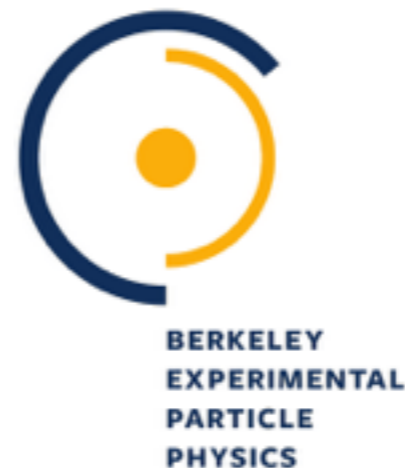# Learning to Classify from Impure Samples with High-Dimensional Data

based on Phys. Rev. D 98, 011502(R) [published yesterday!]

Patrick Komiske, Eric Metodiev,
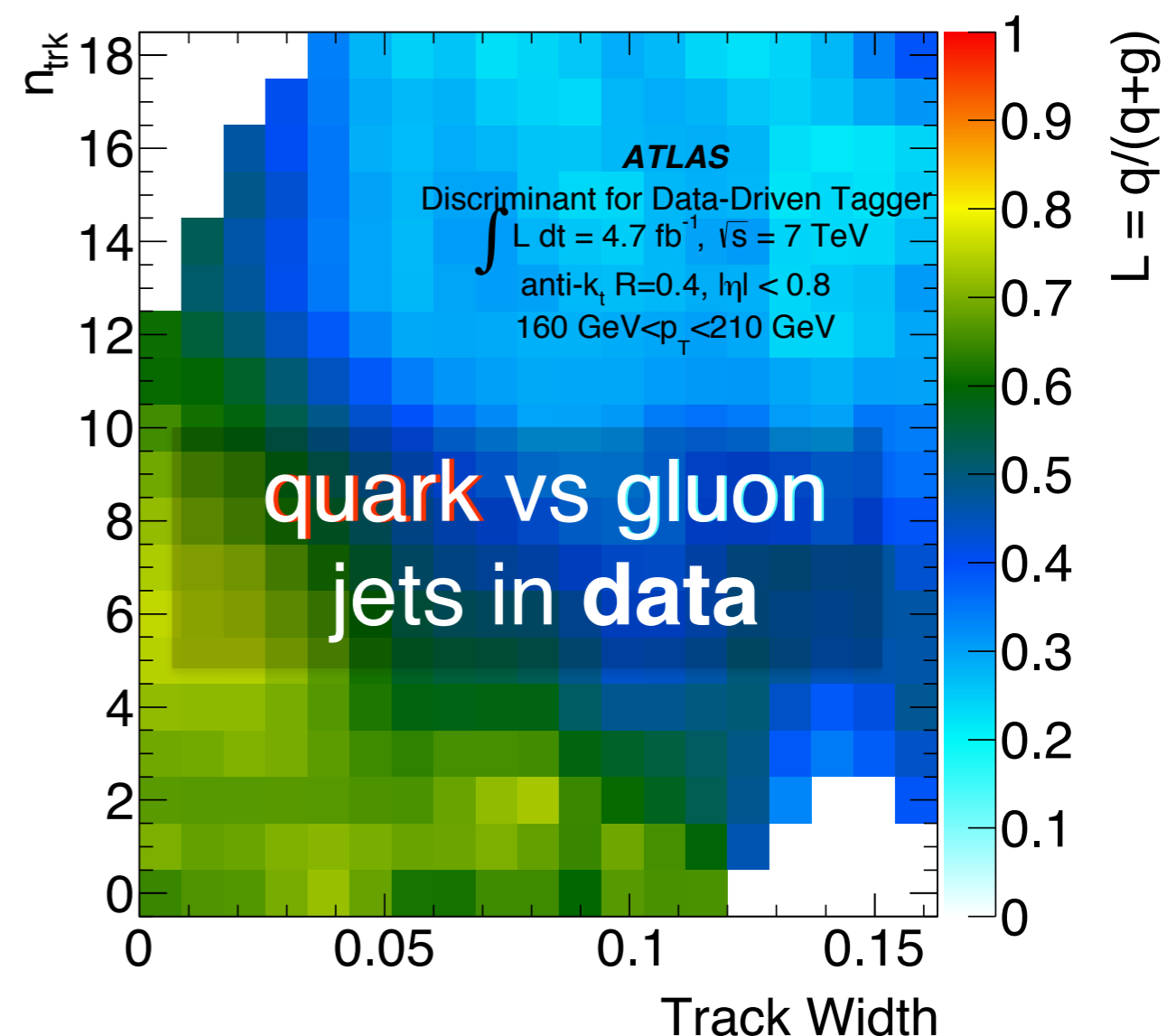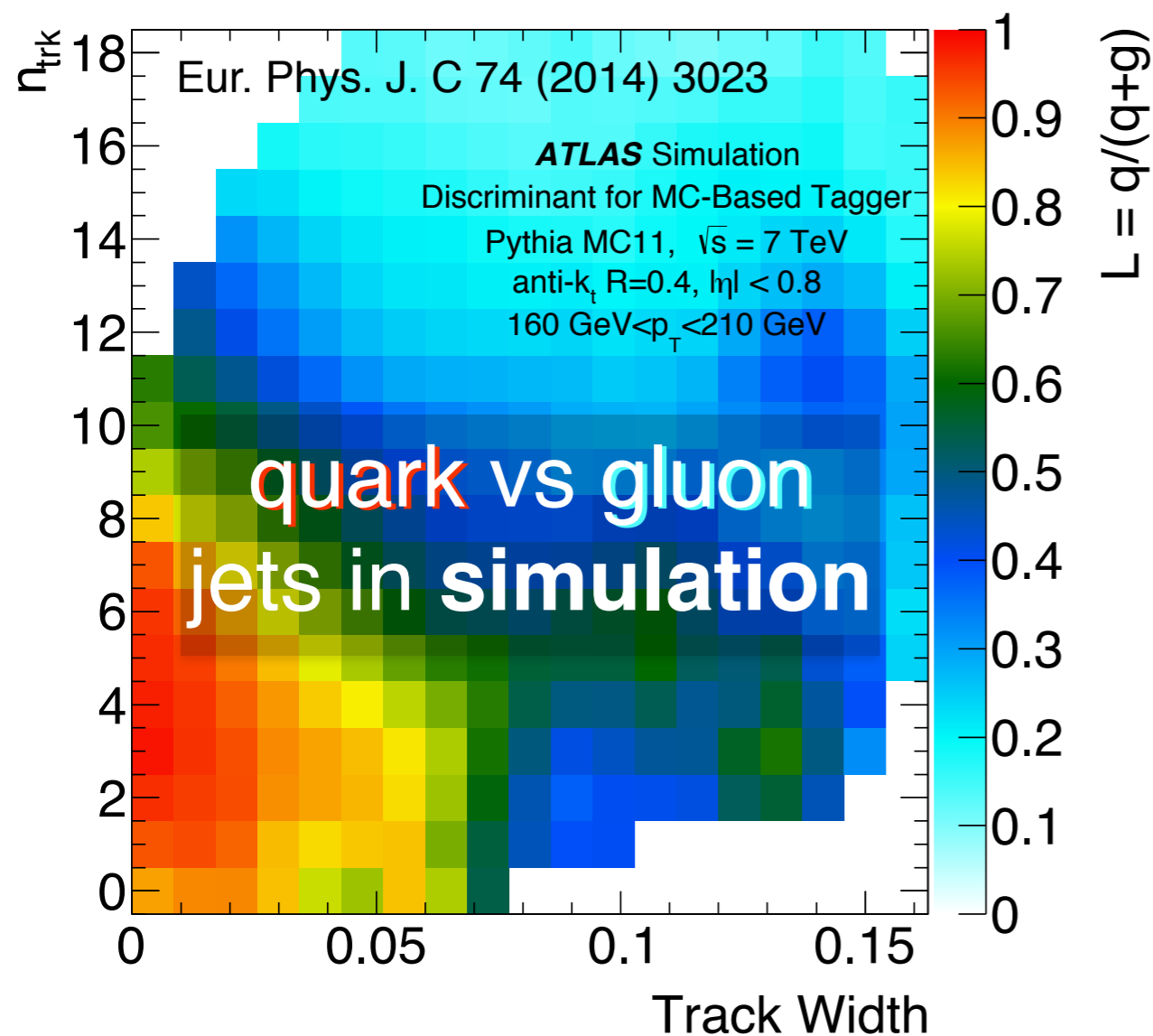Benjamin Nachman, Matt Schwartz

*...building on work also in collaboration with Lucio Dery, Francesco Rubbo, Ariel Schwartzman, Jesse Thaler*

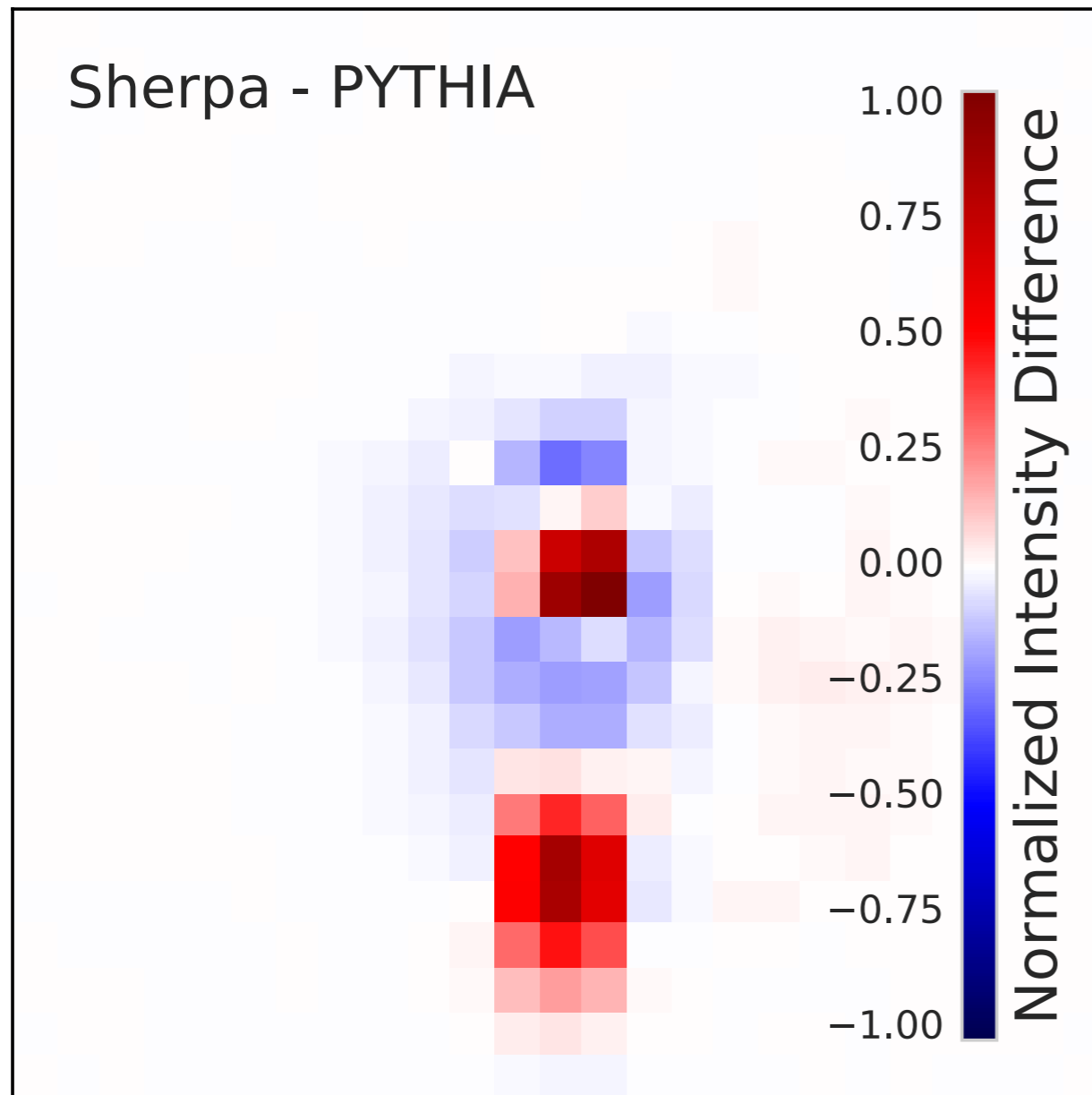Usual paradigm: train in simulation, test on data.



If data and simulation differ, this is **sub-optimal**!

Sherpa - PYTHIA

Normalized Intensity Difference

1.00
0.75
0.50
0.25
0.00
−0.25
−0.50
−0.75
−1.00

Especially important for **deep learning** using subtle features → hard to model!

*W boson radiation pattern - same physics, different simulators!*

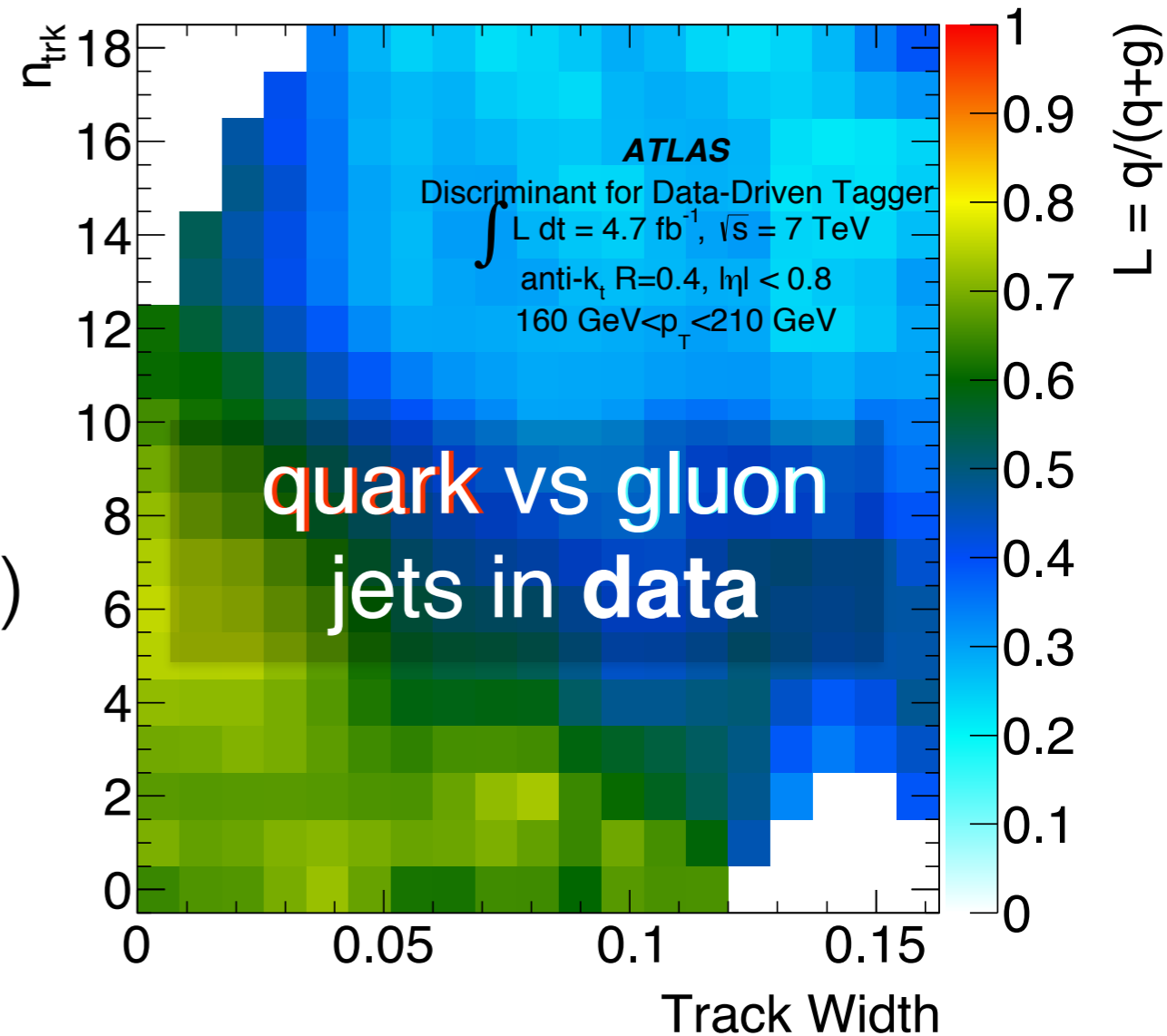J. Barnard, E. Dawe, M. Dolan, N. Rajcic, Phys. Rev. D 95 (2017) 014018

How did we make this plot?

dijets = $f_q$ × Q + (1-$f_q$) × G

Z+jets = $g_q$ × Q + (1-$g_q$) × G

two equations, two unknowns (Q, G)

We often know f, g
(from ME + PDF) much better than
full radiation pattern inside jets.



This doesn't work well when you have more than 2
observables because the templates become sparse.

$$f_{\text{full}} = \text{argmin}_{f':\mathbb{R}^n \to \{0,1\}} \sum_{i=1}^{N} \ell\left(f'(x_i) - t_i\right)$$

*loss fcn.*

*labels*

## LoLiProp

*Learning from Label Proportions*

Solution: Train using class proportions. Work "on average"

Mixed Sample 1

Mixed Sample 2

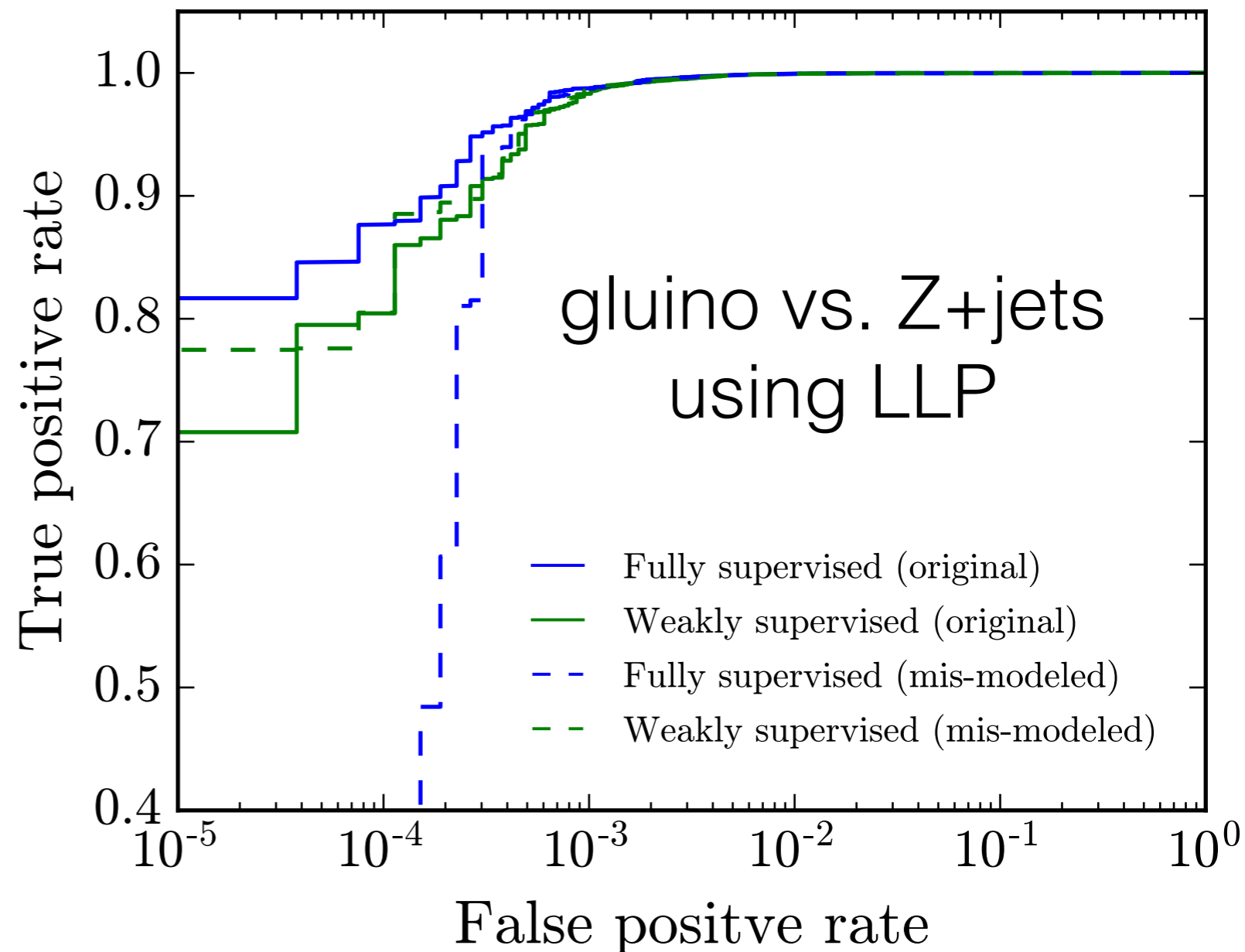$$f_{\text{weak}} = \text{argmin}_{f':\mathbb{R}^n \to [0,1]} \ell\left(\sum_{i=1}^{N} \frac{f'(x_i)}{N} - y\right)$$
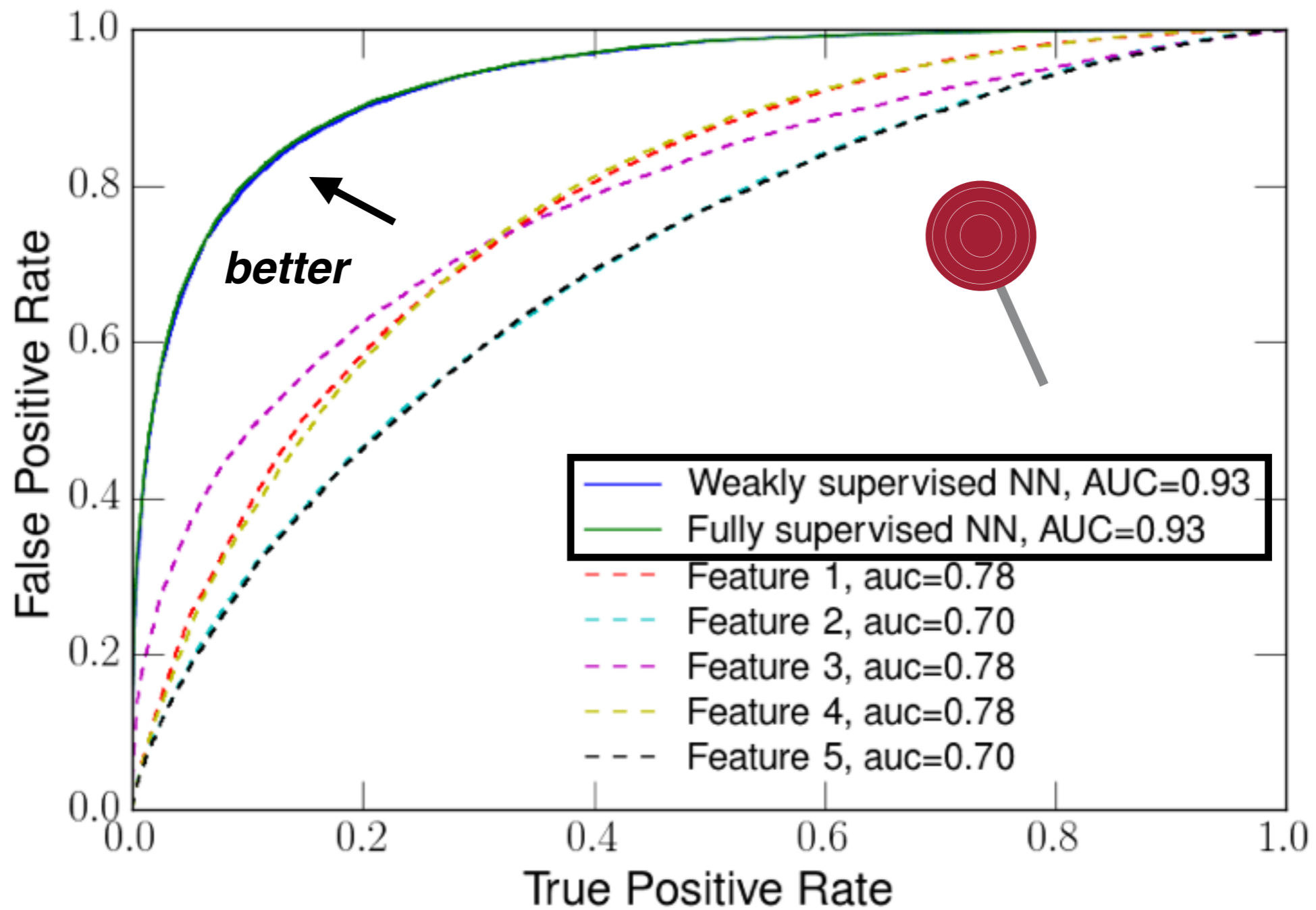
*proportions*

L. Dery, **BPN**, F. Rubbo, A. Schwartzman, JHEP 05 (2017) 145

Even though the proportions are required as input, if they are slightly wrong, you can end up with the correct classifier.



gluino vs. Z+jets
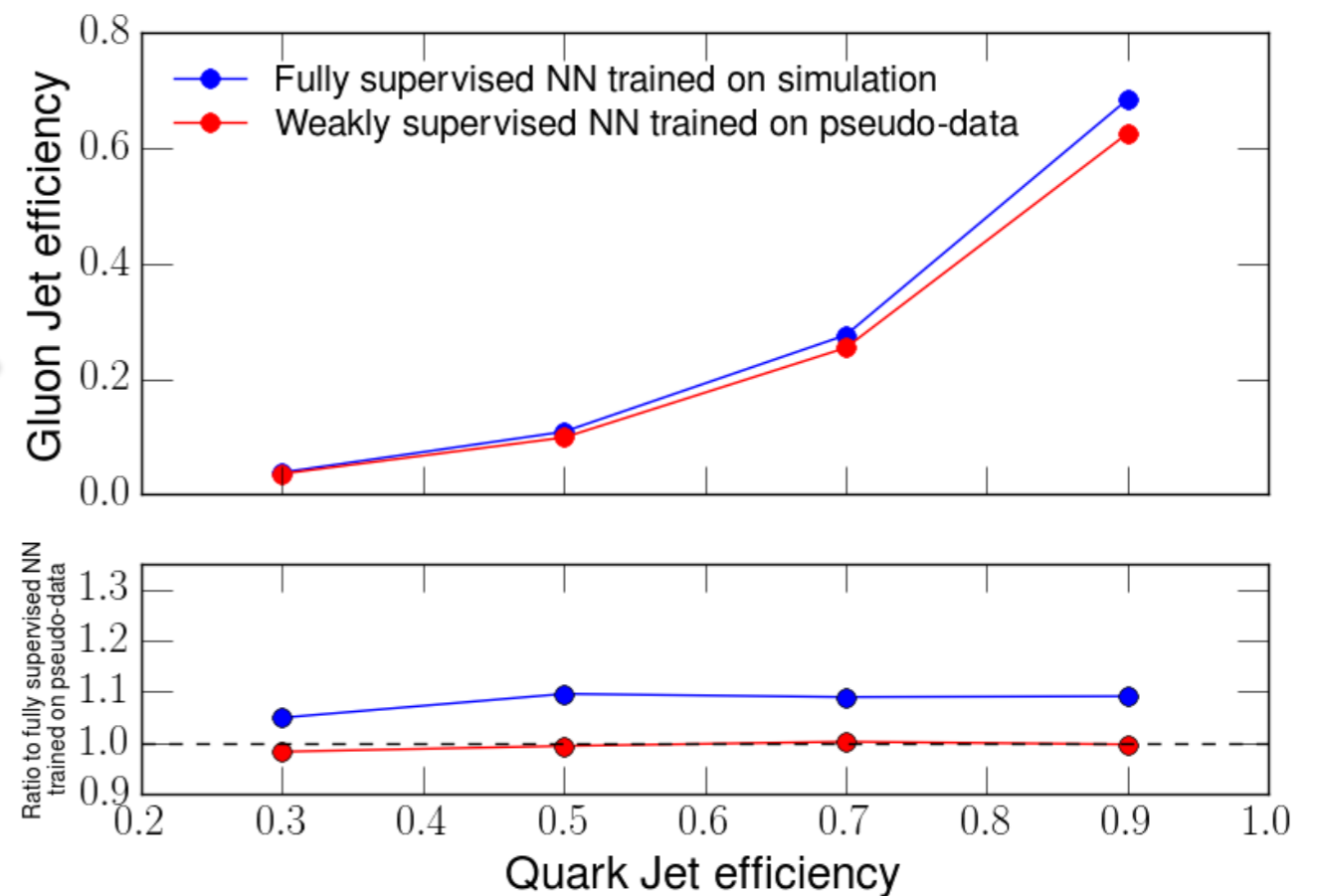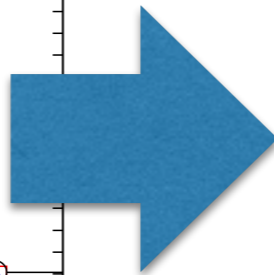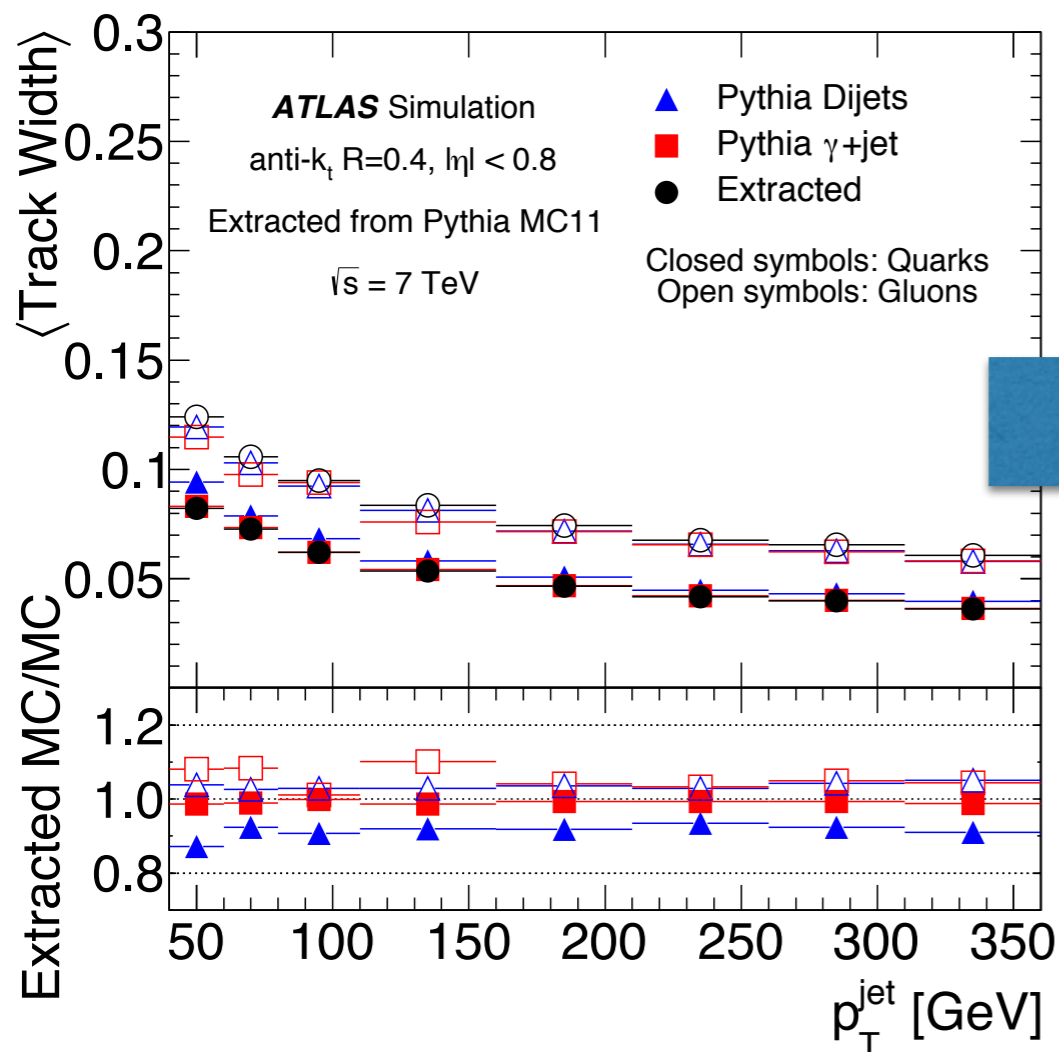using LLP

—— Fully supervised (original)
—— Weakly supervised (original)
- - - Fully supervised (mis-modeled)
- - - Weakly supervised (mis-modeled)
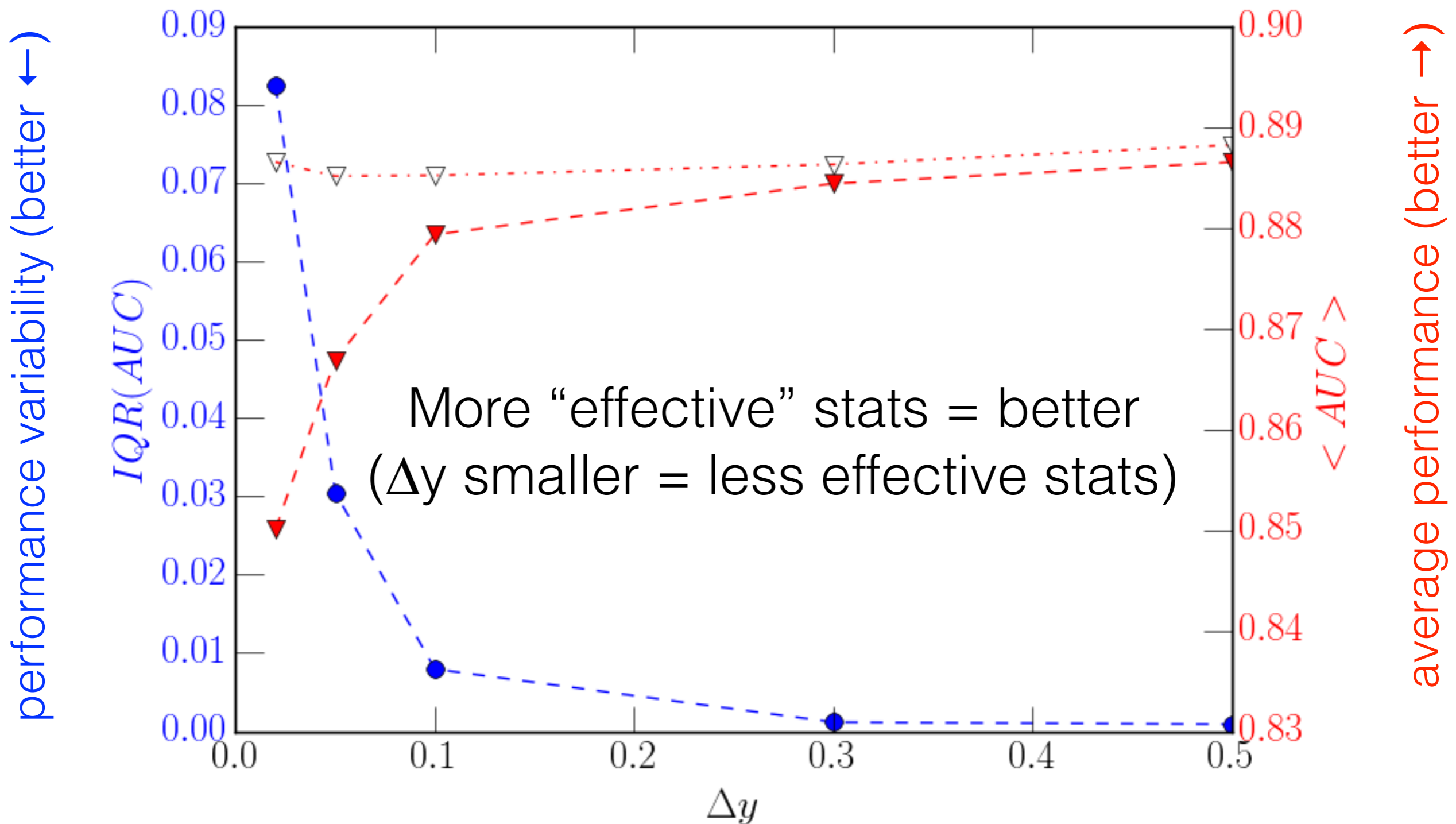
T. Cohen, M. Freytsis, B. Ostdiek, https://arxiv.org/abs/1706.09451

Given the data/MC disagreement from the first slide, this is what you might expect in terms of the performance difference.

More "effective" stats = better
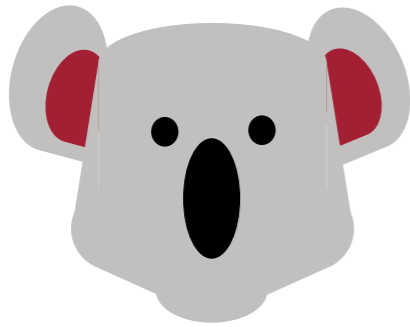(Δy smaller = less effective stats)

how different are the proportions for the two mixed samples

**CWoLa**

*Classification Without Labels*

Solution: Train **directly on data** using mixed samples



E. Metodiev, **BPN**, J. Thaler, JHEP 10 (2017) 51

As with LLP, need sufficient effective statistics

Can't learn when the two proportions are the same.

| Property | LLP | CWoLa |
|---|:---:|:---:|
| Compatible with any trainable model | ✓ | ✓ |
| No training modifications needed | ✗ | ✓ |
| Training does not need fractions | ✗ | ✓ |
| Smooth limit to full supervision | ✗ | ✓ |
| Works for $> 2$ mixed samples | ✓ | ? |

There are many O(1)-dimensional ML problems for jets, but since the full radiation pattern is higher dimensional, need to go to bigger!

We'll use jet images as a testing ground, still focusing on quarks versus gluons.

after Pixel Standardization

gluons

Translated Azimuthal Angle $\phi$

Translated Pseudorapidity $\eta$

after Pixel Standardization

quarks

Translated Azimuthal Angle $\phi$
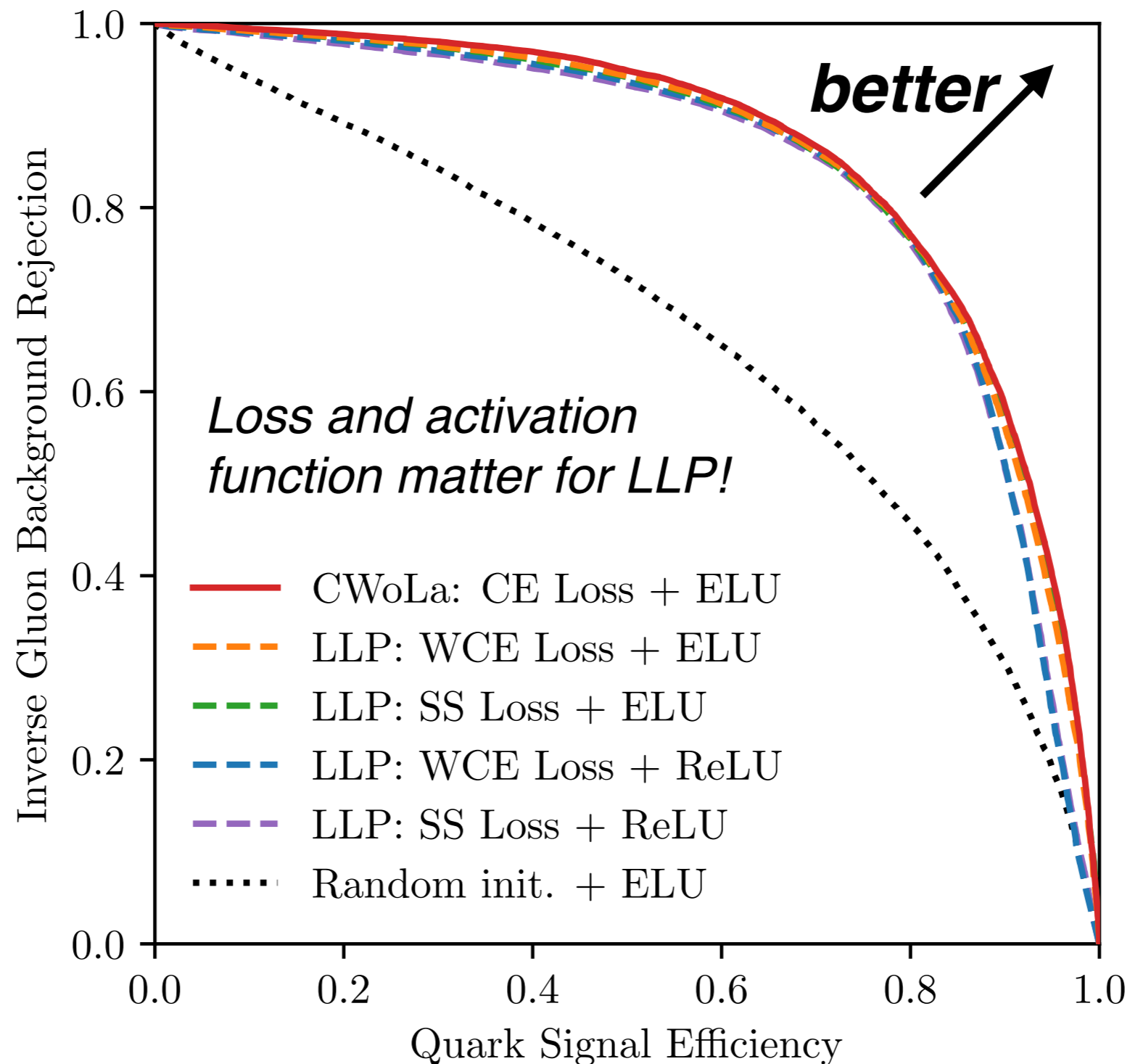
Translated Pseudorapidity $\eta$

The CWoLa approach works out-of-the box - can use well-tested CNN architecture with usual cross-entropy loss.

On the other hand, LLP requires significant work on the technical implementation / optimization.

$$\ell_{\mathrm{WMSE}} = \sum_a \left( f_a - \frac{1}{N} \sum_{i=1}^{N} h(\mathbf{x}_i) \right)^2 \qquad \ell_{\mathrm{WCE}} = \sum_a \mathrm{CE}\left( f_a, \frac{1}{N} \sum_{i=1}^{N} h(\mathbf{x}_i) \right)$$
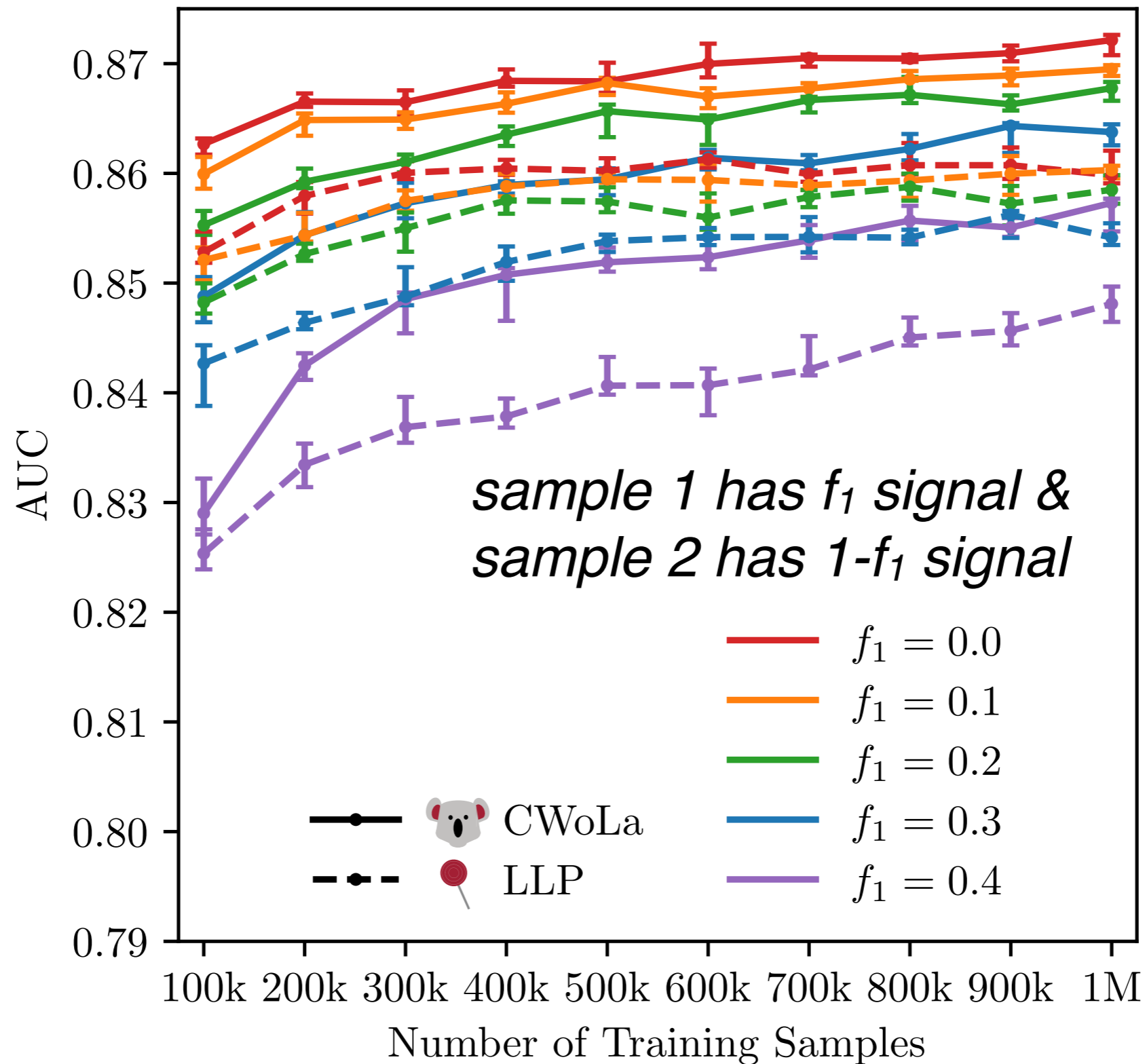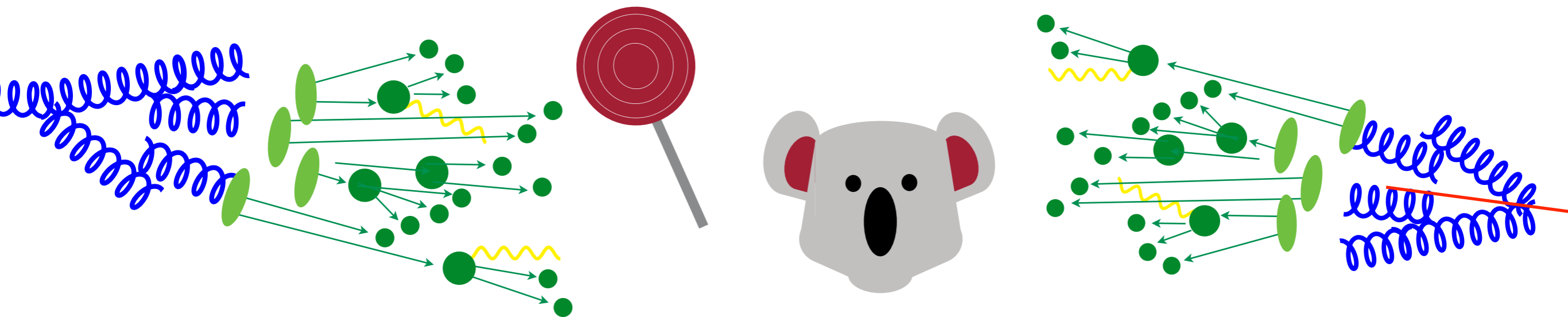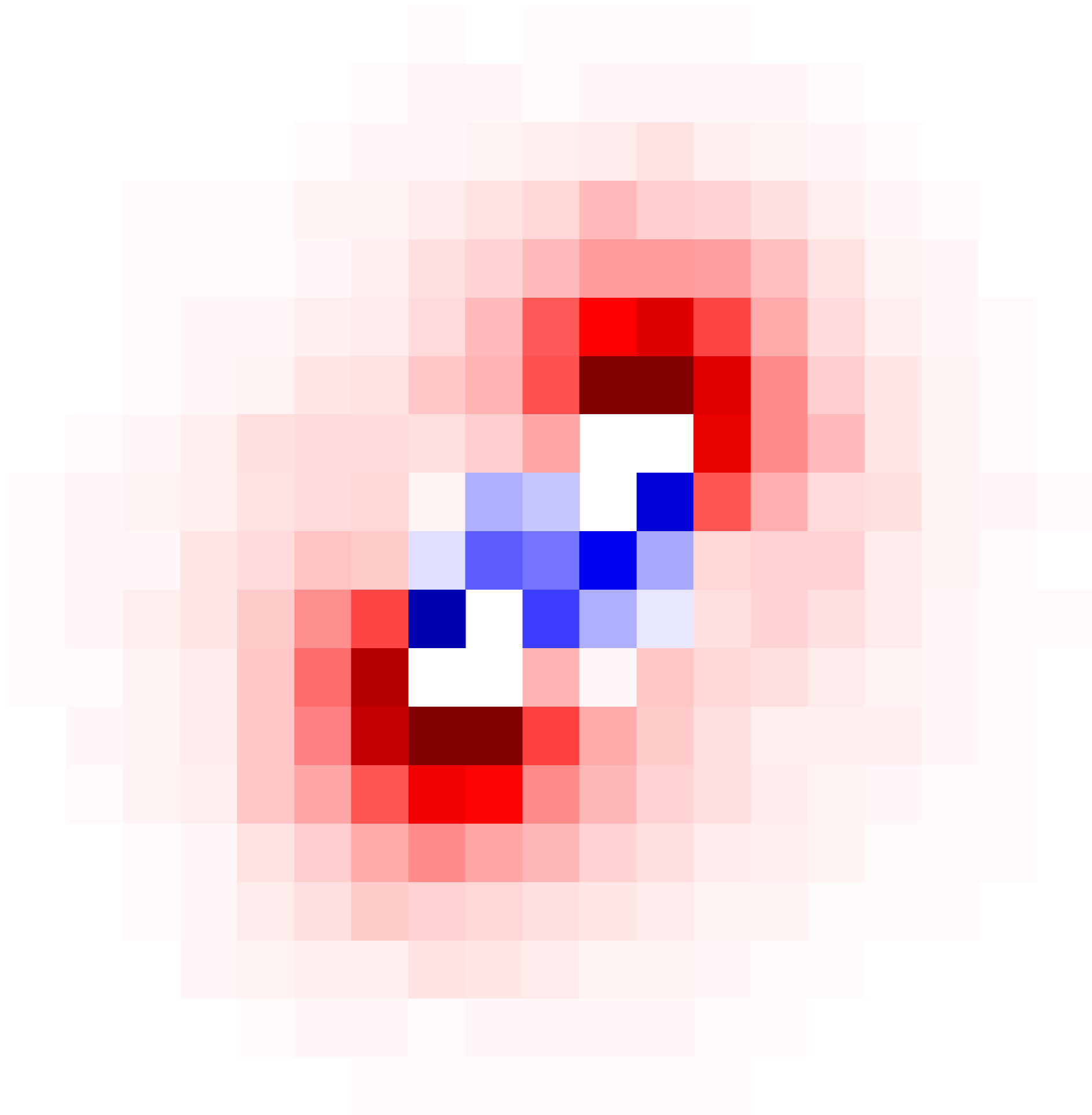
Weak supervision is a new & exiting paradigm for training classifiers. We can learn directly from nature!

This is particularly important for jet physics, where there are concerns about mis-modeling. We have shown that the methods work even for high-dimensional data.



To see how else these ideas could be used, see Jack's CWoLa hunting talk and Eric's Jet Topics talk

Fin.

| Learning | Sample | AUC |
|----------|--------|-----|
| CWoLa | $Z$+jet vs. dijets | $0.8626 \pm 0.0020$ |
|  | Artificial $Z + q/g$ | $0.8621 \pm 0.0019$ |
| LLP | $Z$+jet vs. dijets | $0.8544 \pm 0.0019$ |
|  | Artificial $Z + q/g$ | $0.8549 \pm 0.0018$ |