

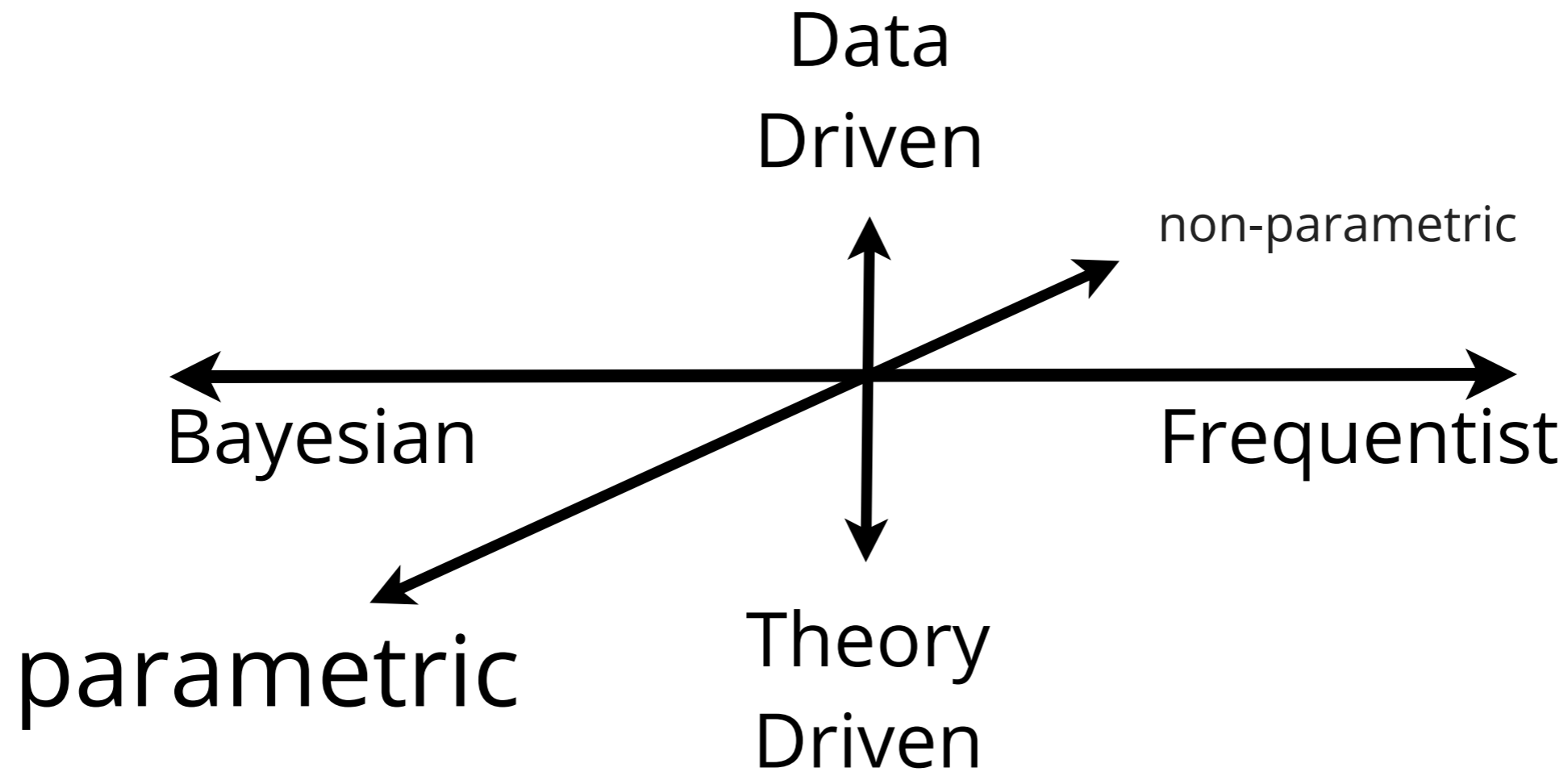
Data-Driven Astronomical Inference

Josh Bloom, UC Berkeley



@profjsb

Inference Space



Hardware	laptops → NERSC
Software	Python/Scipy, R, ...
Carbonware	astro grad students, postdocs

Bayesian Distance Ladder

Pulsational Variables: Period-Luminosity Relation

$$m_{ij} = \mu_i + M_{0j} + \alpha_j \log_{10} (P_i/P_0) + E(B - V)_i \times [R_V \times a (1/\lambda_i) + b (1/\lambda_i)] + \epsilon_{ij}$$

i indexes over individual stars
j indexes over wavebands
a and *b* are fixed constants at each waveband

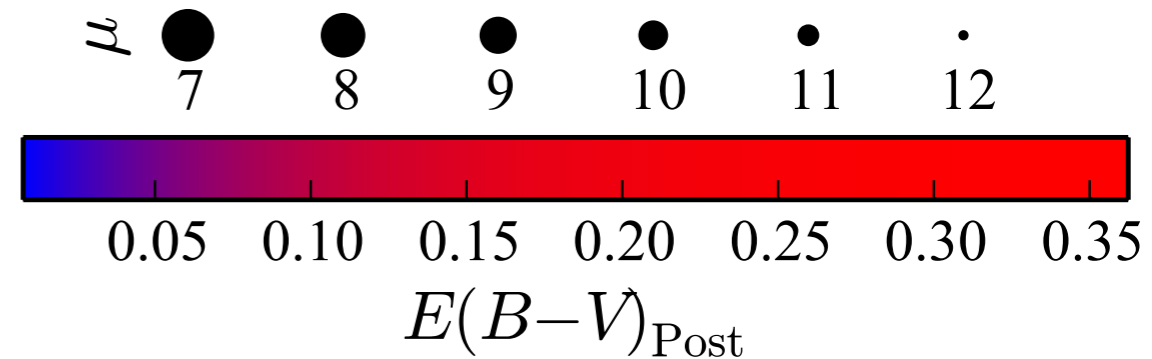
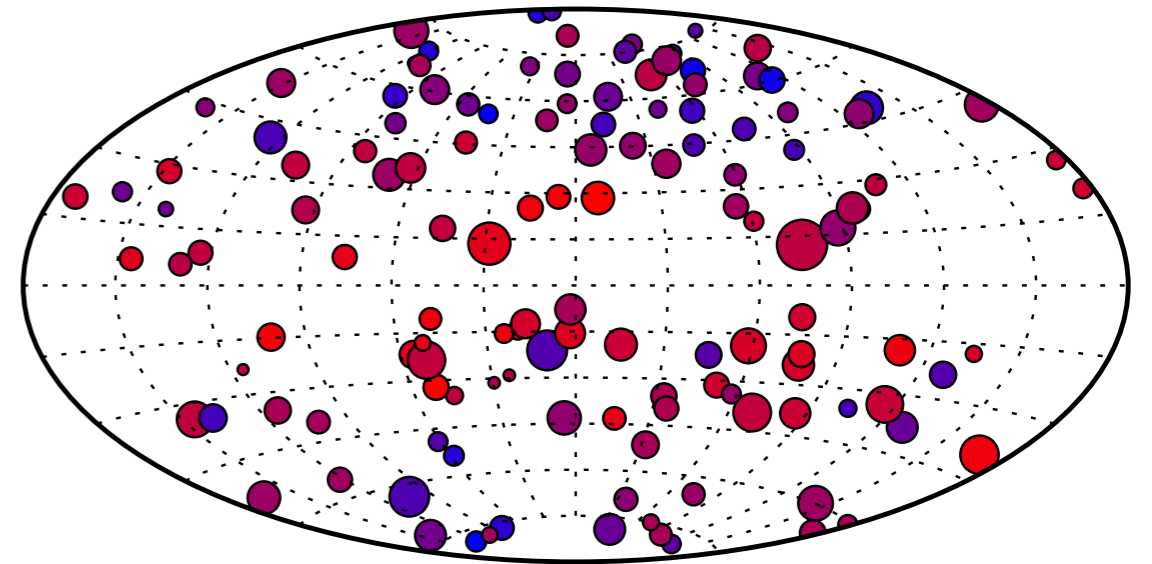
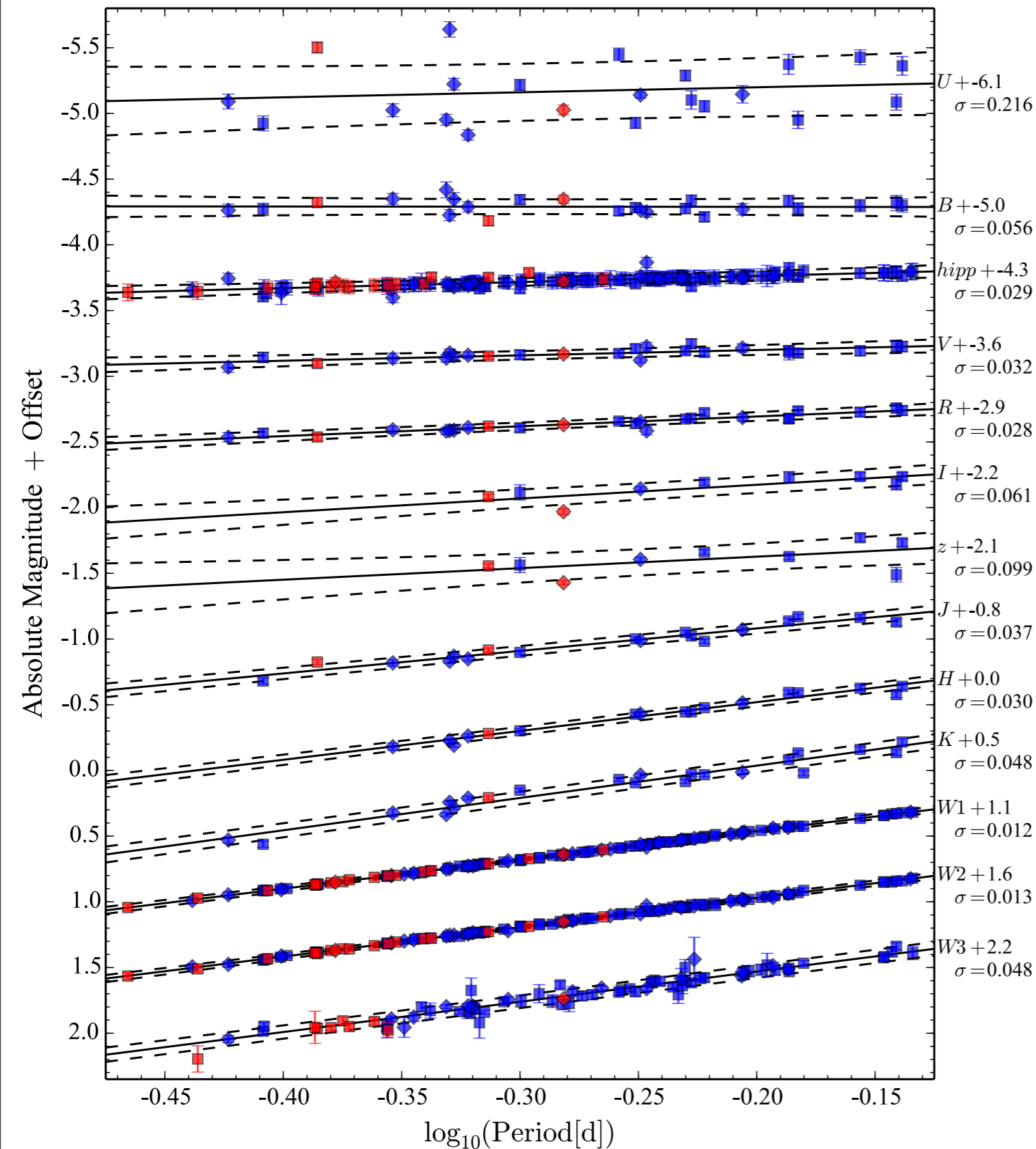
Data 134 RR Lyrae (WISE, Hipparcos, UVIORJHK)

Fit 307 dimensional model parameter inference

- deterministic MCMC model
- ~6 days for a single run (one core)
- parallelism for convergence tests

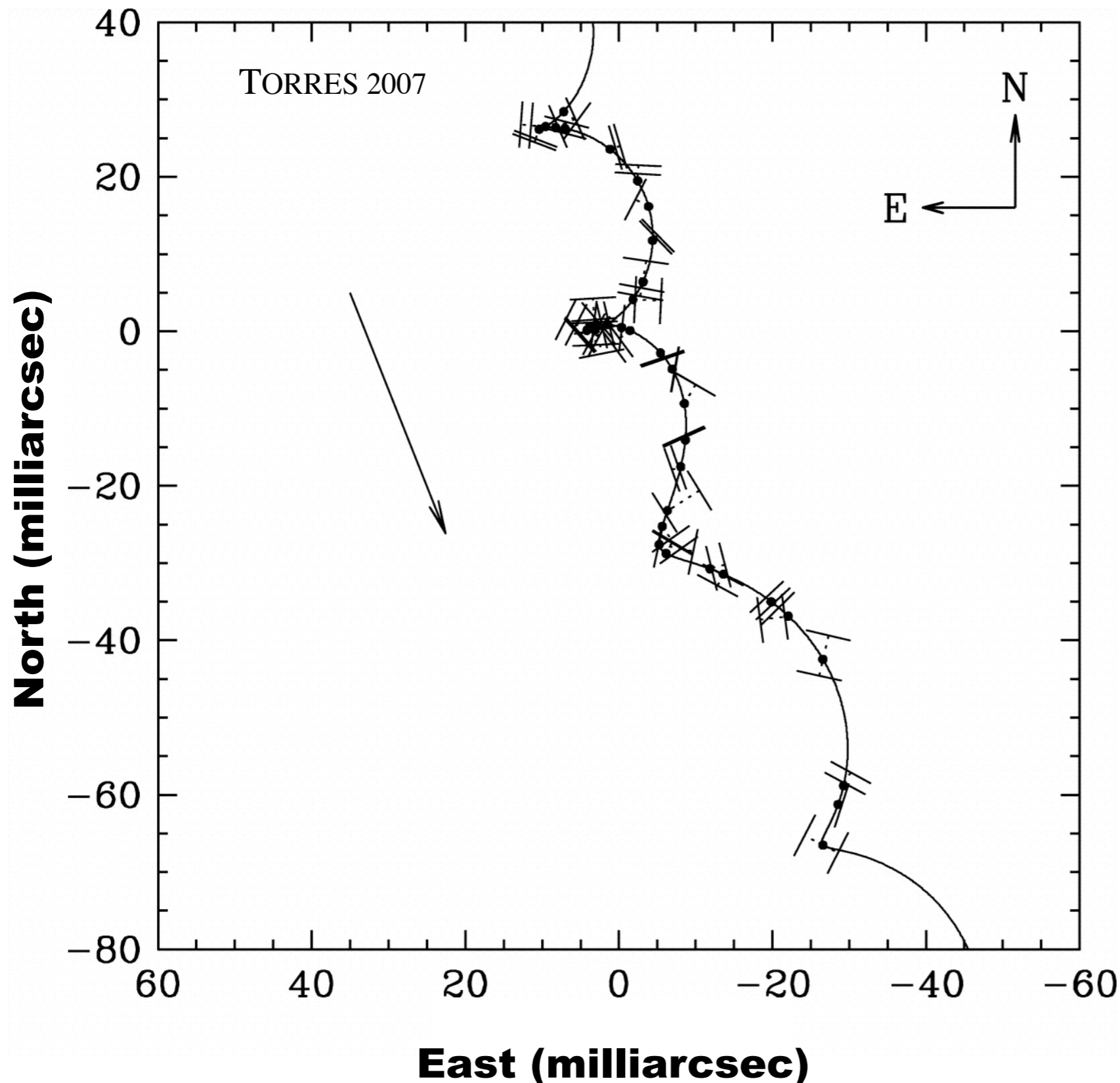
Klein+12; Klein,JSB+14,

Bayesian Distance Ladder



- Approaching 1% distance uncertainty
- Precision 3D dust measurements

Bayesian Astrometry

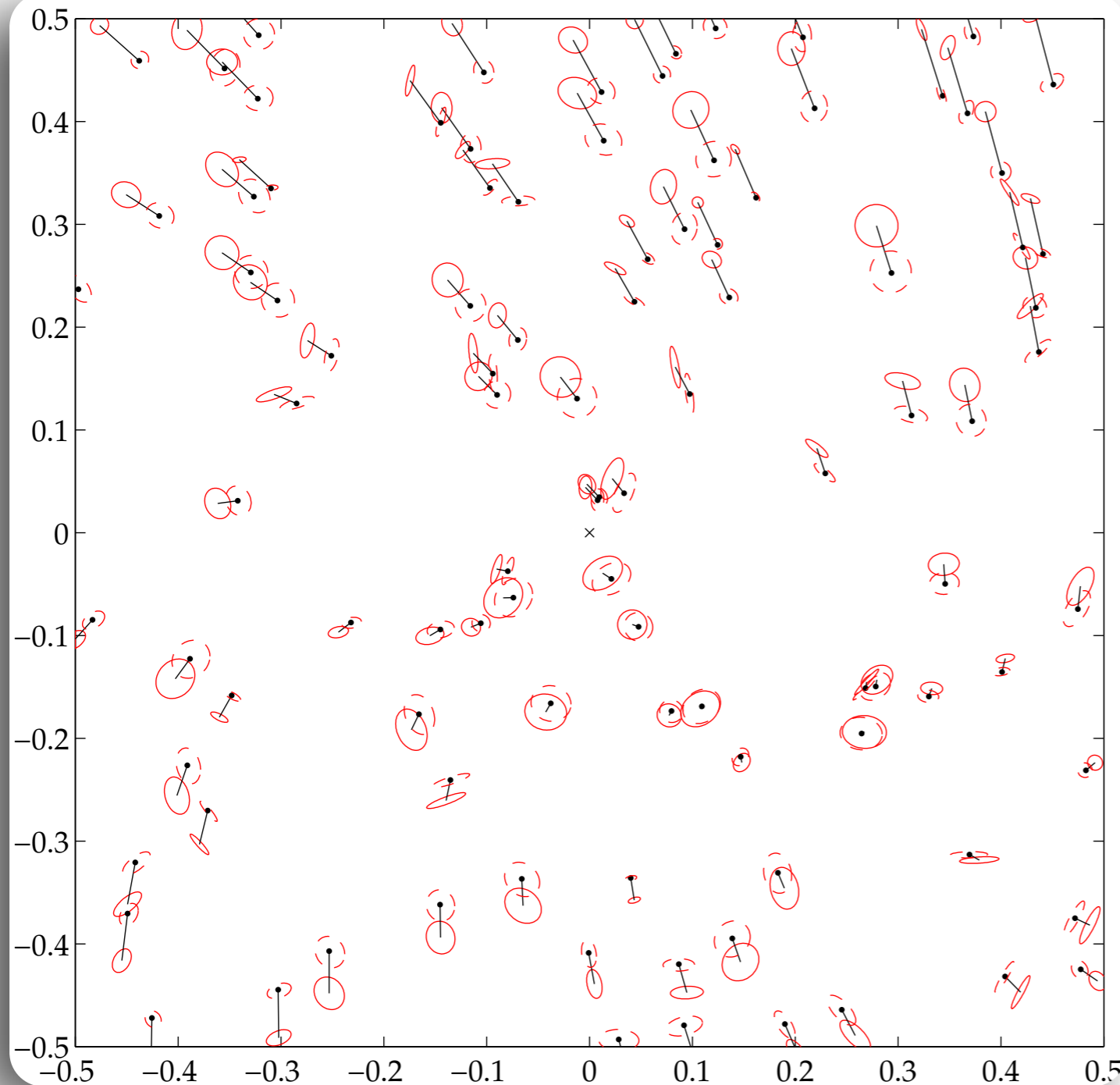


Fitting for:
Parallax, Proper
Motion, Binary
Parameters,
Microlensing...

Hipparcos: 10^6 stars
Gaia: $\sim 10^9$ stars

Bayesian Astrometry

Fork me on GitHub



Step 1:

Regress 7-d parametric affine transformation (scale, rotation, shear, etc.)

Step 2:

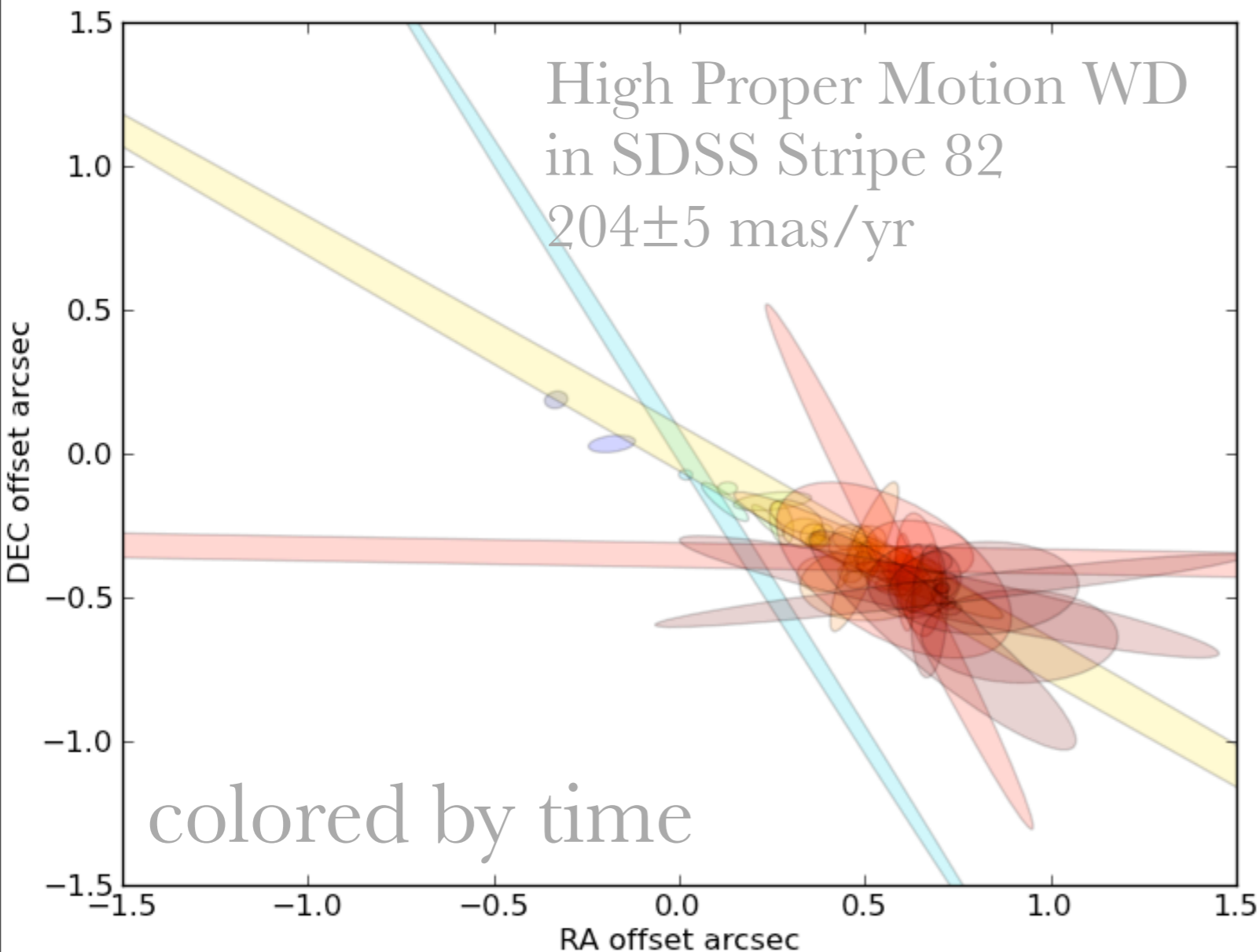
Learn a non-parametric distortion map with Gaussian processes

<http://berianjames.github.com/pyBAST/>

James, JSB+14

Bayesian Astrometry

Fork me on GitHub



Some Clear Benefits

- covariate uncertainties in celestial coordinates
- mapping observed points can incorporate variance throughout image, extending even to highly non-trivial distortion effects
- astrometry can be treated as *Bayesian updating*, allowing incorporation of prior knowledge about proper motion & parallax

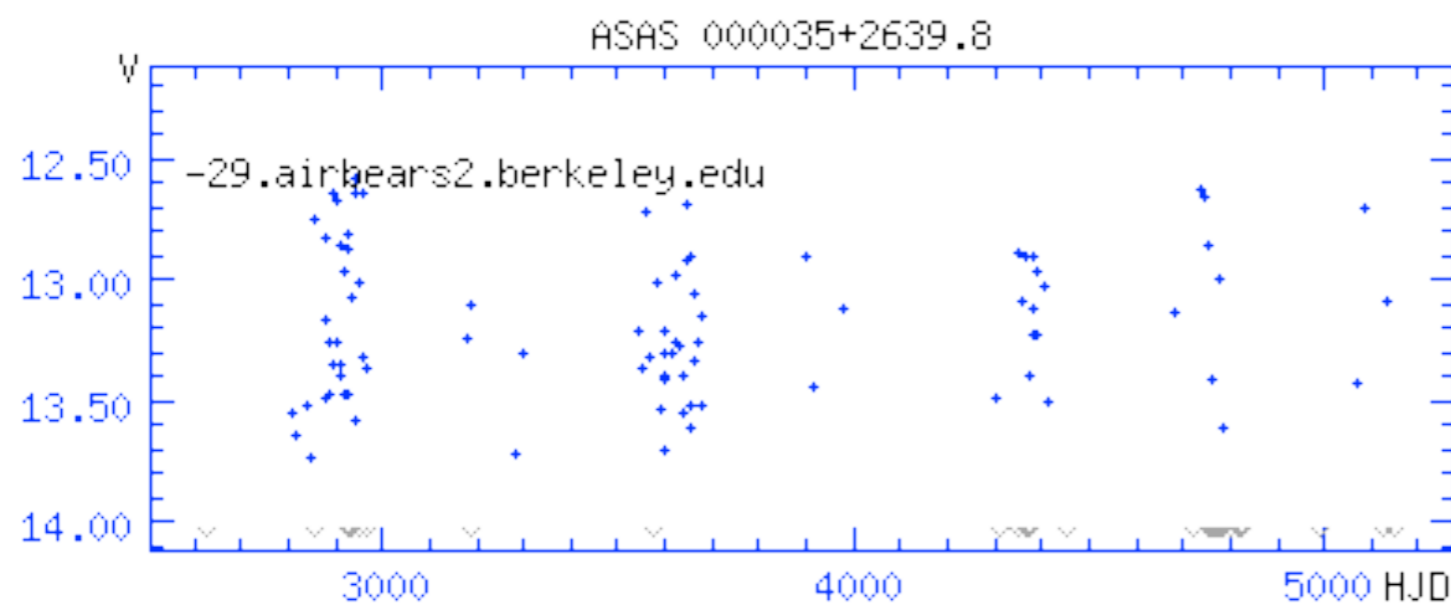
non-parallel Cholesky + MCMC: ~hour for 71 observations

<http://berianjames.github.com/pyBAST/>

Machine Learned Classification

25-class variable star

Data: 50k from ASAS, 810 with known labels
(timeseries, colors)

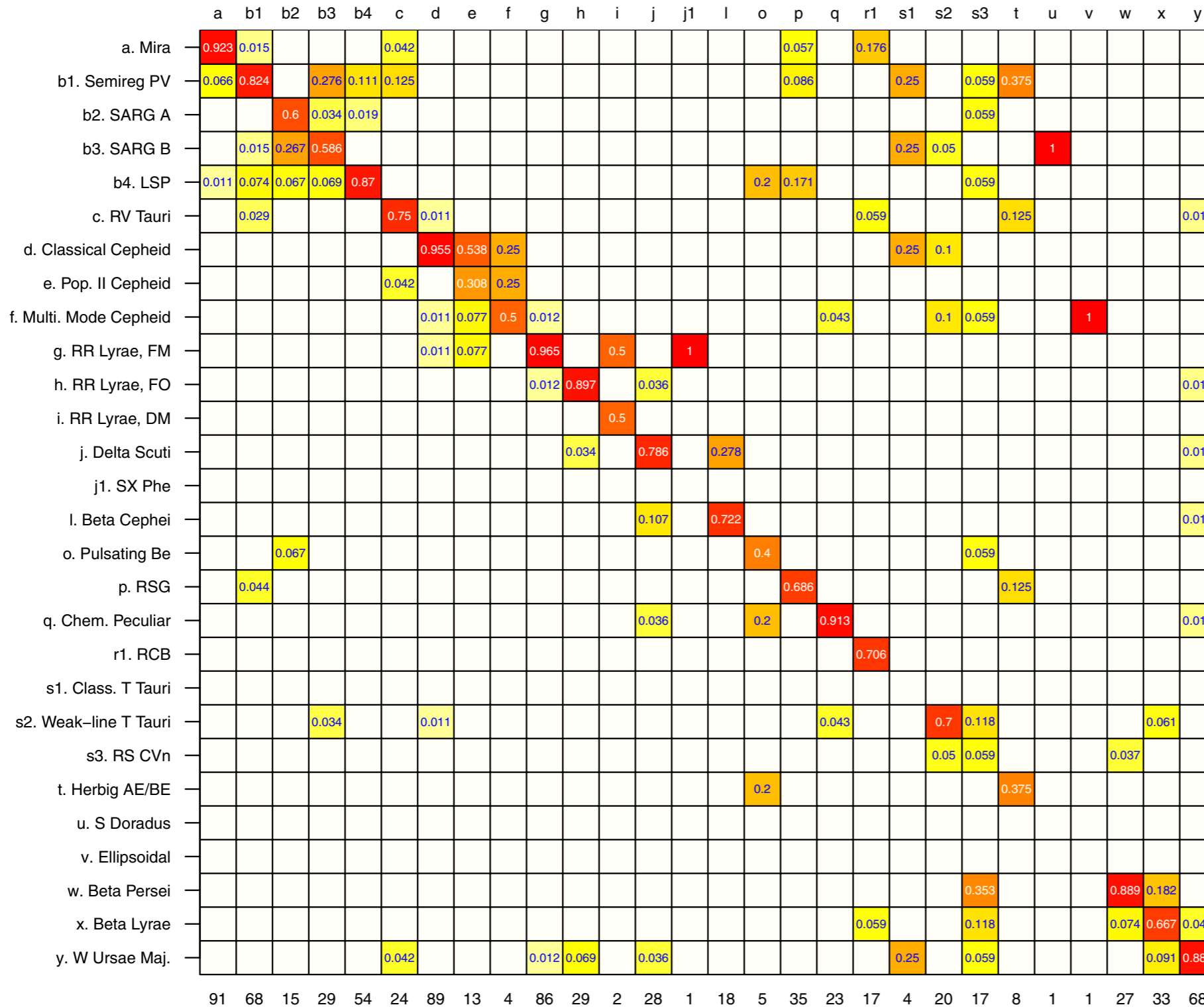


→ $P_{\text{RRL}} = 0.94$

Richards+12

Machine Learned Classification

True Class



74
dimensional
feature set for
learning

featurization is
the bottleneck
(but
embarrassingly
parallel)

Richards+12

Predicted Class

SELECT CLASS:

All

Specify RA/Dec

Customize Table

Export Query Results

Displaying sources 1 to 500

Previous 500

Next 500

MACC (50124)

Rotational (335)

Eruptive (2727)

Binary (11236)

Pulsating (35826)

ASAS_ID	dotAstro_ID	RA	DEC	Class	P_Class	Anomaly	ACVS_Class	Train_Class	P	P_signif	N_epochs	V	c
080940-3810.5	227867	122.415885	-38.174595	Mira	████████	0.040	MIRA=SR	Mira	328.386	20.859	456	7.8	
115501-5915.2	236087	178.762155	-59.258671	Mira	████████	0.075	MIRA		200.167	21.082	339	8.2	
132500-6439.8	238210	201.24495	-64.663232	Mira	████████	0.066	MIRA		350.507	23.531	503	10.34	
161441-3223.5	244080	243.671565	-32.391181	Mira	████████	0.037	MIRA		358.451	16.717	563	9.99	
165413-5615.9	245810	253.55541	-56.265033	Mira	████████	0.033	MIRA	Mira	286.697	16.654	497	9.46	
165538-4506.2	245884	253.907535	-45.102913	Mira	████████	0.040	MIRA		316.996	23.126	402	7.82	
194952+0923.8	258863	297.46893	9.401204	Mira	████████	0.042	MIRA		287.103	12.820	470	10.06	
221800-2936.2	263989	334.501365	-29.604124	Mira	████████	0.055	MIRA		293.459	16.501	415	9.17	
235627-4947.2	265240	359.121555	-49.786453	Mira	████████	0.066	MIRA		266.197	16.322	408	8.69	
044030-3814.2	218251	70.125405	-38.235209	Mira	████████	0.029	MIRA	Mira	390.997	20.138	382	8.98	
070729+0459.1	224085	106.87182	4.986432	Mira	████████	0.050	MIRA		262.462	13.700	480	10.02	
091646-0435.2	230899	139.19409	-4.585538	Mira	████████	0.050	MIRA		268.569	20.205	907	9.56	
094755-6726.9	232018	146.98116	-67.451082	Mira	████████	0.092	MIRA		338.995	18.072	290	10.62	
103823-8046.8	233741	159.6159	-80.7821	Mira	████████	0.082	MIRA		30.522	23.139	823	10.1	
120517-5511.2	236361	181.322	-55.1783	Mira	████████	0.030	MIRA		18.638	18.337	574	9.75	
121938-1915.3	236646	184.9078	-19.27056	Mira	████████	0.075	MIRA		17.856	18.677	615	7.83	



MACC ABOUT CONTACT FAQ



Machine-learned varstar catalog: <http://bigmacc.info>

Doing Science with Probabilistic Catalogs

Demographics (with little followup):

trading high purity at the cost of lower efficiency
e.g., using RRL to find new Galactic structure

Novelty Discovery (with lots of followup):

trading high efficiency for lower purity
e.g., discovering new instances of rare classes

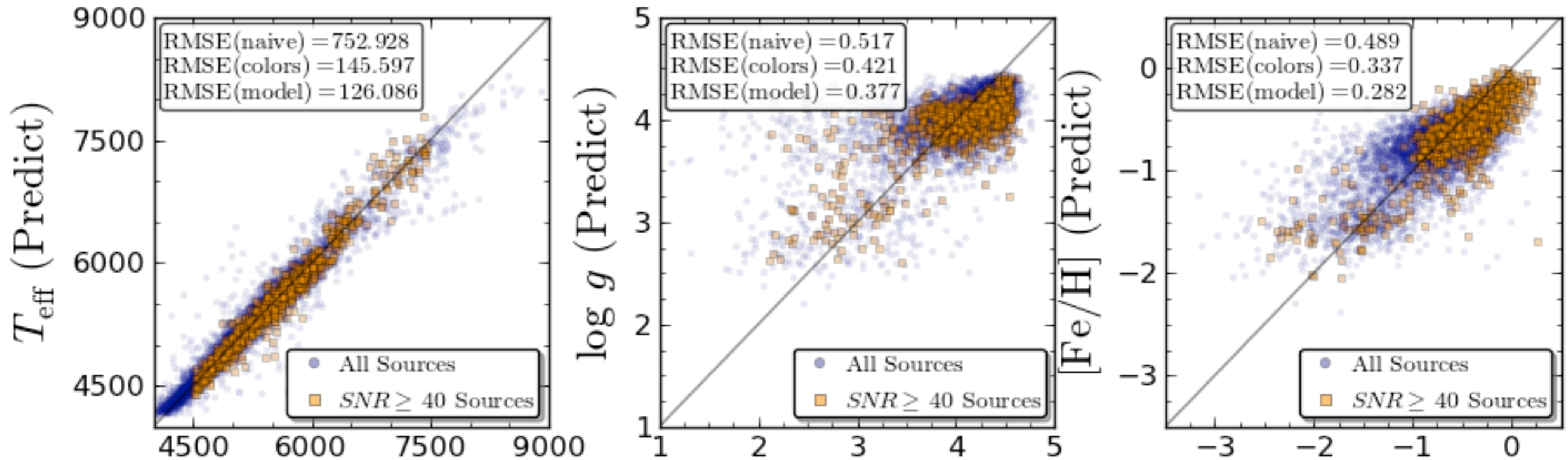
**Discovery of Bright Galactic R Coronae Borealis and DY Persei
Variables: Rare Gems Mined from ASAS**

A. A. Miller^{1,*}, J. W. Richards^{1,2}, J. S. Bloom¹, S. B. Cenko¹, J. M. Silverman¹,
D. L. Starr¹, and K. G. Stassun^{3,4}

[arXiv.org > astro-ph > arXiv:1204.4181](https://arxiv.org/astro-ph/1204.4181)

Turning Imagers into Spectrographs

Time variability + colors \rightarrow fundamental stellar parameters



Data: 5000 variables in SDSS Stripe 82 with spectra
 \sim 80 dimensional regression with Random Forest

Miller, JSB+14

Big Data Challenge: Time & Resources

Large Synoptic Survey Telescope (LSST) - 2018

Light curves for 800M sources every 3 days
 10^6 supernovae/yr, 10^5 eclipsing binaries
3.2 gigapixel camera, 20 TB/night

LOFAR & SKA

150 Gps (27 Tflops) → 20 Pps (~100 Pflops)

Gaia space astrometry mission - 2013

1 billion stars observed ~70 times over 5 years
Will observe 20K supernovae

Many other astronomical surveys are already producing data:
SDSS, iPTF, CRTS, Pan-STARRS, Hipparcos, OGLE, ASAS,
Kepler, LINEAR, DES etc.,

Big Data Challenge: Time & Resources

Large Synoptic Survey Telescope (LSST) - 2018

Light curves for 800M sources every 3 days

10^6 supernovae/yr, 10^5 eclipsing binaries

3.2 gigapixel camera, 20 TB/night

**How do we do discovery,
follow-up, and inference
when the data rates (&
requisite timescales)
preclude human
involvement?**

LOFAR & SKA

150 Gps (100 PB/yr)

Gaia space

1 billion stars of 1000 colors

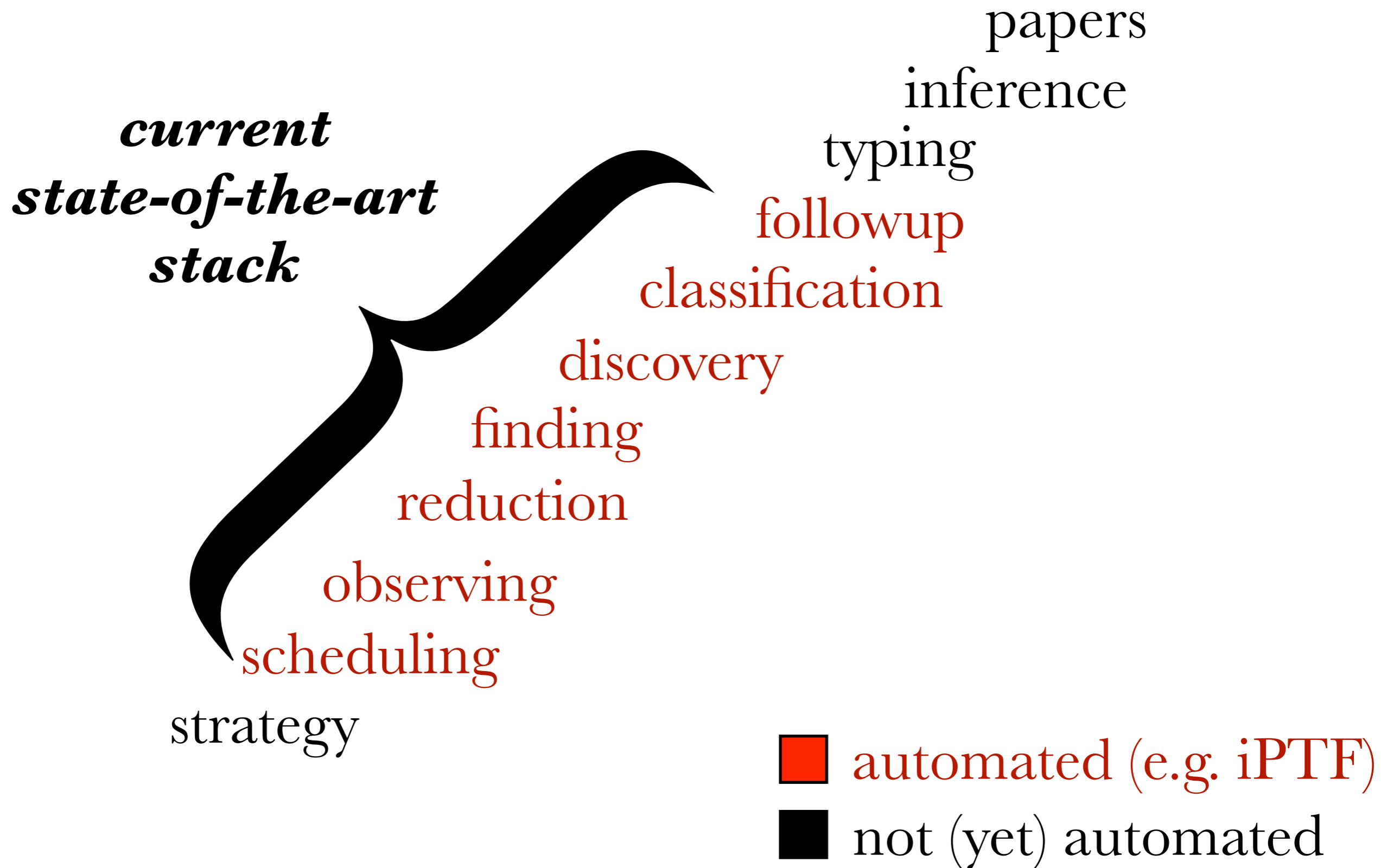
Will observe 20K supernovae

Many other astronomical surveys are already producing data:

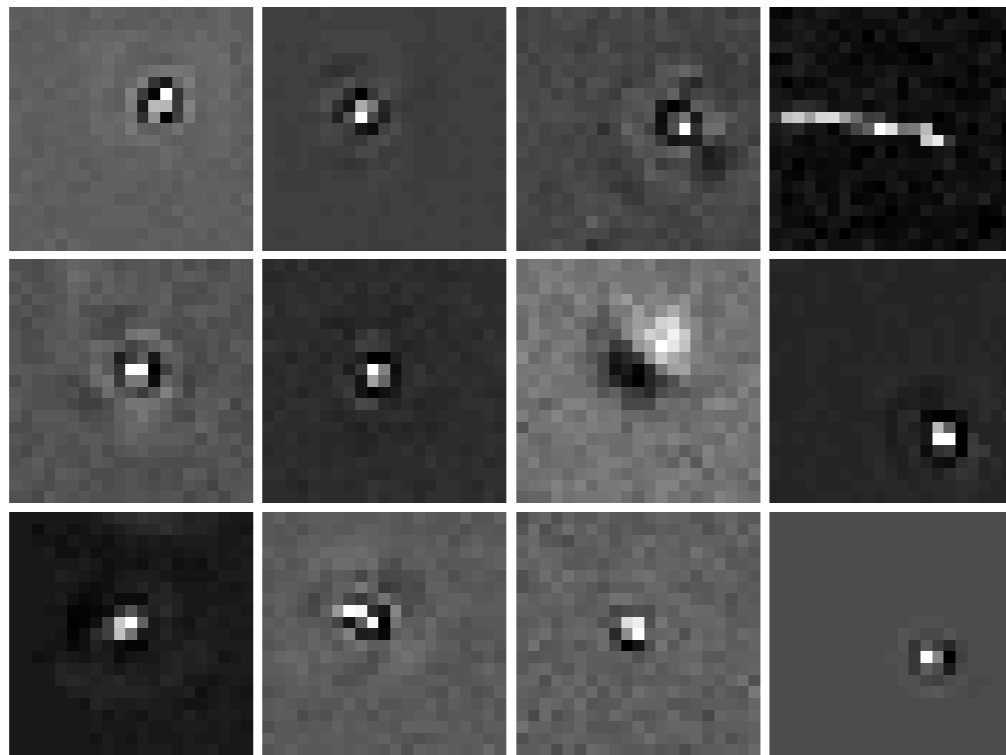
SDSS, IPTF, CRTS, Pan-STARRS, Hipparcos, OGLE, ASAS,

Kepler, LINEAR, DES etc.,

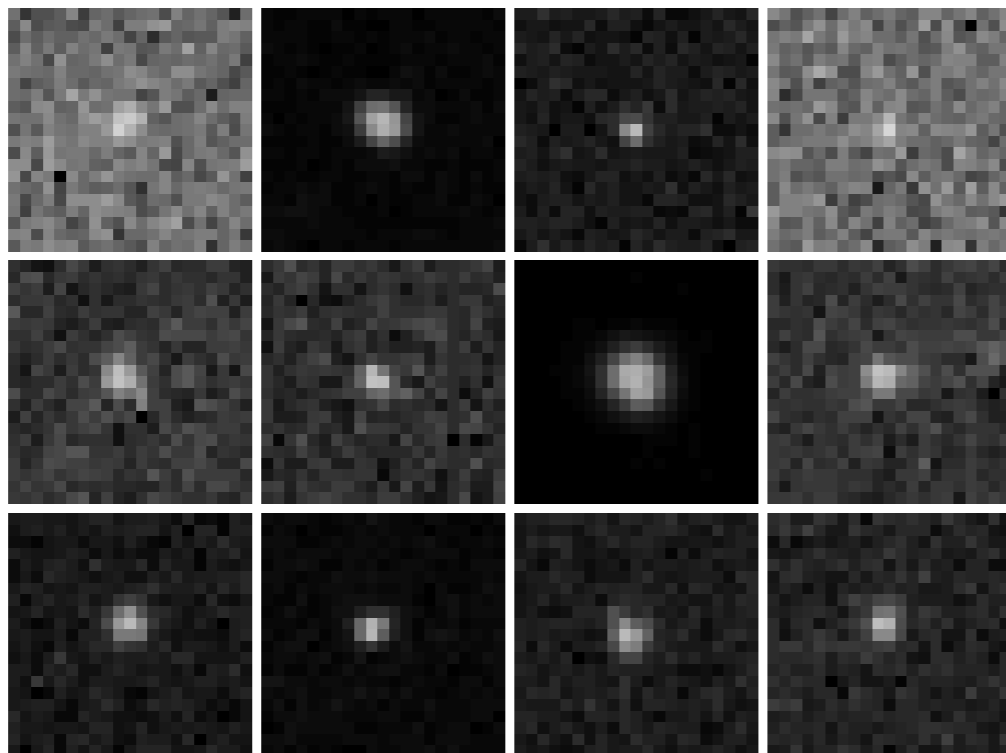
Towards a Fully Automated Scientific Stack for Transients



“bogus”



“real”



PTF subtractions

Goal:

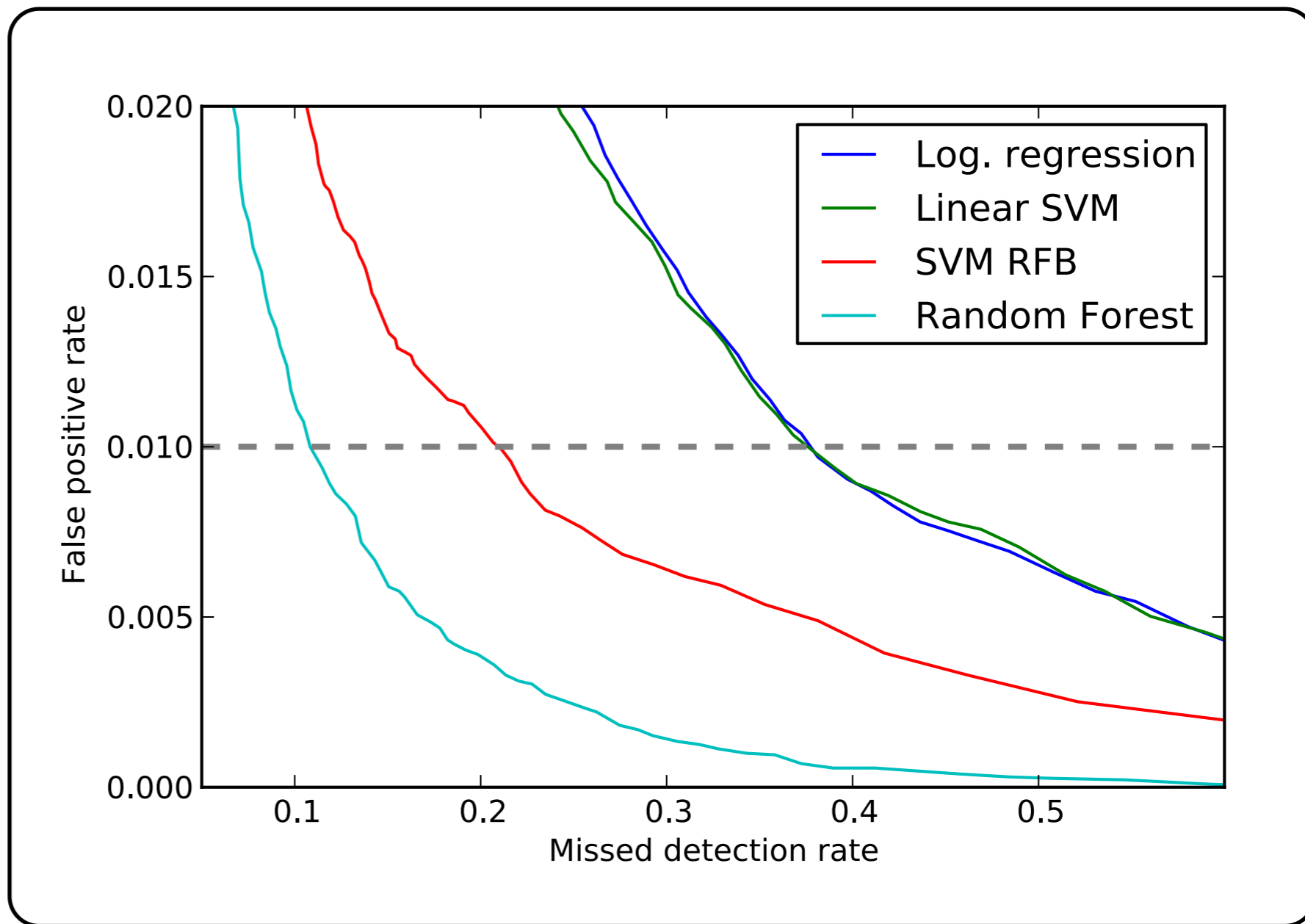
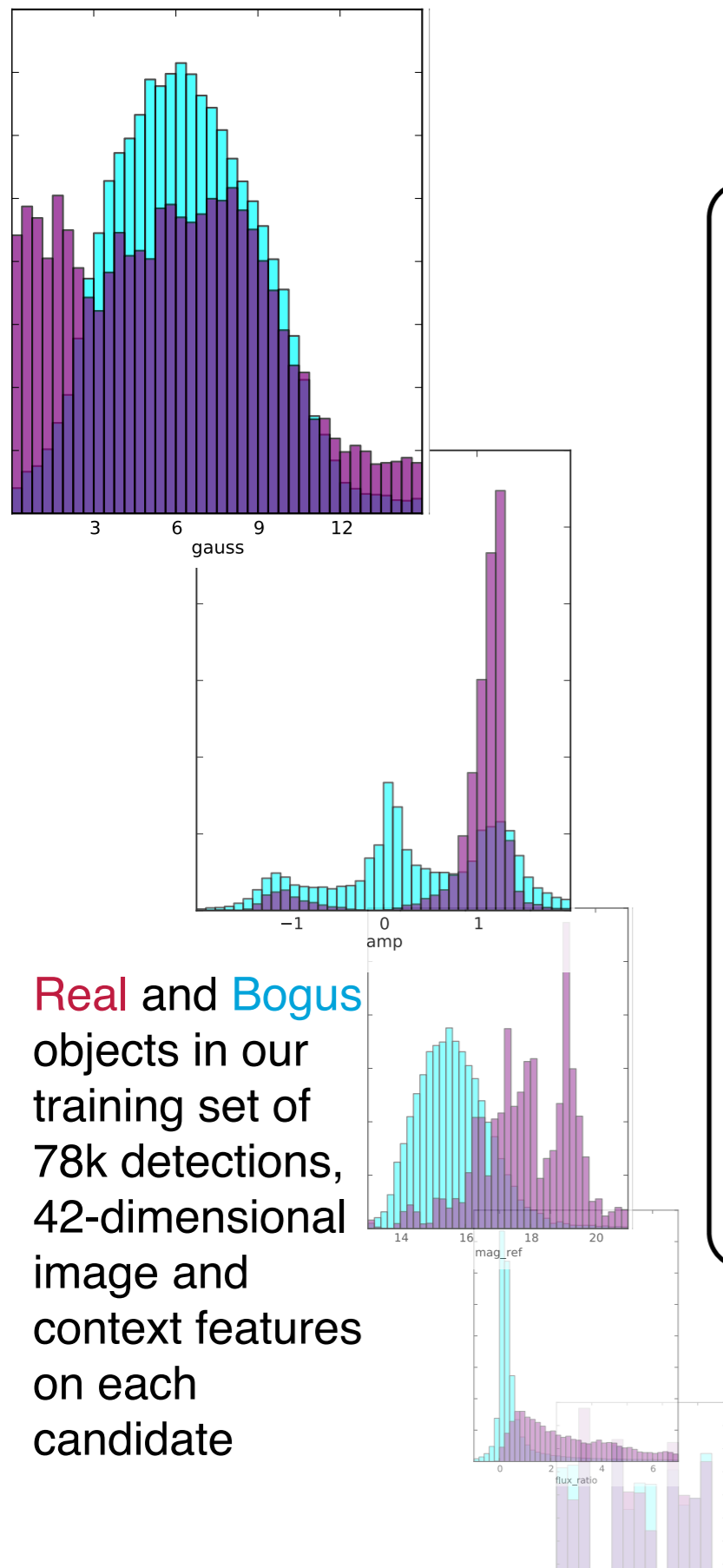
build a framework to
discover variable/
transient sources
without people

- fast (compared to people)
- parallelizable
- transparent
- deterministic
- versionable

*1000 to 1 needle in the
haystack problem*

“Discovery” is Imperfect

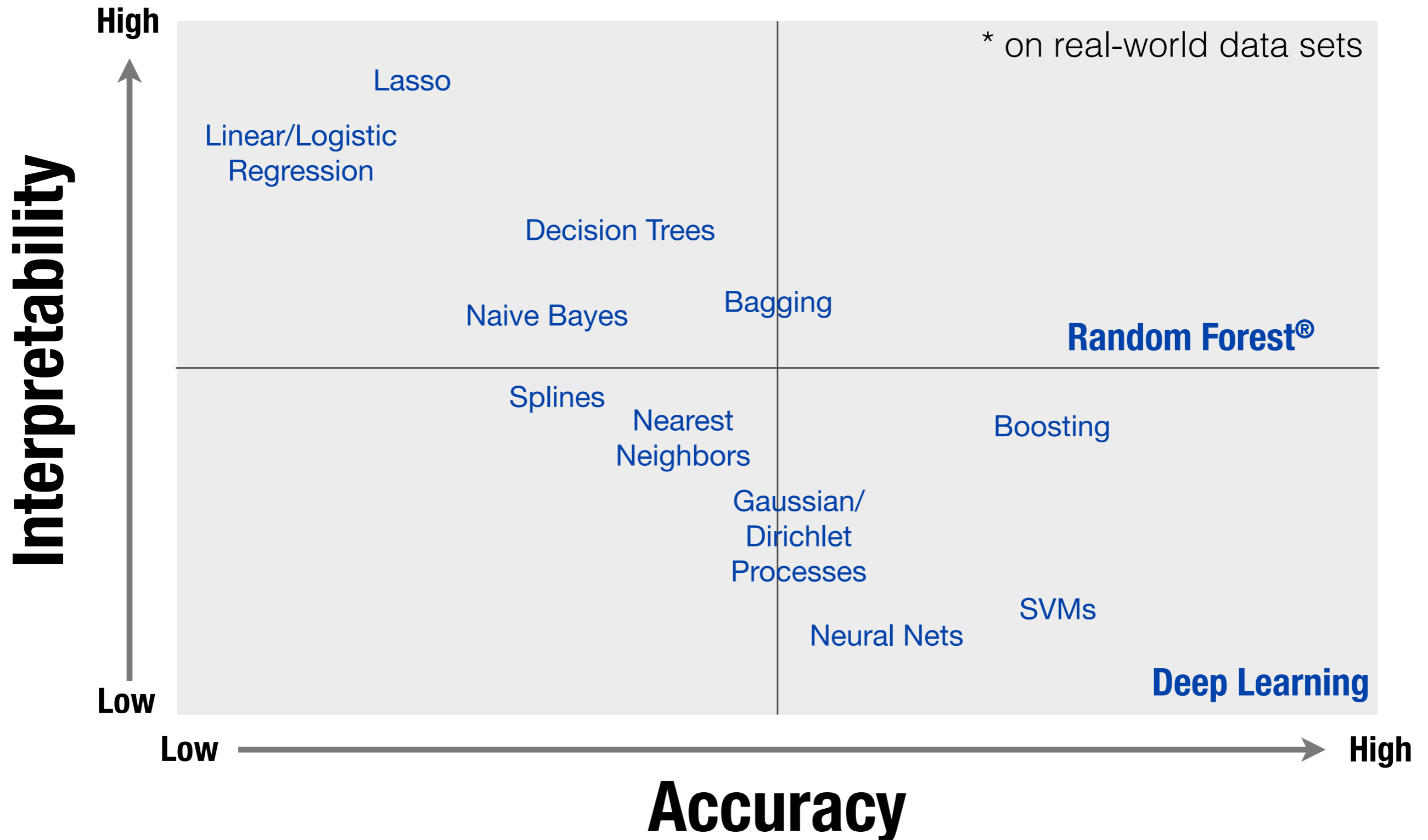
but some classifiers work better than others



Brink+2012

[arXiv.org > astro-ph > arXiv:1209.3775](https://arxiv.org/astro-ph/1209.3775)

ML Algorithmic Trade-Off



Random Forest is a trademark of Salford Systems, Inc.



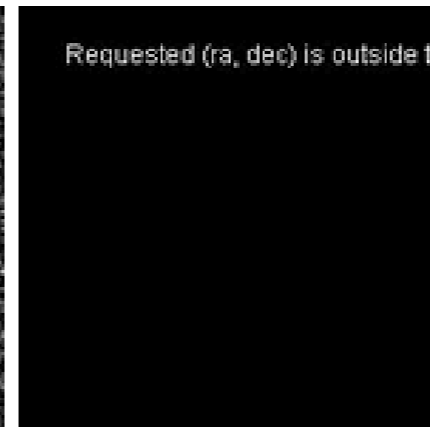
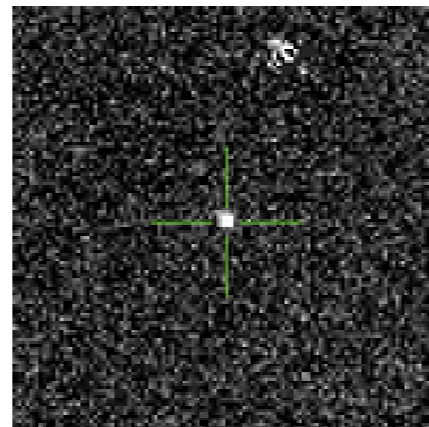
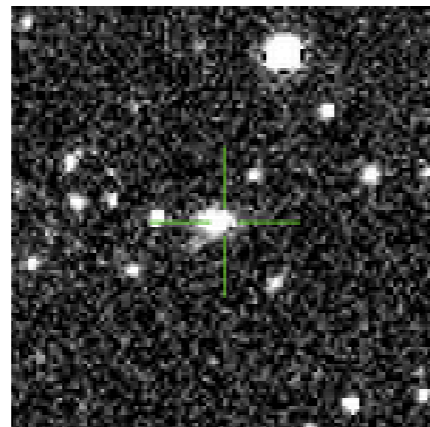
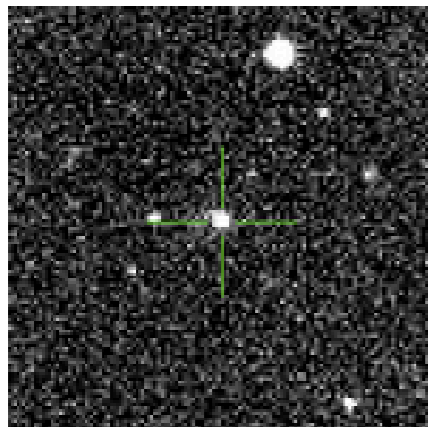
OVERVIEW PHOTOMETRY SPECTROSCOPY FOLLOWUP OBSERVABILITY FINDING CHART SCANNING

NEW

REF

SUB

SDSS



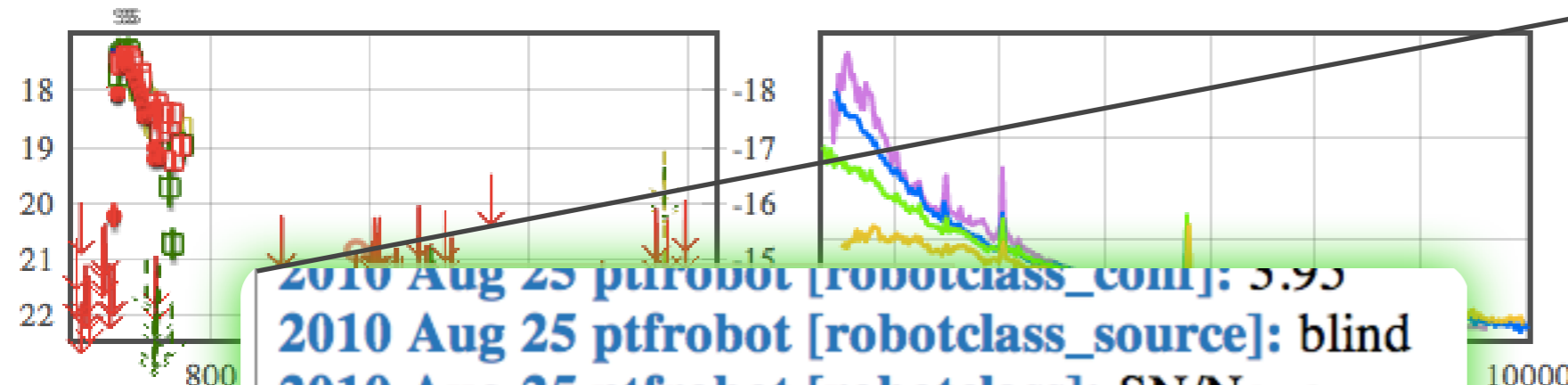
COMMENTS

- 2010 Oct 16 ptfrobot [robotclass_conf]: 3.94
- 2010 Sep 28 ptfrobot [robotclass_conf]: 3.96
- 2010 Sep 21 ptfrobot [robotclass_conf]: 3.97
- 2010 Sep 18 ptfrobot [robotclass_conf]: 4.00
- 2010 Aug 26 ahowell [info]: Redshift. Notice emiss Na ID. [view attachment]
- 2010 Aug 26 ahowell [redshift]: 0.035
- 2010 Aug 26 ahowell [classification]: SN II_n
- 2010 Aug 26 ahowell [comment]: A type II_n. Reduced Daniel Murray's pyraf script! [view attachment]
- 2010 Aug 25 ptfrobot [robotclass_conf]: 3.95
- 2010 Aug 25 ptfrobot [robotclass_source]: blind
- 2010 Aug 25 ptfrobot [robotclass]: SN/Nova
- 2010 Aug 25 ptfrobot [type]: Transient

Mentioned in 7 Email(s)

Add a Comment:

Attach File: No file chosen



2010 Aug 25 ptfrobot [robotclass_conf]: 3.95
 2010 Aug 25 ptfrobot [robotclass_source]: blind
 2010 Aug 25 ptfrobot [robotclass]: SN/Nova
 2010 Aug 25 ptfrobot [type]: Transient

Mentioned in 7 Email(s)

ADDITIONAL IN

- IPAC
- NED
- DSS
- SIMBAD
- HEASARC
- WISE
- SkyView
- PyMP
- PTFInfoBot
- Extinction

Add to Cart



OVERVIEW PHOTOMETRY SPECTROSCOPY FOLLOWUP OBSERVABILITY FINDING CHART SCANNING

NEW

REF

SUB

SDSS

COMMENTS

LETTER

doi:10.1038/nature11877

An outburst from a massive star 40 days before a supernova explosion

E. O. Ofek¹, M. Sullivan^{2,3}, S. B. Cenko⁴, M. M. Kasliwal⁵, A. Gal-Yam¹, S. R. Kulkarni⁶, I. Arcavi¹, L. Bildsten^{7,8}, J. S. Bloom^{4,9}, A. Horesh⁶, D. A. Howell^{8,10}, A. V. Filippenko⁴, R. Laher¹¹, D. Murray¹², E. Nakar¹³, P. E. Nugent^{4,9}, J. M. Silverman^{4,14}, N. J. Shaviv¹⁵, J. Surace¹¹ & O. Yaron¹

2010 Aug 25 ptfrobot [robotclass_source]: Blind

2010 Aug 25 ptfrobot [robotclass]: SN/Nova

2010 Aug 25 ptfrobot [type]: Transient

DM (approximate) = 35.88

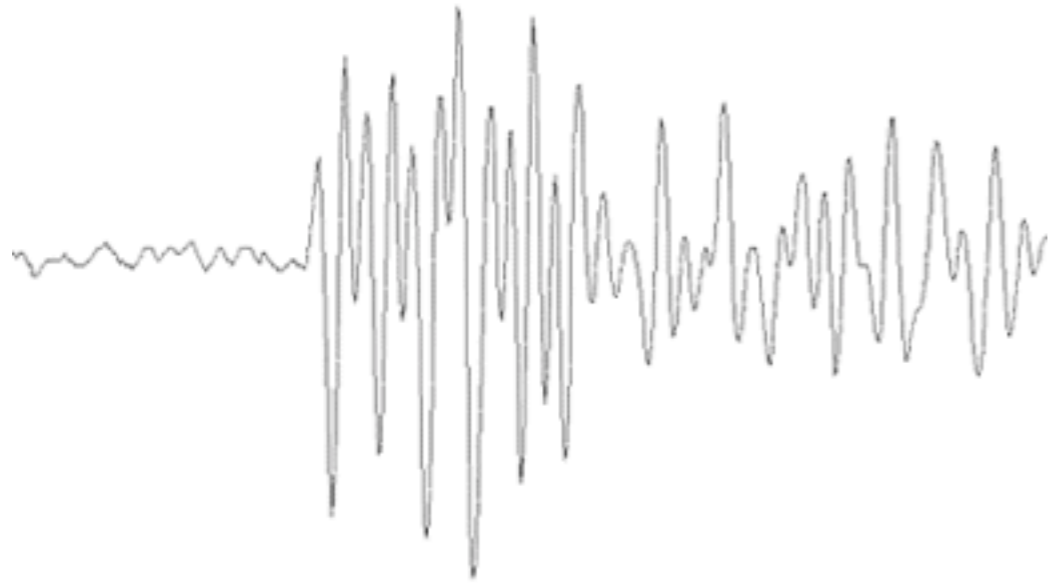
Mentioned in 7 Email(s)

ADDITIONAL INFO

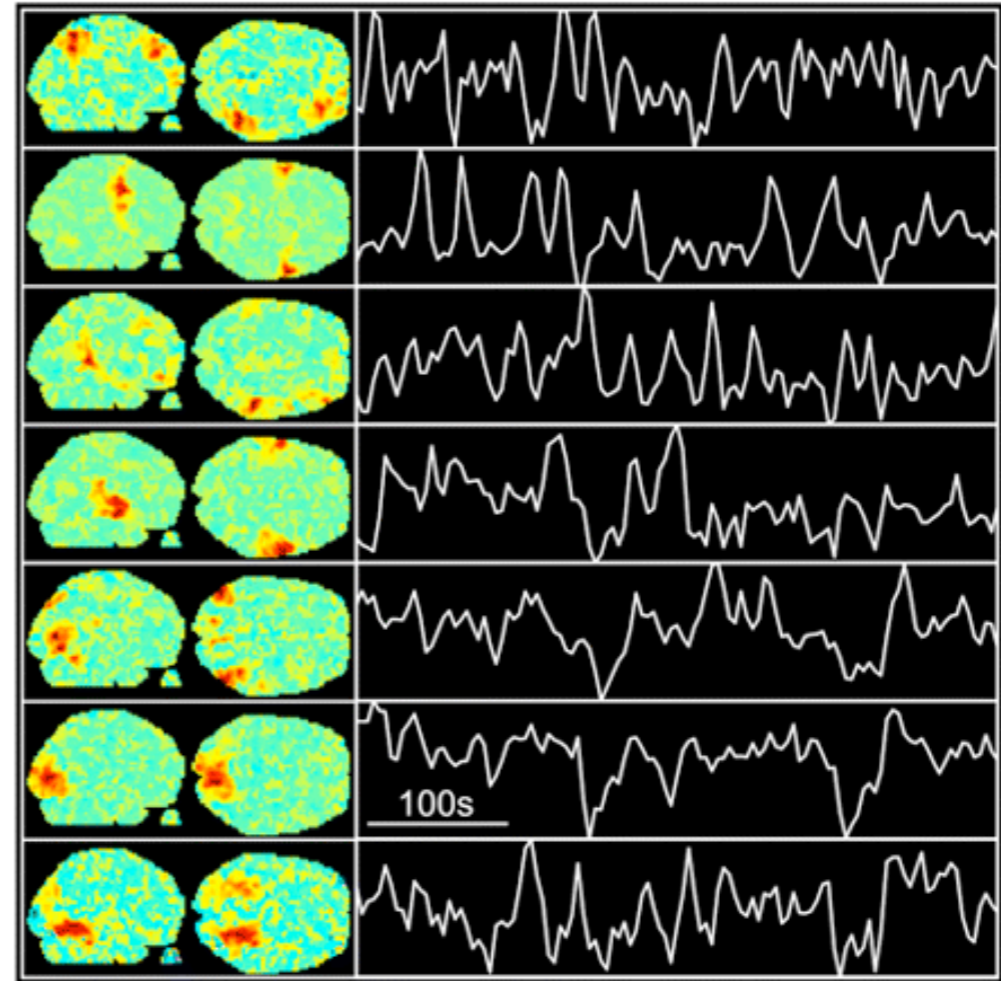
- IPAC
- NED
- DSS
- SIMBAD
- HEASARC
- WISE
- SkyView
- PyMP
- PTFInfoBot
- Extinction

Attach File: No file chosen

Add to Cart



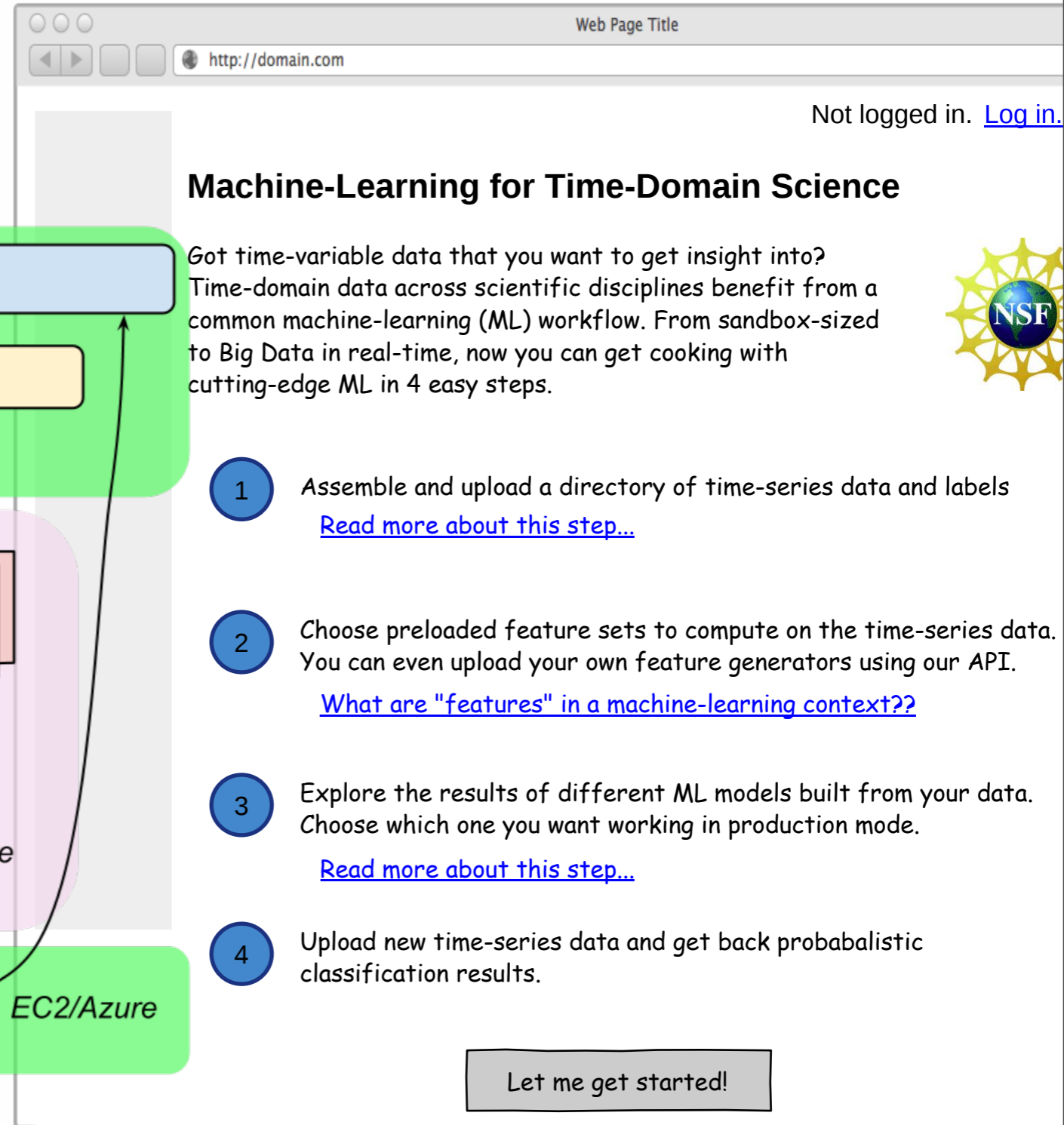
Seismology



Neuroscience

Classification Platform for Novel Scientific Insight on Time-Series Data

NSF/BIGDATA project




Web Page Title

http://domain.com

Not logged in. [Log in.](#)

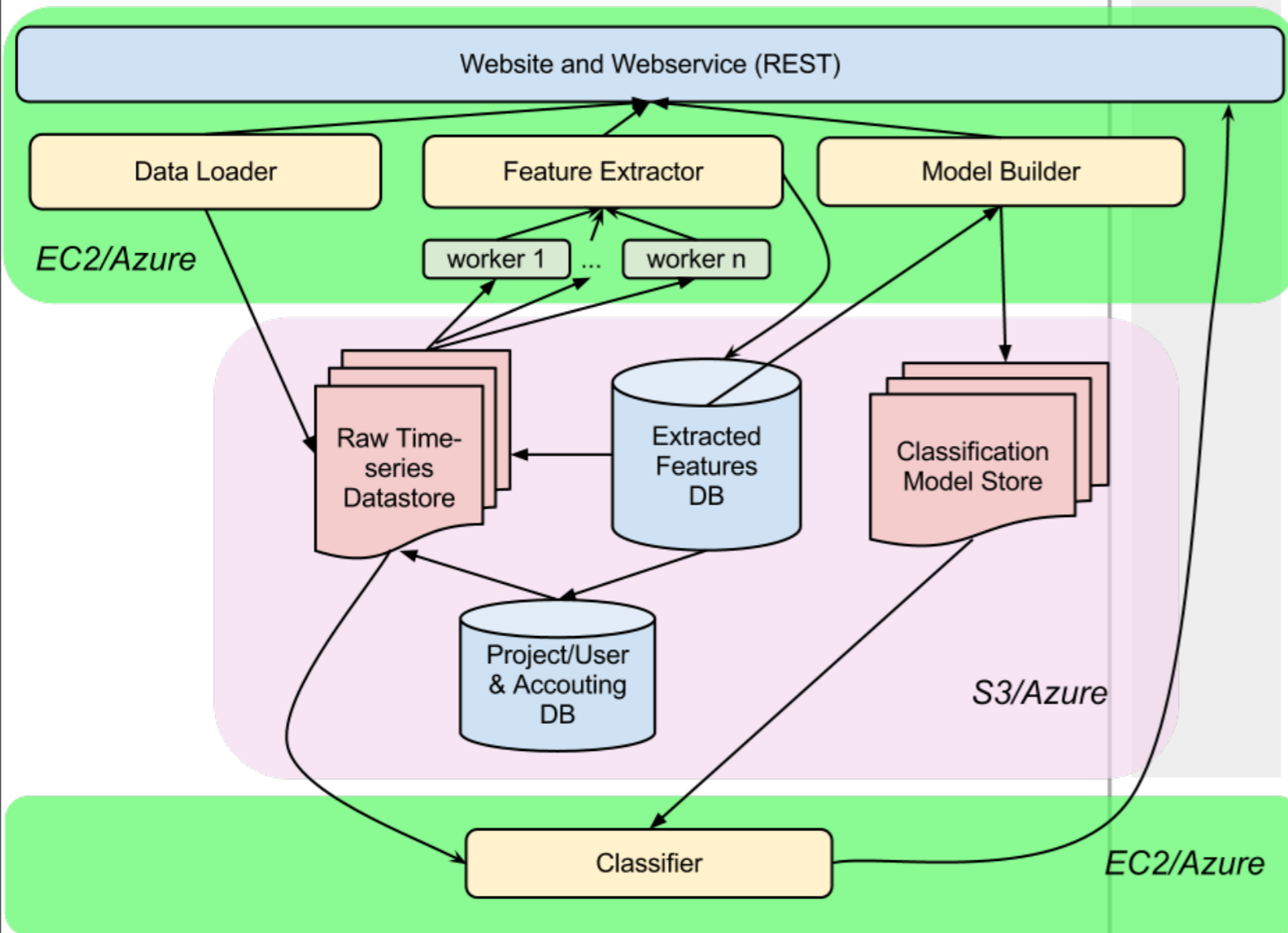
Machine-Learning for Time-Domain Science

Got time-variable data that you want to get insight into? Time-domain data across scientific disciplines benefit from a common machine-learning (ML) workflow. From sandbox-sized to Big Data in real-time, now you can get cooking with cutting-edge ML in 4 easy steps.



- 1 Assemble and upload a directory of time-series data and labels
[Read more about this step...](#)
- 2 Choose preloaded feature sets to compute on the time-series data. You can even upload your own feature generators using our API.
[What are "features" in a machine-learning context??](#)
- 3 Explore the results of different ML models built from your data. Choose which one you want working in production mode.
[Read more about this step...](#)
- 4 Upload new time-series data and get back probabilistic classification results.

Let me get started!



Parting Thoughts

- Astronomy's data deluge demands an abstraction of the traditional roles in the scientific process. Despite automation, crucial (insightful) roles remain for people
- Machine learning is an emerging & useful tool
- Deterministic prediction with verifiable uncertainties is crucial to maximize scientific impact under real-world resource constraints
- Major Challenge: Training & access to great learning frameworks
- In the time-domain, machine-learning ***prediction is the gateway to understanding***