

# Experimental design

---

Alex Malz  
LINCC@CMU

DSFP Session 16



# Overview

The purpose of this lecture is to introduce a sort of metascience:

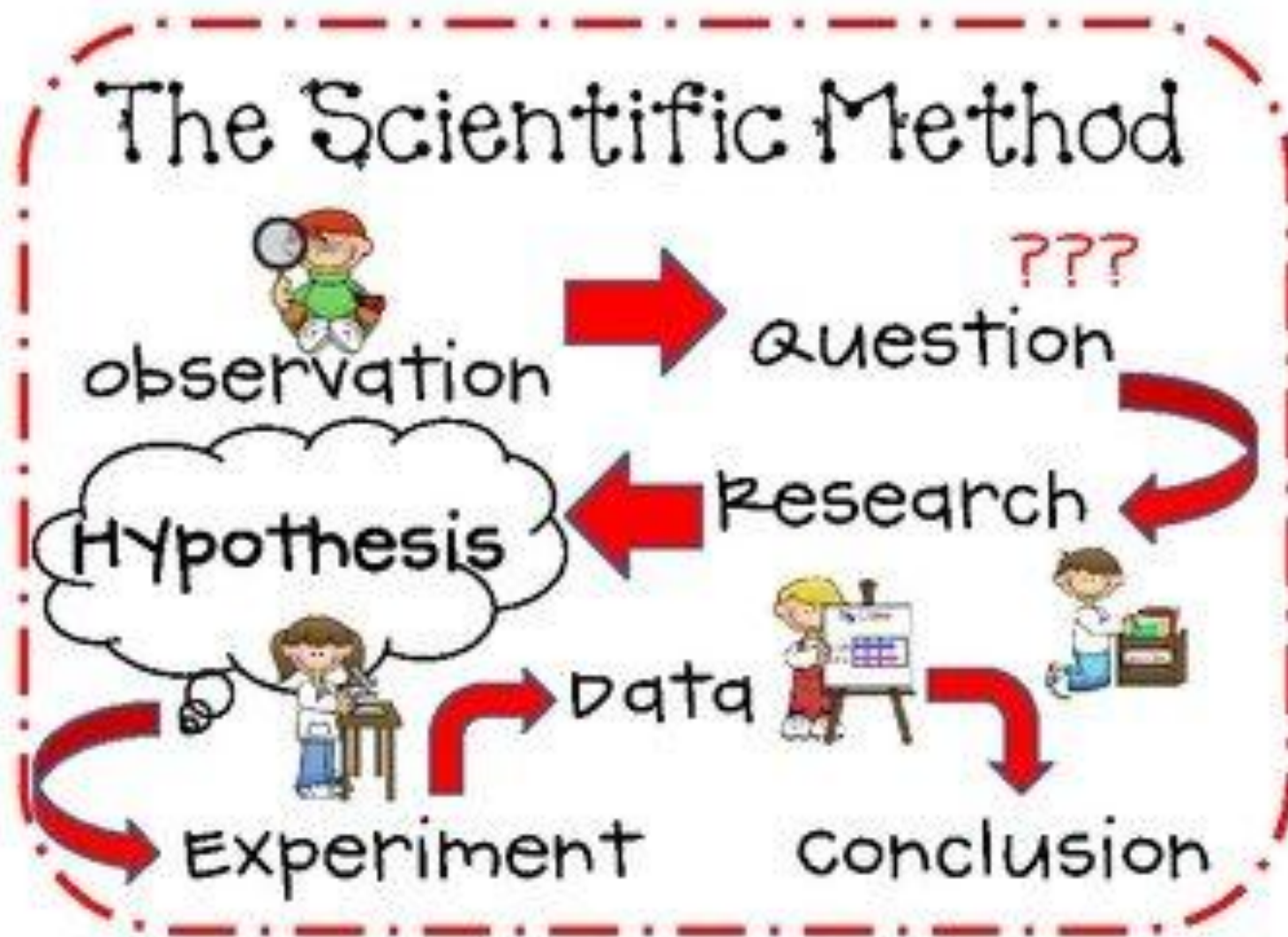
**How do we apply the scientific method to how we do science?**

Probabilistic graphical models are pretty much my favorite thing!

But my secret agenda is for you to incorporate this way of thinking into your hacks, for the experiments in which you prove how great your hierarchical models are.

# What does “experimental design” actually mean?

---



How does observational astronomy  
compare with a physics lab experiment?

Lab physics

vs.

Observational astro

# How does observational astronomy compare with a physics lab experiment?

## Lab physics

vs.

## Observational astro

- Assume physical model  $M$  (and other assumptions  $I$ )
- Seek to constrain physical parameters  $\theta$
- Know controlled experimental conditions  $\phi$
- Make many independent observations  $\{x_i\}$
- ...

Estimate  $p(x | \phi)$   
**likelihood**

- Assume physical model  $M$  (and other assumptions  $I$ )
- Seek to constrain physical parameters  $\theta$
- Learn unknown initial conditions  $\phi$
- Gather data  $x$  from only one observable universe
- ...

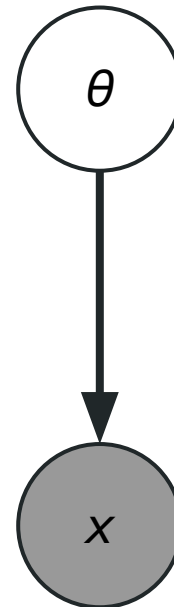
Estimate  $p(\phi | x)$   
**posterior**

# Conditional probability & forward models

The universe has true physical parameters  $\theta$ .

There is a causal relationship  $p(x | \theta)$  between data and the physical parameters; the parameters  $\theta$  determine the data  $x$ .

We observe instances of data  $x$  generated by that model.



# Conditional probability & forward models

The universe has true physical parameters  $\theta$ .

There is a causal relationship  $p(x | \theta)$  between data and the physical parameters; the parameters  $\theta$  determine the data  $x$ .

We observe instances of data  $x$  generated by that model.

We want to constrain the physical parameters  $\theta$  that determined the observed data  $x$ , i.e.  $p(\theta | x)$ .



# Review

---

Let's be critical of  
what we saw  
ereyesterday!



# What do these metrics miss?

Intrinsic scatter

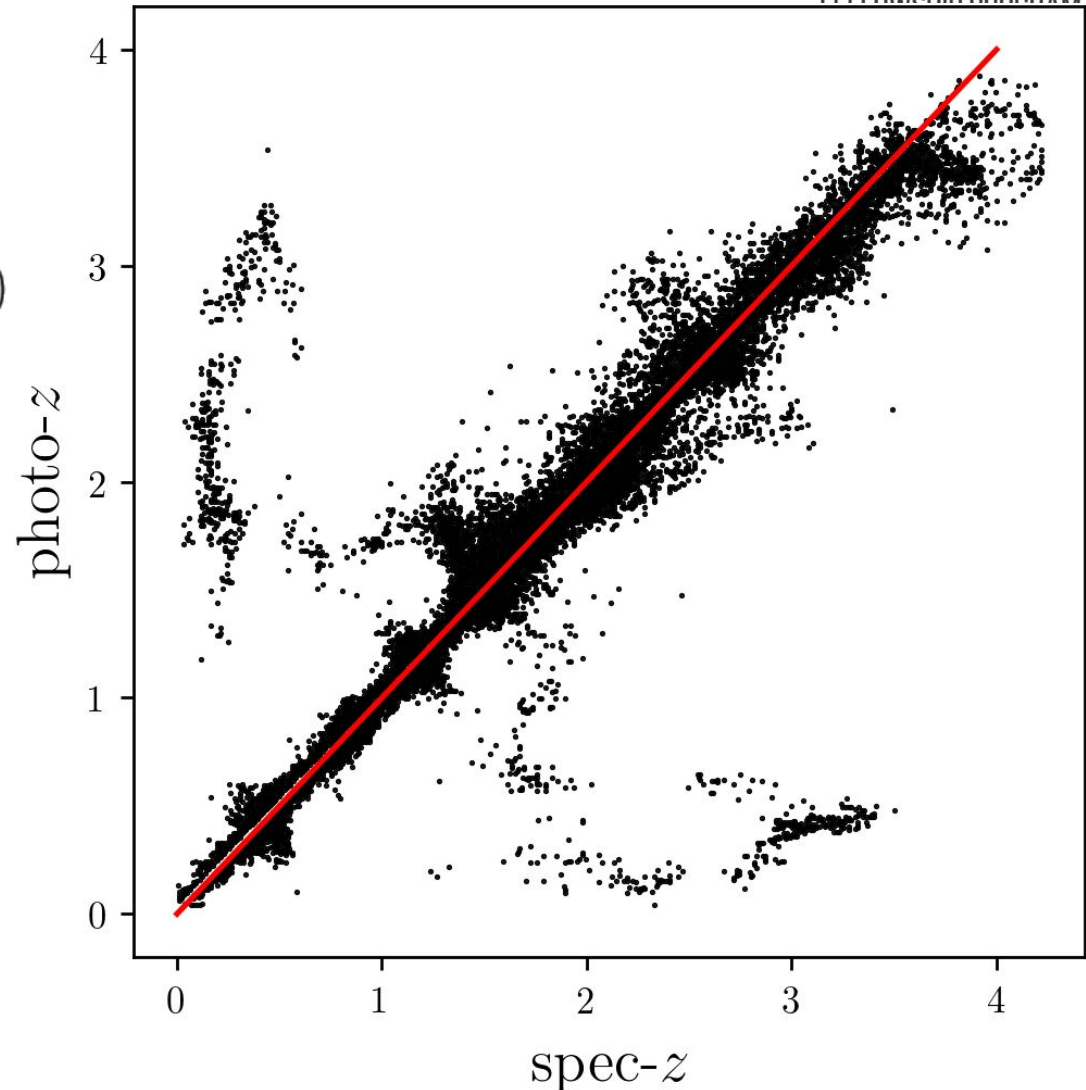
$$\sigma_z < 0.02(1 + z)$$

Bias

$$\langle |z - \hat{z}| \rangle < 0.003(1 + z)$$

Catastrophic outlier rate

$$N_{|z - \hat{z}| > 3\sigma_z} < 0.1 N_{\text{LSST}}$$







# The LSST-DESC PZ DC1 experiment



Motivation: identify the best  
photo-z posterior code  
for LSST-DESC

Data: cosmological redshifts &  
photometry catalog painted  
on N-body simulation

Control: idealized, shared  
prior information



# The LSST-DESC PZ DC1 experiment

Motivation: identify the best  
photo-z posterior code  
for LSST-DESC

Data: cosmological redshifts &  
photometry catalog painted  
on N-body simulation

Control: idealized, shared  
prior information



# Quantitative metrics of 1D PDF ensembles

Root-mean-square Error (RMSE)

$$\text{RMSE} = \sqrt{\int (p_{\text{true}}(z) - \hat{p}_{\text{est}}(z))^2 dz}$$

Kullback-Leibler Divergence (KLD)

$$\text{KLD}[\hat{p}_{\text{est}}(z); p_{\text{true}}(z)] = \int_{-\infty}^{\infty} p_{\text{true}}(z) \log \left[ \frac{p_{\text{true}}(z)}{\hat{p}_{\text{est}}(z)} \right] dz$$

Cumulative Distribution Function (CDF)

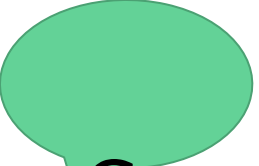
$$\text{CDF}[\hat{p}, z'] \equiv \int_{-\infty}^{z'} \hat{p}(z) dz$$

Probability Integral Transform (PIT)

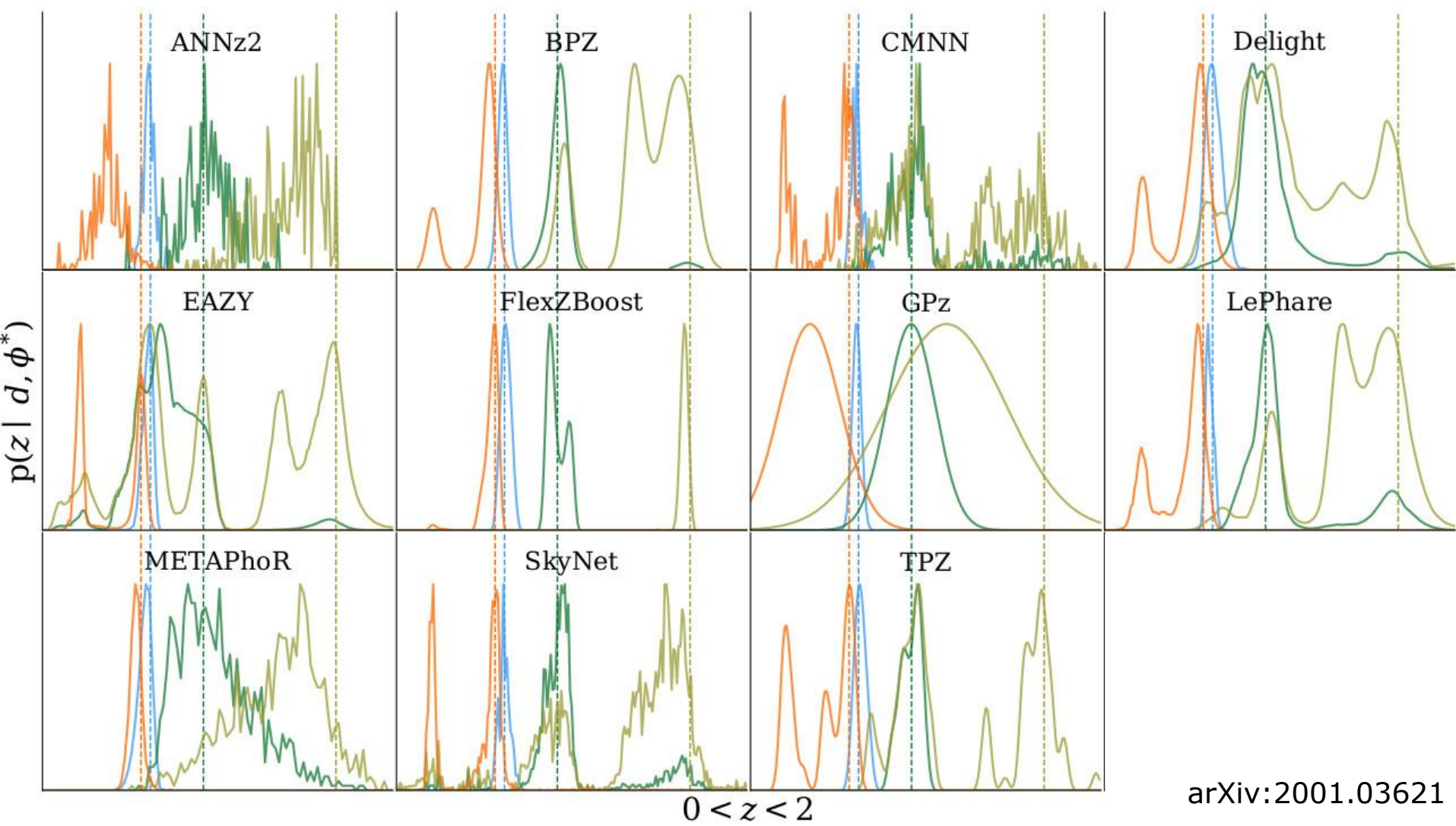
$$P(\text{PIT} \equiv \text{CDF}[\hat{p}, z_{\text{true}}])$$

Quantile-quantile (QQ) Plot  $\sim \int p(\text{PIT}) d\text{PIT}$

**No true  
posteriors  
available!**

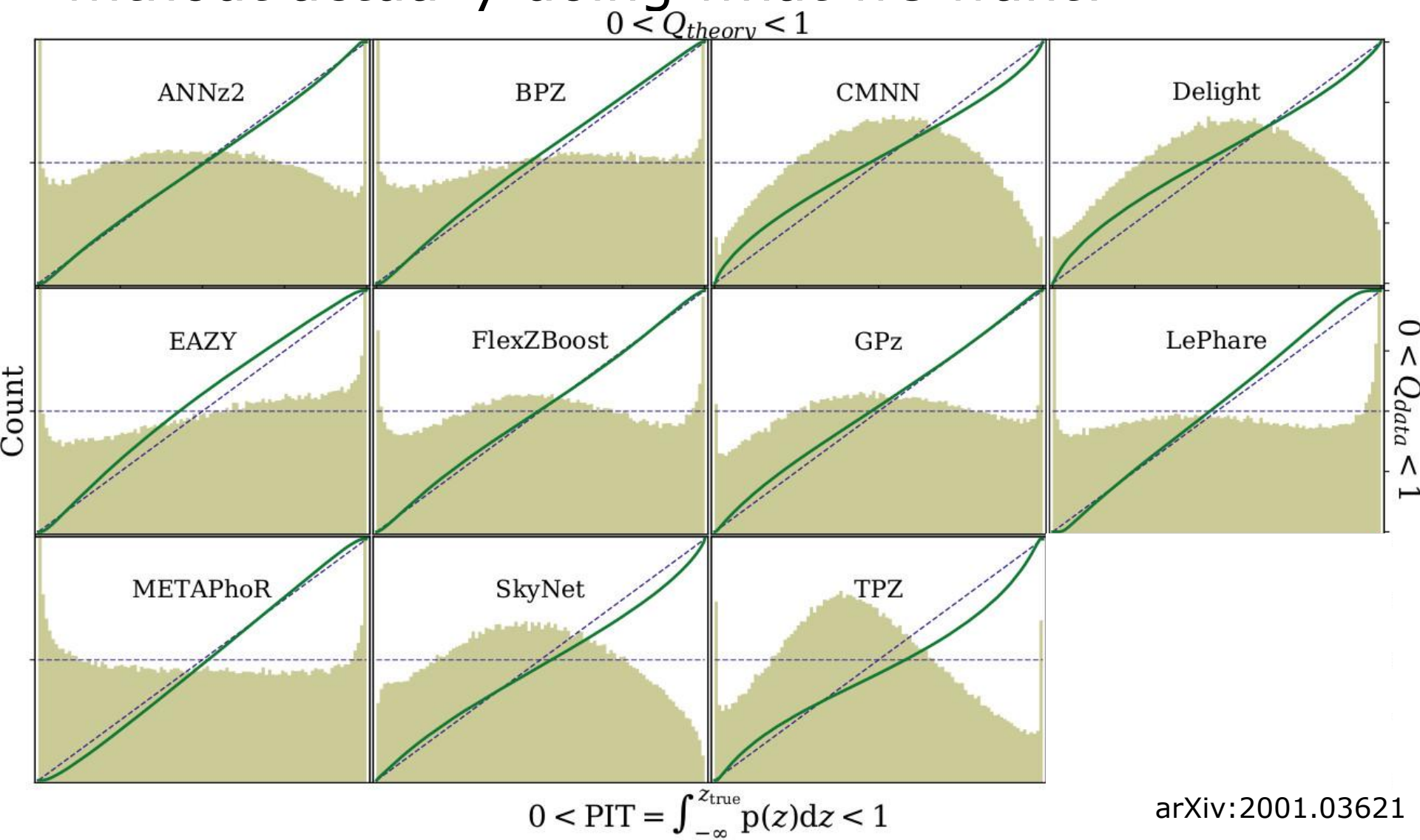


Can we think of an estimation model. . .

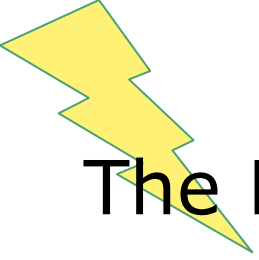




... that does well by the metrics,  
without actually doing what we want?







# The LSST-DESC PZ DC1 experiment

**Wait, what do we want again?**

Motivation: identify the best photo-z posterior code for LSST-DESC

Data: cosmological redshifts & photometry catalog painted on N-body simulation

Control: idealized, shared prior information

# The LSST-DESC PZ DC1 experiment

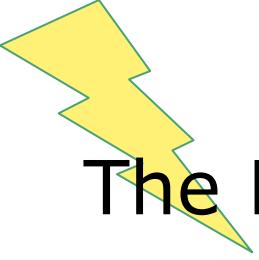
**Wait, what do we want again?**

**“The best” is whatever the metric says it is!**

Motivation: identify the best photo-z posterior code for LSST-DESC

Data: cosmological redshifts & photometry catalog painted on N-body simulation

Control: idealized, shared prior information



# The LSST-DESC PZ DC1 experiment

**Wait, what do we want again?**

Motivation: identify the best photo-z posterior code for LSST-DESC

**“The best” is whatever the metric says it is!**

Data: cosmological redshifts & photometry catalog painted on N-body simulation

**But what do we really want?**

Control: idealized, shared prior information

# The LSST-DESC PZ DC1 experiment

**Wait, what do we want again?**

Motivation: identify the best photo-z posterior code for LSST-DESC

**"The best" is whatever the metric says it is!**

Data: cosmological redshifts & photometry catalog painted on N-body simulation

**But what do we really want?**

Control: idealized, shared prior information

**Here, we want an estimator that extracts redshift information from photometry.**



# Quantitative metrics of 1D PDF ensembles

Root-mean-square Error (RMSE)

$$\text{RMSE} = \sqrt{\int (p_{\text{true}}(z) - \hat{p}_{\text{est}}(z))^2 dz}$$

Kullback-Leibler Divergence (KLD)

$$\text{KLD}[\hat{p}_{\text{est}}(z); p_{\text{true}}(z)] = \int_{-\infty}^{\infty} p_{\text{true}}(z) \log \left[ \frac{p_{\text{true}}(z)}{\hat{p}_{\text{est}}(z)} \right] dz$$

Cumulative Distribution Function (CDF)

$$\text{CDF}[\hat{p}, z'] \equiv \int_{-\infty}^{z'} \hat{p}(z) dz$$

Probability Integral Transform (PIT)

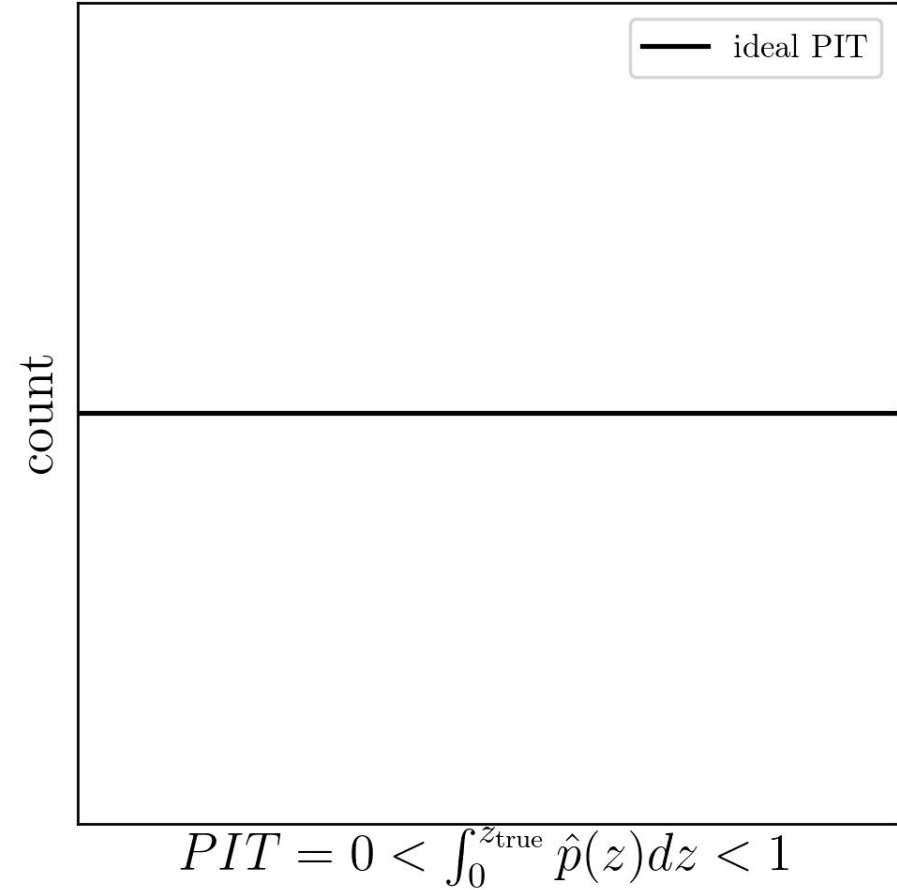
$$P(\text{PIT} \equiv \text{CDF}[\hat{p}, z_{\text{true}}])$$

Quantile-quantile (QQ) Plot  $\sim \int p(\text{PIT}) d\text{PIT}$

**No true  
posteriors  
available!**



Sketch the PIT histogram of an ideal photo-z posterior estimator.



$$\text{CDF}[\hat{p}, z'] \equiv \int_{-\infty}^{z'} \hat{p}(z) dz$$

$$P(\text{PIT} \equiv \text{CDF}[\hat{p}, z_{\text{true}}])$$

# Is the metric inappropriate?

Think *adversarially* to trick the metric into rewarding a wrong answer.

How can you satisfy the metric's criteria without satisfying the criteria you as an experimenter actually care about?

Identify *pathological* cases that fool the metric but are obviously wrong.

Can you think of a null test in the form of a trivially bad answer that performs well by the metric?

# The worst photo- $z$ posterior estimator



# The worst photo- $z$ posterior estimator

1. Make a histogram of training set galaxy redshifts.

# The worst photo- $z$ posterior estimator

1. Make a histogram of training set galaxy redshifts.
2. Read in the test set photometry.

# The worst photo- $z$ posterior estimator

1. Make a histogram of training set galaxy redshifts.
2. Read in the test set photometry.
3. Count the number of test set galaxies.

# The worst photo- $z$ posterior estimator

1. Make a histogram of training set galaxy redshifts.
2. Read in the test set photometry.
3. Count the number of test set galaxies.
4. Discard the test set photometry.

# The worst photo- $z$ posterior estimator

1. Make a histogram of training set galaxy redshifts.
2. Read in the test set photometry.
3. Count the number of test set galaxies.
4. Discard the test set photometry.
5. Return the histogram of training set galaxy redshifts as the photo- $z$  posterior of every test set galaxy.

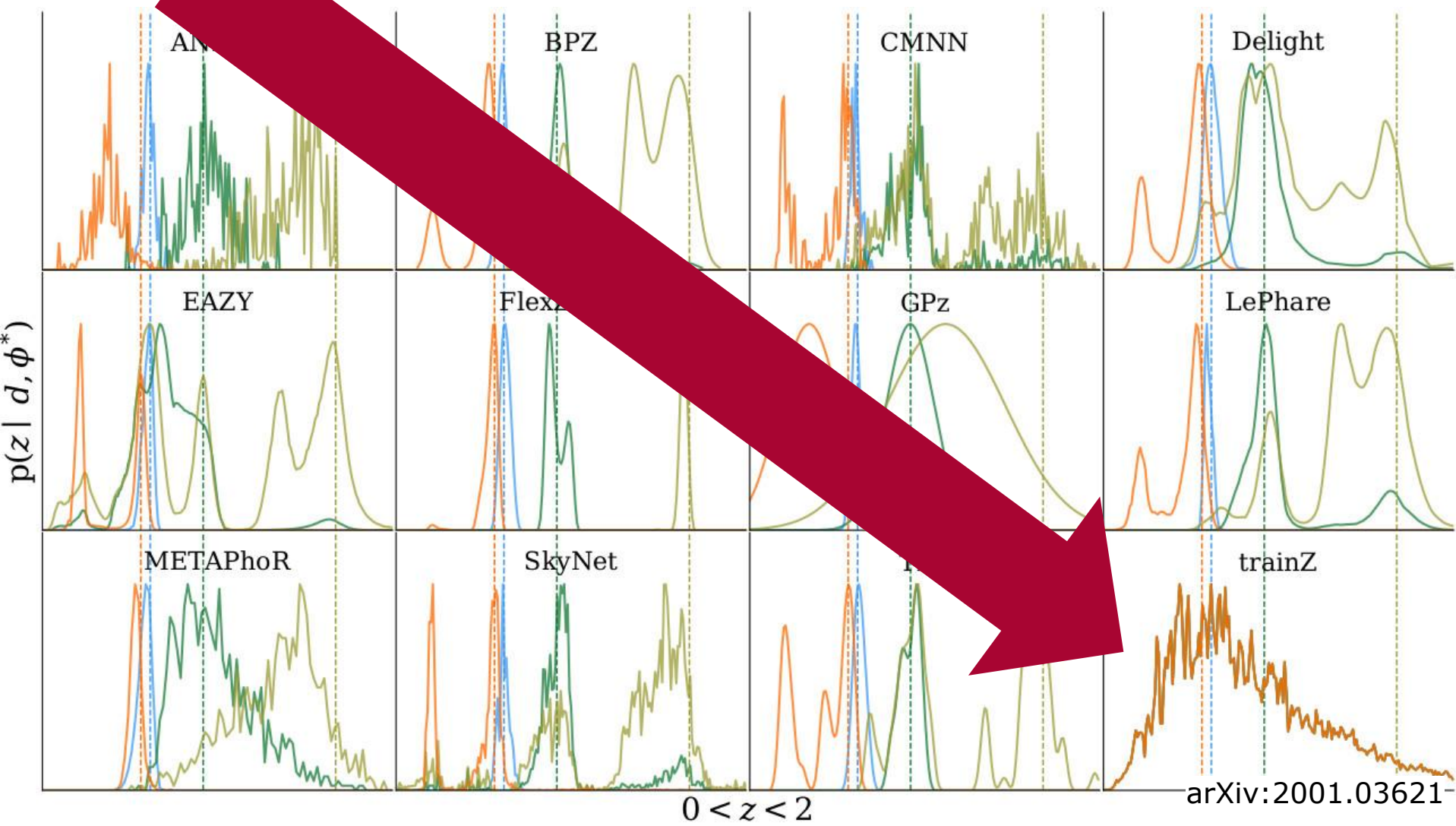
# The worst photo- $z$ posterior estimator

1. Make a histogram of training set galaxy redshifts.
2. Read in the test set photometry.
3. Count the number of test set galaxies.
4. Discard the test set photometry.
5. Return the histogram of training set galaxy redshifts as the photo- $z$  posterior of every test set galaxy.
6. ...

# The worst photo- $z$ posterior estimator

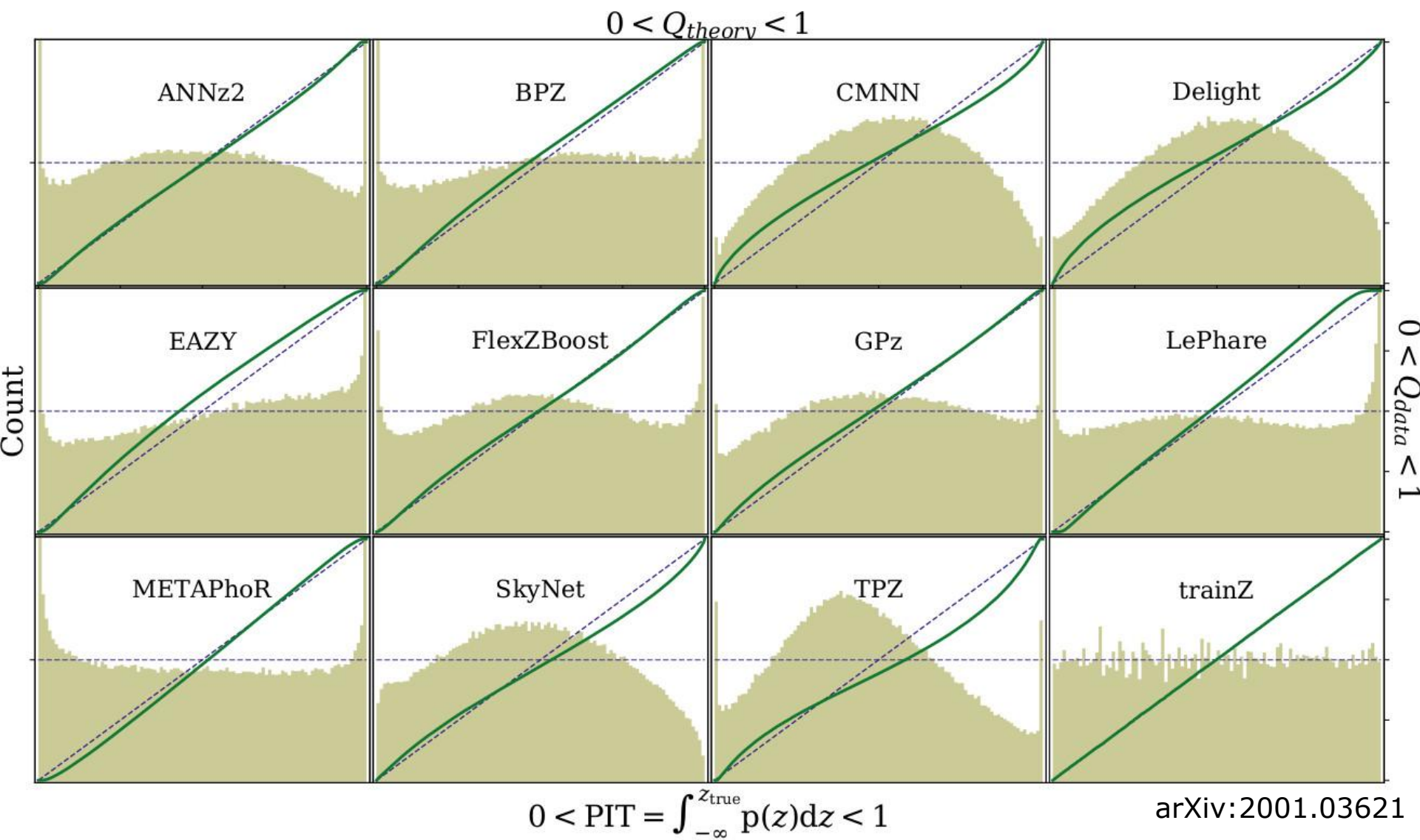
1. Make a histogram of training set galaxy redshifts.
2. Read in the test set photometry.
3. Count the number of test set galaxies.
4. Discard the test set photometry.
5. Return the histogram of training set galaxy redshifts as the photo- $z$  posterior of every test set galaxy.
6. ...
7. Profit!

trainZ is a “control” case





trainZ has nearly perfect PIT & QQ.

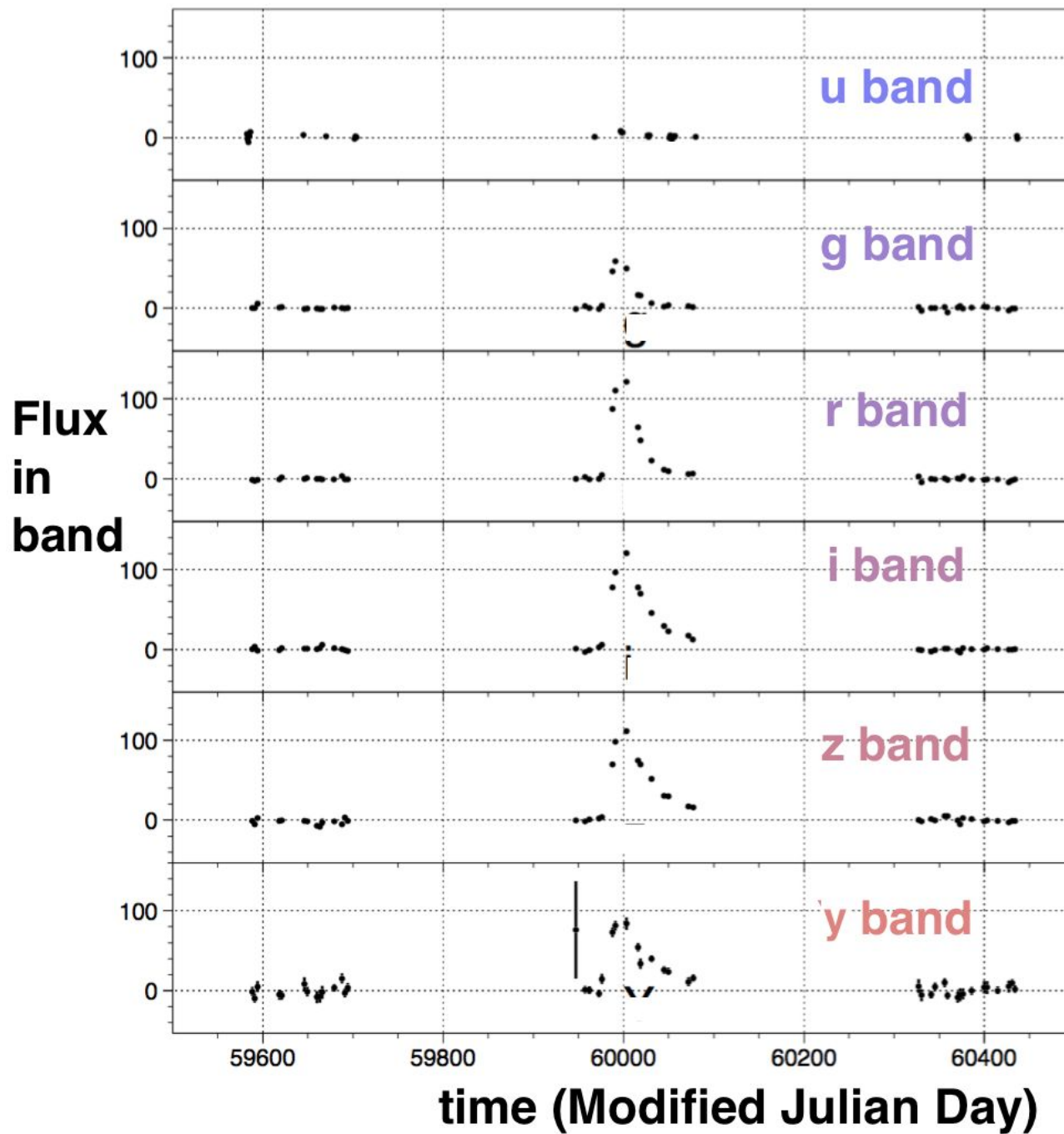


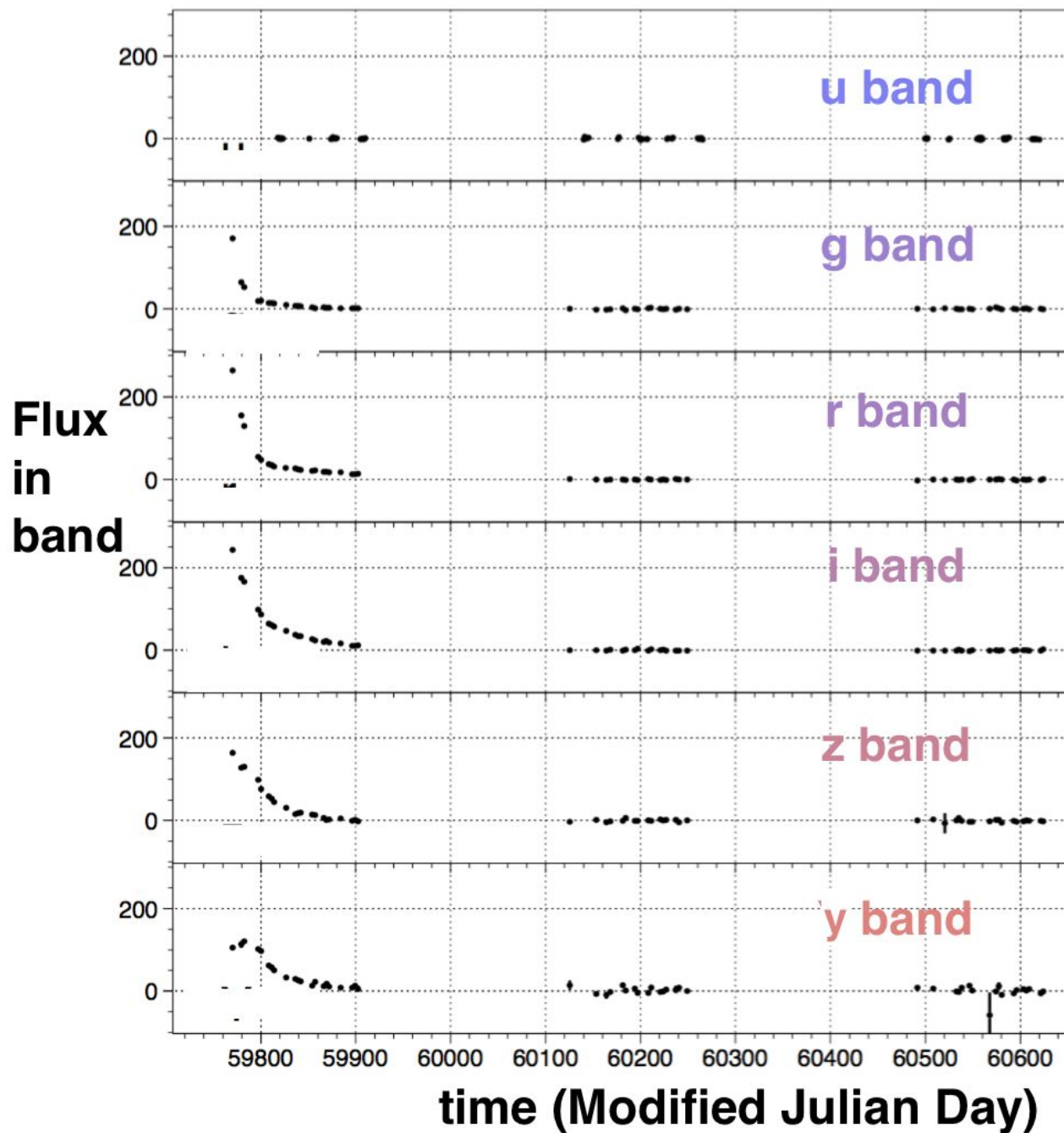
# How does one create a specialized metric for a specific goal?

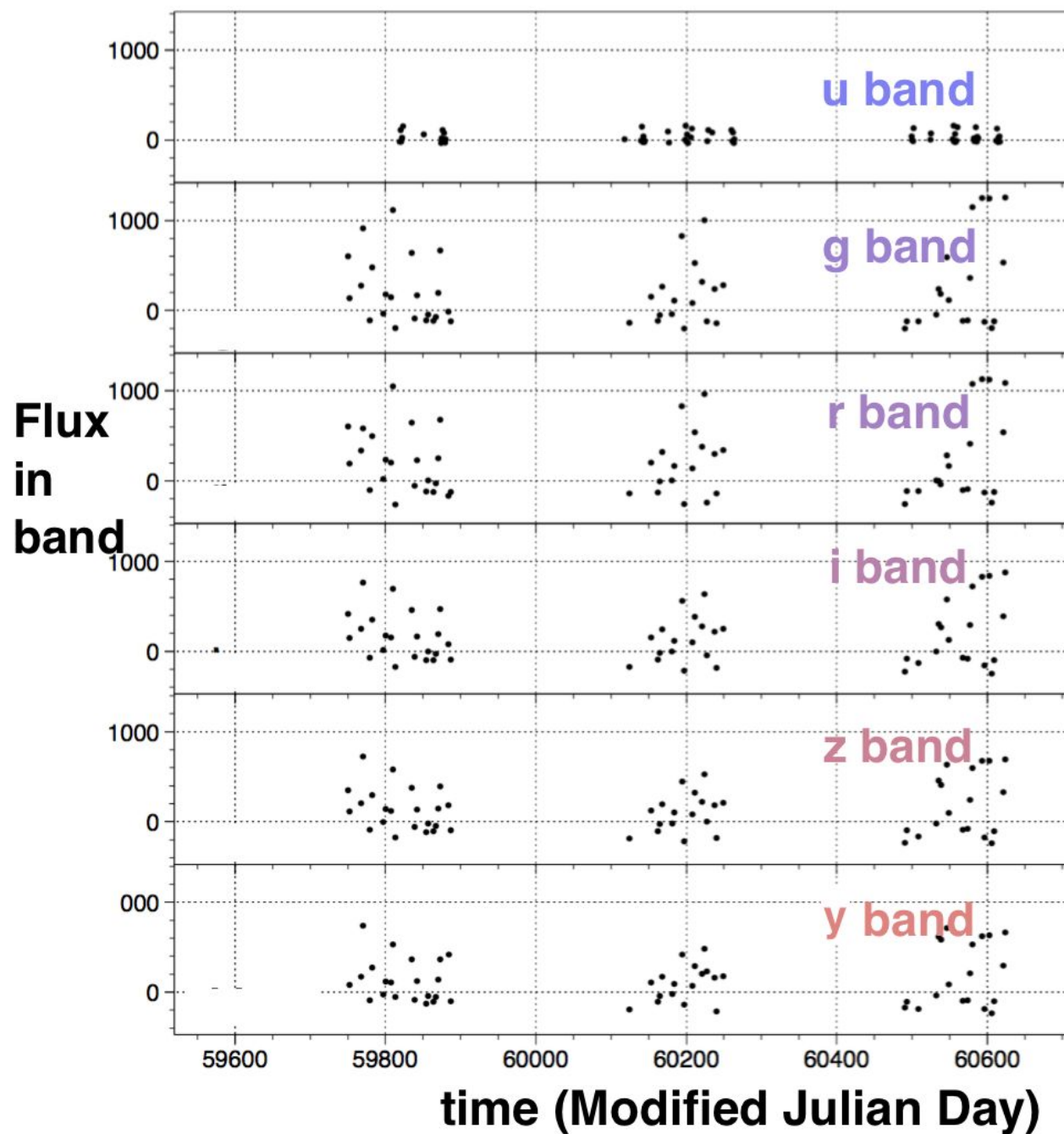
---

Context: The Photometric  
LSST AStronomical  
TIme-series Classification  
Challenge (PLAsTiCC)









# The challenge of the PLAsTiCC metric

## What made PLAsTiCC challenging

- Noisy, incomplete LCs
- Many classes, subclasses
- Diverse science applications

## Why the metric required critical thinking

- Classification PMFs vs. labels
- Kaggle requires single scalar

# Confusion matrix & T/F P/N rates

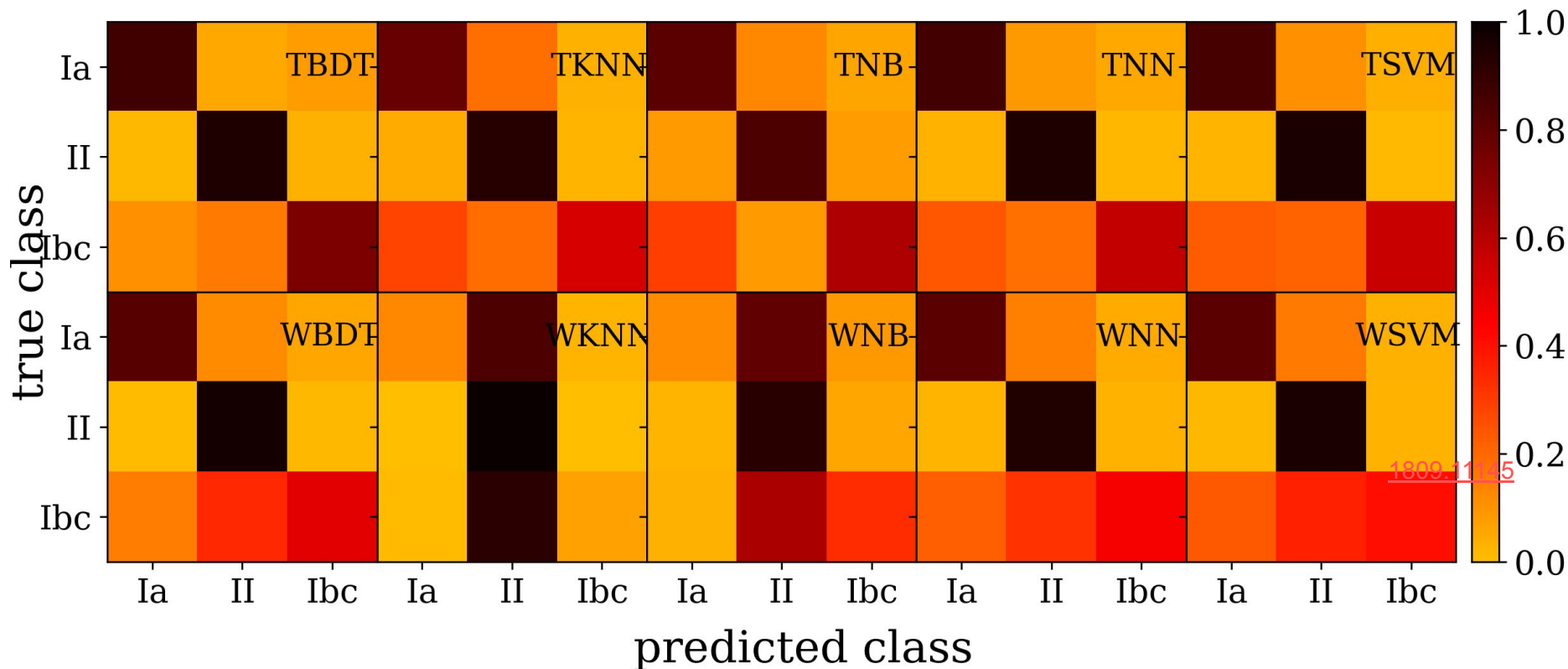
	Actual true	Actual false
Predicted true	True positives (TP)	False positives (FP)
Predicted false	False negatives (FN)	True negatives (TN)

$$\text{Purity} = \text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

$$\text{Efficiency} = \text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$



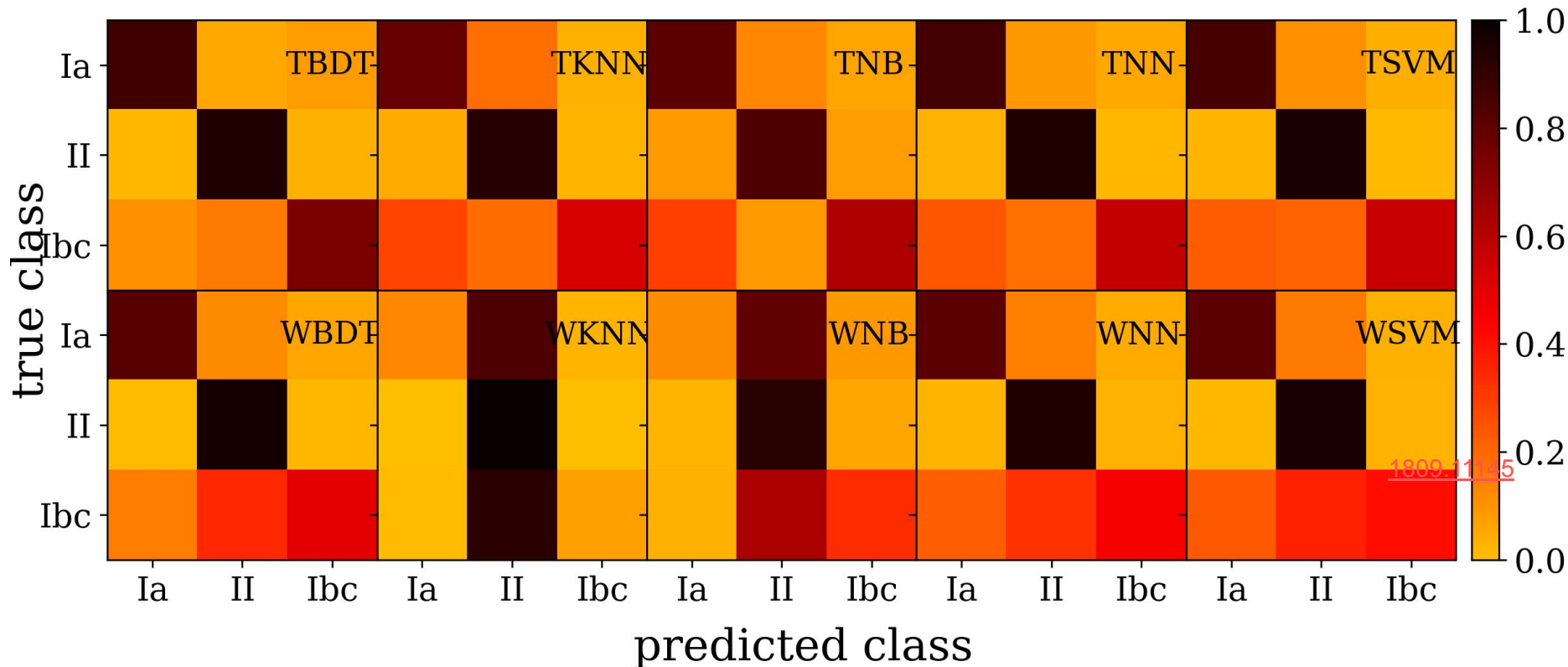
# Confusion matrix for SN classification



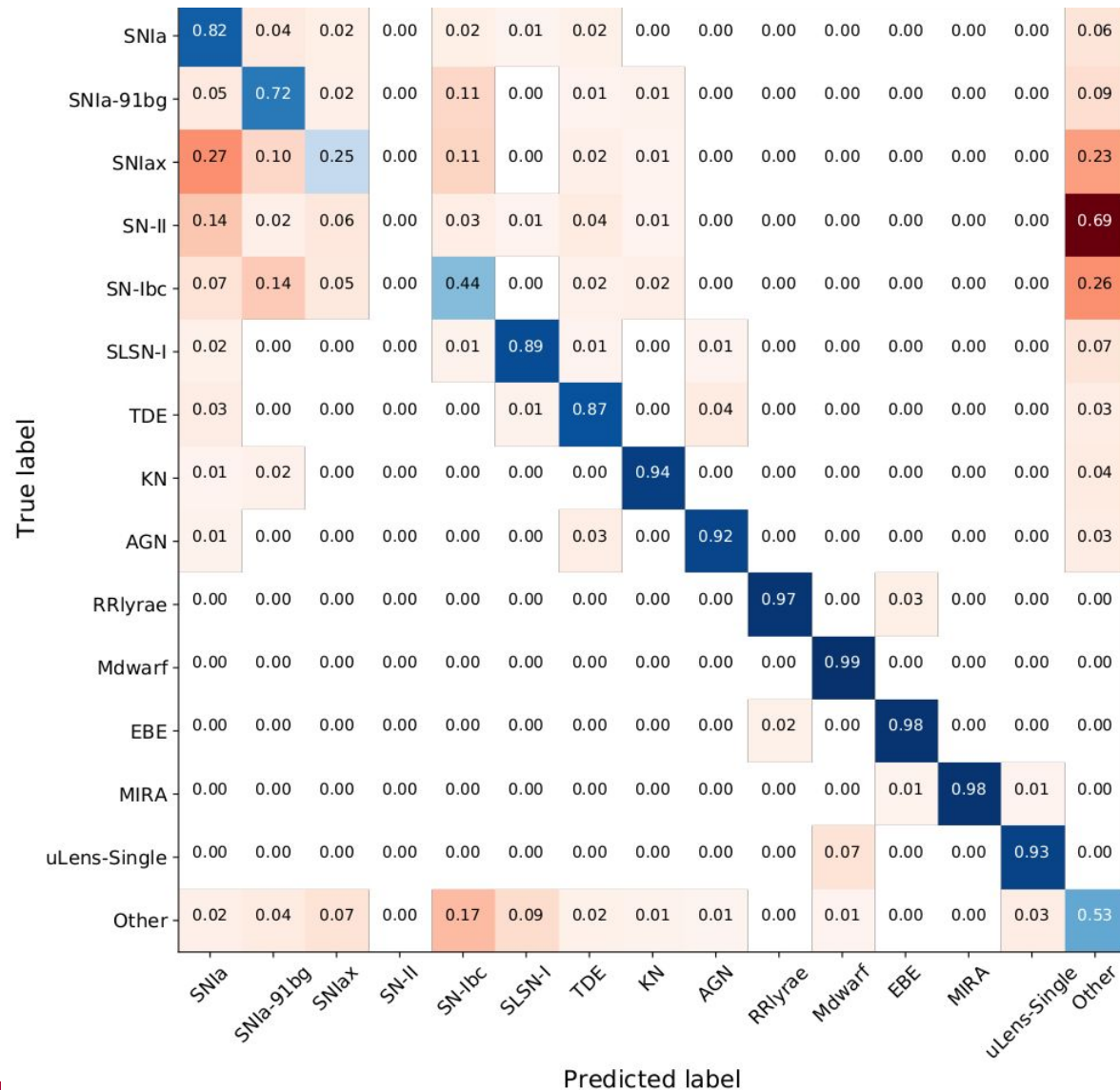





What do we want to reward vs. punish?



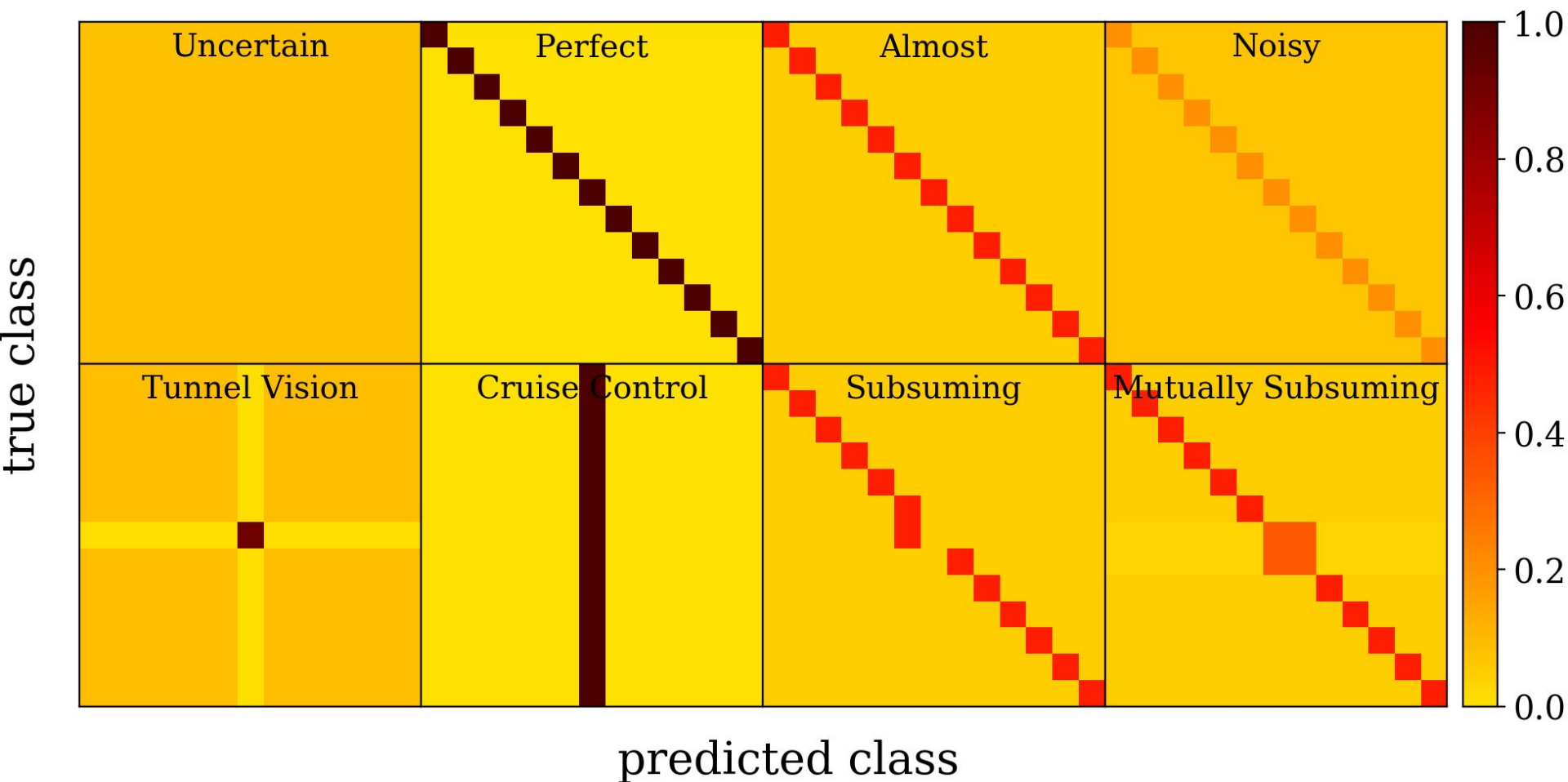
# Many-class confusion matrix





Think adversarially! Devise a pathological classifier and sketch its confusion matrix.

Think adversarially! Devise a pathological classifier and sketch its confusion matrix.



# Quantitative metrics of 1D PDFs

Root-mean-square Error

$$\text{RMSE} = \sqrt{\int \left( P(z) - \hat{P}(z) \right)^2 dz}$$

Kullback-Leibler Divergence

$$\text{KLD} \left[ \hat{P}(z) | P(z) \right] = \int_{-\infty}^{\infty} P(z) \log \left[ \frac{P(z)}{\hat{P}(z)} \right] dz$$

# Quantitative metrics of categorical PMFs

Root mean square Error

**Brier Score**  $B_{LC\ n} \equiv \sum_{\text{class } m=1}^M (\tau_{n,m} - \hat{p}(m|\text{data}_n))^2$

Kullback-Leibler Divergence

**Log-loss**  $L_{LC\ n} \equiv - \sum_{\text{class } m=1}^M \tau_{n,m} \ln[\hat{p}(m|\text{data}_n)]$

What free parameters are chosen by the experimenter?

Root mean square Error

**Brier Score**  $B_{LC\ n} \equiv \sum_{\text{class } m=1}^M (\tau_{n,m} - \hat{p}(m|\text{data}_n))^2$

Kullback-Leibler Divergence

**Log-loss**  $L_{LC\ n} \equiv - \sum_{\text{class } m=1}^M \tau_{n,m} \ln[\hat{p}(m|\text{data}_n)]$

$$Q = \frac{1}{\sum_m W_m} \sum_{m=1}^M W_m \frac{1}{\sum_n w_{n,m}} \sum_{n=1}^N w_{n,m} Q_{n,m}$$

What free parameters are chosen by the experimenter?

Root-mean-square Error

**Brier Score**  $B_{LC\ n} \equiv \sum_{\text{class } m=1}^M (\tau_{n,m} - \hat{p}(m|\text{data}_n))^2$

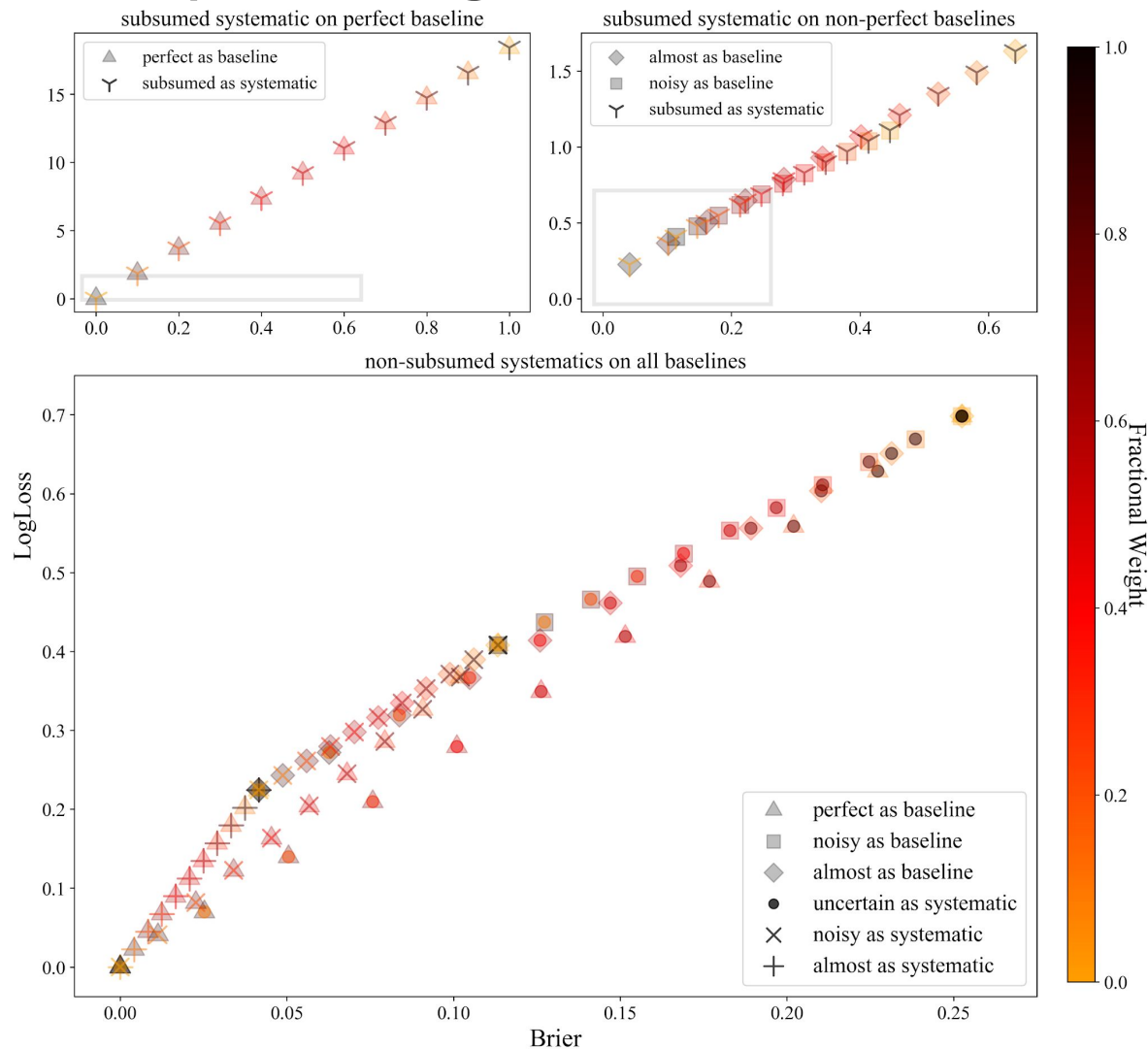
Kullback-Leibler Divergence

**Log-loss**  $L_{LC\ n} \equiv - \sum_{\text{class } m=1}^M \tau_{n,m} \ln[\hat{p}(m|\text{data}_n)]$

$$Q = \frac{1}{\sum_m W_m} \sum_{m=1}^M W_m \frac{1}{\sum_n w_{n,m}} \sum_{n=1}^N w_{n,m} Q_{n,m}$$



tl;dr both metrics considered were robust to pathological classifiers

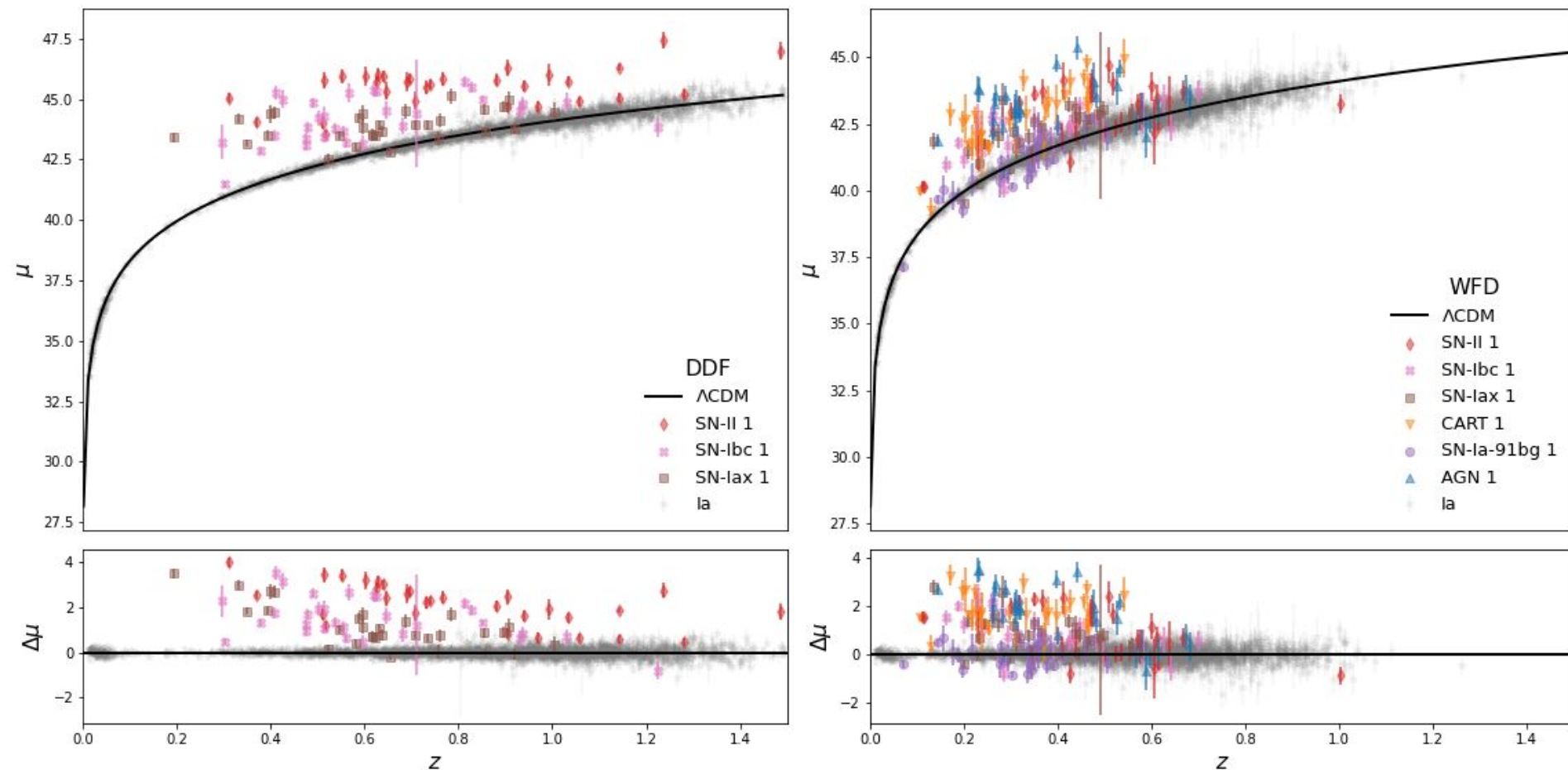


# Another astrophysical example: transient classification for SNIa cosmology

---

How can we choose  
classification metrics  
that reward what we  
really want?

# SN Ia cosmology with contaminants



# SN Ia cosmology with contaminants

Goal: build spectroscopic training set (prior) for classifier (model) to get best constraints on SNIa cosmology given limited follow-up resources

**Wait, what do we want again?**

**“The best” is whatever the metric says it is!**

**But what do we really want?**

**Here, we want a decision metric that identifies most cosmologically impactful light curves for follow-up**

# Classification metrics alone are degenerate

- The *accuracy* is defined as

$$\mathcal{A} = \frac{TP + TN}{N}, \quad (1)$$

where a value closer to unity is more accurate.

- The *purity* (also known as *precision*) is defined as

$$\mathcal{P} = \frac{TP}{TP + FP}, \quad (2)$$

where a value closer to unity is more pure.

- The *efficiency* (also known as *recall*) is defined as

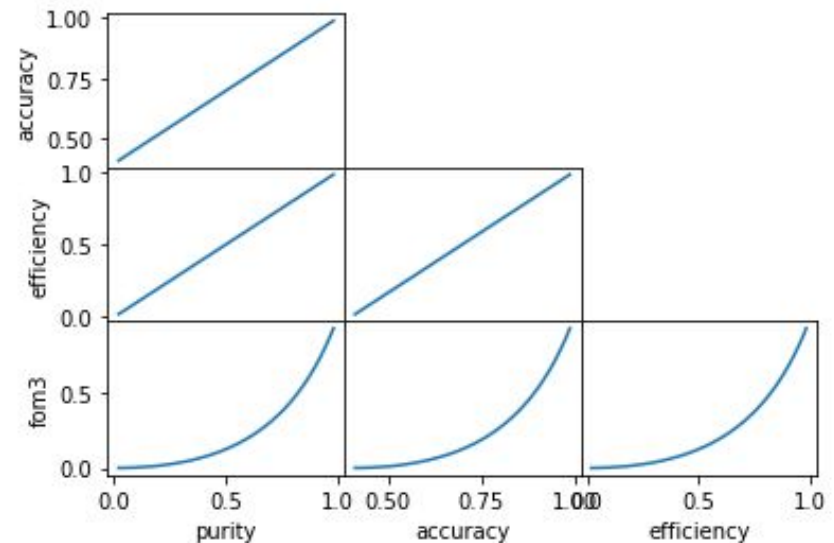
$$\mathcal{R} = \frac{TP}{TP + FN}, \quad (3)$$

where a value closer to unity is more efficient.

- The SNPhotCC defined a *Figure of Merit (FoM)*,

$$\text{FoM}_{W^{\text{false}}} \equiv \text{FoM}(W^{\text{false}}) = \frac{TP}{TP + FN} \times \frac{TP}{TP + W^{\text{false}} \times FP}, \quad (4)$$

where the factor  $W^{\text{false}}$  penalizes false positives. For  $W^{\text{false}} = 1$ ,  $\text{FoM}_1 = \mathcal{R} \times \mathcal{P}$ .<sup>10</sup> We use  $\text{FoM}_3$  in this paper to match the SNPhotCC value of  $W^{\text{false}} = 3$ .



and don't relate directly to cosmology.  
Will using them achieve our goals?

# Classification metrics alone are degenerate

- The *accuracy* is defined as

$$\mathcal{A} = \frac{TP + TN}{N}, \quad (1)$$

where a value closer to unity is more accurate.

- The *purity* (also known as *precision*) is defined as

$$\mathcal{P} = \frac{TP}{TP + FP}, \quad (2)$$

where a value closer to unity is more pure.

- The *efficiency* (also known as *recall*) is defined as

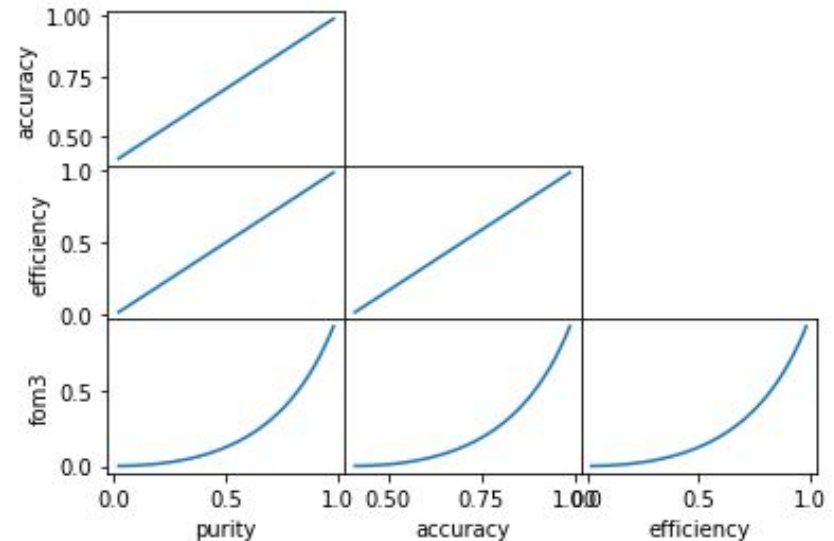
$$\mathcal{R} = \frac{TP}{TP + FN}, \quad (3)$$

where a value closer to unity is more efficient

- The SNPhotCC defined a *Figure of Merit (FoM)*,

$$\text{FoM}_{W^{\text{false}}} \equiv \text{FoM}(W^{\text{false}}) = \frac{TP}{TP + FN} \times \frac{TP}{TP + W^{\text{false}} \times FP}, \quad (4)$$

where the factor  $W^{\text{false}}$  penalizes false positives. For  $W^{\text{false}} = 1$ ,  $\text{FoM}_1 = \mathcal{R} \times \mathcal{P}$ .<sup>10</sup> We use  $\text{FoM}_3$  in this paper to match the SNPhotCC value of  $W^{\text{false}} = 3$ .



and don't relate directly to cosmology.  
 Will using them achieve our goals?

# Is the metric appropriate?

Think *adversarially* to trick the metric into rewarding a wrong answer.

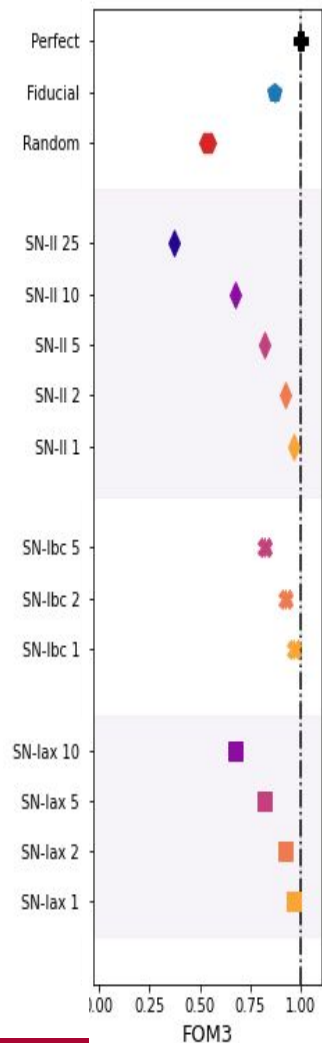
How can you satisfy the metric's criteria without satisfying the criteria you as an experimenter actually care about?

Identify *pathological* cases that fool the metric but are obviously wrong.

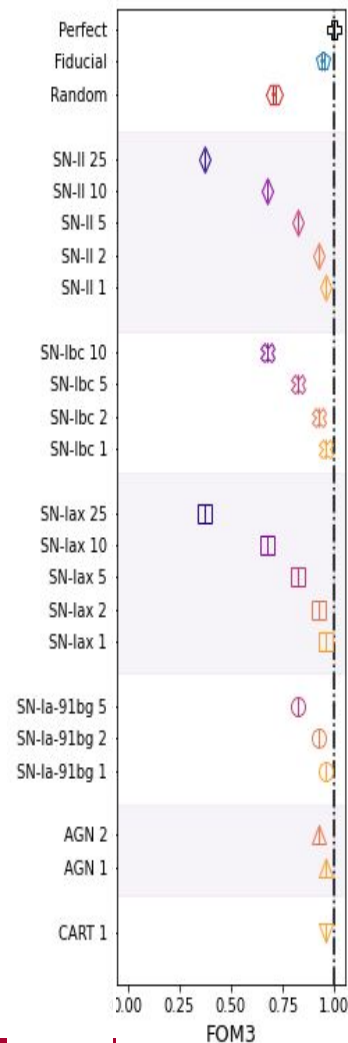
Can you think of a null test in the form of a trivially bad answer that performs well by the metric?

# SNPhotCC FoM can't distinguish contaminant class

## DDF



## WFD





# Quantitative metrics of posterior samples $\{w_i\}$

The *Kullback-Leibler Divergence (KLD)*,

$$KLD = - \int \hat{p}_0(w) \ln \left[ \frac{\hat{p}_{mock}(w)}{\hat{p}_0(w)} \right] dw, \quad (5)$$

is an information theoretic measure of the loss of information due to using an approximation  $\hat{p}_{mock}(w)$  rather than the true distribution  $\hat{p}_0(w)$ ; the KLD has been used before in extragalactic astrophysics (e.g. Malz et al. (2018)).

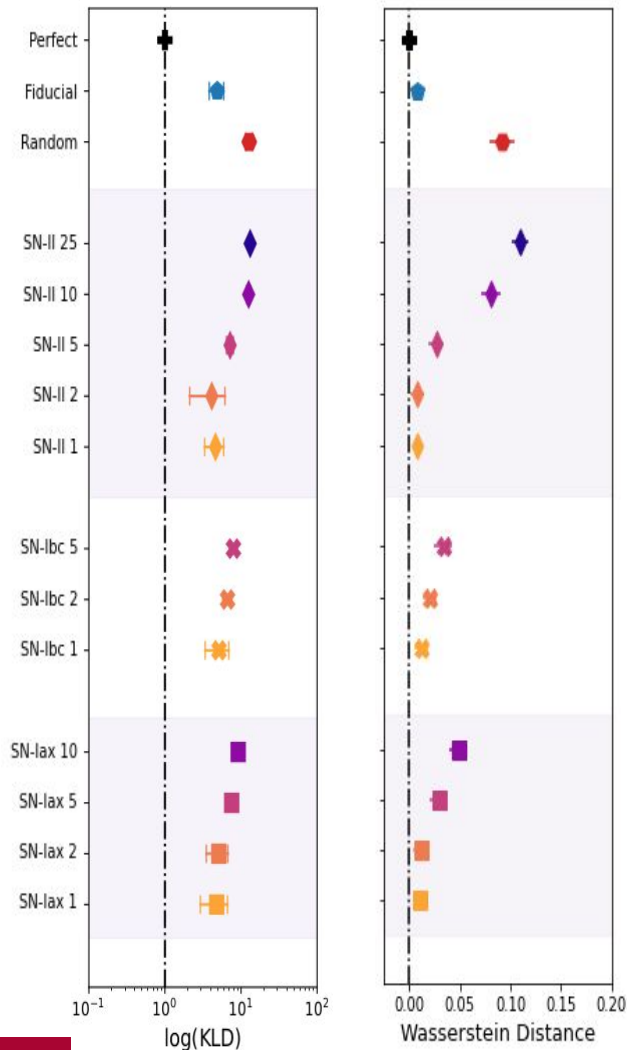
The *Earth-Mover's Distance (EMD)*

$$EMD = \int_{-\infty}^{\infty} \left| \int_{-\infty}^w \hat{p}_0(w') dw' - \int_{-\infty}^w \hat{p}_{mock}(w') dw' \right| dw, \quad (6)$$

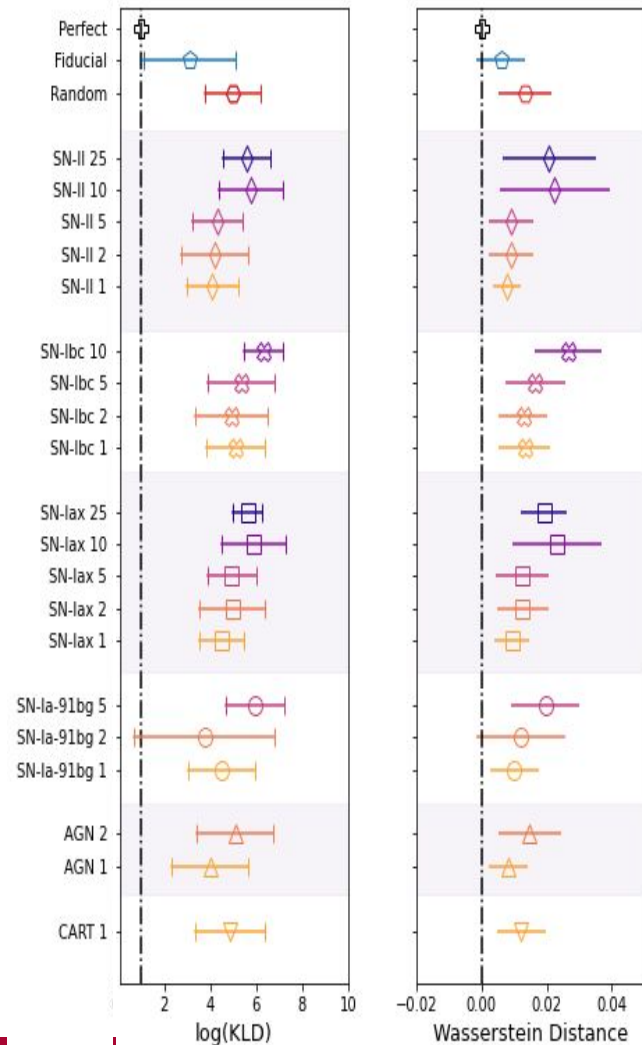
(also known as the first order *Wasserstein metric*) can be intuitively understood as the integrated discrepancy between a pair of PDFs, defined in terms of their cumulative distribution functions (CDFs); the EMD has been used before in cosmology (e.g. Moews et al. 2021).

# Posterior samples' KLD and EMD are sensitive to contaminant class

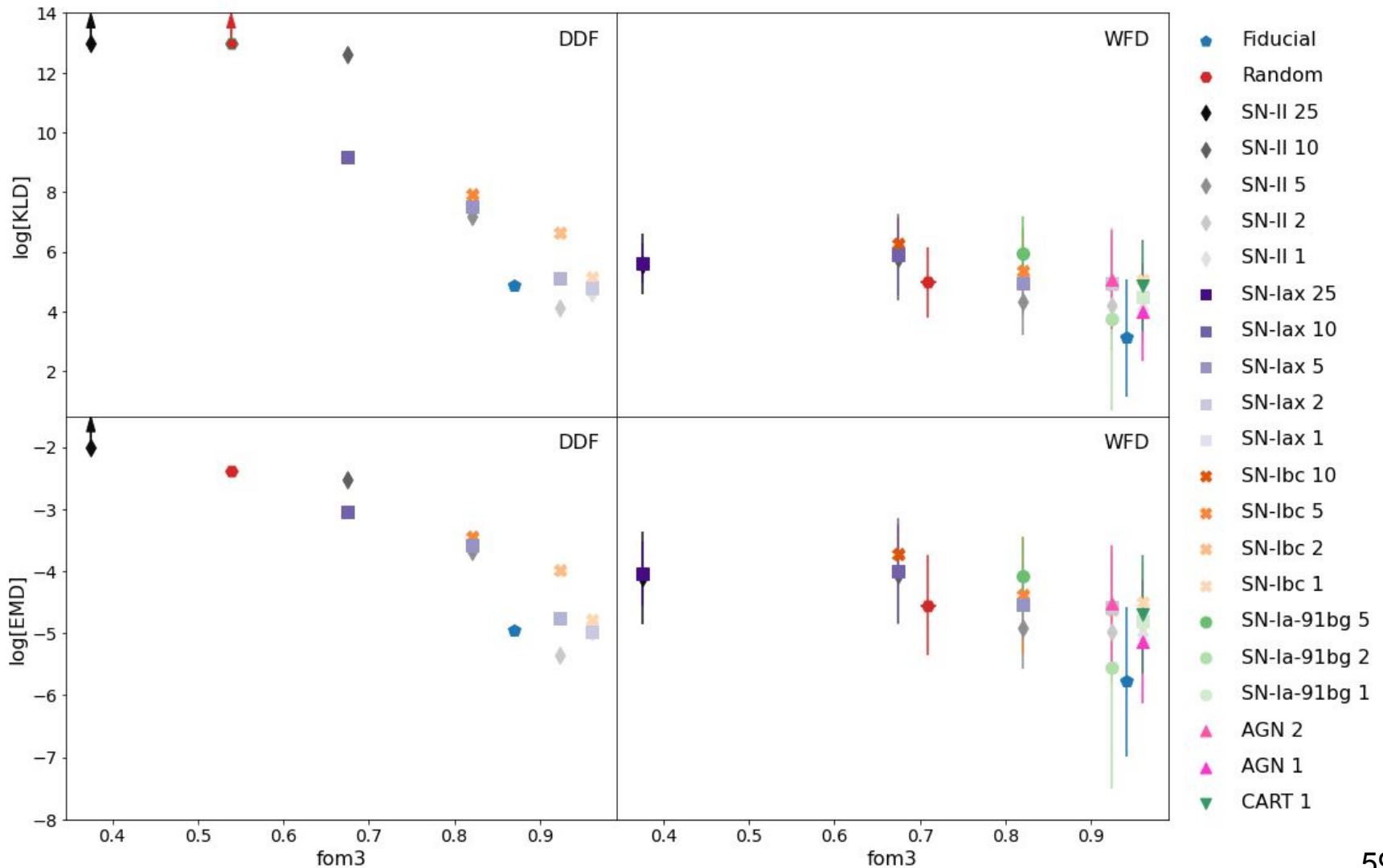
## DDF



## WFD



# Contaminant type matters to cosmology



# Closing thoughts



Probabilistic graphical models are pretty much my favorite thing!  
But my secret agenda is for you to incorporate this way of thinking into your hacks, for the experiments in which you prove how great your hierarchical models are.