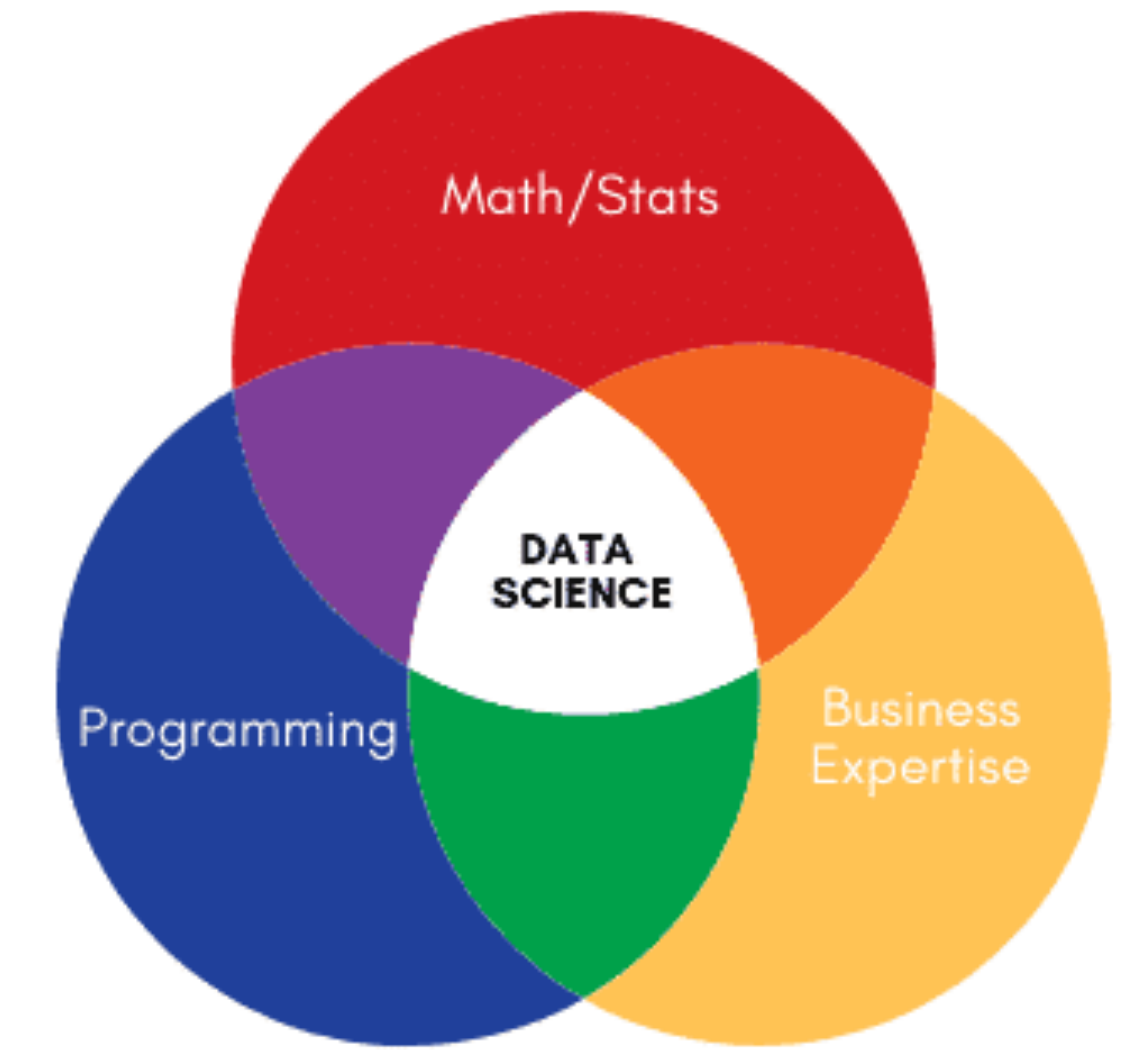# Data Science Challenges in Time Domain Astronomy: Building Methods, Tools and Communities
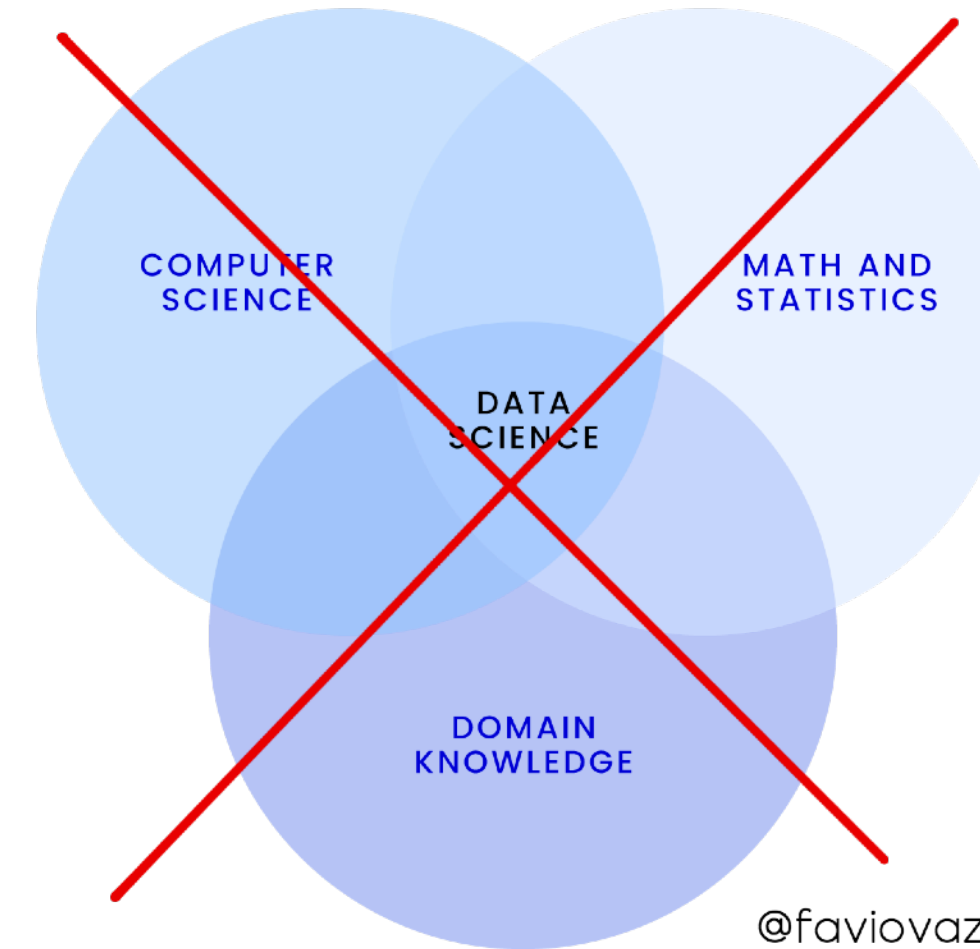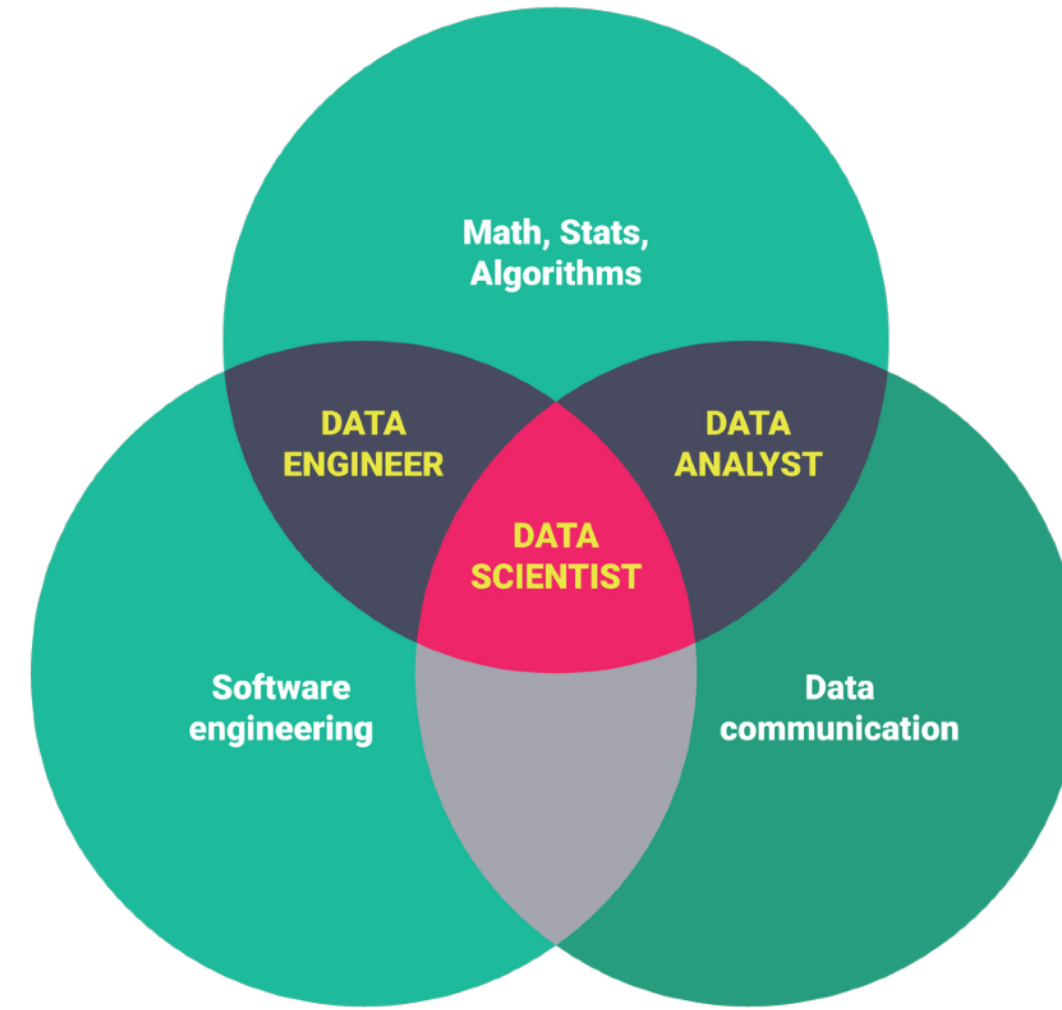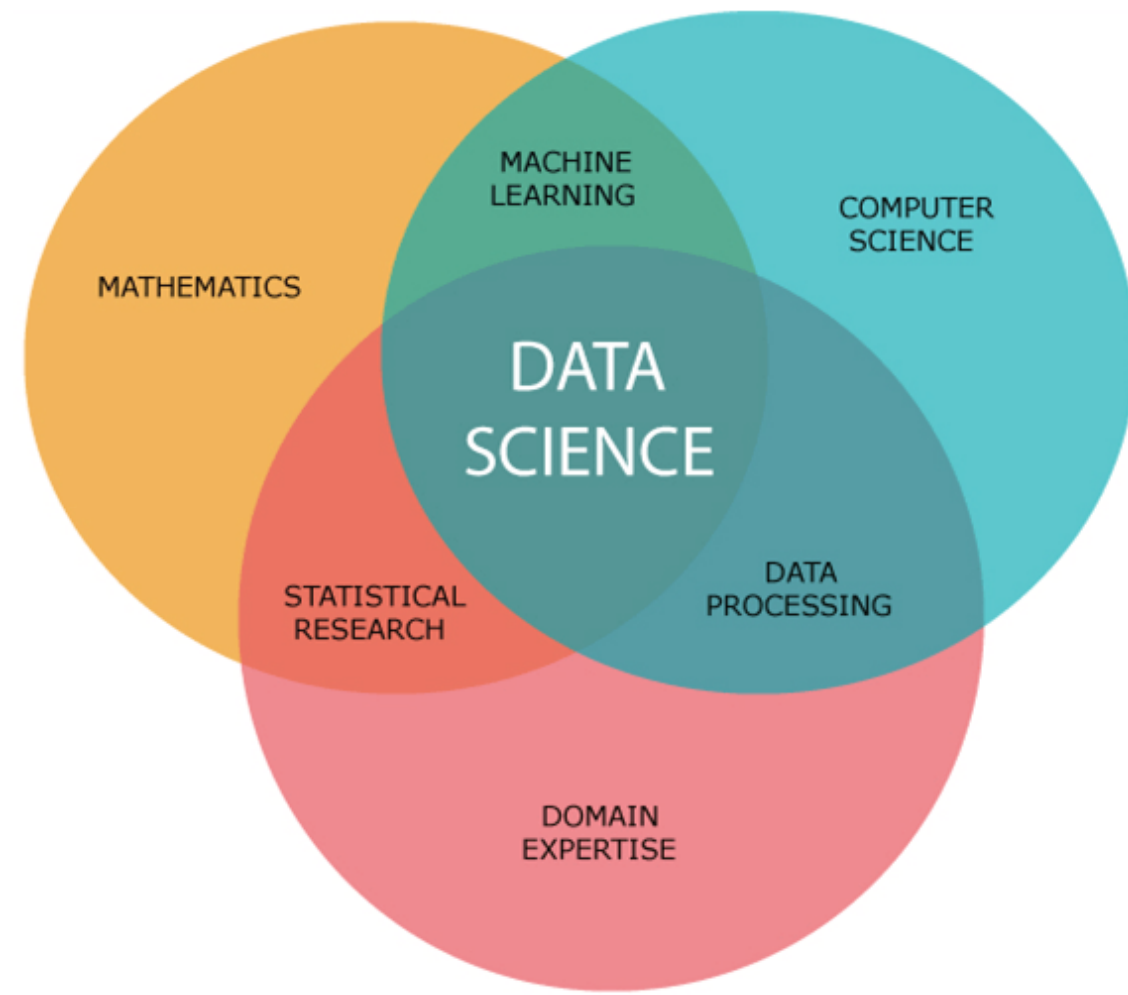
Daniela Huppenkothen

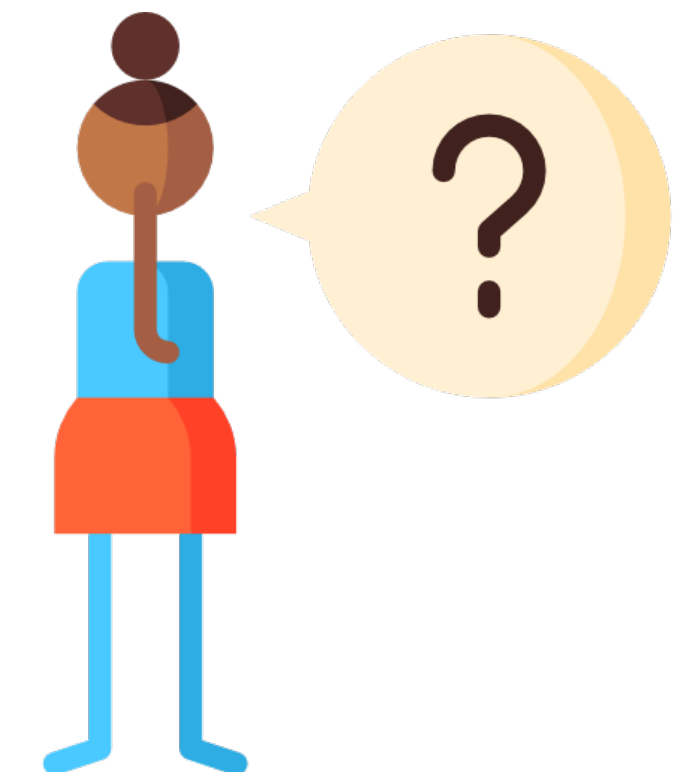DIRAC, University of Washington
eScience Institute, University of Washington

Tiana_Athriel

dhuppenkothen

dhuppenk@uw.edu

# Why data science?

Data science is the art of drawing Venn diagrams?

James Davenport @jradavenport · Jun 7

W/ the help of @jakevdp, here are 200k @ESAGaia + @WISE_Mission stars w/ *predicted* [Fe/H], trained on 30k @APOGEEsurvey 🤩🤩

#GaiaSprint

# Gaia

## 1.7 billion stars

credit: LSST/NOAO

# LSST

40 billion sources
10 million alerts/night

# SKA

160TB raw data/second

# ATHENA

**100x** effective area of Chandra

**7.5** arcmin field of view

... but it's not all **big data!**
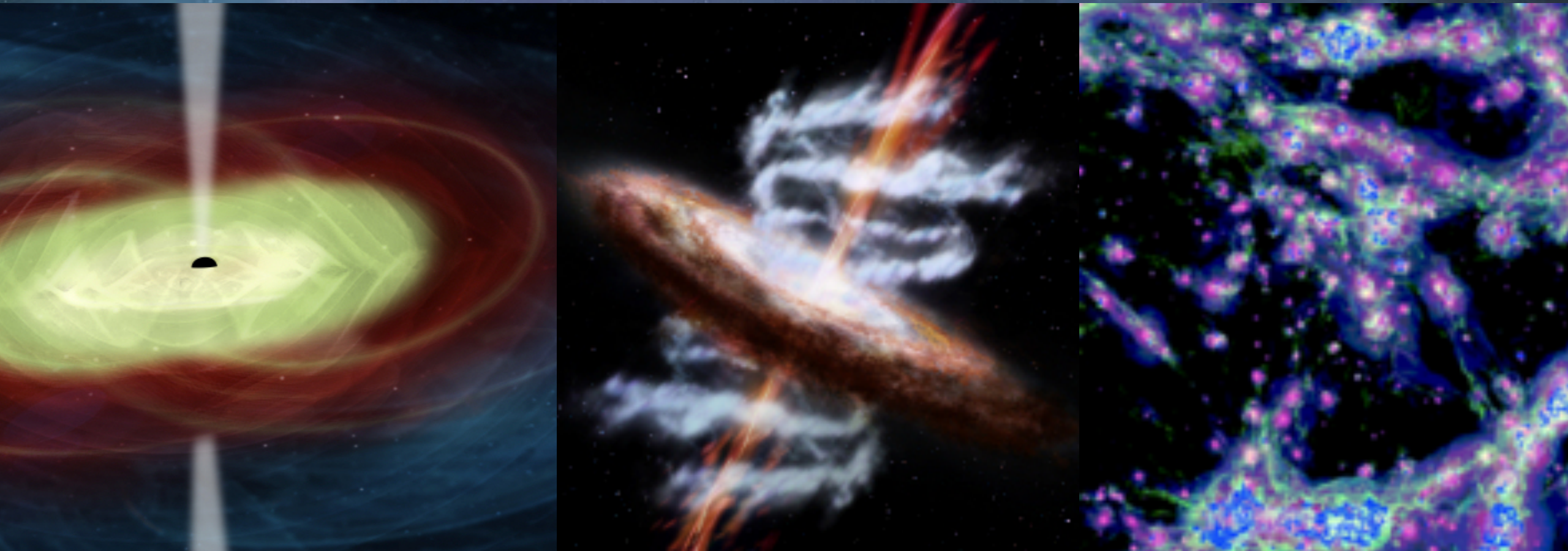
"**Annoyingly not small (and complex) data**"

— Genevieve Graves

# How do we make sense of all this data?

| | Data Science Task | | |
|---|---|---|---|
| | **Description** | **Prediction** | **Causal inference** |
| Example of scientific question | How can we split supernovae into different classes? | What is the probability for a supernovae with a given spectrum/light curve to be a Type a? | Are Type Ia supernovae caused by the explosion of white dwarfs? |
| Data | • Eligibility criteria (is it a supernova?)<br>• features (spectral lines, light curve shape, …) | • Eligibility criteria (is it a supernova?)<br>• Output (classes of supernovae)<br>• Input (spectral lines, light curve shape, …) | • Eligibility criteria (is it a supernova?)<br>• White dwarf explosion models<br>• Features (spectral lines, light curve shape, …) |
| Examples of analytics | Cluster analysis<br>… | Regression<br>Decision trees<br>Random forests<br>Support vector machines<br>Neural networks<br>… | Regression<br>Matching<br>Inverse probability weighting<br>G-formula<br>G-estimation<br>Instrumental variable estimation<br>… |

adapted from Hernan et al (2019)

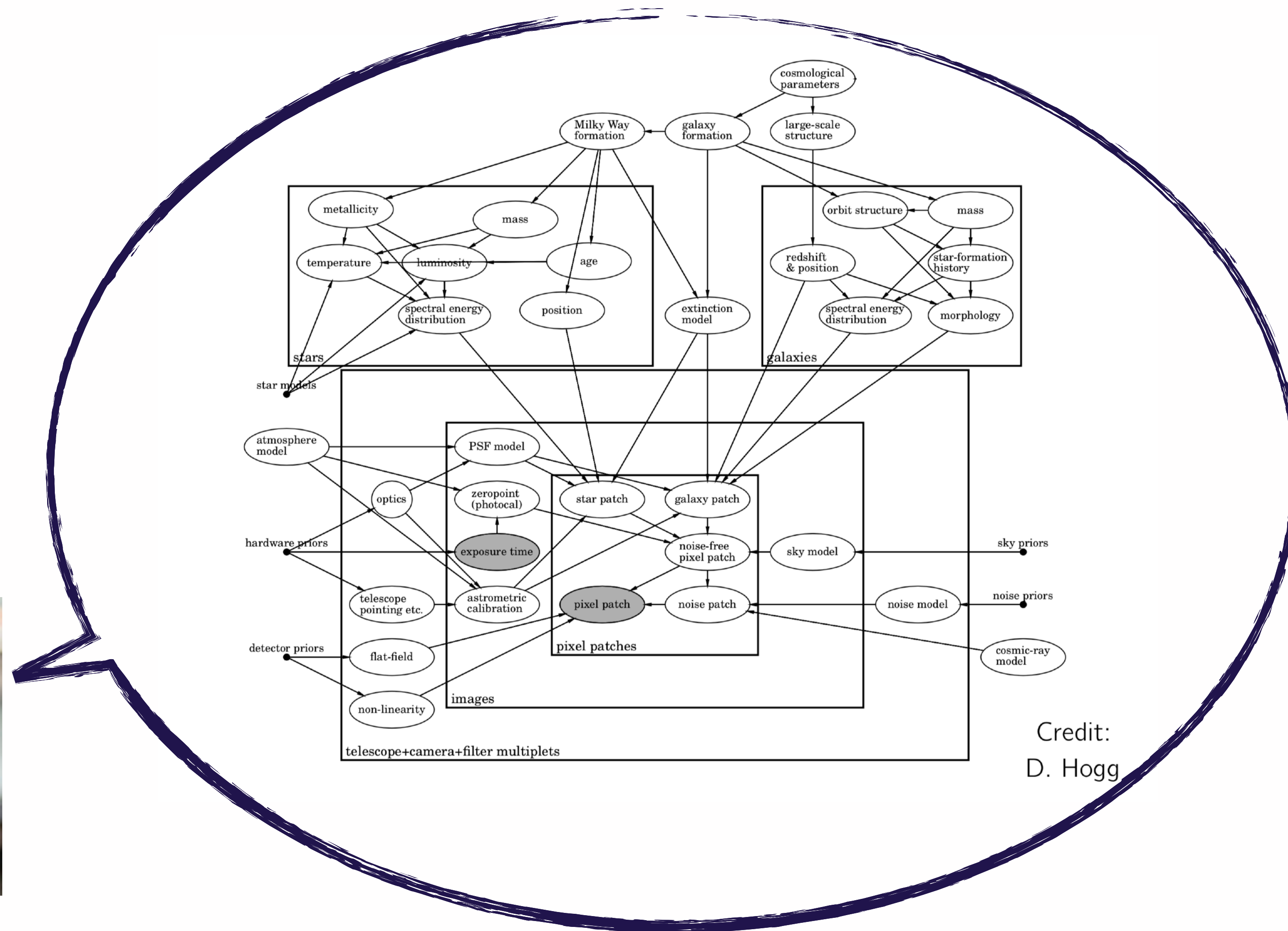| | Data Science Task | | |
|---|---|---|---|
| | **Description** | **Prediction** | **Causal inference** |
| Example of scientific question | How can we split supernovae into different classes? | What is the probability for a supernovae with a given spectrum/light curve to be a Type a? | Are Type Ia supernovae caused by the explosion of white dwarfs? |
| Data | <ul><li>Eligibility criteria (is it a supernova?)</li><li>features (spectral lines, light curve shape, …)</li></ul> | <ul><li>Eligibility criteria (is it a supernova?)</li><li>Output (classes of supernovae)</li><li>Input (spectral lines, light curve shape, …)</li></ul> | <ul><li>Eligibility criteria (is it a supernova?)</li><li>White dwarf explosion models</li><li>Features (spectral lines, light curve shape, …)</li></ul> |
| Examples of analytics | Cluster analysis … | Regression<br>Decision trees<br>Random forests<br>Support vector machines<br>Neural networks<br>… | Regression<br>Matching<br>Inverse probability weighting<br>G-formula<br>G-estimation<br>Instrumental variable estimation<br>… |

adapted from Hernan et al (2019)

DiRAC

## Data Science Task

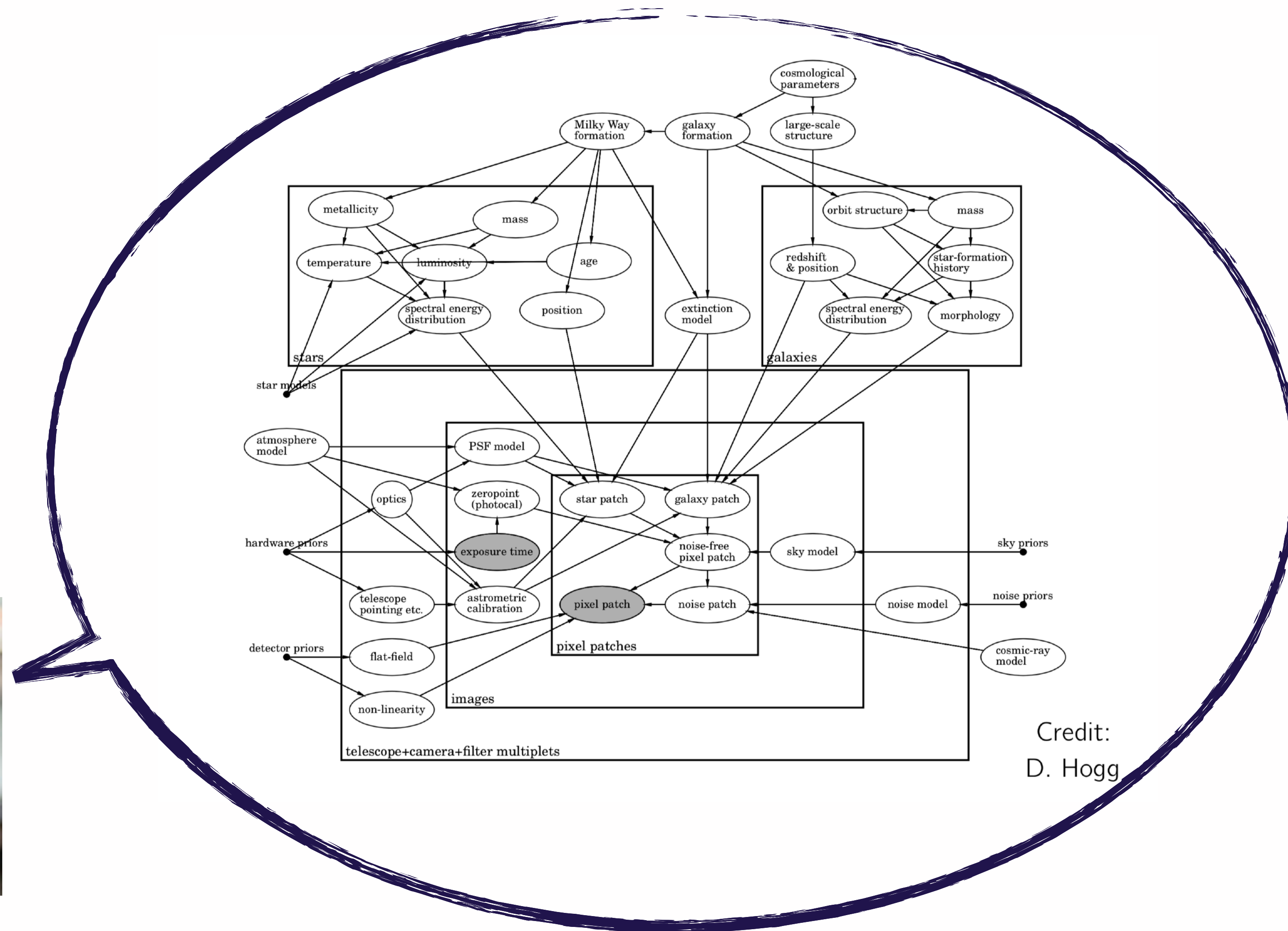| | Description | Prediction | Causal inference |
|---|---|---|---|
| Example of scientific question | How can we split supernovae into different classes? | What is the probability for a supernovae with a given spectrum/light curve to be a Type a? | Are Type Ia supernovae caused by the explosion of white dwarfs? |
| Data | • Eligibility criteria (is it a supernova?)<br>• features (spectral lines, light curve shape, …) | • Eligibility criteria (is it a supernova?)<br>• Output (classes of supernovae)<br>• Input (spectral lines, light curve shape, …) | • Eligibility criteria (is it a supernova?)<br>• White dwarf explosion models<br>• Features (spectral lines, light curve shape, …) |
| Examples of analytics | Cluster analysis … | Regression | Regression … |

**Question to ponder**: where does this distinction **break down**?

DiRAC

Credit:
D. Hogg

Credit: D. Hogg

Maybe I can help?

Credit: D. Hogg

Artificial
Intelligence

Terminator - Rise of The Machines

DiRAC

Are **astronomers going to be** replaced by a **neural network?**

Artificial
Intelligence
Terminator - Rise of The Machines

Are **astronomers** going to be replaced by a **neural network?**

(no)

# PHOTOMETRIC LSST ASTRONOMICAL TIME-SERIES CLASSIFICATION CHALLENGE (PLASTICC)

# PHOTOMETRIC LSST ASTRONOMICAL TIME-SERIES CLASSIFICATION CHALLENGE (PLASTICC)
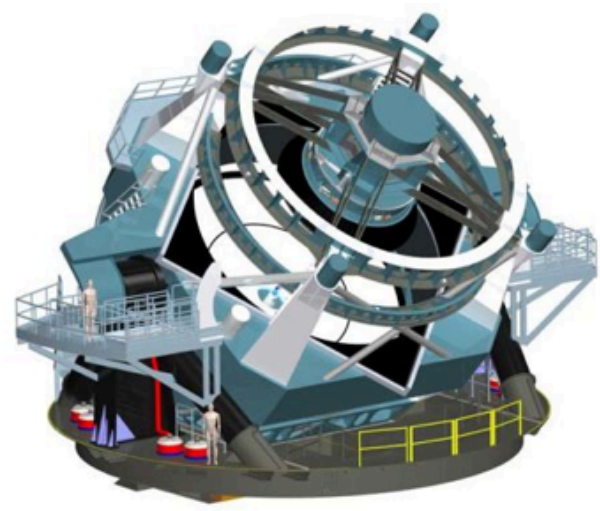
## Overview of 1st place solution

posted in **PLAsTiCC Astronomical Classification** 5 months ago

⬆ 190

**Kyle Boone**
1st place

EDIT: code is now available on my github page at https://github.com/kboone/avocado

First of all, thanks to everyone who participated in this competition! I learned a lot doing it, and I have enjoyed all of the discussions that I had with you. Here is an overview of my model that took 1st place in this competition. I will be releasing the code with a full writeup shortly.

I am an astronomer studying supernova cosmology, so my work mainly focused on trying to tell the different supernova types apart. This ended up working out well because everything else was fairly easy to tell apart. Here is a summary of my solution:

# PHOTOMETRIC LSST ASTRONOMICAL TIME-SERIES CLASSIFICATION CHALLENGE (PLASTICC)

"However, **sound causal inference** requires not only **adequate data analysis techniques** but also **subject-matter expertise** about the causal structure of the problem under study"

(Hernan, 2019: "Comment: Spherical Cows in a Vacuum")

## Overview of 1st place solution
posted in PLAsTiCC Astronomical Classification 5 months ago

⬆ 190

**Kyle Boone**
1st place

EDIT: code is now available on my github page at https://github.com/kboone/avocado

First of all, thanks to everyone who participated in this competition! I learned a lot doing it, and I have enjoyed all of the discussions that I had with you. Here is an overview of my model that took 1st place in this competition. I will be releasing the code with a full writeup shortly.

I am an astronomer studying supernova cosmology, so my work mainly focused on trying to tell the different supernova types apart. This ended up working out well because everything else was fairly easy to tell apart. Here is a summary of my solution:
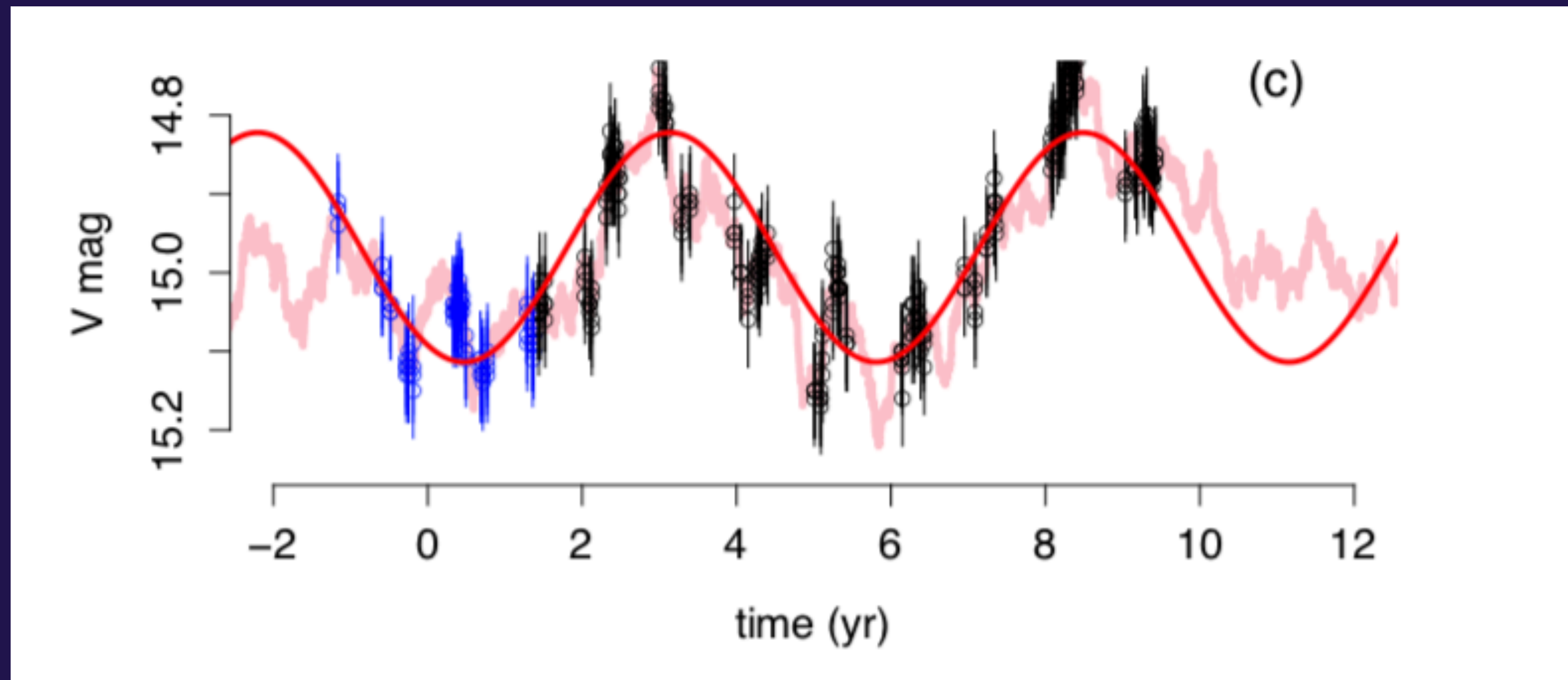
Boone (2019): 1907.04690

# Data Science @ ADASS

- **Yanxia Zhang**: Photometric Redshift Estimation of Quasars by Machine Learning

- **Sweta Singh**: Scientific Visualisation of Extremely Large Distributed Astronomical Surveys

- **Shraddha Surana**: Machine Learning for Scientific Discovery

- **Antonia Rowlinson**: Identifying transient and variable sources in radio images

- **Lightning Session 2**: Automated Bayesian Inference, Supernova detection with deep learning, visualization of virual observatory data, Gaussian process modelling

- **Sessions 5a** and **5b**: data visualization

- Lots of **posters**!

# A Statistical View of Our Data

# Our data is often ...

## stochastic vs deterministic



Vaughan et al (2016)

# Our data is often ...

- **unevenly sampled**
- **heteroscedastic**



Graham et al (2015)

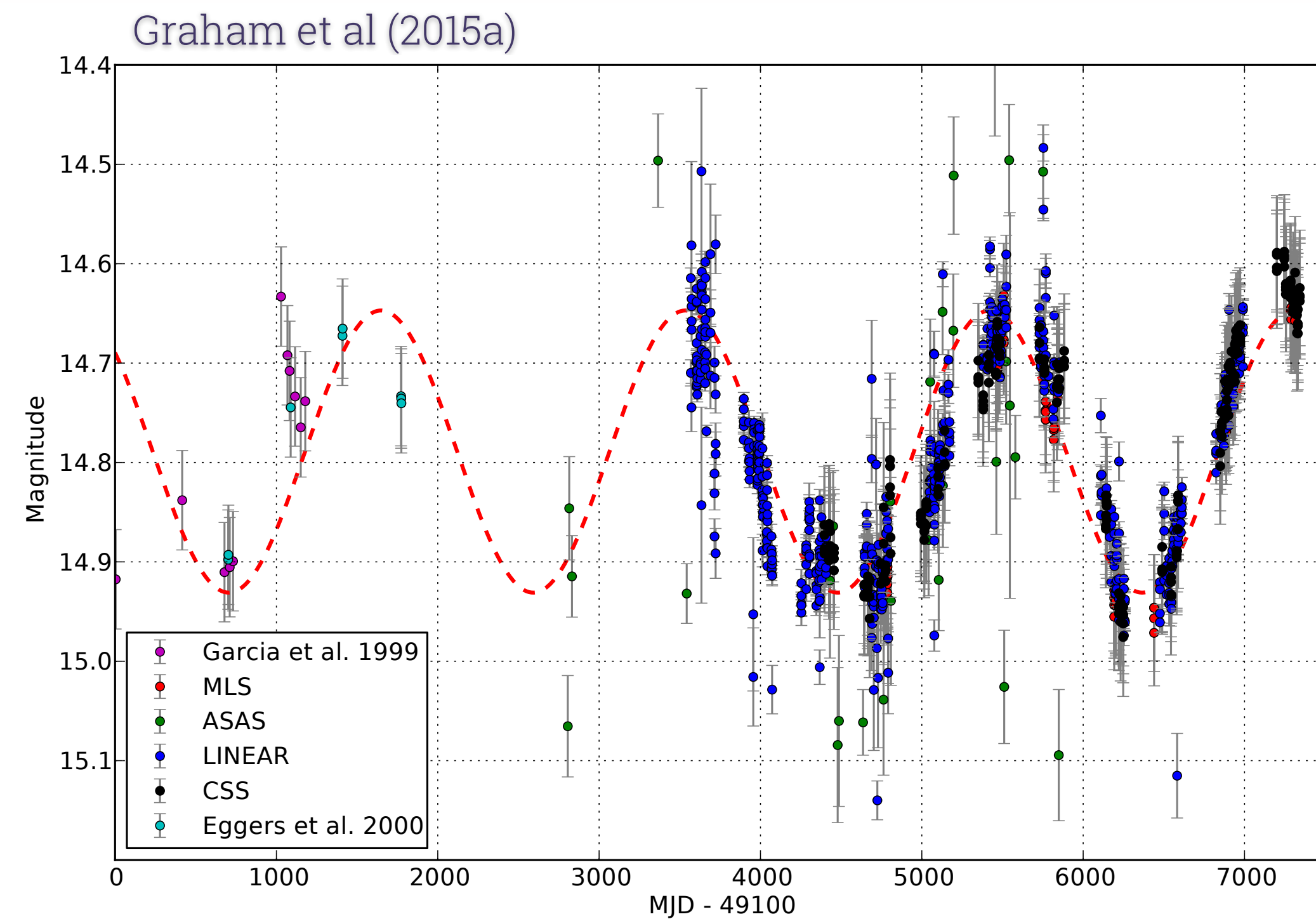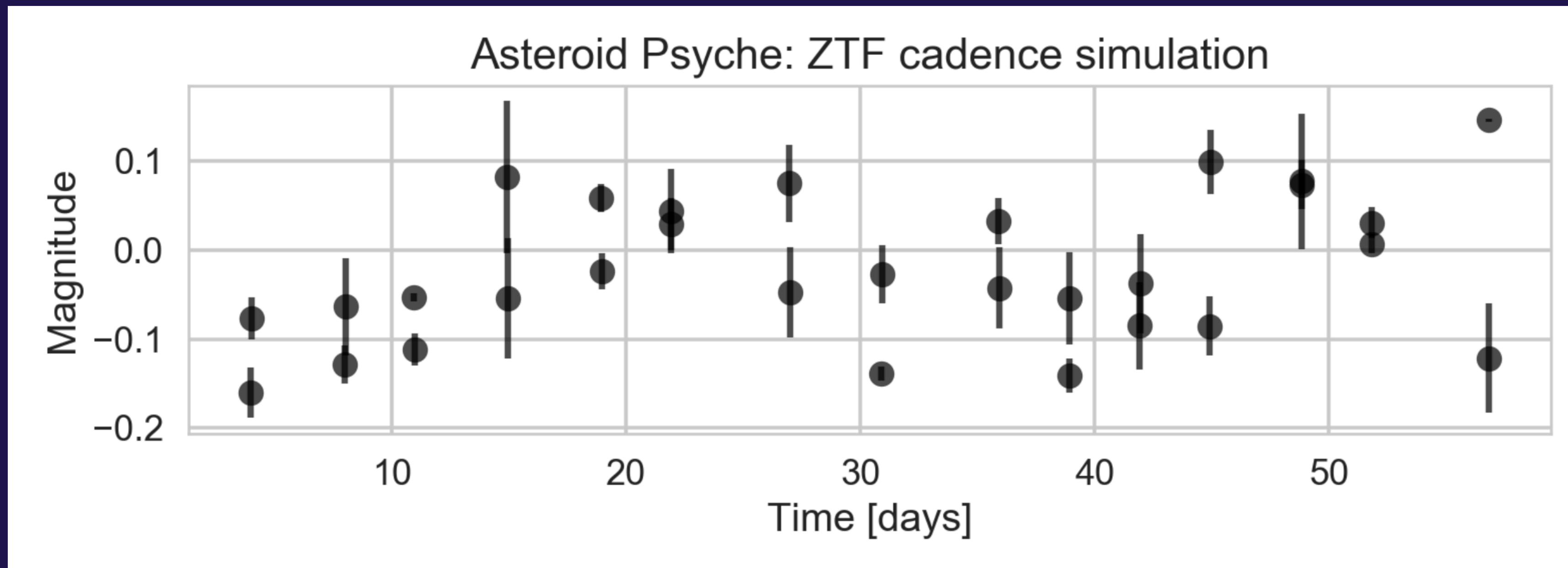# Finding Supermassive Black Hole Binaries: The curious case of PG 1302-102



Graham et al (2015a)

# Finding Supermassive Black Hole Binaries: The curious case of PG 1302-102
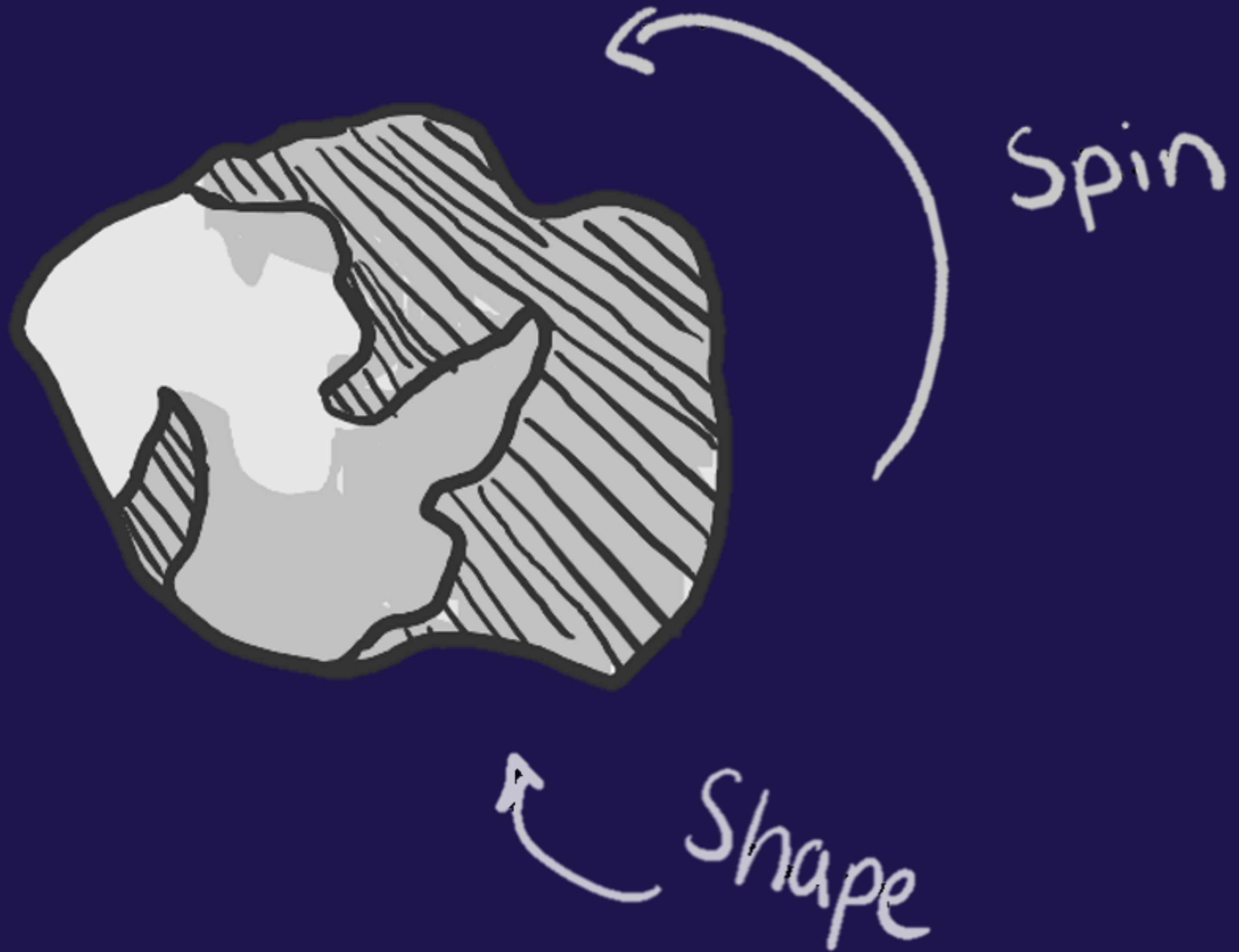
**LSST will extend our baseline by a factor of ~1.5!**

# Our data is often ...

## sparse



Asteroid Psyche: ZTF cadence simulation

SCIENCE

## *An Interstellar Visitor Both Familiar and Alien*

Leer en español

**Dennis Overbye**
**OUT THERE**   NOV. 22, 2017



A Glimpse of Oumuamua
By DENNIS OVERBYE, JONATHAN CORUM and JASON DRAKEFORD

▶   1:36  ————————————●————————————  2:53   HD   🔊   ⛶

Astronomers have discovered a passing rock from another star — the first interstellar asteroid. By DENNIS OVERBYE, JONATHAN CORUM and JASON DRAKEFORD on December 12, 2017. . Watch in Times Video »

Visit the galaxy before the galaxy visits you.

This fall, the galaxy came calling in the form of a small reddish cigar-shaped object named Oumuamua by astronomers based in Hawaii. They discovered it in October, careening through the solar system at
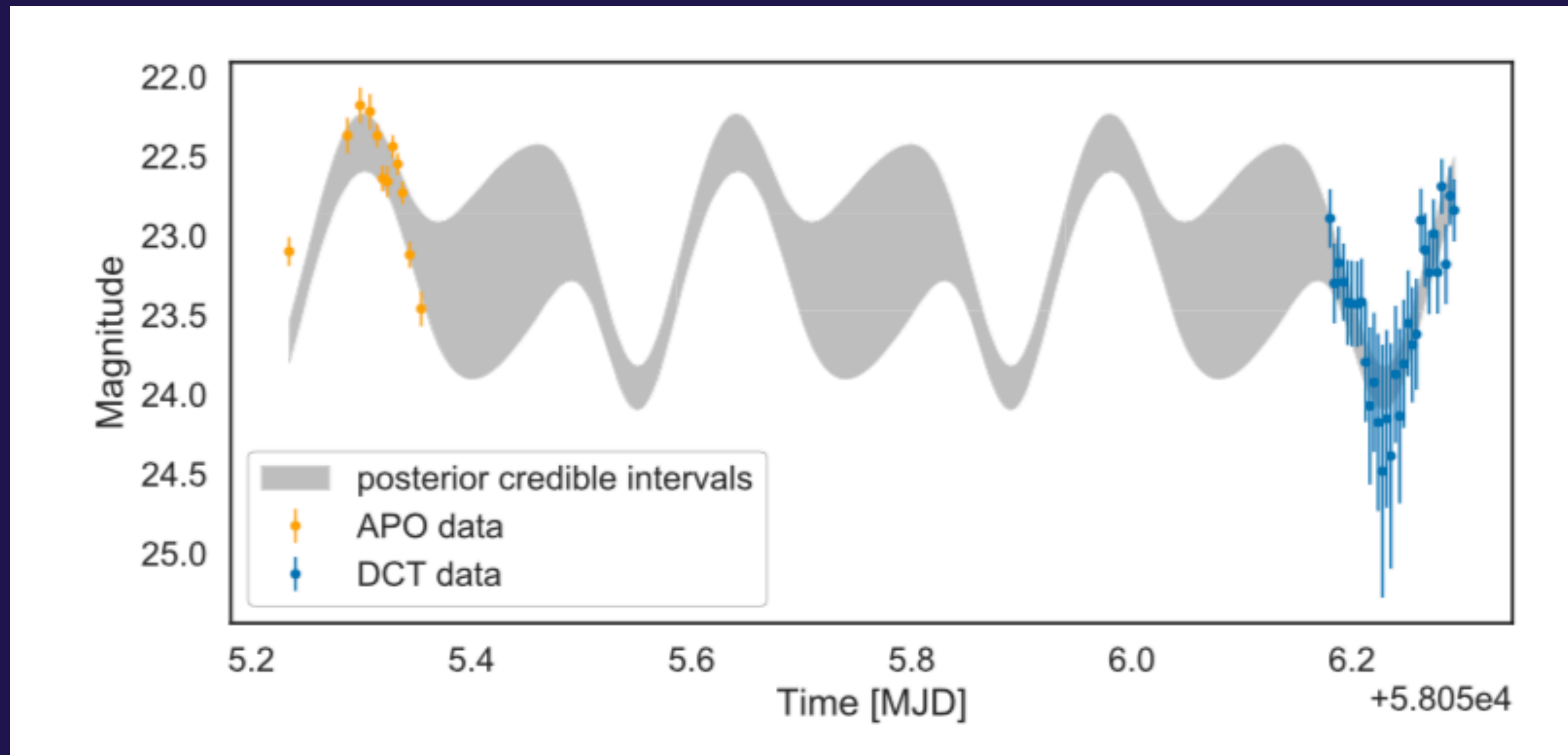
# 1I'Oumuamua



**Discovery Channel Telescope +
Apache Point Observatory 3.5m**

# Model: Gaussian Process with periodic kernel

Bolin et al, incl Huppenkothen (2017)

# Model: Gaussian Process with periodic kernel



Bolin et al, incl Huppenkothen (2017)

# **Model:** Gaussian Process with periodic kernel



better period estimate with less data

result: 1I'Oumuamua is **extremely elongated**, and **probably tumbling**

# Gaussian Processes for Sparse Asteroid Light Curves



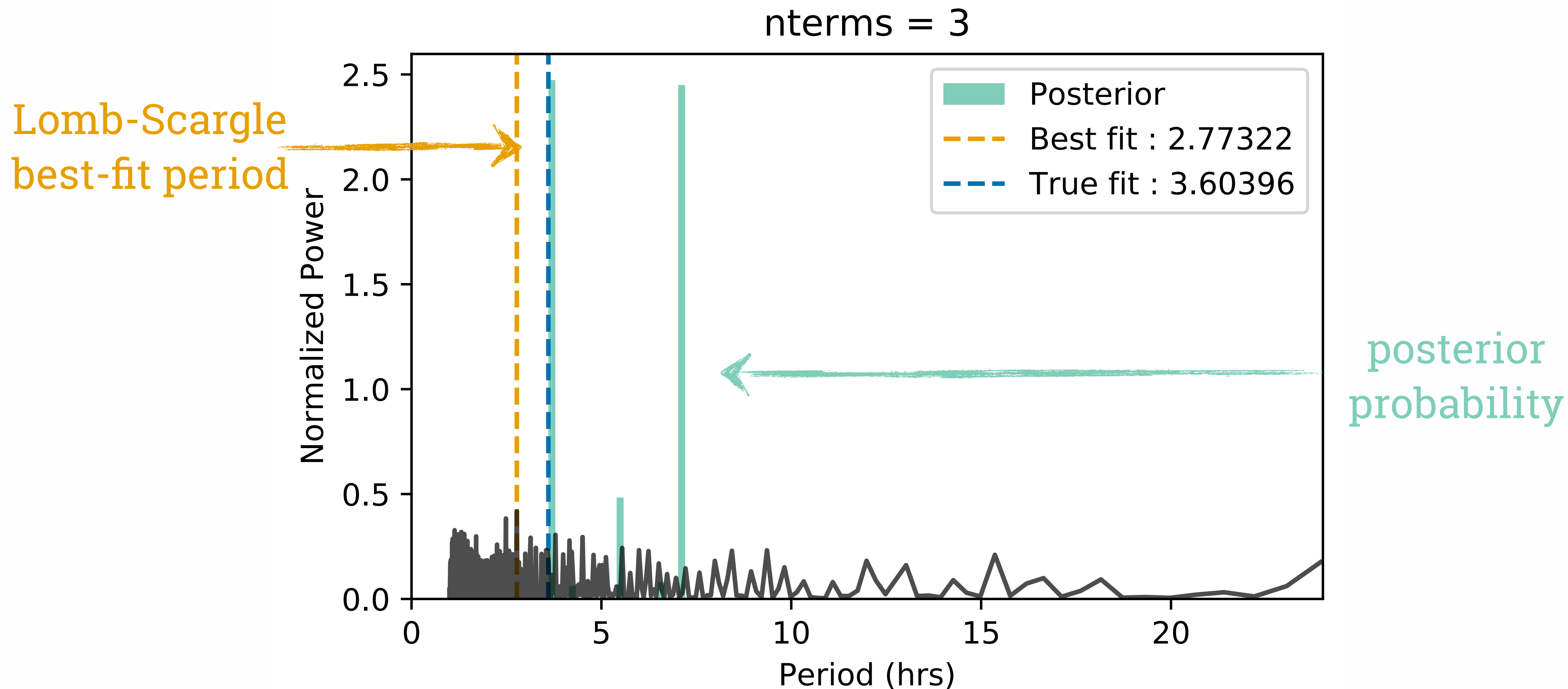nterms = 3

Posterior
Best fit : 2.77322
True fit : 3.60396

Lindberg, Huppenkothen et al (in prep)

# Gaussian Processes for Sparse Asteroid Light Curves



nterms = 3

Lomb-Scargle best-fit period

Posterior
Best fit : 2.77322
True fit : 3.60396

# Gaussian Processes for Sparse Asteroid Light Curves

# Training data sets are biased



PHOTOMETRIC LSST ASTRONOMICAL TIME-SERIES CLASSIFICATION CHALLENGE (PLASTICC)

Typical Supernova

Black Hole Activity

Superluminous Supernova
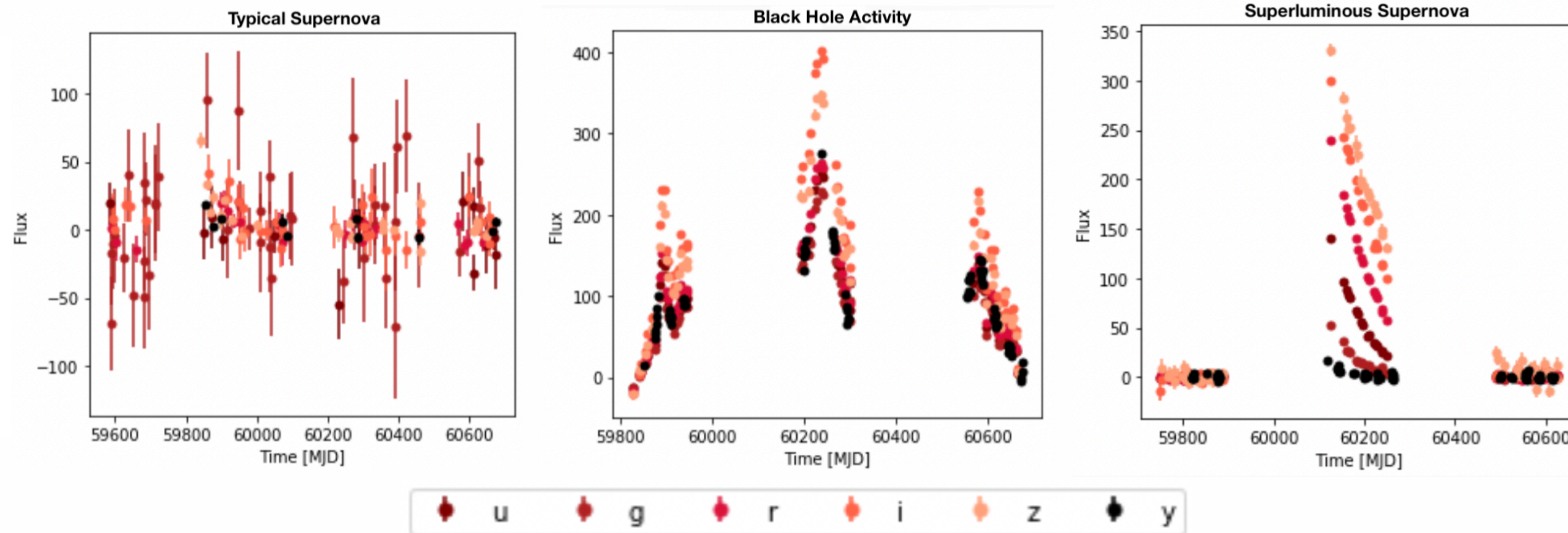
u  g  r  i  z  y

see also: The PLAsTiCC Team (2018)

# Training data sets are biased

- 8000 objects in training data set: bright, low-redshift objects



PHOTOMETRIC LSST
ASTRONOMICAL TIME-
SERIES CLASSIFICATION
CHALLENGE (PLASTICC)
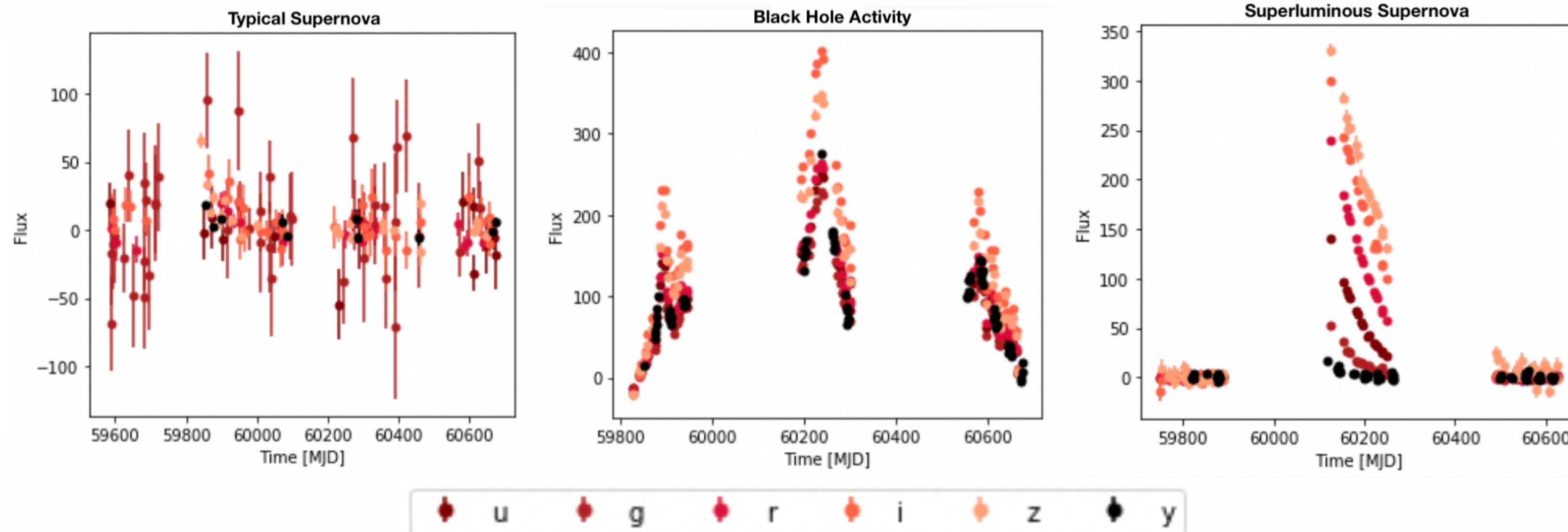
Figure Credit: Leah Fulmer

see also: The PLAsTiCC Team (2018)

# Training data sets are biased

- 8000 objects in training data set: bright, low-redshift objects

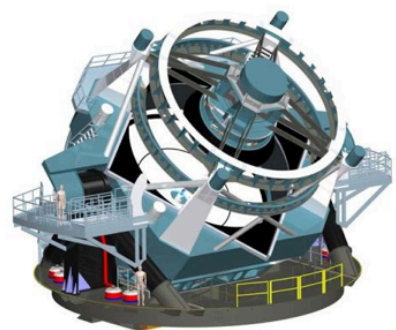- 3.5 million objects in test set: fainter, more distant objects

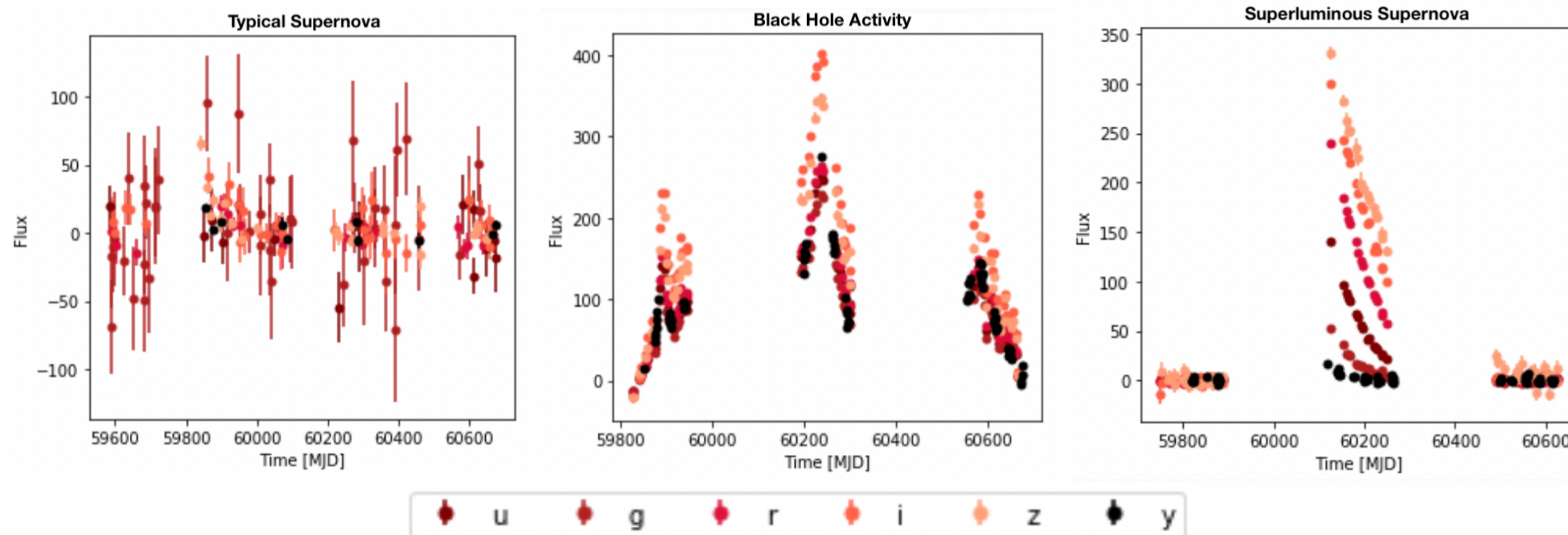**PHOTOMETRIC LSST ASTRONOMICAL TIME-SERIES CLASSIFICATION CHALLENGE (PLASTICC)**



Figure Credit: Leah Fulmer

see also: The PLAsTiCC Team (2018)

# Training data sets are biased

- 8000 objects in training data set: bright, low-redshift objects

- 3.5 million objects in test set: fainter, more distant objects

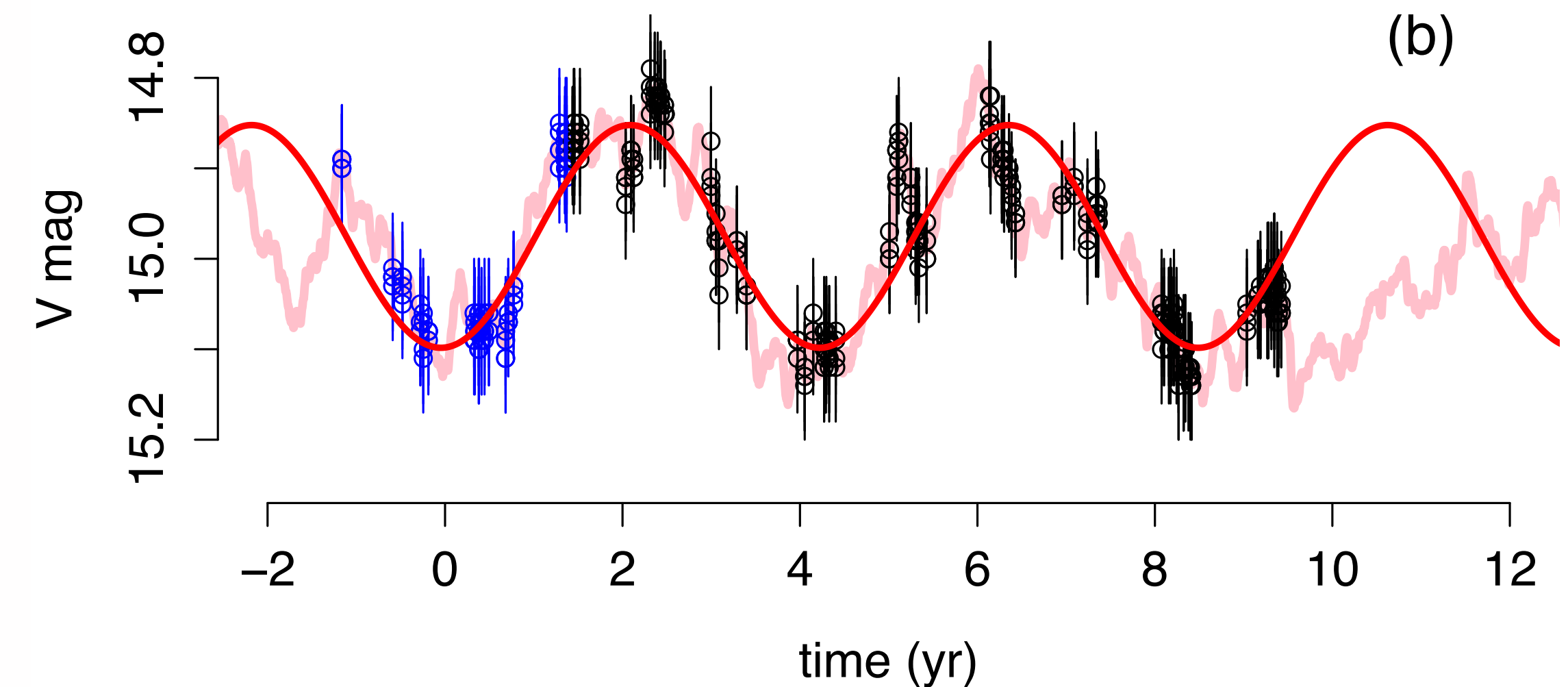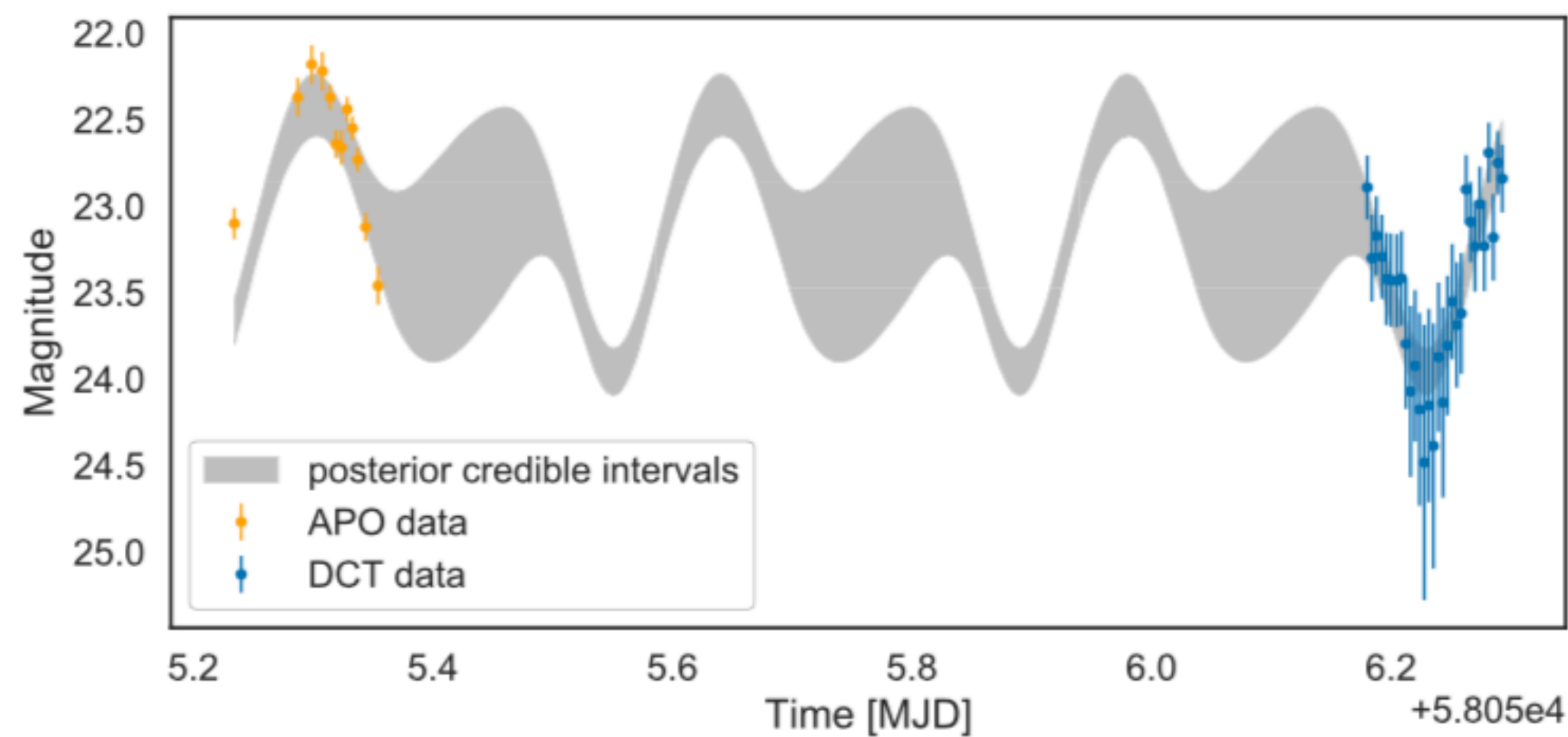- training data properties are non-representative of the test set

PHOTOMETRIC LSST
ASTRONOMICAL TIME-
SERIES CLASSIFICATION
CHALLENGE (PLASTICC)



Figure Credit: Leah Fulmer

see also: The PLAsTiCC Team (2018)

# (Some) Current Major Challenges in Time Domain Astronomy

- uneven sampling

- heteroscedasticity

- non-stationarity

- multi-wavelength data sets

- modeling multiple dimensions simultaneously (time, energy, polarimetry, …)

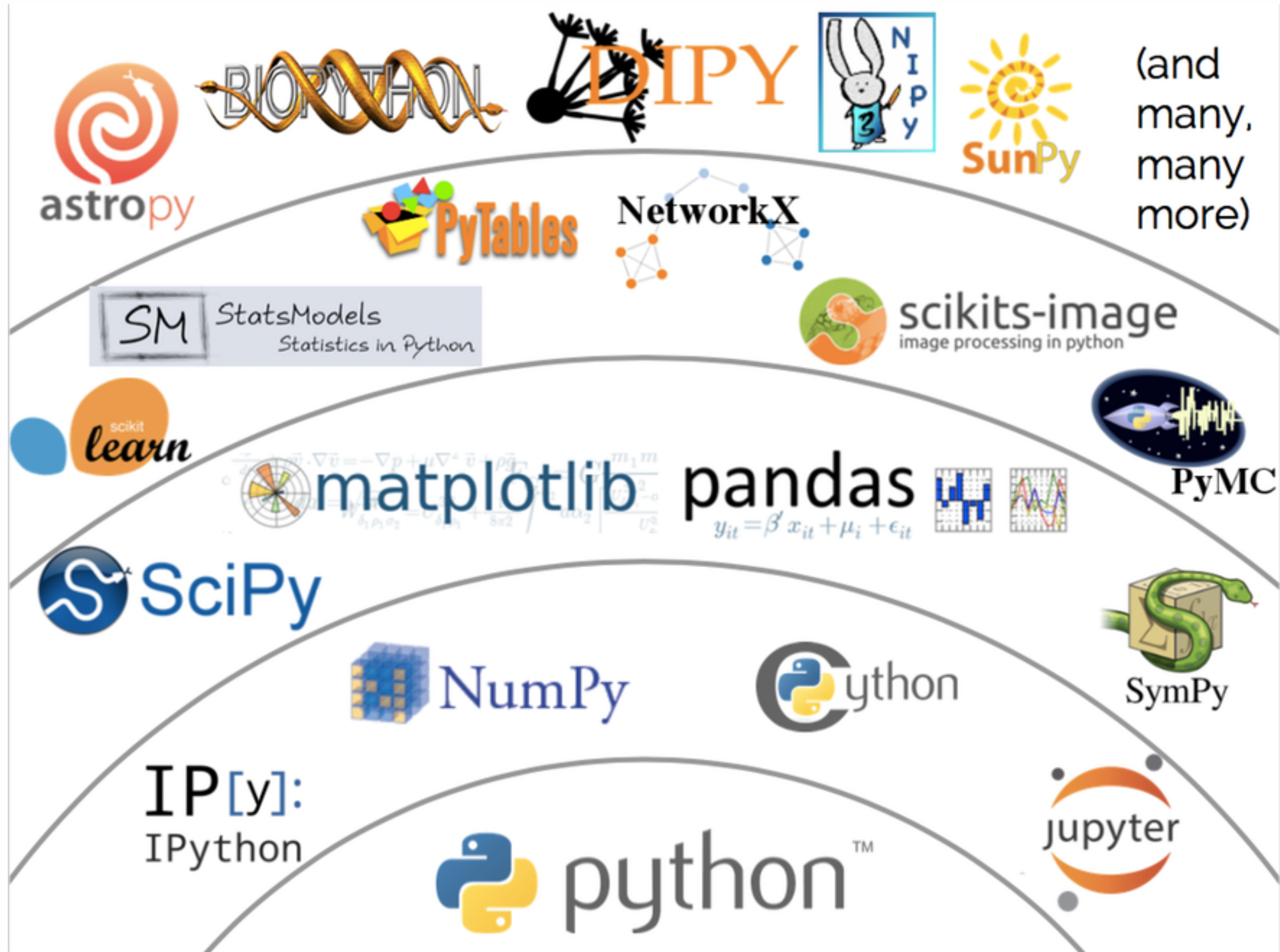- computational scaling

- classification of sources and time series

# Asteroid light curves are a simpler analogue to finding periods in AGN



DiRAC

(a)

V mag

time (yr)

(b)

V mag

time (yr)

(c)

V mag

Vaughan et al, incl Huppenkothen (2016)

posterior credible intervals
APO data
DCT data

Magnitude

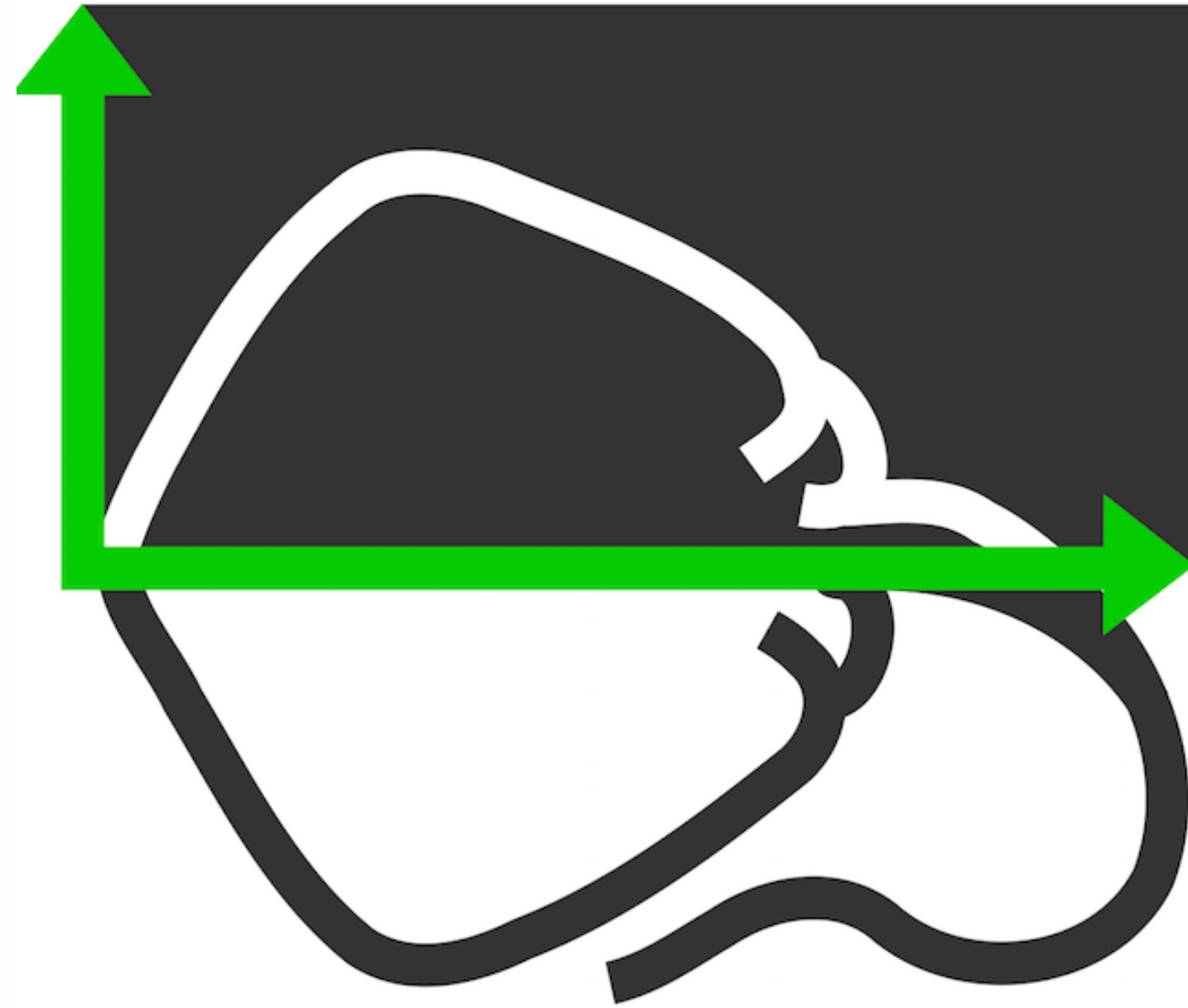Time [MJD]
+5.805e4

Bolin et al, incl Huppenkothen (2017)

**The good news:** many of the really hard problems have **simpler equivalents** in other areas of astronomy or even **other scientific domains**

**Solving data science challenges through communities of practice**

see also: http://msdse.org/reports/

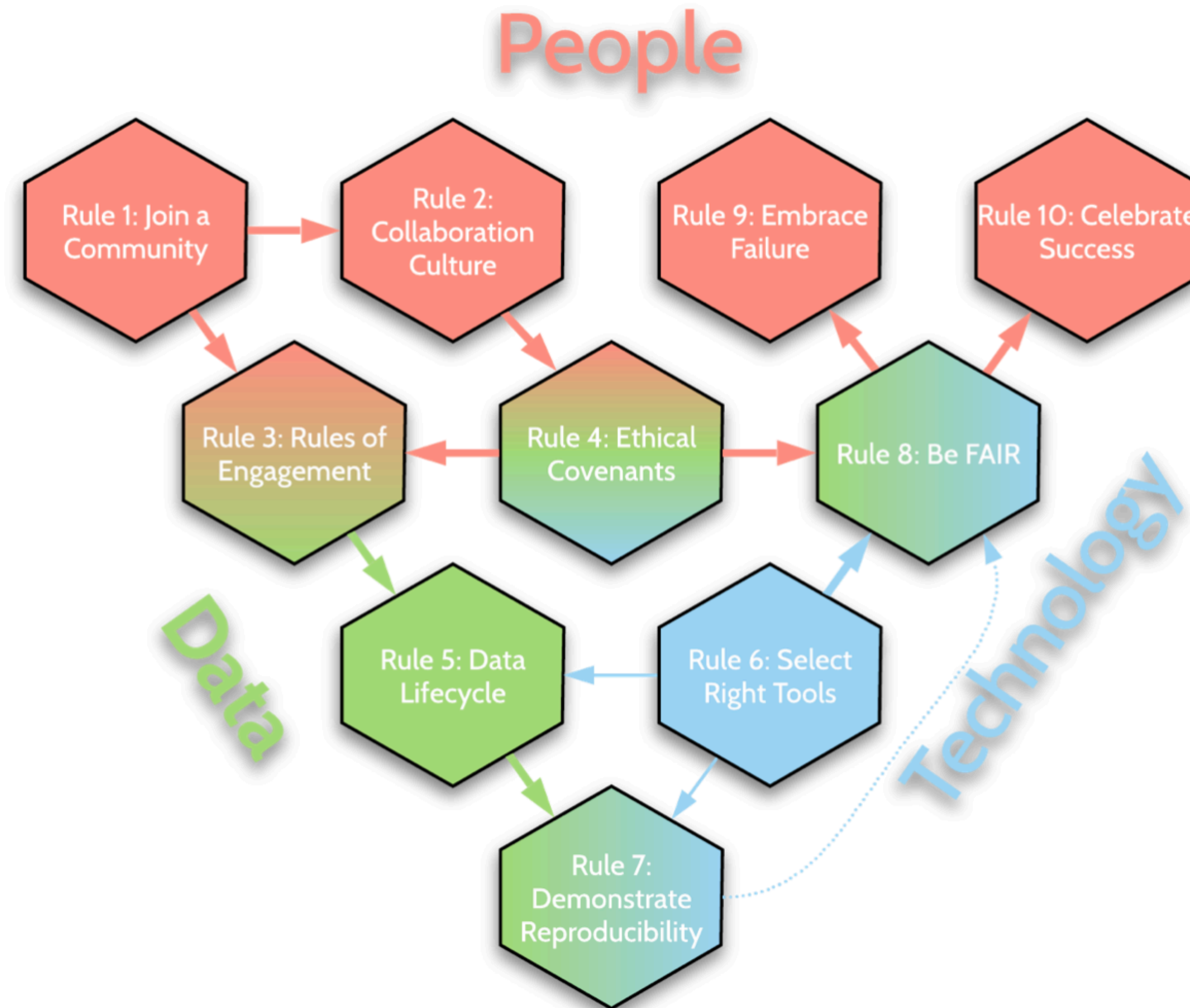**Credit: Jake VanderPlas**

**Stingray**
The Next Generation
Spectral-Timing Software

- 3 lead developers/maintainers (Huppenkothen, Bachetti, Stevens)

- ~10 contributors

- 6 completed Google Summer of Code Projects

- astropy-affiliated project

The largest **data science challenges** will be solved through **collaboration across fields**

# Ten Simple Rules for Researchers Engaging in Data Science and Domain Science Collaboration

# #Astro Hack Week

# #Astro Hack Week

- 5-day workshop
- ~50 participants
- tutorials and break-out sessions
- project work
- Lots of ☕ and 🍪
- participant-driven
- experimental

Hackweek Mission

Collaboration

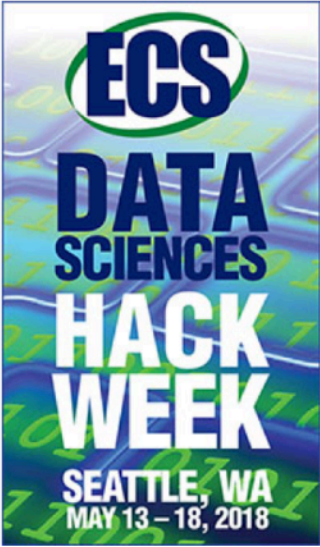Education

Building Software

Networking

Projects

Community

credit: Anthony Arendt

https://geohackweek.github.io

https://neurohackweek.github.io

https://oceanhackweek.github.io

https://www.electrochem.org/233/hack-week

https://waterhackweek.github.io

# Take-Away Lessons

build a **community** first

credit: eScience Institute

build a culture that empowers people to ask fundamental (and trivial) questions

Adapt concepts and ideas to your community's needs

credit: eScience Institute

Evaluate

# Conclusions

# Astronomical time domain data sets are
# complex, unevenly sampled, heteroscedastic, sparse and biased



## This often makes the application of standard tools difficult

# There are many data analysis challenges that are shared across scientific domains

**DiRAC**

## Musical structure analysis

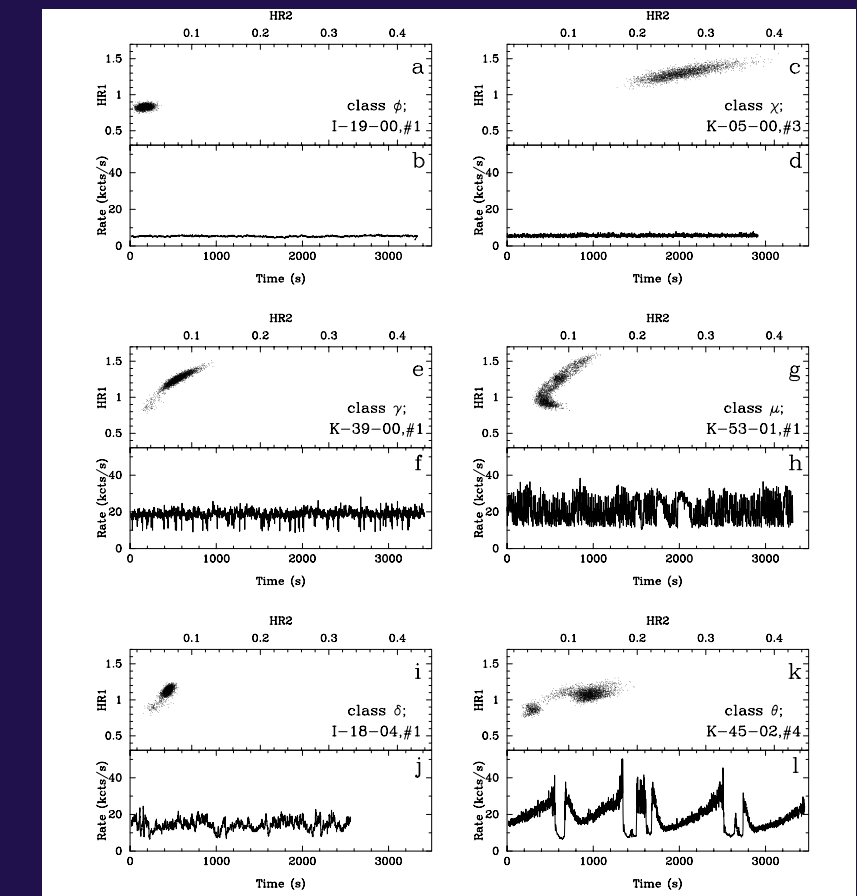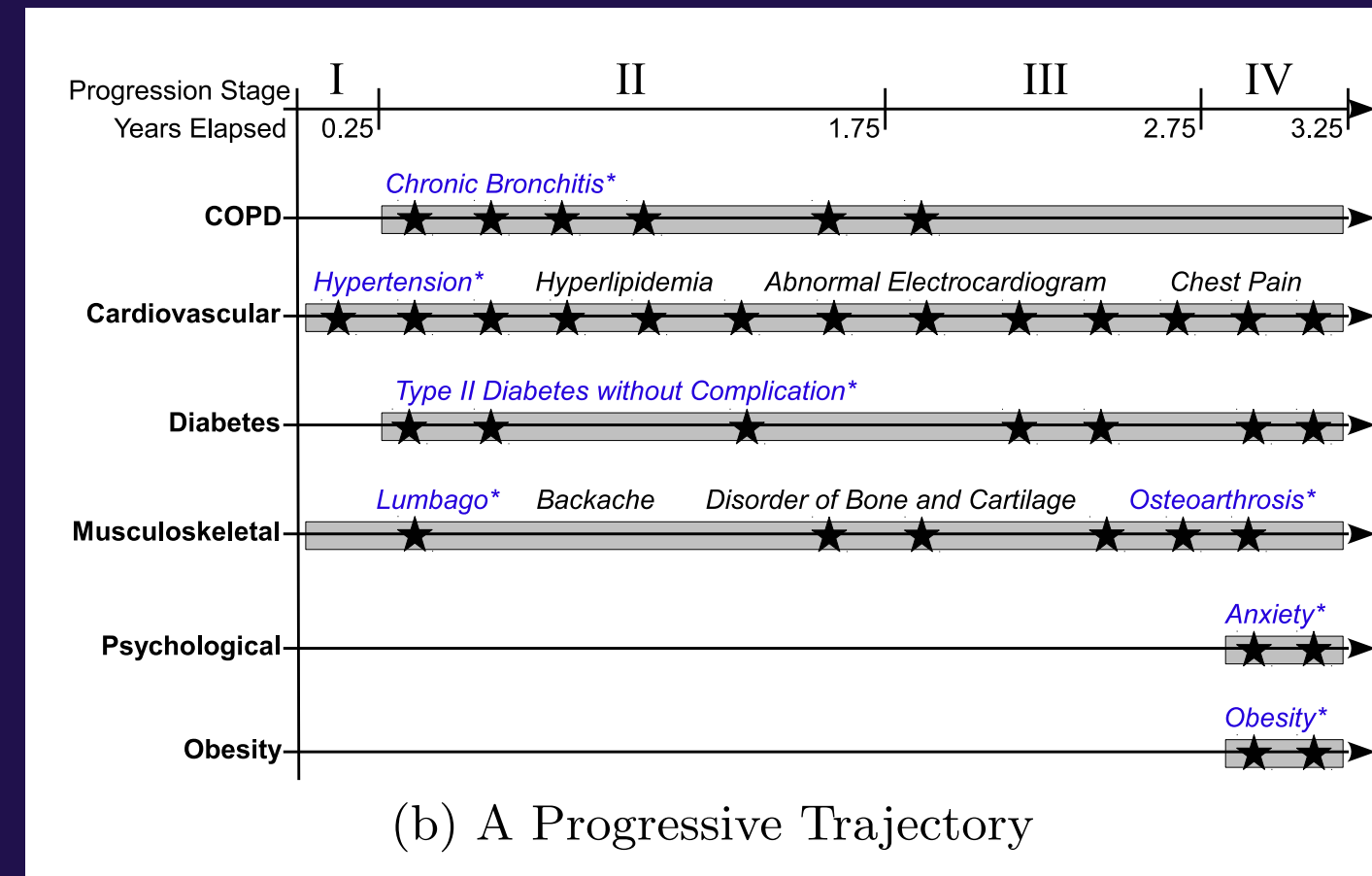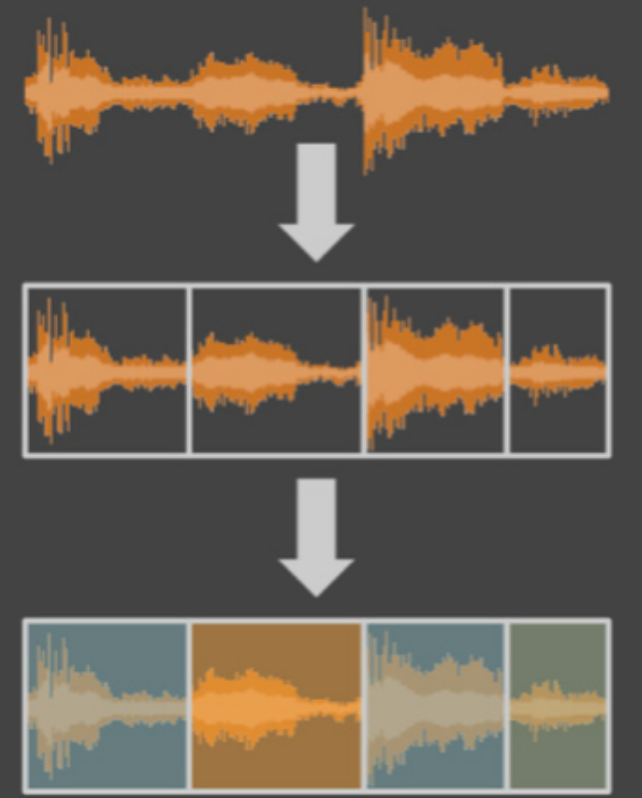1. Detect change-points
   *verse → chorus*

2. Label repeated sections
   ABAC

- ... also *representation* and *visualization*

Progression Stage | I | II | III | IV
Years Elapsed | 0.25 | 1.75 | 2.75 | 3.25

COPD — *Chronic Bronchitis*★

Cardiovascular — *Hypertension*★ Hyperlipidemia Abnormal Electrocardiogram Chest Pain

Diabetes — *Type II Diabetes without Complication*★

Musculoskeletal — *Lumbago*★ Backache Disorder of Bone and Cartilage *Osteoarthrosis*★

Psychological — *Anxiety*★

Obesity — *Obesity*★

(b) A Progressive Trajectory

Data science provides shared venues and a common language to solve these problems across domains

One of the **major challenges** of interdisciplinary research is the **language barrier** between fields



We need people who can translate across disciplinary jargon

**Interdisciplinary and data science research requires community building first and foremost**

**community building ≠ putting people in the same room**

# What can we learn from this community? What mistakes should we avoid making?