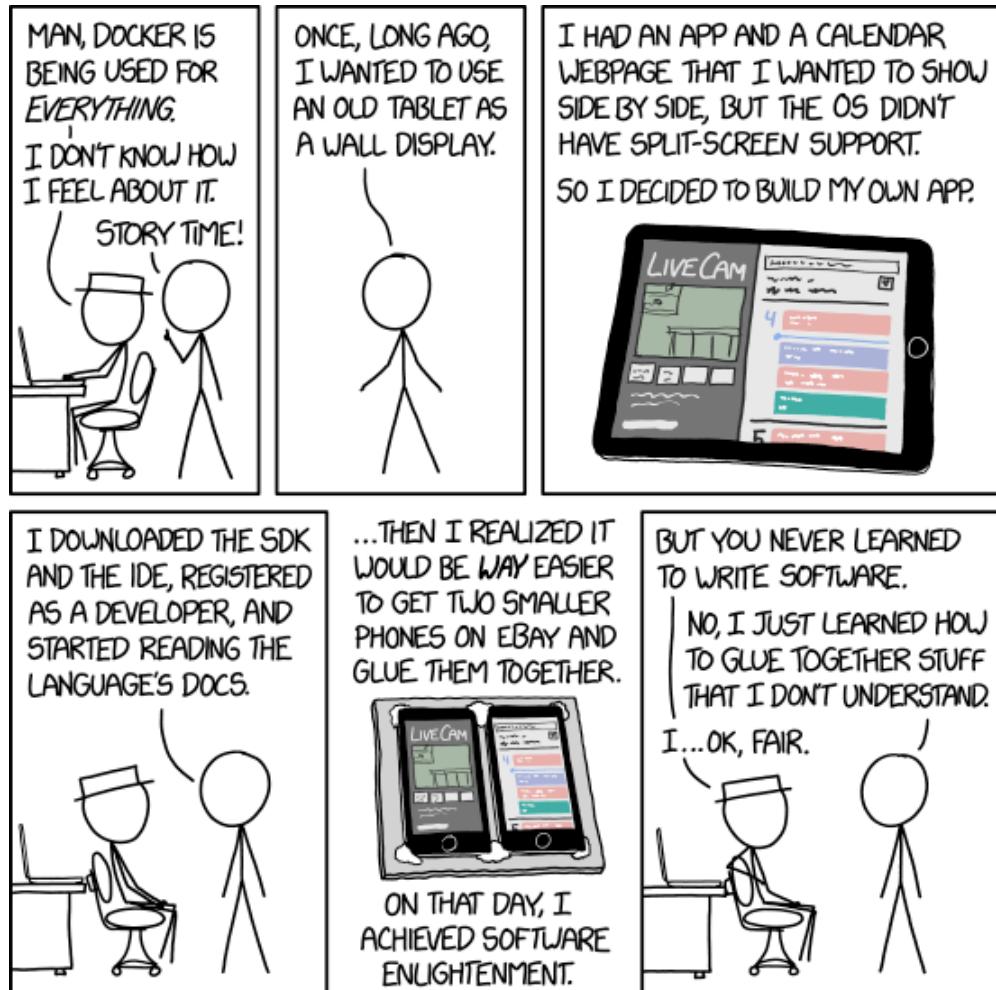


Software Carpentry



xkcd – there may be a lot of this...

Matt Hilton
Astrophysics & Cosmology Research Unit, UKZN

Software carpentry?

- ... as opposed to software engineering???
- An organisation that has been teaching computer skills to scientists since 1998: <https://software-carpentry.org/>
(Disclaimer: I'm not from them... nor am I a software engineer...)
- Topics covered on the UK DISCnet course:
 - Automating tasks with the Unix shell
 - Python basics: building programs with Python
 - Writing robust code and unit testing with Python
 - Version control with git
- We have only today... so:
 - This talk – some edited highlights / discussion / tips
 - Later – a tutorial on Git and setting up a Python application to find extrasolar planets in Kepler space telescope data

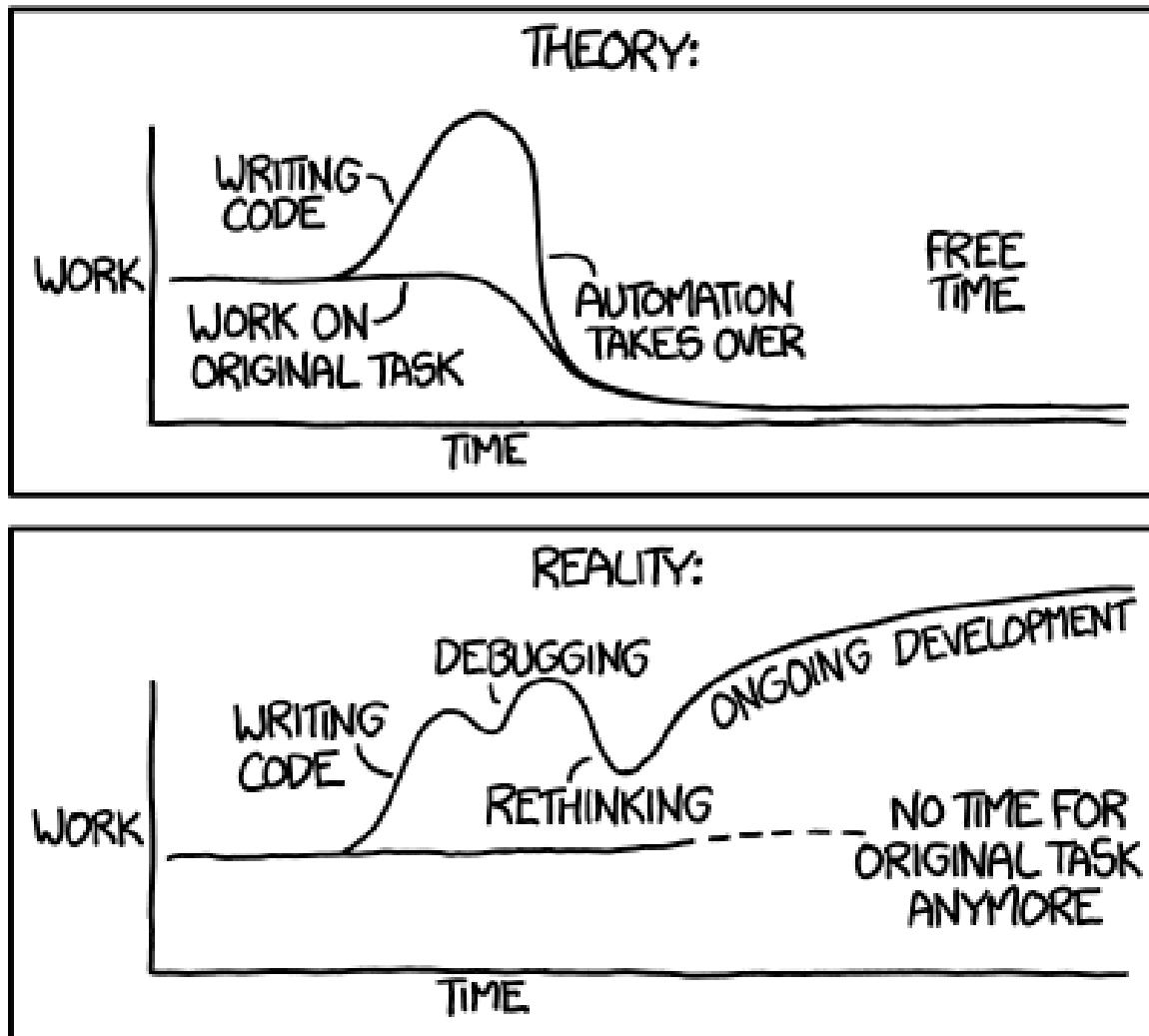


- What is it?
 - An interpreted, cross-platform high-level programming language
- Why use it?
 - Open Source (<https://opensource.org/>)
 - Not as fast as compiled languages like C etc., but much faster to write/debug (higher productivity)
 - Nice, friendly syntax (if you like indenting)
 - The standard library is amazing (with great documentation)...
 - ... numpy / scipy / matplotlib etc. are amazing too
- What are the alternatives?
 - R? (<https://www.r-project.org/>)
 - Perl? (<https://www.perl.org/>)
 - Octave? (<https://www.gnu.org/software/octave/>)
 - Julia? (looks interesting; <https://julialang.org/>)

Automation

- More xkcd

"I SPEND A LOT OF TIME ON THIS TASK.
I SHOULD WRITE A PROGRAM AUTOMATING IT!"



Design

- What do we mean by “well designed software”?
- A good design ought to break your code up into the smallest units that are practical (i.e., functions should be short, modular, reusable), and should be organised in a sensible way (e.g., put related functions in the same module; use classes where it makes sense)
- If you are just doing some data analysis task on your own, you perhaps don’t need to worry about design much...
- On the other hand, you DO need to worry about design if you are writing code that is going to be used by others. You will need to think more deeply about the data structures, the tools needed, and especially the interfaces between them.
- You don’t want to have to keep revising interfaces (this is annoying and breaks other people’s code) – if you design it right the first time, you won’t have to.
- If your code has a good design, it should be easy to maintain, and for others to understand

Coding style

- Consistency, and readability, is key – particularly within a package or module
- You can find the style guide for Python and its Standard Library here:

<https://www.python.org/dev/peps/pep-0008/>

(certainly read this if you haven't done much coding before)

- Never use tab to indent in Python – use spaces instead
- Use comments and docstrings (but keep them up-to-date)
- A couple of naming conventions:
 - Module names, function names, variables:
lowercase_with_underscores
 - Class names: CapitalizedWords or CamelCase
- You'll see in the tutorial code that I use mixedCase a lot out of habit (I don't like all the underscores)... this is fairly common elsewhere (see, e.g., the Qt Framework)

Debugging

- The nice thing about Python is that it always tells you exactly what went wrong, and where...
... if this is a problem in someone else's code: remember to give them the entire error message ('it crashed', 'there is a bug' contains no useful information)

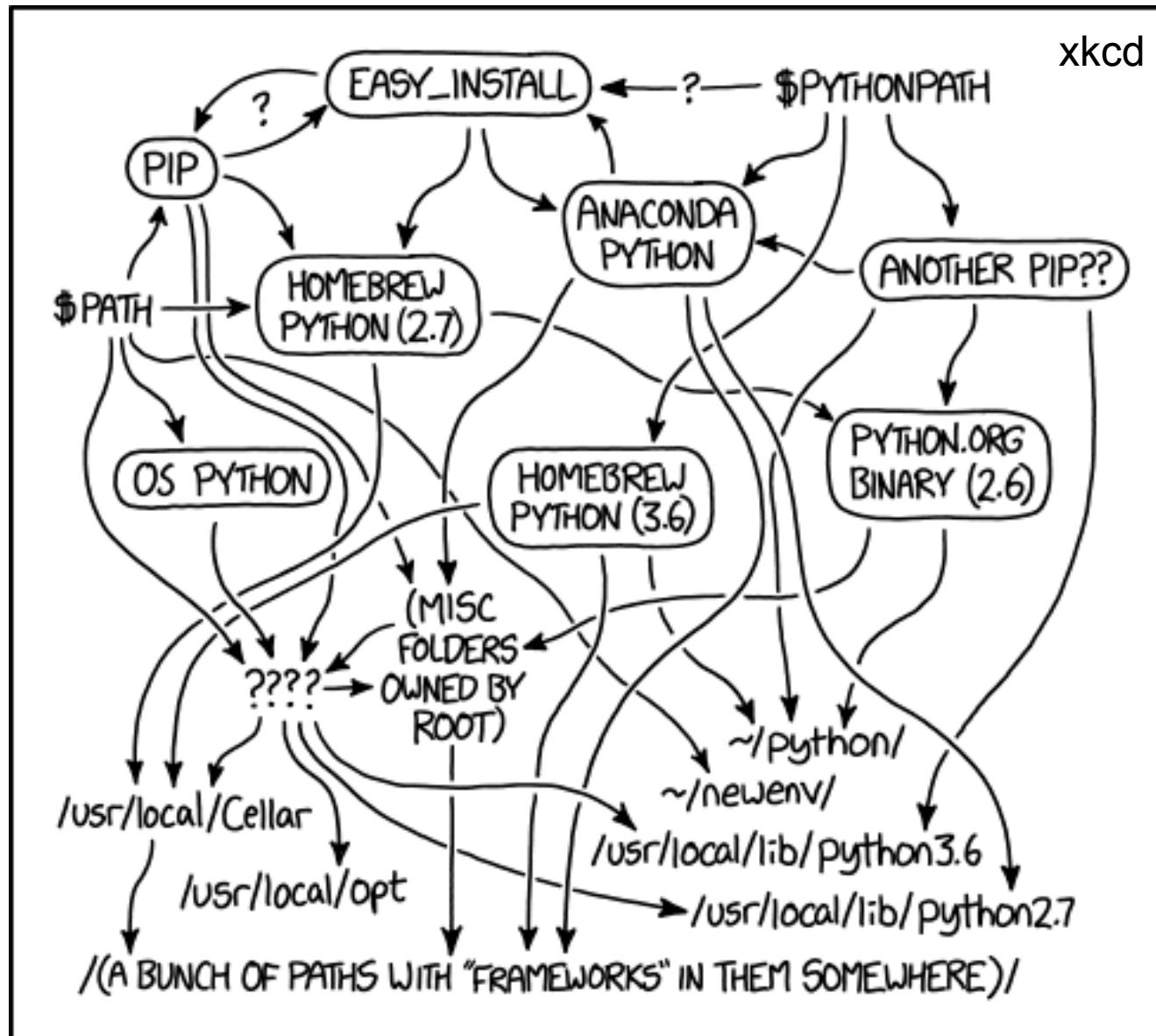
```
Ignoring a divide by zero error ...
Traceback (most recent call last):
  File "debuggingWithIPython.py", line 34, in <module>
    c=divide(a, b)
  File "debuggingWithIPython.py", line 17, in divide
    return a/b
TypeError: unsupported operand type(s) for /: 'float' and 'str'
```

Debugging

```
1 """
2
3 Simple example of using IPython for debugging
4
5 """
6
7 import sys
8 import IPython
9
10 #-----
11 def divide(a, b):
12     """Divide two numbers
13
14     Returns a float
15
16     """
17     return a/b
18
19 #-----
20 # Main
21
22 a=5.0
23 b=2.0
24
25 c=divide(a, b)
26
27 b=0.0
28 try:
29     c=divide(a, b)
30 except ZeroDivisionError:
31     print("Ignoring a divide by zero error ...")
32
33 b="a"
34 try:
35     c=divide(a, b)
36 except:
37     print("This would have triggered a TypeError exception")
38     IPython.embed()
39     sys.exit()
40
```

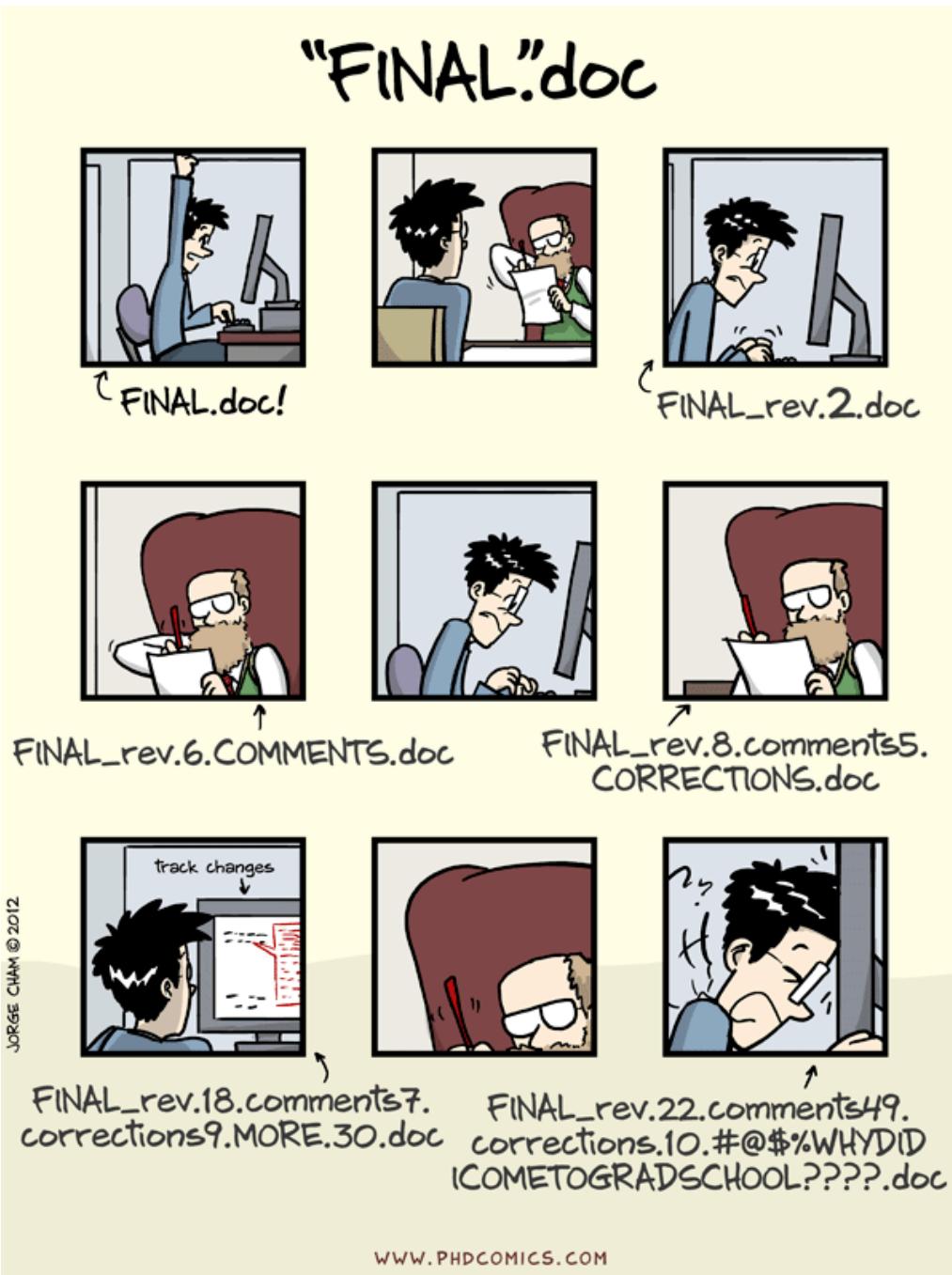
- Because Python is an interpreted language, you can use IPython to “drop in” to a program at any time, inspect the values and types of variables, and figure out why it fell over

Python proliferation



MY PYTHON ENVIRONMENT HAS BECOME SO DEGRADED
THAT MY LAPTOP HAS BEEN DECLARED A SUPERFUND SITE.

Version control



- There are many choices...
... but just use git
- Don't use git to manage a repository with lots of large data files though (keep data and code separate wherever possible)
- Version control helps to keep things running smoothly if you're working in a big collaboration
- Tagging versions / releases of code used in particular papers is a good idea...

Testing, testing, testing...

- You should test your codes frequently
- Some widely-used frameworks exist for Python:
 - The unittest module in the Standard Library:
<https://docs.python.org/3/library/unittest.html>
 - The Nose framework:
<https://nose.readthedocs.io/en/latest/>
- The general idea is to embed test cases for each individual routine / class etc. into your source code, so that you can automate the testing process on “units” that are as small as practical
- Confession:

I've not used either... I just have a bunch of standard things that I run to check that I'm getting the output I expect for a given package (and I can also use version control to check what changed, if things are different)

Refactoring

- You need to maintain your code... from time to time, you might want to re-write / restructure some of it to make it easier to maintain (refactoring)
- In general: don't break interfaces when you do this – this annoys people who use your code (which may include yourself!)
- Sometimes you may think it would be nice to completely re-write something from scratch...
... but ask yourself: is it worth it?

... also remember: the computer does not care if your code is beautiful or not, if it works

“Perfect is the enemy of good enough”

Optimisation

- Yet more xkcd

HOW LONG CAN YOU WORK ON MAKING A ROUTINE TASK MORE EFFICIENT BEFORE YOU'RE SPENDING MORE TIME THAN YOU SAVE?
(ACROSS FIVE YEARS)

		HOW OFTEN YOU DO THE TASK					
		50/DAY	5/DAY	DAILY	WEEKLY	MONTHLY	YEARLY
HOW MUCH TIME YOU SHAVE OFF	1 SECOND	1 DAY	2 HOURS	30 MINUTES	4 MINUTES	1 MINUTE	5 SECONDS
	5 SECONDS	5 DAYS	12 HOURS	2 HOURS	21 MINUTES	5 MINUTES	25 SECONDS
	30 SECONDS	4 WEEKS	3 DAYS	12 HOURS	2 HOURS	30 MINUTES	2 MINUTES
	1 MINUTE	8 WEEKS	6 DAYS	1 DAY	4 HOURS	1 HOUR	5 MINUTES
	5 MINUTES	9 MONTHS	4 WEEKS	6 DAYS	21 HOURS	5 HOURS	25 MINUTES
	30 MINUTES	6 MONTHS	5 WEEKS	5 DAYS	1 DAY	2 HOURS	
	1 HOUR	10 MONTHS	2 MONTHS	10 DAYS	2 DAYS	5 HOURS	
	6 HOURS			2 MONTHS	2 WEEKS	1 DAY	
	1 DAY				8 WEEKS	5 DAYS	

Optimisation

- Speed tips: always use numpy array operations where possible, and avoid huge for loops
- There are some fancy ways to do profiling, but in my experience, using `time.time()` is enough to help you find the slow parts of your code
- Once found, concentrate on those parts (e.g., if a single routine takes 100 sec, and another takes 0.1 sec, which will you focus on? It depends if you're calling the 0.1 sec routine a million times and the 100 sec routine once)
- For calculations that run once and not often (e.g., in part of a pipeline) cache the result somewhere and load it from disk for subsequent re-runs
- If you really need to do something element-by-element on a huge array, here are a couple of possible options
 - Cython (<http://cython.org/>) - I've used this several times
 - Numba (<http://numba.pydata.org/>)

```
1 """
2
3 Example of using time.time() for checking which parts of the code
4 take the longest to run.
5 """
6
7
8 import time
9 import numpy as np
10
11 t0=time.time()
12 A=np.random.uniform(0, 1, [1000, 2000])
13 B=np.random.uniform(0, 1, [A.shape[1], 300])
14 t1=time.time()
15 print("Making matrices: %.3f sec" % (t1-t0))
16
17 C=np.dot(A, B)
18 t2=time.time()
19 print("Matrix multiplication: %.3f sec" % (t2-t1))
20
21 pInv=np.linalg.pinv(A)
22 t3=time.time()
23 print("Pseudo-inverse: %.3f sec" % (t3-t2))
24
25 u, s, vh=np.linalg.svd(A)
26 t4=time.time()
27 print("SVD: %.3f sec" % (t4-t3))
28
29 print("All: %.3f sec" % (t4-t0))
30
31 # Output when running this code on my laptop:
32 #
33 # Making matrices: 0.054 sec
34 # Matrix multiplication: 0.094 sec
35 # Pseudo-inverse: 1.642 sec
36 # SVD: 1.931 sec
37 # All: 3.721 sec
38 #
39 # These are not super-accurate benchmarks, but the idea is to look at the
40 # order-of-magnitude, if you are deciding which parts of the code could use
41 # a speed-up (if possible)
42 #
43 # For individual routines like this, you could use %timeit in IPython instead
44 # (which will give accurate benchmarks)|
```

Parallelisation

- Do you need it? Optimise your code first
- There is the multiprocessing module in the Standard Library:
<https://docs.python.org/3.3/library/multiprocessing.html>
... but I find to use it generally requires major surgery on my code
- MPI4Py:
<http://mpi4py.readthedocs.io/en/stable/>
for me is the easiest / quickest way – and even if you're running on a multicore laptop and not a cluster, you can still use it (and MPI is very easy to install on Ubuntu)
e.g., `mpirun --np 4 your_python_program`
- It is also possible to use threads, rather than processes (advantage: less memory usage; disadvantage: more complicated than MPI, depending on what you are trying to do)

Some tools for collaboration

- Git:
 - Distributed version control: allows many users to work on the same code, and has tools to merge / fork / compare code pretty easily; Open Source
- Github (<https://github.com/>):
 - A web service that hosts Git source-code repositories, with a nice, no-nonsense interface. Free for public, Open Source repositories
 - Alternative services are available (e.g., Sourceforge, Gitlab) that can also host Git repositories
- Slack (<https://slack.com/>):
 - Seems to have become the leading real-time messaging app for code development
 - Free (if you only keep ~10,000 messages) and cross-platform but proprietary software

Github (or Sourceforge or...)

- The most popular site for hosting code repositories in the world is Github (Microsoft bought it the other week)

The screenshot shows a GitHub repository page for 'ACTCollaboration/nemo'. The repository is private, has 168 commits, 2 branches, 1 release, and 3 contributors. The latest commit was 6 days ago. The repository contains files like README.md, setup.py, README.in, MANIFEST.in, examples, and bin.

ACTCollaboration / nemo Private

Code Issues 0 Pull requests 0 Projects 0 Wiki Insights Settings

Map filtering and SZ cluster detection and characterization pipeline Edit

Add topics

168 commits 2 branches 1 release 3 contributors

Branch: master ▾ New pull request Create new file Upload files Find file Clone or download ▾

File	Description	Time
bin	Added mapInterpolator option, and updated examples	13 days ago
examples	Fix for non-tileDeck runs to work with 3D Enki maps	7 days ago
nemo	Fix for overlap regions in tileDeck mode when no surveyMask given	6 days ago
MANIFEST.in	Switched another atpy import for astropy. Added PIL to dependencies l...	a year ago
README.md	Added more selection function plots	15 days ago
setup.py	Switched from pyfits to astropy.io.fits for handling images, as	3 months ago

README.md

Releasing your code

- Where?
- What license?
- Is it easy to use?
- Did you write some nice documentation?

* * * * * Main Page * Download * Installation * Documentation * Bugs and Feature Requests *

astLib Python astronomy modules

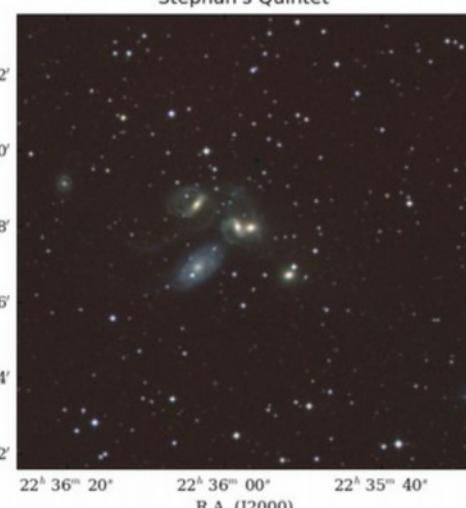
astLib is a set of Python modules that provides some tools for research astronomers. It can be used for astronomical plots, some statistics, common calculations, coordinate conversions, and manipulating FITS images with World Coordinate System (WCS) information through PyWCSTools - a simple wrapping of WCSTools by Jessica Mink. PyWCSTools is distributed (and developed) as part of astLib.

Version 0.10.0 released
Mon, 03/05/2018 - 07:49 — Matt Hilton

Version 0.10.0 of astLib has been released, and is available from both SourceForge and PyPI.

This release breaks compatibility with pyfits (which is deprecated by STScI), and now requires astropy. If you still require a version that works with pyfits, please continue to use release 0.9.3.

Stephan's Quintet

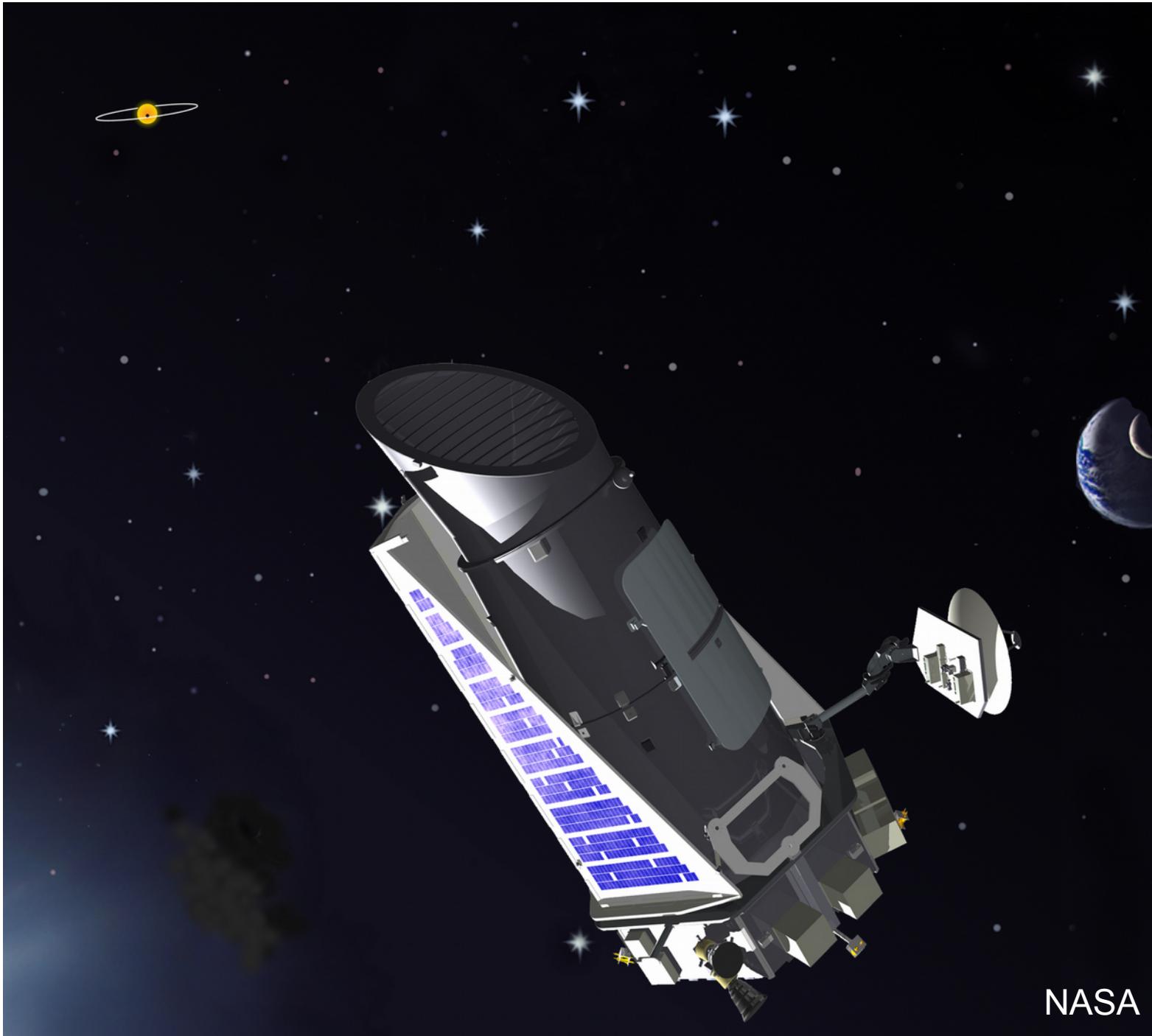


A grayscale astronomical image showing the Stephan's Quintet galaxy cluster. The cluster consists of five galaxies of varying sizes and luminosities. The image is overlaid with coordinate axes: Dec. (J2000) on the left and R.A. (J2000) at the bottom. The Dec. axis ranges from +33°52' to +34°02' with labels every 2'. The R.A. axis ranges from 22°36'20" to 22°35'40" with labels every 20". The title "Stephan's Quintet" is centered above the image.

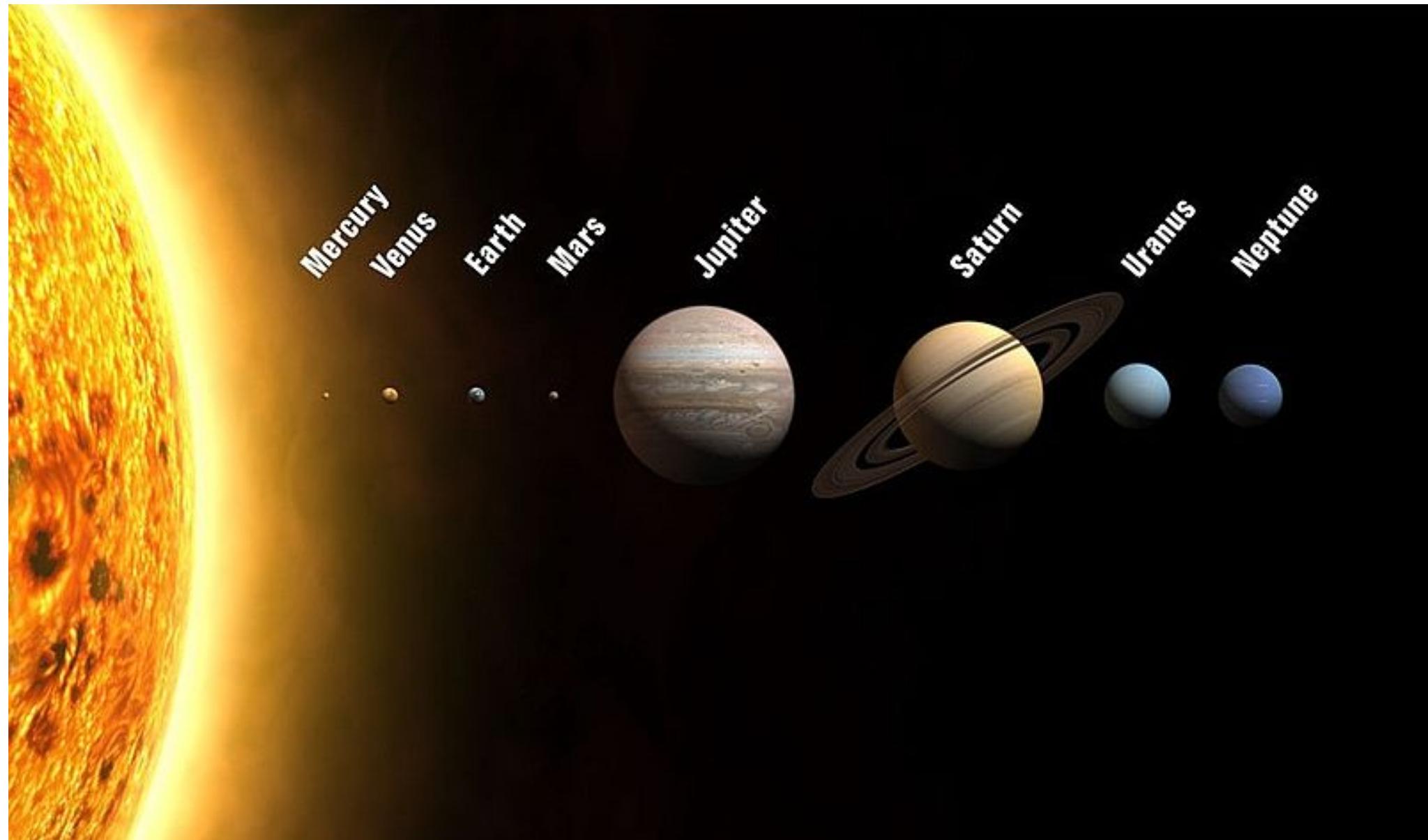
Example output from the `astPlots.ImagePlot` class in version 0.3.1. Image data taken from the Digitized Sky Survey, using the SkyView interface.

<http://astlib.sourceforge.net/> - more than 10 years old...

... and now for some science ...



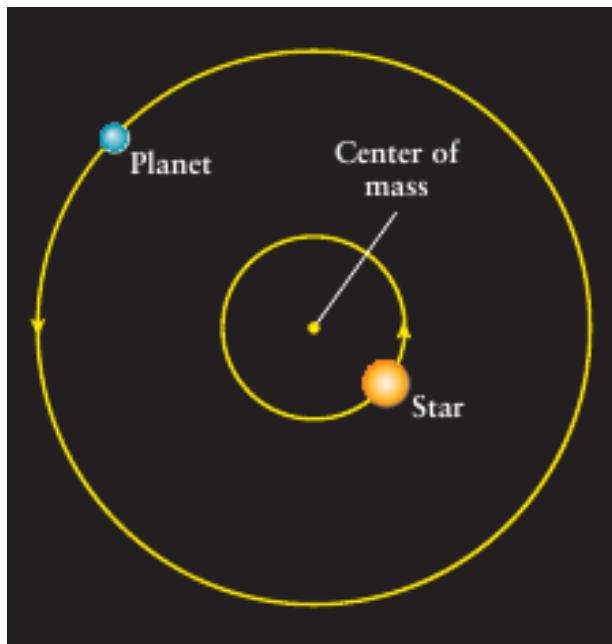
Our Solar System



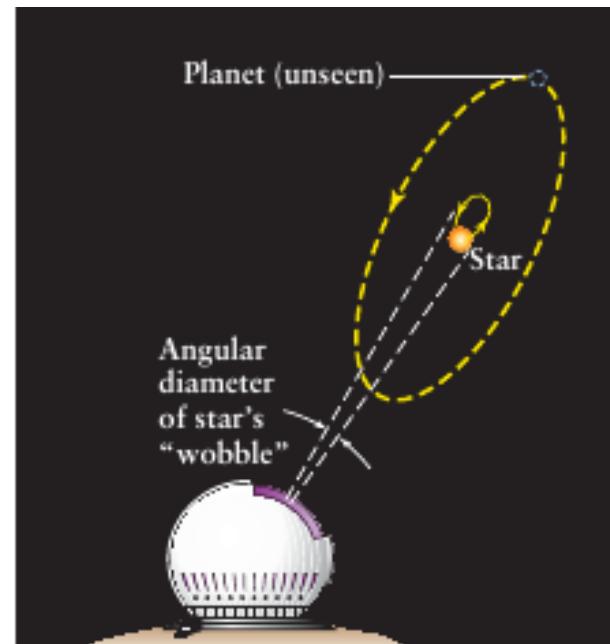
Sizes of planets are to scale – the distances between them are not!

Extrasolar planets

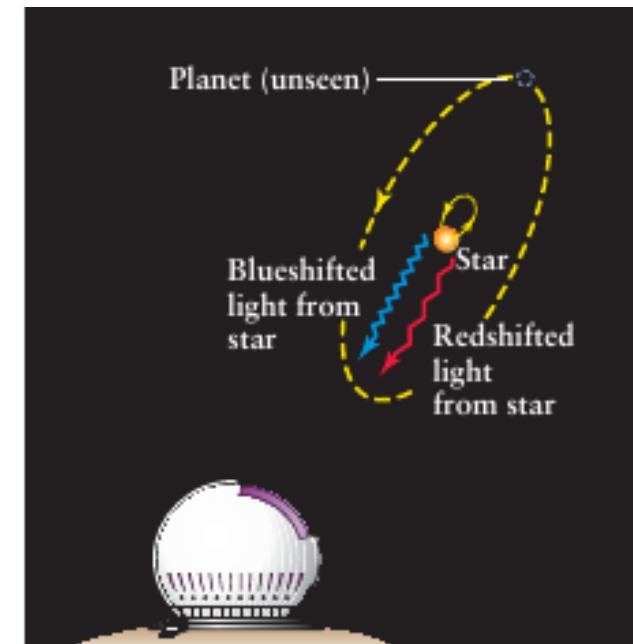
- Since 1995, 3594 planets have been discovered around stars other than the Sun (exoplanet.eu)
- This is very challenging observationally, because stars are \sim 1 billion times brighter than planets
- One way to find extrasolar planets is to look for the wobble of their parent star – this is the **radial velocity method**



(a) A star and its planet



(b) The astrometric method



(c) The radial velocity method

Radial velocity method

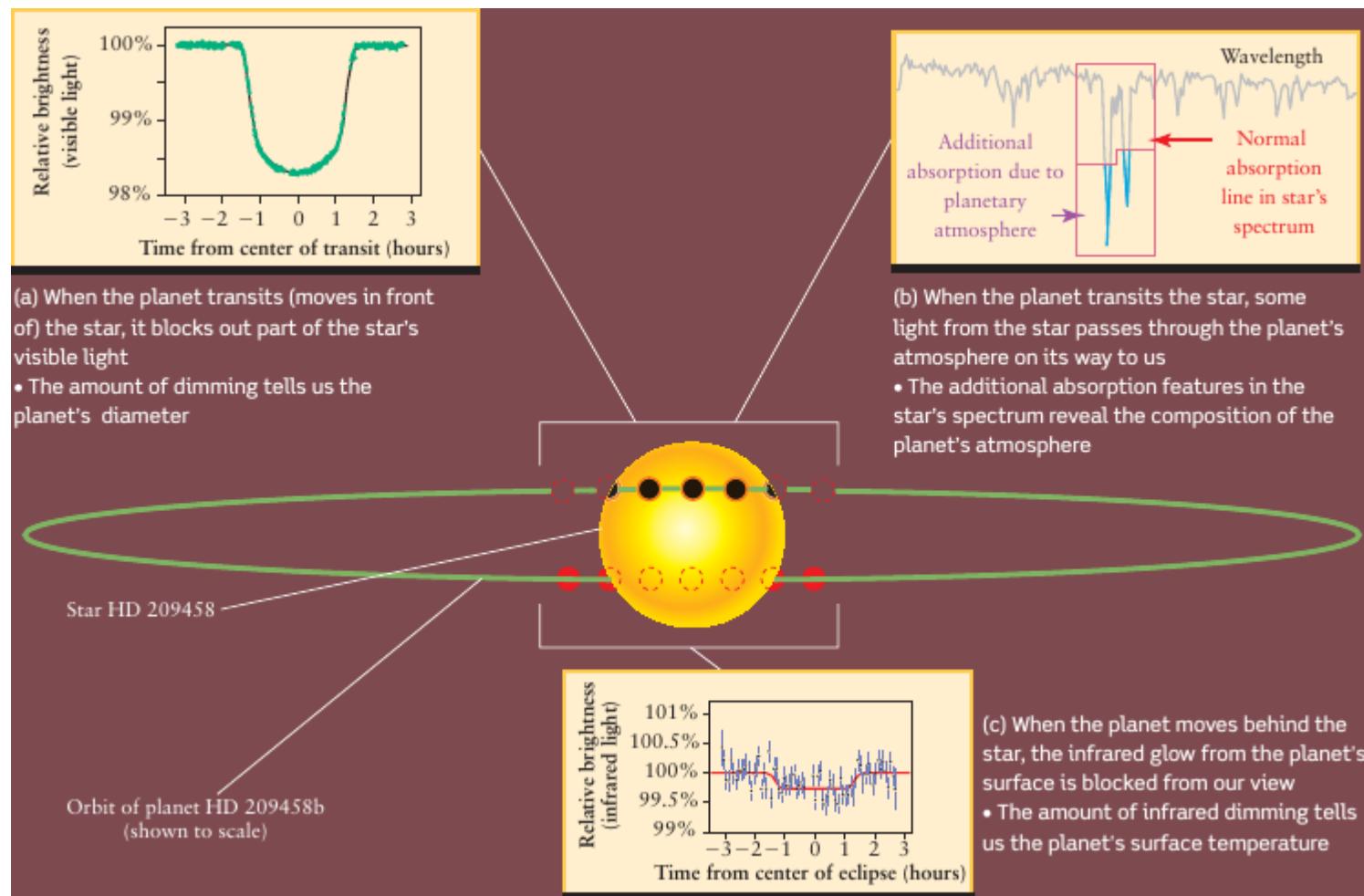


PLANETQUEST
THE SEARCH FOR ANOTHER EARTH

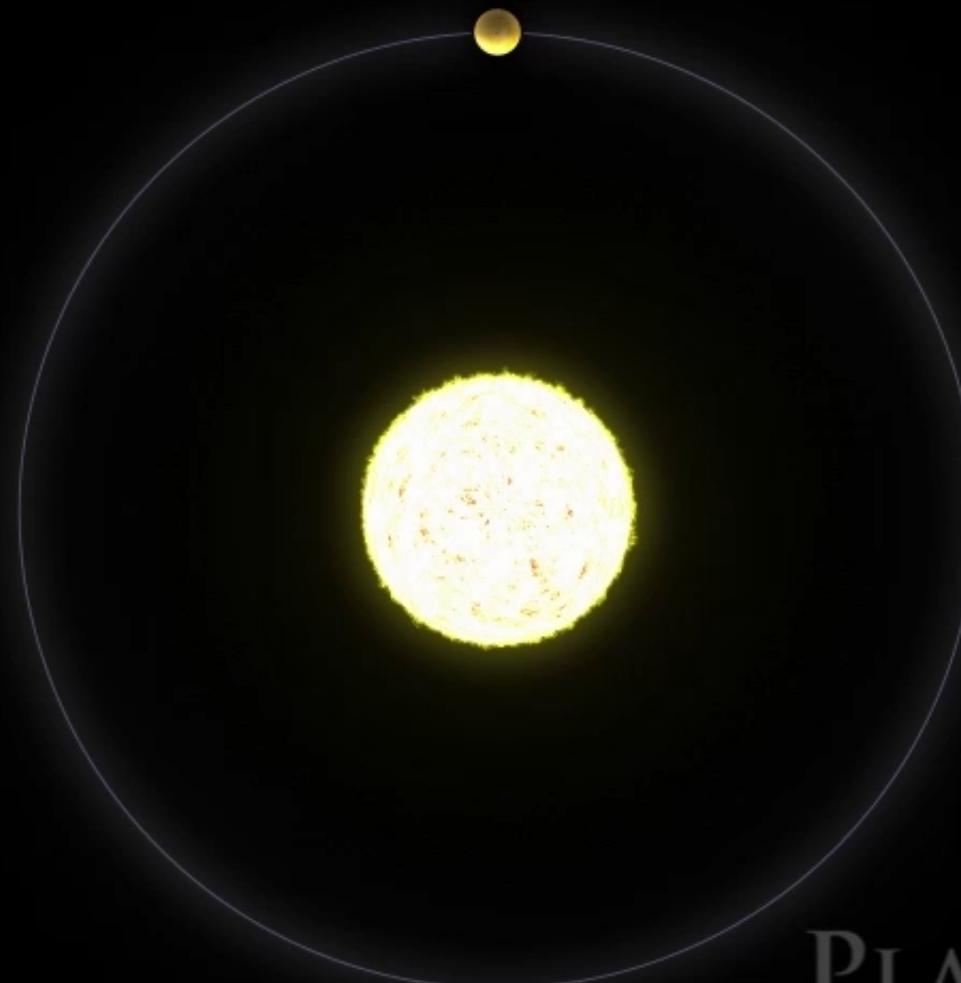
exoplanets.nasa.gov

Transit searches

- Radial velocity searches give an estimate of the mass of an exoplanet, but not its size – for that, **transit searches** can be used
- When an exoplanet transits in front of its parent star, it dims its light for a short period



Transit searches

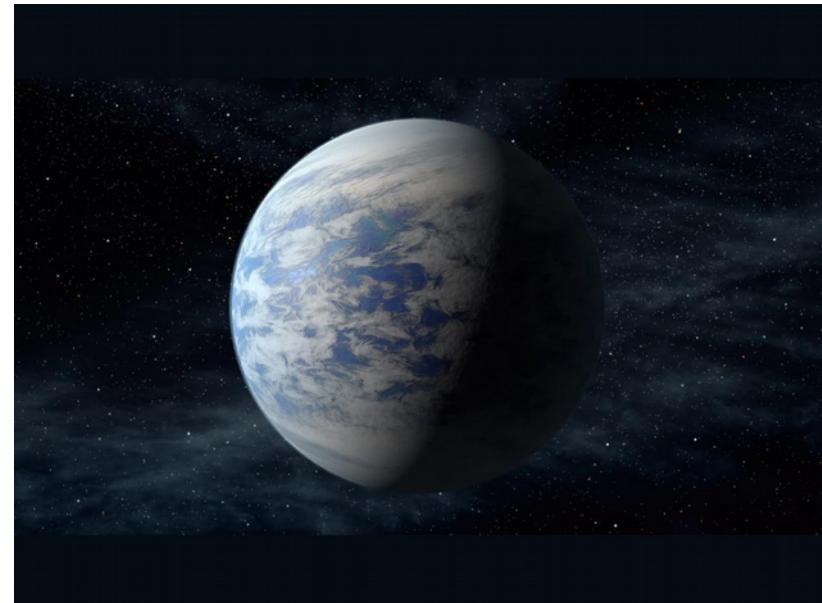


PLANETQUEST
THE SEARCH FOR ANOTHER EARTH

exoplanets.nasa.gov

Properties of exoplanets

- The combination of radial velocity measurements of planets detected using the transit technique gives us both the mass and size of planets, and hence some idea of their average density
- One surprising finding is that many extrasolar planets have masses similar to Jupiter and are orbiting very close to their parent stars
- These objects are called “**hot Jupiters**” – they presumably must have migrated inwards from the edges of their solar systems
- One goal of exoplanet searches is to find planets similar to the Earth to look for signs of **extraterrestrial life**...



Kepler/NASA/JPL

The tutorial...

- We'll go through the process of setting up a Python application to play with Kepler light curve data – this will take you through:
 - Basic Git usage
 - Setting up a pure-Python package (with modules and an executable script)
 - Adding a nice command-line interface
 - Loading and processing Kepler light curve data
- It's open ended, because I leave it to you to figure out how to fit the data...
- You can find it here:

<https://github.com/mattyowl/AIMS-DISCnet-2018>