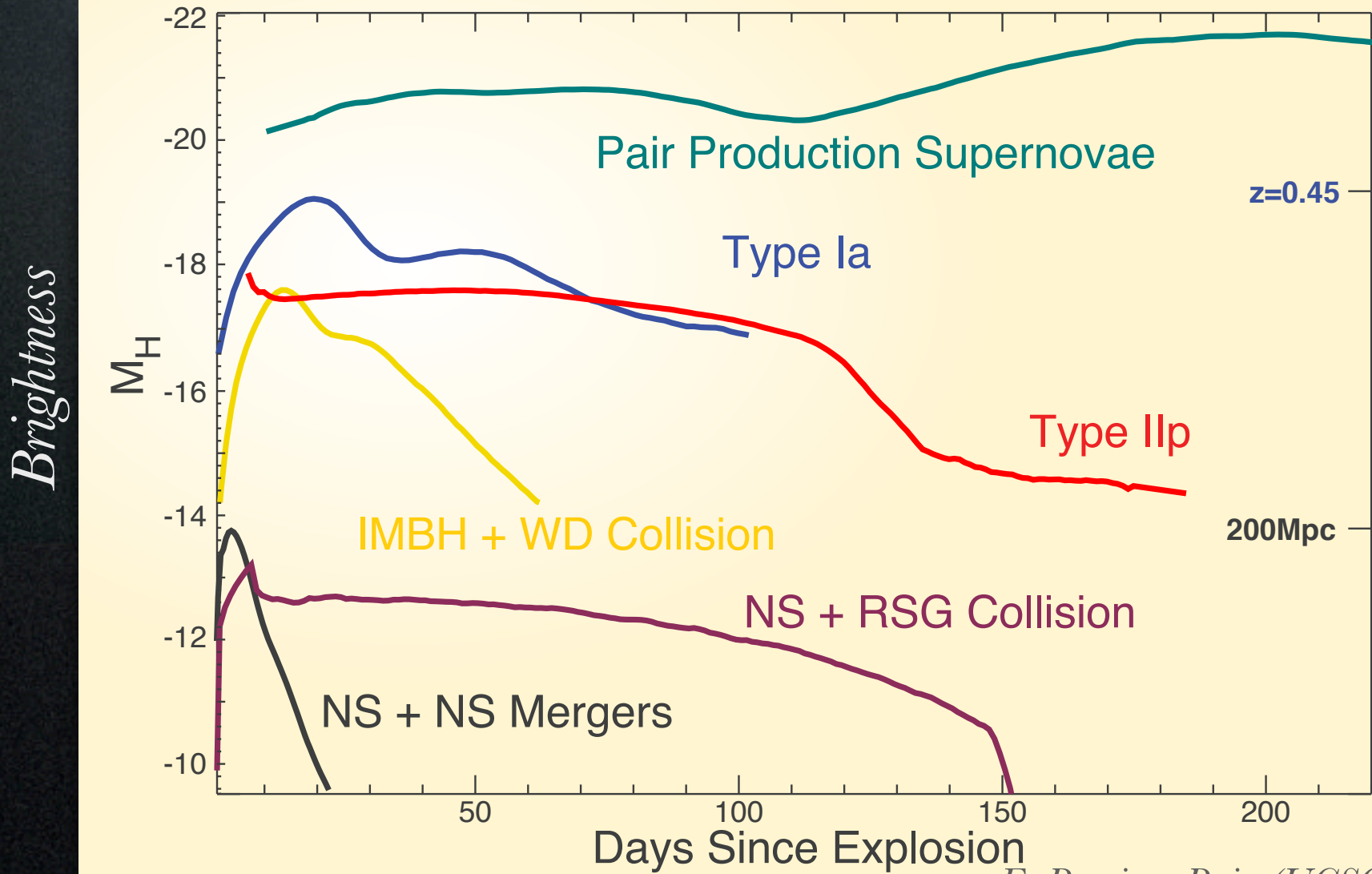January 7 - March 2, 1610

# Real-Time Knowledge Extraction from Massive Time-Series Datastreams

**Josh Bloom**

*Astronomy Department*

jbloom@astro.berkeley.edu

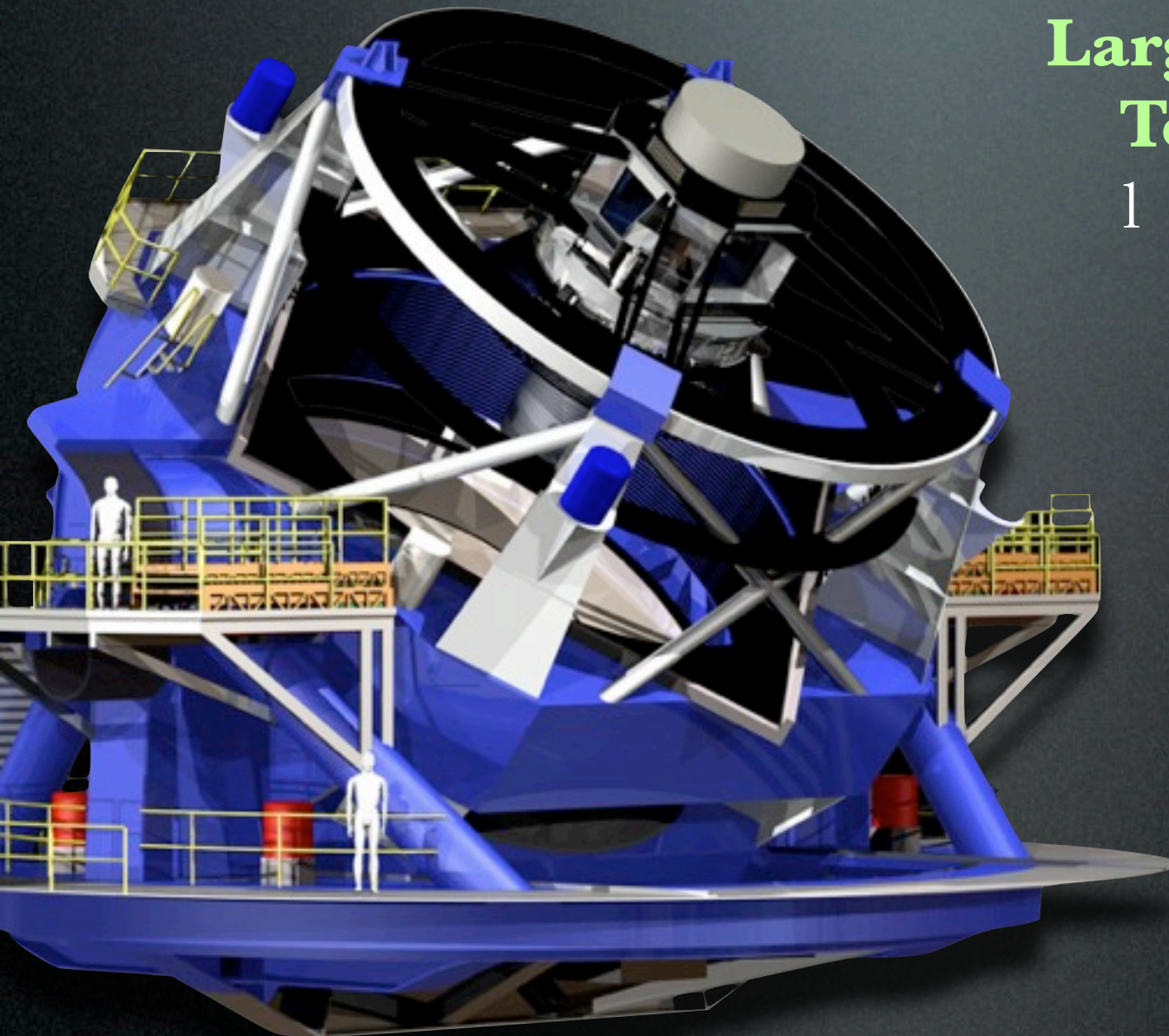# Extragalactic Transient Universe: Explosive Systems

Brightness — $M_H$ vs Days Since Explosion

- Pair Production Supernovae
- Type Ia
- Type IIp
- IMBH + WD Collision
- NS + RSG Collision
- NS + NS Mergers

z=0.45

200Mpc

E. Ramirez-Ruiz (UCSC)

**"Bad" News:**
**Discoveries Swamp Followup Resources**

**Large Synoptic Survey Telescope (LSST):**
1 Gb every 2 seconds

*$10^6$ supernovae/yr*
*$10^5$ eclipsing systems*
*$10^7$ asteroids...*

light curves of 800 million sources every 3 days

# *Transients Classification Project*

***Berkeley Astronomy*:**
      Dan Starr, Dovi Poznanski, Maxime Rischard, Nat Butler,
      Chris Klein, Rachel Kennedy, Justin Huggins, Adam
      Morgan, Adam Miller, JSB
***San Francisco State University*:**
      John M. Brewer
***Berkeley Statistics*:**
      Noureddine El Karoui, John Rice
***Berkeley CS*:**
      Martin Wainwright, Masoud Nikravesh
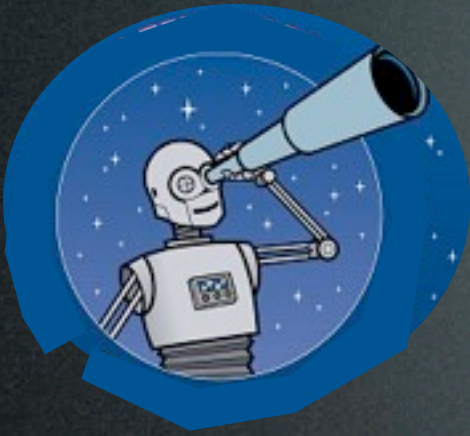***Lawrence Berkeley Lab*:**
      Peter Nugent, Horst Simon
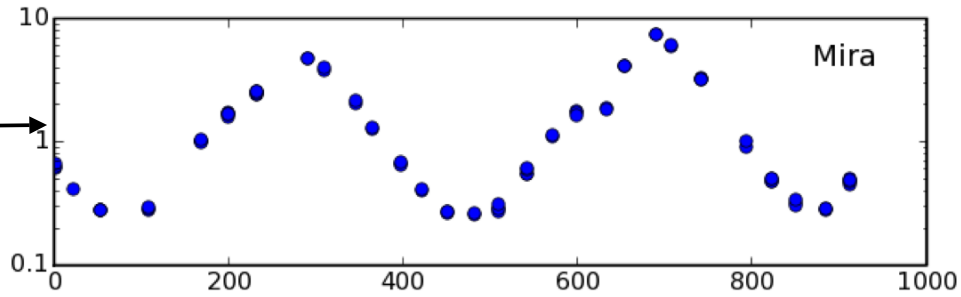***Los Alamos Nat. Lab. / UC Santa Cruz*:**
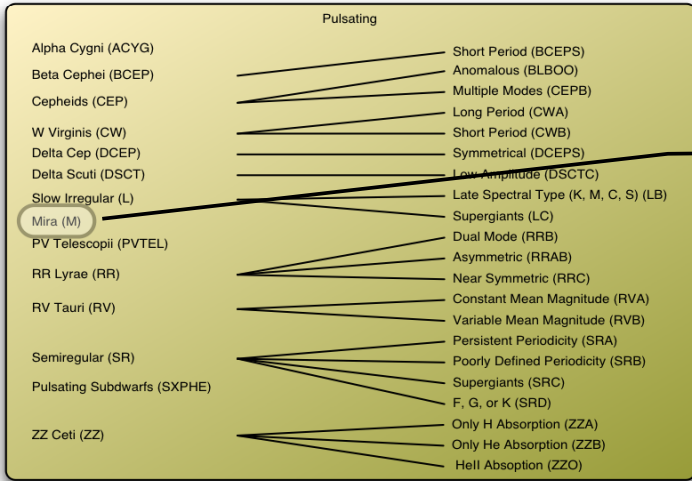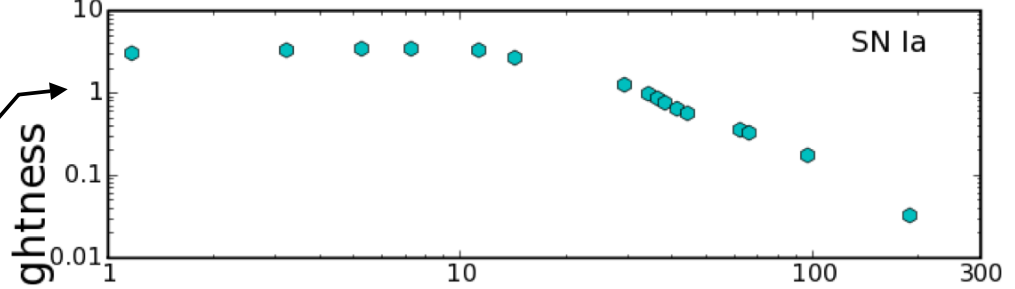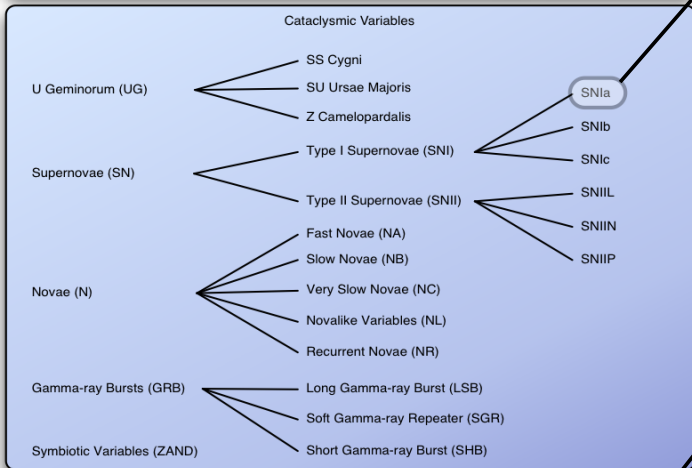      Damian Eads

**Goal: Autonomous creation of new knowledge, that itself spurs further resource allocation & inquiry**

- Generate **probabilistic statements** about the nature of events (ie. classification)

- Provide push/pull **access** to current & past events

- (bootstrap) Learning from feedback

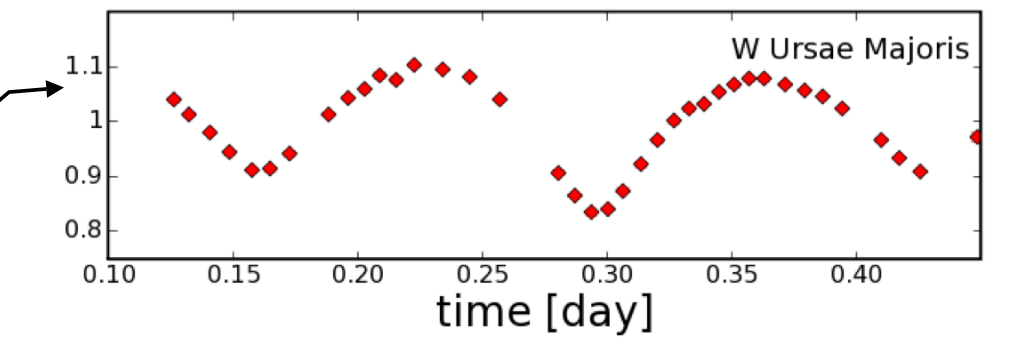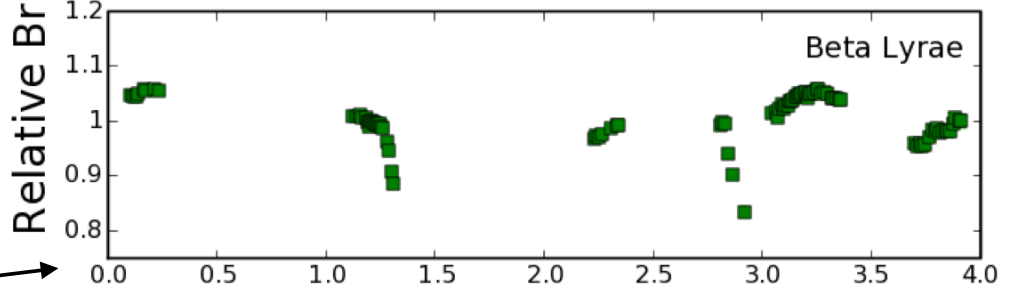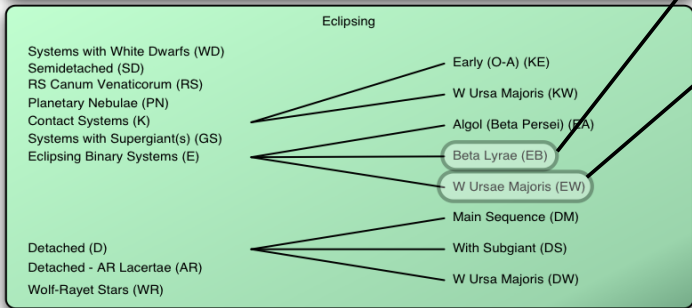- Operate at sufficient & **scalable** rates

**Pulsating Stars**

Pulsating

Alpha Cygni (ACYG)
Beta Cephei (BCEP)
Cepheids (CEP)
W Virginis (CW)
Delta Cep (DCEP)
Delta Scuti (DSCT)
Slow Irregular (L)
Mira (M)
PV Telescopii (PVTEL)
RR Lyrae (RR)
RV Tauri (RV)
Semiregular (SR)
Pulsating Subdwarfs (SXPHE)
ZZ Ceti (ZZ)

Short Period (BCEPS)
Anomalous (BLBOO)
Multiple Modes (CEPB)
Long Period (CWA)
Short Period (CWB)
Symmetrical (DCEPS)
Low Amplitude (DSCTC)
Late Spectral Type (K, M, C, S) (LB)
Supergiants (LC)
Dual Mode (RRB)
Asymmetric (RRAB)
Near Symmetric (RRC)
Constant Mean Magnitude (RVA)
Variable Mean Magnitude (RVB)
Persistent Periodicity (SRA)
Poorly Defined Periodicity (SRB)
Supergiants (SRC)
F, G, or K (SRD)
Only H Absorption (ZZA)
Only He Absorption (ZZB)
HeII Absoption (ZZO)

**Cataclysmic Variables**

Cataclysmic Variables

U Geminorum (UG)
Supernovae (SN)
Novae (N)
Gamma-ray Bursts (GRB)
Symbiotic Variables (ZAND)

SS Cygni
SU Ursae Majoris
Z Camelopardalis
Type I Supernovae (SNI)
Type II Supernovae (SNII)
Fast Novae (NA)
Slow Novae (NB)
Very Slow Novae (NC)
Novalike Variables (NL)
Recurrent Novae (NR)
Long Gamma-ray Burst (LSB)
Soft Gamma-ray Repeater (SGR)
Short Gamma-ray Burst (SHB)

SNIa
SNIb
SNIc
SNIIL
SNIIN
SNIIP

**Eclipsing Systems**

Eclipsing

Systems with White Dwarfs (WD)
Semidetached (SD)
RS Canum Venaticorum (RS)
Planetary Nebulae (PN)
Contact Systems (K)
Systems with Supergiant(s) (GS)
Eclipsing Binary Systems (E)
Detached (D)
Detached - AR Lacertae (AR)
Wolf-Rayet Stars (WR)

Early (O-A) (KE)
W Ursa Majoris (KW)
Algol (Beta Persei) (EA)
Beta Lyrae (EB)
W Ursae Majoris (EW)
Main Sequence (DM)
With Subgiant (DS)
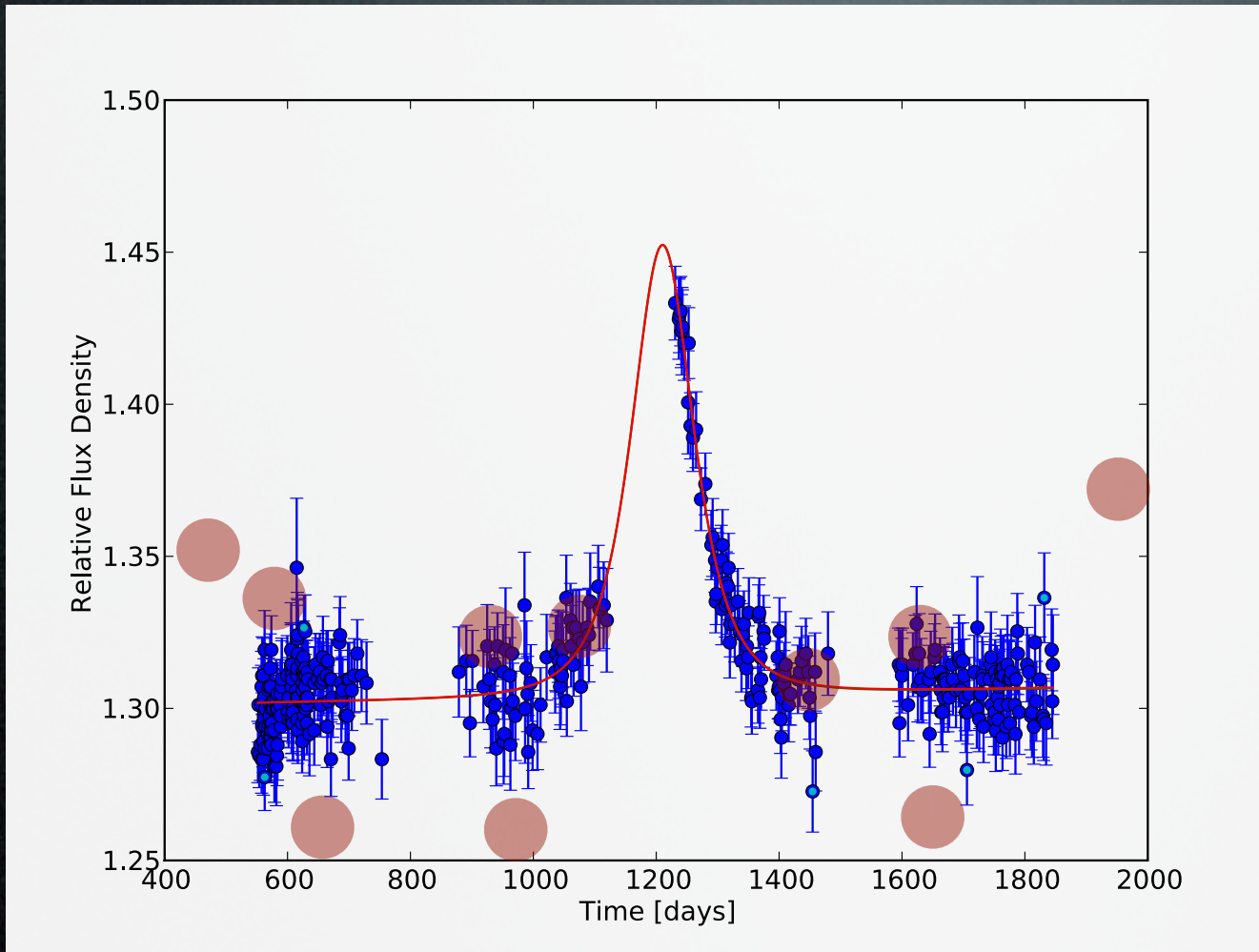W Ursa Majoris (DW)

Mira

SN Ia

Beta Lyrae

W Ursae Majoris

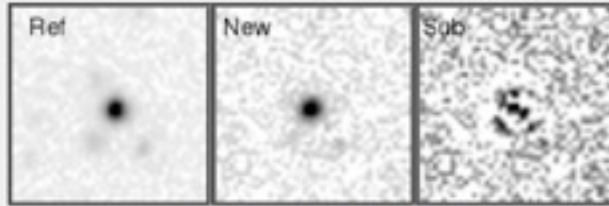Relative Brightness

time [day]

# Considerable Complications with Time Series Data



- noisy, irregularly sampled

- spurious data

- telltale signature event may not have happened yet

class: *microlensing*

2D image classification: Machine-Learning with human input

>1000:1 rejection of bogus candidates (prelim. cuts + machine learning)

most subtractions are bogus...

...*but a long tail of astrophysical goodness*

275,000 very likely real

10M PTF subtractions (1 month of data)

# Major Challenge:

how do we use **domain knowledge** & **known ("labelled") instances** to create a classifier?

*traditional fitting, machine learning, ...*
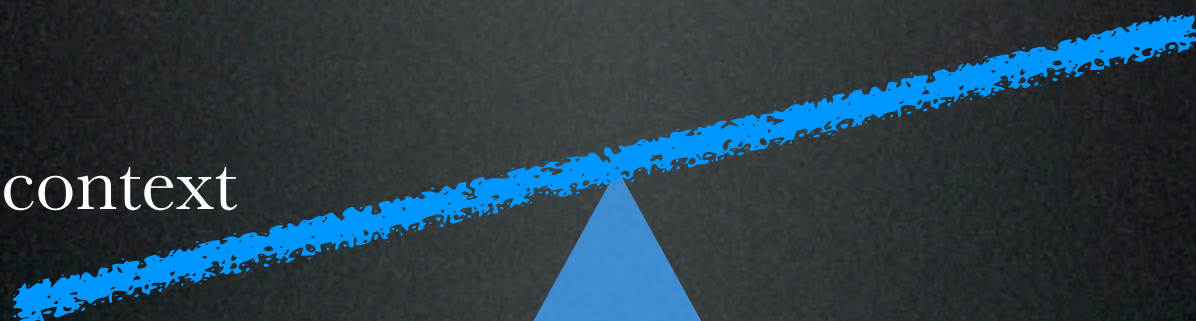
# Machine-Learning Approach to Classification

| Data | Utility for Classification |
|---|---|
| **Time Series** (e.g. color, brightness change, etc.) | • comparison to previously observed sources, & theoretical/ numerical models <br> •historical images: extend time baseline |
| **Context** (e.g. sky location, nearest galaxy type) | situational awareness: expectations of different classes |

*less data regime*

context

time-series

*more data regime*

# Feature Extraction: Homogenizing Heterogenous Data

"Features": real-number metrics that describe the time-domain characteristics & context of a source.

**variability metrics:**
e.g. Stetson indices, $\chi^2/\text{dof}$ (constant hypothesis)

**shape analysis**
e.g. skewness, kurtosis, Gaussianity

**periodic metrics:**
e.g. dominant frequencies in Lomb-Scargle, phase offsets between periods

**context metrics**
e.g. distance to nearest galaxy, type of nearest galaxy, location in the ecliptic plane
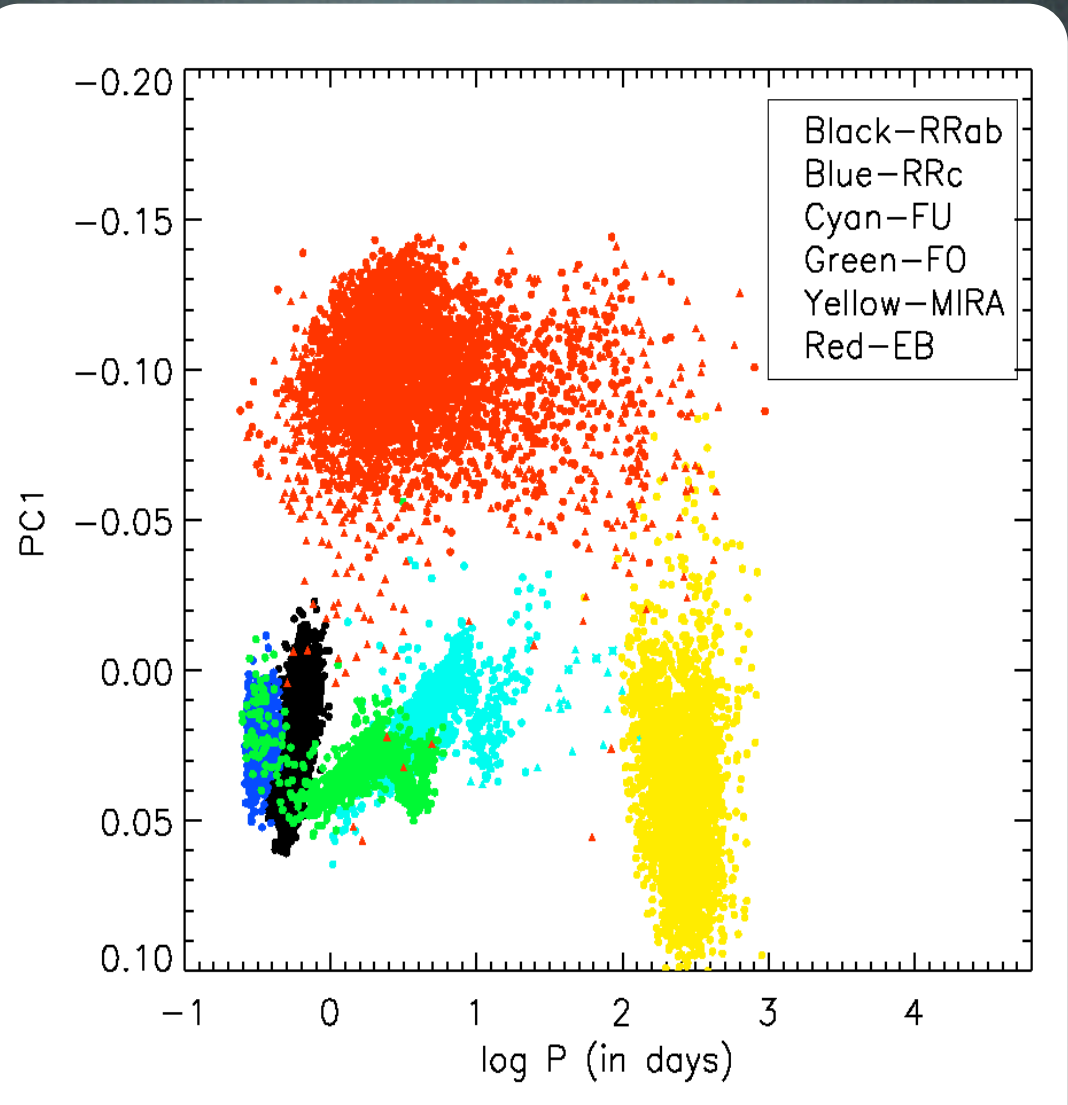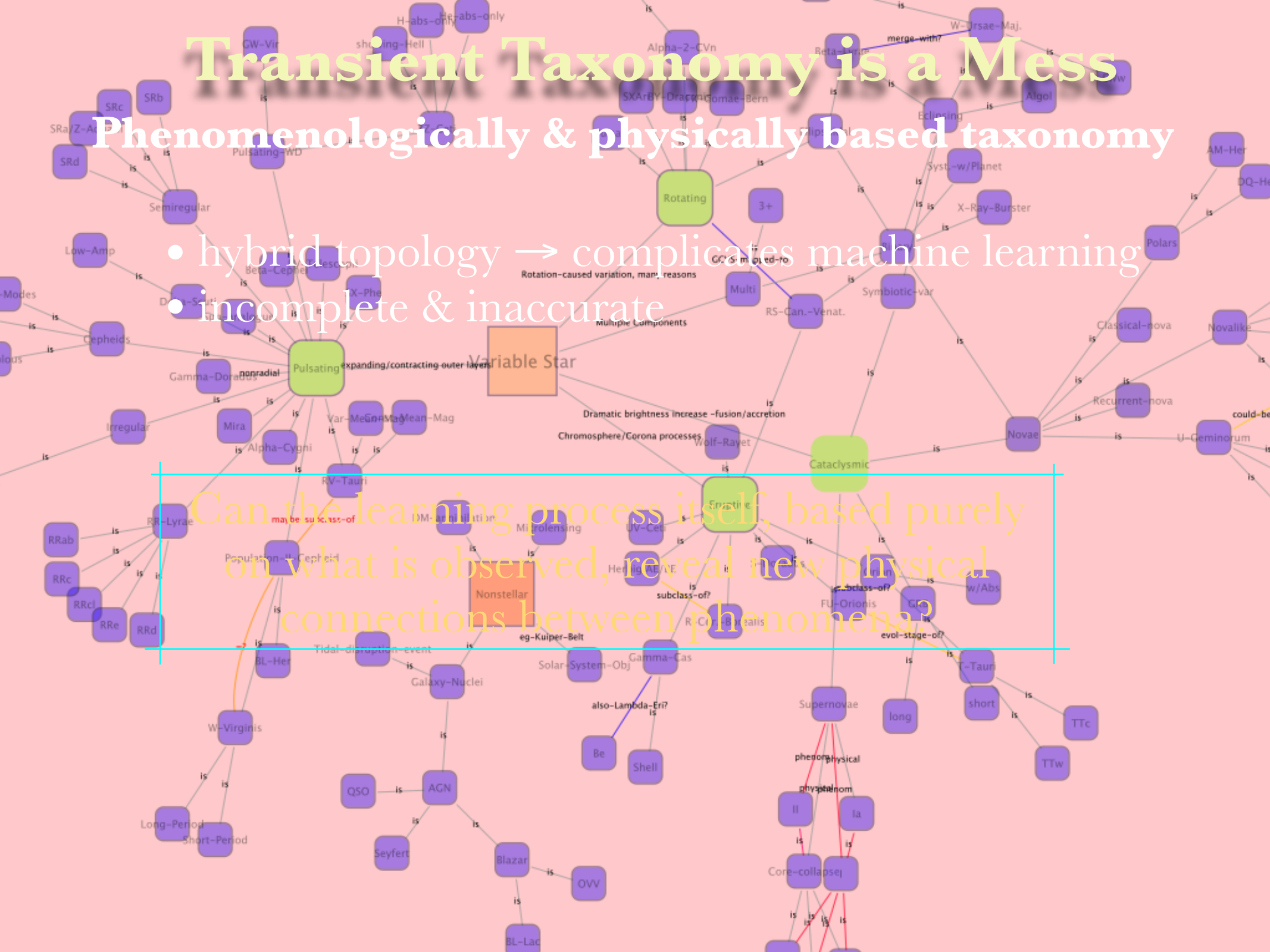
**Fig. 11.** The classification based on PC1 obtained from PCA of 100 interpolated magnitudes for the phase from 0 to 1 in steps of 0.01.

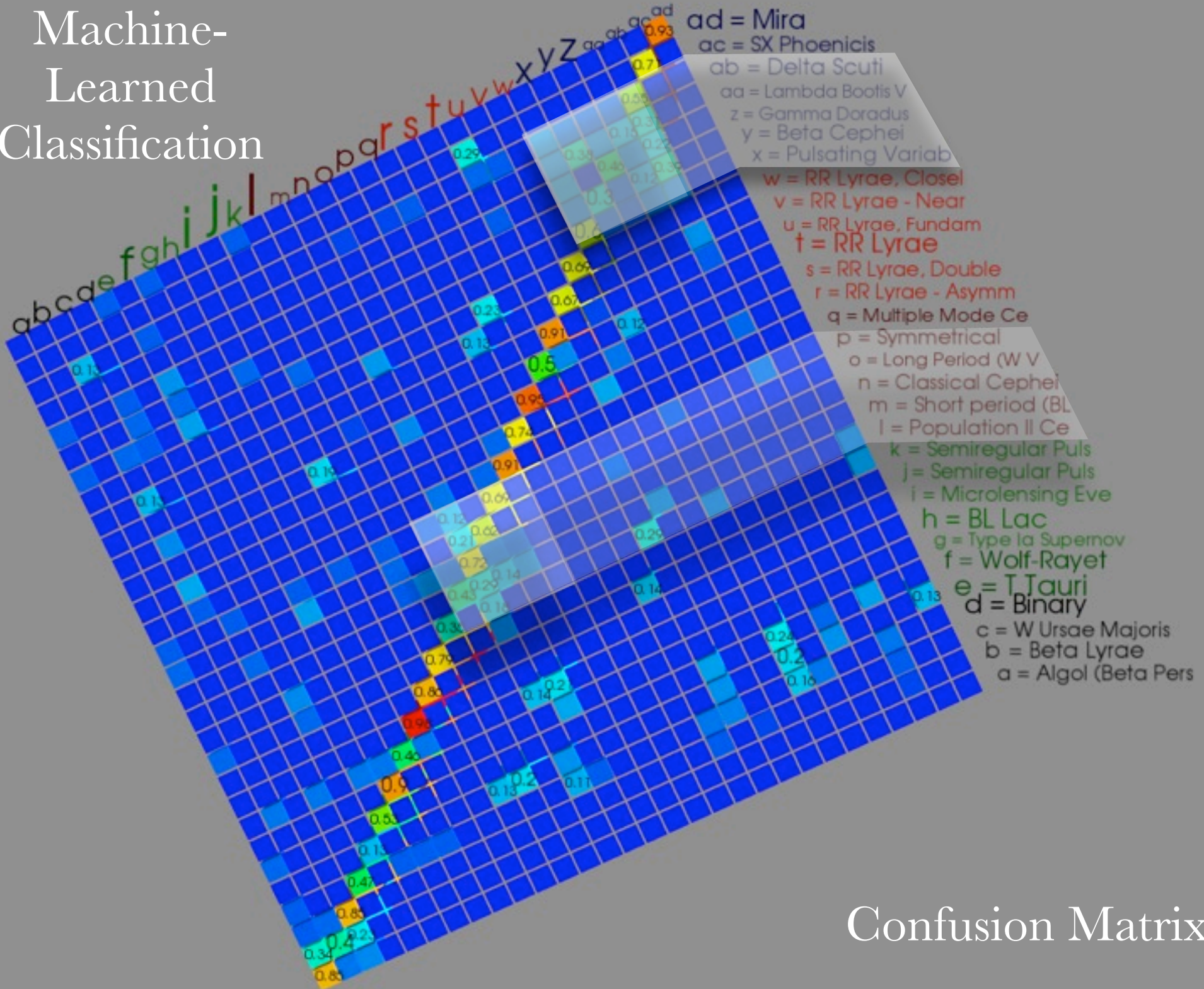Deb & Singh+09

# Transient Taxonomy is a Mess

**Phenomenologically & physically based taxonomy**

- hybrid topology → complicates machine learning
- incomplete & inaccurate

Can the learning process itself, based purely on what is observed, reveal new physical connections between phenomena?

Machine-Learned Classification

Confusion Matrix

ad = Mira
ac = SX Phoenicis
ab = Delta Scuti
aa = Lambda Bootis V
z = Gamma Doradus
y = Beta Cephei
x = Pulsating Variab
w = RR Lyrae, Closel
v = RR Lyrae - Near
u = RR Lyrae, Fundam
t = RR Lyrae
s = RR Lyrae, Double
r = RR Lyrae - Asymm
q = Multiple Mode Ce
p = Symmetrical
o = Long Period (W V
n = Classical Cephei
m = Short period (BL
l = Population II Ce
k = Semiregular Puls
j = Semiregular Puls
i = Microlensing Eve
h = BL Lac
g = Type Ia Supernov
f = Wolf-Rayet
e = T Tauri
d = Binary
c = W Ursae Majoris
b = Beta Lyrae
a = Algol (Beta Pers

# 1. Parallelize the Learning Phase of Machine Learning

**Problem:**
frameworks like Weka (`http://www.cs.waikato.ac.nz/ml/weka/`) are not natively parallel. We will need to burst out training requests on specific time/observation vectors & classify quickly with the results

**Solution:**
build a parallel platform for weka
(GridWeka, Weka-parallel etc. are out of date & probably not elegant)
- develop/adapt Mahout (http://lucene.apache.org/mahout/), ML for Hadoop

http://userweb.port.ac.uk/~khusainr/weka

# 1. Parallelize the Learning Phase of Machine Learning

## Problem:

we have errors on our data (both training sets and instances) & we dont know how to deal with them
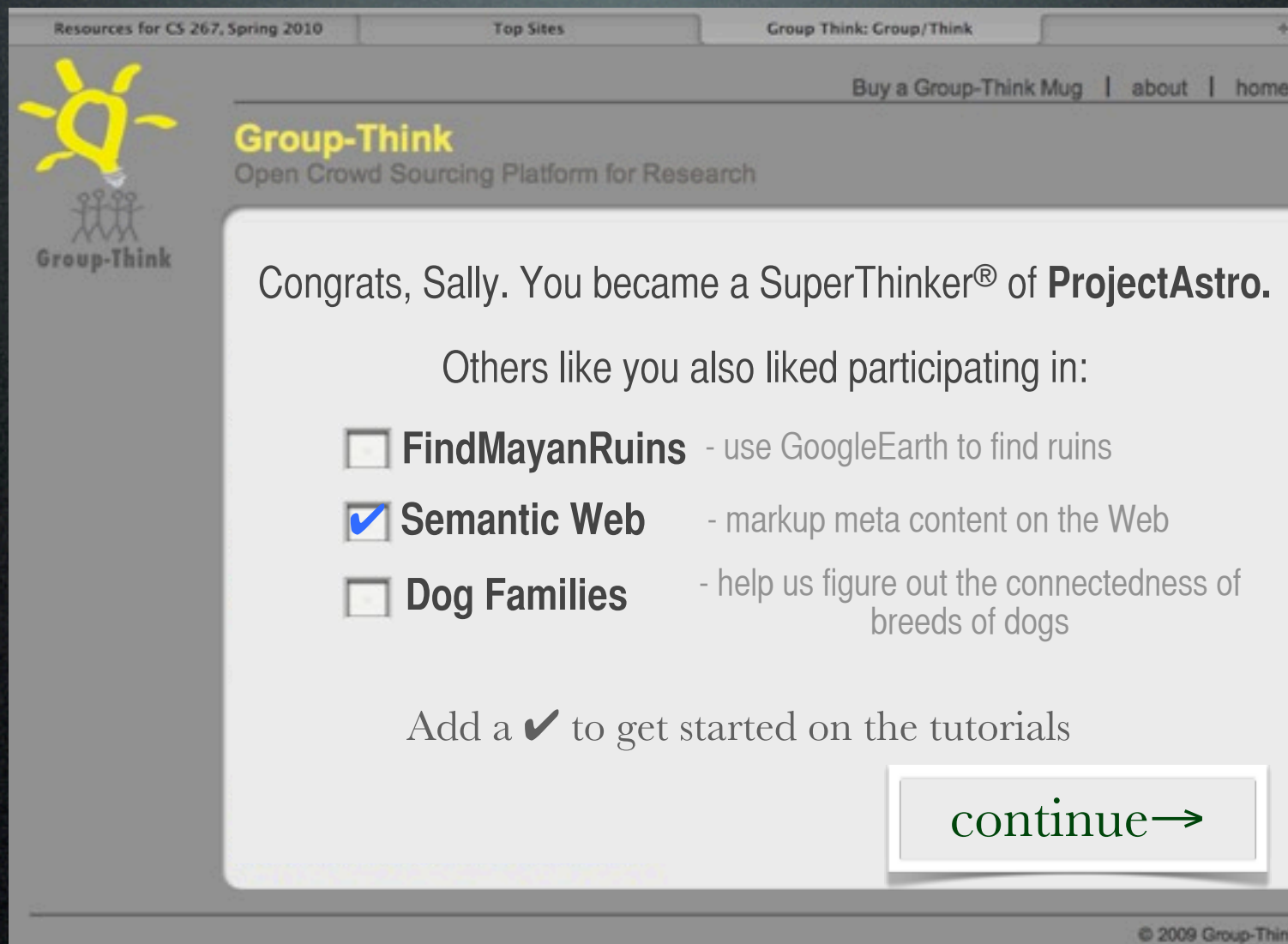
## Sledgehammer Solution:

use a parallel platform to generate distribution of trained models & apply to distribution of instance-based sets

# 1. Parallelize the Learning Phase of Machine Learning

flux

time

68% confidence interval

time

time

time

"fastest rise" feature extractor

0.2

0.5

0.32

# 2. Build a General Crowdsourcing Platform (GroupThink2.0)

- production scale site (GoogleAppEngine or elsewhere), allowing interconnection of projects

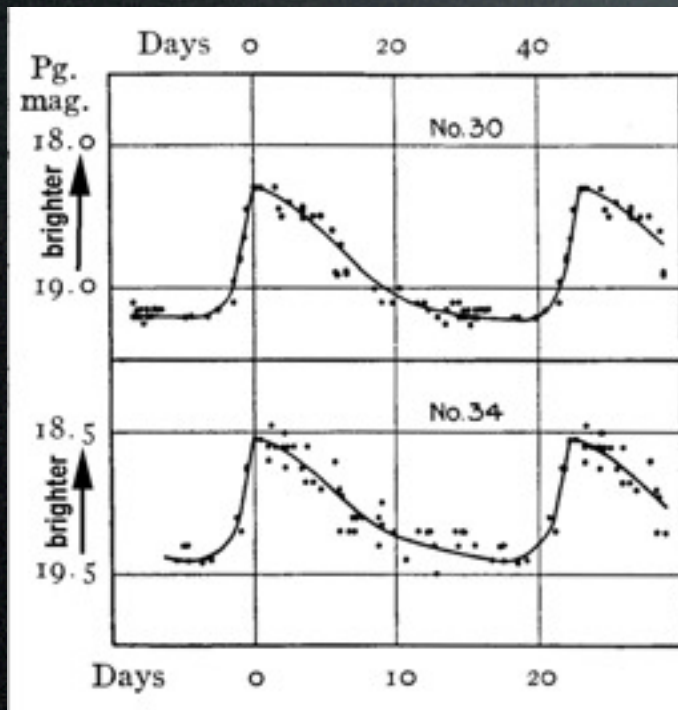## 2. Build a General Crowdsourcing Platform (GroupThink2.0)

- build innovative analytics plugins for projects;
- could require grid/cloud-based analysis for on-the-fly results
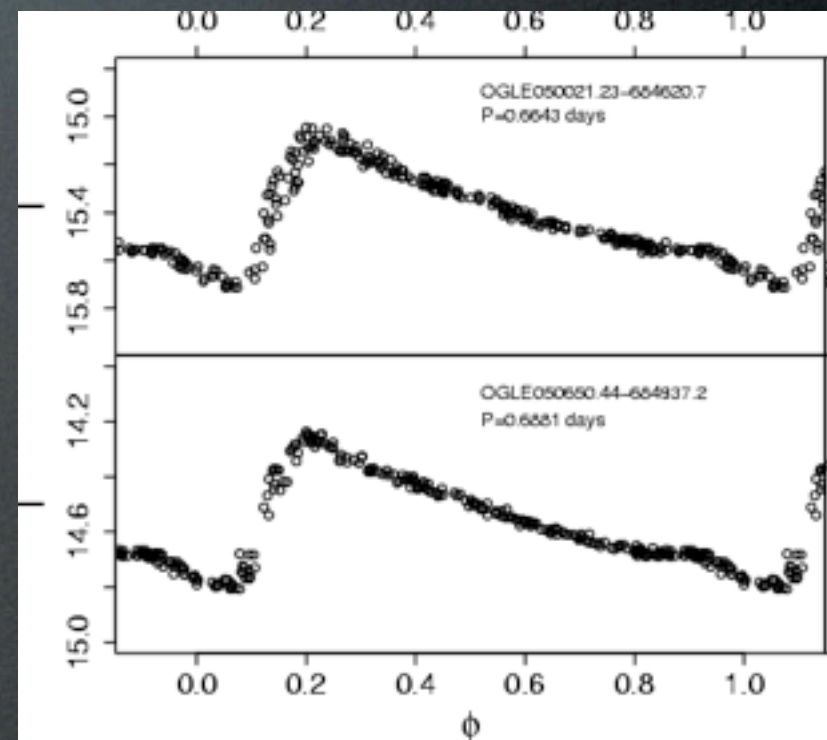
# 3. Parallelized Genetic Programming for Feature Discovery

Instead of handcoding "features" for ML,
using GP (in parallelized environment) to
***discover features*** which give the best
classification



Cepheid

vs.

RR Lyrae

# 4. Parallelized Visual Exploration Tool

allow the armchair astronomer to ask complex
questions of the databases & visualize and interact
with the results (100M+ rows)

| sdss + simbad positions ⬍ | Select query from history... ⬍ | ◀ |

```
select jsb_source.ptfname, oon.val, oan.node_name,sdss.bestz,sdss.bestz_err,sdss.dered_r,jsb_cand.mag_ref, oar_ann.val as cat_offset,oa.val as
sdss_offset from oar_node
join jsb_source on jsb_source.jsb_source_id = oar_node.jsb_source_id
join oar_ann on oar_ann.oar_node_id = oar_node.oar_node_id
join jsb_cand on jsb_cand.lbl_id = jsb_source.initial_lbl_cand_id
join sdss on sdss.jsb_source_id = jsb_source.jsb_source_id
left join oar_ann as oa on ( (oa.jsb_source_id = oar_node.jsb_source_id) and oa.key = 'host_distance_arcsec_sdss')
left join oar_node as oan on ((oan.jsb_source_id =  oar_node.jsb_source_id) and oan.class_type = 'simbad')
left join oar_ann as oon on oon.oar_node_id = oan.oar_node_id
where oar_node.class_type = "sdss"
 and oar_ann.key = "host_distance_arcsec_cat"
```
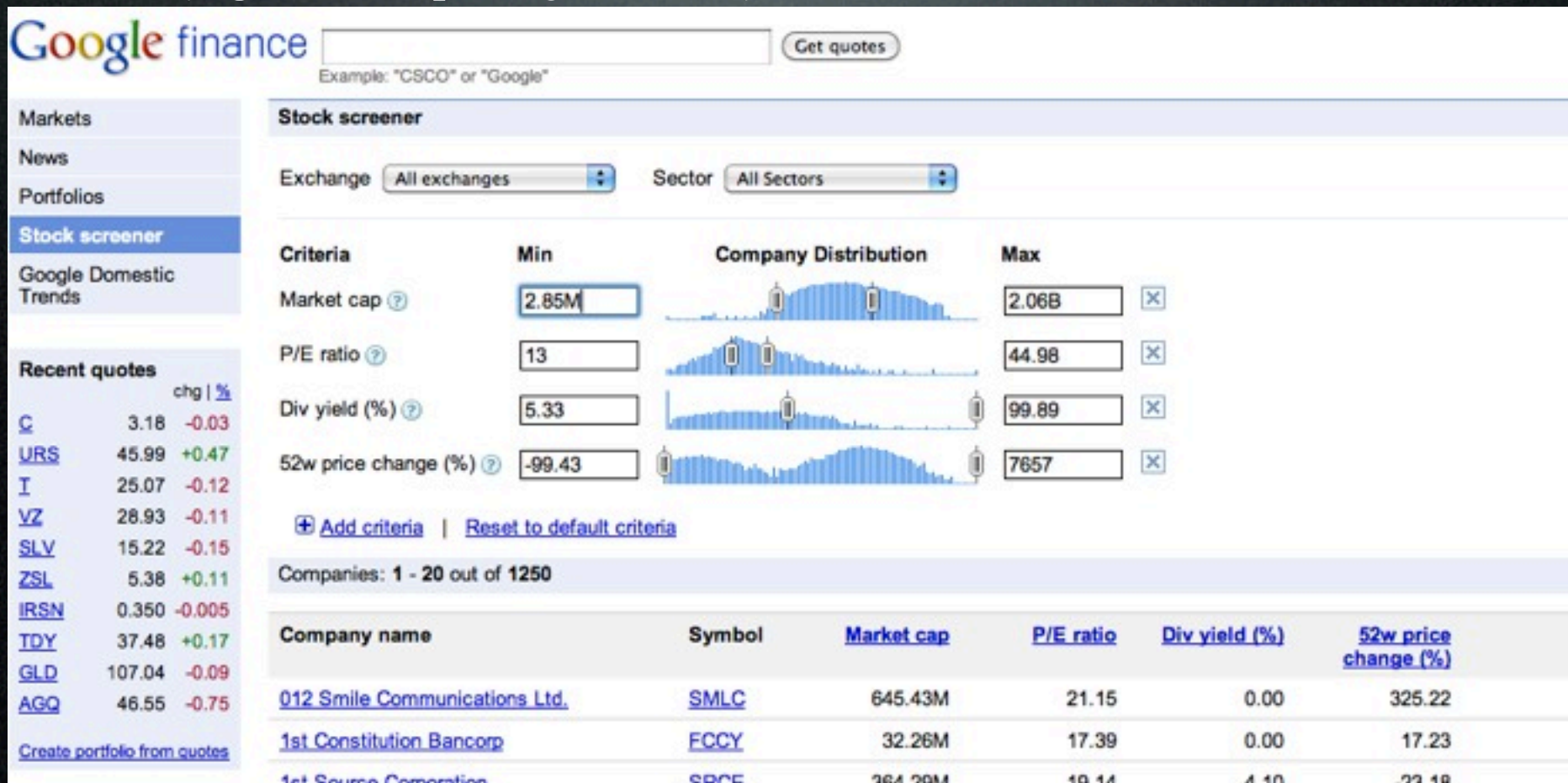
Run

| PTFname | val | node_name | bestz | bestz_err | dered_r | mag_ref |
|---------|-----|-----------|-------|-----------|---------|---------|
| 10bgh | NUL | NUL | 0.1771 | 0.0108 | 17.3009 | 17.263 |
| 10bgb | NUL | NUL | 0.6984 | 0.1571 | 21.7929 | 18.156 |
| 10bfg | NUL | NUL | 0.075 | 0.0294 | 17.9426 | 18.076 |
| 10bea | QSO | extragalactic | 1.2264 | 0.0018 | 17.736 | 17.684 |
| 10bea | NUL | qso | 1.2264 | 0.0018 | 17.736 | 17.684 |
| 10bdv | QSO | extragalactic | 0.5409 | 0.0011 | 18.2925 | 18.302 |

# 4. Parallelized Visual Exploration Tool

allow the armchair astronomer to ask complex questions of the databases & visualize and interact with the results

- parallel database calls with embedded custom code (e.g. Hadoop SQL "hive")

# Resources

1. dotastro.org

2. Harvard TimeSeries Center:
   http://timemachine.iic.harvard.edu/

3. "The Fourth Paradigm: Data-Intensive Scientific Discovery"
   http://research.microsoft.com/en-us/collaboration/fourthparadigm/