

ASTR 635: Exoplanetary Astrophysics

Course Notes

Spring 2024

University of Maryland, College Park

Assistant Professor Tad Komacek

Contents

1 Syllabus Overview, Stars Activity, Formation of the Solar System, Orbit Primer	1
1.1 Course Overview	1
1.2 Stellar physics and radiation fundamentals activity	2
1.3 Formation of the Solar System recap	2
1.3.1 Formation of the Solar System activity	5
1.4 Elliptical orbits primer	5
2 Detecting exoplanets: radial velocity	8
2.1 Radial velocity: notes	8
2.1.1 Doppler shift	8
2.1.2 Radial velocity equation for circular orbits	9
2.1.3 Group activity: deriving the radial velocity semi-amplitude for circular orbits	10
2.1.4 Detecting planets via radial velocity in practice	11
2.2 Radial velocity in practice: group activity	13
3 Detecting exoplanets: astrometry	15
3.1 Astrometry: notes	15
3.1.1 Method, historical and modern observations	15
3.1.2 Astrometric wobble due to a companion planet	18
3.2 Astrometry: group activity	20
4 Detecting exoplanets: transits	21
4.1 Transits: notes	21
4.1.1 Transit depth, probability, and duration	21
4.1.2 Transit geometry, impact parameter	23
4.1.3 Measuring stellar density via transits	25
4.1.4 Transit method in practice	25
4.2 Transits: group activities	27
4.2.1 Calculating transit depth and probability	27
4.2.2 Drawing transits	28
5 Detecting exoplanets: timing	29
5.1 Timing: notes	29
5.1.1 Transit timing variations	29
5.1.2 Principles of detecting planets via timing	32

5.2	Finding planets via pulsar timing: group activity	33
6	Detecting exoplanets: microlensing	34
6.1	Microlensing: notes	34
6.1.1	Lens solution, Einstein radius	34
6.1.2	Peak magnification	36
6.1.3	Planetary perturbation	37
6.1.4	Event length	39
6.1.5	Microlensing in practice	39
7	Detecting exoplanets: direct imaging	41
7.1	Direct imaging intro activity	41
7.2	Direct imaging: notes	41
7.2.1	Planet-star contrast	41
7.2.2	Technological challenges	43
8	Detecting exoplanets: inter-comparison of detection techniques (Day 8)	46
8.1	Activity: Wright & Gaudi for the modern day	46
8.2	Strengths and biases of each detection method	47
8.2.1	Radial velocity	48
8.2.2	Transit	48
8.2.3	Direct imaging	49
8.2.4	Microlensing	49
8.2.5	Astrometry	49
8.3	Key findings from each detection method	50
8.3.1	Radial velocity	52
8.3.2	Transit	53
8.3.3	Direct imaging	54
8.3.4	Microlensing	54
9	Detecting exoplanets: occurrence rates	57
9.1	Occurrence rates	57
9.1.1	General principles	57
9.1.2	Early results	58
9.1.3	Example of deriving occurrence rates: Radius gap	59
10	Planet formation: disk structure	64
10.1	Vertical disk structure	64
10.1.1	Hydrostatic equilibrium	64
10.1.2	Activity: Derive disk density profiles in small groups!	66
10.1.3	Disk flaring	66
10.2	Activity: estimating disk temperatures	66

11 Planet formation: disk thermal structure, dynamics	68
11.1 Disk thermal structure	68
11.1.1 Flared disks	68
11.2 Disk temperature activity: condensation points and ice lines	69
11.3 Disk dynamics	70
11.4 Momentum balance	70
11.4.1 Effect of viscosity	71
11.4.2 Viscosity activity	71
11.4.3 Shakura-Sunyaev disks	72
12 Planet formation: dust and pebble motions	73
12.1 Dust motions	73
12.1.1 Epstein and Stokes drag regimes	73
12.1.2 Dust coupling and settling	74
12.1.3 Radial drift: derivation activity	74
12.1.4 The “meter-size barrier”	76
13 Planet formation: from pebbles to planets	78
13.1 Radial drift activity	78
13.2 Accretion of planetesimals	79
13.2.1 Gravitational focusing	79
13.2.2 Hill radius	80
13.2.3 Isolation mass	80
14 Planet formation: accretion, orbital migration and evolution	82
14.1 Gravitational focusing activity	82
14.2 Steps in the formation of terrestrial planets	82
14.3 Formation of gas giant planets	83
14.3.1 Gravitational instability	83
14.3.2 Core accretion	84
14.4 Migration	85
14.4.1 Type I migration	85
14.4.2 Type II migration	86
14.4.3 Planetesimal disk migration	87
14.5 Models for Solar System evolution	87
15 Exoplanet atmospheres: structure, composition, chemistry, loss	89
15.1 Hydrostatic equilibrium	89
15.2 Atmospheric thermodynamics	90
15.2.1 First law of thermodynamics	90
15.3 Specific heats	91
15.3.1 Convective instability	92
15.3.2 Condensation, clouds and the moist adiabat	95
15.4 Radiative relaxation	96
15.4.1 Radiative timescale activity	97

15.5	Atmospheric composition	98
15.5.1	Compositional diversity	98
15.5.2	Equilibrium chemistry	99
15.5.3	Disequilibrium chemistry and mixing	101
15.6	Atmospheric loss	102
15.6.1	Energetic considerations	102
15.6.2	The cosmic shoreline	104
16	Exoplanet interiors: giant planets	105
16.1	Phases of H/He in giant planets	105
16.2	Interior structures of Solar System giant planets	107
16.3	Hydrostatic equilibrium	108
16.3.1	Central pressure activity	109
16.4	Equations of planetary structure	109
16.4.1	Heat transport in planetary interiors	110
16.4.2	Radius inflation of hot Jupiters	111
16.4.3	Radius evolution, Kelvin-Helmholtz Timescale	112
17	Planetary habitability	113
17.1	The habitable zone	113
17.1.1	Classic 1D framework, carbonate-silicate weathering	113
17.1.2	Clouds and 3D effects	115
17.2	Biosignatures	116
17.2.1	Oxygen and ozone	116
17.2.2	Disequilibrium due to life	117
17.2.3	Biosignature false positives	119
17.3	Discussion activity	121
17.4	Decadal Survey Recommendations	122
17.5	Habitable Worlds Observatory	124
17.5.1	Detecting a sample of potentially habitable ExoEarths with direct imaging	125
17.5.2	Characterizing ExoEarths: reflectance spectra, rotational mapping	128
17.6	Prediction activity!	129
18	Exoplanet characterization: transmission spectroscopy	130
18.1	Fundamentals of transmission spectroscopy	130
18.1.1	Qualitative description	130
18.1.2	Transmission flux ratio	130
18.1.3	Beer's law	131
18.1.4	Application to observed spectra: example of WASP-43b	133
18.2	Highlights of transmission spectroscopy	134
19	Exoplanet interiors: terrestrial planets	138
19.1	Earth's internal structure	138
19.2	Heat transfer via conduction	139

19.2.1	Cooling timescale of Earth activity	140
19.2.2	Historical background: Kelvin’s folly	141
19.3	Convective heat transport	141
19.3.1	Rayleigh-Bernard convection	141
19.3.2	Boundary layer convection	143
19.4	Rocky planet mass-radius relationships	145
20	Exoplanet characterization: emission spectroscopy	146
20.1	Secondary eclipse depth	146
20.2	Linking emission to thermal structure	147
20.2.1	Solutions of the radiative transfer equations	147
20.2.2	Photosphere pressure	148
20.2.3	Absorption and emission features	149
20.3	Emission spectra in practice: WASP-18b with JWST	150
20.4	Emission spectra activity!	151
21	Exoplanet characterization: phase curves	152
21.1	Phase curve fundamentals	152
21.1.1	Orbital phase curves: close-in exoplanets	152
21.1.2	Contribution from reflected light	154
21.1.3	Rotational phase curves: brown dwarfs, wide-separation giant planets	155
21.2	Phase curve theory for tidally locked exoplanets	156
21.2.1	A simple coupled scaling theory for heat transport and winds	156
21.2.2	Comparisons of this simple theory with observations	158
21.2.3	Activity: predicting the day-night contrasts and wind speeds of various exoplanets	161
22	Exoplanet characterization: Atmospheric dynamics	162
22.1	Scale analysis of the momentum equation and basic force balances	162
22.1.1	The momentum equation in Cartesian coordinates	162
22.1.2	Vertical force balance: hydrostatic equilibrium	162
22.1.3	Horizontal force balance: Rossby number, geostrophy	163
22.2	Python activity: Grid of hot Jupiter GCMs	165

1 Syllabus Overview, Stars Activity, Formation of the Solar System, Orbit Primer

Hi all! Welcome to ASTR 635. For our first class meeting, our agenda is the following:

1. Review the course format and structure (15 min)
2. Stars group activity (30 min)
3. Formation of the Solar System: brief recap and group activity (30 min)

I will make an agenda and notes handout for each class, listing the plans for the day as well as providing some further information that you can use for studying. These notes will not necessarily be comprehensive, and will instead be a summary of what is presented in class (sometimes it'll be the opposite, and they'll contain more information than we end up covering). I recommend both taking notes in class and taking notes from the assigned reading for each day. Note that the agenda is always tentative, and I will strive to go at the pace the class would like to rather than rushing through material.

For all future classes, there will be a corresponding reading and an ELMS assignment in which you must submit a question about the reading by 9 am on the day that class takes place. I will use the questions to better orient each lecture to the material that would be most impactful to our learning, and to answer specific thoughtful questions from each of you.

1.1 Course Overview

This course will provide an introduction to the current astrophysical study of exoplanets at the level to prepare undergraduate students to get involved in current research in the field. This course will survey the broad range of exoplanet science, and as such will be split into 3 parts with a coda:

1. Exoplanet detection methods (lectures 1-8).
2. Exoplanet demographics and planet formation (lectures 9-15)
3. Exoplanet atmospheres, interiors, and observational characterization (lectures 16-23).

Each part will end with a midterm exam – there will be no final exam. Instead, the final exam slot will consist of presentations as part of a final project, in which you will conduct novel research in the exoplanet field.

As discussed above, please do the assigned reading (see Table 1 in the Syllabus for assignments) before class and post a thoughtful question on the ELMS assignment for that day by 9 am. Class will include regular activities, usually focused on group problem solving to apply the knowledge we learn from each day's lecture. There will be three problem sets and three group mini-projects, one for each segment in the course triad. The grade distribution of the class is as follows:

Mid-term 1: 10%, **Mid-term 2:** 10%, **Mid-term 3:** 10%.

Problem sets: 15%.

Group mini-projects: 30%.

In-class participation: 5%. The lowest in-class assignment will be dropped.

Pre-class reading questions: 5%. The lowest pre-class assignment will be dropped.

Final project: 15%. The written component will comprise 2/3 of the project grade, and the oral presentation will comprise 1/3 of the project grade.

1.2 Stellar physics and radiation fundamentals activity

Exoplanet science relies on stellar physics, famously stated as “know thy star, know thy planet.” Let’s do a group activity both to get used to interacting in small groups in the classroom and to refresh our memory of fundamental stellar radiation. I’ll distribute markers and tabletop whiteboards for ya’ll to solve these problems on in groups of 2-3!

1. Consider a nearby M-dwarf star with an effective temperature of 2557 K, radius of 0.1234 Solar radii, and parallax of 0.0802 arcseconds.
 - (a) What is the spectral class of this star? Be as specific as possible.
 - (b) What is the bolometric luminosity of this star, in Solar luminosities (i.e., L/L_{\odot})? Note that $R_{\odot} = 6.96 \times 10^8$ m, $L_{\odot} = 3.83 \times 10^{26}$ W, and $T_{\text{eff},\odot} = 5777$ K. The Stefan-Boltzmann constant is 5.67×10^{-8} W m⁻² K⁻⁴.
 - (c) How far away is this star from Earth, in parsecs?
 - (d) At what wavelength does the blackbody spectrum of this star have its maximum? Note that for the Sun, $\lambda_{\text{max}} = 0.502$ μm .
 - (e) This star has seven nearly Earth-sized planets around it that were discovered in 2016 and 2017. What exoplanet system is this?
 - (f) Planet e in this system has a semi-major axis of 0.0293 au. Estimate the radiative equilibrium temperature of this planet, in Kelvin.
 - (g) If you wanted to observe the thermal emission of planet e, approximately what wavelength would you observe in? If you wanted to observe the transmission of light from the host star through the limb of planet e, at what approximate wavelength would you observe in? What current observational facilities might be suitable for each of these observations?

1.3 Formation of the Solar System recap

Our Milky Way Galaxy is teeming with planets. To date, astronomers have discovered 5,569 exoplanets (see Figure 1.1, value updated as of Jan. 12). These planets present an opportunity to understand how planetary systems form, determine whether our Solar System is special, and better understand the physics and chemistry that sets the present day state of planetary interiors, atmospheres, and surfaces. Over the next few weeks, we will first study the various methods by which exoplanets can be detected – each shown by the different colors in Figure 1.1. In today’s class, we will briefly recap the current picture for how planetary systems (including our Solar System) form before studying each detection method in more detail.

Mass – Period Distribution

10 Jan 2024
exoplanetarchive.ipac.caltech.edu

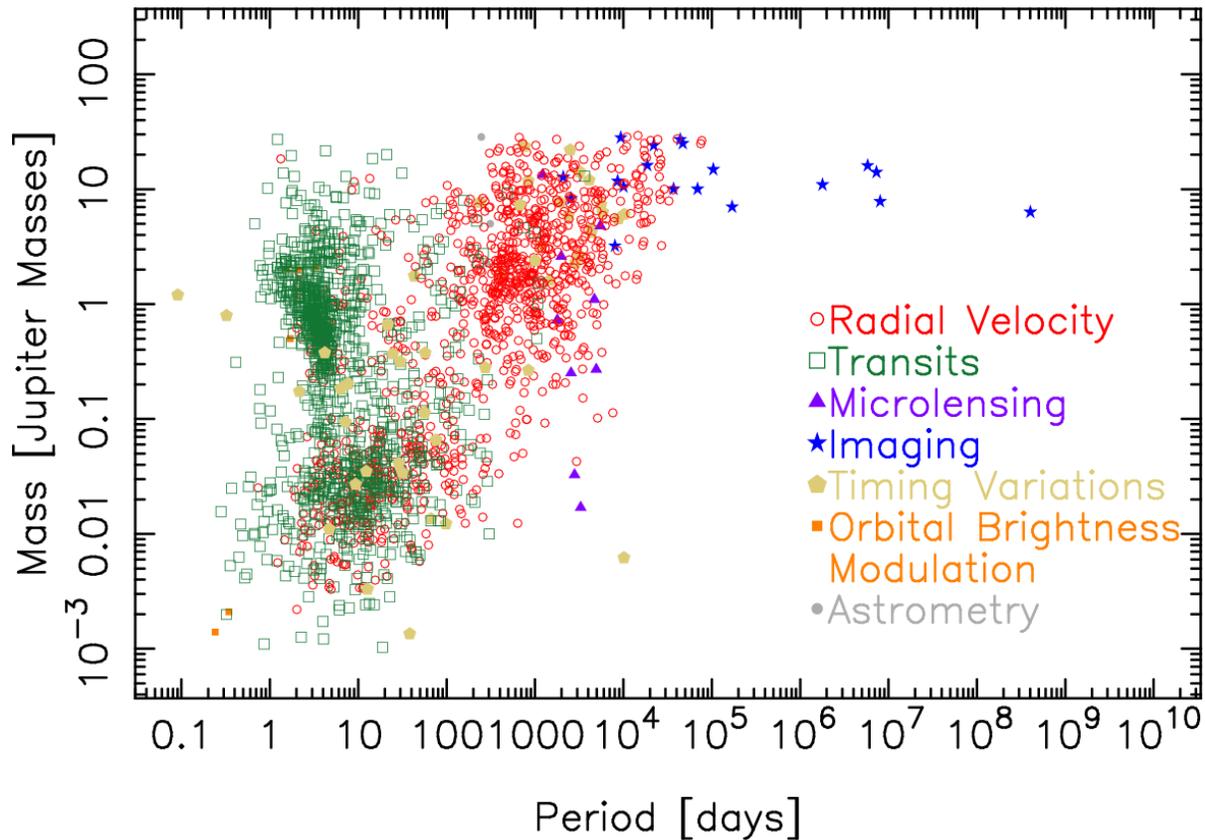


Figure 1.1: The exoplanet census as of January 2024. We'll get to all those detection methods (shown by the different colored points) in the next seven lectures.

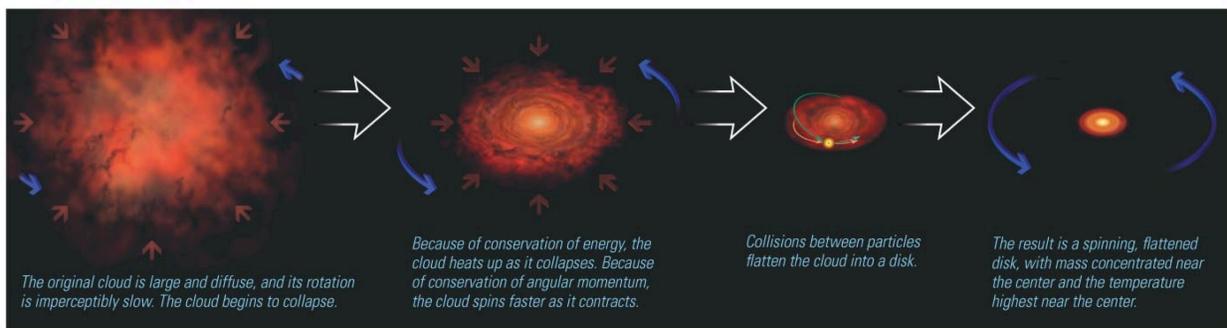


Figure 1.2: Artist's illustration of formation of a protoplanetary disk and nascent Solar System through collapse of a molecular cloud. Literally from my ASTR 100 slides.

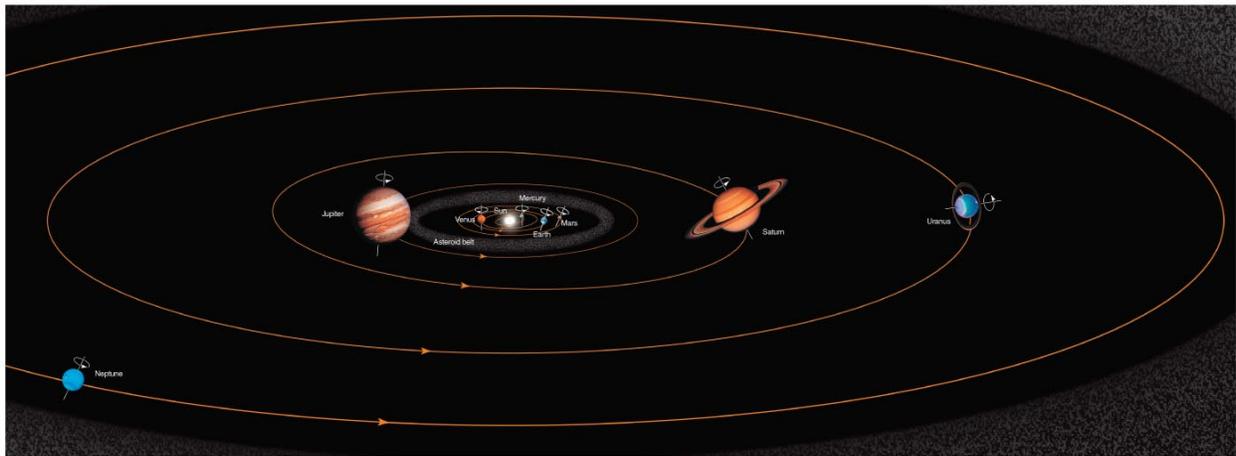
Planetary systems form due to the gravitational collapse of molecular cloud cores, through a step-wise process illustrated in Figure 1.2. This process begins when the molecular cloud becomes dense enough that it reaches a critical density (which can also be re-framed as a

mass or radius, often termed the Jeans mass or Jeans length) such that the internal thermal energy (i.e., gas pressure) is less than the gravitational potential energy and thus the cloud cannot remain in a state of hydrostatic equilibrium. We will derive this criterion when we cover planet formation in the second part of this course, but the Jeans mass can be approximated by (Carroll & Ostlie, 2017)

$$M_J \approx \left(\frac{5kT}{G\mu m_H} \right)^{3/2} \left(\frac{3}{4\pi\rho_0} \right)^{1/2}. \quad (1.1)$$

In Equation (1.1), k is the Boltzmann constant, T is temperature, G is the gravitational constant, μ is the mean molecular weight ($\mu \approx 2.3$ for a cold H/He mixture at Solar composition), m_H is the mass of hydrogen, and ρ_0 is the density of the cloud. The Jeans lengths of approximately Solar-mass clouds are on the order of tens of thousands of astronomical units, and we'll use this in our group activity.

As described in Figure 1.2, clouds spin up and flatten into disks as they collapse due to angular momentum conservation. Angular momentum conservation also implies that all of the planets should have the same sense of orbital revolution, that this should be the same direction as the Sun's rotation around its axis, and that the planets should further rotate around their own axes in the same direction. Figure 1.3 shows a diagram of the rotational and orbital spin vectors for each object in the Solar System. Indeed, as



© 2010 Pearson Education, Inc.

Figure 1.3: Diagram of Solar System (not to scale) showing rotational and orbital spin vectors for each object.

we expect from conservation of angular momentum the revolution of all planets around the Sun is all in the same direction (counter-clockwise), which is the same direction the Sun is rotating on its axis. This is also the same direction as the rotation of most planets around their axes, though the “obliquity,” or tilt of many planets with respect to their orbital plane around the Sun differs drastically (with Mercury having an obliquity of 0.034° and Uranus famous for being on its side, with an obliquity of 97.8° , see <https://nssdc.gsfc.nasa.gov/planetary/factsheet/>). However, Uranus and Venus (obliquity of 177.4°) both have obliquities that lie between 90° and 180° – this means that

their direction of rotation is actually *opposite* to their direction of orbital motion. This is not expected from conservation of angular momentum, and the root cause of the present-day obliquities of both Venus and Uranus is still an active area of research.

1.3.1 Formation of the Solar System activity

This group activity will provide you with some intuition for the partitioning of angular momentum in the present day Solar System, which is linked to that of its birth environment. We'll again use the portable white boards for this! Please get together in groups of 2-3.

Angular momentum in the Solar System is not equitably partitioned between our Sun and the planets. This has strong consequences for our understanding of the formation of our Solar System. The following questions will walk you through this.

1. Calculate the rotational angular momentum of the Sun, assuming it is a uniform density sphere (...this is not the case, but we just need order of magnitude accuracy). Note that the radius of the Sun $R_{\odot} = 6.96 \times 10^8$ m, the mass of the Sun $M_{\odot} = 1.989 \times 10^{30}$ kg, and the rotation period of the Sun $P_{\odot} = 24.5$ days.
2. Calculate the orbital angular momentum of Jupiter. First, write down or derive the orbital velocity of an object around the Sun as a function of the mass of the Sun M_{\odot} , gravitational constant G , and semi-major axis a . Then use the traditional formula for angular momentum to write an expression for the angular momentum of Jupiter's orbit around the Sun. Then plug in, using $M_{\text{Jup}} = 1.898 \times 10^{27}$ kg, $a_{\text{Jup}} = 5.204$ au.
3. (If time remains) Calculate the rotational angular momentum of a Solar-mass molecular cloud with a temperature of 10 K at the Jeans mass (we'll get to that later, for now I'll give you the relevant values). This cloud has a radius of $\approx 16,500$ au, and an angular velocity of $0.03 \text{ km s}^{-1} \text{ pc}^{-1}$ ($1 \text{ pc} = 3.086 \times 10^{13} \text{ km}$). As for the Sun, you can assume it's a uniform density sphere.
4. There are stark decreases in the amount of angular momentum from the molecular cloud, to the orbits of the planets, and further to the rotation of the Sun. Discuss with your peers where all of this angular momentum might have went, and what processes may have led to this "lost" angular momentum.

1.4 Elliptical orbits primer

Ch. 2 of the Exoplanet Handbook textbook and Ch. 1 of the Tremaine Dynamics of Planetary Systems textbook cover elliptical orbits in 3D, and this will be on your Problem Set 1. This section is a quick recap of the salient points that you can refer to for this class. We'll cover this in class if time permits and/or if people want me to (please tell me if you haven't seen this before), otherwise you'll cover it on your own in a fair bit of detail driving the radial velocity shift caused by a planet on an eccentric orbit in Problem Set 1.

Figure 1.4 shows the geometry of a 3D orbit in another planetary system, with the observer looking at the system from the top-down. The orbital plane can be defined with respect to an arbitrary reference plane – for exoplanet systems, we define the orbital plane with respect to the plane of the sky that is perpendicular to the direction toward the observer.

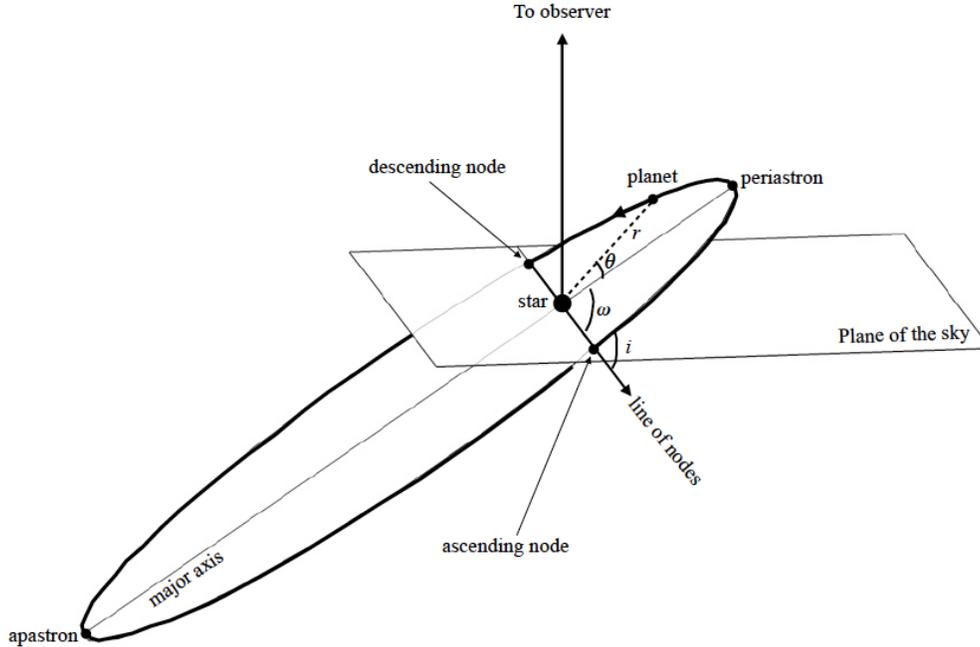


Figure 1.4: Geometry of a 3D orbit. Terms and variables are defined in the text and in Problem Set 1, Question 3. Figure courtesy Prof. Eliza Kempton.

Note that for the Solar System itself, the reference plane is usually different, instead often set to be the ecliptic (Earth’s orbital plane about the Sun).

The primary orbital elements are the eccentricity e , the semi-major axis a , and the inclination i . The semi-major axis a is half of the long axis of the ellipse. The eccentricity $e \equiv \sqrt{1 - b^2/a^2}$, where b is the semi-minor axis of the ellipse (i.e., half of the short axis of the ellipse). The inclination, i , of an orbit is the angle between the orbital plane and the reference plane. The inclination can range from $0^\circ - 180^\circ$, where inclinations from $0^\circ - 90^\circ$ correspond to prograde orbits (in the same direction as the rotation of the primary) while inclinations from $90^\circ - 180^\circ$ correspond to retrograde orbits (opposite to the rotation of the primary). The line of nodes is the intersection between the orbital and reference plane, and the point in the orbit where the planet passes upwards through the line of nodes is the ascending node while the point where the planet passes downwards is the descending node. The angle from a fixed zero point in the orbit to the ascending node is called the longitude of the ascending node, often denoted by Ω .

There are two other important angular orbital elements. The first is the angle between the line to the ascending node and the line toward periapsis, which is termed the argument of periapsis (ω). The second is the true anomaly, $\theta(t)$ in Figure 1.4 ($\nu(t)$ in the Exoplanet Handbook textbook), which is the angle between the periapsis and the planet’s actual time-dependent position as it orbits. Putting these together, we have six orbital elements that specify the location of a planet in its orbit: $a, e, i, \Omega, \omega, \theta$. In practice, astronomers often specify the location of the planet along its orbit using the mean motion

$$n \equiv \frac{2\pi}{P}, \quad (1.2)$$

where P is the orbital period, and the mean longitude (related to the mean anomaly)

$$\lambda = n(t - t_p) + \Omega + \omega = n(t - t_p) + \varpi, \quad (1.3)$$

where t_p is the time of pericenter passage and $\varpi = \Omega + \omega$.

The planet-star distance at any point in its orbit can be expressed as

$$r = \frac{a(1 - e^2)}{1 + e \cos\theta(t)}, \quad (1.4)$$

which is Kepler's first law. Further, note that Kepler's second law (equal areas in equal time, where $A = \pi a^2 \sqrt{1 - e^2}$) states

$$\frac{dA}{dt} = \frac{r^2 d\theta}{2 dt} = \text{constant}. \quad (1.5)$$

However, often the eccentric anomaly $E(t)$ is used rather than the true anomaly $\theta(t)$ to specify the location of the planet along its orbit. The eccentric anomaly is the angle inscribed within an auxiliary circle of the orbital ellipse (i.e., the circle that would be made if you take the radius of the circle to be the semi-major axis of the ellipse, see textbook Fig. 2.1). The true and eccentric anomalies are related as

$$\cos\theta(t) = \frac{\cos E(t) - e}{1 - e \cos E(t)}. \quad (1.6)$$

Further, the mean anomaly

$$M(t) \equiv n(t - t_p) = \lambda - \varpi \quad (1.7)$$

is related to the eccentric anomaly by Kepler's equation

$$M(t) = E(t) - e \sin E(t). \quad (1.8)$$

To practically calculate the position of a planet in its orbit, one can calculate the mean anomaly with Equation (1.7), then solve for the eccentric anomaly iteratively using Equation (1.8), and finally use Equation (1.6) to determine the true anomaly.

2 Detecting exoplanets: radial velocity

Our agenda for Day 2 is the following:

1. Review the concept of the Doppler shift (5 min)
2. One-slide intro to the radial velocity method (5 min)
3. Group activity: derive the radial velocity equation for circular orbits (20 min)
4. Learn in practice how astronomers detect planets via radial velocity (25 min)
5. Group activity: Calculate the radial velocity semi-amplitude due to Earth and Jupiter around the Sun and compare to current astronomical capabilities (20 min)

Today's reading is from our textbook, Ch. 2.1-2.4, Ch. 2.6-2.7 and/or from the Lovis & Fischer handout (which is from the Exoplanets book edited by S. Seager). These readings will cover the fundamentals of orbits (which we'll start discussing this class and build upon in following classes), the principles of radial velocity measurements, modern radial velocity instruments, and some current results for radial velocity observations.

2.1 Radial velocity: notes

2.1.1 Doppler shift

The standard way to measure motion along the line of sight in astronomy is by leveraging the change in the apparent wavelength of emitted light due to this motion. A change in the apparent wavelength due to motion along the line of sight is a fundamental property of waves (both transverse and longitudinal, covering light and sound alike), and this shift is termed the Doppler shift. Astronomers use the Doppler shift of stars induced by unseen planets to infer the presence of exoplanets via the “radial velocity” method.

For the purposes of measuring radial velocity of a star, we'll express the Doppler shift as a difference in the wavelength of a spectral line observed from the star relative to the wavelength of a spectral line emitted from that star, i.e., $\Delta\lambda = \lambda_{\text{obs}} - \lambda_{\text{rest}}$, where λ_{obs} is the observed wavelength of the spectral line and λ_{rest} is the wavelength at which that spectral line would lie if it were emitted at rest (i.e., in a laboratory).

The full (relativistic) Doppler shift is

$$\lambda_{\text{obs}} = \lambda_{\text{rest}} \frac{1 + \frac{v}{c} \cos\theta}{\sqrt{1 - \left(\frac{v}{c}\right)^2}}, \quad (2.1)$$

where v is the magnitude of velocity, c is the speed of light, and θ is the angle of motion relative to the line between the observer and star. For the purposes of detecting exoplanets, we can ignore relativistic effects ($v/c \ll 1$), which reduces the equation to

$$\lambda_{\text{obs}} = \lambda_{\text{rest}} \left(1 + \frac{v}{c} \cos\theta\right). \quad (2.2)$$

Rearranging, we can write the equation for radial velocity $v_r = v \cos\theta$ as

$$v_r = \frac{\Delta\lambda}{\lambda_{\text{rest}}} c. \quad (2.3)$$

As a result, to measure the motion of a star along the line of sight, in principle all we need to measure is the difference in position of a spectral line compared to what it would be in the laboratory ($\Delta\lambda$) and use Equation (2.3) to solve for v_r .

2.1.2 Radial velocity equation for circular orbits

Next we need to link the observable (radial velocity) to fundamental properties of the star-planet system. If you’ve ever learned the fundamental physics of binary star systems, exoplanet detection via radial velocity is fundamentally the same as characterizing the orbits of spectroscopic binary systems. However, unlike in binary systems where the Doppler shift of light from each star can be measured, in exoplanet systems only the Doppler shift of the brighter host star is detectable (in most cases – we’ll cover the utility of planetary Doppler shifts for characterizing atmospheric circulation in a couple months). This means that the properties of the unseen planet are determined solely by studying the apparent motion of the host star.

To estimate the observable radial velocity shift of a star due to an unseen planet, we need to calculate the orbital velocity of this star around the center of mass of the star-planet system. The left-hand side of Figure 2.1 shows the geometry of this “binary” system, with the star and planet on opposite sides of their common center of mass.

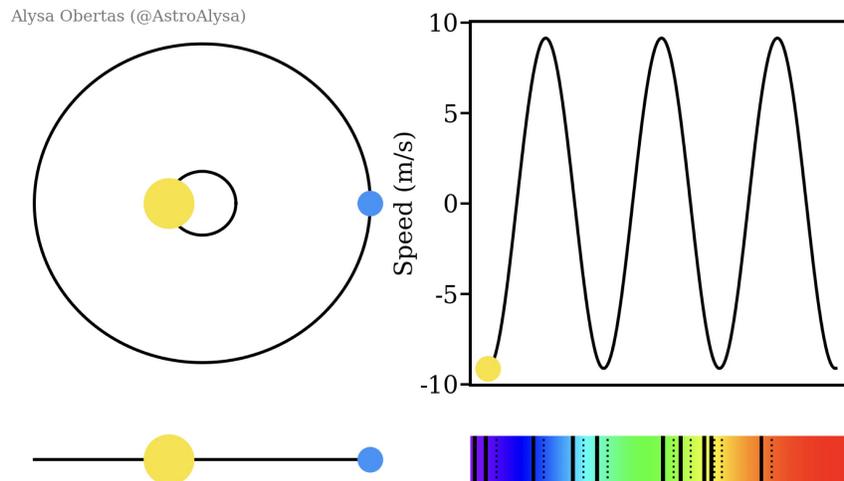


Figure 2.1: The geometry (left) and observables (right) of the radial velocity method. The top-left shows the face-on view of a star-planet system, while the bottom-left shows an edge-on view. The top right shows the resulting radial velocity curve (where, at the moment, the star is at the yellow dot) and the bottom right shows the spectral lines (solid) compared to their rest position (dashed). Figure from <https://astrobites.org/2019/10/16/the-nobel-winning-discovery-of-51-pegasi-b/>.

To start, let us call the mutual orbital period of the planet and star P , the separation between the planet and the center of mass of the system r_p , the separation between the star and the center of mass r_* , and the velocities of the planet and star v_p and v_* , respectively. Next, we can write down expressions for the velocities of the planet and star assuming

circular orbits:

$$\begin{aligned} v_p &= \frac{2\pi r_p}{P}, \\ v_\star &= \frac{2\pi r_\star}{P}. \end{aligned} \tag{2.4}$$

Note that real orbits aren't always aligned with our line of sight, so there is a projection effect that causes the maximum observed velocities to be smaller, depending on their inclination i :

$$\begin{aligned} v_{o,p} &= v_p \cos(90^\circ - i) = v_p \sin(i), \\ v_{o,\star} &= v_\star \sin(i). \end{aligned} \tag{2.5}$$

Let's now take the ratio of $v_{o,p}$ and $v_{o,\star}$:

$$\frac{v_{o,p}}{v_{o,\star}} = \frac{v_p \sin(i)}{v_\star \sin(i)} = \frac{2\pi r_p}{P} \frac{P}{2\pi r_\star} = \frac{r_p}{r_\star}. \tag{2.6}$$

Along with the definition of center of mass ($M_p r_p = M_\star r_\star$), this provides a useful set of relationships between velocity, mass, and separation for (circular) star-planet systems:

$$\frac{v_p}{v_\star} = \frac{r_p}{r_\star} = \frac{M_\star}{M_p}. \tag{2.7}$$

2.1.3 Group activity: deriving the radial velocity semi-amplitude for circular orbits

Recall Kepler's 3rd law,

$$\frac{a^3}{P^2} = \frac{GM_{\text{tot}}}{4\pi^2}, \tag{2.8}$$

where the total mass $M_{\text{tot}} = M_\star + M_p$, and a is the separation between the two objects, which can be related to their combined velocities as

$$a = r_p + r_\star = \frac{P}{2\pi} (v_p + v_\star). \tag{2.9}$$

From these and the discussion above, derive the radial velocity semi-amplitude, $K \equiv v_{o,\star}$, assuming circular orbits:

$$\boxed{K = \left(\frac{2\pi G}{P} \right)^{1/3} \frac{M_p \sin(i)}{(M_\star + M_p)^{2/3}}}. \tag{2.10}$$

Note that we'll derive the full radial velocity equation (i.e., Equation 2.27 from our textbook) as part of our homework that will be assigned next class. However, it's not too different from what we derived in class, just with an extra factor of $(1 - e^2)^{-1/2}$:

$$K = \left(\frac{2\pi G}{P} \right)^{1/3} \frac{M_p \sin(i)}{(M_\star + M_p)^{2/3}} \frac{1}{\sqrt{1 - e^2}}. \tag{2.11}$$

2.1.4 Detecting planets via radial velocity in practice

The radial velocity method was used to find the first exoplanet around a Sun-like star, 51 Pegasi b, by Mayor & Queloz in 1995 (Mayor & Queloz, 1995). Figure 2.2 shows their observed radial velocity curve from observations at Observatoire de Haute-Provence with the ELODIE spectrograph, covering $0.389 \mu\text{m} - 0.6815 \mu\text{m}$ with a radial velocity precision of $\approx 7 \text{ m s}^{-1}$. Using Equation (2.3), we can infer that this 7 m s^{-1} precision equates to a

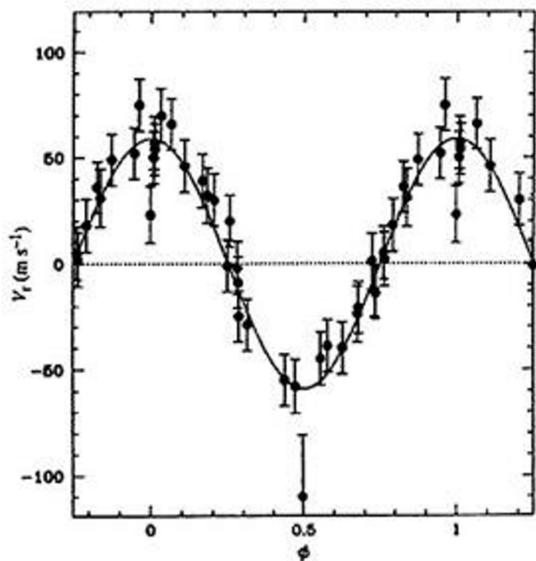


Figure 2.2: Radial velocity curve of 51 Pegasi b, the first exoplanet found via the radial velocity method. The x-axis is orbital phase (this is phase-folded data), and the y-axis is the radial velocity in m/s.

precision of $\Delta\lambda \sim 1.5 \times 10^{-5} \text{ nm}$ given a spectral line rest wavelength of 656.279 nm for the H- α line (note that in practice, radial velocity observations use many lines to calculate the Doppler shift). These and follow-up observations have measured the orbital period of 51 Pegasi b to be 4.23 days and the *minimum mass* of the planet to be $0.468 M_{\text{Jup}}$.

Note that the **radial velocity method alone does not allow a mass to be measured directly, it only places a lower limit on the mass**. This is because M_p and $\sin(i)$ are degenerate in Equation (2.11), only the combination is measured directly. In order to break this degeneracy, the inclination must be inferred through other means. The easiest way to break this is if the planet is also transiting (which allows i to be directly measured), which in turn infers one to measure the mass and radius directly and thus measure the density. We'll discuss this further next week.

Radial velocity measurements require very high spectral resolution $R \equiv \lambda/\Delta\lambda \sim 10^5$ – measuring these small wavelength shifts requires an instrument that measures the higher diffraction orders. The type of instrument regularly used for radial velocity instruments is an echelle spectrometer, which feeds light through both a low dispersion (standard) grating and then a specialized echelle grating that separates the high diffraction orders. The main challenge of radial velocity in its infancy was the wavelength calibration – advances in wavelength calibration using iodine gas cells led to the first radial velocity exoplanet detections. Today, most instruments either use emission lamps (e.g., Th-Ar, U-Ne) or laser frequency combs for wavelength calibration, in principle allowing measurements down to $K \approx 0.01 \text{ m s}^{-1}$.

However, astronomers cannot measure radial velocities down to the cm/s level at present

– the key challenge now is mitigating the effects of the star, which can overwhelm planetary signals or even masquerade as one. There are three main effects that stars have on radial velocity measurements: 1) stellar granulation, 2) stellar oscillations, 3) stellar activity (starspots, plages). These sum to cause a stellar impact in the radial velocity signal at the $\sim 1 \text{ m s}^{-1}$ level or higher.

1. Stellar granulation is the surface representation of small-scale convection cells in the envelope of a star. The radial velocity variability induced is on the order of $\approx 1 \text{ m s}^{-1}$. However, convection is fundamentally well-understood, and there are efforts to model stellar granulation with multi-dimensional stellar atmosphere models to remove the RV signals.
2. Stellar oscillations are due to pressure waves (P-modes) propagating inside stars and causing them to oscillate on a fairly regular timescale. These oscillations are on the order of $\approx 1 \text{ m s}^{-1}$, and they are studied in detail for specific stars with short-cadence measurements through the study of asteroseismology. They are generally dealt with by integrating over the oscillation period, which is on the order of tens of minutes. Importantly, the oscillation period depends on stellar mass because the timescale of oscillation scales as $\tau \propto \sqrt{\rho}$, where ρ is density, so lower-mass stars have longer oscillation timescales (and thus, require longer integrations).
3. Star spots (cool regions of a star) and plages (hotter than average regions) are perhaps the most pernicious stellar RV jitter. The amplitude of these can be up to 100 m s^{-1} for active stars. This often causes active stars to be simply left out of large-scale radial velocity surveys, with astronomers using activity indicators (e.g., Ca II H and K lines) to determine if stars are too active to study with RV. Additionally, stars like our Sun undergo long-term activity cycles that correlate with the migration of starspots from the mid-latitudes to the equator (a “Butterfly” pattern, see Figure 2.3). This means

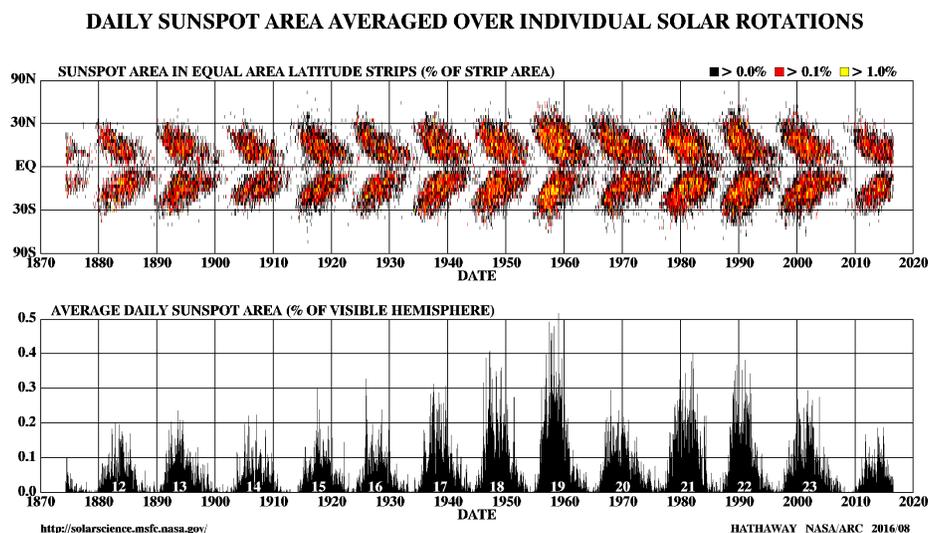


Figure 2.3: Maunder’s butterfly plot showing the Solar activity cycle.

that long-term radial velocity observations could see a periodic trend in the stellar “motion” that is due to stellar activity rather than an unseen planet.

Though it may seem challenging to overwhelm the impacts of stellar activity, astronomers can now regularly push instruments to RV precisions of $\lesssim 0.3 \text{ m s}^{-1}$, and the radial velocity method has detected 1,075 planets to date (https://exoplanetarchive.ipac.caltech.edu/docs/counts_detail.html). Figure 2.4 shows radial velocity curves of two interesting systems, HD 80606 (Naef et al., 2001) and 55 Cancri (Fischer et al., 2008). HD 80606b is

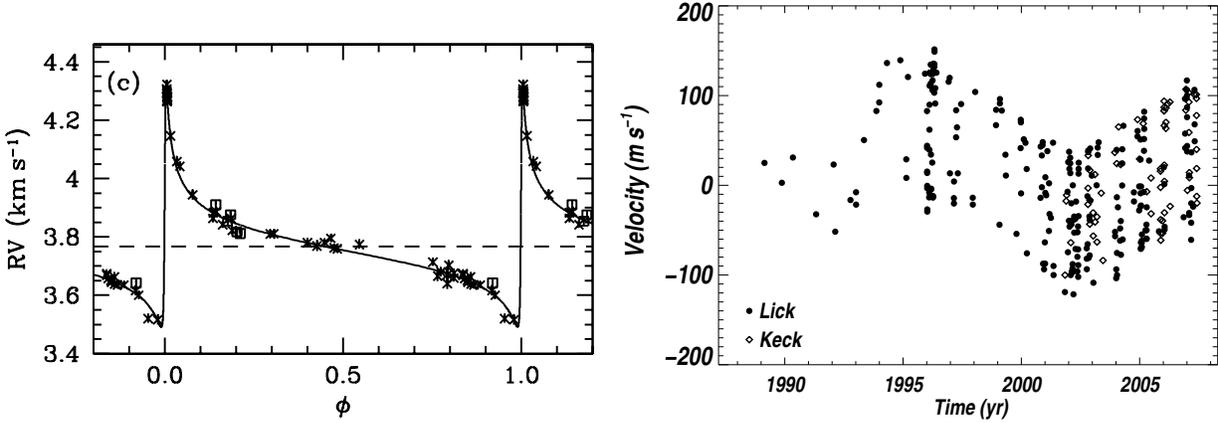


Figure 2.4: Radial velocity curves of HD 80606 (left) and 55 Cancri (right).

one of the most eccentric exoplanets known, with $e = 0.927$. Note how the high eccentricity causes a rapid shift in the stellar RV near periape. 55 Cancri is a system with five known planets, and the RV curve displays this complexity – you can easily see a 14-year period from the outer planet, but there are four other planets with periods of 2.8, 14.7, 44.3, 260.8 days as well. To fit these multi-planet RV curves one usually ignores the effect of planet-planet gravitational interactions, allowing one to model the radial velocity curve of the star as the linear superposition of the radial velocity curve of each individual planet. By subtracting the radial velocity curves of known planets, one can then search for other periodic signals in the RV curve using a periodogram (i.e., Lomb-Scargle) analysis and/or a more statistically robust method (e.g., MCMC).

2.2 Radial velocity in practice: group activity

Using our simplified formula for radial velocity assuming a circular orbit (Equation 2.10), let’s estimate the radial velocity semi-amplitude that planets in our Solar System would cause around another Sun-like star.

1. Calculate the radial velocity semi-amplitude caused by Jupiter orbiting our Sun, assuming that it is viewed edge-on (i.e., $i = 90^\circ$). You may use $M_{\text{Jup}} = 1.898 \times 10^{27} \text{ kg}$, $P_{\text{Jup}} = 4331 \text{ days}$, and $M_{\odot} = 1.989 \times 10^{30} \text{ kg}$.
2. Calculate the radial velocity semi-amplitude again for a Jupiter-mass planet around a Sun-like star, but now with the orbital period of 51 Pegasi b ($P = 4.23 \text{ days}$). Discuss what effects together make it easier to detect planets via radial velocity that are both more massive and closer in to the host star (hint: it’s more than just a larger K).

3. Calculate the radial velocity semi-amplitude caused by Earth orbiting our Sun, again assuming it is viewed edge-on. You may use $M_{\oplus} = 5.97 \times 10^{24}$ kg, and $P_{\oplus} = 365.2$ days. How does this compare to the radial velocity accuracy that modern instruments have of $\approx 0.3 \text{ m s}^{-1}$ (textbook, Ch. 2.4)?

3 Detecting exoplanets: astrometry

Our agenda for Day 3 is the following:

1. Finish up/recap radial velocity (10 min)
2. Cover the fundamental method and history of astrometry (15 min)
3. Derive the astrometric wobble of a star with a companion planet (15 min)
4. Group activity: calculate the astrometric wobble and compare it to historical “detections” (25 min)
5. Cover modern space-based astrometry (10 min)

Today’s reading is from our textbook, Ch. 3.1-3.9. This will cover what sets the apparent size of the astrometric “wobble,” our current ground- and space-based observational capabilities, and inferring planetary system properties with astrometry.

3.1 Astrometry: notes

3.1.1 Method, historical and modern observations

Astrometry is perhaps the most intuitive exoplanet detection method: observing the gravitational influence of an unseen companion planet by studying the changing position of its host star on the sky. As a result of its requirement only to measure precise positions of stars, it has a long historical record dating back to William Herschel claiming a stellar or planetary companion to 70 Ophiuchi in 1779.

The observable for the astrometric detection of exoplanets is simply the angular shift of the location of a star in the sky as it orbits the center of mass of a star-planet system. The top-left hand panel of Figure 2.1 (and the gif from the class slides) demonstrate the orbit of a star around the common center of mass of a face-on star-planet system. However, the astrometric shifts due to the planet (generally $\lesssim 1$ milli-arcsecond) are much smaller than other motions of the star on the sky due to proper motion and parallax – as a result, these effects need to be accounted for to determine the astrometric motion due to the planet.

The proper motion is the apparent shift in angular location of a star as it moves across the sky. Proper motion is caused by the tangential motion of a star on-sky, which is in the direction orthogonal to the radial motion (which as we discussed last class, can be constrained by RV measurements). Figure 3.1 shows the proper motion of Barnard’s star over the course of eight years, which is the highest proper motion of any star (and which you will calculate in our group activity). Proper motion is defined as

$$\mu \equiv \frac{d\theta}{dt} = \frac{v_\theta}{d} \quad (3.1)$$

where θ is the angular position, v_θ is the transverse (tangential) velocity, and d is the distance to the star. As you can see, proper motion is the equivalent to the angular speed of a star across the sky in the transverse direction parallel to the plane of the sky.

Parallax is the apparent shift in position of a star on sky (relative to distant background stars) due to the motion of Earth around the Sun. Figure 3.2 shows a schematic of how



Figure 3.1: Barnard’s star has a high proper motion which is clearly detectable from year-to-year.

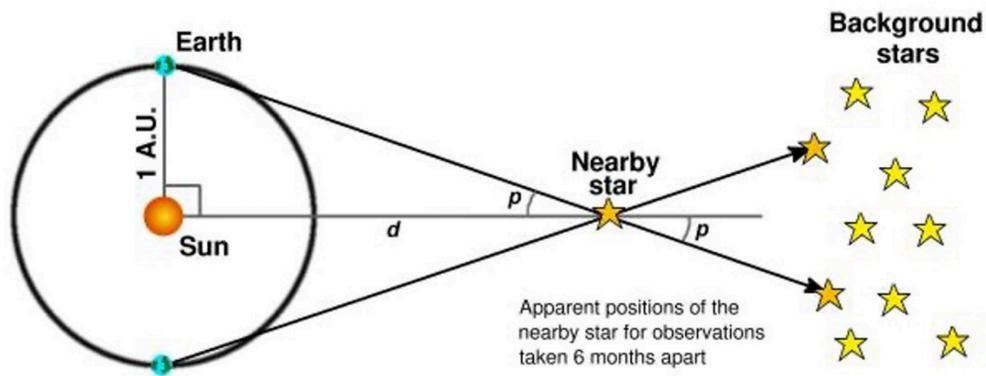


Figure 3.2: Schematic showing how trigonometric parallax relates to distance and the Earth-Sun separation.

the parallax of a star p is related to the Earth-Sun distance (1 au) and the distance to the star d . Using the small-angle approximation, $1 \text{ au} = pd$, and thus the distance to the star is $d = 1 \text{ au}/p$. Astronomers traditionally measure parallax in arcseconds, causing the distance to a star with a parallax of 1 arcsecond to be $d = 206265 \text{ au}$. This distance scale has in turn been defined as the “parsec” or pc (distance with a PARallax of one arcSECond), leading to the traditional formula for parallactic distance that we used on Day 1:

$$d(\text{pc}) = \frac{1}{p(\text{arcsec})}. \quad (3.2)$$

For reference, the parallax of Barnard’s star is 0.545 arcseconds, which implies that its distance is 1.84 pc. As is discussed in Section 3.2, astrometric measurements from the ground are fraught due to atmospheric turbulence limiting seeing, with typical astrometric precision of most ground-based surveys ≈ 100 milli-arcseconds for targets with $m_v = 10$.

However, for large (10 meter-class) telescopes like the VLT, astrometric precision can reach the $\approx 300 \mu\text{as}$ level for bright targets in good seeing conditions. The majority of astrometric measurements of parallax come from space-based measurements with the Hipparcos (1989-1993) and Gaia (2013-) observatories. Hipparcos had a characteristic astrometric accuracy of 1 mas for $m_v = 10$, while Gaia has a significantly improved characteristic accuracy of $10 \mu\text{as}$. Figure 3.3 shows a typical series of astrometric measurements over the course of three years from Hipparcos, showing the proper motion (movement from bottom left to upper right) along with parallax (loops along this movement, one per year).

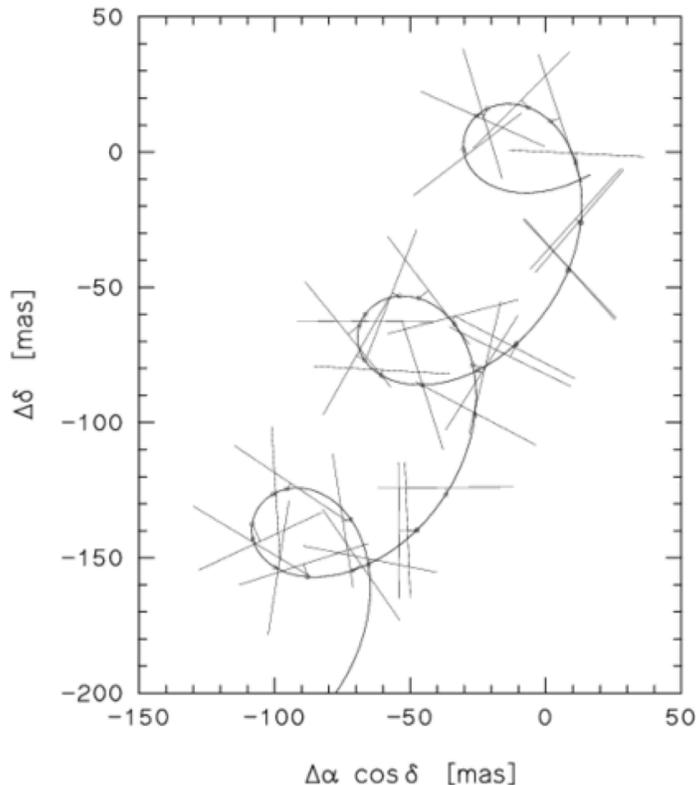


Figure 3.3: Example of the measurement of the path in the sky of a Hipparcos star.

Now that we've covered what needs to be removed to isolate the astrometric signal due to an unseen planet, let's preface our discussion of actual astrometric planetary detections with a parable. Prof. Peter van de Kamp at Swarthmore College observed the astrometric positions of Barnard's star on their 24" telescope starting in 1938. Figure 3.4 shows his observed changes in right ascension of Barnard's star over observations spanning a 31-year baseline from 1938 to 1969. Van de Kamp inferred from the R.A. changes of Barnard's star that there were two planets with masses of 1.1 and 0.8 Jupiter masses orbiting the star with periods of 26 and 12 years (van de Kamp, 1969). However, his own successor at Swarthmore, Wulff Heintz, showed that these changes were not due to stellar motion but due to aberrations on the photographic plates used to record the changing position of the star. This displays the challenges of ground-based astrometry, and as we'll discuss next the first promising astrometric detection came from space-based observations.

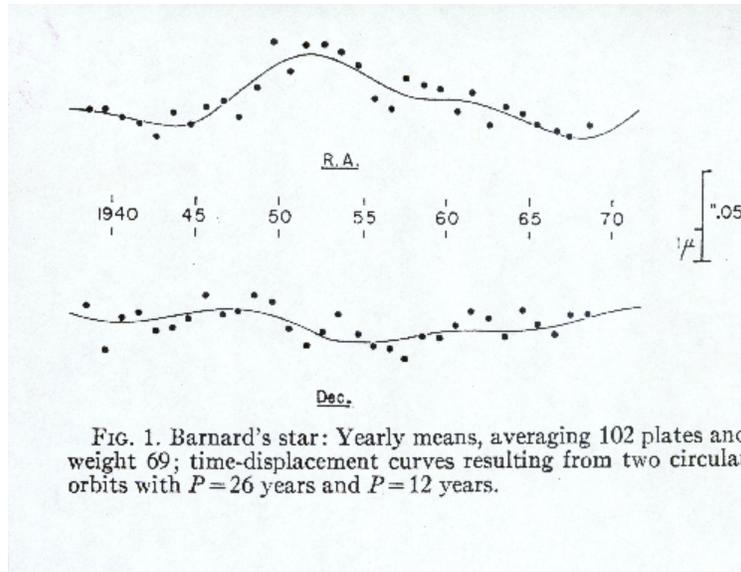


Figure 3.4: Van de Kamp’s observed “astrometric” motion of Barnard’s star, along with fits for two planets with orbital periods of 26 and 16 years.

The astronomy community had to wait until 2002 for an undisputed *detection* of planets via astrometry using Hubble Space Telescope observations of the GJ 876 system (Benedict et al., 2002). In this case (and almost all other cases of planets detected via astrometry), the planets in the GJ 876 system had already been *discovered* via the radial velocity method. There are only three cases of planets discovered by astrometry (according to the NASA Exoplanet Archive, https://exoplanetarchive.ipac.caltech.edu/docs/counts_detail.html), but the upcoming Gaia DR4 will likely change that – Gaia astrometry is expected to discover tens of thousands of planets (Perryman et al., 2014)

Combining planet detections with RV (or direct imaging, as we’ll discuss in two weeks) with follow-up astrometry has provided one method of breaking the $M\sin(i)$ degeneracy. Figures 3.5 and 3.6 show the example of the ν And system, which has three planets that were previously discovered by RV, the outer two of which cause the detectable astrometric motion of ν And A (McArthur et al., 2010). In this case, combining astrometry and RV together allows direct measurements of inclination because astrometry directly constrains a, e, ω , and t_p (where ω is the argument of pericenter and t_p is the time of pericenter passage) while RV directly constrains e, P, t_p, ω along with the combination of a, e, P, i through the RV semi-amplitude K . As a result of astrometry providing a and RV providing P , we can now also solve Kepler’s 3rd law ($a^3/P^2 \propto M_\star + M_p$) to provide another independent determination of M_p and thus $\sin(i)$ combining the two methods. In practice, astronomers do a joint fit to both the astrometry and the RV data (see Figure 3.6) in order to directly constrain M through a joint minimization process (see textbook Eq. 3.25) rather than iterating between the two solutions.

3.1.2 Astrometric wobble due to a companion planet

In the most deceptively simple derivation we’ll do this semester, the maximum astrometric spatial shift of a star due to an unseen planet can be determined from the definition of center

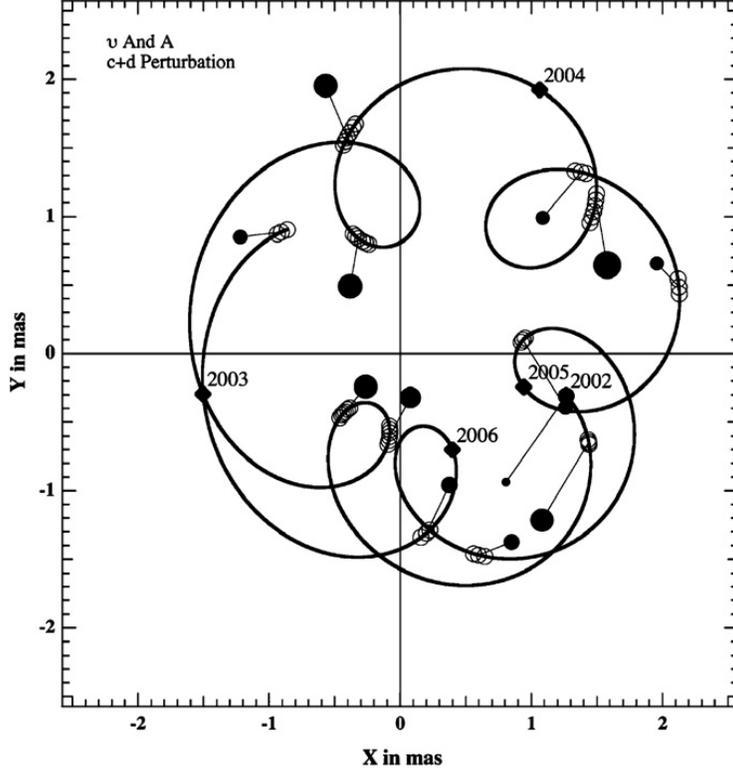


Figure 3.5: On-sky astrometric reflex motion of ν And A due to its companion planets c and d. Measured “normal points” (averages of individual data points) are shown in black filled circles, while the time of observations are shown in the open circles. The curved line shows the best-fit model, and the smaller straight lines show residuals.

of mass of the system that we used previously to derive Equation (2.7):

$$M_{\star}r_{\star} = M_p r_p . \quad (3.3)$$

Defining the orbital semi-major axis $a = r_p + r_{\star}$, we can substitute for r_p

$$r_{\star} = \frac{M_p}{M_{\star}}(a - r_{\star}) \quad (3.4)$$

and re-arrange to find

$$r_{\star} = a \frac{M_p}{M_{\star} + M_p} \approx a \frac{M_p}{M_{\star}} , \quad (3.5)$$

which is equivalent to textbook Eq. 3.1. Importantly, Equation (3.5) is the motion in projected distance, while the observable is the angular shift. Similar to parallax, we can again use the small-angle equation $r_{\star} = d\alpha$ to solve for the angular astrometric shift α

$$\alpha = \frac{a}{d} \frac{M_p}{M_{\star} + M_p} \approx \frac{a}{d} \frac{M_p}{M_{\star}} , \quad (3.6)$$

which can be written in scaled power-law form in units of arcsec (textbook Eq. 3.2):

$$\alpha \approx \left(\frac{M_p}{M_{\star}} \right) \left(\frac{a}{1 \text{ au}} \right) \left(\frac{d}{1 \text{ pc}} \right)^{-1} \text{ arcsec} . \quad (3.7)$$

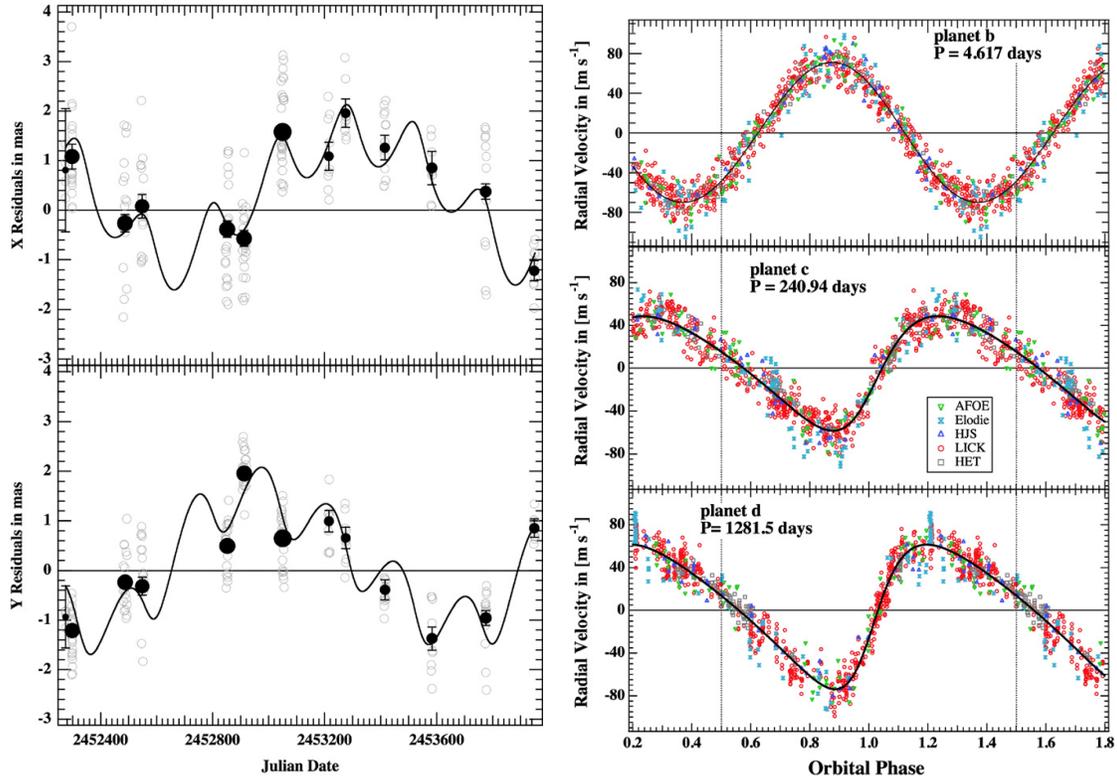


Figure 3.6: Astrometric fit (left) and RV fit (right) to ν And A. The top left panel shows the residual motion due to planet c, and the bottom left planet d, with solid black points showing averaged data and open circles showing individual observations. Planet b is fit for in the RVs, but does not impact the astrometry significantly due to its close-in orbit.

3.2 Astrometry: group activity

Peter Van de Kamp is (in)famous for his early claims of two approximately Jupiter-mass planets orbiting the nearby Barnard’s star ($d = 1.84$ pc, $M = 0.14 M_{\odot}$). This exercise will help you infer whether he even had the ability to detect planets around Barnard’s star via astrometry.

1. Calculate the amplitude of the Sun’s angular astrometric wobble (in units of milli-arcseconds) due to Jupiter if it were viewed from a distance of 10 pc. Note $a_{\text{Jup}} = 5.23$ au.
2. Determine the amplitude of the astrometric wobble of Barnard’s star due to the hypothetical Jovian-mass planet on a 26-year orbital period.
3. The tangential velocity (perpendicular to the radial velocity) of Barnard’s star is $\approx 90 \text{ km s}^{-1}$. Calculate the proper motion of Barnard’s star in units of arcsec/year, and compare this to your answer in part (b).

4 Detecting exoplanets: transits

Our agenda for Day 4 is the following:

1. One-slide intro to the transit method (5 min)
2. Derive the transit depth and transit probability (15 min)
3. Group activity: calculate the transit depth and probability for HD 209458b and Earth (20 min)
4. Mathematics of transits: impact parameter (15 min)
5. Group activity: draw some transits! (20 min, till end of class)
6. Mathematics of transits: stellar density (10 min, if time)
7. Transit method in practice: Kepler, TESS, ground-based surveys (15 min, if time)

Note that we'll be finishing transits during the next lecture, given that timing is a relatively short topic. Today's reading is from the textbook, Ch. 6.1-6.6 and 6.13, and/or the Winn handout on ELMS. This will cover the fundamentals of the transit method, previous transit searches from the ground and space and notable discoveries, as well as modeling transit light curves.

4.1 Transits: notes

4.1.1 Transit depth, probability, and duration

The transit method detects planets through the small dip in observed starlight that occurs when the planet passes between the star and the observer's point of view. The relative fraction of sky that a given object subtends can be quantified by the solid angle

$$\Omega = \frac{A}{4\pi d^2}, \quad (4.1)$$

where A is the projected area of the object and d is the distance to the object. For a (spherical) star or planet, $A = \pi R^2$. The transit depth is then simply the relative fraction of the star's area that the planet covers (blocks)

$$\frac{\Omega_p}{\Omega_\star} \equiv \delta = \frac{\pi R_p^2}{4\pi d^2} \frac{4\pi(d+a)^2}{\pi R_\star^2} \approx \left(\frac{R_p}{R_\star}\right)^2, \quad (4.2)$$

using the valid assumption that the distance between the star and planet $a \ll d$. Then, the total flux from the system during a perfectly edge-on transit event can be related to the unocculted stellar flux F_\star as

$$F = F_\star \left[1 - \left(\frac{R_p}{R_\star}\right)^2 \right], \quad (4.3)$$

where note that F_\star can itself vary over the course of a transit. This is because regions near the edges of the stellar disk appear dimmer by a factor that scales with $\mu = \cos\theta$, where

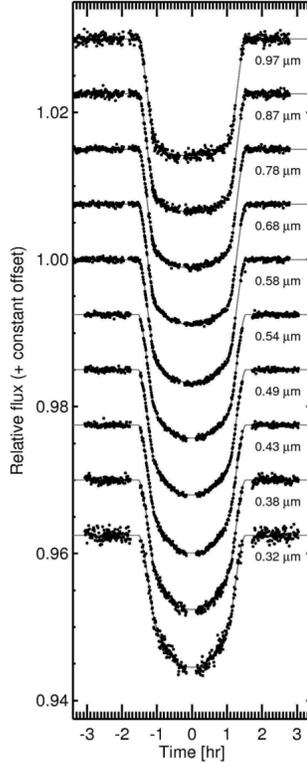


Figure 4.1: Example of limb darkening in an exoplanet transit. Note how the shape of the transit is chromatic, as limb darkening has a larger effect in visible wavelengths.

θ is the angle between the direction of the observer and the location of the stellar surface. This physically occurs because the effective photosphere of the star becomes shallower (at higher altitudes and lower pressures) near the limb due to the enhanced optical path of a light ray to escape the limb relative to the center of the disk. Usually, stellar atmospheres near the photosphere decrease in temperature with increasing height, so emission from these higher regions near the limb is cooler and redder, causing the limb to appear darker – that is why this is termed “limb darkening.” As a result, the decrement in flux is smaller near the edge of the stellar disk and increases towards the center, resulting in transits with a U-shape rather than a flat bottom (see Figure 4.1). This effect is chromatic, such that limb darkening affects bluer wavelengths more, because these wavelengths have a larger change in photosphere pressure from the center to limb of the stellar disk.

For a circular orbit, the transit probability can be calculated by considering the range of angles at which an observer at infinity would see the planet occult the stellar disk. To derive this, we can define an angle $\theta_0 = 90^\circ - i$, which is the angle of the planet’s orbit with respect to the observer. We can then integrate to find the probability that the planet lies within the maximum possible angle from our line of sight such that it transits its host star (which is simple for a circular orbit, but more complex for an elliptical orbit, see textbook section 6.13.6):

$$P \approx \frac{\int_0^{\theta_0} \cos\theta d\theta}{\int_0^{\pi/2} \cos\theta d\theta} = \frac{\sin\theta|_0^{\theta_0}}{\sin\theta|_0^{\pi/2}} = \frac{\sin\theta_0}{1} = \frac{R_\star}{a}. \quad (4.4)$$

As a result, the transit probability does *not* depend on planet radius, and is larger for planets with smaller orbital semi-major axes. The full expression for transit probability on

an eccentric orbit, accounting for the possibility of grazing transits, is

$$P = \frac{R_\star \pm R_p}{a(1 - e^2)}. \quad (4.5)$$

Lastly, the duration of a transit can be simply estimated as

$$\tau \approx \frac{2R_\star}{v_p}, \quad (4.6)$$

where v_p is the orbital velocity of the planet. For a circular orbit, we find

$$\tau \approx 2R_\star \sqrt{\frac{a}{GM_\star}}, \quad (4.7)$$

which can be re-written in terms of common quantities as (textbook, Eq. 6.11):

$$\tau = 13 \text{ hr} \left(\frac{M_\star}{M_\odot} \right)^{-1/2} \left(\frac{a}{1 \text{ au}} \right)^{1/2} \left(\frac{R_\star}{R_\odot} \right). \quad (4.8)$$

4.1.2 Transit geometry, impact parameter

Figure 4.2 shows the detailed transit geometry for a single planet transiting a single star. Assuming the orbit is circular and the planet mass is much less than the stellar mass, there are five equations that can altogether specify the system. The first is simply the transit depth

$$\delta \equiv \left(\frac{R_p}{R_\star} \right)^2, \quad (4.9)$$

which we previously derived.

The second is the transit duration, or t_t , which is the time between the first and fourth contacts. Figure 4.3 shows schematics of the transit duration with respect to the full orbit as well as the disk of the star. Given the triangle made between the line through which the planet crosses the disk and the center of the stellar disk, we can relate half of the length of the chord that the planet traverses to the planet and star radii as $\sqrt{(R_\star + R_p)^2 - a^2 \cos^2(i)}$. Using the left hand side of Figure 4.3, we can then note that the angle of the full orbit that the planet sweeps out during transit is $\sin^{-1}(\sqrt{(R_\star + R_p)^2 - a^2 \cos^2(i)}/a)$. We can thus relate this angle as a fraction of the full orbit (2π) to determine the transit duration

$$t_t = \frac{P}{2\pi} 2 \sin^{-1} \left(\frac{\sqrt{(R_\star + R_p)^2 - a^2 \cos^2(i)}}{a} \right), \quad (4.10)$$

Which is often re-written in the form of Seager & Mallén-Ornelas (2003):

$$t_t = \frac{P}{\pi} \sin^{-1} \left(\frac{R_\star}{a} \left[\frac{(1 + (R_p/R_\star))^2 - [(a/R_\star)\cos i]^2}{1 - \cos^2 i} \right]^{1/2} \right). \quad (4.11)$$

The third is the transit shape, which is the ratio of the duration of the flat bottom of the transit t_f to the full transit t_t . In the case of the flat bottom, half of the chord that the

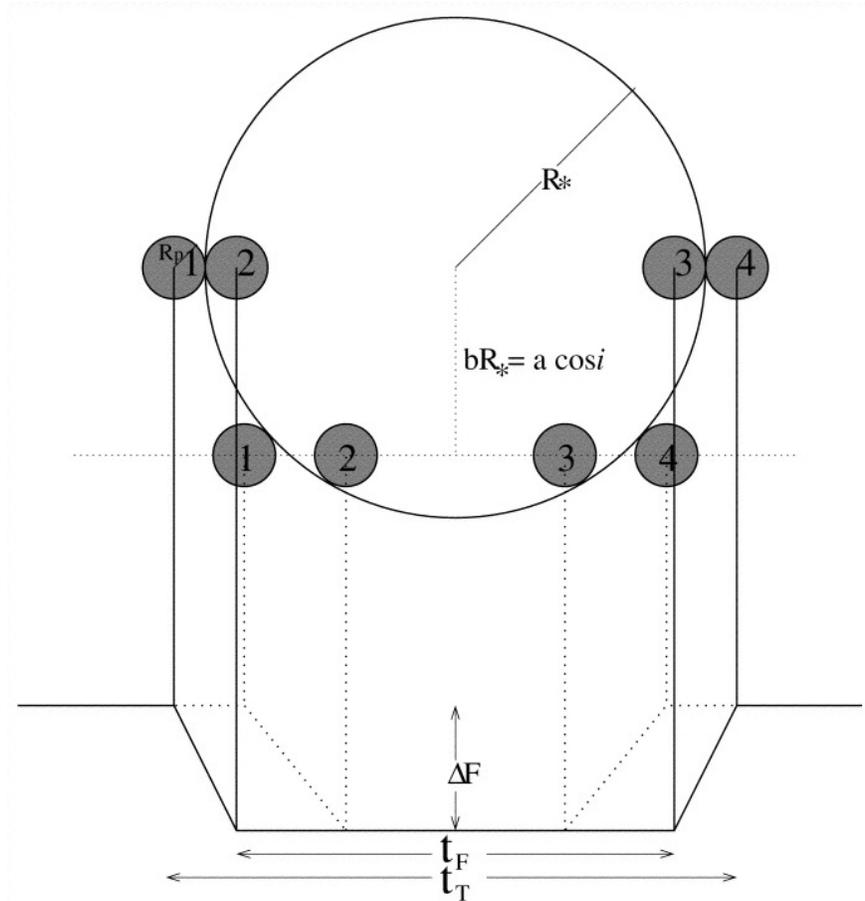


Figure 4.2: Schematic of the transit geometry, showing the 1st, 2nd, 3rd, and 4th contacts, total transit duration t_t , and full occultation duration t_f , along with the definition of impact parameter b . Figure adapted from Seager & Mallén-Ornelas (2003).

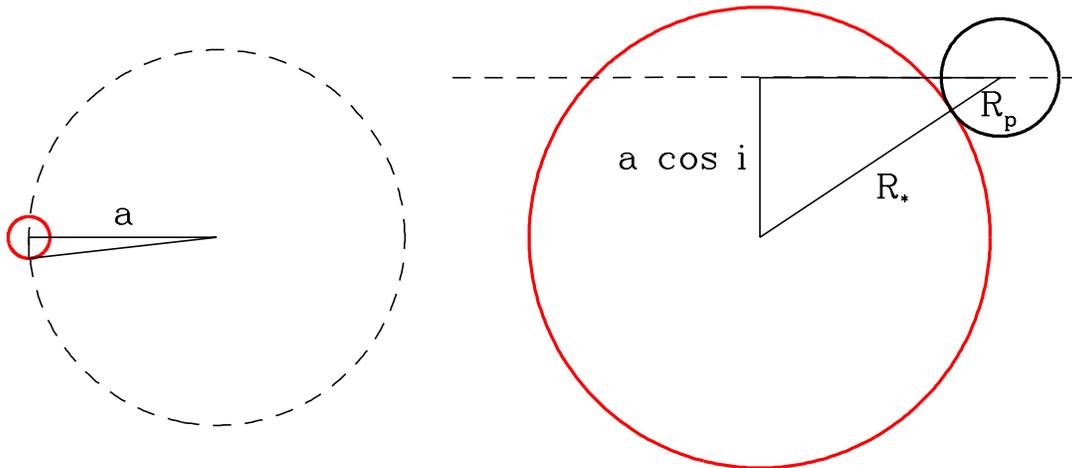


Figure 4.3: Schematic of the transit duration as a fraction of the total orbit (left) for which some portion of the planet occults the star (right). Adapted from Sackett (1999).

planet traverses is $\sqrt{(R_\star - R_p)^2 - a^2 \cos^2(i)}$ (the difference with t_t being the minus sign in $R_\star - R_p$). Thus, the ratio of the durations can be expressed as

$$\frac{\sin(t_f \pi / P)}{\sin(t_t \pi / P)} = \frac{\sqrt{[1 - (R_p / R_\star)]^2 - [(a / R_\star) \cos(i)]^2}}{\sqrt{[1 + (R_p / R_\star)]^2 - [(a / R_\star) \cos(i)]^2}}. \quad (4.12)$$

The combination $(a / R_\star) \cos(i) \equiv b$, which is known as the impact parameter. The impact parameter is the projected distance between the planet and star centers during mid-transit in units of R_\star , and varies from $b = 0$ for a transit that crosses the mid-plane of the stellar disk to $b = \pm 1$ for a grazing transit where at most half the planet occults the stellar disk (note that b can also be larger than 1 for a barely grazing transit).

The fourth equation that specifies a transit is well-known to us – Kepler’s third law:

$$P^2 = \frac{4\pi^2 a^3}{G(M_\star + M_p)}. \quad (4.13)$$

The fifth and final equation is an assumed power-law mass-radius relationship for the host stars

$$R_\star = k M_\star^x, \quad (4.14)$$

where k is constant for each luminosity class (main-sequence, giant stars, etc.) and x describes the power-law relationship for that sequence. For Sun-like stars, $x \approx 0.8$.

4.1.3 Measuring stellar density via transits

The stellar density ρ_\star is the only parameter directly constrained from a transit observation – note that the planetary radius R_p is dependent on the (a priori unknown) R_\star , and thus R_p is usually dependent on our model uncertainty for R_\star . As a result, the direct measurement of ρ_\star is critical to better predicting the stellar radius by benchmarking stellar models.

To derive the stellar density, we start by re-arranging Equation (4.11) to solve for the ratio a / R_\star :

$$\frac{a}{R_\star} = \sqrt{\frac{(1 + \sqrt{\delta})^2 - b^2[1 - \sin^2(t_t \pi / P)]}{\sin^2(t_t \pi / P)}}. \quad (4.15)$$

Denoting the right-hand-side of the preceding equation as $f(\delta, b, t_t, P)$, we can note that $R_\star = a / f(\delta, b, t_t, P) = f(\delta, b, t_t, P)^{-1} [GM_\star P^2 / (4\pi^2)]^{1/3}$. Noting that $\rho_\star = M_\star / R_\star^3$, we can find $\rho_\star = 4\pi^2 / (GP^2) f(\delta, b, t_t, P)^3$ – which does not depend on M_\star or R_\star , and instead only on measurable quantities from a transit observation. The full expression for ρ_\star is

$$\rho_\star = \left(\frac{4\pi^2}{GP^2} \right) \left(\frac{(1 + \sqrt{\delta})^2 - b^2[1 - \sin^2(t_t \pi / P)]}{\sin^2(t_t \pi / P)} \right)^{3/2} \quad (4.16)$$

4.1.4 Transit method in practice

Up until the launch of CoRoT, all detections of exoplanets were from ground-based surveys. After the initial set of surveys done in parking lots with small telescopes (i.e., the Charbonneau transit detection of HD 209458b) that followed up RV detections, a set of more complete wide-field surveys were developed. These surveys were designed to have a

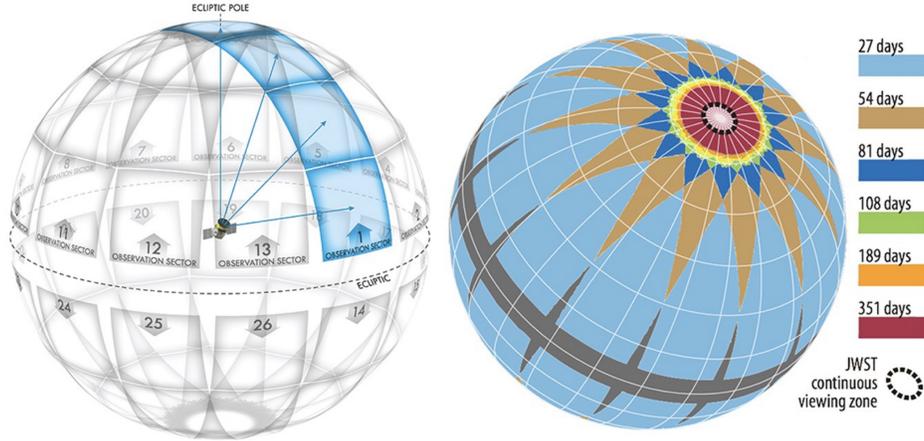


Figure 4.4: The observation “sectors” of TESS (left) and the resulting viewing zones of the mission (right).

large etendue $E = A\Omega$, where A is effective aperture area and Ω is the solid angle on the sky imaged in a single exposure, and short cadences of a few minutes. Long-lasting early surveys include (Super)WASP (2005-), HAT(NET) (2003-), and KELT (2005-) – these early surveys have detected hundreds of planets, including some of the best studied objects (e.g., WASP-39b, the target of the recent JWST ERS transmission spectroscopy program). The recent focus of ground-based surveys has been to detect rocky planets orbiting small M-dwarf stars, and these surveys of nearby M dwarfs include MEarth (2008-), TRAPPIST (2010-), and its successor SPECULOOS (2017-).

The first space-borne transit observatory was CoRoT, which launched in 2006 and detected ≈ 32 planets. Kepler was the first truly transformative exoplanet mission, which was launched in 2009 and detected 2,778 confirmed planets (which is a sea change given that less than 500 planets were known at the time of launch). Kepler was so successful due to its pointing stability, broad field of view, and focus on only a single field of 150,000 stars near the constellation of Cygnus with 30 minute cadence. The current space-based exoplanet detection workhorse is the Transiting Exoplanet Survey Satellite (TESS), which launched in 2018. TESS is the first all-sky exoplanet survey, and it observes space in sectors, which it observes for a little less than a month (27 days) at a time. This results in the ecliptic poles being continuously observable by TESS when it is observing that hemisphere (the sector paths flip from N to S hemisphere, 13 in each), which is critical because the JWST continuous viewing zone is also at the ecliptic poles. As a result of its focus on finding nearby transiting planets that can be studied with follow-up, TESS is often considered a “finder scope” for JWST.

One drawback of transit observations is that they have a somewhat high false positive rate of $\approx 10\%$, given that other astrophysical phenomena (especially eclipsing binaries) can cause transit-like signals. As a result, it is critical to combine transit observations with radial velocity measurements to confirm that the transit signal is indeed due to a planet. Importantly, putting together both transit and RV measurements measure the planet density, as transit measures R_p and i while RV measures the combination $M_p \sin(i)$.

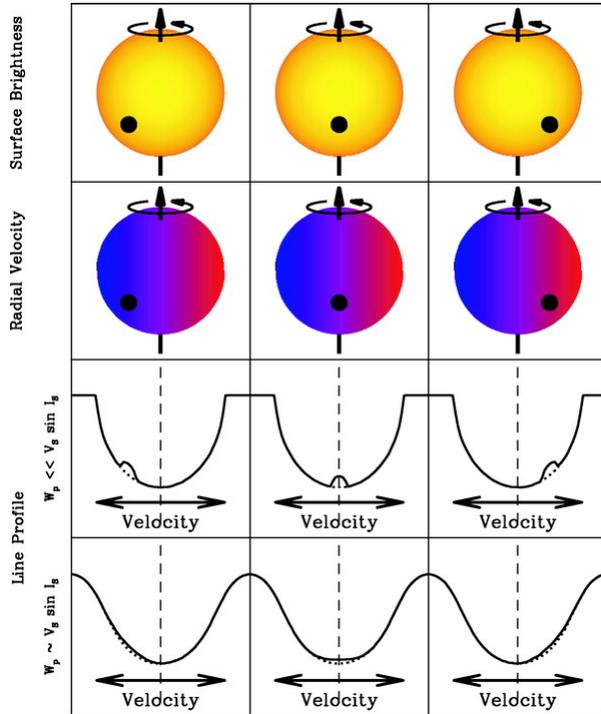


Figure 4.5: The Rossiter-McLaughlin effect is caused by the planet blocking the emission from different regions of the star as it crosses the transit chord, changing the stellar spectra due to the differential blocking of blueshifted vs. red-shifted regions caused by the stellar rotation.

The combination of transit and RV is also a powerful way to determine the angle between the stellar spin and planetary orbit. Figure 4.5 demonstrates how the “Rossiter-McLaughlin effect” caused by a planet blocking emission from different regions of the stellar disk during transit can cause changes in the line profile of the star due to the planet blocking blueshifted or redshifted regions caused by the rotation of the star Doppler shifting/broadening the spectrum. If the planet is aligned with the stellar spin, the variation will be such that the velocity of the star appears to shift back and forth. However, if the planet is in a near-polar orbit, the effect on the Doppler shifting/broadening of the lines will be smaller.

4.2 Transits: group activities

4.2.1 Calculating transit depth and probability

Recall that the transit depth $\delta = (R_p/R_\star)^2$ and the transit probability $p = R_\star/a$. Also, note that the transit duration is (textbook, Eq. 6.11)

$$\tau = 13 \text{ hr} \left(\frac{M_\star}{M_\odot} \right)^{-1/2} \left(\frac{a}{1 \text{ au}} \right)^{1/2} \left(\frac{R_\star}{R_\odot} \right). \quad (4.17)$$

HD 209458b was the first planet discovered via transits. It orbits a Sun-like star with a semi-major axis of 0.048 au, and has a radius of $1.38 R_{\text{Jup}}$.

1. Calculate the transit depth, probability, and duration of HD 209458b. Assume that HD 209458 has the same properties as the Sun. (Groups 1-2: calculate the transit depth. Groups 3-4: calculate the probability. Groups 5-6: calculate the duration.)
2. Calculate the transit depth, probability, and duration of Earth around the Sun.

(Groups 1-2: calculate the transit depth. Groups 3-4: calculate the probability. Groups 5-6: calculate the duration.)

3. Photometric measurements capable of measuring the transit dip of HD 209458b were available more than a decade prior to its detection in 1999. If this was the case, why did it take so long to find the first transiting planet?
4. As you found, detecting Earth around the Sun with the transit is challenging. How would you go about designing a survey to find a copy of Earth with the transit method?

4.2.2 Drawing transits

Let's gain some conceptual understanding by drawing idealized transit events, for 5 different scenarios:

1. Planet with $i = 90^\circ$ transits across the center of the star (i.e., a baseline transit event – use this as the reference for all your other drawings).
2. Smaller radius planet transits the same host star.
3. Planet transiting with $i < 90^\circ$, but not in a grazing configuration.
4. Planet with an impact parameter of $b = 1$.
5. Longer period planet with $i = 90^\circ$.

5 Detecting exoplanets: timing

Our agenda for Day 5 is the following:

1. Recap transit geometry (5 minutes)
2. Measuring stellar density via transits (5 minutes)
3. Group activity: draw some transits! (20 minutes)
4. Transit detections in practice (15 minutes)
5. Transit timing variations (10 minutes)
6. Principles of detecting planets via timing (10 min)
7. Group activity: detecting pulsar planets (if time, if not either start next class with this or skip)

We'll finish covering transits in full (see Day 4 notes) before moving on to timing. Timing on its own has only been used to find 7 confirmed planets around pulsars and 2 orbiting pulsating variable stars, while transits have found 4153 planets, so fractional to its detection count we're giving timing plenty of attention! Today's reading is from the textbook, Chs. 6.20, 4.1-4.2 or the Agol & Fabrycky handout (I highly recommend the latter). This will cover transit timing variations and applying timing to find planets orbiting pulsars.

5.1 Timing: notes

5.1.1 Transit timing variations

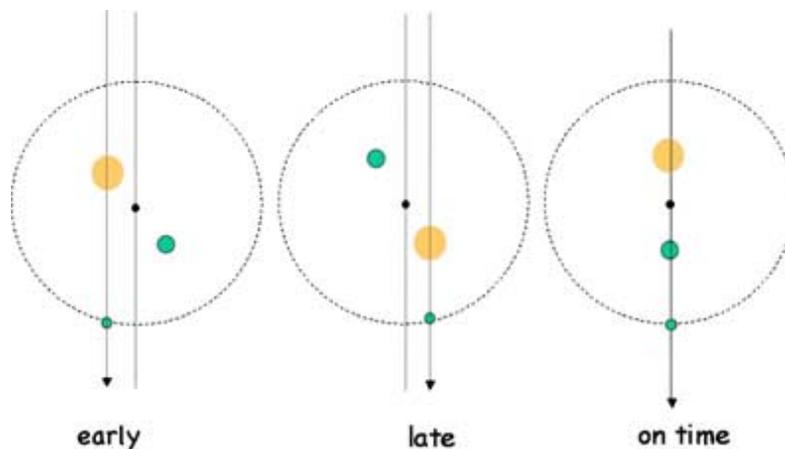


Figure 5.1: Diagram showing how TTVs on an outer planet can be caused by the gravitational influence of an inner planet on the host star. Adapted from Agol et al. (2005).

Transit timing variations (TTVs) are deviations from the regular transit times expected for a single planet on a Keplerian orbit around its host star. Transit timing variations are

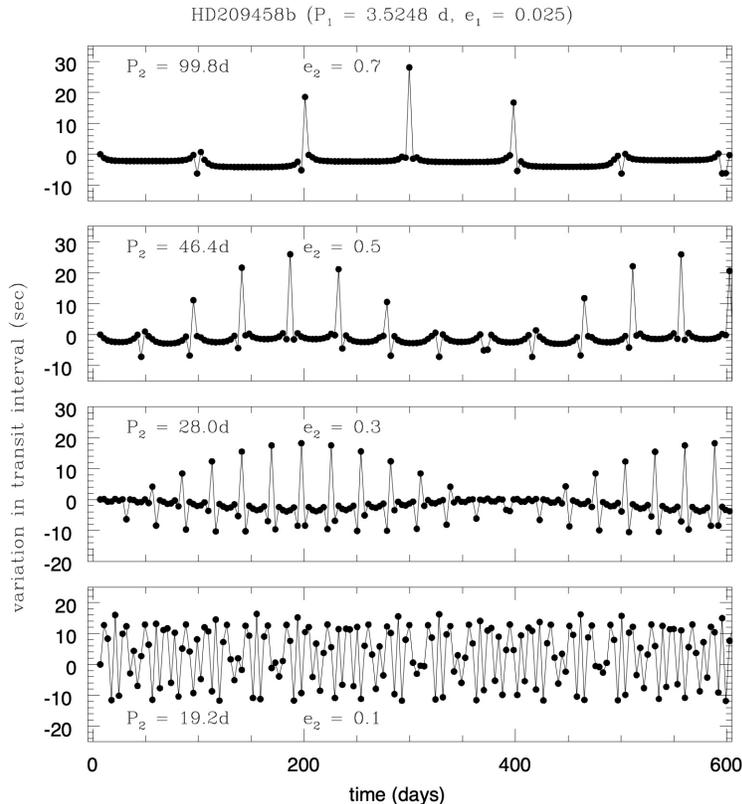


Figure 5.2: Diagram showing how TTVs of hypothetical planet c’s with various orbital periods and eccentricities can affect the transit timing of HD 209458b. Adapted from Holman & Murray (2005).

caused by the presence of additional bodies in the system, causing the orbit of the planet and star to become non-Keplerian and the transit timing to become aperiodic. Figure 5.1 shows an example of TTVs driven by the gravitational influence of an unseen inner planet on the orbital location of the host star relative to the barycenter. This causes transit times to vary between that expected from the (linear) Keplerian ephemeris. Importantly, TTVs can be caused by the interactions of a planet with another planet, which itself does not necessarily need to be transiting – so TTV provides a way to find additional non-transiting planets in a system with known transiting planets (note RV can also do so).

TTVs can thus be used to infer the presence of additional planets in a system from the transit observations of the transiting planet(s) in that system. Figure 5.2 shows an example of the TTVs that would be induced on the hot Jupiter HD 209458b by hypothetical unseen planet c’s with periods of 19.2 – 99.8 days and eccentricities from 0.1 – 0.7. The TTV amplitude is on the order of tens of seconds, which is measurable with a sufficient number of transits. To date, TTVs have been used to detect 28 confirmed planets (https://exoplanetarchive.ipac.caltech.edu/docs/counts_detail.html).

TTVs caused by planet-planet gravitational perturbations are largest for systems near an orbital resonance. The strongest mean motion orbital resonances occur when the ratio of orbital periods of planets are in a ratio $N : (N \pm 1)$, where N is an integer. This causes conjunctions between the planets to always occur at the same orbital phase, giving each

planet a gravitational “kick.” The resulting timing variation scales as

$$\text{TTV} \sim \frac{P}{4.5N} \frac{M_{\text{pert}}}{M_{\text{pert}} + M_{\text{trans}}}, \quad (5.1)$$

where M_{pert} is the mass of the perturber and M_{trans} the mass of the transiting planet. TTVs can be as large as tens of minutes for planets in closely packed resonant chains, easily detectable for transit photometers like TESS with cadences of seconds to minutes¹. Most

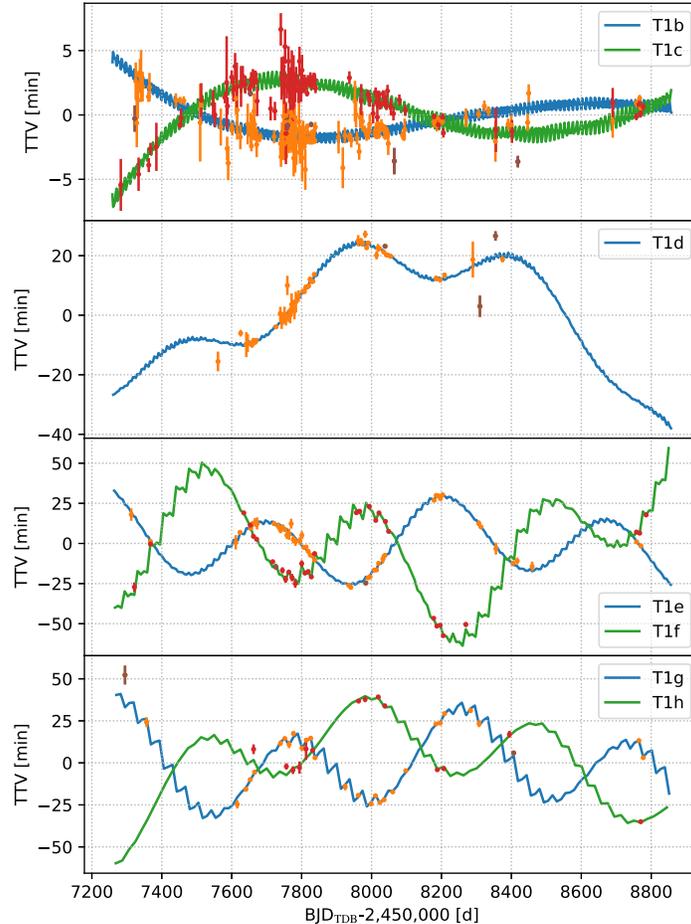


Figure 5.3: Observed TTVs of planets in the TRAPPIST-1 system, along with best fit dynamical models for the contribution from each planet. Adapted from Agol et al. (2021).

notably, the near-resonance of the seven planets in the TRAPPIST-1 system causes large TTVs (see Figure 5.3). This enables the masses of each planet to be measured (to the few % level) by observing their TTVs over long timescales and fitting them with N-body integrations Agol et al. (2021). A particular “chopping” pattern appears in models of TTVs in these near-resonant systems due to the gravitational perturbations being largest at conjunction of planet pairs. The typical period of chopping is the period between conjunctions, known as

¹TESS has an exposure cadence of 2 seconds and postage stamp cadences of 20 seconds and 2 minutes, see <https://tess.mit.edu/science/observations/>.

the synodic period (Agol & Fabrycky, 2018), which for a pair of planets with orbital periods P_1 and P_2 is equal to

$$P_{\text{syn}} = (P_1^{-1} - P_2^{-1})^{-1} . \quad (5.2)$$

5.1.2 Principles of detecting planets via timing

Timing is another indirect detection method that has been used to find 9 planets to date on its own (not combined with transits, given that as you saw above timing is a more useful tool for transiting exoplanets). Timing is useful as a detection method when the host star has some time-periodic signature that would repeat perfectly regularly ad infinitum if a planet were not orbiting it. In the presence of a planet, the motion of the host star around the common center of mass causes this timing to be aperiodic, with an amplitude related to the effect of light travel time on the otherwise periodic signature.

The timing offset due to light travel time τ in the direction of the line-of sight is

$$\tau = \frac{r_\star}{c} , \quad (5.3)$$

where c is the speed of light and r_\star is the separation of the star to the star-planet barycenter (center of mass). In general, the orbit may be offset from our line of sight, requiring the usual factor of $\sin(i)$. The expression for τ from an arbitrary viewing orientation is

$$\tau = \frac{r_\star \sin(i)}{c} . \quad (5.4)$$

For a circular orbit, we previously related $r_\star/r_p = M_p/M_\star$ (see Equation 2.7). We can thus provide a general solution for τ given circular orbits

$$\tau = \frac{r_p M_p \sin(i)}{c M_\star} \approx \frac{a M_p \sin(i)}{c M_\star} , \quad (5.5)$$

where we have made the approximation $r_p \approx a$ in the second expression. This is equivalent to Eq. 4.1 of the textbook.

Most planets found via timing have been found orbiting radio-loud millisecond pulsars. Assuming a circular, edge-on orbit and a pulsar mass of $1.35 M_\odot$ (close to the Chandrasekhar mass), the timing signature of a planet orbiting a pulsar is (textbook, Eq. 4.4)

$$\tau_p \approx 1.2 \text{ ms} \left(\frac{M_p}{M_\oplus} \right) \left(\frac{P}{1 \text{ yr}} \right)^{2/3} , \quad (5.6)$$

where we have used Kepler's third law to relate a and P . Such small timing variations are detectable due to the extreme regularity of millisecond pulsars (with spin-down rates of only $\approx 10^{-19} \text{ s}^{-1}$). However, note that the timing signature scales linearly with the planet's separation from the system barycenter (or $P^{2/3}$), so for planets with wider orbits the timing signature can be as large as a few seconds. This allows for the detection of planets via measuring the change in the timing of regularly pulsating variable stars with short pulsation periods of minutes to hours.

5.2 Finding planets via pulsar timing: group activity

The first confirmed² exoplanetary system (ever!) was found in 1992 (Wolszczan & Frail, 1992) around the pulsar PSR B1257+12. This is a system of three planets, as confirmed two years after the initial discovery of the system (Wolszczan, 1994). Planet b has a mass of $0.022 M_{\oplus}$ and an orbital period of 25.26 days. Planet c has a mass of $4.13 M_{\oplus}$ and a period of 66.54 days. Planet d has a mass of $3.82 M_{\oplus}$ and a period of 98.21 days. Please split into 6 groups – Groups 1-2 will calculate properties of planet b, Groups 3-4 will study planet c, and Groups 5-6 will study planet d.

1. First, calculate the amplitude of the timing signature caused by your planet in ms using Equation (5.6), and compare it to the 6.2 ms period of this millisecond pulsar. One of these three planets was discovered two years after the others – try to determine if your group’s planet could be the one discovered last.
2. Pulsar planets have never been found using another detection method (e.g., transits, RV, imaging, astrometry). This and the next part of the problem will help us understand why. Calculate the astrometric wobble of the host star due to your group’s planet (recall Equation 3.7 from the Day 3 notes). To do this, assume that the stellar mass is $1.35 M_{\odot}$ (which we assumed in Equation 5.6), and use the known distance to the system of 710 pc. Compare this to the Gaia detectability threshold of 0.01 mas.
3. Now calculate the calculate the transit probability of your group’s planet. Note that the stellar radius is tiny (~ 10 km), so the transit probability expression reduces to $p \approx R_p/a$. Assume that the planetary densities are 5.5 g cm^{-3} , similar to Earth. Given that there are only seven confirmed pulsar planets in the NASA exoplanet archive, discuss whether it is likely that a pulsar planet transit would be caught.
4. If time remains, discuss with your group why RV and imaging are also poor detection methods for pulsar planets like those found by Wolszczan (1994). No calculations needed, just qualitatively discuss – even though we haven’t covered direct imaging, consider whether the small separation of the planet could be angularly resolved at the distance of the system.

²As discussed on the slides, confirmation is key – the first reported pulsar planet was simply an alias, see Bailes et al., 1991 and the retraction in Lyne & Bailes, 1992. Thanks to Chris Barnet for mentioning this!

6 Detecting exoplanets: microlensing

Our agenda for Day 6 is the following:

1. Finish up timing, timing activity (20 min)
2. Principles of microlensing (15 min)
3. Microlensing derivations: Einstein radius, magnification, event length (30 min)
4. Microlensing in practice (10 min)

Today's reading is from the textbook, Ch. 5.1-5.4, or from the Gaudi review chapter. This will cover the principles of microlensing as well as the practical interpretation of microlensing light curves.

Our mid-term is coming up rapidly (it's 3 classes from now!). Note that you are allowed to bring a two-sided 8.5 by 11 inch note sheet to the exam. Everything on this sheet must be hand-written, and it must be turned in with your exam. The sheet doesn't need to include constants and Solar System planetary properties, I'll provide them. The exam will cover exoplanet detection methods, everything through what we learn this week in class (i.e., material up to and including direct imaging).

6.1 Microlensing: notes

Gravitational lensing is caused by the fact that matter distorts spacetime, affecting the trajectories of light waves as they propagate across space. If there is (near) alignment of a background light source and a massive object, the massive object can cause the formation images of the sources that are distorted ("lensed") by the object ("lens"). Gravitational lensing is commonly known through the effect of massive galaxies and galaxy clusters on radiation, leading to "strong lensing" effects (e.g., visible Einstein rings and Einstein crosses). For planets, we are concerned with "microlensing" effects where the image itself is not resolved, only the magnification of the source due to gravitational lensing.

6.1.1 Lens solution, Einstein radius

We will derive the lens solution following Paczynski (1996), see also textbook Ch. 5.2.1. Figure 6.1 shows the geometry of the lensing problem for a single lens (e.g., a single star), where S is the source, M is the lens mass, O is the observer, D_s is the distance from the observer to the source, D_d is the distance from the observer to the lens, D_{ds} is the distance from the plane of the lens to the plane of the source, A is the point at which a light ray from the source intersects the plane of the lens, R is the distance from the lens to A , R_s is the distance from the lens to the line connecting the source and observer, I is the position at which the line of sight to the image would be present on the source plane if there were no light deflection, and α is the angular deflection of the light ray as a consequence of general relativity. We can further define θ_s as the angle between the line of sight to the lens and the line of sight to the source, and θ_I as the angle between the line of sight to the lens and the light of sight to the image.

The deflection angle of the light ray α is given from general relativity as

$$\alpha = \frac{4GM}{Rc^2} , \quad (6.1)$$

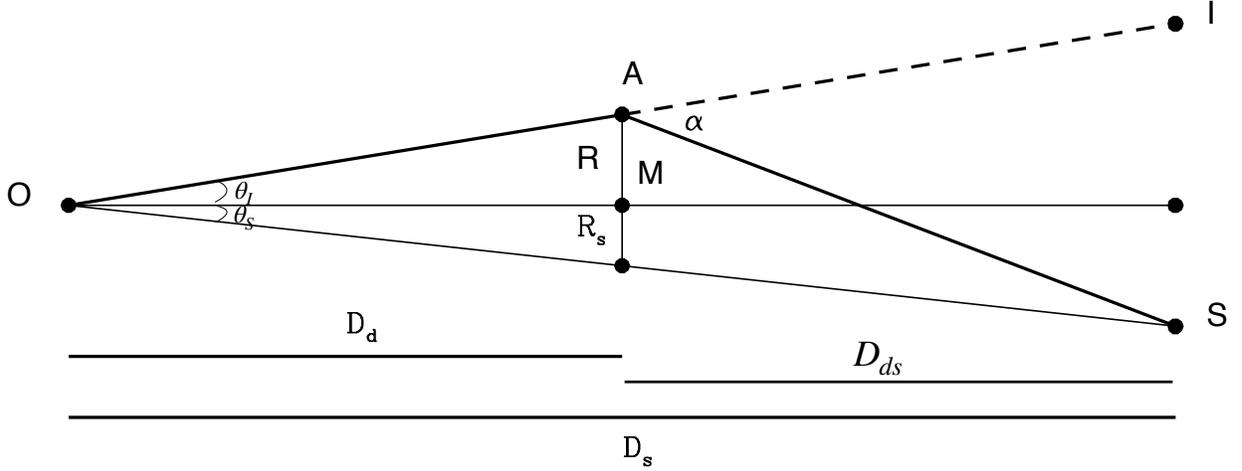


Figure 6.1: Geometry of gravitational lensing. The observer is to the left (O), the source is at the right (S), and the lensing mass is M. Adapted from Paczynski (1996).

where M is the mass of the lens and c is the speed of light. Note that the Schwarzschild radius $R_{Sc} = 2GM/c^2$, so $\alpha = 2R_{Sc}/R$. We can further define the lens (M) as having an angular position (x_m, y_m) on the sky, and the observer as looking at the sky in the direction of angular coordinates (x, y) . Doing so, we can write the position of M on the lens plane as $(X_m = x_m D_d, Y_m = y_m D_d)$, and the position of M on the source plane as $(X_{M,s} = x_m D_s, Y_{M,s} = y_m D_s)$. We can also write the location at which the line of sight of the observer intersects the lens plane at point A as $(X_A = x D_d, Y_A = y D_d)$, and the coordinates for point I on the source plane as $(X_I = x D_s, Y_I = y D_s)$. With this, we can then break the angle of deflection α into two components

$$\alpha_x = \alpha \frac{X_A - X_M}{R}, \alpha_y = \alpha \frac{Y_A - Y_M}{R}. \quad (6.2)$$

From these components, we can determine the coordinates of S on the source plane as

$$X_s = X_I - \alpha_x D_{ds}, Y_s = Y_I - \alpha_y D_{ds}. \quad (6.3)$$

We can now use the small angle equation to relate the distances on the lens plane to the distances on the source plane:

$$\frac{R + R_s}{D_d} = \frac{\sqrt{(X_s - X_I)^2 + (Y_s - Y_I)^2}}{D_s}. \quad (6.4)$$

We can now use Equation (6.3) to relate the distances on the lens plane to the distances to the lens and source

$$R + R_s = \alpha D_{ds} \frac{D_d}{D_s} = \frac{4GM}{Rc^2} \frac{D_{ds} D_d}{D_s}. \quad (6.5)$$

This can be recast as a quadratic lens equation

$$\frac{R_s}{R_E} = -\frac{R}{R_E} + \frac{R_E}{R} \rightarrow R^2 + R_s R - R_E^2 = 0, \quad (6.6)$$

where the linear Einstein ring radius

$$R_E = \sqrt{2R_{Sc}D} = \sqrt{\frac{4GM}{c^2} \frac{D_{ds}D_d}{D_s}} \quad (6.7)$$

with R_{Sc} the Schwarzschild radius and the effective lens distance $D = D_{ds}D_d/D_s$. Using the quadratic formula, we can write the solutions for Equation (6.6)

$$R_{\pm} = 0.5 \left[R_s \pm \sqrt{R_s^2 + 4R_E^2} \right], \quad (6.8)$$

where note that there are two solutions that correspond to two images of the same source located on opposite sides of the lens at angular distances of R_+/D_d and R_-/D_d .

Note that the derivation in the textbook is in angular position rather than physical position (i.e., for θ_S and θ_I). The textbook's version of the lens equation is analogous to the above,

$$\theta_S = \theta_I - 2R_{Sc} \frac{D_{ds}}{D_d D_s} \frac{1}{\theta_I} \rightarrow \theta_I^2 - \theta_S \theta_I - \theta_E^2 = 0, \quad (6.9)$$

where now $\theta_E = \sqrt{2R_{Sc}D_{ds}/(D_d D_s)}$, which is related to R_E by $R_E = \theta_E D_d$. One can also write down expressions for R_E and θ_E in terms of relevant numerical quantities,

$$R_E \approx 8.1 \text{ au} \left(\frac{M}{M_{\odot}} \right)^{1/2} \left(\frac{D_s}{8 \text{ kpc}} \right)^{1/2} \left(\frac{D_d D_{ds}}{D_s^2} \right)^{1/2}, \quad (6.10)$$

$$\theta_E \approx 1.0 \text{ mas} \left(\frac{M}{M_{\odot}} \right)^{1/2} \left(\frac{D_d}{8 \text{ kpc}} \right)^{-1/2} \left(\frac{D_{ds}}{D_s} \right)^{1/2}. \quad (6.11)$$

These angular scales are too small to resolve with most ground-based instruments, which led to the nomenclature ‘‘microlensing’’ due to the images not being resolved. However, these spatial scales of 5-10 au are prime for detecting planets near the ice lines of the protoplanetary disks from which they formed, allowing microlensing to probe a novel region of parameter space (especially relative to transits, which is highly biased toward closer-in orbits). Figure 6.2 shows the resulting appearance of the images of a lensed source as it passes through the Einstein ring of a point mass. Note the two image paths that stay on opposite sides of the source as it is differentially imaged while it passes near the line of sight to the lens.

6.1.2 Peak magnification

The microlensing event causes a magnification of the source due to the two lensed images being brighter than the source itself. This magnification is time-dependent, and peaks at the time of closest angular separation of the source to the lens (see Figure 6.3). Gravitational lensing conserves surface brightness, so the ratio of the image to source intensity (i.e., magnification) is given by the ratio of the area of the image to the area of the source on the lens plane. The magnification A for each image is related to the image position and source position on the lens plane and the derivative of the image position with respect to the source position as

$$A_{\pm} = \left| \frac{R_{\pm}}{R_s} \frac{dR_{\pm}}{dR_s} \right| = \frac{u^2 + 2}{2u\sqrt{u^2 + 4}} \pm 0.5, \quad (6.12)$$

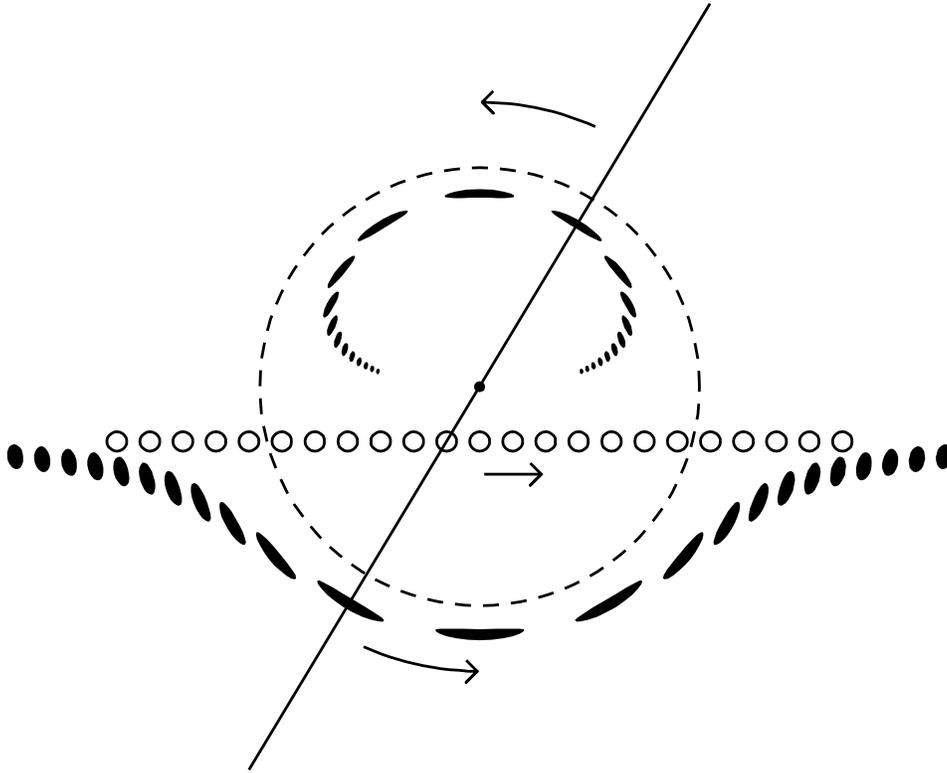


Figure 6.2: The geometry of gravitational lensing, showing the trajectory of the source (open circles) and the images (filled ellipses). Adapted from Paczynski (1996).

Where the derivative of Equation (6.8) was taken using the chain rule, and $u \equiv R_s/R_E$. The total magnification is

$$A = A_+ + A_- = \frac{u^2 + 2}{u\sqrt{u^2 + 4}}, \quad (6.13)$$

which is always larger than one. Also note that the difference between the magnification of the two images is constant, $A_+ - A_- = 1$. For $u \ll 1$, $A \approx u^{-1}$, and for $u \gg 1$, $A \approx 1$. As a result, the magnification drops to 1 (no magnification) far from the Einstein ring, with the magnification during a microlensing event scaling inversely with the ratio of the separation of the source from the Einstein ring – as a result, events with separations near the Einstein ring have the largest detectable microlensing magnification. Note that though the magnification can be formally infinite for a source that has an angular separation of zero from the lens, in practice the magnification is always finite.

6.1.3 Planetary perturbation

Planets cause their own microlensing event that is imprinted upon the larger magnification (and longer duration, see next section) event of their host star. If the planet perturbation is before or after the source crosses the Einstein ring, the planet causes a single trough (if before stellar Einstein ring crossing, due to reducing the brightness of the image within the Einstein ring), or peak (if after stellar Einstein ring crossing, due to splitting the image outside the Einstein ring), see Figure 6.4. However, if the source crosses the planet within the

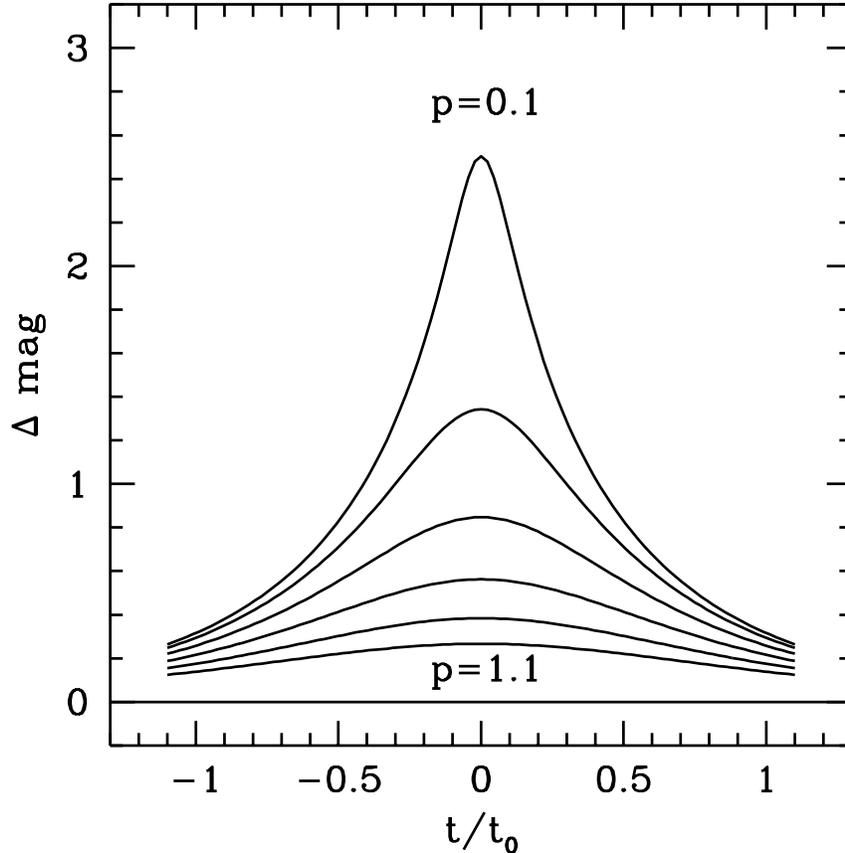


Figure 6.3: The magnification due to a point lensing event as a function of time, normalized by the time it takes the source to move across the Einstein ring t_0 . Here p is the impact parameter, which is the smallest angular distance between the source and the lens (R_s) in units of the angular Einstein ring radius. Adapted from Paczynski (1996).

Einstein ring, the behavior can be more complex, resulting in a multiple-peaked structure of the magnification curve. This is because the planet distorts the magnification field of the host star, and if the source crosses a caustic (region of distortion) the planet effectively induces an astigmatism in the lensing pattern. Caustic crossings can occur for star-planet separations between $0.5 - 2\theta_E$.

One common point of confusion is that we detect planets lensing the far background source, *not* their own host star. This is because the Einstein radius $R_E \propto \sqrt{D_{ds}}$, and so if the planet and star are very nearby the Einstein radius of the planet lensing its host star is very small. The magnification then scales as R_E/R_s , so the magnification due to the planet lensing event of the host star is undetectable.

Lastly, note that planets without a host star (“rogue planets”) can also cause microlensing events. These are short-duration and small-magnification events that look like standard single microlensing events but with masses that are clearly below the Deuterium burning limit. As a result, microlensing is the only detection method that allows the study of planets without either directly seeing their radiation or inferring their presence from their indirect effects on a host star.

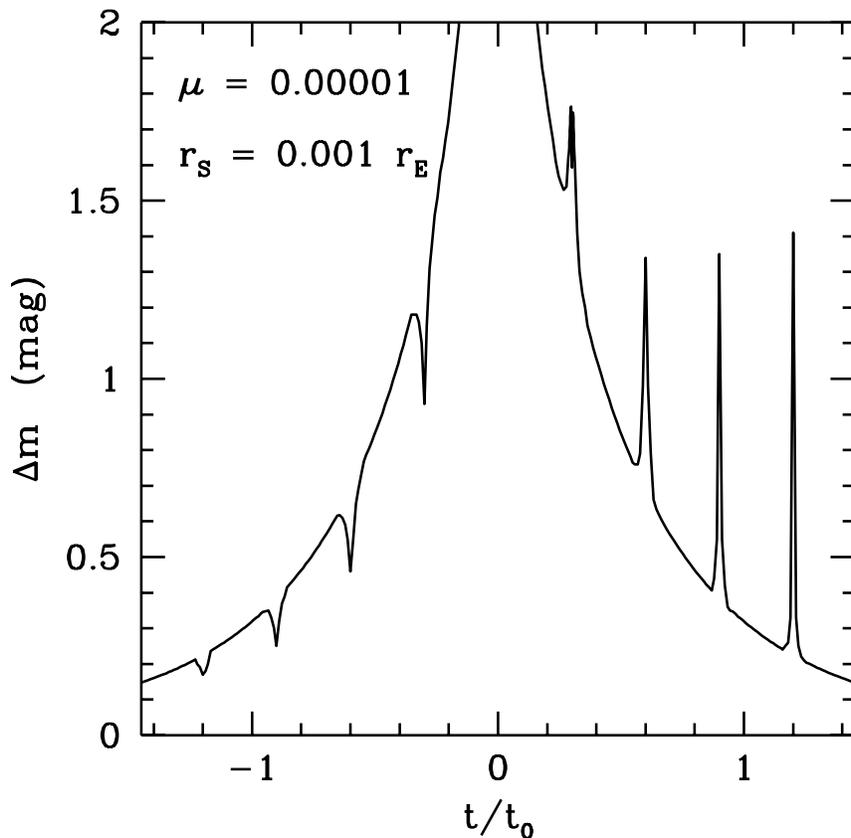


Figure 6.4: The magnification of a hypothetical planetary system with eight equal-mass planets located in a straight line. Note that the planetary perturbations create local minima before Einstein ring crossing and maxima after Einstein ring crossing. Adapted from Paczynski (1996).

6.1.4 Event length

The typical timescale for a microlensing event is the time it takes a typical star in the bulge to cross the Einstein ring. We can express this as $\tau_E = R_E/v_t$, where v_t is the tangential velocity of the star (i.e., on the plane of the sky). Using our expression for the Einstein ring radius, Equation (6.10), we can write this as

$$\tau_E = 0.214 \text{ yr} \left(\frac{M}{M_\odot} \right)^{1/2} \left(\frac{D_d}{10 \text{ kpc}} \right)^{1/2} \left(1 - \frac{D_d}{D_s} \right)^{1/2} \left(\frac{200 \text{ km s}^{-1}}{v_t} \right). \quad (6.14)$$

Importantly, $\tau_E \propto \sqrt{M}$, which means that planetary event durations are much shorter than stellar microlensing events. Planetary events typically last less than a day, and for Earth-mass planets the event durations are on the order of 3 – 5 hr. This is why microlensing searches typically look for longer-term increases in light due to a stellar microlensing event and then use high cadence observations with multiple facilities to search for a planetary companion.

6.1.5 Microlensing in practice

To date, microlensing has discovered 210 planets, with the first discovery of a planet via microlensing (OGLE-2003-BLG-235L b) in 2004. This makes gravitational microlensing the third-most prolific exoplanet discovery method after transit and radial velocities. All planets discovered to date have been detected from the ground, primarily from 3 surveys,

OGLE (1992-), KMT (2009-), and MOA (2006-2014). Each of these surveys use a wide etendue to search the central galactic bulge through Baade’s window. These surveys cover $\approx 8 \text{ deg}^2$ of sky and probe thousands of events per year, most of them stellar without a planetary perturbation. An early microlensing discovery of a planet (OGLE-2005-BLG-390b) is shown in Figure 6.5. Note that the full lensing event lasts over 50 days, which allows for follow-up from a wide range of ground-based observatories to characterize the shape of the stellar and planetary components of the source magnification. These surveys have also

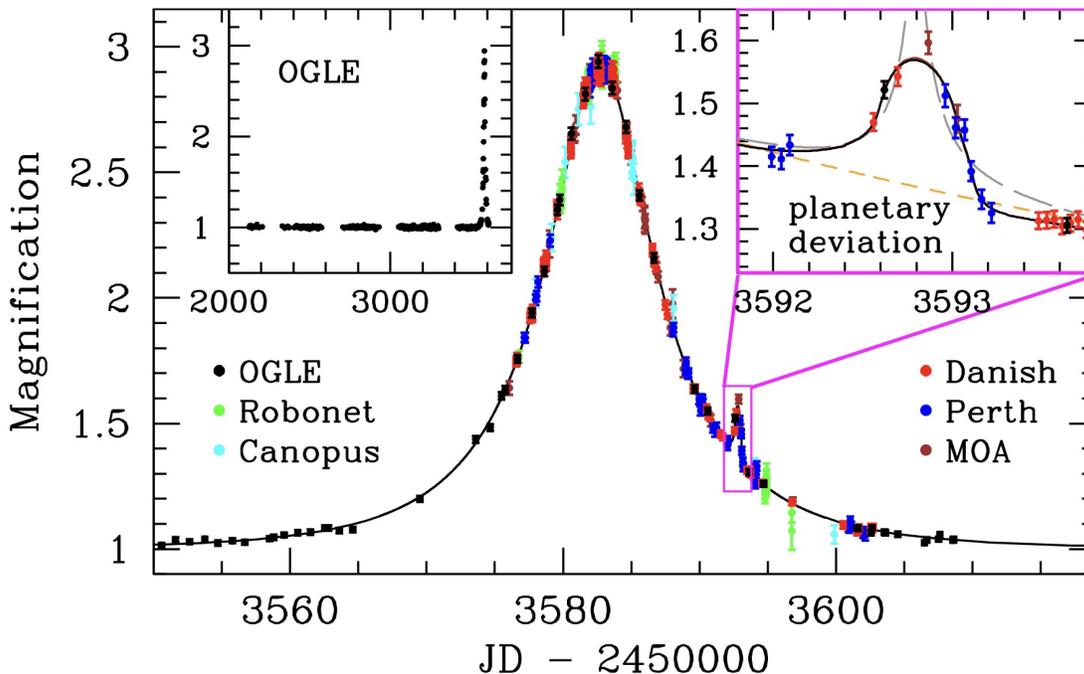


Figure 6.5: The discovery light curve of OGLE-2005-BLG-390 (January 2006). Here the planet causes a small perturbation well after the primary microlensing event. The planet has a mass of $\approx 5.5 M_{\oplus}$ and a orbital period of ≈ 3500 days, and its host star lies at a distance of 6.6 kpc.

found several free-floating planetary candidates, with an expected abundance of $\sim 1 - 2$ per main-sequence star – potentially implying that each forming planetary system leads to the loss of approximately one planet to interstellar space.

Space-based observatories have largely been used to follow up ground-based microlensing signals to date, with few search campaigns that have not found clear evidence of new planetary candidates. However, the Nancy Grace Roman Space Telescope (NGRST, formerly WFIRST) is an upcoming wide-field space-based survey that may revolutionize the field of microlensing. Roman will launch in 2027 to L2 with a 2.4 m mirror, comparable to Hubble but with a 100x larger field of view (0.28 deg^2). Roman is expected to find $\gtrsim 1500$ planets via a wide-field microlensing survey as part of its science objectives.

7 Detecting exoplanets: direct imaging

Our agenda for Day 7 is the following:

1. Microlensing: recap (5 min)
2. Microlensing: magnification derivation (in groups, 20 min)
3. Microlensing in practice, event length, Roman (10 min)
4. Direct imaging intro activity (10 min, skip if past 40 minutes in)
5. Direct imaging: contrast (10 min)
6. Direct imaging in practice (20 min)

Today's reading is textbook Ch. 7.1-7.5 and/or the Traub & Oppenheimer handout. This will cover the fundamentals of direct imaging as well as the practicalities of how direct imaging is conducted using adaptive optics and coronagraphic masks.

7.1 Direct imaging intro activity

This activity is meant to demonstrate why the current observational characterization of exoplanets via direct imaging is limited to young, massive planets orbiting at wide separations from their host star.

1. Groups 1, 2: Calculate the angular separation between Earth and the Sun. Group 1 - calculate this at a distance of 10 pc. Group 2 - calculate this at a distance of 100 pc.
Groups 3, 4: Calculate the angular separation between the Sun and Jupiter. Group 3 - calculate this at a distance of 10 pc. Group 4 - calculate this at a distance of 100 pc.
Groups 5, 6: Calculate the angular separation between HR 8799 and HR 8799b, which has a semi-major axis of 71.6 au. Group 5 - calculate this at a distance of 10 pc. Group 6 - calculate this at a distance of 100 pc.
2. For each of your determined angular separations, calculate the approximate size of a telescope required to detect the planet at a wavelength of $0.6 \mu\text{m}$ around the star at the given distance (assuming diffraction-limited observations).
3. Estimate the visible light star-planet brightness ratio for your planetary system. To do so, assume that the planet has an albedo of 1 (i.e., that it is perfectly reflective), and that we are observing the planet at full phase (i.e., when the full disk is illuminated). For the HR 8799 group, assume that HR 8799b has the same radius as Jupiter.

7.2 Direct imaging: notes

7.2.1 Planet-star contrast

There are two primary components of the light we observe from any planet: the thermal emitted light from the planet itself, and reflected light from the host star. Detections of directly imaged planets to date have been in thermal emission, with these detections finding young giant planets at wide separations from their host stars. This is because giant planets

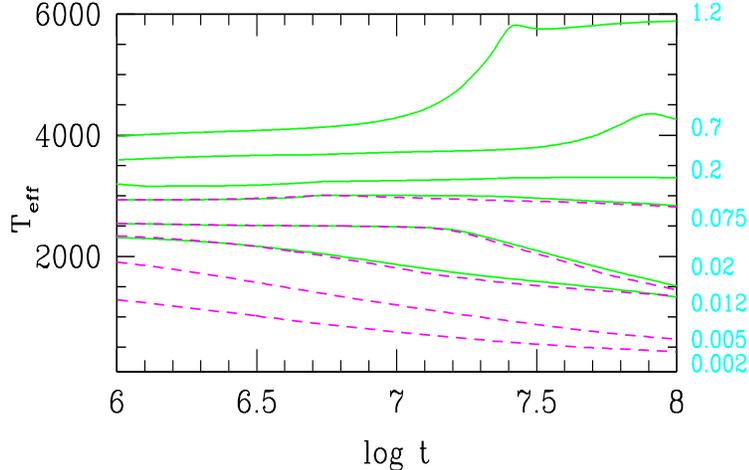


Figure 7.1: Cooling curves over the first 100 Myr of evolution of stars, brown dwarfs and planets. The time is in seconds, temperature in K, and mass of objects in the light blue on the right in M_{\odot} . Adapted from Baraffe et al. (2002).

form with high effective temperatures, and cool from this hot initial condition over time (see Figure 7.1).

The primary quantity to characterize the detectability of a planet via direct imaging is planet-to-star flux ratio, or “contrast” of the planet with respect to the star. The wavelength-dependent contrast can be expressed as

$$f_{\text{em}} = \left(\frac{R_p}{R_{\star}} \right)^2 \frac{B_{\lambda}(T_p)}{B_{\lambda}(T_{\star})} \Phi_{\text{em}}(\lambda, \alpha), \quad (7.1)$$

where B_{λ} is the Planck function and Φ_{em} is the phase function for emission, which depends on wavelength and the star-planet-observer “phase angle” α . For a circular orbit, $\cos \alpha = \sin(\theta + \omega) \sin(i)$. Many observations of planets via direct imaging are in the Rayleigh-Jeans tail of the Planck function, simplifying the contrast to

$$f_{\text{em}} \approx \left(\frac{R_p}{R_{\star}} \right)^2 \frac{T_p}{T_{\star}} \Phi_{\text{em}}(\lambda, \alpha). \quad (7.2)$$

For young planets, the effective temperature of the planet is much greater than the effective temperature the planet would have if it were in thermal equilibrium with the instellation it receives from the host star (Figure 7.1). For mature planets that have cooled off from formation (or more generally, planets that receive much more incident stellar power than their intrinsic cooling luminosity), one can approximate the planetary temperature by the equilibrium temperature

$$T_{\text{eq}} = T_{\star} [f(1 - A_B)]^{1/4} \sqrt{\frac{R_{\star}}{a}}. \quad (7.3)$$

In Equation (7.3), f is a factor that accounts for redistribution of the received heat from the star across the planet by e.g., atmospheric (and/or oceanic) circulation, and is $f = 1/4$ for full redistribution of the incident stellar radiation (if the thermal energy is equally radiated to space over the entire surface of the planet). A_B is the Bond albedo, which is the fraction of total energy incident on the planet that is not absorbed and re-radiated (i.e., the amount that is scattered/reflected by clouds, haze, ice, and gas).

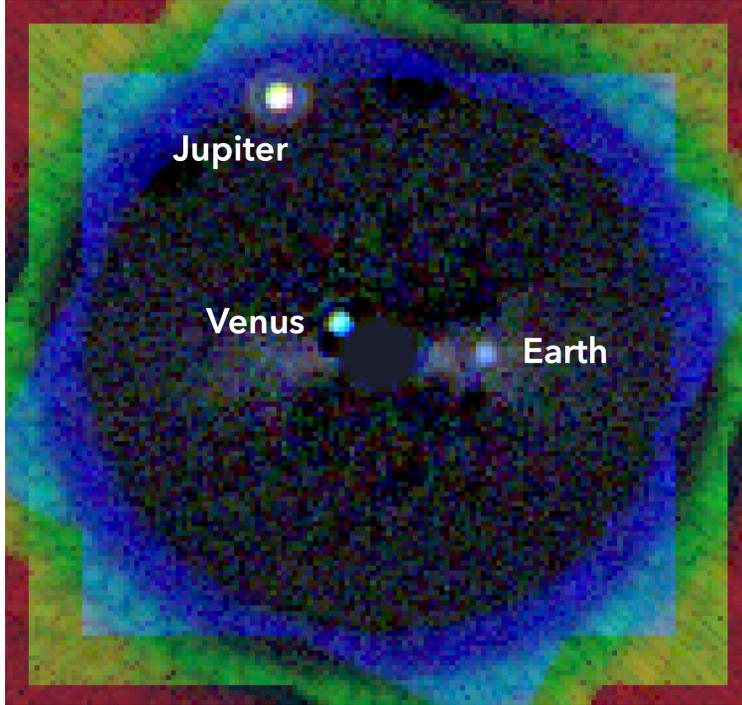


Figure 7.2: Simulated image of our Solar System as viewed in reflected light from a distance of 10 pc by a Habitable Worlds Observatory-like mission. The specific simulations here are for the LUVOIR-A mission concept.

In order to directly image planets that are not young, observations have to probe the reflected light from the planet that originates from the star. Planets have not yet been directly imaged in reflected light, but Figure 7.2 shows an example image of our Solar System in the optical from a distance of 10 pc. The wavelength-dependent contrast in reflected light can be written similarly to that in emitted light as

$$f_{\text{ref}} = \left(\frac{R_p}{a} \right)^2 A_g(\lambda) \Phi_{\text{ref}}(\lambda, \alpha) , \quad (7.4)$$

where A_g is the wavelength-dependent geometric albedo and Φ_{ref} is the phase function in reflected light.

The typical planet-to-star contrasts in thermal emission for young giant planets $f_{\text{em}} \sim 10^{-4} - 10^{-6}$, which is currently achievable with ground-based telescopes that use adaptive optics and coronagraphy (see the next section). Those in reflected light are much smaller – for Jupiter around the Sun, the reflected light contrast is $f_{\text{ref}} \sim 10^{-8}$, and for Earth the contrast is $f_{\text{ref}} \sim 10^{-9}$.

7.2.2 Technological challenges

In order to access the emitted or reflected light from a star, an optical device must be used to suppress the light from that star. The mask that is placed in the focal plane of the telescope is called a coronagraph, given that early development of such a device in the late 1800s and early 1900s was for the purpose of blocking light from the Solar photosphere to reveal the corona. Figure 7.3 shows the optics of a Lyot coronagraph, which was originally implemented in 1931 by Lyot to observe the Sun. This uses a system of three lenses, along with two optical masks. The first objective lens forms an image of the star, and an occulting mask then blocks the disk of the star. There is still significant diffracted light after the

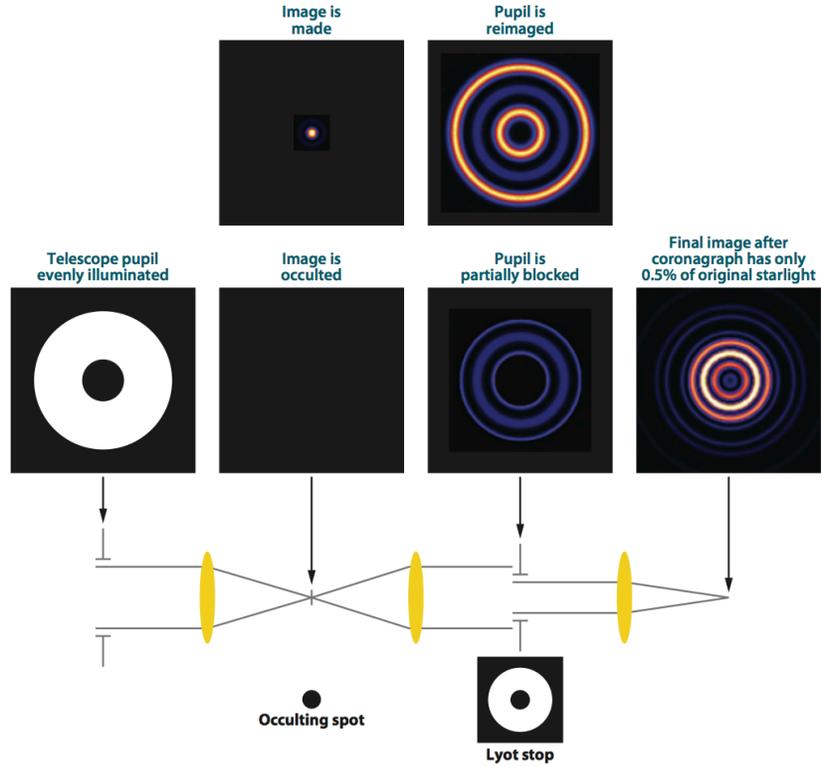


Figure 7.3: Schematic of coronagraphic optics with an occulting spot and internal Lyot stop.

image passes through the occulting mask, and so the light is then put through a field lens that re-images the diffraction pattern, with a Lyot stop then intercepting the diffraction ring while allowing light from the rest of the image to pass through. Finally, the second objective lens places the image on the plane of the detector. The coronagraph suppresses the stellar light within a given region, beyond which the noise is small and planets with a given contrast are detectable. This angular separation between the host star and the region where planets are detectable is termed the “inner working angle” (IWA) of the instrument.

This basic coronagraphic method has been improved with more detailed versions of the Lyot coronagraph, including using pupil apodisation to modify the shape of the point spread function upon entering the optics, using interferometric coronagraphs that remove the diffraction pattern by destructive interference, and using phase masks to shift the light in the focal plane and lead to destructive interference. More recent developments include the vortex coronagraph, which is a high performance phase mask that phase shifts light by transforming its wavefront from planar to helical, resulting in total destructive interference in a dark central core. Potential future developments in coronagraphy include the starshade concept, where an external occulter with a size of ~ 50 m formation flies at a separation of $\sim 75,000$ km from a space telescope (values are approximately those for the HabEx mission concept) and with very precise tolerances (± 1 m position, < 1 mm shape) can occult light from nearby stars to reveal companion planets with contrasts down to 10^{-10} .

Even with a coronagraph, there are still distortions (“speckles”) due to either the at-

mosphere (for ground-based direct imaging) and/or deformations in the mirror and optical system (relevant for both ground- and space-based observations) that must be removed in order to isolate the planetary signal. Astronomers use two primary methods in conjunction to deal with these: 1) Adaptive optics to move the mirror and compensate for phase fluctuations, 2) Differential imaging techniques to remove the noise pattern and clean the image. Adaptive optics is a technique that couples actuators to the telescope mirror itself in order to deform the mirror in a way that compensates for atmospheric turbulence. These are very rapid adjustments that are constantly being made on ~ 1 ms timescales across the mirror, with the number of actuators required scaling with $(D/r_o)^2$, where D is the telescope diameter and r_o is the atmospheric wavefront coherence length (typically ~ 0.2 m in the visible, ~ 1 m in the NIR). As a result, current ground-based direct imaging surveys use $\sim 10^3$ actuators, while future surveys with the ELTs will use $\sim 10^4$ actuators. These adaptive imaging systems are often focused using laser guide stars, which send a laser beam pulse from the ground to the upper atmosphere that either excite Na in the mesosphere or use shorter-wavelength Rayleigh scattering to make an artificial “star” that can be used to focus the adaptive optics system and account for atmospheric turbulence.

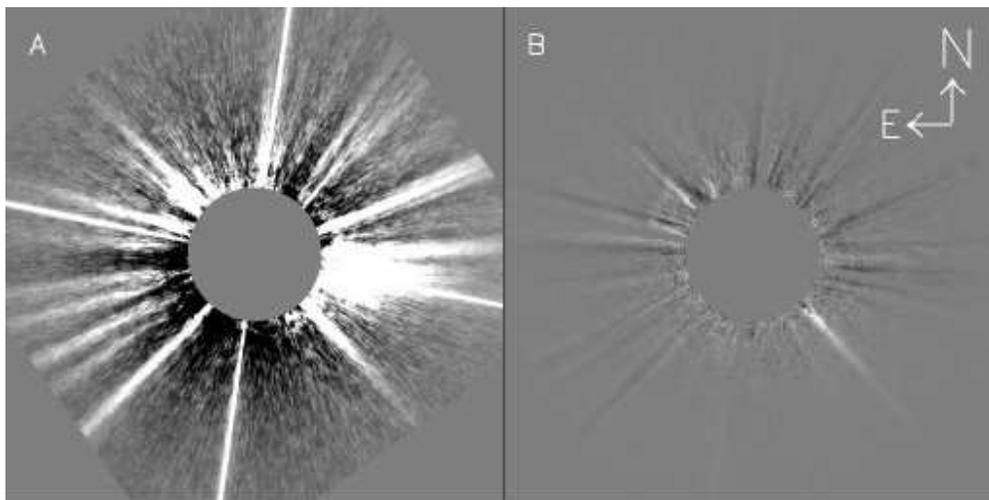


Figure 7.4: A) Image of Vega after flat field normalization, bad pixel correction, distortion correction, but before ADI. B) Image of Vega after a single ADI difference. Adapted from Marois et al. (2006).

The noise pattern is removed for both ground-based and space-based imaging by differential imaging techniques. The most common is angular differential imaging (ADI), which rotates the image and uses the fact that the noise pattern is set by the instrument and optics to then subtract out the noise when combining rotated images. Angular differential imaging from the ground uses the Earth’s rotation over the course of the night to combine images, while space-based observations physically roll the telescope in order to rotate the image. Other forms of differential imaging include reference differential imaging (RDI) where observations of a reference star (ideally without companions) are subtracted, or spectral differential imaging (SDI) where the speckle pattern is suppressed by separating the light from a planetary absorption or emission feature from the stellar spectrum.

8 Detecting exoplanets: inter-comparison of detection techniques (Day 8)

Our agenda for Day 8 is the following:

1. Wright & Gaudi for the modern day activity (30 minutes)
2. Strengths, weaknesses, and resulting biases of each detection method (30 minutes)
3. Open time for questions recapping planet detection and the upcoming midterm (5-10 minutes)
4. Highlights of exoplanet detections (if time remains, finish next class)

Today's reading is the Wright & Gaudi handout. These will cover the population of detections from each detection method as of 2011 (the handout is dated), as well as the strengths and weaknesses of each detection method. The handout is quite long, so it's okay if you read just Chapters 1 and 2 for this class (we'll also cover chapter 3 if you have time to read it).

8.1 Activity: Wright & Gaudi for the modern day

Figure 8.1 below shows my own version of the mass-semimajor axis distribution of exoplanets discovered by transit, RV, microlensing, imaging, astrometry, and timing. In this activity, we'll derive the appropriate sensitivity curves for each technique and manually overplot them on this diagram, provided I can work the projector properly. Please split into five

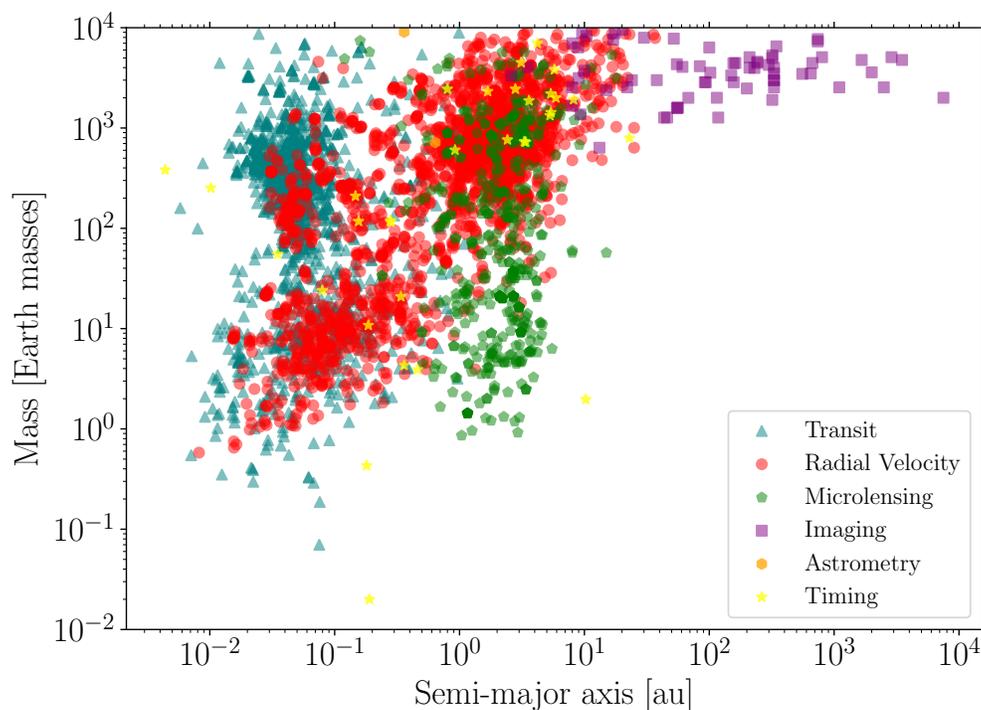


Figure 8.1: Mass vs. semi-major axis of exoplanets discovered by various techniques. Data from the NASA Exoplanet Archive.

groups (3 people each). Each of these five groups will work on one of the five key detection techniques in order to derive approximate mass-semimajor axis detection limits from each method. When you're done deriving your detection limit, please plot it neatly on the blackboard and be prepared to describe how you derived it to the class on the whiteboard. If you finish before another group, please go help that group (a couple of these are harder than others).

1. **Radial velocity.** From the expression in the notes for the radial velocity semi-amplitude (Equation 2.10), derive an equation for the minimum planet mass as a function of semi-major axis that can be detected for planets orbiting Sun-like stars assuming that the minimum detectable RV semi-amplitude is ~ 1 m/s. Over-plot this line on the class' plot of planet mass vs. semi-major axis, and label it.
2. **Transit.** From the transit depth (Equation 4.9 in the notes), write down a scaling expression for the minimum detectable planet mass detectable via transit. To do this, assume that the masses and radii of stars are linearly proportional to one another, and that all exoplanets have a fixed density. Note that the transit SNR approximately scales as the transit depth times the square root of the number of photons obtained during transit events over the survey duration. Then over-plot and label a line on our class plot for Kepler-like transit surveys on the plot assuming that these surveys have a threshold capability to detect the transits of Earth-sized planets orbiting Sun-like stars with semi-major axes of 1 au.
3. **Astrometry.** Starting from the expression for astrometric wobble in Equation (3.6) of the notes, derive the minimum planet mass that is detectable with Gaia as a function of semi-major axis. Assume that Gaia has a sensitivity of 0.01 mas for the brightest stars, and further assume that the host stars are Solar-type and that their typical distances are 10 pc. Over-plot this line on the class' graph and label it.
4. **Microlensing.** Present-day microlensing surveys can find planets with masses of $\approx 1 M_{\oplus}$ around low-mass stars. Assuming that the typical lens star has a mass of $0.5 M_{\odot}$ and a distance of 4 kpc, with source stars typically located at 6 kpc, determine at which semi-major axis there should be a minimum in the microlensing sensitivity curve. Then, assuming that the microlensing sensitivity drops to zero within a factor of ten in semi-major axis in both directions, over-plot and label the microlensing sensitivity curve on the class' plot.
5. **Direct imaging.** Use the population of discovered exoplanets with ground-based direct imaging to motivate a region of parameter space in which direct imaging is sensitive. Specifically, choose a threshold mass above which directly imaged self-luminous planets are detectable, and a semi-major axis that corresponds to the current inner working angle of ground-based observatories. Plot a line (or two) on our graph that boxes in this region where direct imaging from the ground can find planets, and label it.

8.2 Strengths and biases of each detection method

The following discussion is based on Wright & Gaudi (2013).

8.2.1 Radial velocity

The radial velocity semi-amplitude is (Equation 2.10)

$$K = \left(\frac{2\pi G}{P} \right)^{1/3} \frac{M_p \sin(i)}{(M_\star + M_p)^{2/3}}. \quad (8.1)$$

This semi-amplitude is also the signal that we measure in radial velocity, thus the radial velocity signal-to-noise ratio scales as

$$(S/N)_{\text{RV}} \propto P^{-1/3} M_p M_\star^{-2/3}, \quad (8.2)$$

assuming that $M_p \ll M_\star$. Using Kepler's third law ($P \propto a^{3/2} M^{-1/2}$), the radial velocity signal-to-noise ratio scales as $(S/N)_{\text{RV}} \propto a^{-1/2} M_p M_\star^{-1/2}$. We can invert this to find the scaling of minimum mass that can be found for a given signal-to-noise ratio

$$M_{p,\text{min}} \propto a^{1/2} M_\star^{1/2}, \quad (8.3)$$

which demonstrates that (ignoring stellar noise) radial velocity is sensitive to lower-mass planets at smaller separations from lower-mass stars. In reality, M dwarf stars are very noisy due to stellar activity, leading G and K dwarfs to be the optimal stellar types to search for planets around with RV surveys.

8.2.2 Transit

The signal for the transit method is the transit depth, as in Equation (4.9),

$$\delta = \left(\frac{R_p}{R_\star} \right)^2, \quad (8.4)$$

where a greater number of transits over a given observation timescale increases the signal relative to the noise. Thus, the transit signal scales as

$$(S/N)_{\text{tr}} \propto \frac{R_p^2 a}{R_s^2 P} \propto \frac{R_p^2 a M_\star^{1/2}}{M_\star^2 a^{3/2}} \propto R_p^2 a^{-1/2} M_\star^{-3/2}, \quad (8.5)$$

where we have used Kepler's third law and assumed that $R_\star \propto M_\star$, which is valid for $M \lesssim M_\odot$. As a result, the minimum detectable planet radius

$$R_{p,\text{min}} \propto a^{1/4} M_\star^{3/4}. \quad (8.6)$$

Additionally, the transit probability scales as $p \propto M_\star a^{-1}$. Together, these imply that smaller planets are more readily detectable via transit at closer separations around less massive stars, with a slightly weaker dependence on separation and sharper dependence on stellar mass than radial velocities. However, the transit probability strongly biases the population toward small semi-major axes due to its inverse dependence on separation. One other important factor in broad surveys is simply the population of nearby stars, which is heavily weighted toward low-mass M dwarfs (which make up $\sim 70\%$ of the population), further enabling surveys to find planets around small, cool stars.

8.2.3 Direct imaging

The imaging signal-to-noise ratio scales directly with the planet-to-star contrast, which we previously stated for reflected light (Equation 7.4) and thermal emission (Equation 7.2)

$$(S/N)_{di,ref} \propto R_p^2 a^{-2}, \quad (8.7)$$

$$(S/N)_{di,em} \propto R_p^2 R_\star^{-2} T_p T_\star^{-1}, \quad (8.8)$$

where the latter assumes that observations are in the Rayleigh-Jeans tail where $B \propto T$. Assuming that the planet is in radiative equilibrium, we can scale the planetary temperature with the equilibrium temperature $T_{eq} \propto T_\star R_\star^{1/2} a^{-1/2}$, implying that the emission SNR scales as

$$(S/N)_{di,em} \propto R_p^2 R_\star^{-3/2} a^{-1/2}. \quad (8.9)$$

Direct imaging also requires that the planet is beyond the inner working angle of the instrument, with a separation $a > \theta_{IWA} d^{-1}$, where d is the distance to the system. Thus, in thermal emission hotter and larger planets around smaller stars at wider separations and closer distances are more detectable. In reflected light, the inverse-square law causes larger planets closer to the star (but still beyond the inner working angle) to be more detectable. There is also an effect from the host star type for reflected light, as planets will have less reflected light at the same separation around smaller stars. Putting these together, current surveys are only sensitive to planets with $M \gtrsim M_{Jup}$ at wide ($\gtrsim 10$ au) separations around young nearby stars.

8.2.4 Microlensing

Microlensing is most sensitive to planets that have semi-major axes near the Einstein ring radius of their host stars, which is (Equation 6.10):

$$R_E \approx 8.1 \text{ au} \left(\frac{M}{M_\odot} \right)^{1/2} \left(\frac{D_s}{8 \text{ kpc}} \right)^{1/2} \left(\frac{D_d D_{ds}}{D_s^2} \right)^{1/2}. \quad (8.10)$$

As a result, the optimal separation of microlensing scales as

$$a_{ml,opt} \propto M_\star^{1/2}, \quad (8.11)$$

but there is no simple way to write down a signal-to-noise scaling as for other methods. The sensitivity to the host star mass is largely dependent on the event rate, which in turn is related to how many lenses there are in the line of sight from the observer to the bulge of the Milky Way, weighted by the host star Einstein ring radius (which scales as $M_\star^{1/2}$). As a result, microlensing is most sensitive to planets around stars that are less massive than the Sun (because they are more numerous), with a peak at $\approx 0.5 M_\odot$.

8.2.5 Astrometry

The angular astrometric shift of a host star due to an unseen companion planet is (Equation 3.6)

$$\alpha = \frac{a}{d} \frac{M_p}{M_\star + M_p} \approx \frac{a}{d} \frac{M_p}{M_\star}, \quad (8.12)$$

and thus the signal-to-noise ratio of astrometric detections scales as

$$(S/N)_{as} \propto ad^{-1}M_pM_\star^{-1} . \quad (8.13)$$

Inverting this for the minimum detectable planet mass, we find

$$M_{p,\min} \propto a^{-1}d M_\star . \quad (8.14)$$

This implies that astrometry is able to find smaller planets at wider separations and smaller distances from the observer around less massive stars. There is also a hidden trade-off in noise, given that less massive stars (which have greater astrometric shifts for a given planet mass) emit fewer photons, increasing the noise for a given distance.

8.3 Key findings from each detection method

To date (2/20/24), there are 5573 discovered exoplanets, with 4153 discovered by transit, 1075 by radial velocity, 210 by microlensing, 68 by imaging, 54 by timing variations, and 3 by astrometry (see Figure 8.2). There has been a steady increase in RV detections over time,

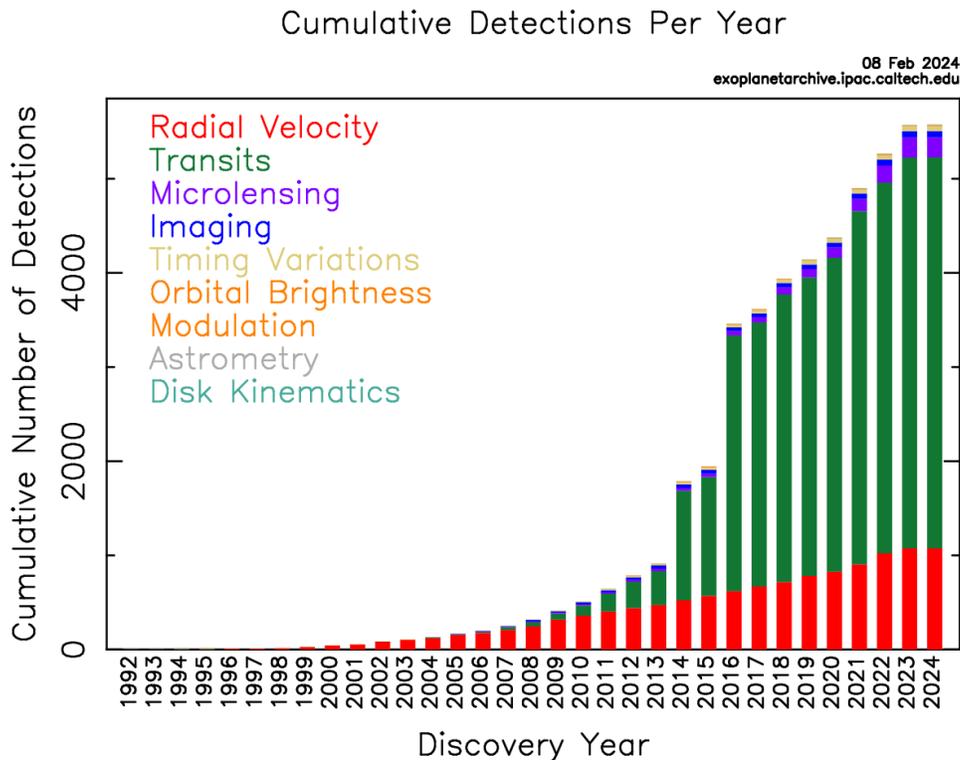


Figure 8.2: Cumulative detections of exoplanets since 1992, colored by detection method. Figure from the NASA Exoplanet Archive.

with transits becoming the dominant method upon the end of the primary Kepler mission in 2013. There are clear jumps in the transit population between 2013-2014 and 2015-2016 – these correspond to Kepler data releases from their data processing and validation pipeline. Both microlensing and imaging have steadily found planets since 2004 and 2005, respectively, with other methods playing a relatively minor role. Note that here “timing

variations” corresponds both to pulsar timing and transit/eclipse timing, of which the latter has been more productive (with 45 planets found by transit/eclipse timing and 7 from pulsar timing).

These detections span a broad range of mass and period space, as previously shown in Figure 1.1. Figure 8.3 shows the population of planets as of a decade ago, also as a function of mass but plotted against the semi-major axis normalized to the snow line, where $a_{sl} = 2.7 \text{ au } M_*/M_\odot$. This normalizes detections as a function of stellar type, enabling more

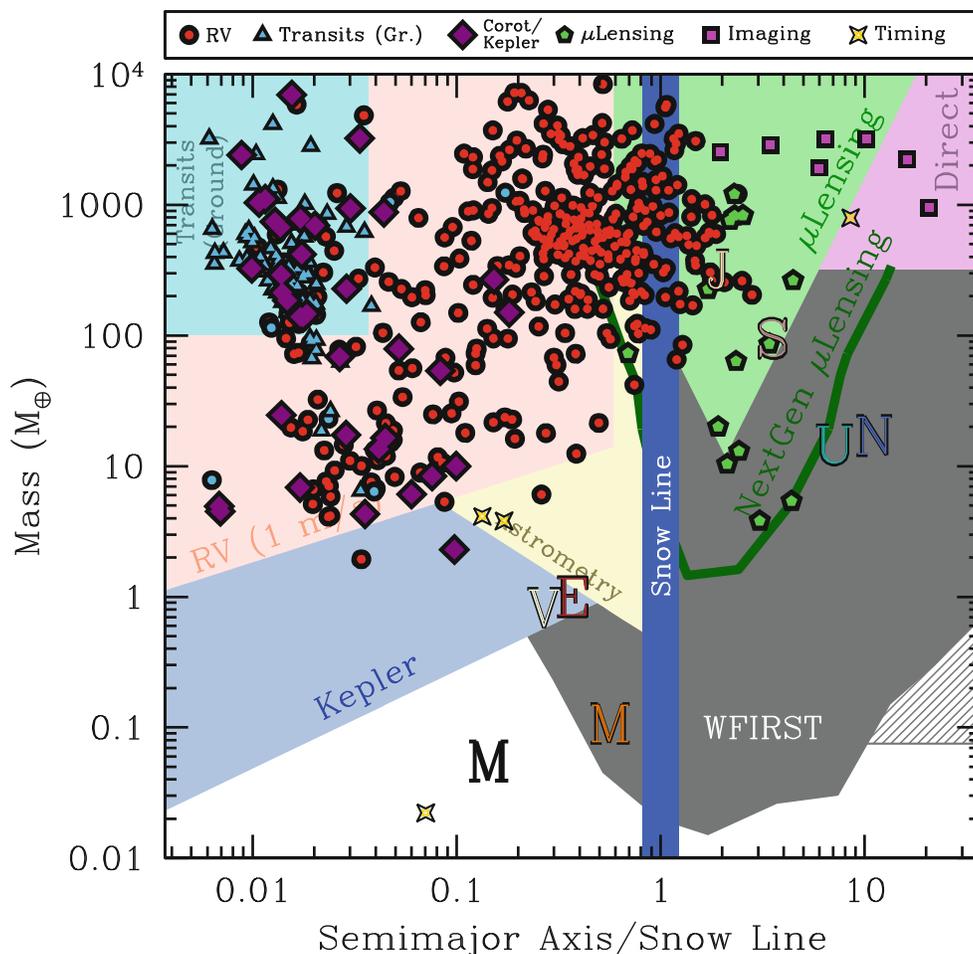


Figure 8.3: Detections of planets as of 2011 compared to sensitivity curves from various detection methods and missions. From Wright & Gaudi (2013).

direct inter-comparison of the entire population of exoplanets. The sensitivity of various detection methods is over-plotted by the shaded regions, displaying the biases of each detection method. For instance, transit is more favorable at small separations, RV is biased towards large masses, astrometric shifts are higher for wide orbits, microlensing preferentially finds planets with projected separations near the Einstein ring radius, and direct imaging can only currently find massive planets beyond the large inner working angles of ground-based coronagraphs. We will dig deeper into the sensitivity of each method in Section 8.2, but let’s first recap the highlights of each of the four most prolific detection methods (RV, transit,

microlensing, and imaging) to date.

8.3.1 Radial velocity

There are three key inferences from the radial velocity detections of hot Jupiters in the late 1990s and early 2000s that have stood the test of time. First, hot Jupiters were inferred to be intrinsically rare, with occurrence rates of $\sim 1\%$ or less – long-baseline RV surveys have proven this to be the case (see next paragraph). Second, gas giant planets are more common around massive stars (Lovis & Mayor, 2007), which agrees with the basic expectation that protoplanetary disk masses should scale with stellar mass, allowing more mass to be incorporated in planets. Finally, as shown in Figure 8.4, gas giants are more common around stars with a higher metallicity ($[\text{Fe}/\text{H}]$). This finding agrees with the basic

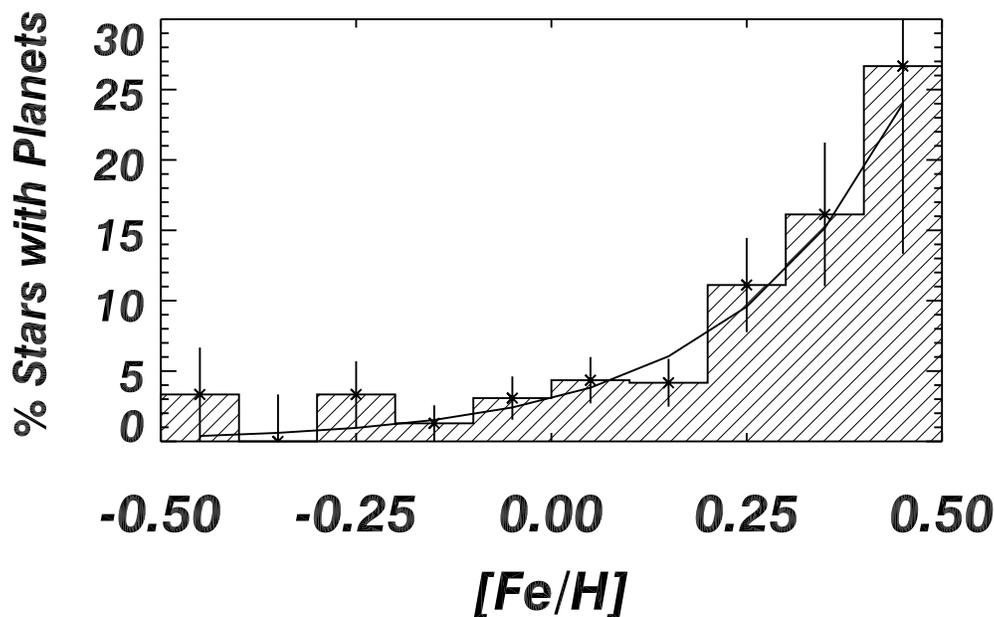


Figure 8.4: The planet-metallicity correlation: stars with a higher metallicity have a greater number of gas giant companions. Adapted from Fischer & Valenti (2005).

expectations of the core accretion hypothesis for giant planet formation, as stars with a greater metallicity should have protoplanetary disks with more metals (including dust and ices) that can be incorporated as the building blocks of the cores of giant planets. This will enable proto-giant planets in more metal-rich disks to build more massive cores, allowing them to be more likely to reach the critical mass to accrete the surrounding H/He gas and form a giant planet via the core accretion instability. We'll discuss core accretion in more detail when we cover planet formation in the following two weeks.

Radial velocity is no longer limited to finding planets with short orbital periods – instead, now it has found (massive) planets out to orbital periods of $\sim 10^5$ days. Figure 8.5 shows the occurrence rate derived from the ≈ 24 year California Legacy radial velocity survey of 719 stars. There is a paucity of hot and warm Jupiters at small semi-major axes $\lesssim 0.5$ au, but a significant increase in the number of gas giant planets near and beyond the snow line ($a \gtrsim 1$ au). This implies that Jupiter and Saturn-like gas giant planets are relatively

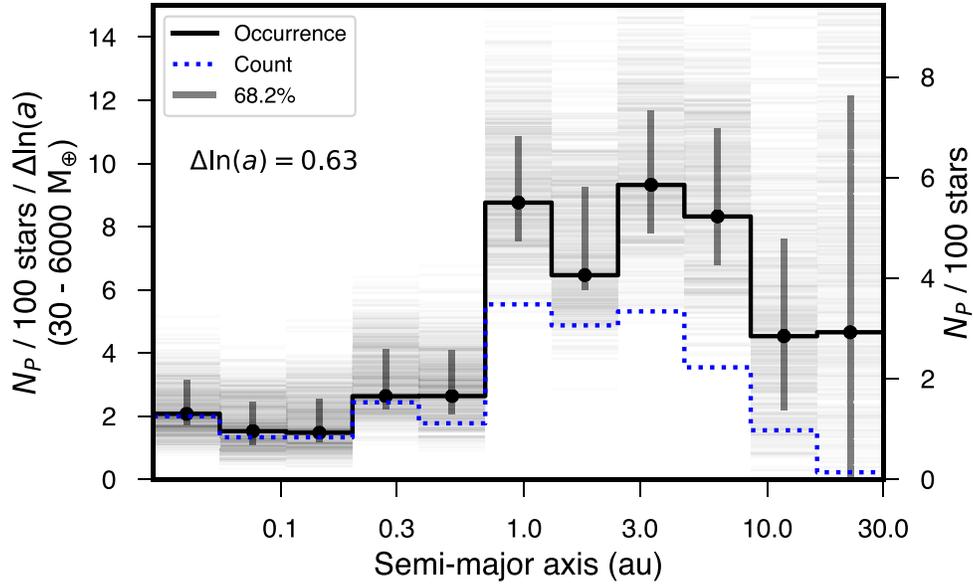


Figure 8.5: Occurrence rate of gas giant planets from the California Legacy RV Survey. There is a clear increase in the number of gas giants at wide separations, with hot Jupiters being relatively uncommon. Adapted from Fulton et al. (2021).

common.

8.3.2 Transit

Transits have provided a wealth of information on a range of planets, from hot gaseous planets orbiting main-sequence Sun-like stars down to temperate rocky planets orbiting M dwarfs. Figure 8.6 shows a summary plot of the transit radii of gas giant planets as a function of equilibrium temperature or, equivalently, incident stellar flux. The radii of warm gas giant planets are relatively independent of incident stellar flux, but those of hot Jupiters with equilibrium temperatures in excess of ≈ 990 K generally increase with equilibrium temperature (Laughlin et al., 2011). Additionally, the radii of many gas giants is larger than theoretical expectations, a “radius inflation” problem that is still unsolved (Fortney et al., 2021). We will discuss the mechanisms that set the radii of hot Jupiters further when we cover the internal structure of gas giants in the third part of this course.

Transit measurements, especially with Kepler, TESS, and targeted ground-based surveys, have provided a wealth of information about planets orbiting stars cooler than our Sun. This includes the detection of the seven-planet TRAPPIST-1 system (Gillon et al., 2017), along with a broad range of population statistics that we will cover in Day 10 when we discuss occurrence rates in more detail. One especially interesting finding of transit observations is that multi-planet systems appear to pack planets together with similar sizes and orbital spacing as their neighbors – like peas in a pod. Figure 8.7 shows this “peas in a pod” pattern for Kepler multi-planet systems around stars that are less massive than the Sun, though note that this trend continues to multi-planet systems around more massive host stars. The prevalence of closely packed planets with similar sizes in many systems may imply that planets migrate inward in the disks in which they form from their formation

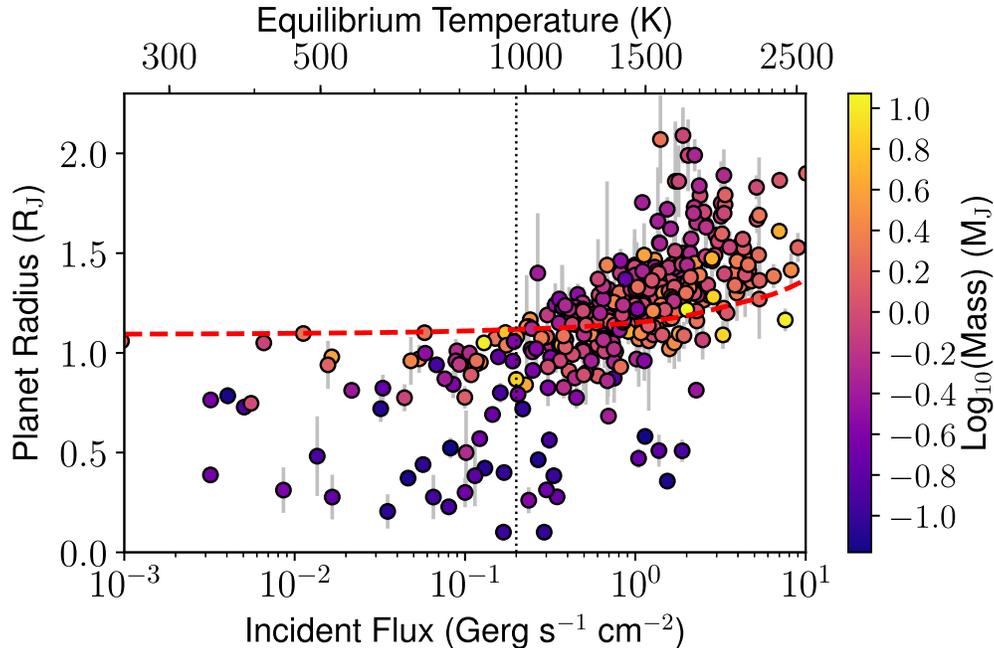


Figure 8.6: The radii of hot gas giant planets increase as a function of incident stellar flux. Additionally, the hottest gas giant planets (with equilibrium temperatures $\gtrsim 990$ K, see vertical dotted line) have radii that can be larger than those predicted from standard evolution models (red dashed line). Adapted from Thorngren & Fortney (2018).

locations. We’ll discuss this “Type 1 migration” within the protoplanetary disk (Type 2 migration is for gas giant planets that can open a gap in the disk) in the coming weeks.

8.3.3 Direct imaging

Direct imaging has been the only method to find massive, wide-separation planets and planetary-mass objects³ since the initial detection of 2M J1207b in 2005. The primary efforts in direct imaging in the past decade have been large (hundreds of stars) ground-based surveys, which have provided information on the statistics of giant planets at wide separations. One key detection is 51 Eri b (see Figure 8.8), which is the lowest-mass directly imaged planet ($2 M_{\text{Jup}}$). Nielsen et al. (2019) conducted an occurrence rate analysis on the GPI survey, finding that directly imaged planets that are still hot from formation are significantly more common around massive stars with $M > 1.5 M_{\odot}$, while the occurrence distribution around lower-mass stars is consistent with zero. There is also evidence for a decrease in the number of gas giant planets with increasing separation from $10 \text{ au} \lesssim a \lesssim 100 \text{ au}$, in tentative agreement with some radial velocity surveys.

8.3.4 Microlensing

Since the initial microlensing discovery of a $\approx 1.5 M_{\text{Jup}}$ planet at $\sim 3 \text{ au}$ (Bond et al., 2004), microlensing has found a wealth of massive planets at intermediate separations, along with approximately a dozen $0.5 - 2 M_{\oplus}$ planets. Figure 8.9 shows the light curve and caustic

³Planetary-mass objects are planetary mass ($M \lesssim 13 M_{\text{Jup}}$) but might have formed like a brown dwarf, via gravitational collapse, rather than bottom-up like a planet.

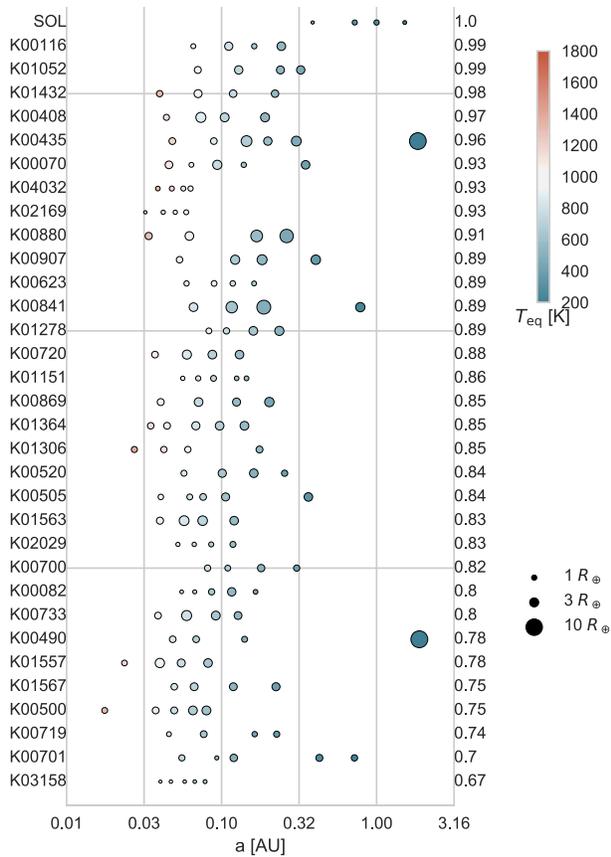


Figure 8.7: Kepler multi-planet systems show visible (and statistical) uniformity in radius of neighboring planets as well as period spacing between neighboring planets. This has been termed the “peas-in-a-pod” pattern of close-in multi-planet systems. Figure adapted from a larger figure in Weiss et al. (2018) that also shows systems with $M_\star > M_\odot$.

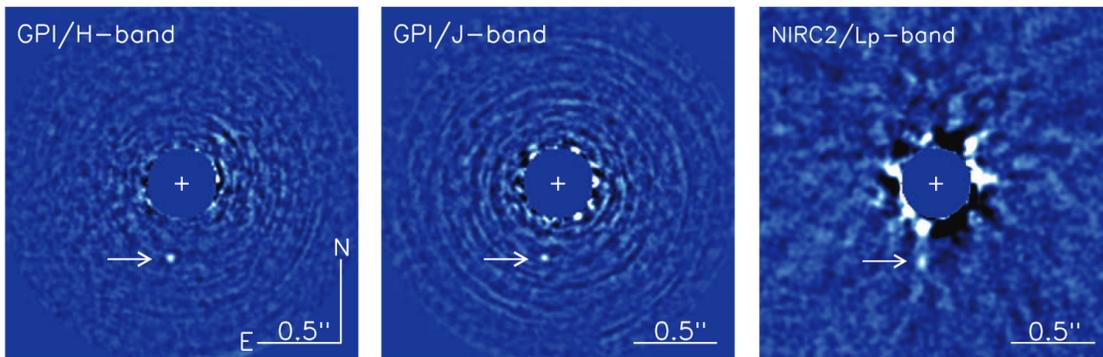


Figure 8.8: Detection images of 51 Eridiani b, the lowest-mass planet (2 Jupiter masses) found via direct imaging, discovered by the Gemini Planet Imager (GPI) Survey. Adapted from Macintosh et al. (2015).

structure of one of the most important early microlensing events, the discovery of a Jupiter-Saturn analogue pair orbiting a $\sim 0.5 M_\odot$ star at a distance of ~ 1.5 kpc with masses of $0.71 M_{\text{Jup}}$ and $0.27 M_{\text{Jup}}$ separations of 2.3 au and 4.6 au. This discovery clearly displayed the strength of microlensing, that it probes regions near where the gas giants in our own Solar System lie. Given that this region is also where the ice lines in protoplanetary disks are expected to be, there is significant hope for future space-based microlensing with Roman and beyond.

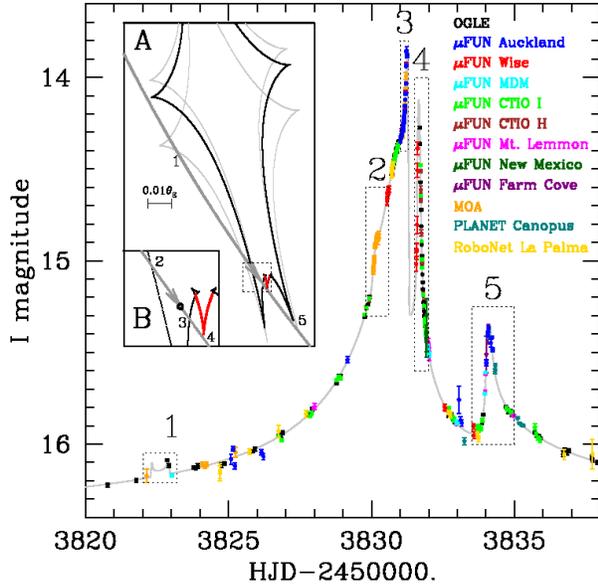


Figure 8.9: Microlensing discovery of two gas giants orbiting a $\sim 0.5 M_{\odot}$ star with masses of 0.71 and 0.27 Jupiter masses at separations of 2.3 and 4.6 au – a Jupiter-Saturn analog system. Adapted from Gaudi et al. (2008).

Microlensing has also found a population of ≈ 30 total free-floating, or “rogue” planets. Figure 8.10 shows one recent discovery, of a free-floating Neptune-mass planet. There has

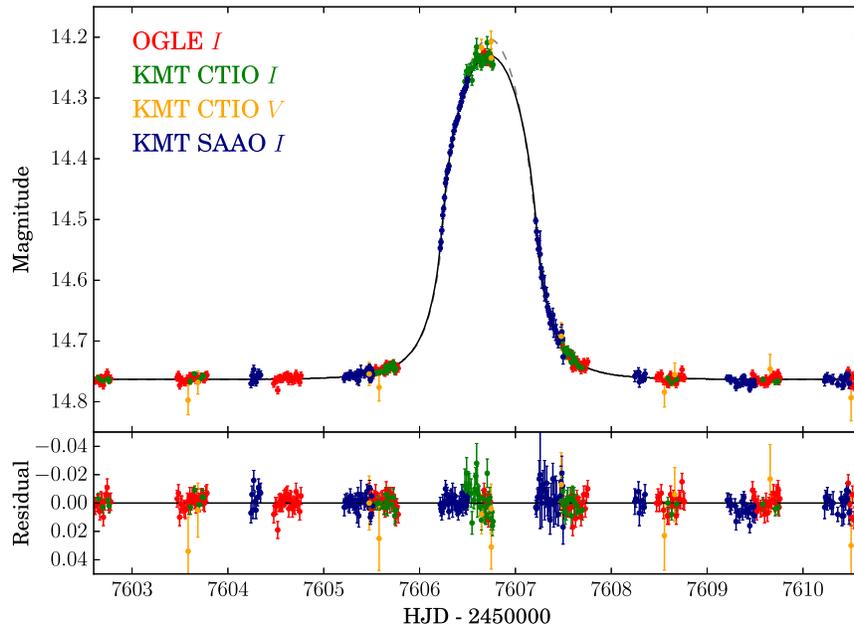


Figure 8.10: Discovery of a free-floating Neptune-mass planet, showcasing the short duration of single “rogue” planet lens events. Adapted from Mróz et al. (2018).

been a more recent discovery of a 41.5 minute microlensing event of a Mars-to-Earth mass rogue planet (Mróz et al., 2020), showcasing the technical capability of current microlensing surveys. Efforts are ongoing to study the statistics of free-floating planets, with some tentative evidence for a gap in the occurrence of microlensing planets with angular Einstein ring radii of 8.8 - 26 micro-arcseconds (Gould et al., 2022).

9 Detecting exoplanets: occurrence rates

Our agenda for Day 9 is the following:

1. Midterm recap (10 minutes)
2. Recap updates to exoplanet sample in the decade since Wright & Gaudi (10 minutes)
3. Deriving occurrence rates from biases: example of the transit method (30 minutes)
4. Highlights of exoplanet detections (remainder of time, mostly for fun)

Today's reading is Sections 2.10-2.12, 5.10, 6.25-6.26, and 7.10 of the second edition of our textbook. If you don't have the second edition, please read Fulton et al. (2017) (<https://ui.adsabs.harvard.edu/abs/2017AJ...154..109F/abstract>) instead, as we'll use that as an example of how to derive occurrence rates from transit surveys.

Our learning goals for today are:

1. Identify strengths and limitations of our current exoplanet sample.
2. Understand how occurrence rates are derived from a uniform but biased sample of planets.
3. Become aware of some key moments in exoplanet discovery from the past two decades, as well as important trends that have been discovered.

9.1 Occurrence rates

9.1.1 General principles

The observed distribution of planets (as shown in Figure 1.1) is strongly biased by the fact that each detection method is more sensitive to a given region of planetary mass/radius, separation, and stellar mass parameter space, among other parameters including age and distance. We previously discussed the general trends of these biases for each detection method in Section 8.2. These biases must be taken into account in order to back out the true underlying distribution of planets.

Occurrence rates have been derived individually for each method, some of which (e.g., for RV, see Figure 8.5) were discussed previously. In this section, we will focus only on occurrence rates via transit, but the general principles of backing out occurrence rates from the observed exoplanet distribution is the same for each detection method. There are three main steps to deriving occurrence rates from a survey with any given detection method:

1. **Sample Selection:** Cull the sample of observations in order to limit observational biases. These include magnitude cuts to only stars that are bright enough to detect planets around, and cuts in stellar properties and planet properties in order to limit outliers and/or systems that were included for reasons going beyond having a uniform survey.

2. **Survey sensitivity:** Perform tests on the data in order to quantify the fraction of planets with given properties (e.g., mass, radius, separation – all of which are connected to the signal-to-noise ratio of a given method) that are recovered. Use this to derive the detectability of the survey as a function of the planetary and/or stellar parameters of interest.
3. **Calculate occurrence:** Weight the actual detections in the survey by the detectability derived in Step (2) in order to calculate the true expected occurrence from the detections at hand. Analyze these occurrence rates as a function of the planetary and/or stellar properties of interest.

Generally, reliable occurrence rates are derived from samples that consist of hundreds to thousands of target stars. In the remainder of this section, we will study occurrence rates derived from the deepest transit survey done to date, from the Kepler primary mission.

9.1.2 Early results

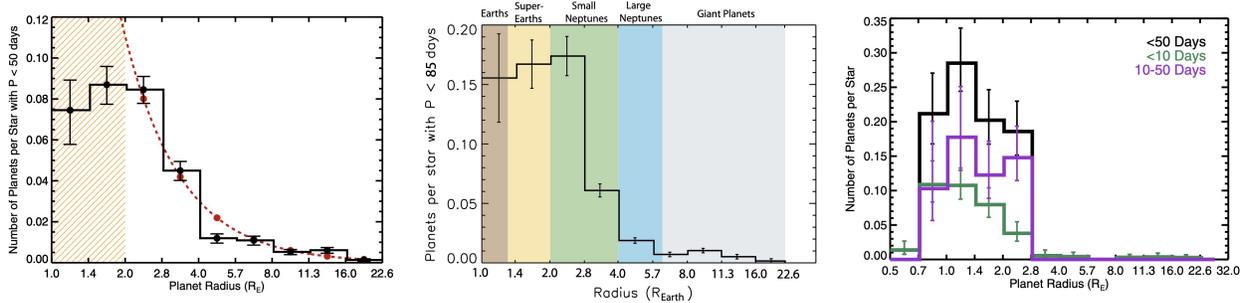


Figure 9.1: Comparison of Kepler occurrence rates derived from Howard et al. (2012) (left, FGK stars), Fressin et al. (2013) (middle, FGK stars), and Dressing & Charbonneau (2013) (right, for $M < M_{\odot}$). All results agree that planets with $R \lesssim 3 R_{\oplus}$ are the most common. Note that the samples are incomplete at the smallest radii (shown by the hatched region in the left-hand panel).

Figure 9.1 compares the occurrence rate of planets (as the number of planets per star) as a function of planet radius derived from Kepler observations. This compares the results for planets orbiting FGK stars to periods of 50 days (Howard et al., 2012), FGK stars to periods of 85 days (Fressin et al., 2013), and small stars with $M < M_{\odot}$ (Dressing & Charbonneau, 2013). The most striking finding is that planet occurrence is much larger for smaller planets relative to gas giant planets. This is not what you would expect from simply looking at the detections of planets (Figure 1.1) by eye, showcasing how by incorporating the sensitivity of a given survey one can reveal the true distribution of planets. Note that the Kepler data shown in Figure 9.1 is incomplete for small planets ($R \lesssim 1 R_{\oplus}$), so though there is a clear increase in planets at small radii it is unknown how this extends down to planets with masses similar to those of Mars or Mercury. Additionally, as we will discuss in the next section further characterization of the host stars in the Kepler survey has allowed more detailed constraints to be placed on the occurrence rate distribution as a function of planet radius.

9.1.3 Example of deriving occurrence rates: Radius gap

The key limiting factor in the dependence of occurrence rate on planet radius is the uncertainties in the host star radii, as the observable linked to planet radius (transit depth) $\propto R_{\star}^{-2}$. Fulton et al. (2017) improved the uncertainties on the stellar radius using the California-Kepler Survey sample of 2025 host star spectra, which allowed a reduction in the stellar uncertainty from typical values of $\approx 25\%$ to $\approx 11\%$. We'll next walk through the specific steps in the occurrence rate analysis of Fulton et al. (2017).

The first step for any occurrence rate analysis is to cull the sample to limit bias. Figure 9.2 shows the cuts made by Fulton et al. (2017) to their sample of host stars. The first cut is

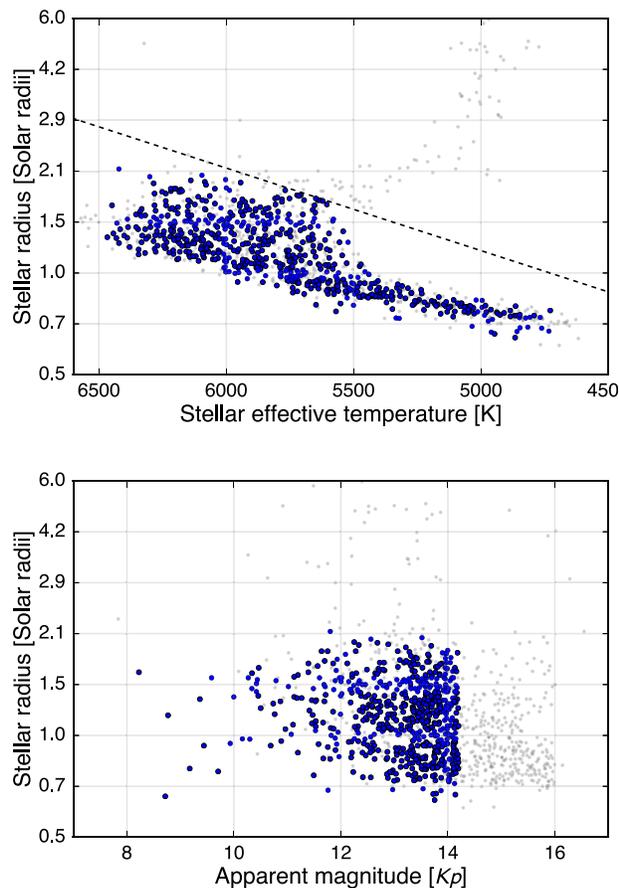


Figure 9.2: The sample of stars used by Fulton et al. (2017) in their analysis (blue) along with the stars removed due to either magnitude cuts or radius cuts at a given effective temperature (gray). The dashed line in the H-R diagram (top panel) shows the radius cut-off used as a function of effective temperature.

to add a filter on stellar radius as a function of temperature (the dashed line in the top panel of Figure 9.2) in order to remove giants from the survey. The second cut to the stellar distribution is to remove stars with Kepler magnitudes $K_p > 14.2$ in order to only include stars that can be well-characterized via follow-up. The third cut to the stellar distribution is to only include stars with $4700 \text{ K} < T_{\text{eff}} < 6500 \text{ K}$, as this is the range of stars where the spectroscopic stellar follow-up provides precise stellar parameters.

On top of these cuts to the stellar distribution, the authors also conduct cuts to the planet candidate distribution (and thus planet host star distribution) as well. The first is to remove signals that are deemed to be false positives from the sample. The second is to remove grazing transits with $b > 0.7$, as the properties of these planets become degenerate

with the limb darkening parameterization used, and are thus less precise. The last cut to the planet sample is to limit the orbital period to only $P < 100$ d in order to ensure a reasonable signal-to-noise for each candidate. Putting all these cuts, 3 for the stellar distribution and 3 for the planet distribution, together reduces the initial sample of 2025 stars in the spectroscopic Kepler follow-up (CKS) observations down to 900 stars that are analyzed to derive occurrence rates.

Step 2 in the occurrence rate calculation is to determine the sensitivity of the survey. Fulton et al. (2017) use results from a previous injection-recovery study of Kepler host stars in order to quantify the fraction of signals that are recovered as a function of their signal-to-noise, which is

$$(S/N)_i = \left(\frac{R_p^2}{R_{\star,i}^2} \right) \sqrt{\frac{T_i}{P}} \frac{1}{N_i}, \quad (9.1)$$

where R_p and P are the radius and period of a given injected planet, $R_{\star,i}$ is the stellar radius for a star in the Kepler catalog, T_i is the observation duration for that star, and N_i is the photometric noise for each star given the transit duration included in the injection. The resulting dependence of the recovery fraction of injected signals on signal-to-noise ratio is shown in Figure 9.3, along with a best-fit Γ distribution function (which the authors call C) to this histogram. The recovery fraction increases with SNR as expected, with injections

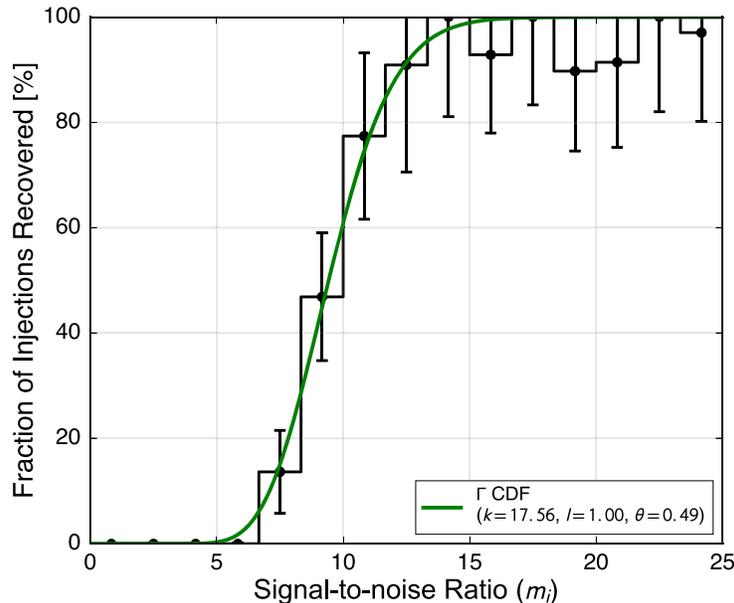


Figure 9.3: The results of injection-recovery tests for the fraction of injected planet signals recovered as a function of the signal-to-noise ratio of the planet signal used in Fulton et al. (2017). The green line is a Γ cumulative distribution function (CDF) fit to the recovery fraction dependence on SNR.

not recovered for $(S/N) < 5$ and injections generally recovered for $(S/N) > 15$.

The next step is to convert these recovery probabilities into a completeness fraction of the survey for a given planet radius and orbital period. To do so, the authors calculate the

fraction of stars where a transiting planet with a given SNR would be detected

$$p_{\text{det}} = \frac{1}{N_{\star}} \sum_i^{N_{\star}} C . \quad (9.2)$$

This quantity p_{det} is also called the “completeness” of the pipeline, and is shown in the top panel of Figure 9.4. As expected, the completeness is one for large planets at short orbital

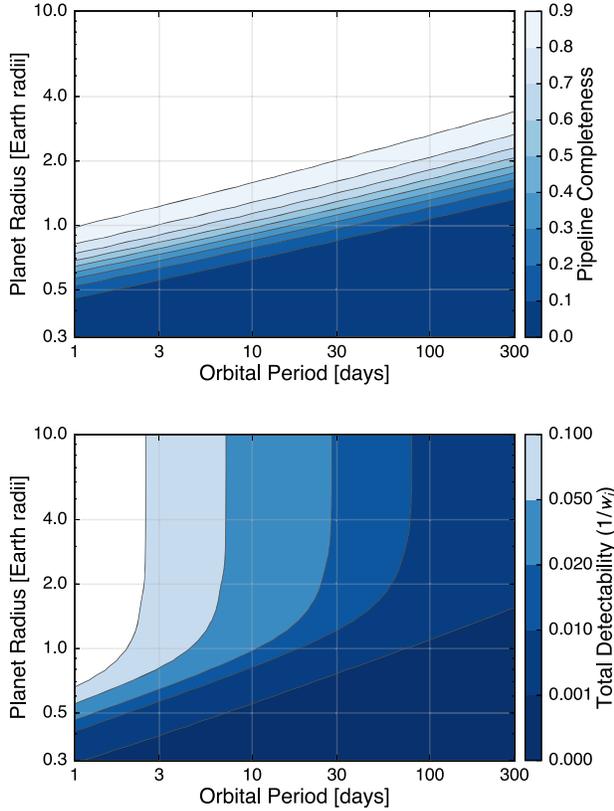


Figure 9.4: Top: The completeness (p_{det}) of the analysis pipeline of Fulton et al. (2017) as a function of radius and period. Bottom: the detectability ($1/w_i$) as a function of radius and period, folding in both the completeness and transit probability.

periods, and drops as the SNR decreases with decreasing radius and increasing orbital period. In order to then calculate the total detectability as a function of planet radius and orbital period, the authors also need to fold in the transit probability

$$p_{\text{tr}} = 0.7 \frac{R_{\star}}{a} , \quad (9.3)$$

where the extra factor of 0.7 comes from the $b < 0.7$ cut used in the sample selection. The authors combine the transit probability and the completeness as a weighting function

$$w_i = \frac{1}{p_{\text{det}} p_{\text{tr}}} , \quad (9.4)$$

which is applied to each planet detection. This weighting function is the “total detectability” plotted in the bottom panel of Figure 9.4, which differs from the completeness due to the lower transit probabilities at longer orbital periods.

Finally, to calculate the occurrence rate in terms of the number of planets per star in a bin of a given planet radius and orbital period, the authors simply take the sum of the weights divided by the number of stars in the sample:

$$f_{\text{bin}} = \frac{1}{N_{\star}} \sum_{i=1}^{n_{\text{pl,bin}}} w_i . \quad (9.5)$$

Figure 9.5 shows a histogram of the resulting occurrence rate distribution of planetary radii. The occurrence rate distribution in Figure 9.5 is similar to that shown in Figure 9.1 in that

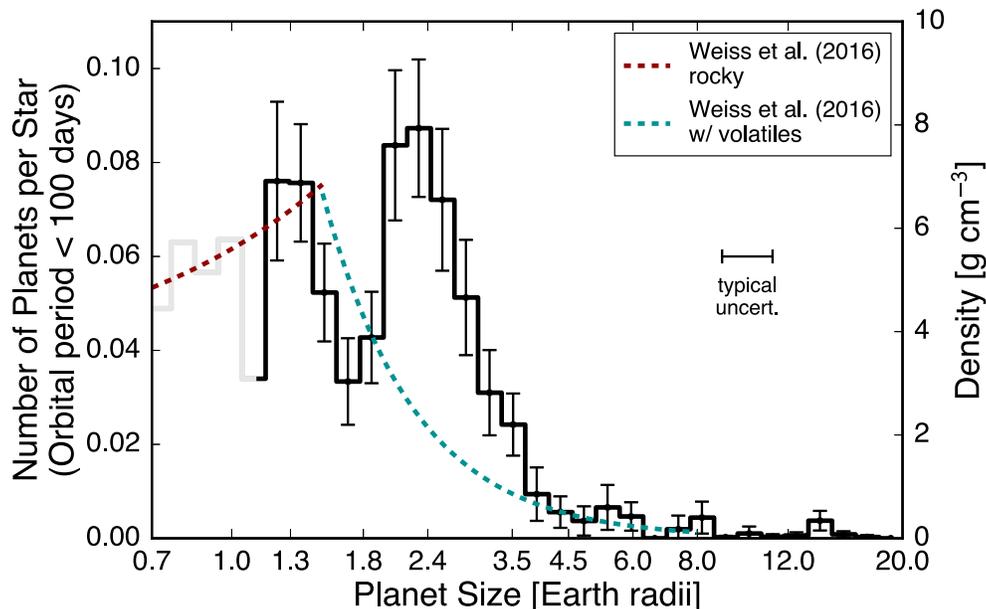


Figure 9.5: The resulting occurrence rates from Fulton et al. (2017) as a function of planet size. Over-plotted on this occurrence rate distribution is the density of planets as a function of radius from Weiss et al. (2017), which shows a peak near the super-Earth occurrence rate peak – implying that the occurrence rate gap can be interpreted as a transition from rocky to gas-rich planets.

the number of planets with $R \lesssim 4 R_{\oplus}$ is much greater than the number of larger gaseous planets. However, with the improved stellar parameters and corresponding reduction in planetary radii, there is a clear bi-modality, with two peaks at radii of $1 - 1.5 R_{\oplus}$ and $2 - 3 R_{\oplus}$, along with a gap in occurrence rate between these peaks. Over-plotted on the histogram of occurrence is the inferred density distribution from the sample of RV and TTV mass and transit radius measurements (Weiss et al., 2017), which shows a peak in density that aligns well with the smaller-radius peak in occurrence. This implies that planets that are above this smaller-radius peak in occurrence have a significant volatile component, while smaller planets are largely rocky. As a result, this divide can be interpreted as a divide between rocky and gaseous planets, with “super-Earths” at radii smaller than the radius gap (but larger than Earth), and “sub-Neptunes” at radii above the gap. Specifically, Fulton et al. (2017) define super-Earths as having radii $1 R_{\oplus} < R < 1.75 R_{\oplus}$ and sub-Neptunes as having radii of $1.75 R_{\oplus} < R < 3.5 R_{\oplus}$.

Given that the derived occurrence rates are a function of both radius and orbital period, the authors can then convert this distribution to show two-dimensional occurrence rate maps. Figure 9.6 shows this in radius-instellation space in order to link to theories for what mechanism could drive this gap in occurrence rate. The existence of a radius gap in the

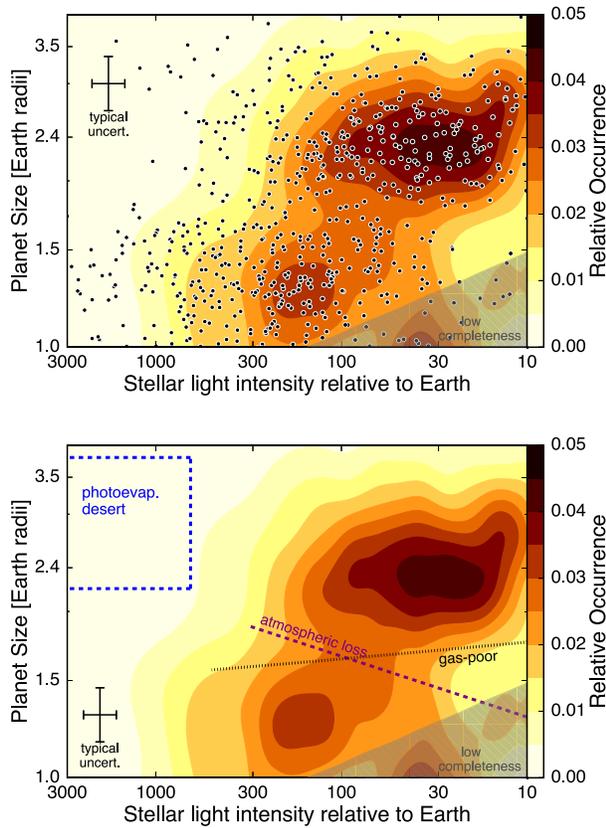


Figure 9.6: The two-dimensional occurrence from Fulton et al. (2017) as a function of planet radius and instellation. The bottom panel shows the scaling relation expected for atmospheric loss via photoevaporation, which is the best explanation for the gap in radii between super-Earths and sub-Neptunes.

Kepler population was actually predicted 4 years before its discovery by Owen & Wu (2013). The fundamental mechanism that is expected to cause this gap is atmospheric loss, where low-mass planets more readily lose their primordial H/He envelopes. Sub-Neptunes above the gap can hold onto their H/He, while super-Earths lose any primary H/He atmosphere. Importantly, the gap appears to move to smaller radii at lower instellations, which agrees with the basic expectation that the mechanism driving this loss is linked to the irradiation – planets that receive more irradiation receive more high-energy stellar photons that can drive atmospheric loss, and planets that are hotter will have more extended atmospheres that can more easily be lost to space. At present, there are two competing mechanisms for the atmospheric loss that carves the radius gap – photoevaporation (loss due to high-energy stellar photons driving atmospheric escape) and core-powered mass loss (outflows driven by the cooling of the interior from formation). A key way to differentiate between these is that core-powered mass loss should not be dependent on the host star type, while photoevaporation is dependent on the host star type through the stellar spectrum. There is tentative evidence for a weak dependence of the radius gap on host star type (Berger et al., 2023), which may point toward core-powered mass loss as the dominant mechanism – regardless, both mechanisms certainly play a role in the evolution of young, low-mass planets with H/He envelopes.

10 Planet formation: disk structure

Our agenda for Day 10 is the following:

1. Star formation recap (5 minutes)
2. Vertical disk structure from hydrostatic equilibrium (40 minutes). As part of this, do small group derivations to get from hydrostatic balance to disk density profile.
3. Disk flaring (10 minutes)
4. Activity: Estimating disk temperatures (20 minutes)
5. (if time) Start disk thermal structure

We'll start reading the Armitage lecture notes for today, and will continue for the next two weeks. Today's reading is Ch. II A-B of the lecture notes, which covers protoplanetary disk structure. The reading for next class is Ch. II C-E. Then, next week we'll cover Ch. III A (Tuesday) and Ch. III B-C (Thursday), and the following week we'll finish with Ch. IV.

10.1 Vertical disk structure

10.1.1 Hydrostatic equilibrium

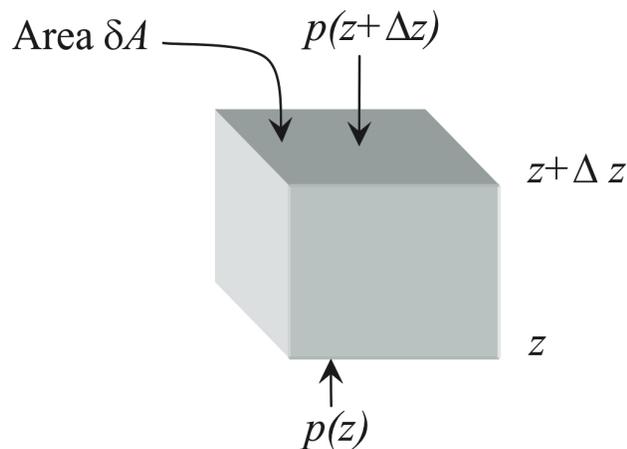


Figure 10.1: Schematic of a parcel of gas in hydrostatic equilibrium.

Figure 10.1 shows a parcel of gas held between z and $z + \Delta$ at pressure $p(z)$ at the bottom and $p(z + \Delta z)$ at the top. Newton's second law ($F = ma$) implies that the change in pressure across Δz is the weight per area of the gas:

$$[p(z + \Delta z) - p(z)]\delta A = -\rho g \Delta z \delta A . \quad (10.1)$$

Thus,

$$\frac{p(z + \Delta z) - p(z)}{\Delta z} = -\rho g , \quad (10.2)$$

and given that

$$\lim_{\Delta z \rightarrow 0} \frac{p(z + \Delta z) - p(z)}{\Delta z} = \frac{dp}{dz} , \quad (10.3)$$

we can write the expression for hydrostatic balance

$$\frac{dp}{dz} = -\rho g . \quad (10.4)$$

Let's now apply this to calculate the vertical density structure of a protoplanetary disk. Figure 10.2 shows the direction of the vertical component of gravity in a geometrically thin

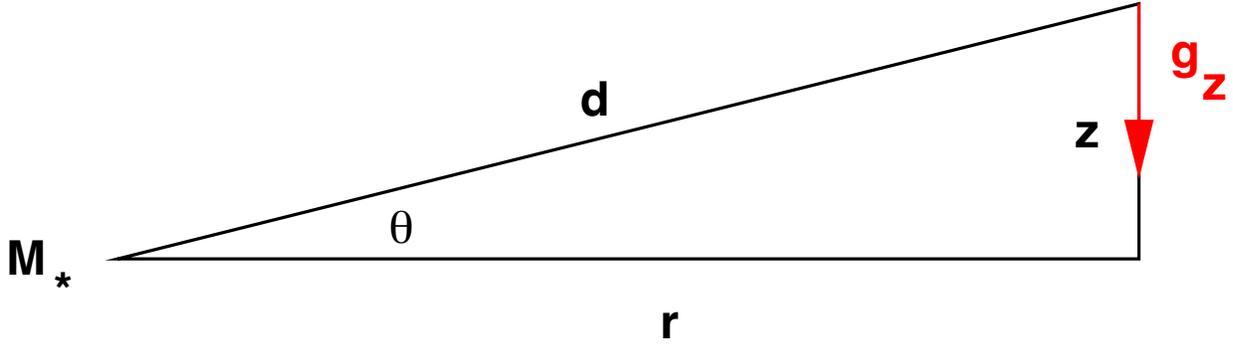


Figure 10.2: Schematic we'll use to calculate the vertical component of gravity g_z in a protoplanetary disk from hydrostatic equilibrium. Adapted from Armitage (2007).

disk. In this case, hydrostatic equilibrium in the vertical direction of the disk is

$$\frac{dp}{dz} = \rho g_z , \quad (10.5)$$

where

$$g_z = g \sin\theta , \quad (10.6)$$

with

$$g = \frac{GM_\star}{d^2} = \frac{GM_\star}{(r^2 + z^2)} . \quad (10.7)$$

We can thus write the vertical component of gravity g_z as

$$g_z = \frac{GM_\star}{(r^2 + z^2)} \frac{z}{\sqrt{r^2 + z^2}} = \frac{GM_\star}{d^3} z . \quad (10.8)$$

We can now use the ideal gas law to relate pressure and density

$$P = \frac{\rho k_B T}{\mu m_p} = \rho c_s^2 , \quad (10.9)$$

where

$$c_s^2 = \frac{k_B T}{\mu m_p} \quad (10.10)$$

is the isothermal sound speed, $\mu \approx 2.3$ is the mean molecular weight and m_p is the proton mass.

10.1.2 Activity: Derive disk density profiles in small groups!

We can now derive the density profiles for disks in hydrostatic equilibrium. Follow these steps, and check in with me after each one.

1. Substitute the ideal gas law and g_z into your expression for vertical hydrostatic equilibrium (Equation 10.5) to express hydrostatic equilibrium as a function of density ρ instead of pressure p . Assume $z \ll r$ to relate $d\rho/dz$ to the sound speed c_s , height z , density ρ , and Keplerian orbital frequency $\Omega = \sqrt{GM_\star/r^3}$.
2. Integrate your expression of vertical hydrostatic equilibrium from $z' = 0$ to $z' = z$ to derive the dependence of density on height, assuming that the disk is isothermal (and thus c_s is constant) and given a mid-plane density at $z = 0$ of ρ_0 . Re-write your expression as a function of the disk density scale height $h = c_s/\Omega$.
3. Integrate your expression for the vertical density profile of the disk from $z' = -\infty$ to $z' = +\infty$ to derive the total surface density of the disk (i.e., the integrated mass in a 2D column of the disk), which we call Σ . Use this to relate the mid-plane density ρ_0 to the surface density Σ and scale height h .

We'll have groups come up and derive each part of the solution for the class, and I'll post the full solutions below (they're also just commented out on the Overleaf if you want to check).

10.1.3 Disk flaring

The shape of the disk depends on the aspect ratio, h/r . This can be related to the sound speed and orbital velocity as

$$\frac{h}{r} = \frac{c_s}{\Omega r} = \frac{c_s}{v_K} = \text{Ma}^{-1}, \quad (10.11)$$

where $v_K = \sqrt{GM_\star/r}$ is the Keplerian orbital velocity and Ma is the Mach number of the disk at the orbital velocity. If we then assume that the sound speed scales as $c_s \propto r^{-\beta}$, then

$$\frac{h}{r} \propto r^{-\beta} r^{1/2} \propto r^{1/2-\beta}. \quad (10.12)$$

Thus, the aspect ratio h/r will increase (i.e., the disk will “flare”) with r if $\beta < 1/2$. Given that $c_s \propto \sqrt{T}$, then the disk will flare if the temperature dependence with separation is $T \propto r^{-1}$ or shallower. As we will see in our activity, this is expected to be the case, and so protoplanetary disks are expected to flare.

10.2 Activity: estimating disk temperatures

Protoplanetary disk temperature (and surface density) profiles are often parameterized as a power-law with separation from the host star r . This activity will walk us through one limiting case for the dependence of disk temperature on r .

1. Split into 5 groups of 3 people each. Then, calculate the equilibrium temperature of a dust grain that lies at a given separation r from a Sun-like host star with a radius R_\star and a temperature T_\star . Assume that the dust grain is a sphere with zero albedo and the same temperature across its surface. Group 1: calculate this value T_{eq} at 0.01 au. Group 2: calculate this at 0.1 au. Group 3: calculate this at 1 au. Group 4: calculate this at 10 au. Group 5: calculate this at 100 au.

2. Have one group member post your answer on the log-log plot on the board of dust grain temperature as a function of separation.
3. Derive a scaling for the dependence of dust grain temperature on semi-major axis. Compare this to that found for a flared disk thermal profile in Equation (51) of the Armitage reading.
4. Roughly estimate the condensation temperatures for various species in a disk, first for a volatile like water (H_2O) and then for refractory species like perovskite (CaTiO_3) and/or silicates like enstatite (e.g., MgSiO_3). From this, estimate where in the disk these species would be found in solid vs. gaseous form.

11 Planet formation: disk thermal structure, dynamics

Our agenda for Day 11 is the following:

1. Flared disk structure (15 minutes)
2. Activity: Estimating flared disk temperatures (15 minutes)
3. Disk dynamics, effective viscosity (30 minutes)
4. Viscosity activity (15 minutes, if time, if not do next class)

Today’s reading is Ch. II C-E of the Armitage notes, which covers condensation and ice lines, dynamics, and the effective viscosity and angular momentum transport within disks.

11.1 Disk thermal structure

11.1.1 Flared disks

Protoplanetary disks are “flared” (have an aspect ratio that increases with separation from the host star) due to instellation puffing up the outer regions of the disk. Figure 11.1 shows a schematic of a flared disk.

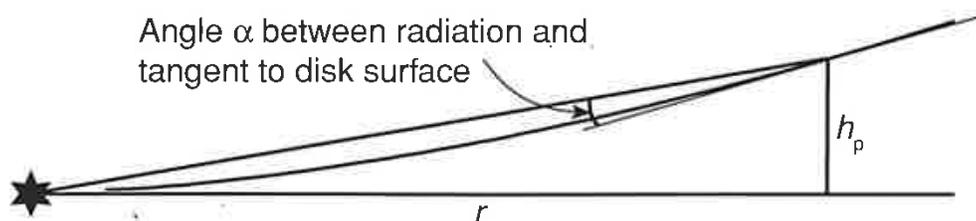


Fig. 2.4. Geometry for calculation of the radial temperature profile of a flared protoplanetary disk. At distance $r \gg R_*$, radiation from the star is absorbed by the disk at height h_p above the mid-plane. The angle between the tangent to the disk surface and the radiation is α .

Figure 11.1: Schematic showing a flared disk along with the flaring angle α . From Armitage (2013). I’ve left his caption in on purpose.

The flaring angle of a disk,

$$\alpha = \frac{dH}{dr} - \frac{H}{r}, \quad (11.1)$$

is the angle between the path of incident stellar radiation and the tangent to the disk surface at a given radius⁴, where H is the height of the disk from the mid-plane to the disk surface

⁴ α can be derived by drawing two right triangles at a given location, one that has a hypotenuse which follows the tangent to the disk surface and one that follows the path of incident stellar radiation. The opening angle of the smaller triangle is dH/dr , while the opening angle of the larger triangle is H/r , where H is the height of the disk at a given location. Thus, the angle between h and the hypotenuse in the smaller triangle is $\pi/2 - dH/dr$, and that for the larger triangle is $\pi/2 - H/r$. Subtracting the smaller angle from the larger angle, we find $\alpha = -H/r + dH/dr$, equivalent to Equation (11.1).

that intercepts starlight (equivalent to h_p in the Armitage (2013) notation, but in this part of the notes I use lowercase h to mean disk scale height). Assuming that the disk is locally in radiative equilibrium with incoming starlight onto an area A and that the flaring angle α is small,

$$\begin{aligned}\frac{L_\star}{4\pi r^2} A \sin\alpha &= A\sigma T_d^4 \\ \frac{L_\star}{4\pi r^2} \alpha &= \sigma T_d^4,\end{aligned}\tag{11.2}$$

where T_d is the temperature of the disk. Rearranging, we find

$$T_d = \left(\frac{L_\star \alpha}{4\pi \sigma r^2} \right)^{1/4}.\tag{11.3}$$

Using the Stefan-Boltzmann law to substitute $L_\star = 4\pi R_\star^2 \sigma T_\star^4$, we can write

$$T_d = \sqrt{\frac{R_\star}{r}} \alpha^{1/4} T_\star.\tag{11.4}$$

Thus, the dependence of the disk temperature is very similar in r, R_\star, T_\star to the equilibrium temperature we previously derived and (meant to) apply in last class' activity, with $T_d \propto r^{-1/2}$. However, there is an additional dependence on α , as disks that are more flared have a greater surface area to intercept incident starlight.

11.2 Disk temperature activity: condensation points and ice lines

This activity will enable us to estimate the effect of flaring on disk temperature and the separations where various species can condense out of the gas and form planetary building blocks.

1. Split into 5 groups of 3 people each. Each group corresponds to a separation r – Group 1 is 0.01 au, Group 2 is 0.1 au, Group 3 is 1 au, Group 4 is 10 au, and Group 5 is 100 au. Calculate the temperature of the disk at this semi-major axis around a Sun-like star with the effect of flaring. Assume $\alpha = 0.05$.
2. Have one group member post your answers on the log-log plot on the board of dust grain temperature as a function of separation, both with and without the effects of flaring.
3. Assuming a flared disk model with $\alpha = 0.05$, calculate the separation from the host star at which rocky material that forms meteorites and the building blocks of rocky planets can condense. Assume that the condensation temperature of rocky (silicate) material is 1500 K in typical disk conditions.
4. Assuming a flared disk model with $\alpha = 0.05$, calculate the separation from the host star at which water vapor can condense (the “ice line”). Assume that the condensation temperature of water is 170 K in typical disk conditions.

5. All of the derivations and calculations above are for the irradiated *surface* of the disk. However, as you'll find in problem set 2, the disk *midplane* is where we expect most of the dust to pile up. Do you expect the temperature in the disk interior to be hotter or colder than the surface, and why? How does this impact what materials can condense at the midplane?

11.3 Disk dynamics

11.4 Momentum balance

The radial force balance for the gas in a protoplanetary disk is largely between three forces: the centrifugal force, pressure gradients, and gravity. We can express this as

$$\frac{v_{\phi,g}^2}{r} = \frac{GM_\star}{r^2} + \rho^{-1} \frac{dp}{dr}, \quad (11.5)$$

where $v_{\phi,g}$ is the orbital velocity of the gas. Importantly, $dp/dr < 0$, which implies that the velocity of the gas will be sub-Keplerian. Assuming that the pressure in the disk follows a power-law as

$$p = p_0 \left(\frac{r}{r_0} \right)^{-n}, \quad (11.6)$$

where $p_0 = \rho_0 c_s^2$, we can substitute this into the force balance to find

$$\frac{v_{\phi,g}^2}{r} = \frac{v_K^2}{r} - \rho^{-1} \rho_0 c_s^2 n \frac{r^{-n-1}}{r_0^{-n}}. \quad (11.7)$$

Multiplying through by r and assuming that r and ρ are taken near the fiducial radius and density r_0 and ρ_0 , we find

$$v_{\phi,g}^2 = v_K^2 - n c_s^2, \quad (11.8)$$

and rearranging

$$v_{\phi,g} = v_K \sqrt{1 - n \frac{c_s^2}{v_K^2}}. \quad (11.9)$$

Note that $c_s = h\Omega$, so $n c_s^2 / v_K^2 = n h^2 \Omega^2 / v_K^2 = n h^2 / r^2$, and as a result

$$v_{\phi,g} = v_K \sqrt{1 - n \frac{h^2}{r^2}}. \quad (11.10)$$

For a typical disk with $h/r = 0.05$ and $n = 3$, $v_{\phi,g} = 0.996 v_K$. Thus, the effect of pressure gradients on the gas motion is relatively small and can be neglected for many of our purposes. However, we will come back to this next week, as it will be critical for determining the motions of intermediate-sized dust grains in the disk, as they move at Keplerian speeds while the gas is slightly sub-Keplerian.

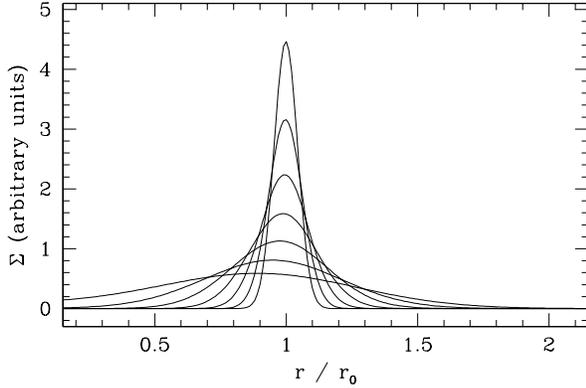


Figure 11.2: Pringle solution for the diffusive spreading of a disk with a constant viscosity outward from a thin ring. Adapted from Armitage (2007).

11.4.1 Effect of viscosity

Viscosity is a measure of the resistance of a fluid that is being deformed by stresses. We have ignored viscosity thus far, but it is of critical importance for the evolution of disks. This is because viscosity leads to angular momentum transport that redistributes the angular momentum in the disk (as some parcels of gas fall into the star, and some move outward) and contributes to disk evolution and eventual dissipation (see Figure 11.2). Importantly, this dissipation is measured to take $\lesssim 10$ Myr, with a mean disk lifetime of 3 – 5 Myr.

Molecular collisions generate a viscosity in a flow due to the finite mean free path of the gas. We can estimate the molecular viscosity of a fluid as

$$\nu_m = \lambda c_s , \quad (11.11)$$

where

$$\lambda = \frac{1}{n\sigma_{\text{mol}}} , \quad (11.12)$$

is the mean free path with n the number density and σ_{mol} the cross-section for molecular collisions. Given that we can re-write the number density $n = \rho/(\mu m_p)$ and have previously shown that the typical density $\rho_0 = \Sigma/(\sqrt{2\pi}h)$, we can write the mean free path as

$$\lambda = \frac{\mu m_p}{\rho\sigma} \approx \frac{\sqrt{2\pi}\mu m_p h}{\Sigma\sigma_{\text{mol}}} . \quad (11.13)$$

We can then write the molecular viscosity as

$$\nu_m \approx \frac{\sqrt{2\pi}\mu m_p}{\Sigma\sigma_{\text{mol}}} h c_s , \quad (11.14)$$

11.4.2 Viscosity activity

This activity will help us understand whether molecular viscosity can serve as the mechanism for angular momentum transport in protoplanetary disks. First, note that the typical timescale for viscous transport can be written as

$$\tau_\nu = \frac{r^2}{\nu} . \quad (11.15)$$

1. Calculate the viscous timescale for a disk from Equation (11.14) at 10 au, assuming it is comprised of a molecular solar composition gas ($\mu \approx 2.3$). Use the Hayashi (1981) prescription for surface density of an MMSN,

$$\Sigma = 1.7 \times 10^3 \text{ g cm}^{-2} \left(\frac{r}{1 \text{ au}} \right)^{-3/2}. \quad (11.16)$$

Further assume that the sound speed at 10 au is $c_s \approx 0.5 \text{ km s}^{-1}$, and that the collision cross section is $\sigma_{\text{mol}} \approx 2 \times 10^{-15} \text{ cm}^2$. Compare the timescale you calculate to the typical timescale for disk evolution of a few Myr.

2. Clearly we need a higher level of viscosity to cause disks to evolve on reasonable timescales. Calculate the required minimum ν to lead to disk evolution at 10 au, assuming a viscous evolution timescale of 10 Myr.
3. Let's assume that our expression for ν_m was incorrect. Instead, let's re-write it as

$$\nu = \alpha h c_s, \quad (11.17)$$

where α is a free parameter that replaces the combination of quantities $\left(\frac{\sqrt{2\pi\mu m_p}}{\Sigma\sigma_{\text{mol}}} \right)$ in Equation (11.14). Calculate the required α to reach a viscous evolution timescale of 10 Myr at 10 au using your solution from part (2).

11.4.3 Shakura-Sunyaev disks

It's clear from the activity that molecular viscosity cannot lead to disk evolution. One important consequence of a small molecular viscosity is a large Reynolds number,

$$\text{Re} = \frac{UL}{\nu}, \quad (11.18)$$

where U is a typical velocity scale and L is a typical length scale. If we take $U \approx c_s$ and $L \approx h$ and $h/r = 0.05$, we find that the typical Reynolds number at 10 au is $\sim 10^{10}$. This is a very large number, which implies that the flow will be very turbulent, with molecular viscosity largely irrelevant to the bulk of the fluid motions.

One possibility for the source of viscosity in disks is that it is caused in itself by the turbulence of the disk, which leads to mixing of neighboring fluid elements that acts as an effective viscosity. We can estimate that this effective viscosity will be limited to velocities smaller than the sound speed c_s (higher velocities lead to shocks and resulting dissipation), and scales less than the disk scale height h (which is generally the smallest physical scale in the disk). Thus, we can write the turbulent viscosity as

$$\nu = \alpha c_s h, \quad (11.19)$$

where α is a dimensionless value, known as the Shakura-Sunyaev α parameter (note that this is not the disk flaring angle⁵). Note that this expression is equivalent to the modified version of Equation (11.14) we used in our activity. The value of α is a priori unknown, with typical values ranging over orders of magnitude from $10^{-6} - 10^{-1}$ depending on the source of turbulence in the disk.

⁵Sorry, I'm sticking with traditional notation, where they are both α . So it goes.

12 Planet formation: dust and pebble motions

Our agenda for Day 12 is the following:

1. Recap viscosity, viscosity activity (25 minutes, see Day 11 notes)
2. Dust drag regimes, coupling (15 minutes)
3. Radial drift intro (10 minutes)
4. Radial drift derivation activity (remainder of class)

Today's reading is Ch. III A of the Armitage notes, which covers planet formation, starting with the interaction of dust and gas in the protoplanetary disk.

12.1 Dust motions

12.1.1 Epstein and Stokes drag regimes

Dust particles moving in a disk feel an aerodynamic force from the gas that opposes the motion of the dust particles. This force is directly proportional to the cross-sectional area of the dust grain πs^2 as well as the relative velocity of the dust grain to the gas disk \mathbf{v} ,

$$\mathbf{F}_D = -\frac{1}{2}C_D\pi s^2\rho\mathbf{v}^2, \quad (12.1)$$

where ρ is the gas density, C_D is a drag coefficient, and boldface stands for a vector.

For small particles with a size less than the mean free path of gas molecules ($s \lesssim \lambda$), the coefficient for this “Epstein drag” regime is

$$C_D = \frac{8v_{\text{th}}}{3v}, \quad (12.2)$$

where the mean thermal velocity in the gas

$$v_{\text{th}} = \sqrt{\frac{8k_B T}{\pi\mu m_p}} = \sqrt{\frac{8}{\pi}}c_s. \quad (12.3)$$

Thus, the drag force in the Epstein regime is

$$\mathbf{F}_D = -\frac{4\pi}{3}\rho s^2 v_{\text{th}}\mathbf{v}. \quad (12.4)$$

Conversely, particles with a size much greater than the mean free path ($s \gtrsim \lambda$) lie in a “Stokes drag” regime. In this case, the drag coefficient C_D depends on the Reynolds number of the particle,

$$\text{Re} = \frac{2sv}{\nu_m}, \quad (12.5)$$

which expresses the ratio between the advection of the particle by fluid motions and diffusion by molecular viscosity. Practically, the Stokes drag coefficient can be written as a piecewise function of the Reynolds number

$$\begin{aligned} C_D &\approx 24\text{Re}^{-1}, \text{Re} < 1 \\ C_D &\approx 24\text{Re}^{-0.6}, 1 < \text{Re} < 800 \\ C_D &\approx 0.44, \text{Re} > 800. \end{aligned} \quad (12.6)$$

Thus, the Stokes drag can be determined for a given Re by inserting this piecewise formulation for the drag coefficient in Equation (12.1). Note that the Epstein and Stokes drag coefficients are equivalent for $s = 9\lambda/4$, which is the transition size between regimes (Epstein drag at smaller s , and Stokes drag at larger s).

12.1.2 Dust coupling and settling

Small dust particles are tightly coupled to the gas, while very large particles are not affected by gas drag. We can quantify how tightly coupled the gas and dust in the disk is by defining a friction time scale for a dust particle of mass m

$$t_{\text{fric}} = \frac{mv}{F_D}, \quad (12.7)$$

which in the Epstein drag regime for a given density $\rho_m = m/(4/3\pi s^3)$ corresponds to

$$t_{\text{fric}} = \frac{\rho_m s}{\rho v_{\text{th}}}. \quad (12.8)$$

Assuming $\rho = 10^{-9} \text{ g cm}^{-3}$, $\rho_m = 3 \text{ g cm}^{-3}$, $v_{\text{th}} = 10^3 \text{ m s}^{-1}$ (appropriate for $a = 1 \text{ au}$), and $s = 1 \text{ }\mu\text{m}$, we find $t_{\text{fric}} = 3 \text{ s}$, implying that micron-sized dust particles are very closely coupled to the gas.

In reality, though small dust particles are tightly coupled to the gas in their horizontal motions, they are not perfectly coupled to the gas in the vertical direction due to the vertical component of the stellar gravitational force. Figure 12.1 shows a calculation of vertical settling (top) and coagulation (bottom) of a particle as it sinks toward the midplane and grows by coagulating with other particles. Typical settling velocities for small micron-sized particles are $\sim 10^5 \text{ yr}$, but the settling velocity scales with particle size – thus, larger particles settle toward the midplane more quickly. We’ll derive this settling velocity and settling timescale dependence on particle size (and other properties) in problem set 2.

12.1.3 Radial drift: derivation activity

Dust in the disk will drift radially due to interactions with gas. Small particles ($s \lesssim 1 \text{ cm}$) are well-coupled to the gas, so they orbit the star at a velocity slightly smaller than the Keplerian velocity (recall Equation 11.10),

$$v_{\phi,g} = v_K \sqrt{1 - n \frac{h^2}{r^2}} = v_K \sqrt{1 - \eta}, \quad (12.9)$$

where $\eta = nc_s^2/v_K^2$, with n the power-law exponent for the radial pressure dependence. Because small particles orbit at this slower velocity, they will not be in radial force balance like the gas because the gas has pressure support (i.e., a pressure gradient term in the force balance) while the dust does not. This causes small dust grains to spiral in toward the host star at a radial terminal velocity. Conversely, large “rocks” with $s \gtrsim 1 \text{ m}$ feel gas drag because their orbits are Keplerian, while those of the gas are sub-Keplerian. This causes an effective headwind that removes angular momentum from the orbit of the rock, causing it to drift inward.

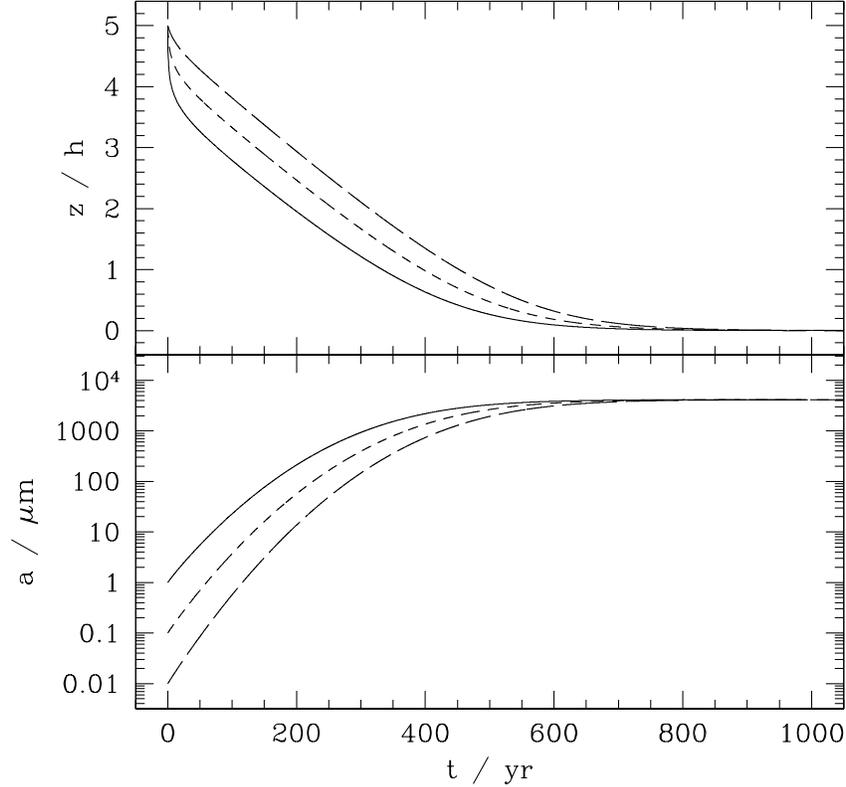


Figure 12.1: Calculation of vertical settling and resulting growth of particles in the protoplanetary disk. The solid line corresponds to $s = 1 \mu\text{m}$, the dashed line to $s = 0.1 \mu\text{m}$, and the long dashed line to $s = 0.01 \mu\text{m}$. Adapted from Armitage (2007).

We can formalize these two concepts by writing a specific force balance for particles (of any size, with a given stopping timescale t_{fric}) in the radial and azimuthal directions. In the radial direction,

$$\frac{dv_r}{dt} = \frac{v_\phi^2}{r} - \Omega_K^2 r - \frac{1}{t_{\text{fric}}}(v_r - v_{r,\text{gas}}), \quad (12.10)$$

where v_r is the radial velocity of the particle, the first term on the RHS is the centrifugal force, the second term on the RHS is gravity, and the third term on the RHS represents gas drag. The azimuthal force balance is only dependent on gas drag,

$$\frac{d(rv_\phi)}{dt} = -\frac{r}{t_{\text{fric}}}(v_\phi - v_{\phi,\text{gas}}), \quad (12.11)$$

where v_ϕ is the azimuthal velocity of the particle.

We can now use these statements of radial and azimuthal force balance for dust particles to derive the radial velocity (toward/away from the star, not to/from Earth!) of dust particles in the disk. Please do so in small groups of 2-3, following these steps:

1. First, simplify the azimuthal equation (Equation 13.3) by assuming that the particle

spirals in through a succession of nearly Keplerian orbits, i.e.,

$$\frac{d(rv_\phi)}{dt} \approx v_r \frac{d(rv_K)}{dr} = \frac{v_r v_K}{2}. \quad (12.12)$$

Plug this into the expression for azimuthal force balance to find an expression for $(v_\phi - v_{\phi,\text{gas}})$.

2. Simplify the radial equation (Equation 13.2) by substituting in $v_K^2 = v_{\phi,\text{gas}}^2 + \eta v_K^2$ and $\Omega_K^2 r = v_K^2/r$. There will be two terms on the right hand side of this equation relating to the azimuthal velocity, make the following first-order assumption that the gas and dust motions are similar

$$\frac{v_\phi^2}{r} - \frac{v_{\phi,\text{gas}}^2}{r} = \frac{(v_\phi + v_{\phi,\text{gas}})(v_\phi - v_{\phi,\text{gas}})}{r} \approx \frac{2v_K(v_\phi - v_{\phi,\text{gas}})}{r}, \quad (12.13)$$

to write the radial velocity equation for the dust to first-order accuracy.

3. Assume that there is no radial acceleration of the dust ($\frac{dv_r}{dt} \approx 0$) and use $(v_\phi - v_{\phi,\text{gas}})$ from the azimuthal force balance to derive the dependence of the radial velocity on $r, v_K, t_{\text{fric}}, v_{r,\text{gas}}, \eta$.
4. Re-cast the radial velocity by defining a dimensionless stopping time $\tau_{\text{fric}} \equiv t_{\text{fric}} \Omega_K = t_{\text{fric}} v_K/r$ to write v_r as a function of $\tau_{\text{fric}}, v_{r,\text{gas}}, v_K, \eta$.
5. The radial drift velocity peaks at $\tau_{\text{fric}} \approx 1$. What is the peak radial velocity, as a function of η, v_K ?
6. Assuming Epstein drag, what typical particle sizes does the peak radial velocity of a dust grain correspond to?

Solutions to your derived expression for particle radial velocity as a function of dimensionless stopping time are shown in Figure 12.2.

12.1.4 The “meter-size barrier”

The fast radial drift velocities of 10 cm - 1 m particles in the protoplanetary disk corresponds to an inward drift timescale of $\sim 10 - 1000$ years, increasing with radial separation from the host star. This implies that grains of this size are lost to radial drift over very short timescales – which necessitates mechanisms that rapidly concentrate meter-sized grains to form planetesimals. It also implies that radial drift will be common, and can potentially lead to build-up of material in the disk if pressure maxima occur, causing both radial inward drift from farther separations and outward drift interior to the pressure bump (see Figure 12.3). In practice, it is expected that dust rapidly coagulates into planetesimals through two-fluid (dust-gas) instabilities, for example the “streaming instability.” In this instability, dust forms a thin, dense mid-plane layer with a dust density comparable to the gas density, leading to clumping of particles that grow and collapse under self-gravity to form planetesimals. This mechanism (and related fluid instability mechanisms) can allow for rapid growth of grains from the tens of cm to hundreds of m scales over short timescales required to bypass the meter-size barrier for planet formation.

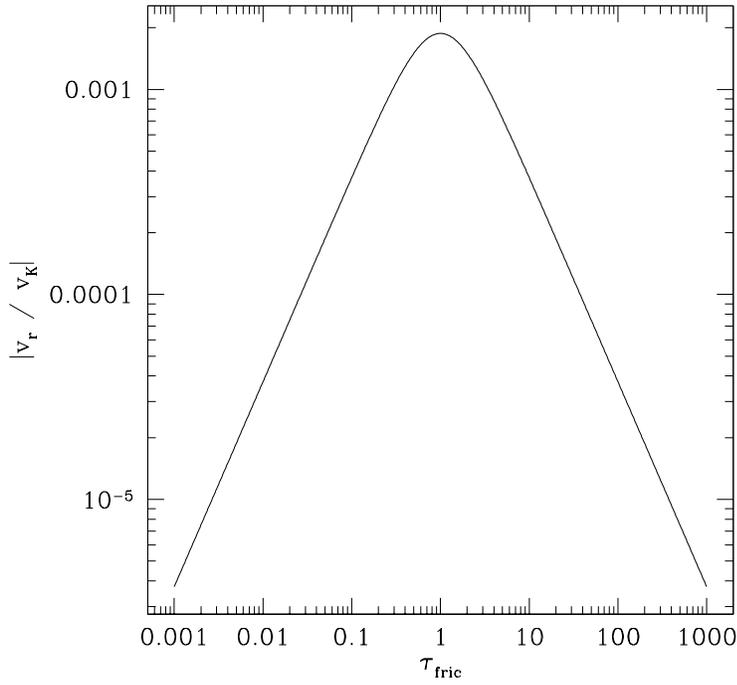


Figure 12.2: Radial drift velocity of particles at the midplane as a function of the dimensionless stopping time. The most rapid inward drift occurs for $\tau_{\text{fric}} \sim 1$, which corresponds to a stopping time of Ω_K^{-1} – particles in the 10 cm - 1 m size range. Adapted from Armitage (2007).

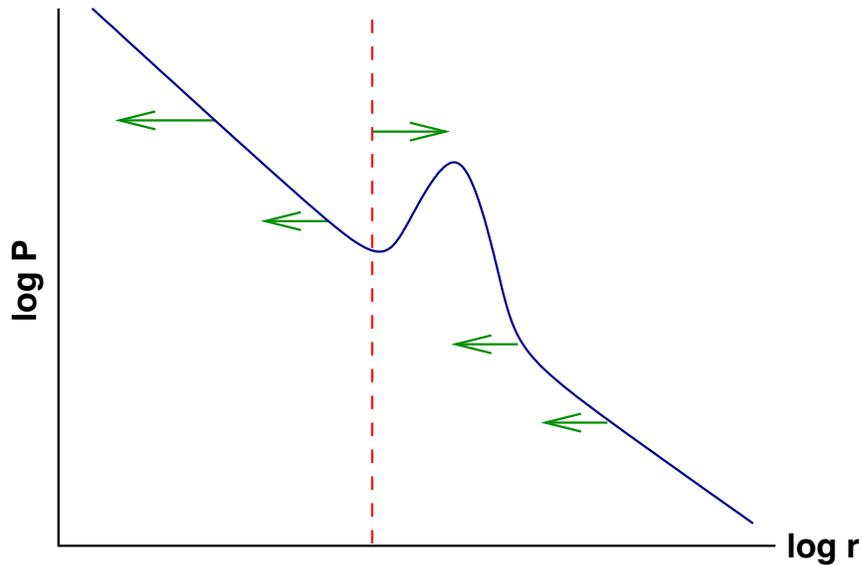


Figure 12.3: Schematic of how pressure maxima in a protoplanetary disk can cause the formation of dust rings due to radial drift and resulting pile-up of material. Adapted from Armitage (2007).

13 Planet formation: from pebbles to planets

Our agenda for Day 13 is the following:

1. Radial drift derivation activity (45 minutes)
2. Planetary accretion: Hill radius, isolation mass (30 minutes)
3. Gravitational focusing intro (if time)

Today's reading is Ch. III B-C of the Armitage notes, which covers planetary accretion from planetesimals and gas giant formation.

13.1 Radial drift activity

We'll start today's class by doing the promised activity solving for the radial drift of dust in the disk. I'm reproducing the text here so you don't have to look at the notes for Day 12, but please refer to them for more information.

Dust in the disk will drift radially due to interactions with gas. Small particles ($s \lesssim 1$ cm) are well-coupled to the gas, so they orbit the star at a velocity slightly smaller than the Keplerian velocity (recall Equation 11.10),

$$v_{\phi,g} = v_K \sqrt{1 - n \frac{h^2}{r^2}} = v_K \sqrt{1 - \eta}, \quad (13.1)$$

where $\eta = nc_s^2/v_K^2$, with n the power-law exponent for the radial pressure dependence. Because small particles orbit at this slower velocity, they will not be in radial force balance like the gas because the gas has pressure support (i.e., a pressure gradient term in the force balance) while the dust does not. This causes small dust grains to spiral in toward the host star at a radial terminal velocity. Conversely, large “rocks” with $s \gtrsim 1$ m feel gas drag because their orbits are Keplerian, while those of the gas are sub-Keplerian. This causes an effective headwind that removes angular momentum from the orbit of the rock, causing it to drift inward.

We can formalize these two concepts by writing a specific force balance for particles (of any size, with a given stopping timescale t_{fric}) in the radial and azimuthal directions. In the radial direction,

$$\frac{dv_r}{dt} = \frac{v_\phi^2}{r} - \Omega_K^2 r - \frac{1}{t_{\text{fric}}}(v_r - v_{r,\text{gas}}), \quad (13.2)$$

where v_r is the radial velocity of the particle, the first term on the RHS is the centrifugal force, the second term on the RHS is gravity, and the third term on the RHS represents gas drag. The azimuthal force balance is only dependent on gas drag,

$$\frac{d(rv_\phi)}{dt} = -\frac{r}{t_{\text{fric}}}(v_\phi - v_{\phi,\text{gas}}), \quad (13.3)$$

where v_ϕ is the azimuthal velocity of the particle.

We can now use these statements of radial and azimuthal force balance for dust particles to derive the radial velocity (toward/away from the star, not to/from Earth!) of dust particles in the disk. Please do so in small groups of 2-3, following these steps:

1. First, simplify the azimuthal equation (Equation 13.3) by assuming that the particle spirals in through a succession of nearly Keplerian orbits, i.e.,

$$\frac{d(rv_\phi)}{dt} \approx v_r \frac{d(rv_K)}{dr} = \frac{v_r v_K}{2}. \quad (13.4)$$

Plug this into the expression for azimuthal force balance to find an expression for $(v_\phi - v_{\phi,\text{gas}})$.

2. Simplify the radial equation (Equation 13.2) by substituting in $v_K^2 = v_{\phi,\text{gas}}^2 + \eta v_K^2$ and $\Omega_K^2 r = v_K^2/r$. There will be two terms on the right hand side of this equation relating to the azimuthal velocity, make the following first-order assumption that the gas and dust motions are similar

$$\frac{v_\phi^2}{r} - \frac{v_{\phi,\text{gas}}^2}{r} = \frac{(v_\phi + v_{\phi,\text{gas}})(v_\phi - v_{\phi,\text{gas}})}{r} \approx \frac{2v_K(v_\phi - v_{\phi,\text{gas}})}{r}, \quad (13.5)$$

to write the radial velocity equation for the dust to first-order accuracy.

3. Assume that there is no radial acceleration of the dust ($\frac{dv_r}{dt} \approx 0$) and use $(v_\phi - v_{\phi,\text{gas}})$ from the azimuthal force balance to derive the dependence of the radial velocity on $r, v_K, t_{\text{fric}}, v_{r,\text{gas}}, \eta$.
4. Re-cast the radial velocity by defining a dimensionless stopping time $\tau_{\text{fric}} \equiv t_{\text{fric}} \Omega_K = t_{\text{fric}} v_K / r$ to write v_r as a function of $\tau_{\text{fric}}, v_{r,\text{gas}}, v_K, \eta$.
5. The radial drift velocity peaks at $\tau_{\text{fric}} \approx 1$. What is the peak radial velocity, as a function of η, v_K ?
6. Assuming Epstein drag, what typical particle sizes does the peak radial velocity of a dust grain correspond to?

13.2 Accretion of planetesimals

Once planetesimals of sizes of hundreds of meters to hundreds of km form, they grow to form terrestrial planets and giant planet cores via accretion of material. The gas disk no longer regulates the radial motion of these planetesimals, and instead the physics of accretion is purely Newtonian. Effectively, the formation of planets from planetesimals requires studying the process of an up to hundred million year long “cascade” of pairwise accretion of solid bodies.

13.2.1 Gravitational focusing

A massive object will deflect the paths of other bodies toward it, increasing its effective cross-section for collisions. Figure 13.1 shows how gravitational focusing of two objects can cause them to collide, even from trajectories which otherwise would not collide without the effects of mutual gravity. It can be shown (activity in the next class) that the effective cross-section for collisions of two bodies with mass m is

$$\Gamma = \pi R_s^2 \left(1 + \frac{v_{\text{esc}}^2}{\sigma^2} \right) = \pi R_s^2 (1 + \theta), \quad (13.6)$$

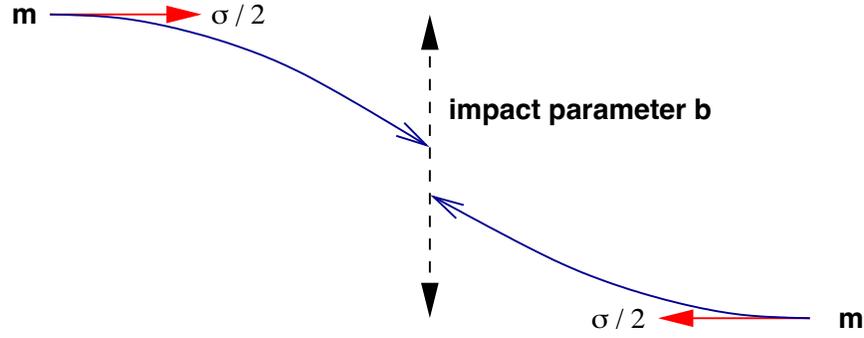


Figure 13.1: Schematic of how gravitational focusing enhances accretion. Objects with a mass m and velocity $\sigma/2$ are deflected by gravity with an impact parameter b , causing accretion. Adapted from Armitage (2007).

where R_s is the sum of the radii of the two objects, v_{esc} is the escape velocity of the two objects when they contact (i.e., $v_{\text{esc}}^2 = 4Gm/R_s$), and σ is the relative velocity of the two objects before they begin to gravitationally interact. $\theta \equiv v_{\text{esc}}^2/\sigma^2$ is the Safronov number, which determines the extent of gravitational focusing of the two bodies.

13.2.2 Hill radius

The Hill sphere is the region within which the gravitational force of a planet (or proto-planet) dominates over the tidal gravitational field of the star. The radius of the Hill sphere is

$$r_H = a \left(\frac{M_p}{3M_\star} \right)^{1/3}, \quad (13.7)$$

where a is semi-major axis and M_p and M_\star are the mass of the (proto)planet and star, respectively. The Hill sphere roughly demarcates the region around the (proto)planet within which it can gravitationally attract particles and accrete them.

13.2.3 Isolation mass

We can use the Hill sphere to estimate the mass of a protoplanet that has accreted all of the planetesimals in its vicinity. This is equivalently known as an object growing to its “isolation mass.” An object can only accrete planetesimals whose orbits lie within its feeding zone Δa , which extends some multiple of Hill radii Cr_h from the planet,

$$\Delta a = Cr_h. \quad (13.8)$$

A typical value of $C = 2\sqrt{3}$, derived from the maximum separation for which collisions between a planetesimal and protoplanet are possible in the three-body problem (Armitage, 2013). The isolation mass is then the mass of planetesimals in the feeding zone of the protoplanetary disk

$$M_{\text{iso}} = 2\pi a \cdot 2\Delta a \Sigma_p = 4\pi a^2 C \left(\frac{M_{\text{iso}}}{3M_\star} \right)^{1/3} \Sigma_p, \quad (13.9)$$

where we have used $\delta a \approx r_H$ and Σ_p is the column mass of planetesimals in the disk, approximately 0.01Σ (i.e., 1/100th of the total column mass). Solving for the isolation mass, we find

$$M_{\text{iso}} = \frac{8}{\sqrt{3}} \pi^{3/2} C^{3/2} M_{\star}^{-1/2} \Sigma_p^{3/2} a^3 . \quad (13.10)$$

For typical disk values, at 1 au $M_{\text{iso}} \sim 0.1 M_{\oplus}$, and at 5 au $M_{\text{iso}} \sim 10 M_{\oplus}$. This implies that terrestrial planets must grow to Earth-like masses through accretion between planetary embryos (each of which are $\sim 0.1 M_{\oplus}$), while the isolation mass for Jupiter corresponds to the approximate expected mass of a giant planet core.

14 Planet formation: accretion, orbital migration and evolution

Our agenda for Day 14 is the following:

1. Gravitational focusing activity (30 minutes)
2. Terrestrial planet formation (15 minutes)
3. Giant planet formation (15 minutes)
4. Nice and Grand Tack Models (15 minutes)

Today's reading is Ch. IV of the Armitage notes, which covers the evolution of planetary systems and the Nice model.

14.1 Gravitational focusing activity

Please work on this activity in small groups of 2-3, and be prepared to write your solutions on the board. Figure 13.1 shows a schematic of the gravitational focusing of two masses with mass m that collide from an initial impact parameter b and initial velocities each of $\sigma/2$.

1. Assuming energy conservation, write down an expression that equates the initial kinetic energy of the objects with the sum of their kinetic energy and gravitational potential energy at closest approach. Assume that their velocity at closest approach is v_c and their separation at closest approach is R_c .
2. Assuming angular momentum conservation, derive an expression for v_c as a function of b and R_c .
3. Substitute your expression for v_c into part (a) to derive an expression for the largest impact parameter b that will lead to a collision between the objects, given that the sum of the radii of the two objects is R_s .
4. Re-arrange this expression to determine the cross-section for collisions $\Gamma = \pi b^2$ as a function of R_s , the mutual escape velocity of the objects at the point of contact v_{esc} , and the sum of the initial velocities σ .
5. Estimate the effective cross-section for collisions of two embryos with masses and radii equal to that of Mars ($M \approx 0.11 M_{\oplus}$, $R \approx 0.53 R_{\oplus}$), assuming a relative velocity at infinity $\sigma = 100 \text{ m s}^{-1}$. Compare this to the physical cross section of each object.

14.2 Steps in the formation of terrestrial planets

The formation of terrestrial planets can be separated into five main stages:

1. The agglomeration of small (starting with sub-micron-sized) dust particles to form cm-m sized "pebbles." The coagulation of dust is mediated by electrostatic forces, which allows dust to grow to pebbles via pairwise collisions.

2. The growth of dust and pebbles to planetesimals. This requires a bypass of the meter-sized barrier, which likely occurs through some form of gravitational instability of solid material in the disk, perhaps mediated by the streaming instability. Once these planetesimals form, they can rapidly grow through accretion of pebbles due to radial drift and gravitational focusing.
3. Runaway growth of the largest planetesimals to become planetary embryos, with a resulting phase of oligarchic growth where embryos grow more slowly until they reach the isolation mass. At this point, each embryo has accreted all material in its feeding zone.
4. Collisions between planetary embryos result in growth of planets to their final masses. The final giant impact between planet and embryo is the point at which the formation of the planet has ceased, and evolution has begun.

14.3 Formation of gas giant planets

Forming a gas giant planet requires an enormous amount of pairwise accretion of rocky bodies and later gas accretion. Figure 14.1 summarizes the challenges in the formation of gas giant planets from micron-sized dust bunnies – the total mass growth is a factor of $\sim 10^{42}$!

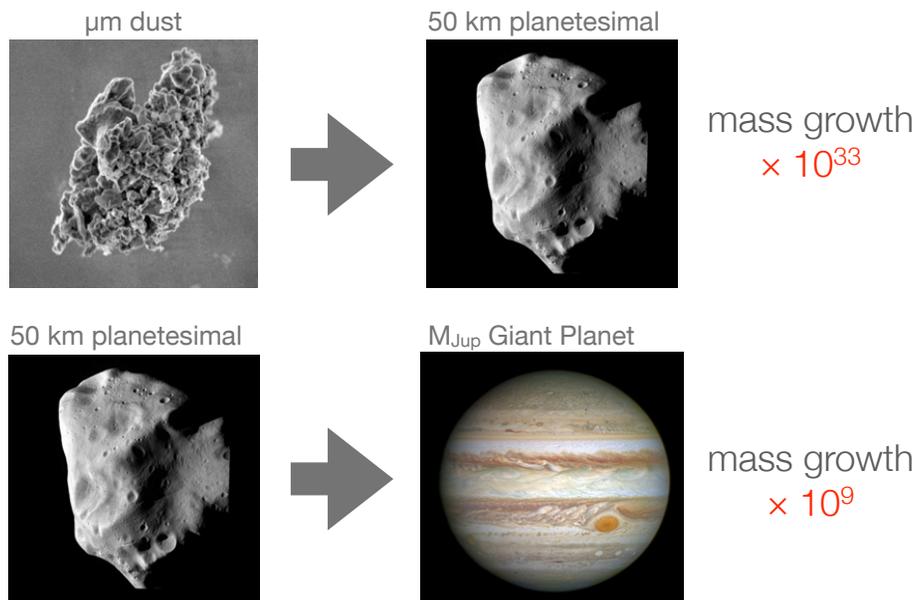


Figure 14.1: Stages in the growth of planets from micron-sized dust bunnies to Jupiter-like masses. Figure courtesy Andrew Youdin.

14.3.1 Gravitational instability

One possible mechanism to form a giant planet is through local collapse of a gravitationally unstable disk. This mechanism can only work in massive and/or cold protoplanetary disks, which are gravitationally unstable on a large scale, leading to instabilities that result in clumping of material and resulting collapse (see Figure 14.2). Gravitational instability of

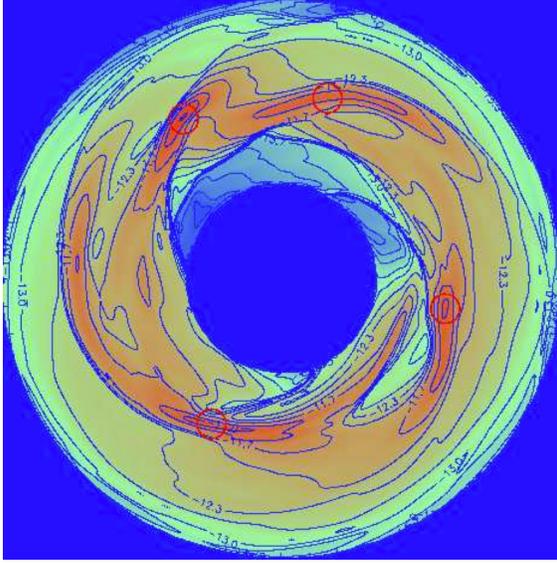


Figure 14.2: Clumping of a gravitationally unstable protoplanetary disk. Figure adapted from Boss (2011).

a disk is only possible if the Toomre Q parameter is sufficiently low

$$Q \equiv \frac{c_s \Omega}{\pi G \Sigma} \lesssim 1, \quad (14.1)$$

where the dependence on c_s/Σ implies that only cold, massive disks will have $Q < 1$. A second criterion for gravitational collapse to form a planet (rather than just a gas and dust clump) is that the collapsing clump be able to cool on a timescale shorter than the collapse timescale

$$\tau_{\text{cool}} < \tau_{\text{ff}}. \quad (14.2)$$

This latter criterion is challenging to satisfy, as the radiative cooling timescale scales as Σ/T^3 – the opposite dependence of the Toomre Q parameter on Σ and T ! Thus, if a disk is sufficiently massive and cool to be gravitationally unstable, it is also likely to be locally too massive and cool to efficiently radiatively cool. This will then prevent a clump from collapsing and forming a planet. Because cooling is required to accrete, it is expected that gravitational instability generally forms objects more massive than planets (i.e., brown dwarfs), and it is not expected to be the dominant gas giant formation mechanism.

14.3.2 Core accretion

Core accretion is the widely accepted theory of gas giant formation at present. Core accretion is a bottom-up mechanism to form gas giant planets, with three stages outlined in Figure 14.3. The first stage of core accretion is for a massive rocky core to form near its isolation mass, similar to the process of terrestrial planet growth. This core will grow massive enough that it will then accrete some gas from the protoplanetary disk, and further grow via accretion of planetesimals and pebbles. Once this core hits a “critical core mass,” it will then undergo runaway gas accretion from the protoplanetary disk, rapidly growing to a Jupiter-like mass. The critical core mass is expected to be $\sim 10 M_{\oplus}$ for typical disk conditions. This is significantly larger than the masses that rocky objects can grow to in the MMSN, resulting in the need for incorporation of ice in the formation of a giant planet core.

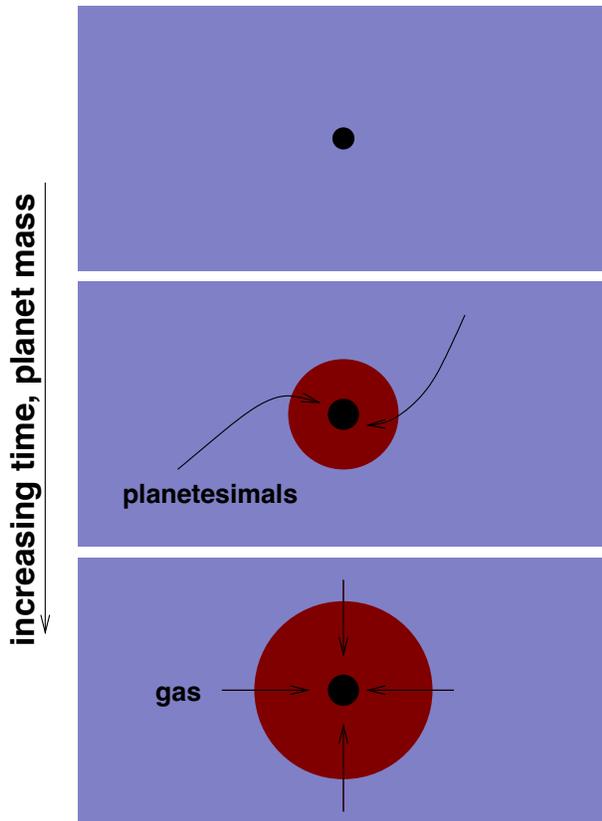


Figure 14.3: Schematic of the stages of the core accretion model for the formation of gas giant planets. Adapted from Armitage (2007).

As a result, the standard model for giant planet formation includes growth of a core outside the snow line in the protoplanetary disk.

Figure 14.4 shows the classic model of core accretion from Pollack et al. (1996). This simulation finds the same three main stages of gas giant formation: initial growth of a core, hydrostatic accretion of gas and planetesimals, and then runaway (non-hydrostatic) growth by accretion of gas from the protoplanetary disk. The key challenge in the core accretion paradigm is reaching the critical core mass quickly, as the gas disk is expected to be lost by 3-5 Myr in most systems. This can be accomplished by reducing the opacity of the gas, which reduces the needed critical core mass for runaway. Another possibility is pebble accretion, which rapidly grows the core mass through radial drift of grains.

14.4 Migration

There are three main mechanisms through which planets can undergo orbital migration (i.e., have a changing semi-major axis with time) during the epoch of planet formation.

14.4.1 Type I migration

Type I migration causes the radial motion of approximately Earth-mass planets through protoplanetary disks. In Type I migration, the planet exerts a negligible influence on the gas in the disk, with torques from the disk controlling the motion of the planet. Type I migration arises due to net torques that arise from Lindblad resonances risen on the planet from the disk interior and the disk exterior to the planet. These Lindblad resonances occur

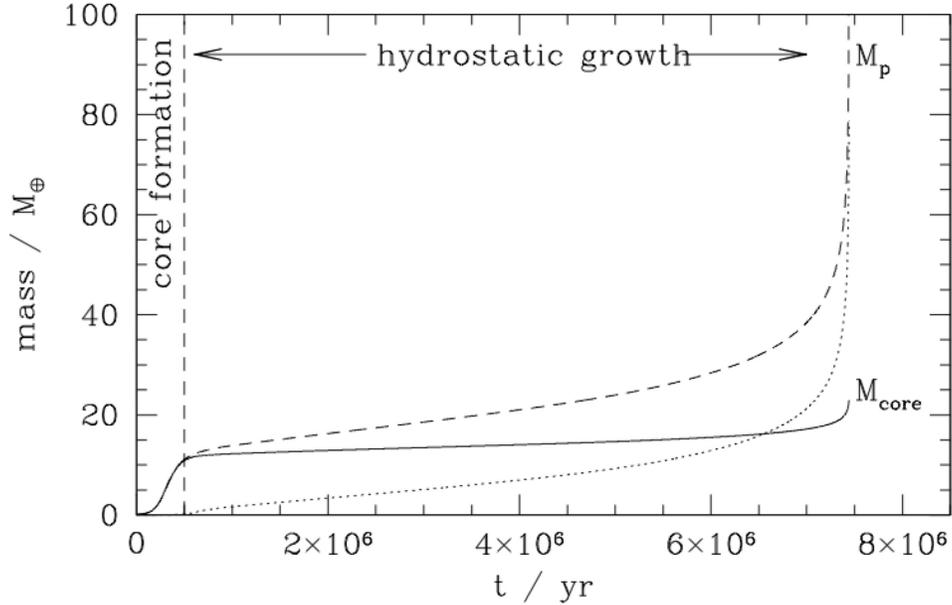


Figure 14.4: Simulation of the growth of a giant planet by core accretion. The solid line shows the evolution of the mass of the core, the dashed line the mass of gas, and the dotted line the total planet mass. Adapted from Pollack et al. (1996).

(for Keplerian orbits) when the angular frequency of the gas is at some integer multiple of the difference between the angular frequency of the gas and the angular frequency of the planet, resulting in the planet gaining angular momentum from interior Lindblad resonances (driving the planet outward) and losing angular momentum from outer Lindblad resonances (driving the planet inward). Generally, the Lindblad resonances exterior to the planet are dominant, resulting in a net torque that causes inward migration of the planet.

14.4.2 Type II migration

Type II migration occurs when a planet is sufficiently massive that it opens a gap in the disk. There are two conditions for gap opening, first that the Hill radius of the planet is greater than the scale height of the disk ($r_H \gtrsim h$), and secondly that the torques removing gas from the disk are faster than the action of viscosity to diffuse the gas back into the disk, which requires a planet-to-star mass ratio $q \gtrsim 10^{-4}$. Figure 14.5 shows the planet-disk interaction in the regime where the planet has opened a gap, with streams of gas that are flowing onto the planet from the inner and outer regions of the disk. Type II migration occurs from planet-disk interactions because the orbital evolution of the planet is directly coupled to the evolution of gas in the disk, the latter of which is controlled by viscosity. This causes the massive planet to follow the sense of gas motion, which is generally inward (toward the star) in regions of the inner disk near the snow line where giant planets are expected to form. However, unlike Type I migration, Type II migration need not be inward – at sufficient distances from the star, it may lead to outward migration of gas giant planets.

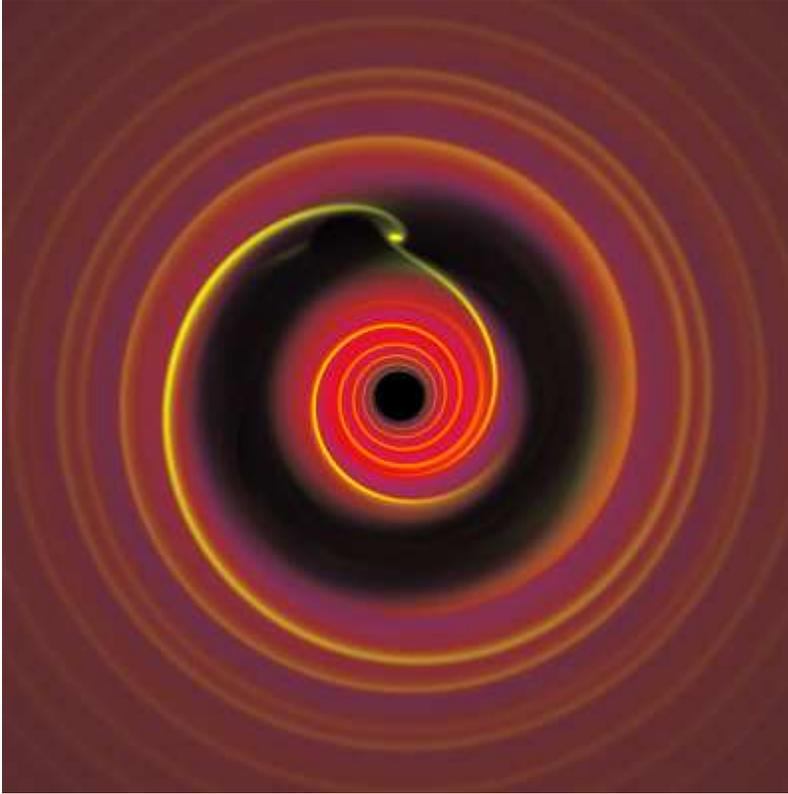


Figure 14.5: Simulation of a gas giant planet that has opened a gap in the protoplanetary disk, with resulting interaction with the disk causing it to migrate inward. Adapted from Armitage (2007).

14.4.3 Planetesimal disk migration

Even after planets form and migrate within the disk through Type I or II migration in a given system, it is likely that planetesimals will remain in the system that have not been incorporated into planets. This planetesimal disk can gravitationally interact with the planet, leading to individual planetesimals being scattered either outward or inward by a planet. Each individual scattering event must conserve angular momentum, so outward scattering of a planetesimal leads to inward migration of a planet, and vice versa. Importantly, the process of planetesimal disk migration can last far longer than the processes of Type I and II migration, which are limited to occur only during the lifetime of the gaseous disk. Planetesimal disk migration is expected to have occurred in our Solar System, leading to a dynamical instability of the gas giant planets and resulting disruption of our system's nascent planetesimal belt.

14.5 Models for Solar System evolution

The Nice model (named after Nice, a seaside city in southern France) is an umbrella term for a variety of dynamical models which predict that our Solar System underwent a large-scale dynamical instability a few hundred Myr after its formation. These models assume that the orbits of the gas and ice giants in the Solar System began more compact than today, with Uranus and Neptune notably at semi-major axes approximately half their present values. Due to Type II migration, Jupiter and Saturn migrate to be in a near-resonant configuration. Then, due to interaction with a planetesimal disk beyond Neptune (encompassing and slightly inward of the current location of Kuiper Belt Objects), the

resonant chain between Jupiter and Saturn (and perhaps Uranus and Neptune) is broken, leading to scattering between the giant planets. This drives rapid outward orbital evolution of Uranus and Neptune, scattering the planetesimal disk, leading to the ejection of a significant portion of this planetesimal disk, with some of it scattered inward, perhaps leading to some of the last large impact basins on the Moon. The remaining objects in the planetesimal disk comprise our current suite of Kuiper Belt Objects, including Pluto, Charon, and Arrokoth.

There are variants of the Nice model that incorporate additional processes which are of interest. One variant of the Nice model includes a fifth gaseous planet in the outer Solar System, approximately the mass of Uranus or Neptune. This additional ice giant is theorized to have been lost from the Solar System due to planet-planet scattering and would now be part of the population of free-floating planets. Another variant is named the “Grand Tack” model, and assumes that Jupiter underwent inward Type II migration before this late-stage scattering event, moving inward to ~ 2 au and then back outward. The Grand Tack model improves upon the Nice model by further explaining the deficit of material in the asteroid belt as well as the low mass of Mars, which is effectively a stranded planetary embryo.

15 Exoplanet atmospheres: structure, composition, chemistry, loss

Our agenda for Day 16 is the following:

1. Recap atmospheric structure, dry and moist adiabats (15 minutes)
2. Stratified atmospheres: radiative relaxation, radiative timescale activity (30 minutes)
3. Atmospheric composition and chemistry (20 minutes)
4. Atmospheric loss, cosmic shoreline (10 minutes)

Today's reading is Sections 3 and 5 of the Zhang Atmospheres on Exoplanets and Brown Dwarfs review paper. This will cover atmospheric loss (Section 3) and atmospheric composition (Section 5). I'm including material that Dr. Lothringer covered in Day 15 in these notes as well for completeness.

15.1 Hydrostatic equilibrium

Recall from our discussion of disks that hydrostatic balance, where pressure gradients balance gravity, implies that the variation of pressure with height is (Equation 10.4):

$$\frac{dp}{dz} = -\rho g . \quad (15.1)$$

This balance applies equally well to disks and atmospheres. Note that if we integrate this equation vertically assuming constant gravity, we find that the surface density (i.e., mass per area) in a given atmosphere column increases with pressure and decreases with gravity

$$\frac{\text{mass}}{\text{area}} = \Sigma = \frac{p}{g} . \quad (15.2)$$

We can use the ideal gas equation to relate p to ρ , here in a format often used in atmospheric science

$$p = \rho RT , \quad (15.3)$$

where $R = R_u/(\mu m_p)$ is the *specific* gas constant, which depends on the atmospheric species of interest, with $R_u = 8.3145 \text{ J mol}^{-1} \text{ K}^{-1}$ the *universal* gas constant. Substituting this into the expression for hydrostatic equilibrium and integrating, we find

$$\begin{aligned} \frac{1}{p} \frac{dp}{dz} &= -\frac{g}{RT} \\ \ln\left(\frac{p}{p_0}\right) &= -\int_{z_0}^z \frac{g}{RT} dz' \\ p &= p_0 \exp\left(-\int_{z_0}^z \frac{g}{RT} dz'\right) = p_0 \exp\left(-\int_{z_0}^z \frac{dz'}{H}\right) , \end{aligned} \quad (15.4)$$

where $H = RT/g$ is the pressure scale height. For an isothermal atmosphere, the expression for hydrostatic equilibrium simplifies to

$$p = p_0 e^{-z/H} , \quad (15.5)$$

and thus the pressure (and density, for an isothermal atmosphere) decreases with increasing height over a characteristic e-folding distance of H . Figure 15.1 shows that as expected, in Earth's atmosphere the dependence of pressure and density on height is approximately exponential.

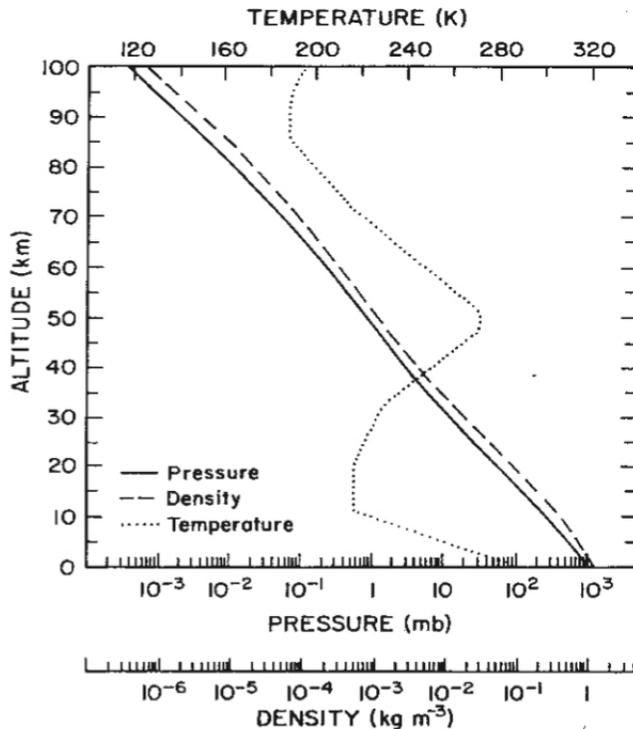


Figure 15.1: U.S. standard atmosphere of Earth. Shown are pressure, density, and temperature profiles for typical Earth climate conditions.

Global-mean pressure (solid), density (dashed), and temperature (dotted), as functions of altitude.
Source: U.S. Standard Atmosphere (1976).

However, the dependence of temperature on height is not so straightforward – for Earth there are multiple atmospheric levels, defined by the temperature gradient. The region near the surface where the temperature decreases with increasing height is known as the *troposphere*, with the region overlying that where the temperature increases with height known as the *stratosphere*. Near-surface tropospheres and overlying stratospheres are ubiquitous in our Solar System, as shown in Figure 15.2. We'll next dive into the thermodynamic properties that control these vertical temperature profiles.

15.2 Atmospheric thermodynamics

15.2.1 First law of thermodynamics

The first law of thermodynamics is a statement of energy conservation – namely, that the rate of change in the internal energy of a system is the sum of the working and heating rates

$$\frac{dU}{dt} = Q + W, \quad (15.6)$$

where U is the internal energy, Q is the heating rate (which represents Q interactions with environments at different temperature through conduction and radiation), and W is the

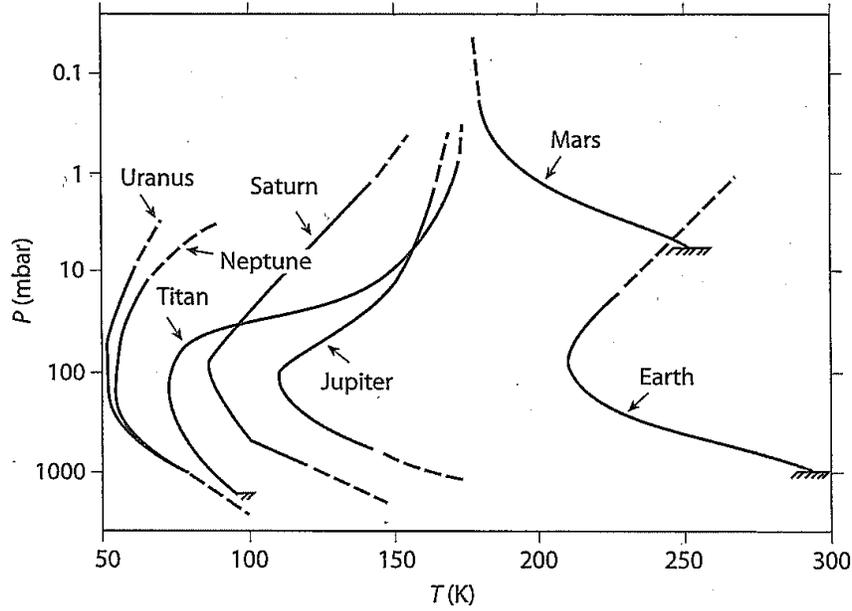


Figure 15.2: Atmospheric temperature-pressure profiles of planets in our Solar System with significant atmospheres (not including Venus), along with Titan. Figure adapted from Seager (2010).

working rate (which represents a mechanical exchange of energy with the environment). For an atmosphere, the working rate is dominated by the expansion work done *on* the system

$$W = -p \frac{dV}{dt}, \quad (15.7)$$

where V is the volume of the system. Thus, for an atmosphere we can write the first law of thermodynamics as

$$\frac{dU}{dt} = Q - p \frac{dV}{dt}. \quad (15.8)$$

For studies of atmospheres, an alternate version of the first law of thermodynamics is often used with enthalpy $H = U + pV$ instead of internal energy,

$$\frac{dH}{dt} = Q + V \frac{dp}{dt}. \quad (15.9)$$

15.3 Specific heats

The chain rule can be used to re-state the change in internal energy as

$$\frac{dU}{dt} = \frac{\partial U}{\partial T} \frac{dT}{dt} + \frac{\partial U}{\partial V} \frac{dV}{dt}. \quad (15.10)$$

We can define the heat capacity at constant volume C_v to be

$$C_v \equiv \left(\frac{\partial U}{\partial T} \right)_v, \quad (15.11)$$

which is a *specific heat* (c_p) when defined per mass or per mole. Note that for an ideal gas changes in internal energy due to changes in volume can be neglected, so we can write the change in internal energy as

$$\frac{dU}{dt} = C_v \frac{dT}{dt}. \quad (15.12)$$

We can equivalently perform the chain rule on the change in enthalpy with time to express

$$\frac{dH}{dt} = \frac{\partial H}{\partial T} \frac{dT}{dt} + \frac{\partial H}{\partial p} \frac{dp}{dt}. \quad (15.13)$$

We can similarly define the heat capacity at constant pressure C_p to be

$$C_p \equiv \left(\frac{\partial H}{\partial T} \right)_p. \quad (15.14)$$

For an ideal gas, the enthalpy changes due to pressure changes are negligible, so we can write the change of enthalpy as

$$\frac{dH}{dt} = C_p \frac{dT}{dt}. \quad (15.15)$$

For an ideal gas, the heat capacity and specific heat are related as

$$\begin{aligned} C_p &= C_v + NR_u, \\ c_p &= c_v + R_u, \end{aligned} \quad (15.16)$$

where N is the number of moles in the system. The heat capacity C_v per molecule is $k_B/2$ per degree of freedom, and note that the universal gas constant is related to the Boltzmann constant as $R_u = N_A k_B$ with $N_A = 6.022 \times 10^{23} \text{ mol}^{-1}$ Avogadro's number. Thus, the specific heat capacity c_v can also be expressed as $R_u/(2)$ per degree of freedom. For diatomic molecules, there are 3 translational and 2 rotational degrees of freedom, leading to $c_v = 5/2 R_u$ and $c_p = 7/2 R_u$. Two common combinations of specific heats are

$$\begin{aligned} \gamma &\equiv \frac{c_p}{c_v}, \\ \kappa &\equiv \frac{R}{c_p}, \end{aligned} \quad (15.17)$$

where for Earth air $\gamma \approx 1.4$ and $\kappa \approx 0.286$. The specific heat capacities and specific heat ratios of relevant atmospheric gases are shown in Figure 15.3.

15.3.1 Convective instability

For adiabatic processes in an ideal gas, $Q = 0$ in the first law of thermodynamics. If we can further consider the specific heat capacity and gas constant to be constant, this leads to an expression for the change in temperature with height of an adiabatic atmosphere (as derived by Prof. Lothringer):

$$\frac{dT}{dz} = -\frac{g}{c_p}. \quad (15.18)$$

Thus, the lapse rate of an adiabatic parcel depends on both the gravity of the planet and the composition of the gas through c_p .

	H_2O	CH_4	CO_2	N_2	O_2	H_2	He	NH_3
Crit. point T	647.1	190.44	304.2	126.2	154.54	33.2	5.1	405.5
Crit. point p	221.e5	45.96e5	73.825e5	34.0e5	50.43e5	12.98e5	2.28e5	112.8
Triple point T	273.15	90.67	216.54	63.14	54.3	13.95	2.17	195.4
Triple point p	611.	.117e5	5.185e5	.1253e5	.0015e5	.072e5	.0507e5	.061e5
L vap(b.p.)	22.55e5	5.1e5	–	1.98e5	2.13e5	4.54e5	.203e5	13.71e5
L vap(t.p.)	24.93e5	5.36e5	3.97e5	2.18e5	2.42e5	??	??	16.58e5
L fusion	3.34e5	.5868e5	1.96e5	.2573e5	.139e5	.582e5	??	3.314e5
L sublimation	28.4e5	5.95e5	5.93e5	2.437e5	2.56e5	??	??	19.89e5
ρ liq(b.p.)	958.4	450.2	1032.	808.6	1141.	70.97	124.96	682.
ρ liq(t.p.)	999.87	??	1110.	??	1307.	??	??	734.2
ρ solid	917.	509.3	1562.	1026.	1351.	88.	200.	822.6
$c_p(0C/1bar)$	1847.	2195.	820.	1037.	916.	14230.	5196.	2060.
$\gamma(c_p/c_v)$	1.331	1.305	1.294	1.403	1.393	1.384	1.664	1.309

Table 2.1: Thermodynamic properties of selected gases. Latent heats of vaporization are given at both the boiling point (the point where saturation vapor pressure reaches $1bar$) and the triple point. Liquid densities are given at the boiling point and the triple point. For CO_2 the 'boiling point' is undefined, so the liquid density is given at $253K/20bar$ instead. Note that the maximum density of liquid water is $1000.00kg/m^3$ and occurs at $-4C$. Densities of solids are given at or near the triple point. All units are mks , so pressures are quoted as Pa with the appropriate exponent. Thus, $1bar$ is written as $1e5$ in the table.

Figure 15.3: Relevant thermodynamic data for common atmospheric gases. Table adapted from Pierrehumbert (2010).

Atmospheres with temperature profiles that decrease more sharply with height than the adiabatic profile are unstable to convection. This is because if we perturb a parcel at the same density as the surroundings, the parcel's density will change according to the adiabatic relation while keeping a pressure close to that of the surroundings. If the density is less than the environmental density for upward displacements (or greater for downward displacements), then the parcel will accelerate away from its initial position – this will initiate convection. Figure 15.4 shows a schematic of convectively unstable and stable temperature profiles compared to the adiabatic lapse rate. If the parcel is hotter than the surrounding air than it will be less dense (and if it is cooler, then it will be more dense). The atmosphere will then be unstable to convection if temperature decreases with height faster than the adiabatic lapse rate $-g/c_p$ and stable if temperature decreases with height more slowly than the adiabatic lapse rate. Thus, we can express the convective stability of the atmosphere as

$$\begin{aligned}
\frac{dT}{dz} &< -\frac{g}{c_p} \text{ unstable,} \\
\frac{dT}{dz} &= -\frac{g}{c_p} \text{ marginally stable,} \\
\frac{dT}{dz} &> -\frac{g}{c_p} \text{ stable.}
\end{aligned}
\tag{15.19}$$

Atmospheres that are forced to be unstable will adjust via convection to have lapse rates that are approximately stable, with $dT/dz \approx -g/c_p$. In general, planetary atmospheres are

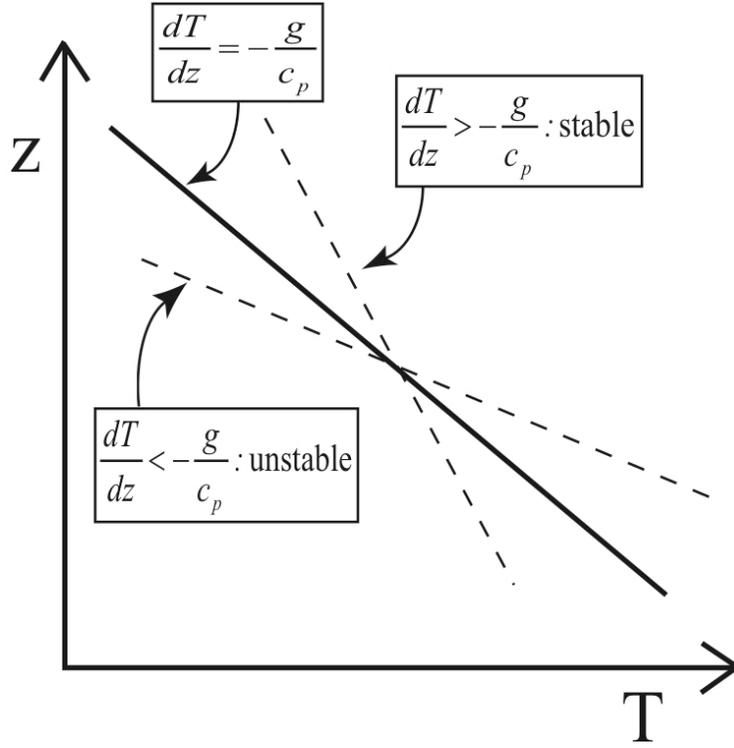


Figure 15.4: Schematic showing the convective lapse rate along with convectively unstable and stable temperature profiles. Figure courtesy Adam Showman.

stable to convection except at deep levels, either near the surface for rocky planets or within the deep envelope and interior in gas giant planets.

Note that potential temperature is a very useful quantity to determine stability of a planetary atmosphere. The potential temperature is defined such that

$$d\ln\theta \equiv d\ln T - \frac{R}{c_p} d\ln p, \quad (15.20)$$

and is equivalent to entropy in an atmospheric context. Integrating for constant $R/c_p \equiv \kappa$, we find

$$\theta = T \left(\frac{p_0}{p} \right)^\kappa. \quad (15.21)$$

We can further relate the potential temperature variation with height $d\theta/dz$ to the lapse rate dT/dz as

$$\begin{aligned} \frac{1}{\theta} \frac{d\theta}{dz} &= \frac{1}{T} \frac{dT}{dz} - \frac{\kappa}{p} \frac{dp}{dz}, \\ \frac{1}{\theta} \frac{d\theta}{dz} &= \frac{1}{T} \frac{dT}{dz} + \frac{\kappa}{p} \rho g, \\ \frac{1}{\theta} \frac{d\theta}{dz} &= \frac{1}{T} \left(\frac{dT}{dz} + \frac{g}{c_p} \right), \end{aligned} \quad (15.22)$$

where we have used hydrostatic equilibrium in step 2 and the ideal gas law in step 3. Note that this implies that if $d\theta/dz = 0$, the lapse rate $dT/dz = -g/c_p$ (i.e., a dry adiabat). Thus, we can write an equivalent expression to Equation (15.19) using potential temperature:

$$\begin{aligned} \frac{d\theta}{dz} &< 0 \text{ unstable,} \\ \frac{d\theta}{dz} &= 0 \text{ marginally stable,} \\ \frac{d\theta}{dz} &> 0 \text{ stable.} \end{aligned} \tag{15.23}$$

15.3.2 Condensation, clouds and the moist adiabat

Clouds can form in planetary atmospheres where species are thermodynamically favored to undergo a phase transition from vapor to solid phases. Condensation and deposition are the process of conversion of species from gaseous to liquid, or gaseous to solid, respectively – generally speaking, cloud formation (both condensation and deposition) are often referred to simply as “condensation” into liquid or solid phases. Cloud condensation is thermodynamically allowed when the partial pressure of a species exceeds its saturation vapor pressure. The saturation vapor pressure is set by the Clausius-Clapeyron relationship, which for an ideal gas can be expressed as

$$\frac{dp_{\text{sat}}}{dT} = \frac{1}{T} \frac{L}{\rho_{\text{vap}}^{-1} - \rho_{\text{cond}}^{-1}}. \tag{15.24}$$

If $\rho_{\text{cond}} \gg \rho_{\text{vap}}$ (as is the case for most relevant atmospheric condensibles), then (using the ideal gas law) Clausius-Clapeyron simplifies to

$$\frac{dp_{\text{sat}}}{dT} = \frac{p_{\text{sat}}L}{RT^2}, \tag{15.25}$$

where p_{sat} is the saturation vapor pressure and L is the latent heat of condensation or fusion of the species of interest. This can be directly integrated to obtain the saturation vapor pressure curve

$$p_{\text{sat}} = p_0 \exp \left[\frac{L}{R} \left(\frac{1}{T_0} - \frac{1}{T} \right) \right], \tag{15.26}$$

implying that the saturation vapor pressure scales exponentially with $-1/T$. As a result, the saturation vapor pressure of species decreases with decreasing temperature, as shown in Figure 15.5.

We can equivalently state that cloud condensation is thermodynamically allowed at a given pressure when $T < T_{\text{cond}}$, which then occurs when temperature profiles cross the condensation curves in Figure 15.5. This is why clouds generally form at high altitudes in planetary atmospheres – because the temperature decreases with height in the troposphere, clouds form at altitudes where the temperature is cool enough to allow for condensation. The exception is when the local saturation vapor pressure is enhanced (e.g., due to mixing), increasing the partial pressure of a species to allow for condensation.

Note that for an atmosphere composed of a single (condensible) component, the Clausius-Clapeyron relationship can be expressed as a lapse rate

$$\frac{d \ln T}{d \ln p} = \frac{RT}{L}, \tag{15.27}$$

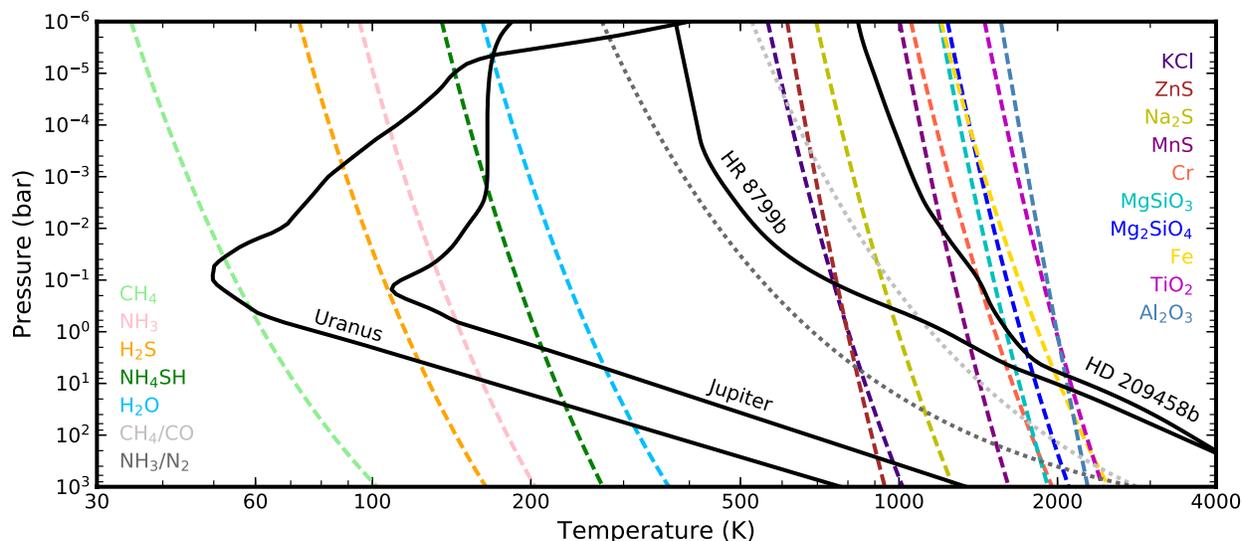


Figure 15.5: Temperature profiles of Solar System giant planets and exoplanets (black lines) compared with condensation curves of various species (dashed lines) and equilibrium chemistry equivalency curves (dotted lines). Figure adapted from Gao et al. (2021).

which is the single component moist adiabat. Given that the dry adiabatic gradient is $d\ln T/d\ln p = R/c_p$, the moist adiabatic gradient is smaller than the dry adiabatic gradient if $L > c_p T$ – which is almost always satisfied for most condensible species of interest, including water, carbon dioxide, methane, and nitrogen. The more general form of the moist adiabat for a dilute condensible species can be expressed as

$$\frac{d\ln T}{d\ln p} = \frac{R_u + \frac{L\xi}{T}}{c_p + \frac{L^2\xi}{R_u T^2}}, \quad (15.28)$$

where $\xi = p_{\text{sat}}/p_d$ is the mixing ratio of the condensible gas, i.e., the ratio of the saturation vapor pressure to the partial pressure of the dry gas component p_d .

15.4 Radiative relaxation

In atmospheres that are convectively stable, radiative heat transport controls the temperature profile. We’ll study the temperature profiles of radiative atmospheres in more detail in upcoming lectures, especially when we cover observational characterization via emission spectroscopy. First, let’s simply calculate how close to a state of radiative equilibrium an atmosphere will be in by estimating its radiative timescale, i.e., the time over which the atmosphere will adjust back to a state of radiative equilibrium.

First, if the atmosphere’s cooling is dominated by radiation, the rate of change of enthalpy must be equal to the outgoing longwave radiation of the planet

$$\frac{dH}{dt} = A\sigma T^4, \quad (15.29)$$

where A is surface area. We can re-write this in terms of specific enthalpy h as

$$m \frac{dh}{dt} = A\sigma T^4, \quad (15.30)$$

and given hydrostatic balance $m/A = \Sigma = p/g$ and assuming an ideal gas where $h = c_p T$ (along with ignoring the dependence of c_p on temperature) we can write

$$c_p \frac{p}{g} \frac{dT}{dt} = \sigma T^4. \quad (15.31)$$

Now, we can determine how long it takes the atmosphere to adjust back to radiative equilibrium from a small temperature perturbation δT . We can write the temperature as

$$T = T_0 + \delta T, \quad (15.32)$$

where T_0 is the radiative equilibrium temperature. Inserting this into our expression, we can write

$$c_p \frac{p}{g} \frac{dT_0}{dt} + c_p \frac{p}{g} \frac{d(\delta T)}{dt} = \sigma (T_0 + \delta T)^4. \quad (15.33)$$

Given that the perturbation is small, we can first-order Taylor expand the right hand side to write

$$c_p \frac{p}{g} \frac{dT_0}{dt} + c_p \frac{p}{g} \frac{d(\delta T)}{dt} = \sigma T_0^4 + 4T_0^3 \delta T. \quad (15.34)$$

Given that the radiative equilibrium state is not changing and that we only consider terms related to the perturbation, we can simplify this to

$$c_p \frac{p}{g} \frac{d(\delta T)}{dt} = 4T_0^3 \delta T. \quad (15.35)$$

Then, we can scale this equation using $d(\delta T)/dt \approx \delta T/\tau_{\text{rad}}$ to write an expression for the radiative timescale

$$\tau_{\text{rad}} \approx \frac{p}{g} \frac{c_p}{4\sigma T^3}. \quad (15.36)$$

The radiative timescale is thus shorter for hotter, thinner atmospheres that have a higher gravity, and vice versa. We'll next apply this to consider how close to radiative equilibrium different planetary atmospheres should be.

15.4.1 Radiative timescale activity

Let's calculate the radiative timescale for different planetary atmospheres to get a sense for how close each atmosphere is to a state of radiative equilibrium. Split into 6 groups, and depending on your group you'll calculate the radiative timescale of different atmospheres. We'll then compare our results.

1. (Groups 1-2) Calculate the radiative timescale at the surface of Earth, assuming that Earth has a temperature equal to its zero-albedo full-redistribution equilibrium temperature and an atmosphere comprised entirely of N_2 .
2. (Groups 3-4) Calculate the radiative timescale at the 1 bar (i.e., 10^5 Pa) level in Jupiter's atmosphere, assuming that Jupiter has a temperature equal to its zero-albedo full-redistribution equilibrium temperature and an atmosphere comprised entirely of H_2 .

- (Groups 5-6) Calculate the radiative timescale at the 1 bar level of a 51 Peg b-like hot Jupiter that orbits a Sun-like star with a separation of 0.05 au. Assume that the 1 bar temperature is equal to the zero-albedo full-redistribution equilibrium temperature and that the planet has an atmosphere comprised entirely of H_2 .

15.5 Atmospheric composition

15.5.1 Compositional diversity

There are two primary ways of measuring the composition of exoplanets. The first is derived from astrophysical measurements of the “metallicity” of stars, and measures the heavy element abundance relative to hydrogen (M/H), in turn relative to our own Sun. Metallicity is best used for gaseous planets that inherited significant amount of their mass from the protoplanetary disk, which is in turn expected to have compositional similarities to the host star. Figure 15.6 shows the metallicity of Solar System ice and gas giants, exoplanets, and brown dwarfs as a function of planet mass. There is an intriguing trend of

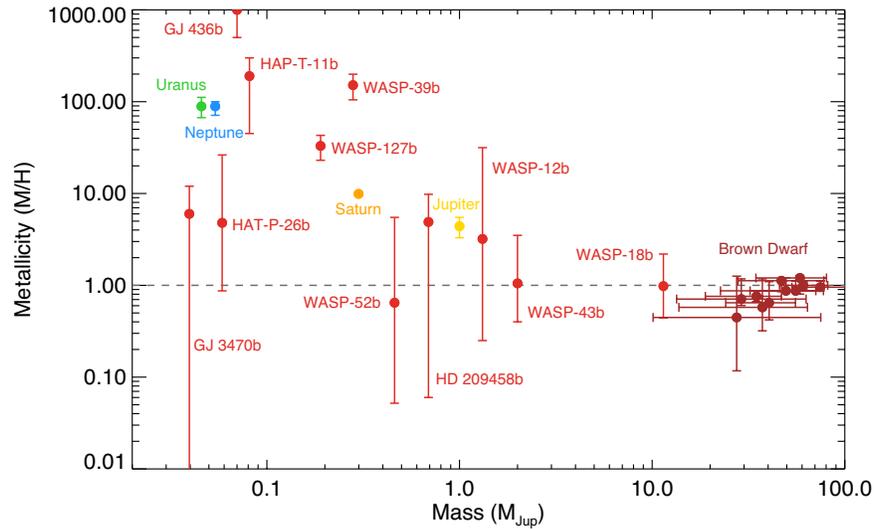


Figure 15.6: Metallicity of Solar System (colors), exoplanets (red), and brown dwarf (maroon) atmospheres as a function of mass. A metallicity of 1 corresponds to the metallicity of our Sun. The Solar System objects show decreasing metallicity with increasing mass, which serves as a first-order expectation for exoplanets that is yet to be discerned. Figure adapted from Zhang (2020).

decreasing metallicity with increasing planet mass for the Solar System gas and ice giants, with Uranus and Neptune having metallicities nearly 100x Solar, Saturn having a metallicity approximately 10x Solar, and Jupiter having a metallicity around 3x Solar. Meanwhile, brown dwarfs (which are expected to form similar to stars) all have low metallicity, implying that their atmospheres do not have an ice or refractory component significantly enhanced from stars. At this moment, exoplanets are effectively a scatter plot in mass-metallicity space – current and future work with JWST and ground-based telescopes to measure the metallicity of exoplanet atmospheres may better discern whether the mass-metallicity “trend” from our Solar System extends outward to exoplanets.

The second way to measure bulk composition is to consider the individual bulk elemental composition. This is normally done by measuring elemental ratios, with the most common being the C/O ratio, given that it is expected to be linked to the relative gaseous vs. refractory composition of the protoplanetary disk (Öberg et al., 2011, see Figure 15.7). Figure

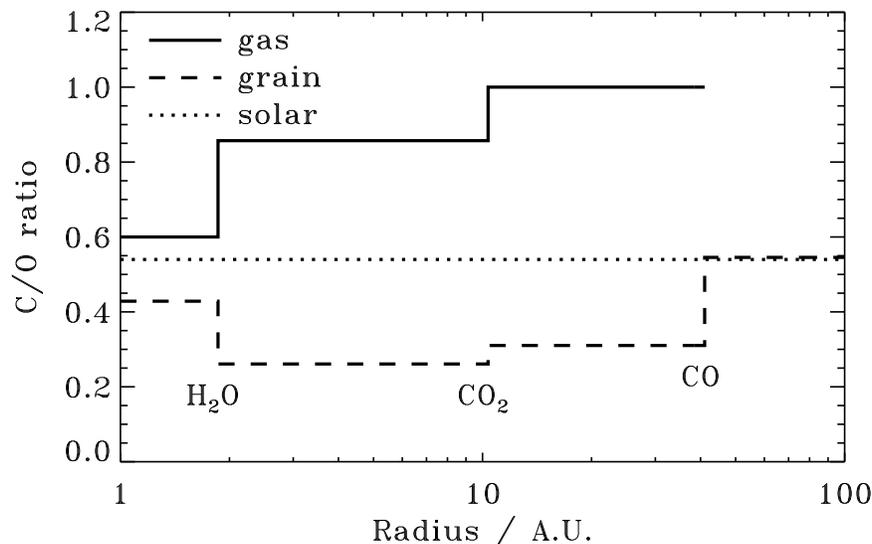


Figure 15.7: C/O ratio of gas (solid) grains (dashed) relative to the Solar C/O ratio (dotted) in a simple model for a protoplanetary disk. The C/O ratio of gas is enriched due to condensation of H₂O and CO₂ at the water and carbon dioxide ice lines, which in turn depletes the C/O ratio of grains that form rocky cores. Figure adapted from Öberg et al. (2011).

15.8 shows a ternary diagram of atmospheric composition as a function of the relative abundance of hydrogen, oxygen, and carbon. The end-members of atmospheric composition are standard gas giant atmospheres (H₂), oxygen-dominated atmospheres, and graphite/carbon monoxide-dominated atmospheres. Generally, planets are expected to have intermediate C/O ratios, with the specific combination of C/O ratio and metallicity determining whether the object is CO or CO₂ dominated.

For rocky planets, the atmospheric abundances are set by the competition of atmospheric loss and outgassing from the interior (e.g., via volcanism). Most Earth-like rocky planets are expected to have lost a primary envelope of H accreted from the protoplanetary disk, and thus their atmospheres are “secondary,” and obtained via outgassing from the solid interior. Thus, the bulk composition of rocky planets is expected to roughly relate to their atmospheric composition, with fractionation from outgassing leading to an atmosphere comprised of volatile species.

15.5.2 Equilibrium chemistry

If an atmosphere is in thermochemical equilibrium, the temperature, pressure, and metallicity (i.e., bulk composition) alone set the abundance distribution of each individual chemical species within the atmosphere. Following Visscher & Moses (2011), we consider a simple

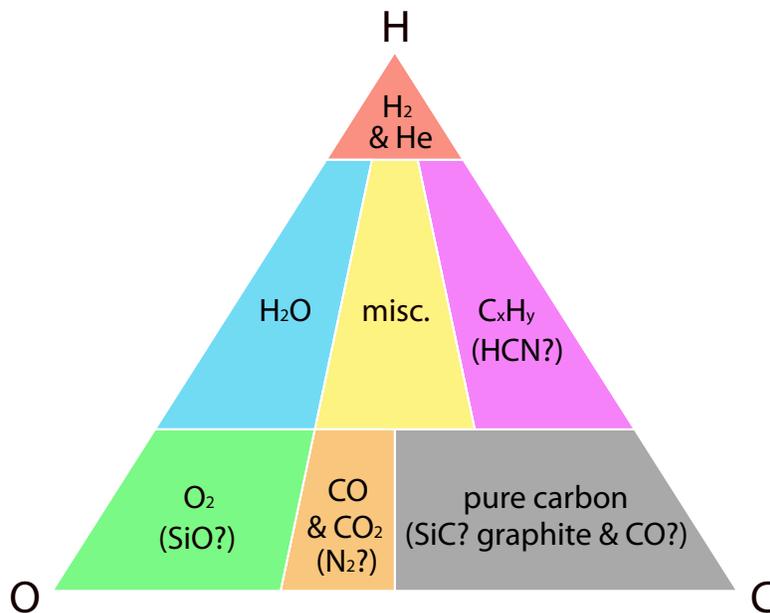


Figure 15.8: A ternary diagram of atmospheric compositions in H-C-O space. Gas giants lie on the top, low C/O planets on the bottom left, and high C/O planets on the bottom right. Figure adapted from Zhang (2020).

balanced gas-phase reaction



where the lowercase letter represents the number of molecules of the uppercase letter. The equilibrium constant of this reaction can be expressed as

$$K_{\text{eq}} = \frac{[C]^c [D]^d}{[A]^a [B]^b}, \quad (15.38)$$

where brackets represent the number density of each species. This equilibrium constant can alternately be written in terms of partial pressures p ,

$$K_p = \frac{p_C^c p_D^d}{p_A^a p_B^b}. \quad (15.39)$$

The equilibrium constant is then related to the standard-state Gibbs free energy change $\Delta_r G^\circ = \Delta_f G^\circ(\text{products}) - \Delta_f G^\circ(\text{reactants})$ and temperature as

$$K_p = \exp\left(-\frac{\Delta_r G^\circ}{RT}\right). \quad (15.40)$$

Thus, to determine the relative abundances of species for a net chemical reaction at a given temperature and pressure, one only needs to obtain the change in Gibbs free energy between products and reactants. The standard-state Gibbs free energy change for formation can be calculated from the standard state enthalpy $\Delta_f H^\circ$ and entropy $\Delta_f S^\circ$ as

$$\Delta_f G^\circ = \Delta_f H^\circ - T\Delta_f S^\circ, \quad (15.41)$$

where the standard state enthalpy of formation and entropy are compiled in standard chemistry reference databases, e.g., <https://webbook.nist.gov/chemistry/>.

Figure 15.5 shows the equilibrium chemistry equivalency curves of methane-carbon monoxide and ammonia-nitrogen, which individually summarize the thermochemical carbon and nitrogen cycles. The net thermochemical carbon cycle can be written as



where CH_4 has a higher abundance at cooler temperatures and higher pressures and CO has a higher abundance at higher temperatures and lower pressures. The net thermochemical nitrogen cycle is



with NH_3 having a higher abundance at low temperatures and high pressures and N_2 having a higher abundance at high temperatures and low pressures.

15.5.3 Disequilibrium chemistry and mixing

Chemical species can be mixed by fluid motions, causing the abundance of a given species to be different than the expected from considerations of thermochemical equilibrium. Such “disequilibrium” states occur when the typical mixing timescale τ_{mix} is shorter than the chemical timescale τ_{chem}

$$\begin{aligned} \tau_{\text{mix}} \ll \tau_{\text{chem}} & \text{ disequilibrium,} \\ \tau_{\text{mix}} \gg \tau_{\text{chem}} & \text{ equilibrium.} \end{aligned} \quad (15.44)$$

Though large-scale vertical motions are not diffusive (rather, they are the combination of many small-scale advective motions), in order to concoct a one-dimensional picture of mixing we can define a vertical diffusion coefficient, often termed K_{zz} . K_{zz} is analogous to the effective viscosity from our discussion of protoplanetary disks, but here limited to vertical transport alone (hence the subscript – it’s really one component of a larger diffusion tensor). We can relate the vertical mixing timescale to K_{zz} as

$$\tau_{\text{mix}} \sim \frac{H^2}{K_{zz}}, \quad (15.45)$$

where $H = RT/g$ is the (isothermal) pressure scale height. The eddy diffusivity is expected to increase with decreasing pressure in planetary atmospheres due to the combination of increased radiative forcing and higher wave amplitude at lower pressures (see Figure 15.9). For radiatively dominated isothermal atmospheres, $K_{zz} \propto p^{-1/2}$. This implies that mixing timescales will decrease with decreasing pressure in most atmospheres. Meanwhile, the chemical timescales in planetary atmospheres are expected to increase with decreasing pressure (at least in the troposphere), as chemical reaction rates drop with decreasing temperature. As a result, there is expected to be “quench point” in planetary atmospheres where $\tau_{\text{mix}} = \tau_{\text{chem}}$, and at lower pressures (higher altitudes) than the quench point $\tau_{\text{mix}} < \tau_{\text{chem}}$. Thus, the quench point is the location in the atmosphere where above the atmosphere is in a state of chemical disequilibrium, and below it is in a state of chemical equilibrium. This quench point is species-dependent, as it depends itself on the chemical timescale, which depends

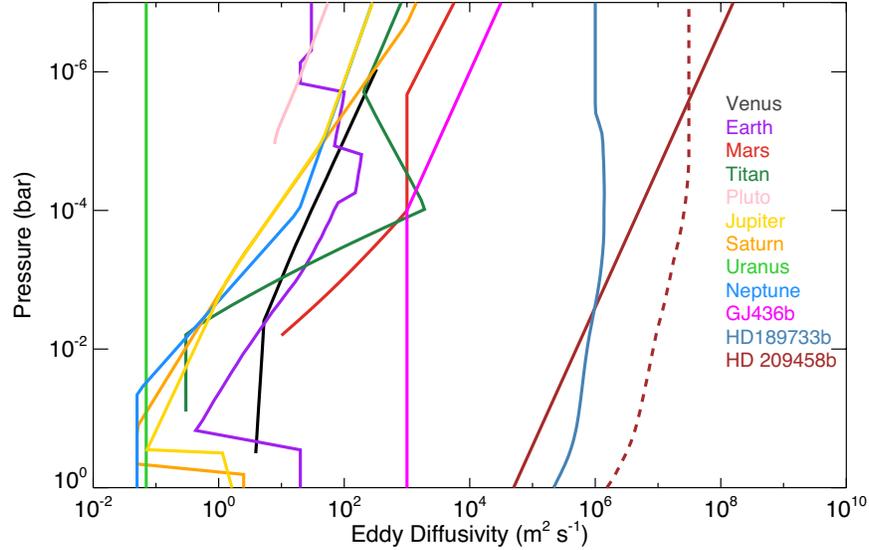


Figure 15.9: Eddy diffusion coefficient K_{zz} for various Solar System objects as well as exoplanets. Figure adapted from Zhang (2020).

on the chemical reaction of interest. Figure 15.10 shows example abundance profiles from a disequilibrium chemistry calculation for the hot Jupiter HD 189733b with the VULCAN code, displaying how species quench at different locations at depth (compare the solid lines to the dotted lines). Note that the process driving chemical disequilibrium will change at pressures lower than the quench point for irradiated atmospheres, as photochemistry becomes the dominant disequilibrium process rather than mixing at low pressures where UV radiation is absorbed (e.g., in Earth’s stratosphere).

15.6 Atmospheric loss

15.6.1 Energetic considerations

Atmospheres can be lost to space through two main types of mechanisms: thermal and non-thermal processes. Thermal atmospheric escape corresponds to cases where the upper atmospheric temperature is high enough that the thermal velocity of the gas approximately exceeds the escape velocity of the planet (see this problem set for a more accurate estimate), implying that the gas is not gravitationally bound to the planet. The potential for thermal escape can be described by the Jeans parameter

$$\lambda = \frac{E_{\text{grav}}}{E_{\text{therm}}} = \frac{GM_p \mu m_p}{k_B T R_p}, \quad (15.46)$$

which is the ratio of gravitational to thermal energy in the upper atmosphere of a planet. Non-thermal escape occurs through processes that are not related to the temperature of the gas, usually related instead to electrical interactions such as stellar wind interactions with ions. Generally, non-thermal escape is not expected to cause total atmospheric loss (unless the host star is very highly active), while thermal escape is able to completely remove atmospheric envelopes especially for close-in, hot, low-mass planets.

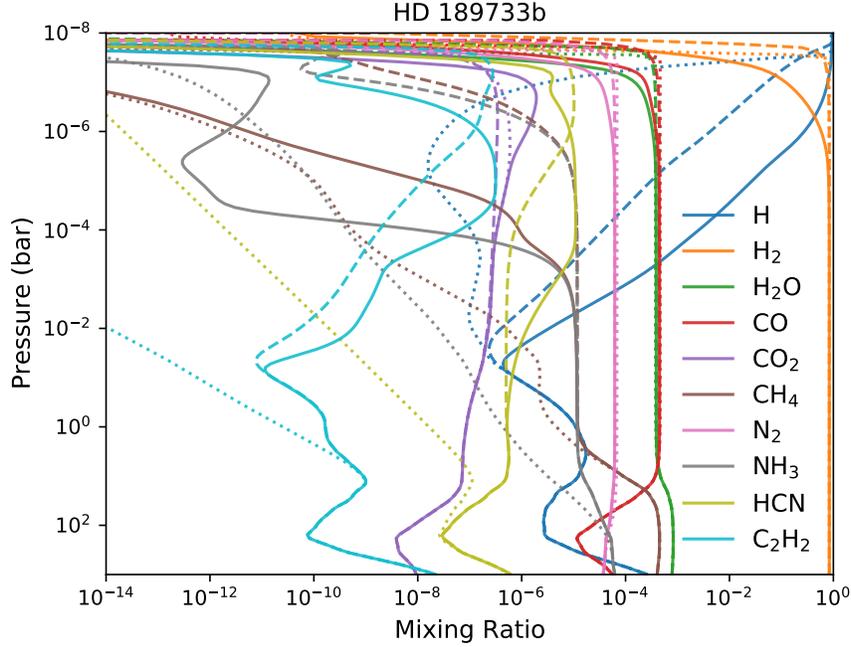


Figure 15.10: Predicted chemical profiles from full chemistry models (solid lines), models without photochemistry (dashed lines), and models assuming equilibrium chemistry (dotted lines) for HD 189733b. Figure adapted from Tsai et al. (2021).

Within thermal escape, there are two further categories: supply-limited and energy-limited escape. Supply-limited escape is escape that is limited by the supply of a low mean molecular weight species, for example hydrogen, and is what regulated the escape of Earth’s primary hydrogen atmosphere. In supply-limited escape, molecules such as water are broken down into atoms by chemical reactions (e.g., photodissociation driven by the host star’s UV radiation), and the lighter atoms segregate to lower pressures due to their larger scale height in the upper atmosphere, causing loss to space. What sets the eventual rate of escape is then the transport of these light atoms to high altitudes from where they can then be lost to space. Energy-limited escape is conversely when the energy available regulates the escape rate. The energy-limited mass loss rate can be roughly estimated as

$$\dot{M} \sim \eta \frac{L_{\text{XUV}} R_p^3}{4GM_p a^2}, \quad (15.47)$$

where L_{XUV} is high-energy portion of the stellar luminosity (with XUV the dominant region of the spectrum for high-energy photons) and η is the efficiency of mass loss through a hydrodynamic wind, which must be calculated from numerical simulations and is $\eta \sim 0.1 - 0.2$. Energy-limited mass loss is larger for planets around more active stars (higher L_{XUV}) with closer in orbits (lower a), and lower masses and larger radii (and thus slower escape velocities).

15.6.2 The cosmic shoreline

A foundational observation was made by Zahnle & Catling (2017) (along with the authors' preceding work) that there appears to be a “cosmic shoreline” in irradiation-escape velocity space, where planets with high instellation and low escape velocities do not have significant atmospheres while planets with high escape velocities and low instellation hold onto thick atmospheres. An updated version of their empirical observation is shown in Figure 15.11, where the cyan line represents the $I \propto v_e^4$ curve that marks the shoreline. Though

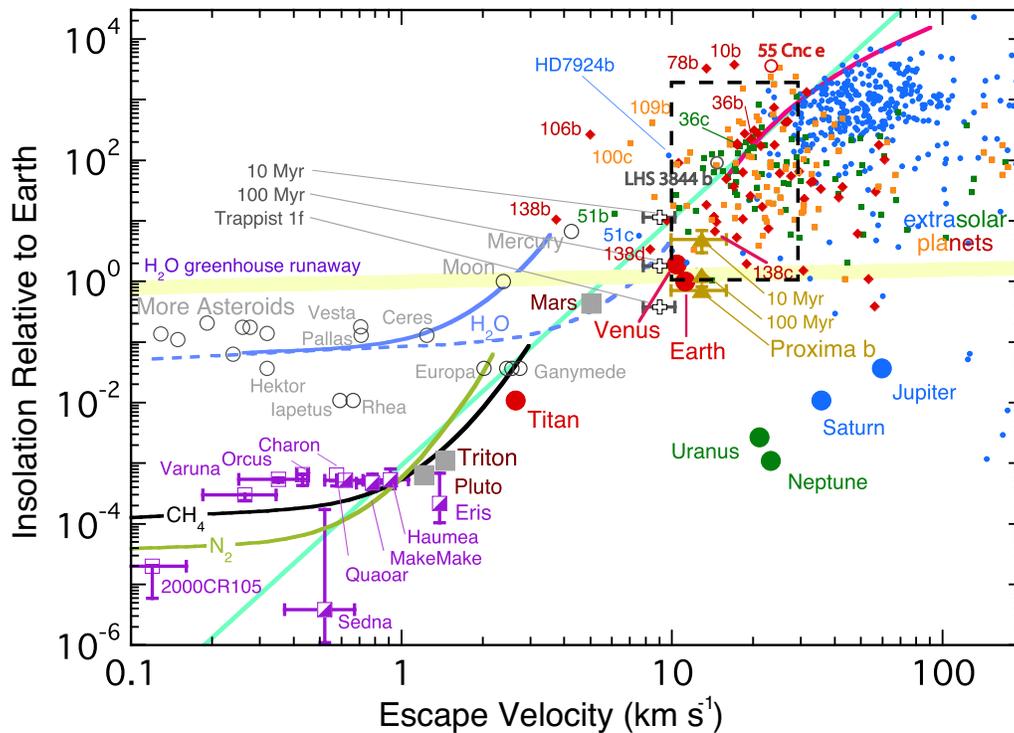


Figure 15.11: The “cosmic shoreline,” i.e. the empirically determined level of instellation that planets below a given escape velocity cannot hold onto thick atmospheres. Shown is a scatterplot of the level of instellation and escape velocity of various planets and moons in our Solar System and beyond, compared with the empirical cosmic shoreline (cyan) and hydrodynamic thermal (energy-limited) escape curves for various species (methane, nitrogen, water – black, gold, blue colored lines). Figure adapted from Zhang (2020), in turn modified from Zahnle & Catling (2017).

this shoreline nicely describes the observed prevalence of atmospheres, it does not correspond one-to-one with existing theoretical predictions. Energy-limited mass-loss would imply that the shoreline would scale as $I \propto v_e^3 \sqrt{\rho}$, too shallow to explain the empirical relationship. Additionally, supply-limited loss cannot alone explain the trend. As a result, Zahnle & Catling (2017) proposed that impact erosion (the loss of atmospheres through energy released via impacts from comets and asteroids) could potentially shape the cosmic shoreline. Observations of the presence/absence of atmospheres on exoplanets (through the transmission and emission spectroscopy methods that we’ll discuss in two weeks) are likely required to provide a firm empirical basis on which to test our fundamental understanding of atmospheric loss.

16 Exoplanet interiors: giant planets

Our agenda for Day 17 is the following:

1. Phase diagram of hydrogen, structure of our gas and ice giants (15 minutes)
2. Hydrostatic equilibrium (for the third time!), central pressure activity (40 minutes)
3. Equations of planetary structure, energy transport and Schwarzschild criterion (20 minutes)

Today's reading is the Fortney Giant Planet Interior Structure and Thermal Evolution review paper. This will provide a comprehensive overview of the current study of the interiors of both Jupiter and Saturn as well as gas giant exoplanets.

16.1 Phases of H/He in giant planets

Hydrogen lies in two main phases in the interiors of gas giant planets: molecular hydrogen (H_2) in the outer envelope and atmosphere, and metallic hydrogen (H^+) in the interior. Metallic hydrogen forms due to pressure ionization of molecular hydrogen, which turns hydrogen into a dense lattice of protons, with a distance between protons in the lattice that is the same as the distance between protons in a hydrogen molecule. The electrons in this lattice are delocalized, thus causing metallic hydrogen to have high thermal and electrical conductivities. As shown in Figure 16.1, metallic hydrogen is expected to form at high densities of $\sim 1 \text{ g cm}^{-3}$, corresponding to pressures between approximately 0.1 – 3 Mbar. As

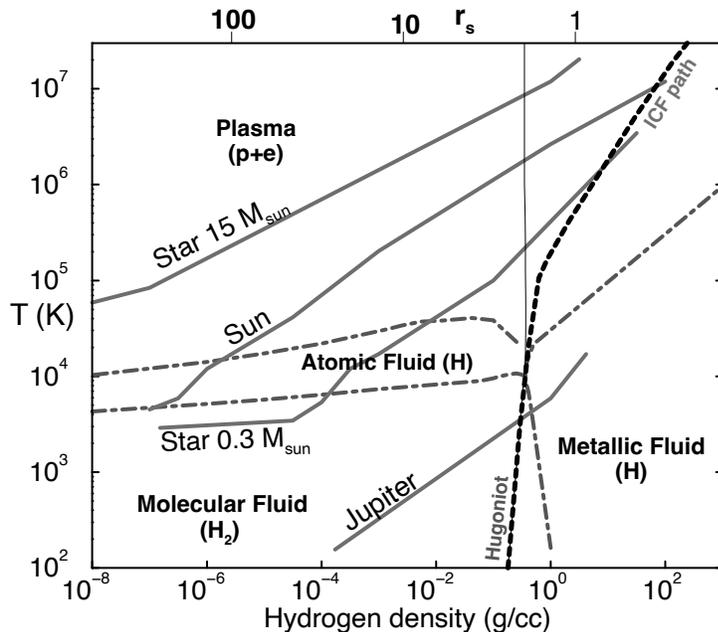


Figure 16.1: Phase diagram of hydrogen, showing the molecular, atomic, metallic, and plasma regimes. Over-plotted in gray are temperature-density profiles of Jupiter and stars with $M = 0.3, 1, 15 M_{\odot}$.

a result, both Jupiter and Saturn are expected to have outer layers of molecular H_2 , with a transition in the deep interior to metallic hydrogen. Note that even hotter gas giant planets (e.g., some hot Jupiters) can have a transition from molecular hydrogen to hydrogen plasma in the interior due to temperature-driven (rather than pressure-driven) ionization, which is

governed by the Saha equation and is what sets the ionization state of stars. In both cases, the ionization (pressure or temperature driven) in the interiors of gas giant planets is high enough that the ratio of the electrostatic potential energy to the thermal energy

$$\Gamma = \frac{E_{\text{coul}}}{E_{\text{th}}} = \frac{e^2}{dk_B T} \sim 1, \quad (16.1)$$

where for an ideal gas $\Gamma \approx 0$. This implies that the equation of state of the interiors of gas giant planets is very far from an ideal gas, and pressure is no longer a linear function of density and temperature.

Helium is more difficult to ionize than hydrogen because it has two electrons rather than one. As a result, He is neutral until pressures of $\gtrsim 50$ Mbar, causing it to never transition to a metallic or plasma form in typical giant planet interiors. Though helium does not transition, when hydrogen undergoes its molecular to metallic transition the hydrogen and helium fluid together undergoes a phase change from being well-mixed and homogeneous at low pressures (where hydrogen is molecular) to being demixed at high pressures (where hydrogen is metallic). Figure 16.2 shows this demixing boundary both in temperature-pressure space as well as physically within Jupiter's interior. The demixing of hydrogen and

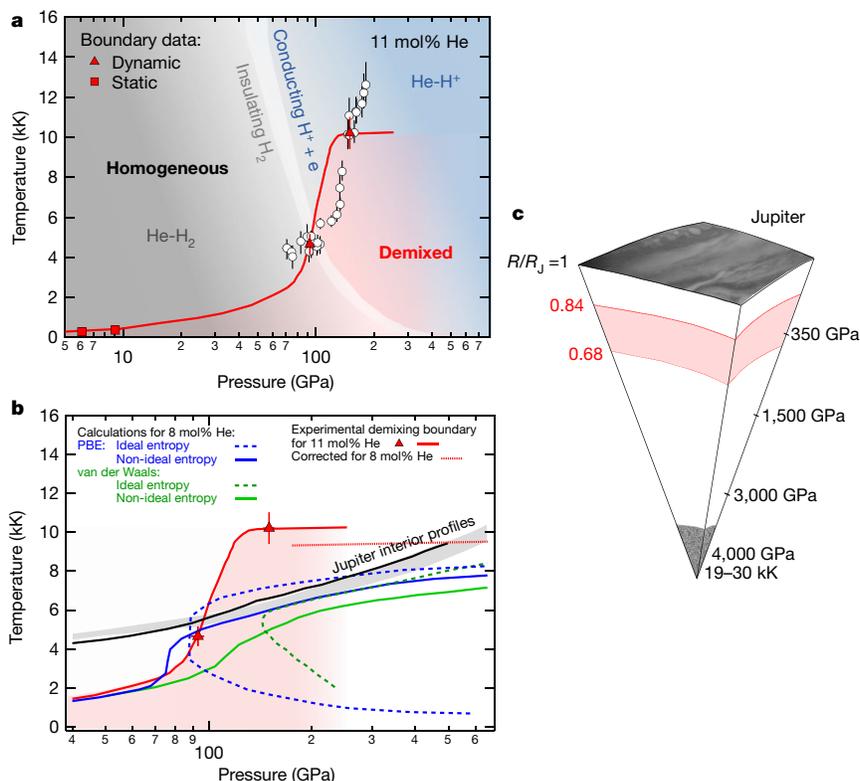


Figure 16.2: Experimental data for the phase diagram of the H-He mixture (a), comparison of experimental data (red) with both previous calculations (green, blue lines) and the Jupiter interior temperature profile (black) (b), and inferred regime of H-He demixing in Jupiter's structure (c). Figure adapted from Brygoo et al. (2021).

helium at ~ 1 Mbar in Jupiter's interior causes helium to phase separate from hydrogen and

“rain” out, falling deeper within the planet. This then causes a compositional gradient in the interior of Jupiter (and Saturn) which has been demonstrated to impact their evolution.

16.2 Interior structures of Solar System giant planets

The fundamental interior structures of Jupiter and Saturn are broadly similar, as shown in Figure 16.3. Both planets have exterior envelopes of molecular hydrogen that transition to

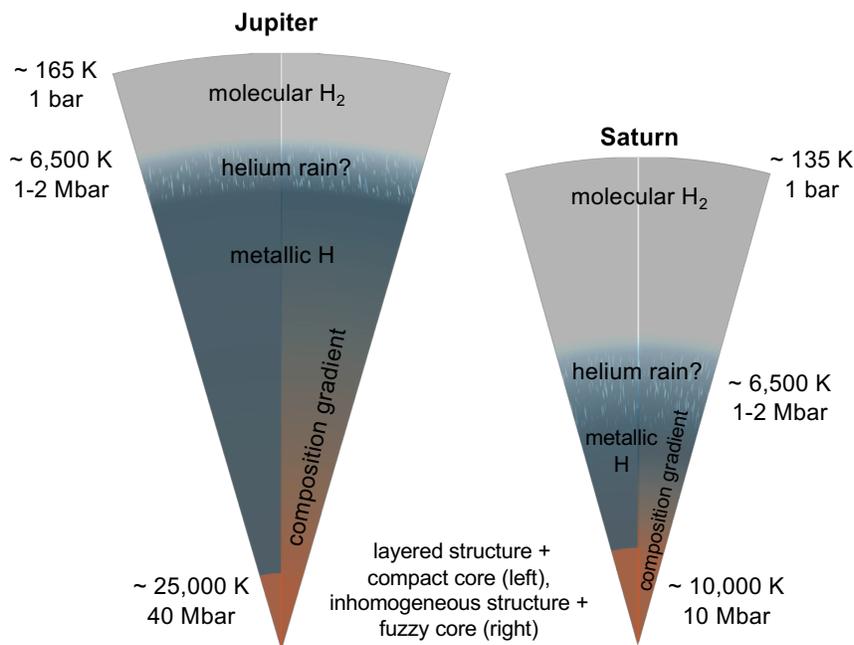


Figure 16.3: Pie slice view of Jupiter and Saturn, showing layers of molecular H₂, helium demixing, metallic hydrogen, and possible structures of the deep interior. Figure courtesy Jonathan Fortney.

metallic hydrogen at pressures of ~ 1 Mbar, with H-He demixing at this region that causes a compositional gradient. From matching internal structure models to precision gravity measurements of the Juno (for Jupiter) and Cassini grand finale (for Saturn) missions, both planets are inferred to have extended heavy element “cores,” which are an extended mix of metal, rock, and H/He. This differs somewhat to the typical expectation from core accretion, where a planet forms a (solid, singular) core and then accretes gas onto of the core from the protoplanetary disk. Three possibilities for this diluted core are: 1) that it occurred as a product of planetesimal accretion during formation, as planetesimals burned up in the envelope before reaching the core; 2) that the core was “dredged” up by convective motions into the envelope; 3) that Jupiter underwent a giant impact early in its evolution that destroyed the core and mixed it upward into the envelope. The fact that both Jupiter and Saturn show evidence for diffuse cores implies that the mechanism may be ubiquitous and thus linked to formation.

Compared to Jupiter and Saturn, relatively little is known about the interiors of Uranus and Neptune. This is because both planets have only been studied by a single flyby mission, Voyager 2, while both Jupiter and Saturn have been characterized in detail by orbiters.

Figure 16.4 shows schematics of the possible interior structures of Uranus and Neptune, where our lack of detailed information prevents a detailed picture of the differences in internal structure between the two planets. In general, models predict that the outer layers of both

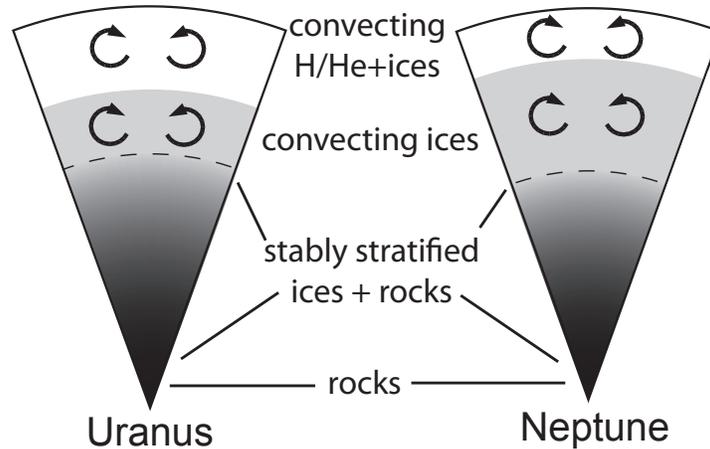


Figure 16.4: Pie slice views of the possible interior structures of Uranus and Neptune, from Fortney et al. (2010). These planets have an envelope dominated by H/He overlying an icy layer of volatiles, which then overlies a deep interior of ice, rock, and metals.

Uranus and Neptune is dominantly (molecular) hydrogen and helium. This H/He envelope then overlies a fluid layer of astrophysical ices (e.g., water, methane), which then overlies an interior of ices and heavy elements. However, it is not known what the mass fraction of various ices and rock/metal is in the interiors of Uranus and Neptune. Additionally, though these planets are often termed “ice giants” due to their composition, the physical state of these high pressure ices are a hot, dense, partially ionized (and conductive) fluid. At the interior conditions of Uranus and Neptune, water and hydrogen are miscible, and thus can be treated as one fluid (Soubiran & Militzer, 2015). Similarly, at high pressure in the interiors of Uranus and Neptune rock and ices may be similarly miscible, leading to a diffuse heavy element interior similar to Jupiter and Saturn. Detailed gravity observations of Uranus and Neptune with orbiters are required to better constrain their interior structure.

16.3 Hydrostatic equilibrium

Recall that hydrostatic equilibrium can be expressed as

$$\frac{dp}{dr} = -\rho g = -\rho \frac{Gm}{r^2}, \quad (16.2)$$

where r is radius from the center of the planet, and m is the mass coordinate (i.e., the enclosed mass). Note that the mass of any given shell of the planet can be related to the density and radius as $dm = 4\pi r^2 \rho dr$. Substituting in for dr , we can express hydrostatic equilibrium in terms of dp/dm as

$$\frac{dp}{dm} = -\frac{Gm}{4\pi r^4}. \quad (16.3)$$

As a result, the pressure must increase as the mass coordinate decreases, going toward the center of the planet. We can now use this to determine the pressure at the center of any given planet, as we'll do in our activity.

16.3.1 Central pressure activity

Use the previously derived expression of hydrostatic equilibrium in mass coordinates to solve the following problems in groups of 2-3.

1. Derive an approximate expression for the pressure at the center of a planet as a function of its mass and radius. To do so, integrate Equation (16.3) from the center to the surface of the planet, assuming constant density and that the surface of the planet is at zero pressure.
2. Use your expression to estimate the pressure at the center of Jupiter, in Mbar (where 1 bar = 10^5 Pa). Compare this to the value shown in Figure 16.3. Describe why your estimate might be different than the exact value that includes density variations with mass coordinate.
3. Given that hydrogen metallizes at pressures $\gtrsim 1$ Mbar, determine whether hydrogen changes phase within Neptune's interior. Note that Neptune has a magnetic field driven by an internal dynamo – ionization of what species could lead to a dynamo that generates Neptune's magnetic field?

16.4 Equations of planetary structure

Five equations fully represent the internal structure of a planet – these are the equations of planetary structure, which are equivalent to the equations of stellar structure but without including nuclear burning. The first two of these have previously been introduced this chapter, and are the equations of mass conservation

$$\frac{dm}{dr} = 4\pi r^2 \rho, \quad (16.4)$$

and hydrostatic equilibrium

$$\frac{dp}{dm} = -\frac{Gm}{4\pi r^4}. \quad (16.5)$$

The third equation of planetary structure is a statement of energy conservation,

$$\frac{dL}{dm} = \epsilon_{\text{grav}} = -T \frac{dS}{dt}, \quad (16.6)$$

where L is the outgoing luminosity at mass coordinate m and S is the entropy, the loss of which drives gravitational cooling ϵ_{grav} and contraction. The fourth equation describes energy transport

$$\frac{dT}{dm} = -\frac{GmT}{4\pi r^2 p} \nabla, \quad (16.7)$$

where $\nabla = d\ln T/d\ln p$ is the logarithmic temperature gradient (equivalent to lapse rate from our discussion of atmospheres). We'll discuss what sets ∇ in the following section.

The final equation of planetary structure is the equation of state, which relates pressure, density, and temperature and depends on the composition of the planet. The equation of state for gas giant planets is non-trivial, because due to the high densities in the interior of the planet quantum mechanics must be taken into account due to the uncertainty principle and Pauli exclusion principle. The pressure in the interior of a giant planet is the sum of the degeneracy pressure and thermal pressure, i.e., $p = p_{\text{deg}} + p_{\text{th}}$. In general, the ratio of thermal pressure to total pressure (at high pressures relevant to the deep interiors of gas giants) is $p_{\text{th}}/p \approx 1 \times 10^{-5}T$. Thus, at typical gas giant central temperatures of 10^4 K, the contribution of thermal pressure to the total pressure is only $\sim 10\%$ – meaning that these objects are highly degenerate! In practice, equations of state for giant planet interiors are tabulated based on numerical quantum mechanics simulations that are benchmarked with high-pressure experiments using either diamond anvil cells or laser compression to reach \sim Mbar pressures.

16.4.1 Heat transport in planetary interiors

The equation for energy transport in planetary interiors can alternately be written as a change in temperature with radius

$$\frac{dT}{dr} = \frac{dp}{dr} \frac{T}{P} \nabla. \quad (16.8)$$

There are two primary ways heat can be transported through planetary interiors: radiation and convection. Radiative energy transport is determined by the rate of diffusion of photons in a random-walk process. The mean free path of photons depends on the number density n and cross section σ , or equivalently mass density ρ and opacity κ , as

$$\lambda = \frac{1}{n\sigma} = \frac{1}{\rho\kappa}. \quad (16.9)$$

The radiative lapse rate in a planetary interior can be related to the opacity, luminosity, pressure, mass coordinate, and temperature as

$$\nabla_{\text{rad}} = \frac{3}{64\pi\sigma G} \frac{\kappa L p}{m T^4}. \quad (16.10)$$

Conversely, the adiabatic lapse rate ∇_{ad} is set by the thermodynamic properties of the planet, and in the envelope (as for an atmosphere) it is $\nabla_{\text{ad}} = R/c_p$.

Whether heat transport and the resulting lapse rate is set by convection or radiation can be determined by the Schwarzschild criterion, which sets the temperature gradient to the smaller of the adiabatic gradient ∇_{ad} or the radiative gradient ∇_{rad} . We can express this as

$$\begin{aligned} \nabla_{\text{ad}} < \nabla_{\text{rad}} & \text{ convection,} \\ \nabla_{\text{ad}} > \nabla_{\text{rad}} & \text{ radiation.} \end{aligned} \quad (16.11)$$

Because the radiative gradient increases with pressure and opacity, generally planets transition from having radiative exteriors to having convective interiors. In some cases, radiative “windows” (Guillot et al., 1995) appear at shallow regions of the otherwise convective envelope due to sharp decreases in the opacity due to compositional variations or changes in the temperature and outgoing luminosity.

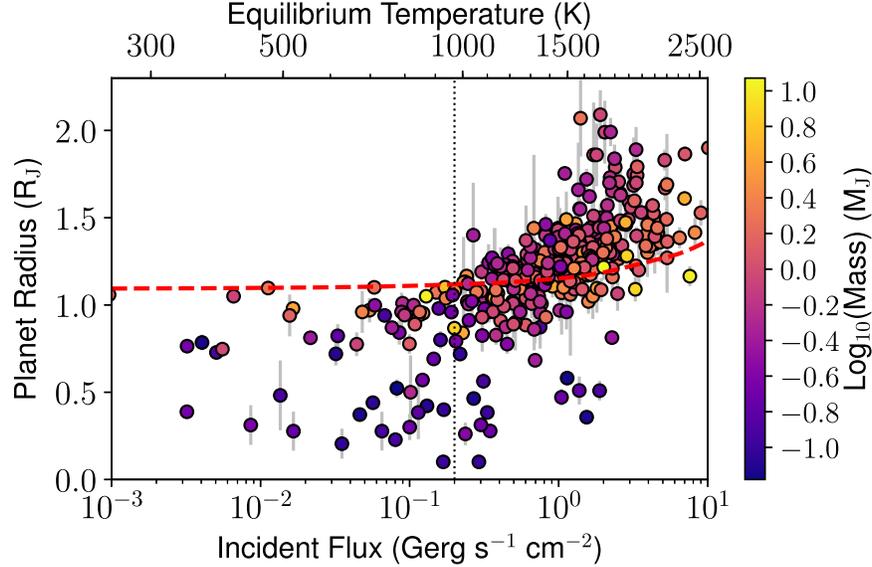


Figure 16.5: The radii of warm and hot Jupiters as a function of the amount of incident flux they receive, with points colored by planet mass. Notably, hot Jupiters with $T_{\text{eq}} \gtrsim 1000$ K have radii that can be larger than standard evolutionary model predictions (red dashed line).

16.4.2 Radius inflation of hot Jupiters

Many hot Jupiters have radii larger than expected from solving the equations of planetary structure, even including an additional atmospheric heating term due to the instellation that the planet receives. Figure 16.5 demonstrates this issue of “radius inflation,” where hot Jupiters with $T_{\text{eq}} \gtrsim 1000$ K have radii that can be larger than standard solutions of the equations of planetary structure, while warm Jupiters always have radii at or below the expected curve. The observed planets with radii smaller than the curve can be explained by adding additional heavy elements into the interior, which increases the bulk density and thus reduces the radius for a given mass. However, explaining the radii above the standard model curve requires some additional physical mechanisms, the causes of which are still active areas of research (for a recent review, see Fortney et al., 2021).

The radius inflation mechanisms can be broken down into two main categories: mechanisms that directly slow the cooling rate of the planet, and mechanisms that offset the cooling of the planet by adding additional heat into the planetary interior. Because (as shown in Figure 16.5) the observed radii of hot Jupiters appear to correlate with the level of incident flux they receive, it is expected that the mechanism(s) that cause the bloated radii of hot Jupiters are linked to the incident flux, with some fraction of the incident stellar power $\gamma = \Gamma/L_{\star}$ being deposited in the interior of the planet. Possible sources of this deep heat include tidal dissipation (Bodenheimer et al., 2001), Ohmic dissipation (Batygin & Stevenson, 2010), or atmospheric winds (Guillot & Showman, 2002). To date, there is no “smoking gun” for which mechanism causes radius inflation. The most powerful statistical result is that the heating efficiency γ peaks at intermediate values of T_{eq} and decreases toward both warm and ultra-hot Jupiters (Thorngren & Fortney, 2018; Sarkis et al., 2021). This implies that the mechanism that causes radius inflation is self-limiting at high temperatures, which aligns

well with mechanisms that are regulated by the feedback of magnetic fields onto motions in the planetary envelope. This is because hotter planets will have more conductive envelopes, which then interact more strongly with planetary magnetic fields, leading to Lorentz forces which generally act against the flow and limit the level of dissipation (Ohmic or mechanical). However, future work studying the timescale of inflation during both the main-sequence and post-main-sequence is required to conclusively identify the mechanism leading to radius inflation (or determine if there are multiple mechanisms).

16.4.3 Radius evolution, Kelvin-Helmholtz Timescale

Recently-formed, hot, young planets have their cooling dominated by gravitational energy loss, with a cooling luminosity of $L \sim -dE_g/dt$. We can scale this expression to derive a characteristic Kelvin-Helmholtz (thermal) timescale

$$\tau_{\text{KH}} \sim \frac{E_g}{L} \sim \frac{GM_p^2}{R_p L}. \quad (16.12)$$

For a typical Jupiter-mass giant planet, the initial Kelvin-Helmholtz timescale is ~ 10 Myr. This implies that the present-day “inflated” radii of hot Jupiters is not solely the consequence of age, as most hot Jupiters are found around main-sequence (Gyr-old) stars. Instead, some process must halt or slow the cooling of the planet, keeping radii large out to late times. Figure 16.6 shows radius evolution curves from the planetary structure model predictions of Komacek & Youdin (2017) for an HD 209458b-like planet with varying depths of deposited heat. The cases with shallow heating show perpetual cooling over time, demonstrating that

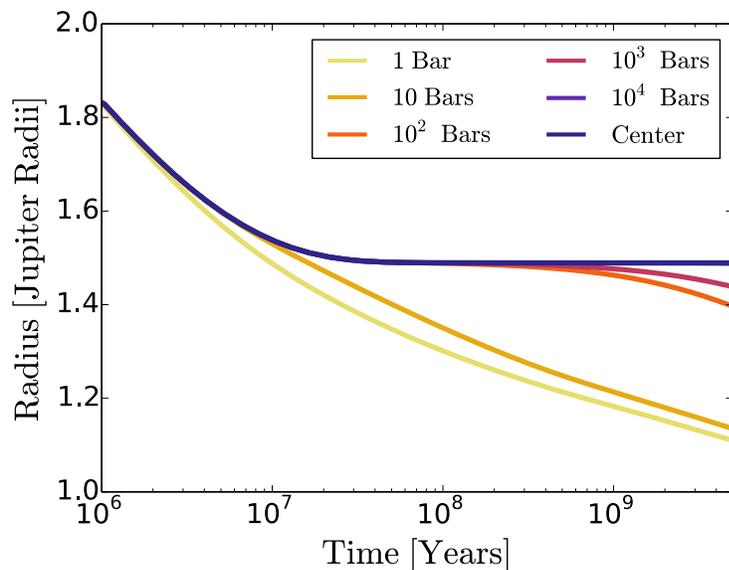


Figure 16.6: Radius evolution for hot Jupiters undergoing deposited heating centered at different depths (from 1 bar to the center), with a fixed heating efficiency (fraction of incident stellar flux converted to heat) of $\gamma = 1\%$. Figure modified from Komacek & Youdin (2017).

continued Kelvin-Helmholtz contraction cannot explain the radii of inflated hot Jupiters. Meanwhile, central heating leads to an equilibrium state where the deposited heating replaces the cooling of the interior, leading to zero net change in the central temperature (and thus radius) with time. As a result, the mechanism that causes the radius inflation of hot Jupiters must be linked toward deep deposition of heat into the planet interior, at least sufficiently deep to slow cooling over Gyrs of evolution.

17 Planetary habitability

Our agenda for our special lecture on habitability is the following:

1. Overview of the habitable zone concept (25 minutes)
2. Biosignatures: oxygen, ozone, chemical disequilibrium (15 minutes)
3. Discussion: which technique would you use to search for biosignatures? (15 minutes)
4. Decadal survey, prediction activity (20 minutes)

Today's reading is a review by Meadows et al., which describes how oxygen can be used as a biosignature, as well as potential false positives for oxygen biosignatures and thus the need to characterize planetary environments in order to discern if an observed biosignature is linked to life. Note that the material in today's class will *not* be covered on the last mid-term.

17.1 The habitable zone

17.1.1 Classic 1D framework, carbonate-silicate weathering

The habitable zone is the region around any given star at which water can reside in liquid form at the surface of a planet that is (roughly) equivalent to Earth in its mass, radius, atmospheric composition, and atmospheric surface pressure. The habitable zone is often collapsed to only be a function of the host star type and instellation (i.e., incident stellar flux) onto the top-of-atmosphere of its companion planet. However, in reality the habitable zone is multi-dimensional, as it critically depends on the age of the system (because Sun-like stars brighten and M-dwarf stars dim over time) and thus the evolutionary history of the planet, along with perturbations of planetary parameters (e.g., mass, radius, atmospheric composition) slightly away from modern Earth values.

The classic model of the habitable zone was first developed by Kasting et al. (1993). The critical improvement of the Kasting model over previous approaches is that Kasting took into account the impact of the carbonate-silicate weathering feedback on climate evolution. Figure 17.1 shows a schematic of this process, which occurs on all Earth-like planets with active (plate) tectonics, surface liquid water (and thus rain and oceans), and silicate rock. The carbonate-silicate cycle begins with the weathering of exposed calcium/magnesium bearing silicate rock ((Ca,Mg)SiO₃) by rain, which causes a chemical reaction by which CO₂ is removed from the atmosphere, producing calcium bicarbonate ions (see the reaction under “Land” in Figure 17.1). Then, the calcium and bicarbonate ions are transported (e.g., by flowing water) to the ocean, and organisms in the ocean use these ions to make calcium carbonate (CaCO₃, see the reaction under “Ocean” in Figure 17.1) a fraction of which is then deposited on the seafloor after these organisms die, forming carbonate sediments (limestone). This calcium carbonate is then subducted into the interior of Earth, where metamorphism due to increasing pressures and temperatures during subduction releases CO₂ that can be degassed via volcanism (see the “metamorphosis” reaction in Figure 17.1). The net carbonate-silicate weathering reaction (Kasting et al., 1993) is



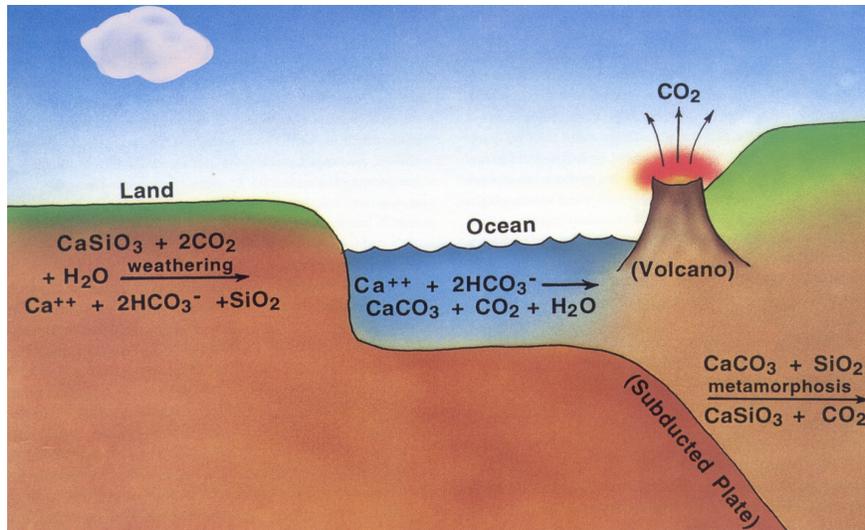


Figure 17.1: Schematic of the carbonate-silicate weathering cycle. Carbon is ingassed into the interior via silicate weathering trapping CO_2 in bicarbonate ions, which are transported into the ocean and then locked up in carbonate minerals. These carbonate species are then subducted to the interior. Carbon dioxide is then returned back into the atmosphere (outgassed) via volcanism.

Importantly, the carbonate-silicate cycle runs faster in hotter climates, as the weathering process that removes CO_2 vapor from the atmosphere is both a temperature-dependent reaction and because it generally rains more in hotter climates, with both increasing the weathering rate. Meanwhile, volcanism is relatively independent of the surface temperature. As a result, in hotter climates carbon is more efficiently removed from the atmosphere to the interior, and vice versa (in colder climates weathering is reduced).

As a result, the carbonate-silicate weathering feedback is a negative (stabilizing) feedback that helps maintain habitable conditions on Earth-like planets that have active tectonics. Note that the typical timescale of the carbonate-silicate weathering feedback is $\gtrsim 1 - 10$ Myr, as it is set by the typical time to remove rock from the seafloor into the interior via subduction. The classical model of the habitable zone includes this negative feedback, which reduces the amount of CO_2 in the atmospheres of planets with high instellation and increases the amount of CO_2 in the atmospheres of planets at low instellation, leading to a greater range of instellations in which there can be habitable surface conditions.

The solid lines in Figure 17.2 show expectations for the habitable zone from a 1D (classical) habitable zone model. The “inner edge” of the habitable zone is the highest incident stellar flux at which the planet can maintain surface liquid water. The inner edge is set by the instellation limit at which the planet loses its water to space by photolysis, with the resulting hydrogen escaping to space (often termed the “moist greenhouse limit”). The “outer edge” of the habitable zone is the lowest instellation at which the planet can maintain surface liquid water. This outer edge is set by the formation of carbon dioxide clouds, which increase the albedo and cool the surface, further promoting CO_2 condensation. Classic 1D models nicely reproduce Earth’s habitability, as well as the potential for Mars to be habitable if it had a thicker Earth-like atmosphere. However, Earth is *very* close to the

inner edge of the habitable zone – just a small perturbation in instellation (or a factor of ~ 10 increase in the CO_2 partial pressure) could cause Earth to reach a moist greenhouse. Though there is little concern that anthropogenic climate change will cause Earth to become uninhabitable over Myr timescales thanks to the silicate-weathering feedback, this is still a useful reminder of the fragility of the habitability of Earth. One consequence of this is that due to the brightening Sun, Earth will begin its transition to a moist greenhouse in ≈ 1.99 Gyr, transforming Earth into a Dune planet (Wolf & Toon, 2015). Earth will then continue to warm due to the lack of an active carbonate-silicate weathering feedback, leading to an eventual runaway greenhouse transition and buildup of CO_2 , through which Earth’s climate will become analogous to present-day Venus.

17.1.2 Clouds and 3D effects

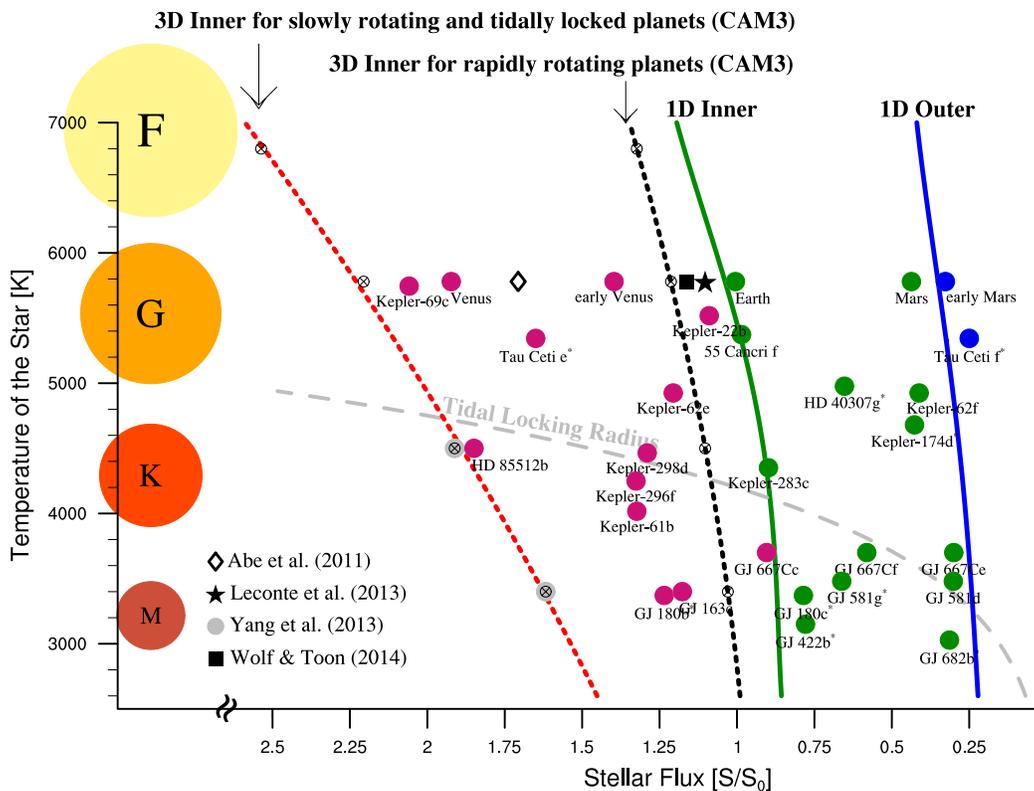


Figure 17.2: The inner and outer edges of the habitable zone from 1D models (solid lines) and 3D GCMs (dashed lines) as a function of stellar effective temperature and incident stellar flux. There are two separate GCM predictions, one for slowly rotating planets and one for rapidly rotating planets. The dayside cloud coverage on slowly rotating planets increases planetary albedo, moving the inner edge closer in. Figure adapted from Yang et al. (2014).

Figure 17.2 also shows two other predictions for the inner edge of the habitable zone, both from three-dimensional climate models (GCMs). These GCMs are similar in their fundamentals to the hot Jupiter GCMs we have discussed previously, but in this case they are tailored to the study of Earth-like atmospheres, most importantly including condensation of liquid water and the resulting formation of water clouds. Both of the GCM predictions for

the inner edge of the habitable zone are closer-in than the classic 1D model predictions. This is because of the cloud formation in the GCM, as liquid water clouds increase the albedo of a planet, causing less insolation to reach the surface from the top-of-atmosphere, and resulting in cooling of the surface relative to a cloud-free state.

There is further a large difference in Figure 17.2 between GCM predictions for slowly rotating and rapidly rotating planets. On slowly rotating planets (where $Ro \gg 1$), cloud formation and dissipation can occur on much shorter timescales than the rotation period. As a result, for moist atmospheres near the inner edge of the habitable zone, there is a persistent deck of clouds on the dayside, greatly increasing the albedo of the planet. In the limiting case of a tidally locked planet, this dayside cloud deck is confined near the substellar point, right at the location of maximum top-of-atmosphere downwelling shortwave radiation. As a result, slowly rotating (i.e., Venus-like) and tidally locked planets have an inner edge of the habitable zone that is predicted to be significantly closer-in than for rapidly rotating planets. This may allow for tidally locked rocky planets around M dwarf stars to maintain habitable surfaces even at close separations, and also may have enabled Venus to have liquid water until ~ 700 Ma (Way et al., 2016).

17.2 Biosignatures

A biosignature is a sign of life on a planet that is remotely detectable, most commonly through spectra of the planetary atmosphere. Biosignatures must reliably point toward inhabited planets, and be detectable through telescopic observations. Life can impact its environment in myriad ways, but the general effect is for life to push the chemistry of its environment away from a state of chemical equilibrium. As a result, we can observationally search for the presence of disequilibrium, either redox disequilibrium due to the prevalence of oxygen and ozone (for modern Earth-like life), or a more general state of disequilibrium by comparing the abundances of multiple species.

17.2.1 Oxygen and ozone

Earth's atmosphere has been oxygenated ever since the Great Oxidation Event (GOE) that occurred at the boundary between the Archean and Proterozoic eons 2.5 Ga. Figure 17.3 shows a timeline of the oxygen content in Earth's atmosphere, along with geochemical records from which this is inferred. The rise in oxygen in Earth's atmosphere corresponds to the onset of life that generates oxygen via photosynthesis, like due to early cyanobacteria (algae). The "smoking gun" of the GOE from the geochemical record is the end of the mass-independent fractionation of sulfur isotopes. The bottom panel of Figure 17.3 shows this geochemical record for both sulfur and carbon isotopes. The $\Delta_{33}S$ shows the difference between the sulfur isotope ratio and that expected from mass-dependent fractionation (the expected way that isotopes are fractionated). The non-zero $\Delta_{33}S$ in the Archean implies a mass-independent fractionation process, which is likely due to sulfur photochemistry as stellar ultraviolet rays could penetrate deep into an anoxic (and thus ozone-free) atmosphere. As oxygen built up, so did ozone, shutting off sulfur photochemistry and the resulting mass-independent fractionation of sulfur. Another key piece of evidence of an early anoxic atmosphere are banded iron formations, which begin to appear at the end of the Archean and form from the precipitation of oxidized iron.

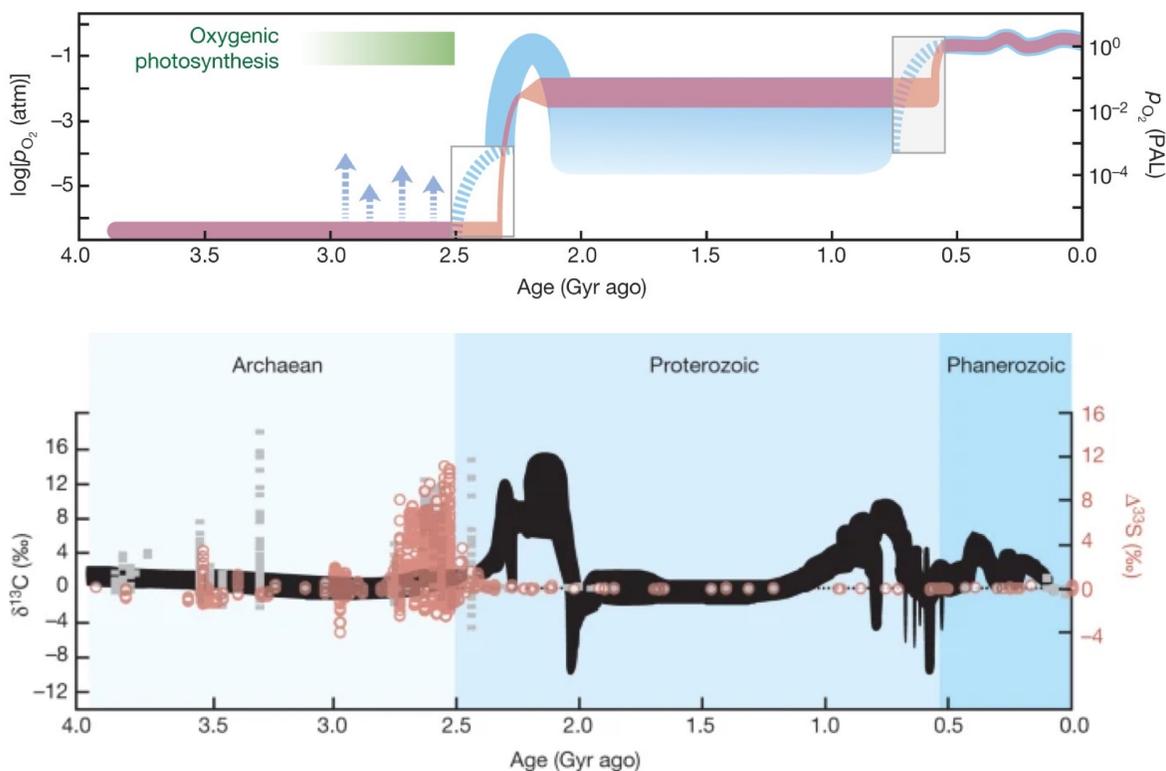


Figure 17.3: The oxygenation of Earth’s atmosphere over time, showing the great oxidation event (GOE) at ~ 2.5 Ga and the transition to complex life ~ 538 Ma (top). The bottom panel shows the carbon (black) and sulfur (red) isotope fractionation over time, showing a sharp change in the sulfur fractionation at the onset of the GOE as mass-independent fractionation of sulfur ceased.

17.2.2 Disequilibrium due to life

The rise of oxygen in Earth’s atmosphere led to methane oxidation through the following reaction



which significantly changed the atmospheric composition of Earth by oxidizing carbon from CH_4 to CO_2 . However, even still, Earth’s atmosphere has a non-zero methane abundance (≈ 1.9 ppm) that is produced by life (most famously, by cows). This implies that this methane oxidation reaction (and other chemical reactions including nitrogen and water) is not in a state of chemical equilibrium – rather, life has driven Earth’s atmosphere and ocean into a state of chemical disequilibrium.

The level of chemical disequilibrium in the atmosphere-ocean system of a planet can be quantified as the “available Gibbs free energy,” which is the difference in Gibbs free energy from the actual state to that in chemical equilibrium. Recall that the Gibbs free energy is related to the equilibrium constant of a reaction K as

$$\Delta_r G = -RT \ln(K), \quad (17.3)$$

and from Section 15.5.2 the equilibrium constant is in turn related to the ratio of partial

pressures of the products to the reactants. Thus, the Gibbs free energy is readily calculated both for a system in thermochemical equilibrium (just based on the temperature and pressure conditions) and from the actual state (just by determining the partial pressures of species for gas, or activity for aqueous species). Figure 17.4 shows a plot of this available Gibbs free energy over Earth history, determined by Krissansen-Totton et al. (2018). The shaded

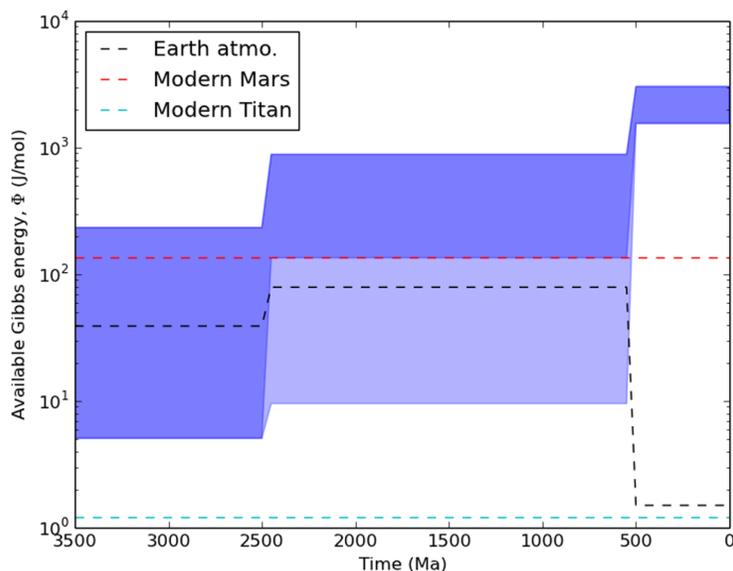


Figure 17.4: The available Gibbs free energy, a metric for disequilibrium, over time for Earth ocean and atmosphere (blue shaded regions) and Earth’s atmosphere (black dashed line) compared with modern Mars (red dashed line) and modern Titan (blue dashed line). The level of disequilibrium in Earth’s biosphere has increased over time, due to the increasing complexity of life. Figure adapted from Krissansen-Totton et al. (2018).

regions show the available Gibbs free energy in Earth’s atmosphere-ocean system – you can see that it increases over time, with two characteristic jumps. One corresponds to the GOE (at 2.5 Ga) and one corresponds to the Cambrian explosion (at 0.53 Ga). Both of these increases correspond directly to increases in the oxygen content of Earth’s atmosphere as shown in Figure 17.3. As a result, the oxygenation of Earth’s atmosphere by life has driven it to a state of disequilibrium.

Thus, if we know that life produces a state of disequilibrium, we can search exoplanetary atmospheres for disequilibrium chemistry and study the environment of the planet to determine if this disequilibrium could be produced by life. Such a search for disequilibrium biosignatures may be feasible with JWST for rocky planets orbiting M dwarf stars. Figure 17.5 shows transmission spectra from JWST NIRSpec/PRISM for TRAPPIST-1e assuming an Archean-like or modern Earth-like atmospheric composition. Due to the strong near-infrared spectral features of CO_2 and CH_4 , searching for life in chemical disequilibrium by constraining the potential partial pressures of carbon dioxide and methane and comparing to chemical equilibrium expectations is more imminently feasible than searching for the spectral features of O_2 and O_3 , which are relatively weak in the near-infrared. Instead, future missions (e.g., Habitable Worlds Observatory) that focus on studying Earth-like planets in the optical are likely required to search for an oxygenic biosignature.

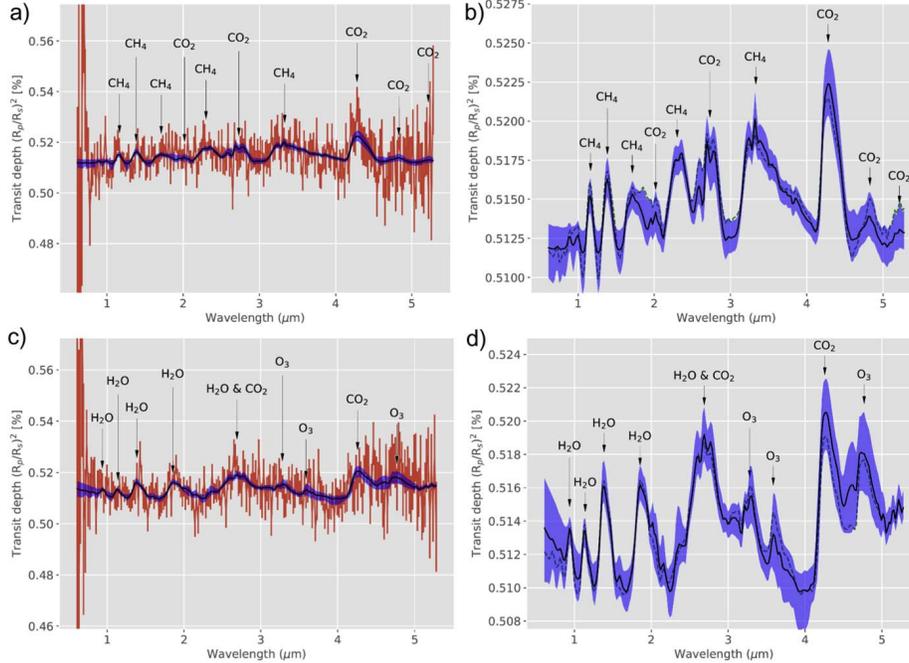


Figure 17.5: Transmission spectra for Archean Earth-like (top) and modern Earth-like (bottom) atmospheres on TRAPPIST-1e as observed with 10 transits by JWST NIR-Spec/PRISM. Disequilibrium biosignatures may be detectable on Archean Earth, and oxygen/ozone could be detectable with sufficient transits. Figure adapted from Krissansen-Totton et al. (2018).

17.2.3 Biosignature false positives

Oxygen is not always produced by life. If we detect oxygen and attribute it to be due to life when it is not, that would be a “false positive” biosignature detection. Astronomers are expected to be skeptical, and thus we must rule out all alternative explanations before claiming a biosignature detection.

One of the most prominent false positive oxygen/ozone biosignatures is the production of abiotic oxygen from water loss. Water high in the atmosphere of an exoplanet can be photodissociated by incident ultraviolet light and broken up into hydrogen and oxygen. The hydrogen is then lost to space due to its low atomic mass, while the oxygen atoms can stick around and form oxygen and ozone. This process (analogous to the moist greenhouse) is expected to be especially common on planets orbiting M dwarf stars, as M dwarfs produce more XUV radiation relative to their full bolometric flux compared to Sun-like stars (see Figure 17.6). As a result, the process of water photodissociation will be more common on water-rich planets orbiting M dwarf stars than Sun-like stars. This will then cause the loss of water on planets orbiting M dwarf stars through the loss of hydrogen and build-up of oxygen – in some (water-rich) model predictions, this process can cause the loss of tens of Earth oceans of water and the build-up of hundreds of bars of oxygen. As a result, we expect that abiotic oxygen production is common in the atmospheres of rocky planets orbiting M dwarf stars.

Figure 17.7 shows a summary of some possible false positive oxygen biosignatures that

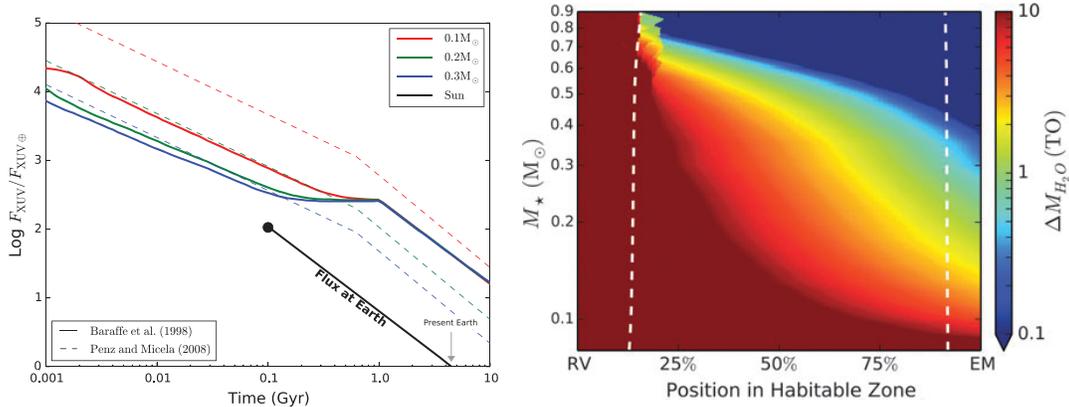


Figure 17.6: Evolution of the XUV irradiation from M dwarf stars of various masses (left-hand panel) from stellar evolutionary models. The right-hand panel shows the effects of stellar XUV on water loss as a function of stellar mass and position in the habitable zone, where the color bar shows the number of oceans lost in a nominal evolutionary model. Figure adapted from Luger & Barnes (2015).

would occur on planets that are not Earth-like. The water loss scenario just described is listed as “ocean loss,” but there are a variety of other possibilities for the production of oxygen false positives. These include thin atmospheres where water can be easily transported to low pressures and photodissociated, leading to the build-up of oxygen (“low non-condensable gas”). These also include CO₂-rich atmospheres, which can lead to photodissociation of

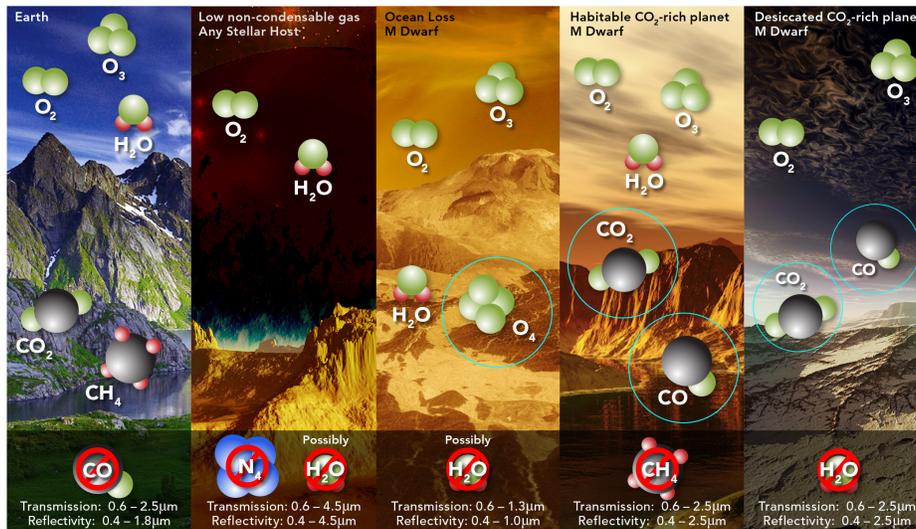


Figure 17.7: Various scenarios to cause false-positive biosignatures for oxygen, compared with an Earth-like case with a robust oxygen biosignature. Figure adapted from Meadows et al. (2018).

carbon dioxide that similarly leads to the build-up of oxygen.

One might then wonder – with all these false positive possibilities, what is the path toward detecting a robust oxygenic biosignature? Figure 17.8 lays out a flowchart that

shows the requirements to do so. The fundamental challenge is that solely a detection of

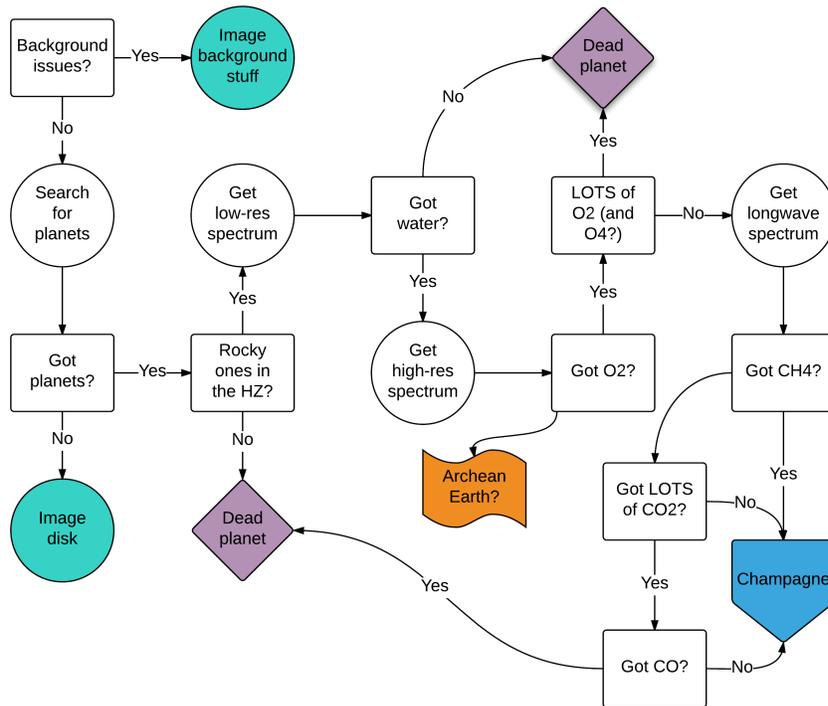


Figure 17.8: A flowchart showing the steps required to observationally determine if a planet has the potential to host Earth-like life. Figure adapted from Meadows et al. (2018).

oxygen is not sufficient – we also need to characterize the environment of this biosignature, including detections and non-detections of other species. For Earth-like life, oxygen (but not too much) must be present in an atmosphere with methane, with some (but not too much) carbon dioxide, and without carbon monoxide. Otherwise, it is impossible to discriminate between the oxygen and a likely false positive. We will discuss NASA’s preferred path toward searching for this series of biosignatures on Earth-like planets around Sun-like stars next. However, first we’ll do a group discussion activity for each of us to think about the best way to search for biosignatures in the atmospheres of exoplanets in our lifetimes.

17.3 Discussion activity

I will assign groups of ~ 3 for this activity. Spend ten minutes discussing the optimal observation strategy to search for biosignatures and understand the planetary environment. As part of your discussion, choose one preferred observational technique, e.g., transit spectroscopy, secondary eclipse spectroscopy, phase curves, direct imaging, or another method. Then, come up with a < 1 minute “elevator” pitch for why this technique is optimal to determine if an exoplanet is inhabited. Each group will present their elevator pitch in front of the class, and we’ll finish by voting on our preferred observational characterization method to search for life.

17.4 Decadal Survey Recommendations

The Astro2020 Decadal survey was released in 2021 (due to pandemic-related delays), and addresses the pathways for astrophysics research in the U.S.A. through the 2020s (and beyond). The Decadal surveys are published by the National Academies of Sciences, Engineering, and Medicine, and provide recommendations for U.S. funding agencies (e.g., NASA, NSF) and the community as a whole for the science to focus on in the coming decade. Astronomy and Planetary Science have separate Decadal surveys, but exoplanets almost entirely falls under the umbrella of Astronomy. The Decadal Surveys are written by a Steering Committee (co-chairs for 2020 were Fiona Harrison and Robert Kennicutt), with input from Science Panels on each relevant science topic (the chair of the Exoplanets, Astrobiology, and Solar System Panel was Vikki Meadows from the University of Washington).

There is myriad recommended astrophysics science in the Decadal survey. Most salient to our class and to the future of exoplanet science is the recommendation that NASA construct a large IR/O/UV space mission to directly image Earth-size exoplanets at approximately 1 au separations from Sun-like stars. This observatory would both directly detect these planets in reflected light, and then follow up to study the reflected light spectrum of the planet from the UV to the IR to search for signs of life. Figure 17.9 is a summary of the evolution of Earth's atmosphere, showing both the abundance (in column mass) of CO₂, CH₄, O₂, O₃, and H₂O over time, along with the observable UV-near IR spectra of Earth during the Archean (4 - 2.5 Ga), Proterozoic (2.5 - 0.53 Ga), and Modern/Phanerozoic (0.53 Ga - present) eons. The recommended large IR/O/UV mission would have broad UV-near IR wavelength coverage in order to identify spectral signatures of the biosignatures oxygen, ozone, and methane as well as the habitability indicator water and the greenhouse gas carbon dioxide. Importantly, the detectability of these species is strongly dependent on abundance, and Earth is our only guideline for the oxygenation of a habitable planet atmosphere. Because Earth has only had a modern level of oxygen for ~ 538 Myr, this implies that we likely need to detect many potentially inhabited Earth-like planets in order to measure a modern-Earth like amount of oxygen. As we'll discuss, this drives the requirements of the recommended large IR/O/UV mission.

Of course, the large IR/O/UV flagship mission is not the only relevant astrophysics project that was recommended by the decadal survey. This pioneer exoplanet mission was recommended alongside a wide range of interesting astrophysics, including continuing to push on gravitational wave transients and multi-messenger astronomy, building the next generation VLA and next generation IceCube, and exploring possibilities for far-infrared and X-ray probe missions. Figure 17.10 shows the approximate timeline of all recommended programs, listed in order of their science category. The large IR/O/UV flagship mission has by far the latest expected date of all of these programs, with an anticipated launch date of the mid-2040s. This was somewhat of a departure from previous Decadal surveys, with the Astro2020 survey not just dictating science through the 2020s but effectively through the next three decades (and beyond). The long wait time for the large IR/O/UV mission is due to a combination of technological challenges related to coronagraphy as well as budgetary demands. As a result, the Decadal also described the importance of the existing HST and JWST observatories on the pathway to characterizing exoplanet atmospheres, as summarized in Figure 17.11. The 2020s and 2030s will likely be focused on characterizing the atmospheres

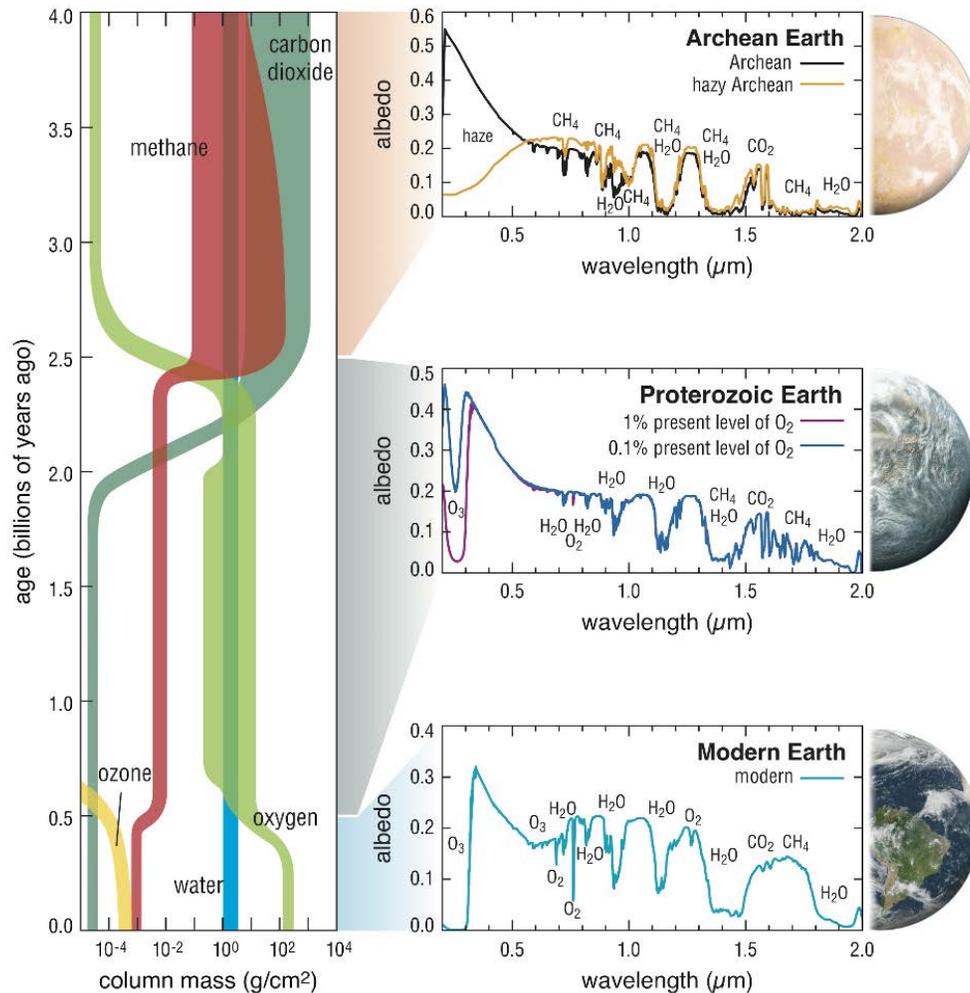


Figure 17.9: Earth’s atmospheric composition through time (left) and resulting spectra (right) for the Archean Earth (4 - 2.5 Ga), Proterozoic Earth (2.5 - 0.5 Ga), and Phanerozoic (Modern) Earth (0.5 Ga - present). Figure adapted from the Astro2020 Decadal survey.

of gaseous planets, sub-Neptunes, and rocky planets orbiting M dwarf stars. Only with a large IR/O/UV flagship can exoplanet science push toward directly imaging Earth-Sun twins in reflected light.

The Decadal specifically recommended that the large IR/O/UV flagship have an approximate inscribed mirror diameter of 6 meters. This is slightly smaller than JWST’s 6.5m mirror, but the effective aperture size is similar due to the recommended use of an off-axis secondary mirror. This 6m inscribed diameter was chosen in order to allow for the detection of ≈ 25 Earth-sized planets in the habitable zones of Sun-like stars, and corresponds to the red dot in Figure 17.12. As we’ll discuss, this choice of a 6m mirror falls between other proposed mission concepts (termed LUVOIR and HabEx), and is meant to be a middle-ground between cost and performance in terms of the number of Earth-like planets that can be characterized.

Even though the large IR/O/UV mission recommended by the decadal is a compromise,

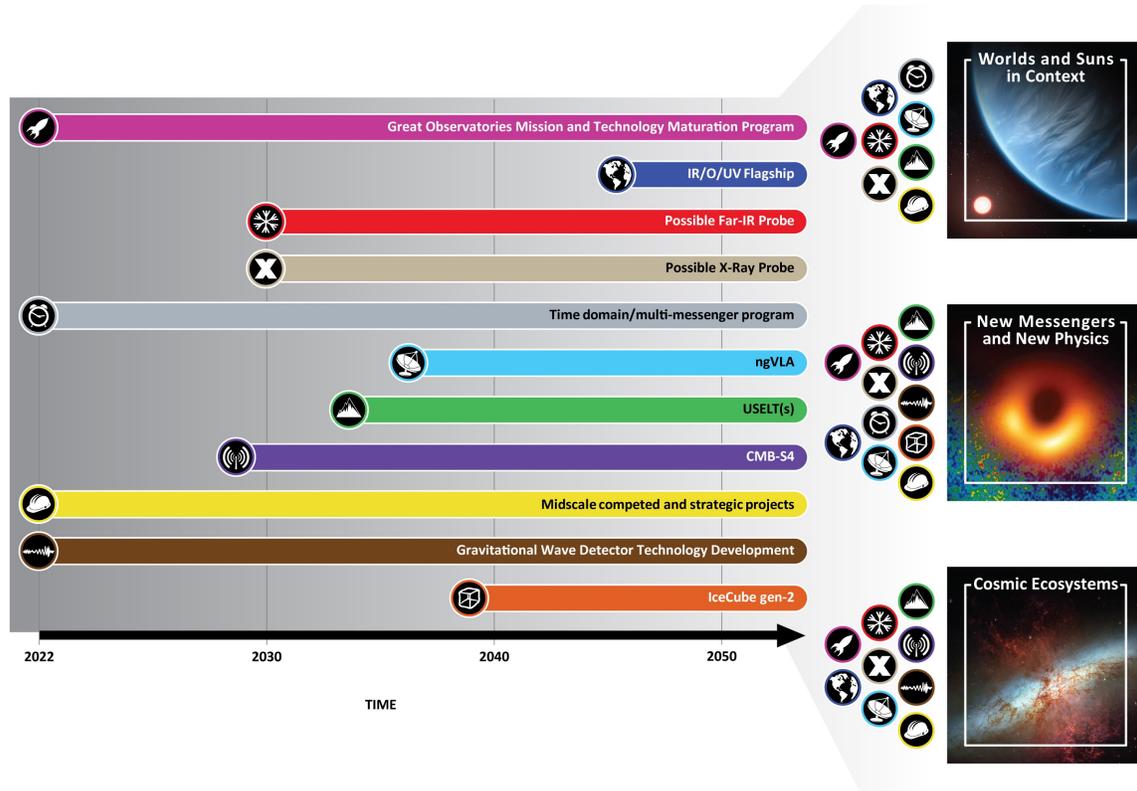


Figure 17.10: Recommended timeline of medium and large missions from the Astro2020 Decadal survey. The IR/O/UV Flagship on the second row is now known as the NASA Habitable Worlds Observatory (HWO) mission concept. Figure adapted from the Decadal survey.

it would still represent an enormous investment from NASA and the U.S. Government. Figure 17.13 shows budget projections of each of the major recommendations from the Astro2020 Decadal (shaded regions), along with the NASA expected budget projection (blue line). Even though the funding for the large IR/O/UV observatory is not meant to officially begin until 2026, it would represent the bulk of NASA’s budget integrated over the 2035-2045 timeframe (the decade before launch). The yearly cost of this flagship would far exceed NASA’s expected budget (by > 500 million USD). As a result, the mission would likely require special budgetary approval by Congress for it to be feasible in the planned timeframe.

17.5 Habitable Worlds Observatory

NASA recently dubbed the large IR/O/UV flagship mission recommended by the Decadal the Habitable Worlds Observatory (HWO). NASA is now beginning a multi-year process of fleshing out this mission concept, including determining potential science goals and assessing risks. This mission would aim to detect multiple ExoEarths with direct imaging, as well as characterize these planets to search for biosignatures and associated environmental context.

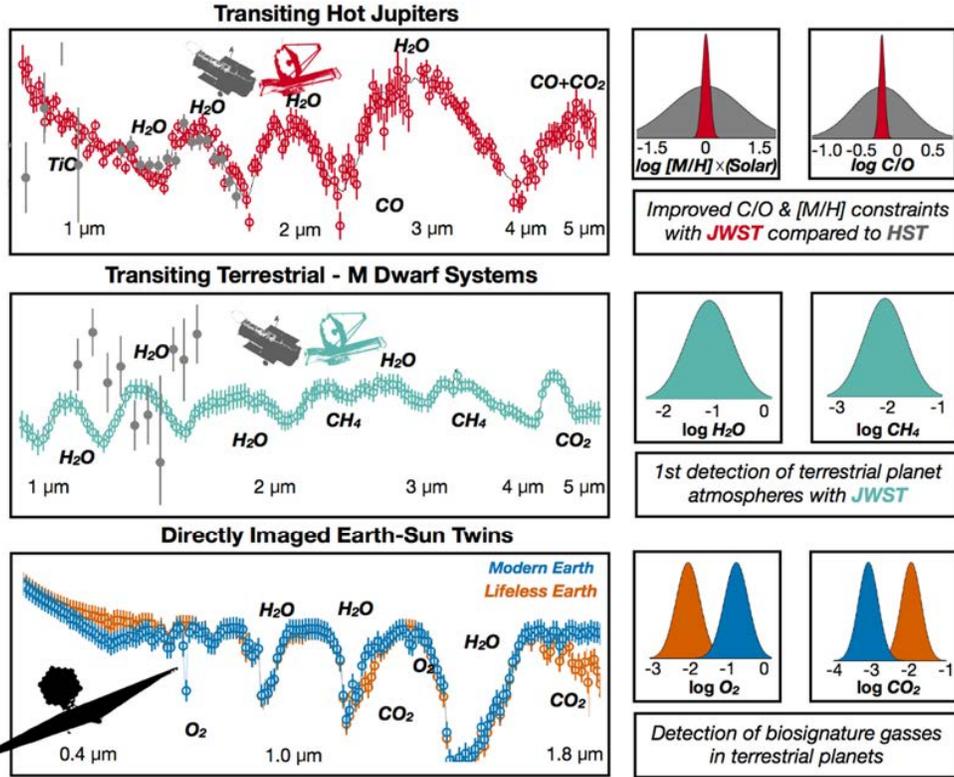


Figure 17.11: Examples of upcoming spectral characterization of exoplanets with HST, JWST, and HWO. The top and middle panels show simulated transmission spectra for hot Jupiters and terrestrial planets with a combination of HST and JWST. The bottom panel shows a simulated reflectance spectrum of an Earth twin with HWO. Figure adapted from the Decadal survey.

17.5.1 Detecting a sample of potentially habitable ExoEarths with direct imaging

The HWO mission concept and the recommendation from the Decadal are the result of decades of work by other mission concept teams. The two most prominent mission concepts to directly image Earth-sized planets orbiting in the habitable zones of Sun-like stars pre-Decadal were LUVOIR and HabEx. LUVOIR itself had two different configurations, LUVOIR-A and LUVOIR-B, with A having a large on-axis mirror and B having a smaller off-axis mirror. Figure 17.14 shows the predicted ExoEarth Candidate (EEC) yield for these two mission configurations. Both of the possible LUVOIR configurations would have used a coronagraph hosted in the optics of the telescope itself. The HabEx mission concept, meanwhile, proposed the use of an external occulter (“starshade”) that would formation fly with a ≈ 4 m diameter telescope to block the light from ExoEarth host stars. The HabEx projections for EEC yield are shown by the yellow line in Figure 17.14. As you can see, the number of habitable planets that would be detected by HabEx would be lower due to the need to physically move the starshade to align it with the telescope and observe a different stellar system. However, HabEx’s strengths are that the mirror would be smaller (and thus

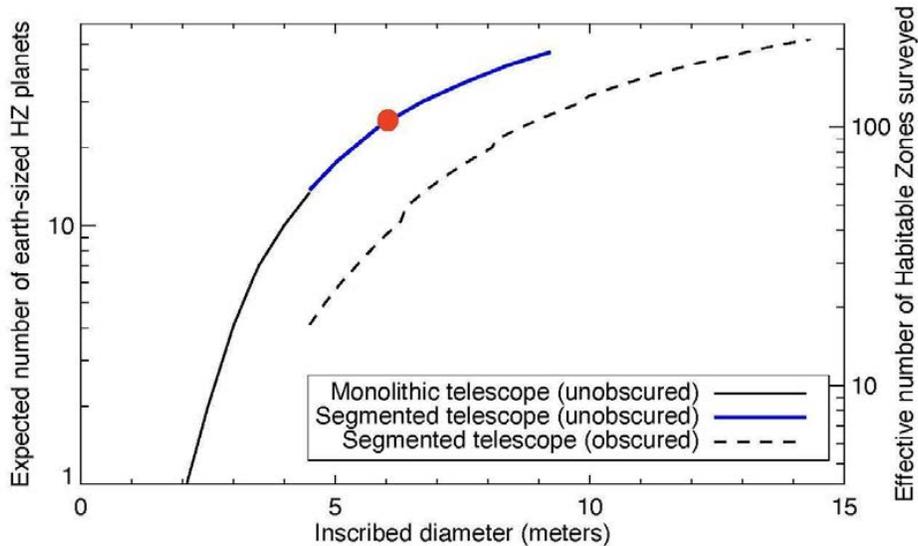


Figure 17.12: Expected number of Earth-sized habitable zone planets as a function of inscribed diameter of HWO assuming $\eta_{\oplus} = 0.24$. The red dot shows the recommend mirror diameter for a large IR/O/UV mission from the Decadal of ≈ 6 m. Figure adapted from the Decadal survey.

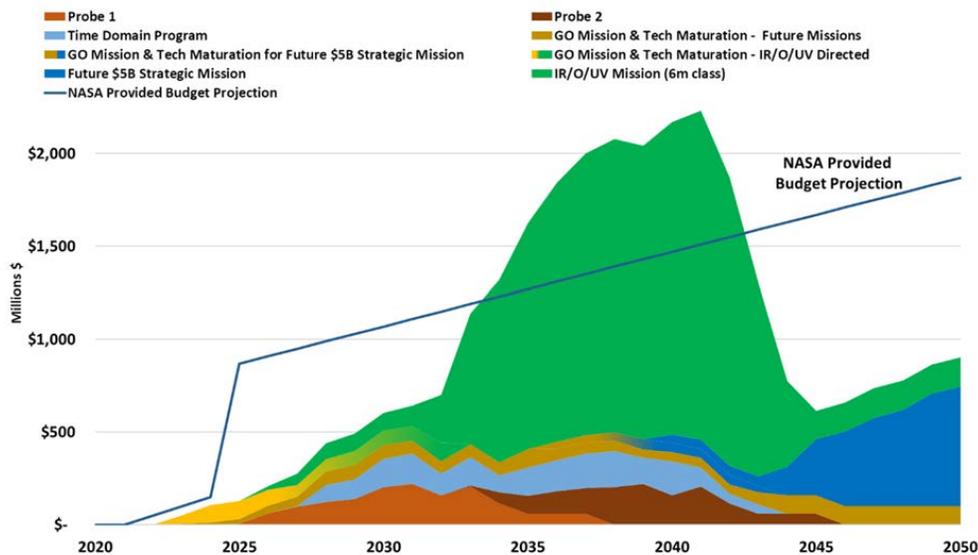


Figure 17.13: Budget projections for recommended programs from the Decadal (filled colors) compared to the NASA overall budget projection (solid blue line). The large IR/O/UV mission (now HWO) is the green shaded region that goes well above the NASA budget projection. Figure adapted from the Decadal survey.

cost likely lower) and that the starshade may more efficiently reduce the contrast down to the $\sim 10^{-10}$ level required to detect an Earth-sized planet around a Sun-like star at 1 au in reflected light.

HWO is planned to have a broad wavelength range covering from the near-UV to the

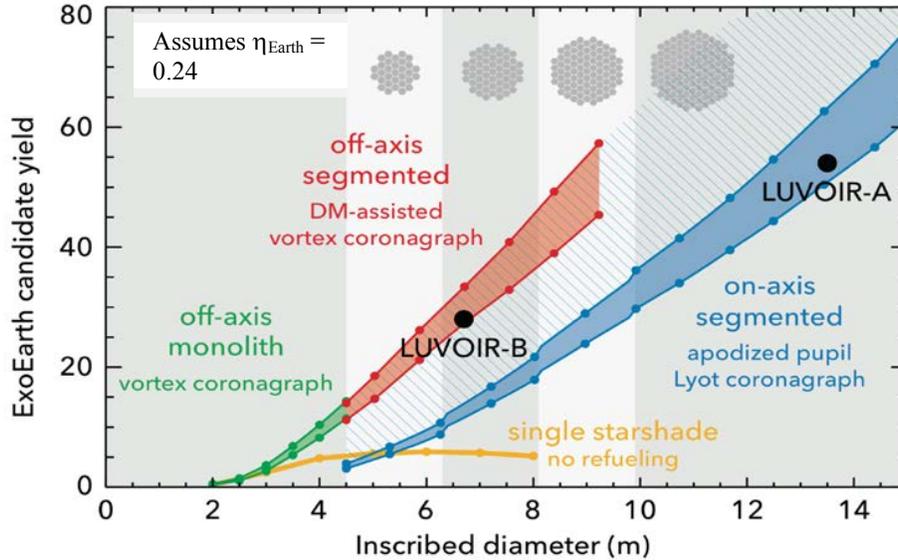


Figure 17.14: Expected ExoEarth candidate yield for the LUVUOIR-A and LUVUOIR-B configurations assuming $\eta_{\oplus} = 0.24$. Shaded regions correspond to different mission architectures, including segmented or monolithic mirrors, on-axis or off-axis secondaries, and including the possibility of a starshade. Figure adapted from the Decadal survey, originally from the LUVUOIR final report.

near-IR. This was also planned in each of the LUVUOIR and HabEx mission concepts, and is motivated in order to detect myriad biosignatures, habitability indicators, and other species that better provide environmental context for a purported biosignature. Figure 17.15 shows a list of various biosignatures and false positive discriminants, as well as their relevant absorption features in the UV-Visible and near-IR. The key motivation for including the UV

Table 3-1. *Desired spectral features for biosignature assessment.*

Biosignatures & False Positive Discriminants (indicated with *)		
Molecules/Feature	UV-VIS wavelengths (0.2–1.0 μm)	NIR wavelengths (1.0–2.0 μm)
O_2	0.2, 0.63, 0.69, 0.76 (strong)	1.27
O_3	0.2–0.35 (strong), 0.5–0.7	
O_4 (O_2 - O_2)*	0.345, 0.36, 0.38, 0.45, 0.48, 0.53, 0.57, 0.63	1.06, 1.27 (strong)
CO^*		1.6
CO_2^*		1.05, 1.21, 1.44, 1.59
CH_4	0.6, 0.79, 0.89, 1.0	1.1, 1.4, 1.7
N_2O		1.5, 1.7, 1.78, 2.0
Organic haze	< 0.5	
Vegetation red edge	0.6 (halophile), 0.7 (photosynthesis)	

Figure 17.15: Potential biosignature spectral features of interest for LUVUOIR, along with their associated wavelengths in the UV-visible and near-infrared. Figure adapted from the LUVUOIR final report.

is to include the strong Hartley-Huggins bands of ozone between $0.2 - 0.35 \mu\text{m}$. The visible has a range of O_2 and O_3 spectral features, as well as methane features, and the possibility of detecting evidence for a “red edge” due to photosynthetic vegetation on the surface. The near-IR is required largely to detect habitability indicators and false positive discriminants, including water, carbon dioxide and carbon monoxide, and methane. The exact planned wavelength coverage and instrument modes of HWO are currently being determined as part of its mission concept phase.

17.5.2 Characterizing ExoEarths: reflectance spectra, rotational mapping

Figure 17.16 shows a simulated spectrum which approximates that expected from HWO for an Earth-twin around a Solar-twin. There is a clear ozone feature in the UV, multiple

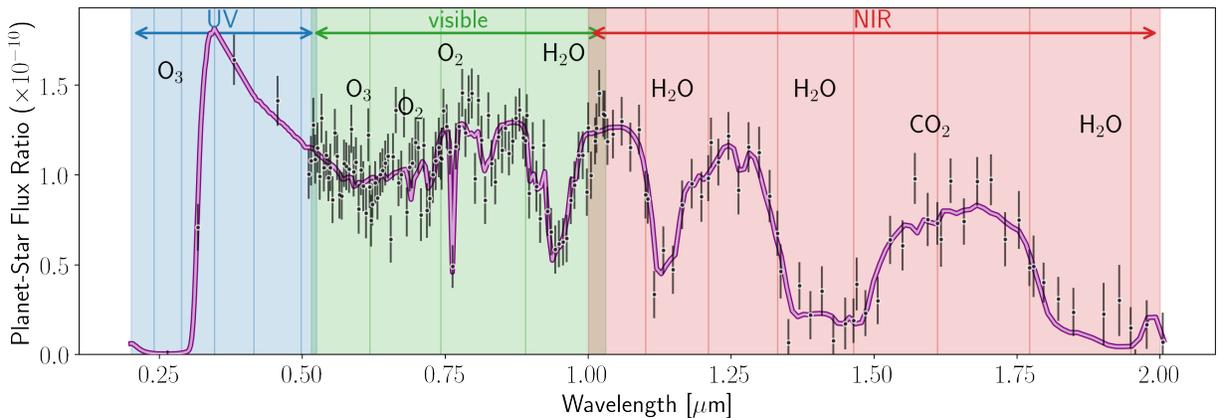


Figure 17.16: Simulated reflectance spectrum of Earth around the Sun. Clear absorption features of ozone, oxygen, and water are seen in the UV-visible, and signatures of water and carbon dioxide occur in the near-infrared. Figure adapted from the Decadal and LUVOIR final report.

oxygen features in the visible, water features in the visible/near-IR, and carbon dioxide in the near-IR. The need to both detect biosignatures, signs of water vapor (and thus a potential surface ocean) and characterize the atmospheric properties through constraining the level of carbon dioxide is what drives the simulated wavelength range. Generally, it is expected that the approximate wavelength range of HWO will be $0.2 - 2 \mu\text{m}$, with wiggle room depending on instrument and detector design along with science requirements to rule out possible biosignature false positives.

Beyond the time-integrated reflectance spectrum alone, there is a range of other science that can be done for Earth-like planets with HWO. One of the most compelling is rotational phase mapping of these planets in order to study how their albedo (reflectance) varies with rotational phase. An example of this using Earth itself is shown in Figure 17.17. From these reflectance rotational phase curves, the rotation period of the planet can be constrained if there are spatial albedo variations. Additionally, a rough surface map of the planet can be made by matching the albedos of various surfaces. Figure 17.17 shows the specific example of using ocean glint during the crescent phase to identify the presence of a liquid water ocean, and that the location of the liquid water ocean approximately matches with the Atlantic and

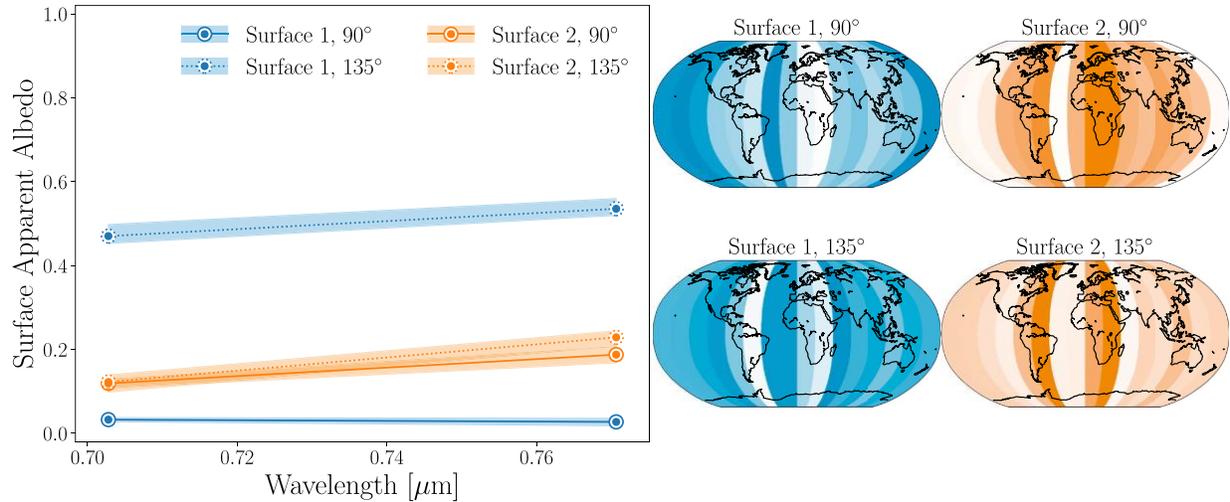


Figure 17.17: Results from simulated mapping of the surface of Earth in reflected light with a HWO-like mission. Multiple surface features with different albedos are inferred with longitudinal variation, corresponding to continents and ocean. A liquid ocean is inferred through the presence of glint, which significantly increases its apparent albedo near crescent phase. Figure adapted from Lustig-Yaeger et al. (2019).

Pacific ocean basins for a simulated Earth test case. As a result, rotational phase mapping could directly probe the habitability of rocky exoplanets by searching for the presence of surface liquid water.

17.6 Prediction activity!

Given its technological complexity, the Decadal survey only provides a rough timeline for the IR/O/UV Flagship (now known as Habitable Worlds Observatory). Let's get in small groups to discuss and make some predictions to then see where we, as a class, stand in terms of our expectations for the long-term future of the search for life on exoplanets. Discuss the following questions in groups of ~ 3:

1. What year do you think the Habitable Worlds Observatory will be launched in?
2. Do you think a robust biosignature will be detected before the Habitable Worlds Observatory is launched? If so, what method (and/or observatory) do you think will detect this biosignature?
3. How many Earth-like exoplanets orbiting Sun-like stars do you think the Habitable Worlds Observatory will detect?
4. Do you think that Habitable Worlds Observatory will detect life on an exoplanet? If so, how many inhabited planets do you think it will find?
5. Assuming the Habitable Worlds Observatory does find Earth-like planets with a sign of life, what type of signature do you think the Habitable Worlds Observatory will detect? Options include but are not limited to oxygen with environmental context, disequilibrium biosignature, technosignature, something unexpected, etc.

18 Exoplanet characterization: transmission spectroscopy

Enclosed are notes that you may find helpful to review before or after Dr. Munazza Alam’s lecture on transmission spectroscopy. Today’s reading is the Kreidberg review chapter on transmission spectroscopy. This will detail the principles of transmission spectroscopy as well as how it can be used to characterize exoplanet atmospheres.

18.1 Fundamentals of transmission spectroscopy

18.1.1 Qualitative description

Transmission spectroscopy probes the atmospheres of exoplanets by studying the transmission (or filtering) of light from the host star through the limb of an exoplanet that appears to occult the host star. Figure 18.1 shows the geometry of a transit event, along with the geometry of the secondary eclipse that is used to derive planetary emission spectra (as we’ll discuss in the next class). Transmission spectra can thus only be observed for transiting

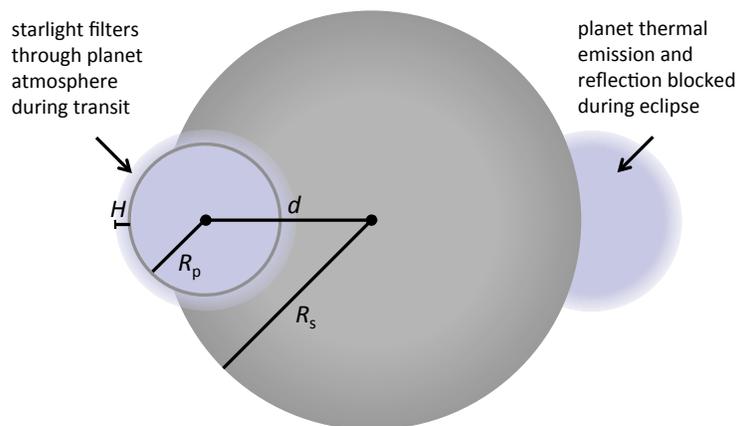


Figure 18.1: Geometry of transmission and emission observations. Here the stellar radius is R_s , planet radius is R_p , the separation of centers of planet and star is d , and the pressure scale height of the planet’s atmosphere is H . Figure adapted from Kreidberg (2017).

exoplanets, and are generally a small effect (at most on the orders of parts per thousand) on top of the larger transit signal. Similarly to the transit method itself, transmission spectral signatures are larger for planets orbiting smaller stars. However, there is also a strong dependence of the transmission spectra signal on the atmospheric composition, temperature, and gravity of the orbiting planet, which we derive below.

18.1.2 Transmission flux ratio

The transmission spectral features are a small atmospheric contribution imprinted on top of the total transit depth. Recall that the transit depth for the solid-body planet is $\delta = (R_p/R_s)^2$. The transit depth including the atmospheric contribution is then

$$\delta_\lambda = \frac{(R_p + A_{H,\lambda})^2}{R_s^2}, \quad (18.1)$$

where $A_{H,\lambda}$ is the apparent atmospheric height which is being probed. We can write the apparent atmospheric height as a number of scale heights, i.e.,

$$A_{H,\lambda} = nH = n\frac{RT}{g} \quad (18.2)$$

where H is the (isothermal) pressure scale height. As a result, we can isolate the transmission spectral feature amplitude (i.e., the contribution due to the atmosphere alone) as

$$\delta_{\lambda,atm} = \frac{(R_p + nH)^2}{R_\star^2} - \frac{R_p^2}{R_\star^2} \approx \frac{2R_p nH}{R_\star^2} \approx \frac{2R_p nRT}{gR_\star^2}. \quad (18.3)$$

Typically, $n \approx 2$ for low spectral resolution observations of cloud-free atmospheres (Kreidberg, 2017). As a result, we expect the transmission spectral feature amplitudes to be larger for hotter and lower-gravity planets, which have a larger scale height. We also expect transmission spectral features to be larger for low mean molecular weight atmospheres, as decreasing mean molecular weight increases R , which also increases the scale height. Note that our derivation does not include the effects of clouds, which also mute spectral features by increasing the optical depth of the atmosphere at low pressures – see the next section.

18.1.3 Beer’s law

Consider radiation with an initial radiance at a given wavelength I_λ that impinges upon an absorbing and emitting slab that has mass density ρ , absorption coefficient k_λ , and thermal radiance B_λ . Lambert’s law states that this slab absorbs radiation, causing a decrease in radiance leaving the slab (at distance dl) of

$$dI_\lambda = -I_\lambda k_\lambda \rho dl \text{ Lambert's law.} \quad (18.4)$$

Similarly, Kirchhoff’s law states that substances in thermodynamic equilibrium emit as efficiently as they absorb, so the change in emitted radiation must be

$$dI_\lambda = B_\lambda k_\lambda \rho dl \text{ Kirchhoff's law.} \quad (18.5)$$

Putting these together, we can write Schwarzschild’s equation for radiative transfer ignoring scattering,

$$\frac{dI_\lambda}{dl} = \rho k_\lambda (B_\lambda - I_\lambda) \text{ Schwarzschild's equation.} \quad (18.6)$$

This equation tells us how the radiance from the slab is affected by a combination of absorption and emission. If we further define the optical path τ as

$$\tau_\lambda = \int_0^l \rho k_\lambda dl', \quad (18.7)$$

we can re-write Schwarzschild’s equation as

$$\frac{dI_\lambda}{d\tau} = B_\lambda - I_\lambda. \quad (18.8)$$

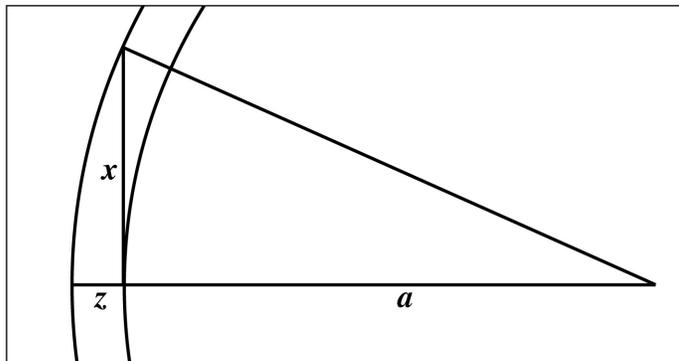


Figure 18.2: Geometry used to derive the slant optical path of transmission through an atmosphere. Figure adapted from Fortney (2005).

In the case where the medium is purely absorbing, $B_\lambda = 0$ and we can integrate to solve for the radiance as a function of optical path as

$$I_\lambda = I_\lambda(0)e^{-\tau} \text{ Beer's law.} \quad (18.9)$$

Beer's law can be applied to transmission geometry, as in most cases we can ignore the atmospheric emission contribution to the transit depth. This then allows us to relate the observed flux deficit to an optical path, and thus a region in the atmosphere that is probed via transmission spectroscopy, given that regions with $\tau > 1$ are opaque to transmitted starlight.

Note that here τ is the slant optical path through the limb of the planet, rather than the optical depth derived previously for light travelling in the vertical direction toward and away from the planetary surface – this slant optical path is significantly larger than the typical optical depth, which is why transit spectra probe relatively low pressures compared to emission spectra. Figure 18.2 shows the geometry of the slant optical path (in the direction of x) relative to the vertical optical depth (in the direction of z). We will now follow the derivation in Fortney (2005). Recall that an isothermal atmosphere in hydrostatic equilibrium has a dependence of pressure p and equivalently number density n on height above a given level z_0 as

$$\begin{aligned} p(z) &= p(z_0) \exp\left(-\frac{(z - z_0)}{H}\right), \\ n(z) &= n(z_0) \exp\left(-\frac{(z - z_0)}{H}\right), \end{aligned} \quad (18.10)$$

where $H = RT/g$ is the pressure scale height, which is equal to the density scale height for an isothermal atmosphere. From Figure 18.2, we can note that

$$a^2 + x^2 = (a + z)^2 = a^2 + 2az + z^2. \quad (18.11)$$

If $2az \gg z^2$, then $z \approx x^2/2a$. As a result, we can write the dependence of number density on x as

$$n(x) = n_0 \exp\left(-\frac{x^2}{2aH}\right). \quad (18.12)$$

If we integrate this dependence over x to find the integrated density N_H from horizon to horizon, we find

$$N_H = \int_{-\infty}^{\infty} n(x)dx = n_0 \frac{1}{2} \sqrt{\pi} \sqrt{2aH} \operatorname{erf} \left(\frac{x}{\sqrt{2aH}} \right) \Big|_{-\infty}^{\infty} = n_0 \sqrt{2\pi aH}. \quad (18.13)$$

The ratio of the horizontally integrated density to the vertically integrated density is equal to the ratio of the horizontal to vertical optical path. This is

$$\frac{N_H}{N_V} = \frac{\tau_H}{\tau_V} = \sqrt{\frac{2\pi a}{H}}. \quad (18.14)$$

Because $a \gg H$, this value is always much larger than one – for Earth it is ~ 75 , and for Jupiter it is ~ 128 . Generally speaking, due to the high slant optical path, transit spectra probe low pressures of ≈ 1 mbar, with higher resolution observations probing even lower pressures in absorption lines. This also causes condensate clouds and hazes to have a large impact on the depth of transmission spectral features, as they move the continuum (deepest region that can be probed) to lower pressures (Fortney, 2005).

18.1.4 Application to observed spectra: example of WASP-43b

Figure 18.3 shows a transmission spectrum observed with the Hubble Space Telescope/Wide Field Camera 3 (HST/WFC3) (bottom) of the hot Jupiter WASP-43b. This

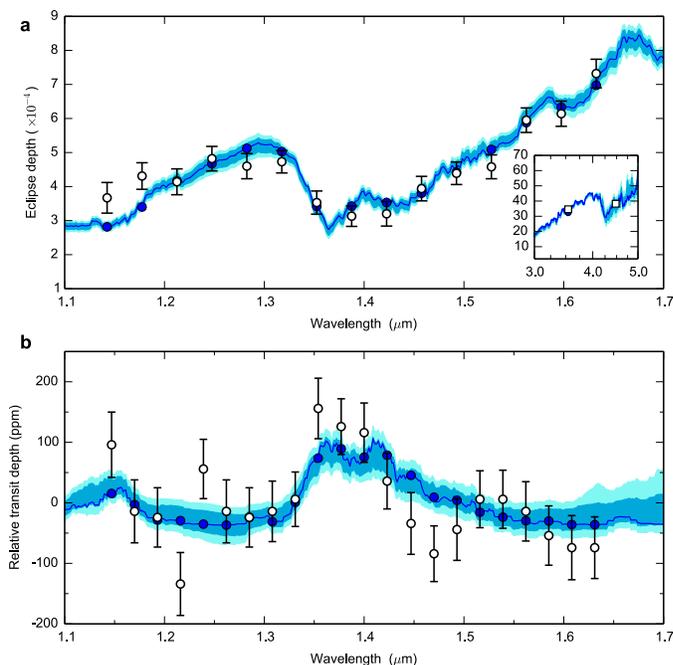


Figure 18.3: Emission spectrum (top) and transmission spectrum (bottom) of WASP-43b observed with the Hubble Space Telescope Wide Field Camera 3 instrument (inset on top shows observations with the Spitzer Space Telescope). Note the spectral feature (bump) in the transmission spectrum at 1.4 μm due to water vapor absorption, which also causes an absorption feature (dip) in the emission spectrum. Figure adapted from Kreidberg (2017).

shows the change in transit depth as a function of wavelength, i.e., isolating only the atmospheric contribution A_H to the total transit depth. There is a bump in the spectrum at a wavelength of 1.4 μm . This bump corresponds to an increase in the effective transit depth, i.e., an increase in the effective size of the planet at 1.4 μm . As a result, this is an absorption feature, as the planet absorbs more light at 1.4 μm than at other wavelengths in

the instrumental range. Note that this bump in the transmission spectrum at $1.4 \mu\text{m}$ aligns well with a dip in the emission spectrum (top) at the same wavelength. We'll cover emission spectra soon, but this dip in emission spectra similarly corresponds to an absorption feature due to the same species causing the bump in the transmission spectrum.

Given a transmission spectral feature at a given wavelength, astronomers can then infer the species that is causing this feature by comparing with model atmospheres. The wavelength-dependence of the atmospheric absorption is determined by the absorption cross section of a given species σ , which (as we discussed in the gas giant interiors class) is related to opacity κ as

$$n\sigma_{\lambda} = \kappa_{\lambda}\rho, \quad (18.15)$$

where subscripts indicate a wavelength-dependence. Figure 18.4 shows the dominant cross-sections in a typical hot Jupiter atmosphere in chemical equilibrium at Solar metallicity, a temperature of 1500 K, and a pressure of 300 mbar. The strongest feature near $1.4 \mu\text{m}$

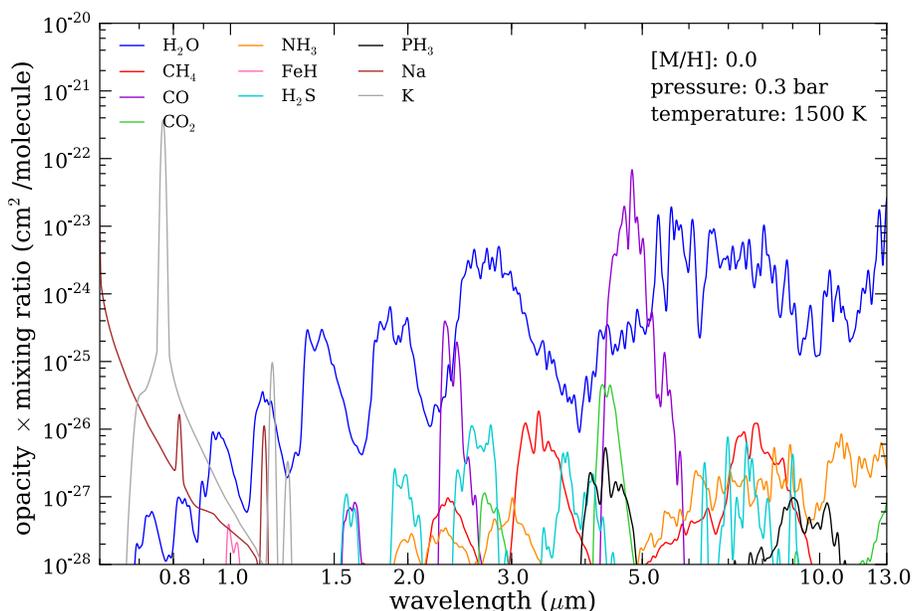


Figure 18.4: Cross sections of various species in a Solar composition atmosphere of a hot Jupiter at a pressure of 300 mbar. Note the strong features due to water vapor and carbon monoxide in the near-infrared. Figure adapted from Kreidberg (2017).

is due to water vapor absorption, which we can identify as the absorber in the WASP-43b HST/WFC3 transmission spectrum. Note that the species which dominates absorption cross-sections is strongly dependent on wavelength, with water vapor, carbon dioxide, and methane (greenhouse gases that reduce Earth's outwelling radiation) being most important in the near-infrared, and sodium and potassium having significant contributions in the visible.

18.2 Highlights of transmission spectroscopy

The first detection of species in the atmosphere of an exoplanet (and thus, detection of an atmosphere itself) via transmission spectroscopy was for the exoplanet HD 209458b with the Hubble Space Telescope/Space Telescope Imaging Spectrograph (HST/STIS) instrument

(Charbonneau et al., 2002). Figure 18.5 shows this detection, in terms of the number of photoelectrons received from the spectrograph as a function of wavelength. Note that this

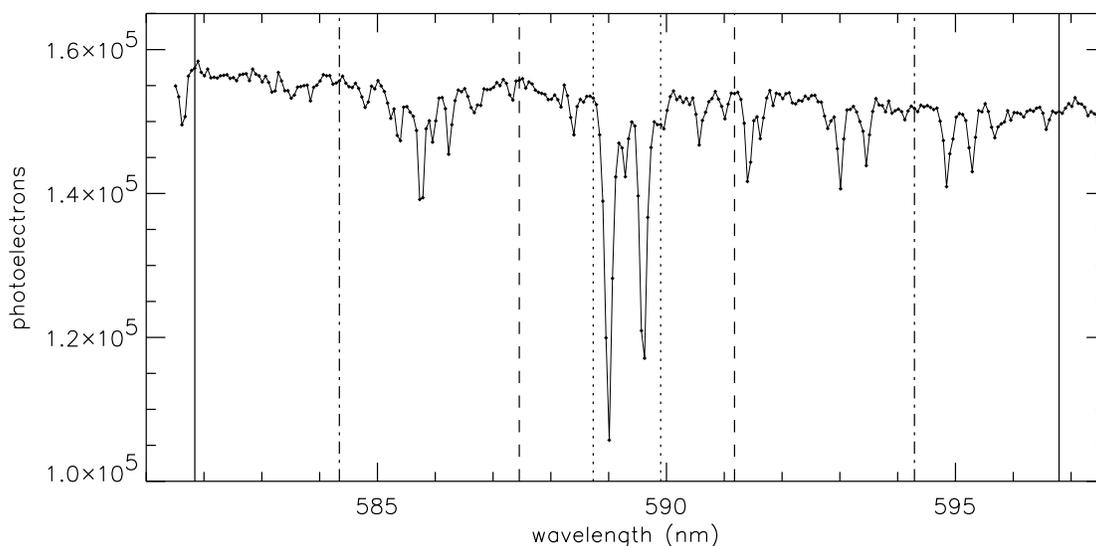


Figure 18.5: First detection of an exoplanet atmosphere. Shown are the transmission spectral features of the Na doublet at 589 nm in the atmosphere of HD 209458b as observed with HST/STIS. Figure adapted from Charbonneau et al. (2002).

y-axis goes in the opposite direction of the current convention for transmission spectroscopy – here absorption is a dip rather than a bump in the spectrum. This HST/STIS spectrum shows clear absorption due to the Na doublet centered at 589 nm, and thus an indication of an absorbing atmosphere of this hot Jupiter. Notably, the depth of these spectral features was smaller than expected for a clear atmosphere. This implies that aerosols increase the height of the atmosphere’s continuum, reducing the depth of these Na spectral features.

The most precise transmission spectrum of a sub-Neptune is the HST/WFC3 spectrum for GJ 1214b (Kreidberg et al., 2014). Figure 18.6 shows the spectro-photometric transit observations of this planet from 1.15 – 1.63 μm , which is directly measured, here from co-adding 15 transit observations together. The relative transit depth as a function of wavelength is then measured from these transit observations at varying wavelengths, and turned into the transmission spectrum shown in Figure 18.7. The measured transmission spectrum of GJ 1214b with HST/WFC3 is consistent with a flat line, with no evidence for absorption by either Solar composition atmosphere (top panel) or even atmospheres comprised of high mean molecular weight species like water, methane, or carbon dioxide. This implies that GJ 1214b has a high-altitude aerosol layer that is roughly wavelength-independent (“gray”) in the WFC3 bandpass that prevents transmission spectral observations from probing the absorption of the gaseous species in the atmosphere. This is one of many instances of clouds impacting the transmission spectra of exoplanets, and generally speaking transmission spectra have found that clouds are ubiquitous in exoplanet atmospheres.

The state of the art of transmission spectral observations is with JWST, as evidenced by the early release science (ERS) observations of the hot Jupiter WASP-39b. Figure 18.8 shows the NIRSpec/PRISM transmission spectrum of WASP-39b (with a single transit!)

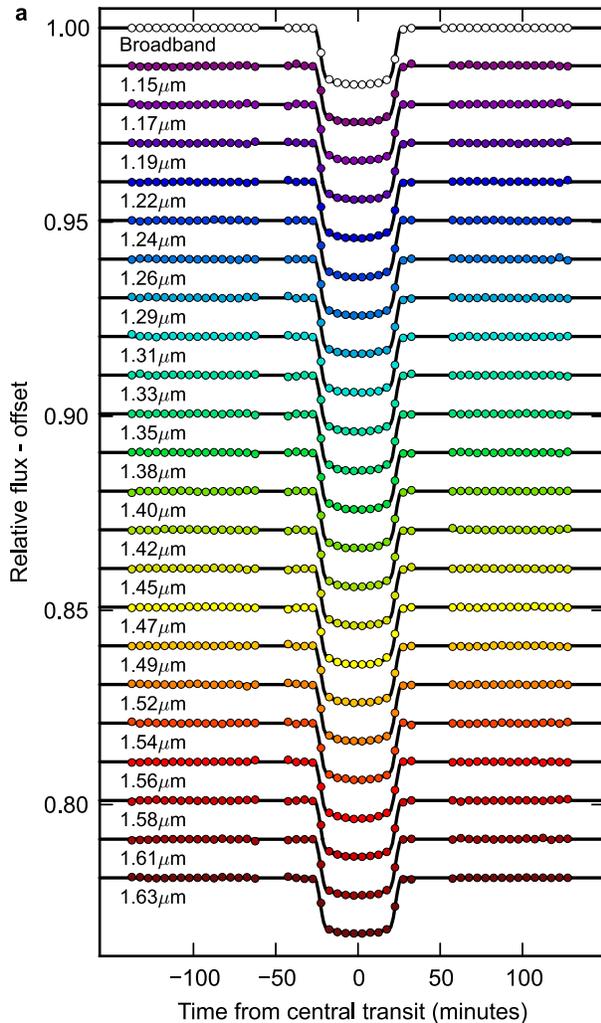


Figure 18.6: Spectro-photometric transit observations of GJ 1214b with HST/WFC3 from 1.1 – 1.7 μm . Figure adapted from Kreidberg et al. (2014).

by the JWST ERS team (Rustamkulov et al., 2023). This is a broadband visible-NIR transmission spectrum from 0.6–5.3 μm , significantly expanding on the wavelength coverage of the HST/WFC3 instrument. Given the broad wavelength coverage, detection of a variety of species is enabled, including a large carbon dioxide feature centered at 4.2 μm , clear water absorption features at multiple wavelengths, sodium in the visible, and SO_2 at 4 μm . Notably, the SO_2 can only be produced by photochemistry (i.e., it does not occur in chemical equilibrium for this planet) – this is the first direct evidence of photochemistry in an exoplanet atmosphere. There is also evidence for a high continuum level of the spectrum, due to an aerosol (cloud) deck preventing transmission to deeper levels.

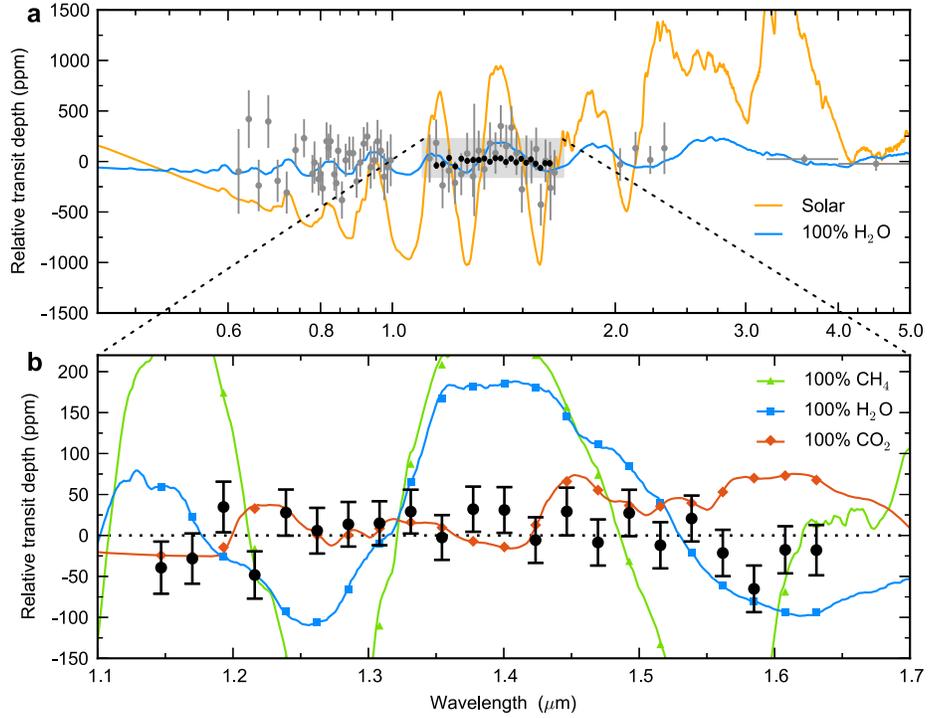


Figure 18.7: The transmission spectrum of the sub-Neptune GJ 1214b as measured through 15 transits with HST/WFC3. Figure adapted from Kreidberg et al. (2014).

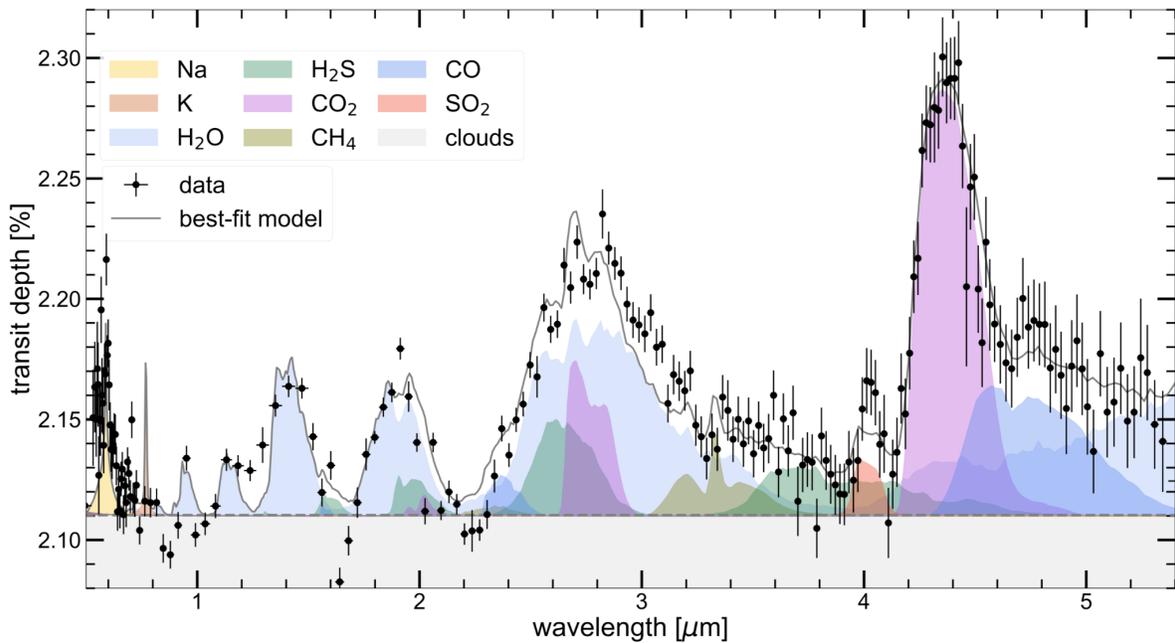


Figure 18.8: The transmission spectrum of the hot Jupiter WASP-39b observed with one transit from JWST NIRSpec/PRISM. There are clear detections of Na, H₂O, CO₂, CO, and SO₂, along with a cloud deck. Figure adapted from Rustamkulov et al. (2023).

19 Exoplanet interiors: terrestrial planets

Our agenda for Day 19 is the following:

1. Heat transfer via conduction (15 minutes)
2. Cooling timescale of Earth activity (25 minutes)
3. Heat transport in rocky planet interiors: Rayleigh-Bernard convection (25 minutes)
4. Exoplanet mass-radius relationships: dependence on composition (10 minutes)

Today's reading is the Sotin review chapter on terrestrial planet interiors. This will provide a comprehensive overview of expectations for the interior structure and heat transport of rocky exoplanets.

19.1 Earth's internal structure

A schematic of the interior structure of Earth is shown in Figure 19.1. Earth's solid

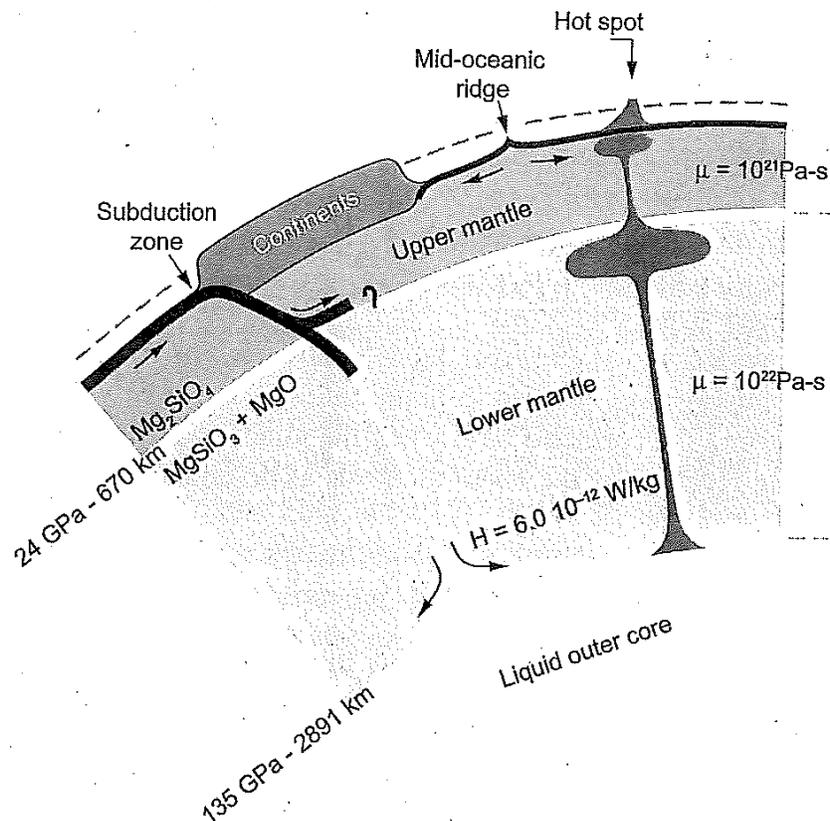


Figure 19.1: Schematic of the regions of Earth's interior that behave like a fluid, on which the continents and oceanic crust float isostatically. Figure adapted from Sotin et al. (2010).

component can be broken into five primary layers, from the surface inward: 1) crust, both continental and oceanic; 2) an upper mantle, comprised dominantly of olivine and enstatite;

3) the lower mantle, comprised dominantly of perovskite and periclase, with a transition to post-perovskite near the bottom of the mantle; 4) a liquid metal outer core; 5) a solid inner core. The relative extent and characteristic density profile of each of these layers are shown in Figure 19.2, which displays the Preliminary Reference Earth Model interior structure derived from seismology. Note the sharp transitions in density between each of

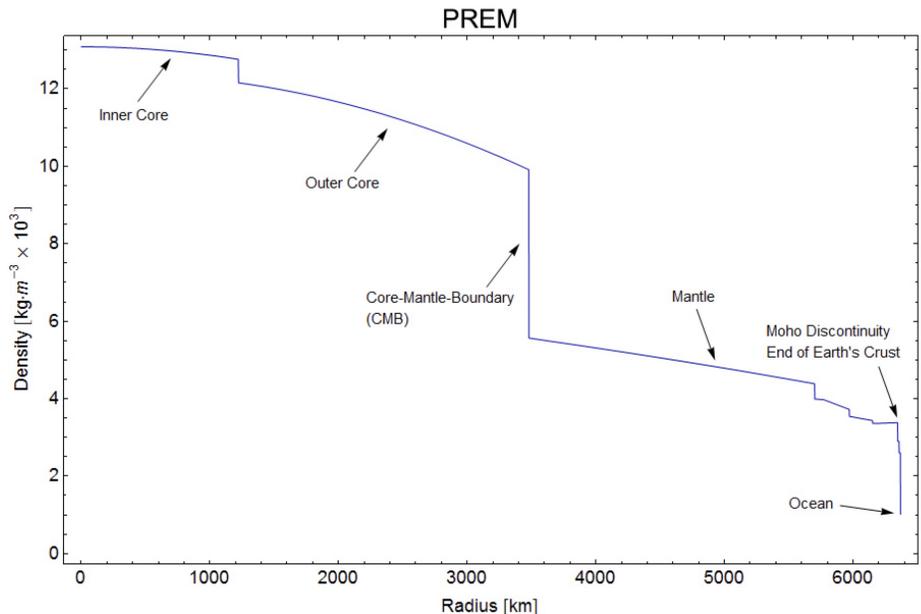


Figure 19.2: Preliminary Reference Earth Model (PREM) of the interior density distribution of Earth. Note the sharp discontinuities in density at the boundary between the solid inner core and fluid outer core, core-mantle boundary, and mantle-crust boundary (Mohorovičić discontinuity).

these layers, including the Mohorovičić discontinuity between the crust and mantle, CMB between mantle and core, and liquid-solid phase transition between the outer and inner core. These transitions occur because Earth’s interior is differentiated, a process which is expected to be ubiquitous for rocky planets due to gravitational segregation during the early stages of planet formation, as they form hot and cool over time.

19.2 Heat transfer via conduction

In solid rocks, conduction is the transfer of thermal energy by vibrations in the lattice of the solid material, with heat transported from regions of high temperature to regions of low temperature. Note that *no bulk movement* of material occurs during conduction, instead heat flows down-gradient from small-scale thermal interactions via a random-walk (i.e., diffusive) process. Fourier’s law of conduction states the the heat flux \mathbf{F} carried by conduction is related to the thermal conductivity k and temperature gradient ∇T as

$$\mathbf{F} = -k\nabla T, \quad (19.1)$$

where the thermal conductivity k has SI units of W/(mK), with typical values of ~ 1 W/(mK) for rock.

Heating via conduction can be related to the divergence of the conductive heat flux as

$$\begin{aligned}\rho c_p \frac{\partial T}{\partial t} &= -\nabla \cdot \mathbf{F} + Q, \\ \rho c_p \frac{\partial T}{\partial t} &= \nabla \cdot (k \nabla T) + Q,\end{aligned}\tag{19.2}$$

where in the latter expression we have substituted Fourier's law. Dividing through by ρc_p , assuming constant k , and expressing the thermal diffusivity

$$\kappa = \frac{k}{\rho c_p},\tag{19.3}$$

we can write a diffusion equation for heat conduction as

$$\frac{\partial T}{\partial t} = \kappa \nabla^2 T + \frac{Q}{\rho c_p}.\tag{19.4}$$

As a result, when ignoring sources, heat conduction is a purely diffusive process. The characteristic diffusion timescale of thermal conduction is governed by the thermal diffusivity

$$\tau \sim \frac{h^2}{\kappa},\tag{19.5}$$

where h is the thickness of the cooling structure. We can estimate the thermal diffusivity as a product of a velocity and a mean-free path

$$\kappa \sim vl,\tag{19.6}$$

where the speed of lattice vibrations (the diffusing quantity) is the speed of sound, and the mean free paths are comparable to the inter-atomic spacing.

19.2.1 Cooling timescale of Earth activity

Let's estimate the typical thermal diffusivity of rock to see how long it takes heat to escape Earth's interior by conduction through the oceanic crust. If we completely ignore heat production in the interior (see next section), in principle this could constrain the cooling age of the interior of Earth, but as you'll find there must be some other mechanism evicting heat from Earth's interior.

1. Assume that the speed of sound in rock in Earth's mantle is 4 km s^{-1} , and that the mean free path of lattice vibrations are 3 \AA . Estimate the thermal diffusivity of Earth mantle rock κ in m^2/s .
2. Earth's continental crust is $\sim 40 \text{ km}$ thick. Estimate the timescale for heat to diffuse through Earth's crust using your calculated thermal diffusivity. Compare this to the age of the Earth.
3. Now calculate the cooling timescale to transport heat from the center of Earth to the surface. Compare this to the age of the universe.

19.2.2 Historical background: Kelvin’s folly

In the 19th century, Lord Kelvin used conduction to estimate the age of the Earth, by calculating its characteristic Kelvin-Helmholtz cooling timescale from the heat flux that can be carried by conduction. He estimated the age by assuming that the Earth formed at a uniform hot temperature T_i and that its surface is maintained at a lower temperature T_0 , and that heat conduction was transferred through a thin near-surface boundary-layer at a thermal gradient $(dT/dz)_0$ that can be measured by studying the thermal gradient below the surface of Earth. Using these, he estimated the age of Earth as

$$t_0 = \frac{(T_i - T_0)^2}{\pi \kappa (dT/dz)_0^2} \sim 65 \text{ Myr.} \quad (19.7)$$

Of course, this turned out to be erroneous, as the age of Earth is ≈ 4.5 Gyr. This is because of two reasons: 1) a significant fraction of Earth’s internal heat budget is caused by radiogenic heat production (see Figure 19.3), which he did not know about at the time, 2) the interior of Earth transports heat by a combination of conduction and convection, which significantly changes the heat flux out of the interior. We’ll next turn to develop a more realistic model

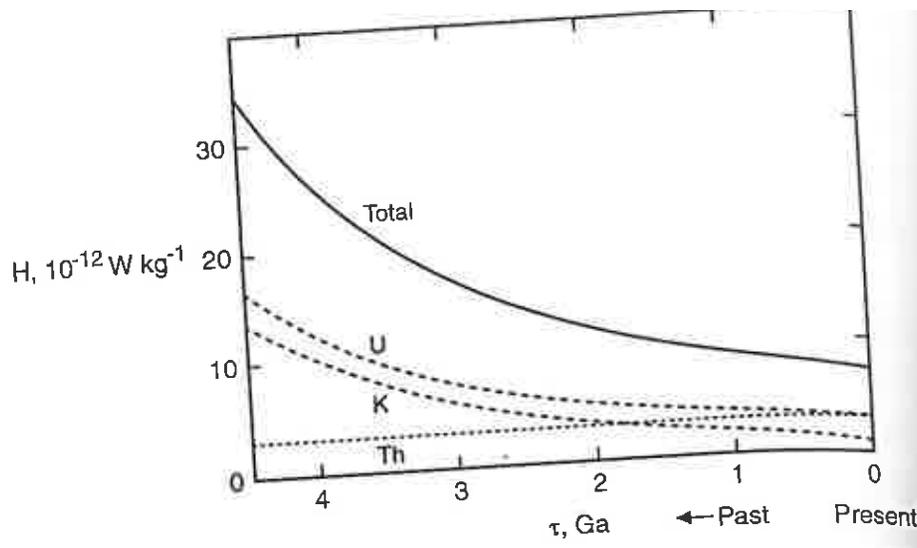


Figure 19.3: Total heat production rate in the mantle of Earth (solid line), along with contributions from uranium, potassium, and thorium (dashed lines) as a function of time before present. Figure poorly photocopied from Turcotte & Schubert (2002).

of the internal heat flux of Earth where heat is transported by convection in the mantle, but the heat evicted out of the interior is set by conduction through a “boundary layer” of crust.

19.3 Convective heat transport

19.3.1 Rayleigh-Bernard convection

Convection can occur in a fluid if it is sufficiently buoyant that the overturning timescale, τ_{over} due to buoyancy-driven motion is shorter than the timescale required for the heat to

thermally diffuse (conduct), τ_{diff} . That is, if

$$\begin{aligned}\tau_{\text{over}} &< \tau_{\text{diff}} \text{ convection,} \\ \tau_{\text{over}} &> \tau_{\text{diff}} \text{ conduction,}\end{aligned}\tag{19.8}$$

which is notably analogous to the Schwarzschild criterion for heat transport in gaseous planets (Equation 16.11), but here thermal conduction through lattice vibrations takes the role of diffusion of photons via radiative heat transport. The overturning timescale can be related to the ratio of the dynamic viscosity $\eta = \nu\rho$ (where ν is the kinematic viscosity, which we discussed previously when covering disks) and buoyancy as

$$\tau_{\text{over}} \sim \frac{\eta}{\Delta\rho gh},\tag{19.9}$$

where $\Delta\rho$ represents the difference in density over a height of thermal perturbations in the fluid h . The time to erase thermal anomalies by diffusion is

$$\tau_{\text{diff}} \sim \frac{h^2}{\kappa},\tag{19.10}$$

which you'll note is equivalent to the conduction timescale discussed previously. As a result, we can write that for thermal convection to occur

$$\frac{\eta}{\Delta\rho gh} < \frac{h^2}{\kappa}.\tag{19.11}$$

If we let $\Delta\rho = \alpha\rho\Delta T$ where α is the thermal expansivity and ΔT is the temperature drop across the height h , we can re-write the criterion for convection as

$$\frac{\alpha\rho g\Delta Th^3}{\eta\kappa} > 1 \text{ for convection.}\tag{19.12}$$

The expression on the left hand side is the definition of the non-dimensional fluid Rayleigh number Ra, which is the ratio of the diffusion to overturning timescale

$$\text{Ra} \equiv \frac{\alpha\rho g\Delta Th^3}{\eta\kappa}.\tag{19.13}$$

In reality, our simple analysis requires an additional numerical factor, such that convection only occurs for Rayleigh numbers above a critical value

$$\text{Ra} > \text{Ra}_{\text{crit}} \sim 10^3 \text{ for convection.}\tag{19.14}$$

We can further define one other non-dimensional number that will be important for characterizing convection, the Nusselt number. The Nusselt number is the ratio of heat flux to the heat flux that would be transported by conduction alone, i.e.,

$$\text{Nu} \equiv \frac{F}{F_{\text{cond}}} = h \frac{F}{k\Delta T},\tag{19.15}$$

the latter expression of which is valid if the thermal conductivity k is a constant. For Earth's mantle, $\text{Nu} \approx 30$.

19.3.2 Boundary layer convection

In reality, the heat flux out of the mantle to the surface of Earth is not evicted by convection, but via conduction through a thin boundary layer, Earth's crust. Similarly, heat transport from the core to the mantle is not regulated by convection, but through conduction through the bottom of the mantle. As a result, Earth's mantle undergoes Rayleigh-Bernard convection in a layer sandwiched between a hot thermal boundary layer at the bottom and a cold thermal boundary layer at the top, as shown in Figure 19.4.

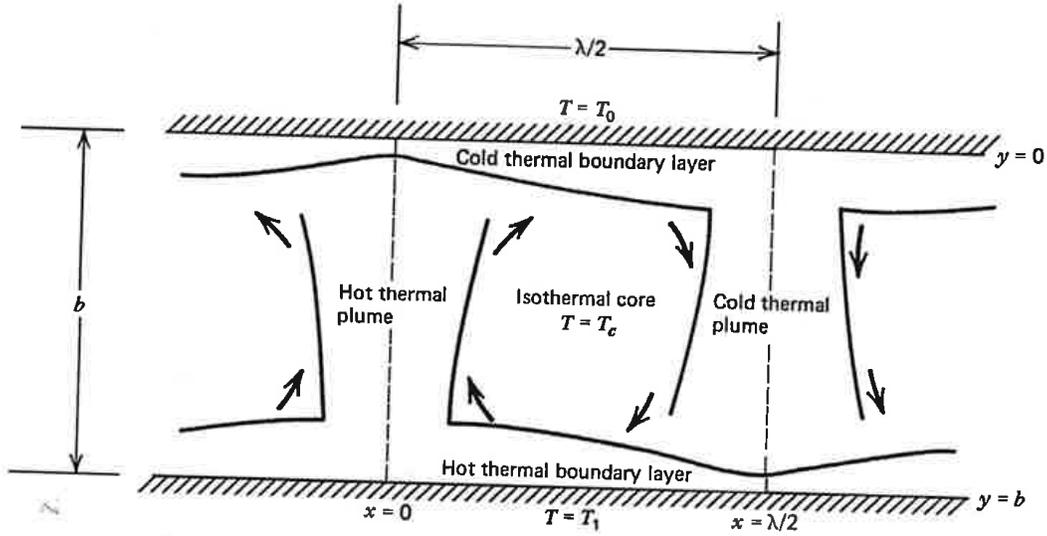


Figure 19.4: Schematic of boundary layer convection, with a (Rayleigh-Bernard) convective interior sandwiched between a hot bottom conductive boundary layer and a cold top conductive boundary layer. Figure adapted from Turcotte & Schubert (2002).

As a result, to order-of-magnitude, the flux transported out of the interior is that conducted across the top thermal boundary layer,

$$F \sim \frac{k\Delta T}{\delta}, \quad (19.16)$$

where δ is the thickness of the boundary layer. To determine the boundary layer thickness, we can estimate the maximum thickness over which it will not convect – i.e., where the diffusion time δ^2/κ will be less than the overturning time of the boundary layer $\eta/(\Delta\rho g\delta)$. This is the same as the criterion for Rayleigh-Bernard convection, but now with a local Rayleigh number required for boundary-layer peel away:

$$\begin{aligned} \text{Ra}_{\text{local}} &> \text{Ra}_{\text{crit}} \\ \frac{\alpha\rho g\Delta T\delta^3}{\eta\kappa} &> \text{Ra}_{\text{crit}}. \end{aligned} \quad (19.17)$$

We can then realize that boundary-layer peel away is also required for convection to occur, so the condition for boundary-layer peel away is also approximately the same criterion for

convection to occur in the first place. As a result, boundary-layer convection occurs when

$$\text{Ra}_{\text{crit}} < \frac{\alpha \rho g \Delta T \delta^3}{\eta \kappa} = \text{Ra} \frac{\delta^3}{h^3}. \quad (19.18)$$

We can re-arrange to solve for the boundary-layer thickness as

$$\delta \sim h \left(\frac{\text{Ra}_{\text{crit}}}{\text{Ra}} \right)^{1/3}, \quad (19.19)$$

the resulting heat flux transported through the boundary layer to the surface

$$F \sim k \frac{\Delta T}{\delta} \sim \frac{k \Delta T}{h} \left(\frac{\text{Ra}}{\text{Ra}_{\text{crit}}} \right)^{1/3}, \quad (19.20)$$

and the Nusselt number

$$\text{Nu} = \frac{F}{F_{\text{cond}}} \sim \left(\frac{\text{Ra}}{\text{Ra}_{\text{crit}}} \right)^{1/3}. \quad (19.21)$$

More generally, depending on the numerical simulation and details of the system of interest, a power-law relationship is found between the Nusselt number and Rayleigh number,

$$\text{Nu} = a \text{Ra}^\beta, \quad (19.22)$$

with $\beta \approx 0.1 - 0.4$. When coupling this boundary-layer convection model to a thermal evolution model (beyond the scope of this class, but see Komacek & Abbot, 2016 for an example including volatile cycling), the resulting temperature structure of the interior of a planet contains three regions: a hot and a cold boundary layer with strong temperature gradients, and a nearly isothermal adiabatic interior. This thermal structure that results from boundary-layer convection is shown in Figure 19.5, which provides a rough first-order model of the geotherm.

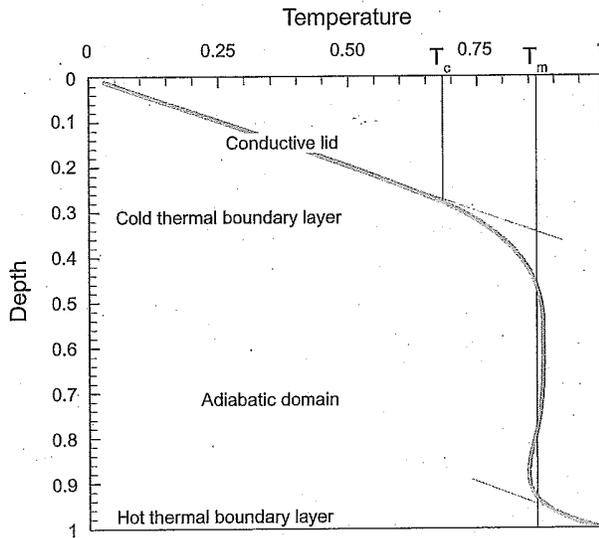


Figure 19.5: Temperature profile resulting from a boundary layer convection model. Note the small temperature gradient in the convective interior due to the relatively small adiabatic lapse rate compared to the conductive lapse rates in both the hot bottom and cold top boundaries. Figure adapted from Sotin et al. (2010).

19.4 Rocky planet mass-radius relationships

Due to the radius gap between sub-Neptunes and super-Earths that we discussed previously (Fulton et al., 2017), it is expected that planets with radii $\lesssim 1.6 R_{\oplus}$ are mostly rocky, while planets that have larger radii may be either rocky or volatile-rich (icy), and likely host a gaseous envelope. Figure 19.6 shows the mass-radius diagram of observed exoplanets with $M < 10 M_{\oplus}$, along with the terrestrial Solar System planets Mars, Venus, and Earth. Inter-

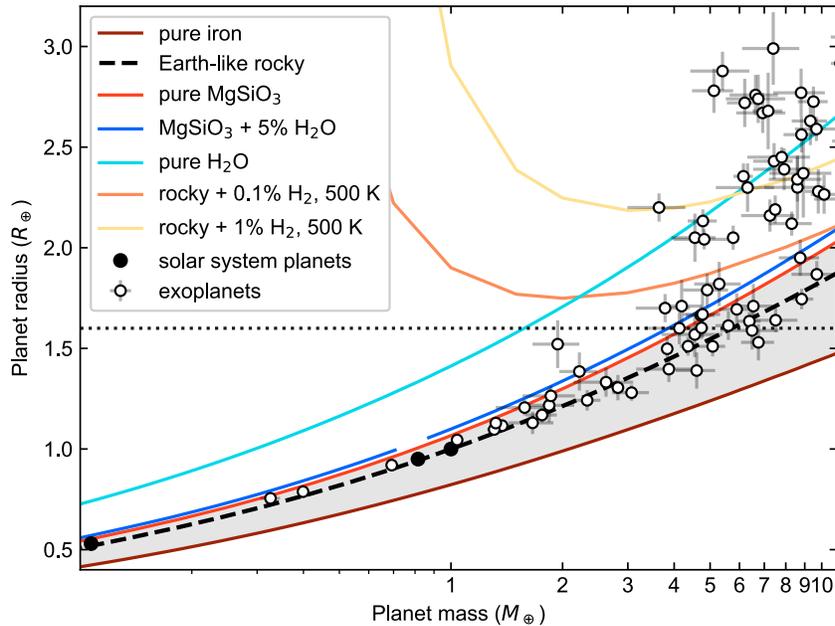


Figure 19.6: Mass and radius measurements for exoplanets (open circles) and Solar System planets (filled circles), compared with mass-radius curves for varying planetary compositions. Only exoplanets with 5σ mass measurements are plotted. The horizontal dotted line marks the radius gap above which planets generally retain significant hydrogen envelopes. Figure adapted from Wordsworth & Kreidberg (2022).

estingly, the vast majority of rocky planets with well-measured masses and radii less than $1.6 R_{\oplus}$ lie fairly close to an Earth-like compositional mass-radius relationship, with only a few planets to date requiring less metal or an envelope of ices to explain their radius. However, as we go to larger radii greater than $1.6 R_{\oplus}$ there are a diversity of planetary densities, with some (large super-Earths) still being well-explained by an Earth-like bulk composition and others (sub-Neptunes) being explained by either a pure water or icy composition (known as “waterworlds”) or a combination of a rocky core and a hydrogen envelope. As a result, the composition of sub-Neptune sized exoplanets is strongly degenerate, with many objects being equally well explained by either a waterworld or a rocky core with a hydrogen/helium envelope. Observational characterization is required (e.g., via transmission and emission spectroscopy with JWST) to determine the planetary atmospheric composition and thus better constrain the bulk composition in order to test these two compositional hypotheses to determine if one is dominant or if sub-Neptunes have a bi-modal distribution of bulk compositions.

20 Exoplanet characterization: emission spectroscopy

Our agenda for Day 20 is the following:

1. Secondary eclipse depth (10 minutes)
2. Thermal structure, formation of absorption and emission features (30 minutes)
3. Example: WASP-18b with JWST (5 minutes)
4. Emission spectra activity on WASP-18b with JWST (30 minutes) – be sure to bring your computer and download the .zip file from ELMS in advance!

Today’s reading is the Deming review paper on how to characterize the atmosphere of a transiting exoplanet. This will both introduce the concept of emission spectroscopy via secondary eclipse measurements as well as how to combine measurements of emission and transmission spectroscopy to better constrain the atmospheric properties of exoplanets.

20.1 Secondary eclipse depth

The secondary eclipse occurs when the planet is occulted by the star. By measuring the differential flux between when the planet is not occulted by the star and we see $F = F_\star + F_p$ and when the planet is occulted and we only see F_\star , we can measure the planetary flux alone. The eclipse depth is then

$$\delta_{\text{ecl}} = \frac{F_p}{F_p + F_s} \approx \frac{F_p}{F_\star}, \quad (20.1)$$

where we have assumed that the planet is much dimmer than the star in the wavelength region of interest. Given that the flux of the planet and star is $F = B(T, \lambda)R^2$, where $B(T, \lambda)$ is the Planck function, we can express the secondary eclipse depth as

$$\delta_{\text{ecl}} \approx \frac{F_p}{F_\star} = \frac{B_p(T, \lambda)R_p^2}{B_\star(T, \lambda)R_\star^2} = \frac{R_p^2 e^{hc/\lambda k T_\star} - 1}{R_\star^2 e^{hc/\lambda k T_p} - 1}. \quad (20.2)$$

Note that the ratio of exponentials is always < 1 (because stars are hotter than their planets), so the secondary eclipse is always smaller than the primary eclipse. At long wavelengths in the Rayleigh-Jeans tail, we can approximate the eclipse depth as $F_p/F_\star \approx T_p/T_\star R_p^2/R_\star^2$.

First, note that the eclipse depth is a wavelength-dependent differential measurement that depends very strongly on the relative radii and temperatures of the planet and star. Thus, eclipse depths will be much larger for hotter and larger planets – for hot Jupiters, they can be on the order of percent, but for cooler planets they can be on the order of a few ppm. Additionally, F_p/F_s typically increases with wavelength due to the shorter wavelength peak of the stellar blackbody, implying that longer infrared wavelengths are preferred to search for secondary eclipses.

The first measured secondary eclipse was by UMD faculty member Drake Deming in 2005 for the planet HD 209458b (Deming et al., 2005, see Figure 20.1), coincident with the detection of a secondary eclipse for TrES-1b (Charbonneau et al., 2005). Both papers even have the same name – “Detection of Thermal Emission from an Extrasolar Planet”! Since then, wavelength-dependent eclipse depths (or “emission spectra”) have been used to characterize planets with a range of observatories.

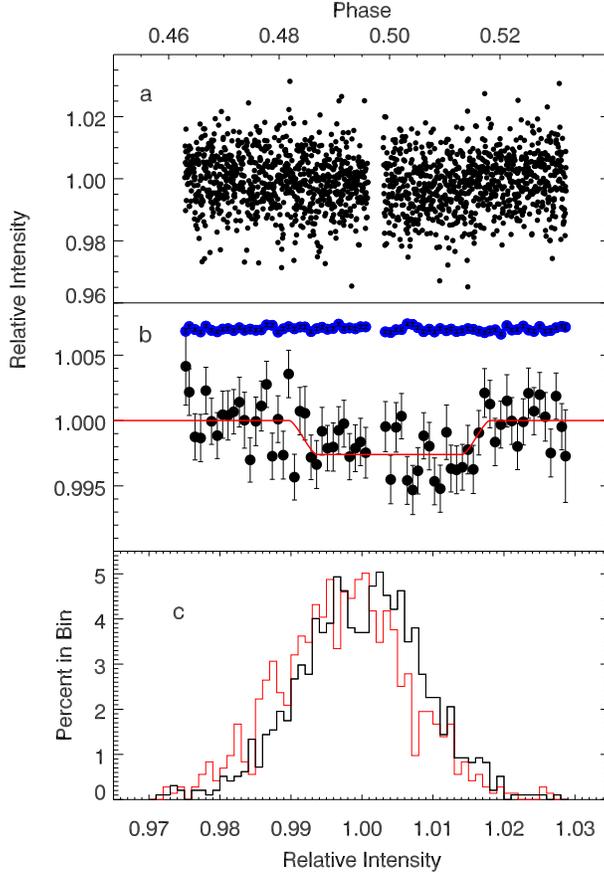


Figure 20.1: Detection of a secondary eclipse on HD 209458b with Spitzer at 24 μm by Deming et al. (2005). The top panel is the raw data, the middle panel shows binned data with best-fit eclipse model, and the bottom panel shows histograms of the binned intensity.

20.2 Linking emission to thermal structure

20.2.1 Solutions of the radiative transfer equations

Recall from our transmission spectra notes that we can write Schwarzschild's equation for radiative transfer ignoring scattering as

$$\frac{dI_\lambda}{dl} = \rho k_\lambda (B_\lambda - I_\lambda) \text{ Schwarzschild's equation} \quad (20.3)$$

in length coordinates, or equivalently in optical path coordinates as

$$\frac{dI_\lambda}{d\tilde{\tau}} = I_\lambda - B_\lambda, \quad (20.4)$$

or optical depth coordinates as

$$\mu \frac{dI_\lambda}{d\tau} = I_\lambda - B_\lambda, \quad (20.5)$$

where $\mu = \cos\theta$, with θ the zenith angle between the light path and the vertical in the atmosphere. Moving forward, we'll use optical path coordinates and drop the tilde for

simplicity. We previously considered atmospheres with no emission to derive Beer's law for transmission, but now we will consider atmospheres with emission.

First, let us consider an isothermal medium where B_λ is constant. We can change variables to simplify our derivation, using

$$\begin{aligned}\chi &= I_\lambda - B_\lambda \\ d\chi &= dI_\lambda.\end{aligned}\tag{20.6}$$

Thus, we can write Schwarzschild's equation with the change of variable as

$$\frac{d\chi}{d\tau} = -\chi.\tag{20.7}$$

Integrating, we find

$$\chi = \chi_0 e^{-\tau}.\tag{20.8}$$

Plugging back in for our change of variables, we find the isothermal solution of Schwarzschild's equation:

$$I_\lambda = B_\lambda + (I_\lambda(0) - B_\lambda) e^{-\tau}.\tag{20.9}$$

Let's consider some limits of this isothermal solution. When the optical path is very small ($\tau \rightarrow 0$), the radiance is the same before and after entering the slab. When the optical path is very large ($\tau = \infty$), $I_\lambda = B_\lambda$ – that is, an infinitely opaque isothermal medium acts as a blackbody. Finally, if $I_\lambda(0) = 0$, then $I_\lambda = B_\lambda(1 - e^{-\tau})$, and the radiance e-folds toward a blackbody with increasing optical thickness.

There is also a general solution to Schwarzschild's equation ignoring scattering. To derive it, we can use an integrating factor e^τ such that

$$\frac{d}{d\tau}(e^\tau I_\lambda) = e^\tau B_\lambda,\tag{20.10}$$

and thus

$$e^\tau I_\lambda|_0^\tau = \int_0^\tau e^{\tau'} B_\lambda d\tau'.\tag{20.11}$$

One can then rearrange to find the general solution to Schwarzschild's equation,

$$I_\lambda(\tau) = I_\lambda(0)e^{-\tau} + \int_{\tau=0}^\tau e^{\tau'-\tau} B_\lambda d\tau'.\tag{20.12}$$

Note that the first half of the right hand side is simply Beer's law, and the second half is the additional contribution from thermal emission.

20.2.2 Photosphere pressure

It's clear from the above derivations that the optical depth (equivalent to the optical path if $\mu = 1$) is a critical parameter for estimating radiative properties of planetary atmospheres. If the optical depth $\tau \gg 1$, then the region is optically thick and $I_\lambda \approx B_\lambda$. If the optical depth is $\tau \ll 1$, the region is optically thin and $I_\lambda \approx I_\lambda(0)$. However, if $\tau \approx 1$, then both incoming starlight is absorbed and radiation escapes to space – this special region of the atmosphere is called the “photosphere.”

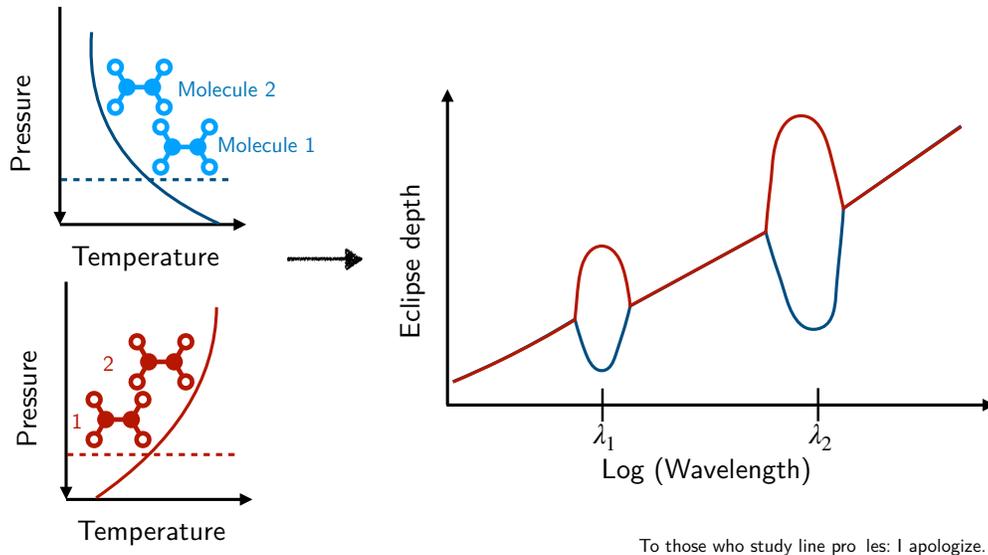


Figure 20.2: Cartoon demonstrating how absorption and emission features arise from non-inverted and inverted temperature pressure profiles, respectively.

We can estimate the photosphere using the definition of optical depth,

$$\tau \approx 1 = \int_z^\infty \kappa_\lambda \rho dz = \kappa_\lambda \frac{p}{g}. \quad (20.13)$$

If we assume that the opacity is constant and use hydrostatic equilibrium ($-\rho dz = dp/g$), we can estimate the photosphere pressure

$$p_{\tau=1} \approx \frac{g}{\kappa_\lambda}. \quad (20.14)$$

For a typical hot Jupiter with $g = 10 \text{ m s}^{-2}$ and $\kappa_{\text{IR}} = 4 \times 10^{-3} \text{ cm}^2 \text{ g}^{-1}$, $p_{\text{IR}} \approx 0.25 \text{ bars}$. Typically the photosphere is at the $\sim 100 \text{ mbar} - 1 \text{ bar}$ pressure level for irradiated planets.

20.2.3 Absorption and emission features

The temperature-pressure profile can be directly probed by measuring spectral features in secondary eclipse spectra. Figure 20.2 demonstrates how the temperature-pressure profile is coupled to spectral features. In non-inverted atmospheres where the temperature decreases with increasing pressure, spectral features appear as dips in the eclipse depth at the wavelengths where the opacity from absorbers is large. This is because the radiatively active molecules absorb the thermal emission from below and re-emit it at their (cooler) local temperature, causing a reduction in the planetary flux at that wavelength. In this case, molecule 2 has a larger relative dip (i.e., a larger absorption feature) because it is optically thick at lower pressures than molecule 1.

Conversely, in thermally inverted atmospheres there is a bump, or emission feature, rather than an absorption feature. This is because the molecules now are re-emitting the radiation they absorb at hotter local temperatures, causing an increase in the planet flux in the wavelengths at which they are opaque. Thus, one can estimate the temperature-pressure profile and chemical abundances together from an observed secondary eclipse observation, as you'll see in our activity.

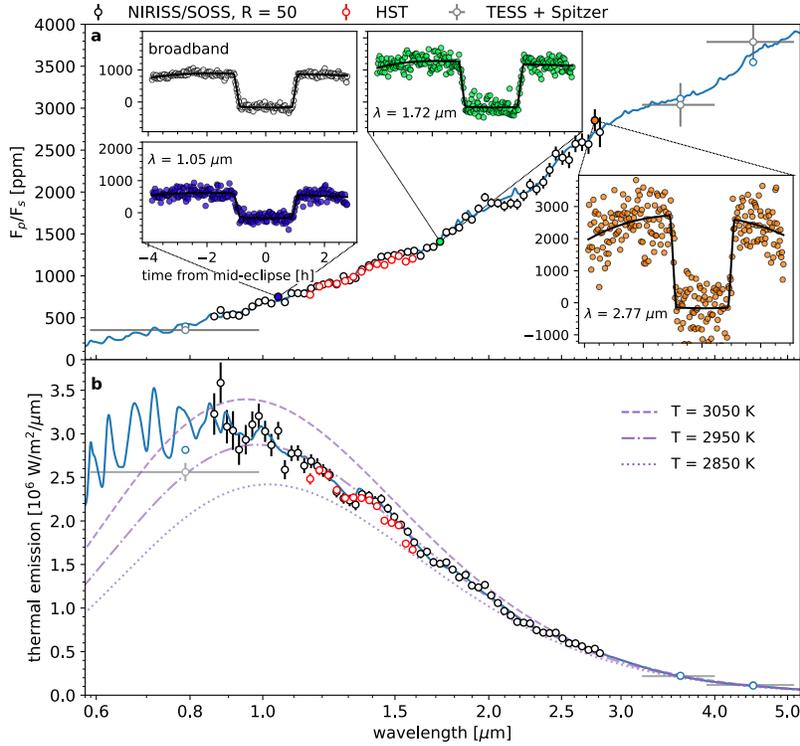


Figure 20.3: Top: secondary eclipse time-series at various wavelengths and secondary eclipse spectrum in F_p/F_* for WASP-18b as measured by JWST NIRISS/SOSS, TESS, and Spitzer. Bottom: The thermal emission from WASP-18b alone, using a PHOENIX stellar model to convert F_p/F_* to F_p . Figure from Coulombe et al. (2023).

20.3 Emission spectra in practice: WASP-18b with JWST

Before we dive in to an activity with real JWST data, let's briefly walk through an example of the type of data you'll be comparing thermal emission models to. Figures 20.3 - 20.5 show the secondary eclipse observations and spectrum, brightness temperature spectrum, and retrieval results for WASP-18b as observed with JWST NIRISS/SOSS by Coulombe et al. (2023). First, note that the spectrum in Figure 20.3 looks by eye nearly featureless. This was already expected from previous HST and Spitzer observations, as WASP-18b is an ultra-hot Jupiter with a dayside temperature of ~ 2900 K. At these high temperatures, molecules begin to thermally dissociate, and thus there will be a reduced amount of infrared opacity compared to cooler planets. When multiplying out a model stellar spectrum to isolate the planet thermal emission, some small features can be discerned, but it's clear they only correspond to thermal variations on the order of ~ 100 K in brightness temperature.

Figure 20.4 shows a spectrum of WASP-18b in brightness temperature rather than flux (i.e., inverting for the blackbody temperature at each wavelength) in order to more clearly display the small wavelength-dependent variations in the spectrum. Given the dominant absorbers overlaid to guide the eye, it is clear that water vapor spectral features are dominating the spectral variation in the near-infrared. There is additionally an increase in the brightness temperature going to shorter wavelengths due to a combination of optical absorbers,

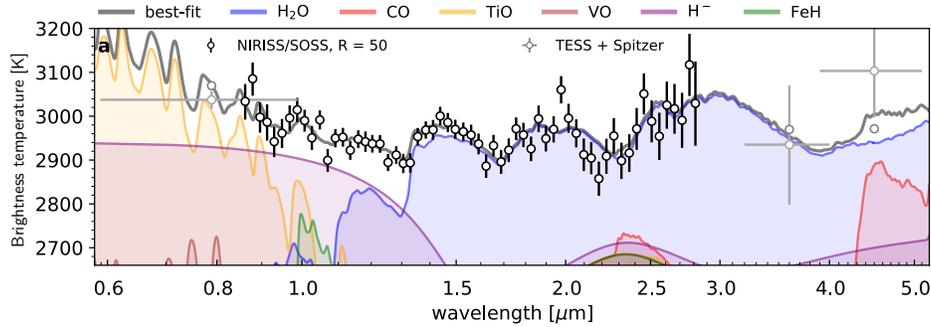


Figure 20.4: Brightness temperature spectrum of WASP-18b with key absorbers shown in the shaded regions. Figure adapted from Coulombe et al. (2023).

including TiO, VO, FeH, and H^- continuum opacity.

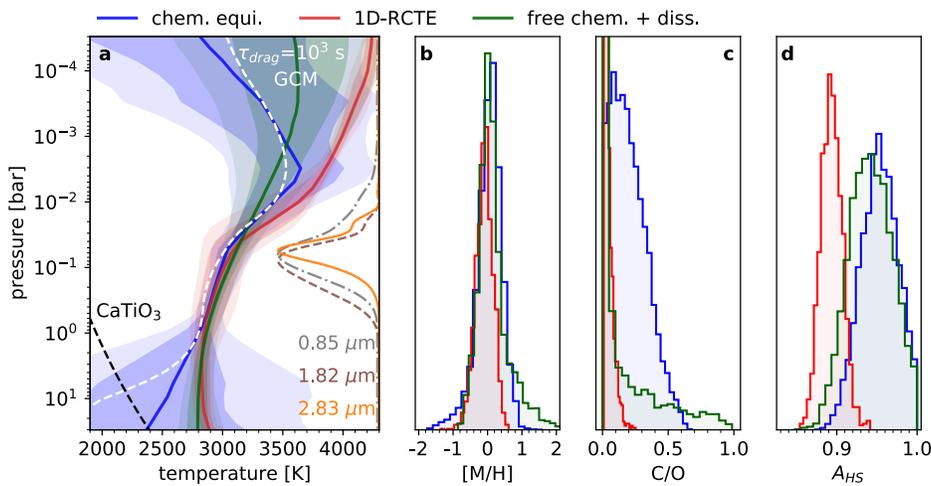


Figure 20.5: Retrieved temperature-pressure profiles (a) compared to condensation curves and GCM results, retrieved metallicity (b), retrieved C/O ratio (c), and retrieved dilution parameter (d) for WASP-18b. Figure from Coulombe et al. (2023).

The spectral features in Figure 20.4 correspond to an increase in temperature, implying that they are emission rather than absorption features. The retrieved temperature-pressure profile is shown in Figure 20.5, and compared to 3D GCM predictions, both of which indicate that the atmosphere is thermally inverted at the photosphere. In addition, the retrievals enable constraints on the atmospheric metallicity and C/O ratio, the former of which is roughly Solar.

20.4 Emission spectra activity!

We'll now do an activity using a Jupyter notebook prepared by Dr. Anjali Piette that will walk you through fitting model emission spectra to the WASP-18b JWST NIRISS/SOSS dataset. Please get in small groups of 2-3 and open the Jupyter notebook on your laptop and walk through the notebook. I'll walk around the room to make sure everyone is able to make progress, and we'll check in and compare our answers at the end.

21 Exoplanet characterization: phase curves

Our agenda for Day 21 is the following:

1. Introduction to phase curves (15 minutes)
2. Overview of the state of phase curve observations (15 minutes)
3. Interpreting phase curves with analytic theory (25 minutes)
4. Activity: predicting day-night temperature contrasts (20 minutes)

Today's reading is Sections 4.3-4.4 of the Zhang review article on exoplanet and brown dwarf atmospheres. These sections will explain the current state of observations and theoretical interpretation of orbital phase curves of exoplanets as well as rotational phase curves of brown dwarfs. You may also want to read the Heng & Showman review on exoplanet phase curves.

21.1 Phase curve fundamentals

Phase curves are measurements of the light from an object as a function of either orbital phase or rotational phase. The light measured can include both thermal emitted light as well as reflected light from a companion object (e.g., a star). Phase curves measured at present for exoplanets are orbital phase curves, where the planet flux is a small modulation onto the (nearly constant) stellar flux. Rotational phase curves, meanwhile, can be presently measured for brown dwarfs and wide-separation giant planets. Phase curves critically provide a light curve measurement of the object's flux in the time-domain, which in turn can be translated to make a (crude) map of the brightness of the planet as a function of longitude, latitude, and/or pressure (where the latter can be inferred only if the measurement is spectroscopic).

21.1.1 Orbital phase curves: close-in exoplanets

Figure 21.1 shows an example phase curve taken by the Spitzer Space Telescope of the hot Jupiter HD 189733b, which is a photometric observation centered at a wavelength of $3.6 \mu\text{m}$. Labeled on the top half of this diagram (which shows the full y-scale) are the secondary eclipse and transit where the planet is occulted by and occults the star, respectively. If you look closely, you can see small changes in the total system flux – the star does not vary significantly over the 2.2 day orbital period of the planet, so these variations are due to the planet. The bottom zooms in to show the effect of the variation in thermal emission from the planet HD 189733b on the phase curve. To first order, this variation is sinusoidal, and in fact phase curves are often fit with a double sinusoid in the literature.

Phase curves are often characterized by two key features: their offset (i.e., time or phase shift in the peak flux from the time of the center of secondary eclipse), and their amplitude (i.e., relative difference between maximum and minimum flux). Phase curve offsets are usually measured in degrees of orbital phase. The sign of the phase curve offset is such that if the phase curve peaks before secondary eclipse, it is positive, and if it peaks after secondary eclipse, the phase curve offset is negative. This convention is chosen such that the phase curve offset matches with the sign of the brightness spot offset in longitude, assuming that the planet is tidally locked to its host star, such that a positive phase curve offset

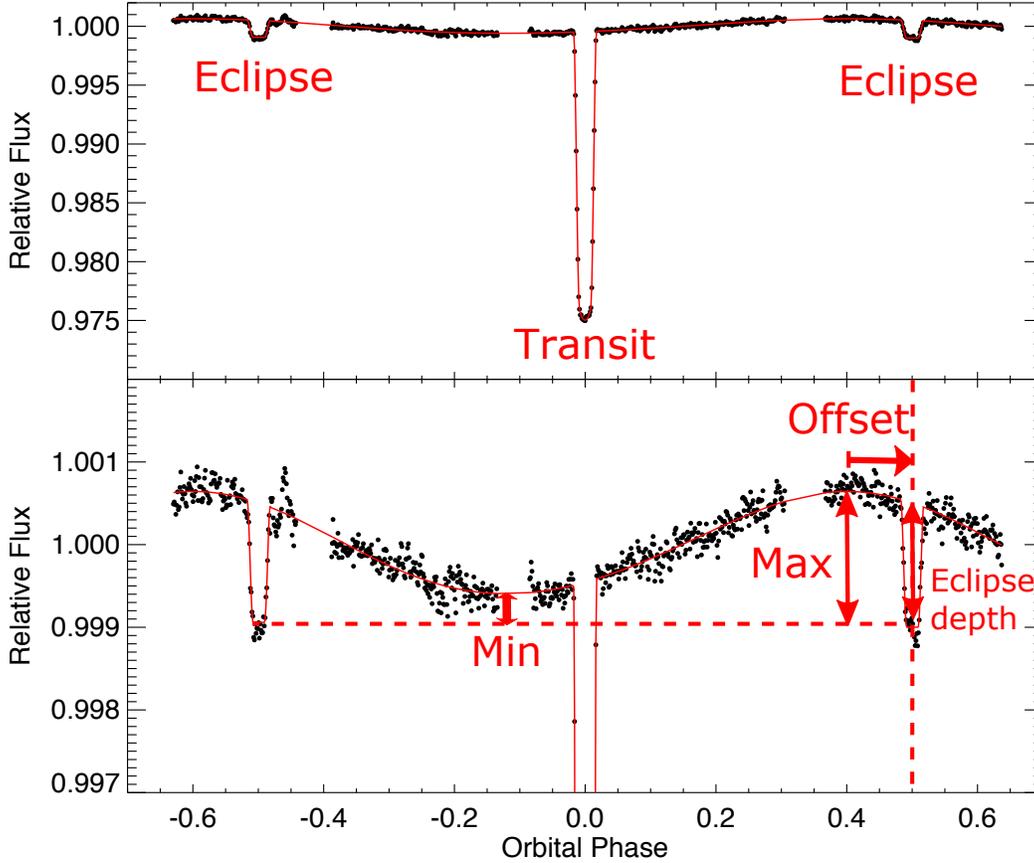


Figure 21.1: Example orbital phase curve observation of the hot Jupiter HD 189733b with the Spitzer Space Telescope at $3.6 \mu\text{m}$. Figure adapted from Parmentier & Crossfield (2018).

corresponds to an eastward bright spot shift and a negative phase curve offset corresponds to a westward bright spot shift. Phase curve amplitudes are usually measured as

$$A_F = \frac{(F_{p,\text{max}} - F_{p,\text{min}})}{F_{p,\text{max}}}, \quad (21.1)$$

where $F_{p,\text{max}}$ and $F_{p,\text{min}}$ are the maximum and minimum flux, respectively. As a result, the phase curve amplitude is $A_F = 1$ when the minimum flux goes to zero (e.g., in the case of no day-night heat transport), and $A_F = 0$ when the planet has a uniform brightness distribution.

A broad range of information about the planet’s thermal structure (and indirectly, the atmospheric circulation) can be inferred from phase curve observations of exoplanets and other complimentary observations. Figure 21.2 shows a summary of the types of measurements that can be made at present for hot Jupiters. From population-level measurements of dayside and nightside flux, one can then translate these fluxes at specific wavelengths into brightness temperatures to constrain the day-to-night temperature contrast and how it depends on planetary parameters, for instance equilibrium temperature (see also Figure 21.6). From spectroscopic phase curve measurements (e.g., with HST or JWST), the changing shapes of spectral features can provide constraints on the temperature profiles as a function

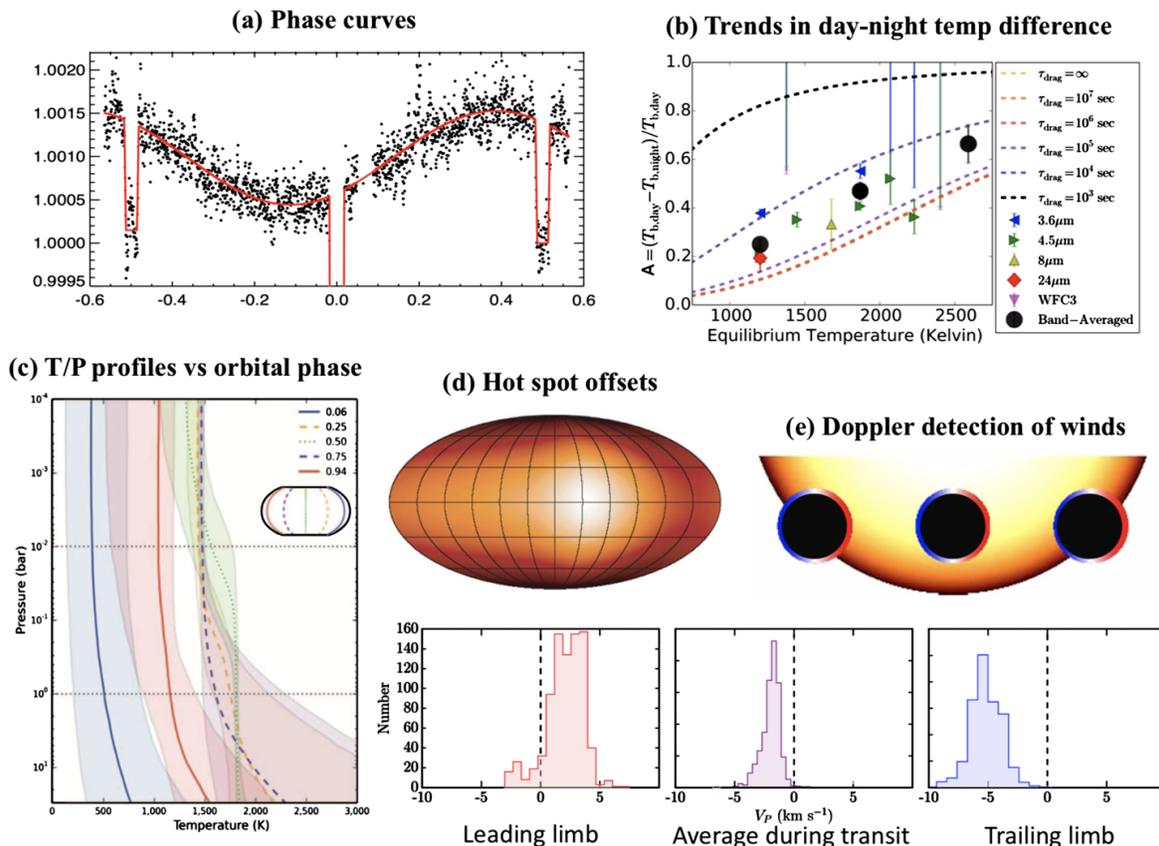


Figure 21.2: Examples of the wealth of information that can be inferred from observations of hot Jupiters, largely with phase curves. These include the day-night brightness temperature contrast, temperature profiles as a function of orbital phase, and hot spot offsets. With high spectral resolution observations of transit, eclipse, or phase variations, Doppler shifts due to winds can also be inferred. Figure adapted from Showman et al. (2020).

of orbital phase. Measurements of phase curve offsets can be translated to brightness temperature offsets in longitude assuming tidal locking, which then allows for an inference to be made on wind patterns that cause this phase curve offset. Finally, by observing phase curves and transmission spectra in high spectral resolution (which requires instrumentation on ground-based telescopes), the wind speeds and patterns in hot exoplanet atmospheres can be inferred via Doppler shifts of spectral lines. Note that wind speeds and rotation are degenerate with one another, so wind speeds can only be inferred under the assumption that the planet is tidally locked (or if the rotation speed is otherwise known).

21.1.2 Contribution from reflected light

Though the majority of phase curves of exoplanets are measured in the infrared with HST, Spitzer, and JWST, and thus largely probe thermal emission, phase curve observations in visible wavelengths with Kepler, TESS, and CHEOPS have a significant reflected light component. As shown in Figure 21.3, numerical models of tidally locked hot Jupiters predict that the emitted light component will lead to a phase shift that is opposite of that

caused by reflected light. This is because the reflected light contribution is dominated by

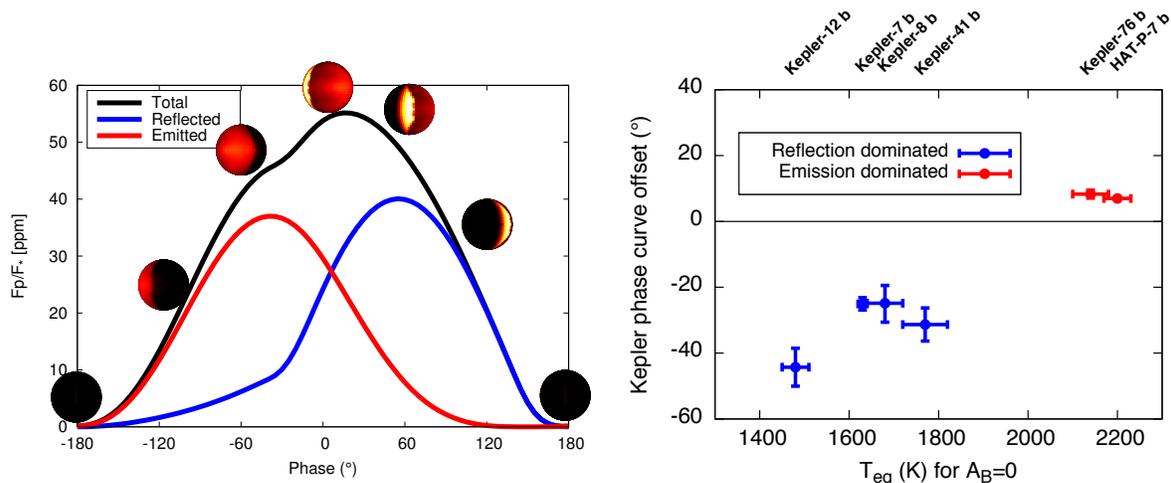


Figure 21.3: Simulated phase curve including both thermal emission and reflected light contributions (left), and observed Kepler phase curve offsets (right). There appears to be a transition in Kepler observations from reflected light setting the phase curve offset at low temperature to thermal emission dominating at higher temperatures. Figure adapted from Parmentier & Crossfield (2018).

scattering, which has a greater contribution from cloudy regions that are at lower temperatures. Meanwhile, the thermal emitted light contribution is larger at clearer, hotter regions, which generally occur on the eastern hemisphere of tidally locked planets. As a result, for tidally locked hot Jupiters, the reflected light component is expected to cause westward bright spot offsets while thermal emission will cause eastward bright spot offsets.

The relative contribution of reflected light vs. thermal emission varies with planetary temperature, as hotter planets have a greater emitted component that will begin to dominate over reflected light (especially in the visible due to the shift of the Planck function to shorter wavelengths). Thus, theoretical predictions expect westward phase curve offsets in the visible for cooler hot Jupiters (due to reflected light) that transition to eastward offsets for hotter planets (due to thermal emission). This is exactly what has been observed in Kepler phase curves, as shown in the right-hand panel of Figure 21.3.

21.1.3 Rotational phase curves: brown dwarfs, wide-separation giant planets

Light curve variability over rotational phase has been measured for many brown dwarfs, both with ground- and space-based observations. Two examples of measured rotational phase curve variability are shown in Figure 21.4. The left-hand panel shows a light curve of SIMP0136, displaying very short-timescale variability that is likely due to the inherent patchiness of the atmosphere that changes due to rotation changing the observable hemisphere with time. The patchiness in brown dwarf atmospheres is largely expected to be due to a cloud coverage and the resulting effective temperature variations (Tan & Showman, 2021) – as a result, the changing light curve with time is likely due to changes in the cloud cover and/or atmospheric circulation of the object. The right-hand panel shows inferred surface maps of the closest brown dwarf to Earth, Luhman 16B, via Doppler imaging. This

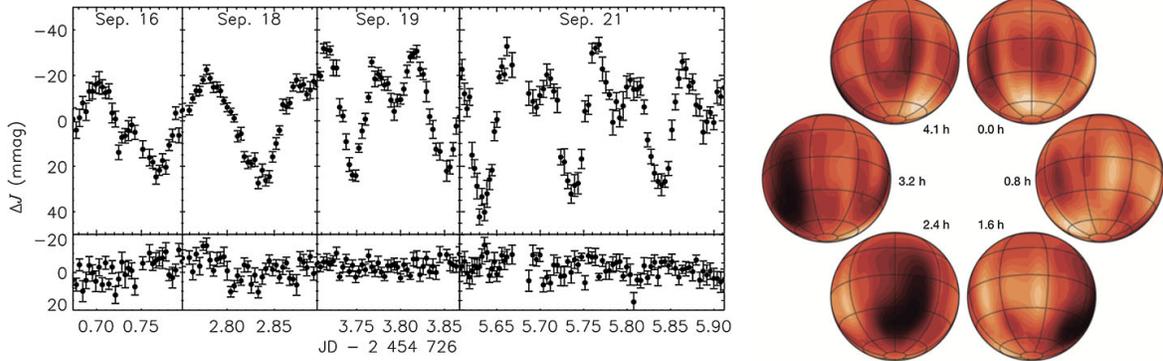


Figure 21.4: Example of a brown dwarf rotational light curve (left, $1.2 \mu\text{m}$ observations of SIMP0136) and using rotational phase variations to infer surface patchiness on a brown dwarf (right, Doppler imaging observations of Luhman 16B). Figure adapted from Showman et al. (2020).

technique shows that the surface is patchy, with clear variations as a function of longitude in the circulation of the object. TESS observations of Luhman 16B have similarly found variability due to both rotation and atmospheric circulation, along with a long-term trend in the brightness of the brown dwarf (Apai et al., 2021).

21.2 Phase curve theory for tidally locked exoplanets

In general, phase curves require detailed three-dimensional atmospheric circulation models, often termed General Circulation Models or GCMs, in order to interpret fully. We will discuss and do an activity to study the detailed output of GCMs in the next class. In order to build intuition, in this class we'll derive a simple first-principles scaling theory for the characteristic day-night temperature contrast and wind speeds of tidally locked gaseous planets.

21.2.1 A simple coupled scaling theory for heat transport and winds

Let's now derive simple analytic predictions for the day-night temperature contrast and wind speeds of hot Jupiters, following a simplified version of the derivations in Perez-Becker & Showman (2013); Komacek & Showman (2016); Zhang & Showman (2017); Zhang (2020). To do so, we'll start by scaling the equations of momentum conservation, hydrostatic equilibrium, continuity equation, and energy conservation. We'll then solve them in one limit, geostrophic balance, which corresponds to when the pressure gradient and Coriolis terms balance in the momentum equation. The full solutions are in Equations (13) and (14) of the Zhang reading.

First, the conservation of momentum can be written

$$\frac{d\mathbf{v}}{dt} = -\frac{\nabla p}{\rho} + \mathbf{g} - 2\boldsymbol{\Omega} \times \mathbf{v} + \mathcal{F}, \quad (21.2)$$

where the first term on the right is the pressure gradient force, the second gravity, the third Coriolis force, and the last one additional dissipation (e.g., frictional drag). Let's scale the horizontal component of each term individually, noting that there is no horizontal

contribution from gravity. First, the horizontal advection term in steady-state can be scaled as

$$\frac{d\mathbf{v}}{dt} = \frac{\partial\mathbf{v}}{\partial t} + \mathbf{v} \cdot \nabla\mathbf{v} \sim \frac{U^2}{L}, \quad (21.3)$$

where U and L are characteristic velocity and length scales, respectively. Next, using the ideal gas law, the pressure gradient term can be written and scaled as

$$-\frac{\nabla p}{\rho} \approx \frac{\rho R \Delta T}{\rho L} = \frac{R \Delta T}{L}, \quad (21.4)$$

where ΔT is a characteristic planetary-scale (day-night) temperature contrast. Finally, the Coriolis force can simply be scaled as

$$2\Omega \times \mathbf{v} \sim \Omega U. \quad (21.5)$$

In geostrophic balance, the Coriolis force and pressure gradient balance, and thus our final scaled momentum equation in geostrophy is

$$\begin{aligned} \frac{R \Delta T}{L} &= \Omega U \\ \Delta T &= \frac{\Omega U L}{R}. \end{aligned} \quad (21.6)$$

This is one equation for two unknowns: ΔT and U . Thus, we need an additional constraint, from energy conservation. This is just the first law of thermodynamics,

$$c_p \frac{dT}{dt} = q + \frac{1}{\rho} \frac{dp}{dt}, \quad (21.7)$$

which we can alternately write as

$$\frac{dT}{dt} - \frac{\omega}{\rho c_p} = \frac{q}{c_p}, \quad (21.8)$$

where note that the vertical pressure velocity $\omega = dp/dt$. We can then make the ansatz that the heating/cooling rate can be prescribed as a linear relaxation of the temperature back to an equilibrium value (T_{eq}) over the radiative timescale τ_{rad} ,

$$\frac{q}{c_p} = -\frac{(T - T_{\text{eq}})}{\tau_{\text{rad}}}, \quad (21.9)$$

and further note that the material derivative of temperature can be written as

$$\frac{dT}{dt} = \frac{\partial T}{\partial t} + \mathbf{v} \cdot \nabla T + \omega \frac{\partial T}{\partial p}. \quad (21.10)$$

In the limit where we ignore vertical motions (a vast simplification from Komacek & Showman, 2016, but congruent with our simplifications of the momentum equation), and assume a steady state, we can then write the thermodynamic energy equation as

$$\mathbf{v} \cdot \nabla T = \frac{(T_{\text{eq}} - T)}{\tau_{\text{rad}}}. \quad (21.11)$$

Further assuming that $T_{\text{eq}} - T \approx \Delta T_{\text{eq}} - \Delta T$, where ΔT is the planetary-scale (day-night) temperature contrast and ΔT_{eq} is the temperature contrast in radiative equilibrium, then we can scale this equation to relate the wind speeds and day-night temperature contrast as

$$\begin{aligned} \frac{U\Delta T}{L} &\sim \frac{(\Delta T_{\text{eq}} - \Delta T)}{\tau_{\text{rad}}} \\ \frac{\Delta T}{\Delta T_{\text{eq}}} &\sim \left(1 + \frac{U\tau_{\text{rad}}}{L}\right)^{-1} \sim \left(1 + \frac{\tau_{\text{rad}}}{\tau_{\text{dyn}}}\right)^{-1}. \end{aligned} \quad (21.12)$$

In the equation above,

$$\tau_{\text{dyn}} \sim \frac{L}{U}, \quad (21.13)$$

where $L \approx a$ is the characteristic length scale of the circulation and U is the characteristic wind speed. Combining with Equation (21.6), we can separately solve to find a quadratic expression for ΔT

$$\frac{R\tau_{\text{rad}}\Delta T^2}{\Omega L^2} + \Delta T - \Delta T_{\text{eq}} \sim 0. \quad (21.14)$$

In the non-linear regime where $\frac{R\tau_{\text{rad}}}{\Omega L^2} \Delta T_{\text{eq}} \gg 1$, we can approximate

$$\Delta T \sim \sqrt{\frac{\Omega L^2 \Delta T_{\text{eq}}}{R\tau_{\text{rad}}}}. \quad (21.15)$$

We can then plug this approximate expression for the day-night contrast into Equation (21.6) to solve for the characteristic wind speed

$$U \sim \sqrt{\frac{R\Delta T_{\text{eq}}}{\Omega\tau_{\text{rad}}}}. \quad (21.16)$$

For a typical hot Jupiter, $R \approx 3600 \text{ J kg}^{-1} \text{ K}^{-1}$, $\Delta T_{\text{eq}} \approx 1000 \text{ K}$, $\Omega = 2\pi/P_{\text{rot}} \approx 3.6 \times 10^{-5} \text{ s}^{-1}$, $\tau_{\text{rad}} \approx 10^5 \text{ s}$, $L \approx a \approx 1 R_{\text{Jup}}$. Using these values, we find a typical hot Jupiter day-night temperature contrast of $\Delta T \sim 700 \text{ K}$ and a typical wind speed of $U \sim 1000 \text{ m s}^{-1}$.

21.2.2 Comparisons of this simple theory with observations

Figure 21.5 compares the predictions for the day-night temperature contrast (top) and phase curve offset (bottom) derived from the full solution of this simple scaling theory (see Equations 14 and 15 of the Zhang review article) with the state of the art of observations prior to JWST⁶. Note that the day-night temperature contrast is directly predicted from the theory, and is here plotted as a fractional contrast

$$A_T = \frac{(T_{\text{day}} - T_{\text{night}})}{T_{\text{day}}}. \quad (21.17)$$

Meanwhile, the phase curve offset can be estimated from the ratio of radiative to dynamical timescales,

$$\delta \sim \tan^{-1} \left(\frac{\tau_{\text{rad}}}{\tau_{\text{dyn}}} \right), \quad (21.18)$$

⁶At the time of writing there are only two published JWST phase curves, but a statistical sample of JWST phase curve observations is likely forthcoming.

where τ_{rad} was previously derived (see Equation 15.36 in the notes). Note that the ratio $\tau_{\text{rad}}/\tau_{\text{dyn}}$ decreases with increasing equilibrium temperature, causing the predicted phase offset to decrease with T_{eq} .

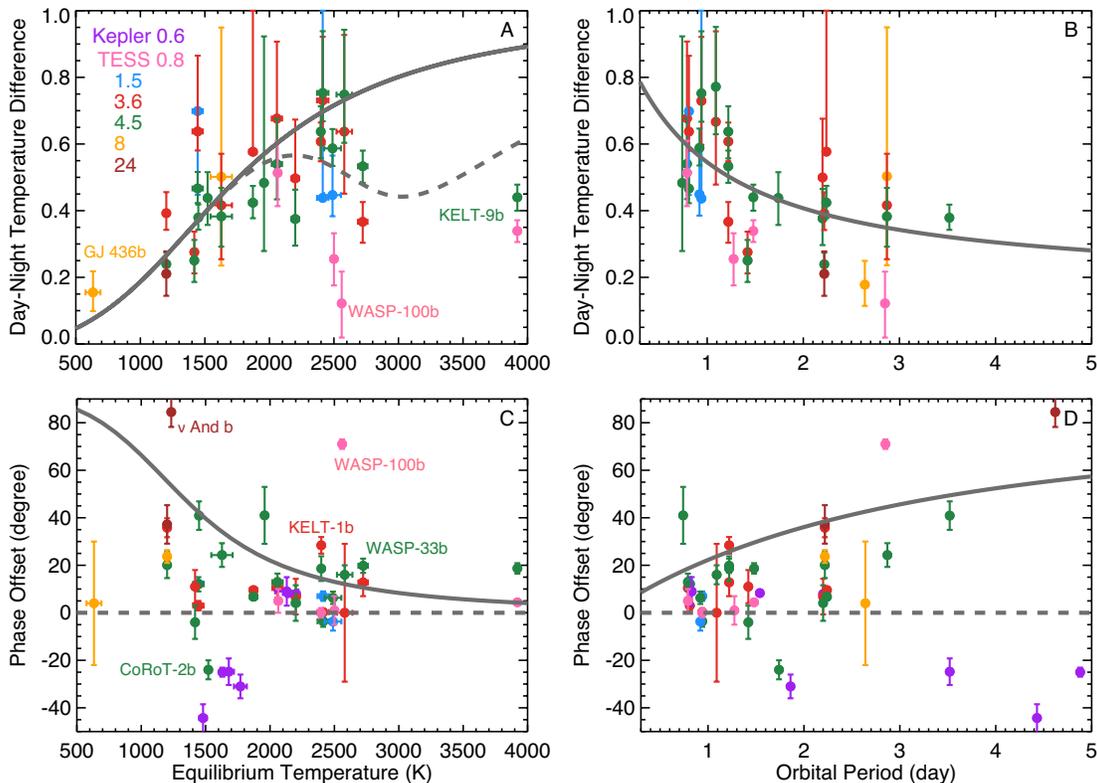


Figure 21.5: Observed day-night brightness temperature difference (top) and phase curve offset (bottom) as a function of equilibrium temperature (left) and orbital period (right), compared with analytic theory (lines). Figure adapted from Zhang (2020).

There is clearly a wide range of scatter in the observed day-night temperature contrast and phase offset (especially in the latter), and there are no statistically robust observational trends from Hubble and Spitzer observations. However, the theory roughly captures the general behavior of the dependence of day-night temperature contrast on equilibrium temperature and rotation period, and incorporating the uncertain strength of atmospheric frictional drag as an effective model uncertainty allows for a better match (Komacek et al., 2017). Meanwhile, the scatter in the observed phase offset is much larger than that predicted by analytic theory. Interestingly, the phase offset of the planetary thermal emission has been measured to be negative with Spitzer phase curves (famously for CoRoT-2b, Dang et al., 2018), which is not expected from standard theoretical models. The analytic theory presented here assumes that the phase offset is set by a kinematic competition between radiation and advection, but in reality it is set by winds Doppler-shifting a planetary-scale wave pattern (Hammond & Pierrehumbert, 2018). Further work is needed to understand the conditions in which this Doppler shifting could lead to westward phase curve offsets.

Note that the dashed line in the day-night temperature difference plot shows modifications

to this scaling theory which include the thermodynamic effects of hydrogen dissociation and recombination (Komacek & Tan, 2018; Tan & Komacek, 2019). Hydrogen begins to thermally dissociate near the photospheres of hot Jupiters with $T \gtrsim 2200$ K, and thus the state of hydrogen transitions from (partially) atomic on the daysides to (largely) molecular form on the nightsides on these “ultra-hot” Jupiters. Energy is required to be input to break the hydrogen bond, and as a result dissociation leads to cooling – conversely, energy is released when hydrogen recombines (analogous to a latent heat), and so recombination leads to heating. This in turn reduces day-night contrasts compared to theoretical expectations that do not include this effect (compare the dashed to the solid lines at $T_{\text{eq}} \gtrsim 2000$ K in Figures 21.5 and 21.6).

One interesting trend from Spitzer phase curve observations is the apparent “flat” nightsides of hot Jupiters (Keating et al., 2019; Beatty et al., 2019). Figure 21.6 shows the day-side and nightside temperatures measured with Spitzer (see also Bell et al., 2021) compared with predictions from the analytic scaling theory above including (solid) and not including (dashed) the thermodynamic effect of hydrogen dissociation and recombination. The night-

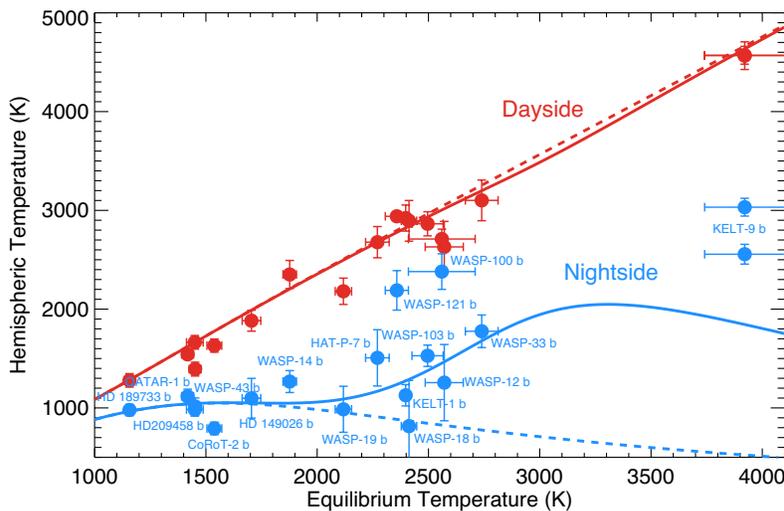


Figure 21.6: Observed dayside (red points) and nightside (blue points) brightness temperatures compared with analytic theory (lines, solid lines include the effects of hydrogen dissociation and recombination while dashed lines do not). Figure adapted from Zhang (2020).

side temperatures of hot Jupiters appear to be roughly constant from 1000 K – 2500 K, even though the dayside temperature roughly scales with the level of irradiation. This may be explained with dynamics, as heat transport from day to night becomes less efficient with increasing equilibrium temperature, increasing the fractional day-night temperature contrast as shown previously in Figure 21.5. However, an alternate explanation of these cold nightsides is a cloud deck that keeps a constant temperature at the cloud top with varying equilibrium temperature – implying that the cloud deck physically moves to higher altitudes with increasing temperature. We’ll explore this possibility in more detail in our hands-on activity in the next class, where we look at the outputs of GCMs that include clouds.

21.2.3 Activity: predicting the day-night contrasts and wind speeds of various exoplanets

We'll do this activity if time allows. Please split up into five groups – each group will study one of the following five planets: KELT-9b, WASP-18b, WASP-43b, GJ 1214b, and K2-18b. Using our derived scaling theory, predict the fractional day-to-night temperature contrast $\Delta T/T_{\text{eq}}$ and characteristic wind speeds U for your planet. To do so, use the NASA Exoplanet Archive (<https://exoplanetarchive.ipac.caltech.edu/>) to find relevant data (e.g., orbital period, semi-major axis, stellar effective temperature, stellar radius) that allows you to calculate quantities needed to estimate the day-night temperature contrast and wind speed. Make sure to explicitly calculate the full-redistribution equilibrium temperature. Report your answers for T_{eq} , $\Delta T/T_{\text{eq}}$, U as plots on the board so we can determine how the atmospheric circulation of tidally locked planets depends on the level of irradiation that they receive.

22 Exoplanet characterization: Atmospheric dynamics

Our agenda for Day 22 is the following:

1. Scale the equations of motion to derive basic force balances and the Rossby number (30 minutes)
2. Activity: Analyze the predictions for atmospheric circulation of hot Jupiters from state-of-the-art GCMs (remaining time)

Today's reading is Section 6 of the Zhang review article on exoplanet and brown dwarf atmospheres. This section will describe our current understanding of the atmospheric circulation of exoplanets derived from a combination of General Circulation Models and theoretical first principles.

22.1 Scale analysis of the momentum equation and basic force balances

22.1.1 The momentum equation in Cartesian coordinates

Recall that the full vector momentum equation is (Equation 21.2)

$$\frac{d\mathbf{v}}{dt} = -\frac{\nabla p}{\rho} + \mathbf{g} - 2\boldsymbol{\Omega} \times \mathbf{v} + \mathcal{F}, \quad (22.1)$$

where the left-hand side represents advection and the terms on the right-hand side (in order) are the pressure gradient, gravity, and Coriolis force, along with additional forces (e.g., drag and resulting dissipation of momentum). In Cartesian coordinates, with x the east-west direction and u the eastward wind, y the north-south direction and v the northward wind, and z the vertical direction and w the vertical (upward) wind, we can write the components of the full vector momentum equation as

$$\frac{du}{dt} = -\frac{1}{\rho} \frac{\partial p}{\partial x} + 2\Omega v \sin\phi - 2\Omega w \cos\phi + \mathcal{F}_x, \quad (22.2)$$

$$\frac{dv}{dt} = -\frac{1}{\rho} \frac{\partial p}{\partial y} - 2\Omega u \sin\phi + \mathcal{F}_y, \quad (22.3)$$

$$\frac{dw}{dt} = -\frac{1}{\rho} \frac{\partial p}{\partial z} - g + 2\Omega u \cos\phi + \mathcal{F}_z, \quad (22.4)$$

where ϕ is latitude, which comes into the equations through the projection of the Coriolis force onto each plane as $\Omega_x = 0$, $\Omega_y = \Omega \cos\phi$, and $\Omega_z = \Omega \sin\phi$.

22.1.2 Vertical force balance: hydrostatic equilibrium

Let's now study the characteristic values of each term in the vertical momentum equation. Re-writing it using the material derivative and dropping \mathcal{F}_z , we can express the vertical momentum equation as

$$\frac{\partial w}{\partial t} + \mathbf{v} \cdot \nabla w = -\frac{1}{\rho} \frac{\partial p}{\partial z} - g + 2\Omega u \cos\phi. \quad (22.5)$$

We can write down a scaled version of this equation as

$$\frac{W}{\tau} + \frac{UW}{L} \sim -\frac{p}{\rho H} - g + \Omega U. \quad (22.6)$$

For characteristic hot Jupiter atmospheric conditions, we can estimate $T \sim 1000$ K, $H \sim RT/g \sim 3600 \text{ J kg}^{-1} \text{ K}^{-1} 1000 \text{ K} / 25 \text{ m s}^{-2} \sim 144$ km, $U \sim 1$ km s⁻¹, $L \sim a \approx 7 \times 10^4$ km, $g = GM/a^2 \sim 25$ m s⁻², $\Omega \sim 3.6 \times 10^{-5}$ s⁻¹, $\tau \sim 10^5$ s, $\Delta p = \rho R \Delta T$ with $\Delta T \sim 1000$ K, and $W \sim UH/L \sim 2$ m s⁻¹. Plugging these in, we find that the approximate scaling for each term (each in m s⁻²) is

$$\begin{aligned} \frac{W}{\tau} + \frac{UW}{L} &\sim -\frac{p}{\rho H} - g + \Omega U, \\ 2 \times 10^{-5} + 3 \times 10^{-5} &\sim -25 - 25 + 0.036. \end{aligned} \quad (22.7)$$

Thus, we can see that the dominant two terms are the pressure gradient and gravity – as expected, hydrostatic balance holds on a large scale. Note that this also holds on a local scale – if we conduct the same exercise as above, but study the local changes in the pressure gradient and gravity due to the circulation, we still find that the pressure gradient and gravity terms are orders of magnitude larger than other terms. As a result, to good approximation we can consider the atmosphere to be in a state of vertical hydrostatic equilibrium.

22.1.3 Horizontal force balance: Rossby number, geostrophy

We can similarly write a full version of the y-component of the horizontal momentum equation neglecting \mathcal{F}_y ,

$$\frac{\partial v}{\partial t} + \mathbf{v} \cdot \nabla v = -\frac{1}{\rho} \frac{\partial p}{\partial y} - 2\Omega u \sin \phi. \quad (22.8)$$

As above, scaling this expression and plugging in the same characteristic values, we now find for the horizontal momentum balance

$$\begin{aligned} \frac{U}{\tau} + \frac{U^2}{L} &\sim -\frac{R\Delta T}{L} - \Omega U, \\ 10^{-2} + 1.4 \times 10^{-2} &\sim -5.1 \times 10^{-2} - 3.2 \times 10^{-2}. \end{aligned} \quad (22.9)$$

This is clearly far trickier! Each of these terms are comparable to one another for hot Jupiters (for Earth, the pressure gradient and Coriolis forces are in balance, with advection playing a minimal role except near the equator). Assuming that the atmosphere is in steady-state ($\partial v/\partial t \approx 0$), we have a three-way force balance between advection, pressure gradients, and Coriolis forces that set the dynamics:

$$\mathbf{v} \cdot \nabla v \approx -\frac{1}{\rho} \frac{\partial p}{\partial y} - 2\Omega u \sin \phi. \quad (22.10)$$

Irradiated atmospheres will always have pressure and temperature gradients due to the radiative forcing contrast between the more irradiated and less irradiated regions. Thus, what determines the dynamical regime of an atmosphere (i.e., which term balances the

pressure gradient) is the ratio of the advective term to the Coriolis term in the momentum equation. This is the Rossby number,

$$\text{Ro} \equiv \frac{\text{advection}}{\text{Coriolis}} = \frac{U}{fL}, \quad (22.11)$$

where $f = 2\Omega\sin\phi$ is the Coriolis parameter. Depending on the Rossby number, an atmosphere can be in one of two regimes

$$\begin{aligned} \text{Ro} < 1 & \text{ geostrophic balance,} \\ \text{Ro} > 1 & \text{ advection dominated.} \end{aligned} \quad (22.12)$$

On Earth, geostrophic balance applies in the mid-latitudes, so these regimes are often called “extra-tropical” ($\text{Ro} < 1$) and “tropical” ($\text{Ro} > 1$) dynamics, respectively. Note that in

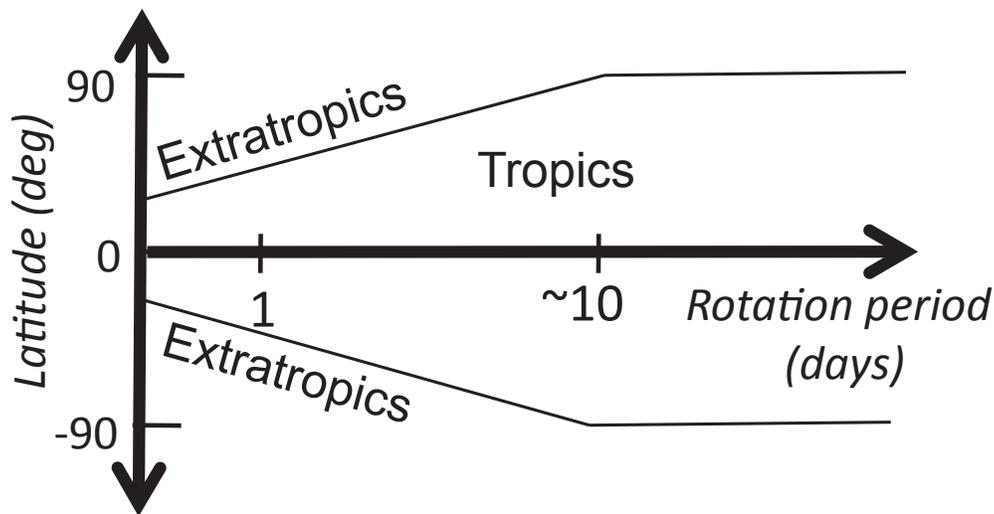


Figure 22.1: Schematic plot showing the latitudinal extent of tropical and extratropical regions on an Earth-sized planet as a function of rotation period. While Earth has both tropical and extratropical dynamical regimes near the equator and mid-latitudes, Earth-sized planets with similar wind speeds and rotation periods longer than 10 days would be “all-tropics” worlds. Figure adapted from (Showman & Kaspi, 2013).

our estimate for a typical hot Jupiter, we expect $\text{Ro} \approx 0.44$. For Earth, $\text{Ro} \approx 0.1$ in the mid-latitudes. Thus, geostrophic balance (“geostrophy”) often holds on the global scale for planets (note that it does not for very slowly rotating planets, like Venus). The horizontal momentum balance in geostrophy can be written as

$$\begin{aligned} fu &\approx -\frac{1}{\rho} \frac{\partial p}{\partial y}, \\ fv &\approx \frac{1}{\rho} \frac{\partial p}{\partial x}. \end{aligned} \quad (22.13)$$

Because geostrophic winds are at right angles to pressure gradients, geostrophy implies that the circulation will follow isobars (contours of constant pressure) – that is, horizontal winds

will flow parallel to isobars. Figure 22.2 shows a weather map near the 500 mbar pressure level over North America. The wind barbs generally point parallel to the isobars, indicating

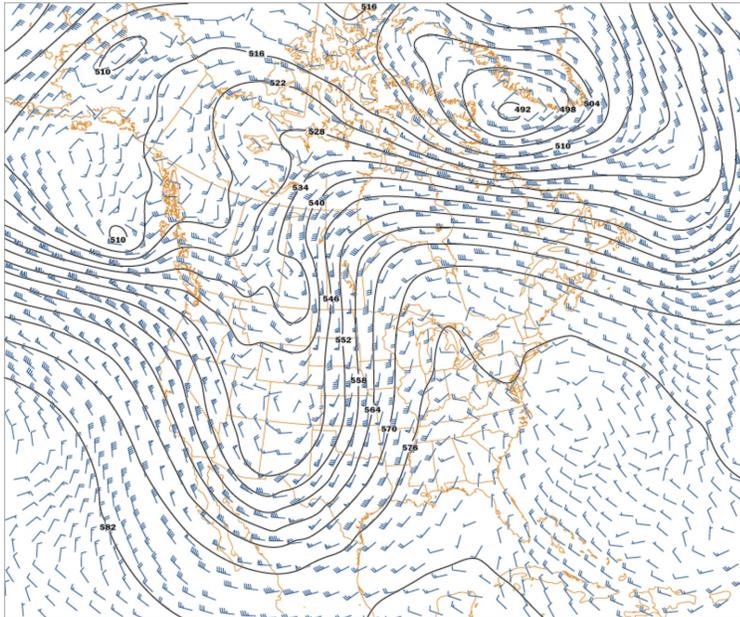


Figure 22.2: Isobars near 500 mbars over North America, along with wind barbs. Note that the wind direction is generally parallel to the isobars, a result of geostrophic balance. Figure adapted from <https://www.weather.gov/jetstream/500mb>.

that the atmosphere in Earth’s mid-latitudes is near geostrophic balance.

22.2 Python activity: Grid of hot Jupiter GCMs

Now that we’ve done some basic scaling, let’s get a sense for what is physically happening in the atmospheres of hot Jupiters by looking at the results of numerical simulations. We’ll specifically look at General Circulation Model simulations, which solve an equation set that includes momentum conservation (force balance), mass conservation, local hydrostatic balance, thermal energy balance (first law of thermodynamics), and an equation of state (ideal gas law). To recap, these are, in order,

$$\frac{d\mathbf{v}}{dt} = -\frac{\nabla p}{\rho} + \mathbf{g} - 2\boldsymbol{\Omega} \times \mathbf{v} + \mathcal{F}, \quad (22.14)$$

$$\frac{d\rho}{dt} + \rho \nabla \cdot \mathbf{v} = 0, \quad (22.15)$$

$$\frac{\partial p}{\partial z} = -\rho g, \quad (22.16)$$

$$c_p \frac{dT}{dt} = q + \frac{1}{\rho} \frac{dp}{dt}, \quad (22.17)$$

$$p = \rho RT. \quad (22.18)$$

Specifically, GCMs solve the “primitive equations” of meteorology – a reduced form of the full Navier-Stokes equations applied to atmospheres on a rotating sphere (which are somewhat different than those above, but the differences are irrelevant for understanding the general scaling of atmospheric dynamics with planetary properties).

The .zip file on ELMS contains GCM output and a Python-based plotting tool to analyze these (kindly provided by Alex Roth and Vivien Parmentier). These GCMs are similar to those in Parmentier et al. (2021), and conducted with the SPARC/MITgcm (Showman et al., 2009). However, this is an updated and novel grid that covers a much larger parameter space over varying planetary equilibrium temperature, host star type, planetary gravity, and planetary atmospheric metallicity, recently posted on arXiv as Roth et al. (2024)⁷.

First, we'll start with the cloudless GCM models. Run the python script "uber-grid_sliders_update.py" and you should see a pop-up window with sliders that allow you to vary planetary parameters and see how it affects spectra, phase curves, and the temperature map of the planet. Use this plotting tool to answer the following questions.

1. Let's start by studying how varying planetary properties affects near-infrared phase curves measured with HST. To select HST for the phase curve, move the "bandpass" slider under the phase curve to "5." While keeping metallicity, stellar mass, and $\log(g)$ at their default values, vary the equilibrium temperature from 1000 to 2400 K. How does the varying equilibrium temperature affect the phase curve amplitude and offset? To better understand these dependencies, look at how the temperature map changes with temperature.

2. We can use the radiative timescale to help interpret what's going on. As a reminder, this can be written as

$$\tau_{\text{rad}} = \frac{p}{g} \frac{c_p}{4\sigma T^3}, \quad (22.19)$$

where p/g is the atmospheric mass, c_p is the specific heat capacity, and T is the atmospheric temperature. We can re-write this expression using opacity rather than p/g to think about the effect of atmospheric metallicity on the dynamics and phase curve. Recall that the optical depth $d\tau = \kappa\rho dl$. Given that the mass per area is $\rho z = p/g$, we can express $\tau \sim \kappa p/g$. If we assume that the emission from the planet comes from the $\tau = 1$ level, then we can write $p/g = 1/\kappa$. Thus, we can re-write the radiative timescale as

$$\tau_{\text{rad}} = \frac{c_p}{4\sigma\kappa T^3}, \quad (22.20)$$

where κ is the opacity, which increases with increasing metallicity. Now keep $T_{\text{eq}} = 1600$ K but vary the metallicity from 0 to 1.5. How do the phase curve offset and phase curve amplitude vary with metallicity?

3. Now keep $T_{\text{eq}} = 1600$ K and $\log(M/H) = 0.0$ but vary the gravity ($\log g$) from 0.8 - 1.8. How do the phase curve offset and phase curve amplitude vary with gravity?
4. Using the expression for radiative timescale, come up with an explanation for what causes the trends in phase curve offset and amplitude with metallicity, gravity and equilibrium temperature.
5. For the case with $T_{\text{eq}} = 2400$ K, look at the phase curve and temperature map with and without TiO/VO (e.g., change the TiO/VO slider from 0 to 1). Specifically, compare

⁷Note that the activity was made with these GCMs before they were published, but the same basic results hold in the final version of the GCM grid.

how the phase curve amplitude changes between the bolometric phase curve (bandpass 12) and the HST/WFC3 phase curve (bandpass 5). Can you explain the difference between the bandpasses?

Now let's consider the effect of clouds in our GCM models. To do so, run the python script "cloud_sliders_update.py" to open the plotting tool. These plots show results of GCMs that contain MnS clouds of a range of prescribed particle sizes (0.1, 1, 10 μm) over varying equilibrium temperature. Use this script to answer the following questions.

6. Keep $T_{\text{eq}} = 1400$ K and reduce the particle size to 0.1 μm . Then, look at the phase curve in bands 2, 4, and 6, which correspond to HST/STIS2, HST/WFC3, and Spitzer Channel 2. For which of these bands is the phase curve dominated by reflected light? Which of the bands are dominated by thermal emission?
7. From comparing the phase curve in band 2 (HST/STIS2) with the cloud and temperature maps at various pressures, estimate the range of pressure levels (in bars) that are probed by the phase curve.
8. Now continuing to study the band 2 phase curve, vary the equilibrium temperature from 1000 - 2000 K while keeping the particle size fixed at 0.1 μm . What is the range of temperatures where the phase curve peaks after secondary eclipse? Use the cloud map to understand what sets this temperature range.
9. Now choose band 6 (Spitzer Ch. 2, which is a wavelength of 4.5 μm), fix $T_{\text{eq}} = 1400$ K, and vary the cloud particle size from 0.1 - 10 μm . Compare the cloudy and no clouds phase curves to determine the effect of clouds on the phase curve. How does the effect of clouds vary with particle size? Why does changing the particle size change the effect of the clouds on the phase curve? Hint: this is related to Mie scattering.

If you've finished early, you can also take a look at the grid of models with a prescribed nightside cloud deck ("NSonly_cloud_sliders_update.py"). These are more idealized, and clouds are placed only on the nightside.

10. For band 5 (Spitzer Ch. 1) in the baseline case, determine how nightside clouds affect the temperature map. Link this to the influence of clouds on the phase curve.
11. Use the temperature map to attempt to guess what pressures are probed in the phase curve. The temperature map at which pressure level provides best fits to the cloudy phase curve? What about the pressure level that best fits the clear phase curves? Are the cloudy and cloud-free photospheres at similar or different pressures, and why?
12. How do nightside clouds affect the dayside spectra (phase 0)? What about the nightside spectra (phase 180)?

References

- [1] Agol, E., & Fabrycky, D. C. 2018, in Handbook of Exoplanets, ed. H. J. Deeg & J. A. Belmonte, 7
- [2] Agol, E., Steffen, J., Sari, R., & Clarkson, W. 2005, MNRAS, 359, 567
- [3] Agol, E., Dorn, C., Grimm, S. L., et al. 2021, , 2, 1
- [4] Apai, D., Nardiello, D., & Bedin, L. R. 2021, ApJ, 906, 64
- [5] Armitage, P. J. 2007, arXiv e-prints, astro
- [6] —. 2013, Astrophysics of Planet Formation
- [7] Bailes, M., Lyne, A. G., & Shemar, S. L. 1991, Nature, 352, 311
- [8] Baraffe, I., Chabrier, G., Allard, F., & Hauschildt, P. H. 2002, A&A, 382, 563
- [9] Batygin, K., & Stevenson, D. 2010, The Astrophysical Journal Letters, 714, L238
- [10] Beatty, T. G., Marley, M. S., Gaudi, B. S., et al. 2019, The Astronomical Journal, 158, 166
- [11] Bell, T. J., Dang, L., Cowan, N. B., et al. 2021, MNRAS, 504, 3316
- [12] Benedict, G. F., McArthur, B. E., Forveille, T., et al. 2002, ApJL, 581, L115
- [13] Berger, T. A., Schlieder, J. E., Huber, D., & Barclay, T. 2023, arXiv e-prints, arXiv:2302.00009
- [14] Bodenheimer, P., Lin, D., & Mardling, R. 2001, The Astrophysical Journal, 548, 466
- [15] Bond, I. A., Udalski, A., Jaroszyński, M., et al. 2004, ApJL, 606, L155
- [16] Boss, A. P. 2011, ApJ, 731, 74
- [17] Brygoo, S., Loubeyre, P., Millot, M., et al. 2021, Nature, 593, 517
- [18] Carroll, B. W., & Ostlie, D. A. 2017, An introduction to modern astrophysics, Second Edition
- [19] Charbonneau, D., Brown, T., Noyes, R., & Gilliland, R. 2002, The Astrophysical Journal, 568, 377
- [20] Charbonneau, D., Allen, L. E., Megeath, S. T., et al. 2005, ApJ, 626, 523
- [21] Coulombe, L.-P., Benneke, B., Challener, R., et al. 2023, Nature, 620, 292
- [22] Dang, L., Cowan, N., Schwartz, J., & et al. 2018, Nature Astronomy, 2, 220
- [23] Deming, D., Seager, S., Richardson, L. J., & Harrington, J. 2005, Nature, 434, 740

- [24] Dressing, C. D., & Charbonneau, D. 2013, *ApJ*, 767, 95
- [25] Fischer, D. A., & Valenti, J. 2005, *ApJ*, 622, 1102
- [26] Fischer, D. A., Marcy, G. W., Butler, R. P., et al. 2008, *ApJ*, 675, 790
- [27] Fortney, J., Dawson, R., & Komacek, T. 2021, *Journal of Geophysical Research: Planets*, 126, e06629
- [28] Fortney, J., Shabram, M., Showman, A., et al. 2010, *The Astrophysical Journal*, 709, 1396
- [29] Fortney, J. J. 2005, *MNRAS*, 364, 649
- [30] Fressin, F., Torres, G., Charbonneau, D., et al. 2013, *ArXiv: 1301.0842*, arXiv:1301.0842
- [31] Fulton, B., Petigura, E., Howard, A., et al. 2017, *The Astronomical Journal*, 154, 109
- [32] Fulton, B. J., Petigura, E. A., Howard, A. W., et al. 2017, *AJ*, 154, 109
- [33] Fulton, B. J., Rosenthal, L. J., Hirsch, L. A., et al. 2021, *ApJS*, 255, 14
- [34] Gao, P., Wakeford, H., Moran, S., & Parmentier, V. 2021, *Journal of Geophysical Research: Planets*, 126, e06655
- [35] Gaudi, B. S., Bennett, D. P., Udalski, A., et al. 2008, *Science*, 319, 927
- [36] Gillon, M., Triaud, A., Demory, B., et al. 2017, *Nature*, 542, 456
- [37] Gould, A., Jung, Y. K., Hwang, K.-H., et al. 2022, *Journal of Korean Astronomical Society*, 55, 173
- [38] Guillot, T., Chabrier, G., Gautier, D., & Morel, P. 1995, *ApJ*, 450, 463
- [39] Guillot, T., & Showman, A. 2002, *Astronomy and Astrophysics*, 385, 156
- [40] Hammond, M., & Pierrehumbert, R. 2018, *The Astrophysical Journal*, 869, 65
- [41] Holman, M. J., & Murray, N. W. 2005, *Science*, 307, 1288
- [42] Howard, A. W., Marcy, G. W., Bryson, S. T., et al. 2012, *ApJS*, 201, 15
- [43] Kasting, J., Whitmire, D., & Reynolds, R. 1993, *Icarus*, 101, 108
- [44] Keating, D., Cowan, N., & Dang, L. 2019, *Nature Astronomy*, 3, 1092
- [45] Komacek, T., & Showman, A. 2016, *The Astrophysical Journal*, 821, 16
- [46] Komacek, T., Showman, A., & Tan, X. 2017, *The Astrophysical Journal*, 835, 198
- [47] Komacek, T., & Tan, X. 2018, *Research notes of the AAS*, 2, 36
- [48] Komacek, T., & Youdin, A. 2017, *The Astrophysical Journal*, 844, 94

- [49] Komacek, T. D., & Abbot, D. S. 2016, *ApJ*, 832, 54
- [50] Kreidberg, L. 2017, *Exoplanet Atmosphere Measurements from Transmission Spectroscopy and Other Planet Star Combined Light Observations*, 100
- [51] Kreidberg, L., Bean, J. L., Désert, J.-M., et al. 2014, *Nature*, 505, 69
- [52] Krissansen-Totton, J., Garland, R., Irwin, P., & Catling, D. 2018, *The Astronomical Journal*, 156, 114
- [53] Krissansen-Totton, J., Olson, S., & Catling, D. C. 2018, *Science Advances*, 4, eaao5747
- [54] Laughlin, G., Crismani, M., & Adams, F. 2011, *The Astrophysical Journal Letters*, 729, L7
- [55] Lovis, C., & Mayor, M. 2007, *A&A*, 472, 657
- [56] Luger, R., & Barnes, R. 2015, *Astrobiology*, 15, 119
- [57] Lustig-Yaeger, J., Meadows, V., & Lincowski, A. 2019, *The Astronomical Journal*, 158, 27
- [58] Lyne, A. G., & Bailes, M. 1992, *Nature*, 355, 213
- [59] Macintosh, B., Graham, J. R., Barman, T., et al. 2015, *Science*, 350, 64
- [60] Marois, C., Lafrenière, D., Doyon, R., Macintosh, B., & Nadeau, D. 2006, *ApJ*, 641, 556
- [61] Mayor, M., & Queloz, D. 1995, *Nature*, 378, 355
- [62] McArthur, B. E., Benedict, G. F., Barnes, R., et al. 2010, *ApJ*, 715, 1203
- [63] Meadows, V., Reinhard, C., Arney, G., & et al. 2018, *Astrobiology*, 18, 630
- [64] Mróz, P., Ryu, Y. H., Skowron, J., et al. 2018, *AJ*, 155, 121
- [65] Mróz, P., Poleski, R., Gould, A., et al. 2020, *ApJL*, 903, L11
- [66] Naef, D., Latham, D. W., Mayor, M., et al. 2001, *A&A*, 375, L27
- [67] Nielsen, E. L., De Rosa, R. J., Macintosh, B., et al. 2019, *AJ*, 158, 13
- [68] Öberg, K. I., Murray-Clay, R., & Bergin, E. A. 2011, *ApJL*, 743, L16
- [69] Owen, J. E., & Wu, Y. 2013, *ApJ*, 775, 105
- [70] Paczynski, B. 1996, *ARA&A*, 34, 419
- [71] Parmentier, V., & Crossfield, I. 2018, *The Exoplanet Handbook (Springer)*, 116

- [72] Parmentier, V., Showman, A., & Fortney, J. 2021, *Monthly Notices of the Royal Astronomical Society*, 501, 78
- [73] Perez-Becker, D., & Showman, A. 2013, *The Astrophysical Journal*, 776, 134
- [74] Perryman, M., Hartman, J., Bakos, G. Á., & Lindegren, L. 2014, *ApJ*, 797, 14
- [75] Pierrehumbert, R. T. 2010, *Principles of Planetary Climate*
- [76] Pollack, J. B., Hubickyj, O., Bodenheimer, P., et al. 1996, *Icarus*, 124, 62
- [77] Roth, A., Parmentier, V., & Hammond, M. 2024, arXiv e-prints:2404.09626
- [78] Rustamkulov, Z., Sing, D. K., Mukherjee, S., et al. 2023, *Nature*, 614, 659
- [79] Sackett, P. D. 1999, in *NATO Advanced Study Institute (ASI) Series C, Vol. 532, Planets Outside the Solar System: Theory and Observations*, ed. J. M. Mariotti & D. Alloin, 189
- [80] Sarkis, P., Mordasini, C., Henning, T., Marleau, G., & Mollière, P. 2021, *Astronomy & Astrophysics*, 645, A79
- [81] Seager, S. 2010, *Exoplanet Atmospheres: Physical Processes*
- [82] Seager, S., & Mallén-Ornelas, G. 2003, *ApJ*, 585, 1038
- [83] Showman, A., & Kaspi, Y. 2013, *The Astrophysical Journal*, 776, 85
- [84] Showman, A., Tan, X., & Parmentier, V. 2020, *Space Science Reviews*, 216, 139
- [85] Showman, A. P., Fortney, J. J., Lian, Y., et al. 2009, *ApJ*, 699, 564
- [86] Sotin, C., Jackson, J., & Seager, S. 2010, *Exoplanets*, ed. S. Seager (University of Arizona Press), 375–395
- [87] Soubiran, F., & Militzer, B. 2015, arXiv e-prints
- [88] Tan, X., & Komacek, T. 2019, *The Astrophysical Journal*, 886, 26
- [89] Tan, X., & Showman, A. 2021, *Monthly Notices of the Royal Astronomical Society*, 502, 678
- [90] Thorngren, D., & Fortney, J. 2018, *The Astronomical Journal*, 155, 214
- [91] Tsai, S., Malik, M., Kitzmann, D., et al. 2021, *The Astrophysical Journal*, 923, 264
- [92] Turcotte, D., & Schubert, G. 2002, *Geodynamics* (New York, NY: Cambridge University Press)
- [93] van de Kamp, P. 1969, *AJ*, 74, 757
- [94] Visscher, C., & Moses, J. 2011, *The Astrophysical Journal*, 738, 72

- [95] Way, M., Del Genio, A., Kiang, N., et al. 2016, *Geophysical Research Letters*, 43, 8376
- [96] Weiss, L. M., Deck, K. M., Sinukoff, E., et al. 2017, *AJ*, 153, 265
- [97] Weiss, L. M., Marcy, G. W., Petigura, E. A., et al. 2018, *AJ*, 155, 48
- [98] Wolf, E., & Toon, O. 2015, *Journal of Geophysical Research: Atmospheres*, 120, 5775
- [99] Wolszczan, A. 1994, *Science*, 264, 538
- [100] Wolszczan, A., & Frail, D. A. 1992, *Nature*, 355, 145
- [101] Wordsworth, R., & Kreidberg, L. 2022, *ARA&A*, 60, 159
- [102] Wright, J. T., & Gaudi, B. S. 2013, in *Planets, Stars and Stellar Systems. Volume 3: Solar and Stellar Planetary Systems*, ed. T. D. Oswalt, L. M. French, & P. Kalas, 489
- [103] Yang, J., Boue, G., Fabrycky, D., & Abbot, D. 2014, *The Astrophysical Journal Letters*, 787, L2
- [104] Zahnle, K., & Catling, D. 2017, *The Astrophysical Journal*, 843, 122
- [105] Zhang, X. 2020, *Research in Astronomy and Astrophysics*, 20, 099
- [106] Zhang, X., & Showman, A. 2017, *The Astrophysical Journal*, 836, 73