# Time series analysis methods
# Some alternatives

Jaan Pelt

Tartu Observatory

Space Climate School

March 30 – April 3, 2016

Levi, Finnish Lapland
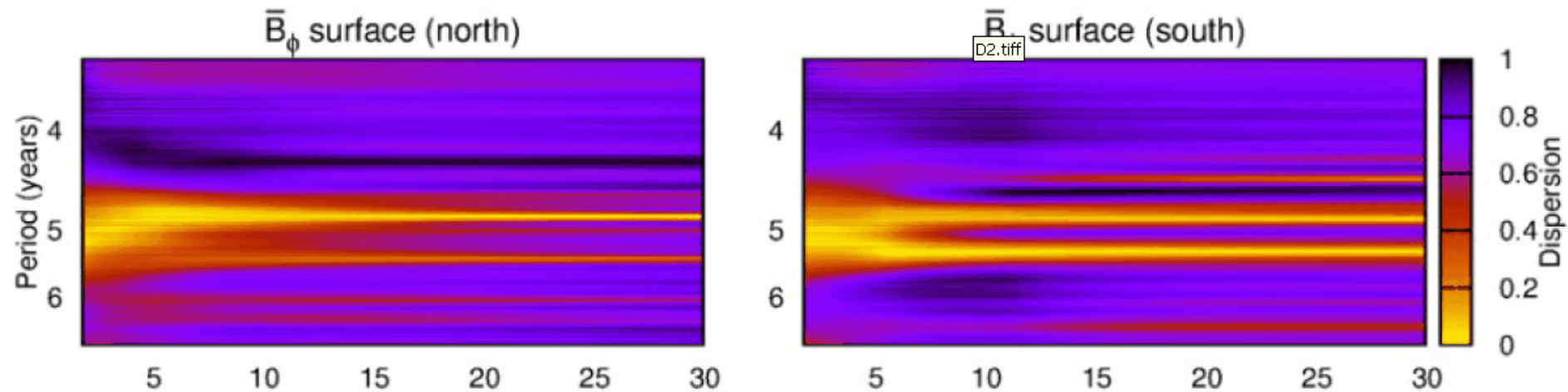


ReSoLVE
CENTRE OF EXCELLENCE

# From harmonics to cycles

- Harmonics
- Periods
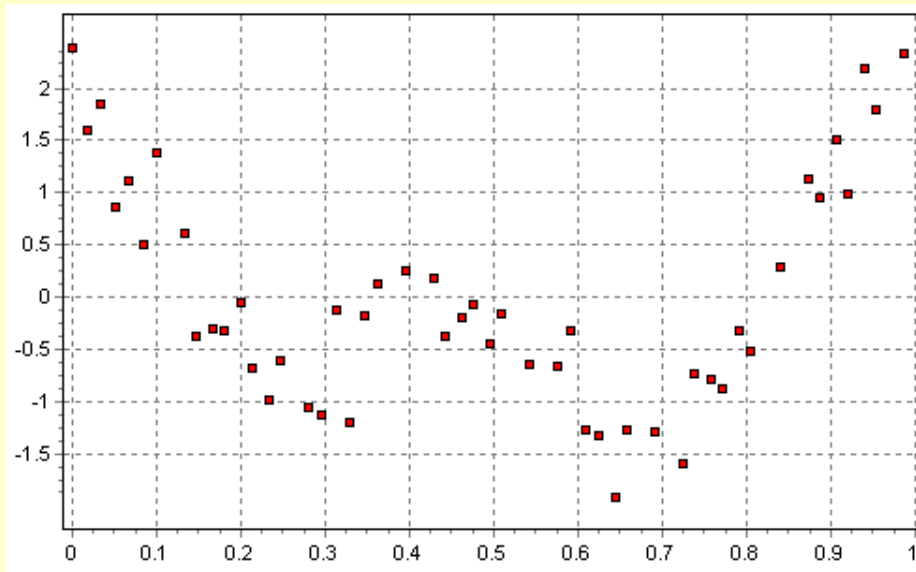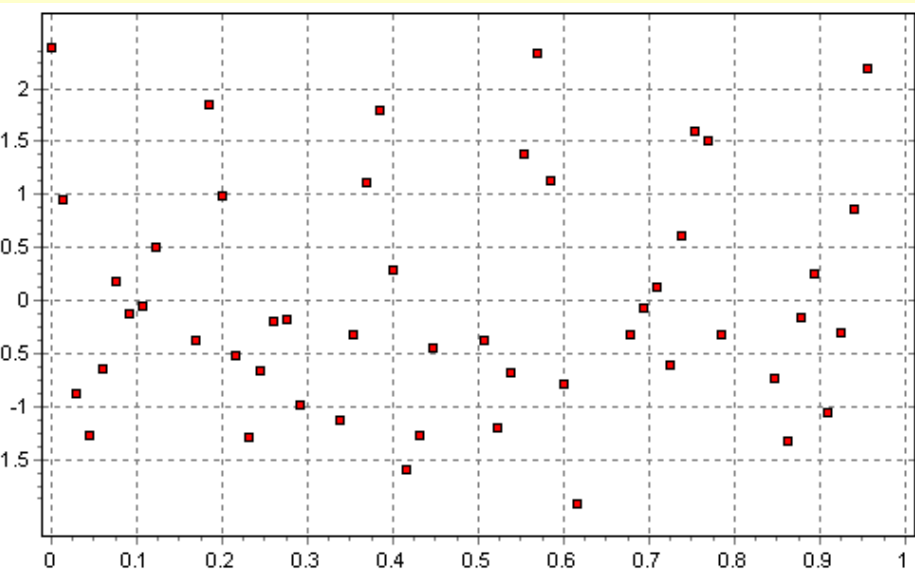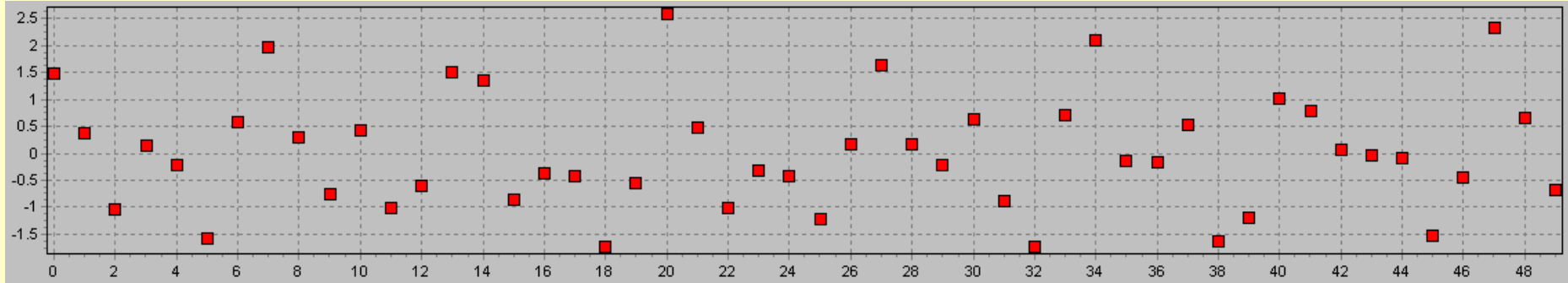- Multiperiodicity
- Cyclicity

# Plan

- What is $D^2$ method?

- What is CF method?

- What is FDC method?.

# What is D² method?



- D² analysis (correlating only phases over a certain coherence time, not over the full time span) reveals the *dominating "solar butterfly" –like period of roughly 5.35 years.*
- This is roughly 4 times shorter than the solar cycle, but if scaled back to solar time units, the simulation length would be roughly 2 millennia.
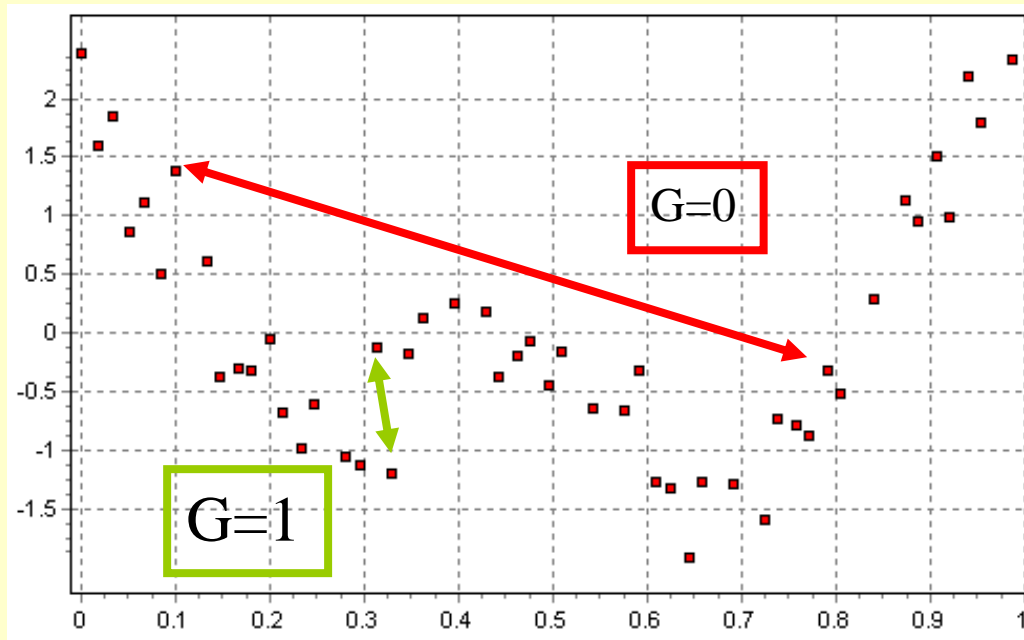
# Phase-process diagram (folding)



$$t_i, f(t_i), i = 1, 2, \ldots, N$$

$$\varphi_P(t_i) = \mathrm{Frac}(t_i P^{-1})$$

$$D^2(\nu) = \frac{\sum_{i=1}^{N-1} \sum_{j=i+1}^{N} G(t_i, t_j, \nu) \left[f(t_i) - f(t_j)\right]^2}{2\sum_{i=1}^{N-1} \sum_{j=i+1}^{N} G(t_i, t_j, \nu)}$$

Weights G are larger than zero when phases of two points in pair are similar, or:

$$t_i - t_j \approx \frac{k}{\nu}, \, k = 0, \pm 1, \pm 2, \ldots,$$

# Smoother periodograms

$$D^2(\nu) = \frac{\sum_{i=1}^{N-1}\sum_{j=i+1}^{N} G^*(t_i, t_j, \nu)\left[f(t_i) - f(t_j)\right]^2}{2\sum_{i=1}^{N-1}\sum_{j=i+1}^{N} G^*(t_i, t_j, \nu)}$$

$$G^*(t_i, t_j, \nu) = G(t_i, t_j, \nu)W(t_i, t_j)$$

$$W(t_i, t_j) = \begin{cases} 1, & |t_i - t_j| \leq t_{\max} \\ 0, & \text{otherwise,} \end{cases}$$

# How to compute?

$$t_i - t_j \approx k\delta t, \text{ for some } k$$

$$C_k = \sum_{t_i - t_j \approx k\delta t} \left[ f(t_i) - f(t_j) \right]^2$$

$$D^2(\nu) = \frac{\sum_{i=1}^{N-1} \sum_{j=i+1}^{N} G(t_i, t_j, \nu) \left[ f(t_i) - f(t_j) \right]^2}{2 \sum_{i=1}^{N-1} \sum_{j=i+1}^{N} G(t_i, t_j, \nu)}$$

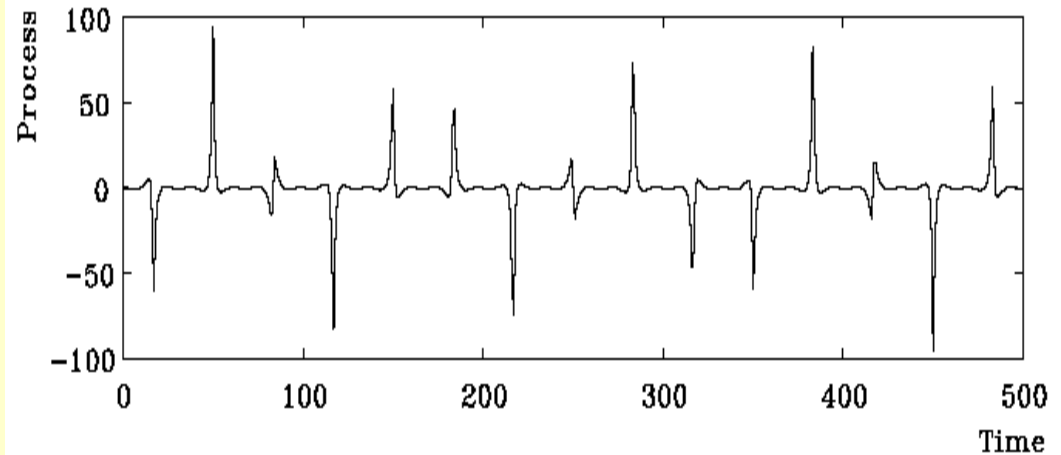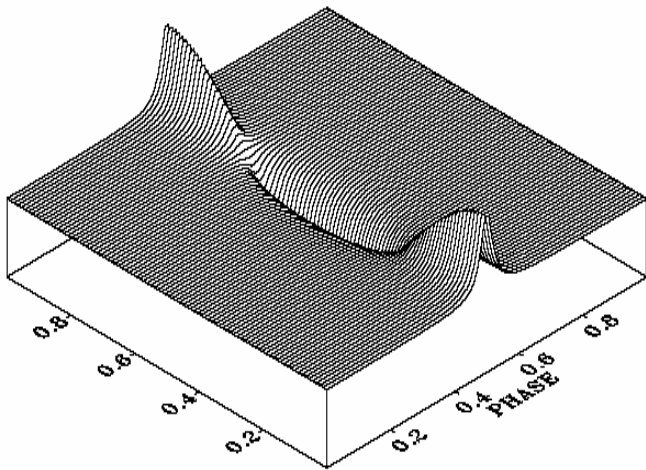$$S_k = \sum_{t_i - t_j \approx k\delta t} 1$$

$$G(t_i, t_j, \nu) = d\big( (t_i - t_j)\nu \big) \approx d(k\delta t\nu)$$

$$D^2(\nu) = \frac{\sum_{k=0}^{K} d(k\delta t\nu) C_k}{2 \sum_{k=0}^{K} d(k\delta t\nu) S_k}$$

$$d(k\delta t\nu) = \sum_{r=0}^{\infty} d_r \cos\big( 2\pi r k\delta t\nu \big)$$

# Multiperiodic oscillations

$$f(t) = \sum_{k_1=-\infty}^{\infty} \cdots \sum_{k_R=-\infty)}^{\infty} \left( a_{k_1,\ldots,k_R} \cos(2\pi \sum_{r=1}^{R} k_r \bar{\nu}_r t) + \right.$$

$$\left. + b_{k_1,\ldots,k_R} \sin(2\pi \sum_{r=1}^{R} k_r \bar{\nu}_r t) \right) + \varepsilon(t).$$
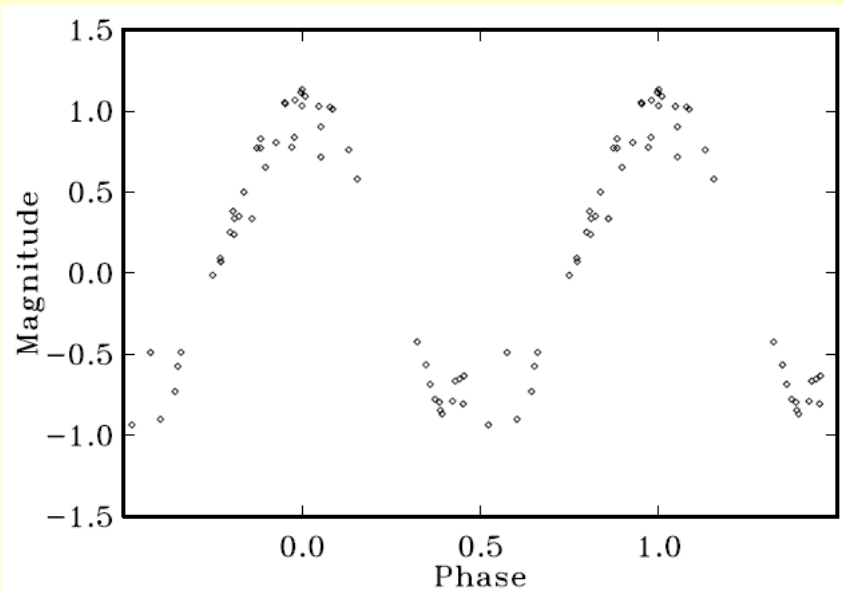
# Phase dispersion for multiperiodic processes

$$\frac{\sum_{i=1}^{N-1} \sum_{j=i+1}^{N} G(t_i, t_j, \nu_1, \nu_2) \big[f(t_i) - f(t_j)\big]^2}{2\sum_{i=1}^{N-1} \sum_{j=i+1}^{N} G(t_i, t_j, \nu_1, \nu_2)}$$
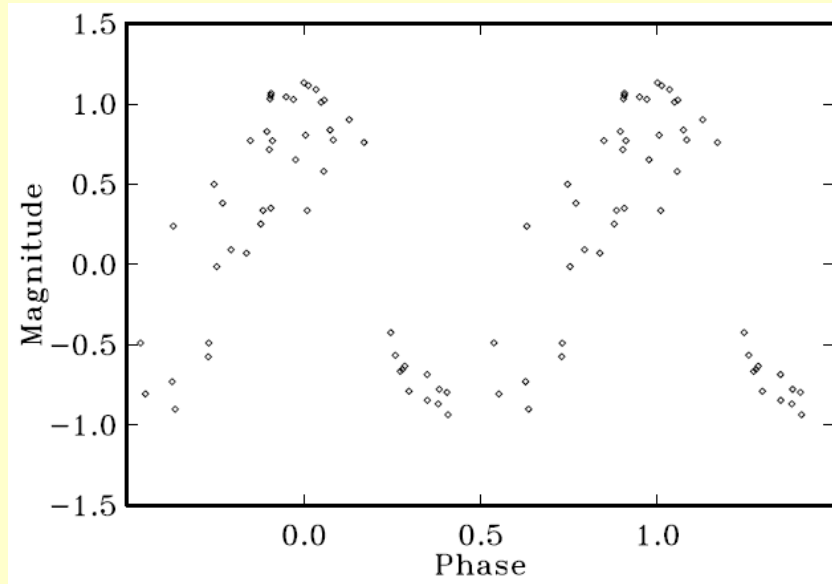
# An example

$$f(t) = 0.55073 \cos(2\pi t/0.53289 + 0.09111) +$$

$$+ 0.58717 \cos(2\pi t/7.83453 + 0.91564).$$
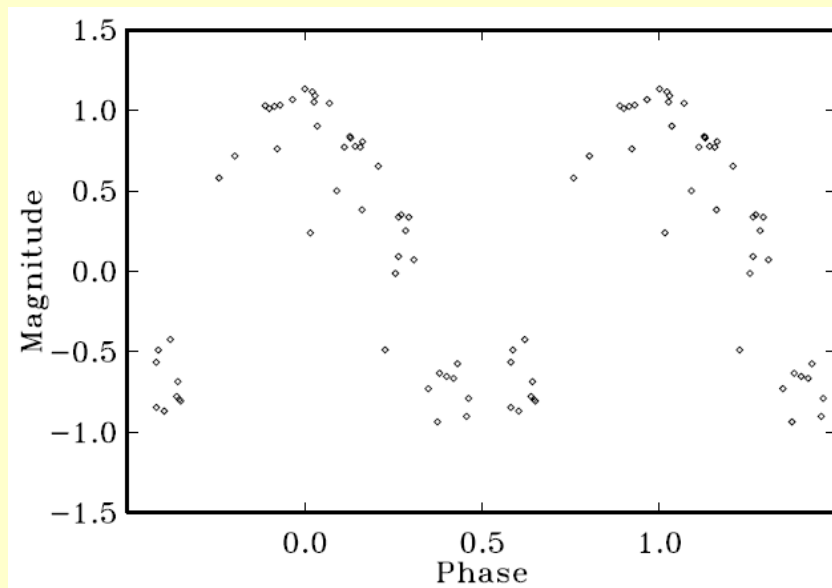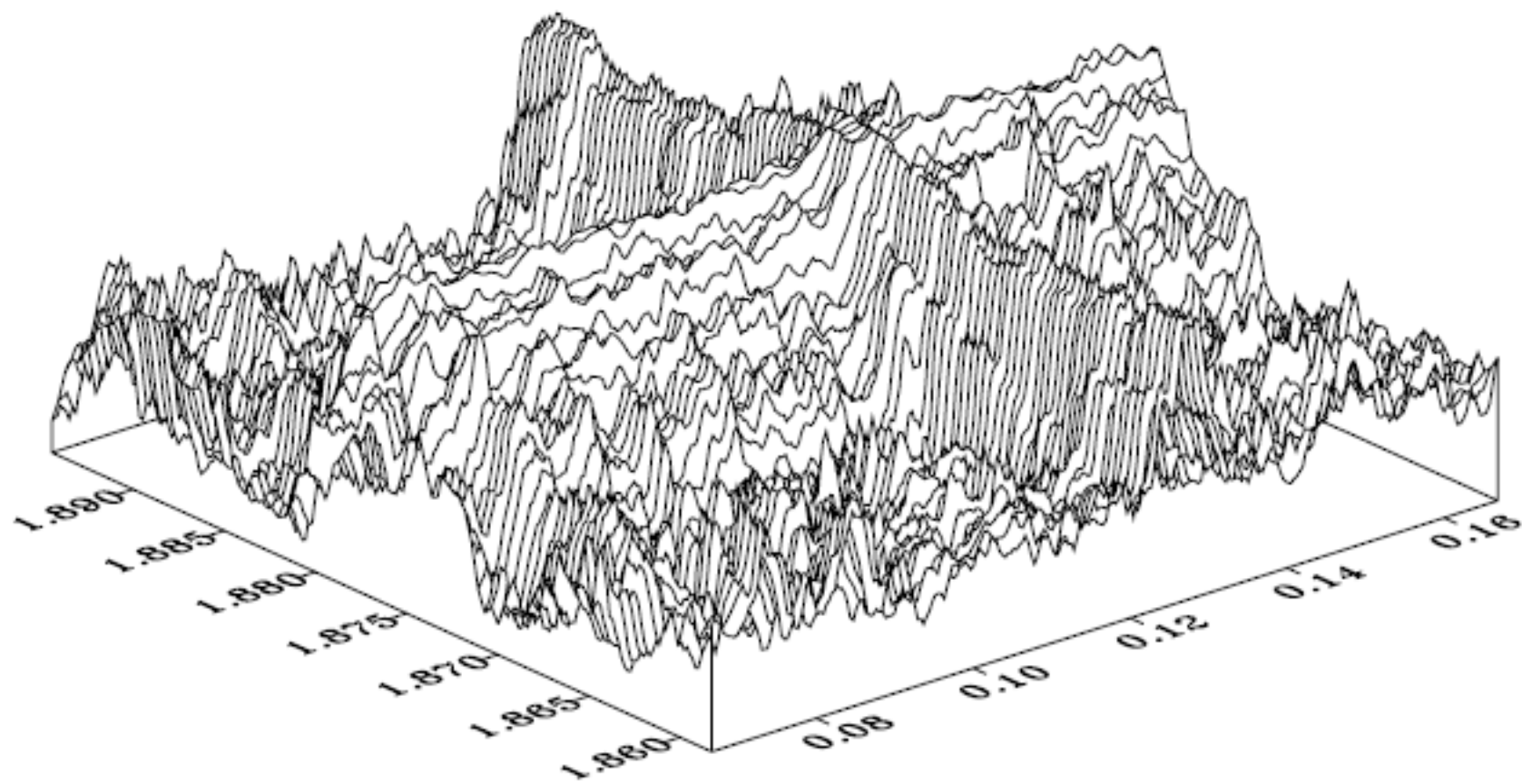
???

# Why?



deepest peak in spectrum $P = 1.143455$



first real period $P = 0.5328956$



second real period $P = 7.8345346$

Two-dimensional "power spectrum" $1 - D(\nu_1, \nu_2)$

# Multichannel search



$$t_{ij} = |t_i - t_j|,$$

$$w_{ij}^c = \frac{1}{\sigma^c(t_i)^2 + \sigma^c(t_j)^2} = \frac{w_i^c \cdot w_j^c}{w_i^c + w_j^c},$$

$$y_{ij}^c = w_{ij}^c \cdot |y_i^c - y_j^c|^2.$$

# Many channels in LSQ format

Dispersions in separate channels

$$WRSS^c(P) = \sum_{i=1}^{N} w_i^c [y_i^c - M(t_i, \beta^c(P))]^2$$

Total dispersion

$$WRSS(P) = \sum_{c=1}^{C} WRSS^c(P)$$

# Total dispersion

$$D(P) = \frac{\sum_{c=1}^{C} \sum_{i=1}^{N-1} \sum_{j=i+1}^{N} g(t_{ij}, P) L(t_{ij}) y_{ij}^{c}}{\sum_{c=1}^{C} \sum_{i=1}^{N-1} \sum_{j=i+1}^{N} g(t_{ij}, P) L(t_{ij}) w_{ij}^{c}}.$$
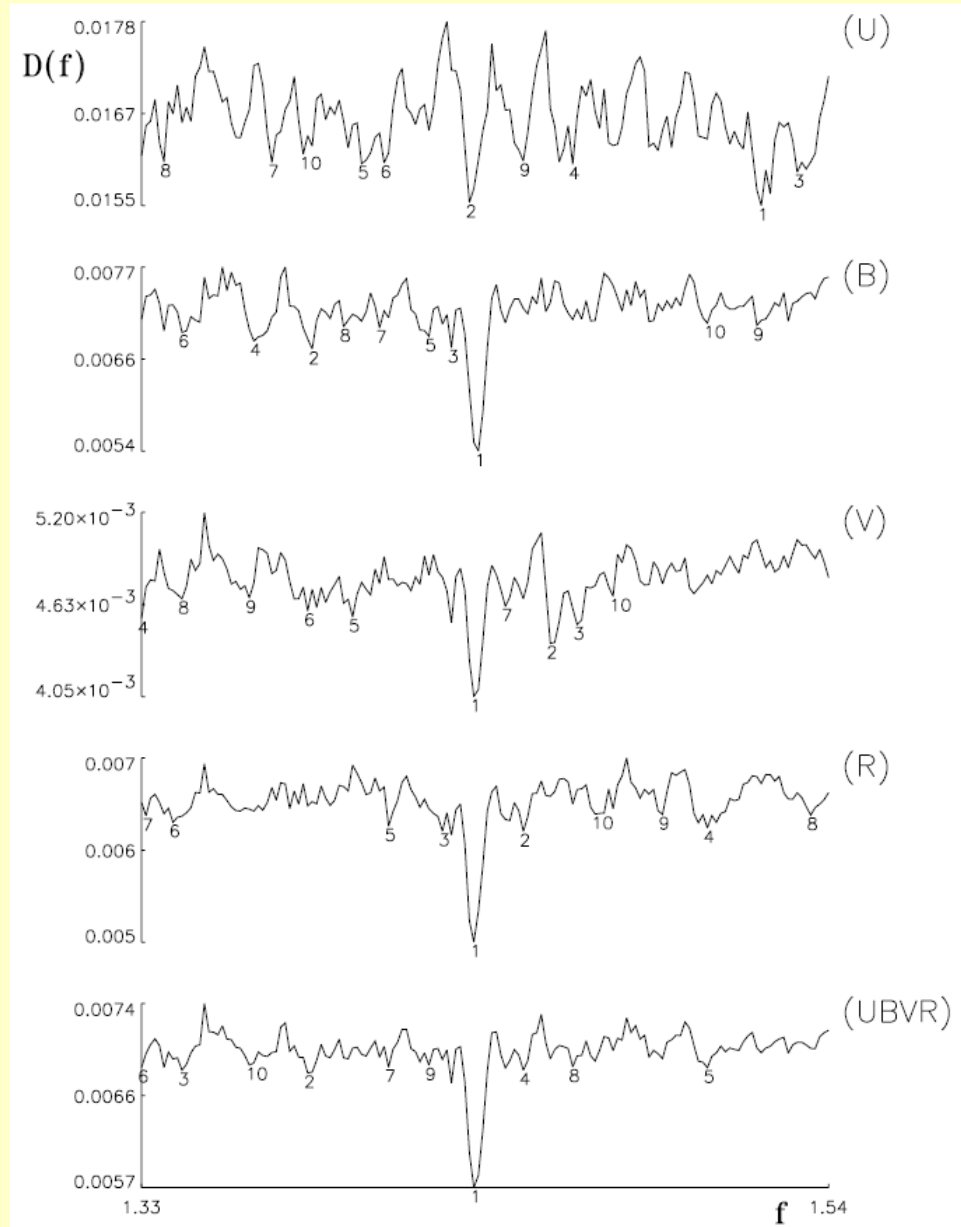
Phases

$$\phi_{ij}(P) = \mathrm{Frac}\frac{t_i - t_j}{P}.$$

must be near by

$$g(t_{ij}, P) = \begin{cases} 1, & \phi_{ij}(P) \leq \tau \quad \text{or} \\ & \phi_{ij}(P) > 1 - \tau \\ 0, & \text{otherwise.} \end{cases}$$

Resolution can be controlled

$$L(t_{ij}) = \begin{cases} 1, & D_{\min} \leq |t_i - t_j| < D_{\max} \\ 0, & \text{otherwise.} \end{cases}$$

$$D^2(P, \Delta t) = \frac{\sum\limits_{i=1}^{N-1} \sum\limits_{j=i+1}^{N} g(t_i, t_j, P, \Delta t)[f(t_i) - f(t_j)]^2}{2\sigma^2 \sum\limits_{i=1}^{N-1} \sum\limits_{j=i+1}^{N} g(t_i, t_j, P, \Delta t)},$$

$$(B.1)$$

where $f(t_i), i = 1, \ldots, N$ is the input time series, $\sigma^2$ is its variance, $g(t_i, t_j, P, \Delta t)$ is the selection function, which is significantly greater than zero only when

$$t_j - t_i \approx kP, k = \pm 1, \pm 2, \ldots \quad \text{and} \qquad (B.2)$$

$$|t_j - t_i| \lesssim \Delta t, \qquad (B.3)$$

where $P$ is the trial period and $\Delta t$ is the so-called coherence time, which is the measure of the width of the sliding time window wherein the data points are taken into account

$$g(t_i, t_j, P, \Delta t) = g_1(t_i, t_j, P) \cdot g_2(t_i, t_j, \Delta t), \tag{B.4}$$

$$g_1(t_i, t_j, P) = \frac{1}{2}\left(\cos\left(2\pi \cdot \mathrm{frac}\left(\frac{t_j - t_i}{P}\right)\right) + 1\right), \tag{B.5}$$

$$g_2(t_i, t_j, \Delta t) = \exp\left(-\ln 2\left(\frac{t_j - t_i}{\Delta t}\right)^2\right), \tag{B.6}$$

where $\mathrm{frac}((t_j - t_i)/P)$ removes the integer part of $(t_j - t_i)/P$.
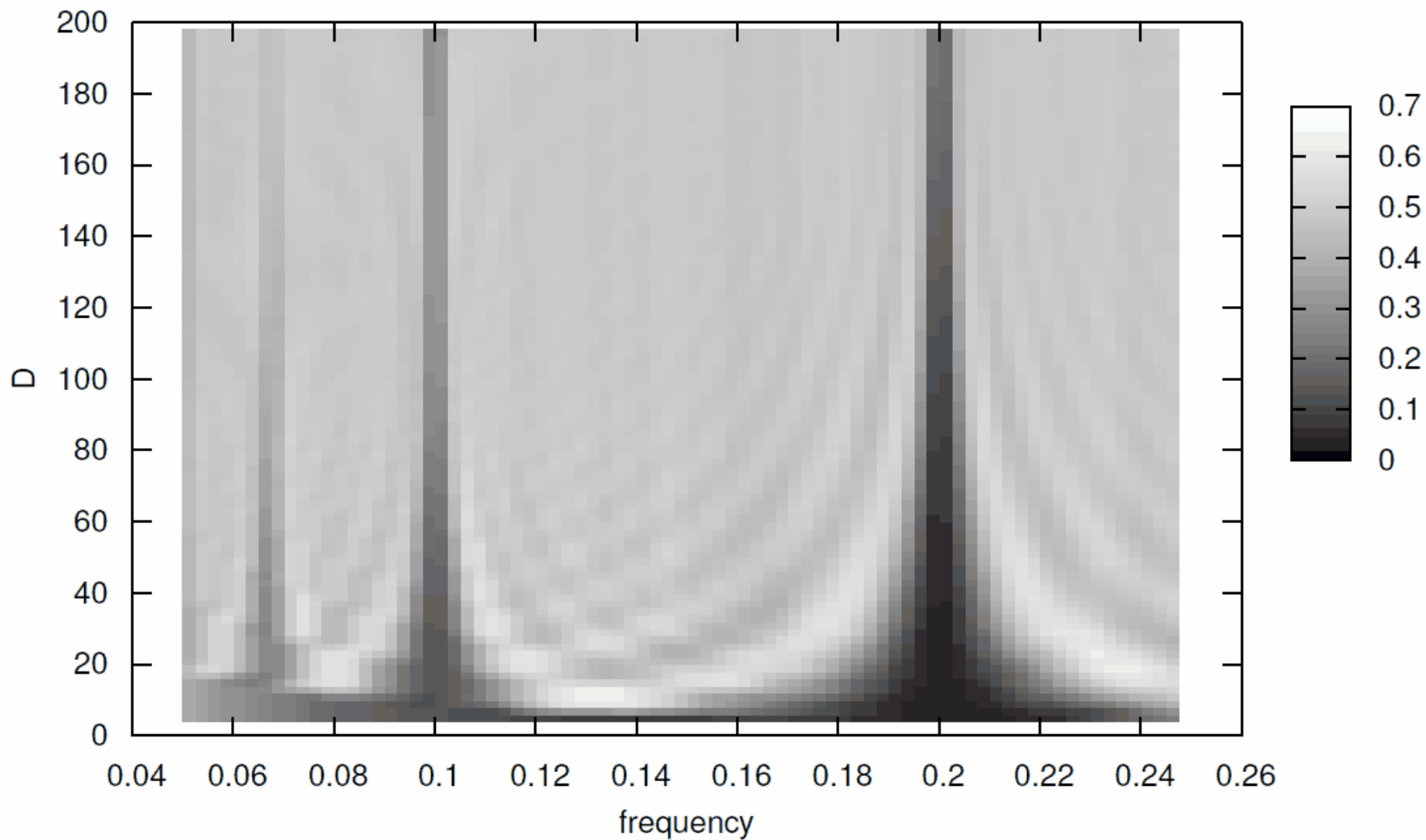
Phase weighting and coherence weighting

Figure 1: $(\nu, D)$-spectrum for data with constant frequency

**Fig. 6.** Phase dispersion analysis results for the components of $\overline{B}$. The calculations were done for $\overline{B}_\phi$ at latitude $\pm 22°$ and radius $0.94 R_\odot$, for $\overline{B}_r$ at latitude $\pm 66°$ and radius $0.94 R_\odot$, and for $\overline{B}_\theta$ at latitude $\pm 49°$ and radius $0.82 R_\odot$. Left (right) refers to the northern (southern) hemisphere.

From photometric channels to $N^2$ or more time dependent pixels.

Distances need to be computed only once!

# $D^2$ Summary

- $D^2$ method can be used for search of periods in irregularly spaced data
- It is usable also in context where periodicity is only approximate (process is cyclic)
- Method can be easily generalized for different geometries
- $D^2$ method allows to classify data curves in frequency and coherence terms

# What is CF method?

## Basic model

Model with single harmonic carrier:

$$f(t) = a(t)\cos(2\pi\nu_0 t) + b(t)\sin(2\pi\nu_0 t),$$

where $\nu_0$ is so called **carrier frequency.**

More general waveform:

$$f(t) = \sum_{k=1}^{K} a_k(t)\cos(2\pi k\upsilon_0 t) + b_k(t)\sin(2\pi k\upsilon_0 t)$$

The coefficients themselves are modeled using harmonic expansions:

$$a_k(t) = \sum_{l}^{L} a_{kl}^{(a)} \cos(2\pi l \upsilon_m t) + b_{kl}^{(a)} \sin(2\pi l \upsilon_m t),$$

$$b_k(t) = \sum_{l}^{L} a_{kl}^{(b)} \cos(2\pi l \upsilon_m t) + b_{kl}^{(b)} \sin(2\pi l \upsilon_m t).$$

# What to do with modulating curves?

For our special case ( due to the theorem of Bedrosian) we can define *Hilbert transform* in the following way—if the original function is:

$$u(t) = a(t)\cos(2\pi f_0 t) + b(t)\sin(2\pi f_0 t),$$

then its Hilbert transformed image is:

$$v(t) = a(t)\sin(2\pi f_0 t) - b(t)\cos(2\pi f_0 t).$$

This is true when frequency domain support for $a(t), b(t)$ do not overlap with carrier frequency $f_0$. The good news is, that (as we saw before) the pair of functions $u(t), v(t)$ allows us to define *instantanous frequency, phase and amplitude* for the every time point of the model curve.

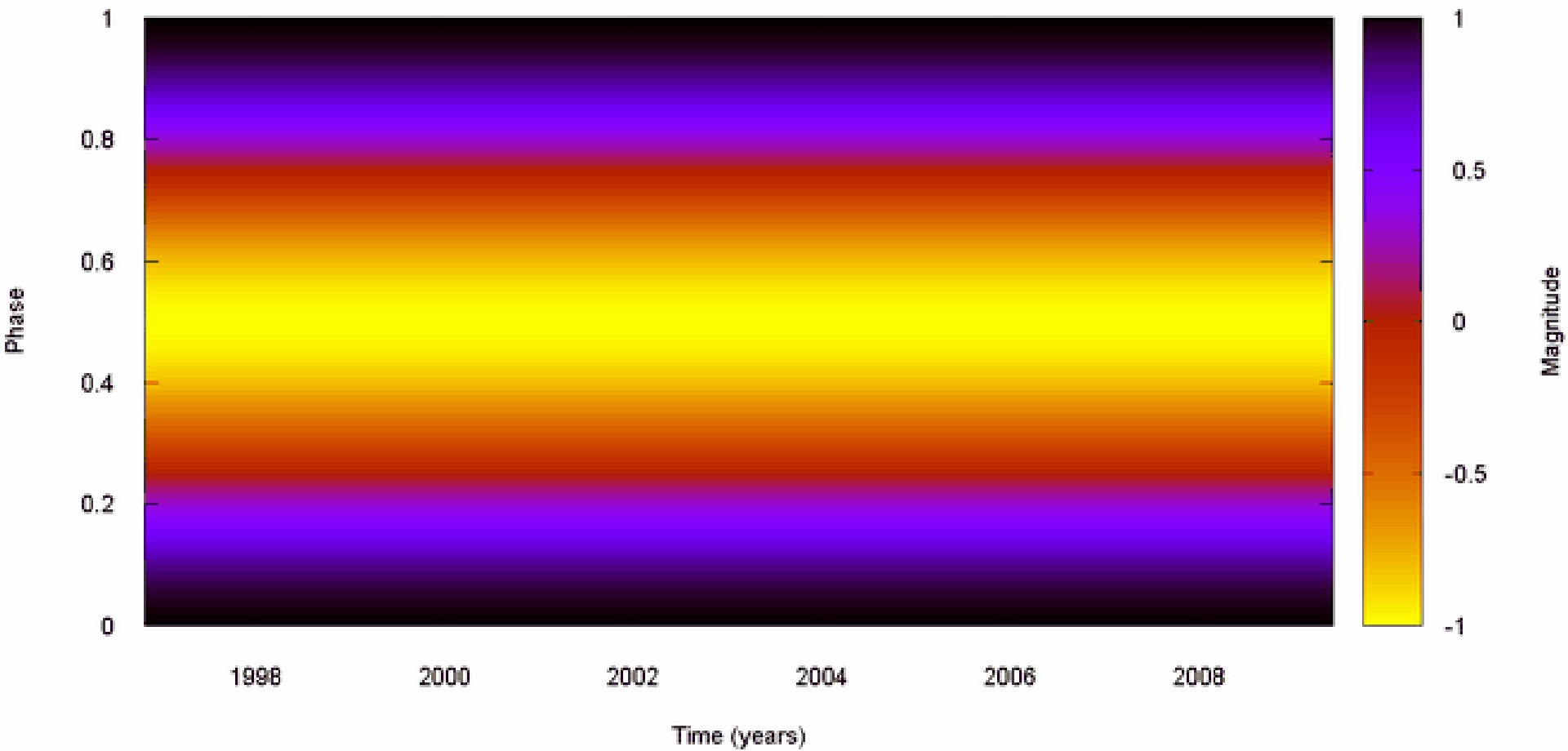# Instantaneous values

Amplitude:

$$A(t) = \sqrt{u^2(t) + v^2(t)},$$

frequency:

$$\nu(t) = \frac{v'(t)u(t) - u'(t)v(t)}{u^2(t) + v^2(t)},$$

phase:

$$\phi(t) = \arctan\left(\frac{u(t)}{v(t)}\right).$$
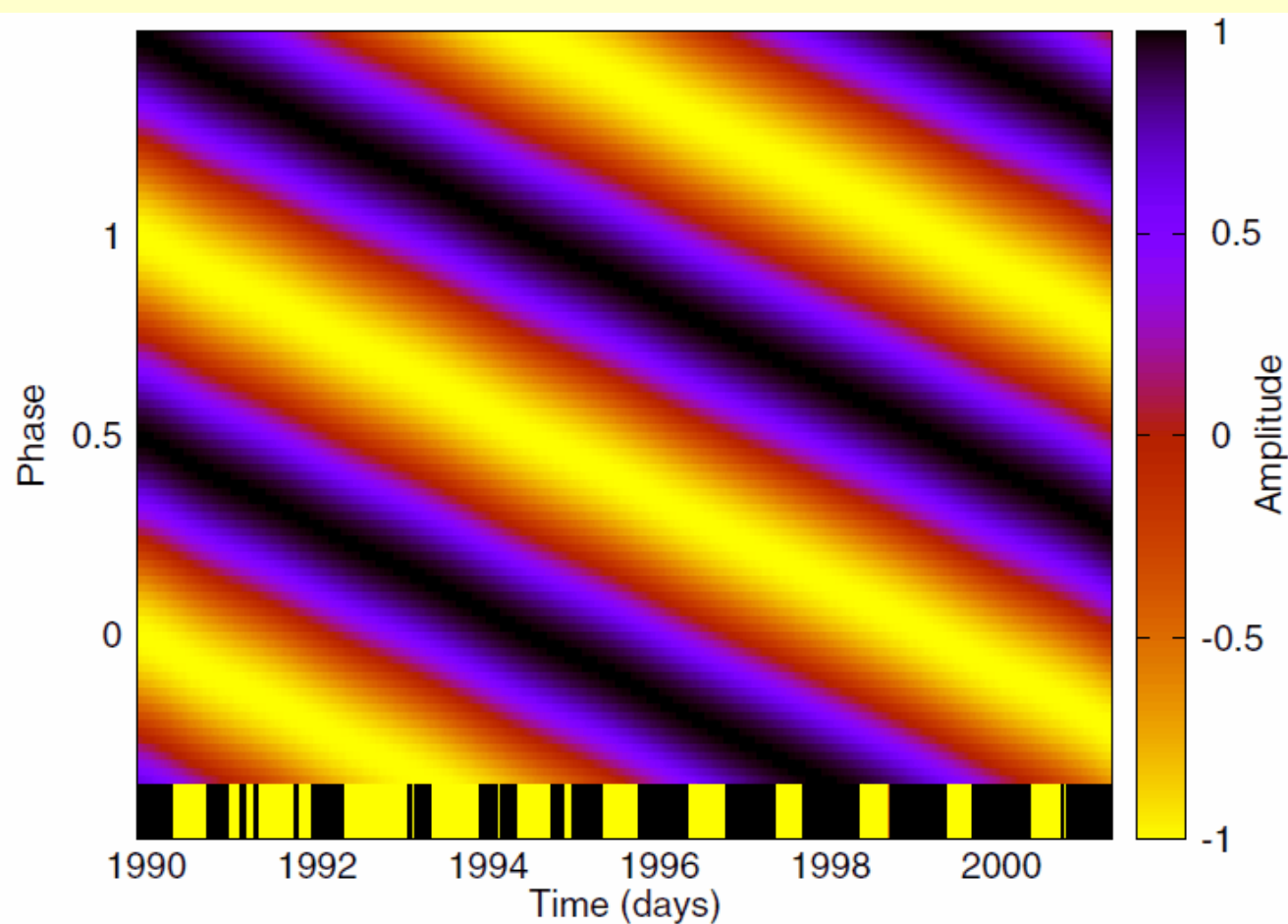
Simple sinusoid

**Fig. 4.** Time-dependent phase diagram, i.e. the light curve amplitude profile over phase ($y$-axis) plotted as function of time ($x$-axis), for the mismatched carrier frequency example presented in Sect. 3.2. In this plot, a slightly too low carrier frequency of $\nu_0 = 0.148714$ is used, with $K = 1$ model and $L = 8$ modulation harmonics. Time points are obtained from real $V$-band photometric observations of LQ Hya; the bar-code in the bottom of the plot indicates when data has been available (black) and the gaps (bright). For visualization details see Sect. 2.6.
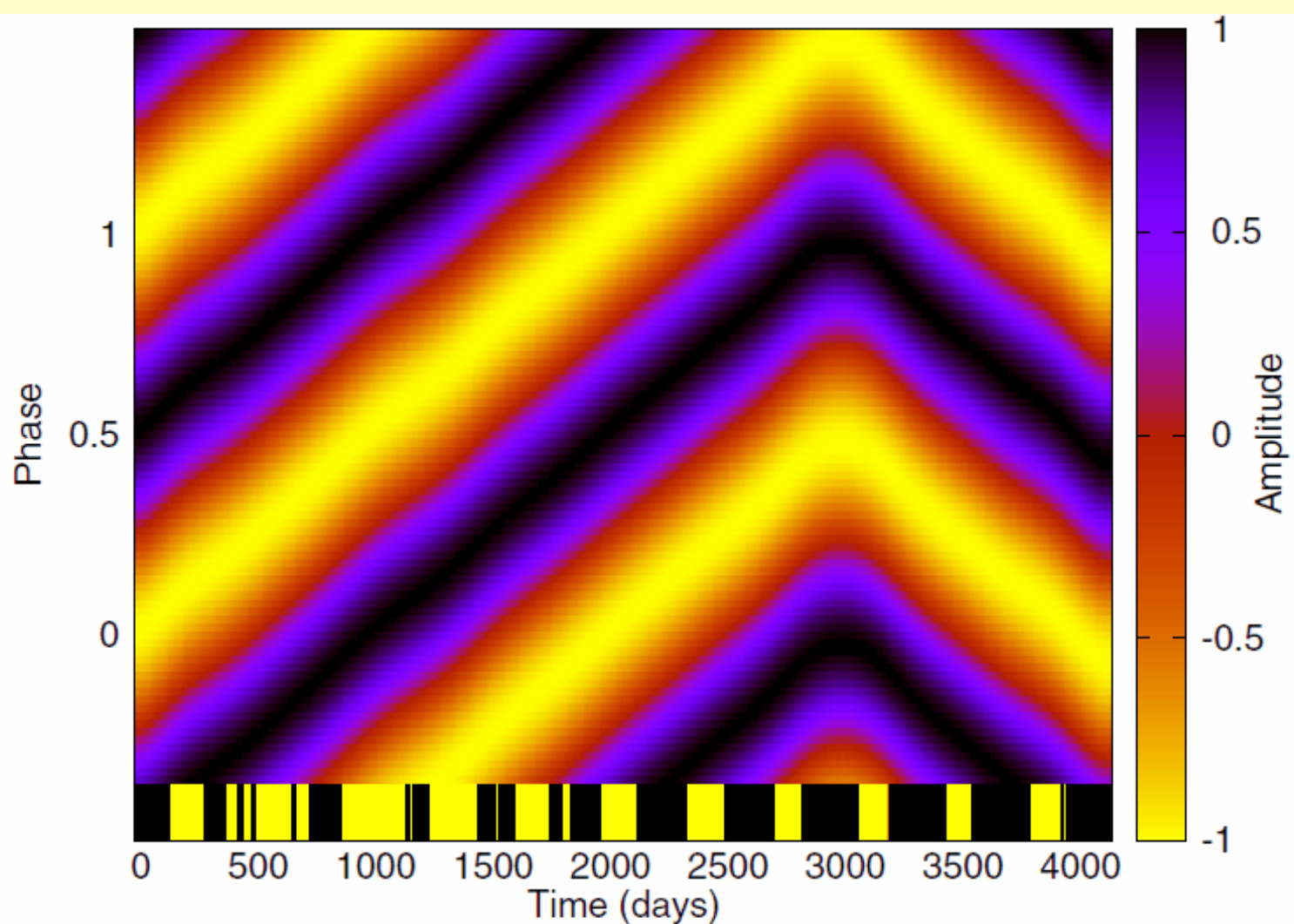
**Fig. 7.** The effect of an abrupt change in the period (Sect. 3.3) in the time-dependent phase diagram. The period before the jump is $P_1 = 6.747017$, and after it $P_2 = 6.7018004$. The used carrier is $P_0 = 6.724333$, the number of model harmonics, $K = 1$, and the number of modulation harmonics, $L = 8$.
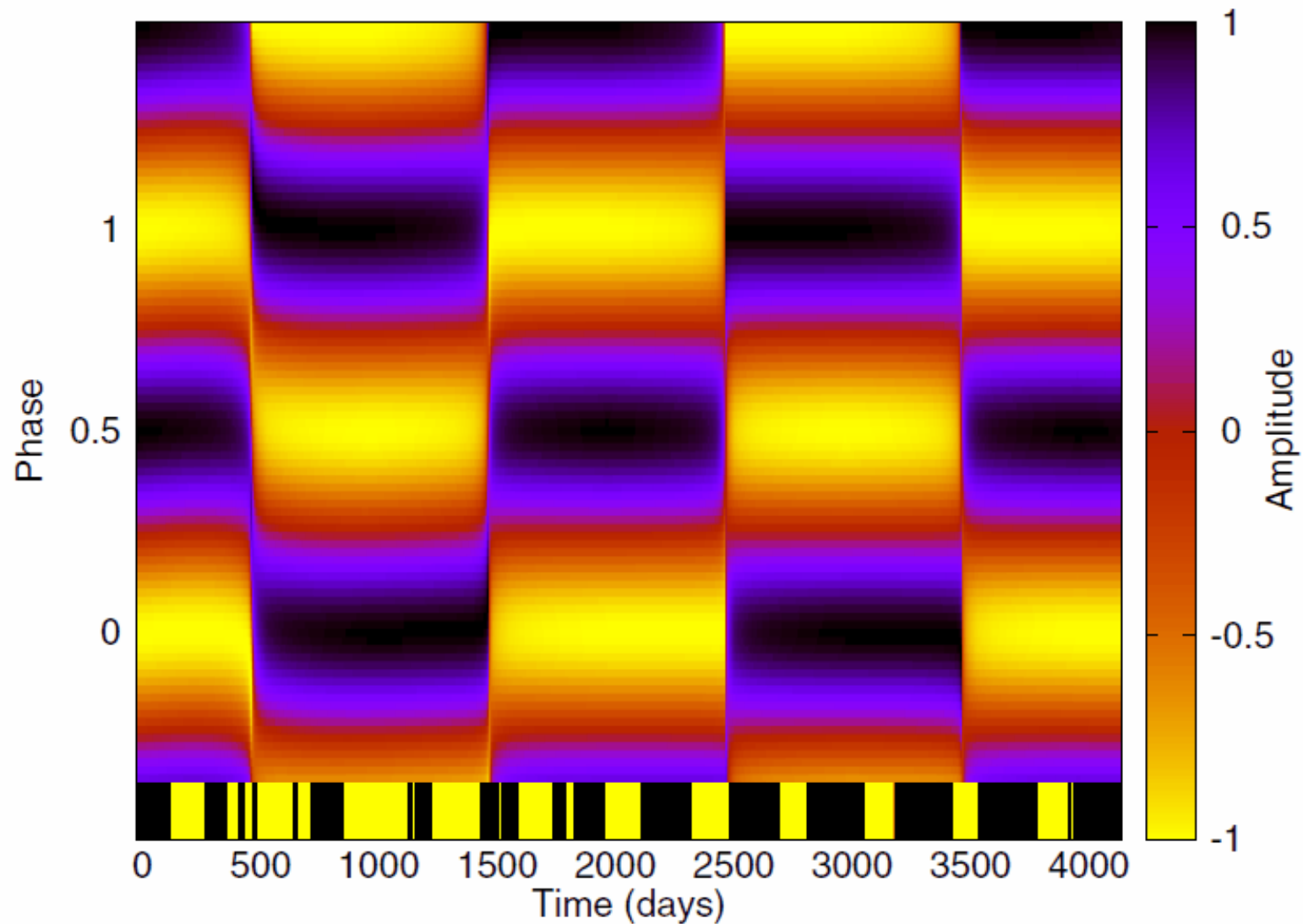
**Fig. 9.** The effect of two beating periods, $P_1 = 6.747017$ and $P_2 = 6.701800$, in the time-dependent phase diagram. The carrier used is $P_0 = 6.724333$, the number of model harmonics $K = 1$, and the number of modulation harmonics $L = 8$. See Sect. 3.4.

**Fig. 12.** The effect of a linear trend in frequency, as described in Sect. 3.5, in the time-dependent phase diagram. The used model parameters for the carrier fit read $\nu_0 = 0.148714$, $K = 1$, and $L = 8$.

FK Com

FK Com detail

# CF Summary

- CF method allows to analyze and visualize nearly periodic processes

- It allows to locate special features in data curves (flip-flops, trends etc)

- Method is of exploratory type and must be supplemented by other techniques.

# What is FDC method?

(a) H1

10 ms light travel time

L1

Test Mass

$L_y = 4$ km

Test Mass

Power Recycling

Beam Splitter

$L_x = 4$ km

Laser Source

20 W

100 kW Circulating Power

Signal Recycling

Test Mass

Test Mass

Photodetector

(b)

H1
L1

Strain noise (Hz$^{-1/2}$)

$10^{-21}$

$10^{-22}$

$10^{-23}$

Frequency (Hz)

20          100          1000

LIGO h = $10^{-21}$
4000 * $10^{-21}$ = 0.4 * $10^{-17}$ m

PROOTONI DIAMEETER 1.67-1.74 * $10^{-15}$ m
JUUKSED 0.17-1.81 * $10^{-4}$

PROOTON/LIGO = 400
(Mujal 1000-10000)

AU 0.15 * $10^{12}$ * $10^{-21}$ = 0.15 * $10^{-9}$ m
Proxima Centauri 0.4 $10^{17}$ * $10^{-21}$ = 0.4 * $10^{-4}$
Foonikiirguse tee 0.127 * $10^{27}$ * $10^{-21}$ = 127km
Tallinn-Põltsamaa





$$g' = \frac{\Delta g}{d} = \frac{\text{change in gravity}}{\text{displacement}}$$

$$h = \frac{2\Delta d}{d} = 2 \times \frac{\text{change in displacement}}{\text{displacement}}$$

Inspiral      Merger   Ring-down

Three targets – TSI Composites

Proxies

**Fig. 1.** Three different transmission curves for the specific width of the time domain filter.

$$T(\nu) = \exp^{-(\nu W)^2}$$

Transmission curve for Gaussian filter

Data -> FFT -> Multiple by TC -> FFT$^{-1}$ =Smoothed data

LF part = smoothed data
HF part = data – smoothed data

To compare different daily data sets we:

A. Cut off parts for which both data exist.

B. Compute $R_c = \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$

|              | PMOD   | RMIB   | ACRIM  | SATIRE-S |
| ------------ | ------ | ------ | ------ | -------- |
| Original     |        |        |        |          |
| PMOD         | 1.0000 | **0.9322** | 0.8553 | 0.8969   |
| RMIB         | **0.9322** | 1.0000 | 0.9167 | 0.8586   |
| ACRIM        | 0.8553 | 0.9167 | 1.0000 | *0.8105* |
| SATIRE-S     | 0.8969 | 0.8586 | *0.8105* | 1.0000 |
| LF           |        |        |        |          |
| PMOD         | 1.0000 | 0.9191 | 0.8422 | **0.9446** |
| RMIB         | 0.9191 | 1.0000 | 0.9022 | 0.8404   |
| ACRIM        | 0.8422 | 0.9022 | 1.0000 | *0.7962* |
| SATIRE-S     | **0.9446** | 0.8404 | *0.7962* | 1.0000 |
| HF           |        |        |        |          |
| PMOD         | 1.0000 | 0.9391 | 0.8939 | 0.8583   |
| RMIB         | 0.9391 | 1.0000 | **0.9424** | 0.8777   |
| ACRIM        | 0.8939 | **0.9424** | 1.0000 | *0.8266* |
| SATIRE-S     | 0.8583 | 0.8777 | *0.8266* | 1.0000 |

**Table 2.** Correlation matrices for input target sets.

|  | PSI | SA | SN | RADIO | MGII | LYMAN |
|---|---|---|---|---|---|---|
| Original | | | | | | |
| PSI | 1.0000 | 0.9424 | 0.8521 | 0.8722 | 0.7726 | *0.7709* |
| SA | 0.9424 | 1.0000 | 0.8786 | 0.9041 | 0.8001 | 0.8116 |
| SN | 0.8521 | 0.8786 | 1.0000 | 0.9466 | 0.9237 | 0.9116 |
| RADIO | 0.8722 | 0.9041 | 0.9466 | 1.0000 | 0.9616 | 0.9552 |
| MGII | 0.7726 | 0.8001 | 0.9237 | 0.9616 | 1.0000 | **0.9735** |
| LYMAN | *0.7709* | 0.8116 | 0.9116 | 0.9552 | **0.9735** | 1.0000 |
| LF | | | | | | |
| PSI | 1.0000 | **0.9978** | 0.9890 | 0.9943 | 0.9885 | 0.9896 |
| SA | **0.9978** | 1.0000 | 0.9880 | 0.9890 | *0.9822* | 0.9822 |
| SN | 0.9890 | 0.9880 | 1.0000 | 0.9915 | 0.9858 | 0.9841 |
| RADIO | 0.9943 | 0.9890 | 0.9915 | 1.0000 | 0.9968 | 0.9930 |
| MGII | 0.9885 | *0.9822* | 0.9858 | 0.9968 | 1.0000 | 0.9900 |
| LYMAN | 0.9896 | 0.9822 | 0.9841 | 0.9930 | 0.9900 | 1.0000 |
| HF | | | | | | |
| PSI | 1.0000 | **0.8939** | 0.7410 | 0.8172 | 0.5715 | 0.5852 |
| SA | **0.8939** | 1.0000 | 0.7452 | 0.8315 | *0.5527* | 0.5909 |
| SN | 0.7410 | 0.7452 | 1.0000 | 0.8020 | 0.7054 | 0.6663 |
| RADIO | 0.8172 | 0.8315 | 0.8020 | 1.0000 | 0.8047 | 0.8014 |
| MGII | 0.5715 | *0.5527* | 0.7054 | 0.8047 | 1.0000 | 0.8725 |
| LYMAN | 0.5852 | 0.5909 | 0.6663 | 0.8014 | 0.8725 | 1.0000 |

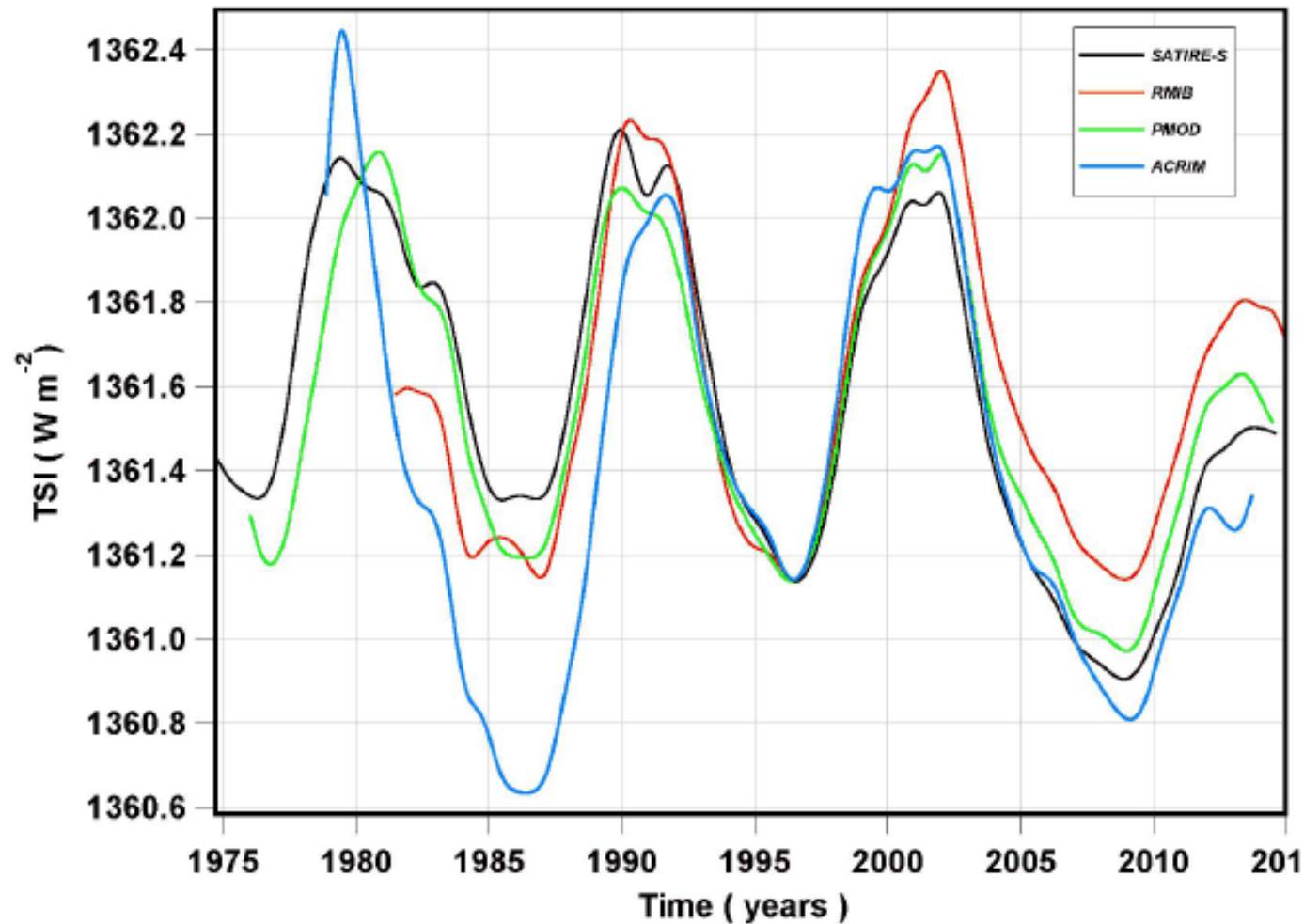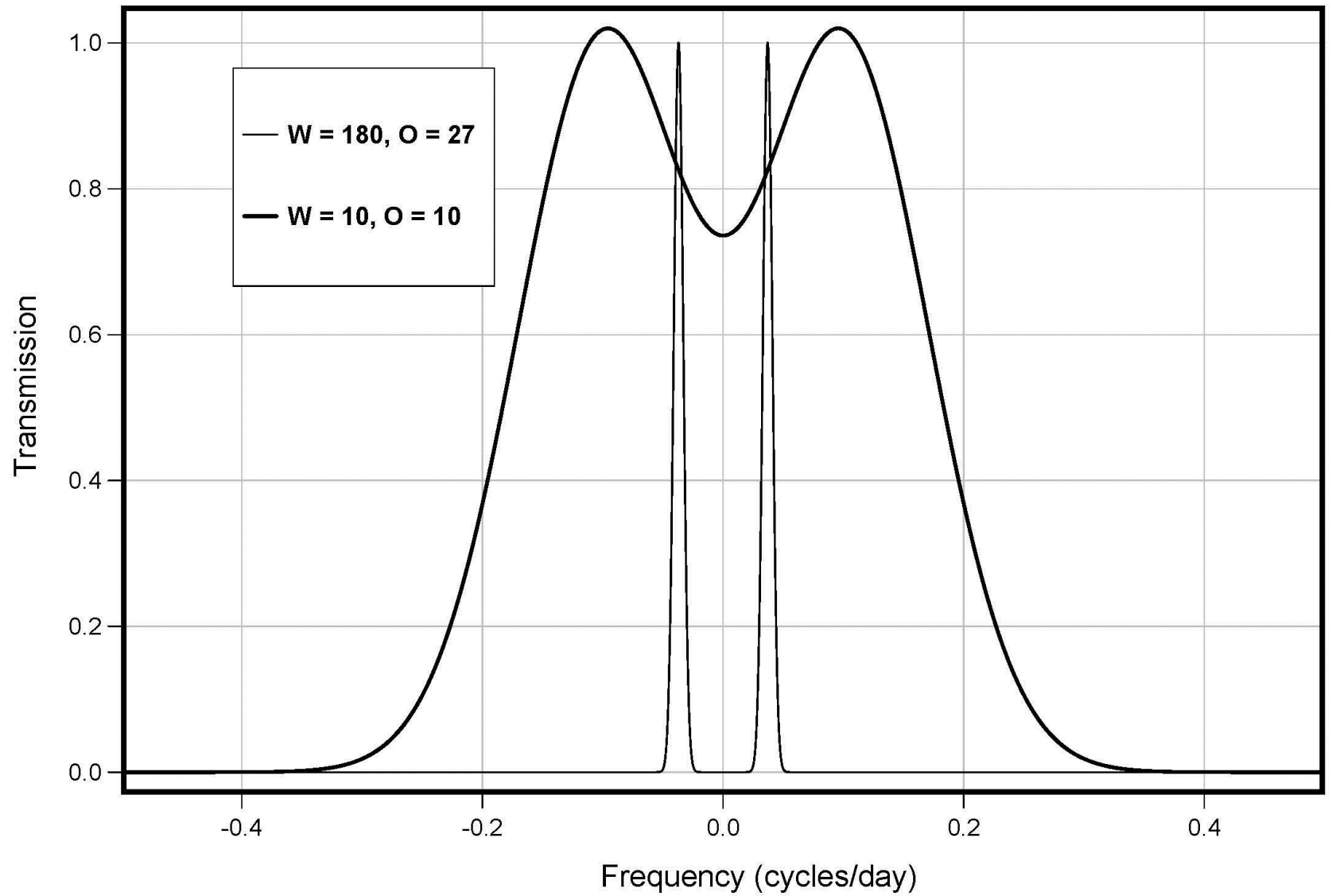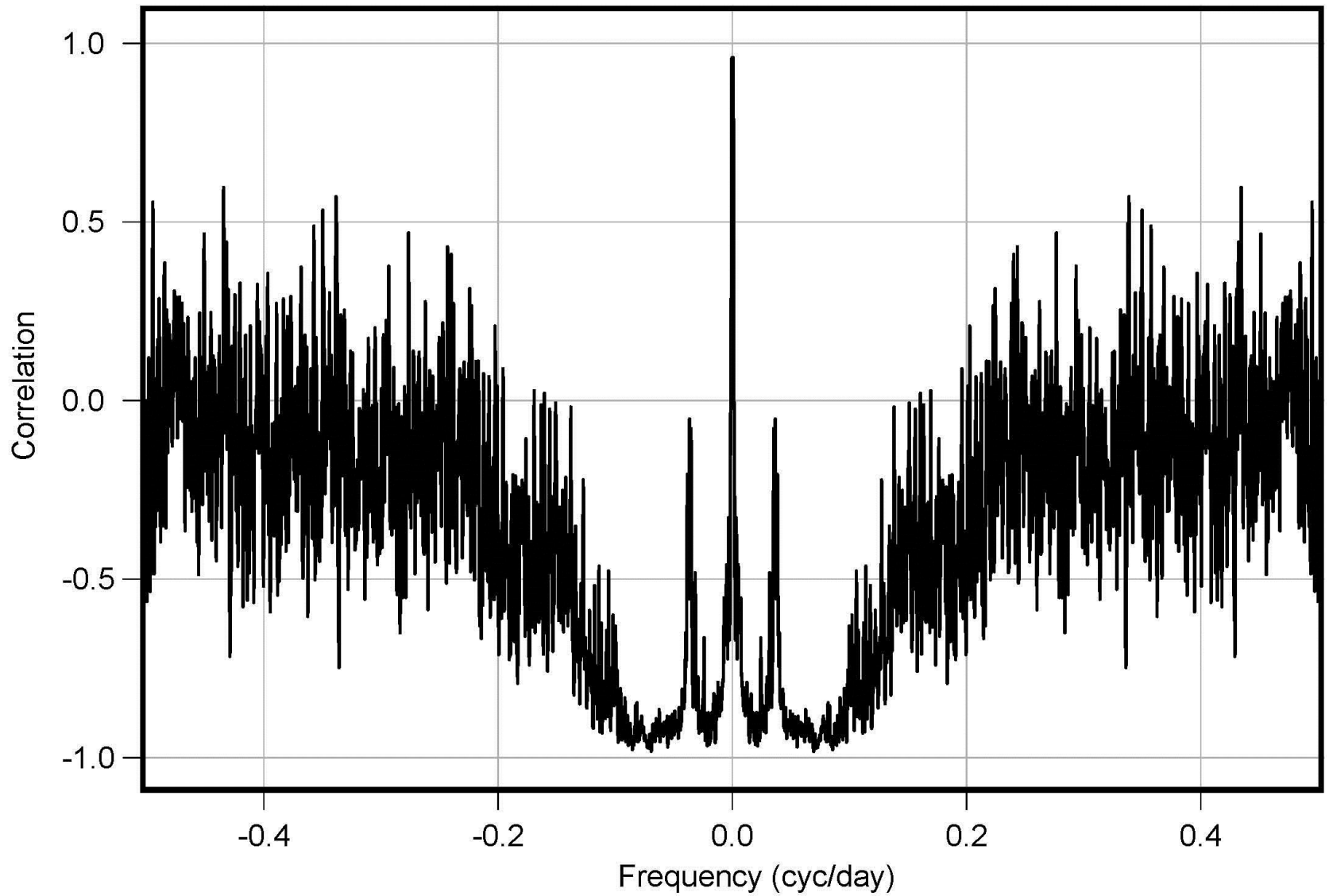**Table 1.** Correlation matrices for input proxy data sets.

**Fig. 6.** The major controversy. To recover TSI for the past dates we need current estimates for calibration. But which one to choose? LF parts ($W = 750$ days) of four target TSI composites. The input curves are shifted to common level at 1996.465.

PMOD vs PSI

W = 2000 days, varying offsets

W = 300, O = 27  (days)

W = 1000, O = 27 (days)

$$SS = \sum_{i=1}^{N} \left( y(t_i) - \sum_{l=0}^{L} a_l E_l(...)[t_i] \right)^2$$

- Start model building from including into it constant level and the input curve as components.
- In each additional step try all the variants from the full library of filtered or otherwise processed components to be as possible next components.
- Evaluate prospective candidates by correlating obtained model with actually observed signal.
- Finally, use the obtained model to hindcast unobserved irradiance values.

| Coefficent | Predictor | Value |
|---|---|---|
| $a_0$ | 1.0 | 1365.551 |
| $a_1$ | $E(0,0)$ | -8.030 |
| $a_2$ | $E(766.229, 0.0)$ | 19.049 |
| $R_c$ | 0.8597 | |

**Table 4.** Summary of least squares fit results with one smoothed component.

**Fig. 7.** Common parts of the PMOD and PSI are practically uncorrelated ($R_c$ = 0.0300). Here we use predicted-observed crossplot as we will do below for more general than trivial linear regression model.

**Fig. 9.** The observations and simple model with smoothed component is much more strongly correlated ($R_c = 0.8597$).

**Fig. 10.** Reconstruction of the TSI using the HF part of PSI combined with LF part of PMOD. Reconstruction is plotted in red, target in black.

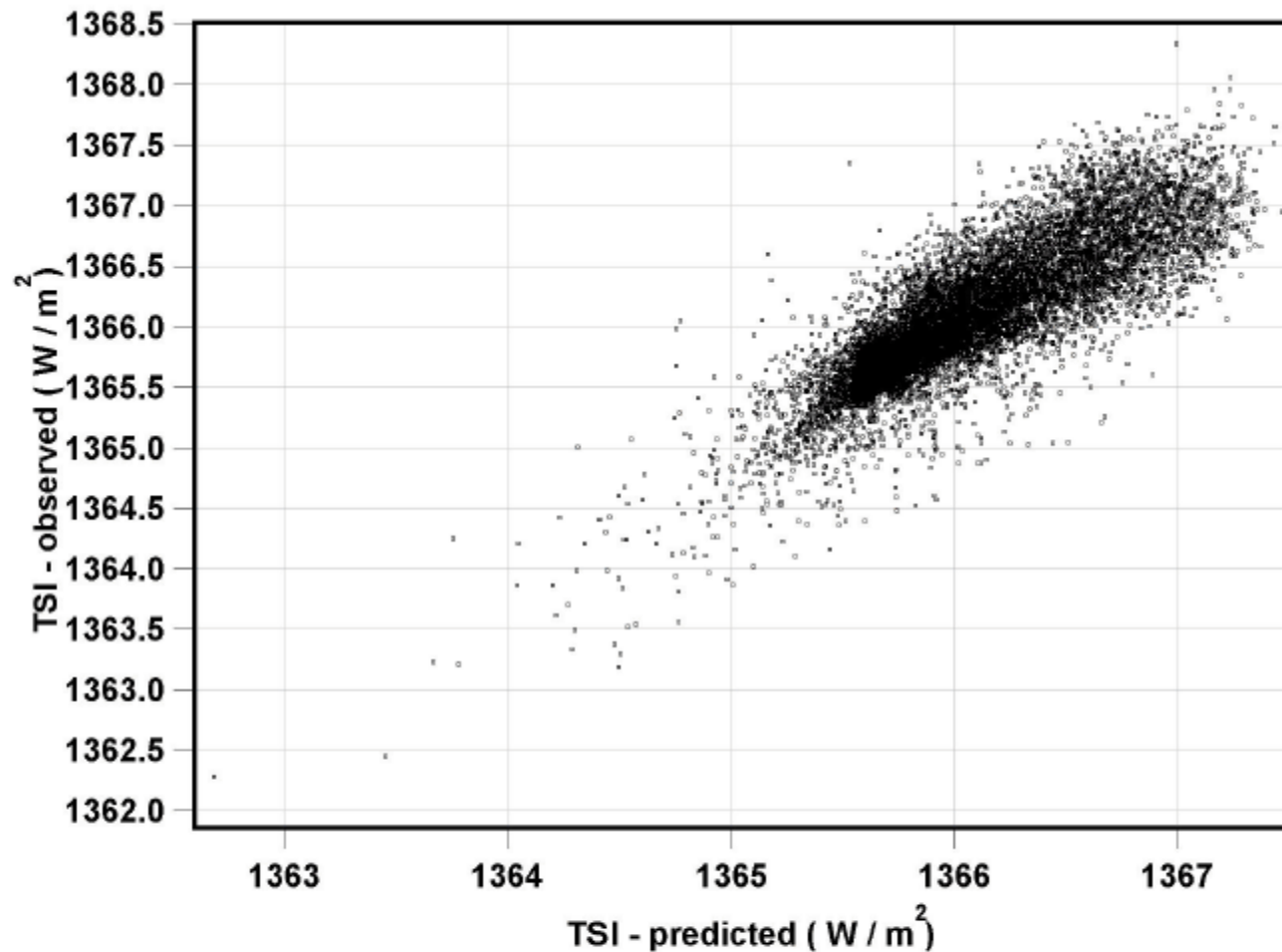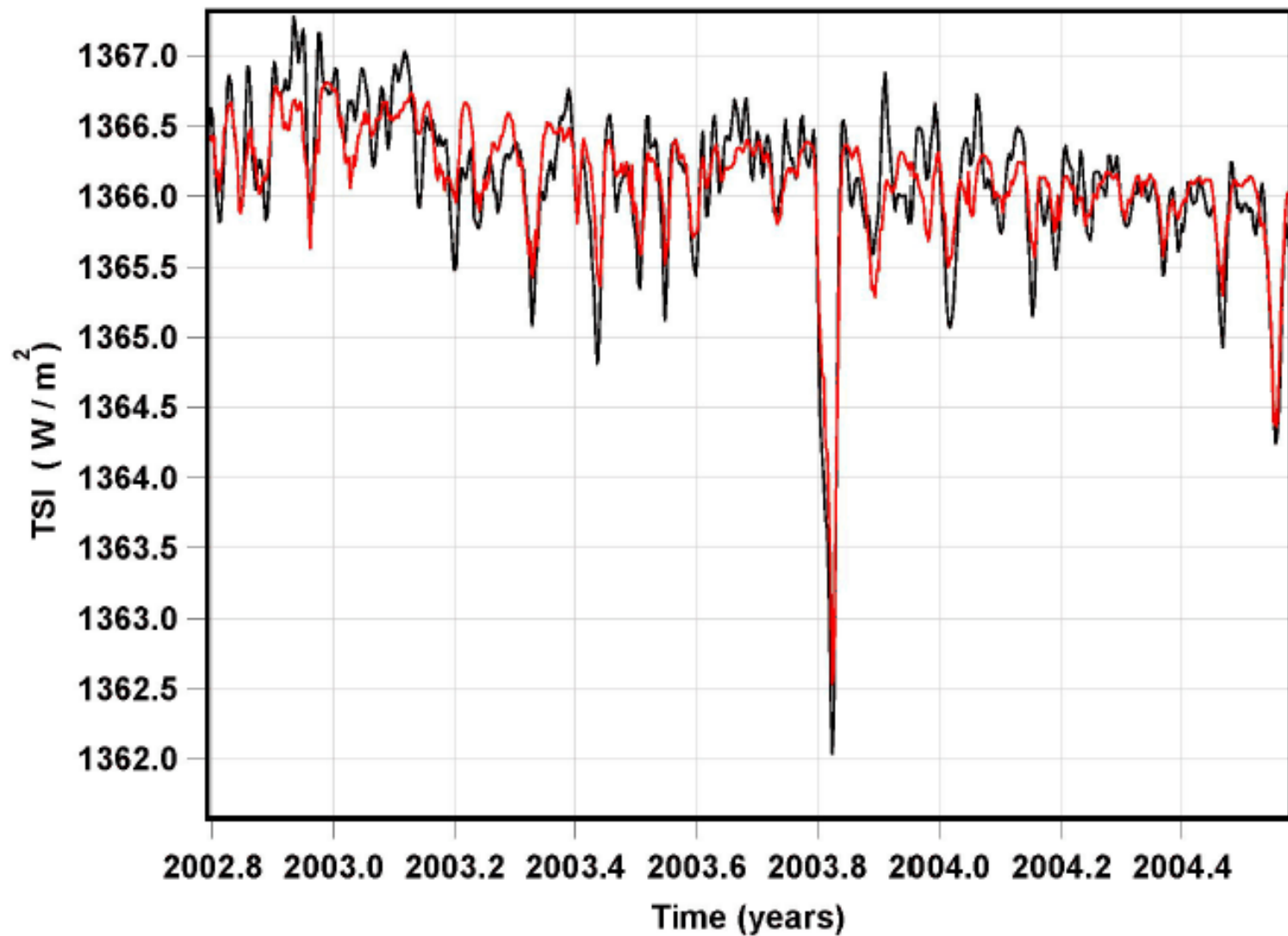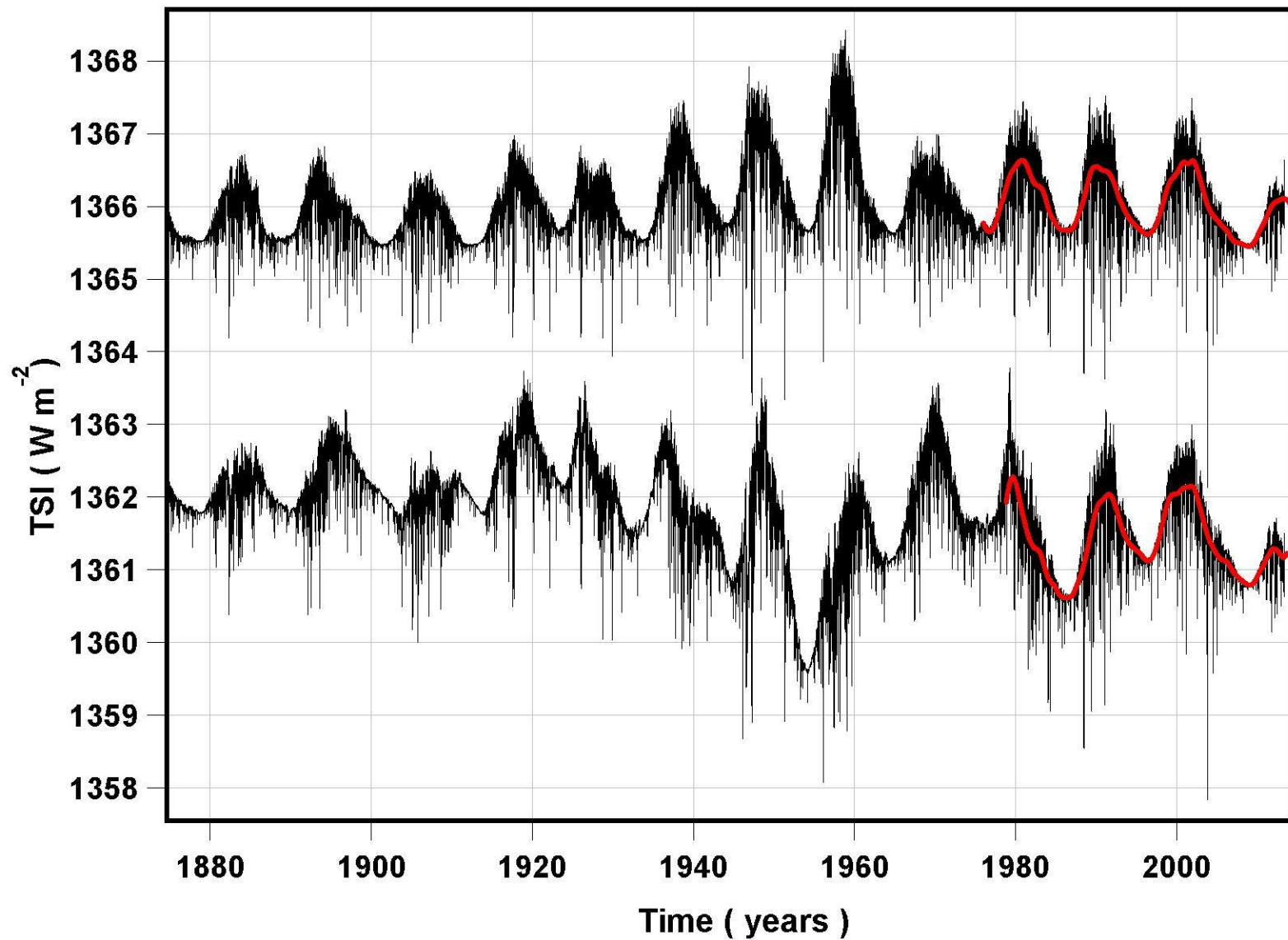| N | Type | Width | 1/Offset | $R_c$ |
|---|---|---|---|---|
| 1 | $E$ | 766.229110 | - | 0.85973057 |
| 2 | $E_t^+$ | 399.116082 | 27.103323 | 0.87599719 |
| 3 | $E^-$ | 2413.855149 | 1811.476298 | 0.88710035 |
| 4 | $E$ | 10.083440 | 10.148801 | 0.89328020 |
| 5 | $E^-$ | 137.021917 | 102.081918 | 0.89713987 |
| 6 | $E^+$ | 3509.347248 | 12.475225 | 0.90071914 |
| 7 | $E^+$ | 2474.591733 | 20.816371 | 0.90294868 |
| 8 | $E^-$ | 11.057883 | 10.679024 | 0.90478582 |
| 9 | $E_t^-$ | 12.699491 | 11.198220 | 0.90733729 |
| 10 | $E_t^+$ | 133.649130 | 188.570202 | 0.90864964 |
| 11 | $E_t^+$ | 189.738179 | 198.084112 | 0.91081665 |
| 12 | $E_t^-$ | 901.197452 | 28.196168 | 0.91161825 |
| 13 | $E_t^+$ | 5170.005637 | 140.675553 | 0.91222860 |
| 14 | $E_t^+$ | 2544.129259 | 31.414278 | 0.91284601 |
| 15 | $E_t^-$ | 33.610573 | 9.976299 | 0.91336121 |

**Table 8.** Modeling PMOD using PSI data. First 15 components from the greedy search for regression components. Parameters $W$ and $O$ are selected without restrictions.

| N | Type | Width | 1/Offset | Correlation |
|---|------|-------|----------|-------------|
| 1 | $E_t^+$ | 5170.005661 | 4419.832453 | 0.78193389 |
| 2 | $E_t^+$ | 2519.302633 | 1630.468830 | 0.83406728 |
| 3 | $E^+$ | 1503.913332 | 289.744282 | 0.85372950 |
| 4 | $E^-$ | 3406.542205 | 148.810590 | 0.87693631 |
| 5 | $E$ | 362.381918 | 27.341738 | 0.88562065 |
| 6 | $E_t^+$ | 430.606068 | 2183.735031 | 0.89144597 |
| 7 | $E^-$ | 2340.484301 | 300.893810 | 0.89771429 |
| 8 | $E^-$ | 165.079222 | 93.264999 | 0.90193772 |
| 9 | $E$ | 9.637961 | 10.600217 | 0.90548052 |
| 10 | $E^+$ | 2570.699109 | 27.380182 | 0.90766399 |
| 11 | $E_t^+$ | 15.521919 | 34.876082 | 0.90919779 |
| 12 | $E^-$ | 1931.046907 | 289.584617 | 0.91083821 |
| 13 | $E_t^+$ | 124.463532 | 112.081094 | 0.91196667 |
| 14 | $E_t^+$ | 95.268221 | 105.427923 | 0.91497974 |
| 15 | $E_t^+$ | 153.399257 | 109.674353 | 0.91650578 |

**Table 9.** Modeling ACRIM using PSI data. First 15 components from the greedy search for regression components. Parameters $W$ and $O$ are selected without restrictions.
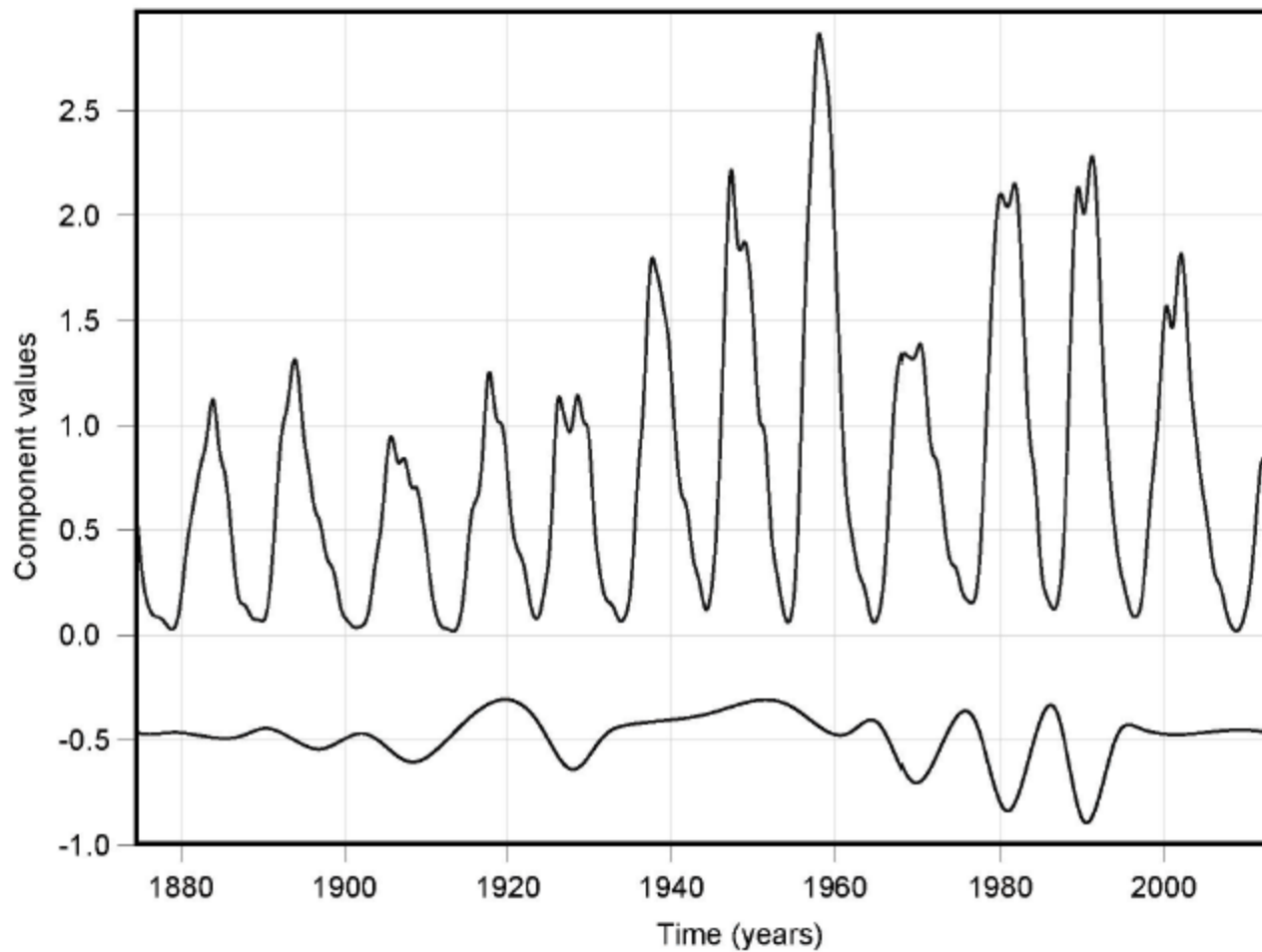
**Fig. 14.** The first (upper) and third (lower) components of the full scale PSI to PMOD model.

# Cross-prediction and sieving

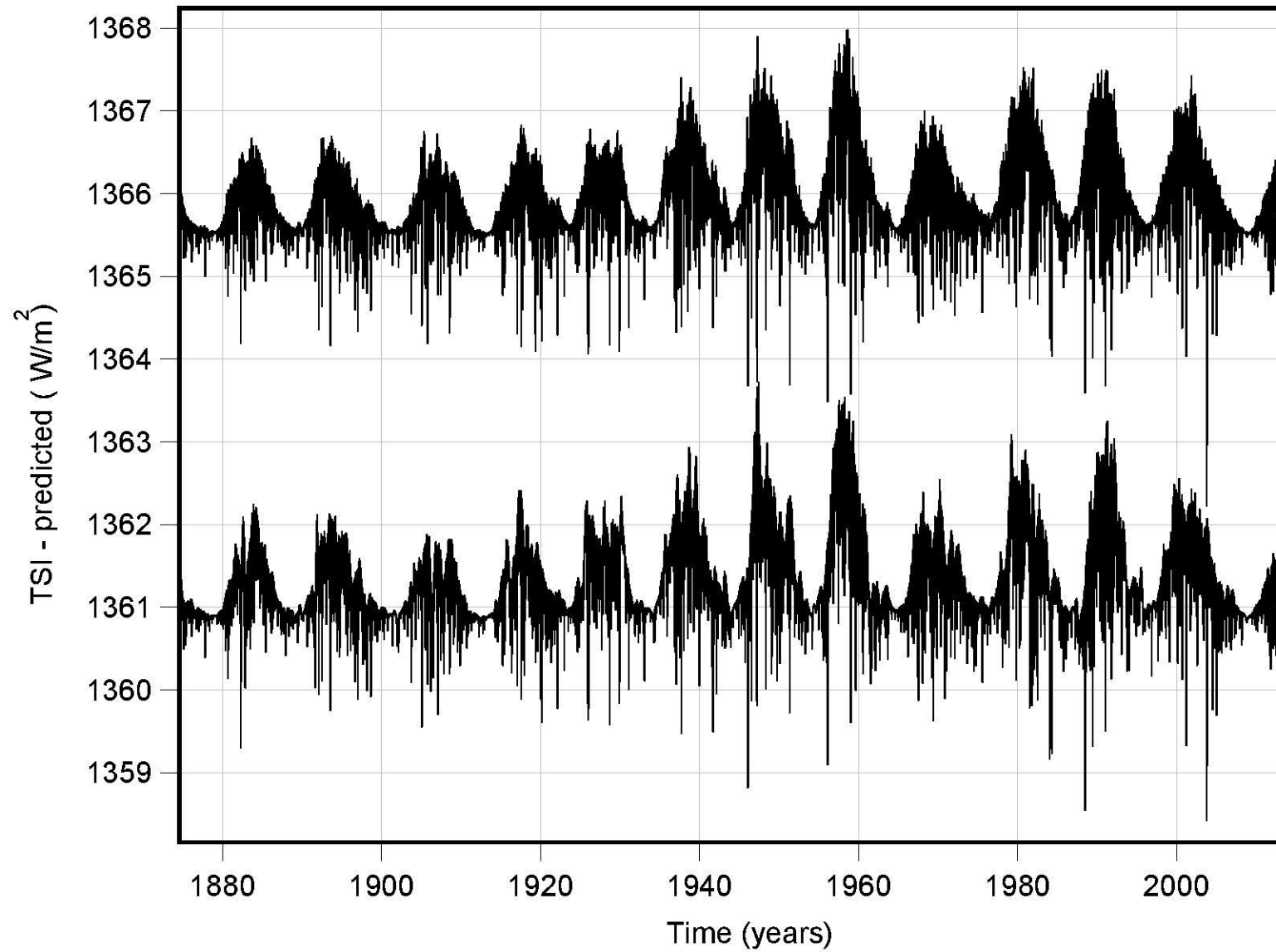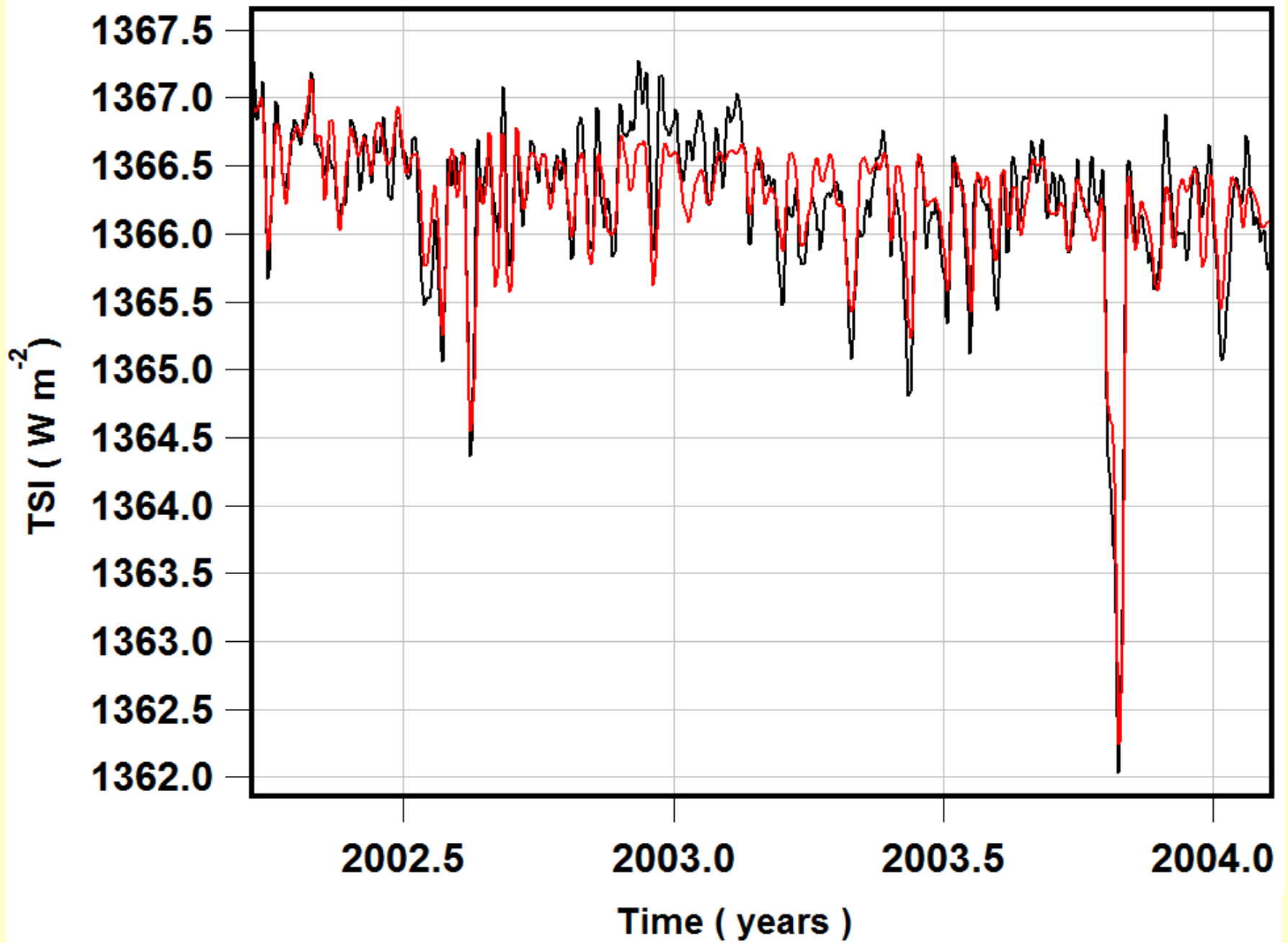$$R_c = \max_{W,O,M} \min(R_c^{I \to II}(W, O, M), R_c^{II \to I}(W, O, M))$$

W<2000
1/O<500
Check for colinearity

| Component | Type | Width | 1/Offset | Correlation |
|---|---|---|---|---|
| 1 | $E$ | 766.203360 | - | 0.85971218 |
| 2 | $E$ | 368.344694 | 27.205664 | 0.87183687 |
| 3 | $E$ | 11.076777 | 10.507722 | 0.87798356 |
| 4 | $E$ | 510.918696 | - | 0.87798440 |
| 5 | $E_t^-$ | 126.141442 | 85.206093 | 0.88078618 |
| 6 | $E^-$ | 137.407018 | 9.624734 | 0.88213771 |
| 7 | $E_t^-$ | 8.636770 | 8.672939 | 0.88481412 |
| 8 | $E^+$ | 356.511603 | 61.922806 | 0.88584926 |
| 9 | $E$ | 194.211239 | 8.970496 | 0.88652112 |
| 10 | $E^+$ | 124.335750 | 112.065571 | 0.88693209 |
| 11 | $E_t^+$ | 145.786806 | 194.401278 | 0.88695375 |
| 12 | $E^+$ | 145.491923 | 169.169762 | 0.88977830 |
| 13 | $E^+$ | 815.678995 | 82.630138 | 0.89115533 |
| 14 | $E_t^+$ | 99.363542 | 89.764059 | 0.89158876 |
| 15 | $E$ | 173.673683 | 77.786024 | 0.89236231 |

**Table 10.** Modeling PMOD data using PSI as a proxy. First 15 components of the regression model where all kind of components were allowed.
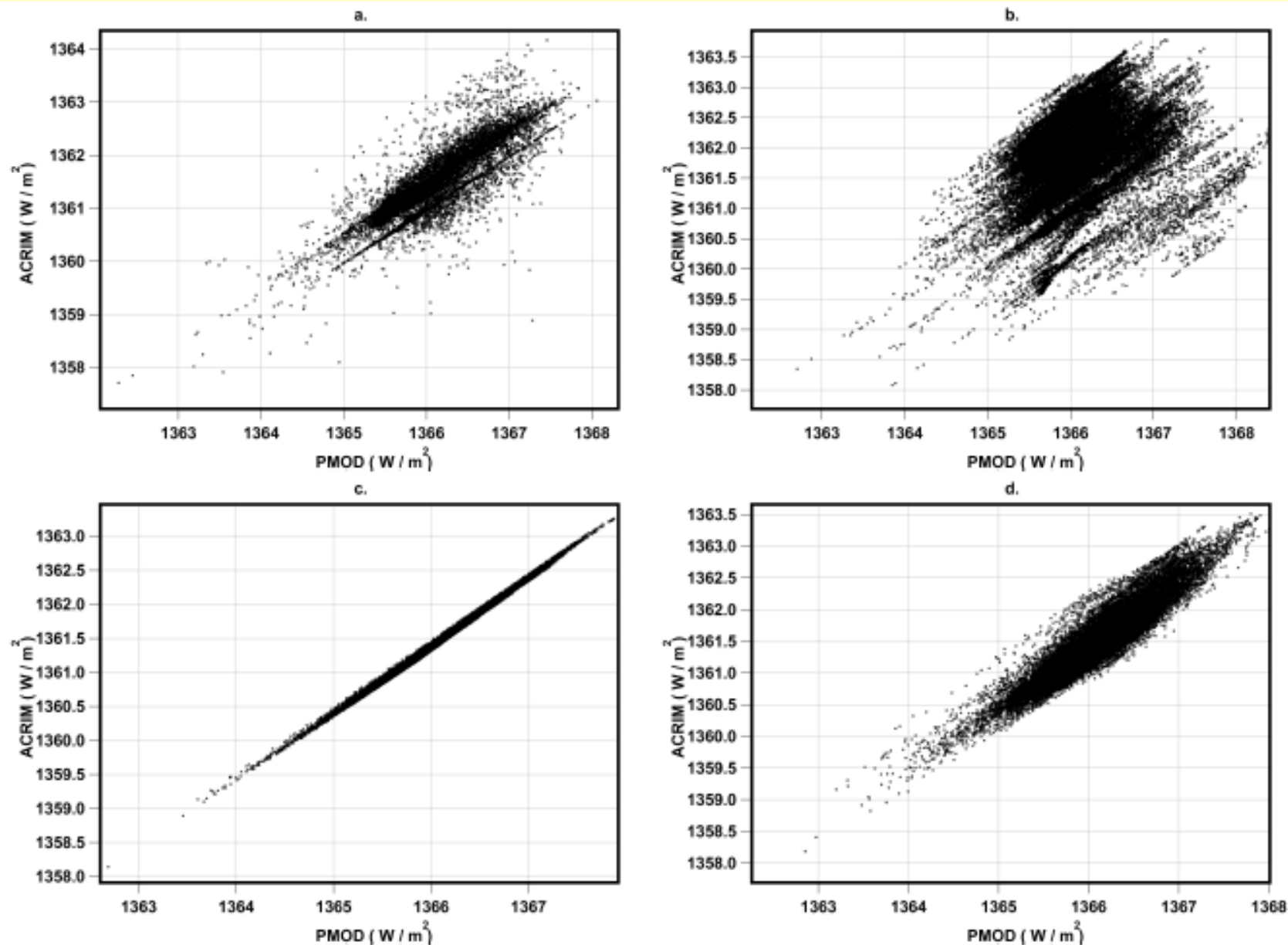
**Fig. 10.** Correlations between PMOD and ACRIM based predictions. Upper panel: a. crossplot of the original target data sets PMOD versus ACRIM ($R_c = 0.8553$), b. crossplot for solution without restrictions ($R_c = 0.3521$, Fig. 11); lower panel: c. crossplot of the predictions based on simplest models ($R_c = 0.9985$), d. crossplot between combined models ($R_c = 0.9384$).

# FDC Summary

- The correlations between time series depend on frequency band involved.

- The models based on smoothed, enveloped etc components can be used for prediction, interpolation and stitching of different fragments.

- Current state of affairs with TSI measurement and composite building is very complicated.

- Probably the best data product which can be given to climatologists is nearly primitive model based only on small number of components. It will with high probability correlate with true TSI at the level of $R_c$=0.85 or so.

# Final words

- $D^2$ method - periodicity, multiperiodicity, period stability, coherence length.

- CF method - tranients in phase, modulations etc.

- FDC method – filter and correlate!