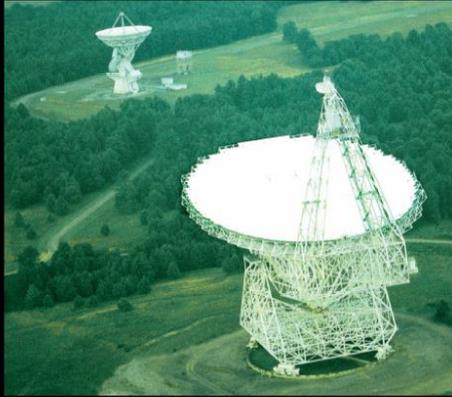


PHIL GREGORY

**Bayesian Logical  
Data Analysis  
for the Physical Sciences**

A Comparative Approach with  
*Mathematica* Support



CAMBRIDGE

# **Bayesian Planet Searches for the 10 cm/s Radial Velocity Era**

**Phil Gregory  
University of British Columbia  
Vancouver, Canada**

**Aug. 4, 2015**

**IAU Honolulu  
Focus Meeting 8  
On Statistics and Exoplanets**

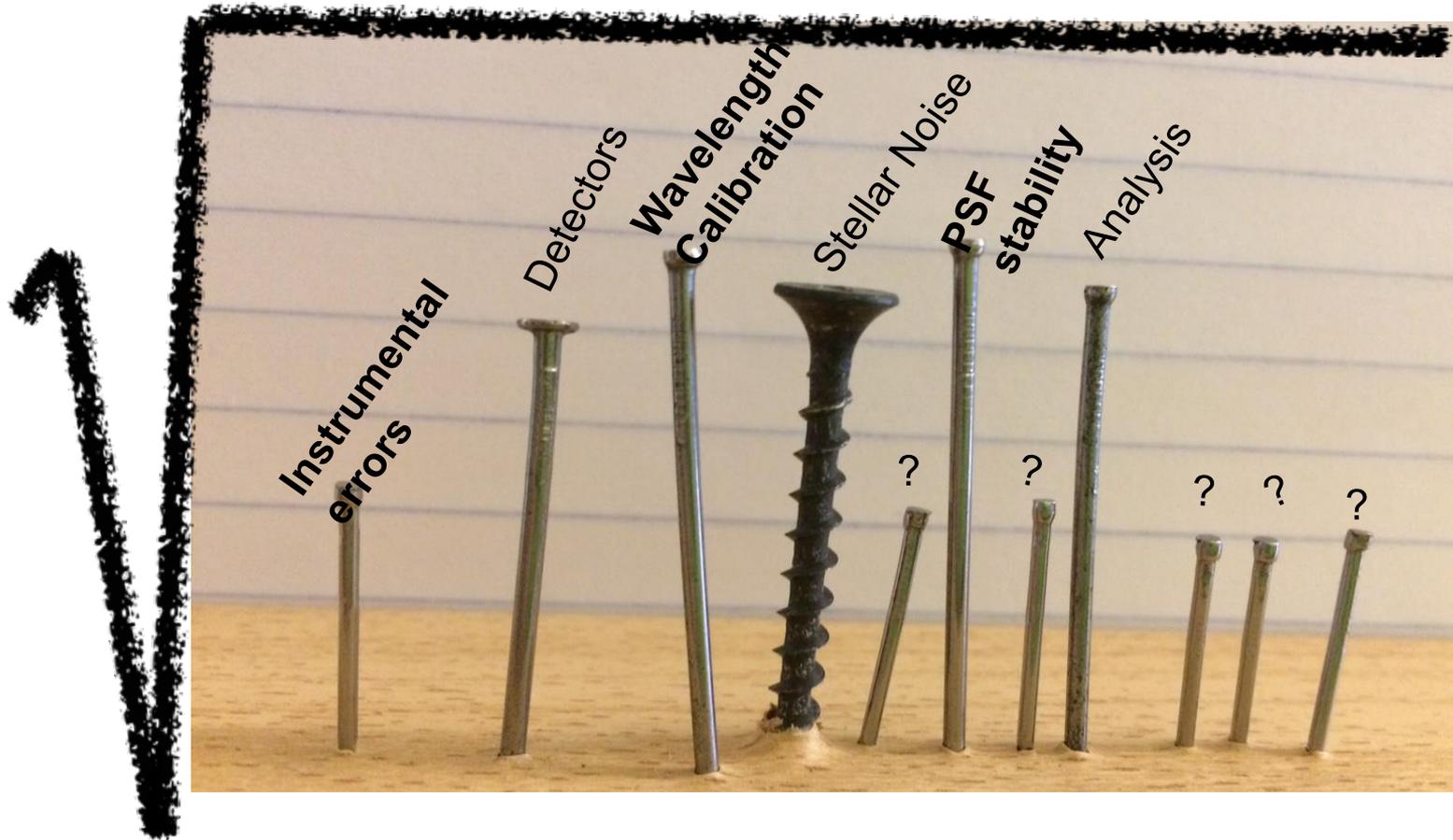
# **Bayesian planet searches for the 10 cm/s radial velocity era**

**Intrinsic stellar variability has become the main limiting factor for planet searches in both transit and radial velocity (RV) data. New spectrographs are under development like ESPRESSO and EXPRES that aim to improve RV precision by a factor of approximately 10 over the current best spectrographs, HARPS and HARPS-N. This will greatly exacerbate the challenge of distinguishing planetary signals from stellar activity induced RV signals.**

**At the same time good progress has been made in simulating stellar activity signals. At the Porto 2014 meeting, “Towards Other Earths II,” Xavier Dumusque challenged the community to a large scale blind test using the simulated RV data at the 1 m/s level of precision, to understand the limitations of present solutions to deal with stellar signals and to select the best approach.**

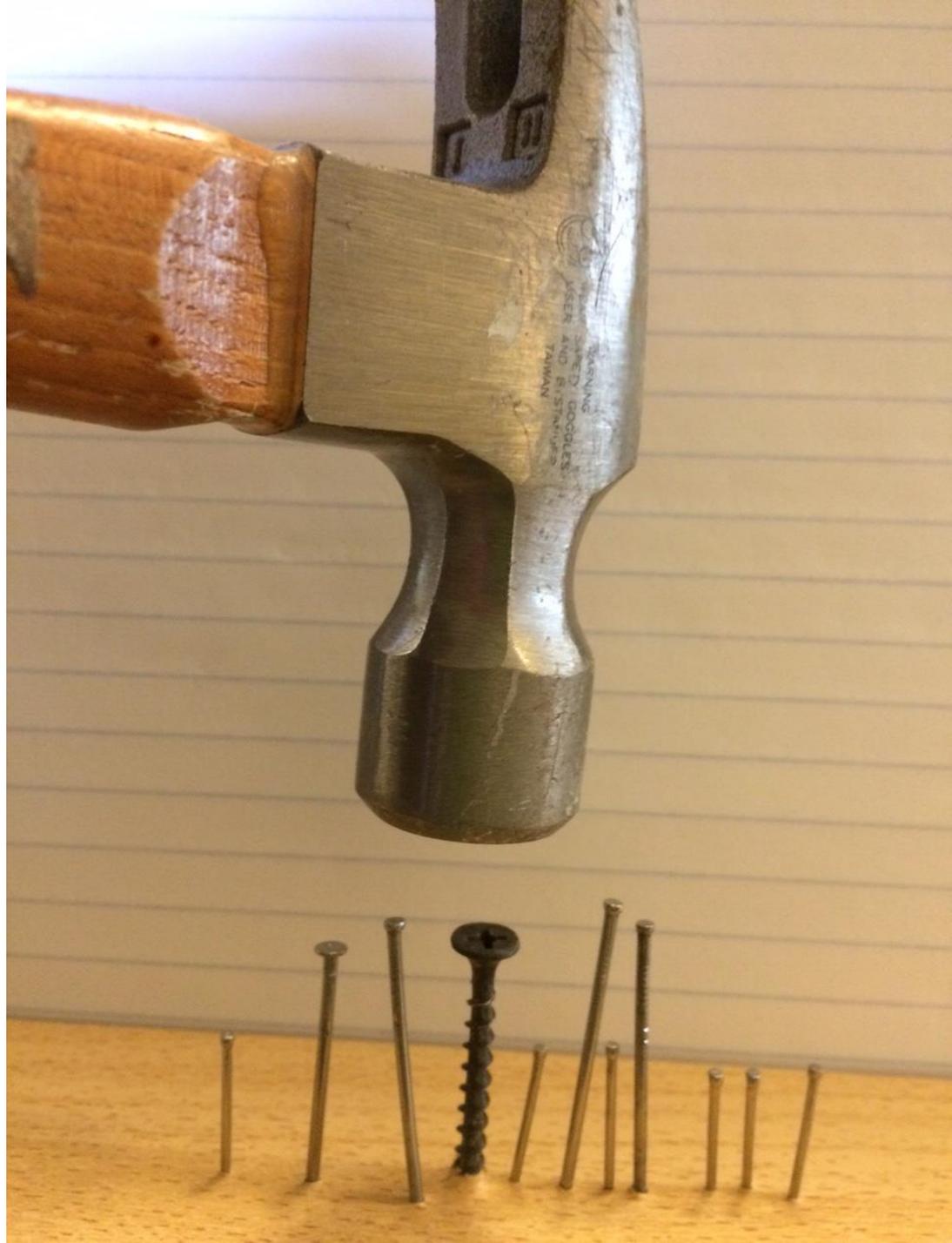
**My talk will focus on some of the statistical lesson learned from this challenge with an emphasis on Bayesian methodology.**

# This is how Debra Fischer portrayed the problem at the recent “Extreme Precision Radial Velocity” meeting at Yale (2015)



We have worked hard over the past 2 decades to improve RV precision. Now seem to be at a point where the largest terms in the error budget are similar magnitude. As we push down, we may encounter new surprises.

**Need to use the right tool**



**Debra Fischer**



**If we eliminate all other error sources except stellar noise, we won't see significant precision gains. We'll be... well... screwed.**

**A key challenge for statistical analysis is to separate planetary signals from stellar activity induced signals.**



**Debra Fischer**

# Stellar activity

Time Scale	Vel. noise	Type of activity	Partial solutions
~ 10 years	1 – 20 m/s	Magnetic cycle	correlation
10 – 50 d	few m/s	Active regions spots and plages	a) correlation b) FF' analysis + Gaussian process
15 min – 2 d	few m/s	Granulations	ave. 3x10 min/night reduce to ~ 0.5 m/s
~ 1 hr	< 1 m/s	Flares	
< 15 min	few m/s	Oscillations	ave. for 15 min reduce to ~ 0.2 m/s

# THE KEPLER FITTING CHALLENGE

(<https://rv-challenge.wikispaces.com>)  
(Google rv challenge wikispace)

IDEA Porto  
September 14

SIMULATION Stellar signals  
Planets

CHALLENGE Can we find  
the planets ?

# THE DATASET

(<https://rv-challenge.wikispaces.com>)

(Google rv challenge wikispace)

## DATA

15 data sets

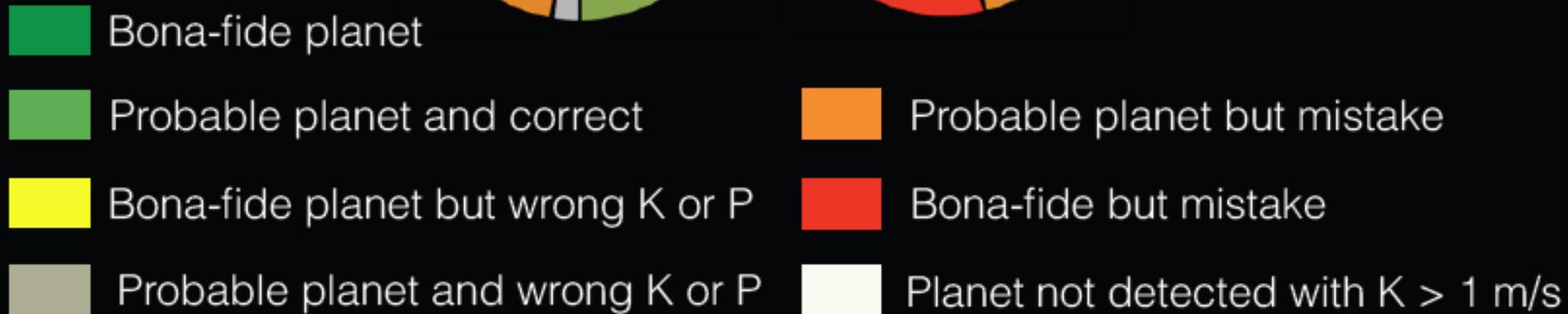
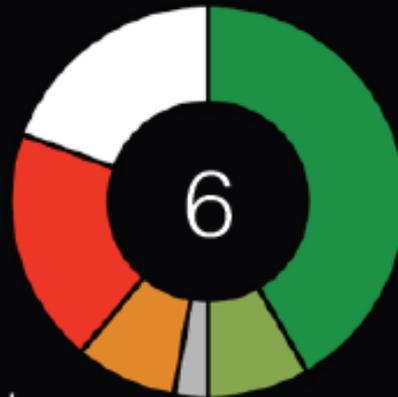
## OBSERVABLES

RV, BIS SPAN, FWHM,  $\text{Log}(R'_{hk})$

# ANALYSIS

- 24 groups where interested (~55 persons)
- Results from 8 groups
  - 5 groups analyzed the 15 RV curves
  - 2 groups only looked at the first 5 systems
  - 1 group analyzed only 2 systems

# PLANET DETECTION



# PODIUM

3RD



PHIL GREGORY

2ND



INAF TORINO

Mario Damasso, Aldo Bonomo, Paolo Giacobbe,  
Raphaëlle Haywood, Matteo Pinamonti, Alessandro Sozzetti

1ST



MIKKO TUOMI &  
GUILLEM ANGLADA ESCUDE

# TAKE AWAY MESSAGE

BEST TECHNIQUES  
MOVING AVERAGE

GAUSSIAN PROCESS

APODIZED KEPLERIAN

RED NOISE MODELS

BAYESIAN FRAMEWORK

# TAKE AWAY MESSAGE

## SIMULATED DATA

Are not bad, similarity with real datasets

$K > 1 \text{ M/S}$

90% of signal recovered

$K < 1 \text{ M/S}$

10% of signal recovered

Developed a new approach for  
the RV challenge based on  
Apodized Keplerian Models

# The Apodized Kepler (AK) model approach

The Kepler radial velocity parameter  $K$  is multiplied by an apodization term of the form

$$\exp\left[-\frac{(t_i - t_a)^2}{2\tau^2}\right]$$

Since a true planetary signal spans the duration of the data the apodization time,  $\tau$ , will be large while a stellar activity induced signal will generally have a small  $\tau$  value.

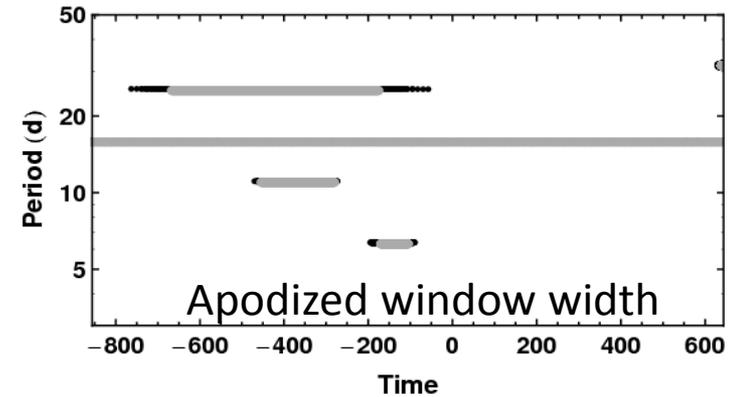
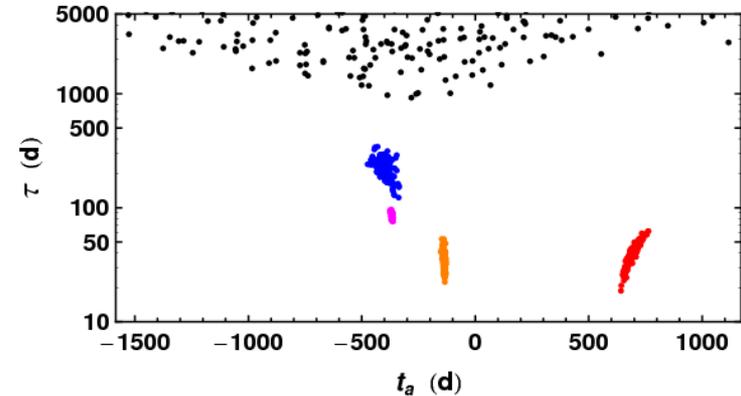
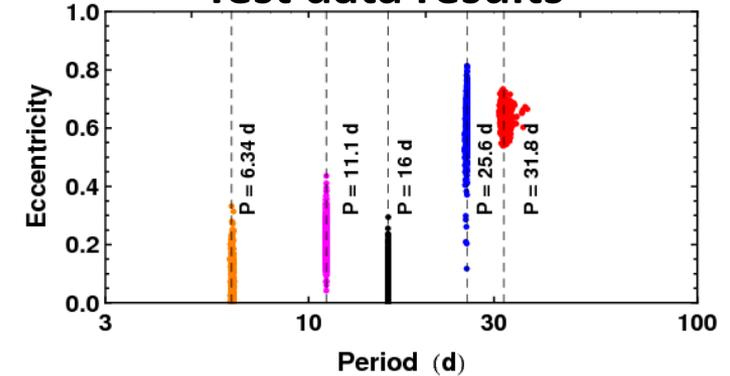
Each model also included a correlation term between RV and the stellar activity diagnostic  $\log(R'_{hk})$  and an extra Gaussian noise term.

The model parameters were explored using my fusion MCMC code and a differential version of the Generalized Lomb-Scargle algorithm.

The figure shows plots of MCMC parameter estimates for a 5 signal model fit to the test data, known to have one planet with a period of 16 d.

Phil Gregory (July 2015)  
University of British Columbia

## Test data results



# Radial velocity model for $m$ signals (planets + stellar activity) plus $\ln(R' hk)$ linear regression term

$$v(t_i) = V + \sum_{j=1}^m \left[ K_j \exp\left[-\frac{(t_i - t_{ja})^2}{2\tau_j^2}\right] \times (\cos\{\theta_j(t_i + \chi_j P_j) + \omega_j\} + e_j \cos \omega_j) \right] + \beta \times \ln[R' hk](t_i)$$

$m$  = the number of apodized Kepler (AK) signals in model.

Linear regression term  $\beta$  is just another fit parameter in the MCMC.

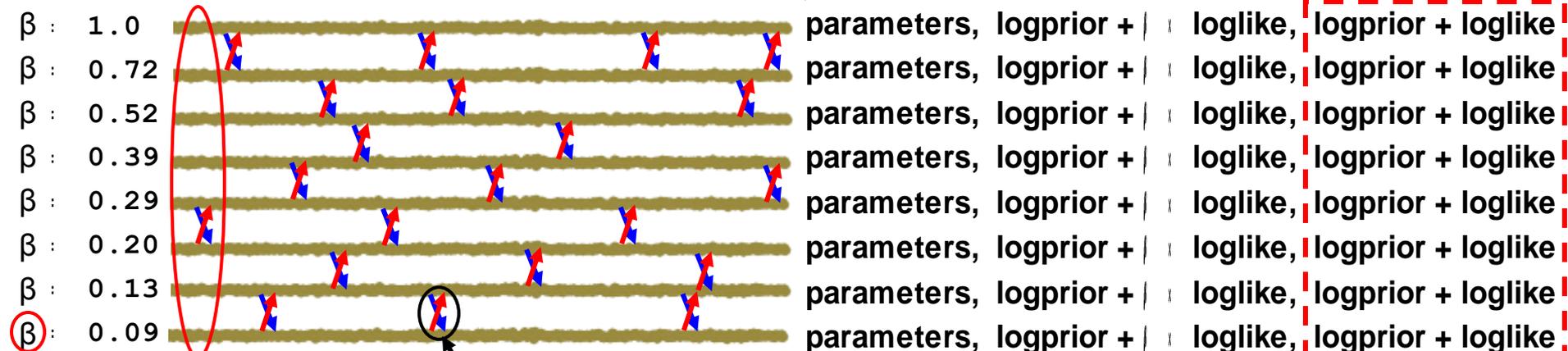
The AK models were explored using an automated fusion MCMC algorithm (FMCMC), a general purpose tool for nonlinear model fitting and regression analysis (Gregory 2013). The AK models combined with the FMCMC algorithm constitute a multi-signal AK periodogram.

**Current analysis assumes multiple independent Keplerian orbits which breaks down for near resonant orbits.**

# Fusion MCMC with Automatic proposal scheme

8 parallel tempering Metropolis chains

Output at each iteration



parameters, logprior +	loglike,	logprior + loglike
parameters, logprior +	loglike,	logprior + loglike
parameters, logprior +	loglike,	logprior + loglike
parameters, logprior +	loglike,	logprior + loglike
parameters, logprior +	loglike,	logprior + loglike
parameters, logprior +	loglike,	logprior + loglike
parameters, logprior +	loglike,	logprior + loglike
parameters, logprior +	loglike,	logprior + loglike



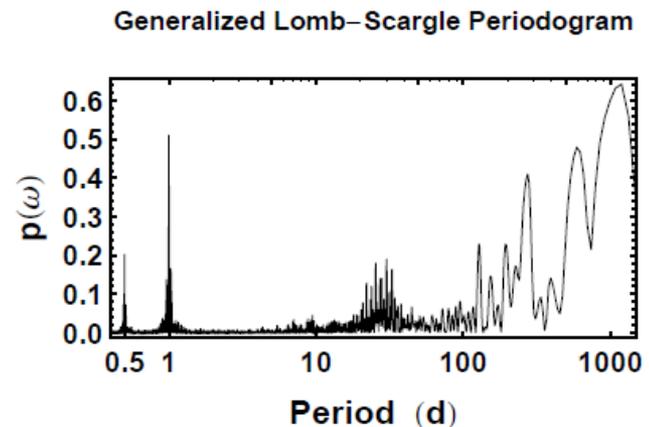
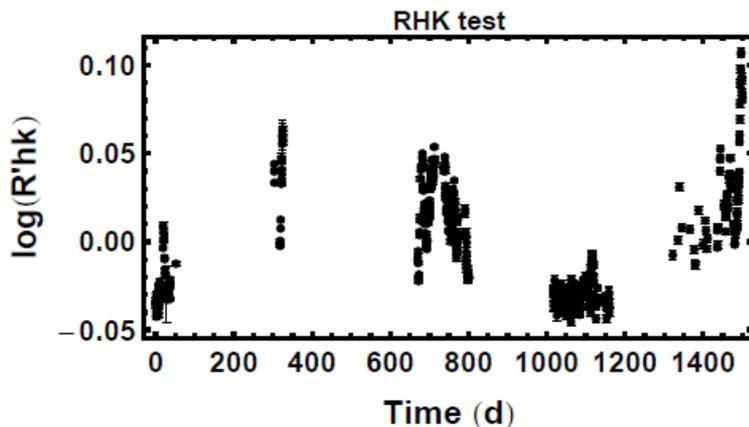
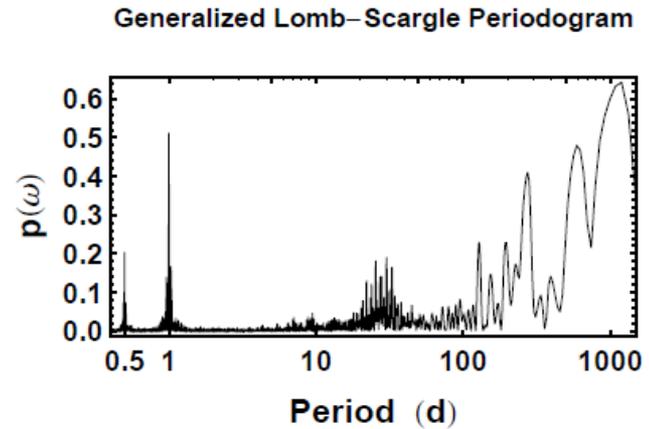
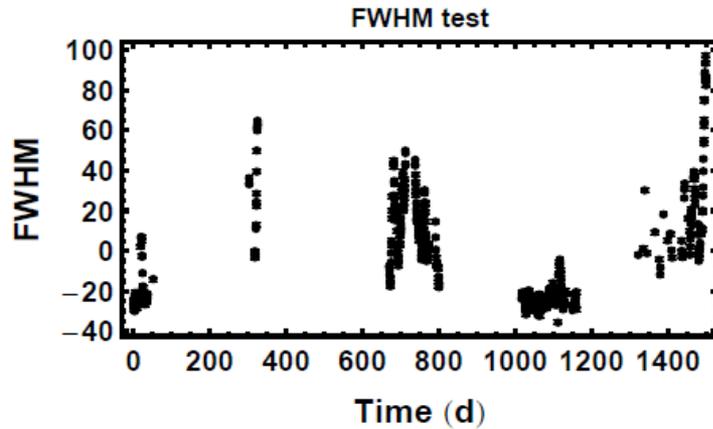
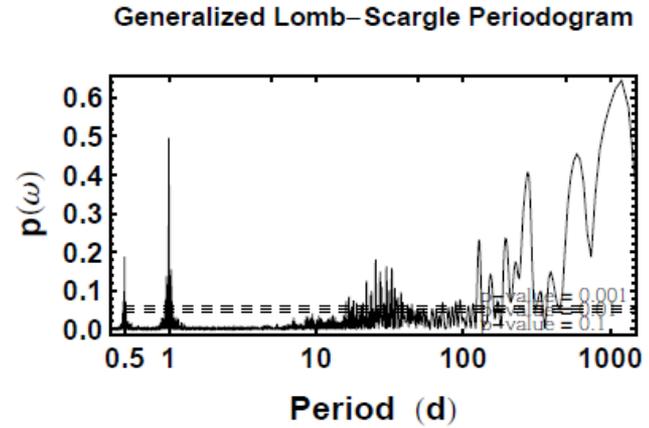
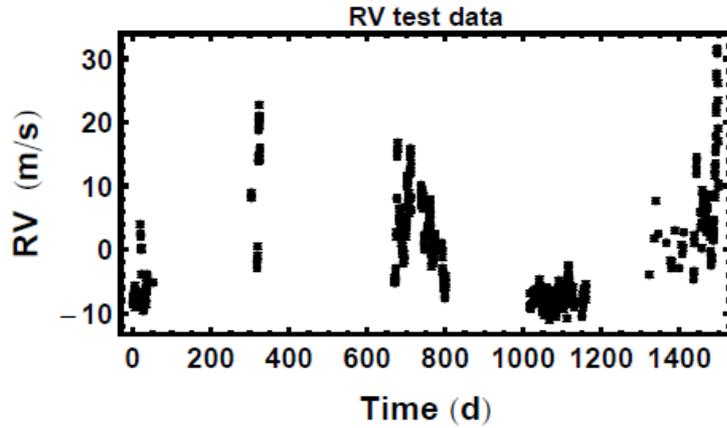
'I' proposals  
Independent Gaussian proposal scheme employed 50% of the time

'C' proposals  
Proposal distribution with built in param. correlations used 50% of the time

Genetic algorithm  
Every 40<sup>th</sup> iteration perform gene swapping operation to breed a more probable parameter set.

MCMC adaptive control system

# Raw RV and the FWHM and $\ln(R'hk)$ diagnostics for Test data set

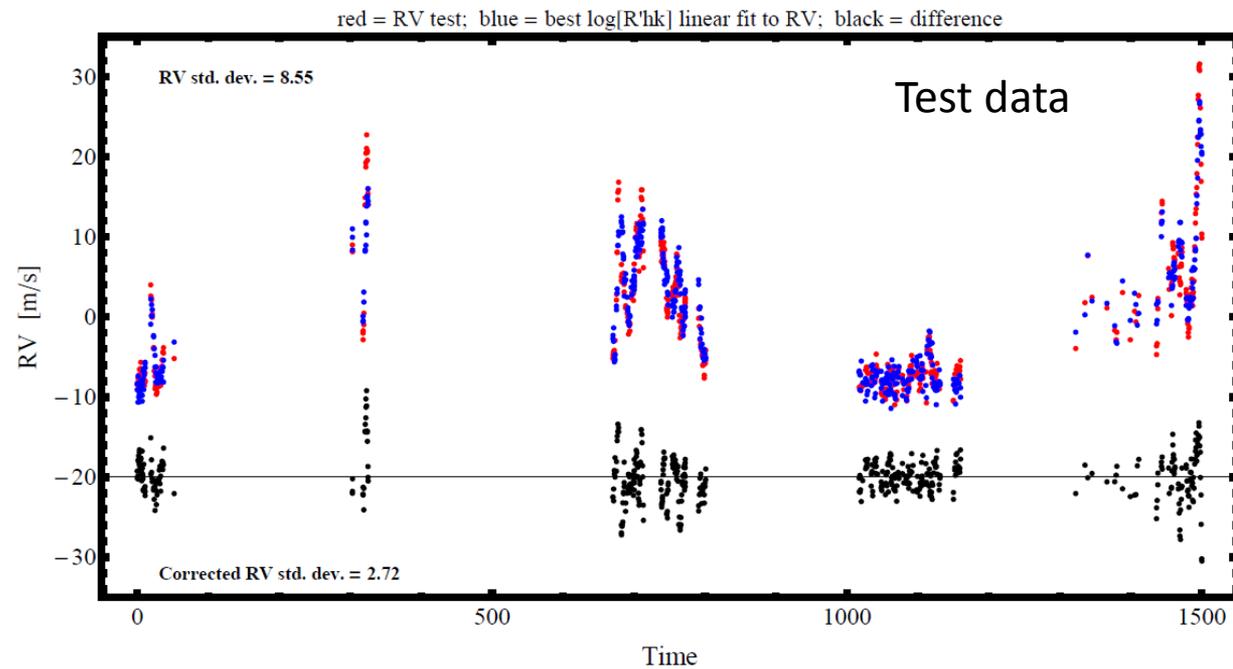


## Top panel

**Red** points shows the raw RV test data,

**Blue** points show the best  $\log(R'hk)$  linear regression fit to the RV data, and

**Black** points = the difference. (call this RV (rhk corrected))

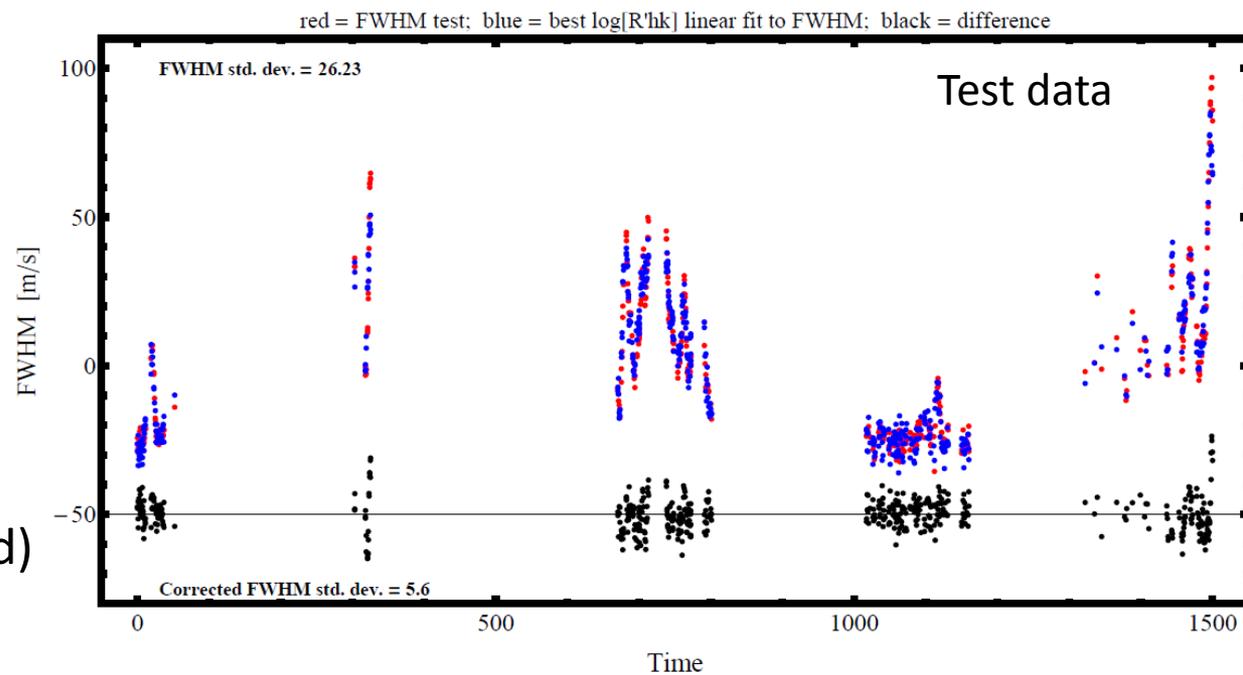


## Bottom panel

**Red** points shows the raw FWHM test data,

**Blue** points show the best  $\log(R'hk)$  linear regression fit to the FWHM data, and

**Black** points = the difference. (call this FWHM (rhk corrected) which is used as a control.)

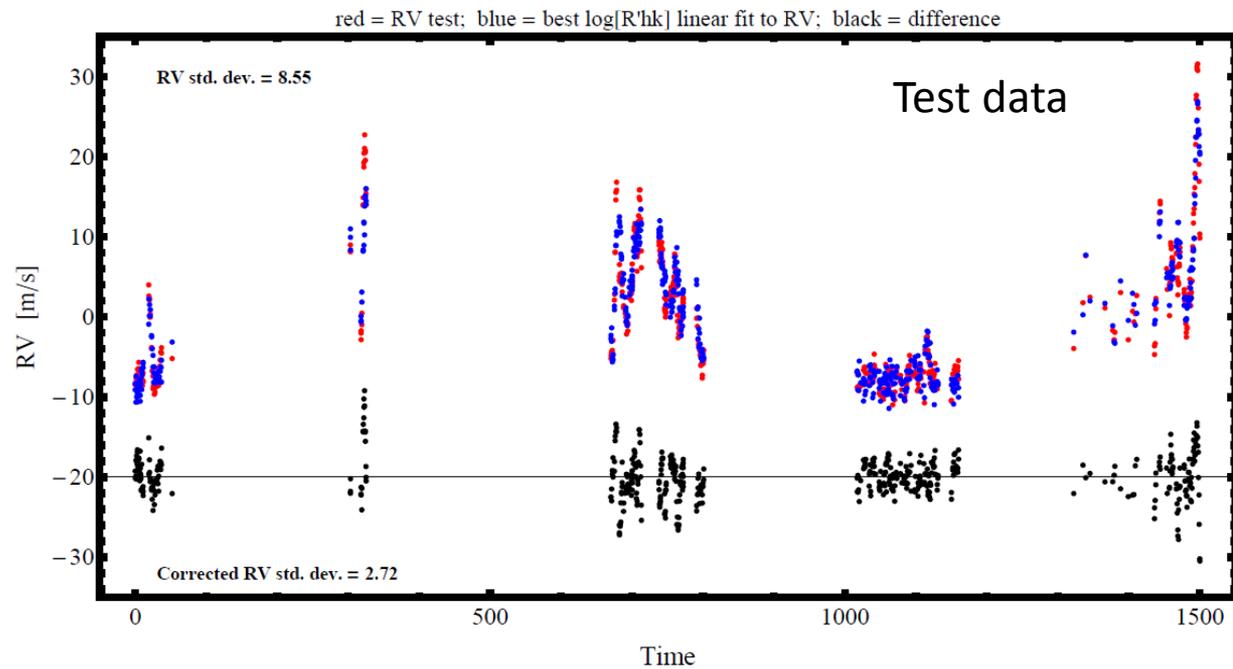


## Top panel

**Red** points shows the raw RV test data,

**Blue** points show the best  $\log(R'hk)$  linear fit to the RV data, and

**Black** points = the difference.  
(call this RV (rhk corrected))

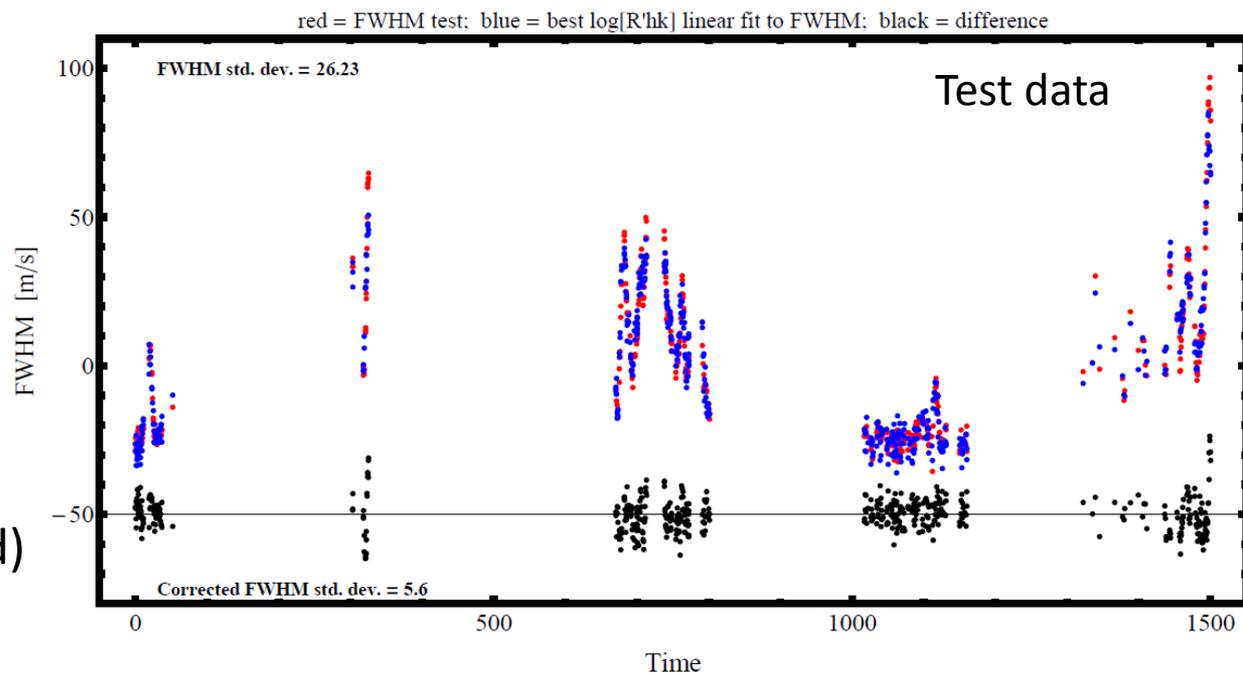


## Bottom panel

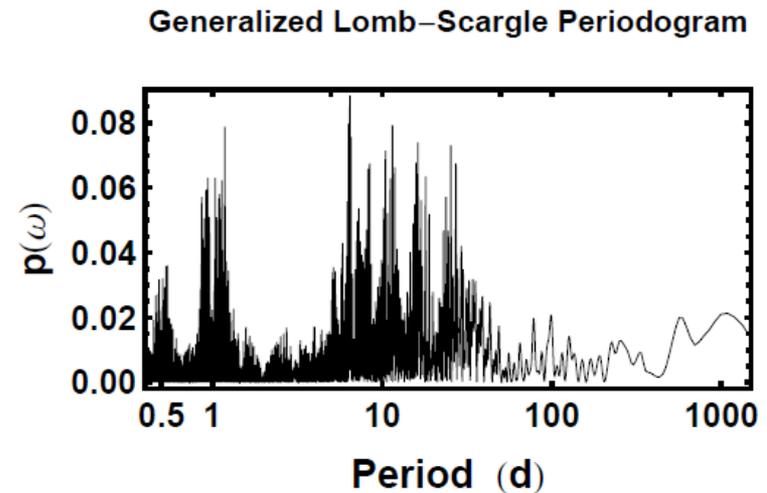
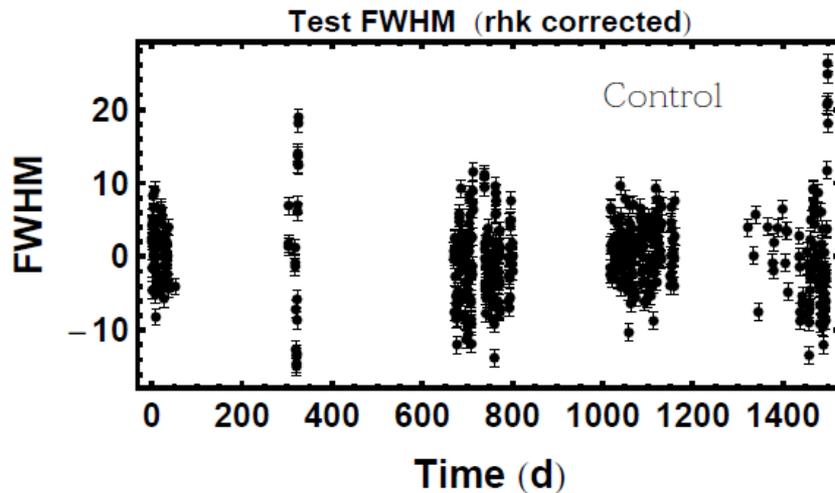
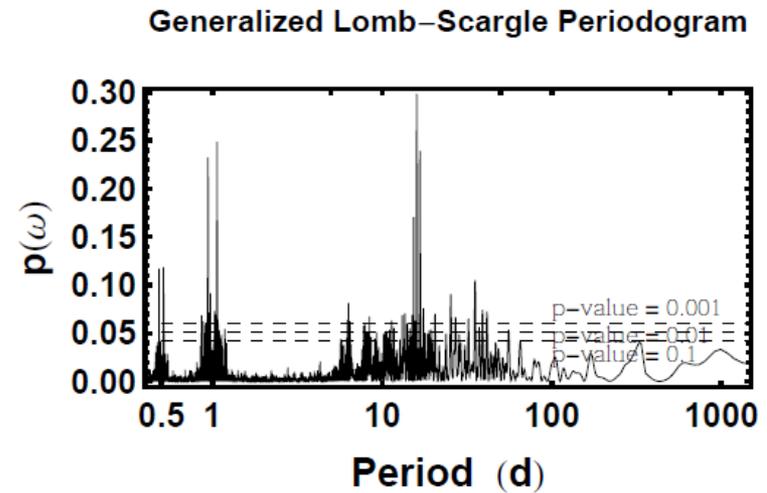
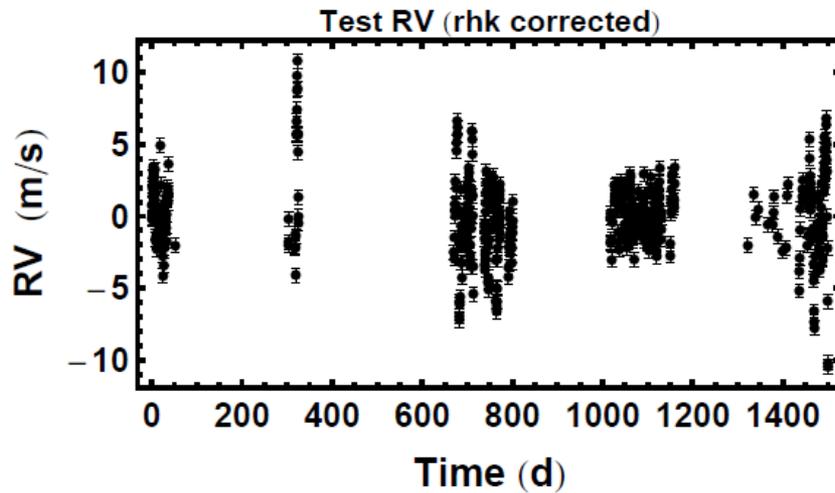
**Red** points shows the raw FWHM test data,

**Blue** points show the best  $\log(R'hk)$  linear fit to the FWHM data, and

**Black** points = the difference.  
(Call this FWHM (rhk corrected)  
which is used as a control.)



# Generalized Lomb-Scargle (GLS) periodogram of RV and FWHM (both rhk corrected).



The GLS periodogram measures the relative  $\chi^2$ -reduction,  $p(\omega)$ , as a function of frequency  $\omega$  and is normalised to unity by  $\chi^2_0$  (the  $\chi^2$  for the weighted mean of the data).

[New: a Bayesian version of GLS now available \(Mortier et al., arXiv:1412.0467.pdf\)](#)

# GLS Spectral difference

of significant  
spectral regions

Black = RV (rhk corr.)

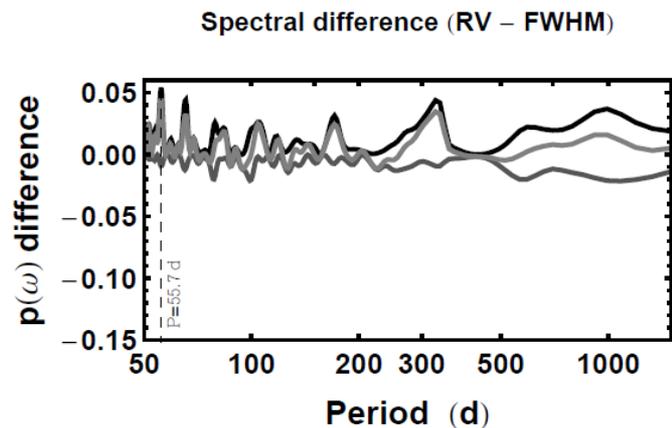
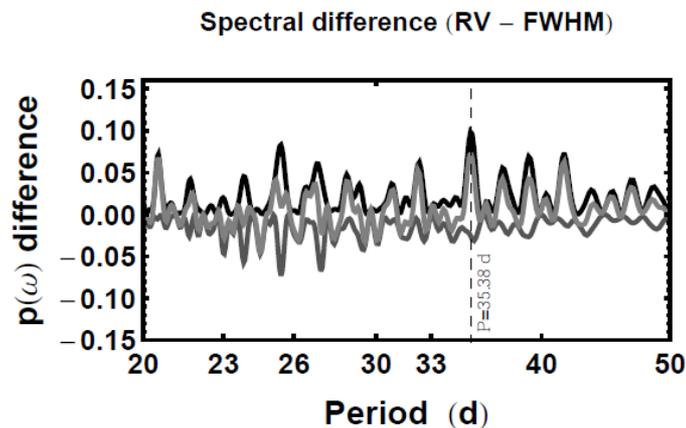
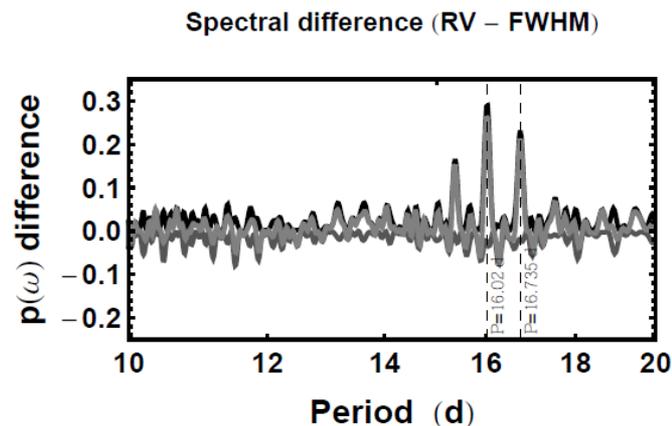
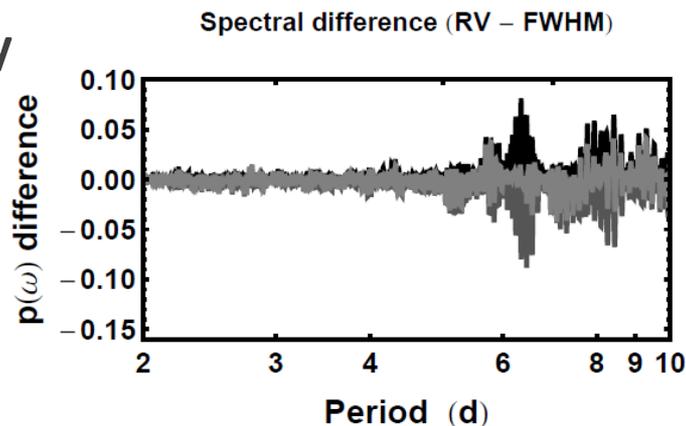
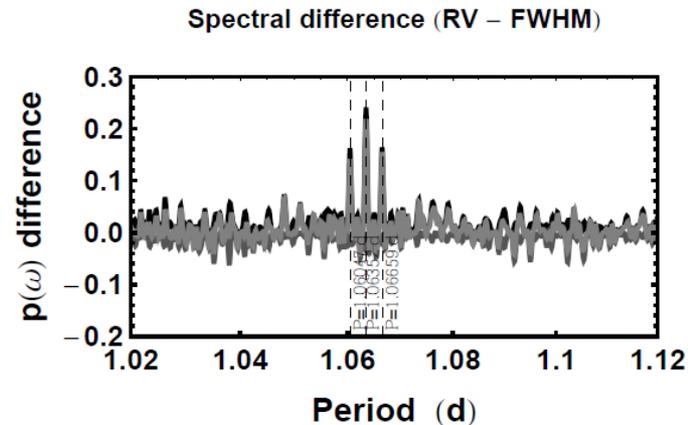
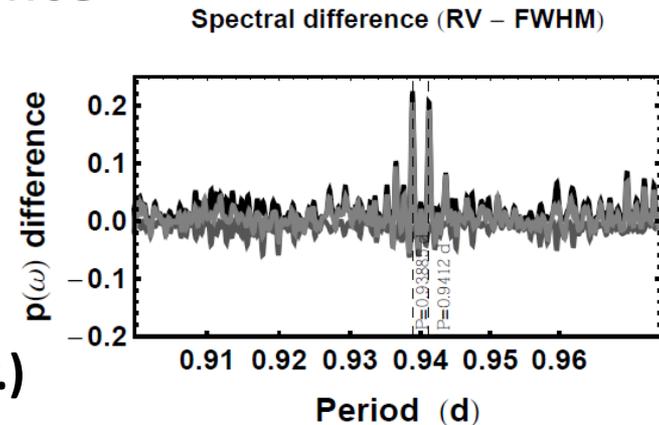
Gray = - FWHM (rhk corr.)

Light Gray = Black + Gray

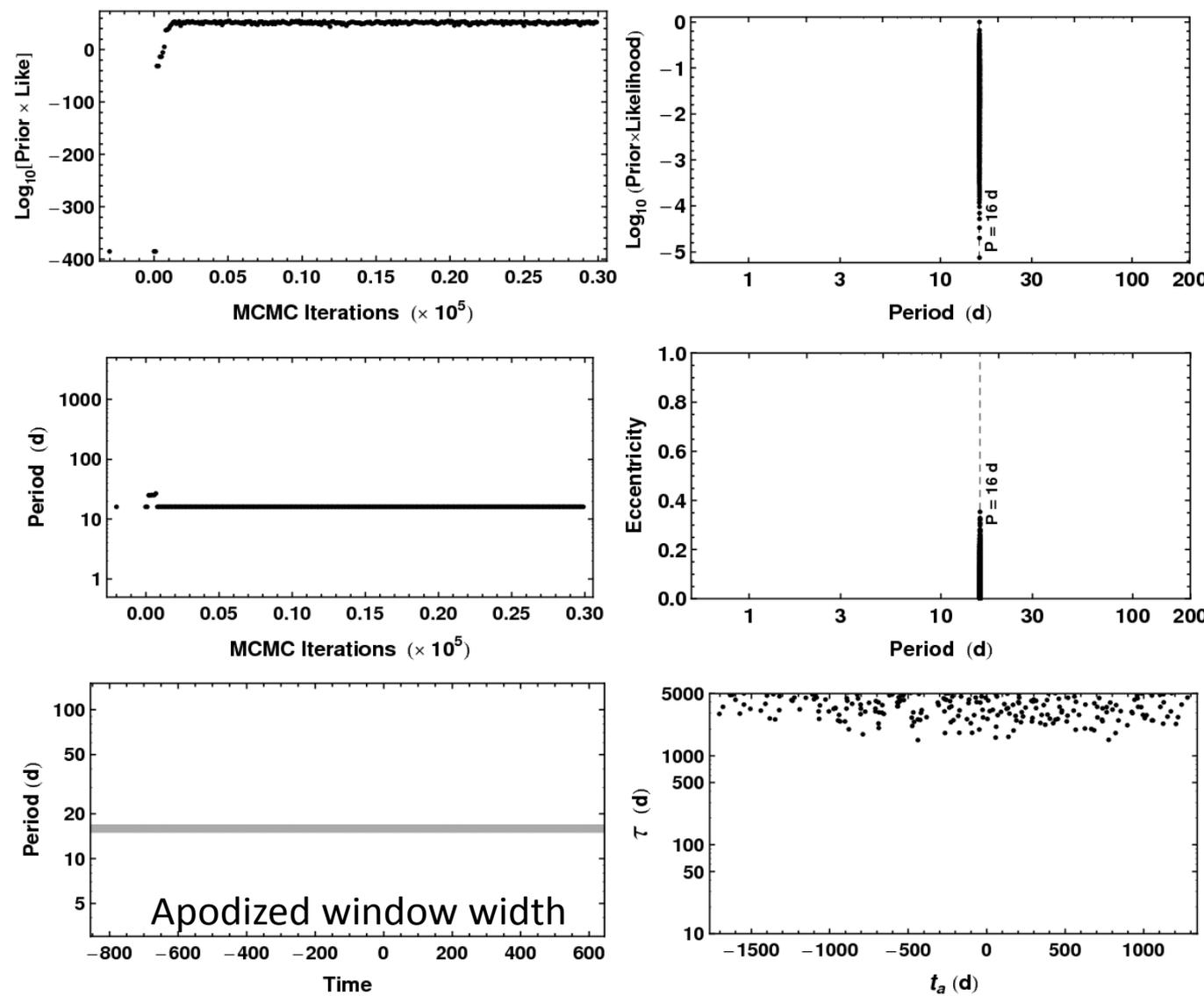
Signals in common to  
both indicate stellar  
activity. Gray trace acts  
as a control.

Dominant 16 d signal  
clearly visible. The next  
big peak on either side  
is a 1 yr alias.

Solar and sidereal day  
aliases seen near  
 $P = 0.94$  &  $1.06$  d.



# Model: 1 apodized Kepler signal + log(R' hk) regression fit (Test data)



The model parameters explored using fusion MCMC.

The figure shows Various plots of the MCMC parameter estimates.

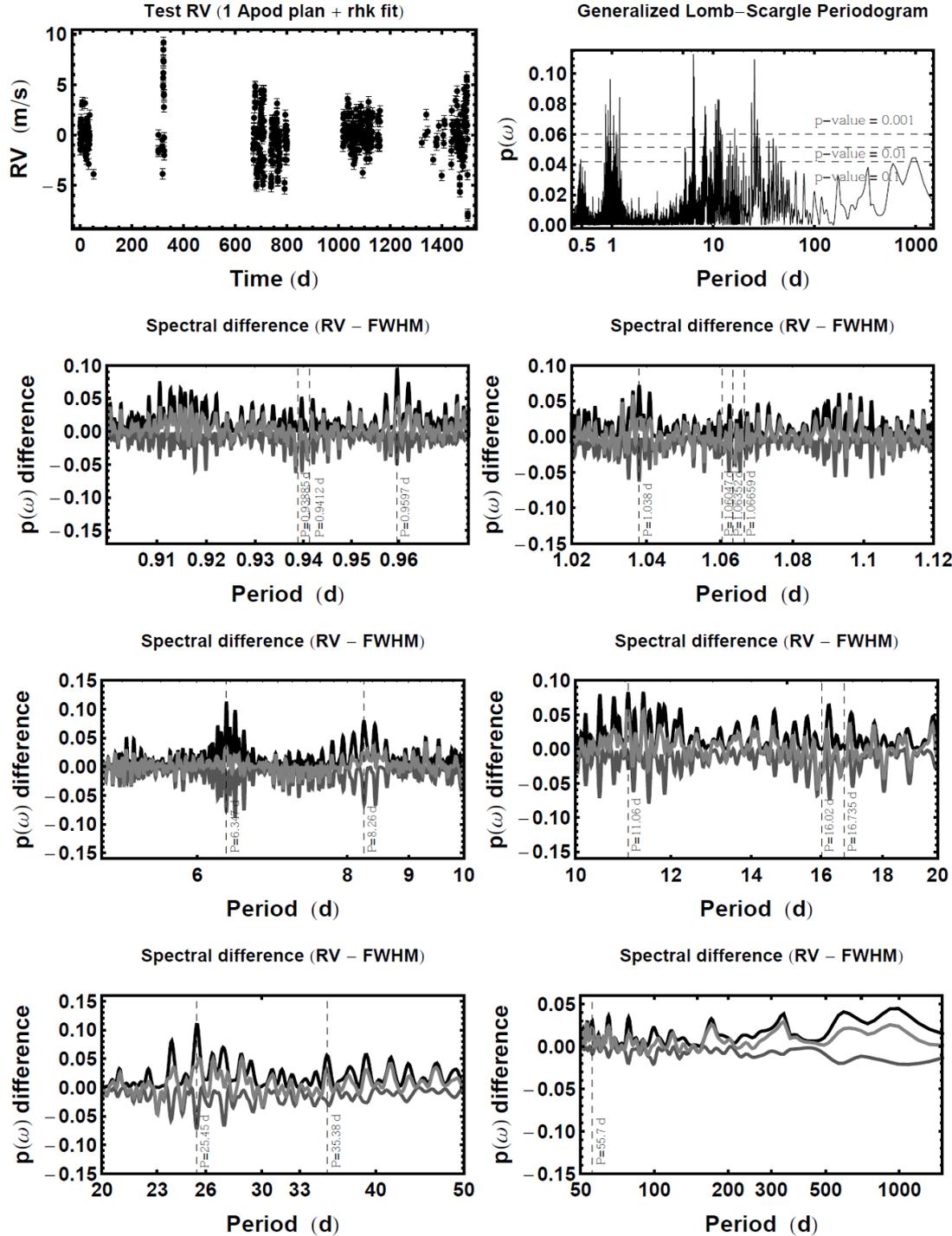
Lower left panel: apodization interval for each signal shown by gray trace for MAP values of  $\tau$  and  $t_a$ . Lower right panel: apodization time constant,  $\tau$ , versus  $t_a$  for the 16 d signal.

# GLS & Spectral difference of residuals from 1 apodized Kepler + rkh fit

Dominant 16 d signal and  
aliases have been removed  
including those near  
 $P = 0.94$  &  $1.06$  d.

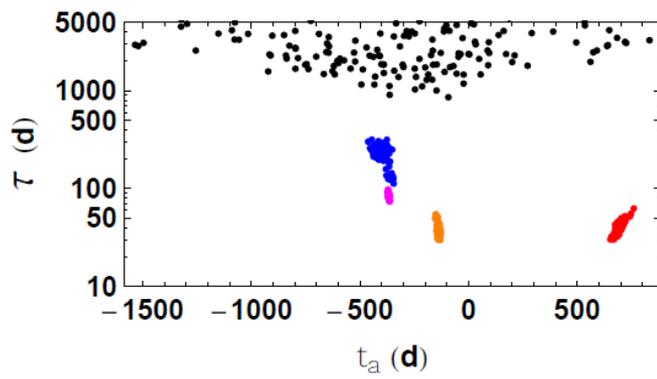
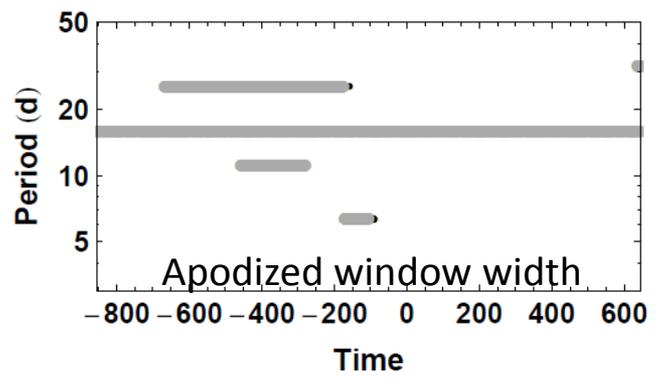
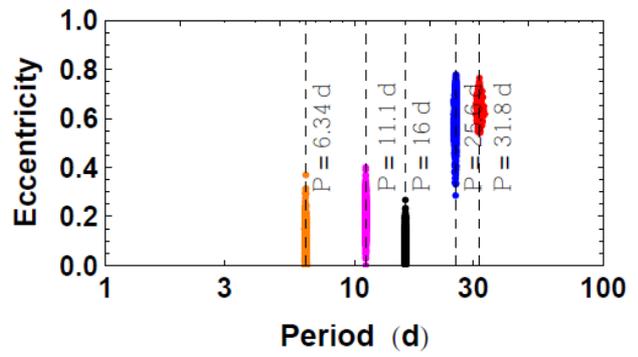
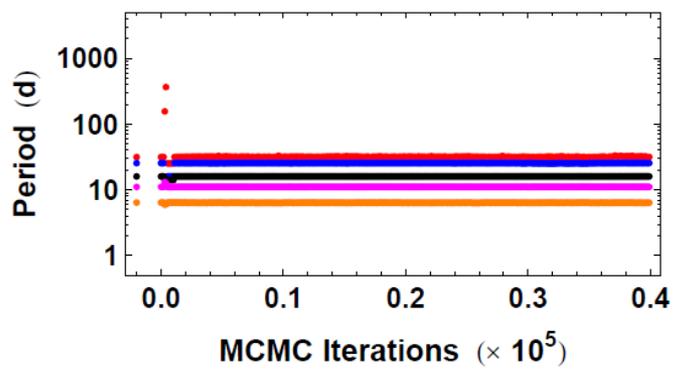
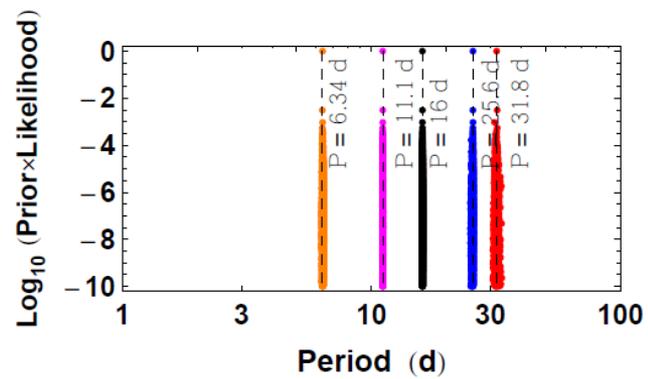
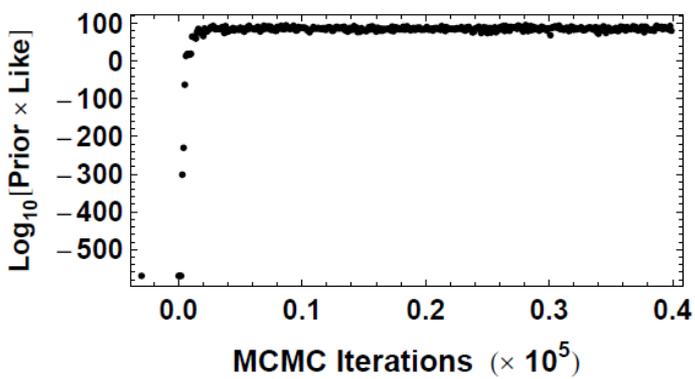
Largest GLS residual peak  
at  $P = 6.3$  d  
has  $p\text{-value} \ll 0.001$

Note: the FWHM control  
indicates 6.3 d is stellar activity



# Model: 5 apodized Kepler signals + log(R'hk) regression fit (Test data)

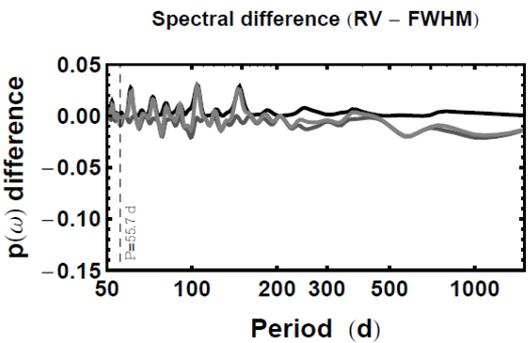
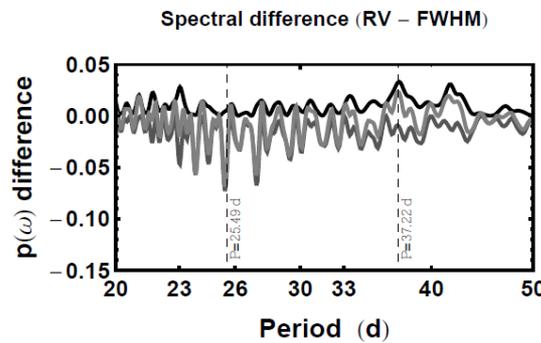
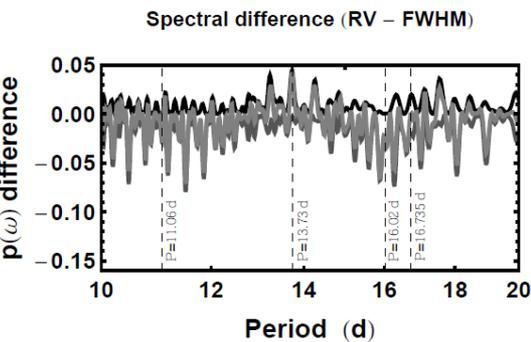
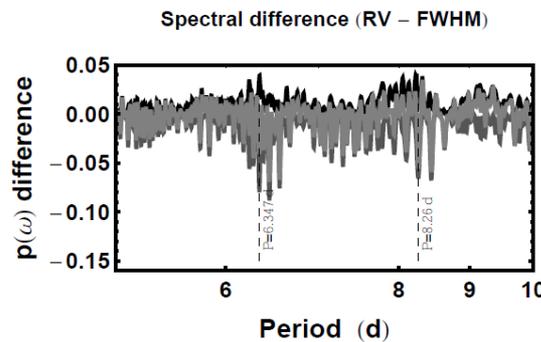
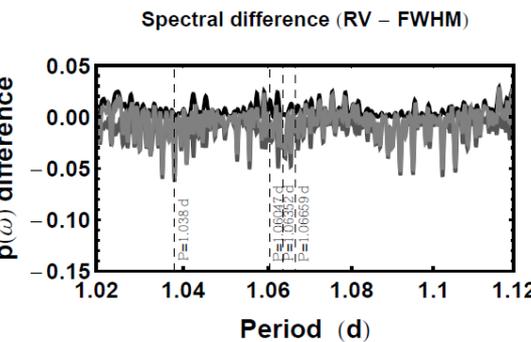
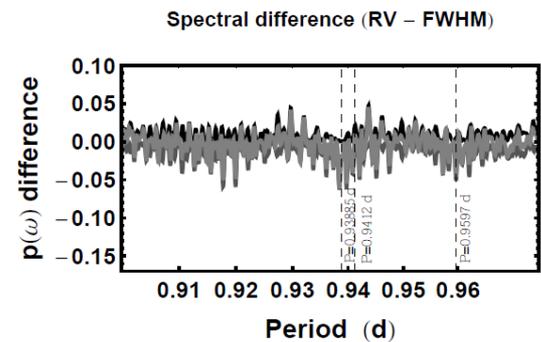
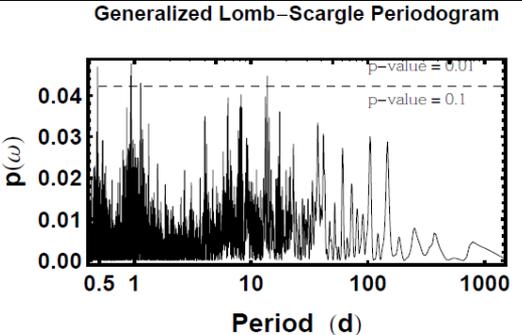
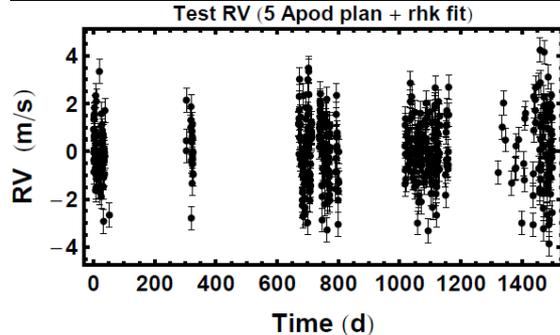
**Only the 16 d signal has an apodization time constant  $\tau$  (d) consistent with a planet.**



Free *Mathematica* fusion MCMC code for simple 2 planet Kepler model and program details available under resources at: <http://www.cambridge.org/pl/academic/subjects/statistics-probability/statistics-physical-sciences-and-engineering/bayesian-logical-data-analysis-physical-sciences-comparative-approach-mathematica-support>

# GLS & Spectral difference of residuals from 5 apodized Kepler + rhk fit

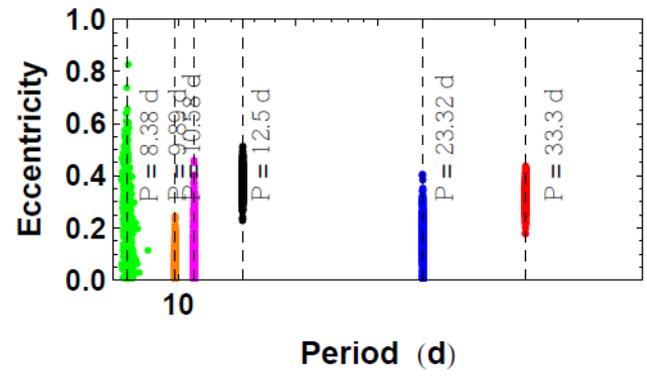
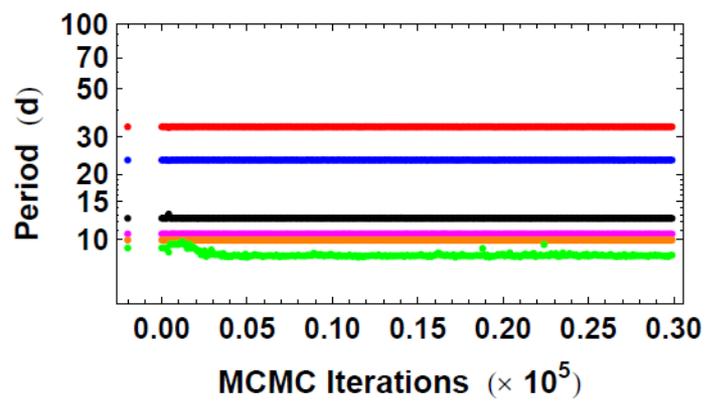
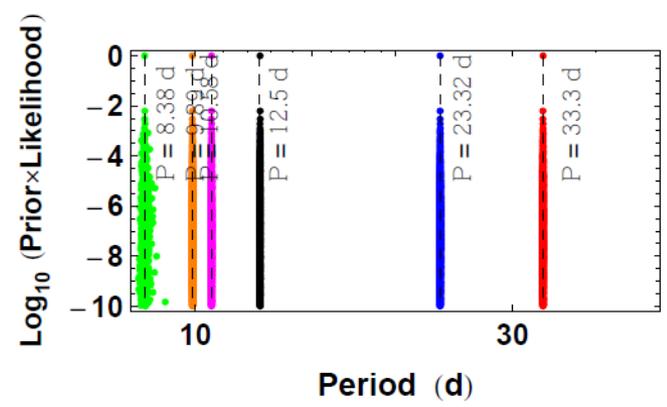
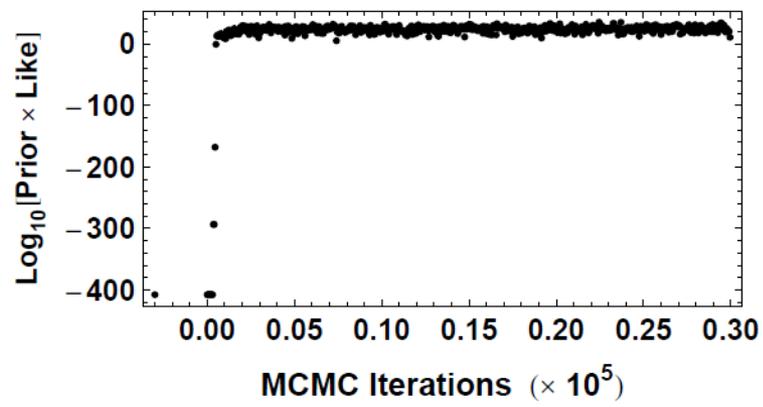
Largest GLS residual peak  
has p-value between  
**0.1 & 0.01**



# RV 1 Results

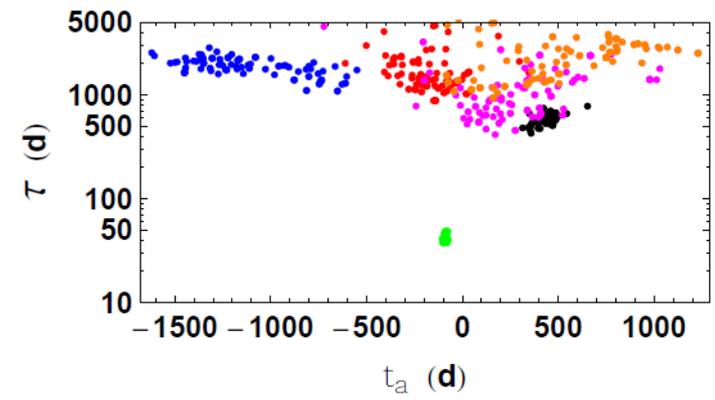
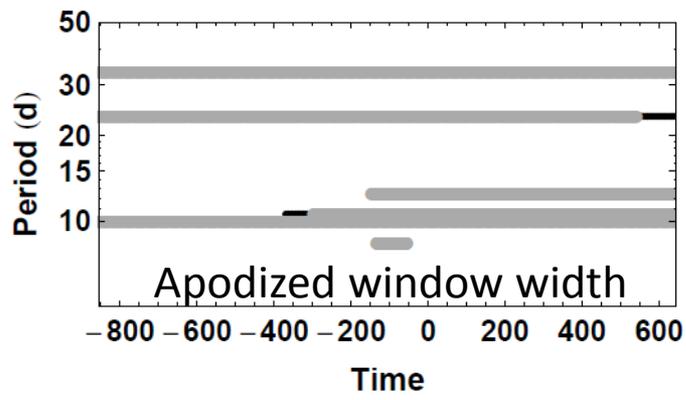
# RV 1 Model: 6 apodized Kepler signals + log(R'hk) regression fit

**Results indicate**  
 3 planets with  
 $P = 9.89, 23.4, 33.3$  d  
 +  
 3 stellar activity (SA)  
 signals

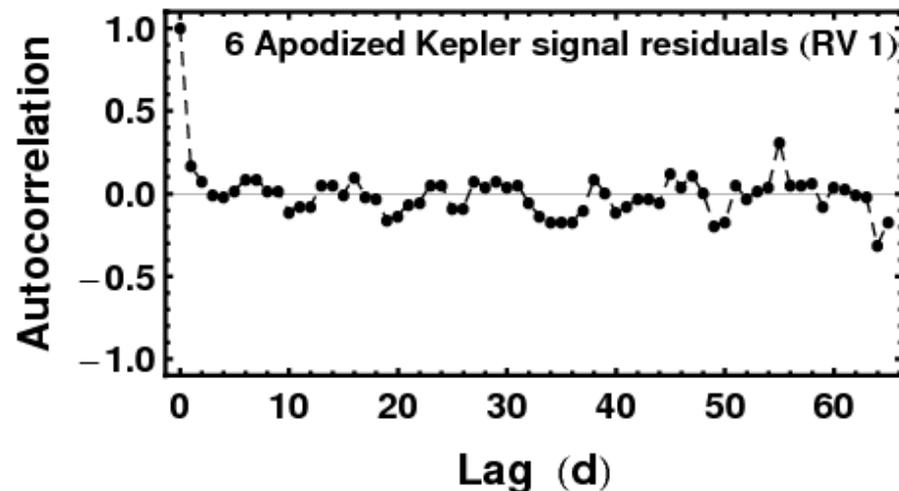
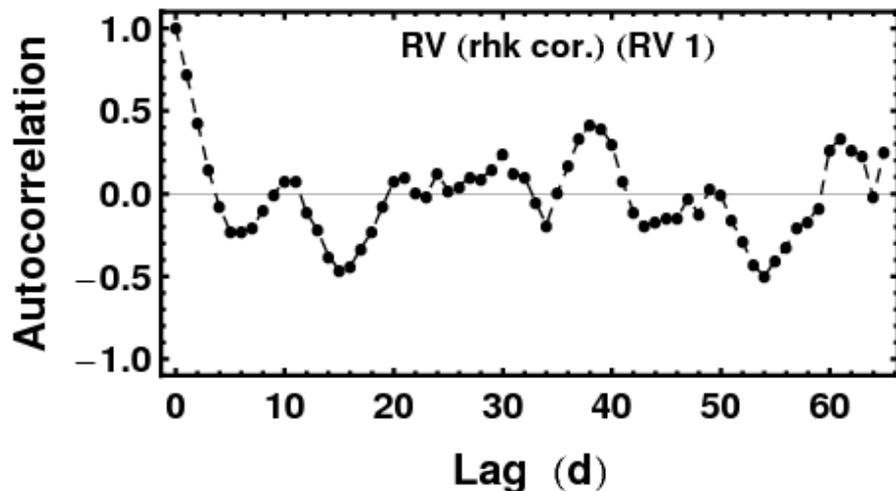
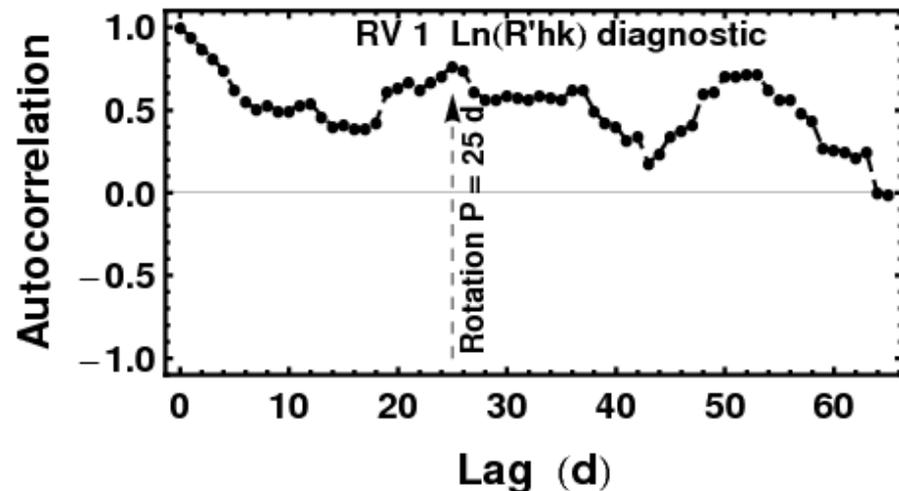
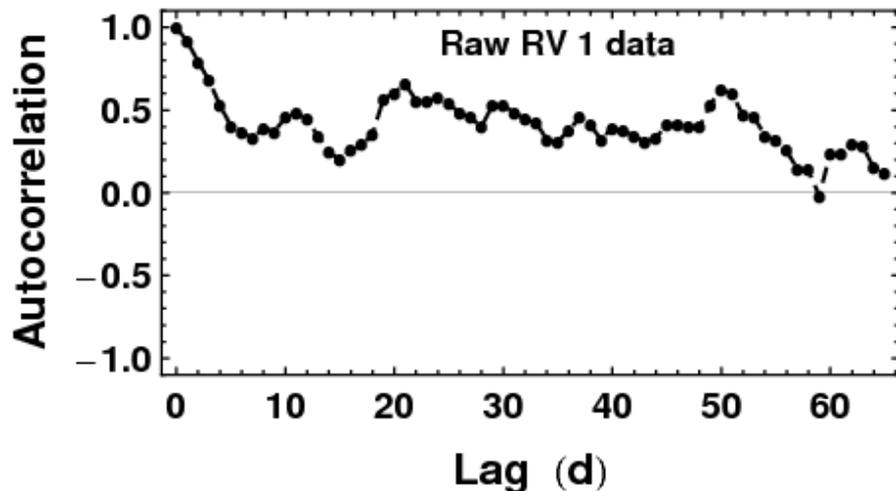


## True planets signals

P (d)	ecc	K (m/s)
9.89	0.1	1.45
23.4	0.12	1.67
33.3	0.08	2.05
112.5	0.21	0.38
273.2	0.16	0.22



# Correlated Noise



By the time the 6 apodized Kepler signals and Log(R'hk) regression are removed, the autocorrelation of the residuals is looking close to white noise.

# RV 2 Results

# RV 2 Model: 8 apodized Kepler signals + log(R'hk) regression fit

## Results indicate

3 planets

$P = 3.77, 10.6, 75.5$  d

(10.6 d listed as a probable due to many nearby SA signals.)

+

5 SA signals

## True planets signals

$P$  (d)     $ecc$      $K$  (m/s)

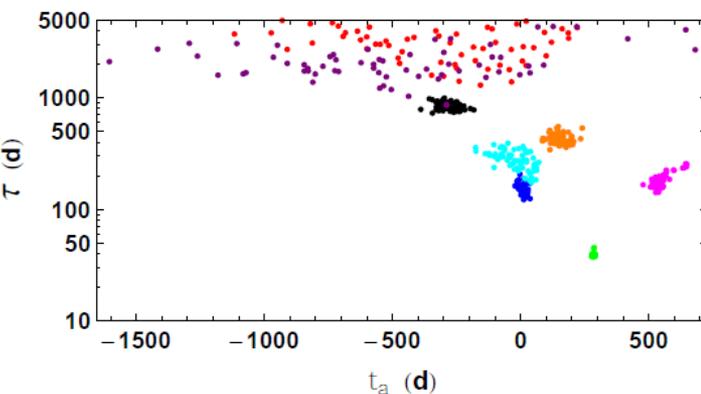
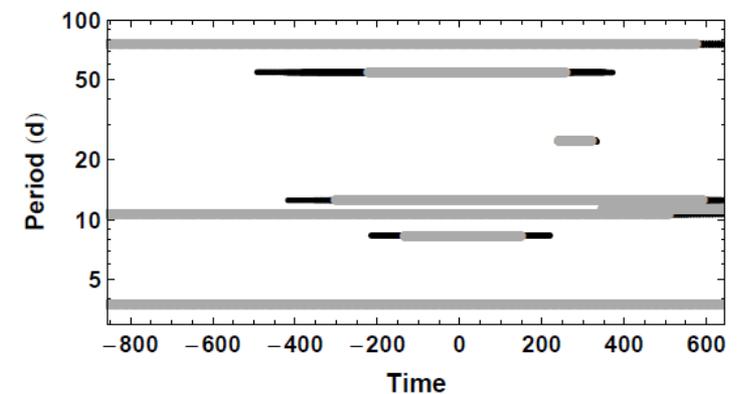
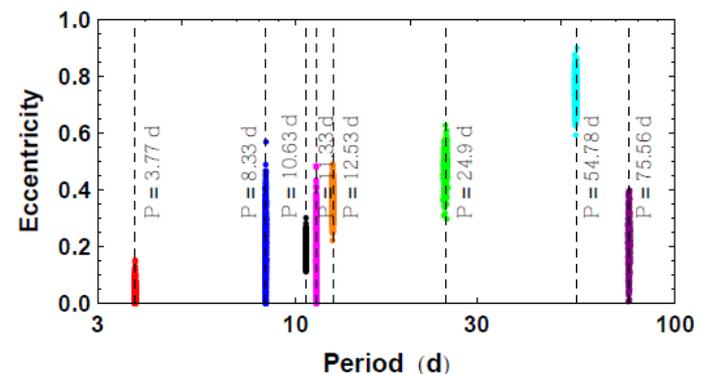
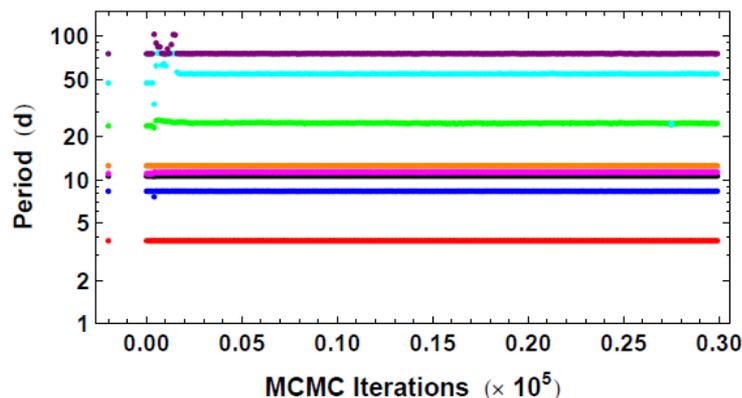
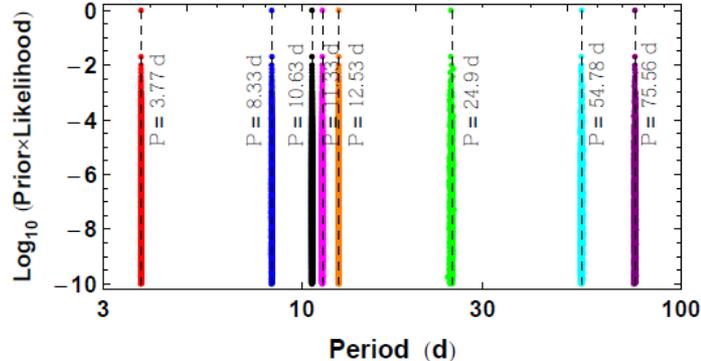
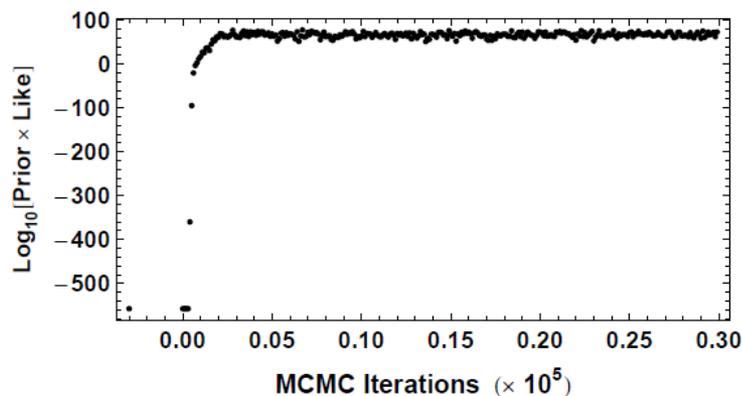
3.77    0.05    2.75

5.79    0.11    0.27

10.6    0.14    2.85

20.2    0.08    0.34

75.3    0.19    1.35



# RV 3 Results

# RV 3 Models

# 6 apodized Kepler signals

## Results indicate

3 planets with  
 $P = 17, 48.8, \sim 1100d$   
 (17 d listed probable  
 due to weak signature  
 in FWHM control)  
 (1100 d credited as  
 harmonic of 2315)

+

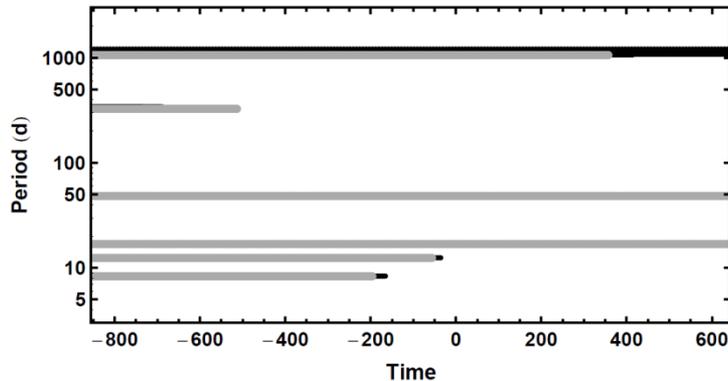
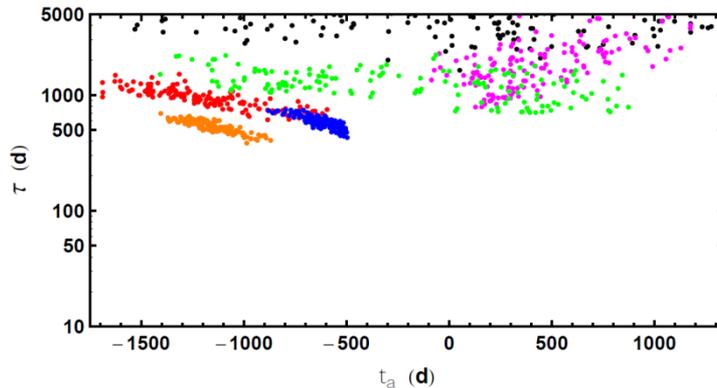
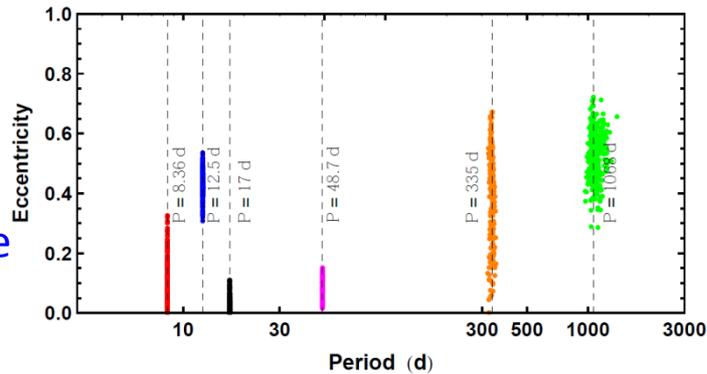
3 SA signals

## True planets signals

P (d)    ecc    K (m/s)

---

1.12	0.0	0.96
17.0	0.15	3.68
26.3	0.08	0.38
48.7	0.06	5.14
201.5	0.2	0.42
596	0.13	1.91
2315	0.15	3.87

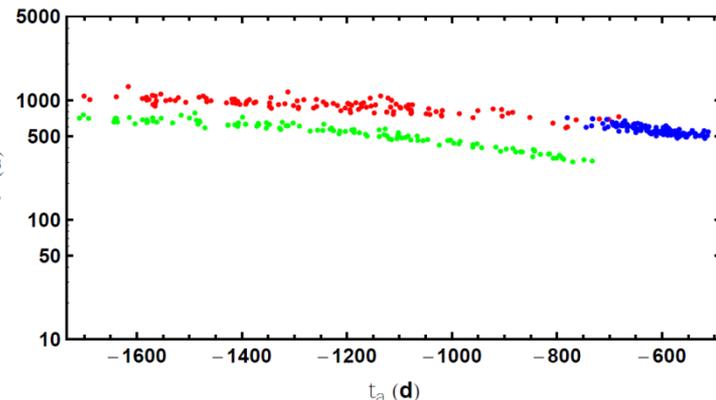
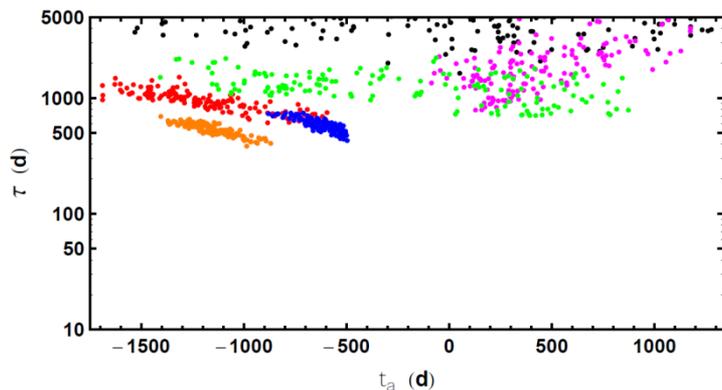
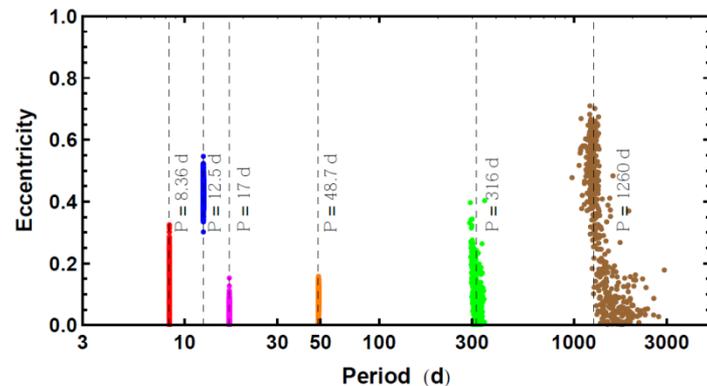
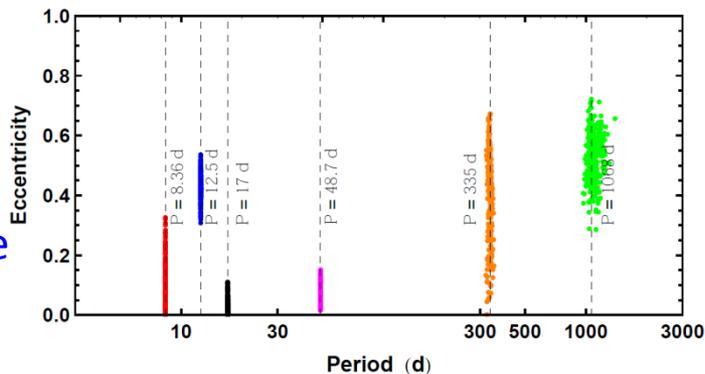


# RV 3 Models

# 6 apodized Kepler signals

# 3 apodized Kepler signals + 3 straight Kepler signals

**Results indicate**  
 3 planets with  
 $P = 17, 48.8, \sim 1100d$   
 (17 d listed probable  
 due to weak signature  
 in FWHM control)  
 (1100 d credited as  
 harmonic of 2315)  
 +  
 3 SA signals

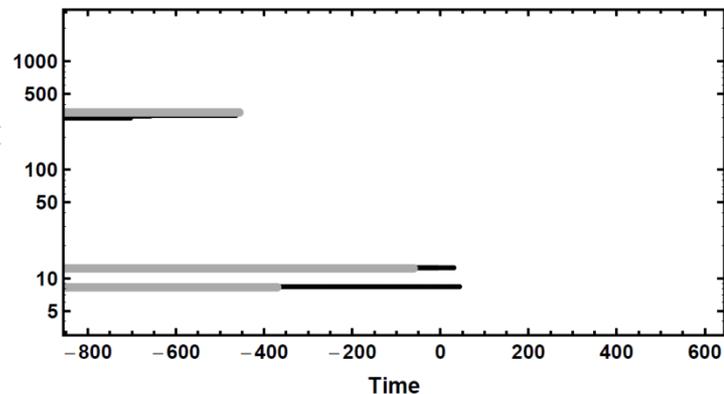
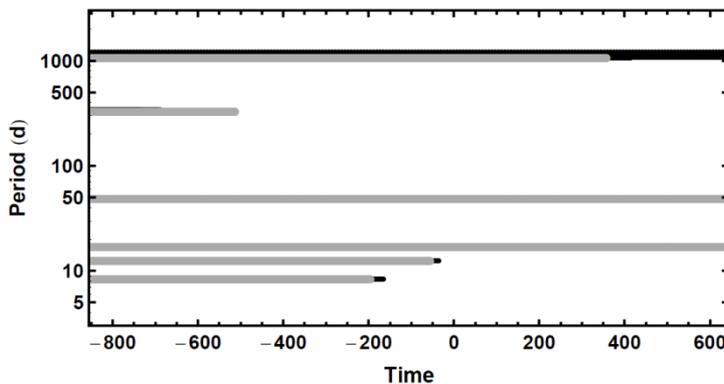


## True planets signals

P (d)    ecc    K (m/s)

---

1.12	0.0	0.96
17.0	0.15	3.68
26.3	0.08	0.38
48.7	0.06	5.14
201.5	0.2	0.42
596	0.13	1.91
2315	0.15	3.87



# RV 4 Results

# RV 4 Model: 8 apodized Kepler signals + log(R'hk) regression fit

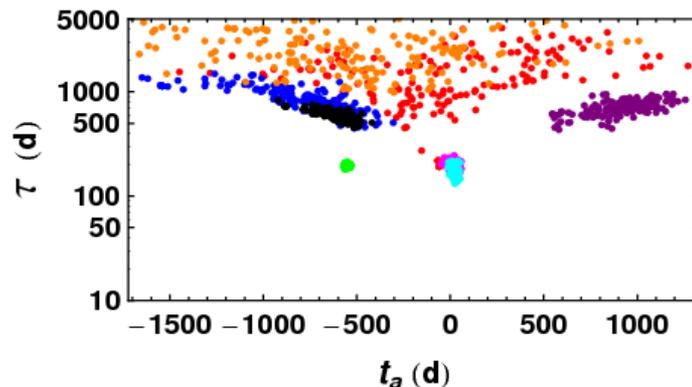
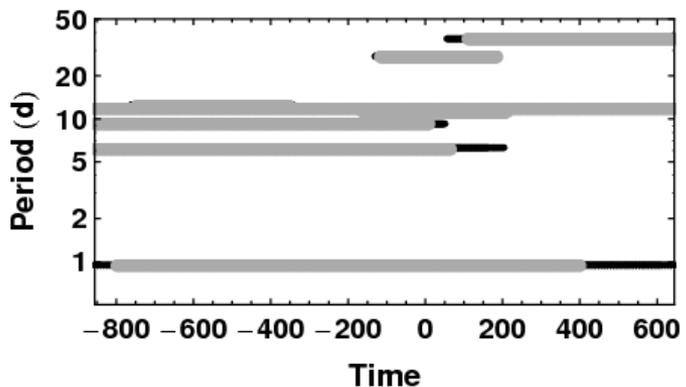
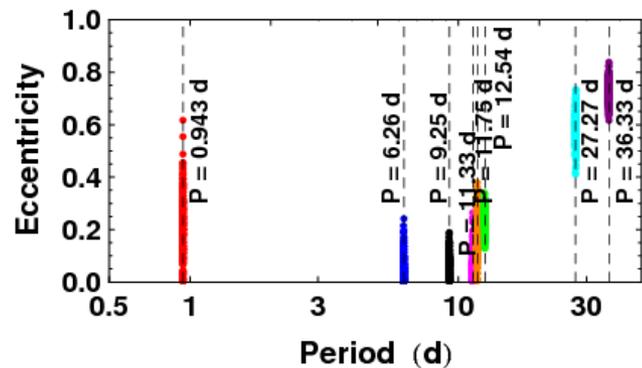
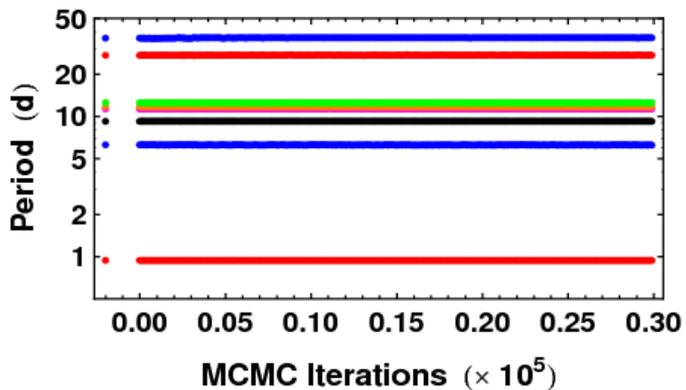
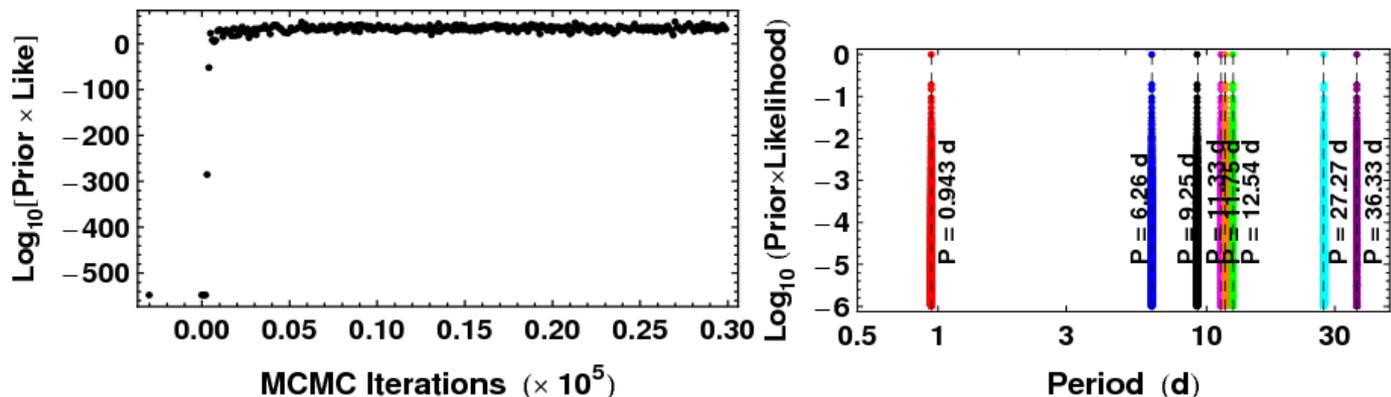
**No definite planets**

Possible planets at  $P = 0.946$  &  $11.75$  d based on apodization.

Bayes factor finds against a real  $P = 0.946$  d planet.

$P = 11.75$  only a possible because of weak FWHM

Control counterpart, see differential GLS periodogram.

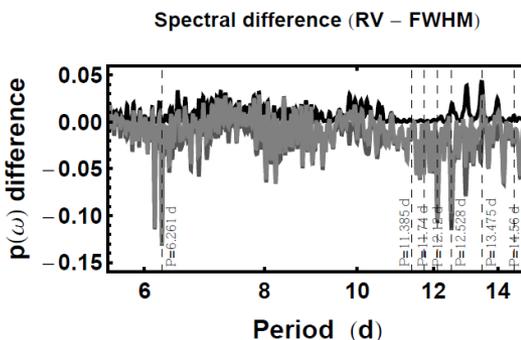
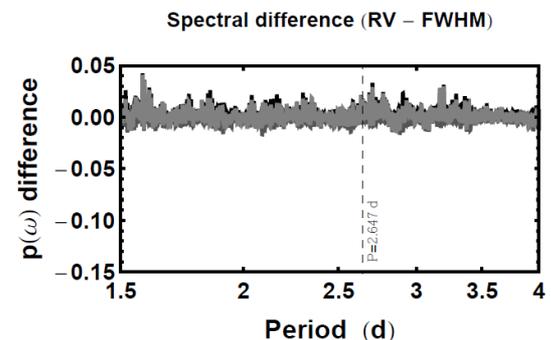
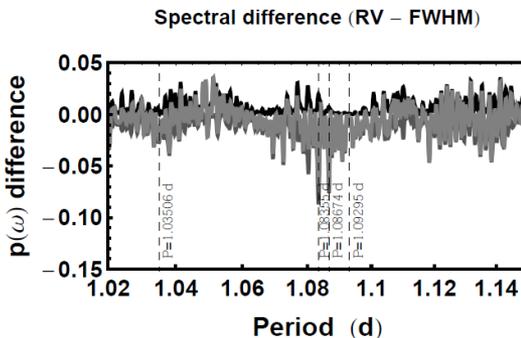
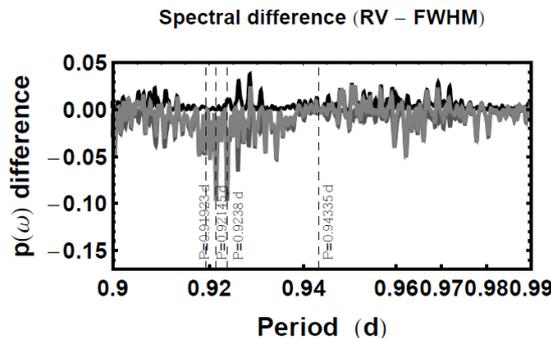
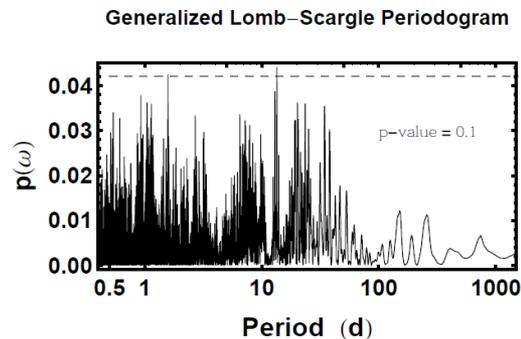
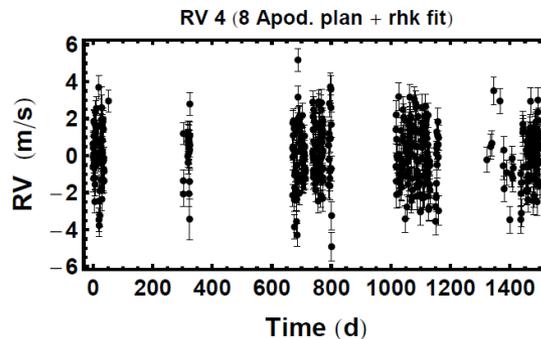


**True planets signals**

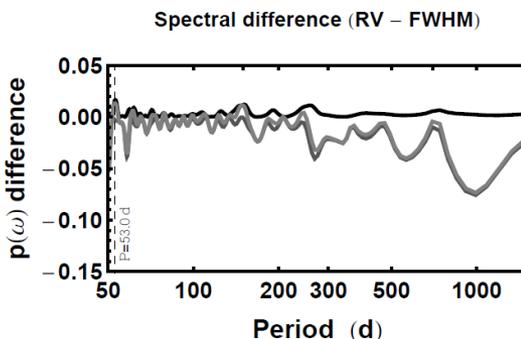
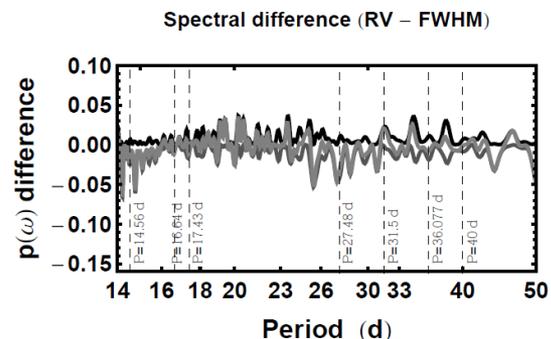
P (d)    ecc    K (m/s)

-----  
None

# GLS & Spectral difference of residuals from 8 apodized Kepler + rhk fit



Significant power at  $P = 11.75$  d  
in FWHM (rhk corr.) control



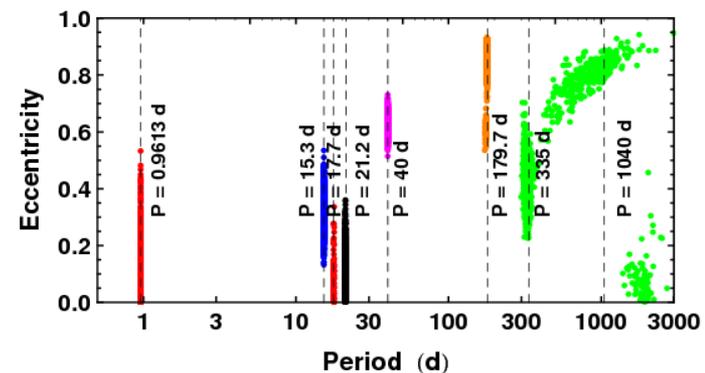
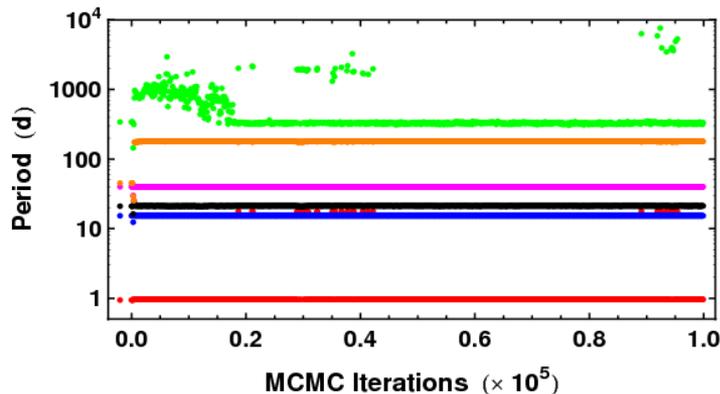
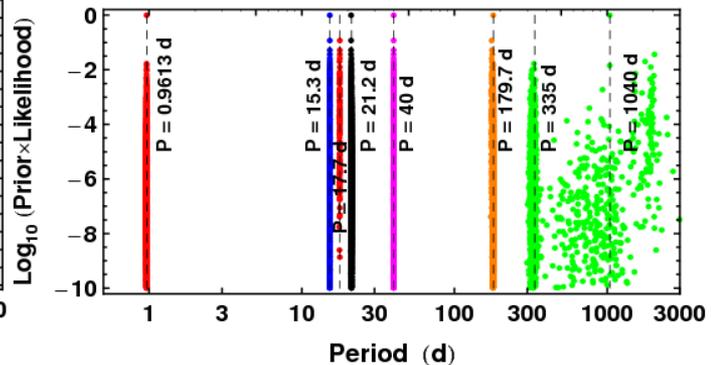
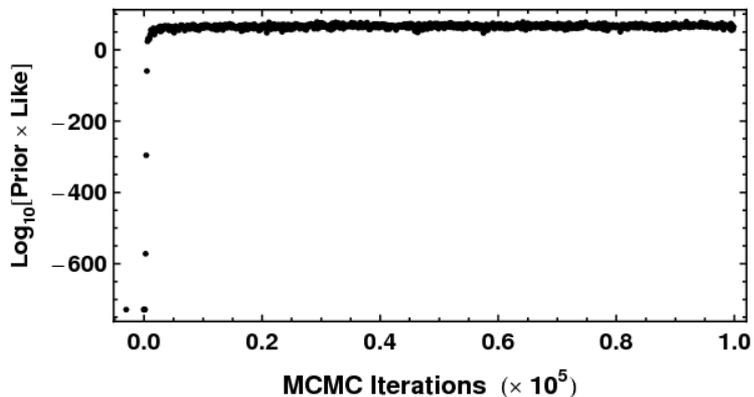
# RV 5 Results

# RV 5 Model: 6 apodized Kepler signals + log(R'hk) regression fit

**No definite planets**

Possible planet at  $P = 0.96$  d based on apodization width.

Bayes factor finds against a real  $P = 0.96$  d planet.

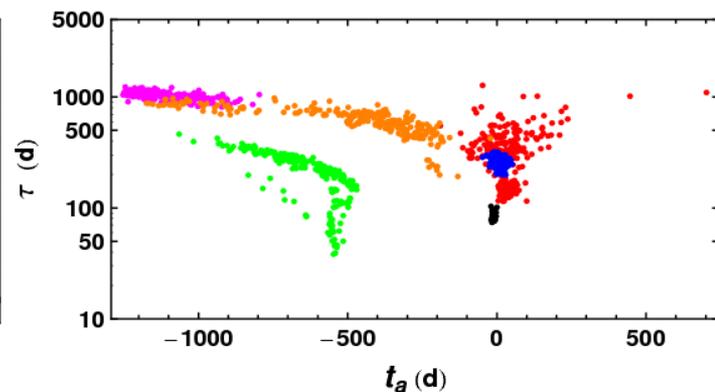
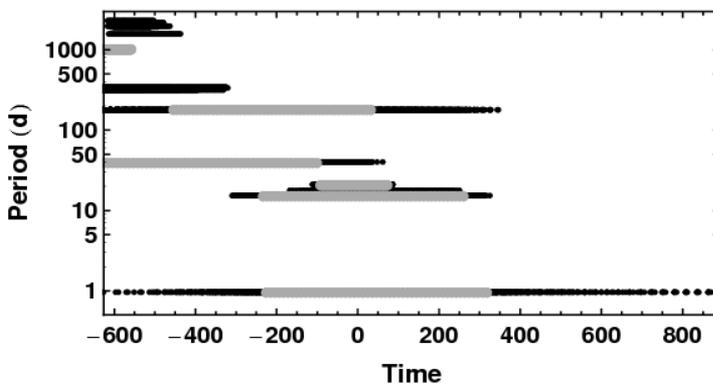


**True planets signals**

P (d)    ecc    K (m/s)

---

14.7	0.17	0.65
26.2	0.25	0.44
34.7	0.03	0.69
173.2	0.05	0.59
283.1	0.3	0.41
616.3	0.03	0.55



# Summary Statistics

Data set	# Signals extracted	Initial $\sigma_I$ (m/s)	Regression residual (m/s)	Final Residual $\sigma_R$ (m/s)	Mean meas. error (m/s)	$\sigma_I/\sigma_R$
Test	5	8.55	2.7	1.35	0.5	6.3
RV 1	6	5.6	3.0	1.44	0.67	3.8
RV 2	8	8.58	4.0	1.42	0.67	6.0
RV 3	6	10.85	5.5	1.79	0.67	6.1
RV 4	8	8.3	3.3	1.52	0.67	5.5
RV 5	6	8.92	2.6	1.18	0.67	7.6

**Conclusion: we are able to dig into the effective noise level set by stellar activity by a factor of ~ 6.  
Still have a long way to go!!**

# Conclusions on Apodized Kepler model

- 1. Conceptually simple approach based on assumption that stellar activity signals vary on time scales shorter than the duration of the data set.  
For very short data sets this assumption would break down.**
- 2. Relatively fast to compute (15 min for a one apodized Kepler model implemented in *Mathematica* and scales linearly with number of signals.)**
- 3. Performed well for  $K > 1$  m/s and resulted in no false detections.**
- 4. Can be employed with other likelihood models (like Student's t) to help with outliers.**
- 5. Next step to see if some combination of the 3 best techniques performs better and try out other apodization functions.**