

# Kernels for Periodic Time Series Arising in Astronomy

Gabriel Wachman

Tufts University

Roni Khardon

Tufts University

Pavlos Protopapas

Harvard CfA and IIC

Charles A. Alcock

Harvard CfA

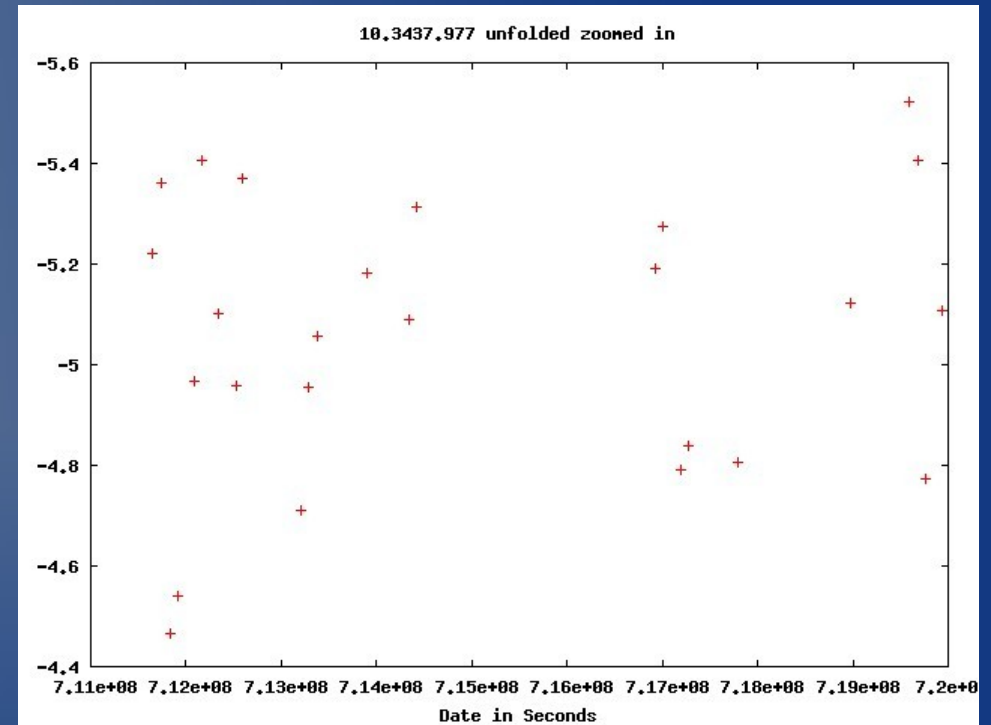
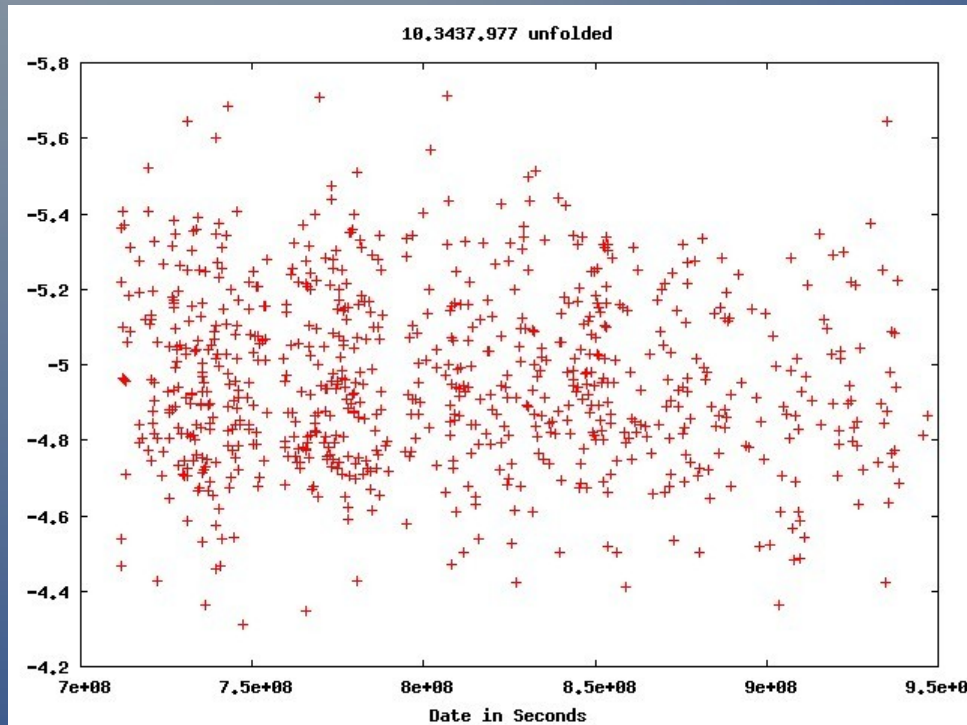
# Acknowledgments

- Initiative for Innovative Computing
- NSF grant IIS-080340
- Odyssey cluster supported by the FAS Research Computing Group at Harvard
- Tufts High-performance computing research cluster

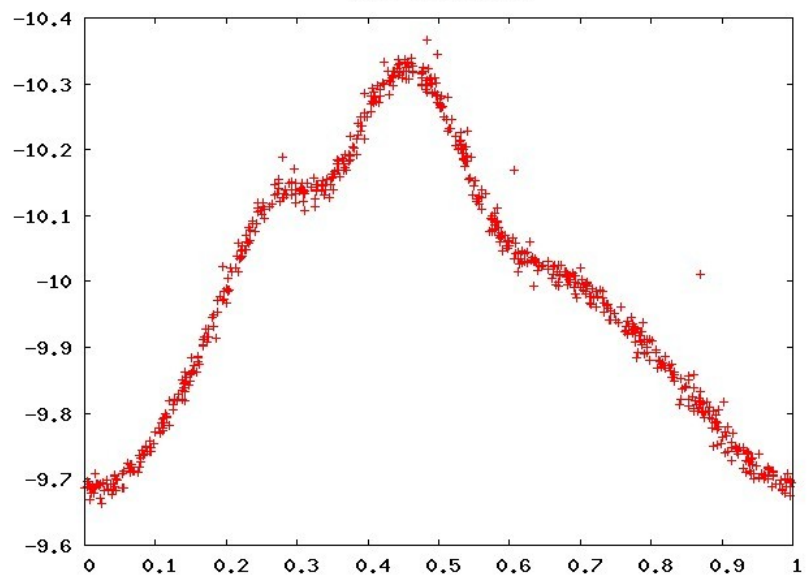
# Motivation

- Astronomy surveys contain millions of objects
- Significant time requirement to do manual classification
- Many stars remain unclassified
- Stars are represented as a *time series*

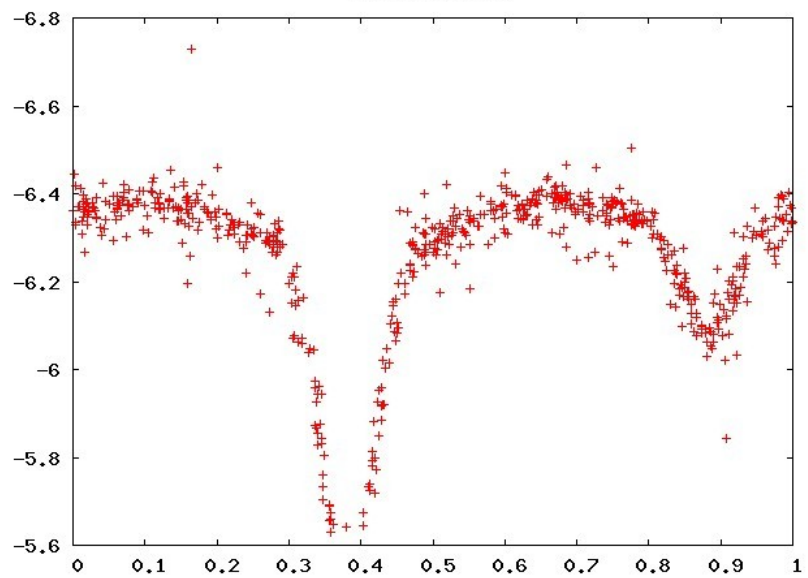
# Periodic Variable Stars



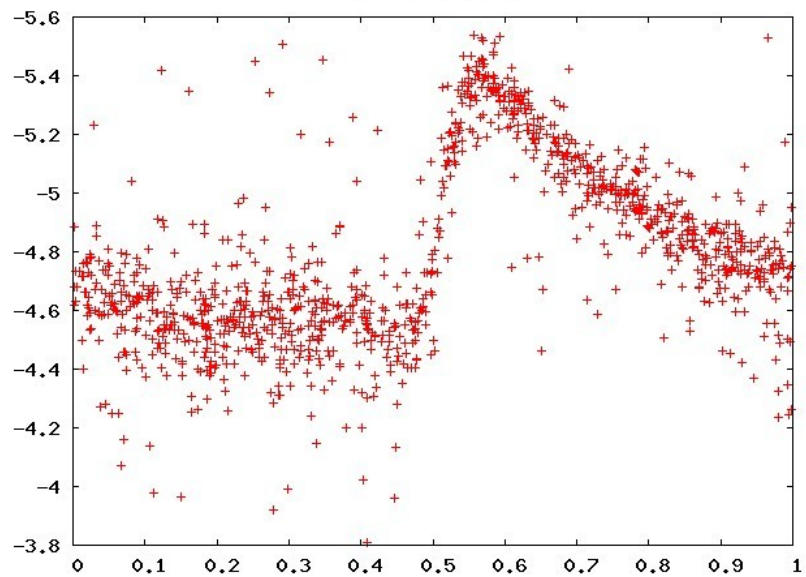
CEPH F1.3441.15



EB F1.3442.233

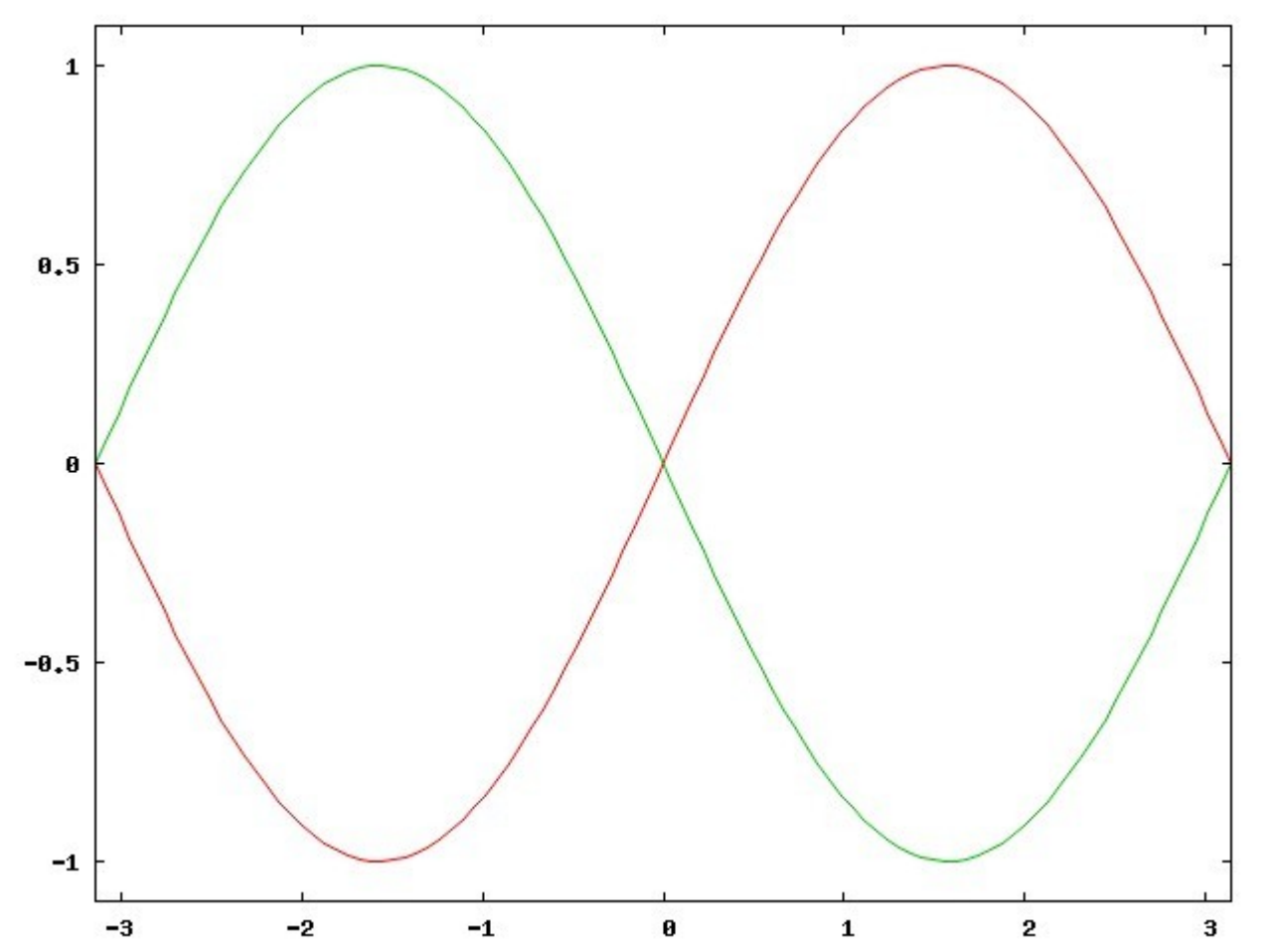


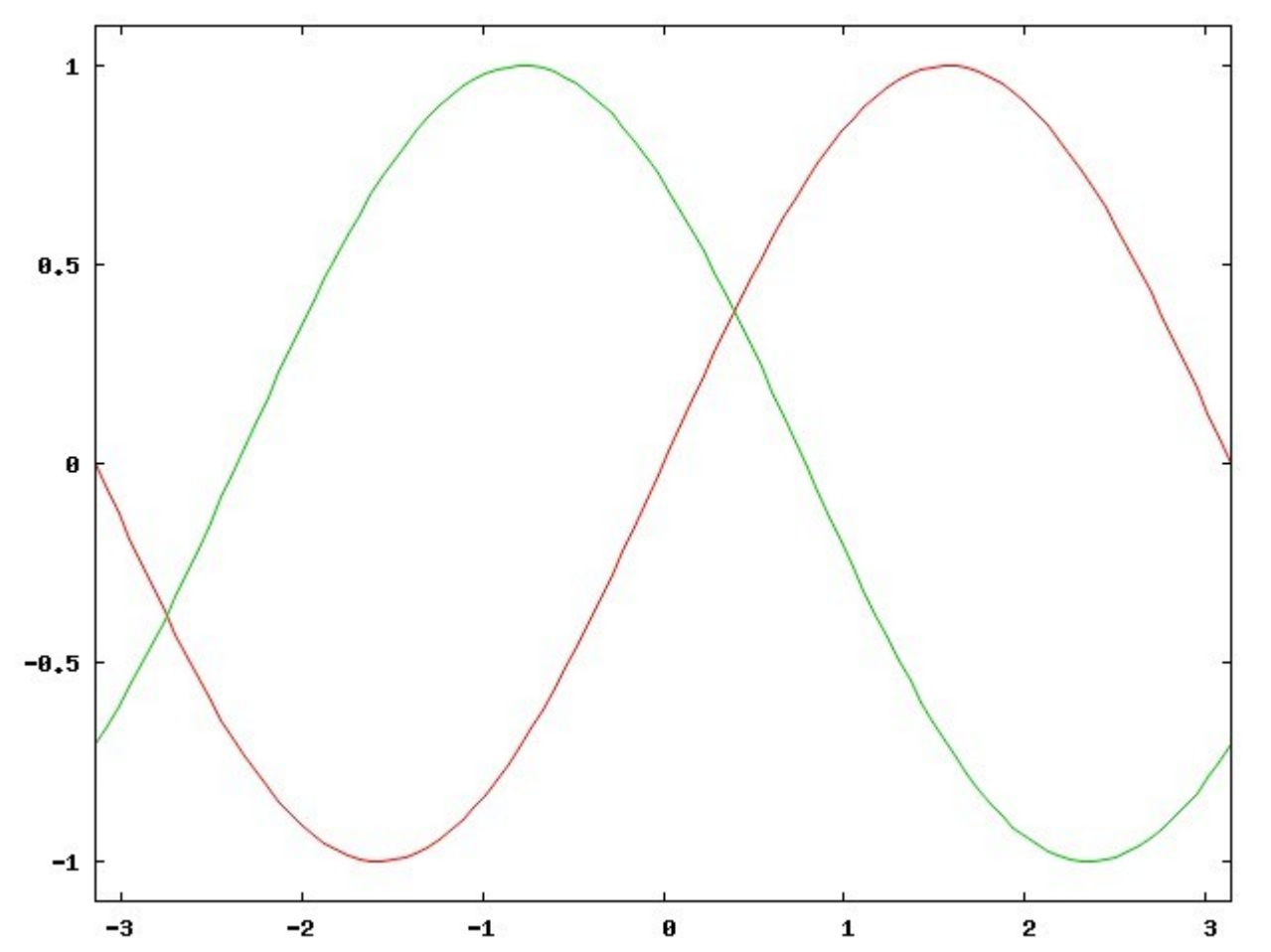
RRL F1.3449.1327



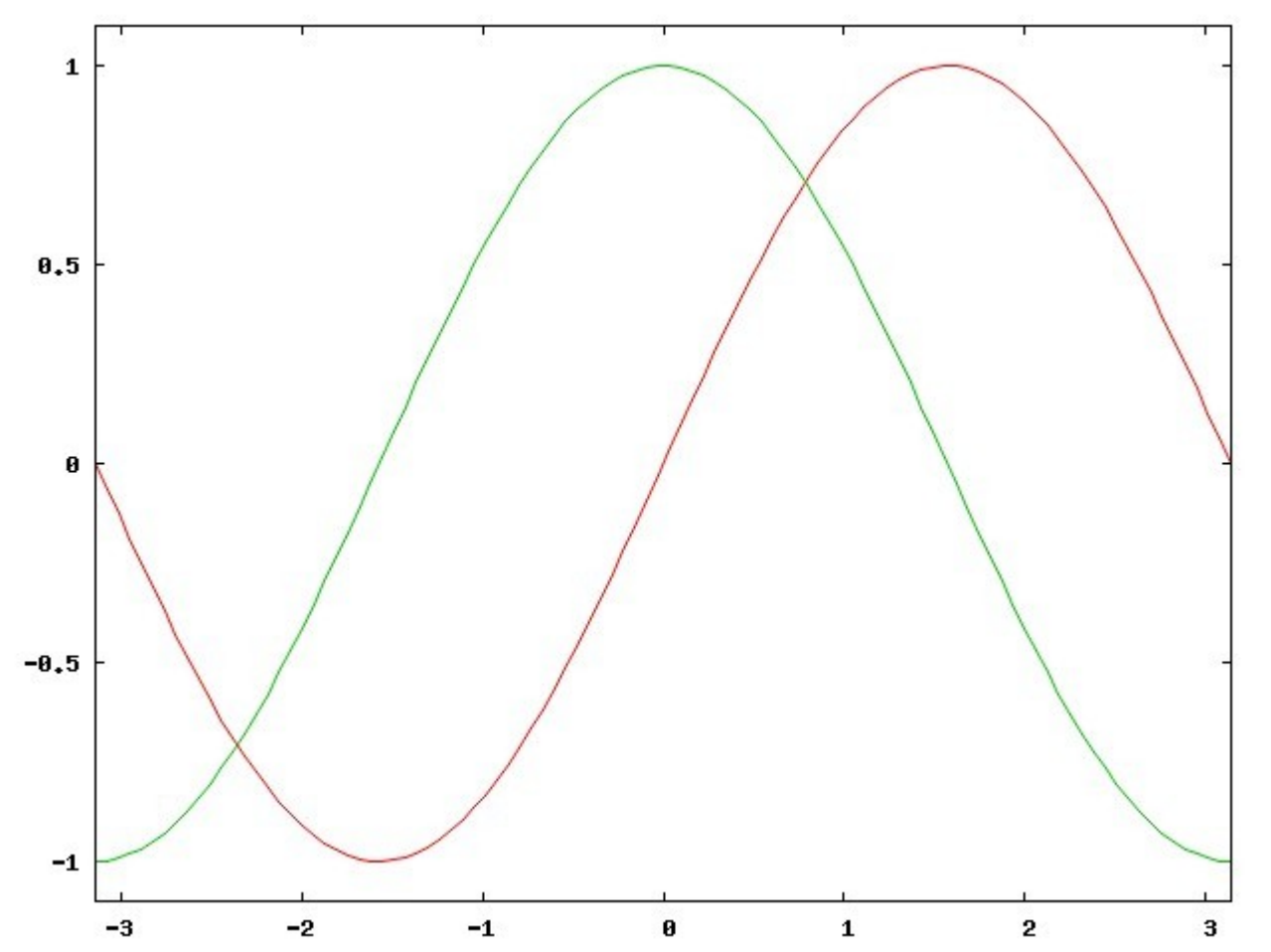
# Agenda

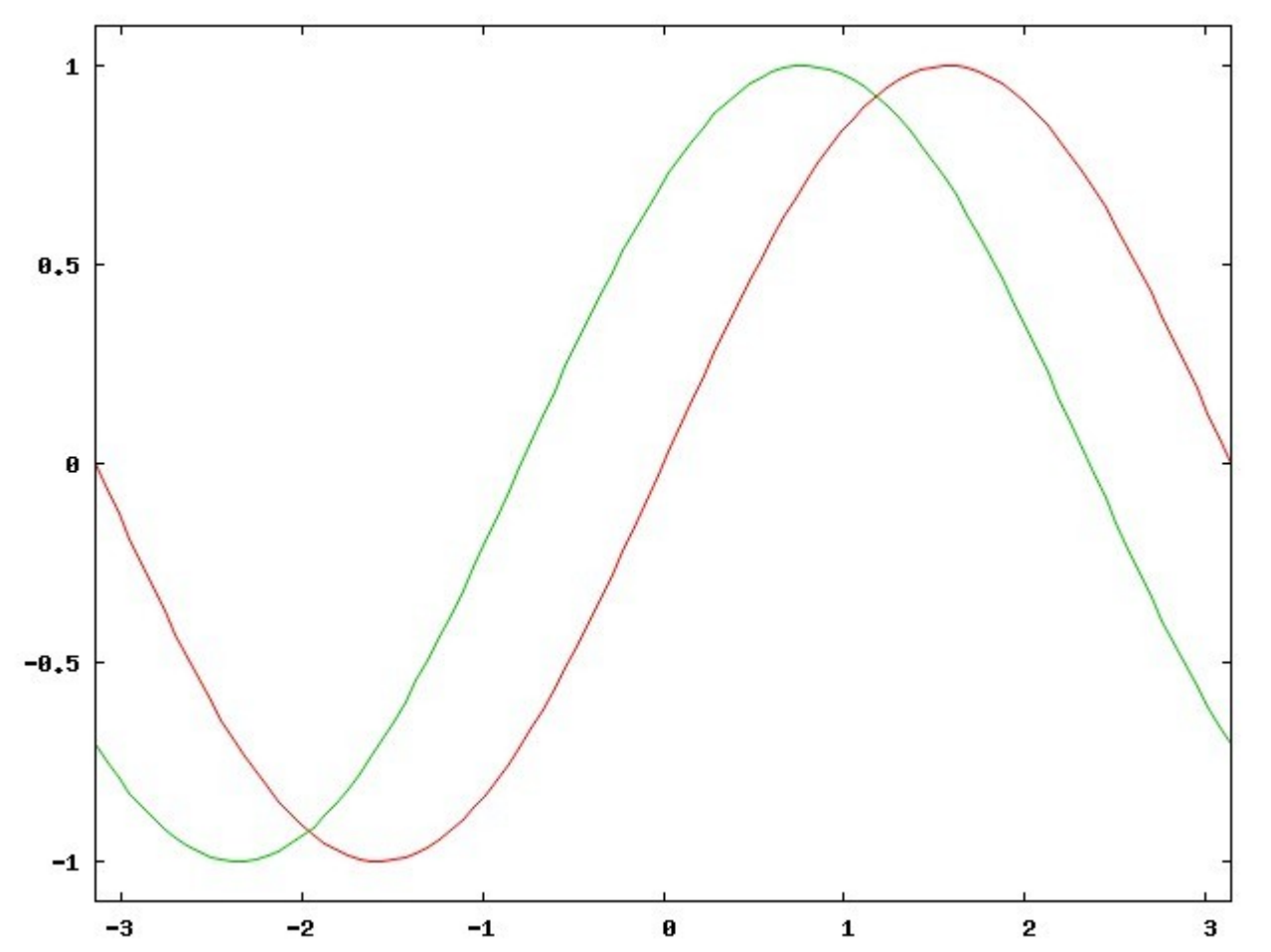
- Introduction and Motivation
- Similarity Measure for Time Series
- Kernel for Periodic Time Series
- Experiments and Results
- Big Picture
- Summary

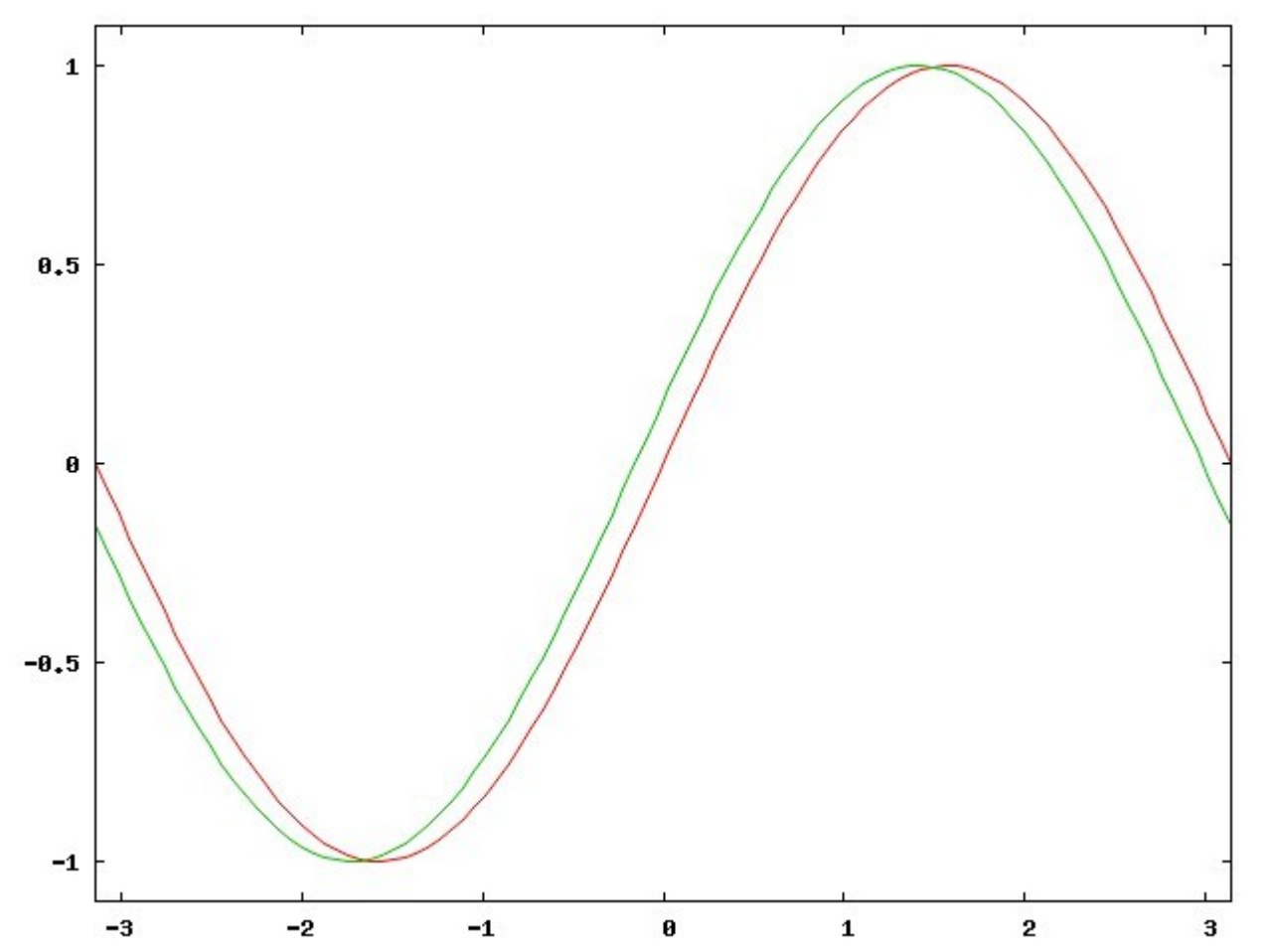












# Kernel for Periodic Time Series

- We want to capture what it means for two time series to have the same “shape”
- S1 [Protopapas06]:

$$\max_s \langle x, y_{+s} \rangle$$

- Does exactly what we want
- Can compute using FFT in  $O(n \log n)$
- Is it a kernel?

# Properties of $S_1$

*Theorem 1:*  $S_1$  satisfies the Cauchy Schwartz inequality.

*Theorem 2:* Can construct a distance measure using  $S_1$  that satisfies triangle inequality.

*Theorem 3:* Any  $3 \times 3$  Gram matrix of  $S_1$  is positive semidefinite.

*Theorem 4:*  $S_1$  is not positive semidefinite.

# Kernel for Periodic Time Series

- K1:

$$\sum_{s=1}^n e^{\gamma \langle x, y_{+s} \rangle}$$

- Intuitively approximates maximum alignment
- Works as well in practice

# Kernel for Time Series

- *Theorem 5:*  $K_1$  is positive semidefinite
- *Theorem 6:*  $K_1$  is computable in  $O(n \log n)$

# Shape Data

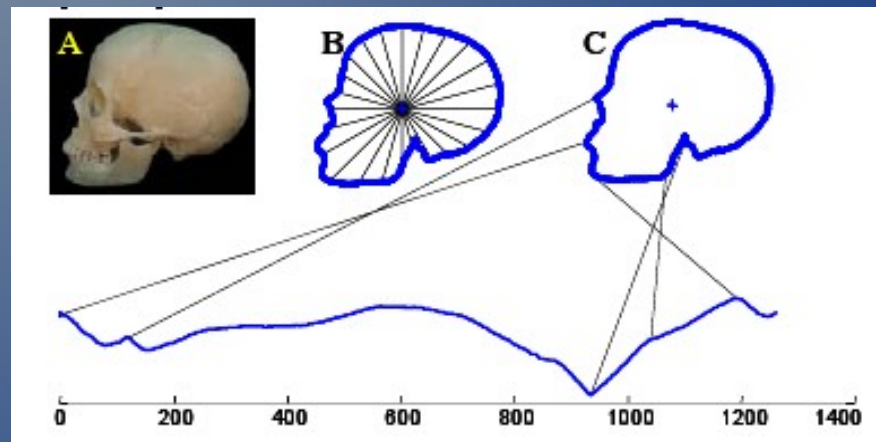


Image from **E. Keogh, L. Wei, X. Xi, S. Lee, M. Vlachos. LB\_Keogh Supports Exact Indexing of Shapes under Rotation Invariance with Arbitrary Representations and Distance Measures. VLDB 2006.**



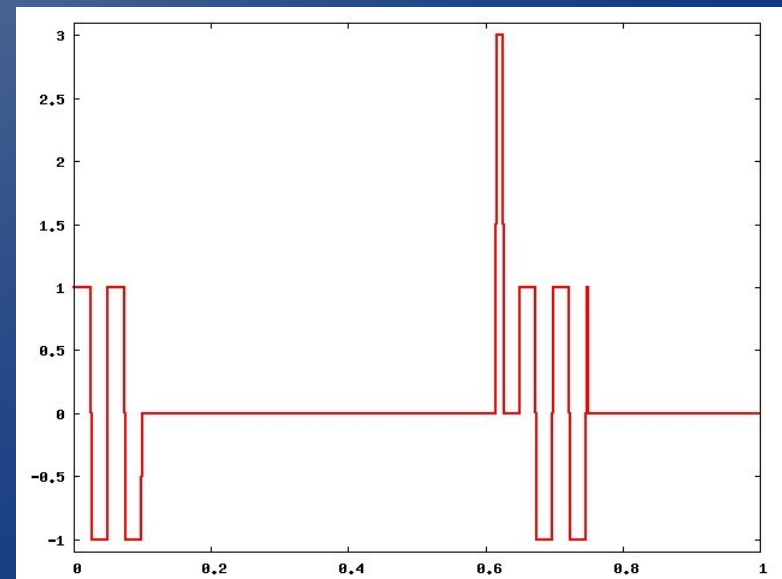
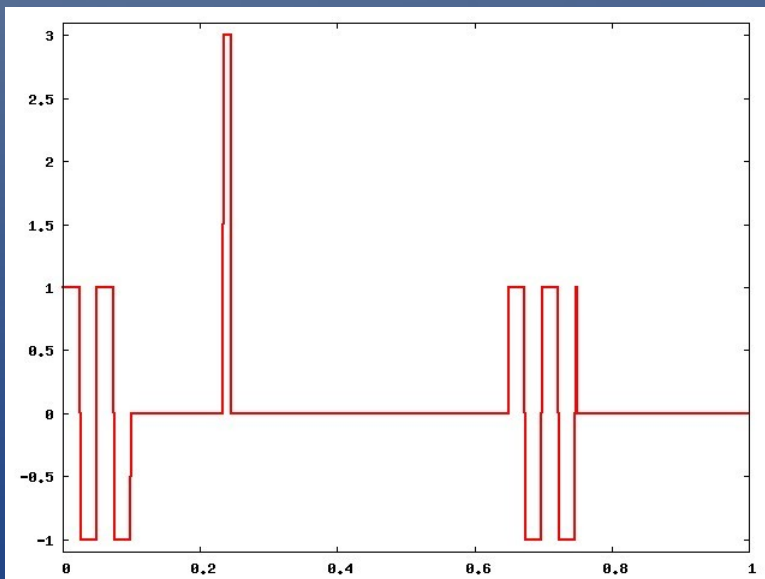
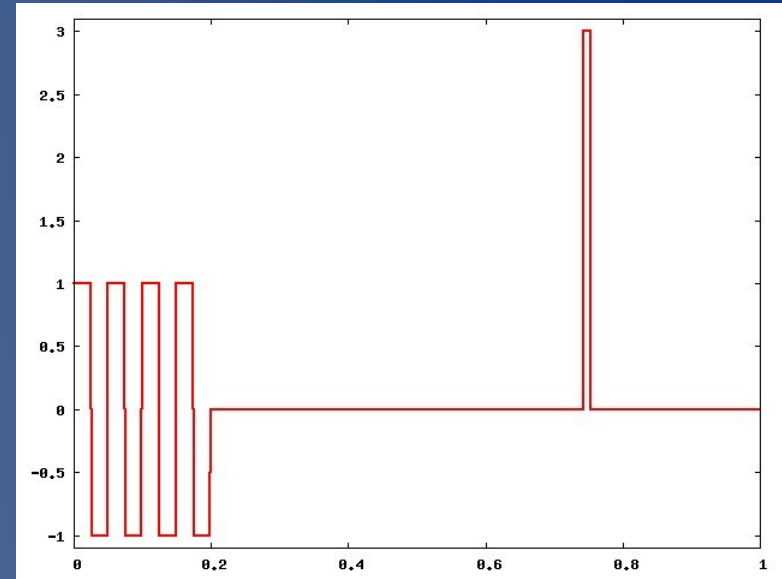
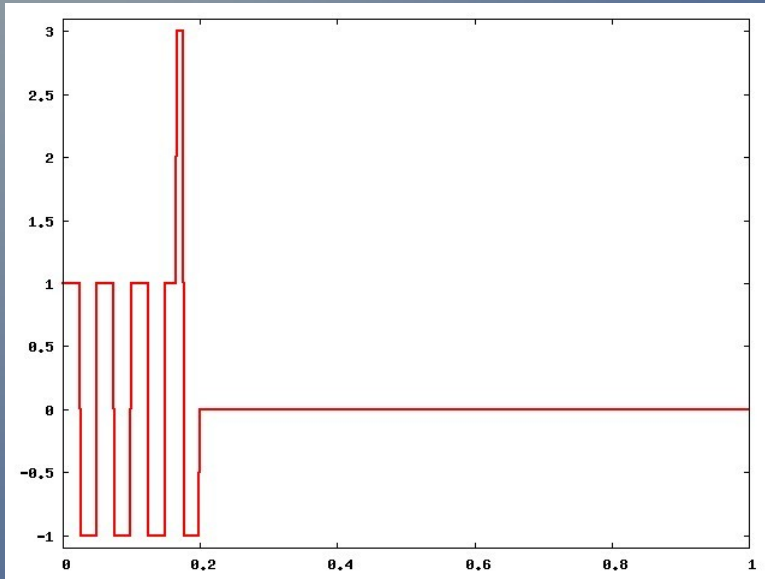
# Experiments on Shape Data

	1-NN S1	1-NN K1	1-NN DTW	1-NN UP	SVM ED	SVM UP	SVM K1
A	$0.54 \pm 0.06$	$0.54 \pm 0.08$	$0.33 \pm 0.06$	$0.49 \pm 0.05$	$0.20 \pm 0.05$	$0.41 \pm 0.05$	$0.63 \pm 0.04$
B	$0.73 \pm 0.04$	$0.73 \pm 0.04$	$0.59 \pm 0.08$	$0.70 \pm 0.07$	$0.40 \pm 0.10$	$0.65 \pm 0.08$	$0.76 \pm 0.08$
I	$0.98 \pm 0.01$	$0.98 \pm 0.01$	$0.84 \pm 0.03$	$0.97 \pm 0.02$	$0.47 \pm 0.03$	$0.80 \pm 0.02$	$0.91 \pm 0.02$

# Experiments on Shape Data

	1-NN S1	1-NN K1	1-NN DTW	1-NN UP	SVM ED	SVM UP	SVM K1
A	$0.54 \pm 0.06$	$0.54 \pm 0.08$	$0.33 \pm 0.06$	$0.49 \pm 0.05$	$0.20 \pm 0.05$	$0.41 \pm 0.05$	$0.63 \pm 0.04$
B	$0.73 \pm 0.04$	$0.73 \pm 0.04$	$0.59 \pm 0.08$	$0.70 \pm 0.07$	$0.40 \pm 0.10$	$0.65 \pm 0.08$	$0.76 \pm 0.08$
I	$0.98 \pm 0.01$	$0.98 \pm 0.01$	$0.84 \pm 0.03$	$0.97 \pm 0.02$	$0.47 \pm 0.03$	$0.80 \pm 0.02$	$0.91 \pm 0.02$

# Breaking Universal Phasing



# Experiments on Artificial Data

	1-NN S1	1-NN K1	1-NN UP	SVM UP	SVM K1
Artificial	$0.99 \pm 0.02$	$1.00 \pm 0.00$	$0.65 \pm 0.04$	$0.50 \pm 0.07$	$0.997 \pm 0.00$
Artificial w/Noise	$0.84 \pm 0.14$	$0.84 \pm 0.12$	$0.61 \pm 0.09$	$0.53 \pm 0.12$	$0.90 \pm 0.05$

# Astronomy Data

## OGLE[Ud97,So03]

- 14087 periodic variables of type Cepheid, EB, RRL
- Periods given
- *Features* are average brightness, average color, and period

# Experiments on Astronomy Data

	1-NN	SVM		1-NN	SVM
S1	$0.844 \pm 0.011$	$0.680 \pm 0.011$	features + S1	$0.991 \pm 0.002$	$0.998 \pm 0.001$
K1	$0.901 \pm 0.008$	$0.947 \pm 0.005$	features + K1	$0.992 \pm 0.002$	$0.998 \pm 0.001$
UP	$0.827 \pm 0.010$	$0.851 \pm 0.006$	features + UP	$0.991 \pm 0.002$	$0.997 \pm 0.001$
			features	$0.938 \pm 0.006$	$0.974 \pm 0.004$

# Producing a Clean Catalog

	Ceph	EB	RRL
Ceph	3416	1	13
EB	0	3389	0
RRL	9	0	7259

	Ceph	EB	RRL
Ceph	3382	1	3
EB	0	3364	0
RRL	1	0	7195

	Ceph	EB	RRL
Ceph	3352	1	0
EB	0	3312	0
RRL	0	0	7138

# Test Survey: MACHO

- ~25 million stars
- ~50,000 are periodic variables
- Two primary issues:
  - Eliminating *just* the other 24,950,000 stars
  - Finding the periods



# Astronomy Classification

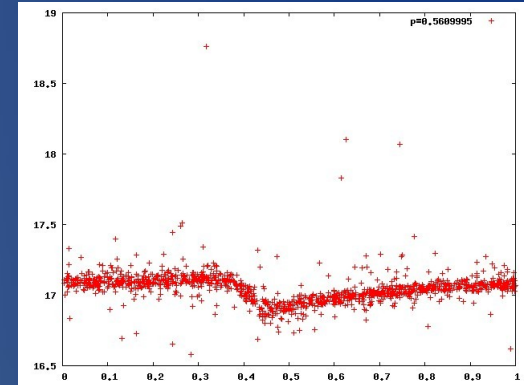
- Eliminate non-variables, non-periodic, and any stars not of type Cepheid, RRL, or EB
- Find periods automatically
- Scalable approach

# Period Estimation

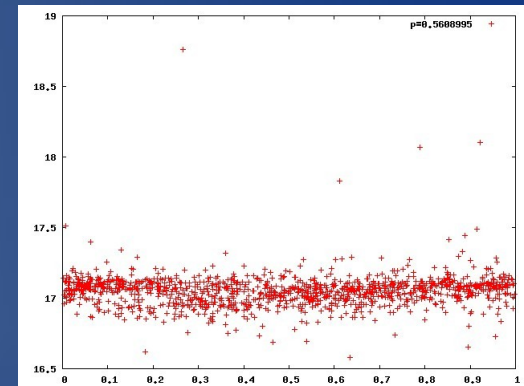
- Challenges

- Sensitivity of period
- Range of potential periods
  - 0.25d – 20d
- Computational cost

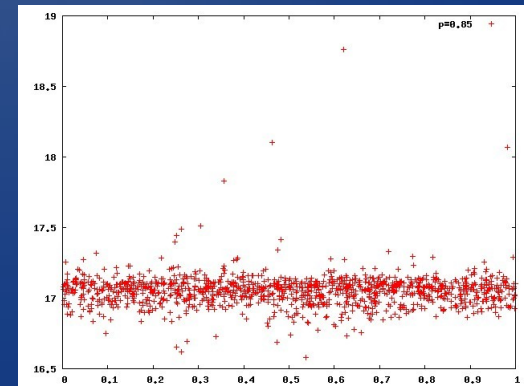
0.5609995



0.5608998



0.85



# Summary

- Insight into cross-correlation
- Kernel for periodic time series
  - Similar to cross-correlation
  - Efficiently computable
  - Works well in multiple domains
- Integral piece of larger astronomy project
- Application of maximum alignment approximation to other domains

# Questions

[gabriel.wachman@tufts.edu](mailto:gabriel.wachman@tufts.edu)