



Human in the loop:
**Active Learning in
Astronomy**

*Artificial Intelligence in Astronomy - ESO, Germany
22 June 2019*

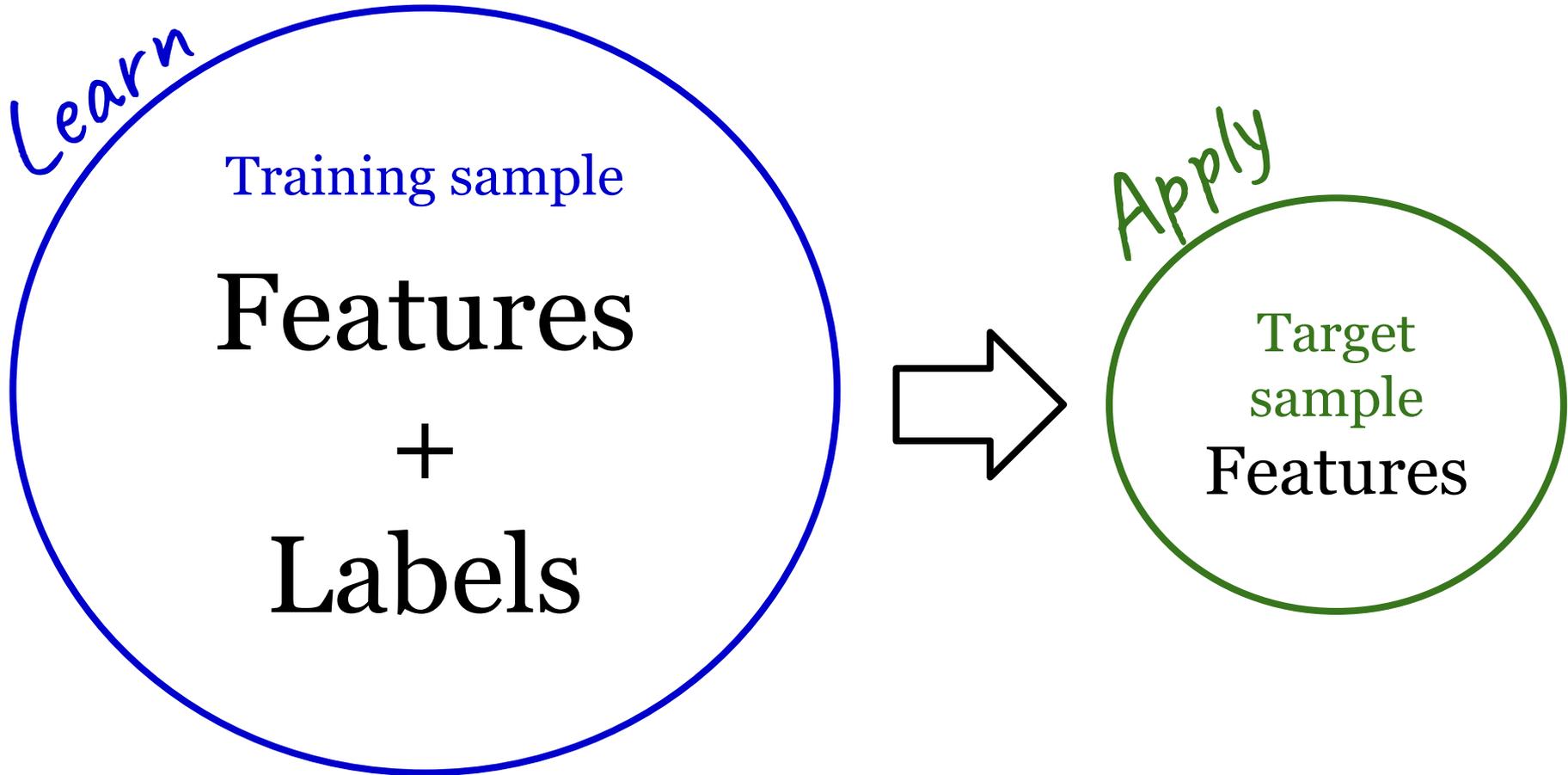
Emille E. O. Ishida

*Laboratoire de Physique de Clermont - Université Clermont-Auvergne
Clermont Ferrand, France*



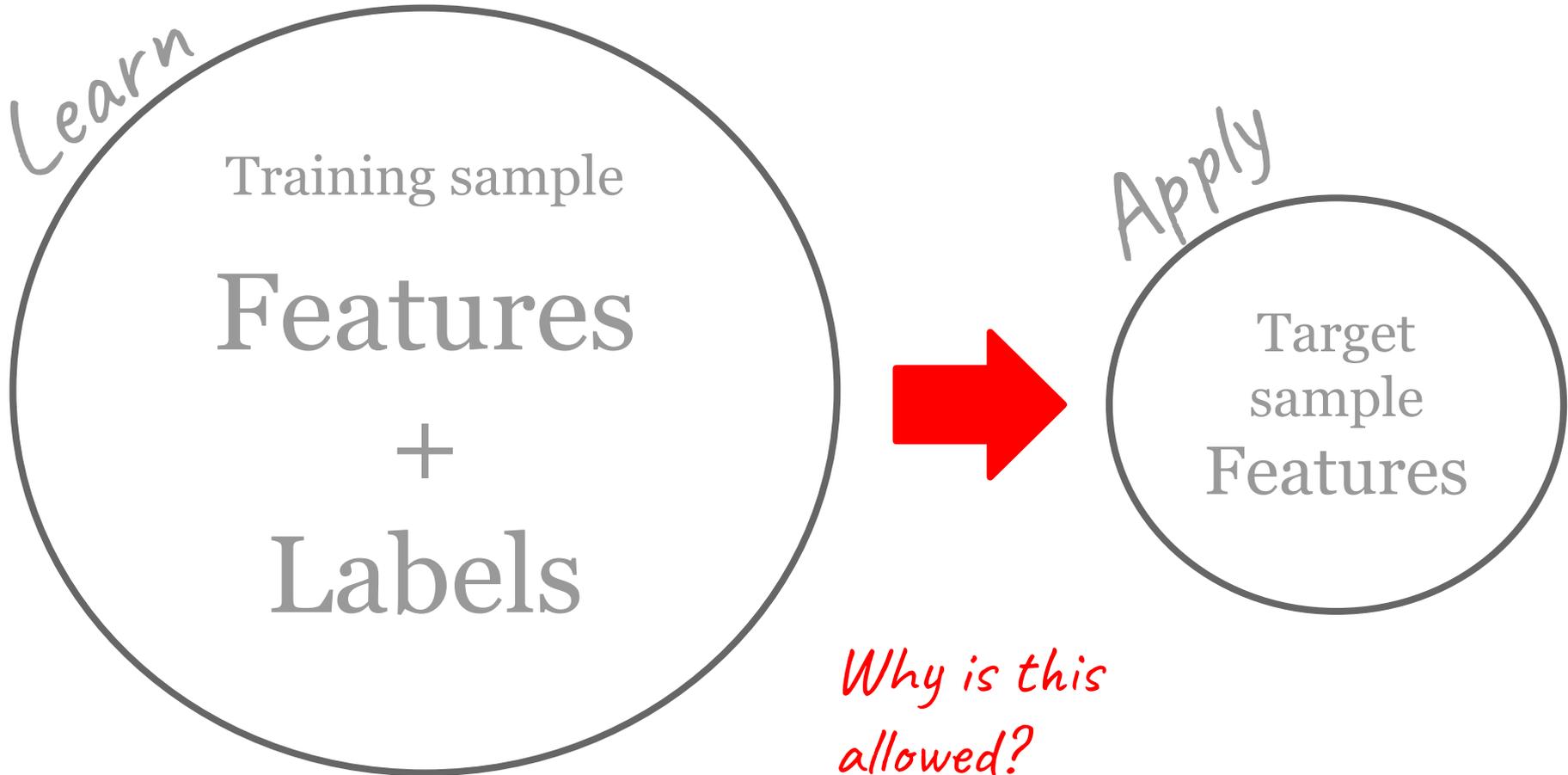
Supervised Learning

Ideal data scenario



Supervised Learning

Learn by example



Supervised ML model

data	training, target
\mathcal{X}	set of all samples, x
\mathcal{Y}	set of possible labels, y
h_{train}	learner: $y_{est;i} = h_{train}(x_i)$
L	Loss function

Hypothesis:
Training is
representative
of target

Data generation model:

$$x_i \sim P_{\mathcal{X}}$$

$f \rightarrow$ true labeling function, $y_i = f(x_i)$

$$L_{data,f}(h) \equiv P_{x \sim data} (h_{train}(x) \neq f(x))$$

Supervised ML model

data training, target
 X set of all samples, x
 Y set of possible labels, y

Machine Learning algorithm

Hypothesis:
training is
representative
of target

h_{train} learner: $y_{est;i} = h_{train}(x_i)$
 L Loss function

Data generation model:

$$x_i \sim P_X$$

$f \rightarrow$ true labeling function, $y_i = f(x_i)$

$$L_{data,f}(h) \equiv P_{x \sim data}(h_{train}(x) \neq f(x))$$

Supervised ML model

data	training, target
\mathcal{X}	set of all samples, x
\mathcal{Y}	set of possible labels, y
h_{train}	learner: $y_{est;i} = h_{train}(x_i)$
L	Loss function

Hypothesis:
Training is
representative
of target

Data generation model:

$$x_i \sim P_{\mathcal{X}}$$

$f \rightarrow$ true labeling function, $y_i = f(x_i)$

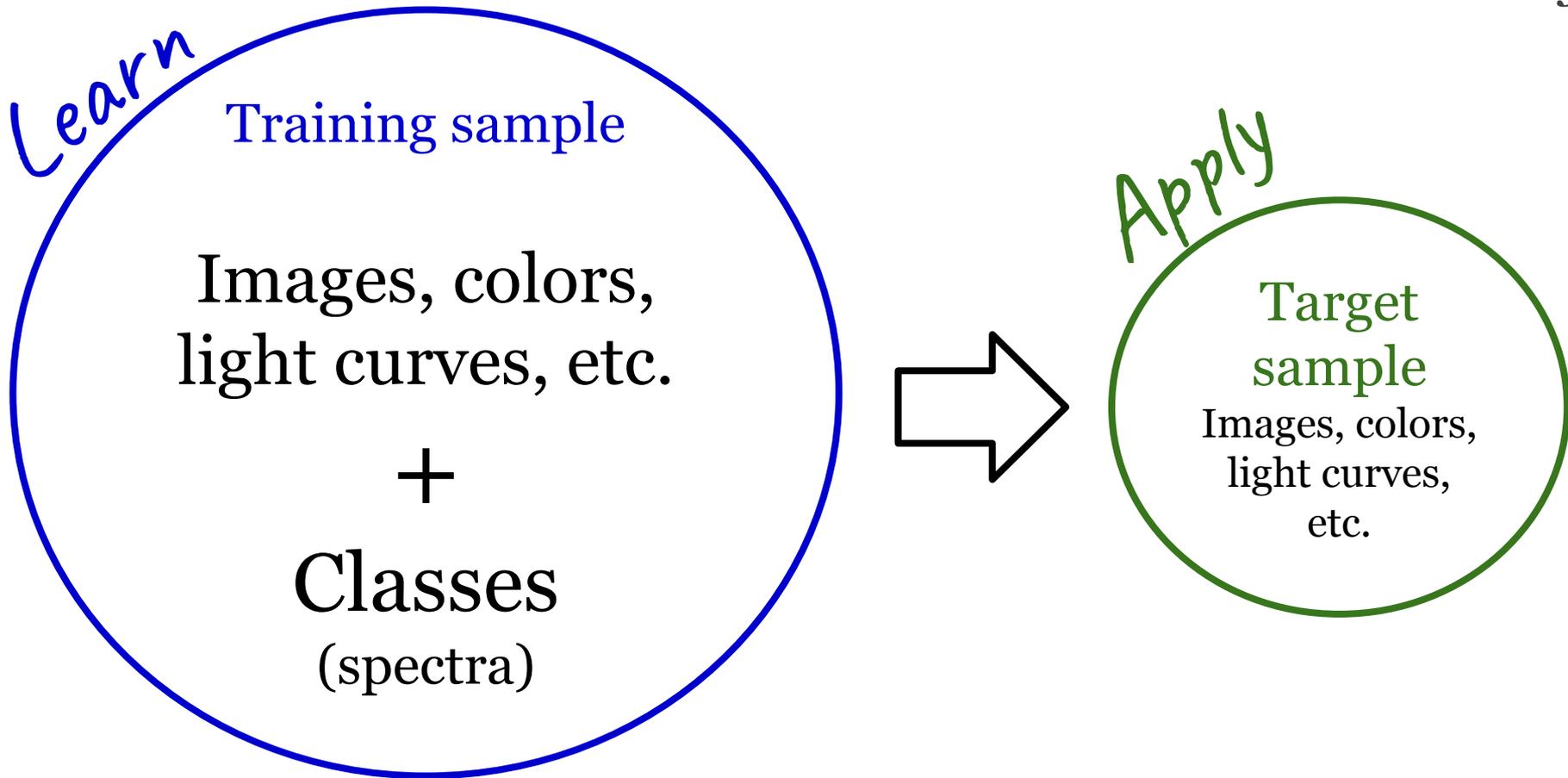
$$L_{data,f}(h) \equiv P_{x \sim data} (h_{train}(x) \neq f(x))$$



How often
does your
data fulfill
these
requirements?

Ideal Supervised learning situation

In astronomy



Astro: supervised learning situation

In astronomy, labels \Rightarrow spectra



Similar examples

Labels are often far too expensive!



amazon

35% OF AMAZON'S REVENUE ARE GENERATED BY IT'S RECOMMENDATION ENGINE.

NETFLIX

75% OF USERS SELECT MOVIES BASED ON NETFLIX'S RECOMMENDATIONS.



Similar examples

Labels are often far too expensive!



Given limited resources, we need recommendation systems!



amazon

35% OF AMAZON'S REVENUE ARE GENERATED BY IT'S RECOMMENDATION ENGINE.

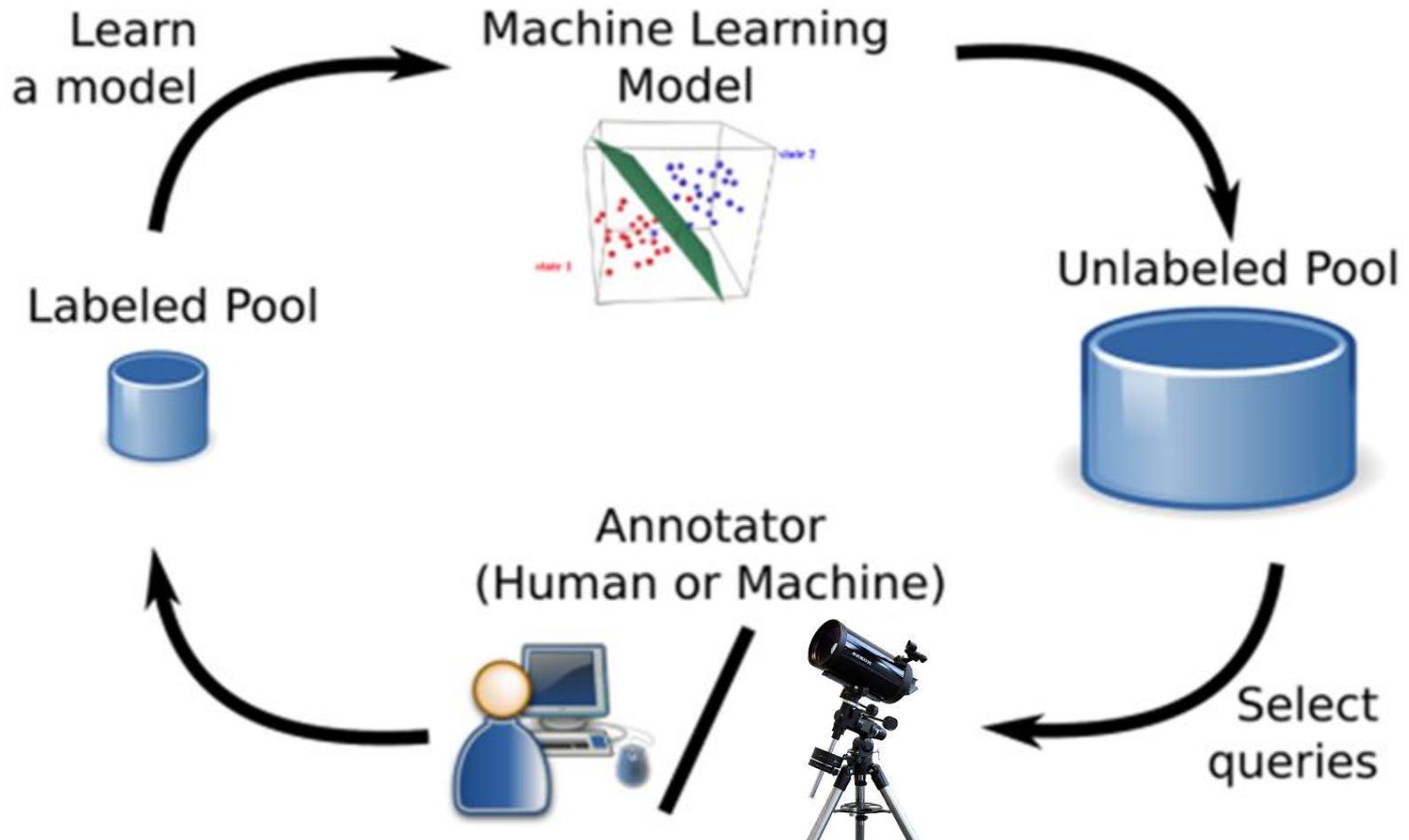
NETFLIX

75% OF USERS SELECT MOVIES BASED ON NETFLIX'S RECOMMENDATIONS.



Active Learning

Optimal classification, minimum training



Optimal Experiment Design

In Statistics literature

$$PQ_{data,f}(x) \equiv P_{x \sim data}(h_{train}(x) \neq f(x) \mid \textit{previous results})$$

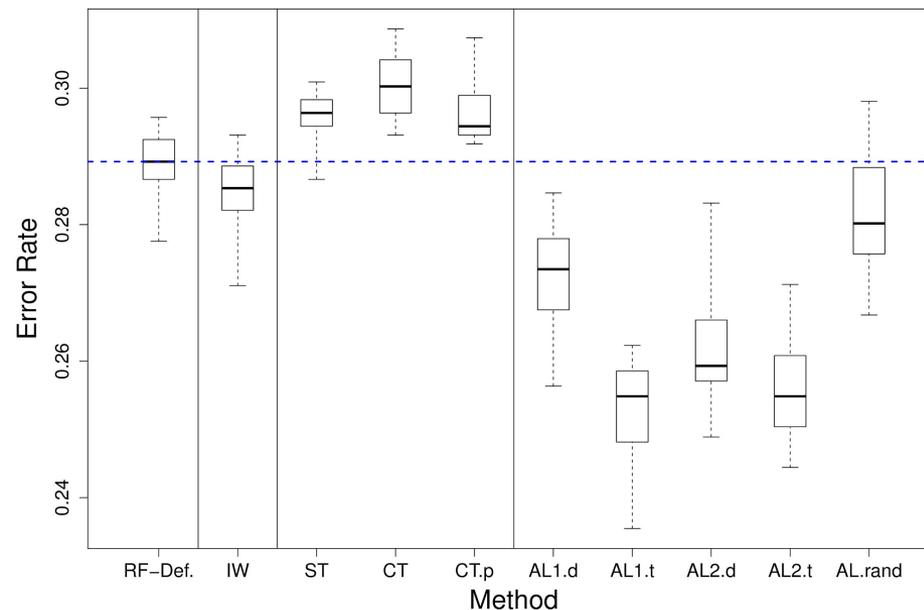
- Pool based
- Generative
- Sequential

Estimation of stellar population parameters

Simulated catalogs - Solorio et al, 2005

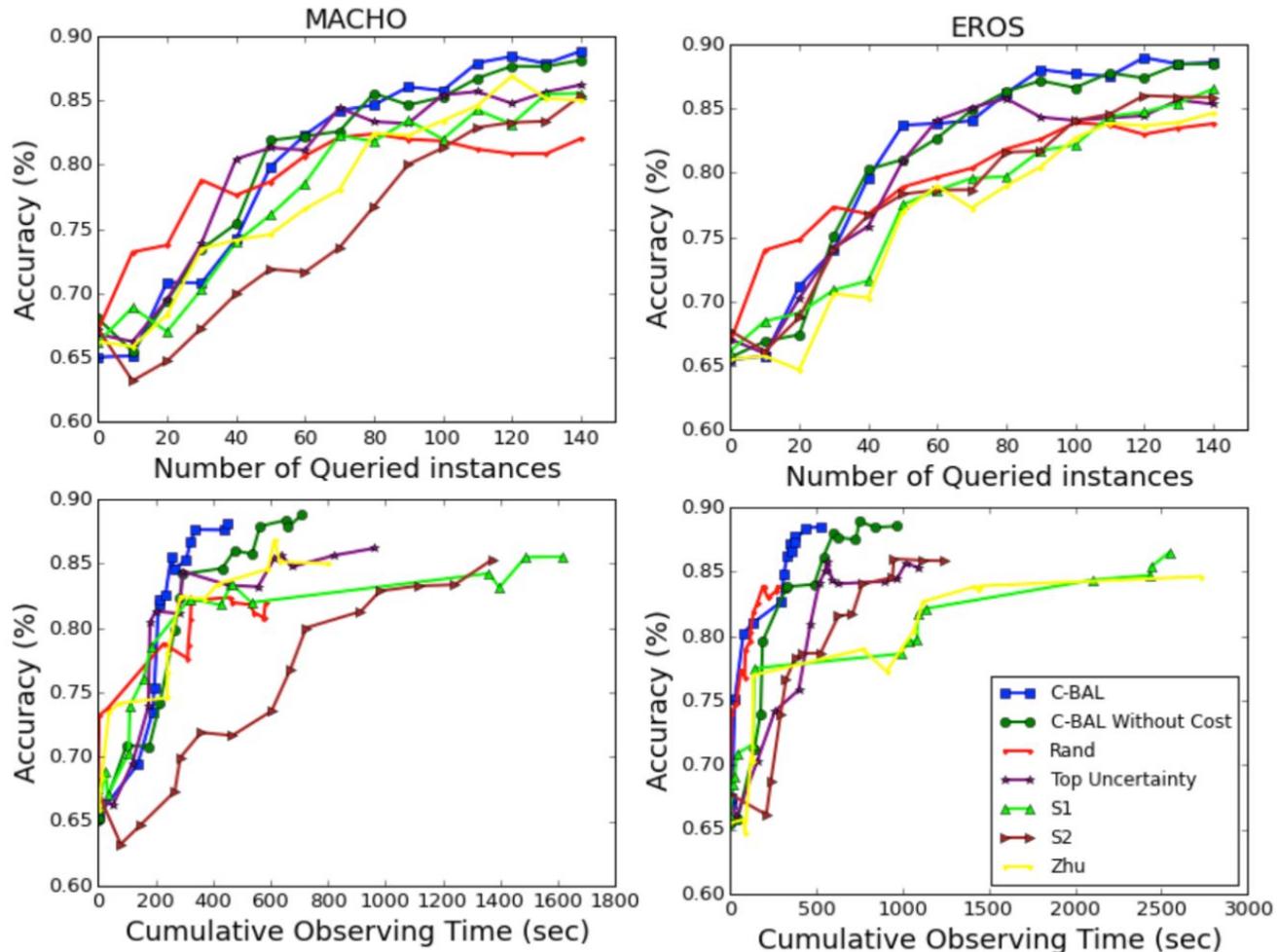
Classification of variable stars

Real data, expected error change - *Richards et al., 2012*



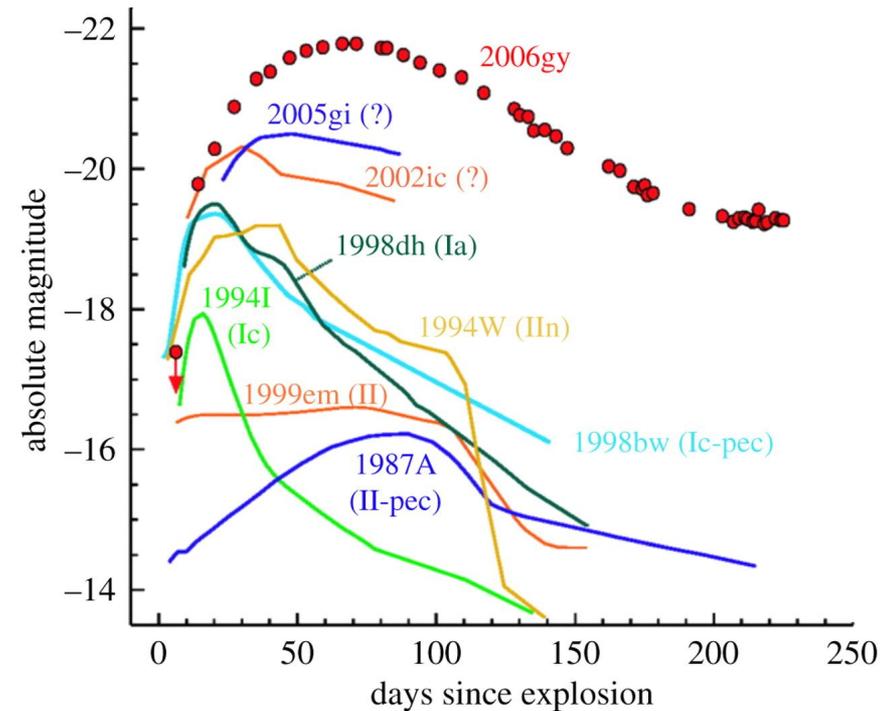
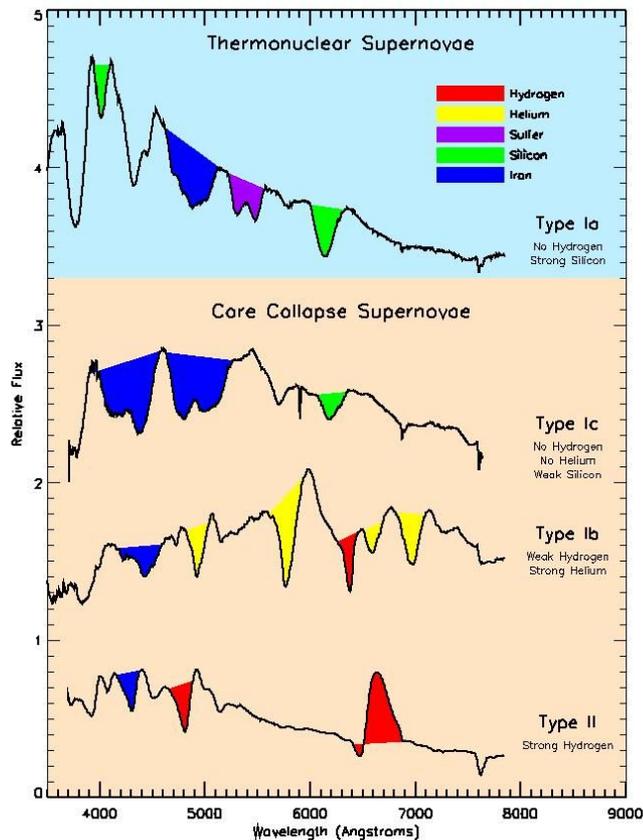
Choosing where to point the telescope

Real catalog - cost sensitive - *Xia et al., 2016*



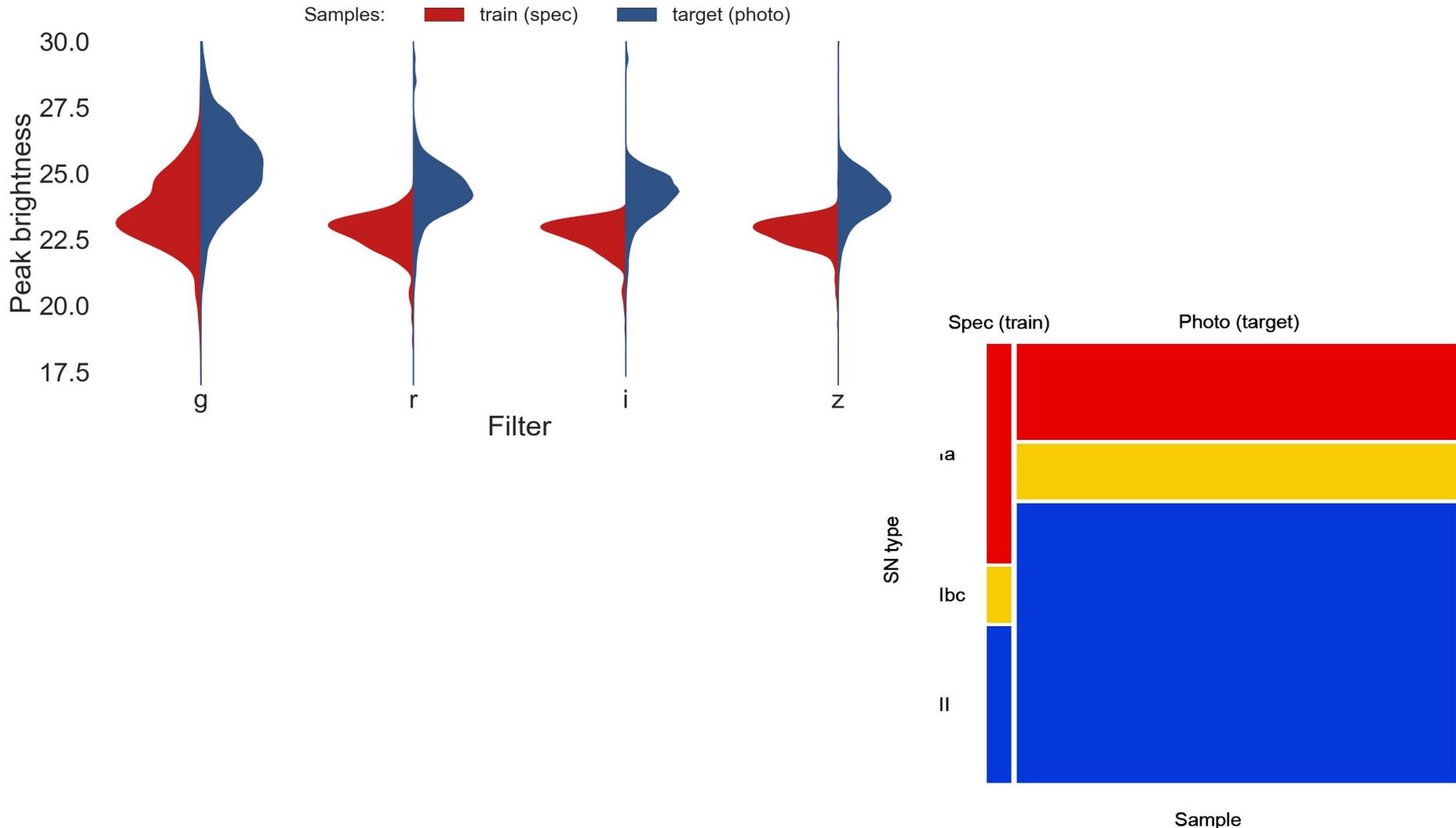
Example of application to astro:

Supernova photometric classification



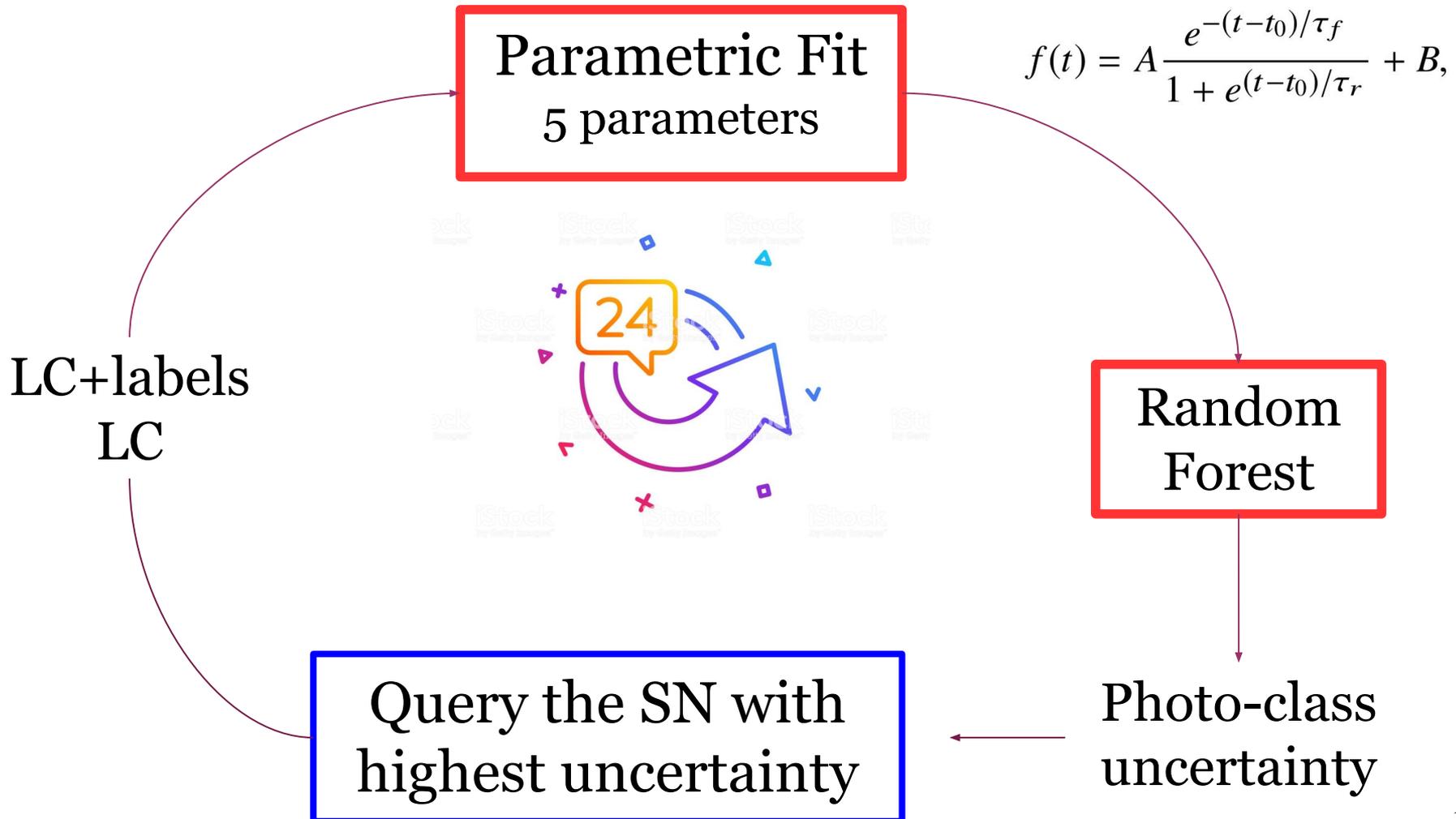
For more on SN classification see Anais Moller's talk this afternoon!

Representativeness



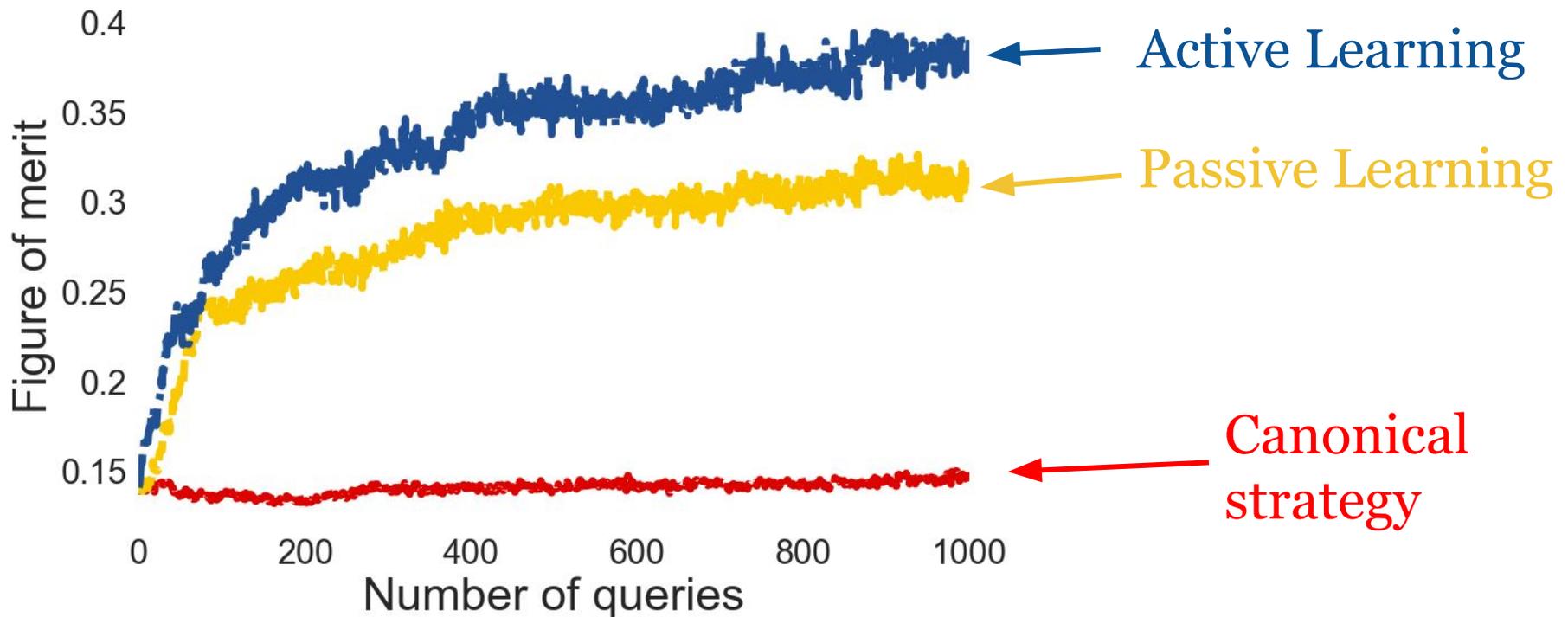
AL for Supernova classification

A strategy



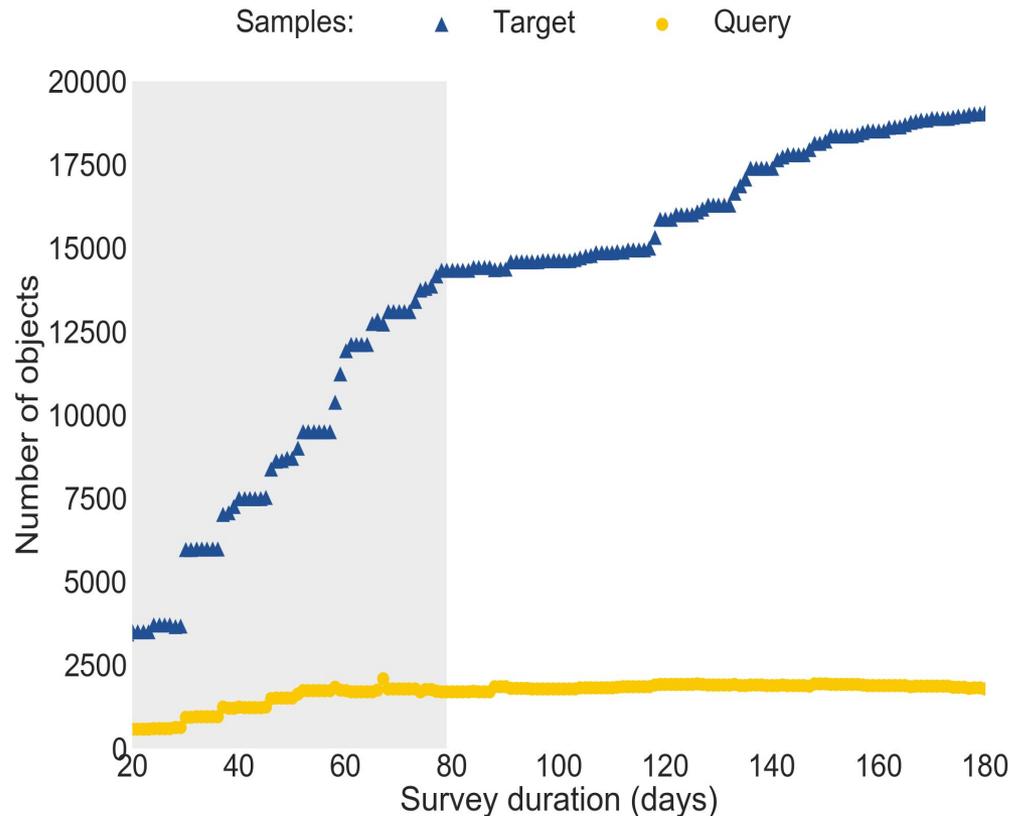
AL for SN classification

Static results



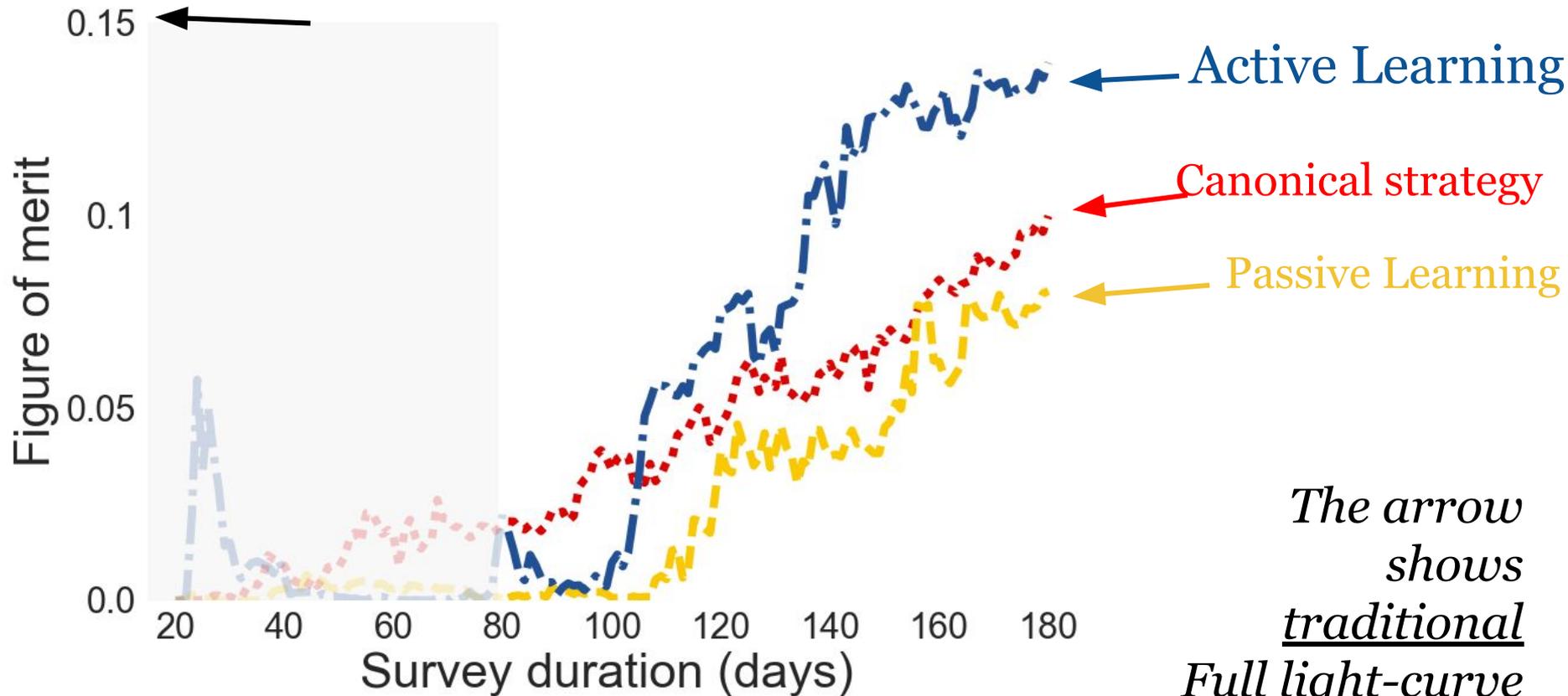
Complications: SNe are transients

Not everything is available for labelling



1. Feature extraction done daily **with available observed epochs until then.**
2. Query sample is also re-defined daily: objects with **r-mag < 24**

Do we even need a training set?



For more on The Cosmostatistics Initiative see Rafael de Souza's talk on Friday!

The arrow shows traditional Full light-curve results with full SNPCC spec 1103 spectra

What comes next?

The Large Synoptic Survey Telescope

Photometric obs:

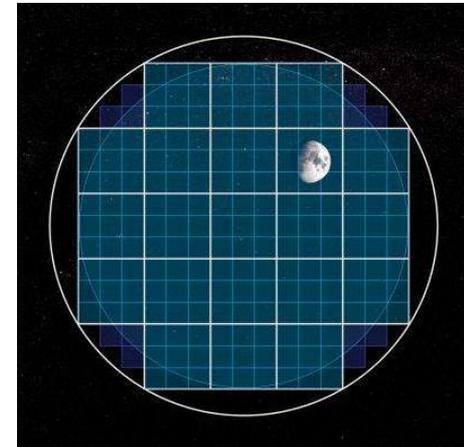
~minute

Spectroscopic obs:

≥ 1 hour (e.g. SDSS)

Multi-fiber spec.

Pointing is not trivial



Camera: **3.2 Giga** pixels and 1.65m

Primary mirror: **8.4m**

Field of view: **3.5 deg**, 40x full moon

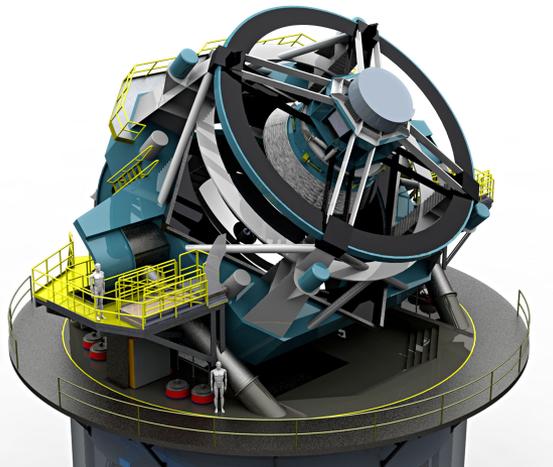
Data production :**15 TB/night**

(3yr LSST=internet today)

~10 million alerts/night

30.000 type Ia SN/yr (today ~1000)

Expected ~ **1000 spectra/yr** (~ 3%)



<https://www.kaggle.com/c/PLAsTiCC-2018>

Featured Prediction Competition

PLAsTiCC Astronomical Classification

Can you help make sense of the Universe?

LSST Project · 1,078 teams · 2 days to go

\$25,000
Prize Money

1,093	1,382	22,430
Teams	Competitors	Entries

[Overview](#) [Data](#) [Kernels](#) [Discussion](#) [Leaderboard](#) [Rules](#) [Host](#)

Overview Edit

- Description**
- Evaluation
- Prizes
- Timeline
- PLAsTiCC's Team

[+ Add Page](#)

Help some of the world's leading astronomers grasp the deepest properties of the universe.

The human eye has been the arbiter for the classification of astronomical sources in the night sky for hundreds of years. But a new facility -- the [Large Synoptic Survey Telescope \(LSST\)](#) -- is about to revolutionize the field, discovering 10 to 100 times more astronomical sources that vary in the night sky than we've ever known. Some of these sources will be completely unprecedented!

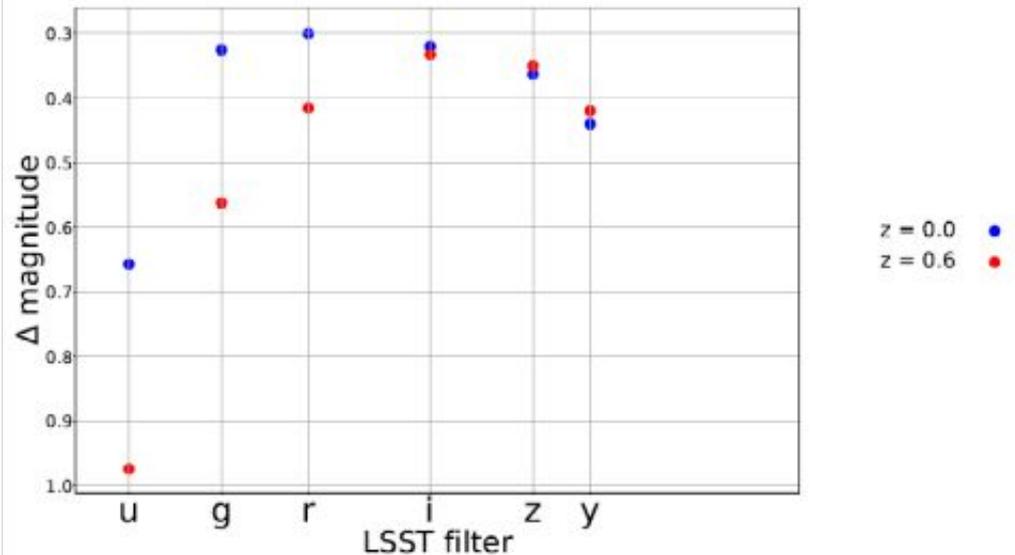
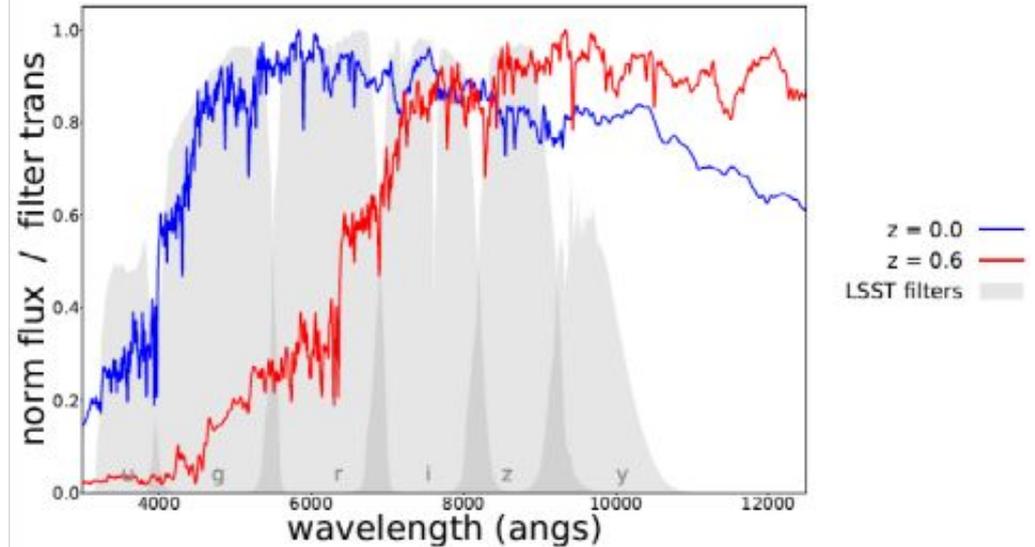
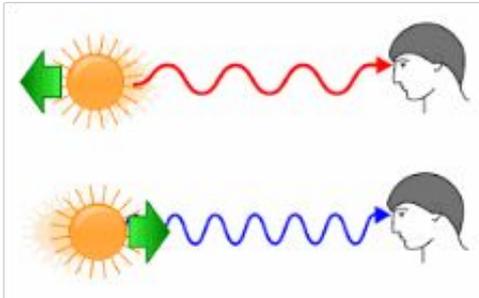


For more on PLAsTiCC see Mi Dai's talk on Wednesday!

Example of application to astro:

Photometric redshift

Idea



AL for Photo-Z

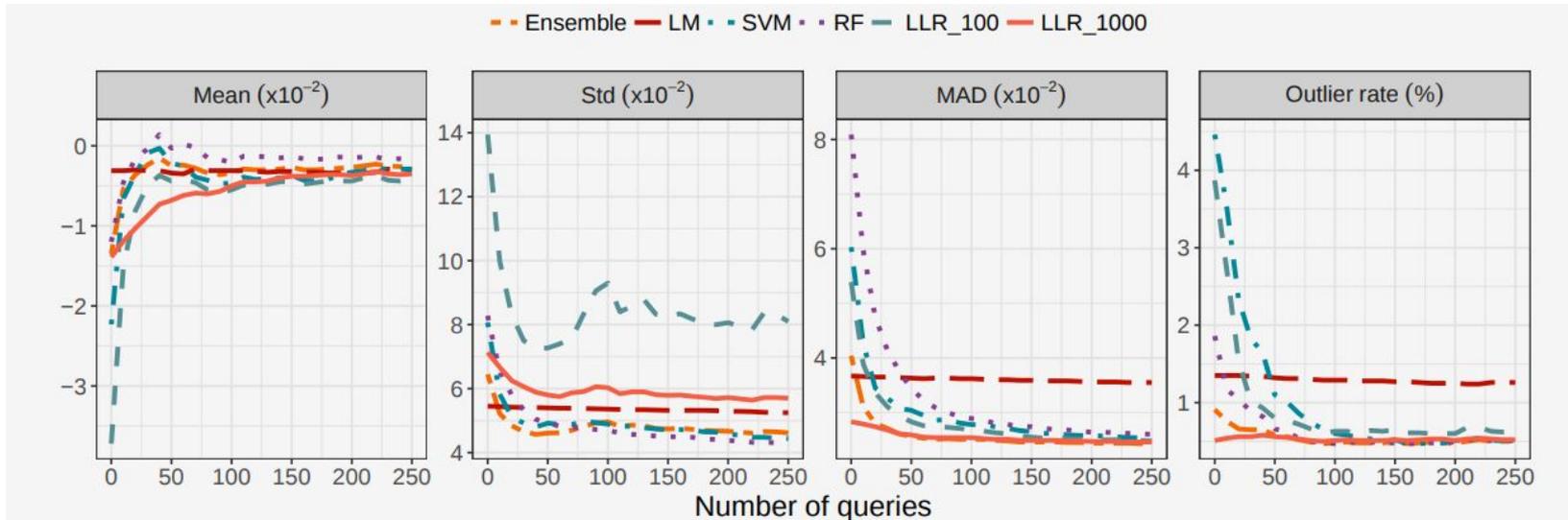
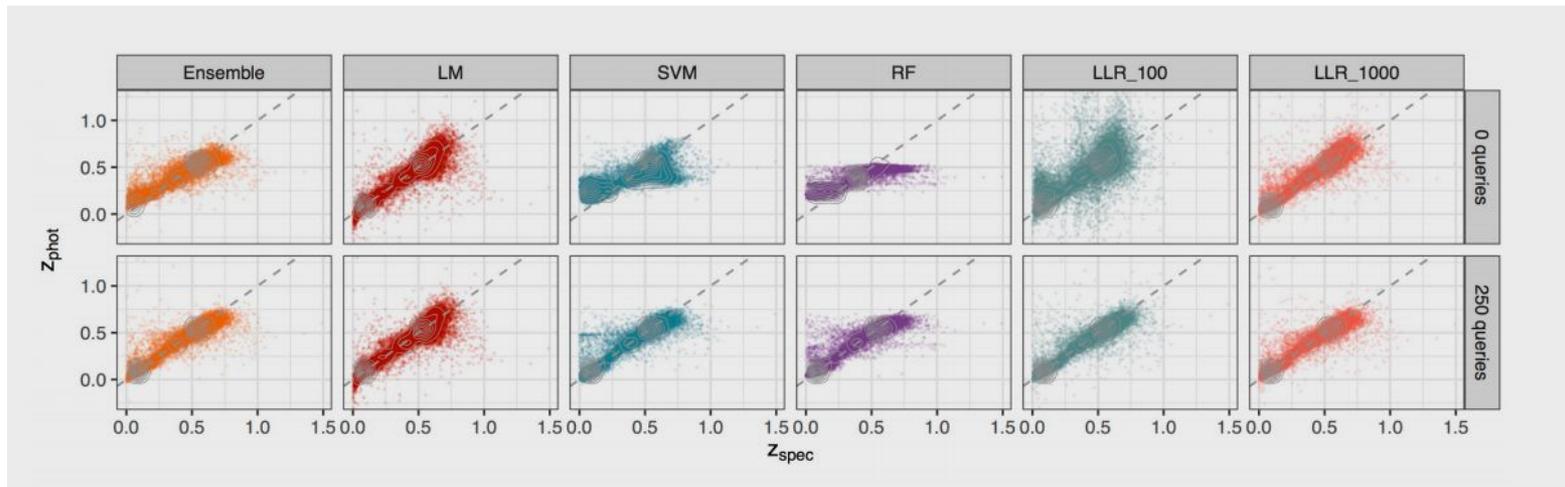
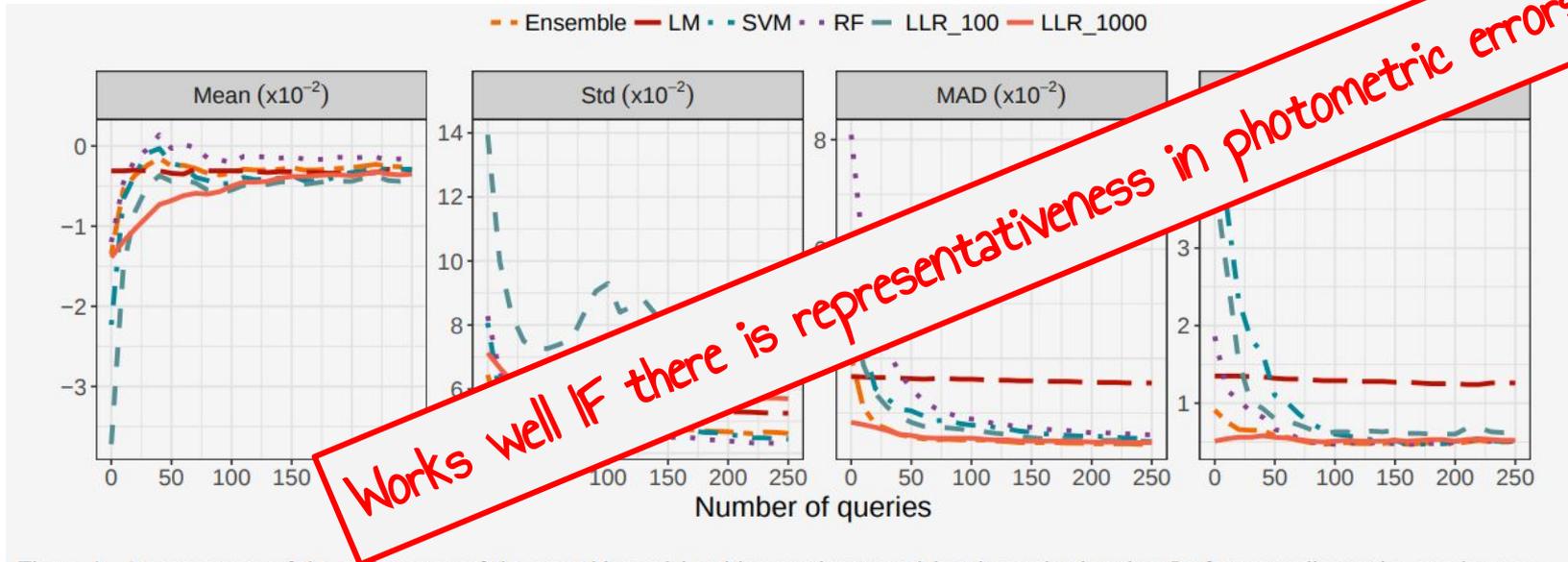
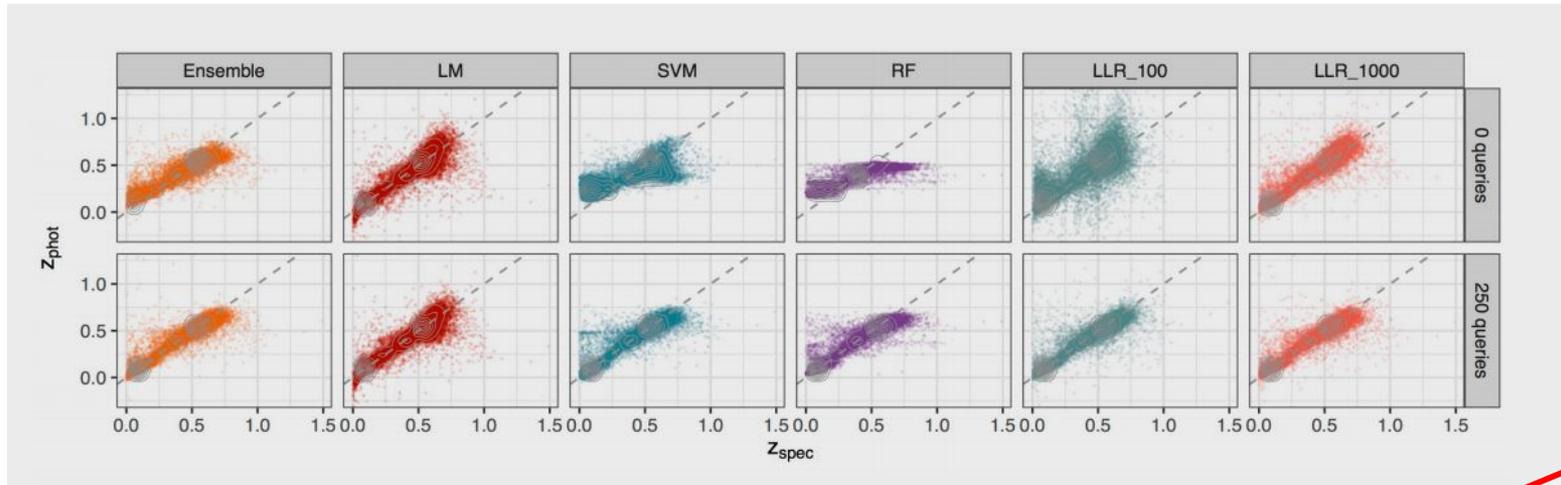


Figure 4. An assessment of the performance of the ensemble model and its constituent models using active learning. Performance diagnostics are shown as a function of the number of queries.

AL for Photo-Z



Works well IF there is representativeness in photometric errors!

Figure 4. An assessment of the performance of the ensemble model and its constituent models using active learning. Performance diagnostics are shown as a function of the number of queries.

Take home message 1:

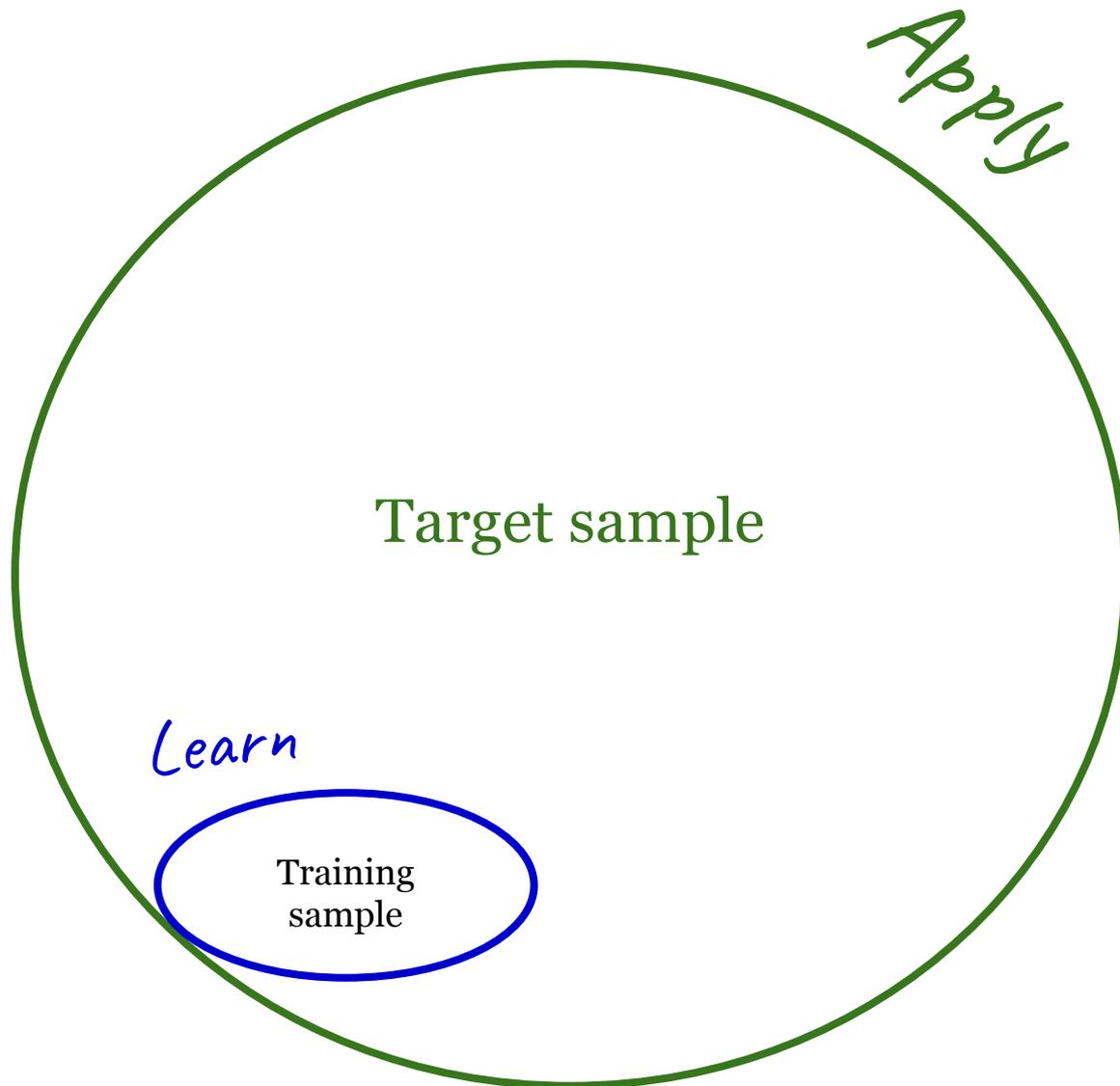
Astronomy needs
optimized samples and
algorithms for
Machine Learning
applications

This means

Interdisciplinarity is the key

Text-book machine learning methods must be adapted
to the peculiarities of astronomical data

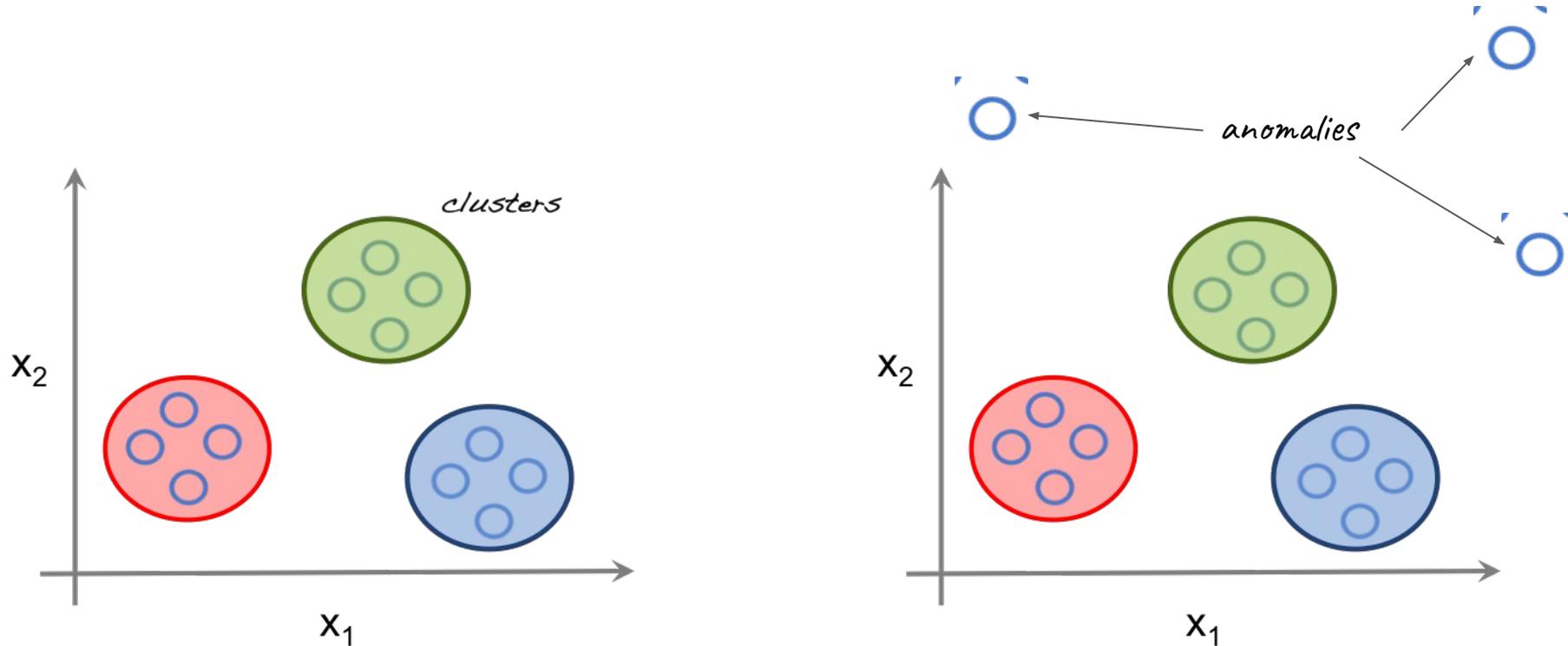
Summary: Supervised Learning



“How do we optimize machine learning results with a minimum number of labeled training instances?”

**Adaptive
Learning
designed for
astronomical
data**

Clustering and Anomaly Detection

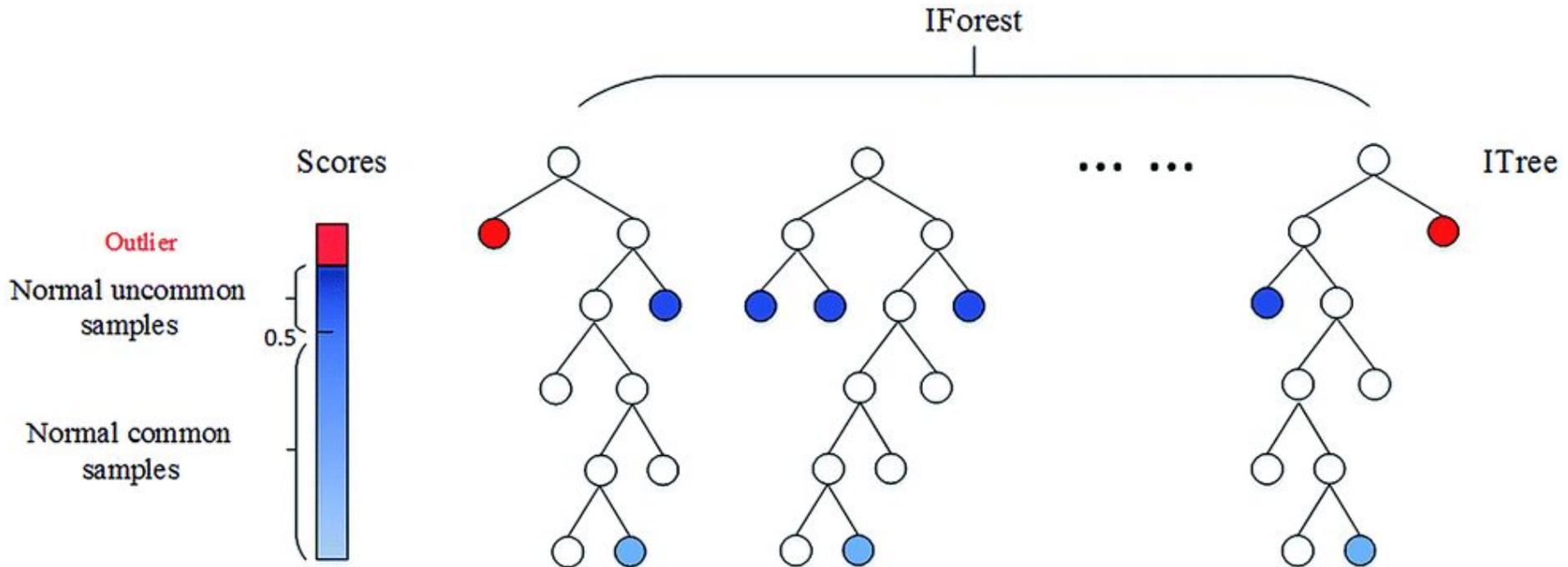


"An anomaly is an observation which deviates so much from the other observations as to arouse suspicions that it was generated by a different mechanism"

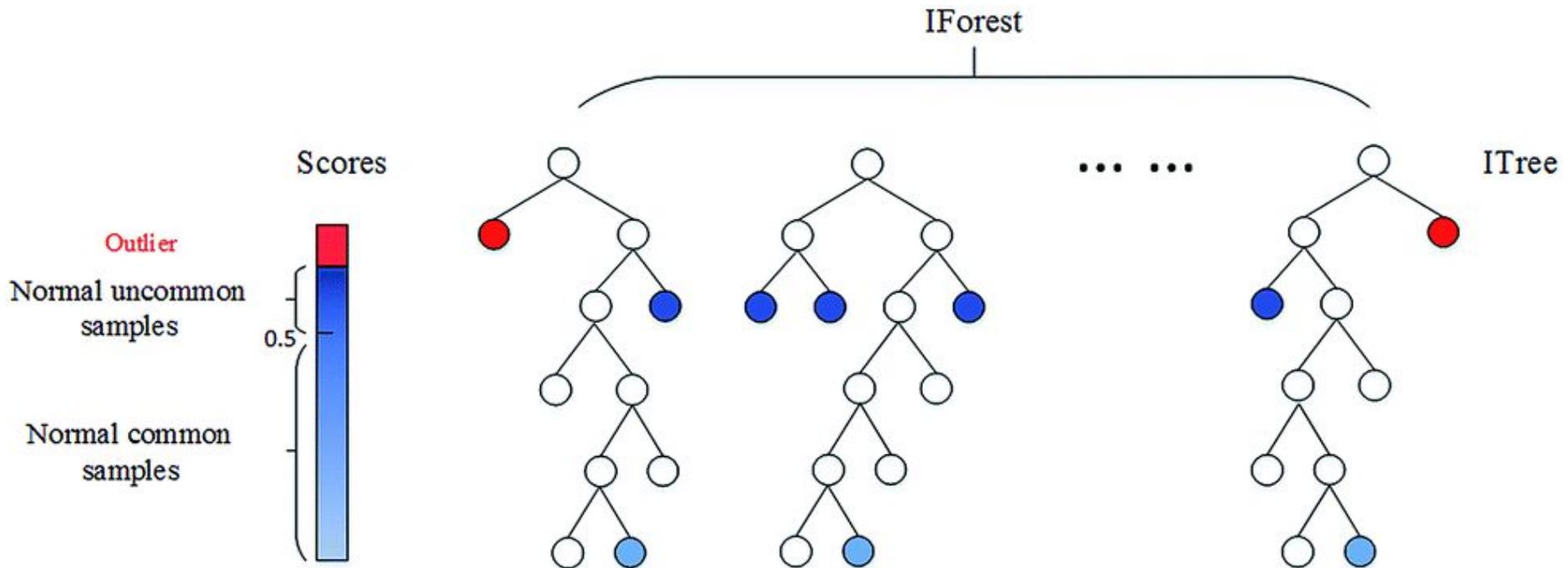
For more on unsupervised methods see Alberto Krone-Martins' talk on Tuesday and Dalya Baron on Wednesday!

Hawkins, 1980

Isolation Forest



Isolation Forest



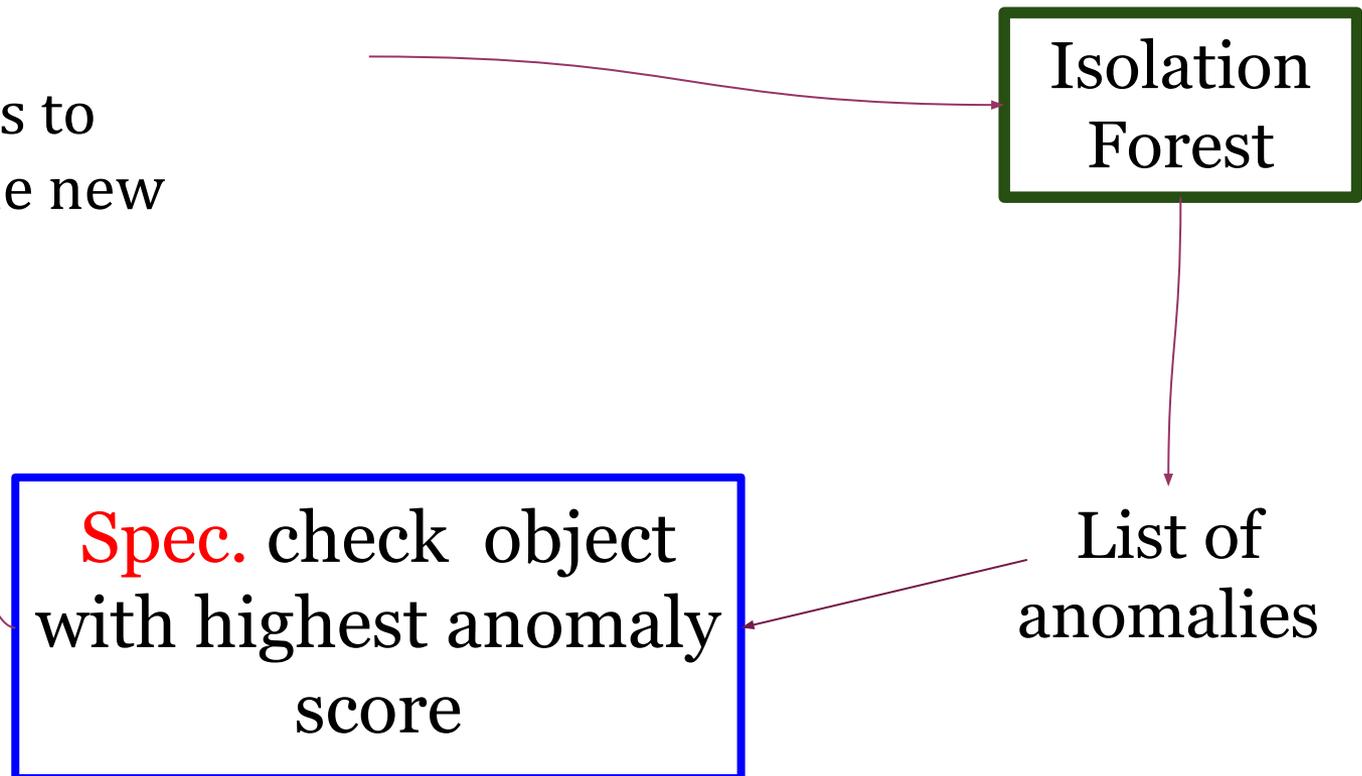
Problem: high occurrence of false positives!

Active Anomaly Detection

A strategy

If yes: check next obj in the anomaly score board

If no: update hyperparameters to accommodate the new information



Does this solve the
problem completely?

No, it is just the best you can do!

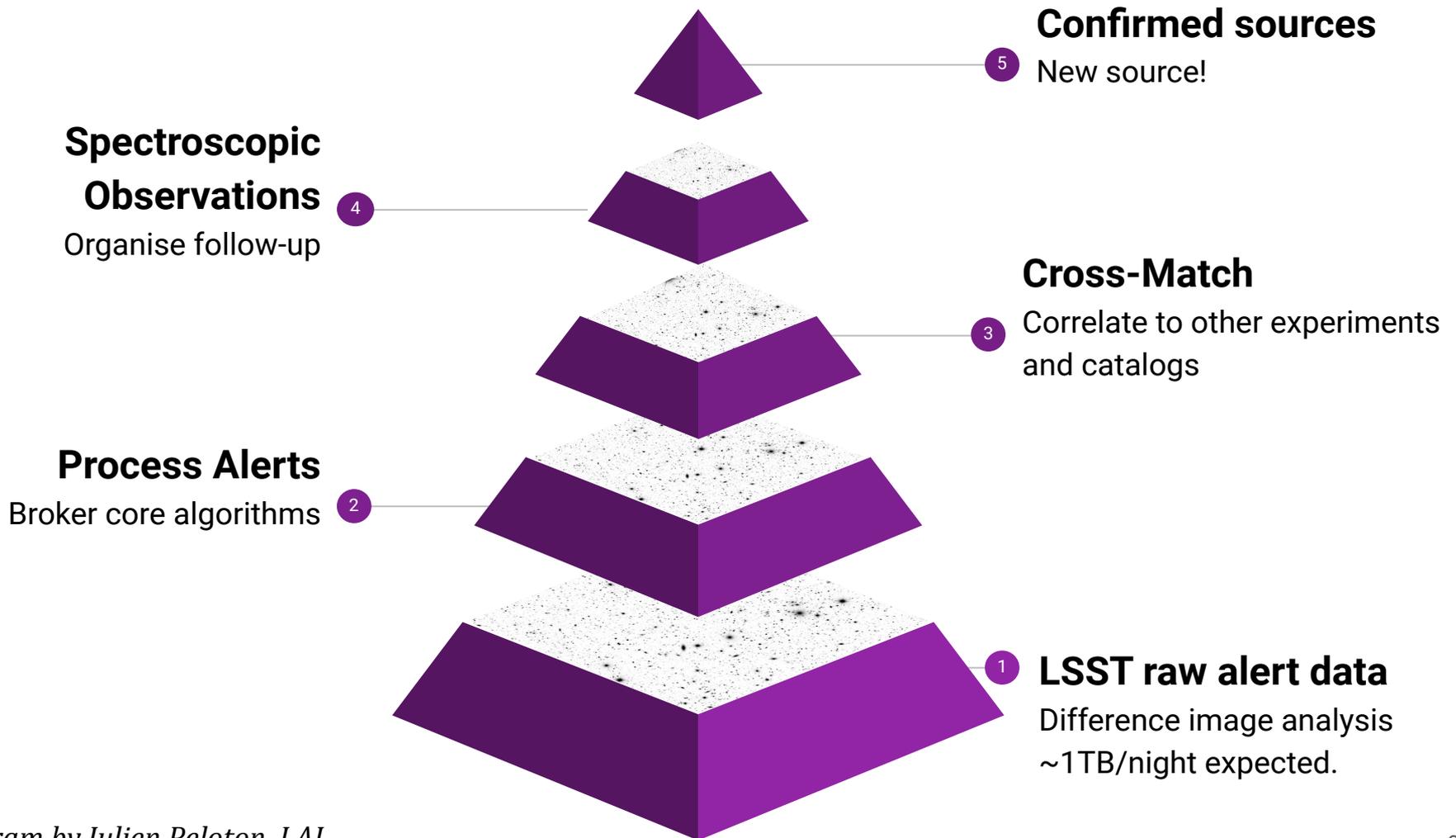
Does this solve the
problem completely?

No, it is just the best you can do!

Is this adaptable to the
upcoming generation of large
scale surveys?

We have to check!

The LSST alert stream



What comes next?

Fink: a community broker based on Active Learning, BNN and Spark

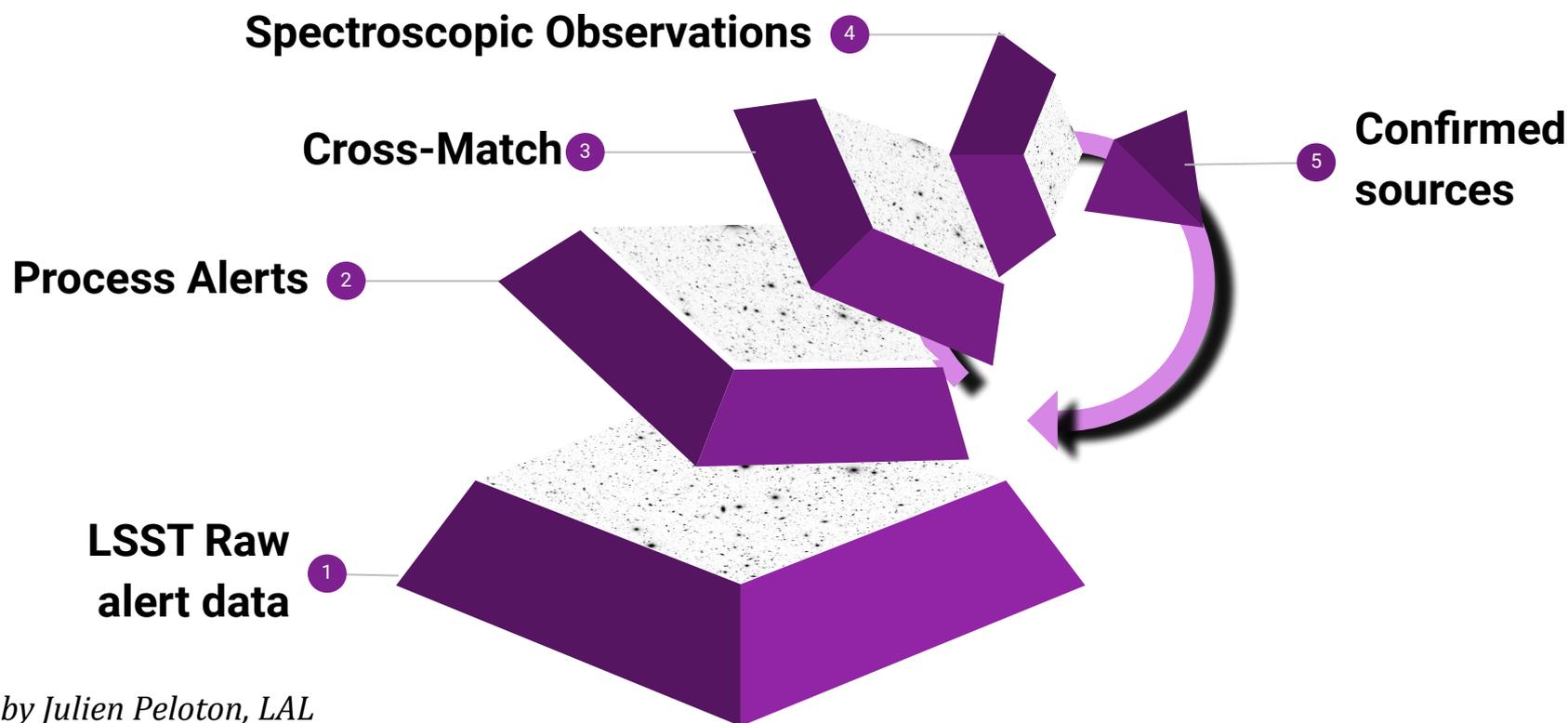


Diagram by Julien Peloton, LAL

<https://fink-broker.readthedocs.io/en/latest/>

Take home message 2:

Serendipitous discoveries will only get more difficult with the next generation of large scale surveys

This means

We need to plan for the unknown

Adaptable algorithms are one possible way to systematically search for new physics - we should think of/try others

THANK
YOU



Extra slides

The queried sample

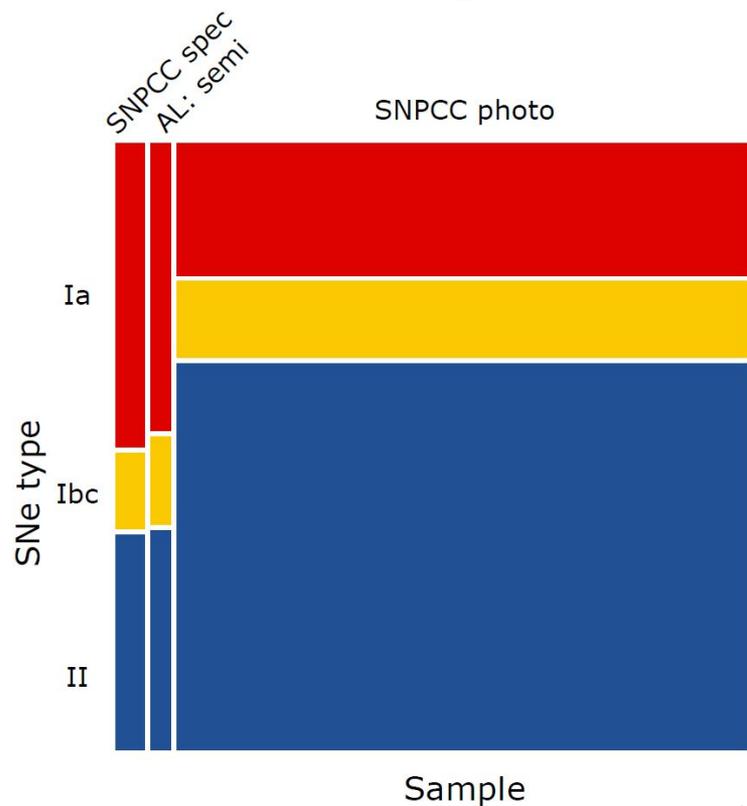
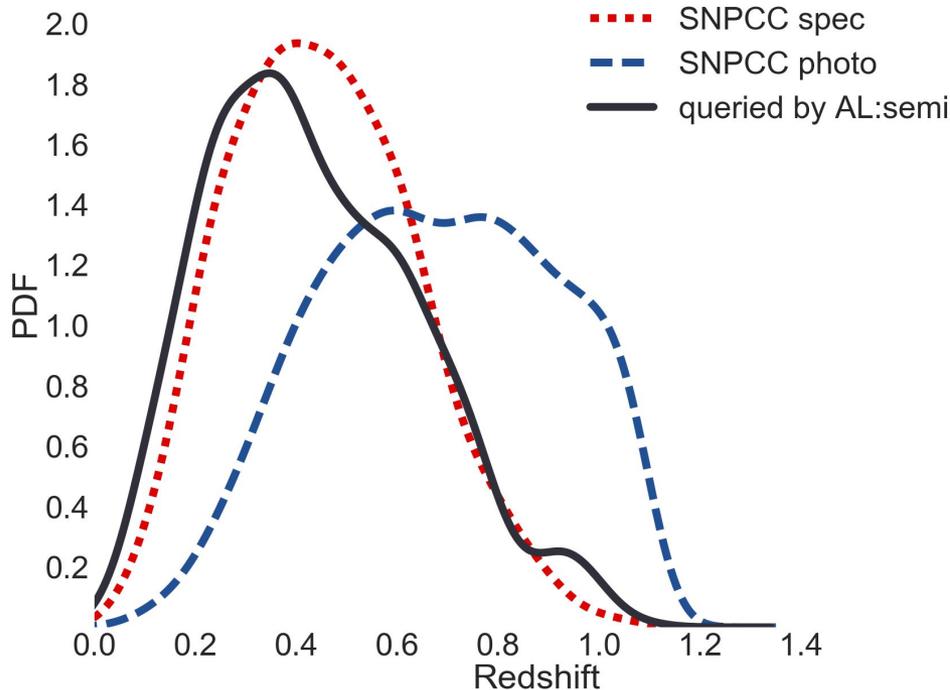
Partial LC, no training, time domain, batch

SNPCC spec:
1103 objects

Queried sample:
800 objects

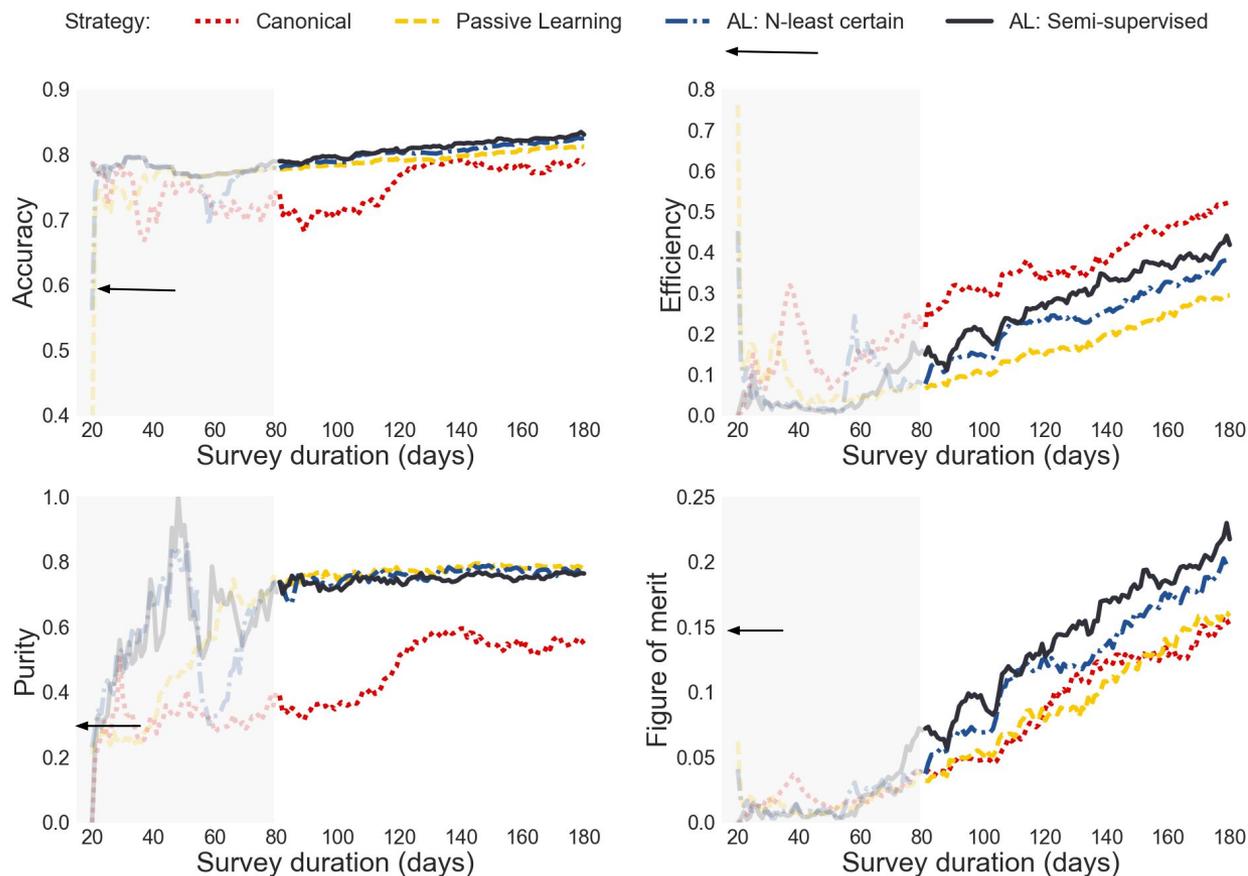
Telescope time:
Queried/spec =

0.999



Batch Mode

Partial LC, no initial training, time domain



The arrow shows traditional Full light-curve results with full SNPCC spec

Happy catalogue

*The effect of coverage +
photometric errors*

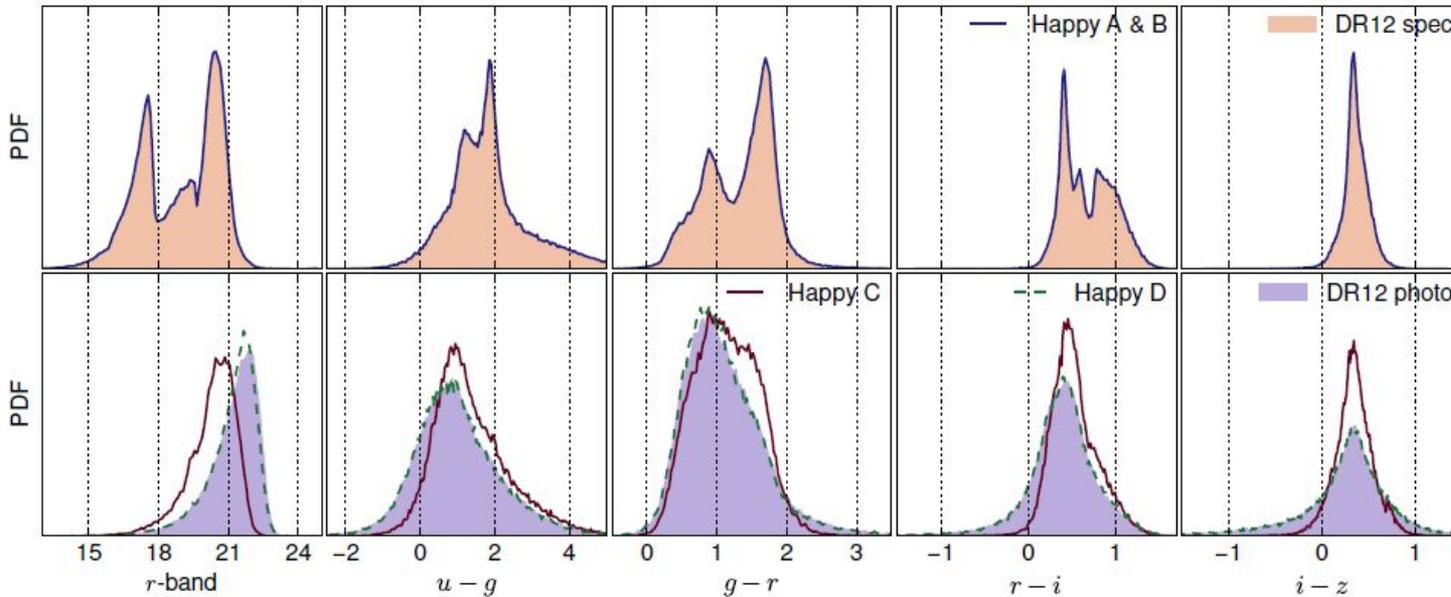
Photometric redshifts

Photometry from SDSS

Spec-z from many different surveys leads to larger
photometric errors and consequently wide domain in r-band and color

From CRP#3 - Beck et al., astro-ph:1701.08748, MNRAS 2017

Each sample has
~ 75000 lines
5 features + errors



Sample A → training

Sample B → target sample, ideal case

Sample C → target sample, data quality similar to A

Sample D → target sample, realistic case

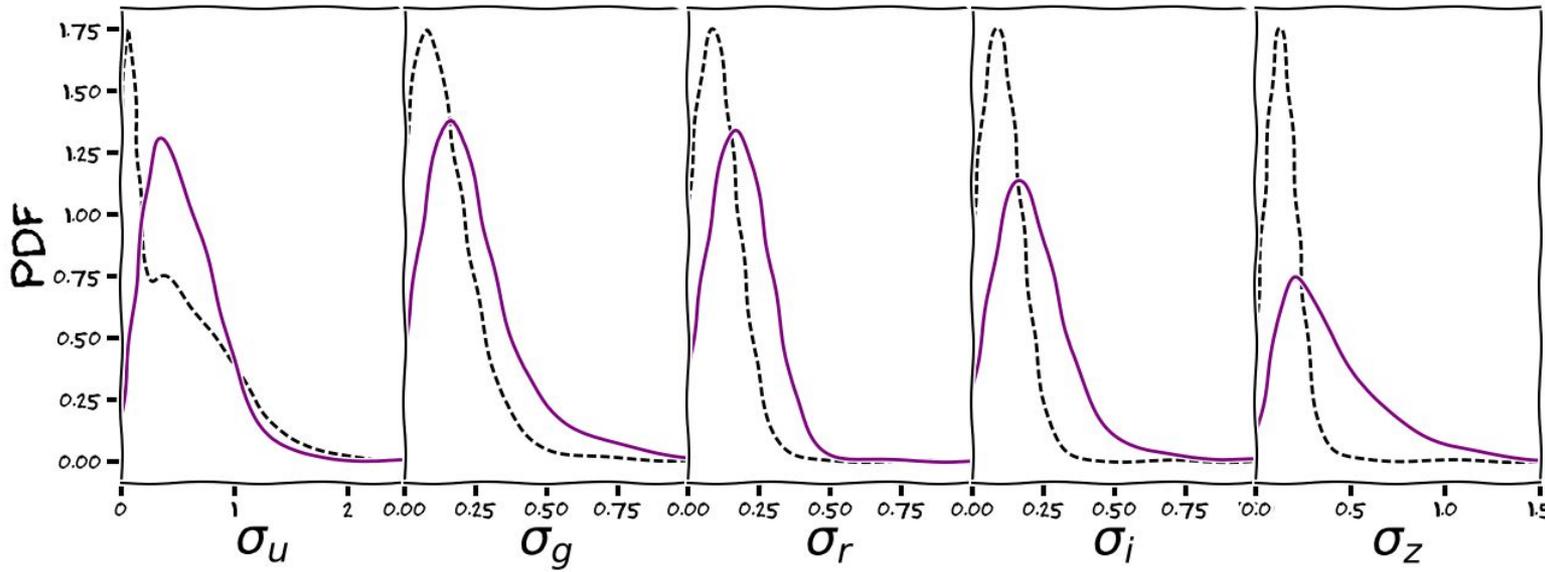
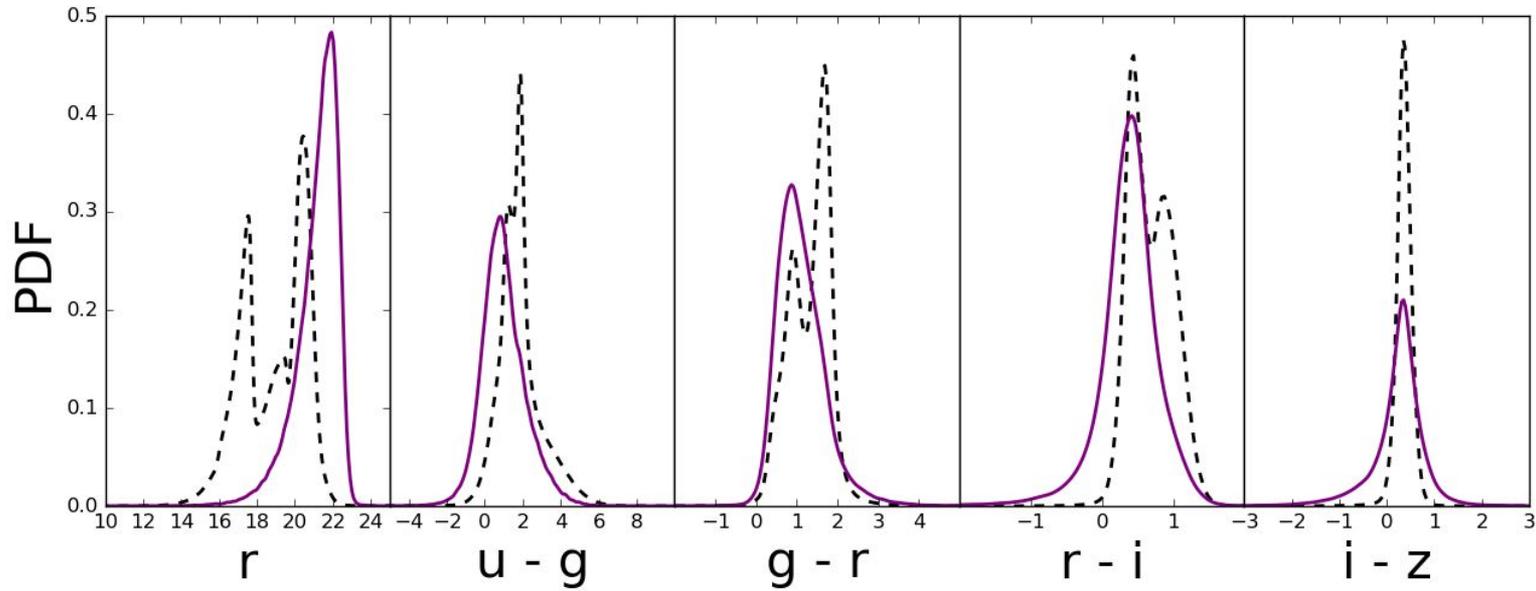
Larger errors when compared to A

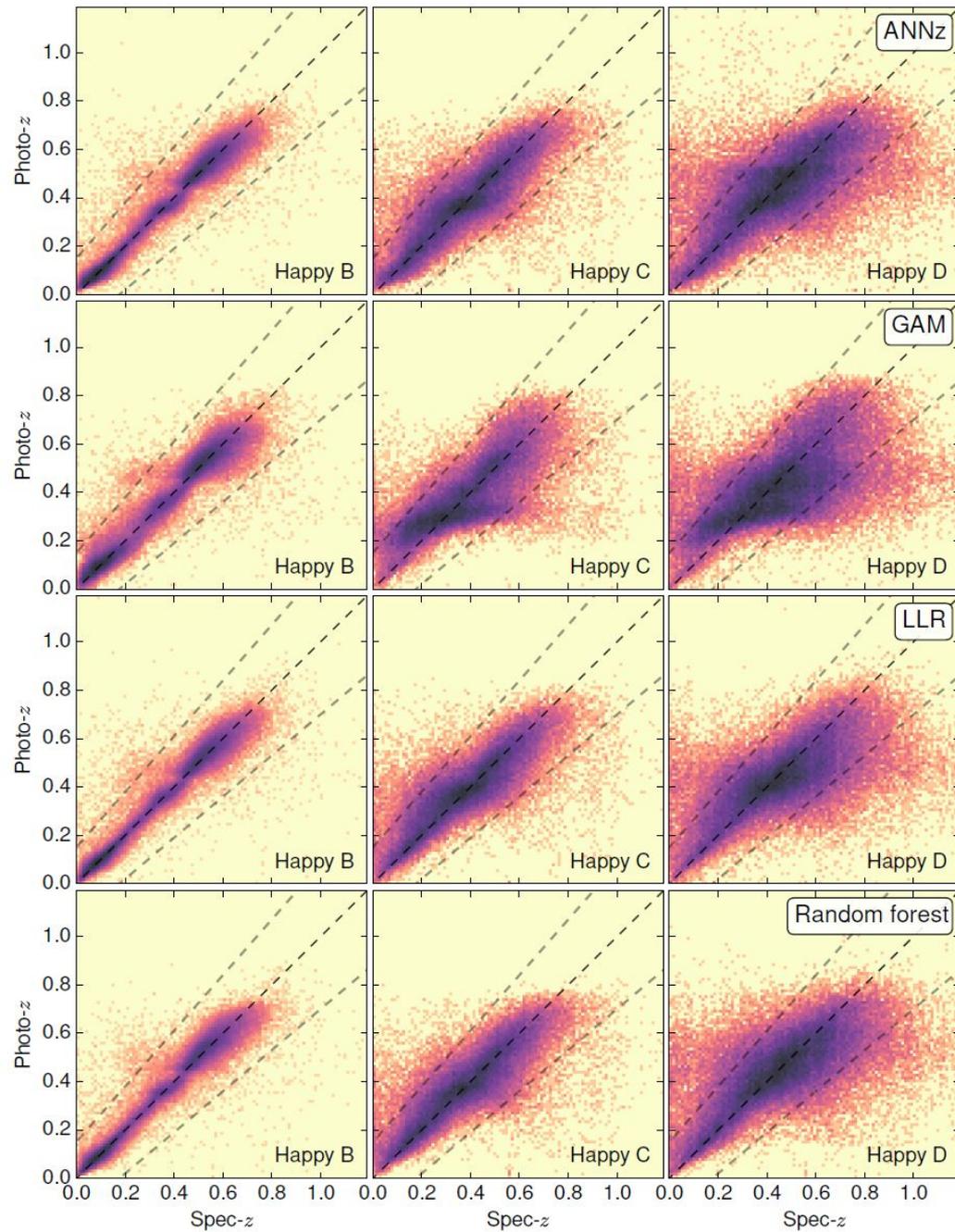
Happy



Redshifts \rightarrow The feature space

----- Training (spec)
— Target (photo)





Happy catalogue

*The effect of coverage +
photometric errors*

Beck et al., astro-ph:1701.08748, MNRAS 2017

