition with wavelength 2.6 mm (BM §8.1.4). The CO luminosity is believed to trace the total mass in molecular gas (mostly H_2). The narrow spiral arms delineated by CO coincide closely with the dust lanes (Rand & Kulkarni 1990), as we would expect if the CO arms trace high-density gas. The enhanced CO emission arises not only because the gas density is higher, but also because the higher density both enhances the rate of collisional excitation and promotes the formation of new molecules.

(d) Neutral atomic gas Spiral structure is also seen in the surface density of neutral atomic hydrogen (HI), as measured by the 21-cm hyperfine transition (BM §8.1.4)—see Plates 5, 6, and 20. For nearby galaxies, the

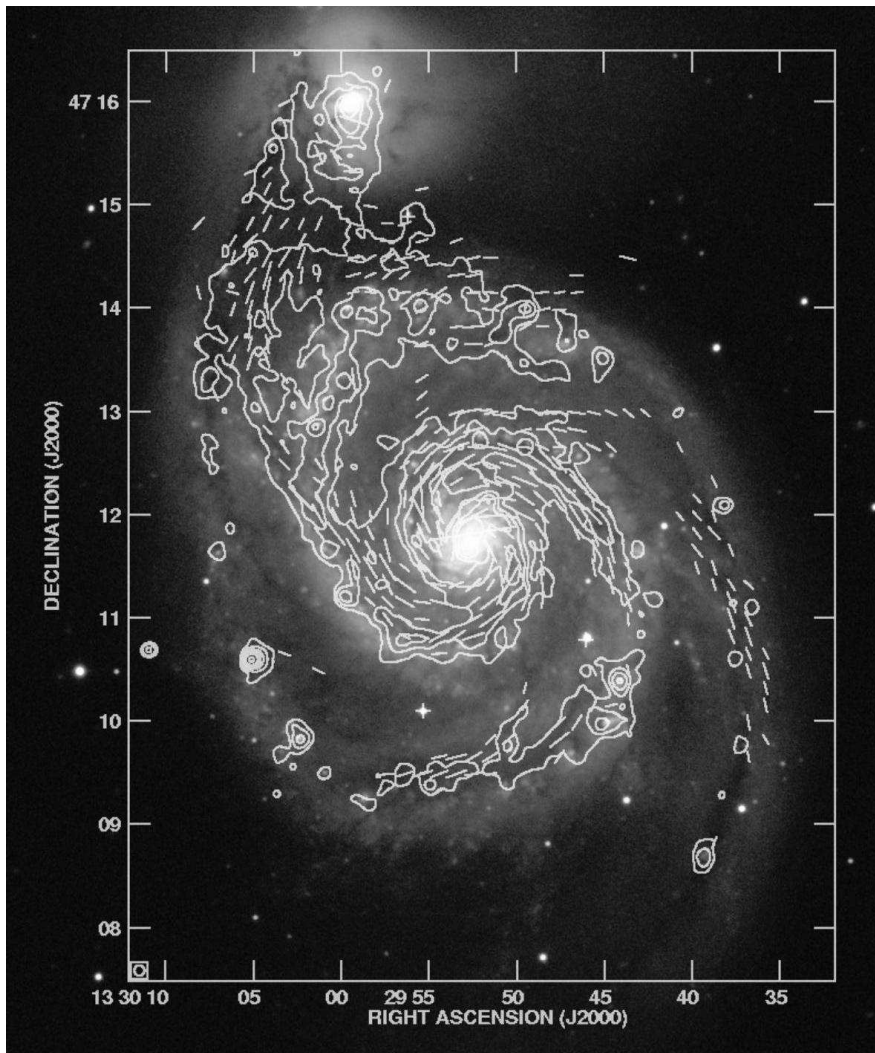


Figure 6.2 Radio emission in M51 at 6 cm wavelength, arising mainly from synchrotron radiation. The map has a resolution of 8 arcsec and is based on combined observations with the Very Large Array and the Effelsberg telescope. The vectors give the orientation of the interstellar magnetic field. The map is overlaid on an optical image of the galaxy. Credit: A. Fletcher and R. Beck (Max Planck Institute for Radio Astronomy, Bonn, Germany). See also Beck et al. (1996).

HI emission can also provide a detailed kinematic map of the mean gas velocity as a function of position in the galaxy disk, which shows how the gas is deflected or shocked by the gravitational field from the spiral arms,

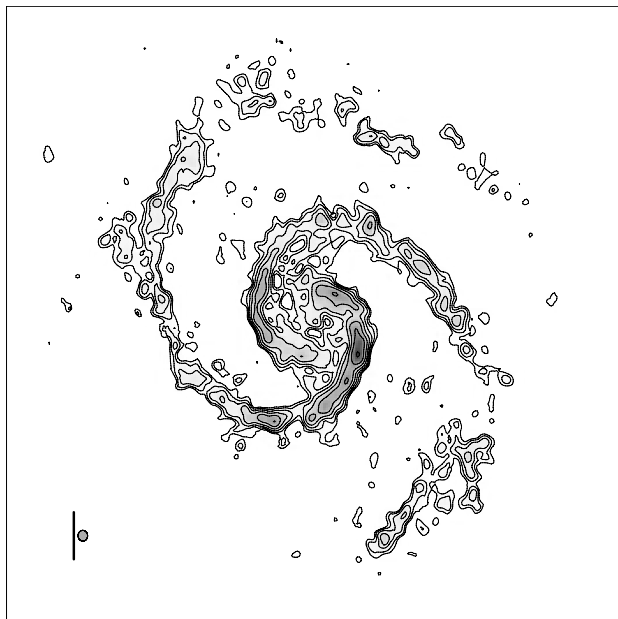


Figure 6.3 Emission from molecular gas in M51, as traced by the 2.6-mm CO line. The contours are logarithmically spaced. The vertical bar represents a distance of 1 kpc, and the adjacent small circle shows the resolution of the map. The companion galaxy NGC 5195 is just off the figure to the top. From Regan et al. (2001), reproduced by permission of the AAS.

and thus provides a further probe of the dynamics of spiral structure. The best-studied example is M81 (Visser 1980; Kaufman et al. 1989; Adler & Westpfahl 1996).

In several grand-design spirals, the HI arms coincide with the bright-star arms. This finding suggests that the enhanced density of HI in the arms does not arise from the compression of the gas, but rather from dissociation of molecular hydrogen by ultraviolet radiation from the young stars formed in the arm (Tilanus & Allen 1989). Thus in some cases the HI arms may be a product, rather than a precursor, of rapid star formation.

(e) HII regions These are regions of ionized gas that surround hot stars (see BM §8.1.3). HII regions are visible in the $H\alpha$ emission line of the Balmer series at 656 nm, which arises when a proton recombines with a free electron to form neutral hydrogen in an excited state, which then cascades down to the ground state by emitting a series of photons. HII regions are also visible at radio frequencies from their bremsstrahlung and radio recombination lines (see BM §8.1.4). The HII regions generally trace the bright-star spiral arms; groups of them are seen in several of the color plates as pinkish knots.

To summarize, in M51 the spiral arms seen in dust, non-thermal radio

continuum emission, and CO all coincide, and presumably mark the high-density gas arms caused by orbit crowding. The spiral arms in HI, young stars, and HII regions are also coincident, and lie outside the gas arms. All of these features can be explained qualitatively by density-wave theory, in which the displacement of the bright-star arms from the gas arms marks the time lag required for star formation. If the spiral is trailing (as described in the next subsection), then this displacement implies that the wave must be moving more slowly than the disk material. However, these observations do not necessarily show that the spiral structure is stationary over many orbital periods—for example, we shall argue below (§6.4.3) that the grand-design spiral in M51 is probably a transitory event, due to a recent encounter with the companion galaxy NGC 5195.

Many of the features seen in M51 are shared by other nearby grand-design spirals such as M81 and M100 (Kaufman et al. 1989; Knapen & Beckman 1996). It is far from clear, however, that the conclusions we have drawn about grand-design spirals can be applied to galaxies with intermediate-scale or flocculent spiral structure. For example, in flocculent spirals the spiral pattern in the old stars is weak or non-existent (Elmegreen & Elmegreen 1984; Thornley 1996), and the distribution of CO is more uniform and less concentrated in narrow arms (Regan et al. 2001).

We close this discussion by repeating a prescient comment made by Oort (1962) on the relation between spiral structure and the interstellar gas. The principal features that distinguish lenticular or S0 galaxies from spirals are the low density of cold interstellar gas, the absence of young stars, and the absence of spiral arms. Only a tiny fraction of gas-poor disk galaxies exhibit spiral arms, and most of these may be spirals that have recently been stripped of gas (Strom, Jensen, & Strom 1976; Kennicutt & Edgar 1986; Yamauchi & Goto 2004). Thus, even though spiral structure is present in the old disk stars, *interstellar gas is essential for persistent spiral structure.*

6.1.3 The geometry of spiral arms

(a) The strength and number of arms Disk galaxies are oriented at random to the line of sight. However, because they are thin, flat, and approximately axisymmetric, we can use the surface-brightness distribution in the sky plane to deduce the surface brightness that would be seen in a face-on view. The two principal uncertainties are the effects of dust obscuration and deviations from axisymmetry (BM §4.4), but these are usually significant only for galaxies that are seen nearly edge-on.

Consider a disk galaxy with face-on surface brightness $I(R, \phi)$, where (R, ϕ) are the usual polar coordinates in the disk plane, centered on the galactic center. If the surface brightness distribution is unchanged under a rotation through $2\pi/m$ radians, $I(R, \phi + 2\pi/m) = I(R, \phi)$, the galaxy is said to have m -fold rotational symmetry and m arms ($m > 0$).

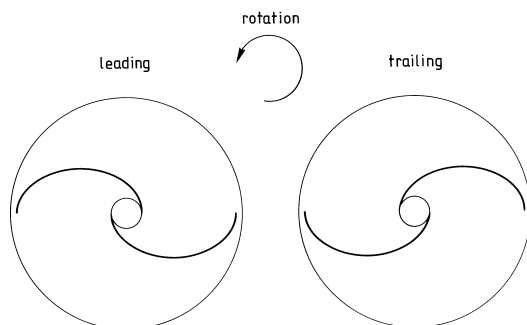


Figure 6.4 Leading and trailing arms.

The strength of the spiral structure can be parametrized by the amplitude of its Fourier components, defined by expressing the surface brightness as a Fourier series (eq. B.66),

$$\frac{I(R, \phi)}{\bar{I}(R)} = 1 + \sum_{m=1}^{\infty} A_m(R) \cos m[\phi - \phi_m(R)] \quad (A_m(R) > 0). \quad (6.1)$$

Here $\bar{I}(R) \equiv (2\pi)^{-1} \int_0^{2\pi} d\phi I(R, \phi)$ is the azimuthally averaged surface brightness at radius R , and A_m and ϕ_m are the amplitude and phase of the m th Fourier component.

If a single Fourier component m dominates the spiral structure, the strength can also be parametrized by the arm-interarm surface-brightness ratio K , which is related to A_m by

$$K = \frac{1 + A_m}{1 - A_m}. \quad (6.2)$$

Most grand-design spiral galaxies have two arms and approximate two-fold rotational symmetry. In near-infrared light, which traces the surface density, the amplitude of the arms lies in the range $0.15 \lesssim A_2 \lesssim 0.6$ (Rix & Zaritsky 1995), corresponding to arm-interarm ratios of $1.4 \lesssim K \lesssim 4$. Grand-design spirals with $m \neq 2$ are rare, although a significant fraction of disk galaxies exhibit lopsided distortions ($A_1 \gtrsim 0.2$) in their outer parts, and careful Fourier decomposition occasionally reveals three-armed spiral patterns (Rix & Zaritsky 1995). The dominance of two-armed patterns in grand-design spirals is a striking observational fact that demands explanation in a successful theory of spiral structure.

(b) Leading and trailing arms Spiral arms can be classified by their orientation relative to the direction of rotation of the galaxy. A **trailing** arm is one whose outer tip points in the direction opposite to galactic rotation, while the outer tip of a **leading** arm points in the direction of rotation (see Figure 6.4).

It is not easy to determine observationally whether the arms of a given galaxy are leading or trailing. In face-on galaxies we cannot determine the

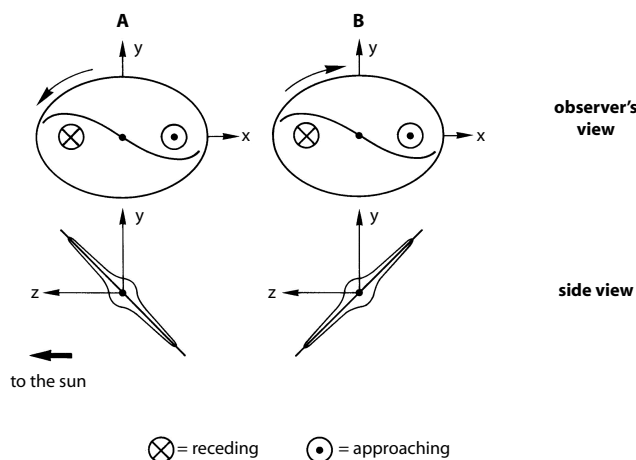


Figure 6.5 The appearance of leading and trailing arms. Galaxy A has leading arms, while galaxy B has trailing arms, but both exhibit the same image on the sky and the same radial-velocity field.

direction of rotation from radial velocities, and in edge-on galaxies we cannot see the spiral arms. Even in galaxies with intermediate inclinations the task is difficult, as we now show. Consider the two galaxies A and B in Figure 6.5. In both cases the (x, y) plane is the celestial sphere and the z axis points towards the Sun. Galaxy A is inclined so that the side nearest the Sun is in the half plane $y > 0$, while galaxy B is closer in the half plane $y < 0$. We have marked a spiral pattern and a rotation direction on both galaxies; the spiral in A is leading and in B is trailing. Despite this difference the appearance of both galaxies as seen from the Sun is the same; Figure 6.5 shows that in both systems the spiral pattern curves in an anti-clockwise direction as one moves out from the center, and the side with $x > 0$ has radial velocity towards the Sun. Thus radial-velocity measurements cannot by themselves distinguish leading and trailing spirals in thin disks.

To determine whether a given galaxy leads or trails, we must determine which side of the galaxy is closer to us. A variety of clues can be used to do this. If the inner disk is dusty, so it absorbs a significant fraction of the starlight passing through it, then the surface brightness of the bulge at a given distance along its apparent minor axis will be lower on the near side, and the number density of any population with a spheroidal distribution (such as globular clusters) will also be lower (see Figure 6.6). Similarly, dust filaments in the inner disk obscure a larger fraction of the bulge light on the near side, and hence are more prominent.

In almost all cases in which the answer is unambiguous, the spiral arms trail (Hubble 1943; de Vaucouleurs 1959; Pasha 1985). The arms in our own Galaxy trail as well (BM §4.3). There are occasional reports of galaxies with

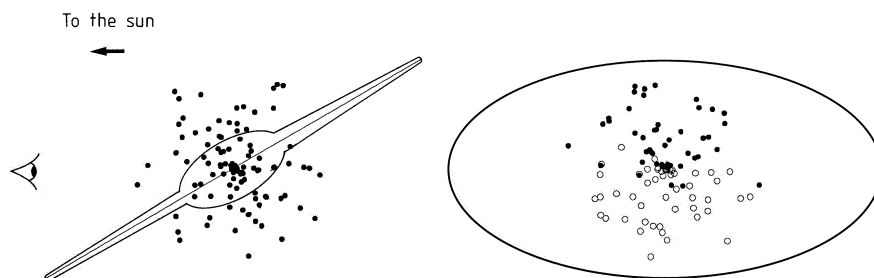


Figure 6.6 Distinguishing near and far sides of a disk galaxy. The dots represent objects such as novae or globular clusters. There is an obscuring dust layer in the central plane of the disk which is shown as a line in the side view at left. In the observer's view, at right, objects behind the dust layer (open circles) are fainter and hence fewer are present in a flux-limited survey.

leading arms (Pasha 1985; Buta, Byrd, & Freeman 2003), and transitory one-armed leading spirals can be produced by plausible dynamical processes, for example encounters with companion galaxies on retrograde orbits. Nevertheless, in the vast majority of cases *spiral arms are trailing*.

(c) The pitch angle and the winding problem The **pitch angle** α at any radius R is the angle between the tangent to the arm and the circle $R = \text{constant}$ (see Figure 6.8); by definition $0 < \alpha < 90^\circ$.

It is useful to think of the center of each arm as a mathematical curve in the plane of the galaxy, which we write in the form $\phi + g(R, t) = \text{constant}$ where t is the time. Suppose that the galaxy has m -fold rotational symmetry, that is, the arm pattern is unchanged if we rotate the galaxy by $2\pi/m$ radians ($m > 0$). Then a more convenient expression, which defines the locations of all m arms, is

$$m\phi + f(R, t) = \text{constant} \pmod{2\pi}, \quad (6.3)$$

where $f(R, t) \equiv mg(R, t)$ is the **shape function**. It is also useful to introduce the **radial wavenumber**

$$k(R, t) \equiv \frac{\partial f(R, t)}{\partial R}. \quad (6.4)$$

The sign of k determines whether the arms are leading or trailing. If, as we shall always assume, $m > 0$ and the galaxy rotates in the direction of increasing ϕ , then

$$\text{leading arms} \Leftrightarrow k < 0 \quad ; \quad \text{trailing arms} \Leftrightarrow k > 0. \quad (6.5)$$

The pitch angle is given by

$$\cot \alpha = \left| R \frac{\partial \phi}{\partial R} \right|, \quad (6.6)$$

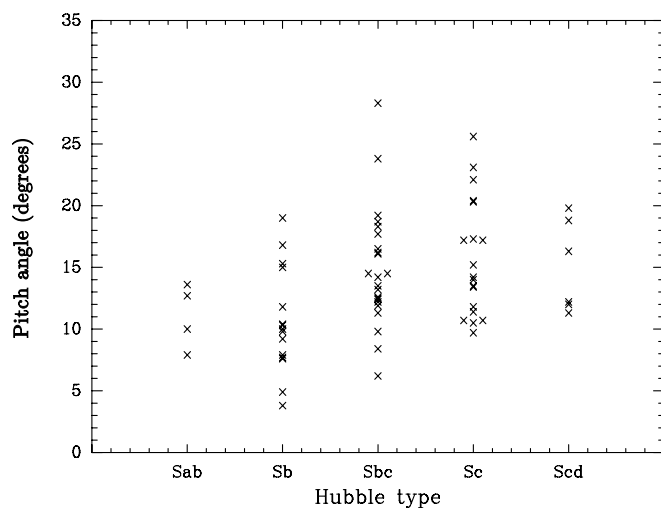


Figure 6.7 Pitch angles (eq. 6.6) in degrees, as a function of Hubble type (Ma 2002).

where the partial derivative is evaluated along the curve (6.3). Thus

$$\cot \alpha = \left| \frac{kR}{m} \right|. \quad (6.7)$$

Figure 6.7 shows the pitch angle as a function of Hubble type for a sample of spiral galaxies. The pitch angle is correlated with Hubble type—as it should be, since openness of the spiral arms is one of the criteria in Hubble’s classification scheme—but there is substantial scatter. The typical spiral has $\alpha \simeq 10^\circ\text{--}15^\circ$.

We now conduct a simple thought experiment. At some initial time $t = 0$ we paint a narrow stripe or arm radially outward across the disk of a galaxy. The initial equation of the stripe is $\phi = \phi_0$, where ϕ is the azimuthal angle (Figure 6.8). The disk rotates with an angular speed $\Omega(R)$, where R is the distance from the center of the disk. The disk is said to be in **differential rotation** if $\Omega(R)$ is not independent of R . When the disk is in differential rotation the arm does not remain radial as the disk rotates. The location of the arm $\phi(R, t)$ is described by the equation

$$\phi(R, t) = \phi_0 + \Omega(R)t. \quad (6.8)$$

The pitch angle is given by equation (6.6),

$$\cot \alpha = Rt \left| \frac{d\Omega}{dR} \right|. \quad (6.9)$$

For a galaxy with a flat circular-speed curve, $R\Omega(R) = v_c = 200 \text{ km s}^{-1}$, $R = 5 \text{ kpc}$, and $t = 10 \text{ Gyr}$, the pitch angle would now be $\alpha = 0.14^\circ$, far

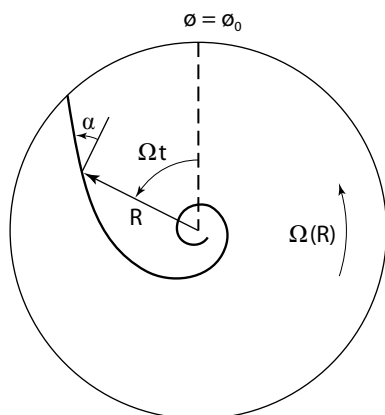


Figure 6.8 How a material arm winds up in a differentially rotating disk. The rotation law is $\Omega(R) \propto R^{-1}$.

smaller than observed pitch angles. This discrepancy is called the **winding problem**: if the material originally making up a spiral arm remains in the arm, the differential rotation of the galaxy winds up the arm in a time short compared with the age of the galaxy. A remarkably clear statement of the winding problem was given over a century ago by Wilczynski (1896).

There are several possible ways to resolve the winding problem:

- (i) It may be that the spiral pattern is statistically in a steady state, but that any individual spiral arm is quite young. If we continuously dribble cream into a freshly stirred cup of coffee, each droplet briefly takes on a spiral form before it is stretched out and disappears. Similarly, if localized luminosity features are continuously produced in a galactic disk (say, by the collapse of a gravitationally unstable patch, leading to a burst of star formation) each feature will be sheared out into a spiral, which winds up more and more but lasts only until the luminous young stars die off. This model, which we discuss further in §6.4.3, is plausible for flocculent galaxies but cannot explain grand-design spirals.
- (ii) As we have discussed, the Lin–Shu hypothesis is that spiral structure is a stationary density wave in the stellar density and gravitational potential of the disk, and hence not subject to the winding problem.²
- (iii) The spiral pattern may be a temporary phenomenon resulting from a recent violent disturbance such as a close encounter with another galaxy.

²Other types of stationary wave are possible. New stars that explode as supernovae induce further star formation in adjacent regions—direct evidence of this is seen in supershells, complex HI structures with radii of several hundred parsecs that appear to be formed by a sequence of supernovae, and are often bordered by young stars. Mueller & Arnett (1976) and Gerola & Seiden (1978) have suggested that star formation induced in this way could lead to a **detonation wave** of star formation that propagates across the entire galaxy like an infectious disease, and takes the form of a grand-design spiral. However, a wave of this kind cannot explain the presence of spiral structure in the old disk stars, and fine tuning is needed to ensure that induced star formation is reliable enough to propagate across the whole galaxy, but not so strong that it consumes all of the interstellar gas.

Encounters may indeed be responsible for many of the most striking grand-design spirals. For example, the prototypical grand-design spiral M51=NGC 5194 has recently suffered an encounter with its companion galaxy NGC 5195 (Plate 1 and Figure 6.25).

(d) The pattern speed According to the Lin–Shu hypothesis, spiral structure is a wave pattern that rotates rigidly. In this case we may define the pattern speed Ω_p to be the angular speed of rotation of the spiral wave as viewed from an inertial frame. If the amplitude of the spiral pattern is small, the material in the galaxy travels in nearly circular orbits at an angular speed that varies with radius, $\Omega(R)$, which we can assume to be positive. The radius at which $\Omega_p = \Omega(R)$, the corotation radius or corotation resonance, was already introduced in §3.3.3 in the context of weak bars. Since the material angular speed $\Omega(R)$ is a decreasing function of radius for almost all galaxies, a spiral pattern at radius R with $\Omega(R) > \Omega_p$ is said to lie “inside corotation,” while a pattern with $\Omega(R) < \Omega_p$ is “outside corotation.”

As we discussed in §6.1.2a, in grand-design spirals the dust lanes usually lie inside the arms traced by luminous stars, and this displacement may reflect the time lag between the maximum compression of the gas and the formation of stars. This explanation requires that the gas and stars approach the density wave from the inside. Since spiral arms are usually trailing, the gas must rotate faster than the spiral pattern, rather than vice versa; in other words, *the pattern in grand-design spirals is usually inside corotation.*

Several methods of varying reliability are used to estimate the pattern speed:

(i) We can fit the observed spiral pattern to analytic or numerical models of the dynamics of the galactic disk that predict the pitch angle as a function of pattern speed and radius, such as the dispersion relation that we derive in §6.2.2 below. This method requires that we understand completely the complex dynamics of the disk, that we have accurate estimates of the disk surface density, velocity dispersion, and circular-speed curve, and that a well-defined pattern speed exists.

(ii) We can use the surface-density distribution measured from the near-infrared light to calculate the gravitational potential in the disk, including the non-axisymmetric potential due to the spiral arms. We then compute the streamlines for cold gas moving in this potential, as a function of pattern speed, and find the pattern speed at which the model line-of-sight velocities best match the velocities determined from the spectral lines of HI and CO (Kranz, Slyz, & Rix 2003).

(iii) In §3.3.3 on weak bars, we introduced the Lindblad radii, at which the apparent frequency of an m -armed perturbation with a fixed pattern speed, $m(\Omega - \Omega_p)$, resonates with the epicycle frequency κ of a star. At these radii, cold gas will experience strong radial forcing, leading to shocks and thus (perhaps) to rapid star formation. These regions of intense star formation may be visible as rings (Buta 1995). The pattern speed can be

determined from the Lindblad radii once the circular-speed curve is known. Unfortunately, the epicycle frequency depends on the first derivative of the angular speed (eq. 3.80), so the result is sensitive to uncertainties in the circular-speed curve. See §6.5.2d for a discussion of rings in barred galaxies.

(iv) Suppose that there is some component of the galaxy that satisfies the continuity equation, such as old disk stars, and that we can determine the mean line-of-sight velocity v_{\parallel} and surface number density or surface brightness Σ of this population at every point on the disk. Let the (x, y) plane coincide with the disk, and let (x', y') be coordinates that rotate with the pattern speed $\Omega_p = \Omega_p \hat{\mathbf{e}}_z$. The disk is stationary in the rotating frame, so in that frame the continuity equation (F.4) reads

$$\frac{\partial}{\partial x'} (\Sigma v'_x) + \frac{\partial}{\partial y'} (\Sigma v'_y) = 0. \quad (6.10)$$

The velocity in an inertial frame is $\mathbf{v} = \mathbf{v}' + \Omega_p \times \mathbf{x}'$ (see §3.3.2), so

$$\begin{aligned} & \frac{\partial}{\partial x'} [\Sigma(v_x + \Omega_p y')] + \frac{\partial}{\partial y'} [\Sigma(v_y - \Omega_p x')] \\ &= \frac{\partial}{\partial x} (\Sigma v_x) + \frac{\partial}{\partial y} (\Sigma v_y) + \Omega_p y \frac{\partial \Sigma}{\partial x} - \Omega_p x \frac{\partial \Sigma}{\partial y} = 0, \end{aligned} \quad (6.11)$$

where in the last line we have dropped the primes on x and y since we can choose the rotating and non-rotating coordinates to coincide at any given instant.

We now integrate this equation over x from $-\infty$ to ∞ . The first and third terms yield vanishing integrals because $\Sigma \rightarrow 0$ as $|x| \rightarrow \infty$. Thus we are left with

$$\frac{\partial}{\partial y} \left(\int_{-\infty}^{\infty} dx \Sigma v_y - \Omega_p \int_{-\infty}^{\infty} dx \Sigma x \right) = 0. \quad (6.12)$$

This result implies that the quantity in parentheses is independent of y ; since it must vanish as $|y| \rightarrow \infty$ it must be zero, and we have (Tremaine & Weinberg 1984a)

$$\Omega_p = \frac{\int_{-\infty}^{\infty} dx \Sigma v_y}{\int_{-\infty}^{\infty} dx \Sigma x}. \quad (6.13)$$

If the x axis is chosen to be the line of nodes of the disk (the intersection of the disk plane with the sky plane), then $v_y = v_{\parallel} / \sin i$ where v_{\parallel} is the line-of-sight velocity relative to the center of the galaxy, and i is the inclination, the angle between the disk normal and the line of sight. Thus the quantities on the right side of equation (6.13) can all be measured. Each value of y yields an independent measurement of the pattern speed, and if all is well these measurements will be consistent.

The crucial assumptions in this method are that there *is* a well-defined pattern speed—that is, that the surface density Σ and velocity field \mathbf{v} are

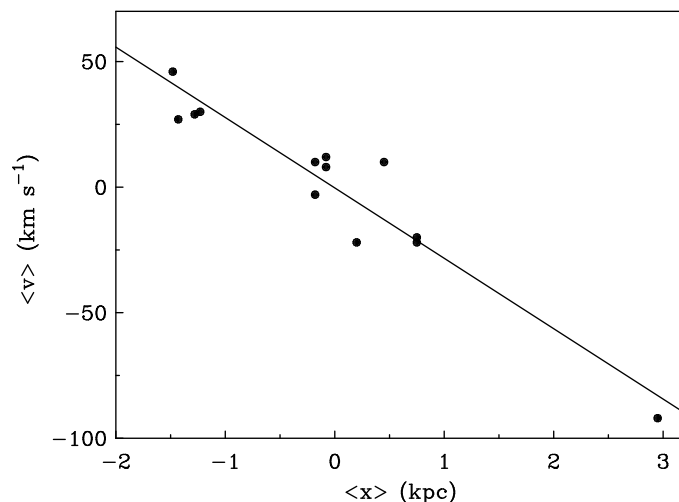


Figure 6.9 Determination of the pattern speed in M100 (Plate 17), using emission from molecular gas. The numerator (vertical axis) and denominator (horizontal axis) of the expression in equation (6.13), plotted for each of 13 integration paths parallel to the major axis of the galaxy. The slope of the line from the origin to each data point gives an independent measurement of the pattern speed. The best-fit slope, shown by the line, is $\Omega_p = 28 \text{ km s}^{-1} \text{ kpc}^{-1}$. From Rand & Wallin (2004).

stationary in time except for an overall rotation—and that the population we are examining satisfies the continuity equation—that is, its members are neither created nor destroyed. The second condition is satisfied by old disk stars (cf. §4.1.1a) but not by luminous young stars, which are short-lived, or by HI gas, since atomic gas can be converted into molecular gas and vice versa (page 467).

So far the method has mostly been used to measure the pattern speeds of bars in SB0 and SBa galaxies with little or no recent star formation (§6.5.1); it is harder to apply to spiral galaxies because the luminous young stars obscure the contribution of the old stars to the spectrum. Rand & Wallin (2004) have applied this method to spirals in which the interstellar gas is dominated by the molecular phase, arguing that in this case a large fraction of the molecular gas cannot be rapidly created or destroyed so the continuity equation should apply. They derive pattern speeds for five galaxies including M100, for which they find $\Omega_p = (28 \pm 5) \text{ km s}^{-1} \text{ kpc}^{-1}$, consistent with estimates from other methods (Figure 6.9).

(v) The pattern speeds of bars are easier to measure than the pattern speeds of spirals (§6.5). Since the spiral arms usually appear to emerge from the tips of the bar, and since strong bars are expected to drive spiral structure, it is natural to assume that the pattern speeds of the spiral and the

bar are equal. Unfortunately this appealing argument leads to an unsettling result: in general the corotation radius for a bar lies close to the end of the bar, which implies that most of the spiral pattern will lie *outside* corotation. However, the dust lanes in barred spirals lie inside the arms traced by luminous stars, so, as we argued at the start of this section, the spirals should lie *inside* corotation. Sellwood & Sparke (1988) suggest that the flaw in this chain of argument is that in fact the spirals have a much smaller pattern speed than the bar. They illustrate from N-body simulations that the spiral appears to be connected to the tips of the bar for most of the beat period between the two pattern speeds, even when they are quite different.

All of these methods assume the stationary spiral structure hypothesis, that is, that the spiral pattern actually *has* a well-defined pattern speed. The agreement of the points in Figure 6.9 with the dashed line provides some reassurance that this assumption is correct in M100. However, a secure confirmation of the stationary spiral structure hypothesis would require accurate, consistent measurements of the spiral pattern speed at a variety of radii, using several different methods. We are far from achieving this goal with any galaxy, including our own.

6.1.4 The anti-spiral theorem

Newton's equations of motion and law of gravitation are time-reversible. If we make a movie of the trajectories of N point masses interacting through their mutual gravity, the trajectories seen when we run the movie backwards are also dynamically possible. Thus, if the dynamics of a spiral galaxy is governed by Newton's equations, and the galaxy is in a stationary state—that is, if there is some frame rotating at a pattern speed Ω_p in which the galaxy's DF is time-independent—then a time-reversed movie of the galaxy should also represent a possible steady-state solution of the equations. However, time reversal changes trailing spirals into leading spirals, by changing the sign of all velocities without changing the instantaneous luminosity density. This argument is the basis of the **anti-spiral theorem**, which states that *if a stationary solution of a time-reversible set of equations has the form of a trailing spiral, then there must be an identical solution in the form of a leading spiral* (Lynden-Bell & Ostriker 1967).

The anti-spiral theorem implies that spiral galaxies cannot be understood simply as steady-state solutions of the collisionless Boltzmann equation and Newton's law of gravity—the prevalence of trailing spirals demands an additional ingredient in the physics. The most likely explanations are that either (i) the spirals are not in a steady state (for example, they are a result of a recent disturbance), or are somehow continuously regenerated; or (ii) the spiral form is influenced by processes that are not time-reversible, such as dissipation in the interstellar gas or absorption at a Lindblad resonance.

6.1.5 Angular-momentum transport by spiral-arm torques

Spiral structure produces a spiral gravitational field, which exerts torques and transfers angular momentum from one part of the disk to another. We now investigate the strength of these torques and their consequences.

Consider a galaxy with disk density $\rho_d(\mathbf{x}, t)$ and potential $\Phi(\mathbf{x}, t)$. We assume that the disk is symmetric about the plane $z = 0$, and focus on the disk material outside a cylinder of radius R_0 whose symmetry axis is the z axis. The torque per unit mass exerted by the gravitational potential is $-\partial\Phi/\partial\phi$, so the z -component of the torque exerted on the disk material outside this cylinder is

$$C_G(R_0) = - \int_{R_0}^{\infty} dR R \int_0^{2\pi} d\phi \int_{-\infty}^{\infty} dz \rho_d \frac{\partial\Phi}{\partial\phi}. \quad (6.14)$$

In principle the torque due to a given spiral density field can be found directly from this formula, but there is a more instructive approach. We assume that the non-axisymmetric component of the potential is generated by the disk, so $\partial\Phi/\partial\phi$ can be replaced by $\partial\Phi_d/\partial\phi$, where Φ_d is the disk potential, and ρ_d is related to Φ_d by Poisson's equation (2.10). Using equation (B.52) for ∇^2 in cylindrical coordinates, we have

$$\begin{aligned} C_G(R_0) = & -\frac{1}{4\pi G} \int_{R_0}^{\infty} dR R \int_0^{2\pi} d\phi \int_{-\infty}^{\infty} dz \\ & \times \left[\frac{1}{R} \frac{\partial}{\partial R} \left(R \frac{\partial\Phi_d}{\partial R} \right) + \frac{1}{R^2} \frac{\partial^2\Phi_d}{\partial\phi^2} + \frac{\partial^2\Phi_d}{\partial z^2} \right] \frac{\partial\Phi_d}{\partial\phi}. \end{aligned} \quad (6.15)$$

The contribution of the second term in square brackets vanishes, since

$$\int_0^{2\pi} d\phi \frac{\partial^2\Phi_d}{\partial\phi^2} \frac{\partial\Phi_d}{\partial\phi} = \frac{1}{2} \int_0^{2\pi} d\phi \frac{\partial}{\partial\phi} \left(\frac{\partial\Phi_d}{\partial\phi} \right)^2 = 0. \quad (6.16)$$

The remaining terms can be integrated by parts with respect to R and z :

$$\begin{aligned} C_G(R_0) = & -\frac{1}{4\pi G} \int_0^{2\pi} d\phi \int_{-\infty}^{\infty} dz R \frac{\partial\Phi_d}{\partial R} \frac{\partial\Phi_d}{\partial\phi} \Big|_{R=R_0}^{R=\infty} \\ & - \frac{1}{4\pi G} \int_{R_0}^{\infty} dR R \int_0^{2\pi} d\phi \frac{\partial\Phi_d}{\partial z} \frac{\partial\Phi_d}{\partial\phi} \Big|_{z=-\infty}^{z=\infty} \\ & + \frac{1}{4\pi G} \int_{R_0}^{\infty} dR R \int_0^{2\pi} d\phi \int_{-\infty}^{\infty} dz \left(\frac{\partial\Phi_d}{\partial R} \frac{\partial^2\Phi_d}{\partial\phi\partial R} + \frac{\partial\Phi_d}{\partial z} \frac{\partial^2\Phi_d}{\partial\phi\partial z} \right). \end{aligned} \quad (6.17)$$

The final line can be shown to vanish, using manipulations similar to those in (6.16). The terms that involve the potential at $z \rightarrow \pm\infty$ or $R \rightarrow \infty$ also vanish, since the potential decays to zero at large distances. Hence

$$C_G(R_0) = \frac{R_0}{4\pi G} \int_0^{2\pi} d\phi \int_{-\infty}^{\infty} dz \frac{\partial\Phi_d}{\partial R} \frac{\partial\Phi_d}{\partial\phi} \Big|_{R=R_0}. \quad (6.18)$$

This formula can be evaluated explicitly for tightly wound spiral structure, in which the pitch angle is small or the radial wavenumber k satisfies $|kR| \gg 1$ (we shall discuss this approximation in greater detail in §6.2.2). Let us assume that the disk is razor-thin, and that the non-axisymmetric component of the surface density has the form

$$\Sigma_d(R, \phi) = \Sigma_1(R) \cos[m\phi + f(R)] \quad (m > 0), \quad (6.19)$$

where $\Sigma_1(R)$ is a slowly varying function of radius, $f(R)$ is the shape function (eq. 6.3), $df/dR = k$ (eq. 6.4), and $|kR| \gg 1$; we may neglect any time dependence since the torque is determined by the instantaneous surface density and potential. The corresponding potential is found by inspecting the real parts of equations (6.29) and (6.30) below:

$$\Phi_d(R, \phi, z) = \Phi_1(R) e^{-|kz|} \cos[m\phi + f(R)], \quad \text{where} \quad \Phi_1 = -\frac{2\pi G \Sigma_1}{|k|}. \quad (6.20)$$

The torque (6.18) is then

$$C_G(R_0) = \text{sgn}(k) \frac{m R_0 \Phi_1^2}{4G} = \text{sgn}(k) \frac{\pi^2 m R_0 G \Sigma_1^2}{k^2}, \quad (6.21)$$

where $\text{sgn}(k) = \pm 1$ depending on the sign of k . In deriving this result we neglect derivatives of $\Sigma_1(R)$ in comparison to radial derivatives of $\cos[m\phi + f(R)]$, since $|kR| \gg 1$ while $\Sigma_1(R)$ is slowly varying.

Since $m > 0$, the sign of the torque depends only on the sign of the wavenumber k and hence on whether the spiral arms are leading or trailing: trailing arms ($k > 0$) exert positive gravitational torque on the outer parts of the disk, and thereby transport angular momentum from the inner to the outer disk, while leading arms transport angular momentum inward.

To estimate the strength of this torque, consider a Mestel disk, in which the unperturbed surface density is $\Sigma_0(R) = v_0^2/(2\pi GR)$ (eq. 4.158), and suppose that the amplitude of the spiral surface density pattern is a fraction A_m of the axisymmetric density, $\Sigma_1(R) = A_m \Sigma_0(R)$. The total angular momentum contained in the unperturbed disk inside radius R_0 is

$$L(R_0) = 2\pi \int_0^{R_0} dR \Sigma_0(R) v_0 R^2 = \frac{v_0^3 R_0^2}{2G}. \quad (6.22)$$

After a time t , the fraction of the angular momentum initially inside R_0 that has been transported outside R_0 by the spiral gravitational torque is

$$\begin{aligned} \frac{C_G t}{L(R_0)} &= \frac{1}{2} m \text{sgn}(k) \frac{A_m^2}{(kR_0)^2} \frac{v_0 t}{R_0} \\ &= \text{sgn}(k) \frac{\tan^2 \alpha}{m} \pi A_m^2 N_{\text{rot}}, \end{aligned} \quad (6.23)$$

where N_{rot} is the age in units of the local rotation period $2\pi R_0/v_0$, and equation (6.7) has been used to write kR_0 in terms of the pitch angle α .

We have seen in §6.1.3 that typical grand-design spirals have surface-density amplitudes A_2 in the range 0.15–0.6, pitch angles of 10° – 15° , corresponding to $\tan\alpha \simeq 0.2$ – 0.3 , and ages $N_{\text{rot}} \simeq 50$ – 100 . For these values the ratio of the time-integral of the torque to the reservoir of angular momentum in equation (6.23) varies from about 0.05 to 5. We conclude that in some cases spiral gravitational torques might significantly rearrange the angular-momentum distribution in the galactic disk over the galactic age of 10 Gyr, if the spiral pattern were permanently present. A similar conclusion was reached by Gnedin, Goodman, & Frei (1995) for the grand-design spiral M100 (Plate 17). These rather large values of $C_G t/L(R)$ for strong, open, grand-design spirals suggest that such patterns cannot be permanent, since they would dramatically alter the galaxy’s angular-momentum distribution on a timescale much less than its age. In other words, it is likely that prominent grand-design spirals persist for only $\lesssim 1$ Gyr.

Gravitational torques tell only part of the story. Just as ordinary sound waves carry momentum and energy by advection (Appendix F.3.1), spiral waves carry angular momentum and energy (see Appendix J). The advective transport is generally comparable in magnitude to transport by gravitational torques. The rate at which angular momentum is transported out of the region inside R_0 is the sum of the gravitational torque $C_G(R_0)$ and the advective current of angular momentum $C_A(R_0)$ (see §6.2.6).

Angular-momentum transport by spiral arms is an example of **secular evolution** in galaxies, slow changes due to internal dynamical processes. Other examples of secular evolution include gas inflow resulting from dissipation in the interstellar gas, the excitation of random velocities of disk stars by the gravitational fields from molecular clouds or spiral arms (often called disk “heating”; see §8.4), the conversion of the inner parts of disks into thick structures that resemble bulges (**pseudobulges**), and the slowing of bar pattern speeds due to dynamical friction from the dark halo (§8.1.1d). Early in the history of the universe, most galaxy evolution was due to external processes such as mergers (§§9.2 and 9.3). As a result of the aging and expansion of the universe, the rate of infall and mergers has declined steadily with time (Figure 9.13). Thus galaxy evolution is experiencing a gradual transition from an early phase dominated by rapid, violent, external events such as mergers to a late phase dominated by slower, more gradual internal processes (Kormendy & Kennicutt 2004).

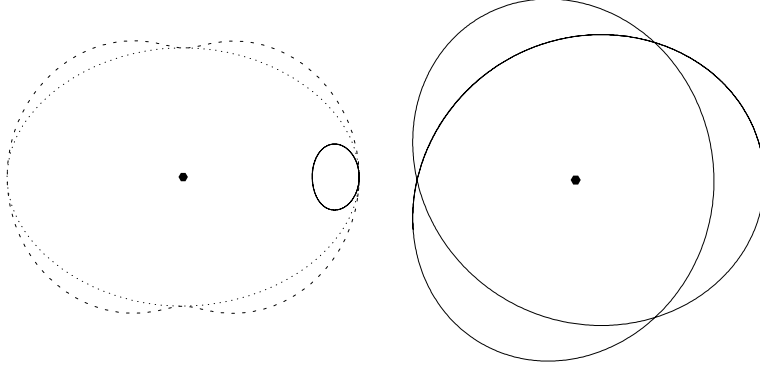


Figure 6.10 The appearance of elliptical orbits in a frame rotating at $\Omega_p = \Omega - n\kappa/m$. Left: $(n, m) = (0, 1)$, solid line; $(1, 2)$, dotted line; $(1, -2)$, dashed line. Right: $(n, m) = (2, 3)$.

6.2 Wave mechanics of differentially rotating disks

6.2.1 Preliminaries

(a) Kinematic density waves The galactocentric distance of a particle that orbits in the equatorial plane of an axisymmetric galaxy is a periodic function of time with period T_r (see eq. 3.17). During the interval T_r the azimuthal angle increases by an amount $\Delta\phi$ (eq. 3.18b). These quantities are related to the radial and azimuthal oscillation frequencies $\Omega_r = 2\pi/T_r$ and $\Omega_\phi = \Delta\phi/T_r$. In general, $\Delta\phi/(2\pi)$ is irrational, so the orbit forms a rosette figure such as the one shown in Figure 3.1.

Now suppose that we view the orbit from a frame that rotates at angular speed Ω_p . In this frame, the azimuthal angle is $\phi_p = \phi - \Omega_p t$, which increases in one radial period by $\Delta\phi_p = \Delta\phi - \Omega_p T_r$. Therefore we can choose Ω_p so that the orbit is closed; in particular, if $\Delta\phi_p = 2\pi n/m$, where m and n are integers, the orbit closes after m radial oscillations. In this case

$$\Omega_p = \Omega_\phi - \frac{n\Omega_r}{m} \simeq \Omega - \frac{n\kappa}{m}, \quad (6.24)$$

where in the last equality we have approximated Ω_ϕ and Ω_r by their values for nearly circular orbits, the circular frequency Ω and the epicycle frequency κ (see eqs. 3.79). The appearance of the closed orbits in the rotating frame is shown in Figure 6.10.

In general $\Omega(R) - n\kappa(R)/m$ will be a function of radius, so no single choice for Ω_p can ensure that orbits at all radii are closed. In Figure 6.11 we show the behavior of $\Omega - n\kappa/m$ for several values of m and n . The curves are

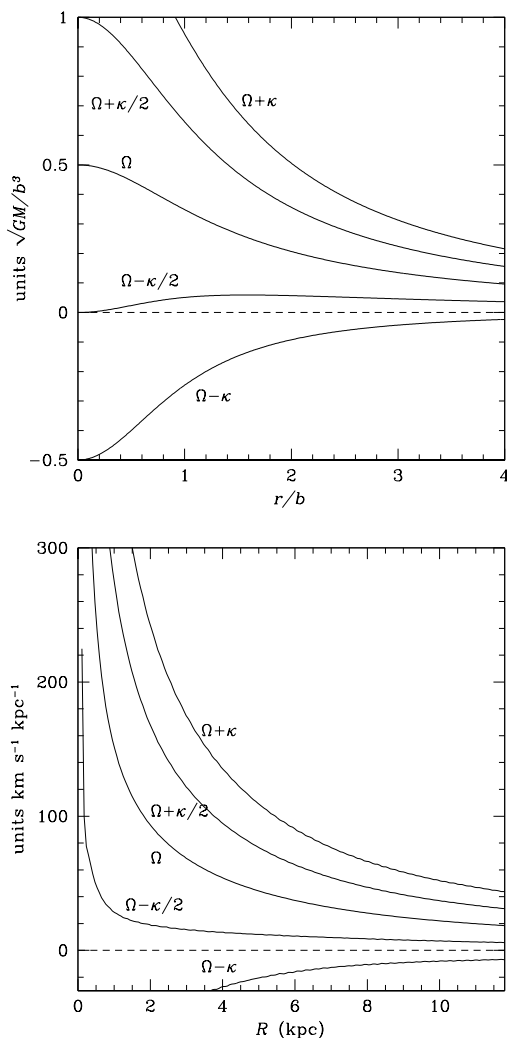


Figure 6.11 Behavior of $\Omega - n\kappa/m$ in: (top) the isochrone potential (eq. 2.47); (bottom) Model I for our Galaxy, described in §2.7.

plotted for two representative galactic circular-speed curves, the isochrone potential (eq. 2.47) and Model I for our Galaxy, as described in Table 2.3.

This diagram exhibits an intriguing fact noticed by Lindblad many decades ago: while most of the $\Omega - n\kappa/m$ curves vary rapidly with radius, the curve for $n = 1, m = 2$ (or $n = 2, m = 4$, etc.) is relatively constant across much of the galaxy.³ To understand the significance of Lindblad's re-

³ This result is related to the shape of galaxy circular-speed curves in their inner parts. In most galaxies the circular speed rises linearly from the center with a steep slope. Thus, both Ω and κ are large at small radii, so for most values of m and n , $|\Omega - n\kappa/m|$ is much larger near the center than at large radii. However, in the central region where the

mark, let us suppose for the moment that $\Omega - \frac{1}{2}\kappa$ were *exactly* constant, and equal to some number Ω_p . Then in a frame rotating at Ω_p the orbits of the type shown as a dotted line in the left panel of Figure 6.10a would be exactly closed at every radius. Hence we could set up a nested, aligned set of these orbits covering a range of radii, as shown in Figure 6.12a. If we fill up these orbits with stars we create a bar-like pattern, which is stationary in the rotating frame and appears as a density wave rotating at the pattern speed Ω_p in the inertial frame. By rotating the axes of the ellipses we can create leading or trailing spiral density waves as in Figure 6.12b and c.

In a real galaxy $\Omega - \frac{1}{2}\kappa$ is not exactly constant. Hence, no matter what the value of Ω_p , most orbits are not exactly closed. The orientations of different orbits drift at slightly different speeds, so the pattern tends to twist or wind up. This is a modified version of the winding problem which we have already discussed—but now applied to density waves rather than material arms—and the rate of winding can be calculated in a similar way. Let $\phi_p(R, t)$ be the angle of the major axis of the pattern, as viewed in the frame rotating at the pattern speed. Let us suppose that the major axes are aligned at time $t = 0$; thus $\phi_p(R, 0) = \phi_0$. The drift rate is $\partial\phi_p/\partial t = \Omega - \frac{1}{2}\kappa - \Omega_p$; thus

$$\phi_p(R, t) = \phi_0 + [\Omega(R) - \frac{1}{2}\kappa(R) - \Omega_p]t \quad (6.25)$$

(cf. eq. 6.8). Equation (6.6) now gives the pitch angle as

$$\cot \alpha = Rt \left| \frac{d(\Omega - \frac{1}{2}\kappa)}{dR} \right|. \quad (6.26)$$

In Model I for the Galactic potential of §2.7, the average of $|R d(\Omega - \frac{1}{2}\kappa)/dR|$ is about $7 \text{ km s}^{-1} \text{ kpc}^{-1}$ between 5 and 10 kpc, and after $t = 10 \text{ Gyr}$ the pitch angle in this region is about $\alpha = 0.8^\circ$. For comparison we computed after equation (6.9) that a material arm would have $\alpha < 0.2^\circ$ in a galaxy with a similar circular-speed curve. Thus, the wave pattern winds up more slowly than the material arm by a factor of five or so. Although the pitch angle is still too small by a factor 10–20, we have come some way towards resolving the winding problem. We conclude that in galaxies with circular-speed curves similar to our own, $n = 1, m = 2$ *density waves can resist the winding process much better than material arms*. This result suggests a natural explanation for the prevalence of two-armed spirals, providing we can find a way to adjust the slow drift rates of all the orbits to a common standard.

Density waves of the type described above are called **kinematic density waves** because they involve only the kinematics of orbits in an axisymmetric potential. In a galaxy, the orbits will deviate from the paths we have assumed

circular speed is roughly proportional to radius, $\Omega \simeq \frac{1}{2}\kappa$ (see eq. 3.80). Thus $\Omega - \frac{1}{2}\kappa$ is much smaller near the center than $\Omega - n\kappa/m$ for other values of n and m .

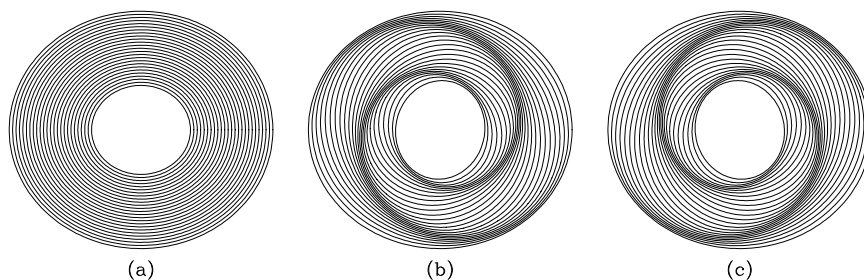


Figure 6.12 Arrangement of closed orbits in a galaxy with $\Omega - \frac{1}{2}\kappa$ independent of radius, to create bars and spiral patterns (after Kalnajs 1973b).

because the spiral pattern itself produces a non-axisymmetric component of the gravitational field. A major goal of spiral-structure theory is to determine whether the non-axisymmetric gravitational field due to the spiral itself can coordinate the drift rates of the orbits in such a way as to produce long-lived spiral patterns.

(b) Resonances Orbits, like springs and drums, have natural resonant frequencies. If the gravitational field generated by spiral structure perturbs an orbit near one of its resonant frequencies, then the response of the orbit is strong, even when the perturbing field is weak. To investigate the response of a stellar disk to non-axisymmetric forces, an essential first step is to locate the resonant orbits.

A gravitational potential that is stationary in a rotating frame can be written in the form $\Phi_1(R, \phi, t) = \Phi(R, \phi - \Omega_p t)$, where Ω_p is the pattern speed of the potential. Examples of systems that generate potentials of this form include the rotating bars seen at the centers of many disk galaxies, a satellite galaxy on a circular orbit in the disk plane, and any stationary spiral structure pattern (i.e., any structure with a well-defined pattern speed). More complicated potentials can be regarded as superpositions of potentials with different pattern speeds.

We now examine the effect of a weak potential of this form on a disk composed of stars on circular or near-circular orbits. Since the potential is periodic in $\phi - \Omega_p t$, it can be decomposed into a series of terms proportional to $\cos[m(\phi - \Omega_p t) + f_m(R)]$, where $m \geq 0$ is an integer. We studied orbits in potentials of this form in §3.3.3, and found that resonances occurred when the circular frequency Ω and the epicycle frequency κ in the unperturbed axisymmetric potential satisfied one of three conditions: $\Omega = \Omega_p$ (corotation resonance), $m(\Omega - \Omega_p) = \kappa$ (inner Lindblad resonance), or $m(\Omega - \Omega_p) = -\kappa$ (outer Lindblad resonance). These resonances occur at specific radii in a differentially rotating disk. The location and even the existence of these radii, called the corotation and Lindblad radii, depend on the circular-speed curve and the pattern speed. For example, inspection of Figure 6.11 shows that the isochrone potential has zero or two inner Lindblad radii with $m = 2$,

depending on the value of Ω_p , while Model I for the Galaxy has one inner Lindblad radius if $\Omega_p > 0$, and 0 otherwise.

Note that the condition (6.24) for a stationary kinematic density wave with $n = \pm 1$ is identical to the condition for a Lindblad resonance. This is to be expected: near resonance a weak perturbation with m -fold rotational symmetry can produce a strong response with the same symmetry, hence we expect that when the Lindblad resonance condition is satisfied exactly, a stationary m -fold wave pattern can persist even in the absence of any perturbing force.

6.2.2 The dispersion relation for tightly wound spiral arms

We analyze the behavior of density waves in disks by a three-step process. First, we use Poisson's equation to calculate the gravitational potential of an assumed surface-density pattern. Second, we determine how this potential affects the stellar orbits and thus alters the surface density in the galaxy. Finally, we match this response surface density to the input surface density to obtain a self-consistent density wave.

We approximate the disk as razor-thin so we can work in two spatial dimensions, rather than three. We also assume that the density perturbations imposed on the original axisymmetric disk are small, so we can analyze the dynamics using linear perturbation theory. Then self-consistent density waves are simply linear modes of the galactic disk, which can be computed using the Kalnajs matrix method that we introduced for spherical stellar systems in §5.3.2.

Calculating the modes of a stellar disk is a difficult task. The modes can be found analytically only for the Kalnajs disks (§5.6.2), and numerical mode calculations have been done for only a handful of models, based on simple potentials such as the isochrone disk (eq. 2.47), the Kuzmin disk (eq. 2.68), and power-law disks (Zang 1976; Kalnajs 1978; Vauterin & Dejonghe 1996; Pichon & Cannon 1997; Evans & Read 1998; Jalali & Hunter 2005).

(a) The tight-winding approximation One of the principal difficulties in mode calculations is that gravity is a long-range force, so perturbations in all parts of the system are coupled. In the early 1960s a number of workers, notably A. J. Kalnajs, C. C. Lin, and A. Toomre, realized that for tightly wound density waves (waves whose radial wavelength is much less than the radius), the long-range coupling is negligible, the response is determined locally, and the relevant solutions are analytic. As we shall see, this **tight-winding, short-wavelength, or WKB approximation**⁴ is an indispensable tool for understanding the properties of density waves in differentially rotating disks.

How tightly wound *are* spiral arms in real galaxies? The radial separation between adjacent arms at a given azimuth is ΔR , where $|f(R + \Delta R, t) -$

⁴Named after the Wentzel–Kramers–Brillouin approximation of quantum mechanics.

$f(R, t) = 2\pi$, and $f(R, t)$ is the shape function (eq. 6.3). If the arms are tightly wound, we may replace $f(R + \Delta R, t)$ by $f(R, t) + (\partial f / \partial R)\Delta R$; using equations (6.4) and (6.7) we can replace $\partial f / \partial R$ by the wavenumber k , to find

$$\Delta R = \frac{2\pi}{|k|} = \frac{2\pi R}{m} \tan \alpha. \quad (6.27)$$

Thus for tightly wound arms the radial wavelength is $2\pi/|k|$, consistent with the usage of the term “wavenumber” in other branches of physics. Figure 6.7 shows that typical pitch angles in spiral galaxies are between 10° and 15° . For two-armed spirals, this implies $|kR| \simeq 7$ –11. The WKB approximation requires $|kR| \gg 1$, although in many situations it works fairly well even for $|kR|$ as small as unity. Thus the WKB approximation is satisfied by most spiral galaxies, but not by a very comfortable margin.

Evidently the results of analyses based on the WKB approximation must be applied with caution: they should be used more as an aid to interpreting numerical mode calculations and N-body simulations than as a definitive theory.

Although the terms “tightly wound” and “short-wavelength” are often applied interchangeably to density waves, there are situations in which these are not equivalent: (i) Axisymmetric waves ($m = 0$) have zero pitch angle (eq. 6.7) and hence are always tightly wound, even though their wavelengths may not be short. The analysis we describe below is valid for axisymmetric waves only if they also have short wavelengths, $|kR| \gg 1$. (ii) If the azimuthal wavenumber is large, waves may have short wavelength without being tightly wound, if $|kR| \sim m \gg 1$. In this case the analysis below fails, and we must resort to a quite different approximation, the sheared sheet of §§8.3.2 and 8.4.2 (Goldreich & Lynden-Bell 1965a; Julian & Toomre 1966; Goldreich & Tremaine 1978).

(b) Potential of a tightly wound spiral pattern The surface density of a zero-thickness disk can be represented mathematically as the sum of an axisymmetric or unperturbed surface density $\Sigma_0(R)$, and a perturbed surface density $\Sigma_1(R, \phi, t)$ which represents the spiral pattern. For a tightly wound spiral it is convenient to write Σ_1 in a form that separates the rapid variations in density as one passes between arms from the slower variation in the strength of the spiral pattern as one moves along an arm. We may accomplish this by writing

$$\Sigma_1(R, \phi, t) = H(R, t)e^{i[m\phi + f(R, t)]}, \quad (6.28)$$

where $f(R, t)$ is the shape function of equation (6.3), and $H(R, t)$ is a slowly varying function of radius that gives the amplitude of the spiral pattern. As usual, the physical surface density is given by the real part of equation (6.28). This expression assumes that the surface-density variation is approximately sinusoidal in radius, which proves to be correct for the linear perturbation

theory that we examine below; more complicated surface-density variations can be Fourier decomposed into a sum of sinusoids.

The next step is to determine the gravitational potential due to the pattern (6.28). Since the perturbed surface density oscillates rapidly around zero mean, there will be nearly complete cancellation of the contribution from the distant parts of the pattern to the local potential; in other words, the perturbed potential at a given location will be almost entirely determined by the properties of the pattern within a few wavelengths of that location. Thus, to determine the potential in the neighborhood of a point (R_0, ϕ_0) , we may replace the shape function $f(R, t)$ by the first two terms in its Taylor series, $f(R_0, t) + k(R_0, t)(R - R_0)$. Hence

$$\Sigma_1(R, \phi, t) \simeq \Sigma_a e^{ik(R_0, t)(R - R_0)}, \text{ where } \Sigma_a = H(R_0, t)e^{i[m\phi_0 + f(R_0, t)]}. \quad (6.29)$$

We have neglected variations with angle ϕ since these are much slower than radial variations when the wave is tightly wound. Equations (6.29) show that in the vicinity of (R_0, ϕ_0) , the spiral wave closely resembles a plane wave with wavevector $\mathbf{k} = k\hat{\mathbf{e}}_R$. The potential of a plane wave in a razor-thin disk was determined in §5.6.1. According to equation (5.161),

$$\Phi_1(R, \phi, z, t) \simeq \Phi_a e^{ik(R_0, t)(R - R_0) - |k(R_0, t)z|}, \text{ where } \Phi_a = -\frac{2\pi G \Sigma_a}{|k|}. \quad (6.30)$$

We are now free to set $R = R_0$, $\phi = \phi_0$, and $z = 0$, thereby obtaining our final result for the potential in the plane due to the surface density (6.28),

$$\Phi_1(R, \phi, t) = -\frac{2\pi G}{|k|} H(R, t) e^{i[m\phi + f(R, t)]}. \quad (6.31)$$

The fractional error in this result is $O(|kR|^{-1})$. An alternative expression can be obtained by differentiating equation (6.31) with respect to radius and neglecting the derivative of $H(R, t)$ compared to the derivative of $f(R, t)$, which also involves an error of order $|kR|^{-1}$ and hence does not degrade the accuracy of the approximation any further. We find

$$\Sigma_1(R, \phi, t) = \frac{i \operatorname{sgn}(k)}{2\pi G} \frac{\partial}{\partial R} \Phi_1(R, \phi, t), \quad (6.32)$$

again with fractional error $O(|kR|^{-1})$. Since this result involves the wavenumber k only through its sign, it is valid for any tightly wound spiral whose Fourier decomposition involves predominantly either leading ($kR \ll -1$) or trailing ($kR \gg 1$) waves. Shu (1970) gives a more accurate version whose fractional error is only $O(|kR|^{-2})$:

$$\Sigma_1(R, \phi, t) = \frac{i \operatorname{sgn}(k)}{2\pi G \sqrt{R}} \frac{\partial}{\partial R} \left[\sqrt{R} \Phi_1(R, \phi, t) \right]. \quad (6.33)$$

(c) The dispersion relation for fluid disks We now determine the response of the galactic disk to a potential perturbation. Since fluid disks are simpler than stellar disks, we consider spiral structure in fluid disks before we tackle the more complicated stellar disks. Many of the results that we obtain presage similar features of stellar disks.

We shall begin without using the tight-winding approximation, only invoking it after equation (6.50) when it is needed to avoid numerical calculations. We again neglect the thickness of the disk, and assume that the pressure p acts only in the disk plane. Thus the motion is confined to the plane $z = 0$. Using equation (B.56) for $(\mathbf{v} \cdot \nabla)\mathbf{v}$ in cylindrical coordinates, Euler's equations (F.10) can be written

$$\begin{aligned} \frac{\partial v_R}{\partial t} + v_R \frac{\partial v_R}{\partial R} + \frac{v_\phi}{R} \frac{\partial v_R}{\partial \phi} - \frac{v_\phi^2}{R} &= -\frac{\partial \Phi}{\partial R} - \frac{1}{\Sigma_d} \frac{\partial p}{\partial R} \\ \frac{\partial v_\phi}{\partial t} + v_R \frac{\partial v_\phi}{\partial R} + \frac{v_\phi}{R} \frac{\partial v_\phi}{\partial \phi} + \frac{v_\phi v_R}{R} &= -\frac{1}{R} \frac{\partial \Phi}{\partial \phi} - \frac{1}{\Sigma_d R} \frac{\partial p}{\partial \phi}, \end{aligned} \quad (6.34)$$

where we have replaced the volume density ρ by the disk surface density Σ_d since we are dealing with a two-dimensional disk. We choose a simple equation of state, namely

$$p = K \Sigma_d^\gamma. \quad (6.35)$$

With this equation of state, sound waves in a disk with surface density Σ_0 propagate at a speed v_s given by equation (F.50),

$$v_s^2(\Sigma_0) = \left(\frac{dp}{d\Sigma} \right)_{\Sigma_0} = \gamma K \Sigma_0^{\gamma-1}. \quad (6.36)$$

The equations of motion (6.34) are simplified if we replace p by the specific enthalpy (see eq. F.29)

$$h = \frac{\gamma}{\gamma-1} K \Sigma_d^{\gamma-1}. \quad (6.37)$$

For example, the right side of the first of equations (6.34) then becomes

$$-\frac{\partial \Phi}{\partial R} - \frac{1}{\Sigma_d} \frac{\partial p}{\partial R} = -\frac{\partial \Phi}{\partial R} - \gamma K \Sigma_d^{\gamma-2} \frac{\partial \Sigma_d}{\partial R} = -\frac{\partial}{\partial R}(\Phi + h), \quad (6.38)$$

with a similar simplification in the second of equations (6.34).

We now assume that the spiral wave is only a small perturbation on a steady-state axisymmetric disk, so we can linearize the equations of motion. Denoting quantities in the unperturbed disk by the subscript "0" we have $v_{R0} = 0$ and $\partial \Phi_0 / \partial \phi = \partial p_0 / \partial \phi = 0$. Euler's equations for the unperturbed disk become simply

$$\frac{v_{\phi 0}^2}{R} = \frac{d}{dR}(\Phi_0 + h_0) = \frac{d\Phi_0}{dR} + v_s^2 \frac{d}{dR} \ln \Sigma_0, \quad (6.39)$$

which equates the centripetal acceleration $v_{\phi 0}^2/R$ on the left to the gravitational and pressure force per unit mass on the right. In the cases of interest to us, the sound speed v_s is much smaller than the rotation speed $v_{\phi 0}$ —for example, in the Galactic interstellar gas $v_s \simeq 10 \text{ km s}^{-1}$ while $v_{\phi 0} \simeq 220 \text{ km s}^{-1}$ —and thus the term $v_s^2 d \ln(\Sigma_0)/dR$ in equation (6.39) can be neglected. Hence

$$v_{\phi 0} \simeq \sqrt{R \frac{d\Phi_0}{dR}} = R\Omega(R), \quad (6.40)$$

where $\Omega(R)$ is the circular frequency.

We now write $v_R = \epsilon v_{R1}$, $v_\phi = v_{\phi 0} + \epsilon v_{\phi 1}$, $h = h_0 + \epsilon h_1$, $\Sigma_d = \Sigma_0 + \epsilon \Sigma_{d1}$, $\Phi = \Phi_0 + \epsilon \Phi_1$, where $\epsilon \ll 1$ and the quantities with subscript 1 are of the same order of magnitude as the quantities with subscript 0. Keeping only terms that are first order in ϵ , equations (6.34) yield

$$\begin{aligned} \frac{\partial v_{R1}}{\partial t} + \Omega \frac{\partial v_{R1}}{\partial \phi} - 2\Omega v_{\phi 1} &= -\frac{\partial}{\partial R}(\Phi_1 + h_1), \\ \frac{\partial v_{\phi 1}}{\partial t} + \left[\frac{d(\Omega R)}{dR} + \Omega \right] v_{R1} + \Omega \frac{\partial v_{\phi 1}}{\partial \phi} &= -\frac{1}{R} \frac{\partial}{\partial \phi}(\Phi_1 + h_1). \end{aligned} \quad (6.41)$$

The square bracket in the second equation may be written $-2B(R)$, where $B(R)$ is defined in equation (3.83) and is related to the epicycle frequency $\kappa(R)$ by equation (3.84).

Any solution of equations (6.41) is a sum of terms of the form

$$\begin{aligned} v_{R1} &= \text{Re}[v_{Ra}(R)e^{i(m\phi - \omega t)}] \quad ; \quad v_{\phi 1} = \text{Re}[v_{\phi a}(R)e^{i(m\phi - \omega t)}], \\ \Phi_1 &= \text{Re}[\Phi_a(R)e^{i(m\phi - \omega t)}] \quad ; \quad h_1 = \text{Re}[h_a(R)e^{i(m\phi - \omega t)}], \\ \Sigma_{d1} &= \text{Re}[\Sigma_{da}(R)e^{i(m\phi - \omega t)}]. \end{aligned} \quad (6.42)$$

where $m \geq 0$ is an integer, and the perturbation has m -fold rotational symmetry. Substituting these definitions into equations (6.41) and solving for v_{Ra} and $v_{\phi a}$ we find

$$\begin{aligned} v_{Ra}(R) &= \frac{i}{\Delta} \left[(\omega - m\Omega) \frac{d}{dR}(\Phi_a + h_a) - \frac{2m\Omega}{R}(\Phi_a + h_a) \right], \\ v_{\phi a}(R) &= -\frac{1}{\Delta} \left[2B \frac{d}{dR}(\Phi_a + h_a) + \frac{m(\omega - m\Omega)}{R}(\Phi_a + h_a) \right], \end{aligned} \quad (6.43)$$

where

$$\Delta \equiv \kappa^2 - (\omega - m\Omega)^2, \quad (6.44)$$

and κ , Ω , Δ are all functions of radius, as are Φ_a and h_a . By equations (6.36) and (6.37), the linearized version of the equation of state is

$$h_a = \gamma K \Sigma_0^{\gamma-2} \Sigma_{da} = v_s^2 \Sigma_{da} / \Sigma_0, \quad (6.45)$$

where the sound speed $v_s(\Sigma_0)$ is a function of the unperturbed density.

If ω is real, there may be radii at which $\Delta = 0$ and equations (6.43) diverge. The origin of this singularity can be clarified by writing the exponent in equations (6.42) as $i(m\phi - \omega t) = im(\phi - \Omega_p t)$, where $\Omega_p = \omega/m$ is the pattern speed. Thus the singularity in the perturbed velocities arises when

$$\Omega_p = \Omega \pm \frac{\kappa}{m}, \quad (6.46)$$

which we recognize as the condition for a Lindblad resonance, or a kinematic density wave with $n = \pm 1$ (eq. 6.24).⁵ The singularity arises because non-axisymmetric disturbances can be sustained at these resonances even when the forcing functions in square brackets in equations (6.43) are zero. Our analysis breaks down near these resonances and a separate, more careful treatment is required (Goldreich & Tremaine 1979).

The perturbed surface density is related to the perturbed velocities by the equation of continuity (F.3). Keeping only terms linear in the small quantity ϵ , we have

$$\frac{\partial \Sigma_{d1}}{\partial t} + \nabla \cdot (\Sigma_{d1} \mathbf{v}_0) + \nabla \cdot (\Sigma_0 \mathbf{v}_1) = 0, \quad (6.47)$$

which we write in cylindrical coordinates using equation (B.47):

$$\frac{\partial \Sigma_{d1}}{\partial t} + \Omega \frac{\partial \Sigma_{d1}}{\partial \phi} + \frac{1}{R} \frac{\partial}{\partial R} (R v_{R1} \Sigma_0) + \frac{\Sigma_0}{R} \frac{\partial v_{\phi 1}}{\partial \phi} = 0. \quad (6.48)$$

With equations (6.42) this becomes

$$-i(\omega - m\Omega)\Sigma_{da} + \frac{1}{R} \frac{d}{dR} (R v_{Ra} \Sigma_0) + \frac{im\Sigma_0}{R} v_{\phi a} = 0. \quad (6.49)$$

Equations (6.43), (6.45), and (6.49) provide four constraints on the five variables Σ_{da} , v_{Ra} , $v_{\phi a}$, h_a , and Φ_a . Thus, they determine the dynamical *response* Σ_{da} of the disk to an *imposed* potential Φ_a . If this potential is generated through Poisson's equation by a surface density $\Sigma_1 = \text{Re}[\Sigma_a(R)e^{i(m\phi - \omega t)}]$, then we may formally write

$$\Sigma_{da}(R) = \int dR' \tilde{P}_m(R, R', \omega) \Sigma_a(R'), \quad (6.50)$$

where $\tilde{P}_m(R, R', \omega)$ is the polarization function that relates the response density $\Sigma_{da}(R)$ to the total density $\Sigma_a(R)$ (cf. eq. 5.7). A self-consistent density wave requires that $\Sigma_a = \Sigma_{da}$. Thus equations (6.43), (6.45), (6.49),

⁵ Resonances with $|n| > 1$ are unimportant in fluid disks because they correspond to orbits that are self-intersecting, as in Figure 6.10b, and hence cannot be present in a fluid.

and Poisson's equation can be solved numerically to yield the shapes and frequencies of the modes in a given disk (Hunter 1965; Bardeen 1975; Aoki et al. 1979).

In this section we concentrate on a simpler task: we use the WKB approximation to obtain analytic *local* solutions for density waves. As usual, the potential of a tightly wound wave can be written in the form

$$\Phi_a(R) = F(R)e^{if(R)} = F(R)e^{i\int^R k dR}, \quad (6.51)$$

where $k = df(R)/dR$ and $|kR| \gg 1$. The potential and surface density are related by Poisson's equation (6.31), which holds with fractional error $O(|kR|^{-1})$. The disk response Σ_{da} and—through equation (6.45)— h_a share with Φ_a the factor $\exp[if(R)]$, which varies rapidly with radius. Hence, in equations (6.43), the terms proportional to $(\Phi_a + h_a)/R$ are smaller than those involving $d(\Phi_a + h_a)/dR$ by a factor $\sim kR$, and we can neglect them without degrading the accuracy of our work. We may also write $d(\Phi_a + h_a)/dR = ik(\Phi_a + h_a)$ to the same level of accuracy. Thus equations (6.43) simplify to

$$v_{Ra} = -\frac{(\omega - m\Omega)}{\Delta}k(\Phi_a + h_a) \quad ; \quad v_{\phi a} = -\frac{2iB}{\Delta}k(\Phi_a + h_a). \quad (6.52)$$

Similarly, in equation (6.49) we replace $d(R\Sigma_0 v_{Ra})/dR$ by $ikR\Sigma_0 v_{Ra}$. Since equations (6.52) show that v_{Ra} and $v_{\phi a}$ are of the same order, the second term dominates over the third in equation (6.49) by $O(|kR|)$, and hence we drop the latter term. The continuity equation thus has the form

$$-(\omega - m\Omega)\Sigma_{da} + k\Sigma_0 v_{Ra} = 0. \quad (6.53)$$

We eliminate v_{Ra} using the first of equations (6.52), eliminate h_a using (6.45), and eliminate Φ_a using (6.30). We find

$$\begin{aligned} \Sigma_{da}(R) &= \tilde{P}_m(k, R, \omega)\Sigma_a(R), \quad \text{where} \\ \tilde{P}_m(k, R, \omega) &= \frac{2\pi G\Sigma_0|k|}{\kappa^2 - (\omega - m\Omega)^2 + v_s^2 k^2} \end{aligned} \quad (6.54)$$

is the polarization function for tightly wound density waves. The tight-winding approximation has allowed us to reduce the integral equation (6.50) to a simple analytic expression for the polarization function.

For self-consistent density waves, $\tilde{P}_m(k, R, \omega)$ must be unity, yielding the dispersion relation for a fluid disk in the tight-winding limit:

$$(\omega - m\Omega)^2 = \kappa^2 - 2\pi G\Sigma|k| + v_s^2 k^2, \quad (6.55)$$

where for brevity we have dropped the subscript in the equilibrium surface density Σ_0 . In the case of $m = 0$ waves in a disk with uniform rotation

(constant angular speed, for which $\kappa = 2\Omega$), equation (6.55) reduces to the dispersion relation for a rotating sheet, equation (5.164).

We have derived the dispersion relation (6.55) by dropping all terms in Euler's equation, the continuity equation, and Poisson's equation that are smaller than the dominant terms by $O(|kR|^{-1})$ or more. At the cost of more algebra, we can drop only those terms that are smaller than the dominant terms by $O(|kR|^{-2})$ or more (this requires using the more accurate version 6.33 of Poisson's equation). This procedure leads to the additional constraint

$$\frac{d}{dR} \left[R |\Phi_a|^2 \left(\frac{v_s^2 |k|}{\pi G \Sigma} - 1 \right) \right] = 0. \quad (6.56)$$

In Appendix J we show that this constraint has a simple physical interpretation. As we have seen in §6.1.5, spiral waves transport angular momentum. The rate at which angular momentum is transported outward through radius R is the **angular-momentum current**, which is composed of the gravitational current $C_G(R)$ due to torques exerted by the inner disk on the outer disk (eqs. 6.21 or J.23), and the advective current $C_A(R)$ (eq. J.22). Equation (6.56) is simply the condition that the total angular-momentum current $C_G(R) + C_A(R)$ carried by the density wave is conserved (eq. J.24), as must be true in a steady state.

Before discussing the implications of the WKB dispersion relation (6.55), we shall derive the analogous result for collisionless stellar-dynamical disks.

(d) The dispersion relation for stellar disks The WKB dispersion relation for a disk of stars may be calculated by the same principles that we used to obtain the dispersion relation for a fluid disk: we use the equations of motion to calculate the surface-density perturbation Σ_{d1} arising from a potential perturbation Φ_1 of the form (6.42), and then require that Σ_{d1} and Φ_1 be related by Poisson's equation. The hardest step in this calculation is determining the perturbation \bar{v}_{R1} in the mean radial velocity of the stars that is induced by Φ_1 at a given point (R, ϕ) . If the disk were quite cold, that is, if the unperturbed orbits were all circular, we could obtain \bar{v}_{R1} from equation (6.52) with $h_a = 0$, because a cold stellar disk is dynamically equivalent to a fluid disk with zero pressure. Thus for a cold stellar disk,

$$\bar{v}_{Ra} = -\frac{\omega - m\Omega}{\Delta} k \Phi_a, \quad (6.57)$$

where Δ is defined by equation (6.44).

This expression is accurate if the disk is cool enough that the typical epicycle amplitude is much smaller than the wavelength $2\pi/k$ of the imposed spiral pattern: if this condition is not fulfilled, stars passing through a given location (R, ϕ) at a given time, which come from a range of radii of width equal to twice the epicycle amplitude, will have sampled entirely different parts of the spiral potential. Consequently, the effects of the spiral potential

on the mean velocity perturbation will partially cancel. Formally, we can allow for this cancellation by rewriting equation (6.57) in the form

$$\bar{v}_{Ra} = -\frac{\omega - m\Omega}{\Delta} k \Phi_a \mathcal{F}, \quad (6.58)$$

where $\mathcal{F} \leq 1$ is the **reduction factor**, the factor by which the response of the disk to a given spiral perturbation is reduced below the value for a cold disk. We temporarily defer the task of computing \mathcal{F} .

Once we have \bar{v}_{Ra} , it is straightforward to calculate the response density Σ_{da} , since the Jeans equation (4.204) is identical to the continuity equation of the fluid disk. Thus the analog of equation (6.53) is

$$-(\omega - m\Omega)\Sigma_{da} + k\Sigma_0\bar{v}_{Ra} = 0. \quad (6.59)$$

The final step is to eliminate \bar{v}_{Ra} between equations (6.58) and (6.59), and combine the resulting equation with the WKB form of Poisson's equation (6.30) to obtain the polarization function for stellar disks (cf. eq. 6.54),

$$\tilde{P}_m(k, R, \omega) = \frac{2\pi G \Sigma_0 |k| \mathcal{F}}{\kappa^2 - (\omega - m\Omega)^2}. \quad (6.60)$$

Self-consistent density waves require that $\tilde{P}_m(k, R, \omega) = 1$, which yields the dispersion relation for a stellar disk in the tight-winding limit,

$$(\omega - m\Omega)^2 = \kappa^2 - 2\pi G \Sigma |k| \mathcal{F}, \quad (6.61)$$

the analog of equation (6.55) for a fluid disk.

In Appendix K we evaluate the reduction factor \mathcal{F} for a razor-thin disk having the Schwarzschild DF discussed in §4.4.3, which has the form

$$f_0(R, v_R, v_\phi) = \frac{\gamma \Sigma(R)}{2\pi \sigma_R^2(R)} \exp \left[-\frac{v_R^2 + \gamma^2 \tilde{v}_\phi^2}{2\sigma_R^2(R)} \right], \quad (6.62)$$

where $\Sigma(R)$ is the surface density, $\tilde{v}_\phi = v_\phi - v_c(R)$, $v_c(R) = R\Omega(R)$ is the circular speed, $\sigma_R(R)$ is the radial velocity dispersion, and $\gamma(R) = 2\Omega(R)/\kappa(R)$. As required by the epicycle approximation, we assume that $\sigma_R \ll v_c$. For this DF, the reduction factor can be written in the form (eqs. K.25 and K.21)

$$\begin{aligned} \mathcal{F} \left(\frac{\omega - m\Omega}{\kappa}, \frac{\sigma_R^2 k^2}{\kappa^2} \right) &\equiv \mathcal{F}(s, \chi) = \frac{2}{\chi} (1 - s^2) e^{-\chi} \sum_{n=1}^{\infty} \frac{I_n(\chi)}{1 - s^2/n^2} \\ &= \frac{1 - s^2}{\sin \pi s} \int_0^\pi d\tau e^{-\chi(1 + \cos \tau)} \sin s\tau \sin \tau. \end{aligned} \quad (6.63)$$

Here $I_n(\chi)$ is a modified Bessel function (Appendix C.7). These formulae were derived independently by Lin & Shu (1966) and Kalnajs (1965). Note that the definition of \mathcal{F} implies that $\mathcal{F}(s, 0) = 1$.

The complicated arguments of the reduction factor have a simple physical interpretation. The first, $(\omega - m\Omega)/\kappa$, is simply the ‘‘Doppler-shifted’’ forcing frequency ω as viewed from a star orbiting at the circular frequency Ω , divided by the radial frequency κ . The second is proportional to the square of the ratio of the typical epicycle size σ_R/κ to the spiral wavelength $2\pi/|k|$; we expect that the response will be small when this term is large.

The dispersion relations (6.55) and (6.61) are the key equations in the analytic study of density waves in disks. As we have seen, the conditions for the validity of the WKB approximation are satisfied by only a modest margin for the spiral structure in typical galactic disks. Nevertheless, when properly buttressed by numerical work, these dispersion relations provide an invaluable guide to the behavior of density waves in galaxies. Like any local dispersion relations, they establish the relation between wavenumber and frequency that is satisfied by a traveling wave as it propagates across the disk. They do *not* show that a permanent standing wave pattern can be set up in the disk—this requires more input physics, including the boundary conditions at the center and outer edge of the disk, and an understanding of how the wave behaves at the Lindblad and corotation resonances, where the local dispersion relations break down.

6.2.3 Local stability of differentially rotating disks

The dispersion relations (6.55) and (6.61) for fluid and stellar disks can be used to determine whether a disk is locally stable to axisymmetric perturbations.⁶ All of the analysis we have done so far is for tightly wound non-axisymmetric disturbances, that is, for $|kR/m| \gg 1$. However, it is easy to see that the dispersion relations (6.55) and (6.61) also hold for axisymmetric disturbances ($m = 0$) so long as $|kR| \gg 1$.

Consider first the case of a cold disk. A cold fluid disk has $v_s = 0$, so for axisymmetric disturbances, equation (6.55) becomes

$$\omega^2 = \kappa^2 - 2\pi G\Sigma|k|. \quad (6.64)$$

For a cold stellar disk, $\sigma_R = 0$, and since $\mathcal{F}(s, 0) = 1$ the dispersion relation (6.61) also reduces to equation (6.64)—as expected, since a cold fluid disk is equivalent to a cold stellar disk.

⁶ Local stability to non-axisymmetric perturbations is more complicated. Goldreich & Lynden-Bell (1965a) and Julian & Toomre (1966) show that non-axisymmetric disturbances in fluid and stellar disks always wind up, just like the kinematic density waves of §6.2.1a, and hence appear to be stable. However, unstable non-axisymmetric modes can occur if there is feedback from trailing to leading waves, as we describe in §6.3.2.

Since the quantities on the right side of equation (6.64) are real, ω^2 must also be real. If $\omega^2 > 0$, then ω is real and the disk is stable. If, on the other hand, $\omega^2 < 0$, say $\omega^2 = -p^2$, then $\omega = \pm ip$, and $\exp(-i\omega t) = \exp(\pm pt)$. Hence for $\omega^2 < 0$, there is a perturbation whose amplitude grows exponentially, and the disk is unstable. Thus all perturbations with wavenumber $|k| < k_{\text{crit}}$ or wavelength $\lambda > \lambda_{\text{crit}}$ are unstable, where

$$k_{\text{crit}} \equiv \frac{\kappa^2}{2\pi G\Sigma} \quad ; \quad \lambda_{\text{crit}} \equiv \frac{2\pi}{k_{\text{crit}}} = \frac{4\pi^2 G\Sigma}{\kappa^2}. \quad (6.65)$$

Moreover, the instability is a violent one: as the wavelength of the disturbance shrinks to zero, the growth rate $p = (4\pi^2 G\Sigma/\lambda - \kappa^2)^{1/2}$ grows without limit—a cold, zero-thickness disk disintegrates on small scales in an arbitrarily short time!

Next consider a fluid disk with non-zero sound speed. For axisymmetric disturbances, equation (6.55) reads

$$\omega^2 = \kappa^2 - 2\pi G\Sigma|k| + v_s^2 k^2. \quad (6.66)$$

Once again, the disk is unstable if and only if $\omega^2 < 0$, and the line of neutral stability is

$$\kappa^2 - 2\pi G\Sigma|k| + v_s^2 k^2 = 0. \quad (6.67)$$

The fluid disk is stable if there is no solution of equation (6.67) for positive $|k|$. Since the equation is quadratic in $|k|$, it is easily solved, and we find that axisymmetric stability requires

$$Q \equiv \frac{v_s \kappa}{\pi G\Sigma} > 1 \quad (\text{for fluids}). \quad (6.68)$$

The line of neutral stability defined by equation (6.67) is drawn in Figure 6.13 in terms of the dimensionless ratios Q and $\lambda/\lambda_{\text{crit}}$.

The stability criterion (5.166) for a uniformly rotating sheet is a special case of (6.68) when $\kappa = 2\Omega$. Of course, in a general disk v_s , κ , and Σ are all functions of radius, so Q is also a function of radius. In this case $Q(R) < 1$ implies only *local* axisymmetric instability near radius R , in the sense that a short-wavelength traveling wave that crosses a region with $Q(R) < 1$ will be amplified while it is in this region.

The analysis of the stability of a stellar disk is similar. By analogy we expect that the boundary between stable and unstable axisymmetric waves is given by $\omega = 0$, just as in the case of a fluid disk. Thus from equation (6.61) the stability boundary is given by

$$\kappa^2 = 2\pi G\Sigma|k|\mathcal{F}(0, \sigma_R^2 k^2/\kappa^2), \quad (6.69)$$

or, using the first of equations (6.63) and the identity (C.67),

$$\frac{|k|\sigma_R^2}{2\pi G\Sigma} = \left[1 - e^{-\sigma_R^2 k^2/\kappa^2} I_0 \left(\frac{\sigma_R^2 k^2}{\kappa^2} \right) \right], \quad (6.70)$$

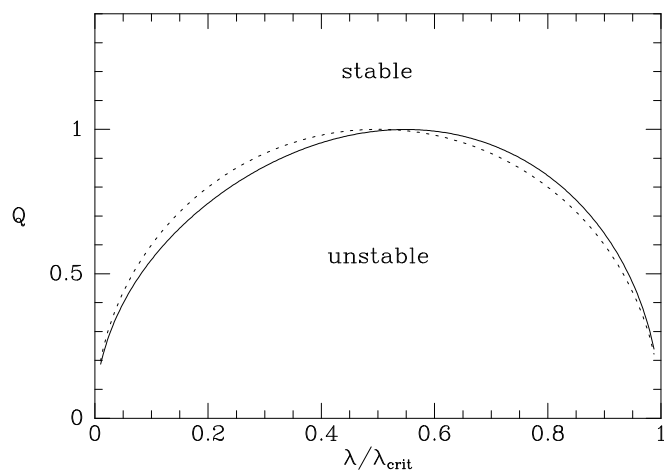


Figure 6.13 Neutral stability curves for tightly wound axisymmetric perturbations in a fluid disk (dashed line, from eq. 6.67) and a stellar disk (solid line, from eq. 6.70).

a relation first derived by Toomre (1964).

There is no solution to (6.70), and thus the stellar disk is stable, if

$$Q \equiv \frac{\sigma_R \kappa}{3.36 G \Sigma} > 1 \quad (\text{for stars}). \quad (6.71)$$

The stability boundary (6.70) is plotted in Figure 6.13 as a function of the dimensionless ratios Q and $\lambda/\lambda_{\text{crit}}$. Note the close analogy between fluid and stellar disks: the dashed (fluid) and solid (stellar) stability curves in Figure 6.13 almost coincide, and the stability criterion for stellar disks (6.71) is obtained from the criterion for fluid disks (6.68) simply by replacing the sound speed v_s by the radial velocity dispersion σ_R , and the coefficient $\pi \simeq 3.14$ by 3.36. The inequality (6.68) or (6.71) is known as **Toomre’s stability criterion**;⁷ its physical interpretation is discussed in the context of the uniformly rotating sheet in §5.6.1. Toomre’s Q can be thought of as a temperature scale for galactic disks. “Hot” disks have large velocity dispersion and high Q , while “cool” disks have low dispersion and Q , and “cold” disks have zero dispersion and $Q = 0$.

As Q drops below unity, instability first appears at a single wavelength, which we may write as $\lambda(\text{most unstable}) \equiv p\lambda_{\text{crit}}$, where the constant p is 0.5 or 0.55 for zero-thickness fluid or stellar disks, respectively. Toomre’s stability criterion is reliable only if $\lambda(\text{most unstable})$ is short compared to the size of the system (since we have used the WKB approximation), and long compared to the thickness of the disk (since we have modeled the disk as razor-thin), but in practice it often works reasonably well somewhat outside the regime in which it is strictly justified.

⁷ An approximate version of equation (6.68) dates back to Safronov (1960).

The most reliable evidence on the value of Q in a real galactic disk comes from the solar neighborhood. The solar neighborhood contains significant amounts of both stars and gas. The total surface density in stars is $\Sigma_\star \simeq (36 \pm 5) \mathcal{M}_\odot \text{pc}^{-2}$, from Table 1.1. From Table 1.2, the epicycle frequency $\kappa_0 \simeq (37 \pm 3) \text{km s}^{-1} \text{kpc}^{-1}$. We adopt the radial velocity dispersion of red main-sequence stars ($B - V > 0.6$), since stars in this color range are mostly old, and dominate the total stellar mass in the solar neighborhood; from Table 1.2 we have $\sigma_R = (38 \pm 2) \text{km s}^{-1}$. With these parameters, the stars alone give $Q_\star = 2.7 \pm 0.4$ (eq. 6.71). The interstellar gas is a complex multi-phase medium but for our purposes we may treat it as a fluid with surface density $\Sigma_g \simeq 13 \mathcal{M}_\odot \text{pc}^{-2}$ (Table 1.1) and sound speed $v_s = 7 \text{km s}^{-1}$. With these parameters, the gas alone has $Q_g = 1.5$ (eq. 6.68). Of course, the perturbations in the stars and gas are coupled through gravity, so the separate analyses we have carried out for the stability of fluid and stellar disks must be combined. Although the surface density of the gas is much smaller than that of the stars, it turns out to have a strong destabilizing effect on the combined system because it is cold ($v_s \ll \sigma_R$). Using parameters similar to the ones here, Rafikov (2001) finds that the solar neighborhood is stable, but not by a large margin—only 15% changes in the surface density or sound speed of the gas component can lead to instability. The most unstable wavelength is $\sim 2 \text{kpc}$, so the WKB approximation is fairly accurate.

There are at least two major uncertainties in this estimate of the local stability of the solar neighborhood. (i) We have oversimplified the interstellar gas by treating it as a single barotropic fluid. (ii) The stellar surface density in spiral galaxies varies with azimuth due to spiral structure, by factors of 1.4–4 (eq. 6.2); we do not understand well how the local stellar surface density is related to the azimuthally averaged surface density at the solar radius, or how such variations affect stability.

6.2.4 Long and short waves

The dispersion relations (6.55) and (6.61) for fluid and stellar disks relate the wavenumber k or wavelength $2\pi/|k|$ to the frequency ω or pattern speed $\Omega_p = \omega/m$. We have plotted these relations in Figure 6.14, assuming for simplicity that Toomre’s Q (eqs. 6.68 or 6.71) is the same at all radii. We use the dimensionless wavenumber k/k_{crit} (eq. 6.65) as the horizontal coordinate, and the dimensionless frequency

$$s \equiv \frac{m(\Omega_p - \Omega)}{\kappa} \quad (6.72)$$

as the vertical coordinate. The corotation resonance is at $s = 0$, and the Lindblad resonances are at $s = \pm 1$. With these coordinates, Figure 6.14 applies to any disk with constant Q —only the relation between s and radius depends on the circular-speed curve. For example, an m -armed spiral in a

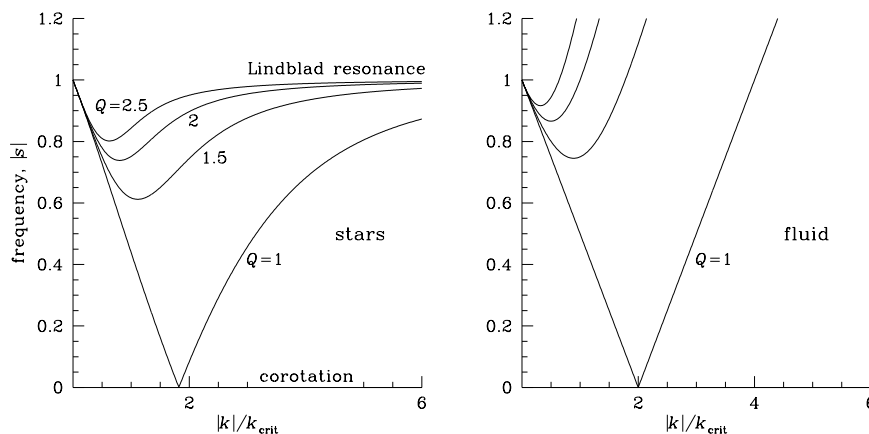


Figure 6.14 The dispersion relation for tightly wound disturbances in stellar (left panel) and fluid (right panel) disks. The horizontal coordinate is the wavenumber k in units of the critical wavenumber k_{crit} (eq. 6.65), and the vertical coordinate is the dimensionless frequency s (eq. 6.72). The curves shown are for $Q = 1, 1.5, 2,$ and 2.5 . Since only $|s|$ and $|k|$ are shown, there is no distinction between leading and trailing waves, or waves inside and outside corotation. Secondary branches of the dispersion relation (not shown) occupy narrow regions near $|s| = 2, 3, \dots$ (Lovelace, Jore, & Haynes 1997).

Mestel disk (constant rotation speed) with corotation radius R_{CR} has $s = m(R/R_{\text{CR}} - 1)/\sqrt{2}$ and the Lindblad resonances are at $R = R_{\text{CR}}(1 \pm \sqrt{2}/m)$.

These figures illustrate the following:

- (i) The dispersion relations for leading ($k < 0$) and trailing ($k > 0$) waves (eq. 6.5) are identical. This is consistent with the anti-spiral theorem of §6.1.4. The dispersion relations are also even functions of s . Thus if $(s, k/k_{\text{crit}})$ is a solution of the dispersion relation, so are $(-s, k/k_{\text{crit}})$, $(-s, -k/k_{\text{crit}})$, and $(s, -k/k_{\text{crit}})$.
- (ii) In disks with $Q > 1$ there is a region around corotation in which the dispersion relation has no real solutions. In this “forbidden” region, tightly wound density waves are **evanescent**, that is, the wave has a complex wavenumber k and thus decays or grows exponentially with radius. The width of the forbidden region increases as Q increases.
- (iii) In the region that is outside the forbidden region but between the Lindblad resonances $s = \pm 1$, there are two branches of the dispersion relation. The **long-wave** branch (the one with larger λ or smaller $|k|$) begins at $k = 0$, $|s| = 1$, and $|s|$ decreases as $|k|$ increases. Sufficiently near the Lindblad resonances the long-wave dispersion relation is independent of Q and is the same for fluid and stellar disks:

$$\frac{|k|}{k_{\text{crit}}} = 1 - s^2. \quad (6.73)$$

The behavior of fluid and stellar disks in this regime is similar because the dispersion relation is determined solely by the self-gravity of the disk, which dominates at large wavelengths, not by the pressure or velocity dispersion, which dominate at short wavelengths.

- (iv) Along the **short-wave** branch (the one with smaller λ or larger $|k|$), $|s|$ increases as $|k|$ increases. In stellar disks, $|k| \rightarrow \infty$ as $|s| \rightarrow 1$, so the short-wave branch terminates at the Lindblad resonance $|s| = 1$. In fluid disks the short-wave branch passes smoothly through the Lindblad resonance, crossing at wavenumber $|k| = 4k_{\text{crit}}/Q^2 = 2\pi G\Sigma/v_s^2$.
- (v) As $Q \rightarrow \infty$, the forbidden region grows to cover the entire region $|s| < 1$ between the Lindblad resonances. Thus in a stellar disk with $Q \gg 1$, a stationary disturbance can be present only near the Lindblad resonances at $s = \pm 1$ —not a surprising result, since we have seen that stationary kinematic waves can exist at the Lindblad resonances in the absence of *any* collective effects. A fluid disk with $Q \gg 1$ sustains density waves only outside the Lindblad resonances, $|s| > 1$, satisfying the dispersion relation $s^2 = 1 + v_s^2 k^2 / \kappa^2$. These are sound waves modified by Coriolis and centrifugal forces, in which the self-gravity of the disk plays no role.

There is a simple physical basis for much of this behavior. The dimensionless frequency s represents the ratio of the forcing frequency seen by a particle, $m(\Omega_p - \Omega)$, to its natural radial frequency, κ . If no perturbing forces are present, a steady wave can be set up only if the two frequencies are equal, which occurs at the Lindblad resonances, $s = \pm 1$. When Q is of order unity, the self-gravity of the disk becomes important. Self-gravity is attractive, and reduces the natural radial frequency below κ , so waves can be present when $|s| < 1$. As the wavenumber $|k|$ increases, self-gravitational forces become more important, so the natural frequency $|s|$ decreases with increasing $|k|$ on the long-wave branch. Eventually the repulsive pressure forces in a fluid disk, or the reduction factor \mathcal{F} in a stellar disk, begin to dominate over self-gravity, so $|s|$ increases again on the short-wave branch. In a fluid disk, the pressure forces can actually increase the natural radial frequency above κ , so waves can also be present when $|s| > 1$. A stellar disk has no pressure forces and therefore cannot support waves with $|s| > 1$.⁸

6.2.5 Group velocity

The waves described by Figure 6.14 are traveling waves, and they propagate with a group velocity (Appendix F.4). The group velocity of a wave packet in a homogeneous dispersive medium is $v_g = d\omega(k)/dk$. Similarly, when the medium is inhomogeneous and the frequency of a wave of given k depends

⁸ Except in narrow regions near the resonances at $|s| = 2, 3, \dots$ which are unlikely to be of much practical importance for galaxy disks (Lovelace, Jore, & Haynes 1997).

on position, $\omega = \omega(k, R)$, the group velocity at radius R is (Whitham 1974)

$$v_g(R) = \frac{\partial \omega(k, R)}{\partial k}, \quad (6.74)$$

so long as the distance over which ω varies is much larger than the wavelength. Toomre (1969) first pointed out that this relation could be applied to Figure 6.14 to determine the evolution of a tightly wound density wave packet. We show in the next subsection that the group velocity also determines the direction and rate of angular momentum and energy transport in the disk.

The dispersion relation for a fluid disk (6.55) yields a group velocity

$$v_g(R) = \text{sgn}(k) \frac{|k|v_s^2 - \pi G\Sigma}{\omega - m\Omega}. \quad (6.75)$$

In this equation v_s , Σ , and Ω are functions of radius determined by the unperturbed disk, the frequency ω is a property of the wave packet, and k is determined from these by the dispersion relation. A wave packet localized around a radius R propagates radially outward if $v_g(R) > 0$, and inward if $v_g(R) < 0$.

The group velocity in stellar disks is derived from (6.74) and the WKB dispersion relation (6.61) (Toomre 1969):

$$v_g(R) = -\frac{\kappa}{k} \frac{1 + 2 \frac{\partial \ln \mathcal{F}(s, \chi)}{\partial \ln \chi}}{\frac{\partial}{\partial s} \ln \frac{\mathcal{F}(s, \chi)}{1 - s^2}}. \quad (6.76)$$

The group velocity can be obtained by a graphical construction from Figure 6.14. The axes are $x \equiv k/k_{\text{crit}}$ and $s = (\omega - m\Omega)/\kappa$ (to within a sign). Thus $d\omega|_R = \kappa ds$ and $dk|_R = k_{\text{crit}} dx$; hence $v_g = (\partial\omega/\partial k)_R = (\kappa/k_{\text{crit}})(ds/dx)$. In other words, to within a sign the group velocity is simply the slope of the curves in Figure 6.14 times the characteristic velocity $\kappa/k_{\text{crit}} = 2\pi G\Sigma/\kappa$. The sign of the group velocity is simply the sign of this slope times the sign of ks .

In the solar neighborhood, $\kappa/k_{\text{crit}} \simeq 36 \text{ km s}^{-1}$ (using $\Sigma = 49 \mathcal{M}_\odot \text{ pc}^{-2}$, $\kappa = 37 \text{ km s}^{-1} \text{ kpc}^{-1}$ from Tables 1.1 and 1.2). It is instructive to estimate the time required for a wave packet to propagate across the Galaxy. Let us take 0.3 as a typical absolute value of the slopes of the curves in Figure 6.14. Then the group velocity is $0.3\kappa/k_{\text{crit}} \simeq 12 \text{ km s}^{-1}$, and the time required to propagate 5 kpc is 400 Myr. For comparison, the rotation period at the solar radius is $2\pi/\Omega \simeq 220 \text{ Myr}$ (Table 1.2). It is evident that any tightly wound wave packet will propagate across a substantial fraction of the galactic disk within at most a few rotation times; more precisely, it will propagate into

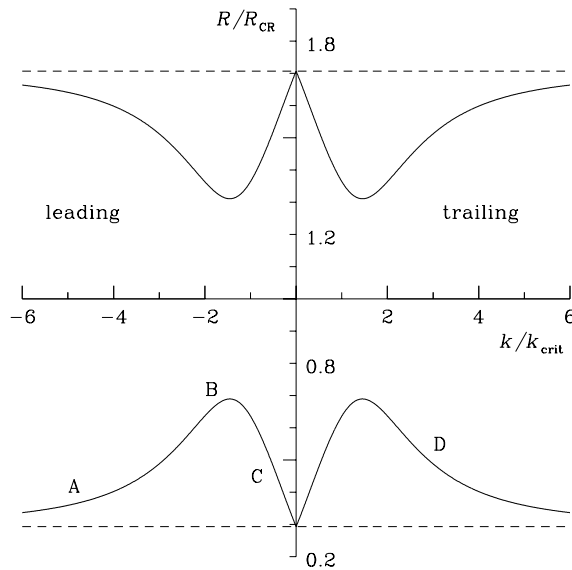


Figure 6.15 Dispersion relation in the form of wavenumber versus radius for a tightly wound $m = 2$ wave in a stellar Mestel disk with $Q = 1.2$. The radial scale is in units of the corotation radius R_{CR} . The inner and outer Lindblad resonances, at $R = 0.293R_{\text{CR}}$ and $R = 1.707R_{\text{CR}}$, are marked by dashed lines.

either a Lindblad resonance or the forbidden region around the corotation resonance. What is the fate of a density wave packet at these resonances?

The life story of a density wave can be deduced with the help of Figure 6.15, which shows the relation between wavenumber and radius in a stellar Mestel disk with constant $Q = 1.2$. Let us begin with a wave packet localized near point A, which can be characterized as a leading ($k < 0$) short-branch wave. The group velocity is positive and the packet propagates outward. As the radius increases, $|k|$ decreases. Eventually, at point B, the edge of the forbidden region, the group velocity changes sign, and the packet begins to propagate inward as a leading wave on the long branch (point C). This reversal in direction can be thought of as a reflection of the packet off the forbidden region, in the same way that wave packets in quantum mechanics reflect off potential barriers.

On the long branch, the tight-winding approximation is suspect, because $|k| \rightarrow 0$ as the group velocity carries the wave packet towards the Lindblad resonance at $k = 0$. The accuracy of the tight-winding approximation for a

given disk can be characterized by the parameter

$$X \equiv \frac{k_{\text{crit}}R}{m} = \frac{\kappa^2 R}{2\pi G \Sigma m}; \quad (6.77)$$

equation (6.7) shows that X is simply the cotangent of the pitch angle for waves of the critical wavenumber k_{crit} . So long as $X \gg 1$ —so long as the tight-winding approximation is satisfied at the critical wavelength—the WKB approximation is valid throughout most of the long branch, and when the wave reaches the Lindblad resonance it simply reflects off this resonance and propagates away on the long trailing branch (Goldreich & Tremaine 1978, 1979). At the forbidden region it is reflected again, this time into a short-branch trailing wave, and then propagates inward (point D). From now on its wavelength becomes shorter and shorter without limit; in a stellar disk the wave is eventually absorbed at the Lindblad resonance, by a process akin to Landau damping (Mark 1974), while in a fluid disk the wave propagates straight through the resonance and on towards the center of the disk. The evolution of waves outside corotation follows a similar sequence. The behavior of waves when $X \lesssim 1$ is described in §6.3.2.

In the evolution that we have described here, the wavenumber k increases monotonically: a tightly wound leading wave becomes first a loosely wound leading wave, then a loosely wound trailing wave, and finally a trailing wave that becomes ever more tightly wound. This behavior is reminiscent of the winding up of material spiral arms. The rate of winding of tightly wound density waves can be worked out quantitatively using the dispersion relation and the group velocity. For the simple example of a fluid, $Q = \text{constant}$ Mestel disk, we show in Problem 6.5 that a trailing wave winds up at a rate

$$\frac{d}{dt}(\cot \alpha) = \Omega_p. \quad (6.78)$$

For comparison, according to equation (6.9), a material arm in a Mestel disk winds up at the rate

$$\frac{d}{dt}(\cot \alpha) = R \left| \frac{d\Omega}{dR} \right| = \Omega, \quad (6.79)$$

where the last equation follows because $\Omega \propto R^{-1}$ in a Mestel disk. Also for comparison, an $m = 2$ kinematic density wave winds up at the rate (eq. 6.26)

$$\frac{d}{dt}(\cot \alpha) = \left| R \frac{d(\Omega - \frac{1}{2}\kappa)}{dR} \right| = \left(1 - \frac{1}{\sqrt{2}} \right) \Omega, \quad (6.80)$$

using the relation $\kappa = \sqrt{2}\Omega$ for a Mestel disk.

In this example, wave packets wind up at a rate comparable to material arms or kinematic density waves: waves outside corotation ($\Omega < \Omega_p$) wind up

faster than material arms, while waves inside corotation ($\Omega > \Omega_p$) wind up slower. Kinematic density waves wind up about 0.3 times as fast as material arms. This conclusion is a serious blow to the Lin–Shu hypothesis: our hope at the outset of this analysis was that non-axisymmetric gravitational forces within the disk could eliminate the winding problem and enable the disk to sustain a long-lived spiral pattern. Instead, we have seen that including these forces actually *increases* the winding rate over much of the disk in comparison to the winding rate for kinematic density waves.

However, the behavior of galactic disks is more complex than this simple example would suggest, mainly because the tight-winding approximation often fails badly for waves on the long branch of the dispersion relation. For example, in the solar neighborhood the parameter X of equation (6.77) equals 4.2 for a two-armed spiral, which is not large enough that the tight-winding approximation is trustworthy on the long branch. As a consequence, we shall find that the simple evolutionary path that was described above—in which short leading waves reflect off the forbidden region into long leading waves, which then reflect off the Lindblad resonance into long trailing waves, etc.—requires major modifications in realistic galactic disks.⁹ To investigate the fate of leading disturbances and their influence on disk stability, we must turn to numerical experiments, as described in §6.3.

6.2.6 Energy and angular momentum in spiral waves

We have seen that tightly wrapped density waves in a rotating fluid or stellar disk obey a WKB dispersion relation and propagate radially at the group velocity $v_g(R)$. Waves in a dispersive medium transport energy and momentum through the medium at the group velocity—we show this for ordinary sound waves in Appendix F.3.1—and thus we expect that tightly wrapped density waves do the same, except that in an axisymmetric system they transport angular momentum rather than linear momentum.

The total angular-momentum current—the rate at which angular momentum is transferred outward across a cylinder of radius R —due to tightly wrapped waves in a fluid disk without viscosity is given by equation (J.24):

$$C_L(R) = \text{sgn}(k) \frac{mR|\Phi_a|^2}{4G} \left(\frac{v_s^2|k|}{\pi G\Sigma} - 1 \right). \quad (6.81)$$

We may define the **angular-momentum density** $L_w(R)$ in the spiral wave such that the wave angular momentum between R and $R + dR$ is $L_w(R)dR$. Then if the current C_L arises from the transport of angular-momentum density at the group velocity, we have

$$L_w(R) = \frac{C_L(R)}{v_g(R)} = \frac{m^2 R(\Omega_p - \Omega)|\Phi_a|^2}{4\pi G^2 \Sigma}. \quad (6.82)$$

⁹ There are, of course, disk systems other than galaxies in which the tight-winding approximation works extremely well. For example, in Saturn’s rings, $X \approx 10^7/m$.

Notice that the wave can have either positive or negative angular-momentum density: waves inside corotation have $L_w < 0$ while those outside corotation have $L_w > 0$.¹⁰

Similarly, the angular-momentum current and angular-momentum density in a stellar disk can be derived from equations (J.25) and (6.76):

$$\begin{aligned} C_L(R) &= -\text{sgn}(k) \frac{mR|\Phi_a|^2}{4G} \left(1 + 2 \frac{\partial \ln \mathcal{F}(s, \chi)}{\partial \ln \chi} \right) \\ L_w(R) &= \frac{C_L(R)}{v_g(R)} = \frac{mR|\Phi_a|^2 |k|}{4G\kappa} \frac{\partial \mathcal{F}(s, \chi)}{\partial s} \frac{1}{1-s^2}. \end{aligned} \quad (6.83)$$

It is straightforward to show from equation (6.63) for $\mathcal{F}(s, \chi)$ that the sign of L_w is the sign of $\Omega_p - \Omega$, so tightly wound waves in a stellar disk, like those in a fluid disk, have positive or negative angular-momentum density according to whether they are outside ($\Omega < \Omega_p$) or inside ($\Omega > \Omega_p$) corotation.

The energy current and energy density of the wave can be determined from C_L and L_w by the following heuristic argument. Imagine that we excite a wave in the disk by the imposition of a potential $\Phi_{\text{ext}}(\mathbf{R}, t)$. The torque per unit mass exerted on the disk is $-\partial\Phi_{\text{ext}}/\partial\phi$ so the rate of change of the disk angular momentum is

$$\dot{L} = - \int d^2\mathbf{R} \Sigma \frac{\partial\Phi_{\text{ext}}}{\partial\phi}. \quad (6.84)$$

The rate of change of the disk energy is (cf. eq. D.10)

$$\dot{E} = \int d^2\mathbf{R} \Sigma \frac{\partial\Phi_{\text{ext}}}{\partial t}. \quad (6.85)$$

To excite a wave with pattern speed Ω_p , we use an external potential that has the same pattern speed, that is, $\Phi_{\text{ext}}(\mathbf{R}, t) = \Phi_{\text{ext}}(R, \phi - \Omega_p t)$. In this case, $\partial\Phi_{\text{ext}}/\partial t = -\Omega_p \partial\Phi_{\text{ext}}/\partial\phi$, so $\dot{E} = \Omega_p \dot{L}$. If this energy and angular momentum all goes into the wave then the wave energy and angular-momentum densities are related by

$$E_w(R) = \Omega_p L_w(R), \quad (6.86)$$

a result that can be confirmed by more rigorous but longer routes. Similarly, the energy current $C_E(R)$ is just Ω_p times the angular-momentum current, since both energy and angular momentum are transported at the group velocity.

¹⁰ A simple device to remember the sign of this result is to think of the wave as having negative angular momentum if its “speed” Ω_p is less than the disk speed Ω .

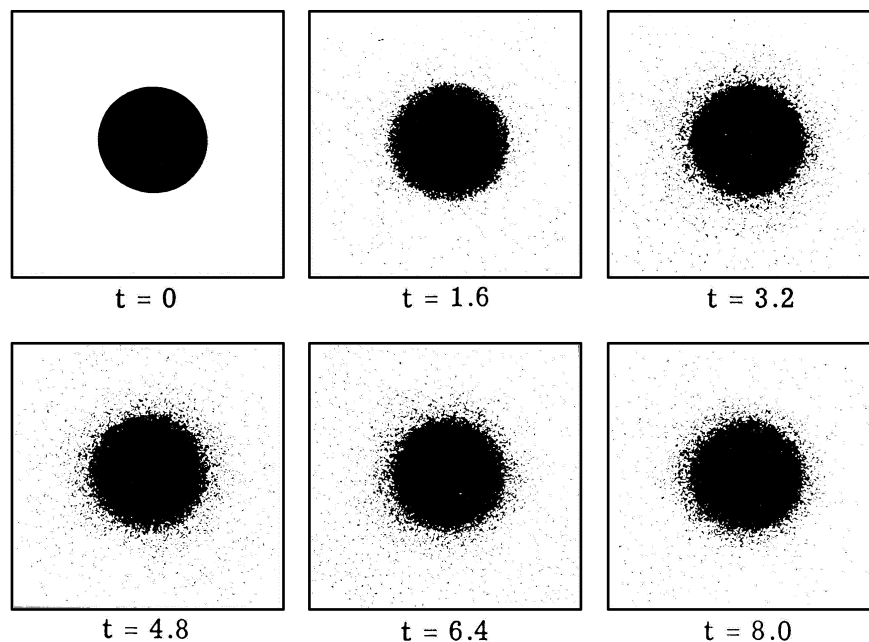


Figure 6.16 Initial evolution of a uniformly rotating disk of 10^5 stars with $Q = 1$ that is constrained to remain axisymmetric. The time unit is the constant rotation period of the stars in the initial disk, $2\pi/\Omega_0$. From Hohl (1971), reproduced by permission of the AAS.

6.3 Global stability of differentially rotating disks

In the last section we described how tightly wound spiral disturbances propagate through galactic disks. Unfortunately, the WKB analysis of tightly wound waves does not give a complete picture of disk dynamics, because it does not apply to loosely wound structures. Hence we now turn to a description of numerical experiments on disk dynamics. We shall find that many of the results of these experiments can be understood by treating the disk as a resonant cavity, within which tightly wound disturbances rattle to and fro.

6.3.1 Numerical work on disk stability

One of the earliest and most influential studies of disk stability was carried out by Hohl (1971), who followed 10^5 bodies using a particle-mesh Poisson solver (§2.9.3) on a 256×256 grid. The initial potential and surface density were those of a Kalnajs disk (eqs. 4.166 and 4.167). The stellar velocities were assigned from the Schwarzschild velocity distribution (4.157), with initial radial velocity dispersion chosen so that the disk was marginally stable in the WKB approximation ($Q = 1$, eq. 6.71).

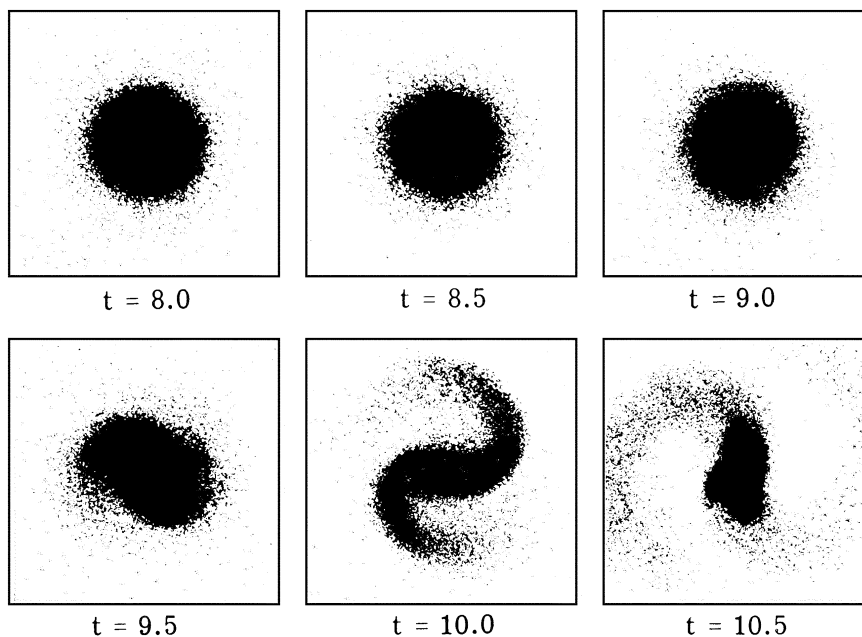


Figure 6.17 Further evolution of the disk in Figure 6.16, after removing the constraint that the disk remain axisymmetric.

These initial conditions do not correspond to an exact stationary solution of the collisionless Boltzmann equation (the DF that *is* a stationary solution for the chosen surface density, the Kalnajs DF of eq. 4.168, is somewhat unrealistic because it has an integrable singularity). To obtain a stationary DF, Hohl simply ran his program for several orbital times while constraining the gravitational field to remain azimuthally symmetric. The resulting evolution is shown in Figure 6.16. Apart from some blurring of the sharp outer edge of the initial distribution, there is rather little change, indicating that the disk has settled into equilibrium. The velocity dispersion of the stars does not rise noticeably, Q remains near unity, and there is no sign of any instability. These results suggest that in this case Toomre's *local* stability criterion $Q > 1$ is also sufficient for *global* stability to axisymmetric modes.

With the disk now in equilibrium, Hohl removed the constraint that the gravitational field should remain axisymmetric. The resulting evolution, shown in Figure 6.17, is dramatically different. In less than two rotations, the disk evolves into a bar-like structure with trailing spiral arms. At later times the bar gradually dissolves, leaving behind a disk with large random velocities ($Q \approx 2-4$) surrounding a slowly rotating oval structure. There is evidently a strong $m = 2$ or **bar instability**, which was not predicted by the local analysis of the previous section.

Hohl's results are consistent with the analytic study of the modes of the Kalnajs disks reported in §5.6.2. We found there that the Kalnajs disks are bar-unstable if $\Omega/\Omega_0 > 0.507$ (eq. 5.189), and in Hohl's axisymmetric disk $\Omega/\Omega_0 = 0.81$, well above the instability threshold.¹¹ The investigations by Kalnajs and Hohl were almost simultaneous and mutually reinforcing: Kalnajs's analytic work showed that the strong bar instability found by Hohl was not a numerical artifact, while Hohl's work showed that the instability was not due to the singularity in the Kalnajs DF.

The bar instability is equally strong in differentially rotating disks, as illustrated by the simulation in Figure 6.18.

By now the stability of a wide range of disk models has been investigated, both by N-body simulations (Zang & Hohl 1978; Sellwood 1981, 1985; Sellwood & Moore 1999) and linear mode calculations using the methods described in §5.3.2 (Zang 1976; Kalnajs 1978; Sawamura 1988; Vauterin & Dejonghe 1996; Pichon & Cannon 1997; Evans & Read 1998).¹² The results largely confirm the two main conclusions of Hohl's classic paper:

- (i) Toomre's local stability criterion $Q > 1$ is a fairly accurate predictor of stability to axisymmetric modes of all wavelengths.
- (ii) If most of the kinetic energy of a disk is in rotational motion, then the disk is usually strongly unstable to a large-scale bar-like mode.

The second conclusion leads to the important question: why are disk galaxies apparently stable? Ostriker & Peebles (1973) were the first to stress the grave consequences of the bar instability for our own Galaxy and other disk galaxies. They argued that the most plausible way to stabilize the Galaxy was to add a massive dark halo, possibly containing even more mass than the visible disk within the solar radius. The dark halo provides part of the equilibrium gravitational field, thereby reducing the required disk mass and the destabilizing effect of the disk's self-gravity. For example, Kalnajs disks embedded in a rigid halo are stable to the bar mode so long as the disk contributes less than $f_d = \frac{2}{3}$ of the equilibrium radial force (Problem 5.18).

The need to stabilize disks against bar formation was one of the influential early arguments for the presence of large amounts of dark matter in galaxies, although we shall see below that there are other ways to stabilize galactic disks, and the actual contribution of dark matter to the total mass inside the solar radius may be as small as 10–20%.

¹¹ Axisymmetric ($m = 0$) instabilities are also present in some of the Kalnajs disks, as seen in Figure 5.5. In rapidly rotating disks these arise because $Q < 1$, while in slowly rotating Kalnajs disks they are probably connected with the integrable singularity in the DF; thus neither instability is expected to be present in Hohl's model.

¹² An important element of these investigations is the comparison of the results from these two methods, which has been done for the Kalnajs disks (Sellwood 1983), the Kuzmin disk (Sellwood & Athanassoula 1986), the isochrone disk (Earn & Sellwood 1995), and the Mestel disk (Sellwood & Evans 2001).

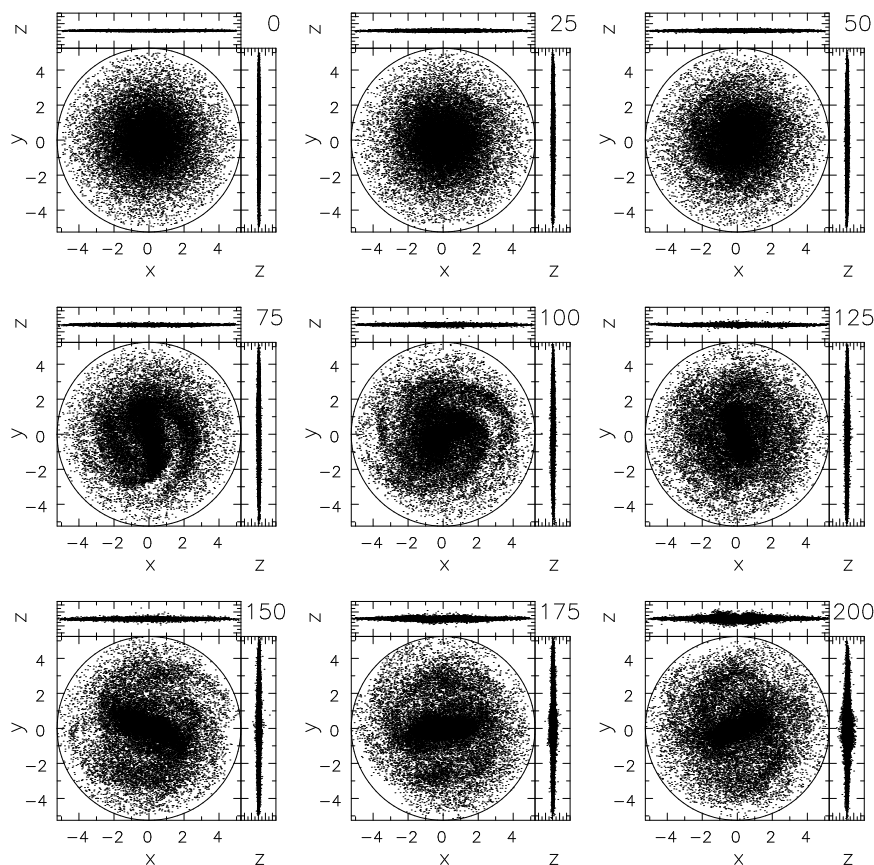


Figure 6.18 The development of the bar instability in a differentially rotating disk. The initial surface-density distribution in the disk is exponential (eq. 2.162) with unit scale length R_d , unit total mass, and $G = 1$. The initial radial velocity dispersion in the disk is chosen so that $Q = 1.2$ (eq. 6.71). The disk is embedded in an approximately spherical dark halo; thus at radius $2R_d$ the disk contributes only $f_d = 0.53$ of the total radial force. The gravitational potential is softened with a softening length $\epsilon = 0.1$ (eq. 2.226). The disk and halo are modeled by 0.5×10^6 and 2.5×10^6 particles, respectively. Each panel shows three orthogonal views of the disk, with a time label at the upper right. Courtesy of J. Sellwood.

6.3.2 Swing amplifier and feedback loops

Many features of the bar instability can be understood by augmenting the WKB dispersion relations with two new physical concepts.

(a) The swing amplifier In §6.2.5 we argued that any leading disturbance in a disk inevitably unwinds. If the parameter X defined by equation (6.77) is of order unity or smaller, the WKB approximation is unable

to follow the evolution of an unwinding leading wave after it reaches the long branch of the dispersion relation. We may use numerical experiments to study what actually happens next.

Following Toomre and Zang (Toomre 1981) we consider a stellar Mestel disk with constant Q . We assume that the disk is embedded in a rigid, fixed “halo” component, such that the disk contributes only a fraction $f_d < 1$ of the total radial force in the equilibrium system. If the disk surface density is written as $\Sigma(R) = v_0^2/(2\pi GR)$ (eq. 4.158), the disk contributes a radial force v_0^2/R , and the circular frequency Ω is therefore given by

$$\Omega^2 = \frac{v_0^2}{f_d R^2}. \quad (6.87)$$

Using the relation $\kappa = \sqrt{2}\Omega$ (see eq. 3.80), we have from equation (6.77),

$$X = \frac{2}{f_d m}. \quad (6.88)$$

Toomre and Zang used linear perturbation theory to follow numerically the evolution of a leading wave packet with $m = 2$ in a disk with $Q = 1.5$, $f_d = \frac{1}{2}$, and $X = 2$. The results are shown in Figure 6.19. As expected, within a few rotation periods the wave unwinds into a relatively open pattern (frame 3) and then into a trailing pattern that becomes more and more tightly wound (frame 9). The striking feature of these calculations is that the amplitude of the trailing wave in frame 9 is about twenty times larger than the amplitude of the initial leading wave in frame 1, and that at intermediate stages (frames 4, 5, and 6) an even stronger *transient* spiral pattern is formed.

These results are a manifestation of **swing amplification**, a phenomenon that is not captured by the WKB approximation.¹³ The following heuristic explanation of the basic mechanism of the amplification is given by Toomre (1981). Consider a material arm described by equation (6.8). For the purposes of this argument we redefine the quadrant of the pitch angle so that $0 < \alpha < 90^\circ$ for trailing arms and $90^\circ < \alpha < 180^\circ$ for leading arms. Thus equation (6.9) is written as

$$\cot \alpha = -Rt \frac{d\Omega}{dR} = 2At, \quad (6.89)$$

where A , defined in equation (3.83), is a measure of the shear in the disk. The rate of change of pitch angle is

$$\frac{d\alpha}{dt} = -\frac{2A}{1 + 4A^2 t^2}. \quad (6.90)$$

¹³ Weaker amplification, up to a factor of two, is present in the WKB approximation if Q is very close to unity. This process has an elegant physical interpretation in terms of tunneling across the forbidden region at corotation (Mark 1976; Goldreich & Tremaine 1978, 1979).

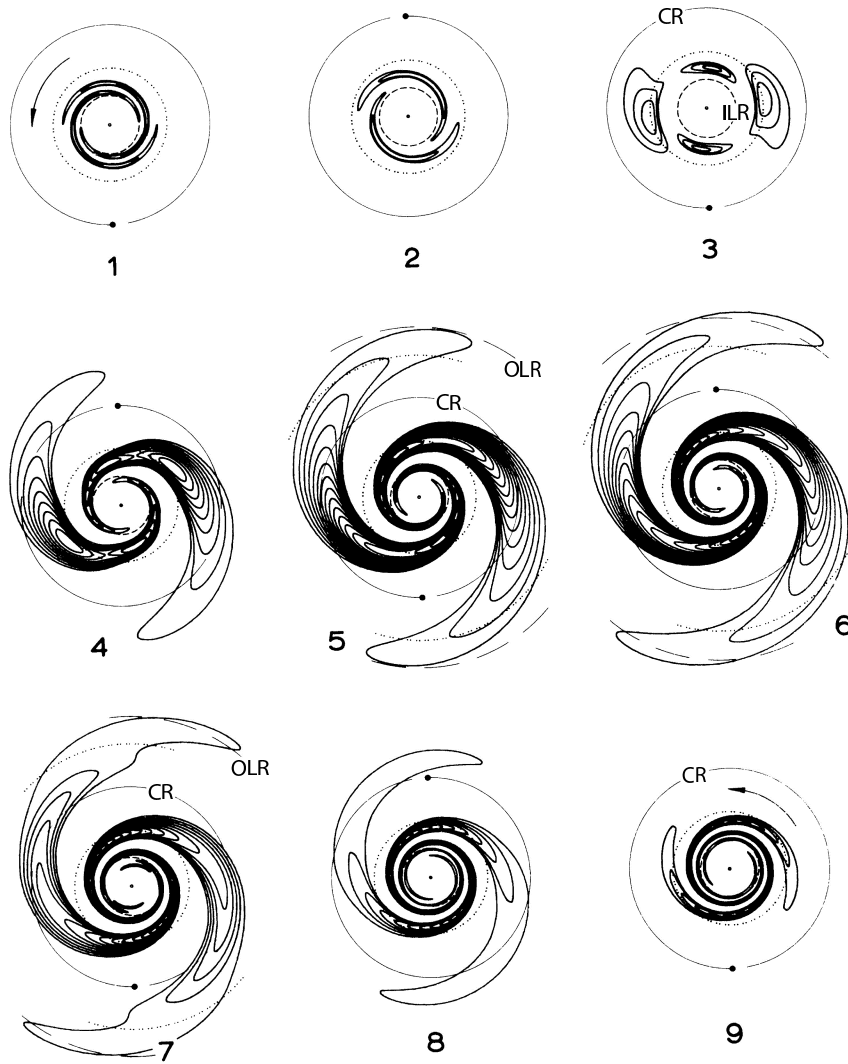


Figure 6.19 Evolution of a packet of leading waves in a Mestel disk with $Q = 1.5$ and $f_d = 1/2$ (equal contributions from the disk and the rigid halo to the flat circular-speed curve). Contours represent fixed fractional excess surface densities; since the calculations are based on linear perturbation theory, the amplitude normalization is arbitrary. Contours in regions of depleted surface density are not shown. The time interval between diagrams is one-half of a rotation period at corotation. ILR, CR, and OLR denote the radii of the inner Lindblad resonance, the corotation resonance, and the outer Lindblad resonance. From Toomre (1981), © Cambridge University Press 1981. Reprinted by permission of Cambridge University Press.

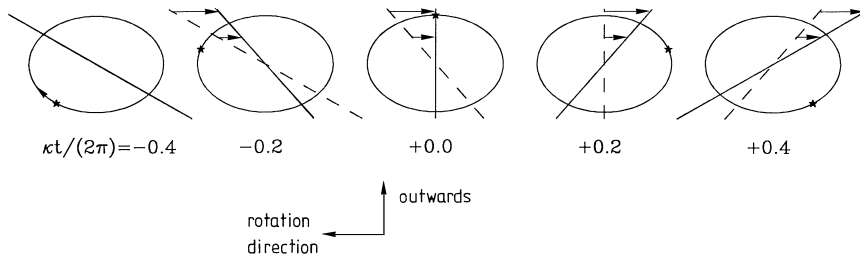


Figure 6.20 Schematic diagram showing the reason for swing amplification. The panels show the motion of a star in its epicycle and the motion of an unwinding material arm in a Mestel disk. The dashed and solid straight lines show the movement of the material arm between panels. The interval between panels is 0.2 times the epicycle period. Note the temporary similarity of angular speeds of the arm and the star: between $\kappa t/(2\pi) \simeq -0.2$ and $+0.2$ the arm and the star swing around at roughly the same rate, so the gravitational field of the arm can steadily attract the star.

When the arm is tightly wound, its rotation rate $d\alpha/dt$ is slow, but as it swings from leading to trailing it reaches a maximum rotation rate of $2A$. This maximum is comparable to the average angular speed of stars around their epicycles, κ (for a Mestel disk, $2A = \Omega$ and $\kappa = \sqrt{2}\Omega$). Moreover, both the unwinding of the arm and the rotation of the stars around their epicycles are in the same sense, opposite to the direction of rotation (see Figure 6.20). Thus there is a temporary near-match between the epicyclic motion and the rotating spiral feature, which enhances the effect of the gravitational force from the spiral on the stellar orbit—and the contribution of the star’s own gravity to the spiral perturbation. This enhancement can lead to rapid growth in the strength of the arm over the interval of about one radian when the arm is most open.

Swing amplification is effective when Q exceeds unity, but not by too much—so the disk is stable, but still responds strongly to gravitational perturbations. Numerical experiments show that the amplification factor during the swing from leading to trailing is extremely sensitive to the value of Q when Q is near unity. For example, in Mestel disks there is an increase of a factor of five in the amplification as Q is decreased from 1.5 to 1.2. Strong amplification also requires that X is not too large—if $X \gg 1$, the wave will propagate away from corotation on the long branch of the dispersion relation before swinging through $\alpha = 90^\circ$, so corotating stars are not coupled to the swinging spiral field. Figure 6.21 shows the gain of the swing amplifier as a function of the parameters X and Q , in disks with a flat circular-speed curve. This figure suggests that $1 \lesssim X \lesssim 3$, $Q \lesssim 2$ are necessary and sufficient conditions for swing amplification factors of more than a few in such disks; and when $X \simeq 2$ and $Q = 1.2$ the amplification factor can easily be

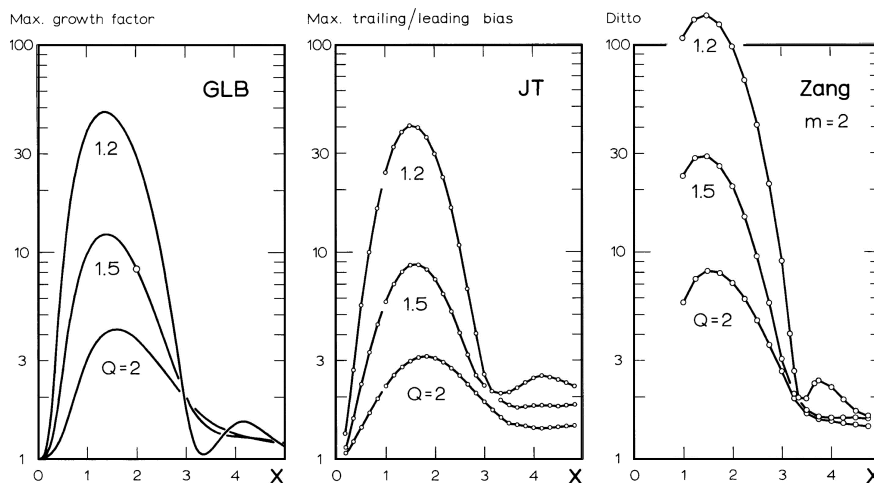


Figure 6.21 Gain of the swing amplifier as a function of X and Q , for a disk with a flat circular-speed curve. The first two panels show calculations based on the sheared sheet (§8.3.2); the left panel is for a fluid (Goldreich & Lynden-Bell 1965a) and the middle panel is for a stellar system (Julian & Toomre 1966). The right panel shows calculations using the Zang (1976) model of a Mestel disk. Note that all three models give similar results, that the gain is small for $X \gtrsim 3$, and that the vertical scale is logarithmic. From Toomre (1981), © Cambridge University Press 1981. Reprinted by permission of Cambridge University Press.

30–100.¹⁴

(b) Feedback loops Swing amplification of a single leading disturbance is not sufficient by itself to destabilize a galactic disk. However—as anyone who has set up a public address system knows—an amplifier together with positive feedback from output to input can give rise to instability. Thus, any mechanism that turns trailing waves into leading waves is liable to initiate instability in a disk with a strong swing amplifier. For example:

- (i) Suppose that the disk has a sharp outer edge that lies outside the forbidden region around corotation but inside the outer Lindblad resonance. Trailing waves that approach this edge can reflect off it, in the same way that waves reflect off the end of a hanging chain or an organ pipe. The reflection reverses the sign of the wavenumber k and hence reflects trailing waves into leading waves with the same wavelength. This is a simple example of feedback, but it is probably not present in most galactic disks because their outer edges are not sufficiently sharp.
- (ii) In some cases the disk may have no inner Lindblad resonance, because the maximum value of $\Omega - \frac{1}{2}\kappa$ is less than the pattern speed. For ex-

¹⁴ The gain of the swing amplifier in disks with other circular-speed curves is discussed by Athanassoula (1984) and Fuchs (2001).

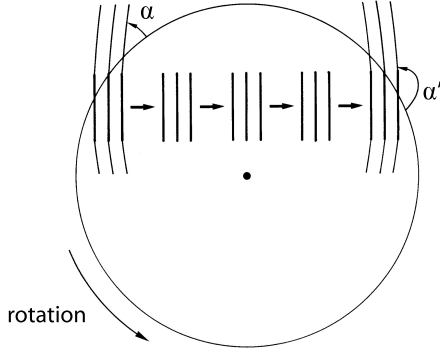


Figure 6.22 A graphical argument that suggests why trailing waves that propagate through the center of a disk emerge as leading waves. A small patch of three incoming trailing waves with pitch angle $\alpha < 90^\circ$ is shown on the left. The patch propagates through the center as a plane wave and emerges with a pitch angle $\alpha' = 180^\circ - \alpha$. Since $\alpha' > 90^\circ$ the emerging wave is leading.

ample, in the isochrone potential (eq. 2.47), there is no inner Lindblad resonance when $\Omega_p > 0.0593(GM/b^3)^{1/2}$ (Figure 6.11a). In this case, trailing waves can propagate right in to the center of the disk, from where they emerge as leading waves propagating outward (see Figure 6.22).

- (iii) Nonlinear interactions of trailing spiral waves can generate leading spirals. To illustrate how this happens, we consider two tightly wrapped trailing spirals with surface density of the form

$$\Sigma_i(R, \phi, t) = H_i(R) \cos[f_i(R) + m_i(\phi - \Omega_{pi}t)] \quad (i = 1, 2). \quad (6.91)$$

In the collisionless Boltzmann equation, the gravitational field due to one of these waves interacts with the perturbed DF of the other to generate a perturbation to the DF whose spatial form is proportional to $\Sigma_1 \Sigma_2$. The rapidly varying component of this perturbation is proportional to $\cos[f_1(R) + m_1(\phi - \Omega_{p1}t)] \cos[f_2(R) + m_2(\phi - \Omega_{p2}t)]$. Using the formula $\cos x \cos y = \frac{1}{2} \cos(x + y) + \frac{1}{2} \cos(x - y)$, it is easy to show that this corresponds to spiral waves with shape function $f_\pm(R) = f_1(R) \pm f_2(R)$, azimuthal wavenumber $m_\pm = m_1 \pm m_2$, and pattern speed $\Omega_{p\pm} = (m_1 \Omega_{p1} \pm m_2 \Omega_{p2}) / (m_1 \pm m_2)$. If the radial wavenumber $k_\pm = df_\pm/dR$ is negative and satisfies the WKB dispersion relation with the corresponding azimuthal wavenumber and pattern speed, then the nonlinear interaction will launch a leading wave that can later be amplified by the swing amplifier (Sygnet et al. 1988; Fuchs, Dettbarn, & Tsuchiya 2005).

(c) Physical interpretation of the bar instability These results yield a simple physical interpretation of the bar instability: any minor leading disturbance unwinds and is then swing amplified into a short trailing disturbance, which propagates through the disk center and emerges as a short leading disturbance, which then unwinds and is amplified further.

One clue that supports this point of view is shown in Figure 6.23. The figure shows an unstable $m = 2$ mode for a disk with surface density $\Sigma(R) = \Sigma_0 \exp(-\frac{1}{2}R^2/R_0^2)$. Only a fraction f_d of the radial force in the

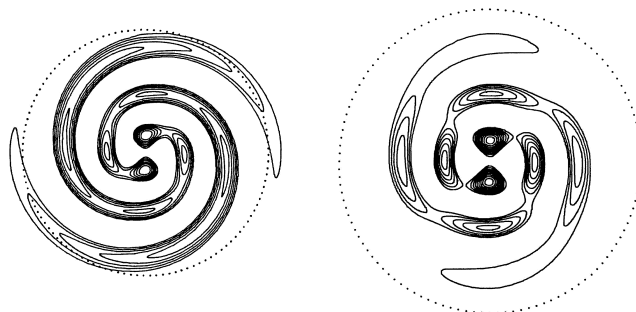


Figure 6.23 Shapes of unstable modes in a Gaussian disk (Toomre 1981). Left: no rigid halo. Right: one-third of the central force arises from a rigid halo. The growth rate is 17% of the pattern speed in the first case but only 3% of the pattern speed in the second. This calculation used a cold disk and softened gravity to model a disk with non-zero Q (see Problem 6.3). Dots mark the corotation circle. From Toomre (1981), © Cambridge University Press 1981. Reprinted by permission of Cambridge University Press.

unperturbed system arises from the disk; the remainder comes from a rigid “halo” component. The fraction $f_d = 1$ in the model on the left (no halo) and $f_d = \frac{2}{3}$ on the right. The details of the disk model are described in Toomre (1981); the main feature to notice here is the “lumpy” structure of the mode, which is much more prominent in the right panel. This lumpy structure is naturally explained as the result of interference between leading and trailing waves of nearly equal amplitude, propagating through the disk center. In the model on the left, which has a large growth rate, the swing-amplified trailing waves have larger amplitude than the leading waves, so the mode looks like a trailing spiral even though a weak leading wave is present. In the model on the right, which has a much lower growth rate because of the rigid halo, the amplitudes of the leading and trailing waves are nearly equal, so the interference pattern is more pronounced. This interpretation also predicts, correctly, that growing modes in a differentially rotating disk are always trailing.

These arguments imply that *any* simple stability criterion based on the global properties of the galaxy cannot accurately predict stability in all cases, because stability depends on the efficiency of the trailing→leading feedback through the center, which in turn depends sensitively on the properties of the disk near the center. A disk can be stabilized by a small readjustment of the inner mass distribution that increases the central angular speed and thus creates an inner Lindblad resonance, cutting off the propagation of density waves.

6.3.3 The maximum-disk hypothesis

A major uncertainty in understanding the dynamics of disk galaxies is the fractional contribution of the halo to the total mass within the outer radius R_* of the stellar disk (of course, the dark halo contributes most of the mass at radii $\gg R_*$). In particular, disks in which the fractional contribution of the dark halo is small are far more susceptible to the bar instability. The circular-speed curve strongly constrains the overall mass distribution, but the relative contributions of the disk and halo to this mass are poorly determined, because the mass-to-light ratio Υ_d of the disk is uncertain.

Some constraints on Υ_d can be obtained by fitting models of stellar populations (see BM §5.4 and Bell & de Jong 2001) to the observed colors or spectra of galaxy disks. However, these models depend on uncertain assumptions about the distribution of stellar ages and masses—in particular, low-mass stars contribute little or no light but could have a substantial contribution to Υ_d .

The circular-speed curves of most spiral galaxies out to $\sim R_*$ can be fitted rather well by simple models in which the mass-to-light ratio of the disk is independent of radius and there is little or no dark mass (Palunas & Williams 2000; Sancisi 2004). Dark matter is *required* only to fit HI rotation curves, which extend to several times R_* . Advocates of the **maximum-disk hypothesis** argue that the success of these simple models would be an improbable coincidence unless the disk contributed most of the mass inside R_* . This argument is appealing but far from rigorous, since (i) the radial profile of the dark matter is unknown and may be similar to that of the disk, and (ii) disks have radial color gradients, which suggest that the mass-to-light ratio of the disk is likely to depend on radius.

It is worthwhile to define the maximum-disk hypothesis more precisely (e.g., Sackett 1997). Assume that the disk and bulge have mass-to-light ratios Υ_d and Υ_b that are independent of radius, so we can derive their contribution to the circular-speed curve once we know these two parameters. Assume that the dark halo has a given functional form such as an NFW profile (§2.2.2g), which depends on two parameters ρ_0 and a . Then find the combination of Υ_d , Υ_b , ρ_0 and a that fits the circular-speed curve while minimizing the halo mass inside R_* ; the resulting disk parameters define the “maximum disk.” The maximum-disk hypothesis is that most galaxies have nearly maximum disks.

Since dark halos are required to fit the circular-speed curves outside R_* , and the density of the halo increases inwards, even maximum disks have *some* halo contribution to the radial force inside R_* . Typically, a maximum disk contributes a fraction $f_d = 0.55$ – 0.9 of the radial force at a radius of $2.2R_d$ (here R_d is the scale length of the exponential disk, eq. 1.7); since the circular speed at a given radius is proportional to the square root of the force, maximum disks contribute 75–95% of the circular speed. For comparison, in the two models of the Milky Way described in §2.7, the disk contributes

$f_d = 0.74$ (Model I) and $f_d = 0.32$ (Model II) of the force at $2.2R_d$. Thus Model I has close to a maximum disk, but Model II does not.

Many indirect arguments have been used to constrain the relative disk and halo contributions to the gravitational field. These include:

- (i) Fitting the two-dimensional velocity field in barred or unbarred spirals to the gravitational field generated by the disk stars and dark halo (Weiner et al. 2001; Kranz, Slyz, & Rix 2003); the basic assumptions here are that the mass-to-light ratio of the stars is independent of position and that the halo, in contrast to the stars, has no small-scale non-axisymmetric structure that could perturb the velocities.
- (ii) Modeling the mass distribution in the Galaxy as arising from stars and a dark halo, where the distribution of stars is inferred from the near-infrared light distribution (Englmaier & Gerhard 1999; Klypin, Zhao, & Somerville 2002). These mass models must explain the kinematics of HI and CO gas, including non-circular motions induced by the Galactic bar, as well as the optical depth to gravitational microlensing in the Galactic bulge (BM §10.2.2, Binney & Evans 2001; Sumi et al. 2003).
- (iii) The pattern speed of a bar embedded in a massive halo is expected to decay rapidly due to dynamical friction, as described in §8.1.1d. Observations show, however, that the pattern speeds of bars are generally high, which limits the halo contribution to the mass in the inner galaxy (see §6.5.2e).
- (iv) The Tully–Fisher law (eq. 1.24) implies that the maximum rotation speed in disk galaxies is tightly correlated with the luminosity of the galaxy. If the halo contribution to the disk rotation is small and the disk mass-to-light ratio Υ_d is independent of radius, then the maximum rotation speed should be $v_c = 0.62(G\Upsilon_d L/R_d)^{1/2}$ for an exponential disk of scale length R_d (see Figure 2.17). Since disks of a given luminosity L have a range in scale lengths, the residuals in v_c from the Tully–Fisher relation should be anti-correlated with the scale length if the maximum-disk hypothesis is correct. No such correlation is observed (Courteau & Rix 1999).
- (v) In disks with a flat circular-speed curve, swing amplification is strongest when the parameter X (eq. 6.77) is approximately 1.5 (Figure 6.21). Given the circular-speed curve, which determines the epicycle frequency $\kappa(R)$, and an assumed disk mass-to-light ratio, which determines the surface density $\Sigma(R)$, we can therefore compute the azimuthal wavenumber m_{\max} that is likely to be subject to the strongest swing amplification at a given radius. In grand-design spirals, where there is usually a prominent two-armed spiral, m_{\max} should be 2, and this requirement constrains the mass-to-light ratio of the disk and hence the halo fraction (Athanasoula, Bosma, & Papaioannou 1987).

These arguments have not yet provided clear and consistent evidence on the dark halo contribution to the mass within the visible disk radius R_* .

Preliminary conclusions are that (i) high-luminosity spirals with large rotation velocities appear to have near-maximum disks, while the dark halo contributes half or more of the total mass within R_* in low-luminosity spirals (Kranz, Slyz, & Rix 2003); (ii) the disk contributes $f_d \simeq 0.5\text{--}0.9$ of the total radial force at the solar radius in the Milky Way, or at $2.2R_d$ in similar spiral galaxies. Correspondingly, the halo contributes a fraction $f_h \equiv 1 - f_d$ of the local force in the range 0.1–0.5.

Maximum disks are difficult to reconcile with the standard cosmological model of halo formation from cold dark matter (page 728), which produces dark halos with strong central concentration (Navarro, Frenk, & White 1997).

6.3.4 Summary

Analyses based on the WKB approximation show that disks are stable to short-wavelength perturbations if Toomre's $Q > 1$ (eqs. 6.68 and 6.71). Numerical experiments suggest that $Q > 1$ implies stability to long-wavelength axisymmetric perturbations as well, but not to non-axisymmetric ones.

In many disks, the strongest non-axisymmetric instability is an $m = 2$ bar-like instability. The instability appears to be due to a feedback loop in which trailing waves propagate through the galactic center, emerging as leading waves that are swing amplified into stronger trailing waves.

The bar instability can be stabilized by the following mechanisms:

- (i) Increase Q . Large random motions reduce the susceptibility of a disk to gravitational instability and thus inhibit the swing amplifier. Figure 6.21 shows that $Q \lesssim 2$ is necessary for strong swing amplification in a disk with a flat circular-speed curve. This mechanism is unlikely to be important for most galaxies, because the measured values of Q in the solar neighborhood (Rafikov 2001) and in other disk galaxies (Bottema 1993) are too small to stabilize the disk.
- (ii) Increase X (eq. 6.77). The presence of a halo or some other hot component that does not respond to the perturbations will lead to more tightly wound waves. For example, Figure 6.21 shows that $X \lesssim 3$ is needed for strong swing amplification in a disk with a flat circular-speed curve. Equation (6.88) shows that in such a disk, an $m = 2$ wave has $X < 3$ if more than $\frac{1}{3}$ of the equilibrium radial force arises from the disk. This is how Ostriker & Peebles (1973) were able to stabilize the disk by adding a massive halo. However, massive halos may provide short-term stability at the cost of long-term instability: resonant interactions with halo stars can drain angular momentum from the bar, leading to slow but steady growth of the bar over Gyr timescales (Athanasoula 2002). Massive halos are probably responsible for stabilizing low-luminosity disk galaxies, since their circular-speed curves imply that most of the mass inside the outer disk radius R_* is dark.
- (iii) Cut off the feedback. The commonest feedback loop involves the propagation of waves through the disk center, which requires that there is

no inner Lindblad resonance. Thus a cool disk can be stabilized by a readjustment of the mass distribution in the inner disk so as to increase the central angular speed and thus create an inner Lindblad resonance (see Sellwood & Moore 1999 and Sellwood & Evans 2001 for examples). This is the mechanism that is likely responsible for the stability of most luminous spiral galaxies; there may be help from mechanism of (ii), but a massive halo is not *required* for stability.

6.4 Damping and excitation of spiral structure

6.4.1 Response of the interstellar gas to a density wave

Only disk galaxies containing substantial quantities of cool interstellar gas exhibit spiral structure. To explain this empirical fact, and to interpret the rich phenomenology of spiral structure as seen in a variety of gas tracers (HI, CO, magnetic fields, dust, etc.), we need to understand how interstellar gas responds to the gravitational forces from a density wave.

In general, calculations of spiral structure in the interstellar gas must be done numerically, since the linear theory we have used so far is invalid for the strong density contrasts seen in the gas. However, with some approximations we can reduce the problem to a simple analog that is easy to understand. The interstellar gas generally contains only a small fraction of the total surface density in a galactic disk—Table 1.1 states that in the solar neighborhood this fraction is about 25%. Hence to a first approximation we can assume that the gas responds to the potential of the stars alone. We write this potential as the sum of the unperturbed axisymmetric potential $\Phi_0(R)$ and the perturbed potential due to a tightly wound stellar density wave,

$$\Phi_1(R, \varphi) = F \cos(kR + m\varphi), \quad (6.92)$$

where φ is now the azimuthal angle in a frame rotating at the pattern speed Ω_p . We shall concentrate our attention on interstellar clouds (see BM §9.6) since these contain the bulk of the atomic and molecular gas. In a first approximation, the clouds can be regarded as test particles moving in the potential of the stellar disk. Their motion can be obtained from the equations of motion for a star in a weak non-axisymmetric field (§3.3.3). Adapting equations (3.142) to this case, we find

$$\begin{aligned} \ddot{R}_1 + \left(\frac{d^2\Phi_0}{dR^2} - \Omega^2 \right) R_1 - 2R_0\Omega\dot{\varphi}_1 &= kF \sin(kR + m\varphi), \\ \ddot{\varphi}_1 + 2\Omega\frac{\dot{R}_1}{R_0} &= \frac{mF}{R_0^2} \sin(kR + m\varphi), \end{aligned} \quad (6.93)$$

where Ω , R_1 , and φ_1 are defined by equations (3.137) and (3.139).

We shall focus on clouds in a small patch of the galactic disk, and we may choose the origin of time so that their unperturbed azimuthal trajectory is $\varphi_0(t) = (\Omega - \Omega_p)t$. Since the perturbing force is weak, we can replace φ on the right side by its unperturbed value $\varphi_0(t)$. However, we do not replace R by its unperturbed value R_0 —the wavenumber k is large because the density wave is tightly wound, and we want to include the possibility that kR_1 is of order unity, even though R_1 is small.

Since the waves are tightly wound, the radial force, which is the right side of the first of equations (6.93), is larger by a factor kR_0/m than the azimuthal force, which is R_0 times the right side of the second of equations (6.93). Hence the right side of the second equation can be dropped, and the remainder of that equation integrates immediately to an approximate statement of the conservation of angular momentum

$$\dot{\varphi}_1 + \frac{2\Omega R_1}{R_0} = \text{constant}. \quad (6.94)$$

Since a readjustment of R_0 simply changes R_1 by a constant, we can always choose R_0 so that the constant in (6.94) is zero; then, eliminating φ_1 from (6.93) and using equation (3.146b), we obtain

$$\ddot{R}_1 + \kappa^2 R_1 = kF \sin[k(R_0 + R_1) + m(\Omega - \Omega_p)t], \quad (6.95)$$

where κ is the epicycle frequency.

There is a simple analog system that obeys equation (6.95). Consider the endless row of identical pendulums shown in Figure 6.24. Each pendulum has length L and therefore its natural oscillation frequency at small amplitude is $\kappa \equiv (g/L)^{1/2}$ where g is the acceleration due to gravity. Let the horizontal position of the support at the top of a given pendulum be x_0 , and the position of the bob at the bottom be $x_0 + x_1(x_0, t)$. Finally, suppose that each pendulum is subjected to a horizontal force per unit mass $kF \sin(kx + \omega t)$. Then the equation of motion is

$$\frac{\partial^2 x_1}{\partial t^2} + \kappa^2 x_1 = kF \sin[k(x_0 + x_1) + \omega t], \quad (6.96)$$

which is the same as (6.95) if x replaces R and ω replaces $m(\Omega - \Omega_p)$.

If the forcing F is sufficiently strong, adjacent pendulums may collide. A collision takes place if the bobs of two pendulums whose supports are separated by a small value Δx_0 cross, that is, if $x_0 + x_1(x_0, t) > x_0 + \Delta x_0 + x_1(x_0 + \Delta x_0, t)$. Letting Δx_0 shrink to zero, we find that the condition for collision is $\partial x_1 / \partial x_0 < -1$. To obtain a crude estimate of the level of forcing that leads to collisions, we revert for the moment to the linear

approximation, where kx_1 can be dropped from the argument of the sine. In this approximation we find

$$x_1 = \frac{kF}{\kappa^2 - \omega^2} \sin(kx_0 + \omega t). \quad (6.97)$$

Hence in the linear approximation collisions occur if

$$f \equiv \frac{k^2|F|}{|\kappa^2 - \omega^2|} > 1. \quad (6.98)$$

This condition has not been derived self-consistently, since it requires that $|kx_1| > 1$ at some point in the cycle, so the approximation of dropping kx_1 from the argument of the sine is not valid. Nevertheless, numerical experiments with the full nonlinear equation (6.95) show that in many cases the crude criterion $f \gtrsim 1$ provides a fairly accurate prediction of when collisions occur. Thus f proves to be a useful dimensionless parameter for describing the effects of spiral waves on cold interstellar gas.

Figure 6.24, taken from Toomre (1977a), shows the results of numerical integrations of equation (6.95) with $\omega = 0.75\kappa$, assuming that any collisions are completely inelastic. In the top panel, where $f = 0.5$, no collisions occur, and the linear approximation is fairly accurate. As the forcing is increased, collisions first occur at $f = 0.98$. The bottom panel shows the behavior of the row of pendulums at even stronger forcing, $f = 1.5$. The most prominent features are the “traffic jams” where several pendulums are in contact. Each pendulum lingers for about 20% of the cycle in the traffic jam and then swings away to the left until it enters the next jam.

The motion of this row of pendulums provides a simple model for the motion of clouds in a spiral density wave. Each pendulum can be regarded as a cloud. The traffic jams correspond to narrow regions of shocked high-density gas. Equation (6.95) and Figure 6.24 show that dense, narrow fluid arms are a natural consequence of the response of inelastic interstellar gas to a tightly wound density wave. The figure also shows that the traffic jams are located near the minima of the forcing potential; in the language of §6.1.2, the gas arm coincides with the potential arm.

The picture we have just described is based on a grossly oversimplified model of the interstellar gas. In reality, most of interstellar space is occupied by a hot ($T \approx 10^6$ K) ionized medium, which hardly responds at all to the spiral gravitational field of the stars because its internal pressure is so high. The interstellar clouds occupy only a few percent of the disk volume. An additional complication is that the clouds have random velocities $v_{\text{RMS}} \approx 6 \text{ km s}^{-1}$; hence the cloud motion is not perfectly ordered and clouds can collide even outside the traffic jams (the typical time between collisions is $t_{\text{coll}} \approx 10 \text{ Myr}$). A better model would treat the clouds as atoms of an imperfect fluid with mean free path $v_{\text{RMS}} t_{\text{coll}} \approx 100 \text{ pc}$. In these more realistic

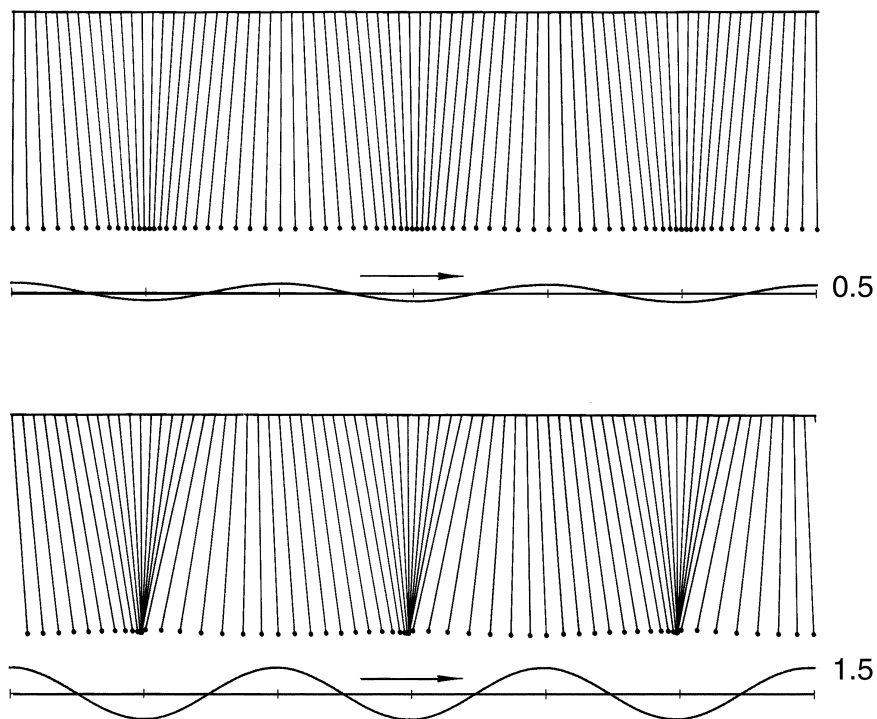


Figure 6.24 A simple model for the response of the interstellar gas. The diagram depicts a row of pendulums exposed to a horizontal force field proportional to $\sin(kx + \omega t)$. The sinusoids below the pendulums represent the forcing potential. The natural small-amplitude frequency of the pendulums is κ , so the equation of motion is (6.96). The motion is shown for $\omega = 0.75\kappa$ and forcing amplitudes of $0.5f$ (top) and $1.5f$ (bottom; see eq. 6.98). All collisions are assumed to be inelastic. From Toomre (1977a), based on unpublished work by A. J. Kalnajs. Reprinted, with permission, from the *Annual Review of Astronomy and Astrophysics*, **15** ©1977 by Annual Reviews (www.annualreviews.org).

models the traffic jam is spread out into a high-density region of small but finite width. An additional complication arises from the self-gravity of the gas. Although it is reasonable to neglect the self-gravity of the gas in linear theory, because the fractional surface density in gas is small, in the traffic jams the gas and star densities may be comparable and the self-gravity of the gas then plays an important role (Balbus 1988).

To sum up, interstellar clouds respond in a strongly nonlinear fashion to the imposed spiral field from a density wave. The resulting high-density traffic jams occur near the minima of the spiral potential, and can be identified with the gas arms seen in grand-design spiral galaxies. It is likely that star formation will proceed much more rapidly in these high-density regions, thus producing the young stars that delineate the bright-star arms. The displacement of the bright-star arms inside the gas and dust arms reflects the

time interval required for star formation (§6.1.3d).

6.4.2 Response of a density wave to the interstellar gas

The shocks or “traffic jams” induced by spiral arms in the interstellar gas dissipate energy. The source of this energy is the spiral density wave in the stars, and here we investigate whether and how this dissipation damps the density wave. In particular, equations (6.82)–(6.86) show that waves with pattern speed in the range $0 < \Omega_p < \Omega$ have negative energy density, so one might expect that such waves are actually *amplified* by energy dissipation—but this expectation turns out to be incorrect (Kalnajs 1972c).

We consider a fluid disk orbiting in a tightly wrapped spiral potential. The flow field and properties of the fluid, and the potential, are assumed to be stationary as viewed in a frame rotating with some pattern speed Ω_p . The amplitude of the spiral is assumed sufficiently large to induce shocks in the fluid, but sufficiently small that the fluid elements still travel on nearly circular orbits. We assume that the fluid is cold, in the sense that its sound speed is small compared to the circular speed. We allow for sources and sinks of mass, energy, and angular momentum at the inner and outer edges of the disk that may be required to maintain a steady state in the presence of dissipation.

We assume for simplicity that the fluid is adiabatic (Appendix F.1.4), so the specific entropy is conserved away from shocks, and increases discontinuously when the fluid element crosses a shock. In a steady-state system this means that the fluid streamlines cannot be closed, since the entropy per unit mass of the fluid on any closed streamline that crosses a shock must grow continually, and hence cannot be stationary. Thus the fluid must slowly spiral in or out. This steady radial drift of the fluid leads to a mass current $C_M(R)$, defined to be the rate of transfer of fluid mass outward across radius R , which must be independent of radius in a steady state.

Since the fluid is cold, its internal energy is much smaller than its orbital energy, and since the fluid orbits are nearly circular, the orbital energy per unit mass is $\epsilon(R) = \frac{1}{2}\Omega^2 R^2 + \Phi(R)$. The radial drift of the fluid adds energy to the annulus $(R, R+dR)$ at a rate $C_M[\epsilon(R) - \epsilon(R+dR)] = -C_M(d\epsilon/dR) dR$. Shocks dissipate energy, so their contribution to the energy balance in the annulus is a negative quantity that we shall call $\dot{E}_s dR$. Finally, gravitational forces between the spiral wave and the fluid transfer energy from the wave to the fluid. The energy density of the spiral wave, E_w , is given by equations (6.83) and (6.86), and the rate at which energy is transferred from the wave to the fluid in the annulus $(R, R+dR)$ may thus be written $-\dot{E}_w dR$. In a steady state, the sum of these three contributions must vanish, so we have

$$\dot{E}_w + C_M \frac{d\epsilon}{dR} = \dot{E}_s. \quad (6.99)$$

The equation for angular-momentum balance is similar, except that momentum is conserved across a shock front so the analog to the dissipation rate \dot{E}_s is zero. Thus we have

$$\dot{L}_w + C_M \frac{d\ell}{dR} = 0, \quad (6.100)$$

where $\ell(R) = \Omega R^2$ is the angular momentum per unit mass in circular orbits.

The energy and angular-momentum density of the spiral wave in the stellar disk are related by $E_w = \Omega_p L_w$ (eq. 6.86), so we may solve equations (6.99) and (6.100) to find

$$\dot{E}_w = \frac{\Omega_p d\ell/dR}{\Omega_p d\ell/dR - d\epsilon/dR} \dot{E}_s \quad ; \quad C_M = \frac{1}{d\epsilon/dR - \Omega_p d\ell/dR} \dot{E}_s. \quad (6.101)$$

Using the relations $d\epsilon/dR = \frac{1}{2}R\kappa^2$ and $d\ell/dR = \frac{1}{2}R\kappa^2/\Omega$ (Problem 3.32), we have

$$\dot{E}_w = \frac{\Omega_p}{\Omega_p - \Omega} \dot{E}_s \quad ; \quad C_M = \frac{2\Omega}{R\kappa^2(\Omega - \Omega_p)} \dot{E}_s. \quad (6.102)$$

Since \dot{E}_s , the rate of energy dissipation in shocks per unit radius, is negative, we conclude that waves outside corotation ($0 < \Omega < \Omega_p$) have $\dot{E}_w < 0$ and therefore lose energy to the fluid. Since these waves have positive energy density (from the discussion in §6.2.6), they are consequently damped. Waves inside corotation ($0 < \Omega_p < \Omega$) have $\dot{E}_w > 0$ and gain energy from the fluid, but have negative energy density and therefore are also damped. In short, *tightly wound spiral density waves are always damped by shocks in the interstellar fluid.*

Equations (6.102) also show that the fluid inside corotation drifts inward ($C_M < 0$), while fluid outside corotation drifts outward. Thus, *shocks in the interstellar fluid repel the fluid from corotation.*

These simple arguments lead to an apparent paradox: interactions between the density wave in the stars and the interstellar gas damp the density wave, yet the observations show that disk galaxies exhibit spiral structure if and only if they contain interstellar gas. The reason why galaxies *without* gas have no spiral structure is straightforward. As packets of density waves in a stellar disk wind up, they are absorbed at the Lindblad resonances (§6.2.5). The energy contained in the waves is transferred to random motions of the stars, causing the disk to heat and Toomre's Q to grow (see §8.4.2 below for more on this process). As the disk heats, the efficiency of the swing amplifier declines, reducing the susceptibility of the disk to further spiral-making. In gas-free disks, this process eventually leads to a hot, quiescent, axisymmetric disk like those seen in lenticular galaxies. Thus the growth of strong spirals, like the growth of ethanol-producing yeasts, is self-limiting because the heating produced by the spirals kills off the swing amplifier, just as fermentation of wine and beer is terminated when the ethanol poisons the yeast cells.

The reason why galaxy disks *with* gas always exhibit spiral structure is less clear. A related, but even harder question is what determines the properties of the spiral structure in a given galaxy—amplitude, number of arms, stationary or transient, grand-design or flocculent, etc. These issues are discussed in the next subsection, although many of the answers are tentative or incomplete.

6.4.3 Excitation of spiral structure

The presence of *some* form of spiral structure in galaxy disks that contain interstellar gas is not hard to understand. Star formation caused by any local gravitational instability in the interstellar gas creates a patch of new stars, which is sheared by differential rotation into an appearance of spirality. As time passes, the arm is sheared more and more while the most luminous stars die. Both effects lead to the gradual disappearance of the arm fragment. Meanwhile new arms form elsewhere. Such processes are almost certainly responsible for flocculent spirals such as M63 (Plate 9), in which the structure is patchy, without any long-range order or any corresponding spiral pattern in the old stellar disk (Elmegreen & Elmegreen 1984; Thornley 1996).

A more difficult and interesting question is the origin of intermediate-scale and grand-design spiral structure. The distinction between flocculent and grand-design spiral structure was concisely stated by Oort (1962): “In systems with strong differential rotation . . . spiral features are quite natural. Every structural irregularity is likely to be drawn out into part of a spiral. But *this* is not the phenomenon we must consider. We must consider a spiral structure extending over the whole galaxy, from the nucleus to its outermost part, and consisting of two arms starting from diametrically opposite points. Although this structure is often hopelessly irregular and broken up, the general form . . . can be recognized in many [galaxies].” In the remainder of this subsection we sketch some of the likely (and less likely) ways in which intermediate-scale and grand-design spirals can be excited.

(a) Excitation by companion galaxies There are some cases in which the cause of grand-design spiral structure is obvious. In particular, some of the most beautiful grand-design spirals have nearby companion galaxies. The best example is M51 (Plate 1), whose companion galaxy NGC 5195 is located near the tip of one of its two main spiral arms. Another example is M81 (Plate 8), which is interacting with its companion galaxies M82 and NGC 3077 (Yun, Ho, & Lo 1994). In a classic paper, Toomre & Toomre (1972) showed that the large-scale morphology of the spiral structure in M51 could be explained as a result of a recent encounter with NGC 5195 (Figure 6.25). The long HI spiral arm extending far beyond the stellar disk on the opposite side from NGC 5195 (Plate 5) is almost certainly also a product of this encounter. Other grand-design spirals with companions can be successfully modeled as well (Figure 6.26). Even in cases where no obvious companion is

present, it is possible that a grand-design spiral has been excited by a recent encounter with substructure in the dark-matter halo.

(b) Excitation by bars Since grand-design spirals can be excited by exterior gravitational perturbations from companions, it is natural to ask whether they can also be excited by interior perturbations, in particular by the bars that are found at the centers of many disks (§6.5). Here the chain of cause and effect is less clear. The natural first impression from examining images of barred spirals is that bars drive spiral structure, since the spiral arms in galaxies such as NGC 1300 (Plate 10) appear to emerge from the tips of the bar. Moreover the fraction of grand-design spirals appears to be much larger among barred galaxies than unbarred ones (Kormendy & Norman 1979; Elmegreen & Elmegreen 1982). However, this assumption leads to a contradiction: as we described in §6.1.3d, the pattern speed of the spiral structure in barred galaxies appears to be much smaller than the pattern speed of the bar; thus the spiral cannot be driven directly by the bar (Sellwood & Sparke 1988). If indeed the bar and the spiral are unrelated dynamical phenomena, is the spiral transient or stationary? And why does the bar appear to enhance the likelihood of a grand-design pattern? The answers are unknown.

(c) Stationary spiral structure The Lin–Shu hypothesis states that spiral structure consists of a stationary density wave, which remains unchanged except for an overall rotation over many crossing times. This hypothesis does not address the origin of the spiral structure, but it is natural to assume the following sequence of events: the stellar disk is unstable, and a growing spiral pattern forms, dominated by the most unstable mode; as the wave amplitude builds, the interstellar gas begins to shock, damping the density wave in the stars; eventually the wave reaches a stable, finite amplitude at which the damping and growth rates are equal, and this is the state in which the galaxy is found today. If this sequence is correct, the shape of a grand-design spiral wave should be similar to that of the most unstable mode.

This scenario successfully explains—and indeed predicted—several features seen in real galaxies:

- Spirals are trailing because unstable modes are usually trailing (§6.3.2c). The anti-spiral theorem (§6.1.4) does not apply because the steady state is established through dissipative shocks in the interstellar gas.
- As we have seen, grand-design spiral galaxies exhibit smooth, broad arms in infrared light, with relative amplitudes of up to 60% (§6.1.3a). Since the old disk stars that dominate the light in this waveband are also the dominant contributors to the mass of the disk, the presence of a spiral pattern in red light is direct evidence that a real *density* wave is present—the whole mass of the disk participates in the spiral wave—as predicted by the Lin–Shu hypothesis.
- The Lin–Shu hypothesis successfully explains why the interstellar gas defines narrower and stronger arms than the old stars (§6.3.4); why the

luminous young stars define narrow arms (the high density in the narrow gas arms triggers rapid star formation, and the lifetimes of the most luminous stars are so short that they cannot drift far from their formation sites before they die); and why the bright-star arms are downstream from the gas arms (the displacement represents the time required for the compressed gas clouds to collapse).

These observations provide persuasive evidence that grand-design spiral galaxies contain density waves, and that the behavior of the interstellar gas is roughly consistent with the assumption that the pattern is a stationary density wave. However, they do not establish that the Lin–Shu hypothesis is correct. Other mechanisms can excite transient trailing spirals in a stellar disk, including encounters with companion galaxies and swing amplification of generic initial disturbances (Figure 6.19). The compression of the interstellar gas and consequent star formation in a spiral arm occur on a timescale that is short compared to the rotation period of the galaxy, and thus transitory density waves with lifetimes of less than a rotation period can produce equally good agreement with the observations of gas and bright-star arms.

Indeed, there is strong circumstantial evidence that most spiral structure is *not* stationary, at least in unbarred galaxies:

- Even though the most regular and well-formed grand-design spirals are generally chosen for testing the Lin–Shu hypothesis, it has proved to be frustratingly difficult to determine the basic parameters of steady-state models, such as the location of the Lindblad and corotation resonances.
- N-body simulations of galactic disks do not generally develop stationary, grand-design spiral structure, so if the Lin–Shu hypothesis is correct, some fundamental physics must be missing from these simulations.
- The high rates of angular-momentum transport in strong, open grand-design spirals that we found in §6.1.5 rearrange the galaxy’s angular-momentum distribution on a timescale considerably less than 10 Gyr, and hence suggest that the strongest spiral patterns last for much less than that time.
- Scattering of disk stars by transient, rather than stationary, spirals provides promising explanations of the relation between age and velocity dispersion (§8.4.2) and age and metallicity (Sellwood & Binney 2002) in the solar neighborhood.

(d) Excitation of intermediate-scale structure The swing amplifier causes stellar disks to respond vigorously to a wide variety of gravitational disturbances, and in particular to the clumpy mass distribution in the gas disk (Julian & Toomre 1966). Moreover the characteristic scale of this predominantly trailing response, determined by the critical wavelength λ_{crit} of the stars (eq. 6.65), is much larger than the scale of the gas disturbances that excited it—using the parameters from §6.2.3, $\lambda_{\text{crit}} \simeq 4.5$ kpc in the solar neighborhood—so the stars enhance both the amplitude and the scale of the noise in the gas disk. In more picturesque language, “[the stellar disk] is

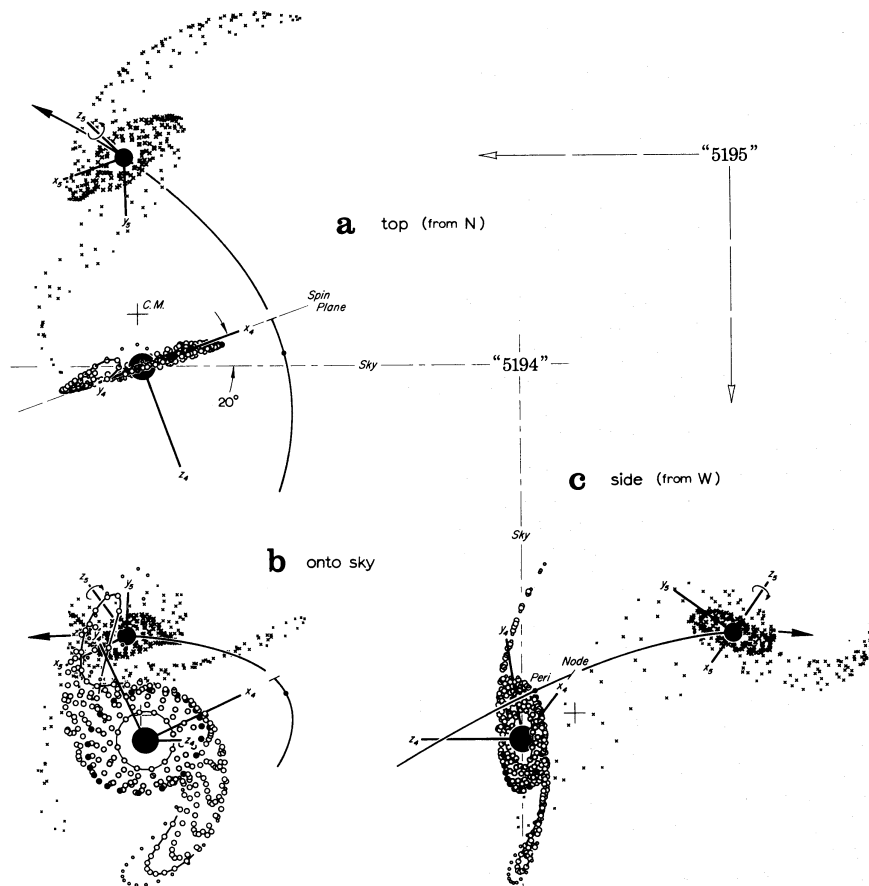


Figure 6.25 Model of the encounter between M51 and NGC 5195, shown in three orthogonal views. The lower left view can be compared with Figure 6.1 or Plate 1. Note that the low-density tidal tail at the 2 o'clock position relative to the center of M51 is similar to a low surface-brightness feature in Figure 6.1. In this pioneering experiment the galaxies were represented as point masses (the filled circles) surrounded by disks of orbiting test particles. From Toomre & Toomre (1972), reproduced by permission of the AAS.

to disk galaxies what a soundboard is to a piano. It organizes and augments the chaotic aspects of spiral galaxies . . . whenever the stellar disk is presented with a relatively flat spectrum of gravitational noise from the gas clouds, it picks out and augments the spatial frequencies which it prefers. And . . . it is this bias which leads to pictures that human astronomers happen to prefer as well" (Toomre & Kalnajs 1991).

At best, this is an appealing scenario for intermediate-scale spiral structure, rather than a real theory. Many questions remain unanswered: what determines whether a galaxy has flocculent, intermediate-scale or grand-design

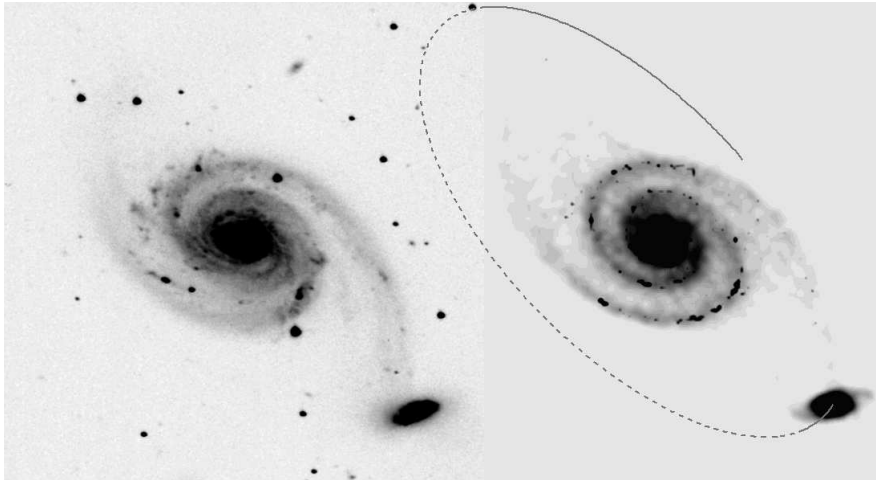


Figure 6.26 Model of the interacting galaxy pair NGC 7753, a large grand-design Sb spiral, and NGC 7752 (the small compact companion at lower right). The left panel shows a negative V -band image and the right panel shows an N-body simulation. The orbit of the companion is marked with a solid line above the disk plane and a dashed line below. The two galaxies are separated by 60 kpc. From Salo & Laurikainen (1993), reproduced by permission of the AAS.

structure? What are the relative roles played by inhomogeneities in the gas disk and substructure in the dark halo? Can nonlinear interactions between waves help to organize or generate spiral structure on intermediate scales? And what, if anything, does understanding such spiral structure teach us about the structure and evolution of galaxies?

6.5 Bars

6.5.1 Observations

Reviews of the properties of bars in disk galaxies are given in BM (§4.4.7), Sellwood & Wilkinson (1993) and Buta et al. (1996). Images of barred galaxies are shown in Plate 10 and Figures 6.27 and 6.28. Here we focus on those properties that are most relevant to interpreting the dynamics.

Bars vary from those that dominate the appearance of the disk, such as the ones shown in the images, to weak oval distortions that are visible only in careful Fourier decompositions of the light distribution. Thus the fraction of disk galaxies that are barred depends on the selection criterion. Classification by eye—still the most reliable method—shows that about 30% of spiral galaxies are strongly barred in optical light; the fraction rises to 50%

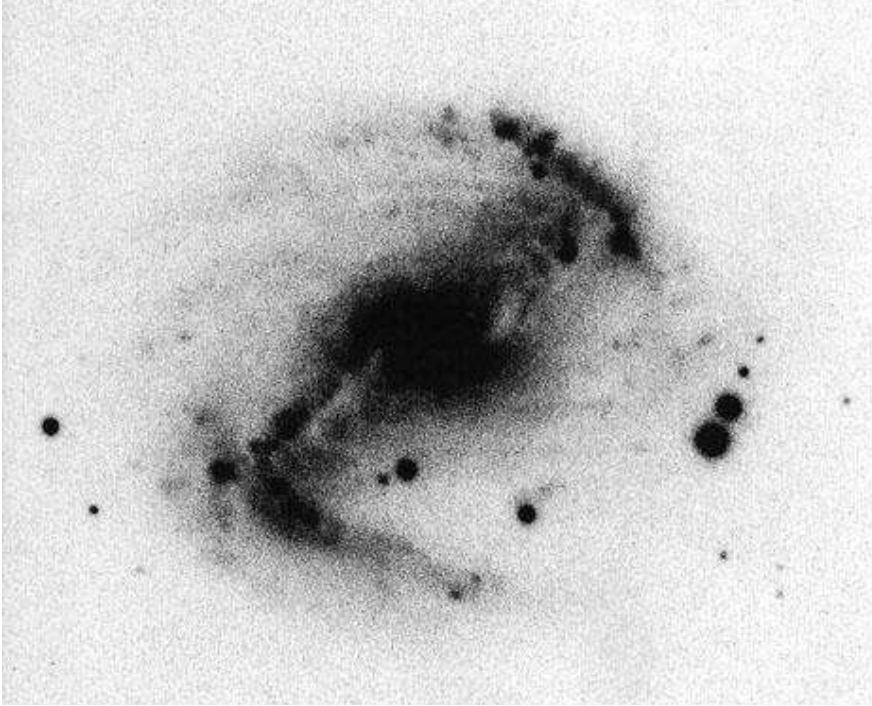


Figure 6.27 NGC 5383, a barred spiral (Hubble classification SBb). Note the two nearly straight dust lanes parallel to the bar, which appear as light streaks in this negative image. From Sandage & Bedke (1994), courtesy of the Observatories of the Carnegie Institution of Washington.

or more if weak bars are included. Bars appear even more prominent in near-infrared images (Eskridge et al. 2000). Since the near-infrared light traces the disk mass (§6.1.2), strong bars represent a substantial non-axisymmetric distortion of the mass distribution of the disk.

Our Galaxy is the nearest barred spiral (see §2.7e), although this is far from obvious because the characteristic non-axisymmetric structure of a bar cannot be seen in any edge-on galaxy. The presence of a bar at the center of our Galaxy was suggested long ago (Johnson 1957; de Vaucouleurs 1964), but this insight was not widely accepted for several decades, until overwhelming evidence had accumulated from several lines of investigation. Modern observational probes of the Galactic bar include the kinematics of HI and molecular gas in the central few kpc, near-infrared surface photometry, gravitational microlensing, and star counts (Gerhard 2002; Merrifield 2003).

Two of the Galaxy's satellites, the Large and Small Magellanic Clouds, are barred irregular galaxies. The nearest giant spiral galaxy, M31, contains an oval distortion which can be interpreted as a bar (Stark 1977), although most observers still classify M31 as a normal rather than a barred spiral.

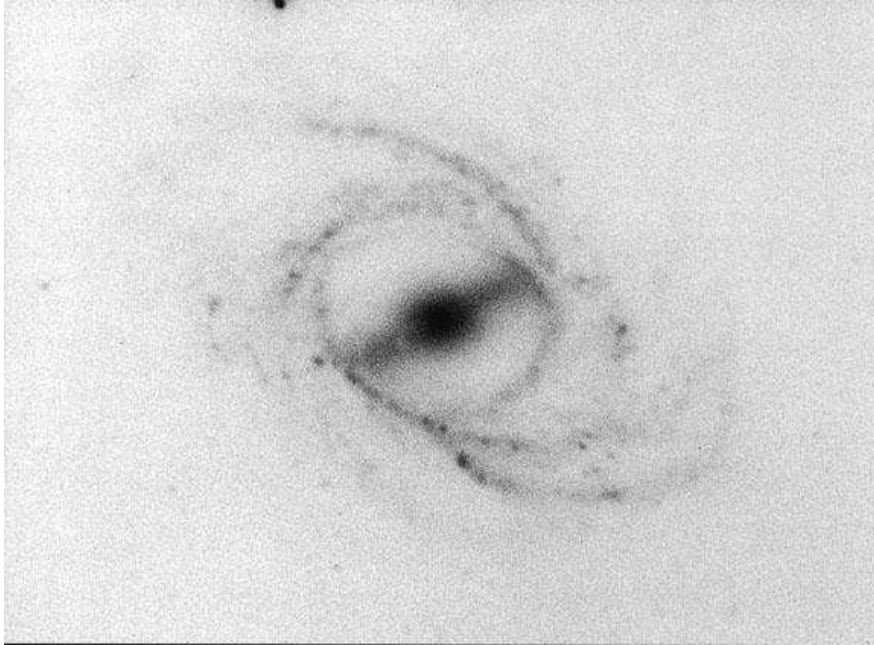


Figure 6.28 NGC 2523, an SBb spiral with a prominent ring. Note the narrowness of the spiral arms and how one forks into two branches just outside the ring. From Sandage & Bedke (1994), courtesy of the Observatories of the Carnegie Institution of Washington.

Bar properties depend on the Hubble type of the host galaxy (see Elmegreen 1996 for a review). Bars in galaxies with early Hubble types (SBa–SBbc; see page 29) are not centrally condensed—the surface brightness along the major axis is nearly flat with a sharp cutoff at the end of the bar—while later Hubble types (SBbc–SBm) contain bars with exponential surface-brightness profiles similar to that of the disk. The bars in early-type galaxies are also larger relative to the size of the disk. The bar contains up to a third of the galaxy’s total luminosity in early-type disk galaxies, and a smaller fraction in late-type galaxies. Since the bar is concentrated towards the center of the galaxy, it makes an even larger relative contribution to the luminosity and mass of the inner region. In contrast to elliptical galaxies, the isophotes of bars are not even approximately elliptical—typically they have shapes that are intermediate between ellipses and rectangles.

Bars are generally quite elongated. For SB galaxies, the median axis ratio in the equatorial plane is about 2:1, and for the Galaxy this ratio is 3:1. The strength of the bar can be characterized by the bar-interbar contrast ratio (see eq. 6.2 for the analogous quantity for spiral arms) which is typically $K \simeq 3\text{--}6$, comparable to the arm-interarm contrast of the strongest spirals (Elmegreen et al. 1996). The thickness of bars is hard to measure since it

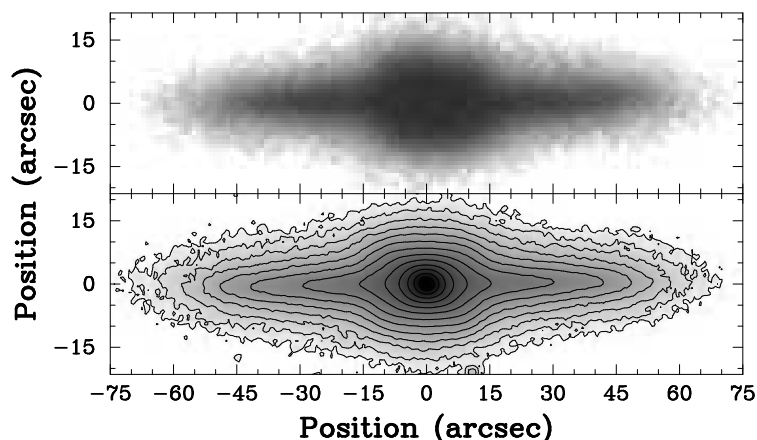


Figure 6.29 A boxy bulge in the edge-on galaxy NGC 1381. The top panel shows a near-infrared image and the bottom panel shows the isophotal contours for this image, spaced by $0.5 \text{ mag arcsec}^{-2}$. From Bureau et al. (2006), by permission of Blackwell Publishing.

is difficult to recognize an edge-on barred galaxy. However, there is good evidence that the bulges with boxy or “peanut-shaped” isophotes that are seen in $\sim 40\%$ of disk galaxies that are viewed nearly edge-on (see Figure 6.29 or Plate 15) are really bars whose long axis is perpendicular to the line of sight (BM §§4.4.7 and 11.3; Kuijken & Merrifield 1995; Bureau & Freeman 1999), while the bulges with elliptical isophotes that are seen in the remainder of edge-on galaxies represent either unbarred galaxies or bars whose long axis is along the line of sight. We shall amplify this point in §6.5.2c.

Prominent dust lanes are found in many bars. The lanes are displaced towards the leading side of the bar, i.e., they are displaced from the center line of the bar in the direction of rotation (this conclusion is based on the assumption that the spiral arms are trailing; see §6.1.3b). Often the lanes are remarkably straight, as in NGC 1300 and NGC 5383 (Plate 10 and Figure 6.27). At small radii the straight, offset lanes often appear to curl around to form an inner or nuclear ring. It is natural to assume that these dust lanes, like the ones in spirals, are regions of highly compressed or shocked gas. This conjecture is supported by the detection of radio emission from some of these dust lanes, just as is seen in spirals (§6.1.2).

Many barred spirals exhibit clusters of young stars and HII regions near the tips of the bars, which presumably arise from rapid star formation in high-density gas compressed by the bar potential, just as rapid star formation is seen in the dense gas in spiral arms.

Apart from the presence of the bar and its associated effects, there are few systematic differences between barred disk galaxies and their unbarred counterparts.

The pattern speed We have seen that there is little direct evidence to

support the Lin–Shu hypothesis that spiral structure is a stationary density wave with a well-defined pattern speed (§6.4.3). In contrast, there are strong reasons to believe that bars *do* have a well-defined pattern speed—most obviously, bars are straight rather than trailing, so it is most unlikely that they are transitory structures that are being wound up by differential rotation. The bar pattern speed Ω_b is usually parametrized by the ratio

$$\mathcal{R} = R_{\text{CR}}/a_b \quad (6.103)$$

of the corotation radius (§6.1.3d) to the bar semi-major axis.¹⁵

Dynamical arguments show that weak bars, at least, must have $\mathcal{R} > 1$; that is, weak bars cannot extend beyond corotation (§6.5.2a). Thus, bars are often said to be “fast” if $\mathcal{R} \approx 1$ and “slow” if $\mathcal{R} \gg 1$.

The most straightforward way to measure the pattern speed is from the flow pattern of a tracer population that obeys the continuity equation, such as old disk stars; in this case Ω_b is given by equation (6.13). This approach has been used to measure the bar pattern speed in over a dozen lenticular and early-type spiral galaxies, with a typical uncertainty of about 30%. The method works less well for late-type spirals, in which the kinematic and photometric data are contaminated by young stars that may be formed in shocks associated with the bar: these do not satisfy the continuity equation and therefore violate the fundamental premise on which equation (6.13) rests. Figure 6.30 shows the bar semi-major axis a_b and corotation radius R_{CR} for 19 barred galaxies examined by this and other methods. The ratio $\mathcal{R} = R_{\text{CR}}/a_b$ is constant along the dashed lines. Within the errors, all bars have $0.9 \lesssim \mathcal{R} \lesssim 1.3$ and thus are “fast.”

Two other lines of evidence also indicate that bars are fast:

- (i) Numerical models of gas flow in barred galaxies show that $\mathcal{R} \simeq 1.2 \pm 0.2$ is required to reproduce the straight, offset dust lanes found in many bars (see §6.5.2d).
- (ii) By assuming a constant mass-to-light ratio for the stars, one can determine the gravitational field of a barred galaxy from optical or near-infrared photometry. One then finds the steady-state flow pattern of cold gas in this field for a range of pattern speeds, and matches this flow to the observed gas kinematics and density field to estimate the pattern speed. This procedure has been applied with varying success to about ten barred galaxies, and the more reliable results, shown in Figure 6.30, are consistent with $\mathcal{R} \simeq 1$.

A less reliable way to estimate the pattern speed is to assign specific features in the photometry and kinematics (e.g., rings, reversals in the relative phase

¹⁵ The parameter \mathcal{R} is somewhat imprecise, both because the bar does not have a sharp edge and because the corotation radius is difficult to determine observationally, even if the pattern speed is known. Furthermore, associating a single radius with corotation is valid only when the non-axisymmetric forces from the bar are weak.

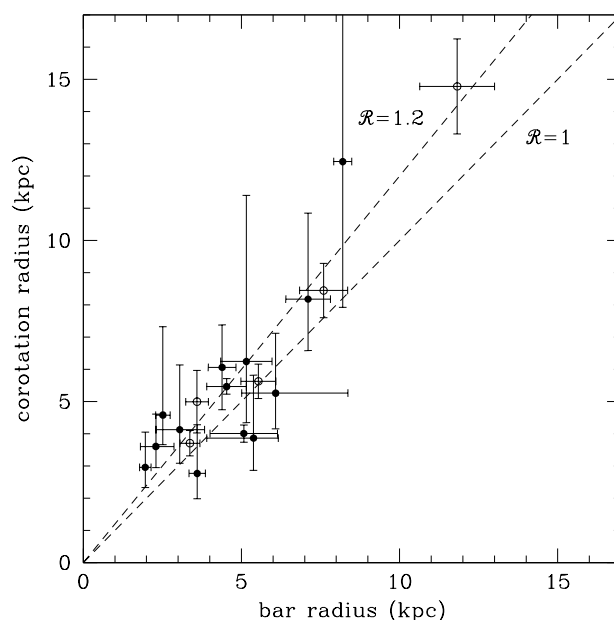


Figure 6.30 Measurements of bar semi-major axes a_b (horizontal axis) and corotation radii R_{CR} (vertical axis). Uncertainties of $\pm 10\%$ are assigned when errors are not quoted in the original papers. The Hubble constant is assumed to be $70 \text{ km s}^{-1} \text{ Mpc}^{-1}$. The ratio $\mathcal{R} = R_{CR}/a_b$ is constant on the dashed lines. Filled circles denote measurements using equation (6.13) (Merrifield & Kuijken 1995; Gerssen, Kuijken, & Merrifield 1999; Debattista & Williams 2001; Debattista, Corsini, & Aguerri 2002; Gerssen, Kuijken, & Merrifield 2003; Aguerri, Debattista, & Corsini 2003); open circles are derived by fitting gas kinematics to dynamical models (Hunter et al. 1988; England, Gottesman, & Hunter 1990; Lindblad, Lindblad, & Athanassoula 1996; Laine, Shlosman, & Heller 1998; Weiner et al. 2001; Weiner & Sellwood 1999; Debattista, Gerhard, & Sevenster 2002). All measured bar pattern speeds lie near $\mathcal{R} = 1$.

of gas and potential arms, etc.) to the Lindblad or corotation resonances (Patsis, Skokos, & Athanassoula 2003). An example of a prominent ring is seen in NGC 2523 in Figure 6.28.

In summary, most bar pattern speeds appear to lie in the range $0.9 \lesssim \mathcal{R} \lesssim 1.3$ and thus bars are fast. This is a striking result, since we shall find that galaxies in which the mass within the disk radius is dominated by a dark halo are expected to have slow bars (see §8.1.1d).

6.5.2 Dynamics of bars

Bars, like many elliptical galaxies, are triaxial stellar systems. However, the dynamical structures of these two types of triaxial system are quite different,

since the pattern speeds of ellipticals are believed to be small or zero, $\mathcal{R} \gg 1$, whereas the pattern speeds of bars are high, $\mathcal{R} \simeq 1$. Bars are easier to study observationally than triaxial elliptical galaxies, both because the directions of the principal axes are known (two of the axes lie in the surrounding disk) and because the flow pattern of cold gas can be used to map the gravitational field.

Like some spiral structure, a bar can be thought of as a density wave. Most bars are so strong that the linear perturbation theory that we used to study spiral structure must be supplemented by other tools for describing bar dynamics, but the artificial case of weak bars provides an introduction to many of the dynamical phenomena seen in their stronger cousins.

(a) Weak bars In a nearly axisymmetric potential, almost all of the stars are on loop orbits parented by the nearly circular, closed, loop orbit described by equation (3.148a) with $C_1 = 0$, or, in a continuum description, by equations (6.43) with $h_a = 0$. Equations (6.43) can be combined with Poisson's equation and the linearized continuity equation (6.49) to construct self-consistent weak bars. However, we can obtain some insight by a simpler route. The long axis of the bar lies along $\varphi = 0, \pi$ (see eq. 3.143 and the subsequent discussion), so the closed loop orbit is elongated along the bar whenever the quantity C_2 defined by equation (3.148b) is positive. Stars on orbits that are elongated along the bar generally contribute to the bar-like nature of the overall gravitational field, whereas stars on orbits that are elongated perpendicular to the bar tend to cancel the gravitational field of the bar. Hence *any self-consistent weak bar must be composed mainly of orbits with $C_2 > 0$* .¹⁶ Let us examine the behavior of C_2 as a function of radius, considering first the region inside corotation, $R < R_{CR}$. Since the angular speed decreases outward in almost all galaxies, inside corotation we have $\Omega > \Omega_b$ so the coefficient of the term Φ_b/R in C_2 is $2\Omega/(\Omega - \Omega_b) > 2$. For most reasonable bar potentials, the term proportional to Φ_b swamps the term proportional to $d\Phi_b/dR$. Since $\Phi_b < 0$ by definition (see the discussion following eq. 3.143), we find that $C_2 > 0$ if and only if $\Delta > 0$ or $\Omega_b > \Omega - \frac{1}{2}\kappa$. Outside corotation, $\Omega < \Omega_b$ so the term $2\Omega\Phi_b/[R(\Omega - \Omega_b)]$ has the opposite sign to Φ_b , as does the term $d\Phi_b/dR$ since the amplitude of the bar potential $|\Phi_b|$ normally decreases outward. Since $\Phi_b < 0$, we conclude that $C_2 > 0$ if and only if $\Delta < 0$, which in turn requires $\Omega_b > \Omega + \frac{1}{2}\kappa$. We conclude that *self-consistent weak bars can exist only between the inner Lindblad resonance and corotation, or outside the outer Lindblad resonance*. Bars of the second kind are not relevant since they cannot extend continuously from the origin, as the observations require. Thus (i) weak bars must rotate sufficiently rapidly to avoid the inner Lindblad resonance, and (ii) they must end before corotation, that is, $\mathcal{R} > 1$ (Contopoulos 1980).

¹⁶ In the notation of §3.3.2, most of the bar orbits must be parented by the long-axis orbits, a statement that remains valid even for strong bars.

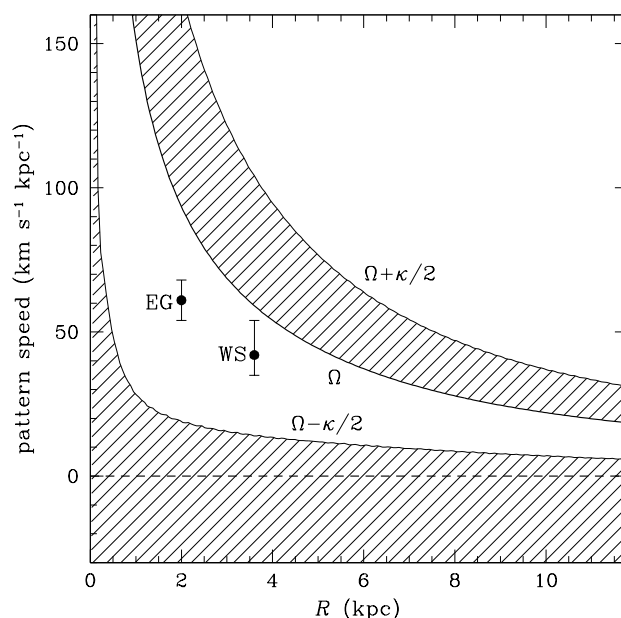


Figure 6.31 This plot shows Ω and $\Omega \pm \frac{1}{2}\kappa$ for Model I of the Galaxy described in §2.7 (cf. Figure 6.11). The shading denotes the regions in which a self-consistent weak bar cannot exist. The points mark the semi-major axis and pattern speed for the Galactic bar derived by Englmaier & Gerhard (1999) and Weiner & Sellwood (1999).

The corresponding constraints on the pattern speed of the Galactic bar are shown in Figure 6.31. These constraints are approximate, because they are based on both linear perturbation theory and a crude treatment of the relation between orbit shape and the gravitational field. In particular, they apply only to the main body of the bar, containing the bulk of its mass, and do not exclude the presence of an inner Lindblad resonance near the center of the bar, which is often required to match the gas flow—see part (d) below.

(b) Strong bars Constructing self-consistent dynamical models for strong bars is a difficult task, which has been approached with a variety of theoretical tools. Each of these tools has helped to illuminate different aspects of bar structure.

The only exact, self-consistent models of bars were constructed by Freeman (1966) (Box 4.2). Unfortunately, the gravitational potential in the Freeman bars is quadratic in the coordinates, so the potential and the properties of the orbits are rather different from real bars.

Numerical integrations of orbits in realistic barred potentials suggest that most of the stars in the bar must be on prograde orbits parented by the

closed long-axis orbits (the sequence x_1 of §3.3.2), since this is the only major orbit family that is elongated in the same sense as the potential (Contopoulos & Papayannopoulos 1980; Teuben & Sanders 1985; Athanassoula 1992a).¹⁷

The simplest way to construct self-consistent bar models is by N-body simulations of disks that are susceptible to the bar instability (Sparke & Sellwood 1987; O’Neill & Dubinski 2003). These simulations confirm that bars are composed mainly of orbits associated with the x_1 sequence; in addition, they find that the bar shape is intermediate between an ellipse and a rectangle and the end of the bar is close to corotation, just as in real bars. Of course, there is no guarantee that bars are formed from an initially axisymmetric disk through the bar instability (Sellwood 2000b), so bars constructed in this way may be only a subset of the bar models that could be constructed by more general procedures such as Schwarzschild’s method (Pfenniger 1984; Zhao 1996). Nevertheless, the properties of bars formed by the bar instability agree remarkably well with observed bars (O’Neill & Dubinski 2003).

(c) The vertical structure of bars Much of our understanding of disk and bar dynamics has been developed through the study of razor-thin analytic and numerical models of stellar systems. Three-dimensional models of axisymmetric disks generally behave in about the same way as their two-dimensional analogs. Thus it was surprising when three-dimensional simulations of bar-unstable disks showed that newly formed bars often bend out of the disk plane. The vertical bending motion rapidly loses coherence and is transformed into random vertical motions, leaving a bar that is substantially thicker than the surrounding disk (Raha et al. 1991). When viewed from the side (i.e., from a line of sight in the disk plane that is perpendicular to the long axis of the bar), the thickened bar appears boxy or peanut-shaped (Combes & Sanders 1981; Combes et al. 1990), and thus it is likely that the boxy bulges seen in some edge-on disk galaxies are really edge-on bars. The structure of the orbit families of three-dimensional bars is considerably more complicated than that of planar bars (Pfenniger & Friedli 1991; Skokos, Patsis, & Athanassoula 2002).

The dynamics of the buckling instability in simplified stellar systems is described in §6.6.2.

(d) Gas flow in bars Numerical models of gas flow through a rotating bar can be compared with observations to constrain the pattern speed and other properties of bars in real galaxies. The straight, offset dust lanes seen in barred galaxies such as NGC 1300 (Plate 10) and NGC 5383 (Figure 6.27), like the dust lanes in spiral galaxies, presumably mark the location of traffic jams containing a greatly enhanced density of interstellar material, which arise when the trajectories of interstellar clouds intersect. In hydrodynamic models, the traffic jams appear as shock waves in the fluid.

¹⁷ Our earlier finding that most of the stars in weak bars must be on orbits with $C_2 > 0$ is the restriction of this statement to weakly non-axisymmetric potentials.

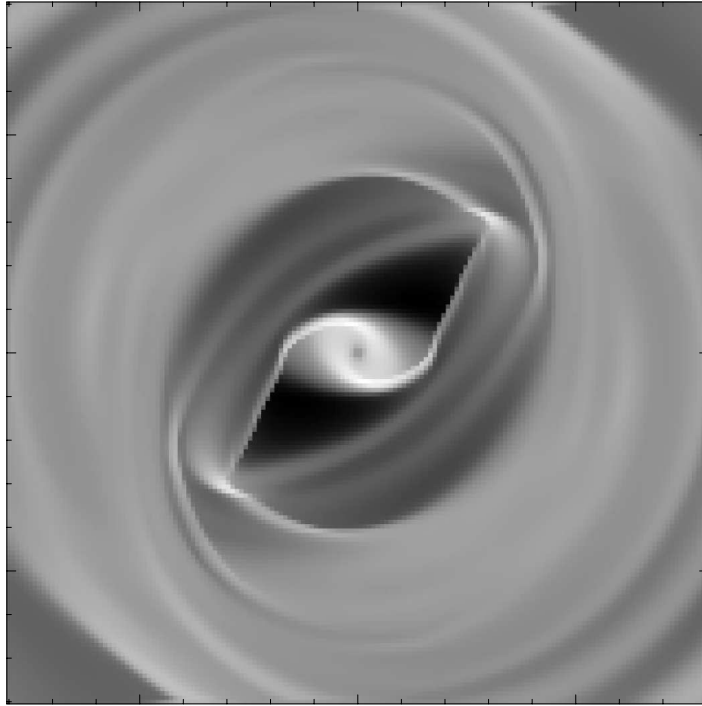


Figure 6.32 A simulation of gas flow in a barred spiral. The brightness is proportional to the local gas density. The bar (not shown) is oriented at 45° to the horizontal, and rotates clockwise. The bar is 10 kpc long, in comparison to the box size of 16 kpc. The thin, straight, bright lines mark shocks, which are offset from the leading edge of the bar and angled by 20° to the bar major axis. The rotation parameter $\mathcal{R} = 1.2$ (eq. 6.103). From Athanassoula (1992b), by permission of Blackwell Publishing.

Athanassoula (1992b) finds that reproducing these shocks requires that the rotation parameter defined by equation (6.103) satisfies $\mathcal{R} \simeq 1.2 \pm 0.2$: for smaller values of \mathcal{R} the dust lanes are straight but centered on the bar's major axis, not offset, while for larger values of \mathcal{R} the shocks are not straight (see Figure 6.32). Offset shocks appear if and only if the potential supports two families of periodic orbits, the x_1 family at larger radii elongated parallel to the bar's long axis and the x_2 family at smaller radii elongated perpendicular to the long axis (§3.3.2), and this in turn requires that the bar pattern has at least one inner Lindblad resonance (Sanders & Tubbs 1980; Athanassoula 1992b). Loosely speaking, the offset shocks arise as the gas attempts to make the transition from one family to the other without orbit crossing.

The dust lanes along the bar are not sites of rapid star formation, possibly because the velocity shear downstream is so large that the gas remains gravitationally stable. Rapid star formation is observed only at the ends of the bar, and it is encouraging that numerical simulations exhibit strong den-

sity enhancements with low shear near the end of the bar, suggesting that these are fertile sites for star formation.

Gas-flow models such as the one shown in Figure 6.32 exhibit strong non-circular motions, even for rather weak bars. For example, the strongly perturbed velocity field in Figure 6.32 results from a non-axisymmetric force that is nowhere greater than about 20% of the axisymmetric force. This strong response arises because most of the gas is close to a Lindblad resonance and hence responds strongly to even a weak imposed force. These non-circular motions give rise to characteristic kinematic features that can be used to identify edge-on barred galaxies (see BM Figure 4.60, and Bureau & Freeman 1999).

The net mass-inflow rate, given by the mass-weighted azimuthally averaged radial velocity, is much smaller than the RMS radial velocity, but not zero—typically it is $\lesssim 1 \text{ km s}^{-1}$. Even this small average velocity leads to substantial inflow over the age of the galaxy, since 1 km s^{-1} is equivalent to 10 kpc in 10 Gyr. Thus strong bars can transport most of the gas in a galaxy to its center within the lifetime of the galaxy.

The gas is transported inward until it reaches radii where the shocks peter out, inside the outermost inner Lindblad resonance, and gathers on nearly circular orbits of the x_2 sequence. There it forms a ring, which probably can be identified with the **nuclear rings** seen in many barred galaxies (Knapen 1999; Pérez-Ramírez et al. 2000). These typically have radii of a few hundred pc and are often sites of vigorous star formation. Nuclear rings are promising reservoirs for the gas that feeds the accretion disks in active galactic nuclei and builds the massive black holes that are present in the centers of most galaxies; unfortunately, the dynamical processes by which the gas migrates from the nuclear ring ($r \sim 100 \text{ pc}$) inward by five orders of magnitude in radius to the accretion disk ($r \sim 10^{-3} \text{ pc}$) remain unclear (Phinney 1994). Possible mechanisms include nested bars and spirals, mergers with galaxy satellites, magnetic fields, or gravitational instability.

Other ring structures are also present in barred galaxies. Many barred galaxies contain elliptical **inner rings** that are aligned with the bar and touch its end (see Figure 6.28 for a particularly clear example). A smaller fraction of barred galaxies contain an **outer ring** with radius of about 2.2 ± 0.5 times that of the inner ring (see BM §§4.1.1 and 4.4.7, and Buta 1995). Rings of gas and young stars form naturally at resonances, where distinct families of periodic orbits intersect, leading to orbit crossing. It is tempting to identify outer rings with the outer Lindblad resonance of the bar (Schwarz 1981; Kalnajs 1991), just as nuclear rings are identified with the inner Lindblad resonance. Inner rings are sometimes identified with either the corotation resonance, or the $m = 4$ inner Lindblad resonance, or the $m = 2$ **ultraharmonic resonance** which is caused by nonlinear effects from the $m = 2$ perturbing potential (Patsis, Skokos, & Athanassoula 2003). Such claims must be regarded as plausible speculations, since there is little

direct observational confirmation.

(e) Slow evolution of bars N-body simulations of bars usually focus on their formation from an axisymmetric disk by the bar instability, which takes only a few crossing times. However, most bars are $\gtrsim 10^2$ crossing times old, and slow evolutionary processes over this much longer time can thus change the properties of bars, or even destroy them.

Bars can be destroyed by the growth of relatively small central mass concentrations, such as massive black holes or nuclear star clusters. The growth of a central mass reduces the volume of phase space occupied by the orbits parented by the x_1 family of periodic orbits, which provides the main support for the bar, leading to the dissolution of the bar when the central mass exceeds a few percent of the bar mass (Hasan & Norman 1990; Norman, Sellwood, & Hasan 1996).¹⁸ The central mass concentration may arise from the inward transport of gas by the bar, as outlined in topic (d) above; just as with humans, the dissipative lifestyle of a bar can lead to its early demise. However, it remains unclear how effective such processes are in practice: massive black holes are probably too small to destroy bars, and central gas concentrations may be too diffuse (Shen & Sellwood 2004).

Other important evolutionary processes involving bars include (i) transfer of angular momentum between the bar and dark halo by dynamical friction, described in §8.1.1d, or between the bar and the disk; (ii) the modification or destruction of bars by mergers with small galaxies or halo substructure, and (iii) the role of bars in forming bulges, lenses, and other apparently distinct components of disk galaxies (Kormendy & Kennicutt 2004).

6.6 Warping and buckling of disks

6.6.1 Warps

In visible light, disk galaxies are remarkably thin and flat. However, the HI disk, which usually extends to larger radii than the visible disk, frequently shows a noticeable warp. Warps are most clearly visible in edge-on disk galaxies (BM Figure 8.30) but can often be detected in—or inferred from the velocity fields of—galaxies that are inclined to the line of sight (Plate 6 and BM §8.2.4). There is also a pronounced warp in the HI disk of our own Galaxy (BM §9.2.6). Warped disks are common: for example, *all* of the spiral galaxies in the Local Group (our own Galaxy, M31, and M33) are warped. Usually, a warped galaxy is twisted upward at one side and downward at

¹⁸ A similar process was discussed in §3.8.3: central mass concentrations in non-rotating triaxial potentials can destroy the triaxiality, in this case by eroding the population of box orbits that supports the triaxiality.

the other; these are sometimes called “integral sign” warps since such a disk resembles an integral sign when viewed edge-on.

Explanations of the origin of warps can be divided into two broad classes: (i) those in which the warp is generated by hydrodynamic or magnetohydrodynamic forces on the HI gas; in this case the modest warps seen in some stellar disks arise from the response of the stars to the gravitational field from the warped gas; (ii) those in which the warp arises from gravitational forces that affect the stars and gas equally; in this case the warp is more prominent in the gas because it extends to larger radii. We shall concentrate on theories of the second kind mainly because they appear to be simpler and more plausible. For general reviews of warps, see Binney (1992) and Nelson & Tremaine (1996).

(a) Kinematics of warps Warps in a stellar disk can be regarded as vertical oscillations or bending waves of the disk, analogous to the horizontal oscillations or density waves that we examined earlier in this chapter. In the usual cylindrical coordinates (R, ϕ, z) , a warped but razor-thin disk can be described by its height $z(R, \phi, t)$ above the galactic plane at position (R, ϕ) and time t . Suppose that initially the warp has the form

$$z_d(R, \phi, t = 0) = z_0(R) \cos m\phi \quad (m \geq 0), \quad (6.104)$$

and the vertical velocity at every point is zero. An integral-sign warp like those seen in most warped galaxies has $m = 1$. We assume for the moment that the outer disk has negligible mass, so its motion is determined by a fixed, azimuthally symmetric potential $\Phi(R, z)$ due to the halo and the massive, flat, inner disk. We now ask how the warp evolves with time. We shall use the epicycle approximation, but neglect the epicycle motion in R and ϕ , an approximation that is valid if the horizontal scale of the warp is much larger than a typical epicycle amplitude X (eq. 3.91), as is usually the case. Thus by equation (3.90) a particle labeled by i has coordinates

$$R_i = \text{constant} \quad ; \quad \phi_i(t) = \Omega(R_i)t + \phi_{0i} \quad ; \quad z_i(t) = Z_i \cos[\nu(R_i)t + \zeta_i], \quad (6.105)$$

where Ω is the circular frequency and the vertical frequency ν is defined by equation (3.79b).

In order to match the assumed initial conditions—vertical displacement (6.104) and zero vertical speed at $t = 0$ —we must have $\zeta_i = 0$ and $Z_i = z_0(R_i) \cos m\phi_{0i}$. Therefore the disk height at time t is given by

$$z_d[R_i, \phi_i(t), t] = z_0(R_i) \cos(m\phi_{0i}) \cos[\nu(R_i)t]. \quad (6.106)$$

Replacing ϕ_{0i} by $\phi_i - \Omega(R_i)t$ and dropping the subscripts, we have

$$z_d(R, \phi, t) = z_0(R) \cos[m(\phi - \Omega t)] \cos(\nu t), \quad (6.107)$$

where Ω and ν are evaluated at R . Using the identity $\cos a \cos b = \frac{1}{2} \cos(a + b) + \frac{1}{2} \cos(a - b)$, we find

$$z(R, \phi, t) = \frac{1}{2} z_0(R) \left\{ \cos[m\phi - (m\Omega - \nu)t] + \cos[m\phi - (m\Omega + \nu)t] \right\}. \quad (6.108)$$

We see that the warp runs around the circle $R = \text{constant}$ as two waves. These are analogous to the kinematic density waves discussed in §6.2.1a, and are called kinematic bending waves. The angular phase velocity or pattern speed Ω_p of the two waves is

$$\Omega_p = \begin{cases} \Omega(R) - \nu(R)/m & \text{("slow" wave),} \\ \Omega(R) + \nu(R)/m & \text{("fast" wave).} \end{cases} \quad (6.109)$$

In a flattened galaxy $\nu > \Omega$, so the pattern speed of the $m = 1$ slow wave is retrograde, $\Omega_p < 0$. From now on, we consider only the slow wave, since the fast wave will disappear rapidly through the winding process discussed below.

If only the slow $m = 1$ wave is present, the warped disk attains its maximum height at $\phi(R, t) = [\Omega(R) - \nu(R)]t$. Because $\Omega - \nu$ is not exactly constant with radius, the warp will tend to wind up. Physically, the winding occurs because each ring in the warped disk precesses at a different rate. Our analysis of the winding problem in spiral arms from §6.1.3c can be applied without change, and we find that the pitch angle of the line of maximum height is (cf. eq. 6.6)

$$\cot \alpha = Rt \left| \frac{d}{dR} (\Omega - \nu) \right|. \quad (6.110)$$

Thus kinematic warps suffer from the same winding problem as kinematic density waves.¹⁹

The severity of the winding problem depends on the circular and vertical frequencies Ω and ν and their dependence on radius. In a spherical galaxy, in which $\Phi = \Phi(\sqrt{R^2 + z^2})$, $\Omega = \nu$ (Problem 3.14) so no winding occurs. Physically, the warp does not wind up because an inclined ring does not precess in a spherical field.

To examine the winding rate in a more realistic flattened galaxy, let us continue to neglect the gravitational field from the disk, and consider the halo to be described by the logarithmic potential of equation (2.71a),

$$\Phi_L = \frac{1}{2} v_0^2 \ln \left(R_c^2 + R^2 + \frac{z^2}{q_\Phi^2} \right), \quad (6.111)$$

¹⁹ One difference is that in the usual case where $\nu > \Omega$ and both frequencies decrease outward, bending waves wind up in the opposite sense to density waves, that is, trailing bending waves wind into leading ones.

where we shall set the core radius R_c to zero for simplicity. Then

$$\Omega = \frac{v_0}{R} \quad ; \quad \nu = \frac{v_0}{q_\Phi R}, \quad (6.112)$$

and hence equation (6.110) yields

$$\cot \alpha = \frac{v_0 t}{R} \left| \frac{1}{q_\Phi} - 1 \right|. \quad (6.113)$$

To model the implications of this result for our own Galaxy, we use $v_0 = 220 \text{ km s}^{-1}$, $R = 12 \text{ kpc}$, and $t = 10 \text{ Gyr}$. Maps of Galactic HI show no observable twist in the warp, suggesting that conservatively we may assume that the pitch angle $\alpha \gtrsim 45^\circ$. Then equation (6.113) yields $|q_\Phi^{-1} - 1| < 0.005$. It is highly improbable that the Galactic halo is so accurately spherical. Moreover, accounting for the disk potential will only increase the winding rate for a kinematic bending wave. Thus we require a mechanism either to excite fresh warps or to maintain a warp against winding.

(b) Bending waves with self-gravity In light of the analogy between bending and density waves, it is natural to ask whether the winding problem could be solved by the self-gravity of the disk. As a first step, it is useful to examine the behavior of tightly wound bending waves, even though observed warps show little or no sign of spirality in their inner parts.

Once the disk has non-zero mass, the definition of the circular and vertical frequencies, $\Omega(R)$ and $\nu(R)$, must be made more precise. We write

$$\Omega^2 = \frac{1}{R} \frac{\partial(\Phi_h + \Phi_d)}{\partial R} \Big|_{z=0} \quad ; \quad \nu^2 = \frac{\partial^2 \Phi_h}{\partial z^2} \Big|_{z=0}. \quad (6.114)$$

Here $\Phi_h(R, z)$ is the potential due to all components of the galaxy other than the disk and $\Phi_d(R, z)$ is the potential due to the unperturbed flat disk. The reason for this distinction is that the vertical restoring force due to the self-gravity of the disk depends entirely on the shape of the bending wave (as we shall see below) and hence cannot be included in the unperturbed potential.

We consider first the behavior of bending waves in an infinite, self-gravitating thin sheet of constant surface density Σ . The sheet initially occupies the plane $z = 0$, at the minimum of a ‘‘halo’’ potential $\frac{1}{2}\nu^2 z^2$. The sheet is uniform in the x - and y -directions, and we neglect the effects of rotation. Assume that the height of the sheet $z_s(x, t)$ oscillates according to

$$z_s(x, t) = z_0 e^{i(kx - \omega t)}. \quad (6.115)$$

When we neglect any horizontal motions, the equation of motion is simply

$$\frac{\partial^2 z_s(x, t)}{\partial t^2} = g_z(x, t) - \nu^2 z_s(x, t), \quad (6.116)$$

where g_z is the vertical acceleration due to the self-gravity of the sheet. We now evaluate g_z .

Consider an infinite straight wire with mass per unit length ζ . The gravitational potential at a distance R from the wire is $\Phi(R)$, and Gauss's theorem (2.12) implies that $4\pi G\zeta = 2\pi d\Phi/dR$, so $\Phi(R) = 2G\zeta \ln R + \text{constant}$. Thus, if we consider the sheet as a collection of such wires, the potential due to the sheet is

$$\Phi(x, z, t) = 2G\Sigma \int_{-\infty}^{\infty} dx' \ln \sqrt{(x-x')^2 + [z-z_s(x', t)]^2} + \text{constant}. \quad (6.117)$$

The vertical acceleration is

$$g_z(x, t) = -\left. \frac{\partial \Phi}{\partial z} \right|_{z=z_s} = -2G\Sigma \wp \int_{-\infty}^{\infty} dx' \frac{z_s(x, t) - z_s(x', t)}{(x-x')^2 + [z_s(x, t) - z_s(x', t)]^2}, \quad (6.118)$$

where \wp denotes the Cauchy principal value (eq. C.6). If the amplitude of the bending waves is small in the sense that $|kz_0| \ll 1$, we can drop the term involving $z_s(x, t) - z_s(x', t)$ in the denominator. Substituting equation (6.115) for the bending wave and replacing the dummy variable x' by $u \equiv x' - x$, we have

$$\begin{aligned} g_z(x, t) &= -2G\Sigma z_0 e^{i(kx - \omega t)} \wp \int_{-\infty}^{\infty} du \frac{1 - e^{iku}}{u^2} \\ &= -2\pi G\Sigma |k| z_s(x, t). \end{aligned} \quad (6.119)$$

Finally, substituting this result and equation (6.115) into the equation of motion (6.116) we obtain the dispersion relation

$$\omega^2 = \nu^2 + 2\pi G\Sigma |k|. \quad (6.120)$$

Notice that ω^2 is positive definite, so bending waves in the infinite sheet are always stable.

This result can be applied to tightly wound bending waves of the form

$$z(R, \phi, t) = \text{Re}[z_a(R) e^{im(\phi - \Omega_p t)}], \quad \text{where } z_a(R) = Z(R) e^{i \int k dR}, \quad (6.121)$$

$m \geq 0$, and $|kR| \gg 1$. Just as in the case of tightly wound density waves, the effect of self-gravity can be determined by noting that the tightly wound bending waves locally resemble the planar bending waves in the infinite sheet. Thus the vertical restoring force due to self-gravity is given by equation (6.119) as $g_z = -2\pi G\Sigma |k| z$. This extra force increases the total restoring force from $-\nu^2 z$ to $-(\nu^2 + 2\pi G\Sigma |k|)z$ and raises the vertical frequency from ν to $(\nu^2 + 2\pi G\Sigma |k|)^{1/2}$. It follows that the pattern speeds Ω_p of the fast and slow waves (eq. 6.109) change from $\Omega \pm \nu/m$ to $\Omega \pm (\nu^2 + 2\pi G\Sigma |k|)^{1/2}/m$.

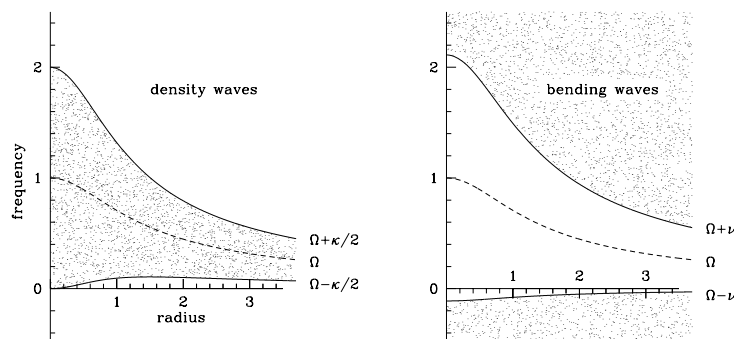


Figure 6.33 Stippling marks the regions where tightly wound $m = 2$ density waves (left) and $m = 1$ bending waves (right) exist in a flattened logarithmic potential. The potential has the form (6.111) with $v_0 = R_c = 1$ and $q_\Phi = 0.9$. The solid curves mark the Lindblad resonances on the left panel and the vertical resonances on the right panel, and the dashed curve marks the corotation resonance. For density waves, we assume that $Q = 1$ so there is no forbidden region around corotation (cf. Figure 6.14).

Using the relation $\omega = m\Omega_p$ that connects the frequency of a wave to its pattern speed, this result can be expressed as the dispersion relation

$$(\omega - m\Omega)^2 = \nu^2 + 2\pi G\Sigma|k|. \quad (6.122)$$

There are **vertical resonances** at $m(\Omega - \Omega_p) = \pm\nu$, which are analogous to the Lindblad resonances at $m(\Omega - \Omega_p) = \pm\kappa$.

Note the similarity to the dispersion relation for tightly wound density waves in a cold disk (eq. 6.64). The only differences are that (i) the epicycle frequency κ is replaced by the vertical frequency ν , and (ii) the minus sign in front of the term $2\pi G\Sigma|k|$ in the density-wave dispersion relation is replaced by a plus sign. The physical origin of this change of sign is that self-gravity increases the stiffness of the disk to bending waves, while it decreases the stiffness to density waves. The sign change implies that the regions in which tightly wound bending and density waves can propagate are quite different (see Figure 6.33).

Since the right side of equation (6.122) is always positive, ω is always real, and we conclude that *tightly wound bending waves in a cold disk are always stable*. This behavior is in sharp contrast to tightly wound density waves in a cold disk, which are violently unstable (page 495).

(c) The origin of warps There are many theories for the origin of warps:

Could galaxies be susceptible to a bending instability analogous to the bar instability which we have found to be so common? The answer seems to be “no.” A razor-thin, cold, self-gravitating disk embedded in an axisymmetric halo is stable to all $m = 0$ and $m = 1$ vertical perturbations, at least

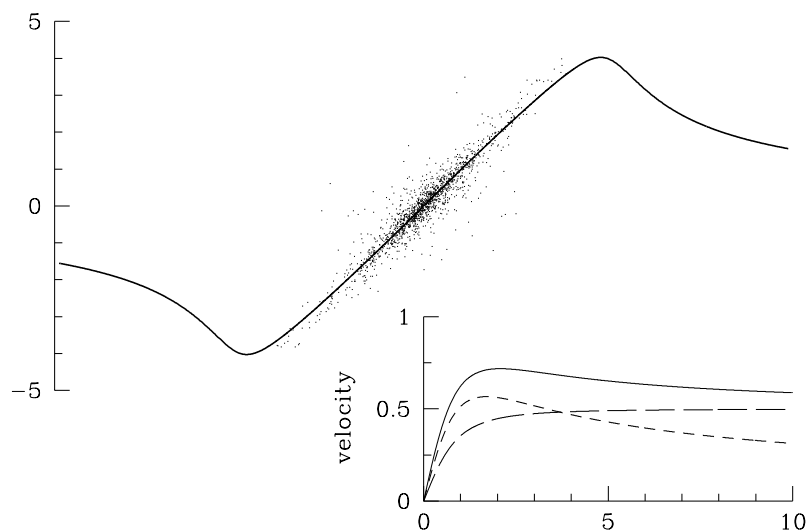


Figure 6.34 The Laplacian surface for a disk embedded in a halo. The dots mark the distribution of stars in a Miyamoto–Nagai disk (eq. 2.69) with mass $M = 1$, scale length $a = 1$, and $b/a = 0.2$. The disk is inclined by 45° to a dark halo having a logarithmic potential (eq. 2.71a) with axis ratio $q_\Phi = 0.9$, core radius $R_c = 1$, and asymptotic circular speed $v_0 = 0.5$. At small radii the Laplacian surface almost coincides with the midplane of the disk, while at large radii it follows the equatorial plane of the halo. The insert at bottom right shows the circular-speed curve that would obtain if the disk and halo were aligned: the short- and long-dashed curves show the circular speed due to the disk and halo potentials, and the solid curve shows the circular speed in the combined potential.

if the halo is spherical or flattened (Hunter & Toomre 1969). Since observed warps have $m = 1$ they should be stable.

Do warps represent a stable $m = 1$ bending mode, perhaps excited early in the lifetime of the galaxy? Numerical calculations of bending modes of galactic disks (Hunter & Toomre 1969; Sparke 1984) show that in most cases they damp out in a few rotation times, essentially because bending waves, like density waves, have a group velocity that carries them to the edge of the disk, where they are absorbed.²⁰ Just as in the case of density waves, the effects of self-gravity do not solve the winding problem for bending waves.

Consider next a disk that is embedded in a flattened axisymmetric halo, and imagine that the symmetry axis of the disk is tilted with respect to the symmetry axis of the halo. Let us assume for the moment that the halo is rigid. The inner part of the disk is so tightly coupled gravitationally that it also acts like a rigid body, so the inner disk remains flat, and its

²⁰ More precisely, the retrograde $m = 1$ waves propagate inward to the inner vertical resonance, where they reflect into leading waves, which then propagate out to the disk edge. See Problem 6.7.

Box 6.2: The Laplacian surface

Here we determine the shape of a cold gas disk orbiting in the combined gravitational field of a disk galaxy and a dark halo with misaligned symmetry axes. We approximate the disk galaxy and the dark halo as rigid bodies, and replace the gas disk by a system of concentric circular rings of negligible mass. The plane of the ring at radius r is inclined to the disk plane by an angle $i(r)$. Thus if $\hat{\mathbf{e}}(r)$ and $\hat{\mathbf{e}}_d$ are unit vectors pointing along the symmetry axes of the ring and the disk, we have $\cos i(r) = \hat{\mathbf{e}}(r) \cdot \hat{\mathbf{e}}_d$. The disk galaxy exerts a torque per unit mass $\mathbf{N}_d(r)$ on each ring, the magnitude of which depends on the radius r of the ring and the inclination $i(r)$, and the direction of which must be normal to both symmetry axes $\hat{\mathbf{e}}(r)$ and $\hat{\mathbf{e}}_d$. Thus we can write

$$\mathbf{N}_d = w_d(r, \hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_d) \hat{\mathbf{e}} \times \hat{\mathbf{e}}_d, \quad (1)$$

where $w_d(r, \hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_d)$ is straightforward to compute if we know the disk potential $\Phi_d(R, z)$. Similarly, the torque per unit mass from the dark halo on the ring at radius r has the form

$$\mathbf{N}_h = w_h(r, \hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_h) \hat{\mathbf{e}} \times \hat{\mathbf{e}}_h. \quad (2)$$

The ring rotates with angular speed equal to the circular frequency at that radius, $\Omega(r)$, so its angular momentum per unit mass is $\mathbf{L}(r) = \Omega(r)r^2\hat{\mathbf{e}}(r)$. The disk and halo torques cause it to precess at a rate

$$\frac{d\mathbf{L}}{dt} = L \frac{d\hat{\mathbf{e}}}{dt} = \mathbf{N}_d + \mathbf{N}_h = \hat{\mathbf{e}} \times [\hat{\mathbf{e}}_d w_d(r, \hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_d) + \hat{\mathbf{e}}_h w_h(r, \hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_h)]. \quad (3)$$

A coherent cold disk can survive only if all the rings precess at the same rate. This goal can be achieved if the net torque on each ring vanishes, $\mathbf{N}_d + \mathbf{N}_h = 0$. It is straightforward to show that this requires that $\hat{\mathbf{e}}$, $\hat{\mathbf{e}}_d$, and $\hat{\mathbf{e}}_h$ all lie in the same plane, and equation (3) can then be solved numerically for the orientation of the ring at radius r . The superposition of these rings at various radii yields the **Laplacian surface** (see Figure 6.34 and Problem 6.6). At small radii the gravitational field from the disk dominates, so $|w_d| \gg |w_h|$ and the Laplacian surface almost coincides with the disk midplane, while at large radii $|w_h| \gg |w_d|$ so the surface lies in the equatorial plane of the halo.

This simple derivation neglects several complicating factors. (i) We have also neglected the precession of the disk galaxy due to torques from the halo. (ii) The behavior of a fluid element is not the same as a rigid ring, although this approximation appears to work so long as the halo torques are not too strong. (iii) The most serious concern, as described in the text, is our approximation that the halo is a rigid body.

symmetry axis simply precesses around the symmetry axis of the halo at fixed inclination. The outer part of the disk, however, is only weakly coupled to the inner part, and therefore tends to settle into the symmetry plane of the halo (see Figure 6.34). The resulting surface, called the Laplacian surface (Box 6.2), can closely resemble observed warps (Toomre 1983; Sparke 1984).

A misaligned disk of this kind can be regarded as a special case of a bending wave, in the following sense. An isolated disk, or a disk embedded in a spherical halo, always has a trivial bending wave (the “tilt” mode), in which the symmetry axis of the disk is perturbed but the disk remains planar (in eq. 6.121, $z_a \propto R$, $m = 1$, $\Omega_p = 0$). If the halo is then slowly flattened, the tilt mode acquires a non-zero pattern speed and a warped shape, which is simply the Laplacian surface shown in Figure 6.34. In contrast to other bending modes, the tilt mode can survive indefinitely in a rigid halo.

The survival of the tilt mode in a realistic dark halo is a much more complicated issue. A halo composed of collisionless particles responds strongly to an embedded precessing disk. In particular, if the plane of the disk is initially inclined with respect to the equatorial plane of the halo, the inner halo quickly realigns itself so that its equatorial plane coincides with the disk. This rearrangement dramatically reduces the torques that the halo exerts on each ring of the disk, with the result that the disk warp associated with the tilt mode in a rigid halo is no longer a mode—the warp exhibits strong differential precession and rapidly winds up (Binney, Jiang, & Dutta 1998). On the other hand, if the tilt mode is driven by the torque that a misaligned outer halo exerts on the combined disk plus inner halo, the mode will only gradually decay by a combination of phase mixing and energy transfer to the halo (Shen & Sellwood 2006).

We conclude that the survival of warps depends strongly on the unknown properties of the surrounding halo, and that in many cases warps will damp or wind up rapidly. Although warps are likely to be short-lived, they are also relatively easy to excite: for a disk embedded in a spherical halo the tilt mode is neutrally stable, so this mode is likely to be nearly neutral—and hence excited to a large amplitude by small perturbations—even in a moderately flattened halo. This argument suggests that warps can be excited by almost any large-scale gravitational field that continually or repeatedly perturbs the disk, even if it is relatively weak.

Some warps might represent the response of the disk to a recent close encounter with a companion galaxy. For our Galaxy the natural candidates are the Large Magellanic Cloud (LMC) and the Sagittarius galaxy (§1.1.3). The orbits of both galaxies are known reasonably well from the kinematics of their tidal streamers (see BM §8.4.1 and §8.1.1c for the LMC, and §8.3.3 for Sagittarius). The LMC is more massive, but Sagittarius is closer, so their tidal fields are comparable in strength. The tidal fields may be amplified by collective effects within the Milky Way’s dark halo (Weinberg 1998). A traditional concern with this mechanism is that many warped galaxies do

not have nearby companions, but substructure in the dark halo could also be responsible.

An alternative model is that warps arise from the changing orientation of the dark halo (Ostriker & Binney 1989; Debattista & Sellwood 1999). Simulations of halo formation show that most halos are triaxial (§9.3.3). The stellar disk may be assumed to lie in one of the principal planes of the triaxial halo, normally the one perpendicular to the smallest axis, which we denote by $\hat{\mathbf{e}}_z$. The orientation of the principal axes is continually changing in response to the tidal torques and the infall of new halo material (§9.3.3). If the angular speed of the principal-axis frame relative to an inertial frame is $\boldsymbol{\Omega}$, then in the principal-axis frame the disk is subject to a Coriolis force $-2\boldsymbol{\Omega} \times \mathbf{v}$ (eq. 3.117). The z -component of this force, $-2v_c(R)(\boldsymbol{\Omega} \cdot \hat{\mathbf{e}}_R)\hat{\mathbf{e}}_z$, distorts the disk into an $m = 1$ warp. The amplitude of the warp depends on the restoring force from the disk and halo, and thus on the location in the disk and the flattening of the halo. For plausible rotation rates $|\boldsymbol{\Omega}|$ of a radian or so per 10 Gyr, warps are produced with properties similar to those observed (Jiang & Binney 1999).

The distinction between these mechanisms of damping and excitation is blurry: the dark halo of a galaxy is a complex, time-dependent system, containing both baryonic substructure (satellite galaxies) and dark-matter substructure, which is continually being erased by tidal disruption and replenished by new infalling substructure (§9.3.3c), while the orientation of the large-scale principal axes of the halo is continuously changing in response to infall and external torques. We know that the halo is a noisy environment and that a warp is the dominant response of a galaxy disk to a wide variety of gravitational excitations; but determining which particular source of gravitational noise is responsible for a given warp may be a difficult and unrewarding task.

6.6.2 Buckling instability

We stated in the last subsection that razor-thin, cold, self-gravitating disks are stable to all $m = 0$ and $m = 1$ bending waves, and to tightly wound bending waves with any m . Nevertheless, an out-of-plane or buckling instability can arise in some disks, and this instability is worth investigating both for insight into disk dynamics and as a possible explanation for the vertical structure of bars (§6.5.2c). The buckling instability can be analyzed exactly in the Kalnajs disks (Polyachenko 1977), but we shall take a more physical approach.

As a first step, we examine the behavior of an infinite razor-thin plane sheet of surface density Σ , which we take to be located initially in the plane $z = 0$ at the minimum of an external “halo” potential $\frac{1}{2}\nu^2 z^2$. The system is similar to the one analyzed in §6.6.1b, except that now we shall assume that the stars have a Maxwellian distribution of velocities in the plane $z = 0$, that

is, the velocity DF is $f(v_x, v_y) \propto \exp[\frac{1}{2}(v_x^2 + v_y^2)/\sigma^2]$. We now distort the sheet into a planar bending wave traveling in the x direction, of the form (6.115). We shall assume that the cohesive forces in the sheet are sufficiently strong that all of the stars, whatever their velocity, remain confined to the warped sheet, i.e., they have the same coordinate $z_s(x, t)$. In practice this vertical cohesion is provided by the self-gravity of the sheet, and requires that the frequency at which the stars rattle vertically within the sheet, $\sim (4\pi G\rho)^{1/2}$ where ρ is the mass density within the sheet, exceed the frequency at which the stars encounter the bending oscillations, $\omega - kv_x$. If the thickness h of the sheet is small enough, this condition is satisfied for the overwhelming majority of stars, since $\rho \sim \Sigma/h$.

The velocity v_y has no effect on the interaction of the star with a wave traveling in the x -direction, so we can neglect motion in the y -direction. Stars with different velocities v_x will experience different histories of vertical acceleration as they traverse the corrugations caused by the wave. The equation of motion of a star in this system is the same as equation (6.116), except that the Eulerian derivative $\partial/\partial t$ is replaced by the convective derivative $d/dt = \partial/\partial t + v_x\partial/\partial x$ (eq. F.8); thus,

$$\left(\frac{\partial^2}{\partial t^2} + 2v_x \frac{\partial^2}{\partial x \partial t} + v_x^2 \frac{\partial^2}{\partial x^2} \right) z_s(x, t) = g_z(x, t) + g_i - \nu^2 z_s(x, t), \quad (6.123)$$

where g_i is the internal cohesive acceleration felt by each star due to the gravity of the others. Newton's third law tells us that the sum over stars of the accelerations g_i must vanish, so we can eliminate this term by averaging the equation over all of the stars at a given location; denoting this average by $\langle \cdot \rangle$, we have $\langle g_i \rangle = 0$, $\langle v_x \rangle = 0$, $\langle v_x^2 \rangle = \sigma^2$, so

$$\left(\frac{\partial^2}{\partial t^2} + \sigma^2 \frac{\partial^2}{\partial x^2} \right) z_s(x, t) = g_z(x, t) - \nu^2 z_s(x, t). \quad (6.124)$$

Substituting from equations (6.115) and (6.119), we obtain the dispersion relation (Toomre 1966; Kulsrud, Mark, & Caruso 1971; Fridman & Polyachenko 1984)

$$\omega^2 = \nu^2 + 2\pi G\Sigma|k| - \sigma^2 k^2. \quad (6.125)$$

This dispersion relation is an interesting contrast to the WKB dispersion relation for axisymmetric density waves (eq. 6.66), which reads

$$\omega^2 = \kappa^2 - 2\pi G\Sigma|k| + v_s^2 k^2. \quad (6.126)$$

The vertical frequency due to the halo ν replaces the epicycle frequency κ , and the signs in front to the remaining two terms are reversed. Consequently the self-gravity of the disk tends to destabilize density waves but stabilizes bending waves, while horizontal dispersion or sound speed stabilizes density waves but destabilizes bending waves!

The dispersion relation (6.125) implies that bending waves are always unstable if the wavelength is sufficiently short ($|k|$ sufficiently large). The instability sets in when the centrifugal force that arises as stars pass over the corrugations exceeds the gravitational restoring force. This **buckling instability** is closely related to the **fire-hose instability**, in which a hanging hose begins to oscillate when a strong flow of water passes through it, and to the **Kelvin–Helmholtz instability**, which arises when two incompressible fluids, one floating atop the other, are in relative horizontal motion. For an isolated disk ($\nu^2 = 0$), all wavelengths $\lambda < \lambda_b \equiv \sigma^2/(G\Sigma)$ are unstable to bending oscillations. However, our analysis is valid only for wavelengths that are long compared to the actual thickness of the sheet that we have approximated as razor-thin. Thus an instability is present only if $k_b z_0 \lesssim 1$ where $k_b = 2\pi/\lambda_b$ and z_0 is a measure of the thickness. For an isothermal stellar disk with vertical velocity dispersion σ_z we may take $z_0 = \sigma_z^2/(2\pi G\Sigma)$ from equation (4.302c), so $k_b z_0 = (\sigma_z/\sigma)^2$. Thus we expect that isolated disks are unstable to buckling if the ratio σ_z/σ is small enough, that is, if the random velocities within the plane are much larger than those in the perpendicular direction.

This expectation can be tested by generalizing the analysis above from an infinite sheet of zero thickness to an infinite slab of stars with non-zero thickness. To be specific, we consider the system with equilibrium DF

$$f(z, \mathbf{v}) = \frac{\rho_0}{(2\pi)^{3/2} \sigma \sigma_z^{1/2}} \exp \left[-\frac{v_x^2 + v_y^2}{2\sigma^2} - \frac{v_z^2}{2\sigma_z^2} - \frac{\Phi(z)}{\sigma_z^2} \right]. \quad (6.127)$$

It is easy to see that the DF (6.127) is a function only of the integrals of motion— v_x , v_y , and $E_z = \frac{1}{2}v_z^2 + \Phi(z)$ —and thus, by the Jeans theorem, the DF is a stationary solution of the collisionless Boltzmann equation. The corresponding density distribution $\rho(z)$ is derived in Problem 4.21. The stability of a stellar system with this DF has been determined by Toomre (1966) and by Araki (1987), who find that the system is stable to buckling modes if and only if $\sigma_z/\sigma > 0.293$.

Merritt & Sellwood (1994) suggest that this stability criterion is not stringent enough when applied to realistic axisymmetric disks rather than infinite slabs, because the vertical restoring force due to self-gravity is smaller for long wavelengths. They argue that $\sigma_z/\sigma \gtrsim 0.6$ is a more accurate requirement for isolated disks. A disk embedded in a massive halo will be more stable than an isolated disk.

Is the solar neighborhood stable against buckling? The simple models we have discussed neglect rotation and assume that the velocity-dispersion tensor is isotropic in the disk plane, and neither assumption is valid in the solar neighborhood. However, if we simply take the RMS value of the dispersions in the radial and azimuthal directions to represent the horizontal dispersion σ , then the ratio $\sigma_z/\sigma \simeq 0.6$ in the solar neighborhood (BM Table

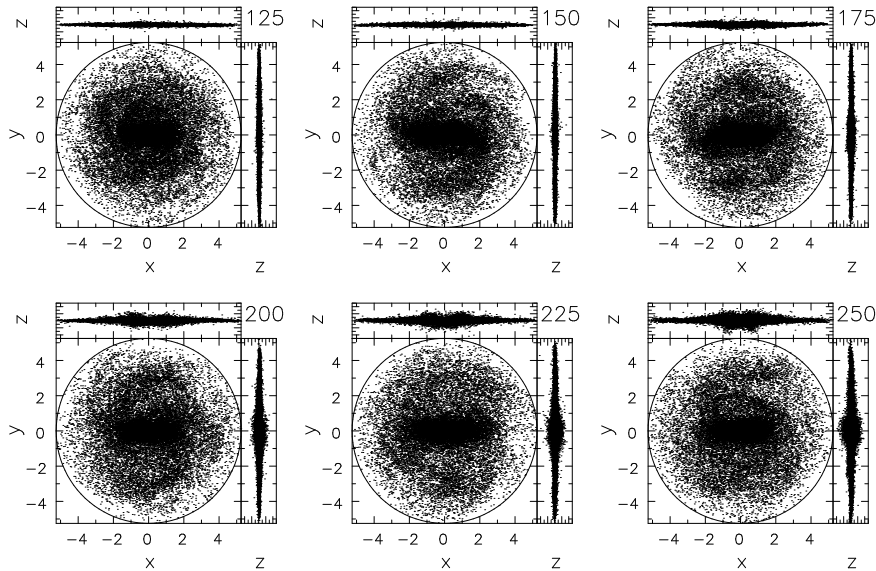


Figure 6.35 The buckling instability in a differentially rotating disk. This is the same simulation as in Figure 6.18, but shown at later times since the buckling instability develops after the bar instability. Each image has been rotated to align the bar with the x axis. Compare the boxy shape of the thickened disk in the edge-on views at late times with the boxy bulge of NGC 1381 in Figure 6.29. Courtesy of J. Sellwood.

10.2), close to the value that Merritt & Sellwood (1994) estimate is required for stability in an isolated disk. Since the (uncertain) halo contribution to the solar neighborhood dynamics (§6.3.4) also helps to stabilize the disk, it is likely that the solar neighborhood is safely stable to buckling.

The buckling instability does appear to play a role in the evolution of other stellar systems:

- (i) Bars formed in thin disks rapidly evolve into thick structures that appear boxy when viewed edge on (Figure 6.35 and §6.5.2c), probably because they are unstable to buckling. The rapid streaming motion along the major axis of a bar increases the effective horizontal dispersion σ and thus makes the bar more susceptible to the buckling instability than the axisymmetric disk from which it arose.
- (ii) Hot stellar systems are unstable to buckling if they are too flat. Since luminous elliptical galaxies are supported mainly by random motions rather than by rotation (§4.4.2c and BM §11.2.1) the buckling instability limits the flattening of these galaxies. Observations show that few, if any, elliptical galaxies have isophotal axis ratios less than 0.4 (i.e., there are almost no ellipticals flatter than E6 in Hubble's classification scheme; see BM §4.1.1), and this limit probably arises because flatter galaxies are unstable to buckling (Fridman & Polyachenko 1984).

Problems

6.1 [2] Show that the mass arm (the maximum of the surface density) in a tightly wound spiral leads the potential arm (the minimum of the gravitational potential) at a given radius by an angle

$$\Delta\phi = \frac{1}{km} \frac{d}{dR} \ln R^{1/2} |\Phi(R)|, \quad (6.128)$$

where $|\Phi(R)|$ is the amplitude of the spiral potential and the result has fractional error $O(|kR|^{-1})$. Thus the mass arm can either lead or lag the potential arm, depending on the radial dependence of the strength of the spiral. Hint: use equation (6.33).

6.2 [2] Consider a razor-thin disk containing some tracer population that satisfies the continuity equation. Assume that the surface density and mean line-of-sight velocity of the tracers are known at every point, that the inclination of the disk to the line of sight is known, and that the disk has a well-defined pattern speed Ω_p . Then Ω_p can be determined from equation (6.13). Show that the mean velocity of the tracers normal to the line of sight can also be determined at every point (Sridhar & Sambhus 2003).

6.3 [2] A useful model for exploring the properties of differentially rotating disks employs the softened point-mass potential (eq. 2.226). In this model the usual gravitational potential due to a particle of mass m at distance d , $-Gm/d$, is replaced by $-Gm/(d^2 + \epsilon^2)^{1/2}$, where ϵ is the softening length (Miller 1971).

(a) Consider a surface-density distribution in the $z = 0$ plane,

$$\Sigma_1(R, \phi) = \Sigma_a e^{i[m\phi + f(R)]}. \quad (6.129)$$

Argue that the softened potential at $z = 0$ due to this surface density distribution is equal to the usual Newtonian potential created by the same density distribution at a height $z = \epsilon$. Hence, if the density distribution is tightly wound, $|R \partial f / \partial R| \gg 1$, show by an extension of the arguments given in §6.2.2b that the softened potential due to the surface density (6.129) is

$$\Phi_\epsilon(R, \phi, z = 0) = -\frac{2\pi G e^{-|k|\epsilon}}{|k|} \Sigma_a e^{i[m\phi + f(R)]}, \quad \text{where } k = \frac{\partial f}{\partial R}. \quad (6.130)$$

(b) Show that the WKB dispersion relation for a cold disk with softened gravity is

$$(m\Omega - \omega)^2 = \kappa^2 - 2\pi G \Sigma |k| \mathcal{F}, \quad \text{where } \mathcal{F} = e^{-|k|\epsilon}. \quad (6.131)$$

The reduction factor \mathcal{F} due to softened gravity mimics the reduction factor due to velocity dispersion (eq. 6.61). Thus, cold disks with softened gravity provide close analogs to stellar disks that are much easier to investigate numerically.

(c) Show that a cold disk is stable to short-wavelength axisymmetric disturbances if

$$\epsilon > \frac{2\pi G \Sigma}{\kappa^2 e}. \quad (6.132)$$

6.4 [1] Show that the group velocity of tightly wound density waves in a fluid disk with $Q = 1$ is equal (within a sign) to the sound speed.

6.5 [2] For theoretical analyses, it is useful to modify the definition (6.7) of the pitch angle α to

$$\cot \alpha = \frac{kR}{m}; \quad (6.133)$$

thus trailing waves have pitch angle in the range $0 < \alpha < 90^\circ$ and leading waves have $90^\circ < \alpha < 180^\circ$. The rate of change of pitch angle in a wave packet is determined by the equation

$$\frac{d}{dt} \cot \alpha = \left. \frac{\partial \cot \alpha}{\partial R} \right|_\omega v_g = \left. \frac{\partial(kR)}{\partial R} \right|_\omega \frac{v_g}{m}, \quad (6.134)$$

where v_g is the group velocity. For a fluid Mestel disk with $Q = \text{constant}$, show that

$$\frac{d}{dt}(\cot \alpha) = \Omega_p. \quad (6.135)$$

6.6 [3] The goal of this problem is to compute the shape of the Laplacian surface for a simple model of a disk-halo system.

(a) We model the halo potential Φ_h by the logarithmic potential of equation (2.71) in the limit of large core radius R_c ; thus

$$\Phi_h = \frac{v_0^2}{2R_c^2} \left(R^2 + \frac{z^2}{q_\Phi^2} \right) + \text{constant}. \quad (6.136)$$

Show that the torque per unit mass from the halo on a ring of radius r with symmetry axis $\hat{\mathbf{e}}$ is given by equation (2) of Box 6.2, with

$$w_h(r, \hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_h) = \frac{v_0^2 r^2}{2R_c^2} \left(\frac{1}{q_\Phi^2} - 1 \right) \hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_h. \quad (6.137)$$

(b) We model the disk galaxy as an exponential disk, having mass M and scale length R_d . At distances exceeding a few scale lengths, the disk potential can be approximated by its monopole and quadrupole terms, given in equation (2.255). In this approximation, show that the torque from the disk on the ring is given by equation (1) of Box 6.2, with

$$w_d(r, \hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_d) = \frac{9GM R_d^2}{2r^3} \hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_d. \quad (6.138)$$

(c) Assume that the disk symmetry axis is tipped relative to the halo symmetry axis by an inclination I . Show that the inclination i of the Laplacian surface at radius r , relative to the disk, is given by

$$\frac{v_0^2 r^5}{9GM R_c^2 R_d^2} \left(\frac{1}{q_\Phi^2} - 1 \right) = \frac{\sin 2i}{\sin 2(I - i)}. \quad (6.139)$$

(d) Draw the Laplacian surface or a cross-section of it.

6.7 [1] Show from equation (6.122) that the group velocity of bending waves with pattern speed Ω_p is

$$v_g = -\text{sgn}(k) \frac{\pi G \Sigma}{m(\Omega - \Omega_p)}. \quad (6.140)$$

Hence show that $m = 1$ retrograde bending waves propagate inward when they are trailing, and outward when they are leading.

7

Kinetic Theory

So far we have concentrated on collisionless systems, in which the constituent particles move under the influence of the smoothed-out gravitational field generated by all the other particles. This approximation is not completely accurate. As described in §1.2, individual stellar encounters¹ gradually perturb stars away from the trajectories they would have taken if the gravitational field were perfectly smooth; in effect the stars diffuse in phase space away from their original orbits. After many such encounters the star eventually loses its memory of the original orbit, and finds itself on a wholly unrelated one. The characteristic time over which this loss of memory occurs is called the relaxation time t_{relax} ; over timescales exceeding t_{relax} the approximation of a smooth gravitational potential is incorrect.

The collisionless Boltzmann equation, which has been our main tool so far, is not valid when encounters are important. Thus we begin this chapter by reviewing general results about stellar systems that hold in the presence of encounters (§7.2 and §7.3). The equations that describe the behavior of stellar systems in the presence of encounters are derived in §7.4, and these are used to investigate the evolution of spherical stellar systems in §7.5.

¹We generally use the term “encounter” to denote the gravitational perturbation of the orbit of one star by another, and “collision” to denote actual physical contact between stars. However, to conform with common use, we use the terms “collisional” or “collisionless” to describe stellar systems in which encounters do or do not play a role.

For other discussions of the topics in this chapter, see Hénon (1973b), Spitzer (1987), and Heggie & Hut (2003).

7.1 Relaxation processes

The relaxation time is of order

$$t_{\text{relax}} \approx \frac{0.1N}{\ln N} t_{\text{cross}}, \quad (7.1)$$

where t_{cross} is the crossing time and N is the number of stars in the system (eq. 1.38). The relaxation time exceeds the crossing time if $N \gtrsim 40$. Galaxies typically have $N \approx 10^{11}$ and $t_{\text{cross}} \approx 100$ Myr, so the effects of stellar encounters can be ignored over a galaxy's lifetime of 10 Gyr. However, encounters have played a central role in determining the present structure of many other stellar systems, such as globular clusters ($N \approx 10^5$, $t_{\text{cross}} \approx 10^5$ yr, lifetime 10 Gyr), open clusters ($N \approx 10^2$, $t_{\text{cross}} \approx 1$ Myr, lifetime 100 Myr), the central parsec of galaxies ($N \approx 10^6$, $t_{\text{cross}} \approx 10^4$ yr, lifetime 10 Gyr), and the centers of clusters of galaxies ($N \approx 10^3$, $t_{\text{cross}} \approx 1$ Gyr, lifetime 10 Gyr).

The fundamental equations describing motion in a collisionless system of N stars of mass m are the collisionless Boltzmann and Poisson equations (eqs. 4.7 and 2.10),

$$\frac{\partial f}{\partial t} + [f, H] = 0 \quad ; \quad \nabla^2 \Phi(\mathbf{x}, t) = 4\pi GmN \int d^3\mathbf{v} f(\mathbf{x}, \mathbf{v}, t). \quad (7.2)$$

Here the Hamiltonian $H(\mathbf{x}, \mathbf{v}, t) = \frac{1}{2}v^2 + \Phi(\mathbf{x}, \mathbf{v}, t)$ and the DF $f(\mathbf{x}, \mathbf{v}, t)$ represents the probability that a given star is found in unit phase-space volume near the phase-space position (\mathbf{x}, \mathbf{v}) . In Chapter 4 we developed models of stellar systems by solving these equations exactly. For example, in spherical models such as the Hernquist model, the gravitational field is precisely time-independent and spherical, so each star conserves its energy and angular momentum. However, in any stellar system with finite N , the energy and angular momentum of individual stars are not precisely conserved, because each star is subject to fluctuating forces from encounters with its neighbors. Therefore the collisionless Boltzmann equation does not provide a complete description of the dynamics of stellar systems with finite N .

Encounters drive the evolution of a stellar system by several distinct mechanisms:

(a) Relaxation Each star slowly wanders away from its initial orbit. As a result of this phase-space diffusion, the entropy of the stellar system increases, and its structure becomes less sensitive to its initial conditions. We have seen in §4.10.1 that the high-entropy states of a self-gravitating gas are

very inhomogeneous, with a dense central core and an extended halo. Thus we expect that relaxation will drive stellar systems towards configurations having small, dense cores and large, low-density halos.

(b) Equipartition A typical stellar system contains stars with a wide range of masses. From elementary kinetic theory we know that encounters tend to produce equipartition of kinetic energy: on average, particles with large kinetic energy $\frac{1}{2}mv^2$ lose energy to particles with less kinetic energy. In an ordinary gas, this process leads to a state in which the mean-square velocity of a population of particles is inversely proportional to mass. By contrast, in a stellar system, massive stars that lose kinetic energy fall deeper into the gravitational potential well, and may even increase their kinetic energy as a result, just as an Earth satellite speeds up as it loses energy from atmospheric drag. Conversely, less massive stars preferentially diffuse towards the outer parts of the stellar system, where the velocity dispersion may be smaller.

(c) Escape From time to time an encounter gives a star enough energy to escape from the stellar system. Thus there is a slow but irreversible leakage of stars from the system, so stellar systems gradually evolve towards a final state consisting of only two stars in a Keplerian orbit, all the others having escaped to infinity. The timescale over which the stars “evaporate” in this way can be related to the relaxation timescale by the following simple argument (Ambarzumian 1938; Spitzer 1940). From equation (2.31) the escape speed v_e at \mathbf{x} is given by $v_e^2(\mathbf{x}) = -2\Phi(\mathbf{x})$. The mean-square escape speed in a system whose density is $\rho(\mathbf{x})$ is therefore

$$\langle v_e^2 \rangle = \frac{\int d^3\mathbf{x} \rho(\mathbf{x}) v_e^2(\mathbf{x})}{\int d^3\mathbf{x} \rho(\mathbf{x})} = -2 \frac{\int d^3\mathbf{x} \rho(\mathbf{x}) \Phi(\mathbf{x})}{M} = -\frac{4W}{M}, \quad (7.3)$$

where M and W are the total mass and potential energy of the system (eq. 2.18). According to the virial theorem (4.250), $-W = 2K$, where $K = \frac{1}{2}M\langle v^2 \rangle$ is the total kinetic energy. Hence

$$\langle v_e^2 \rangle^{1/2} = 2\langle v^2 \rangle^{1/2}; \quad (7.4)$$

in words, the RMS escape speed is just twice the RMS speed. The fraction of particles in a Maxwellian distribution that have speeds exceeding twice the RMS speed is $\gamma = 7.38 \times 10^{-3}$ (Problem 4.18). We can crudely assume that evaporation removes a fraction γ of the stars every relaxation time. Then the rate of loss is

$$\frac{dN}{dt} = -\frac{\gamma N}{t_{\text{relax}}} \equiv -\frac{N}{t_{\text{evap}}}, \quad (7.5)$$

where the **evaporation time**, the characteristic time in which the system’s stars are lost, is $t_{\text{evap}} = t_{\text{relax}}/\gamma \simeq 140 t_{\text{relax}}$. Thus we expect that any

stellar system will lose a substantial fraction of its stars in about $10^2 t_{\text{relax}}$ (see §7.5.2).

(d) Inelastic encounters So far we have treated stars as point masses, but in dense stellar systems we must consider the possibility that two stars occasionally pass so close that they raise powerful tides on one another or even suffer a physical collision. Energy dissipation in near-collisions reduces the total kinetic energy of the system and can lead to the formation of binary stars. Head-on or nearly head-on collisions can result in the coalescence of the colliding stars, leading to the otherwise unexpected presence of massive, short-lived stars in an old stellar system.

The characteristic timescale on which a star suffers a collision is given approximately by

$$t_{\text{coll}} \approx (n\Sigma v)^{-1}, \quad (7.6)$$

where n is the number density of stars, Σ is the collision cross-section, and v is the RMS stellar velocity. We may write $n \approx N/r^3$, where r is the radius of the system, and $\Sigma \approx \pi(2r_*)^2$, where r_* is the stellar radius (neglecting gravitational focusing; a more precise result is given in eq. 7.194). In terms of the crossing time $t_{\text{cross}} \approx r/v$,

$$\frac{t_{\text{coll}}}{t_{\text{cross}}} \approx \frac{r^2}{4\pi N r_*^2}. \quad (7.7)$$

From the virial theorem we have $v^2 \approx GNm/r$ where m is the stellar mass; it proves convenient to use this relation to eliminate r in favor of v . We also eliminate r_* in favor of the escape speed from the stellar surface, $v_* = \sqrt{2Gm/r_*}$ ($v_* = 618 \text{ km s}^{-1}$ for the Sun). Then

$$\frac{t_{\text{coll}}}{t_{\text{cross}}} \approx 0.02N \left(\frac{v_*}{v}\right)^4. \quad (7.8)$$

In terms of the relaxation time (eq. 7.1),

$$\frac{t_{\text{coll}}}{t_{\text{relax}}} \approx 0.2 \left(\frac{v_*}{v}\right)^4 \ln N. \quad (7.9)$$

For systems in which the escape speed from individual objects is much larger than the RMS orbital velocity (such as open and globular clusters, and most galaxies), we have $t_{\text{coll}} \gg t_{\text{relax}}$, so inelastic encounters play only a minor role in determining the overall structure of the stellar system. However, such encounters can occasionally produce exotic single or binary stars, which provide direct evidence of recent non-gravitational interactions.

(e) Binary formation by triple encounters A binary star cannot form in an isolated encounter of two point masses, because the relative motion is always along a hyperbola. However, an encounter involving three stars can

leave two of the participants in a bound Keplerian orbit. It is simple to estimate the rate of formation of binaries by this process. We showed in equation (1.30) that the velocity perturbation in an encounter of two stars of mass m and relative velocity v is $\delta v \approx Gm/bv$, where b is the distance of closest approach. We may rewrite this as

$$\frac{\delta v}{v} \approx \frac{b_{90}}{b}, \quad \text{where} \quad b_{90} \approx \frac{Gm}{v^2} \quad (7.10)$$

is the impact parameter at which the relative velocity is deflected by 90° in the encounter (see eq. 3.51 for a precise definition). If three stars approach one another within a distance b , we expect the velocity perturbations to be of similar magnitude. Thus, to form a binary by a triple encounter, we must have $\delta v \approx v$, which requires $b \approx b_{90}$. For a given star, the time interval between encounters with other stars at separation b_{90} or less is of order $(nb_{90}^2 v)^{-1}$ (eq. 7.6). In each such encounter, there is a probability nb_{90}^3 that a third star will also lie within a distance b_{90} . Hence the time t_3 required for a given star to suffer a triple encounter at separation less than b_{90} is $t_3 \approx (n^2 b_{90}^5 v)^{-1}$. Substituting for b_{90} from equation (7.10), we find the time required for a given star to become part of a binary by a triple encounter to be (Goodman & Hut 1993)

$$t_3 \approx \frac{v^9}{n^2 G^5 m^5}. \quad (7.11)$$

Using the virial theorem, $v^2 \approx GNm/r$, we may express t_3 in terms of the relaxation time (eq. 7.1):

$$\frac{t_3}{t_{\text{relax}}} \approx 10N^2 \ln N. \quad (7.12)$$

Hence the total number of binaries formed per relaxation time is only

$$\frac{Nt_{\text{relax}}}{t_3} \approx \frac{0.1}{N \ln N}. \quad (7.13)$$

Since the system dissolves after the evaporation time of about $10^2 t_{\text{relax}}$, the rate of binary formation by triple encounters is negligible if N is much larger than 10. We discuss binary formation and evolution further in §7.5.7.

(f) Interactions with primordial binaries The many binary stars found in the solar neighborhood were produced when their component stars were formed, rather than by subsequent triple or inelastic encounters. It is likely that binary stars are similarly produced during the formation of globular and open clusters. These are called **primordial binary stars** to distinguish them from binaries formed by dynamical processes long after

their constituent stars. Gravitational forces during encounters transfer energy between the orbits of primordial binaries and other cluster stars. Such energy exchange can dramatically alter the energy balance in the cluster, even if binaries are rare, because the binding energy in the binary orbit can be much larger than the kinetic energy of a typical cluster star. Consider, for example, a globular cluster with mass $M = 10^5 \mathcal{M}_\odot$ and RMS velocity $\langle v^2 \rangle^{1/2} = 10 \text{ km s}^{-1}$. From the virial theorem, its binding energy is $-\frac{1}{2}M\langle v^2 \rangle = 10^{50}$ erg. A binary star consisting of two $1 \mathcal{M}_\odot$ stars with a separation of $2 R_\odot$ has a binding energy of 1×10^{48} erg. Thus, 100 such binaries contain as much binding energy as the whole cluster of 10^5 stars.

7.2 General results

7.2.1 Virial theorem

In Chapter 4 we used the collisionless Boltzmann equation to prove the tensor virial theorem (eqs. 4.241 and 4.247),

$$\frac{1}{2} \frac{d^2 I_{jk}}{dt^2} = 2K_{jk} + W_{jk}, \quad (7.14)$$

which relates the tensor I_{jk} of an isolated stellar system to the kinetic- and potential-energy tensors, K_{jk} and W_{jk} . We now show that with slight modifications this result also holds for collisional systems.

Proof: Consider a system of particles with masses m_α and positions \mathbf{x}_α , $\alpha = 1, \dots, N$. We define the tensor (cf. eq. 4.243)

$$I_{jk} \equiv \sum_{\alpha=1}^N m_\alpha x_{\alpha j} x_{\alpha k}, \quad (7.15)$$

where $x_{\alpha j}$ is the j th Cartesian component of the vector \mathbf{x}_α . The second time derivative of I_{jk} is

$$\frac{d^2 I_{jk}}{dt^2} = \sum_{\alpha=1}^N m_\alpha (x_{\alpha j} \ddot{x}_{\alpha k} + 2\dot{x}_{\alpha j} \dot{x}_{\alpha k} + \ddot{x}_{\alpha j} x_{\alpha k}). \quad (7.16)$$

The acceleration of particle α is

$$\ddot{x}_{\alpha j} = \sum_{\substack{\beta=1 \\ \beta \neq \alpha}}^N \frac{Gm_\beta (x_{\beta j} - x_{\alpha j})}{|\mathbf{x}_\beta - \mathbf{x}_\alpha|^3}; \quad (7.17)$$

substituting this result and a similar formula for $\dot{x}_{\alpha k}$ into equation (7.16) we find

$$\begin{aligned} \frac{d^2 I_{jk}}{dt^2} &= 2 \sum_{\alpha=1}^N m_{\alpha} \dot{x}_{\alpha j} \dot{x}_{\alpha k} \\ &+ \sum_{\substack{\alpha, \beta=1 \\ \beta \neq \alpha}}^N \frac{G m_{\alpha} m_{\beta}}{|\mathbf{x}_{\alpha} - \mathbf{x}_{\beta}|^3} [x_{\alpha j} (x_{\beta k} - x_{\alpha k}) + x_{\alpha k} (x_{\beta j} - x_{\alpha j})]. \end{aligned} \quad (7.18)$$

By analogy with equation (4.240b), we define the kinetic-energy tensor for a system of point particles to be

$$K_{jk} \equiv \frac{1}{2} \sum_{\alpha=1}^N m_{\alpha} \dot{x}_{\alpha j} \dot{x}_{\alpha k}. \quad (7.19)$$

By analogy with equations (2.21a) and (2.22), we define the potential-energy tensor for a system of point particles as²

$$\begin{aligned} W_{jk} &= G \sum_{\substack{\alpha, \beta=1 \\ \beta \neq \alpha}}^N m_{\alpha} m_{\beta} \frac{x_{\alpha j} (x_{\beta k} - x_{\alpha k})}{|\mathbf{x}_{\alpha} - \mathbf{x}_{\beta}|^3} \\ &= -\frac{1}{2} G \sum_{\substack{\alpha, \beta=1 \\ \beta \neq \alpha}}^N m_{\alpha} m_{\beta} \frac{(x_{\alpha j} - x_{\beta j})(x_{\alpha k} - x_{\beta k})}{|\mathbf{x}_{\alpha} - \mathbf{x}_{\beta}|^3}, \end{aligned} \quad (7.20)$$

where the second line is obtained by interchanging the indices α and β in the first line and averaging this result with the first line. From the second line we conclude that \mathbf{W} is symmetric, that is, $W_{jk} = W_{kj}$. The second term on the right side of equation (7.18) is just $W_{jk} + W_{kj} = 2W_{jk}$, and the first term is $4K_{jk}$, so we have successfully arrived at equation (7.14). \triangleleft

The most useful form of the virial theorem is obtained by taking the trace of the tensor \mathbf{I} , $I \equiv \text{trace}(\mathbf{I}) \equiv \sum_{j=1}^3 I_{jj}$. Furthermore, we assume that the system is in a steady state, so $d^2 I / dt^2 = 0$. The trace of equation (7.18) then becomes the scalar virial theorem $2K + W = 0$ (eq. 4.248), where now

$$K = \text{trace}(\mathbf{K}) = \frac{1}{2} \sum_{\alpha=1}^N m_{\alpha} v_{\alpha}^2 \quad ; \quad W = \text{trace}(\mathbf{W}) = -\frac{1}{2} \sum_{\substack{\alpha, \beta=1 \\ \alpha \neq \beta}}^N \frac{G m_{\alpha} m_{\beta}}{|\mathbf{x}_{\alpha} - \mathbf{x}_{\beta}|} \quad (7.21)$$

² We can justify these analogies, at least formally, by replacing the continuous density $\rho(\mathbf{x})$ in (2.21a) and (2.22) by a sum of delta functions: $\rho(\mathbf{x}) = \sum_{\alpha=1}^N m_{\alpha} \delta(\mathbf{x} - \mathbf{x}_{\alpha})$.

are the total kinetic and potential energies.

The only approximation involved in deriving the scalar virial theorem is that I is time-independent. This is a good approximation for equilibrium stellar systems with $N \gg 1$, but in a system with a small number of particles there are statistical fluctuations in I , so the scalar virial theorem holds only for the time-averaged values of K and W .

7.2.2 Liouville's theorem

We have argued that the collisionless Boltzmann equation cannot provide a complete description of the dynamics of a stellar system with finite N . We now discuss a generalization of the collisionless Boltzmann equation that remedies this shortcoming, at least formally. We represent the state of a system of N stars by a point in a $6N$ -dimensional space, called Γ -**space**, whose coordinates are the positions and velocities of all the stars. This state is sometimes called a **microstate** and its representative point a Γ -**point**. In practice, we do not have—and do not want—the detailed information that is required to specify a microstate. We are concerned only with the “average” behavior of the macroscopic properties of the system (density distribution, velocity distribution at a given position, fraction of binary stars, etc.). Thus it is useful to imagine that at some initial time we are given the probability that a system is found in each small volume in Γ -space, and to follow the evolution of this probability distribution, rather than the evolution of a single Γ -point. There is an obvious analogy to the methods of Chapter 4, where we found it simpler to follow the evolution of the probability density in six-dimensional phase space, rather than the orbits of individual stars.

Denote the position and velocity of the α th particle by the canonical coordinates $\mathbf{q}_\alpha, \mathbf{p}_\alpha$, where $\alpha = 1, \dots, N$ (normally \mathbf{q}_α and \mathbf{p}_α are the position and velocity, but they could be any canonical coordinates and momenta). Then the six-dimensional vector $\mathbf{w}_\alpha \equiv (\mathbf{q}_\alpha, \mathbf{p}_\alpha)$ denotes the location of a particle in phase space. The Γ -point of a system in the $6N$ -dimensional Γ -space is determined by the collection of N six-vectors $\mathbf{w}_1, \dots, \mathbf{w}_N$. The probability that a Γ -point is found in a unit volume of Γ -space at time t is denoted by $f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N, t)$; since the probability density integrates to unity we have

$$\int d^6\mathbf{w}_1 \cdots d^6\mathbf{w}_N f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N, t) = 1, \quad \text{where} \quad d^6\mathbf{w}_\alpha \equiv d^3\mathbf{q}_\alpha d^3\mathbf{p}_\alpha. \quad (7.22)$$

The function $f^{(N)}$ is the **N-body distribution function** or N-body DF.

For the sake of simplicity, we shall usually assume that all the particles are identical (same mass, composition, etc.)—it is straightforward to modify the derivations below when several different kinds of particle are present.

Since the particles are identical, the N-body DF can be taken to be a symmetric function of $\mathbf{w}_1, \dots, \mathbf{w}_N$. In other words,

$$f^{(N)}(\dots, \mathbf{w}_\alpha, \dots, \mathbf{w}_\beta, \dots) = f^{(N)}(\dots, \mathbf{w}_\beta, \dots, \mathbf{w}_\alpha, \dots) \quad \text{for all } \alpha, \beta. \quad (7.23)$$

The equation governing the evolution of $f^{(N)}$ is analogous to the collisionless Boltzmann equation governing the evolution of the phase-space density f (§4.1). In fact, to derive the equation for $f^{(N)}$ we need only reinterpret the 3-dimensional vectors \mathbf{q} and \mathbf{p} in that section as $3N$ -dimensional vectors $(\mathbf{q}_1, \dots, \mathbf{q}_N)$, $(\mathbf{p}_1, \dots, \mathbf{p}_N)$. Then the analogs of equations (4.7)–(4.10) are

$$\frac{\partial f^{(N)}}{\partial t} + \sum_{\alpha=1}^N \left(\dot{\mathbf{q}}_\alpha \cdot \frac{\partial f^{(N)}}{\partial \mathbf{q}_\alpha} + \dot{\mathbf{p}}_\alpha \cdot \frac{\partial f^{(N)}}{\partial \mathbf{p}_\alpha} \right) = 0; \quad (7.24)$$

$$\frac{\partial f^{(N)}}{\partial t} + [f^{(N)}, H_N] = 0; \quad (7.25)$$

$$\frac{d f^{(N)}}{dt} = 0; \quad (7.26)$$

where d/dt is the convective derivative in Γ -space, and $[\cdot, \cdot]$ denotes the Poisson bracket in Γ -space. In other words the flow of Γ -points through Γ -space is incompressible: the probability density of Γ -points $f^{(N)}$ around the Γ -point of a given system always remains constant. This is **Liouville's theorem**, and equations (7.24)–(7.26) are **Liouville's equation**.³

If (i) we work in an inertial frame, (ii) we choose our canonical coordinates and momenta to be the positions \mathbf{x}_α and velocities \mathbf{v}_α , and (iii) our particles have mass m and interact only through their mutual gravitation, then Liouville's equation can be written in the form

$$\frac{\partial f^{(N)}}{\partial t} + \sum_{\alpha=1}^N \left(\mathbf{v}_\alpha \cdot \frac{\partial f^{(N)}}{\partial \mathbf{x}_\alpha} - \sum_{\substack{\beta=1 \\ \beta \neq \alpha}}^N \frac{\partial \Phi_{\alpha\beta}}{\partial \mathbf{x}_\alpha} \cdot \frac{\partial f^{(N)}}{\partial \mathbf{v}_\alpha} \right) = 0, \quad (7.27)$$

where $\Phi_{\alpha\beta} = -Gm/|\mathbf{x}_\alpha - \mathbf{x}_\beta|$.

Any N-body DF of the form

$$f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N) = f[H_N(\mathbf{w}_1, \dots, \mathbf{w}_N)] \quad (7.28)$$

³ We adopt the convention that the collisionless Boltzmann equation applies to 6-dimensional phase space and Liouville's equation applies to $6N$ -dimensional Γ -space, although some authors use the term Liouville's equation in both cases. With our convention, Liouville's equation is actually not due to Liouville. It was first explicitly derived by Gibbs (1884), two years after Liouville's death. Gibbs was also the first to recognize its importance in astronomy. It might therefore be better called the Gibbs equation.

is a solution of Liouville's equation. The proof is an obvious extension of the Jeans theorem (§4.2). In thermal equilibrium, we would have

$$f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N) = C \exp[-\beta H_N(\mathbf{w}_1, \dots, \mathbf{w}_N)], \quad (7.29)$$

where C and β are positive constants. Thermal equilibrium cannot be achieved in a gravitational N-body system because the normalization condition (7.22) cannot be satisfied for a DF of the form (7.29).⁴

7.2.3 Reduced distribution functions

We now investigate how the N-body DF $f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N, t)$ is related to the usual DF in six-dimensional phase space, $f(\mathbf{w}, t)$ (§4.1). We introduce first the **reduced** or **K-body distribution function**, which is obtained by integrating the N-body DF over $N - K$ of the six-vectors \mathbf{w}_α . Since $f^{(N)}$ is a symmetric function of the \mathbf{w}_α (eq. 7.23), without loss of generality we may choose the integration variables to be $\mathbf{w}_{K+1}, \dots, \mathbf{w}_N$. Thus we define

$$f^{(K)}(\mathbf{w}_1, \dots, \mathbf{w}_K, t) \equiv \int d^6 \mathbf{w}_{K+1} \cdots d^6 \mathbf{w}_N f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N, t). \quad (7.30)$$

From equation (7.22), the normalization of the K -body DF is simply

$$\int d^6 \mathbf{w}_1 \cdots d^6 \mathbf{w}_K f^{(K)}(\mathbf{w}_1, \dots, \mathbf{w}_K, t) = 1. \quad (7.31)$$

The one-body DF is

$$f^{(1)}(\mathbf{w}_1, t) \equiv \int d^6 \mathbf{w}_2 \cdots d^6 \mathbf{w}_N f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N, t). \quad (7.32)$$

The one-body DF describes the probability of finding a particular star in a unit volume of phase space centered on \mathbf{w}_1 . This is the same as the definition of the phase-space DF in §4.1, and therefore we are free to simplify our notation by writing

$$f(\mathbf{w}, t) = f^{(1)}(\mathbf{w}, t). \quad (7.33)$$

In many situations, it is useful to write the two-body DF in the form

$$f^{(2)}(\mathbf{w}_1, \mathbf{w}_2, t) = f(\mathbf{w}_1, t)f(\mathbf{w}_2, t) + g(\mathbf{w}_1, \mathbf{w}_2, t). \quad (7.34)$$

⁴ The integral diverges at both large and small scales. When the particles are separated by large distances, $f^{(N)}$ depends on velocity but is independent of position, so the spatial integral diverges. When two particles α and β approach one another, $\Phi_{\alpha\beta}$ diverges so $\exp(-\beta H_N)$ becomes extremely large.

The function g is called the **two-body correlation function**; the terminology is borrowed from probability theory, where the variables x and y are said to be uncorrelated if the joint probability $p(x, y)$ can be factored into a product of the form $p_x(x)p_y(y)$. Loosely speaking, the two-body correlation function measures the excess probability of finding a particle at \mathbf{w}_1 due to the presence of a particle at \mathbf{w}_2 . A more precise statement can be derived from the laws of conditional probability (eq. B.85), which state that the probability that a star is located in a unit volume of phase space centered on \mathbf{w}_1 , given that a star is known to be located at \mathbf{w}_2 , is

$$f(\mathbf{w}_1|\mathbf{w}_2) = \frac{f^{(2)}(\mathbf{w}_1, \mathbf{w}_2)}{\int d^6\mathbf{w}'_1 f^{(2)}(\mathbf{w}'_1, \mathbf{w}_2)} = \frac{f(\mathbf{w}_1)f(\mathbf{w}_2) + g(\mathbf{w}_1, \mathbf{w}_2)}{f(\mathbf{w}_2) + \int d^6\mathbf{w}'_1 g(\mathbf{w}'_1, \mathbf{w}_2)}. \quad (7.35)$$

In particular, if the correlation function $g(\mathbf{w}_1, \mathbf{w}_2) = 0$, then $f(\mathbf{w}_1|\mathbf{w}_2) = f(\mathbf{w}_1)$; in other words, the presence of a star at \mathbf{w}_2 has no effect on the probability of finding a star near \mathbf{w}_1 .

The use of reduced DFs can be illustrated by computing the expectation value of the kinetic and potential energy for a stellar system. From equation (7.21), the expectation of the kinetic energy is

$$\langle K \rangle = \frac{1}{2}m \int d^6\mathbf{w}_1 \cdots d^6\mathbf{w}_N f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N, t) \sum_{\alpha=1}^N v_\alpha^2; \quad (7.36)$$

since the stars are identical, this simplifies to

$$\langle K \rangle = \frac{1}{2}Nm \int d^6\mathbf{w}_1 f(\mathbf{w}_1, t) v_1^2. \quad (7.37)$$

Similarly, *any* observable that involves only quantities that depend additively on the phase-space coordinates of single stars can be expressed in terms of the one-body DF. Such observables include density, surface brightness, line-of-sight velocity distribution, metallicity distribution, etc.

The expectation of the potential energy is

$$\langle W \rangle = -\frac{1}{2} \int d^6\mathbf{w}_1 \cdots d^6\mathbf{w}_N f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N, t) \sum_{\substack{\alpha, \beta=1 \\ \alpha \neq \beta}}^N \frac{Gm^2}{|\mathbf{x}_\alpha - \mathbf{x}_\beta|}. \quad (7.38)$$

Since the stars are identical, and there are $N(N-1)$ ways in which we can choose two distinct stars α and β from N , this simplifies to

$$\langle W \rangle = -\frac{1}{2}Gm^2N(N-1) \int d^6\mathbf{w}_1 d^6\mathbf{w}_2 \frac{f^{(2)}(\mathbf{w}_1, \mathbf{w}_2, t)}{|\mathbf{x}_2 - \mathbf{x}_1|}. \quad (7.39)$$

Thus the potential energy depends only on the two-body DF. If the correlation function is small, that is, if $|g(\mathbf{w}_1, \mathbf{w}_2, t)| \ll f(\mathbf{w}_1)f(\mathbf{w}_2)$, then for $N \gg 1$ the potential energy simplifies to

$$W = -\frac{1}{2}Gm^2N^2 \int d^6\mathbf{w}_1 d^6\mathbf{w}_2 \frac{f(\mathbf{w}_1, t)f(\mathbf{w}_2, t)}{|\mathbf{x}_2 - \mathbf{x}_1|} = \frac{1}{2} \int d^3\mathbf{x} \rho(\mathbf{x})\Phi(\mathbf{x}), \quad (7.40)$$

which is the expression we have used in prior chapters (eq. 2.18).

7.2.4 Relation of Liouville's equation to the collisionless Boltzmann equation

The N-body DF is said to be **separable** if it is simply the product of one-body DFs, that is, if

$$f^{(N)}(\mathbf{w}_1, \dots, \mathbf{w}_N, t) = \prod_{\beta=1}^N f(\mathbf{w}_\beta, t). \quad (7.41)$$

As we have seen, this assumption implies that the positions of stars are uncorrelated, in the sense that the probability of finding a star near any phase-space position \mathbf{w}_1 is unaffected by the presence or absence of stars at nearby points. We now assume that the N-body DF is separable, and ask for the equation governing the evolution of the one-body DF f . To find this, we integrate Liouville's equation (7.27) over $d^6\mathbf{w}_2 \cdots d^6\mathbf{w}_N$. The term involving $\partial f^{(N)}/\partial t$ simply yields $\partial f(\mathbf{w}_1, t)/\partial t$. The term involving $\partial f^{(N)}/\partial \mathbf{x}_\alpha$ yields zero if $\alpha = 2, \dots, N$ because $\int d^3\mathbf{x}_\alpha \partial f^{(N)}/\partial \mathbf{x}_\alpha = 0$ so long as $f^{(N)} \rightarrow 0$ sufficiently fast as $|\mathbf{x}_\alpha| \rightarrow \infty$. The integration of the term involving $\partial f^{(N)}/\partial \mathbf{v}_\alpha$ yields zero if $\alpha = 2, \dots, N$ for a similar reason. Thus we obtain

$$\begin{aligned} & \frac{\partial f(\mathbf{w}_1, t)}{\partial t} + \mathbf{v}_1 \cdot \frac{\partial f(\mathbf{w}_1, t)}{\partial \mathbf{x}_1} \\ & - \frac{\partial f(\mathbf{w}_1, t)}{\partial \mathbf{v}_1} \cdot \sum_{\beta=2}^N \int d^6\mathbf{w}_2 \cdots d^6\mathbf{w}_N \frac{\partial \Phi_{1\beta}}{\partial \mathbf{x}_1} \prod_{\alpha=2}^N f(\mathbf{w}_\alpha, t) = 0. \end{aligned} \quad (7.42)$$

Each term in the sum is identical, and $\int d^6\mathbf{w} f(\mathbf{w}, t) = 1$, so this becomes

$$\frac{\partial f(\mathbf{w}_1, t)}{\partial t} + \mathbf{v}_1 \cdot \frac{\partial f(\mathbf{w}_1, t)}{\partial \mathbf{x}_1} - (N-1) \frac{\partial f(\mathbf{w}_1, t)}{\partial \mathbf{v}_1} \cdot \int d^6\mathbf{w}_2 \frac{\partial \Phi_{12}}{\partial \mathbf{x}_1} f(\mathbf{w}_2, t) = 0. \quad (7.43)$$

The expectation value of the gravitational potential at \mathbf{x}_1 is

$$\bar{\Phi}(\mathbf{x}_1, t) = N \int d^6\mathbf{w}_2 \Phi_{12} f(\mathbf{w}_2, t), \quad (7.44)$$

so equation (7.43) simplifies to

$$\frac{\partial f(\mathbf{w}, t)}{\partial t} + \mathbf{v} \cdot \frac{\partial f(\mathbf{w}, t)}{\partial \mathbf{x}} - \frac{N-1}{N} \frac{\partial \bar{\Phi}(\mathbf{x}, t)}{\partial \mathbf{x}} \frac{\partial f(\mathbf{w}, t)}{\partial \mathbf{v}} = 0. \quad (7.45)$$

In the limit $N \rightarrow \infty$ this becomes the collisionless Boltzmann equation (4.11). Thus we have shown that *the collisionless Boltzmann equation results from the Liouville equation when $N \gg 1$ and the N -body DF is separable.*

If the DF is *not* separable and $N \gg 1$, it is straightforward to show that equation (7.45) must be replaced by

$$\frac{df}{dt} = \Gamma[f], \quad (7.46)$$

where $\Gamma[f]$ is the **encounter operator**, given by

$$\Gamma[f(\mathbf{w}_1, t)] \equiv N \int d^6 \mathbf{w}_2 \frac{\partial \Phi_{12}}{\partial \mathbf{x}_1} \cdot \frac{\partial g(\mathbf{w}_1, \mathbf{w}_2, t)}{\partial \mathbf{v}_1}, \quad (7.47)$$

and g is the two-body correlation function (eq. 7.34). Thus the correlations between particles in phase space, as measured by $g(\mathbf{w}_1, \mathbf{w}_2, t)$, drive the rate of change of the phase-space density around a given star, given by $\Gamma[f]$.

We can determine the encounter operator $\Gamma[f]$ in two ways. The first approach is through the correlation function. Just as we derived equation (7.46) for the one-body DF by integrating Liouville's equation over $d^6 \mathbf{w}_2 \cdots d^6 \mathbf{w}_N$, we can derive an equation for the correlation function—or, what is equivalent, the two-body DF $f^{(2)}(\mathbf{w}_1, \mathbf{w}_2, t)$ —by integrating Liouville's equation over $d^6 \mathbf{w}_3 \cdots d^6 \mathbf{w}_N$. Unfortunately, just as equation (7.46) for the one-body DF depends on the two-body DF through the appearance of $g(\mathbf{w}_1, \mathbf{w}_2, t)$ on the right side, the equation for the two-body DF depends on the three-body DF.⁵ However, in the limit where the number of stars $N \rightarrow \infty$, while the total mass mN remains constant, we can neglect the contribution of the three-body correlation function to the equation governing the two-body DF. The resulting equation can be solved to determine the two-body correlation function, and this can be substituted into equation (7.47) to determine the encounter operator (Gilbert 1968; Lifshitz & Pitaevskii 1981).

A more physical approach, which we take in §7.4, is to ask how encounters between stars modify the one-body DF. This approach is only practical when the encounters can be approximated as localized in both time and space; fortunately, we shall see that this approximation is remarkably accurate for most stellar systems.

⁵ Continuing in this way, we would obtain a sequence of equations of rapidly increasing complexity, each expressing the rate of change of $f^{(n)}$ in terms of $f^{(n+1)}$. This sequence is known as the **BBGKY hierarchy**, after N. N. Bogoliubov, M. Born and H. S. Green, J. G. Kirkwood, and J. Yvon, who all discovered the equations independently between 1935 and 1946.

7.3 The thermodynamics of self-gravitating systems

7.3.1 Negative heat capacity

By analogy with an ideal gas, we can define the temperature T of a self-gravitating system at position \mathbf{x} through the relation (eq. F.33)

$$\frac{1}{2}m\overline{v^2} = \frac{3}{2}k_{\text{B}}T, \quad (7.48)$$

where m is the stellar mass and k_{B} is Boltzmann's constant. In general the mean-square velocity $\overline{v^2}$ and hence the temperature depend on position. The mass-weighted mean temperature is $\overline{T} \equiv \int d^3\mathbf{x} \rho(\mathbf{x})T / \int d^3\mathbf{x} \rho(\mathbf{x})$, where $\rho(\mathbf{x})$ is the density, and the total kinetic energy of a system of N identical stars is therefore

$$K = \frac{3}{2}Nk_{\text{B}}\overline{T}. \quad (7.49)$$

If the system is stationary, then the virial theorem (4.250) states that the total energy including gravitational potential energy is $E = -K$, so

$$E = -\frac{3}{2}Nk_{\text{B}}\overline{T}. \quad (7.50)$$

The heat capacity of the system is

$$C \equiv \frac{dE}{d\overline{T}} = -\frac{3}{2}Nk_{\text{B}}, \quad (7.51)$$

which is *negative*: by losing energy the system grows hotter (Lynden–Bell & Wood 1968; Lynden–Bell 1999).

This apparently paradoxical result is not restricted to stellar dynamics. *Any* bound, finite system in which the dominant forces are gravitational has negative heat capacity. In fact, the stability of nuclear burning in the cores of stars like the Sun is a consequence of this property: the reaction rates are strongly increasing functions of density and temperature, so if the reactions proceed so fast that the heat generated cannot be conducted away, thermal energy accumulates in the core, which expands and cools, bringing the reaction rates back towards equilibrium.

The thermodynamics of a system with negative heat capacity is quite different from that of normal laboratory systems. For example, suppose that we place a self-gravitating fluid or stellar system in contact with a heat bath (ignoring the practical difficulties of doing this). Assume that initially the bath and the system have the same temperature T . If a small amount of heat $dQ > 0$ is transferred to the bath, the temperature of the system will change to $\overline{T} - dQ/C > \overline{T}$. The self-gravitating system is now hotter than the bath, and heat continues to flow from hot to cold, that is, from the system to the bath. Therefore, the temperature of the system rises without limit. Similarly, if heat begins to flow from the bath to the system, the temperature of the system will decrease to zero. Thus any system with negative heat capacity in contact with a heat bath is thermodynamically unstable.

Box 7.1: Statistical mechanics of stellar systems

Many familiar results from statistical mechanics do not apply to systems with long-range forces such as stellar systems (for reviews see Padmanabhan 1990; Lynden–Bell 1999; Katz 2003). For example:

Energy is not extensive. In most laboratory systems, the total energy is an extensive property; that is, if the system is divided into parts then its energy is the sum of the energies of the parts. However, in a gravitating system of particles with mass m and number density n the potential energy between a particle and its neighbors within a sphere of radius R is $W = -Gm^2n \int d^3\mathbf{x}/r = -4\pi Gm^2n \int_0^R dr r = -2\pi Gm^2nR^2$. Thus most of the contribution to the potential comes from distant particles, so the energy is not extensive.

The microcanonical ensemble does not exist. The probability distribution of configurations of a closed system at fixed energy is derived by assuming that all states on the energy hypersurface in phase space have equal probability. However, the energy hypersurface of an isolated self-gravitating system is unbounded (except for $N = 2$, $E < 0$), so the microcanonical probability distribution cannot be defined.

The canonical ensemble does not exist. The probability distribution of states of a system that is in contact with a heat bath at temperature T is derived by maximizing the entropy at fixed mean energy. However, we have shown in §4.10.1 that an isolated self-gravitating system has no maximum-entropy state.

The heat capacity at constant volume is not positive. The usual “proof” of this result relies on the relation $\langle(\Delta T)^2\rangle = T^2/C_V$ (Landau & Lifshitz 1980) between the mean-square fluctuations in temperature and the heat capacity at constant volume; this formula relies on the canonical distribution and so cannot be applied to stellar systems.

7.3.2 The gravothermal catastrophe

To investigate the thermodynamic behavior of self-gravitating systems more rigorously, we consider a self-gravitating ideal gas of N point particles⁶ of mass m , having total mass $M = Nm$. The gas is enclosed by a spherical container of radius r_b . The gas is assumed to be thermally conducting, so heat can flow from one part of the gas to another in response to temperature gradients. The equilibrium state of this system is therefore isothermal, and in fact is precisely the isothermal sphere described in §4.3.3b, truncated at the radius of the wall r_b .⁷ Hence the density distribution of the gas in the

⁶ We assume that the particles have no internal degrees of freedom, and that binary systems and other short-range correlations between the particles are unimportant.

⁷ By Newton’s first theorem the absence of material outside the container does not

container, $\rho(r)$, is given by Figure 4.6 for $r < r_b$, while the pressure is related to the density by the ideal gas law (F.31),

$$p(r) = \frac{\rho(r)k_B T}{m} = \frac{\rho(r)}{m\beta}, \quad \text{where} \quad \beta \equiv \frac{1}{k_B T} \quad (7.52)$$

is called the **inverse temperature** (though it has dimensions of inverse energy).

This system provides a highly idealized model of a stellar system; its principal advantage, in addition to mathematical simplicity, is that it does *not* contain most of the processes that drive the evolution of stellar systems—there is no equipartition, escape, inelastic encounters, binary formation, etc.—so we can isolate the effects of relaxation.

We choose the arbitrary constant in the definition of the gravitational potential $\Phi(r)$ so that $\Phi(r_b) = -GM/r_b$; thus if the container has negligible mass, $\Phi \rightarrow 0$ as $r \rightarrow \infty$. The total energy of the gas is $E = K + W$, where the kinetic energy $K = \frac{3}{2}Nk_B T = \frac{3}{2}M/(m\beta)$ (eq. 7.49) and the potential energy $W = 2\pi \int_0^{r_b} dr r^2 \rho(r) \Phi(r)$ (eq. 2.23). We could determine W by evaluating this integral using the density distribution of the isothermal sphere, but this task can be simplified by a trick. In Problem 4.33 we derived a form of the virial theorem valid for a collisionless system confined to a spherical container. The same theorem holds for a gaseous system and reads

$$2K + W = E + K = 4\pi r_b^3 p(r_b). \quad (7.53)$$

Hence

$$E = 4\pi r_b^3 p(r_b) - K = \frac{4\pi r_b^3 \rho(r_b)}{m\beta} - \frac{3M}{2m\beta}. \quad (7.54)$$

We eliminate the inverse temperature β in favor of the King radius r_0 and central density ρ_0 of the isothermal sphere, using equations (4.100) and (4.106), which yield

$$\beta = \frac{9}{4\pi G m \rho_0 r_0^2}. \quad (7.55)$$

We then eliminate the central density in favor of the mass M by the relation

$$M = 4\pi \rho_0 r_0^3 \int_0^{\tilde{r}_b} d\tilde{r} \tilde{r}^2 \tilde{\rho}(\tilde{r}) \equiv 4\pi \rho_0 r_0^3 \tilde{M}(\tilde{r}_b), \quad (7.56)$$

where $\tilde{r} \equiv r/r_0$, $\tilde{r}_b \equiv r_b/r_0$, $\tilde{\rho} \equiv \rho/\rho_0$, and the second equality in equation (7.56) defines a dimensionless mass \tilde{M} . We can therefore rewrite the inverse temperature in the form

$$\beta = \frac{9r_0}{Gm} \frac{\tilde{M}(\tilde{r}_b)}{M} = 9 \frac{r_b}{GMm} \frac{\tilde{M}(\tilde{r}_b)}{\tilde{r}_b}. \quad (7.57)$$

affect the gravitational field inside the container.

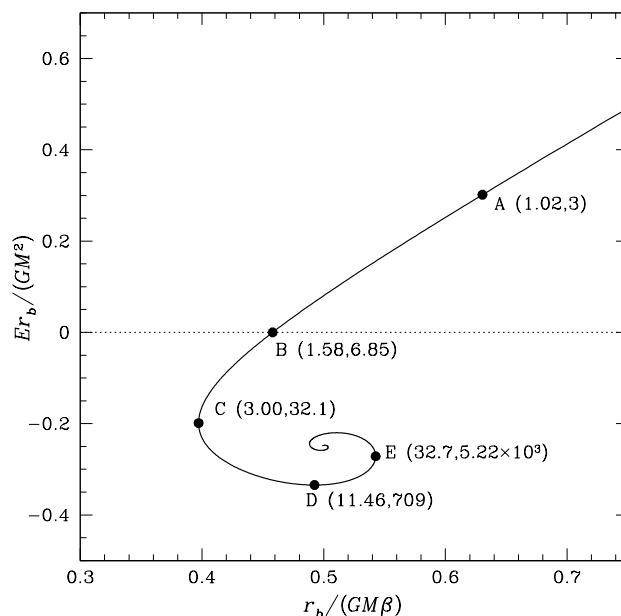


Figure 7.1 The dimensionless temperature $r_b/(GMm\beta) = k_B T r_b/(GMm)$ and dimensionless energy Er_b/GM^2 for a mass M of isothermal gas in a spherical container of radius r_b . The points A–E are labeled by the dimensionless box radius $\tilde{r}_b = r_b/r_0$ and by the density contrast $\mathcal{R} = \rho_0/\rho(r_b)$ (see Problem 7.6). The curve spirals inward to the point $(\frac{1}{2}, -\frac{1}{4})$ corresponding to the singular isothermal sphere.

Using this result to eliminate β from equation (7.54), we have

$$E = \frac{GM^2}{r_b} \left[\frac{\tilde{r}_b^4 \tilde{\rho}(\tilde{r}_b)}{9M^2(\tilde{r}_b)} - \frac{\tilde{r}_b}{6M(\tilde{r}_b)} \right], \quad (7.58)$$

where $\tilde{\rho}(\tilde{r})$ is determined by the differential equation (4.107a).

For a container of a given radius r_b containing a given mass M , we can use equations (7.57) and (7.58) to determine the inverse temperature β and the energy E as functions of the parameter \tilde{r}_b and hence as functions of one another. The result is shown in Figure 7.1. We have plotted the dimensionless ratios $r_b/(GMm\beta)$ and Er_b/GM^2 so the graph can be used for containers of any radius r_b containing any gas mass M . The numbers by the labels A–E give the dimensionless box radius \tilde{r}_b as well as

$$\mathcal{R} \equiv \frac{\rho_0}{\rho(r_b)} = \frac{1}{\tilde{\rho}(\tilde{r}_b)}, \quad (7.59)$$

which measures the density contrast between the center and the edge of the container.

We now perform a thought experiment. Suppose that the walls of the container conduct heat, and that the container is surrounded by a heat bath at very high temperature, so the gravitational potential energy of the gas is small compared with the kinetic energy in random motions (far above point A on the curve in Figure 7.1). Under these conditions, the gas behaves almost like an ideal gas with no self-gravity: the energy $E \simeq \frac{3}{2}M/(m\beta) = \frac{3}{2}Mk_{\text{B}}T/m$, the heat capacity $C = dE/dT \simeq \frac{3}{2}Mk_{\text{B}}/m$ is positive, and the gas is nearly homogeneous, so the density contrast \mathcal{R} is near unity. We now slowly reduce the temperature of the heat bath. The gas remains isothermal, energy flows from the gas to the bath, and the system moves down along the curve in Figure 7.1 to lower temperature and energy. As it passes point B ($\mathcal{R} = 6.8$) its total energy passes through zero and becomes negative, although its heat capacity, measured by the slope of the curve in Figure 7.1, remains positive. As we continue to reduce the temperature of the heat bath, the system continues to lose energy and cool, its heat capacity $C = dE/dT$ becoming larger and larger, and finally at point C ($\mathcal{R} = 32.125$) the heat capacity becomes infinite. There is *no* equilibrium state cooler than $T(C) = 0.40 GMm/(k_{\text{B}}r_{\text{b}})$. Systems between points C and D ($32.125 < \mathcal{R} < 708.61$) have negative heat capacity and are unstable for the reasons given at the end of §7.3.1: if at any instant the gas is hotter than the heat bath, energy flows from the gas to the bath; because of its negative heat capacity the gas becomes hotter as it loses energy; the increased temperature difference leads to even faster energy loss to the heat bath; and the gas is heated without limit. Similarly, if the gas is momentarily cooler than the bath, then energy flows into the gas, which cools without limit as a result.

In systems between points D and E ($708.61 < \mathcal{R} < 5221.5$) the heat capacity is once again positive. Nevertheless, it turns out that these systems are also unstable when in contact with a heat bath (Horwitz & Katz 1978; Katz 1978, 1979). In fact, *an isothermal gas in a conducting spherical container in contact with a heat bath is unstable whenever the density contrast between center and edge exceeds $\mathcal{R} = 32.125$.*

Now perform a second thought experiment. Suppose that our container is surrounded by a thermally insulating wall (i.e., the stars bounce off the wall without loss or gain of energy). Initially the gas has energy E , mass M , and radius $r_{\text{b}0}$. Suddenly the container is expanded to a radius r_{b} . Since the expansion is sudden, the gas does no work on the container, and thus E is constant. After the expansion, when the gas has again settled into equilibrium, its temperature can be determined from Figure 7.1. If the energy of the gas is positive, then the value of the vertical coordinate Er_{b}/GM^2 is increased by a factor $r_{\text{b}}/r_{\text{b}0}$, so the system moves along the curve towards the upper right. The value of \mathcal{R} decreases towards unity and the gas is therefore more homogeneous. However, if the energy is negative, then Er_{b}/GM^2 becomes more negative by a factor of $r_{\text{b}}/r_{\text{b}0}$. If the final value of Er_{b}/GM^2 is more negative than -0.335 (the value at point D in Figure 7.1), then no equilibrium is possible after the expansion.

Point D has an even deeper significance than this. Stability analysis shows that a thermally isolated ($E = \text{constant}$) sphere is unstable at all points in the equilibrium sequence of Figure 7.1 beyond point D (Antonov 1962a; Lynden-Bell & Wood 1968; Horwitz & Katz 1978; Katz 1978, 1979). In other words, *an isothermal gas in an insulating spherical container is unstable if the density contrast between center and edge exceeds $\mathcal{R} = 708.61$* . It can be shown that this instability, named the **gravothermal catastrophe** by Lynden-Bell & Wood (1968), arises because the isothermal sphere with $\mathcal{R} > 708.61$ is a local entropy *extremum* (i.e., a saddle point) but not a *maximum* at fixed E , M , and r_b . Thus the unstable system can reach states of higher entropy by evolving away from isothermality.

The onset of instability can be heuristically explained in terms of the virial theorem. The halo has positive heat capacity C_h since it is not strongly influenced by self-gravity, while the core, which is confined primarily by self-gravity, has negative heat capacity C_c . If the core momentarily becomes hotter than the halo, heat flows from the core to the halo, and the temperatures of *both* the core and halo rise. If $C_h < |C_c|$, the halo temperature rises more than the core temperature and the heat flow is shut off. If $C_h > |C_c|$, the halo has so much thermal inertia that it cannot heat up as fast as the core, and the temperature difference between core and halo grows. Of course, the division into a separate core and halo is artificial, but this argument captures the essence of the instability that sets in at $\mathcal{R} = 708.61$.⁸

Is there a gravothermal catastrophe in isothermal stellar systems as well as in gaseous ones? The answer is yes. The velocity distribution in an isothermal stellar system is the same as in an isothermal gas (§4.3.3b), and hence the entropy of an isothermal stellar system is the same as that of an isothermal gas with the same temperature and density distribution. If we can approximate the stellar encounters as local and instantaneous (see §7.4.2 for a discussion of the validity of this approximation), then Boltzmann's H-theorem tells us that the entropy of the stellar system cannot decrease with time (Lifshitz & Pitaevskii 1981). Hence instability will arise when the equilibrium state becomes a local entropy extremum other than a maximum, just as it does for a gaseous system. The gravothermal catastrophe in a gas develops through heat conduction and hence the growth time is comparable to the thermal diffusion time. The analog of the diffusion time in a stellar system is the relaxation time t_{relax} , so the gravothermal catastrophe in stellar systems should develop on a timescale of order t_{relax} . These arguments have been confirmed by numerical calculations of the normal modes of an isothermal stellar system (Inagaki 1980; Ipser & Kandrup 1980).

⁸ A crude version of this argument was given already by Landau (1932), who argued that quantum-mechanical effects led eventually to the existence of a stable equilibrium, and in this way derived the maximum mass limit for a degenerate star (the Chandrasekhar limit). Toy models that illustrate many features of the gravothermal catastrophe are described by Aronson & Hansen (1972), Lynden-Bell & Lynden-Bell (1977), and Padmanabhan (1990).

To investigate the relevance of these thermodynamic considerations to stellar systems, we must first develop tools that enable us to follow the evolution of a stellar system over timescales longer than the relaxation time.

7.4 The Fokker–Planck approximation

7.4.1 The master equation

Under the influence of the smooth potential $\Phi(\mathbf{x})$, the DF $f(\mathbf{x}, \mathbf{v}, t)$ obeys the collisionless Boltzmann equation $df/dt = 0$, where the derivative is taken along the phase-space path of a given star (eq. 4.10); in words, the phase-space probability density of stars around a given star always remains the same. When encounters are taken into account, the phase-space density around a star changes with time, at a rate determined by the encounter operator $\Gamma[f]$ (eq. 7.46). A mathematically precise—but complicated—expression for the encounter operator is given by equation (7.47); in this section we use physical arguments to derive simpler expressions that are more powerful tools for studying the evolution of stellar systems.

Let $\Psi(\mathbf{w}, \Delta\mathbf{w})d^6(\Delta\mathbf{w})\Delta t$ be the probability that a star with the phase-space coordinates $\mathbf{w} = (\mathbf{q}, \mathbf{p})$ is scattered to the volume of phase space $d^6(\Delta\mathbf{w})$ around $\mathbf{w} + \Delta\mathbf{w}$ during the short time interval Δt . The **transition probability** Ψ includes the effects of encounters with other stars but not acceleration by the smooth potential of the stellar system, since the latter is accounted for already in the collisionless Boltzmann equation. To distinguish the star whose trajectory we are following from the stars doing the scattering, we follow the convention of §1.2.1 and call the former the subject star and the latter the field stars.

As a result of encounters, subject stars are scattered out of a unit volume of phase space centered on \mathbf{w} at a rate

$$\left. \frac{\partial f(\mathbf{w})}{\partial t} \right|_- = -f(\mathbf{w}) \int d^6(\Delta\mathbf{w}) \Psi(\mathbf{w}, \Delta\mathbf{w}). \quad (7.60)$$

Encounters also scatter subject stars into this volume, at a rate

$$\left. \frac{\partial f(\mathbf{w})}{\partial t} \right|_+ = \int d^6(\Delta\mathbf{w}) \Psi(\mathbf{w} - \Delta\mathbf{w}, \Delta\mathbf{w}) f(\mathbf{w} - \Delta\mathbf{w}). \quad (7.61)$$

The sum $(\partial f/\partial t)_- + (\partial f/\partial t)_+$ equals the encounter operator $\Gamma[f]$. Hence we arrive at the **master equation**

$$\frac{df}{dt} = \Gamma[f] = \int d^6(\Delta\mathbf{w}) [\Psi(\mathbf{w} - \Delta\mathbf{w}, \Delta\mathbf{w}) f(\mathbf{w} - \Delta\mathbf{w}) - \Psi(\mathbf{w}, \Delta\mathbf{w}) f(\mathbf{w})]. \quad (7.62)$$

This simple derivation masks a subtle point. The master equation is *not* time-reversible: a DF that is localized near a single point in phase space spreads outward under the influence of the encounter operator, but an extended distribution cannot shrink to a point. On the other hand, Liouville’s equation (7.26), which provides an exact description of the stellar system, *is* time-reversible. Irreversibility sneaks into our derivation through Boltzmann’s celebrated assumption of “molecular chaos”—the assumption that the distributions of subject and field stars are statistically independent. This is equivalent to the assumption that the transition probability $\Psi(\mathbf{w}, \Delta\mathbf{w})$ and the DF $f(\mathbf{w})$ are statistically independent, so the scattering rates in equations (7.60) and (7.61) can be written as products of Ψ and f .

7.4.2 Fokker–Planck equation

In the crude estimate of the relaxation time presented in §1.2.1, we found that, per crossing time, encounters give rise to a mean-square velocity perturbation (eq. 1.36)

$$\Delta v^2 \approx \frac{8v^2}{N} \ln \Lambda, \quad \text{where} \quad \Lambda \approx \frac{R}{b_{90}}. \quad (7.63)$$

Here N is the number of stars in the system and v is the star’s velocity. This result arose from integrating over all impact parameters between the 90° deflection radius $b_{90} \approx Gm/v^2$ (eq. 7.10) and the system’s characteristic radius R . The contribution to Δv^2 from impact parameters in any interval (b_1, b_2) can be obtained by simply replacing $\ln(R/b_{90})$ in (7.63) by $\ln(b_2/b_1)$. Thus *equal octaves of impact parameter contribute equally to Δv^2* ; in other words, encounters with impact parameters in the range b_{90} to $2b_{90}$, $2b_{90}$ to $4b_{90}$, and so forth, right up to the interval $\frac{1}{2}R$ to R , are all of equal importance for the relaxation process.

The virial theorem implies that $v^2 \approx GmN/R$ where m is the stellar mass; thus $R/b_{90} \approx N$. For the systems considered in this book, N is generally large— 10^5 for a globular cluster and 10^{10} for an elliptical galaxy, corresponding to 17 to 33 octaves in R/b_{90} . This large number enables us to make several critical approximations that dramatically simplify the study of the evolution of stellar systems due to encounters:

(a) Weak encounters The fractional velocity change in an encounter is $\delta v/v \approx b_{90}/b$ (eq. 7.10). Since equal octaves contribute equally to the scattering, and most octaves between b_{90} and R have $b \gg b_{90}$, it follows that *most of the scattering is due to weak encounters, that is, ones with $\delta v \ll v$* .

The dominance of weak encounters allows us to simplify the encounter operator. For weak encounters, $|\Delta\mathbf{w}|$ is small, and we can expand the first

term of equation (7.62) in a Taylor series

$$\begin{aligned} \Psi(\mathbf{w} - \Delta\mathbf{w}, \Delta\mathbf{w})f(\mathbf{w} - \Delta\mathbf{w}) &= \Psi(\mathbf{w}, \Delta\mathbf{w})f(\mathbf{w}) \\ &- \sum_{i=1}^6 \Delta w_i \frac{\partial}{\partial w_i} [\Psi(\mathbf{w}, \Delta\mathbf{w})f(\mathbf{w})] \\ &+ \frac{1}{2} \sum_{i,j=1}^6 \Delta w_i \Delta w_j \frac{\partial^2}{\partial w_i \partial w_j} [\Psi(\mathbf{w}, \Delta\mathbf{w})f(\mathbf{w})] + O(\Delta\mathbf{w}^3). \end{aligned} \quad (7.64)$$

The **Fokker–Planck approximation** consists of truncating this series after the second-order terms, so (7.62) becomes

$$\Gamma[f] = - \sum_{i=1}^6 \frac{\partial}{\partial w_i} \{D[\Delta w_i]f(\mathbf{w})\} + \frac{1}{2} \sum_{i,j=1}^6 \frac{\partial^2}{\partial w_i \partial w_j} \{D[\Delta w_i \Delta w_j]f(\mathbf{w})\}, \quad (7.65)$$

where $D[\Delta w_i]$ denotes the expectation of the change in w_i per unit time,

$$D[\Delta w_i] \equiv \int d^6(\Delta\mathbf{w}) \Psi(\mathbf{w}, \Delta\mathbf{w}) \Delta w_i. \quad (7.66)$$

Similarly,

$$D[\Delta w_i \Delta w_j] \equiv \int d^6(\Delta\mathbf{w}) \Psi(\mathbf{w}, \Delta\mathbf{w}) \Delta w_i \Delta w_j. \quad (7.67)$$

The quantities $D[\Delta w_i]$ and $D[\Delta w_i \Delta w_j]$ are known as **diffusion coefficients** since they characterize the rate at which stars diffuse through phase space as a result of encounters. The use of square brackets in the notation is a reminder that the diffusion coefficient $D[\Delta w_i \Delta w_j]$ is *not* a function of the variable $\Delta w_i \Delta w_j$. Rather, it is an average of $\Delta w_i \Delta w_j$ over $\Delta\mathbf{w}$, and a function of the position in phase space where the average is taken.⁹

Equation (7.65) can be extended to a higher order of approximation by including diffusion coefficients arising from the third- and higher-order terms of the Taylor series (7.64), but these are generally smaller, by the factor $\ln \Lambda$ defined in equation (7.63), and may be neglected (see Appendix L for more detail). It might also be thought that $D[\Delta w_i \Delta w_j]$ is much less than $D[\Delta w_i]$, but Δw_i fluctuates in sign, whereas $(\Delta w_i)^2$ is always positive, and in fact these two coefficients are generally comparable in size (§7.4.3).

The second-order diffusion coefficient $D[\Delta w_i \Delta w_j]$ governs the rate at which the subject star executes a random walk in phase space, analogous to the Brownian motion of microscopic particles. The first-order diffusion

⁹Note also that $D[(\Delta\mathbf{w})^2]$ is *not* the same as $D[\Delta(\mathbf{w}^2)]$. Since $\Delta(\mathbf{w}^2) = (\mathbf{w} + \Delta\mathbf{w})^2 - \mathbf{w}^2 = 2\mathbf{w} \cdot \Delta\mathbf{w} + (\Delta\mathbf{w})^2$, the two diffusion coefficients are related by $D[\Delta(\mathbf{w}^2)] = 2\mathbf{w} \cdot D[\Delta\mathbf{w}] + D[(\Delta\mathbf{w})^2]$.

coefficient $D[\Delta w_i]$ represents a steady drift through phase space, rather than a random walk.

Equations (7.46) and (7.65) together constitute the **Fokker–Planck equation**. The Fokker–Planck equation has the virtue that all of the dependence on the field-star DF is contained in the diffusion coefficients, which are functions only of the phase-space coordinates of the subject star. Thus the Fokker–Planck equation is a differential equation, rather than an integro-differential equation like the master equation, and hence is much easier to solve. For this reason, the Fokker–Planck equation has become the principal tool for the study of encounters in stellar systems.¹⁰

The Fokker–Planck equation is reminiscent of, but distinct from, the **diffusion equation** in phase space, which reads

$$\frac{\partial f}{\partial t} = \sum_{i,j=1}^6 \frac{\partial}{\partial w_i} \left(C_{ij} \frac{\partial}{\partial w_j} f(\mathbf{w}) \right), \quad (7.68)$$

where the tensor $C_{ij}(\mathbf{w})$ is also usually called the diffusion coefficient.

(b) Local encounters Equation (7.63) shows that equal octaves in impact parameter contribute equally to gravitational scattering. Since most octaves between b_{90} and R have $b \ll R$, it follows that *most of the scattering is due to short-range or local encounters, that is, ones with $b \ll R$.*

The dominance of local encounters helps to justify the assumption of molecular chaos discussed in the preceding subsection: two stars that have suffered a local encounter with impact parameter $b \ll R$ are unlikely to have another encounter with $b \ll R$ in the lifetime of the stellar system. Even if they do have another local encounter, this will occur only after many orbits and perturbations by many other stars, and all memory of their previous encounter will have been erased.

Other important consequences follow from the dominance of local encounters: (i) since the encounter time $\sim b/v$ is short—much less than the crossing time—the encounter affects the velocity only, not the position, of the interacting stars; (ii) during the encounter the stars may be assumed to move on Keplerian hyperbolae, unaffected by the large-scale potential of the stellar system; (iii) the effects of stellar encounters on a star at \mathbf{x} can be calculated as if the star were embedded in an infinite homogeneous medium in which the DF is everywhere equal to the DF at \mathbf{x} .

Because local encounters affect velocity only, the Fokker–Planck encounter operator (7.65) is simplified if we choose the canonical phase-space coordinates \mathbf{w} to be Cartesian coordinates (\mathbf{x}, \mathbf{v}) . Then $\Psi(\mathbf{w}, \Delta \mathbf{w})$ is zero unless $\Delta \mathbf{x} = \mathbf{0}$, and as a consequence any diffusion coefficient of the form

¹⁰ Studies that do not rely on the Fokker–Planck approximation are described in Goodman (1983) and references therein.

$D[\Delta x_i]$, $D[\Delta x_i \Delta x_j]$ or $D[\Delta x_i \Delta v_j]$ is zero. Thus the encounter operator simplifies to

$$\Gamma[f] = - \sum_{i=1}^3 \frac{\partial}{\partial v_i} \{D[\Delta v_i]f(\mathbf{w})\} + \frac{1}{2} \sum_{i,j=1}^3 \frac{\partial^2}{\partial v_i \partial v_j} \{D[\Delta v_i \Delta v_j]f(\mathbf{w})\}. \quad (7.69)$$

Notice that $D[\Delta v_i]$ is a component of a vector and $D[\Delta v_i \Delta v_j]$ is a component of a tensor.

(c) Orbit-averaging For $N \gg 1$ the relaxation time in stellar systems is much larger than the crossing time. Thus, changes in the DF caused by encounters are expected to be small over a single orbital period. Hence it is useful to separate the slow changes in the phase-space coordinates caused by encounters from the rapid changes associated with orbital motion in the smooth potential; in effect, we **orbit-average** the Fokker–Planck equation.

Orbit averaging is most easily understood by working in angle-action variables. The Fokker–Planck equation in angle-action coordinates can be written as

$$\frac{\partial f}{\partial t} + j_i \frac{\partial f}{\partial J_i} + \dot{\theta}_i \frac{\partial f}{\partial \theta_i} = \Gamma[f]; \quad (7.70)$$

here, and henceforth, we adopt the summation convention (page 772). The time derivatives \dot{J}_i and $\dot{\theta}_i$ refer to motion in the smooth potential (that is, neglecting encounters), and the encounter operator $\Gamma[f]$ is given by equation (7.65), where we have chosen as phase-space coordinates $\mathbf{w} = (\boldsymbol{\theta}, \mathbf{J})$.

Since the evolution due to encounters is slow (i.e., since $t_{\text{relax}} \gg t_{\text{cross}}$), the DF is always approximately in a steady state, so the smooth potential of the stellar system evolves slowly with time. In most cases of interest (e.g., spherical systems) the potential at any instant will be integrable, so the corresponding Hamiltonian depends on the actions only, not the angles, $H(\boldsymbol{\theta}, \mathbf{J}, t) \simeq H(\mathbf{J}, t)$. Similarly, the strong Jeans theorem tells us that the DF depends only on the actions, $f(\boldsymbol{\theta}, \mathbf{J}, t) \simeq f(\mathbf{J}, t)$. Since the Hamiltonian depends only on the actions, $\dot{J}_i = -\partial H / \partial \theta_i$ vanishes, and $\dot{\theta}_i = \partial H / \partial J_i = \text{constant}$.

We now orbit average, by operating on (7.70) with $(2\pi)^{-3} \int d^3\boldsymbol{\theta}$. Since all quantities are periodic in $\boldsymbol{\theta}$, all terms involving $\partial / \partial \theta_i$ in the encounter operator vanish, as does the term $\dot{\theta}_i (\partial f / \partial \theta_i)$. Thus equation (7.70) can be written as

$$\begin{aligned} \frac{\partial f}{\partial t} &= \frac{1}{(2\pi)^3} \int d^3\boldsymbol{\theta} \Gamma[f] \\ &= - \frac{\partial}{\partial J_i} \left\{ f \overline{D}[\Delta J_i] \right\} + \frac{1}{2} \frac{\partial^2}{\partial J_i \partial J_j} \left\{ f \overline{D}[\Delta J_i \Delta J_j] \right\} \end{aligned} \quad (7.71a)$$

where the **orbit-averaged diffusion coefficients** are

$$\overline{D}[\Delta J_i] = \frac{1}{(2\pi)^3} \int d^3\boldsymbol{\theta} D[\Delta J_i], \quad (7.71b)$$

with a similar definition for $\overline{D}[\Delta J_i \Delta J_j]$.

An alternative form of equation (7.71) provides additional insight. Let \mathcal{V} be some volume in action space, and \mathcal{S} its surface. The number of stars in \mathcal{V} is $N_{\mathcal{V}}(t) = N \int_{\mathcal{V}} d^3 \mathbf{J} \int d^3 \boldsymbol{\theta} f(\mathbf{J}, t) = (2\pi)^3 N \int_{\mathcal{V}} d^3 \mathbf{J} f(\mathbf{J}, t)$, and we have

$$\frac{dN_{\mathcal{V}}}{dt} = (2\pi)^3 \int_{\mathcal{V}} d^3 \mathbf{J} \frac{\partial f}{\partial t} = - \int_{\mathcal{V}} d^3 \mathbf{J} \frac{\partial \mathcal{F}_i}{\partial J_i}, \quad (7.72)$$

where

$$\mathcal{F}_i \equiv (2\pi)^3 \left(f \overline{D}[\Delta J_i] - \frac{1}{2} \frac{\partial}{\partial J_j} \left\{ f \overline{D}[\Delta J_i \Delta J_j] \right\} \right). \quad (7.73)$$

The physical meaning of \mathcal{F}_i is seen by applying the divergence theorem—equation (B.43), except in action space rather than real space—to equation (7.72):

$$\frac{dN_{\mathcal{V}}}{dt} = - \oint_{\mathcal{S}} d^2 S_i \mathcal{F}_i, \quad (7.74)$$

where $d^2 \mathbf{S}$ is an element of the surface \mathcal{S} in action space and $d^2 S_i$ is one of its components. Thus \mathcal{F}_i a component of a vector that represents the flux of stars in action space due to encounters.

The principal advantage of orbit averaging is that the domain of the Fokker–Planck equation is reduced from six phase-space coordinates plus time to three actions plus time.

7.4.3 Fluctuation-dissipation theorems

Although we have shown that isolated stellar systems in thermal equilibrium do not exist, many galaxies and star clusters have DFs that are not far from the Maxwellian DF that is characteristic of thermodynamic equilibrium. Thus it is useful to examine the properties that the diffusion coefficients would have in this state.

Consider a stellar system containing two types of star: subject stars of mass m and field stars of mass m_a . In thermodynamic equilibrium, the velocity dispersions σ and σ_a of the two types of stars are related by energy equipartition; thus

$$m\sigma^2 = m_a\sigma_a^2 \equiv \beta^{-1}, \quad (7.75)$$

where β is the inverse temperature (eq. 7.52). In thermodynamic equilibrium, the one-body DF of the subject stars is Maxwellian,

$$f(\mathbf{w}) \propto e^{-\beta m H(\mathbf{J})}, \quad (7.76)$$

where H is the Hamiltonian for a particle of unit mass, and the system does not evolve, so the flux of stars through action space due to encounters must vanish. Thus equation (7.73) yields

$$f \overline{D}[\Delta J_i] = \frac{1}{2} \frac{\partial}{\partial J_j} \left(f \overline{D}[\Delta J_i \Delta J_j] \right); \quad (7.77)$$

as usual, there is an implicit summation over the index j from 1 to 3. After substituting equation (7.76) for the DF, we have

$$\overline{D}[\Delta J_i] = \frac{1}{2} \frac{\partial}{\partial J_j} \overline{D}[\Delta J_i \Delta J_j] - \frac{1}{2} m \beta \Omega_j \overline{D}[\Delta J_i \Delta J_j], \quad (7.78)$$

where $\Omega(\mathbf{J}) \equiv \partial H / \partial \mathbf{J}$ (eq. 3.190).

Since the diffusion coefficients depend on the actions of the subject star and the DF of the field stars, but not the DF of the subject stars, this relation must hold *whether or not* the subject stars are actually in thermal equilibrium, so long as the field stars have a Maxwellian DF with inverse temperature β .

The relation (7.78) is an example of a **fluctuation-dissipation theorem**; these theorems can take a wide variety of forms and are a powerful tool in statistical mechanics (see, for example, Pathria 1972; Landau & Lifshitz 1980; Nelson & Tremaine 1999). Such relations are useful in part because $\overline{D}[\Delta J_i \Delta J_j]$ can often be calculated using only first-order perturbation theory, while $\overline{D}[\Delta J_i]$ can require analyzing the dynamics to second-order.

In the limit where the subject star has zero mass, equation (7.78) simplifies to

$$\overline{D}[\Delta J_i] = \frac{1}{2} \frac{\partial}{\partial J_j} \overline{D}[\Delta J_i \Delta J_j]. \quad (7.79)$$

Note that this result is independent of the inverse temperature β . Since the diffusion coefficients are linear functions of the field star DF, we can superimpose the diffusion coefficients for Maxwellian DFs with different temperatures β to determine the diffusion coefficient for any ergodic DF, that is, any DF of the form $f(H)$. Since the relation (7.79) holds for each of the Maxwellians, it must hold for their sum. Thus *equation (7.79) relates the orbit-averaged diffusion coefficients for a zero-mass particle in any stellar system with an ergodic DF*. Moreover, in this case the orbit-averaged Fokker–Planck equation (7.71) simplifies to

$$\frac{\partial f}{\partial t} = \frac{1}{2} \frac{\partial}{\partial J_i} \left(\overline{D}[\Delta J_i \Delta J_j] \frac{\partial f}{\partial J_j} \right), \quad (7.80)$$

which is simply the diffusion equation (7.68) in action space.

In an infinite homogeneous system, the Jeans swindle (§5.2.2) allows us to assume that the unperturbed orbits of stars have constant velocity. In this case the Hamiltonian $H = \frac{1}{2} v^2$ and the velocity vector \mathbf{v} behaves like an action vector \mathbf{J} in a general potential, in that both are constant on unperturbed orbits. Moreover orbit-averaging is not needed because the medium is homogeneous. Thus the analog of the fluctuation-dissipation theorem (7.78) for a Maxwellian distribution of field stars is

$$D[\Delta v_i] = \frac{1}{2} \frac{\partial}{\partial v_j} D[\Delta v_i \Delta v_j] - \frac{1}{2} m \beta v_j D[\Delta v_i \Delta v_j]; \quad (7.81)$$

and the analog of (7.79) for zero-mass subject stars is

$$D[\Delta v_i] = \frac{1}{2} \frac{\partial}{\partial v_j} D[\Delta v_i \Delta v_j], \quad (7.82)$$

which holds for any field-star DF that is isotropic in velocity space.

7.4.4 Diffusion coefficients

We now evaluate the diffusion coefficients $D[\Delta v_i]$ and $D[\Delta v_i \Delta v_j]$ in the approximation that encounters are local. The diffusion coefficients represent mean changes per unit time due to a large number of encounters of the subject star with the field stars. Each encounter is assumed to be independent of all the others, and to involve only a single pair of stars (i.e., triple and multiple encounters are neglected). The effect of the overall gravitational potential of the stellar system is neglected, so the relative orbit of the two stars is a hyperbola, whose shape is determined by the relative velocity V_0 at large separations and the impact parameter b . All of these approximations are justified because the scattering is dominated by local encounters, as discussed in §7.4.2.

The diffusion coefficients are functions of the position \mathbf{x} and velocity \mathbf{v} of the subject star. In Appendix L we show that they can be written as

$$\begin{aligned} D[\Delta v_i] &= 4\pi G^2 m_a (m + m_a) \ln \Lambda \frac{\partial}{\partial v_i} h(\mathbf{x}, \mathbf{v}), \\ D[\Delta v_i \Delta v_j] &= 4\pi G^2 m_a^2 \ln \Lambda \frac{\partial^2}{\partial v_i \partial v_j} g(\mathbf{x}, \mathbf{v}). \end{aligned} \quad (7.83a)$$

Here m and m_a are the masses of the subject star and field stars, and the field-star DF $f_a(\mathbf{x}_a, \mathbf{v}_a)$ is normalized so that $\int d^3 \mathbf{v}_a f_a(\mathbf{x}_a, \mathbf{v}_a) = n$, where n is the number density of field stars. The functions $g(\mathbf{x}, \mathbf{v})$, $h(\mathbf{x}, \mathbf{v})$ are Rosenbluth potentials (Rosenbluth, MacDonald, & Judd 1957),

$$h(\mathbf{x}, \mathbf{v}) = \int d^3 \mathbf{v}_a \frac{f_a(\mathbf{x}, \mathbf{v}_a)}{|\mathbf{v} - \mathbf{v}_a|} \quad ; \quad g(\mathbf{x}, \mathbf{v}) = \int d^3 \mathbf{v}_a f_a(\mathbf{x}, \mathbf{v}_a) |\mathbf{v} - \mathbf{v}_a|. \quad (7.83b)$$

The factor $\ln \Lambda$, called the Coulomb logarithm, is often found in formulae for scattering rates from $1/r$ potentials, such as the gravitational potential of a point mass or the electrostatic potential of a point charge. In the present context we have (eq. L.15)

$$\Lambda = \frac{b_{\max} v_{\text{typ}}^2}{G(m + m_a)} \quad (7.84)$$

where v_{typ} is a typical velocity of stars in the system, and b_{max} is the maximum impact parameter considered. Numerical experiments show that the appropriate value for b_{max} is roughly the orbital radius R (Farouki & Salpeter 1994; Theuns 1996). The lack of a precise definition for b_{max} is a blemish on the derivation of the diffusion coefficients, but only a minor one: we have already shown in §7.4.2 that $\Lambda \approx N$ in typical systems containing N stars, so the fractional uncertainty in the Coulomb logarithm caused by, for example, a factor of two uncertainty in b_{max} is only $\sim \ln 2 / \ln N$, which is 0.06 for a globular cluster ($N \approx 10^5$) and only 0.03 for a galaxy ($N \approx 10^{11}$). For more discussion of these issues, see Weinberg (1989, 1993). For numerical evaluation of the diffusion coefficients, a useful rule is that $\Lambda = \lambda N$ with $\lambda \simeq 0.1$ for clusters composed of stars of a single mass (see page 588 and Giersz & Heggie 1994).

It is important to remember that this approach does not properly represent the effects of either very close or very distant encounters. Close encounters, those with impact parameter less than the 90° deflection radius $b_{90} \approx Gm/v^2$, have $\delta v/v$ of order unity and hence violate the Fokker–Planck approximation. Distant encounters, with impact parameters of order the system size R , cannot be treated using the local approximation, which is valid only for $b \ll R$. Nevertheless, the Fokker–Planck plus local approximations yield satisfactory results whenever $\ln \Lambda$ is significantly larger than unity. This is because, as we saw at the beginning of §7.4.2, equal octaves of impact parameter contribute equally to the relaxation process. When $\ln \Lambda \approx \ln(R/b_{90})$ is large, there are many octaves in impact parameter that contribute to the relaxation, and failure of the approximations for an octave or two at each end of this range does not lead to significant error.

When the mass m of the subject star is zero, the expressions (7.83) for the diffusion coefficients can be shown to be consistent with the fluctuation-dissipation theorem (7.82) (the proof uses the identity $\nabla^2 |\mathbf{x}| = 2/|\mathbf{x}|$).

The diffusion coefficients are simplified if the field star DF $f_a(\mathbf{x}, \mathbf{v}_a)$ is isotropic, that is, if it depends on velocity only through $v_a \equiv |\mathbf{v}_a|$. In a spherical system this condition is satisfied if f_a is ergodic, that is, if it depends only on the Hamiltonian H and not the angular momentum L . If f_a is isotropic, the only preferred direction in velocity space is defined by the velocity of the subject star \mathbf{v} , and therefore it is natural to choose a coordinate system in which $\hat{\mathbf{e}}_z$ is parallel to \mathbf{v} , and $\hat{\mathbf{e}}_x, \hat{\mathbf{e}}_y$ are perpendicular to \mathbf{v} . Then the symmetry of the problem demands that

$$D[(\Delta v_x)^2] = D[(\Delta v_y)^2], \quad (7.85)$$

and that

$$D[\Delta v_x] = D[\Delta v_y] = D[\Delta v_x \Delta v_y] = D[\Delta v_x \Delta v_z] = D[\Delta v_y \Delta v_z] = 0. \quad (7.86)$$

Hence there are only three independent diffusion coefficients, which we may write as

$$\begin{aligned} D[\Delta v_{\parallel}] &\equiv D[\Delta v_z], \\ D[(\Delta v_{\parallel})^2] &\equiv D[(\Delta v_z)^2], \\ D[(\Delta \mathbf{v}_{\perp})^2] &\equiv 2D[(\Delta v_x)^2] = 2D[(\Delta v_y)^2]. \end{aligned} \quad (7.87)$$

Here $\Delta v_{\parallel} \equiv \Delta \mathbf{v} \cdot \hat{\mathbf{v}}$, where $\hat{\mathbf{v}} \equiv \mathbf{v}/|\mathbf{v}|$ is a unit vector pointing along \mathbf{v} , and $\Delta \mathbf{v}_{\perp} \equiv \Delta \mathbf{v} - \Delta v_{\parallel} \hat{\mathbf{v}}$ is the component of $\Delta \mathbf{v}$ that is perpendicular to \mathbf{v} . The evaluation of these diffusion coefficients, described in Appendix L, yields:

$$\begin{aligned} D[\Delta v_{\parallel}] &= -\frac{16\pi^2 G^2 m_a (m + m_a) \ln \Lambda}{v^2} \int_0^v dv_a v_a^2 f_a(\mathbf{x}, v_a), \\ D[(\Delta v_{\parallel})^2] &= \frac{32\pi^2 G^2 m_a^2 \ln \Lambda}{3} \left[\int_0^v dv_a \frac{v_a^4}{v^3} f_a(\mathbf{x}, v_a) + \int_v^{\infty} dv_a v_a f_a(\mathbf{x}, v_a) \right], \\ D[(\Delta \mathbf{v}_{\perp})^2] &= \frac{32\pi^2 G^2 m_a^2 \ln \Lambda}{3} \\ &\quad \times \left[\int_0^v dv_a \left(\frac{3v_a^2}{v} - \frac{v_a^4}{v^3} \right) f_a(\mathbf{x}, v_a) + 2 \int_v^{\infty} dv_a v_a f_a(\mathbf{x}, v_a) \right]. \end{aligned} \quad (7.88)$$

The relation of these coefficients to the diffusion coefficients in an arbitrary Cartesian coordinate system is given by equation (L.24):

$$\begin{aligned} D[\Delta v_i] &= \frac{v_i}{v} D[\Delta v_{\parallel}], \\ D[\Delta v_i \Delta v_j] &= \frac{v_i v_j}{v^2} \left\{ D[(\Delta v_{\parallel})^2] - \frac{1}{2} D[(\Delta \mathbf{v}_{\perp})^2] \right\} + \frac{1}{2} \delta_{ij} D[(\Delta \mathbf{v}_{\perp})^2], \end{aligned} \quad (7.89)$$

where $\delta_{ij} = 1$ if $i = j$, and 0 otherwise.

The rate of change of the kinetic energy of the subject star is

$$\begin{aligned} D[\Delta E] &= m \sum_{i=1}^3 (v_i D[\Delta v_i] + \frac{1}{2} D[\Delta v_i \Delta v_i]) \\ &= m (v D[\Delta v_{\parallel}] + \frac{1}{2} D[(\Delta v_{\parallel})^2] + \frac{1}{2} D[(\Delta \mathbf{v}_{\perp})^2]) \\ &= 16\pi^2 G^2 m m_a \ln \Lambda \left[m_a \int_v^{\infty} dv_a v_a f_a(v_a) - m \int_0^v dv_a \frac{v_a^2}{v} f_a(v_a) \right]. \end{aligned} \quad (7.90)$$

The first of the integrals in the last line describes the growth of the kinetic energy of the subject star due to gravitational encounters with the field stars (“heating”), while the second describes the cooling effect of the first-order diffusion coefficient $D[\Delta v_{\parallel}]$. Only field stars moving faster than the subject star ($v_a > v$) contribute to the heating, while only stars moving slower than the subject star contribute to the cooling. The ratio of heating to cooling at a given speed v is proportional to m_a/m . The energy of a subject star reaches

equilibrium when the two terms balance—this is the classical phenomenon of **equipartition of energy**.

When the subject star is much more massive than the field stars ($m \gg m_a$), the first-order diffusion coefficient $D[\Delta v_{\parallel}]$ is larger than the second-order coefficients $D[(\Delta v_{\parallel})^2]/v$ and $D[(\Delta \mathbf{v}_{\perp})^2]/v$ by a factor of order m/m_a . In an isotropic field-star distribution, the effect of this diffusion coefficient on the subject star is identical to that of a force per unit mass equal in magnitude to $|D[\Delta v_{\parallel}]|$ and directed opposite to the subject star's velocity. The frictional forces on a body such as a ball flying through the air are also directed opposite to the velocity, so the force exerted on a massive subject star as a result of encounters with field stars is known as **dynamical friction** (Chandrasekhar 1943a). See §8.1 for a thorough discussion of dynamical friction and its astrophysical applications.

The diffusion coefficients can be explicitly evaluated if the field star DF $f_a(v_a)$ is known. The most important case is when the DF is Maxwellian,

$$f_a(v_a) = \frac{n}{(2\pi\sigma^2)^{3/2}} e^{-v_a^2/2\sigma^2}, \quad (7.91)$$

where n is the number density and σ is the one-dimensional velocity dispersion of the field stars. Evaluating the integrals of equation (7.88) we find

$$\begin{aligned} D[\Delta v_{\parallel}] &= -\frac{4\pi G^2 \rho (m + m_a) \ln \Lambda}{\sigma^2} G(X), \\ D[(\Delta v_{\parallel})^2] &= \frac{4\sqrt{2}\pi G^2 \rho m_a \ln \Lambda}{\sigma} \frac{G(X)}{X}, \\ D[(\Delta \mathbf{v}_{\perp})^2] &= \frac{4\sqrt{2}\pi G^2 \rho m_a \ln \Lambda}{\sigma} \left[\frac{\operatorname{erf}(X) - G(X)}{X} \right], \end{aligned} \quad (7.92)$$

where $\rho = m_a n$ is the density of field stars, $X \equiv v/(\sqrt{2}\sigma)$, $\operatorname{erf}(X)$ is the error function (Appendix C.3), and

$$G(X) \equiv \frac{1}{2X^2} \left[\operatorname{erf}(X) - X \frac{d \operatorname{erf}(X)}{dX} \right] = \frac{1}{2X^2} \left[\operatorname{erf}(X) - \frac{2X}{\sqrt{\pi}} e^{-X^2} \right]. \quad (7.93)$$

The behavior of $G(X)$ and the related functions occurring in equations (7.92) is shown in Figure 7.2.

It is straightforward to verify that the diffusion coefficients defined by equations (7.89) and (7.92) satisfy the fluctuation-dissipation theorem in the form (7.81), with the inverse temperature $\beta = 1/(m_a \sigma^2)$.

Heating of the Galactic disk by MACHOs We may use these results to estimate the rate at which the velocity dispersion of stars in the solar neighborhood grows from interactions with a hypothetical population of non-luminous, compact objects such as black holes (MACHOs) in the dark halo

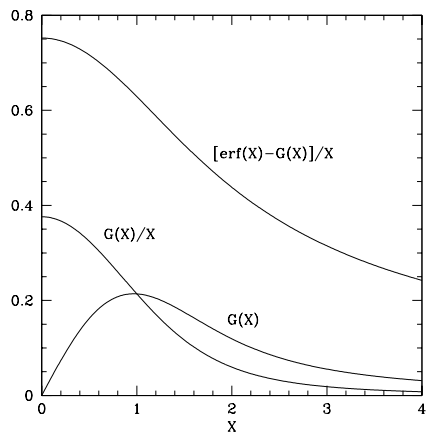


Figure 7.2 The function $G(X)$ (eq. 7.93) that appears in the diffusion coefficients (7.92) for a Maxwellian field-star DF. The functions $G(X)/X$ and $[\text{erf}(X) - G(X)]/X$ are also shown.

(page 16). The gradual growth of the random velocities of disk stars—a process called “disk heating”—can be caused by a number of mechanisms that are investigated in much more detail in §8.4. Heating by MACHOs is the simplest of these mechanisms (though probably far from the dominant one) and thus provides a good illustration of the use of the diffusion coefficients.

We begin by constructing a simple model for the DF of MACHOs in the dark halo. We shall assume that the DF is Maxwellian. The circular speed of the Galaxy is approximately constant near the solar radius, at $v_c = (220 \pm 20) \text{ km s}^{-1}$ (Table 1.2). If the Galaxy were spherical with a flat circular-speed curve, its density would be $\rho(r) = v_c^2 / (4\pi G r^2)$ (eq. 4.103), which is the singular isothermal sphere generated by a Maxwellian DF with dispersion $\sigma = v_c / \sqrt{2}$ (eq. 4.104). If a fraction f_h of this density were contributed by MACHOs, their density at the solar radius R_0 would be

$$\rho = f_h \frac{v_c^2}{4\pi G R_0^2} = 0.014 f_h \mathcal{M}_\odot \text{ pc}^{-3} \left(\frac{v_c}{220 \text{ km s}^{-1}} \frac{8 \text{ kpc}}{R_0} \right)^2. \quad (7.94)$$

Dynamical models of the Galaxy suggest that $f_h \simeq 0.1$ – 0.5 if the dark halo is made entirely of MACHOs (see §1.1.2 and §6.3.3); this is an upper limit to the actual value of f_h since part or all of the dark halo could be in some other form such as WIMPs (page 16).

Now we examine the evolution of the dispersion of the disk stars in the direction normal to the Galactic plane, which we label by z . The energy per unit mass arising from motion in this direction is $E_z \equiv \frac{1}{2}v_z^2 + \Phi_z(z)$ (eq. 3.74). Encounters with MACHOs cause this energy to increase at a rate

$$\begin{aligned} D[\Delta E_z] &= D[\Delta(\tfrac{1}{2}v_z^2)] = v_z D[\Delta v_z] + \tfrac{1}{2} D[(\Delta v_z)^2] \\ &= \frac{v_z^2}{v} D[\Delta v_\parallel] + \frac{v_z^2}{2v^2} \{ D[(\Delta v_\parallel)^2] - \tfrac{1}{2} D[(\Delta \mathbf{v}_\perp)^2] \} + \tfrac{1}{4} D[(\Delta \mathbf{v}_\perp)^2], \end{aligned} \quad (7.95)$$

where we have used equation (7.89). Since the disk stars travel on nearly circular orbits, we may replace v by the circular speed v_c , and neglect terms that are smaller than the dominant terms by a factor of $(v_z/v_c)^2$; thus all terms except the last one in the second line of equation (7.95) can be dropped. Then from equation (7.92) we have

$$D[\Delta E_z] = \frac{1}{4}D[(\Delta \mathbf{v}_\perp)^2] = \frac{\sqrt{2}\pi G^2 \rho m_a \ln \Lambda}{\sigma} \left[\frac{\operatorname{erf}(X) - G(X)}{X} \right], \quad (7.96)$$

where $X = v/(\sqrt{2}\sigma)$. Since $v \simeq v_c$ and we have argued above that $\sigma = v_c/\sqrt{2}$, we have $X = 1$ so the quantity in square brackets is 0.6289. Then by inserting equation (7.94) we have

$$D[\Delta E_z] = 0.314 \frac{G f_h m_a v_c \ln \Lambda}{R_0^2}. \quad (7.97)$$

In a self-gravitating disk, the virial theorem implies that $E_z = 3K_z = \frac{3}{2}\sigma_z^2$, where K_z is the kinetic energy per unit mass in z -motions (Problem 7.3) and σ_z is the RMS z -velocity in the disk. Thus $d\sigma_z^2/dt = \frac{2}{3}D[\Delta E_z]$; we integrate this equation to obtain

$$\begin{aligned} \sigma_z(t) &= 0.46 \left(\frac{G f_h m_a v_c t \ln \Lambda}{R_0^2} \right)^{1/2} \\ &= 13 \text{ km s}^{-1} \frac{8 \text{ kpc}}{R_0} \left(\frac{f_h}{0.5} \frac{m_a}{10^6 \mathcal{M}_\odot} \frac{v_c}{220 \text{ km s}^{-1}} \frac{\ln \Lambda}{10} \frac{t}{10 \text{ Gyr}} \right)^{1/2}. \end{aligned} \quad (7.98)$$

To estimate the size of the Coulomb logarithm, we may use equation (7.84) with $v_{\text{typ}} \simeq 200 \text{ km s}^{-1}$, $b_{\text{max}} \simeq 5 \text{ kpc}$, $m + m_a \simeq 10^6 \mathcal{M}_\odot$, which yields $\ln \Lambda \simeq 10.7$; as usual, the results are insensitive to the precise parameter choices since they enter only in the argument of the logarithm.

Before discussing the implications of this result, we examine the evolution of the velocities parallel to the plane. Replacing the radial Hamiltonian in the epicycle approximation, equation (3.102), by the radial energy E_R , and writing the radial velocity as $v_R = \dot{x}$, the corresponding diffusion coefficient can be written as

$$\begin{aligned} D[\Delta E_R] &= v_R D[\Delta v_R] + \frac{1}{2} D[(\Delta v_R)^2] \\ &\quad + \frac{1}{2} \gamma^2 \{ 2(v_\phi - v_c) D[\Delta v_\phi] + D[(\Delta v_\phi)^2] \}. \end{aligned} \quad (7.99)$$

Using equation (7.89) and observing that in a thin, cold disk $|v_R|, |v_\phi - v_c| \ll v$ and $v \simeq v_c$, this result simplifies to

$$D[\Delta E_R] = \frac{1}{4} D[(\Delta \mathbf{v}_\perp)^2] + \frac{1}{2} \gamma^2 D[(\Delta v_\parallel)^2]. \quad (7.100)$$

With the same assumptions about the DF of the MACHOs that we used to derive equation (7.97), we find

$$D[\Delta E_R] = \frac{G f_h m_a v_c \ln \Lambda}{R_0^2} (0.314 + 0.855 \Omega^2 / \kappa^2). \quad (7.101)$$

To find the rate of growth of the dispersions, we note that E_R is the energy of a harmonic oscillator, which is twice the mean value of the kinetic energy of the oscillator. The mean value of E_R in an ensemble of oscillators is thus $\langle E_R \rangle = \langle \dot{x}^2 \rangle = \sigma_R^2$, since \dot{x} is the radial velocity. Thus $d\sigma_R^2/dt = D[\Delta E_R]$, so

$$\sigma_R(t) = \left(\frac{G f_h m_a v_c t \ln \Lambda}{R_0^2} \right)^{1/2} (0.314 + 0.855 \Omega^2 / \kappa^2)^{1/2}. \quad (7.102)$$

Combining this result with equation (7.98), we find that in a galaxy with a flat circular-speed curve, in which $\kappa = \sqrt{2}\Omega$,

$$\frac{\sigma_z}{\sigma_R} = 0.53 \quad (7.103)$$

if the disk heating is due to MACHOs. This result is consistent with the observed ratio, $\sigma_z/\sigma_R \simeq 0.5$ (Figure 8.11 or BM Table 10.2). However, if MACHOs are the only heat source, both σ_R and σ_z should grow as $t^{0.5}$; the data in Figure 8.11 show that although this relation is approximately correct for σ_z , σ_R grows more slowly, as $t^{0.3}$. Thus it is unlikely that MACHOs are the dominant cause of disk heating, especially since we shall show in §8.4 that the effects of molecular clouds and spiral arms can explain the observed age-velocity dispersion relation without any contribution from MACHOs.

We can use these results to set an upper limit to the MACHO mass: since other mechanisms also heat the disk, and the dispersion of the oldest stars in the solar neighborhood is $\sigma_z \lesssim 30 \text{ km s}^{-1}$, equation (7.98) implies that the MACHO mass is (Lacey & Ostriker 1985)

$$m \lesssim 6 \times 10^6 \mathcal{M}_\odot \frac{0.5}{f_h}. \quad (7.104)$$

Other constraints on the MACHO mass are discussed in §8.2.2e.

7.4.5 Relaxation time

We can use the diffusion coefficients to improve the crude estimate of the relaxation time that we made in §1.2. We shall base our estimate on $D[(\Delta v_{\parallel})^2]$, and define the relaxation time of a subject star to be

$$t_{\text{relax}} \equiv \frac{v^2}{D[(\Delta v_{\parallel})^2]}. \quad (7.105)$$

To measure the characteristic relaxation time for a population of identical stars, we set $m_a = m$ and assume that the velocity distribution of the field stars is Maxwellian with dispersion σ (eq. 7.91). If the speed of the subject star were equal to the RMS speed of the field stars, $v = \sqrt{3}\sigma$, the value of X in equation (7.92) would be $(\frac{3}{2})^{1/2} = 1.225$. Since $G(X)/X$ is not a rapidly varying function of X (Figure 7.2), we simply set $X = 1.225$ in $D[(\Delta v_{\parallel})^2]$ and replace the factor v^2 in (7.105) by σ^2 , which is the mean-square velocity in any one direction. We find

$$\begin{aligned} t_{\text{relax}} &= 0.34 \frac{\sigma^3}{G^2 m \rho \ln \Lambda} \\ &= \frac{18 \text{ Gyr } 1 \mathcal{M}_{\odot}}{\ln \Lambda} \frac{10^3 \mathcal{M}_{\odot} \text{ pc}^{-3}}{m \rho} \left(\frac{\sigma}{10 \text{ km s}^{-1}} \right)^3. \end{aligned} \quad (7.106)$$

This definition is somewhat arbitrary, but the main role of our definition is simply to provide a fiducial timescale for parametrizing the speed of relaxation processes.

A natural first application of this formula is to estimate the relaxation time in the solar neighborhood. We take $\rho = 0.041 \mathcal{M}_{\odot} \text{ pc}^{-3}$ from Table 1.1. The velocity DF in the solar neighborhood is not isotropic, but we neglect this complication and simply identify the RMS velocity of nearby old stars, 50 km s^{-1} according to Table 1.2, with $\sqrt{3}\sigma$, which yields $\sigma \simeq 30 \text{ km s}^{-1}$. We set $\Lambda \simeq h\sigma^2/(G\mathcal{M}_{\odot})$ (eq. 7.84), where $h \simeq 500 \text{ pc}$ is a measure of the disk thickness, so $\ln \Lambda \simeq 18.5$. If the typical stellar mass is $m = 1 \mathcal{M}_{\odot}$, we find $t_{\text{relax}} \simeq 6 \times 10^{14} \text{ yr}$, almost five orders of magnitude longer than the age of the Galaxy. We conclude that relaxation due to stellar encounters is completely unimportant in the solar neighborhood.

One shortcoming of equation (7.106) for other applications is that the density ρ is often not directly observable. To derive a more convenient—but more approximate—formula for spherical stellar systems, we may assume that the stellar system is a singular isothermal sphere (§4.3.3b), consistent with our assumption that the DF is Maxwellian. Then equation (4.103) relates the density and velocity dispersion, so equation (7.106) becomes

$$\begin{aligned} t_{\text{relax}}(r) &= 2.1 \frac{\sigma r^2}{Gm \ln \Lambda} \\ &= \frac{5 \text{ Gyr } 1 \mathcal{M}_{\odot}}{\ln \Lambda} \frac{\sigma}{m} \frac{1}{10 \text{ km s}^{-1}} \left(\frac{r}{1 \text{ pc}} \right)^2. \end{aligned} \quad (7.107)$$

This result is strictly valid only for a singular isothermal sphere, but provides a useful first approximation to the relaxation time in many equilibrium stellar systems, so long as the system is not highly flattened or rapidly rotating.

The relaxation time can vary by several orders of magnitude in different regions of a single system. For reference purposes it is often useful

to characterize a system by a single measure of the relaxation time. To this end, we replace the density ρ in equation (7.106) by the mean density inside the system's half-mass radius r_h , which is just $\frac{1}{2}M/(\frac{4}{3}\pi r_h^3)$, and replace $3\sigma^2$ by the mean-square speed of the system's stars $\langle v^2 \rangle$. By equation (4.249b), $\langle v^2 \rangle \simeq 0.45GM/r_h$. To evaluate the Coulomb logarithm $\ln \Lambda$, we take equation (7.84), replace v_{typ}^2 by $\langle v^2 \rangle$, and replace the maximum impact parameter b_{max} by the half-mass radius r_h . These substitutions give $\Lambda = r_h \langle v^2 \rangle / (2Gm) = \lambda N$, with $\lambda \simeq 0.2$. Comparison of the relaxation rates in systems with different N suggests that $\lambda \simeq 0.1$ provides a better fit (Giersz & Heggie 1994), and this is the value we shall adopt. Then equation (7.106) yields the **half-mass relaxation time** (Spitzer 1969)

$$\begin{aligned} t_{\text{rh}} &= \frac{0.17N}{\ln(\lambda N)} \sqrt{\frac{r_h^3}{GM}} \\ &= \frac{0.78 \text{ Gyr}}{\ln(\lambda N)} \frac{1 \mathcal{M}_{\odot}}{m} \left(\frac{M}{10^5 \mathcal{M}_{\odot}} \right)^{1/2} \left(\frac{r_h}{1 \text{ pc}} \right)^{3/2}. \end{aligned} \quad (7.108)$$

The unshaded histogram in Figure 7.3 shows the distribution of t_{rh} for Galactic globular clusters; in almost all cases $t_{\text{rh}} \lesssim 10$ Gyr, confirming that substantial relaxation has occurred in these systems. The unshaded histogram shows the relaxation times for nearby galaxies, at the smallest radii that can be resolved by optical telescopes. Almost all of the relaxation times are much longer than 10 Gyr, so relaxation due to stellar encounters plays little or no role in determining the observable structure of galaxies.

7.4.6 Numerical methods

(a) Fluid models Throughout this book, we have often used fluids as analogs of stellar systems satisfying the collisionless Boltzmann equation. We can stretch this analogy even further, by using thermal conduction in fluids as an analog of relaxation. The basis for this analogy is that both are slow diffusive processes that cause entropy to increase. An obvious limitation of the analogy is that in fluids the mean free path is short, so conduction is a local process, while stars in stellar systems travel many times the system radius during a relaxation time.

The heat flux in a fluid is $\mathbf{q} = -\kappa \nabla T$ (eq. F.24), where T is the temperature and κ is the thermal conductivity. We can make a simple order-of-magnitude estimate of the thermal conductivity of a fluid. Let us suppose that the mean free path between collisions is λ and the mean time between collisions is τ . Now assume that there is a temperature gradient dT/dx in the x -direction only. The flow of particles from right to left across a unit area in the y - z plane at $x = 0$ is $\approx n\lambda/\tau$, where n is the number density (there are $\sim n\lambda$ particles per unit area within one mean free path of $x = 0$, and roughly half of these will cross $x = 0$ within one collision time). On average,

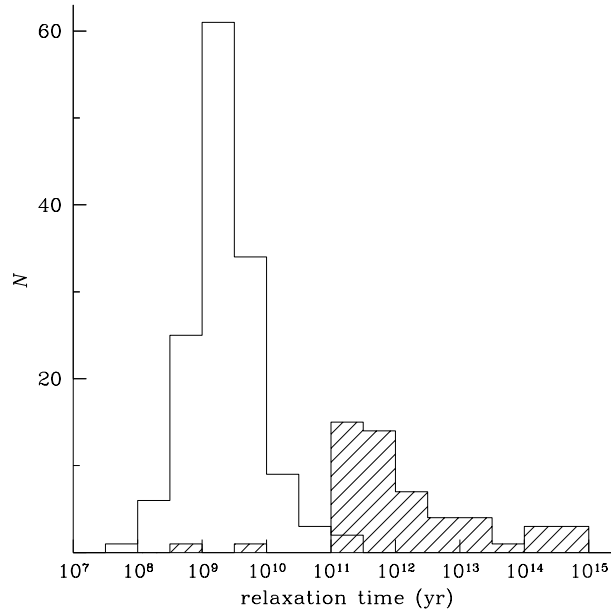


Figure 7.3 The distribution of relaxation times for 141 Galactic globular clusters (unshaded histogram) and 53 nearby galaxies (shaded histogram). The results for globular clusters are based on the half-mass relaxation time (eq. 7.108), and the catalog of Harris (1996), assuming a stellar mass $m = 0.7\mathcal{M}_\odot$ and mass-to-light ratio $\Upsilon = 2\Upsilon_\odot$. The results for galaxies are based on equation (7.106) and the galaxy sample of Faber et al. (1997); the relaxation times are evaluated at a radius corresponding to 0.1 arcsec, about the limiting resolution of the best telescopes. The two galaxies with the shortest relaxation times are M32 and M31 in the Local Group.

each molecule carries an energy $\approx k_B T_+$, where T_+ is the temperature at the location of its last collision and k_B is Boltzmann’s constant. Roughly, the last collision occurred one mean free path to the right of the wall, so $T_+ \approx T + \lambda(dT/dx)$, where T is the temperature at $x = 0$. The corresponding heat flux—energy per unit time crossing unit area at $x = 0$ from right to left—is

$$q_+ \approx -\frac{n\lambda k_B T_+}{\tau} \approx -\frac{n\lambda k_B T}{\tau} - \frac{n\lambda^2 k_B}{\tau} \frac{dT}{dx}. \quad (7.109)$$

Similarly, particles crossing $x = 0$ from left to right carry an energy of order $k_B T_- \approx k_B(T - \lambda dT/dx)$, and yield a heat flux

$$q_- \approx +\frac{n\lambda k_B T_-}{\tau} \approx \frac{n\lambda k_B T}{\tau} - \frac{n\lambda^2 k_B}{\tau} \frac{dT}{dx}. \quad (7.110)$$

The net heat flux is $q = q_+ + q_-$, and this must equal $-\kappa(dT/dx)$. Thus

$$\kappa \approx \frac{n\lambda^2 k_B}{\tau}. \quad (7.111)$$

In fluids, the mean free path and collision time are related by $\lambda/\tau \approx v_{\text{typ}}$, where v_{typ} is the typical velocity of a molecule. Thus we can eliminate λ to obtain

$$\kappa \approx nv_{\text{typ}}^2 k_B \tau. \quad (7.112)$$

Now let us pretend that a spherical stellar system like a globular cluster has a thermal conductivity κ , and use similar arguments to estimate what it would be. The “collision time” is simply the relaxation time, since stellar encounters change the stellar orbits substantially on this timescale. In spherical systems, we are interested in heat transport only in the radial direction; in typical spherical stellar systems, most stars oscillate in and out by a significant fraction of their mean orbital radius on a timescale much shorter than the relaxation time, so the appropriate “mean free path” is approximately the orbital radius r . (Note that in contrast to fluid systems the mean free path is not the total distance traveled between collisions, that is, the relation $\lambda/\tau \approx v_{\text{typ}}$ does *not* hold.) Then from equation (7.111) the thermal conductivity is (Lynden–Bell & Eggleton 1980)

$$\kappa \approx \frac{nr^2 k_B}{t_{\text{relax}}}. \quad (7.113)$$

Notice the curious fact that the thermal conductivity of a fluid *increases* with increasing collision time (eq. 7.112), while the thermal conductivity of a stellar system *decreases* with increasing relaxation time (eq. 7.113). The reason for the difference is that in a stellar system the mean free path in the radial direction is $\lambda \approx r$ and therefore is independent of the collision time, while in a fluid the mean free path grows with the collision time, $\lambda \approx v_{\text{typ}}\tau$.

We now derive the equations that govern the slow evolution of a spherical fluid as heat flows through it. Because the system evolves slowly, it is close to static at all times. Thus it must satisfy the equation of hydrostatic equilibrium (F.12),

$$\frac{\partial p}{\partial r} = -\rho \frac{\partial \Phi}{\partial r}, \quad (7.114)$$

where $p(r, t)$ is the pressure, $\Phi(r, t)$ is the gravitational potential, and we have used partial derivatives since the variables depend on both radius and time in an evolving system. The system must also satisfy Poisson’s equation in the form

$$\frac{\partial \Phi}{\partial r} = \frac{GM}{r^2} \quad ; \quad \frac{\partial M}{\partial r} = 4\pi r^2 \rho. \quad (7.115)$$

The heat flux is related to the change in specific entropy by equations (F.8), (F.21) and (B.53),

$$\begin{aligned} T \frac{ds}{dt} &= T \left(\frac{\partial s}{\partial t} + v \frac{\partial s}{\partial r} \right) \\ &= -\frac{1}{\rho r^2} \frac{\partial}{\partial r} (r^2 q) + \epsilon, \end{aligned} \quad (7.116)$$

where ϵ is the rate of energy production per unit mass (for example, by binary stars; see eq. 7.182), $d/dt = \partial/\partial t + v\partial/\partial r$ is a convective derivative, which measures the time derivative following a given mass shell rather than one at constant radius, and $v(r, t)$ is the (small) radial velocity associated with the slow evolution of the mass distribution through the continuity equation (F.3),

$$\frac{\partial \rho}{\partial t} + \frac{1}{r^2} \frac{\partial}{\partial r} (r^2 \rho v) = 0. \quad (7.117)$$

In practice, it is often simpler to use the mass $M(r)$ rather than the radius r as an independent variable; in this case the convective derivative and partial derivative are the same, $d/dt = (\partial/\partial t)_M$.

To relate the fluid system governed by these equations to a stellar system, we replace the temperature T by the one-dimensional velocity dispersion σ , using the relation $k_B T = m\sigma^2$, where m is the mass of the stars. The radial heat flux $q = -\kappa \partial T/\partial r$, so with equation (7.113) we have

$$q = -\frac{c_1 r^2 \rho}{t_{\text{relax}}} \frac{\partial \sigma^2}{\partial r}, \quad (7.118)$$

where $\rho = mn$ is the mass density and c_1 is a constant of order unity. Using equation (7.106) for the relaxation time, this result can be restated as

$$q = -\frac{c_2 r^2 G^2 \rho^2 m \ln \Lambda}{\sigma^3} \frac{\partial \sigma^2}{\partial r}, \quad (7.119)$$

where $\ln \Lambda$ is the usual Coulomb logarithm and c_2 is also of order unity. A simpler form is obtained if we use equation (7.107) for the relaxation time,

$$q = -c_3 G m \rho \ln \Lambda \frac{\partial \sigma}{\partial r}. \quad (7.120)$$

Finally, the pressure and entropy of our “ideal gas” of stars are related to the density and velocity dispersion by equations (F.31) and (F.43),

$$p = \rho \sigma^2 \quad ; \quad s = \frac{k_B}{m} \ln \left(\frac{\sigma^3}{\rho} \right) + \text{constant}. \quad (7.121)$$

Equations (7.114)–(7.121), together with a prescription for the rate of internal energy generation ϵ , constitute a complete set of partial differential

equations in two variables (r or $M(r)$, and t) that can be solved to follow the evolution of a stellar system due to relaxation. Notice that so long as the energy generation rate ϵ is negligible, the uncertain coefficients c_2 and c_3 in equations (7.119) and (7.120) affect only the timescale, not the outcome, of the evolution, since they can be removed by rescaling the time variable t to $c_2 t$ or $c_3 t$. Fluid models of this kind are described by Hachisu et al. (1978), Lynden–Bell & Eggleton (1980), and Giersz & Spurzem (1994).

Fluid models are much simpler than models based on the Fokker–Planck approximation. They provide extremely valuable insights into the evolution of N-body systems but these should be used with caution and confirmed using the more accurate methods described below.

A closely related approach is to take moments of the Fokker–Planck equation over velocity space, thereby deriving collisional analogs of the Jeans equations (§4.8). For example, Larson (1970a) investigated the evolution of a spherically symmetric stellar system by assuming that the DF can be written as a Maxwellian times a power series in v_r and $v_\theta^2 + v_\phi^2$. By truncating the power series at terms of order v^4 , he was able to obtain a closed set of differential equations for the moments of the DF. With this code he carried out some of the first numerical simulations of the gravothermal catastrophe and the evolution of star clusters (Larson 1970b, 1970a). Thorough treatments of the relation of moment equations to fluid models are given by Bettwieser (1983) and Louis (1990).

(b) Monte Carlo methods This approach can be regarded as a hybrid between direct N-body integrations and numerical solutions of the Fokker–Planck equation. In any stellar system with $N \gg 1$, there is a range of time intervals Δt that are short compared to the crossing time, but long enough that a subject star experiences a large number of encounters with passing stars. We do not care about the details of the velocity kicks associated with each encounter, but only about the statistical properties of the total velocity kick $\Delta \mathbf{v}$ at the end of the interval Δt . Fortunately, these are easy to determine: by the central limit theorem (Appendix B.10), the probability distribution of $\Delta \mathbf{v}$ is Gaussian, and hence is determined entirely by the mean $\mu_i = \langle \Delta v_i \rangle$ and the covariance matrix $C_{ij} = \langle \Delta v_i \Delta v_j \rangle$ (cf. eqs. B.98 and B.99). These can be computed from the diffusion coefficients (7.83) since

$$\mu_i = D[\Delta v_i] \Delta t \quad ; \quad C_{ij} = D[\Delta v_i \Delta v_j] \Delta t. \quad (7.122)$$

In Monte Carlo solutions of the Fokker–Planck equations, we choose a random sample of $p \ll N$ subject stars from the stellar system. The orbit of each subject star is integrated in the gravitational potential of the system, and in addition it is perturbed by velocity kicks at time intervals Δt , chosen by the prescription in the preceding paragraph. In the case of a spherical star cluster, it is convenient to think of each subject star as representing N/p real stars, each with the same pericenter and apocenter but randomly distributed orbital planes. The introduction of these shells of stars,

sometimes called “superstars,” forces the calculation to maintain spherical symmetry and thus greatly simplifies the calculation of the gravitational potential (see Problem 2.24).

A method of this kind is sometimes called an “orbit-following” Monte Carlo solution of the Fokker–Planck equation. An alternative approach is to use an “orbit-averaged” Monte Carlo method. These are based on the observation that encounters affect the orbit only over a time t_{relax} that is much longer than the crossing time t_{cross} . Thus there is little to be gained by integrating the orbits of the superstars. Rather, it is sufficient to monitor only the energy E and angular momentum L (or apocenter and pericenter, or radial and azimuthal action) of each superstar and not its instantaneous radius. At each timestep—which can now be much longer than t_{cross} —the superstar is assigned to a randomly chosen phase of its orbit and E and L are perturbed by random kicks, whose probability distribution is determined using the central limit theorem and the diffusion coefficients at the current radius.

Orbit-averaged methods are much faster than orbit-following methods, since the timestep is a fraction of the relaxation time rather than the crossing time. However, despite their greater computational cost, orbit-following methods are useful because they provide a more accurate treatment of the evaporation process (see §7.5.2), and they can be used to follow processes such as violent relaxation and tidal shocks that occur on a crossing time or less. Orbit-following methods were pioneered by Spitzer and his collaborators (Spitzer 1987). Orbit-averaged methods were introduced by Hénon (1972, 1973b) and have been employed by a number of groups to study the evolution of globular clusters and the dense centers of galaxies (Stodólkiewicz 1986; Giersz 1998; Freitag, Rasio, & Baumgardt 2006).

A drawback of Monte Carlo methods is that properties such as the density distribution must be estimated by counting particles, and thus are subject to statistical fluctuations that obscure small-scale features.

(c) Numerical solution of the Fokker–Planck equation For a spherical system, the orbit-averaged Fokker–Planck equation is a partial differential equation in three variables: time t , radial action J_r , and angular momentum L . The solution is complicated by the fact that the density and thus the potential of the system are slowly changing as the system evolves; thus each time the DF $f(J_r, L, t)$ is updated, we must recompute the potential using Poisson’s equation, and in turn the change in potential affects the DF. One attractive feature of writing the DF in terms of the actions J_r and L is that they are adiabatic invariants (§3.6); thus the DF $f(J_r, L, t)$ is invariant under slow changes in the potential.

Direct solution of the Fokker–Planck equation yields an estimate of the DF that is free from the statistical noise inherent in Monte Carlo methods. However, it is generally easier to incorporate additional effects such as stellar evolution, stellar collisions, external tidal fields, binary formation, rotation, large-angle scattering, etc. into a Monte Carlo simulation.

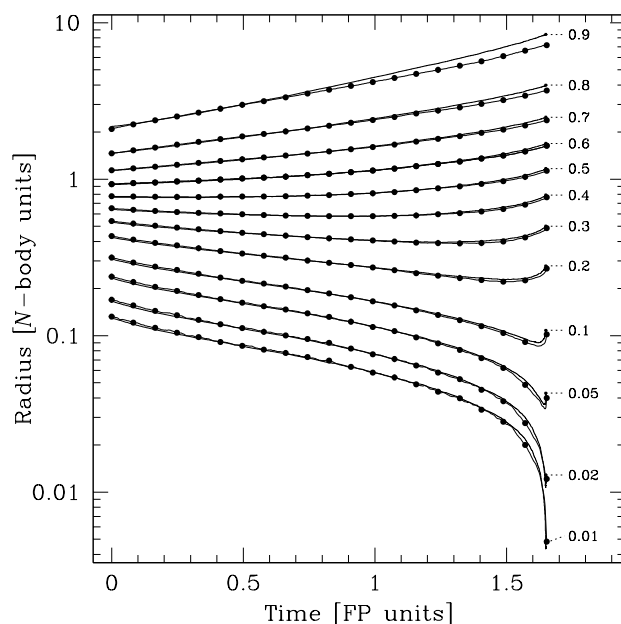


Figure 7.4 The evolution of an ergodic Plummer model, according to an orbit-averaged Monte-Carlo solution of the Fokker–Planck equation with $N = 3 \times 10^5$ superstars (heavy solid lines) and an N-body simulation with $N = 65\,536$ (dots, connected by light solid lines). The radii containing a given fraction of the total mass (1%, 2%, 5%, ..., 90%) are shown as a function of time. The plot is based on **N-body units**, in which $G = 1$ and the mass and total energy of the initial system are given by $M = 1$ and $-4E = 1$. In N-body units, the units of length and time are $L_0 = GM^2/(-4E)$ and $T_0 = GM^{5/2}(-4E)^{-3/2}$, and the scale length of the initial Plummer model is $b = 3\pi L_0/16$. In the figure, the radii are shown in units of L_0 and the time in units of $T_R \equiv T_0 N / \ln(\lambda N)$, $\lambda \simeq 0.1$, so T_R is comparable to the relaxation time (eq. 7.108). From Freitag, Rasio, & Baumgardt (2006).

The details of the numerical methods used to solve the orbit-averaged Fokker–Planck equation (7.71) are described by Cohn (1979, 1980), Takahashi (1995) and Drukier et al. (1999).

(d) N-body integrations The most accurate, and most expensive, way to follow the evolution of a star cluster is through direct N-body integrations (Aarseth 2003; Heggie & Hut 2003).

Simulations of a massive globular cluster with $N \approx 10^6$ are still not practical. The largest calculations done so far have $N \approx 10^5$, and these can require months of time on a special-purpose computer (Baumgardt & Makino 2003). Extrapolating the behavior seen in small- N simulations to larger N requires great care, because, as we have seen in §7.1, many physical processes contribute to the evolution of globular clusters—relaxation, ejection,

tion, evaporation, mass loss from stellar evolution, formation of binary stars, etc.—and these depend on N in different ways (e.g., Aarseth & Heggie 1998; Fukushige & Heggie 2000).

Although N-body simulations are simple in principle, many sophisticated techniques are needed to make them efficient enough to be useful. (i) The central core of the cluster is much denser than its outer parts, so the timestep required for stars in the core must be much shorter than the timestep for most other stars. This problem is solved by assigning each star its own timestep (§3.4.6). (ii) Binary and triple stars can be formed with orbital periods far shorter than the crossing time. These require special treatment, by a combination of analytical solutions of the Kepler problem (§3.1b) and other techniques (Mikkola & Aarseth 1996, 1998). (iii) The computation of the force by direct summation (§2.9.1) requires of order N^2 operations, while advancing the orbits requires only of order N operations. Thus the force computation takes far more time than the orbit calculation. This problem can be addressed in either hardware or software. The hardware solution is to develop special-purpose computers that parallelize and pipeline the force calculation, leaving the much easier task of advancing the orbits to a general-purpose host computer (Makino 2001). The software solution is to separate the rapidly varying forces from the small number of nearby particles and the slowly varying forces from the large number of distant ones, using a neighbor scheme (Ahmad & Cohen 1973) or a tree code (§2.9.2). (iv) Regularization (§3.4.7) dramatically reduces the numerical errors that arise in following close encounters between two stars.

(e) Checks and comparisons It is extremely important to compare the results given by the four different methods described in this section. Many such comparisons have been made, and the agreement has steadily improved as the codes have become more sophisticated (e.g., Giersz & Spurzem 1994; Spurzem & Takahashi 1995; Heggie et al. 1998).

Figure 7.4 shows the results of one such comparison. The initial state was an isolated Plummer model with an ergodic DF (§4.3.3a). All of the stars had equal masses. The radii containing various fractions of the total mass are plotted as a function of time. The results from an N-body integration with $N = 65\,536$ particles (Baumgardt et al. 2003) and an orbit-averaged Monte Carlo simulation (Freitag, Rasio, & Baumgardt 2006) are quite similar.

In addition to checking numerical methods, such comparisons confirm the validity of the basic assumptions of the Fokker–Planck and local approximations, that diffusion due to weak, local, two-body encounters is the principal source of relaxation in stellar systems.

7.5 The evolution of spherical stellar systems

In this section we describe the evolution of a spherical stellar system due to encounters between stars. For concreteness we shall usually assume that the system in question is a globular cluster, although applications to the centers of galaxies are discussed briefly in §7.5.9.

The formation of globular clusters is not well understood, and thus we have only a crude idea of their typical states just after they have settled into dynamical equilibrium. Fortunately for our purposes here, relaxation tends to erase a cluster's memory of its initial state, so the numerical experiments described below give very similar results for a wide range of initial conditions.

Since the relaxation time is inversely proportional to density, evolution due to relaxation proceeds most rapidly in the dense central regions of the cluster. Within the relaxed central region, the DF $f(E, L)$ and the density distribution $n(r)$ should approximate the isothermal distribution of §4.3.3b, i.e., the DF should be approximately Maxwellian at energies well below the escape energy.

In the outer parts of the cluster (the “halo”) the relaxation time is long, and encounters have relatively little effect. However, as relaxation proceeds, the halo population is augmented by stars that were originally in the relaxed central region but have now diffused to higher energies as a result of encounters. Although the apocenters of such orbits may lie far out in the halo, their pericenters must still lie in the relaxed center—an orbit that remains bound can never be expelled completely from the relaxed region by encounters, since encounters are effective only in this region. These *nouveau riche* halo members, who have risen in the world in consequence of a series of profitable encounters with their less fortunate neighbors, ultimately overwhelm the original halo members, who acquired their wealth at birth. Thus after a few central relaxation times, the properties of the halo are determined by relaxation, rather than by the initial conditions.

Many of the properties of this halo can be derived analytically (Spitzer & Shapiro 1972). Let $n(E, t)dE$ be the number of halo stars with energy between E and $E+dE$ at time t ; as argued above, these will have pericenters in the relaxed central region and apocenters in the halo. The Fokker–Planck equation for $n(E, t)$ is derived from a master equation by arguments precisely analogous to those used in §7.4.2, and reads

$$\frac{\partial n(E, t)}{\partial t} = -\frac{\partial}{\partial E} \{n(E, t)D[\Delta E]\} + \frac{1}{2} \frac{\partial^2}{\partial E^2} \{n(E, t)D[(\Delta E)^2]\}, \quad (7.123)$$

where $D[\Delta E]$ and $D[(\Delta E)^2]$ are the mean and mean-square changes in energy per unit time. This equation can be rewritten as a continuity equation for the flux $\mathcal{F}_E(E, t)$, which is the rate at which stars diffuse across a constant-energy surface in phase space,

$$\frac{\partial n}{\partial t} + \frac{\partial \mathcal{F}_E}{\partial E} = 0 \quad (7.124a)$$

where

$$\mathcal{F}_E(E, t) = n(E, t)D[\Delta E] - \frac{1}{2} \frac{\partial}{\partial E} \{n(E, t)D[(\Delta E)^2]\}. \quad (7.124b)$$

Now let Φ_c be a typical value of the potential in the central region. Halo stars have energy $|E| \ll |\Phi_c|$, that is, they are much closer to the escape energy $E = 0$ than stars in the relaxed center. Thus the speed of a halo star passing through the relaxed center is $v = [2(E - \Phi)]^{1/2} \simeq (-2\Phi_c)^{1/2}$, approximately independent of the energy. This means that the mean and mean-square energy changes during one pericenter passage through the relaxed region should not depend strongly on the energy. However, the radial period *does* depend strongly on the energy: since the halo stars spend most of their time outside most of the cluster mass, their orbital period is given approximately by the Kepler formula

$$T_r(E) = \frac{2\pi GM}{(-2E)^{3/2}}, \quad (7.125)$$

where M is the cluster mass (eqs. 3.31 and 3.32). Thus the diffusion coefficients, which characterize the orbit-averaged energy changes per unit time, must have the form

$$D[\Delta E] = \frac{c'_1}{T_r(E)} = c_1(-E)^{3/2} \quad ; \quad D[(\Delta E)^2] = \frac{c'_2}{T_r(E)} = c_2(-E)^{3/2}, \quad (7.126)$$

where the c 's are constants. The flux (7.124b) becomes

$$\mathcal{F}_E(E, t) = c_1(-E)^{3/2}n(E, t) - \frac{1}{2}c_2 \frac{\partial}{\partial E} \left[(-E)^{3/2}n(E, t) \right]. \quad (7.127)$$

In a steady state, the flux is independent of E and t , and determined by the rate at which the relaxed region feeds stars into the halo. If we take \mathcal{F}_E to be constant, and set $y(E) = (-E)^{3/2}n(E)$, we can solve the differential equation (7.127) to obtain

$$y(E) = \frac{\mathcal{F}_E}{c_1} \left(1 - c_3 e^{2c_1 E/c_2} \right), \quad (7.128)$$

where c_3 is an integration constant.

To proceed further, we need to specify the boundary condition on the differential equation (7.127) at the outer edge of the cluster. For an isolated cluster, we may assume that stars escape for $E > 0$, so the appropriate boundary condition is $y(0) = 0$. When the cluster orbits within a galaxy, stars at large radii may be torn off by the tidal field of the galaxy, even though they have bound energies, $E < 0$. This process will be investigated

in §8.3; for our present purposes, it is sufficient to assume that stars are lost from the cluster if their energy exceeds $E_t \equiv -GM/r_t$, where r_t is usually called the “tidal radius.” Note that the outer radius of King models of stellar systems has the same name (see §4.3.3 and eq. 4.113), as does the sharp outer boundary of globular clusters seen by observers—the underlying assumption behind this confusing practice is that all three measures of the maximum radius of a cluster reflect the same underlying physics.

The appropriate boundary condition is then $y(E_t) = 0$; applying this boundary condition to eliminate the integration constant, we have

$$n(E) = \frac{\mathcal{F}_E}{c_1(-E)^{3/2}} \left[1 - e^{2c_1(E-E_t)/c_2} \right] \quad (E < E_t), \quad (7.129)$$

and zero otherwise. Since the constants c_1 and c_2 are determined by the relaxation rate in the central region, their ratio $c_1/c_2 = D[\Delta E]/D[(\Delta E)^2]$ must be of order Φ_c^{-1} . Since our analysis is valid only for halo energies $|E| \ll |\Phi_c|$, we suffer no additional loss in accuracy if we replace the exponential e^x in equation (7.129) by $1 + x$. Thus

$$n(E) = \frac{2\mathcal{F}_E}{c_2(-E)^{3/2}}(E_t - E) \quad (E < E_t), \quad (7.130)$$

and zero otherwise. At fixed angular momentum, the DF is related to $n(E)$ by (see Problem 4.8)

$$f(E) \propto \frac{n(E)}{T_r} \propto E_t - E; \quad (7.131)$$

thus the DF goes linearly to zero at the edge of the cluster.¹¹

A simple argument now yields the halo density $n(r)$. Since the eccentricities of stars in the halo are large, we can approximate the orbits as radial. Then the fraction of the time that a star of energy E spends in the interval from r to $r + dr$ is $p(r|E)dr$, where

$$p(r|E)dr = \frac{2}{T_r} \frac{dr}{|v_r|} = \frac{2(-E)^{3/2}dr}{\pi GM \sqrt{E + GM/r}}; \quad (7.132)$$

here $v_r = \pm\sqrt{2(E + GM/r)}$ is the radial velocity, and we have used equation (7.125). Thus with equation (7.130) the density of halo stars is

$$n(r) = \frac{1}{4\pi r^2} \int dE p(r|E)n(E) = \frac{\mathcal{F}_E}{\pi^2 c_2 GM r^2} \int_{-GM/r}^{E_t} dE \frac{(E_t - E)}{\sqrt{E + GM/r}}. \quad (7.133)$$

¹¹ This treatment fails when $|E| \lesssim \epsilon$, where ϵ is the RMS energy change per orbit. See Spitzer & Shapiro (1972) for a more accurate analysis.

The integration is straightforward, and yields

$$n(r) = \frac{4\mathcal{F}_E(GM)^{1/2}}{3\pi^2 c_2 r^{7/2}} \left(1 - \frac{r}{r_t}\right)^{3/2}, \quad (7.134)$$

in which we have replaced the escape energy E_t by the tidal radius $r_t = -GM/E_t$. Thus, for example, the halo density in an isolated cluster with $E_t = 0$ is $n(r) \propto r^{-7/2}$.

This result does not hold at arbitrarily large radius, for at least three reasons: (i) Close to the tidal radius r_t , the approximation that the orbits are precisely radial fails, because the kinetic energy associated with the tangential motion, $\frac{1}{2}L^2/r^2$, becomes comparable to the total kinetic energy $E - GM/r$; in this region the density falls as $n(r) \propto (1 - r/r_t)^{5/2}$ (Problem 4.19). (ii) Escaping stars contribute a density $n(r) \propto r^{-2}$ (the velocity v of these stars is constant, and the total flux through a given radius, $4\pi r^2 n(r)v$, must also be constant), which eventually dominates at sufficiently large radii (§8.3.3). (iii) Stars with $E > E_t$ may remain close to the cluster for many orbital periods. A more careful treatment of escaping stars in tidally limited clusters is given by Fukushige & Heggie (2000).

The relatively simple analytic arguments given in this section yield a rich harvest of constraints on the properties of globular clusters, which should hold in the interval after a few central relaxation times and before the late stages of core collapse (§7.5.3): (i) within the relaxed central region, the DF should be approximately isothermal, $f(E, L) \propto \exp(-E/\sigma^2)$, at energies well below the escape energy; (ii) within the relaxed region, the density distribution $n(r)$ should be approximately that of an isothermal sphere; (iii) there should be relatively few stars with angular momenta greater than a cutoff L_0 corresponding to the angular momentum of a nearly unbound orbit that grazes the relaxed region; (iv) the DF should tend to zero near the escape energy as $f \propto E_t - E$ (eq. 7.131); (v) there should be an extended region in which the number density tends to zero as $n(r) \propto r^{-7/2}$ in an isolated cluster, and as $n(r) \propto r^{-7/2}(1 - r/r_t)^{3/2}$ in a cluster with tidal radius r_t ; (vi) the radial and tangential velocity dispersions should be the same in the relaxed central region (since the DF is ergodic there), while in the halo the velocity ellipsoid should become more and more radial, with the tangential dispersion falling off as r^{-1} (since $L = rv_t \lesssim L_0$ so $v_t \lesssim L_0/r$) and the radial dispersion in an isolated system falling off more slowly as $r^{-1/2}$.

A DF that satisfies all of these criteria is the Michie DF of equation (4.117). Thus, the Michie DF provides a good empirical model for the DFs of globular clusters and other partially relaxed stellar systems.

With this framework in place, we now review some of the most important processes that determine the collisional evolution of a globular cluster. The reader should bear in mind that in fact all of these processes occur simultaneously, and it may be difficult in practice—and sometimes even in principle—to isolate the effects of any single process on cluster evolution.

7.5.1 Mass loss from stellar evolution

Stars often eject mass from their surfaces near the ends of their lives. If the mass-losing star is located in a globular cluster, the ejected gas is likely to escape the cluster, either because the ejection velocity exceeds the escape speed from the cluster or because intracluster gas is regularly swept out by galactic gas when the cluster passes through the disk. Thus the cluster mass declines as stars evolve.

The evolution timescale of a typical population of stars is usually much longer than the crossing time in the cluster. Thus the adiabatic invariants of the stellar orbits are conserved as the cluster loses mass (§3.6). For example, consider a star orbiting in the outer parts of the cluster. Here the potential is close to Keplerian, so the azimuthal and radial actions may be written as $J_2 = (GMa)^{1/2}(1 - e^2)^{1/2}$ and $J_3 = (GMa)^{1/2}[1 - (1 - e^2)^{1/2}]$, where M is the total cluster mass, a is the semi-major axis and e is the eccentricity (Table E.1). If the total mass changes by a factor ψ , conservation of J_2 and J_3 requires that e remains constant and $a \propto 1/\psi$. Thus the orbits expand, but retain the same shape. Models of the stellar mass distribution and its evolution suggest $\psi \simeq 0.7$ for an old globular cluster (Box 7.2).

This expansion has important consequences for the cluster, which we may investigate using a simple model. Let us assume that the cluster has a Plummer density profile, equation (2.44b). Then the mean density interior to radius r is

$$\bar{\rho}(r) = \frac{M(r)}{\frac{4}{3}\pi r^3} = \frac{3M}{4\pi b^3} \frac{1}{(1 + r^2/b^2)^{3/2}}, \quad (7.135a)$$

where M is the mass and b is the scale length of the Plummer model. If the mass of the stars is slowly reduced by a factor ψ , all of the orbits will expand adiabatically by a factor $1/\psi$, so the density profile will be that of a Plummer model with mass $M' = \psi M$ and scale length $b' = b/\psi$. The mean density will then be

$$\bar{\rho}'(r) = \frac{3M'}{4\pi b'^3} \frac{1}{(1 + r^2/b'^2)^{3/2}} = \frac{3\psi^4 M}{4\pi b^3} \frac{1}{(1 + \psi^2 r^2/b^2)^{3/2}}. \quad (7.135b)$$

Note that when the mass falls by a factor ψ , the central density falls by the much greater factor ψ^4 .

We shall show in §8.3.1 that tidal forces from the host galaxy impose an outer limit to a globular cluster at a radius r_J such that $\bar{\rho}(r_J) \sim \bar{\rho}_h$ (eq. 8.92), where $\bar{\rho}_h$ is the mean density of the host galaxy inside the orbital radius of the cluster. If we make the crude approximations that (i) the tidal forces simply cut off the cluster at r_J without affecting the density distribution inside that radius, and (ii) the orbital radius and hence $\bar{\rho}_h$ remain constant, then we can relate the limiting radii r_J and r'_J before and after mass loss by setting $\bar{\rho}(r_J) = \bar{\rho}'(r'_J)$ in equations (7.135). Thus we find

$$r'_J{}^2 = \psi^{2/3} r_J^2 + b^2 \left(\psi^{2/3} - \frac{1}{\psi^2} \right). \quad (7.136)$$

Box 7.2: The initial mass function and the initial-final mass function

The mass lost by a cluster through stellar evolution depends mainly on two functions. The **initial mass function** or IMF $\xi(m)$ specifies the distribution of masses of stars just after they have formed (BM §5.1.9); thus, immediately after a burst of star formation in some region the number of stars with masses in the range $(m, m + dm)$ is

$$dn \propto \xi(m)dm, \quad (1)$$

where for our purposes the normalization of $\xi(m)$ is arbitrary. The historic Salpeter (1955) IMF is $\xi(m) = m^{-2.35}$, which remains a reasonable approximation for $m \gtrsim 0.5 \mathcal{M}_\odot$ but overestimates the number of stars at lower masses. A better representation of the data is (Kroupa 2002)

$$\xi(m) = \begin{cases} m^{-0.3} & 0.01 < m / \mathcal{M}_\odot < 0.08 \\ c_1 m^{-1.3} & 0.08 < m / \mathcal{M}_\odot < 0.5 \\ c_2 m^{-2.3} & 0.5 < m / \mathcal{M}_\odot < 1 \\ c_3 m^{-2.7} & 1 < m / \mathcal{M}_\odot, \end{cases} \quad (2)$$

where the constants c_i are chosen so that $\xi(m)$ is continuous. The data are generally consistent with the hypothesis that the IMF is universal—the same in any environment in which stars are formed.

The second important function is the **initial-final mass function** $\mu(m, t)$: if a star's initial mass is m , $\mu(m, t)$ is its mass after time t . For old globular clusters, the time t can be taken to be the age of the universe t_0 , and we write $\mu(m) \equiv \mu(m, t_0)$. This function can be estimated from theoretical calculations of mass loss in the late stages of stellar evolution, the distribution of masses of the nuclei of planetary nebulae, or the ages of white dwarfs in clusters (Han, Podsiadlowski, & Eggleton 1994; Hurley, Pols, & Tout 2000). A simple approximation is

$$\mu(m) = \begin{cases} m & m < 1 \text{ (main-sequence star)} \\ 0.5 + 0.13(m - 1) & 1 < m < 8 \text{ (white dwarf)} \\ 1.3 + 0.025(m - 8) & 8 < m \text{ (neutron star or black hole)}, \end{cases} \quad (3)$$

where m is in solar masses and the brackets indicate the type of remnant.

The ratio of the mass of the cluster at t_0 to its initial mass is then

$$\psi = \frac{\int dm m \xi(m)}{\int dm \mu(m) \xi(m)}. \quad (4)$$

For the IMF and initial-final mass function given in equations (2) and (3), $\psi = 0.72$. This estimate is uncertain because of the uncertain shape of the IMF at large masses—half of the mass loss comes from $m > 3 \mathcal{M}_\odot$.

If the initial tidal radius is large, $r_J \gg b$, then $r'_J = \psi^{1/3} r_J$, so the shrinkage in the limiting radius is modest, only a factor of 0.89 if $\psi = 0.7$. However, if the tidal forces are stronger, so the initial limiting radius is smaller, the effects of mass loss are more severe. In particular, if $r_J/b < 1.26$ then for $\psi \leq 0.7$ the cluster is completely disrupted (i.e., there is no solution for r'_J).

These results suggest that clusters with low central concentration (low ratio of central density to mean density; see page 30) are particularly susceptible to disruption by tidal forces as their stars lose mass. This conclusion is confirmed by more accurate models. For example, Chernoff & Weinberg (1990) consider mass loss in King models, in which the concentration is specified by $c = \log_{10}(r_t/r_0)$ where r_0 and r_t are the cluster's King radius and tidal radius (eq. 4.114). They find that models with $c < 0.5$ are completely disrupted if 30% of the mass is lost ($\psi < 0.7$). Clusters with larger values of the concentration can survive much larger fractional mass loss, although at the expense of losing most of their extended envelopes.

The median concentration of Galactic globular clusters is $c \simeq 1.5$; with this concentration, they are disrupted only if $\psi \lesssim 0.08$ (Chernoff & Weinberg 1990), so mass loss due to stellar evolution plays only a minor role in their dynamical evolution. Nevertheless, the fragility of low-concentration clusters suggests that the initial cluster population may have been much larger than the present one, containing many low-concentration clusters that are no longer with us (see §7.5.6).

7.5.2 Evaporation and ejection

Stars can escape from a cluster by two conceptually distinct mechanisms: (i) A single close encounter with another star can produce a velocity change comparable with the initial velocities of the two stars, thereby leaving one of the stars with a speed exceeding the local escape speed; we call this process **ejection**. (ii) A series of weaker, more distant encounters can gradually increase the energy of a star, until a final weak encounter gives the star slightly positive energy and it escapes; we call this process **evaporation** to suggest its more gradual nature.

Hénon (1960, 1969) has calculated the ejection rate for an isolated system with a Plummer density distribution and an ergodic DF. If all the stars have the same mass—we shall call this a **single-mass cluster**—the ejection rate is

$$\frac{dN}{dt} = -1.05 \times 10^{-3} \frac{N}{t_{\text{rh}} \ln(\lambda N)}, \quad (7.137)$$

where we have expressed the result in units of the half-mass relaxation time t_{rh} (eq. 7.108). The Coulomb logarithm $\ln(\lambda N)$ is present only to cancel the dependence of t_{rh} on $\ln(\lambda N)$ —since the Coulomb logarithm arises from the cumulative effect of distant encounters, it does not appear in the ejection rate, which is due to close encounters.

From equation (7.137) we can define an ejection time

$$t_{\text{ej}} = - \left(\frac{1}{N} \frac{dN}{dt} \right)^{-1} = 1 \times 10^3 \ln(\lambda N) t_{\text{rh}}. \quad (7.138)$$

For typical values of the Coulomb logarithm, $\ln(\lambda N) \approx 10$, we shall find that t_{ej} is much longer than the evaporation time due to weak encounters. Hence for most purposes we can neglect ejection relative to evaporation.

Evaporation is a more complicated process than ejection: myriads of weak encounters cause the star to diffuse at random through phase space, and some of the most energetic stars wander into the unbound region of phase space. Stars on high-energy orbits that lie entirely within the (low-density) halo experience very few encounters. Thus the evaporation rate is dominated by stars on highly elongated orbits, which are buffeted by encounters as they plunge through the dense cluster center, as discussed at the beginning of this section. As the energy of such a star approaches escape energy, the apocenter grows and the orbital period becomes longer, but the pericenter tends to remain at roughly the same distance. Thus the RMS energy change per orbit is approximately constant, and we denote this constant by ϵ_2 . When the energy of the halo star is within ϵ_2 of the escape energy E_t , there is a substantial chance that it will escape after its next passage through pericenter.

The behavior described above has important consequences for the orbit-averaged Fokker–Planck equation. Orbit averaging is valid only so long as the fractional changes in the orbital parameters are small in a single orbit. However, a star can escape from an isolated cluster only when its binding energy $E_t - E$ is comparable to the RMS energy change per orbit ϵ_2 . Hence *an orbit-averaged calculation cannot accurately predict the rate of escape from an isolated cluster*. Hénon (1960, 1969) has given a striking example of the problems that can arise in this way. Consider an isolated cluster, in which the escape energy $E_t = 0$. Let us calculate a lower limit to the escape time by assuming that the energy change per orbit is always positive and equal to the RMS change ϵ_2 , rather than being positive or negative with nearly equal probability. Then in the orbit-averaged approximation, the rate of change of a star's energy is

$$\frac{dE}{dt} = \frac{\epsilon_2}{T_r} = \frac{\sqrt{2}\epsilon_2(-E)^{3/2}}{\pi GM}, \quad (7.139)$$

where the radial period T_r is given by (7.125). Integrating this differential equation with the initial condition $E = E_0 < 0$ at $t = 0$, we find

$$E(t) = \frac{E_0}{\left(1 + \frac{|E_0|^{1/2} \epsilon_2 t}{\sqrt{2}\pi GM} \right)^2}. \quad (7.140)$$

Thus the star achieves escape energy $E = 0$ only as $t \rightarrow \infty$; in words, a star can never escape from an isolated cluster in the orbit-averaged approximation. This is **Hénon's paradox**.¹²

More accurate Fokker–Planck calculations, which do not use orbit-averaging, show that in an isolated single-mass cluster the evaporation rate is given by

$$t_{\text{evap}} = -N(dN/dt)^{-1} = ft_{\text{rh}}, \quad (7.141)$$

where t_{rh} is the half-mass relaxation time (eq. 7.108) and $f \approx 300$ (Spitzer 1987). We shall find below that core collapse in such clusters occurs after about 16 times the initial half-mass relaxation time $t_{\text{rh},i}$; thus we expect that only a few percent of the stars in a cluster will escape before core collapse. After core collapse, an isolated cluster starts to expand, the relaxation rate slows, and the rate of escape therefore slows as well. A completely isolated globular cluster would require much longer than the age of the universe to evaporate completely—isolated single-mass clusters lose about 75% of their mass only after $10^7 t_{\text{rh},i}$, almost independent of N (Baumgardt, Hut, & Heggie 2002).

These considerations suggest that tidal forces from the host galaxy play a central role in determining the evaporation rate from a cluster. Truncation of the cluster at the tidal radius (§8.3.1) eliminates Hénon's paradox, because the period of an orbit at the escape energy is finite. Thus the evaporation rate *can* be determined from orbit-averaged Fokker–Planck calculations: for clusters with tidal radii and other properties similar to those observed in Galactic globular clusters, the coefficient f in equation (7.141) lies in the range (Spitzer 1987; Gnedin, Lee, & Ostriker 1999)

$$f \approx 20\text{--}60. \quad (7.142)$$

In contrast to isolated clusters, the expansion of the cluster following core collapse accelerates evaporation, by spilling stars over the tidal radius. Thus the evaporation is slowest at the start, so the timescale required for the cluster to evaporate is given approximately by equations (7.141) and (7.142), with t_{rh} replaced by the initial half-mass relaxation time.

These considerations are based on a simple model of tidal forces, in which stars that wander outside the tidal radius are instantaneously lost from the cluster. Greater accuracy requires the inclusion of tidal shocks (§7.5.6), as well as a more careful treatment of the dynamical process by which stars leak through the tidal boundary and then drift slowly away from the cluster (Fukushige & Heggie 2000; Baumgardt 2001).

In the light of this discussion, it is interesting to re-examine Figure 7.3, which shows the distribution of t_{rh} in Galactic globular clusters. There are almost no globulars with t_{rh} less than about 100 Myr, or about 1% of the age

¹² As Hénon says, “This paradox is somewhat reminiscent of Achilles and the tortoise; unfortunately the conclusion is different.”

of the Galaxy. This observation strongly suggests that there once were many globulars with shorter relaxation times, which all evaporated after $\approx 10^2 t_{\text{rh}}$. Once again, this explanation hints that the number of globular clusters that existed shortly after the Galaxy formed may have been much larger than the present population (see §7.5.6).

The maximum lifetime of a stellar system Galaxies often contain concentrations of 10^6 – $10^9 \mathcal{M}_\odot$ of material in their central few parsecs. These are sometimes called **massive dark objects** because they emit no radiation and are detected only by their effects on the kinematics of nearby stars and gas. Massive dark objects are usually assumed to be black holes, but we now ask whether they could instead be clusters of low-luminosity stars (perhaps brown dwarfs or neutron stars) or stellar-mass black holes.

We have argued in equation (7.141) that an isolated single-mass star cluster will lose a substantial fraction of its mass by evaporation in roughly $300t_{\text{rh}}$, where t_{rh} is the current half-mass relaxation time. This is an upper limit to the lifetime of a cluster, since most of the effects neglected in this calculation speed up the evaporation rate: systems with a range of stellar masses or significant tidal forces have shorter evaporation times; physical collisions and inelastic encounters between the stars (§7.5.8) dissipate energy and hence accelerate cluster evolution; in the densest clusters gravitational radiation is an additional drain on the cluster energy that accelerates the evolution, etc. (Maoz 1991). Unless the cluster formed recently—an improbable coincidence—we must have $300t_{\text{rh}} \gtrsim 10 \text{ Gyr}$, which together with equation (7.108) implies

$$r_{\text{h}} \gtrsim 0.01 \text{ pc} (\ln \lambda M/m)^{2/3} \left(\frac{m}{\mathcal{M}_\odot} \right)^{2/3} \left(\frac{10^8 \mathcal{M}_\odot}{M} \right)^{1/3}. \quad (7.143)$$

Here M and r_{h} are the mass and half-mass radius of the cluster, m is the mass of a typical member, and $\lambda \simeq 0.1$.

The strongest constraint on the properties of a dark central cluster comes from observations of stellar orbits in the Galactic center (Schödel et al. 2002), which show that the central mass $M = 4 \times 10^6 \mathcal{M}_\odot$ is contained within a radius of $\lesssim 0.001 \text{ pc}$. This is consistent with the constraint (7.143) only if its constituents have mass $m \lesssim 2 \times 10^{-4} \mathcal{M}_\odot$. Normal objects of this mass (brown dwarfs or planets) are not allowed, because their collision time would be far shorter than the age of the Galaxy; black holes are allowed but low-mass black holes of this kind are not formed by any known process. Similar constraints can be derived for the massive dark objects in the galaxies NGC 4258 and M31 (Maoz 1995; Bender et al. 2005).

Arguments such as these strongly suggest that the massive dark objects in galactic centers are single (or perhaps binary) black holes, a result that is also implied by the demography of active galactic nuclei (§1.1.6).

7.5.3 Core collapse

The evolution of the mass distribution in an isolated cluster that began as a Plummer model is shown in Figure 7.4. The outer half of the cluster expands, due to the gradual growth of the halo as core stars diffuse towards the escape energy. At the same time, the center contracts; the central 1% of the mass contracts by a factor $k \gtrsim 30$, corresponding to an increase in the central density by a factor $k^3 \gtrsim 3 \times 10^4$. This contraction process, known as **core collapse**, leads to such dramatic growth in the central density that early calculations of the evolution of spherical star clusters, which did not include energy exchange with binary stars (page 559), culminated in an apparent singularity in the central density. This singularity—at which the early numerical codes crashed—occurred at about $16t_{\text{rh},i}$, where $t_{\text{rh},i}$ denotes the initial half-mass relaxation time.

A more accurate calculation of core collapse is shown in Figure 7.5 (Takahashi 1995), based on solving the orbit-averaged Fokker–Planck equation for an isolated, spherical, single-mass cluster without binaries, again starting from a Plummer model. Direct solution of the Fokker–Planck equation allows far greater dynamic range than Monte Carlo or N-body methods, and Takahashi was able to follow the evolution of the central density over a factor exceeding 10^{13} .

The density profiles of these models are shown in the top panel of Figure 7.5. As the cluster evolves, the core radius shrinks and the central density grows. The density profile outside the core approaches a power law $\rho \propto r^{-2.23}$; by the end of the computation this power law extends over more than six orders of magnitude in radius.

The bottom panel shows the behavior of the anisotropy parameter $\beta = 1 - \overline{v_\theta^2}/\overline{v_r^2}$. At large radii, $\beta \simeq 1$, indicating that the orbits are nearly radial, as we expect from the model for the cluster halo developed at the start of this section. At the smallest radii, inside the continually shrinking core, we find $\beta \simeq 0$, indicating that the velocity distribution is isotropic. In the radius range in which the density profile in the top panel is a power law, there is a constant small radial anisotropy, $\beta \simeq 0.08$, or $\overline{v_\theta^2}/\overline{v_r^2} \simeq 0.92$.

To understand the dynamics of core collapse, we shall first investigate the evolution of a system in which the density profile outside the core varies as a power law in radius all the way to infinity. In this limit the density profile evolves self-similarly, that is, profiles at different times differ only in normalization and scale. A self-similar solution of this kind can be written in the form

$$\rho(r, t) = \rho_0(t)\rho_\star(r_\star), \quad (7.144)$$

where

$$r_\star = \frac{r}{r_0(t)}, \quad (7.145)$$

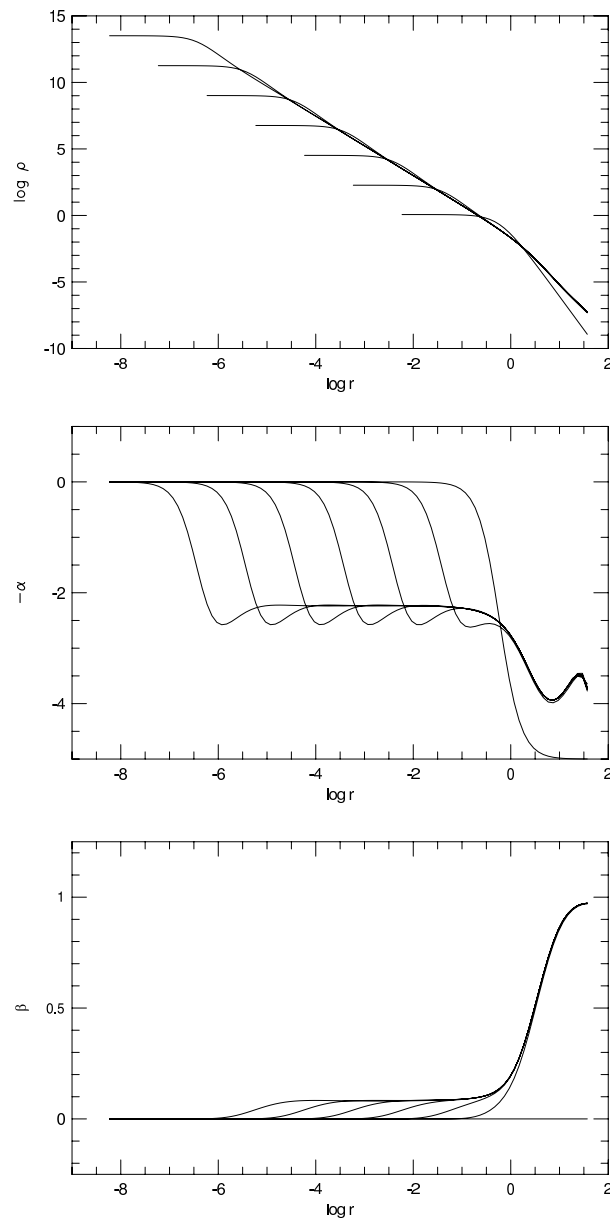


Figure 7.5 An orbit-averaged Fokker-Planck calculation of core collapse. The initial state was an isolated, single-mass Plummer model. Top: evolution of the density profile with time. The central density increases and the core radius decreases with time. Middle: the logarithmic density gradient $d \ln \rho / d \ln r \equiv -\alpha$. The region outside the core takes on a power-law profile with $\alpha = 2.23$. Bottom: evolution of the velocity anisotropy parameter, $\beta = 1 - \overline{v_\theta^2} / \overline{v_r^2}$ (eq. 4.61). From Takahashi (1995), by permission of the Astronomical Society of Japan.

and $\rho_0(t)$ and $r_0(t)$ are the central density and some measure of core radius, evaluated at time t . We shall choose r_0 to be the King radius (eq. 4.106)

$$r_0 \equiv \sqrt{\frac{9\sigma^2}{4\pi G\rho_0}}, \quad (7.146)$$

which is nearly equal to the core radius for the isothermal sphere and many other stellar systems; here $3\sigma^2$ is the mean-square speed at the center.

The top panel of Figure 7.5 shows that the density is nearly independent of time at radii that lie well outside the core, and hence

$$0 = \frac{\partial\rho(r,t)}{\partial t} = \frac{d\rho_0}{dt}\rho_* - \rho_0 \frac{d\rho_*}{dr_*} \frac{r}{r_0^2} \frac{dr_0}{dt} \quad \text{for } r \gg r_0. \quad (7.147)$$

Thus

$$\frac{r_0}{\rho_0} \frac{dt}{dr_0} \frac{d\rho_0}{dt} = \frac{r_*}{\rho_*} \frac{d\rho_*}{dr_*} \quad \text{for } r_* \gg 1. \quad (7.148)$$

Since the left side is a function only of t , and the right side is a function only of r_* , both must be equal to a constant, which we shall call $-\alpha$. Hence

$$\begin{aligned} \rho_*(r_*) &\propto r_*^{-\alpha} \quad \text{for } r_* \gg 1, \\ \rho_0(t) &\propto r_0^{-\alpha}(t). \end{aligned} \quad (7.149)$$

Fitting this form to the slope of the power law in Figure 7.5, we find

$$\alpha = 2.23. \quad (7.150)$$

Larson (1970a) first derived this behavior using fluid models, and found $\alpha \simeq 2.4$; Cohn (1980) found $\alpha = 2.23$ from a solution of the isotropized Fokker–Planck equation; and Lynden–Bell & Eggleton (1980) found $\alpha = 2.21$ from fluid models.

Let us define the core mass to be $M_0(t) \equiv \rho_0(t)r_0^3(t)$. From the second of equations (7.149) we find $M_0(t) \propto r_0^{3-\alpha}$. According to equation (7.146) the central velocity dispersion is $\sigma \propto \rho_0^{1/2}r_0 \propto r_0^{1-\alpha/2}$. From equation (7.106), the relaxation time in the core is $t_{\text{relax}} \propto \sigma^3/\rho_0 \propto r_0^{3-\alpha/2}$ (neglecting variations in $\ln\Lambda$). Since core collapse is driven by relaxation, we expect the characteristic timescale for changes in the core to be comparable to the relaxation time in the core. Thus we may write

$$\frac{1}{r_0} \frac{dr_0}{dt} \propto \frac{1}{t_{\text{relax}}} \propto r_0^{\alpha/2-3}. \quad (7.151)$$

This equation is easily solved to yield

$$r_0(t) \propto (t_{\text{cc}} - t)^{2/(6-\alpha)} \propto \tau^{0.53}, \quad (7.152)$$

where t_{cc} is a constant of integration representing the moment of collapse, $\tau = t_{\text{cc}} - t$ is the time remaining until collapse, and in the second proportionality we have inserted the value of α from equation (7.150). Similarly,

$$\begin{aligned} \rho_0(t) &\propto \tau^{-2\alpha/(6-\alpha)} \propto \tau^{-1.18} & ; & \quad \sigma^2(t) \propto \tau^{(4-2\alpha)/(6-\alpha)} \propto \tau^{-0.12}, \\ M_0(t) &\propto \tau^{(6-2\alpha)/(6-\alpha)} \propto \tau^{0.41} & ; & \quad t_{\text{relax}}(r=0) \propto \tau. \end{aligned} \quad (7.153)$$

The last equation implies that the time to core collapse is always a fixed multiple of the central relaxation time. The proportionality constant is

$$\tau \simeq 300 t_{\text{relax}}(r=0). \quad (7.154)$$

For typical clusters, this simple relation holds only during the late stages of core collapse, after the central density has grown by a factor 10^3 or more; at earlier times $\tau/t_{\text{relax}}(r=0)$ is generally smaller (Quinlan 1996a).

Core collapse is a manifestation of the gravothermal catastrophe of §7.3.2. In most clusters the velocity dispersion decreases outwards; that is, the inner parts of the cluster are hotter than the outer parts. Therefore star-star encounters transfer energy outwards. Since the cluster core has negative heat capacity, it grows hotter as it loses energy. As it grows hotter, the core shrinks in size and mass and grows in density; if no other process intervenes, the core collapses to zero radius and infinite density in a finite time.¹³

7.5.4 After core collapse

What eventually halts core collapse? We remark first that despite the formal density singularity, core collapse is a relatively unspectacular process. To see this, let us employ the similarity relations to follow core collapse in a typical globular cluster. Eliminating the time τ between the relations for r_0 , M_0 , and σ_0 in equations (7.152) and (7.153), we have

$$r_0 \propto M_0^{1/(3-\alpha)} \propto M_0^{1.30} \quad ; \quad \sigma_0 \propto M_0^{(2-\alpha)/(6-2\alpha)} \propto M_0^{-0.15}. \quad (7.155)$$

Let us assume that core collapse starts when the cluster has parameters given by Table 1.3. Thus, initially, the central density $\rho_0 \simeq 1 \times 10^4 \mathcal{M}_\odot \text{pc}^{-3}$, the velocity dispersion $\sigma \simeq 6 \text{ km s}^{-1}$, and the core or King radius $r_0 \simeq 1 \text{ pc}$. The initial core mass is thus $M_0 = \rho_0 r_0^3 \simeq 1 \times 10^4 \mathcal{M}_\odot$; assuming a typical stellar mass $m = 0.7 \mathcal{M}_\odot$, the initial number of stars in the core is $M_0/m \simeq 1.4 \times 10^4$. Then equation (7.155) yields

$$r_0 \simeq 4 \times 10^{-6} N_0^{1.30} \text{ pc} \quad ; \quad \sigma \simeq 25 N_0^{-0.15} \text{ km s}^{-1}, \quad (7.156)$$

¹³ Models in which the inner parts are colder than the outer parts, such as the Hernquist model shown in Figure 4.4, exhibit core expansion rather than core collapse (Hachisu et al. 1978; Quinlan 1996a).

where N_0 is the number of stars remaining in the core. Thus, even when there is only a handful of stars left in the core (say, $N_0 = 10$), $r_0 \simeq 2 \times 10^9$ km (somewhat larger than the distance between the Sun and Saturn) and $\sigma \simeq 20$ km s⁻¹. Clearly, the statistical approximation that $N \gg 1$ fails long before the density is high enough for stellar coalescence, relativistic effects, or other exotic phenomena to become important.

It turns out that core collapse is eventually halted by binary stars. In most clusters these will be primordial (§7.1f), but even if no primordial binaries are present, we now show that binaries are formed by the core-collapse process.¹⁴ We have seen that the number of binaries formed by three-body encounters is about $0.1(N_0 \ln N_0)^{-1}$ per relaxation time (eq. 7.13). Since the characteristic evolution time during core collapse is $\approx 300t_{\text{relax}}(N_0)$ (eq. 7.154), the interval in which the core has about N_0 stars is roughly $300t_{\text{relax}}(N_0)$, and during this time we expect $\approx 30(N_0 \ln N_0)^{-1}$ binary stars to be formed. Thus the first binary forms by a three-body encounter when $N_0 \ln N_0 \approx 30$, corresponding to $N_0 \approx 12$ stars left in the core. These results suggest that core collapse is halted when the core contains only 10–20 stars, or earlier if primordial binaries are present.

How does binary formation halt core collapse? Consider a three-body interaction of stars with initial kinetic energies K_i , $i = 1, 2, 3$. We assume that after the interaction, stars 1 and 2 form a binary; the kinetic energy of the center of mass of the binary is \bar{K}_b , the internal energy of the binary is $E_b < 0$, and the kinetic energy of the third star is K'_3 . Conservation of energy requires that $\bar{K}_b + E_b + K'_3 = K_1 + K_2 + K_3$, so $\bar{K}_b + K'_3 > K_1 + K_2 + K_3$. Thus the kinetic energy stored in the centers of mass of the single star and the binary is larger than the initial kinetic energy of the three stars: in other words, *the formation of binary stars provides a heat source for the cluster*. Core collapse occurs because the inner parts of the cluster have negative heat capacity, and evolve by losing energy to the halo and thereby growing hotter; any heat source such as binary-star formation adds energy to the core, cooling it until the temperature gradient declines to zero, thus halting the collapse. Later, we shall derive a quantitative formula (7.184) for the rate of energy generation per unit mass due to binary formation.

A second heat source is interactions of field stars with primordial binaries. These interactions generate heat at a rate per unit mass given by equation (7.182).

How does the cluster evolve after core collapse? The binaries formed in core collapse reside close to the cluster center, and therefore pump kinetic energy into single stars that pass near the center. This energy is then shared among the other cluster stars through two-body relaxation. The resulting

¹⁴ Remarkably, two decades before the first reliable post-collapse simulations, Hénon (1961) already recognized that the boundary conditions for cluster evolution required an energy source at $r = 0$ and argued that the required energy might be produced by the formation of binary and multiple stars.

evolution of the cluster is self-similar, with the cluster radius, mass, mean density, and velocity dispersion evolving as (Problem 7.12b)

$$R \propto \tilde{\tau}^{2/3} \quad ; \quad M \propto \text{constant} \quad ; \quad \rho \propto \tilde{\tau}^{-2} \quad ; \quad v \propto \tilde{\tau}^{-1/3}, \quad (7.157)$$

where $\tilde{\tau}$ is the time since core collapse (Goodman 1984).

The cluster energy grows during this expansion as $E \approx -GM^2/R \propto -\tilde{\tau}^{-2/3}$. On the other hand, the rate of heat generation by binaries depends strongly on the density in the small central region where they are concentrated. Thus the central density adjusts itself so that the rate of energy generation matches the demands of the similarity solution (Hénon 1961). A similar situation occurs in stars, where the luminosity is determined by the efficiency of radiative and convective energy transport to the photosphere, and the core temperature rises or falls so that the rate of energy generation by nuclear reactions matches this demand.

This simple picture is incomplete. It turns out that the self-similar post-collapse evolution described by equations (7.157) is unstable, at least for $N \gtrsim 8000$ in the case of a single-mass cluster. The instability leads to **gravothermal oscillations**, in which the central density and core radius oscillate, in some cases by several orders of magnitude. These oscillations are shown in Figure 7.6, which was computed by solving the Fokker–Planck equation including a heat source due to binaries. All of the panels show identical evolution up to core collapse, which occurs at $t/t_{\text{rh},i} \simeq 16$. After core collapse, the model with $N = 8000$ displays the self-similar behavior $\rho_0 \propto \tilde{\tau}^{-2}$ predicted by equation (7.157). However, all of the models with $N \gtrsim 8000$ exhibit dramatic gravothermal oscillations superimposed on this self-similar decline in central density. The oscillations have sharp spikes in the central density, suggesting that the core is re-collapsing at each oscillation.

As N increases, the character of gravothermal oscillations changes in several ways: (i) their amplitude grows, to over four orders of magnitude in central density by $N = 10^5$ and six orders of magnitude by $N = 10^6$; (ii) the oscillations become increasingly chaotic, and (iii) the fraction of each oscillation that is spent in the high-density state becomes smaller and smaller. One consequence of the last of these is that we are much more likely to observe a post-collapse cluster in a low-density state than in a high-density state; thus, although many of the globular clusters in the Galaxy have suffered core collapse, the chance of finding one with near-zero core radius is small.

Gravothermal oscillations were first discovered and analyzed using fluid models (Sugimoto & Bettwieser 1983; Goodman 1987), and by now have been found in Monte Carlo simulations, numerical solutions of the Fokker–Planck equation, and even N-body integrations (Makino 1996). The unstable oscillations set in at around $N \simeq 8000$ in Figure 7.6; however, in N-body simulations the oscillations are present at somewhat smaller N because the energy input from binaries is stochastic, leading to noise that drives oscillations even when they are weakly stable. A more detailed description of gravothermal oscillations is given by Heggie & Hut (2003).

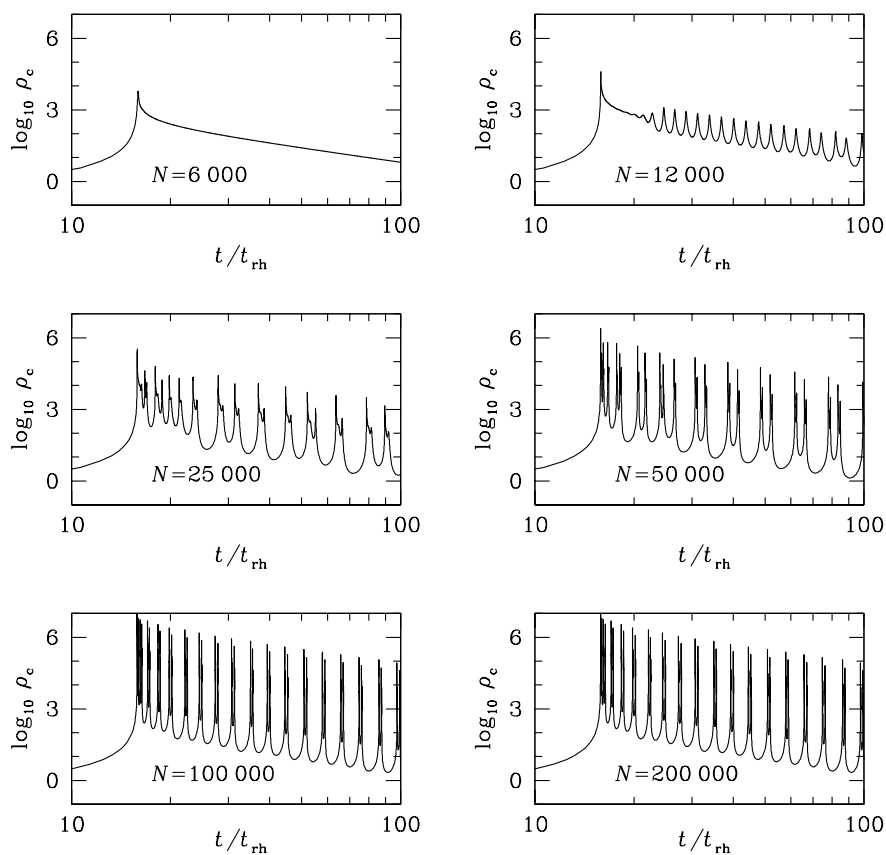


Figure 7.6 Gravothermal oscillations in a star cluster (Breedon & Cohn 1995). The central density is plotted against time in units of the initial half-mass relaxation time $t_{rh,i}$ (eq. 7.108). Each panel is labeled by the number of stars in the cluster, which determines the relative rates of relaxation and binary heating. The initial spike in the central density represents core collapse. The post-collapse evolution of clusters with $N \lesssim 8000$ is stable, and follows the prediction of equation (7.157). The oscillations in unstable clusters become stronger as N increases.

7.5.5 Equipartition

So far, our analysis has been based on idealized single-mass models. If a range of stellar masses is present, encounters tend to establish equipartition of energy. The more massive stars lose energy and sink towards the center, while the lighter stars gain energy, so their orbits expand. The equipartition

timescale is comparable to t_{rh} . Understanding the degree of segregation of stars by mass is critical for the construction of any model of a globular cluster that is to be compared to observations.

The familiar concept of “energy equipartition” must be applied with care in self-gravitating systems. Consider two populations of stars, with masses m_1 and $m_2 > m_1$, in a star cluster with potential $\Phi(\mathbf{x})$. The mean energy per star in population i may be written $\langle E \rangle_i = m_i \langle \frac{1}{2}v^2 + \Phi(\mathbf{x}) \rangle_i$, where the average is over the stars in population i . Encounters do *not* establish equipartition in the sense that $\langle E \rangle_1 \rightarrow \langle E \rangle_2$; indeed, this relation is physically meaningless since it can always be satisfied when $m_1 \neq m_2$ by adding an appropriate constant to $\Phi(\mathbf{x})$, a change that has no effect on the dynamics. Encounters *do* tend to establish equipartition of kinetic energy between the two populations at a given location, so if population 2 has a larger value of $m_i \langle v^2 \rangle$ at a given location than population 1, then energy is transferred from population 2 to population 1. However, as orbital energy is shorn from the heavier stars they sink lower in the potential $\Phi(\mathbf{x})$, where the orbital speeds are likely to be higher: thus equipartition can actually lead to a growth in the kinetic energy of the more massive stars.

It is instructive to explore these arguments using a simplified analysis of equipartition in a spherical cluster. The cluster is assumed to contain just two types of stars, with masses m_1 and $m_2 \gg m_1$. Moreover the total mass of the heavy stars, M_2 , is assumed to be much smaller than the core mass of the system of lighter stars, which we write as $\rho_{c1} r_{c1}^3$. Equipartition causes the heavy stars to sink to the center of the core of light stars, where they form a small, dense subsystem. The virial theorem for the heavy stars may be written in the form (Problem 7.4)

$$2K_2 + W_2 - 4\pi G \int_0^\infty dr r^2 \rho_2(r) \frac{M_1(r)}{r} = 0. \quad (7.158)$$

Here $\rho_1(r)$ and $\rho_2(r)$ are the densities of the light and heavy stars, and $M_1(r) = 4\pi \int_0^r dr r^2 \rho_1(r)$ is the mass of light stars interior to r . Since the total mass in heavy stars is small, the heavy system will not strongly perturb the core of light stars, and $\rho_1(r)$ will be approximately constant, $\rho_1(r) \equiv \rho_{c1}$. We may also write $W_2 = -fGM_2^2/r_{h2}$, where r_{h2} is the half-mass radius of the heavy stars and f is a dimensionless constant that is approximately 0.45 for many systems (see eq. 4.249b). Thus

$$\langle v^2 \rangle_2 = f \frac{GM_2}{r_{h2}} + \frac{4\pi G \rho_{c1}}{3} \langle r^2 \rangle_2 = f \frac{GM_2}{r_{h2}} + \frac{4\pi}{3} g^2 G \rho_{c1} r_{h2}^2; \quad (7.159)$$

here $\langle \cdot \rangle_i$ denotes an average over population i , and $\langle r^2 \rangle_2$, the mean-square radius of the heavy stars, has been written as $g^2 r_{h2}^2$ with g a dimensionless constant of order unity.

In equipartition, the mean kinetic energy of the two populations must be the same; thus $m_2 \langle v^2 \rangle_2 = m_1 \langle v^2 \rangle_1 = 3m_1 \sigma^2$, where σ represents the

one-dimensional dispersion in the core of light stars. Using equation (7.146) to express σ in terms of the central density ρ_{c1} and King radius r_{c1} of the light stars, we find

$$\frac{4\pi}{3} \frac{m_1}{m_2} G \rho_{c1} r_{c1}^2 = f \frac{GM_2}{r_{h2}} + \frac{4\pi}{3} G \rho_{c1} g^2 r_{h2}^2. \quad (7.160)$$

The right side, considered as a function of r_{h2} , has a minimum value equal to $(9\pi f^2 g^2 G^3 M_2^2 \rho_{c1})^{1/3}$, which occurs at $r_{h2} = (3fM_2/8\pi\rho_{c1}g^2)^{1/3}$. Hence equipartition cannot be achieved unless the value of the left side exceeds this minimum, which leads to the inequality

$$\frac{M_2}{\rho_{c1} r_{c1}^3} \leq \frac{1.61}{fg} \left(\frac{m_1}{m_2} \right)^{3/2}. \quad (7.161)$$

There is a simple physical explanation of why equipartition cannot be achieved when this inequality is violated. If the mass in heavy stars is too large, they form an independent self-gravitating system at the center of the core of light stars. Encounters cause the system of heavy stars to lose energy to the light stars. According to the virial theorem, this energy loss causes the velocity dispersion of the heavies to increase, so they evolve away from, not towards, equipartition, and the process of energy loss, heating, and contraction of the heavy system must continue indefinitely. This phenomenon is sometimes called the **equipartition instability** (Spitzer 1969, 1987).

In realistic systems with a distribution of stellar masses, the effect of this instability is to produce a dense central core of heavy stars, which contracts independently from the rest of the core. As this core becomes denser and denser, the gravothermal instability eventually dominates over the equipartition instability, and the core collapses in much the same way as the core in a single-component stellar system. The equipartition instability thus accelerates core collapse, so the time required for core collapse is much shorter in a system with a distribution of stellar masses than in a single-mass system. For typical mass distributions the time to core collapse is only $2-4t_{\text{rh},i}$, compared to $16t_{\text{rh},i}$ in a single-mass cluster, where $t_{\text{rh},i}$ is the initial half-mass relaxation time of equation (7.108).

Equipartition also causes the lighter stars in the cluster to evaporate more quickly. Consequently, relaxed globular clusters are expected to be depleted in stars with masses less than a few tenths of a solar mass. Since these stars have high mass-to-light ratios (see Table 3.13 or Figure 5.11 of BM), their loss tends to lower the mass-to-light ratio of the remaining cluster. This process may help to explain why the mass-to-light ratios of globular clusters ($\Upsilon \simeq 2\Upsilon_{\odot}$, see Table 1.3) are lower than those of other old stellar systems ($\Upsilon \simeq 10\Upsilon_{\odot}$ in the centers of ellipticals and the bulges of spiral galaxies).

7.5.6 Tidal shocks and the survival of globular clusters

A tidal shock is a rapidly changing external gravitational field that accelerates stars in the outer parts of the cluster, leading to the expansion of the cluster and escape of some of its members. The dynamics of tidal shocks is discussed in §8.2. Tidal shocks are important for a wide variety of satellite stellar systems. In galaxies like our own, the strongest tidal shocks on globular clusters are generated as the cluster crosses the galactic disk or passes through pericenter on a highly eccentric orbit (§8.2.2f).

Tidal shocks speed up core collapse and shorten cluster lifetimes (e.g., Aguilar, Hut & Ostriker 1988; Murali & Weinberg 1997a, 1997b; Gnedin, Lee, & Ostriker 1999). To illustrate the effect of tidal shocks on the globular cluster population, we may parametrize clusters by their total mass M , half-mass radius r_h , and orbital radius R , and examine their destruction by three mechanisms: (i) evaporation, which occurs in a time $ft_{\text{rh},i}$, where $t_{\text{rh},i}$ is the initial half-mass relaxation time and $f \approx 20\text{--}60$ (eqs. 7.108 and 7.142); (ii) disk and bulge shocks; (iii) dynamical friction, a process introduced on page 583 and described fully in §8.1, which causes massive clusters to spiral in towards the galactic center at a rate determined by their mass and orbital radius (eq. 8.24).

Figure 7.7 shows the contours in the (r_h, M) plane at which the lifetime from the combined effects of all three destruction mechanisms equals 10 Gyr, for several values of the orbital radius (Tremaine 1975; Fall & Rees 1977; Gnedin & Ostriker 1997; Fall & Zhang 2001). Clusters inside the roughly triangular **survival area** are expected to survive to the present time. The figure also plots the locations of the Galactic globular clusters. Overall, the survival area contains the observed clusters rather well. The presence of a handful of clusters outside the survival area is not surprising, since we expect to find some clusters nearing the ends of their lives if destruction is a continuous process. The figure also shows that the absence of clusters with $M \gtrsim 3 \times 10^6 M_\odot$ is not due to any known destruction process, and hence presumably reflects the initial mass distribution of the clusters. The fact that the bulk of the clusters are crowded near the bottom apex of the survival area suggests that a much more extensive initial population of low-mass globular clusters has gradually been whittled away by relaxation and tidal shocks. Observational support for this view comes from the large numbers of short-lived young globular clusters found in interacting and merging galaxies (Whitmore & Schweizer 1995; Goudfrooij et al. 2004; Larsen 2004); presumably these clusters formed during the merger and have not yet been destroyed by dynamical processes. Two corollaries are that the properties of the observed cluster distribution are largely determined by the constraint that clusters must survive for 10 Gyr, and that much of the stellar halo may consist of debris from disrupted globular clusters.

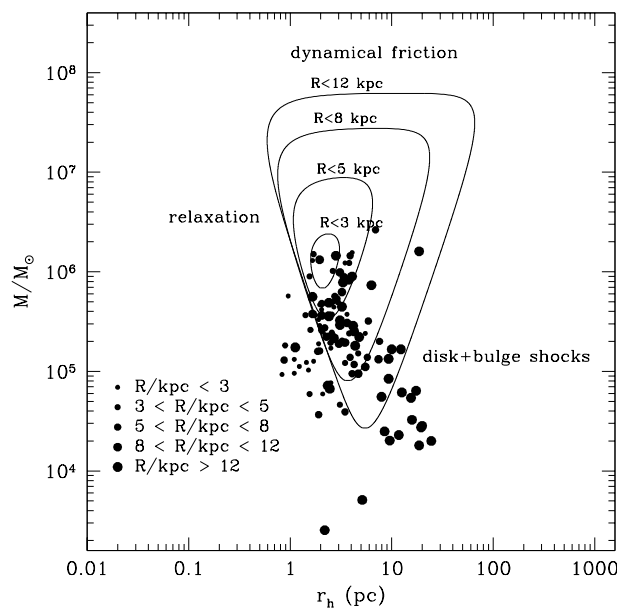


Figure 7.7 The distribution of Galactic globular clusters in mass M , half-mass radius r_h , and orbital radius R (kpc). The different dot sizes denote clusters in different radial ranges. The contours represent locations at which the expected lifetime of a cluster is 10 Gyr at a given orbital radius. The destruction mechanisms are evaporation, dynamical friction, and tidal shocks. From Gnedin & Ostriker (1997).

7.5.7 Binary stars

In this subsection we investigate the formation, evolution, and destruction of binary stars as a result of gravitational encounters with other cluster stars. Our principal interest is in the role played by binaries in the evolution of stellar systems (Hut et al. 1992). For a more thorough analysis see Valtonen & Karttunen (2006).

Let us consider a homogeneous stellar system consisting of single (or “field” stars) and binary stars. The binaries form and dissolve through gravitational interactions with the field stars. For simplicity we shall assume that all the stars have the same mass m . The field stars have density ρ and velocity dispersion σ , with a Maxwellian velocity distribution (eq. 7.91). We assume that the velocity distribution of the centers of mass of the binary stars is also Maxwellian, with velocity dispersion σ_b . If there is equipartition of kinetic energy between the field stars and the binary centers of mass, $\frac{1}{2}m\sigma^2 = \frac{1}{2}(2m)\sigma_b^2$, then $\sigma_b = \sigma/\sqrt{2}$.

Let $\mathbf{x} \equiv \mathbf{x}_1 - \mathbf{x}_2$ be the separation vector between the components of a bi-

nary, and let $\mathbf{V} = \dot{\mathbf{x}}$ be their relative velocity. According to equation (D.34), the internal energy of the binary (i.e., its total energy minus the kinetic energy of the center-of-mass motion) is

$$\tilde{E} = \frac{1}{2}\mu V^2 - \frac{Gm^2}{r}, \quad (7.162)$$

where $\mu = \frac{1}{2}m$ is the reduced mass. The tilde on \tilde{E} is a reminder that in this section “energy” denotes a quantity having units mass \times (velocity)², rather than the specific energy (energy per unit mass) with units (velocity)² that is common elsewhere in the book.

The separation vector \mathbf{x} satisfies the equation of motion of a fictitious “reduced particle” of mass μ orbiting in the potential $-Gm^2/r$ (Appendix D.1). Hence it follows a Keplerian ellipse with semi-major axis a , where

$$\tilde{E} = -\frac{Gm^2}{2a}. \quad (7.163)$$

For a circular orbit the relative speed is $V = (2Gm/a)^{1/2}$.

A binary is called **soft** if $|\tilde{E}|/m\sigma^2 < 1$ and **hard** if $|\tilde{E}|/m\sigma^2 > 1$.¹⁵ The behaviors of soft and hard binaries are quite different, and they will be analyzed separately below. We shall often concentrate on the behaviors of very soft or very hard binaries, which are simpler to understand, and extrapolate our results to the transition region near $|\tilde{E}| \approx m\sigma^2$.

In §7.1 we made a rough estimate of the rate of binary formation by three-body encounters (eq. 7.11). We shall now be more precise. Let us assume that there is equipartition of kinetic energy between the field stars and the binary centers of mass, and let $C(\tilde{E})d\tilde{E}$ be the rate per unit volume at which binaries are created with energy in the range $(\tilde{E}, \tilde{E} + d\tilde{E})$, $\tilde{E} < 0$. Since formation requires a three-body encounter, we expect that $C(\tilde{E}) \propto \rho^3$, where $\rho = mn$ is the density of field stars. Similarly, let $B(\tilde{E})$ be the destruction rate for a binary of energy \tilde{E} ; since binaries are destroyed by encounters with field stars, we expect that $B(\tilde{E}) \propto \rho$. We show in Appendix M.4 that these rates are related by the principle of detailed balance, which requires that

$$C(\tilde{E}) = \frac{\pi^{3/2}G^3\rho^2m^{5/2}}{4\sigma^3|\tilde{E}|^{5/2}}e^{|\tilde{E}|/m\sigma^2}B(\tilde{E}). \quad (7.164)$$

Encounters with field stars tend to make the distribution of binaries in phase space uniform on a given constant-energy surface. If this uniformity

¹⁵ These definitions need to be modified if one or both of the binary components is much more massive than the field stars. See Quinlan (1996b).

is achieved, the probability that a binary star has eccentricity in the range $(e, e + de)$ is equal to $2e de$ (Problem 4.8).

(a) Soft binaries The evolution of binaries is easiest to describe in the limit in which they are very soft, $|\tilde{E}| \ll m\sigma^2$. Consider an encounter of star 1 with a field star, at an impact parameter that is much less than the binary separation. Then the encounter changes the velocity of star 2 much less than the velocity of star 1. We may neglect the orbital motion of the binary during the encounter—the binary is soft, so the relative speed of the binary components is much less than their velocity relative to the field star—and write the change of internal energy in the encounter as

$$\Delta\tilde{E} = \frac{1}{2}\mu\Delta(V^2) = \frac{1}{2}m[\Delta\mathbf{v}_1 \cdot (\mathbf{v}_1 - \mathbf{v}_2) + \frac{1}{2}(\Delta\mathbf{v}_1)^2]. \quad (7.165)$$

The mean value of $\Delta\tilde{E}$ per unit time can be obtained by replacing $\Delta\mathbf{v}_1$ and $\Delta\mathbf{v}_1^2$ by the diffusion coefficients (7.89). Thus we set

$$D[\Delta\mathbf{v}_1] = \frac{\mathbf{v}_1}{v_1}D[\Delta v_{\parallel}] \quad ; \quad D[(\Delta\mathbf{v}_1)^2] = D[(\Delta\mathbf{v}_{\perp})^2] + D[(\Delta v_{\parallel})^2]. \quad (7.166)$$

Substituting into (7.165) and using equation (7.92), we find

$$D[\Delta\tilde{E}] = \frac{2\pi G^2 m^2 \rho \ln \Lambda}{\sigma} \left[-2 \frac{\mathbf{v}_1 \cdot (\mathbf{v}_1 - \mathbf{v}_2)}{v_1 \sigma} G(X_1) + \frac{1}{\sqrt{2}} \frac{\text{erf}(X_1)}{X_1} \right], \quad (7.167)$$

where $X_1 = v_1/(\sqrt{2}\sigma)$, $\ln \Lambda$ is the Coulomb logarithm, and the function $G(X_1)$ is defined by equation (7.93). For a very soft binary $|\mathbf{v}_1 - \mathbf{v}_2| \ll \sigma$, so the first term may be neglected in comparison with the second; moreover, we can replace v_1 by the speed of the binary center of mass v_{cm} . Finally we double the value of $D[\Delta\tilde{E}]$ to account for encounters of field stars with star 2. Thus for very soft binaries

$$D[\Delta\tilde{E}] = \frac{2^{3/2}\pi G^2 m^2 \rho \ln \Lambda}{\sigma} \frac{\text{erf}(X_{\text{cm}})}{X_{\text{cm}}}, \quad (7.168)$$

where $X_{\text{cm}} = v_{\text{cm}}/\sqrt{2}\sigma$. We now average over v_{cm} , assuming for simplicity that the binary center-of-mass motion is in equipartition with the field stars, so $\sigma_b = \sigma/\sqrt{2}$. Using the relation

$$\begin{aligned} \left\langle \frac{\text{erf } X_{\text{cm}}}{X_{\text{cm}}} \right\rangle &\equiv \frac{\int d^3\mathbf{v}_{\text{cm}} \exp(-\frac{1}{2}v_{\text{cm}}^2/\sigma_b^2) \text{erf } X_{\text{cm}}/X_{\text{cm}}}{\int d^3\mathbf{v}_{\text{cm}} \exp(-\frac{1}{2}v_{\text{cm}}^2/\sigma_b^2)} \\ &= \frac{\int_0^\infty X \text{erf}(X) e^{-2X^2} dX}{\int_0^\infty dX X^2 e^{-2X^2}} = \sqrt{\frac{8}{3\pi}}, \end{aligned} \quad (7.169)$$

we obtain

$$\langle D[\Delta\tilde{E}] \rangle = 8\sqrt{\frac{\pi}{3}} \frac{G^2 m^2 \rho}{\sigma} \ln \Lambda. \quad (7.170)$$

To ensure that our expressions for the diffusion coefficients are valid, we must have $1 \ll \Lambda = b_{\max} v_{\text{typ}}^2 / (2Gm)$ (eq. 7.84). For binary stars, we set the maximum impact parameter equal to half the binary semi-major axis, $b_{\max} = \frac{1}{2}a$, since our derivation is valid only for encounters that are much closer to one star than the other (encounters with impact parameter $b \gg a$ affect the center-of-mass velocity of the binary much more strongly than the relative velocity). We also replace v_{typ} by 2.1σ , the RMS relative velocity between the binary center of mass and the field stars. After these approximations $\Lambda \approx 1.1a\sigma^2 / (Gm)$. Using equation (7.163) we can write $\Lambda \approx 0.6m\sigma^2 / |\tilde{E}|$; thus the condition that the binary is very soft guarantees that $\Lambda \gg 1$.

Equation (7.170) implies that on average soft binaries gain energy from encounters with field stars; in other words, *soft binaries become softer*. This result can be interpreted in terms of energy equipartition. Since the binary is very soft, its internal kinetic energy is much less than the kinetic energy of the field stars. The growth of its internal energy is a manifestation of the tendency of the system to evolve towards equipartition, although the negative specific heat of the binary thwarts this tendency.

We now examine how soft binaries are disrupted. We argued in §7.5.2 that stars escape from a cluster by two distinct mechanisms, ejection and evaporation. The same two processes describe the destruction of soft binaries by encounters with field stars: in ejection—often called **ionization** in the context of binary stars—a single close encounter with a field star leaves the binary with positive internal energy, while in evaporation a series of more distant encounters gradually increase the internal energy until it is positive.

The ejection or ionization rate for very soft binaries is (Heggie 1975)

$$B(\tilde{E}) = \frac{40\sqrt{\pi}}{3\sqrt{3}} \frac{G^2 m^2 \rho}{\sigma |\tilde{E}|}, \quad (7.171)$$

where once again we have assumed that all masses are equal and that the binary center-of-mass motion is in equipartition with the field stars. The expected lifetime before ejection is

$$t_{\text{ej}} = \frac{1}{B(\tilde{E})} = 0.037 \frac{\sigma}{G\rho a}. \quad (7.172)$$

The evaporation rate can be determined from the energy diffusion coefficient $D[\Delta\tilde{E}]$ (eq. 7.170); since this is independent of \tilde{E} except for a weak dependence on the Coulomb logarithm, the lifetime of a very soft binary of energy \tilde{E} before it evaporates is simply

$$t_{\text{evap}} = \frac{|\tilde{E}|}{\langle D[\Delta\tilde{E}] \rangle} = \frac{\sqrt{3}\sigma}{16\sqrt{\pi}G\rho a \ln \Lambda} = 0.061 \frac{\sigma}{G\rho a \ln \Lambda}. \quad (7.173)$$

The evaporation time is shorter than the ejection time (7.172) by of order the Coulomb logarithm $\ln \Lambda \approx \ln(0.6m\sigma^2/|\tilde{E}|)$; thus evaporation dominates for very soft binaries with $|\tilde{E}| \ll m\sigma^2$, while ejection and evaporation occur at comparable rates when $|\tilde{E}| \sim m\sigma^2$.

The evaporation time t_{evap} is also much shorter than the local relaxation time t_{relax} (eq. 7.106) for very soft binaries, since the latter is the average time required for changes in energy of order $m\sigma^2$, while disruption of a very soft binary requires only changes of order $|\tilde{E}| \ll m\sigma^2$.

In relaxed clusters the rate of evaporation of soft binaries is balanced by the rate of soft binary formation due to three-body encounters. It turns out that the equilibrium number of very soft binaries resulting from these processes is of order unity, independent of the number of stars in the cluster (Appendix M). Thus, soft binaries play no significant role in the evolution of stellar systems. Nevertheless, the process of disruption of soft binary stars is important to understand in other contexts, for example to constrain the properties of MACHOs in the local dark halo (§8.2.2e).

The disruption of soft binaries by much more massive field objects is discussed in §8.2.2d.

(b) Hard binaries The interaction of a hard binary with a field star can be extremely complex. The reason for this complexity is that the binding energy $|\tilde{E}|$ of a hard binary is larger than the typical kinetic energy $\sim m\sigma^2$ of the binary center of mass and the field star. Hence the three stars can temporarily form a bound three-body system, a possibility that is not available to soft binaries.

Figure 7.8 shows an interaction of a hard binary and a field star. This encounter is an example of an **exchange**: the original binary containing stars 1 and 2 is dissociated, and stars 2 and 3 join to produce a new binary, while star 1 escapes to infinity. The encounter is impossible to describe simply; most of the time the system displays a hierarchy in which one star travels in an elongated Keplerian ellipse with a tightly bound binary at one focus. At each pericenter passage of the outermost star, the three stars interact strongly, one of the three stars is flung into an elongated orbit, and the process repeats. Ultimately one star escapes from the triple, and the three-body interaction ends. It should be stressed that the trajectory in Figure 7.8 is actually *less* complicated than the typical interaction between a very hard binary ($|\tilde{E}| \sim 50m\sigma^2$) and a field star (Hut & Bahcall 1983).

Another possibility is a **flyby**, in which the outgoing state is the same as the incoming state; for example, a binary containing stars 1 and 2 approaches a single star 3, and after the interaction stars 1 and 2 are still bound and 3 is still single, although with different kinetic energies. The third and final possibility is ionization, in which all three stars leave the encounter as single stars; ionization is negligible for very hard binaries since it requires that the total center-of-mass energy of the hard binary plus the field star is positive (see eq. 7.174 below).

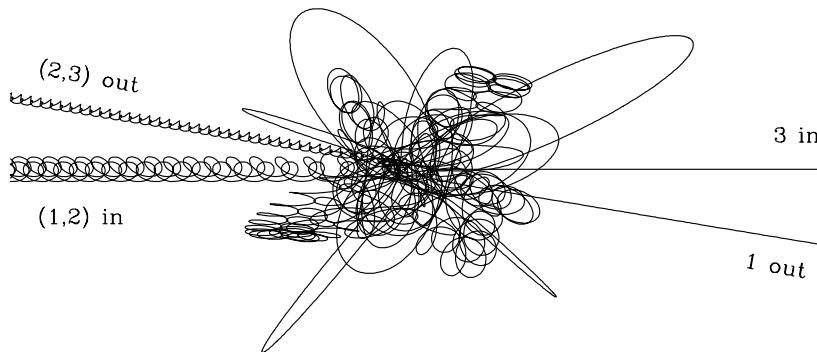


Figure 7.8 An interaction between a hard binary and a field star. All three stars have equal mass and the orbits are plotted in the center-of-mass frame. The binary, containing stars 1 and 2, enters from the left; the single star (labeled 3) enters from the right. After a complicated interaction, star 1 escapes, leaving 2 and 3 behind as a newly formed binary. After Hut & Bahcall (1983).

A heuristic argument can be used to determine the effect of encounters on the internal energy $\tilde{E} < 0$ of hard binaries. The initial relative velocity of the binary and the field star is of order the cluster dispersion σ , which is much less than the orbital speed in a very hard binary. After a complex exchange interaction like the one shown in Figure 7.8, the ejected star will typically have a speed of order the orbital speed of the initial binary. Hence the outgoing single star will have a higher speed than the incoming single star. By energy conservation, the internal energy \tilde{E} of the binary must therefore decrease, that is, $|\tilde{E}|$ increases; in other words, *hard binaries become harder*.

Combining this result with our conclusions on soft binaries, we arrive at **Heggie’s law**: on average, *hard binaries get harder and soft binaries get softer* (Heggie 1975; Hills 1975). There is a “watershed” energy, near $-m\sigma^2$, at which the average rate of energy input from encounters is zero.

(c) Reaction rates Quantitative estimates of the rate of formation, evolution, and destruction of hard binaries have been made by many authors (Heggie 1975; Hut & Bahcall 1983; Sigurdsson & Phinney 1993; Valtonen & Karttunen 2006). In most cases these are based on analytic formulae that are empirically adjusted to match numerical experiments. We shall now list some of these results; we continue to restrict ourselves to the case in which all of the stars have a single mass m and the velocity dispersion of the center of mass of the binaries has the equipartition value, $1/\sqrt{2}$ times the dispersion of the field stars.

We first consider the ionization rate, in which an encounter with a field

star destroys the binary to leave three single stars. Hard binaries are much more difficult to ionize than soft binaries: ionization requires that the total energy of the binary plus the field star in the center-of-mass frame is positive, and for very hard binaries this is true only for the very few field stars whose speed is much larger than the velocity dispersion σ .

The probability per unit time that a given binary is ionized by a single close encounter may be written as

$$B(\tilde{E}) = \frac{8\sqrt{\pi}G^2m^3\rho\sigma}{3^{3/2}|\tilde{E}|^2} \left(1 + \frac{m\sigma^2}{5|\tilde{E}|}\right)^{-1} \left[1 + \exp\left(\frac{|\tilde{E}|}{m\sigma^2}\right)\right]^{-1}, \quad (7.174)$$

where as usual we assume that all masses are equal and the velocity distributions of the binary and field stars are Maxwellian and in equipartition. This reduces to the ionization rate for soft binaries, equation (7.171), in the limit $|\tilde{E}| \ll m\sigma^2$. Equation (7.174) was obtained by Hut & Bahcall (1983) from numerical experiments; the functional form is an interpolation based on formulae derived by Heggie (1975) for very hard and very soft binaries. In principle, the ionization rate depends on the initial eccentricity of the binary, but the numerical experiments show that the formula is accurate to within $\pm 20\%$ at all eccentricities. Equation (7.174) shows that the lifetime of a hard binary against disruption becomes exponentially long as the binary becomes very hard, and thus that primordial hard binaries, unlike soft ones, can survive from the birth of a cluster to the present day.

The formation rate per unit volume is obtained from the principle of detailed balance (7.164),

$$C(\tilde{E}) = \frac{2\pi^2G^5m^{11/2}\rho^3}{3^{3/2}\sigma^2|\tilde{E}|^{9/2}} \left(1 + \frac{m\sigma^2}{5|\tilde{E}|}\right)^{-1} \left[1 + \exp\left(-\frac{|\tilde{E}|}{m\sigma^2}\right)\right]^{-1}. \quad (7.175)$$

The total formation rate of hard binaries is

$$C_h = \int_{m\sigma^2}^{\infty} d|\tilde{E}| C(\tilde{E}) = 0.74 \frac{G^5\rho^3m^2}{\sigma^9}, \quad (7.176)$$

consistent with the rough estimate n/t_3 obtained from equation (7.11).

The formation rate of soft binaries with binding energy exceeding $|\tilde{E}|$ is

$$\int_{|\tilde{E}|}^{m\sigma^2} d|\tilde{E}'| C(\tilde{E}') \rightarrow 3.8 \frac{G^5\rho^3m^2}{\sigma^9} \left(\frac{m\sigma^2}{|\tilde{E}|}\right)^{5/2} \quad \text{as } |\tilde{E}| \rightarrow 0. \quad (7.177)$$

This creation rate diverges as $|\tilde{E}| \rightarrow 0$, but such wide binaries remain bound for only a fraction of a period before being disrupted by passing stars or tidal

forces: they can be regarded as pairs of unassociated particles that simply happen to have nearly the same velocity and position.

The rate at which hard binaries become harder is (Heggie 1975; Heggie & Hut 1993)

$$\langle \dot{\tilde{E}} \rangle = -3.8 \frac{G^2 \rho m^2}{\sigma}. \quad (7.178)$$

Using equation (7.163), this result can be expressed in terms of the binary semi-major axis,

$$\left\langle \frac{d}{dt} \frac{1}{a} \right\rangle = 7.6 \frac{G\rho}{\sigma}. \quad (7.179)$$

Notice that the average rate of change of energy is independent of the energy; in other words, the rate of hardening is independent of the hardness.

The following heuristic derivation of equation (7.178) illuminates the physics of binary hardening. Consider a field star that is traveling towards a hard binary with initial velocity V_0 and impact parameter b , both taken relative to the center of mass of the binary. Its initial angular momentum is $L = bV_0$ (eq. 3.46). The field star will interact strongly with the binary only if it approaches to within a distance $\sim a$ of the binary's center of mass. At this point its velocity is similar to the orbital speed $v \simeq (Gm/a)^{1/2}$ of the binary, since it has fallen into the same potential well that governs the orbital motion of the binary. Its angular momentum is of order av and this cannot be very different from the initial angular momentum L ; thus only stars with impact parameter $b \lesssim b_0 \equiv av/V_0$ can interact strongly with the binary. The number of strongly interacting field stars per unit time is thus $F \sim nV_0 b_0^2 \sim na^2 v^2 / \sigma$, where n is the number density of stars surrounding the binary and we have replaced V_0 by the velocity dispersion σ . Each of these strongly interacting stars will follow a complicated trajectory like the one in Figure 7.8, but ultimately will be ejected from the binary. The ejection velocity is expected to be of order v , since this is the relative velocity of the field star and the binary components when they are interacting strongly. Thus, on average, each field star robs the binary of energy $\sim mv^2$. The rate of energy loss is thus $Fmv^2 \sim \rho a^2 v^4 / \sigma \sim G^2 \rho m^2 / \sigma$, where $\rho = nm$ is the mass density of field stars near the binary. This is an approximate statement of equation (7.178).

The rate of energy loss (7.178) can be simply expressed in terms of the local relaxation time (eq. 7.106): taking $\ln \Lambda \simeq 10$, the average energy change per relaxation time is

$$\langle \Delta \tilde{E} \rangle \simeq -0.13 m \sigma^2. \quad (7.180)$$

As the binary hardens and its semi-major axis shrinks, close encounters with field stars become rarer and more violent, even though the average hardening rate remains constant. The recoil velocity given to the binary in a close encounter also becomes larger and larger as it hardens, and eventually the process terminates when the recoil from some encounter is so large that

the binary is ejected from the cluster. For typical cluster potentials, a hard binary survives until $\tilde{E} \approx -km\sigma^2$ where $k \approx 100$ (Goodman & Hut 1993; Heggie & Hut 2003). Thus the lifetime of a typical hard binary is

$$t_{\text{hb}} \approx k \frac{m\sigma^2}{-\langle \tilde{E} \rangle} = \frac{k}{3.8} \frac{\sigma^3}{G^2 \rho m} \simeq \frac{k}{0.13} t_{\text{relax}}. \quad (7.181)$$

After a few relaxation times, the distribution of soft and moderately hard binaries approaches a steady state. The net energy released by interactions between binaries and single stars then arises solely from the flux of hard binaries towards larger and larger values of $|\tilde{E}|$, which continues until the recoil of the binary from an encounter ejects it from the cluster. Thus the rate of energy generation per unit mass due to encounters between binary and single stars is just

$$\epsilon = -\langle \tilde{E} \rangle \frac{n_{\text{b}}}{\rho} = 3.8 \frac{G^2 m^2 n_{\text{b}}}{\sigma}, \quad (7.182)$$

where n_{b} is the number density of hard binaries. This number density is determined by the competition between the formation of moderately hard binaries by three-body encounters (eq. 7.176) and the ejection of very hard binaries from the cluster; thus $n_{\text{b}} \simeq C_{\text{h}} t_{\text{hb}}$ where C_{h} is the hard-binary formation rate of equation (7.176). We have

$$n_{\text{b}} \simeq 0.20k \frac{G^3 \rho^2 m}{\sigma^6}. \quad (7.183)$$

Substituting this result into equation (7.182), we obtain the equilibrium energy generation rate per unit mass due to binaries,

$$\epsilon \simeq C_{\text{b}} \frac{G^5 \rho^2 m^3}{\sigma^7} \quad (7.184)$$

with $C_{\text{b}} = 0.74k$, or $C_{\text{b}} \approx 75$ for the value $k \approx 100$ estimated above. Numerical experiments yield a similar value, $C_{\text{b}} \simeq 80$ (Goodman & Hut 1993; Heggie & Hut 2003).

The initial evolution of globular clusters is strongly affected by the number of primordial binaries in the cluster. Spectroscopic observations suggest that the fractional abundance of primordial hard binaries in globular clusters is 10–20% (Hut et al. 1992; Rubenstein & Bailyn 1997; Albrow et al. 2001). The heating of the cluster by primordial hard binaries delays core collapse, just as the burning of nuclear fuel halts the collapse of stars. Eventually, as we have described, each binary becomes so hard that an encounter ejects it from the cluster. Once most of the hard binaries have fled, core collapse proceeds just as in a cluster without primordial binaries.

The results described here are restricted to single-mass clusters; more comprehensive results for unequal masses are given in Sigurdsson & Phinney (1993) and Valtonen & Karttunen (2006).

Encounters of binaries with single stars is not the only relevant dynamical process involving binaries: during the late stages of collapse of a binary-rich core, binary-binary encounters are likely to predominate over binary-single encounters (Valtonen & Mikkola 1991).

7.5.8 Inelastic encounters

So far we have approximated stars as point masses. However, in the densest stellar systems, such as the center of a globular cluster in the late stages of core collapse or the cusp of stars surrounding the central black hole in a galaxy (§7.5.9), stars may pass so close to one another that they collide. Even a near-miss can have important consequences, because the relative kinetic energy of the passing stars is dissipated by the tides the stars raise in each other. The loss of this energy hastens core collapse, and the close encounters or collisions can form unusual objects such as blue stragglers (see page 628), or close binary stars that are visible as cataclysmic variables or X-ray sources (Pooley & Hut 2006).

To investigate these effects we compute the **collision time** t_{coll} , where $1/t_{\text{coll}}$ is the collision rate, that is, the average number of physical collisions that a given star suffers per unit time. For simplicity we again restrict ourselves to a cluster in which all stars have the same mass m . Consider an encounter with initial relative velocity \mathbf{V}_0 and impact parameter b . The angular momentum per unit mass of the reduced particle is $L = bV_0$ (eq. 3.46). At the distance of closest approach, which we denote by r_{coll} , the radial velocity is zero, so the angular momentum is $L = r_{\text{coll}}V_{\text{max}}$, where V_{max} is the relative speed at r_{coll} . From the energy equation (7.162), the energy in the center-of-mass frame is $\tilde{E} = \frac{1}{2}\mu V^2 - Gm^2/r$, where $\mu = \frac{1}{2}m$ is the reduced mass (eq. D.32). Equating the energy at r_{coll} and $r \rightarrow \infty$, we have

$$\frac{1}{4}mV_0^2 = \frac{1}{4}mV_{\text{max}}^2 - \frac{Gm^2}{r_{\text{coll}}}. \quad (7.185)$$

Since L is conserved, we can eliminate V_{max} to obtain

$$b^2 = r_{\text{coll}}^2 + \frac{4Gmr_{\text{coll}}}{V_0^2}. \quad (7.186)$$

If we set r_{coll} equal to the sum of the radii of the two stars, then a collision will occur if and only if the impact parameter is less than the value of b determined by equation (7.186).

Let $f(\mathbf{v}_a)d^3\mathbf{v}_a$ be the number of stars per unit volume with velocities in the range \mathbf{v}_a to $\mathbf{v}_a + d^3\mathbf{v}_a$. The number of encounters with these stars per

unit time with impact parameter less than b that are suffered by a given star is just $f(\mathbf{v}_a) d^3\mathbf{v}_a$ times the volume of an annulus with radius b and length $V_0 = |\mathbf{v} - \mathbf{v}_a|$, where \mathbf{v} is the velocity of the subject star. The collision rate is then

$$R_{\text{coll}} = \int d^3\mathbf{v}_a f(\mathbf{v}_a) \pi b^2 |\mathbf{v} - \mathbf{v}_a|. \quad (7.187)$$

We now average this rate over \mathbf{v} by multiplying (7.187) by $f(\mathbf{v})/n$, where $n = \int d^3\mathbf{v} f(\mathbf{v})$ is the number density of stars, and integrating over $d^3\mathbf{v}$. This gives an estimate of the collision time,

$$\frac{1}{t_{\text{coll}}} \equiv \langle R_{\text{coll}} \rangle = \frac{\pi}{n} \int d^3\mathbf{v} d^3\mathbf{v}_a f(\mathbf{v}) f(\mathbf{v}_a) b^2 |\mathbf{v} - \mathbf{v}_a|. \quad (7.188)$$

To evaluate the integrals we assume that the DF is Maxwellian with dispersion σ (eq. 7.91). Substituting for b^2 from equation (7.186) we obtain

$$\frac{1}{t_{\text{coll}}} = \frac{n}{8\pi^2\sigma^6} \int d^3\mathbf{v} d^3\mathbf{v}_a e^{-(v^2+v_a^2)/2\sigma^2} \left(r_{\text{coll}}^2 |\mathbf{v} - \mathbf{v}_a| + \frac{4Gmr_{\text{coll}}}{|\mathbf{v} - \mathbf{v}_a|} \right). \quad (7.189)$$

We now replace the dummy variables \mathbf{v} and \mathbf{v}_a by the relative velocity $\mathbf{V} = \mathbf{v} - \mathbf{v}_a$ and the center-of-mass velocity $\mathbf{v}_{\text{cm}} = \frac{1}{2}(\mathbf{v} + \mathbf{v}_a)$. We have

$$\mathbf{v} = \mathbf{v}_{\text{cm}} + \frac{1}{2}\mathbf{V} \quad ; \quad \mathbf{v}_a = \mathbf{v}_{\text{cm}} - \frac{1}{2}\mathbf{V}. \quad (7.190)$$

The Jacobian determinant of the map between $(\mathbf{v}, \mathbf{v}_a)$ and $(\mathbf{V}, \mathbf{v}_{\text{cm}})$ satisfies $\partial(\mathbf{v}, \mathbf{v}_a)/\partial(\mathbf{v}_{\text{cm}}, \mathbf{V}) = 1$, and hence the velocity-space volume element $d^3\mathbf{v} d^3\mathbf{v}_a$ can be replaced by $d^3\mathbf{v}_{\text{cm}} d^3\mathbf{V}$. In addition $v^2 + v_a^2 = 2v_{\text{cm}}^2 + \frac{1}{2}V^2$. We find

$$\frac{1}{t_{\text{coll}}} = \frac{n}{8\pi^2\sigma^6} \int d^3\mathbf{v}_{\text{cm}} d^3\mathbf{V} e^{-(v_{\text{cm}}^2+V^2/4)/\sigma^2} \left(r_{\text{coll}}^2 V + \frac{4Gmr_{\text{coll}}}{V} \right). \quad (7.191)$$

The integral over \mathbf{v}_{cm} is

$$\int d^3\mathbf{v}_{\text{cm}} e^{-v_{\text{cm}}^2/\sigma^2} = \pi^{3/2}\sigma^3. \quad (7.192)$$

Substituting this result in equation (7.191) and integrating over all directions of \mathbf{V} , we have

$$\frac{1}{t_{\text{coll}}} = \frac{\pi^{1/2}n}{2\sigma^3} \int_0^\infty dV e^{-V^2/4\sigma^2} (V^3 r_{\text{coll}}^2 + 4GmV r_{\text{coll}}). \quad (7.193)$$

The remaining integrals are easy, and yield

$$\frac{1}{t_{\text{coll}}} = 4\sqrt{\pi}n\sigma \left(r_{\text{coll}}^2 + \frac{Gm}{\sigma^2} r_{\text{coll}} \right). \quad (7.194)$$

The second term represents the enhancement in the collision rate due to **gravitational focusing**, that is, the deflection of trajectories by the gravitational attraction of the two stars.

If r_* is the stellar radius, we may set $r_{\text{coll}} = 2r_*$. It is convenient to introduce the escape speed from the stellar surface, $v_* = \sqrt{2Gm/r_*}$, and to write equation (7.194) as

$$\frac{1}{t_{\text{coll}}} = 16\sqrt{\pi}n\sigma r_*^2 \left(1 + \frac{v_*^2}{4\sigma^2}\right) = 16\sqrt{\pi}n\sigma r_*^2(1 + \Theta), \quad (7.195a)$$

where we have introduced the **Safronov number**¹⁶

$$\Theta \equiv \frac{v_*^2}{4\sigma^2} = \frac{Gm}{2\sigma^2 r_*} = 9.54 \frac{m}{\mathcal{M}_\odot} \frac{R_\odot}{r_*} \left(\frac{100 \text{ km s}^{-1}}{\sigma}\right)^2. \quad (7.195b)$$

Numerically, we have

$$t_{\text{coll}} = \begin{cases} 6.8 \times 10^3 \text{ Gyr} \frac{10^5 \text{ pc}^{-3}}{n} \frac{100 \text{ km s}^{-1}}{\sigma} \left(\frac{R_\odot}{r_*}\right)^2 & \Theta \ll 1; \\ 7.1 \times 10^2 \text{ Gyr} \frac{10^5 \text{ pc}^{-3}}{n} \frac{\sigma}{100 \text{ km s}^{-1}} \frac{R_\odot}{r_*} \frac{\mathcal{M}_\odot}{m} & \Theta \gg 1. \end{cases} \quad (7.196)$$

These formulae underestimate the collision rate in clusters with a substantial population of hard binary stars, since the typical close encounter of a single star with a hard binary (Figure 7.8) forms a temporary triple system, which offers many possible opportunities for a collision.

The ratio of the collision time to the relaxation time (eq. 7.106) has the simple form

$$\frac{t_{\text{coll}}}{t_{\text{relax}}} = 0.4 \ln \Lambda \frac{\Theta^2}{1 + \Theta}, \quad (7.197)$$

in approximate agreement with our crude estimate (7.9) for $\Theta \ll 1$. Core collapse occurs in about 300 central relaxation times (eq. 7.154), and hence collisions can have a substantial influence on core collapse if $t_{\text{coll}}/t_{\text{relax}} \lesssim 300$. Let us now evaluate this ratio in some typical cases. The escape speed from the Sun is $v_* = 618 \text{ km s}^{-1}$. In a typical globular cluster, where $\sigma \simeq 10 \text{ km s}^{-1}$, the Safronov number for solar-type stars is $\Theta \simeq 1 \times 10^3$. Hence $t_{\text{coll}} \simeq 4 \times 10^3 t_{\text{relax}}$ for $\ln \Lambda \approx 10$. Thus, in the early stages of core collapse, the effects of collisions of solar-type stars, or indeed of any type of stars, are negligible compared to the effects of relaxation.¹⁷ However, collisions *can*

¹⁶ After V. S. Safronov, who introduced Θ in studies of the collision of planetesimals in the early solar system.

¹⁷ Red giants have much larger radii (up to $\approx 30 R_\odot$) but are less abundant by about a factor of 100. Since the Safronov number is large in a globular cluster, the collision rate is proportional to nr_{coll} , so a main-sequence star collides with giants even less often than it collides with other main-sequence stars.

be important in the late stages of core collapse, since the velocity dispersion grows: equation (7.156) predicts that the dispersion grows by a factor of 4 or so as the core mass falls by a factor 10^4 , so the Safronov number may become as small as $\Theta \approx 60$, in which case $t_{\text{coll}} \simeq 20t_{\text{relax}}$ for $\ln \Lambda \approx 1$. Collisions also play an important role in the central parsec of a galaxy, where $\sigma \approx 100\text{--}200 \text{ km s}^{-1}$, so the Safronov number is $\Theta \approx 2.4\text{--}10$, and $t_{\text{coll}}/t_{\text{relax}} \approx 10\text{--}50$ for $\ln \Lambda \approx 14$.

A head-on collision of two main-sequence stars in a globular cluster produces a luminous, hot remnant ($L \approx 50 L_{\odot}$, $T \approx 8000 \text{ K}$), which cools and contracts to the main sequence in a few million years (Sills et al. 2002). The product of an off-axis collision is less clear: the remnant has a large spin, and as it contracts it spins up by conservation of angular momentum, until eventually it reaches break-up speed. The fate of off-axis merger remnants depends strongly on whether mechanisms such as a magnetic wind or circumstellar disk allow the remnant to shed angular momentum.

The remnants of stellar collisions may be detectable as **blue stragglers**, stars that are located in the globular-cluster color-magnitude diagram on the main sequence but beyond the turnoff point (Figure 1.2; BM §6.1; Bailyn 1995). This location implies that blue stragglers have lifetimes shorter than the cluster age. Collisions occurring long after the cluster formed provide a natural way to create such stars, but there are alternative mechanisms such as mass transfer between or the coalescence of the two components of a close binary star, and the importance of collisions in forming blue stragglers remains unclear.

Some process unique to high-density stellar systems is also likely to be responsible for the formation of low-mass X-ray binaries in globular clusters, since these binaries are $\sim 10^3$ times more abundant in the densest globular clusters than in the Galaxy as a whole (Pooley & Hut 2006). The most promising candidate process is collisions between neutron stars and red-giant stars (Ivanova et al. 2005). Similarly, millisecond pulsars are $\sim 10^2$ times more abundant in dense clusters. These are old, quiescent neutron stars that have been spun up by an accretion disk, enabling them to shine again as **recycled pulsars**. The source of the disk material could be residual material left over after an off-center inelastic encounter, or a nearby companion star (Bailyn 1995; Phinney 1996).

To close this discussion, we briefly mention the effects of near-collisions. A close encounter of two stars raises violent tides on the surface of each star. The energy that powers these disturbances comes from the relative kinetic energy of the stars, and the loss of this energy in a close encounter can leave the two stars with negative orbital energy—that is, as a bound binary system. This effect is especially important in systems such as globular clusters where the Safronov number $\Theta \gg 1$, since to form a binary the stars only need to lose their orbital energy ($\sim m\sigma^2$), which is small compared to their internal binding energies ($\sim mv_{\star}^2$). When tidal dissipation just manages to bind two stars, the resulting binary is highly eccentric. The evolution of

such binaries is still not well understood (Kochanek 1992; Mardling 1995; Mardling & Aarseth 2001). It is unclear whether the tidal disturbances excited in the first near-collision will damp by the time of the next pericenter passage (Kumar & Goodman 1996): if so, then the binary orbit will steadily decay and circularize; if not, the tidal forces on subsequent passages can either add or remove energy from the orbit, depending on the phase of the tidal oscillations at the time of each periastron passage. In the latter case, the binary orbital energy will random walk, and the binary may even be disrupted if the random walk achieves positive orbital energies. An additional complication is that the energy dissipated in the stars by the tides will heat and expand their outer parts, perhaps causing them to lose mass, merge, or even be disrupted.

If the stars survive, and remain bound but do not merge, the binary will eventually settle into a circular orbit as more and more energy is dissipated. Since the initial orbital angular momentum was small, and total angular momentum is conserved, the size of the resulting binary orbit will be only a few stellar radii (Problem 7.16). Such binaries are sometimes called **tidal-capture binaries**. The fraction of near-collisions that achieve this final state remains quite uncertain.

7.5.9 Stellar systems with a central black hole

The power emitted by quasars and other active galactic nuclei is believed to arise from the consumption of gas and stars by a black hole of mass $M_{\bullet} \approx 10^6\text{--}10^9 M_{\odot}$ (§1.1.6 and Krolik 1999). These black holes should normally be located at the center of the galaxy, since the galactic orbit of any object of such large mass decays rapidly due to dynamical friction (§8.1.1a). Observations of the centers of nearby galaxies show that most contain massive dark objects which are almost certainly black holes (page 605); these are likely to have been active galactic nuclei early in the lifetime of the galaxy, but are now dark because they are starved of fuel (Kormendy 2004).

Smaller black holes ($M_{\bullet} \approx 10^2\text{--}10^4 M_{\odot}$) might also be present at the centers of some globular clusters.

(a) Consumption of stars by the black hole One obvious effect of a black hole on the surrounding system is that it swallows any star whose orbit carries it within the black-hole horizon. In the Schwarzschild metric, the radial coordinate of the horizon of a non-rotating black hole is

$$r_{\bullet} = \frac{2GM_{\bullet}}{c^2} = 2.95 \text{ km} \frac{M_{\bullet}}{M_{\odot}} = 9.57 \times 10^{-8} \text{ pc} \frac{M_{\bullet}}{10^6 M_{\odot}}, \quad (7.198)$$

where c is the speed of light. Now consider a test particle that approaches the black hole from infinity on an orbit with angular momentum L and zero energy. It can be shown that the orbit crosses the horizon if $L < L_{\bullet} \equiv$

$4GM_{\bullet}/c$ (Novikov & Frolov 1989). It proves convenient to express this result in terms of the pericenter q_K of a parabolic Keplerian orbit with the same angular momentum. Since $L = (2GM_{\bullet}q_K)^{1/2}$ for Keplerian motion, the star crosses the horizon if $q_K < q_{\bullet} \equiv 4r_{\bullet} = 8GM_{\bullet}/c^2$.

Tidal forces from the black hole can also destroy stars, often at distances much larger than the horizon. Let us consider a star on a parabolic orbit around the black hole, with pericenter distance q_K that is much larger than q_{\bullet} , so the orbit is Keplerian. At the instant of pericenter passage, the force per unit mass on a fluid element on the surface of the star facing the black hole is $\mathbf{F}_{\bullet,i} = -GM_{\bullet}\hat{\mathbf{e}}_r/(q_K - r_{\star})^2$, where $\hat{\mathbf{e}}_r$ is the unit vector pointing away from the black hole and r_{\star} is the radius of the star. Similarly, the force on a surface element facing away from the black hole is $\mathbf{F}_{\bullet,o} = -GM_{\bullet}\hat{\mathbf{e}}_r/(q_K + r_{\star})^2$. The difference between these forces is

$$\Delta\mathbf{F}_{\bullet} = \mathbf{F}_{\bullet,i} - \mathbf{F}_{\bullet,o} = -GM_{\bullet}\hat{\mathbf{e}}_r \left(\frac{1}{(q_K - r_{\star})^2} - \frac{1}{(q_K + r_{\star})^2} \right) \simeq -\frac{4GM_{\bullet}r_{\star}}{q_K^3}\hat{\mathbf{e}}_r, \quad (7.199a)$$

where we have assumed that $r_{\star} \ll q_K$. The force differential $\Delta\mathbf{F}_{\bullet}$ tends to pull the star apart, but this tendency is countered by the gravitational force from the star itself, $\mathbf{F}_{\star,i} = Gm\hat{\mathbf{e}}/r_{\star}^2$ on the inner surface and $\mathbf{F}_{\star,o} = -\mathbf{F}_{\star,i}$ on the outer surface. Thus

$$\Delta\mathbf{F}_{\star} = \mathbf{F}_{\star,i} - \mathbf{F}_{\star,o} = \frac{2Gm}{r_{\star}^2}\hat{\mathbf{e}}_r. \quad (7.199b)$$

The star is likely to be disrupted at pericenter if the magnitude of the force differential in (7.199a) exceeds (7.199b), that is, if

$$q_K < q_{\text{dis}} \equiv gr_{\star} \left(\frac{M_{\bullet}}{m} \right)^{1/3}, \quad (7.200)$$

where g is of order unity. Numerical calculations of the disruption of stars on parabolic orbits confirm that equation (7.200) yields the correct scaling of the disruption radius with the properties of the black hole and the star, and show that g varies between 0.9 and 1.7 for a plausible range of stellar density distributions (Lai, Rasio, & Shapiro 1994; Diener et al. 1995).

This Newtonian analysis can be extended to parabolic orbits with pericenter distances comparable to the horizon, using numerical models for relativistic tidal disruption. The results can still be expressed using equation (7.200), but now q_K must be interpreted as a function of the orbital angular momentum $L = (2GM_{\bullet}q_K)^{1/2}$ rather than as the pericenter distance, and the dimensionless coefficient g depends on both the stellar model and the black-hole mass. If we define ‘‘tidal disruption’’ to mean the loss of half of the stellar mass, then for stars of solar mass and radius, modeled as $n = 1.5$ polytropes, Ivanov & Chernyakova (2006) find that g varies from 1.1 to 2.4 as M_{\bullet} varies from $10^6 M_{\odot}$ to $4 \times 10^7 M_{\odot}$.

If $q_\bullet \gtrsim q_{\text{dis}}$ stars are swallowed whole by the black hole, while if $q_\bullet \lesssim q_{\text{dis}}$ stars are disrupted by the tidal field of the hole before reaching the horizon. Thus, stars are swallowed whole from nearly parabolic orbits if

$$M_\bullet \gtrsim 1.4 \times 10^7 M_\odot g^{3/2} \left(\frac{r_\star}{R_\odot} \right)^{3/2} \left(\frac{M_\odot}{m} \right)^{1/2}. \quad (7.201)$$

We shall say that a star has been “eaten” by the black hole if it is either tidally disrupted or swallowed whole, and thus we define q_{eat} to be the larger of q_{dis} and q_\bullet .

In a spherical system, orbits having pericenter distance $q_K < q_{\text{eat}}$ occupy a specific region in energy-angular momentum space, which is called the **loss cone**.¹⁸ For highly eccentric orbits the loss cone consists of the region with angular momentum $L < L_{\text{lc}} \equiv (2GM_\bullet q_{\text{eat}})^{1/2}$.

If the gravitational potential from the black hole and the surrounding stars were precisely spherical, then the loss cone at a given energy would be emptied in one orbital period and no further stars would be eaten (except for dwarf stars just outside the loss cone that begin ascending the giant branch, thereby expanding in radius so the loss cone grows to engulf them). In realistic systems, however, there is a steady supply of fresh stars into the loss cone from two dynamical mechanisms:

- (i) Through encounters with other stars, the star may diffuse in energy and angular momentum until it enters the loss cone, at which point it is disrupted or swallowed in less than an orbital time (Frank & Rees 1976; Lightman & Shapiro 1977).
- (ii) If the galaxy is non-spherical, torques from its overall mass distribution can carry stars into the loss cone (Magorrian & Tremaine 1999).

As a result of these two mechanisms, a star will be disrupted every 10^4 – 10^5 yr in a typical galaxy (Magorrian & Tremaine 1999). The gas released by tidally disrupted stars may re-accrete onto the black hole over a timescale of months to years, producing a characteristic X-ray flare that signals the presence of a black hole (Komossa et al. 2004).

(b) The effect of a central black hole on the surrounding stellar system Let us imagine that a black hole of mass M_\bullet sits at the center of a spherical stellar system. At large radii, the gravitational field is dominated by the stellar mass distribution and is approximately equal to $-GM_\star(r)/r^2$, where $M_\star(r)$ is the mass of stars inside radius r . At small radii, the gravitational field is dominated by the contribution of the black hole, $-GM_\bullet/r^2$. The transition between these two limits, where the central force from the

¹⁸ The term “cone” reflects the geometry in velocity space, not action space or energy-angular momentum space: at a given position, the stars with pericenter distance $< q_{\text{eat}}$ have velocity vectors that fill a cone in velocity space, with the symmetry axis of the cone pointing towards the black hole.

black hole equals the force from the stars, occurs at the dynamical radius r_g of the black hole (page 353).

The formula for the dynamical radius depends on the density distribution in the galaxy. In a spherical galaxy with a constant-density central core, such as a non-singular isothermal sphere or a King model (§4.3.3b,c), the field due to the stars is $-\frac{4}{3}\pi G\rho_0 r$ where ρ_0 is the central density. Comparing this to the field from the black hole, we have

$$r_g = \left(\frac{3M_\bullet}{4\pi\rho_0}\right)^{1/3} = \left(\frac{GM_\bullet r_0^2}{3\sigma^2}\right)^{1/3}. \quad (7.202)$$

Here r_0 is the King radius (4.106).

Inside the dynamical radius the black hole dominates the dynamics, while outside the self-gravity of the stars determines the dynamics. The dynamical radius for Hernquist models with central black holes is shown in Figure 4.20.

The effect of a central black hole on its host stellar system depends on how it forms. If the black hole grows slowly—perhaps by accretion of gas over Myr or longer timescales—it will compress the surrounding stellar orbits to form a cusp, as described in §4.6.1a. On the other hand, if the hole grows by merging with other black holes that spiral to the center by dynamical friction, as described in §8.1.1a, the energy transferred to the stars from the inspiraling black hole causes the stellar orbits to expand, thereby reducing the central density of the host system (**core scouring**). Without knowing the history of the black hole, we generally cannot predict its influence on the surrounding stellar system.

However, there is one case in which the properties of the stellar system surrounding a black hole *can* be predicted without knowledge of how the system formed: if its age is much larger than the central relaxation time, a steady-state distribution of stars is established in which there is a balance between the tendency of the system to come into thermal equilibrium, with $f \propto \exp(-H/\sigma^2)$, and the continuous depopulation of orbits with pericenters smaller than q_{eat} .

This steady-state density distribution can be determined by a simple but rather subtle argument (see Problem 7.19, Bahcall & Wolf 1976, and Shapiro & Lightman 1976). Let us assume that the density near the hole is a power law, $n(r) \propto r^{-s}$. The Jeans equations dictate that the mean-square velocity at any radius should be of order $\langle v^2 \rangle \simeq GM_\bullet/r$. Hence the local relaxation time (7.106) is $t_{\text{relax}} \approx \langle v^2 \rangle^{3/2}/(G^2 m^2 n) \propto r^{s-3/2}$ (neglecting the Coulomb logarithm). A star that is eaten by the hole from radius r_\bullet has energy of order $E(r) = -GM_\bullet m/r_\bullet$, and in a steady state this loss of negative energy implies a flow of positive energy out through the cusp. Relaxation among the $N(r)$ cusp stars interior to r can carry energy of order $N(r)E(r)$ through the shell at radius r per relaxation time; since $N(r) \propto r^{3-s}$ the flow of energy through radius r is $N(r)E(r)/t_{\text{relax}} \propto r^{-2s+7/2}$. In a steady state the

energy flow must be independent of radius and hence $s = \frac{7}{4}$. This scaling law is verified by numerical solutions of the Fokker–Planck equation (Cohn & Kulsrud 1978).

Fokker–Planck, N-body, and orbit-averaged Monte Carlo codes that follow the evolution of dense stellar systems containing a black hole are described by Cohn & Kulsrud (1978), Freitag & Benz (2001), and Baumgardt, Makino, & Ebisuzaki (2004).

7.6 Summary

Several distinct processes drive the evolution of an isolated self-gravitating system: (i) escape of stars to infinity; (ii) core collapse; (iii) equipartition; (iv) formation of hard binaries; (v) gravitational interactions of hard binaries with single stars and one another. Both escape and core collapse occur on a timescale proportional to and somewhat larger than the relaxation time t_{rh} , but core collapse is generally the faster process. In single-mass systems, the timescale for core collapse is $16t_{\text{rh}}$.

Hard binaries—both primordial binaries and binaries formed in the late stages of core collapse—act as energy sources, much like the nuclear reactions occurring in the center of a star, and halt core collapse. Remarkably, it appears that even a single hard binary is able to halt the core collapse of a cluster of 10^6 stars.

After core collapse, the inner parts of the cluster re-expand, and may experience gravothermal oscillations in which the central density varies by several orders of magnitude. It is likely that many globular clusters have undergone core collapse, but as yet there is no clear observational signature that distinguishes post-collapse from pre-collapse clusters.

Isolated post-collapse clusters expand and evolve more and more slowly. In contrast, if the cluster is truncated by a tidal field the post-collapse expansion spills mass over the tidal radius, leading to rapid dissolution.

Stellar collisions are relatively rare in globular clusters because the escape speed from the surface of the stars is much higher than the velocity dispersion in the cluster. So although collisions and near-collisions can produce exotic stellar species (blue stragglers, cataclysmic variables, X-ray binaries, recycled pulsars, etc.), they are unlikely to influence the evolution of the cluster as a whole.

Problems

7.1 [1] Consider a system in which the interparticle potential energy has the form $\Phi_{\alpha\beta} = -C|\mathbf{x}_\alpha - \mathbf{x}_\beta|^{-p}$, where p and C are positive constants.

(a) Show that the scalar virial theorem has the form

$$2K + pW = 0, \quad (7.203)$$

where K is the kinetic energy and W is the potential energy.

(b) For what values of p does the system have negative heat capacity, in the sense of equation (7.51)?

7.2 [1] Prove that a system of N self-gravitating point masses with positive energy must disrupt, in the sense that at least one star must escape. Hint: use the virial theorem, and prove that the moment of inertia must increase without limit.

7.3 [2] A simple model for a galactic disk consists of N infinite, parallel sheets, each having surface density σ . Let z be the coordinate perpendicular to the sheets and label the position of the i th sheet by z_i .

(a) Show that the equation of motion is

$$\ddot{z}_i = 2\pi G\sigma(N_{+i} - N_{-i}), \quad (7.204)$$

where N_{+i} is the number of sheets with $z > z_i$ and N_{-i} is the number with $z < z_i$.

(b) Show that the Hamiltonian of the system can be written as

$$H = \frac{1}{2\sigma} \sum_{i=1}^N p_i^2 + \pi G\sigma^2 \sum_{\substack{i,j=1 \\ i \neq j}}^N |z_i - z_j|, \quad (7.205)$$

where p_i is the momentum conjugate to z_i .

(c) If the first and second terms in the Hamiltonian are identified with the kinetic energy K and potential energy W per unit area, show that the virial theorem has the form

$$2K = W \quad \text{or} \quad E \equiv K + W = 3K. \quad (7.206)$$

7.4 [1] Consider a stellar system composed of two types of stars, with density distributions $\rho_1(\mathbf{x})$ and $\rho_2(\mathbf{x})$ and corresponding potentials $\Phi_1(\mathbf{x})$ and $\Phi_2(\mathbf{x})$. Show that in a steady state, the scalar virial theorem for component 2 may be written in the form

$$2K_2 + W_2 - \int d^3\mathbf{x} \rho_2(\mathbf{x}) \mathbf{x} \cdot \nabla \Phi_1(\mathbf{x}) = 0, \quad (7.207)$$

where K_2 is the kinetic energy of component 2, and W_2 is the potential energy due to the mutual interaction of the stars of component 2. Hint: use Problem 4.38.

7.5 [2] The **Klimontovich distribution function** for an N -body system of identical point masses m is

$$f(\mathbf{x}, \mathbf{v}, t) = \sum_{\alpha=1}^N \delta[\mathbf{x} - \mathbf{x}_\alpha(t)] \delta[\mathbf{v} - \mathbf{v}_\alpha(t)], \quad (7.208)$$

where δ denotes the three-dimensional delta function (Appendix C.1). The functions $\mathbf{x}_\alpha(t)$, $\mathbf{v}_\alpha(t)$ describe the trajectories of the N bodies, which satisfy the equations

$$\dot{\mathbf{x}}_\alpha = \mathbf{v}_\alpha \quad ; \quad \dot{\mathbf{v}}_\alpha = Gm \sum_{\beta \neq \alpha} \frac{\mathbf{x}_\beta - \mathbf{x}_\alpha}{|\mathbf{x}_\alpha - \mathbf{x}_\beta|^3}, \quad (7.209)$$

(a) Prove that the Klimontovich DF is an exact solution of the collisionless Boltzmann equation (4.11) for the one-body Hamiltonian

$$H(\mathbf{x}, \mathbf{v}, t) = \frac{1}{2}v^2 - \sum_{\alpha=1}^N \frac{Gm}{|\mathbf{x}_\alpha(t) - \mathbf{x}|}. \quad (7.210)$$

The Klimontovich DF offers an exact formal description of the N-body stellar system, and provides an alternative to the BBGKY hierarchy for the systematic derivation of kinetic equations to describe stellar systems (Nishikawa & Wakatani 2000).

(b) We have argued that the collisionless Boltzmann equation cannot account for relaxation due to encounters between individual stars. How is this consistent with our conclusion that the Klimontovich DF is an exact solution of the N-body equations of motion?

7.6 [2] The object of this problem is to determine the behavior of the curve in Figure 7.1 in the limit as the central concentration $\mathcal{R} \rightarrow \infty$.

(a) The density of an isothermal sphere satisfies equation (4.107a),

$$\frac{d}{d\tilde{r}} \left(\tilde{r}^2 \frac{d \ln \tilde{\rho}}{d\tilde{r}} \right) = -9\tilde{r}^2 \tilde{\rho}. \quad (7.211)$$

As the dimensionless radius $\tilde{r} \rightarrow \infty$, the solutions of equation (7.211) approach the singular isothermal sphere $\tilde{\rho}_S(\tilde{r}) = 2/(9\tilde{r}^2)$. To determine the asymptotic behavior more accurately, define new variables u and $z(u)$ by $u = 1/\tilde{r}$ and $\tilde{\rho} \equiv \tilde{\rho}_S(\tilde{r}) \exp(z)$. Show that equation (7.211) becomes

$$u^2 \frac{d^2 z}{du^2} + 2(e^z - 1) = 0. \quad (7.212)$$

(b) By linearizing equation (7.212) for small z , show that the asymptotic behavior of the density $\tilde{\rho}$ is described by the equation (Chandrasekhar 1939)

$$\tilde{\rho}(\tilde{r}) \simeq \tilde{\rho}_S(\tilde{r}) \left[1 + \frac{A}{\tilde{r}^{1/2}} \cos \left(\frac{1}{2} \sqrt{7} \ln \tilde{r} + \phi \right) \right], \quad (7.213)$$

where A and ϕ are constants determined by the boundary conditions at small radii. Thus, at large \tilde{r} , the density of the isothermal sphere *oscillates* around the singular solution, with fractional amplitude decreasing as $\tilde{r}^{-1/2}$.

(c) Now consider an isothermal gas enclosed in a spherical box of radius r_b , with inverse temperature β (cf. eq. 7.55). As $\tilde{r}_b \rightarrow \infty$, show that the mass M of gas can be written in the form

$$x \equiv \frac{r_b}{GMm\beta} \simeq \frac{1}{2} - \frac{A}{8\tilde{r}_b^{1/2}} \left[\cos \left(\frac{1}{2} \sqrt{7} \ln \tilde{r}_b + \phi \right) + \sqrt{7} \sin \left(\frac{1}{2} \sqrt{7} \ln \tilde{r}_b + \phi \right) \right]. \quad (7.214)$$

Hence, argue that when the central concentration $\mathcal{R} = 1/\tilde{\rho}(\tilde{r}_b)$ is large, the curve in Figure 7.1 becomes vertical at successive values of \mathcal{R} that are in the ratio $\exp(4\pi/\sqrt{7}) = 115.54$.

7.7 [1] A particle undergoes a one-dimensional random walk, defined as follows. If the particle is at position x , then during any short time interval Δt the mean-square change in position is $\langle (\Delta x)^2 \rangle = D(x)\Delta t$, while the mean change is $\langle \Delta x \rangle = 0$. Let $p(x, t)dx$ be the probability that the particle is found in the interval $(x, x + dx)$ at time t . What is the partial differential equation governing $p(x, t)$?

7.8 [2] Consider a D -dimensional stellar system containing N identical stars that interact by inverse-square forces ($D = 2$, flat disk; $D = 3$, sphere, etc.). Show that the relaxation time and the crossing time in such a system are related by

$$t_{\text{relax}} \approx \begin{cases} t_{\text{cross}}, & D = 2, \\ N \ln N t_{\text{cross}}, & D = 3, \\ N t_{\text{cross}}, & D > 3. \end{cases} \quad (7.215)$$

7.9 [2] A subject mass M is embedded in an infinite homogeneous sea of field stars, with mass $m \ll M$ and isotropic DF $f(v)$. Using the Fokker–Planck equation and the diffusion coefficients (7.88), show that when the subject mass is in thermal equilibrium with the field stars its DF is Maxwellian, with velocity dispersion

$$\sigma_M^2 = \frac{m}{M} \frac{\int_0^\infty dv v f(v)}{f(0)}. \quad (7.216)$$

Show that when $f(v)$ is Maxwellian, this condition reduces to the requirement of energy equipartition between the subject mass and the field stars.

7.10 [2] A subject mass M is embedded in an infinite homogeneous sea of field stars with number density $n(m)dm$ in the mass range $(m, m+dm)$. The field stars of all masses have a Maxwellian DF with dispersion σ . Thus the field stars are not in thermal equilibrium, which would require that σ^2 is inversely proportional to m . Let us assume, however, that the subject star is in thermal equilibrium with the field stars (this could occur, for example, if $M \gg m$).

(a) Show that the dispersion of the subject star is

$$\sigma_M^2 = \sigma^2 \frac{\int dm m^2 n(m)}{M \int dm m n(m)}. \quad (7.217)$$

(b) Show that the DF of the subject star is Maxwellian.

7.11 [2] The diffusion coefficients for a Maxwellian field star distribution depend on the functions $G(X)$ and $\text{erf}(X) - G(X)$, where $G(X)$ is defined by equation (7.93) and $X = v/(\sqrt{2}\sigma)$. Show that

$$\lim_{X \rightarrow 0} \frac{\text{erf}(X)}{X} = \frac{2}{\sqrt{\pi}} \quad ; \quad \lim_{X \rightarrow 0} \frac{G(X)}{X} = \frac{2}{3\sqrt{\pi}}. \quad (7.218)$$

Thus show that as the velocity of the subject star $v \rightarrow 0$, the diffusion coefficients of equation (7.92) satisfy $D[(\Delta v_{\parallel})^2] = \frac{1}{2}D[(\Delta \mathbf{v}_{\perp})^2]$. Explain physically why this must be so.

7.12 [2] Suppose that a spherical cluster evolves self-similarly as a result of relaxation. In this case its evolution can be described by two functions $M(t)$ and $R(t)$, the mass and characteristic radius as functions of time. Since the evolution is driven by relaxation, we expect that

$$\frac{1}{M} \frac{dM}{dt} = \frac{C_M}{t_{\text{relax}}} \quad ; \quad \frac{1}{R} \frac{dR}{dt} = \frac{C_R}{t_{\text{relax}}} \quad (7.219)$$

where C_M and C_R are constants of order unity. Neglecting changes in the Coulomb logarithm, the relaxation time $t_{\text{relax}} \propto R^{3/2} M^{1/2}$ (eq. 7.108).

(a) If the evolution is dominated by evaporation, we expect that $M(t)$ declines with time while the cluster energy $E \propto GM^2/R$ remains constant, since evaporating stars leave with nearly zero energy. In this case show that

$$M(t) \propto \tau^{2/7} \quad ; \quad R(t) \propto \tau^{4/7}, \quad (7.220)$$

where τ is the time remaining until the cluster disappears.

(b) After core collapse, the evolution of a cluster is dominated by the energy input from binary stars at the cluster center, so the cluster energy E grows but the mass M remains approximately constant. In this case show that

$$R(t) \propto \tau^{2/3}, \quad (7.221)$$

where τ is the time elapsed since core collapse (Goodman 1984).

7.13 [1] Using equation (7.173) for the evaporation time of soft binaries, estimate the maximum semi-major axis of a primordial soft binary that could survive for 10 Gyr in the solar neighborhood. Assume that the DF in the solar neighborhood is isotropic and Maxwellian, with RMS velocity 50 km s^{-1} , that all stars have mass $1 \mathcal{M}_\odot$, and that the stellar density is $\rho = 0.04 \mathcal{M}_\odot \text{ pc}^{-3}$ (from Tables 1.1 and 1.2).

7.14 [3] A population of very hard binaries, each with total mass m_b and internal energy \tilde{E} , is embedded in a distribution of field stars of mass m . The velocity distributions of the field stars and of the centers of mass of the binaries are Maxwellian, with dispersions σ and σ_b respectively. Show that the disruption rate of the binaries contains an exponential factor

$$\exp \left[-\frac{(m + m_b)|\tilde{E}|}{mm_b(\sigma^2 + \sigma_b^2)} \right], \quad (7.222)$$

which reduces to the exponential factor in equation (7.174) when $m_b = 2m$, $\sigma_b^2 = \frac{1}{2}\sigma^2$, and $|\tilde{E}| \gg m\sigma^2$.

7.15 [1] What is the closest approach that a star is likely to have made to the Sun during its lifetime of 4.5 Gyr, assuming that the Sun's environment has always been similar to the present solar neighborhood? Use the same parameters for the solar neighborhood as in Problem 7.13.

7.16 [1] A tidal-capture binary is formed as a result of a close encounter of two stars of equal mass m . The minimum separation during the encounter is d_{\min} , and the orbital energy dissipated in the encounter is $\Delta E \ll Gm^2/d_{\min}$. Once the binary has formed, more energy is dissipated in each successive orbit, until eventually the binary orbit is circularized. If the spin angular momentum of the stars is negligible compared to the orbital angular momentum, show that the radius of the final circular orbit is $2d_{\min}$.

7.17 [2] This problem investigates the distribution of nearest neighbors of stars and the forces from them. Consider an infinite, homogeneous system of stars of mass m and number density n , and assume that the DF is separable, that is, that the two-body correlation function is negligible (§7.2.4).

(a) Show that the probability that the nearest neighbor of a star lies within distance r is

$$1 - e^{-4\pi nr^3/3}. \quad (7.223)$$

(b) Show that the probability that the force per unit mass exerted on a star by its nearest neighbor lies in the range $(F, F + dF)$ is

$$dp_F = \frac{dF}{F_0} W(F/F_0), \quad \text{where } F_0 = Gmn^{2/3} \quad (7.224)$$

and

$$W(\xi) = \frac{2\pi}{\xi^{5/2}} \exp \left(-\frac{4\pi}{3\xi^{3/2}} \right). \quad (7.225)$$

(c) The probability distribution of the total force per unit mass exerted on a star by *all* of its neighbors can be shown to be given by equation (7.224) with

$$W(\xi) = \frac{2}{\pi\xi} \int_0^\infty dx x \sin x e^{-(sx/\xi)^{3/2}}, \quad \text{where } s = 2\pi \left(\frac{4}{15} \right)^{2/3}. \quad (7.226)$$

This is the **Holtmark distribution** (Chandrasekhar 1943b). Using numerical or analytic methods, show that the expressions (7.225) and (7.226) agree for large ξ .

(d) The total force is the sum of a large number of random variables (the forces from individual neighbor stars). Why then is the probability distribution (7.226) not Gaussian, as implied by the central limit theorem (Appendix B.10)? See Feller (1971) for a thorough discussion of the relation between the Holtmark and Gaussian distributions.

7.18 [2] A black hole of mass M_\bullet is embedded in the center of an infinite, homogeneous, three-dimensional sea of test particles. Far from the hole, the test particles have a Maxwellian velocity distribution (7.91) with number density n_0 and velocity dispersion σ . Show that the density distribution of test particles that are not bound to the hole is

$$n(r) = n_0 \left\{ 2\sqrt{\frac{r_F}{\pi r}} + e^{r_F/r} \left[1 - \operatorname{erf} \left(\sqrt{\frac{r_F}{r}} \right) \right] \right\}, \quad (7.227)$$

where the error function $\operatorname{erf}(x)$ is defined in Appendix C.3, and $r_F = GM_\bullet/\sigma^2$. Show that close to the hole, $r \ll r_F$, $n(r) \propto r^{-1/2}$. Thus there is a weak density cusp around the hole, similar in structure to the cusps seen in luminous elliptical galaxies (Nakano & Makino 1999).

7.19 [3] In §7.5.9b we derived the steady-state density distribution of stars around a central black hole, in the case where the relaxation time is shorter than the age of the system. Consider the following alternate derivation. Assume that the density near the hole is a power law, $n(r) \propto r^{-s}$. The mean-square velocity at any radius should be of order $\langle v^2 \rangle \simeq GM_\bullet/r$, so the local relaxation time is $t_{\text{relax}} \approx \langle v^2 \rangle^{3/2} / (G^2 m^2 n) \propto r^{s-3/2}$. Relaxation among the $N(r)$ cusp stars interior to r can lead to a flow of stars through the shell at radius r that is of order $N(r)/t_{\text{relax}} \sim n(r)r^3/t_{\text{relax}} \sim r^{-2s+9/2}$. In a steady state, this flow must be independent of radius, so $s = \frac{9}{4}$. This differs from the (correct) result $s = \frac{7}{4}$ derived in §7.5.9b. What is wrong with the argument presented here?

8

Collisions and Encounters of Stellar Systems

Our Galaxy and its nearest large neighbor, the spiral galaxy M31, are falling towards one another and will probably collide in about 3 Gyr (see Plate 3 and Box 3.1).

A collision between our Galaxy and M31 would have devastating consequences for the gas in both systems. If a gas cloud from M31 encountered a Galactic cloud, shock waves would be driven into both clouds, heating and compressing the gas. In the denser parts of the clouds, the compressed post-shock gas would cool rapidly and fragment into new stars. The most massive of these would heat and ionize much of the remaining gas and ultimately explode as supernovae, thereby shock-heating the gas still further. Depending on the relative orientation of the velocity vectors of the colliding clouds, the post-collision remnant might lose much of its orbital angular momentum, and then fall towards the bottom of the potential well of the whole system, thereby enhancing the cloud-collision and star-formation rates still further. We do not yet have a good understanding of this complex chain of events, but there is strong observational evidence that collisions between gas-rich galaxies like the Milky Way and M31 cause the extremely high star-formation rates observed in starburst galaxies (§8.5.5).

In contrast to gas clouds, stars emerge unscathed from a galaxy collision.

To see this, consider what would happen to the solar neighborhood in a collision with the disk of M31. According to Table 1.1, the surface density of visible stars in the solar neighborhood is $\simeq 30 \mathcal{M}_\odot \text{pc}^{-2}$. Assuming that most of these are similar to the Sun, the number density of stars is $N \simeq 30 \text{pc}^{-2}$ and the fraction of the area of the galactic disk that is filled by the disks of these stars is of order $N\pi R_\odot^2 \approx 5 \times 10^{-14}$. Thus even if M31 were to score a direct hit on our Galaxy, the probability that even one of the 10^{11} stars in M31 would collide with any star in our Galaxy is small.¹

However, the distribution of the stars in the two galaxies would be radically changed by such a collision, because the gravitational field of M31 would deflect the stars of our Galaxy from their original orbits and vice versa for the stars of M31. In this process, which is closely related to violent relaxation (§4.10.2), energy is transferred from ordered motion (the relative motion of the centers of mass of the two galaxies) to random motion. Thus the collision of two galaxies is inelastic, just as the collision of two lead balls is inelastic—in both cases, ordered motion is converted to random motion, of the stars in one case and the molecules in the other (Holmberg 1941; Alladin 1965). Of course, since stars move according to Newton’s laws of motion, the total energy of the galactic system is strictly conserved, in contrast to the lead balls where the energy in random motion of the molecules (i.e., heat) is eventually lost as infrared radiation.

A consequence of this inelasticity is that galaxy collisions often lead to **mergers**, in which the final product of the collision is a single merged stellar system. In fact, we believe that both galaxies and larger stellar systems such as clusters of galaxies are created by a hierarchical or “bottom-up” process in which small stellar systems collide and merge, over and over again, to form ever larger systems (§9.2.2).

The most straightforward way to investigate what happens in galaxy encounters is to simulate the process using an N-body code. Figure 8.1 shows an N-body simulation of the collision of the Galaxy and M31. This is an example of a **major merger**, in which the merging galaxies have similar masses, and the violently changing gravitational field leads to a merger remnant that looks quite different from either of its progenitors. In contrast, **minor mergers**, in which one of the merging galaxies is much smaller than the other, leave the larger galaxy relatively unchanged.

Not every close encounter between galaxies leads to a merger. To see this, let v_∞ be the speed at which galaxy A initially approaches galaxy B and consider how the energy that is gained by a star in galaxy A depends on v_∞ . As we increase v_∞ , the time required for the two galaxies to pass through one another decreases. Hence the velocity impulse $\Delta \mathbf{v} = \int dt \mathbf{g}(t)$ due to the gravitational field $\mathbf{g}(t)$ from galaxy B decreases, and less and less energy is transferred from the relative orbit of the two galaxies to the random motions

¹ We have neglected gravitational focusing, which enhances the collision probability by a factor of about five but does not alter this conclusion (eq. 7.194).

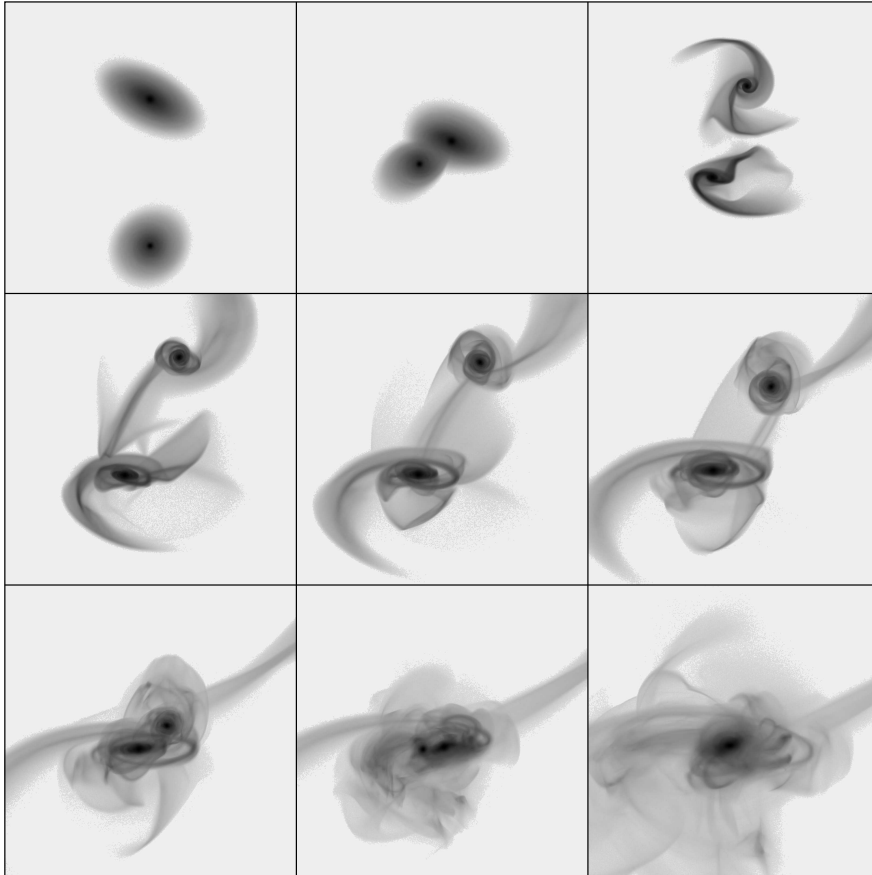


Figure 8.1 An N-body simulation of the collision between the Galaxy (bottom) and M31 (top) which is expected to occur roughly 3 Gyr from now. The simulation follows only the evolution of the stars in the two galaxies, not the gas. Each galaxy is represented by roughly 10^8 stars and dark-matter particles. The viewpoint is from the north Galactic pole. Each panel is 180 kpc across and the interval between frames is 180 Myr. After the initial collision, an open spiral pattern is excited in both disks and long tidal tails are formed. The galaxies move apart by more than 100 kpc and then fall back together for a second collision, quickly forming a remnant surrounded by a complex pattern of shells. The shells then gradually phase mix, eventually leaving a smooth elliptical galaxy. Image provided by J. Dubinski (Dubinski, Mihos, & Hernquist 1996; Dubinski & Farah 2006).

of their stars; in fact, when v_∞ is large, $|\Delta\mathbf{v}| \propto 1/v_\infty$. Thus, when v_∞ exceeds some critical speed v_f , the galaxies complete their interaction with sufficient orbital energy to make good their escape to infinity. If $v_\infty < v_f$, the galaxies merge, while if $v_\infty \gg v_f$ the encounter alters both the orbits and

the internal structures of the galaxies only slightly.² This simple argument explains why most galaxies in rich clusters have not merged: although the density of galaxies in the clusters is high, so collisions are frequent, the random velocities of cluster galaxies are so high that the loss of orbital energy in a collision is negligible—the galaxies simply pass through one another, like ghosts.

Until the 1970s, most astronomers believed that collisions between galaxies were negligible, except in high-density regions such as clusters. This belief was based on the following argument. The velocities of galaxies are the sum of the Hubble velocity (eq. 1.13) appropriate to that galaxy's position, and a residual, or **peculiar velocity**. Typical peculiar velocities are $v_p \approx 100 \text{ km s}^{-1}$ (Willick et al. 1997). The number density of galaxies is described by the Schechter law (eq. 1.18), so the density of luminous galaxies ($L \gtrsim L_*$) is $n \approx \phi_* \approx 10^{-2} \text{ Mpc}^{-3}$. Most of the stars in a typical luminous galaxy are contained within a radius $R \approx 10 \text{ kpc}$, so the collision cross-section between two such galaxies is $\Sigma \approx \pi(2R)^2$. If the positions and velocities of the galaxies are uncorrelated, the rate at which an L_* galaxy suffers collisions with similar galaxies is then expected to be of order $n\Sigma v_p \approx 10^{-6} \text{ Gyr}^{-1}$, so only about one galaxy in 10^5 would suffer a collision during the age of the universe. Such arguments led astronomers to think of galaxies as island universes that formed and lived in isolation.

This estimate of the collision rate turns out to be far too low, for two reasons. (i) The stars in a galaxy are embedded in a dark halo, which can extend to radii of several hundred kpc. Once two dark halos start to merge, their high-density centers, which contain the stars and other baryonic matter, experience a drag force from dynamical friction (§7.4.4) as they move through the common halo. Dynamical friction causes the baryon-rich central regions to spiral towards the center of the merged halo, where they in turn merge. Thus the appropriate cross-section is proportional to the square of the dark-halo radius rather than the square of the radius of the stellar distribution. (ii) As we describe in §9.1, the departures of the matter distribution in the universe from exact homogeneity arose through gravitational forces, and in particular the peculiar velocities of galaxies relative to the Hubble flow are caused by gravitational forces from nearby galaxies. Consequently, the peculiar velocities of nearby galaxies are correlated—nearby galaxies are falling towards one another, just like our Galaxy and M31—so the collision rate is much higher than it would be if the peculiar velocities were randomly oriented. In §8.5.6 we show that the merger rate for L_* galaxies is $\sim 0.01 \text{ Gyr}^{-1}$, 10^4 times larger than our naïve estimate.

When two dark halos of unequal size merge, the smaller halo orbits within the larger one, on a trajectory that steadily decays through dynamical friction. As the orbit decays, the satellite system is subjected to disruptive

² Thus, galaxies behave somewhat like the toy putty that is elastic at high impact speeds, but soft and inelastic at low speeds.

processes of growing strength. These include steady tidal forces from the host galaxy, and rapidly varying forces as the smaller halo passes through the pericenter of its orbit. As stars are lost from the satellite, they spread out in long, thin tidal streamers that can provide vivid evidence of ongoing disruption. Eventually the satellite is completely disrupted, and its stars and dark matter phase-mix with those of the host system.

These processes, which we examine in this chapter, are common to a wide variety of astrophysical systems. Dynamical friction (§8.1) drives the orbital evolution not only of satellite galaxies, but also black holes and globular clusters near the centers of galaxies, and bars in barred spiral galaxies. Tidal forces erode satellite galaxies, globular clusters, and galaxies in clusters, and also determine the lifetimes of star clusters and wide binary stars. We shall focus on the effects of tidal forces in two extreme and analytically tractable limits: §8.2 is devoted to impulsive tides, which last for only a short time, while §8.3 examines the effects of static tides. §8.4 describes the dynamics of encounters in galactic disks, and their effect on the kinematics of stars in the solar neighborhood. Finally, in §8.5 we summarize and interpret the observational evidence for ongoing mergers between galaxies, and estimate the merger rate.

8.1 Dynamical friction

A characteristic feature of collisions of stellar systems is the systematic transfer of energy from their relative orbital motion into random motions of their constituent particles. This process is simplest to understand in the limiting case of minor mergers, in which one system is much smaller than the other.

We consider a body of mass M traveling through a population of particles of individual mass $m_a \ll M$. Following §1.2.1 we call M the subject body and the particles of mass m_a field stars. The subject body usually is a small galaxy or other stellar system and thus has non-zero radius, but we shall temporarily assume that it is a point mass. The field stars are members of a much larger host system of mass $\mathcal{M} \gg M$, which we assume to be so large that it can be approximated as infinite and homogeneous. The influence of encounters with the field stars on the subject body can then be characterized using the diffusion coefficients derived in §7.4.4. Because the test body is much more massive than the field stars, the first-order diffusion coefficients $D[\Delta v_i]$ are much larger than the second-order coefficients $D[\Delta v_i \Delta v_j]/v$ (cf. eqs. 7.83 with $m \gg m_a$). Thus the dominant effect of the encounters is to exert dynamical friction (page 583), which decelerates the subject body at a rate

$$\frac{d\mathbf{v}_M}{dt} = D[\Delta \mathbf{v}] = -4\pi G^2 M m_a \ln \Lambda \int d^3 \mathbf{v}_a f(\mathbf{v}_a) \frac{\mathbf{v}_M - \mathbf{v}_a}{|\mathbf{v}_M - \mathbf{v}_a|^3}, \quad (8.1a)$$

where

$$\Lambda \approx \frac{b_{\max}}{b_{90}} \approx \frac{b_{\max} v_{\text{typ}}^2}{GM} \gg 1 \quad (8.1b)$$

and b_{90} is the 90° deflection radius defined in equation (3.51). Here we have used equations (7.83), assuming $M \gg m_a$ and adjusting the notation appropriately. The field-star DF $f(\mathbf{x}, \mathbf{v}_a)$ is normalized so $\int d^3\mathbf{v}_a f(\mathbf{x}, \mathbf{v}_a) = n(\mathbf{x})$, where n is the number density of field stars in the vicinity of the subject body.

We now estimate the typical value of the factor Λ in the Coulomb logarithm. When a subject body of mass M orbits in a host system of mass $\mathcal{M} \gg M$ and radius \mathcal{R} , the typical relative velocity is given by $v_{\text{typ}}^2 \approx GM/\mathcal{R}$. To within a factor of order unity, the maximum impact parameter $b_{\max} \approx R$, where R is the orbital radius of the subject body. Then $\Lambda \approx (M/\mathcal{M})(R/\mathcal{R})$, which is large whenever $M \ll \mathcal{M}$, unless the subject body is very close to the center of the host.

If the subject body has a non-zero radius, the appropriate value for the Coulomb logarithm is modified to

$$\ln \Lambda = \ln \left(\frac{b_{\max}}{\max(r_h, GM/v_{\text{typ}}^2)} \right), \quad (8.2)$$

where r_h is the half-mass radius of the subject system (see Problem 8.2).

If the field stars have an isotropic velocity distribution,³ equation (7.88) yields a simpler expression for the dynamical friction,

$$\frac{d\mathbf{v}_M}{dt} = -16\pi^2 G^2 M m_a \ln \Lambda \left[\int_0^{v_M} dv_a v_a^2 f(v_a) \right] \frac{\mathbf{v}_M}{v_M^3}; \quad (8.3)$$

thus, only stars moving slower than M contribute to the friction. Like an ordinary frictional drag, the force described by equation (8.3) always opposes the motion ($d\mathbf{v}_M/dt$ is anti-parallel to \mathbf{v}_M). Equation (8.3) is usually called **Chandrasekhar's dynamical friction formula** (Chandrasekhar 1943a).

If the subject mass is moving slowly, so v_M is sufficiently small, we may replace $f(v_a)$ in the integral of equation (8.3) by $f(0)$ to find

$$\frac{d\mathbf{v}_M}{dt} \simeq -\frac{16\pi^2}{3} G^2 M m_a \ln \Lambda f(0) \mathbf{v}_M \quad (v_M \text{ small}). \quad (8.4)$$

Thus at low velocity the drag is proportional to v_M —just as in Stokes's law for the drag on a marble falling through honey. On the other hand, for sufficiently large v_M , the integral in equation (8.3) converges to a definite limit equal to the number density n divided by 4π :

$$\frac{d\mathbf{v}_M}{dt} = -4\pi G^2 M m_a n \ln \Lambda \frac{\mathbf{v}_M}{v_M^3} \quad (v_M \text{ large}). \quad (8.5)$$

³ See Problem 8.3 for the case of an ellipsoidal velocity distribution.

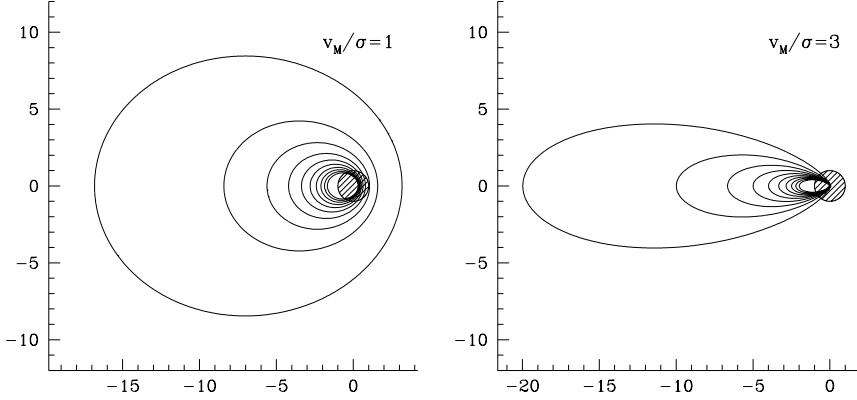


Figure 8.2 A mass M travels from left to right at speed v_M , through a homogeneous Maxwellian distribution of stars with one-dimensional dispersion σ . Deflection of the stars by the mass enhances the stellar density downstream, and the gravitational attraction of this wake on M leads to dynamical friction. The contours show lines of equal stellar density in a plane containing the mass M and the velocity vector \mathbf{v}_M ; the velocities are $v_M = \sigma$ (left panel) and $v_M = 3\sigma$ (right panel). The fractional overdensities shown are 0.1, 0.2, ..., 0.9, 1. The unit of length is chosen so that $GM/\sigma^2 = 1$. The shaded circle has unit radius and is centered at M . The overdensities are computed using equation (8.148), which is based on linear response theory; for a nonlinear treatment see Mulder (1983).

Thus the frictional force *falls* like v_M^{-2} —in contrast to the motion of solid bodies through fluids, where the drag force *grows* as the velocity increases.

If $f(\mathbf{v}_a)$ is Maxwellian with dispersion σ , then equation (8.3) becomes (cf. eqs. 7.91–7.93)

$$\frac{d\mathbf{v}_M}{dt} = -\frac{4\pi G^2 M n m \ln \Lambda}{v_M^3} \left[\operatorname{erf}(X) - \frac{2X}{\sqrt{\pi}} e^{-X^2} \right] \mathbf{v}_M, \quad (8.6)$$

where $X \equiv v_M/(\sqrt{2}\sigma)$ and erf is the error function (Appendix C.3). This important formula illustrates two features of dynamical friction:

- (i) The frictional drag is proportional to the mass density nm of the stars being scattered, but independent of the mass of each individual star. In particular, if we replace nm in equation (8.6) by the overall background density ρ , we obtain a formula that is equally valid for a host system containing a spectrum of different stellar masses:

$$\frac{d\mathbf{v}_M}{dt} = -\frac{4\pi G^2 M \rho \ln \Lambda}{v_M^3} \left[\operatorname{erf}(X) - \frac{2X}{\sqrt{\pi}} e^{-X^2} \right] \mathbf{v}_M. \quad (8.7)$$

- (ii) The frictional acceleration is proportional to M and thus the frictional force must be proportional to M^2 . It is instructive to consider why

this is so. Stars are deflected by M in such a way that the density of background stars behind M is greater than in front of it (see Figure 8.2 and Problem 8.4). The amplitude of this density enhancement or wake is proportional to M and the gravitational force that it exerts on M is proportional to M times its amplitude. Hence the force is proportional to M^2 .

The validity of Chandrasekhar’s formula Although Chandrasekhar’s formula (8.3) was derived for a mass moving through an infinite homogeneous background, it can be (and usually is) employed to estimate the drag on a small body traveling through a much larger host system. In such applications we replace $f(v)$ by the value of the DF in the vicinity of the small body, v_{typ} by the local velocity dispersion, and b_{max} by the distance of the subject body from the center of the host. When employed in this way, Chandrasekhar’s formula suffers from several internal inconsistencies:

- (i) The choices of b_{max} and v_{typ} are rather arbitrary.
- (ii) We have neglected the self-gravity of the wake. Thus equation (8.3) takes into account the mutual attraction of M and the background stars, but neglects the attraction of the background stars for each other.
- (iii) We obtained equation (8.3) in the approximation that stars move past M on Keplerian hyperbolae. Orbits in the combined gravitational fields of M and the host system would be more complex.

These deficiencies become especially worrisome when M is so large as to be comparable to the mass of the host system that lies interior to M ’s orbit. Nevertheless, N-body simulations and linearized response calculations show that Chandrasekhar’s formula provides a remarkably accurate description of the drag experienced by a body orbiting in a stellar system, usually within a factor of two and often considerably better (Weinberg 1989; Fujii, Funato, & Makino 2006).

The fundamental reasons for this success were discussed in the derivation of the Fokker–Planck equation in §7.4.2, and derive from the large ratio between the maximum and minimum impact parameters that contribute to the Coulomb logarithm $\ln \Lambda = \ln(b_{\text{max}}/b_{90})$. Consider, for example, a black hole of mass $M = 10^6 \mathcal{M}_{\odot}$, orbiting at radius 1 kpc in a galaxy with velocity dispersion 200 km s^{-1} . Then we may set $b_{\text{max}} \approx 1 \text{ kpc}$ and $V_0 \approx 200 \text{ km s}^{-1}$, so $b_{90} = 0.1 \text{ pc}$ and $\ln \Lambda = 9.2$. To address the seriousness of problem (i) above, suppose that we have overestimated b_{max} by a factor of two, so the correct value is only half the orbital radius or 0.5 kpc; then $\ln \Lambda = 8.5$, a change of less than 10%. In words, the drag force is insensitive to changes of order unity in b_{max} and v_{typ} , because $\ln \Lambda$ is large. To address problems (ii) and (iii) we note the effects of self-gravity are important only on scales comparable to the Jeans length, which in turn is comparable to b_{max} . Thus the effects of self-gravity are negligible, and the approximation of a Keplerian hyperbola should be valid, for encounters with impact parameter much less than b_{max} . Suppose then that we consider only encounters with $b < 100 \text{ pc}$

or 10% of the orbital radius. The contribution to the Coulomb logarithm from these encounters is $\ln(100 \text{ pc}/b_{90}) = 5.5$, a difference of only 25% from our original estimate. In words, most of the total contribution to the drag comes from encounters with sufficiently small impact parameters that the neglect of self-gravity and the approximation of Keplerian orbits introduce negligible errors.

A more sophisticated treatment of dynamical friction that avoids the inconsistencies of Chandrasekhar's formula requires the machinery of linear response theory that was developed in §5.3. The subject body is regarded as an external potential $\Phi_e(\mathbf{x}, t)$ that excites a response density in the host system, governed by the response function $R(\mathbf{x}, \mathbf{x}', \tau)$. We then solve Poisson's equation to determine the gravitational potential generated by the response density, and the force exerted on the subject body by this response potential is dynamical friction (Weinberg 1986, 1989).

Like Landau damping, dynamical friction illustrates the curious fact that irreversible processes can occur in a system with reversible equations of motion. We have seen in §5.5.3 that Landau damping in spherical stellar systems arises from resonances between the oscillations of the system and the orbital frequencies of individual stars. Similarly, dynamical friction can be shown to arise from resonances between the orbital frequencies of the subject body and the stars (Tremaine & Weinberg 1984b). The rich Fourier spectrum of the gravitational potential from an orbiting point mass ensures that many orbital resonances contribute to the drag force, and the cumulative effect of these many weak resonances gives rise to the Coulomb logarithm in Chandrasekhar's formula.

8.1.1 Applications of dynamical friction

(a) Decay of black-hole orbits The centers of galaxies often contain black holes with masses 10^6 – $10^9 M_\odot$ (§1.1.6). It is natural to ask whether such objects could also be present at other locations in the galaxy, where they would be even harder to find. To investigate this question, we imagine a black hole of mass M on a circular orbit of radius r , and ask how long is needed for dynamical friction to drag the black hole to the galaxy center.

The flatness of many observed rotation curves suggests that we approximate the density distribution by a singular isothermal sphere (eq. 4.103),

$$\rho(r) = \frac{v_c^2}{4\pi G r^2}, \quad (8.8)$$

where $v_c = \sqrt{2}\sigma$ is the constant circular speed (eq. 4.104). The DF of the isothermal sphere is Maxwellian, so equation (8.7) gives the frictional force

$\tilde{F} = M|d\mathbf{v}_M/dt|$ on the black hole:

$$\begin{aligned}\tilde{F} &= \frac{4\pi G^2 M^2 \rho(r) \ln \Lambda}{v_c^2} \left[\operatorname{erf}(X) - \frac{2X}{\sqrt{\pi}} e^{-X^2} \right] \\ &= 0.428 \ln \Lambda \frac{GM^2}{r^2},\end{aligned}\quad (8.9)$$

where $X = v_c/(\sqrt{2}\sigma) = 1$.

This force is tangential and directed opposite to the velocity of the black hole, causing it to lose angular momentum \tilde{L} at a rate

$$\frac{d\tilde{L}}{dt} = -\tilde{F}r \simeq -0.428 \ln \Lambda \frac{GM^2}{r}. \quad (8.10)$$

Thus the black hole spirals towards the center of the galaxy, while remaining on a nearly circular orbit. Since the circular-speed curve of the singular isothermal sphere is flat, the black hole continues to orbit at speed v_c as it spirals inward, so its angular momentum at radius r is $\tilde{L} = Mrv_c$. Substituting the time derivative of this expression into equation (8.10), we obtain

$$r \frac{dr}{dt} = -0.428 \ln \Lambda \frac{GM}{v_c} = -0.302 \ln \Lambda \frac{GM}{\sigma}. \quad (8.11)$$

If we neglect the slow variation of $\ln \Lambda$ with radius, we can solve this differential equation subject to the initial condition that the radius is r_i at zero time. We find that the black hole reaches the center after a time⁴

$$t_{\text{fric}} = \frac{1.65}{\ln \Lambda} \frac{r_i^2 \sigma}{GM} = \frac{19 \text{ Gyr}}{\ln \Lambda} \left(\frac{r_i}{5 \text{ kpc}} \right)^2 \frac{\sigma}{200 \text{ km s}^{-1}} \frac{10^8 \mathcal{M}_\odot}{M}. \quad (8.12)$$

This equation can be cast into a simpler form using the crossing time $t_{\text{cross}} = r_i/v_c$, the time required for the subject body to travel one radian,

$$t_{\text{fric}} = \frac{1.17}{\ln \Lambda} \frac{\mathcal{M}(r)}{M} t_{\text{cross}}, \quad (8.13)$$

where $\mathcal{M}(r) = v_c^2 r/G$ is the mass of the host galaxy contained within radius r . This result is approximately correct even for mass distributions other than the singular isothermal sphere; in words, if the ratio of the mass of the subject body to the interior mass of the host is $\mu \ll 1$, then the subject body spirals to the center of the host in roughly $1/(\mu \ln \Lambda)$ initial crossing times.

For characteristic values $b_{\text{max}} \approx 5 \text{ kpc}$, $M = 10^8 \mathcal{M}_\odot$, and $v_{\text{typ}} \approx \sigma = 200 \text{ km s}^{-1}$, we have by equation (8.1b) that $\ln \Lambda \simeq 6$. Thus for the standard

⁴Equation (8.8) overestimates the density inside the galaxy's core, but this leads to a negligible error in the inspiral time, since the decay is rapid at small radii anyway.

parameters in equation (8.12), the inspiral time t_{fric} is only 3 Gyr. Black holes on eccentric orbits have even shorter inspiral times than those on circular orbits with the same mean radius, since the eccentric orbit passes through regions of higher density where the drag force is stronger. We conclude that any $10^8 M_{\odot}$ black hole that is formed within ~ 10 kpc of the center of a typical galaxy will spiral to the center within the age of the universe. Thus massive black holes should normally be found at the center of the galaxy, unless they are far out in the galactic halo.

(b) Galactic cannibalism Most large galaxies are accompanied by several satellite galaxies, small companion galaxies that travel on bound orbits in the gravitational potential of the larger host. The satellites of our own Milky Way galaxy include the Sagittarius dwarf galaxy, the Large and Small Magellanic Clouds (§1.1.3 and Plate 11), and several dozen even smaller galaxies at distances of ~ 100 – 300 kpc. Two satellite galaxies of the nearby disk galaxy M31 appear in Plate 3.

Satellites orbiting within the extended dark halo of their host experience dynamical friction, leading to orbital decay. As the satellite orbit decays, tidal forces from the host galaxy (§8.3) strip stars from the outer parts of the satellite, until eventually the entire satellite galaxy is disrupted—this process, in which a galaxy consumes its smaller neighbors, is an example of a minor merger, or, more colorfully, **galactic cannibalism**.

The rate of orbital decay for a satellite of fixed mass M is described approximately by equation (8.12). This formula does not, however, allow for mass loss due to tidal stripping as the satellite spirals inward. To account crudely for this process, we shall refer forward to §8.3, in which we show that the outer or tidal radius of a satellite is given approximately by its Jacobi radius r_J , defined by equation (8.91). Once again we assume that the host galaxy is a singular isothermal sphere, so its mass interior to radius r is $\mathcal{M}(r) = v_{\mathcal{M}}^2 r / G = 2\sigma_{\mathcal{M}}^2 r / G$, where $v_{\mathcal{M}}$ and $\sigma_{\mathcal{M}} = v_{\mathcal{M}} / \sqrt{2}$ are the circular speed and velocity dispersion of the host; in this case equations (8.91) and (8.108) yield

$$r_J = \left(\frac{M}{2\mathcal{M}(r)} \right)^{1/3} \quad r = \left(\frac{GM r^2}{4\sigma_{\mathcal{M}}^2} \right)^{1/3}. \quad (8.14)$$

We shall assume that the satellite galaxy is also a singular isothermal sphere, but one that is sharply truncated at r_J . Thus the total mass of the satellite is $M = 2\sigma_s^2 r_J / G$, where σ_s is its velocity dispersion. (A truncated isothermal sphere is not a self-consistent solution of the collisionless Boltzmann and Poisson equations, so this should be regarded as a fitting formula without much dynamical significance.) Equation (8.14) then yields

$$r_J = \frac{\sigma_s}{\sqrt{2}\sigma_{\mathcal{M}}} r \quad \text{which implies that} \quad M = \frac{\sqrt{2}\sigma_s^3 r}{G\sigma_{\mathcal{M}}}. \quad (8.15)$$

Substituting into equation (8.11), we obtain the rate of orbital decay,

$$\frac{dr}{dt} = -0.428 \ln \Lambda \frac{\sigma_s^3}{\sigma_{\mathcal{M}}^2}, \quad (8.16)$$

and neglecting the slow variation in $\ln \Lambda$ with radius, we find the inspiral time from radius r_i to be

$$\begin{aligned} t_{\text{fric}} &= \frac{2.34 \sigma_{\mathcal{M}}^2}{\ln \Lambda \sigma_s^3} r_i \\ &= \frac{2.7 \text{ Gyr}}{\ln \Lambda} \frac{r_i}{30 \text{ kpc}} \left(\frac{\sigma_{\mathcal{M}}}{200 \text{ km s}^{-1}} \right)^2 \left(\frac{100 \text{ km s}^{-1}}{\sigma_s} \right)^3. \end{aligned} \quad (8.17)$$

To evaluate the Coulomb logarithm, we use equation (8.2). The half-mass radius r_h of the satellite is half of its Jacobi radius, and the typical velocity may be taken to be $v_{\text{typ}} = \sigma_{\mathcal{M}}$. Then the two quantities in the denominator of equation (8.2) are given by equations (8.15),

$$r_h = \frac{\sigma_s}{2^{3/2} \sigma_{\mathcal{M}}} r \quad ; \quad \frac{GM}{v_{\text{typ}}^2} = \frac{\sqrt{2} \sigma_s^3}{\sigma_{\mathcal{M}}^3} r. \quad (8.18)$$

The velocity dispersion of a galaxy is correlated with its mass through the Faber–Jackson law (1.21). Satellite galaxies have smaller luminosities than their hosts, and hence smaller dispersions. If $\sigma_s \lesssim 0.5 \sigma_{\mathcal{M}}$, the first term in equation (8.18) is larger than the second, so the argument of the Coulomb logarithm is $\Lambda = b_{\text{max}}/r_h$; setting $b_{\text{max}} = r$ we have finally $\Lambda = 2^{3/2} \sigma_{\mathcal{M}}/\sigma_s$. Thus, for example, equation (8.17) implies that in a host galaxy with dispersion 200 km s^{-1} , a satellite galaxy with dispersion $\sigma \gtrsim 50 \text{ km s}^{-1}$ will merge from a circular orbit with radius 30 kpc within 10 Gyr.

(c) Orbital decay of the Magellanic Clouds In general, the orbits of satellites of the Milky Way cannot be determined, because their velocities perpendicular to the line of sight are either unknown or have large observational uncertainties. However, much more information is available for the Large and Small Magellanic Clouds. Not only do we have good estimates for their velocities perpendicular to the line of sight (Kallivayalil et al. 2006), but the correct Cloud orbits must be able to reproduce the dynamics of the Magellanic Stream, a narrow band of neutral hydrogen gas that extends over 120° in the sky and is believed to have been torn off the Small Cloud by the gravitational field of the Galaxy about 1–1.5 Gyr ago. (See BM §8.4.1 and Putman et al. 2003 for a description of the observations.)

Several groups have modeled the dynamics of the Magellanic Stream and the resulting constraints on the Cloud orbits (Murai & Fujimoto 1980; Lin & Lynden–Bell 1982; Gardiner, Sawa, & Fujimoto 1994; Connors et al. 2004). They find that the orbital plane of the Clouds is nearly perpendicular

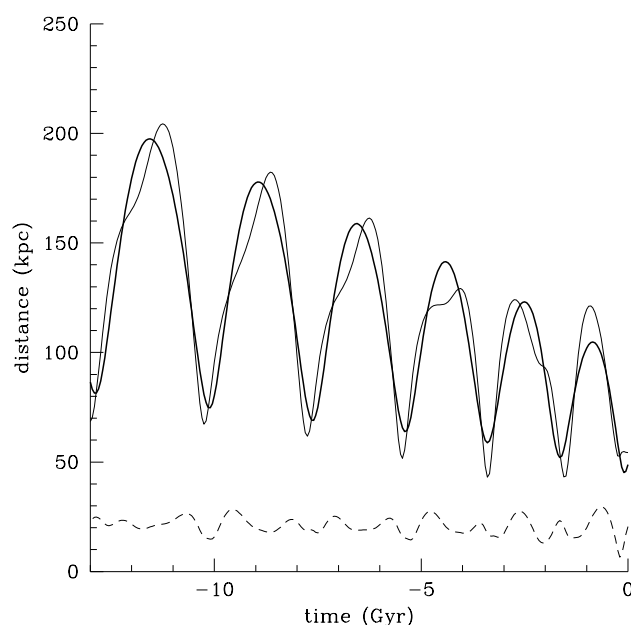


Figure 8.3 The decay of the orbits of the Magellanic Clouds around our Galaxy. The upper curves show the radius of the Clouds from the Galactic center (thick line for the Large Cloud and thin line for the Small Cloud), and the lower, dashed curve shows the distance between the Large and Small Cloud. The Galaxy potential is that of a singular isothermal sphere with circular speed $v_c = 220 \text{ km s}^{-1}$, and the drag force is computed using Chandrasekhar's formula (8.7). The initial conditions at $t = 0$ are chosen to reproduce the observed distances and radial velocities of the Clouds and the kinematics of the Magellanic Stream (Gardiner, Sawa, & Fujimoto 1994).

to the Galactic plane; the sense of the orbit is such that the Clouds are approaching the Galactic plane with the Magellanic Stream trailing behind; the orbit is eccentric (the apocenter/pericenter distance is $\gtrsim 2$); and the Clouds are presently near pericenter (Figure 8.3). As seen in the figure, the orbits of the Magellanic Clouds are decaying due to dynamical friction. The ongoing mass loss from the Clouds that generates the Magellanic Stream provides circumstantial evidence that the orbit is continuing to shrink.

In this model the Clouds merge with the Milky Way in about 6 Gyr, although the model is unrealistic beyond about 3 Gyr in the future, when the Galaxy experiences a much more violent merger with M31 (Box 3.1).

(d) Dynamical friction on bars Dynamical friction can be generated by any time-varying large-scale gravitational field. An important example is the interaction between a bar in a disk galaxy and the surrounding dark halo. As a first approximation, let us think of the bar as a rigid body, consisting

of two masses M at the ends of a rod of length $2r$ that revolves around its center at the bar pattern speed Ω_b . In a strong bar M would not be much smaller than the mass of the galaxy interior to r . In this circumstance equation (8.13) suggests that the bar should lose its angular momentum in a few crossing times, which is much shorter than the age of the galaxy. Thus we might expect that bars in disk galaxies with massive halos would have zero angular momentum and zero pattern speed.

Improving on this crude model is a challenging analytic task, for several reasons: first, the gravitational potential of a bar is more complicated than the potential from a point mass; second, in contrast to most orbiting bodies, bars are extended objects, so the friction is not dominated by local encounters; third, dynamical friction exerts a torque on the bar but we do not understand the reaction of the bar to that torque: does its pattern speed increase or decrease? does the bar grow stronger or weaker? etc.

Accurate analytic determinations of the frictional torque on a bar from the dark halo can be derived using perturbation theory (Weinberg 1985), but N-body simulations can be more informative because they determine both the torque on the bar and its resulting evolution. Simulations confirm that the halo exerts a strong frictional torque on the bar, and show that in response the bar pattern speed rapidly decays but the bar remains intact (Sellwood 1980; Hernquist & Weinberg 1992; Debattista & Sellwood 1998, 2000).

These theoretical results imply that if massive dark halos are present in the inner parts of barred galaxies, bars should rotate slowly. However, this conclusion is inconsistent with observations: as we saw in §6.5.1a the ratio \mathcal{R} of the corotation radius to the bar semi-major axis (eq. 6.103) generally lies in the range 0.9–1.3, where $\mathcal{R} \simeq 1$ is the maximum allowed rotation rate for a weak bar. This problem can be resolved if spiral galaxies have maximum disks (§6.3.4), for then the halo mass is relatively small in the inner few kpc, where interactions with the bar are strongest.

(e) Formation and evolution of binary black holes Since most galaxies contain black holes at their centers, it is natural to ask what happens to the black holes when a satellite galaxy merges with a larger host.

As the satellite's orbit decays, its stars are stripped by tidal forces that become stronger and stronger as the orbit shrinks (eq. 8.15), until eventually only its central black hole is left. The orbit of the black hole continues to decay from dynamical friction, although at a slower rate since the mass of the black hole is only a small fraction of the mass of the original satellite galaxy. Assuming that the host galaxy also contains a central black hole, we expect that eventually the two black holes will form a bound binary system.

After the black-hole binary is formed, its orbit continues to decay by dynamical friction. Equation (8.3) still describes the drag acting on each black hole, with the maximum impact parameter b_{\max} appearing in the Coulomb logarithm now equal to the binary semi-major axis a .

As the binary orbit shrinks, the relative orbital velocity v of the two black holes grows. Eventually the orbital velocity greatly exceeds the velocity dispersion σ of the stars in the galaxy. For a circular orbit, this occurs when the binary semi-major axis a satisfies

$$\frac{G(M_1 + M_2)}{a} \gg \sigma^2 \quad \text{or} \quad a \ll 10 \text{ pc} \frac{M_1 + M_2}{10^8 \mathcal{M}_\odot} \left(\frac{200 \text{ km s}^{-1}}{\sigma} \right)^2, \quad (8.19)$$

where M_1 and M_2 are the masses of the black holes. Following the discussion in §7.5.7, we shall say that the black-hole binary is hard when $v > \sigma$.

For hard binaries Chandrasekhar's dynamical friction formula is no longer valid, but an approximate formula for the rate of orbital decay can be derived by arguments similar to those used to derive the hardening rate for binary stars in equation (7.179). These yield (Quinlan 1996b)

$$\frac{d}{dt} \left(\frac{1}{a} \right) = -C \frac{G\rho}{\sigma}, \quad C = 14.3, \quad (8.20)$$

where ρ is the density of stars in the vicinity of the binary. This result is almost independent of the eccentricity of the binary and depends only weakly on the mass ratio M_2/M_1 so long as it is not too far from unity.⁵

Under the assumption that the galaxy has a constant-density core, we can integrate equation (8.20) to obtain $1/a = \text{constant} - CG\rho t/\sigma$. Choosing the origin of time so that the constant is zero, we obtain $a = \sigma/(CG\rho t)$. The King radius of the galaxy, r_0 , is related to ρ and σ via $4\pi G\rho r_0^2 = 9\sigma^2$ (eq. 4.106). Eliminating ρ with the help of this equation, we have finally

$$a(t) = \frac{4\pi r_0^2}{9C\sigma t} = 0.005 \text{ pc} \frac{200 \text{ km s}^{-1}}{\sigma} \left(\frac{r_0}{100 \text{ pc}} \right)^2 \frac{\text{Gyr}}{t}. \quad (8.21)$$

Thus interactions with stars in the host galaxy can drive the black-hole binary to semi-major axes as small as a few milliparsecs; the corresponding relative speed for a circular orbit is

$$v = \sqrt{\frac{G(M_1 + M_2)}{a}} = 2.1 \times 10^4 \text{ km s}^{-1} \left(\frac{M_1 + M_2}{10^8 \mathcal{M}_\odot} \frac{10^{-3} \text{ pc}}{a} \right)^{1/2}. \quad (8.22)$$

There is an important case in which this analysis fails. Only stars with angular momentum $L \lesssim [G(M_1 + M_2)a]^{1/2}$ interact strongly with the binary, and if the binary semi-major axis a is much smaller than the King radius r_0 then this is much smaller than the typical angular momentum $L \sim r_0\sigma$ of

⁵ The numerical coefficient differs from the one in equation (7.179) because here the binary components are much more massive than the field stars, while in equation (7.179) the binary components and the field stars all have the same mass.

stars in the core. Thus, only a small fraction of the stars in the core interact strongly with the black holes. The region in phase space with such small angular momentum is called the loss cone, by analogy with the loss cone from which stars are consumed by a single black hole (§7.5.9). The mass of stars in the loss cone shrinks as the binary semi-major axis decreases, and eventually may become smaller than the black-hole mass. In this case the binary can empty the loss cone, and the shrinkage of the semi-major axis will stall. Once the loss cone has been emptied, the rate of continued evolution is much less certain, being determined by the rate at which the loss cone is slowly refilled by two processes: diffusion of angular momentum due to two-body relaxation (Chapter 7), or torques from the host galaxy, if its overall mass distribution is non-spherical (Yu 2002; Makino & Funato 2004).

If the binary semi-major axis shrinks far enough, gravitational radiation takes over as the dominant cause of orbital decay. A binary black hole on a circular orbit with semi-major axis a will coalesce under the influence of gravitational radiation in a time (Peters 1964)

$$\begin{aligned}
 t_{\text{gr}} &= \frac{5c^5 a^4}{256G^3 M_1 M_2 (M_1 + M_2)} \\
 &= 5.81 \text{ Myr} \left(\frac{a}{0.01 \text{ pc}} \right)^4 \left(\frac{10^8 \mathcal{M}_\odot}{M_1 + M_2} \right)^3 \frac{(M_1 + M_2)^2}{M_1 M_2}.
 \end{aligned} \tag{8.23}$$

The characteristic decay time due to gravitational radiation therefore scales as a^4 . In contrast, the decay time $(d \ln a / dt)^{-1}$ due to dynamical friction varies as $1/a$. Consequently, the actual decay time, which is set by the more efficient of the two processes, has a maximum at the semi-major axis where the two decay times are equal (Begelman, Blandford & Rees 1980). This radius is referred to as the **bottleneck radius**, and lies between 0.003 pc and 3 pc depending on the galaxy density distribution and black-hole masses (Yu 2002). The bottleneck radius is where binary black holes are most likely to be found.

The decay time at the bottleneck is quite uncertain, since it depends on both the extent to which the loss cone is depopulated, and the possible contribution of gas drag. If the bottleneck decay time exceeds 10 Gyr, most galaxies should contain binary black holes at their centers. If the decay time is less than 10 Gyr, most black-hole binaries will eventually coalesce. Coalescing black holes are of great interest because they generate strong bursts of gravitational radiation that should ultimately be detectable, even at cosmological distances, and thus provide a unique probe of both galaxy evolution and general relativity.

(f) Globular clusters These systems may experience significant orbital decay from dynamical friction. The rate of decay and inspiral time can be described approximately by equations (8.11) and (8.12). For a typical cluster

mass $M = 2 \times 10^5 M_\odot$ (Table 1.3) the inspiral time from radius r_i is

$$t_{\text{fric}} = 64 \text{ Gyr} \frac{\sigma}{200 \text{ km s}^{-1}} \left(\frac{r_i}{1 \text{ kpc}} \right)^2, \quad (8.24)$$

where σ is the velocity dispersion of the host galaxy and we have assumed $\ln \Lambda = 5.8$, from equation (8.2) with $b_{\text{max}} = 1 \text{ kpc}$ and $r_h = 3 \text{ pc}$ (Table 1.3). Orbital decay is most important for low-luminosity host galaxies, which have small radii and low velocity dispersions. Many dwarf elliptical galaxies exhibit a deficit of clusters near their centers and compact stellar nuclei, which may arise from the inspiral and merger of these clusters (Lotz et al. 2001). A puzzling exception is the Fornax dwarf spheroidal galaxy, a satellite of the Milky Way, which contains five globular clusters despite an estimated inspiral time of only $t_{\text{fric}} \sim 1 \text{ Gyr}$ (Tremaine 1976a). Why these clusters have not merged at the center of Fornax remains an unsolved problem.

8.2 High-speed encounters

One of the most important classes of interaction between stellar systems is high-speed encounters. By “high-speed” we mean that the duration of the encounter—the interval during which the mutual gravitational forces are significant—is short compared to the crossing time within each system. A typical example is the collision of two galaxies in a rich cluster of galaxies (§1.1.5). The duration of the encounter is roughly the time it takes the two galaxies to pass through one another; given a galaxy size $r \sim 10 \text{ kpc}$ and the typical encounter speed in a rich cluster, $V \approx 2000 \text{ km s}^{-1}$, the duration is $r/V \approx 5 \text{ Myr}$. For comparison the internal dispersion of a large galaxy is $\sigma \approx 200 \text{ km s}^{-1}$ so the crossing time is $r/\sigma \approx 50 \text{ Myr}$, a factor of ten larger.

As we saw at the beginning of this chapter, the effect of an encounter on the internal structure of a stellar system *decreases* as the encounter speed increases. Hence high-speed encounters can be treated as small perturbations of otherwise steady-state systems.

We consider an encounter between a stellar system of mass M_s , the subject system, and a passing perturber—a galaxy, gas cloud, dark halo, black hole, etc.—of mass M_p . At the instant of closest approach, the centers of the subject system and the perturber are separated by distance b and have relative speed V . If the relative speed is high enough, then:

- (i) The kinetic energy of relative motion of the two systems is much larger than their mutual potential energy, so the centers travel at nearly uniform velocity throughout the encounter.
- (ii) In the course of the encounter, the majority of stars will barely move from their initial locations with respect to the system center. Thus the

gravitational force from the perturber can be approximated as an impulse of very short duration, which changes the velocity but not the position of each star. A variety of analytic arguments (see page 658) and numerical experiments (Aguilar & White 1985) suggest that this **impulse approximation** yields remarkably accurate results, even when the duration of the encounter is almost as long as the crossing time.⁶

We now ask how the structure of the subject system is changed by the passage of the perturber. We work in a frame that is centered on the center of mass of the subject system before the encounter. Let m_α be the mass of the α th star of the subject system, and let $\dot{\mathbf{v}}'_\alpha$ be the rate of change in its velocity due to the force from the perturber. We break $\dot{\mathbf{v}}'_\alpha$ into two components. The component that reflects the rate of change of the center-of-mass velocity of the subject system is

$$\dot{\mathbf{v}}_{\text{cm}} \equiv \frac{1}{M_s} \sum_{\beta} m_{\beta} \dot{\mathbf{v}}'_{\beta}, \quad \text{where} \quad M_s \equiv \sum_{\beta} m_{\beta} \quad (8.25a)$$

is the mass of the subject system. The component

$$\dot{\mathbf{v}}_{\alpha} \equiv \dot{\mathbf{v}}'_{\alpha} - \dot{\mathbf{v}}_{\text{cm}} \quad (8.25b)$$

gives the acceleration of the α th star with respect to the center of mass.

Let $\Phi(\mathbf{x}, t)$ be the gravitational potential due to the perturber. Then

$$\dot{\mathbf{v}}'_{\alpha} = -\nabla\Phi(\mathbf{x}_{\alpha}, t), \quad (8.26)$$

and equation (8.25b) can be written

$$\dot{\mathbf{v}}_{\alpha} = -\nabla\Phi(\mathbf{x}_{\alpha}, t) + \frac{1}{M_s} \sum_{\beta} m_{\beta} \nabla\Phi(\mathbf{x}_{\beta}, t). \quad (8.27)$$

In the impulse approximation, \mathbf{x}_{α} is constant during an impulsive encounter, so

$$\Delta\mathbf{v}_{\alpha} = \int_{-\infty}^{\infty} dt \dot{\mathbf{v}}_{\alpha} = \int_{-\infty}^{\infty} dt \left[-\nabla\Phi(\mathbf{x}_{\alpha}, t) + \frac{1}{M_s} \sum_{\beta} m_{\beta} \nabla\Phi(\mathbf{x}_{\beta}, t) \right]. \quad (8.28)$$

The potential energy of the subject system does not change during the encounter, so in the center-of-mass frame the change in the energy, \tilde{E} , is simply

⁶ Condition (ii) almost always implies condition (i), but condition (i) need not imply condition (ii): for example, a star passing by the Sun at a relative velocity $v \simeq 50 \text{ km s}^{-1}$ and an impact parameter $b \simeq 0.01 \text{ pc}$ will hardly be deflected at all and hence satisfies condition (i). However, the encounter time $b/v \simeq 200 \text{ yr}$ is much larger than the orbital period of most of the planets so the encounter is adiabatic, rather than impulsive.

the change in the internal kinetic energy, \tilde{K} . Here tildes on the symbols are a reminder that these quantities have units of mass \times (velocity)², in contrast to the usual practice in this book where E and K denote energy per unit mass. We have

$$\begin{aligned}\Delta\tilde{E} = \Delta\tilde{K} &= \frac{1}{2} \sum_{\alpha} m_{\alpha} [(\mathbf{v}_{\alpha} + \Delta\mathbf{v}_{\alpha})^2 - \mathbf{v}_{\alpha}^2] \\ &= \frac{1}{2} \sum_{\alpha} m_{\alpha} [|\Delta\mathbf{v}_{\alpha}|^2 + 2\mathbf{v}_{\alpha} \cdot \Delta\mathbf{v}_{\alpha}].\end{aligned}\tag{8.29}$$

In any static axisymmetric system, $\sum_{\alpha} m_{\alpha} \mathbf{v}_{\alpha} \cdot \Delta\mathbf{v}_{\alpha} = 0$ by symmetry (see Problem 8.5). Thus the energy changes that are first-order in the small quantity $\Delta\mathbf{v}$ average to zero, and the change of internal energy is given by the second-order quantity

$$\Delta\tilde{E} = \Delta\tilde{K} = \frac{1}{2} \sum_{\alpha} m_{\alpha} |\Delta\mathbf{v}_{\alpha}|^2.\tag{8.30}$$

This simple derivation masks several subtleties:

(a) Mass loss Equation (8.29) shows that the encounter redistributes a portion of the system's original energy stock: stars in which $\mathbf{v}_{\alpha} \cdot \Delta\mathbf{v}_{\alpha} > 0$ gain energy, while those with $\mathbf{v}_{\alpha} \cdot \Delta\mathbf{v}_{\alpha} < 0$ may lose energy. The energy gained by some stars may be so large that they escape from the system, and then the overall change in energy of the stars that remain bound can be negative. Thus the energy per unit mass of the bound remnant system may decrease (become more negative) as the result of the encounter, even though the encounter always adds energy to the original system.

(b) Return to equilibrium After the increments (8.28) have been added to the velocities of all the stars of the subject system, it no longer satisfies the virial theorem (4.250). Hence the encounter initiates a period of readjustment, lasting a few crossing times, during which the subject system settles to a new equilibrium configuration.

If the perturbation is weak enough that no stars escape, some properties of this new equilibrium can be deduced using the virial theorem. Let the initial internal kinetic and total energies be \tilde{K}_0 and \tilde{E}_0 , respectively. Then the virial theorem implies that

$$\tilde{K}_0 = -\tilde{E}_0.\tag{8.31}$$

Since the impulsive encounter increases the kinetic energy by $\Delta\tilde{K}$ and leaves the potential energy unchanged, the final energy is

$$\tilde{E}_1 = \tilde{E}_0 + \Delta\tilde{K}.\tag{8.32}$$

Once the subject system has settled to a new equilibrium state, the final kinetic energy is given by the virial theorem,

$$\tilde{K}_1 = -\tilde{E}_1 = -(\tilde{E}_0 + \Delta\tilde{K}) = \tilde{K}_0 - \Delta\tilde{K}. \quad (8.33)$$

Thus if the impulsive encounter *increases* the kinetic energy by $\Delta\tilde{K}$, the subsequent relaxation back to dynamical equilibrium *decreases* the kinetic energy by $2\Delta\tilde{K}$!

(c) Adiabatic invariance The impulse approximation is valid only if the encounter time is short compared to the crossing time. In most stellar systems the crossing time is a strong function of energy or mean orbital radius, so the impulse approximation is unlikely to hold for stars near the center. Indeed, sufficiently close to the center, the crossing times of most stars may be so short that their orbits deform adiabatically as the perturber approaches (§3.6.2c). In this case, changes that occur in the structure of the orbits as the perturber approaches will be reversed as it departs, and the encounter will leave most orbits in the central region unchanged.

If we approximate the potential near the center of the stellar system as that of a harmonic oscillator with frequency Ω , then the energy change imparted to the stars in an encounter of duration τ is proportional to $\exp(-\alpha\Omega\tau)$ for $\Omega\tau \gg 1$, where α is a constant of order unity (see §3.6.2a). However, this strong exponential dependence does not generally hold in realistic stellar systems. The reason is that some of the stars are in resonance with the slowly varying perturbing force, in the sense that $\mathbf{m} \cdot \boldsymbol{\Omega} \simeq 0$ where the components of $\boldsymbol{\Omega}(\mathbf{J})$ are the fundamental frequencies of the orbit (eq. 3.190), and \mathbf{m} is an integer triple. At such a resonance, the response to a slow external perturbation is large—in the language of Chapter 5, the polarization matrix diverges (eq. 5.95). A careful calculation of the contribution of both resonant and non-resonant stars shows that the total energy change in an encounter generally declines only as $(\Omega\tau)^{-1}$ for $\Omega\tau \gg 1$, rather than exponentially (Weinberg 1994a).

8.2.1 The distant-tide approximation

The calculation of the effects of an encounter simplifies considerably when the size of the subject system is much less than the impact parameter.

Let $\Phi(\mathbf{x}, t)$ be the gravitational potential of the perturber, in a frame in which the center of mass of the subject system is at the origin. When the distance to the perturber is much larger than the size of the subject system, the perturbing potential will vary smoothly across it, and we may therefore expand the field $-\nabla\Phi(\mathbf{x}, t)$ in a Taylor series about the origin:

$$-\frac{\partial\Phi}{\partial x_j}(\mathbf{x}, t) = -\Phi_j(t) - \sum_k \Phi_{jk}(t)x_k + O(|\mathbf{x}|^2), \quad (8.34a)$$

where $\mathbf{x} = (x_1, x_2, x_3)$ and

$$\Phi_j \equiv \left. \frac{\partial \Phi}{\partial x_j} \right|_{\mathbf{x}=0} ; \quad \Phi_{jk} \equiv \left. \frac{\partial^2 \Phi}{\partial x_j \partial x_k} \right|_{\mathbf{x}=0}. \quad (8.34b)$$

Dropping the terms $O(|\mathbf{x}|^2)$ constitutes the **distant-tide approximation**. Encounters for which the both the distant-tide and impulse approximations are valid are often called **tidal shocks**.

Substituting into equations (8.25a) and (8.26), we find that Φ_{jk} does not contribute to the center-of-mass acceleration $\dot{\mathbf{v}}_{\text{cm}}$, because the center of mass is at the origin so $\sum_{\beta} m_{\beta} \mathbf{x}_{\beta} = 0$. Similarly, substituting (8.34a) into (8.27), we find that Φ_j does not contribute to $\dot{\mathbf{v}}_{\alpha}$ because $\sum_{\beta} m_{\beta} = M_s$. Thus

$$\dot{\mathbf{v}}_{\alpha} = - \sum_{j,k=1}^3 \hat{\mathbf{e}}_j \Phi_{jk} x_{\alpha k}. \quad (8.35)$$

If the perturber is spherical and centered at $\mathbf{X}(t)$, we may write $\Phi(\mathbf{x}, t) = \Phi(|\mathbf{x} - \mathbf{X}(t)|)$ and (cf. Box 2.3)

$$\Phi_j = -\Phi' \frac{X_j}{X} ; \quad \Phi_{jk} = \left(\Phi'' - \frac{\Phi'}{X} \right) \frac{X_j X_k}{X^2} + \frac{\Phi'}{X} \delta_{jk}, \quad (8.36)$$

where $X = |\mathbf{X}|$ and all derivatives of Φ are evaluated at X .

An important special case occurs when the impact parameter is large enough that we may approximate the perturber as a point mass M_p . Then $\Phi(X) = -GM_p/X$ and we have

$$\Phi_j = -\frac{GM_p}{X^3} X_j ; \quad \Phi_{jk} = \frac{GM_p}{X^3} \delta_{jk} - \frac{3GM_p}{X^5} X_j X_k. \quad (8.37)$$

Thus the equation of motion (8.35) becomes

$$\dot{\mathbf{v}}_{\alpha} = -\frac{GM_p}{X^3} \mathbf{x}_{\alpha} + \frac{3GM_p}{X^5} (\mathbf{X} \cdot \mathbf{x}_{\alpha}) \mathbf{X}. \quad (8.38)$$

We argued at the beginning of this section that in the impulse approximation, the orbit of the perturber can be assumed to have constant relative velocity \mathbf{V} . We align our coordinate axes so that \mathbf{V} lies along the z axis, and the perturber's orbit lies in the yz plane, and choose the origin of time to coincide with the point of closest approach. Then $\mathbf{X}(t) = (0, b, Vt)$, where b is the impact parameter, and we have

$$\dot{\mathbf{v}}_{\alpha} = -\frac{GM_p \mathbf{x}_{\alpha}}{[b^2 + (Vt)^2]^{3/2}} + \frac{3GM_p (y_{\alpha} b + z_{\alpha} Vt)}{[b^2 + (Vt)^2]^{5/2}} (b \hat{\mathbf{e}}_y + Vt \hat{\mathbf{e}}_z). \quad (8.39)$$

In the impulse approximation, \mathbf{x}_α is constant during the encounter, so

$$\begin{aligned}\Delta \mathbf{v}_\alpha &= \int_{-\infty}^{\infty} dt \dot{\mathbf{v}}_\alpha \\ &= GM_p \int_{-\infty}^{\infty} dt \left\{ -\frac{(x_\alpha, y_\alpha, z_\alpha)}{[b^2 + (Vt)^2]^{3/2}} + 3(0, b, Vt) \frac{y_\alpha b + z_\alpha Vt}{[b^2 + (Vt)^2]^{5/2}} \right\} \\ &= \frac{GM_p}{b^2 V} \left(-x_\alpha \int_{-\infty}^{\infty} \frac{du}{(1+u^2)^{3/2}}, y_\alpha \int_{-\infty}^{\infty} du \frac{2-u^2}{(1+u^2)^{5/2}}, \right. \\ &\quad \left. z_\alpha \int_{-\infty}^{\infty} du \frac{2u^2-1}{(1+u^2)^{5/2}} \right),\end{aligned}\tag{8.40}$$

where we have made the substitution $u = Vt/b$. Evaluating the integrals in equation (8.40), we obtain finally

$$\Delta \mathbf{v}_\alpha = \frac{2GM_p}{b^2 V} (-x_\alpha, y_\alpha, 0).\tag{8.41}$$

The error introduced in this formula by the distant-tide approximation is of order $|\mathbf{x}|/b \ll 1$. The velocity increments tend to deform a sphere of stars into an ellipsoid whose long axis lies in the direction of the perturber's point of closest approach. This distortion is reminiscent of the way in which the Moon raises tides on the surface of the oceans.

By equations (8.30) and (8.41) the change in the energy per unit mass in the distant-tide approximation is (Spitzer 1958)

$$\Delta \tilde{E} = \frac{2G^2 M_p^2}{V^2 b^4} \sum_{\alpha} m_{\alpha} (x_{\alpha}^2 + y_{\alpha}^2).\tag{8.42}$$

If the subject system is spherical, then $\sum m_{\alpha} x_{\alpha}^2 = \sum m_{\alpha} y_{\alpha}^2 = \frac{1}{3} M_s \langle r^2 \rangle$, where $\langle r^2 \rangle$ is the mass-weighted mean-square radius of the stars in the subject system. In this case equation (8.42) simplifies to

$$\Delta \tilde{E} = \frac{4G^2 M_p^2 M_s}{3V^2 b^4} \langle r^2 \rangle.\tag{8.43}$$

Equation (8.43) shows that for large impact parameter b the energy input in tidal shocks varies as b^{-4} . Thus the encounters that have the strongest effect on a stellar system are those with the smallest impact parameter b , which unfortunately are also those for which the approximation of a point-mass perturber is invalid. Fortunately, it is a straightforward numerical task to generalize these calculations to a spherical perturber with an arbitrary mass distribution, using equations (8.30), (8.35), and (8.36) (Aguilar & White 1985; Gnedin, Hernquist, & Ostriker 1999). Let $U(b/r_h)$ be the ratio

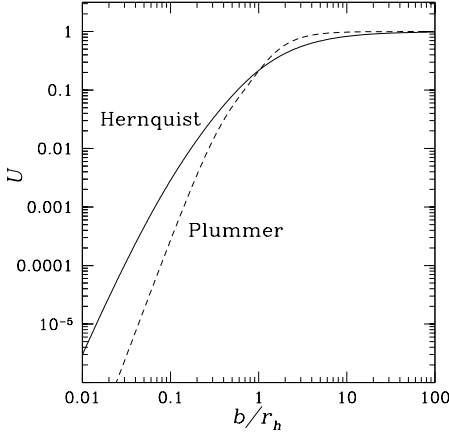


Figure 8.4 Energy input in a tidal shock due to a perturber with a Plummer or Hernquist mass distribution (eqs. 2.44b and 2.67). Here b is the impact parameter, r_h is the half-mass radius of the Plummer or Hernquist model, and U is the ratio of the energy input to that caused by a point mass perturber (eq. 8.44). The integral $W = \int dx U(x)/x^3 = 0.5675$ for the Plummer model and 1.239 for the Hernquist model (eq. 8.52).

of the impulsive energy change caused by a perturber of half-mass radius r_h to the input from a point of the same total mass, which is given by (8.43). Then we have

$$\Delta \tilde{E} = \frac{4G^2 M_p^2 M_s}{3V^2 b^4} U(b/r_h) \langle r^2 \rangle. \quad (8.44)$$

Figure 8.4 shows $U(x)$ for the Plummer and Hernquist mass distributions.

8.2.2 Disruption of stellar systems by high-speed encounters

In many cases we are interested in the cumulative effect of encounters on a stellar system that is traveling through a sea of perturbers. Let us assume that the perturbers have mass M_p and number density n_p , and a Maxwellian DF with velocity dispersion σ_p in one dimension. Similarly, we assume that the subject system is a randomly chosen member of a population having a Maxwellian velocity distribution with dispersion σ_s .

Consider the rate at which the subject system encounters perturbers at relative speeds in the range $(V, V + dV)$ and impact parameters in the range $(b, b + db)$. With our assumption of Maxwellian velocity distributions, the distribution of relative velocities of encounters is also Maxwellian, with dispersion (Problem 8.8)

$$\sigma_{\text{rel}} = (\sigma_s^2 + \sigma_p^2)^{1/2}. \quad (8.45)$$

Thus the probability that the subject system and perturber have relative speed in the given range is

$$dP = \frac{4\pi V^2 dV}{(2\pi\sigma_{\text{rel}}^2)^{3/2}} \exp\left(-\frac{V^2}{2\sigma_{\text{rel}}^2}\right), \quad (8.46)$$

and the average rate at which a subject system encounters perturbers with speed V and impact parameter b is

$$\dot{C} = n_p V 2\pi b db dP = \frac{2\sqrt{2\pi}n_p b db}{\sigma_{\text{rel}}^3} \exp\left(-\frac{V^2}{2\sigma_{\text{rel}}^2}\right) V^3 dV. \quad (8.47)$$

In the distant-tide approximation, the energy input to a star is proportional to the square of its radius (eq. 8.42). So we focus on stars in the outer parts of the subject system, for which we can assume that the gravitational potential of the subject system is Keplerian. In a Keplerian potential, the time averaged mean-square radius of an orbit with semi-major axis a and eccentricity e is $(1 + \frac{3}{2}e^2)a^2$ (Problem 3.9). If the DF of the subject system is isotropic in velocity space, the average of e^2 over all the stars with a given semi-major axis or energy is $\frac{1}{2}$ (Problem 4.8). Thus, if we average over stars with different orbital phases and eccentricities but the same energy, $\langle r^2 \rangle = \frac{7}{4}a^2$, and equation (8.44) yields an average change in energy per unit mass of

$$\langle \Delta E \rangle = \frac{\langle \Delta \tilde{E} \rangle}{M_s} = \frac{7G^2 M_p^2 a^2}{3V^2 b^4} U(b/r_h). \quad (8.48)$$

From Figure 8.4 we see that for $b \gg r_h$, $U \simeq 1$ so $\langle \Delta E \rangle \propto b^{-4}$ is a steeply declining function of impact parameter. The frequency of encounters with impact parameters in the range $(b, b+db)$ is proportional to $b db$ so the rate at which energy is injected by encounters in this range decreases with increasing b as db/b^3 . On the other hand, U rapidly decreases with decreasing b once $b \lesssim r_h$, with the result that encounters with impact parameters $b \sim r_h$ inflict the most damage.

If the damage from a single encounter with $b \sim r_h$ is not fatal for the system, we say that we are in the **diffusive regime** because the effects from a whole sequence of encounters will accumulate, as in the diffusive relaxation processes that we discussed in Chapter 7. If, by contrast, a single encounter at impact parameter $b \sim r_h$ will shatter the system, the damage sustained by the system will be small until it is disrupted by a single, closest encounter, and we say that we are in the **catastrophic regime**.

(a) The catastrophic regime We first determine the largest impact parameter parameter b_1 at which a single encounter can disrupt the system. Since we are in the catastrophic regime, we may assume that $b_1 \gtrsim r_h$ so the energy per unit mass injected by an encounter at impact parameter b_1 is given by equation (8.48) with $U(b/r_h) \simeq 1$. Equating this to the absolute value of the energy of an individual star $E = -GM_s/2a$ (eq. 3.32), we obtain

$$1 = \frac{\langle \Delta E \rangle}{|E|} = \frac{14GM_p^2 a^3}{3M_s V^2 b_1^4} \quad \text{so} \quad b_1(V) = 1.5 \left(\frac{GM_p^2 a^3}{M_s V^2} \right)^{1/4}. \quad (8.49)$$

The rate at which disruptive encounters occur is then given by equation (8.47):

$$\begin{aligned}\mathcal{R} &\equiv \frac{2\sqrt{2\pi}n_p}{\sigma_{\text{rel}}^3} \int_0^\infty dV V^3 \exp\left(-\frac{V^2}{2\sigma_{\text{rel}}^2}\right) \int_0^{b_1(V)} b db \\ &= \sqrt{\frac{14}{3}} \pi \frac{G^{1/2} n_p M_p a^{3/2}}{M_s^{1/2}}.\end{aligned}\quad (8.50)$$

The disruption time of a subject system with semi-major axis a is

$$t_d \simeq \mathcal{R}^{-1} \simeq k_{\text{cat}} \frac{1}{G\rho_p} \left(\frac{GM_s}{a^3}\right)^{1/2}, \quad (8.51)$$

where $\rho_p \equiv M_p n_p$ is the mass density of perturbers and k_{cat} is of order unity. Our analytic treatment yields $k_{\text{cat}} = 0.15$ but Monte-Carlo simulations of catastrophic disruption suggest that $k_{\text{cat}} \simeq 0.07$ (Bahcall, Hut, & Tremaine 1985). It is remarkable that the disruption time in the catastrophic regime is independent of both the velocity dispersion σ_{rel} and the mass of individual perturbers, so long as their overall mass density ρ_p is fixed.

(b) The diffusive regime In this regime, each encounter imparts a velocity impulse $\Delta \mathbf{v}$ (eq. 8.28) that satisfies $|\Delta \mathbf{v}| \ll |\mathbf{v}|$. The corresponding change in the energy per unit mass is $\Delta E = \mathbf{v} \cdot \Delta \mathbf{v} + \frac{1}{2} |\Delta \mathbf{v}|^2$. Thus the **diffusion term** $\mathbf{v} \cdot \Delta \mathbf{v}$ is much larger than the **heating term** $\frac{1}{2} |\Delta \mathbf{v}|^2$. On the other hand the direction of the velocity impulse, which depends on the relative orientation of the star and the perturber, is usually uncorrelated with the direction of the velocity \mathbf{v} of the star relative to the center of mass of the subject system, which depends on the orbital phase of the star. Thus the average of the diffusion term over many encounters is zero, while the heating term systematically increases the energy.⁷ An accurate description of the evolution of the energy under the influence of many high-speed encounters requires the inclusion of both terms, using the Fokker–Planck equation that we described in §7.4.2. Nevertheless, for the sake of simplicity, and since our estimates will be crude anyway, we focus our attention exclusively on the heating term.

Combining equations (8.47) and (8.48), we find that the rate of energy increase for stars with semi-major axis a is

$$\begin{aligned}\dot{E} &= \dot{C} \langle \Delta E \rangle \\ &= \frac{14}{3} \sqrt{2\pi} \frac{G^2 M_p^2 n_p a^2}{\sigma_{\text{rel}}^3} \int_0^\infty dV V \exp\left(-\frac{V^2}{2\sigma_{\text{rel}}^2}\right) \int \frac{db}{b^3} U(b/r_h) \\ &= \frac{14}{3} \sqrt{2\pi} \frac{G^2 M_p^2 n_p a^2}{\sigma_{\text{rel}}^2 r_h^2} W, \quad \text{where } W \equiv \int \frac{dx}{x^3} U(x).\end{aligned}\quad (8.52)$$

⁷This argument is similar to, but distinct from, the argument leading from equation (8.29) to equation (8.30), which involved an average over the effects of a single collision on many stars rather than an average over the effect of many collisions on a single star.

In general, W must be evaluated numerically for a given mass model. For a Plummer model, $W = 0.5675$ and for a Hernquist model $W = 1.239$ (Figure 8.4).

For point-mass perturbers, $U(x) = 1$, and the heating rate is

$$\dot{E} = \frac{14}{3} \sqrt{2\pi} \frac{G^2 M_p^2 n_p a^2}{\sigma_{\text{rel}}} \int \frac{db}{b^3}. \quad (8.53)$$

This integral over impact parameter diverges at small b . In practice, the distant-tide approximation fails when the impact parameter is comparable to the size of the subject system, so the integration should be cut off at this point.

Comparing the heating rate (8.52) to the energy of an individual star $E = -\frac{1}{2}GM_s/a$, we obtain the time required for the star to escape:

$$t_d \simeq \frac{|E|}{\dot{E}} \simeq \frac{0.043}{W} \frac{\sigma_{\text{rel}} M_s r_h^2}{GM_p^2 n_p a^3}. \quad (8.54)$$

For point-mass perturbers, we use equation (8.53), with the integration over impact parameter cut off at b_{min} :

$$t_d \simeq \frac{|E|}{\dot{E}} \simeq 0.085 \frac{\sigma_{\text{rel}} M_s b_{\text{min}}^2}{GM_p^2 n_p a^3}. \quad (8.55)$$

These are only approximate estimates. A more accurate treatment would employ the Fokker–Planck equation (7.123); in this equation the diffusion coefficient $D[\Delta E]$ is the quantity here called \dot{E} , and the diffusion coefficient $D[(\Delta E)^2]$ would be computed similarly as the rate of change of the mean-square energy. Generally this treatment gives a half-life for a star with a given semi-major axis that is a few times shorter than the estimate (8.54).

(c) Disruption of open clusters The masses of open clusters lie in the range $10^2 \mathcal{M}_\odot \lesssim M_c \lesssim 10^4 \mathcal{M}_\odot$, and their half-mass radii and internal velocity dispersions are $r_{h,c} \approx 2 \text{ pc}$ and $\sigma_c \approx 0.3 \text{ km s}^{-1}$ (Table 1.3). The crossing time at the half-mass radius is $r_{h,c}/\sigma_c \approx 10 \text{ Myr}$. Much of the interstellar gas in our Galaxy is concentrated into a few thousand **giant molecular clouds** of mass $M_{\text{GMC}} \gtrsim 10^5 \mathcal{M}_\odot$ and radius $r_{h,\text{GMC}} \approx 10 \text{ pc}$. Both open clusters and molecular clouds travel on nearly circular orbits through the Galactic disk, with random velocities of order 7 km s^{-1} ; thus the dispersion in relative velocity is $\sigma_{\text{rel}} \simeq \sqrt{2} \times 7 \text{ km s}^{-1} \simeq 10 \text{ km s}^{-1}$ (eq. 8.45). The duration of a cluster-cloud encounter with impact parameter $b > r_{\text{GMC}}$ is then $b/\sigma_{\text{rel}} \simeq (b/10 \text{ pc}) \text{ Myr}$, which is shorter than the cluster crossing time for $b \lesssim 100 \text{ pc}$. Thus we may use the impulse approximation to study the effect of close encounters with molecular clouds on open clusters.

The impact parameter at which a typical encounter with a point-mass perturber would disrupt the cluster is given by equation (8.49); identifying

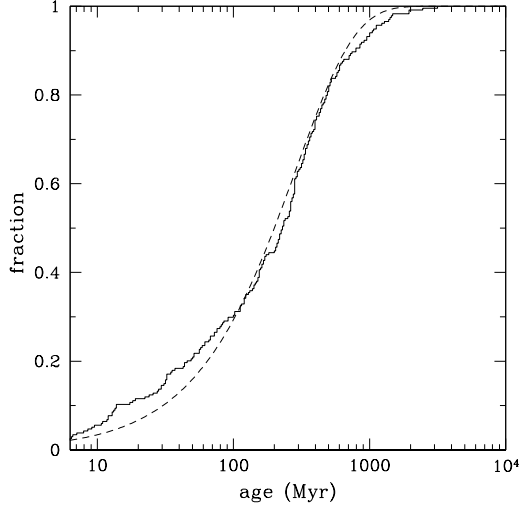


Figure 8.5 The fraction of nearby open clusters younger than a given age. The cluster sample is from Piskunov et al. (2007). The curve is derived from a simple theoretical model in which clusters are born at a constant rate and the probability that a cluster survives for time t is $\exp(-t/\tau)$ with $\tau \simeq 300$ Myr. A Kolmogorov–Smirnov test (Press et al. 1986) shows that the two distributions are statistically indistinguishable.

the open cluster with the subject system and the molecular cloud with the perturber we obtain

$$b_1(\sigma_{\text{rel}}) = 15 \text{ pc} \left(\frac{M_{\text{GMC}}}{10^5 \mathcal{M}_{\odot}} \right)^{1/2} \left(\frac{300 \mathcal{M}_{\odot}}{M_c} \right)^{1/4} \times \left(\frac{a}{2 \text{ pc}} \right)^{3/4} \left(\frac{10 \text{ km s}^{-1}}{\sigma_{\text{rel}}} \right)^{1/2}. \quad (8.56)$$

Since this distance is larger than the cloud size $r_{\text{h,GMC}} \approx 10 \text{ pc}$, even when the semi-major axis is as small as the typical cluster half-mass radius of 2 pc, the encounters are in the catastrophic regime. Hence the disruption time is given by equation (8.51):

$$t_d \simeq 250 \text{ Myr} \frac{k_{\text{cat}}}{0.07} \frac{0.025 \mathcal{M}_{\odot} \text{ pc}^{-3}}{\rho_{\text{GMC}}} \left(\frac{M_c}{300 \mathcal{M}_{\odot}} \right)^{1/2} \left(\frac{2 \text{ pc}}{a} \right)^{3/2}, \quad (8.57)$$

where the mean density of gas in molecular clouds is taken to be about half of the total gas density in the solar neighborhood (see Table 1.1).

This result is quite uncertain, not only because the derivation of equation (8.57) is highly idealized, but also because of uncertainties in the molecular cloud parameters and the large dispersion in open-cluster parameters. Nevertheless, the available data suggest that the median lifetime of open clusters is remarkably close to this simple estimate (Figure 8.5). In fact, it was the observation that there are few open clusters with ages $\gtrsim 500$ Myr that prompted Spitzer (1958) to argue that clusters might be dissolved by the very clouds that bring them into the world.

(d) Disruption of binary stars Binary stars can be thought of as clusters with just two members and, like clusters, they can be disrupted by

encounters with passing perturbers. Obviously the vulnerability of a binary to disruption is an increasing function of the semi-major axis a of its components. Binary semi-major axes are usually measured in terms of the **astronomical unit**, $1 \text{ AU} = 1.496 \times 10^{11} \text{ m} = 4.848 \times 10^{-6} \text{ pc}$ (approximately the mean Earth-Sun distance; see Appendix A).

First we consider disruption of binaries in the solar neighborhood by passing stars. We focus on stars—both the binary components and their perturbers—that have ages comparable to the age of the Galaxy and masses comparable to that of the Sun, since these contain most of the stellar mass in the solar neighborhood. The velocity distribution of such stars is triaxial, but we may approximate this distribution by an isotropic Maxwellian with a one-dimensional dispersion $\sigma_* \simeq 30 \text{ km s}^{-1}$ ($1/\sqrt{3}$ of the RMS velocity, from Table 1.2). The velocity distribution is the same for single and binary stars, so the relative dispersion is $\sigma_{\text{rel}} = \sqrt{2} \times 30 \text{ km s}^{-1} \simeq 40 \text{ km s}^{-1}$ (eq. 8.45).

According to equation (8.49), the maximum impact parameter for a catastrophic encounter between a binary star of total mass M_b and a passing star of mass M_p is

$$\begin{aligned} b_1(\sigma_{\text{rel}}) &\simeq 1.5a \left(\frac{GM_p^2}{M_b \sigma_{\text{rel}}^2 a} \right)^{1/4} \\ &\simeq 0.11a \left(\frac{2 \mathcal{M}_\odot}{M_b} \frac{10^4 \text{ AU}}{a} \right)^{1/4} \left(\frac{M_p}{1 \mathcal{M}_\odot} \frac{40 \text{ km s}^{-1}}{\sigma_{\text{rel}}} \right)^{1/2}. \end{aligned} \quad (8.58)$$

Unless the semi-major axis is so small that the probability of a close encounter is negligible, this result shows that $b_1 \lesssim a$ for solar-type stars in the solar neighborhood, and thus that the encounters are in the diffusive regime. The disruption time is given by equation (8.55), setting $b_{\text{min}} \sim a$, where the distant-tide approximation fails; thus (Öpik 1932; Heggie 1975)

$$t_d \simeq k_{\text{diff}} \frac{\sigma_{\text{rel}} M_b}{GM_p^2 n_p a}, \quad (8.59)$$

where $k_{\text{diff}} \equiv 0.085(b_{\text{min}}/a)^2$. We can refine this estimate by recalling the discussion of the disruption of soft binaries in §7.5.7a; equation (7.173) in that section describes the disruption time in the diffusive regime for the case in which the component stars of the binary have the same mass as the perturbing stars, so $M_b = 2M_p$, and the velocity dispersion σ of the perturbers and the binaries is the same, so $\sigma_{\text{rel}} = \sqrt{2}\sigma$. Equating the two expressions, we find $k_{\text{diff}} \simeq 0.022/\ln \Lambda$, where $\Lambda \approx \sigma_{\text{rel}}^2 a / (GM_p)$. For binaries with $a \sim 10^4 \text{ AU}$ in the solar neighborhood, this formula yields $k_{\text{diff}} \simeq 0.002$, and Monte Carlo simulations yield a similar value (Bahcall, Hut, & Tremaine 1985). The rather small value of k_{diff} arises in part because \dot{E} grows as a^2 , so the heating rate accelerates as the binary gains energy, and in part

because close encounters with either member of the binary contribute to the disruption rate, an effect not accounted for in equation (8.53).

In the solar neighborhood, equation (8.59) yields

$$t_d \simeq 15 \text{ Gyr} \frac{k_{\text{diff}}}{0.002} \frac{\sigma_{\text{rel}}}{40 \text{ km s}^{-1}} \frac{M_b}{2 \mathcal{M}_\odot} \left(\frac{1 \mathcal{M}_\odot}{M_p} \right)^2 \frac{0.05 \text{ pc}^{-3}}{n_p} \frac{10^4 \text{ AU}}{a}. \quad (8.60)$$

Thus the upper limit to the semi-major axes of old binary stars in the solar neighborhood is $a \simeq 2 \times 10^4 \text{ AU}$.

Now consider the effects of molecular clouds. Replacing the perturber mass M_p in equation (8.58) by the typical cloud mass $M_{\text{GMC}} \approx 10^5 \mathcal{M}_\odot$, we find that the maximum impact parameter for impulsive disruption is

$$b_1(\sigma_{\text{rel}}) \simeq 1.9 \text{ pc} \left(\frac{2 \mathcal{M}_\odot}{M_b} \right)^{1/4} \left(\frac{a}{10^4 \text{ AU}} \right)^{3/4} \left(\frac{M_{\text{GMC}}}{10^5 \mathcal{M}_\odot} \frac{30 \text{ km s}^{-1}}{\sigma_{\text{rel}}} \right)^{1/2}. \quad (8.61)$$

We have used a fiducial value $\sigma_{\text{rel}} = 30 \text{ km s}^{-1}$, which is the sum in quadrature of the dispersions of the stars, $\sigma_* \simeq 30 \text{ km s}^{-1}$, and the clouds, $\sigma_{\text{GMC}} \simeq 7 \text{ km s}^{-1}$. Since b_1 is smaller than the cloud radius $r_{\text{h,GMC}} \approx 10 \text{ pc}$, the encounters are in the diffusive regime. The disruption time can be estimated from equation (8.54), using the value $W = 0.5675$ appropriate for a Plummer model of the cloud's density distribution:

$$t_d \simeq 0.075 \frac{\sigma_{\text{rel}} M_b r_{\text{h,GMC}}^2}{G M_{\text{GMC}}^2 n_{\text{GMC}} a^3}. \quad (8.62)$$

The cloud parameters M_{GMC} , n_{GMC} , and $r_{\text{h,GMC}}$ are all poorly known. Fortunately, they enter this equation in terms of the observationally accessible combinations $\Sigma_{\text{GMC}} \equiv (M/\pi r_{\text{h}}^2)_{\text{GMC}}$, the mean surface density of a cloud, and $\rho_{\text{GMC}} = (Mn)_{\text{GMC}}$, the mean density of molecular gas. We adopt $\Sigma_{\text{GMC}} \simeq 300 \mathcal{M}_\odot \text{ pc}^{-2}$ and $\rho_{\text{GMC}} \simeq 0.025 \mathcal{M}_\odot \text{ pc}^{-3}$ (Hut & Tremaine 1985). Thus

$$t_d \simeq 380 \text{ Gyr} \frac{M_b}{2 \mathcal{M}_\odot} \left(\frac{10^4 \text{ AU}}{a} \right)^3 \frac{\sigma_{\text{rel}}}{30 \text{ km s}^{-1}}. \quad (8.63)$$

Although this result is subject to substantial uncertainties, together with equation (8.60) it implies that binaries with semi-major axes $\gtrsim 2 \times 10^4 \text{ AU} \simeq 0.1 \text{ pc}$ cannot survive in the solar neighborhood for its lifetime of $\sim 10 \text{ Gyr}$, due to the combined effects of high-speed encounters with molecular clouds and other stars.

The widest known binary stars in the disk do indeed have separations of about 0.1 pc (Chanamé & Gould 2004); however, there is little evidence for or against the cutoff in the binary distribution that we have predicted at this separation (Wasserman & Weinberg 1987). Binary stars in the stellar halo appear to exist with even larger separations; such binaries can survive

because their velocity σ_{rel} relative to the disk is much higher, and because they spend only a fraction of their orbit in the disk, so the disruptive effects from disk stars and molecular clouds are much weaker (Yoo, Chanamé, & Gould 2004).

(e) Dynamical constraints on MACHOs One possible constituent of the dark halo is MACHOs, compact objects such as black holes or non-luminous stars (§1.1.2). Suppose that MACHOs contribute a fraction $f_{\text{h}} \gtrsim 0.5$ of the radial force in the solar neighborhood; this is close to the maximum allowed since the disk contributes a fraction $f_{\text{d}} = 1 - f_{\text{h}} \gtrsim 0.4$ (§6.3.3). Then limits on the optical depth of the dark halo to gravitational lensing (Alcock et al. 2001; Tisserand et al. 2007) imply that the MACHO mass

$$m \lesssim 10^{-7} \mathcal{M}_{\odot} \quad \text{or} \quad m \gtrsim 30 \mathcal{M}_{\odot}. \quad (8.64)$$

In §7.4.4 we showed that encounters between MACHOs and disk stars add kinetic energy to the disk stars and thereby increase both the velocity dispersion and the disk thickness; even if this is the only mechanism that heats the disk—and we shall see in §8.4 that it is not—the observed dispersion requires that $m \lesssim (5\text{--}10) \times 10^6 \mathcal{M}_{\odot}$ (eq. 7.104). We now investigate what additional constraints can be placed on the MACHO mass by the effect of high-speed encounters of MACHOs on binary stars.

We write the number density of MACHOs as $n = \rho/m$. If the MACHO mass is small enough, disruption is in the diffusive regime, and we can use equation (8.59) to estimate the disruption time:

$$t_{\text{d,diff}} \simeq k_{\text{diff}} \frac{\sigma_{\text{rel}} M_{\text{b}}}{G m \rho a} \quad (k_{\text{diff}} \approx 0.002), \quad (8.65a)$$

where M_{b} is the mass of the binary. In the catastrophic regime, the disruption time is given by equation (8.51):

$$t_{\text{d,cat}} \simeq k_{\text{cat}} \frac{M_{\text{b}}^{1/2}}{G^{1/2} \rho a^{3/2}} \quad (k_{\text{cat}} \approx 0.07). \quad (8.65b)$$

The transition between these two regimes occurs when the critical impact parameter $b_1(\sigma_{\text{rel}})$ (eq. 8.49) is of order the binary semi-major axis a ; however, a more accurate way to determine the transition is to set the actual disruption time to

$$t_{\text{d}} = \min(t_{\text{d,diff}}, t_{\text{d,cat}}) \quad (8.66)$$

and identify the transition with the MACHO mass m_{crit} at which $t_{\text{d,diff}} = t_{\text{d,cat}}$. Thus we find

$$m_{\text{crit}} = \frac{k_{\text{diff}}}{k_{\text{cat}}} \left(\frac{\sigma_{\text{rel}}^2 M_{\text{b}} a}{G} \right)^{1/2} \quad (k_{\text{diff}} k_{\text{cat}} \approx 0.03). \quad (8.67)$$

Notice that for $m > m_{\text{crit}}$, the disruption time $t_{\text{d,cat}}$ depends on the overall density contributed by the MACHOs but not their individual masses. Thus the survival of a given type of binary system either rules out *all* MACHO masses above m_{crit} and *some* masses below m_{crit} (if the system's age exceeds $t_{\text{d,cat}}$) or does not rule out *any* masses (if its age is less than $t_{\text{d,cat}}$).

To plug in numbers for the solar neighborhood, we use the simple model for the DF of MACHOs in the dark halo that we described on page 584. In this model the local density of MACHOs is given by equation (7.94), and the relative dispersion between MACHOs is $\sigma_{\text{rel}} = \sqrt{2}\sigma = v_c$, where v_c is the circular speed (eq. 8.45)—this is also roughly the dispersion between the MACHOs and stars, whether they belong to the disk or the stellar halo. Then

$$m_{\text{crit}} \simeq 30 \mathcal{M}_{\odot} \frac{k_{\text{diff}}/k_{\text{cat}}}{0.03} \left(\frac{M_b}{2 \mathcal{M}_{\odot}} \frac{a}{10^4 \text{ AU}} \right)^{1/2} \frac{v_c}{220 \text{ km s}^{-1}}. \quad (8.68)$$

To evaluate the disruption time in the catastrophic regime, $m > m_{\text{crit}}$, we use equation (8.65b), and take the local MACHO density from equation (7.94). Assuming the solar radius $R_0 = 8 \text{ kpc}$ and the solar circular speed $v_c = v_0 = 220 \text{ km s}^{-1}$, we have

$$t_{\text{d,cat}} \simeq 20 \text{ Gyr} \frac{0.5}{f_h} \frac{k_{\text{cat}}}{0.07} \left(\frac{10^4 \text{ AU}}{a} \right)^{3/2}. \quad (8.69)$$

For dark-halo fractions $f_h \simeq 0.5$, the disruption time $t_{\text{d,cat}}$ is larger than 10 Gyr for semi-major axes $a \lesssim 1.6 \times 10^4 \text{ AU}$. In the diffusive regime, the disruption time is even longer. Disk binaries with semi-major axes larger than this limit are likely to be disrupted by encounters with other disk stars (eq. 8.59) and so we cannot probe the MACHO mass with disk binaries. Halo binaries are much less susceptible to other stars and molecular clouds, because they spend only a small fraction of their time in the disk, and therefore might be present with semi-major axes large enough to provide useful constraints on the MACHO population. Thus, if a population of halo binaries with $a \gtrsim 2 \times 10^4 \text{ AU}$ were discovered, we could rule out a substantial contribution to the local gravitational field for all MACHO masses exceeding $30 \mathcal{M}_{\odot}$ (eq. 8.68). Yoo, Chanamé, & Gould (2004) offer evidence that halo binaries exist with semi-major axes as large as $a \sim 10^5 \text{ AU}$. Together with the microlensing constraint (8.64) this conclusion, if verified by larger samples, would virtually rule out MACHOs as a significant constituent of the dark halo in the solar neighborhood.

(f) Disk and bulge shocks Globular clusters in disk galaxies pass through the disk plane twice per orbit. As they cross the plane, the gravitational field of the disk exerts a compressive gravitational force which is superposed on the cluster's own gravitational field, pinching the cluster briefly along the normal to the disk plane. Repeated pinching at successive passages

through the disk can eventually disrupt the cluster. This process is known as **disk shocking** (Ostriker, Spitzer, & Chevalier 1972).

Let $Z \equiv Z_{\text{cm}} + z$ be the height above the disk midplane of a cluster star, with $Z_{\text{cm}}(t)$ the height of the cluster's center of mass. Then so long as the cluster is small compared to the disk thickness, we may use the distant-tide approximation, and equation (8.35) yields

$$\dot{v}_z = - \left(\frac{\partial^2 \Phi_d}{\partial Z^2} \right)_{\text{cm}} z, \quad (8.70)$$

where $v_z = \dot{z}$ is the z -velocity of the star relative to the cluster center.

The gravitational potential arising from a thin disk of density $\rho_d(R, z)$ is $\Phi_d(R, Z)$, where (eq. 2.74)

$$\frac{d^2 \Phi_d}{dZ^2} = 4\pi G \rho_d. \quad (8.71)$$

Thus

$$\dot{v}_z = -4\pi G \rho_d(R, Z_{\text{cm}}) z, \quad (8.72)$$

where R is the radius at which the cluster crosses the disk.

If the passage of the cluster through the disk is sufficiently fast for the impulse approximation to hold, z is constant during this passage, and the velocity impulse is

$$\Delta v_z = \int dt \dot{v}_z = -4\pi G z \int dt \rho_d[R, Z_{\text{cm}}(t)]. \quad (8.73)$$

To a good approximation we can assume that the velocity of the center of mass of the cluster is constant as it flies through the disk, so $Z_{\text{cm}}(t) = V_z t + \text{constant}$, where V_z is the Z -velocity of the cluster; eliminating the dummy variable t in favor of Z_{cm} we have

$$\Delta v_z = -\frac{4\pi G z}{|V_z|} \int dZ_{\text{cm}} \rho_d(R, Z_{\text{cm}}) = -\frac{4\pi G \Sigma_d(R) z}{|V_z|}, \quad (8.74)$$

where $\Sigma_d(R) \equiv \int dZ \rho_d(R, Z)$ is the surface density of the disk.

From equation (8.30), the energy per unit mass gained by the cluster in a single disk passage is

$$\Delta E = \frac{1}{2} \langle (\Delta v_z)^2 \rangle = \frac{8\pi^2 G^2 \Sigma_d^2}{V_z^2} \langle z^2 \rangle. \quad (8.75)$$

If the cluster is spherically symmetric, the average value of z^2 for stars at a given radius r is $\frac{1}{3}r^2$. As shown on page 662, if the cluster has an ergodic DF

the average value of r^2 for stars with a given semi-major axis is $\frac{7}{4}a^2$. Thus the energy gain is

$$\Delta E = \frac{14\pi^2 G^2 \Sigma_d^2 a^2}{3V_z^2}. \quad (8.76)$$

The cluster passes through the disk twice in each orbital period T_ψ , so the disruption time is

$$t_d \simeq \frac{1}{2} T_\psi \frac{|E|}{\Delta E} = 0.005 \frac{M_{\text{gc}} V_z^2 T_\psi}{G \Sigma_d^2 a^3}, \quad (8.77)$$

where we have set $E = \frac{1}{2} GM_{\text{gc}}/a$, since the potential is Keplerian in the outer parts of the cluster, where the effect of disk shocking is strongest.⁸

In the solar neighborhood the Galactic disk has a midplane volume density $\rho \simeq 0.10 \mathcal{M}_\odot \text{pc}^{-3}$, and surface density $\Sigma_d \simeq 50 \mathcal{M}_\odot \text{pc}^{-2}$ (Table 1.1). The effective thickness of the disk is $h \equiv \Sigma_d/\rho \simeq 500 \text{pc}$. If we approximate the potential of the Milky Way as spherically symmetric, with circular speed v_c at all radii, then the mean-square speed of a collection of test particles such as clusters is $\langle V^2 \rangle = v_c^2$ (Problem 4.35), so if the cluster distribution is spherical, we expect that $\langle V_z^2 \rangle = \frac{1}{3} v_c^2$; thus $\langle V_z^2 \rangle^{1/2} \simeq 130 \text{km s}^{-1}$ for $v_c \simeq 220 \text{km s}^{-1}$. Equation (8.77) can be rewritten

$$t_d \simeq 340 \text{Gyr} \frac{M_{\text{gc}}}{2 \times 10^5 \mathcal{M}_\odot} \frac{T_\psi}{200 \text{Myr}} \times \left(\frac{V_z}{130 \text{km s}^{-1}} \right)^2 \left(\frac{50 \mathcal{M}_\odot \text{pc}^{-2}}{\Sigma_d} \right)^2 \left(\frac{10 \text{pc}}{a} \right)^3. \quad (8.78)$$

This result is based on the impulse approximation, whose validity we must check. The duration of the encounter of the cluster with the disk is $\tau \approx h/V_z \simeq 4 \text{Myr}$ for $V_z \simeq 130 \text{km s}^{-1}$ and $h \simeq 500 \text{pc}$. The crossing time in the outer parts of the cluster is roughly the inverse of the orbital frequency, $(a^3/GM_{\text{gc}})^{1/2} \simeq 1 \text{Myr} (a/10 \text{pc})^{3/2}$ for a cluster mass of $2 \times 10^5 \mathcal{M}_\odot$. Thus the impulse approximation is valid only in the outer parts of the cluster, $a \gtrsim 30 \text{pc}$, and at these semi-major axes the disruption time is $\lesssim 10 \text{Gyr}$. We conclude that disk shocks can lead to substantial erosion of the outermost stars in a typical globular cluster orbiting at the solar radius. For clusters orbiting at smaller radii, disk shocks are even more important, since the orbital time is shorter and the disk surface density is larger.

Bulge shocking is a closely related process. Here the rapidly changing external gravitational field arises as a cluster on a highly eccentric orbit

⁸ A subtle but important assumption in deriving this result is that the orbital phase and thus the value of the height z is uncorrelated between successive disk passages. This assumption is plausible because the interval between successive disk passages is likely to vary considerably for an eccentric, inclined cluster orbit in a realistic disk galaxy potential.

plunges through the bulge of a disk galaxy or the dense center of an elliptical galaxy. In this case the use of the term “shock” is less apt, since the duration of the encounter is not very short compared to the crossing time, even in the outer parts of the cluster, and the impulse approximation is not strictly valid. Nevertheless, the results are similar: the encounters systematically pump energy into the stars in the outer cluster, leading to the escape of stars and the eventual dissolution of the cluster.

The evolution of the globular-cluster population under the influence of disk and bulge shocks is described in §7.5.6.

(g) High-speed interactions in clusters of galaxies The study of galaxies in clusters provides unique insights into galaxy formation and evolution, not only because many dynamical processes are stronger and more obvious in the high-density cluster environment, but also because clusters can be detected at high redshift, enabling the evolution of the galaxy population to be studied directly.

The relative velocities between galaxies in a rich cluster, $\sim 2000 \text{ km s}^{-1}$, are so large that collisions of galaxies last only a few Myr, far less than the crossing time in the galaxy, so they can be treated by the impulse approximation. We model the cluster as a singular isothermal sphere, with density $\rho(r) = \sigma^2/(2\pi Gr^2)$ (eq. 4.103). In clusters the galaxies and dark matter have a similar distribution, so it is reasonable to assume that the ratio $M_\star \equiv \rho/n$ of the mass density to the galaxy number density—in other words the mass per galaxy—is constant. (Note that this is not necessarily the mass *of* the galaxy, since most of the mass is probably spread uniformly through the cluster and is not associated with any individual galaxies.) We focus on galaxies with the characteristic luminosity $L_\star = 2.9 \times 10^{10} L_\odot$ in the R band (eq. 1.18). The mass-to-light ratio in rich clusters is $\Upsilon_R \approx 200\Upsilon_\odot$ (eq. 1.25) so the mass associated with each L_\star galaxy is

$$M_\star = \Upsilon_R L_\star \approx 6 \times 10^{12} \mathcal{M}_\odot. \quad (8.79)$$

Consequently, the number density of L_\star galaxies is

$$n(r) = \frac{\rho(r)}{M_\star} = \frac{\sigma^2}{2\pi GM_\star r^2}. \quad (8.80)$$

Now let us estimate the rate at which a given galaxy encounters other galaxies. The velocity dispersion in clusters is so high that gravitational focusing is negligible, so if galaxies are deemed to collide when their centers come within a collision radius r_{coll} , then from equation (7.194) the collision time is given by

$$\frac{1}{t_{\text{coll}}} = 4\sqrt{\pi}n\sigma r_{\text{coll}}^2 \simeq \frac{2}{\sqrt{\pi}} \frac{\sigma^3 r_{\text{coll}}^2}{GM_\star r^2}. \quad (8.81)$$

After replacing r_{coll} with $2r_{\text{h}}$, where r_{h} is the galaxy radius, equation (8.81) yields

$$t_{\text{coll}} \simeq 0.2 \frac{GM_{\star} r^2}{\sigma^3 r_{\text{h}}^2} \simeq 6 \text{ Gyr} \left(\frac{800 \text{ km s}^{-1}}{\sigma} \right)^3 \left(\frac{20 \text{ kpc}}{r_{\text{h}}} \frac{r}{0.5 \text{ Mpc}} \right)^2 \frac{M_{\star}}{5 \times 10^{12} \mathcal{M}_{\odot}}. \quad (8.82)$$

Thus the stellar component of a galaxy in the central 0.5 Mpc of a rich cluster is likely to have suffered at least one close encounter with another galaxy. What are the consequences of such encounters?

As we have seen, high-speed collisions between galaxies have only a small effect on the distribution of stars, but if both galaxies contain gas disks the gas will suffer a violent collision and be lost from the galaxies. Spitzer & Baade (1951) suggested that collisions might transform spiral galaxies into gas-free lenticular galaxies, thereby explaining the observation that spirals are replaced by lenticulars in high-density environments such as clusters (§1.1.3). An alternative and more likely explanation (Gunn & Gott 1972) is that ram pressure, heating, and other interactions with hot intergalactic gas in the cluster have gradually eroded the gas disks of spiral galaxies (see van Gorkom 2004 for a review).

Isolated galaxies have dark halos that extend to several hundred kpc. In clusters, encounters strip off the outer parts of these halos, so we expect that cluster galaxies will have much smaller and less massive halos than galaxies in low-density environments. We can make a crude estimate of this effect using equation (8.54). For this purpose, we assume that the subject system and the perturber are identical. Thus $M_{\text{s}} = M_{\text{p}} = M$ and $\sigma_{\text{rel}} = \sqrt{2}\sigma$ where σ is the velocity dispersion in the cluster. We set the dimensionless parameter $W \approx 1$, and set $r_{\text{h}} = a$ to estimate the disruption time for stars at the half-mass radius. Thus

$$t_{\text{d}} \simeq 0.06 \frac{\sigma}{GMnr_{\text{h}}}. \quad (8.83)$$

We eliminate the number density n using equation (8.80), and write the mass M of the galaxy in terms of r_{h} and its velocity dispersion $\sigma_{\text{s}}^2 = \frac{1}{3}\langle v^2 \rangle$, using the virial theorem in the form (4.249b); the dispersion for an L_{\star} galaxy is $\sigma_{\text{s}} \simeq 200 \text{ km s}^{-1}$ (eq. 1.21). Thus

$$\begin{aligned} t_{\text{d}} &\simeq 0.056 \frac{GM_{\star} r^2}{\sigma \sigma_{\text{s}}^2 r_{\text{h}}^2} \\ &\simeq 3.8 \text{ Gyr} \frac{M_{\star}}{5 \times 10^{12} \mathcal{M}_{\odot}} \frac{800 \text{ km s}^{-1}}{\sigma} \left(\frac{200 \text{ km s}^{-1}}{\sigma_{\text{s}}} \right)^2 \left(\frac{r}{0.5 \text{ Mpc}} \frac{50 \text{ kpc}}{r_{\text{h}}} \right)^2. \end{aligned} \quad (8.84)$$

Note that the disruption time is related to the collision time (eq. 8.82) by the simple formula $t_{\text{d}}/t_{\text{coll}} \simeq 0.25(\sigma/\sigma_{\text{s}})^2$; as the velocity dispersion σ of the

cluster increases relative to the dispersion σ_s of the galaxies, gravitational interactions become less and less important relative to physical collisions.

This result shows that encounters with other galaxies will erode the dark halos of galaxies residing in the inner 0.5 Mpc of a rich cluster to a radius $r_h \lesssim 50$ kpc. Most of the dark-halo mass has therefore been stripped from the individual galaxies, and is now smoothly distributed throughout the cluster (Richstone 1976). Since the disruption time is proportional to M_s/a^3 (eq. 8.54) and thus to the mean density of the subject system, low-density galaxies are more severely affected by encounters, and can be completely disrupted near the cluster center. The large-scale static tidal field of the cluster also strips the outer halo of cluster galaxies, a process that we shall investigate in the next section.

Clusters form hierarchically from smaller systems that resemble groups of galaxies (§§1.1.5 and 9.2.2). Groups have velocity dispersions of only ~ 300 km s⁻¹, so encounters in groups occur at lower speeds and have stronger effects—they frequently lead to mergers—and dynamical friction is more powerful. Most of the galaxy evolution that we see in the centers of rich clusters may thus be due to “pre-processing” in the groups that later merged to form the cluster. For example, the exceptionally luminous brightest cluster galaxies or cD galaxies (§1.1.3) that are found at the centers of clusters probably arise from the merger of galaxies in precursor groups (Dubinski 1998).

8.3 Tides

In the last section we examined how tidal shocks from high-speed encounters heat stellar systems and erode their outer parts. We now consider the opposite limiting case of a static tidal field. The simplest example of a static tide occurs when a satellite travels on a circular orbit in the gravitational field of a much larger spherical host system. In this case, the satellite experiences no shocks—in fact, in the frame rotating with the satellite the external tidal field is stationary—so in the absence of other relaxation effects, a sufficiently small system could survive indefinitely. However, a static tidal field prunes distant stars from the satellite system, thereby enforcing an upper limit on its size. Observationally, globular clusters and other satellite systems often show a fairly sharp outer boundary, which is called the tidal radius on the assumption that it is caused by this process (see §4.3.3c, BM §6.1.10, and King 1962).

In the following subsections, we analyze the effect of a tidal field on a satellite in a circular orbit using two complementary approaches.

8.3.1 The restricted three-body problem

Let us suppose that the host and satellite systems are point masses M and m , traveling at separation R_0 in a circular orbit around their mutual center of mass. The **restricted three-body problem** is to find the trajectory of a massless test particle that orbits in the combined gravitational field of these two masses (Szebehely 1967; Hénon 1997; Valtonen & Karttunen 2006). Solutions of this problem provide a good approximation to the motion of stars in the outer parts of a satellite stellar system that is on a circular orbit near or beyond the outer edge of a spherical host system.

The two masses orbit their common center of mass with angular speed

$$\Omega = \sqrt{\frac{G(M+m)}{R_0^3}}, \quad (8.85)$$

so the gravitational field is stationary when referred to a coordinate system centered on the center of mass that rotates at speed Ω . We orient this coordinate system so that the centers of the satellite and host systems are at $\mathbf{x}_m = [MR_0/(M+m), 0, 0]$ and $\mathbf{x}_M = [-mR_0/(M+m), 0, 0]$, and the angular speed is $\boldsymbol{\Omega} = (0, 0, \Omega)$. In §3.3.2 we showed that on any orbit in such a system, the Jacobi integral

$$\begin{aligned} E_J &= \frac{1}{2}v^2 + \Phi(\mathbf{x}) - \frac{1}{2}|\boldsymbol{\Omega} \times \mathbf{x}|^2 \\ &= \frac{1}{2}v^2 + \Phi_{\text{eff}}(\mathbf{x}) \end{aligned} \quad (8.86)$$

is conserved (eq. 3.113). Since $v^2 \geq 0$, a star with Jacobi integral E_J can never trespass into a region where $\Phi_{\text{eff}}(\mathbf{x}) > E_J$. Consequently, the surface $\Phi_{\text{eff}}(\mathbf{x}) = E_J$, the zero-velocity surface for stars of Jacobi integral E_J , forms an impenetrable wall for such stars. Figure 8.6 shows contours of constant Φ_{eff} in the equatorial plane of two orbiting point masses; the Lagrange points are the extrema (maxima and saddle points) of this surface. It is instructive to compare these contours to those in a bar-like potential, shown in Figure 3.14.⁹ The stability of orbits near the Lagrange points in the restricted three-body problem is discussed in Problem 3.25.

From the figure we see that the zero-velocity surfaces near each body are centered on it, but farther out the zero-velocity surfaces surround both bodies. Hence, at the critical value of Φ_{eff} corresponding to the last zero-velocity surface to enclose only one body, there is a discontinuous change in the region confined by the Jacobi integral. The last closed zero-velocity surface surrounding a single body is called its **tidal** or **Roche surface**; since this surface touches the Lagrange point L_3 that lies between the two masses

⁹Note that different authors use different conventions for the numbering of the Lagrange points L_1, L_2, L_3 .

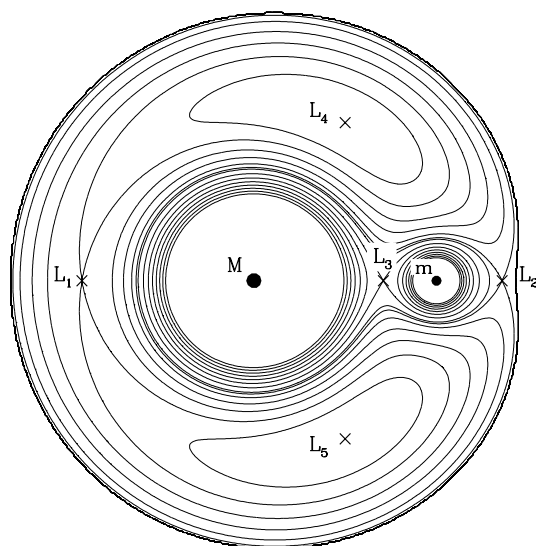


Figure 8.6 Contours of equal effective potential Φ_{eff} defined by equation (8.88) for two point masses in a circular orbit. The mass ratio $m/M = \frac{1}{9}$. The points L_1, \dots, L_5 are the Lagrange points. The L_4 and L_5 points form an equilateral triangle with the two masses (Problem 3.25).

on the line connecting them, it is natural to identify the outermost radius of orbits bound to m as the distance r_J between m and L_3 .

We may evaluate r_J by noticing that at $(x_m - r_J, 0, 0)$ the effective potential has a saddle point, so

$$\left(\frac{\partial \Phi_{\text{eff}}}{\partial x} \right)_{(x_m - r_J, 0, 0)} = 0. \quad (8.87)$$

For two point masses a distance R_0 apart, equations (8.85) and (8.86) imply

$$\Phi_{\text{eff}}(\mathbf{x}) = -G \left[\frac{M}{|\mathbf{x} - \mathbf{x}_M|} + \frac{m}{|\mathbf{x} - \mathbf{x}_m|} + \frac{M+m}{2R_0^3} (x^2 + y^2) \right]. \quad (8.88)$$

At a point between the two masses, (8.87) is satisfied if

$$0 = \frac{1}{G} \left(\frac{\partial \Phi_{\text{eff}}}{\partial x} \right)_{(x_m - r_J, 0, 0)} = \frac{M}{(R_0 - r_J)^2} - \frac{m}{r_J^2} - \frac{M+m}{R_0^3} \left(\frac{MR_0}{M+m} - r_J \right). \quad (8.89)$$

This equation leads to a fifth-order polynomial whose roots give r_J . In general these roots must be found numerically. However, if the satellite is small, $m \ll M$, then $r_J \ll R_0$, and we can expand $(R_0 - r_J)^{-2}$ in powers of r_J/R_0 to find

$$0 = \frac{M}{R_0^2} \left(1 + \frac{2r_J}{R_0} + \dots \right) - \frac{m}{r_J^2} - \frac{M}{R_0^2} + \frac{M+m}{R_0^3} r_J \simeq \frac{3Mr_J}{R_0^3} - \frac{m}{r_J^2}. \quad (8.90)$$

Truncating the series in this way is none other than the distant-tide approximation. Then to first order in r_J/R_0 ,

$$r_J = \left(\frac{m}{3M} \right)^{1/3} R_0. \quad (8.91)$$

We call the radius r_J the **Jacobi radius** of the mass m ; alternative names are the **Roche** or **Hill radius**. The Jacobi radius of an orbiting stellar system is expected to correspond to the observational tidal radius, the maximum extent of the satellite system. However this correspondence is only approximate, for several reasons:

- (i) The Roche surface is not spherical (see Problem 8.11), so it cannot be fully characterized by a single radius.
- (ii) All we have established is that a test particle can never cross the Roche surface if it lies inside the Roche surface and has a velocity (in the rotating frame) small enough that $E_J < \Phi_{\text{eff}}(L_3)$. Stars with larger velocities may or may not escape from the satellite; conversely, stars that lie outside the Roche surface can, in some cases, remain close to the satellite for all future times (see Problem 8.13 and Hénon 1970). The approximate correspondence between the Jacobi radius and the observational tidal radius arises because the fraction of velocity space occupied by orbits that remain close to the satellite diminishes rapidly beyond r_J .
- (iii) In most applications, the satellite system is not on a circular orbit. When m is on an eccentric orbit, there is no reference frame in which the potential experienced by a test particle is stationary, and no analog of the Jacobi integral exists.¹⁰ Thus no direct generalization of our derivation of the Jacobi radius to the case of non-circular satellite orbits is possible. King (1962) and others have argued that if the satellite is on a non-circular orbit, the tidal radius is still given by equation (8.91), but with R_0 replaced by the pericenter distance (we used an analogous argument to describe the tidal disruption of stars orbiting a massive black hole; see eq. 7.200). A more accurate approach is to recognize that the effect of tidal fields on satellites in non-circular orbits is intermediate between tidal radii—a concept that applies to circular orbits—and tidal

¹⁰ Although analogs of the Lagrange points can exist (Szebehely 1967).

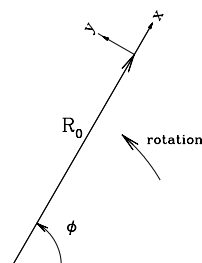


Figure 8.7 The rotating (x, y) coordinate system used in Hill's approximation.

shocks—which apply to high-velocity or plunging orbits. The tidal radius limits the satellite at a fixed size, no matter how many orbits it travels, while tidal shocks prune the satellite more and more at every pericenter passage.

- (iv) Stars are usually lost from the satellite as a result of weak perturbations, such as two-body relaxation, that drive E_J slightly higher than $\Phi_{\text{eff}}(L_3)$. Such stars drift slowly away from the satellite and thus can remain close to the satellite for many orbital periods, thereby contributing to the star counts even though they are no longer bound to the satellite (Fukushige & Heggie 2000).
- (v) In many cases, the satellite orbits within the body of the host system, so the point-mass approximation used in deriving equation (8.91) is not accurate. This defect, at least, is easy to remedy—see equation (8.106) below.

Tidal forces can be thought of as imposing a limit on the density of a satellite. Let $\bar{\rho} \equiv m/(\frac{4}{3}\pi r_J^3)$ be the mean density of the satellite within a distance r_J , and $\bar{\rho}_h \equiv M/(\frac{4}{3}\pi R_0^3)$ be the mean density of the host inside the orbital radius R_0 . Then equation (8.91) states that

$$\bar{\rho} = 3\bar{\rho}_h; \quad (8.92)$$

to within a factor of order unity, a satellite is pruned by tidal forces until its mean density equals the mean density of its host interior to its orbital radius.

8.3.2 The sheared-sheet or Hill's approximation

When the satellite is much smaller than the distance to the center of the host system, we can use the distant-tide approximation for the host's gravitational field (§8.2.1). We consider a spherically symmetric host system with potential $\Phi(R)$ at a distance R from its center; here we do *not* assume that the host is a point mass, so $\Phi(R)$ is not necessarily the Keplerian potential $-GM/R$.

We assume that the satellite travels on a circular orbit at distance R_0 from the center of the host. We work in a frame with origin at the center of mass of the satellite, in which the x - y plane coincides with the orbital plane of the satellite, $\hat{\mathbf{e}}_x$ points directly away from the center of the host system, and $\hat{\mathbf{e}}_y$ points in the direction of the orbital motion of the satellite (see Figure 8.7). This frame rotates with the circular frequency $\boldsymbol{\Omega}_0 \equiv \Omega_0 \hat{\mathbf{e}}_z$, so the acceleration of a particle in the satellite is given by equation (3.116),

$$\frac{d^2 \mathbf{x}}{dt^2} = -\nabla \Phi - 2\boldsymbol{\Omega}_0 \times \frac{d\mathbf{x}}{dt} - \boldsymbol{\Omega}_0 \times (\boldsymbol{\Omega}_0 \times \mathbf{x}), \quad (8.93)$$

where

$$\nabla \Phi = \nabla \Phi_s + \sum_{k=1}^3 \Phi_{jk} x_k. \quad (8.94)$$

Here $\Phi_s(\mathbf{x})$ is the gravitational potential from the satellite, and the second term arises from the distant-tide approximation (8.35). In our coordinate system, the center of the host is located at $\mathbf{X} = (-R_0, 0, 0)$ and from equation (8.36):

$$\Phi_{xx} = \Phi''(R_0) \quad ; \quad \Phi_{yy} = \Phi_{zz} = \frac{\Phi'(R_0)}{R_0} \quad ; \quad \Phi_{xy} = \Phi_{xz} = \Phi_{yz} = 0. \quad (8.95)$$

The equations of motion (8.93) read

$$\begin{aligned} \ddot{x} &= 2\Omega_0 \dot{y} + [\Omega_0^2 - \Phi''(R_0)] x - \frac{\partial \Phi_s}{\partial x}; \\ \ddot{y} &= -2\Omega_0 \dot{x} + \left[\Omega_0^2 - \frac{\Phi'(R_0)}{R_0} \right] y - \frac{\partial \Phi_s}{\partial y}; \\ \ddot{z} &= -\frac{\Phi'(R_0)}{R_0} z - \frac{\partial \Phi_s}{\partial z}. \end{aligned} \quad (8.96)$$

Using the relation $\Phi'(R_0) = R_0 \Omega_0^2$ we see that the term in square brackets in the second line vanishes. Moreover we can rewrite $\Omega_0^2 - \Phi''(R_0)$ as $-2R_0 \Omega_0 \Omega'(R_0)$ and this in turn can be rewritten as $4\Omega_0 A_0$ where $A_0 = A(R_0)$ is given by equation (3.83). Thus

$$\ddot{x} - 2\Omega_0 \dot{y} - 4\Omega_0 A_0 x = -\frac{\partial \Phi_s}{\partial x} \quad ; \quad \ddot{y} + 2\Omega_0 \dot{x} = -\frac{\partial \Phi_s}{\partial y} \quad ; \quad \ddot{z} + \Omega_0^2 z = -\frac{\partial \Phi_s}{\partial z}. \quad (8.97)$$

These are the equations of motion in the **sheared sheet** or **Hill's approximation**, named after the mathematician G. W. Hill, who used this approach to study the motion of the Moon in the nineteenth century (Murray & Dermott 1999).

(a) The epicycle approximation and Hill's approximation We first consider the trajectories of test particles in the absence of a satellite, $\Phi_s(\mathbf{x}) = 0$. The simplest solutions of equations (8.97) have the form

$$x(t) = x_g = \text{constant} \quad ; \quad y(t) = -2A_0 x_g t + \text{constant} \quad ; \quad z(t) = 0. \quad (8.98)$$

These are the analogs of circular orbits in the host system. The general solution is

$$\begin{aligned} x(t) &= x_g + X \cos(\kappa_0 t + \alpha), \\ y(t) &= y_g(t) - Y \sin(\kappa_0 t + \alpha), \quad \text{where } y_g(t) = y_{g0} - 2A_0 x_g t, \\ z(t) &= Z \cos(\Omega_0 t + \alpha_z), \end{aligned} \quad (8.99)$$

where x_g , y_{g0} , X , Z , α and α_z are arbitrary constants and

$$\kappa_0^2 = 4\Omega_0(\Omega_0 - A_0) = -4\Omega_0 B_0 \quad ; \quad \frac{X}{Y} = \frac{\kappa_0}{2\Omega_0}. \quad (8.100)$$

Here $B_0 = A_0 - \Omega_0$ (eq. 3.84). Thus we have re-derived the epicycle approximation of §3.2.3, in particular the relation between the epicycle frequency κ_0 and Oort's constants (eq. 3.84) and the ratio of the axes of the epicycle (eq. 3.95). The difference between the two derivations is that §3.2.3 described an *approximate* solution of the *exact* equations of motion for a particle on a nearly circular orbit, while here we have found an *exact* solution of Hill's *approximate* equations of motion. Note that in Hill's approximation all particles have the same epicycle frequency.

It is straightforward to verify that when $\Phi_s = 0$, the following expressions are integrals of the motion:

$$E_{\parallel} \equiv \frac{1}{2}(\dot{x}^2 + \dot{y}^2 - 4\Omega_0 A_0 x^2) \quad ; \quad E_{\perp} \equiv \frac{1}{2}(z^2 + \Omega_0^2 z^2) \quad ; \quad L \equiv \dot{y} + 2\Omega_0 x; \quad (8.101)$$

$E_{\parallel} + E_{\perp}$ and $R_0 L$ differ from the Jacobi integral and the angular momentum by constant terms and terms of order $O(x^3, y^3)$. These expressions are related to the constants in the orbit solutions (8.99) by

$$E_{\parallel} = 2A_0 B_0 x_g^2 + \frac{1}{2}\kappa_0^2 X^2 \quad ; \quad E_{\perp} = \frac{1}{2}\Omega_0^2 Z^2 \quad ; \quad L = -2B_0 x_g. \quad (8.102)$$

A circular orbit has $E_{\parallel} = \frac{1}{2}A_0 L^2/B_0$; hence it is natural to define the **epicycle energy** E_x as the difference

$$\begin{aligned} E_x &\equiv E_{\parallel} - \frac{A_0 L^2}{2B_0}, \\ &= \frac{1}{2}[\dot{x}^2 + \kappa_0^2(x - x_g)^2], \\ &= \frac{1}{2}\kappa_0^2 X^2, \\ &= \frac{1}{2}\dot{x}^2 + \frac{2\Omega_0^2}{\kappa_0^2}(\dot{y} + 2A_0 x)^2. \end{aligned} \quad (8.103)$$

Some of these results, derived in other ways and with slightly different notation, have already appeared as equations (3.86) and (3.102).

(b) The Jacobi radius in Hill's approximation If a satellite is present, with potential $\Phi_s(\mathbf{x})$, the integrals in equation (8.101) are no longer conserved; the only remaining classical integral is (Problem 8.14)

$$E \equiv \frac{1}{2}(\dot{x}^2 + \dot{y}^2 + \dot{z}^2 - 4\Omega_0 A_0 x^2 + \Omega_0^2 z^2) + \Phi_s(\mathbf{x}). \quad (8.104)$$

This integral is the analog of the Jacobi integral (8.86).

Now let us imagine that the satellite potential Φ_s arises from a mass m that is located at $\mathbf{x} = \mathbf{0}$. The equations of motion (8.97) become

$$\ddot{x} = 2\Omega_0\dot{y} + 4\Omega_0A_0x - \frac{Gmx}{r^3}; \quad \ddot{y} = -2\Omega_0\dot{x} - \frac{Gmy}{r^3}; \quad \ddot{z} = -\Omega_0^2z - \frac{Gmz}{r^3}, \quad (8.105)$$

where $r^2 = x^2 + y^2 + z^2$. The test particle remains stationary ($\ddot{x} = \dot{x} = \ddot{y} = \dot{y} = \ddot{z} = 0$) if and only if $y = z = 0$ and $4\Omega_0A_0 = Gm/|x|^3$. These conditions are satisfied for the points on the x axis with

$$x = \pm r_J, \quad \text{where} \quad r_J \equiv \left(\frac{Gm}{4\Omega_0A_0} \right)^{1/3}. \quad (8.106)$$

These stationary points are analogs to the Lagrange points L_2 and L_3 in the restricted three-body problem (Figure 8.6). If the host is a point mass $M \gg m$, then $\Omega(R) = (GM/R^3)^{1/2}$ so $A_0 = \frac{3}{4}\Omega_0$ and

$$r_J = \left(\frac{m}{3M} \right)^{1/3} R_0. \quad (8.107)$$

Thus we recover expression (8.91) for the Jacobi radius. For a spherical host with mass $M(R)$ interior to radius R , it is straightforward to show that this expression is modified by replacing M by $M(R_0)$ and multiplying the Jacobi radius by a factor

$$f = \left(1 - \frac{1}{3} \frac{d \ln M}{d \ln R} \right)^{-1/3}. \quad (8.108)$$

The factor f is unity for a point mass and 1.145 for a singular isothermal sphere ($M \propto R$). For a homogeneous sphere ($M \propto R^3$) f diverges, so there is no Jacobi radius: in this case the host potential $\Phi(R) = \frac{1}{2}\Omega_0^2R^2$ and Oort's constant $A_0 = 0$, so the tidal field $4\Omega_0A_0x$ in the equations of motion (8.97) is absent. Physically, there is no Jacobi radius because all stars in this potential have the same orbital period: thus, even if the satellite mass were zero, stars in nearly circular orbits with similar radii and azimuths will continue to have similar radii and azimuths at all future times.

8.3.3 Tidal tails and streamers

We now investigate what happens to stars after they are stripped from a satellite by tidal forces, with the help of the angle-action variables described in §3.5 (Helmi & White 1999; Tremaine 1999). Consider a satellite of mass m orbiting a host that has mass $M \gg m$ interior to the satellite orbit. At its pericenter, a distance R from the center of the host, the satellite is pruned by tidal forces to a radius $r \approx R(m/M)^{1/3}$. Its velocity at pericenter is $V \approx$

$(GM/R)^{1/2}$. To a first approximation, we may assume that stars lost from the satellite no longer feel its gravitational field, and follow orbits determined solely by the field of the host. On such orbits, the actions \mathbf{J} are constant and the angles $\boldsymbol{\theta}$ increase linearly with time, at a rate $\dot{\boldsymbol{\theta}} = \boldsymbol{\Omega} = \partial H / \partial \mathbf{J}$ where H is the Hamiltonian corresponding to the host potential. The stripped stars have a range of actions and angles, which we write as $\mathbf{J}_0 \pm \Delta \mathbf{J}$, $\boldsymbol{\theta}_0 \pm \Delta \boldsymbol{\theta}$. The mean actions \mathbf{J}_0 and the mean angles $\boldsymbol{\theta}_0$ at the time the satellite passes through pericenter are very nearly the actions and angles of the satellite at that time, since the tidal forces are symmetric about its center of mass. The spread in actions and angles in the stripped stars arises from two effects: (i) the stars are lost from both the inner and outer edge of the satellite (near the Lagrange points L_3 and L_2), and (ii) the stars have a range of velocities, roughly equal to the velocity dispersion σ of the satellite. These effects lead to a fractional spread $r/R \sim (m/M)^{1/3}$ in position and σ/V in velocity. Since $\sigma \approx (Gm/r)^{1/2} \sim V(m/M)^{1/2}(R/r)^{1/2} \sim V(m/M)^{1/3}$ the two effects yield approximately the same fractional spread. Thus, the stripped stars are initially distributed through ranges in action and angle given by

$$\frac{\Delta J_i}{J_i}, \Delta \theta_i \sim \left(\frac{m}{M}\right)^{1/3}. \quad (8.109)$$

The spread in actions leads to a spread in orbital frequencies

$$\Delta \Omega_i \sim \sum_{j=1}^3 H_{ij} \Delta J_j, \quad \text{where} \quad D_{ij} \equiv \frac{\partial^2 H}{\partial J_i \partial J_j} \quad (8.110)$$

is the Hessian of the Hamiltonian. The spread in angles grows linearly with time, such that

$$\Delta \boldsymbol{\theta}(t) = \Delta \boldsymbol{\theta}(0) + \Delta \boldsymbol{\Omega} t, \quad (8.111)$$

where $t = 0$ is the time at which the stars were stripped. At large times the second term dominates, so we have

$$\Delta \boldsymbol{\theta}(t) \simeq t \mathbf{D} \cdot \Delta \mathbf{J}. \quad (8.112)$$

Since the matrix \mathbf{D} is symmetric, it is diagonalizable, that is, there exists an orthogonal matrix \mathbf{A} such that

$$\mathbf{A} \mathbf{D} \mathbf{A}^T = \tilde{\mathbf{D}}, \quad (8.113)$$

where $\mathbf{A}^T = \mathbf{A}^{-1}$ is the transpose of \mathbf{A} ($A_{jk}^T = A_{kj}$), and $\tilde{\mathbf{D}}$ is the diagonal matrix formed by the eigenvalues λ_i of \mathbf{D} . We now make a canonical transformation to new angle-action variables $(\boldsymbol{\theta}', \mathbf{J}')$ using the generating function $S(\boldsymbol{\theta}, \mathbf{J}') = \mathbf{J}' \cdot \mathbf{A} \cdot \boldsymbol{\theta}$ (eq. D.93); thus

$$\boldsymbol{\theta}' = \frac{\partial S}{\partial \mathbf{J}'} = \mathbf{A} \cdot \boldsymbol{\theta} \quad ; \quad \mathbf{J} = \frac{\partial S}{\partial \boldsymbol{\theta}} = \mathbf{A}^T \cdot \mathbf{J}'. \quad (8.114)$$

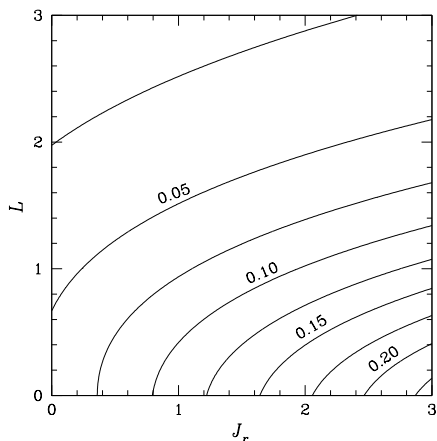


Figure 8.8 The ratio of the two largest eigenvalues of \mathbf{D} , the Hessian of the Hamiltonian, for the isochrone potential (see §3.1c and Problem 3.41). The axes are the radial action J_r and the angular momentum L . When this ratio is small compared to unity, tidally stripped stars form a one-dimensional filament or tidal streamer.

In terms of the new variables, equation (8.112) becomes

$$\Delta\theta'(t) \simeq t\tilde{\mathbf{D}} \cdot \Delta\mathbf{J}' \quad \text{or} \quad \Delta\theta'_i(t) \simeq t\lambda_i\Delta J'_i \quad (\text{no summation over } i). \quad (8.115)$$

This result shows that the cloud of escaped stars spreads into three of the six phase-space dimensions, at rates determined by the initial spread in actions and the eigenvalues λ_i of the matrix \mathbf{D} . Small satellites have a smaller spread in actions and so disperse more slowly. If one of the three eigenvalues is zero, or at least much smaller than the other two, the cloud will expand in two dimensions in phase space, creating a sheet; this is the situation for tidal streamers in a spherical host galaxy. If two of the three eigenvalues are zero, the cloud will expand in one dimension to produce a filament; this is the situation in a Keplerian potential. Even when two or more of the eigenvalues of \mathbf{D} are non-zero, usually one is large enough compared to the others that the disrupted stars form a relatively thin tail, which is called a **tidal streamer** or **tail** (see Figure 8.8)—usually the term “tail” is reserved for the long, prominent, massive streamers formed in major mergers of two disk galaxies.

Known tidal streamers are associated with the Magellanic Clouds (the Magellanic Stream, already described in §8.1.1c), the globular cluster Pal 5 (Figure 8.9), and the Sagittarius galaxy (Figure 8.10).

Unlike comet tails, tidal streamers are symmetrical structures that both lead and lag the satellite along its orbit. For example, in Figure 8.9 the upper streamer is made up of stars that have longer orbital periods than the cluster, and hence trail behind it; conversely, the streamer at lower right contains stars that are on shorter-period orbits, and race ahead of it.

In Chapter 9 we shall argue that galaxies form by hierarchical merging of smaller subunits. In the merging process, these subunits are disrupted by tidal forces, and the debris—both stars and dark matter—forms a vast web

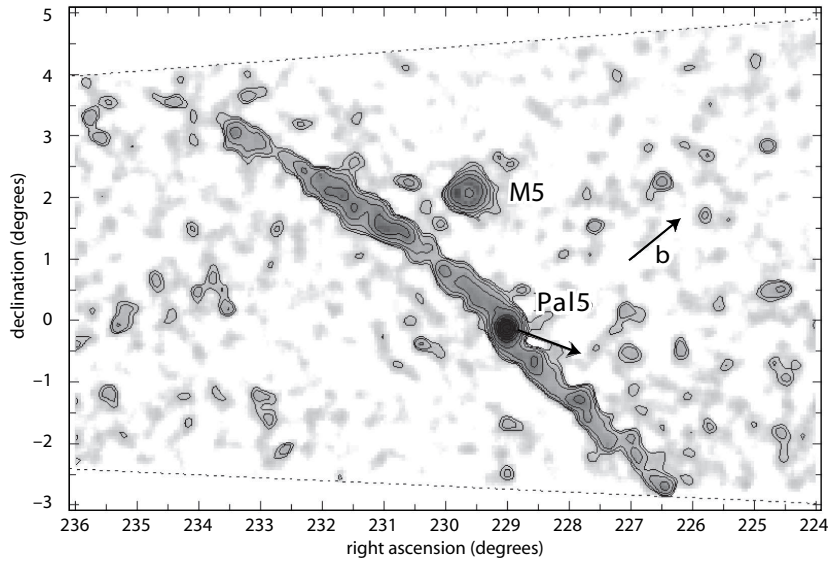


Figure 8.9 Tidal streamers emerging from the globular cluster Pal 5. The plot shows the surface density of stars whose distances are consistent with the cluster distance. The clump labeled “M5” is a residual feature from the unrelated cluster M5. The arrow from Pal 5 indicates the direction of its motion on the sky, and the arrow labeled “b” shows the direction of increasing Galactic latitude. The dotted lines mark the borders of the field. See Grillmair & Dionatos (2006) for maps of the streamers at even larger distances from the cluster. From Odenkirchen et al. (2003), by permission of the AAS.

of tidal streamers. The number of streamers per unit volume and the corresponding degree of irregularity in the mass distribution of the halo depend on the distance from the center of the galaxy: at small radii, the galaxy is hundreds or thousands of crossing times old and the tidal streamers are thoroughly phase-mixed, while at large radii subunits are falling in for the first time and the substructure will be much more prominent (Helmi, White, & Springel 2003). At any given radius, the substructure is likely to be stronger in the baryons (stars and gas) than in the dark matter, since the baryons are concentrated in the dense centers of the dark-matter halos and thus are less susceptible to tidal forces. Efforts to detect and disentangle this web are still in their infancy.

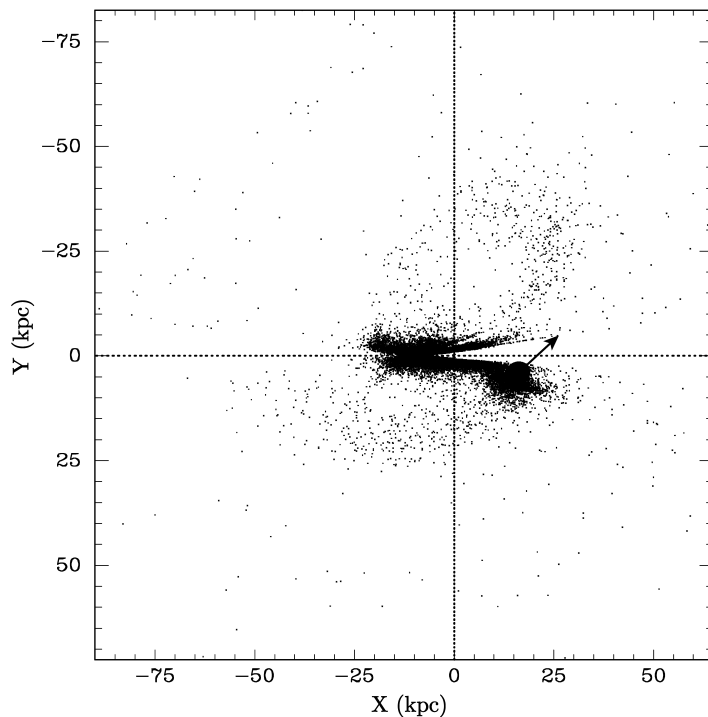


Figure 8.10 The distribution of M-giant stars lying within 7 kpc of the orbital plane of the Sagittarius dwarf galaxy. The figure is a projection onto this orbital plane, which is tipped by 77° from the Galactic plane. The Galactic disk lies along $Y = 0$, the Galactic center is at the origin, and the Sun is at $X \simeq -8$ kpc, $Y \simeq 0$ (by coincidence, the Sun lies nearly in the Sagittarius orbital plane). Stars that are highly reddened have been removed, which creates the wedge-shaped gap stretching right from the Sun. The Sagittarius galaxy is located at $X \simeq 15$ kpc, $Y \simeq 5$ kpc, and the arrow extending from it indicates the direction of its velocity vector. Tidal debris from the galaxy is evident as the prominent arc passing through $(X, Y) \simeq (25 \text{ kpc}, -30 \text{ kpc})$ above the Galactic plane, and through $(X, Y) \simeq (-15 \text{ kpc}, 15 \text{ kpc})$ below the plane. Most of the width of the arcs is probably due to errors in the stellar distances. From Majewski et al. (2003).

8.4 Encounters in stellar disks

The velocity distribution of stars in the solar neighborhood is approximately described by the Schwarzschild distribution introduced in §4.4.3 (see also Problem 8.16). In this DF, the number of stars with velocity \mathbf{v} in a small range $d^3\mathbf{v}$ is

$$f(\mathbf{v})d^3\mathbf{v} = \frac{n_0 d^3\mathbf{v}}{(2\pi)^{3/2}\sigma_R\sigma_\phi\sigma_z} \exp\left[-\left(\frac{v_R^2}{2\sigma_R^2} + \frac{v_\phi^2}{2\sigma_\phi^2} + \frac{v_z^2}{2\sigma_z^2}\right)\right]. \quad (8.116)$$

Here n_0 is the number of stars per unit volume, σ_R , σ_ϕ , and σ_z are the velocity dispersions along the three axes of a cylindrical coordinate system

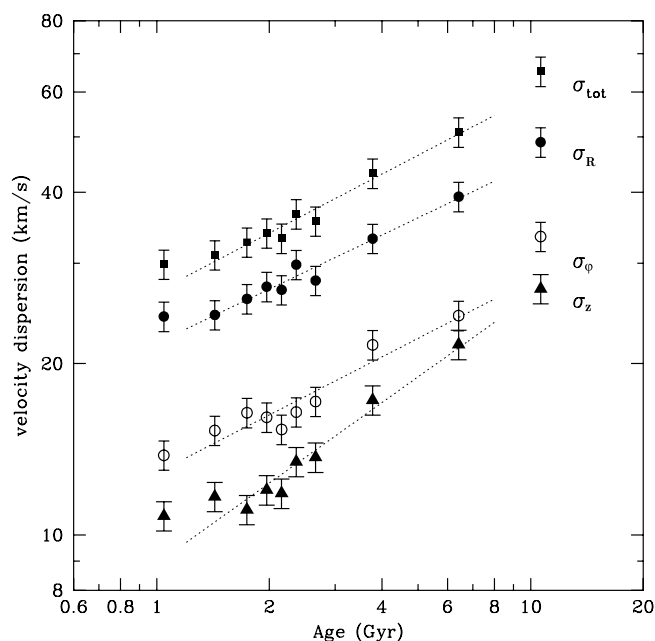


Figure 8.11 The velocity dispersion of stars in the solar neighborhood as a function of age, from Nordström et al. (2004). From bottom to top, the plots show the vertical dispersion σ_z , the azimuthal dispersion σ_ϕ , the radial dispersion σ_R , and the RMS velocity $(\sigma_R^2 + \sigma_\phi^2 + \sigma_z^2)^{1/2}$. The lines show fits of the form $\sigma_i \propto t^\alpha$ where t is the age; from bottom to top the best-fit exponents α are 0.47, 0.34, 0.31, and 0.34.

centered on the center of the Galaxy, and \mathbf{v} is the velocity relative to the velocity of a circular orbit passing through the solar neighborhood (the Local Standard of Rest or LSR; see §1.1.2). The radial and azimuthal dispersions are approximately related by

$$\frac{\sigma_\phi}{\sigma_R} = \frac{\kappa}{2\Omega}, \quad (8.117)$$

where κ and Ω are the epicycle frequency and the azimuthal frequency for nearly circular orbits in the solar neighborhood (see eq. 3.100 and Problem 4.43). Equation (8.116) states that the density of stars in velocity space is constant on ellipsoids with principal axes σ_R , σ_ϕ and σ_z , called velocity ellipsoids in §4.1.2.

Although the shape of the velocity ellipsoid is approximately the same for different types of stars, its size is not: the dispersions σ_i ($i = R, \phi, z$) of cool, red stars are almost three times as large as those of hot, blue stars (BM Figure 10.12 and Tables 10.2 and 10.3). Since blue stars are young,

while red stars are a mixture of mostly old and a few young stars, this trend suggests that stars are born on nearly circular or “cold” orbits, and as a stellar population ages it “heats up” in the sense that the dispersions σ_i increase. This hypothesis can be confirmed by measuring the age-velocity dispersion relation for nearby stars (Figure 8.11). These observations show that $\sigma_i \propto t^\alpha$, where $\alpha \simeq 0.5$ for σ_z and $\alpha \simeq 0.3$ for both σ_R and σ_ϕ —the exponent is necessarily the same for these two dispersions, because they are related by equation (8.117).

We refer to the steady increase of these dispersions with time as **disk heating**, and in this section we investigate the dynamics of this process. A natural first step in this investigation is to wonder whether disk heating can be due to the accumulation of small velocity kicks from passing stars. This process was described briefly in §1.2.1 and more thoroughly in §7.4. In particular, in the discussion following equation (7.106) we saw that encounters between stars in the solar neighborhood have a negligible effect on their velocities over the age of the Galaxy. Thus we must seek other explanations.

The simplest mechanism for disk heating is encounters with hypothetical massive objects in the dark halo, or MACHOs. This process was investigated in §7.4.4, where we found that the predicted time dependence of the velocity dispersion σ_R is incorrect. Moreover, the required MACHO mass appears to be incompatible with observations of wide binary stars in the halo (§8.2.2e). We therefore examine other possibilities.

8.4.1 Scattering of disk stars by molecular clouds

Long before molecular clouds were detected, Spitzer & Schwarzschild (1951, 1953) suggested that encounters between disk stars and massive gas clouds might be responsible for the random velocities of old disk stars. In Figure 8.12 we illustrate how a molecular cloud or other mass m on a circular orbit in a disk affects the orbits of nearby stars. Since the cloud mass is $\lesssim 10^{-5}$ times the mass of the Galaxy, we may use Hill’s approximation (§8.3.2), in which the cloud is at rest at the origin of a rotating Cartesian coordinate system, with the x axis pointing radially outward and the y axis in the direction of rotation. For simplicity we neglect motion perpendicular to the x - y plane. The stellar trajectories are given by the equations of motion (8.97), where the cloud potential $\Phi_s = -Gm/(x^2 + y^2)^{1/2}$. In the figure, we write the distances in terms of the Jacobi radius of the cloud (eq. 8.106),

$$r_J = \left(\frac{Gm}{4\Omega_0 A_0} \right)^{1/3} = 52 \text{ pc} \left(\frac{m}{10^5 \mathcal{M}_\odot} \frac{\Omega_0}{A_0} \right)^{1/3} \left(\frac{220 \text{ km s}^{-1}}{v_c} \frac{R_0}{8 \text{ kpc}} \right)^{2/3}. \quad (8.118)$$

In these units, the equations of motion are independent of m , so Figure 8.12 applies to clouds of any mass.

The figure shows only stars on initially circular orbits that are larger than the cloud’s orbit. The behavior of orbits that are smaller than the

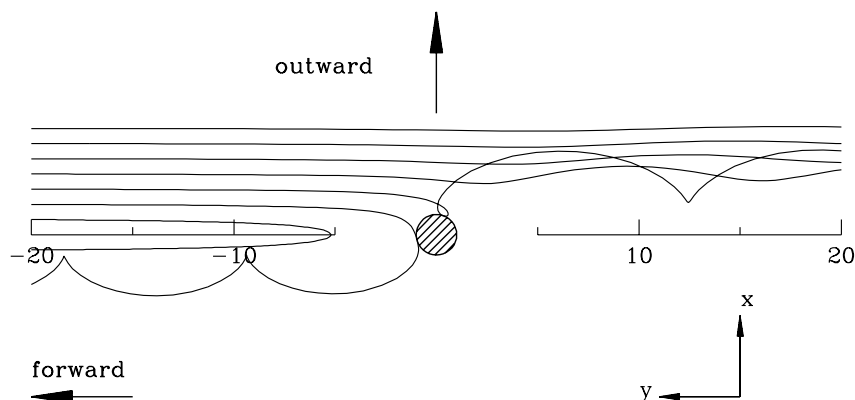


Figure 8.12 The trajectories of stars near a mass point in a disk. The orbit of the mass point is circular, as are the initial orbits of the stars. The coordinate frame co-rotates with the mass point, which is therefore fixed at the origin; in this expanded view of the area near the mass point, upwards is radially outward and the mass point is traveling to the left in an inertial frame. All orbits are restricted to the $z = 0$ plane. Circular orbits appear as straight horizontal lines. A sequence of seven orbits is shown, all initially circular with radius slightly larger than that of the point mass, so it overtakes them. The behavior of orbits inside the point mass is given by reflecting these orbits through the origin. The trajectories of the stars are described by equations (8.97). The disk is assumed to have a flat circular-speed curve, so Oort's constant $A_0 = \Omega_0/2$ and the epicycle frequency $\kappa_0 = \sqrt{2}\Omega_0$. Distances are measured in units of the Jacobi radius r_J of the mass point (eq. 8.106), which is also the radius of the circle representing its location.

cloud's can be deduced by reflecting the orbits shown through the origin of the figure. The initial orbits shown have angular speeds that are smaller than the cloud's, so the cloud overtakes them (i.e., they move to the right in the cloud frame of reference that is used in the figure). If the initial difference in orbital radii $\Delta r \lesssim r_J$, the encounter simply reverses the direction of the orbit relative to the cloud, without imparting any significant epicycle motion.¹¹ If $\Delta r \sim r_J$, the encounter imparts significant epicycle motion—the epicycle amplitude is comparable to Δr and the encounter may or may not reverse the overall direction of motion of the orbit relative to the cloud. For $\Delta r \gtrsim r_J$, the star passes the cloud and acquires a small epicyclic motion. It is this excitation of epicycle motion by encounters with clouds that warms the disk.

We may estimate the efficiency of this process by using the impulse approximation to find the radial velocity acquired by a star that is initially on a circular orbit of radius R . If R_c is the radius of the cloud orbit, then in the sheared sheet the star's initial orbit is $x = R - R_c \equiv b = \text{constant}$,

¹¹ This is an example of the donkey effect, described in Box 3.3. No epicycle motion is excited because the star approaches the cloud slowly, so its eccentricity is an adiabatic invariant. For a comprehensive discussion of the trajectories in this problem, see Petit & Hénon (1986).

$y(t) = \text{constant} - 2A_0bt$ (eq. 8.99). Integrating the gravitational attraction of the potential $\Phi_s(\mathbf{x})$ along this trajectory yields

$$\dot{x} = - \int_{-\infty}^{\infty} dt \frac{Gmb}{[b^2 + y^2(t)]^{3/2}} = - \frac{Gm}{A_0b^2}. \quad (8.119)$$

In the impulse approximation, immediately after the encounter we have

$$x = b \quad ; \quad \dot{x} = - \frac{Gm}{A_0b^2} \quad ; \quad \dot{y} = -2A_0b, \quad (8.120)$$

and the corresponding epicycle energy and amplitude are given by equations (8.103),

$$\Delta E_x = \frac{1}{2}\dot{x}^2 = \frac{f^2}{2} \left(\frac{Gm}{A_0b^2} \right)^2 \quad ; \quad X = f \frac{Gm}{A_0\kappa_0b^2}, \quad (8.121a)$$

where $f = 1$. We have introduced the correction factor f because the impulse approximation is not accurate: since the radial velocity oscillates with the epicycle frequency κ_0 , the impulse approximation requires that the duration of the encounter is much less than the epicycle period. In fact, the duration is $\approx b/\dot{y} = (2A)^{-1}$, which is comparable to the epicycle period $2\pi/\kappa_0$. Hence the impulse approximation makes an error of order unity in the epicycle energy. The correct derivation (Julian & Toomre 1966 and Problem 8.19) yields

$$f = \frac{\Omega_0}{A_0} K_0 \left(\frac{\kappa_0}{2A_0} \right) + \frac{\kappa_0}{2A_0} K_1 \left(\frac{\kappa_0}{2A_0} \right), \quad (8.121b)$$

where K_ν is a modified Bessel function (Appendix C.7). For a flat circular-speed curve, $A = \frac{1}{2}\Omega$, $\kappa = \sqrt{2}\Omega$, and $f = 0.923$; for a Keplerian curve $f = 1.680$.

This result is based on linear perturbation theory, and hence is valid only when the epicycle amplitude X induced by the encounter is much smaller than the impact parameter b . Requiring that $X \lesssim b$ implies that

$$\frac{b}{r_J} \gtrsim \left(\frac{4f\Omega_0}{\kappa_0} \right)^{1/3} \approx 1, \quad (8.122)$$

where we have used equation (8.118) for the Jacobi radius. Thus, equations (8.121) are valid for encounters with impact parameters that are large compared to the Jacobi radius, and they show that the epicycle energy excited in an encounter falls off as b^{-4} . On the other hand, Figure 8.12 shows that encounters at small impact parameters, $b \lesssim r_J$, simply switch the star from one circular orbit to another, with no sensible increase in the star's random velocity. These considerations imply that the strongest encounters

have $b \sim r_J$, and hence justify our treatment of the molecular cloud as a point mass: since the typical cloud radius $R \sim 10$ pc is much smaller than its Jacobi radius (8.118), the non-zero cloud size has little influence on the rate of disk heating.

The speeds with which the stars in Figure 8.12 approach the scattering cloud are due entirely to the differential rotation of the galactic disk, $\dot{y} = -2A_0b$. Once the stars have acquired non-zero epicycle energy, we have to consider two types of encounter. For large impact parameters b or small epicycle amplitudes X , the approach speed is still dominated by the contribution from differential rotation (**shear-dominated** encounters), but at impact parameters $b \lesssim \kappa_0 X/A_0$, the encounter geometry will be determined mainly by the star's epicyclic motion (**dispersion-dominated** encounters).

In contrast to shear-dominated encounters, dispersion-dominated encounters usually *can* be treated using the impulse approximation. In this approximation, the magnitude of the velocity change $\Delta \mathbf{v}$ in a single encounter is proportional to v^{-1} where v is the encounter velocity. Thus the average change in epicycle energy E_x is proportional to v^{-2} , and since the number of encounters per unit time is proportional to v and $v^2 \sim E_x$ we expect (cf. Problem 8.20)

$$\frac{dE_x}{dt} \propto v^{-1} \propto \frac{1}{\sqrt{E_x}}. \quad (8.123)$$

Integrating this result we find that $E_x \propto t^{2/3}$ so the velocity dispersion $v \propto t^\alpha$ with $\alpha = \frac{1}{3}$. This simple calculation somewhat overestimates the rate of growth of the dispersion, since the thickness of the stellar disk is larger than the thickness of the cloud layer, so stars spend a smaller and smaller fraction of their time in the cloud layer as the vertical dispersion, and the resulting thickness of the stellar disk, continue to grow. The numerical calculations described below are consistent with this argument, suggesting that $\alpha \simeq 0.2$ – 0.25 for heating by molecular clouds. This value is too low to match the observations shown in Figure 8.11—just the opposite problem from MACHO-dominated heating, which gives an exponent that is too large (eq. 7.102).

We have shown that encounters with clouds “heat” the disk, in the sense that the mean epicycle energy increases with time. It is instructive to ask where this energy comes from, since the total energy or Jacobi integral of the star (eq. 8.104) is conserved during each encounter. The first of equations (8.103) shows that in a razor-thin disk the difference in epicycle energy E_x before and after the encounter is equal and opposite to the difference in $\frac{1}{2}A_0L^2/B_0$; in most galactic potentials $A_0/B_0 < 0$ so we conclude that an increase in E_x is accompanied by an increase in $|L|$ or in $|x_g|$ (eq. 8.102), where x_g is the difference in radius between the guiding center of the stellar orbit and the orbital radius of the molecular cloud. In words, the gravitational interaction with the cloud *repels* the stars, in the sense that their mean

orbital radius is shifted away from the cloud. Thus the energy to heat the disk comes from a redistribution of the surface mass density of the stars in the disk; the cloud acts as a catalyst to expedite this redistribution of energy but does not contribute any of its own energy to the disk heating.

Numerous authors have estimated the rate at which star-cloud encounters heat disks (Spitzer & Schwarzschild 1951, 1953; Jenkins 1992; Hänninen & Flynn 2002). The best estimate of the number density and masses of molecular clouds in the solar neighborhood leads to a rate of velocity-dispersion growth that is too small by a factor of two or more; but the heating rate is likely to be enhanced by the swing-amplified response or spiral wake induced in the stellar disk by the gravitational field from the molecular cloud, which can be several times more massive than the cloud itself (Julian & Toomre 1966; Julian 1967).

These studies also show that the predicted ratio σ_z/σ_R of the vertical and radial dispersions is $\simeq 0.6$ (Ida, Kokubo, & Makino 1993), not far from the observed ratio of 0.5. However, the predicted age-velocity dispersion relation is approximately a power law, $\sigma_i \propto t^\alpha$, with exponent $\alpha \simeq 0.2$ – 0.25 . This is significantly lower than the observed exponent, which is 0.3 for σ_R and σ_ϕ and even larger for σ_z (Figure 8.11). This result suggests that molecular clouds are unlikely to be the primary cause of disk heating.

8.4.2 Scattering of disk stars by spiral arms

The disks of spiral galaxies are far from smooth. Gas, dust, and young stars are always concentrated into spiral arms. Spiral features are also found in the old stars that make up most of the mass of galactic disks (§6.1.2), so it is natural to ask whether the gravitational fields of spiral features, like the fields from molecular clouds, are able to heat galactic disks (Barbanis & Woltjer 1967).

Consider a weak spiral potential with pattern speed Ω_p ,

$$\Phi_s(R, \phi, t) = \epsilon F(R) \cos[f(R) + m\phi - \Omega_p t], \quad (8.124)$$

where $\epsilon \ll 1$. To illustrate the effect of this potential on a stellar orbit, we shall make two assumptions that simplify the algebra but still retain most of the important dynamics: (i) we work in the sheared-sheet approximation (§8.3.2); (ii) we consider only tightly wound spirals, for which the wavenumber $k \equiv df/dR$ is large compared to $1/R$ (eq. 6.4).

The sheared-sheet approximation is valid in a neighborhood of the disk centered at a point $[R_0, \phi_0(t)]$ that rotates at the circular angular speed $\dot{\phi}_0 = \Omega_0 = \Omega(R_0)$. We expand the spiral potential in a Taylor series around this point, using the coordinates $x = R \cos(\phi - \phi_0) - R_0 \simeq R - R_0$ and $y = R \sin(\phi - \phi_0) \simeq R_0(\phi - \phi_0)$. In this neighborhood, we can approximate

the shape function as $f(R) \simeq f(R_0) + kx$. Since the amplitude $F(R)$ varies slowly, it can be replaced by a constant, $F_0 \equiv F(R_0)$. Thus we have

$$\begin{aligned}\Phi_s(x, y, t) &= \epsilon F_0 \cos[f(R_0) + kx + my/R_0 + m\phi_0 - \Omega_p t] \\ &= \epsilon F_0 \cos[kx + my/R_0 + m(\Omega_0 - \Omega_p)t + \text{constant}].\end{aligned}\quad (8.125)$$

We now substitute this potential into the equations of motion (8.97) of the sheared sheet, neglecting motion in the z -direction perpendicular to the disk plane:

$$\begin{aligned}\ddot{x} - 2\Omega_0 \dot{y} - 4\Omega_0 A_0 x &= \epsilon k F_0 \sin[kx + my/R_0 + m(\Omega_0 - \Omega_p)t + \text{constant}]; \\ \ddot{y} + 2\Omega_0 \dot{x} &= \frac{\epsilon m}{R_0} F_0 \sin[kx + my/R_0 + m(\Omega_0 - \Omega_p)t + \text{constant}].\end{aligned}$$

Since the wave is assumed to be tightly wound, its pitch angle is small so $|k| \gg m/R_0$ (eq. 6.7). Thus the right side of the second equation is much smaller than the corresponding term in the first, and can be neglected. The second equation can then be integrated to yield $\dot{y} + 2\Omega_0 x = \text{constant}$, and this can be substituted into the first equation to give

$$\ddot{x} + \kappa_0^2 x + \text{constant} = \epsilon k F_0 \sin[kx + my/R_0 + m(\Omega_0 - \Omega_p)t + \text{constant}]; \quad (8.126)$$

here κ_0 is the epicycle frequency (8.100). The constant on the left side can be dropped, since it can be absorbed by a shift in the origin of the x -coordinate.

In the absence of a spiral ($F_0 = 0$) the solution to this equation is given by equations (8.99); we shall assume that the unperturbed motion is circular, so the trajectory is $\mathbf{x}_0(t) = (x_g, y_{g0} - 2A_0 x_g t)$. Now consider how this motion is modified by the weak spiral potential on the right side of equation (8.126). We write the trajectory as $\mathbf{x}(t) = \mathbf{x}_0(t) + \epsilon \mathbf{x}_1(t)$, where $\epsilon \mathbf{x}_1(t)$ is the perturbation induced by the spiral. Then the terms of order ϵ in equation (8.126) yield

$$\begin{aligned}\ddot{x}_1 + \kappa_0^2 x_1 &= k F_0 \sin[kx_g + m(y_{g0} - 2A_0 x_g t)/R_0 + m(\Omega_0 - \Omega_p)t + \text{constant}] \\ &= k F_0 \sin(kx_g + \omega t + c).\end{aligned}\quad (8.127)$$

In the last expression we have absorbed y_{g0} in the constant c , and set $\omega = m(\Omega_0 - 2A_0 x_g/R_0 - \Omega_p)$; this is the frequency at which the unperturbed orbit encounters successive crests of the spiral potential.

This equation can be solved to yield

$$x_1(t) = \frac{k F_0}{\kappa_0^2 - \omega^2} \sin(kx_g + \omega t + c). \quad (8.128)$$

The solution diverges when $\omega = \pm \kappa_0$. These points can be thought of as the Lindblad resonances of the sheared sheet: at these locations, like the Lindblad resonances in a disk, the frequency of excitation by the spiral potential

coincides with the frequency κ_0 of the particle's natural radial oscillation. This result is a close analog of equation (3.148), which was derived in the context of weak bars.

Equation (8.128) shows that the spiral potential imposes a forced radial oscillation on the star but does not lead to any steady growth in the radial oscillation $x_1(t)$. In other words, *a spiral potential with a fixed pattern speed cannot heat the disk*, except perhaps at the Lindblad resonances where our simple derivation fails.

This result implies that disk heating requires *transitory* rather than steady spiral patterns. To illustrate this, let us multiply the potential (8.124) or (8.125) by a Gaussian function of time, $p(t) = (2\pi s^2)^{-1/2} \exp(-\frac{1}{2}t^2/s^2)$. Equation (8.127) is thereby modified to read

$$\ddot{x}_1 + \kappa_0^2 x_1 = kF_0 p(t) \sin(kx_g + \omega t + c), \quad (8.129)$$

which has the solution

$$x_1(t) = X_1 \cos(\kappa_0 t + \alpha_1) + \frac{kF_0}{\kappa_0} \int_{-\infty}^t dt' p(t') \sin(kx_g + \omega t' + c) \sin[\kappa_0(t - t')], \quad (8.130)$$

where X_1 and α_1 are arbitrary constants. Inserting the chosen form for $p(t)$ and setting the amplitude X_1 of the free oscillation to zero, we obtain

$$x_1(t \rightarrow \infty) = \frac{kF_0}{2\kappa_0} \left[\cos(\kappa_0 t - kx_g - c) e^{-s^2(\omega + \kappa_0)^2/2} - \cos(\kappa_0 t + kx_g + c) e^{-s^2(\omega - \kappa_0)^2/2} \right]. \quad (8.131)$$

Thus the transitory spiral pattern has induced a permanent epicyclic oscillation. When the characteristic duration of the transient, s , is much greater than the orbital period, the induced epicycle amplitude is strongly peaked near the Lindblad resonances $\omega = \pm\kappa_0$. On the other hand, when the duration of the transient is short, the arguments of the exponential are small and epicycle motion is induced over a wide range of radii in the disk.

This example shows that the ability of spiral structure to heat the disk is strongly dependent on its temporal structure. According to the Lin–Shu hypothesis (§6.1), in which spiral structure is a stationary wave with a single, well-defined pattern speed, disk heating is negligible except at the Lindblad resonances. In such models the disk can be heated over a wide range of radii only if the pattern speed evolves with time, so the Lindblad resonances slowly sweep across most of the disk. On the other hand, if the spiral structure is transient, the whole disk can be heated—this situation is likely to occur in flocculent spirals, intermediate-scale spirals, or grand-design spirals excited by recent encounters.

Let us suppose that a given star is subjected to N independent transient perturbations. Each transient induces an epicyclic motion whose radial

component can be written in the form $x_1(t) = a_i \cos(\kappa_0 t + \alpha_i)$, $i = 1, \dots, N$, where a_i and α_i are given by equations similar to (8.131). After N transients,

$$\begin{aligned} x_1(t) &= \sum_{i=1}^N a_i \cos(\kappa_0 t + \alpha_i) \\ &= \left(\sum a_i \cos \alpha_i \right) \cos \kappa_0 t - \left(\sum a_i \sin \alpha_i \right) \sin \kappa_0 t \\ &\equiv a_f \cos(\kappa_0 t + \alpha_f) \end{aligned} \quad (8.132)$$

where

$$a_f^2 = \left(\sum a_i \cos \alpha_i \right)^2 + \left(\sum a_i \sin \alpha_i \right)^2 = \sum_{i,j=1}^N a_i a_j \cos(\alpha_i - \alpha_j). \quad (8.133)$$

Since the transients are uncorrelated, the phases of the epicyclic oscillations that they induce are also uncorrelated. Hence on average $\cos(\alpha_i - \alpha_j)$ will be zero when $i \neq j$, and the only terms in the sum that contribute to the final amplitude a_f will be those with $i = j$. Thus $a_f^2 \simeq N \langle a^2 \rangle$, where $\langle a^2 \rangle$ is the mean-square amplitude induced by a single transient. If the rate of occurrence and the strength of new transients are independent of time, we conclude that a_f^2 , and hence the squared velocity dispersion v^2 , should grow linearly with time. In other words, $v \propto t^\alpha$, where $\alpha = 0.5$. This behavior holds only so long as a_f is not too large: once the epicycle size becomes comparable to the radial wavelength of the spiral arms, the effects of the spiral tend to average out over the epicycle period, so the heating is weaker—this is the same effect that leads to the reduction factor in the WKB dispersion relation for spiral waves, as described in §6.2.2d. Estimates of the heating rate at larger amplitudes can be obtained using the Fokker–Planck equation (Jenkins & Binney 1990; Jenkins 1992) or numerical simulations (De Simone, Wu, & Tremaine 2004); these calculations show that α can vary between 0.25 and 0.5 depending on the properties of the spiral transients (duration, strength, pitch angle, etc.). This range of α is nicely consistent with the observed exponent for the growth of the radial dispersion, $\alpha \simeq 0.3$, and provides a substantially better fit than the values predicted for heating by molecular clouds.

The radial and azimuthal velocity dispersions are related by equation (8.117), so the exponent in the age-velocity dispersion relation must be the same for these two axes of the velocity ellipsoid. However, spiral structure cannot excite velocities in the z -direction effectively, since its spatial and temporal scales are much larger than the amplitude or period of oscillations perpendicular to the disk plane. Thus, scattering by spiral arms cannot explain the relation between age and the z -velocity dispersion σ_z . Probably gravitational scattering by molecular clouds redistributes the radial and azimuthal velocities into the direction perpendicular to the plane (Carlberg

1987; Jenkins & Binney 1990). Thus molecular clouds are responsible for the shape, but not the size, of the velocity ellipsoid.

Transient spiral arms have other interesting consequences for the distribution of disk stars. Strong transients can produce long-lived clumps of stars in velocity space, sometimes called star streams or moving groups (see page 327 and De Simone, Wu, & Tremaine 2004). Spiral waves also redistribute the angular momenta of disk stars, leading to substantial inward and outward migration of individual stars over the lifetime of the Galaxy (Sellwood & Binney 2002).

8.4.3 Summary

There is little doubt that irregularities in the Galaxy's gravitational field heat the disk and thereby determine the velocity distribution of disk stars. It is less clear *which* irregularities dominate this process. We have discussed the influence of hypothetical massive objects in the dark halo (MACHOs), molecular clouds, and transient spiral arms. Other possibilities include merging satellite galaxies (Walker, Mihos, & Hernquist 1996; Velázquez & White 1999), substructure in the dark halo (Benson et al. 2004), or the Galactic bar (Kalnajs 1991; Dehnen 2000a). The simplest explanation that appears to be consistent with most of the observations is the combined effects of spiral transients and molecular clouds.

8.5 Mergers

So far we have investigated galaxy mergers and encounters through limiting cases that are analytically tractable. For example, minor mergers occur through dynamical friction (§8.1), which leads to gradual orbital decay, and as the orbit shrinks tidal forces and tidal shocks (§§8.2 and 8.3) become stronger and stronger, until either the small galaxy is completely disrupted or its core comes to rest at the center of the larger galaxy.

In major mergers the physical processes are qualitatively similar, but harder to quantify. The relative velocity of the centers of mass of the two galaxies is converted into randomly directed velocities of their individual stars—the same process as dynamical friction—but the conversion is so rapid that the galaxies merge into a single steady-state system within a few crossing times. Thus, numerical simulations such as the one shown in Figure 8.1, rather than analytic arguments, are the primary tool for understanding major mergers.

In this section we shall focus on features of major mergers that have direct observational consequences; these are important because they provide the “smoking gun” that enables us to identify galaxies that are participating in ongoing mergers, and thus to explore the physics of mergers. Reviews

of interacting and merging galaxies are given by Barnes & Hernquist (1992) and Kennicutt, Schweizer, & Barnes (1998).

8.5.1 Peculiar galaxies

A small fraction of galaxies are found in a highly disturbed state. The importance of these puzzling systems was emphasized by Arp (1966), who compiled an *Atlas of Peculiar Galaxies* containing over 300 such objects (see also Arp & Madore 1987). Arp argued that “if we could analyze a galaxy in a laboratory, we would deform it, shock it, probe it, in order to discover its properties” and that the peculiarities of the galaxies in his atlas offered a range of experiments on galaxies furnished to us by nature, which we should learn from. At one time it was widely believed that unusual systems of this kind were exploding galaxies or galaxies with very strong magnetic fields, but by the early 1970s it became clear that most are actually colliding systems, and that many of these collisions will result in mergers.

Figure 8.13 shows the pair of galaxies NGC 4038/4039 from the Arp atlas. This system consists of overlapping blobs of light from which two curved tails of much lower surface brightness emerge, giving rise to its common name “the Antennae.” From end to end, the tails span over 100 kpc. Can this striking morphology be the signature of a merger? In a classic paper, Toomre & Toomre (1972) showed that this is indeed the case. The Toomres studied encounters between disks of massless particles orbiting around point masses: even with this grossly oversimplified model of a galaxy—the disk has no self-gravity, there is no massive halo, and the disk circular-speed curve is Keplerian rather than flat—they were able to show that for a suitable choice of initial conditions, it is possible to find a pair of colliding stellar systems that is remarkably similar to Figure 8.13. We show their model in Figure 8.14. More accurate models that include the self-gravity of the disk and a massive halo largely confirm the Toomres’ conclusions (Barnes 1988; Dubinski, Mihos, & Hernquist 1999).

The Toomres’ model predicted the line-of-sight velocity at each point in the system. The observed velocities were found to be in complete agreement with the model, and show that the point of closest approach of the two galaxies, when the tails were launched, occurred 0.5 Gyr ago (Hibbard et al. 2001).

The tails seen in the Antennae differ in one important respect from the tidal streamers discussed in §8.3.3. The streamers discussed in that section are composed of stars stripped from small satellites of much larger stellar systems; the streamers are narrow because the satellite is small. In contrast, the two merging systems in the Antennae have comparable size; the tails are narrow because the stars come from cold stellar systems—the disks of the two merging spiral galaxies—so all the stars near a given location have nearly the same initial velocity. Mergers of hot stellar systems of comparable size do not generate narrow tidal tails.

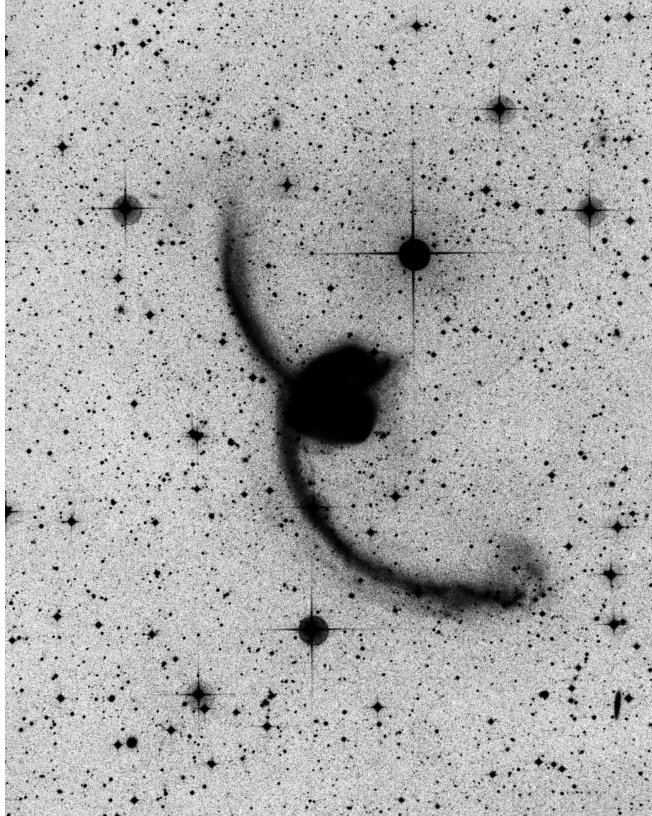


Figure 8.13 The interacting galaxies NGC 4038 and NGC 4039, the “Antennae.” This is an overexposed image to emphasize the low surface-brightness tidal tails. The distance from the overlapping blobs at the center to the bright star above and to the right of them is 40 kpc. Courtesy of D. F. Malin and the Anglo–Australian Telescope Board.

Another route to the same conclusion is through the collisionless Boltzmann equation, which shows that the density of stars in phase space is conserved (eq. 4.10). A long-lived tidal tail or streamer must have high phase-space density, since the spatial density must be high if the tail is to be visible against the background galaxy, and the velocity dispersion must be low if it is not to disperse quickly. Thus the progenitor of the tidal tail or streamer must have high phase-space density, a condition that is satisfied by both satellite stellar systems (because their spatial densities are high and their velocity dispersions are low compared to the larger host galaxy) and disks (because the velocity dispersion is low).

Another galaxy with prominent tidal tails that is almost certainly an ongoing merger is NGC 4676 (“the Mice”), shown in Figure 8.15.

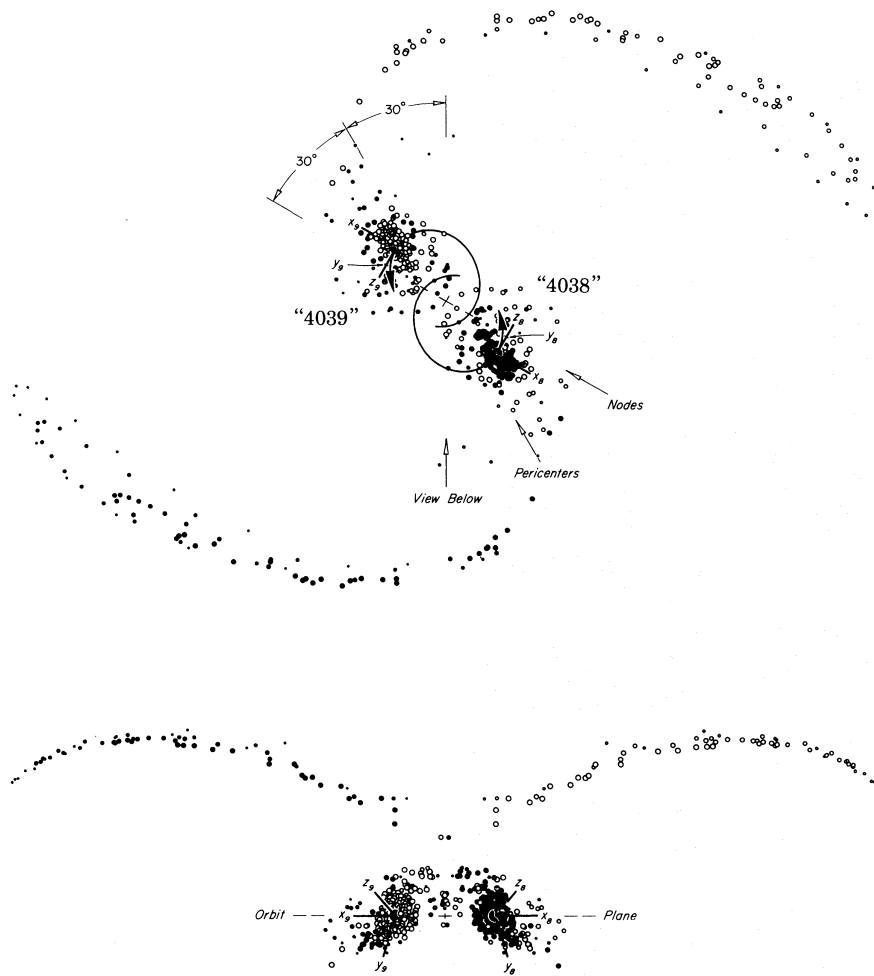


Figure 8.14 A model of the NGC 4038/4039 pair by Toomre & Toomre (1972). Reproduced by permission of *The Astrophysical Journal*.

8.5.2 Grand-design spirals

We have seen in Chapter 6 that grand-design spirals such as M51 (Plate 1) or M81 (Plate 8) often have companion galaxies nearby, and that the gravitational forces from an encounter with a companion can excite a strong but transitory spiral pattern (Figure 6.26). In most cases the orbit of the companion galaxy that excited the spiral will decay by dynamical friction, so the two galaxies are likely to merge in the future. Thus many of the most beautiful and striking spiral galaxies in the sky are likely to be the product of major mergers.

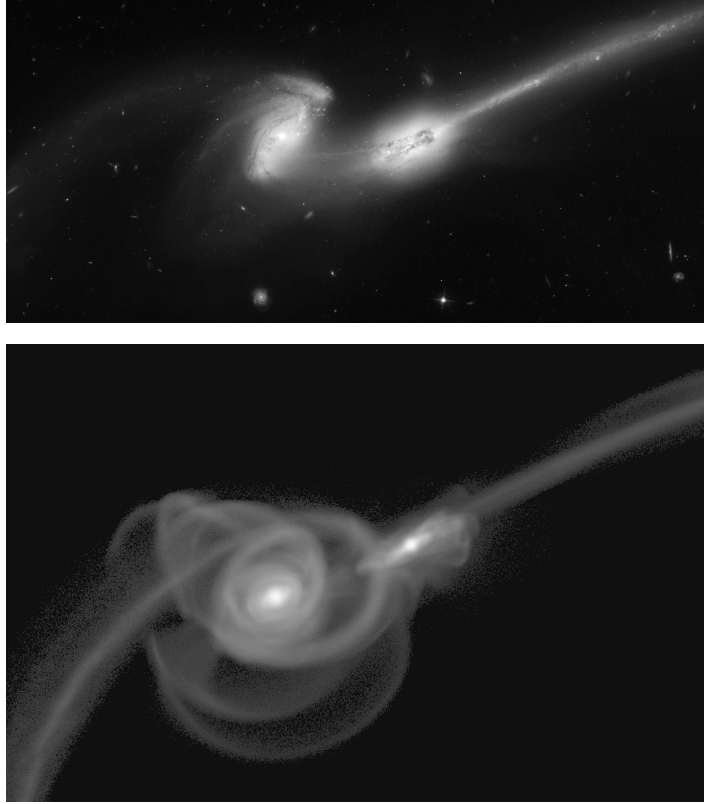


Figure 8.15 The Mice, NGC 4676, a pair of interacting galaxies at a distance of 95 Mpc. Top: optical image from the Hubble Space Telescope. Bottom: an N-body model. Credit for HST image: NASA, H. Ford (JHU), G. Illingworth (UCSC/LO), M. Clampin (STScI), G. Hartig (STScI), the ACS Science Team, and ESA. Credit for N-body model: J. Dubinski (Dubinski & Farah 2006).

8.5.3 Ring galaxies

A handful of galaxies exhibit a distinctive morphology consisting of a luminous ring of young stars that is both rotating and expanding, usually with one or more compact companion galaxies nearby. Figure 8.16 shows one example, the “Cartwheel Galaxy.” These remarkable systems are known as **ring galaxies** or sometimes **collisional ring galaxies** (Appleton & Struck-Marcell 1996).¹²

Ring galaxies form when a disk galaxy collides head-on with another system (Lynds & Toomre 1976). The collision excites a radially expanding

¹² These are distinct from the prominent rings that are seen in some barred galaxies, which are thought to arise from rapid star formation in gas that is in resonance with the bar (see §6.5.2d and Buta 1995).

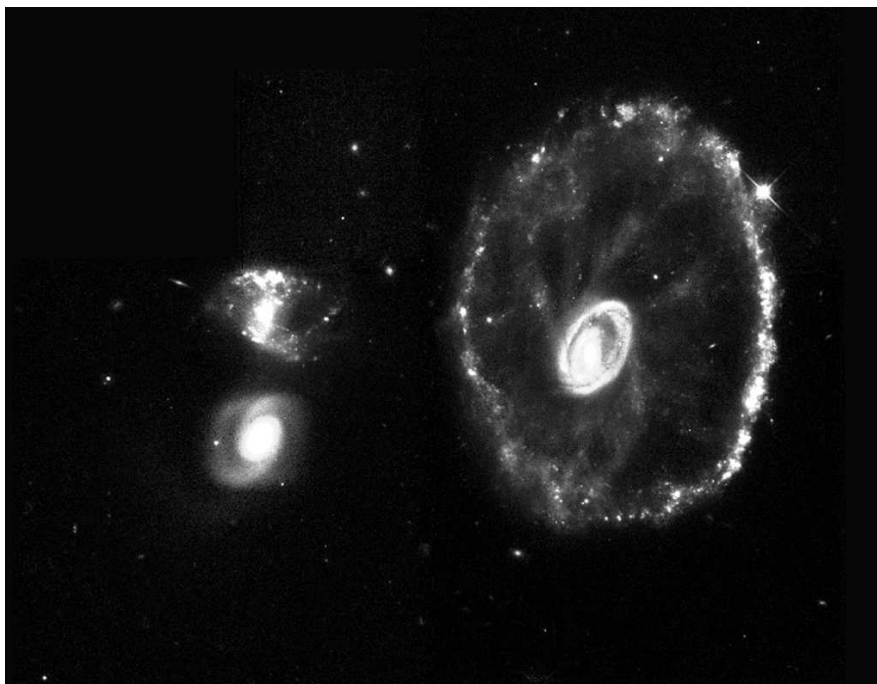


Figure 8.16 The Cartwheel galaxy, the prototypical ring galaxy, at a distance of 125 Mpc. The ring diameter is about 45 kpc. The lower of the two compact stellar systems to the left of the ring—both members of the same group of galaxies as the Cartwheel—is probably responsible for the ring structure. Credit: K. Borne (George Mason University) and NASA.

density wave that triggers star formation in the disk as it passes. The compact systems are the surviving central cores of the colliding galaxies. Ring galaxies are rare—about one in 10^4 galaxies—because they are short-lived, and because they are produced only in collisions with near-zero impact parameter.

We can use the impulse approximation to develop an instructive model of this process, even though this approximation may not hold for all ring galaxies. Consider a singular isothermal sphere that contains a rotating, cold, disk of test particles in the plane $z = 0$, and suppose that it collides with a second singular isothermal sphere having the same circular speed v_c , traveling along the z axis with relative speed $V \gg v_c$. The gravitational potential of each sphere is $\Phi(r) = v_c^2 \ln r$, and in Problem 8.7 it is shown that the change Δv_R in the velocity of a disk star at initial radius R is then

$$\Delta v_R = -2R \frac{v_c^2}{V} \int_R^\infty \frac{dr}{r\sqrt{r^2 - R^2}} = -\frac{\pi v_c^2}{V}. \quad (8.134)$$

If V/v_c is sufficiently large, the velocity impulse $\Delta v_R/v_c$ will be small, so

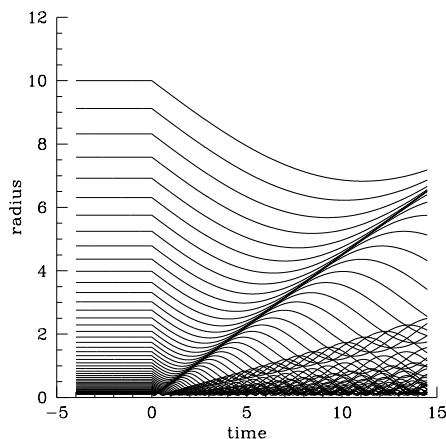


Figure 8.17 The evolution of the radii of particles in a disk after a head-on encounter, as described by equation (8.135), following Toomre (1978). The ratio $V/v_c = 7$.

(i) we may neglect the changes in the target's potential that are generated by the collision; (ii) we can describe the subsequent motion of the disk stars using the epicycle approximation. If the collision is assumed to occur at time $t = 0$, the radius of a star that is initially at R_0 is given by the solution of equation (3.78a) that satisfies the initial conditions $x(0) = R - R_0 = 0$, $\dot{x}(0) = \Delta v_R$:

$$R(R_0, t) = R_0 + \frac{\Delta v_R}{\kappa_0} \sin(\kappa_0 t) \quad (t > 0). \quad (8.135)$$

Here κ_0 is the epicycle frequency at R_0 , which is given by $\kappa_0^2 = 2v_c^2/R_0^2$ (eq. 3.79a).

The evolution of the radii of particles in the disk is shown in Figure 8.17. The crowding of particle orbits gives rise to strong axisymmetric density waves that propagate out through the disk. The point of maximum compression of the particle orbits is likely to be a region of enhanced star formation, which we identify with the ring of luminous young stars. The outward propagation of the density waves implies that the region inside the ring should contain older, redder stars that were formed when the ring was smaller, and such radial color gradients are indeed observed in several ring galaxies.

In practice, the collision of two galaxies is never precisely along the z axis of the disk, as assumed in this simple model. However, numerical experiments such as those shown in Figure 8.18 show that whenever an intruder passes close to the center of the target disk on a trajectory that is angled by less than about 30° from the symmetry axis of the disk, a striking ring is generated. If the intruder misses the center of the target just slightly, the dense center of the target galaxy is displaced from the center of the ring, as is observed in the Cartwheel Galaxy.

8.5.4 Shells and other fine structure

Figure 8.19 shows images of NGC 3923 and NGC 1344, which exhibit arclike

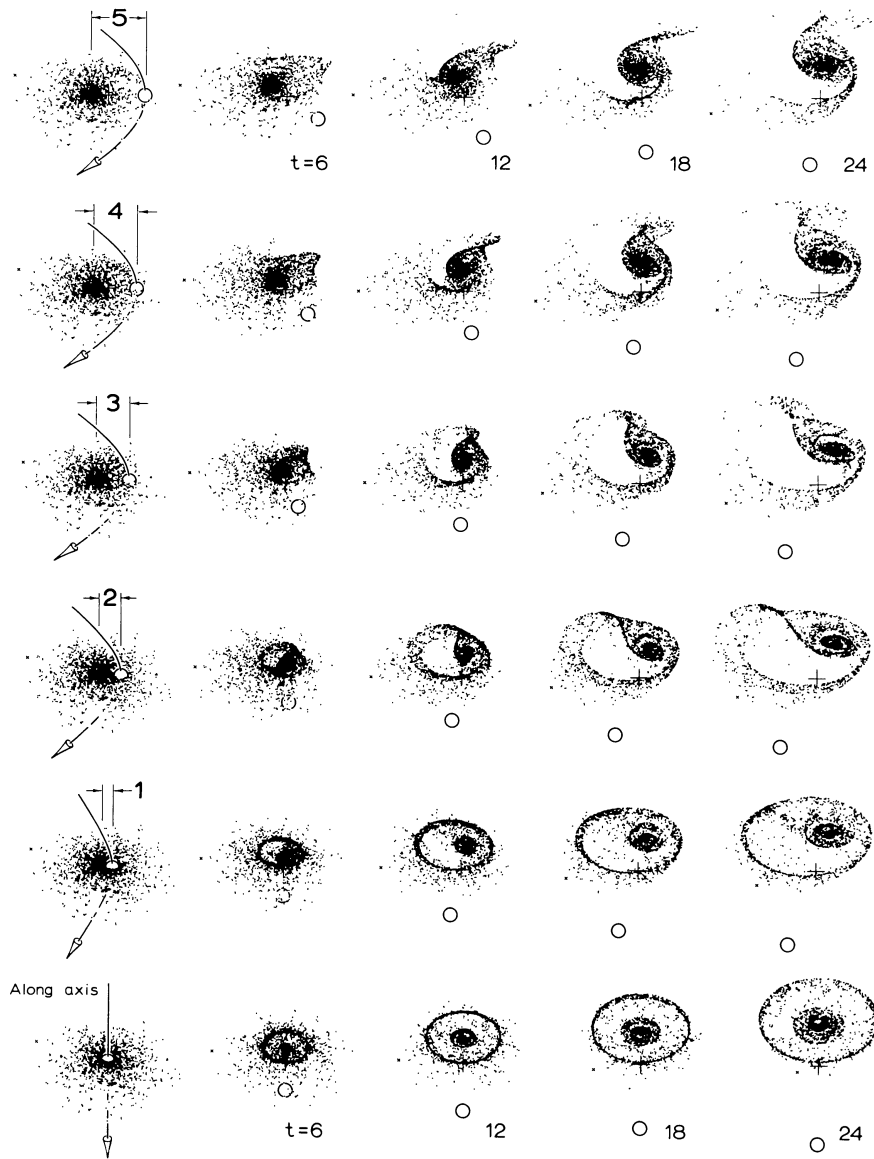


Figure 8.18 Six encounters between a disk of test particles orbiting a point mass and an intruder of half the mass, marked by an open circle. The relative orbit is parabolic, and the system is viewed from 45° above the disk. A ring is generated when the impact parameter is small compared to the size of the disk (bottom three rows). From Toomre (1978), with permission of Springer Science and Business Media.

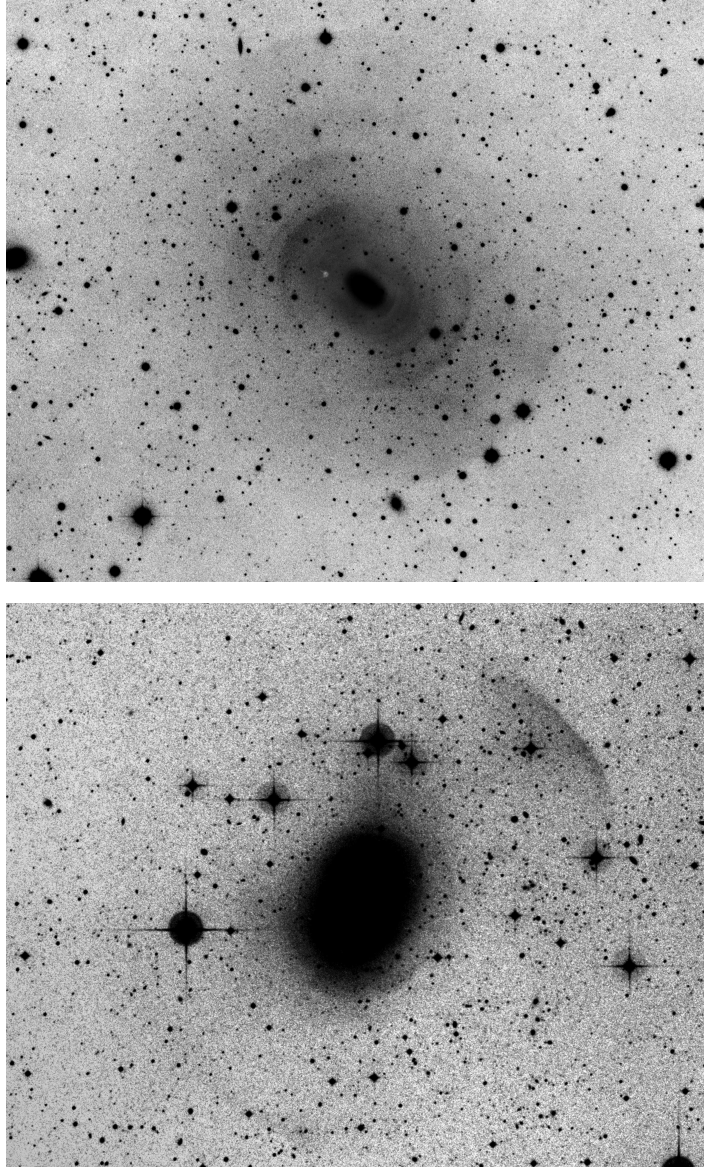


Figure 8.19 The giant elliptical galaxies NGC 3923 (top) and NGC 1344 (bottom) are surrounded by faint shells. The images have been processed to accentuate the shells using a high-pass spatial filter. Courtesy of David Malin, © Anglo–Australian Observatory.

shells in the surface brightness on both sides of the galaxy; careful analysis reveals over 20 such shells in this galaxy. The fraction of otherwise smooth

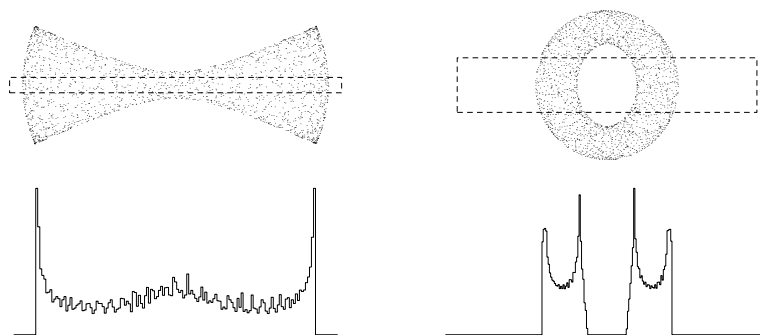


Figure 8.20 The box and loop orbits shown in Figure 3.8. The upper figures show stars at randomly chosen phases on the orbits to give a visual impression of the density distribution in the orbit, and the lower figures show the number of stars as a function of horizontal position within the boxes marked by dashed lines. Note the sharp cusps at the turning points of the orbits.

galaxies (ellipticals and lenticulars) that exhibit shells can be as large as 30–50%, depending on how closely one looks. Shells may also be present in spiral galaxies, but are camouflaged by spiral structure, dust, and irregular star formation in the disk. Spectra show that the shells are composed of stars, not gas. As is often the case in astronomy, the most famous examples of this phenomenon are atypical. The shells in NGC 3923 are exceptionally sharp and numerous, aligned with the major axis of the galaxy, and interleaved in radius (the shells from one side alternate in radius with those from the other side). In contrast, most shell galaxies contain $\lesssim 3$ detectable shells, and these are fainter, more diffuse, and have a less regular geometrical structure.

Other types of fine structure are also seen in elliptical galaxies (Kennicutt, Schweizer, & Barnes 1998), and are given names such as “loops,” “ripples,” “plumes,” “jets,” “X-structures,” etc. The tidal streamers described in §8.3.3 are also a kind of fine structure.

Most fine structure in galaxies is formed by the same process that forms tidal streamers, namely the disruption of a stellar system that has high phase-space density—either a small, hot galaxy or a large, cold one.

For example, consider the fate of the stars in a small satellite that is disrupted by a host galaxy. Initially the disrupted stars will form a tidal streamer, but eventually the streamer will disperse. If the host potential is regular, the stars will finally spread into a cloud of particles that have similar actions but uniformly distributed angles. Such a cloud gives rise to surface-density distributions such as those shown in Figure 8.20 for box and loop orbits. A box orbit produces an X-shaped structure, while a loop orbit produces an annulus with sharp edges at its inner and outer turning points.

In projection, the edges of these particle distributions can appear as shells, as indicated in the figure.

Another simple example is the disruption of a disk galaxy that is on a radial orbit in a spherical potential (Quinn 1984). This process can be explored by releasing a cloud of test particles in a fixed potential (Figure 8.21). Since the angular momentum of the test particles is nearly zero, the motion can be followed in the two-dimensional (r, v_r) plane. Shells are formed at the turning points of the orbits, and they are interleaved in radius like the shells observed in NGC 3923. The rather special circumstances of the encounter (radial orbit, spherical host potential) are consistent with the observation that most shell galaxies do not exhibit the regular geometrical structure seen in this example.

The much more common case of the disruption of a galaxy on a non-radial orbit can also produce shells, such as those shown in Figure 8.22, but now the shells display the more complex geometry that is encountered in most shell galaxies (Hernquist & Quinn 1988, 1989).

More generally, shells arise when stars are confined to a subspace of lower dimensionality than the full six-dimensional phase space. The projection of this smooth manifold onto the two-dimensional plane of the sky gives rise to caustics, which can be classified using catastrophe theory (Tremaine 1999).

8.5.5 Starbursts

So far we have focused on the effects of mergers on a galaxy's stars, but the effects on its gas—if the galaxy has a gas disk—are even more dramatic. As Toomre & Toomre (1972) wrote, “Would not the violent mechanical agitation of a close tidal encounter—let alone an actual merger—already tend to bring *deep* into a galaxy a fairly *sudden* supply of fresh fuel in the form of interstellar material?” The Toomres' prescient question was answered by Larson & Tinsley (1978), who showed that many merger remnants had anomalous blue colors consistent with young, massive stars formed in a recent **starburst**—a short, intense period of rapid star formation at a rate far exceeding that of a normal galaxy. Since that time a wide variety of observations has confirmed that vigorous star formation occurs in merging galaxies. Among the most striking of these observations is the discovery of almost 10^3 blue objects in the Antennae (Figure 8.13), which appear to be young globular clusters formed in the merger (Whitmore & Schweizer 1995).

The observational link between mergers and star formation was cemented by the discovery of **starburst galaxies**. These are among the most luminous galaxies known, emitting up to $10^{12.5} L_{\odot}$, mostly at infrared wavelengths. This intense emission comes from young stars shrouded in dust and concentrated near the center of the galaxy. The emission is powered by extremely high star-formation rates, as large as $10^3 M_{\odot} \text{ yr}^{-1}$, compared to a few $M_{\odot} \text{ yr}^{-1}$ in galaxies like the Milky Way. Starburst galaxies usually

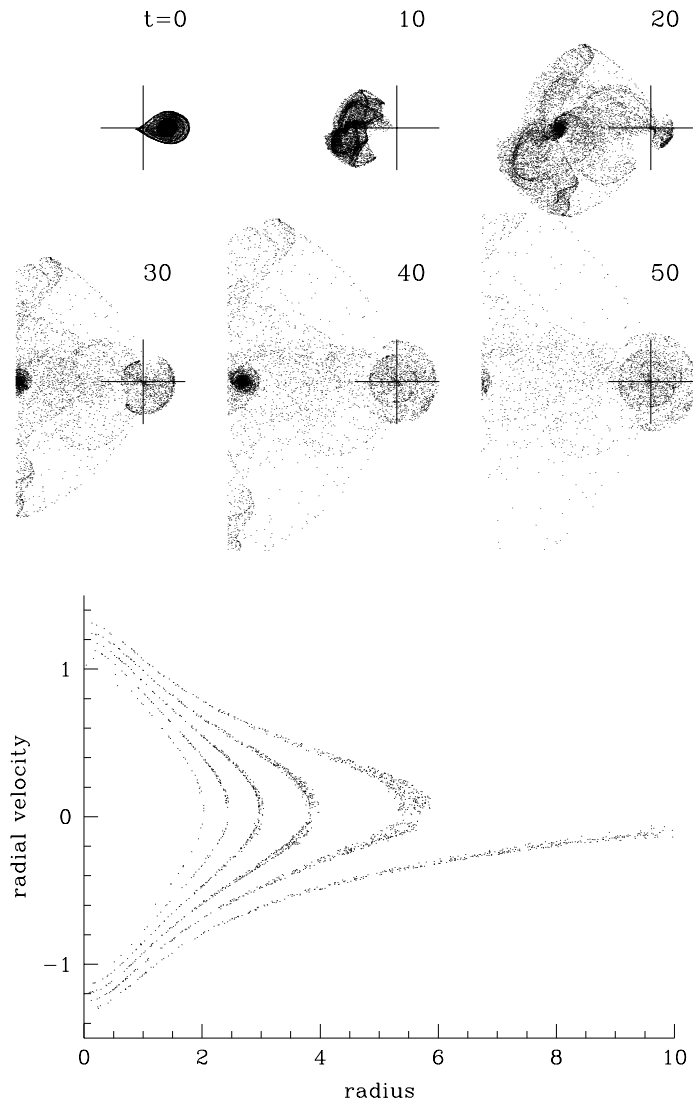


Figure 8.21 The disruption of a disk galaxy on a radial orbit in a spherical potential. The host galaxy is represented by a rigid, fixed Plummer potential (2.44a) with mass $M = 1$ and scale length $b = 1$, and the satellite has mass $m = 0.1$ and is modeled by a rigid Kuzmin-disk potential (2.68a) with scale length $a = 0.5b$, containing 10 000 test particles on initially circular orbits. Top: The cross marks the center of the host galaxy and the length of each arm of the cross is $5b$. The evolution is viewed from a direction normal to the plane of the disk galaxy (the x - y plane). The distribution of test particles is first shown just before the satellite reaches the center of the host, falling in from infinity along the positive x axis, and at intervals of 10 time units thereafter. Bottom: the projection of the test particles onto the radius-radial velocity plane at time $t = 50$ (Quinn 1984).

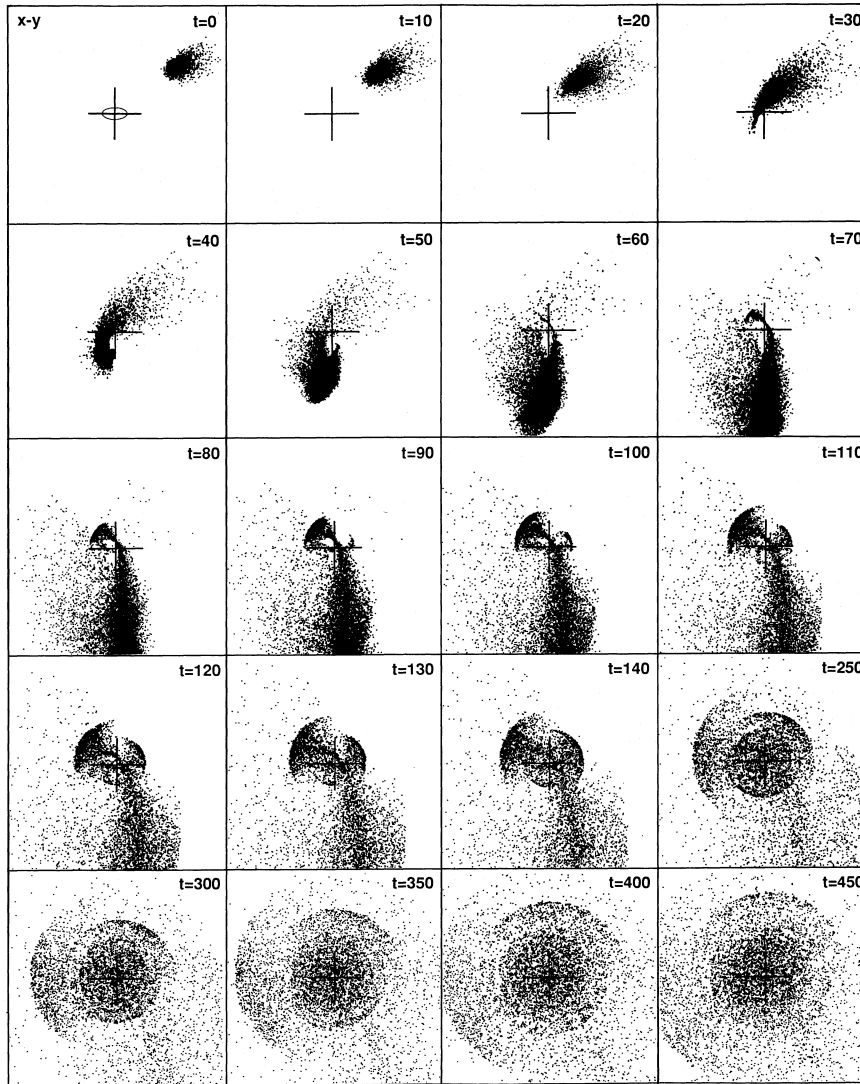


Figure 8.22 The formation of shells in the disruption of a small spherical galaxy on a non-radial orbit. The small ellipse in the first frame represents the approximate location and projected shape of the larger galaxy in the encounter. The simulation uses 20 000 test particles. From Hernquist & Quinn (1989), reproduced by permission of the AAS.

exhibit tidal streamers or other optical features indicative of a recent collision or merger (Sanders & Mirabel 1996; Kennicutt, Schweizer, & Barnes 1998; Kennicutt 1998). The presence of these short-lived features, and the rapid consumption of gas required by the high star-formation rate, imply

that starbursts last for only a few tens of Myr.

The reason for these very high star-formation rates is suggested by numerical simulations of mergers of disk galaxies that contain both gas and stars (Noguchi 1988; Mihos & Hernquist 1996). During the merger, both the gas and stars form strong bars. The gas bar leads the stellar bar, so the gravitational torque from the stars rapidly drains angular momentum from the gas. Remarkably, in a typical major merger the gas can lose up to 90% of its angular momentum in a fraction of an orbital period, thus settling into a dense rotating disk $\lesssim 0.5$ kpc across, in which star formation is likely to be extremely rapid.

8.5.6 The merger rate

The rate at which galaxies merge is a fundamental point of comparison between observations and models of structure formation. The merger rate can be determined from observations using two quite different methods.

The first method is to count the fraction of galaxies showing obvious features of an ongoing or recent major merger, such as tidal tails or starbursts, and combine this fraction with an estimate of how long such features last to determine the merger rate per galaxy. The first attempt of this kind was made by Toomre (1977b), who pointed out that about 10 of the ~ 4000 NGC galaxies¹³ show prominent tidal tails, and that these tails probably last no more than ~ 0.5 Gyr. Thus the rate of major mergers is probably about $10/4000/0.5 \text{ Gyr} \simeq 0.005 \text{ Gyr}^{-1}$ for a luminous galaxy. This estimate neglects two important biases, of opposite sign: first, not all major mergers yield visible tidal tails; second, mergers enhance the star-formation rate, so galaxies experiencing mergers are more luminous than quiescent galaxies, and hence will be over-represented in a flux-limited catalog like the NGC. A crude assumption is that these two biases cancel, leaving the original estimate of 0.005 major mergers per Gyr approximately correct.

A second approach is to count the fraction of galaxies with companions within a given radius, and combine this fraction with estimates of the rate of decay of the companion orbit by dynamical friction to obtain the merger rate (Tremaine 1981). The distribution of companions is described by the **galaxy-galaxy correlation function** $\xi(r)$, defined so that the probability of finding two galaxies in the volumes $d^3\mathbf{r}_1$ and $d^3\mathbf{r}_2$ separated by $r \equiv |\mathbf{r}_1 - \mathbf{r}_2|$ is

$$d^2p = n_0^2[1 + \xi(r)] d^3\mathbf{r}_1 d^3\mathbf{r}_2, \quad (8.136)$$

where $n_0 d^3\mathbf{r}$ is the probability of finding one galaxy in the volume $d^3\mathbf{r}$. Over a wide range of separations and luminosities, the galaxy-galaxy correlation

¹³ NGC stands for “New General Catalog,” a catalog of galaxies, nebulae, and clusters compiled by Dreyer in 1888, which was a revision and expansion of Herschel’s “General Catalog” of 1864. See BM Appendix B.

function can be described by a power law,

$$\xi(r) \simeq \left(\frac{r_0}{r}\right)^\gamma, \quad (8.137)$$

where (Hawkins et al. 2003)

$$r_0 = (7.2 \pm 0.4)h_7^{-1} \text{ Mpc} \quad ; \quad \gamma = 1.67 \pm 0.03, \quad (8.138)$$

over the range $1 \lesssim \xi \lesssim 10^3$. Thus the number density of companions at distance $r \ll r_0$ from a primary galaxy is

$$n(r) = n_0 \left(\frac{r_0}{r}\right)^\gamma. \quad (8.139)$$

The rate of decay of the orbital radius of a companion due to dynamical friction is given approximately by equation (8.16). Assuming that the primary galaxy and the companion have similar luminosity and velocity dispersion, σ , we have

$$\frac{dr}{dt} \simeq -0.4f\sigma, \quad (8.140)$$

where we have set the Coulomb logarithm $\ln \Lambda \simeq 1$ according to the arguments that follow equation (8.17), and $f \gtrsim 1$ is a correction factor that arises because the relative orbit is likely to be elongated rather than circular, which accelerates the decay. For each primary galaxy, the current of merging companions through radius r is then

$$\dot{N}(r) \simeq \frac{1}{2} \times 4\pi r^2 n(r) \left| \frac{dr}{dt} \right| \simeq 0.8\pi f n_0 \sigma r_0^\gamma r^{2-\gamma}; \quad (8.141)$$

the factor $\frac{1}{2}$ is needed to avoid counting each galaxy twice, once as a primary galaxy and once as a companion. In a steady state, the current should be independent of r , and this constant number would represent the merger rate; the actual weak dependence on r seen in equation (8.141), $\dot{N} \propto r^{2-\gamma} \propto r^{0.3}$, probably arises because our approximation that both galaxies are isothermal spheres is not very accurate. We shall equate the merger rate \dot{N}_{merge} to $\dot{N}(r_{\text{min}})$, where $r_{\text{min}} = 20h_7^{-1} \text{ kpc}$ is roughly the radius at which the stellar distributions of two luminous galaxies begin to merge.

The derived merger rate depends on the minimum luminosity of the companion galaxies that we consider—clearly, if we count minor mergers, the merger rate will be higher than if we count only major mergers. For the present estimate we shall consider only mergers of companions with luminosity $L > L_*$, where L_* is the characteristic Schechter luminosity defined by equation (1.18). Using that equation, the average number density of galaxies more luminous than L_* is

$$n_0(L > L_*) = \int_{L_*}^{\infty} dL \phi(L) = \phi_* \int_1^{\infty} dx x^\alpha e^{-x}. \quad (8.142)$$

For the parameters given after equation (1.18), we find $n_0(L > L_*) = 0.21\phi_* \simeq 1.0 \times 10^{-3} h_7^3 \text{Mpc}^{-3}$. According to the Faber–Jackson law (1.21), the dispersion of an L_* galaxy is $\sigma_* \simeq 200 \text{ km s}^{-1}$. Thus

$$\begin{aligned} \dot{N}_{\text{merge}}(L > L_*) &\simeq 0.8\pi f n_0(L > L_*) \sigma_* r_0^\gamma r_{\text{min}}^{2-\gamma} \\ &\approx 0.008 h_7 \frac{f}{2} \text{Gyr}^{-1}. \end{aligned} \quad (8.143)$$

The approximate agreement of this estimate with the rate 0.005 Gyr^{-1} obtained by Toomre (1977b) provides encouraging evidence that our understanding of the merging process is sound. More recent determinations of the merger rate (Conselice 2006) are also roughly consistent with these crude estimates.

Problems

8.1 [1] Two identical galaxies are initially at rest, at a large distance from one another. They are spherical, composed solely of identical stars, and their light distributions obey the Sérsic law (1.17) with Sérsic index m and effective radius R_e . The galaxies fall together and merge. If the merger product also satisfies the Sérsic law with the same index, what is its effective radius?

8.2 [1] The derivation of the dynamical friction formula (8.1) assumes that the subject system is a point mass, but in many cases of interest the subject system is an extended body, such as a star cluster or satellite galaxy, characterized by a half-mass radius r_h . If the point of closest approach of the field star to the center of the subject body is $\lesssim r_h$ then the deflection of the field-star orbit, and its contribution to the drag force, will be smaller than if the subject body were a point of the same total mass.

(a) Argue that the total drag force is largely unaffected by the non-zero size of the subject body if $r_h \lesssim b_{90}$, where b_{90} is given by equation (3.51).

(b) If $r_h \gtrsim b_{90}$, argue that encounters with impact parameter $\lesssim r_h$ make a negligible contribution to the total drag force. Using the first of equations (L.11), argue that in this case the argument of the Coulomb logarithm is given by $\Lambda \simeq b_{\text{max}}/r_h$.

(c) Combine these conclusions to argue that the correct value of the argument of the Coulomb logarithm for a subject body of half-mass radius r_h is approximately

$$\Lambda = \frac{b_{\text{max}}}{\max(r_h, GM/v_{\text{typ}}^2)}, \quad (8.144)$$

and that the fractional error in $\ln \Lambda$ that arises from using this expression is of order $(\ln \Lambda)^{-1}$.

8.3 [3] In the core of a certain flattened elliptical galaxy, the mean stellar velocity vanishes and the velocity distribution is Gaussian, with dispersion σ_z parallel to the galaxy's symmetry axis $\hat{\mathbf{e}}_z$, and dispersion $\sigma_\perp = \sigma_z/\sqrt{1-e^2} > \sigma_z$ in directions orthogonal to $\hat{\mathbf{e}}_z$. A massive body moves through the core at velocity $\mathbf{v} = v_z \hat{\mathbf{e}}_z + v_\perp \hat{\mathbf{e}}_\perp$, where $\hat{\mathbf{e}}_\perp \cdot \hat{\mathbf{e}}_z = 0$. Show that the frictional drag on the body may be written $\mathbf{F} = -\gamma_z v_z \hat{\mathbf{e}}_z - \gamma_\perp v_\perp \hat{\mathbf{e}}_\perp$, where

$$1 < \frac{\gamma_z}{\gamma_\perp} = \frac{I(1, \frac{3}{2})}{I(2, \frac{1}{2})} ; \quad I(\mu, \nu) \equiv \int_1^\infty \frac{d\lambda}{\lambda^\mu (\lambda - e^2)^\nu} \exp \left[-\frac{1}{2} \sigma_\perp^{-2} \left(\frac{v_\perp^2}{\lambda} + \frac{v_z^2}{\lambda - e^2} \right) \right]. \quad (8.145)$$

Hint: use the analogy between the Rosenbluth potential $h(\mathbf{v})$ and the gravitational potential (see discussion following eq. L.19) and equation (2.125).

8.4 [3] Chandrasekhar's dynamical friction formula can be derived using the linear response theory developed in §5.2.4 (Marochnik 1967; Kalnajs 1972b).

(a) Consider a point mass M traveling on the straight-line trajectory $\mathbf{x}_M(t) = \mathbf{v}_M t$ through a uniform stellar system. Show that the spatial Fourier transform (eq. 5.26) of the density response is given by

$$\bar{\rho}_{s1}(\mathbf{k}, t) = M \int dt' \bar{R}(\mathbf{k}, t - t') e^{-i\mathbf{k} \cdot \mathbf{v}_M t'}, \quad (8.146)$$

where $\bar{R}(\mathbf{k}, \tau)$ is the response function (eq. 5.27).

(b) As we showed in §5.2.4, an infinite homogeneous stellar system is unstable, so to avoid an infinite response we must suppress the self-gravity of the system when evaluating equation (8.146). The justification for this neglect is that the instability arises on scales comparable to the Jeans length, while the dominant contribution to dynamical friction comes from encounters at much smaller distances (page 576), for which the effects of self-gravity are small. To remove self-gravity, we simply replace the response function R in equation (8.146) by the polarization function P , which measures the response to a given total potential rather than a given external potential. With this substitution, use equation (5.55) to show that if the stellar system has a Maxwellian DF, the Fourier transform of the density response is

$$\bar{\rho}_{s1}(\mathbf{k}, t) = 4\pi G M \rho e^{-i\mathbf{k} \cdot \mathbf{v}_M t} \int_0^\infty d\tau \tau e^{i\mathbf{k} \cdot \mathbf{v}_M \tau - (k\sigma\tau)^2/2}, \quad (8.147)$$

where ρ and σ are the density and velocity dispersion of the host system.

(c) Take the inverse Fourier transform of $\bar{\rho}_{s1}$, and evaluate the resulting integrals to show that the density response is

$$\rho_{s1}(\mathbf{x}, t) = \frac{GM\rho}{\sigma^2 r} \exp\left(-\frac{v_M^2 \sin^2 \theta}{2\sigma^2}\right) \left[1 - \operatorname{erf}\left(\frac{v_M \cos \theta}{\sqrt{2}\sigma}\right)\right], \quad (8.148)$$

where erf denotes the error function (Appendix C.3), $\mathbf{r} = \mathbf{x} - \mathbf{x}_M$, $r = |\mathbf{r}|$, and θ is the angle between \mathbf{v}_M and \mathbf{r} . Hint: carry out the integral over \mathbf{k} first, using polar coordinates in \mathbf{k} -space with the polar axis along the vector $\mathbf{r} + \mathbf{v}_M \tau$; then evaluate the integral over τ after transforming to the variable $u = 1/\tau$.

(d) Show that the gravitational force exerted on M by this density distribution is

$$\mathbf{F} = 2\pi \frac{G^2 M^2 \rho}{\sigma^2} \frac{\mathbf{v}_M}{v_M} \int \frac{dr}{r} \int_{-1}^1 d\mu \mu \exp\left[-\frac{v_M^2(1-\mu^2)}{2\sigma^2}\right] \left[1 - \operatorname{erf}\left(\frac{v_M \mu}{\sqrt{2}\sigma}\right)\right], \quad (8.149)$$

where $\mu = \cos \theta$.

(e) The upper limit to the integral over radius should be of order the size R of the host system, while the lower limit should be roughly the 90° deflection radius $b_{90} \approx GM/\sigma^2$ (eq. 3.51), since interior to this radius the perturbations to the orbits of passing stars are so large that linear response theory is invalid. With these limits, show that evaluation of the integrals in equation (8.149) yields the standard dynamical friction formula (8.7), with $\Lambda = R/b_{90}$.

8.5 [1] At some initial time the stellar streaming velocity $\bar{\mathbf{v}}(\mathbf{x})$ within an axisymmetric galaxy of density $\rho(R, z)$ constitutes circular rotation at angular frequency $\omega(R, z)$. The galaxy is then perturbed by the high-speed passage of a massive system. Show that within the impulse approximation the instantaneous change $\Delta \mathbf{v}$ in \mathbf{v} that is produced by the encounter satisfies

$$\int d^3 \mathbf{x} \rho (\bar{\mathbf{v}} \cdot \Delta \mathbf{v}) = 0. \quad (8.150)$$

Hint: write $\bar{\mathbf{v}} = \omega R \hat{\mathbf{e}}_\phi$ and exploit the fact that $\Delta \mathbf{v}$ can be derived from a potential.

8.6 [2] Reproduce Figure 8.4.

8.7 [2] Consider a high-speed head-on encounter at relative velocity V . Assume that the perturber is spherical, with gravitational potential $\Phi(r)$, and let (R, z) be cylindrical coordinates such that the z axis coincides with the perturber's trajectory (i.e., the trajectory is $R = 0, z = Vt$).

(a) Show that the only non-zero component of the impulse to a star at (R, z) is

$$\Delta v_R = -\frac{2R}{V} \int_R^\infty \frac{dr}{\sqrt{r^2 - R^2}} \frac{d\Phi}{dr}. \quad (8.151)$$

(b) If the perturber is a Plummer model, $\Phi = -GM/\sqrt{r^2 + b^2}$, of mass M and scale length b (§2.2c), show that the impulse is

$$\Delta v_R = -\frac{2GMR}{V(R^2 + b^2)}. \quad (8.152)$$

(c) If the perturber and the perturbed system are identical Plummer models, show that the energy per unit mass gained by each system in the encounter is

$$\Delta E = \frac{G^2 M^2}{3V^2 b^2}. \quad (8.153)$$

8.8 [1] Show that the probability $P(V) dV$ that two stars drawn from Maxwellian distributions with one-dimensional dispersions σ_1 and σ_2 have relative speed in the interval $(V, V + dV)$ is

$$P(V) dV = (2\pi\sigma^2)^{-3/2} \exp\left(-\frac{V^2}{2\sigma^2}\right) V^2 dV, \quad \sigma^2 = \sigma_1^2 + \sigma_2^2. \quad (8.154)$$

In words, the relative speed distribution is Maxwellian, with squared dispersion equal to the sum of the squared dispersions of the two populations.

8.9 [1] Is there more angular momentum in the orbit of the Magellanic Clouds around our Galaxy or in the spin of the disk of our Galaxy?

8.10 [3] A frictionless railroad crosses a valley that separates two flat plateaus of equal height. At $t = 0$, two cars, each of mass m , are sent off with the same speed v and separation d from the horizontal stretch of track on one side of the valley. Show that when the cars emerge onto the horizontal stretch of track on the other side of the valley, they have zero relative velocity and their separation is unchanged.

Discuss the relation between this system and disk shocking of globular clusters; in particular, why does passage through the disk heat the cluster but leave the relative velocity of the cars unchanged? Hint: consider adding a spring of rest length d and stiffness ω^2/m between the two cars.

8.11 [2] A satellite system of mass m is in a circular orbit around a point-mass host $M \gg m$. Let (x, y, z) be Cartesian coordinates with $\hat{\mathbf{e}}_x$ pointing along the line joining the two masses and $\hat{\mathbf{e}}_z$ normal to the orbital plane. The distance of the tidal surface from m along the x axis is r_J (eq. 8.91). Show that the distance of this surface from m along the y - and z -axes is $\frac{2}{3}r_J$ and $(3^{2/3} - 3^{1/3})r_J$, respectively. Thus, the Roche surface is not spherical.

8.12 [2] In the distant-tide approximation, the tidal field around a freely falling satellite of a host galaxy can be written in the form $-\nabla\Phi_t = -\sum_{i,j=1}^3 \hat{\mathbf{e}}_i \Phi_{ij} x_j$, where $\{x_j\}$ are non-rotating Cartesian coordinates centered on the body (see §8.2.1). The tidal field is said to be **compressive** along axis i if $\hat{\mathbf{e}}_i \cdot \nabla\Phi_t > 0$, that is, if the tidal force points towards the center of the satellite.

(a) If the host galaxy is spherical with density $\rho_h(R)$ at a distance R from its center, prove that the tidal force is compressive along all three axes if and only if $\rho_h > \frac{2}{3}\bar{\rho}_h$, where $\bar{\rho}_h$ is the mean density of the host interior to R (cf. eq. 8.92).

(b) If the density of the host is $\rho(R) \propto R^{-\gamma}$, prove that the tidal force is compressive in all directions if and only if $\gamma < 1$. In a host with this property tidal disruption cannot occur, no matter how small the mass of the satellite may be. How is this result consistent with the discussion of tidal disruption in §8.3?

8.13 [3] This problem investigates how orbits that lie far beyond the Jacobi radius can remain bound to a satellite. We consider a satellite on a circular orbit, using Hill's approximation (§8.3.2) and restrict our attention to the orbital plane of the satellite, $z = 0$. Since the orbits in question are much larger than the Jacobi radius, the gravitational field of the satellite is weak. Thus we may assume that the orbit is described approximately by the solution (8.99) over timescales of order the epicycle period $T_r = 2\pi/\kappa_0$, with constants of motion x_g, y_{g0}, X, Y , and ψ that change slowly due to perturbations from the satellite.

(a) Show that the guiding-center radius x_g changes at a rate

$$\dot{x}_g = \frac{1}{2B_0} \frac{\partial \Phi_s}{\partial y}, \tag{8.155}$$

where $\Phi_s(\mathbf{x}) = -Gm/(x^2 + y^2)^{1/2}$ is the potential from the satellite. Hint: use equations (8.97), (8.101), and (8.102).

(b) If the perturbations from the satellite are weak, and $|x_g| \ll |y_g|$ (assumptions we will justify below) then the term $\partial \Phi_s / \partial y$ in equation (8.155) can be replaced by its average over an epicycle period at fixed values of the constants of motion; that is

$$\frac{\partial \Phi_s}{\partial y} \Rightarrow \left\langle \frac{\partial \Phi_s}{\partial y} \right\rangle \equiv \frac{Gm}{2\pi} \int_0^{2\pi} d\tau \frac{y_g - Y \sin \tau}{[(x_g + X \cos \tau)^2 + (y_g - Y \sin \tau)^2]^{3/2}}. \tag{8.156}$$

In the limit where $|x_g| \ll |y_g| \ll X$, that is, where the distance of the guiding center from the satellite is much less than the epicycle size, show that

$$\left\langle \frac{\partial \Phi_s}{\partial y} \right\rangle = -\frac{Gm y_g}{X^3} W \left(\frac{\Omega_0}{\kappa_0} \right), \quad \text{where } W(u) \equiv \frac{2}{\pi} \int_0^{\pi/2} d\tau \frac{8u^2 \sin^2 \tau - \cos^2 \tau}{[\cos^2 \tau + 4u^2 \sin^2 \tau]^{5/2}}. \tag{8.157}$$

Hint: expand equation (8.156) in a Taylor series, and use equation (8.100). The function $W(u)$ varies from 0.10032 for $u = 1$ (Keplerian orbits) to 0.22662 for $u = 2^{-1/2}$ (flat circular-speed curve) to 0.5 for $u = \frac{1}{2}$ (harmonic oscillator).

(c) Differentiating the equation for $y_g(t)$ in (8.99) and using the assumption that the time derivatives of the constants of motion are small yields $\dot{y}_g = -2A_0 x_g$. Using this result and equation (8.155) show that the equation of motion for the guiding center is

$$\ddot{y}_g = -\frac{A_0}{B_0} \left\langle \frac{\partial \Phi_s}{\partial y} \right\rangle. \tag{8.158}$$

Interpret this result in terms of the “effective mass” introduced in Box 3.3.

(d) Show that the motion of the guiding center is given by

$$x_g(t) = X_g \cos(\omega t + \alpha) \quad ; \quad y_g(t) = Y_g \sin(\omega t + \alpha), \tag{8.159}$$

where X_g and α are arbitrary and

$$\frac{X_g}{Y_g} = -\frac{\omega}{2A_0} \quad ; \quad \omega^2 = \frac{Gm A_0}{-B_0 X^3} W \left(\frac{\Omega_0}{\kappa_0} \right) = \frac{4\Omega_0 A_0^2}{-B_0} \left(\frac{r_J}{X} \right)^3 W \left(\frac{\Omega_0}{\kappa_0} \right); \tag{8.160}$$

the final form has been derived with the use of equation (8.106). In most galactic potentials $A_0 > 0$ and $B_0 < 0$, so $\omega^2 > 0$ and ω is real. Thus the guiding center oscillates around the satellite; if the orbit lies far outside the tidal radius ($X \gg r_J$) then (i) the oscillation is slow in the sense that $\omega \ll \Omega_0$, and (ii) the excursions in y_g are much larger than the excursions in x_g , consistent with the assumptions we made in deriving this result. The approximations we have used also require that $Y_g \ll Y$, that is, the amplitude of the guiding-center oscillations must be smaller than the epicycle amplitude.

(e) In Problem 5.1, we showed that a solid ring orbiting a planet is unstable. This calculation neglected tidal forces. Would a solid ring that is much larger than the Jacobi radius be stable?

8.14 [1] (a) Derive the energy integral (8.104) for the sheared sheet in two ways, first by multiplying the equations of motion (8.97) by \dot{x} , \dot{y} , \dot{z} respectively, adding, and integrating; second by finding the Lagrangian and Hamiltonian that yield the equations of motion.

(b) Assume that we impose periodic boundary conditions on the sheared sheet, by identifying $y + 2\pi R$ with y . Find the angle-action variables for the case $\Phi_s = 0$, and relate these to the angle-action variables in the epicycle approximation (§3.5.3b).

8.15 [1] A spherical host galaxy contains two small satellites having masses m_1 and m_2 . The satellites travel on nearly circular orbits with nearly the same orbital radius and plane.

(a) Argue that their interactions can be described using Hill's approximation (8.97) in the form

$$\ddot{x}_1 - 2\Omega_0\dot{y}_1 - 4\Omega_0 A_0 x_1 = -\frac{\partial\Phi_{12}}{\partial x_1}; \quad \ddot{y}_1 + 2\Omega_0\dot{x}_1 = -\frac{\partial\Phi_{12}}{\partial y_1}; \quad \ddot{z}_1 + \Omega_0^2 z_1 = -\frac{\partial\Phi_{12}}{\partial z_1}, \quad (8.161)$$

where $\Phi_{12} = -Gm_2/|\mathbf{x}_1 - \mathbf{x}_2|$. Here $\mathbf{x}_i \equiv (x_i, y_i, z_i)$ is the position of satellite i , $i = 1, 2$. The equation of motion for satellite 2 is obtained by interchanging the indices 1 and 2.

(b) In this approximation, what is the trajectory of the center of mass of the two satellites, $\mathbf{x}_{\text{cm}} \equiv (m_1\mathbf{x}_1 + m_2\mathbf{x}_2)/(m_1 + m_2)$?

(c) Show that determining the motion of the two satellites can be reduced to solving the equation of motion for a single particle with position $\mathbf{x} \equiv \mathbf{x}_2 - \mathbf{x}_1$.

8.16 [2] This problem analyzes the sheared sheet (§8.3.2) as a model for the kinematics of the solar neighborhood or other stellar disks. For simplicity, we restrict ourselves to a two-dimensional disk, ignoring motion in the z -coordinate, although the results are easily generalized to three-dimensional disks.

(a) Show that in the absence of local mass concentrations (that is, if the satellite potential $\Phi_s = 0$) the equations of motion (8.97) are invariant under the transformation

$$x \rightarrow x + \Delta x \quad ; \quad y \rightarrow y - 2A_0\Delta x t. \quad (8.162)$$

Describe the physical meaning of this symmetry.

(b) According to the Jeans theorem, the equilibrium DF $f(x, y, \dot{x}, \dot{y})$ can depend only on the integrals of motion E_{\parallel} and L (eq. 8.101). Show that the only combination of these integrals that is invariant under the transformation (8.162) is the epicycle energy defined in equation (8.103). Thus argue that if the disk is smooth on small scales, the DF must have the form $f(E_x)$.

(c) For a DF of this form, show that the surface density is independent of position, the mean radial velocity vanishes, the mean azimuthal or y -velocity is $-2A_0x$, and the ratio of the dispersions in the azimuthal and radial directions is

$$\frac{\sigma_y^2}{\sigma_x^2} = \frac{\int d\dot{x}d\dot{y} f(E_x)(\dot{y} + 2A_0x)^2}{\int d\dot{x}d\dot{y} f(E_x)\dot{x}^2} = \frac{\kappa_0^2}{4\Omega_0^2}, \quad (8.163)$$

a result already derived by different methods in equation (3.100).

(d) Show that if $f(E_x) \propto \exp(-E_x^2/\sigma_0^2)$ then

$$f(x, y, \dot{x}, \dot{y}) = \frac{\Sigma}{\pi\sigma_0^2} \frac{\Omega_0}{\kappa_0} \exp\left[-\frac{\dot{x}^2}{2\sigma_0^2} - \frac{2\Omega_0^2(\dot{y} + 2A_0x)^2}{\kappa_0^2\sigma_0^2}\right]. \quad (8.164)$$

Show that this is the analog of the Schwarzschild DF introduced in §4.4.3.

(e) Does this DF exhibit asymmetric drift (§4.8.2a)?

8.17 [2] Assume that the Sun travels in a circular orbit in the Galactic plane. Let (x, y, z) be rotating Cartesian coordinates centered on the Sun, with $\hat{\mathbf{e}}_x$ pointing away from the Galactic center and $\hat{\mathbf{e}}_z$ pointing to the north Galactic pole.

(a) Show that the zero-velocity surfaces in the combined gravitational field of the Sun and the Galaxy are given by

$$2A(B - A)x^2 + (A^2 - B^2 + 2\pi G\rho_0)z^2 - \frac{GM_\odot}{r} = \text{constant}, \quad (8.165)$$

where ρ_0 is the density in the solar neighborhood, A and B are Oort's constants, and $r^2 = x^2 + y^2 + z^2$. Hint: see Problem 3.18.

(b) Let x_J, y_J, z_J be the intersections of the Sun's Roche surface with the coordinate axes. Evaluate these quantities in parsecs, using the parameters in Tables 1.1 and 1.2.

8.18 [2] Reproduce Figure 8.8.

8.19 [3] The goal of this problem is to determine the epicycle amplitude induced in a star as it passes a molecular cloud, in the shear-dominated regime. We use the equations of motion (8.97) and neglect motion perpendicular to the x - y plane. We assume that the cloud is at the origin and that the star is initially on a circular orbit with impact parameter b , so $\mathbf{x}(t) = (b, -2A_0bt)$. If the cloud potential is $\phi(\mathbf{x}) = -Gm/(x^2 + y^2)^{1/2}$ and its mass is sufficiently small that the right sides of the equations of motion can be evaluated along the unperturbed stellar orbit, show that after the encounter the epicycle amplitude is (Julian & Toomre 1966)

$$X = \frac{Gm\Omega_0}{\kappa_0 A_0^2 b^2} \left[K_0 \left(\frac{\kappa_0}{2A_0} \right) + \frac{\kappa_0}{2\Omega_0} K_1 \left(\frac{\kappa_0}{2A_0} \right) \right], \quad (8.166)$$

where K_ν is a modified Bessel function (Appendix C.7). Thus, derive the correction factor f in equation (8.121b).

8.20 [2] In this problem we estimate the rate of growth of epicycle energy in the dispersion-dominated regime. Consider a star traveling on a nearly circular orbit in the equatorial plane of a razor-thin galaxy. At time zero, the star is instantaneously deflected by the gravitational field from a nearby molecular cloud that is itself on a perfectly circular orbit. The star is traveling at speed v with respect to the cloud, and the encounter deflects it through an angle η onto a new, nearly circular orbit within the galactic plane. Show that the deflection changes the star's epicycle energy by an amount

$$\Delta E_x = E_x(\gamma^2 - 1) \left[\sin^2 \eta (\sin^2 \alpha - \gamma^{-2} \cos^2 \alpha) - \frac{1}{2\gamma} \sin 2\eta \sin 2\alpha \right], \quad (8.167)$$

where $\gamma = 2\Omega/\kappa$ and α is the epicycle phase (see eqs. 3.91 and 3.93, or 8.99). At the radius of the star's orbit, there are n clouds per unit area, each having mass m . The mass distribution in the clouds can be represented by a Plummer model with scale length b , which is much smaller than the star's epicycle radius. Using the impulse approximation, assuming that the relative velocity is dominated by the velocity dispersion of the stars, and assuming that the deflection angle η is small, show that the expectation value of the rate of change of epicycle energy is

$$\dot{E}_x = \frac{\sqrt{2}G^2 m^2 n}{b\sqrt{E_x}} (\gamma^2 - 1) \int_0^{\pi/2} d\alpha \frac{\sin^2 \alpha - \gamma^{-2} \cos^2 \alpha}{(\sin^2 \alpha + \gamma^{-2} \cos^2 \alpha)^{3/2}}. \quad (8.168)$$

Verify that $\dot{E}_x > 0$ for $\gamma > 1$. What happens to stars in a galaxy with $\gamma < 1$?

9

Galaxy Formation

Why is the universe populated by galaxies, rather than a uniform sea of stars? Why are most stars in galaxies with luminosities near $L_\star \simeq 3 \times 10^{10} L_\odot$ (eq. 1.18)? What is the physical origin of the fundamental plane of elliptical galaxies (eq. 1.20) and the Tully–Fisher law for disk galaxies (eq. 1.24)? These are the kinds of questions that a complete theory of galaxy formation should answer. The answers to these questions are still incomplete, and what answers we do have are based on a wide range of physics that extends well beyond stellar dynamics. Nevertheless, stellar dynamics does play a central role in the formation of galaxies and in determining their characteristics, and that role is the subject of this chapter.

In §1.3 (eq. 1.73) we saw that only 15% of the cosmic matter density is contributed by baryonic matter, the remaining $\sim 85\%$ of the density being non-baryonic matter that has no strong or electromagnetic interactions. Thus it is a useful first approximation to neglect the contribution of the baryonic matter to the gravitational forces involved in galaxy formation. In this approximation the dark matter is not affected by the baryons. Without the complications of baryonic physics (radiation, fluid dynamics, star formation, etc.) the dynamics of dark matter becomes a well-defined—though difficult—stellar-dynamical problem that has been extensively studied, and largely solved, by cosmologists over the last several decades. In this chapter we review the results of this study.

The emergence of stars, galaxies, groups and clusters of galaxies, and

even larger inhomogeneities from the nearly homogeneous early universe is referred to as **structure formation**. In §9.1, which is concerned with the earliest phase in the emergence of structure, we shall derive results that apply to baryons and dark matter alike. In §§9.2 and 9.3 we shall restrict ourselves to discussion of the dynamics of dark matter in order to avoid the complexities of nonlinear gas dynamics. This restriction is unfortunate since the sophisticated theoretical framework that emerges describes only the distribution of matter that we cannot see, but the results are essential for understanding the behavior of the baryons that we do see. Finally, in §9.4 we discuss in a speculative vein the complex dynamics of ordinary matter that leads to the formation of visible stars and galaxies near the centers of dark halos.

9.1 Linear structure formation

As we described in §1.3 the universe is homogeneous and isotropic on large scales, and therefore described by the metric (1.44) with the evolution of the scale factor $a(t)$ given by the Friedmann–Robertson–Walker (FRW) equations (1.47)–(1.50). A fundamental assumption of modern cosmology is that early in its history the universe was almost perfectly homogeneous, and that stars, galaxies, clusters, and other large-scale structures developed by the growth of gravitationally unstable fluctuations in the density $\rho(\mathbf{x})$ of the baryonic and non-baryonic matter. Rather than working directly with $\rho(\mathbf{x})$, we use the dimensionless **overdensity**

$$\delta(\mathbf{x}) \equiv \frac{\rho(\mathbf{x})}{\rho_0} - 1, \quad (9.1)$$

where ρ_0 denotes the average matter density over a volume V that is sufficiently large that the universe can be considered homogeneous. We shall usually assume that \mathbf{x} is a comoving coordinate, so the physical distance corresponding to the comoving distance x is $a(t)x$, where $a(t)$ is the cosmic scale factor (§1.3.1). In the linear regime $|\delta| \ll 1$.

In the real universe δ has a well-defined value at each location \mathbf{x} . However, we do not know what this value was early in the life of the universe, and it is natural to imagine it to be the outcome of pure chance. That is, we imagine $\delta(\mathbf{x})$ to be a random variable, and we say that the function δ defines a random field. Thus we replace the concept of a homogeneous universe with that of a statistically homogeneous universe: one in which the statistical properties of $\rho(\mathbf{x})$ in any volume V of given size and shape are independent of the location of its centroid.

If δ is to be a continuous function, the (random) values that it takes at two nearby points \mathbf{x}' and $\mathbf{x}' + \mathbf{x}$ must be correlated: given that the field is

continuous and we know the value of $\delta(\mathbf{x}')$, then the uncertainty in $\delta(\mathbf{x}' + \mathbf{x})$ must become smaller and smaller as $\mathbf{x} \rightarrow 0$. The degree to which $\delta(\mathbf{x}')$ and $\delta(\mathbf{x}' + \mathbf{x})$ are mutually dependent on one another is quantified by the **correlation function**

$$\xi(\mathbf{x}) \equiv \langle \delta(\mathbf{x}')\delta(\mathbf{x}' + \mathbf{x}) \rangle. \quad (9.2)$$

Here the angle brackets imply that the expectation value is to be taken. The assumption of statistical homogeneity allows us to interpret this as the average over all points \mathbf{x}' . If the universe is statistically isotropic, ξ cannot depend on the direction of the displacement \mathbf{x} , but only on its magnitude $x = |\mathbf{x}|$. Thus we can assume that the correlation function has the form

$$\xi(x) \equiv \langle \delta(\mathbf{x}')\delta(\mathbf{x}' + \mathbf{x}) \rangle. \quad (9.3)$$

For simplicity, we restrict our averages to a large but finite volume V , which we take to be a cube of side $V^{1/3}$. We apply periodic boundary conditions on this cube—since we can make V very large, this assumption has no effect on the answer. Then the overdensity is periodic and can be expanded as a Fourier series (Appendix B.4),

$$\delta(\mathbf{x}) = \frac{1}{V} \sum_{\mathbf{k}} \delta_{\mathbf{k}} e^{i\mathbf{k}\cdot\mathbf{x}} \quad \text{where} \quad \delta_{\mathbf{k}} = \int_V d^3\mathbf{x} \delta(\mathbf{x}) e^{-i\mathbf{k}\cdot\mathbf{x}}, \quad (9.4)$$

$\mathbf{k} = 2\pi\mathbf{n}/V^{1/3}$ with $\mathbf{n} = (n_1, n_2, n_3)$ a triple of integers, and the sum is over all \mathbf{n} . Note that $\delta_{\mathbf{0}} = 0$ because $\delta(\mathbf{x})$ has zero mean by definition. Like $\delta(\mathbf{x})$, $\delta_{\mathbf{k}}$ is a random variable with zero mean. The reality of $\delta(\mathbf{x})$ is assured by including both \mathbf{k} and $-\mathbf{k}$ in the sum, with

$$\delta_{-\mathbf{k}} = \delta_{\mathbf{k}}^*. \quad (9.5)$$

When we insert the Fourier expansion (9.4) into the definition (9.3) of the correlation function, we find

$$\xi(x) = \frac{1}{V^2} \sum_{\mathbf{k}'\mathbf{k}} \langle \delta_{\mathbf{k}'}\delta_{\mathbf{k}} \rangle e^{i(\mathbf{k}'+\mathbf{k})\cdot\mathbf{x}'} e^{i\mathbf{k}\cdot\mathbf{x}}. \quad (9.6)$$

Since the right side cannot depend on \mathbf{x}' , we infer that

$$\langle \delta_{\mathbf{k}'}\delta_{\mathbf{k}} \rangle = 0 \quad \text{if} \quad \mathbf{k}' \neq -\mathbf{k}. \quad (9.7)$$

Hence equation (9.6) simplifies to

$$\xi(x) = \frac{1}{V^2} \sum_{\mathbf{k}} \langle \delta_{-\mathbf{k}}\delta_{\mathbf{k}} \rangle e^{i\mathbf{k}\cdot\mathbf{x}} = \frac{1}{V^2} \sum_{\mathbf{k}} \langle |\delta_{\mathbf{k}}|^2 \rangle e^{i\mathbf{k}\cdot\mathbf{x}} \equiv \frac{1}{V} \sum_{\mathbf{k}} P(\mathbf{k}) e^{i\mathbf{k}\cdot\mathbf{x}}, \quad (9.8)$$

where $P(\mathbf{k}) \equiv \langle |\delta_{\mathbf{k}}|^2 \rangle / V$ is the **power spectrum** of the random field. Equation (9.8) states that the power spectrum is the Fourier transform of the correlation function. Since the universe is isotropic, the power spectrum can only depend on $k = |\mathbf{k}|$. The variance of the overdensity is

$$\langle \delta^2(\mathbf{x}) \rangle = \xi(0) = \frac{1}{V} \sum_{\mathbf{k}} P(k). \quad (9.9)$$

9.1.1 Gaussian random fields

We say that $\delta(\mathbf{x})$ is a **Gaussian random field** if, except for the obvious constraint implied by the reality condition (9.5), the $\delta_{\mathbf{k}}$ are independent random variables.¹ Equation (9.4) can be written

$$\delta(\mathbf{x}) = \frac{1}{V} \sum'_{\mathbf{k}} (\delta_{\mathbf{k}} e^{i\mathbf{k}\cdot\mathbf{x}} + \delta_{-\mathbf{k}} e^{-i\mathbf{k}\cdot\mathbf{x}}), \quad (9.10)$$

where \sum' is the sum over only half of \mathbf{k} -space.² In a Gaussian field, each term $\delta_{\mathbf{k}} \exp(i\mathbf{k}\cdot\mathbf{x}) + \delta_{-\mathbf{k}} \exp(-i\mathbf{k}\cdot\mathbf{x})$ is an independent real random variable. According to the central limit theorem (Appendix B.10), a sum of a large number of independent random variables has a Gaussian distribution: that is, the probability that $\delta(\mathbf{x})$ lies in the small interval $(\delta, \delta + d\delta)$ is

$$dp = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{\delta^2}{2\sigma^2}\right) d\delta, \quad (9.11)$$

where $\sigma^2 = \langle \delta^2(\mathbf{x}) \rangle$ is the position-independent variance of the overdensity field. It is from this property that Gaussian random fields derive their name.

Any quantity that is a linear function of the values taken by the overdensity field at several positions will also have a Gaussian distribution. For example, the gradient of the overdensity field is $\nabla\delta = iV^{-1} \sum_{\mathbf{k}} \delta_{\mathbf{k}} \mathbf{k} \exp(i\mathbf{k}\cdot\mathbf{x})$, and the central limit theorem can be applied to this sum to show that its distribution is Gaussian.

The central limit theorem can also be used to establish a much more powerful result on the joint distribution of the overdensities at several locations. Let $\mathbf{u} \equiv (\delta_1, \dots, \delta_D)$ be the vector describing the overdensity $\delta_i \equiv \delta(\mathbf{x}_i)$ at the D positions $\mathbf{x}_1, \dots, \mathbf{x}_D$. Then, by an extension of equation (9.10),

$$\mathbf{u} = \frac{1}{V} \sum'_{\mathbf{k}} [(e^{i\mathbf{k}\cdot\mathbf{x}_1}, \dots, e^{i\mathbf{k}\cdot\mathbf{x}_D}) \delta_{\mathbf{k}} + (e^{-i\mathbf{k}\cdot\mathbf{x}_1}, \dots, e^{-i\mathbf{k}\cdot\mathbf{x}_D}) \delta_{-\mathbf{k}}]. \quad (9.12)$$

Each term in the sum is an independent vector random variable, so by equation (B.99) the joint probability distribution of $\delta_1, \dots, \delta_D$ is given by

$$dp = \frac{d\delta_1 \cdots d\delta_D}{(2\pi)^{D/2} |\mathbf{C}|^{1/2}} \exp\left(-\frac{1}{2} \sum_{i,j=1}^D \delta_i C_{ij}^{-1} \delta_j\right), \quad (9.13)$$

¹ Mathematically we require that if $\delta_{\mathbf{k}} = a + ib$, with a, b real, and $p_{\mathbf{k}}(\delta_{\mathbf{k}}) da db$ is the probability that $\delta_{\mathbf{k}}$ lies in the infinitesimal area of the complex plane $da db$, then the probability that $\delta_{\mathbf{k}}$ lies in this area *and* $\delta_{\mathbf{k}'}$ lies in $da' db'$ is $p_{\mathbf{k}}(\delta_{\mathbf{k}}) p_{\mathbf{k}'}(\delta_{\mathbf{k}'}) da db da' db'$ so long as $\mathbf{k} \neq \pm\mathbf{k}'$.

² For example the volume $k_1 > 0$ plus the half-plane $k_1 = 0, k_2 > 0$, plus the half-line $k_1 = k_2 = 0, k_3 > 0$. Recall that $\delta_0 = 0$.

where \mathbf{C}^{-1} is the inverse of the matrix \mathbf{C} defined by

$$C_{ij} = \langle \delta(\mathbf{x}_i) \delta(\mathbf{x}_j) \rangle = \xi(|\mathbf{x}_j - \mathbf{x}_i|), \quad (9.14)$$

$|\mathbf{C}|$ is the determinant of this matrix, and $\xi(x)$ is defined by equation (9.2). Since the matrix \mathbf{C} is determined by the correlation function $\xi(x)$, which in turn is determined by the power spectrum $P(\mathbf{k})$ through equation (9.8), *all of the statistical properties of a Gaussian random field are determined by its power spectrum* (e.g., Bardeen et al. 1986).

(a) Filtering To understand the physical meaning of the power spectrum, imagine smoothing the density field with a filter that averages out all density fluctuations on scales less than the **smoothing length** L . More precisely, we convolve the overdensity field $\delta(\mathbf{x})$ with a filter or **window function** $W(\mathbf{x})$ to generate a smoothed field $\delta_L(\mathbf{x})$:

$$\delta_L(\mathbf{x}) \equiv \int d^3\mathbf{x}' W(\mathbf{x} - \mathbf{x}') \delta(\mathbf{x}'), \quad (9.15)$$

where $\int d^3\mathbf{x} W(\mathbf{x}) = 1$ and $W(\mathbf{x}) \simeq 0$ for $|\mathbf{x}| \gtrsim L$. In terms of its Fourier components, the smoothed overdensity is $\delta_L(\mathbf{x}) = V^{-1} \sum_{\mathbf{k}} \delta_{L,\mathbf{k}} \exp(i\mathbf{k} \cdot \mathbf{x})$. It is straightforward to show that the relation between the smoothed and unsmoothed Fourier components (eq. 9.4) is

$$\delta_{L,\mathbf{k}} = \widetilde{W}(\mathbf{k}) \delta_{\mathbf{k}} \quad \text{where} \quad \widetilde{W}(\mathbf{k}) \equiv \int d^3\mathbf{x} W(\mathbf{x}) e^{-i\mathbf{k} \cdot \mathbf{x}}. \quad (9.16)$$

Thus smoothing reduces the Fourier amplitudes by a factor $|\widetilde{W}(\mathbf{k})|$, which is unity for $\mathbf{k} = \mathbf{0}$ and tends to zero for $k \gtrsim K \equiv 2\pi/L$. For brevity we shall say that a filter of this kind has scale K^{-1} although the actual smoothing length is closer to $2\pi/K$.

Although several different filters are used in cosmology, we shall work with fields that have been smoothed by simply setting to zero all amplitudes with k greater than K —that is, we shall work with the **sharp k-space filter**. This filter's window function $W_K(\mathbf{x})$ is determined by taking the inverse of the filter's Fourier transform, $\widetilde{W}_K(\mathbf{k}) = H(1 - k/K)$, where H is the step function (Appendix C.1 and eq. B.68). Denoting $x = |\mathbf{x}|$ we find that

$$W_K(\mathbf{x}) = \frac{1}{2\pi^2 x^3} (\sin Kx - Kx \cos Kx). \quad (9.17)$$

Any given smoothing filter defines a characteristic mass that is equal to the mean density times the volume contained within one smoothing length. One natural way to define this mass is

$$M \equiv \rho_0 \int d^3\mathbf{x} \frac{W(\mathbf{x})}{W(0)} = \frac{\rho_0}{W(0)}. \quad (9.18)$$

For the sharp k -space filter $W_K(0) = K^3/(6\pi^2)$, so

$$M_K = \frac{6\pi^2\rho_0}{K^3}. \quad (9.19)$$

By equations (9.9) and (9.16), the variance of an overdensity field that has been smoothed with the filter W_K is

$$\sigma_K^2 = \frac{1}{V} \sum_{\mathbf{k}} \widetilde{W}_K^2(\mathbf{k}) P(k) = \frac{1}{V} \sum_{|\mathbf{k}| < K} P(k). \quad (9.20)$$

If we now proceed to the limit $V \rightarrow \infty$ of a large box, the interval between values of $k_i = 2\pi n_i/V^{1/3}$ over which we are summing becomes infinitesimal, and we have

$$\frac{1}{V} \sum_{\mathbf{k}} \rightarrow \frac{1}{V} \int d^3\mathbf{n} = \frac{1}{(2\pi)^3} \int d^3\mathbf{k}. \quad (9.21)$$

Then equation (9.20) becomes

$$\sigma_K^2 = \frac{1}{2\pi^2} \int_0^K dk k^2 P(k). \quad (9.22)$$

(b) The Harrison–Zeldovich power spectrum The simplest power spectra are power laws,

$$P(k) \propto k^n. \quad (9.23)$$

For a power spectrum of this form,

$$\sigma_K^2 \propto \int_0^K dk k^{2+n} \propto K^{3+n}. \quad (9.24)$$

For example, suppose that the density field consists of points laid down at random, with average density n ; in other words, the density field is a Poisson process (see Appendix B.8). In a box of side K^{-1} , the mean number of points will be $N = nK^{-3}$, the variance in this number will be $\sigma_N^2 = N$ (eq. B.84), the fractional standard deviation will be $\sigma_N/N = N^{-1/2} \propto K^{3/2}$, and this is just σ_K . Comparing this result with equation (9.24) shows that the power spectrum of a Poisson process has $n = 0$. Power spectra with $n > 0$ are smoother than Poisson processes on large scales.

The most important power-law spectrum in cosmology is the **Harrison–Zeldovich** or **scale-invariant spectrum**, which has $n = 1$. Different theories of inflation predict spectra that deviate little from the Harrison–Zeldovich spectrum but Zeldovich’s³ interest in this particular case was motivated by a much simpler argument. In general relativity, the fluctuations

³ Ya. B. Zeldovich (1914–1987) received his Ph.D. for work in chemistry, and then made important contributions to the theory of combustion, detonation, nuclear chain reactions and particle physics. From 1965 to the end of his life he worked on astrophysics, including the CMB and galaxy formation. Prior to 1963 he was a key figure in the Soviet nuclear-weapons program and he was never permitted to travel outside the Soviet bloc.

in the metric tensor are of order Φ/c^2 , where Φ is the gravitational potential arising from the density fluctuations. If the metric-tensor fluctuations are too large on large scales, the universe will not be homogeneous and isotropic; if too large on small scales, small structures will be relativistic. Therefore it is natural to be interested in a spectrum for which the RMS fluctuation in Φ is independent of scale. From equation (9.24), the RMS density fluctuation on a scale K^{-1} is $\rho_{\text{RMS}}(K) \approx \rho_0 \sigma_K \propto K^{(3+n)/2}$. The RMS mass fluctuation is $M_{\text{RMS}}(K) \approx \rho_{\text{RMS}}(K) K^{-3} \propto K^{(n-3)/2}$, and the corresponding fluctuation in the gravitational potential is $\Phi_{\text{RMS}}(K) \approx GM_{\text{RMS}}(K)K \propto K^{(n-1)/2}$. Thus the RMS potential fluctuation is independent of scale if and only if $n = 1$.

9.1.2 Gravitational instability in the expanding universe

At early times $|\delta(\mathbf{x})| \ll 1$ everywhere and we can derive linear equations for the evolution of the overdensity. For the moment we assume that the cosmic material behaves like a fluid, so consider the form taken by the continuity and Euler equations (eqs. F.3 and F.10) in comoving coordinates \mathbf{x} . In principle this can be found by use of the chain rule, but in practice it is simpler to go back to the underlying physics. There are two limiting cases to consider, depending on whether relativistic or non-relativistic fluid dominates the inhomogeneities.

(a) Non-relativistic fluid Let V be a volume that is fixed in comoving coordinates. With ρ the usual fluid density, the mass contained in V is $M = a^3 \int_V d^3\mathbf{x} \rho$. With $\mathbf{v} = \dot{\mathbf{x}}$ the rate of change of the comoving coordinate of the fluid element that is at \mathbf{x} , the physical velocity of the fluid relative to the boundary of V is $a\mathbf{v}$, so the rate at which mass flows out of V is $a^2 \oint_V d^2\mathbf{S} \cdot (a\mathbf{v}) \rho$, where $a^2 d^2\mathbf{S}$ is the vector whose magnitude is the area of an element of the surface of V , and whose direction is the outward normal to the surface. Equating this to minus the derivative of the integral that gives M , and bearing in mind that V is arbitrary, we conclude that

$$\frac{\partial}{\partial t}(a^3 \rho) + a^3 \nabla \cdot (\rho \mathbf{v}) = 0, \quad (9.25)$$

where $\nabla = \partial/\partial\mathbf{x}$. Expanding the derivative on the left and identifying $H(t) = \dot{a}/a$ as the Hubble parameter at time t , the continuity equation in comoving coordinates becomes

$$\frac{\partial \rho}{\partial t} + 3H\rho + \nabla \cdot (\rho \mathbf{v}) = 0. \quad (9.26)$$

Consider next the form of the Euler equation in comoving coordinates. The kinetic energy per unit mass of a particle that is located at \mathbf{x} and has comoving velocity $\mathbf{v} \equiv d\mathbf{x}/dt$ is $\frac{1}{2}(\dot{a}\mathbf{x} + a\mathbf{v})^2$, so the particle's Lagrangian is

$$\mathcal{L} = \frac{1}{2}(\dot{a}\mathbf{x} + a\mathbf{v})^2 - \Phi. \quad (9.27)$$

From equation (D.48) we deduce that its equation of motion is

$$\frac{d}{dt} [(\dot{a}\mathbf{x} + a\mathbf{v})a] - (\dot{a}\mathbf{x} + a\mathbf{v})\dot{a} + \nabla\Phi = 0. \quad (9.28)$$

We can rearrange this to

$$\frac{d\mathbf{v}}{dt} + 2H\mathbf{v} + \frac{\ddot{a}}{a}\mathbf{x} + \frac{1}{a^2}\nabla\Phi = 0. \quad (9.29)$$

If we now consider the particle to be the element of fluid located at \mathbf{x} , and bear in mind that in addition to the gravitational force this element may be subject to a pressure force, we see that in comoving coordinates Euler's equation (F.10) reads

$$\frac{\partial\mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla)\mathbf{v} + 2H\mathbf{v} + \frac{\ddot{a}}{a}\mathbf{x} = -\frac{1}{a^2} \left(\frac{1}{\rho}\nabla p + \nabla\Phi \right). \quad (9.30)$$

We now linearize the continuity equation (9.26) around the undisturbed Hubble flow in which $\mathbf{v} = 0$ and ρ is independent of \mathbf{x} . Using subscripts 0 and 1 to denote the undisturbed and perturbed parts of quantities, we have

$$\frac{\partial\rho_0}{\partial t} + 3H\rho_0 = 0 \quad ; \quad \frac{\partial\rho_1}{\partial t} + 3H\rho_1 + \rho_0\nabla \cdot \mathbf{v} = 0. \quad (9.31)$$

The first of these simply implies that $\rho_0 \propto a^{-3}$ (eq. 1.59 with $w = 0$). Dividing the second equation by $\rho_0(t)$ yields

$$\begin{aligned} 0 &= \frac{\partial\delta}{\partial t} + \frac{\partial \ln \rho_0}{\partial t} \delta + 3H\delta + \nabla \cdot \mathbf{v} \\ &= \frac{\partial\delta}{\partial t} + \nabla \cdot \mathbf{v}, \end{aligned} \quad (9.32)$$

where the second equality uses the first of equations (9.31).

Similarly linearizing equation (9.30) we obtain

$$\ddot{a}\mathbf{x} = -\nabla\Phi_0 \quad ; \quad \frac{\partial\mathbf{v}}{\partial t} + 2H\mathbf{v} = -\frac{1}{a^2} \left(\frac{1}{\rho_0}\nabla p_1 + \nabla\Phi_1 \right). \quad (9.33)$$

In the first equation, $\nabla\Phi_0 = -GM\mathbf{x}/(a^2x^3)$, where M is the gravitational mass within a sphere of comoving radius x , and this equation is equivalent to equation (1.47). Taking the curl of the second equation, we find that the vorticity $\boldsymbol{\omega} \equiv \nabla \times \mathbf{v}$ satisfies $\partial\boldsymbol{\omega}/\partial t = -2H\boldsymbol{\omega}$, so any vorticity initially present decays with time. Hence we may assume that $\boldsymbol{\omega} = \mathbf{0}$. Since any vector field that has vanishing curl can be written as the gradient of a scalar

field, we can write $\mathbf{v} = \nabla\psi$. When we use this expression to eliminate \mathbf{v} from (9.33) we obtain

$$\nabla \left[\frac{\partial\psi}{\partial t} + 2H\psi + \frac{1}{a^2} \left(\frac{p_1}{\rho_0} + \Phi_1 \right) \right] = 0. \quad (9.34)$$

Evidently the expression in square brackets is independent of \mathbf{x} , so it can only depend on t . Nothing of physical significance is changed if we add to ψ any function of t only, and we exploit this freedom to ensure that the square bracket always vanishes, so we have

$$\frac{\partial\psi}{\partial t} + 2H\psi + \frac{1}{a^2} \left(\frac{p_1}{\rho_0} + \Phi_1 \right) = 0. \quad (9.35)$$

Taking the Fourier transform of equation (9.32) yields $\partial\delta_{\mathbf{k}}/\partial t = k^2\psi_{\mathbf{k}}$. Similarly Fourier transforming equation (9.35) and eliminating $\psi_{\mathbf{k}}$ between these equations we obtain

$$\frac{\partial^2\delta_{\mathbf{k}}}{\partial t^2} + 2H\frac{\partial\delta_{\mathbf{k}}}{\partial t} + \frac{k^2}{a^2} \left(\frac{p_{1\mathbf{k}}}{\rho_0} + \Phi_{1\mathbf{k}} \right) = 0. \quad (9.36)$$

To proceed further we require a relation between $p_{1\mathbf{k}}$ and $\delta_{\mathbf{k}}$. Simple inflationary theories predict that the cosmic fluid is everywhere on the same adiabat, so $p_{1\mathbf{k}} = v_s^2\rho_0\delta_{\mathbf{k}}$ (Peacock 1999), where v_s is the sound speed (eq. F.50). Poisson's equation $a^{-2}\nabla^2\Phi_1 = 4\pi G\rho_1$ links Φ_1 to the perturbed density. Taking the Fourier transform of Poisson's equation we eliminate $\Phi_{1\mathbf{k}}$ from equation (9.36) and have finally

$$\frac{\partial^2\delta_{\mathbf{k}}}{\partial t^2} + 2H\frac{\partial\delta_{\mathbf{k}}}{\partial t} + \left(\frac{k^2}{a^2}v_s^2 - 4\pi G\rho_0 \right) \delta_{\mathbf{k}} = 0. \quad (9.37)$$

In this equation ρ_0 is the undisturbed value of the density of the non-relativistic fluid, whose disturbed value is measured by $\delta_{\mathbf{k}}$. Our derivation has excluded relativistic fluids, which can couple to the non-relativistic fluid, either through Φ_1 or through electromagnetic interactions such as Thomson scattering. However, the derivation of equation (9.37) is valid in the presence of vacuum energy, which does not cluster, if ρ_0 is interpreted as the matter density ρ_m rather than the total density $\rho_m + \rho_\Lambda$. In any given cosmological model H , v_s and ρ_0 are known functions of time, so we can solve this equation for the evolution of the overdensity.

We first examine the solution in the important case of large-scale, long-wavelength perturbations, for which k is small. Then the term that is proportional to k^2 can be neglected. In this case the solutions are given in Problem 9.1 for a universe that contains only non-relativistic matter and vacuum energy. In the special case of matter domination, which holds over

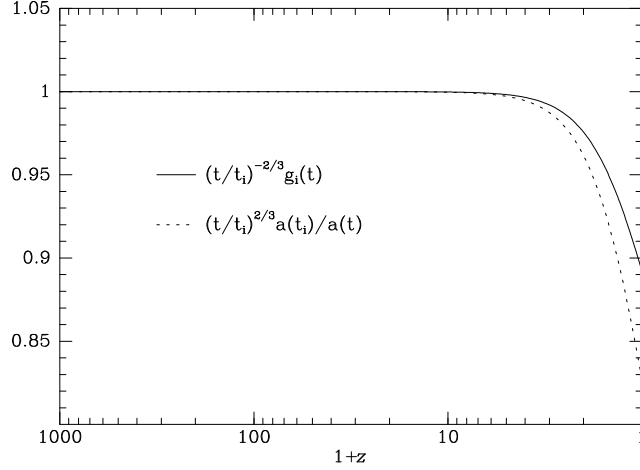


Figure 9.1 Full curve: $(t/t_i)^{-2/3}$ times the growth factor $g_i(t)$ defined by equation (9.40) with t_i corresponding to $1+z = 1000$. The values of g_i were obtained by numerical evaluation of the integral (9.92) for a flat universe with $\Omega_{m0} = 1 - \Omega_{\Lambda 0} = 0.27$ (eq. 1.73) Dashed curve: $(t/t_i)^{2/3} a(t_i)/a(t)$ for the same cosmology.

the interval $3100 \gtrsim z \gtrsim 0.5$ (eqs. 1.70 and 1.74) in which most structure forms, $H = 2/(3t)$ and $4\pi G\rho_0 = 2/(3t^2)$ (eq. 1.63). The solutions of (9.37) are then linear combinations of the growing and decaying power laws,

$$\delta_{\mathbf{k}} \propto t^{2/3} \quad \text{and} \quad \delta_{\mathbf{k}} \propto t^{-1}. \quad (9.38)$$

After some time the growing solution will dominate, so we may neglect the decaying one. In the current and future epoch of domination by vacuum energy, H is a constant and $4\pi G\rho_0 = \frac{3}{2}H^2\Omega_m$ (cf. eq. 1.65) rapidly becomes negligible, so equation (9.37) with $k = 0$ may be approximated by

$$\frac{\partial^2 \delta_{\mathbf{k}}}{\partial t^2} + 2H \frac{\partial \delta_{\mathbf{k}}}{\partial t} = 0. \quad (9.39)$$

The general solution of this equation is $\delta_m(t) = \alpha + \beta \exp(-2Ht)$, where α and β are constants; in words, $\delta_{\mathbf{k}}$ is asymptotically constant. To summarize, density fluctuations grow as $t^{2/3}$ when the universe is matter-dominated, but this growth freezes out once the universe becomes dominated by vacuum energy at $z_{m\Lambda} \simeq 0.5$.

Let t_i be a time early in the era of matter domination, for example the time of decoupling. Then we define the **growth factor** $g_i(t)$ to be the growth in the amplitude of long-wavelength perturbations between t_i and a subsequent time t :

$$g_i \equiv \frac{\delta_{\mathbf{0}}(t)}{\delta_{\mathbf{0}}(t_i)}. \quad (9.40)$$

The full curve in Figure 9.1 is a plot of $(t/t_i)^{-2/3}g_i(t)$ for a flat universe similar to our own (eq. 1.73), while the dashed curve shows $(t/t_i)^{2/3}a(t_i)/a(t)$. Throughout the matter-dominated era both plotted quantities are indistinguishable from unity, attesting to the accuracy of the approximations $g_i(t) \propto t^{2/3}$ and $a(t) \propto t^{2/3}$. The curves turn below unity around $z = 2$, as g_i starts to grow more slowly, and a grows faster, than in a matter-dominated cosmology. However, even at the present epoch g_i is only 10% below the value it would have in the absence of vacuum energy.

(b) Relativistic fluid Equation (9.37) is invalid during the radiation-dominated era, $z \gtrsim z_{\gamma m} \simeq 3100$ (eq. 1.70), because the mass-energy density of radiation does not satisfy the continuity equation (9.26): during expansion the total energy of a body of radiation decreases. However there is a different conservation law we can use. The mean free path of photons is short, so each volume of the photon fluid conserves its entropy. Thus the radiation's entropy density *does* satisfy equation (9.26). The energy and entropy densities of black-body radiation of temperature T are proportional to T^4 and T^3 , respectively (Problem 9.2), so from equation (9.26) we conclude that conservation of entropy implies that

$$\frac{\partial \rho^{3/4}}{\partial t} + 3H\rho^{3/4} + \nabla \cdot (\rho^{3/4}\mathbf{v}) = 0. \quad (9.41)$$

After multiplying through by $\frac{4}{3}\rho^{1/4}$ this can be rewritten

$$\frac{\partial \rho}{\partial t} + 4H\rho + \mathbf{v} \cdot \nabla \rho + \frac{4}{3}\rho \nabla \cdot \mathbf{v} = 0. \quad (9.42)$$

Linearizing, we find that for radiation the analog of equation (9.32) is

$$\frac{\partial \delta}{\partial t} + \frac{4}{3}\nabla \cdot \mathbf{v} = 0. \quad (9.43)$$

Taking the Fourier transform of this equation we find that $\partial \delta_{\mathbf{k}}/\partial t = \frac{4}{3}k^2\psi_{\mathbf{k}}$. Before proceeding, we need to modify the analysis for non-relativistic matter in two ways. First, the linearized Euler equation (9.35) is modified by special relativity to (Problem 9.4)

$$\frac{\partial \psi}{\partial t} + 2H\psi + \frac{1}{a^2} \left(\frac{p_1}{\rho_0 + p_0/c^2} + \Phi_1 \right) = 0. \quad (9.44)$$

For radiation $p = \frac{1}{3}\rho c^2$, so $p_{1\mathbf{k}}/(\rho_0 + p_0/c^2) = \frac{1}{4}c^2\delta_{\mathbf{k}}$. Second, according to general relativity, in Poisson's equation we must replace the density ρ by $\rho + 3p/c^2$ (Problem 9.5). For radiation this implies that $(k/a)^2\Phi_{1\mathbf{k}} = -8\pi G\rho_0\delta_{\mathbf{k}}$. When these relations are used to eliminate $\psi_{\mathbf{k}}$, $p_{1\mathbf{k}}$ and $\Phi_{1\mathbf{k}}$ from the Fourier transform of equation (9.44), we obtain

$$\frac{\partial^2 \delta_{\mathbf{k}}}{\partial t^2} + 2H\frac{\partial \delta_{\mathbf{k}}}{\partial t} + \left(\frac{k^2 c^2}{3a^2} - \frac{32}{3}\pi G\rho_0 \right) \delta_{\mathbf{k}} = 0. \quad (9.45)$$

In the radiation-dominated era we have $2H = 1/t$ and $\frac{32}{3}\pi G\rho_0 = 1/t^2$ (eq. 1.64), so equation (9.45) becomes

$$\frac{\partial^2 \delta_{\mathbf{k}}}{\partial t^2} + \frac{1}{t} \frac{\partial \delta_{\mathbf{k}}}{\partial t} + \left(\frac{k^2 c^2}{3a^2} - \frac{1}{t^2} \right) \delta_{\mathbf{k}} = 0. \quad (9.46)$$

For $k/a \gtrsim 1/ct$ the first term in the bracket dominates the second, and causes $\delta_{\mathbf{k}}$ to oscillate: radiation pressure stabilizes perturbations with wavelengths $2\pi a/k$ smaller than the horizon size $\sim ct$, just as pressure was found to stabilize perturbations with wavelength less than the Jeans length in §5.2: the short-wavelength fluctuations are simply sound waves. On the other hand, when the wavelength is significantly larger than the horizon size, we can neglect this term and the solutions are linear combinations of the power laws

$$\delta_{\mathbf{k}}(t) \propto t \quad \text{and} \quad \delta_{\mathbf{k}}(t) \propto t^{-1}. \quad (9.47)$$

These results tell us how fluctuations in the radiation density evolve in the radiation-dominated era. During this epoch a small fraction of the total density is contributed by non-relativistic matter, in the form of baryons and collisionless dark matter. This density is too small to affect the evolution of perturbations in the relativistic fluid. However, we need to understand how perturbations in these trace constituents evolve in the radiation-dominated era.

The baryon fluid is fully ionized and its free electrons scatter radiation efficiently. Consequently, the baryon fluid is forced to move with the same velocity as the radiation fluid. From equations (9.32) and (9.43), it follows that $\delta(\text{matter}) = \frac{3}{4}\delta(\text{radiation})$, so $\delta(\text{matter}) \propto t$ for wavelengths greater than the horizon size.

Next consider the evolution of fluctuations in the collisionless dark matter. In a collisionless fluid, random particle velocities cause overdensities to disperse, just as density perturbations in a stellar system with wavelength shorter than the Jeans wavelength are damped (§5.2.4). This process erases fluctuations on scales $\lesssim \sigma t$ where σ is the velocity dispersion of the particles at time t . At early times most candidate particles for the dark matter are relativistic,⁴ so $\sigma \sim c$ and the fluctuations are damped on all scales less than the horizon $a(t)x_h(t) \sim ct$ (eq. 1.68). Once the dark-matter particles become non-relativistic, their velocity dispersion declines and the damping scale becomes smaller than the current horizon. Comparison with observations shows no evidence of suppression of the density fluctuations on small scales, which implies that the unknown dark-matter particles must have mass $\gtrsim 1$ keV (Blumenthal et al. 1984). Such particles become non-relativistic at

⁴ It has been conjectured that there is a field whose excitations, **axions**, would be non-relativistic from their instant of formation, when the background temperature was ~ 100 MeV (Bond, Szalay, & Turner 1982).

$z \gtrsim 5 \times 10^6$ and their random velocities at decoupling are $\lesssim 50 \text{ km s}^{-1}$; dark matter satisfying these constraints is known as **cold dark matter (CDM)**.

Precise calculations analogous to the simple ones that we have done can be carried out for a realistic universe that contains baryons, radiation and collisionless dark matter. The results of these calculations are encapsulated in the **transfer function**

$$T^2(k) = \frac{\langle \delta_k^2 \rangle_{z=0}}{\langle \delta_k^2 \rangle_{z \rightarrow \infty}} \bigg/ \frac{\langle \delta_0^2 \rangle_{z=0}}{\langle \delta_0^2 \rangle_{z \rightarrow \infty}}. \quad (9.48)$$

Here δ_0 is the overdensity on large scales, which obeys equations such as (9.37) and (9.46) with $k = 0$. Thus $T(k)$ represents the factor by which the linear fluctuations with wavenumber k are enhanced or suppressed relative to large-scale fluctuations, which are easily calculated for any given cosmological model—see Problem 9.1. By definition, $T(0) = 1$.

The transfer function is plotted in Figure 9.2 for a universe with parameters given by equation (1.73). Its structure is straightforward to understand qualitatively. While the universe is radiation-dominated, fluctuations in both matter and radiation initially grow as $\delta \propto t$ (eq. 9.47), so long as the comoving wavelength of the fluctuation is larger than the growing horizon size ($a/k \gtrsim ct$). Once the fluctuation is contained within the horizon, its growth is stopped or greatly slowed: fluctuations in the radiation density cannot grow because the wavenumber exceeds the Jeans wavenumber as described after (9.46), and fluctuations in the baryons cannot grow because they are strongly coupled to the radiation. The growth of fluctuations in the non-baryonic dark matter also slows, but for a different reason: because the universe is radiation-dominated, the characteristic expansion time $a/\dot{a} \sim (G\rho_\gamma)^{-1/2}$ is much shorter than the characteristic growth time for gravitational instability in the matter, $\sim (G\rho_m)^{-1/2}$, so the instability does not have time to grow.

Eventually, at redshift $z_{\gamma m} \simeq 3100$ (eq. 1.70 and Problem 1.13), the universe becomes matter-dominated, and the dark-matter fluctuations begin to grow as $\delta \propto t^{2/3}$ (eq. 9.38). The baryons remain frozen to the radiation, and thus their fluctuations cannot grow, until decoupling at redshift $z_d \simeq 1100$ (eq. 1.71). During decoupling, the population of free electrons fades and the mean free path of photons increases. Hence there is an interval in which the baryons are still coupled to the photons, but the mean free path of the photons is not small. In this interval the photon-baryon fluid has large coefficients of viscosity and thermal conductivity, so fluctuations with scales smaller than the horizon, which are sound waves, are rapidly damped—this process is called **Silk damping** (Silk 1968). In fact, they are essentially eliminated on all scales smaller than that of rich clusters of galaxies $\approx 10 \text{ Mpc}$. If there were no dark matter, structures with scales smaller than 10 Mpc could not form, so dark matter is *required* for galaxy

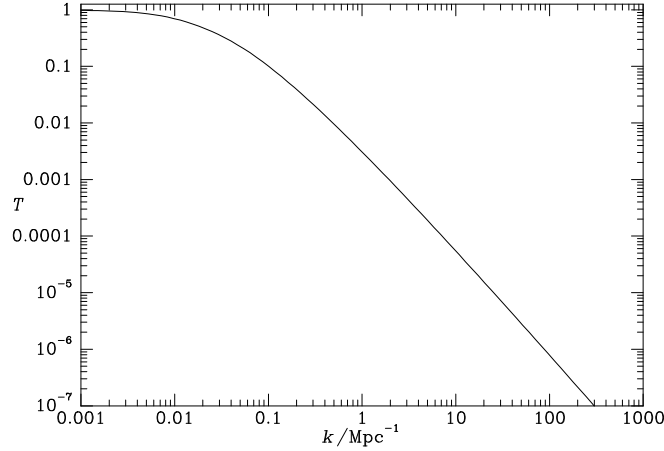


Figure 9.2 The transfer function for a universe with parameters given by equation (1.73) using the approximation (9.50).

formation.⁵ However, small-scale fluctuations in the dark matter persist after decoupling, so the baryons fall into the potential wells that they generate, and from then on the dark-matter and baryon fluctuations evolve together. Thus notwithstanding Silk damping, the transfer function is essentially the same for dark matter and baryons in a given cosmological model.

Let $t_{\gamma m}$ and $x_{\gamma m} = x_h(t_{\gamma m})$ be the time and the comoving horizon size when the universe becomes matter-dominated; for the model of equation (1.73), $t_{\gamma m} \simeq 6 \times 10^4$ yr (eq. 1.70) and $x_{\gamma m} \simeq 125$ Mpc (Problem 1.13). Fluctuations with comoving wavelength $2\pi/k$ larger than $x_{\gamma m}$ grow as t for $t \lesssim t_{\gamma m}$ and as $t^{2/3}$ later, as predicted by equations (9.38) and (9.47); thus the transfer function $T(k)$ is near unity for $kx_{\gamma m} \lesssim 1$. For shorter wavelengths, the growth $\delta \propto t$ is suppressed by a factor $t_h/t_{\gamma m}$, where t_h is the time when this wavelength enters the horizon, defined by $kx_h(t_h) \sim 1$. Since $x_h \propto t^{1/2}$ in the radiation-dominated era (eq. 1.68), $t_h/t_{\gamma m} \sim 1/(kx_{\gamma m})^2$. Thus the transfer function is roughly

$$T(k) \approx \begin{cases} 1 & (kx_{\gamma m} \lesssim 1) \\ (kx_{\gamma m})^{-2} & (kx_{\gamma m} \gtrsim 1). \end{cases} \quad (9.49)$$

The exact transfer function for a given cosmological model can be calculated using publicly available codes such as CMBFAST (Seljak & Zaldarriaga 1996). For the important case of a flat universe containing cold dark matter, in

⁵ If we were to dispense with dark matter, the rotation curves of galaxies would oblige us to subscribe to a theory of gravity that differs materially from that used here, and in such a theory galaxies might be able to form without dark matter.

which the baryon density is negligible ($\Omega_{\text{b}0} \ll \Omega_{\text{m}0}$), the transfer function can be fitted by the formula (Bardeen et al. 1986)

$$T(k) = \frac{\ln(1 + 2.34q)}{2.34q [1 + 3.89q + (16.1q)^2 + (5.46q)^3 + (6.71q)^4]^{1/4}}, \quad (9.50)$$

where $q \equiv \frac{2.04k}{\Omega_{\text{m}0} h_7^2 \text{Mpc}^{-1}}$.

If linear theory were valid to the present, the current matter power spectrum $P(k)$ would be proportional to the product of the **primordial power spectrum** $P_{\text{prim}}(k)$ established by inflation and the square of the transfer function,

$$P(k) \propto T^2(k) P_{\text{prim}}(k). \quad (9.51)$$

The primordial power spectrum encapsulates the physics of inflation, and the transfer function encapsulates the subsequent physics. The solid line in Figure 9.3 shows the theoretical power spectrum for a Harrison–Zeldovich primordial spectrum (eq. 9.23 with $n = 1$) and the transfer function that is plotted in Figure 9.2. Using equation (9.49), we have

$$P(k) \propto \begin{cases} kx_{\gamma\text{m}} & (kx_{\gamma\text{m}} \lesssim 1) \\ 1/(kx_{\gamma\text{m}})^3 & (kx_{\gamma\text{m}} \gtrsim 1). \end{cases} \quad (9.52)$$

The comparison of observations with the linear theory of structure formation is complicated by the fact that on small scales $\lesssim 10$ Mpc the fluctuations are currently nonlinear. Thus a theoretical model of nonlinear structure formation (§9.2) is required to convert measurements of structure at $z = 0$ to the linear fluctuations from which these structures arose. Let z_i be a redshift large enough that fluctuations on the scales of interest are linear, but small enough for the transfer function between that time and the present to be unity. Then if the power spectrum and variance at z_i were $P_i(k)$ and σ_{ik}^2 , and linear theory continued to be valid to the present epoch, the power spectrum and RMS density fluctuation would now be

$$P^{\text{L}}(k) \equiv g_i^2(t_0) P_i(k) \propto T^2(k) P_{\text{prim}}(k) \quad ; \quad \sigma_k^{\text{L}} \equiv g_i(t_0) \sigma_{ik}. \quad (9.53)$$

These are the quantities plotted in Figure 9.3, while Figure 9.4 is a plot of σ_k^{L} against M_k (eq. 9.19) rather than k .

The predicted matter power spectrum can be compared to a wide variety of observations, including measurements of: fluctuations in the temperature of the cosmic microwave background (CMB—see §1.3.5); the clustering of galaxies on large scales; the distribution of clusters of galaxies, which lie at the rare high peaks in the overdensity field; weak gravitational lensing, which measures the distortion of images of distant galaxies by the gravitational field of intervening density fluctuations (Schneider 2006); and the Lyman- α

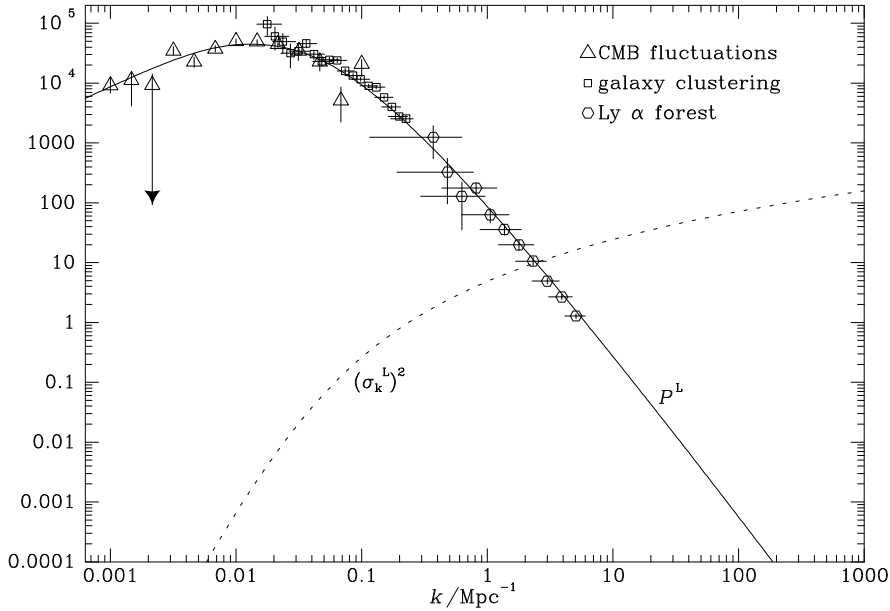


Figure 9.3 The matter power spectrum $P^L(k)$ (eqs. 9.8 and 9.53) at zero redshift, in units of Mpc^3 . The solid line is the theoretical spectrum for a flat FRW model with a Harrison–Zeldovich initial spectrum [$P_i(k) \propto k$] and other parameters from equation (1.73). Triangles show results from the CMB (Figure 11 of Spergel et al. 2007), squares show results from galaxies in the Sloan Digital Sky Survey (Table 3 of Tegmark et al. 2004) scaled to the Hubble constant of equation (1.73), and hexagons (from Tegmark et al. 2004) show results from an analysis of the Lyman- α forest (Gnedin & Hamilton 2002). The dotted line shows the variance $(\sigma_k^L)^2$ of the theoretical overdensity field smoothed on the scale k^{-1} with a sharp k -space filter (eq. 9.22).

forest, the rich set of absorption lines in the spectra of distant quasars. The relation of these measurements to the actual matter power spectrum is not always straightforward. For example, galaxy clustering measures the power spectrum $P_g(k)$ of galaxies, not the overdensity power spectrum $P(k)$. The **bias factor** $b \equiv P_g(k)/P(k)$ is likely to depend on k , galaxy luminosity, Hubble type, etc. The hope is that any bias arises from the complicated processes of galaxy formation, which should have only a limited spatial range (at most the size of the largest clusters, a few Mpc), so any bias on larger scales should be small. The data points in Figure 9.3 show that estimates of the matter power spectrum obtained from the CMB, galaxy clustering, and the Lyman- α forest agree remarkably well with the theoretical curve for a Harrison–Zeldovich primordial spectrum.

The normalization of the power spectrum cannot be predicted theoretically and is fitted to the observational measurements. Historically, the normalization has been parameterized by σ_8 , the RMS overdensity when filtered

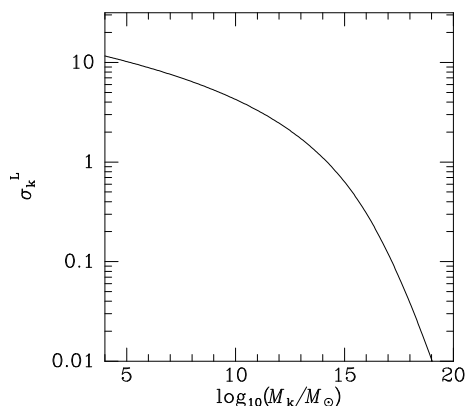


Figure 9.4 The RMS overdensity fluctuation σ_k^L as a function of the mass scale M_k defined by (9.19) for the standard Λ CDM model.

by a particular window function in (9.15), namely the **top-hat function**

$$W_L(\mathbf{x}) = \begin{cases} (\frac{4}{3}\pi L^3)^{-1} & \text{for } |\mathbf{x}| < L \equiv 11.4h_7^{-1} \text{ Mpc} \\ 0 & \text{otherwise.} \end{cases} \quad (9.54)$$

This particular value of L , which is 8 Mpc for $H_0 = 100 \text{ km s}^{-1} \text{ Mpc}^{-1}$, is chosen because then $\sigma_8 \simeq 1$. In fact, current data yield (Spergel et al. 2007)

$$\sigma_8 = 0.76 \pm 0.06. \quad (9.55)$$

This normalization yields the values of the power spectrum and RMS overdensity that are plotted in Figures 9.3 and 9.4.

The success of these observational tests yields the following conclusions: (i) on large scales, the universe is homogeneous and isotropic, as described by the FRW metric (1.44); (ii) the geometry of the universe is flat, as predicted by inflation; (iii) the present densities in ordinary or baryonic matter, dark matter, and vacuum energy are given approximately by equation (1.73); (iv) the dark matter is cold in the sense described on page 728; (v) the initial density fluctuations were small ($|\delta| \ll 1$) and described by a Gaussian random field; (vi) the initial power spectrum of the density fluctuations was approximately the Harrison-Zeldovich spectrum (eq. 9.23 with $n = 1$). The cosmological model with these properties is called the **standard Λ CDM model** after its two principal constituents, vacuum energy and cold dark matter.

9.2 Nonlinear structure formation

The density in a luminous galaxy at a radius of a few kpc is some 10^5 times larger than the critical density ρ_c (eq. 1.55). Galaxy formation therefore involves highly nonlinear density fluctuations, and our investigation of the linear power spectrum must be supplemented by approximate analytic arguments and numerical simulations to follow structure formation into the nonlinear regime. We start by considering a highly idealized model that enables us to introduce the principal concepts.

9.2.1 Spherical collapse

We assume that the background FRW model is matter-dominated and flat; this is a reasonable approximation between $z_{\gamma m} \simeq 3100$ (eq. 1.70) and $z_{m\Lambda} \simeq 0.5$ (eq. 1.74), and it is in this interval that most structures form. Suppose that at some initial time t_i early in this redshift interval, there is a spherically symmetric density fluctuation such that the volume-averaged overdensity within a sphere of radius r_i is $\delta_i \ll 1$. Far outside the sphere the density is that of the FRW model (eq. 1.63),

$$\rho_m(t) = \frac{1}{6\pi G t^2}. \quad (9.56)$$

The total mass in the sphere is

$$M = \frac{4}{3}\pi(1 + \delta_i)\rho_m(t_i)r_i^3. \quad (9.57)$$

We now follow the evolution of the radius $r(t)$ of material at initial radius r_i . The gravitational acceleration of this material is determined only by the total interior mass M , which is constant—the matter is assumed to be cold so there is no flow of material across the shell $r(t)$. Thus Newton's laws imply that

$$\frac{d^2 r(t)}{dt^2} = -\frac{GM}{r^2(t)}. \quad (9.58)$$

This is the equation of motion of a projectile launched vertically from the surface of a spherical body of mass M . Since we are interested in a perturbation that eventually collapses to form a galaxy, we shall assume that the projectile has too little energy to escape. The solution to equation (9.58) can then be written parametrically (eq. 3.28),

$$r = a(1 - \cos \eta) \quad ; \quad t = \sqrt{\frac{a^3}{GM}}(\eta - \sin \eta) + t', \quad (9.59)$$

where $r_{\max} = 2a$ is the **turnaround radius** at which the expansion halts and collapse commences, which occurs at $\eta = \pi$. We may set $t' = 0$ by

arguing that the radius $r(t)$ should be zero at $t = 0$, the time of the Big Bang.⁶ The average density inside the sphere is $\rho_s(t) \equiv M/\frac{4}{3}\pi r^3(t)$, and dividing this by equation (9.56) we obtain the average overdensity inside $r(t)$ as

$$\delta(t) \equiv \frac{\rho_s(t)}{\rho_m(t)} - 1 = \frac{9}{2} \frac{(\eta - \sin \eta)^2}{(1 - \cos \eta)^3} - 1. \quad (9.60)$$

At turnaround the overdensity is⁷

$$\delta_{\max} \equiv \delta(\eta = \pi) \quad \text{or} \quad \delta_{\max} = \frac{9\pi^2}{16} - 1 = 4.55. \quad (9.61)$$

For small density contrast, $\eta \ll 1$, and we can expand (9.60) as a power series to obtain $\delta(t) = \frac{3}{20}\eta^2 + O(\eta^4)$. Thus the initial condition $\delta(t_i) = \delta_i \ll 1$ implies $\eta_i^2 = \frac{20}{3}\delta_i$. From a similar expansion of the first of equations (9.59), we have that $a \simeq 2r_i/\eta_i^2 \simeq \frac{3}{10}r_i/\delta_i$. When we use this relation to eliminate r_i from (9.57), we find that the turnaround radius is

$$r_{\max} = 2a \simeq \left(\frac{243}{250}\right)^{1/3} \frac{(GMt_i^2)^{1/3}}{\delta_i}. \quad (9.62)$$

The turnaround time is

$$t_{\max} = \pi \sqrt{\frac{a^3}{GM}} = \pi \left(\frac{243}{2000}\right)^{1/2} \frac{t_i}{\delta_i^{3/2}} = 1.095 \frac{t_i}{\delta_i^{3/2}}. \quad (9.63)$$

The larger the initial fluctuation δ_i , the sooner the collapse commences.

In this oversimplified model, the collapse reaches singular density at $\eta = 2\pi$ or $t = 2t_{\max}$. In practice, density fluctuations are neither spherically symmetric nor isolated, and the collapsing dark matter will undergo violent relaxation and phase mixing (§§4.10.2 and 9.2.2 below) and settle into an equilibrium configuration, called a **halo**. This collapse, mixing and relaxation process is sometimes called **virialization**, since the halo should satisfy the virial theorem (4.248) after it is complete. A reasonable approximation is that a realistic halo virializes at about the same time that our idealized model collapses to a singular density, that is, at $2t_{\max}$.

Once a halo has settled to approximate equilibrium, its radius can be estimated from the virial theorem. Let us assume that the density fluctuation

⁶Our approximation that the universe is matter-dominated is invalid for $t < t_{\gamma m} \simeq 6 \times 10^4$ yr (eq. 1.70). However, this time is negligibly small.

⁷This analysis neglects the influence of vacuum energy on the dynamics of the collapse. This is justifiable providing the density in the collapsing object is always larger than the vacuum density, that is $\rho_m(1 + \delta) > \rho_\Lambda$. Since ρ_Λ is constant, this condition is always satisfied if it is satisfied at turnaround. Therefore we require $\rho_\Lambda < 5.55\rho_m$ at turnaround. At the present time, $\rho_\Lambda \simeq 2.7\rho_m$ (eq. 1.73), and is smaller at earlier times. Therefore this condition is satisfied for any structure that has already turned around.

inside $r(t)$ is nearly homogeneous at turnaround. Then the potential energy of the material inside the turnaround radius r_{\max} is $W_{\max} = -\frac{3}{5}GM^2/r_{\max}$ (eq. 2.41). Moreover let us assume that all of the material in this region turns around at the same time, so the kinetic energy is zero at t_{\max} and the total energy $E = W_{\max}$. After virialization, which conserves the total energy, the potential energy $W = 2E$ (eq. 4.250), so using $r_h \simeq 0.45GM^2/|W|$ (eq. 4.249b) we find that the half-mass radius of the relaxed system is

$$r_h \simeq 0.375r_{\max}, \quad (9.64)$$

In words, the half-mass radius of the virialized system is about one-third of the turnaround radius. After virialization the mean-square speed of its particles is

$$\langle v^2 \rangle = |W|/M = \frac{6}{5}GM/r_{\max}. \quad (9.65)$$

Since half of the halo mass is inside the half-mass radius, the mean density inside this radius is $\rho_h = \frac{1}{2}M/(\frac{4}{3}\pi r_h^3)$. Assuming once again that the halo becomes virialized at $t = 2t_{\max}$, equation (9.56) implies that the ratio of this density to the unperturbed or background density at this time is

$$\frac{\rho_h}{\rho_m(2t_{\max})} = \frac{9GMt_{\max}^2}{r_h^3} = \frac{9\pi^2}{8} \left(\frac{r_{\max}}{r_h} \right)^3 \simeq 200; \quad (9.66)$$

here we have used equations (9.63) and (9.64) and the relation $r_{\max} = 2a$. This result suggests that regions in which the density exceeds a few hundred times the background density should be part of a halo; although the derivation of equation (9.66) is valid only for a flat, matter-dominated universe, this statement remains valid for most cosmological models of interest. Sometimes the **virial radius** r_{200} of a halo is defined to be the radius at which the density equals 200 times the critical density; inside r_{200} the halo is assumed to be in virial equilibrium, and the mass inside r_{200} is used as a measure of the total mass of the halo.

9.2.2 The cosmic web

We now turn from studying the idealized collapse of a spherical overdensity to what really happens during structure formation. Throughout the period of linear growth, the over-density field $\delta(\mathbf{x})$ retains the same shape; only its amplitude increases as $t^{2/3}$. Thus, if we plot a contour map of $\delta(\mathbf{x})$ at a series of times, we can make the contours at different times coincide by an appropriate rescaling of the plotted contour levels. Moreover, if $\delta(\mathbf{x})$ is a Gaussian field as we have assumed, the properties of the contours for negative values of δ , which enclose underdense regions, are statistically indistinguishable from the contours for $\delta > 0$ that enclose overdense regions.⁸ These two symmetries are broken once the nonlinear regime is entered.

⁸This symmetry follows from the fact that if $\delta(\mathbf{x})$ is a Gaussian field, then so is $-\delta(\mathbf{x})$.

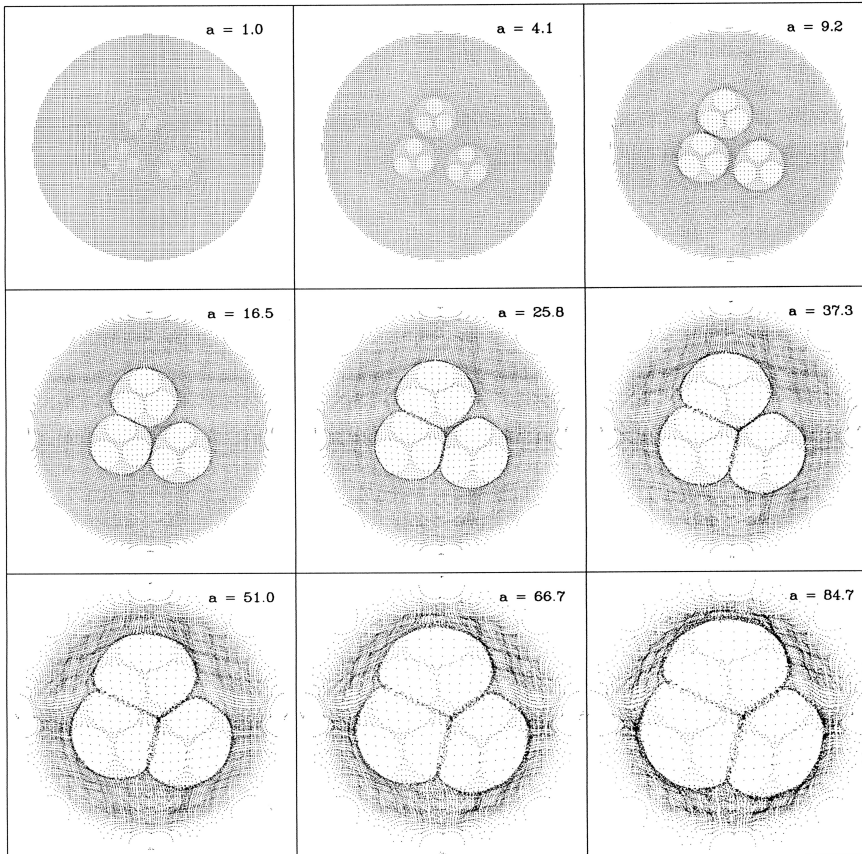


Figure 9.5 The interaction of voids. Top left: near the center of an expanding, self-gravitating sphere are placed three regions of slightly reduced density, and inside each of these are three spheres of even lower density. The low-density spheres expand faster than the region around them, so they run into each other and merge. In the last frame there is a single void within which a web of higher-density regions traces the locations of the original low-density spheres. From Dubinski et al. (1993), by permission of the AAS.

Imagine a roughly spherical region that is underdense. The recession of matter from the center of this region is slowed by gravity, but less so than the general cosmic expansion, so the underdense region expands more rapidly than the universe as a whole. This relative growth of underdense regions, or **voids**, inevitably causes neighboring voids to collide with one another as shown in Figure 9.5. The expanding voids shepherd matter into high-density **sheets** that separate them. When the sheets surrounding three voids meet, a **filament** of high-density material forms. The gravitational attraction of this filament draws in material from the adjoining sheets.

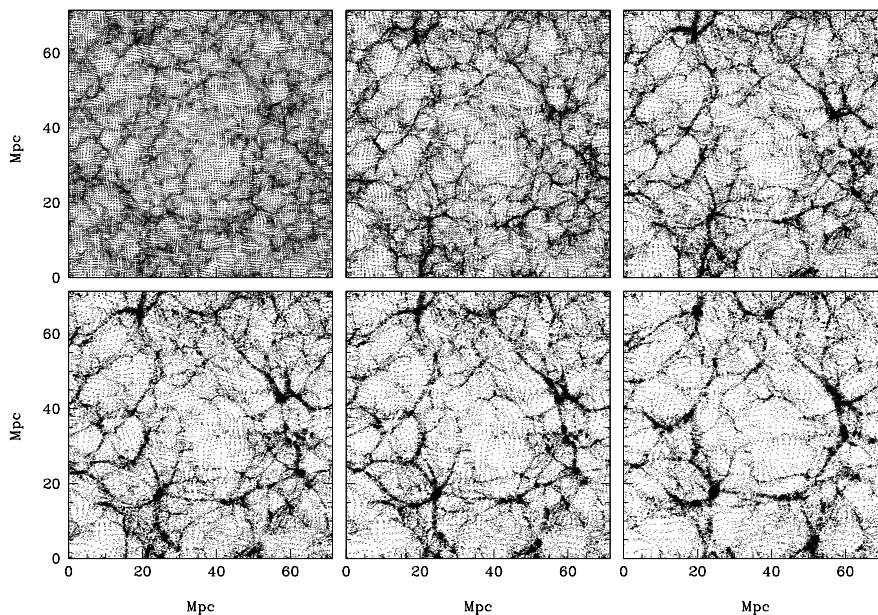


Figure 9.6 Six snapshots of the evolution of a region of a simulated universe: from top left to bottom right $a(t)/a(t_0) = 0.1, 0.2, 0.3, 0.35, 0.4,$ and 0.5 . In the earlier frames many small voids are visible. By the last frame many of these have been squashed into the wall that surrounds the central giant void. From Sheth & van de Weygaert (2004), by permission of the AAS.

Just as the volume of underdense regions grows, the volume occupied by overdense material diminishes because gravity slows the expansion of overdense regions more than the universe as a whole. In fact, matter becomes more and more strongly confined to a **cosmic web** of thin dense sheets that enclose voids (Figure 9.6). The sheets feed matter into a network of filaments at the intersections of the sheets, and these filaments in turn drain into **nodes** at which a number of filaments intersect. Thus the simple picture of spherical collapse is replaced by a model in which material falls first along sheets into filaments and then along filaments into nodes, which develop into virialized halos.

Let us assume that $P(k) \simeq 0$ for $k > K_{\max}$, where K_{\max} is determined by the mass of whatever particle provides the dark matter (see the discussion on page 728). Thus the overdensity field $\delta(\mathbf{x})$ will be smooth on scales $\lambda \lesssim K_{\max}^{-1}$. The nonlinear regime will first be entered where the overdensity is initially largest, and it will be at these locations that the first virialized objects will form through collapse onto a node of the developing cosmic web. The contribution to the variance in the overdensity from an octave in wavenumber is $\sim k^3 P(k)$ (eq. 9.22). Since equation (9.52) shows that $P(k) \propto k^{-3}$ on comoving scales $\lesssim x_{\gamma m} \sim 100$ Mpc, all octaves satisfying

$x_{\gamma\text{m}}^{-1} \lesssim k \lesssim K_{\text{max}}$ contribute equally to the overdensity at a typical point. Thus $\delta(\mathbf{x})$ will be largest at locations that happen to lie near the crests of a number of waves of very different wavelengths. Hence we expect the first objects to form near the crests of both long- and short-wavelength waves, and they will be part of a web that has a characteristic scale of order λ . Gradually the web spreads out from these high-density locations.

Now consider the structure of the overdensity field smoothed on some scale K^{-1} in the range $(\lambda, x_{\gamma\text{m}})$. At a given time t_i this smoothed overdensity field will have smaller variance than the unsmoothed overdensity field. Consequently, the smoothed overdensity field will start attaining the critical density for collapse only after the first objects formed. However, at this time the structure of the smoothed overdensity field will be very similar to the structure that the unsmoothed overdensity field had just before the first objects formed—both have $P(k) \propto k^{-3}$ up to some maximum value of k , beyond which $P(k) \simeq 0$. In particular there will be voids bounded by sheets, like those present when the first structures formed, but with an increase in the linear scale from λ to $\sim K^{-1}$. Thus the cosmic web is constantly regenerated, on ever larger scales; the principal difference between the first appearance of the web and its subsequent reincarnations is that “particles” in the later reincarnations are halos rather than primordial fluid. Since $\sigma_8 \simeq 1$, the typical radius of an over-dense sphere that is currently collapsing is ~ 10 Mpc. Figure 9.7 shows the cosmic web as it appears in the 2dF galaxy-redshift survey, and its appearance is dominated by voids that currently have radii ~ 20 Mpc, but have smaller comoving radii because they have been expanding faster than a region of equal mass in the underlying homogeneous universe.

We argued above that the first objects formed near the crests of long waves, but these are precisely the nodes of the current cosmic web. So we expect to find the oldest objects at the nodes of the current cosmic web, which are the centers of rich clusters of galaxies.

Thus structure formation is **hierarchical** or “bottom-up” in the sense that smaller structures form before larger ones. An ever larger-scale cosmic web is constantly forming from the debris of previous smaller-scale webs. The old concept of galaxies as island universes, which form and evolve without interactions with their neighbors, is turned on its head: galaxies evolve mainly through interactions with their neighbors, by a hierarchical process of halos merging.

At any given redshift halos span a wide range in mass. Some mergers involve halos of comparable mass, and these are called major mergers, while minor mergers involve a small halo falling into a much larger neighbor. An individual halo may contain substructure consisting of smaller halos that it acquired earlier in its life, and that have not yet been disrupted by the processes described in Chapter 8. We define a **primary halo** to be one surrounded by infalling matter, while one that orbits inside a more massive halo is a **subhalo**.

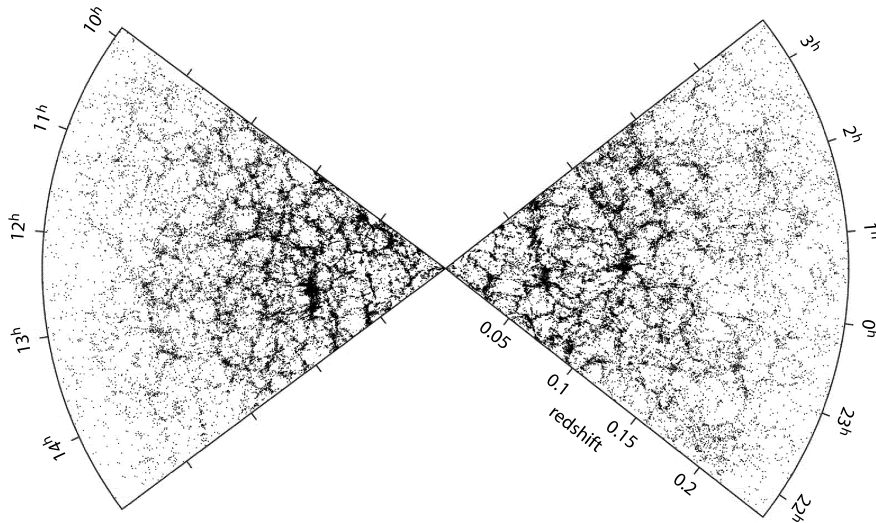


Figure 9.7 The redshift distributions of galaxies along two strips on the sky, each 3° wide. The density of galaxies falls away at large redshifts because at large distances only the most luminous galaxies have been observed. Many voids of diameter $\delta z \simeq 0.01$ ($40h_7^{-1}$ Mpc) can be seen. The data were taken at the Anglo-Australian Observatory as part of the 2dF Galaxy redshift Survey. Credit: Matthew Colless and the 2dF Galaxy Redshift Survey team.

9.2.3 Press–Schechter theory

We now put our treatment of the development of the cosmic web on a quantitative basis by developing **extended Press–Schechter theory**, which has evolved from a seminal paper by Press & Schechter (1974)—more detail can be found in Bond et al. (1991) and Lacey & Cole (1993). We focus on estimating the expected number of halos of a given mass per unit volume (the **mass function**), and the rate at which halos of various masses merge with one another.

Our work on the spherical-collapse model suggests that protohalos can be approximately identified long before they collapse, as regions in which the average overdensity exceeds some **critical overdensity** δ_c . To estimate δ_c at a time t_i much less than the collapse time t , we use equation (9.63); assuming that the halos virialize at $t = 2t_{\text{max}}$,

$$\delta_c(t_i, t) = \left(\frac{2 \times 1.095 t_i}{t} \right)^{2/3} = 1.686 (t_i/t)^{2/3}. \quad (9.67)$$

We have derived the coefficient 1.686 assuming a particular cosmology and, unrealistically, homogeneous collapse. However numerical simulations show

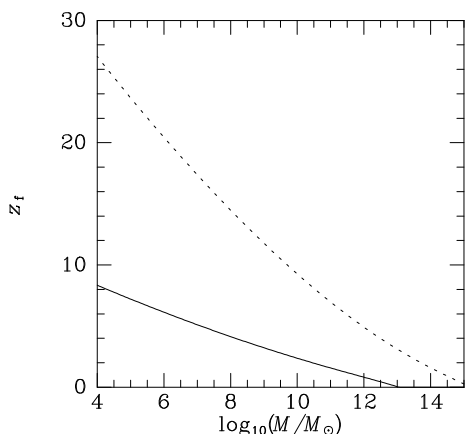


Figure 9.8 The full curve shows the redshift z_f at which a 1σ overdensity of mass M collapsed, while the dashed curve shows the collapse redshifts of 3σ overdensities as a function of mass. These predictions are based on the spherical-collapse model and the fluctuation spectrum plotted in Figure 9.3.

that when this coefficient is used in formulae that we derive below, it yields remarkably good fits to the properties of halos for a wide range of cosmological models (Jenkins et al. 2001).

Figure 9.4 is a plot of the linear growth factor $g_i(t_0)$ times the RMS overdensity fluctuation $\sigma_K(t_i)$ within a spherical region at t_i , as a function of the mean mass M_K contained within such a region (eqs. 9.19 and 9.22). This quantity, $\sigma_K^L = g_i(t_0)\sigma_K(t_i)$ (eq. 9.53), is independent of t_i and would be the RMS overdensity on mass-scale M_K at the present time if linear theory were still valid. According to (9.67), a typically overdense region of mass M_K will have virialized at the time t_f or redshift z_f when $g_i(t_f)\sigma_K(t_i) \simeq 1.686$. The full curve in Figure 9.8 plots z_f as a function of M_K . We see that objects bigger than a large star cluster have typically formed at $z \lesssim 10$.

While a typical overdense region of mass M will not collapse until redshift $z_f(M)$ given by the full curve in Figure 9.8, many halos of this mass will form earlier, from rarer regions of higher overdensity. For example, a $3\sigma_K$ peak in the overdensity field collapsed at $t_f/3^{3/2}$. The dotted line in Figure 9.8 shows the corresponding collapse redshift as a function of M .

Equation (9.65) gives the mean-square velocity of a halo's particles in terms of the radius r_{\max} at which a perturbation turns around. Equation (9.62) connects r_{\max} to $\delta_i/t_i^{2/3}$. If we estimate this ratio as $\sigma_K(t_i)/t_i^{2/3}$, we obtain the estimated RMS velocity of particles in halos as a function of halo mass M_K that is plotted in Figure 9.9. The RMS velocity reaches the values $\sim 100 \text{ km s}^{-1}$ characteristic of galaxies like the Milky Way for halo masses $M \sim 10^{12} \mathcal{M}_\odot$ that are comparable to the masses estimated observationally, for example from the dynamics of the Local Group (Box 3.1), providing a welcome confirmation that we are on the right track.

We focus on a given position, \mathbf{x} , at the initial time t_i and follow the evolution of $\delta(\mathbf{x})$ as we add in the waves that make up the Fourier sum (9.4). We start with the longest wavelengths (smallest values of $|\mathbf{k}|$). At each stage in the process we add to $\delta(\mathbf{x})$ a quantity Δ_K that is the overall perturbation

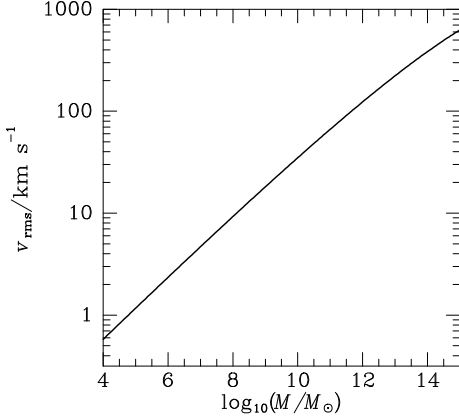


Figure 9.9 The typical RMS internal velocity of a halo. This prediction is based on the spherical-collapse model and the fluctuation spectrum plotted in Figure 9.3.

contributed by all waves that have wavevectors in a spherical shell in \mathbf{k} -space, that is

$$\Delta_K = \frac{1}{V} \sum_{K \leq |\mathbf{k}| < K+dK} \delta_{\mathbf{k}} e^{i\mathbf{k} \cdot \mathbf{x}}. \tag{9.68}$$

Since Δ_K is the sum of a large number of statistically independent random variables $\delta_{\mathbf{k}}$, by the central limit theorem its probability distribution is Gaussian.

Let σ_K^2 be the variance of $\delta(\mathbf{x})$ when we have included all waves with $|\mathbf{k}| \leq K$ (eq. 9.22). Evidently, σ_K^2 will grow monotonically from zero as we add more shells. In fact, since Δ_K is uncorrelated with the value reached by $\delta(\mathbf{x})$ before Δ_K is added, σ_K^2 will increase by $\langle \Delta_K^2 \rangle$ at each step. By contrast, $|\delta(\mathbf{x})|$ does not necessarily increase at a particular step because Δ_K may have the opposite sign to the current value of $\delta(\mathbf{x})$. In fact, if we plot the values of $\delta(\mathbf{x})$ at each stage against the corresponding values of σ_K^2 , $\delta(\mathbf{x})$ will undergo a random walk (Figure 9.10), starting from the origin, in which σ_K^2 plays the role of step number. The RMS distance that a particle executing a random walk moves in n steps is proportional to \sqrt{n} . In our context, the “distance” walked is just $|\delta(\mathbf{x})|$ and the constant of proportionality between the mean-square of the distance traveled and the step number— σ_K^2 —is unity because $\langle |\delta^2| \rangle = \sigma_K^2$ by definition. This result is manifestly independent of the power spectrum, which is why it is advantageous to use σ_K^2 as the independent variable rather than K .

In the spherical-collapse model, a region will collapse and virialize by time t if its overdensity at the initial time $t_i < t$ exceeds the critical overdensity $\delta_c(t_i, t)$ given by (9.67). The dashed line in Figure 9.10 shows this critical overdensity for $t/t_i = 196$. If the random walk first carries $\delta(\mathbf{x})$ over this line at ordinate σ_K^2 , we deduce that \mathbf{x} belongs to a region of scale K^{-1} that will collapse at time t . As t increases, the dashed line in the figure sinks, and the point at which the full curve first crosses the dashed line moves to the left. This leftward movement is jerky because the full curve is jagged,

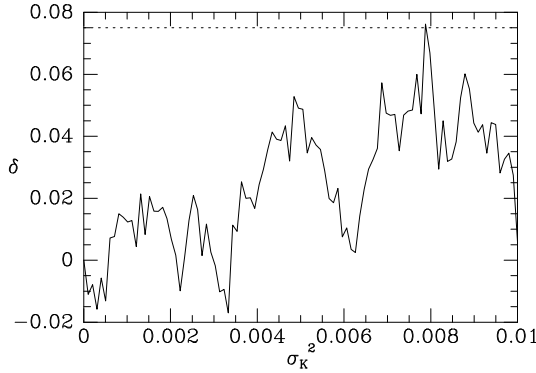


Figure 9.10 The overdensity $\delta(\mathbf{x})$ at some location smoothed by including only waves with wavenumbers $< K$, versus the corresponding variance $\sigma_K^2 = \langle \delta^2(\mathbf{x}) \rangle$ (solid line). The vertical coordinate executes a random walk, with σ_K^2 playing the role of the number of steps. The dashed line shows the critical overdensity δ_c (eq. 9.67) for fluctuations that will collapse by $t_f = 196t_i$.

and each jerk to the left is associated with a discrete increase in the length scale K^{-1} and mass M_K (eq. 9.19) of the largest halo that contains our mass element. Most such leftward movements are small and correspond to minor accretion events which together make up a steady rain or **infall** of material onto the halo. At certain times the dashed line just touches the top of an isolated peak in the full curve, and the point of first intersection leaps to the left; this situation corresponds to major mergers in which the halo that incorporates our mass element merges with a similar or even larger halo, so K^{-1} and M_K increase by large amounts.

For $\delta > 0$ the equation $\delta_c(t_i, t) = \delta$ gives the collapse time t of a region that has overdensity δ at time t_i . Thus for $\delta > 0$ the vertical axis in Figure 9.10 can be regarded as a time axis as shown in Figure 9.11—smaller values of δ correspond to larger times. The point marked A in Figure 9.11 corresponds to $t/t_i = t_A/t_i = 104$ and $\sigma_K^2 = \sigma_{K_A}^2 = 0.0079$. Similarly, the point B corresponds to time $t_B > t_A$ and scale $K_B^{-1} > K_A^{-1}$. At t_A the mass element at \mathbf{x} becomes part of a halo of scale K_A^{-1} , and at t_B it becomes part of a more massive halo of scale K_B^{-1} . At other times the scale of the halo in which our mass element is embedded can be read off from the horizontal position of the heavy line in Figure 9.11. The horizontal sections of that line correspond to mergers, in which the scale and mass of the halo increase discontinuously. The mass of the halo is related to its scale by equation (9.19).

We have described the evolution of the overdensity at \mathbf{x} in one particular realization of the early universe. Actually our interest is in the statistics of the early universe, so we want to know the probability that at a given stage σ_K^2 in the addition process the overdensity at a randomly chosen point \mathbf{x} lies in the interval $(\delta, \delta + d\delta)$. We denote this probability by $p_K(\delta) d\delta$ and note that through its subscript it is a function of σ_K^2 . If we consider a large number N of locations \mathbf{x} , the number of locations at which the overdensity takes a value in the range $(\delta, \delta + d\delta)$ will be $Np_K(\delta) d\delta$. As we add more waves, and σ_K^2 increases, the overdensity at some locations will decrease, while that at other locations will increase, with the consequence that the

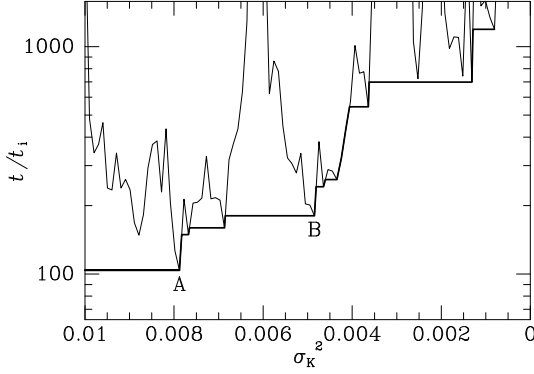


Figure 9.11 The light curve is as in Figure 9.10 but with δ replaced by $t/t_1 = (1.686/\delta)^{3/2}$ and with the ordering of the horizontal axis reversed so M_K rather than σ_K^2 increases to the right. The heavy curve shows how M_K increases discontinuously with time.

cloud of points will diffuse along the line of δ -values. Thus the density of the cloud, $Np_K(\delta)$, must satisfy the diffusion equation (cf. eq. 7.68) with σ_K^2 playing the role of time:

$$\frac{\partial p_K}{\partial \sigma_K^2} = C \frac{\partial^2 p_K}{\partial \delta^2}. \tag{9.69}$$

Here C is a diffusion coefficient, which we determine by noting that $\sigma_K^2 \equiv \langle \delta^2 \rangle = \int d\delta \delta^2 p_K(\delta)$. Differentiating both sides of this equation with respect to σ_K^2 , carrying the derivative inside the integral on the right, and then using equation (9.69), we have

$$1 = C \int d\delta \delta^2 \frac{\partial^2 p_K}{\partial \delta^2}. \tag{9.70}$$

Two integrations by parts and the normalization condition $\int d\delta p_K = 1$ enable us to evaluate the integral and infer that $C = \frac{1}{2}$, so

$$\frac{\partial p_K}{\partial \sigma_K^2} = \frac{1}{2} \frac{\partial^2 p_K}{\partial \delta^2}. \tag{9.71}$$

We seek a solution of equation (9.71) that satisfies the normalization condition $\int d\delta p_K(\delta) = 1$ and the boundary value that at $\sigma_K^2 = 0$, $p_K(\delta) = 0$ for $\delta \neq 0$. The required solution is (Problem 9.9)

$$p_K(\delta) = \frac{1}{(2\pi\sigma_K^2)^{1/2}} \exp\left(-\frac{\delta^2}{2\sigma_K^2}\right). \tag{9.72}$$

A related solution of equation (9.71) is

$$p_K(\delta) = \frac{1}{(2\pi\sigma_K^2)^{1/2}} \left[\exp\left(-\frac{\delta^2}{2\sigma_K^2}\right) - \exp\left(-\frac{(\delta - 2\delta_c)^2}{2\sigma_K^2}\right) \right], \tag{9.73}$$

where δ_c is a constant. This is the difference between the probability density of particles that are initially at the origin and ones that are initially at

$\delta = 2\delta_c$. Since (9.71) is linear, this difference of solutions is itself a solution. By symmetry it vanishes at $\delta = \delta_c$ for all σ_K^2 . Let us now restrict our attention to the region $\delta < \delta_c$. In this region (9.73) satisfies the boundary condition $p_K(\delta_c) = 0$, and as $\sigma_K^2 \rightarrow 0$ it vanishes for all $\delta \neq 0$, and (ii) it satisfies the normalization condition $\int_{-\infty}^{\delta_c} d\delta p_K(\delta) = 1$. That is, it gives us the probability density of particles that are released from the origin and diffuse towards an absorbing barrier at $\delta = \delta_c$.

(a) The mass function It will prove useful to calculate the probability $p_1(K, t)d\sigma_K^2$ that at time t the primary (i.e., largest) halo that contains some given mass element has scale K^{-1} , where K^{-1} is the radius of the sphere associated with σ_K^2 . This is the probability that the random walk of δ first moves above $\delta_c(t_i, t)$ at the step $\sigma_K^2 \rightarrow \sigma_K^2 + d\sigma_K^2$. Our solution (9.73) to the absorbing-barrier problem enables us to calculate $p_1(K, t)$ as follows. The probability that $\delta < \delta_c$ is $\int_{-\infty}^{\delta_c} d\delta p_K$. If we evaluate this integral with p_K given by (9.73), we obtain the probability that $\delta(\mathbf{x})$ has *not* reached δ_c . Minus the rate of change of this probability is the probability density that we seek. Thus

$$p_1(K, t) = -\frac{\partial}{\partial \sigma_K^2} \int_{-\infty}^{\delta_c} d\delta p_K(\delta) = -\frac{1}{2} \left. \frac{\partial p_K}{\partial \delta} \right|_{-\infty}^{\delta_c}, \quad (9.74)$$

where the last equality follows from (9.71). With equation (9.73) we have

$$p_1(K, t) = \frac{\delta_c(t_i, t)}{(2\pi)^{1/2} \sigma_K^3} \exp\left(-\frac{\delta_c^2}{2\sigma_K^2}\right). \quad (9.75)$$

The quantity $p_1(K, t) d\sigma_K^2$ is the probability that at time t a randomly chosen mass element is part of a halo of scale K^{-1} that is not part of a larger halo. Let $f(M_K, t) dM_K$ be the probability that the largest halo that contains a given mass element has a mass in the interval $(M_K, M_K + dM_K)$. Then $f(M_K, t) dM_K = p_1(K, t) d\sigma_K^2$ so

$$f(M_K, t) = \frac{\delta_c/\sigma_K}{(2\pi)^{1/2} M_K} \exp\left(-\frac{\delta_c^2}{2\sigma_K^2}\right) \left. \frac{d \ln \sigma_K^2}{d \ln M_K} \right|, \quad (9.76)$$

where the right side is a function of M_K and t because δ_c is a function of t and K is a function of M_K (eqs. 9.67 and 9.19). At time t let there be dn halos per unit volume with masses in the interval $(M, M + dM)$. Then the total mass per unit volume in these halos is $M dn$ and this must equal $\rho_0 f(M, t) dM$, so we have

$$\left. \frac{dn}{d \ln M} \right|_t = \rho_0 f(M, t). \quad (9.77)$$

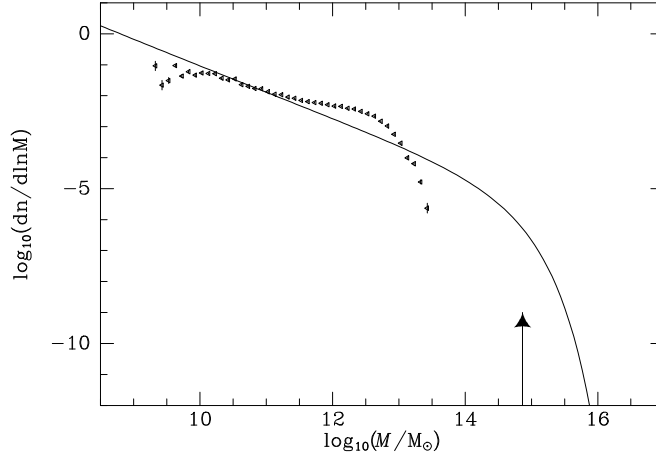


Figure 9.12 The mass function of halos in the standard standard Λ CDM model from equation (9.77) (curve) and the galaxy luminosity function (symbols) scaled to mass-to-light ratio $\Upsilon_R = 220 M_\odot / L_\odot$ (eq. 1.76). The luminosity function is taken from Blanton et al. (2005). The error bars on most points are too small to be seen. The vertical arrow marks the value of M_c .

The full curve in Figure 9.12 is a plot of $dn/d\ln M$ at the present epoch from (9.77) for the standard Λ CDM model described at the end of §9.1. The exponential in (9.76) rapidly diminishes $dn/d\ln M$ for values of M large enough that $\sigma_K \ll \delta_c$. The **characteristic halo mass** $M_c(t)$ is the mass at which $\sigma_K^2(t_i) = \frac{1}{6}\delta_c^2(t_i, t)$ above which $dn/d\ln M$ declines rapidly— $M_c(t_0) = 7.4 \times 10^{14} M_\odot$. The behavior of $dn/d\ln M$ for halos with mass much smaller than M_c depends on the shape of the power spectrum. If in some range of wavenumbers we approximate the power spectrum by the power law $P(k) \propto k^n$, the variance on length scale K^{-1} is given by equation (9.24), so the variance on mass scale $M \propto K^{-3}$ is $\sigma_K^2 \propto M^{-1-n/3}$. In this case equation (9.77) predicts $dn/d\ln M \propto M^{(n-3)/6}$ for $M \ll M_c$. According to equation (9.52) and Figure 9.3, $n \simeq -3$ for all galaxy-sized masses. Thus equation (9.77) predicts $dn/d\ln M \propto M^{-1}$ for $M \ll M_c$, in agreement with Figure 9.12. The symbols in Figure 9.12 show the galaxy luminosity function for one particular scaling $L = M/\Upsilon$ between mass and luminosity. Clearly no scaling could make the luminosity function coincide with the full curve: the luminosity of a halo cannot be proportional to its mass. We shall return to this point at the end of §9.4.

If we return to (9.75), we can derive an equation that provides a valuable way to test extended Press–Schechter theory with N-body experiments (§9.3.1). In a given simulation at time t let dn be the number density of halos that have scales K^{-1} in the interval in which the variance σ_K^2 of the initial density field changes by $d\sigma_K^2$. Then the contribution $M_K dn$ of these

halos to the total mass density ρ_0 must be $\rho_0 p_1(K, t) d\sigma_K^2$, so

$$F \equiv \frac{M_K}{\rho_0} \left| \frac{dn}{d \ln \sigma_K^2} \right| = \sigma_K^2 p_1(K, t) = \frac{\delta_c / \sigma_K}{(2\pi)^{1/2}} \exp\left(-\frac{\delta_c^2}{2\sigma_K^2}\right). \quad (9.78)$$

Since the right side of this equation depends only on δ_c / σ_K , extended Press–Schechter theory predicts that the left side, which can be determined by counting halos, should have this property also. From equation (9.67) we have that

$$\frac{\delta_c}{\sigma_K} = 1.686 \frac{(t_i/t)^{2/3}}{\sigma_K} = \frac{1.686}{\tilde{\sigma}_K} \quad \text{where} \quad \tilde{\sigma}_K \equiv (t/t_i)^{2/3} \sigma_K. \quad (9.79)$$

Thus according to extended Press–Schechter theory, the quantity F depends on t and σ_K only in the combination $\tilde{\sigma}_K$. This is a strong prediction that we shall test in §9.3.1. The quantity $\tilde{\sigma}_K$ has a simple physical interpretation: it is the RMS fluctuation in M_K at time t that linear theory predicts for a strictly matter-dominated universe.

(b) The merger rate The rate at which halos merge can be deduced from Press–Schechter theory in three steps.

(i) Consider a mass element that is within a halo of mass M_2 at time t_2 . What is the probability $f_1(M_1, t_1 | M_2, t_2) dM_1$ that this element belonged to a primary halo with mass in the range $(M_1, M_1 + dM_1)$ at time $t_1 < t_2$? As above, each mass M_j is associated with a wavenumber K_j through (9.19), and each time t_j is associated with a value of the critical overdensity $\delta_j = \delta_c(t_i, t_j)$. Then $f_1 dM_1$ is simply the probability that the random walk, after reaching δ_2 at “time” $\sigma_{K_2}^2$, will first reach $\delta_1 > \delta_2$ in the “time” interval $(\sigma_{K_1}^2, \sigma_{K_1}^2 + d\sigma_{K_1}^2)$. This is the problem addressed in equation (9.75), except that the starting point is $(\sigma_{K_2}^2, \delta_2)$ instead of $(0, 0)$, so the answer can be obtained by replacing σ_K^2 by $\sigma_{K_1}^2 - \sigma_{K_2}^2$ and δ_c by $\delta_1 - \delta_2$:

$$\begin{aligned} f_1(M_1, t_1 | M_2, t_2) dM_1 \\ = \frac{\delta_1 - \delta_2}{(2\pi)^{1/2} (\sigma_{K_1}^2 - \sigma_{K_2}^2)^{3/2}} \exp\left(\frac{-(\delta_1 - \delta_2)^2}{2(\sigma_{K_1}^2 - \sigma_{K_2}^2)}\right) d\sigma_{K_1}^2 \begin{pmatrix} \sigma_{K_1}^2 > \sigma_{K_2}^2 \\ \delta_1 > \delta_2 \end{pmatrix}. \end{aligned} \quad (9.80)$$

(ii) Consider now the probability $f_2(M_2, t_2 | M_1, t_1) dM_2$ that a halo of mass M_1 at time t_1 has been subsumed in a halo that at time $t_2 > t_1$ has mass in the interval $(M_2, M_2 + dM_2)$. Again associating K_j with M_j and δ_j with t_j , we apply Bayes’s theorem (eq. B.87) to obtain

$$\begin{aligned} f_2(M_2, t_2 | M_1, t_1) &= \frac{f_1(M_1, t_1 | M_2, t_2) f(M_2, t_2)}{f(M_1, t_1)} \\ &= \frac{(\delta_1 - \delta_2)(\delta_2/\delta_1)(\sigma_{K_1}/\sigma_{K_2})^3}{(2\pi)^{1/2} (\sigma_{K_1}^2 - \sigma_{K_2}^2)^{3/2}} \left| \frac{d\sigma_{K_2}^2}{dM_2} \right| \\ &\quad \times \exp\left(-\frac{(\delta_1 - \delta_2)^2}{2(\sigma_{K_1}^2 - \sigma_{K_2}^2)} - \frac{\delta_2^2}{2\sigma_{K_2}^2} + \frac{\delta_1^2}{2\sigma_{K_1}^2}\right), \end{aligned} \quad (9.81)$$

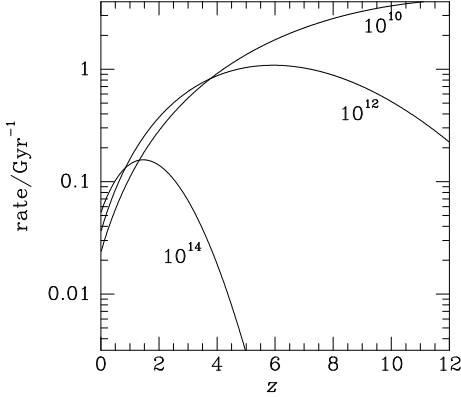


Figure 9.13 Merger rate as given by equation (9.82) for equal-mass mergers. The curves are labeled by the pre-merger mass M_1 in solar masses. The figure is for the standard Λ CDM model.

where $\sigma_{K_1}^2 > \sigma_{K_2}^2$, $\delta_1 > \delta_2$, and $f(M, t)$, the probability that a mass element belongs to a halo of mass M , is given by equation (9.76).

(iii) Now let t_2 approach t_1 , that is write $t_2 = t_1 - dt$ and $\delta_2 = \delta_1 - d\delta_c$. In this limit the probability of multiple mergers is negligible, so $f_2(M_2, t_1 - dt|M_1, t_1)dM_2$ is the probability that a halo with mass M_1 experiences a single merger that boosts its mass to the interval $(M_2, M_2 + dM_2)$. Thus the probability per unit time that a halo of mass M_1 will experience a merger that increases its mass to M_2 is given by

$$\frac{d^2p}{d \ln M_2 dt} = \frac{M_2}{(2\pi)^{1/2}} \left| \frac{d\delta_c}{dt} \frac{d\sigma_{K_2}^2}{dM_2} \right| \left[\frac{\sigma_{K_1}^2/\sigma_{K_2}^2}{\sigma_{K_1}^2 - \sigma_{K_2}^2} \right]^{3/2} \exp \left[-\frac{\delta_c^2}{2} \left(\frac{1}{\sigma_{K_2}^2} - \frac{1}{\sigma_{K_1}^2} \right) \right]. \tag{9.82}$$

The right side depends on time only through the critical overdensity $\delta_c(t_1, t)$. Figure 9.13 shows the merger rate $d^2p/d \ln M_2 dt$ as a function of redshift for $M_1 = 10^{10}, 10^{12}$ and $10^{14} \mathcal{M}_\odot$ and $M_2 = 2M_1$ (equal-mass mergers). At a given mass the rate rises at high redshift to a peak around the redshift at which 3σ peaks are virializing (Figure 9.8) and then declines as more and more halos of the given mass are subsumed in larger halos; for all three masses plotted, the merger rate is predicted to be steeply declining at the current epoch.

In §8.5.6 a mixture of observational data and dynamical arguments led us to estimate that the merger rate per halo of L_* galaxies is $\simeq 0.008 \text{ Gyr}^{-1}$. The halos of these galaxies have masses $\approx 10^{12} \mathcal{M}_\odot$, and for this mass Figure 9.13 predicts a merger rate per halo that is four times higher. Is this a worrying discrepancy? In fact the merger rate we have calculated here is for mergers of *primary* halos, whereas our earlier rate includes mergers of subhalos. A large fraction of halos of mass $10^{12} \mathcal{M}_\odot$ are now subhalos—recall from Figure 9.12 that the characteristic halo mass for clustering is currently $M_c = 7.4 \times 10^{14} \mathcal{M}_\odot$. Once a halo has fallen in to a significantly more massive halo, it is unlikely to merge with one of its peers. Hence when, as in §8.5.6, the merger rate per halo is averaged over galaxies that are hosted by

both primary halos and subhalos, we expect to obtain a value that is smaller than the merger rate per halo when only primary halos are considered, as here.

Many properties of the hierarchical growth of dark halos can be investigated by constructing Monte Carlo realizations of the evolution of the dark-halo population (a **merger tree**). Techniques for constructing merger trees using extended Press–Schechter theory are described by Sheth & Lemson (1999) and Somerville & Kolatt (1999).

Extended Press–Schechter theory has several limitations:

- (i) It uses the arbitrary prescription (9.19) to relate the mass of a halo to the scale K^{-1} . In numerical experiments the mass of a halo is measured, arbitrarily, within the virial radius r_{200} (see discussion following eq. 9.66).
- (ii) It cannot address the fate of halos once they become subhalos by falling into a larger halo, which is unfortunate given that the halos of satellite galaxies and galaxies in groups and clusters are all subhalos rather than primary halos.
- (iii) It provides no information about the spatial or velocity distribution of halos.

Despite these drawbacks, extended Press–Schechter theory is a powerful tool for understanding nonlinear structure formation. Its limitations are less remarkable than its successes.

9.2.4 Collapse and virialization in the cosmic web

In §9.2.2 we saw that the cosmic web largely comprises sheets into which material is falling on both sides. We can explore this process by idealizing the flow as perfectly one-dimensional and neglecting motion parallel to the plane of a sheet. Then we need to follow a set of sheets of material as they move along their common normal. We discretize the problem and suppose that we have $2n+1$ sheets, all with the same surface density, that are all perpendicular to the x axis and move parallel to this axis. The equations of motion of the sheets, which are given in Problem 7.3, can be solved analytically between sheet crossings. Hence a computer can integrate these equations exactly, apart from roundoff error (e.g., Yamashiro, Gouda, & Sakagami 1992).

We assume that all the sheets were at $x = 0$ for fixed time τ in the past (the Big Bang). They then expanded with the Hubble flow, were slowed by their mutual self-gravity and began to collapse. We start the integration just after the time of maximum expansion, before any sheets have crossed. For any given spatial distribution of the sheets, their velocities follow from the requirement that all sheets were located at the origin at the Big Bang. To introduce a small degree of initial inhomogeneity, we choose the locations of the sheets so that the coarse-grained density is proportional to

$[1 - 0.3 \cos(0.3\pi\xi)]^{-1}$, where ξ varies from -1 to 1 as one crosses the set of sheets.

Figure 9.14 shows, from top left to bottom right, six stages in the virialization of 401 sheets from these initial conditions. The elapsed time in units of τ is given in the top left corner of each panel. The negative slope in the top left panel implies that at the start of the simulation, the system is already collapsing. In the next panel ($t = 0.41$) the vertical orientation of the distribution near the origin indicates that the center has already finished collapsing. By the time of the third panel, the center has expanded and collapsed once more. The edge, by contrast, is only just starting to collapse for the second time. In this simple system, the phase of a particle's oscillation about the spatial origin is related to be the polar angle in the phase plane of Figure 9.14. The panels for $t > 2$ show that the phase lag $\Delta\psi$ between the particles near the center and particles near the outside grows rapidly. By the time of the last panel, $\Delta\psi$ has become so large near the center that it is hard to follow the spiral of phase points.

The two processes discussed in §4.10.2 are manifest in Figure 9.14: (i) the winding up of the line of phase points into an ever tighter spiral is phase mixing, and (ii) violent relaxation causes the edge of the occupied part of the phase plane to move outwards, and the points near the center to move towards the origin as the system's time-varying gravitational field transfers energy from the central to the peripheral sheets. For example, between the top left and top right panels, the sheets that are at $|v_x| \lesssim 0.2$ in the second panel have fallen together under their mutual self gravity, and are beginning to cross one another. Sheets further out have not begun to cross. However, when the inner sheets expand, the outer sheets are falling in past them, and the inner sheets have to climb out of a deeper well than they fell into. Conversely, the outer sheets fall into a well that has significantly weakened by the time they rise up the other side.

This transfer of energy from the inner to the outer sheets increases the density contrast between the center and the outside. Since the frequency of a sheet's oscillations through the center scales as the square root of the mean density $\bar{\rho}$ interior to it (eq. 2.40), the energy transfer enhances the rate at which the phase lag $\Delta\psi$ between the inner and outer particles grows. The energy transfer works most effectively between groups of sheets with phase lag $\Delta\psi \sim \pm\pi/4$. So, as the panels of Figure 9.14 clearly show, the characteristic distance scale of the transfer rapidly decreases. In an oversimplified picture, the first collapse transfers energy between the inner and the outer half of the sheets. The second collapse at the center transfers energy from the first quartile to the second quartile, and a little later energy is transferred from the second to the third quartile, and later still on out to the fourth quartile. By this time the additional central collapse pictured in the fourth panel of Figure 9.14 is transferring energy from the innermost $\frac{1}{8}$ of the sheets, and so on. At each stage the transfers become smaller because they involve a smaller and smaller fraction of the mass, but Figure 9.15 shows that they

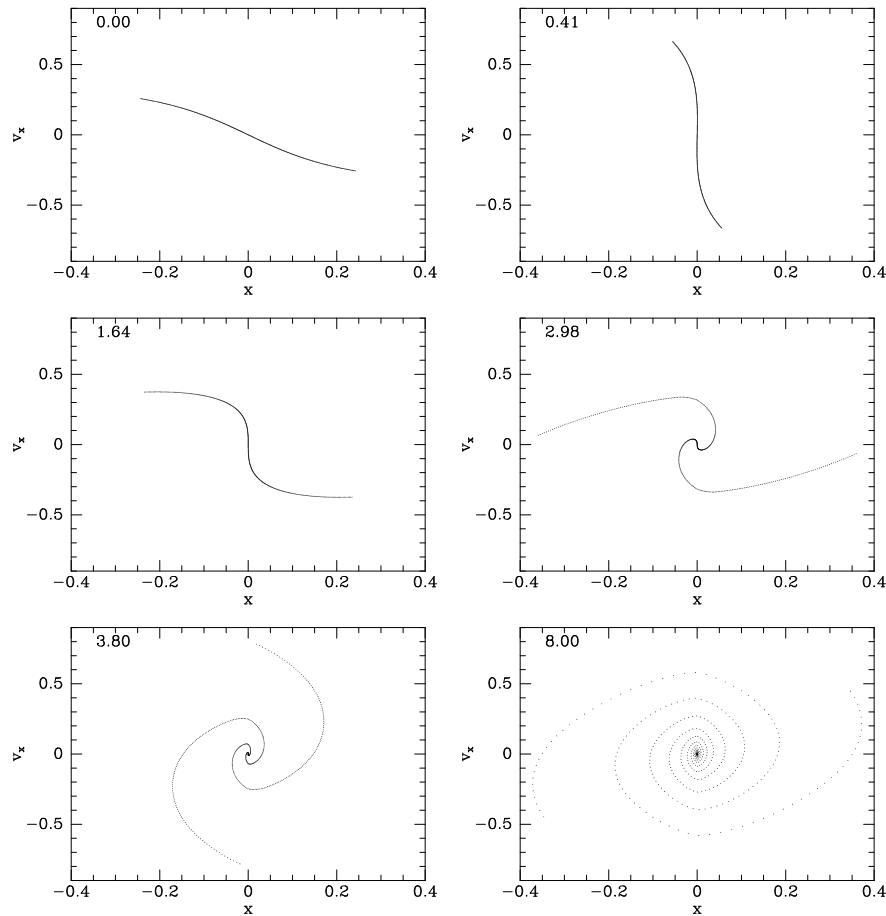


Figure 9.14 From top left to bottom right, six stages in the virialization of 401 sheets. The number in the top left corner of each panel is the time since the initial conditions were imposed. The units of time are the turnaround time τ of the underlying homogeneous distribution of sheets.

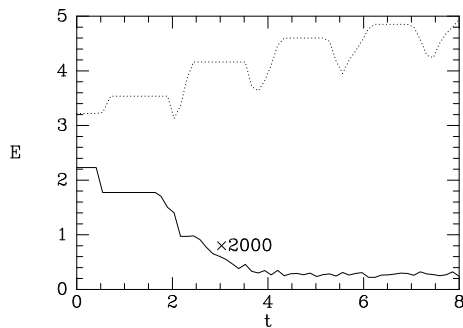


Figure 9.15 Energy of the innermost 12 sheets (full curve) and outermost 12 sheets (dotted curve) as a function of time. For clarity the energy of the inner sheets has been multiplied by 2000. The energy losses of the innermost 12 sheets are suppressed by the discreteness of the system from about $t = 4$. The energy is defined as in Problem 7.3.

remain significant to remarkably late times for the outer sheets. When only a finite number of sheets are used in the simulation, the loss of energy by the central sheets ceases once the phase difference between successive sheets exceeds $\sim \pi/4$ at the center, and the system settles to a configuration in which the coarse-grained central density is finite. If a continuum of sheets of infinitesimal surface density were used, the central sheets would continue to lose energy for longer, and the numerical experiments suggest that the final relaxed density profile would have a cusp in which the density would diverge as $x \rightarrow 0$ as $\rho \sim x^{-1/2}$ (Binney 2004b). This divergence—the formation of a central cusp in the density—is possible only with initially *cold* matter: the cusp has infinite phase-space density, and since phase mixing can only reduce the maximum coarse-grained phase-space density (§4.10.2a) a cusp can form only if the initial distribution also involves infinite phase-space density.

In our idealized model we have neglected motion parallel to the sheets of the cosmic web. In reality these motions will become significant soon after virialization perpendicular to the sheets is complete. The dominant effect is streaming within sheets towards the filaments at which several sheets join, and along the filaments to the nodes at which filaments meet. Halos grow at the nodes by accreting a stream from each filament. Since by this time the sheets have singular central densities, some of the infalling matter is exceedingly dense. We shall see below that three-dimensional simulations of the clustering process indicate that some of this dense material streams to the center of the halo and forms a cusp in which the density diverges as radius $r \rightarrow 0$.

The simulations we have just described provide a model for what happened immediately after the cosmic web first formed. As we saw in §9.2.2, ever larger-scale cosmic webs subsequently formed as the longer waves become nonlinear. In these subsequent stages of structure formation, the cosmic web is not made up of a smooth collisionless fluid of elementary particles, but is the lumpy aggregation of relaxed halos that formed as an earlier web collapsed. Hence we must imagine sheets of cuspy halos falling together and the virialization of the system will involve many mergers of objects that have cuspy density profiles.

9.3 N-body simulations of clustering

We now summarize the most important insights into the structure and dynamics of halos that have emerged from numerical simulations of the cosmological clustering of collisionless matter. Such simulations have been the primary theoretical tool for investigating dark-matter structure for a quarter of a century. Over this period the scale of these simulations has increased enormously—by 2005 state-of-the-art simulations contained 10^{10} particles in a box $\simeq 700$ Mpc on a side, in which particles interacted with a softening length of $\lesssim 2$ kpc. Nevertheless, the questions that can be definitively

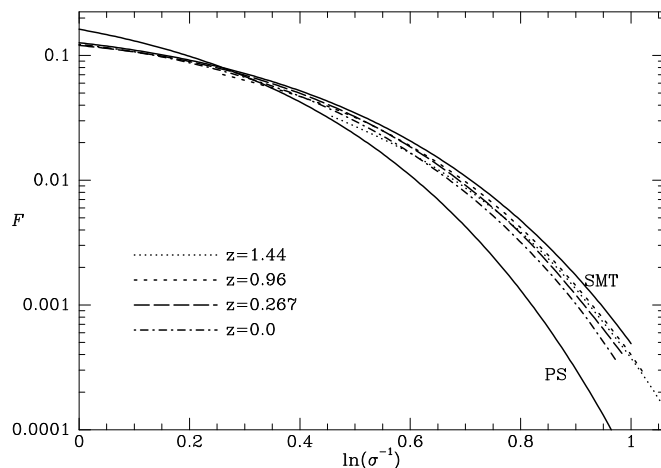


Figure 9.16 Mass functions of dark-matter halos from numerical simulations (broken curves) and from theory (full curves). On the vertical axis we plot $(M_K/\rho_0)|dn/d \ln \sigma_K^2|$ (eq. 9.78), while the horizontal axis shows values of σ_K at the relevant redshift, which in equation (9.79) is denoted $\tilde{\sigma}_K$. The curve marked PS is the prediction (9.78) of extended Press–Schechter theory, while the curve marked SMT is the prediction (9.83) of Sheth, Mo, & Tormen (2001). After Jenkins et al. (2001).

answered by simulations are still limited by finite resolution, both in mass (determined by the mass of an individual particle) and in distance (determined by a combination of particle density and softening length). These limitations must be borne in mind when examining small-scale features in simulations.

9.3.1 The mass function of halos

In §9.2.3 we obtained a prediction (9.77) for the number density of primary halos as a function of mass, and it is important to compare this with the results obtained from simulations. A key prediction of extended Press–Schechter theory was that the quantity $F \equiv (M_K/\rho_0)|dn/d \ln \sigma_K^2|$ should be independent of time t when plotted as a function of the dispersion at t , $\tilde{\sigma}_K$ (eq. 9.78). The broken curves in Figure 9.16 show for four redshifts the values of F that Jenkins et al. (2001) determined by counting primary halos in numerical simulations of the standard Λ CDM cosmology. To the degree that the curves from different redshifts overlaid one another, the prediction of extended Press–Schechter theory is vindicated. On the other hand, the theory’s functional form $F(\tilde{\sigma}_K)$, which is given by the extreme right side of equation (9.78) and in Figure 9.16 is shown by the curve labeled “PS”, is clearly an inadequate approximation to the form of F that emerges from the simulations. By generalizing the spherical-collapse model of §9.2.1 to the

case of collapsing ellipsoidal bodies, Sheth, Mo, & Tormen (2001) obtained a slightly different form of F :

$$F(\tilde{\sigma}_K^2) = 0.322 \left[1 + \left(\frac{\tilde{\sigma}_K}{\delta_1} \right)^{0.3} \right] \frac{\delta_1 / \tilde{\sigma}_K}{(2\pi)^{1/2}} \exp\left(-\frac{\delta_1^2}{2\tilde{\sigma}_K^2}\right) \quad (\delta_1 \equiv 1.418). \quad (9.83)$$

This prediction, shown by the curve marked “SMT” in Figure 9.16, is in much better agreement with the simulations.

The simulations show that halos contain many subhalos. Subhalos were not counted in either the analytical work of §9.2.3 or in the simulation results shown in Figure 9.16, but they are important for astronomy. The simplest heuristic model for subhalos is that the distribution of normalized subhalo masses m/M , where m and M are the masses of the subhalo and primary halo, is the same for all primary halos. That is, the number of subhalos per primary halo in some mass interval is $dn = f(m/M)d(m/M)$. Obviously, $f(m/M)$ must turn steeply down as m/M approaches unity, and moreover it is natural to assume that $f(m/M)$ is a power law when $m/M \ll 1$. These considerations suggest the simple form for the subhalo mass function

$$dn = c \left(\frac{m}{M} \right)^a \exp(-km/M) \frac{dm}{M}, \quad (9.84)$$

where c , a and k are constants (note the analogy to the Schechter law 1.18). This formula provides a remarkably good fit to N -body simulations of structure formation by Gao et al. (2004). The constant $a \simeq -2$, so the overall shape of the subhalo mass function (9.84) is similar to the mass distribution of primary halos, shown in Figure 9.12. The normalizing constant c is a weak function of primary halo mass, growing by a factor of two as the primary halo mass grows by 10^3 . Equation (9.84) predicts that the fraction of a primary halo’s mass that is contained in subhalos with m/M exceeding some small constant ϵ is independent of the mass of the primary halo; the simulations indicate that for $\epsilon = 10^{-4}$ this fraction is 5–10%.

9.3.2 Radial density profiles

At the smallest resolved radii, the densities of halos can be approximated by a power law $\rho \propto r^{-\alpha}$, where $\alpha \approx 1$ (Diemand, Moore, & Stadel 2004; Hayashi et al. 2004). Measurements of gravitational lensing of background galaxies suggest that in galaxy clusters the dark-matter density profile is indeed at least as steep as $\rho \propto r^{-1}$ (Sand et al. 2004). However, there is much debate about whether such steep density profiles are consistent with data for much smaller halos. In particular, it is not clear that such cusps are consistent with the circular-speed curves of low-luminosity, dark-matter dominated galaxies (McGaugh, Barker, & de Blok 2003), with the mass

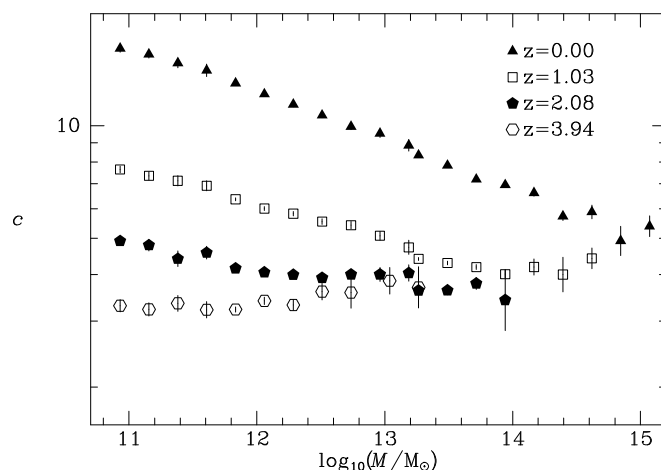


Figure 9.17 Median concentrations $c = r_{200}/a$ of halos of various masses at four redshifts (after Zhao et al. 2003).

distribution of the Milky Way inside the solar circle (Binney & Evans 2001; Klypin, Zhao, & Somerville 2002), with the pattern speeds of galactic bars (§8.1.1d; Debattista & Sellwood 2000), or with the evidence for maximum disks in galaxies (§6.3.3).

Outside the central cusp, the density profile of a halo depends on whether it is a primary halo or a subhalo. The density profile of a primary halo steepens to $\rho \propto r^{-3}$ at large r . Remarkably, it appears that to a good approximation all primary halos have a universal density profile, independent of cosmology, which can be fitted by the NFW model introduced in §2.2.2g. If $\rho \propto r^{-3}$, the mass within radius r diverges like $\ln r$ at large r (eq. 2.66). In reality, the halo ceases to be an equilibrium structure, and the NFW model no longer applies, around the virial radius r_{200} (§9.2.1). Figure 9.17 shows for halos of various masses at four redshifts the median concentrations $c \equiv r_{200}/a$, where a is the NFW scale radius. The filled triangles show that at the current epoch the concentration of a typical halo decreases with increasing halo mass. The pentagons and hexagons show that at earlier epochs the dependence of c on mass was weaker. It is easy to understand this result qualitatively: we have seen that only low-mass halos are formed at high redshift (Figure 9.8). If such a halo accretes relatively little since its formation and is still a low-mass halo, its inner structure, and therefore its scale radius a , will not have changed much. However, its virial radius r_{200} will have grown as the mean cosmic density fell, so c will have grown. If, by contrast, an initially low-mass halo becomes a massive halo, it will have experienced a succession of mergers with objects of comparable mass, giving violent relaxation an opportunity to increase a , so higher masses correspond

to lower concentrations.

Subhalos have density profiles that usually consist of a single power-law segment with slope $\alpha \simeq 1$ followed by an abrupt cutoff. Such a profile arises naturally through the tidal truncation of an NFW profile by the gravitational field of the primary halo (§8.3; Kazantzidis et al. 2004).

The origin of the universal density profile of primary dark halos is not well understood. Halos settle to this profile regardless of the cosmology and the power spectrum $P(k)$ (Navarro, Frenk, & White 1996). In particular, if the initial power spectrum is truncated, so $P(k) = 0$ for $k \geq K$, low-mass halos are eliminated but the density profiles of the primary halos that do form remain essentially unchanged (Huss, Jain, & Steinmetz 1999; Moore et al. 1999). In such a simulation, the cuspy central density profiles of the smallest halos cannot be fossils of small-scale structure in the initial conditions; rather, they must have arisen through virialization along the lines discussed in §9.2.4. By contrast, in simulations that start from a power spectrum that has substantial small-scale power, halos form largely through merging and accretion of smaller halos, yet they have the same density profile. Thus the explanation we seek must come in two parts: (i) why does virialization produce something like an NFW profile? and (ii) why does this profile reproduce itself when systems with NFW profiles merge?

The following argument offers an explanation of item (ii) above (cf. Dekel et al. 2003). Consider what happens when a satellite on a nearly circular orbit spirals into a much larger system. Assume that inside the Jacobi or tidal radius the satellite's density varies as $l^{-\beta}$, where l is distance from the satellite's center, and that the host's density varies with radius as $r^{-\alpha}$. Dynamical friction drains energy from the satellite's orbit, causing it to shrink (§8.1). Simultaneously, tides strip mass from the satellite (§8.3); this stripped material subsequently contributes to the density of the host around the current radius r of the satellite orbit. The masses interior to radius r in the host and radius l in the satellite are

$$M(r) = K_a r^{3-\alpha} \quad ; \quad m(l) = K_b l^{3-\beta}, \quad (9.85)$$

where K_a and K_b are unimportant constants. From equations (8.91) and (8.108) the radius to which the satellite has been stripped is therefore

$$l(r) = f \left(\frac{K_b}{3K_a} \right)^{1/\beta} r^{\alpha/\beta}, \quad (9.86)$$

where f is a factor of order unity that depends on α . The mass dm that is stripped as the orbit sinks through a radial interval dr is $(dm/dl)(dl/dr)dr$ and we should compare this with $(dM/dr)dr$, which is the mass that the host originally had in that radial range. If the ratio of these masses increases as r decreases, stripped material will be becoming more important as the center is approached, and the accretion event will steepen the host's density

profile, while this profile will be flattened in the opposite case. The crucial ratio of masses is

$$\mu \equiv \frac{dm}{dl} \frac{dr}{dM} \frac{dl}{dr} = \frac{\alpha(3-\beta)K_b l^{2-\beta}}{\beta(3-\alpha)K_a r^{2-\alpha}} f\left(\frac{K_b}{3K_a}\right)^{1/\beta} r^{\alpha/\beta-1} \propto r^{3(\alpha/\beta-1)}. \quad (9.87)$$

Thus if $\beta > \alpha$, and the satellite is the cuspier system, μ will increase as r decreases and the merger will steepen the host's density profile. Conversely, if $\beta < \alpha$, the merger will flatten the host's density profile. Thus halos that arise from the merging of many small halos will acquire the density profile of their small progenitors. Simulations of virialization tend to produce inner profiles with $\alpha \simeq 1$, as in an NFW model (Navarro, Frenk, & White 1995; Moore et al. 1999, 2004). Hence the universality of the NFW profile probably arises because (i) it is produced during virialization of the first halos and (ii) it survives subsequent mergers.

9.3.3 Internal dynamics of halos

Halos prove to be slowly rotating, triaxial stellar systems that are dynamically similar to luminous elliptical galaxies.

(a) The shapes of halos The shapes of halos are important for understanding many phenomena, including tidal streamers, warps, X-ray halos, and weak gravitational lensing. They have been investigated by many authors—Bailin & Steinmetz (2005) provide a useful summary of this work.

Halos are roughly ellipsoidal objects in the sense that the principal-axis directions determined by considering only matter that lies within distance r of the center, depend at most weakly on r . Figure 9.18 shows that halos are strongly triaxial: the left panel shows that the ratio c/a of the lengths of the shortest and longest principal axes clusters around 0.6, while the lower right panel plots the distribution of the **triaxiality parameter**

$$T \equiv \frac{a^2 - b^2}{a^2 - c^2}, \quad (9.88)$$

where b is the length of the intermediate axis. T varies from zero for oblate spheroids to unity for prolate halos and the right panel of Figure 9.18 shows that halos with $0.5 < T < 0.85$ are common while spheroidal halos, either oblate or prolate, are exceedingly rare. The scarcity of prolate halos is emphasized by the upper right panel, which shows the distribution of triaxiality parameters obtained when the axis lengths are chosen randomly between 0 and 1. This distribution peaks at $T = 1$ because when $b^2, c^2 \ll a^2$, $T \simeq 1$ regardless of the values taken by b and c . By contrast the distributions of triaxiality parameters of N-body halos is sharply depressed at $T = 1$.

If a halo still lies in a sheet of the cosmic web, the direction of the smallest principal axis has a strong tendency to be the normal to that sheet.

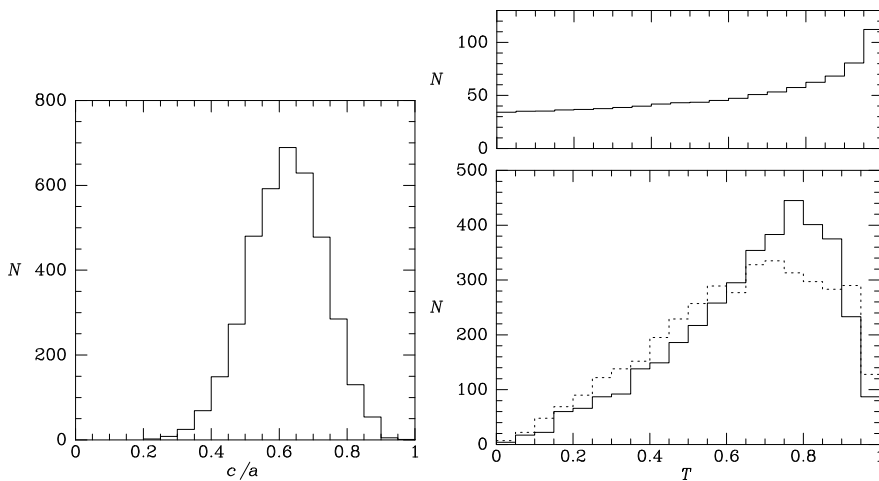


Figure 9.18 Left panel: distribution of the axis ratio c/a for halos that form in numerical cosmological N-body simulations. Lower right panel: the distribution of the triaxiality parameters (eq. 9.88) of these halos measured at radius $r = 0.12r_{200}$ (full line) and at $r = r_{200}$ (dashed curve). Upper right panel: the distribution of triaxiality parameters when the axis lengths are each chosen with uniform probability on $(0, 1)$ (after Bailin & Steinmetz 2005).

It is likely that the shape of a halo reflects the last major merger that was involved in its formation. For example, prolate shapes arise when halos of comparable mass merge from a nearly radial orbit, while oblate halos form when similar halos merge from a nearly circular orbit (Moore et al. 2004).

(b) Rotation of halos The dynamical importance of rotation for a self-gravitating system is quantified by the **spin parameter**

$$\lambda = \frac{J|E|^{1/2}}{GM^{5/2}}, \quad (9.89)$$

where M is the system's mass, E is its energy and J is its spin angular momentum. The spin parameter is a dimensionless number (Problem 9.7) that increases linearly with the rotation speed when the system's energy and shape are unchanged⁹ and it is constructed from quantities that are constant for an isolated system. Thus λ provides a measure of how rapidly a system is rotating that is independent of the system's mass and can be evaluated even before the system has virialized. For a non-rotating system $\lambda = 0$, while for a cold, razor-thin, self-gravitating exponential disk $\lambda = 0.4255$ (Problem 9.10).

⁹ We saw in §4.2.1b that for axisymmetric systems the part of the DF, f_- , that is odd in L_z contributes neither to the density nor to E , so such changes can be made by changing f_- .

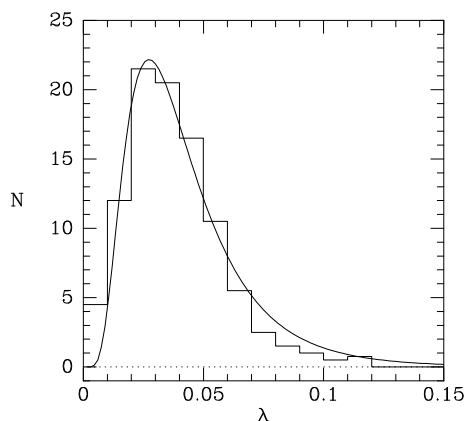


Figure 9.19 A histogram of the values of λ for halos in a simulated standard Λ CDM universe. The curve shows the log-normal distribution $dP(\lambda) \propto \exp[-\ln^2(\lambda/\lambda_0)/2\sigma^2]d\ln\lambda$ with $\lambda_0 = 0.037$ and $\sigma = 0.55$ (after Bullock et al. 2001).

The principal limitation of λ is that the values of J and M are both sensitive to the definition of the halo's outer boundary.

Figure 9.19 shows the distribution of λ -values for halos formed in a simulation of the standard Λ CDM cosmology. The distribution can be fitted by a log-normal distribution with median value $\lambda = 0.037$. This distribution appears to be independent of halo mass or environment (Barnes & Efstathiou 1987; Warren et al. 1992). Since the typical λ -value of a halo is much less than 0.5, we conclude that halos are not significantly flattened by rotation.

Bullock et al. (2001) investigated the radial distribution of angular momentum within halos that form in a standard Λ CDM cosmology. They found that the fraction $f(j)$ of the mass that has specific angular momentum less than j can be fitted by the simple formula

$$f(j) = \frac{\mu j}{j + j_0}, \quad (9.90)$$

where j_0 and μ are parameters. The function $f(j)$ is proportional to j for small j and increases more slowly at larger j . The maximum specific angular momentum j_{\max} is such that $f(j_{\max}) = 1$, so $j_{\max} = j_0/(\mu - 1)$. Typically $\mu \simeq 1.25$ and $j_0 \simeq 1.4\lambda r_{200}v_{200}$, where v_{200} is the circular speed at r_{200} , so $j_{\max} \simeq 5.6\lambda r_{200}v_{200}$.

The mean angular-momentum vector of material near the center generally lies within $\approx 25^\circ$ of the halo's shortest principal axis, but Bailin & Steinmetz (2005) find that there is very little correlation between the directions of the angular-momentum vectors of material near the center and on the periphery of a typical halo. These findings are consistent with the prediction of linear theory that the direction of the angular momentum that is being accreted by a halo varies dramatically over cosmic time (Quinn & Binney 1992). On account of tidal torques and infall, the total angular momentum of a halo varies strongly over cosmic time, both in direction and magnitude. Collisionless relaxation processes within a halo work constantly to align the

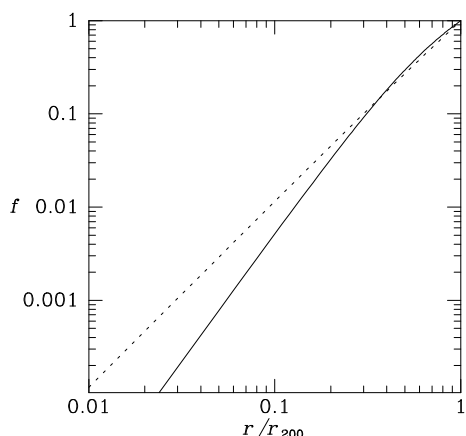


Figure 9.20 The full curve is a fit from Gao et al. (2004) to the fraction of subhalos contained within radius r , while the dashed line shows the fraction of the mass of an NFW halo that is inside r . The halo has been assumed to have concentration $c = 8$.

angular-momentum vectors of different mass shells, and they work fastest in the inner regions where the crossing time is shortest. Thus the overall effect is to produce a coherent inner portion that is constantly slewing its spin axis under the influence of torques produced by the less coherent outer regions (Binney & May 1986; Ostriker & Binney 1989). As we described in §6.6.1, disk warps may be a manifestation of this process.

(c) Dynamics of halo substructure As we saw in §9.3.1, primary halos have subhalos orbiting within them that contain up to 10% of the primary halo’s mass. Subhalos are constantly eroded by tides (§8.3) and tidal shocks (§8.2), and dragged inwards by dynamical friction (§8.1). The consequences of these processes are evident in subhalo statistics compiled by Gao et al. (2004):

- (i) A subhalo is much more likely to have been accreted recently than the halo’s matter as a whole. About 70% of subhalos fell into their primary halo since $z = 0.5$, and $\approx 90\%$ since $z = 1$.
- (ii) The ratio of the present mass of a subhalo to the mass it had when it fell in to the primary halo at redshift z decreases rapidly with increasing z , independent of the original subhalo mass.
- (iii) The spatial distribution of subhalos is more extended than that of all the halo’s mass (Figure 9.20). This finding presumably arises because subhalos at smaller radii are subject to stronger disruptive processes.

There is an important corollary of these results for the relation between the luminosity of a satellite galaxy or a galaxy in a cluster and the mass of its surrounding subhalo. The stellar parts of galaxies are much less affected by tides and other disruptive dynamical processes than their host halos because they are confined to the central regions of their halos. Consider two subhalos that start out with the same mass and galaxy luminosity, one at small radius and one at large radius in the primary halo. The subhalo at smaller radius will be eroded much faster than the one at larger radius, but the luminosities of the galaxies they contain will remain the same. So eventually, the relation

between galaxy luminosity and subhalo mass will vary strongly with position within the primary halo. In other words, there is no simple relation between the present distribution of subhalo masses and galaxy luminosities.

9.4 Star formation and feedback

Thus far we have been exclusively concerned with dark matter, which has simple physics but is so far observable only through its gravitational field. All our knowledge of the universe has been gleaned by studying particles, such as photons, neutrinos and high-energy nuclei, that have been emitted by “ordinary,” or “baryonic,” matter, such as the matter we and the Earth are made of. We now discuss what baryonic matter did as clustering developed.

(a) Reionization As we saw in §9.1, baryonic matter was tightly coupled to the photons of the CMB until the epoch of decoupling at $z \simeq 1000$, so baryonic fluctuations could not grow at a time when the dark matter was already clustering. Since the dark matter had this head start, after decoupling the baryons simply fell into the potential wells associated with pre-existing dark-matter overdensities.

Because the mass in baryons is much smaller than the mass in dark matter, the presence of the baryons probably had a very small impact on the development of the universe from decoupling until the first regions started to collapse, probably no earlier than $z \approx 30$ (Figure 9.8). The behavior of baryons during collapse was much richer and more dramatic than the response of the dark matter. Collisions between the neutral hydrogen and helium atoms that formed at decoupling both heated the gas and excited its atoms, causing them to emit photons and free electrons. The latter catalyzed the formation of hydrogen molecules, which can radiate at lower temperatures than atomic hydrogen. On account of the propensity of excited atoms and molecules to radiate, much of the gravitational energy released by the collapsing gas was lost, and gas sank towards the bottom of the potential well of whatever halo it was in. On account of this dissipation of energy—a process not available to dark matter—soon gas dominated the mass density at the centers of halos, even though the total halo mass was dominated by dark matter.

Continued radiation caused the density to run away to very high values in restricted volumes. Eventually these regions became optically thick to the photons radiated by gas atoms. Once photons became trapped, further collapse simply heated the gas. As the temperature rose, dynamical equilibrium became possible, and stars formed. We are not sure what these first stars looked like—the structure of present-day stars is heavily influenced by the trace amounts of heavy elements such as carbon and oxygen that they contain, while the first stars contained no such pollutants. Moreover, we do not know what the typical mass of a first-generation star was. However, the

CMB provides evidence that between $z \approx 20$ and 10 a great deal of ionizing radiation had been released by these stars, which suggests that many were massive objects ($M \gtrsim 10 M_{\odot}$).

As growing numbers of massive stars poured ionizing photons into the surrounding gas, the fraction of the universe containing ionized rather than neutral gas rapidly increased, until by $z \simeq 6$ almost all the volume of the universe was filled with ionized gas—we know this by studying lines in the spectra of distant quasars that are due to absorption of photons by neutral hydrogen and helium atoms in the intergalactic medium (Problem 1.14). This process is called **reionization** since the gas had previously been ionized until decoupling.

(b) Feedback When hydrogen and helium are photo-ionized, electrons are ejected at speeds characteristic of temperatures $\approx 10^4$ K. Hence, as the universe was reionized, the intergalactic medium, which still contained nearly all the baryons, was heated to $T \approx 10^4$ K.

At a temperature of 10^4 K the sound speed in hydrogen is $\simeq 10$ km s $^{-1}$, so gas at this temperature would not have fallen into gravitational potential wells that had escape speeds $v_e \lesssim 10$ km s $^{-1}$, which is the case for halos less massive than $\approx 10^8 M_{\odot}$ (Figure 9.9). Thus a halo less massive than $M \approx 10^8 M_{\odot}$ that had not already formed stars before the universe reionized will have missed its opportunity, and will not subsequently form stars. This argument suggests that there may be “dark-dark” halos that contain no stars (Dekel & Woo 2003).

We have only a limited understanding of star formation, in part because it is a consequence of the microstructure of the interstellar medium (ISM). Both observationally and theoretically, this is exceedingly hard to study, and we have only intimations of what the ISM may look like on the smallest scales. However, observations clearly show that (i) stars form from very cold ($T \lesssim 20$ K) molecular gas, and (ii) when a mass of cold gas gathers, the ensuing starburst (§8.5.5) converts only a small fraction of it into stars, blowing most of it away. The propensity of star formation to inject significant quantities of energy into the ISM is called **feedback**—this is negative feedback since the injected energy tends to quench subsequent star formation. A number of lines of observational evidence suggest that feedback is a remarkably effective process: for every unit of mass that is turned into stars, at least another unit of mass is driven out of a typical starbursting galaxy.

Supernovae can heat gas to temperatures $T \gtrsim 10^6$ K at which it can escape from halos with circular speeds $\lesssim 100$ km s $^{-1}$ (Dekel & Silk 1986). This energy input drives bulk motion of the gas at speeds of hundreds of km s $^{-1}$ on galactic scales (a **galactic wind**). The hot, fast-moving gas of the wind entrains colder, denser gas, which can be detected either through its emission lines (Strickland et al. 2004), or through its absorption of light from a background quasar (Pettini et al. 2001). Thus a starburst-driven wind is associated with an **outflow** that moves much more gas out of the galaxy than is directly involved in the wind.

The dynamics of the Local Group provides evidence for the importance of outflows. In Box 3.1 we concluded that the Local Group has mass $M \simeq 5 \times 10^{12} \mathcal{M}_\odot$, and in §1.3.5 we saw that $\simeq 20\%$ of matter is baryonic (eq. 1.75), so the Local Group should contain $1 \times 10^{12} \mathcal{M}_\odot$ of baryons. The Galaxy contains $\simeq 5 \times 10^{10} \mathcal{M}_\odot$ of baryonic mass (Table 1.2) and accounts for about a third of the Local Group's luminosity (BM Table 4.3). If the other Local-Group galaxies contribute a similar amount of light per unit baryonic mass, the galaxies of the Local Group must contain $\lesssim 2 \times 10^{11} \mathcal{M}_\odot$ of baryonic mass in total, no more than a quarter of the total baryonic mass of the Local Group. Observational constraints on the pressure of hot gas around the Galaxy require that the remaining gas must extend through a volume $\gtrsim 1$ Mpc in size (Problem 9.12). Thus fully three-quarters of the Local Group's baryons lie outside galaxies, probably in hot gas (Pedersen et al. 2006). This prediction can be reconciled with our picture of structure formation only if outflows have driven most of the baryons out of individual halos.

Another telling argument for outflows is that at least half of the heavy elements in rich clusters of galaxies are in the intracluster medium rather than in galaxies. Since the luminosity of the cluster is dominated by the galaxies, these elements must have been made in galaxies, and transferred to the intracluster medium by outflows. The synthesis of heavy elements was dominated by galaxies with luminosities $L \gtrsim L_*$, the characteristic luminosity of the galaxy luminosity function (§1.1.3), and these galaxies have high escape velocities ($v_e \gtrsim 500 \text{ km s}^{-1}$). Therefore massive outflows must be possible from galaxies with deep potential wells.

In the case of cluster galaxies, outflows can be driven in two ways. The starburst-driven winds described above are one mechanism, and the other is **ram-pressure stripping** by intergalactic gas (van Gorkom 2004). A typical cluster galaxy is moving through intracluster gas of density ρ at about $v \approx 1000 \text{ km s}^{-1}$, and the impact of this gas produces a ram pressure of order ρv^2 . This pressure sweeps the galactic gas into the turbulent wake that trails the galaxy. In practice star-formation powered outflows and ram-pressure stripping probably work in tandem: energy injected by star formation produces diffuse, extended gas, which is then easily removed by ram pressure. Note that ram-pressure stripping is effective only when a dense intracluster medium already exists as a result of winds, galaxy collisions or other mechanisms.

Like the thick disk of our own galaxy, the stellar populations of the early-type galaxies that dominate galaxy clusters have high abundances of α nuclides such as ^{16}O and ^{24}Mg relative to ^{56}Fe . As described on page 13, a high α/Fe ratio indicates that the era of star formation ended in $\lesssim 1$ Gyr. The natural explanation for this brief era of star formation is that cold gas was quickly removed by outflows powered by starbursts or ram pressure.

(c) Mergers, starbursts and quiescent accretion When two halos of comparable size merge, the gas in their disks is violently shocked. Consequently, the merger produces a starburst (§8.5.5). Stars that are formed in

such events are unlikely to be neatly arranged on circular, coplanar orbits, since the gas from which they form is not in quiescent rotation. Hence starbursts that occur during mergers probably give rise to bulges rather than stellar disks, and many bulges must have form in this way, and/or through the mixing together of pre-existing smaller disks.

So long as the cold ISM forms a quiescent disk, conservation of angular momentum keeps it away from the central black hole, so the latter cannot grow rapidly by accreting cold, dense gas. During a merger the gravitational field is far from axisymmetric, and by losing all its angular momentum (§8.5.5) some gas is likely to plunge to the center, and feed the black hole. Hence, the black hole is expected to grow rapidly at times of bulge formation; this may explain the observed tight correlation between black-hole mass and bulge luminosity and velocity dispersion (eq. 1.27) and also suggests a close relation between quasars and major mergers.

At least two other mechanisms are capable of forming bulges. First, in §6.6.2 we saw how the buckling instability can fatten a barred disk into a peanut-shaped bulge. Second, bulges can be made out of disks when a satellite galaxy plunges deep into its host. If the intruder is much less massive than the disk, the latter will become a thick disk such as the Milky Way possesses (page 13). If the intruder is massive, the disk will be shattered, and its debris will subsequently be part of a bulge. Thus we suspect that some bulge stars were born on non-circular orbits during merger-driven starbursts, while others were born in quiescent disks and subsequently scattered into their present orbits. The α/Fe ratio mentioned above provides a useful diagnostic, since a bulge formed from a disk that had been quiescent for $\gtrsim 1$ Gyr would have a near-solar α/Fe ratio.

Several lines of evidence indicate that spiral galaxies are continually acquiring gas, from accreted galaxies, and perhaps from the hot intergalactic medium discussed above. On account of the ability of gas to radiate, the incoming gas clouds will soon settle onto closed orbits in one of the galaxy's principal planes. In most cases this plane will coincide with the plane of the pre-existing stellar disk, but there are exceptions. In (rare) polar-ring galaxies (BM §8.2.5) the gas is observed to orbit in a plane perpendicular to the plane of the existing stellar disk. In even rarer cases, for example NGC 4550, the original disk and the accreted gas disk lie in the same plane but rotate in opposite senses. These exotic systems confirm that galaxies can be rejuvenated by acquiring a fresh source of gas, and illustrate the chaotic nature of galaxy formation.

In the inner regions of well-studied galaxies, baryons now dominate the mass budget (§§2.7 and 4.9.2). Consequently, the gravitational field in these regions is expected to differ significantly from that found in simulations that include only dark matter. If the baryons arrived in bursts, as during mergers, simulations that include the complex physics of baryons are required to predict the current gravitational field from cosmological initial conditions. By contrast, if the baryons arrived slowly through steady accretion, we can

determine the current force field by exploiting the adiabatic invariance of actions (§4.6.1b). The simplicity of the second approach has led to its being used more widely than is probably justified.

(d) The role of central black holes Figures 9.8 and 9.9 show that the characteristic mass of dark halos, M_c , and the associated characteristic velocity dispersion σ_c are constantly increasing. Consequently, the temperature that gas must attain if it is to escape from a typical halo, which is proportional to σ_c^2 , also increases. By contrast, the highest temperature to which stellar winds and supernovae can heat gas during a starburst is fixed at $\lesssim 10^7$ K, which yields $v_s \simeq 300$ km s $^{-1}$. Consequently, there is a critical halo mass, $M_{\text{trap}} \simeq 10^{13} M_\odot$ (Figure 9.9), above which gas is trapped, and supernova-heated gas simply accumulates in the halo’s potential well. Atmospheres of trapped hot gas in luminous elliptical galaxies have been detected and extensively studied through their thermal X-ray emission.

As star formation continues in a halo with $M > M_{\text{trap}}$, the density of hot, trapped gas rises. The gas can cool radiatively, and in response to cooling its density rises to maintain pressure balance. The radiative cooling time of the gas is $t_{\text{cool}} = \frac{3}{2} k_B T / \dot{E}$, where \dot{E} is the rate per particle at which the gas radiates. \dot{E} is proportional to the particle density n because radiation is caused by collisions of ions with free electrons, so t_{cool} is shortest at the center. Consequently, the first manifestation of cooling is an increase in the density of the gas at the center. A massive halo contains a black hole at its center (eq. 1.27), and this increased gas density boosts the rate at which the black hole accretes gas, and the rate of release of accretion energy. The details of this energy release are inadequately understood, but observations indicate that a significant fraction, in some cases almost all, of the energy emerges as a collimated outflow that heats the surrounding gas mechanically (e.g., Omma & Binney 2004).

Roughly three-quarters of rich clusters of galaxies host **cooling-flow** X-ray sources in which $t_{\text{cool}} \lesssim 300$ Myr at $r \simeq 10$ kpc. These systems must be several Gyr old, so this **cooling time** is at least an order of magnitude shorter than the lifetime, and even shorter cooling times must occur at smaller radii. Consequently, in the absence of heating these systems would develop infinite central densities within the next $\lesssim 100$ Myr. The X-ray morphologies of the systems are remarkably similar to one another (Donahue et al. 2006), which suggests that they are in approximate steady states, presumably because the central black hole is acting like a thermostatically controlled central-heating system. The cooling-flow phenomenon¹⁰ occurs in individual galaxies as well as in rich clusters of galaxies, but such cooling flows have been less thoroughly studied because they are much harder to observe.

The cooling-flow phenomenon has a major impact on galaxy formation by effectively quenching star formation in halos more massive than M_{trap} .

¹⁰ The name “cooling flow” is a misnomer derived from a defunct model: in these systems gas radiates but in a time-averaged sense neither cools nor flows inwards.

This happens because once a dense atmosphere at $T \gtrsim 10^7$ K has built up, thermal conduction and turbulent mixing carries heat from the hot atmosphere into any infalling clouds of cold gas, causing them to mix in with the hot atmosphere rather than forming a disk of cold gas in which stars can form. Notice that the heat absorbed by infalling cold clouds, like the energy radiated by the hot atmosphere, ultimately comes from accretion onto the central black hole.

(e) Origin of the galaxy luminosity function We finally return to the discrepancy between the galaxy luminosity function, shown by symbols in Figure 9.12, and the mass function of halos, given by the full curve in that figure. In that figure the galaxy luminosity function has been converted into a mass function by assuming that all galaxies have the same mass-to-light ratio, and that all of the mass of the universe is contained in galaxies (eq. 1.76). If the second of these assumptions is incorrect, we should shift the mass function of galaxies horizontally to the left. Whether or not such a shift is made, there are two irreconcilable differences between the mass functions of halos and galaxies: there are too many halos at both the largest and smallest masses.

These discrepancies can be resolved only if the mass-to-light ratio varies with halo mass. In particular, the most massive halos must have abnormally large mass-to-light ratios: in effect there is an upper limit to the luminosity of a galaxy, no matter how massive its halo is. Our discussion of how galaxies form and evolve suggests how this upper limit arises: halos more massive than M_{trap} have luminosities not much above L_* because stars ceased to form in them around the time that their luminosities reached L_* .

The discrepancy at low halo masses can in principle be resolved in two ways. Either low-mass halos have abnormally low mass-to-light ratios so their luminosities cluster around L_* , or they have abnormally high mass-to-light ratios, so they are faint or invisible. Our discussion of how galaxies form and evolve suggests which of these possibilities is likely to be true: low-mass halos have abnormally high mass-to-light ratios because a combination of efficient feedback and early heating of the intergalactic medium has kept gas out of their shallow potential wells. Thus at a qualitative level feedback from stars and active galactic nuclei can explain why the galaxy luminosity function differs strongly from the mass function of halos (Binney 2004a; Cattaneo et al. 2006).

9.5 Conclusions

Galaxies dominate the visual appearance of the universe, but over the last several decades we have concluded that they are only the tips of vast icebergs of dark matter that are in turn embedded in a mysterious sea of vacuum energy. Our primitive ideas as to the nature of dark matter and vacuum energy are speculative and possibly far from the truth, yet we understand

enough about these exotic materials to be able to assemble a coherent and elegant account of the evolution of galaxies and the universe with remarkable predictive power.

A convenient starting point is the redshift $z_{\gamma\text{m}} \simeq 3100$ when matter first dominated the overall mass budget of the universe. The CMB provides a clear window on the universe at the redshift of decoupling $z_{\text{d}} \simeq 1100$, when the scale factor was only three times larger. Through it we have seen the ripples in the cosmic density from which large-scale structure grew, and the power spectrum of these fluctuations is consistent with the theoretical predictions of Harrison and Zeldovich and of inflationary models of the very early universe.

Given these initial conditions, stellar dynamics, the subject of this book, is the only physics needed to follow the evolution of dark matter right up to the present epoch.

The story is more uncertain as regards the baryonic component of the universe, which is insignificant in terms of mass but dominates observational phenomena. The involvement of baryons in galaxy formation implies that a clear understanding of galaxy formation and evolution can be obtained only by complementing stellar dynamics with many other branches of astrophysics, including star formation, stellar evolution, the dynamics of interstellar gas and dust, accretion disks, and even black-hole physics. This is a vast enterprise that is in its infancy. A particularly important and challenging area is gas dynamics, which is involved in both star formation and its aftermath, and imprints on the world of galaxies characteristic scales that have their origin in atomic and nuclear physics.

We have sketched the main points of a theory of galaxy formation that may emerge in the years to come. Only time and much work will show whether our sketch is true to life.

Problems

9.1 [2] (a) Let $u(t)$ be a solution of $\ddot{y} + p(t)\dot{y} + q(t)y = 0$. By writing $y(t) = v(t)u(t)$ show that the general solution is

$$y(t) = Au(t) + Bu(t) \int^t dt' \frac{e^{-P(t')}}{u(t')^2}, \quad (9.91)$$

where A and B are arbitrary constants and $P(t) = \int^t dt' p(t')$.

(b) Show that the decaying solution of equation (9.37) for large-scale ($k \rightarrow 0$) density perturbations in a universe dominated by matter and vacuum energy is $\delta = \dot{a}/a$. If possible, you should derive this simple result using both mathematical and physical arguments. Hint: use equations (1.49) and (1.59).

(c) Show that for $k \rightarrow 0$ the growing solution of equation (9.37) is

$$\delta \propto \frac{\dot{a}}{a} \int_0^a \frac{da'}{a'^3}. \quad (9.92)$$

9.2 [1] The energy density and entropy density of black-body radiation are universal functions of the temperature, which we may write as $u_{\text{BB}}(T)$ and $s_{\text{BB}}(T)$. The pressure exerted by black-body radiation is $p = \frac{1}{3}\rho c^2 = \frac{1}{3}u_{\text{BB}}$ (eq. 1.58). Prove that $u_{\text{BB}}(T) = aT^4$ and $s_{\text{BB}}(T) = \frac{4}{3}aT^3$, where a is a universal constant. Hint: use the thermodynamic relation $dU = T dS - p dV$.

9.3 [2] The Lagrangian of a particle in comoving coordinates is given by equation (9.27).

(a) Show that the corresponding Hamiltonian is

$$H(\mathbf{x}, \mathbf{p}, t) = \frac{\mathbf{p}^2}{2a^2} - \frac{\dot{a}}{a} \mathbf{p} \cdot \mathbf{x} + \Phi(\mathbf{x}, t), \quad (9.93)$$

where $\mathbf{p} \equiv a^2 \dot{\mathbf{x}} + a\dot{a}\mathbf{x}$ is the momentum conjugate to \mathbf{x} .

(b) Show that after a suitable canonical transformation, the motion is described by the Hamiltonian

$$H'(\mathbf{x}, \mathbf{p}', t) = \frac{(\mathbf{p}')^2}{2a^2} + \Phi(\mathbf{x}, t) + \frac{1}{2}a\ddot{a}x^2, \quad (9.94)$$

where $\mathbf{p}' = a^2 \dot{\mathbf{x}}$. Hint: use equations (D.93) and (D.98).

(c) If the universe is homogeneous, show that the potential must have the form

$$\Phi(\mathbf{x}, t) = -\frac{1}{2}a\ddot{a}x^2 + \text{constant}. \quad (9.95)$$

Hint: consider the equations of motion for a particle moving with the Hubble flow.

(d) If the universe contains non-relativistic matter with mean density $\rho_0(t)$ and overdensity $\delta(\mathbf{x}, t)$ (eq. 9.1), argue that the potential must have the form

$$\Phi(\mathbf{x}, t) = -\frac{1}{2}a\ddot{a}x^2 + \phi(\mathbf{x}, t) \quad \text{where} \quad \nabla_{\mathbf{x}}^2 \phi = 4\pi G a^2 \rho_0 \delta \quad \text{with} \quad \nabla_{\mathbf{x}} \equiv \partial/\partial \mathbf{x}. \quad (9.96)$$

(e) Show that the equations of motion in comoving coordinates are

$$\dot{\mathbf{x}} = \frac{\mathbf{p}'}{a^2} \quad ; \quad \dot{\mathbf{p}}' = -\frac{\partial \phi}{\partial \mathbf{x}}. \quad (9.97)$$

(f) Find a symplectic integration algorithm (§3.4) to solve the equations of motion (9.97).

9.4 [3] (a) By considering the form taken by the energy-momentum tensor $T^{\alpha\beta}$ of a perfect fluid in the fluid's rest frame, explain why $T^{\alpha\beta}$ has to take the form $(\rho + p/c^2)u^\alpha u^\beta + pg^{\alpha\beta}$, where $x^0 = ct$, $\mathbf{u}(\mathbf{x})$ is the four-velocity of material at the event \mathbf{x} , the metric is taken to have signature $\mathbf{s} = (-1, 1, 1, 1)$, and ρ and p are the inertial density and pressure.

(b) In special relativity, the equations of hydrodynamics are obtained from the conservation law $\partial_\beta T^{\alpha\beta} = 0$, where $\partial_\beta \equiv \partial/\partial x^\beta$ and the metric $g^{\alpha\beta}$ is replaced by the Minkowski metric $\eta^{\alpha\beta} \equiv \delta_{\alpha\beta}$. For a fluid composed of relativistic particles (so p is of order ρc^2), show that in the absence of gravitational fields Euler's equation takes the form

$$\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} = -\frac{1 - v^2/c^2}{\rho + p/c^2} \left(\nabla p + \frac{\mathbf{v}}{c^2} \frac{\partial p}{\partial t} \right). \quad (9.98)$$

Hence derive equation (9.44).

9.5 [3] In this problem we use notation introduced in Problem 9.4 to determine the gravitational field generated by a slowly evolving distribution of relativistic matter when the field is weak. Since the field is weak, we can use coordinates that are nearly inertial, so we may take the metric to be of the form $g_{\alpha\beta} = \eta_{\alpha\beta} + h_{\alpha\beta}$, where $|h_{\alpha\beta}| \ll 1$.

(a) The four-velocity u^α of a freely-falling particle satisfies the equation of motion¹¹

$$\frac{du^\alpha}{d\tau} + \Gamma_{\mu\nu}^\alpha u^\mu u^\nu = 0, \quad (9.99)$$

where τ is the particle's proper time and the Christoffel symbol $\Gamma_{\mu\nu}^\alpha$ is related to the metric by

$$\Gamma_{\mu\nu}^\alpha = \frac{1}{2}g^{\alpha\beta} (\partial_\mu g_{\beta\nu} + \partial_\nu g_{\mu\beta} - \partial_\beta g_{\mu\nu}). \quad (9.100)$$

Show that for a slow-moving particle in a weak gravitational field the spatial components of this equation may be approximated by Newton's law of motion $\dot{u}_i = -\partial_i \Phi$ provided $h_{00} = -2\Phi/c^2$.

¹¹ Notice the analogy with the special-relativistic motion of a particle of mass m and charge q that moves in the electromagnetic field described by $F_{\mu\nu}$: $du^\alpha/d\tau - (q/m)F^\alpha{}_\mu u^\mu = 0$.

(b) It is often convenient to work with coordinates that satisfy the **harmonic gauge condition**

$$g^{\mu\nu}\Gamma_{\mu\nu}^{\alpha} = 0, \quad (9.101)$$

In such coordinates the Ricci tensor associated with a weak gravitational field may be approximated as $R_{\mu\nu} = \frac{1}{2}\partial_{\alpha}\partial^{\alpha}h_{\mu\nu}$. Given that Einstein's field equations are

$$R_{\mu\nu} = -\frac{8\pi G}{c^4}\left(T_{\mu\nu} - \frac{1}{2}g_{\mu\nu}T^{\gamma}_{\gamma}\right), \quad (9.102)$$

show that in the harmonic gauge the Einstein equations for the weak gravitational field generated by a perfect fluid (Problem 9.4) are

$$\partial_{\alpha}\partial^{\alpha}h_{\mu\nu} = -\frac{16\pi G}{c^4}\left[(\rho + p/c^2)u_{\mu}u_{\nu} + \frac{1}{2}g_{\mu\nu}(\rho c^2 - p)\right]. \quad (9.103)$$

Show that when the field is only slowly varying and we work in the fluid's rest frame, Φ satisfies the modified Poisson equation

$$\nabla^2\Phi = 4\pi G(\rho + 3p/c^2). \quad (9.104)$$

9.6 [2] Consider the gravitational collapse of an initially homogeneous and spheroidal cloud of pressureless fluid. A time t later, let \mathbf{x} be the position vector of the particle that originally had position vector \mathbf{q} . Show that if the components of \mathbf{x} are $x_i(\mathbf{q}, t) = \alpha_i(t)q_i$ (no summation on i), then the cloud remains homogeneous, and that for $\alpha_i(0) = 1$ the density at any time is $\rho(t) = \rho(0)/\prod_i \alpha_i(t)$. Show from the equations of motion of particles that (Lin, Mestel, & Shu 1965)

$$\ddot{\alpha}_i = -2\pi G\rho(0)\frac{A_i\alpha_i}{\alpha_1\alpha_2\alpha_3} \quad (\text{no summation on } i), \quad (9.105)$$

where $A_i(\boldsymbol{\alpha})$ is given by Table 2.1. By numerically integrating these coupled differential equations, show that an initially oblate spheroid with axis ratio 0.95 collapses to a disk whose radius is smaller than the original semi-major axis of the spheroid by a factor 0.0744. What is the connection between the dynamics of this system and the formation of the cosmic web?

9.7 [1] Let $-E$ be the energy and J the spin angular momentum of a system. Show that $GE^{\beta}J^{\gamma}$ cannot be dimensionless. Find values of α , β , and γ for which $GM^{\alpha}E^{\beta}J^{\gamma}$ is dimensionless.

9.8 [3] The original argument that led Press & Schechter (1974) to their formula for the distribution of halo masses was much simpler and more direct than that given in §9.2.3. From the probability that the overdensity smoothed on scale K^{-1} lies near δ , show that the probability $P_c(K, t)$ that at time t a given mass element is part of a collapsed structure of scale K^{-1} is

$$P_c(K, t) = \frac{1}{2} - \frac{1}{2}\operatorname{erf}\left(\frac{\delta_c(t_i, t)}{\sqrt{2}\sigma_K(t_i)}\right), \quad (9.106)$$

where δ_c is defined by (9.67). Hence show that the number density of halos with mass in $(M, M + dM)$ is

$$\frac{dn}{d\ln M} = \frac{\rho_0\delta_c/\sigma_K}{2(2\pi)^{1/2}M}\exp\left(-\frac{\delta_c^2}{2\sigma_K^2}\right)\left|\frac{d\ln\sigma_K^2}{d\ln M}\right|. \quad (9.107)$$

This formula for $dn/d\ln M$ agrees with (9.77) except for a factor of two. Which one is correct, and why?

9.9 [2] We seek solutions of the diffusion equation

$$\frac{\partial P}{\partial t} = D\frac{\partial^2 P}{\partial x^2} \quad \text{subject to the normalization condition} \quad \int_{-\infty}^{\infty} dx P(x, t) = 1.$$

(a) Use dimensional analysis to show that there should be solutions proportional to $Q(\xi)/t^{1/2}$, where $\xi \equiv x^2/(Dt)$ is a dimensionless variable.

(b) Show that Q satisfies

$$4\xi Q'' + (2 + \xi)Q' + \frac{1}{2}Q = 0. \quad (9.108)$$

(c) By obtaining a series solution of this equation, or otherwise, derive the function (9.72).

9.10 [1] Show that a centrifugally supported Kuzmin disk (§2.3.1) has spin parameter (eq. 9.89) $\lambda = \sqrt{8}/5 = 0.566$, and that a centrifugally supported exponential disk has $\lambda = 0.4255$.

9.11 [2] Let the rate per unit volume at which halos of mass M_1 and M_2 merge to form halos of mass $M_f = M_1 + M_2$ be $Q(M_1, M_2)n(M_1)n(M_2) dM_1 dM_2$, where $n(M) dM$ is the number density of halos of mass M . Show that extended Press–Schechter theory predicts that Q is given by

$$Q(M_1, M_2) = \frac{M_2}{\rho_0} \left| \frac{d \ln \delta_c}{dt} \frac{d\sigma_f^2}{dM_f} \frac{dM_2}{d\sigma_2^2} \right| \left(\frac{\sigma_1^2 \sigma_2^2 / \sigma_f^2}{\sigma_1^2 - \sigma_f^2} \right)^{3/2} \exp \left[-\frac{\delta_c^2}{2} \left(\frac{1}{\sigma_f^2} - \frac{1}{\sigma_1^2} - \frac{1}{\sigma_2^2} \right) \right], \quad (9.109)$$

where σ_1^2 is the variance of δ on mass scale M_1 , etc. Show that $Q(M_1, M_2)$ is not a symmetrical function of its arguments. Is this fact worrying? (Benson, Kamionkowski, & Hassani 2005).

9.12 [1] Spitzer¹² (1956) inferred from absorption lines in the spectra of stars at high Galactic latitudes that the Galaxy is surrounded by gas at the virial temperature with $n_e T \simeq 5 \times 10^8 \text{ K m}^{-3}$ near the Sun. Taking the Galactic potential to be spherical with $v_c = 220 \text{ km s}^{-1}$ at all radii, and the gas to be an isothermal mixture of fully ionized hydrogen and helium, with one ion in ten ^4He , show that the gas mass inside radius r is

$$M(r) \simeq 10^8 \mathcal{M}_\odot \left(\frac{10^6 \text{ K}}{T} \right) \int_0^{r/R_0} dx x^{2-\alpha},$$

where $\alpha \simeq 3.63(10^6 \text{ K}/T)$. Hence estimate the radius to which this atmosphere must extend if it is to contain that part of the Galaxy's share of the Local Group's baryons that is not in Galactic stars.

¹² As early as 1946 Lyman Spitzer (1914–1997) wrote a paper that explored the benefits of having a telescope in space. In 1951 he initiated work on magnetically confined nuclear fusion, the initially secret *Project Matterhorn* that later became the Princeton Plasma Physics Laboratory. He pioneered the kinetic theory of plasmas and the study of the interstellar medium, and made important contributions to our understanding of the dynamics of star clusters. From 1962 he led the project that culminated in NASA's *Copernicus* satellite, which from 1972–1981 opened up the far ultraviolet to astronomers. He played a leading role in the planning of and advocacy for the *Hubble Space Telescope*.

Appendices

Appendix A: Useful numbers

Physical constants¹

gravitational constant	$G = 6.6742(10) \times 10^{-11} \text{ m}^3 \text{ kg}^{-1} \text{ s}^{-2}$
speed of light	$c = 2.99792458 \times 10^8 \text{ m s}^{-1}$ (definition)
magnetic constant	$\mu_0 = 4\pi \times 10^{-7} \text{ N A}^{-2}$ $= 1.256637 \dots \times 10^{-6} \text{ N A}^{-2}$ (definition)
electric constant	$\epsilon_0 = (\mu_0 c^2)^{-1}$ $= 8.854188 \dots \times 10^{-12} \text{ F m}^{-1}$ (definition)
Planck constant	$h = 6.6260693(11) \times 10^{-34} \text{ J s}$ $\hbar = h/(2\pi)$ $= 1.05457168(18) \times 10^{-34} \text{ J s}$
Boltzmann's constant	$k_B = 1.3806505(24) \times 10^{-23} \text{ J K}^{-1}$
electron charge	$e = 1.60217653(14) \times 10^{-19} \text{ coulomb}$
proton mass	$m_p = 1.67262171(29) \times 10^{-27} \text{ kg}$
electron mass	$m_e = 9.1093826(16) \times 10^{-31} \text{ kg}$
Stefan–Boltzmann constant	$\sigma = \pi^2 k_B^4 / (60 \hbar^3 c^2)$ $= 5.670400(40) \times 10^{-8} \text{ W m}^{-2} \text{ K}^{-4}$
Thomson cross-section	$\sigma_T = e^4 / (6\pi \epsilon_0^2 m_e^2 c^4)$ $= 6.65245873(13) \times 10^{-29} \text{ m}^2$

¹ Taken from Mohr & Taylor (2005), Yao et al. (2006), Standish (1995), Cox (2000), BM, and Spergel et al. (2007). Numbers given in parentheses indicate one standard deviation uncertainty in the last digits of the preceding number.

Astronomical constants

astronomical unit	$1 \text{ AU} = 1.49597871475(30) \times 10^{11} \text{ m}$
parsec	$1 \text{ pc} = (648\,000/\pi) \text{ AU}$ $= 3.0856775975(6) \times 10^{16} \text{ m}$
year ²	$1 \text{ yr} = 3.15582 \times 10^7 \text{ s}$
solar mass	$\mathcal{M}_\odot = 1.9884(3) \times 10^{30} \text{ kg}$
heliocentric gravitational constant	$G\mathcal{M}_\odot = 1.32712440018(8) \times 10^{20} \text{ m}^3 \text{ s}^{-2}$
solar radius	$R_\odot = 6.9551(3) \times 10^8 \text{ m}$
solar luminosity (bolometric)	$L_\odot = 3.845(8) \times 10^{26} \text{ W}$
escape speed from Sun	$v_\star = (2G\mathcal{M}_\odot/R_\odot)^{1/2}$ $= 617.8 \text{ km s}^{-1}$
solar absolute magnitude	$M_V = +4.83 \quad M_R = +4.42$
Earth mass	$\mathcal{M}_\oplus = 3.0034896(2) \times 10^{-6} \mathcal{M}_\odot$ $= 5.9723(9) \times 10^{24} \text{ kg}$
Hubble constant	$H_0 = 70h_7 \text{ km s}^{-1} \text{ Mpc}^{-1}$ $h_7 = 1.05 \pm 0.05$
Hubble time	$H_0^{-1} = 13.969h_7^{-1} \text{ Gyr}$
age of the universe	$t_0 = (13.73 \pm 0.16) \text{ Gyr}$
critical density	$\rho_{c0} = 3H_0^2/(8\pi G)$ $= 9.2040(14) \times 10^{-27} h_7^2 \text{ kg m}^{-3}$ $= 1.359929 \times 10^{11} h_7^2 \mathcal{M}_\odot \text{ Mpc}^{-3}$

Useful relations

$$1 \text{ km s}^{-1} \simeq 1 \text{ pc per million years (actually 1.023)}$$

$$1 \text{ radian} = 206\,265 \text{ arcsec}$$

Appendix B: Mathematical background

The text presupposes an acquaintance with mathematical physics at the level of Jackson (1999) or Arfken & Weber (2005). This appendix contains a summary of some of the material and formulae that will be needed.

B.1 Vectors

The location of the point with Cartesian coordinates (x, y, z) may be described by a **position vector**,

$$\mathbf{x} = x\hat{\mathbf{e}}_x + y\hat{\mathbf{e}}_y + z\hat{\mathbf{e}}_z, \quad (\text{B.1})$$

where $\hat{\mathbf{e}}_x$, $\hat{\mathbf{e}}_y$, and $\hat{\mathbf{e}}_z$ are fixed unit vectors that point along the x , y and z axes. The distance of the point from the origin is written r or $|\mathbf{x}|$ and is equal to $(x^2 + y^2 + z^2)^{1/2}$.

Similarly, we represent an arbitrary vector \mathbf{A} in component form as

$$\mathbf{A} = A_x\hat{\mathbf{e}}_x + A_y\hat{\mathbf{e}}_y + A_z\hat{\mathbf{e}}_z. \quad (\text{B.2})$$

²In Galactic and extragalactic astronomy, “year” is an approximate time unit used for convenience. We have given the Gaussian year (the period of a test particle orbiting the Sun with a semi-major axis of 1 AU), which is the same as the sidereal year (the period of revolution of the Earth with respect to the fixed stars) at the quoted accuracy.

The magnitude of a vector \mathbf{A} is $A \equiv |\mathbf{A}| \equiv (A_x^2 + A_y^2 + A_z^2)^{1/2}$.

The **scalar** or **dot product** of two vectors \mathbf{A} and \mathbf{B} is

$$\mathbf{A} \cdot \mathbf{B} \equiv |\mathbf{A}||\mathbf{B}| \cos \psi, \quad (\text{B.3})$$

where ψ is the angle between the two vectors, placed tail to tail. Note that $\mathbf{A} \cdot \mathbf{B} = \mathbf{B} \cdot \mathbf{A}$ and $\mathbf{A} \cdot \mathbf{A} = |\mathbf{A}|^2$. Since $\hat{\mathbf{e}}_x \cdot \hat{\mathbf{e}}_x = \hat{\mathbf{e}}_y \cdot \hat{\mathbf{e}}_y = \hat{\mathbf{e}}_z \cdot \hat{\mathbf{e}}_z = 1$, and $\hat{\mathbf{e}}_x \cdot \hat{\mathbf{e}}_y = \hat{\mathbf{e}}_x \cdot \hat{\mathbf{e}}_z = \hat{\mathbf{e}}_y \cdot \hat{\mathbf{e}}_z = 0$, we may write the dot product in component form as

$$\mathbf{A} \cdot \mathbf{B} = \sum_{i=1}^3 A_i B_i, \quad (\text{B.4})$$

where the subscripts 1, 2, and 3 stand for x , y , and z , respectively. For simplicity we generally adopt the **summation convention**: we automatically sum from 1 to 3 over any dummy subscript that appears repeatedly in one term of an equation. Thus equation (B.4) may be written

$$\mathbf{A} \cdot \mathbf{B} = A_i B_i. \quad (\text{B.5})$$

The **vector** or **cross product** of two vectors is

$$\mathbf{A} \times \mathbf{B} \equiv AB \sin \psi \hat{\mathbf{p}}, \quad (\text{B.6})$$

where $\hat{\mathbf{p}}$ is a unit vector that is perpendicular to the plane containing \mathbf{A} and \mathbf{B} and points in the direction of movement of a right-hand screw when \mathbf{A} is rotated about the origin into \mathbf{B} . Note that $\mathbf{A} \times \mathbf{B} = -\mathbf{B} \times \mathbf{A}$, that $\mathbf{A} \times \mathbf{A} = 0$, and that $\hat{\mathbf{e}}_x \times \hat{\mathbf{e}}_y = \hat{\mathbf{e}}_z$, $\hat{\mathbf{e}}_y \times \hat{\mathbf{e}}_z = \hat{\mathbf{e}}_x$, $\hat{\mathbf{e}}_z \times \hat{\mathbf{e}}_x = \hat{\mathbf{e}}_y$. In component form the cross product may be written

$$\mathbf{A} \times \mathbf{B} = \epsilon_{ijk} \hat{\mathbf{e}}_i A_j B_k, \quad (\text{B.7})$$

where a sum over i , j , and k is implied by the summation convention. Here ϵ_{ijk} is the **permutation tensor** which is defined to be zero if any two or more of the indices i , j , and k are equal, +1 if (i, j, k) is an even permutation of $(1, 2, 3)$ [the even permutations are $(2, 3, 1)$ and $(3, 1, 2)$] and -1 if (i, j, k) is an odd permutation of $(1, 2, 3)$.

Three important identities that involve the dot and cross product are:

$$\mathbf{A} \cdot (\mathbf{B} \times \mathbf{C}) = \mathbf{C} \cdot (\mathbf{A} \times \mathbf{B}) = \mathbf{B} \cdot (\mathbf{C} \times \mathbf{A}), \quad (\text{B.8})$$

$$\mathbf{A} \times (\mathbf{B} \times \mathbf{C}) = (\mathbf{A} \cdot \mathbf{C})\mathbf{B} - (\mathbf{A} \cdot \mathbf{B})\mathbf{C}, \quad (\text{B.9})$$

$$(\mathbf{A} \times \mathbf{B}) \cdot (\mathbf{C} \times \mathbf{D}) = (\mathbf{A} \cdot \mathbf{C})(\mathbf{B} \cdot \mathbf{D}) - (\mathbf{A} \cdot \mathbf{D})(\mathbf{B} \cdot \mathbf{C}). \quad (\text{B.10})$$

In proving these identities it is useful to employ the relation

$$\epsilon_{ijk} \epsilon_{klm} = \delta_{il} \delta_{jm} - \delta_{im} \delta_{jl}, \quad (\text{B.11})$$

where the summation convention has been used (over index k) and δ_{pq} is defined to be 1 if $p = q$ and zero otherwise.

The velocity and acceleration of a particle may be written in Cartesian components as

$$\mathbf{v} \equiv \dot{\mathbf{x}} \equiv \frac{d\mathbf{x}}{dt} = \dot{x}\hat{\mathbf{e}}_x + \dot{y}\hat{\mathbf{e}}_y + \dot{z}\hat{\mathbf{e}}_z \quad ; \quad \mathbf{a} \equiv \ddot{\mathbf{x}} \equiv \frac{d^2\mathbf{x}}{dt^2} = \ddot{x}\hat{\mathbf{e}}_x + \ddot{y}\hat{\mathbf{e}}_y + \ddot{z}\hat{\mathbf{e}}_z. \quad (\text{B.12})$$

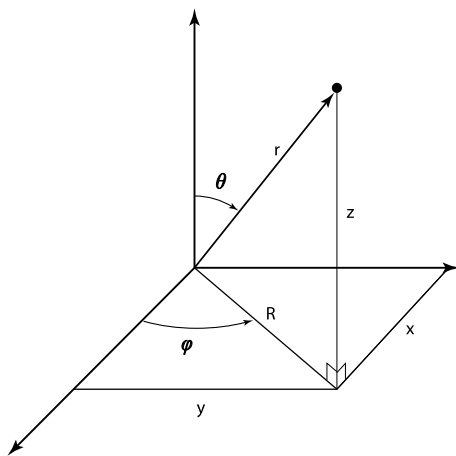


Figure B.1 The three main coordinate systems: Cartesian (x, y, z) , cylindrical (R, ϕ, z) , and spherical (r, θ, ϕ) .

B.2 Curvilinear coordinate systems

Let (q_1, q_2, q_3) denote the coordinates of a point in an arbitrary coordinate system. A fundamental quantity of any coordinate system is the **metric tensor** $h_{ij}(\mathbf{q})$, defined such that the distance ds between the points (q_1, q_2, q_3) and $(q_1 + dq_1, q_2 + dq_2, q_3 + dq_3)$ is given by

$$ds^2 = h_{ij} dq_i dq_j; \tag{B.13}$$

summation over i and j from 1 to 3 is implied by the summation convention. The coordinate systems used in this book are **orthogonal**, that is, $h_{ij} = 0$ if $i \neq j$. In this case we write $h_{ii} \equiv h_i^2$, so

$$ds^2 = h_i^2 dq_i^2. \tag{B.14}$$

The velocity is

$$\dot{\mathbf{x}} = \sum_i h_i \frac{dq_i}{dt} \hat{\mathbf{e}}_i, \tag{B.15}$$

where $\hat{\mathbf{e}}_i$ is a unit vector pointing in the direction from (q_1, q_2, q_3) to $(q_1 + dq_1, q_2, q_3)$, etc. The volume element in orthogonal coordinates is

$$d^3\mathbf{x} = h_1 h_2 h_3 dq_1 dq_2 dq_3. \tag{B.16}$$

For Cartesian coordinates, $(q_1, q_2, q_3) = (x, y, z)$ and $h_1 = h_2 = h_3 = 1$.

Cylindrical coordinate system In this system the location of a point is denoted by the triple (R, ϕ, z) , where R is the perpendicular distance from the z axis to the point, and ϕ is the azimuthal angle between the x axis and the projection of the position vector onto the (x, y) plane (Figure B.1). Thus the relation to Cartesian coordinates is

$$x = R \cos \phi \quad ; \quad y = R \sin \phi \quad ; \quad z = z. \tag{B.17}$$

In cylindrical coordinates the position vector is

$$\mathbf{x} = R \hat{\mathbf{e}}_R + z \hat{\mathbf{e}}_z. \tag{B.18}$$

Any vector may be written $\mathbf{A} = A_R \hat{\mathbf{e}}_R + A_\phi \hat{\mathbf{e}}_\phi + A_z \hat{\mathbf{e}}_z$, where $\hat{\mathbf{e}}_\phi = \hat{\mathbf{e}}_z \times \hat{\mathbf{e}}_R$ and

$$A_x = A_R \cos \phi - A_\phi \sin \phi \quad ; \quad A_y = A_R \sin \phi + A_\phi \cos \phi \quad ; \quad A_z = A_z. \tag{B.19}$$

The expressions for dot and cross products in cylindrical coordinates are simply equations (B.5) and (B.7), with the subscripts (1, 2, 3) denoting (R, ϕ, z) instead of (x, y, z) ; however, the decomposition into components must be carried out at the same position for both vectors in the product, since the directions of $\hat{\mathbf{e}}_R$ and $\hat{\mathbf{e}}_\phi$ depend on position.

The velocity in cylindrical coordinates is

$$\mathbf{v} = \frac{d\mathbf{x}}{dt} = \dot{R}\hat{\mathbf{e}}_R + R\dot{\phi}\hat{\mathbf{e}}_\phi + \dot{z}\hat{\mathbf{e}}_z. \quad (\text{B.20})$$

To compute $\dot{\hat{\mathbf{e}}}_R$ we use equation (B.19) with $A_R = 1$, $A_\phi = 0$, $A_z = 0$. Thus $\hat{\mathbf{e}}_R = \cos\phi\hat{\mathbf{e}}_x + \sin\phi\hat{\mathbf{e}}_y$, and $d\hat{\mathbf{e}}_R = (-\sin\phi\hat{\mathbf{e}}_x + \cos\phi\hat{\mathbf{e}}_y)d\phi$. The expression in parentheses is just $\hat{\mathbf{e}}_\phi$. After carrying out a similar analysis for $\hat{\mathbf{e}}_\phi$ we have

$$\frac{d\hat{\mathbf{e}}_R}{d\phi} = \hat{\mathbf{e}}_\phi \quad ; \quad \frac{d\hat{\mathbf{e}}_\phi}{d\phi} = -\hat{\mathbf{e}}_R, \quad (\text{B.21})$$

and

$$\dot{\hat{\mathbf{e}}}_R = +\dot{\phi}\hat{\mathbf{e}}_\phi \quad ; \quad \dot{\hat{\mathbf{e}}}_\phi = -\dot{\phi}\hat{\mathbf{e}}_R. \quad (\text{B.22})$$

Thus the velocity is

$$\mathbf{v} = \dot{R}\hat{\mathbf{e}}_R + R\dot{\phi}\hat{\mathbf{e}}_\phi + \dot{z}\hat{\mathbf{e}}_z. \quad (\text{B.23})$$

The acceleration is

$$\mathbf{a} = \frac{d\mathbf{v}}{dt} = (\ddot{R} - R\dot{\phi}^2)\hat{\mathbf{e}}_R + (2\dot{R}\dot{\phi} + R\ddot{\phi})\hat{\mathbf{e}}_\phi + \ddot{z}\hat{\mathbf{e}}_z. \quad (\text{B.24})$$

For cylindrical coordinates the metric tensor is given by (B.14) with

$$h_R = 1 \quad ; \quad h_\phi = R \quad ; \quad h_z = 1. \quad (\text{B.25})$$

The volume element is $d^3\mathbf{x} = R dR d\phi dz$.

Spherical coordinate system The position of a point in these coordinates is denoted by (r, θ, ϕ) (Figure B.1). The coordinate r is the radial distance from the origin to the point; θ is the angle between the position vector and the z axis; and ϕ is the same azimuthal angle used in cylindrical coordinates. The relation to Cartesian coordinates is

$$x = r \sin\theta \cos\phi \quad ; \quad y = r \sin\theta \sin\phi \quad ; \quad z = r \cos\theta. \quad (\text{B.26})$$

In spherical coordinates the position vector is simply

$$\mathbf{x} = r\hat{\mathbf{e}}_r. \quad (\text{B.27})$$

Any vector may be written $\mathbf{A} = A_r\hat{\mathbf{e}}_r + A_\theta\hat{\mathbf{e}}_\theta + A_\phi\hat{\mathbf{e}}_\phi$, where $\hat{\mathbf{e}}_\theta = \hat{\mathbf{e}}_\phi \times \hat{\mathbf{e}}_r$ and

$$\begin{aligned} A_x &= A_r \sin\theta \cos\phi + A_\theta \cos\theta \cos\phi - A_\phi \sin\phi, \\ A_y &= A_r \sin\theta \sin\phi + A_\theta \cos\theta \sin\phi + A_\phi \cos\phi, \\ A_z &= A_r \cos\theta - A_\theta \sin\theta. \end{aligned} \quad (\text{B.28})$$

Once again, the expressions for dot and cross products in spherical coordinates are equations (B.5) and (B.7), with the subscripts (1, 2, 3) denoting (r, θ, ϕ) instead of (x, y, z) , and with the understanding that the decomposition into components must be carried out at the same position for both vectors in the product.

The rate of change of the unit vectors is

$$\dot{\hat{\mathbf{e}}}_r = \dot{\theta}\hat{\mathbf{e}}_\theta + \dot{\phi}\sin\theta\hat{\mathbf{e}}_\phi \quad ; \quad \dot{\hat{\mathbf{e}}}_\theta = -\dot{\theta}\hat{\mathbf{e}}_r + \dot{\phi}\cos\theta\hat{\mathbf{e}}_\phi \quad ; \quad \dot{\hat{\mathbf{e}}}_\phi = -\dot{\phi}\sin\theta\hat{\mathbf{e}}_r - \dot{\phi}\cos\theta\hat{\mathbf{e}}_\theta. \quad (\text{B.29})$$

Thus the velocity is

$$\mathbf{v} = \dot{r}\hat{\mathbf{e}}_r + r\dot{\theta}\hat{\mathbf{e}}_\theta + r\sin\theta\dot{\phi}\hat{\mathbf{e}}_\phi. \quad (\text{B.30})$$

The acceleration is

$$\begin{aligned} \mathbf{a} = \frac{d\mathbf{v}}{dt} = & (\ddot{r} - r\dot{\theta}^2 - r\sin^2\theta\dot{\phi}^2)\hat{\mathbf{e}}_r + (2\dot{r}\dot{\theta} + r\ddot{\theta} - r\sin\theta\cos\theta\dot{\phi}^2)\hat{\mathbf{e}}_\theta \\ & + (r\sin\theta\ddot{\phi} + 2\sin\theta\dot{r}\dot{\phi} + 2r\cos\theta\dot{\theta}\dot{\phi})\hat{\mathbf{e}}_\phi. \end{aligned} \quad (\text{B.31})$$

For spherical coordinates the metric tensor is given by

$$h_r = 1 \quad ; \quad h_\theta = r \quad ; \quad h_\phi = r\sin\theta. \quad (\text{B.32})$$

The volume element is $d^3\mathbf{x} = r^2\sin\theta dr d\theta d\phi$.

For brevity the two angular coordinates are sometimes written as $\boldsymbol{\Omega} \equiv (\theta, \phi)$, and the element of area on the unit sphere is written $d^2\Omega \equiv \sin\theta d\theta d\phi$.

B.3 Vector calculus

Gradient In Cartesian coordinates we define the **gradient** of a scalar function of position $f(\mathbf{x})$ to be

$$\boldsymbol{\nabla} f \equiv \hat{\mathbf{e}}_x \frac{\partial f}{\partial x} + \hat{\mathbf{e}}_y \frac{\partial f}{\partial y} + \hat{\mathbf{e}}_z \frac{\partial f}{\partial z} = \hat{\mathbf{e}}_i \frac{\partial f}{\partial x_i}, \quad (\text{B.33})$$

where $(x_1, x_2, x_3) = (x, y, z)$ and we have used the summation convention. The symbol $\boldsymbol{\nabla}$ is called grad, del, or nabla and is considered to be a vector operator defined by

$$\boldsymbol{\nabla} = \hat{\mathbf{e}}_i \frac{\partial}{\partial x_i}. \quad (\text{B.34})$$

Where it is necessary to distinguish the variable used in the gradient, we write

$$\boldsymbol{\nabla}_{\mathbf{x}} \quad \text{or} \quad \frac{\partial}{\partial \mathbf{x}}. \quad (\text{B.35})$$

The change in the value of f between the points \mathbf{x} and $\mathbf{x} + d\mathbf{x}$ is

$$df = \frac{\partial f}{\partial x_i} dx_i = \boldsymbol{\nabla} f \cdot d\mathbf{x}. \quad (\text{B.36})$$

If \mathbf{x} and $\mathbf{x} + d\mathbf{x}$ lie on a surface S on which f is constant, that is, if $f(\mathbf{x}) = f(\mathbf{x} + d\mathbf{x})$, then $df = 0$ and $\boldsymbol{\nabla} f \cdot d\mathbf{x} = 0$, so $\boldsymbol{\nabla} f$ is orthogonal to $d\mathbf{x}$. Since $d\mathbf{x}$ can be chosen to be any infinitesimal vector lying in S , $\boldsymbol{\nabla} f$ must be normal to S itself. Hence the gradient of f is normal to surfaces of constant f .

In cylindrical coordinates the expression for the gradient must be modified, since $d\mathbf{x}$ is not $dR\hat{\mathbf{e}}_R + d\phi\hat{\mathbf{e}}_\phi + dz\hat{\mathbf{e}}_z$ but $dR\hat{\mathbf{e}}_R + R d\phi\hat{\mathbf{e}}_\phi + dz\hat{\mathbf{e}}_z$. Hence for consistency with equation (B.36) we must have

$$\boldsymbol{\nabla} = \hat{\mathbf{e}}_R \frac{\partial}{\partial R} + \frac{\hat{\mathbf{e}}_\phi}{R} \frac{\partial}{\partial \phi} + \hat{\mathbf{e}}_z \frac{\partial}{\partial z}. \quad (\text{B.37})$$

In spherical coordinates

$$\boldsymbol{\nabla} = \hat{\mathbf{e}}_r \frac{\partial}{\partial r} + \frac{\hat{\mathbf{e}}_\theta}{r} \frac{\partial}{\partial \theta} + \frac{\hat{\mathbf{e}}_\phi}{r\sin\theta} \frac{\partial}{\partial \phi}. \quad (\text{B.38})$$

The general form valid in any orthogonal coordinate system is

$$\nabla = \frac{\hat{\mathbf{e}}_i}{h_i} \frac{\partial}{\partial q_i}. \quad (\text{B.39})$$

We shall use the identity

$$\nabla_{\mathbf{x}}[\mathbf{a} \cdot (\mathbf{x} \times \mathbf{b})] = \nabla_{\mathbf{x}}[\mathbf{x} \cdot (\mathbf{b} \times \mathbf{a})] = \mathbf{b} \times \mathbf{a}, \quad (\text{B.40})$$

where \mathbf{a} and \mathbf{b} are constants.

Divergence In Cartesian coordinates we define the divergence of a vector-valued function $\mathbf{F}(\mathbf{x})$ to be the scalar

$$\nabla \cdot \mathbf{F} \equiv \frac{\partial F_x}{\partial x} + \frac{\partial F_y}{\partial y} + \frac{\partial F_z}{\partial z} = \frac{\partial F_i}{\partial x_i}. \quad (\text{B.41})$$

To clarify the physical meaning of the divergence, consider a volume V enclosed by a surface S . For simplicity, assume initially that the volume is cubical, occupying the region $x_{ia} \leq x_i \leq x_{ib}$, $i = 1, 2, 3$. Then

$$\begin{aligned} \int_V d^3\mathbf{x} \nabla \cdot \mathbf{F} &= \int_{x_{1a}}^{x_{1b}} dx_1 \int_{x_{2a}}^{x_{2b}} dx_2 \int_{x_{3a}}^{x_{3b}} dx_3 \left(\frac{\partial F_1}{\partial x_1} + \frac{\partial F_2}{\partial x_2} + \frac{\partial F_3}{\partial x_3} \right) \\ &= \int_{x_{2a}}^{x_{2b}} dx_2 \int_{x_{3a}}^{x_{3b}} dx_3 [F_1(x_{1b}, x_2, x_3) - F_1(x_{1a}, x_2, x_3)] + \text{two similar terms.} \end{aligned} \quad (\text{B.42})$$

This expression may be written more concisely as the area integral $\oint_S d^2\mathbf{S} \cdot \mathbf{F}$, where d^2S is the area of a small element of the surface and $d^2\mathbf{S}$ is a vector normal to the surface, pointing outward, with magnitude d^2S . (Note that in this notation $d^3\mathbf{x}$ is a scalar but $d^2\mathbf{S}$ is a vector.) We may generalize this result to an arbitrary volume by dividing the volume into many small cubes and noting that the surface integrals from the inside faces of the cubes cancel. Hence for an arbitrary volume V ,

$$\int_V d^3\mathbf{x} \nabla \cdot \mathbf{F} = \oint_S d^2\mathbf{S} \cdot \mathbf{F}. \quad (\text{B.43})$$

This result is known as the **divergence theorem** or Gauss's theorem (we reserve the latter term for the application of the divergence theorem to Poisson's equation; see eq. 2.12). In Cartesian coordinates we have

$$\int_V d^3\mathbf{x} \frac{\partial F_j}{\partial x_j} = \oint_S d^2S_j F_j. \quad (\text{B.44})$$

An immediate consequence of the divergence theorem is that for arbitrary scalar and vector functions g and \mathbf{F} ,

$$\int_V d^3\mathbf{x} g \nabla \cdot \mathbf{F} = \oint_S g \mathbf{F} \cdot d^2\mathbf{S} - \int_V d^3\mathbf{x} (\mathbf{F} \cdot \nabla)g, \quad (\text{B.45})$$

which is a three-dimensional analog of integration by parts.

If $\mathbf{F} = f\hat{\mathbf{e}}_i$, then we have $\int_V d^3\mathbf{x} (\partial f/\partial x_i) = \oint_S d^2S_i f$, which can be written more compactly after multiplication by $\hat{\mathbf{e}}_i$ and summing over indices as

$$\int_V d^3\mathbf{x} \nabla f = \oint_S d^2\mathbf{S} f. \quad (\text{B.46})$$

In cylindrical coordinates

$$\nabla \cdot \mathbf{F} = \frac{1}{R} \frac{\partial}{\partial R}(RF_R) + \frac{1}{R} \frac{\partial F_\phi}{\partial \phi} + \frac{\partial F_z}{\partial z}. \quad (\text{B.47})$$

This result can be derived from the divergence theorem, or by writing $\nabla \cdot \mathbf{F} = (\hat{\mathbf{e}}_R \partial / \partial R + \hat{\mathbf{e}}_\phi \partial / R \partial \phi + \hat{\mathbf{e}}_z \partial / \partial z) \cdot (F_R \hat{\mathbf{e}}_R + F_\phi \hat{\mathbf{e}}_\phi + F_z \hat{\mathbf{e}}_z)$ and expanding the expression using equation (B.21). In spherical coordinates

$$\nabla \cdot \mathbf{F} = \frac{1}{r^2} \frac{\partial}{\partial r} (r^2 F_r) + \frac{1}{r \sin \theta} \frac{\partial}{\partial \theta} (\sin \theta F_\theta) + \frac{1}{r \sin \theta} \frac{\partial F_\phi}{\partial \phi}. \quad (\text{B.48})$$

In arbitrary orthogonal coordinates

$$\nabla \cdot \mathbf{F} = \frac{1}{h_1 h_2 h_3} \left[\frac{\partial}{\partial q_1} (h_2 h_3 F_1) + \frac{\partial}{\partial q_2} (h_3 h_1 F_2) + \frac{\partial}{\partial q_3} (h_1 h_2 F_3) \right]. \quad (\text{B.49})$$

Laplacian The divergence of the gradient of a scalar function is called the **Laplacian** of that function. Thus the Laplacian of $F(\mathbf{x})$ is

$$\nabla^2 F \equiv \nabla \cdot (\nabla F). \quad (\text{B.50})$$

In different coordinate systems we have

$$\nabla^2 F = \frac{\partial^2 F}{\partial x^2} + \frac{\partial^2 F}{\partial y^2} + \frac{\partial^2 F}{\partial z^2}; \quad (\text{B.51})$$

$$\nabla^2 F = \frac{1}{R} \frac{\partial}{\partial R} \left(R \frac{\partial F}{\partial R} \right) + \frac{1}{R^2} \frac{\partial^2 F}{\partial \phi^2} + \frac{\partial^2 F}{\partial z^2}; \quad (\text{B.52})$$

$$\nabla^2 F = \frac{1}{r^2} \frac{\partial}{\partial r} \left(r^2 \frac{\partial F}{\partial r} \right) + \frac{1}{r^2 \sin \theta} \frac{\partial}{\partial \theta} \left(\sin \theta \frac{\partial F}{\partial \theta} \right) + \frac{1}{r^2 \sin^2 \theta} \frac{\partial^2 F}{\partial \phi^2}; \quad (\text{B.53})$$

$$\nabla^2 F = \frac{1}{h_1 h_2 h_3} \left[\frac{\partial}{\partial q_1} \left(\frac{h_2 h_3}{h_1} \frac{\partial F}{\partial q_1} \right) + \frac{\partial}{\partial q_2} \left(\frac{h_3 h_1}{h_2} \frac{\partial F}{\partial q_2} \right) + \frac{\partial}{\partial q_3} \left(\frac{h_1 h_2}{h_3} \frac{\partial F}{\partial q_3} \right) \right]. \quad (\text{B.54})$$

Convective operator We shall use the convective operator $(\mathbf{A} \cdot \nabla) \mathbf{B}$. In Cartesian coordinates we have

$$(\mathbf{A} \cdot \nabla) \mathbf{B} = \hat{\mathbf{e}}_i A_j \frac{\partial B_i}{\partial x_j}. \quad (\text{B.55})$$

In cylindrical coordinates

$$\begin{aligned} (\mathbf{A} \cdot \nabla) \mathbf{B} &= \left(A_R \frac{\partial B_R}{\partial R} + \frac{A_\phi}{R} \frac{\partial B_R}{\partial \phi} + A_z \frac{\partial B_R}{\partial z} - \frac{A_\phi B_\phi}{R} \right) \hat{\mathbf{e}}_R \\ &+ \left(A_R \frac{\partial B_\phi}{\partial R} + \frac{A_\phi}{R} \frac{\partial B_\phi}{\partial \phi} + A_z \frac{\partial B_\phi}{\partial z} + \frac{A_\phi B_R}{R} \right) \hat{\mathbf{e}}_\phi \\ &+ \left(A_R \frac{\partial B_z}{\partial R} + \frac{A_\phi}{R} \frac{\partial B_z}{\partial \phi} + A_z \frac{\partial B_z}{\partial z} \right) \hat{\mathbf{e}}_z. \end{aligned} \quad (\text{B.56})$$

In spherical coordinates

$$\begin{aligned} (\mathbf{A} \cdot \nabla) \mathbf{B} &= \left(A_r \frac{\partial B_r}{\partial r} + \frac{A_\theta}{r} \frac{\partial B_r}{\partial \theta} + \frac{A_\phi}{r \sin \theta} \frac{\partial B_r}{\partial \phi} - \frac{A_\theta B_\theta + A_\phi B_\phi}{r} \right) \hat{\mathbf{e}}_r \\ &+ \left(A_r \frac{\partial B_\theta}{\partial r} + \frac{A_\theta}{r} \frac{\partial B_\theta}{\partial \theta} + \frac{A_\phi}{r \sin \theta} \frac{\partial B_\theta}{\partial \phi} + \frac{A_\theta B_r}{r} - \frac{A_\phi B_\phi \cot \theta}{r} \right) \hat{\mathbf{e}}_\theta \\ &+ \left(A_r \frac{\partial B_\phi}{\partial r} + \frac{A_\theta}{r} \frac{\partial B_\phi}{\partial \theta} + \frac{A_\phi}{r \sin \theta} \frac{\partial B_\phi}{\partial \phi} + \frac{A_\phi B_r}{r} + \frac{A_\phi B_\theta \cot \theta}{r} \right) \hat{\mathbf{e}}_\phi. \end{aligned} \quad (\text{B.57})$$

In arbitrary orthogonal coordinates

$$(\mathbf{A} \cdot \nabla)\mathbf{B} = \hat{\mathbf{e}}_j \left[\frac{A_i}{h_i} \frac{\partial B_j}{\partial q_i} + \frac{B_i}{h_i h_j} \left(A_j \frac{\partial h_j}{\partial q_i} - A_i \frac{\partial h_i}{\partial q_j} \right) \right]. \quad (\text{B.58})$$

Green's theorem Let γ be a closed curve in the (x, y) plane, and S the area enclosed by this curve. Then **Green's theorem** states that for any functions $f(x, y)$, $g(x, y)$

$$\oint_{\gamma} [f(x, y) dx + g(x, y) dy] = \iint_S dx dy \left(\frac{\partial g}{\partial x} - \frac{\partial f}{\partial y} \right), \quad (\text{B.59})$$

where the line integral is taken anti-clockwise around the curve. This is demonstrated by first proving the theorem for the small square area $x \rightarrow x + \Delta x$, $y \rightarrow y + \Delta y$, then adding together many of these areas to fill up S .

In particular, let $f(x, y) = -y$, $g(x, y) = x$. Then

$$\oint_{\gamma} (x dy - y dx) = 2 \iint_S dx dy. \quad (\text{B.60})$$

Now¹ $\oint_{\gamma} y dx = -\oint_{\gamma} x dy$ so

$$\oint_{\gamma} x dy = \iint_S dx dy, \quad (\text{B.61})$$

which is the area enclosed by the curve γ .

B.4 Fourier series and transforms

Any function $f(x)$ that is periodic with period L , so that $f(x + L) = f(x)$, can be written as a **Fourier series**,

$$f(x) = \frac{1}{c} \sum_{n=-\infty}^{\infty} F_n \exp\left(\frac{2\pi i n x}{L}\right); \quad (\text{B.62})$$

normally the constant c is set to unity, but for some purposes (see below) it is more convenient to set $c = L$.

To find the coefficients F_n , we multiply this equation by $\exp(-2\pi i m x/L)$, where m is an integer, and integrate from $-L/2$ to $L/2$:

$$\int_{-L/2}^{L/2} dx f(x) \exp\left(-\frac{2\pi i m x}{L}\right) = \sum_{n=-\infty}^{\infty} \frac{F_n}{c} \int_{-L/2}^{L/2} dx \exp\left(\frac{2\pi i (n - m)x}{L}\right). \quad (\text{B.63})$$

The integral on the right side is zero for integer n and m unless $n = m$, in which case it equals L . Thus

$$F_n = \frac{c}{L} \int_{-L/2}^{L/2} dx f(x) \exp\left(-\frac{2\pi i n x}{L}\right). \quad (\text{B.64})$$

If $f(x)$ is real, then $F_n^* = F_{-n}$.

¹ To see this, let t be a parameter that varies from 0 to 1 as we travel around the curve γ , so the curve is defined by $[x(t), y(t)]$, $0 \leq t < 1$. Then integrating by parts we have $\oint_{\gamma} x dy = \int_0^1 dt x(t) dy(t)/dt = -\int_0^1 dt y(t) dx(t)/dt = -\oint_{\gamma} y dx$.

When x is the azimuthal angle ϕ , $L = 2\pi$ and it is convenient to write $F_n \equiv a_n \exp(-in\phi_n)$ where a_n is real and positive. Then for $c = 1$ equation (B.62) becomes

$$f(\phi) = \sum_{n=-\infty}^{\infty} a_n \exp[in(\phi - \phi_n)]. \quad (\text{B.65})$$

If $f(\phi)$ is real, then $a_{-n} = a_n$, $\phi_{-n} = \phi_n$, and

$$f(\phi) = a_0 + 2 \sum_{m=1}^{\infty} a_m \cos m(\phi - \phi_m). \quad (\text{B.66})$$

To derive analogous formulae for non-periodic functions, we consider the limit in which both $|n|$ and L become large. We replace n by a continuous variable $k \equiv 2\pi n/L$, replace LF_n/c by a continuous function $F(k)$, and replace the sum \sum_n in (B.62) by the integral $\int dn = L \int dk/(2\pi)$. Thus we arrive at the **Fourier transform**,

$$F(k) = \int_{-\infty}^{\infty} dx f(x)e^{-ikx} \quad ; \quad f(x) = \int_{-\infty}^{\infty} \frac{dk}{2\pi} F(k)e^{ikx}. \quad (\text{B.67})$$

Here the advantage of the choice $c = L$ is that F_n and $F(k = 2\pi n/L)$ are identical (see, for example, §9.1).

The Fourier transform is easily generalized to D -dimensional vectors:

$$F(\mathbf{k}) = \int d^D \mathbf{x} f(\mathbf{x})e^{-i\mathbf{k}\cdot\mathbf{x}} \quad ; \quad f(\mathbf{x}) = \int \frac{d^D \mathbf{k}}{(2\pi)^D} F(\mathbf{k})e^{i\mathbf{k}\cdot\mathbf{x}}, \quad (\text{B.68})$$

where the integral is over the entire D -dimensional space. In this book, the variable x or \mathbf{x} in the Fourier transform often denotes position, and then k or \mathbf{k} is called the **wavenumber** or **wavevector**.

We shall also use Fourier transforms in the time domain. Here we restrict our attention to functions $f(t)$ that vanish for $t < 0$. The temporal Fourier transform of $f(t)$ is then given by

$$F(\omega) = \int_0^{\infty} dt f(t)e^{i\omega t} \quad ; \quad f(t) = \int_{-\infty}^{\infty} \frac{d\omega}{2\pi} F(\omega)e^{-i\omega t}, \quad (\text{B.69})$$

where ω is the **frequency**. Note that the sign convention differs from (B.67) (there are good reasons for this). Note also that ω can be a complex number, in contrast to k or \mathbf{k} which is real. The use of complex ω enables us to generalize the temporal Fourier transform $F(\omega)$ to a wide class of functions $f(t)$ that diverge as $t \rightarrow \infty$. Suppose there exists a real constant $c > 0$ such that $f(t) \exp(-ct) \rightarrow 0$ as $t \rightarrow \infty$. Set $g(t) = f(t) \exp(-ct)$. Then

$$G(\omega) = \int_0^{\infty} dt f(t)e^{(i\omega - c)t} \quad ; \quad g(t) = \int_{-\infty}^{\infty} \frac{d\omega}{2\pi} e^{-i\omega t} G(\omega). \quad (\text{B.70})$$

Now set $\omega' \equiv \omega + ic$ and $F(\omega') \equiv G(\omega)$. The preceding equation becomes

$$F(\omega') = \int_0^{\infty} dt f(t)e^{i\omega' t} \quad ; \quad f(t) = \int_{ic-\infty}^{ic+\infty} \frac{d\omega'}{2\pi} F(\omega')e^{-i\omega' t}, \quad (\text{B.71})$$

and $F(\omega')$ is defined whenever $\text{Im}(\omega') \geq c$.

B.5 Abel integral equation

Let

$$f(x) = \int_x^\infty \frac{dt g(t)}{(t-x)^\alpha} \quad (0 < \alpha < 1). \quad (\text{B.72a})$$

Then

$$\begin{aligned} g(t) &= -\frac{\sin \pi \alpha}{\pi} \frac{d}{dt} \int_t^\infty \frac{dx f(x)}{(x-t)^{1-\alpha}} \\ &= -\frac{\sin \pi \alpha}{\pi} \int_t^\infty \frac{dx}{(x-t)^{1-\alpha}} \frac{df}{dx}. \end{aligned} \quad (\text{B.72b})$$

The first line can be proved by substituting equation (B.72a) into (B.72b), interchanging the order of integration, and using the integral

$$\int_0^1 \frac{du}{u^\alpha(1-u)^{1-\alpha}} = \frac{\pi}{\sin \pi \alpha}. \quad (\text{B.73})$$

The second line can be proved by replacing the integration variable x in the first expression for $g(t)$ by $u = x - t$ and then carrying out the differentiation with respect to t .

Another useful result, which can be proved similarly, is

$$f(x) = \int_0^x \frac{dt g(t)}{(x-t)^\alpha} \quad (0 < \alpha < 1), \quad (\text{B.74a})$$

$$g(t) = \frac{\sin \pi \alpha}{\pi} \frac{d}{dt} \int_0^t \frac{dx f(x)}{(t-x)^{1-\alpha}} = \frac{\sin \pi \alpha}{\pi} \left[\int_0^t \frac{dx}{(t-x)^{1-\alpha}} \frac{df}{dx} + \frac{f(0)}{t^{1-\alpha}} \right]. \quad (\text{B.74b})$$

B.6 Schwarz's inequality

For any two real functions $A(x)$ and $B(x)$,

$$\int dx A^2 \int dx B^2 \geq \left(\int dx AB \right)^2, \quad (\text{B.75})$$

with equality if and only if $A(x) = cB(x)$ for some constant c .

To prove this, let $C(x) = A(x) - \lambda B(x)$. Then $\int dx C^2$ is non-negative, so

$$\int dx C^2 = \int dx A^2 - 2\lambda \int dx AB + \lambda^2 \int dx B^2 \geq 0. \quad (\text{B.76})$$

Now set $\lambda = \int dx AB / \int dx B^2$. Equation (B.76) becomes

$$\int dx A^2 - \frac{(\int dx AB)^2}{\int dx B^2} \geq 0, \quad (\text{B.77})$$

which proves (B.75).

B.7 Calculus of variations

Consider a curve $\mathbf{x} = \mathbf{x}(t)$, $t_0 \leq t \leq t_1$, which we label by γ . We define a function

$$I(\gamma) \equiv \int_{t_0}^{t_1} dt L[\mathbf{x}(t), \dot{\mathbf{x}}(t), t]. \quad (\text{B.78})$$

Now consider a nearby curve γ' defined by $\mathbf{x} = \mathbf{x}(t) + \epsilon \mathbf{h}(t)$, where $\mathbf{h}(t_0) = \mathbf{h}(t_1) = 0$. As $\epsilon \rightarrow 0$ we have

$$\begin{aligned} I(\gamma') - I(\gamma) &= \int_{t_0}^{t_1} dt [L(\mathbf{x} + \epsilon \mathbf{h}, \dot{\mathbf{x}} + \epsilon \dot{\mathbf{h}}, t) - L(\mathbf{x}, \dot{\mathbf{x}}, t)] \\ &= \epsilon \int_{t_0}^{t_1} dt \left[\mathbf{h} \cdot \frac{\partial L}{\partial \mathbf{x}} + \dot{\mathbf{h}} \cdot \frac{\partial L}{\partial \dot{\mathbf{x}}} \right] + O(\epsilon^2), \end{aligned} \quad (\text{B.79})$$

where the integral on the second line is evaluated along the unperturbed curve γ . On this curve L can be considered to be a function only of time, $L(t) = L[\mathbf{x}(t), \dot{\mathbf{x}}(t), t]$. Hence we may integrate by parts to obtain

$$I(\gamma') - I(\gamma) = -\epsilon \int_{t_0}^{t_1} dt \mathbf{h} \cdot \left[\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{\mathbf{x}}} \right) - \frac{\partial L}{\partial \mathbf{x}} \right] + \epsilon \left(\mathbf{h} \cdot \frac{\partial L}{\partial \dot{\mathbf{x}}} \right)_{t_0}^{t_1} + O(\epsilon^2). \quad (\text{B.80})$$

The boundary term is seen to be zero since $\mathbf{h}(t_0) = \mathbf{h}(t_1) = 0$.

The curve γ is an **extremal** if $I(\gamma) - I(\gamma') = O(\epsilon^2)$ for all variations \mathbf{h} . Extremal curves have the largest or smallest values of the function I on any continuous curves connecting the fixed endpoints $\mathbf{x}(t_0)$ and $\mathbf{x}(t_1)$. At an extremal curve, the integral in equation (B.80) must vanish for all variations $\mathbf{h}(t)$ for which $\mathbf{h}(t_0) = \mathbf{h}(t_1) = 0$; thus the condition for an extremal curve is

$$\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{\mathbf{x}}} \right) - \frac{\partial L}{\partial \mathbf{x}} = 0. \quad (\text{B.81})$$

This is the **Euler–Lagrange equation** for an extremal curve I .

B.8 Poisson distribution

This distribution describes the statistical properties of a wide range of physical phenomena, but for concreteness we shall derive it in the context of the distribution of stars in phase space.

Suppose that the phase-space density of stars, f , is uniform over some phase-space volume \mathcal{V} , and that the distribution of stars is separable (§7.2.4) so the probability of finding a star at any phase-space position is unaffected by the presence or absence of other stars nearby. On average, we expect to find $r = f\mathcal{V}$ stars in this volume, but we want more detailed information: the probability p_n that exactly n stars are present.

Let us divide \mathcal{V} into a large number K of cells, each of volume $\Delta\mathcal{V} = \mathcal{V}/K$. If K is sufficiently large, the probability of finding two or more stars in any cell is negligible, and the probability of finding one star in a cell is just $f\Delta\mathcal{V} = r/K$. Hence the probability of finding n stars in \mathcal{V} is simply the probability that n cells are occupied by a star and $K - n$ are vacant. This is given by the binomial

distribution,

$$\begin{aligned}
 p_n &= \frac{K!}{n!(K-n)!} (f\Delta\mathcal{V})^n (1-f\Delta\mathcal{V})^{K-n} \\
 &= \frac{K(K-1)\cdots(K-n+1)}{n!} \frac{r^n}{K^n} \left(1 - \frac{r}{K}\right)^{K-n} \\
 &= \left(1 - \frac{1}{K}\right) \cdots \left(1 - \frac{n-1}{K}\right) \frac{r^n}{n!} \left(1 - \frac{r}{K}\right)^{K-n}.
 \end{aligned} \tag{B.82}$$

Now let $K \rightarrow \infty$; each of the factors $(1 - j/K)$ approaches unity, $(1 - r/K)^K \rightarrow \exp(-r)$ and $(1 - r/K)^{-n} \rightarrow 1$. Thus

$$p_n = \frac{r^n}{n!} e^{-r}, \tag{B.83}$$

which is the **Poisson distribution**. It is straightforward to show that $\sum_{n=1}^{\infty} p_n = 1$, as it must, and that the mean and variance of the Poisson distribution are

$$\langle n \rangle = \sum_{n=1}^{\infty} n p_n = r \quad ; \quad \langle (n - \mu)^2 \rangle = r. \tag{B.84}$$

More generally, the Poisson distribution describes the probability of obtaining n successes in a **Poisson process**, which is a large number K of trials in which the probability of success in a single trial is r/K , and the probability of success in different trials is independent.

B.9 Conditional probability and Bayes's theorem

The joint probability distribution $p(x, y)$ of two variables x and y is defined so that $p(x, y) dx dy$ is the probability that the first variable is found in the interval $(x, x + dx)$ and the second variable is in the interval $(y, y + dy)$. For many purposes we are interested in the conditional probability distribution $p(x|y)$, where $p(x|y) dx$ is the probability that the first variable is found in the interval $(x, x + dx)$ given that the second variable has the value y .

The relation between these two functions is fundamental to probability theory and statistical inference. Let X be the event "first variable is found in the interval $(x, x + dx)$," and Y the event "second variable is in the interval $(y, y + dy)$." The probability that both events occur is $P_{XY} = p(x, y) dx dy$. The probability that Y occurs, whatever the value of x may be, is $P_Y = dy \int dx' p(x', y) \equiv p_y(y) dy$. The probability that X occurs, given that Y occurs, is $P_{X|Y} = P_{XY}/P_Y$. Thus

$$p(x|y) = \frac{p(x, y)}{\int dx' p(x', y)} = \frac{p(x, y)}{p_y(y)}. \tag{B.85}$$

Similarly,

$$p(y|x) = \frac{p(x, y)}{\int dy' p(x, y')} = \frac{p(x, y)}{p_x(x)}; \tag{B.86}$$

eliminating $p(x, y)$ between equations (B.85) and (B.86) gives

$$p(y|x) = p(x|y) \frac{p_y(y)}{p_x(x)} = \frac{p(x|y)p_y(y)}{\int dy' p(x|y')p_y(y')}. \tag{B.87}$$

This is the celebrated **Bayes's theorem**, proved by Bayes and by Laplace in the eighteenth century. The formula is equally valid if x and y are vectors.

The most important application of Bayes's theorem is in statistical inference. Here y represents a set of model parameters and x represents a set of experimental or observational data. The model to be tested is encapsulated in $p(x|y)$, which predicts the probability of observing a given set of data x given a particular value for the model parameters y . The function $p_y(y)$ is the prior probability of y , which is our best estimate of the probability distribution of the model parameters before the data are taken—usually $p_y(y) \propto \text{constant}$ or $p_y(y) \propto y^{-1}$. Then $p(y|x)$ is the probability distribution of the model parameters given the data, from which the best estimate of y and associated error bars can be derived (e.g., Saha 2003).

B.10 Central limit theorem

This theorem governs the behavior of sums of large numbers of independent random variables. For a thorough treatment see Feller (1971).

Let u_1, u_2, \dots, u_n be independent random variables with distinct probability distributions $p_j(u)$; that is, the probability that u_j lies in the interval $(u_j, u_j + du_j)$ is $p_j(u_j)du_j$. Since $p_j(u)$ is a probability distribution, we must have $p_j(u) \geq 0$ and $\int du p_j(u) = 1$.

The first two moments of the distributions are

$$\int du u p_j(u) \equiv \mu_j \quad ; \quad \int du (u - \mu_j)^2 p_j(u) \equiv \sigma_j^2; \quad (\text{B.88})$$

these define the **mean** μ_j and the **variance** σ_j^2 or the **standard deviation** σ_j . For simplicity we assume that $\mu_j = 0$, although our results are easily generalized to distributions with non-zero mean, by working with the variables $u'_j \equiv u_j - \mu_j$. We define the normalized sum

$$s \equiv \frac{1}{\sqrt{n}}(u_1 + \dots + u_n). \quad (\text{B.89})$$

The central limit theorem states that as $n \rightarrow \infty$ the probability distribution of s approaches

$$g(s) \equiv \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{s^2}{2\sigma^2}\right), \quad \text{where} \quad \sigma^2 \equiv \frac{1}{n} \sum_{j=1}^n \sigma_j^2, \quad (\text{B.90})$$

which is the **normal** or **Gaussian** probability distribution. Thus the sum $\sum_{i=1}^n u_i$ is also normally distributed, with variance $n\sigma^2$ equal to the sum of the individual variances.

Proof: Since the u_j are statistically independent, the probability that they lie in the volume $du_1 \cdots du_n$ around (u_1, \dots, u_n) is $p_1(u_1) \cdots p_n(u_n) du_1 \cdots du_n$. Thus the probability that the normalized sum lies in the interval $(s, s + ds)$ is $q(s)ds$, where

$$q(s) = \int du_1 \cdots du_n p_1(u_1) \cdots p_n(u_n) \delta\left(s - \frac{u_1 + \dots + u_n}{\sqrt{n}}\right); \quad (\text{B.91})$$

here δ denotes the delta function and we have used equation (C.7). Taking the

Fourier transform (B.67), we have

$$\begin{aligned} Q(\omega) &\equiv \int ds g(s) e^{-i\omega s} \\ &= \int du_1 \cdots du_n p_1(u_1) \cdots p_n(u_n) \exp \left[\frac{-i\omega(u_1 + \cdots + u_n)}{\sqrt{n}} \right] \\ &= \prod_{j=1}^n \int du p_j(u) e^{-i\omega u/\sqrt{n}}. \end{aligned} \quad (\text{B.92})$$

As $n \rightarrow \infty$ at fixed u , the exponent in the last equation tends to zero, so we can expand the exponential in a Taylor series,

$$\begin{aligned} \int du p_j(u) e^{-i\omega u/\sqrt{n}} &= \int du p_j(u) \left[1 - i\omega u/\sqrt{n} - \frac{1}{2}(\omega u)^2/n + O(n^{-3/2}) \right] \\ &= 1 - \frac{\omega^2 \sigma_j^2}{2n} + O(n^{-3/2}), \end{aligned} \quad (\text{B.93})$$

where the last equation employs our assumption that the mean $\mu_j = 0$. Substituting equation (B.93) into (B.92) and taking the logarithm, we have

$$\ln Q(\omega) = \sum_{j=1}^n \ln \left[1 - \frac{1}{2} \omega^2 \sigma_j^2 / n + O(n^{-3/2}) \right]. \quad (\text{B.94})$$

We expand the right side in its Taylor series to find

$$\ln Q(\omega) = -\frac{\omega^2}{2n} \sum_{j=1}^n \sigma_j^2 + O(n^{-1/2}) \quad (\text{B.95})$$

so as $n \rightarrow \infty$

$$Q(\omega) = \exp \left(-\frac{\omega^2}{2n} \sum_{j=1}^n \sigma_j^2 \right). \quad (\text{B.96})$$

Taking the inverse Fourier transform (B.67),

$$q(s) = \int \frac{d\omega}{2\pi} Q(\omega) e^{i\omega s} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}s^2/\sigma^2}, \quad (\text{B.97})$$

where $\sigma^2 = n^{-1} \sum_{j=1}^n \sigma_j^2$. This is the normal distribution (B.90), as required. ◁

The central limit theorem can be generalized to vectors. Let $\mathbf{u}^1, \dots, \mathbf{u}^n$ be independent random variables with dimension D and probability distribution $p^j(\mathbf{u}^j)$. The moments of the distribution are

$$\int d^D \mathbf{u} u_i p^j(\mathbf{u}) \equiv \mu_i^j \quad ; \quad \int d^D \mathbf{u} (u_i - \mu_i^j)(u_k - \mu_k^j) p^j(\mathbf{u}) \equiv C_{ik}^j, \quad (\text{B.98})$$

where $i, k = 1, \dots, D$, $j = 1, \dots, n$, and \mathbf{C}^j is the **covariance matrix**.

For simplicity we assume again that the means $\boldsymbol{\mu}^j = 0$. Then the normalized sum $\mathbf{s} \equiv \sum_{j=1}^n \mathbf{u}^j/\sqrt{n}$ approaches the multivariate Gaussian distribution

$$g(\mathbf{s}) = \frac{1}{(2\pi)^{D/2} |\mathbf{C}|^{1/2}} \exp \left(-\frac{1}{2} \mathbf{s} \cdot \mathbf{M} \cdot \mathbf{s} \right), \quad (\text{B.99})$$

where $|\mathbf{C}|$ is the determinant of \mathbf{C} and

$$\mathbf{C} \equiv \frac{1}{n} \sum_{j=1}^n \mathbf{C}^j \quad ; \quad \mathbf{M} = \mathbf{C}^{-1}. \quad (\text{B.100})$$

Appendix C: Special functions

Complete descriptions of the special functions of mathematical physics are given in many books (Morse & Feshbach 1953; Abramowitz & Stegun 1964; Jackson 1999; Gradshteyn & Ryzhik 2000; Arfken & Weber 2005) and on web sites such as mathworld.wolfram.com.

Throughout this appendix, x and z represent real and complex variables respectively.

C.1 Delta function and step function

The **delta function** is a singular function defined by the properties

$$\delta(x) = 0 \quad \text{for } x \neq 0 \quad ; \quad \int_{-\infty}^{\infty} dx f(x) \delta(x) = f(0), \quad (\text{C.1})$$

where $f(x)$ is an arbitrary function continuous at $x = 0$. The delta function is sometimes called the Dirac delta function, after its inventor P.A.M. Dirac. The definition (C.1) implies that

$$\delta[f(x)] = \sum_j \frac{\delta(x - x_j)}{|f'(x_j)|}, \quad (\text{C.2})$$

where x_j are the solutions of $f(x_j) = 0$.

The three-dimensional delta function is a function of $\mathbf{x} = (x, y, z)$ defined by

$$\delta(\mathbf{x}) = \delta(x) \delta(y) \delta(z). \quad (\text{C.3})$$

The delta function can be written as

$$\delta(x) = \frac{1}{\pi} \lim_{\epsilon \rightarrow 0} \frac{|\epsilon|}{\epsilon^2 + x^2} = \frac{1}{\pi} \lim_{a \rightarrow \infty} \frac{\sin ax}{x} = \frac{1}{2\pi} \int_{-\infty}^{\infty} dx e^{ixt}. \quad (\text{C.4})$$

The **Plemelj identity** is

$$\lim_{\epsilon \rightarrow 0} \int_a^b \frac{dx f(x)}{x - y + i|\epsilon|} = \wp \int_a^b \frac{dx f(x)}{x - y} - i\pi f(y), \quad (\text{C.5})$$

where $f(x)$ is an arbitrary continuous function, x and y are real, $a < y < b$, and \wp denotes the Cauchy principal value,

$$\begin{aligned} \wp \int_a^b \frac{dx f(x)}{x - y} &\equiv \lim_{\epsilon \rightarrow 0} \left(\int_a^{y-|\epsilon|} \frac{dx f(x)}{x - y} + \int_{y+|\epsilon|}^b \frac{dx f(x)}{x - y} \right) \\ &= \lim_{\epsilon \rightarrow 0} \int_a^b dx f(x) \frac{x - y}{(x - y)^2 + \epsilon^2}. \end{aligned} \quad (\text{C.6})$$

The Plemelj identity is easy to prove using the last line in (C.6) and the first of the expressions in (C.4).

If $p(x_1, \dots, x_N) dx_1 \cdots dx_N$ is the number of objects with parameters in the range $x_1 \rightarrow x_1 + dx_1, \dots, x_N \rightarrow x_N + dx_N$, and G is some function of x_1, \dots, x_N , then the number of objects with G in the range $g \rightarrow g + dg$ is $n(g) dg$, where

$$n(g) = \int dx_1 \cdots dx_N p(x_1, \dots, x_N) \delta[g - G(x_1, \dots, x_N)]. \quad (\text{C.7})$$

The **step function** is defined by

$$H(x) = \begin{cases} 0 & (x < 0), \\ 1 & (x > 0). \end{cases} \quad (\text{C.8})$$

Obviously,

$$\frac{dH(x)}{dx} = \delta(x). \quad (\text{C.9})$$

C.2 Factorial or gamma function

$$z! \equiv \Gamma(z+1) \equiv \int_0^\infty dt t^z e^{-t} \quad (\text{Re } z > -1). \quad (\text{C.10})$$

By analytic continuation the function $z!$ can be defined for all complex numbers z , except for simple poles at $-1, -2, \dots$. Some useful relations are

$$z! = z(z-1)!, \quad (\text{C.11})$$

$$(z-1)!(-z)! = \pi \csc \pi z, \quad (\text{C.12})$$

$$(2z)! = \frac{2^{2z}}{\sqrt{\pi}} z! (z - \frac{1}{2})!, \quad (\text{C.13})$$

$$z!^* = (z^*)!, \quad (\text{C.14})$$

$$x! = \sqrt{2\pi} x^{x+1/2} e^{-x} [1 + O(x^{-1})] \quad \text{as } x \rightarrow \infty; \quad (\text{C.15})$$

the last of these is **Stirling's approximation** to $x!$ for large x .

Special values of $x!$ include

$$\begin{aligned} n! &= 1 \cdot 2 \cdot 3 \cdots n & ; & \quad 0! = 1! = 1; \\ (-\frac{1}{2})! &= \sqrt{\pi} = 1.77245 & ; & \quad (\frac{1}{2})! = \frac{1}{2}\sqrt{\pi} = 0.88623. \end{aligned} \quad (\text{C.16})$$

C.3 Error function, Dawson's integral, and plasma dispersion function

The **error function** is defined by

$$\text{erf } z \equiv \frac{2}{\sqrt{\pi}} \int_0^z dt e^{-t^2} = \frac{2}{\sqrt{\pi}} \sum_{n=0}^{\infty} \frac{(-1)^n z^{2n+1}}{n!(2n+1)}. \quad (\text{C.17})$$

We have

$$\text{erf}(0) = 0; \quad \text{erf}(\infty) = 1; \quad \text{erf}(-z) = -\text{erf}(z); \quad \text{erf}(z^*) = [\text{erf}(z)]^*. \quad (\text{C.18})$$

$$\text{As } x \rightarrow \infty, \quad 1 - \text{erf } x \rightarrow \frac{e^{-x^2}}{\sqrt{\pi}x} [1 + O(x^{-2})]. \quad (\text{C.19})$$

The error function is closely related to **Dawson's integral**,

$$F_{\pm}(x) \equiv e^{\mp x^2} \int_0^x dy e^{\pm y^2}; \tag{C.20}$$

in particular,

$$F_-(z) = \frac{1}{2}\sqrt{\pi}e^{z^2} \operatorname{erf}(z) \quad ; \quad F_+(z) = -\frac{1}{2}\sqrt{\pi}ie^{-z^2} \operatorname{erf}(iz). \tag{C.21}$$

On the interval $(0, \infty)$, the function $F_+(x)$ has a maximum of 0.54104 at $x = 0.92414$. As $x \rightarrow \infty$, $F_+(x) \rightarrow 1/(2x)$.

The **plasma dispersion function** (Fried & Conte 1961) is

$$Z(z) = i\sqrt{\pi}e^{-z^2} [1 + \operatorname{erf}(iz)]. \tag{C.22}$$

An alternative expression is

$$Z(z) = \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} ds \frac{e^{-s^2}}{s-z} \quad (\operatorname{Im} z > 0), \tag{C.23}$$

and its analytic continuation for $\operatorname{Im} z \leq 0$. For real and imaginary argument,

$$Z(x) = i\sqrt{\pi}e^{-x^2} - 2F_+(x) \quad ; \quad Z(iy) = i\sqrt{\pi}e^{y^2} (1 - \operatorname{erf} y). \tag{C.24}$$

For $y > 0$,

$$Z(x - iy) = Z^*(x + iy) + 2i\sqrt{\pi}e^{-(x-iy)^2}. \tag{C.25}$$

We have also

$$Z(0) = i\sqrt{\pi} \quad ; \quad Z(z^*) = -[Z(-z)]^* \quad ; \quad \frac{dZ}{dz} = -2 - 2zZ(z). \tag{C.26}$$

As $|z| \rightarrow \infty$,

$$Z(z) = i\sqrt{\pi}\sigma e^{-z^2} - \sum_{n=0}^{\infty} \frac{(n - \frac{1}{2})!}{\sqrt{\pi}z^{2n+1}}, \tag{C.27}$$

where $\sigma = 0, 1, 2$ if $\operatorname{Im}(z)$ is positive, zero or negative.

C.4 Elliptic integrals

The incomplete elliptic integrals of the first and second kinds are

$$F(\theta, k) \equiv \int_0^{\theta} \frac{d\phi}{\sqrt{1 - k^2 \sin^2 \phi}} \quad ; \quad E(\theta, k) \equiv \int_0^{\theta} d\phi \sqrt{1 - k^2 \sin^2 \phi}. \tag{C.28}$$

The complete elliptic integrals of the first and second kinds are

$$K(k) \equiv F(\frac{1}{2}\pi, k) = \int_0^{\pi/2} \frac{d\phi}{\sqrt{1 - k^2 \sin^2 \phi}} = \int_0^1 \frac{dt}{\sqrt{(1-t^2)(1-k^2t^2)}}, \tag{C.29}$$

$$E(k) \equiv E(\frac{1}{2}\pi, k) = \int_0^{\pi/2} d\phi \sqrt{1 - k^2 \sin^2 \phi} = \int_0^1 dt \sqrt{\frac{1 - k^2t^2}{1 - t^2}}.$$

Note that $K(0) = E(0) = \pi/2$, $E(1) = 1$, and $K(k)$ diverges logarithmically as $k \rightarrow 1$, in that $K(k) - \frac{1}{2} \ln[16/(1 - k^2)] \rightarrow 0$.

C.5 Legendre functions

The Legendre functions of the first and second kinds, $P_\lambda^\mu(z)$ and $Q_\lambda^\mu(z)$, are linearly independent solutions of the differential equation

$$\frac{d}{dz} \left[(1-z^2) \frac{dw}{dz} \right] - \frac{\mu^2}{1-z^2} w + \lambda(\lambda+1)w = 0. \quad (\text{C.30})$$

For $\text{Re}(\lambda) > 0$, the Legendre functions of the first kind diverge ($\propto z^\lambda$) as $|z| \rightarrow \infty$, while the functions of the second kind vanish [$\propto z^{-(\lambda+1)}$]. As $z \rightarrow 0$

$$\begin{aligned} \left[\frac{d \ln P_\lambda^\mu(z)}{dz} \right]_{z=0} &= 2 \tan\left[\frac{1}{2}\pi(\lambda+\mu)\right] \frac{[\frac{1}{2}(\lambda+\mu)]![\frac{1}{2}(\lambda-\mu)]!}{[\frac{1}{2}(\lambda+\mu-1)]![\frac{1}{2}(\lambda-\mu-1)]!}, \\ \left[\frac{d \ln Q_\lambda^\mu(z)}{dz} \right]_{z=0} &= 2 \exp\left\{\frac{1}{2}\pi i \operatorname{sgn}[\operatorname{Im}(z)]\right\} \frac{[\frac{1}{2}(\lambda+\mu)]![\frac{1}{2}(\lambda-\mu)]!}{[\frac{1}{2}(\lambda+\mu-1)]![\frac{1}{2}(\lambda-\mu-1)]!}. \end{aligned} \quad (\text{C.31})$$

Here $\operatorname{sgn}(x) = +1$ if $x > 0$ and -1 if $x < 0$.

For many applications we are interested in the Legendre functions with real arguments, $z = x$, in the interval $-1 \leq x \leq 1$. Unless μ is an even integer, $P_\lambda^\mu(x+i\epsilon)$ and $P_\lambda^\mu(x-i\epsilon)$ are different for real x and ϵ as $\epsilon \rightarrow 0$. Thus it is conventional to redefine the Legendre functions for $-1 \leq x \leq 1$ by

$$\begin{aligned} P_\lambda^\mu(x) &\equiv \frac{1}{2} \lim_{\epsilon \rightarrow 0} \left[e^{\pi i \mu/2} P_\lambda^\mu(x+i|\epsilon|) + e^{-\pi i \mu/2} P_\lambda^\mu(x-i|\epsilon|) \right] \\ Q_\lambda^\mu(x) &\equiv \frac{1}{2} e^{-i\pi\mu} \lim_{\epsilon \rightarrow 0} \left[e^{-\pi i \mu/2} Q_\lambda^\mu(x+i|\epsilon|) + e^{\pi i \mu/2} Q_\lambda^\mu(x-i|\epsilon|) \right]. \end{aligned} \quad (\text{C.32})$$

For $\mu = 0$ and λ a non-negative integer, the Legendre functions are polynomials given by the formula

$$P_l(x) \equiv P_l^0(z) = \frac{1}{2^l l!} \frac{d^l}{dx^l} (x^2-1)^l. \quad (\text{C.33})$$

These **Legendre polynomials** are also generated by the relation

$$\frac{1}{\sqrt{1-2xt+t^2}} = \sum_{l=0}^{\infty} P_l(x) t^l \quad |t| < 1, \quad |x| \leq 1, \quad (\text{C.34})$$

which leads to an expression for the inverse distance between the points \mathbf{x} and \mathbf{x}' ,

$$\frac{1}{|\mathbf{x} - \mathbf{x}'|} = \sum_{l=0}^{\infty} \frac{r_{<}^l}{r_{>}^{l+1}} P_l(\cos \gamma), \quad (\text{C.35})$$

where $r_{<} = \min(|\mathbf{x}|, |\mathbf{x}'|)$, $r_{>} = \max(|\mathbf{x}|, |\mathbf{x}'|)$, and γ is the angle between the two vectors.

For integer $m > 0$ and integer $l \geq 0$ the Legendre functions are sometimes called **associated Legendre functions**, and are given by¹

$$P_l^m(x) = (-1)^m (1-x^2)^{m/2} \frac{d^m P_l^0(x)}{dx^m} = (-1)^m \frac{(1-x^2)^{m/2}}{2^l l!} \frac{d^{l+m} P_l^0(x)}{dx^{l+m}}. \quad (\text{C.36})$$

¹ Our convention for associated Legendre functions with real x between -1 and $+1$ follows Abramowitz & Stegun (1964), Press et al. (1986), Gradshteyn & Ryzhik (2000), and software such as IDL, Maple, and Mathematica, but differs from Morse & Feshbach (1953), Arfken & Weber (2005), and the first edition of this book by a factor $(-1)^m$.

Note that $P_l^m(x)$ vanishes for $m > l$, and that $P_l^m(x)$ is even in x if $l - m$ is even, and odd if $l - m$ is odd. We have

$$P_l^m(0) = (-1)^{(l+m)/2} \frac{(l+m)!}{2^l [\frac{1}{2}(l-m)]! [\frac{1}{2}(l+m)]!} \quad (l-m \text{ even}), \quad (C.37)$$

and zero if $(l - m)$ is odd. For integer $m > 0$,

$$P_l^{-m}(x) = (-1)^m \frac{(l-m)!}{(l+m)!} P_l^m(x). \quad (C.38)$$

The associated Legendre functions are orthogonal in the sense that

$$\int_{-1}^1 dx P_l^m(x) P_n^m(x) = \frac{2}{2l+1} \frac{(l+m)!}{(l-m)!} \delta_{ln}. \quad (C.39)$$

$$\int_{-1}^1 \frac{dx}{1-x^2} P_l^m(x) P_l^k(x) = \frac{1}{m} \frac{(l+m)!}{(l-m)!} \delta_{mk}. \quad (C.40)$$

The associated Legendre functions can be written most compactly using the substitution $x = \cos \theta$; since $-1 \leq x \leq 1$ we take $0 \leq \theta \leq \pi$ and let $c = \cos \theta$, $s = \sin \theta$:

$$\begin{aligned} P_0(c) &= 1 \\ P_1(c) &= c & P_1^1(c) &= -s \\ P_2(c) &= \frac{1}{2}(3c^2 - 1) & P_2^1(c) &= -3cs & P_2^2(c) &= 3s^2 \\ P_3(c) &= \frac{1}{2}(5c^3 - 3c) & P_3^1(c) &= -\frac{3}{2}s(5c^2 - 1) & P_3^2(c) &= 15cs^2 & P_3^3(c) &= -15s^3. \end{aligned} \quad (C.41)$$

C.6 Spherical harmonics

A spherical harmonic is defined by the expression

$$Y_l^m(\theta, \phi) = \sqrt{\frac{2l+1}{4\pi} \frac{(l-m)!}{(l+m)!}} P_l^m(\cos \theta) e^{im\phi} \quad (m \geq 0). \quad (C.42)$$

For most purposes, the indices of the spherical harmonics can be restricted to $l = 0, 1, 2, \dots$ and $m = -l, -l + 1, \dots, l - 1, l$. The variables lie in the range $0 \leq \theta \leq \pi$ and $0 \leq \phi \leq 2\pi$ and usually represent the angular coordinates in a spherical coordinate system (see Figure B.1). Note that

$$Y_l^{-m}(\theta, \phi) = (-1)^m Y_l^{m*}(\theta, \phi), \quad (C.43)$$

where the asterisk denotes complex conjugation.

The most important feature of the spherical harmonics, which is easily proved using equation (C.39), is that they are orthonormal in the sense that

$$\oint d^2\Omega Y_k^{n*}(\Omega) Y_l^m(\Omega) = \int_0^\pi d\theta \sin \theta \int_0^{2\pi} d\phi Y_k^{n*}(\theta, \phi) Y_l^m(\theta, \phi) = \delta_{kl} \delta_{nm}. \quad (C.44)$$

An arbitrary function of position $f(\mathbf{r})$ can be written in spherical coordinates as a series of spherical harmonics,

$$f(\mathbf{r}) = f(r, \theta, \phi) = \sum_{n=0}^{\infty} \sum_{m=-l}^l f_{lm}(r) Y_l^m(\theta, \phi). \quad (C.45)$$

Multiplying by $Y_n^{k*}(\theta, \phi)$, integrating over solid angle, and using equation (C.44), we find

$$f_{nk}(r) = \int d^2\Omega Y_k^{n*}(\theta, \phi) f(\mathbf{r}). \quad (\text{C.46})$$

The addition theorem for spherical harmonics states that if the directions (θ, ϕ) and (θ', ϕ') are separated by an angle γ , then

$$P_l(\cos \gamma) = \frac{4\pi}{2l+1} \sum_{m=-l}^l Y_l^{m*}(\theta', \phi') Y_l^m(\theta, \phi). \quad (\text{C.47})$$

Together with equation (C.35), this leads to an expression for the inverse distance between the points $\mathbf{x} = (r, \theta, \phi)$ and $\mathbf{x}' = (r', \theta', \phi')$:

$$\frac{1}{|\mathbf{x} - \mathbf{x}'|} = \sum_{l=0}^{\infty} \sum_{m=-l}^l \frac{4\pi}{2l+1} \frac{r_{<}^l}{r_{>}^{l+1}} Y_l^{m*}(\theta', \phi') Y_l^m(\theta, \phi), \quad (\text{C.48})$$

where $r_{<} = \min(r, r')$ and $r_{>} = \max(r, r')$.

Using equations (B.53) and (C.30) we can show that

$$\nabla^2 [f(r) Y_l^m(\theta, \phi)] = \left[\frac{1}{r^2} \frac{d}{dr} \left(r^2 \frac{df}{dr} \right) - \frac{l(l+1)}{r^2} f(r) \right] Y_l^m(\theta, \phi). \quad (\text{C.49})$$

The first few spherical harmonics are:

$$\begin{aligned} Y_0^0(\theta, \phi) &= \frac{1}{\sqrt{4\pi}} \\ Y_1^0(\theta, \phi) &= \sqrt{\frac{3}{4\pi}} \cos \theta & Y_1^{\pm 1}(\theta, \phi) &= \mp \sqrt{\frac{3}{8\pi}} \sin \theta e^{\pm i\phi} \\ Y_2^0(\theta, \phi) &= \sqrt{\frac{5}{16\pi}} (3 \cos^2 \theta - 1) & Y_2^{\pm 1}(\theta, \phi) &= \mp \sqrt{\frac{15}{8\pi}} \sin \theta \cos \theta e^{\pm i\phi} \\ & & Y_2^{\pm 2}(\theta, \phi) &= \sqrt{\frac{15}{32\pi}} \sin^2 \theta e^{\pm 2i\phi}. \end{aligned} \quad (\text{C.50})$$

C.7 Bessel functions

The most complete reference is Watson (1995).

The Bessel functions of the first and second kind, $J_\nu(z)$ and $Y_\nu(z)$, are linearly independent solutions of the differential equation

$$\frac{1}{z} \frac{d}{dz} \left(z \frac{dw}{dz} \right) + \left(1 - \frac{\nu^2}{z^2} \right) w = 0. \quad (\text{C.51})$$

In series form,

$$J_\nu(z) = \sum_{k=0}^{\infty} \frac{(-1)^k}{k! (\nu+k)!} \left(\frac{1}{2} z \right)^{\nu+2k}, \quad (\text{C.52})$$

and $Y_\nu(z)$ is defined by the relation

$$Y_\nu(z) = \frac{\cos \nu\pi J_\nu(z) - J_{-\nu}(z)}{\sin \nu\pi}, \quad (\text{C.53})$$

or by its limiting value if ν is an integer. The function $Y_\nu(z)$ diverges as $z^{-|\nu|}$ when $z \rightarrow 0$. As $x \rightarrow \infty$

$$J_\nu(x) \rightarrow \sqrt{\frac{2}{\pi x}} \cos \left(x - \frac{1}{2} \nu\pi - \frac{1}{4} \pi \right) + O(x^{-1}). \quad (\text{C.54})$$

If $\nu \equiv n$ is an integer

$$J_n(-z) = (-1)^n J_n(z) \quad ; \quad J_{-n}(z) = (-1)^n J_n(z) \quad ; \quad Y_{-n}(z) = (-1)^n Y_n(z), \quad (\text{C.55})$$

$$J_n(z) = \frac{1}{\pi} \int_0^\pi d\theta \cos(z \sin \theta - n\theta). \quad (\text{C.56})$$

If C_ν denotes either J_ν or Y_ν ,

$$C_{\nu-1}(z) + C_{\nu+1}(z) = \frac{2\nu}{z} C_\nu(z) \quad ; \quad C_{\nu-1}(z) - C_{\nu+1}(z) = 2 \frac{dC_\nu(z)}{dz}. \quad (\text{C.57})$$

For $\nu = 0$ these relations imply

$$J'_0(z) = -J_1(z). \quad (\text{C.58})$$

An important integral identity is

$$\int_0^\infty dk k \int_0^\infty dR R F(R) J_\nu(kR) J_\nu(kr) = F(r) \quad (\nu \geq -\frac{1}{2}), \quad (\text{C.59})$$

where $F(R)$ is an arbitrary function. If

$$g(k) = \int_0^\infty dr r f(r) J_\nu(kr) \quad (\text{C.60a})$$

then g is called the **Hankel transform** of f , and equation (C.59) yields

$$f(r) = \int_0^\infty dk k g(k) J_\nu(kr). \quad (\text{C.60b})$$

The Hankel transform is defined for any integer ν and all real $\nu \geq -\frac{1}{2}$.

The **modified Bessel functions** are

$$I_\nu(z) = i^{-\nu} J_\nu(iz) \quad ; \quad K_\nu(z) = K_{-\nu}(z) = \frac{\pi}{2} \frac{I_{-\nu}(z) - I_\nu(z)}{\sin \nu\pi}; \quad (\text{C.61})$$

the second equation is replaced by its limiting value if ν is an integer. As $z \rightarrow 0$,

$$\begin{aligned} I_\nu(z) &\rightarrow \frac{1}{\nu!} \left(\frac{1}{2}z\right)^\nu \quad (\nu \neq -1, -2, \dots); \\ K_\nu(z) &\rightarrow \frac{(\nu-1)!}{2} \left(\frac{1}{2}z\right)^{-\nu} \quad (\nu > 0); \end{aligned} \quad (\text{C.62})$$

At large x ,

$$I_\nu(x) \rightarrow \frac{e^x}{\sqrt{2\pi x}} \quad ; \quad K_\nu(x) \rightarrow \sqrt{\frac{\pi}{2x}} e^{-x}. \quad (\text{C.63})$$

If $\nu \equiv n$ is an integer,

$$I_n(-z) = (-1)^n I_n(z) \quad ; \quad I_n(z) = I_{-n}(z) = \frac{1}{\pi} \int_0^\pi d\theta e^{z \cos \theta} \cos(n\theta). \quad (\text{C.64})$$

If Z_ν denotes either I_ν or $e^{i\pi\nu} K_\nu$,

$$Z_{\nu-1}(z) - Z_{\nu+1}(z) = \frac{2\nu}{z} Z_\nu(z) \quad ; \quad Z_{\nu-1}(z) + Z_{\nu+1}(z) = 2 \frac{dZ_\nu(z)}{dz}; \quad (\text{C.65})$$

for $\nu = 0$ these imply

$$I'_0(z) = I_1(z) \quad ; \quad K'_0(z) = -K_1(z). \quad (\text{C.66})$$

We shall use the results

$$e^{z \cos \theta} = I_0(z) + 2 \sum_{n=1}^{\infty} I_n(z) \cos n\theta = \sum_{n=-\infty}^{\infty} I_n(z) \cos n\theta; \quad (\text{C.67})$$

$$\int_0^{\infty} dx x^{\mu} J_{\nu}(x) = 2^{\mu} \frac{[\frac{1}{2}(\nu + \mu + 3)]!}{[\frac{1}{2}(\nu - \mu + 3)]!} \quad \text{Re}(\nu + \mu) > -1, \text{Re}(\mu) < \frac{1}{2}. \quad (\text{C.68})$$

$$\begin{aligned} \int_u^{\infty} dx \frac{x e^{-\mu x}}{(x^2 - u^2)^{1-\nu}} \\ = \frac{2^{\nu-1/2}(\nu-1)!}{\sqrt{\pi}} \frac{u^{\nu+1/2}}{\mu^{\nu-1/2}} K_{\nu+1/2}(u\mu) \quad (u > 0, \text{Re } \mu, \nu > 0); \end{aligned} \quad (\text{C.69})$$

$$\begin{aligned} \int_0^b dy y^{\nu} (b^2 - y^2)^{\nu-3/2} K_{\nu}(y) = 2^{\nu-3} \sqrt{\pi} (\nu - \frac{3}{2})! b^{2\nu-1} \\ \times [I_{\nu-1}(\frac{1}{2}b) K_{\nu}(\frac{1}{2}b) - I_{\nu}(\frac{1}{2}b) K_{\nu-1}(\frac{1}{2}b)] \quad (\text{Re } \nu > -\frac{1}{2}). \end{aligned} \quad (\text{C.70})$$

Appendix D: Mechanics

We assume a background in classical mechanics at the advanced undergraduate level, including basic Hamiltonian mechanics. Useful reference texts include Landau & Lifshitz (1989), José & Saletan (1998), and Sussman & Wisdom (2001). The most elegant and mathematical treatment of the subject is found in Arnold (1989). This appendix contains a brief summary of the concepts employed in this book.

D.1 Single particles

The momentum of a particle is $\mathbf{p} = m\mathbf{v}$, where m is its mass and \mathbf{v} is its velocity. Its motion is described by **Newton's second law**,

$$\mathbf{F} = \frac{d\mathbf{p}}{dt}, \quad (\text{D.1})$$

where \mathbf{F} is the force acting on the particle. Thus, if the mass of the particle is constant,

$$\frac{d\mathbf{v}}{dt} = \frac{d^2\mathbf{x}}{dt^2} = \frac{\mathbf{F}}{m}. \quad (\text{D.2})$$

The **work** done against the force \mathbf{F} in moving a particle from \mathbf{x}_1 to \mathbf{x}_2 is

$$W_{12} = - \int_{\mathbf{x}_1}^{\mathbf{x}_2} d\mathbf{x} \cdot \mathbf{F}, \quad (\text{D.3})$$

a line integral that is to be taken along the particle's trajectory from \mathbf{x}_1 to \mathbf{x}_2 . The rate at which work is done on the force is

$$\frac{dW}{dt} = - \frac{d}{dt} \int_{\mathbf{x}_1}^{\mathbf{x}_2(t)} d\mathbf{x} \cdot \mathbf{F} = - \frac{d\mathbf{x}}{dt} \cdot \mathbf{F} = -\mathbf{F} \cdot \mathbf{v}, \quad (\text{D.4})$$

evaluated at \mathbf{x}_2 . For a particle of fixed mass,

$$\begin{aligned} W_{12} = -m \int_{\mathbf{x}_1}^{\mathbf{x}_2} d\mathbf{x} \cdot \frac{d^2\mathbf{x}}{dt^2} = -m \int_{\mathbf{x}_1}^{\mathbf{x}_2} dt \frac{d\mathbf{x}}{dt} \cdot \frac{d^2\mathbf{x}}{dt^2} = -m \int_{\mathbf{x}_1}^{\mathbf{x}_2} dt \mathbf{v} \cdot \frac{d\mathbf{v}}{dt} \\ = \frac{1}{2}m[v^2(\mathbf{x}_1) - v^2(\mathbf{x}_2)]. \end{aligned} \quad (\text{D.5})$$

The **kinetic energy** of a particle is $K \equiv \frac{1}{2}mv^2$, so

$$W_{12} = K_1 - K_2. \quad (\text{D.6})$$

Many forces encountered in nature are **conservative**, that is, the work W_{12} is independent of the path taken between the endpoints \mathbf{x}_1 and \mathbf{x}_2 . In this case, we may choose some fixed point \mathbf{x}_0 and define the potential energy $V(\mathbf{x})$ by

$$V(\mathbf{x}) \equiv - \int_{\mathbf{x}_0}^{\mathbf{x}} d\mathbf{x} \cdot \mathbf{F} \quad \text{then} \quad W_{12} = V(\mathbf{x}_2) - V(\mathbf{x}_1), \quad (\text{D.7})$$

A common convention that fixes the otherwise arbitrary zero point of V is to place the point \mathbf{x}_0 at “infinity,” that is, far from all interacting bodies, where \mathbf{F} is negligibly small. All such points yield the same zero point, since the work done in moving from one point at “infinity” to another is negligible.

There are two immediate consequences of equations (D.7). Taking the gradient of V , we obtain

$$\mathbf{F} = -\nabla V(\mathbf{x}) = -\frac{\partial V}{\partial \mathbf{x}}; \quad (\text{D.8})$$

and substituting the second of equations (D.7) into (D.6) yields

$$K_1 + V(\mathbf{x}_1) = K_2 + V(\mathbf{x}_2). \quad (\text{D.9})$$

Thus if the **energy** of a particle is defined to be $E = K + V = \frac{1}{2}mv^2 + V(\mathbf{x})$, we find that *if the forces acting on a particle are conservative, then its energy is conserved*. We shall also encounter cases in which the forces are conservative if the trajectory is traversed instantaneously, but the force at a given position is time-dependent. In these cases we can still write $\mathbf{F} = -\nabla V(\mathbf{x}, t)$, but the energy E is no longer conserved; in fact,

$$\frac{dE}{dt} = \frac{\partial V(\mathbf{x}, t)}{\partial t}. \quad (\text{D.10})$$

Forces due to gravity are conservative (the proof is given at the beginning of Chapter 2); for gravity the potential energy of a particle is proportional to its mass and we may write $V(\mathbf{x}) = m\Phi(\mathbf{x})$, where Φ is the gravitational potential. Thus the energy per unit mass $\frac{1}{2}v^2 + \Phi(\mathbf{x})$ is conserved. Since we deal almost exclusively with gravitational forces in this book, we shall usually shorten the term “energy per unit mass” to simply “energy,” and refer to $\frac{1}{2}v^2$ as the kinetic energy and Φ as the potential energy. When the distinction is not clear from the context, we use the notation E for energy per unit mass and $\tilde{E} = mE$ for energy.

The **angular momentum** of a particle relative to some origin O is

$$\mathbf{L} \equiv \mathbf{x} \times \mathbf{p}, \quad (\text{D.11})$$

where the position vector \mathbf{x} is measured from O . The torque is

$$\mathbf{N} = \mathbf{x} \times \mathbf{F}. \quad (\text{D.12})$$

We have

$$\frac{d\mathbf{L}}{dt} = \frac{d\mathbf{x}}{dt} \times \mathbf{p} + \mathbf{x} \times \frac{d\mathbf{p}}{dt} = \mathbf{v} \times \mathbf{p} + \mathbf{x} \times \mathbf{F}. \quad (\text{D.13})$$

The first term is proportional to $\mathbf{p} \times \mathbf{p} = 0$, and thus

$$\mathbf{N} = \frac{d\mathbf{L}}{dt}; \quad (\text{D.14})$$

in words, *the torque is equal to the rate of change of angular momentum*.

D.2 Systems of particles

Consider an isolated system of N particles of masses m_α and positions \mathbf{x}_α , $\alpha = 1, \dots, N$. The total mass of the system and the total force on particle α are

$$M = \sum_{\alpha=1}^N m_\alpha \quad ; \quad \mathbf{F}_\alpha = \sum_{\substack{\beta=1 \\ \beta \neq \alpha}}^N \mathbf{F}_{\alpha\beta}, \quad (\text{D.15})$$

where $\mathbf{F}_{\alpha\beta}$ is the force exerted on particle α by particle β . According to **Newton's third law**

$$\mathbf{F}_{\alpha\beta} = -\mathbf{F}_{\beta\alpha}. \quad (\text{D.16})$$

The **center of mass** is located at

$$\mathbf{x}_{\text{cm}} = \frac{\sum_{\alpha=1}^N m_\alpha \mathbf{x}_\alpha}{M}. \quad (\text{D.17})$$

Thus

$$\frac{d^2 \mathbf{x}_{\text{cm}}}{dt^2} = \frac{1}{M} \sum_{\alpha=1}^N m_\alpha \frac{d^2 \mathbf{x}_\alpha}{dt^2} = \frac{1}{M} \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \mathbf{F}_{\alpha\beta}. \quad (\text{D.18})$$

The sum over $\mathbf{F}_{\alpha\beta}$ vanishes since each pair $\mathbf{F}_{\alpha\beta} + \mathbf{F}_{\beta\alpha}$ sums to zero. Thus

$$\frac{d^2 \mathbf{x}_{\text{cm}}}{dt^2} = 0, \quad (\text{D.19})$$

and we conclude that *the center of mass of an isolated system moves at uniform velocity.*

Similarly, the total angular momentum is

$$\mathbf{L} = \sum_{\alpha=1}^N \mathbf{x}_\alpha \times \mathbf{p}_\alpha, \quad (\text{D.20})$$

and

$$\frac{d\mathbf{L}}{dt} = \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \mathbf{x}_\alpha \times \mathbf{F}_{\alpha\beta}. \quad (\text{D.21})$$

The right side is a sum of pairs of the form $\mathbf{x}_\alpha \times \mathbf{F}_{\alpha\beta} + \mathbf{x}_\beta \times \mathbf{F}_{\beta\alpha} = (\mathbf{x}_\alpha - \mathbf{x}_\beta) \times \mathbf{F}_{\alpha\beta}$. In the case of gravity and most other phenomena, the interparticle force acts along the line joining the two particles, so this term also vanishes. In words, *the total angular momentum of an isolated system is conserved if the interparticle forces act along the line joining the particles.*

As for a single particle, we compute the work done against the forces \mathbf{F}_α on particle α in moving the system from configuration 1 to configuration 2,

$$W_{12} = - \sum_{\alpha=1}^N \int_1^2 d\mathbf{x}_\alpha \cdot \mathbf{F}_\alpha, \quad (\text{D.22})$$

and as in equation (D.5) we may write

$$W_{12} = K_{\text{tot},1} - K_{\text{tot},2}, \quad (\text{D.23})$$

where $K_{\text{tot}} = \frac{1}{2} \sum_{\alpha=1}^N m_\alpha v_\alpha^2$ is the total kinetic energy of the system. Note that

$$K_{\text{tot}} = \frac{1}{2} M v_{\text{cm}}^2 + \frac{1}{2} \sum_{\alpha=1}^N m_\alpha v_\alpha'^2, \quad (\text{D.24})$$

where $\mathbf{v}'_\alpha = \mathbf{v}_\alpha - \mathbf{v}_{\text{cm}}$ and $\mathbf{v}_{\text{cm}} = \dot{\mathbf{x}}_{\text{cm}}$ is the velocity of the center of mass. Thus the total kinetic energy is the sum of the kinetic energy of motion relative to the center of mass, and the kinetic energy of a single particle of mass M moving at the velocity of the center of mass.

In many cases the interparticle forces are conservative, and can be written as

$$\mathbf{F}_{\alpha\beta} = -\frac{\partial}{\partial \mathbf{x}_\alpha} V(|\mathbf{x}_\alpha - \mathbf{x}_\beta|). \quad (\text{D.25})$$

Note that this form automatically guarantees the validity of Newton's third law, and that the interparticle forces act along the line joining the particles. We have

$$W_{12} = \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} \int_1^2 d\mathbf{x}_\alpha \cdot \frac{\partial}{\partial \mathbf{x}_\alpha} V(|\mathbf{x}_\alpha - \mathbf{x}_\beta|). \quad (\text{D.26})$$

We define the potential energy of the system to be

$$W_{\text{tot}} = \frac{1}{2} \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} V(|\mathbf{x}_\alpha - \mathbf{x}_\beta|). \quad (\text{D.27})$$

Then, as a result of any small change $\mathbf{x}_\alpha \rightarrow \mathbf{x}_\alpha + d\mathbf{x}_\alpha$, $\alpha = 1, \dots, N$,

$$dW_{\text{tot}} = \sum_{\alpha=1}^N \sum_{\beta \neq \alpha} d\mathbf{x}_\alpha \cdot \frac{\partial V(|\mathbf{x}_\alpha - \mathbf{x}_\beta|)}{\partial \mathbf{x}_\alpha}, \quad (\text{D.28})$$

where the factor $\frac{1}{2}$ has disappeared because the term involving a given pair of particles—say, $V(|\mathbf{x}_1 - \mathbf{x}_2|)$ —appears twice in equation (D.27). Thus equations (D.26) and (D.28) yield $W_{12} = W_{\text{tot},2} - W_{\text{tot},1}$, and equation (D.23) yields

$$K_{\text{tot},1} + W_{\text{tot},1} = K_{\text{tot},2} + W_{\text{tot},2}; \quad (\text{D.29})$$

in words, *the total energy of the system $K_{\text{tot}} + W_{\text{tot}}$ is conserved.*

The behavior of an isolated two-body system with conservative interparticle forces is particularly simple. We have

$$m_1 \frac{d^2 \mathbf{x}_1}{dt^2} = \mathbf{F}_{12} = -m_2 \frac{d^2 \mathbf{x}_2}{dt^2}. \quad (\text{D.30})$$

The center of mass is at $\mathbf{x}_{\text{cm}} = (m_1 \mathbf{x}_1 + m_2 \mathbf{x}_2)/M$ and moves at uniform velocity. The relative separation vector $\mathbf{r} = \mathbf{x}_2 - \mathbf{x}_1$ obeys the equation

$$\frac{d^2 \mathbf{r}}{dt^2} = -\frac{\mathbf{F}_{12}}{\mu} = \frac{1}{\mu} \frac{\partial V(|\mathbf{r}|)}{\partial \mathbf{x}_1}, \quad (\text{D.31})$$

where the **reduced mass** is

$$\mu \equiv \frac{m_1 m_2}{m_1 + m_2}. \quad (\text{D.32})$$

Since $\partial V/\partial \mathbf{x}_1 = -\partial V/\partial \mathbf{r}$, we have

$$\frac{d^2 \mathbf{r}}{dt^2} = -\frac{1}{\mu} \frac{\partial V(|\mathbf{r}|)}{\partial \mathbf{r}}, \quad (\text{D.33})$$

which is the equation of motion of a fictitious single particle (the **reduced particle**) of mass μ in the fixed potential $V(|\mathbf{r}|)$. Thus the two-body problem has been reduced to a one-body problem. The total energy is

$$E = \frac{1}{2} M v_{\text{cm}}^2 + \frac{1}{2} (m_1 v_1'^2 + m_2 v_2'^2) + V = \frac{1}{2} M v_{\text{cm}}^2 + \frac{1}{2} \mu \dot{\mathbf{r}}^2 + V(|\mathbf{r}|), \quad (\text{D.34})$$

in which we have used the relations $\mathbf{v}'_1 = -m_2\dot{\mathbf{r}}/M$, $\mathbf{v}'_2 = m_1\dot{\mathbf{r}}/M$.

Rigid bodies Although stellar systems are not rigid bodies, the simple behavior of rigid bodies provides a useful comparison to the behavior of rotating stellar systems.

Consider a rigid body that rotates about its center of mass at $\mathbf{x} = 0$, with angular velocity $\boldsymbol{\Omega}$. Its angular momentum and kinetic energy are given by equations (D.20) and (D.24),

$$\mathbf{L} = \sum_{\alpha=1}^N m_{\alpha} \mathbf{x}_{\alpha} \times \mathbf{v}_{\alpha} \quad ; \quad K = \frac{1}{2} \sum_{\alpha=1}^N m_{\alpha} v_{\alpha}^2, \quad (\text{D.35})$$

where the sum is over all of the particles in the body. Since the body is rigid, the velocity at each point is

$$\mathbf{v}_{\alpha} = \boldsymbol{\Omega} \times \mathbf{x}_{\alpha}. \quad (\text{D.36})$$

Using the vector identity (B.8), we have

$$K = \frac{1}{2} \sum_{\alpha=1}^N m_{\alpha} \mathbf{v}_{\alpha} \cdot (\boldsymbol{\Omega} \times \mathbf{x}_{\alpha}) = \frac{1}{2} \sum_{\alpha=1}^N m_{\alpha} \boldsymbol{\Omega} \cdot (\mathbf{x}_{\alpha} \times \mathbf{v}_{\alpha}) = \frac{1}{2} \boldsymbol{\Omega} \cdot \mathbf{L}. \quad (\text{D.37})$$

Moreover

$$\mathbf{L} = \sum_{\alpha=1}^N m_{\alpha} \mathbf{x}_{\alpha} \times (\boldsymbol{\Omega} \times \mathbf{x}_{\alpha}) = \mathbf{I}' \cdot \boldsymbol{\Omega} \quad \text{or} \quad L_j = \sum_{k=1}^3 I'_{jk} \Omega_k. \quad (\text{D.38})$$

Here \mathbf{I}' is the **moment of inertia tensor**, a symmetric tensor having components¹

$$I'_{jk} \equiv \sum_{\alpha=1}^N m_{\alpha} (|\mathbf{x}_{\alpha}|^2 \delta_{jk} - x_{\alpha j} x_{\alpha k}). \quad (\text{D.41})$$

The kinetic energy is given by

$$K = \frac{1}{2} \boldsymbol{\Omega} \cdot \mathbf{I}' \cdot \boldsymbol{\Omega} = \frac{1}{2} \mathbf{L} \cdot (\mathbf{I}')^{-1} \cdot \mathbf{L}. \quad (\text{D.42})$$

If the coordinate axes coincide with the principal axes of the moment of inertia tensor, then the tensor is diagonal. If, in addition, the angular velocity vector lies in one of the principal axes (say, the z axis), we have

$$\mathbf{L} = L \hat{\mathbf{e}}_z \quad ; \quad \boldsymbol{\Omega} = \Omega \hat{\mathbf{e}}_z \quad ; \quad K = \frac{L^2}{2I}, \quad (\text{D.43})$$

where

$$I \equiv I'_{zz} = \sum_{\alpha=1}^N m_{\alpha} (x_{\alpha}^2 + y_{\alpha}^2). \quad (\text{D.44})$$

¹ Unfortunately, this term is sometimes used to denote the tensor

$$I_{jk} \equiv \sum_{\alpha=1}^N m_{\alpha} x_{\alpha j} x_{\alpha k} \quad (\text{D.39})$$

(cf. eqs. 4.243 and 7.15). The relation between the two definitions is

$$I'_{jk} = \text{trace}(\mathbf{I}) \delta_{jk} - I_{jk}. \quad (\text{D.40})$$

D.3 Lagrangian dynamics

Consider a particle moving in a conservative force field defined by the potential energy $V(\mathbf{x}, t)$. We define the **Lagrangian**

$$\mathcal{L}(\mathbf{x}, \dot{\mathbf{x}}, t) \equiv K - V = \frac{1}{2}m\dot{\mathbf{x}}^2 - V(\mathbf{x}, t). \quad (\text{D.45})$$

The **principle of least action** or **Hamilton's principle** states that *the motion of the particle from time t_0 to t_1 is along a curve $\mathbf{x}(t)$ that is an extremal of the action*

$$I \equiv \int_{t_0}^{t_1} dt \mathcal{L}. \quad (\text{D.46})$$

The proof is simple. According to the Euler–Lagrange equation (B.81), the trajectory is an extremal of I if and only if

$$0 = \frac{d}{dt} \left(\frac{\partial \mathcal{L}}{\partial \dot{\mathbf{x}}} \right) - \frac{\partial \mathcal{L}}{\partial \mathbf{x}} = m\ddot{\mathbf{x}} + \frac{\partial V}{\partial \mathbf{x}}, \quad (\text{D.47})$$

which is simply a restatement of Newton's second law.

The attraction of this approach is that the Lagrangian \mathcal{L} is a *scalar* function. Hence it is straightforward to compute \mathcal{L} as a function $\mathcal{L}(\mathbf{q}, \dot{\mathbf{q}}, t)$ of arbitrary—sometimes called **generalized**—coordinates \mathbf{q} and their time derivatives $\dot{\mathbf{q}}$. Extremizing the action with \mathcal{L} expressed in this form yields **Lagrange's equations**

$$\frac{d}{dt} \left(\frac{\partial \mathcal{L}}{\partial \dot{\mathbf{q}}} \right) - \frac{\partial \mathcal{L}}{\partial \mathbf{q}} = 0, \quad (\text{D.48})$$

which are the equations of motion in the generalized coordinates. This approach avoids the heavy algebra that is often required to express vector equations in curvilinear coordinates.

D.4 Hamiltonian dynamics

For a given set of generalized coordinates \mathbf{q} we define the **generalized momentum** \mathbf{p} to be

$$\mathbf{p} \equiv \left(\frac{\partial \mathcal{L}}{\partial \dot{\mathbf{q}}} \right)_{\mathbf{q}, t}. \quad (\text{D.49})$$

The **Hamiltonian** is

$$H(\mathbf{q}, \mathbf{p}, t) \equiv \mathbf{p} \cdot \dot{\mathbf{q}} - \mathcal{L}(\mathbf{q}, \dot{\mathbf{q}}, t), \quad (\text{D.50})$$

where it is understood that $\dot{\mathbf{q}}$ is to be eliminated in favor of \mathbf{q} , \mathbf{p} , and t using equation (D.49).

D.4.1 Hamilton's equations

The total derivative of the Hamiltonian is

$$dH = \left(\frac{\partial H}{\partial \mathbf{q}} \right)_{\mathbf{p}, t} \cdot d\mathbf{q} + \left(\frac{\partial H}{\partial \mathbf{p}} \right)_{\mathbf{q}, t} \cdot d\mathbf{p} + \left(\frac{\partial H}{\partial t} \right)_{\mathbf{q}, \mathbf{p}} dt. \quad (\text{D.51})$$

Differencing equation (D.50) we also have

$$\begin{aligned} dH &= \mathbf{p} \cdot d\dot{\mathbf{q}} + \dot{\mathbf{q}} \cdot d\mathbf{p} - \left(\frac{\partial \mathcal{L}}{\partial \mathbf{q}} \right)_{\dot{\mathbf{q}}, t} \cdot d\mathbf{q} - \left(\frac{\partial \mathcal{L}}{\partial \dot{\mathbf{q}}} \right)_{\mathbf{q}, t} \cdot d\dot{\mathbf{q}} - \left(\frac{\partial \mathcal{L}}{\partial t} \right)_{\mathbf{q}, \dot{\mathbf{q}}} dt \\ &= \dot{\mathbf{q}} \cdot d\mathbf{p} - \left(\frac{\partial \mathcal{L}}{\partial \mathbf{q}} \right)_{\dot{\mathbf{q}}, t} \cdot d\mathbf{q} - \left(\frac{\partial \mathcal{L}}{\partial t} \right)_{\mathbf{q}, \dot{\mathbf{q}}} dt, \end{aligned} \quad (\text{D.52})$$

where the first and fourth terms in the first line have canceled because of equation (D.49). Since equations (D.51) and (D.52) must be the same, we have

$$\dot{\mathbf{q}} = \left(\frac{\partial H}{\partial \mathbf{p}} \right)_{\mathbf{q}, t}; \quad \left(\frac{\partial H}{\partial \mathbf{q}} \right)_{\mathbf{p}, t} = - \left(\frac{\partial \mathcal{L}}{\partial \mathbf{q}} \right)_{\dot{\mathbf{q}}, t}; \quad \left(\frac{\partial H}{\partial t} \right)_{\mathbf{q}, \mathbf{p}} = - \left(\frac{\partial \mathcal{L}}{\partial t} \right)_{\mathbf{q}, \dot{\mathbf{q}}}. \quad (\text{D.53})$$

Using Lagrange's equations (D.48) and simplifying the notation, the first two of these equations lead us to **Hamilton's equations**

$$\dot{\mathbf{q}} = \frac{\partial H}{\partial \mathbf{p}}; \quad \dot{\mathbf{p}} = - \frac{\partial H}{\partial \mathbf{q}}. \quad (\text{D.54})$$

The **configuration space** of a system is the n -dimensional space with coordinates (q_1, \dots, q_n) . The corresponding **momentum space** has coordinates (p_1, \dots, p_n) . A system with n -dimensional configuration and momentum spaces is said to have n **degrees of freedom**. Phase space is the $2n$ -dimensional space with coordinates $(q_1, \dots, q_n, p_1, \dots, p_n) \equiv (\mathbf{q}, \mathbf{p}) \equiv \mathbf{w}$. Since Hamilton's equations (D.54) are first-order differential equations, if we are given the Hamiltonian and a particle's phase-space coordinates \mathbf{w}_0 at time $t = 0$, we can solve for the coordinates \mathbf{w}_t at any later (or earlier) time t . Thus through each point \mathbf{w}_0 in phase space there passes a unique phase-space trajectory $\mathbf{w}_t(\mathbf{w}_0)$, which gives the future and past phase-space coordinates of the particle that has coordinates \mathbf{w}_0 at $t = 0$. No two of these trajectories can ever intersect: if they did, the trajectory of a particle that started from the intersection point would not be unique. We define the **time-evolution operator** \mathbf{H}_t by

$$\mathbf{w}_t \equiv \mathbf{H}_t(\mathbf{w}_0), \quad (\text{D.55})$$

and say that the operator \mathbf{H}_t defines a **Hamiltonian flow** in phase space.

Along a trajectory $\mathbf{w}(t)$, the Hamiltonian $H[\mathbf{w}(t), t]$ changes at a rate

$$\frac{dH}{dt} = \frac{\partial H}{\partial \mathbf{q}} \cdot \dot{\mathbf{q}} + \frac{\partial H}{\partial \mathbf{p}} \cdot \dot{\mathbf{p}} + \frac{\partial H}{\partial t} = \frac{\partial H}{\partial t}, \quad (\text{D.56})$$

where we have used (D.54) to eliminate $\dot{\mathbf{q}}$ and $\dot{\mathbf{p}}$. Hence, if $\partial \mathcal{L} / \partial t = 0$, it follows from (D.53) that the Hamiltonian is conserved along all dynamical trajectories.

Thus, for example, consider motion in the time-independent potential $V(\mathbf{x})$. If we work in Cartesian coordinates, the Lagrangian $\mathcal{L} = \frac{1}{2}m\dot{\mathbf{x}}^2 - V(\mathbf{x})$ depends only on \mathbf{x} and $\dot{\mathbf{x}}$, so $\partial \mathcal{L} / \partial t = 0$. Hence the Hamiltonian H is conserved. The physical meaning of this conservation law is easily found. We have $\mathbf{p} = \partial \mathcal{L} / \partial \dot{\mathbf{x}} = m\dot{\mathbf{x}}$ and

$$H(\mathbf{x}, \mathbf{p}) = \mathbf{p} \cdot \dot{\mathbf{x}} - \mathcal{L} = \frac{p^2}{2m} + V(\mathbf{x}); \quad (\text{D.57})$$

this is simply the total energy $E = K + V$, which we know to be conserved for motion in a time-independent potential. Thus *for motion in a fixed potential the Hamiltonian is equal to the total energy*.

In most cases, the velocity \mathbf{v} of a particle of mass m is just $1/m$ times the particle's momentum \mathbf{p} , so by a slight abuse of language, we shall also use the term phase space to denote the six-dimensional space with coordinates (\mathbf{x}, \mathbf{v}) .

D.4.2 Poincaré invariants

Let \mathcal{S}_0 be any two-dimensional surface in phase space, and let us specify points on \mathcal{S}_0 by some set of coordinates (u, v) . The time-evolution operator \mathbf{H}_t maps each point of \mathcal{S}_0 into a new surface \mathcal{S}_t and also maps the curves of constant u, v on \mathcal{S}_0 into curves that define (u, v) coordinates on \mathcal{S}_t . With this definition, each pair (u, v) defines a trajectory.

We define

$$\begin{aligned}
 A(t) &\equiv \iint_{\mathcal{S}_t} d\mathbf{q} \cdot d\mathbf{p} = \sum_{i=1}^n \iint_{\mathcal{S}_t} dq_i dp_i \\
 &= \sum_{i=1}^n \iint_{\mathcal{S}_t} dudv \frac{\partial(q_i, p_i)}{\partial(u, v)},
 \end{aligned}
 \tag{D.58}$$

and calculate dA/dt . We set $t' \equiv t + \delta t$, where δt is small, and let $q(u, v), q'(u, v)$ be points on the same trajectory at times t and t' , with similar definitions for p and p' . To first order in the small interval δt , Hamilton's equations (D.54) yield

$$(\mathbf{q}', \mathbf{p}') = \mathbf{H}_{\delta t}(\mathbf{q}, \mathbf{p}) = \left(\mathbf{q} + \frac{\partial H}{\partial \mathbf{p}} \delta t, \mathbf{p} - \frac{\partial H}{\partial \mathbf{q}} \delta t \right).
 \tag{D.59}$$

Differentiating these equations with respect to u and v , we find

$$\begin{aligned}
 \frac{\partial(q'_i, p'_i)}{\partial(u, v)} &= \begin{vmatrix} \frac{\partial q_i}{\partial u} + \frac{\partial^2 H}{\partial u \partial p_i} \delta t & \frac{\partial q_i}{\partial v} + \frac{\partial^2 H}{\partial v \partial p_i} \delta t \\ \frac{\partial p_i}{\partial u} - \frac{\partial^2 H}{\partial u \partial q_i} \delta t & \frac{\partial p_i}{\partial v} - \frac{\partial^2 H}{\partial v \partial q_i} \delta t \end{vmatrix} \\
 &= \frac{\partial(q_i, p_i)}{\partial(u, v)} + \left(\frac{\partial q_i}{\partial v} \frac{\partial^2 H}{\partial u \partial q_i} - \frac{\partial p_i}{\partial u} \frac{\partial^2 H}{\partial v \partial p_i} + \frac{\partial p_i}{\partial v} \frac{\partial^2 H}{\partial u \partial p_i} - \frac{\partial q_i}{\partial u} \frac{\partial^2 H}{\partial v \partial q_i} \right) \delta t + O(\delta t)^2.
 \end{aligned}
 \tag{D.60}$$

Thus

$$\begin{aligned}
 \frac{dA}{dt} &= \lim_{\delta t \rightarrow 0} \frac{1}{\delta t} \iint dudv \sum_i \left[\frac{\partial(q'_i, p'_i)}{\partial(u, v)} - \frac{\partial(q_i, p_i)}{\partial(u, v)} \right] \\
 &= \sum_i \iint dudv \left(\frac{\partial q_i}{\partial v} \frac{\partial^2 H}{\partial u \partial q_i} - \frac{\partial p_i}{\partial u} \frac{\partial^2 H}{\partial v \partial p_i} + \frac{\partial p_i}{\partial v} \frac{\partial^2 H}{\partial u \partial p_i} - \frac{\partial q_i}{\partial u} \frac{\partial^2 H}{\partial v \partial q_i} \right).
 \end{aligned}
 \tag{D.61}$$

One may show that the sum of the brackets in the second of equations (D.61) vanishes, by replacing every occurrence of $\partial/\partial u$ in the second derivatives by $\sum_k \left(\frac{\partial q_k}{\partial u} \frac{\partial}{\partial q_k} + \frac{\partial p_k}{\partial u} \frac{\partial}{\partial p_k} \right)$ and similarly for $\partial/\partial v$. Hence $dA/dt = 0$ and we have:

Poincaré invariant theorem *If $\mathcal{S}(0)$ is any two-surface in phase space, and $\mathcal{S}(t)$ is the surface into which $\mathcal{S}(0)$ is mapped by the time-evolution operator \mathbf{H}_t , then*

$$\iint_{\mathcal{S}(0)} d\mathbf{q} \cdot d\mathbf{p} = \iint_{\mathcal{S}(t)} d\mathbf{q} \cdot d\mathbf{p}.
 \tag{D.62}$$

This conserved quantity is known as a **Poincaré invariant**.

Corollary *If $\gamma(0)$ is any closed path through phase space, and $\gamma(t)$ is the path to which $\gamma(0)$ is mapped by the time-evolution operator, then*

$$\oint_{\gamma(0)} d\mathbf{q} \cdot \mathbf{p} = \oint_{\gamma(t)} d\mathbf{q} \cdot \mathbf{p}.
 \tag{D.63}$$

Proof: By Green's theorem (B.61),

$$\oint_{\gamma(t)} \mathbf{dq} \cdot \mathbf{p} = \sum_i \oint_{\gamma(t)} dq_i p_i = \sum_i \iint_{\mathcal{S}(t)} dq_i dp_i = \iint_{\mathcal{S}(t)} \mathbf{dq} \cdot \mathbf{dp}, \quad (\text{D.64})$$

where $\mathcal{S}(t)$ is any surface that has $\gamma(t)$ as its boundary. The result now follows from the Poincaré invariant theorem. ◁

D.4.3 Poisson brackets

Let $A(\mathbf{w})$ and $B(\mathbf{w})$ be any two scalar functions of the phase-space coordinates. Then the **Poisson bracket** is defined by

$$[A, B] \equiv \frac{\partial A}{\partial \mathbf{q}} \cdot \frac{\partial B}{\partial \mathbf{p}} - \frac{\partial A}{\partial \mathbf{p}} \cdot \frac{\partial B}{\partial \mathbf{q}}. \quad (\text{D.65})$$

An equivalent definition is

$$[A, B] = \sum_{\alpha, \beta=1}^{2n} J_{\alpha\beta} \frac{\partial A}{\partial w_\alpha} \frac{\partial B}{\partial w_\beta}, \quad (\text{D.66})$$

where the **symplectic matrix** is

$$\mathbf{J} \equiv \begin{pmatrix} \mathbf{0} & \mathbf{I} \\ -\mathbf{I} & \mathbf{0} \end{pmatrix}, \quad (\text{D.67})$$

and $\mathbf{0}$ and \mathbf{I} are the $n \times n$ zero and unit matrix. The symplectic matrix has the useful properties

$$\mathbf{J}^{-1} = \mathbf{J}^T = -\mathbf{J} \quad ; \quad \mathbf{J}^2 = -\mathbf{I} \quad ; \quad |\mathbf{J}| = 1, \quad (\text{D.68})$$

where $|\mathbf{J}|$ denotes the determinant of \mathbf{J} , and the superscript “T” denotes the transpose of a matrix, $A_{ij}^T = A_{ji}$.

It is straightforward to verify the following properties of Poisson brackets:

- (i) $[A, B] = -[B, A]$ and $[A + B, C] = [A, C] + [B, C]$;
- (ii) $[AB, C] = A[B, C] + B[A, C]$;
- (iii) $[[A, B], C] + [[B, C], A] + [[C, A], B] = 0$ (**Jacobi identity**);
- (iv) Hamilton's equations may be written

$$\dot{w}_\alpha = [w_\alpha, H] \quad \text{or} \quad \dot{\mathbf{w}} = [\mathbf{w}, H] \quad \text{or} \quad \dot{\mathbf{w}} = \mathbf{J} \cdot \frac{\partial H}{\partial \mathbf{w}}; \quad (\text{D.69})$$

- (v) The phase-space coordinates $\mathbf{w} = (\mathbf{q}, \mathbf{p})$ satisfy the **canonical commutation relations** $[p_i, p_j] = [q_i, q_j] = 0$ and $[q_i, p_j] = \delta_{ij}$, or simply

$$[w_\alpha, w_\beta] = J_{\alpha\beta}. \quad (\text{D.70})$$

D.4.4 Canonical coordinates and transformations

The power of Lagrange's equations (D.48) is that they can describe the motion of a dynamical system in *any* system of coordinates \mathbf{q} . It is therefore natural to ask what are the most general phase-space coordinates in which motion can be described by Hamilton's equations.

Any set of phase-space coordinates $\mathbf{W} \equiv \{W_\alpha, \alpha = 1, \dots, 2n\}$ is said to be **canonical** if

$$[W_\alpha, W_\beta] = J_{\alpha\beta}. \quad (\text{D.71})$$

Equation (D.70) shows that ordinary phase-space coordinates are canonical. If \mathbf{W} and \mathbf{w} are two sets of canonical coordinates, then the function relating them, $\mathbf{W}(\mathbf{w})$, is a **canonical transformation** or **map**.

Let \mathbf{W} be a set of canonical coordinates; then with equation (D.70) and the chain rule we have

$$\begin{aligned} [A, B] &= \sum_{\alpha, \beta=1}^{2n} J_{\alpha\beta} \frac{\partial A}{\partial w_\alpha} \frac{\partial B}{\partial w_\beta} = \sum_{\kappa\lambda} \left(\sum_{\alpha\beta} J_{\alpha\beta} \frac{\partial W_\kappa}{\partial w_\alpha} \frac{\partial W_\lambda}{\partial w_\beta} \right) \frac{\partial A}{\partial W_\kappa} \frac{\partial B}{\partial W_\lambda} \\ &= \sum_{\kappa\lambda} [W_\kappa, W_\lambda] \frac{\partial A}{\partial W_\kappa} \frac{\partial B}{\partial W_\lambda} = \sum_{\kappa\lambda} J_{\kappa\lambda} \frac{\partial A}{\partial W_\kappa} \frac{\partial B}{\partial W_\lambda}. \end{aligned} \quad (\text{D.72})$$

Thus the derivatives involved in the definition (D.65) of the Poisson bracket can be taken with respect to *any* set of canonical coordinates, just as the vector formula $\nabla \cdot \mathbf{a} = \sum_i (\partial a_i / \partial x_i)$ is valid in any Cartesian coordinate system.

The rate of change of an arbitrary canonical coordinate W_α along an orbit is

$$\dot{W}_\alpha = \sum_{\beta=1}^{2n} \frac{\partial W_\alpha}{\partial w_\beta} \dot{w}_\beta = \sum_{\beta\gamma} \frac{\partial W_\alpha}{\partial w_\beta} J_{\beta\gamma} \frac{\partial H}{\partial w_\gamma}, \quad (\text{D.73})$$

where we have used Hamilton's equations (D.69). We now write H in terms of the new coordinates W_α , so

$$\dot{W}_\alpha = \sum_{\beta\gamma\delta} \frac{\partial W_\alpha}{\partial w_\beta} J_{\beta\gamma} \frac{\partial W_\delta}{\partial w_\gamma} \frac{\partial H}{\partial W_\delta} = \sum_{\delta} [W_\alpha, W_\delta] \frac{\partial H}{\partial W_\delta}. \quad (\text{D.74})$$

Since the coordinates are canonical, this simplifies to

$$\dot{W}_\alpha = \sum_{\delta} J_{\alpha\delta} \frac{\partial H}{\partial W_\delta}, \quad (\text{D.75})$$

which is simply Hamilton's equations (D.69) in the coordinate system \mathbf{W} . Thus *Hamilton's equations are valid in any canonical coordinate system*, whatever the Hamiltonian may be.

The Jacobian matrix relating two coordinate systems \mathbf{w} and \mathbf{W} is a $2n \times 2n$ matrix $\mathbf{g}(\mathbf{w})$ defined by

$$g_{\alpha\beta} \equiv \frac{\partial W_\alpha}{\partial w_\beta}. \quad (\text{D.76})$$

In terms of this matrix, equation (D.66) can be written

$$[W_\alpha, W_\beta] = \sum_{\gamma\delta} J_{\gamma\delta} g_{\alpha\gamma} g_{\beta\delta}. \quad (\text{D.77})$$

If the coordinates \mathbf{W} are canonical, $[W_\alpha, W_\beta] = J_{\alpha\beta}$, then this equation can be rewritten in matrix notation as

$$\mathbf{g} \cdot \mathbf{J} \cdot \mathbf{g}^T = \mathbf{J}. \quad (\text{D.78})$$

Conversely, if equation (D.78) holds the coordinates \mathbf{W} are canonical.

An equivalent condition is obtained by left-multiplying this result by \mathbf{J} and using equation (D.68) to obtain $\mathbf{J} \cdot \mathbf{g} \cdot \mathbf{J} \cdot \mathbf{g}^T = -\mathbf{I}$; then left-multiplying by \mathbf{g}^T to obtain $\mathbf{g}^T \cdot \mathbf{J} \cdot \mathbf{g} \cdot \mathbf{J} \cdot \mathbf{g}^T = -\mathbf{g}^T$; then right-multiplying by the inverse of \mathbf{g}^T to obtain $\mathbf{g}^T \cdot \mathbf{J} \cdot \mathbf{g} \cdot \mathbf{J} = -\mathbf{I}$; then right-multiplying by \mathbf{J} to obtain

$$\mathbf{g}^T \cdot \mathbf{J} \cdot \mathbf{g} = \mathbf{J}. \quad (\text{D.79})$$

A matrix \mathbf{g} that satisfies the condition (D.78) or (D.79) is said to be **symplectic**. Thus a transformation $\mathbf{W}(\mathbf{w})$ is canonical if and only if its Jacobian matrix is symplectic. Sometimes the terms “canonical” and “symplectic” are used interchangeably.

Consider a small phase-space volume element, which may be written $d^{2n}\mathbf{w}$ and $d^{2n}\mathbf{W}$ in two different canonical coordinate systems. The relation between these two expressions is

$$d^{2n}\mathbf{W} = \left\| \frac{\partial \mathbf{W}}{\partial \mathbf{w}} \right\| d^{2n}\mathbf{w} = \|\mathbf{g}\| d^{2n}\mathbf{w}, \quad (\text{D.80})$$

where $\|\mathbf{g}\|$ denotes the absolute value of the determinant $|\mathbf{g}|$. From equation (D.78), $|\mathbf{J}| = |\mathbf{g}| \times |\mathbf{J}| \times |\mathbf{g}^T| = |\mathbf{J}| \times |\mathbf{g}|^2$, where we have used the fact that the determinant of a matrix and its transpose are identical. Since $|\mathbf{J}|$ is non-zero (eq. D.68), $|\mathbf{g}| = \pm 1$, so $\|\mathbf{g}\| = 1$ and

$$d^{2n}\mathbf{W} = d^{2n}\mathbf{w}; \quad (\text{D.81})$$

in words, *the phase-space volume element is the same in any canonical coordinates.*

Finally, we consider how the Poincaré invariant is changed by an canonical transformation. In the canonical coordinates \mathbf{W} , the analog of equation (D.58) is

$$\begin{aligned} A'(t) &\equiv \iint_{S_t} d\mathbf{Q} \cdot d\mathbf{P} = \sum_{i=1}^n \iint_{S_t} dudv \frac{\partial(Q_i, P_i)}{\partial(u, v)} \\ &= \sum_{i=1}^n \iint_{S_t} dudv \left(\frac{\partial Q_i}{\partial u} \frac{\partial P_i}{\partial v} - \frac{\partial Q_i}{\partial v} \frac{\partial P_i}{\partial u} \right). \end{aligned} \quad (\text{D.82})$$

The expression in brackets can be rewritten using the symplectic matrix,

$$\begin{aligned} A'(t) &= \sum_{\alpha, \beta=1}^{2n} \iint_{S_t} dudv \frac{\partial W_\alpha}{\partial u} J_{\alpha\beta} \frac{\partial W_\beta}{\partial v} \\ &= \sum_{\alpha, \beta, \gamma, \delta} \iint_{S_t} dudv \frac{\partial W_\alpha}{\partial w_\gamma} \frac{\partial w_\gamma}{\partial u} J_{\alpha\beta} \frac{\partial W_\beta}{\partial w_\delta} \frac{\partial w_\delta}{\partial v} \\ &= \sum_{\gamma, \delta} \iint_{S_t} dudv \frac{\partial w_\gamma}{\partial u} (\mathbf{g}^T \cdot \mathbf{J} \cdot \mathbf{g})_{\gamma\delta} \frac{\partial w_\delta}{\partial v} \\ &= \sum_{\gamma, \delta} \iint_{S_t} dudv \frac{\partial w_\gamma}{\partial u} J_{\gamma\delta} \frac{\partial w_\delta}{\partial v}, \end{aligned} \quad (\text{D.83})$$

where the last line follows from equation (D.79) because \mathbf{W} is canonical. The last line is simply the Poincaré invariant $A(t)$ in the original coordinates \mathbf{w} . We conclude that Poincaré invariants are conserved in a canonical transformation. Thus, if \mathbf{w} and \mathbf{W} are any two sets of canonical coordinates,

$$\iint_S d\mathbf{q} \cdot d\mathbf{p} = \iint_S d\mathbf{Q} \cdot d\mathbf{P} \quad ; \quad \oint_\gamma d\mathbf{q} \cdot \mathbf{p} = \oint_\gamma d\mathbf{Q} \cdot \mathbf{P} \quad (\text{D.84})$$

for all surfaces S and closed curves γ in phase space. The converse argument shows that any transformation that conserves all integrals of the form $\iint d\mathbf{q} \cdot d\mathbf{p}$ and $\oint d\mathbf{q} \cdot \mathbf{p}$ is canonical. Thus a transformation is canonical if and only if it conserves the Poincaré invariants of all surfaces in phase space.

We have shown that the mapping defined by the time-evolution operator \mathbf{H}_t (eq. D.55) conserves Poincaré invariants; hence the transformation defined by a

Hamiltonian flow over any time interval t is canonical. Since canonical transformations preserve the phase-space volume element, the phase-space volume element is conserved in the flow. Since any finite volume is the sum of volume elements, *the volume of any arbitrary region in phase space is conserved by a Hamiltonian flow.* This result leads immediately to the collisionless Boltzmann equation in the form (4.10).

D.4.5 Extended phase space

A Hamiltonian that does not depend explicitly on time is said to be **autonomous**. Although most of the Hamiltonians that we shall deal with in this book are autonomous, it is important to understand the behavior of dynamical systems governed by time-dependent or non-autonomous Hamiltonians $H(\mathbf{q}, \mathbf{p}, t)$. The trajectories in such systems still satisfy Hamilton's equations (D.54), but by (D.56) the energy of the system $E(t) = H[\mathbf{q}(t), \mathbf{p}(t), t]$ is no longer conserved along a trajectory.

Let us define $p_0 \equiv -E$, $q_0 = t$, and an **extended phase space** with coordinates (\mathbf{Q}, \mathbf{P}) , where $\mathbf{Q} \equiv (q_0, q_1, \dots, q_n)$ and $\mathbf{P} \equiv (p_0, p_1, \dots, p_n)$. We introduce a new variable τ called the **fictitious time**, which serves as the time coordinate in the extended phase space. Now consider the trajectories in the extended phase space governed by the Hamiltonian

$$\mathcal{H}(\mathbf{Q}, \mathbf{P}) \equiv H(\mathbf{q}, \mathbf{p}, t) - E = H(\mathbf{q}, \mathbf{p}, q_0) + p_0. \quad (\text{D.85})$$

These are given by Hamilton's equations

$$\frac{dQ_i}{d\tau} = \frac{\partial \mathcal{H}}{\partial P_i} \quad ; \quad \frac{dP_i}{d\tau} = -\frac{\partial \mathcal{H}}{\partial Q_i} \quad i = 0, \dots, n. \quad (\text{D.86})$$

For $i = 0$ the first of these becomes $dt/d\tau = 1$ which implies that the fictitious time and the actual time coincide along a trajectory; and the second becomes $dE/d\tau = \partial H/\partial t$ which is identical to equation (D.56). For $i = 1, \dots, n$ the equations of motion (D.86) are identical to Hamilton's equations in the original phase space. Thus *the behavior of a dynamical system governed by a time-dependent Hamiltonian with n degrees of freedom can be described by an autonomous Hamiltonian with $(n + 1)$ degrees of freedom.*

D.4.6 Generating functions

Consider two points \mathbf{w}_0 and \mathbf{w} in phase space, and two distinct paths γ_0 and γ_1 from \mathbf{w}_0 to \mathbf{w}_1 . Let Γ denote the closed path from \mathbf{w}_0 to \mathbf{w} and back that is created by first traversing γ_0 and then traversing γ_1 in the reverse direction. If (\mathbf{q}, \mathbf{p}) and $(\mathbf{q}', \mathbf{p}')$ are canonical coordinate systems then the conservation of Poincaré invariants in canonical transformations implies that

$$\oint_{\Gamma} (d\mathbf{q} \cdot \mathbf{p} - d\mathbf{q}' \cdot \mathbf{p}') = 0; \quad (\text{D.87})$$

splitting the path Γ into its components we have

$$\int_{\gamma_0} (d\mathbf{q} \cdot \mathbf{p} - d\mathbf{q}' \cdot \mathbf{p}') = \int_{\gamma_1} (d\mathbf{q} \cdot \mathbf{p} - d\mathbf{q}' \cdot \mathbf{p}'). \quad (\text{D.88})$$

We conclude that the integral does not depend on the path of integration, so for a fixed initial point \mathbf{w}_0 we may write

$$S(\mathbf{w}) = \int_{\mathbf{w}_0}^{\mathbf{w}} (d\mathbf{q} \cdot \mathbf{p} - d\mathbf{q}' \cdot \mathbf{p}'), \quad (\text{D.89})$$

or

$$dS = d\mathbf{q} \cdot \mathbf{p} - d\mathbf{q}' \cdot \mathbf{p}', \quad (\text{D.90})$$

where dS is an exact differential. S is called a **generating function** of the canonical transformation from (\mathbf{q}, \mathbf{p}) to $(\mathbf{q}', \mathbf{p}')$.

Now let us assume that we use $(\mathbf{q}, \mathbf{q}')$ as phase-space coordinates instead of (\mathbf{q}, \mathbf{p}) . Then $dS(\mathbf{q}, \mathbf{q}') = d\mathbf{q} \cdot \mathbf{p} - d\mathbf{q}' \cdot \mathbf{p}'$, where \mathbf{p} and \mathbf{p}' are regarded as functions of \mathbf{q} and \mathbf{q}' . Since dS is an exact differential, we must have

$$\mathbf{p} = \frac{\partial S(\mathbf{q}, \mathbf{q}')}{\partial \mathbf{q}} \quad ; \quad \mathbf{p}' = -\frac{\partial S(\mathbf{q}, \mathbf{q}')}{\partial \mathbf{q}'}. \quad (\text{D.91})$$

Every sufficiently smooth and non-degenerate function $S(\mathbf{q}, \mathbf{q}')$ defines a canonical transformation through these relations.

Similarly, let $S_2(\mathbf{q}, \mathbf{p}') \equiv \mathbf{q}' \cdot \mathbf{p}' + S$, where \mathbf{p} and \mathbf{q}' are now regarded as functions of \mathbf{q} and \mathbf{p}' . Then equation (D.90) becomes

$$dS_2(\mathbf{q}, \mathbf{p}') = \mathbf{q}' \cdot d\mathbf{p}' + d\mathbf{q}' \cdot \mathbf{p}' + dS = \mathbf{q}' \cdot d\mathbf{p}' + d\mathbf{q} \cdot \mathbf{p} \quad (\text{D.92})$$

and

$$\mathbf{q}' = \frac{\partial S_2(\mathbf{q}, \mathbf{p}')}{\partial \mathbf{p}'} \quad ; \quad \mathbf{p} = \frac{\partial S_2(\mathbf{q}, \mathbf{p}')}{\partial \mathbf{q}}. \quad (\text{D.93})$$

If $S_3(\mathbf{q}', \mathbf{p}) \equiv \mathbf{q} \cdot \mathbf{p} - S$, where \mathbf{q} and \mathbf{p}' are regarded as functions of \mathbf{q}' and \mathbf{p} , then

$$\mathbf{p}' = \frac{\partial S_3(\mathbf{q}', \mathbf{p})}{\partial \mathbf{q}'} \quad ; \quad \mathbf{q} = \frac{\partial S_3(\mathbf{q}', \mathbf{p})}{\partial \mathbf{p}}. \quad (\text{D.94})$$

Canonical transformations can also be defined by generating functions of the form $S_4(\mathbf{p}, \mathbf{p}')$. Notice that all these generating functions depend on one old and one new variable so the canonical transformation defined by a generating function is, inconveniently, always implicit.

The generating function $S(\mathbf{q}, \mathbf{q}') = \mathbf{q} \cdot \mathbf{q}'$ yields the canonical transformation $(\mathbf{q}', \mathbf{p}') = (\mathbf{p}, -\mathbf{q})$; this transformation simply exchanges position and momentum and highlights the fact that configuration space and momentum space have equal status in Hamiltonian mechanics. The generating function $S_2(\mathbf{q}, \mathbf{p}') = \mathbf{q} \cdot \mathbf{p}'$ yields the identity transformation $(\mathbf{q}', \mathbf{p}') = (\mathbf{q}, \mathbf{p})$.

These results can be extended to time-dependent generating functions $S(\mathbf{w}, t)$, using the extended phase space $(\mathbf{Q}, \mathbf{P}) = (q_0, \mathbf{q}, p_0, \mathbf{p})$ defined earlier. Let us define a canonical transformation in the extended phase space by the generating function

$$S_2(\mathbf{Q}, \mathbf{P}') = q_0 p'_0 + S_2(\mathbf{q}, \mathbf{p}', q_0). \quad (\text{D.95})$$

The analogs to equations (D.93) in the extended phase space become

$$q'_0 = \frac{\partial S_2}{\partial p'_0} = q_0 \quad ; \quad p_0 = \frac{\partial S_2}{\partial q_0} = p'_0 + \frac{\partial S_2}{\partial q_0} \quad ; \quad \mathbf{q}' = \frac{\partial S_2}{\partial \mathbf{p}'} \quad ; \quad \mathbf{p} = \frac{\partial S_2}{\partial \mathbf{q}}. \quad (\text{D.96})$$

The third and fourth equations are simply a restatement of the original canonical transformation between \mathbf{w} and \mathbf{w}' as defined by equations (D.93). The Hamiltonian $\mathcal{H}(\mathbf{Q}, \mathbf{P})$ in the extended phase space (eq. D.85) is transformed to

$$\mathcal{H}'(\mathbf{Q}', \mathbf{P}') = \mathcal{H}(\mathbf{Q}, \mathbf{P}) = H(\mathbf{q}, \mathbf{p}, q_0) + p_0 = H(\mathbf{q}, \mathbf{p}, q_0) + p'_0 + \frac{\partial S_2}{\partial q_0}, \quad (\text{D.97})$$

where q_0 , \mathbf{q} , and \mathbf{p} are functions of q'_0 , \mathbf{q}' , and \mathbf{p}' defined by equations (D.96). Now we transform back from the extended phase space $(\mathbf{Q}', \mathbf{P}')$ to ordinary phase space with coordinates $(\mathbf{q}', \mathbf{p}')$. The time variable is unchanged, since $q_0 = q'_0$. However, the motion is now described by a Hamiltonian

$$H'(\mathbf{q}', \mathbf{p}', t) = H(\mathbf{q}, \mathbf{p}, t) + \frac{\partial S_2}{\partial q_0} = H(\mathbf{q}, \mathbf{p}, t) + \frac{\partial S_2}{\partial t}. \quad (\text{D.98})$$

Similarly, the generating function $S_3(\mathbf{q}', \mathbf{p}, t)$ transforms the Hamiltonian to

$$H'(\mathbf{q}', \mathbf{p}', t) = H(\mathbf{q}, \mathbf{p}, t) - \frac{\partial S_3}{\partial t}. \quad (\text{D.99})$$

Appendix E: Delaunay variables for Kepler orbits

Angle-action variables for Kepler orbits are fundamental tools of celestial mechanics. They follow immediately from formulae in §3.5.2 for the isochrone potential by taking the limit $b \rightarrow 0$. From equations (3.226), the Kepler Hamiltonian and frequencies are

$$H_K = -\frac{(GM)^2}{2(L + J_r)^2} \quad ; \quad \Omega_r = \Omega_\vartheta = \frac{(GM)^2}{(L + J_r)^3} \quad ; \quad \Omega_\phi = \text{sgn}(J_\phi)\Omega_\vartheta. \quad (\text{E.1})$$

Since H_K depends on the actions in the combination $L + J_r$, it is convenient to make this combination an action. **Delaunay variables** (J_a, J_b, J_c) are defined by the generating function

$$S_d = \theta_\phi J_a + \theta_\vartheta (J_b - |J_a|) + \theta_r (J_c - J_b). \quad (\text{E.2})$$

Differentiating with respect to the old angles we discover the connection between the new and old actions:

$$J_\phi = J_a \quad ; \quad J_\vartheta = J_b - |J_a| \quad ; \quad J_r = J_c - J_b. \quad (\text{E.3})$$

Thus the Delaunay actions are $(J_a, J_b, J_c) = (L_z, L, L + J_r)$. In these variables the Kepler Hamiltonian and frequencies take on an extremely simple form:

$$H_K = -\frac{(GM)^2}{2J_c^2} \quad ; \quad \Omega_a = \Omega_b = 0 \quad ; \quad \Omega_c = \frac{(GM)^2}{J_c^3}. \quad (\text{E.4})$$

Since Ω_a and Ω_b both vanish, the Delaunay angles θ_a and θ_b are both integrals of motion. Thus the Kepler potential has five isolating integrals—three actions and two angles. Differentiating S_d with respect to J_a , etc., we find that the Delaunay angles are

$$\theta_a = \theta_\phi - \text{sgn}(J_a)\theta_\vartheta = \theta_1 \quad ; \quad \theta_b = \theta_\vartheta - \theta_r \quad ; \quad \theta_c = \theta_r. \quad (\text{E.5})$$

Physically, the constancy of θ_a and θ_b implies that both the orientation of the orbital plane and the direction of the line of apsides are fixed, or, equivalently, that both the angular-momentum vector \mathbf{L} and the eccentricity vector \mathbf{e} (Box 3.2) are time-independent.¹

A Kepler orbit is conventionally described by its semi-major axis a and eccentricity e , which were defined in §3.1b, and by three angular orbital elements, i ,

¹ Conservation of the three components of \mathbf{L} and the three components of \mathbf{e} yields only five independent constraints because $\mathbf{L} \cdot \mathbf{e} = 0$.

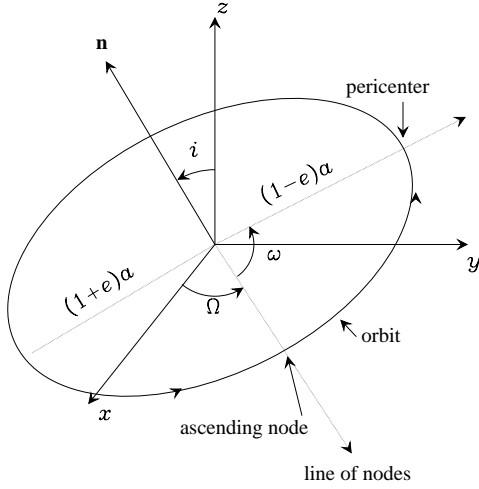


Figure E.1 Definition of Kepler orbital elements. The line of nodes is the intersection of the orbital plane with the xy -plane, and the angle Ω between it and the x axis is the longitude of the ascending node. The pericenter occurs where the orbit crosses a line from the center that makes an angle ω , the argument of pericenter, with the line of nodes.

Ω and ω , as shown in Figure E.1. In the limit $b \rightarrow 0$ the parameter c of equations (3.240) becomes $-GM/(2H)$, so by equation (3.32) it is equal to a . The variable e in equation (3.240) becomes $(1 - L^2/GMa)^{1/2}$ so by (3.25b) it is equal to the Kepler eccentricity having the same symbol. With equations (3.230) we have that a , e , and i are given by the Delaunay actions:

$$a = \frac{J_c^2}{GM} \quad ; \quad e = \sqrt{1 - \frac{J_b^2}{J_c^2}} \quad ; \quad i = \cos^{-1}(J_a/J_b). \quad (\text{E.6})$$

The inverse relations are

$$J_a = J_b \cos i \quad ; \quad J_b = J_c \sqrt{1 - e^2} \quad ; \quad J_c = \sqrt{GMa}. \quad (\text{E.7})$$

In the limit $b \rightarrow 0$, equations (3.240) reduce to equation (3.28) between r and η , so the latter is the eccentric anomaly. By the first of equations (3.241), the angle $\theta_c = \theta_3$ is simply $2\pi t/T_r$, where $T_r = 2\pi/\Omega_c$ is the radial period and t is the time elapsed since the last pericenter passage (eq. 3.28b). In celestial mechanics this angle is known as the mean anomaly and is usually denoted by ℓ . The second of equations (3.241) shows that

$$\theta_2 = \psi + \theta_3 - 2 \tan^{-1} \left(\sqrt{\frac{1+e}{1-e}} \tan\left(\frac{1}{2}\eta\right) \right). \quad (\text{E.8})$$

Equation (3.29) shows that the last term is simply $\psi_0 - \psi$, where ψ_0 is the angle measured in the orbit plane from the ascending node to the pericenter; this angle is the **argument of pericenter** and is usually denoted by ω (Figure E.1). Thus for a Kepler potential

$$\theta_2 = \omega + \theta_3 = \omega + \ell. \quad (\text{E.9})$$

The analogous relations for the other angle-action variables are given in Table E.1.

Table E.1 Angle-action variables in a Kepler potential

actions	$J_\phi = \sqrt{GMa(1-e^2)} \cos i$ $J_\vartheta = \sqrt{GMa(1-e^2)}(1 - \cos i)$ $J_r = \sqrt{GMa(1-\sqrt{1-e^2})}$
angles	$\theta_\phi = \Omega + \text{sgn}(L_z)(\omega + \ell)$; $\theta_\vartheta = \omega + \ell$; $\theta_r = \ell$
Hamiltonian	$-\frac{1}{2}(GM)^2/(J_r + J_\vartheta + J_\phi) = -\frac{1}{2}GM/a$
frequencies	$\Omega_\phi = \text{sgn}(J_\phi)\Omega_\vartheta$ $\Omega_\vartheta = (GM)^2/(J_r + J_\vartheta + J_\phi)^2 = (GM/a^3)^{1/2}$ $\Omega_r = \Omega_\vartheta$
actions	$J_1 = \sqrt{GMa(1-e^2)} \cos i$ $J_2 = \sqrt{GMa(1-e^2)}$ $J_3 = \sqrt{GMa(1-\sqrt{1-e^2})}$
angles	$\theta_1 = \Omega$; $\theta_2 = \omega + \ell$; $\theta_3 = \ell$
Hamiltonian	$-\frac{1}{2}(GM)^2/(J_2 + J_3)^2 = -\frac{1}{2}GM/a$
frequencies	$\Omega_1 = 0$; $\Omega_2 = (GM)^2/(J_2 + J_3)^3 = (GM/a^3)^{1/2}$; $\Omega_3 = \Omega_2$
actions	$J_a = \sqrt{GMa(1-e^2)} \cos i$ $J_b = \sqrt{GMa(1-e^2)}$ $J_c = \sqrt{GMa}$
angles	$\theta_a = \Omega$; $\theta_b = \omega$; $\theta_c = \ell$
Hamiltonian	$-\frac{1}{2}(GM)^2/J_c^2 = -\frac{1}{2}GM/a$
frequencies	$\Omega_a = 0$; $\Omega_b = 0$; $\Omega_c = (GM)^2/J_c^3 = (GM/a^3)^{1/2}$

NOTES: Actions and angles are expressed in terms of the standard Kepler orbital elements: semi-major axis a , eccentricity e , inclination i , longitude of the ascending node Ω , argument of pericenter ω , and mean anomaly ℓ . These are defined in §3.1b and Figure E.1. Unfortunately, Ω is also used for the frequency corresponding to a given action, but in this case it is always accompanied by a subscript.

Appendix F: Fluid mechanics

The basic principles of fluid mechanics play an important role in galaxy dynamics, both because fluid and stellar systems behave in similar ways and because gas dynamics is central to the formation of galaxies. We review here some of the concepts of fluid mechanics that are used in the book. For further reading see Landau & Lifshitz (2000).

F.1 Basic equations

The state of a fluid is specified by its density $\rho(\mathbf{x}, t)$, pressure $p(\mathbf{x}, t)$, and velocity field $\mathbf{v}(\mathbf{x}, t)$, and possibly by other thermodynamic functions such as the temperature $T(\mathbf{x}, t)$ and **specific entropy** or entropy per unit mass $s(\mathbf{x}, t)$.

F.1.1 Continuity equation

Consider an arbitrary closed volume V that is fixed in space and bounded by a surface S . The mass of fluid in this volume is $M(t) = \int_V d^3\mathbf{x} \rho(\mathbf{x}, t)$, and $M(t)$ changes

with time at a rate $dM/dt = \int_V d^3\mathbf{x} \partial\rho/\partial t$. The mass flowing out through the surface area element d^2S per unit time is $\rho\mathbf{v} \cdot d^2\mathbf{S}$, where $d^2\mathbf{S}$ is an outward-pointing vector, normal to the surface, with magnitude d^2S . Thus $dM/dt = -\oint_S d^2\mathbf{S} \cdot (\rho\mathbf{v})$ and hence

$$\int_V d^3\mathbf{x} \frac{\partial\rho}{\partial t} + \oint_S d^2\mathbf{S} \cdot (\rho\mathbf{v}) = 0. \quad (\text{F.1})$$

Using the divergence theorem (B.43),

$$\int_V d^3\mathbf{x} \left[\frac{\partial\rho}{\partial t} + \nabla \cdot (\rho\mathbf{v}) \right] = 0; \quad (\text{F.2})$$

since this result must hold for any volume, we arrive at the **continuity equation**

$$\frac{\partial\rho}{\partial t} + \nabla \cdot (\rho\mathbf{v}) = 0. \quad (\text{F.3})$$

In Cartesian coordinates this can be written

$$\frac{\partial\rho}{\partial t} + \frac{\partial}{\partial x_j}(\rho v_j) = 0, \quad (\text{F.4})$$

where we have used the summation convention (page 772).

In perturbation theory, we sometimes need to know the change in density resulting from a small displacement of the fluid. Let the fluid element at position \mathbf{x} be displaced to $\mathbf{x} + \epsilon\boldsymbol{\xi}(\mathbf{x})$, where ϵ is sufficiently small. If the displacement occurs at $t = t_0$, we may write $\mathbf{v}(\mathbf{x}, t) = \boldsymbol{\xi}(\mathbf{x})\delta(t - t_0)$, where δ denotes the delta function (Appendix C.1). If we now integrate equation (F.3) in time from just before to just after t_0 , we have

$$\rho_1 = -\nabla \cdot (\rho_0\boldsymbol{\xi}), \quad (\text{F.5})$$

where $\epsilon\rho_1(\mathbf{x})$ is the change in density at \mathbf{x} . We have replaced the density by its unperturbed value in the divergence, since it is multiplied by the small quantity $\boldsymbol{\xi}$. For more details see, for example, Chandrasekhar (1969).

F.1.2 Euler's equation

In a fluid that has no viscosity, the total pressure force acting on a volume is $-\oint_S d^2\mathbf{S} p$. In addition, there may be some external force, in particular from a gravitational potential $\Phi(\mathbf{x}, t)$. Thus Newton's second law reads

$$M \frac{d\mathbf{v}}{dt} = -\oint_S d^2\mathbf{S} p - M\nabla\Phi. \quad (\text{F.6})$$

According to the divergence theorem in the form (B.46), $\oint_S d^2\mathbf{S} p = \int_V d^3\mathbf{x} \nabla p$, and since equation (F.6) must hold for any volume V ,

$$\rho \frac{d\mathbf{v}}{dt} = -\nabla p - \rho\nabla\Phi. \quad (\text{F.7})$$

The quantity $d\mathbf{v}/dt$ in equation (F.7) is the acceleration of the fluid element at \mathbf{x} . This is not necessarily the same as $\partial\mathbf{v}/\partial t$, the rate of change of velocity at the point \mathbf{x} , since different fluid elements occupy this point at different times: for example, in a waterfall $\partial\mathbf{v}/\partial t = 0$ but $d\mathbf{v}/dt = \mathbf{g}$, the acceleration due to the Earth's gravity. More generally, if $w(\mathbf{x}, t)$ is any intrinsic property of the fluid—pressure, density, velocity, etc.—we define dw/dt to be the rate of change of w as seen by an observer traveling with the fluid. The quantity dw/dt is sometimes

referred to as the **Lagrangian** or **convective derivative** of w , to distinguish it from the **Eulerian derivative** $\partial w/\partial t$. The relation between these two derivatives is straightforward to derive. The change dw during the interval dt is the sum of the change at a given point in space, $(\partial w/\partial t)dt$, and the difference in w between two points separated by $d\mathbf{x} = \mathbf{v}dt$ at the same instant. The latter change is $\sum_{i=1}^3(\partial w/\partial x_i)dx_i = d\mathbf{x} \cdot \nabla w = \mathbf{v} \cdot \nabla w dt$. Thus

$$\frac{dw}{dt} = \frac{\partial w}{\partial t} + \mathbf{v} \cdot \nabla w. \quad (\text{F.8})$$

If we replace w in this general formula by the three components of the velocity \mathbf{v} , we have

$$\frac{d\mathbf{v}}{dt} = \frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla)\mathbf{v}. \quad (\text{F.9})$$

See equations (B.55) to (B.58) for $(\mathbf{v} \cdot \nabla)\mathbf{v}$ in various coordinate systems.

Combining equations (F.7) and (F.9) we arrive at **Euler's equation**

$$\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla)\mathbf{v} = -\frac{1}{\rho}\nabla p - \nabla\Phi. \quad (\text{F.10})$$

In Cartesian coordinates, using the summation convention, we have

$$\frac{\partial v_i}{\partial t} + v_j \frac{\partial v_i}{\partial x_j} = -\frac{1}{\rho} \frac{\partial p}{\partial x_i} - \frac{\partial \Phi}{\partial x_i}. \quad (\text{F.11})$$

A static fluid has $\mathbf{v} = 0$ and thus obeys the equation of **hydrostatic equilibrium**,

$$\nabla p = -\rho \nabla \Phi. \quad (\text{F.12})$$

Euler's equation can be regarded as a statement of conservation of momentum. The momentum per unit volume of the fluid is $\rho\mathbf{v}$. The rate of change of momentum per unit volume is

$$\frac{\partial}{\partial t}(\rho v_i) = v_i \frac{\partial \rho}{\partial t} + \rho \frac{\partial v_i}{\partial t}. \quad (\text{F.13})$$

Using the continuity and Euler equations (F.4) and (F.11) this can be rewritten as

$$\frac{\partial}{\partial t}(\rho v_i) = -v_i \frac{\partial}{\partial x_j}(\rho v_j) - \rho v_j \frac{\partial v_i}{\partial x_j} - \frac{\partial p}{\partial x_i} - \rho \frac{\partial \Phi}{\partial x_i} = -\frac{\partial \Pi_{ij}}{\partial x_j} - \rho \frac{\partial \Phi}{\partial x_i}, \quad (\text{F.14})$$

where

$$\Pi_{ij} \equiv \rho v_i v_j + p \delta_{ij}, \quad (\text{F.15})$$

and δ_{ij} is 1 if $i = j$ and zero otherwise. To interpret equation (F.14), we integrate over a volume V that is fixed in space and bounded by a surface S . We have

$$\begin{aligned} \frac{\partial}{\partial t} \int_V d^3\mathbf{x} \rho v_i &= - \int_V d^3\mathbf{x} \frac{\partial \Pi_{ij}}{\partial x_j} - \int_V d^3\mathbf{x} \rho \frac{\partial \Phi}{\partial x_i} \\ &= - \oint_S d^2 S_j \Pi_{ij} - \int_V d^3\mathbf{x} \rho \frac{\partial \Phi}{\partial x_i}, \end{aligned} \quad (\text{F.16})$$

where in the last equation we have used the divergence theorem (B.44), except that the vector F_j in that equation has been replaced by the tensor Π_{ij} . The left side is the rate of change of the i th component of the momentum contained in the volume V . The integral over $\partial\Phi/\partial x_i$ on the right is the rate at which momentum is added to this volume from the gravitational acceleration $-\nabla\Phi$. The surface integral on the right side therefore represents the rate at which momentum flows out of the volume through the surface S . In particular, $\sum_{j=1}^3 \Pi_{ij} d^2 S_j$ is the rate at which

the i th component of the momentum flows through the surface element $d^2\mathbf{S}$. The tensor Π_{ij} is therefore called the **momentum flux tensor**.

F.1.3 Energy equation

A fluid carries kinetic energy, internal energy, and gravitational potential energy. We ignore other energy forms, such as magnetic fields (Kulsrud 2005).

The kinetic energy per unit volume is $K(\mathbf{x}, t) = \frac{1}{2}\rho v^2$. Using the continuity equation (F.3), Euler's equation (F.10), and the vector identity $\mathbf{v} \cdot [(\mathbf{v} \cdot \nabla)\mathbf{v}] = v_i v_j (\partial v_i / \partial x_j) = \frac{1}{2}\mathbf{v} \cdot \nabla v^2$, the rate of change of kinetic energy per unit volume can be written

$$\begin{aligned} \frac{\partial K}{\partial t} &= \frac{1}{2} \frac{\partial \rho}{\partial t} v^2 + \rho \mathbf{v} \cdot \frac{\partial \mathbf{v}}{\partial t} \\ &= -\frac{1}{2} v^2 \nabla \cdot (\rho \mathbf{v}) - \frac{1}{2} \rho \mathbf{v} \cdot \nabla v^2 - \mathbf{v} \cdot \nabla p - \rho \mathbf{v} \cdot \nabla \Phi \\ &= -\frac{1}{2} \nabla \cdot (\rho \mathbf{v} v^2) - \mathbf{v} \cdot \nabla p - \rho \mathbf{v} \cdot \nabla \Phi. \end{aligned} \quad (\text{F.17})$$

The internal energy per unit volume is $U(\mathbf{x}, t) = \rho u$, where u is the specific internal energy (internal energy per unit mass), which is determined by the equation of state of the fluid. The rate of change of internal energy per unit volume is

$$\frac{\partial U}{\partial t} = \frac{\partial \rho}{\partial t} u + \rho \frac{\partial u}{\partial t} = -u \nabla \cdot (\rho \mathbf{v}) + \rho \frac{\partial u}{\partial t}. \quad (\text{F.18})$$

The potential energy per unit volume due to an external gravitational field is $W_e = \rho \Phi_e$, and the energy due to the self-gravity of the fluid is $W_s = \frac{1}{2}\rho \Phi_s$, where $\nabla^2 \Phi_s = 4\pi G \rho$ (eq. 2.18). Setting $W \equiv W_e + W_s$ and $\Phi \equiv \Phi_e + \Phi_s$, we have

$$\begin{aligned} \frac{\partial W}{\partial t} &= \frac{\partial}{\partial t} \rho (\Phi_e + \frac{1}{2}\Phi_s) = \frac{\partial \rho}{\partial t} (\Phi_e + \frac{1}{2}\Phi_s) + \rho \frac{\partial (\Phi_e + \frac{1}{2}\Phi_s)}{\partial t} \\ &= -(\Phi_e + \frac{1}{2}\Phi_s) \nabla \cdot (\rho \mathbf{v}) + \rho \frac{\partial (\Phi_e + \frac{1}{2}\Phi_s)}{\partial t}. \end{aligned} \quad (\text{F.19})$$

The total rate of change of energy is

$$\begin{aligned} \frac{\partial}{\partial t} (K + U + W) &= -\nabla \cdot [\rho \mathbf{v} (\frac{1}{2}v^2 + u + \Phi)] - \mathbf{v} \cdot \nabla p + \rho \mathbf{v} \cdot \nabla u \\ &\quad + \rho \frac{\partial u}{\partial t} + \frac{1}{2}\Phi_s \nabla \cdot (\rho \mathbf{v}) + \rho \frac{\partial (\Phi_e + \frac{1}{2}\Phi_s)}{\partial t} \\ &= -\nabla \cdot [\rho \mathbf{v} (\frac{1}{2}v^2 + u + \Phi)] - \mathbf{v} \cdot \nabla p + \rho \frac{du}{dt} + \frac{1}{2}\Phi_s \nabla \cdot (\rho \mathbf{v}) + \rho \frac{\partial (\Phi_e + \frac{1}{2}\Phi_s)}{\partial t}, \end{aligned} \quad (\text{F.20})$$

where we have used equation (F.8) to replace the Eulerian derivative of u by the convective derivative.

To proceed further, we use the second law of thermodynamics, which states that the specific entropy of a fluid element changes if heat flows into or out of it, or if heat is generated within the element by viscous dissipation, nuclear reactions, or other mechanisms. Let \mathbf{q} be the heat flux, so the rate of heat flow out of a volume is $\oint_S d^2\mathbf{S} \cdot \mathbf{q}$, which is equal to $\int_V d^3\mathbf{x} \nabla \cdot \mathbf{q}$ by the divergence theorem, equation (B.43). Let ϵ be the rate of internal energy production per unit mass. Then according to the second law, the specific entropy changes at a rate given by

$$\rho T \frac{ds}{dt} = -\nabla \cdot \mathbf{q} + \rho \epsilon. \quad (\text{F.21})$$

The specific entropy is related to the volume per unit mass $V = 1/\rho$, the internal energy u , and the specific enthalpy $h \equiv u + pV = u + p/\rho$ by

$$du = Tds - pdV = Tds + \frac{p}{\rho^2}d\rho \quad ; \quad dh = Tds + Vdp = Tds + \frac{dp}{\rho}. \quad (\text{F.22})$$

We can now simplify equation (F.20) by using (F.22) to replace du/dt by $Tds/dt + (p/\rho^2)d\rho/dt$ and replacing u in the first term on the right side by $h = u + p/\rho$. After some algebra we obtain

$$\frac{\partial}{\partial t}(K + U + W) = -\nabla \cdot [\rho\mathbf{v}(\frac{1}{2}v^2 + h + \Phi) + \mathbf{q}] + \rho\epsilon + \frac{1}{2}\Phi_s \nabla \cdot (\rho\mathbf{v}) + \rho \frac{\partial(\Phi_e + \frac{1}{2}\Phi_s)}{\partial t}. \quad (\text{F.23})$$

In the simplest case, the heat flux is proportional to the temperature gradient, so we may write

$$\mathbf{q} = -\kappa\nabla T, \quad (\text{F.24})$$

where κ is the thermal conductivity, which in general may depend on T and ρ .

Using the equation of continuity,

$$\begin{aligned} \frac{1}{2}\Phi_s \nabla \cdot (\rho\mathbf{v}) + \frac{1}{2}\rho \frac{\partial\Phi_s}{\partial t} &= \frac{1}{2} \left(\rho \frac{\partial\Phi_s}{\partial t} - \frac{\partial\rho}{\partial t} \Phi_s \right) \\ &= \frac{1}{8\pi G} \nabla \cdot \left(\frac{\partial\Phi_s}{\partial t} \nabla\Phi_s - \Phi_s \nabla \frac{\partial\Phi_s}{\partial t} \right), \end{aligned} \quad (\text{F.25})$$

where the last equality can be verified using Poisson's equation.

Thus we have

$$\begin{aligned} \frac{\partial}{\partial t}(K + U + W) &= \frac{\partial}{\partial t} [\rho(\frac{1}{2}v^2 + u + \Phi_e + \frac{1}{2}\Phi_s)] \\ &= -\nabla \cdot \left[\rho\mathbf{v}(\frac{1}{2}v^2 + h + \Phi) - \kappa\nabla T \right. \\ &\quad \left. + \frac{1}{8\pi G} \left(\Phi_s \nabla \frac{\partial\Phi_s}{\partial t} - \frac{\partial\Phi_s}{\partial t} \nabla\Phi_s \right) \right] + \rho\epsilon + \rho \frac{\partial\Phi_e}{\partial t}, \end{aligned} \quad (\text{F.26})$$

where $\Phi = \Phi_s + \Phi_e$. The terms on the right side of equation (F.26) can be interpreted as follows. The quantity $\rho\mathbf{v}(\frac{1}{2}v^2 + h + \Phi)$ is the **energy flux vector**, which represents the rate at which energy is convected by the fluid. The term proportional to κ represents the conduction of heat. The terms involving Φ_s represent the rate at which energy is transported by self-gravity. The term $\rho\epsilon$ represents the rate of heat generation by internal processes. The term involving the external potential Φ_e reflects changes in the zero-point of the external gravitational potential (cf. eq. D.10).

F.1.4 Equation of state

To relate the pressure and density, we need an **equation of state** $p = p(\rho, s)$ or $p = p(\rho, T)$. For our purposes it is usually sufficient to consider the simple case of a **barotropic** equation of state, where the pressure is determined by the density,

$$p = p(\rho). \quad (\text{F.27})$$

The most important examples of barotropic equations of state arise when the fluid is **isentropic** or **adiabatic**, that is, has constant specific entropy. In an isentropic

fluid the thermodynamic relations (F.22) simplify to

$$du = \frac{p}{\rho^2} d\rho \quad ; \quad dh = \frac{dp}{\rho}, \quad (\text{F.28})$$

so we may write

$$u(\rho) = \int_0^\rho d\rho' \frac{p(\rho')}{\rho'^2} \quad ; \quad h(\rho) = \int_0^\rho \frac{dp(\rho')}{\rho'} = \int_0^\rho \frac{d\rho'}{\rho'} \frac{dp(\rho')}{d\rho'}. \quad (\text{F.29})$$

For a barotropic fluid, Euler's equation (F.10) becomes simply

$$\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} = -\nabla(h + \Phi). \quad (\text{F.30})$$

F.2 The ideal gas

One of the simplest fluids is the ideal gas, whose equation of state is

$$p = \frac{\rho k_B T}{m}, \quad (\text{F.31})$$

where k_B is Boltzmann's constant and m is the mass of a single molecule. The pressure in an ideal gas is related to the velocity dispersion of the molecules in each direction by

$$p = \rho \overline{v_x^2} = \rho \overline{v_y^2} = \rho \overline{v_z^2}, \quad (\text{F.32})$$

and thus

$$\overline{v_x^2} = \overline{v_y^2} = \overline{v_z^2} = \frac{1}{3} \overline{v^2} = \frac{k_B T}{m}. \quad (\text{F.33})$$

The first of equations (F.22) implies that

$$ds = \frac{1}{T} du - \frac{p}{T\rho^2} d\rho. \quad (\text{F.34})$$

Writing the internal energy u as a function of T and ρ , we have

$$ds = \frac{1}{T} \left(\frac{\partial u}{\partial T} \right)_\rho dT + \left[\frac{1}{T} \left(\frac{\partial u}{\partial \rho} \right)_T - \frac{p}{T\rho^2} \right] d\rho. \quad (\text{F.35})$$

If we write the entropy s as a function of T and ρ , we have

$$ds = \left(\frac{\partial s}{\partial T} \right)_\rho dT + \left(\frac{\partial s}{\partial \rho} \right)_T d\rho. \quad (\text{F.36})$$

Since dT and $d\rho$ are arbitrary, their coefficients in the previous two equations must be the same. Hence

$$\left(\frac{\partial s}{\partial T} \right)_\rho = \frac{1}{T} \left(\frac{\partial u}{\partial T} \right)_\rho \quad ; \quad \left(\frac{\partial s}{\partial \rho} \right)_T = \frac{1}{T} \left(\frac{\partial u}{\partial \rho} \right)_T - \frac{p}{T\rho^2}. \quad (\text{F.37})$$

The derivative of the first expression with respect to ρ is $\partial^2 s / \partial \rho \partial T$, which must equal the derivative of the second expression with respect to T . Thus

$$\left(\frac{\partial}{\partial \rho} \right)_T \frac{1}{T} \left(\frac{\partial u}{\partial T} \right)_\rho = \left(\frac{\partial}{\partial T} \right)_\rho \left[\frac{1}{T} \left(\frac{\partial u}{\partial \rho} \right)_T - \frac{p}{T\rho^2} \right]. \quad (\text{F.38})$$

This expression can be simplified to

$$\left(\frac{\partial u}{\partial \rho} \right)_T = \frac{p}{\rho^2} - \frac{T}{\rho^2} \left(\frac{\partial p}{\partial T} \right)_\rho. \quad (\text{F.39})$$

For the ideal gas, the equation of state (F.31) then implies that $(\partial u/\partial \rho)_T = 0$; thus the internal energy per unit mass is a function of temperature alone, $u = u(T)$. For an ideal gas of point particles, the internal energy is simply the kinetic energy associated with random motions,

$$u = \frac{1}{2} \overline{v^2} = \frac{3k_B T}{2m}; \quad (\text{F.40})$$

more generally, if the particles have q internal degrees of freedom,

$$u = \frac{(q+3)k_B T}{2m}. \quad (\text{F.41})$$

Equation (F.35) can now be evaluated as

$$ds = \frac{k_B}{m} \left[\frac{q+3}{2} \frac{dT}{T} - \frac{d\rho}{\rho} \right], \quad (\text{F.42})$$

which can be integrated to yield the specific entropy,

$$s = \frac{k_B}{m} \ln \left(\frac{T^{(q+3)/2}}{\rho} \right) + \text{constant}. \quad (\text{F.43})$$

There are two important special cases in which the ideal gas is barotropic. (i) If the temperature of the gas is fixed everywhere at T_0 , then the fluid is said to be **isothermal**. In this case

$$p = K\rho, \quad (\text{F.44})$$

with $K = k_B T_0/m$. (ii) If the fluid is isentropic, equation (F.43) implies that

$$\rho \propto T^{(q+3)/2} = T^{1/(\gamma-1)}, \quad (\text{F.45})$$

where we have replaced q by a new constant $\gamma \equiv (q+5)/(q+3)$. Using the equation of state (F.31), we then have

$$p = K\rho^\gamma, \quad (\text{F.46})$$

where K is a constant that depends on the specific entropy. A barotropic equation of state with the power-law form (F.46) is known as a **polytropic** equation of state; the isothermal equation of state (F.44) is polytropic with $\gamma = 1$.

F.3 Sound waves

We examine the evolution of small disturbances in a stationary barotropic fluid of constant density ρ_0 . We assume that the gravitational field $\nabla\Phi = 0$. If the fluid is subjected to a small perturbation, we may write

$$\rho(\mathbf{x}, t) = \rho_0 + \epsilon\rho_1(\mathbf{x}, t); \quad h(\mathbf{x}, t) = h_0 + \epsilon h_1(\mathbf{x}, t); \quad \mathbf{v}(\mathbf{x}, t) = \epsilon\mathbf{v}_1(\mathbf{x}, t), \quad (\text{F.47})$$

where $\epsilon \ll 1$ and the quantities with subscripts 0 and 1 are of the same order of magnitude. Substituting (F.47) into equations (F.3), (F.29), and (F.30), we find that the terms that are independent of ϵ vanish, and discarding terms proportional to ϵ^2 we obtain

$$\frac{\partial \rho_1}{\partial t} + \rho_0 \nabla \cdot \mathbf{v}_1 = 0 \quad ; \quad h_1 = \left(\frac{dp}{d\rho} \right)_{\rho_0} \frac{\rho_1}{\rho_0} \quad ; \quad \frac{\partial \mathbf{v}_1}{\partial t} = -\nabla h_1. \quad (\text{F.48})$$

We differentiate the first of these with respect to time and eliminate \mathbf{v}_1 and h_1 to obtain the **wave equation**

$$\frac{\partial^2 \rho_1}{\partial t^2} - v_s^2 \nabla^2 \rho_1 = 0, \quad (\text{F.49})$$

where

$$v_s^2 \equiv \left(\frac{dp}{d\rho} \right)_{\rho_0}. \quad (\text{F.50})$$

The solution of this equation is simplest to understand when ρ_1 depends on only one coordinate, say, x . Then

$$\frac{\partial^2 \rho_1}{\partial t^2} - v_s^2 \frac{\partial^2 \rho_1}{\partial x^2} = 0. \quad (\text{F.51})$$

The general solution is

$$\rho_1 = f_+(x - v_s t) + f_-(x + v_s t), \quad (\text{F.52})$$

where f_+ and f_- are arbitrary functions. This solution consists of two superimposed waves, both traveling at speed v_s —one (f_+) to the right and the other (f_-) to the left. The disturbances are sound waves, and v_s is the **sound speed**.

The simplest examples of sound waves are uniform wavetrains of the form

$$\rho_1 = A \cos(kx - \omega t + \text{constant}), \quad (\text{F.53})$$

which satisfy equation (F.51) when

$$\omega^2 = v_s^2 k^2. \quad (\text{F.54})$$

This is the **dispersion relation** for sound waves.

In general the perturbations in a sound wave are adiabatic, and hence equations (F.31), (F.46), and (F.50) yield for the sound speed in an ideal gas

$$v_s = \sqrt{\frac{\gamma p_0}{\rho_0}} = \sqrt{\frac{\gamma k_B T_0}{m}} = \sqrt{\frac{1}{3} \gamma v^2}, \quad (\text{F.55})$$

where T_0 is the temperature of the equilibrium gas and the second equality follows from (F.33). Hence the speed of sound is close to the RMS speed of the molecules.

F.3.1 Energy and momentum in sound waves

We now compute the density and flux of energy and momentum in a train of sound waves. We assume that the waves are traveling in the x -direction and that $v_y = v_z = 0$. In the absence of a gravitational field, the energy density of the fluid is $E = \rho(\frac{1}{2}v_x^2 + u)$ (eq. F.26; see Problem 5.7 for sound waves with self-gravity). In the undisturbed state $\rho = \rho_0$ is a constant and $v_x = 0$. In the presence of a wave, the value of ρu peaks where ρ is largest and is smallest at troughs of ρ . These extrema of ρu occur where the fluid is stationary and there is no kinetic contribution to E , so E is an oscillating function. We want to identify the amount by which E is changed by the wave when averaged over a cycle, for it is this net change in E that must be supplied to the fluid to establish a wavetrain. Only this average contribution to E is propagated by the wave rather than sloshing backwards and forwards during each cycle. The derivation below follows Lighthill (1978).

We make use of the energy transport equation (F.26), which in the present context reads

$$\frac{\partial E}{\partial t} + \frac{\partial F}{\partial x} = 0, \quad \text{where} \quad F = \rho v_x (\frac{1}{2} v_x^2 + h) \quad (\text{F.56})$$

is the energy flux. We write

$$E = (\rho u)_0 + E' + E_w \quad ; \quad F = F' + F_w, \quad (\text{F.57})$$

where E' and E_w are the non-wave and wave energy densities, and F' and F_w are the non-wave and wave energy fluxes. We demand that the non-wave and wave contributions separately satisfy transport equations analogous to (F.56), as they must if they are to be interpreted as densities and fluxes of conserved quantities.

Since the fluid is assumed barotropic, ρu is a function of only ρ , so it can be expanded in a Taylor series,

$$\rho u = (\rho u)_0 + \left. \frac{d(\rho u)}{d\rho} \right|_0 \Delta\rho + \frac{1}{2} \left. \frac{d^2(\rho u)}{d\rho^2} \right|_0 (\Delta\rho)^2 + O(\Delta\rho^3), \quad (\text{F.58})$$

where $\Delta\rho = \rho - \rho_0$. Equation (F.28) and the definition of specific enthalpy h imply that

$$\left. \frac{d(\rho u)}{d\rho} \right|_0 = u_0 + \frac{p_0}{\rho_0} = h_0 \quad ; \quad \left. \frac{d^2(\rho u)}{d\rho^2} \right|_0 = \frac{1}{\rho_0} \left. \frac{dp}{d\rho} \right|_0 = \frac{v_s^2}{\rho_0}; \quad (\text{F.59})$$

thus

$$\rho u = (\rho u)_0 + h_0 \Delta\rho + \frac{v_s^2}{2\rho_0} (\Delta\rho)^2 + O(\Delta\rho^3). \quad (\text{F.60})$$

We take the non-wave energy and flux to be

$$E' \equiv h_0 \Delta\rho \quad ; \quad F' = h_0 \rho v_x, \quad (\text{F.61a})$$

so

$$\frac{\partial E'}{\partial t} + \frac{\partial F'}{\partial x} = h_0 \left[\frac{\partial \rho}{\partial t} + \frac{\partial(\rho v_x)}{\partial x} \right], \quad (\text{F.61b})$$

which equals zero because of the continuity equation (F.4). Since both the total energy and the non-wave energy satisfy transport equations of the form (F.56), their difference, the wave energy, must do so as well. The wave energy and wave flux are then

$$\begin{aligned} E_w &= E - E' - (\rho u)_0 = \frac{1}{2} \rho v_x^2 + \rho u - (\rho u)_0 - h_0(\rho - \rho_0), \\ F_w &= F - F' = \frac{1}{2} \rho v_x^3 + \rho v_x (h - h_0). \end{aligned} \quad (\text{F.62})$$

Note that the wave and non-wave energy densities are defined exactly, not merely to some order in perturbation theory. However, for most purposes it is sufficient to work in the approximation that the amplitude of the sound wave is small. Let us therefore write the density $\rho(\mathbf{x}, t) = \rho_0 + \rho_1(\mathbf{x}, t) + \rho_2(\mathbf{x}, t) + \dots$, where ρ_n is of order A^n , and A is the small amplitude of the sound wave. We write a similar expansion for v_x , noting that $v_{x0} = 0$ since the equilibrium fluid is stationary. We have

$$\begin{aligned} E_w &= \frac{1}{2} \rho_0 v_{x1}^2 + \frac{v_s^2}{2\rho_0} \rho_1^2 + O(A^3), \\ F_w &= \rho_0 v_{x1} h_1 = v_s^2 \rho_1 v_{x1} + O(A^3). \end{aligned} \quad (\text{F.63})$$

From the linearized continuity equation (F.48), the first-order density and velocity perturbations associated with a wavetrain may be written as

$$\begin{aligned} \rho_1 &= A_\rho \cos(kx - \omega t) \\ v_{x1} &= A_v \cos(kx - \omega t), \end{aligned} \quad \text{where} \quad A_v = \frac{\omega}{k\rho_0} A_\rho. \quad (\text{F.64})$$

The wave energy density and flux can thus be written to second-order in the amplitude as

$$\begin{aligned} E_w &= \frac{1}{2}\rho_0 A_v^2 \cos^2(kx - \omega t) + \frac{v_s^2 A_\rho^2}{2\rho_0} \cos^2(kx - \omega t) \\ &= \frac{\omega^2 + v_s^2 k^2}{2k^2 \rho_0} A_\rho^2 \cos^2(kx - \omega t) = \frac{v_s^2}{\rho_0} A_\rho^2 \cos^2(kx - \omega t), \quad (\text{F.65}) \\ F_w &= v_s^2 \rho_1 v_{x1} = \frac{v_s^2 \omega}{k \rho_0} A_\rho^2 \cos^2(kx - \omega t) = \pm v_s E_w, \end{aligned}$$

where we have used the dispersion relation $\omega = \pm v_s k$ (eq. F.54), and the \pm signs correspond to waves traveling to the right and to the left. We conclude that the energy flux in a plane sound wave equals the energy density times the speed of sound; that is, the energy of the wave propagates at the sound speed.

Averaging over one cycle, we have

$$\langle E_w \rangle = \frac{1}{2}\rho_0 A_v^2 \quad ; \quad \langle F_w \rangle = \frac{\omega}{2k} \rho_0 A_v^2 = \pm v_s \langle E_w \rangle. \quad (\text{F.66})$$

The transport equation (F.56) for the wave energy density and flux E_w and F_w describes conservation of energy to second-order in the wave amplitude with a constant multiple of the continuity equation (F.61b) excluded. This definition of the wave energy eliminates much larger first-order terms that have no bearing on the behavior of the waves.

We may calculate the momentum density and flux of the wave similarly. The momentum density in the fluid containing the wavetrain we are considering is $P = \rho v_x$, and from equation (F.14) the momentum flux (rate of transfer of x -momentum in the x -direction) is $\Pi = \rho v_x^2 + p$. We write

$$P = P' + P_w \quad ; \quad \Pi = p_0 + \Pi' + \Pi_w, \quad (\text{F.67})$$

where $P' = \rho_0 v_x$, $\Pi' = \frac{1}{2}\rho_0 v_x^2 + \rho_0(h - h_0)$, and

$$\begin{aligned} P_w &= (\rho - \rho_0)v_x = \rho_1 v_{x1} + O(A^3) = \pm \frac{v_s}{\rho_0} A_\rho^2 \cos^2(kx - \omega t) + O(A^3) \\ \Pi_w &= (\rho - \frac{1}{2}\rho_0)v_x^2 + (p - p_0) - \rho_0(h - h_0) = \frac{1}{2}\rho_0 v_{x1}^2 + \frac{v_s^2}{2\rho_0} \rho_1^2 + O(A^3) \quad (\text{F.68}) \\ &= \frac{v_s^2}{\rho_0} A_\rho^2 \cos^2(kx - \omega t) + O(A^3) \end{aligned}$$

are the momentum density and momentum flux of the wave. These satisfy the conservation equation

$$\frac{\partial P_w}{\partial t} + \frac{\partial \Pi_w}{\partial x} = 0, \quad (\text{F.69})$$

and $\Pi_w = \pm v_s P_w$, so the momentum of the wave is propagated at the sound speed. The first-order non-wave contributions P' and Π' also satisfy a conservation equation,

$$\frac{\partial P'}{\partial t} + \frac{\partial \Pi'}{\partial x} = 0, \quad (\text{F.70})$$

which is simply ρ_0 times Euler's equation.

The transport of energy and momentum through the fluxes F_w and Π_w is sometimes called **advective** or **Reynolds** transport. Advective transport is transport of a quantity due to bulk motions of the fluid (as opposed, say, to diffusion or

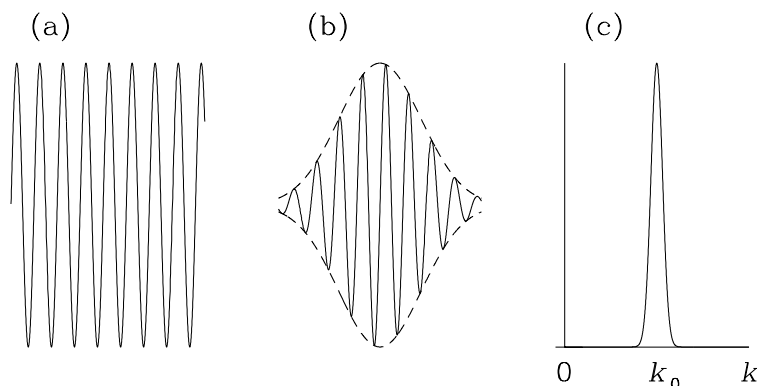


Figure F.1 (a) The real part of a wavetrain with wavenumber k_0 . (b) A typical wave packet. The dashed curve is the envelope or amplitude of the wave packet. (c) The absolute value $|F(k)|$ of the Fourier amplitude of the packet.

other microscopic processes). Reynolds transport is associated with the Reynolds stress $\langle \rho v_i v_j \rangle$. See Whitham (1974) or Lighthill (1978) for general discussions of energy and momentum transport by waves.

F.4 Group velocity

The dispersion relation for sound waves $\omega(k) = \pm v_s k$ (eq. F.54) is linear in the wavenumber k . This is not a general property of all waves; for example, the dispersion relation for sound waves in a medium with self-gravity has the form $\omega^2(k) = v_s^2 k^2 - 4\pi G \rho_0$ (eq. 5.34). If $\omega(k)$ is not linear in k , the medium is said to be **dispersive**.

One example of a wave in a dispersive medium is a uniform wavetrain of the form $\exp\{i[k_0 x - \omega(k_0)t]\}$, as shown in Figure F.1a. However, any wavetrain of this sort is somewhat unphysical, since it extends throughout all space. To describe a spatially localized wave is more difficult, since general solutions of the form (F.52) do not exist for a dispersive medium. Instead, we must construct a **wave packet**, a wave whose amplitude is non-zero over a region whose extent is finite, but changes only on scales much larger than the wavelength. A typical wave packet is shown in Figure F.1b. The wave packet can be represented mathematically by superimposing uniform wavetrains:

$$f(x, t) = \int_{-\infty}^{\infty} \frac{dk}{2\pi} F(k) e^{i[kx - \omega(k)t]}; \quad (\text{F.71})$$

with this definition, $F(k) \exp[-i\omega(k)t]$ is the spatial Fourier transform of $f(x, t)$, as defined in equation (B.67).

In any small region, the wave packet looks like a wavetrain of fixed wavenumber, say, k_0 . Hence $F(k)$ will be non-zero only for values of k near k_0 (see Figure F.1c), and we may therefore expand $\omega(k)$ in a Taylor series, $\omega(k) \simeq$

$\omega(k_0) + (k - k_0)v_g$, where

$$v_g \equiv \left(\frac{d\omega(k)}{dk} \right)_{k_0}. \quad (\text{F.72})$$

When we write $k = k_0 + u$, equation (F.71) becomes

$$f(x, t) \simeq e^{i[k_0x - \omega(k_0)t]} \int_{-\infty}^{\infty} \frac{du}{2\pi} F(k_0 + u) e^{iu(x - v_g t)}. \quad (\text{F.73})$$

We can rewrite this expression as

$$f(x, t) \simeq e^{i[k_0x - \omega(k_0)t]} A(x - v_g t), \quad (\text{F.74})$$

where

$$A(x) \equiv \int_{-\infty}^{\infty} \frac{du}{2\pi} F(k_0 + u) e^{iux}. \quad (\text{F.75})$$

Because $F(k)$ is non-zero only for k near k_0 , the integrand in this equation is non-zero only for $|u| \ll k_0$. Thus $A(x)$ varies much more slowly than $\exp\{i[k_0x - \omega(k_0)t]\}$, and we can picture the wave at any instant as having a rapidly varying phase $\propto \exp(ik_0x)$ modulated by a slowly varying amplitude $|A(x - v_g t)|$ (the dashed line in Figure F.1b). This argument shows that the envelope of the wave propagates with a velocity v_g , which is known as the **group velocity**. The group velocity represents the velocity of any true physical disturbance, since a disturbance arising from any physical source is always localized.

The group velocity is distinct from the **phase velocity** v_p , the velocity of a given crest of the wave. On a crest $kx - \omega t = \text{constant}$; thus

$$v_p = \frac{\omega}{k}. \quad (\text{F.76})$$

The phase velocity and group velocity are equal in a non-dispersive medium, for then $\omega(k)$ is linear in k . It can be shown that the energy of a dispersive wave propagates at the group velocity (Whitham 1974). For the sound waves investigated in Appendix F.3.1, $v_p = v_g = v_s$, the sound speed.

Appendix G: Discrete Fourier transforms

The discrete Fourier transform of a set of K numbers, $\{x_k\}$ ($k = 0, \dots, K-1$), is

$$\hat{x}_m \equiv \sum_{k=0}^{K-1} x_k e^{-2\pi i k m / K} \quad (m = 0, \dots, K-1). \quad (\text{G.1a})$$

We have:

Discrete Fourier transform theorem *If the K numbers \hat{x}_m are defined by equation (G.1a), then*

$$x_{k'} = \frac{1}{K} \sum_{m=0}^{K-1} \hat{x}_m e^{2\pi i k' m / K}. \quad (\text{G.1b})$$

Proof: We multiply both sides of equation (G.1a) by $(1/K)e^{2\pi i k' m / K}$ and sum over m , to obtain

$$\frac{1}{K} \sum_{m=0}^{K-1} \hat{x}_m e^{2\pi i k' m / K} = \frac{1}{K} \sum_{m=0}^{K-1} \sum_{k=0}^{K-1} x_k e^{2\pi i (k' - k) m / K} = \frac{1}{K} \sum_{k=0}^{K-1} x_k \sum_{m=0}^{K-1} e^{2\pi i (k' - k) m / K}. \quad (\text{G.2})$$

The inner sum on the extreme right of equation (G.2) is a geometric series, with sum K if $k' = k$, or (since k and k' are integers and $|k' - k| < K$)

$$\frac{1 - e^{2\pi i(k' - k)}}{1 - e^{2\pi i(k' - k)/K}} = 0 \text{ if } k' \neq k. \tag{G.3}$$

Thus the only contributing term in the outer sum on the right of equation (G.3) is that for which $k = k'$. ◁

In 1965 Cooley and Tukey published an algorithm, known as the **fast Fourier transform**, by which the discrete transform $\{\hat{x}_m\}$ of K numbers x_k can be evaluated in only $\sim K \ln(K)$ multiplications and additions.¹ For large K this is far fewer operations than the K^2 operations required by a direct evaluation of the sums of equation (G.1a) (Press et al. 1986).

Discrete Fourier transforms share many of the properties of continuous Fourier transforms. To demonstrate these properties, it is necessary to define the quantities x_k for k outside the range $(0, K - 1)$ by the rule

$$x_k = x_{k+mK} \quad \text{for all integer } m, \tag{G.4}$$

i.e., to assume that x_k is periodic with period K . Note that \hat{x}_m is already periodic with period K (see eq. G.1a). We may now prove:

Discrete Fourier convolution theorem *If the three sets of K numbers $\{x_k\}$, $\{y_k\}$, $\{z_k\}$ satisfy (G.4) and are related by*

$$z_k = \sum_{k'=0}^{K-1} y_{k-k'} x_{k'}, \quad \text{then } \hat{z}_m = \hat{y}_m \hat{x}_m. \tag{G.5}$$

Proof: We take the discrete Fourier transform of both sides of equation (G.5) and then rearrange the resulting double sum. We have

$$\hat{z}_m = \sum_{k=0}^{K-1} e^{-2\pi i k m / K} \sum_{k'=0}^{K-1} y_{k-k'} x_{k'} = \sum_{k'=0}^{K-1} x_{k'} e^{-2\pi i k' m / K} \sum_{k=0}^{K-1} y_{k-k'} e^{-2\pi i (k-k') m / K}. \tag{G.6}$$

If we now define $k'' \equiv k - k'$, the inner sum in equation (G.6) becomes \hat{y}_m . This is independent of k' , so it may be taken out of the outer sum, which then yields \hat{x}_m . ◁

Application to James's Fourier potential solver Given $K - 1$ real numbers x_k ($k = 1, \dots, K - 1$) we define the sine transform of the set to be the numbers

$$x^\alpha(S) \equiv \sum_{k=1}^{K-1} x_k \sin\left(\frac{\pi \alpha k}{K}\right) \quad (j = 1, \dots, K - 1). \tag{G.7a}$$

Multiplying both sides of this equation by $\sin(\pi \alpha j / K)$ and summing over α we find that the inverse transformation is

$$x_j = \frac{2}{K} \sum_{\alpha=1}^{K-1} x^\alpha(S) \sin\left(\frac{\pi \alpha j}{K}\right). \tag{G.7b}$$

¹ The Cooley–Tukey algorithm requires that $K = p^m$, where m is an integer and p is a prime number; usually $p = 2$.

Analogously we define the cosine transform of a set of K numbers x_k ($k = 0, \dots, K-1$) to be²

$$x^\alpha(C) \equiv \sum_{k=0}^{K-1} x_k \cos\left(\frac{\pi\alpha k}{K}\right). \quad (\text{G.8})$$

There is no elegant formula for recovering the x_k from the $x^\alpha(C)$, but the matrix that is the inverse of $A_{\alpha k} \equiv \cos(\pi\alpha k/K)$ can be calculated numerically if required.

Sine transforms help us to convolve values of the density on a grid with a softening kernel (§2.9). For simplicity we consider the two-dimensional case, but the results generalize straightforwardly to three dimensions. If the density at grid point (i, j) is ρ_{ij} , the potential values are

$$\Phi_{ij} = \sum_{kl} S_{i-k, j-l} \rho_{kl}, \quad (\text{G.9})$$

where S is the softening kernel. Since this is an even function of its indices, i.e., $S_{ij} = S_{-i, j}$ etc., and it will be evaluated only for $|i-j| \leq K$, it can be written as a sum of cosines:

$$S_{ij} = \sum_{\alpha\beta} S_C^{\alpha\beta} \cos\left(\frac{\pi\alpha i}{K}\right) \cos\left(\frac{\pi\beta j}{K}\right). \quad (\text{G.10})$$

Using this expression to eliminate $S_{i-k, j-l}$ from (G.9), we have

$$\begin{aligned} \Phi_{ij} &= \sum_{kl} \sum_{\alpha\beta} S_C^{\alpha\beta} \cos\left(\frac{\pi\alpha(i-k)}{K}\right) \cos\left(\frac{\pi\beta(j-l)}{K}\right) \rho_{kl} \\ &= \sum_{kl} \sum_{\alpha\beta} S_C^{\alpha\beta} \left[\cos\left(\frac{\pi\alpha i}{K}\right) \cos\left(\frac{\pi\alpha k}{K}\right) + \sin\left(\frac{\pi\alpha i}{K}\right) \sin\left(\frac{\pi\alpha k}{K}\right) \right] \\ &\quad \times \left[\cos\left(\frac{\pi\beta j}{K}\right) \cos\left(\frac{\pi\beta l}{K}\right) + \sin\left(\frac{\pi\beta j}{K}\right) \sin\left(\frac{\pi\beta l}{K}\right) \right] \rho_{kl}, \end{aligned} \quad (\text{G.11})$$

where we have twice used $\cos(A-B) = \cos A \cos B + \sin A \sin B$. When we multiply out the two big square brackets, we get a sum of four terms of the type

$$\begin{aligned} \sum_{\alpha\beta} S_C^{\alpha\beta} \sum_{kl} \cos\left(\frac{\pi\alpha i}{K}\right) \cos\left(\frac{\pi\alpha k}{K}\right) \sin\left(\frac{\pi\beta j}{K}\right) \sin\left(\frac{\pi\beta l}{K}\right) \rho_{kl} \\ = \sum_{\alpha\beta} S_C^{\alpha\beta} \cos\left(\frac{\pi\alpha i}{K}\right) \sin\left(\frac{\pi\beta j}{K}\right) \rho^{\alpha\beta}(\text{CS}), \end{aligned} \quad (\text{G.12})$$

²The connection between discrete complex Fourier transforms and discrete sine and cosine transforms is established by: (i) taking the period of the data x_k to be $2K$; (ii) breaking the sum $\hat{x}_m = \sum_{k=0}^{2K-1} x_k e^{-\pi i k m / K}$ into a sum from 0 to $K-1$ and another from K to $2K-1$ and then using the periodicity of the x_k to express the latter as a sum from $-K$ to -1 ; (iii) writing $x_k = x_k^+ + x_k^-$, where $x_k^\pm \equiv \frac{1}{2}(x_k \pm x_{-k})$, and then combining terms associated with k and $-k$ to produce

$$\hat{x}_m = x_0 + (-1)^m x_{-K} + 2 \sum_{k=1}^{K-1} \left[x_k^+ \cos\left(\frac{\pi k m}{K}\right) - i x_k^- \sin\left(\frac{\pi k m}{K}\right) \right].$$

where

$$\rho^{\alpha\beta}(\text{CS}) \equiv \sum_{kl} \cos\left(\frac{\pi\alpha k}{K}\right) \sin\left(\frac{\pi\beta l}{K}\right) \rho_{kl} \quad (\text{G.13})$$

is ρ_{kl} with its first index cosine transformed and its second index sine transformed. When these and the analogous results are inserted into (G.11), we find

$$\begin{aligned} \Phi_{ij} = \sum_{\alpha\beta} S_C^{\alpha\beta} \left[\rho^{\alpha\beta}(\text{CC}) \cos\left(\frac{\pi\alpha i}{K}\right) \cos\left(\frac{\pi\beta j}{K}\right) + \rho^{\alpha\beta}(\text{CS}) \cos\left(\frac{\pi\alpha i}{K}\right) \sin\left(\frac{\pi\beta j}{K}\right) \right. \\ \left. + \rho^{\alpha\beta}(\text{SC}) \sin\left(\frac{\pi\alpha i}{K}\right) \cos\left(\frac{\pi\beta j}{K}\right) + \rho^{\alpha\beta}(\text{SS}) \sin\left(\frac{\pi\alpha i}{K}\right) \sin\left(\frac{\pi\beta j}{K}\right) \right]. \end{aligned} \quad (\text{G.14})$$

Thus, Φ_{ij} is the sum of four contributions, one from each of the four sine or cosine transforms of ρ_{ij} . The generalization to three dimensions is clear: in three dimensions there are eight possible transforms, $\rho^{\alpha\beta\gamma}(\text{CSS})$ etc., and each one will give rise to a contribution to Φ_{ijk} . The necessity of calculating eight sine/cosine transforms of ρ in order to use the convolution theorem to calculate the potential corresponds to the doubling of the range of each index mentioned below equation (2.237) in the context of complex DFTs.

Box 2.5 describes a technique for obtaining the gravitational potential of a general mass distribution without convolving the softening kernel with the full density distribution. Instead a potential is calculated from the triple sine transform of ρ , and a correction to this potential made by convolving the softening kernel with a distribution of masses that is confined to the walls of the box. Here we explain how the hollowness of the mass distribution dramatically simplifies the convolution.

We break the boundary points of the grid into six isolated pieces, being the top and bottom, front and back, etc., of the grid's bounding box. We ensure that each of these bounding planes is a complete $(K+1) \times (K+1)$ plane of grid points by duplicating nodes of the grid that lie along edges of the grid, and splitting into three copies the nodes that lie at the corners of the grid. We divide any mass assigned to these replicated nodes in equal parts to the copies. Now the triple cosine transform of a density distribution that is confined to the top and bottom planes, $k=0$ and $k=K$, is

$$\begin{aligned} \rho_{\text{top/bot}}^{\alpha\beta\gamma}(\text{CCC}) &= \sum_{ij} (\rho_{ij0} + (-1)^\gamma \rho_{ijK}) \cos\left(\frac{\pi\alpha i}{K}\right) \cos\left(\frac{\pi\beta j}{K}\right), \\ &\equiv \rho_0^{\alpha\beta\cdot}(\text{CC}) + (-1)^\gamma \rho_K^{\alpha\beta\cdot}(\text{CC}), \end{aligned} \quad (\text{G.15})$$

where the second line introduces a convenient notation for keeping track of the objects generated. The key point is that the right side of equation (G.15) is just a two-dimensional cosine transform, which is inexpensive to calculate. The other two pairs of bounding planes generate similar contributions to $\rho^{\alpha\beta\gamma}(\text{CCC})$, so for a hollow mass distribution $\rho^{\alpha\beta\gamma}(\text{CCC})$ may be expressed as a sum of two-dimensional cosine transforms. The other transforms required, such as $\rho^{\alpha\beta\gamma}(\text{CCS})$, $\rho^{\alpha\beta\gamma}(\text{CSS})$, and $\rho^{\alpha\beta\gamma}(\text{SSS})$ are even easier to evaluate because, for example, $\sum_k \rho_{ijk} \sin(\pi\gamma k/K)$ vanishes for a mass distribution that is confined to the top and bottom bounding surfaces.

Appendix H: The Antonov–Lebovitz theorem

This theorem states that all non-radial modes of a barotropic star with $dp(\rho)/d\rho > 0$ are stable. The original proofs are by Antonov (1962b) and Lebovitz (1965), but here we follow Aly & Pérez (1992).

To prove the theorem, we must show that $W[\rho_1]$ in equation (5.120) is never negative. We begin by rewriting the second term in using Poisson’s equation (5.115), so

$$W[\rho_1] = \int d^3\mathbf{x} \left| \frac{d\Phi}{d\rho} \right|_0 \rho_1^2(\mathbf{x}) + \int d^3\mathbf{x} \rho_1(\mathbf{x}) \Phi_1(\mathbf{x}). \quad (\text{H.1})$$

Now we apply Schwarz’s inequality (B.75), setting $A = |d\Phi/d\rho|_0^{1/2} \rho_1$ and $B = \Phi_1/|d\Phi/d\rho|_0^{1/2}$. We have

$$W[\rho_1] = \int d^3\mathbf{x} (A^2 + AB) \geq \frac{X^2}{Y} + X = \frac{X}{Y} (X + Y), \quad (\text{H.2})$$

where

$$X \equiv \int d^3\mathbf{x} AB = \int d^3\mathbf{x} \rho_1 \Phi_1 \quad ; \quad Y \equiv \int d^3\mathbf{x} B^2 = \int d^3\mathbf{x} \left| \frac{d\rho}{d\Phi} \right|_0 \Phi_1^2. \quad (\text{H.3})$$

As shown in Box 5.2, any non-radial mode can be chosen to be the real part of a function having angular dependence proportional to $Y_l^m(\theta, \phi)$ with $l \geq 1$. Thus if $\Phi_1(\mathbf{x})$ is the potential of a mode, we may write

$$\Phi_1(\mathbf{x}) = \text{Re}[\Phi_c(\mathbf{x})], \quad \text{where} \quad \Phi_c(\mathbf{x}) \equiv \frac{d\Phi_0}{dr} s_{lm}(r) Y_l^m(\theta, \phi); \quad (\text{H.4})$$

the common factor $d\Phi_0/dr$ has been introduced to simplify later equations. Using equation (C.49), Poisson’s equation can be written

$$\rho_1(\mathbf{x}) = \text{Re}[\rho_c(\mathbf{x})] \quad (\text{H.5a})$$

where

$$\rho_c(\mathbf{x}) = \frac{1}{4\pi G} \nabla^2 \Phi_c = \frac{1}{4\pi G} Y_l^m(\theta, \phi) \left(\frac{1}{r^2} \frac{d}{dr} r^2 \frac{d}{dr} s_{lm} \frac{d\Phi_0}{dr} - \frac{l(l+1)}{r^2} s_{lm} \frac{d\Phi_0}{dr} \right). \quad (\text{H.5b})$$

Then $X = \int d^3\mathbf{x} \rho_1 \Phi_1 = \frac{1}{4} \int r^2 dr d^2\Omega (\rho_c \Phi_c + \rho_c \Phi_c^*) + \text{CC}$, where “CC” denotes the complex conjugate of the preceding terms in a sum. Using the orthonormality of the spherical harmonics (eq. C.44), we have

$$X = \frac{1}{16\pi G} \int_0^\infty dr \frac{d\Phi_0}{dr} s_{lm}^* \left(\frac{d}{dr} r^2 \frac{d}{dr} s_{lm} \frac{d\Phi_0}{dr} - l(l+1) s_{lm} \frac{d\Phi_0}{dr} \right) + \text{CC}. \quad (\text{H.6})$$

We now carry out the derivatives in the first term in square brackets, and integrate the term involving $d^2 s_{lm}/dr^2$ by parts. The boundary terms can be shown to vanish (Φ_1 decays at least as fast as r^{-2} as $r \rightarrow \infty$ because mass conservation dictates that there cannot be any monopole component of Φ_1 ; thus s_{lm} is constant or decaying as $r \rightarrow \infty$). We simplify the result using Poisson’s equation in the form

$$\frac{d^2 \Phi_0}{dr^2} + \frac{2}{r} \frac{d\Phi_0}{dr} = 4\pi G \rho_0, \quad (\text{H.7})$$

which can be differentiated to eliminate the term involving $d^3\Phi_0/dr^3$. After some algebra we find

$$X = \frac{1}{2} \int_0^\infty dr r^2 \frac{d\rho_0}{dr} \frac{d\Phi_0}{dr} |s_{lm}|^2 - \frac{1}{8\pi G} \int_0^\infty dr \left(\frac{d\Phi_0}{dr} \right)^2 \left(\left| r \frac{ds_{lm}}{dr} \right|^2 + (l^2 + l - 2) |s_{lm}|^2 \right). \quad (\text{H.8})$$

Similarly, we use equation (H.4) and the orthonormality of the spherical harmonics to evaluate Y in terms of $s_{lm}(r)$,

$$Y = \frac{1}{2} \int_0^\infty dr r^2 \left| \frac{d\rho}{d\Phi} \right|_0 \left(\frac{d\Phi_0}{dr} \right)^2 |s_{lm}|^2 = -\frac{1}{2} \int_0^\infty dr r^2 \frac{d\rho_0}{dr} \frac{d\Phi_0}{dr} |s_{lm}|^2, \quad (\text{H.9})$$

where we have assumed that $(d\rho/d\Phi)_0 < 0$ since $dp/d\rho > 0$ by assumption and $(dp/d\Phi)_0 = -\rho_0 < 0$ by hydrostatic equilibrium (F.12). Combining equations (H.8) and (H.9), we have

$$X + Y = -\frac{1}{8\pi G} \int_0^\infty dr \left(\frac{d\Phi_0}{dr} \right)^2 \left(\left| r \frac{ds_{lm}}{dr} \right|^2 + (l^2 + l - 2) |s_{lm}|^2 \right). \quad (\text{H.10})$$

Non-radial modes have $l \geq 1$, and for these $l^2 + l - 2 \geq 0$. Thus $X + Y \leq 0$ for non-radial modes. Moreover $Y \geq 0$ and $X \leq 0$; the latter statement is evident from a comparison of equations (2.17) and (2.18). Thus equation (H.2) implies that $W[\rho_1] \geq 0$ for non-radial modes.

Appendix I: The Doremus–Feix–Baumann theorem

This theorem states that any spherical stellar system with an ergodic equilibrium DF that satisfies $f'_0(E) < 0$ is stable to radial perturbations. Proofs are given by Doremus, Feix, & Baumann (1971) and Kandrup & Sygnet (1985).

We introduce polar coordinates (v, η, ψ) in velocity space (eq. 4.63). The volume element in velocity space is

$$d^3\mathbf{v} = v^2 dv \sin \eta d\eta d\psi = v_t dv_t dv_r d\psi, \quad (\text{I.1})$$

where v_r is the radial velocity and $v_t \equiv v \sin \eta = (v_\theta^2 + v_\phi^2)^{1/2}$ is the tangential velocity. The energy is $E = H_0(r, v) = \frac{1}{2}v^2 + \Phi_0(r)$, where $\Phi_0(r)$ is the potential, and the angular momentum is $L = rv \sin \eta = rv_t$. Since both the equilibrium system and the perturbation are spherically symmetric, neither can depend on the angular variables (θ, ϕ) ; moreover we can integrate the DF over the angle ψ , and work in the variables (r, E, L) .¹ In these variables the volume element in velocity space is

$$d^3\mathbf{v} = 2\pi \frac{dE dL^2}{ry(r, E, L)}, \quad \text{where } y(r, E, L) = \{2r^2[E - \Phi_0(r)] - L^2\}^{1/2} \quad (\text{I.2})$$

¹ In principle, the perturbed DF f_1 could depend on ψ , so long as the density $\rho_1 = \int d^3\mathbf{v} f_1$ is spherically symmetric. It is straightforward to generalize the proof to account for this possibility.

if the quantity in braces is positive, and zero otherwise. Note that a factor of two has been included because there are stars with both positive and negative radial velocities at a given value of (r, E, L) .

The proof is based on the variational principle (5.131). The second term of this expression is twice the potential energy of the perturbation. Since the perturbation is spherically symmetric, we can use equation (2.17) to replace this term by

$$W_2[f_1] = -\frac{1}{G} \int dr \left(r \frac{d\Phi_1}{dr} \right)^2, \quad (\text{I.3})$$

where

$$\begin{aligned} \frac{d\Phi_1}{dr} &= \frac{GM_1(r)}{r^2} = \frac{4\pi G}{r^2} \int_0^r dr' r'^2 \int d^3\mathbf{v}' f_1(r', \mathbf{v}') \\ &= \frac{8\pi^2 G}{r^2} \int_0^r dr' r' \int \frac{dE dL^2}{y(r', E, L)} f_1(r', \mathbf{v}'). \end{aligned} \quad (\text{I.4})$$

At this point we change variables from $f_1(r, \mathbf{v})$ to $g(r, E, L)$, where

$$f_1(r, \mathbf{v}) \equiv \frac{f'_0(E)y(r, E, L)}{r} \frac{\partial}{\partial r} [y(r, E, L)g(r, E, L)]. \quad (\text{I.5})$$

Substituting this expression into (I.3), we have

$$\begin{aligned} W_2[f_1] &= -64\pi^4 G \int \frac{dr}{r^2} \left[\int_0^r dr' \int dE dL^2 f'_0(E) \frac{\partial}{\partial r'} y(r', E, L) g(r', E, L) \right]^2 \\ &= -64\pi^4 G \int \frac{dr}{r^2} \left[\int dE dL^2 f'_0(E) y(r, E, L) g(r, E, L) \right]^2; \end{aligned} \quad (\text{I.6})$$

the contribution to the integral from its lower limit $r' = 0$ is zero since $y(0, E, L) = 0$.

Using equations (I.2) and (I.5), and recalling that $f'_0(E) < 0$, the first term of equation (5.131) can be written as

$$W_1[f_1] = -8\pi^2 \int dr dE dL^2 \frac{f'_0(E)y(r, E, L)}{r} \left[\frac{\partial(yg)}{\partial r} \right]^2. \quad (\text{I.7})$$

We now apply Schwarz's inequality (B.75) to the integral over $dE dL^2$ in equation (I.6), with $A = (-f'_0 y)^{1/2}$ and $B = Ag$:

$$\begin{aligned} W_2[f_1] &\geq -64\pi^4 G \int \frac{dr}{r^2} \int dE dL^2 f'_0(E) y(r, E, L) g^2(r, E, L) \\ &\quad \times \int dE_1 dL_1^2 f'_0(E_1) y(r, E_1, L_1). \end{aligned} \quad (\text{I.8})$$

To evaluate the last integral we integrate by parts with respect to E_1 :

$$\int dE_1 dL_1^2 f'_0(E_1) y(r, E_1, L_1) = - \int dE_1 dL_1^2 f_0(E_1) \frac{\partial y}{\partial E_1}; \quad (\text{I.9})$$

the boundary terms vanish because $y = 0$ at the minimum energy, while $f_0(E) = 0$ at the maximum energy. Now $\partial y / \partial E = r^2 / y$, so using equation (I.2) we have

$$\begin{aligned} \int dE_1 dL_1^2 f'_0(E_1) y(r, E_1, L_1) &= -r^2 \int \frac{dE_1 dL_1^2}{y(r, E_1, L_1)} f_0(E_1) \\ &= -\frac{r^3}{2\pi} \int d^3\mathbf{v}_1 f_0\left(\frac{1}{2}v_1^2 + \Phi_0\right) = -\frac{r^3}{2\pi} \rho_0(r), \end{aligned} \quad (\text{I.10})$$

where $\rho_0(r)$ is the density of the equilibrium system. With this result, equations (I.7) and (I.8) simplify to

$$W[f_1] = W_1[f_1] + W_2[f_1] \geq -8\pi^2 \int dE dL^2 f'_0(E) \times \int \frac{dr y(r, E, L)}{r} \left\{ \left[\frac{\partial(yg)}{\partial r} \right]^2 - 4\pi G \rho_0(r) r^2 g^2 \right\}. \quad (\text{I.11})$$

We integrate the component in square brackets by parts. The boundary terms vanish, since $y(r, E, L) = 0$ at $r = 0$ or as $r \rightarrow \infty$. The remaining integral involves several terms, one of which contains $\partial^2 g / \partial r^2$; we integrate this by parts a second time, and the boundary terms vanish for the same reason. The result is

$$W[f_1] \geq -8\pi^2 \int dE dL^2 f'_0(E) \int \frac{dr y(r, E, L)}{r} \times \left\{ \left(y \frac{\partial g}{\partial r} \right)^2 - g^2 \left[\left(\frac{\partial y}{\partial r} \right)^2 + y \frac{\partial^2 y}{\partial r^2} - \frac{y}{r} \frac{\partial y}{\partial r} \right] - 4\pi G \rho_0(r) r^2 g^2 \right\}. \quad (\text{I.12})$$

Using the definition of $y(r, E, L)$, equation (I.2), we have

$$\frac{1}{r} \frac{\partial}{\partial r} \frac{1}{r} \frac{\partial y^2}{\partial r} = -2 \frac{d^2 \Phi_0}{dr^2} - \frac{6}{r} \frac{d\Phi_0}{dr}. \quad (\text{I.13})$$

Using Poisson's equation in the form (H.7) this can be rewritten as

$$\left(\frac{\partial y}{\partial r} \right)^2 + y \frac{\partial^2 y}{\partial r^2} - \frac{y}{r} \frac{\partial y}{\partial r} = -4\pi G \rho_0(r) r^2 - r \frac{d\Phi_0}{dr}, \quad (\text{I.14})$$

so (I.12) simplifies to

$$W[f_1] \geq -8\pi^2 \int dE dL^2 f'_0(E) \int \frac{dr y(r, E, L)}{r} \left[\left(y \frac{\partial g}{\partial r} \right)^2 + g^2 r \frac{d\Phi_0}{dr} \right], \quad (\text{I.15})$$

which is non-negative since $f'_0(E) < 0$ and $d\Phi_0/dr > 0$. Thus the system is stable.

Appendix J: Angular-momentum transport in disks

We have shown in §6.1.5 that the gravitational torques exerted by a trailing spiral pattern transport angular momentum outward in a disk galaxy. In this appendix we present a more complete description of angular-momentum transport in fluid and stellar disks.

J.1 Transport in fluid and stellar systems

The angular momentum per unit volume in a fluid¹ is $\tilde{\mathbf{l}} = \rho \mathbf{x} \times \mathbf{v}$. Its rate of change can be written in Cartesian coordinates as

$$\frac{\partial \tilde{l}_i}{\partial t} = \epsilon_{ijk} x_j \frac{\partial}{\partial t} (\rho v_k), \quad (\text{J.1})$$

¹ The tilde on \mathbf{l} is a flag that in this appendix, the angular momentum and energy of a fluid element of mass m are defined as $m \mathbf{x} \times \mathbf{v}$ and $m(\frac{1}{2}v^2 + \Phi)$, while in most other sections of this book, "angular momentum" and "energy" are shorthand for specific angular momentum (angular momentum per unit mass) and specific energy, $\mathbf{x} \times \mathbf{v}$ and $\frac{1}{2}v^2 + \Phi$.

where the permutation tensor ϵ_{ijk} is defined after equation (B.7) and we have used the summation convention defined just above that equation. Using equation (F.14) this can be rewritten as

$$\frac{\partial \tilde{l}_i}{\partial t} = -\epsilon_{ijk} x_j \left(\frac{\partial \Pi_{km}}{\partial x_m} + \rho \frac{\partial \Phi}{\partial x_k} \right) = -\frac{\partial \Lambda_{im}}{\partial x_m} - \rho (\mathbf{x} \times \nabla \Phi)_i, \quad (\text{J.2})$$

where the momentum flux tensor $\Pi_{km} = \rho v_k v_m + p \delta_{km}$ and

$$\Lambda_{im} = \epsilon_{ijk} x_j \Pi_{km}. \quad (\text{J.3})$$

The second term in equation (J.2) is simply the gravitational torque per unit volume exerted on the fluid. The first arises from the flow or current of angular momentum; specifically, the outward flow of the angular momentum component \tilde{L}_i through a unit area with outward normal parallel to the unit vector $\hat{\mathbf{e}}_m$ is Λ_{im} , and thus Λ_{im} is known as the **angular-momentum flux tensor**. Note that $\Lambda_{im} \neq \Lambda_{mi}$.

The analogous equation for a stellar system can be written down by inspection. The derivation of the momentum flux tensor in Appendix F relies on the equation of continuity and Euler's equation. The analogs to these equations for a stellar system are the Jeans equations (4.204) and (4.209); the only changes are that the velocity v_i is replaced by the mean velocity \bar{v}_i , and the pressure is replaced by the product of the density and the velocity-dispersion tensor, $\rho \sigma_{ij}^2$. Thus the momentum flux tensor for a stellar system must be

$$\Pi_{km} = \rho \bar{v}_k \bar{v}_m + \rho \sigma_{km}^2 = \rho \overline{v_k v_m}; \quad (\text{J.4})$$

the second expression follows from the first by equation (4.26). The angular-momentum flux tensor for a stellar system is then given by equation (J.3), if the momentum flux tensor is given by equation (J.4).

J.2 Transport in a disk with stationary spiral structure

We now specialize to a razor-thin disk in the $z = 0$ plane. We assume that the disk has a stationary spiral pattern with pattern speed Ω_p ; that is, the surface density $\Sigma(R, \phi, t)$ and all other physical variables depend on the azimuthal angle ϕ and time t only in the combination $\phi - \Omega_p t$. We allow for sources and sinks of angular momentum, such as torques from an external potential, but assume that there are no sources or sinks of mass. We shall focus on the flow of angular momentum through the cylinder $R = R_0 = \text{constant}$.

The mass inside R_0 is constant, since the surface density is constant in the frame rotating at the pattern speed. Thus the net flow of mass through R_0 must be zero; that is, $C_M(R_0) = 0$ where

$$C_M(R) = R \int_0^{2\pi} d\phi \Sigma v_R \quad (\text{J.5})$$

is the **mass current**, and v_R is the radial velocity.

The rate of change of the angular momentum \tilde{L} inside R_0 is given by integrating equation (J.2) over this volume:

$$\frac{d\tilde{L}}{dt} = \int_{R < R_0} d^3 \mathbf{x} \frac{\partial \tilde{l}_z}{\partial t} = - \int_{R < R_0} d^3 \mathbf{x} \left(\frac{\partial \Lambda_{zm}}{\partial x_m} + \epsilon_{zjk} \rho x_j \frac{\partial \Phi}{\partial x_k} \right). \quad (\text{J.6})$$

For an arbitrary vector field \mathbf{F} , the divergence theorem (B.43) implies that

$$\int_{R < R_0} d^3\mathbf{x} \nabla \cdot \mathbf{F} = \oint_S d^2\mathbf{S} \cdot \mathbf{F} = R_0 \int_0^{2\pi} d\phi \int_{-\infty}^{\infty} dz F_R, \quad (\text{J.7})$$

where $F_R = \mathbf{F} \cdot \hat{\mathbf{e}}_R$ is the component of \mathbf{F} in the radial direction, and S is the surface of the cylinder $R = R_0$. We can apply this result to the tensor Λ_{zm} by treating the three components Λ_{zm} , $m = 1, 2, 3$ as a vector; thus

$$\int_{R < R_0} d^3\mathbf{x} \frac{\partial \Lambda_{zm}}{\partial x_m} = R_0 \int_0^{2\pi} d\phi \int_{-\infty}^{\infty} dz \Lambda_{zR}. \quad (\text{J.8})$$

We now write $\rho(R, \phi, z, t) = \Sigma(R, z, t)\delta(z)$ and integrate over z to obtain

$$\frac{d\tilde{L}}{dt} = -R_0^2 \int_0^{2\pi} d\phi \Sigma v_\phi v_R - \int_0^{R_0} dR R \int_0^{2\pi} d\phi \Sigma \frac{\partial \Phi}{\partial \phi}. \quad (\text{J.9})$$

This result is for a fluid; for a stellar system we simply replace $v_\phi v_R$ by $\overline{v_\phi v_R}$.

The first integral on the right side of this equation arises from advective or Reynolds transport of angular momentum due to the bulk motion of the fluid (see discussion following eq. F.70),² and leads to a flow of angular momentum through radius R_0 or **advective angular-momentum current**

$$C_A(R_0) \equiv R_0^2 \int_0^{2\pi} d\phi \Sigma v_\phi v_R. \quad (\text{J.10})$$

The second term on the right side of (J.9) represents the gravitational torque exerted on the fluid inside radius R_0 . The gravitational potential exerting this torque can be written $\Phi \equiv \Phi_{\text{ext}} + \Phi_{\text{in}} + \Phi_{\text{out}}$, where Φ_{ext} is the potential from sources external to the disk, Φ_{in} is the potential due to the disk mass inside R_0 , and Φ_{out} is due to the disk mass outside R_0 . The torque due to Φ_{in} must vanish, since the disk interior to R_0 cannot exert any net torque on itself. The torque due to external sources is

$$\tilde{N}(R_0) \equiv - \int_0^{R_0} dR R \int_0^{2\pi} d\phi \Sigma \frac{\partial \Phi_{\text{ext}}}{\partial \phi}. \quad (\text{J.11})$$

The torque due to the disk exterior to R_0 adds angular momentum to the disk inside R_0 and removes it from the disk outside R_0 . Hence this torque redistributes of the angular momentum in the disk; the corresponding **gravitational angular-momentum current** through R_0 is

$$C_G(R_0) \equiv \int_0^{R_0} dR R \int_0^{2\pi} d\phi \Sigma \frac{\partial \Phi_{\text{out}}}{\partial \phi}. \quad (\text{J.12})$$

This quantity was already defined and discussed in §6.1.5.

The equation for angular-momentum conservation inside radius R_0 is

$$\frac{d\tilde{L}}{dt} = \tilde{N}(R_0) - C_A(R_0) - C_G(R_0); \quad (\text{J.13})$$

² The term **lorry transport** was used by Lynden-Bell & Kalnajs (1972), who likened the fluid elements to a fleet of lorries or trucks carrying coal. The trucks travel outward full of coal and return empty, so there is an outward flow of coal but no net flow of trucks. Similarly the fluid elements carry angular momentum outward, deposit it near the apocenter of their orbits, and return to acquire more angular momentum near pericenter, leading to an outward flow of angular momentum with no outward flow of mass.

in words, the rate of change of the angular momentum of the disk interior to R_0 equals the external torque on that region minus the advective and gravitational currents of angular momentum out of that region.

J.3 Transport in perturbed axisymmetric disks

We next consider the case of small perturbations from axisymmetry, so

$$\begin{aligned} v_R &= \epsilon v_{R1} + \epsilon^2 v_{R2} + \dots & ; & & v_\phi &= v_{\phi 0} + \epsilon v_{\phi 1} + \epsilon^2 v_{\phi 2} + \dots, \\ \Sigma &= \Sigma_0 + \epsilon \Sigma_1 + \epsilon^2 \Sigma_2 + \dots & ; & & \Phi &= \Phi_0 + \epsilon \Phi_1 + \epsilon^2 \Phi_2 + \dots, \end{aligned} \quad (\text{J.14})$$

where $\epsilon \ll 1$ and $v_{\phi 0}(R)$, $\Sigma_0(R)$ and $\Phi_0(R)$ are independent of angle. Note that the zero-order terms in radial velocity v_R and external potential Φ_{ext} both vanish—the first since the unperturbed disk is in circular rotation, and the second since the external perturbations are assumed to be weak. To first order in the small quantity ϵ , equations (J.5), (J.10), (J.11), and (J.12) become

$$\begin{aligned} C_{M1}(R_0) &= \epsilon R_0 \int_0^{2\pi} d\phi \Sigma_0 v_{R1}, \\ C_{A1}(R_0) &= \epsilon R_0^2 \int_0^{2\pi} d\phi \Sigma_0 v_{\phi 0} v_{R1}, \\ \tilde{N}_1(R_0) &= -\epsilon \int_0^{R_0} dR R \int_0^{2\pi} d\phi \Sigma_0 \frac{\partial \Phi_{\text{ext},1}}{\partial \phi}, \\ C_{G1}(R_0) &= \epsilon \int_0^{R_0} dR R \int_0^{2\pi} d\phi \Sigma_0 \frac{\partial \Phi_{\text{out},1}}{\partial \phi}. \end{aligned} \quad (\text{J.15})$$

The mass current $C_{M1}(R_0)$ must vanish, since the mass current is zero to all orders in a stationary disk. Since C_{A1} is just $R_0 v_{\phi 0}$ times C_{M1} , C_{A1} must also vanish. Since $\Sigma_0(R)$ is independent of the angle ϕ , the integrals for $\tilde{N}_1(R_0)$ and $C_{G1}(R_0)$ must also vanish. Thus *the angular-momentum transport in an axisymmetric disk due to external torques and non-axisymmetric density waves of strength ϵ is of order ϵ^2* . This result is closely related to the finding that the energy and momentum flux in a sound wave is second-order in the wave amplitude (Appendix F.3.1).

To quadratic order in ϵ we have

$$\begin{aligned} C_{M2}(R_0) &= \epsilon^2 R_0 \int_0^{2\pi} d\phi (\Sigma_1 v_{R1} + \Sigma_0 v_{R2}), \\ C_{A2}(R_0) &= \epsilon^2 R_0^2 \int_0^{2\pi} d\phi (\Sigma_1 v_{\phi 0} v_{R1} + \Sigma_0 v_{\phi 1} v_{R1} + \Sigma_0 v_{\phi 0} v_{R2}), \\ &= \epsilon^2 R_0^2 \int_0^{2\pi} d\phi \Sigma_0 v_{\phi 1} v_{R1} + R_0 v_{\phi 0} C_{M2}(R_0) = \epsilon^2 R_0^2 \Sigma_0 \int_0^{2\pi} d\phi v_{\phi 1} v_{R1}, \\ \tilde{N}_2(R_0) &= -\epsilon^2 \int_0^{R_0} dR R \int_0^{2\pi} d\phi \Sigma_1 \frac{\partial \Phi_{\text{ext},1}}{\partial \phi}, \\ C_{G2}(R_0) &= \epsilon^2 \int_0^{R_0} dR R \int_0^{2\pi} d\phi \Sigma_1 \frac{\partial \Phi_{\text{out},1}}{\partial \phi}. \end{aligned} \quad (\text{J.16})$$

Here the last expression for $C_A(R_0)$ follows because $C_M(R_0)$ is zero to all orders. For stellar systems the advective current is obtained by replacing $v_{\phi 1} v_{R1}$ by $\overline{v_{\phi 1} v_{R1}}$.

Following equation (6.42), we can write the perturbed quantities as a sum over azimuthal wavenumber $m \geq 0$ of terms of the form

$$v_{R1} = \text{Re}[v_{Ra}(R)e^{im(\phi - \Omega_p t)}] = \frac{1}{2}v_{Ra}(R)e^{im(\phi - \Omega_p t)} + \frac{1}{2}v_{Ra}^*(R)e^{-im(\phi - \Omega_p t)}, \quad (\text{J.17})$$

etc. Axisymmetric perturbations ($m = 0$) cannot transport angular momentum, so we may restrict the sum to terms with $m > 0$. Then equation (J.13) can be integrated over ϕ to yield a sum over m of terms of the form

$$\begin{aligned} \frac{d\tilde{L}}{dt} &= \tilde{N}_2(R_0) - C_{A2}(R_0) - C_{G2}(R_0) \\ &= -\frac{1}{2}\pi\epsilon^2 R_0^2 \Sigma_0 v_{\phi a} v_{Ra}^* + \frac{1}{2}\pi m i \epsilon^2 \int_0^{R_0} dR R \Sigma_a (\Phi_{\text{ext},a}^* + \Phi_{\text{out},a}^*) + \text{CC}, \end{aligned} \quad (\text{J.18})$$

where ‘‘CC’’ stands for the complex conjugate of the preceding terms.

Using equations (6.43), the advective current can be written, after some algebra, as

$$C_{A2}(R_0) = \frac{\pi\epsilon^2 R_0 \Sigma_0 i m}{2\Delta} \left[(\Phi_a + h_a)^* \frac{d}{dR} (\Phi_a + h_a) \right] + \text{CC}, \quad (\text{J.19})$$

where $\Delta = \kappa^2 - m^2(\Omega - \Omega_p)^2$.

J.4 Transport in the WKB approximation

The expressions we have derived for the angular-momentum currents can be evaluated in the WKB approximation analyzed in §6.2. In this approximation d/dR can be replaced by ik in equation (J.19), where k is the radial wavenumber, so the advective current simplifies to

$$C_{A2}(R_0) = -\frac{\pi m \epsilon^2 R_0 \Sigma_0 k}{\Delta} |\Phi_a + h_a|^2. \quad (\text{J.20})$$

Using equations (6.30) and (6.45), we can express the perturbed enthalpy in terms of the perturbed potential,

$$h_a = v_s^2 \frac{\Sigma_a}{\Sigma_0} = -\frac{v_s^2 |k|}{2\pi G \Sigma_0} \Phi_a, \quad (\text{J.21})$$

so the advective angular-momentum current is

$$C_{A2}(R_0) = -\frac{\pi m \epsilon^2 R_0 \Sigma_0 k}{\Delta} \left(1 - \frac{v_s^2 |k|}{2\pi G \Sigma_0} \right)^2 |\Phi_a|^2. \quad (\text{J.22})$$

The gravitational current has already been evaluated in equation (6.21); adjusting the notation to the present usage we have

$$C_{G2}(R_0) = \text{sgn}(k) \frac{m \epsilon^2 R_0 |\Phi_a|^2}{4G}; \quad (\text{J.23})$$

so the total angular-momentum current from a tightly wrapped wave in a fluid disk is (Goldreich & Tremaine 1979)

$$C_{\tilde{L}}(R_0) = C_{A2}(R_0) + C_{G2}(R_0) = \text{sgn}(k) \frac{m \epsilon^2 R_0 |\Phi_a|^2}{4G} \left(\frac{v_s^2 |k|}{\pi G \Sigma_0} - 1 \right), \quad (\text{J.24})$$

where Δ has been eliminated using the WKB dispersion relation (6.55). The WKB analysis of disk dynamics in §6.2.2 shows that the total angular-momentum current is conserved, as it must be (eq. 6.56).

For a stellar disk the analogous formula is (Toomre 1969)

$$C_{\tilde{L}}(R_0) = -\text{sgn}(k) \frac{m\epsilon^2 R_0 |\Phi_a|^2}{4G} \left(1 + 2 \frac{\partial \ln \mathcal{F}(s, \chi)}{\partial \ln \chi} \right), \quad (\text{J.25})$$

where $s = m(\Omega_p - \Omega)/\kappa$, $\chi = \sigma_R^2 k^2 / \kappa^2$, and the reduction factor $\mathcal{F}(s, \chi)$ is defined in equation (6.63).

Appendix K: Derivation of the reduction factor

Our goal in this appendix is to derive the reduction factor \mathcal{F} , the factor by which the response of a stellar disk to an imposed potential is reduced below that of a cold disk (eq. 6.58). We assume that the potential is tightly wound, that the orbits in the disk can be described by the epicycle approximation, and that the disk is razor-thin and flat, with the DF having the Schwarzschild form (6.62).

As usual we assume that the perturbations are small, so we can linearize the equations of motion. The linearized collisionless Boltzmann equation (5.15) reads

$$f_1(\mathbf{x}, \mathbf{v}, t) = - \int_{-\infty}^t dt' [f_0, \Phi_1]_{\mathbf{x}', \mathbf{v}', t'} = \int_{-\infty}^t dt' \frac{\partial f_0}{\partial \mathbf{v}'}(\mathbf{x}', \mathbf{v}') \cdot \nabla' \Phi_1(\mathbf{x}', t'); \quad (\text{K.1})$$

here $f_0(\mathbf{x}, \mathbf{v})$ and $f_1(\mathbf{x}, \mathbf{v}, t)$ are the equilibrium and perturbed DFs, $\Phi_1(\mathbf{x}, t)$ is the perturbed gravitational potential, and $[\cdot, \cdot]$ is the Poisson bracket, evaluated along the unperturbed orbit $(\mathbf{x}', \mathbf{v}') \equiv (\mathbf{x}[t'], \mathbf{v}[t'])$ that arrives at \mathbf{x}, \mathbf{v} at time t . In writing equation (K.1) we have assumed that $f_1 \rightarrow 0$ as $t \rightarrow -\infty$, i.e., the perturbation was small in the distant past.

We write the potential perturbation in the disk plane in the form (cf. eqs. 6.42 and 6.51)

$$\Phi_1(\mathbf{x}, t) = \Phi_a(R) e^{i(m\phi - \omega t)} = F(R) e^{i(\int^R k \, dR + m\phi - \omega t)}, \quad (\text{K.2})$$

where as usual only the real part is physical. The component of $\nabla \Phi_1$ in the disk plane is

$$\nabla \Phi_1 = \left[\hat{\mathbf{e}}_R \left(\frac{dF}{dR} + ikF \right) + \hat{\mathbf{e}}_\phi \frac{imF}{R} \right] e^{i(\int^R k \, dR + m\phi - \omega t)}. \quad (\text{K.3})$$

Since the potential is tightly wound, $|kR| \gg 1$, we keep only the term proportional to k :

$$\nabla \Phi_1 \simeq \hat{\mathbf{e}}_R ikF e^{i(\int^R k \, dR + m\phi - \omega t)}. \quad (\text{K.4})$$

Substituting this result and the Schwarzschild DF (6.62) into equation (K.1) yields

$$f_1(R, \phi, v_R, v_\phi, t) = -\frac{i}{2\pi} \int_{-\infty}^t dt' \left(\frac{\gamma k F \Sigma}{\sigma_R^4} \right)_{R'} e^{i(\int^{R'} k'' \, dR'' + m\phi' - \omega t') - w' v'_{R'}}, \quad (\text{K.5})$$

where the unperturbed trajectory at time t' is $(\mathbf{x}', \mathbf{v}') = (R', \phi', v'_{R'}, v'_\phi)$,

$$w' = \frac{v'_{R'}{}^2 + \gamma^2(R') [v'_\phi - v_c(R')]^2}{2\sigma_R^2(R')}, \quad (\text{K.6})$$

and $\gamma(R) = 2\Omega(R)/\kappa(R)$.

The epicycle orbit equations (3.75), (3.91), and (3.97) yield the relations

$$\begin{aligned} R' &= R_g + X \cos(\kappa t' + \psi), \\ v'_R &= -\kappa X \sin(\kappa t' + \psi), \\ v'_\phi - v_c(R') &= -\frac{\kappa}{\gamma} X \cos(\kappa t' + \psi), \end{aligned} \quad (\text{K.7})$$

where the epicycle amplitude X , guiding-center radius R_g , and phase ψ are determined by the boundary conditions that $R' = R$, $v'_R = v_R$, and $v'_\phi = v_\phi$ at $t' = t$. Straightforward manipulation of equations (K.7) yields

$$R' = R - \frac{\gamma}{\kappa}(v_\phi - v_c)(\cos \tau - 1) + \frac{v_R}{\kappa} \sin \tau \quad ; \quad v'_R = v_R \cos \tau + \gamma(v_\phi - v_c) \sin \tau, \quad (\text{K.8})$$

where $\tau = \kappa(t' - t)$ and v_ϕ , v_R and $v_c(R)$ are all evaluated at t .

We are assuming that the epicycle amplitude X is small, so X and $|R' - R|$ are always much less than R . We are also assuming that the spiral wave is tightly wound, so $|kR| \gg 1$. The epicycle amplitude can be either larger or smaller than the radial spacing between arms, so $|kX|$ can be either large or small compared to unity. With this ordering, equation (K.5) can be simplified considerably:

- (i) Since $|R' - R| \ll R$, to a good approximation $k(R')$, $F(R')$, $\Sigma(R')$, $\sigma_R(R')$, and $\gamma(R')$ can all be replaced by their values at R and taken outside the integral.
- (ii) At the same level of approximation we can write

$$e^{i \int^{R'} k'' dR''} \simeq e^{i[\int^R k'' dR'' + k(R)(R' - R)]}. \quad (\text{K.9})$$

- (iii) Since the azimuthal angle ϕ' enters only through the factor $\exp(im\phi')$, and the epicycle amplitude is small, we may neglect the epicyclic motion in ϕ and set $\phi' = \phi + \Omega(t' - t)$. We cannot use the same argument to neglect the epicyclic motion in R since the small factor $(R' - R)$ is multiplied by the large factor k .
- (iv) Equations (K.7) imply that $w' = \kappa^2 X^2 / (2\sigma_R^2)$. Since X is an integral of the motion, w' must therefore be conserved along the trajectory.¹ Thus $w' = w$ and the factor $\exp(-w)$ can be taken out of the integral.

Inserting these approximations into equation (K.5), we find

$$\begin{aligned} f_1(R, \phi, v_R, v_\phi, t) &= -\frac{i\gamma k F \Sigma}{2\pi\kappa\sigma_R^3} e^{i(\int^R k'' dR'' + m\phi - \omega t)} e^{-(u^2 + v^2)/2} \\ &\times \int_{-\infty}^0 d\tau e^{i(k\sigma_R/\kappa)[u \sin \tau - v(\cos \tau - 1)] - is\tau} (u \cos \tau + v \sin \tau). \end{aligned} \quad (\text{K.10})$$

Here we have introduced the notations

$$s = \frac{\omega - m\Omega}{\kappa} \quad ; \quad u = \frac{v_R}{\sigma_R} \quad ; \quad v = \gamma \frac{v_\phi - v_c}{\sigma_R}. \quad (\text{K.11})$$

The mean radial velocity at a given point is

$$\bar{v}_{R1} = \frac{\int d^2\mathbf{v} (f_0 + f_1) v_R}{\int d^2\mathbf{v} (f_0 + f_1)}. \quad (\text{K.12})$$

¹This result also follows from the Jeans theorem: since the unperturbed DF depends on velocity only through w , w must be an integral.

Since $\int d^2\mathbf{v} f_0 v_R = 0$, to first order in the perturbation we have

$$\bar{v}_{R1} = \frac{\int d^2\mathbf{v} f_1 v_R}{\int d^2\mathbf{v} f_0} = \frac{\sigma_R^3 \int dudv f_1 u}{\gamma \Sigma}. \quad (\text{K.13})$$

Writing $\bar{v}_{R1} = \bar{v}_{Ra} \exp[i(m\phi - \omega t)]$ (eq. 6.42) and using equation (K.2) and the definition of the reduction factor \mathcal{F} (eq. 6.58), we find

$$\begin{aligned} \mathcal{F} &= i \frac{1-s^2}{2\pi s} \int_{-\infty}^{\infty} du u \int_{-\infty}^{\infty} dv e^{-(u^2+v^2)/2} \\ &\quad \times \int_{-\infty}^0 d\tau e^{i(k\sigma_R/\kappa)[u \sin \tau - v(\cos \tau - 1)] - is\tau} (u \cos \tau + v \sin \tau). \end{aligned} \quad (\text{K.14})$$

We now use the following integrals:

$$\begin{aligned} \int_{-\infty}^{\infty} dx e^{-x^2/2} e^{i\mu x} &= \sqrt{2\pi} e^{-\mu^2/2} \quad ; \quad \int_{-\infty}^{\infty} dx e^{-x^2/2} x e^{i\mu x} = i\sqrt{2\pi} \mu e^{-\mu^2/2}, \\ \int_{-\infty}^{\infty} dx e^{-x^2/2} x^2 e^{i\mu x} &= \sqrt{2\pi} (1 - \mu^2) e^{-\mu^2/2}. \end{aligned} \quad (\text{K.15})$$

Evaluating the u and v integrals and replacing τ by $-\tau$, we obtain

$$\mathcal{F}(s, \chi) = i \frac{1-s^2}{s} \int_0^{\infty} d\tau e^{is\tau - \chi(1 - \cos \tau)} (\cos \tau - \chi \sin^2 \tau), \quad (\text{K.16})$$

where $\chi = (k\sigma_R/\kappa)^2$. Now write the integral as a sum of integrals from $\tau = 0$ to 2π , 2π to 4π , etc.:

$$\mathcal{F}(s, \chi) = i \frac{1-s^2}{s} \sum_{n=0}^{\infty} e^{2\pi i n s} \int_0^{2\pi} d\tau e^{is\tau - \chi(1 - \cos \tau)} (\cos \tau - \chi \sin^2 \tau). \quad (\text{K.17})$$

The geometric series in (K.17) has the form $\sum_{n=0}^{\infty} p^n$, where $|p| = |\exp(2\pi i s)| = \exp[-2\pi \text{Im}(s)] < 1$, since $\text{Im}(s) > 0$ (we have assumed that the perturbation vanishes as $t \rightarrow -\infty$). Hence the series is convergent, with sum $1/(1-p)$. Thus

$$\mathcal{F}(s, \chi) = i \frac{1-s^2}{s} \frac{1}{1 - e^{2\pi i s}} \int_0^{2\pi} d\tau e^{is\tau - \chi(1 - \cos \tau)} (\cos \tau - \chi \sin^2 \tau). \quad (\text{K.18})$$

Replacing the variable τ by $\tau + \pi$, we have

$$\mathcal{F}(s, \chi) = i \frac{1-s^2}{s} \frac{e^{\pi i s}}{1 - e^{2\pi i s}} \int_{-\pi}^{\pi} d\tau e^{is\tau - \chi(1 + \cos \tau)} (-\cos \tau - \chi \sin^2 \tau). \quad (\text{K.19})$$

We can write $(1 - e^{2\pi i s})/e^{\pi i s} = -2i \sin \pi s$. Also, if we write $e^{is\tau} = \cos s\tau + i \sin s\tau$, only the cosine term will contribute to the integral since the sine produces an integrand that is odd in τ . Thus

$$\mathcal{F}(s, \chi) = \frac{1-s^2}{s \sin \pi s} \int_0^{\pi} d\tau e^{-\chi(1 + \cos \tau)} \cos s\tau (\cos \tau + \chi \sin^2 \tau). \quad (\text{K.20})$$

The integral can be simplified by writing the second term of the integrand as $\chi e^{-\chi(1 + \cos \tau)} \cos s\tau \sin^2 \tau = \cos s\tau \sin \tau d[e^{-\chi(1 + \cos \tau)}]/d\tau$ and integrating by parts:

$$\mathcal{F}(s, \chi) = \frac{1-s^2}{\sin \pi s} \int_0^{\pi} d\tau e^{-\chi(1 + \cos \tau)} \sin s\tau \sin \tau. \quad (\text{K.21})$$

Yet another integration by parts yields a second form:

$$\mathcal{F}(s, \chi) = \frac{1-s^2}{\chi} \left[1 - \frac{s}{\sin \pi s} \int_0^\pi d\tau e^{-\chi(1+\cos \tau)} \cos s\tau \right]. \quad (\text{K.22})$$

We obtain a third form, involving the modified Bessel function I_n , by using the identities (C.64) and (C.67):

$$\mathcal{F}(s, \chi) = \frac{1-s^2}{\chi} \left[1 - \frac{s}{\sin \pi s} e^{-\chi} \sum_{n=-\infty}^{\infty} (-1)^n I_n(\chi) \int_0^\pi d\tau \cos n\tau \cos s\tau \right]. \quad (\text{K.23})$$

It is simple to show that $\int_0^\pi d\tau \cos n\tau \cos s\tau = (-1)^{n+1} s \sin \pi s / (n^2 - s^2)$; hence

$$\mathcal{F}(s, \chi) = \frac{1-s^2}{\chi} \left[1 + s^2 e^{-\chi} \sum_{n=-\infty}^{\infty} \frac{I_n(\chi)}{n^2 - s^2} \right]. \quad (\text{K.24})$$

This equation can be simplified by rewriting the unit term in the square brackets as $1 = e^{-\chi} \sum_{n=-\infty}^{\infty} I_n(\chi)$, an identity that follows from equation (C.67) by setting $\theta = 0$. After some rearrangement, we obtain a final form,

$$\mathcal{F}(s, \chi) = \frac{2}{\chi} (1-s^2) e^{-\chi} \sum_{n=1}^{\infty} \frac{I_n(\chi)}{1-s^2/n^2}. \quad (\text{K.25})$$

Appendix L: The diffusion coefficients

We consider a subject star of mass m and velocity \mathbf{v} moving through a homogeneous sea of field stars of mass m_a and DF $f_a(\mathbf{v}_a)$; the DF is normalized so that $\int d^3\mathbf{v}_a f_a(\mathbf{v}_a) = n$, the number density of field stars. Due to encounters with the field stars, the subject star suffers a gradually accumulating velocity change $\Delta\mathbf{v}$. Our goal is to calculate the diffusion coefficients $D[\Delta v_i]$ and $D[\Delta v_i \Delta v_j]$, which measure the average velocity changes per unit time,

$$D[\Delta v_i] = \frac{\langle \Delta v_i \rangle}{\Delta t} \quad ; \quad D[\Delta v_i \Delta v_j] = \frac{\langle \Delta v_i \Delta v_j \rangle}{\Delta t}. \quad (\text{L.1})$$

The diffusion coefficients were first derived by Landau (1936) in the context of Coulomb interactions in plasmas, but the derivation here follows Rosenbluth, MacDonald, & Judd (1957).

Suppose that in a particular encounter the subject and field stars have velocities \mathbf{v} and \mathbf{v}_a respectively. The relative velocity is $\mathbf{V} = \mathbf{v} - \mathbf{v}_a$. Its initial value is \mathbf{V}_0 and the change in velocity caused by the encounter is $\Delta\mathbf{V}$. The dynamics of the encounter is described in §3.1d and in particular the relation between $\Delta\mathbf{V}$ and the change in velocity of the subject star $\Delta\mathbf{v}$ is given by equation (3.45), which in the present notation reads

$$\Delta\mathbf{v} = \frac{m_a}{m + m_a} \Delta\mathbf{V}. \quad (\text{L.2})$$

We introduce a coordinate system $\hat{\mathbf{e}}'_1, \hat{\mathbf{e}}'_2, \hat{\mathbf{e}}'_3$, such that $\hat{\mathbf{e}}'_1$ is parallel to \mathbf{V}_0 (see Figure L.1). Thus

$$\mathbf{V}_0 \cdot \hat{\mathbf{e}}'_1 = |\mathbf{V}_0| \equiv V_0 \quad ; \quad \mathbf{V}_0 \cdot \hat{\mathbf{e}}'_2 = \mathbf{V}_0 \cdot \hat{\mathbf{e}}'_3 = 0. \quad (\text{L.3})$$

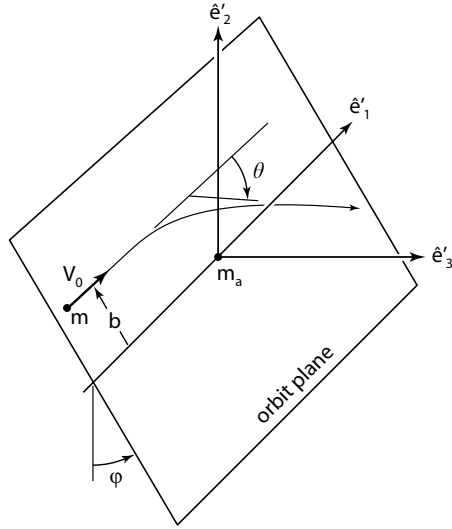


Figure L.1 Geometry of an encounter. The angle θ is the deflection angle θ_{defl} introduced on page 155 and in Figure 3.2.

We denote the angle between the plane of the relative orbit and $\hat{\mathbf{e}}'_2$ by ϕ . The change in \mathbf{v} during the encounter may be written

$$\Delta \mathbf{v} = -\Delta v_{\parallel} \hat{\mathbf{e}}'_1 + \Delta v_{\perp} (-\hat{\mathbf{e}}'_2 \cos \phi + \hat{\mathbf{e}}'_3 \sin \phi), \quad (\text{L.4})$$

where Δv_{\parallel} and Δv_{\perp} are the magnitudes of the components of $\Delta \mathbf{v}$ that are parallel and perpendicular to \mathbf{V}_0 . The signs are chosen to ensure that Δv_{\parallel} and Δv_{\perp} are positive if the interaction between the two particles is attractive.

Using the identity $\Delta \mathbf{v} = \sum_{k=1}^3 (\Delta \mathbf{v} \cdot \hat{\mathbf{e}}'_k) \hat{\mathbf{e}}'_k$ we have

$$\begin{aligned} \Delta v_i &= \sum_{k=1}^3 (\Delta \mathbf{v} \cdot \hat{\mathbf{e}}'_k) (\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_k), \\ \Delta v_i \Delta v_j &= \sum_{k,l=1}^3 (\Delta \mathbf{v} \cdot \hat{\mathbf{e}}'_k) (\Delta \mathbf{v} \cdot \hat{\mathbf{e}}'_l) (\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_k) (\hat{\mathbf{e}}_j \cdot \hat{\mathbf{e}}'_l). \end{aligned} \quad (\text{L.5})$$

Since all angles $0 \leq \phi < 2\pi$ are equally probable, and Δv_{\parallel} and Δv_{\perp} do not depend on ϕ , we can take averages of equations (L.4) and (L.5) over ϕ . Denoting these averages by $\langle \cdot \rangle_{\phi}$, we have $\langle \cos \phi \rangle_{\phi} = \langle \sin \phi \rangle_{\phi} = 0$, $\langle \cos^2 \phi \rangle_{\phi} = \langle \sin^2 \phi \rangle_{\phi} = \frac{1}{2}$. Thus equations (L.4) and (L.5) yield

$$\begin{aligned} \langle \Delta v_i \rangle_{\phi} &= -\Delta v_{\parallel} (\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_1), \\ \langle \Delta v_i \Delta v_j \rangle_{\phi} &= (\Delta v_{\parallel})^2 (\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_1) (\hat{\mathbf{e}}_j \cdot \hat{\mathbf{e}}'_1) \\ &\quad + \frac{1}{2} (\Delta v_{\perp})^2 [(\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_2) (\hat{\mathbf{e}}_j \cdot \hat{\mathbf{e}}'_2) + (\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_3) (\hat{\mathbf{e}}_j \cdot \hat{\mathbf{e}}'_3)]. \end{aligned} \quad (\text{L.6})$$

From equations (3.54) we have

$$\Delta v_{\perp} = \frac{2m_a V_0}{m + m_a} \frac{b/b_{90}}{1 + b^2/b_{90}^2} \quad ; \quad \Delta v_{\parallel} = \frac{2m_a V_0}{m + m_a} \frac{1}{1 + b^2/b_{90}^2}, \quad (\text{L.7})$$

where b is the impact parameter and

$$b_{90} \equiv \frac{G(m + m_a)}{V_0^2} \tag{L.8}$$

is the 90° deflection radius defined in equation (3.51).

We now sum the effects of all the encounters. The number density of field stars in the velocity-space volume $d^3\mathbf{v}_a$ is $f_a(\mathbf{v}_a)d^3\mathbf{v}_a$. The number of encounters in a time Δt with impact parameters between b and $b + db$ is just this density times the volume of an annulus with inner radius b , outer radius $b + db$, and length $V_0\Delta t$, that is,

$$2\pi b db V_0 \Delta t f_a(\mathbf{v}_a) d^3\mathbf{v}_a. \tag{L.9}$$

Thus

$$D[\Delta v_i] = \frac{\langle \Delta v_i \rangle}{\Delta t} = 2\pi \int d^3\mathbf{v}_a db b f_a(\mathbf{v}_a) V_0 \langle \Delta v_i \rangle_\phi, \tag{L.10}$$

with a similar equation for $D[\Delta v_i \Delta v_j]$. We can carry out the integration over the impact parameter b using equations (L.7). The range of integration is from 0 to b_{\max} . The integrals involved are:

$$\begin{aligned} \int_0^{b_{\max}} db b \Delta v_{\parallel} &= \frac{m_a V_0 b_{90}^2}{m + m_a} \ln(1 + \Lambda^2), \\ \int_0^{b_{\max}} db b (\Delta v_{\parallel})^2 &= 2 \left(\frac{m_a V_0 b_{90}}{m + m_a} \right)^2 \left(1 - \frac{1}{1 + \Lambda^2} \right), \\ \int_0^{b_{\max}} db b (\Delta v_{\perp})^2 &= 2 \left(\frac{m_a V_0 b_{90}}{m + m_a} \right)^2 \left(\ln(1 + \Lambda^2) + \frac{1}{1 + \Lambda^2} - 1 \right), \end{aligned} \tag{L.11}$$

where

$$\Lambda \equiv \frac{b_{\max}}{b_{90}}. \tag{L.12}$$

The choice of the appropriate value for b_{\max} is discussed following equation (7.84), where it is argued that b_{\max} is given approximately by the radius of the subject star's orbit. In most applications, Λ is very large, and therefore without loss of accuracy we can discard terms involving $(1 + \Lambda^2)^{-1}$ and replace $\ln(1 + \Lambda^2)$ by $2 \ln \Lambda$ in equation (L.11). Furthermore, although this is a less accurate approximation, we can discard terms of order unity compared to those of order $\ln \Lambda$. In this manner we arrive at a simplified version of equations (L.11), keeping only terms of order $\ln \Lambda$,

$$\begin{aligned} \int_0^{b_{\max}} db b \Delta v_{\parallel} &= 2 \frac{m_a V_0 b_{90}^2}{m + m_a} \ln \Lambda, \\ \int_0^{b_{\max}} db b (\Delta v_{\parallel})^2 &= 0 \quad ; \quad \int_0^{b_{\max}} db b (\Delta v_{\perp})^2 = 4 \left(\frac{m_a V_0 b_{90}}{m + m_a} \right)^2 \ln \Lambda. \end{aligned} \tag{L.13}$$

At this point we investigate whether diffusion tensors of rank higher than two play an important role in gravitational scattering. By analogy with equation (L.6), a diffusion tensor of rank n , $D[\Delta v_{i_1} \cdots \Delta v_{i_n}]$, will involve products of unit vectors with terms such as $(\Delta v_{\parallel})^c (\Delta v_{\perp})^d$ where $c + d = n$. When $b \gg b_{90}$, $\Delta v_{\perp} \propto b^{-1}$ and $\Delta v_{\parallel} \propto b^{-2}$, while for $b \ll b_{90}$, $\Delta v_{\perp} \propto b$ and $\Delta v_{\parallel} \propto \text{constant}$. Therefore $(\Delta v_{\parallel})^c (\Delta v_{\perp})^d \propto b^{-2c-d} = b^{-n-c}$ for $b \gg b_{90}$ and $\propto b^d = b^{n-c}$ for $b \ll b_{90}$. The integral over impact parameter in the analog of equation (L.10) will have the form $\int db b^{1-n-c}$ for $b \gg b_{90}$ and $\int db b^{1+n-c}$ for $b \ll b_{90}$. Since $0 \leq c \leq n$ these integrals

cannot give a divergent factor such as $\ln \Lambda$ for $n > 2$. Thus diffusion tensors of rank three or higher are smaller than the tensors of ranks 1 and 2 by at least $O(\ln \Lambda)^{-1}$, and hence can be dropped in the Fokker–Planck expansion (7.65) (Hénon 1973b). An explicit calculation of the higher-rank tensors is given by Hénon (1960a).

Using equations (L.6), (L.10), and (L.13) we can write the diffusion tensors as

$$\begin{aligned} D[\Delta v_i] &= -4\pi \frac{m_a}{m + m_a} \int d^3 \mathbf{v}_a V_0^2 b_{90}^2 f_a(\mathbf{v}_a) \ln \Lambda (\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_1), \\ D[\Delta v_i \Delta v_j] &= 4\pi \left(\frac{m_a}{m + m_a} \right)^2 \int d^3 \mathbf{v}_a V_0^3 b_{90}^2 f_a(\mathbf{v}_a) \ln \Lambda \\ &\quad \times [(\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_2)(\hat{\mathbf{e}}_j \cdot \hat{\mathbf{e}}'_2) + (\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_3)(\hat{\mathbf{e}}_j \cdot \hat{\mathbf{e}}'_3)]. \end{aligned} \quad (\text{L.14})$$

Since we assume that $\ln \Lambda$ is large, we do not make any significant additional error by replacing the factor V_0 in Λ by some typical stellar speed v_{typ} . Thus

$$\Lambda = \frac{b_{\text{max}} v_{\text{typ}}^2}{G(m + m_a)}, \quad (\text{L.15})$$

which is independent of \mathbf{v}_a so $\ln \Lambda$ may be taken outside the integral. Also, the expressions involving unit vectors can be simplified, since $\sum_{p=1}^3 (\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_p)(\hat{\mathbf{e}}_j \cdot \hat{\mathbf{e}}'_p) = \hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}_j = \delta_{ij}$, and $(\hat{\mathbf{e}}_i \cdot \hat{\mathbf{e}}'_1) = V_{0i}/V_0$. After eliminating b_{90} using equation (L.8), we obtain

$$\begin{aligned} D[\Delta v_i] &= -4\pi G^2 m_a (m + m_a) \ln \Lambda \int d^3 \mathbf{v}_a \frac{f_a(\mathbf{v}_a)}{V_0^3} V_{0i}, \\ D[\Delta v_i \Delta v_j] &= 4\pi G^2 m_a^2 \ln \Lambda \int d^3 \mathbf{v}_a \frac{f_a(\mathbf{v}_a)}{V_0} \left(\delta_{ij} - \frac{V_{0i} V_{0j}}{V_0^2} \right). \end{aligned} \quad (\text{L.16})$$

These equations can be simplified further by noting that $V_0 = [\sum_{i=1}^3 (v_i - v_{ai})^2]^{1/2}$, so

$$\frac{\partial}{\partial v_i} \frac{1}{V_0} = -\frac{V_{0i}}{V_0^3} \quad ; \quad \frac{\partial^2}{\partial v_i \partial v_j} V_0 = \frac{\delta_{ij}}{V_0} - \frac{V_{0i} V_{0j}}{V_0^3}; \quad (\text{L.17})$$

thus we can write

$$\begin{aligned} D[\Delta v_i] &= 4\pi G^2 m_a (m + m_a) \ln \Lambda \frac{\partial}{\partial v_i} h(\mathbf{v}), \\ D[\Delta v_i \Delta v_j] &= 4\pi G^2 m_a^2 \ln \Lambda \frac{\partial^2}{\partial v_i \partial v_j} g(\mathbf{v}), \end{aligned} \quad (\text{L.18})$$

where the **Rosenbluth potentials** are

$$h(\mathbf{v}) \equiv \int d^3 \mathbf{v}_a \frac{f_a(\mathbf{v}_a)}{|\mathbf{v} - \mathbf{v}_a|} \quad ; \quad g(\mathbf{v}) \equiv \int d^3 \mathbf{v}_a f_a(\mathbf{v}_a) |\mathbf{v} - \mathbf{v}_a|. \quad (\text{L.19})$$

To within a proportionality constant, $h(\mathbf{v})$ is simply the gravitational potential at \mathbf{v} generated by a fictitious body in velocity space that has density $f_a(\mathbf{v}_a)$ (cf. eq. 2.3).

When the field star DF is isotropic, so $f_a(\mathbf{v}_a)$ depends only on $v_a = |\mathbf{v}_a|$, more explicit expressions can be obtained. In this case, symmetry dictates that the Rosenbluth potentials depend only on $v = |\mathbf{v}|$. To evaluate $g(v)$ and $h(v)$, we use equation (C.35) to write

$$\frac{1}{|\mathbf{v} - \mathbf{v}_a|} = \sum_{l=0}^{\infty} \frac{v_{<}^l}{v_{>}^{l+1}} P_l(\cos \gamma), \quad (\text{L.20})$$

where $v_<$ and $v_>$ are the smaller and larger of v and v_a , $P_l(x)$ is a Legendre polynomial, and γ is the angle between \mathbf{v} and \mathbf{v}_a . Thus

$$h(v) = 2\pi \sum_{l=0}^{\infty} \int_0^{\infty} dv_a \frac{v_a^2 v_<^l}{v_>^{l+1}} f_a(v_a) \int_0^{\pi} d\gamma \sin \gamma P_l(\cos \gamma). \quad (\text{L.21})$$

Using the relation $\int_{-1}^1 dx P_l(x) = 2\delta_{l0}$ (eq. C.39 with $m = n = 0$) we obtain

$$h(v) = 4\pi \left[\frac{1}{v} \int_0^v dv_a v_a^2 f_a(v_a) + \int_v^{\infty} dv_a v_a f_a(v_a) \right]. \quad (\text{L.22})$$

To evaluate $g(v)$ we write $|\mathbf{v} - \mathbf{v}_a| = (v^2 - 2vv_a \cos \gamma + v_a^2)/|\mathbf{v} - \mathbf{v}_a|$ and use equation (L.20) to expand $1/|\mathbf{v} - \mathbf{v}_a|$. Using the identity $\int_{-1}^1 dx x P_l(x) = \frac{2}{3}\delta_{l1}$ (eq. C.39 with $m = 0, n = 1$) we find that

$$g(v) = \frac{4\pi}{3} \left[\int_0^v dv_a \left(3v_a^2 v + \frac{v_a^4}{v} \right) f_a(v_a) + \int_v^{\infty} dv_a \left(3v_a^3 + v^2 v_a \right) f_a(v_a) \right]. \quad (\text{L.23})$$

Since the Rosenbluth potentials depend only on v , we can use the relation $\partial v / \partial v_i = v_i / v$ to rewrite equations (L.18) in the form

$$\begin{aligned} D[\Delta v_i] &= \frac{v_i}{v} D[\Delta v_{\parallel}], \\ D[\Delta v_i \Delta v_j] &= \frac{v_i v_j}{v^2} \left(D[(\Delta v_{\parallel})^2] - \frac{1}{2} D[(\Delta \mathbf{v}_{\perp})^2] \right) + \frac{1}{2} \delta_{ij} D[(\Delta \mathbf{v}_{\perp})^2], \end{aligned} \quad (\text{L.24})$$

where

$$\begin{aligned} D[\Delta v_{\parallel}] &= 4\pi G^2 m_a (m + m_a) \ln \Lambda h'(v), \\ D[(\Delta v_{\parallel})^2] &= 4\pi G^2 m_a^2 \ln \Lambda g''(v) \quad ; \quad D[(\Delta \mathbf{v}_{\perp})^2] = \frac{8\pi G^2 m_a^2 \ln \Lambda}{v} g'(v). \end{aligned} \quad (\text{L.25})$$

The reason for this notation can be seen by considering the case in which \mathbf{v} lies along one of the coordinate axes, say $\hat{\mathbf{e}}_1$, so $v_1 = v$ and $v_2 = v_3 = 0$. Then $D[(\Delta v_{\parallel})^2] = D[(\Delta v_1)^2]$ is the total diffusion rate parallel to the velocity vector, and $D[(\Delta \mathbf{v}_{\perp})^2] = 2D[(\Delta v_2)^2] = 2D[(\Delta v_3)^2]$ is the total diffusion rate in the two-dimensional plane perpendicular to the velocity vector. Note that here the subscripts “ \parallel ” and “ \perp ” refer to directions parallel and perpendicular to the subject star velocity \mathbf{v} , whereas in equations (L.4) to (L.13) they refer to directions parallel and perpendicular to the relative velocity $\mathbf{V} = \mathbf{v} - \mathbf{v}_a$.

Substituting from equations (L.22) and (L.23) into (L.25), we have

$$\begin{aligned} D[\Delta v_{\parallel}] &= -\frac{16\pi^2 G^2 m_a (m + m_a) \ln \Lambda}{v^2} \int_0^v dv_a v_a^2 f_a(v_a), \\ D[(\Delta v_{\parallel})^2] &= \frac{32\pi^2 G^2 m_a^2 \ln \Lambda}{3} \left[\int_0^v dv_a \frac{v_a^4}{v^3} f_a(v_a) + \int_v^{\infty} dv_a v_a f_a(v_a) \right], \\ D[(\Delta \mathbf{v}_{\perp})^2] &= \frac{32\pi^2 G^2 m_a^2 \ln \Lambda}{3} \\ &\quad \times \left[\int_0^v dv_a \left(\frac{3v_a^2}{v} - \frac{v_a^4}{v^3} \right) f_a(v_a) + 2 \int_v^{\infty} dv_a v_a f_a(v_a) \right]. \end{aligned} \quad (\text{L.26})$$

Notice that $D[\Delta v_{\parallel}]$ depends only on the total number density of field stars traveling slower than the subject star. This is a simple corollary of Newton's first and second theorems from §2.2.1: the Rosenbluth potential $h(\mathbf{v})$ (eq. L.19) is the velocity-space

analog of the definition (2.3) of the gravitational potential, with \mathbf{v} replacing \mathbf{x} and $f_a(\mathbf{v}_a)$ replacing $\rho(\mathbf{x}')$; the expression for $D[\Delta v_i]$ in (L.18) involves the velocity-space gradient of $h(\mathbf{v})$, which is analogous to the force; so Newton's theorems tell us that this gradient depends only on the number of stars inside the velocity-space sphere of radius $|\mathbf{v}|$.

Appendix M: The distribution of binary energies

We have argued in §§4.10.1 and 7.3.2 that a stellar system has no maximum-entropy state and thus cannot be in thermal equilibrium. Nevertheless, for many purposes the distribution of soft binaries can be approximated by an equilibrium distribution, because soft binaries are so weakly bound that they approach equilibrium much faster than the stellar system as a whole.

M.1 The evolution of the energy distribution of binaries

As in §7.5.7, we let $B(\tilde{E})$ denote the destruction rate for a binary of energy \tilde{E} , while $C(\tilde{E}') d\tilde{E}'$ and $Q(\tilde{E}, \tilde{E}') d\tilde{E}'$ are, respectively, the rate of creation of binaries with energy in the range $(\tilde{E}', \tilde{E}' + d\tilde{E}')$ and the rate at which a binary of energy \tilde{E} makes transitions to energies in this range.¹ These functions are related by

$$B(\tilde{E}) = \int_0^\infty d\tilde{E}' Q(\tilde{E}, \tilde{E}'); \quad (\text{M.1})$$

that is, the rate of destruction of binaries is the rate at which they make transitions to positive energy. The number of binaries $n_b(\tilde{E}, t)$ per unit volume and per unit energy must satisfy a master equation (cf. §7.4.1 and Goodman & Hut 1993)

$$\begin{aligned} \frac{\partial n_b}{\partial t} = & C(\tilde{E}) - B(\tilde{E})n_b(\tilde{E}, t) \\ & + \int_{-\infty}^0 d\tilde{E}' n_b(\tilde{E}', t) Q(\tilde{E}', \tilde{E}) - n_b(\tilde{E}, t) \int_{-\infty}^0 d\tilde{E}' Q(\tilde{E}, \tilde{E}'); \end{aligned} \quad (\text{M.2})$$

the third and fourth terms on the right side denote respectively the rate at which binaries transition into a unit energy range centered on \tilde{E} , and the rate at which binaries already in this energy range are scattered to other bound energies. Equation (M.2) assumes that binaries interact only with single field stars, not with one another; this is reasonable if the fraction of stars in binaries is small.

Fitting functions for $B(\tilde{E})$, $C(\tilde{E})$, and $Q(\tilde{E}, \tilde{E}')$ are given by Heggie & Hut (1993) and the steady-state solution for equation (M.2) is described by Goodman & Hut (1993).

¹ The tilde on \tilde{E} is a flag that in this appendix “energy” denotes a quantity having units mass \times (velocity)², rather than the specific energy (energy per unit mass) with units (velocity)² that is common elsewhere in the book.

M.2 The two-body distribution function in thermal equilibrium

Consider a system of N identical stars of mass m , contained in a volume V . We denote the position and velocity of star n by $\mathbf{w}_n = (\mathbf{x}_n, \mathbf{v}_n)$, and invoke the Jeans swindle to neglect large-scale gravitational fields within the system. As in §7.2.3, we define the two-body DF $f^{(2)}$ such that the probability of finding two stars in the phase-space volume elements $d^6\mathbf{w}_1$ and $d^6\mathbf{w}_2$ is $f^{(2)}(\mathbf{w}_1, \mathbf{w}_2) d^6\mathbf{w}_1 d^6\mathbf{w}_2$.

In thermal equilibrium we expect that

$$f^{(2)}(\mathbf{w}_1, \mathbf{w}_2) = ce^{-\beta\tilde{H}}, \quad (\text{M.3})$$

where c and β are positive constants and the Hamiltonian \tilde{H} is given by

$$\tilde{H}(\mathbf{w}_1, \mathbf{w}_2) = \frac{1}{2}m(v_1^2 + v_2^2) - \frac{Gm^2}{|\mathbf{x}_1 - \mathbf{x}_2|}. \quad (\text{M.4})$$

If the separation of stars 1 and 2 is large, their distributions should be independent, so $f^{(2)}$ should be simply the product of one-body DFs. If these are Maxwellian, with dispersion σ , then we have

$$f^{(2)}(\mathbf{w}_1, \mathbf{w}_2) = \frac{1}{(2\pi\sigma^2)^3 V^2} e^{-v_1^2/2\sigma^2} e^{-v_2^2/2\sigma^2} \quad \text{as } |\mathbf{x}_1 - \mathbf{x}_2| \rightarrow \infty. \quad (\text{M.5})$$

Comparison of equations (M.3) and (M.5) permits us to evaluate the constants c and β . We find

$$f^{(2)}(\mathbf{w}_1, \mathbf{w}_2) = \frac{1}{(2\pi\sigma^2)^3 V^2} \exp \left[-\frac{1}{2\sigma^2} \left(v_1^2 + v_2^2 - \frac{2Gm}{|\mathbf{x}_1 - \mathbf{x}_2|} \right) \right]. \quad (\text{M.6})$$

M.3 The distribution of binary energies in thermal equilibrium

In thermal equilibrium, the number density of binaries with internal energies in the range $(\tilde{E}, \tilde{E} + d\tilde{E})$ is $n_{\text{eq}}(\tilde{E})d\tilde{E}$, where

$$n_{\text{eq}}(\tilde{E}) = \frac{N^2}{2V} \int d^6\mathbf{w}_1 d^6\mathbf{w}_2 f^{(2)}(\mathbf{w}_1, \mathbf{w}_2) \delta[\tilde{E}(\mathbf{w}_1, \mathbf{w}_2) - \tilde{E}]. \quad (\text{M.7})$$

In this expression we have used the identity (C.7), and we have divided by a factor of two because each binary pair is counted twice in the integration. The function $\tilde{E}(\mathbf{w}_1, \mathbf{w}_2)$, the internal energy of the binary, is given by equation (7.162).

At this point it is convenient to convert the integration variables from \mathbf{v}_1 and \mathbf{v}_2 to the relative velocity \mathbf{V} and the center of mass velocity $\mathbf{v}_{\text{cm}} = \frac{1}{2}(\mathbf{v}_1 + \mathbf{v}_2)$. Following the arguments after equation (7.189), we have

$$n_{\text{eq}}(\tilde{E}) = \frac{n^2}{16\pi^3\sigma^6 V} \int d^3\mathbf{x}_1 d^3\mathbf{x}_2 d^3\mathbf{v}_{\text{cm}} d^3\mathbf{V} \times \exp \left[-\frac{1}{\sigma^2} \left(v_{\text{cm}}^2 + \frac{1}{4}V^2 - \frac{Gm}{|\mathbf{x}_1 - \mathbf{x}_2|} \right) \right] \delta \left(\frac{1}{4}mV^2 - \frac{Gm^2}{|\mathbf{x}_1 - \mathbf{x}_2|} - \tilde{E} \right), \quad (\text{M.8})$$

where N/V has been replaced by the number density n . The integral over \mathbf{v}_{cm} is given in equation (7.192). Because of the δ function, we can set $\frac{1}{4}V^2 - Gm/|\mathbf{x}_1 - \mathbf{x}_2| = \tilde{E}/m$ in the exponential. We can replace the dummy variable \mathbf{x}_1 by $\mathbf{x} = \mathbf{x}_1 - \mathbf{x}_2$; moreover, since we are primarily interested in binaries whose separation $|\mathbf{x}|$ is much less than the system size, we can just as well extend the integral over \mathbf{x}

to cover all space (i.e., we neglect edge effects). In this approximation, the integral over \mathbf{x}_2 is independent of \mathbf{x} and yields $\int d^3\mathbf{x}_2 = V$. Thus

$$n_{\text{eq}}(\tilde{E}) = \frac{n^2}{16\pi^{3/2}\sigma^3} e^{-\tilde{E}/m\sigma^2} \int d^3\mathbf{x} d^3\mathbf{v} \delta\left(\frac{1}{4}mV^2 - \frac{Gm^2}{|\mathbf{x}|} - \tilde{E}\right). \quad (\text{M.9})$$

Since the integrand depends only on $|\mathbf{x}|$ and $|\mathbf{v}|$,

$$n_{\text{eq}}(\tilde{E}) = \frac{\pi^{1/2}n^2}{\sigma^3} e^{-\tilde{E}/m\sigma^2} \int_0^\infty dV V^2 \int_0^\infty dx x^2 \delta\left(\frac{1}{4}mV^2 - \frac{Gm^2}{x} - \tilde{E}\right). \quad (\text{M.10})$$

Performing the integration over x , replacing the number density n by the mass density $\rho = mn$, and recalling that $\tilde{E} < 0$, we obtain

$$n_{\text{eq}}(\tilde{E}) = \frac{\pi^{1/2}G^3\rho^2m^4}{\sigma^3} e^{|\tilde{E}|/m\sigma^2} \int_0^\infty \frac{dV V^2}{(\frac{1}{4}mV^2 + |\tilde{E}|)^4}. \quad (\text{M.11})$$

Using the result $\int_0^\infty dx x^2/(x^2 + 1)^4 = \pi/32$, we find

$$n_{\text{eq}}(\tilde{E}) = \frac{\pi^{3/2}G^3\rho^2m^{5/2}}{4\sigma^3|\tilde{E}|^{5/2}} e^{|\tilde{E}|/m\sigma^2}. \quad (\text{M.12})$$

The integral $\int d\tilde{E} n_{\text{eq}}(\tilde{E})$ diverges, both as $|\tilde{E}| \rightarrow 0$ and as $|\tilde{E}| \rightarrow \infty$. The divergence as $|\tilde{E}| \rightarrow \infty$ is not worrying because hard binaries are not in thermal equilibrium anyway; rather, as we have seen in §7.5.7, there is a steady flux of hard binaries towards larger and larger values of $|\tilde{E}|$. The divergence as $|\tilde{E}| \rightarrow 0$ arises because of the large volume of phase space available to a loosely bound binary. However, in a realistic stellar cluster such binaries are very short-lived, because they are disrupted by tidal forces or encounters with field stars. In particular, tidal forces—which are not present in this analysis because of the Jeans swindle—impose an upper limit to the semi-major axis of a binary that is roughly the Jacobi radius $r_J \approx R/N^{1/3}$ (eq. 8.91), where $M = Nm$ and R are the mass and radius of the host system. Hence the minimum $|\tilde{E}|$ for a binary is $|\tilde{E}|_{\text{min}} \approx Gm^2/r_J \approx GmM/(N^{2/3}R)$, and using the virial theorem in the form $\sigma^2 \approx GM/R$ we can write

$$|\tilde{E}|_{\text{min}} \approx \frac{m\sigma^2}{N^{2/3}}. \quad (\text{M.13})$$

Thus we expect that the relation (M.12) applies in the energy range $|\tilde{E}|_{\text{min}} \lesssim |\tilde{E}| \lesssim m\sigma^2$ —all soft binaries that can survive tidal disruption—so long as there has been time for thermal equilibrium to be achieved.

We may use these results to estimate the equilibrium total number of soft binaries in a cluster. We integrate equation (M.12) over the range just derived to obtain the number density of soft binaries,

$$n_{\text{soft}} \simeq \frac{\pi^{3/2}G^3\rho^2m^{5/2}}{4\sigma^3} \int_{|\tilde{E}|_{\text{min}}}^{m\sigma^2} d|\tilde{E}| \frac{e^{|\tilde{E}|/m\sigma^2}}{|\tilde{E}|^{5/2}} \simeq \frac{\pi^{3/2}G^3\rho^2m}{4\sigma^6} \int_{N^{-2/3}}^1 dx \frac{e^x}{x^{5/2}}. \quad (\text{M.14})$$

For $N \gg 1$ the integral is dominated by the region $x \ll 1$, so we can replace the exponential by unity to find

$$n_{\text{soft}} \approx \frac{\pi^{3/2}G^3\rho^2mN}{6\sigma^6}. \quad (\text{M.15})$$

The total number of soft binaries is $N_{\text{soft}} \approx n_{\text{soft}} R^3$, and if we set $\rho \approx Nm/R^3$ and $\sigma^2 \approx GNm/R$, we find $N_{\text{soft}} \approx 1$, that is, there is only of order one soft binary in the whole system! Obviously, soft binaries play no role of any consequence in the evolution of stellar systems; even a large primordial population of soft binaries will rapidly be destroyed by encounters with field stars.

M.4 The principle of detailed balance

Let us temporarily assume that a stellar system is in thermal equilibrium. The rate of creation of binaries per unit volume is $C(\tilde{E})$, and in thermal equilibrium this must equal the rate of destruction of binaries per unit volume, which is $B(\tilde{E})n_{\text{eq}}(\tilde{E})$. Similarly, the rate of transition of binaries from energies in the range $(\tilde{E}, \tilde{E} + d\tilde{E})$ to energies in the range $(\tilde{E}', \tilde{E}' + d\tilde{E}')$ is $n_{\text{eq}}(\tilde{E})Q(\tilde{E}, \tilde{E}') d\tilde{E} d\tilde{E}'$, and this must equal the rate of transition in the opposite direction, which is $n_{\text{eq}}(\tilde{E}')Q(\tilde{E}', \tilde{E}) d\tilde{E}' d\tilde{E}$. Thus

$$\begin{aligned} C(\tilde{E}) &= B(\tilde{E})n_{\text{eq}}(\tilde{E}) \\ n_{\text{eq}}(\tilde{E})Q(\tilde{E}, \tilde{E}') &= n_{\text{eq}}(\tilde{E}')Q(\tilde{E}', \tilde{E}), \end{aligned} \quad (\text{M.16})$$

where $n_{\text{eq}}(\tilde{E})$ is given by equation (M.12). It is easy to verify that if these relations are satisfied then so is the master equation (M.2).

Equation (M.12) can be combined with the first of equations (M.16) to yield

$$C(\tilde{E}) = \frac{\pi^{3/2} G^3 \rho^2 m^{5/2}}{4\sigma^3 |\tilde{E}|^{5/2}} e^{|\tilde{E}|/m\sigma^2} B(\tilde{E}). \quad (\text{M.17})$$

The rates $B(\tilde{E})$ and $C(\tilde{E})$ do not depend on the assumption that the binary distribution is in thermal equilibrium, so long as the field-star distribution is Maxwellian. Hence (M.17) must hold whether or not the binary distribution is in thermal equilibrium. This **principle of detailed balance** allows us to calculate both the formation and destruction rates if either one of them is known (Heggie 1975).

References

Most of the articles in this reference list are available on-line from the Smithsonian/NASA Astrophysics Data System (ADS), adsabs.harvard.edu/abstract_service.html. Articles that are more than three years old are usually free of charge. Newer articles are available either (i) to institutional and individual subscribers through a link to the site maintained by the publisher; or (ii) through a link to a pre-publication version of the manuscript maintained on the free arXiv e-print service, arXiv.org.

We use the following abbreviations for astronomy journals, which conform with ADS standards:

A&A	Astronomy and Astrophysics
A&AS	Astronomy and Astrophysics Supplement Series
AcA	Acta Astronomica
AJ	Astronomical Journal
AN	Astronomische Nachrichten
AnAp	Annales d'Astrophysique
ApJ	Astrophysical Journal
ApJS	Astrophysical Journal Supplement Series
ApL	Astrophysical Letters
Ap&SS	Astrophysics and Space Science
ARA&A	Annual Review of Astronomy and Astrophysics
AZh	Astronomicheskii Zhurnal
BAN	Bulletin of the Astronomical Institutes of the Netherlands
CeMec	Celestial Mechanics
CeMDA	Celestial Mechanics and Dynamical Astronomy
Icar	Icarus
MNRAS	Monthly Notices of the Royal Astronomical Society
Nat	Nature
NewA	New Astronomy
PASA	Publications of the Astronomical Society of Australia
PASJ	Publications of the Astronomical Society of Japan
PASP	Publications of the Astronomical Society of the Pacific
PAZh	Pis'ma Astronomicheskii Zhurnal
QJRAS	Quarterly Journal of the Royal Astronomical Society
Sci	Science
SvA	Soviet Astronomy
SvAL	Soviet Astronomy Letters
ZA	Zeitschrift für Astrophysik

- Aarseth, S.J. 2003. *Gravitational N-Body Simulations: Tools and Algorithms* (Cambridge: Cambridge University Press)
- Aarseth, S.J., & Binney, J. 1978. *MNRAS*, **185**, 227
- Aarseth, S.J., & Heggie, D.C. 1998. *MNRAS*, **297**, 794
- Abramowitz, M., & Stegun, I.A., ed. 1964. *Handbook of Mathematical Functions* (New York: Dover). Also www.math.sfu.ca/~cbm/aands/
- Adler, D.S., & Westpfahl, D.J. 1996. *AJ*, **111**, 735
- Aguerri, J.A.L., Debattista, V.P., & Corsini, E.M. 2003. *MNRAS*, **338**, 465
- Aguilar, L., Hut, P., & Ostriker, J.P. 1988. *ApJ*, **335**, 720
- Aguilar, L., & White, S.D.M. 1985. *ApJ*, **295**, 374
- Ahmad, A., & Cohen, L. 1973. *J. Comput. Phys.*, **12**, 389
- Albrow, M.D., Gilliland, R.L., Brown, T.M., Edmonds, P.D., Guhathakurta, P., & Sarajedini, A. 2001. *ApJ*, **559**, 1060
- Alcock, C., et al. 2001. *ApJ*, **550**, L169
- Alladin, S.M. 1965. *ApJ*, **141**, 768
- Allen, S., Schmidt, R.W., & Fabian, A.C. 2002. *MNRAS*, **334**, L11
- Aly, J.-J., & Pérez, J. 1992. *MNRAS*, **259**, 95
- Ambarzumian, V.A. 1938. *Ann. Leningrad State Univ.*, **22** (in Russian). Also Good-

- man & Hut (1985), 521 (in English)
- Antonov, V.A. 1960. *AZh*, **37**, 918 (in Russian). Also *SvA*, **4**, 859 (in English)
- Antonov, V.A. 1962a, *Vestnik Leningrad Univ.*, **7**, 135 (in Russian). Also Goodman & Hut (1985), 525 (in English)
- Antonov, V.A. 1962b, *Vestnik Leningrad Univ.*, **19**, 96 (in Russian). Also de Zeeuw (1987), 531 (in English)
- Antonov, V.A. 1971. *Trudy Astron. Obs. Leningrad*, **28**, 64 (in Russian)
- Antonov, V.A. 1973. In *The Dynamics of Galaxies and Star Clusters*, ed. G.B. Omarov (Alma Ata: Nauka), 139 (in Russian). Also de Zeeuw (1987), 549 (in English)
- Aoki, S., Noguchi, M., & Iye, M. 1979. *PASJ*, **31**, 737
- Appleton, P.N., & Struck-Marcell, C. 1996. *Fund. Cosmic Physics*, **16**, 111
- Araki, S. 1985. Ph.D. thesis, Massachusetts Institute of Technology
- Araki, S. 1987. *AJ*, **94**, 99
- Arfken, G.B., & Weber, H.J. 2005. *Mathematical Methods for Physicists* (6th ed.; Burlington: Elsevier)
- Arnold, V.I. 1989. *Mathematical Methods of Classical Mechanics* (2nd ed.; New York: Springer)
- Aronson, E.B., & Hansen, C.J. 1972. *ApJ*, **177**, 145
- Arp, H. 1966. *Atlas of Peculiar Galaxies* (Washington: Carnegie Institution). Also *ApJS*, **14**, 1
- Arp, H., & Madore, B.S. 1987. *A Catalogue of Southern Peculiar Galaxies and Associations* (New York: Cambridge University Press)
- Ascasibar, Y., & Binney, J. 2005. *MNRAS*, **356**, 872
- Ashman, K.A., & Zepf, S.E. 1998. *Globular Cluster Systems* (Cambridge: Cambridge University Press)
- Athanassoula, E. 1984. *Phys. Rep.*, **114**, 319
- Athanassoula, E. 1992a. *MNRAS*, **259**, 328
- Athanassoula, E. 1992b. *MNRAS*, **259**, 345
- Athanassoula, E. 2002. *ApJ*, **569**, L83
- Athanassoula, E., Bosma, A., & Papaioannou, S. 1987. *A&A*, **179**, 23
- Baade, W. 1963. *Evolution of Stars and Galaxies* (Cambridge, MA: Harvard University Press), 63
- Bahcall, J.N., Hut, P., & Tremaine, S. 1985. *ApJ*, **290**, 15
- Bahcall, J.N., & Wolf, R.A. 1976. *ApJ*, **209**, 214
- Bailin, J., & Steinmetz, M. 2005. *ApJ*, **627**, 647
- Bailyn, C.D. 1995. *ARA&A*, **33**, 133
- Balbus, S.A. 1988. *ApJ*, **324**, 60
- Barbanis, B., & Woltjer, L. 1967. *ApJ*, **150**, 461
- Bardeen, J.M. 1975. In *IAU Symposium 59, Dynamics of Stellar Systems*, ed. A. Hayli (Dordrecht: Reidel), 297
- Bardeen, J.M., Bond, J.R., Kaiser, N., & Szalay, A.S. 1986. *ApJ*, **304**, 15
- Barnes, J. 1985. In Goodman & Hut 1985, 297
- Barnes, J. 1988. *ApJ*, **331**, 699
- Barnes, J., & Efstathiou, G. 1987. *ApJ*, **319**, 575
- Barnes, J., Goodman, J., & Hut, P. 1986. *ApJ*, **300**, 112
- Barnes, J., & Hernquist, L. 1992. *ARA&A*, **30**, 705
- Barnes, J., & Hut, P. 1986. *Nat*, **324**, 446
- Bartholomew, P. 1971. *MNRAS*, **151**, 333
- Basu, S., Mouschovias, T.C., & Paleologou, E. 1997. *ApJ*, **480**, L55
- Battaglia, G., Fraternali, F., Oosterloo, T., & Sancisi, R. 2006. *A&A*, **447**, 49
- Baumgardt, H. 2001. *MNRAS*, **325**, 1323
- Baumgardt, H., Heggie, D.C., Hut, P., & Makino, J. 2003. *MNRAS*, **341**, 247
- Baumgardt, H., Hut, P., & Heggie, D.C. 2002. *MNRAS*, **336**, 1069
- Baumgardt, H., & Makino, J. 2003. *MNRAS*, **340**, 227
- Baumgardt, H., Makino, J., & Ebisuzaki, E. 2004. *ApJ*, **613**, 11331143
- Beck, R., Brandenburg, A., Moss, D., Shukurov, A., & Sokoloff, D. 1996. *ARA&A*, **34**, 155
- Beck, R., Ehle, M., Shoutenkov, V., Shukurov, A., & Sokoloff, D. 1999. *Nat*, **397**, 324
- Begelman, M.C., Blandford, R.D., & Rees, M.J. 1980. *Nat*, **287**, 307
- Bell, E.F., & de Jong, R.S. 2001. *AJ*, **550**, 212
- Belokurov, V., et al. 2007. *ApJ*, **654**, 897
- Bender, R., et al. 2005. *ApJ*, **631**, 280
- Benson, A.J., Kamionkowski, M., & Hassani, S.H. 2005. *MNRAS*, **357**, 847
- Benson, A.J., Lacey, C.G., Frenk, C.S., Baugh, C.M., & Cole, S. 2004. *MNRAS*, **351**, 1215
- Bettwieser, E. 1983. *MNRAS*, **203**, 811
- Binney, J. 1978. *MNRAS*, **183**, 501
- Binney, J. 1982. *MNRAS*, **201**, 1
- Binney, J. 1987. In *The Galaxy*, ed. G. Gilmore & R. Carswell (Dordrecht: Reidel), 399
- Binney, J. 1992. *ARA&A*, **30**, 51
- Binney, J. 2004a. *MNRAS*, **347**, 1093
- Binney, J. 2004b. *MNRAS*, **350**, 939
- Binney, J. 2005. *MNRAS*, **363**, 937

- Binney, J., & Evans, N.W. 2001. *MNRAS*, **327**, L27
- Binney, J., Gerhard, O., & Hut, P. 1985. *MNRAS*, **215**, 59
- Binney, J., Jiang, I.-G., & Dutta, S. 1998. *MNRAS*, **297**, 1237
- Binney, J., & Mamon G.A. 1982. *MNRAS*, **200**, 361
- Binney, J., & May, A. 1986. *MNRAS*, **218**, 743
- Binney, J., & Merrifield, M. 1998. *Galactic Astronomy* (Princeton: Princeton University Press)
- Binney, J., & Spiegel, D. 1982. *ApJ*, **252**, 308
- Binney, J., & Spiegel, D. 1983. In *IAU Colloquium 76, The Nearby Stars and the Stellar Luminosity Function*, ed. A.G.D. Philip & A.R. Uggren (Schenectady: L. Davis Press), 259
- Bissantz, N., & Gerhard, O. 2002. *MNRAS*, **330**, 591
- Blanton, R.B., Lupton, R.H., Schlegel, D.J., Strauss, M.A., Brinkmann, J., Fukugita, M., & Loveday, J. 2005. *ApJ*, **631**, 208
- Blumenthal, G.R., Faber, S.M., Primack, J.R., & Rees, M.J. 1984. *Nat.*, **311**, 517
- Bode, P., & Ostriker, J.P. 2003. *ApJS*, **145**, 1
- Bond, J.R., Cole, S., Efstathiou, G., & Kaiser, N. 1991. *ApJ*, **379**, 440
- Bond, J.R., Szalay, A.S., & Turner, M.S. 1982. *Phys. Rev. Lett.* **48**, 1636
- Bontekoe, T.J.R. 1988. Ph.D. thesis, University of Groningen
- Bottema, R. 1993. *A&A*, **275**, 16
- Brada, R., & Milgrom, M. 1995. *ApJ*, **444**, 71
- Breeden, J.L., & Cohn, H.N. 1995. *ApJ*, **448**, 672
- Brown, W.R., Geller, M.J., Fabricant, D.G., & Kurtz, M.J. 2001. *AJ*, **122**, 714
- Bryan, G.H. 1888. *Phil. Trans. R. Soc. London A*, **180**, 187
- Bullock, J.S. 2002. In *The Shapes of Galaxies and their Dark Halos*, ed. P. Natarajan (Singapore: World Scientific), 109
- Bullock, J.S., Dekel, A., Kolatt, T.S., Kravtsov, A.V., Klypin, A.A., Porciani, C., & Primack, J.R. 2001. *ApJ*, **555**, 240
- Bureau, M., Aronica, G., Athanassoula, E., Dettmar, R.-J., Bosma, A., & Freeman, K.C. 2006. *MNRAS*, **370**, 753
- Bureau, M., & Freeman, K.C. 1999. *AJ*, **118**, 126
- Buta, R. 1995. *ApJS*, **96**, 39
- Buta, R., Byrd, G.G., & Freeman, T. 2003. *AJ*, **125**, 634
- Buta, R., Crocker, D.A., & Elmegreen, B.G., ed. 1996. *ASP Conference Series 91, Barred Galaxies* (San Francisco: Astronomical Society of the Pacific)
- Butcher, J.C. 1987. *The Numerical Analysis of Ordinary Differential Equations* (Chichester: John Wiley & Sons)
- Calzetti, D., et al. 2005. *ApJ*, **633**, 871
- Cappellari, M., et al. 2006. *MNRAS*, **366**, 1126
- Cappellari, M., et al. 2007. *MNRAS*, **379**, 418
- Carlberg, R.G. 1987. *ApJ*, **322**, 59
- Carney, B.W., & Harris, W.E. 2001. In *Saas-Fee Advanced Course 28, Star Clusters*, ed. L. Labhardt & B. Binggeli (Berlin: Springer)
- Carroll, S.M., Press, W.H., & Turner, E.L. 1992. *ARA&A*, **30**, 499
- Case, K.M. 1959. *Ann. Phys.*, **7**, 349
- Cattaneo, A., Dekel, A., Devriendt, J., Guiderdoni, B., & Blaizot, J. 2006. *MNRAS*, **370**, 1651
- Cayrel, R., et al. 2001. *Nat.*, **409**, 691
- Chanamé, J., & Gould, A. 2004. *ApJ*, **601**, 289
- Chandrasekhar, S. 1939. *An Introduction to the Theory of Stellar Structure* (Chicago: University of Chicago Press). Reissued by Dover 1973
- Chandrasekhar, S. 1942. *Principles of Stellar Dynamics* (Chicago: University of Chicago Press). Reissued by Dover 2005
- Chandrasekhar, S. 1943a. *ApJ*, **97**, 255
- Chandrasekhar, S. 1943b. *Rev. Mod. Phys.*, **15**, 1. Reprinted in Wax, N., ed. 1954. *Selected Papers on Noise and Stochastic Processes* (Dover: New York)
- Chandrasekhar, S. 1961. *Hydrodynamic and Hydromagnetic Stability* (Oxford: Oxford University Press). Reissued by Dover 1981
- Chandrasekhar, S. 1963. *ApJ*, **138**, 896
- Chandrasekhar, S. 1964. *ApJ*, **139**, 664
- Chandrasekhar, S. 1969. *Ellipsoidal Figures of Equilibrium* (New Haven: Yale University Press)
- Chernoff, D.F., & Weinberg, M.D. 1990. *ApJ*, **351**, 121
- Ciotti, L. 1991. *A&A*, **249**, 99
- Ciotti, L., & Bertin, G. 1999. *A&A*, **352**, 447
- Cohn, H. 1979. *ApJ*, **234**, 1036
- Cohn, H. 1980. *ApJ*, **242**, 765
- Cohn, H., & Kulsrud, R.M. 1978. *ApJ*, **226**, 1087
- Combes, F., Debbasch, F., Friedli, D., & Pfenniger, D. 1990. *A&A*, **233**, 82
- Combes, F., & Sanders, R.H. 1981. *A&A*, **96**, 164
- Connors, T.W., Kawata, D., Maddison, S.T., & Gibson, B.K. 2004. *PASA*, **21**, 222
- Conselice, C. 2006. *ApJ*, **638**, 686

- Contopoulos, G. 1954. *ZA*, **35**, 67 (in German)
- Contopoulos, G. 1980. *A&A*, **81**, 198
- Contopoulos, G., & Papayannopoulos, Th. 1980. *A&A*, **92**, 33
- Côté, P., McLaughlin, D.E., Cohen, J.G., & Blakeslee, J.P. 2003. *ApJ*, **591**, 850
- Couchman, H.M.P. 1991. *ApJ*, **368**, L23
- Courteau, S., & Rix, H.-W. 1999. *ApJ*, **513**, 561
- Cox, A.N., ed. 2000. *Allen's Astrophysical Quantities* (4th ed.; New York: Springer)
- Crézé, M., Chereul, E., Bienaymé, O., & Pichon, C. 1998. *A&A*, **329**, 920
- Cuddeford, P. 1993. *MNRAS*, **262**, 1076
- Davies, R.L., Efstathiou, G., Fall, S.M., Illingworth, G., & Schechter, P.L. 1983. *ApJ*, **266**, 41
- de Jong, R.S. 1996. *A&A*, **313**, 45
- De Simone, R.S., Wu, X., & Tremaine, S. 2004. *MNRAS*, **350**, 627
- de Vaucouleurs, G. 1948. *AnAp*, **11**, 247 (in French)
- de Vaucouleurs, G. 1959. In *Handbuch der Physik* **53**, ed. S. Flügge (Berlin: Springer), 275
- de Vaucouleurs, G. 1964. In *IAU Symposium 20, The Galaxy and the Magellanic Clouds*, ed. F.J. Kerr & A.W. Rodgers (Canberra: Australian Academy of Science), 195
- de Zeeuw, T. 1985. *MNRAS*, **216**, 273
- de Zeeuw, T., ed. 1987. *IAU Symposium 127, Structure and Dynamics of Elliptical Galaxies* (Dordrecht: Reidel)
- Debattista, V., Corsini, E.M., & Aguerri, J.A.L. 2002. *MNRAS*, **332**, 65
- Debattista, V., Gerhard, O., & Sevenster, M.N. 2002. *MNRAS*, **334**, 355
- Debattista, V., & Sellwood, J.A. 1998. *ApJ*, **493**, L5
- Debattista, V., & Sellwood, J.A. 1999. *ApJ*, **513**, L107
- Debattista, V., & Sellwood, J.A. 2000. *ApJ*, **543**, 704
- Debattista, V., & Williams, T.B. 2001. In *ASP Conference Series 230, Galaxy Disks and Disk Galaxies*, ed. J.G. Funes & E.M. Corsini (San Francisco: Astronomical Society of the Pacific), 553
- Dehnen, W. 1993. *MNRAS*, **265**, 250
- Dehnen, W. 1999a. *AJ*, **118**, 1190
- Dehnen, W. 1999b. *AJ*, **118**, 1201
- Dehnen, W. 2000a. *AJ*, **119**, 800
- Dehnen, W. 2000b. *ApJ*, **536**, L39
- Dehnen, W. 2001. *MNRAS*, **324**, 273
- Dehnen, W. 2005. *MNRAS*, **360**, 892
- Dehnen, W., & Binney, J.J. 1998a. *MNRAS*, **294**, 429
- Dehnen, W., & Binney, J.J. 1998b. *MNRAS*, **298**, 387
- Dejonghe, H. 1986. *Phys. Rep.*, **133**, 217
- Dejonghe, H., & Merritt, D. 1992. *ApJ*, **391**, 531
- Dekel, A., Arad, I., Devor, J., & Birnboim, Y. 2003. *ApJ*, **588**, 680
- Dekel, A., & Silk, J. 1986. *ApJ*, **303**, 39
- Dekel, A., & Woo, J. 2003. *MNRAS*, **344**, 1131
- Diemand, J., Moore, B., & Stadel, J. 2004. *MNRAS*, **353**, 624
- Diener, P., Kosovichev, A.G., Kotok, E.V., Novikov, I.D., & Pethick, C.J. 1995. *MNRAS*, **275**, 498
- Donahue, M., Horner, D.J., Cavagnolo, K.W., & Voit, G.M. 2006. *ApJ*, **643**, 730
- Doremus, J.P., Feix, M.R., & Baumann, G. 1971. *Phys. Rev. Lett.* **26**, 725
- Draine, B.T. 2004. In *Saas-Fee Advanced Course 32, The Cold Universe*, ed. D. Pfenniger & Y. Revaz (Berlin: Springer), 1
- Dressler, A., Lynden-Bell, D., Burstein, D., Davies, R.L., Faber, S.M., Terlevich, R.J., & Wegner, G. 1987. *ApJ*, **313**, 42
- Drukier, G.A., Cohn, H.N., Lugger, P.M., & Yong, H. 1999. *ApJ*, **518**, 233
- Dubinski, J. 1998. *ApJ*, **502**, 141
- Dubinski, J., & Carlberg, R.G. 1991. *ApJ*, **378**, 496
- Dubinski, J., da Costa, L.N., Goldwirth, D.S., Lecar, M., & Piran, T. 1993. *ApJ*, **410**, 458
- Dubinski, J., & Farah, J. 2006. *Gravitas: portraits of a Universe in Motion* (DVD media), www.galaxydynamics.org (ISBN 0-9738457-0-8)
- Dubinski, J., Mihos, J.C., & Hernquist, L. 1996. *ApJ*, **462**, 576
- Dubinski, J., Mihos, J.C., & Hernquist, L. 1999. *ApJ*, **526**, 607
- Dumas, H.S., & Laskar, J. 1993. *Phys. Rev. Lett.* **70**, 2975
- Earn, D.J.D., & Sellwood, J.A. 1995. *ApJ*, **451**, 533
- Eddington, A.S. 1916a. *MNRAS*, **76**, 525
- Eddington, A.S. 1916b. *MNRAS*, **76**, 572
- Efstathiou, G., Davis, M., Frenk, C.S., & White, S.D.M. 1985. *ApJS*, **57**, 241
- Einstein, A. 1921. In *Festschrift der Kaiser-Wilhelm-Gesellschaft zur Förderung der Wissenschaftler zu ihrem zehnjährigen Jubiläum* (Berlin: Springer), 50 (in German). Also in M. Janssen et al. (2002), *The Collected Papers of Albert Einstein*, **7**, The Berlin Years: Writings, 1918–1921 (Princeton: Princeton University Press), 181
- Eisenhauer, F., Schödel, R., Genzel, R., Ott, T., Tecza, M., Abuter, R., Eckart, A., & Alexander, T. 2003. *ApJ*, **597**, L121

- Elmegreen, D.M. 1981. *ApJS*, **47**, 229
- Elmegreen, D.M. 1996. In Buta et al. (1996), 23
- Elmegreen, D.M. 1998. *Galaxies and Galactic Structure* (Upper Saddle River: Prentice Hall)
- Elmegreen, D.M., & Elmegreen, B.G. 1982. *MNRAS*, **201**, 1021
- Elmegreen, D.M., & Elmegreen, B.G. 1984. *ApJS*, **54**, 127
- Elmegreen, D.M., & Elmegreen, B.G. 1987. *ApJ*, **314**, 3
- Elmegreen, D.M., Elmegreen, B.G., Chromey, F.R., Hasselbacher, D.A., & Bissell, B.A. 1996. *AJ*, **111**, 2233
- England, M.N., Gottesman, S.T., & Hunter, J.H. 1990. *ApJ*, **348**, 456
- Englmaier, P., & Gerhard, O. 1999. *MNRAS*, **304**, 512
- Eskridge, P.B., et al. 2000. *AJ*, **119**, 536
- Eskridge, P.B., et al. 2002. *ApJS*, **143**, 73
- Evans, N.W. 1993. *MNRAS*, **260**, 191
- Evans, N.W. 1994. *MNRAS*, **267**, 333
- Evans, N.W., & Collett, J.L. 1994. *ApJ*, **420**, L67
- Evans, N.W., & de Zeeuw, P.T. 1992. *MNRAS*, **257**, 152
- Evans, N.W., & Read, J.C.A. 1998. *MNRAS*, **300**, 106
- Faber, S.M., et al. 1997. *AJ*, **114**, 1771
- Fall, S.M., & Rees, M.J. 1977. *MNRAS*, **181**, 37p
- Fall, S.M., & Zhang, Q. 2001. *ApJ*, **561**, 751
- Famaey, B., & Binney, J. 2005. *MNRAS*, **363**, 603
- Famaey, B., Jorissen, A., Luri, X., Mayor, M., Udry, S., Dejonghe, H., & Turon, C. 2005. *A&A*, **430**, 165
- Farouki, R.T., & Salpeter, E.E. 1994. *ApJ*, **427**, 676
- Feller, W. 1971. *An Introduction to Probability Theory and its Applications* (2nd ed.; New York: Wiley)
- Flynn, C., Holmberg, J., Portinari, L., Fuchs, B., & Jahreiss, H. 2006. *MNRAS*, **372**, 1149
- Forest, E., & Ruth, R.D. 1990. *Physica D*, **43**, 105
- Frank, J., & Rees, M. J. 1976. *MNRAS*, **176**, 633
- Freedman, W.L., et al. 2001. *ApJ*, **553**, 47
- Freeman, K.C. 1966. *MNRAS*, **134**, 15
- Freeman, K.C. 1970. *ApJ*, **160**, 811
- Freeman, K.C. 1993. In *ASP Conference Series 48, The Globular Cluster-Galaxy Connection*, ed. G.H. Smith & J.P. Brodie (San Francisco: Astronomical Society of the Pacific), 608
- Freitag, M., & Benz, W. 2001. *A&A*, **375**, 711
- Freitag, M., Rasio, F.A., & Baumgardt, H. 2006. *MNRAS*, **368**, 121
- Frenk, C.S., et al., 1999. *ApJ*, **525**, 554
- Fricke, W. 1952. *AN*, **280**, 193 (in German)
- Fridman, A.M., & Polyachenko, V.L. 1984. *Physics of Gravitating Systems* (New York: Springer-Verlag)
- Fried, B.D., & Conte, S.D. 1961. *The Plasma Dispersion Function* (New York: Academic Press)
- Fuchs, B. 2001. *A&A*, **368**, 107
- Fuchs, B., Dettbarn, C., & Tsuchiya, T. 2005. *A&A*, **444**, 1
- Fujii, M., Funato, Y., & Makino, J. 2006. *PASJ*, **58**, 743
- Fukugita, M., Hogan, C.J., & Peebles, P.J.E. 1998. *ApJ*, **503**, 518
- Fukugita, M., & Peebles, P.J.E. 2004. *ApJ*, **616**, 642
- Fukushige, T., & Heggie, D.C. 2000. *MNRAS*, **318**, 753
- Fux, R. 2001. *A&A*, **373**, 511
- Gaitskell, R.J. 2004. *Ann. Rev. Nucl. Part. Sci.*, **54**, 315
- Gao, L., White, S.D.M., Jenkins, A., Stoehr, F., & Springel, V. 2004. *MNRAS*, **355**, 819
- Gardiner, L.T., Sawa, T., & Fujimoto, M. 1994. *MNRAS*, **266**, 567
- Gebhardt, K., et al. 2003. *ApJ*, **583**, 92
- Gebhardt, K., & Kissler-Patig, M. 1999. *AJ*, **118**, 1526
- Gerhard, O. 1991. *MNRAS*, **250**, 812
- Gerhard, O. 1993. *MNRAS*, **265**, 213
- Gerhard, O. 2002. In *ASP Conference Series 273, The Dynamics, Structure and History of Galaxies*, ed. G.S. Da Costa & H. Jerjen (San Francisco: Astronomical Society of the Pacific), 73
- Gerhard, O., & Binney, J.J. 1985. *MNRAS*, **216**, 467
- Gerhard, O., Kronawitter, A., Saglia, R.P., & Bender R. 2001. *AJ*, **121**, 1936
- Gerola, H., & Seiden, P.E. 1978. *ApJ*, **223**, 129
- Gerssen, J., Kuijken, K., & Merrifield, M. 1999. *MNRAS*, **306**, 926
- Gerssen, J., Kuijken, K., & Merrifield, M. 2003. *MNRAS*, **345**, 261
- Gibbs, J.W. 1884. *Proc. Am. Assoc. Adv. Sci.*, **33**, 57
- Giersz, M. 1998. *MNRAS*, **298**, 1239
- Giersz, M., & Heggie, D.C. 1994. *MNRAS*, **268**, 257
- Giersz, M., & Spurzem, R. 1994. *MNRAS*, **269**, 241
- Gilbert, I.H. 1968. *ApJ*, **152**, 1043
- Gilmore, G., & Reid, N. 1983. *MNRAS*, **202**, 1025

- Gnedin, N.Y., & Hamilton, A.J.S. 2002. *MNRAS*, **334**, 107
- Gnedin, O.Y., Goodman, J., & Frei, Z. 1995. *AJ*, **110**, 1105
- Gnedin, O.Y., Hernquist, L., & Ostriker, J.P. 1999. *ApJ*, **514**, 109
- Gnedin, O.Y., Lee, H.M., & Ostriker, J.P. 1999. *ApJ*, **522**, 935
- Gnedin, O.Y., & Ostriker, J.P. 1997. *ApJ*, **474**, 223
- Goldreich, P., & Lynden-Bell, D. 1965a. *MNRAS*, **130**, 125
- Goldreich, P., & Lynden-Bell, D. 1965b. *MNRAS*, **130**, 97
- Goldreich, P., & Tremaine, S. 1978. *ApJ*, **222**, 850
- Goldreich, P., & Tremaine, S. 1979. *ApJ*, **233**, 857
- Goldreich, P., & Tremaine, S. 1981. *ApJ*, **243**, 1062
- Goldstein, H., Safko, J.L., & Poole, C.P. 2002. *Classical Mechanics* (3rd ed.; San Francisco: Addison Wesley)
- Goodman, J. 1983. *ApJ*, **270**, 700 (erratum in *ApJ*, **278**, 893)
- Goodman, J. 1984. *ApJ*, **280**, 298
- Goodman, J. 1987. *ApJ*, **313**, 576
- Goodman, J. 1988. *ApJ*, **329**, 612
- Goodman, J., & Binney, J. 1984. *MNRAS*, **207**, 511
- Goodman, J., Heggie, D.C., & Hut, P. 1993. *ApJ*, **415**, 715
- Goodman, J., & Hut, P. 1993. *ApJ*, **403**, 271
- Goodman, J., & Hut, P., ed. 1985. *IAU Symposium 113, Dynamics of Star Clusters* (Dordrecht: Reidel)
- Goudfrooij, P., Gilmore, D., Whitmore, B.C., & Schweizer, F. 2004. *ApJ*, **613**, L121
- Gradshteyn, I.S., & Ryzhik, I.M. 2000. *Table of Integrals, Series, and Products* (6th ed.; San Diego: Academic Press)
- Grillmair, C.J., & Dionatos, O. 2006. *ApJ*, **641**, L37
- Gunn, J.E. 1978. In *Eighth Advanced Course of the Swiss Society of Astronomy and Astrophysics, Observational Cosmology*, ed. A. Maeder, L. Martinet, & G. Tammann (Sauverny: Geneva Observatory), 1
- Gunn, J.E., & Gott, J.R. 1972. *ApJ*, **176**, 1
- Hachisu, I., Nakada, Y., Nomoto, K., & Sugimoto, D. 1978. *Prog. Theor. Phys.*, **60**, 393
- Häfner, R., Evans, N.W., Dehnen, W., & Binney, J. 2000. *MNRAS*, **314**, 433
- Hairer, E., Lubich, C., & Wanner, G. 2002. *Geometric Numerical Integration* (Berlin: Springer)
- Han, Z., Podsiadlowski, P., & Eggleton, P.P. 1994. *MNRAS*, **270**, 121
- Hänninen, J., & Flynn, C. 2002. *MNRAS*, **337**, 731
- Häring, N., & Rix, H.-W. 2004. *ApJ*, **604**, L89
- Harris, W.E. 1996. *AJ*, **112**, 1487. See also physwww.mcmaster.ca/%7Eeharris/Databases.html
- Harris, W.E. 2003. In *A Decade of Hubble Space Telescope Science*, ed. M. Livio, K. Noll, & M. Stiavelli (Cambridge: Cambridge University Press), 78
- Hasan, H., & Norman, C. 1990. *ApJ*, **361**, 69
- Hawkins, E., et al. 2003. *MNRAS*, **346**, 78
- Hayashi, E., et al. 2004. *MNRAS*, **355**, 794
- Heggie, D.C. 1975. *MNRAS*, **173**, 729
- Heggie, D.C., Giersz, M., Spurzem, R., & Takahashi, K. 1998. In *Highlights of Astronomy 11A*, ed. J. Andersen (Dordrecht: Kluwer), 591
- Heggie, D.C., & Hut, P. 1993. *ApJS*, **85**, 347
- Heggie, D.C., & Hut, P. 2003. *The Gravitational Million Body Problem* (Cambridge: Cambridge University Press)
- Helmi, A., & White, S.D.M. 1999. *MNRAS*, **307**, 495
- Helmi, A., White, S.D.M., & Springel, V. 2003. *MNRAS*, **339**, 834
- Hemsendorf, M., & Merritt, D. 2002. *ApJ*, **580**, 606
- Hénon, M. 1959. *AnAp*, **22**, 491 (in French)
- Hénon, M. 1960a. *AnAp*, **23**, 467 (in French)
- Hénon, M. 1960b. *AnAp*, **23**, 474 (in French)
- Hénon, M. 1960. *AnAp*, **23**, 668 (in French) and 1969. *A&A*, **2**, 151
- Hénon, M. 1961. *AnAp*, **24**, 369 (in French)
- Hénon, M. 1970. *A&A*, **9**, 24
- Hénon, M. 1972. In *IAU Colloquium 10, Gravitational N-Body Problem*, ed. M. Lecar (Dordrecht: Reidel), 44, 406. See also *Ap&SS*, **13**, 284 & **14**, 151
- Hénon, M. 1973a. *A&A*, **24**, 229
- Hénon, M. 1973b. In *Third Advanced Course of the Swiss Society of Astronomy and Astrophysics, Dynamical Structure and Evolution of Stellar Systems*, ed. L. Martinet & M. Mayor (Sauverny: Geneva Observatory), 183
- Hénon, M. 1982. *A&A*, **114**, 211
- Hénon, M. 1997. *Generating Families in the Restricted Three-Body Problem* (Berlin: Springer)
- Hernquist, L. 1990. *ApJ*, **356**, 359
- Hernquist, L., & Ostriker, J.P. 1992. *ApJ*, **386**, 375
- Hernquist, L., & Quinn, P.J. 1988. *ApJ*, **331**, 682
- Hernquist, L., & Quinn, P.J. 1989. *ApJ*, **342**, 1

- Hernquist, L., & Weinberg, M.D. 1992. *ApJ*, **400**, 80
- Herrnstein, J.R., Moran, J.M., Greenhill, L.J., & Trotter, A.S. 2005. *ApJ*, **629**, 719
- Heymans, C., et al. 2006. *MNRAS*, **371**, L60
- Hibbard, J.E., van der Hulst, J.M., Barnes, J.E., & Rich, R.M. 2001. *AJ*, **122**, 2969
- Hickson, P. 1997. *ARA&A*, **35**, 357
- Hills, J.G. 1975. *AJ*, **80**, 809
- Hockney, R.W., & Eastwood, J.W. 1988. *Computer Simulation Using Particles* (Bristol: Adam Hilger)
- Hodge, P. 1992. *The Andromeda Galaxy* (Dordrecht: Kluwer)
- Hohl, F. 1971. *ApJ*, **168**, 343
- Hohl, F., & Hockney, R.W. 1969. *J. Comput. Phys.*, **4**, 306
- Hohl, F., & Zang, T.A. 1979. *AJ*, **84**, 585
- Holm, D.D., Marsden, J.E., Ratiu, T., & Weinstein, A. 1985. *Phys. Rep.*, **123**, 1
- Holmberg, E. 1941. *ApJ*, **94**, 385
- Holmberg, J., & Flynn, C. 2000. *MNRAS*, **313**, 209
- Holmberg, J., & Flynn, C. 2004. *MNRAS*, **352**, 440
- Horedt, G.P. 2004. *Polytropes* (Dordrecht: Kluwer)
- Horwitz, G., & Katz, J. 1978. *ApJ*, **222**, 941
- Houghton, R.C.W., Magorrian, J., Sarzi, M., Thatte, N., Davies, R.L., & Krajnović, D. 2006. *MNRAS*, **367**, 2
- Hubble, E. 1930. *ApJ*, **71**, 231
- Hubble, E. 1943. *ApJ*, **97**, 112
- Humphrey, P.J., Buote, D.A., Gastaldello, F., Zappacosta, L., Bullock, J.S., Brighenti, F., & Mathews, W.G. 2006. *ApJ*, **646**, 899
- Hunter, C. 1963. *MNRAS*, **126**, 299
- Hunter, C. 1965. *MNRAS*, **129**, 321
- Hunter, C. 1977. *AJ*, **82**, 271
- Hunter, C. 2001. *MNRAS*, **328**, 829
- Hunter, C., & Qian, E. 1993. *MNRAS*, **262**, 401
- Hunter, C., & Toomre, A. 1969. *ApJ*, **155**, 747
- Hunter, J.H., England, M.N., Gottesman, S.T., Ball, R., & Huntley, J.M. 1988. *ApJ*, **324**, 721
- Hurley, J.R., Pols, O.R., Aarseth, S.J., & Tout, C.A. 2005. *MNRAS*, **363**, 293
- Hurley, J.R., Pols, O.R., & Tout, C.A. 2000. *MNRAS*, **315**, 543
- Huss, A., Jain, B., & Steinmetz M. 1999. *ApJ*, **517**, 64
- Hut, P., & Bahcall, J.N. 1983. *ApJ*, **268**, 319
- Hut, P., et al. 1992. *PASP*, **104**, 981
- Hut, P., & Heggie, D.C. 2002. *J. Stat. Phys.*, **109**, 1017
- Hut, P., & Tremaine, S. 1985. *AJ*, **90**, 1548
- Ida, S., Kokubo, E., & Makino, J. 1993. *MNRAS*, **263**, 875
- Inagaki, S. 1980. *PASJ*, **32**, 213
- Ipser, J.R., & Kandrup, H.E. 1980. *ApJ*, **241**, 1141
- Ivanov, P. B., & Chernyakova, M. A. 2006. *A&A*, **448**, 843
- Ivanova, N., Rasio, F.A., Lombardi, J.C., Dooley, K.L., & Proulx, Z.F. 2005. *ApJ*, **621**, L109
- Jackson, J.D. 1999. *Classical Electrodynamics* (3rd ed.; New York: John Wiley & Sons)
- Jaffe, W. 1983. *MNRAS*, **202**, 995
- Jalali, M.A., & Hunter, C. 2005. *ApJ*, **630**, 804
- James, R.A. 1977. *J. Comput. Phys.*, **25**, 71
- Jarvis, B.J., & Freeman, K.C. 1985. *ApJ*, **295**, 314
- Jeans, J.H. 1902. *Phil. Trans. R. Soc. London*, **199**, 1
- Jeans, J.H. 1915. *MNRAS*, **76**, 70
- Jeans, J.H. 1919. *Phil. Trans. R. Soc. London A*, **218**, 157
- Jenkins, A. 1992. *MNRAS*, **257**, 620
- Jenkins, A., & Binney, J. 1990. *MNRAS*, **245**, 305
- Jenkins, A., Frenk, C.S., White, S.D.M., Colberg, J.M., Cole, S., Evrard, A.E., Couchman, H.M.P., & Yoshida, N. 2001. *MNRAS*, **321**, 372
- Jiang, I.-G., & Binney, J. 1999. *MNRAS*, **303**, L7
- Jørgensen, I., Franx, M., & Kjaergaard, P. 1996. *MNRAS*, **280**, 167
- Johnson, H. M. 1957. *AJ*, **62**, 19
- José, J.V., & Saletan, E.J. 1998. *Classical Dynamics: A Contemporary Approach* (Cambridge: Cambridge University Press)
- Julian, W.H. 1967. *ApJ*, **148**, 175
- Julian, W.H., & Toomre, A. 1966. *ApJ*, **146**, 810
- Kaasalainen, M., & Binney, J. 1994. *Phys. Rev. Lett.* **73**, 2377
- Kahn, F.D., & Woltjer, L. 1959. *ApJ*, **130**, 705
- Kallivayalil, N., van der Marel, R.P., Alcock, C., Axelrod, T., Cook, K.H., Drake, A.J., & Geha, M. 2006. *ApJ*, **638**, 772
- Kalnajs, A.J. 1965. Ph.D. thesis, Harvard University
- Kalnajs, A.J. 1971. *ApJ*, **166**, 275
- Kalnajs, A.J. 1972a. *ApJ*, **175**, 63
- Kalnajs, A.J. 1972b. In *IAU Colloquium 10, Gravitational N-Body Problem*, ed. M. Lecar (Dordrecht: Reidel), 13. See also *Ap&SS*, **13**, 279
- Kalnajs, A.J. 1972c. *ApL*, **11**, 41
- Kalnajs, A.J. 1973a. *ApJ*, **180**, 1023 (erratum in *ApJ*, **185**, 393)

- Kalnajs, A.J. 1973b. *PASA*, **2**, 174
- Kalnajs, A.J. 1976. *ApJ*, **205**, 751
- Kalnajs, A.J. 1977. *ApJ*, **212**, 637
- Kalnajs, A.J. 1978. In *IAU Symposium 77, Structure and Properties of Nearby Galaxies*, ed. E.M. Berkhuijsen & R. Wielebinski (Dordrecht: Reidel), 113
- Kalnajs, A.J. 1991. In *Dynamics of Disc Galaxies*, ed. B. Sundelius (Göteborg: Dept. of Astronomy and Astrophysics, Göteborgs University and Chalmers University of Technology), 323
- Kalnajs, A.J. 1999. In *ASP Conference Series 165, The Galactic Halo*, ed. B.K. Gibson, T.S. Axelrod, & M.E. Putman (San Francisco: Astronomical Society of the Pacific), 325
- Kalnajs, A.J., & Athanassoula-Georgala, E. 1974. *MNRAS*, **168**, 287
- Kandrup, H.E., & Sygnet, J.F. 1985. *ApJ*, **298**, 27
- Kapteyn, J.C. 1922. *ApJ*, **55**, 302
- Katz, J. 1978. *MNRAS*, **183**, 765
- Katz, J. 1979. *MNRAS*, **189**, 817
- Katz, J. 2003. *Found. Phys.*, **33**, 223
- Kaufman, M., Bash, F.N., Hine, B., Rots, A.H., Elmegreen, D.M., & Hodge, P.W. 1989. *ApJ*, **345**, 674
- Kazantzidis, S., Mayer, L., Mastropietro, C., Diemand, J., Stadel, J., & Moore, B. 2004. *ApJ*, **608**, 663
- Kellogg, O.D. 1953. *Foundations of Potential Theory* (New York: Dover)
- Kennicutt, R.C. 1998. *ARA&A*, **36**, 189
- Kennicutt, R.C., & Edgar, B.K. 1986. *ApJ*, **300**, 132
- Kennicutt, R.C., Schweizer, F., & Barnes, J.E. 1998. *Saas-Fee Advanced Course 26 Lecture Notes, Galaxies: Interactions and Induced Star Formation* (Berlin: Springer)
- Kent, S.M., & Gunn, J.E. 1982. *AJ*, **87**, 945
- Kiessling, M.K.-H. 2003. *Adv. Appl. Math.*, **31**, 132
- Kim, W.-T., & Ostriker, E.C. 2002. *ApJ*, **570**, 132
- King, I.R. 1962. *AJ*, **67**, 471
- King, I.R. 1966. *AJ*, **71**, 64
- King, I.R. 1981. *QJRAS*, **22**, 227
- Klypin, A., Zhao, H., & Somerville, R. 2002. *ApJ*, **573**, 597
- Knapen, J.H. 1999. In *ASP Conference Series 187, The Evolution of Galaxies on Cosmological Timescales*, ed. J.E. Beckman & T.J. Mahoney (San Francisco: Astronomical Society of the Pacific), 72
- Knapen, J.H., & Beckman, J.E. 1996. *MNRAS*, **283**, 251
- Knebe, A., Green, A., & Binney, J. 2001. *MNRAS*, **325**, 845
- Kochanek, C.S. 1992. *ApJ*, **385**, 604
- Komossa, S., Halpern, J., Schartel, N., Hasinger, G., Santos-Lleo, M., & Predehl, P. 2004. *ApJ*, **603**, L17
- Kormendy, J. 2004. In *Coevolution of Black Holes and Galaxies*, ed. L.C. Ho (Pasadena: Carnegie Observatories), 1
- Kormendy, J., & Kennicutt, R.C. 2004. *ARA&A*, **42**, 603
- Kormendy, J., & Norman, C. 1979. *ApJ*, **233**, 539
- Kranz, T., Slyz, A., & Rix, H.-W. 2003. *ApJ*, **586**, 143
- Krauss, L.M., & Chaboyer, B. 2003. *Sci*, **299**, 65
- Krolik, J. 1999. *Active Galactic Nuclei* (Princeton: Princeton University Press)
- Kronawitter, A., Saglia, R.P., Gerhard, O., & Bender, R. 2000. *A&AS*, **144**, 53
- Kroupa, P. 2002. *Sci*, **295**, 82
- Kuijken, K., & Gilmore, G. 1989. *MNRAS*, **239**, 605
- Kuijken, K., & Gilmore, G. 1991. *ApJ*, **367**, L9
- Kuijken, K., & Merrifield, M.R. 1995. *ApJ*, **443**, L13
- Kuijken, K., & Tremaine, S. 1991. In *Dynamics of Disc Galaxies*, ed. B. Sundelius (Göteborg: Dept. of Astronomy and Astrophysics, Göteborgs University and Chalmers University of Technology), 71
- Kulsrud, R.M. 2005. *Plasma Physics for Astrophysics* (Princeton: Princeton University Press)
- Kulsrud, R.M., & Mark, J.W.-K. 1970. *ApJ*, **160**, 471
- Kulsrud, R.M., Mark, J.W.-K., & Caruso, A. 1971. *Ap&SS*, **14**, 52
- Kumar, P., & Goodman, J. 1996. *ApJ*, **466**, 946
- Kuzmin, G.G. 1956. *AZh*, **33**, 27 (in Russian)
- Lacey, C., & Cole, S. 1993. *MNRAS*, **262**, 627
- Lacey, C., & Ostriker, J.P. 1985. *ApJ*, **299**, 633
- Lai, D., Rasio, F.A., & Shapiro, S.L. 1994. *ApJ*, **437**, 742
- Laine, S., Shlosman, I., & Heller, C.H. 1998. *MNRAS*, **297**, 1052
- Landau, L. 1932. *Phys. Z. Sowjetunion*, **1**, 285. See also ter Haar, D. 1965. *Collected Papers of L.D. Landau* (New York: Gordon & Breach), 60
- Landau, L. 1936. *Phys. Z. Sowjetunion*, **10**, 154 (in German). See also ter Haar, D. 1965. *Collected Papers of L.D. Landau* (New York: Gordon & Breach), 163 (in English)
- Landau, L.D. 1946. *J. Phys. (USSR)*, **10**, 25.

- See also ter Haar, D. 1965. *Collected Papers of L.D. Landau* (New York: Gordon & Breach), 445
- Landau, L.D., & Lifshitz, E.M. 1980. *Statistical Physics* (3rd ed.; Oxford: Pergamon)
- Landau, L.D., & Lifshitz, E.M. 1989. *Mechanics* (3rd ed.; Oxford: Pergamon)
- Landau, L.D., & Lifshitz, E.M. 1999. *The Classical Theory of Fields* (4th ed.; Oxford: Butterworth-Heinemann)
- Landau, L.D., & Lifshitz, E.M. 2000. *Fluid Mechanics* (2nd ed.; Oxford: Butterworth-Heinemann)
- Laplace, P.S. 1829. *Celestial Mechanics*, trans. N. Bowditch (New York: Chelsea)
- Larsen, S.S. 2004. In *ASP Conference Series 322, The Formation and Evolution of Massive Young Star Clusters*, ed. H.J.G.L.M. Lamers, L.J. Smith, & A. Nota (San Francisco: Astronomical Society of the Pacific), 19
- Larson, R.B. 1970a. *MNRAS*, **147**, 323
- Larson, R.B. 1970b. *MNRAS*, **150**, 93
- Larson, R.B., & Tinsley, B.M. 1978. *ApJ*, **219**, 46
- Laskar, J. 1990. *Icar*, **88**, 266
- Lauer, T.R., et al. 1995. *AJ*, **110**, 2622
- Laval, G., Mercier, C., & Pellat, R. 1965. *Nucl. Fusion*, **5**, 156
- Lebovitz, N.R. 1965. *ApJ*, **142**, 229
- Leeuw, F., Combes, F., & Binney, J. 1993. *MNRAS*, **262**, 1013
- Lichtenberg, A.J., & Lieberman, M.A. 1992. *Regular and Chaotic Dynamics* (2nd ed.; New York: Springer)
- Lifshitz, E.M., & Pitaevskii, L.P. 1981. *Physical Kinetics* (Oxford: Pergamon)
- Lighthill, M.J. 1978. *Waves in Fluids* (Cambridge: Cambridge University Press)
- Lightman, A. P., & Shapiro, S. L. 1977. *ApJ*, **211**, 244
- Lin, C.C., Mestel, L., & Shu, F.H. 1965. *ApJ*, **142**, 1431
- Lin, C.C., & Shu, F.H. 1966. *Proc. Nat. Acad. Sci.*, **55**, 229
- Lin, D.N.C., & Lynden-Bell, D. 1982. *MNRAS*, **198**, 707
- Lindblad, P.A.B., Lindblad, P.O., & Athanassoula, E. 1996. *A&A*, **313**, 65
- Lindblom, L. 1992. *Phil. Trans. R. Soc. London A*, **340**, 353
- López-Corredoira, M., Cabrera-Lavers, A., & Gerhard, O. 2005. *A&A*, **439**, 107
- Lotz, J.M., Telford, R., Ferguson, H.C., Miller, B.W., Stiavelli, M., & Mack, J. 2001. *ApJ*, **552**, 572
- Louis, P.D. 1990. *MNRAS*, **244**, 478
- Lovelace, R.V.E., Jore, K.P., & Haynes, M.P. 1997. *ApJ*, **475**, 83
- Lupton, R., Blanton, M.R., Fekete, G., Hogg, D.W., O'Mullane, W., Szalay, A., & Wherry, N. 2004. *PASP*, **116**, 133
- Lynden-Bell, D. 1960. *MNRAS*, **120**, 204
- Lynden-Bell, D. 1962a. *MNRAS*, **123**, 447
- Lynden-Bell, D. 1962b. *MNRAS*, **124**, 279
- Lynden-Bell, D. 1967. *MNRAS*, **136**, 101
- Lynden-Bell, D. 1969. *MNRAS*, **144**, 189
- Lynden-Bell, D. 1979. *MNRAS*, **187**, 101
- Lynden-Bell, D. 1999. *Physica A*, **263**, 293
- Lynden-Bell, D., & Eggleton, P.P. 1980. *MNRAS*, **191**, 483
- Lynden-Bell, D., & Kalnajs, A.J. 1972. *MNRAS*, **157**, 1
- Lynden-Bell, D., & Lynden-Bell, R.M. 1977. *MNRAS*, **181**, 405
- Lynden-Bell, D., & Ostriker, J.P. 1967. *MNRAS*, **136**, 293
- Lynden-Bell, D., & Pringle, J.E. 1974. *MNRAS*, **168**, 603
- Lynden-Bell, D., & Wood, R. 1968. *MNRAS*, **138**, 495
- Lynds, R., & Toomre, A. 1976. *ApJ*, **209**, 382
- Lyttleton, R.A. 1953. *The Stability of Rotating Liquid Masses* (Cambridge: Cambridge University Press)
- Ma, J. 2002. *A&A*, **388**, 389
- Magorrian, J. 2006. *MNRAS*, **373**, 425
- Magorrian, J., & Tremaine, S. 1999. *MNRAS*, **309**, 447
- Magorrian, S.J. 2007. *astro-ph/0703406*
- Magorrian, S.J., & Binney, J. 1994. *MNRAS*, **271**, 949
- Majewski, S.R., Skrutskie, M.F., Weinberg, M.D., & Ostheimer, J.C. 2004. *ApJ*, **599**, 1082
- Makino, J. 1991. *ApJ*, **369**, 200
- Makino, J. 1996. *ApJ*, **471**, 796
- Makino, J. 2001. In *ASP Conference Series 228, Dynamics of Star Clusters and the Milky Way*, ed. S. Deiters, B. Fuchs, A. Just, R. Spurzem, & R. Wielen (San Francisco: Astronomical Society of the Pacific), 87
- Makino, J., & Aarseth, S.J. 1992. *PASJ*, **44**, 151
- Makino, J., Fukushige, T., Koga, M., & Namura, K. 2003. *PASJ*, **55**, 1163
- Makino, J., & Funato, Y. 2004. *ApJ*, **602**, 93
- Malin, D. 1993. *A View of the Universe* (Cambridge, MA: Sky Publishing)
- Maoz, E. 1991. *ApJ*, **375**, 687
- Maoz, E. 1995. *ApJ*, **447**, L91
- Mardling, R. 1995. *ApJ*, **450**, 722
- Mardling, R., & Aarseth, S.J. 2001. *MNRAS*, **321**, 398
- Mark, J.W.-K. 1971. *ApJ*, **169**, 455
- Mark, J.W.-K. 1974. *ApJ*, **193**, 539

- Mark, J.W.-K. 1976. *ApJ*, **205**, 363
- Marochnik, L.S. 1967. *AZh*, **44**, 5 (in Russian). Also *SvA*, **11**, 873 (in English)
- Marochnik, L.S., & Suchkov, A.A. 1996. *The Milky Way Galaxy* (Amsterdam: Gordon and Breach)
- Mateo, M.L. 1998. *ARA&A*, **36**, 435
- Mathews, W.G., & Brighenti, F. 2003. *ARA&A*, **41**, 191
- Mathur, S. 1990. *MNRAS*, **243**, 529
- McGaugh, S.S., Barker, M.K., & de Blok, W.J.G. 2003. *ApJ*, **584**, 566
- McGlynn, T.A. 1984. *ApJ*, **281**, 13
- Merrifield, M.R. 2003. In *ASP Conference Series 317, Milky Way Surveys: the Structure and Evolution of our Galaxy*, ed. D. Clemens, R.Y. Shah, & T. Brainerd (San Francisco: Astronomical Society of the Pacific), 289
- Merrifield, M.R., & Kuijken, K. 1995. *MNRAS*, **274**, 933
- Merritt, D. 1985. *AJ*, **90**, 1027 and *MNRAS*, **214**, 25p
- Merritt, D. 1999. *PASP*, **111**, 129
- Merritt, D., & Aguilar, L.A. 1985. *MNRAS*, **217**, 787
- Merritt, D., & Fridman, T. 1996. *ApJ*, **460**, 136
- Merritt, D., & Quinlan, G.D. 1998. *ApJ*, **498**, 625
- Merritt, D., & Sellwood, J.A. 1994. *ApJ*, **425**, 551
- Merritt, D., & Valluri, M. 1999. *AJ*, **118**, 1177
- Mestel, L. 1963. *MNRAS*, **126**, 553
- Meza, A., & Zamorano, N. 1997. *ApJ*, **490**, 136
- Michie, R.W. 1963. *MNRAS*, **125**, 127
- Michie, R.W., & Bodenheimer, P.H. 1963. *MNRAS*, **126**, 269
- Mihos, J.C., & Hernquist, L. 1996. *ApJ*, **464**, 641
- Mikkola, S. 1997. *CeMDA*, **68**, 87
- Mikkola, S., & Aarseth, S.J. 1996. *CeMDA*, **64**, 197
- Mikkola, S., & Aarseth, S.J. 1998. *NewA*, **3**, 309
- Miller, R.H. 1964. *ApJ*, **140**, 250
- Miller, R.H. 1971. In *IAU Colloquium 10, Gravitational N-Body Problem*, ed. M. Lecar (Dordrecht: Reidel), 213. Also *Ap&SS*, **14**, 73
- Miller, R.H., & Smith, B.F. 1979. *ApJ*, **227**, 785
- Miyamoto, M., & Nagai, R. 1975. *PASJ*, **27**, 533
- Mohr, P.J., & Taylor, B.N. 2005. *Rev. Mod. Phys.*, **77**, 1. See also physics.nist.gov/cuu/Constants/index.html
- Moore, B., Kazantzidis, S., Diemand, J., & Stadel, J. 2004. *MNRAS*, **354**, 522
- Moore, B., Quinn, T., Governato, F., Stadel, J., & Lake, G. 1999. *MNRAS*, **310**, 1147
- Morse, P.M., & Feshbach, H. 1953. *Methods of Theoretical Physics* (New York: McGraw-Hill)
- Mueller, M.W., & Arnett, W.D. 1976. *ApJ*, **210**, 670
- Mulchaey, J.S., Dressler, A., & Oemler, A., ed. 2004. *Carnegie Observatories Astrophysics Series, 3, Clusters of Galaxies: Probes of Cosmological Structure and Galaxy Evolution* (Cambridge: Cambridge University Press)
- Mulder, W.A. 1983. *A&A*, **117**, 9
- Murai, T., & Fujimoto, M. 1980. *PASJ*, **32**, 581
- Murali, C., & Weinberg, M.D. 1997a. *MNRAS*, **288**, 749
- Murali, C., & Weinberg, M.D. 1997b. *MNRAS*, **291**, 717
- Murray, C.D., & Dermott, S.F. 1999. *Solar System Dynamics* (Cambridge: Cambridge University Press)
- Naab, T., Jesseit, R., & Burkert, A. 2006. *MNRAS*, **372**, 839
- Nagai, R., & Miyamoto, M. 1976. *PASJ*, **28**, 1
- Nakano, T., & Makino, J. 1999. *ApJ*, **525**, L77
- Navarro, J.F., Frenk, C.S., & White, S.D.M. 1995. *MNRAS*, **275**, 720
- Navarro, J.F., Frenk, C.S., & White, S.D.M. 1996. *ApJ*, **462**, 563
- Navarro, J.F., Frenk, C.S., & White, S.D.M. 1997. *ApJ*, **490**, 493
- Nelson, R.W., & Tremaine, S. 1996. In *Gravitational Dynamics*, ed. O. Lahav, E. Terlevich, & R.J. Terlevich (Cambridge: Cambridge University Press), 73
- Nelson, R.W., & Tremaine, S. 1999. *MNRAS*, **306**, 1
- Nishikawa, K., & Wakatani, M. 2000. *Plasma Physics* (3rd ed.; Berlin: Springer)
- Noguchi, M. 1988. *A&A*, **203**, 259
- Nordström, B., et al. 2004. *A&A*, **418**, 989
- Norman, C.A., May, A., & van Albada, T.S. 1985. *ApJ*, **296**, 20
- Norman, C.A., Sellwood, J.A., & Hasan, H. 1996. *ApJ*, **462**, 124
- Novikov, I. D., & Frolov, V. P. 1989. *Physics of Black Holes* (Dordrecht: Kluwer)
- O'Neill, J.K., & Dubinski, J. 2003. *MNRAS*, **346**, 251
- Odenkirchen, M., et al. 2003. *AJ*, **126**, 2385
- Omma, H., & Binney, J. 2004. *MNRAS*, **350**, L13
- Oort, J.H. 1932. *BAN*, **6**, 249

- Oort, J.H. 1962. In *Interstellar Matter in Galaxies*, ed. L. Woltjer (New York: Benjamin)
- Oort, J.H. 1970. In *IAU Symposium 38, The Spiral Structure of Our Galaxy*, ed. W. Becker & G. Contopoulos (Dordrecht: Reidel), 1
- Öpik, E. 1932. *Proc. Am. Acad. Arts. Sci.*, **67**, 169
- Osipkov, L.P. 1979. *PAZh*, **5**, 77 (in Russian). See also *SvAL*, **4**, 42 (in English)
- Ostriker, E.C., & Binney, J.J. 1989. *MNRAS*, **237**, 785
- Ostriker, J.P., & Peebles, P.J.E. 1973. *ApJ*, **186**, 467
- Ostriker, J.P., Spitzer, L., & Chevalier, R.A. 1972. *ApJ*, **176**, L51
- Padmanabhan, T. 1990. *Phys. Rep.*, **188**, 285
- Palmer, P.L. 1994. *Stability of Collisionless Stellar Systems* (Dordrecht: Kluwer)
- Palmer, P.L., & Papaloizou, J. 1987. *MNRAS*, **224**, 1043
- Palmer, P.L., & Papaloizou, J. 1988. *MNRAS*, **231**, 935
- Palunas, P., & Williams, T.B. 2000. *AJ*, **120**, 2884
- Parker, E.N. 1966. *ApJ*, **145**, 811
- Pasha, I.I. 1985. *SvAL*, **11**, 1
- Pasha, I.I. 2002. *Istoriko-Astronomicheskie Issledovaniya*, **27**, 102 (in Russian). See also arXiv.org/abs/astro-ph/0406142 (in English)
- Pasha, I.I. 2004. *Istoriko-Astronomicheskie Issledovaniya*, **29**, 8 (in Russian). See also arXiv.org/abs/astro-ph/0406143 (in English)
- Pathria, R.K. 1972. *Statistical Mechanics* (Oxford: Pergamon)
- Patsis, P.A., Skokos, C., & Athanassoula, E. 2003. *MNRAS*, **346**, 1031
- Peacock, J.A. 1999. *Cosmological Physics* (Cambridge: Cambridge University Press)
- Pedersen, K., Rasmussen, J., Sommer-Larsen, J., Toft, S., Benson, A.J., & Bower, R.G. 2006. *NewA*, **11**, 465
- Peebles, P.J.E. 1993. *Principles of Physical Cosmology* (Princeton: Princeton University Press)
- Peebles, P.J.E. 1996. In *Gravitational Dynamics*, ed. O. Lahav, E. Terlevich, & R.J. Terlevich (Cambridge: Cambridge University Press), 219
- Pérez-Ramírez, D., Knapen, J.H., Peletier, R.F., Laine, S., Doyon, R., & Nadeau, D. 2000. *MNRAS*, **317**, 234
- Perryman, M.A.C., et al. 1995. *A&A*, **304**, 69
- Peters, P.C. 1964. *Phys. Rev.*, **136**, B1224
- Petit, J.-M., & Hénon, M. 1986. *Icar*, **66**, 536
- Pettini, M., Shapley, A.E., Steidel, C.C., Cuby, J.-G., Dickinson, M., Moorwood, A.F.M., Adelberger, K.L., & Giavalisco, M. 2001. *ApJ*, **554**, 981
- Pfenniger, D. 1984. *A&A*, **141**, 171
- Pfenniger, D., & Friedli, D. 1991. *A&A*, **252**, 75
- Phinney, E.S. 1994. In *Mass-Transfer Induced Activity in Galaxies*, ed. I. Shlosman (Cambridge: Cambridge University Press), 1
- Phinney, E.S. 1996. In *ASP Conference Series 90, The Origins, Evolution, and Destinies of Binary Stars in Clusters*, ed. E.F. Milone & J.-C. Mermilliod (San Francisco: Astronomical Society of the Pacific), 163
- Pichon, C., & Cannon, R.C. 1997. *MNRAS*, **291**, 616
- Piskunov, A.E., Schilbach, E., Kharchenko, N.V., Röser, S., & Scholz, R.-D. 2007. *A&A*, **468**, 151
- Plummer, H.C. 1911. *MNRAS*, **71**, 460
- Polyachenko, V.L. 1977. *PAZh*, **3**, 99 (in Russian). Also *SvAL*, **3**, 51 (in English)
- Polyachenko, V.L. 1981. *PAZh*, **7**, 142 (in Russian). Also *SvAL*, **7**, 79 (in English)
- Polyachenko, V.L., & Shukhman, I.G. 1973. *AZh*, **50**, 721 (in Russian). Also *SvA*, **17**, 460 (in English)
- Polyachenko, V.L., & Shukhman, I.G. 1981. *AZh*, **58**, 933 (in Russian). Also *SvA*, **25**, 533 (in English)
- Pooley, D., & Hut, P. 2006. *ApJ*, **646**, L143
- Prada, F., et al. 2003. *ApJ*, **598**, 260
- Prendergast, K.H., & Tomer, E. 1970. *AJ*, **75**, 674
- Press, W.H., Flannery, B.P., Teukolsky, S.A., & Vetterling, W.T. 1986. *Numerical Recipes* (New York: Cambridge University Press).
- Press, W.H., & Schechter, P. 1974. *ApJ*, **187**, 425
- Putman, M.E., Staveley-Smith, L., Freeman, K.C., Gibson, B.K., & Barnes, D.G. 2003. *ApJ*, **586**, 170
- Qian, E. 1992. *MNRAS*, **257**, 581
- Qian, E., de Zeeuw, P.T., van der Marel, R.P., & Hunter, C. 1995. *MNRAS*, **274**, 602
- Quinlan, G. 1996a. *NewA*, **1**, 255
- Quinlan, G. 1996b. *NewA*, **1**, 35
- Quinlan, G., Hernquist, L., & Sigurdsson, S. 1995. *ApJ*, **440**, 554
- Quinn, P.J. 1984. *ApJ*, **279**, 596
- Quinn, T., & Binney, J. 1992. *MNRAS*, **255**, 729
- Rafikov, R.R. 2001. *MNRAS*, **323**, 445
- Raha, N., Sellwood, J.A., James, R.A., & Kahn, F.D. 1991. *Nat*, **352**, 411
- Rand, R.J., & Kulkarni, S.R. 1990. *ApJ*, **349**, L43

- Rand, R.J., & Wallin, J.F. 2004. *ApJ*, **614**, 157
- Regan, M.W., et al. 2001. *ApJ*, **561**, 218
- Reid, M.J., & Brunthaler, A. 2004. *ApJ*, **616**, 872
- Rein, G. 2007. In *Handbook of Differential Equations: Evolutionary Equations 3*, ed. C.M. Dafermos & E. Feireisl (Amsterdam: Elsevier)
- Rhoads, J.E. 1998. *AJ*, **115**, 472
- Richstone, D.O. 1976. *ApJ*, **204**, 642
- Richstone, D.O. 1982. *ApJ*, **252**, 496
- Richstone, D.O., & Potter, M.D. 1982. *ApJ*, **254**, 451
- Richstone, D.O., & Tremaine, S. 1985. *ApJ*, **296**, 370
- Riess, A., et al. 2004. *ApJ*, **607**, 665
- Rix, H.-W., & Rieke, M.J. 1993. *ApJ*, **418**, 123
- Rix, H.-W., & Zaritsky, D. 1995. *ApJ*, **447**, 82
- Roberts, J.A.G., & Quispel, G.R.W. 1992. *Phys. Rep.*, **216**, 63
- Roberts, P.H. 1962. *ApJ*, **136**, 1108
- Romanowsky, A.J., Douglas, N.G., Arnaboldi, M., Kuijken, K., Merrifield, M.R., Napolitano, N.R., Capaccioli, M., & Freeman, K.C. 2003. *Sci*, **301**, 1696
- Rood, H.J., Page, T.L., Kintner, E.C., & King, I.R. 1972. *ApJ*, **175**, 627
- Rosenbluth, M.N., MacDonald, W.M., & Judd, D.L. 1957. *Phys. Rev.*, **107**, 1
- Rots, A.H., Bosma, A., van der Hulst, J.M., Athanassoula, E., & Crane, P.C. 1990. *AJ*, **100**, 387
- Rowley, G. 1988. *ApJ*, **331**, 124
- Rubenstein, E.P., & Bailyn, C.D. 1997. *ApJ*, **474**, 701
- Sackett, P.D. 1997. *ApJ*, **483**, 103
- Safronov, V.S. 1960. *AnAp*, **23**, 979
- Saha, P. 1991. *MNRAS*, **248**, 494
- Saha, P. 1992. *MNRAS*, **254**, 132
- Saha, P. 1993. *MNRAS*, **262**, 1062
- Saha, P. 2003. *Principles of Data Analysis* (Great Malvern: Capella Archive)
- Sakai, S., et al. 2000. *ApJ*, **529**, 698
- Salo, H., & Laurikainen, E. 1993. *ApJ*, **410**, 586
- Salpeter E.E. 1955. *ApJ*, **121**, 161
- Sancisi, R. 2004. In *IAU Symposium 220, Dark Matter in Galaxies*, ed. S.D. Ryder, D.J. Pisano, M.A. Walker, & K.C. Freeman (San Francisco: Astronomical Society of the Pacific), 233
- Sand, D.J., Treu, T., Smith, G.P., & Ellis, R.S. 2004. *ApJ*, **604**, 88
- Sandage, A., & Bedke, J. 1994. *The Carnegie Atlas of Galaxies* (Washington: Carnegie Institution)
- Sanders, D.B., & Mirabel, I.F. 1996. *ARA&A*, **34**, 749
- Sanders, R.H., & Tubbs, A.D. 1980. *ApJ*, **235**, 803
- Sargent, W.L.W., Young, P.J., Bokserberg, A., Shortridge, K., Lynds, C.R., & Hartwick, F.D.A. 1978. *ApJ*, **221**, 731
- Satoh, C. 1980. *PASJ*, **32**, 41
- Sawamura, M. 1988. *PASJ*, **40**, 279
- Schneider, P. 2006. In *Gravitational Lensing: Strong, Weak, and Micro*, ed. G. Meylan, P. Jetzer, & P. North (Berlin: Springer), 273
- Schneider, P., Ehlers, J., & Falco, E.E. 1999. *Gravitational Lenses* (Berlin: Springer)
- Schödel, R., et al. 2002. *Nat*, **419**, 694
- Schuster, A. 1883. *British Assoc. Rep.*, **427**
- Schwarz, M.P. 1981. *ApJ*, **247**, 77
- Schwarzschild, K. 1907. *Göttingen Nachr.*, **614**
- Schwarzschild, M. 1979. *ApJ*, **232**, 236
- Schweizer, F. 1976. *ApJS*, **31**, 313
- Seljak, U., & Zaldarriaga, M. 1996. *ApJ*, **469**, 437
- Sellwood, J.A. 1980. *A&A*, **89**, 296
- Sellwood, J.A. 1981. *A&A*, **99**, 362
- Sellwood, J.A. 1983. *J. Comput. Phys.*, **50**, 337
- Sellwood, J.A. 1985. *MNRAS*, **217**, 127
- Sellwood, J.A. 1999. In *ASP Conference Series 182, Galaxy Dynamics*, ed. D. Merritt, J.A. Sellwood, & M. Valluri (San Francisco: Astronomical Society of the Pacific), 351
- Sellwood, J.A. 2000a. *ApJ*, **540**, L1
- Sellwood, J.A. 2000b. In *ASP Conference Series 197, Dynamics of Galaxies: From the Early Universe to the Present*, ed. F. Combes, G.A. Mamon, & V. Charmandaris (San Francisco: Astronomical Society of the Pacific), 3
- Sellwood, J.A., & Athanassoula, E. 1986. *MNRAS*, **221**, 195
- Sellwood, J.A., & Binney, J.J. 2002. *MNRAS*, **336**, 785
- Sellwood, J.A., & Evans, N.W. 2001. *ApJ*, **546**, 176
- Sellwood, J.A., & McGaugh, S.S. 2005. *ApJ*, **634**, 70
- Sellwood, J.A., & Merritt, D. 1994. *ApJ*, **425**, 530
- Sellwood, J.A., & Moore, E.M. 1999. *ApJ*, **510**, 125
- Sellwood, J.A., & Sparke, L.S. 1988. *MNRAS*, **231**, 25p
- Sellwood, J.A., & Wilkinson, A. 1993. *Rep. Prog. Phys.*, **56**, 173
- Shapiro, S. L., & Lightman, A. P. 1976. *Nat*, **262**, 743
- Shen, J., & Sellwood, J.A. 2004. *ApJ*, **604**, 614

- Shen, J., & Sellwood, J.A. 2006. *MNRAS*, **370**, 2
- Sheth, R.K., & Lemson, G. 1999. *MNRAS*, **305**, 946
- Sheth, R.K., Mo, H.J., & Tormen, G. 2001. *MNRAS*, **323**, 1
- Sheth, R.K., & van de Weygaert, R. 2004. *MNRAS*, **350**, 517
- Shu, F.H. 1969. *ApJ*, **158**, 505
- Shu, F.H. 1970. *ApJ*, **160**, 99
- Shu, F.H. 1982. *The Physical Universe: an Introduction to Astronomy* (Mill Valley: University Science Books)
- Sigurdsson, S., & Phinney, E.S. 1993. *ApJ*, **415**, 631
- Silk, J. 1968. *ApJ*, **151**, 459
- Sills, A., Adams, T., Davies, M.B., & Bate, M.R. 2002. *MNRAS*, **332**, 49
- Skokos, C., Patsis, P.A., & Athanassoula, E. 2002. *MNRAS*, **333**, 847
- Smith, M.C., et al. 2007. *MNRAS*, **379**, 755
- Somerville, R.S., & Kolatt, T.S. 1999. *MNRAS*, **305**, 1
- Sparke, L.S. 1984. *ApJ*, **280**, 117
- Sparke, L.S., & Gallagher, J.S. 2000. *Galaxies in the Universe: an Introduction* (Cambridge: Cambridge University Press)
- Sparke, L.S., & Sellwood, J.A. 1987. *MNRAS*, **225**, 653
- Spergel, D.N., et al. 2007. *ApJS*, **170**, 377
- Spitzer, L. 1940. *MNRAS*, **100**, 396
- Spitzer, L. 1942. *ApJ*, **95**, 329
- Spitzer, L. 1956. *ApJ*, **124**, 20
- Spitzer, L. 1958. *ApJ*, **127**, 17
- Spitzer, L. 1969. *ApJ*, **158**, L139
- Spitzer, L. 1987. *Dynamical Evolution of Globular Clusters* (Princeton: Princeton University Press)
- Spitzer, L., & Baade, W. 1951. *ApJ*, **113**, 413
- Spitzer, L., & Schwarzschild, M. 1951. *ApJ*, **114**, 385
- Spitzer, L., & Schwarzschild, M. 1953. *ApJ*, **118**, 106
- Spitzer, L., & Shapiro, S.L. 1972. *ApJ*, **173**, 529
- Spurzem, R., & Takahashi, K. 1995. *MNRAS*, **272**, 772
- Sridhar, S., & Sambhus, N. 2003. *MNRAS*, **345**, 539
- Sridhar, S., & Touma, J. 1996. *MNRAS*, **279**, 1263
- Standish, E.M. 1995. In *Highlights of Astronomy 10*, ed. I. Appenzeller (Dordrecht: Kluwer), 180. See also ssd.jpl.nasa.gov
- Stark, A.A. 1977. *ApJ*, **213**, 368
- Statler, T. 1987. *ApJ*, **321**, 113
- Stiefel, E.L., & Schiefele, G. 1971. *Linear and Regular Celestial Mechanics* (Berlin: Springer)
- Stix, T.H. 1992. *Waves in Plasmas* (New York: American Institute of Physics)
- Stodólkiewicz, J.S. 1986. *AcA*, **36**, 19
- Strickland, D.K., Heckman, T.M., Colbert, E.J.M., Hoopes, C.G., & Weaver, K.A. 2004. *ApJ*, **606**, 829
- Strom, S.E., Jensen, E.B., & Strom, K.M. 1976. *ApJ*, **206**, L11
- Sugimoto, D., & Bettwieser, E. 1983. *MNRAS*, **204**, 19p
- Sumi, T., et al. 2003. *ApJ*, **591**, 204
- Summers, D., & Thorne, R.M. 1991. *Phys. Fluids B*, **3**, 1835
- Sussman, G.J., & Wisdom, J. 2001. *Structure and Interpretation of Classical Mechanics* (Cambridge, MA: MIT Press)
- Swaters, R.A., Madore, B.F., van den Bosch, F.C., & Balcells, M. 2003. *ApJ*, **583**, 732
- Sweet, P.A. 1963. *MNRAS*, **125**, 285
- Syer, D., & Tremaine, S. 1996. *MNRAS*, **282**, 223
- Sygné, J.F., Tagger, M., Athanassoula, E., & Pellat, R. 1988. *MNRAS*, **232**, 733
- Szebehely, V.G. 1967. *Theory of Orbits* (New York: Academic Press)
- Takahara, F. 1976. *Prog. Theor. Phys.*, **56**, 1665
- Takahashi, K. 1995. *PASJ*, **47**, 561
- Tassoul, J.L. 1978. *Theory of Rotating Stars* (Princeton: Princeton University Press)
- Tegmark, M., et al. 2004. *ApJ*, **606**, 702
- Teuben, P.J., & Sanders, R.H. 1985. *MNRAS*, **212**, 257
- Theuns, T. 1996. *MNRAS*, **279**, 827
- Thornley, M.D. 1996. *ApJ*, **469**, L45
- Tilanus, R.P.J., & Allen, R.J. 1989. *ApJ*, **339**, L57
- Tisserand, P., et al. 2007. *A&A*, **469**, 387
- Toomre, A. 1963. *ApJ*, **138**, 385
- Toomre, A. 1964. *ApJ*, **139**, 1217
- Toomre, A. 1966. In *Geophysical Fluid Dynamics, notes on the 1966 Summer Study Program at Woods Hole Oceanographic Institute*, ref. no. 66-46, 111
- Toomre, A. 1969. *ApJ*, **158**, 899
- Toomre, A. 1977a. *ARA&A*, **15**, 437
- Toomre, A. 1977b. In *The Evolution of Galaxies and Stellar Populations*, ed. B.M. Tinsley & R.B. Larson (New Haven: Yale University Observatory), 401
- Toomre, A. 1978. In *IAU Symposium 79, The Large Scale Structure of the Universe*, ed. M.S. Longair & J. Einasto (Dordrecht: Reidel), 109
- Toomre, A. 1981. In *The Structure and Evolution of Normal Galaxies*, ed. S.M. Fall & D. Lynden-Bell (Cambridge: Cambridge University Press), 111
- Toomre, A. 1982. *ApJ*, **259**, 535

- Toomre, A. 1983. In IAU Symposium 100, Internal Kinematics and Dynamics of Galaxies, ed. E. Athanassoula (Dordrecht: Reidel), 177
- Toomre, A., & Kalnajs, A.J. 1991. In Dynamics of Disc Galaxies, ed. B. Sundelius (Göteborg: Dept. of Astronomy and Astrophysics, Göteborgs University and Chalmers University of Technology), 341
- Toomre, A., & Toomre, J. 1972. *ApJ*, **178**, 623
- Touma, J., & Tremaine, S. 1997. *MNRAS*, **292**, 905
- Tremaine, S. 1975. Ph.D. thesis, Princeton University
- Tremaine, S. 1976a. *ApJ*, **203**, 345
- Tremaine, S. 1976b. *MNRAS*, **175**, 557
- Tremaine, S. 1981. In The Structure and Evolution of Normal Galaxies, ed. S.M. Fall & D. Lynden-Bell (Cambridge: Cambridge University Press), 67
- Tremaine, S. 1999. *MNRAS*, **307**, 877
- Tremaine, S., et al. 2002. *ApJ*, **574**, 740
- Tremaine, S., Hénon, M., & Lynden-Bell, D. 1986. *MNRAS*, **219**, 285
- Tremaine, S., Richstone, D.O., Byun, Y.-I., Dressler, A., Faber, S.M., Grillmair, C., Kormendy, J., & Lauer, T.R. 1994. *AJ*, **107**, 634
- Tremaine, S., & Weinberg, M.D. 1984a. *ApJ*, **282**, L5
- Tremaine, S., & Weinberg, M.D. 1984b. *MNRAS*, **209**, 729
- Tremaine, S., & Yu, Q. 2000. *MNRAS*, **319**, 1
- Valluri, M., & Merritt, D. 2000. In The Chaotic Universe, ed. V.G. Gurzadyan & R. Ruffini (Singapore: World Scientific), 229
- Valtonen, M., & Karttunen, H. 2006. The Three-Body Problem (Cambridge: Cambridge University Press)
- Valtonen, M., & Mikkola, S. 1991. *ARA&A*, **29**, 9
- van Albada, T.S. 1982. *MNRAS*, **201**, 939
- van Albada, T.S., Bahcall, J.N., Begeman, K., & Sancisi, R. 1985. *ApJ*, **295**, 305
- van de Ven, G., Hunter, C., Verolme, E.K., & de Zeeuw, P.T. 2003. *MNRAS*, **342**, 1056
- van den Bergh, S. 2000. The Galaxies of the Local Group (Cambridge: Cambridge University Press)
- van der Marel, R.P. 1999. *AJ*, **117**, 744
- van der Marel, R.P., Alves, D.R., Hardy, E., & Suntzeff, N.B. 2002. *AJ*, **124**, 2639
- van der Marel, R.P., Cretton, N., de Zeeuw, P.T., & Rix, H.-W. 1998. *ApJ*, **493**, 613
- van der Marel, R.P., & Franx, M. 1993. *ApJ*, **407**, 525
- van Gorkom, J.H. 2004. In Clusters of Galaxies: Probes of Cosmological Structure and Galaxy Evolution, ed. J.S. Mulchaey, A. Dressler, & A. Oemler (Cambridge: Cambridge University Press), 305
- van Kampen, N.G. 1955. *Physica*, **21**, 949
- Vandervoort, P.O. 2003. *MNRAS*, **339**, 537
- Vauterin, P., & Dejonghe, H. 1996. *A&A*, **313**, 465
- Velázquez, H., & White, S.D.M. 1999. *MNRAS*, **304**, 254
- Visser, H.C.D. 1980. *A&A*, **88**, 149,159
- Walker, I.R., Mihos, J.C., & Hernquist, L. 1996. *ApJ*, **460**, 121
- Warren, M.S., Quinn, P.J., Salmon, J.K., & Zurek, W.H. 1992. *ApJ*, **399**, 405
- Wasserman, I., & Weinberg, M.D. 1987. *ApJ*, **312**, 390
- Watson, G.N. 1995. A Treatise on the Theory of Bessel Functions (2nd ed.; Cambridge: Cambridge University Press)
- Weinberg, M.D. 1983. *ApJ*, **271**, 595
- Weinberg, M.D. 1985. *MNRAS*, **213**, 451
- Weinberg, M.D. 1986. *ApJ*, **300**, 93
- Weinberg, M.D. 1989. *MNRAS*, **239**, 549
- Weinberg, M.D. 1991a. *ApJ*, **368**, 66
- Weinberg, M.D. 1991b. *ApJ*, **373**, 391
- Weinberg, M.D. 1993. *ApJ*, **410**, 543
- Weinberg, M.D. 1994a. *AJ*, **108**, 1398,1403, 1414
- Weinberg, M.D. 1994b. *ApJ*, **421**, 481
- Weinberg, M.D. 1998. *MNRAS*, **299**, 499
- Weinberg, S. 1972. Gravitation and Cosmology (New York: Wiley)
- Weiner, B.J., & Sellwood, J.A. 1999. *ApJ*, **524**, 112
- Weiner, B.J., Sellwood, J.A., & Williams, T.B. 2001. *ApJ*, **546**, 931
- Whitham, G.B. 1974. Linear and Nonlinear Waves (New York: Wiley)
- Whitmore, B.C., & Schweizer, F. 1995. *AJ*, **109**, 960
- Wilczynski, E.J. 1896. *ApJ*, **4**, 97
- Wilkinson, A., & James, R.A. 1982. *MNRAS*, **199**, 171
- Wilkinson, M.I., & Evans, N.W. 1999. *MNRAS*, **310**, 645
- Willick, J.A., Strauss, M.A., Dekel, A., & Kolatt, T. 1997. *ApJ*, **486**, 629
- Wilson, C.P. 1975. *AJ*, **80**, 175
- Woolley, R., & Dickens, R.J. 1961. *Roy. Greenwich Obs. Bulletin*, **42**, 291
- Yamashiro, T., Gouda, N., & Sakagami, M. 1992. *Prog. Theor. Phys.*, **88**, 269
- Yamauchi, C., & Goto, T. 2004. *MNRAS*, **352**, 815
- Yao, W.-M., et al. 2006. *J. Phys. G*, **33**, 1. See also pdg.lbl.gov
- Yoo, J., Chanamé, J., & Gould, A. 2004. *ApJ*, **601**, 311

- Yoshida, H. 1982. *CeMec*, **28**, 239
Yoshida, H. 1993. *CeMDA*, **56**, 27
Young, P. 1980. *ApJ*, **242**, 1232
Yu, Q. 2002. *MNRAS*, **331**, 935
Yun, M.S., Ho, P.T.P., & Lo, K.Y. 1994. *Nat.*, **372**, 530
Zang, T.A. 1976. Ph.D. thesis, Massachusetts Institute of Technology
Zang, T.A., & Hohl, F. 1978. *ApJ*, **226**, 521
Zhao, D.H., Jing, Y.P, Mo, H.J., & Börner, G. 2003. *ApJ*, **597**, L9
Zhao, H. 1996. *MNRAS*, **283**, 149
Zinn, R. 1985. *ApJ*, **293**, 424
Zoccali, M., et al. 2004. *A&A*, **399**, 931
Zwicky, F. 1933. *Helv. Phys. Acta*, **6**, 110
Zwicky, F. 1955. *PASP*, **67**, 232