

# Astrostatistics

Lecture 02: Monday, 21 January 2019

## **Recommended Reading:**

Feigelson & Babu: Chapters 1-4

Ivezic: Chapters 1, 3-5

C. Schafer article:

“A Framework for Statistical Inference in Astrophysics”

Intro to Statistics in Astronomy

Review of Probability & Statistics Foundations

Classical & Bayesian Statistical Inference

**kmandel@statslab.cam.ac.uk**

<https://github.com/CambridgeAstroStat/PartIII-Astrostatistics-2019>

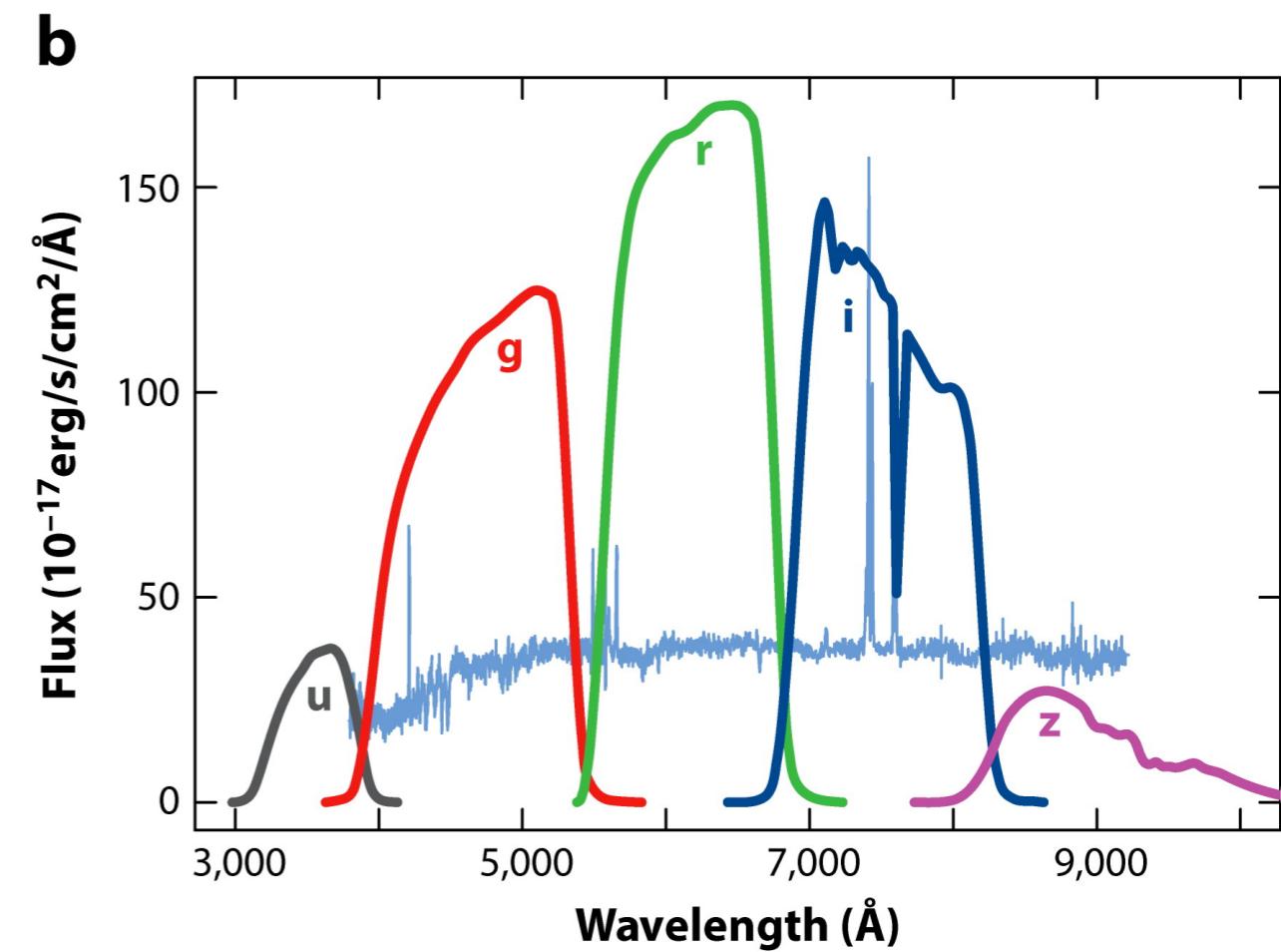
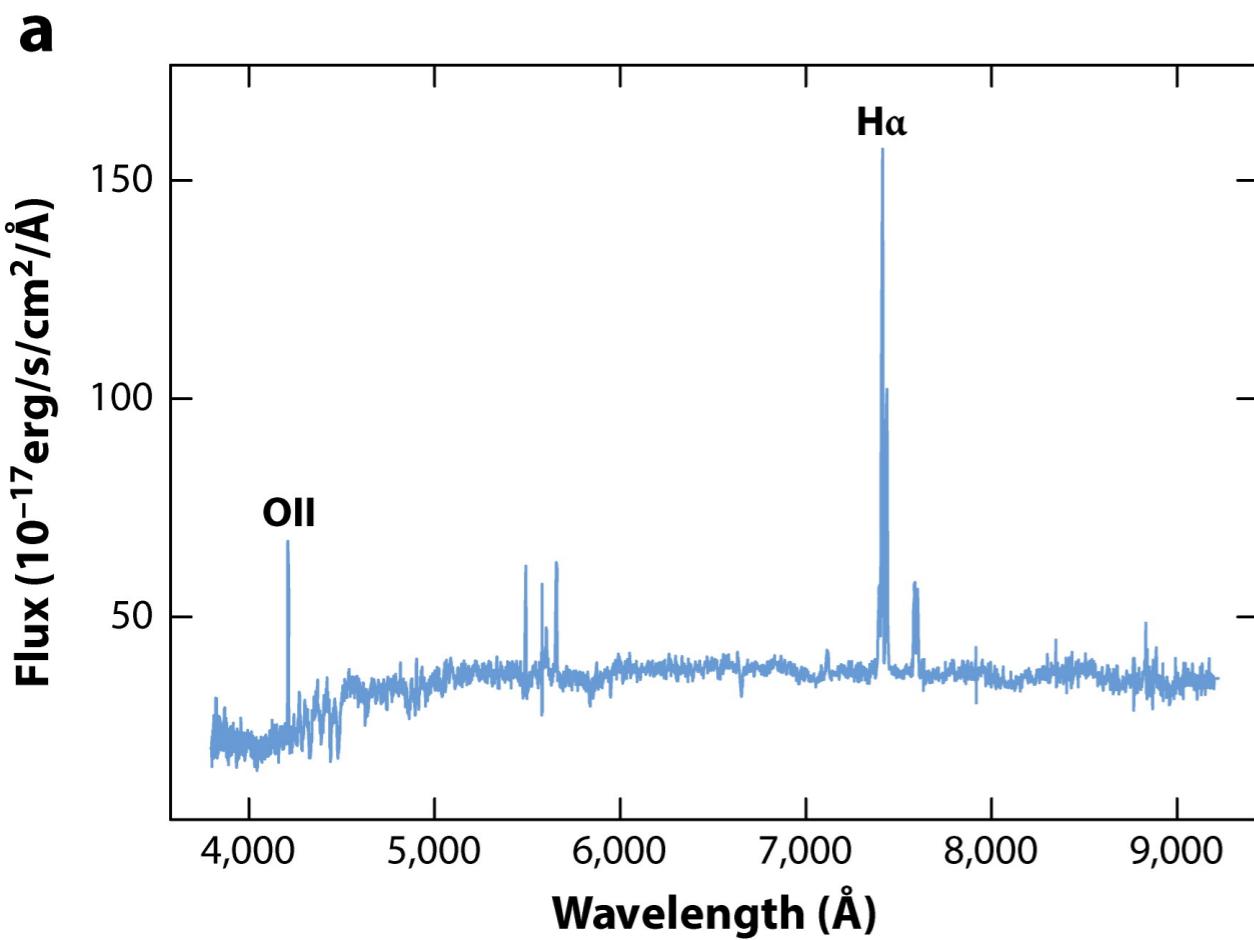
# Today

- Introduction to Astronomical Data Types (for statisticians)
  - See also Schafer article “A Framework for Statistical Inference in Astrophysics”
- Motivational Case Studies:
  - Bayesian Inference of the Milky Way Galactic Mass
  - Gravitationally Lensed Time Delay Estimation
  - **Gaussian Processes for Spectral Time Series (Radial Velocity Analysis)**
  - **Hierarchical Bayes for Supernova Cosmology**

# What astronomers measure

- Astrometry (angular position on sky, e.g. Gaia)
- Photometry (how bright is it?)
  - Flux = photons (or energy) per second per meter<sup>2</sup>
  - (apparent) Magnitude =  $-2.5 \times \log_{10} [\text{Flux}] + \text{const}$ 
    - Absolute Magnitude =  $-2.5 \times \log_{10} [\text{Luminosity}] + \text{const}$   
= apparent magnitude at fixed distance of 10pc
- Spectroscopy (brightness versus wavelength)
- Time Series (light curves): Transients & Variables (e.g. stars, quasars, supernovae, exoplanets), Moving objects (e.g. asteroids)
- Spatial Variation (images, maps) - clustering, spatial correlation functions
- Combinations of the above, e.g.
  - Astrometry vs. Time = Proper Motion (e.g. stars, satellite galaxies)
  - Spectroscopic Time Series (Radial Velocity studies of stars)

# Spectroscopy and Photometry



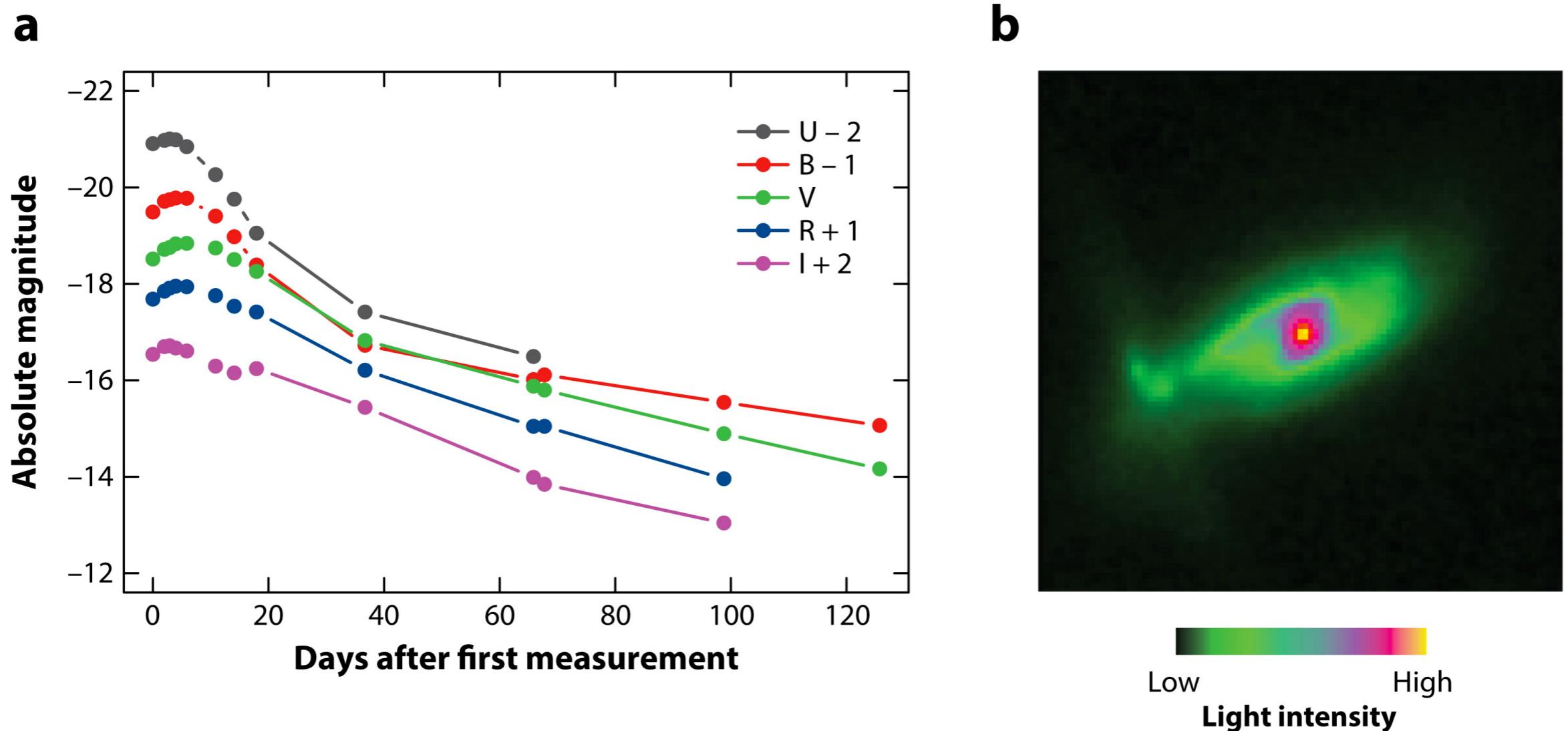
Schafer CM. 2015.

Annu. Rev. Stat. Appl. 2:141–62

Galaxy Spectrum

Galaxy Photometry  
(Brightness: Flux / Magnitude)

# Temporal & Spatial Variation



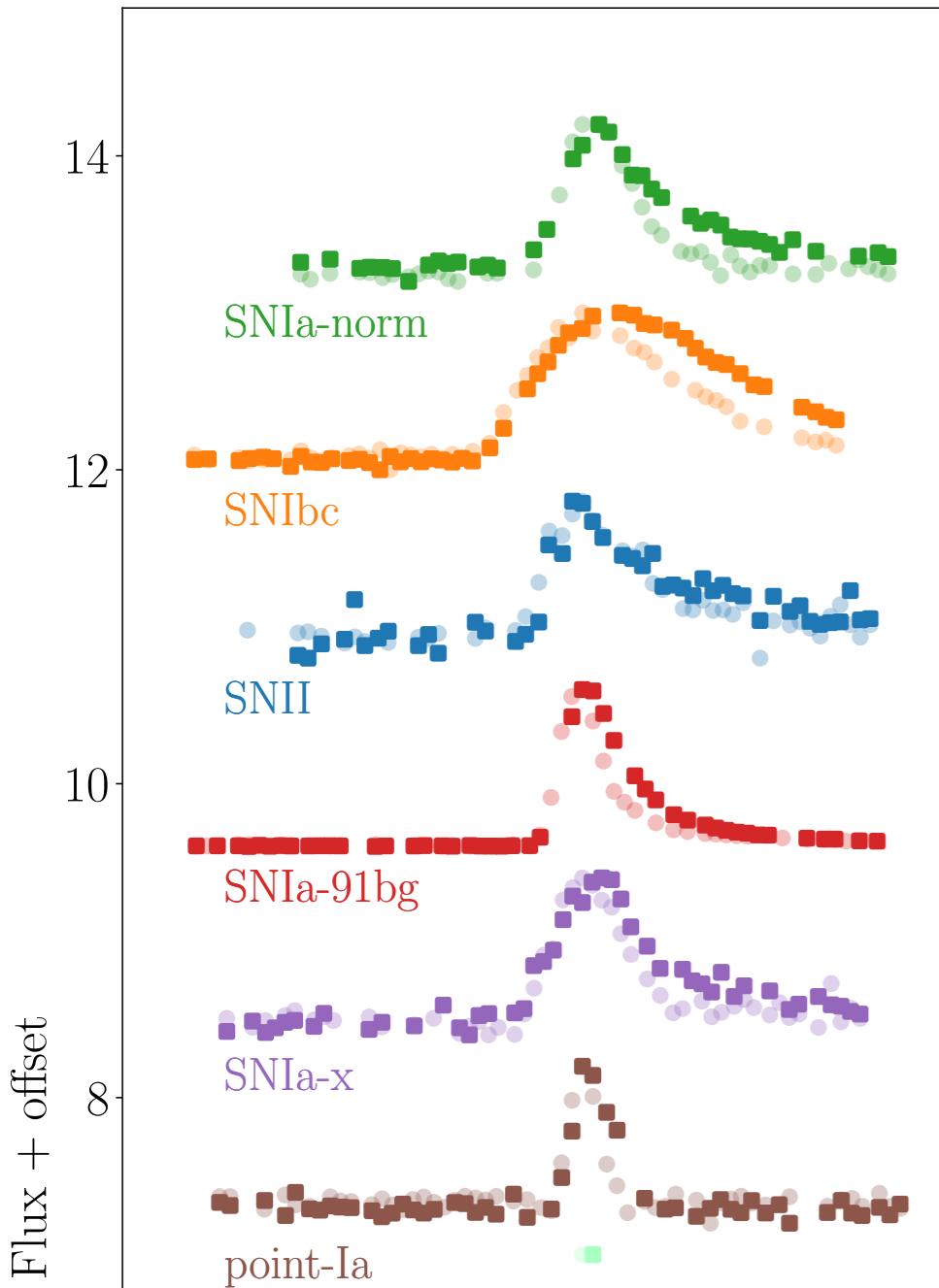
Schafer CM. 2015.

Annu. Rev. Stat. Appl. 2:141–62

Time Series (Light Curve)  
Supernova

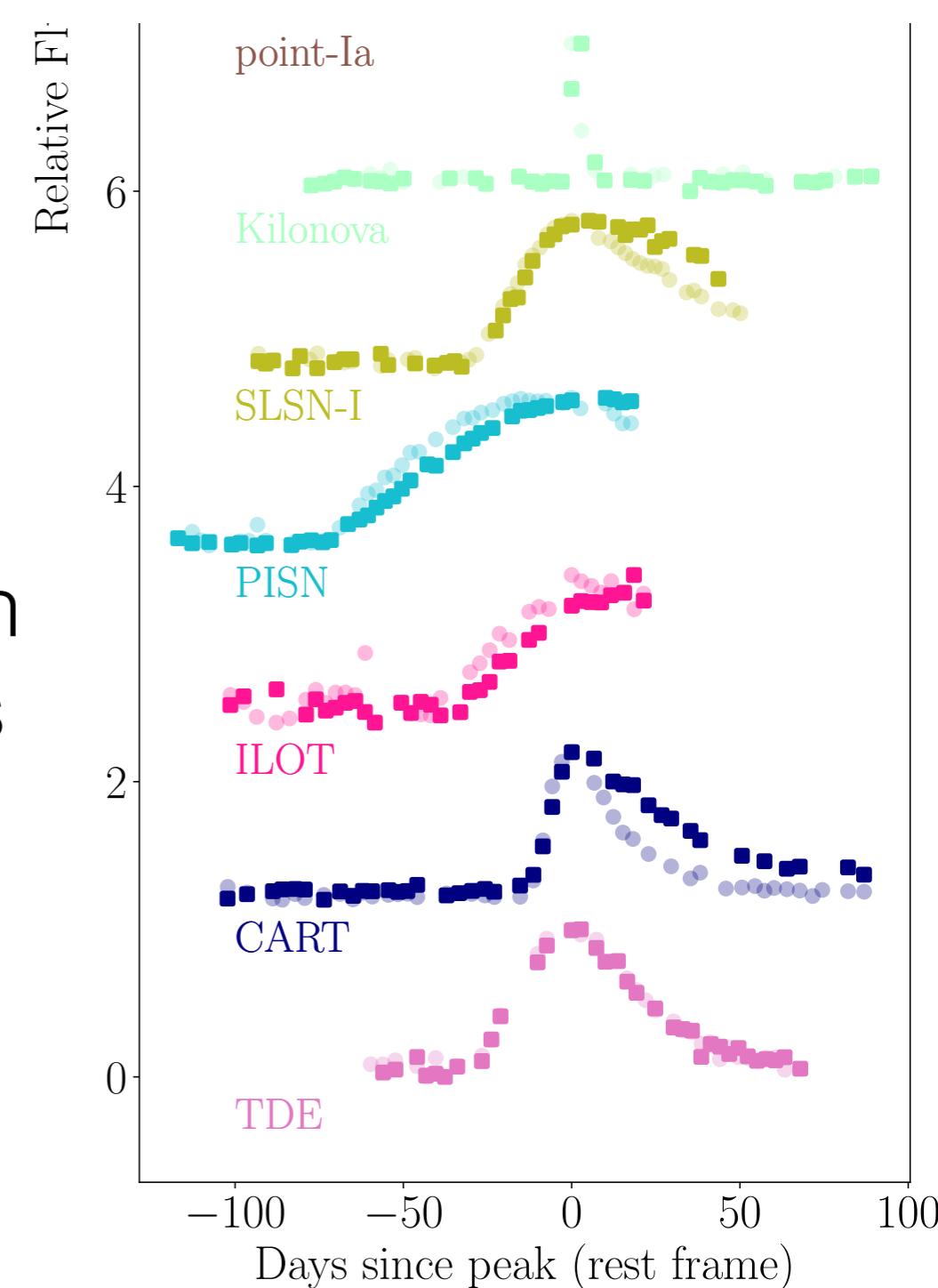
Galaxy Image  
(Intensity Map)

# Time-Series (Time Domain Astronomy)



## Goals:

- Classification
- Astrophysics
- Cosmology



**Figure 2.** The light curves of one example transient from each of the 12 transient classes is plotted with an offset. We have only plotted transients with a high signal-to-noise and with a low simulated host redshift ( $z < 0.2$ ) to facilitate comparison of light curve shape between the classes. The opaque square markers plots the  $r$  band light curves of each transient, while the transparent circle markers are the  $g$  band light curves of each transient.

Muthukrishna et al. 2019  
Deep Learning for  
Transient Classification

# Photometric LSST Astronomical Time Series Classification Challenge (PLAsTiCC):

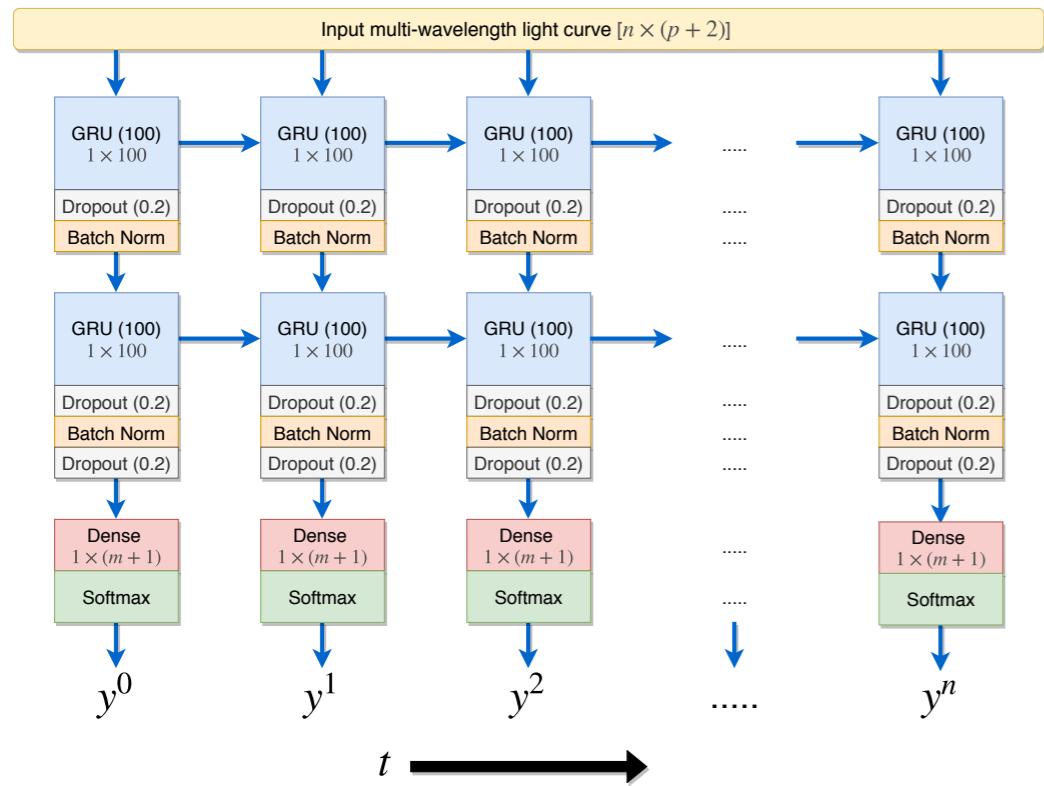
<https://www.kaggle.com/c/PLAsTiCC-2018>

The image shows the landing page for the PLAsTiCC Astronomical Classification competition on Kaggle. The background features a dark blue abstract geometric pattern. At the top left, there is a trophy icon and the text "Featured Prediction Competition". In the center, the competition title "PLAsTiCC Astronomical Classification" is displayed in large white font, with the subtitle "Can you help make sense of the Universe?" below it. To the right, a large "\$25,000 Prize Money" is prominently shown. At the bottom left, there is a small LSST logo and the text "LSST Project · 1,094 teams · a month ago". Below the main title, there is a navigation bar with links: Overview (which is underlined in blue), Data, Kernels, Discussion, Leaderboard, and Rules.

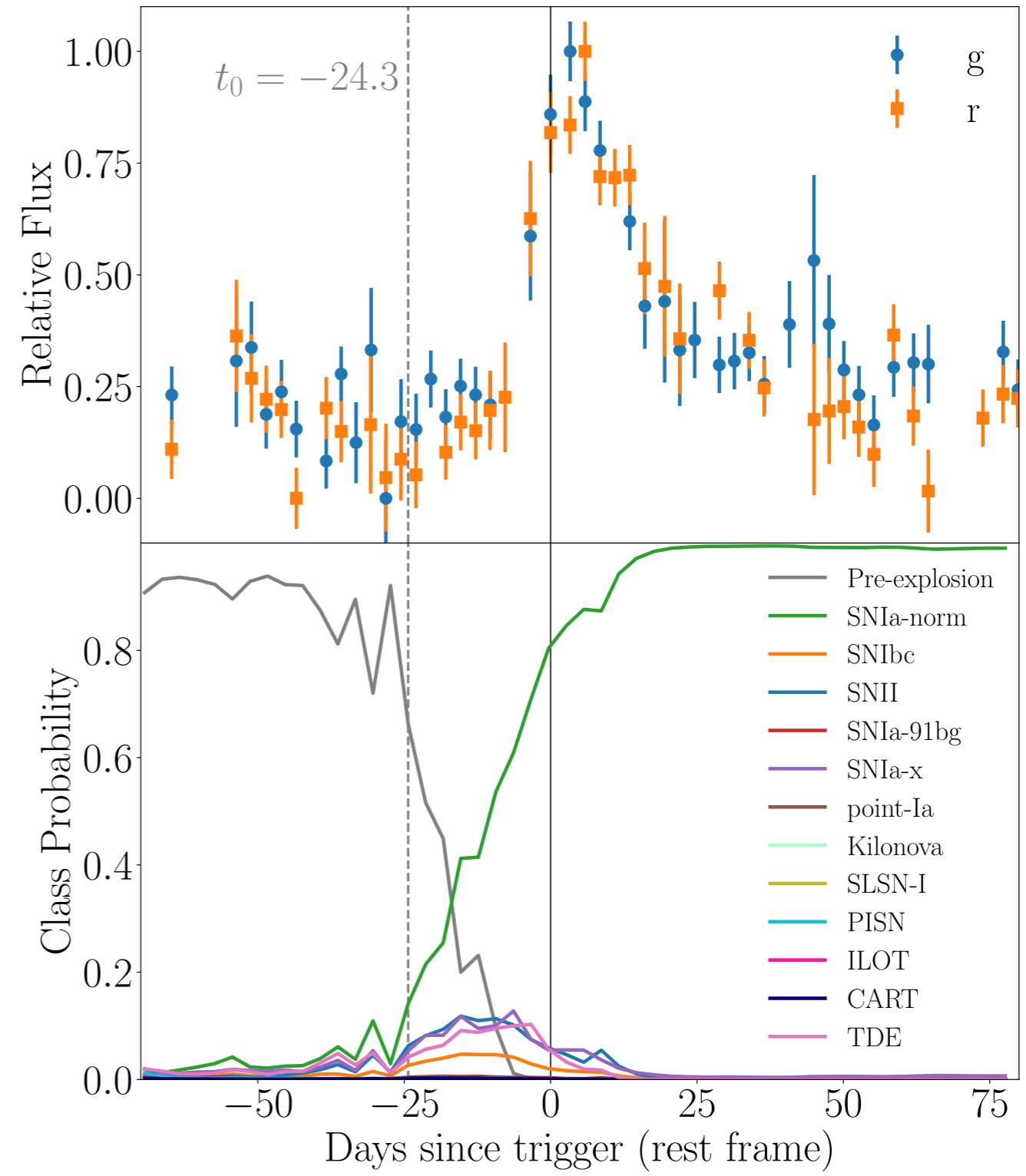
This image shows a detailed view of the competition's "Overview" section. It includes a "Description" box with the text "Help some of the world's leading astronomers grasp the deepest properties of the universe." and an "Evaluation" box. To the right, there is a large, blurry image showing a field of stars or astronomical data. The overall layout is clean and professional, typical of a Kaggle competition page.

# Transient Time Series Classification

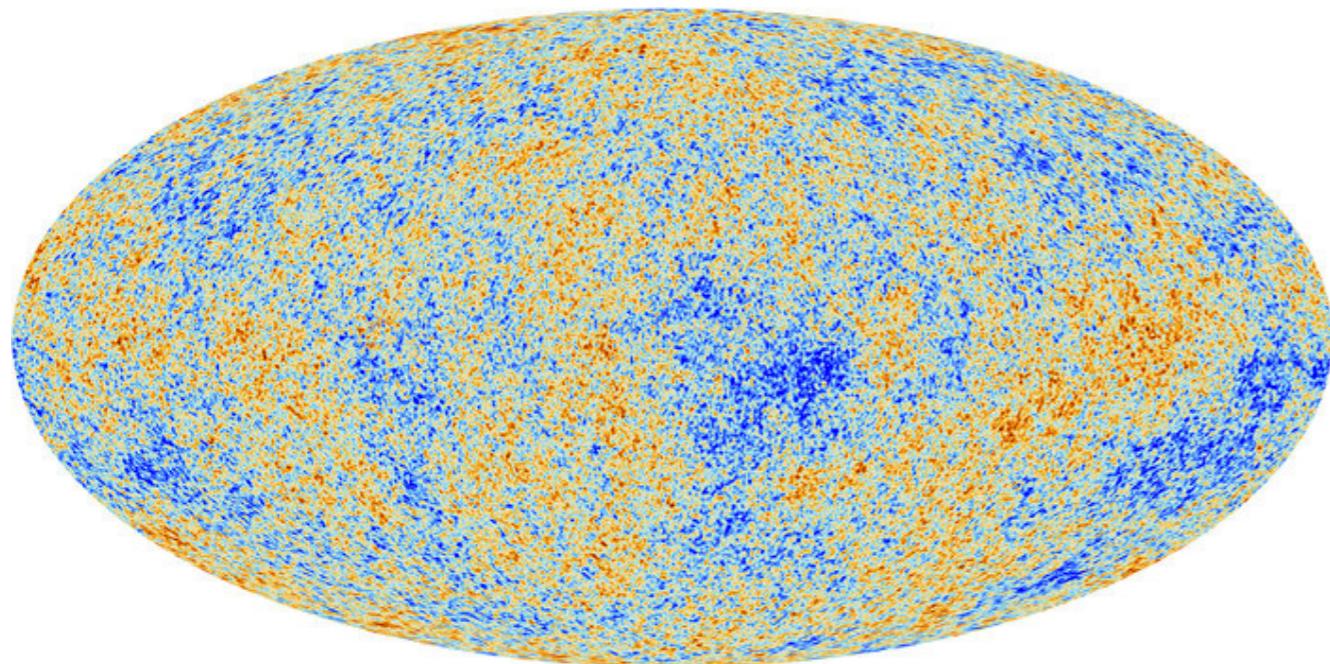
## Deep Recurrent Neural Network



Muthukrishna et al. 2019  
Deep Learning for  
Transient Classification

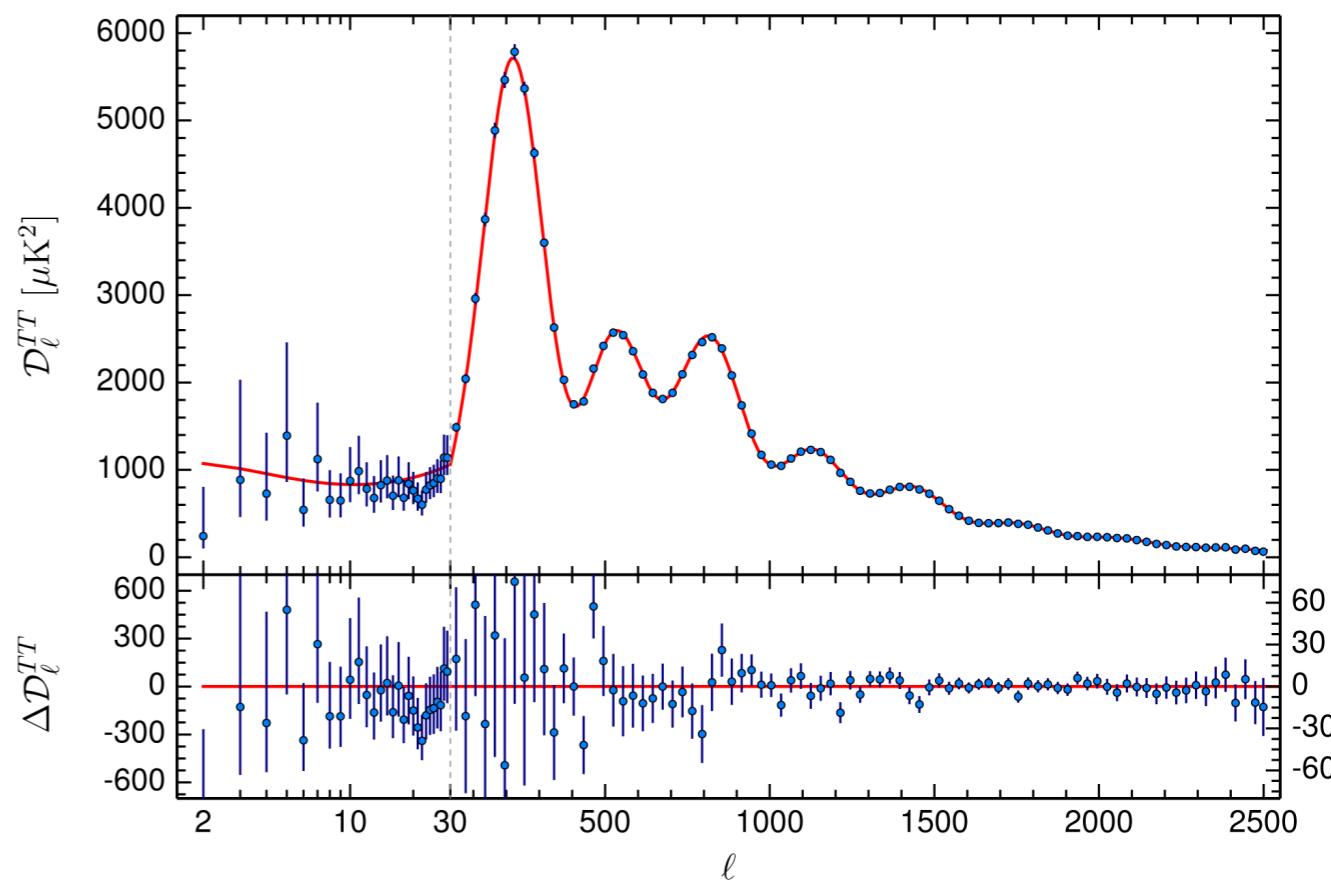


# Spatial Variation of Intensity



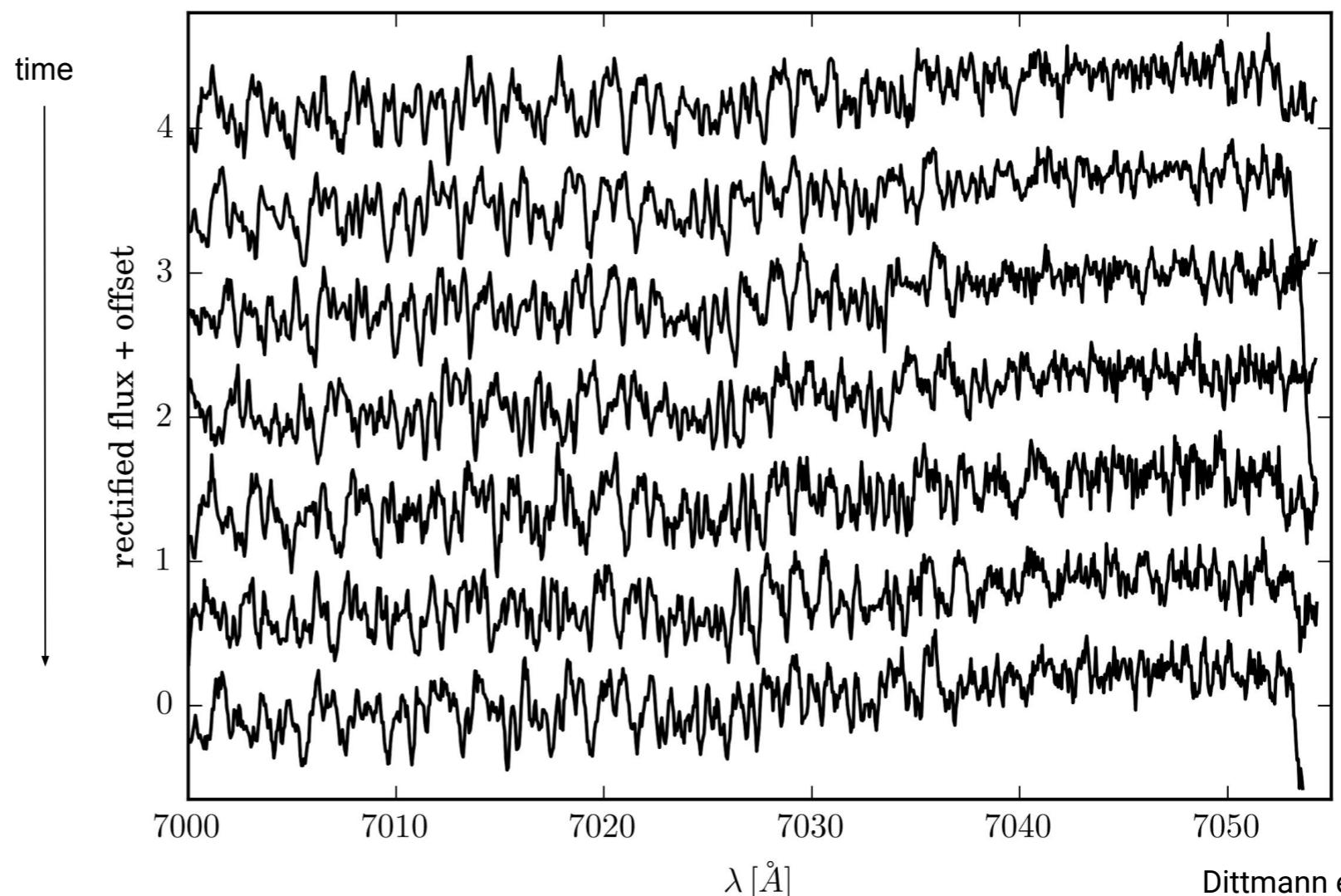
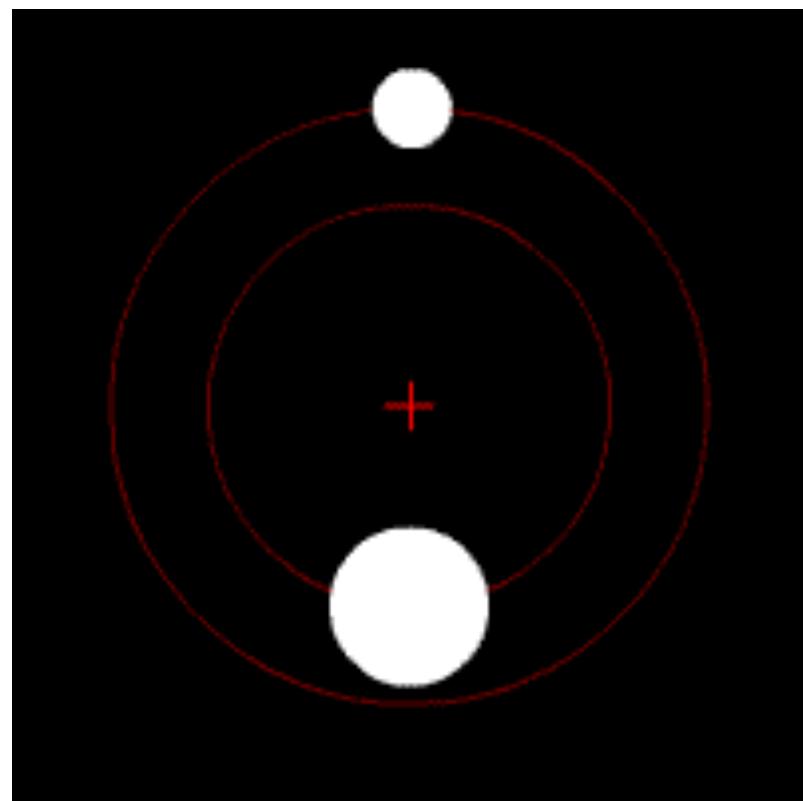
Cosmic Microwave  
Background (Planck)  
~ Gaussian Random Field  
(mean = 2.7 K,  
std dev  $\sim 10^{-5}$ )

Power Spectrum  
(~Fourier Transform of  
Correlation Function)

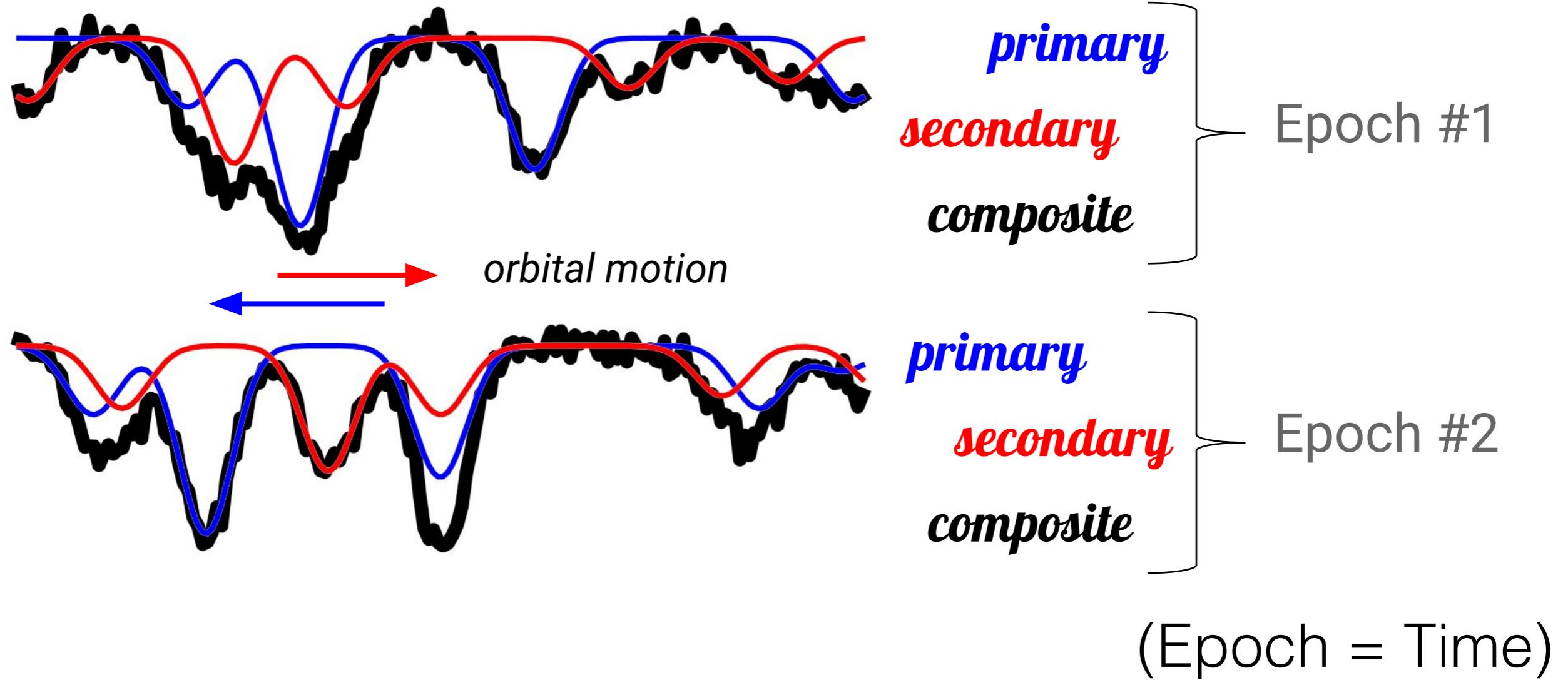


*Astrostatistics Case Studies:*  
Disentangling Time Series Spectra with Gaussian  
Processes: Applications to Radial Velocity Analysis  
(Czekala et al. 2017, ApJ, 840, 49. arXiv:1702.05652)

**Raw Observations of the LP661-13 M4 Binary**



# Spectroscopic Binary Stars



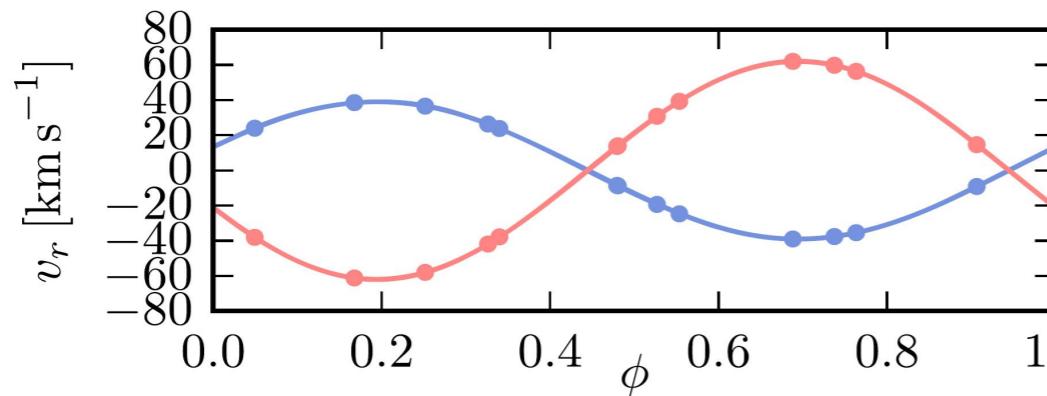
We only observe the “noisy” sum of two (latent) spectra.  
Latent (underlying) spectra are unknown functions  
Observed spectrum = Measured Data

# *Forward Model = Generates Data*

## Problem setup

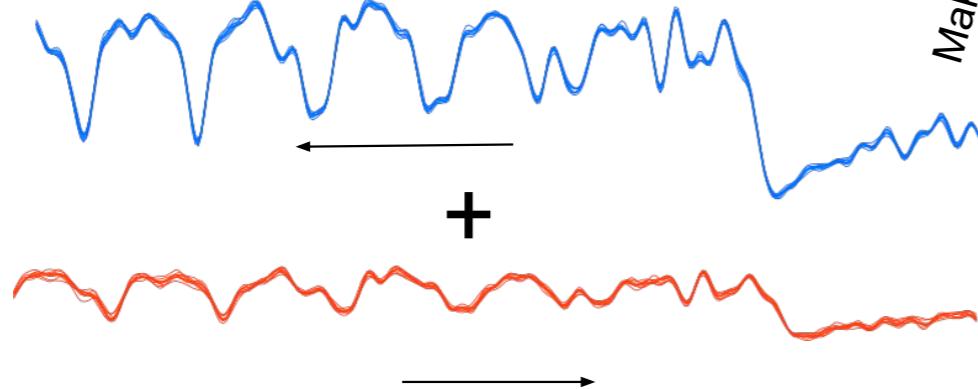
**Orbit:** period,  
eccentricity,  
phase, etc.

?



Model  
spectra

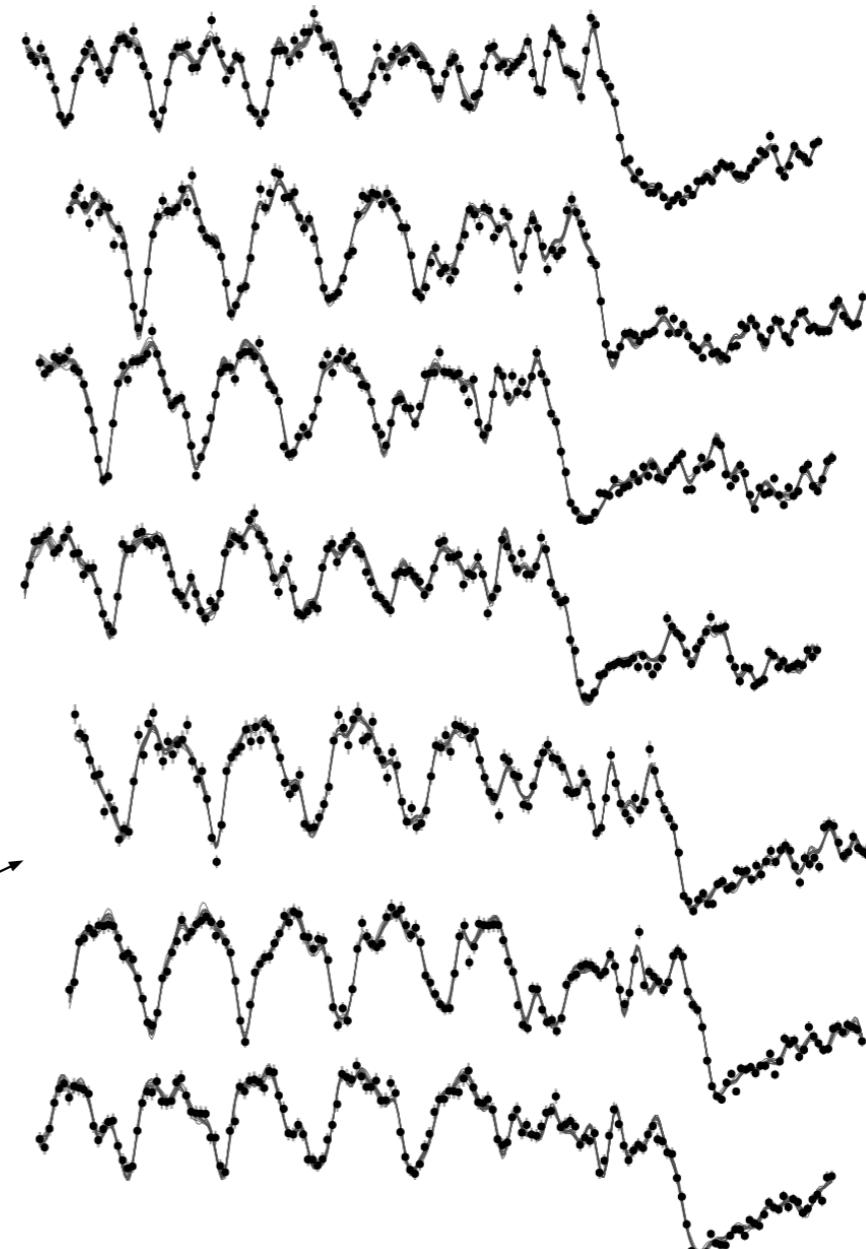
?



Velocity shifts

Make composite spectra

Data spectra

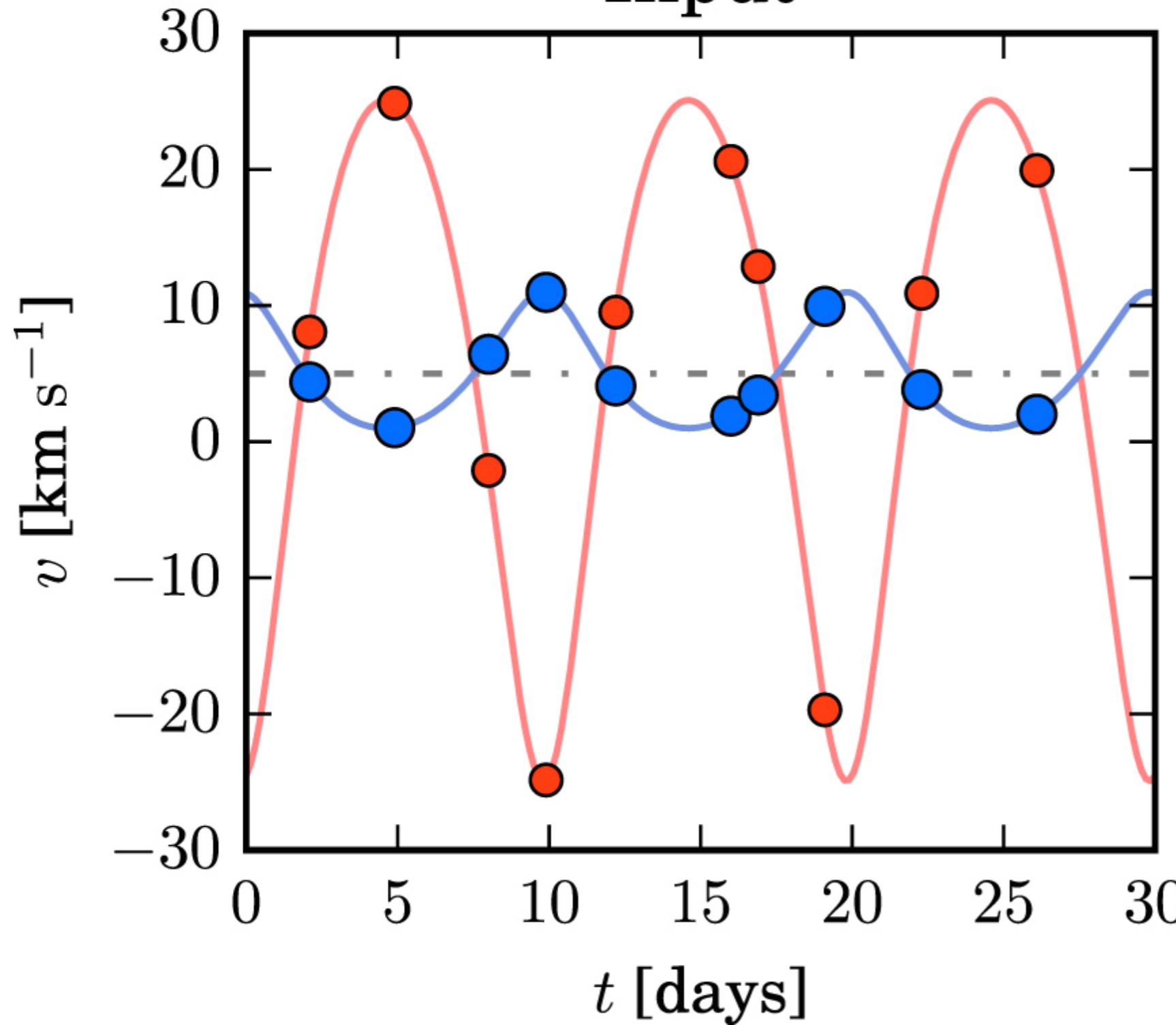


<https://www.youtube.com/watch?v=kHjN42ft6aU>

Goal: Go Backwards and Infer the Component Spectra & Orbital Parameters from noisy, observed (composite) spectra time series

# Orbital Parametric Model

Input



- Seven Parameters:
- Mass Ratio
  - Velocity Amplitude
  - eccentricity
  - Arg of Periastron
  - Epoch of Periastron
  - Orbital Period
  - Systemic Velocity

# Nonparametric Bayes

## Gaussian processes

We will model the latent stellar spectrum  $f_\lambda$  as a Gaussian process

$$f_\lambda \sim \text{GP}(\mu(\lambda), k(\lambda, \lambda'))$$

A function is said to have a Gaussian process if for any collection of inputs the random vector  $\mathbf{f}$  has a multivariate Gaussian distribution with mean  $\mathbf{\mu}$  and covariance matrix given by  $k$  evaluated over ***lambda***

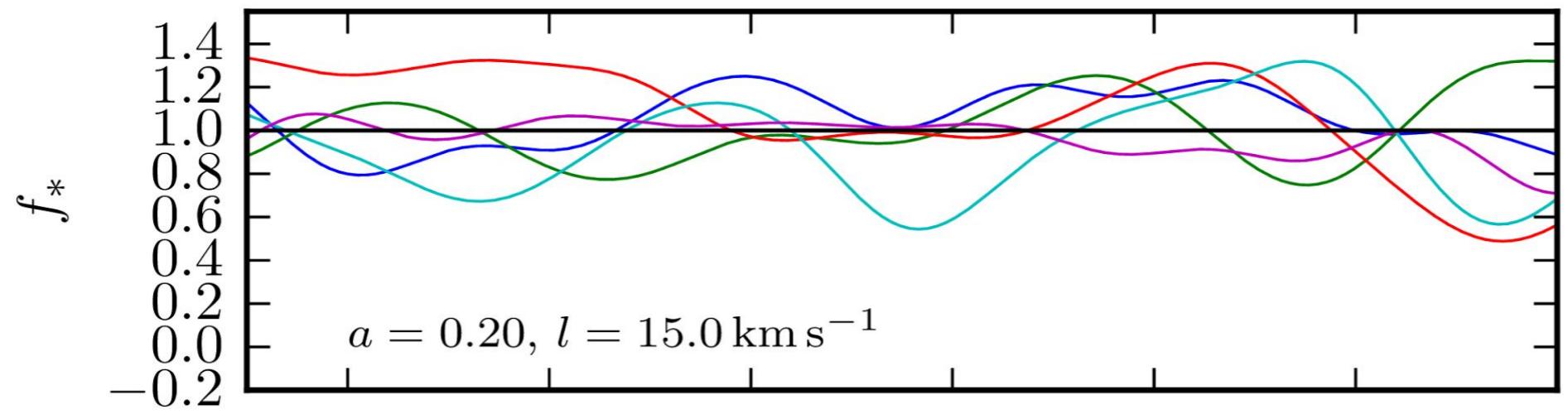
For a covariance kernel, we will use the commonly used squared exponential kernel, which relates pixels in the spectrum based upon their distance in log-wavelength ( $\propto$  velocity)

$$k_{ij}(r_{ij} | a, l) = a^2 \exp\left(-\frac{r_{ij}^2}{2l^2}\right)$$

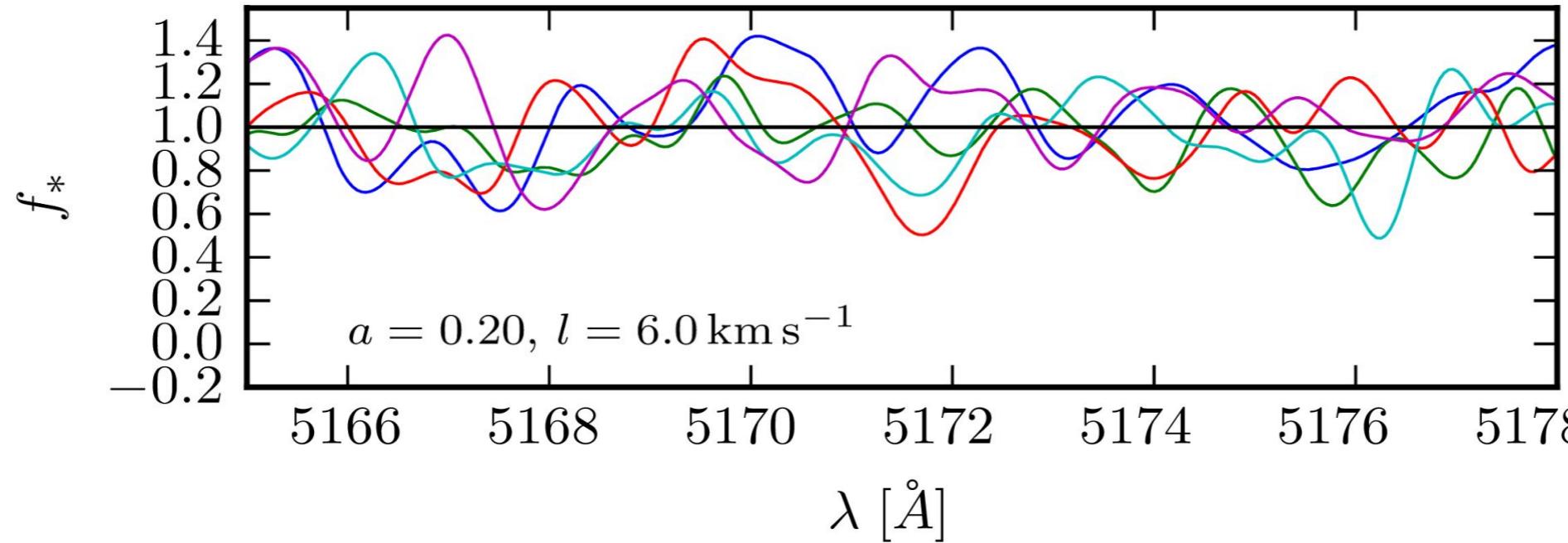
Gaussian Process = a prior on functions (latent spectra)

## **Gaussian Process model for a single, stationary star**

(Zoomed) draws from the prior



$l$



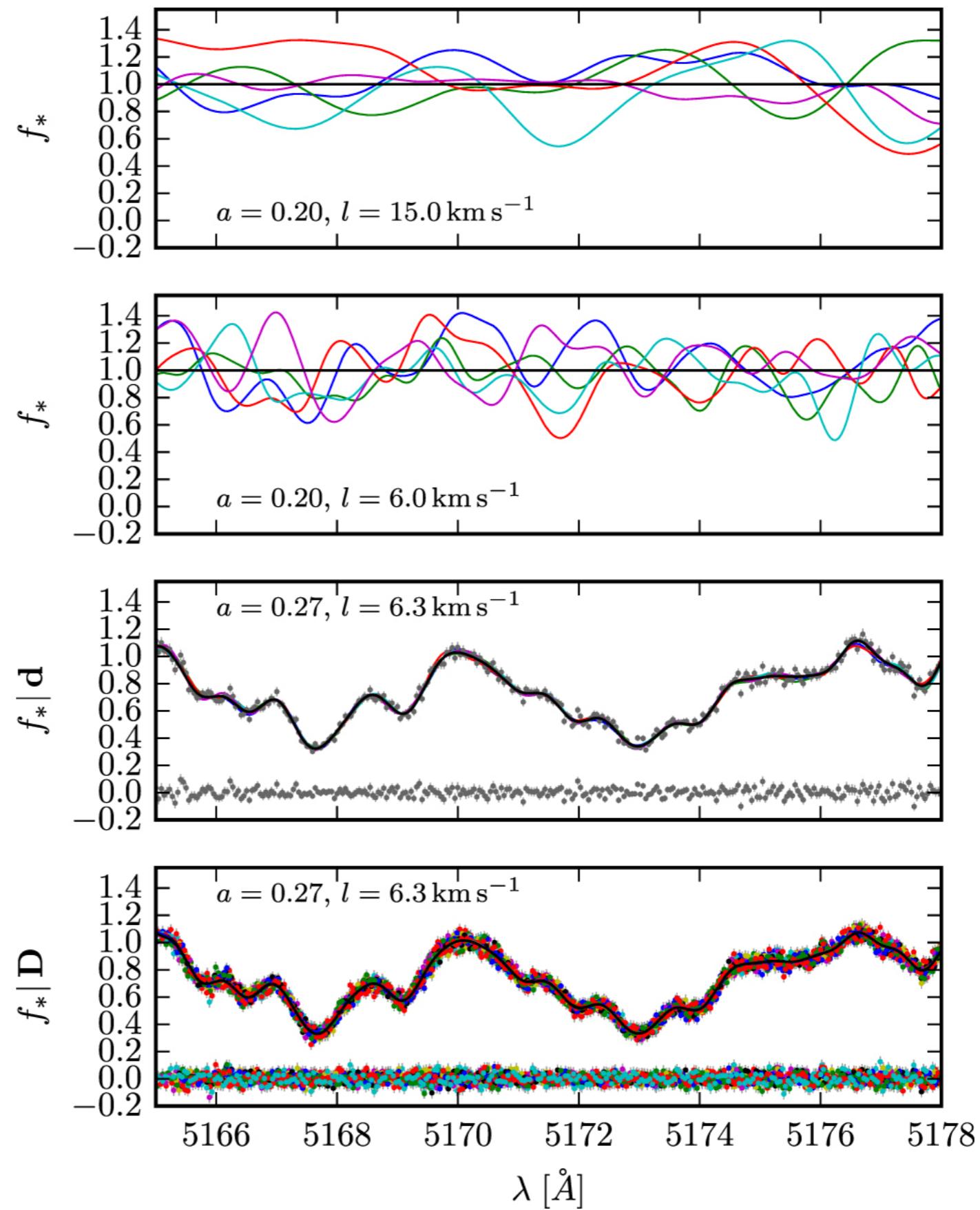
$l$

Inference = Which function is most consistent with the data?

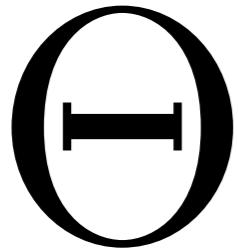
# Gaussian Process: Priors & Posteriors

GP prior  
(long/short length scales)

GP Posterior  
(conditioned on data spectrum  $\mathbf{d}$ )  
Inference of latent spectrum



# Known Unknowns



7-dim Orbital Parameters = Period, Phase, eccentricity, Velocity Amplitude

$f(\lambda), g(\lambda)$

( $\infty$ -dim) Latent Functions = the unobserved component spectra of the primary (f) and secondary (g) stars

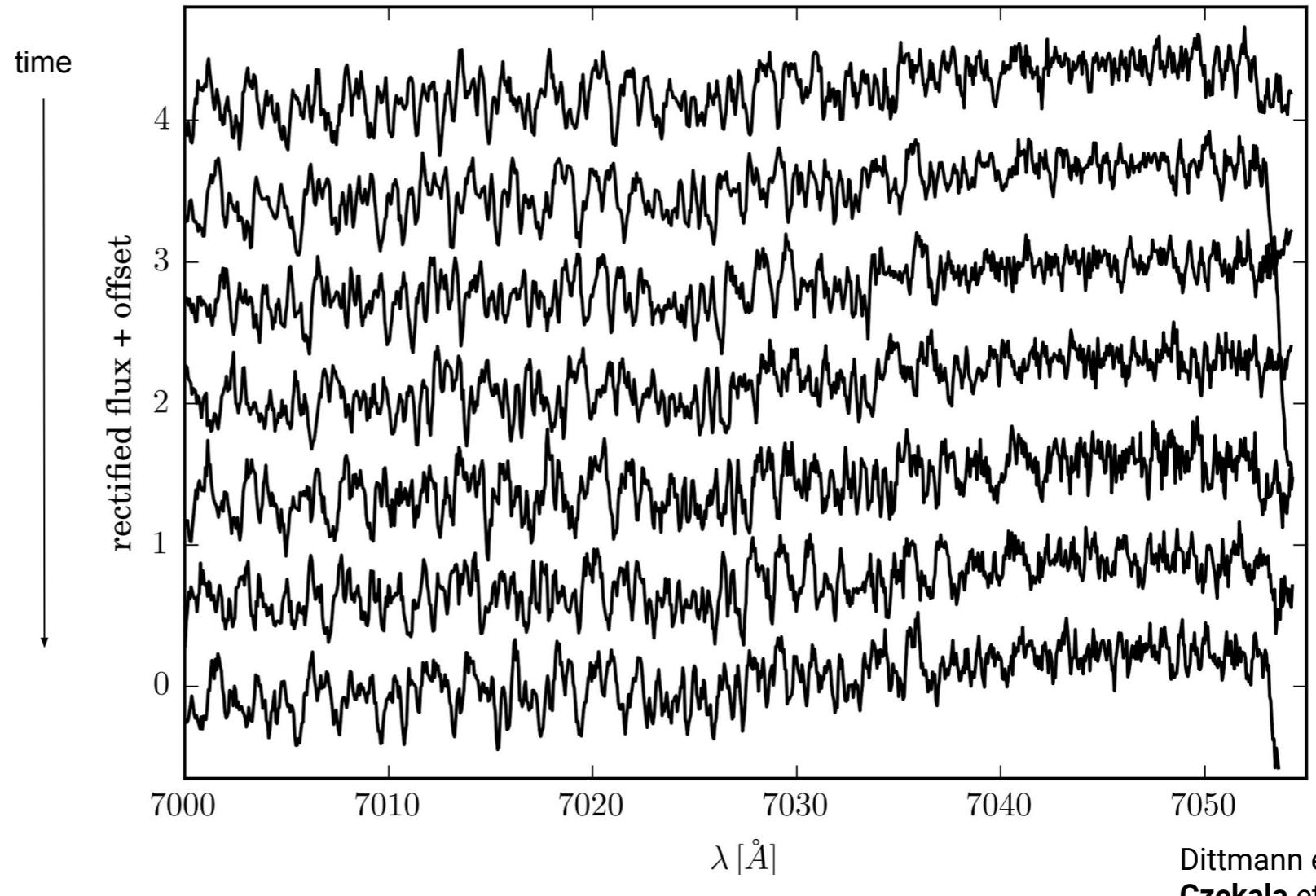
$\alpha =$   
 $(a_f, l_f, a_g, l_g)$

4-dim GP hyperparameters = controlling the amplitude and smoothness of Gaussian Process prior on latent spectra

# Knowns (Data)

Raw Observations of the LP661-13 M4 Binary

**D** =



Dittmann et al. 17  
Czekala et al. 17a

# Bayesian Inference

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

In this case:

$$\begin{aligned} P(\Theta, f, g, \alpha | D) &\propto \\ P(D | \Theta, f, g, \alpha) \times P(\Theta, f, g, \alpha) \end{aligned}$$

a probability density on (4+7+ $\infty$ )-dim parameter space

# Bayesian Computation

1. Run Markov Chain Monte Carlo (MCMC)  
(e.g. *emcee* affine-invariant ensemble sampler)  
on the 4+7 small dimensional marginal posterior

$$P(\Theta, \alpha | D) = \int df \int dg P(\Theta, f, g, \alpha | D)$$

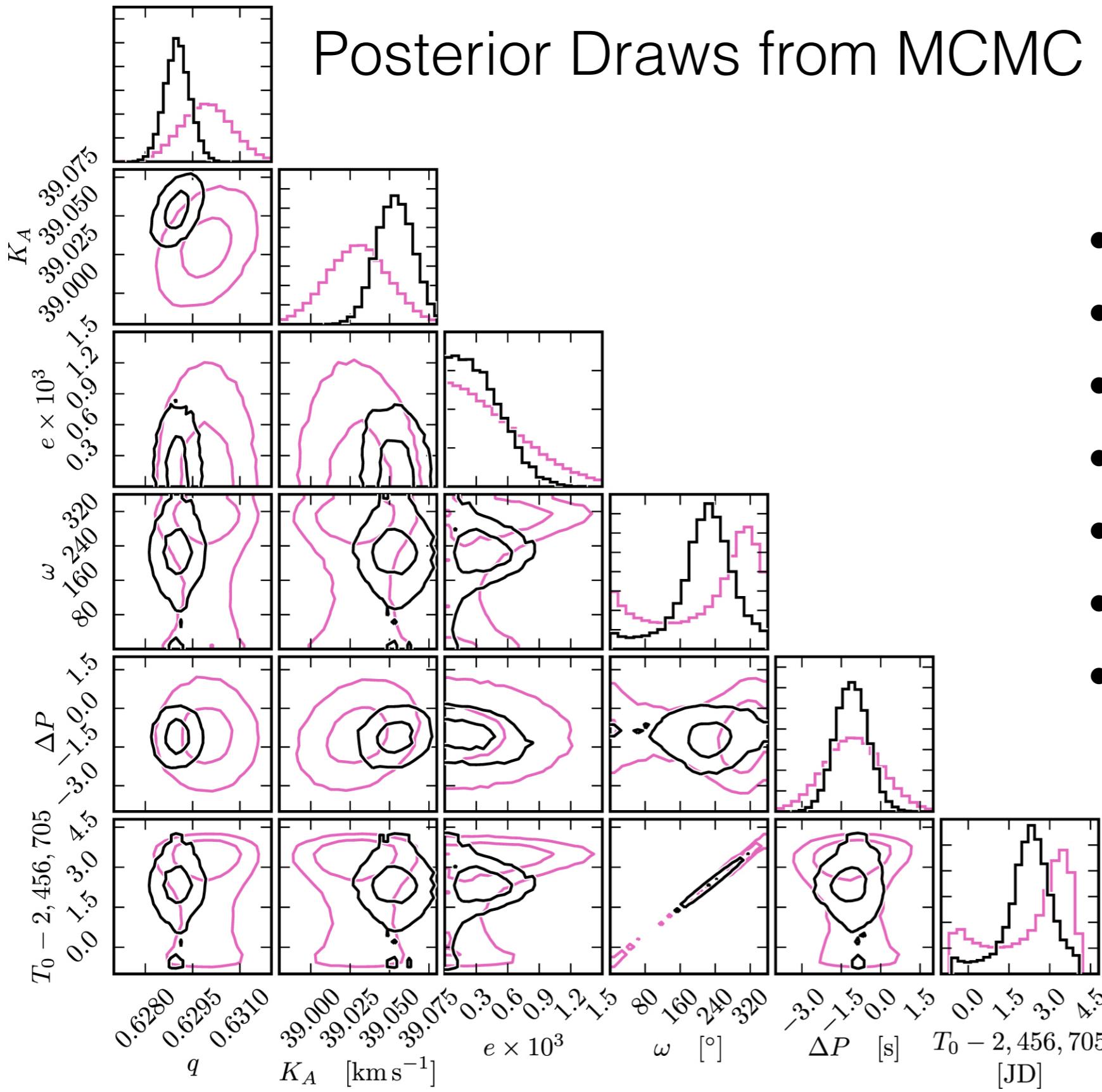
MCMC generates samples:  $\Theta_i, \alpha_i \sim P(\Theta, \alpha | D)$

2. Draw high-dim (**f**, **g**) spectra from the posterior predictive distribution

$$f_i, g_i \sim P(f, g | \Theta_i, \alpha_i, D)$$

# Application to the Mid-M-Dwarf Binary LP661-13

Posterior Draws from MCMC



Seven Orbital  
Parameters:

- Mass Ratio
- Velocity Amplitude
- eccentricity
- Arg of Periastron
- Epoch of Periastron
- Orbital Period
- Systemic Velocity

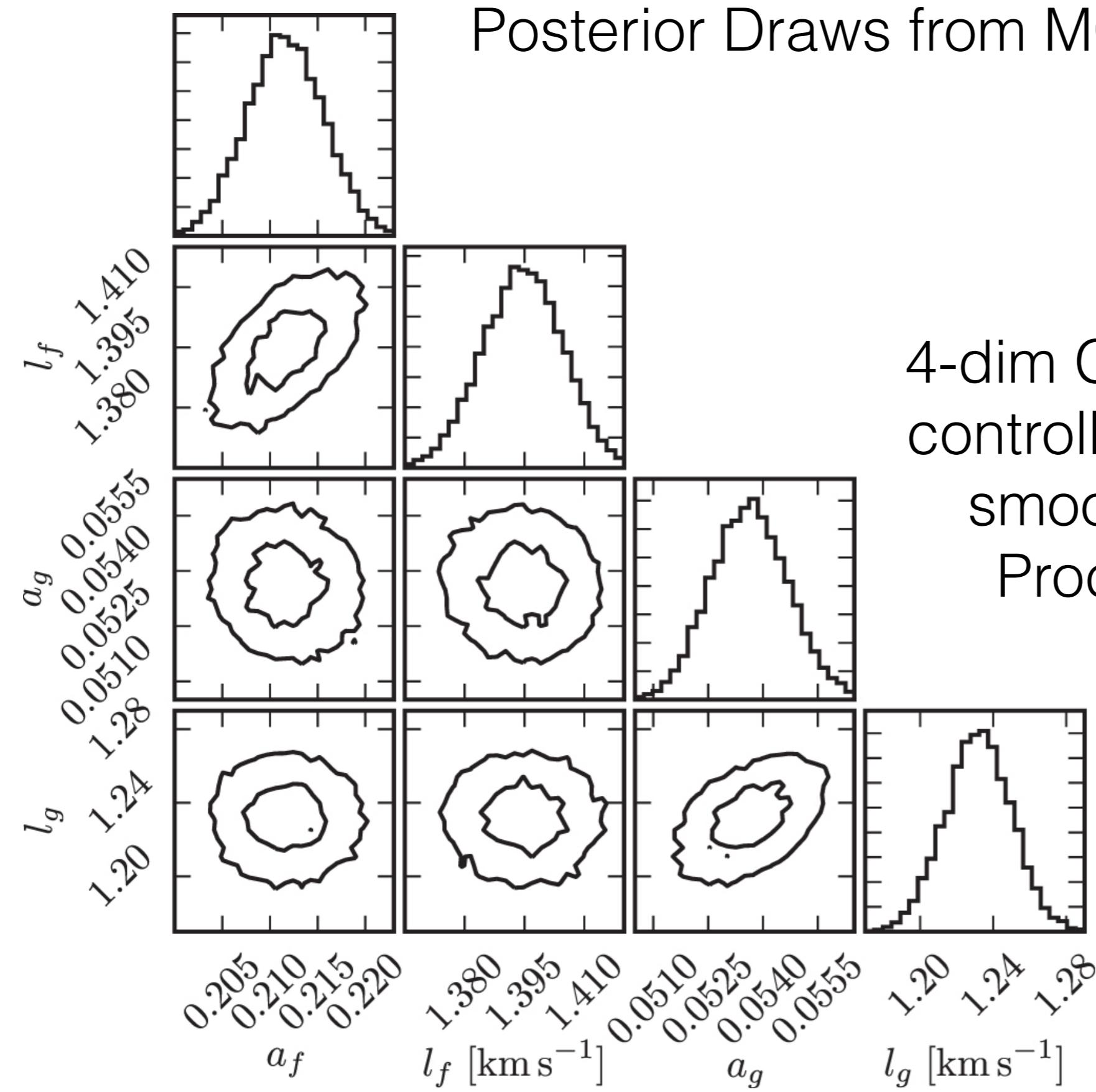
(Purple:  
Conventional  
Analysis)

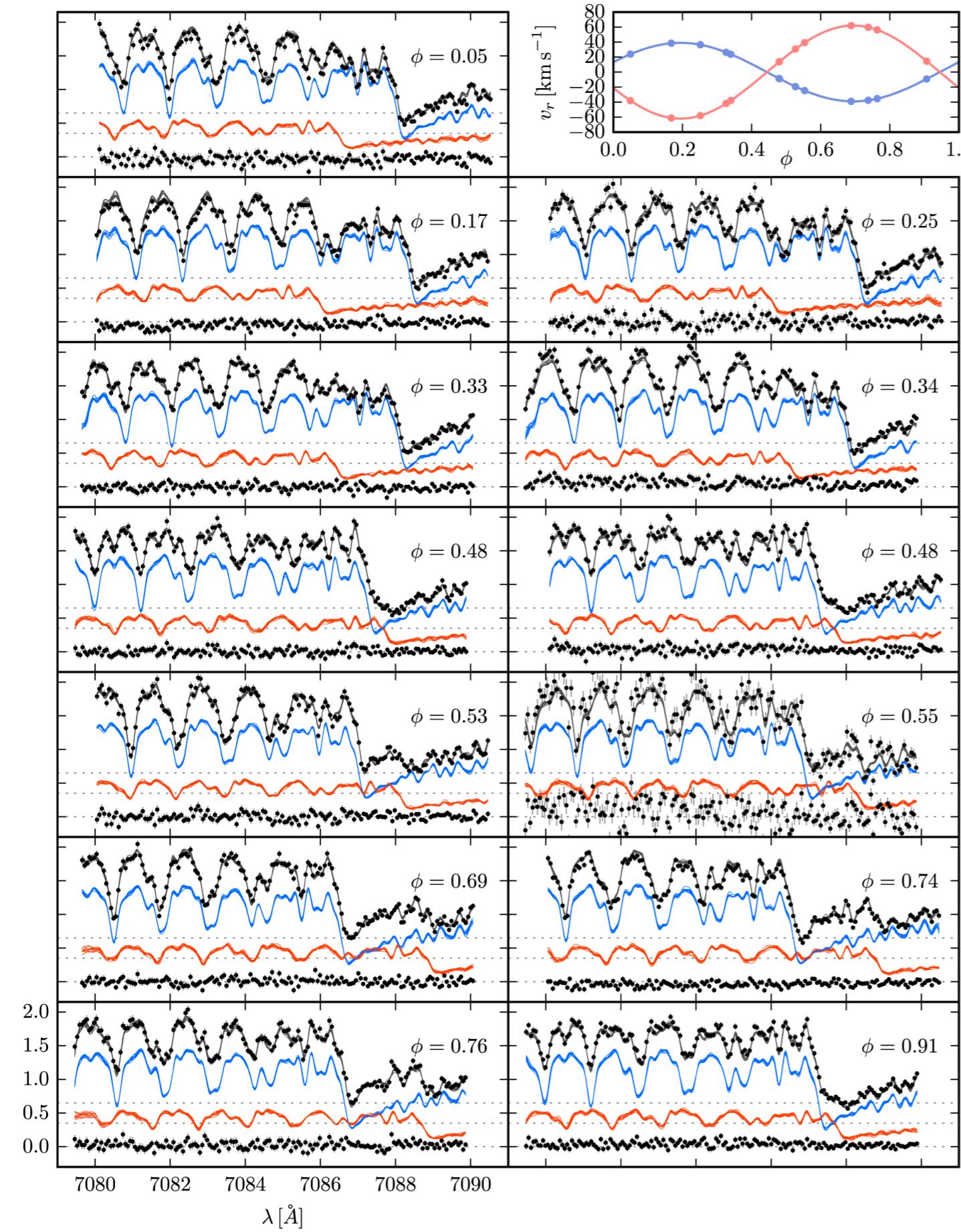
# Application to the Mid-M-Dwarf Binary LP661-13

Posterior Draws from MCMC     $\alpha =$

$$(a_f, l_f, a_g, l_g)$$

4-dim GP hyperparameters =  
controlling the amplitude and  
smoothness of Gaussian  
Process prior on latent  
spectra





Posterior Inference of  
Component Spectra  
**(f, g)**

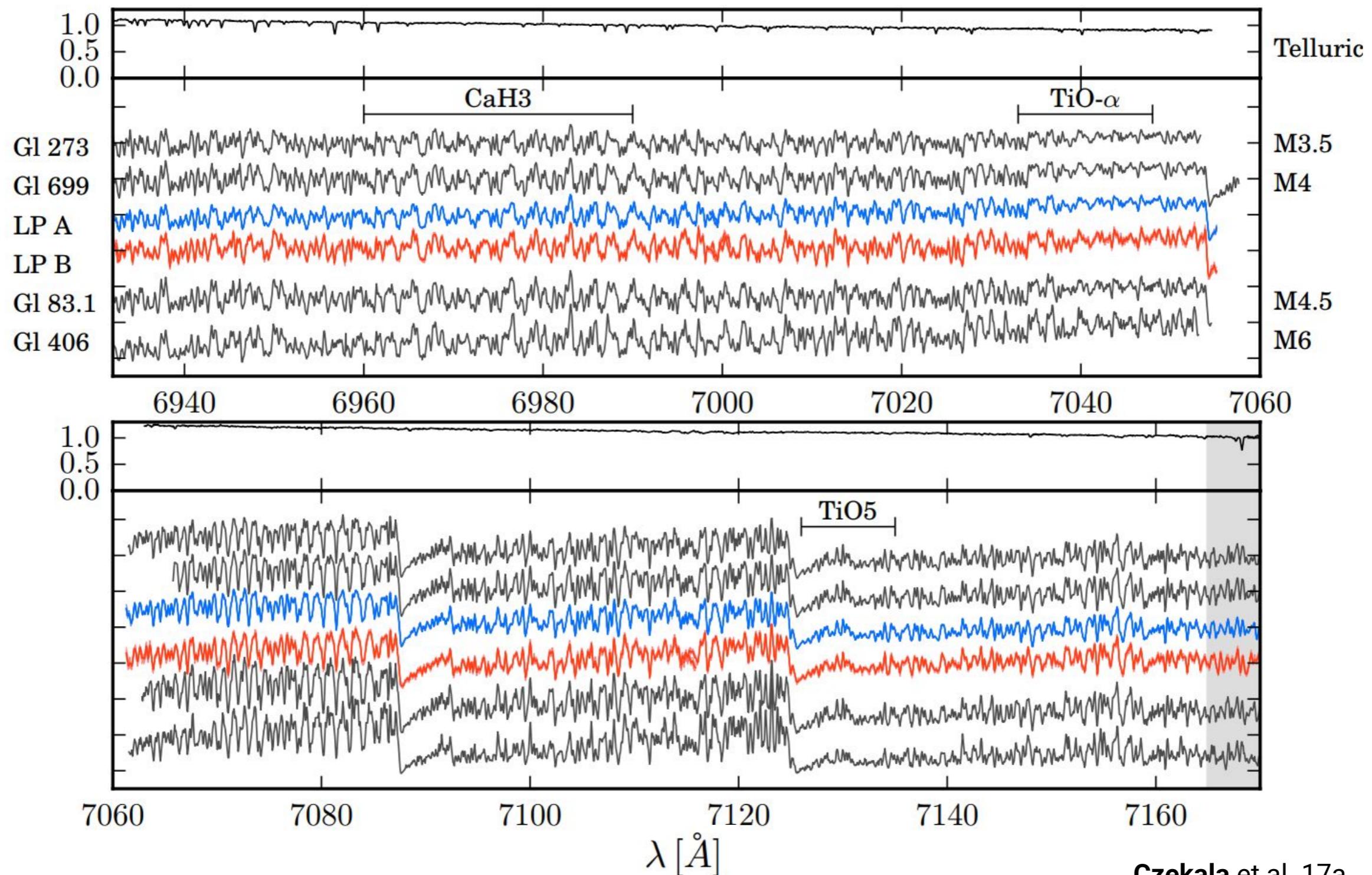
Compared to 10 epochs of  
observed spectra **(data)**

Model Checking!  
Checking Fit against Data

# Model Checking!

## Checking Fit against Domain Knowledge (astrophysics)!

**Disentangled spectra match other single standard stars**



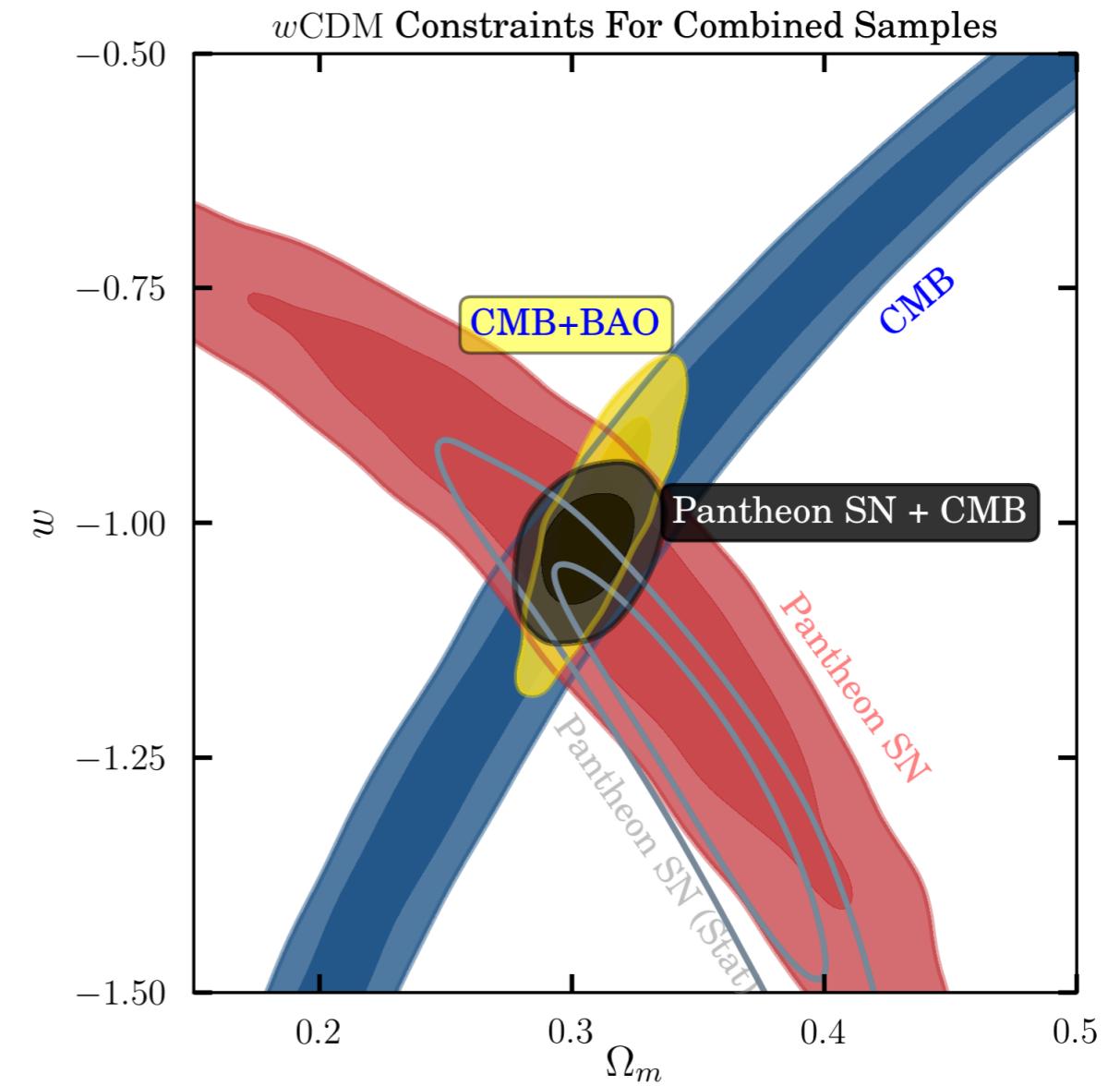
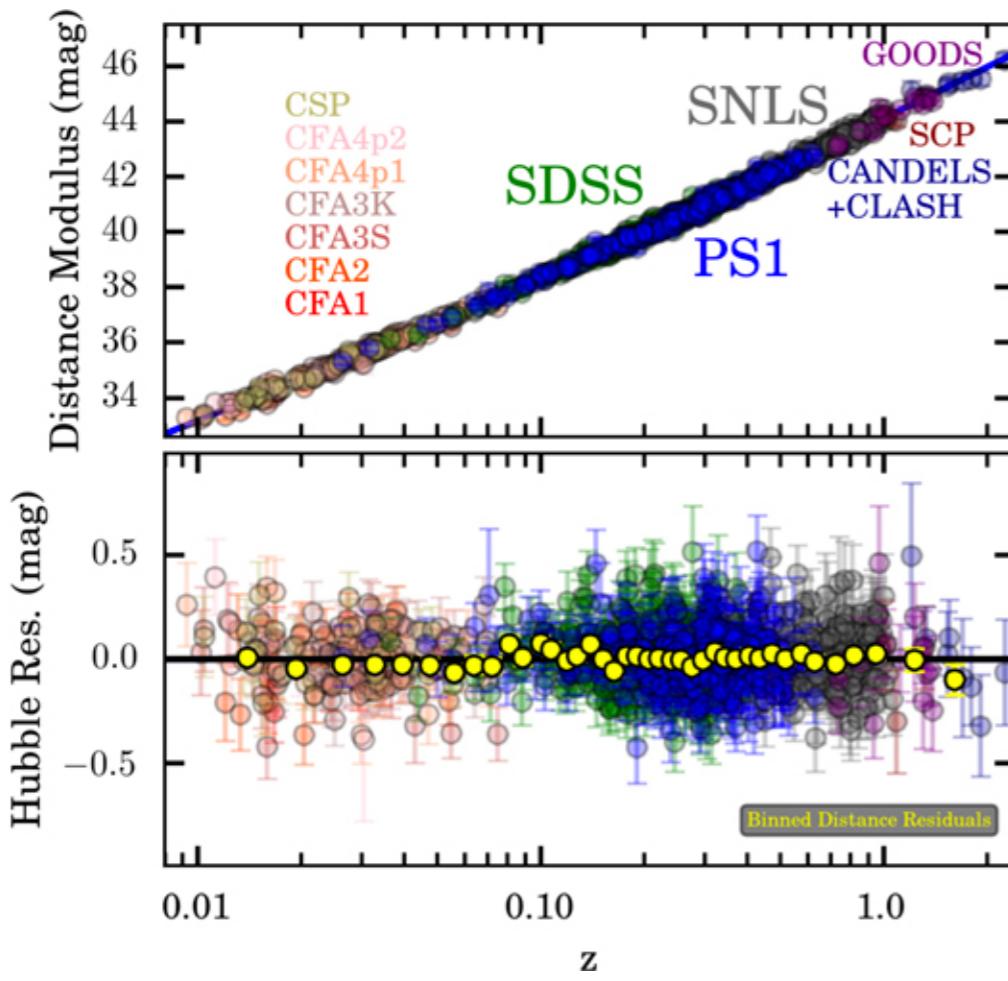
# *Astrostatistics Case Study 1:* Disentangling Time Series Spectra with Gaussian Processes: Applications to Radial Velocity Analysis (Czekala et al. 2017, arXiv:1702.05652)

<http://psoap.readthedocs.io/en/latest/>

- Statistics:
  - Parametric Modelling (Stellar Orbit Parameters)
  - Nonparametric Modelling (Gaussian Process Spectrum)
  - Bayesian Inference (probability of unknowns given data)
  - Markov Chain Monte Carlo (computing posterior probability)
- Astronomy:
  - Applications to Radial Velocity Analysis of Stars/Exoplanets

# Astrostatistics Case Study 4:

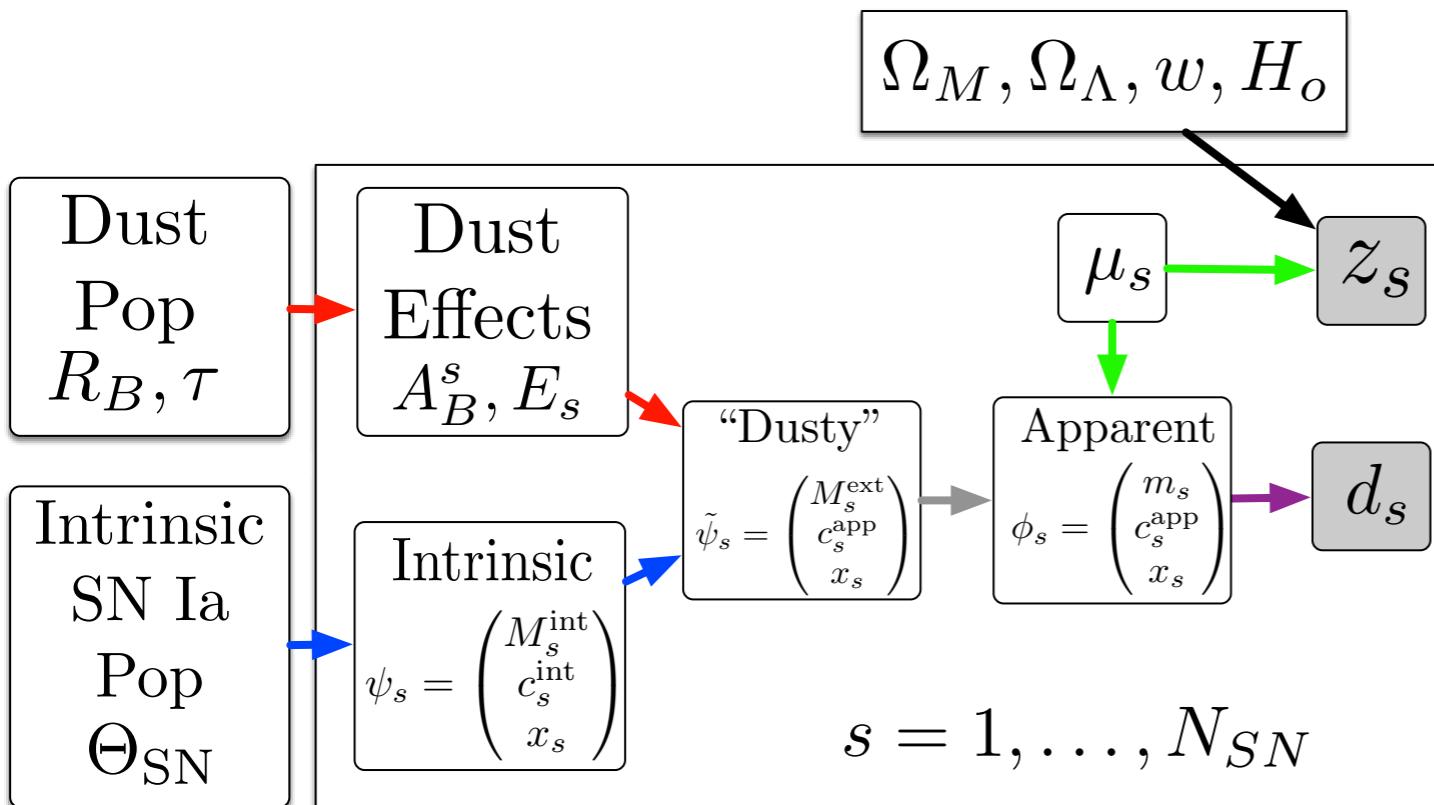
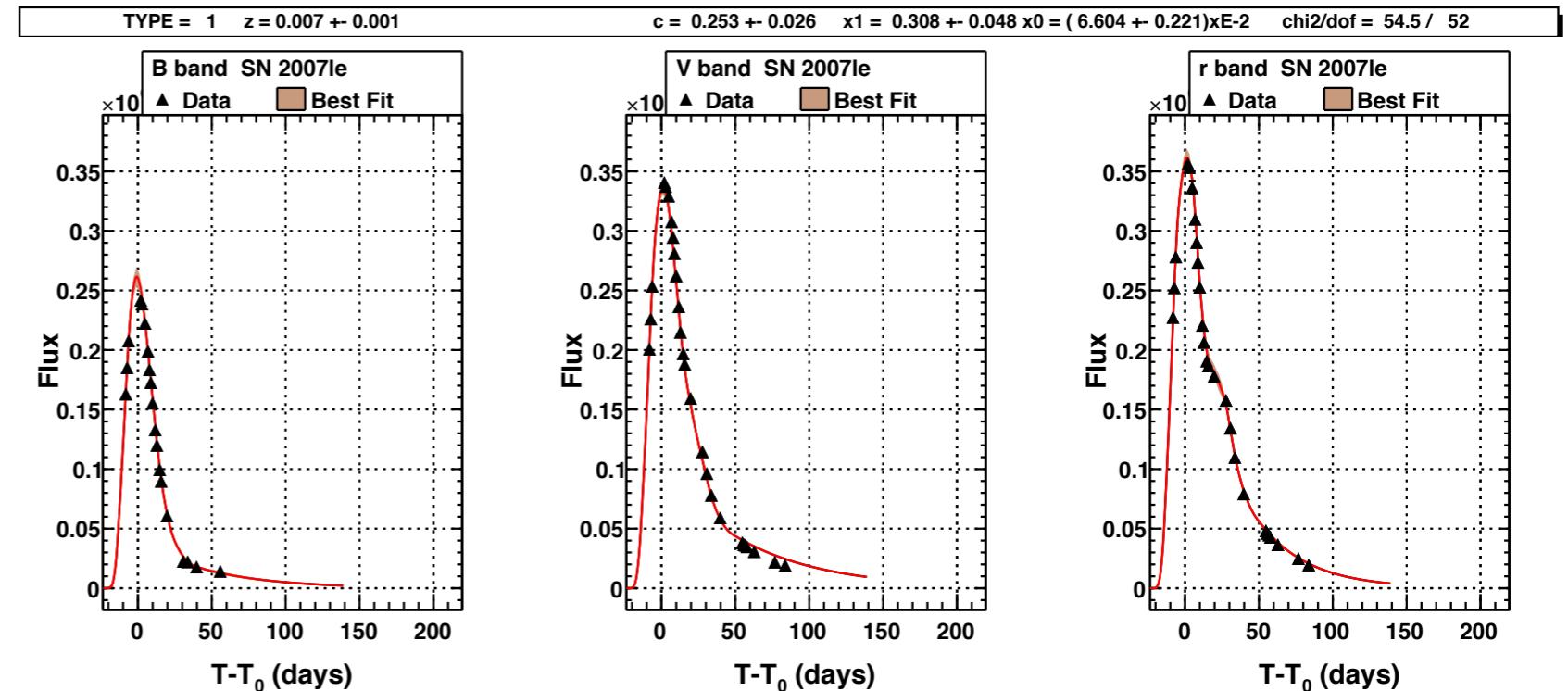
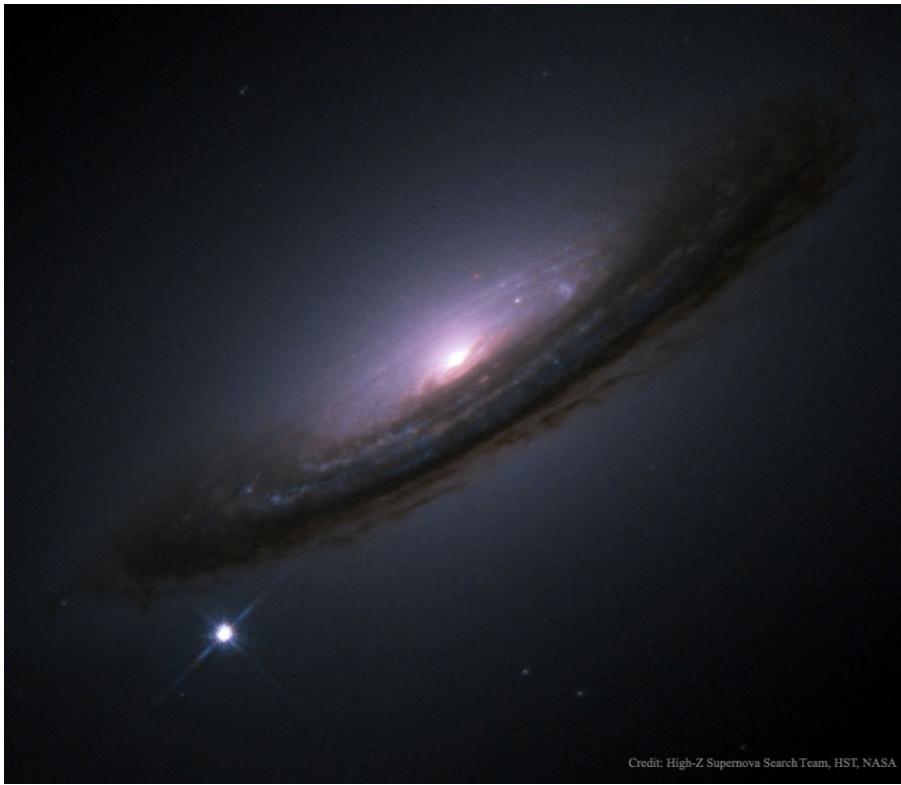
## Supernova Cosmology



**Figure 20.** Confidence contours at 68% and 95% for the  $\Omega_m$  and  $w$  cosmological parameters for the  $w$ CDM model. Constraints from CMB (blue), SN - with systematic uncertainties (red), SN - with only statistical uncertainties (gray-line), and SN+CMB (purple) are shown.

# Astrostatistics Case Study 4:

## Hierarchical Bayesian Models for Supernovae (Mandel et al. 2017)



- Hierarchical Bayes / Multilevel Models
- Probabilistic Graphical Models
- Latent Variable Modelling
- Time Series Analysis
- MCMC