

TIME SERIES ANALYSIS WITH GAUSSIAN PROCESSES AND THE SEARCH FOR EARTH 2.0

Dan Foreman-Mackey

CCPP@NYU // github.com/dfm // [@exoplaneteer](https://twitter.com/exoplaneteer) // dfm.io

I write **code** for good & **astrophysics**.



I study
space.

this isn't what
my data look like



I study
space.



I do **data science**.

not data
science.



Photo credit **James Silvester**
silvesterphoto.tumblr.com

CATERING PREP

data science.



Flickr user Marcin Wichary

data science.



my

GOALS

for today's talk

no convince you that

exoplanets are cool

EXOPLANETS

for today's talk

no convince you that

exoplanets are cool

demonstrate some

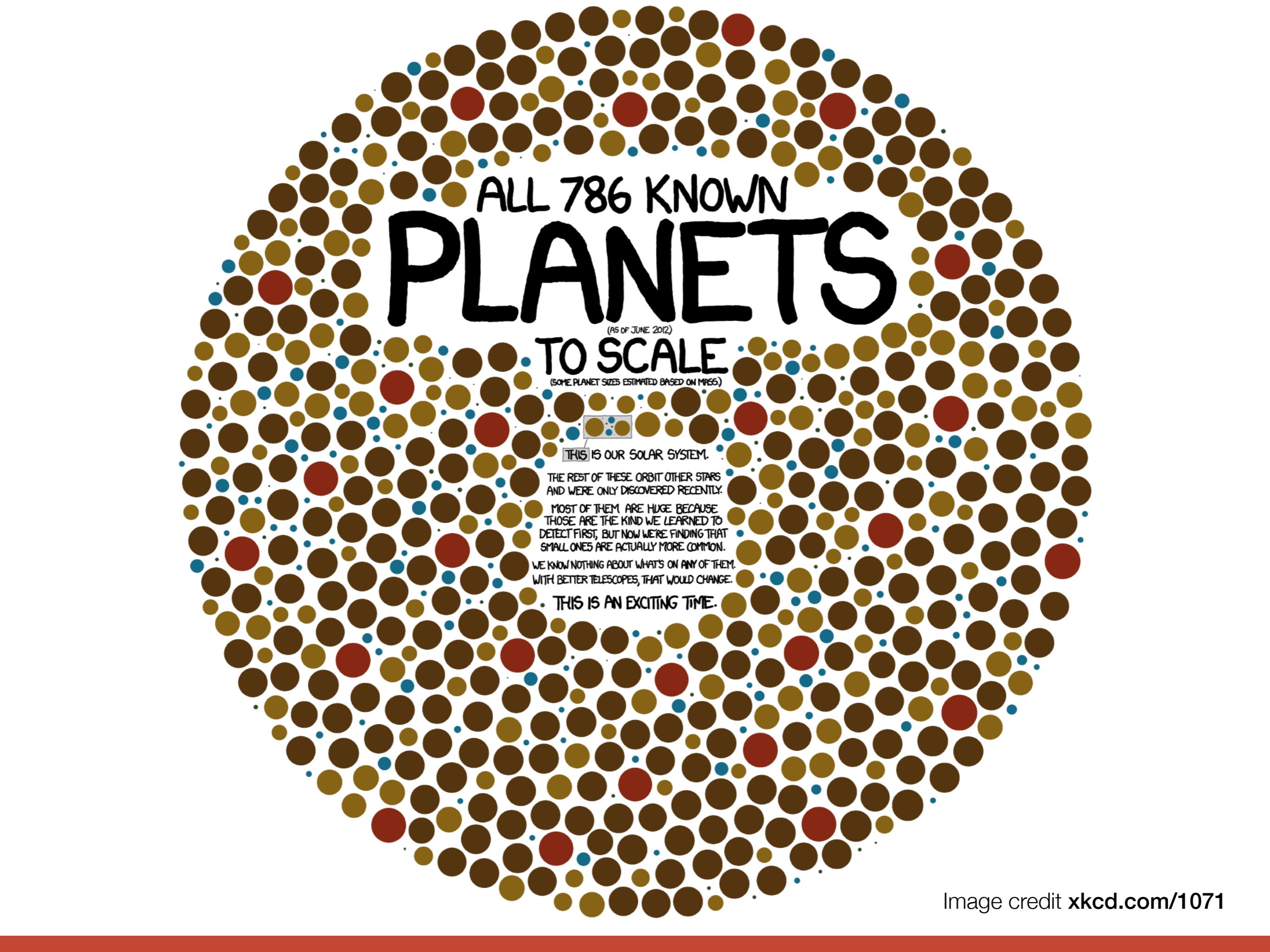
sick Python code

for today's talk

ex·o·plan·et

'eksō,planēt/

noun. a planet that orbits a star outside the solar system.



ALL 786 KNOWN PLANETS

(AS OF JUNE 2012)

TO SCALE

(SOME PLANET SIZES ESTIMATED BASED ON MASS.)

THIS IS OUR SOLAR SYSTEM.

THE REST OF THESE ORBIT OTHER STARS
AND WERE ONLY DISCOVERED RECENTLY.

MOST OF THEM ARE HUGE BECAUSE
THOSE ARE THE KIND WE LEARNED TO
DETECT FIRST, BUT NOW WE'RE FINDING THAT
SMALL ONES ARE ACTUALLY MORE COMMON.

WE KNOW NOTHING ABOUT WHAT'S ON ANY OF THEM.
WITH BETTER TELESCOPES, THAT WOULD CHANGE.

• THIS IS AN EXCITING TIME.

OUR NEIGHBORHOOD

A PORTRAIT OF ALL HABITABLE-ZONE PLANETS WITHIN 60 LIGHT-YEARS OF EARTH
(CONSTRUCTED FROM STATISTICAL DATA ON TYPICAL PLANET SIZES AND ORBITS)

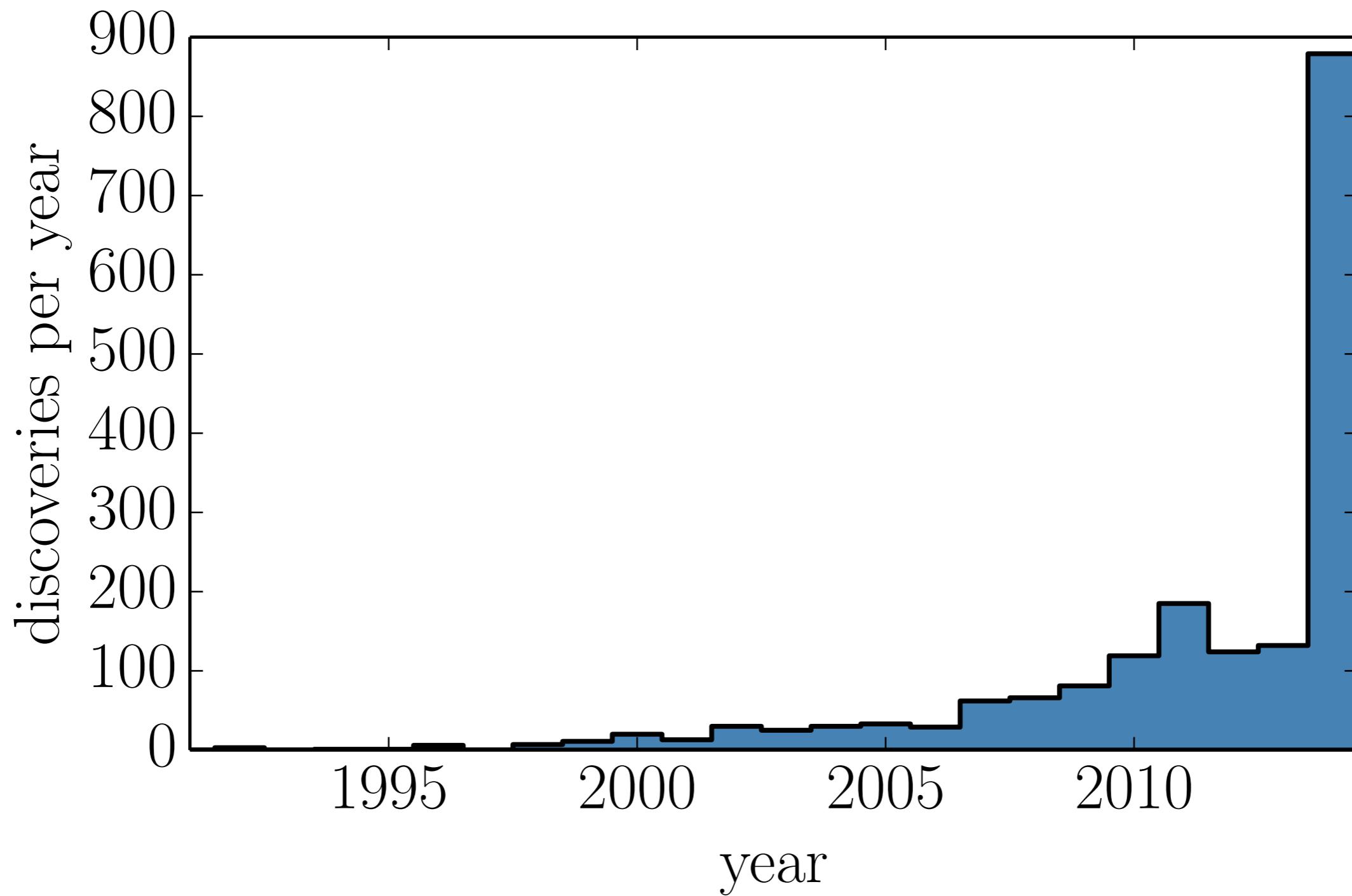
PLANETS AROUND SUN-LIKE STARS



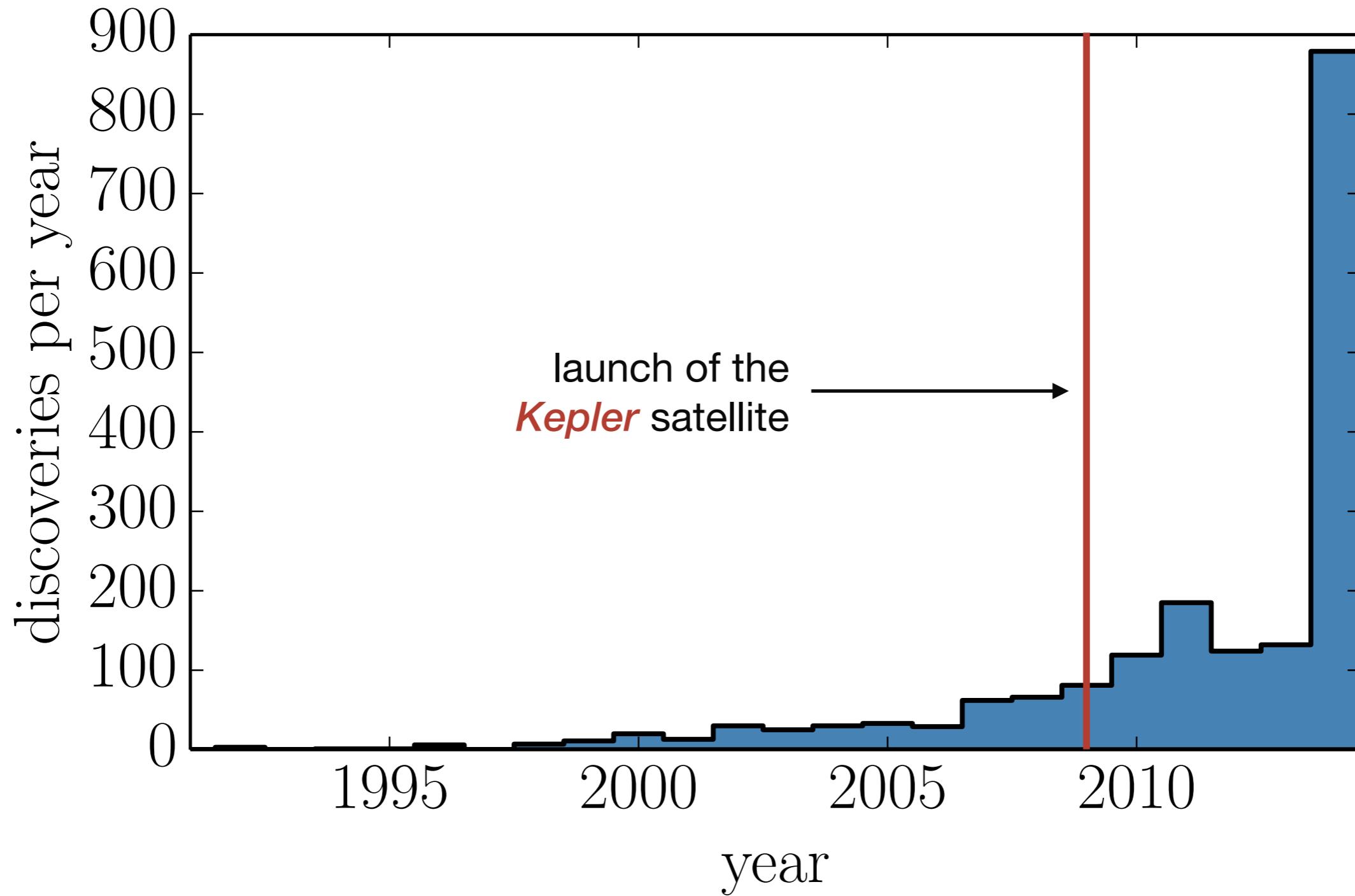
PLANETS AROUND OTHER STARS

EARTH-SIZED PLANETS

EARTH



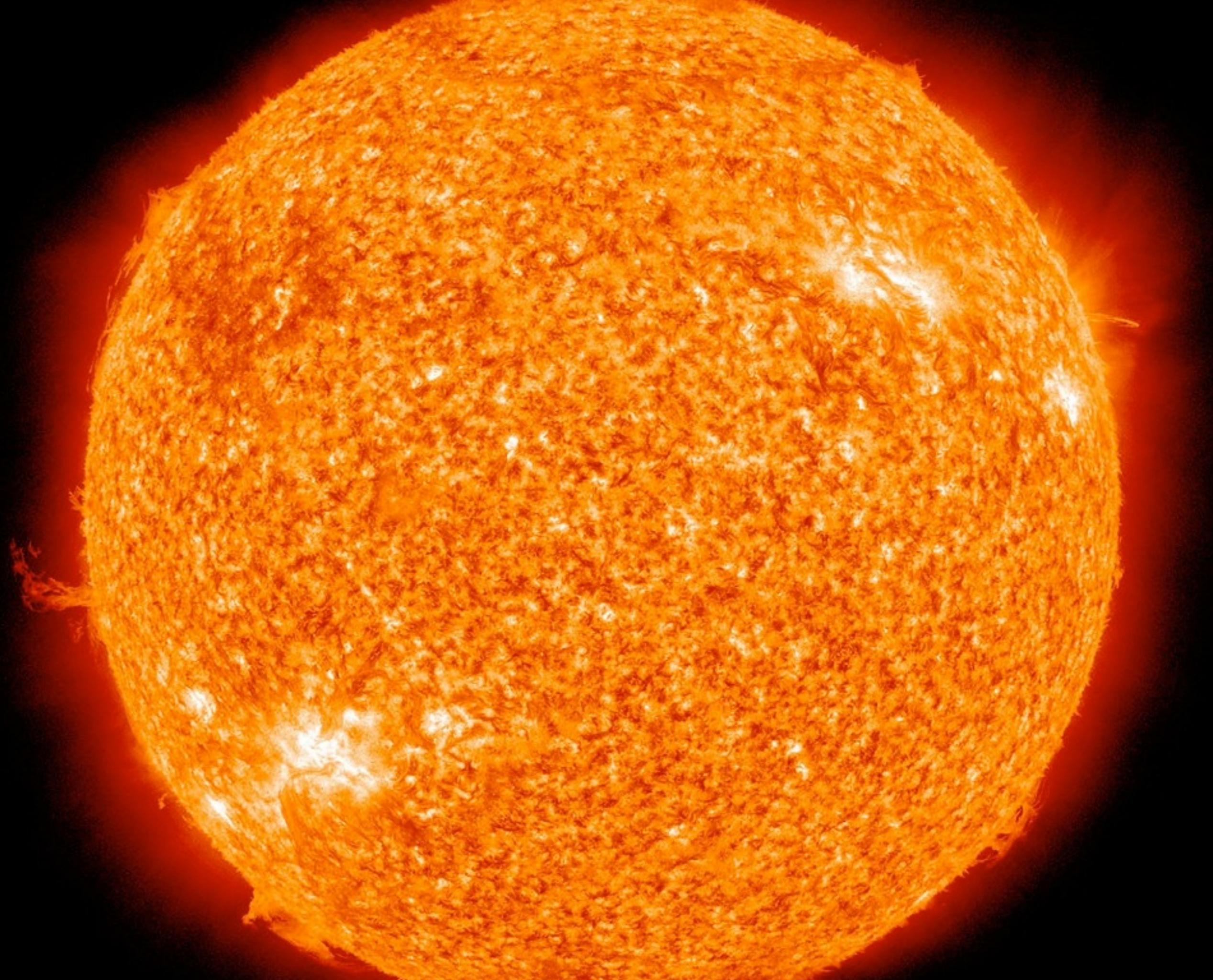
Data from **The Open Exoplanet Catalogue**



Data from [The Open Exoplanet Catalogue](#)

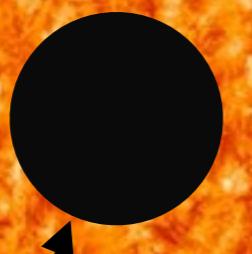
Exoplanets are **hard** to find

The **transit** method

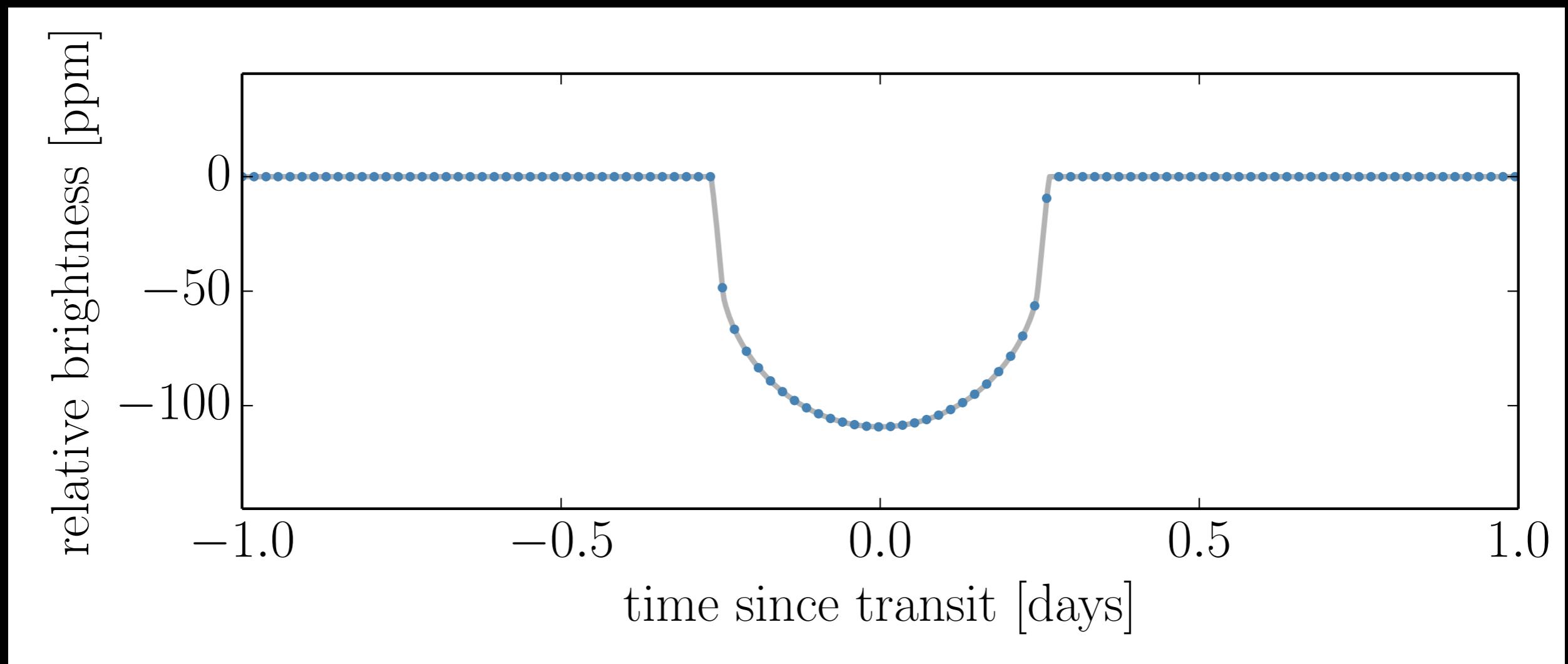
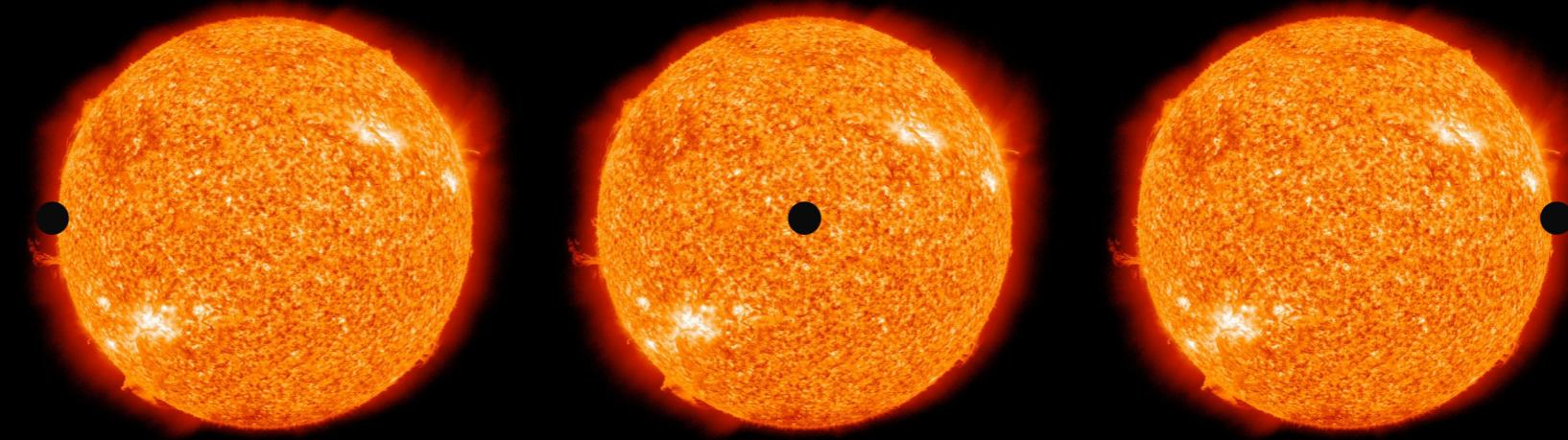




Earth



Jupiter



need to look at **the right place**
at **the right time**

and measure brightness
extremely precisely

Kepler



Image credit **NASA**



Photo credit: NASA



Image credit **Carter Roberts**

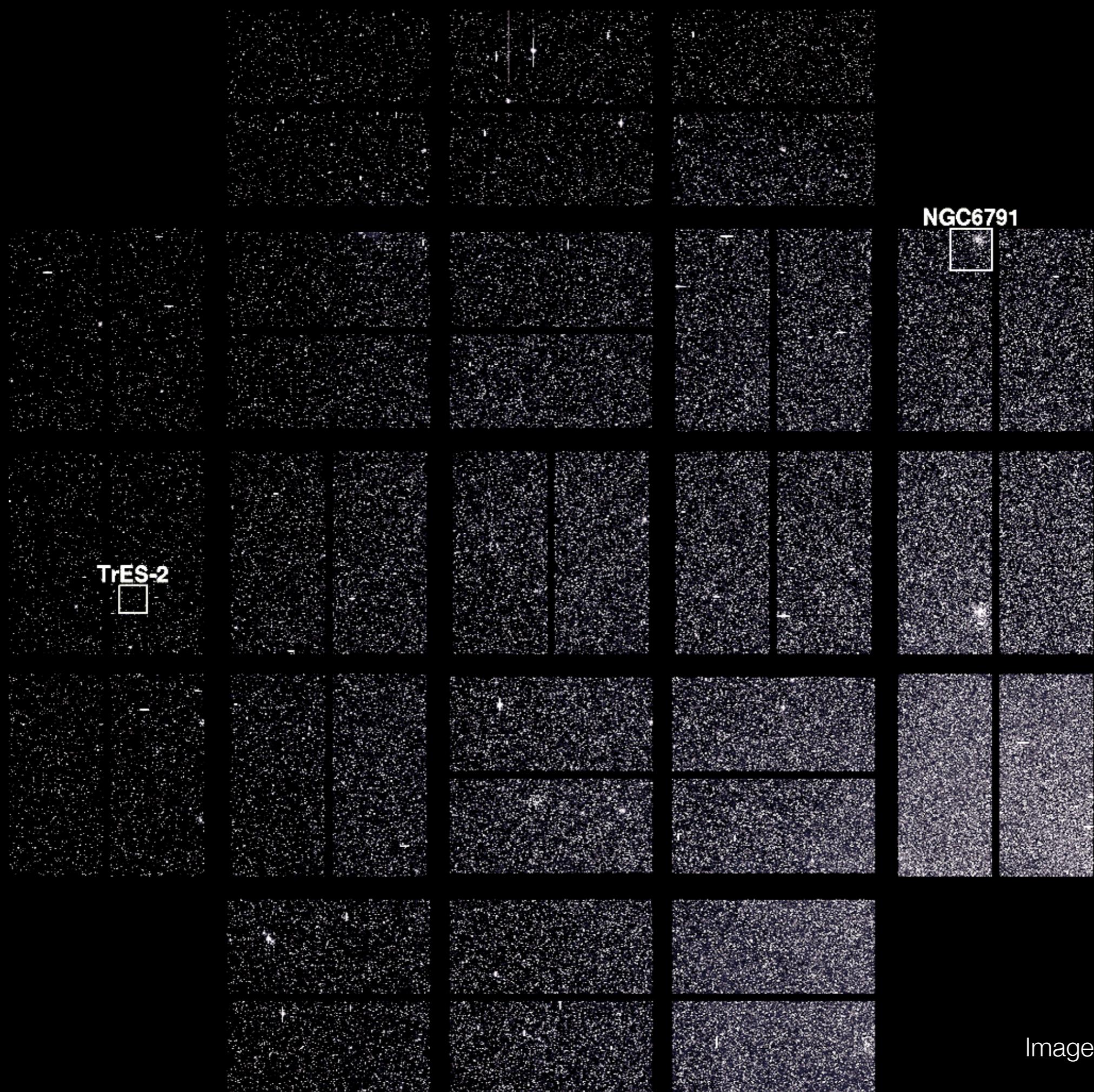
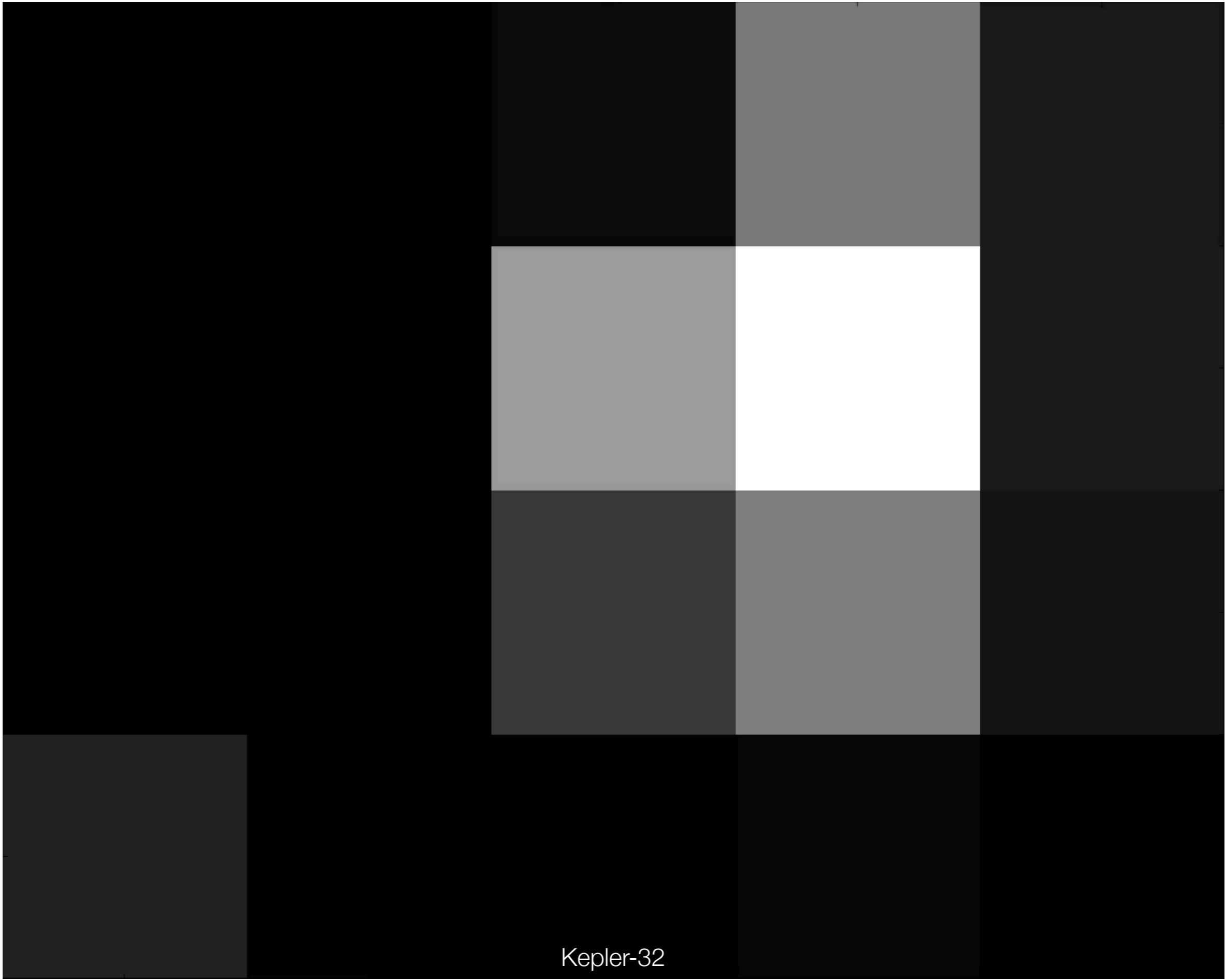
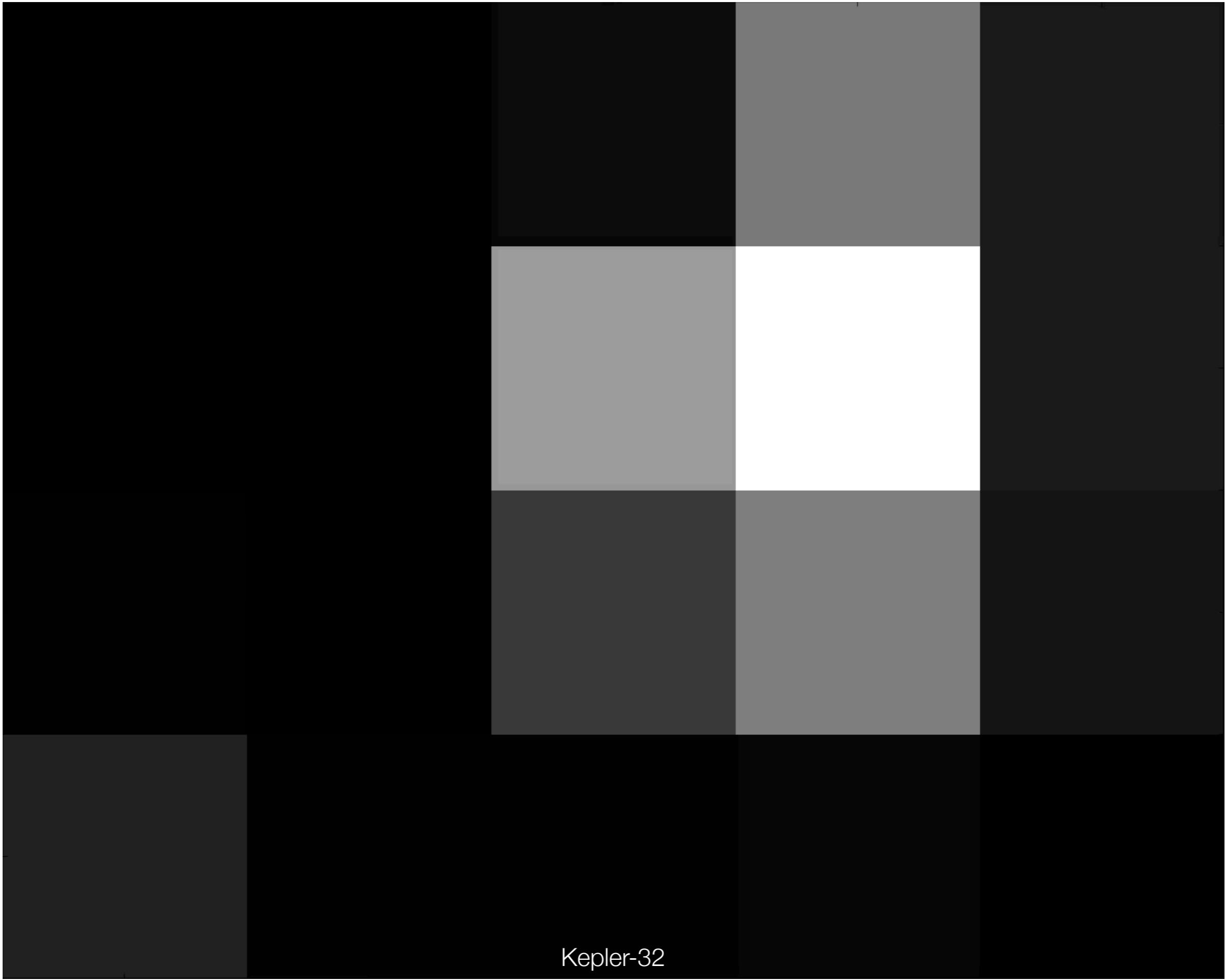


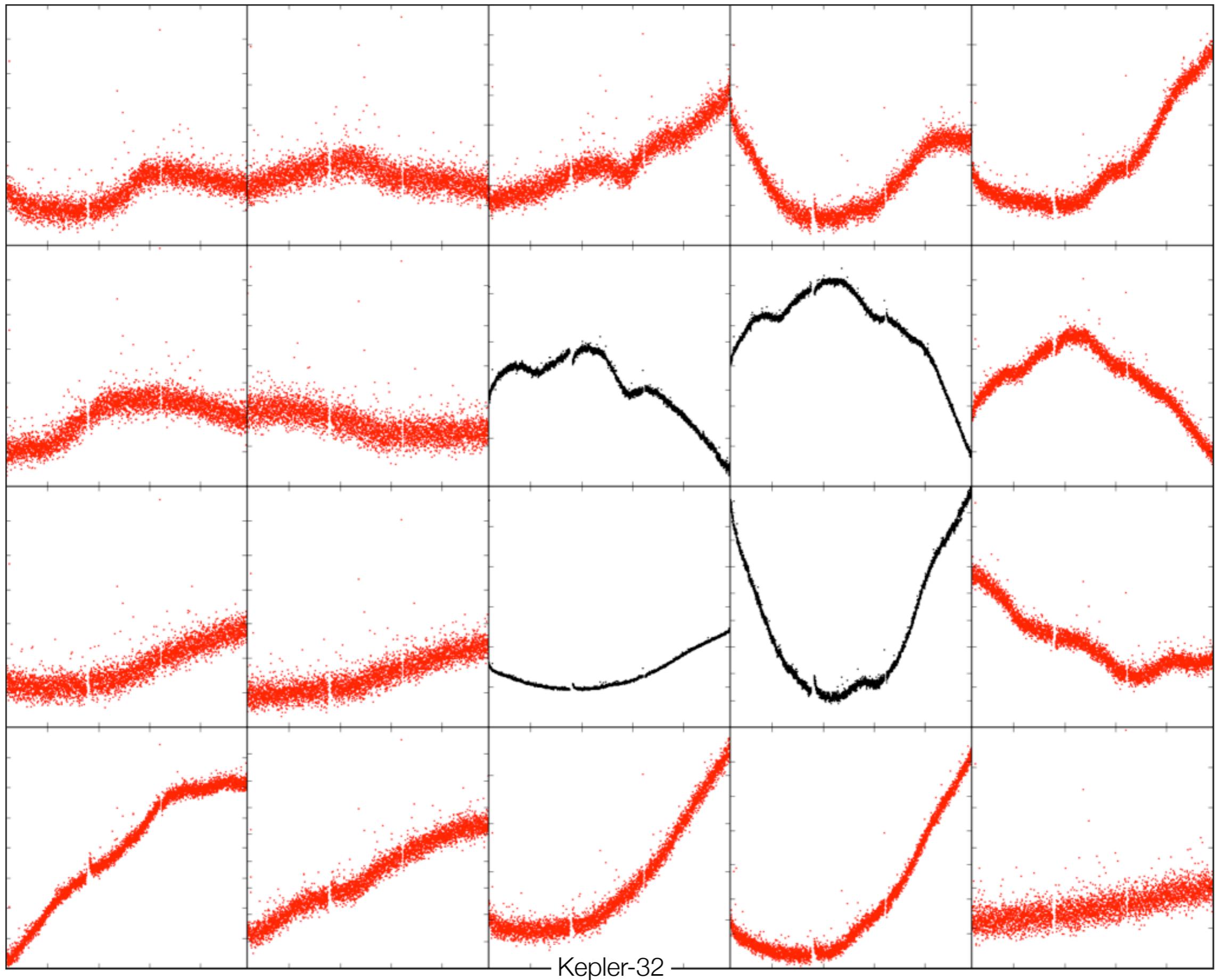
Image credit **NASA**

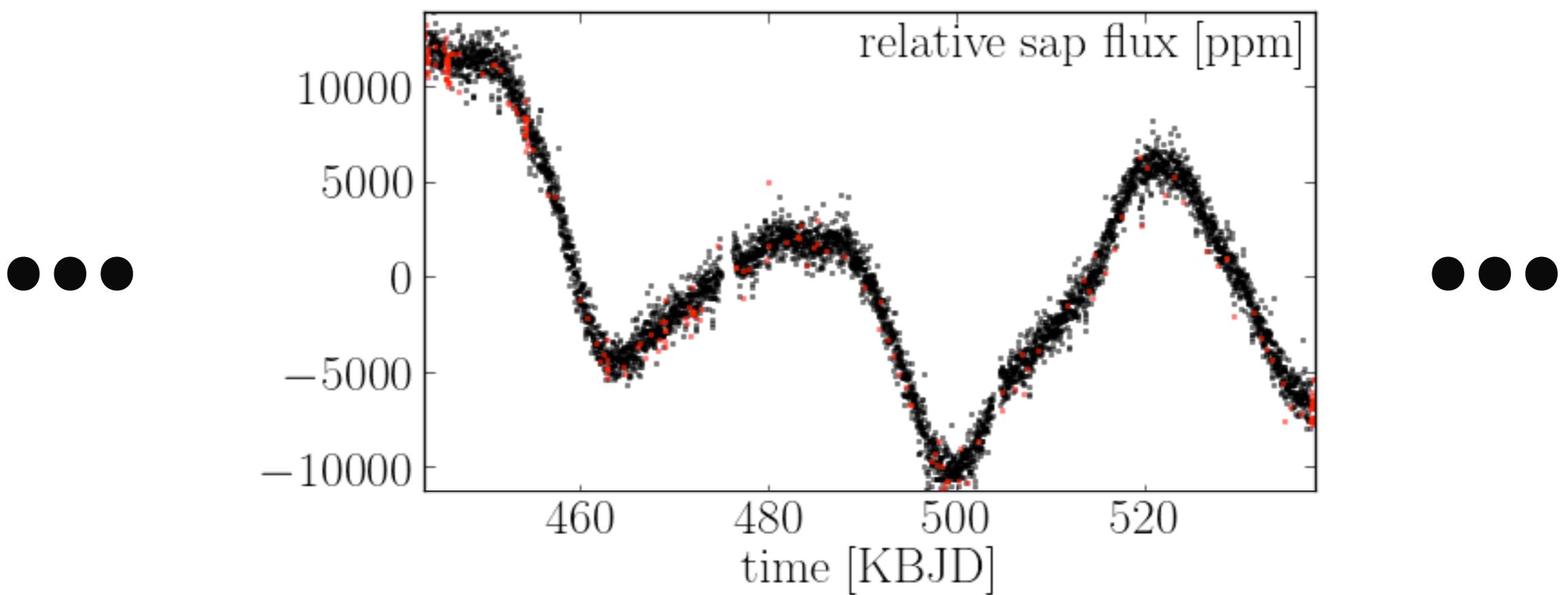


Kepler-32

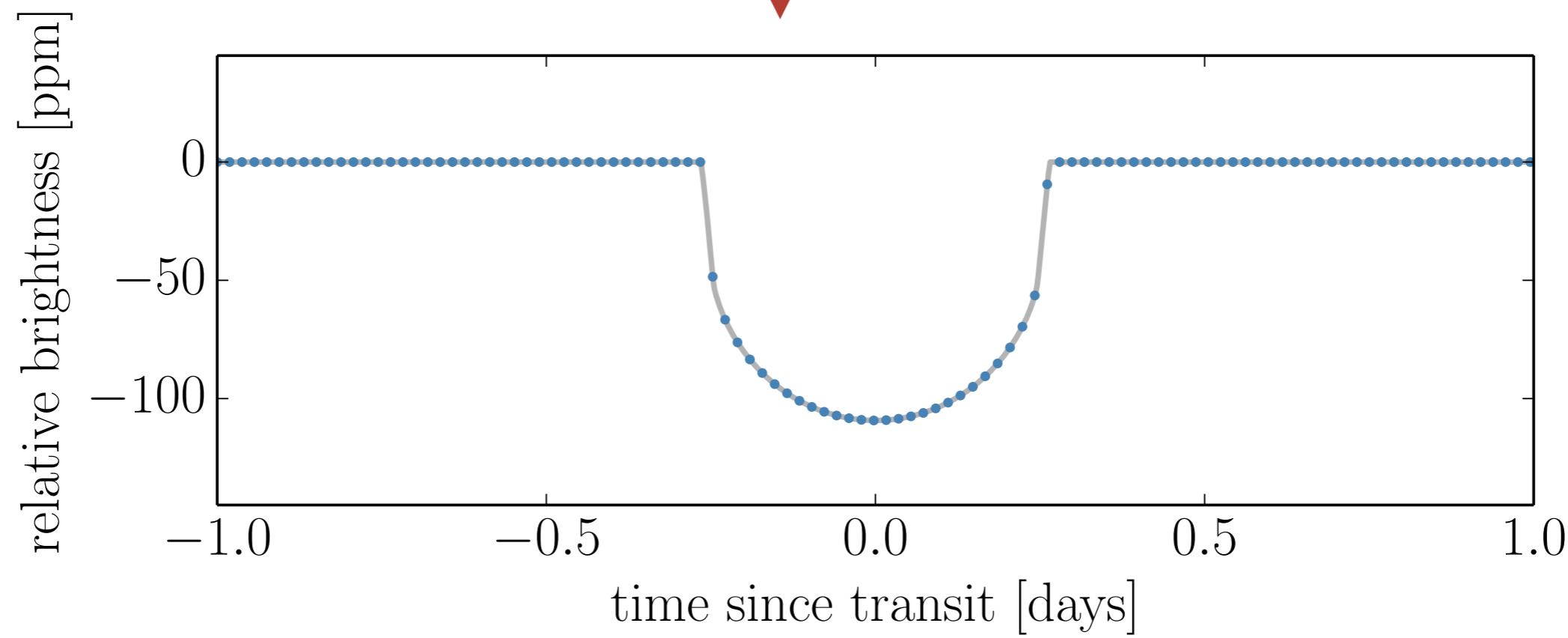
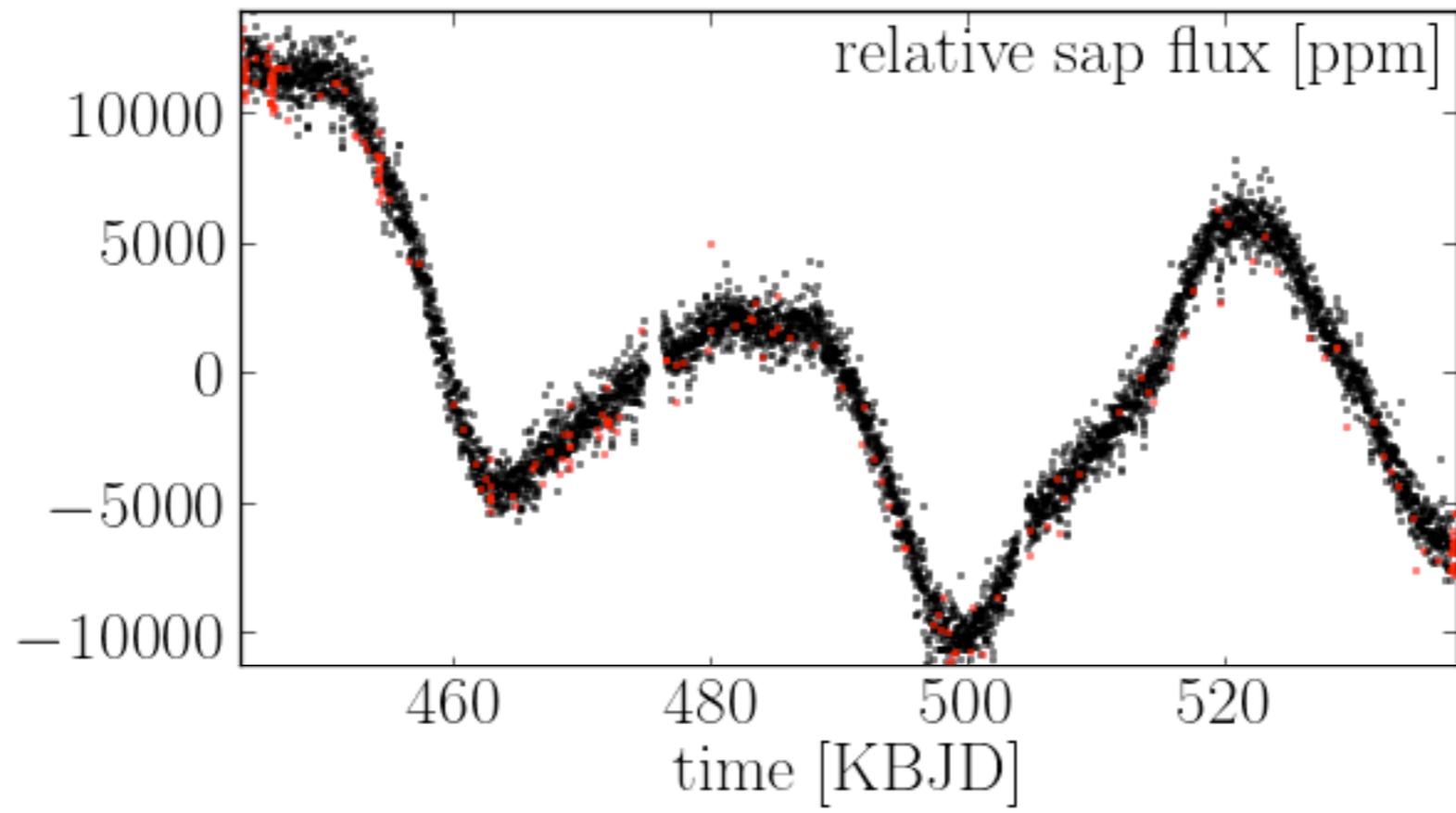


Kepler-32

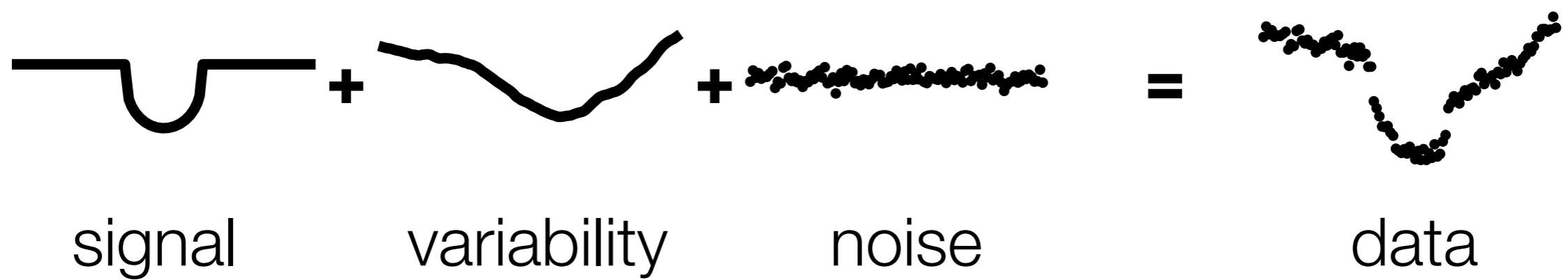




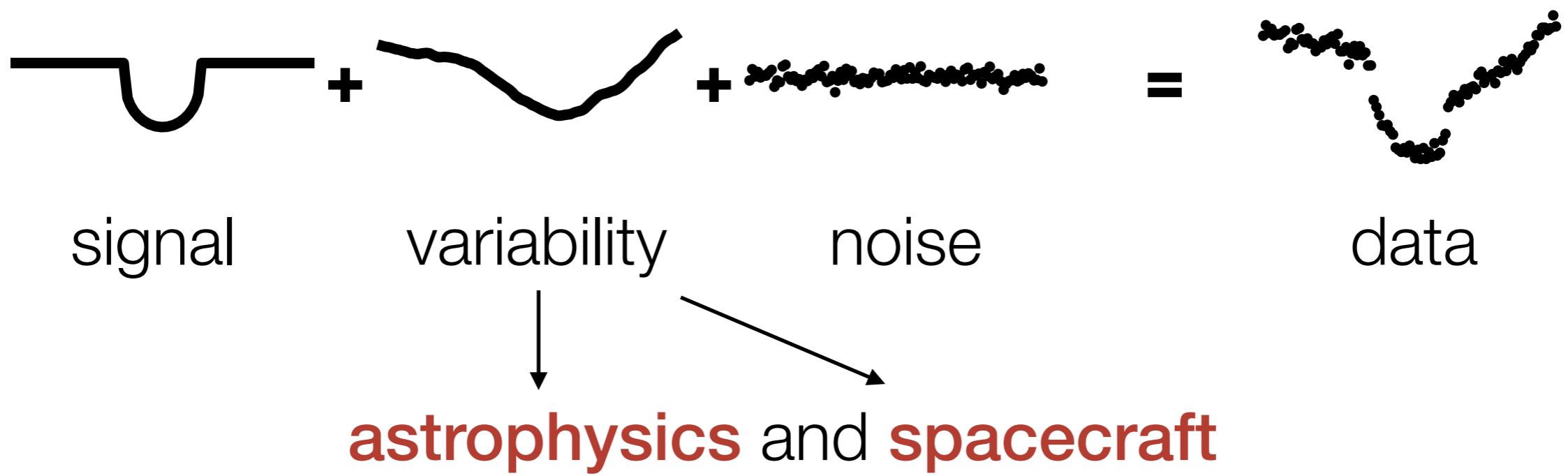
Kepler-32



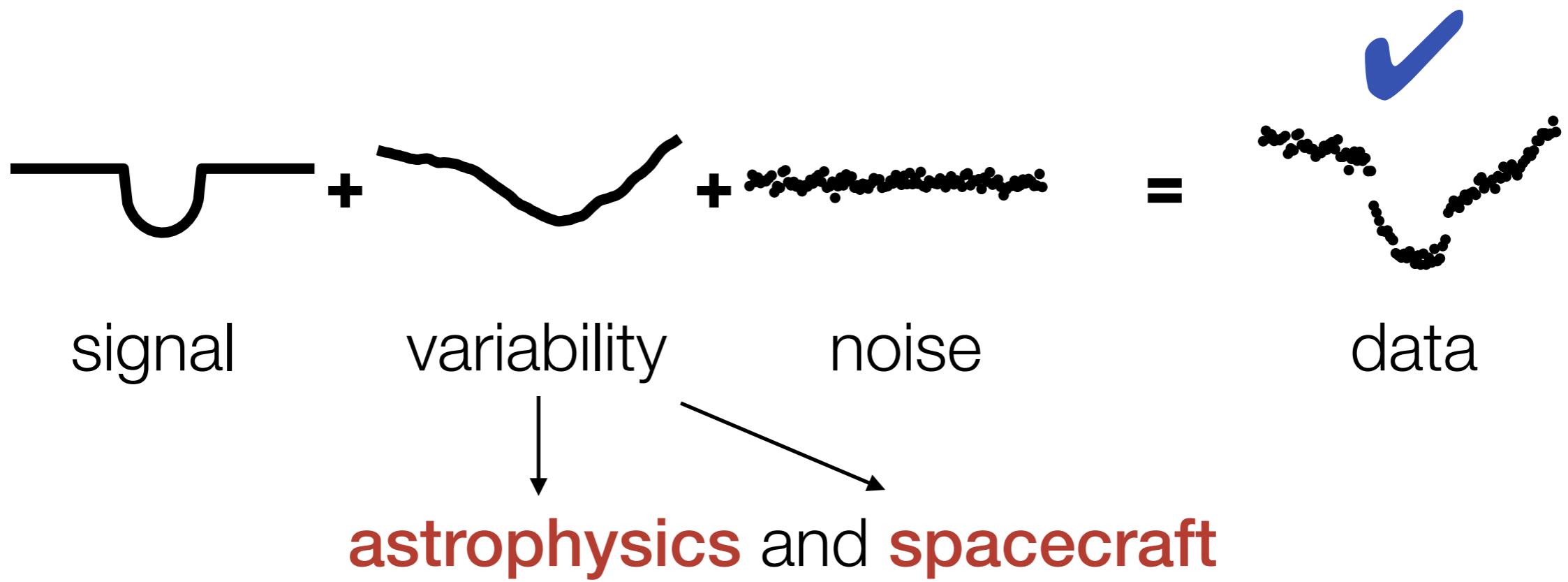
The **anatomy** of a **transit** observation



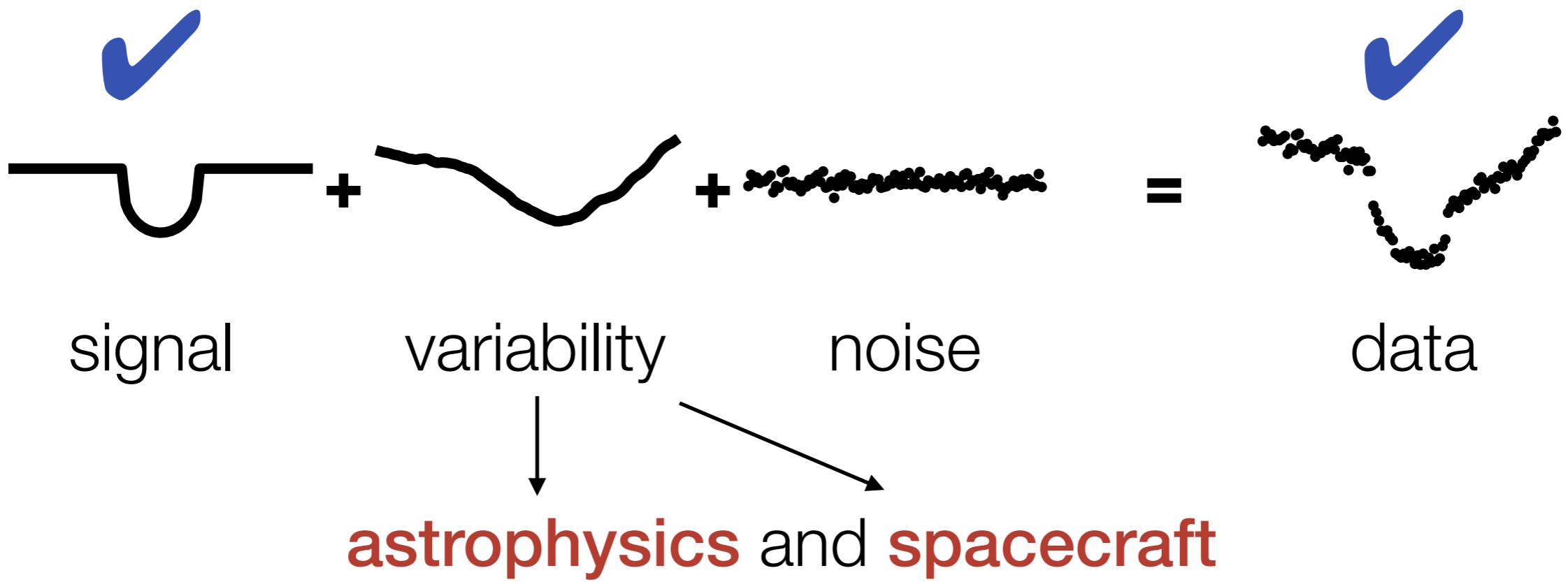
The **anatomy** of a **transit** observation



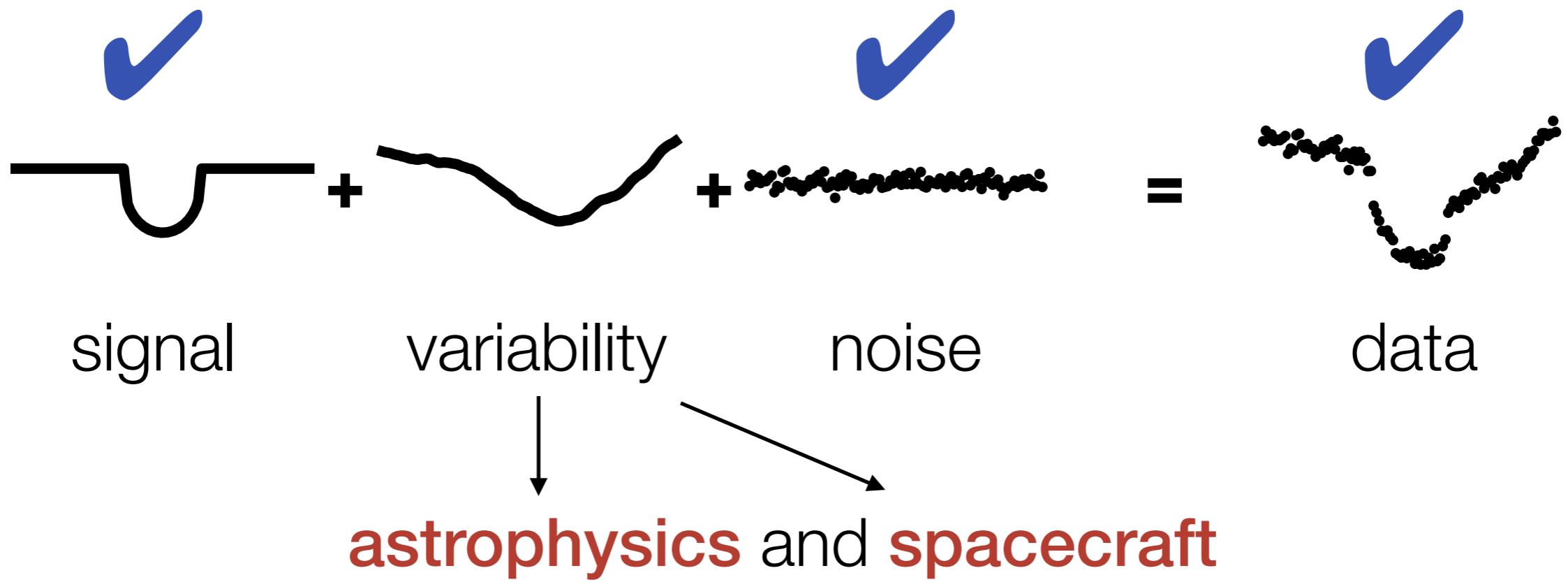
The **anatomy** of a **transit** observation



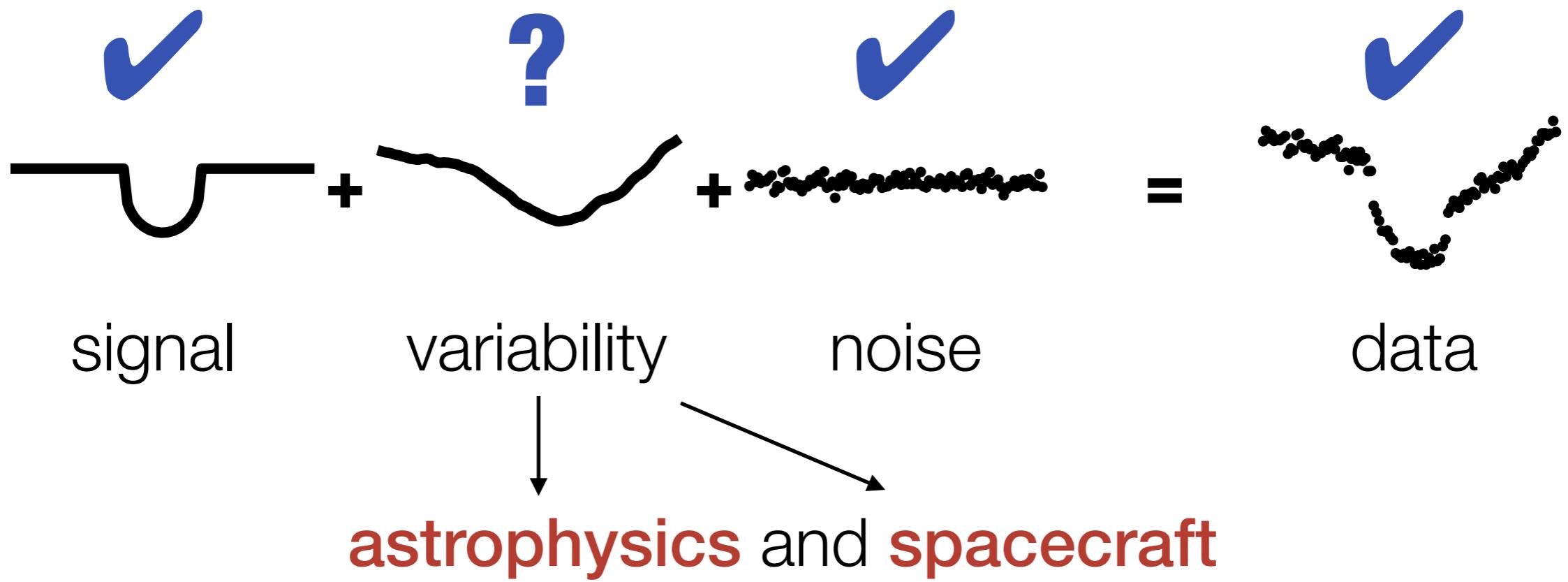
The **anatomy** of a **transit** observation



The **anatomy** of a **transit** observation



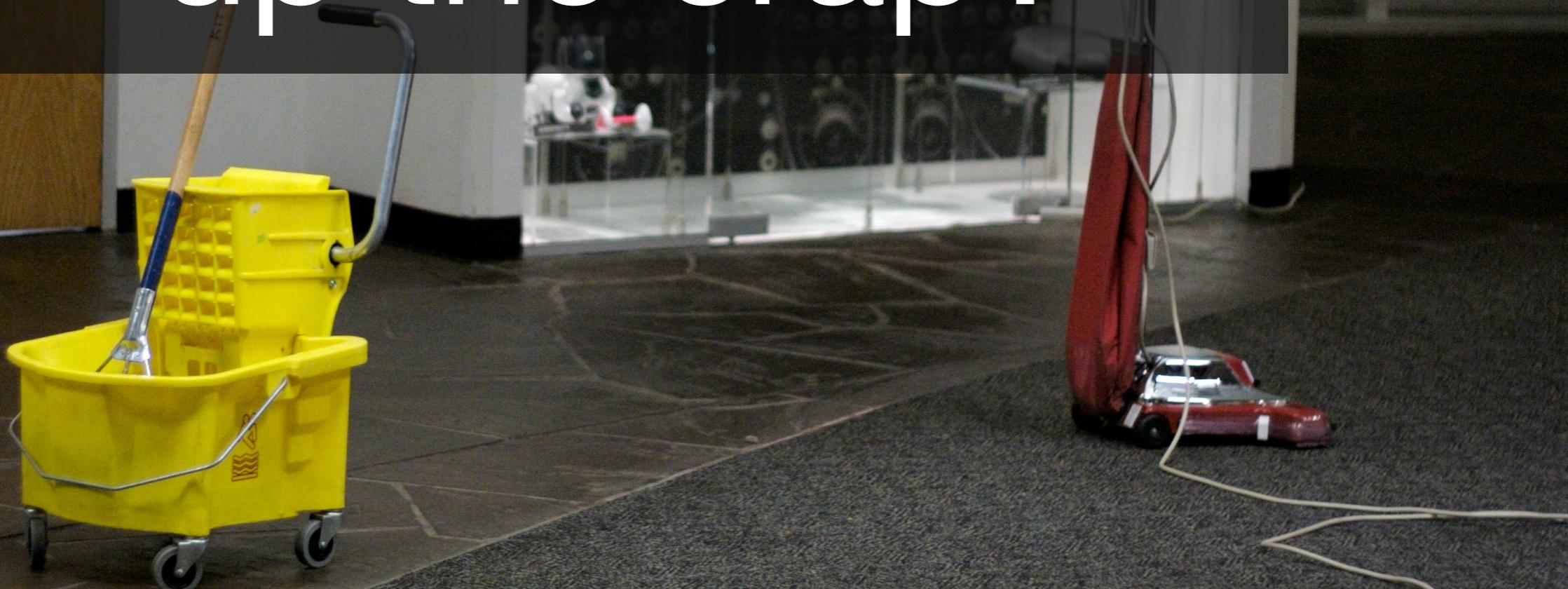
The **anatomy** of a **transit** observation



CATERING PREP

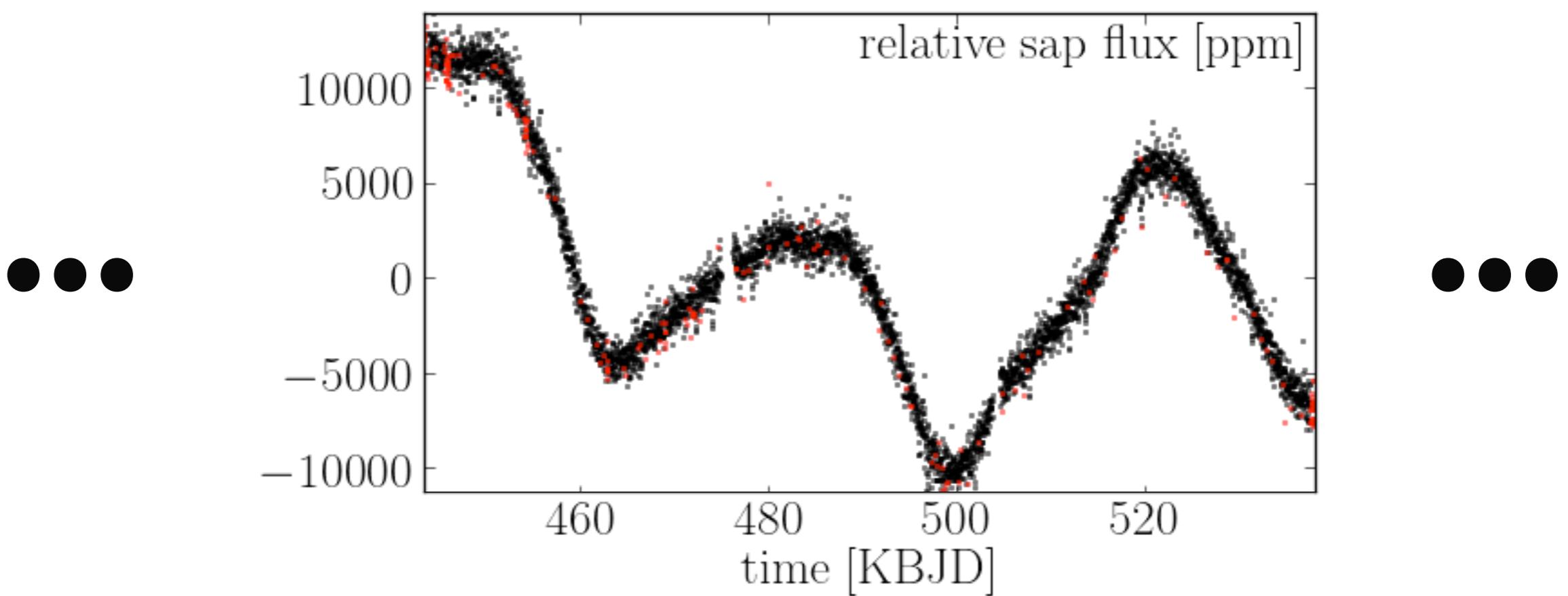
Kitchen

so.
how do we clean
up the crap?

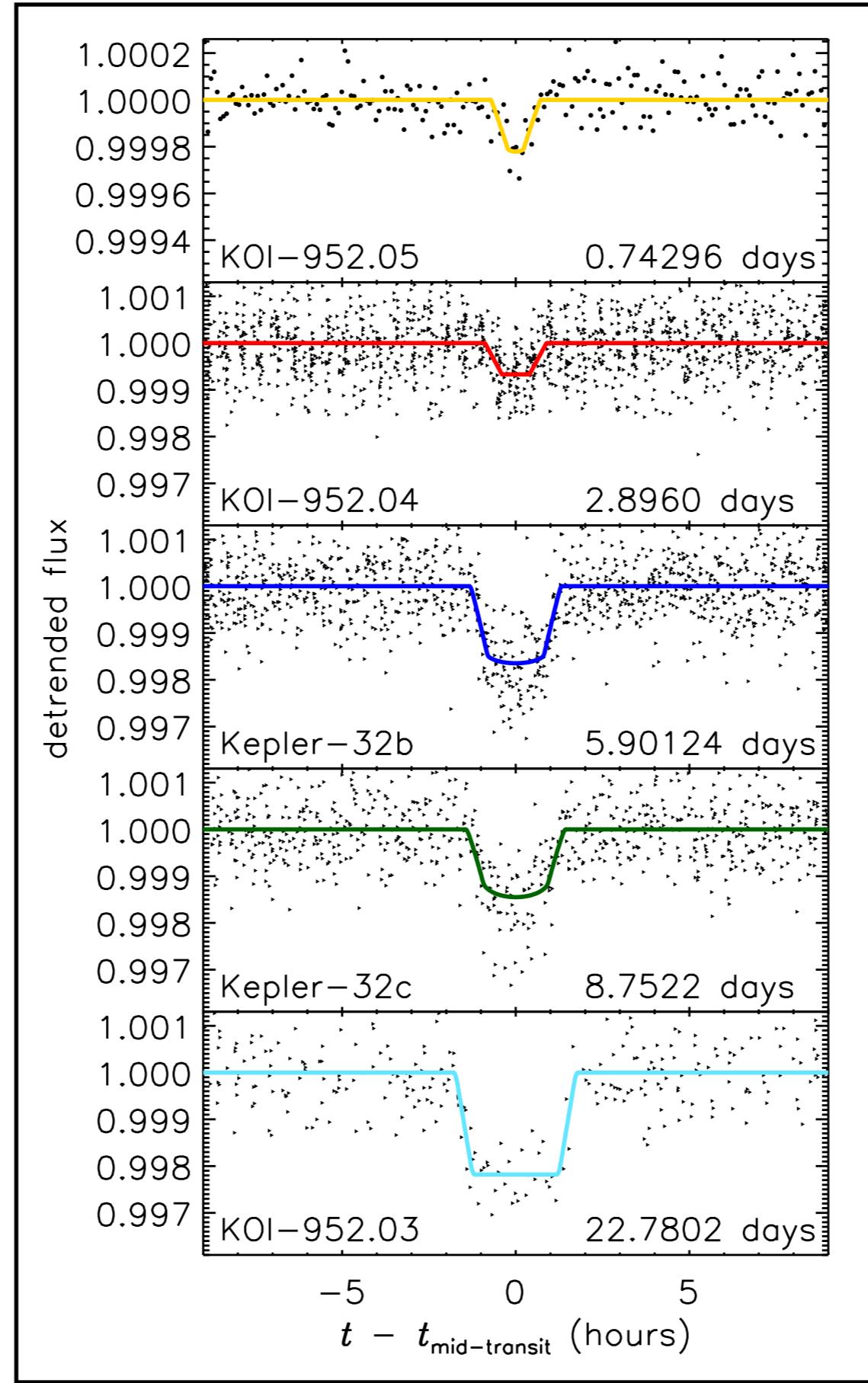


Flickr user Marcin Wichary

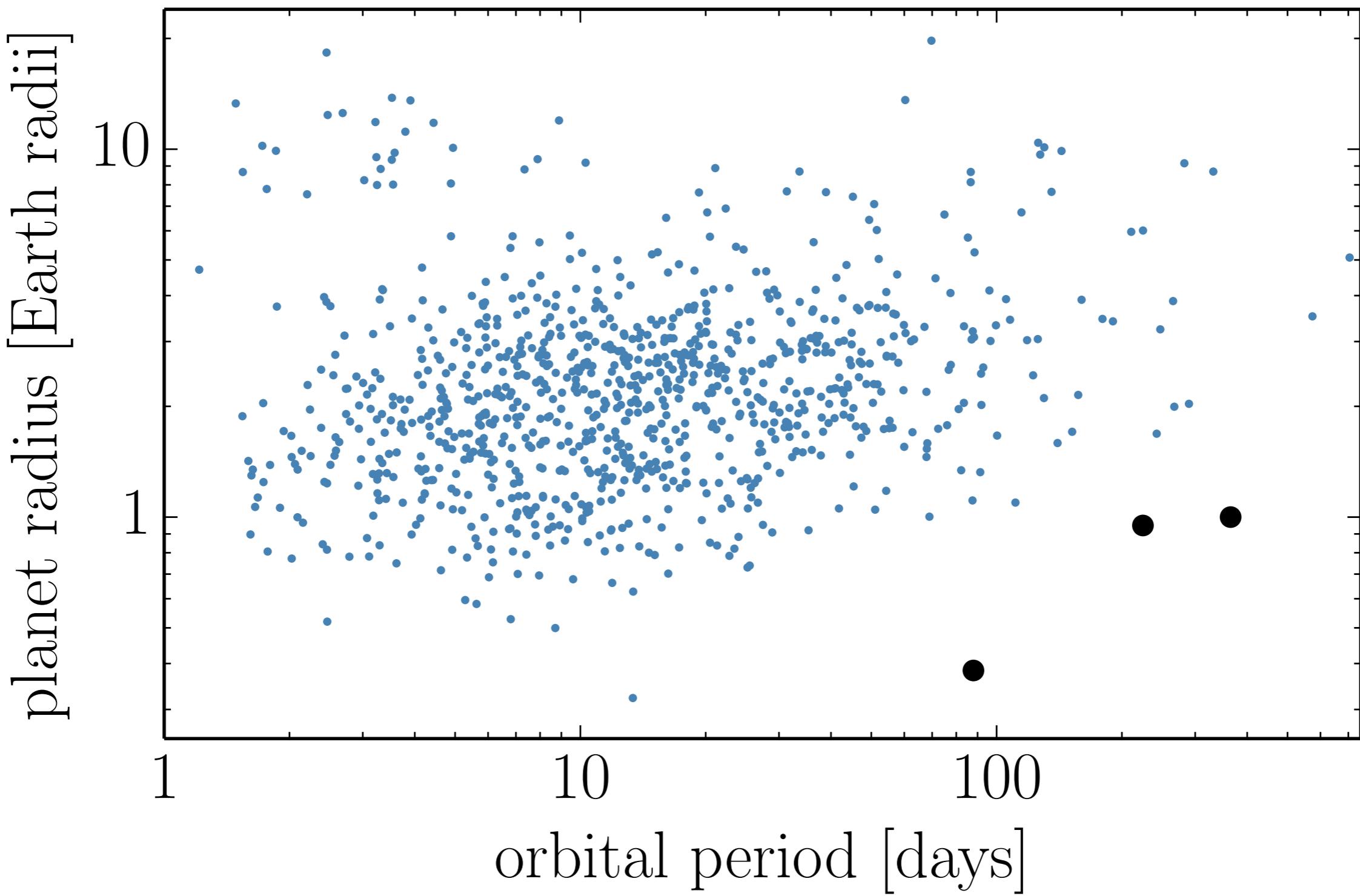
Standard practice: **Filtering**



Kepler-32



Credit **Fabrycky et al. (2012)**



Data from **NASA Exoplanet Archive (11/22)**

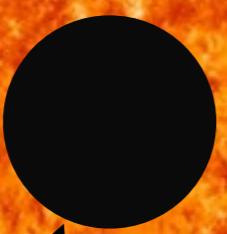
Exoplanets are **hard** to find

Small exoplanets are **harder** to find



Earth

Exoplanets on long orbits are **even harder** to find



Jupiter

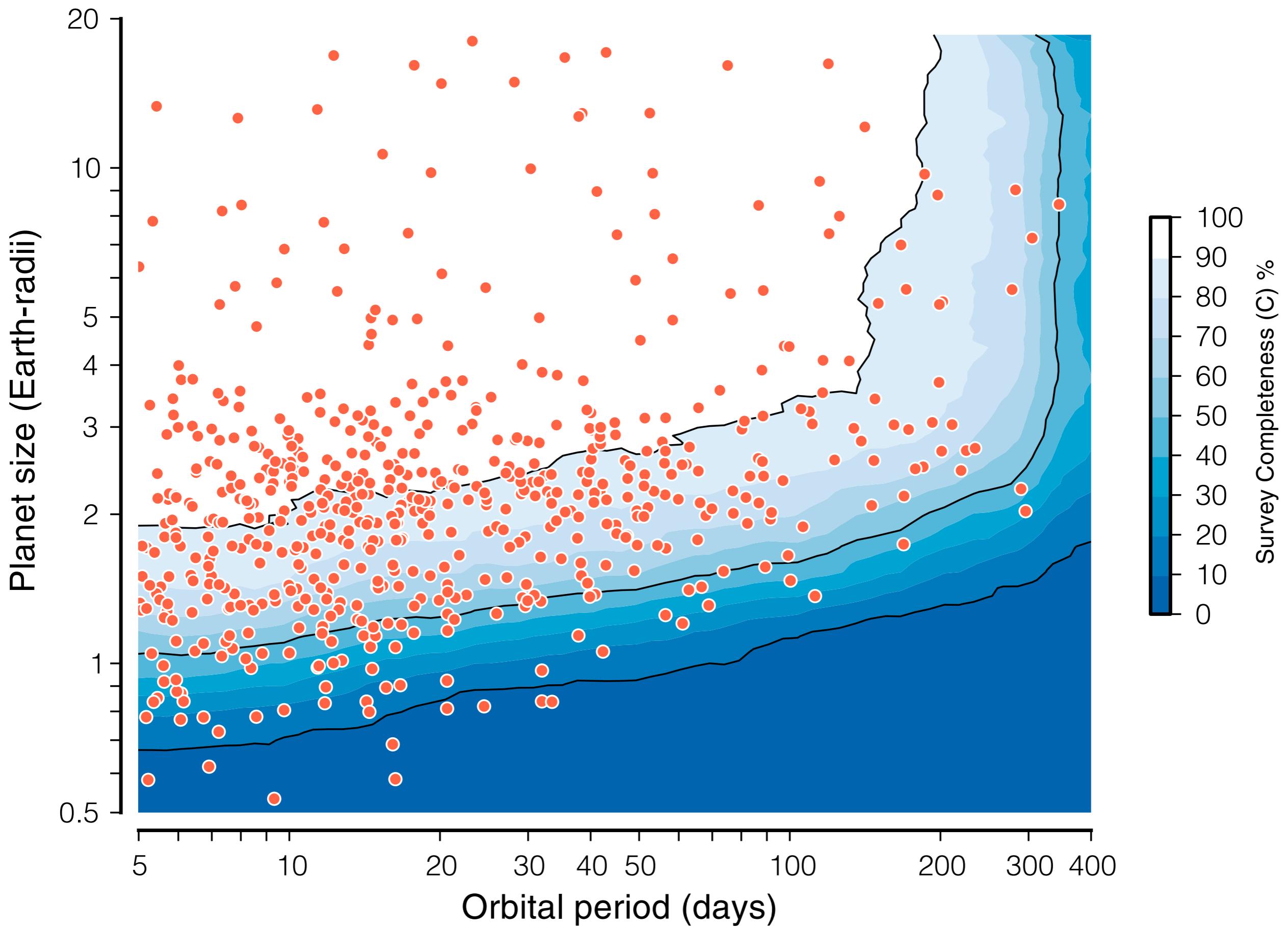


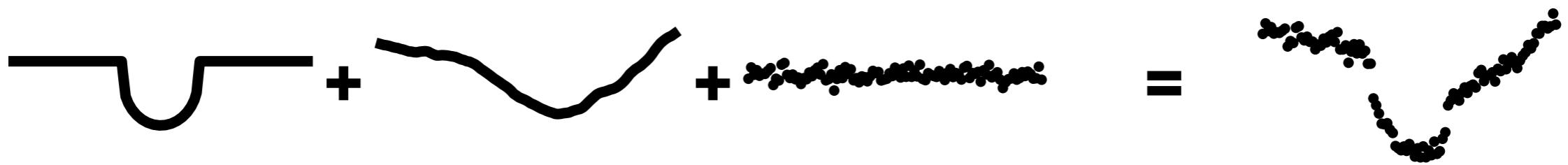
Figure credit: Petigura, Howard & Marcy (2013)

What about Gaussian Processes?

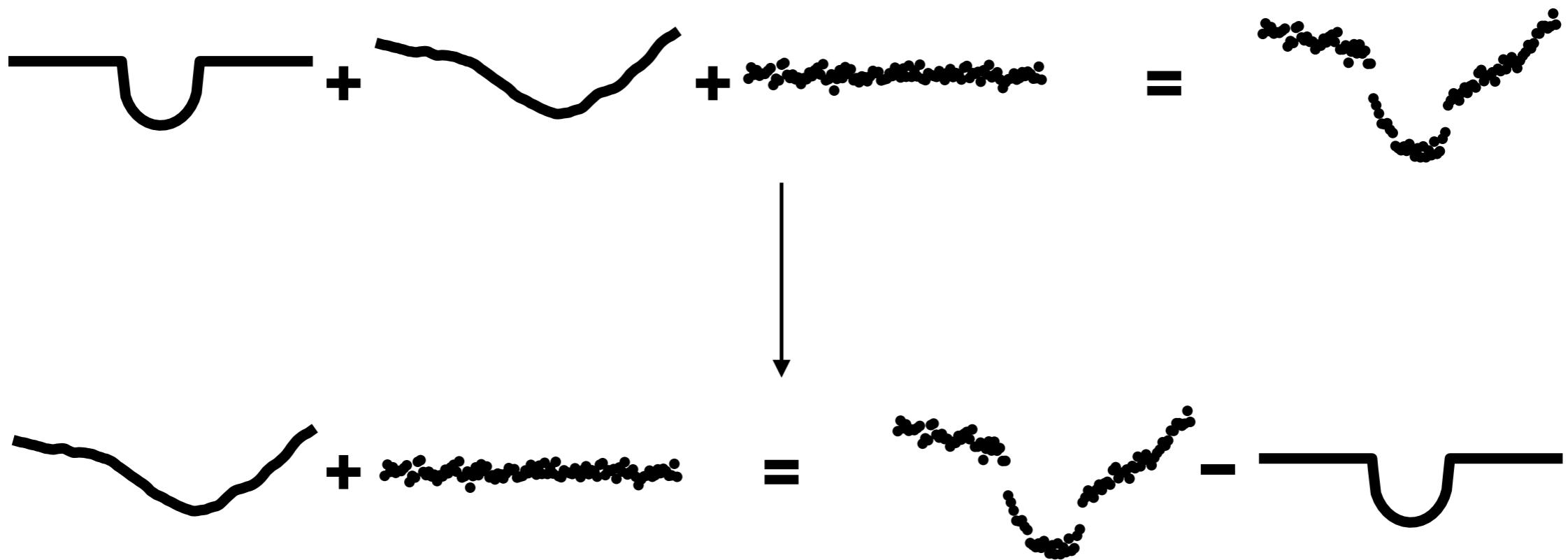
gaussianprocess.org/gpml

Rasmussen & Williams

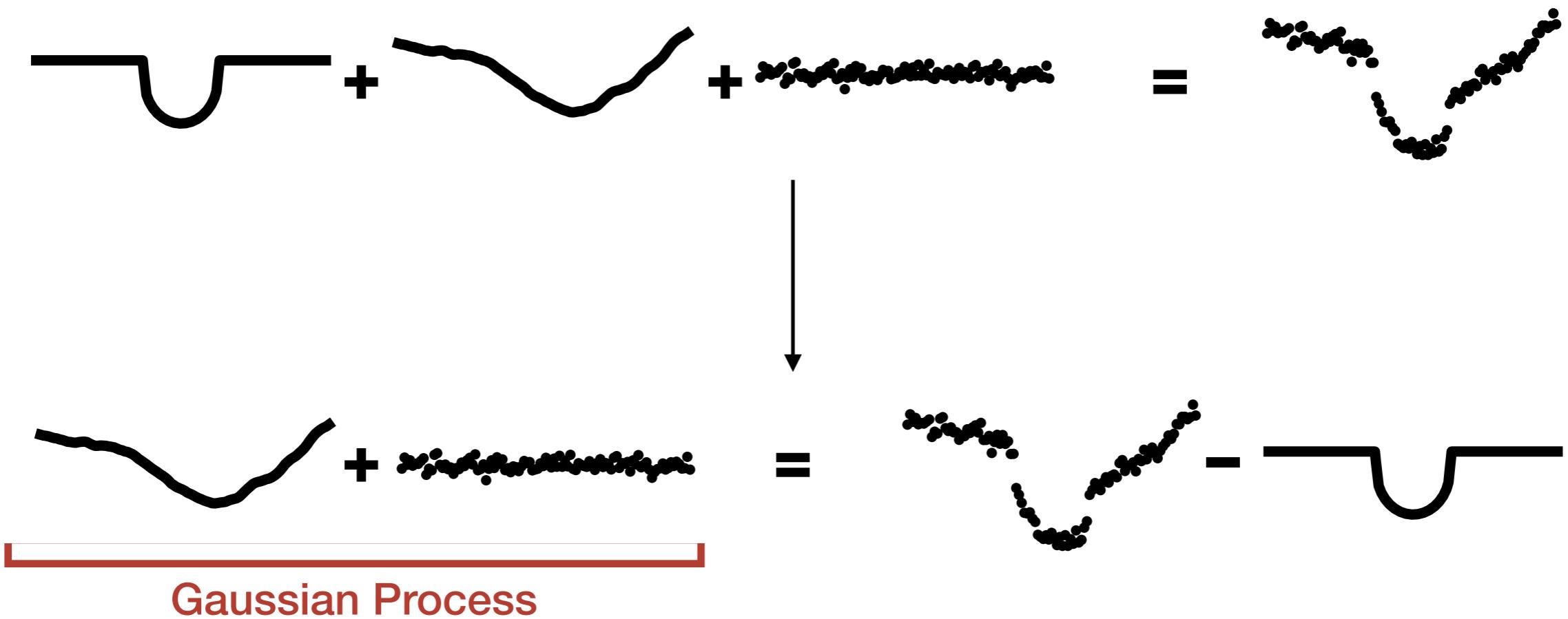
Modeling a **light curve** using a **Gaussian Processes**



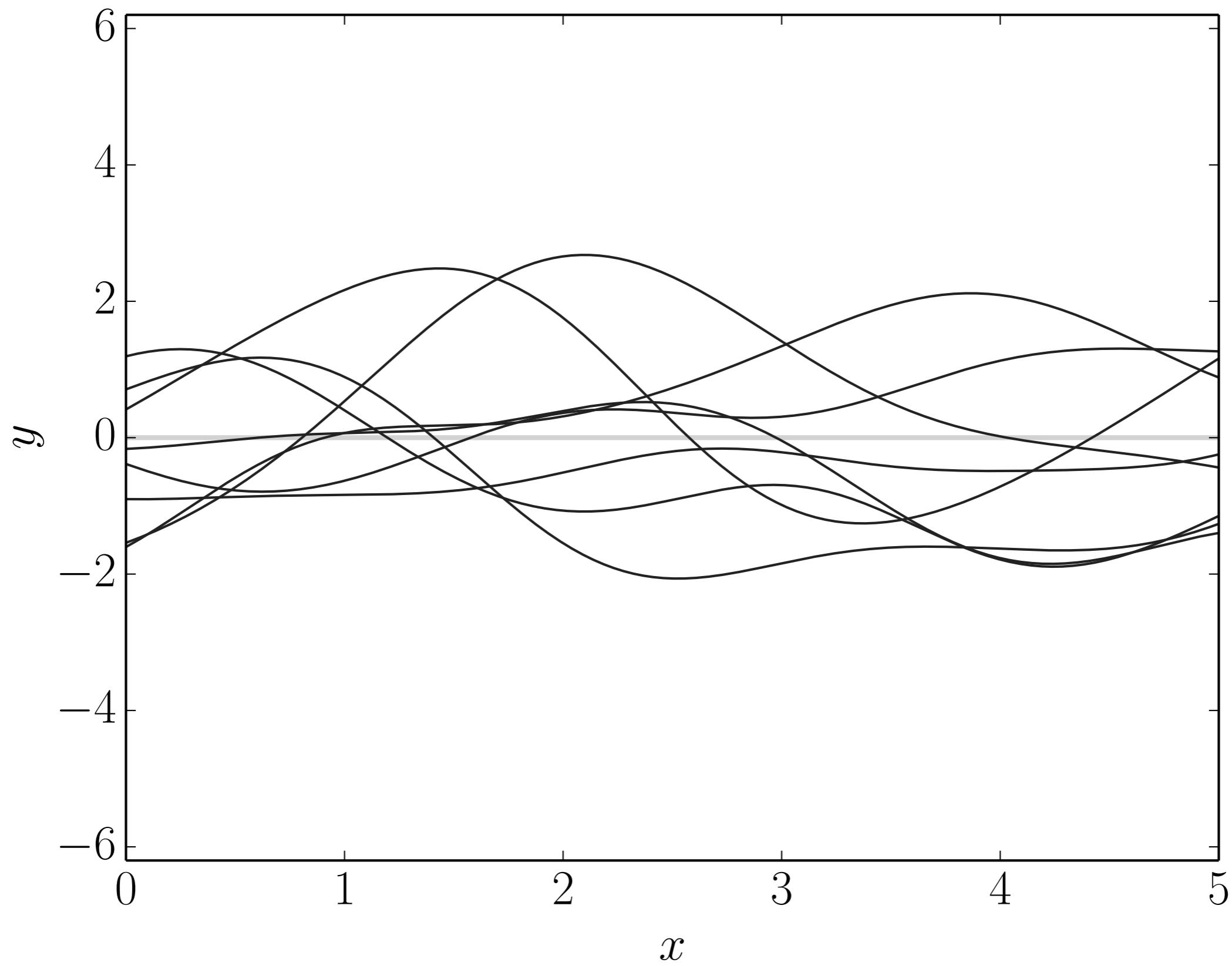
Modeling a **light curve** using a **Gaussian Processes**

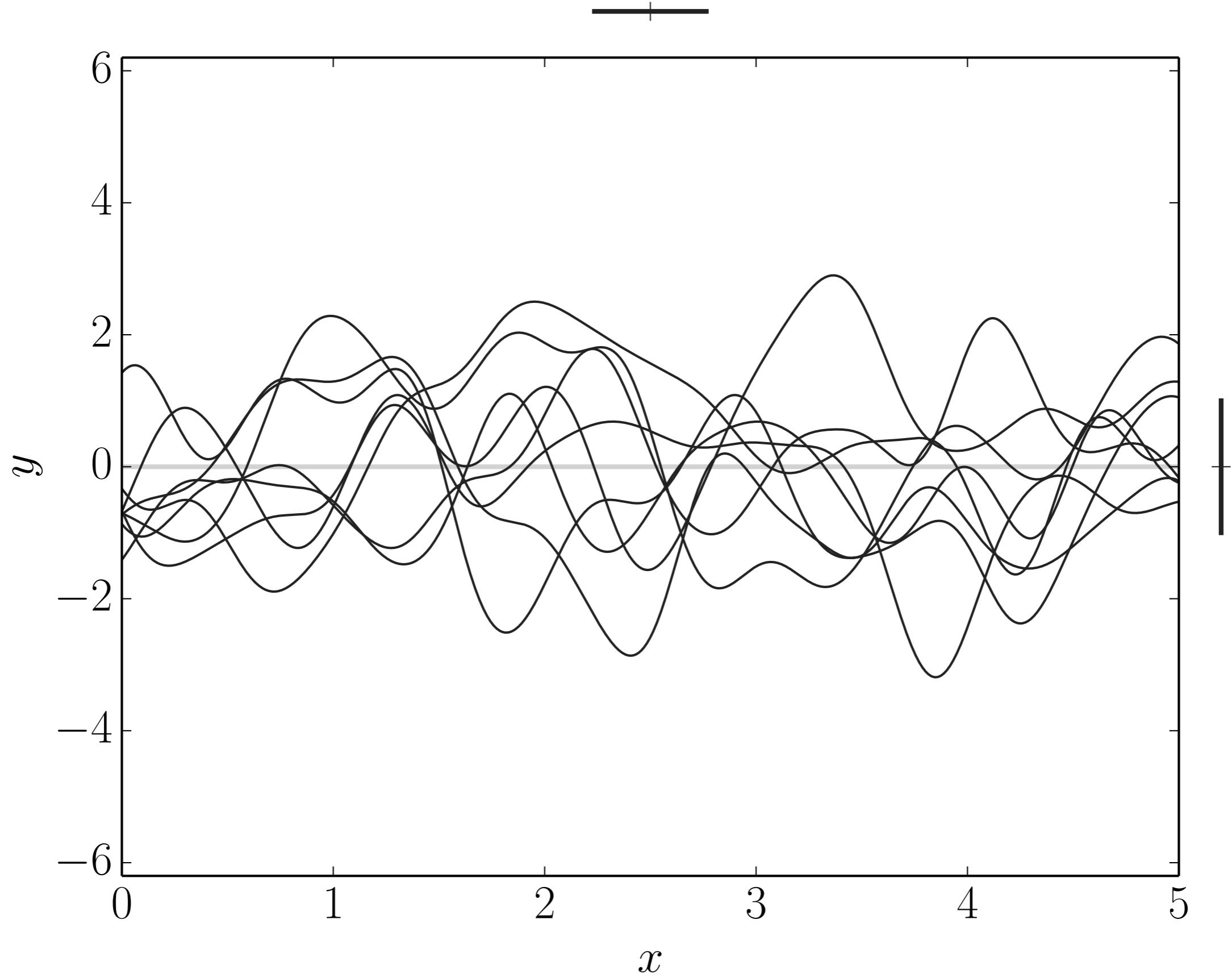


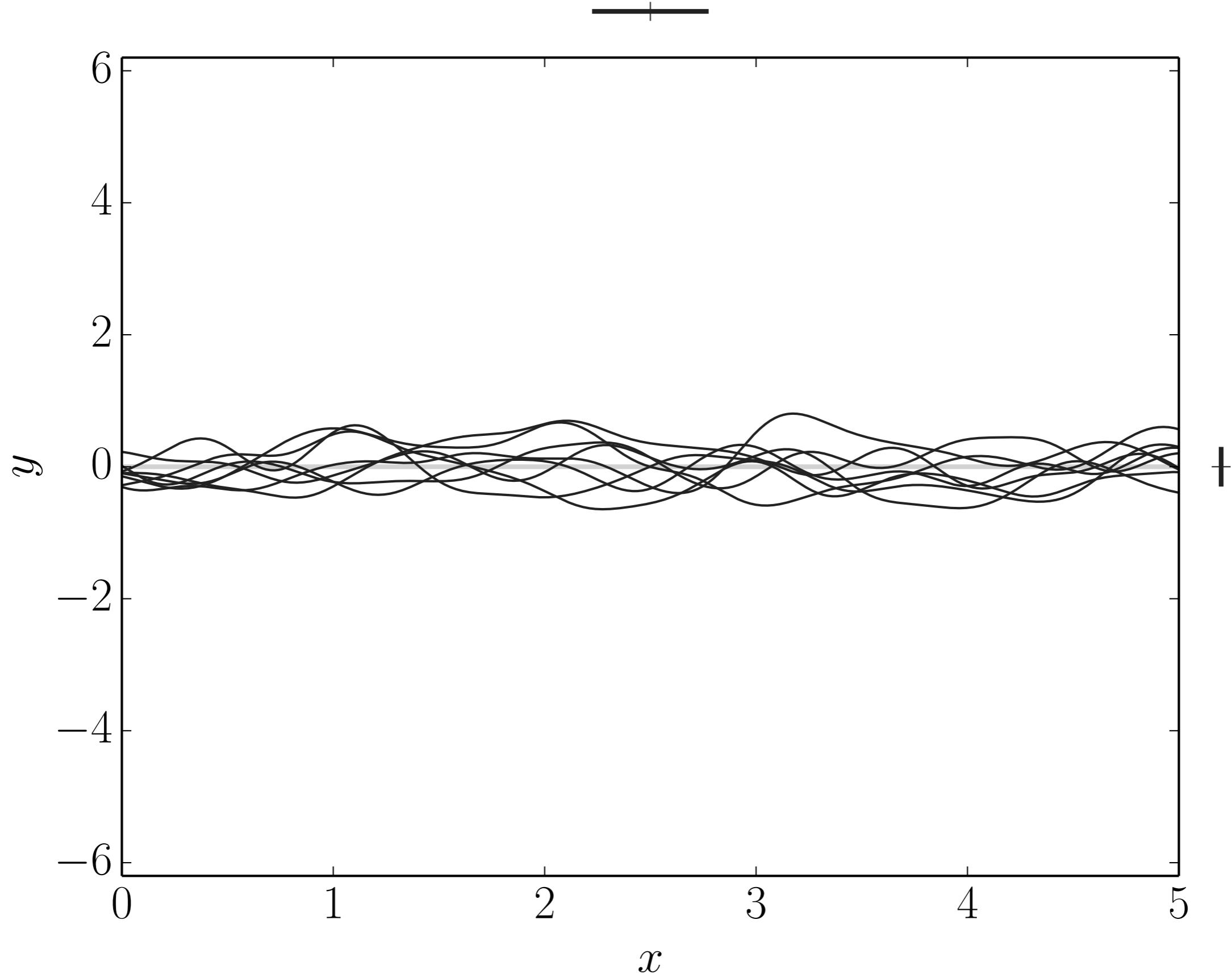
Modeling a light curve using a Gaussian Processes

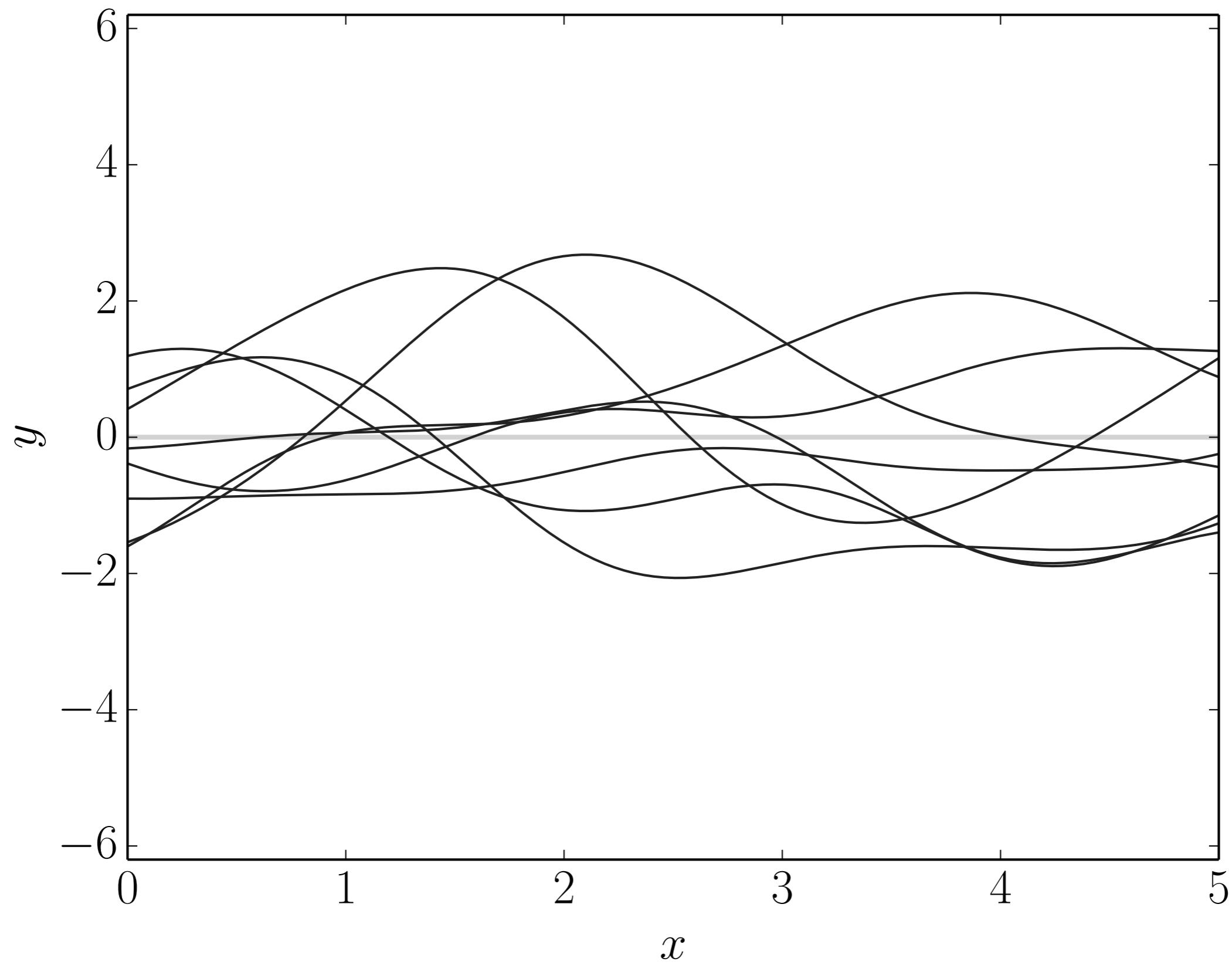


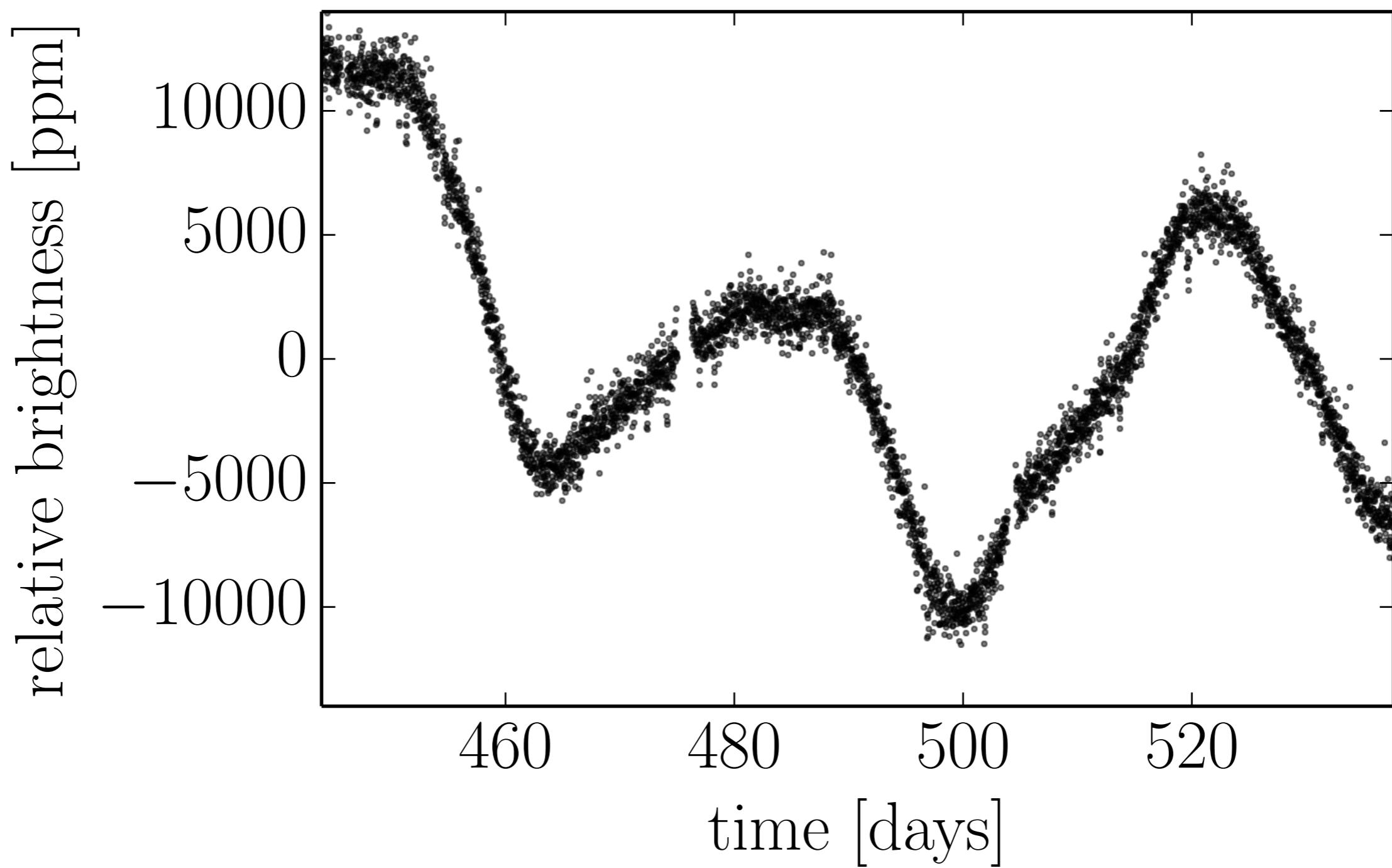
What is a Gaussian Process?

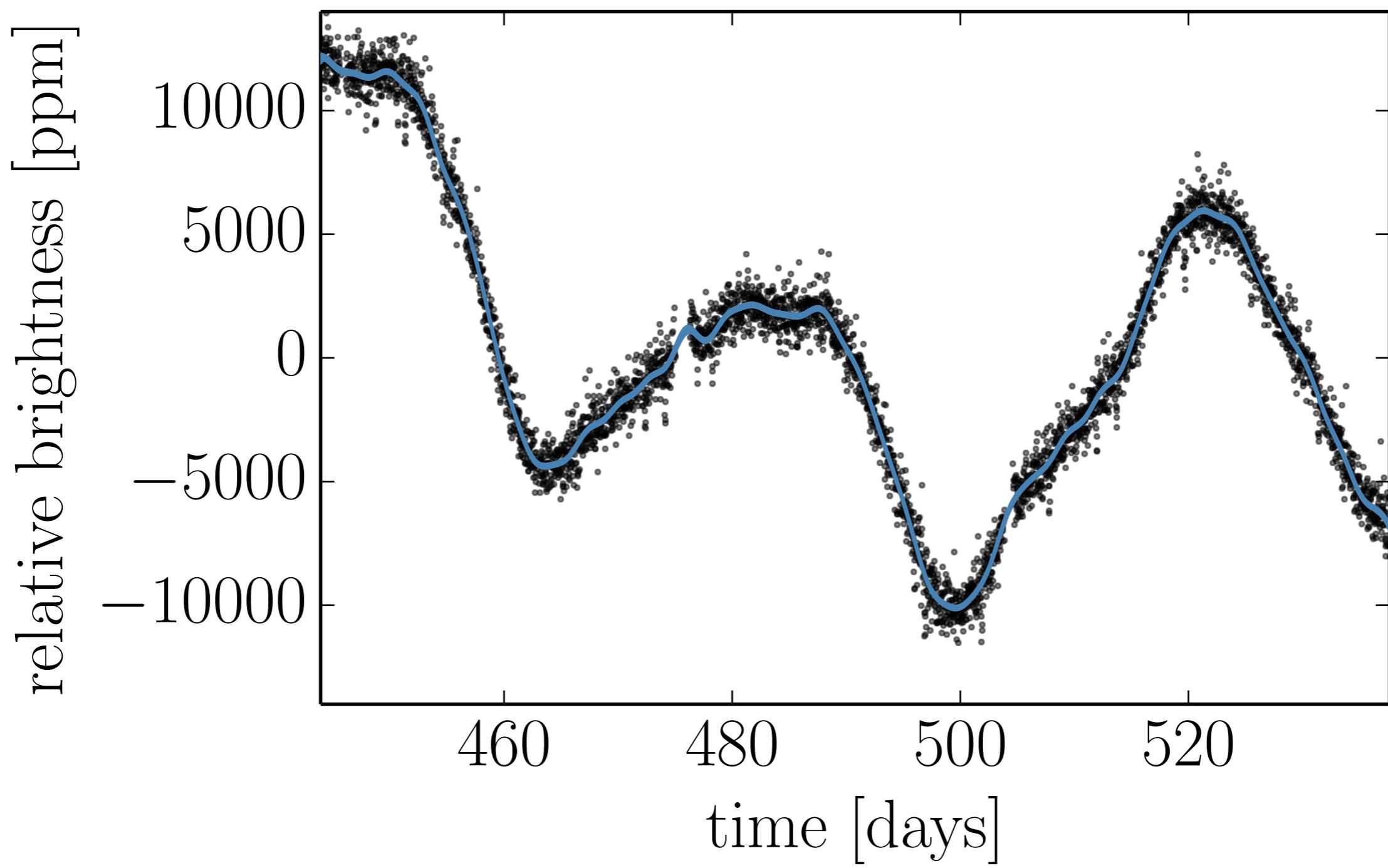












the data are drawn from one

**huge
Gaussian**

*

* the dimension is the number of data points.

The mathematical model

$$y \sim \mathcal{N}(f_{\theta}(x), K_{\alpha}(x, \sigma))$$

where

$$[K_{\alpha}(x, \sigma)]_{ij} = \sigma_i^2 \delta_{ij} + k_{\alpha}(x_i, x_j)$$

The mathematical model

$$\begin{aligned}\log p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\sigma}, \boldsymbol{\theta}, \boldsymbol{\alpha}) = & -\frac{1}{2} [\mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})]^T K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma})^{-1} [\mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})] \\ & -\frac{1}{2} \log \det K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma}) - \frac{N}{2} \log 2 \pi\end{aligned}$$

where

$$[K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma})]_{ij} = \sigma_i^2 \delta_{ij} + k_{\boldsymbol{\alpha}}(x_i, x_j)$$

The mathematical model

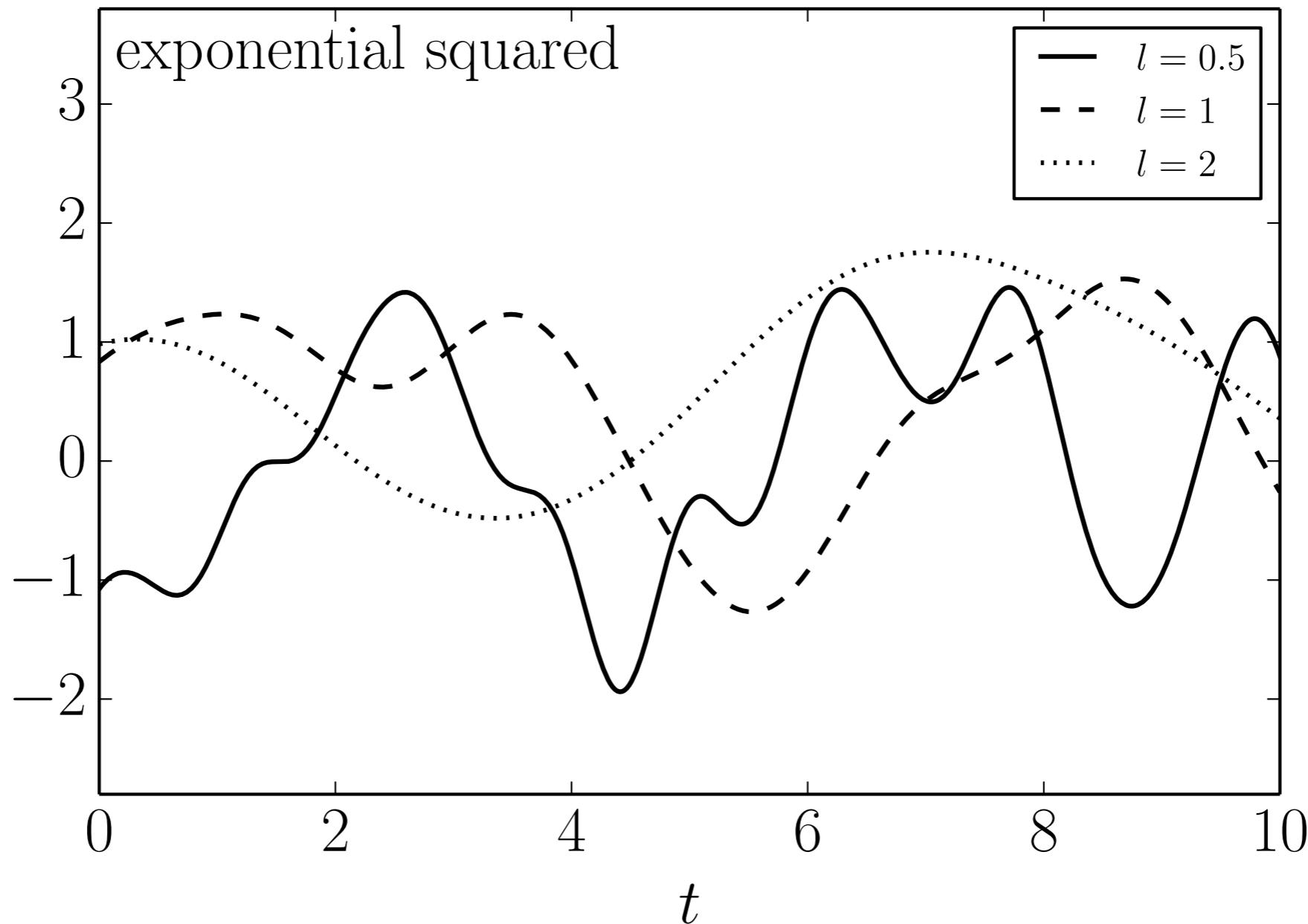
$$\begin{aligned}\log p(\mathbf{y} | \mathbf{x}, \boldsymbol{\sigma}, \boldsymbol{\theta}, \boldsymbol{\alpha}) = & -\frac{1}{2} [\mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})]^T K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma})^{-1} [\mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})] \\ & -\frac{1}{2} \log \det K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma}) - \frac{N}{2} \log 2 \pi\end{aligned}$$

where

$$[K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma})]_{ij} = \sigma_i^2 \delta_{ij} + \underbrace{k_{\boldsymbol{\alpha}}(x_i, x_j)}_{\text{kernel function}} \quad (where \text{ the } magic \text{ happens})$$

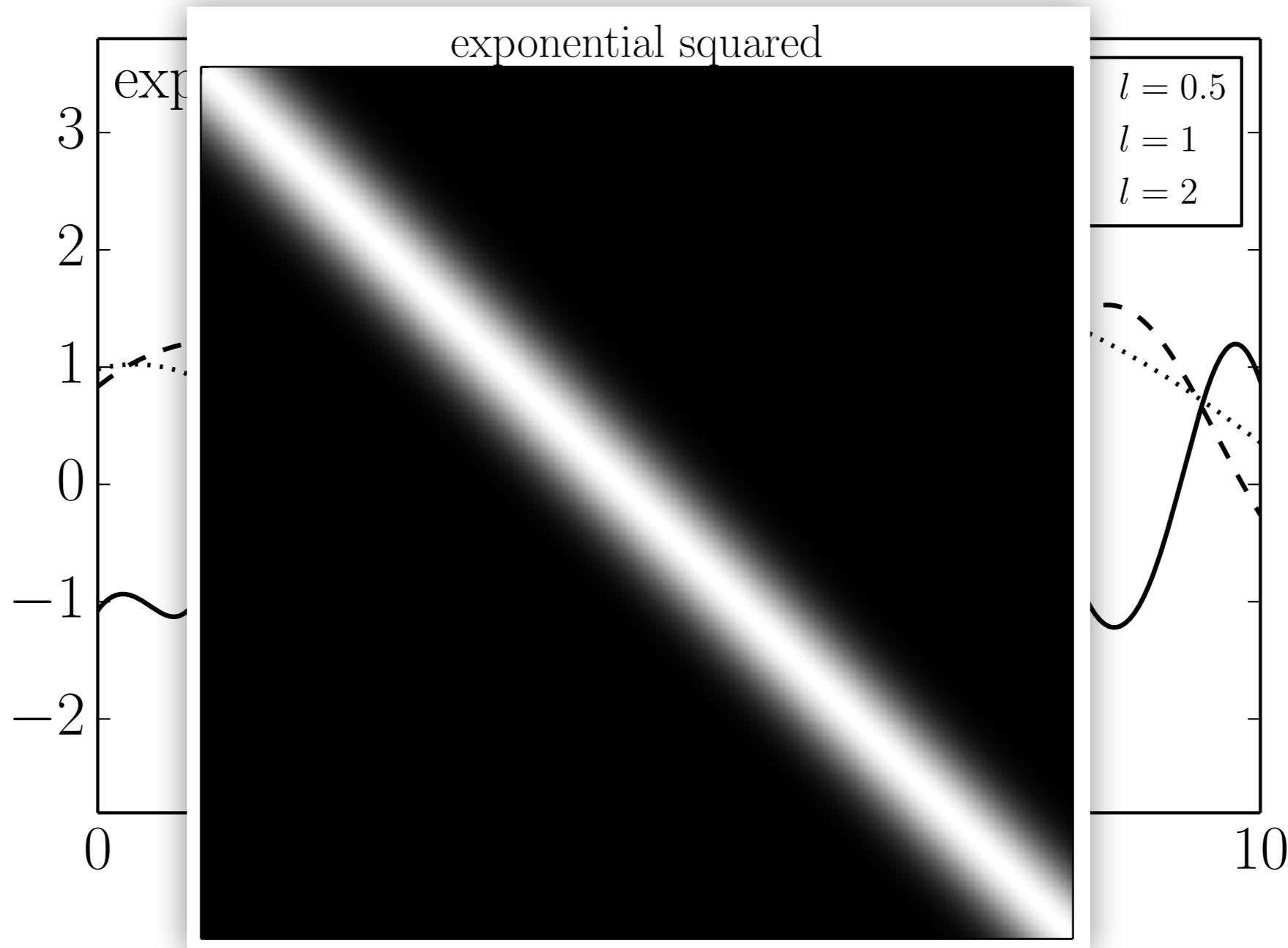
The choice of **kernel**

$$k_{\alpha}(x_i, x_j) = \exp\left(-\frac{[x_i - x_j]^2}{2 \ell^2}\right)$$



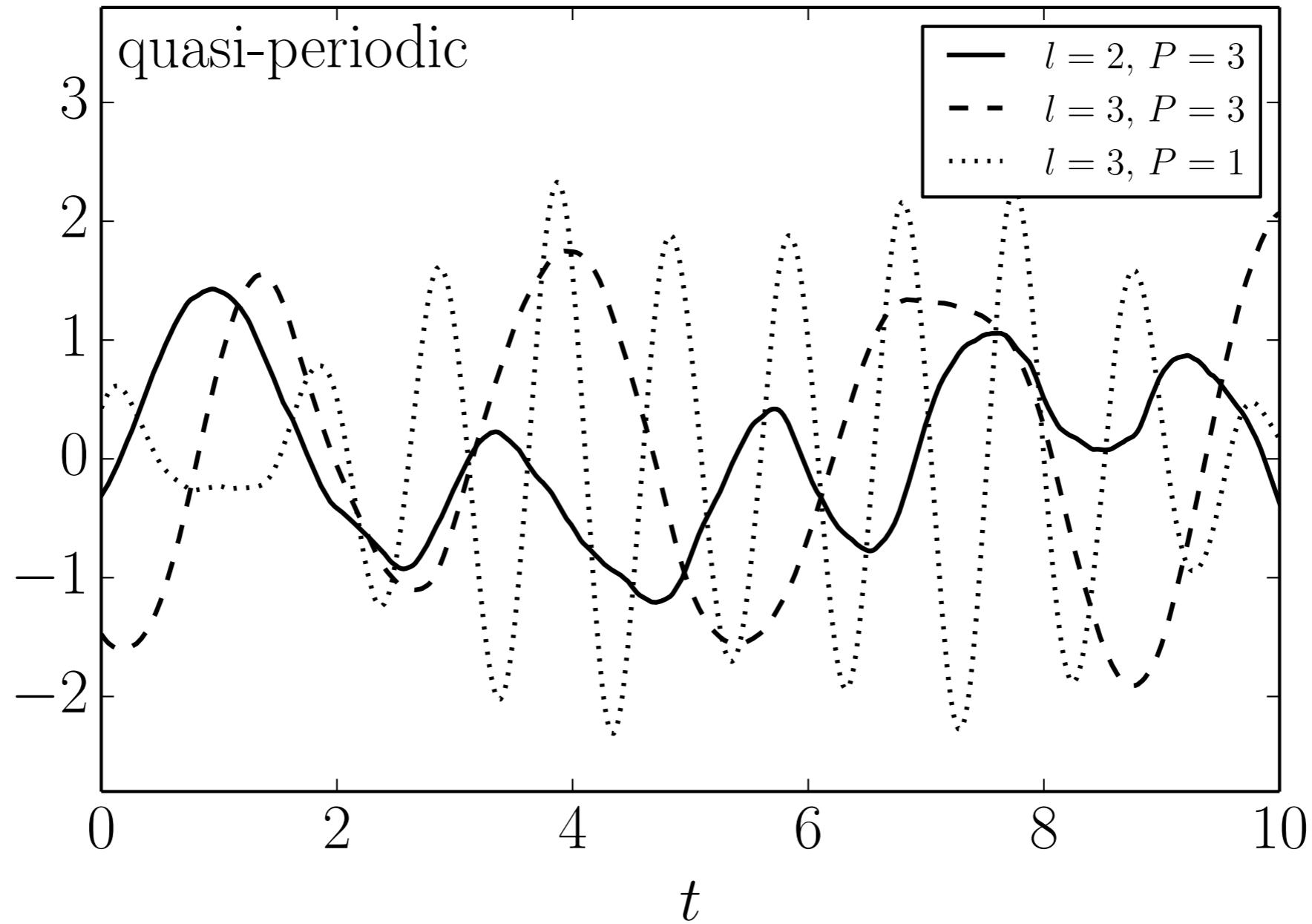
The choice of **kernel**

$$k_{\alpha}(x_i, x_j) = \exp\left(-\frac{[x_i - x_j]^2}{2 \ell^2}\right)$$



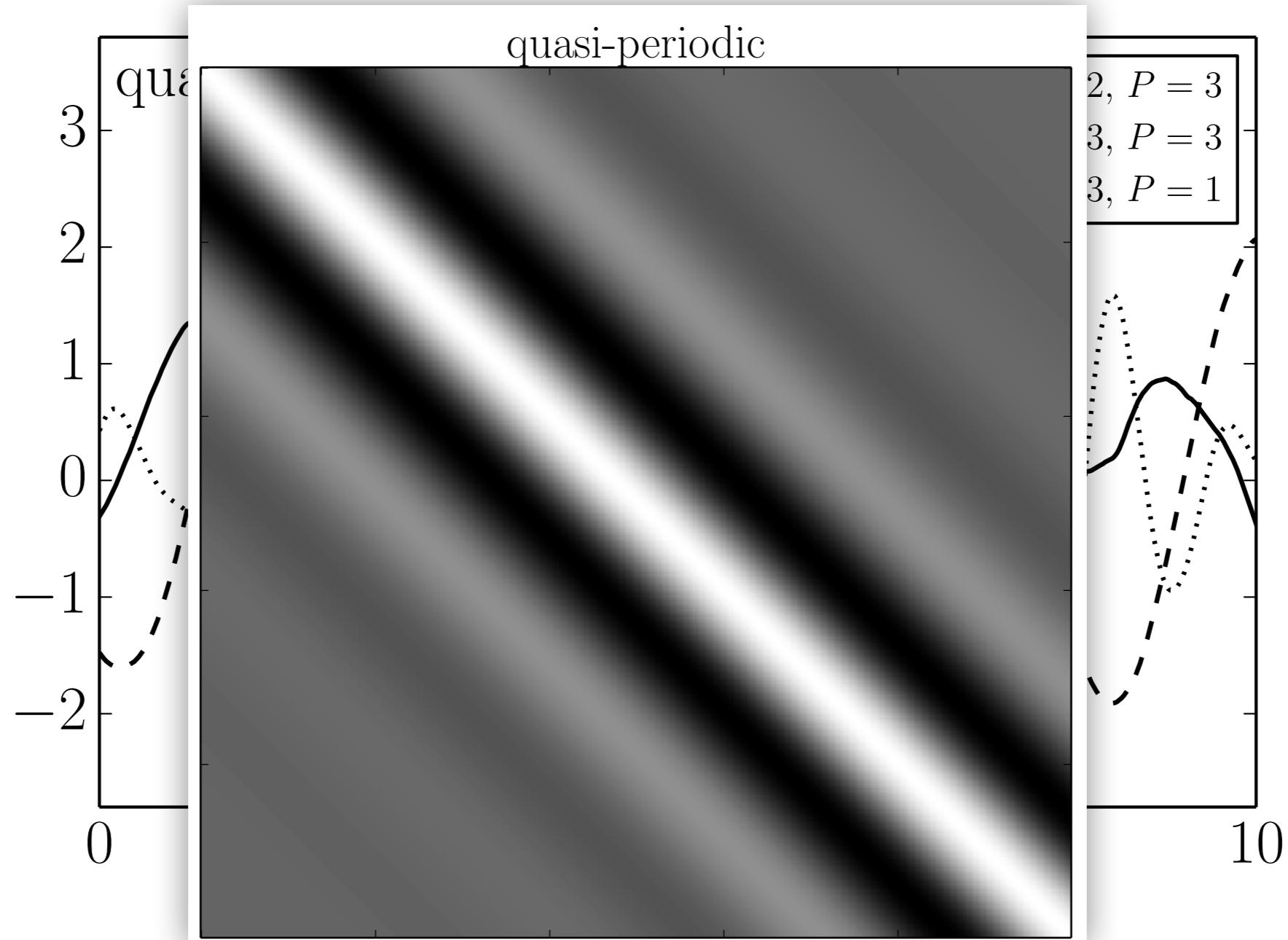
The choice of **kernel**

$$k_{\alpha}(x_i, x_j) = \left[1 + \frac{\sqrt{3} |x_i - x_j|}{\ell} \right] \exp \left(-\frac{|x_i - x_j|}{\ell} \right) \cos \left(\frac{2 \pi |x_i - x_j|}{P} \right)$$

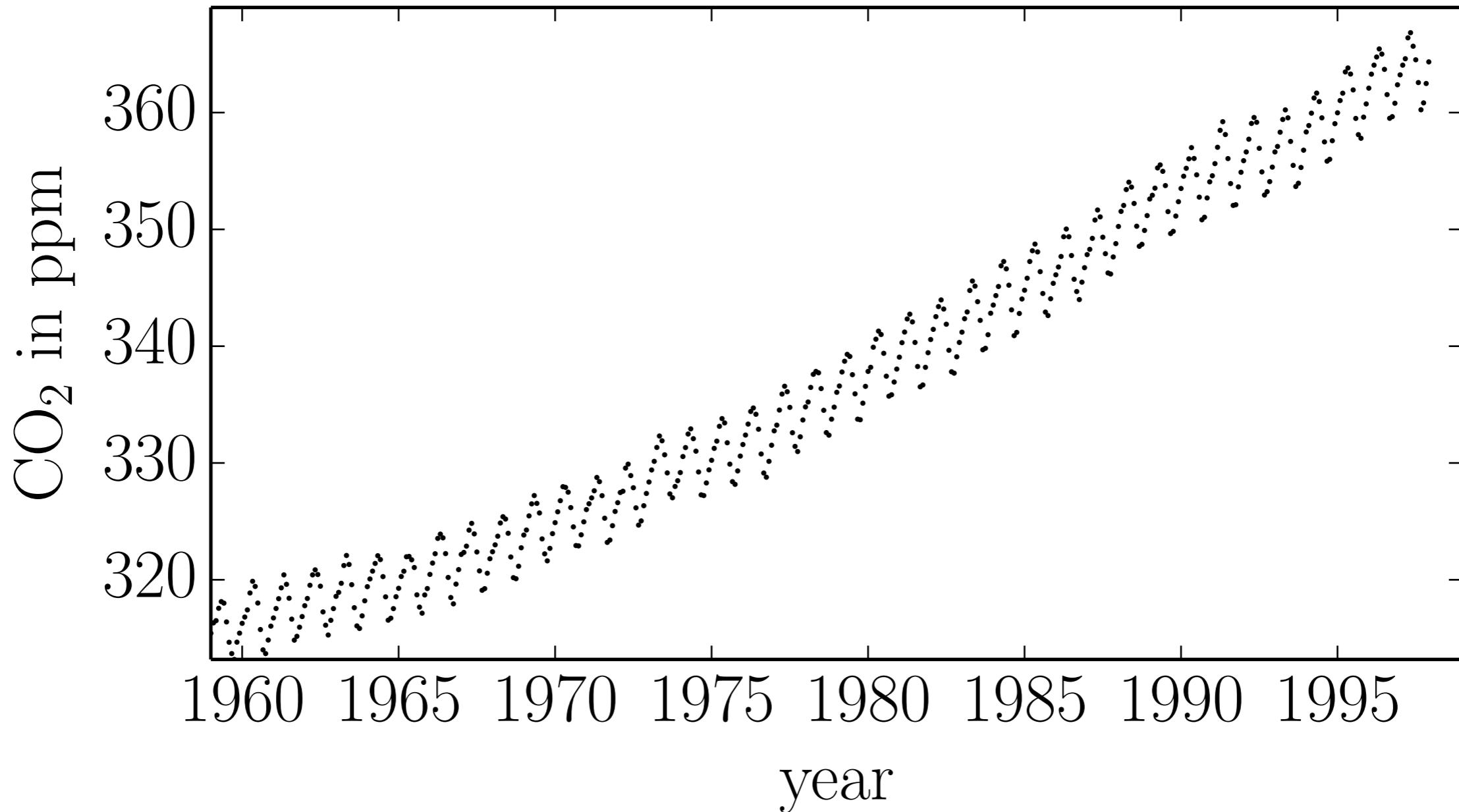


The choice of **kernel**

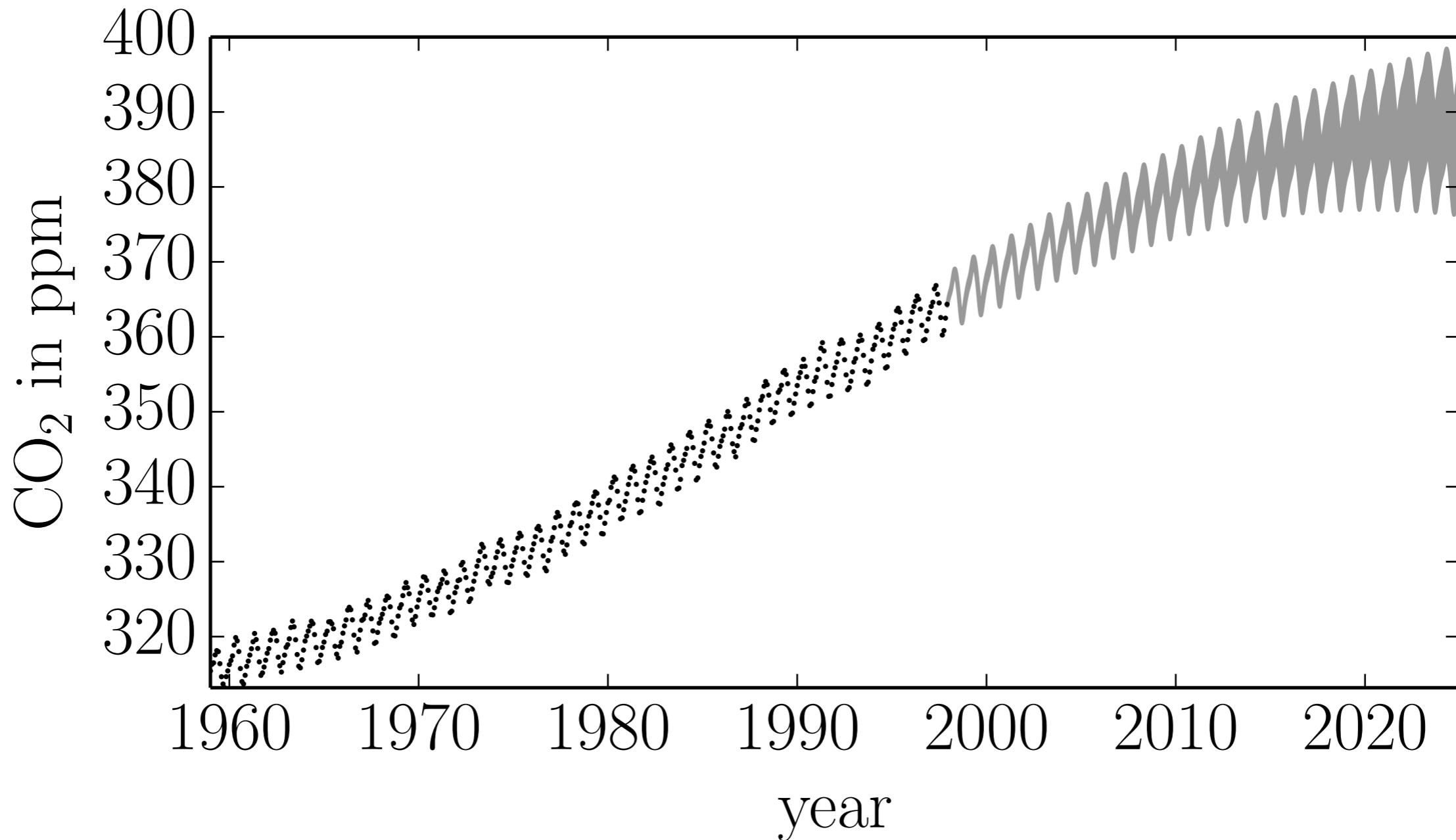
$$k_{\alpha}(x_i, x_j) = \left[1 + \frac{\sqrt{3} |x_i - x_j|}{\ell} \right] \exp \left(-\frac{|x_i - x_j|}{\ell} \right) \cos \left(\frac{2 \pi |x_i - x_j|}{P} \right)$$



The choice of kernel



The choice of **kernel**



The mathematical model

$$\begin{aligned}\log p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\sigma}, \boldsymbol{\theta}, \boldsymbol{\alpha}) = & -\frac{1}{2} [\mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})]^T K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma})^{-1} [\mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})] \\ & -\frac{1}{2} \log \det K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma}) - \frac{N}{2} \log 2 \pi\end{aligned}$$

where

$$[K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma})]_{ij} = \sigma_i^2 \delta_{ij} + \underbrace{k_{\boldsymbol{\alpha}}(x_i, x_j)}_{\text{kernel function}} \quad (where \text{ the } magic \text{ happens})$$

A simple & efficient Python implementation

```
import numpy as np
from scipy.linalg import cho_factor, cho_solve

def kernel(x1, x2):
    # ...

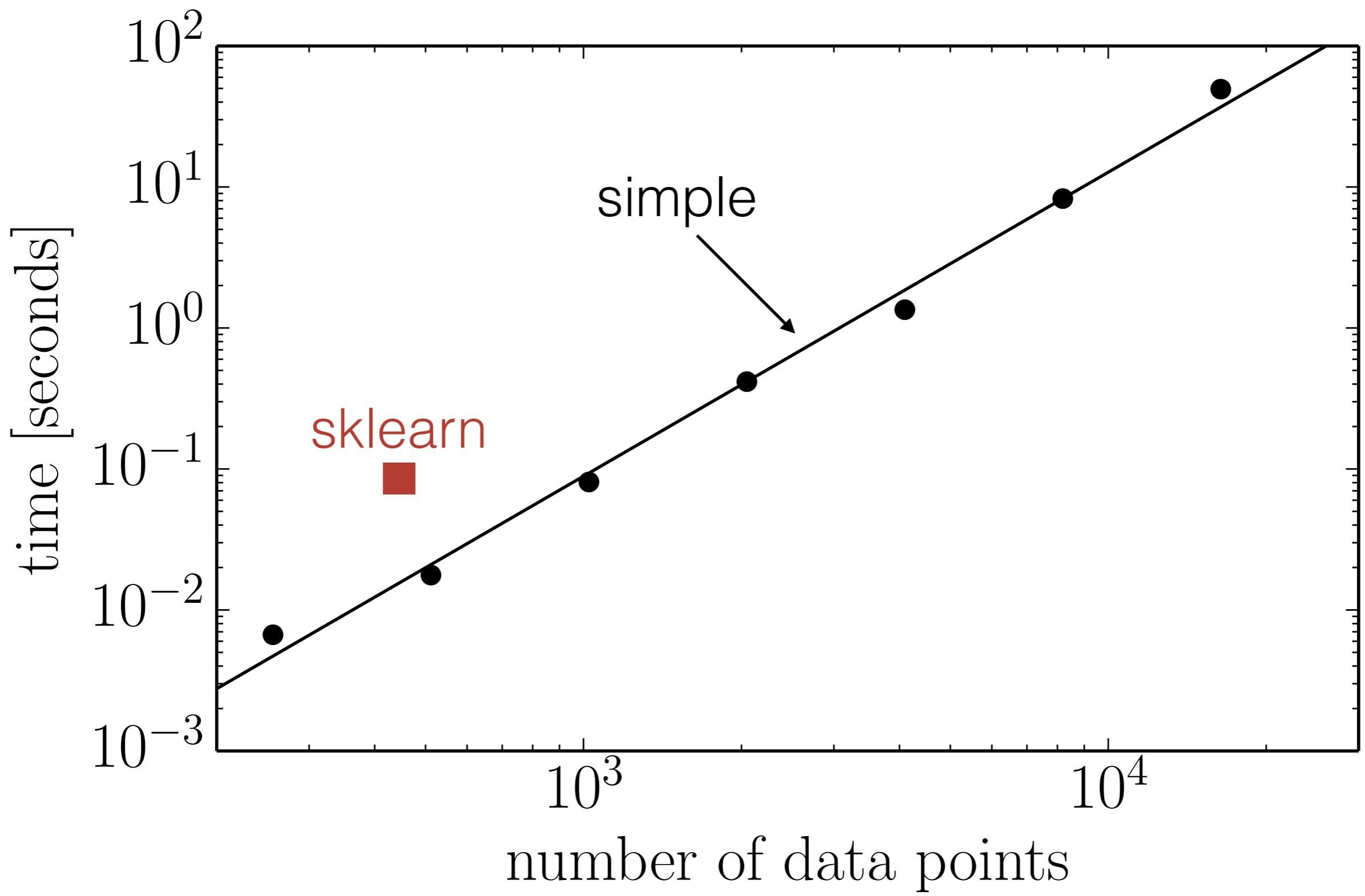
def gp_lnlike(x, y, yerr):
    C = kernel(x[:, None], x[None, :])
    C[np.diag_indices_from(C)] += yerr ** 2
    factor, flag = cho_factor(C)
    logdet = 2*np.sum(np.log(np.diag(factor)))
    return -0.5 * (np.dot(y, cho_solve((factor, flag), y))
                  + logdet + len(x)*np.log(2*np.pi))
```

Using **George**

```
import george
import numpy as np

# kernel = george.kernels...

def george_lnlike(x, y, yerr):
    gp = george.GP(kernel)
    gp.compute(x, yerr)
    return gp.lnlikelihood(y)
```



What's the **catch**?

My Problem

=

Big Data

(by some definition)

Note: I hate myself for this slide too...

Computational complexity.

$$\begin{aligned}\log p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\sigma}, \boldsymbol{\theta}, \boldsymbol{\alpha}) = & -\frac{1}{2} [\mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})]^T K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma})^{-1} [\mathbf{y} - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x})] \\ & -\frac{1}{2} \log \det K_{\boldsymbol{\alpha}}(\mathbf{x}, \boldsymbol{\sigma}) - \frac{N}{2} \log 2 \pi\end{aligned}$$

compute **factorization** // evaluate **log-det** // apply **inverse**

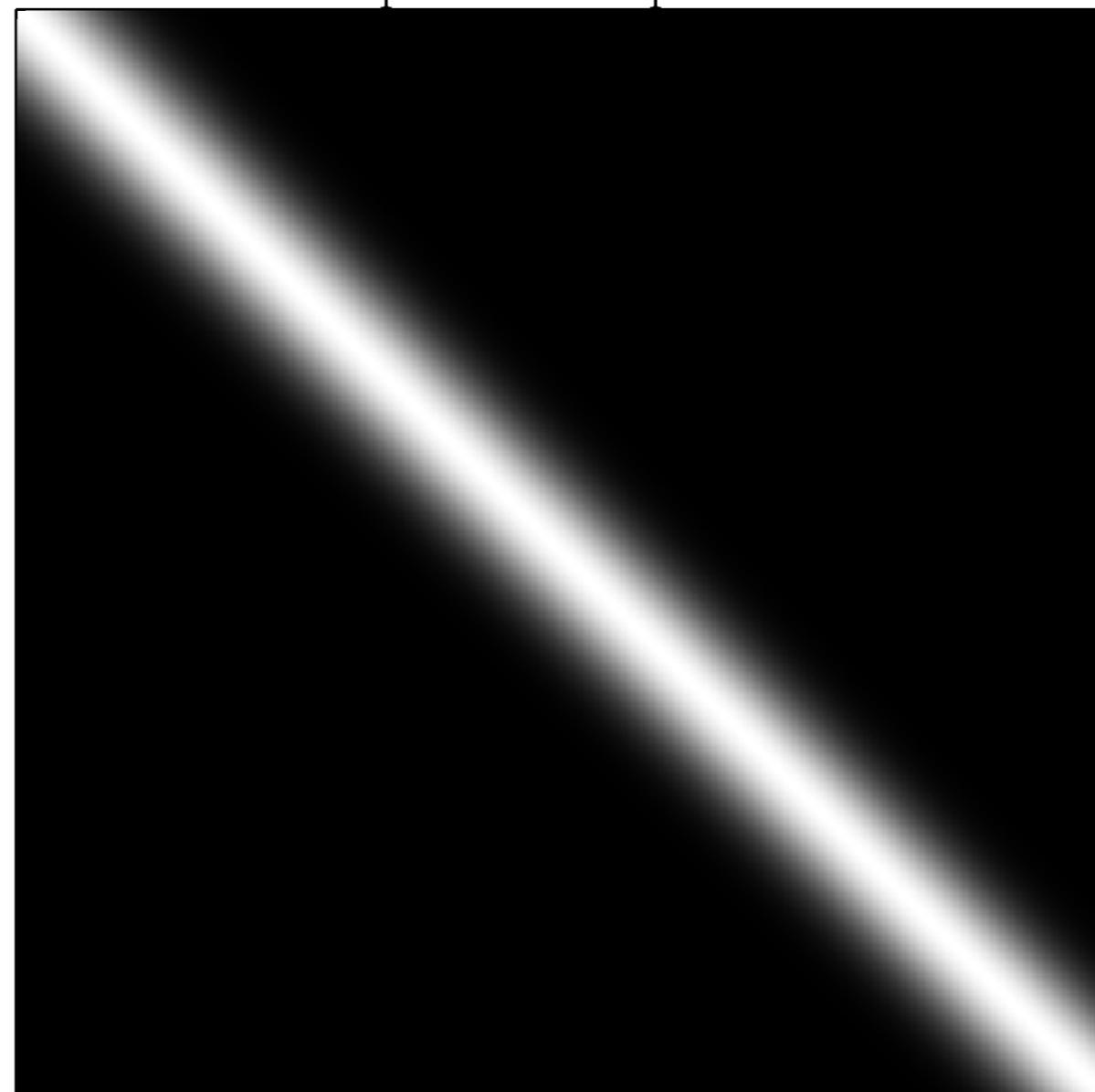
naïvely: $\mathcal{O}(N^3)$

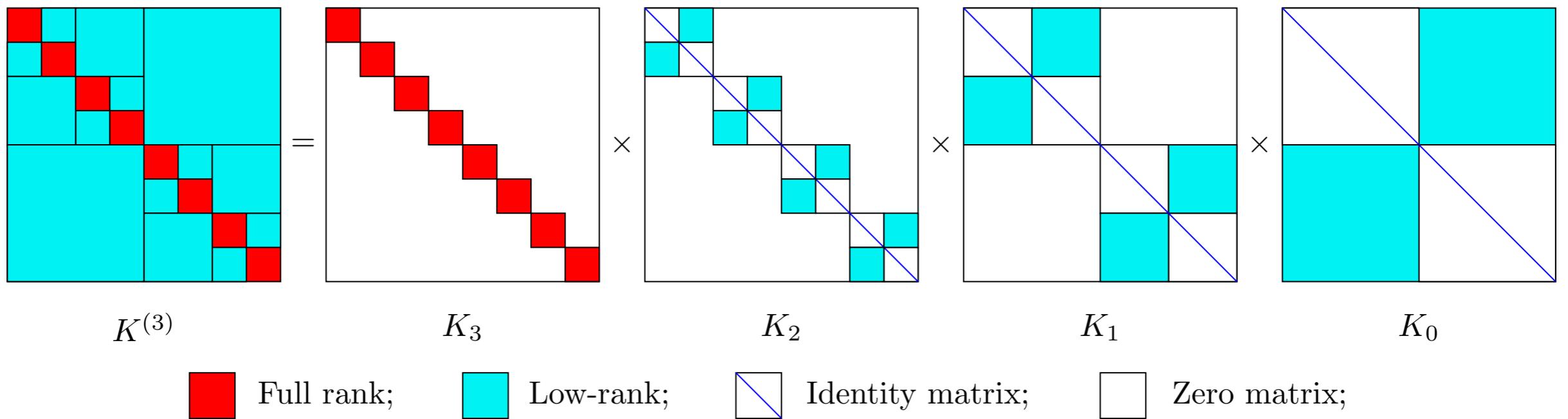


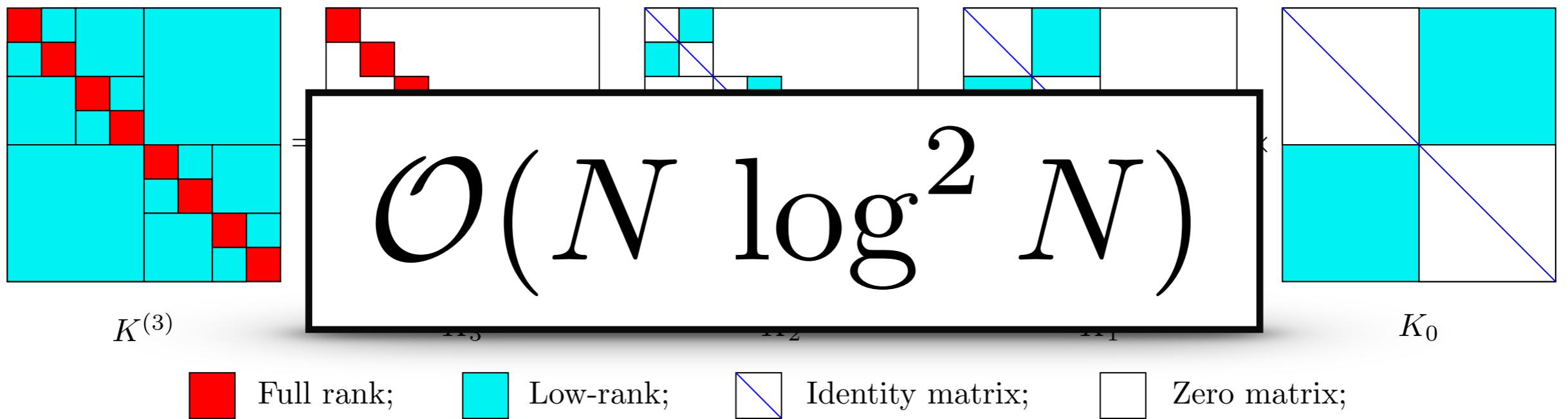
Aren't kernel matrices **Hierarchical Off-Diagonal Low-Rank?**

— not me

exponential squared







github.com/sivaramambikasaran/HODLR

2. dfm@moka | tmux ne...:header (tmux)

```
164 >     >     >     temp.block(nRank[0], 0, nRank[1] , n) => Vinverse[0]*matrix.block(start, 0 , child
165 >
166 >     >     > //> Computes tempSolve => Kinverse\temp-
167 >
168 >     >     > MatrixXd tempSolve => Kinverse.solve(temp);-
169 >
170 >     >     > //> Computes matrix => matrix-Uinverse*tempSolve-
171 >
172 >     >     > matrix.block(start, 0, child[0]->nSize, n) => matrix.block(start, 0, child
173 >     >     > matrix.block(start + child[0]->nSize, 0, child[1]->nSize, n) => matrix.block(sta
174 >     >     }-
175 >     };-
176 >
177 /*!-
178     Computes the determinant of the matrix.-*
179 */
180 > void compute_Determinant() {-
181     if (Kinverse.rows()>0) {           // Check needed when the matrix is predomin
182         MatrixXd LU      = Kinverse.matrixLU();-
183         determinant      = log(fabs(LU(0,0)));-
184         for (int k=1; k<Kinverse.rows(); ++k) {-
185             determinant+=log(fabs(LU(k,k)));-
186         }-
187         // Previous version which had some underflow.-*
188         // determinant=> log(fabs(K.determinant()));-
189     }-
190 }-
HODLR_Node.hpp [cpp]
```

george 1:Vim ✘ 0 ✓ 0 17 Nov 19:22

The HODLR solver from **George**

```
import george
import numpy as np

# kernel = george.kernels...

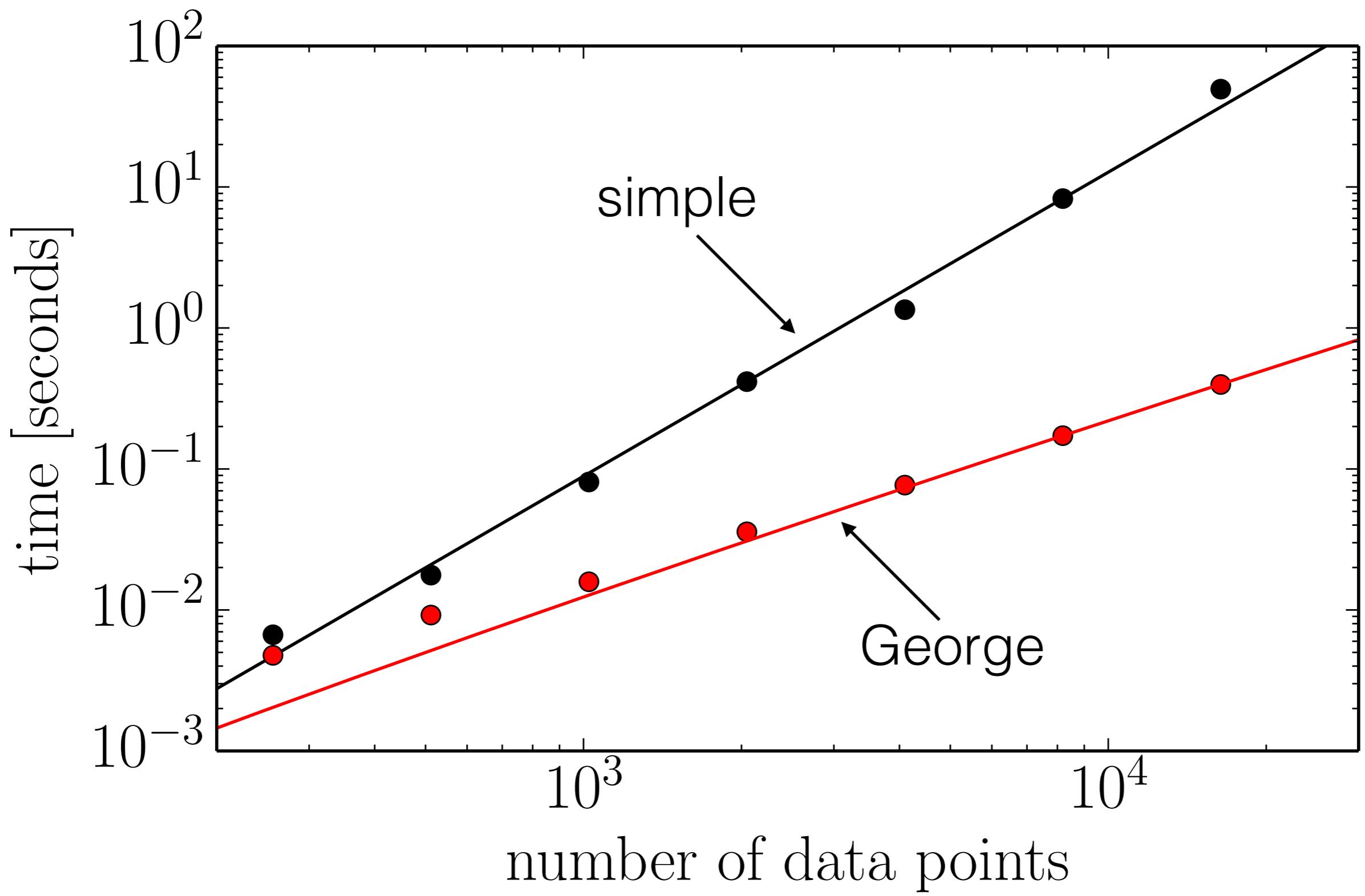
def george_lnlike(x, y, yerr):
    gp = george.GP(kernel)
    gp.compute(x, yerr)
    return gp.lnlikelihood(y)
```

The HODLR solver from George

```
import george
import numpy as np

# kernel = george.kernels...

def george_lnlike(x, y, yerr):
    gp = george.GP(kernel, solver=george.HODLRSolver)
    gp.compute(x, yerr)
    return gp.lnlikelihood(y)
```



References

gaussianprocess.org/gpml

[github.com/dfm/gp
george](https://github.com/dfm/gpgeorge)

danfm@nyu.edu