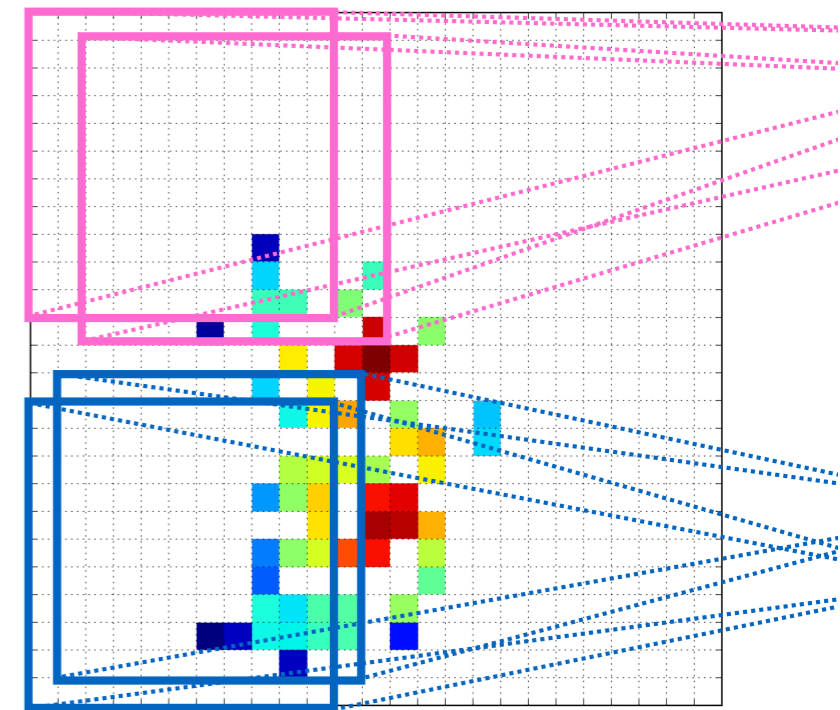


Reducing our dependence on slow simulators with deep learning

Benjamin Nachman

Lawrence Berkeley National Laboratory

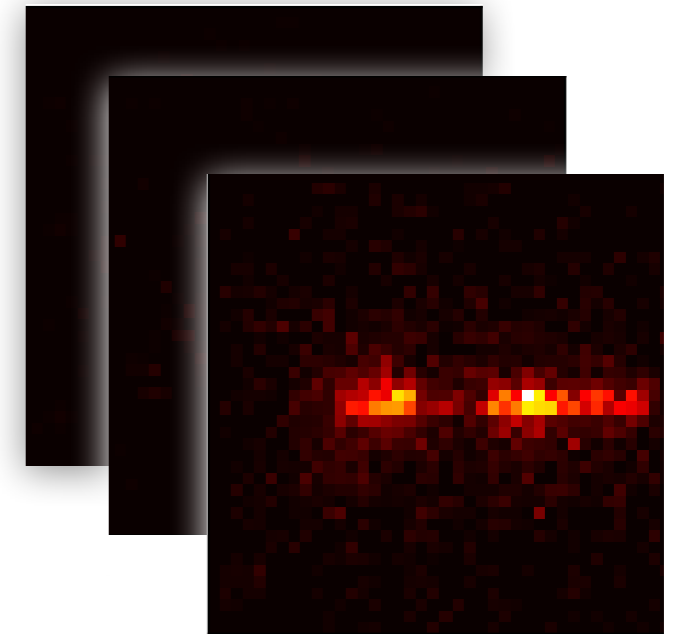


*TRISEP School
August 2, 2019*

ML for classification and beyond



There is a lot more ML can do than classify examples!



Classification

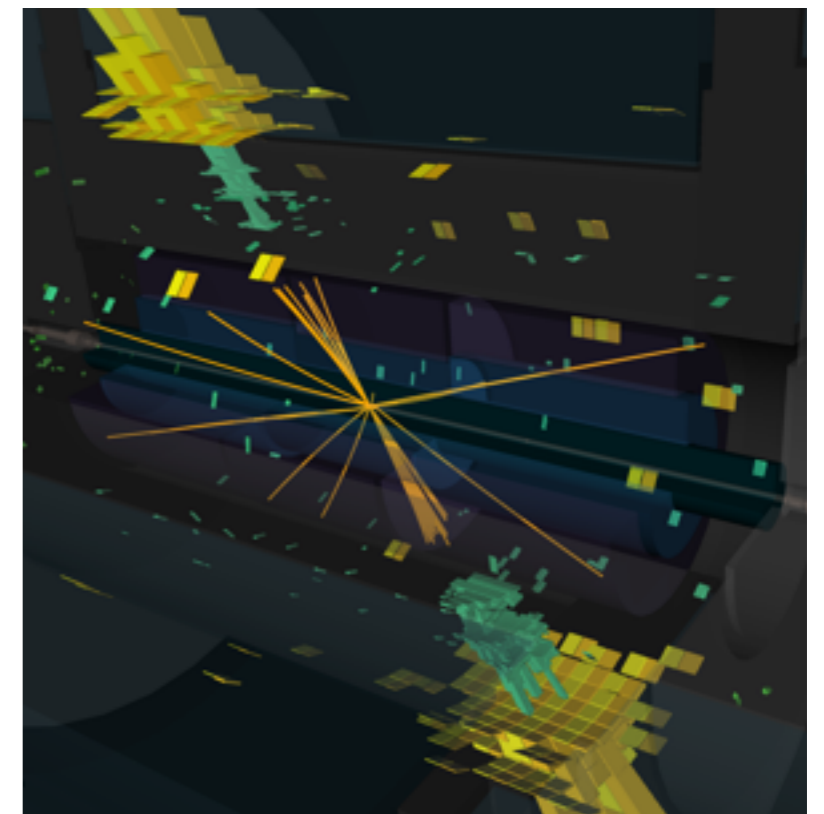
provide examples for training

arbitrarily many categories

Generation

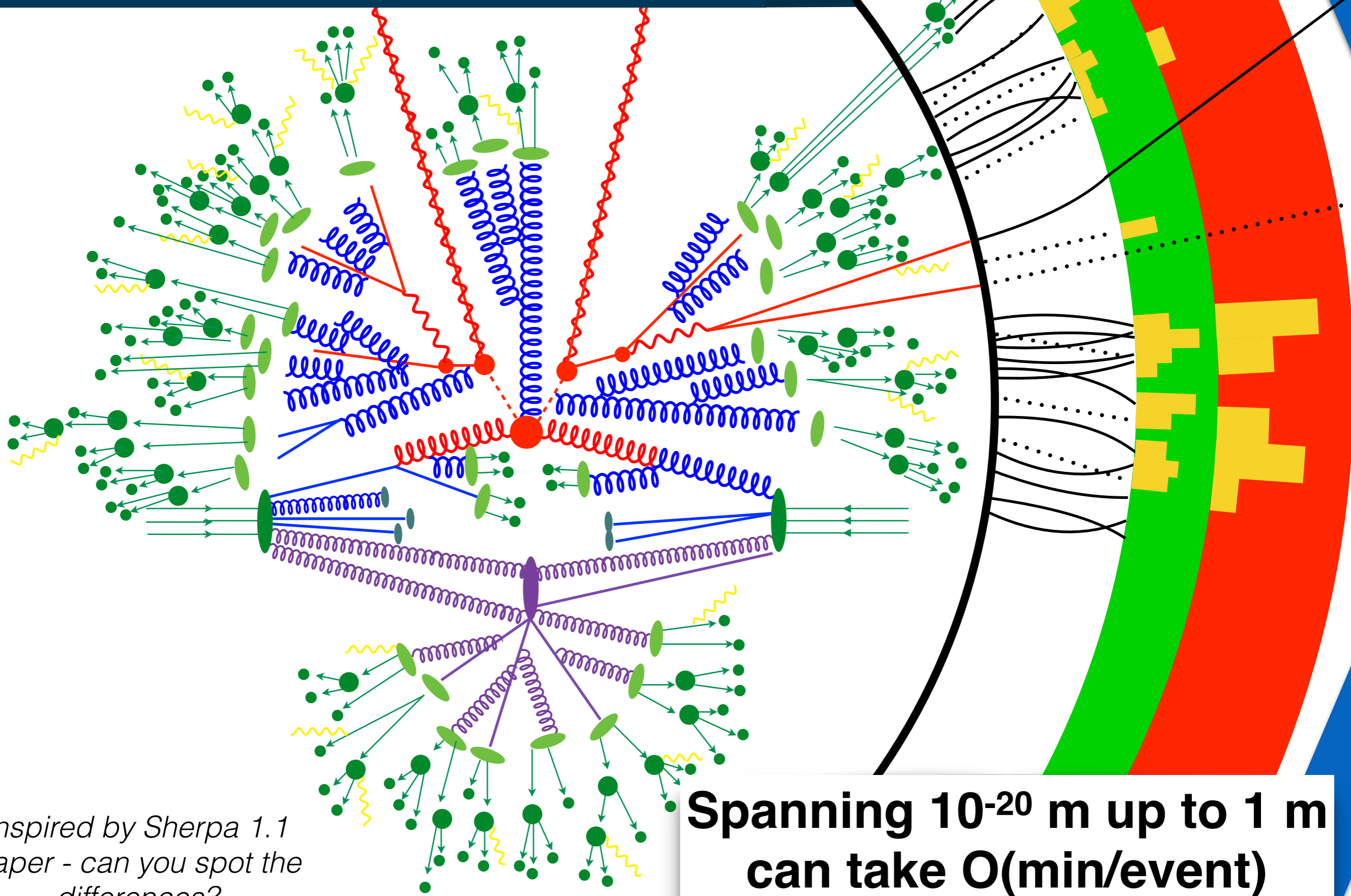
Regression

map noise to structure



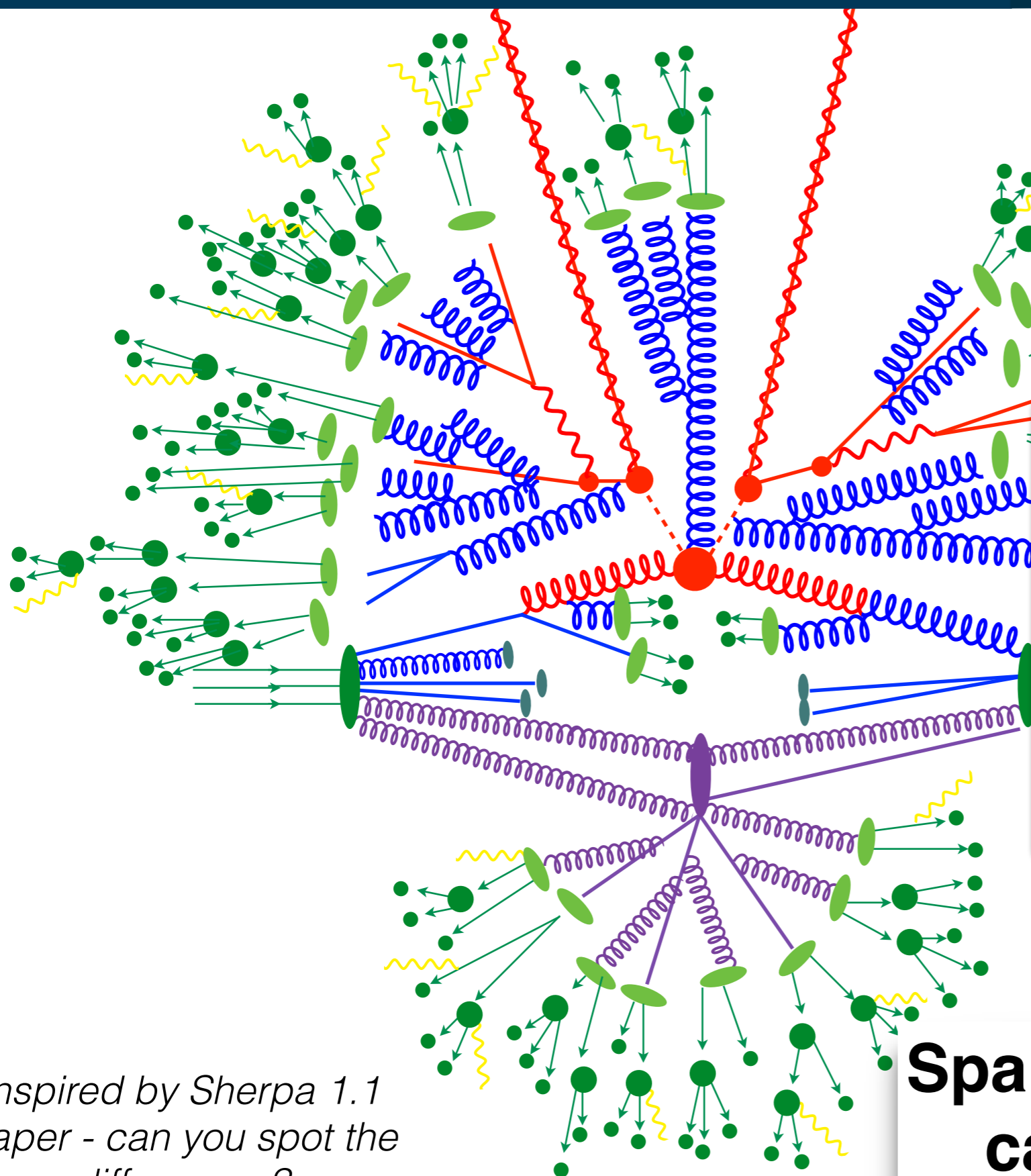
- Simulation dependence in collider HEP
- **Generative models**: accelerate simulations with deep learning
- **Generalized numerical inversion**: machine learn with prior independence
- **Weak supervision**: learn directly from (unlabeled) data
- CWoLa hunting: model agnostic **anomaly detection**

Simulation at the LHC



Inspired by Sherpa 1.1 paper - can you spot the differences?

**Spanning 10^{-20} m up to 1 m
can take O(min/event)**



Part I:
accelerate
the existing
simulation

**Spanning 10^{-20} m up to 1 m
can take $O(\text{min}/\text{event})$**

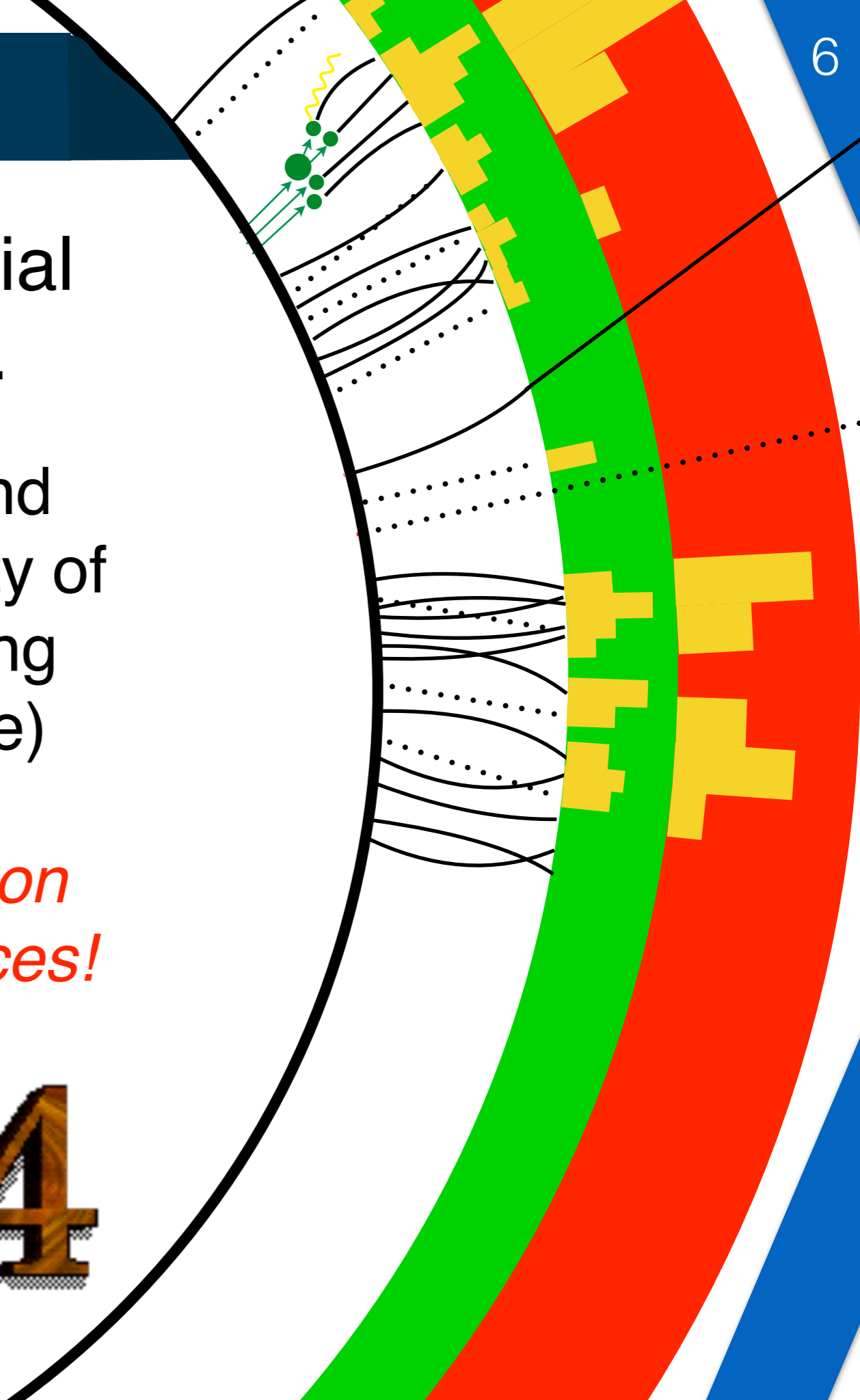
*Inspired by Sherpa 1.1
paper - can you spot the
differences?*

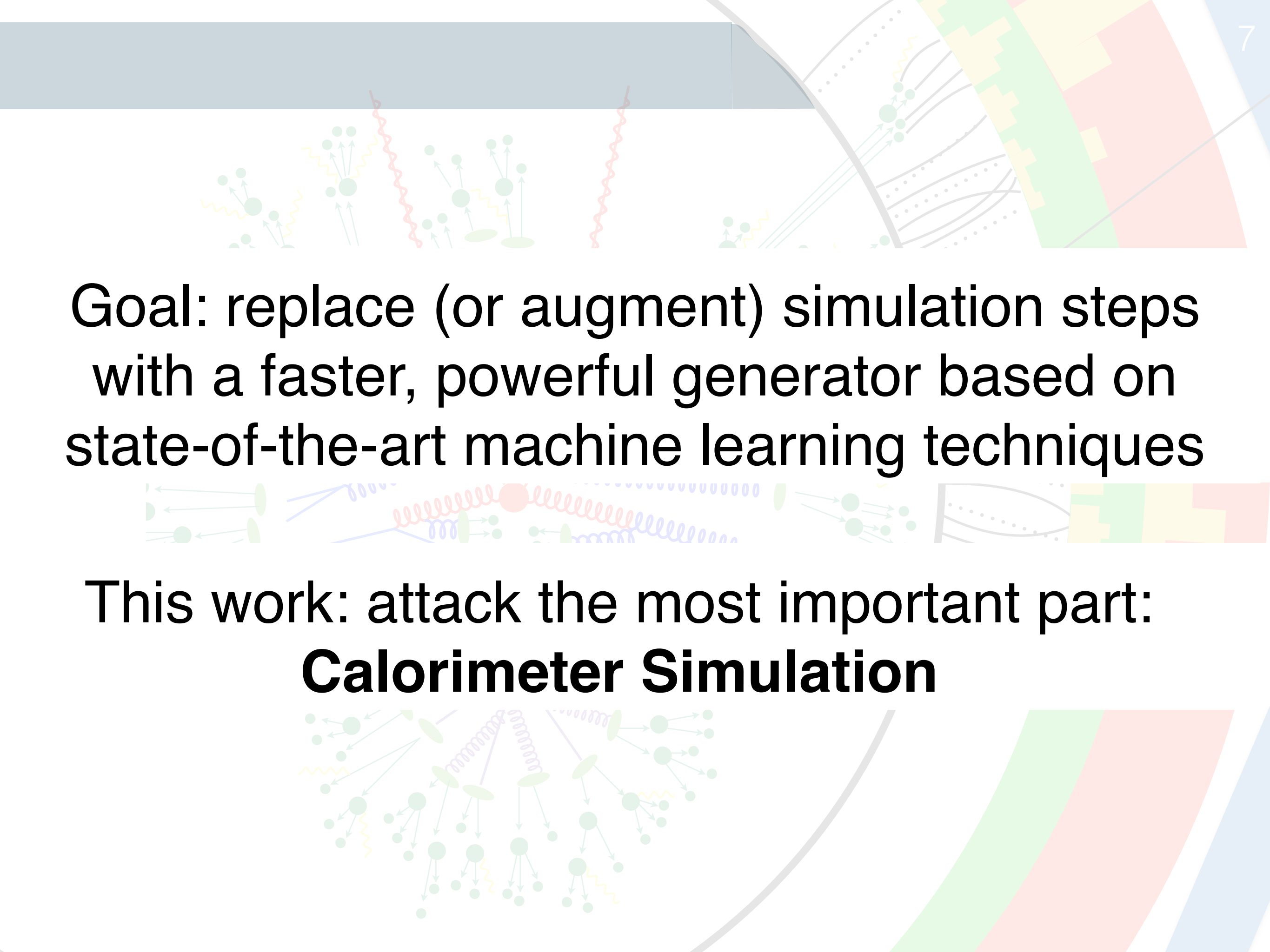
State-of-the-art for material interactions is Geant 4.

Includes electromagnetic and hadronic physics with a variety of lists for increasing/decreasing accuracy (at the cost of time)

This accounts for $O(1)$ fraction of all HEP computing resources!

Geant 4





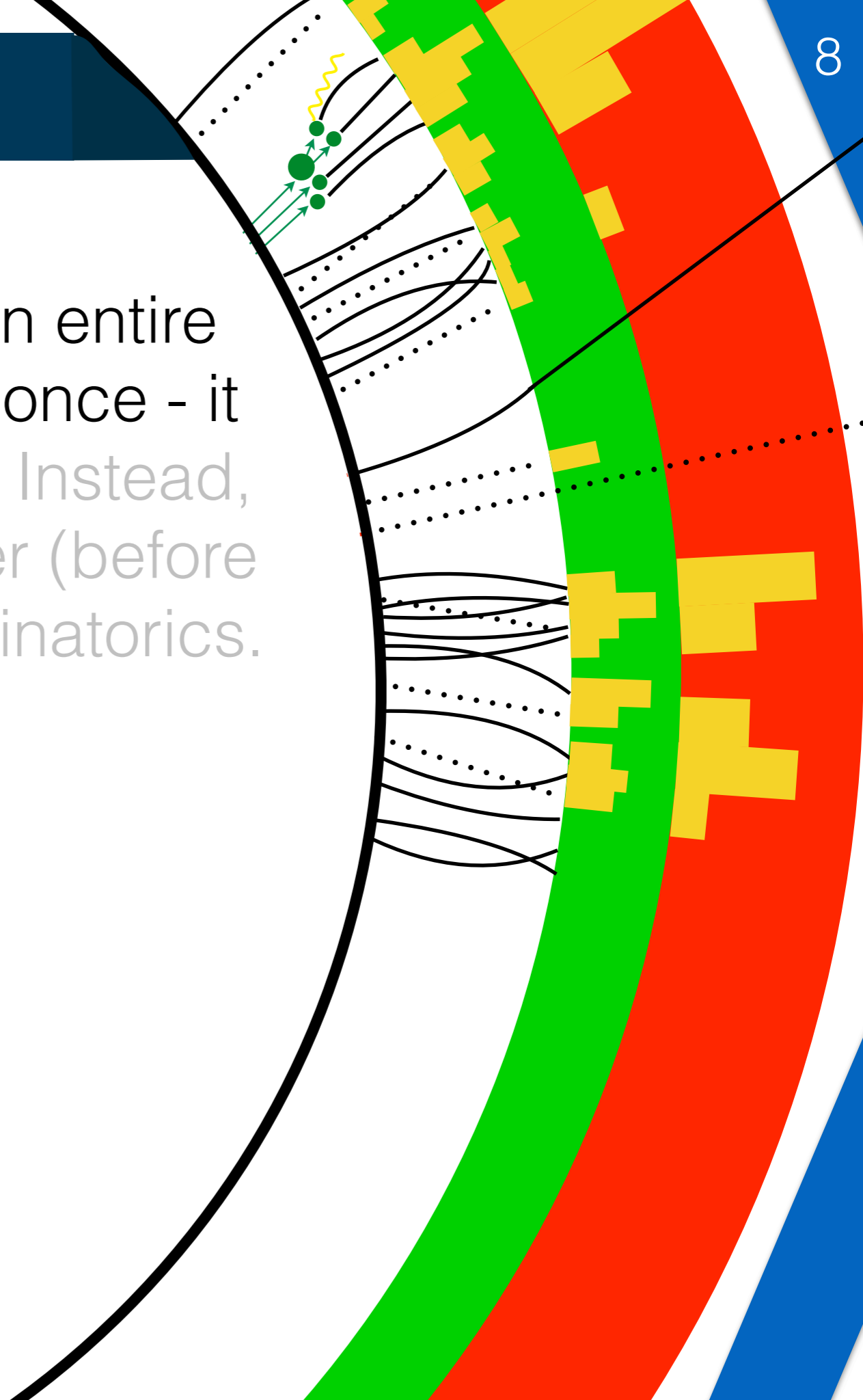
Goal: replace (or augment) simulation steps with a faster, powerful generator based on state-of-the-art machine learning techniques



This work: attack the most important part:
Calorimeter Simulation

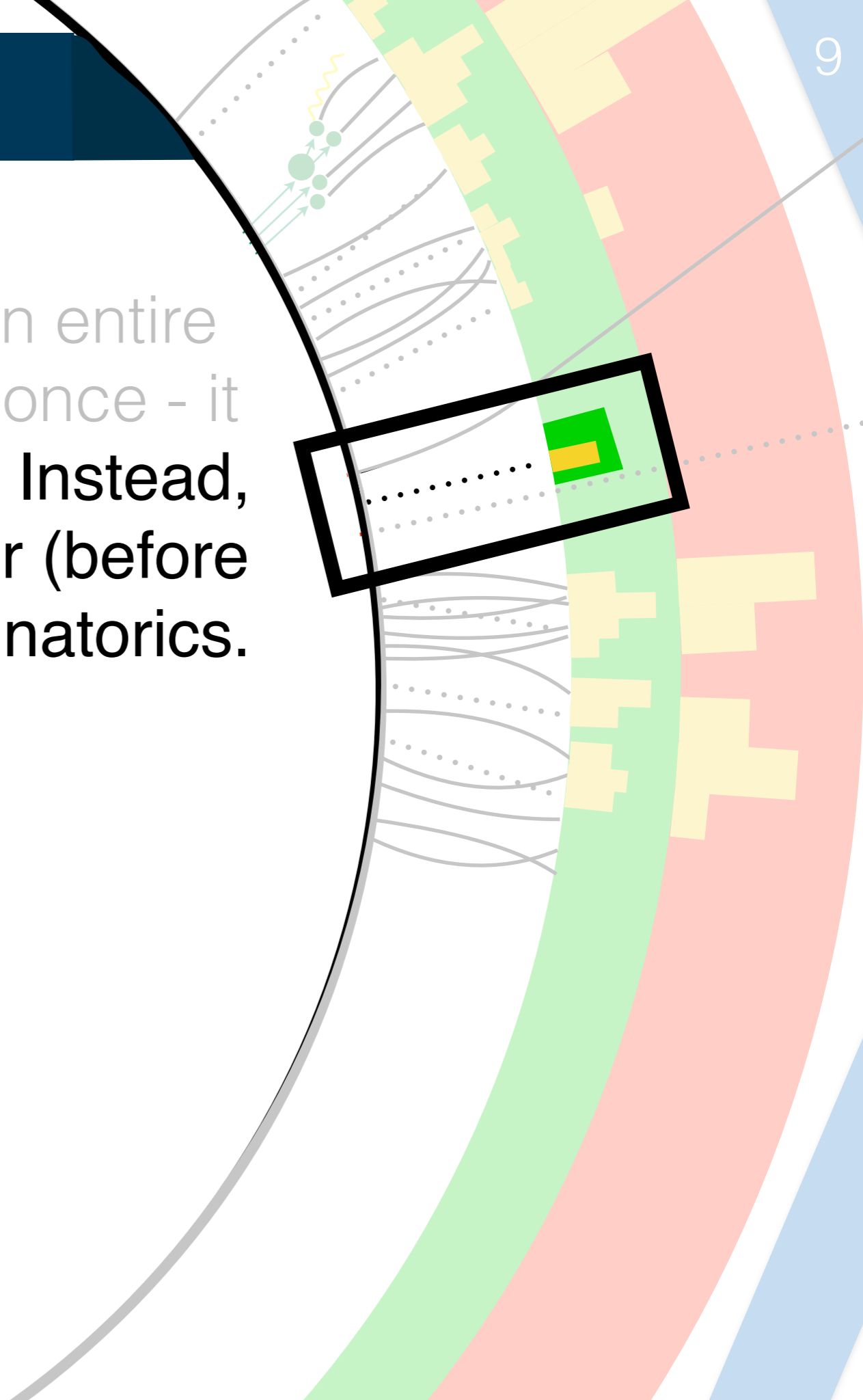
Factorization

We are not trying to generate an entire event ($O(1000)$ particles) all at once - it would be **very hard to validate!** Instead, generate a single particle shower (before electronics) and appeal to combinatorics.



Factorization

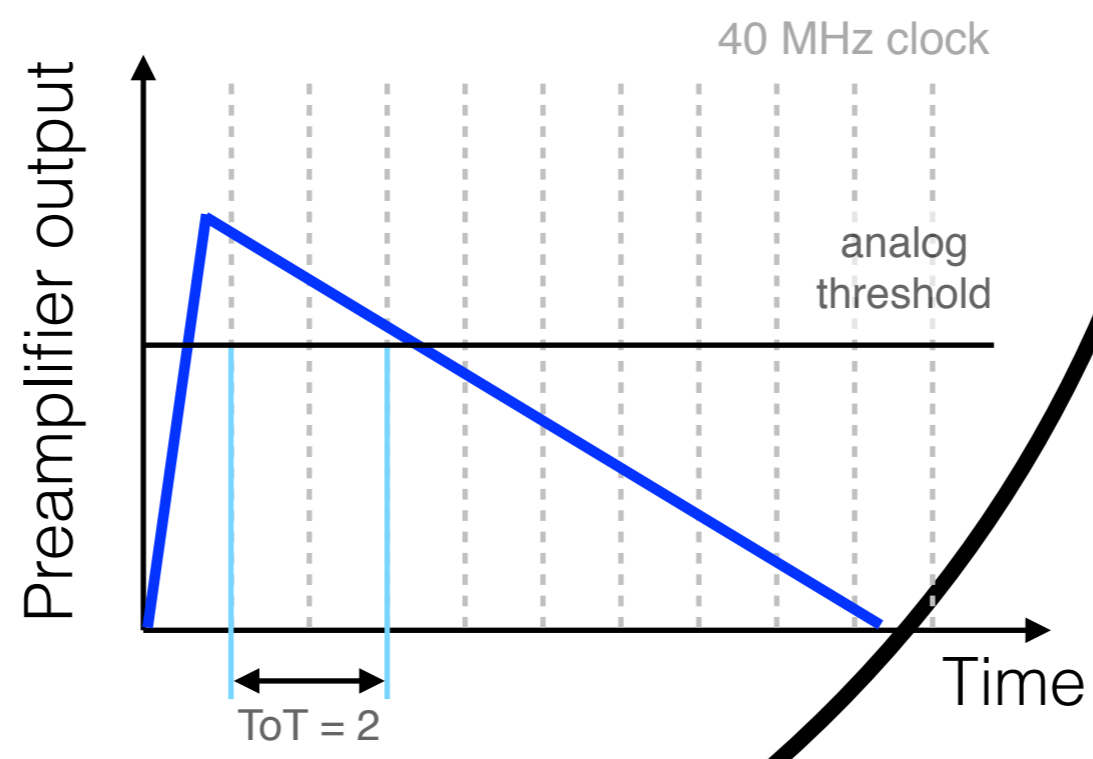
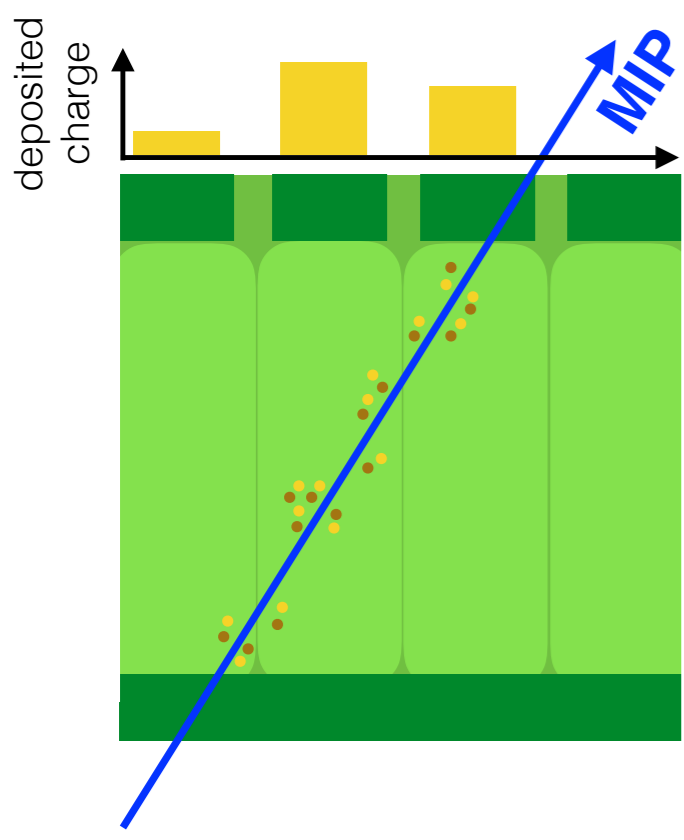
We are not trying to generate an entire event ($O(1000)$ particles) all at once - it would be **very hard to validate!** Instead, generate a single particle shower (before electronics) and appeal to combinatorics.



Factorization

We are not trying to read out every event ($O(1000)$ per event would be **very** hard to generate a single channel of electronics) and a

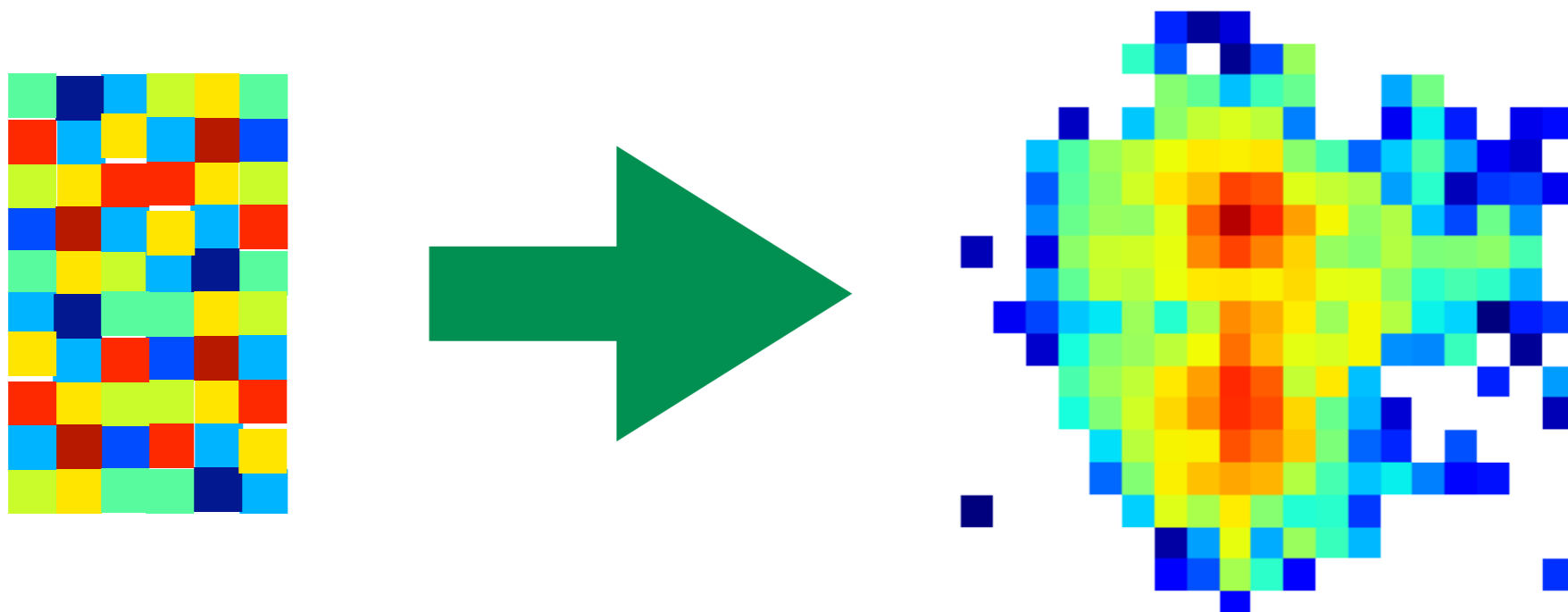
N.B. **calorimeter energy deposits factorize** (sum of the deposits is the deposit of the sum) but **digitization (w/ noise) does not!**



Now to the machine learning

11

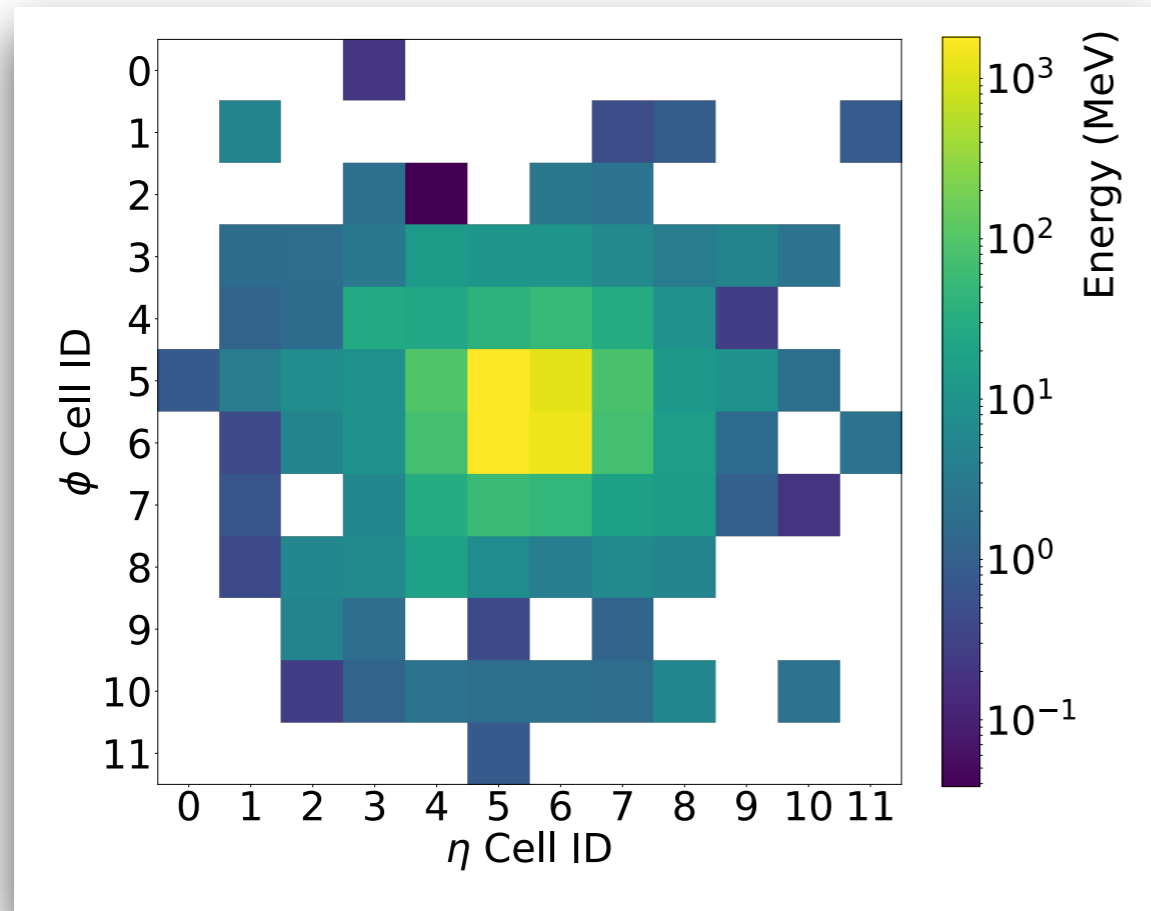
A **generator** is nothing other than a function that maps random numbers to structure.



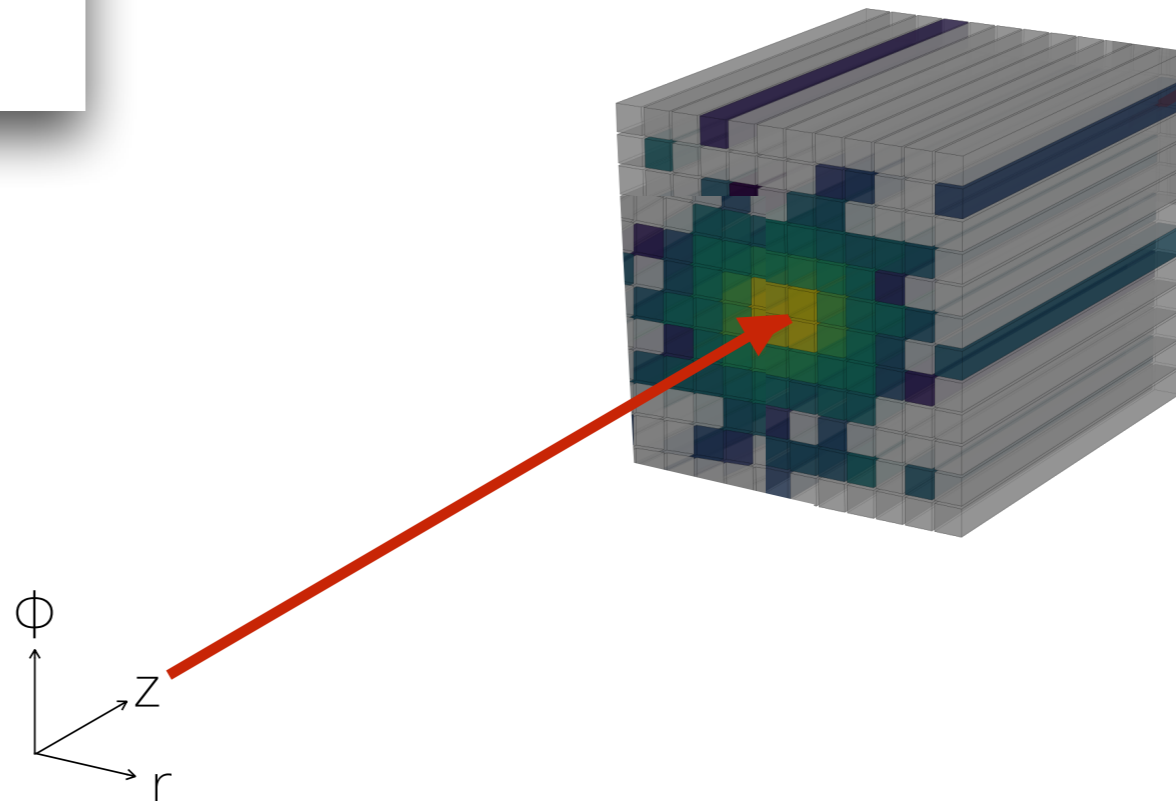
Our structure: **calorimeter images**

Calorimeter images

12



Grayscale images:
Pixel intensity =
energy deposited

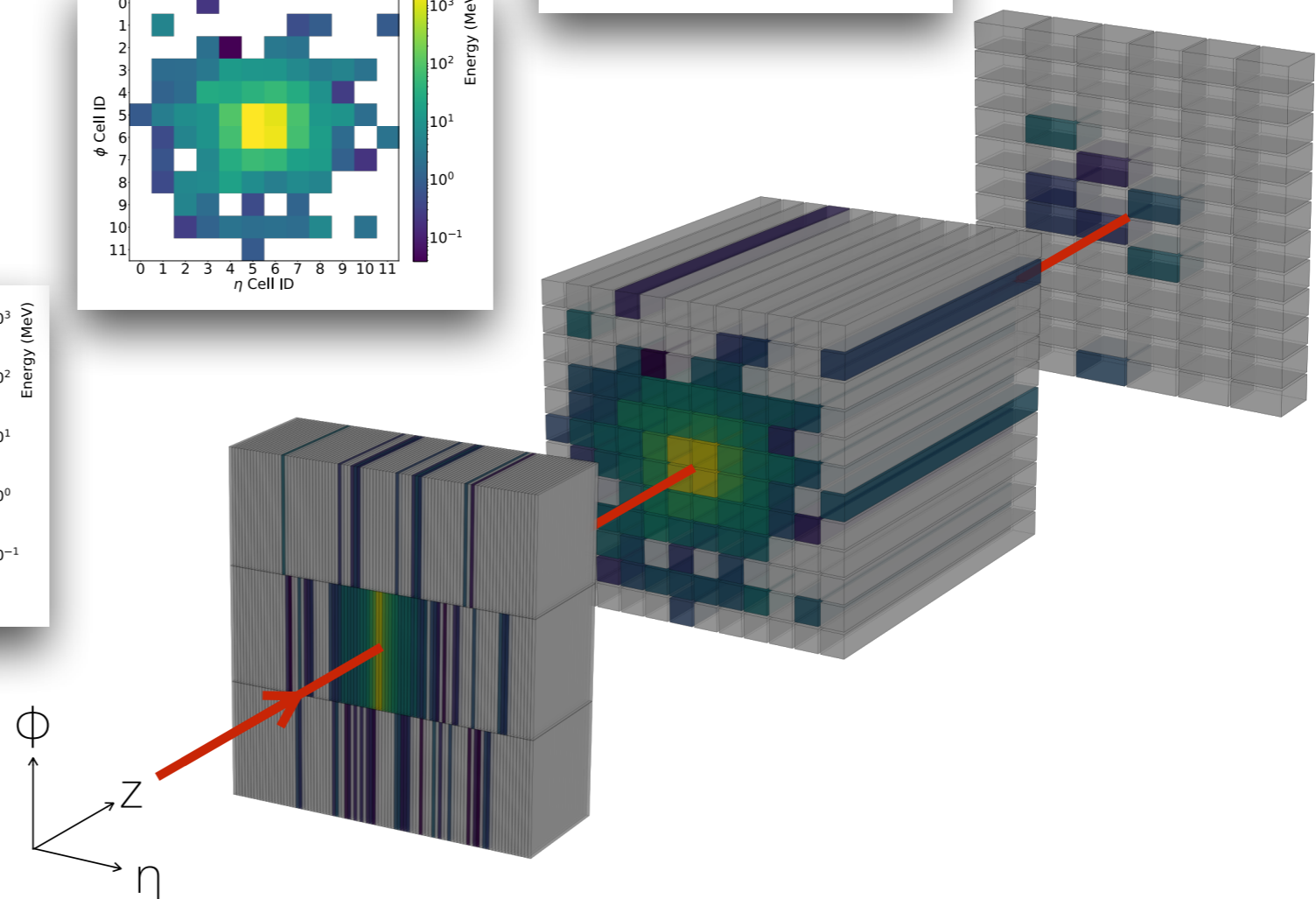
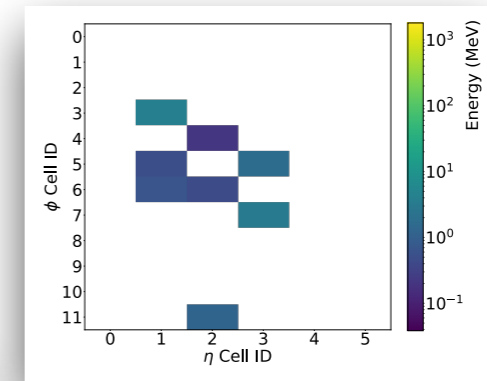
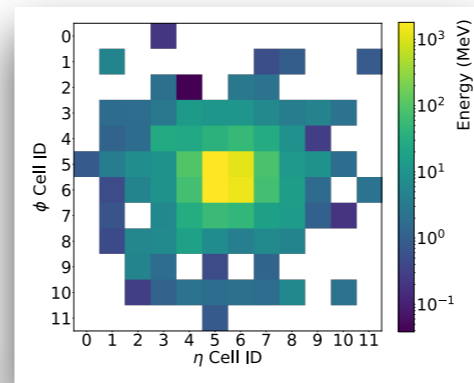
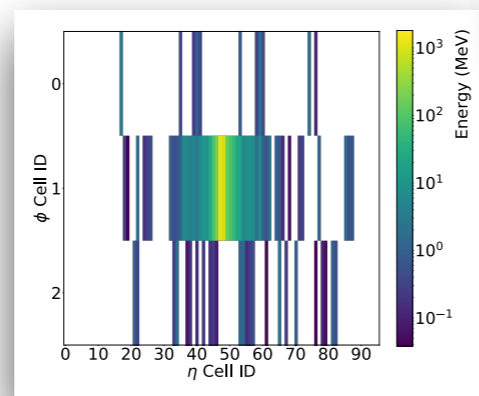


Calorimeter images

13

Challenge: **multiple layers**
with **non-uniform granularity**
and a **causal relationship**?

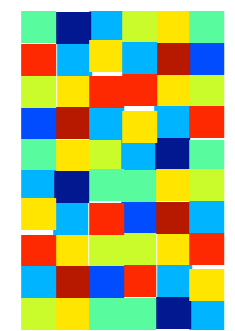
*N.B. images are
 $O(1000)$ dimensional*



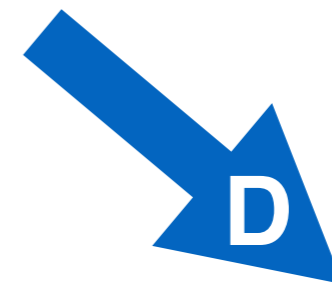
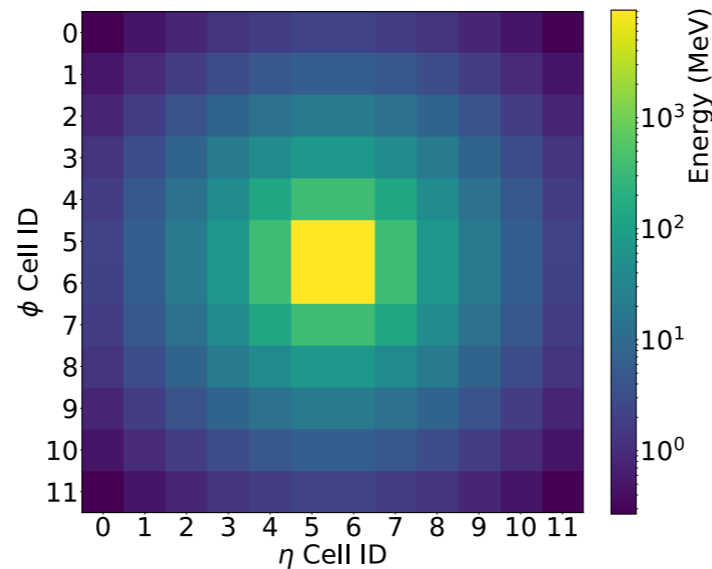
One possibility: GANs

14

Generative Adversarial Networks (GAN):
*A two-network game where one **maps noise to images** and one **classifies images as fake or real**.*

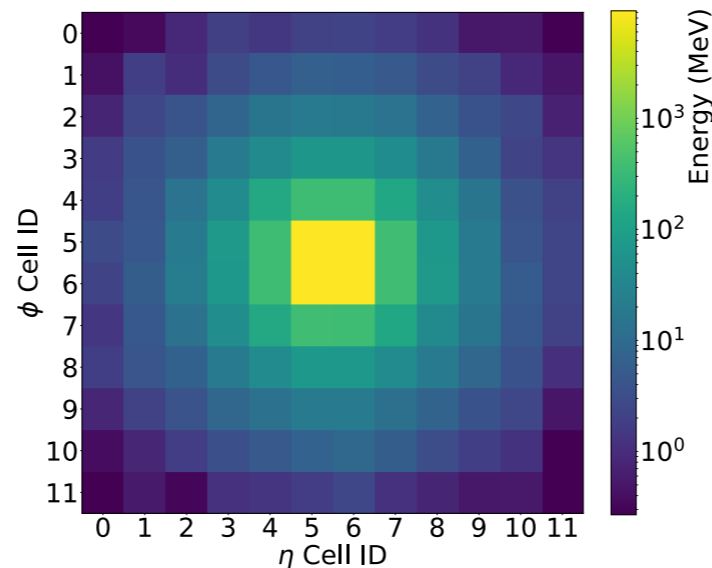


noise



{real, fake}

When **D** is maximally confused, **G** will be a good generator

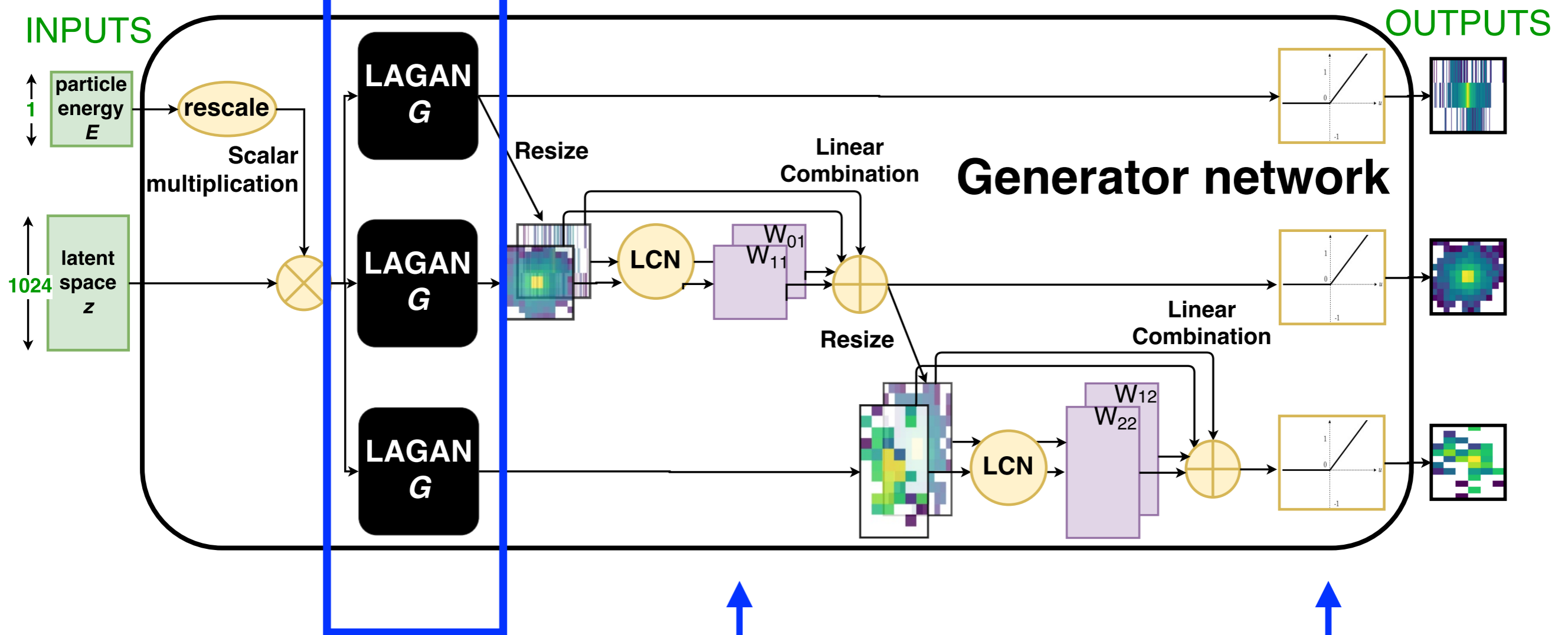


Physics-based simulator

Introducing CaloGAN

One image per calo layer

One network per particle type; input particle energy



use layer i as input to layer $i+1$

ReLU to encourage sparsity

Introducing CaloGAN

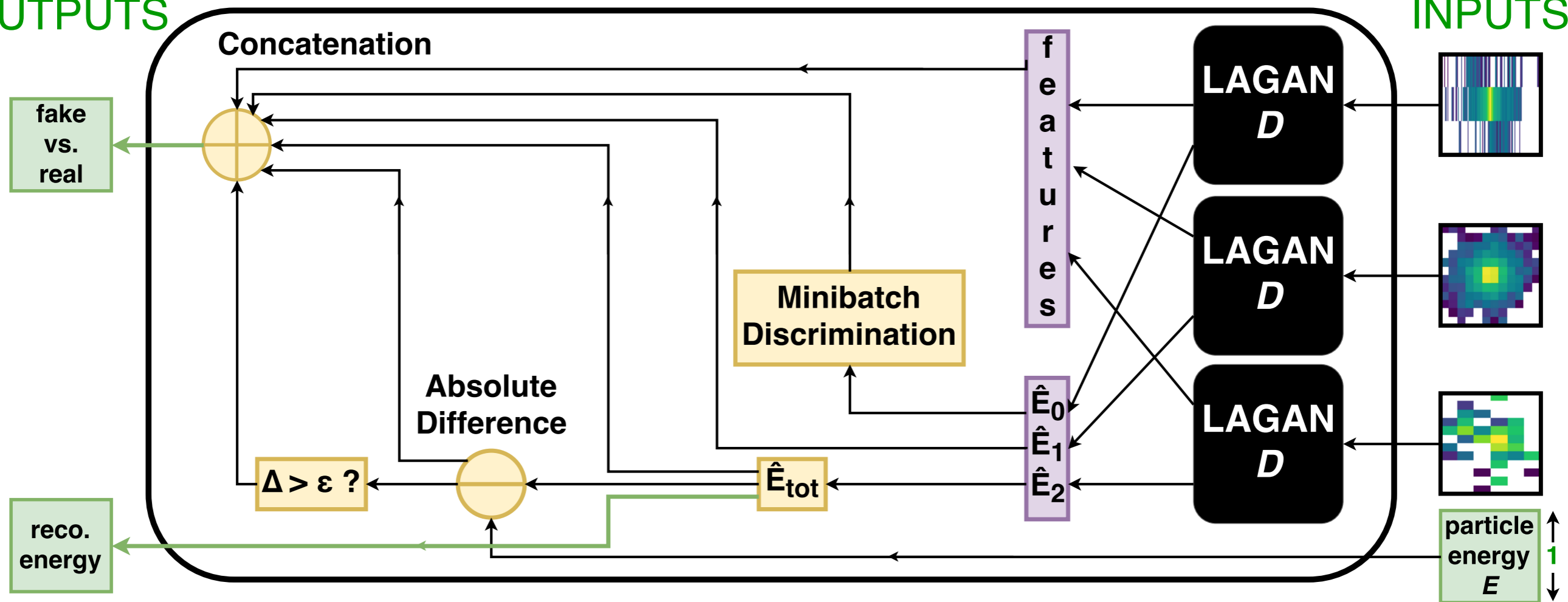
Mode collapse: learns to generate one part of the distribution well, but leaves out other parts.

help avoid 'mode collapse'



OUTPUTS

INPUTS

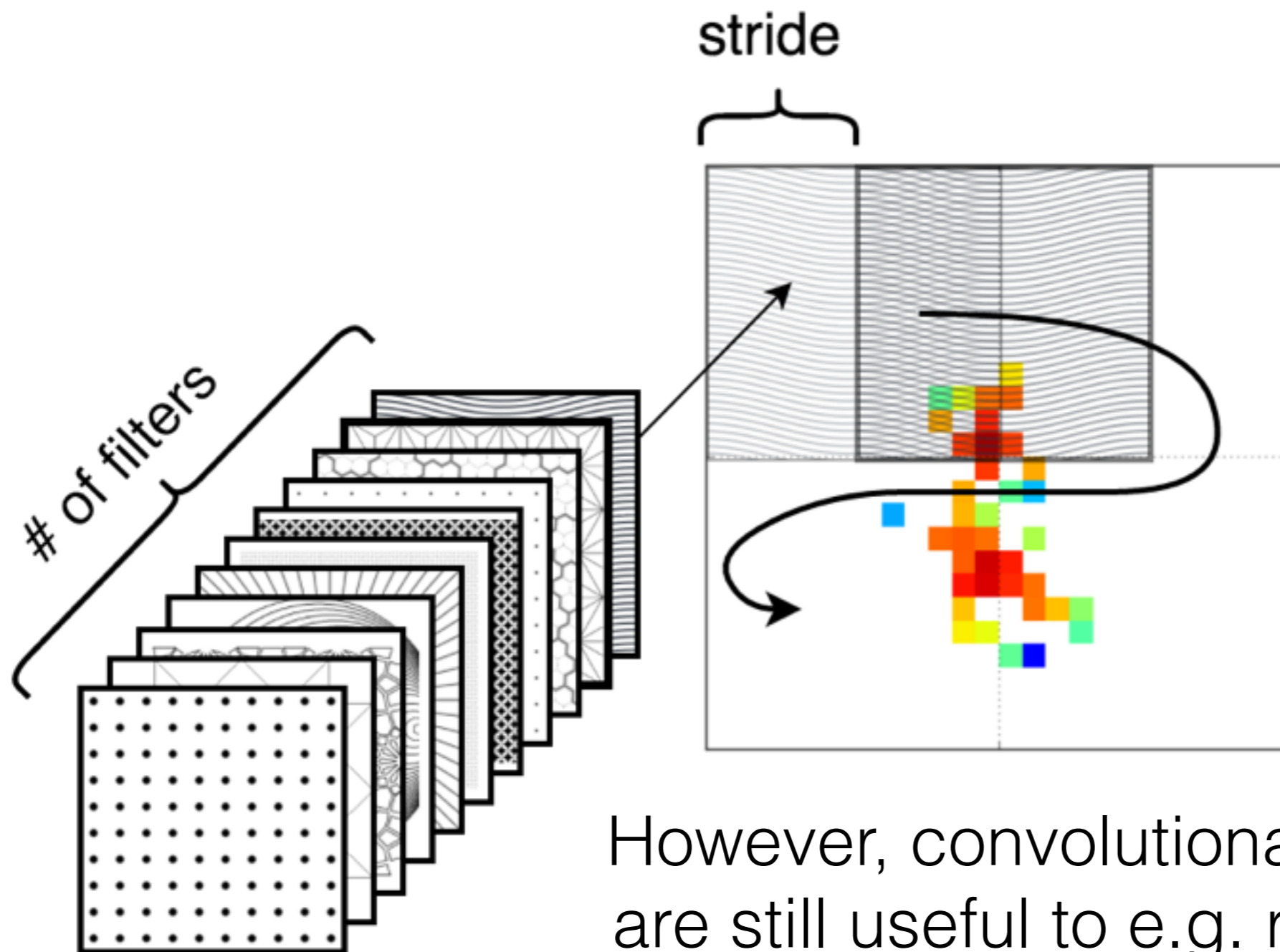


Discriminator network

Locally connected layers

Due to the structure of the problem, we do not have translation invariance.

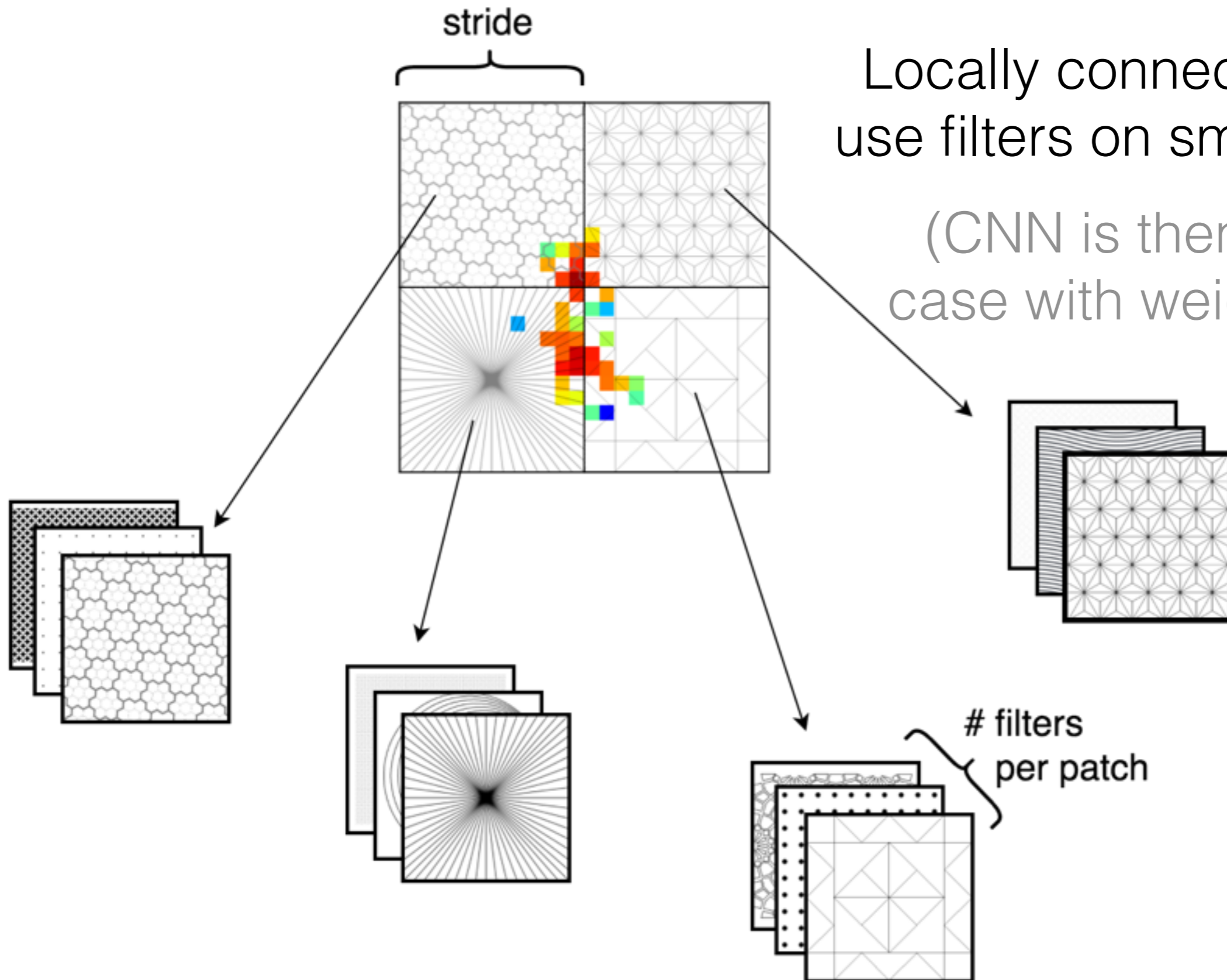
Classification studies found fully connected networks outperformed CNNs



However, convolutional-like architectures are still useful to e.g. reduce parameters

Locally connected layers

18



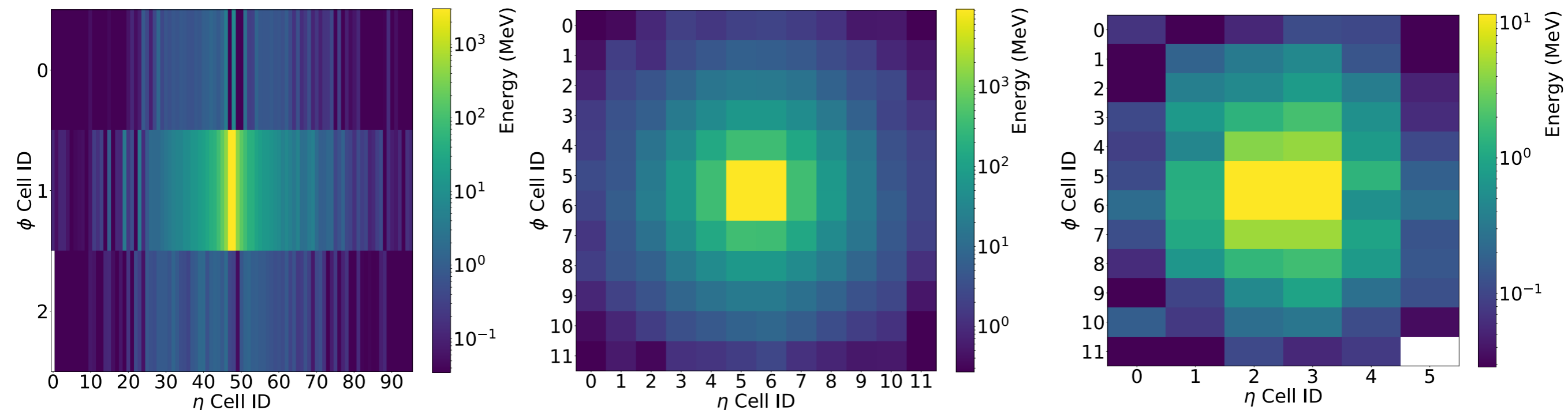
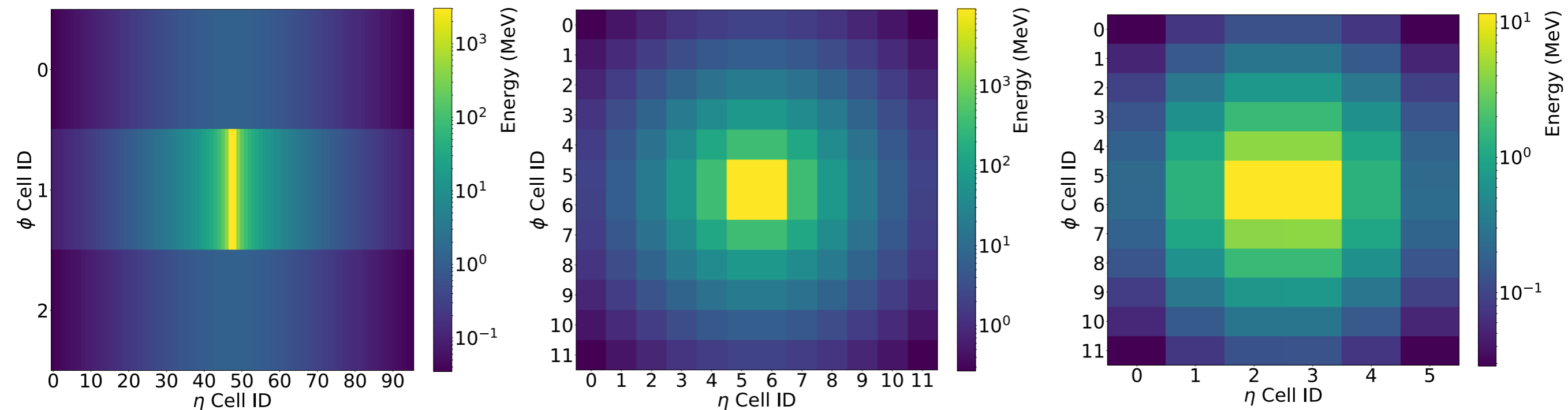
Locally connected layers use filters on small patches

(CNN is then a special case with weight sharing)

filters
per patch

Results: average images

Geant4

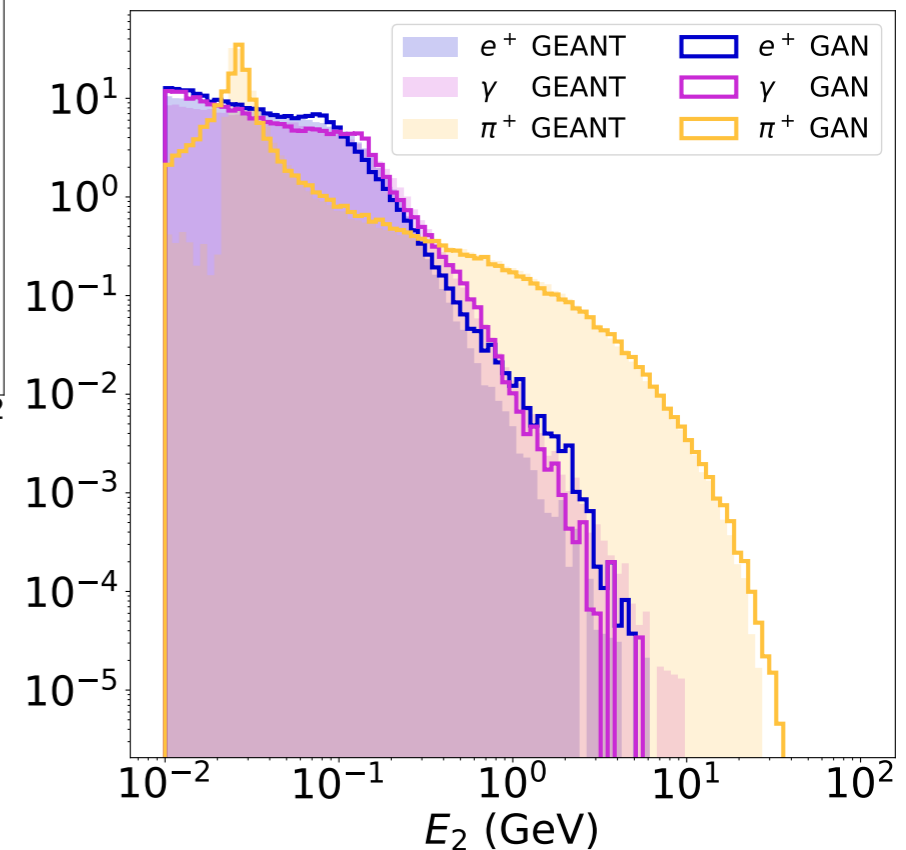
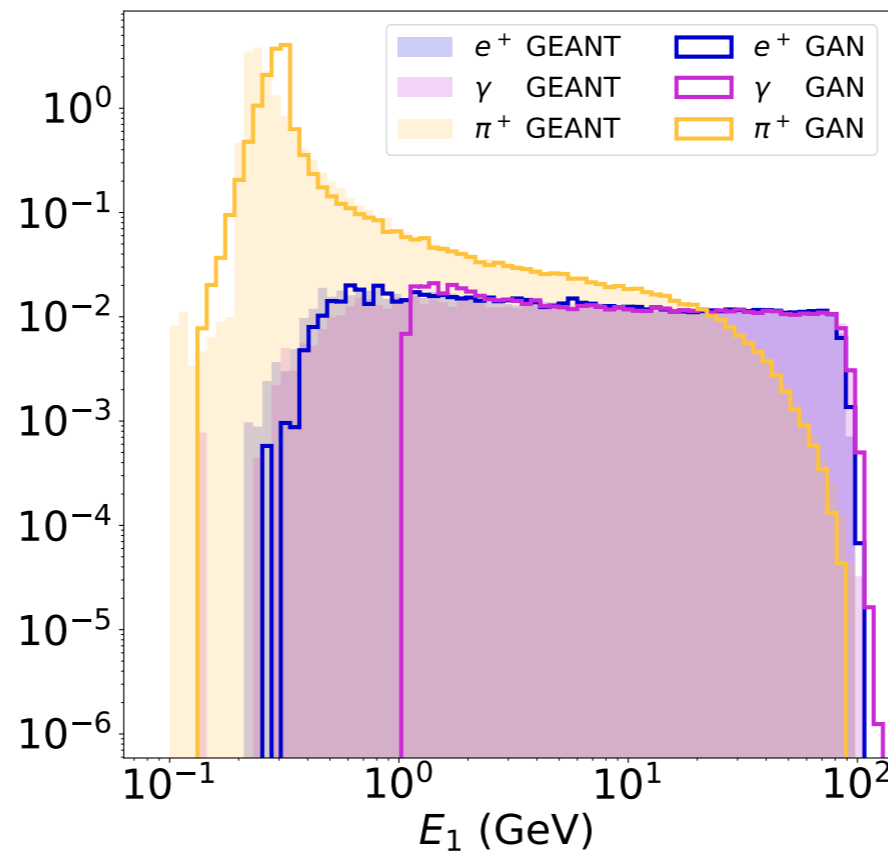
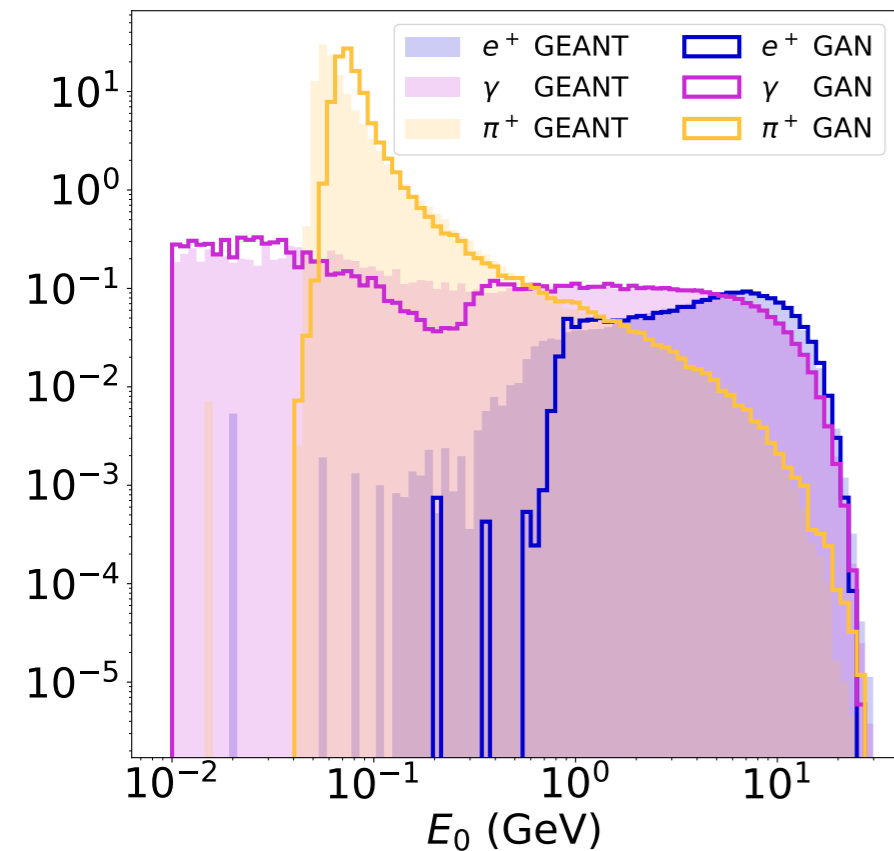


CaloGAN

Energy per layer

20

Pions deposit much less energy in the first layers; leave the calorimeter with significant energy



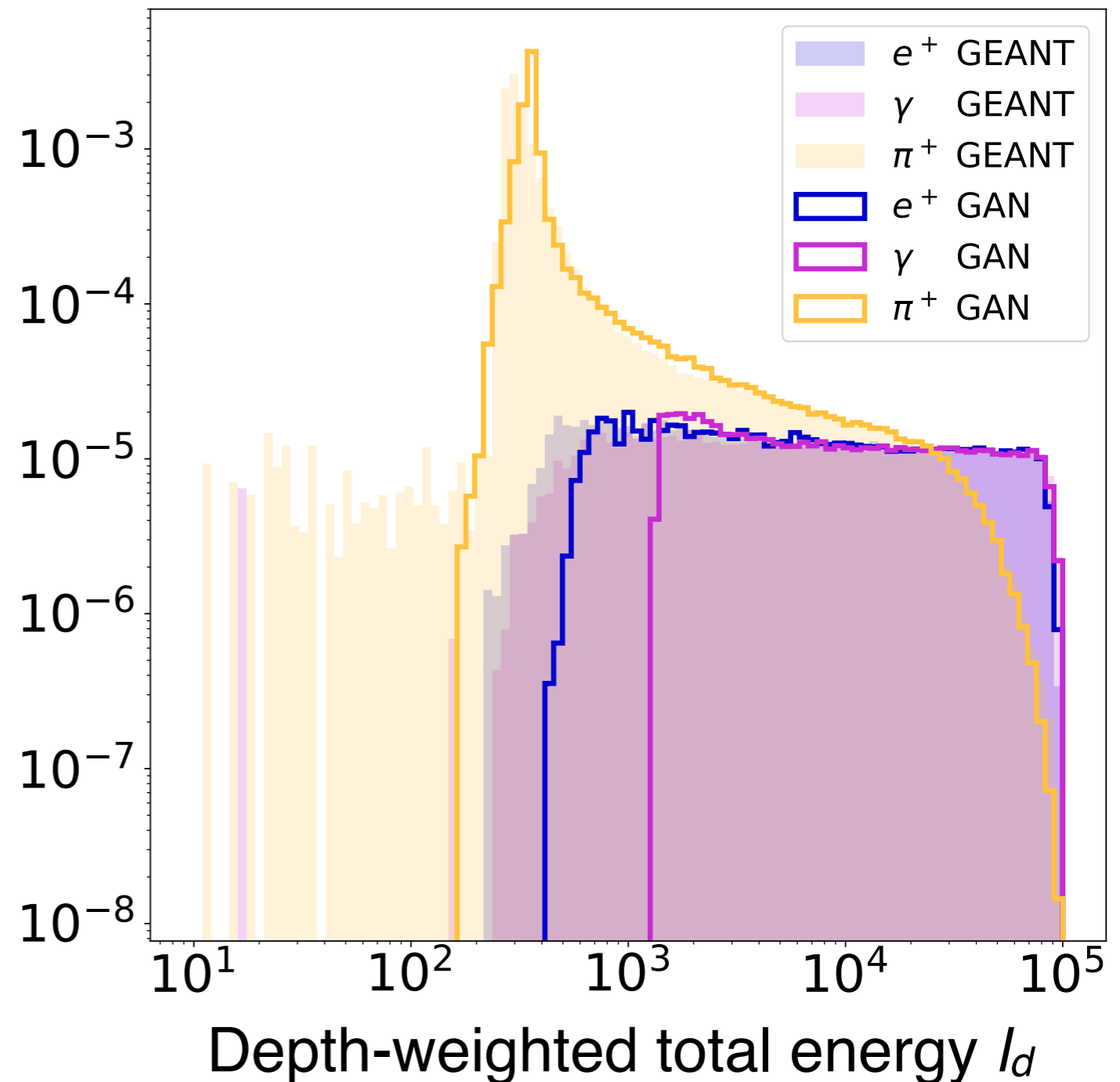
N.B. can always add these (and others) explicitly to the training

Warning: challenge with GANs

21

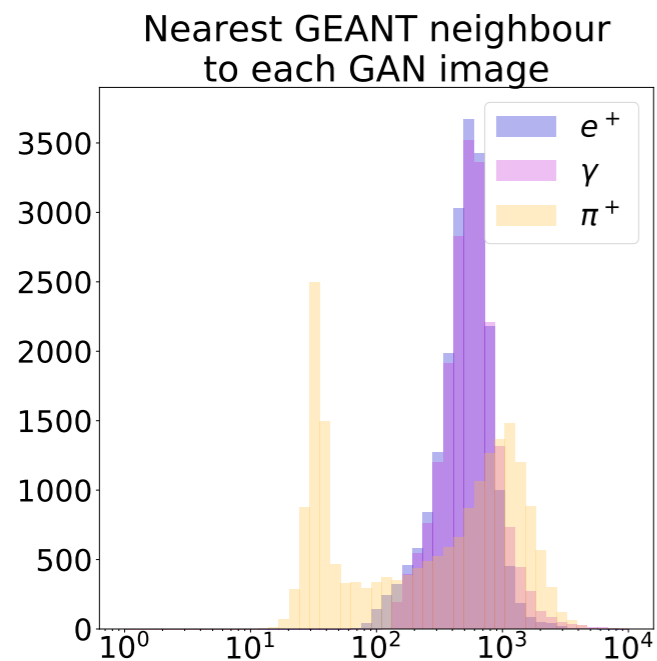
Unlike for classifiers, it is not easy to figure out which GAN is a good GAN - trying to learn a $O(1000)$ generative model and not a single likelihood ratio!

...this is a place where **science applications can make a big impact on ML.**

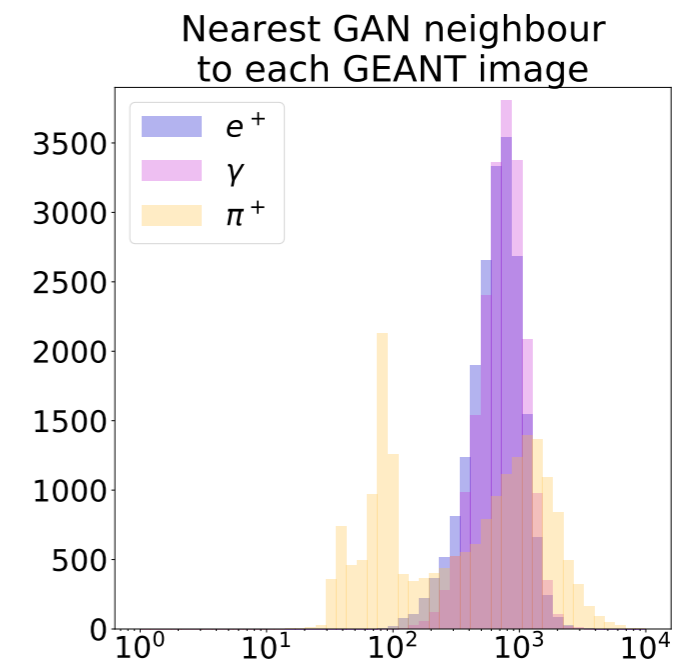


“Overtraining”

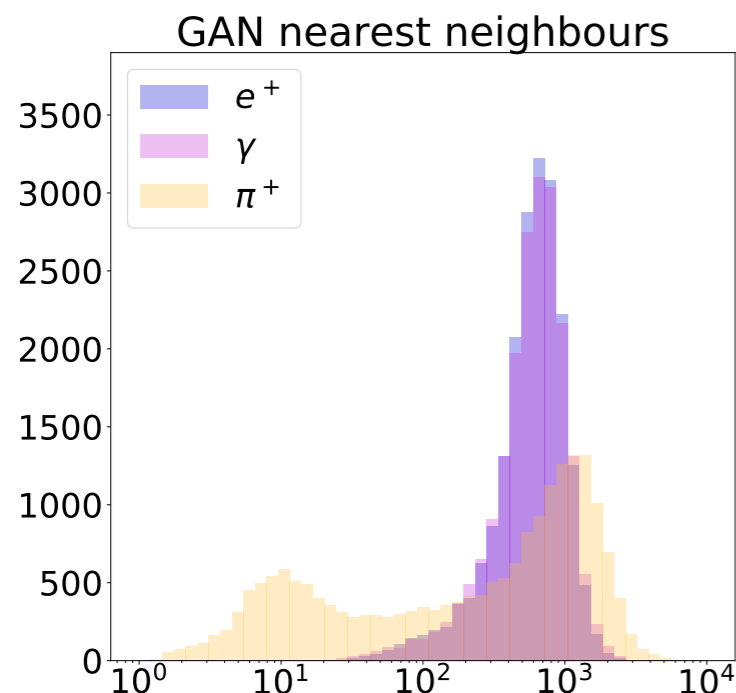
22



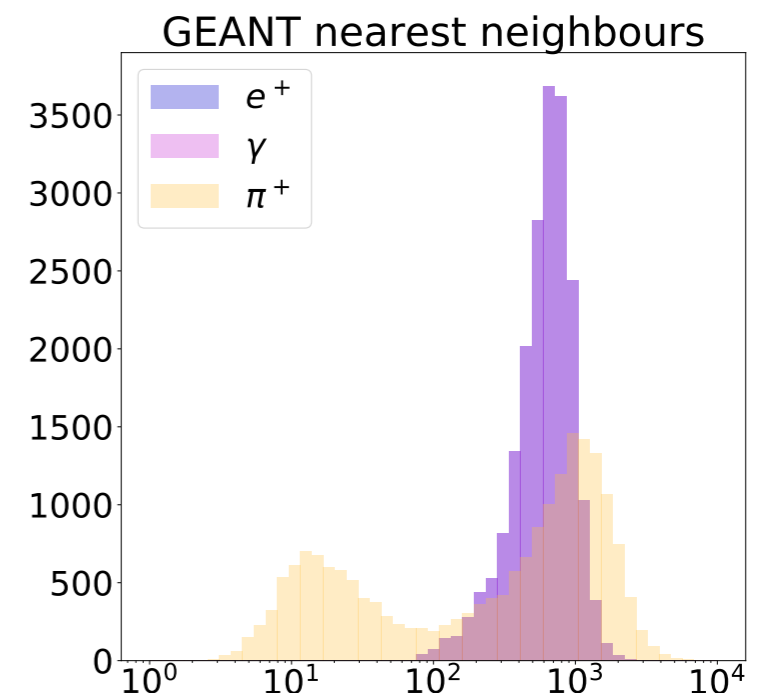
~not
memorizing



A key challenge in training GANs is the diversity of generated images. This does not seem to be a (big) problem for CaloGAN.

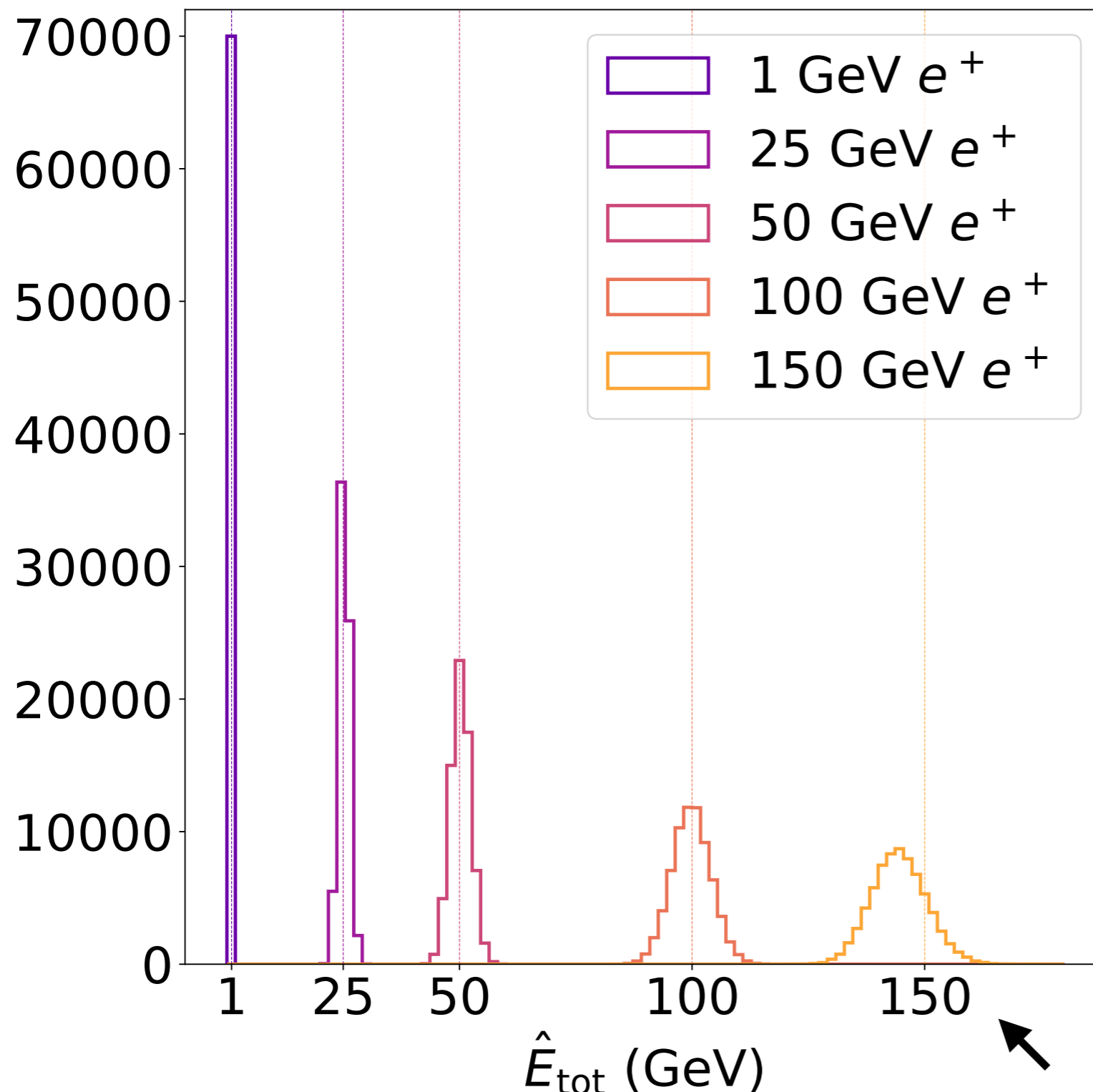


~no mode
collapse



Extrapolating

23



GANs are not designed to extrapolate, but in some cases, they can smoothly go on!

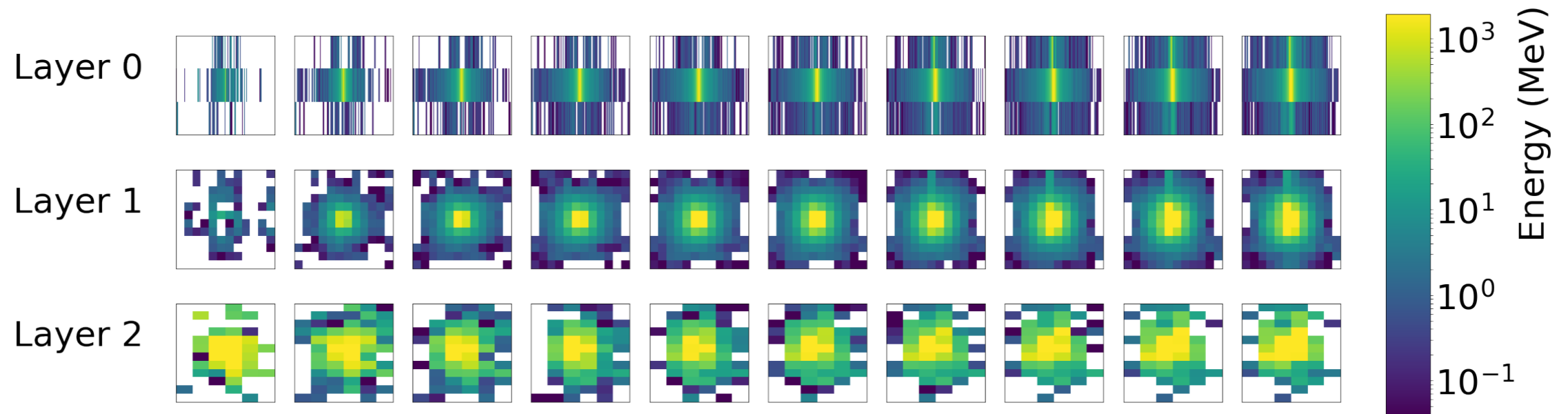
works here until there is no new physical principles which turn on at some energy

Beyond our training sample!

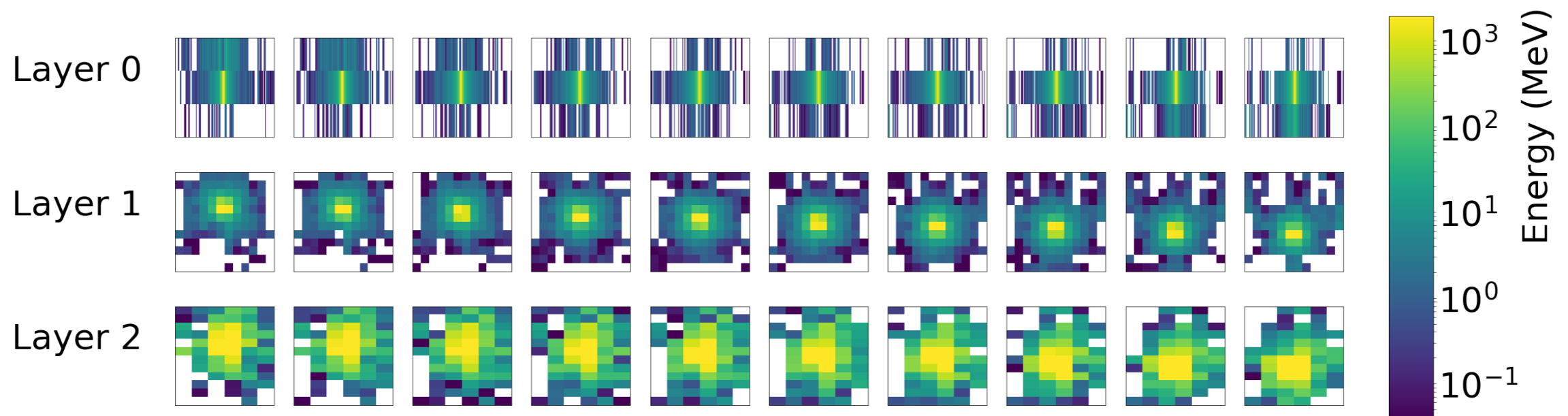
Conditioning

24

Fix noise, scan latent variable corresponding to energy



Fix noise, scan latent variable corresponding to x-position



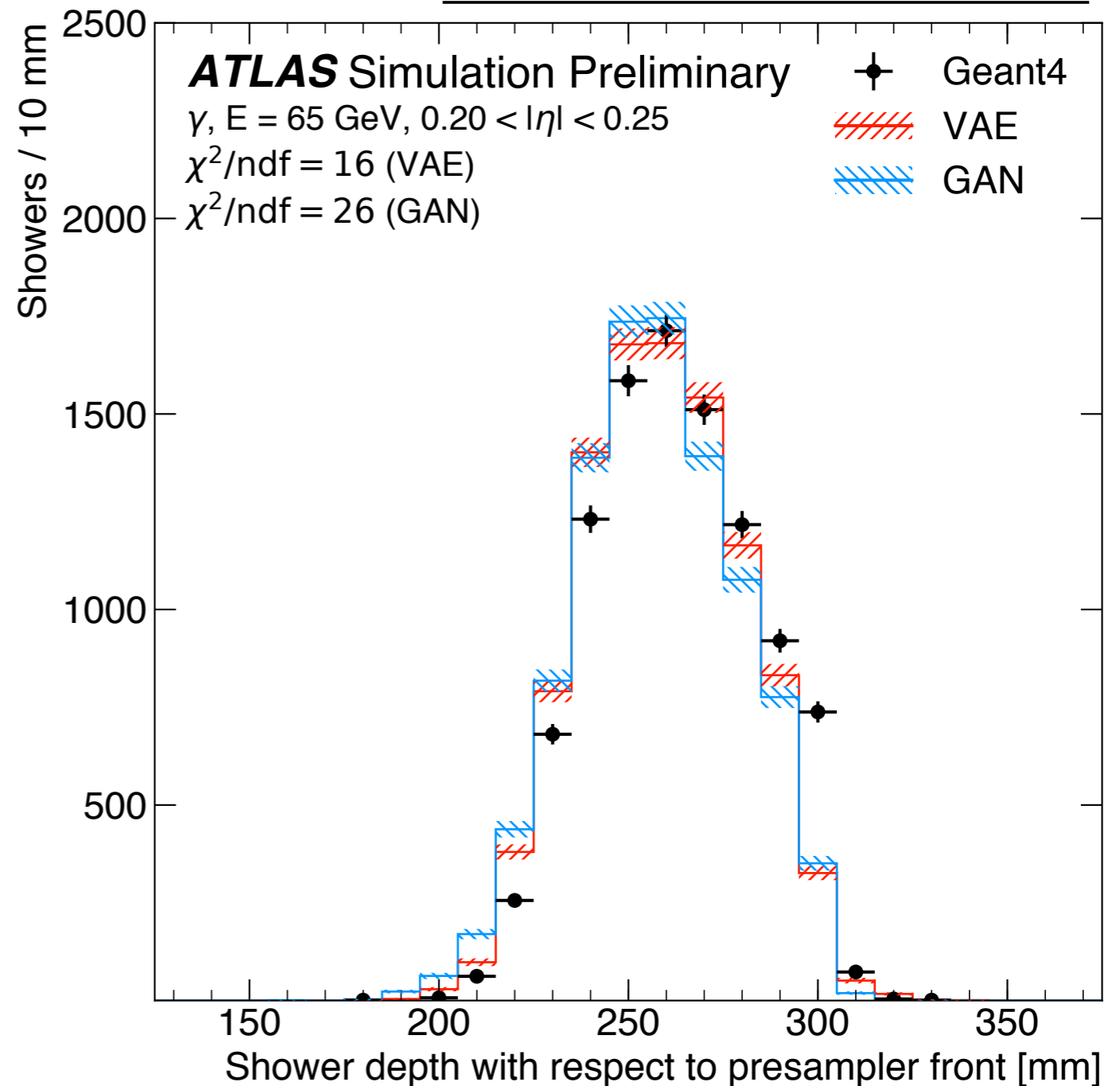
Timing

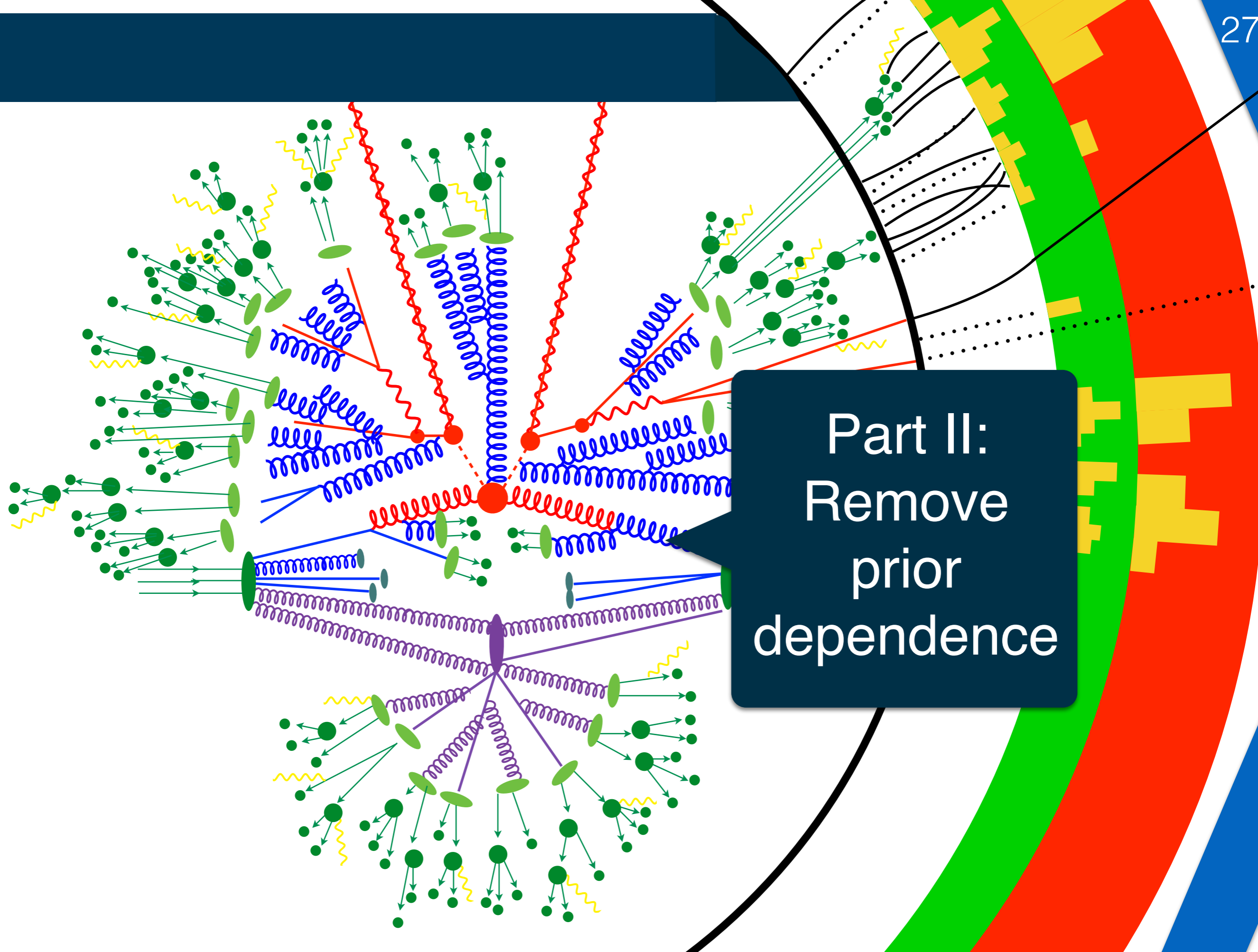
Generation Method	Hardware	Batch Size	milliseconds/shower
GEANT4	CPU	N/A	1772 ←
CALOGAN	CPU <i>Intel Xeon E5-2670</i>	1	13.1
		10	5.11
		128	2.19
		1024	2.03
	GPU <i>NVIDIA K80</i>	1	14.5
		4	3.68
		128	0.021
		512	0.014
		1024	0.012 ←

(clearly these numbers will change as both technologies improve - this is simply meant to be qualitative and motivating!)

Integrating these techniques into a full detector simulation is another layer of complication, but is possible and hopefully worth the effort!

ATL-SOFT-PUB-2018-001





Part II:
Remove
prior
dependence

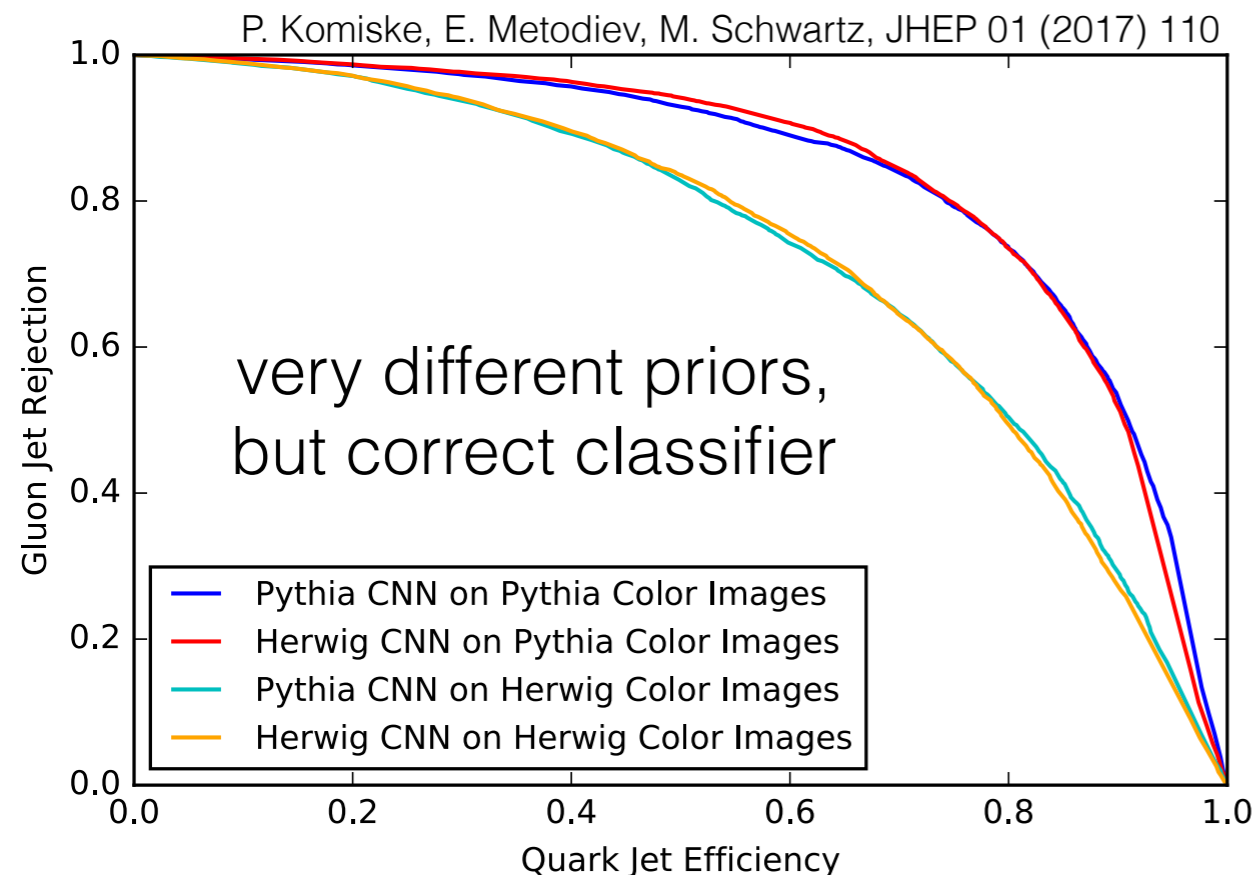
Simulation dependence + regression

28

One source of MC dependence is the same as classification:

→ mis-modeling dependencies between features

However, there is a new source:
dependence on the feature priors.

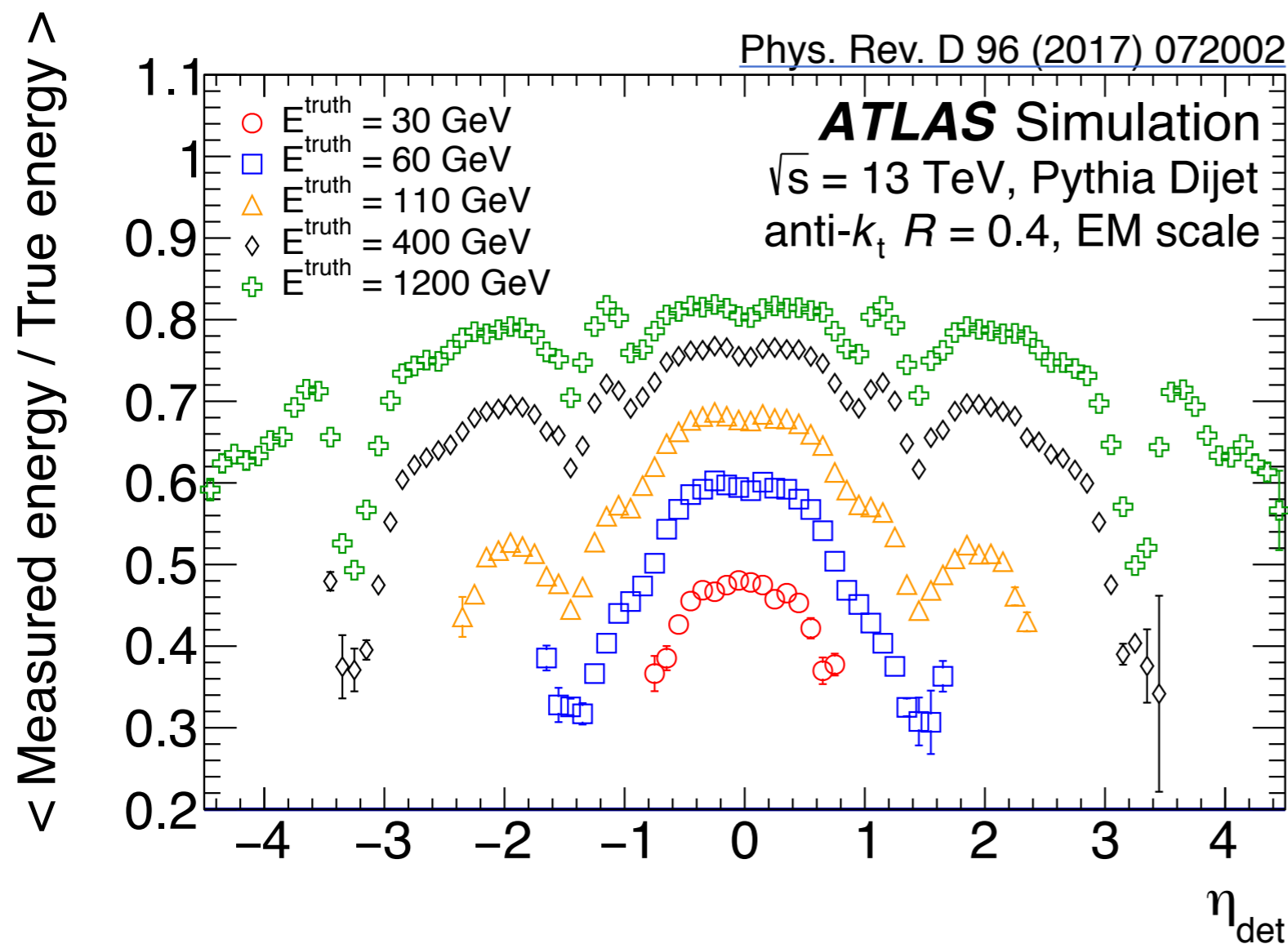


This is really new for regression -
for classification, if $p(x|\text{signal})$ and $p(x|\text{background})$ are mis-modeled,
you get the optimal answer as
long as their ratio is correct.

Simulation dependence + regression

29

An example that you can have in mind is energy calibration.



We want to predict the true energy given the measured energy

(and possibly other features - more on that soon)

...however what I'm about to say applied more generally (though the impact is biggest when the resolution is poorest)

What can go wrong?

30

Suppose you have some features x and you want to predict y .

detector energy

true energy

One way to do this is to find an f that minimizes the mean squared error:

$$f = \operatorname{argmin}_g \sum_i (g(x_i) - y_i)^2$$

Then, $f(x) = E[y|x]$.

*Proof in the lectures
from yesterday!*

Why is this a problem?

What can go wrong?

$$f(x) = E[y|x] = \int dy y p(y|x)$$

$$E[f(x)|y] = \int dx dy' y' p_{\text{train}}(y'|x) p_{\text{test}}(x|y)$$

this need not be y even if $p_{\text{train}} = p_{\text{test}}$ (!)

One solution: Numerical inversion

32

ATLAS and CMS use a trick to be prior-independent:

Numerical inversion *instead of predicting y from x , predict x from y and then invert the function*

... put another way:

learn $f: y \rightarrow x$ and then for a given x , predict $f^{-1}(x)$

by construction, f is independent of $p(y)$ and thus f^{-1} also does not depend on $p(y)$, as desired.

Caveats about numerical inversion

33

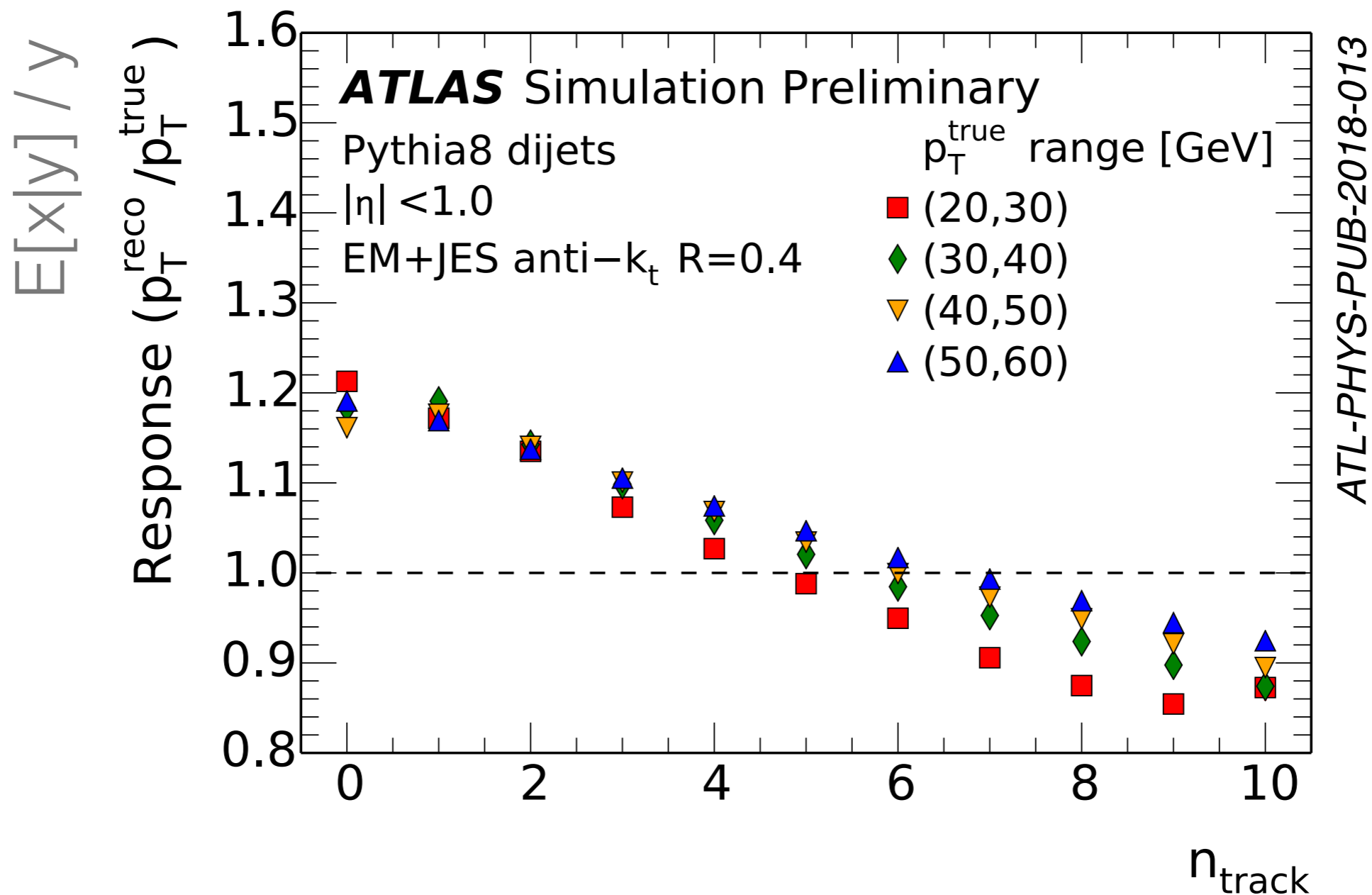
This procedure is independent of the prior $p(y)$ but may not close exactly, i.e. $E[f^{-1}(x)|y]$ may not be y .

...under mild assumptions, it does close for the mean absolute error, but usually has some non-closure for the squared error.

Also, the calibration procedure can distort the underlying distribution, i.e. if you start with a Gaussian, you almost never end up with exactly a Gaussian.

+ more features

The detector response of jets depends on many properties of the jet. Ideally, the calibration can include this!

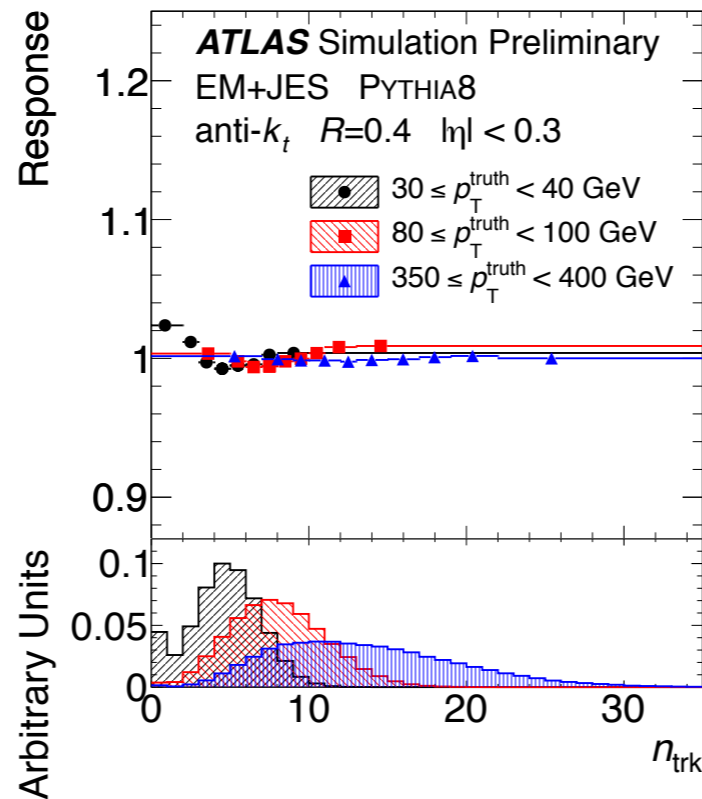
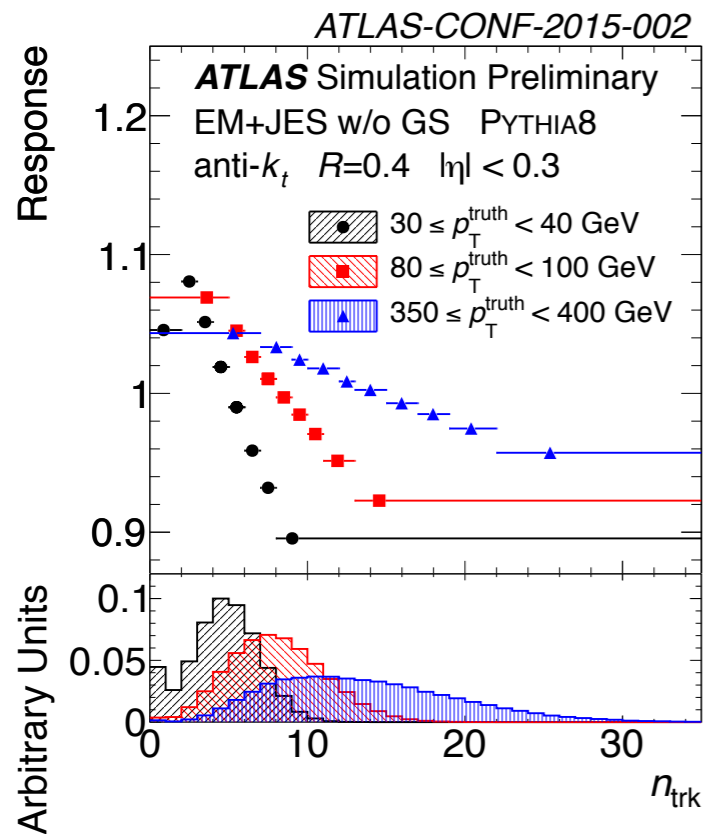


Global sequential calibration

35

The current ATLAS approach to including more features is to repeat NI sequentially:

$$p_T^{\text{reco}} \mapsto \hat{p}_T^{\text{reco}} = f_{\theta_n}^{-1} \left(\cdots f_{\theta_2}^{-1} \left(f_{\theta_1}^{-1} \left(p_T^{\text{reco}} \right) \right) \cdots \right)$$



This works well when the jet response is independent of θ_i given θ_j .

For reasons discussed earlier, we can't include correlations by learning y given x and all the θ 's.

However, it would still be great to use machine learning to automatically and efficiently make use of correlated information.

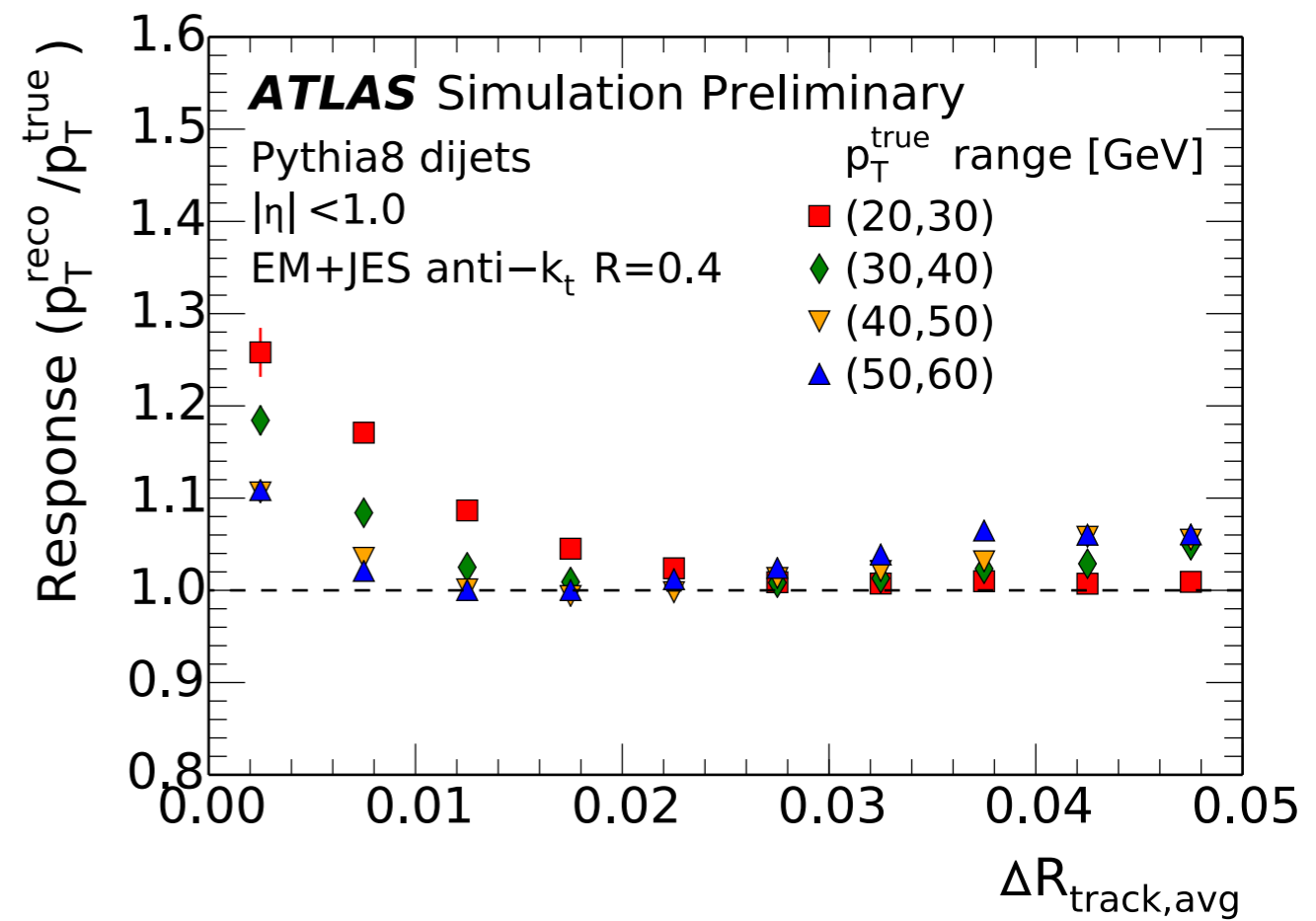
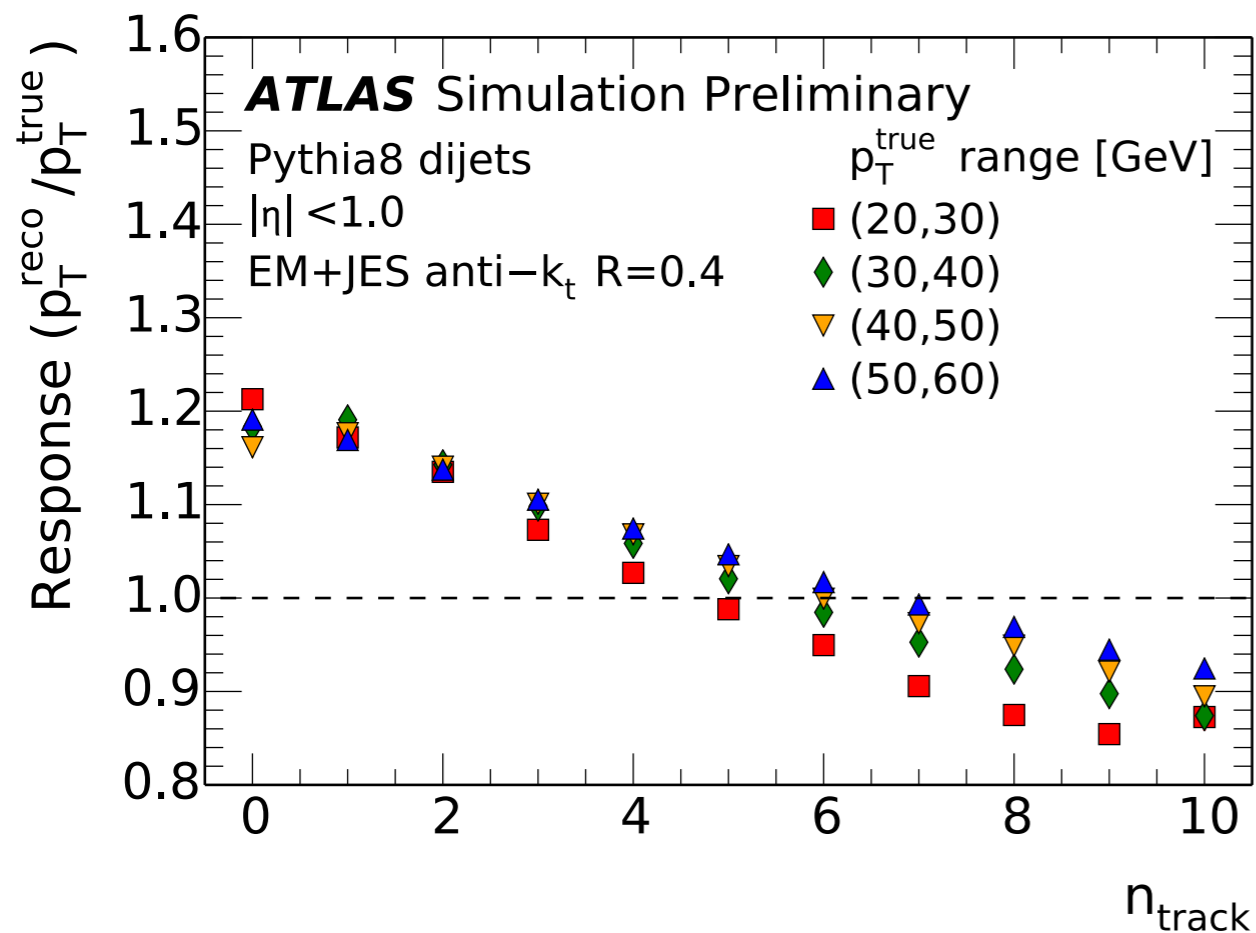
We cannot use numerical inversion out-of-the-box because we now have a many-to-one function.

Since we are not (necessarily) interested in calibrating the θ 's, we can generalize NI as follows:

- (1) Learn a function f to predict x given y and all the θ 's.
- (2) For every combination of θ , invert f .
- (3) Calibrate via $x \rightarrow f_{\theta}^{-1}(x)$

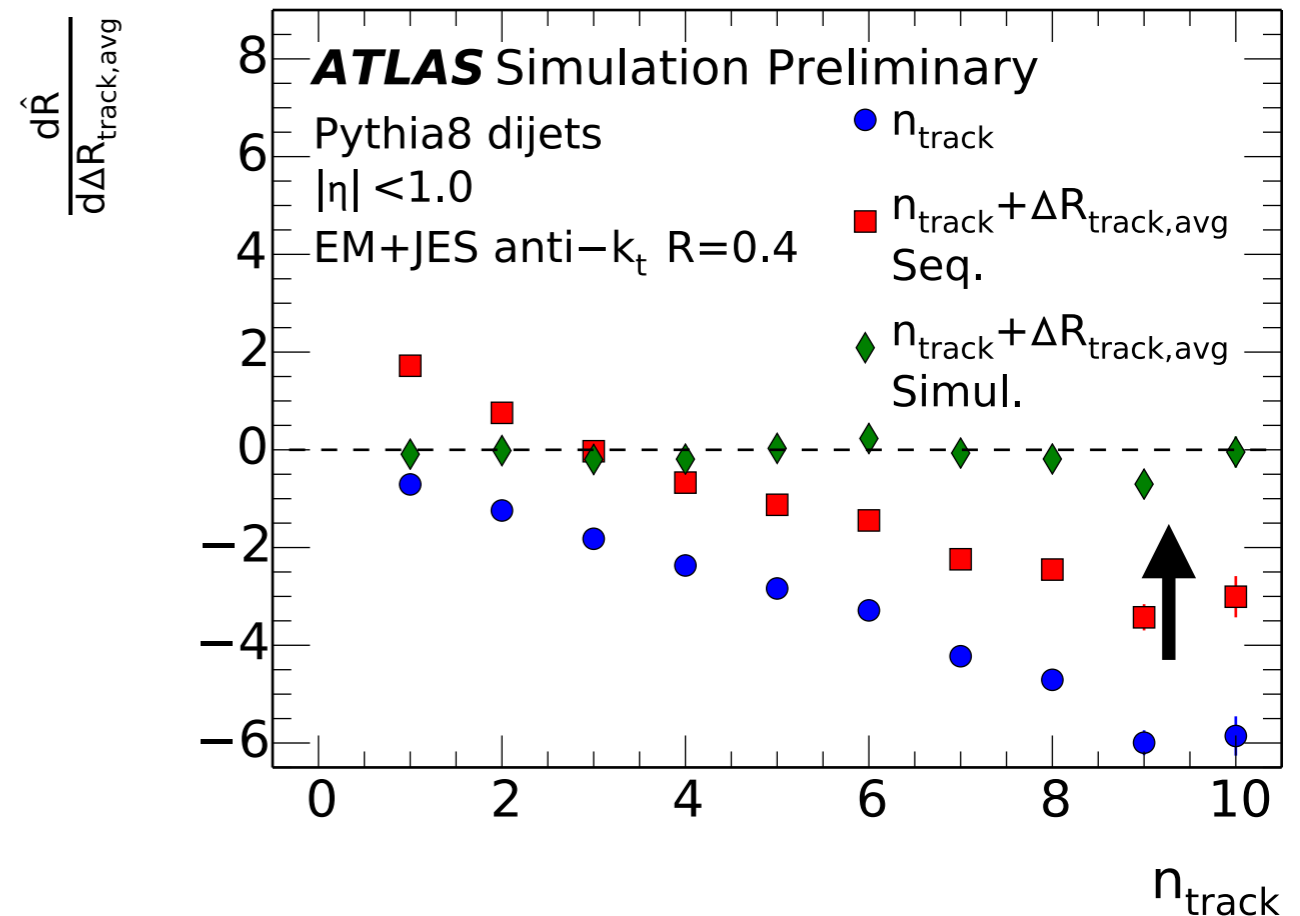
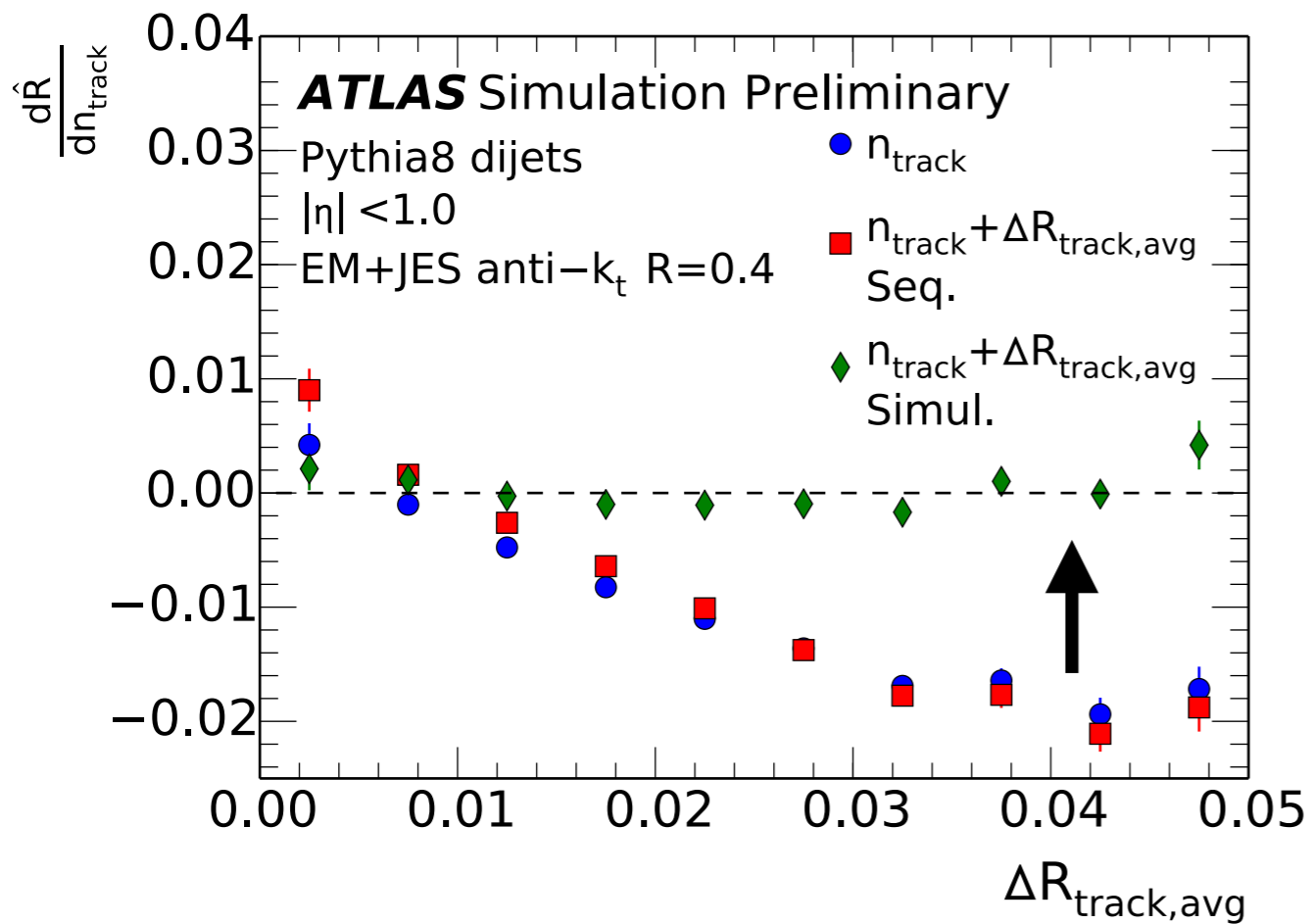
Step (2) is intractable, so replace it with another learning step: predict y given $f(y, \theta)$ and θ .

Consider two features:



average track p_T -weighted
distance from jet center
(doesn't really matter what it is)

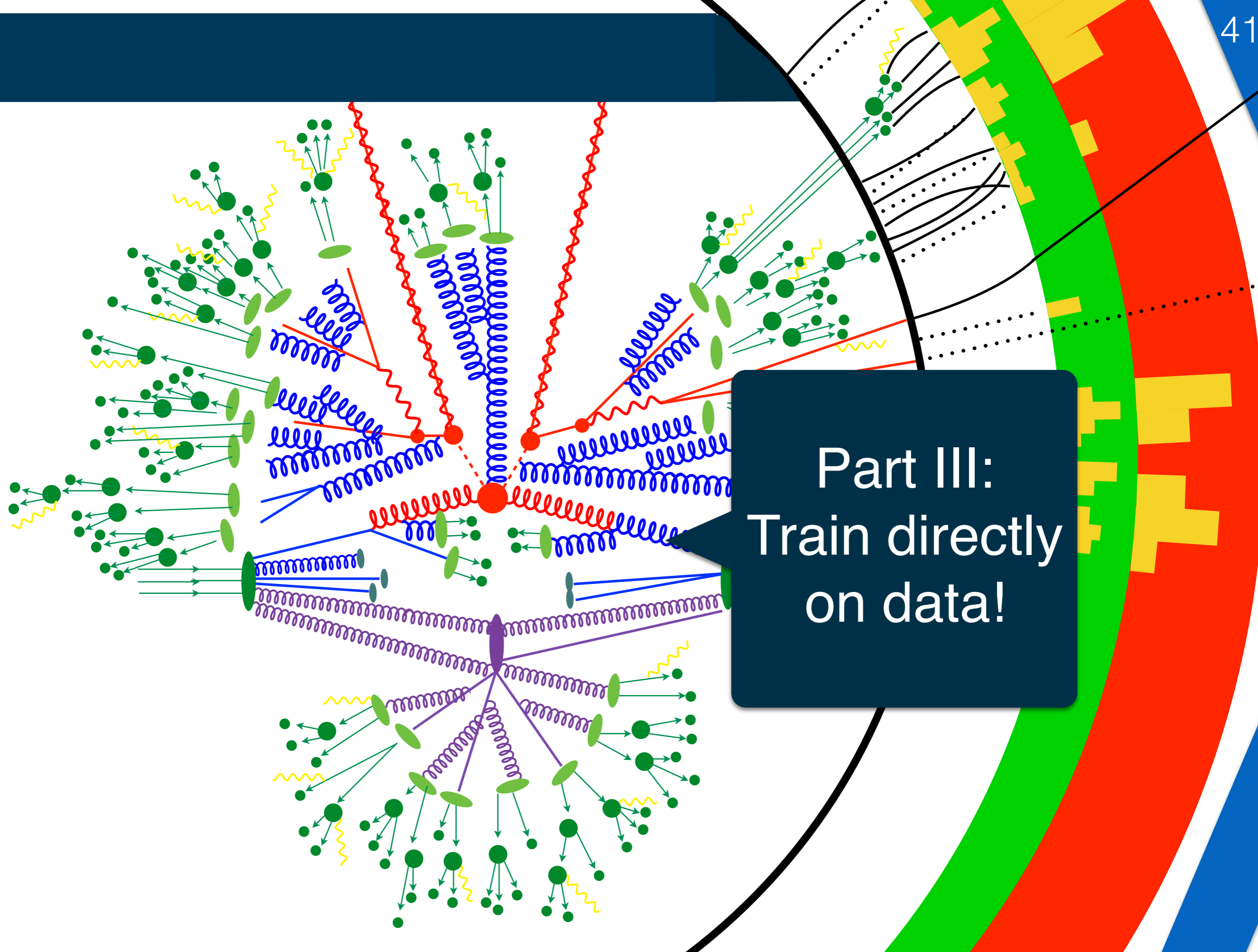
\hat{R} is the calibrated $E[x|y] / y$



Only the simultaneous approach removes the full residual dependence!

Adding more features (with more interdependencies) will lead to more dramatic improvements.

We can also extend this approach to calibrate other observables and even simultaneously calibrate some of the θ 's (even more generalized NI!)

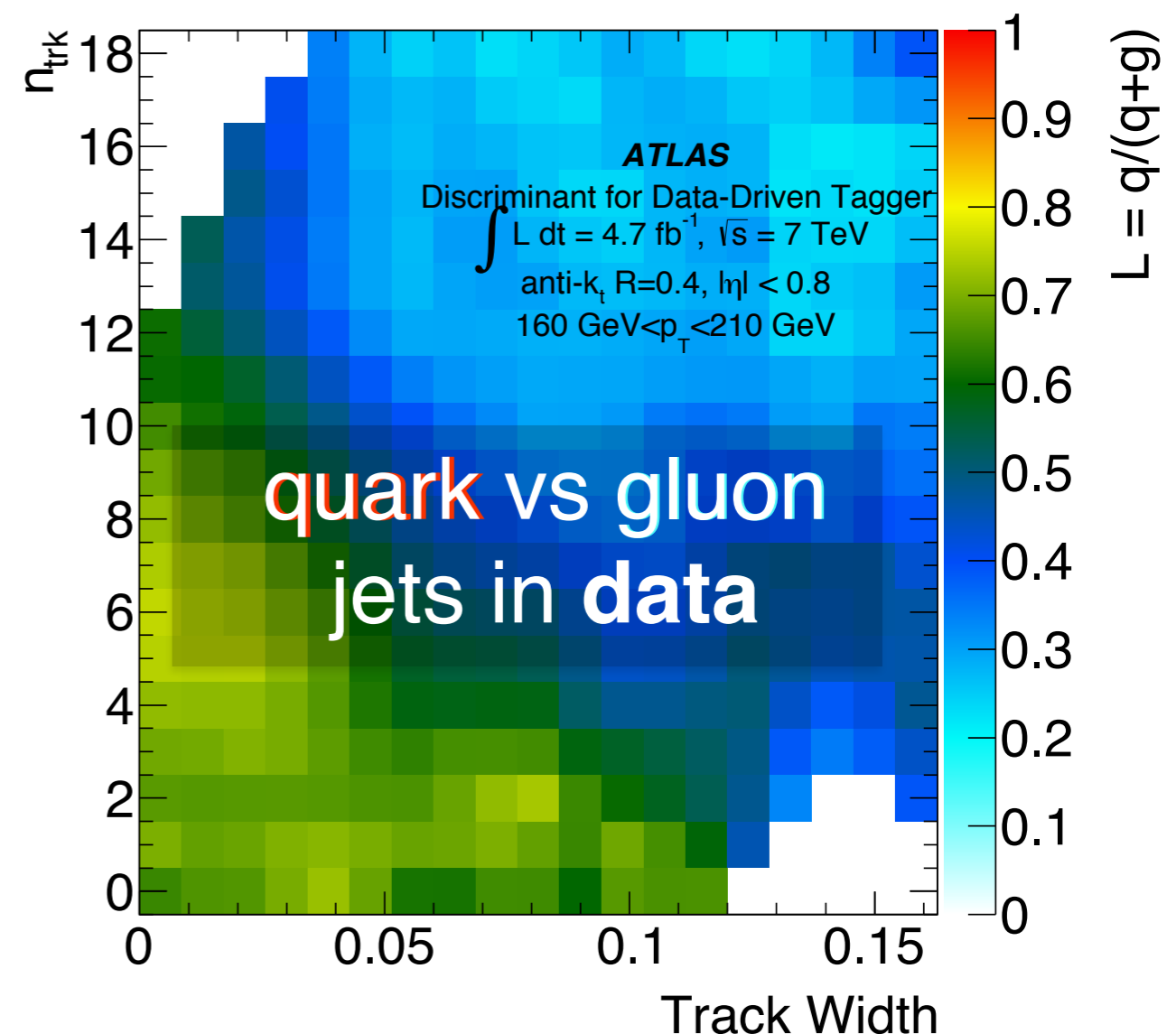
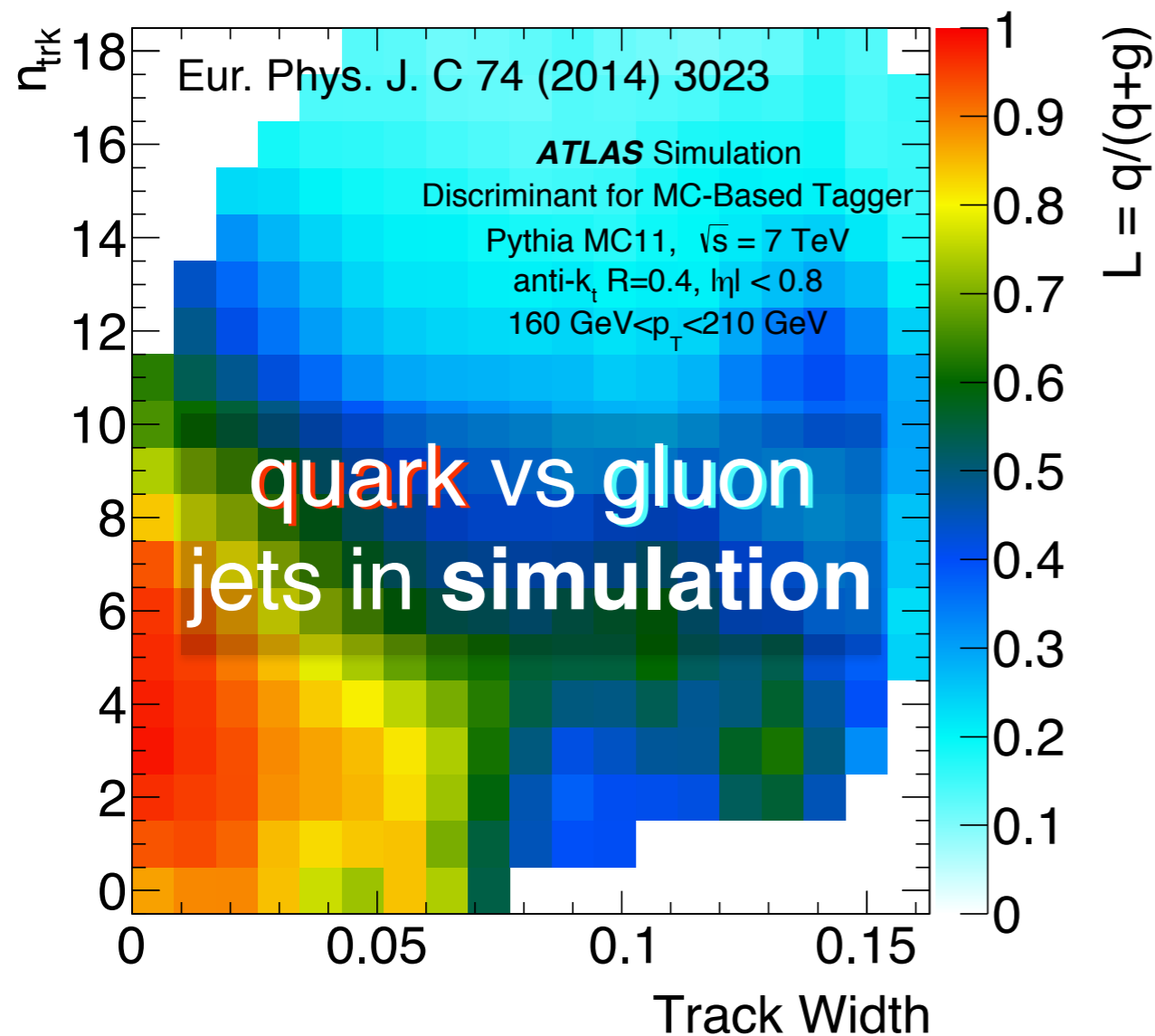


Part II:
Train directly
on data!

Background and Motivation

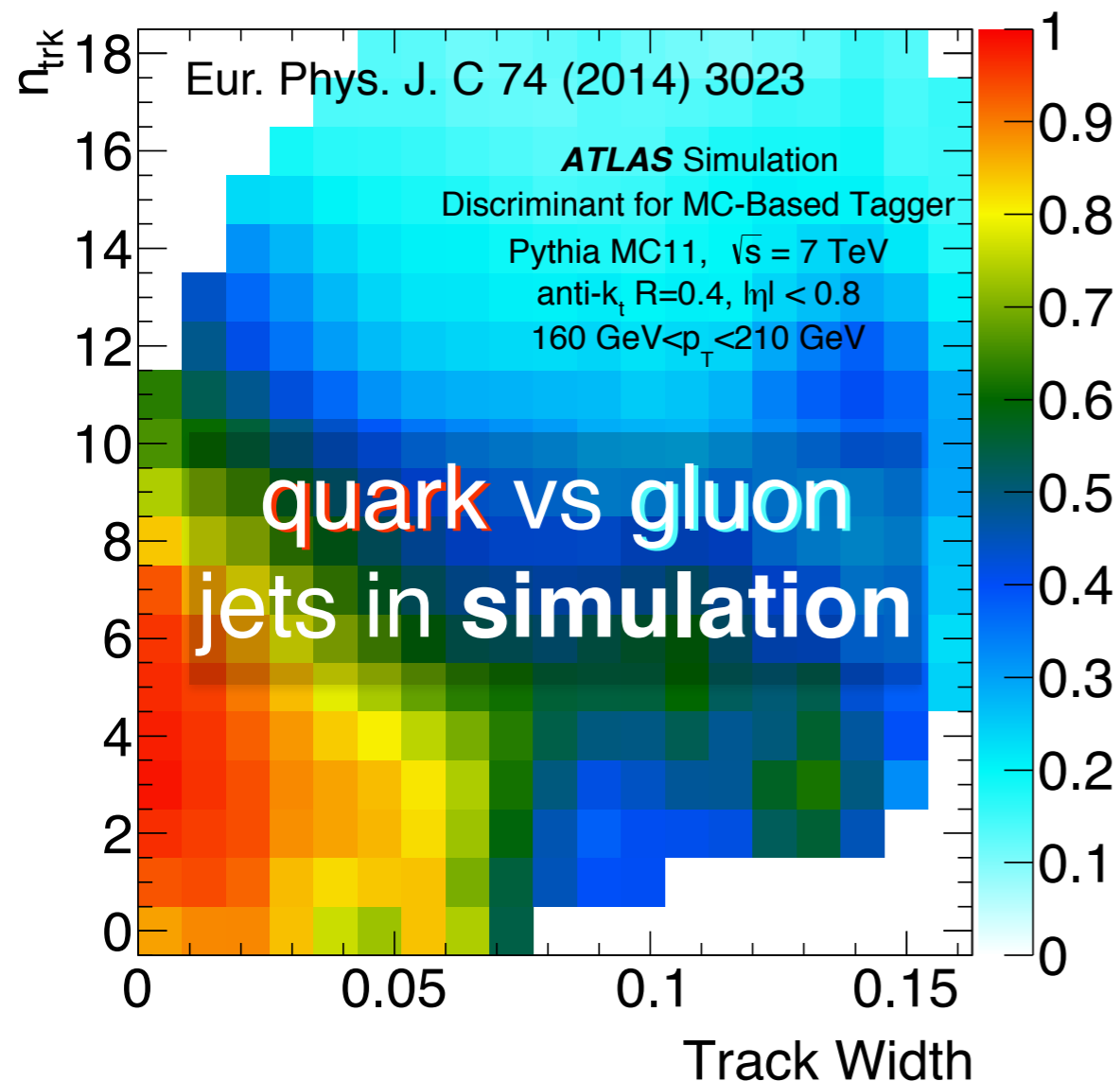
42

Usual paradigm: train in simulation, validate on data, test on data.



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

Usual paradigm: **train in simulation**, validate on data, test on data.



Recall: optimal classifier is a threshold cut on the likelihood ratio.

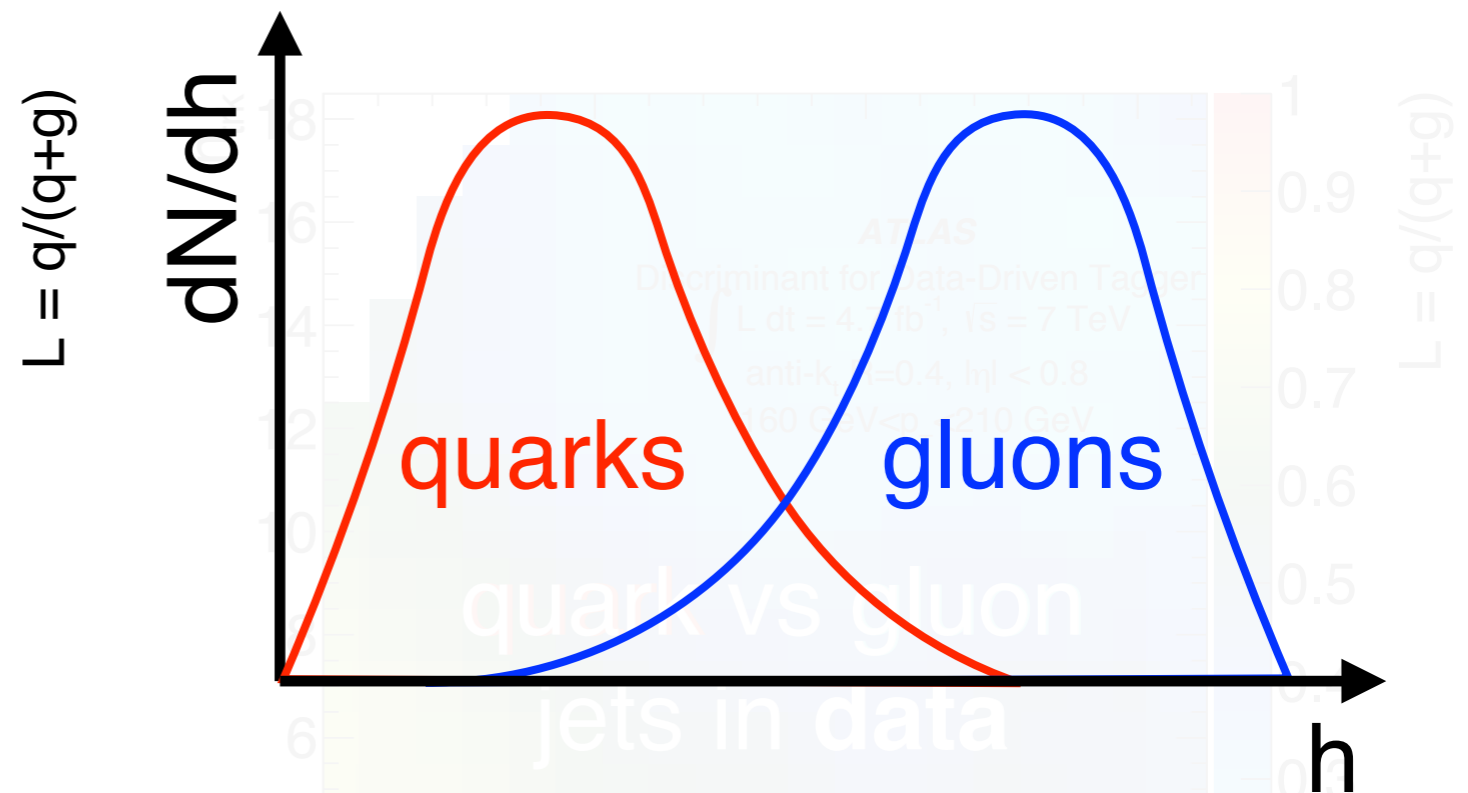
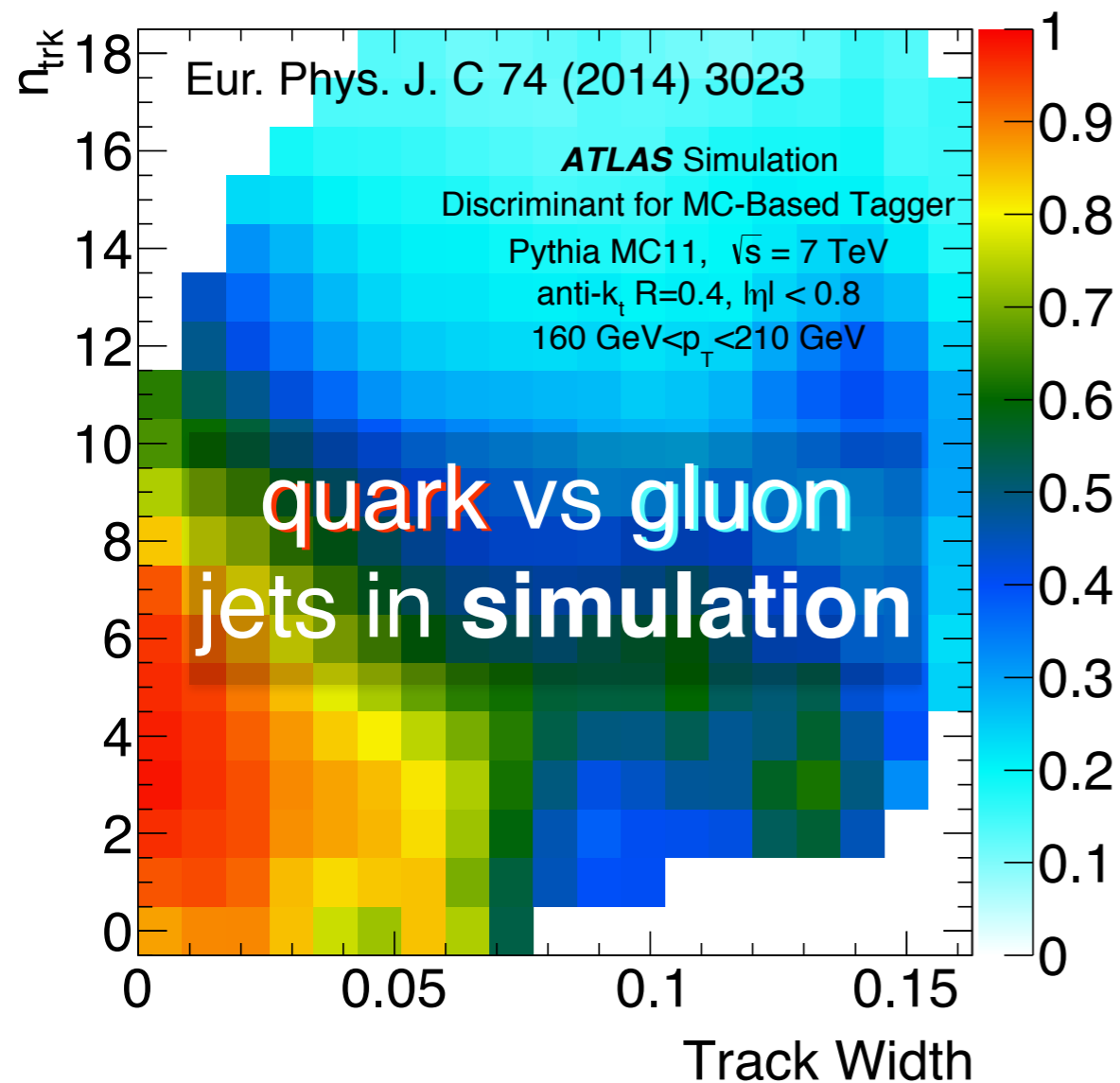
For a 2D feature space, no need for a NN or BDT - can use a histogram to “train” the classifier.

$$h(n_{\text{trk}}, \text{Track Width}) \rightarrow [0, 1]$$

Background and Motivation

44

Usual paradigm: **train in simulation**, validate on data, test on data.



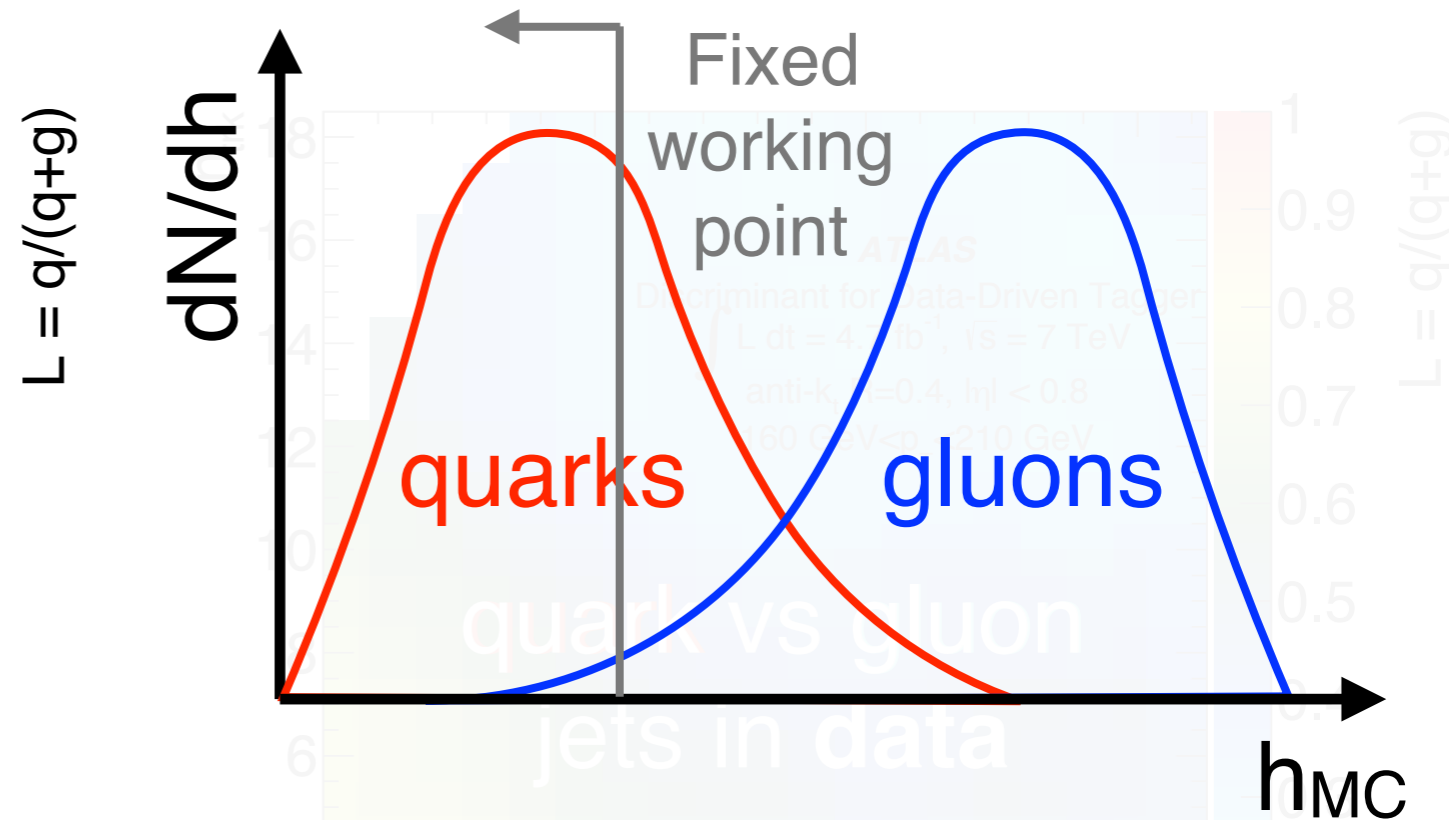
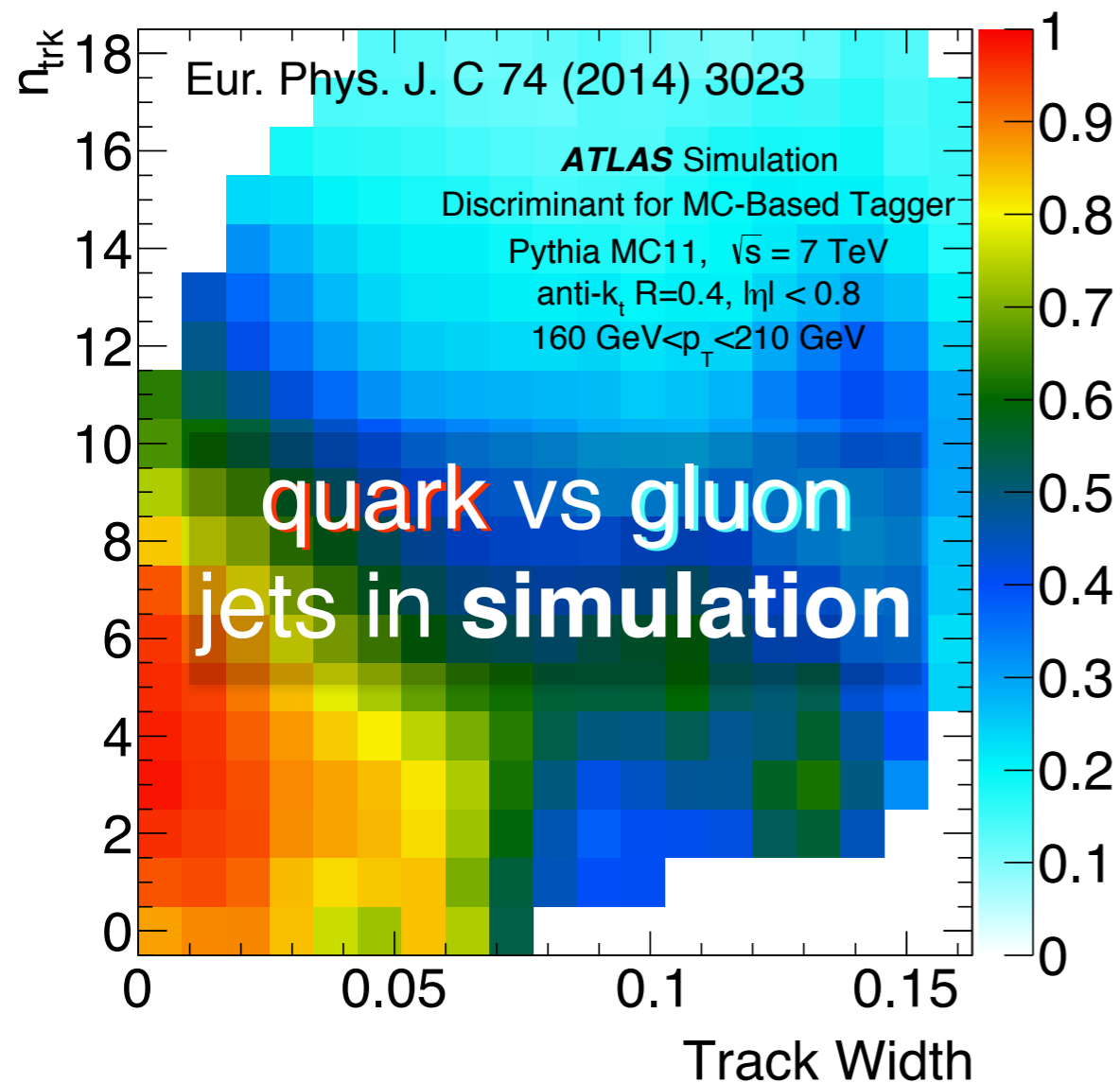
For a 2D feature space, no need for a NN or BDT - can use a histogram to “train” the **classifier**.

$$h(n_{\text{trk}}, \text{Track Width}) \rightarrow [0, 1]$$

Background and Motivation

45

Usual paradigm: **train in simulation**, validate on data, test on data.

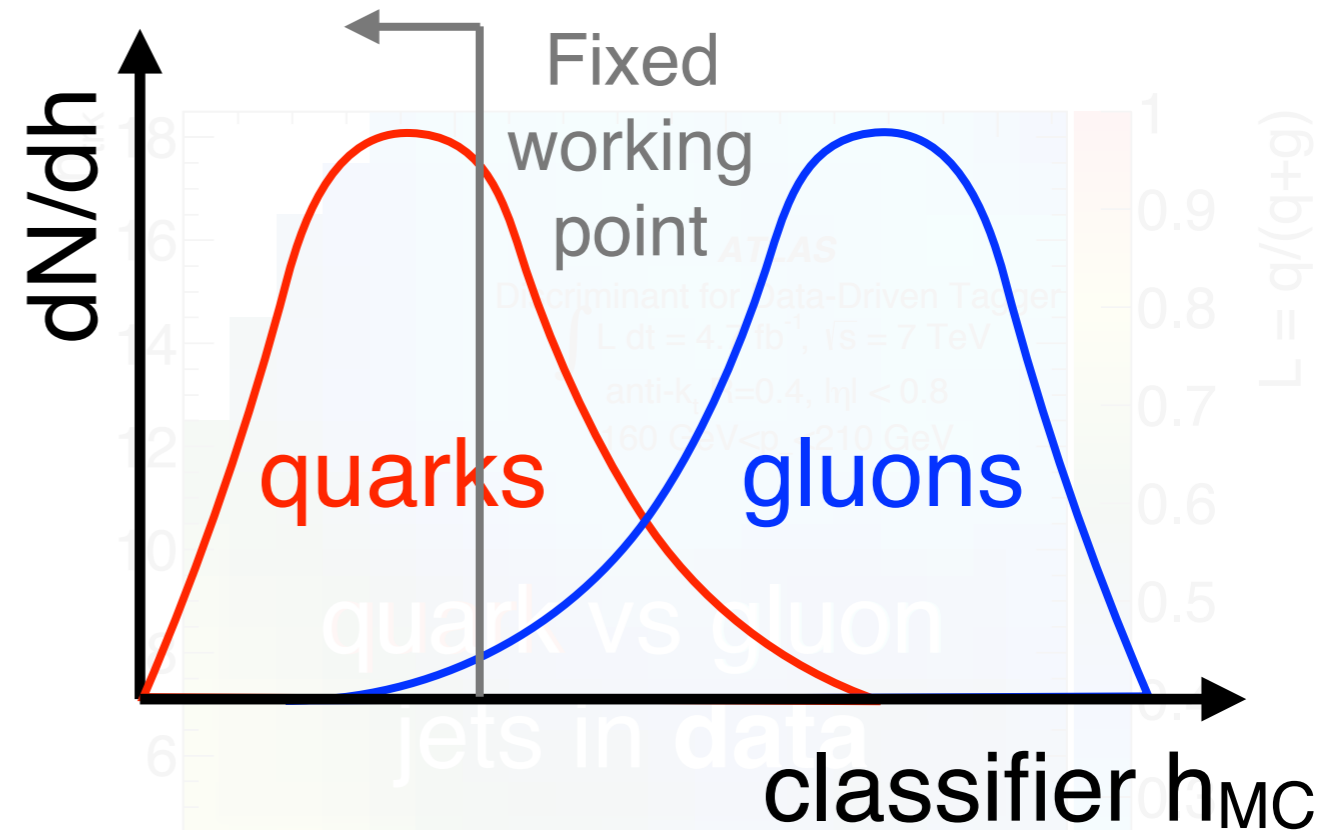
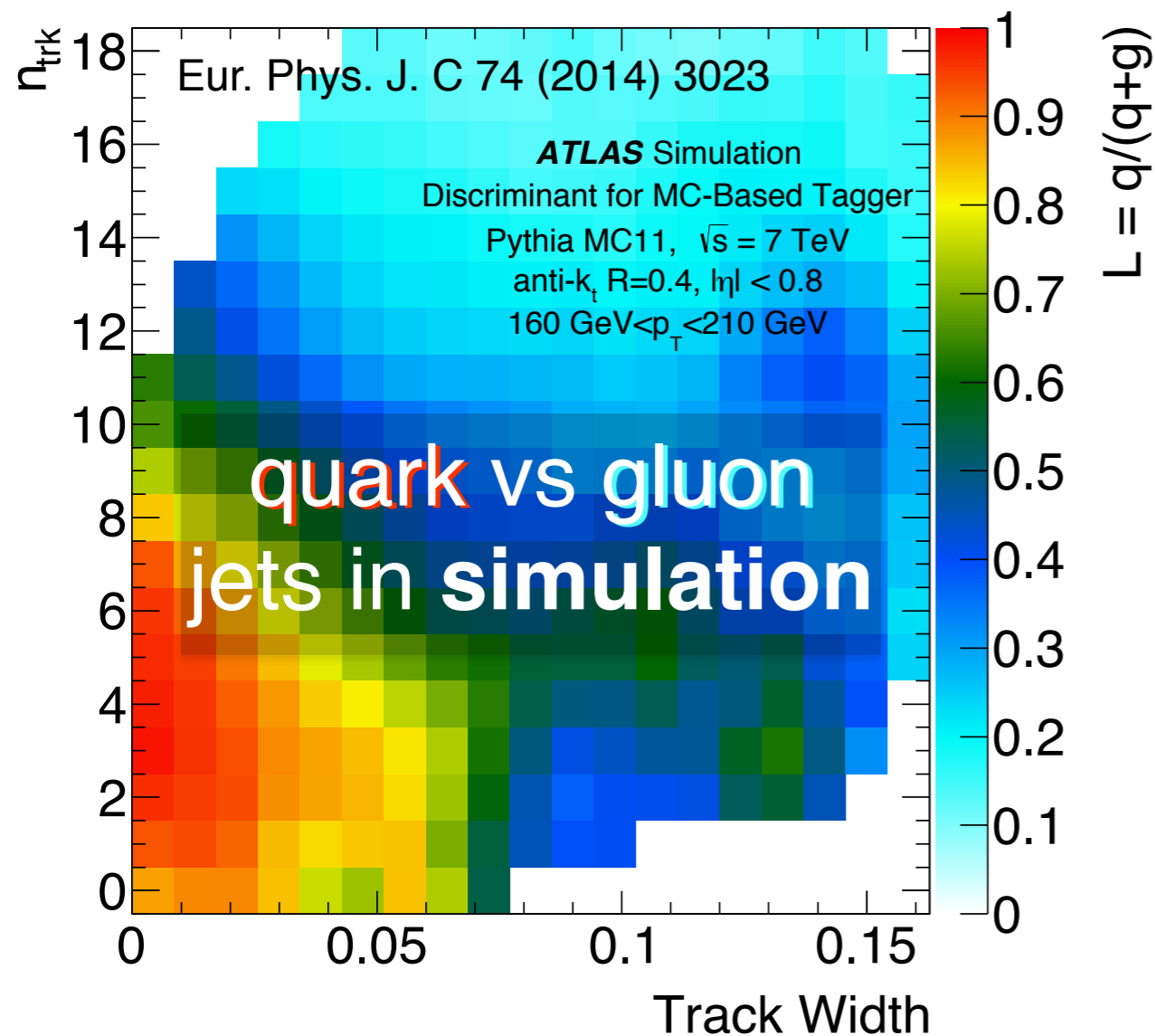


For a 2D feature space, no need for a NN or BDT - can use a histogram to "train" the classifier.

$$h_{\text{MC}}(n_{\text{trk}}, \text{Track Width}) \rightarrow [0, 1]$$

Background and Motivation

Usual paradigm: **train in simulation**, validate on data, test on data.



WP in simulation:
 $\epsilon_{\text{signal,MC}}$, $\epsilon_{\text{back,MC}}$

Background and Motivation

47

Usual paradigm: train in simulation, **validate on data**, test on data.

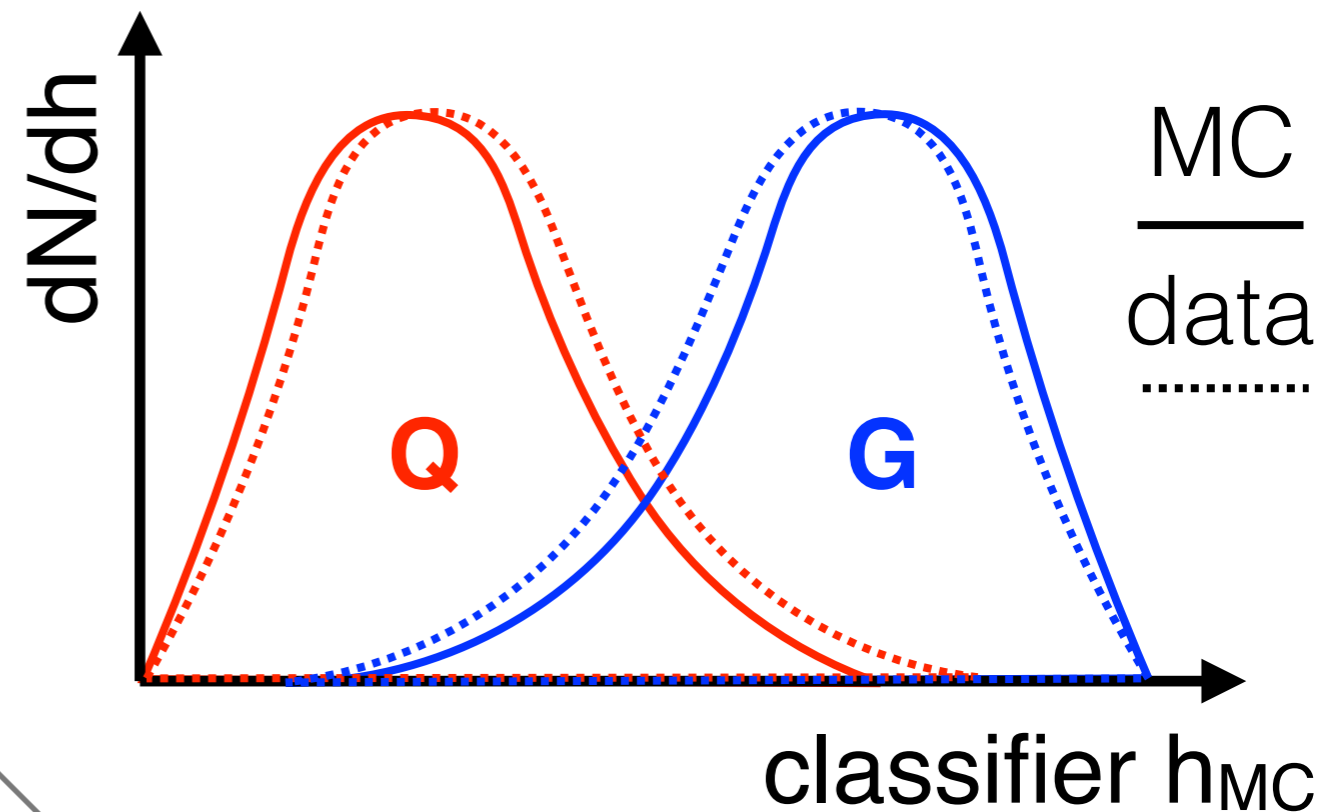
Determine the performance of the WP in data.

How did we get this?

$$\text{dijets} = f_q \times \mathbf{Q} + (1-f_q) \times \mathbf{G}$$

$$\text{Z+jets} = g_q \times \mathbf{Q} + (1-g_q) \times \mathbf{G}$$

2 equations, 2 unknowns (\mathbf{Q} , \mathbf{G})



two event samples with different q/g fractions

(N.B. f & g from simulation and selection can't bias \mathbf{Q} and \mathbf{G} - more on that later)

Background and Motivation

48

Usual paradigm: train in simulation, **validate on data**, test on data.

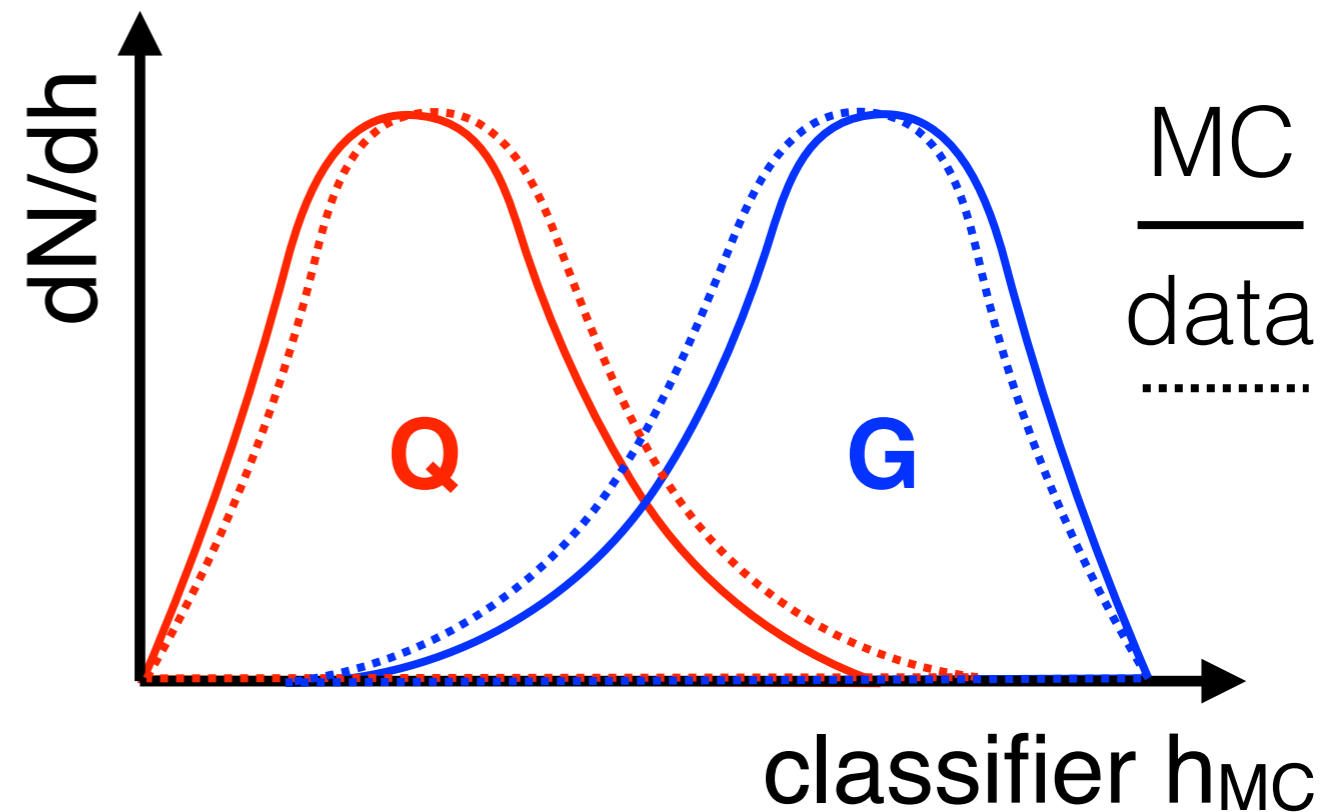
Determine the performance of the WP in data.

How did we get this?

$$\text{dijets} = f_q \times \mathbf{Q} + (1-f_q) \times \mathbf{G}$$

$$\text{Z+jets} = g_q \times \mathbf{Q} + (1-g_q) \times \mathbf{G}$$

2 equations, 2 unknowns (\mathbf{Q} , \mathbf{G})



WP in data:

$\epsilon_{\text{signal,data}}$, $\epsilon_{\text{back,data}}$

Can correct the MC to have the same performance as data.

Usual paradigm: train in simulation, validate on data, **test on data.**

Once we have scale factors (& their uncertainty), we can ensure that our analysis will be accurate.

...so what is the problem?

remember my claim from earlier:

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

This is an accuracy versus precision problem. It is “easy” to achieve accuracy through calibration, but the results may not be the best one possible.

Background and Motivation



In this 2D feature space, we can actually derive h_{data} .

Using the same trick as earlier:

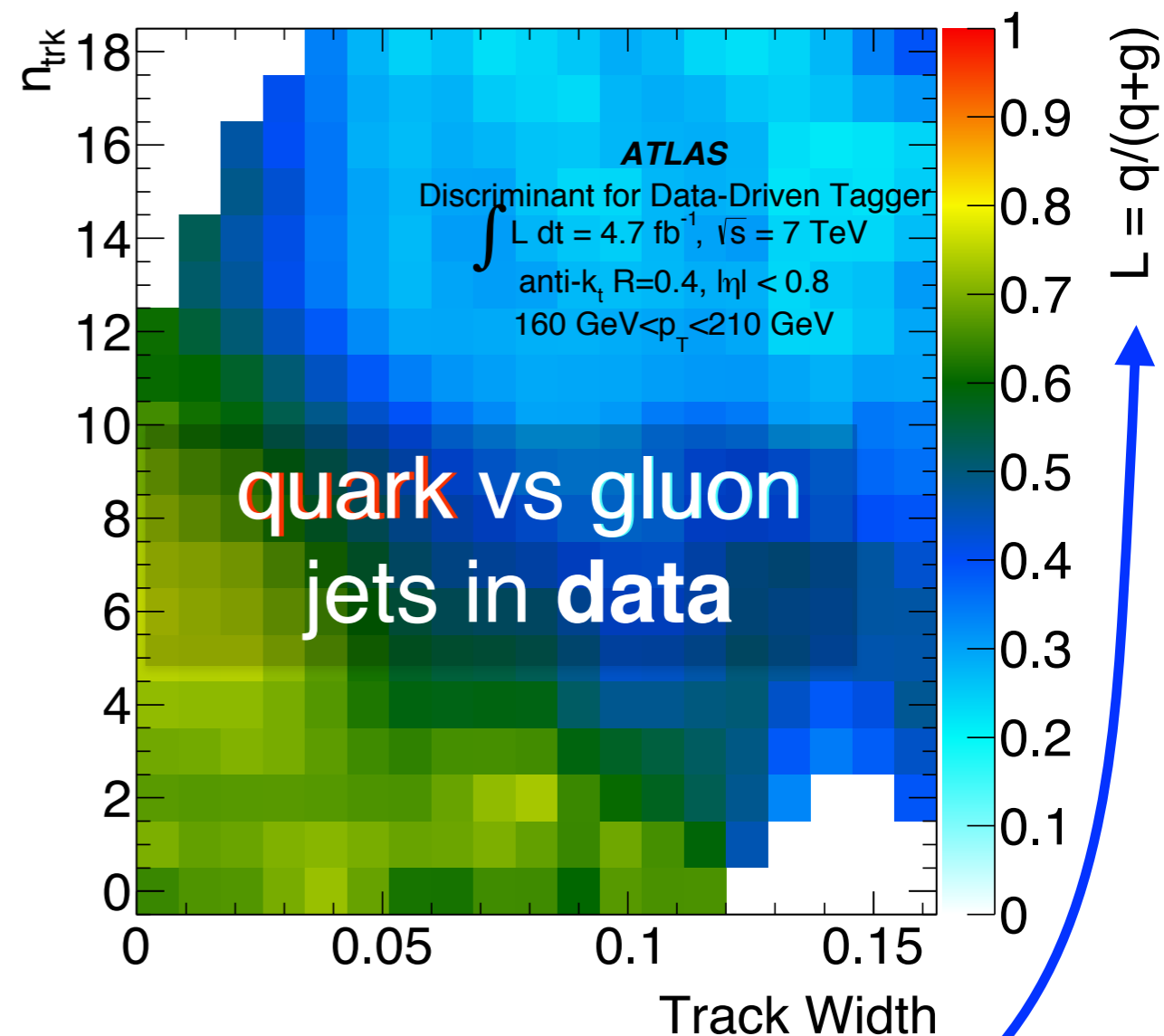
$$\text{dijets} = f_q \times Q + (1-f_q) \times G$$

$$L = q/(q+g)$$

$$\text{Z+jets} = g_q \times Q + (1-g_q) \times G$$

2 equations, 2 unknowns (Q, G)

(now Q and G are 2D histograms)



in general:

$$h_{\text{MC}}(n_{\text{trk}}, \text{Track Width}) \neq h_{\text{data}}(n_{\text{trk}}, \text{Track Width})$$

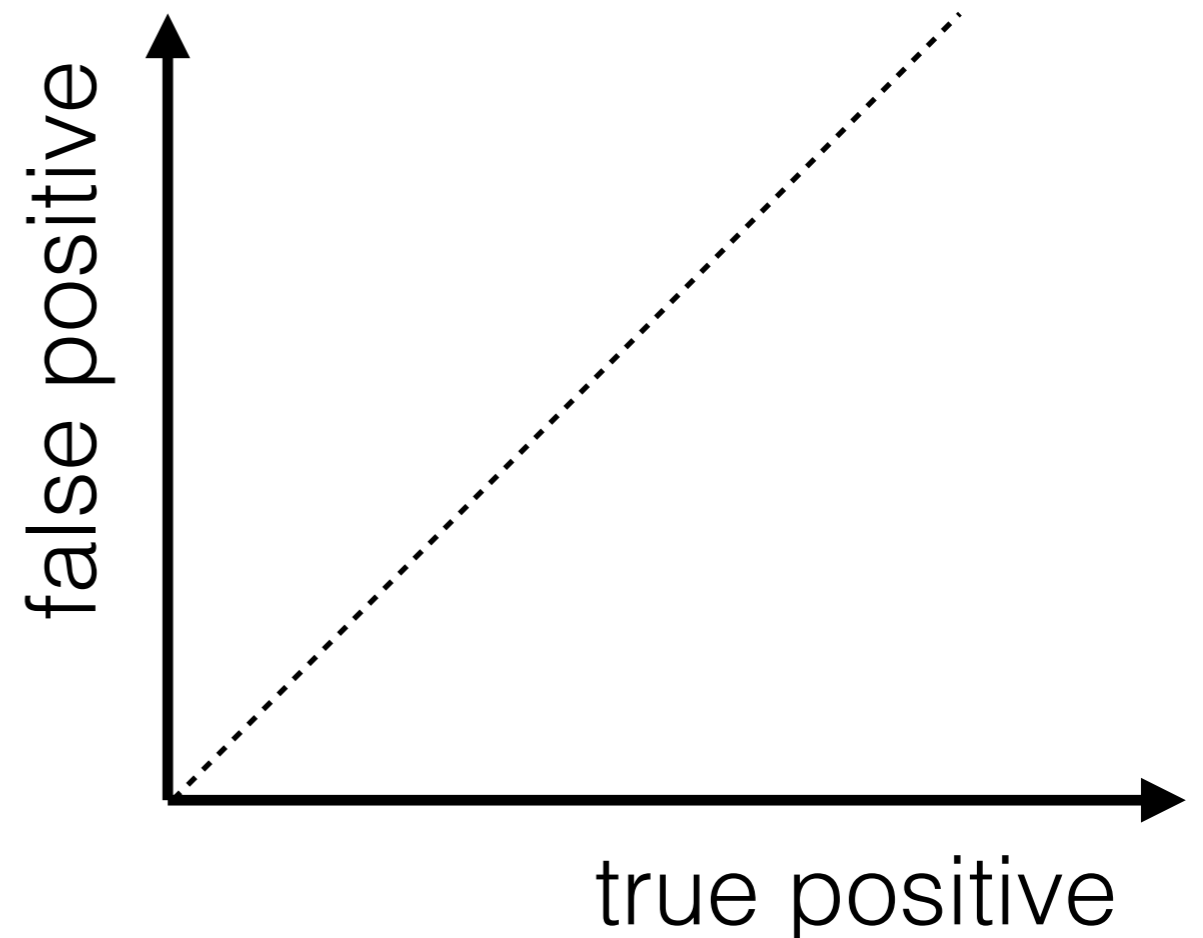
Take it to the extreme

51

To stress this point, suppose that h_{MC} is the random classifier:

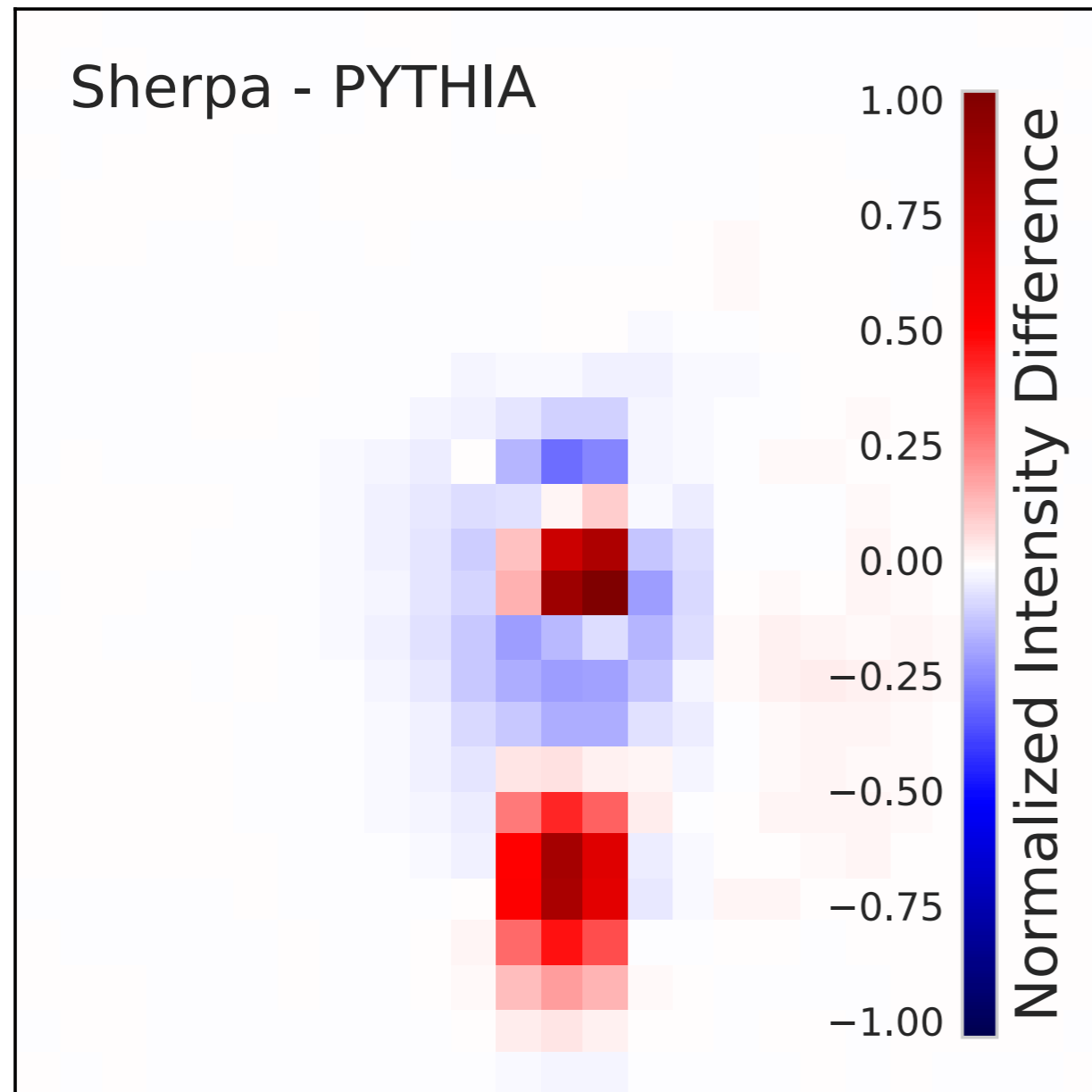
$h_{MC} = 0$ if you pick a random number x in $[0, 1]$ and $x < \epsilon$
1 otherwise

We can calibrate this classifier in data, but clearly, it is sub-optimal !!



One more slide about why it matters

52



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

One solution: Learn from data!

53

One of the biggest challenges with any MC-based method is that it can't use information that the MC doesn't know about.

One solution is to train directly on data !

In general, this is not possible since data are unlabeled. However, in a wide range of cases, it is possible to work with less.

There is an interesting connection between what I'm calling "weak supervision" and the topic of "label noise".

The setup: suppose you have (at least) two mixed samples, each composed of two classes (say q and g).

Requirement:

The two classes are well-defined i.e. q in sample 1 is statistically identical to q in sample 2).

Weak sup. option 1: Use class proportions

55

Remember this plot?

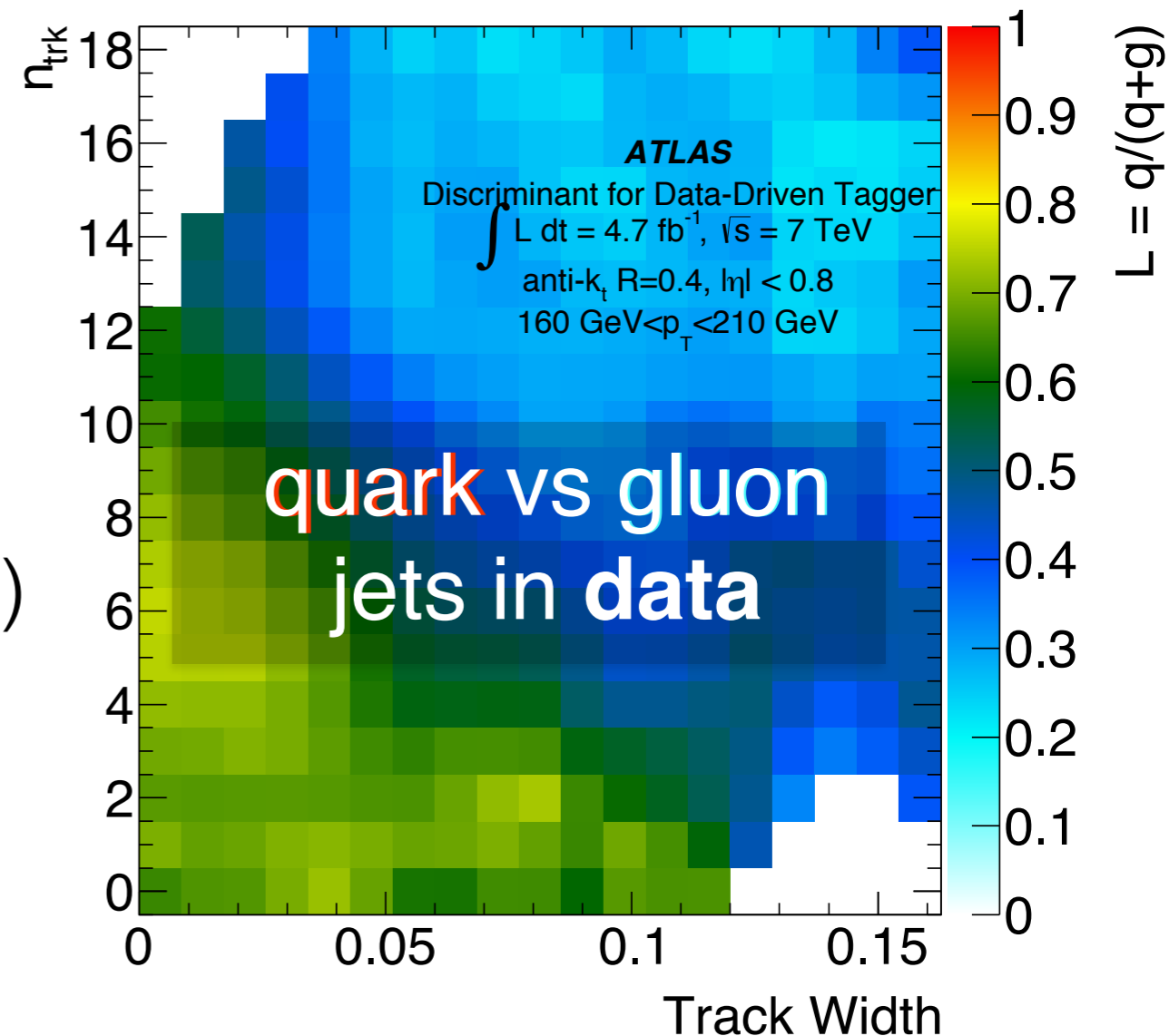
$$\text{dijets} = f_q \times Q + (1-f_q) \times G$$

$$\text{Z+jets} = g_q \times Q + (1-g_q) \times G$$

two equations, two unknowns (Q, G)

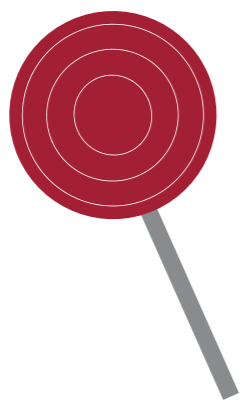
We often know f, g

(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.

Method 1: Learn from Proportions



LoLiProp

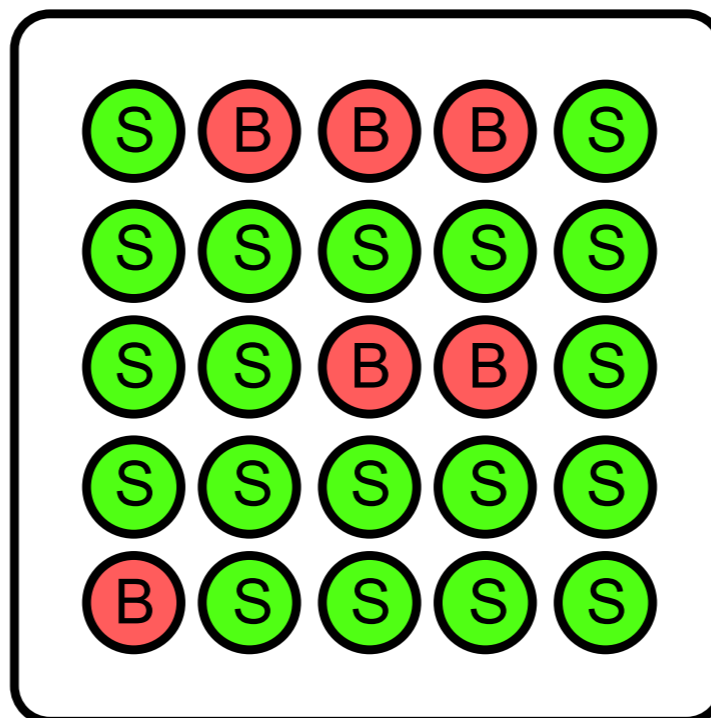
Learning from Label Proportions

Solution: Train using class proportions.
Work “on average”

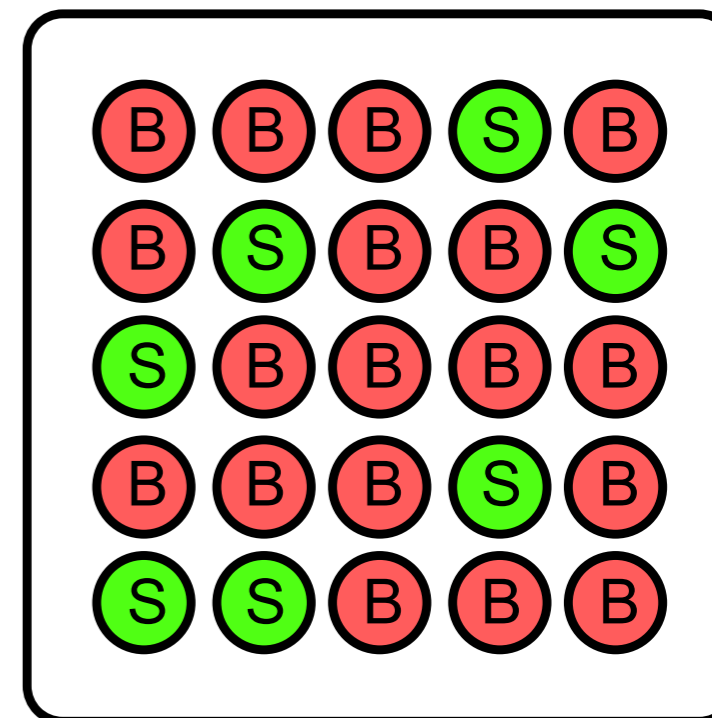
$$f_{\text{full}} = \operatorname{argmin}_{f': \mathbb{R}^n \rightarrow \{0,1\}} \sum_{i=1}^N \ell(f'(x_i) - t_i)$$

ℓ ← loss fcn. t_i ← labels

Mixed Sample 1



Mixed Sample 2



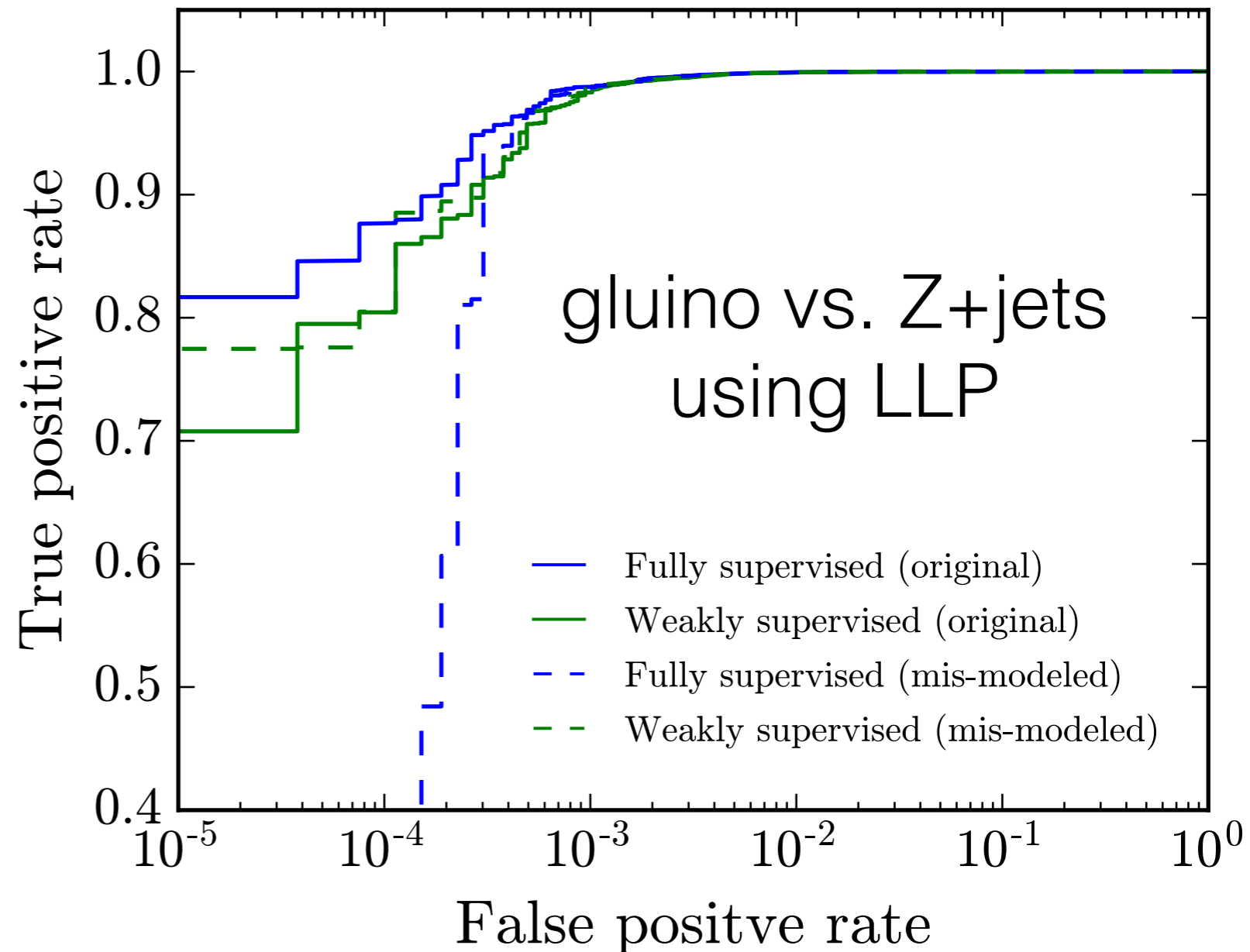
$$f_{\text{weak}} = \operatorname{argmin}_{f': \mathbb{R}^n \rightarrow [0,1]} \ell \left(\sum_{i=1}^N \frac{f'(x_i)}{N} - y \right)$$

$\frac{f'(x_i)}{N}$ ← proportions y ← proportions

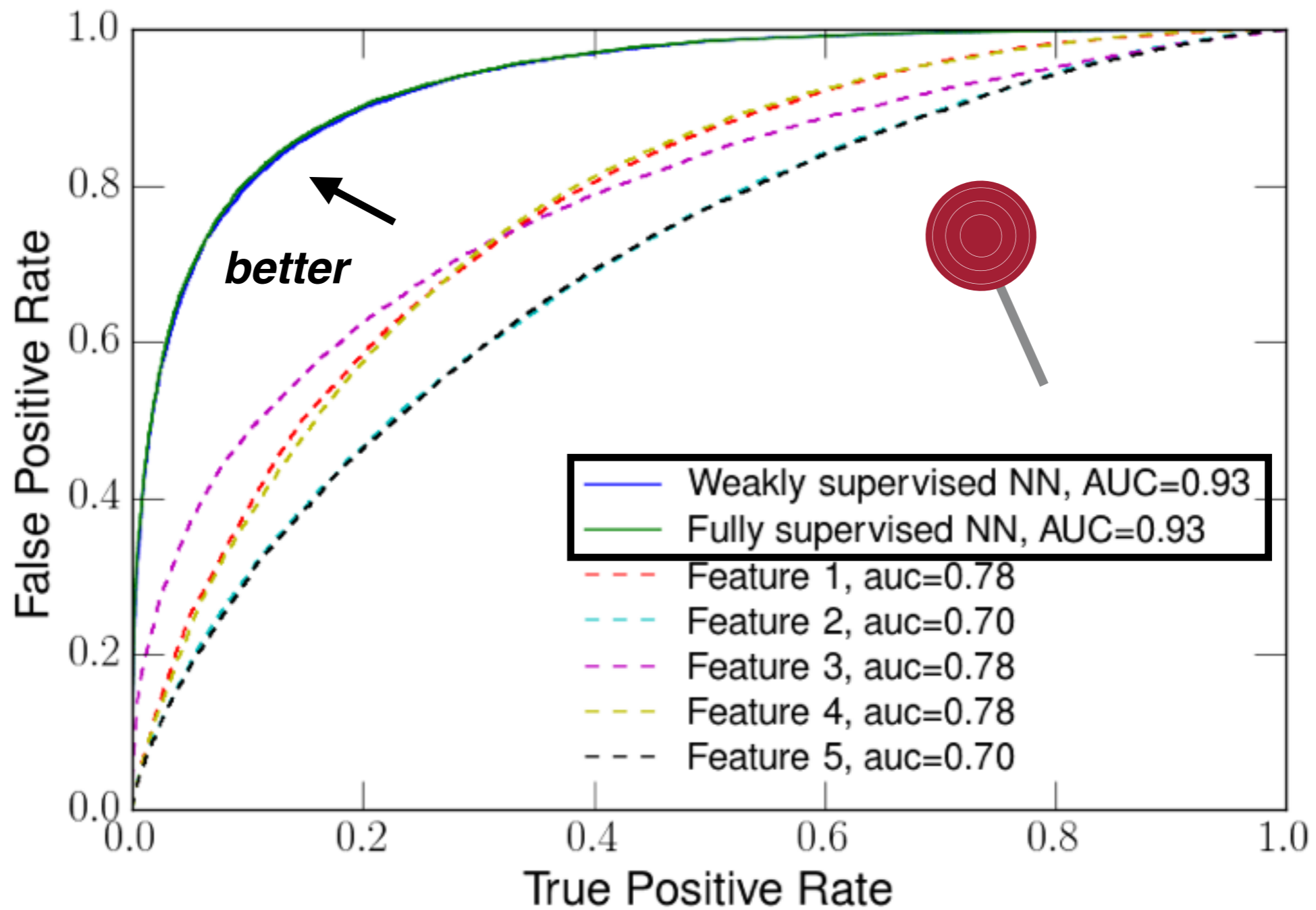
N.B. Don't need 100% fraction accuracy

57

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

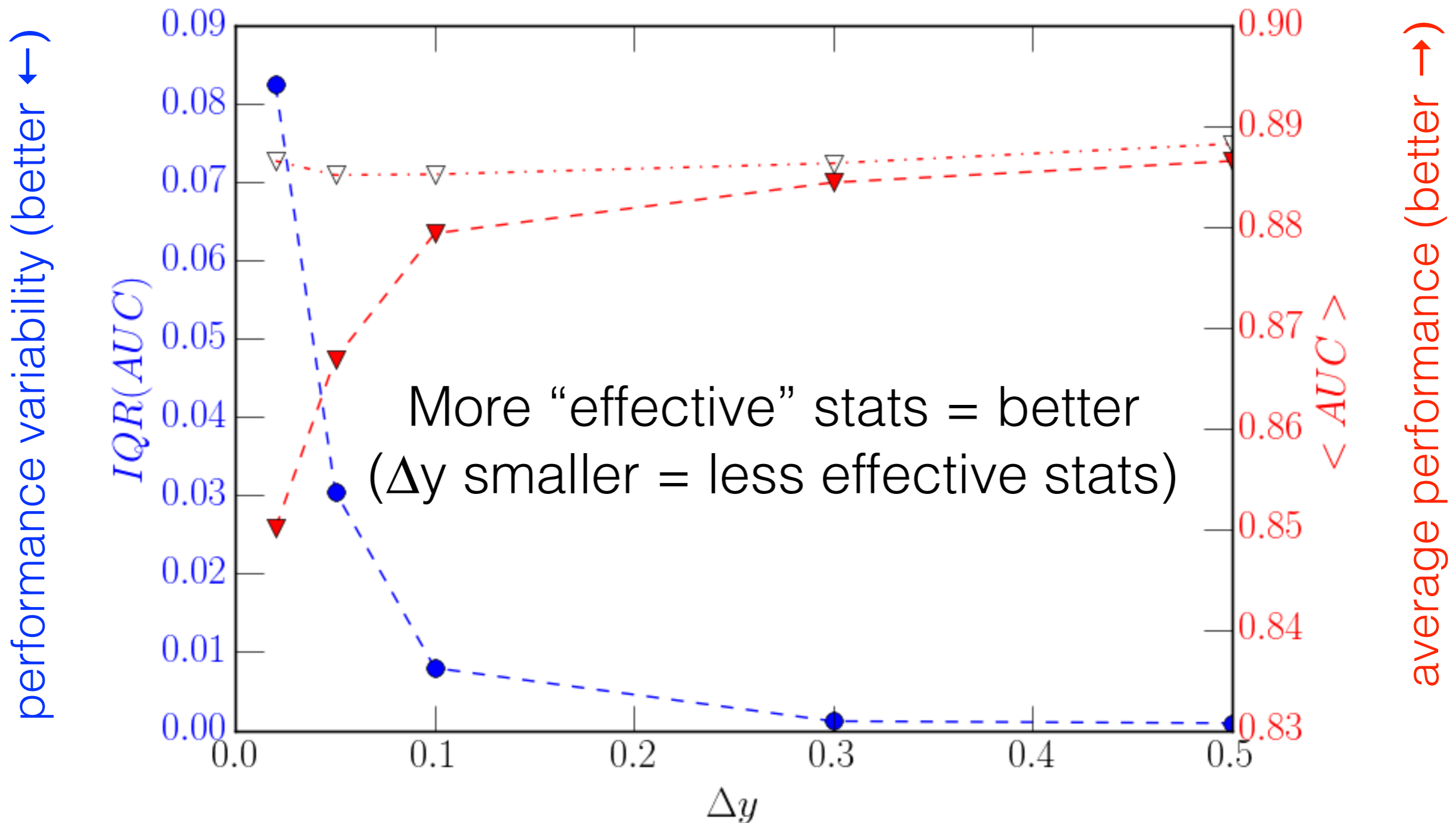


Works in low-dimensions



A note about training statistics

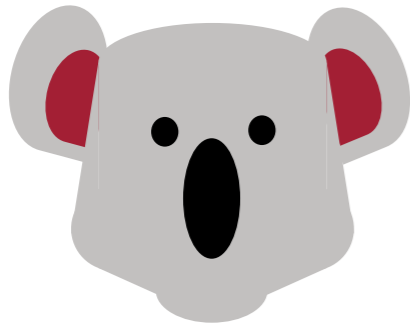
59



how different are the proportions for the two mixed samples

Method 2: Learning without Proportions

60

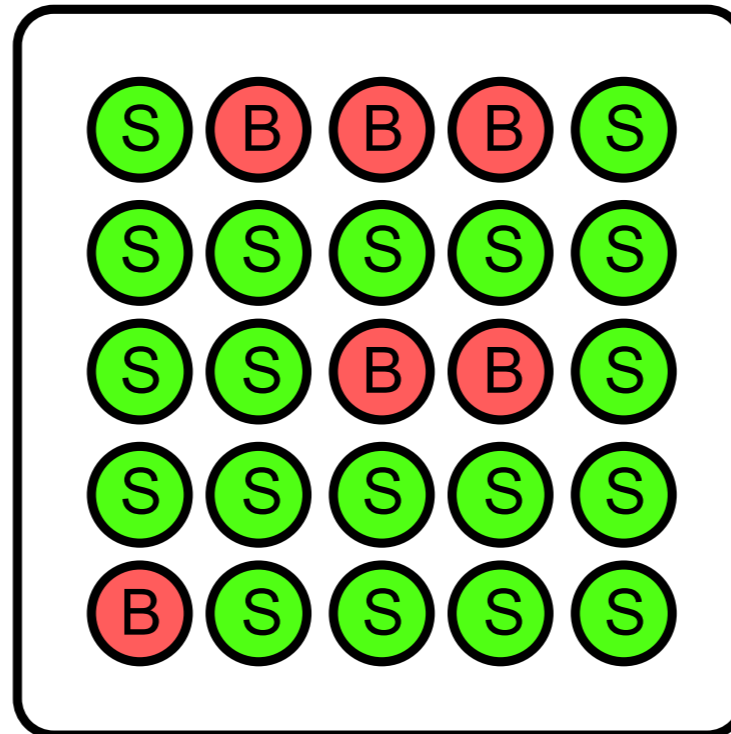


CWoLa

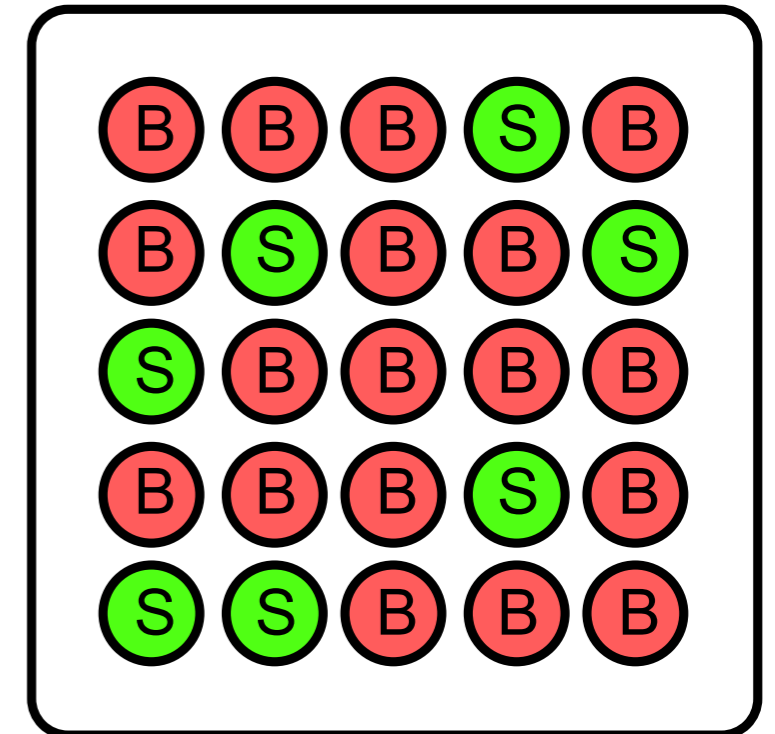
*Classification
Without Labels*

Solution: Train
directly on data using
mixed samples

Mixed Sample 1



Mixed Sample 2



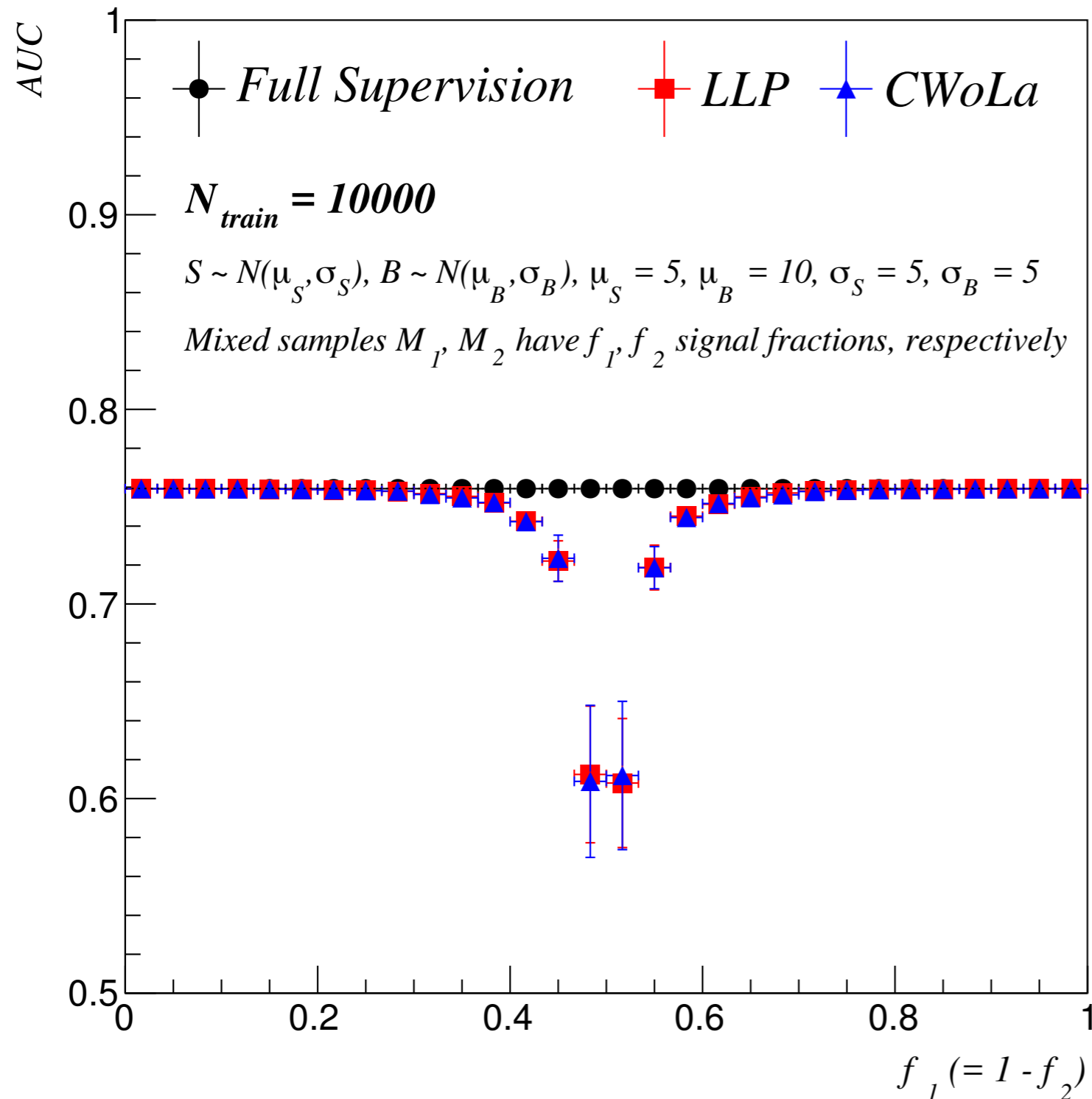
0

1

Classifier



A note about training statistics

62



As with LLP, need sufficient effective statistics

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

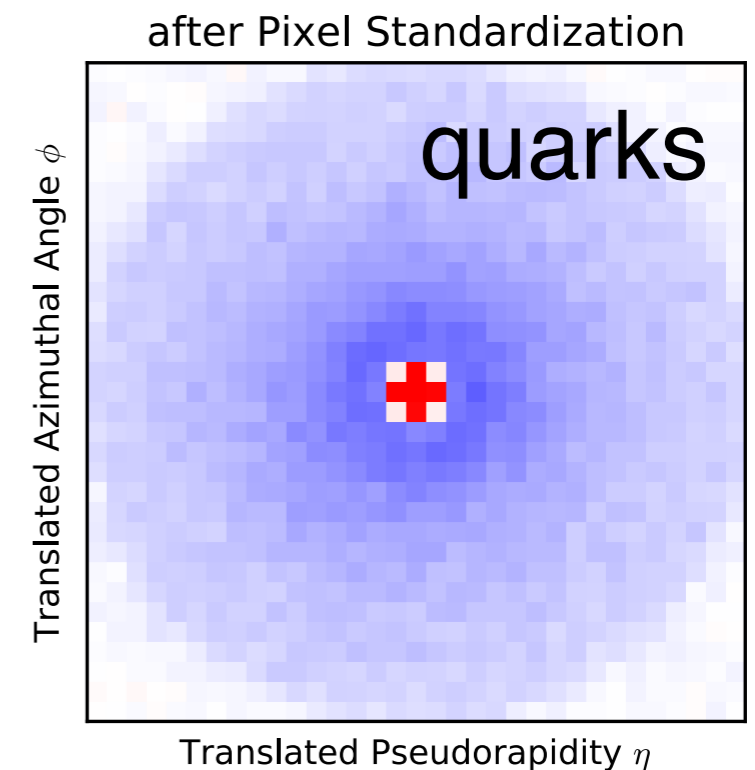
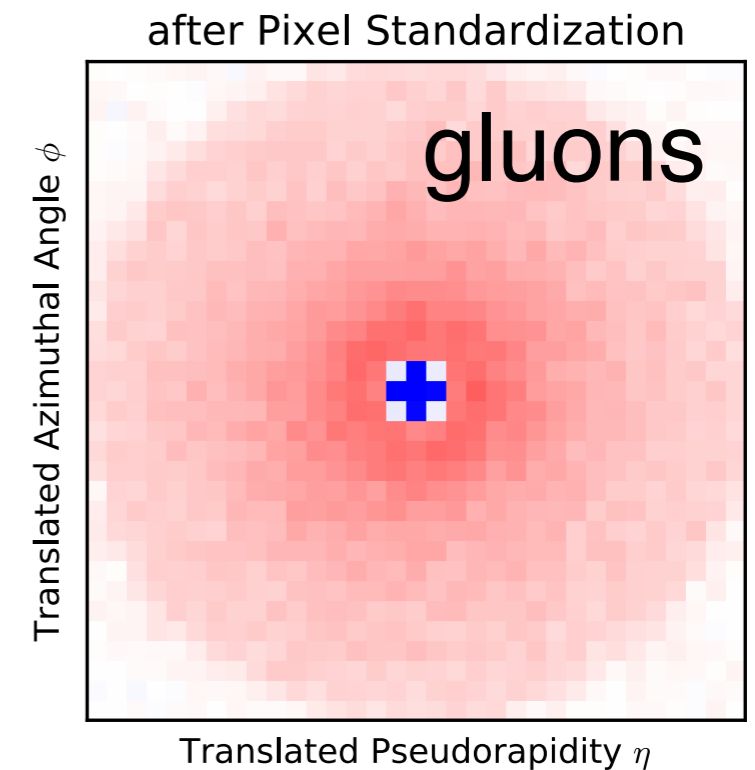
Property	 LLP	 CWoLa
Compatible with any trainable model	✓	✓
No training modifications needed	✗	✓
Training does not need fractions	✗	✓
Smooth limit to full supervision	✗	✓
Works for > 2 mixed samples	✓	?

Next step: what about high dim.?

64

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.

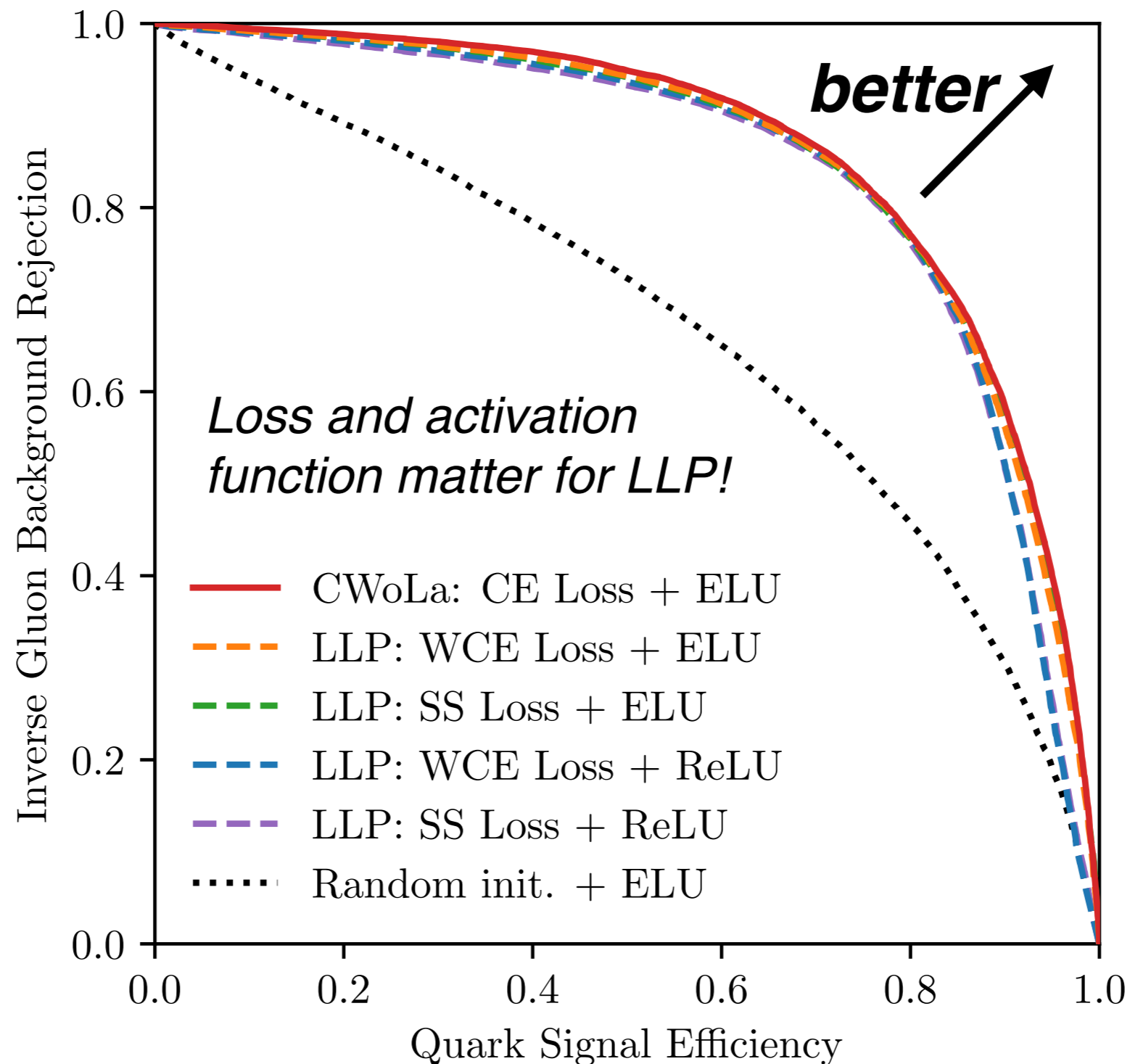


The CWoLa approach works out-of-the box - can use well-tested CNN architecture with usual cross-entropy loss.

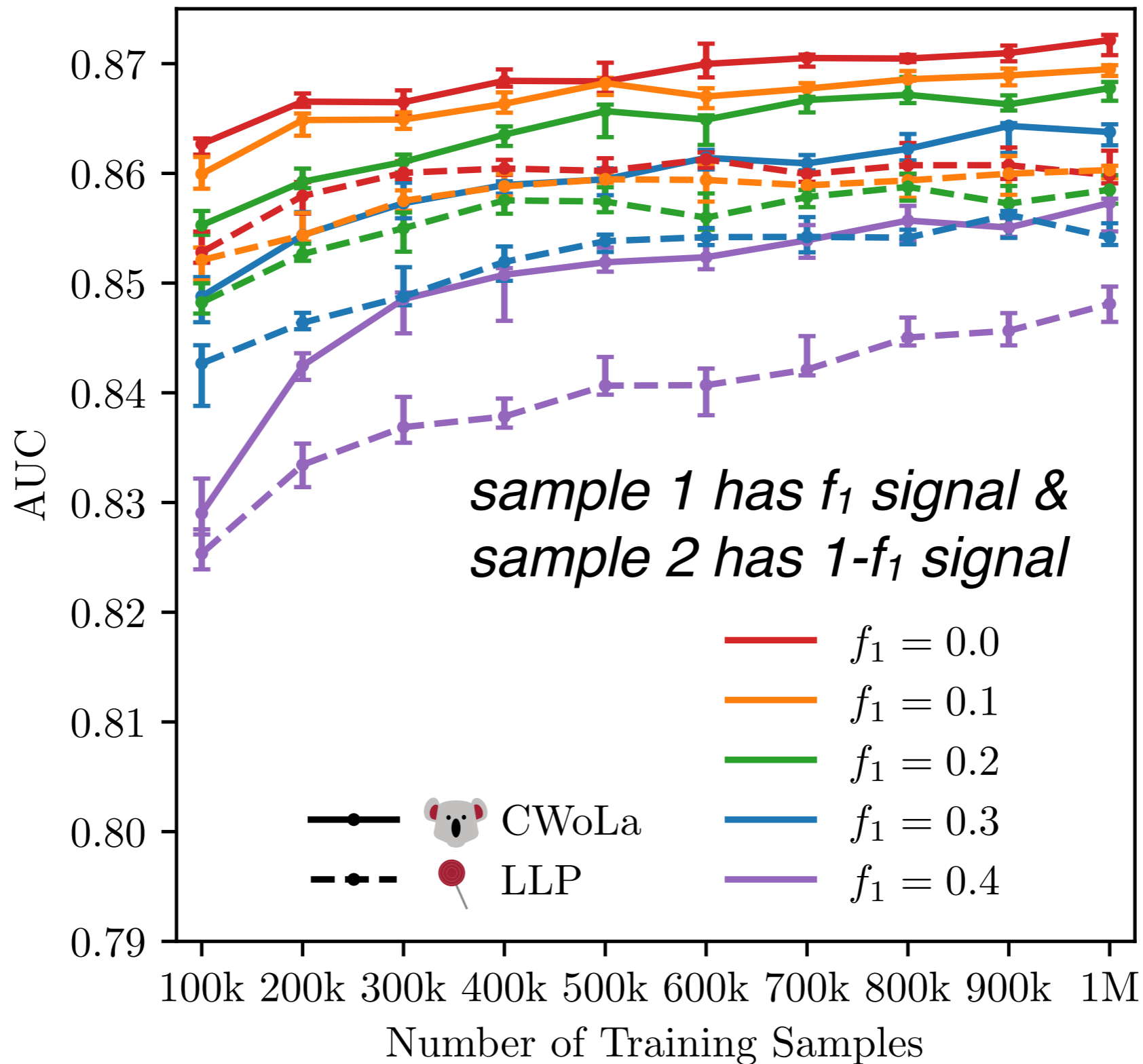
On the other hand, LLP requires significant work on the technical implementation / optimization.

$$\ell_{\text{WMSE}} = \sum_a \left(f_a - \frac{1}{N} \sum_{i=1}^N h(\mathbf{x}_i) \right)^2 \quad \ell_{\text{WCE}} = \sum_a \text{CE} \left(f_a, \frac{1}{N} \sum_{i=1}^N h(\mathbf{x}_i) \right)$$

Works in many-dimensions!



A note about training statistics



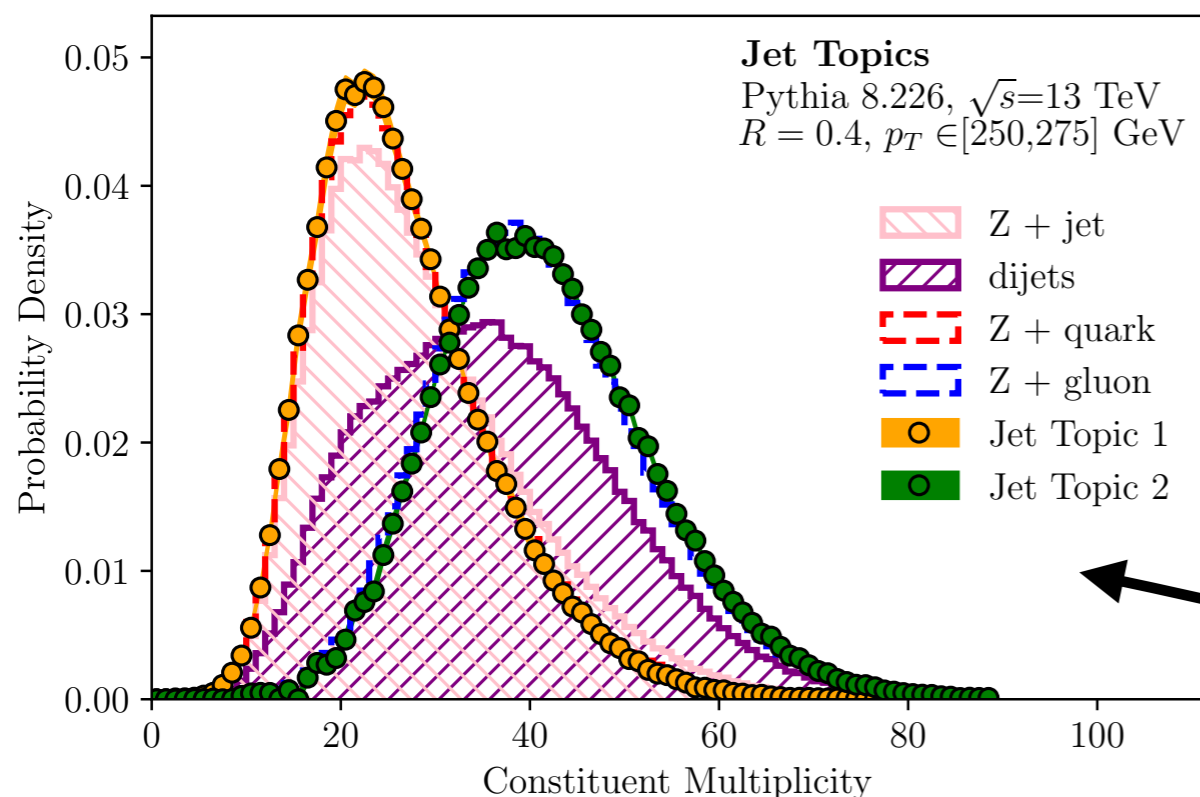
Hybrid Approaches

68

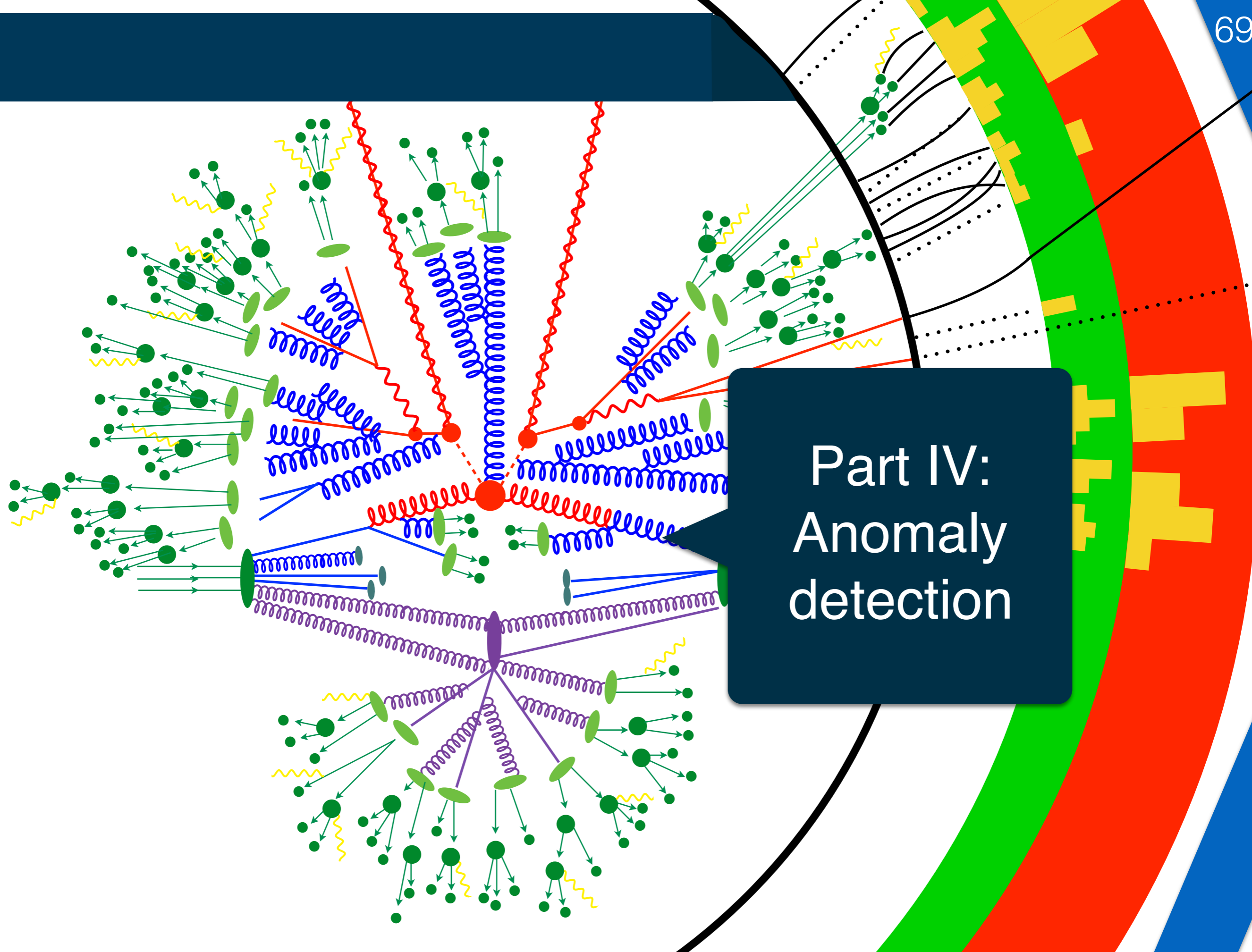
As usual, it is likely that the best approach will use all of the available information, including some input from simulation.

...this could be as simple as pre-training in MC and then running weak supervision

E. Metodiev and J. Thaler, Phys. Rev. Lett. 120 (2018) 241602



Another interesting direction is to push the weak supervision paradigm a step further and **define** the classes so that it works.



Part IV:
Anomaly
detection

I will leave you with one last, but very exciting topic.

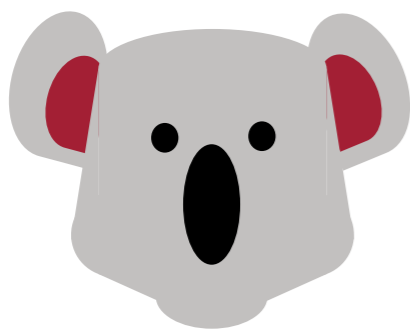
One of the most important goals of HEP is to search for new particles. However, we have not found anything (significantly) unexpected in a while ... we need simulation-independent ways of searching for new particles !

anomalies, i.e. something unexpected



N.B. The approach discussed here is not the only one - see also M. Farina, Y. Nakai, D. Shih, 1808.08992 & T. Heimel, G. Kasieczka, T. Phlen, J. Thompson, 1808.08979 for an alternative approach based on auto-encoders.

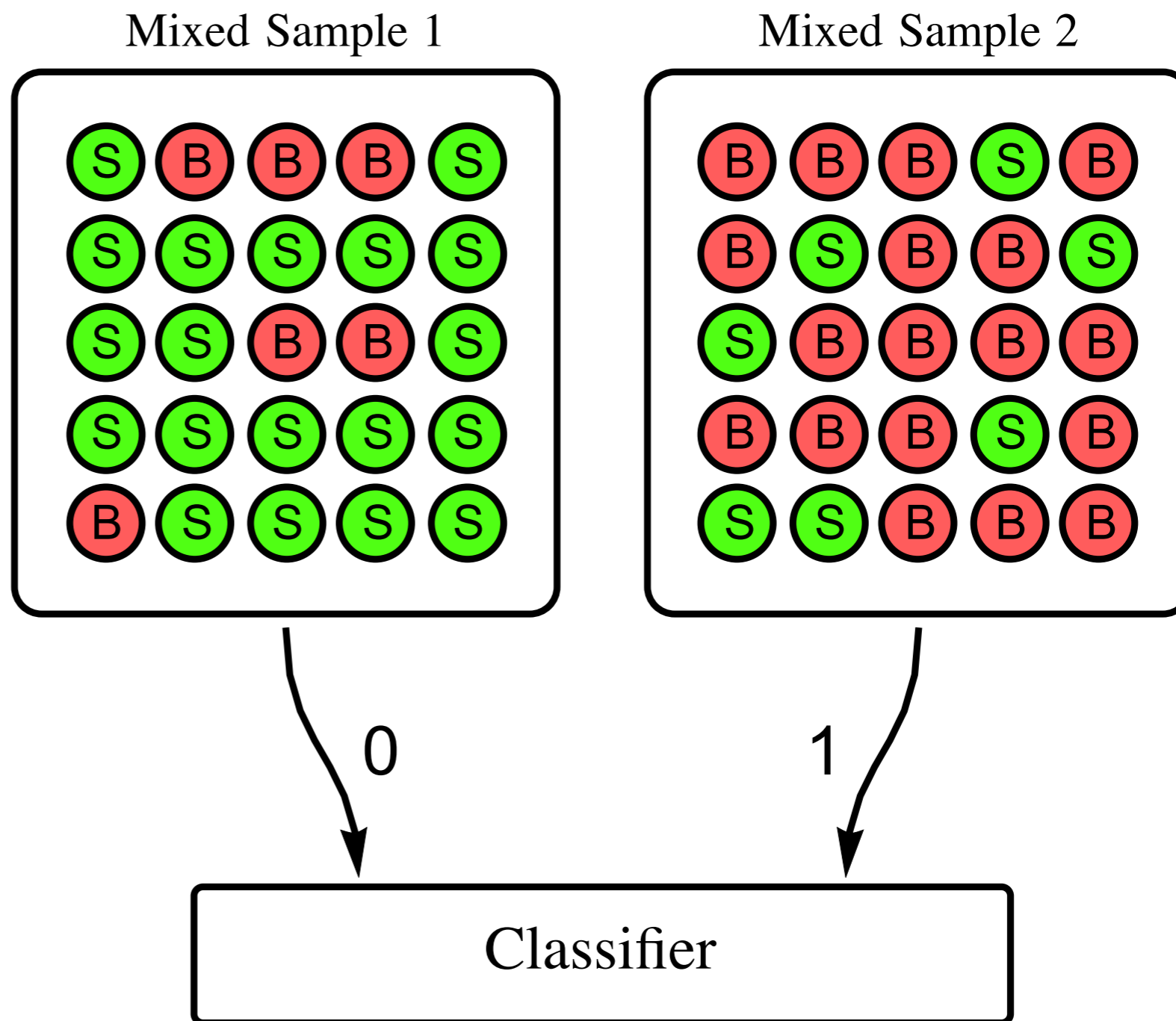
Remember CWoLa ...



CWoLa

*Classification
Without Labels*

Solution: Train
directly on data using
mixed samples



Weak/unsupervised learning for anomalies

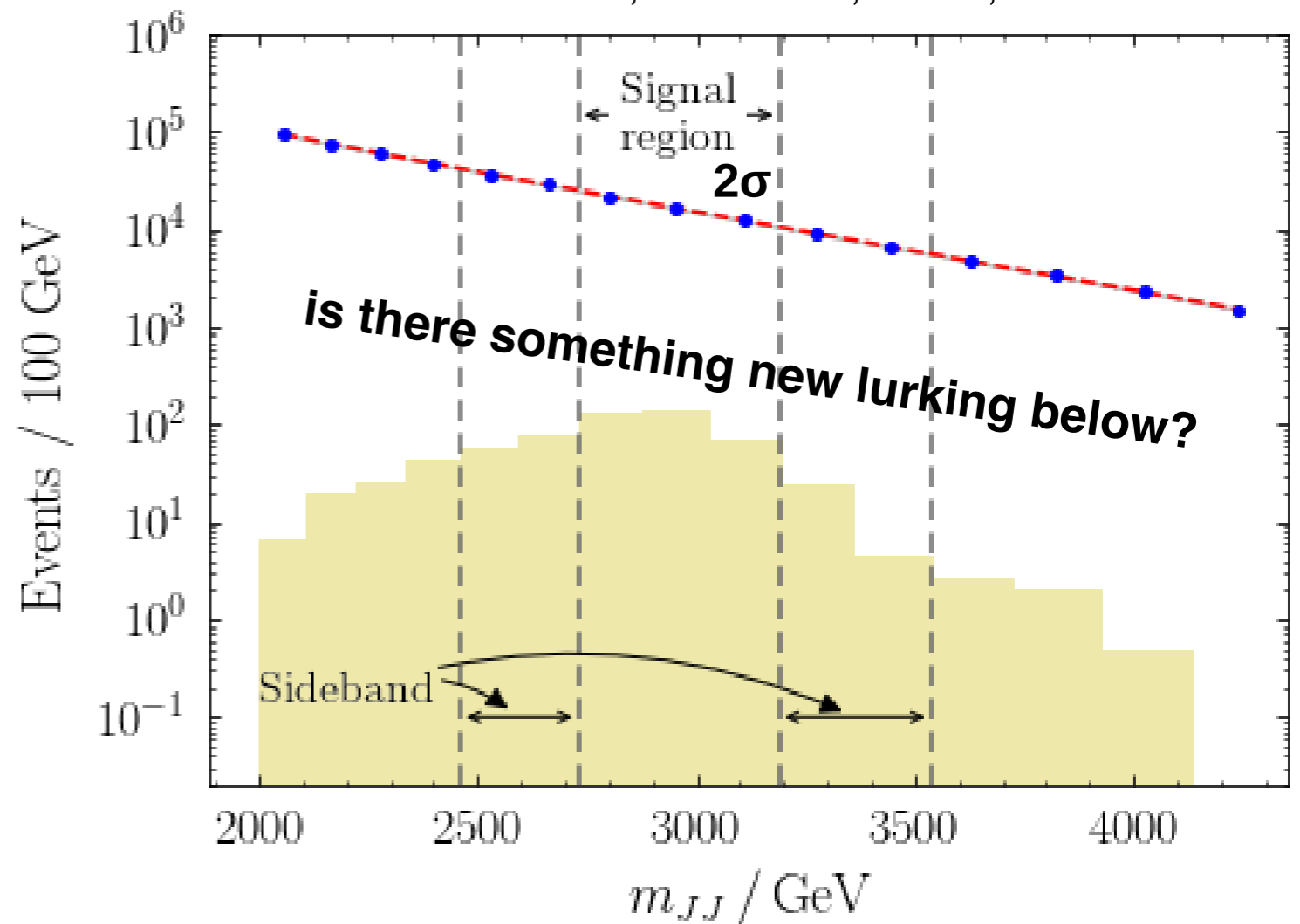


Can we take this idea one step further to look for something unexpected?

= CWoLa Hunting*



J. Collins, K. Howe, **BPN**, 1805.02664



Minimal assumption: if there is a signal, it is localized in one known dimension.

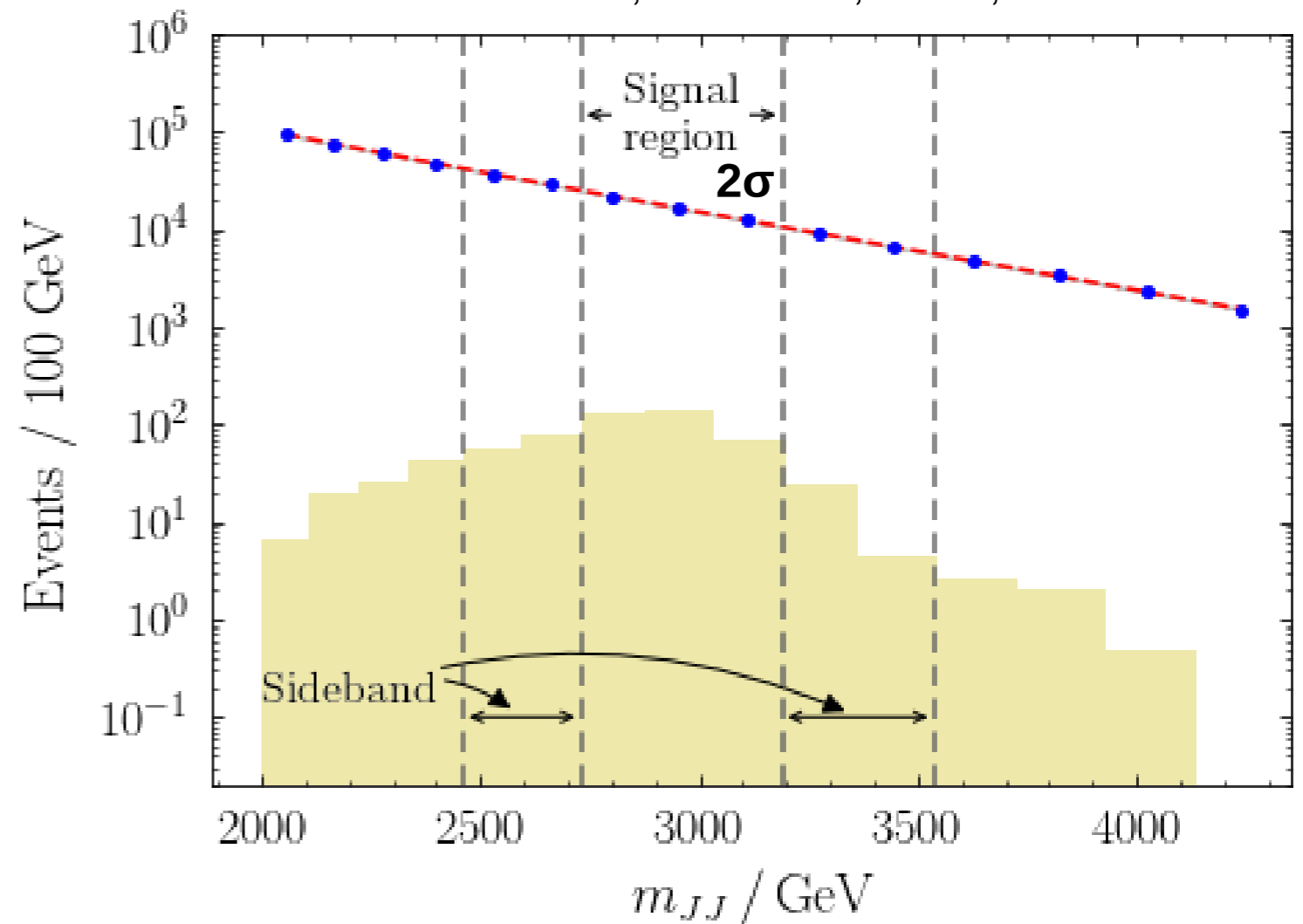
*Image from [this article](#). This Koala is actually being freed - I do not condone violence against these animals!

Mixed sample 1:
signal region

Mixed sample 2:
sideband region

Train a classifier to
distinguish the two
mixed samples.

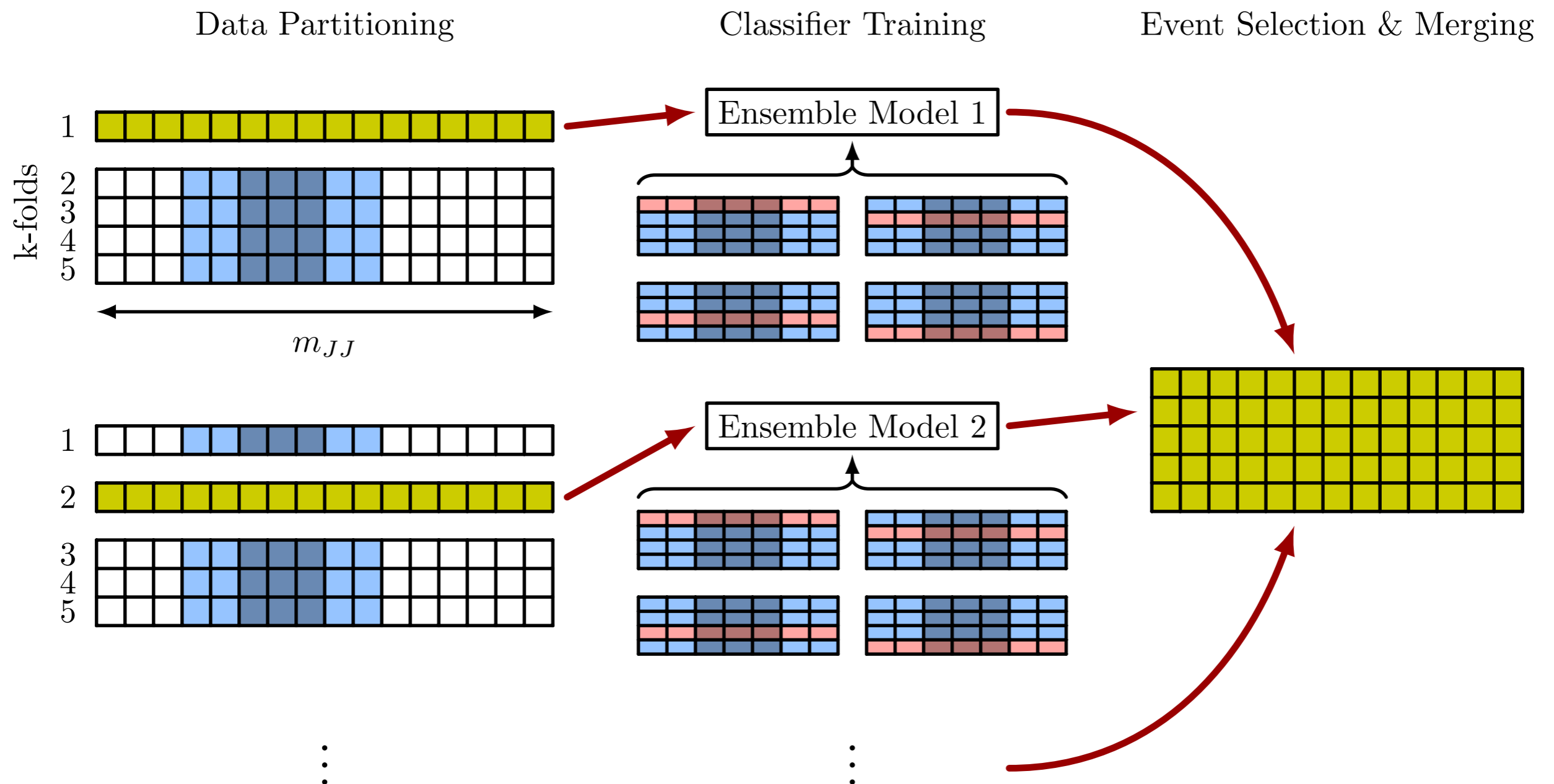
J. Collins, K. Howe, **BPN**, 1805.02664



If there is a signal, there will be something to learn and the signal will be enhanced. If no signal, nothing to learn.

Weak/unsupervised learning for anomalies

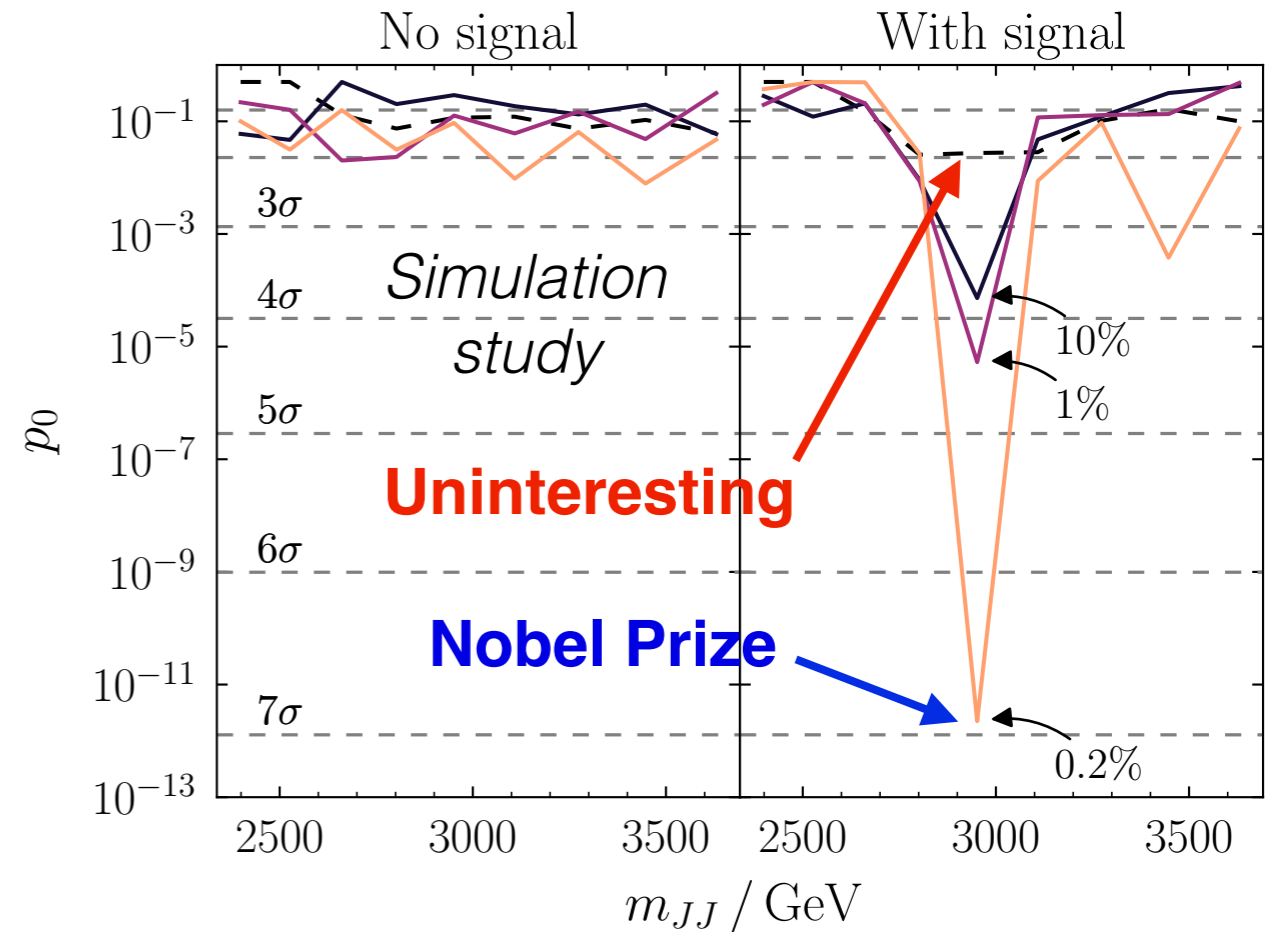
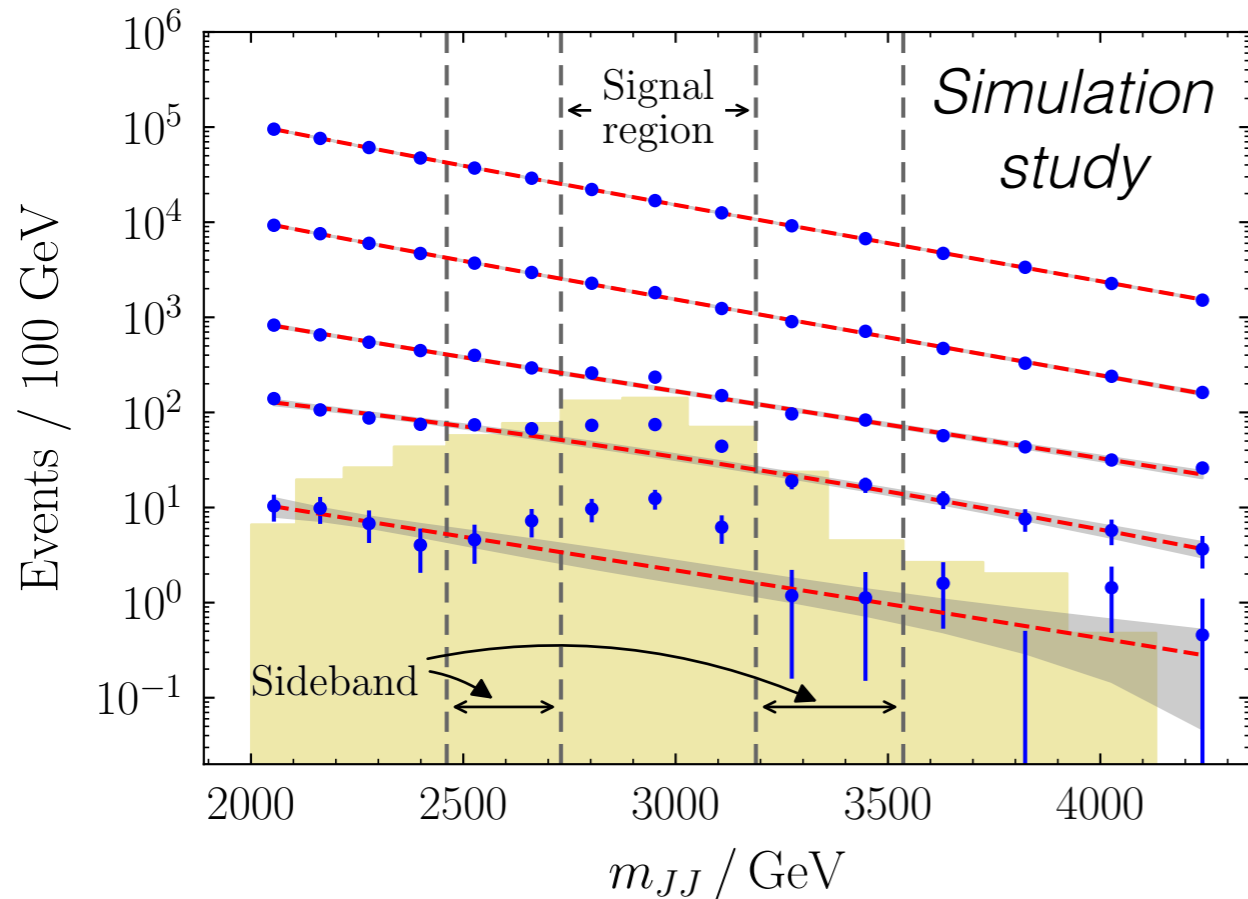
Need to be careful about testing/training on the same data.



Weak/unsupervised learning for anomalies

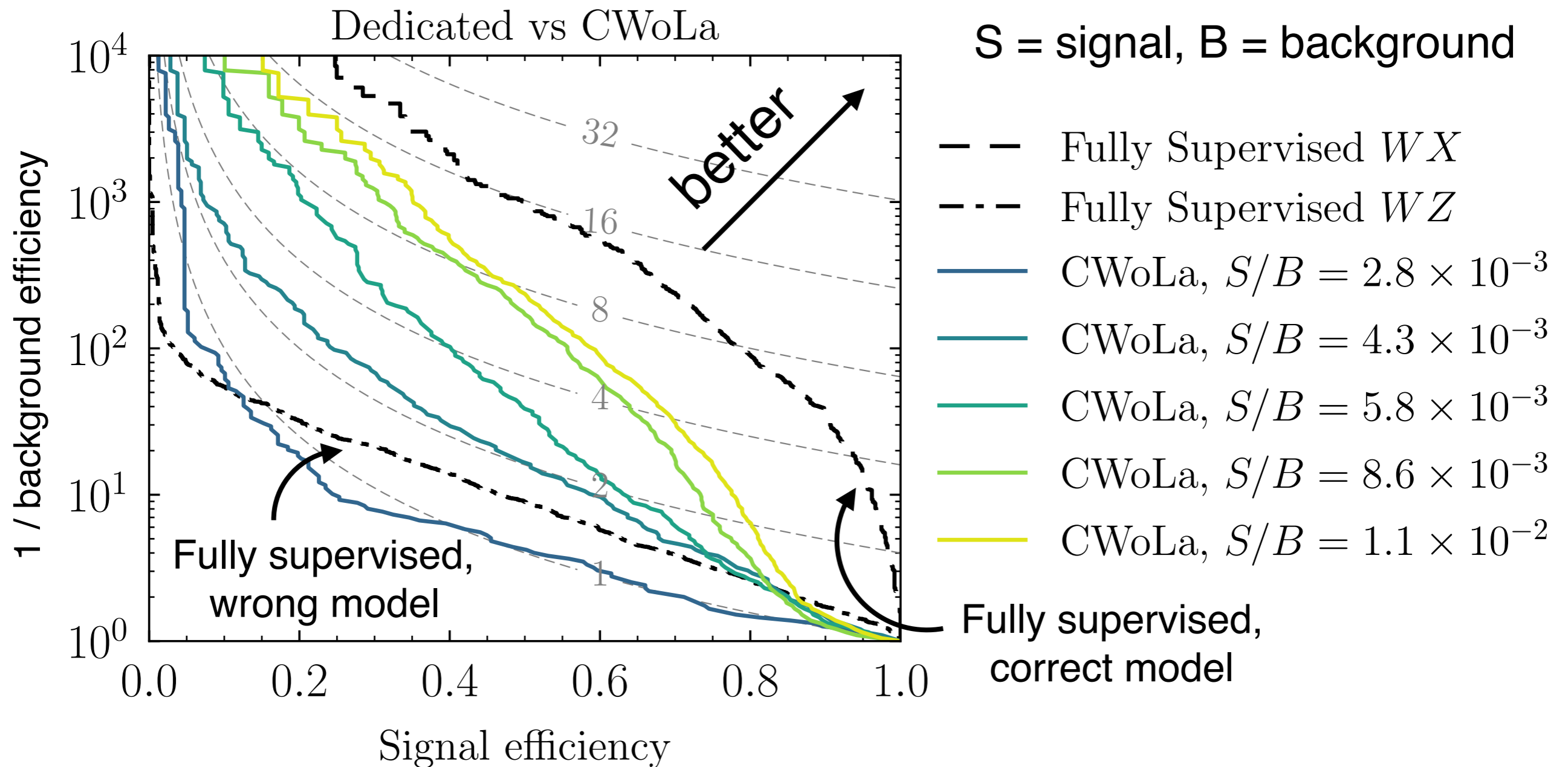


J. Collins, K. Howe, **BPN**, 1805.02664



Using a classifier trained to distinguish a signal region from a sideband, make progressively harsher cuts on the NN output

CWoLa hunting vs. Full Supervision



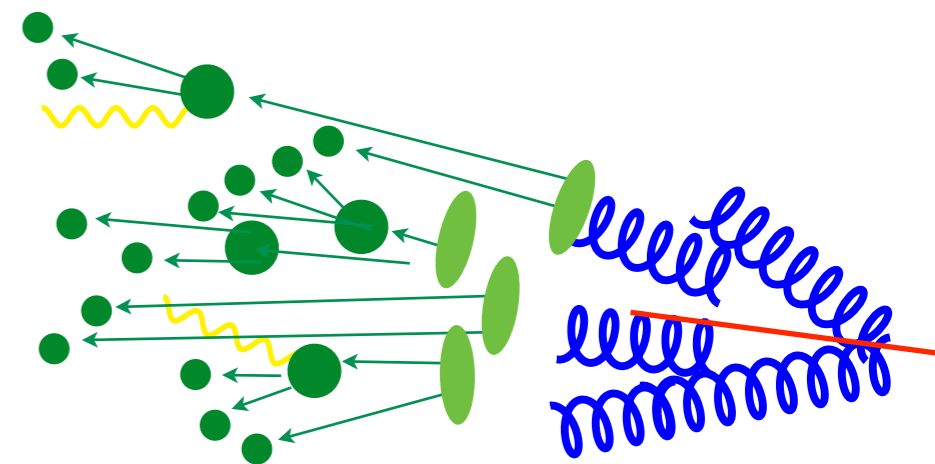
If you know what you are looking for, you should look for it. If you don't know, then CWoLa hunting may be able to catch it!

(1) Need an observable X (e.g. m_{JJ}) where the signal is localized and the background is not.

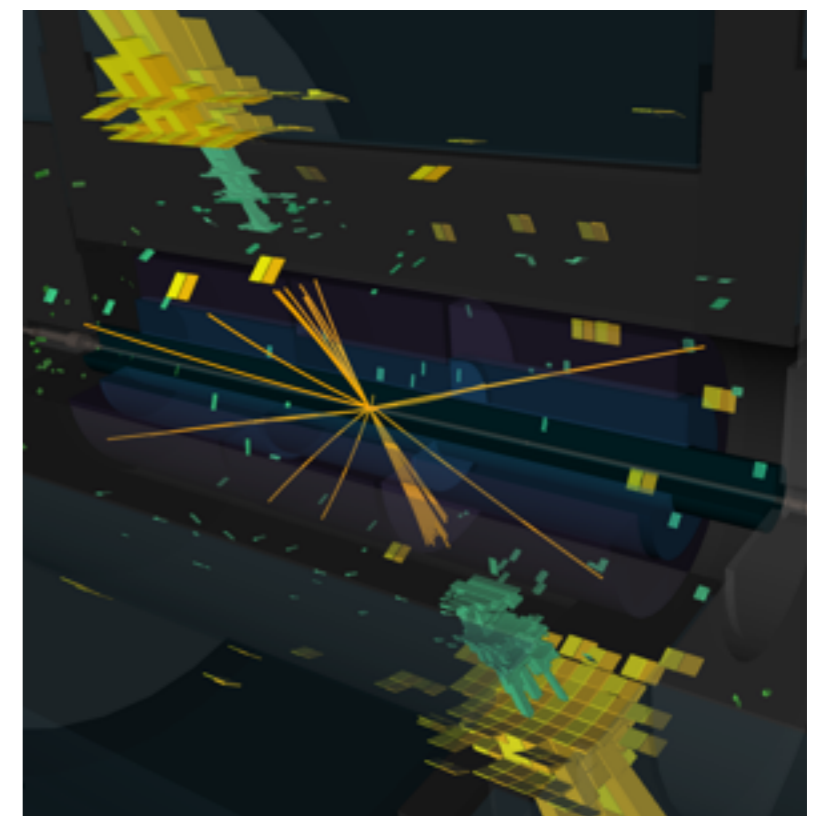
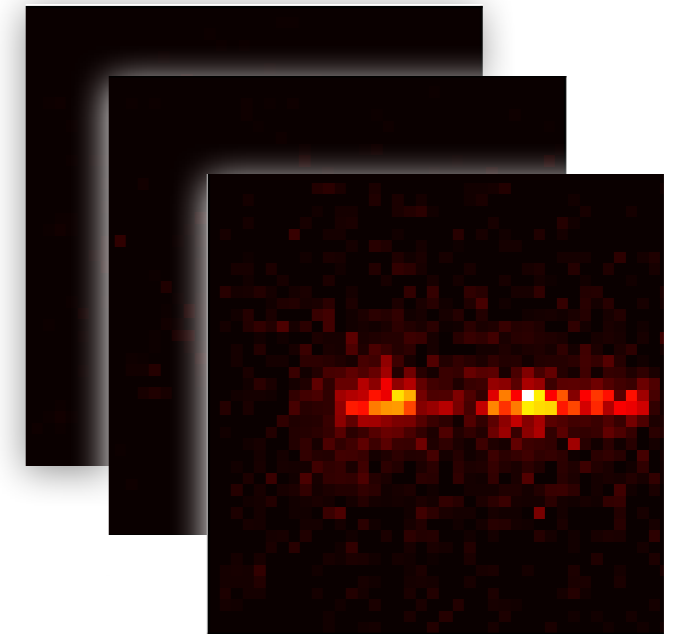
(2) Identify features Y (e.g. jet substructure) that are \sim independent of X , but can be useful for identifying a broad range of new particles.

actually, we don't need independence, we just need them to not allow us to sculpt bumps.

- **Generative models:** accelerate simulations with deep learning
- **Generalized numerical inversion:** machine learn with prior independence
- **Weak supervision:** learn directly from (unlabeled) data
- CWoLa hunting: model agnostic **anomaly detection**



Accelerating simulation/reducing simulation dependence can improve performance and may even help us to understand something new and fundamental about nature!



Classification

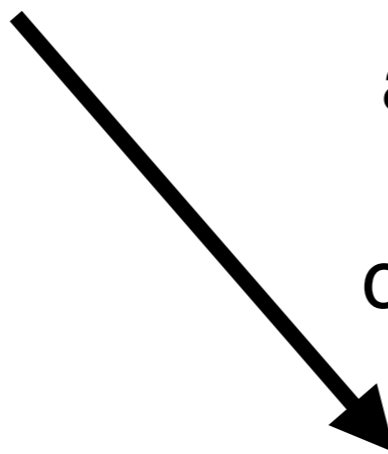
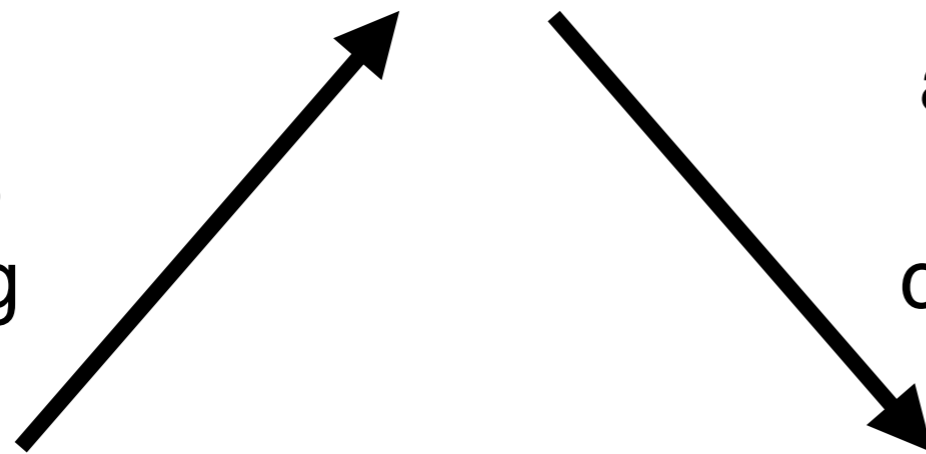
provide examples for training

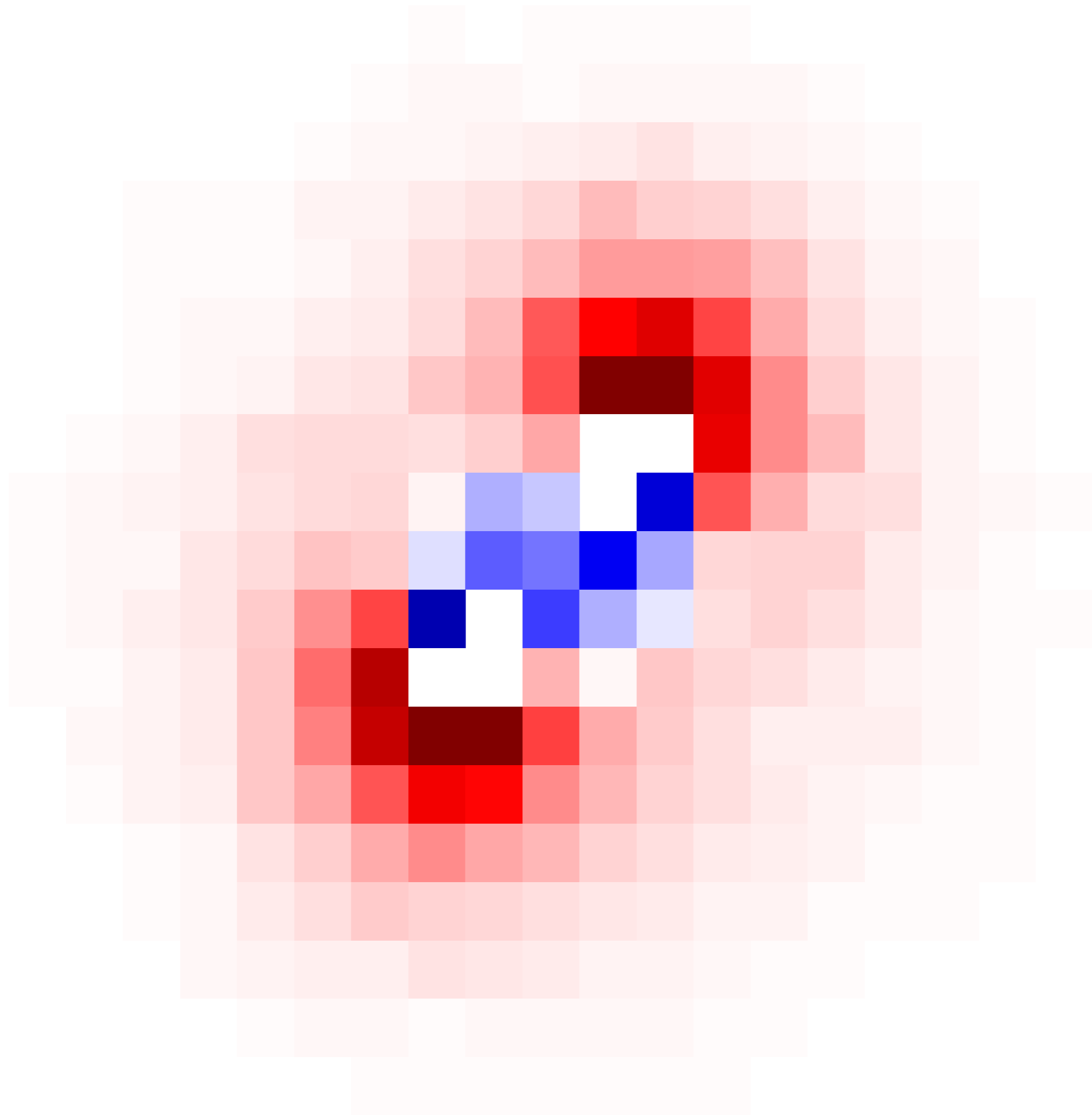
arbitrarily many categories

Generation

Regression

map noise to structure





Fin.