

Astrostatistics: 06 Mar 2019

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

- Example Class 3: Thu, 14 March, 12-2pm MR9
- Example Class 4: Probably ~ April 29-30
- Make-up Lecture:
 - Tue 12 Mar, 2pm MR 14, or Fri 15 Mar 12pm MR13 ?
- Finish Gaussian Processes in Astrophysics
- Case Study: application to gravitationally lensed supernova
- Hierarchical Bayes

Human Learning of Gaussian Processes

- Classic Text: Rasmussen & Williams (2006)
 - “Gaussian Processes for Machine Learning”, Ch 1-2,4-5
 - Free Online: <http://www.gaussianprocess.org/gpml/>
- Ivezic, Sec 8.10 GP Regression, (Ch 8 is Regression)
- **Bishop: Pattern Recognition & Machine Learning, Ch 6**
 - **Also free online:**
<https://www.microsoft.com/en-us/research/people/cmbishop/#!prml-book>
- Gelman, Bayesian Data Analysis 3rd Ed., Chapter 21
- “Practical Introduction to GPs for Astronomy” - D. Foreman-Mackey
 - http://hea-www.harvard.edu/AstroStat/aas231_2018/DForeman-Mackey_20180110_aas231.pdf

Human Learning of Gaussian Processes

Christopher Bishop at Microsoft

Secure | https://www.microsoft.com/en-us/research/people/cmbishop/#!prml-book

Christopher Bishop

Technical Fellow and Laboratory Director, Microsoft Research Cambridge

Contact Info

+44 (0)1223 479993
Email
LinkedIn
Twitter

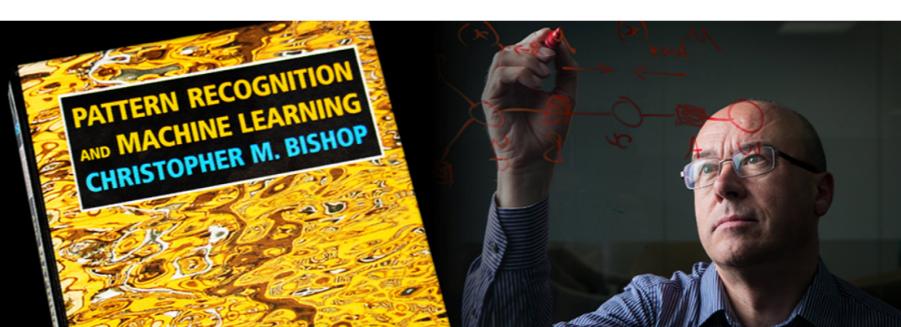
Research areas

Artificial intelligence

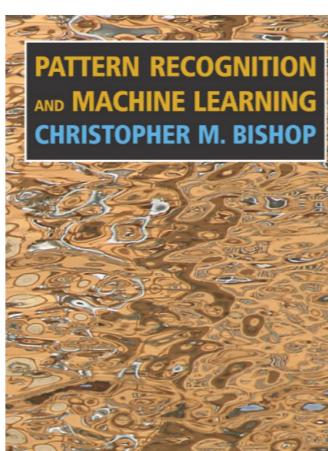
About Publications Videos Downloads PRML Book

Pattern Recognition and Machine Learning by Christopher Bishop

Download FREE PDF



Pattern Recognition and Machine Learning



<https://www.microsoft.com/en-us/research/uploads/prod/2006/01/Bishop-Pattern-Recognition-and-Machine-Learning-2006.pdf> A comprehensive introduction to the fields of pattern recognition and machine learning.

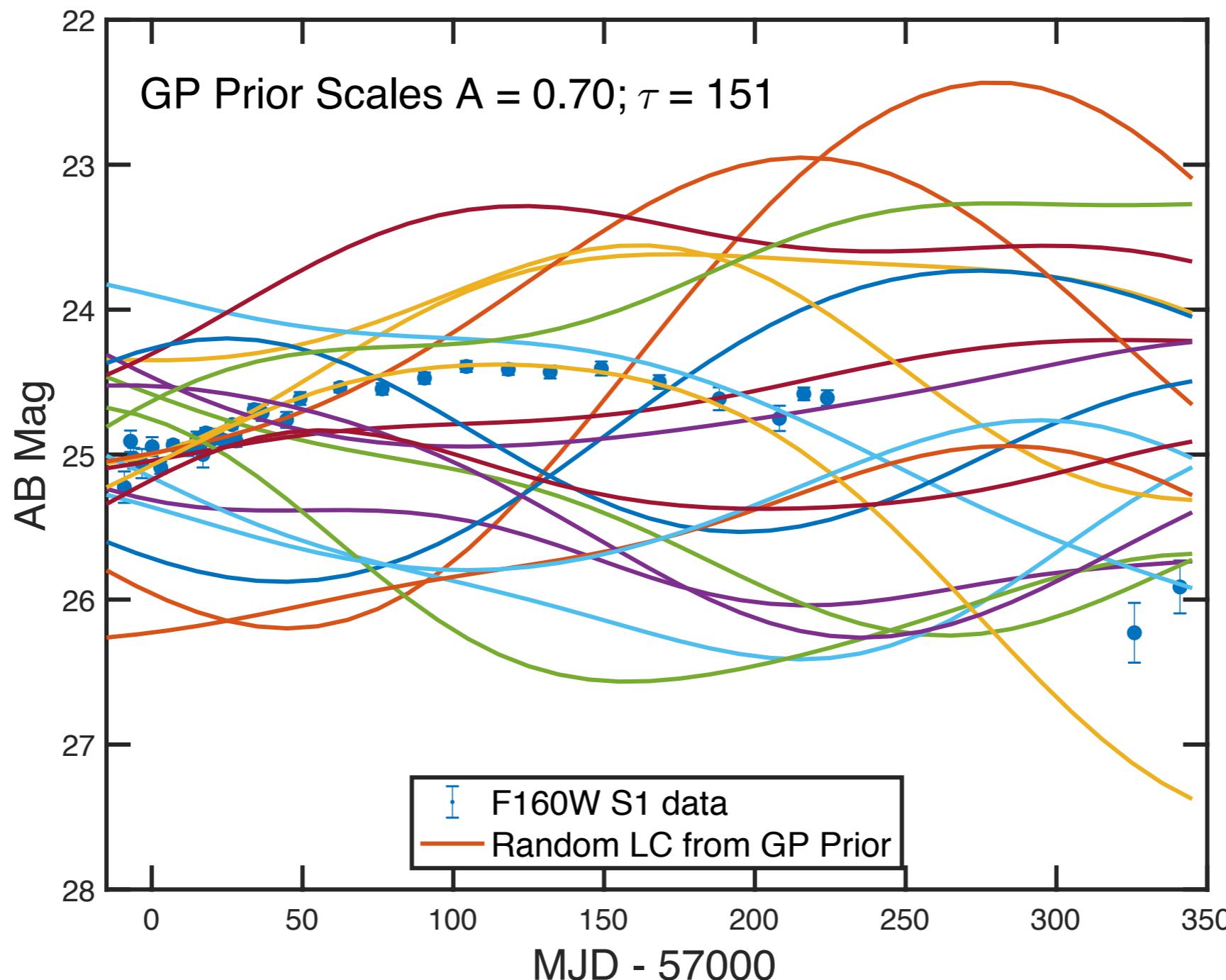
Review: Fitting a GP to data

1. If we knew the characteristic scales of the kernel (A , τ^2), then how do we fit the data at observed times to find the curve for unobserved times? (computing the posterior)
2. How do we fit for the characteristic scales of the kernel (hyperparameters)? (model selection)

1. Which curve from the prior is the best description of the data?

$$\text{Cov}[f(t), f(t')] = k(t, t') = A \exp(-|t - t'|^2/\tau^2)$$

$$f | A, \tau \sim N(1c, K)$$



Posterior Inference with GPs

Estimating the underlying curve:

f_o = observed points at times t_o (training set)

f_* = function at unobserved times t_* (prediction or test set)

Joint:

$$\begin{pmatrix} f_o \\ f_* \end{pmatrix} \sim N \left(\begin{bmatrix} 1c \\ 1c \end{bmatrix}, \begin{bmatrix} K(t_o, t_o) & K(t_*, t_o) \\ K(t_o, t_*) & K(t_*, t_*) \end{bmatrix} \right)$$

Populating the Covariance Matrix

$K(t, t')$ has i,j-th entry = $k(t_i, t'_j)$

Using the assumed kernel function

$$\text{Cov}[f(t), f(t')] = k(t, t') = A \exp(-|t - t'|^2 / \tau^2)$$

Review: Posterior Inference with GPs

Estimating the underlying curve:

f_o = observed points at times t_o

f_* = function at unobserved times t_*

Jointly Gaussian: $\begin{pmatrix} f_o \\ f_* \end{pmatrix} \sim N \left(\begin{bmatrix} 1c \\ 1c \end{bmatrix}, \begin{bmatrix} K(t_o, t_o) & K(t_*, t_o) \\ K(t_o, t_*) & K(t_*, t_*) \end{bmatrix} \right)$

Posterior is also Gaussian $f_*|f_o \sim N(\mathbb{E}[f_*|f_o], \text{Var}[f_*|f_o])$

Posterior Mean:

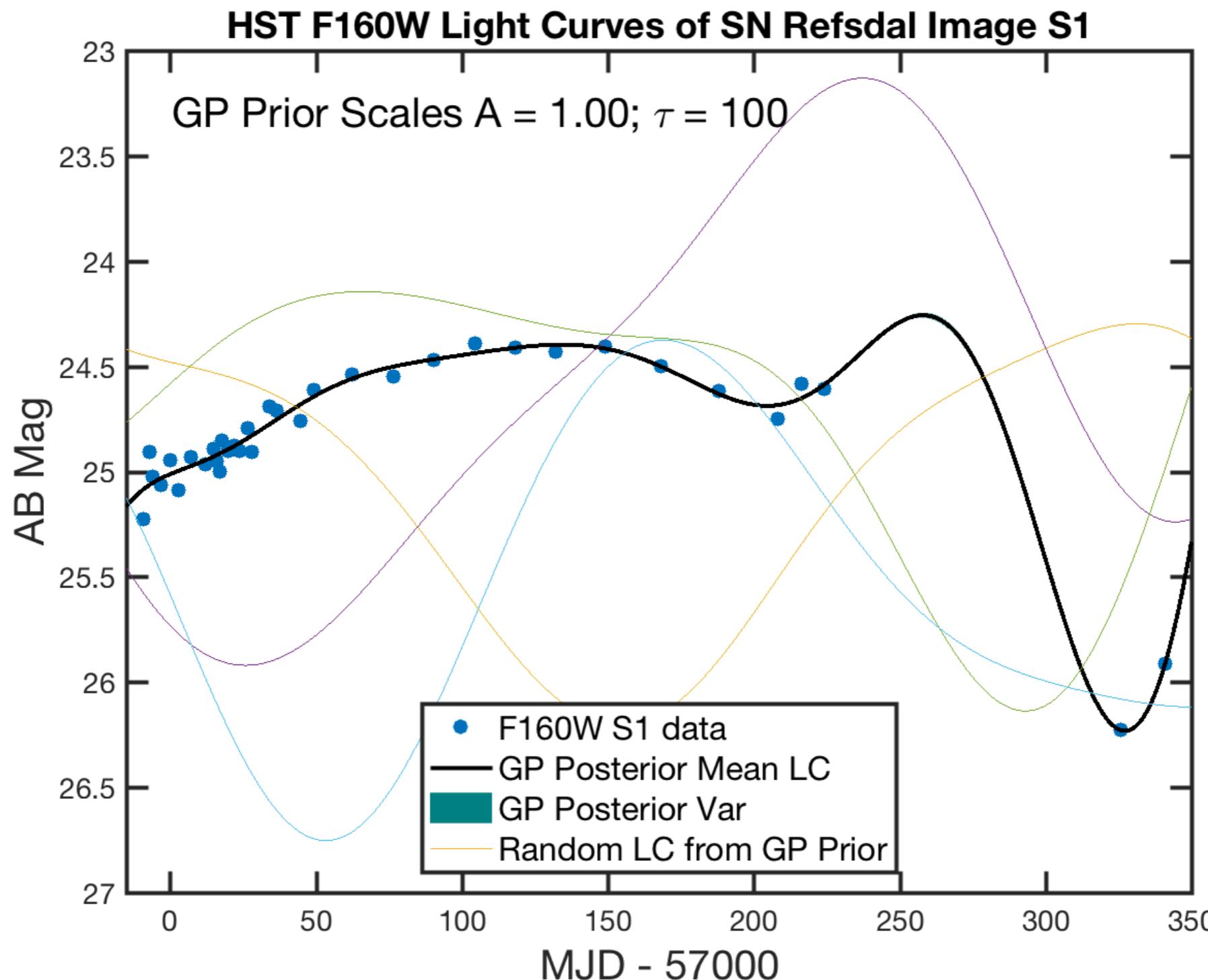
$$\mathbb{E}[f_*|f_o] = 1c + K(t_*, t_o)K(t_o, t_o)^{-1}(f_o - 1c)$$

Posterior Co(variance):

$$\text{Var}[f_*|f_o] = K(t_*, t_*) - K(t_*, t_o)K(t_o, t_o)^{-1}K(t_o, t_*)$$

Posterior Inference with GPs

Estimating the underlying curve:



Accounting for Measurement Error

$$\mathbf{y}_o | \mathbf{f}_o \sim N(\mathbf{f}_o, \mathbf{W})$$

\mathbf{y}_o are measured values of \mathbf{f}_o at times t_o

\mathbf{W} is measurement covariance matrix
(often diagonal for independent noise)

$$W_{ij} = \delta_{ij} \sigma_i^2$$

$$\begin{pmatrix} \mathbf{y}_o \\ \mathbf{f}_* \end{pmatrix} \sim N \left(\begin{bmatrix} \mathbf{1}_c \\ \mathbf{1}_c \end{bmatrix}, \begin{bmatrix} \mathbf{K}(\mathbf{t}_o, \mathbf{t}_o) + \mathbf{W} & \mathbf{K}(\mathbf{t}_*, \mathbf{t}_o) \\ \mathbf{K}(\mathbf{t}_o, \mathbf{t}_*) & \mathbf{K}(\mathbf{t}_*, \mathbf{t}_*) \end{bmatrix} \right)$$

Now can calculate function prediction at unobserved points

$$\mathbf{f}_* | \mathbf{y}_o \sim N(\mathbb{E}[\mathbf{f}_* | \mathbf{y}_o], \text{Var}[\mathbf{f}_* | \mathbf{y}_o])$$

Using conditional properties of Gaussian as before

Accounting for Measurement Error:
Derivation as the sum of two GPs at the observed times

GP of Intrinsic Curve

$$f(t) \sim \mathcal{GP}(m(t) = c, k(t, t'))$$

f_o = function at observed times t_o

$$f_o \sim N[\mathbf{1}_c, \mathbf{K}(t_o, t_o)]$$

GP of Measurement Error

$$\mathbf{y}_o | f_o \sim N(f_o, \mathbf{W})$$

Same as: (mean-zero noise)

$$\mathbf{y}_o = f_o + \boldsymbol{\epsilon} \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{W})$$

Most common case:
heteroskedastic uncorrelated measurement error:
 $\text{Cov}(\epsilon_i, \epsilon_j) \equiv W_{ij} = \delta_{ij} \sigma_i^2$

Accounting for Measurement Error:
Derivation as the sum of two GPs at the observed times

Intrinsic/Latent Process: $f_o \sim N[\mathbf{1}_c, K(t_o, t_o)]$

Measurement Process: $y_o | f_o \sim N(f_o, W)$

$$y_o = f_o + \epsilon \quad \epsilon \sim N(\mathbf{0}, W)$$

$$\text{Cov}(y_o, y_o) = \text{Cov}(f_o, f_o) + \text{Cov}(\epsilon, \epsilon) + 2 \text{Cov}(f_o, \epsilon)$$

$$\text{Cov}(f_o, f_o) = K(t_o, t_o) \quad (\text{GP of intrinsic curve})$$

$$\text{Cov}(\epsilon, \epsilon) = W \quad (\text{measurement noise})$$

(the two processes are uncorrelated)

$$2\text{Cov}[f_o, \epsilon] = 0$$

Therefore: $\text{Cov}[y_o, y_o] = K(t_o, t_o) + W$

Accounting for Measurement Error:
Derivation as the sum of two GPs at the observed times

$$\mathbf{y}_o | \mathbf{f}_o \sim N(\mathbf{f}_o, \mathbf{W})$$

$$\mathbf{y}_o = \mathbf{f}_o + \boldsymbol{\epsilon} \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{W})$$

$$\begin{pmatrix} \mathbf{f}_o \\ \mathbf{f}_* \end{pmatrix} \sim N \left(\begin{bmatrix} 1_c \\ 1_c \end{bmatrix}, \begin{bmatrix} \mathbf{K}(\mathbf{t}_o, \mathbf{t}_o) & \mathbf{K}(\mathbf{t}_*, \mathbf{t}_o) \\ \mathbf{K}(\mathbf{t}_o, \mathbf{t}_*) & \mathbf{K}(\mathbf{t}_*, \mathbf{t}_*) \end{bmatrix} \right)$$

Similar arguments for:

$$\text{Cov}[\mathbf{y}_o, \mathbf{f}_*] = \text{Cov}[\mathbf{f}_o, \mathbf{f}_*] + \text{Cov}[\boldsymbol{\epsilon}, \mathbf{f}_*]$$

$$\text{Cov}[\mathbf{y}_o, \mathbf{f}_*] = \mathbf{K}(\mathbf{t}_o, \mathbf{t}_*) + 0 = \mathbf{K}(\mathbf{t}_o, \mathbf{t}_*)$$

$$\begin{pmatrix} \mathbf{y}_o \\ \mathbf{f}_* \end{pmatrix} \sim N \left(\begin{bmatrix} 1_c \\ 1_c \end{bmatrix}, \begin{bmatrix} \mathbf{K}(\mathbf{t}_o, \mathbf{t}_o) + \mathbf{W} & \mathbf{K}(\mathbf{t}_*, \mathbf{t}_o) \\ \mathbf{K}(\mathbf{t}_o, \mathbf{t}_*) & \mathbf{K}(\mathbf{t}_*, \mathbf{t}_*) \end{bmatrix} \right)$$

Accounting for Measurement Error: Derivation using Conditional/Marginal properties of MV Gaussian

$$\begin{pmatrix} f_o \\ f_* \end{pmatrix} \sim N \left(\begin{bmatrix} 1_c \\ 1_c \end{bmatrix}, \begin{bmatrix} K(t_o, t_o) & K(t_*, t_o) \\ K(t_o, t_*) & K(t_*, t_*) \end{bmatrix} \right)$$

$$\begin{pmatrix} y_o \\ f_* \end{pmatrix} \mid \begin{pmatrix} f_o \\ f_* \end{pmatrix} \sim N \left(\begin{bmatrix} f_o \\ f_* \end{bmatrix}, \begin{bmatrix} W & 0 \\ 0 & 0 \end{bmatrix} \right)$$

$$P(V) \& P(U|V) \longrightarrow P(U)$$

$$\begin{pmatrix} y_o \\ f_* \end{pmatrix} \sim N \left(\begin{bmatrix} 1_c \\ 1_c \end{bmatrix}, \begin{bmatrix} K(t_o, t_o) + W & K(t_*, t_o) \\ K(t_o, t_*) & K(t_*, t_*) \end{bmatrix} \right)$$

Now can calculate function prediction at unobserved points

$$f_* | y_o \sim N(\mathbb{E}[f_* | y_o], \text{Var}[f_* | y_o])$$

Using conditional properties of Gaussian as before

Accounting for Measurement Error:

$$\begin{pmatrix} \mathbf{y}_o \\ f_* \end{pmatrix} \sim N \left(\begin{bmatrix} 1c \\ 1c \end{bmatrix}, \begin{bmatrix} \mathbf{K}(t_o, t_o) + \mathbf{W} & \mathbf{K}(t_*, t_o) \\ \mathbf{K}(t_o, t_*) & \mathbf{K}(t_*, t_*) \end{bmatrix} \right)$$

Now can calculate function prediction at unobserved points

$$f_* | \mathbf{y}_o \sim N(\mathbb{E}[f_* | \mathbf{y}_o], \text{Var}[f_* | \mathbf{y}_o])$$

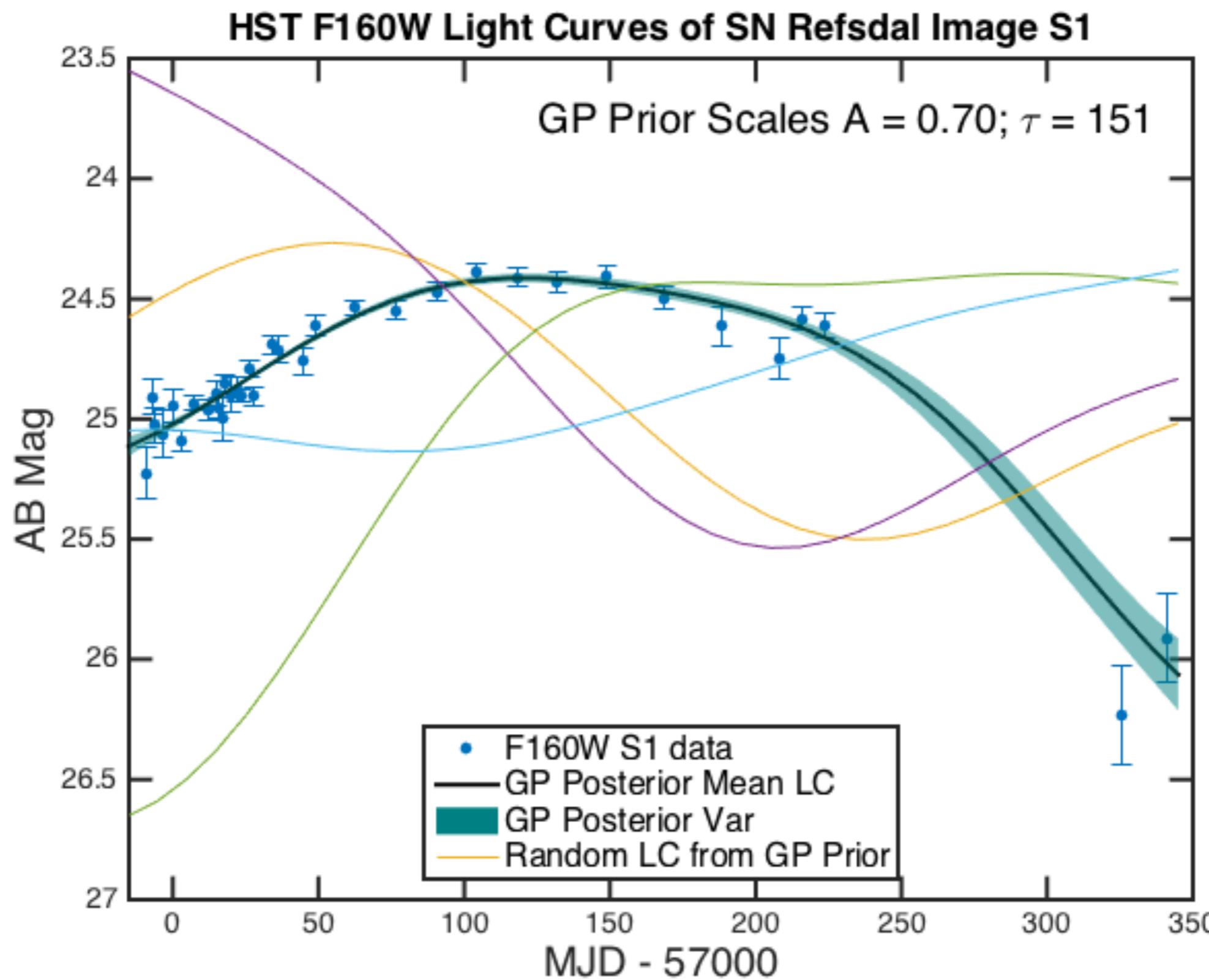
Using Gaussian Conditional Properties:

$$\mathbb{E}[f_* | \mathbf{y}_o] = 1c + \mathbf{K}(t_*, t_o)[\mathbf{K}(t_o, t_o) + \mathbf{W}]^{-1}(\mathbf{y}_o - 1c)$$

$$\text{Var}[f_* | \mathbf{y}_o] = \mathbf{K}(t_*, t_*) - \mathbf{K}(t_*, t_o)[\mathbf{K}(t_o, t_o) + \mathbf{W}]^{-1}\mathbf{K}(t_o, t_*)$$

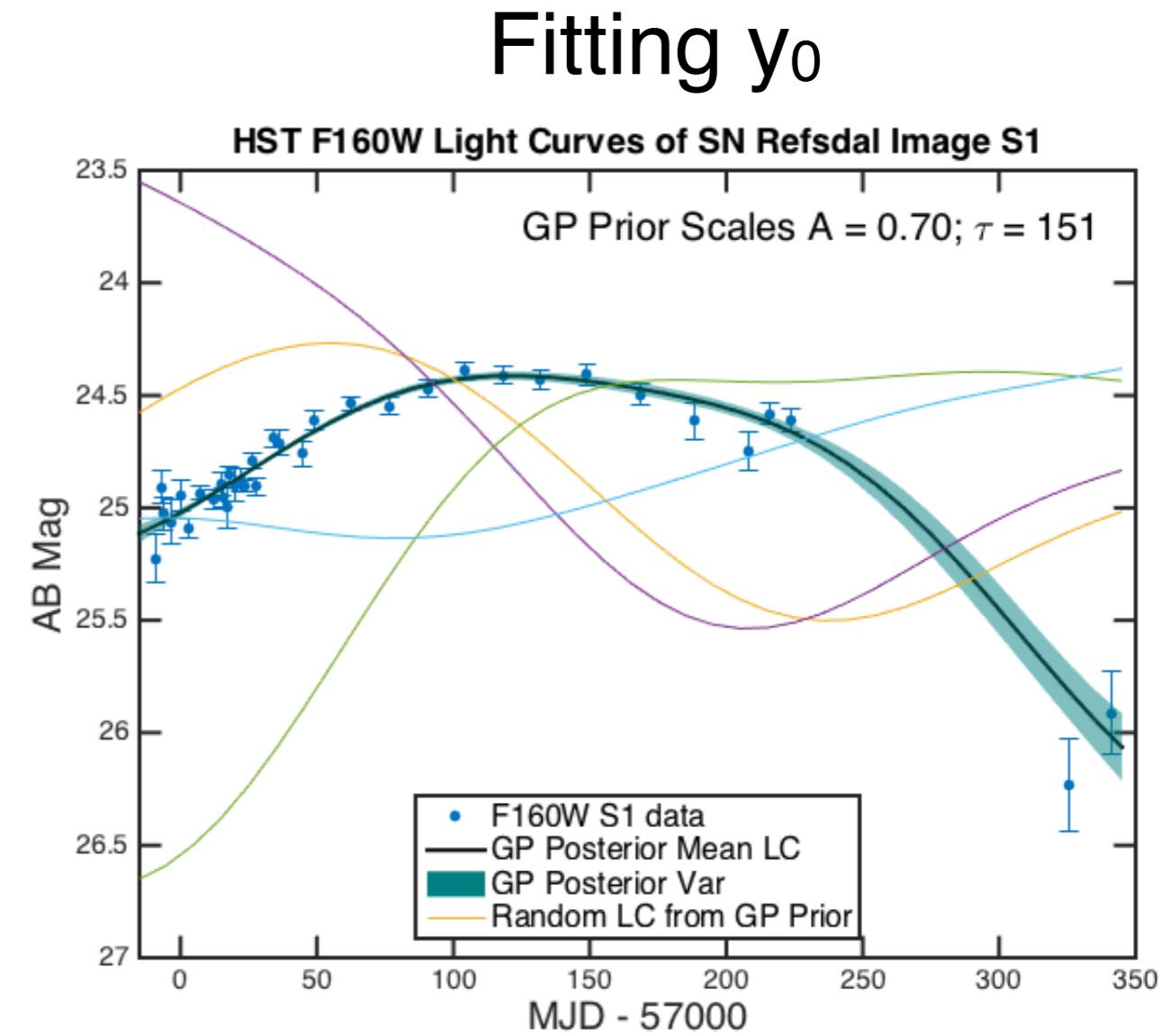
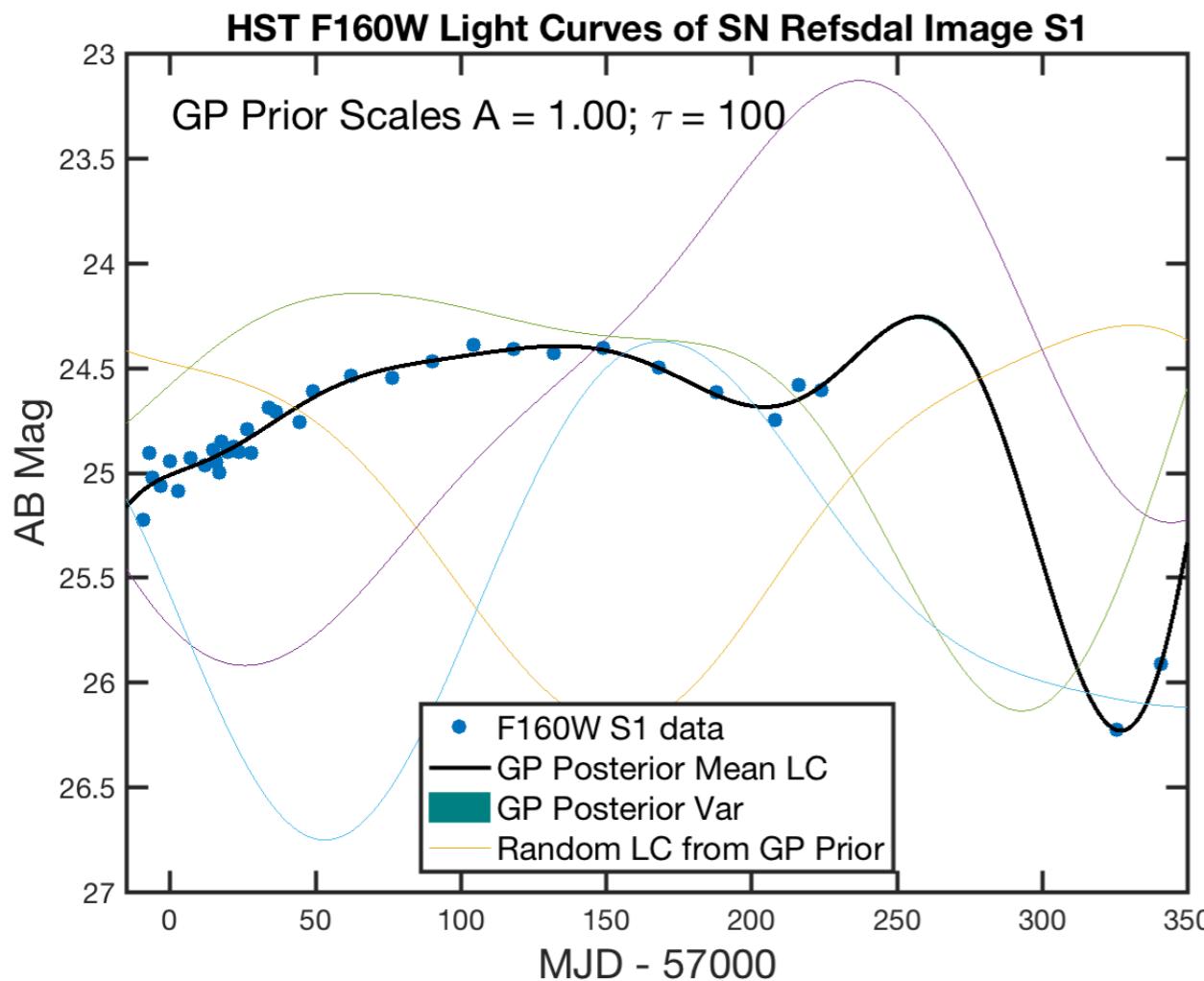
Posterior Inference with GPs

Accounting for Measurement Error



GP Fitting with measurement error

(over)fitting f_0



Actually a small “nugget term” was added to regularise matrix inversion - acts like meas. err.

$$f_0 \sim N(\mathbf{1}_c, \mathbf{K} + \sigma^2 \mathbf{I})$$

$$\sigma = 10^{-3}$$

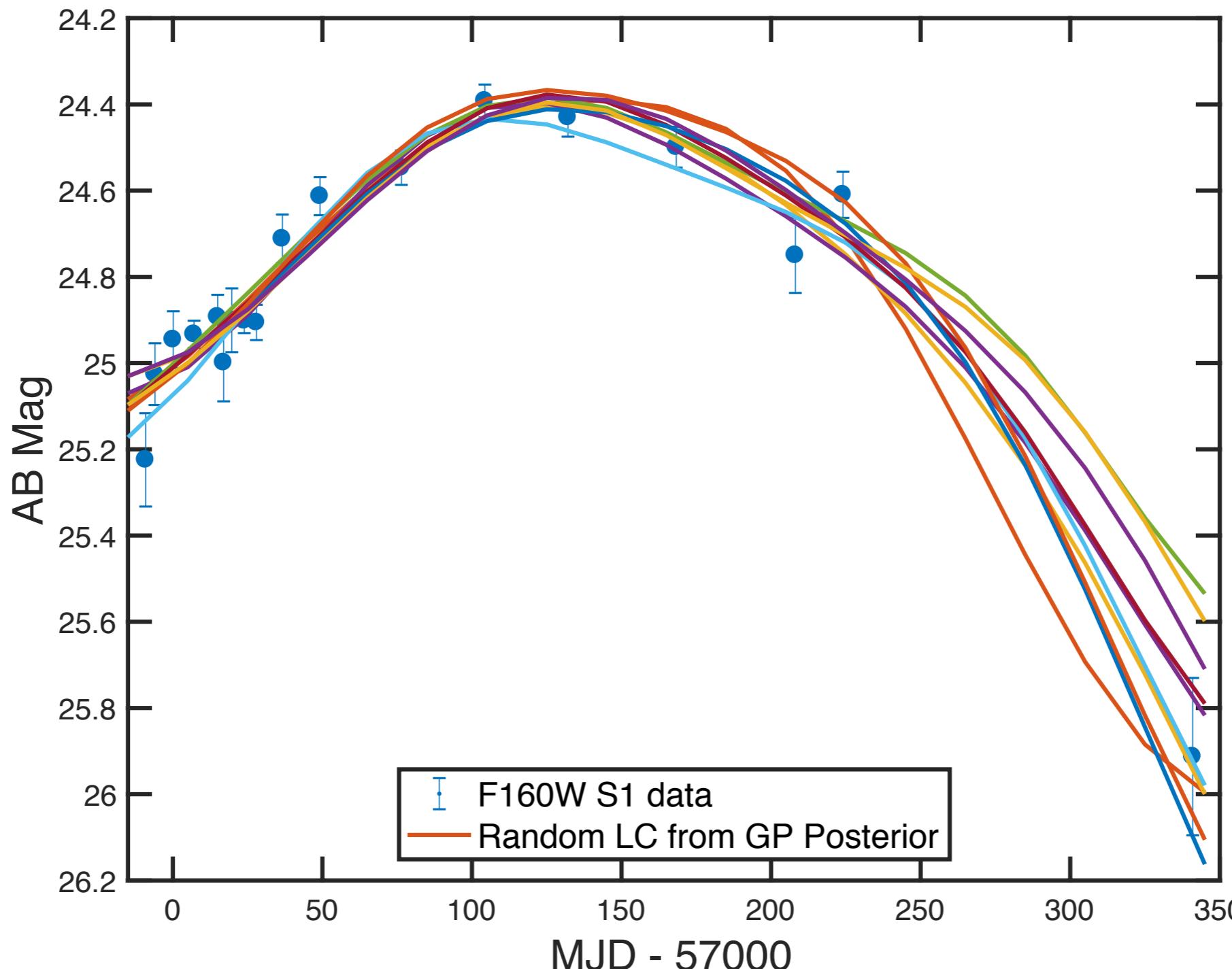
Using actual measurement uncertainties

$$\mathbf{y}_0 \sim N(\mathbf{1}_c, \mathbf{K} + \mathbf{W})$$

$$W_{ij} = \sigma_i^2 \delta_{ij}$$

Posterior Inference with GPs

Accounting for Measurement Error: Random draws from the posterior given noisy data



Fitting a GP to data

1. If we knew the characteristic scales of the kernel (A , τ^2), then how do we fit the data at observed times to find the curve for unobserved times? (computing the posterior)
2. How do we fit for the characteristic scales of the kernel (hyperparameters)? (model selection)

Bayesian Model Selection: tuning the hyperparameters (A , τ)

Integrating out the latent function $f(t)$
gives us the Marginal Likelihood:

$$P(\mathbf{y}_o | A, \tau^2) = \int P(\mathbf{y}_o | f_o) \times P(f_o | A, \tau^2) df_o$$

$$P(\mathbf{y}_o | A, \tau^2) = \int N[\mathbf{y}_o | f_o, \mathbf{W}] \times N[f_o | \mathbf{1}_c, \mathbf{K}(t_o, t_o)] df_o$$

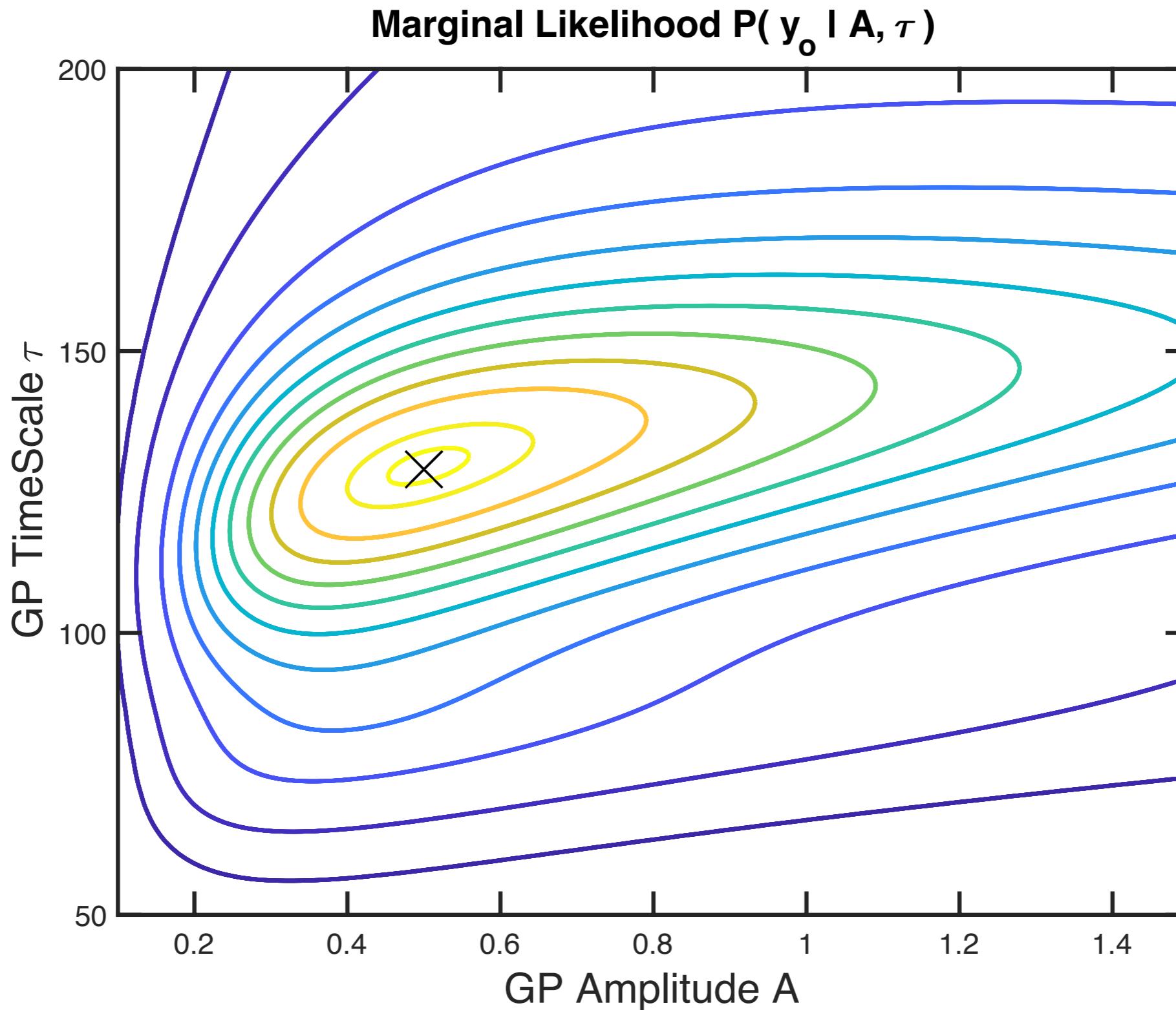
$$L(A, \tau^2) = P(\mathbf{y}_o | A, \tau^2) = N[\mathbf{y}_o | \mathbf{1}_c, \mathbf{K}_{A, \tau^2}(t_o, t_o) + \mathbf{W}]$$

\mathbf{W} = measurement error covariance

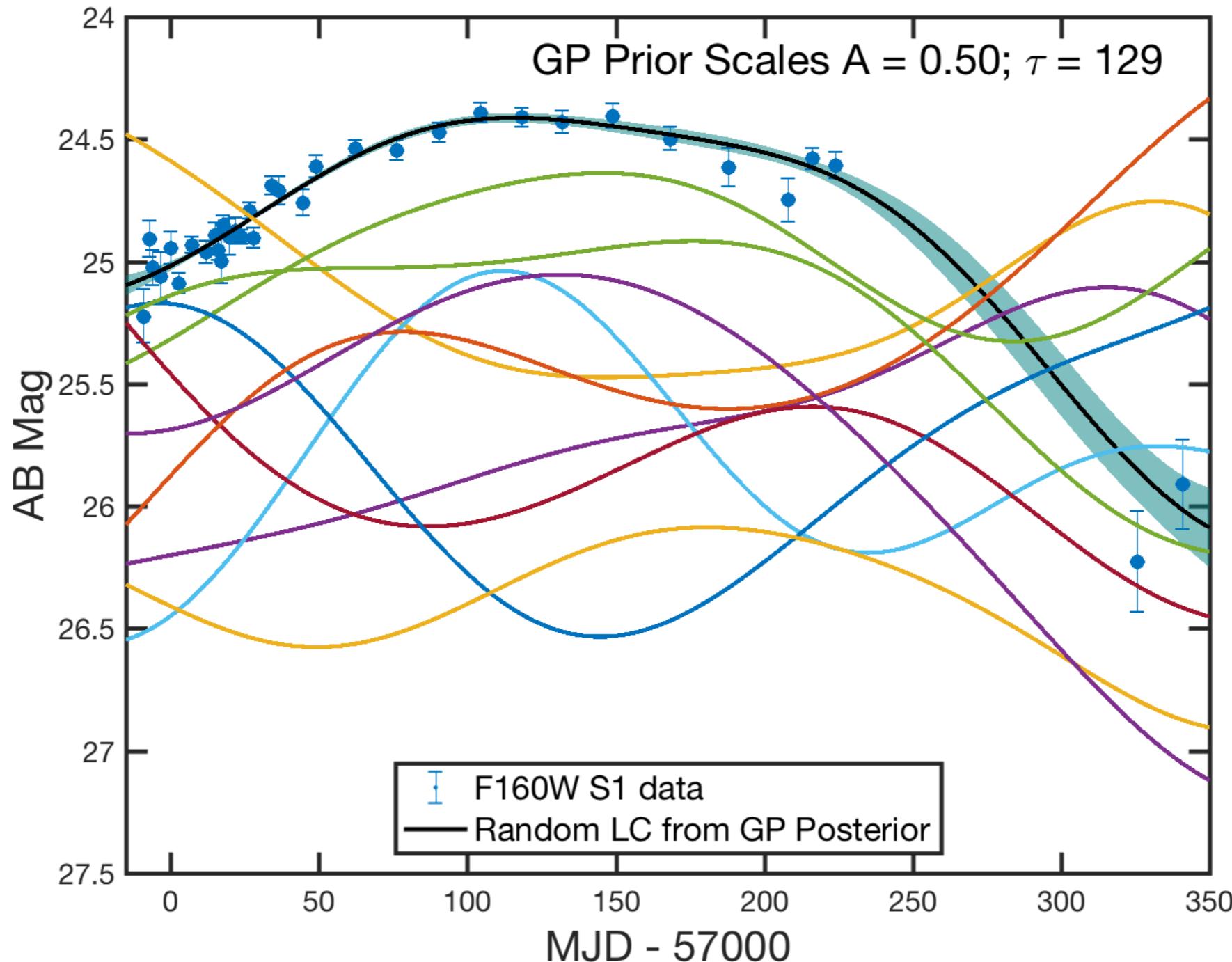
\mathbf{K}_{A, τ^2} = GP covariance

Which we can optimise (max likelihood)
or specify a prior on (A , τ) and sample from posterior

Demonstration

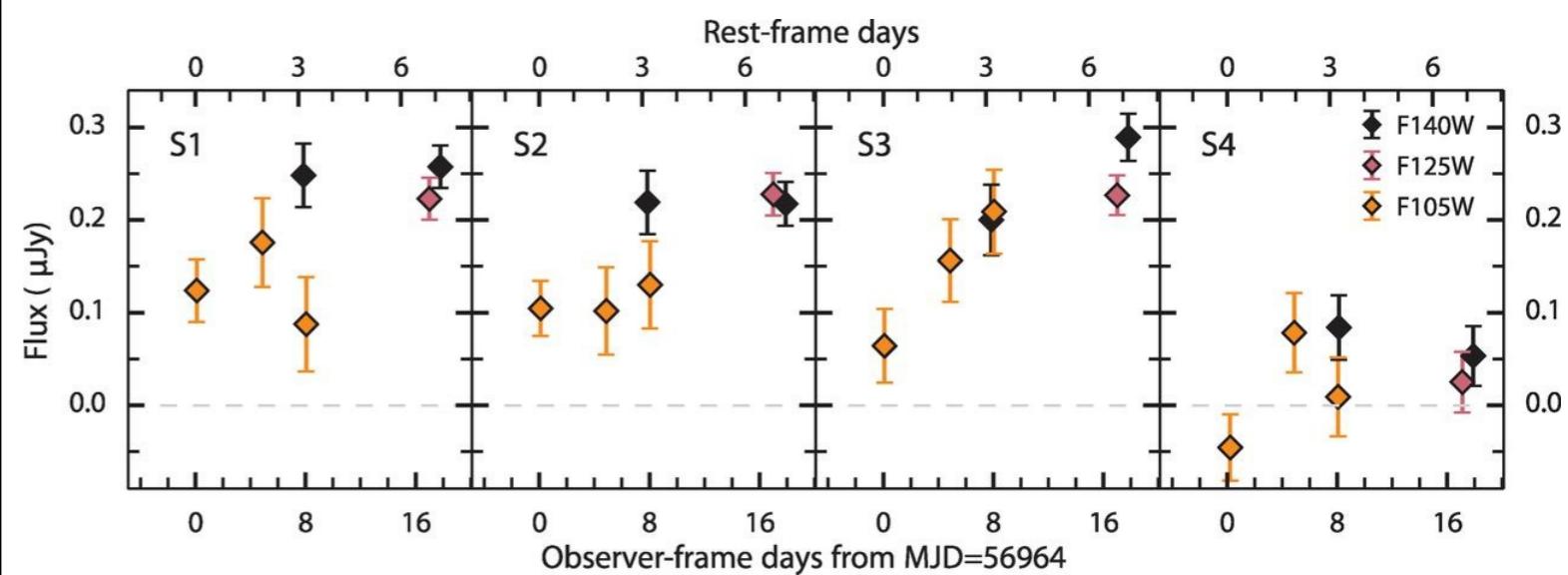
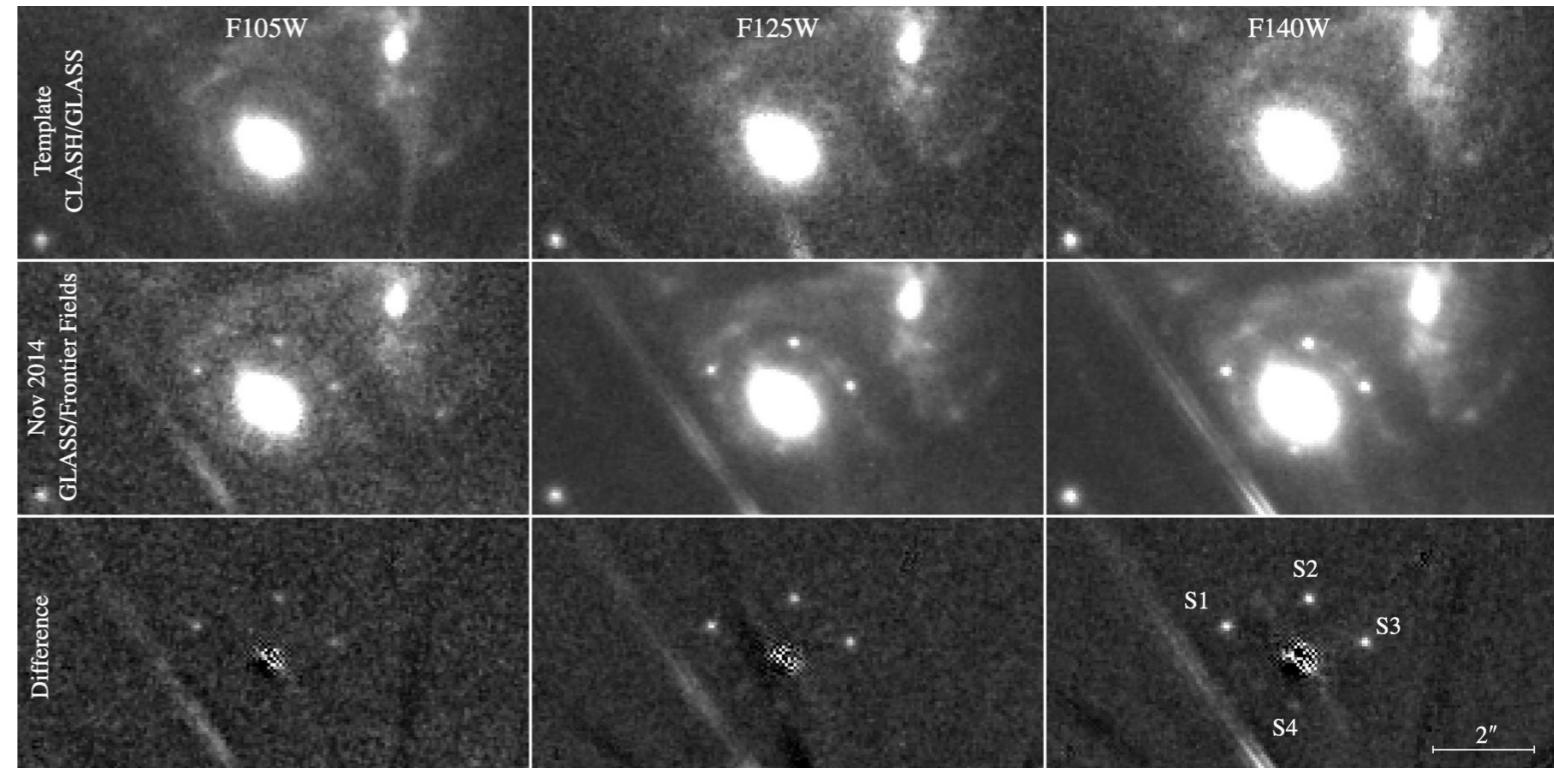
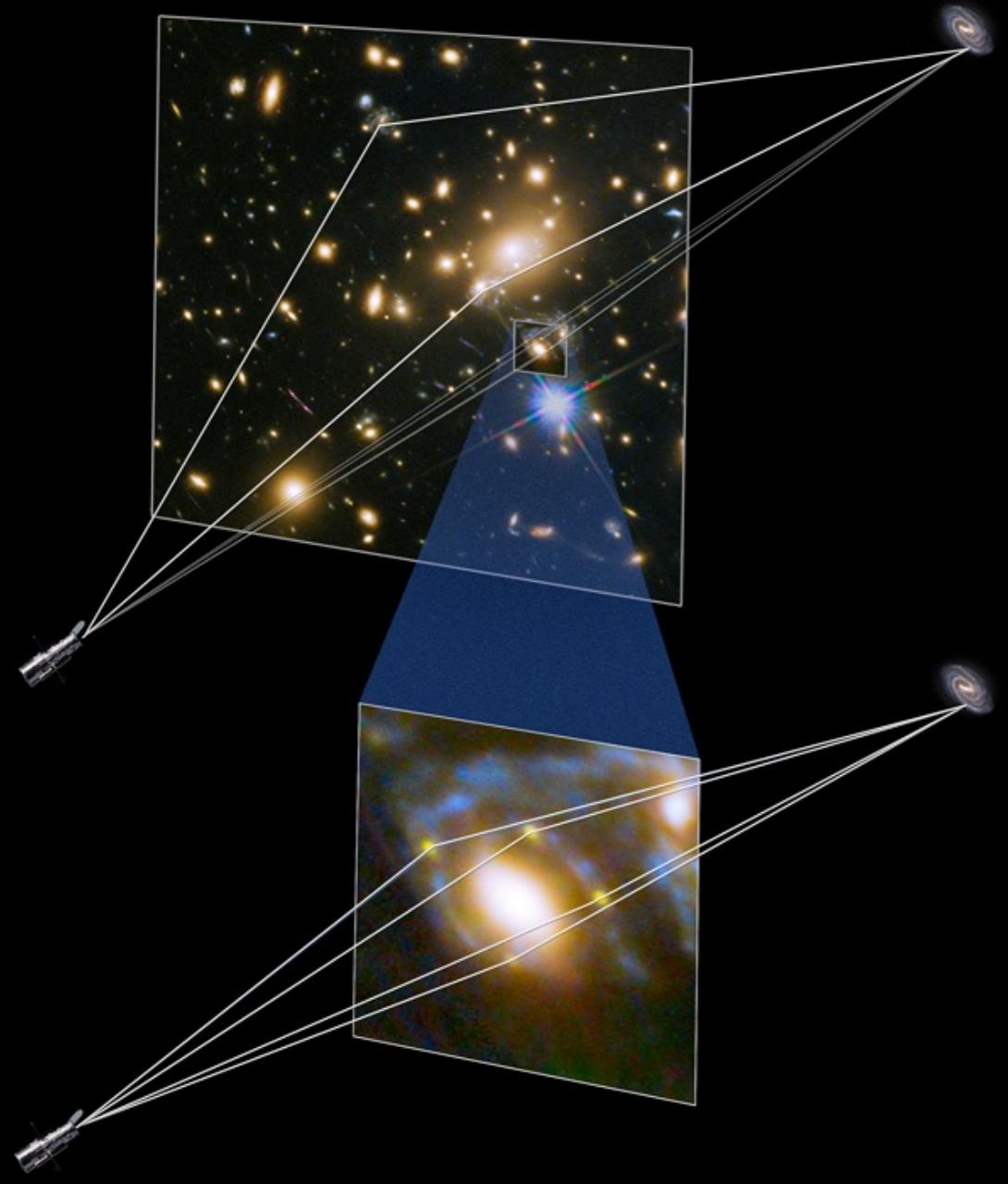


GP Fit with optimised hyperparameters



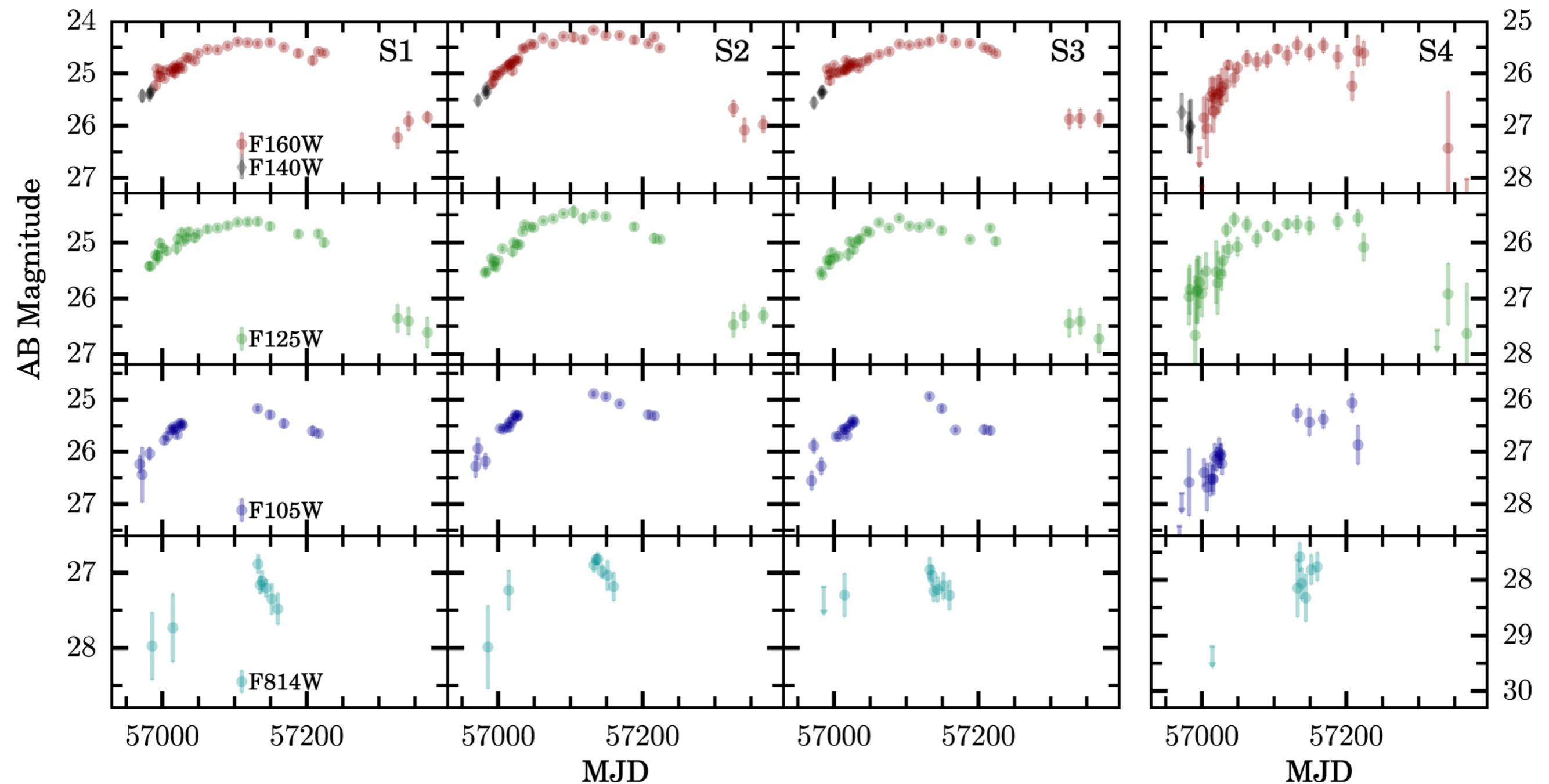
Case Study: SN Refsdal

Hubble Sees Distant Supernova
Multiply Imaged by Foreground Galaxy Cluster



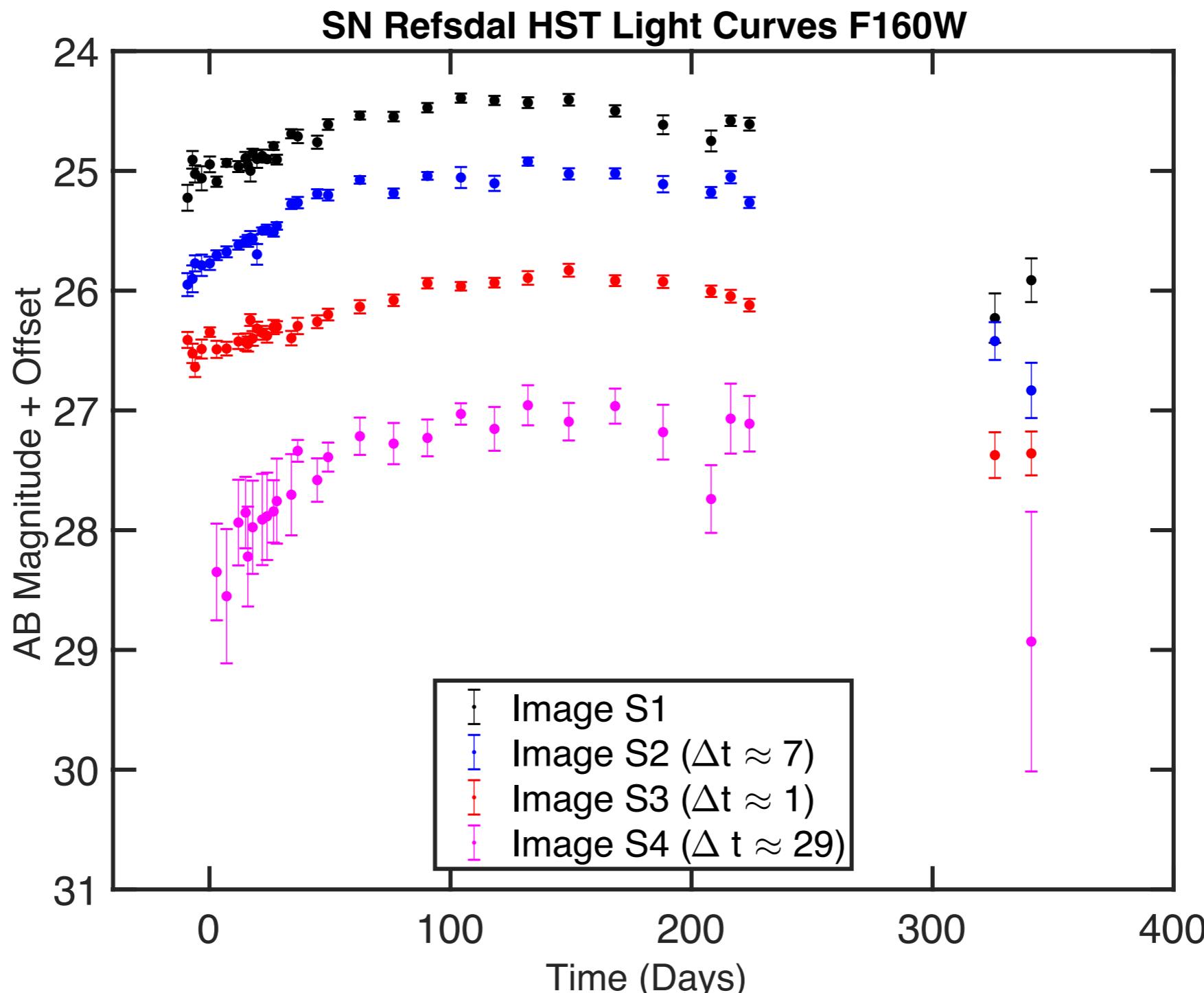
Time Series of SN brightnesses of each image: S1-S4

Hubble Space Telescope time series of SN Refsdal multiple images (Rodney et al. 2016)



Brightness Time Series [MJD = Modified Julian Day]

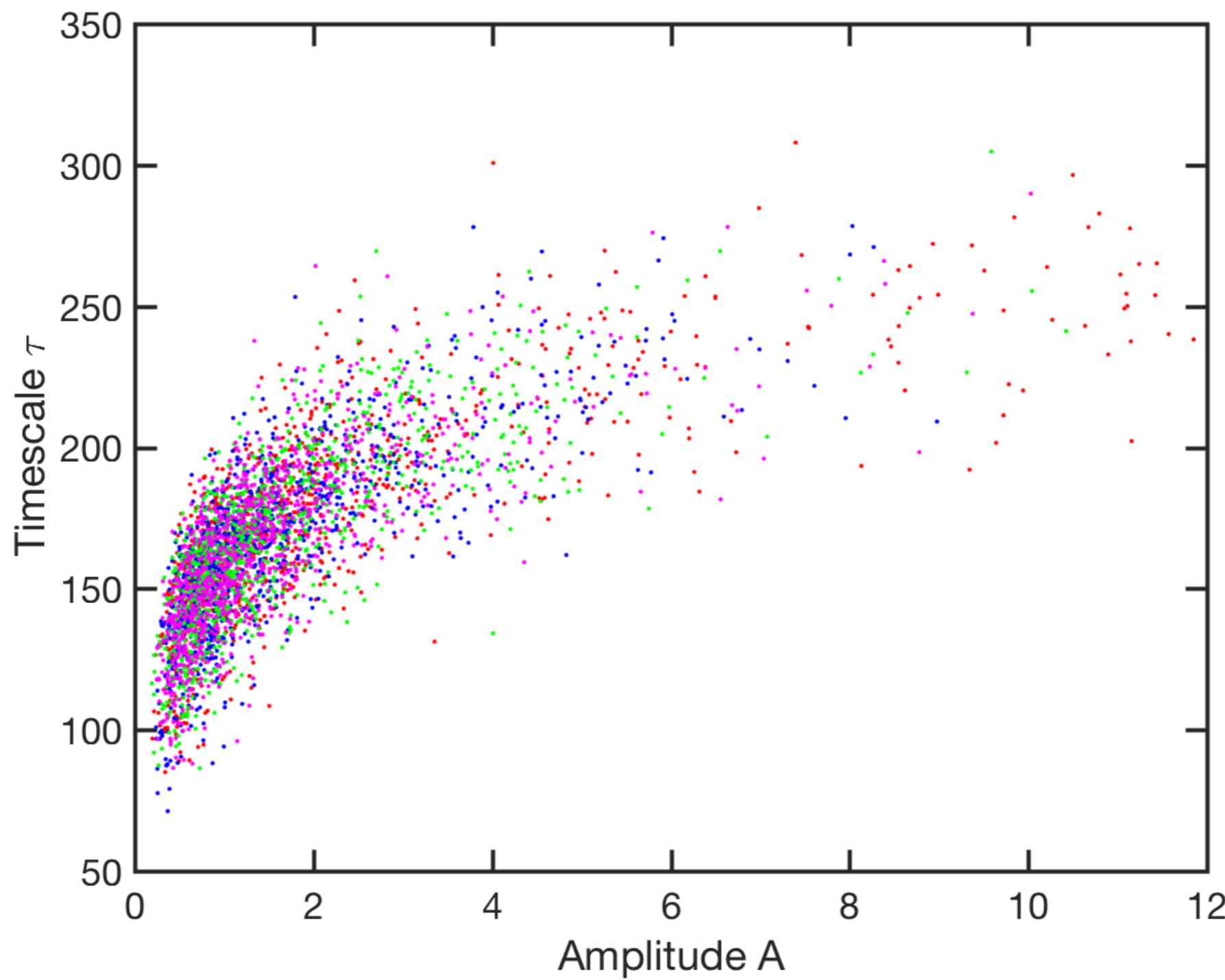
Hubble Space Telescope Time Series (light curves) of SN Refsdal at $\lambda \approx 1.6 \mu\text{m}$



Rodney et al. 2016: Photometry & Time Delay Measurements
of the first Einstein Cross Supernova

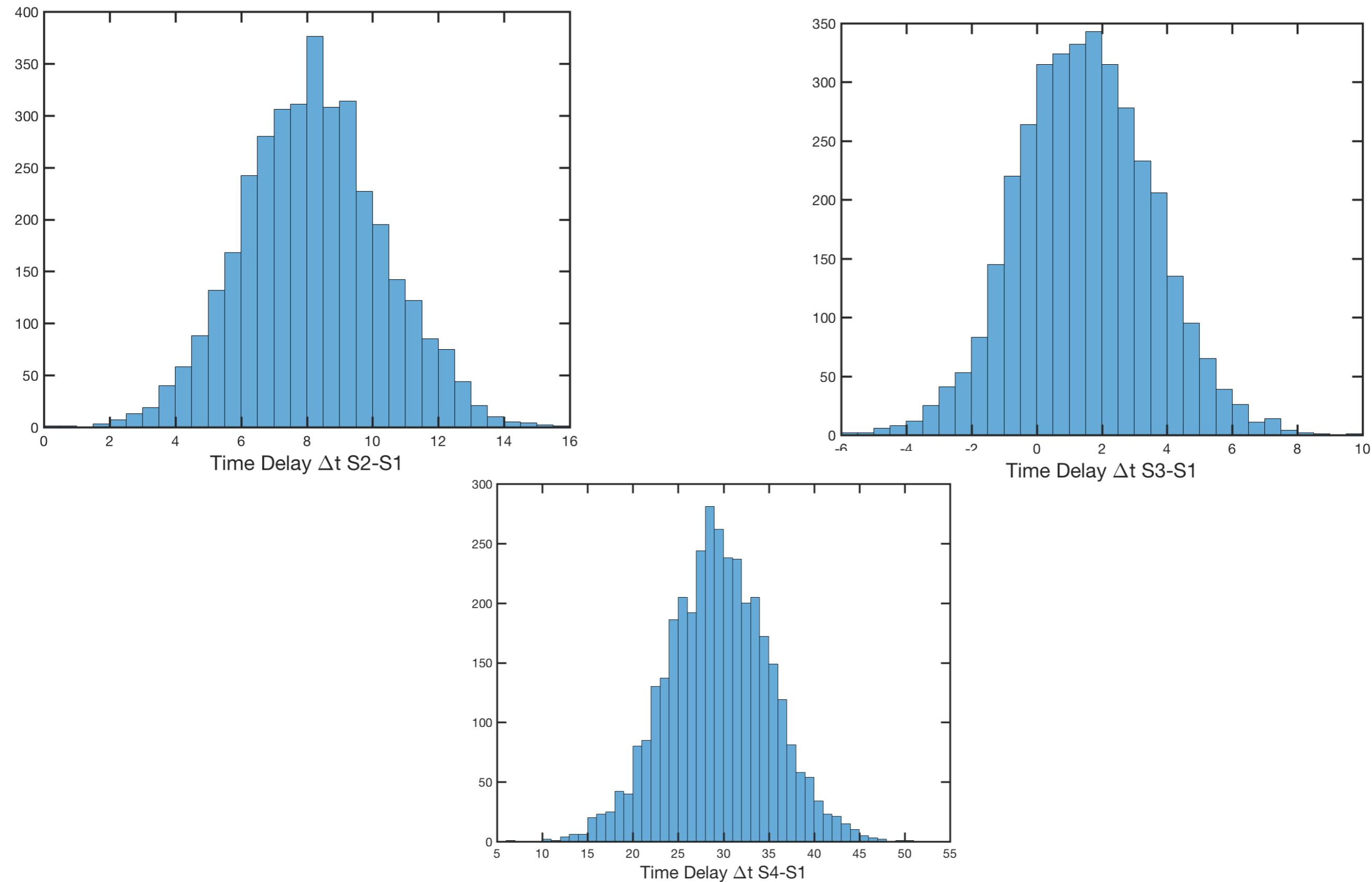
Write down Bayesian Model

Metropolis-within-Gibbs results: 8 parameters, 4 chains GP Hyperparameters



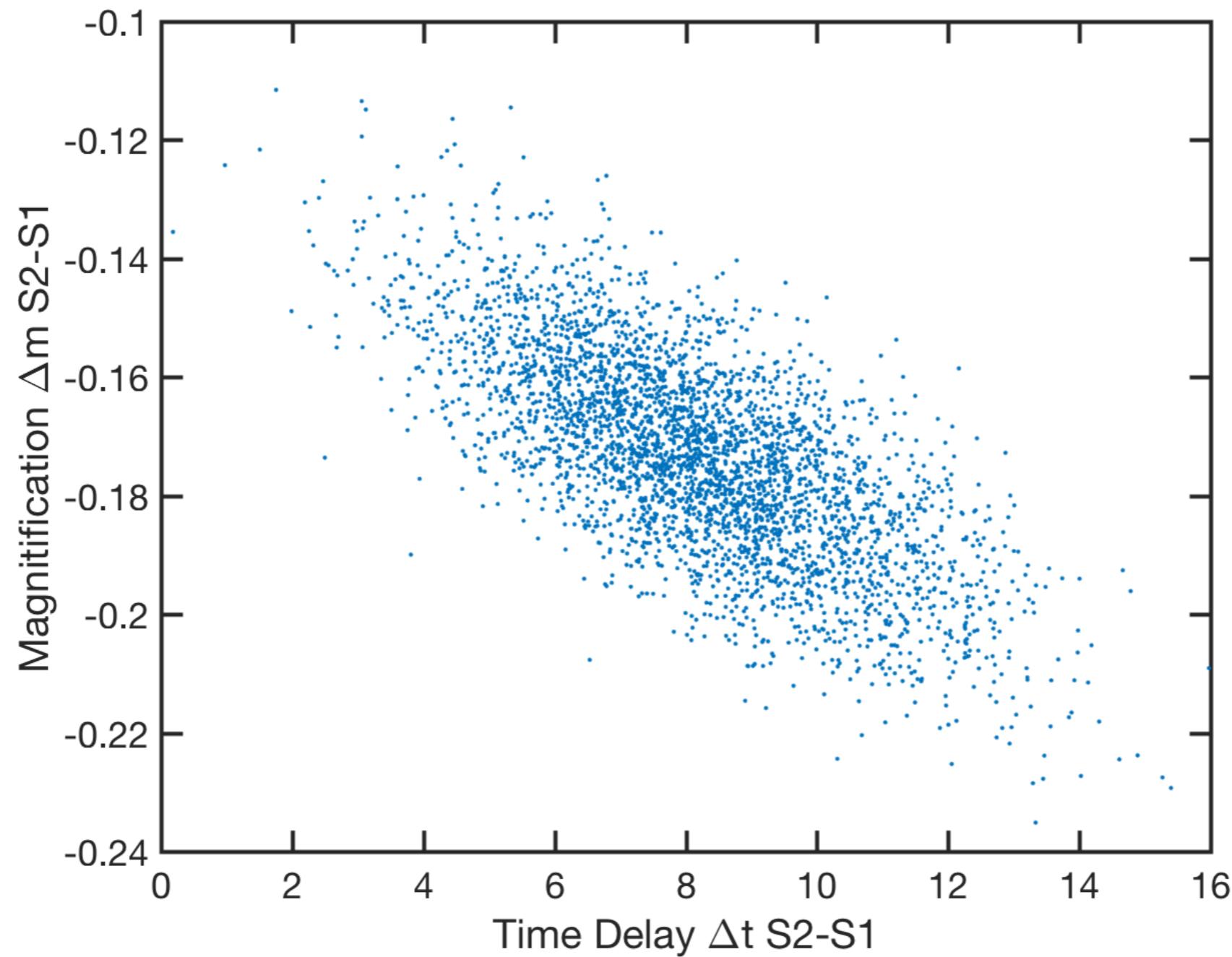
Metropolis-within-Gibbs results: 4 chains

Time Delays



Metropolis-within-Gibbs results: 4 chains

Time Delays vs Magnifications



Metropolis-within-Gibbs results: 4 chains

Deshifting the data

