cylindrical space.

# Machine Learn
# High Energy Ph

Jet Image    Convolution    Max-Pool

## Benjamin Nachman

*Lawrence Berkeley National Laboratory*

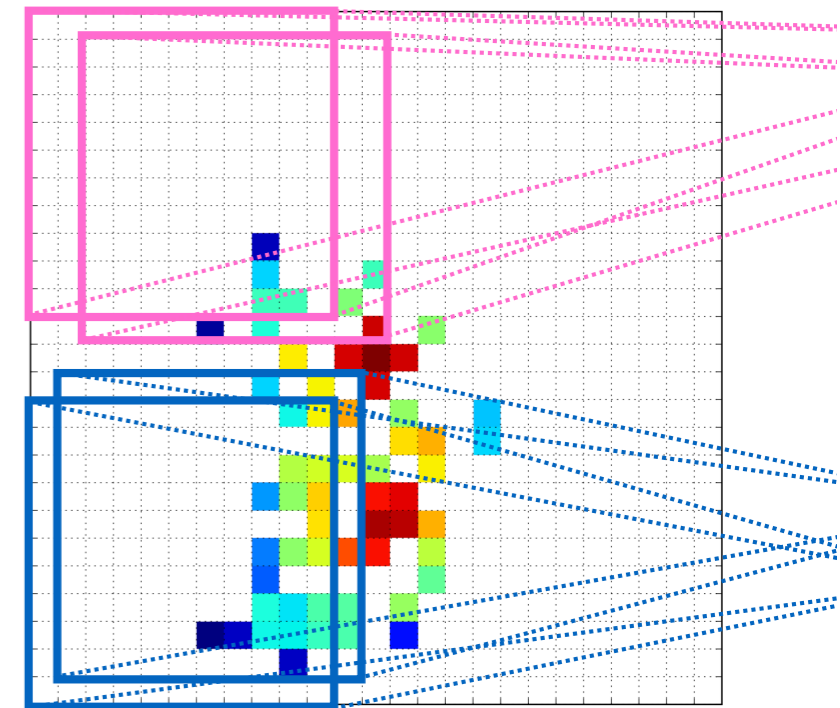cern.ch/bnachman    @bpnachman    bnachman
bpnachman@lbl.gov

**BERKELEY LAB**

**BERKELEY EXPERIMENTAL PARTICLE PHYSICS**

*Lund MCNet ML School*
*June 24, 2020*

Theory of everything

**Fast simulation**

**Parameter estimation / unfolding**

**Physics simulators**

Detector-level observables

↓

Pattern recognition

**Nature**

↓

Experiment

**Online processing & quality control**

↓

Detector-level observables

↓

Pattern recogn

**Data curation**

**Anomaly detection**

calibration
clustering
tracking
noise mitigation
particle identification

...

**Roadmap for this lecture**

Theory of everything

**Nature**

Fast simulation

Parameter estimation / unfolding

**Physics simulators**

Experiment

Online processing & quality control

Detector-level observables

Detector-level observables

Pattern recognition

Pattern recogn

Data curation

Anomaly detection

**calibration**

clustering

tracking

noise mitigation

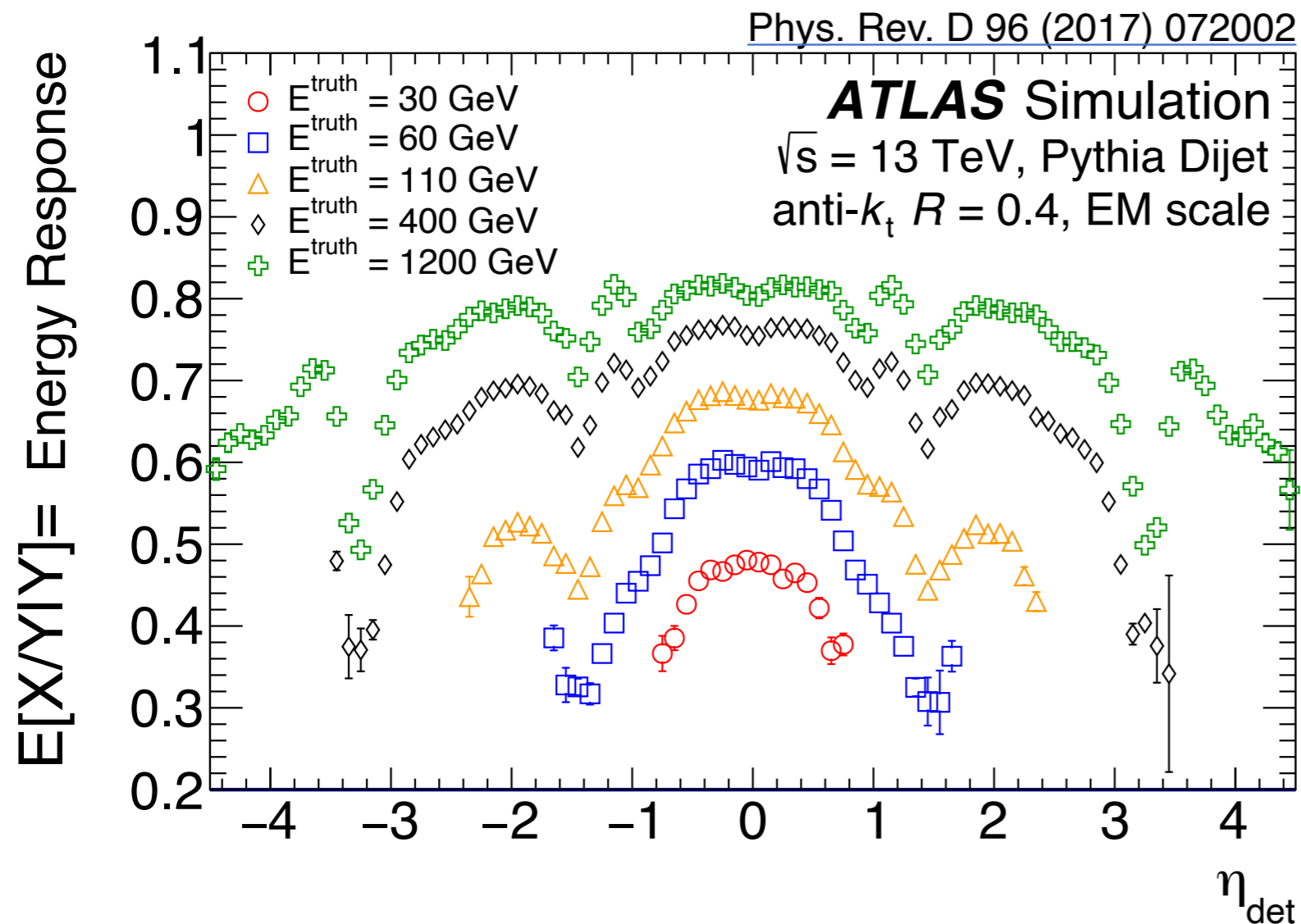particle identification

...

You measure X and want E[X|Y] = Y.

We'll discuss how to perform this calibration using machine learning.

One of the themes throughout this lecture will be mitigating simulation (prior) dependence.

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



Phys. Rev. D 96 (2017) 072002

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)

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 (MSE):

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

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

*Check this using the method we discussed yesterday!*

Could this be a problem?

$$f(x) = E[y|x] = \int \mathrm{d}y \, y \, p(y|x)$$

$$E[f(x)|y] = \int \mathrm{d}x \, \mathrm{d}y' \, y' \, p_{\mathrm{train}}(y'|x) \, p_{\mathrm{test}}(x|y)$$

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

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 $f^{-1}(x)$

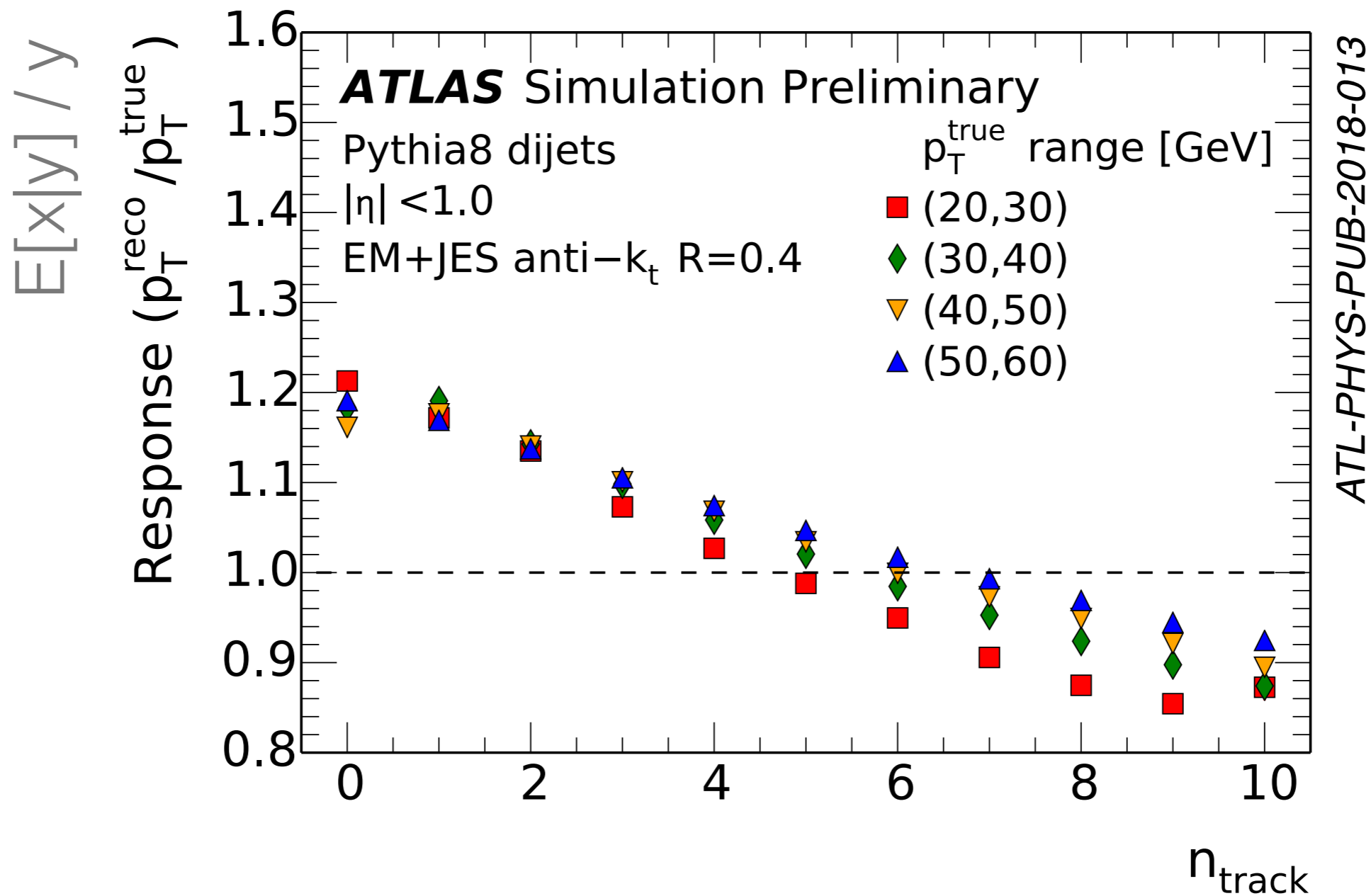by construction, f is independent of p(y) and thus $f^{-1}$ also does not depend on p(y), as desired.

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 MSE.

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.

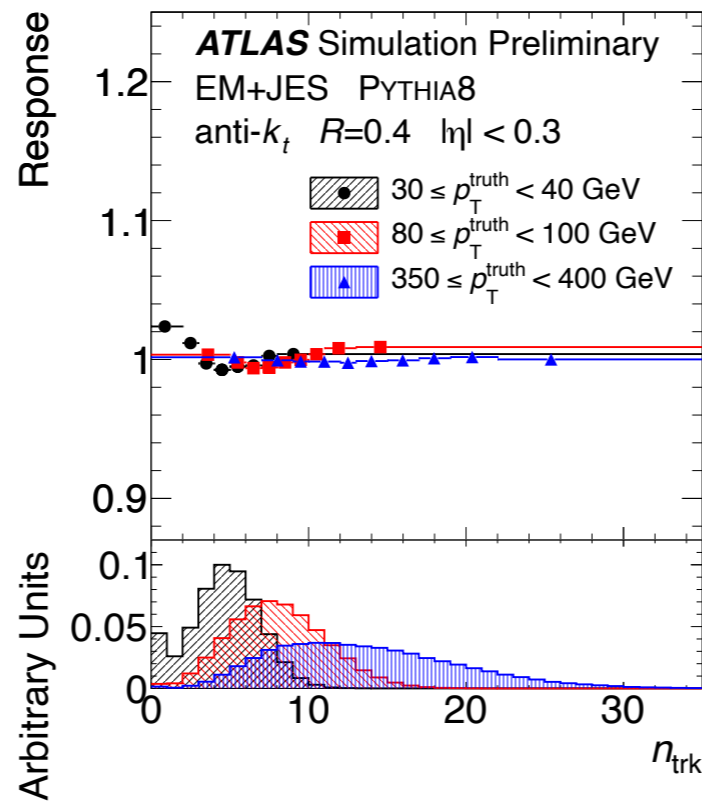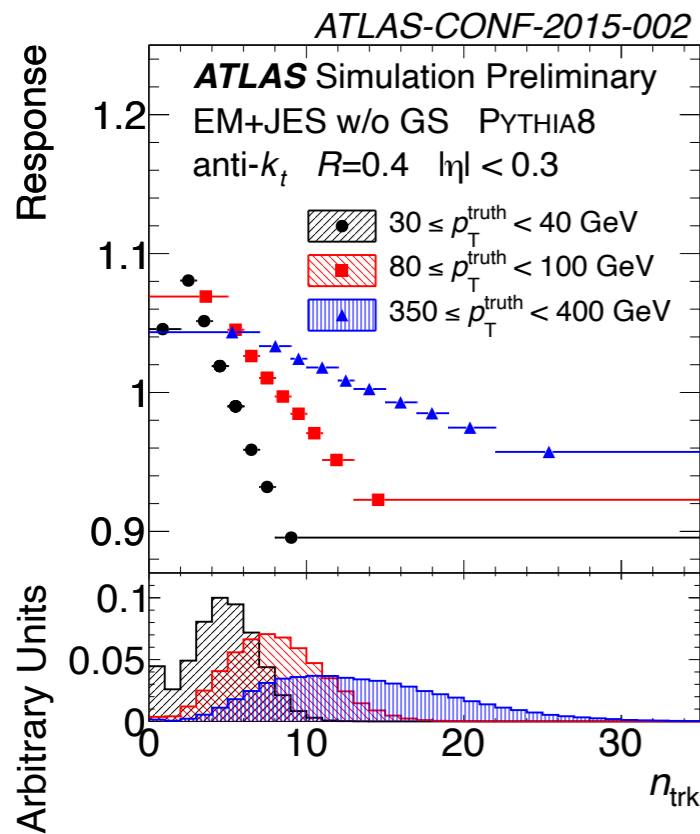The detector response of jets depends on many properties of the jet. Ideally, the calibration can include this!

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

$$p_{\mathrm{T}}^{\mathrm{reco}} \mapsto \hat{p}_{\mathrm{T}}^{\mathrm{reco}} = f_{\theta_n}^{-1}\left(\cdots f_{\theta_2}^{-1}\left(f_{\theta_1}^{-1}\left(p_{\mathrm{T}}^{\mathrm{reco}}\right)\right)\cdots\right)$$



ATLAS-CONF-2015-002

This works well when the jet response is independent of $\theta_i$ given $\theta_j$.

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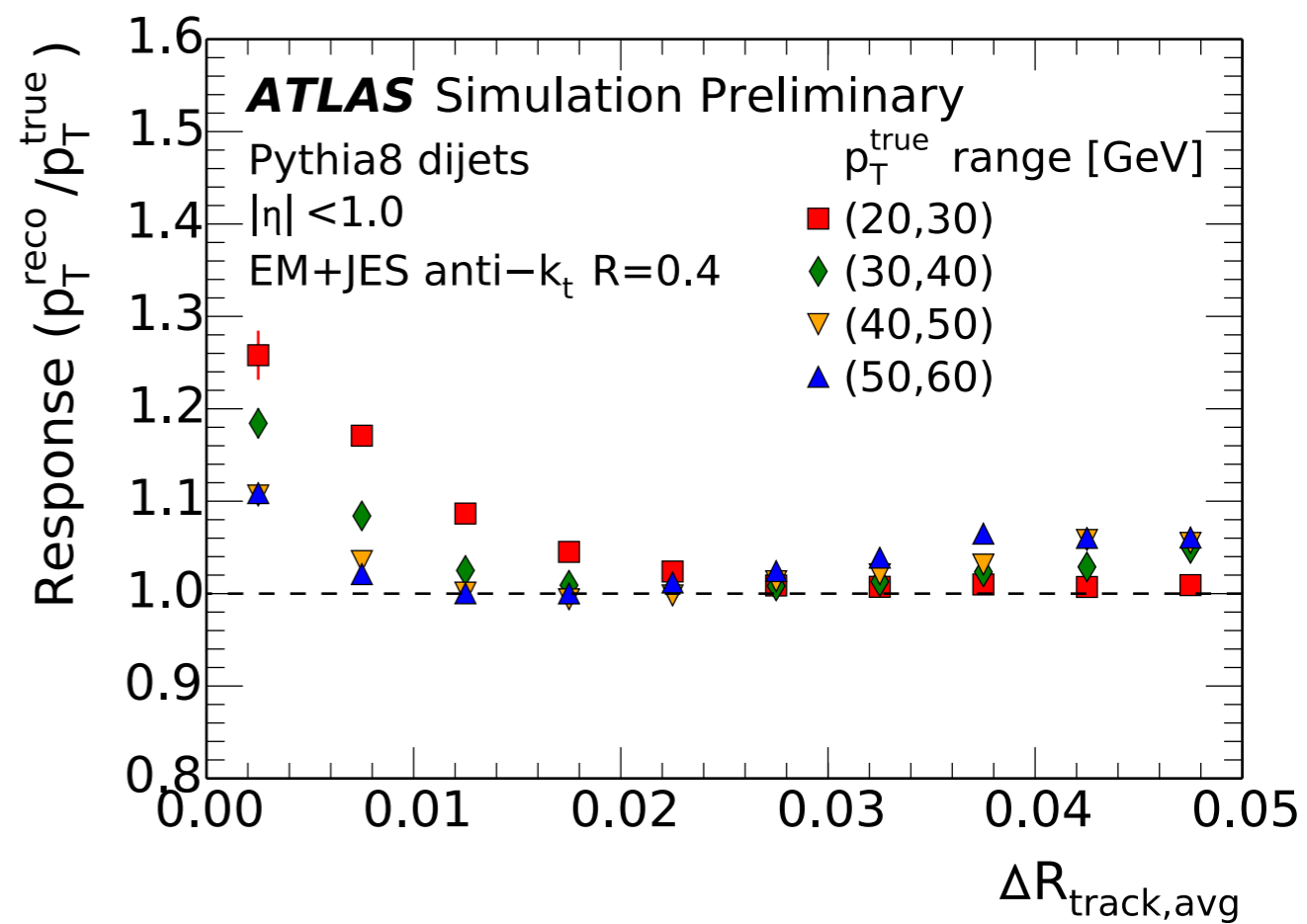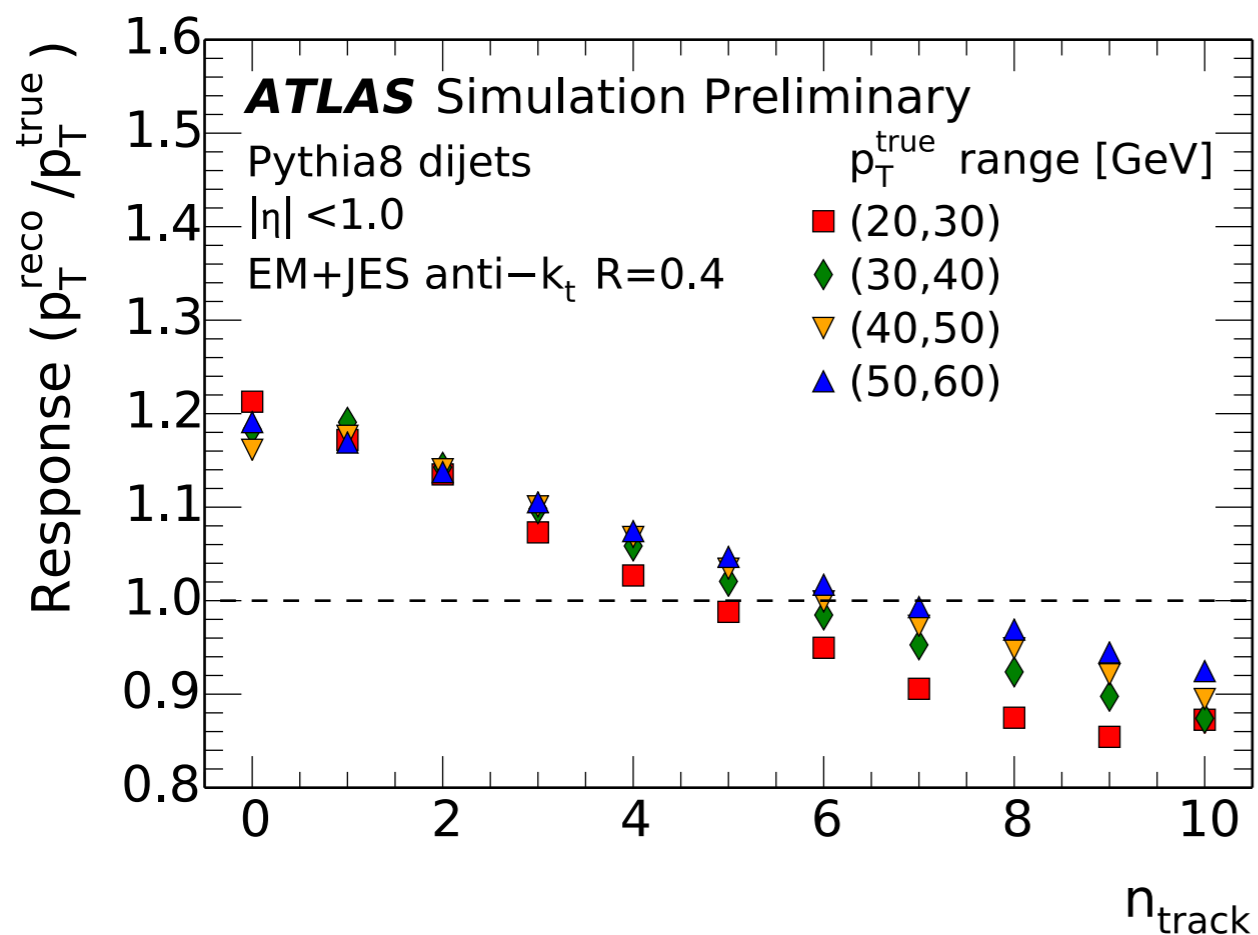
We cannot use numerical inversion out-of-the-box because we now have a many-to-one function.

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

(1) Learn a function f to predict x given y and all the $\theta$'s.

(2) For every combination of $\theta$, invert f.

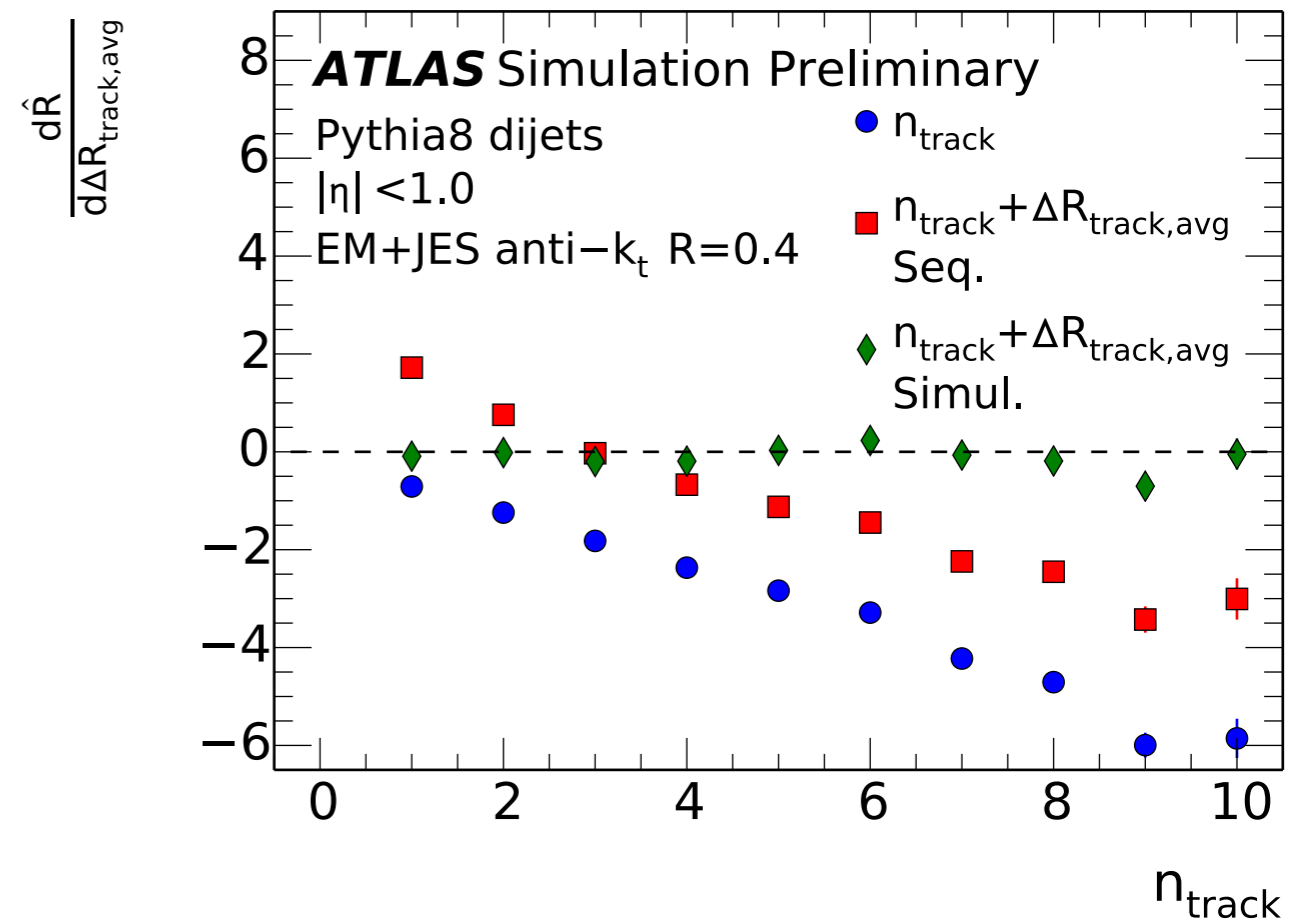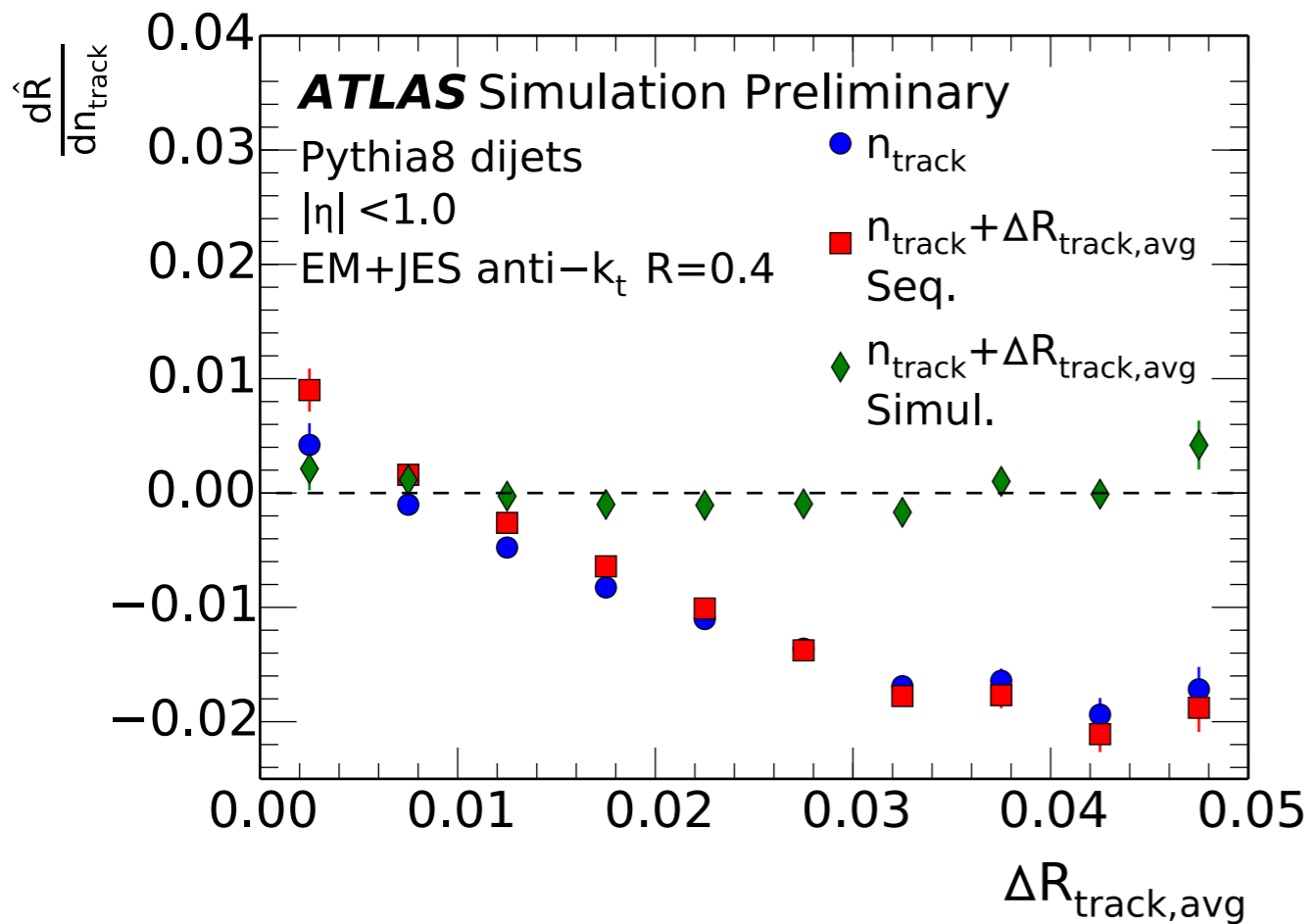(3) Calibrate via x → $f_\theta^{-1}(x)$

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

## Consider two features:



average track p_T-weighted
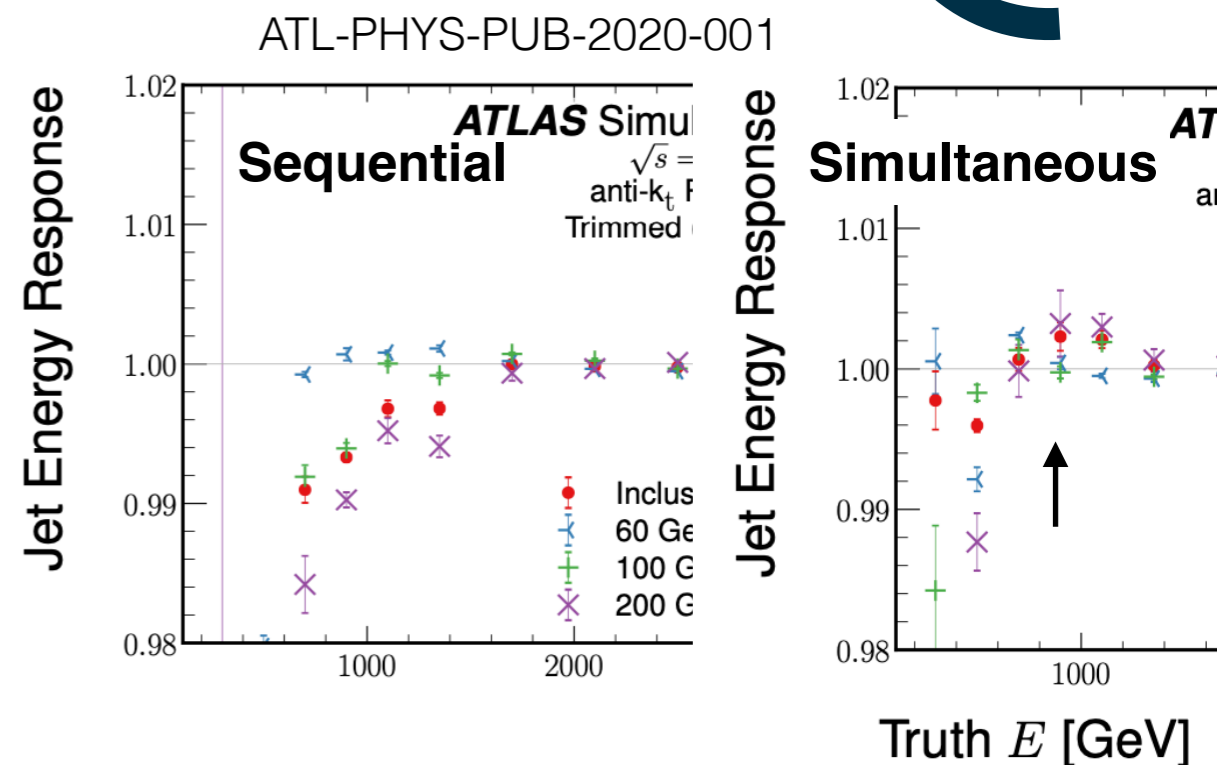distance from jet center

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



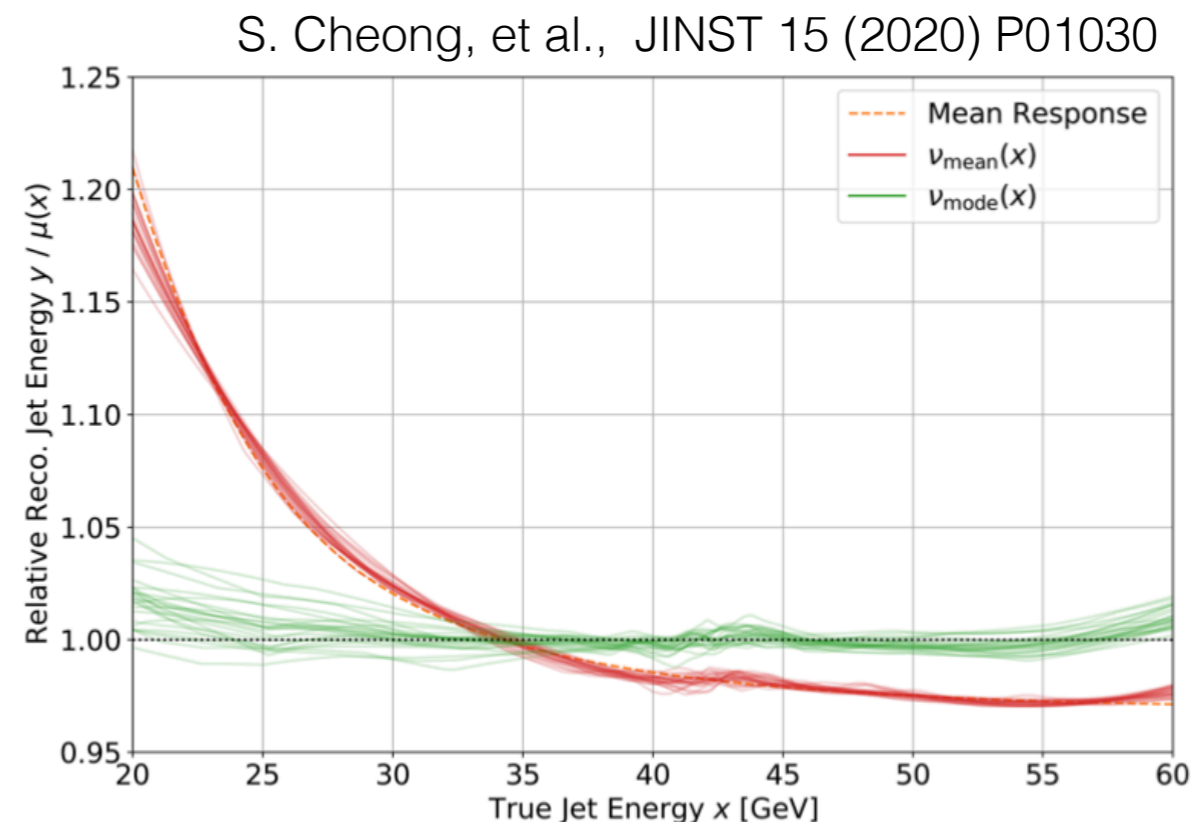Only the **simultaneous approach** removes the full residual dependence!

Can also simultaneously calibrate a subset of the θ's (e.g. jet energy and mass)

In many cases, it is desirable to calibrate the mode and not the mean since p(X|Y) is asymmetric.

*Can achieve this with modified loss function!*

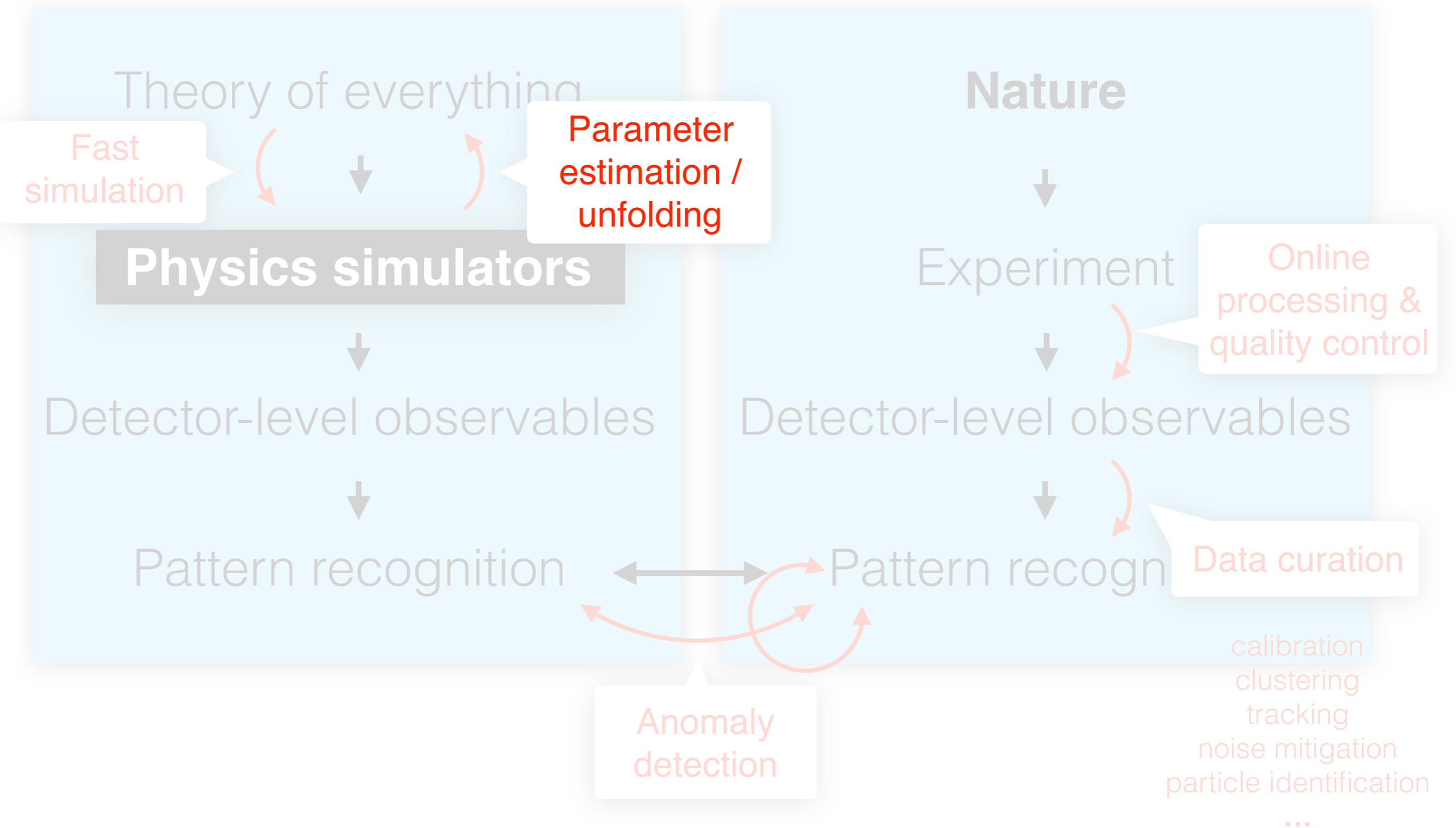ATL-PHYS-PUB-2020-001



S. Cheong, et al.,  JINST 15 (2020) P01030

There are many more applications of regression in HEP, but calibration is a prototypical example.

When building a regression model, it is critical to be wary of prior dependence and to pick the loss function based on what you actually want to learn (mean/median/mode/IQR/etc.)

Theory of everything

**Nature**

Fast simulation

Parameter estimation / unfolding

**Physics simulators**

Experiment

Online processing & quality control

Detector-level observables

Detector-level observables

Pattern recognition

Pattern recogn

Data curation

Anomaly detection

calibration
clustering
tracking
noise mitigation
particle identification
...

Key **challenge** and **opportunity**: *hypervariate phase space*
& *hyper spectral data*

*Image inspired by JHEP 02 (2009) 007*

Typical collision events
at the LHC produce
**O(1000+)** particles

We detect these
particles with
**O(100 M)**
readout channels



Not to scale!

Key **challenge** and **opportunity**: *hypervariate phase space*
& *hyper spectral data*

*Image inspired by JHEP 02 (2009) 007*

Typical collision events at the LHC produce **O(1000+)** particles

We detect these particles with **O(100 M)** readout channels



Not to scale!

**Want this**

**Measure this**



i.e. remove detector distortions

If you know $p(\textbf{\textit{meas.}} \mid \textbf{\textit{true}})$, could do maximum likelihood, i.e.

$$\text{unfolded} = \underset{true}{\text{argmax}}\ p(\textbf{\textit{measured}} \mid \textbf{\textit{true}})$$



**Want this**

**Measure this**

$p(\textbf{\textit{meas.}} \mid \textbf{\textit{true}})$ = "response matrix" or "point spread function"

If you know *p(**meas.** / **true**)*, could do maximum likelihood, i.e.

> *unfolded = argmax p(**measured** / **true**)*
>                     *true*

⚠️ Challenge: **measured** is hyperspectral and **true** is hypervariate … *p(**meas.** | **true**) is **intractable** !*

*p(**meas.** / **true**) = "response matrix" or "point spread function"*

If you know *p(**meas.** | **true**)*, could do maximum likelihood, i.e.

*unfolded = argmax p(**measured** | **true**)*
*true*

⚠️ Challenge: **measured** is hyperspectral and **true** is hypervariate … *p(**meas.** | **true**) is **intractable** !*

However: we have **simulators** that we can use to sample from *p(**meas.** | **true**)*

→ **Simulation-based (likelihood-free) inference**

*p(**meas.** | **true**)* = "response matrix" or "point spread function"

I'll briefly show you one solution to give you a sense of the power of likelihood-free inference.

I'll briefly show you one solution to give you a sense of the power of likelihood-free inference.

The solution will be built on **reweighting**

dataset 1: sampled from **p(x)**
dataset 2: sampled from **q(x)**

Create weights **w(x) = q(x)/p(x)** so that when dataset 1 is weighted by **w**, it is statistically identical to dataset 2.

I'll briefly show you one solution to give you a sense of the power of likelihood-free inference.

The solution will be built on **reweighting**

dataset 1: sampled from **p(x)**
dataset 2: sampled from **q(x)**

Create weights **w(x) = q(x)/p(x)** so that when dataset 1 is weighted by **w**, it is statistically identical to dataset 2.

What if we don't (and can't easily) know **q** and **p**?

**Fact**: Neutral networks learn to approximate the likelihood ratio = $q(x)/p(x)$

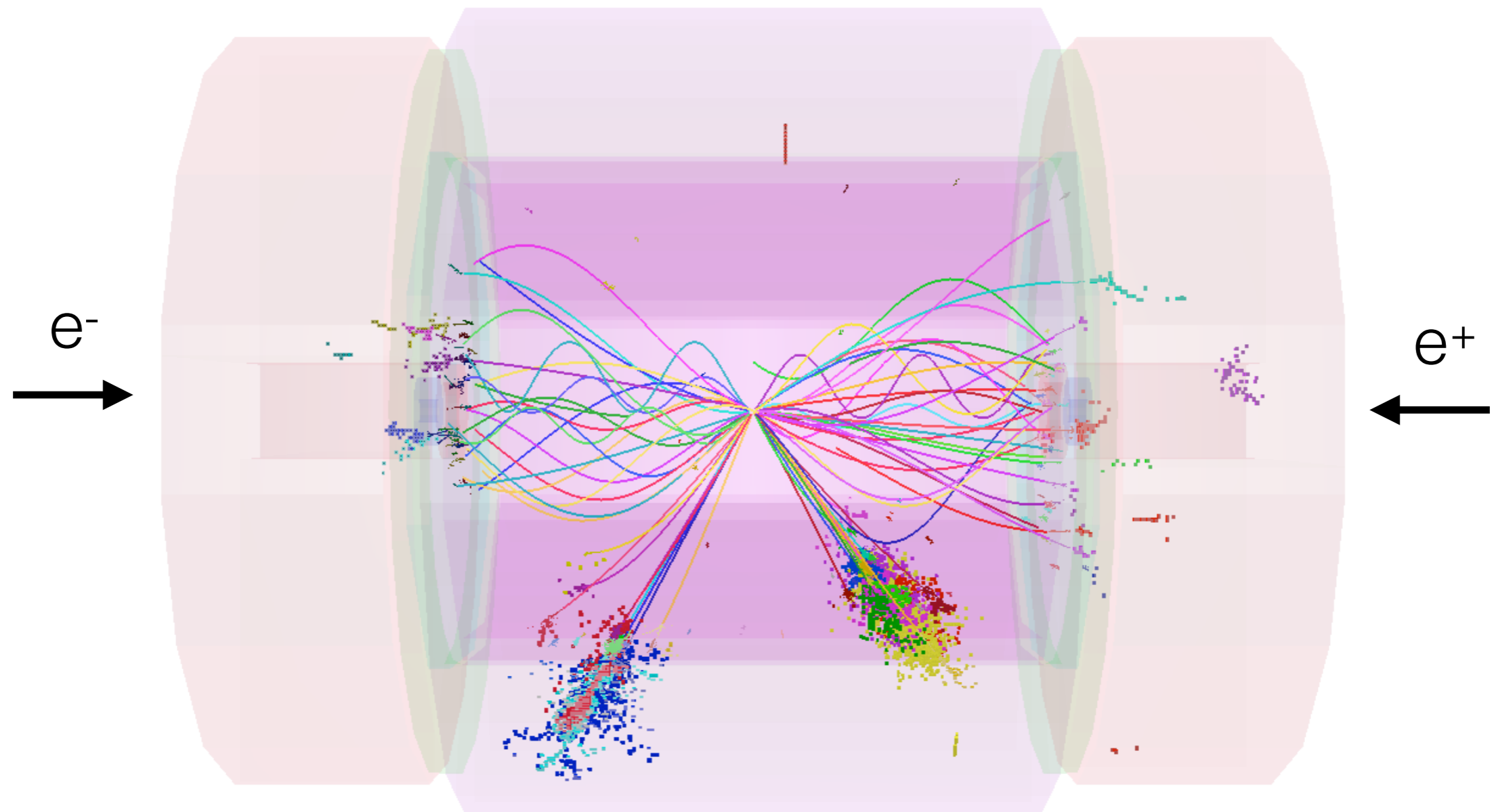(see previous lecture! Can you derive the monotonic relation?)

Solution: train a neural network to distinguish the two datasets!

This turns the problem of **density estimation** (**hard**) into a problem of **classification** (**easy**)

Particularly useful for particle physics, where collisions may produce a variable # of particles which are interchangeable*



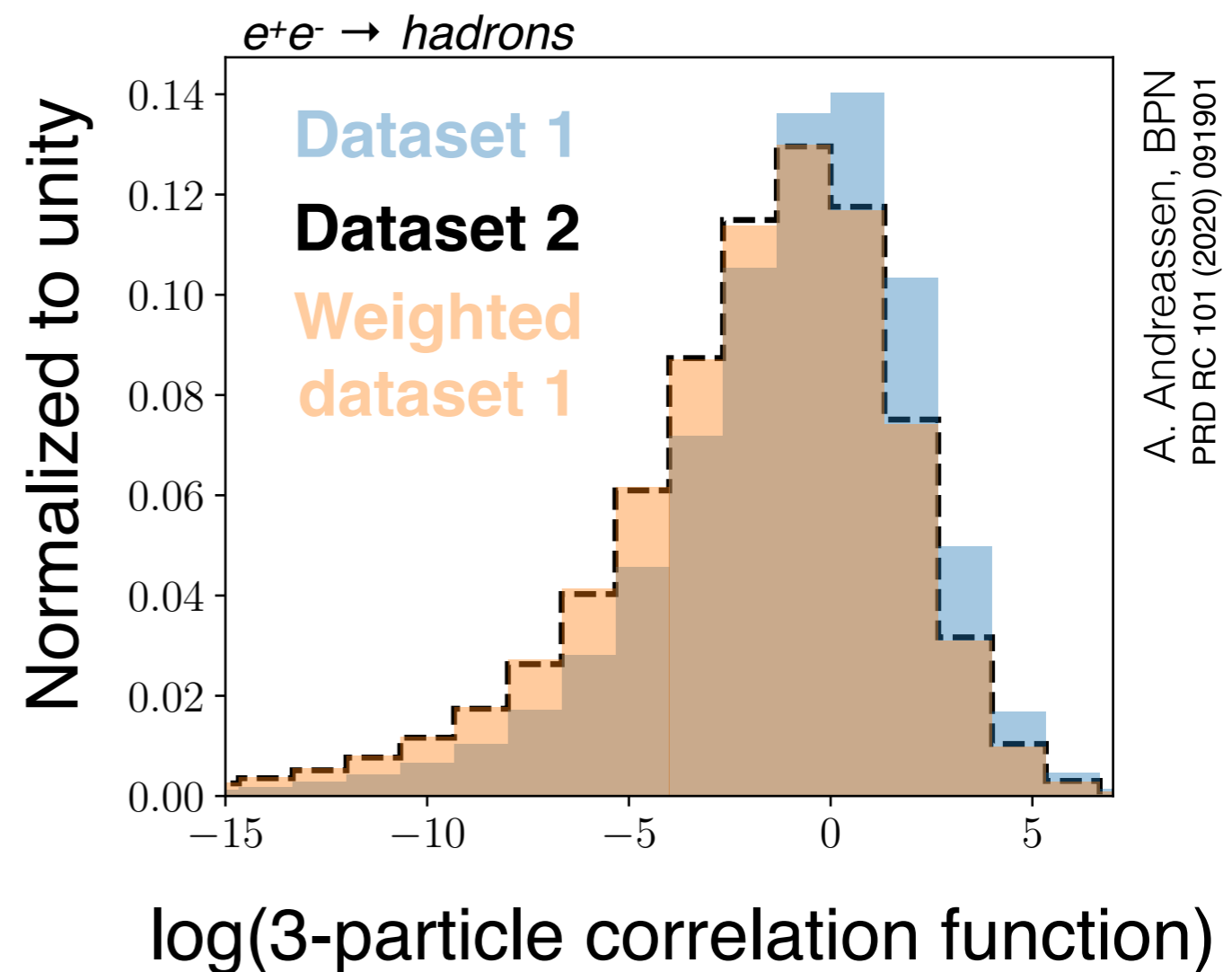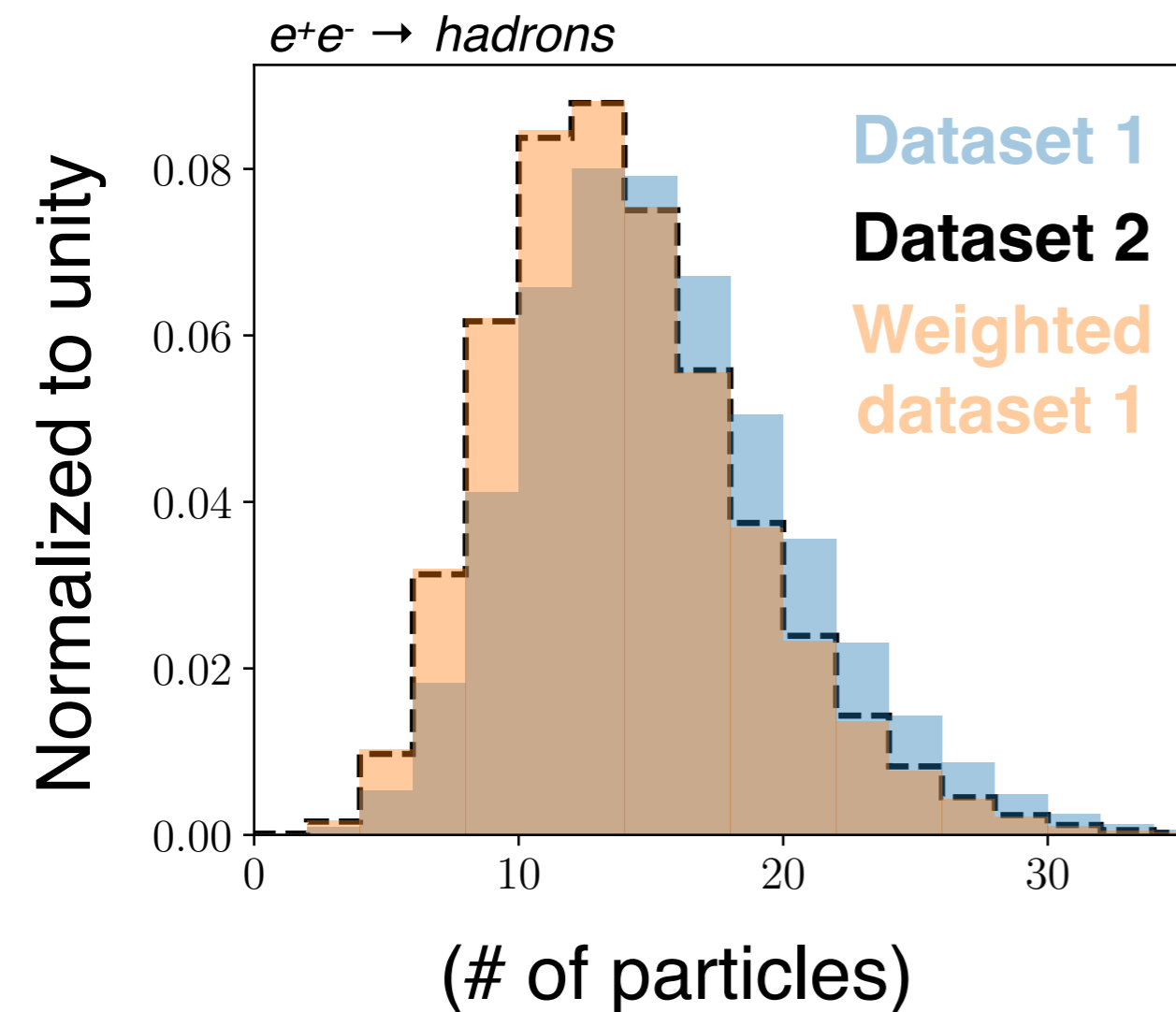$e^-$ →

← $e^+$

*deep learning architecture: deep sets, Zaheer et al., NIPS 2017, Komiske, Metodiev, Thaler, JHEP 01 (2019) 121

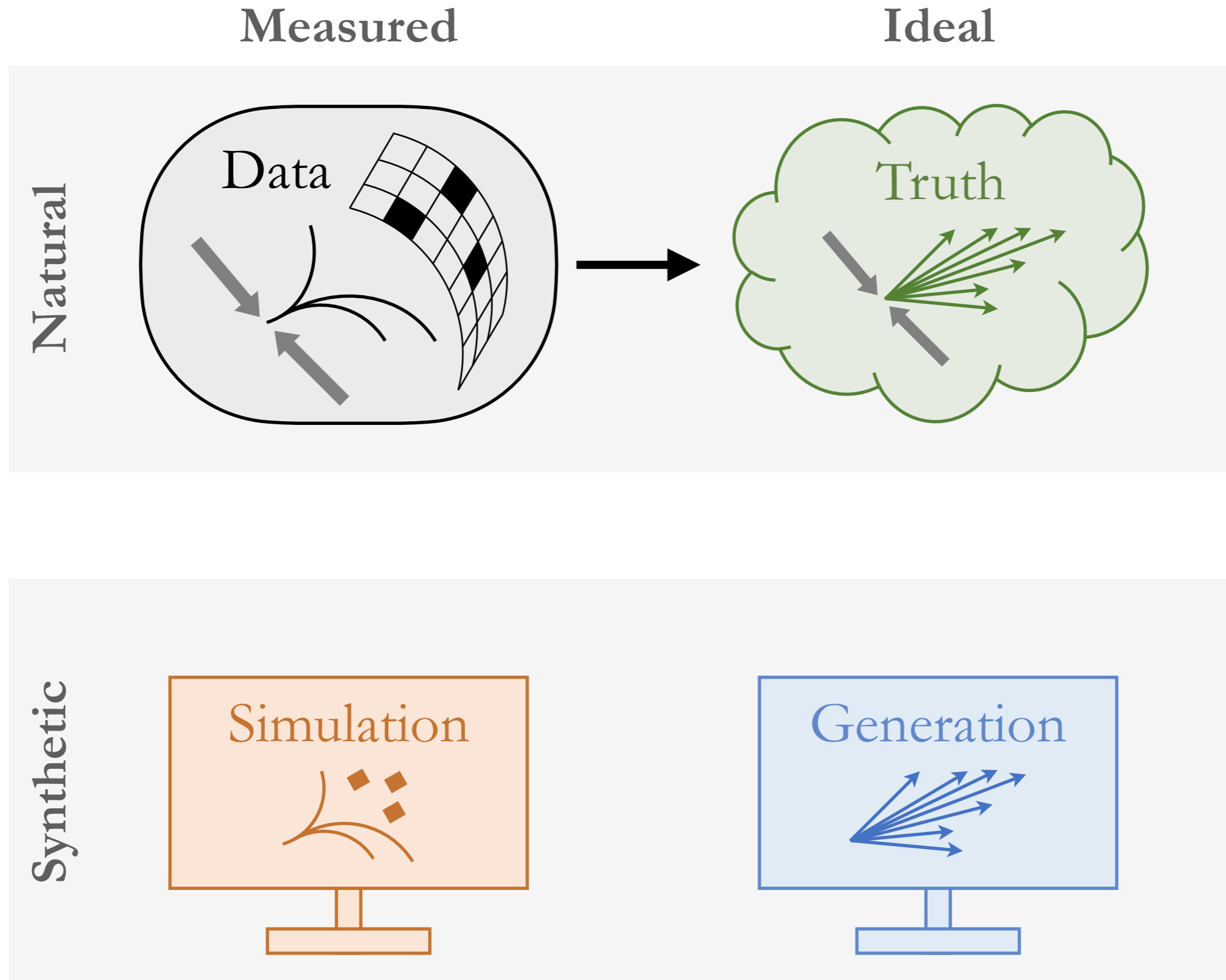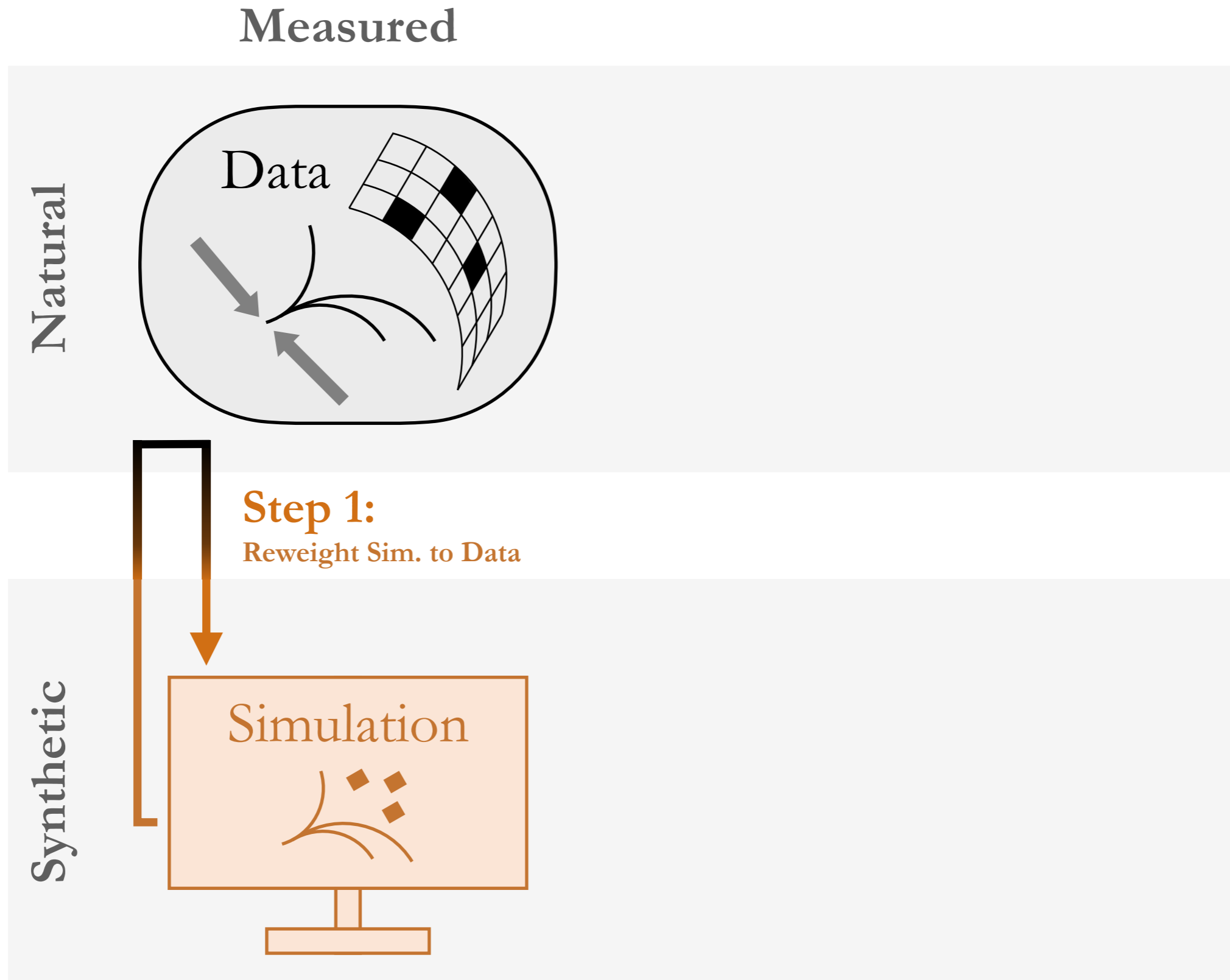Reweight the **full phase space** and then check for various binned 1D observables.



$e^+e^- \rightarrow$ *hadrons*

Normalized to unity

**Dataset 1**
**Dataset 2**
**Weighted dataset 1**

(# of particles)



$e^+e^- \rightarrow$ *hadrons*

Normalized to unity

**Dataset 1**
**Dataset 2**
**Weighted dataset 1**

A. Andreassen, BPN
PRD RC 101 (2020) 091901

log(3-particle correlation function)

We call this Deep neural networks using Classification for Tuning and Reweighting or "**DCTR**"

**Measured**

**Natural**

Data

**Synthetic**

**Step 1:**

**Reweight Sim. to Data**

Simulation

**Measured**

**Ideal**

Natural

Data

Step 1:
Reweight Sim. to Data

Step 2:
Reweight Gen.

Synthetic

Simulation

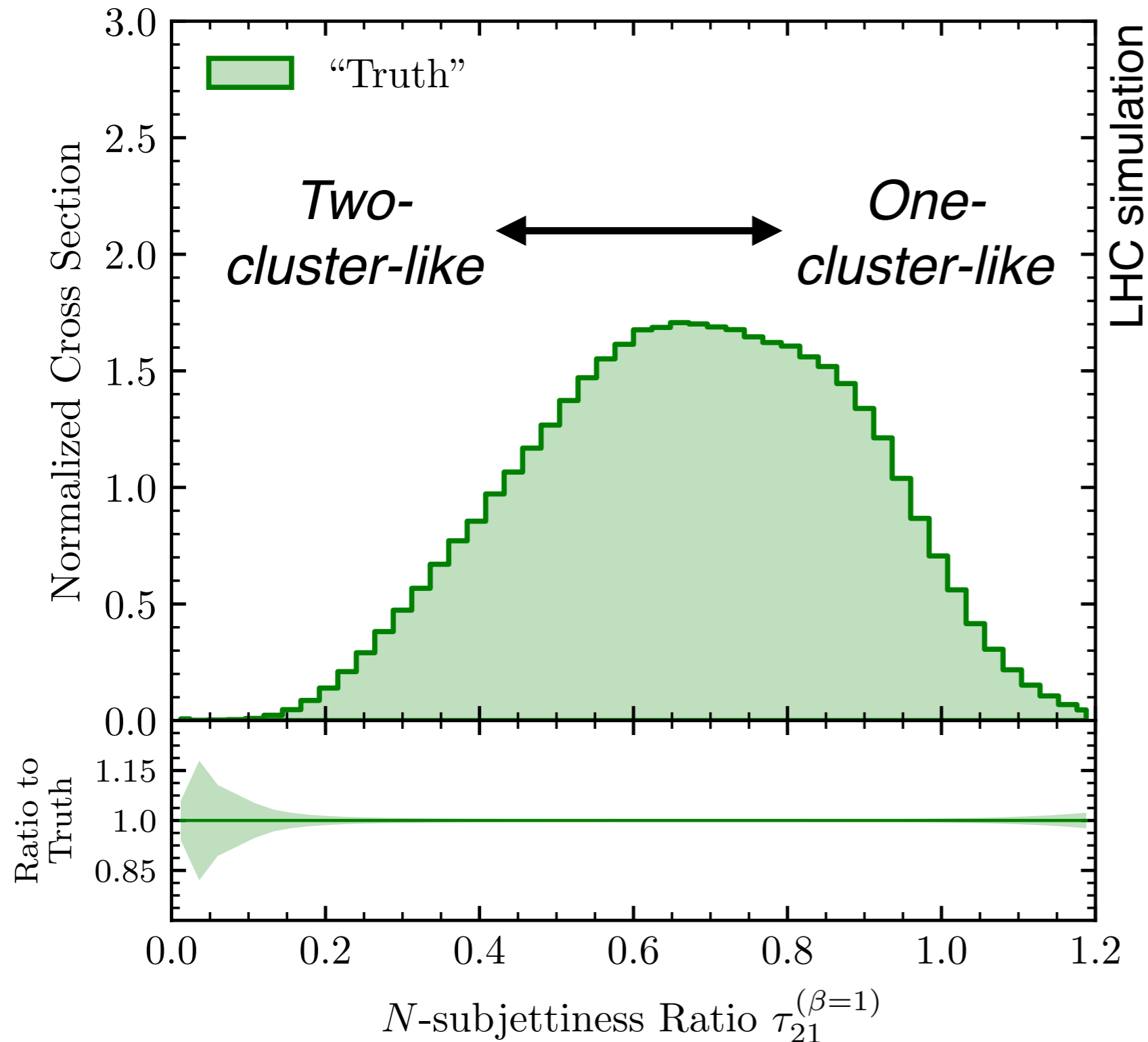Generation

Push Weights

A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, PRL 124 (2020) 182001

Consider this observable, which characterizes the substructure

A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, PRL 124 (2020) 182001



LHC simulation

Consider this observable, which characterizes the substructure

A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, PRL 124 (2020) 182001



*Two-cluster-like*

*One-cluster-like*

"Data"

"Truth"

LHC simulation

Normalized Cross Section

Ratio to Truth

$N$-subjettiness Ratio $\tau_{21}^{(\beta=1)}$

Consider this observable, which characterizes the substructure

Measured

Ideal

Data

Truth

Natural

Synthetic

Simulation

Generation
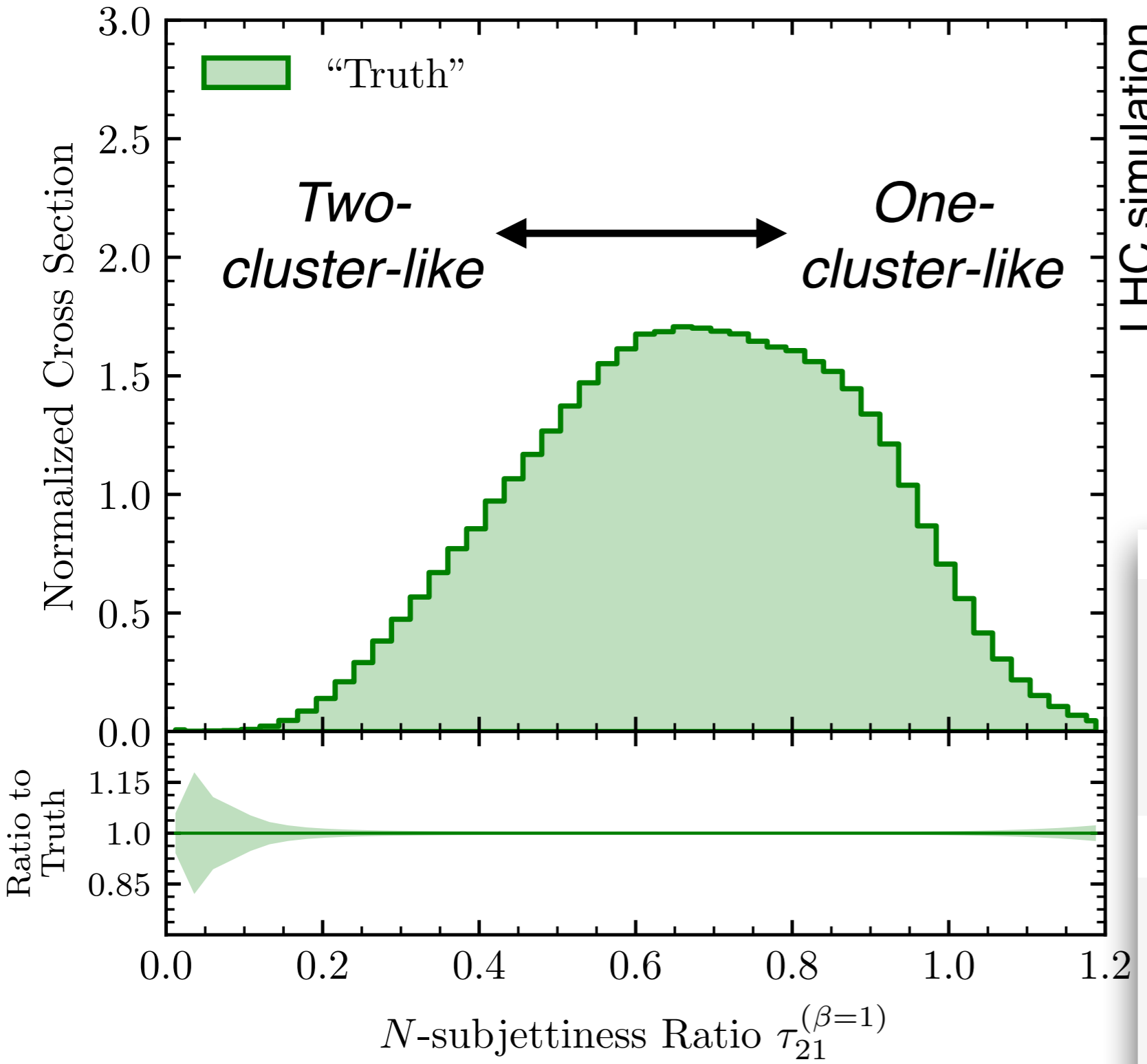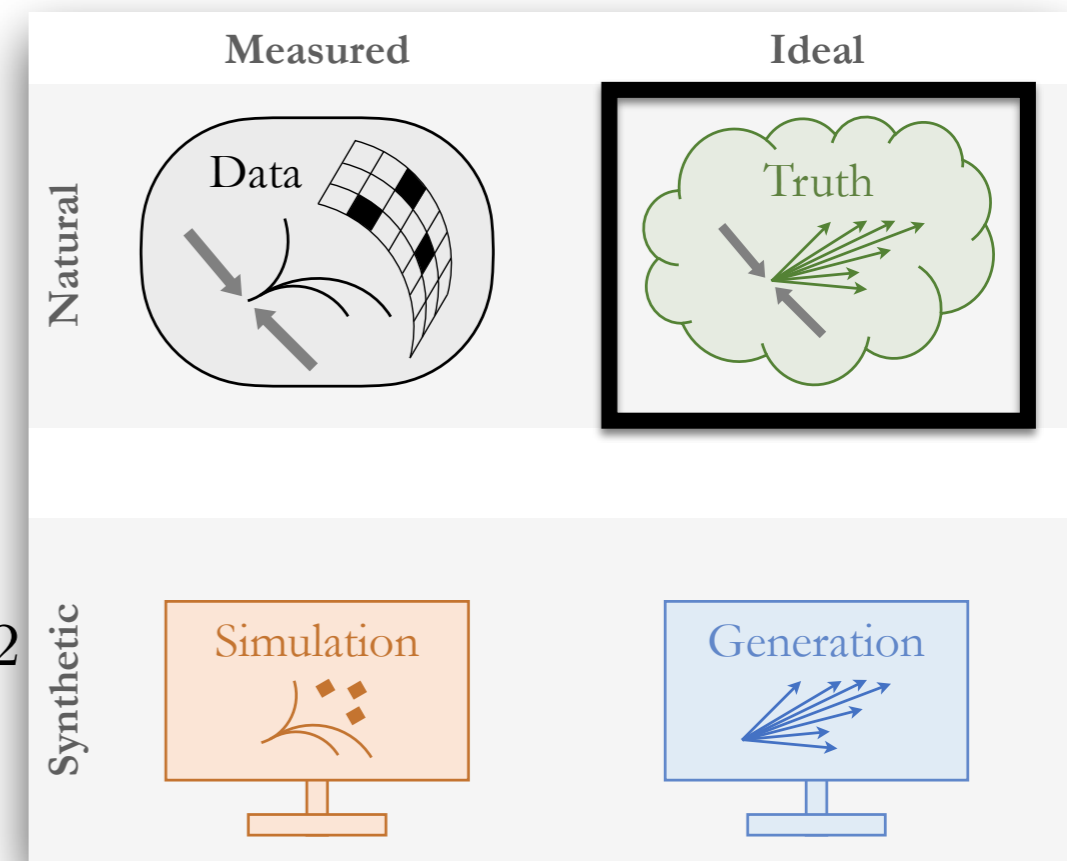
A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, PRL 124 (2020) 182001



Consider this observable, which characterizes the substructure

LHC simulation

A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, PRL 124 (2020) 182001

Consider this observable, which characterizes the substructure

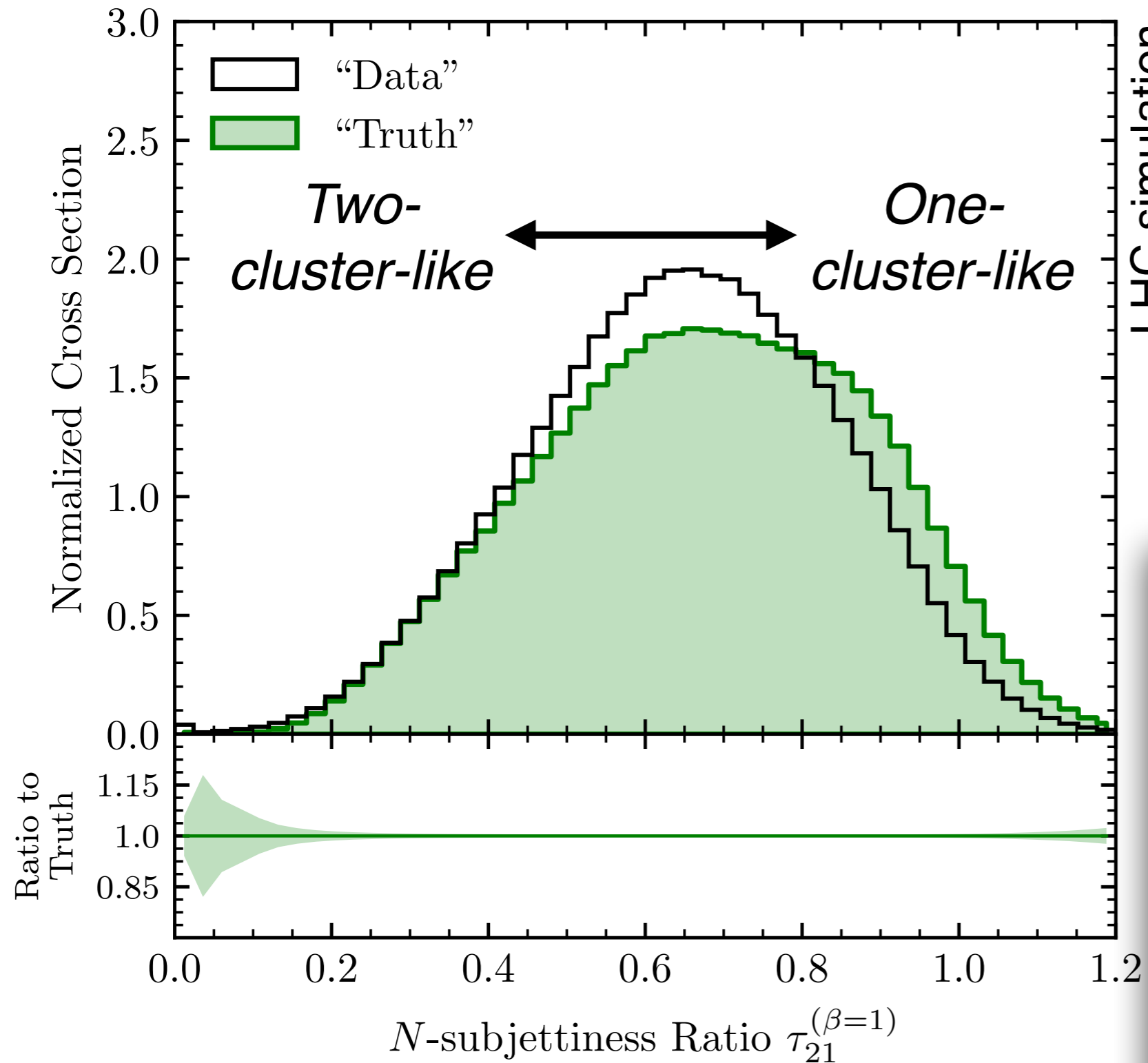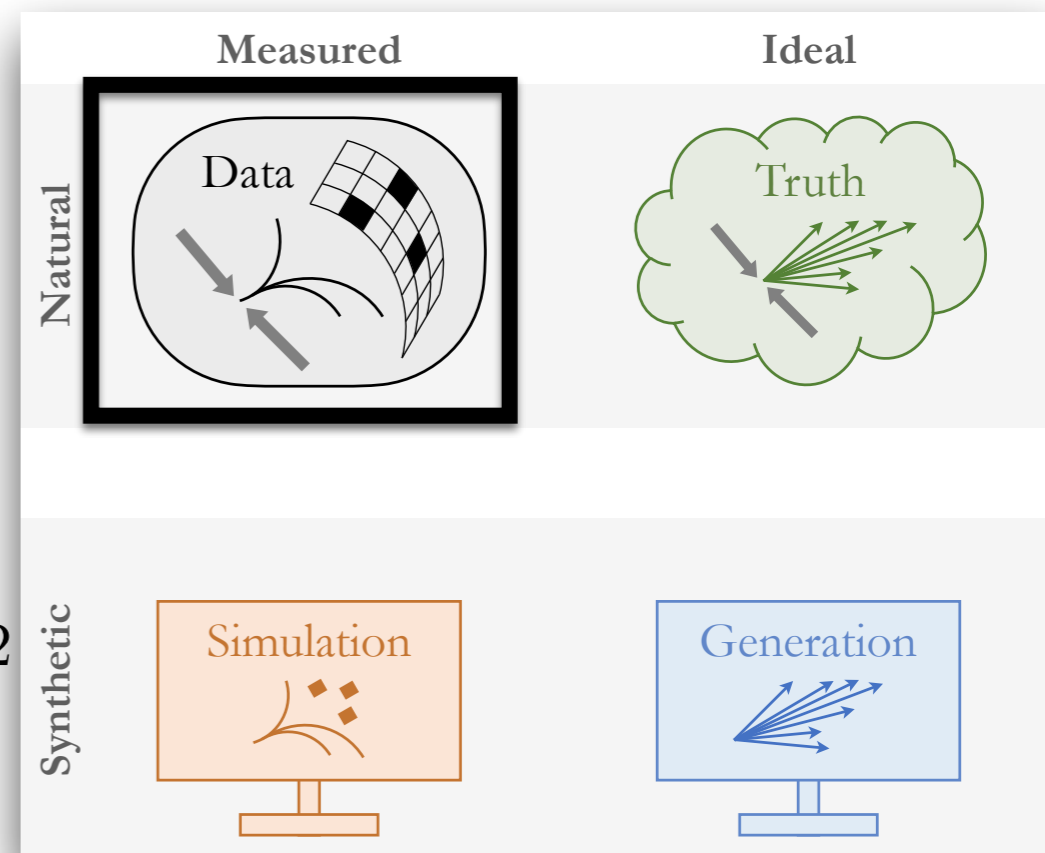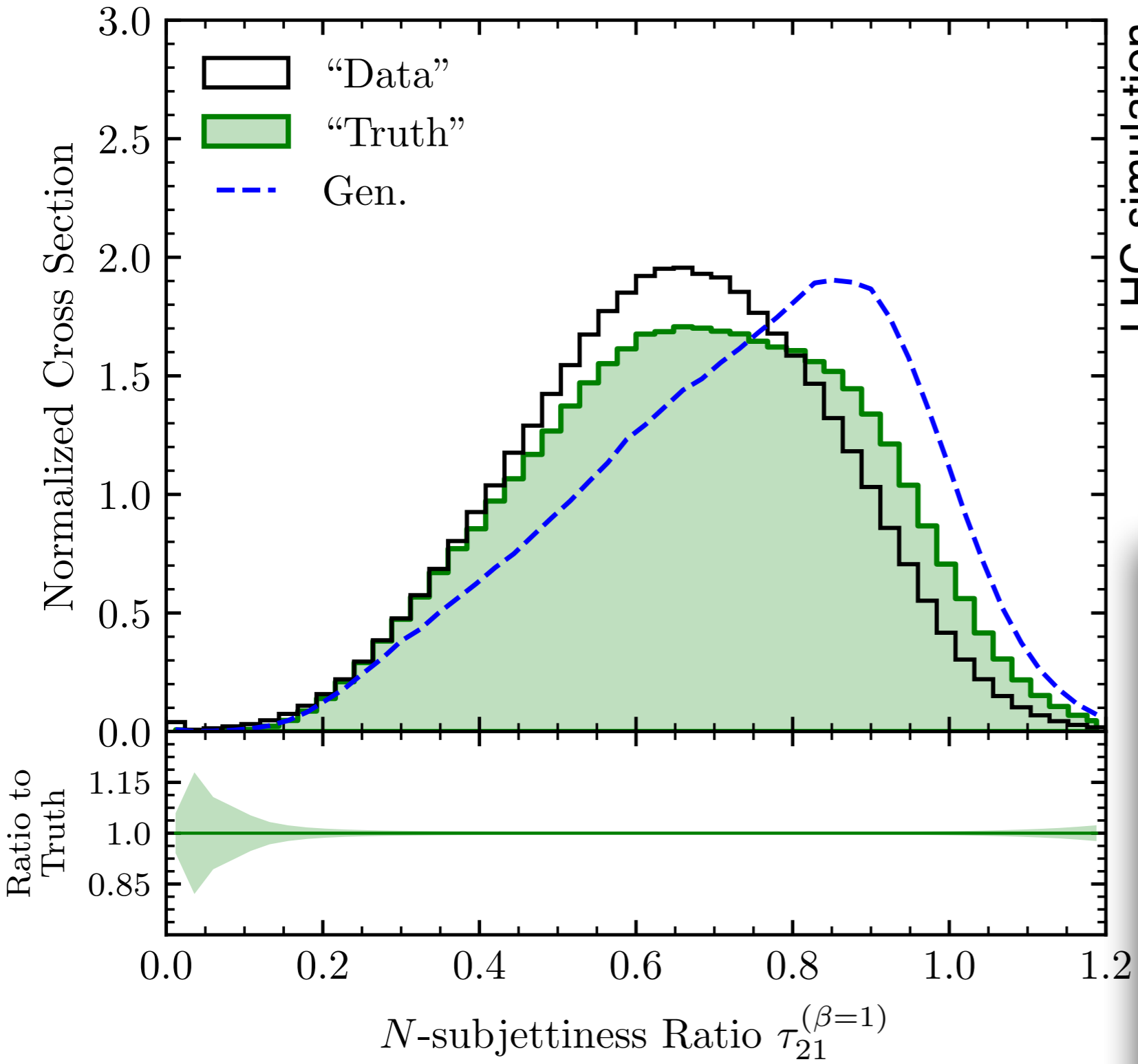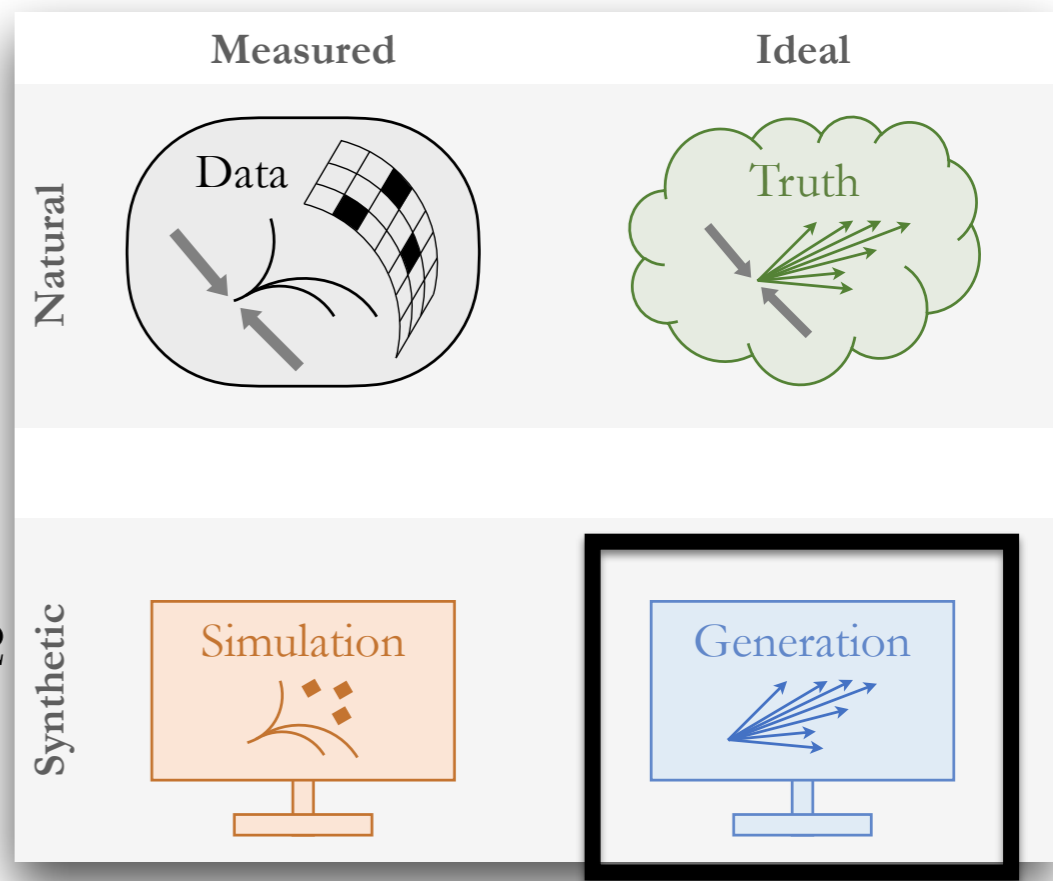A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, PRL 124 (2020) 182001



LHC simulation

IBU is the current standard. It is a 1D binned and iterative approach.

A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, PRL 124 (2020) 182001



LHC simulation

OmniFold outperforms IBU even though it is not tailored to this observable

A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, PRL 124 (2020) 182001

[A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, 1907.08209]



OmniFold is:
- ***Unbinned***
- *Maximum likelihood*
- *Full phase space (compute observables post-facto)*
- *Improves the resolution from auxiliary features*

[A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, 1907.08209]



OmniFold is:
- *Unbinned*
- ***Maximum likelihood***
- *Full phase space (compute observables post-facto)*
- *Improves the resolution from auxiliary features*

[A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, 1907.08209]

OmniFold is:
- *Unbinned*
- *Maximum likelihood*
- ***Full phase space*** *(compute observables post-facto)*
- *Improves the resolution from auxiliary features*

[A. Andreassen, P. Komiske, E. Metodiev, BPN, J. Thaler, 1907.08209]



# OmniFold is:
- *Unbinned*
- *Maximum likelihood*
- *Full phase space (compute observables post-facto)*
- ***Improves the resolution from auxiliary features***

extreme example:

$$\text{measured}|\text{true} = \text{true} + X$$

$$X \sim \mathcal{N}(\mu, \sigma)$$

If you control for X (=auxiliary feature), response is a delta-function!

One of the features of HEP that distinguishes it from other fields is the availability of a high-fidelity simulation (thanks to MCNet collaborators!)

These simulations are usually expensive and non-differentiable.  A variety of ML methods can scaffold on top of our simulators to allow us to use all their physics to extract the most information from our data.

Theory of everything

Nature

Fast simulation

Parameter estimation / unfolding

**Physics simulators**

Experiment

Online processing & quality control

Detector-level observables

Detector-level observables

Pattern recognition ⟷ Pattern recogn

Data curation

Anomaly detection

calibration
clustering
tracking
noise mitigation
particle identification
...

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

This is only possible because of **factorization** (*Markov Property*): given the physics at one energy (~1/length) scale, the physics at the next one is independent of what came before.

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

We begin with a model and ME generators.

$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu}$$
$$+ i\bar{\psi} \slashed{D} \psi$$
$$+ \psi_i y_{ij} \psi_j \phi + \mathrm{h.c.}$$
$$+ |D_\mu \phi|^2 - V(\phi)$$
$$+ ???$$

*A lot of interesting work on efficient phase space generation with ML - see the living review for links*



```
**************************************************************
*                                                            *
*                        W E L C O M E  to                   *
*              M A D G R A P H 5 _ a M C @ N L O             *
*                                                            *
*                                                            *
*                 *                          *               *
*                  *        * *             *                *
*                    * * * * 5 * * * *                       *
*                  *        * *             *                *
*                 *                          *               *
*                                                            *
**************************************************************
```

## Standard is automated
## NLO or LO + matched

*For many cases, this is slow but not limiting (yet)*

# Part II: Fragmentation

Fragmentation uses MCMC;
standard is leading-log.

Not a limiting factor in
terms of computing time.

# Part III: Material Interactions

State-of-the-art for material interactions is Geant4.

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 competing resources!*

Geant 4

# Part IV: Digitization

It is important to mention that **after** Geant4, each experiment has custom code for *digitization*

this can also be slow; but is usually faster than G4 and reconstruction

deposited charge

MIP

Preamplifier output

40 MHz clock

analog threshold

Time

ToT = 2

It is important to mention that **after** Geant4, each experiment has custom code for *digitization*

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

deposited charge

Preamplifier output

analog threshold

Time

ToT = 2

# Factorization

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

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**

# Why should **you** care?
## N.B. ALL jet substructure analyses in ATLAS are forced to use full simulation as current fast sim. is not good enough.



**Standard Model Production Cross Section Measurements**

*Status: July 2017*

Why should **you** care?

N.B. ALL jet substructure analyses in ATLAS are forced to use full simulation as current fast sim. is not good enough.

Standard Model Production Cross Section Measurements

*Status: July 2017*

If we don't do something, the HL-LHC won't be possible. If we do something now, we can save O($10 million/year).

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



Our structure: **calorimeter images**

Grayscale images:
Pixel intensity =
energy deposited

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

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

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

[I. Goodfellow et al., NIPS 2014]

[L. de Oliveira, M. Paganini, BPN, PRL 120 (2018) 042003]

One image per calo. layer

One network per particle type; input particle energy



INPUTS

OUTPUTS

particle energy *E*

1

rescale

Scalar multiplication

latent space *z*

1024

Single Image Gen.

Single Image Gen.

Single Image Gen.

Resize

Linear Combination

NN to learn coefficients

**Generator network**

Resize

Linear Combination

NN to learn coefficients

[L. de Oliveira, M. Paganini, BPN, CSBS 1 (2017) 4]

use layer i as input to layer i+1

ReLU to encourage sparsity

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

help avoid 'mode collapse'



OUTPUTS

Concatenation

fake vs. real

Minibatch Discrimination

Absolute Difference

$\Delta > \epsilon$ ?

$\hat{E}_{tot}$

reco. energy

features

$\hat{E}_0$
$\hat{E}_1$
$\hat{E}_2$

Single Image Network

Single Image Network

Single Image Network

INPUTS

particle energy

$E$

Encourage energy conservation

**Discriminator network**

## Full physics generator (Geant4)



**CaloGAN**

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

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.

Nearest GEANT neighbour to each GAN image

Nearest GAN neighbour to each GEANT image

not memorizing

GAN nearest neighbours

GEANT nearest neighbours

no "mode collapse"

| Generation Method | Hardware | Batch Size | milliseconds/shower |
|---|---|---|---|
| GEANT4 | CPU | N/A | 1772 ⟵ |
| CALOGAN | CPU *Intel Xeon E5-2670* | 1 | 13.1 |
| | | 10 | 5.11 |
| | | 128 | 2.19 |
| | | 1024 | 2.03 |
| | GPU *NVIDIA K80* | 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!)

*If the physics changes, no expectation that the GAN will be accurate.*

Beyond our training sample!

Generative models are becoming more powerful & popular
(not just GANs, but other models like Variational Autoencoders and Normalizing Flows)

Our applications are often more challenging than industry because our data are less "structured" than natural images and we also have a strong requirement of quantitive and not just qualitative quality (e.g. jets versus celebrities)

…you will hear more about GANs in HEP tomorrow!

Theory of everything

**Fast simulation**

**Parameter estimation / unfolding**

**Physics simulators**

Detector-level observables

Pattern recognition

**Nature**

Experiment

**Online processing & quality control**

Detector-level observables

Pattern recogn

**Data curation**

**Anomaly detection**

calibration
clustering
tracking
noise mitigation
particle identification

...

"But what are the uncertainties on the NN"?

- question asked by every reviewer

"But what are the uncertainties on the NN"?

- question asked by every reviewer

**Let's consider this question in the context of a search for new particles in collision events.**

*this is representative for many analyses at the LHC, for example*

1. Train a classifier (in sim.) for signal vs. background.

2. Define a control region (**CR**) and a signal region (**SR**) using (1).

3. Check / modify simulation in CR.

4. Compare data and simulation in SR.

```
Significantly
different? go to
Stockholm : publish
limits.
```



ATLAS Collaboration, 2004.01678

**ATLAS**
$\sqrt{s}$=13 TeV, 139 fb$^{-1}$
$\sigma(H)=\sigma_{SM}(H)$
$B(H\rightarrow Za)=100\%$

• Data
— Background
-- H→Za (0.5 GeV)
⋯ H→Za (1.5 GeV)
-·- H→Za (2.5 GeV)

CR  SR

Events / 0.01

Data / Bkgd

Bkgd MC Stat

NN output

Precision / Optimality

*Bad use of our data, time, money, etc. but **not wrong**.*

Accuracy / Bias

Precision / Optimality: $\mathrm{NN(x)} \neq \frac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$

Optimal by Neyman-Pearson

Accuracy / Bias

*Note that this is not p(x|S) / p(x|B), however the two are monotonically related to each other.*

Precision / Optimality: $\mathrm{NN(x)} \neq \dfrac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$

Accuracy / Bias: $p_{\mathrm{prediction}}(\mathrm{NN}) \neq p_{\mathrm{true}}(\mathrm{NN})$

*The distribution of the (corrected) sim. is not correct.*

Precision / Optimality: $\mathrm{NN(x)} \neq \dfrac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$

$$p_{\mathrm{train}}(x) \neq p_{\mathrm{true}}(x)$$

*inaccurate training data*

*limited training statistics*

$$\mathrm{NN(x)}\big|_{p_{\mathrm{true}}=p_{\mathrm{train}}} \neq \dfrac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$$

*model/optimization flexibility*

**Statistical uncertainty**

**Systematic uncertainty**

*~ aleatoric*

*~ epistemic*

*limited prediction statistics*

$$p_{\mathrm{prediction}}(x) \neq p_{\mathrm{true}}(x)$$

*inaccurate prediction data*

Accuracy / Bias: $p_{\mathrm{prediction}}(\mathrm{NN}) \neq p_{\mathrm{true}}(\mathrm{NN})$

Precision / Optimality: $\mathrm{NN(x)} \neq \frac{p_{\text{true}}(x|\text{S+B})}{p_{\text{true}}(x|\text{B})}$

$p_{\text{train}}(x) \neq p_{\text{true}}(x)$
*inaccurate training data*

*limited training statistics*

$\mathrm{NN(x)}|_{p_{\text{true}}=p_{\text{train}}} \neq \frac{p_{\text{true}}(x|\text{S+B})}{p_{\text{true}}(x|\text{B})}$
*model/optimization flexibility*

**Statistical uncertainty** | **Systematic uncertainty**

*limited prediction statistics*

$p_{\text{prediction}}(x) \neq p_{\text{true}}(x)$
*inaccurate prediction data*

Accuracy / Bias: $p_{\text{prediction}}(\mathrm{NN}) \neq p_{\text{true}}(\mathrm{NN})$

You can always accomplish this by bootstrapping: making pseudo-datasets from resampling and then retraining.

It is important to fix the NN initialization so that you are not also testing your sensitivity to that.

This can be painful because it requires retraining many NNs.

Maybe can accomplish with one Bayesian NN? See e.g. S. Bollweg, et al., SciPost Phys. 8, 006 (2020), 1904.10004 for a particle physics example.

Precision / Optimality: $\mathrm{NN(x)} \neq \frac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$

$p_{\mathrm{train}}(x) \neq p_{\mathrm{true}}(x)$

*inaccurate training data*

*limited training statistics*

$\mathrm{NN(x)}|_{p_{\mathrm{true}}=p_{\mathrm{train}}} \neq \frac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$

*model/optimization flexibility*

Statistical uncertainty      Systematic uncertainty

$p_{\mathrm{prediction}}(x) \neq p_{\mathrm{true}}(x)$

*limited prediction statistics*

*inaccurate prediction data*

Accuracy / Bias: $p_{\mathrm{prediction}}(\mathrm{NN}) \neq p_{\mathrm{true}}(\mathrm{NN})$

As with all systematic uncertainties, this is hard to quantify.

One component is due to the modeling of p(x) - more on this later.

Testing the flexibility of the network requires checking the sensitivity to the architecture (#layers, nodes/layer, etc.), the initialization, the training procedure (#epochs, learning rate, etc.)

Precision / Optimality: $\text{NN(x)} \neq \frac{p_{\text{true}}(x|\text{S+B})}{p_{\text{true}}(x|\text{B})}$

limited training statistics

$p_{\text{train}}(x) \neq p_{\text{true}}(x)$

*inaccurate training data*

$\text{NN(x)}|_{p_{\text{true}}=p_{\text{train}}} \neq \frac{p_{\text{true}}(x|\text{S+B})}{p_{\text{true}}(x|\text{B})}$

*model/optimization flexibility*

Statistical uncertainty    Systematic uncertainty

limited prediction statistics

$p_{\text{prediction}}(x) \neq p_{\text{true}}(x)$

*inaccurate prediction data*

Accuracy / Bias: $p_{\text{prediction}}(\text{NN}) \neq p_{\text{true}}(\text{NN})$

Can be estimated via bootstrapping.  Less painful here because the NN's are fixed.

N.B. it may be possible to design a network that is designed to minimize uncertainty at inference.  This does not work in all cases, but early studies in particle physics seem promising: S. Wunsch et al., 2003.07186, P. da Castro et al., CPC 244 (2019) 170, 1806.04743

Precision / Optimality: $\mathrm{NN(x)} \neq \frac{p_{\text{true}}(x|\text{S+B})}{p_{\text{true}}(x|\text{B})}$

$p_{\text{train}}(x) \neq p_{\text{true}}(x)$

*inaccurate training data*

*limited training statistics*

$\mathrm{NN(x)}|_{p_{\text{true}}=p_{\text{train}}} \neq \frac{p_{\text{true}}(x|\text{S+B})}{p_{\text{true}}(x|\text{B})}$

*model/optimization flexibility*

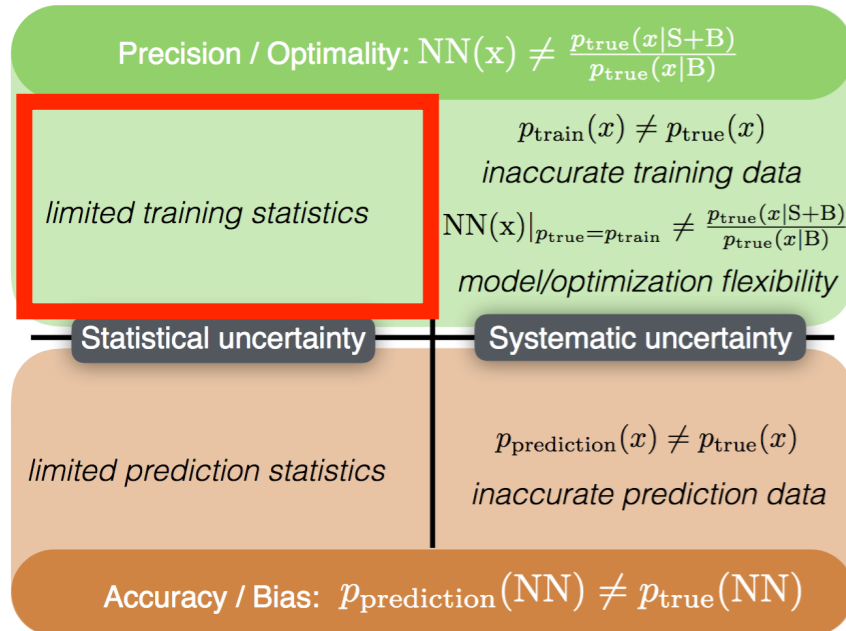Statistical uncertainty — Systematic uncertainty

*limited prediction statistics*

$p_{\text{prediction}}(x) \neq p_{\text{true}}(x)$

*inaccurate prediction data*

Accuracy / Bias: $p_{\text{prediction}}(\mathrm{NN}) \neq p_{\text{true}}(\mathrm{NN})$

**This is the trickiest one…**

*…because we need the uncertainty on the modeling of x and x can be high-dimensional!*

In many cases, the uncertainties factorize, e.g. the uncertainty on two photon energies can be decomposed into the uncertainty on each photon.

However, in many cases, we simply do not know the full uncertainty model (= nuisance parameters and their distribution)

One word of caution: current paradigm for uncertainties may be too naive for high-dimensional analysis!

(truly end-to-end)

e.g. for some uncertainties, we often compare two different models - one nuisance parameter.

How can we even see how sensitive we are to high-dimensional effects?

One word of caution: current paradigm for uncertainties may be too naive for high-dimensional analysis!

(truly end-to-end)

e.g. for some uncertainties, we often compare two different models - one nuisance parameter.

How can we even see how sensitive we are to high-dimensional effects?

**Answer: borrow tools from AI Safety**

There is a vast literature on how easy it is to "attack" a NN.

*They want to know: how subtle can an attack be and still significantly impact the output.*

We know (hope?!) that nature is not evil, but these tools can help us probe the high-dimensional sensitivity of our NNs.

Model Physical Dynamics by Sampling from Distribution

STOP STOP STOP STOP STOP STOP

$f_\theta(x)$ → Output → SPEED LIMIT 45 Target

Stationary + Drive-By Testing

Input → RP₂ → Mask → Perturbed Stop Sign Under Varying Distances/Angles

K. Eykholt et. al, 1707.08945

$J$ = collision event (in all of its high-dimensional glory)

$f$ = fixed classifier for signal vs. background

Loss

$$\mathcal{L}_{\mathrm{sig}} = \log(1 - f(g(J))),$$
$$\mathcal{L}_{\mathrm{bg}} = \lambda_{\mathrm{cls}}(f(J) - f(g(J)))^2$$
$$+ \sum_i \lambda_{\mathrm{obs}}^{(i)}(\mathcal{O}^{(i)}(J) - \mathcal{O}^{(i)}(g(J))^2$$

$g$ is a learned NN that maps $J$ to $J + \delta J$.

$O(J)$ are observables that will be validated in the CR.

"worst-case uncertainty"

Precision / Optimality: $\mathrm{NN(x)} \neq \frac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$

$p_{\mathrm{train}}(x) \neq p_{\mathrm{true}}(x)$

*inaccurate training data*

*limited training statistics*

$\mathrm{NN(x)}|_{p_{\mathrm{true}}=p_{\mathrm{train}}} \neq \frac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$

*model/optimization flexibility*

Statistical uncertainty — Systematic uncertainty

$p_{\mathrm{prediction}}(x) \neq p_{\mathrm{true}}(x)$

*limited prediction statistics*

*inaccurate prediction data*

Accuracy / Bias: $p_{\mathrm{prediction}}(\mathrm{NN}) \neq p_{\mathrm{true}}(\mathrm{NN})$

Train with more events!

Precision / Optimality: $\text{NN}(\text{x}) \neq \frac{p_\text{true}(x|\text{S+B})}{p_\text{true}(x|\text{B})}$

$p_\text{train}(x) \neq p_\text{true}(x)$

*inaccurate training data*

*limited training statistics*

$\text{NN}(\text{x})|_{p_\text{true}=p_\text{train}} \neq \frac{p_\text{true}(x|\text{S+B})}{p_\text{true}(x|\text{B})}$

*model/optimization flexibility*

Statistical uncertainty | Systematic uncertainty

$p_\text{prediction}(x) \neq p_\text{true}(x)$

*limited prediction statistics*

*inaccurate prediction data*

Accuracy / Bias: $p_\text{prediction}(\text{NN}) \neq p_\text{true}(\text{NN})$

Train with more events!

…maybe use NN's to help with that

Might be possible to reduce uncertainties or at least alleviate analysis complexity by making your NN independent of known nuisance parameters*.

…might also be better to explicitly depend on the nuisance parameters and profile them in data.
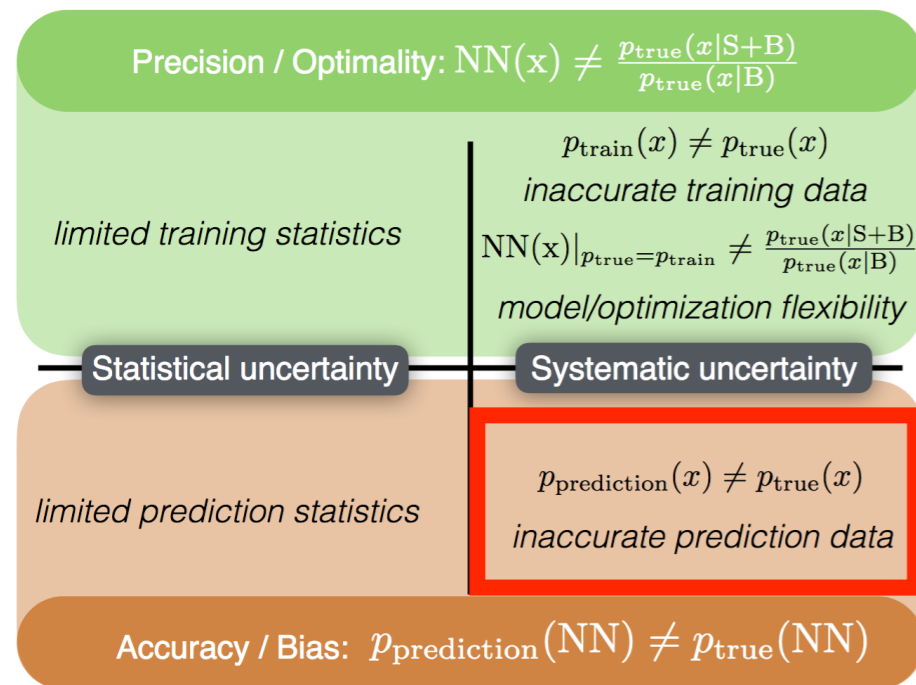
Work hard to understand the true nuisance parameters in the hypervariate parameter space.

Precision / Optimality: $\mathrm{NN(x)} \neq \frac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$

*limited training statistics*

$p_{\mathrm{train}}(x) \neq p_{\mathrm{true}}(x)$
*inaccurate training data*

$\mathrm{NN(x)}|_{p_{\mathrm{true}}=p_{\mathrm{train}}} \neq \frac{p_{\mathrm{true}}(x|\mathrm{S+B})}{p_{\mathrm{true}}(x|\mathrm{B})}$

*model/optimization flexibility*

Statistical uncertainty — Systematic uncertainty

*limited prediction statistics*

$p_{\mathrm{prediction}}(x) \neq p_{\mathrm{true}}(x)$
*inaccurate prediction data*

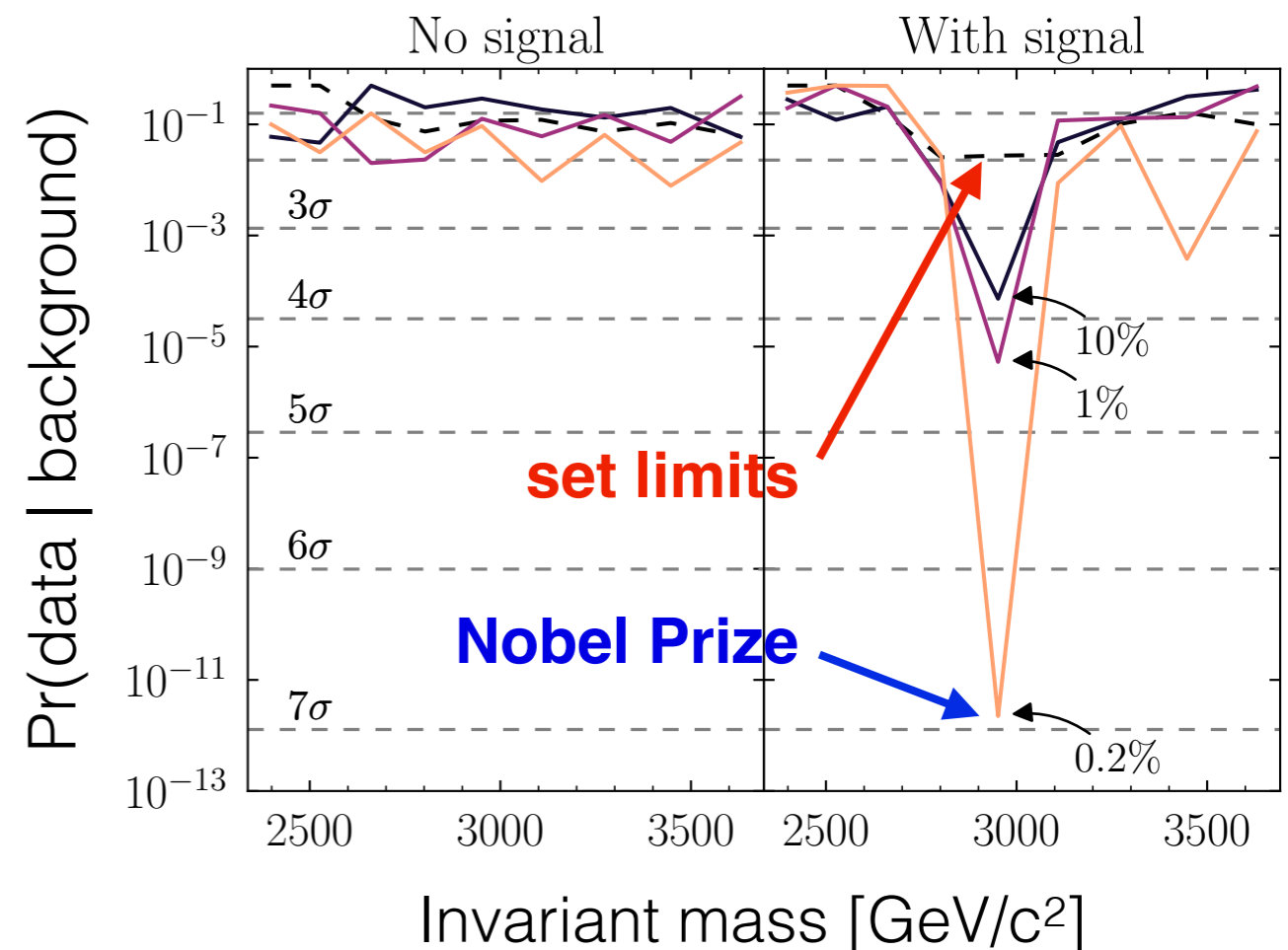Accuracy / Bias: $p_{\mathrm{prediction}}(\mathrm{NN}) \neq p_{\mathrm{true}}(\mathrm{NN})$

In my opinion, this is **THE** biggest challenge with deploying NN-based analyses … solving it will require hard physics work.

Work hard to understand the true nuisance parameters in the hypervariate parameter space.

Don't use simulation!
(not always possible and of course, still has assumptions…)



No signal     With signal

Pr(data | background)

$10^{-1}$

$3\sigma$

$10^{-3}$

$4\sigma$

$10^{-5}$

$5\sigma$

$10^{-7}$

$6\sigma$

$10^{-9}$

$10^{-11}$

$7\sigma$

$10^{-13}$

2500   3000   3500    2500   3000   3500

set limits

Nobel Prize

10%

1%

0.2%

Invariant mass [GeV/c²]

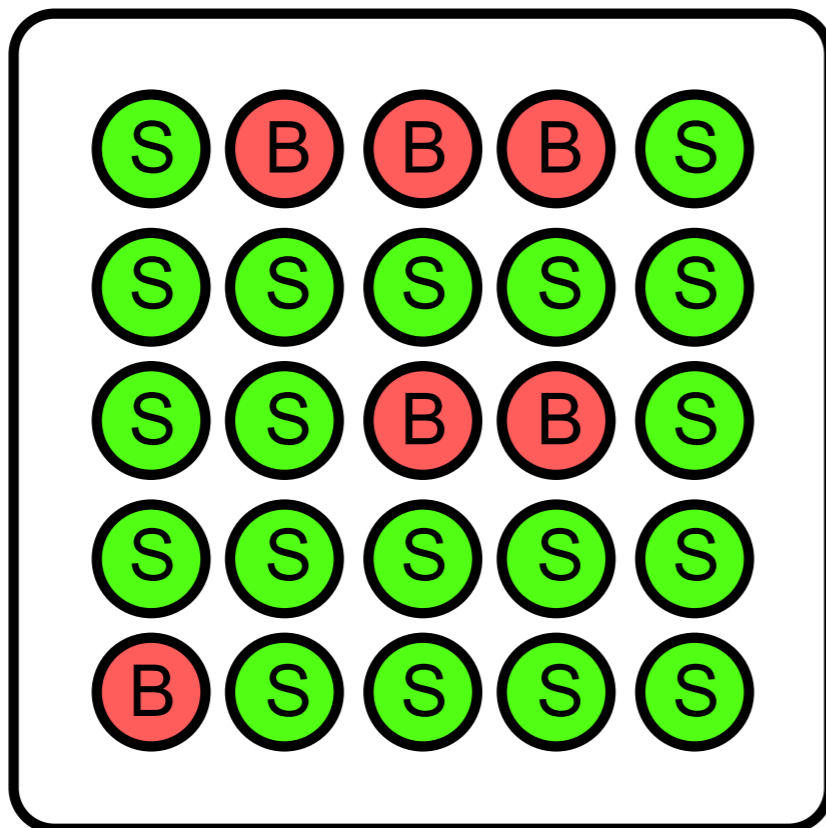Why can't I just pay some physicists to label events and then train a neural network using those labels?



Image credit: pixabay.com

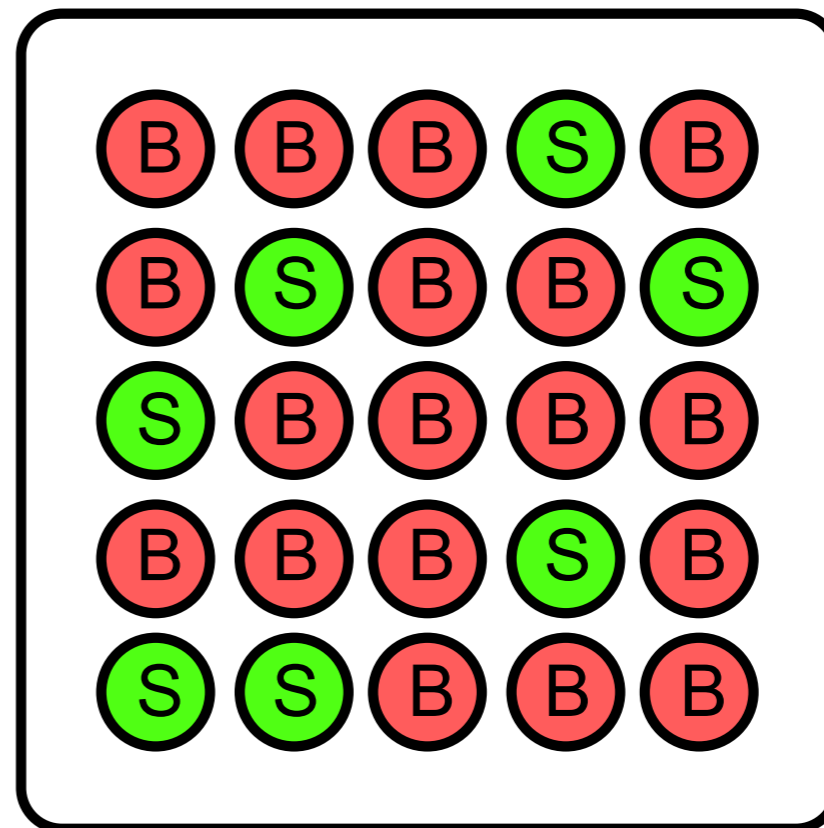Answer: this is not cats-versus-dogs … thanks to quantum mechanics it is **not possible to know** what happened.

The data are unlabeled and in the best case, come to us as mixtures of two classes ("signal" and "background").
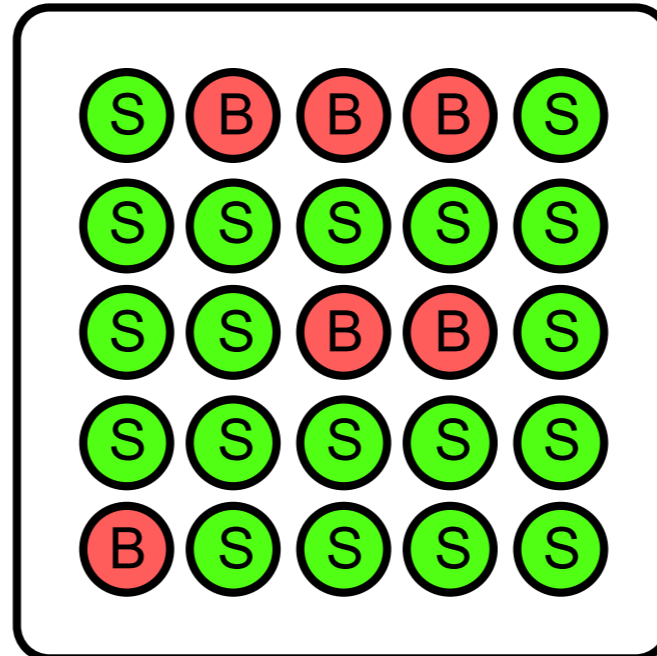


Mixed Sample 1

Mixed Sample 2

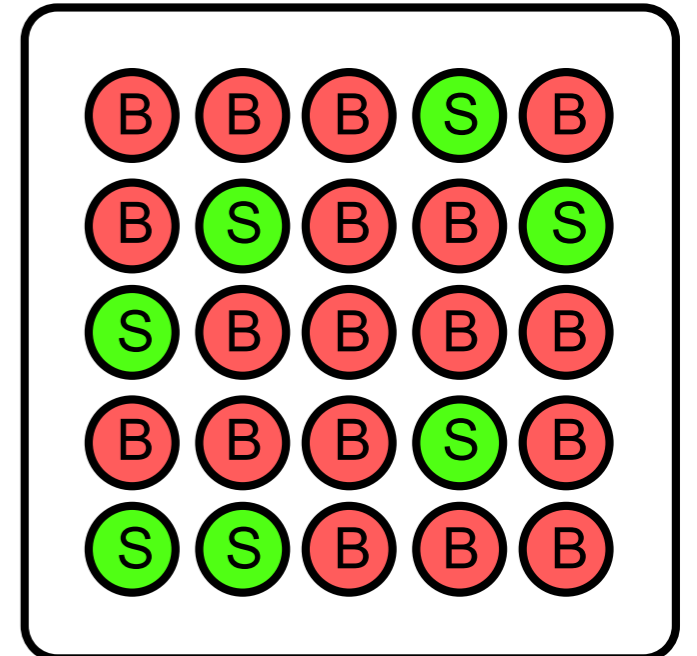(we don't get to observe the color of the circles)

# Weak supervision:
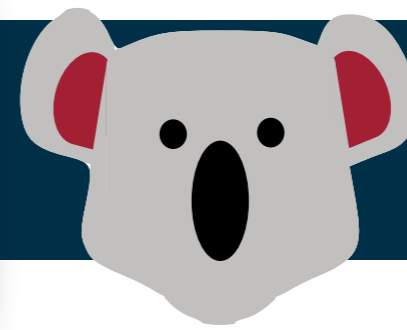## *Classification Without Labels*

Can we learn without any label information?



Mixed Sample 1

Mixed Sample 2

Can we learn without any label information?

**Yes !**

*Training on impure samples is (asymptotically) equivalent to training on pure samples*



Mixed Sample 1

Mixed Sample 2

0

1

Classifier

E. Metodiev, **BPN**, J. Thaler, JHEP 10 (2017) 51
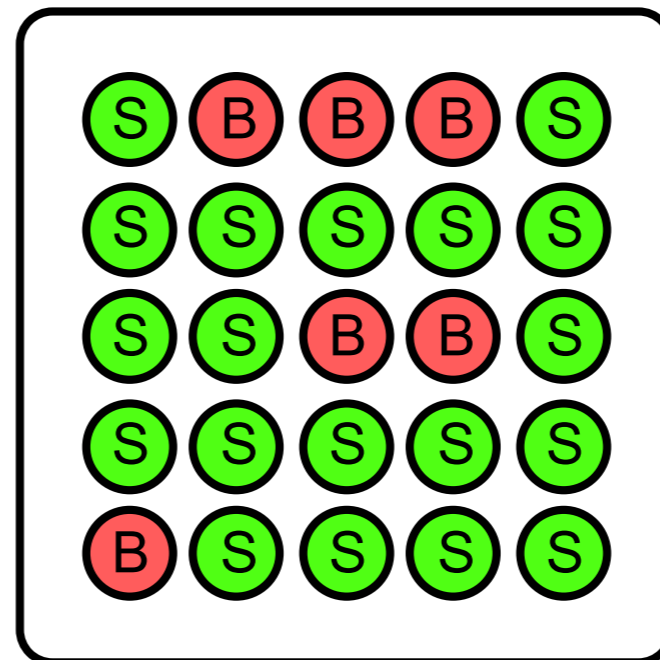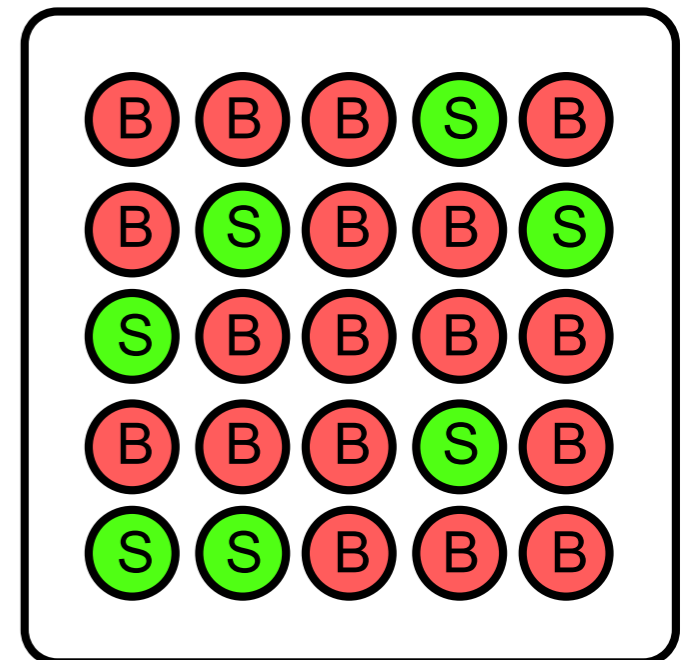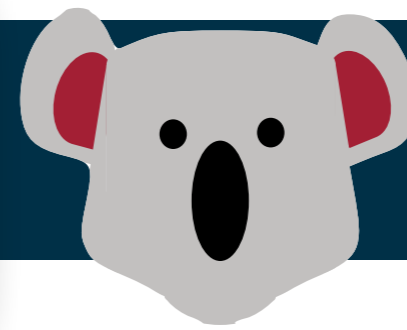
# Weak supervision:
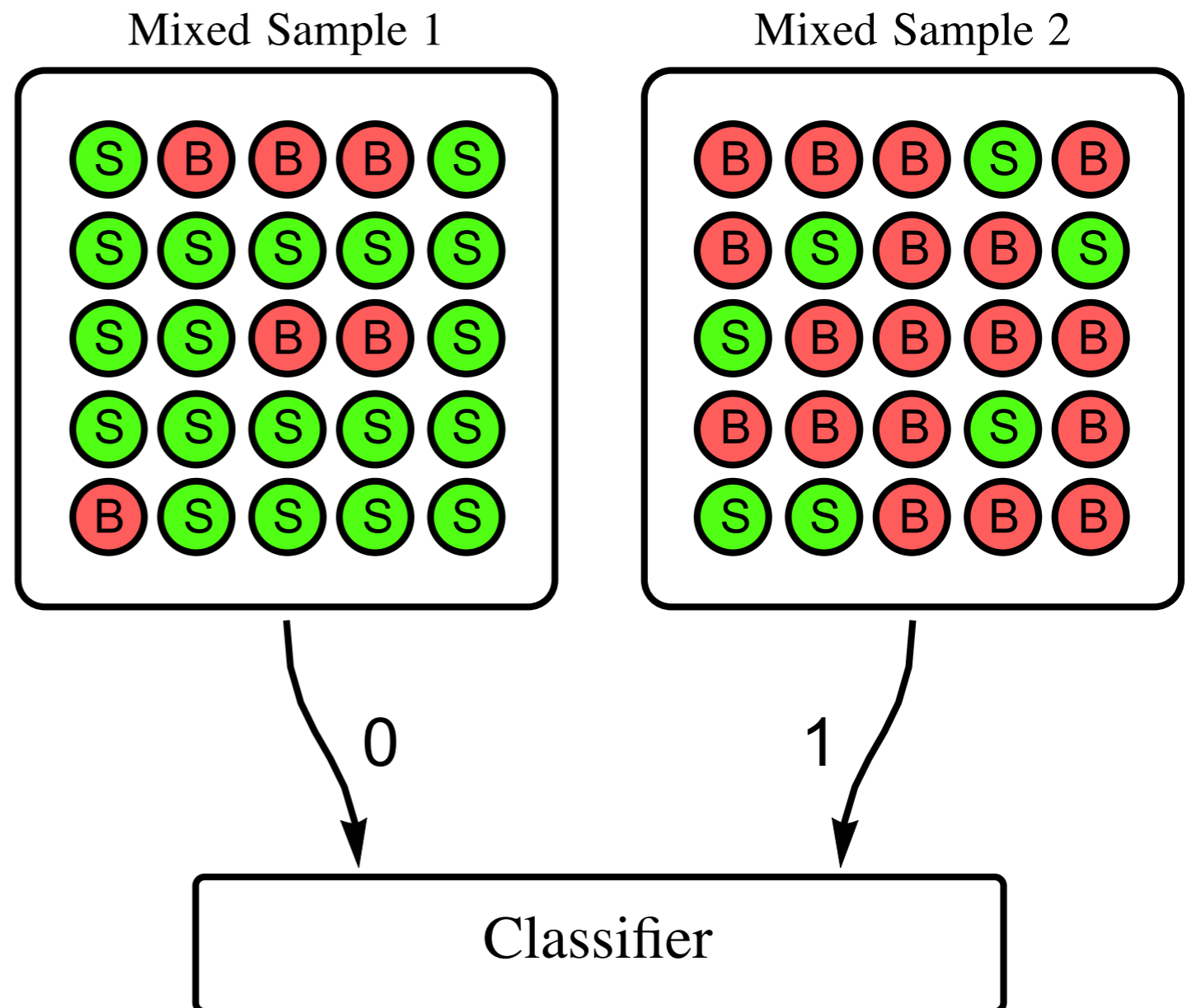## *Classification Without Labels*

Can we learn without any label information?

**Yes !**

*Training on impure samples is (asymptotically) equivalent to training on pure samples*

**Exercise: What does this mean and can you prove it?**



Mixed Sample 1

Mixed Sample 2

0

1

Classifier

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

J. Collins, K. Howe, BPN,
Phys. Rev. Lett. 121 (2018)
241803, 1805.02664

Mixed Sample 2

Mixed Sample 1

$dN/dm_{res}$

background

signal

$m_{res}$

hypervariate
feature space

+ be careful to not pay a big trials factor
(ask if interested)

Example: two "jet" search

Jet 1

p

Features: radiation
pattern inside each jet

p

ATLAS
EXPERIMENT

Run: 302347
Event: 753275626
2016-06-18 18:41:48 CEST

Jet 2

LHC simulation

Events / 100 GeV

$m_{JJ}$ [GeV/c²]

sidebands



Events / 100 GeV

$10^5$

$10^4$

$10^3$

$10^2$

$10^1$

$10^0$

2000    3000    4000

$m_{JJ}$ [GeV/$c^2$]

LHC simulation

sidebands

standard parametric
fit to background



100%

LHC simulation

Events / 100 GeV

$10^5$
$10^4$
$10^3$
$10^2$
$10^1$
$10^0$

2000    3000    4000

$m_{JJ}$ [GeV/$c^2$]

sidebands

standard parametric
fit to background



LHC simulation

100%

10%

1%

0.2%

0.02%

Events / 100 GeV

$10^5$
$10^4$
$10^3$
$10^2$
$10^1$
$10^0$

2000    3000    4000

$m_{JJ}$ [GeV/$c^2$]

- - - -  no cut on NN          ———  most 1% signal-region-like

———  most 10% signal-region-like     ———  most 0.2% signal-region-like

sidebands

standard parametric
fit to background

No signal

Events / 100 GeV

$10^5$

$10^4$

$10^3$

$10^2$

$10^1$

$10^0$

**100%**

**10%**

**1%**

**0.2%**

**0.02%**

LHC simulation

2000          3000          4000

$m_{JJ}$ [GeV/c$^2$]

units of the
proton mass

$-\Pr(\text{data} \,|\, \text{background})$

$10^0$

$10^{-2}$

$10^{-4}$

$10^{-6}$

$10^{-8}$

$10^{-10}$

$10^{-12}$

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

2500          3000          3500

$m_{JJ}$ [GeV/c$^2$]

- - - -    no cut on NN                              most 1% signal-region-like

———    most 10% signal-region-like          most 0.2% signal-region-like

LHC simulation

**100%**

**10%**

1%

0.2%

0.02%

Events / 100 GeV

$m_{JJ}$ [GeV/c$^2$]

No signal

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

$-$Pr(data | background)

$m_{JJ}$ [GeV/c$^2$]

- - - no cut on NN

—— most 10% signal-region-like

—— most 1% signal-region-like

—— most 0.2% signal-region-like

LHC simulation

Events / 100 GeV

$10^5$

$10^4$

$10^3$

$10^2$

$10^1$

$10^0$

**100%**

**10%**

**1%**

**0.2%**

**0.02%**

2000

3000

4000

$m_{JJ}$ [GeV/c$^2$]

No signal

$-$ Pr(data | background)

$10^0$

$10^{-2}$

$10^{-4}$

$10^{-6}$

$10^{-8}$

$10^{-10}$

$10^{-12}$

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

2500

3000

3500

$m_{JJ}$ [GeV/c$^2$]

- - - no cut on NN

—— most 10% signal-region-like

—— most 1% signal-region-like

—— most 0.2% signal-region-like

LHC simulation

100%

10%

1%

0.2%

0.02%

$m_{JJ}$ [GeV/c$^2$]

No signal

$10^0$

$10^{-2}$

$3\sigma$

$10^{-4}$

$4\sigma$

$10^{-6}$

$5\sigma$

$10^{-8}$

$6\sigma$

$10^{-10}$

$10^{-12}$

$7\sigma$

Pr(data | background)

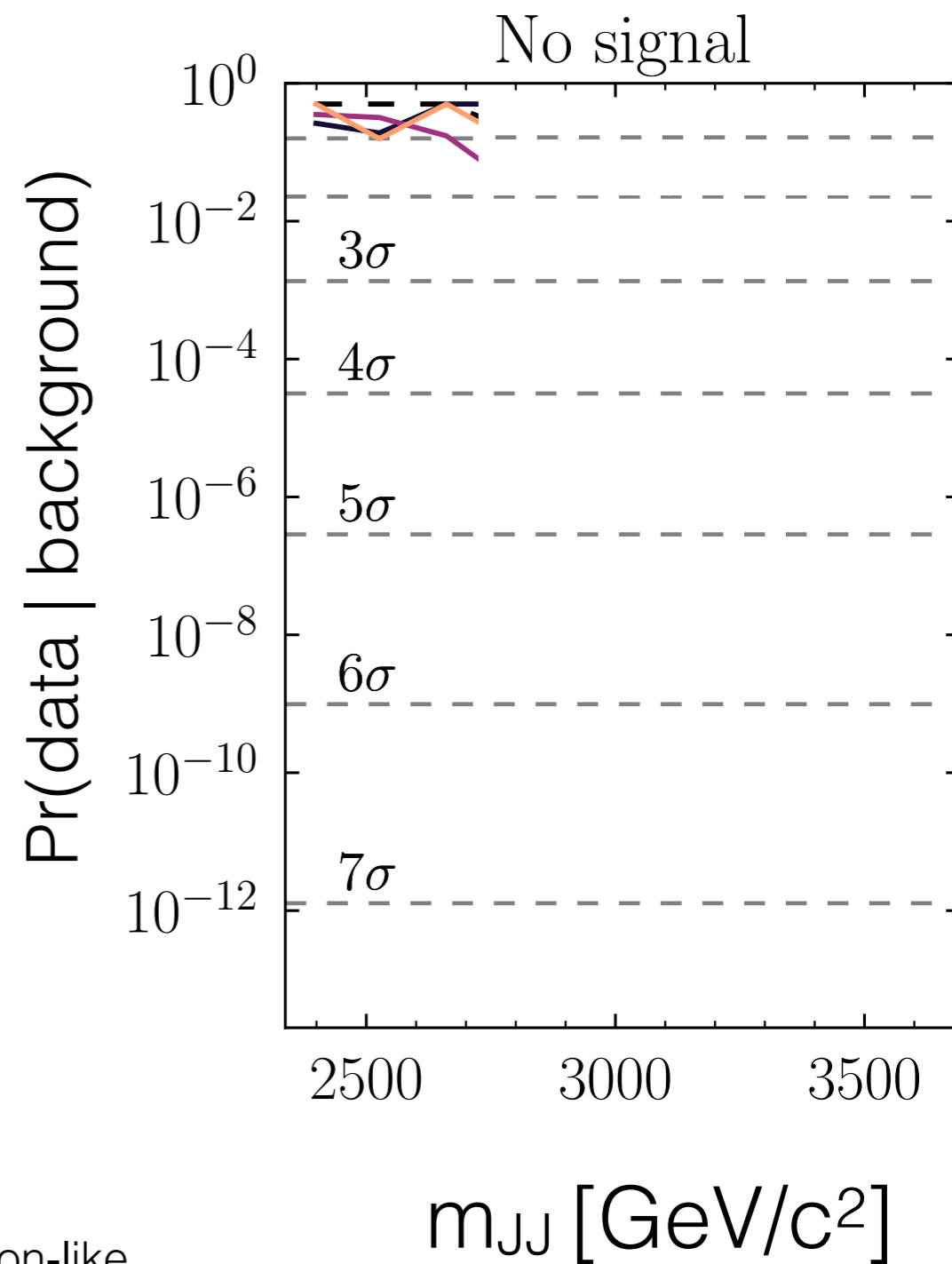2500    3000    3500

$m_{JJ}$ [GeV/c$^2$]

- - - - no cut on NN

——— most 10% signal-region-like

——— most 1% signal-region-like

——— most 0.2% signal-region-like

LHC simulation

100%

10%

1%

0.2%

0.02%

$m_{JJ}$ [GeV/c²]

1000

No signal

$\Pr(\text{data} \mid \text{background})$

$10^0$

$10^{-2}$

$3\sigma$

$10^{-4}$

$4\sigma$

$10^{-6}$

$5\sigma$

$10^{-8}$

$6\sigma$

$10^{-10}$

$7\sigma$

$10^{-12}$

2500    3000    3500

$m_{JJ}$ [GeV/c²]

- - - - no cut on NN      ——— most 1% signal-region-like

——— most 10% signal-region-like      ——— most 0.2% signal-region-like

LHC simulation

**100%**

**10%**

**1%**

**0.2%**

**0.02%**

$m_{JJ}$ [GeV/c$^2$]

1000

No signal

$10^0$

$10^{-2}$

$3\sigma$

$10^{-4}$

$4\sigma$

$10^{-6}$

$5\sigma$

$10^{-8}$

$6\sigma$

$10^{-10}$

$7\sigma$

$10^{-12}$

Pr(data | background)

2500    3000    3500

$m_{JJ}$ [GeV/c$^2$]
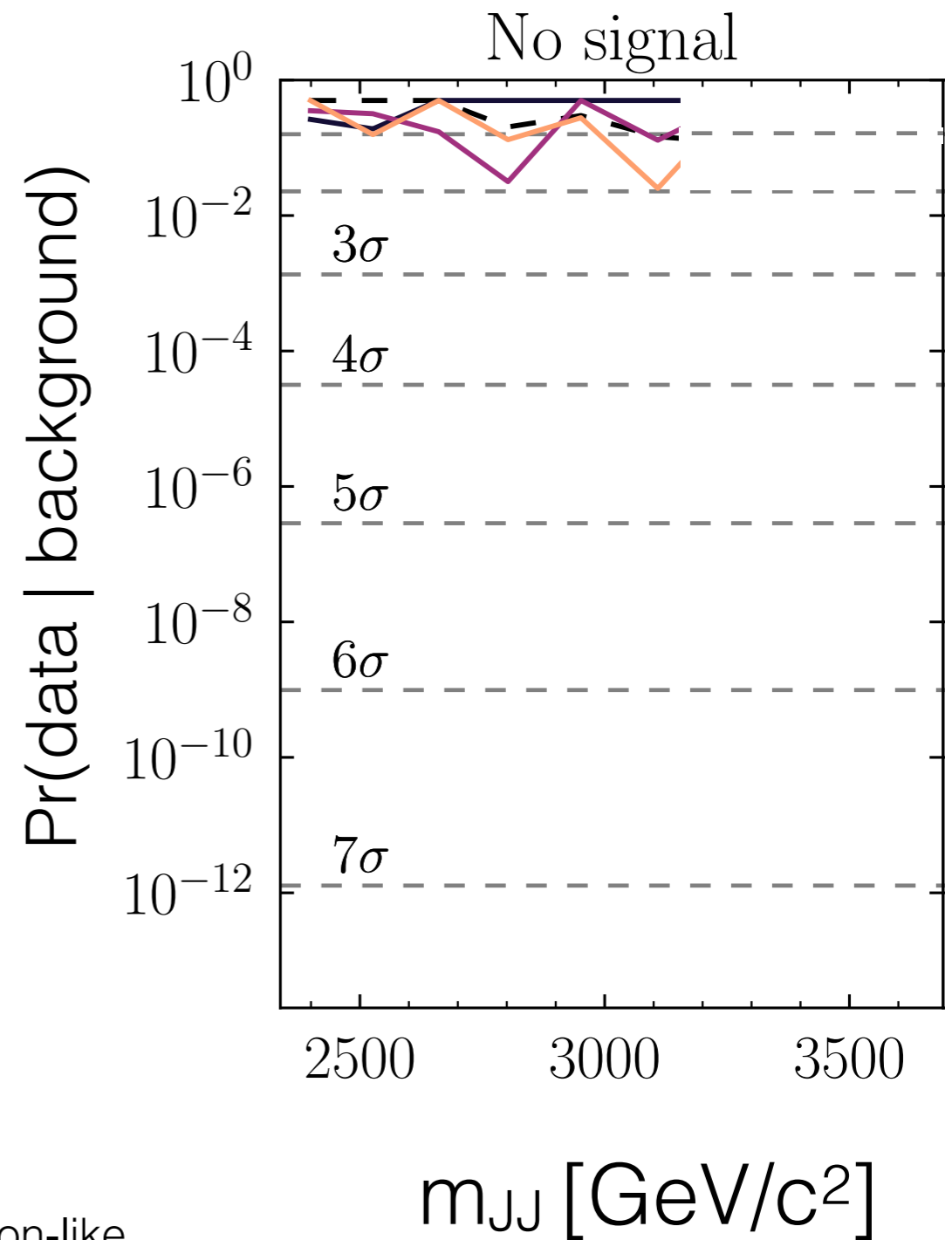
- - - no cut on NN

——— most 1% signal-region-like

——— most 10% signal-region-like

——— most 0.2% signal-region-like

LHC simulation

100%

10%

1%

0.2%

0.02%

$m_{JJ}$ [GeV/c$^2$]

1000

No signal

$10^0$

$10^{-2}$ — $3\sigma$

$10^{-4}$ — $4\sigma$

$10^{-6}$ — $5\sigma$

$10^{-8}$ — $6\sigma$

$10^{-10}$

$10^{-12}$ — $7\sigma$

Pr(data | background)

2500     3000     3500

$m_{JJ}$ [GeV/c$^2$]

- - - -  no cut on NN

——  most 10% signal-region-like

——  most 1% signal-region-like

——  most 0.2% signal-region-like

LHC simulation

**100%**

**10%**

**1%**

**0.2%**

**0.02%**

$m_{JJ}$ [GeV/c²]

No signal

Pr(data | background)

$10^0$

$10^{-2}$

$10^{-4}$

$10^{-6}$

$10^{-8}$

$10^{-10}$

$10^{-12}$

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

2500    3000    3500

$m_{JJ}$ [GeV/c²]

- - - - no cut on NN                     ——— most 1% signal-region-like

——— most 10% signal-region-like     ——— most 0.2% signal-region-like

LHC simulation

100%

10%

1%

0.2%

0.02%

S / 100 GeV

$m_{JJ}$ [GeV/c$^2$]

1000

No signal

$10^0$

$10^{-2}$

$10^{-4}$

$10^{-6}$

$10^{-8}$

$10^{-10}$

$10^{-12}$

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

Pr(data | background)

2500    3000    3500

$m_{JJ}$ [GeV/c$^2$]

- - - -  no cut on NN

———  most 10% signal-region-like

———  most 1% signal-region-like

———  most 0.2% signal-region-like

sidebands

standard parametric
fit to background



LHC simulation

Events / 100 GeV

$10^5$

$10^4$

$10^3$

$10^2$

$10^1$

$10^0$

**100%**

**10%**

**1%**

*0.2%*

*0.02%*

*signal*

2000    3000    4000

$m_{JJ}$ [GeV/c$^2$]

W

W'

X

- - - -  no cut on NN          |          most 1% signal-regi

――――  most 10% signal-reg ... ke    |    ――――  most 0.2% signal-reg

sidebands

standard parametric
fit to background

LHC simulation

100%

$m_{JJ}$ [GeV/c$^2$]

signal-reg

signal-reg

$-$Pr(data | background)

$10^{-2}$

$10^{-4}$

$10^{-6}$

With signal

$3\sigma$

$4\sigma$

$5\sigma$

2500    3000    3500

$m_{JJ}$ [GeV/c$^2$]

signal

LHC simulation

100%

$m_{JJ}$ [GeV/c²]

signal-regio

signal-reg

With signal

$10^{-2}$

$3\sigma$

$10^{-4}$

$4\sigma$

$10^{-6}$

$5\sigma$

$10^{-8}$

$6\sigma$

$10^{-10}$

$10^{-12}$

$7\sigma$

Pr(data | background)

2500   3000   3500

$m_{JJ}$ [GeV/c²]

LHC simulation

*signal*

100%

m_JJ [GeV/c²]

signal-region-like

signal-region-like

With signal

Pr(data | background)

$10^{-2}$

$10^{-4}$

$10^{-6}$

$10^{-8}$

$10^{-10}$

$10^{-12}$

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

2500    3000    3500

m_JJ [GeV/c²]

LHC simulation

signal

100%

m_{JJ} [GeV/c²]

signal-regi

signal-reg

Pr(data | background)

With signal

$10^{-2}$

$3\sigma$

$10^{-4}$

$4\sigma$

$10^{-6}$

$5\sigma$

$10^{-8}$

$6\sigma$

$10^{-10}$

$10^{-12}$

$7\sigma$

2500    3000    3500

m_{JJ} [GeV/c²]

LHC simulation

100%

$m_{JJ}$ [GeV/c²]

signal-regio

signal-regi

With signal

$\Pr(\text{data} \mid \text{background})$

$10^{-2}$

$10^{-4}$

$10^{-6}$

$10^{-8}$

$10^{-10}$

$10^{-12}$

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

2500    3000    3500

$m_{JJ}$ [GeV/c²]

signal

LHC simulation

m_JJ [GeV/c²]

100%

Pr(data | background)

With signal

$10^{-2}$
$10^{-4}$
$10^{-6}$
$10^{-8}$
$10^{-10}$
$10^{-12}$

$3\sigma$
$4\sigma$
$5\sigma$
$6\sigma$
$7\sigma$

1%

0.2%

2500    3000    3500

m_JJ [GeV/c²]

signal-region-like

signal-region-like

LHC simulation

100%

signal

m_JJ [GeV/c²]

signal-region

signal-region

With signal

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

10%

1%

0.2%

$10^{-2}$

$10^{-4}$

$10^{-6}$

$10^{-8}$

$10^{-10}$

$10^{-12}$

Pr(data | background)

2500        3000        3500

$m_{JJ}$ [GeV/c²]

signal

100%

LHC simulation

signal-region

signal-region

m_JJ [GeV/c²]

Pr(data | background)

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

$10^{-2}$

$10^{-4}$

$10^{-6}$

$10^{-8}$

$10^{-10}$

$10^{-12}$

With signal

10%

1%

0.2%

2500     3000     3500

$m_{JJ}$ [GeV/c²]

LHC simulation

$10^5$

$10^4$

$10^3$

100%

10%

1%

0.2%

0.02%

$m_{JJ}$ [GeV/c$^2$]

000

signal-region-like

signal-region-like

With signal

$3\sigma$

$4\sigma$

$5\sigma$

$6\sigma$

$7\sigma$

$10^{-2}$

$10^{-4}$

$10^{-6}$

$10^{-8}$

$10^{-10}$

$10^{-12}$

Pr(data | background)

10%

1%

0.2%

2500        3000        3500

$m_{JJ}$ [GeV/c$^2$]

ATLAS Collaboration, 2005.02983
Analysis Team: A. Cukeriman, BPN



First round, keep it simple: feature space is 2D (jet masses)

ATLAS Collaboration, 2005.02983
Analysis Team: A. Cukeriman, BPN



"Jet" 1

B

A

C

"Jet" 2

ATLAS
$\sqrt{s}$ = 13 TeV, 139 fb$^{-1}$
No Injected Signal

ATLAS
$\sqrt{s}$ = 13 TeV, 139 fb$^{-1}$
✕  Injected Signal
$m_A$ = 3000 GeV
$m_B$ = 400 GeV
$m_C$ = 80 GeV

First round, keep it simple: feature space is 2D (jet masses)

← *Better*

95% CL Exclusion Limit σ(pp →A →BC) [fb]

Four model parameters: $m_A$, $m_B$, $m_C$, σ

We fix the masses and set limits on the cross section σ

"Jet" 1

B

A

C

"Jet" 2

$(m_B, m_C)$ [GeV]

(400,400)  (400,200)  (400,80)  (200,200)  (200,80)  (80,80)

**ATLAS** Collaboration, 2005.02983
Analysis Team: A. Cukeriman, BPN

**Better**

95% CL Exclusion Limit σ(pp →A →BC) [fb]

40
35
30
25
20
15
10
5
0

Four model parameters: $m_A$, $m_B$, $m_C$, σ

We fix the masses and set limits on the cross section σ

*ATLAS*

*This is an official result from the ATLAS Collaboration at CERN!*

B
A
C

(400,400)
(400,200)
(400,80)
(200,200)
(200,80)
(80,80)

$(m_B, m_C)$ [GeV]

**ATLAS** Collaboration, 2005.02983
Analysis Team: A. Cukeriman, BPN

Better

95% CL Exclusion Limit
$\sigma(pp \to A \to BC)$ [fb]

40

35

30

25

20

15

10

5

0

Four model
parameters: $m_A$,
$m_B$, $m_C$, $\sigma$

We fix the
masses and set
limits on the
cross section $\sigma$

**ATLAS**

$\sqrt{s}$ = 13 TeV, 139 fb$^{-1}$

*The data were
collected during the
last run in 2015-2018*

$(m_B, m_C)$ [GeV]

(400,400)
(400,200)
(400,80)
(200,200)
(200,80)
(80,80)

B

A

C

ATLAS Collaboration, 2005.02983
Analysis Team: A. Cukeriman, BPN

*Better* ← (blue arrow)

95% CL Exclusion Limit $\sigma(pp \to A \to BC)$ [fb]

ATLAS

$\sqrt{s}$ = 13 TeV, 139 fb$^{-1}$

$\in$ =0.1, $m_A$ =3000 GeV

*$m_A$ is fixed and we also fix a NN efficiency of 10%*

$(m_B, m_C)$ [GeV]

(400,400)  (400,200)  (400,80)  (200,200)  (200,80)  (80,80)

**ATLAS** Collaboration, 2005.02983
Analysis Team: A. Cukeriman, BPN

B    A    C

Fun fact: this plot required training 10k NNs

*Better*

95% CL Exclusion Limit $\sigma(pp \to A \to BC)$ [fb]

As we learn from data, for each signal cross section, need to retrain!

**ATLAS**

$\sqrt{s}$ = 13 TeV, 139 fb$^{-1}$

$\in$ =0.1, $m_A$ =3000 GeV

■ Observed

--- Expected

±1 σ

±2 σ

$(m_B, m_C)$ [GeV]

(400,400)  (400,200)  (400,80)  (200,200)  (200,80)  (80,80)

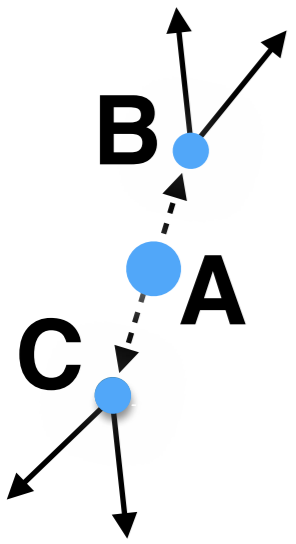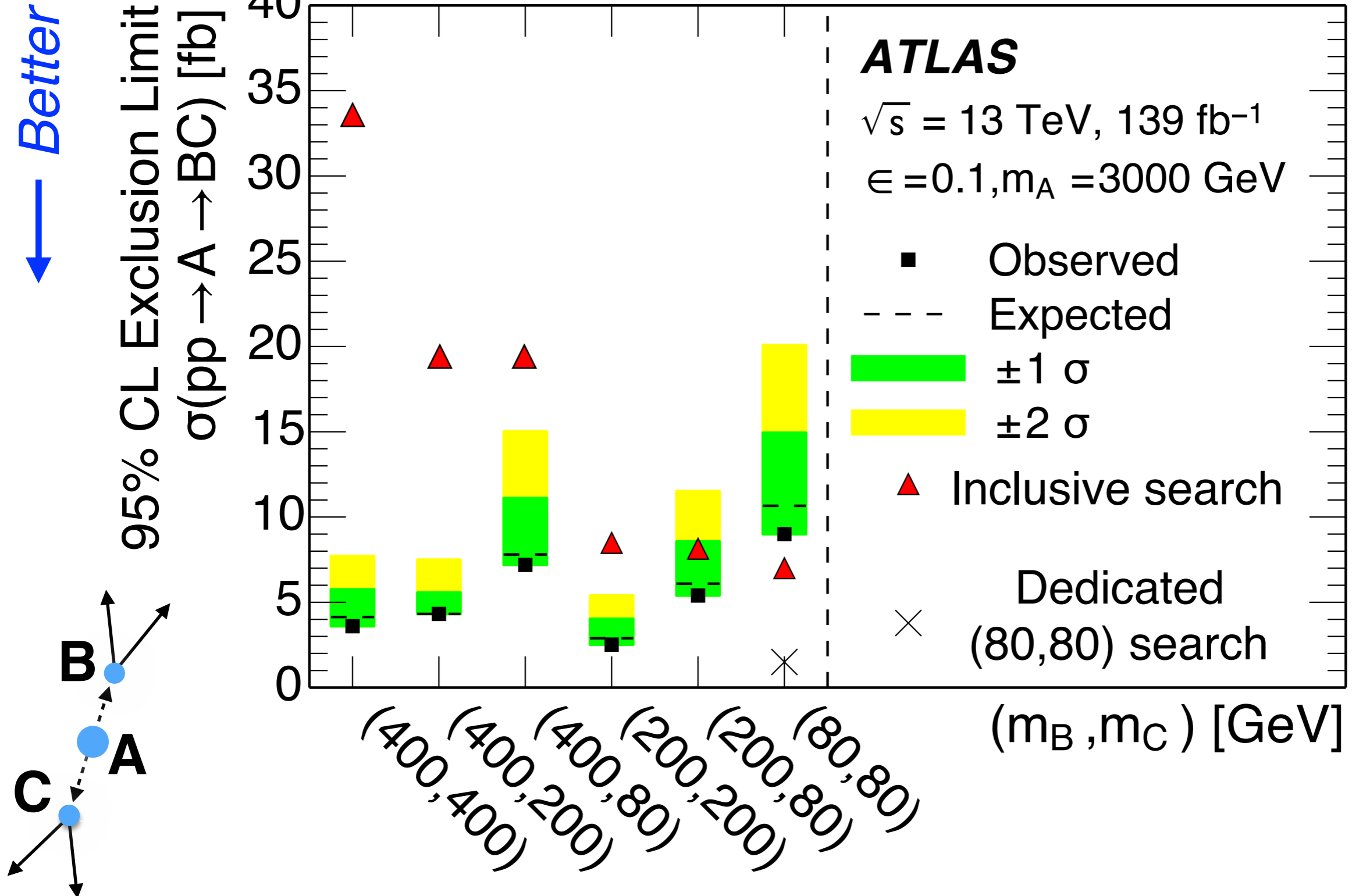ATLAS Collaboration, 2005.02983
Analysis Team: A. Cukeriman, BPN

*Fun fact: this plot required training 10k NNs*

Fun fact: this plot required training 10k NNs

*Deep learning + weak supervision + anomaly detection leading to real physics output!*

**ATLAS**

$\sqrt{s}$ = 13 TeV, 139 fb$^{-1}$

$\in$ =0.1, $m_A$ =3000 GeV

- ■ Observed
- - - - Expected
- ±1 σ
- ±2 σ
- ▲ Inclusive search
- ✕ Dedicated (80,80) search

95% CL Exclusion Limit $\sigma(pp \to A \to BC)$ [fb]

← Better

$(m_B, m_C)$ [GeV]

(400,400) (400,200) (400,80) (200,200) (200,80) (80,80)

B A C

ATLAS Collaboration, 2005.02983
Analysis Team: A. Cukeriman, BPN

new ideas!

background model independence

Some searches (train signal versus data)

autoencoders
LDA
ANODE
CWoLa
SALAD

Most searches (train with simulations)

Some searches (train data versus background simulation)

signal model independence

*Rapidly developing area - LHC Olympics 2020 to help facilitate!*



*G. Kasieczka. BPN, D. Shih*
*https://lhco2020.github.io/homepage/*

***We need your great ideas!***

*A. Andreassen.* **BPN***, D. Shih, PRD 101 (2020) 095004*      *M. Farina, Y. Nakai, D. Shih, PRD 101 (2020) 075021*

**BPN***, D. Shih, PRD 101 (2020) 075042*      *T. Heimel, G. Kasieczka, T. Plehn, J. Thompson, SciPost Phys. 6 (2019) 030*

*J. Collins, K. Howe,* **BPN***, PRL 121 (2018) 241803*      *B. Dillon, D. Faroughy, J. Kamenik, PRD 100 (2019) 056002*
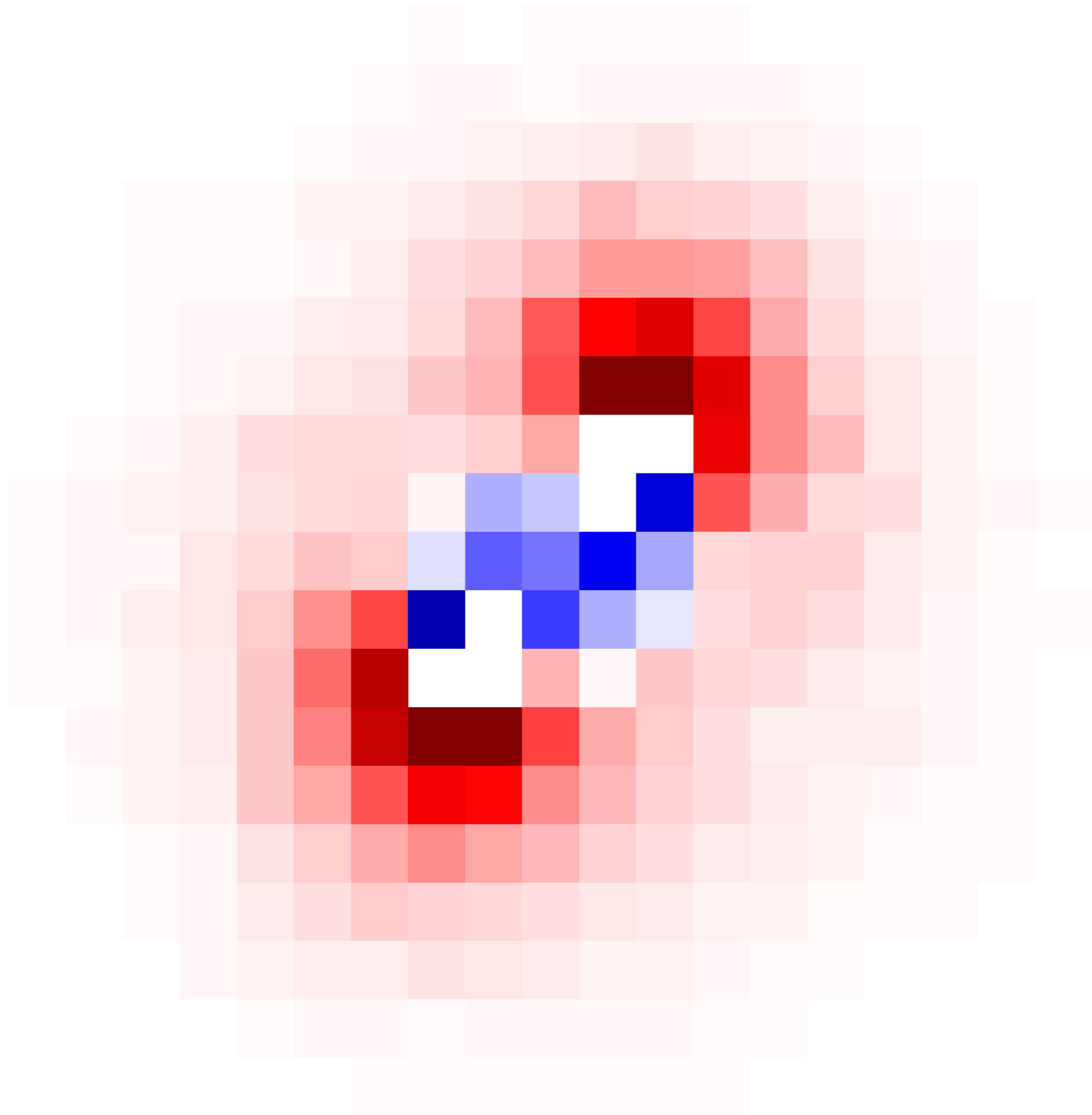
Deep learning has a great potential to **enhance**, **accelerate**, and **empower** HEP analyses.
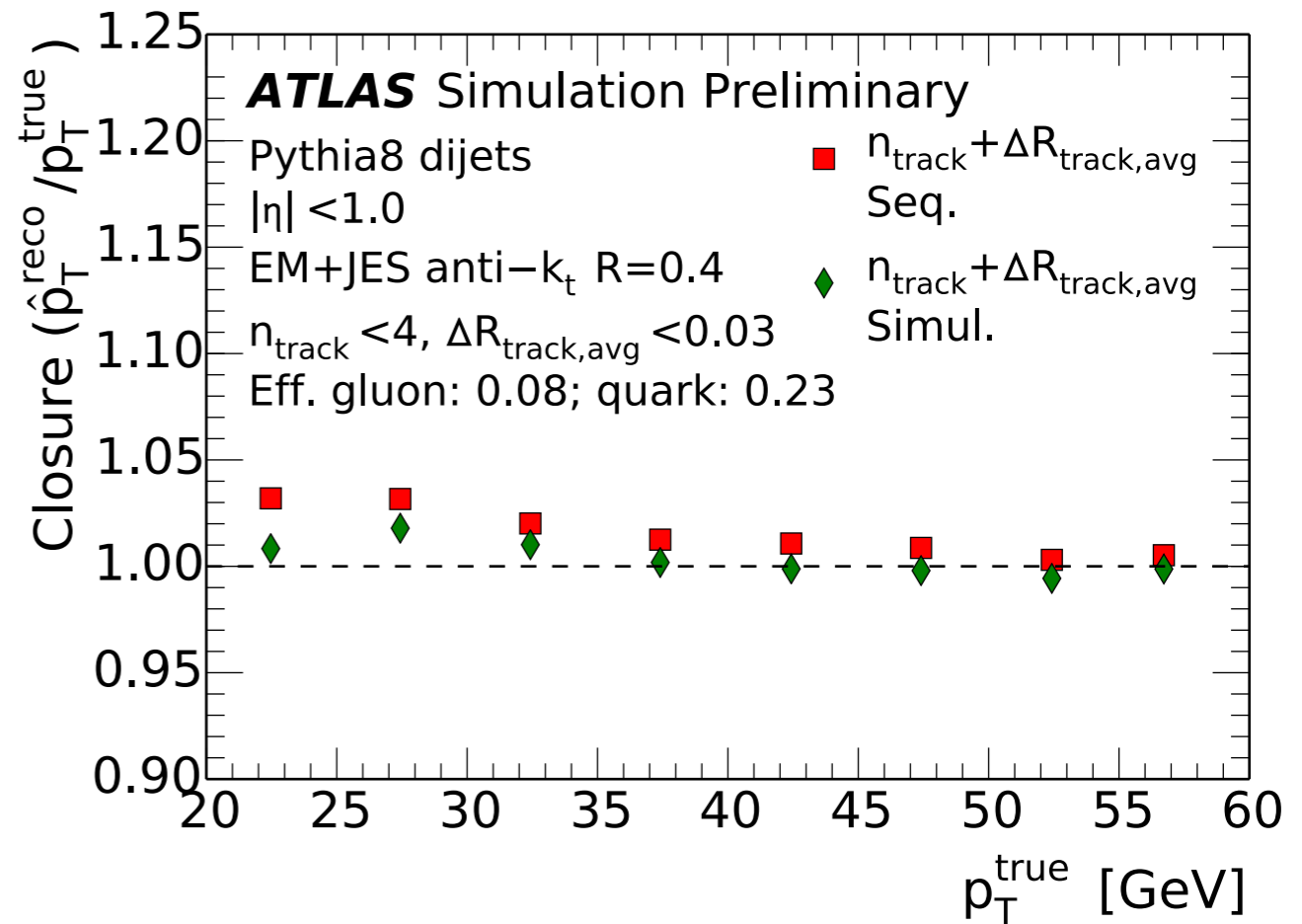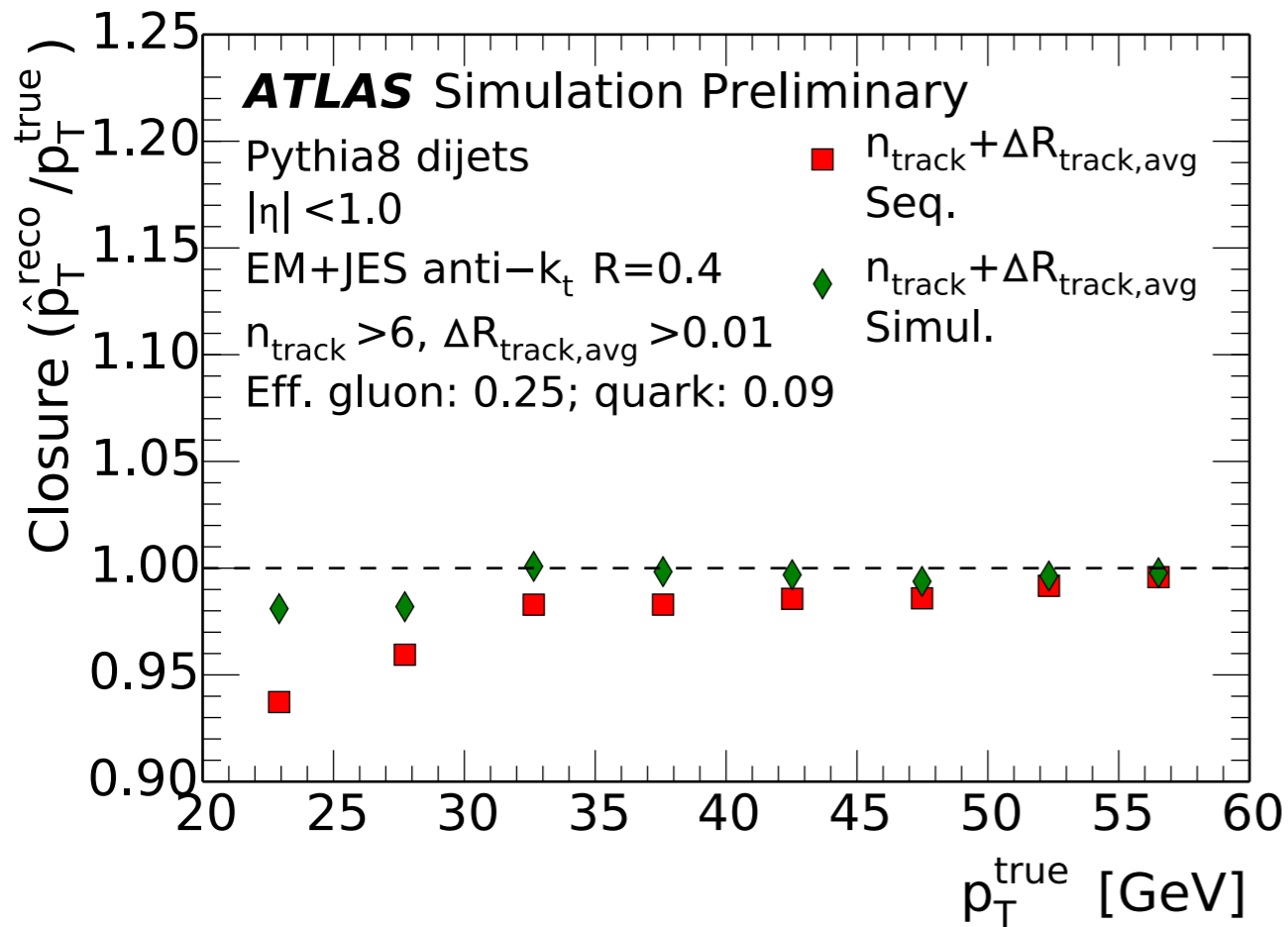
*Disclaimer: I have given you a biased perspective of new developments - there is a growing community within HEP!*



The **full phase space** of our experiments is now explorable and deep learning will allow us this information to discover fundamental properties of nature!
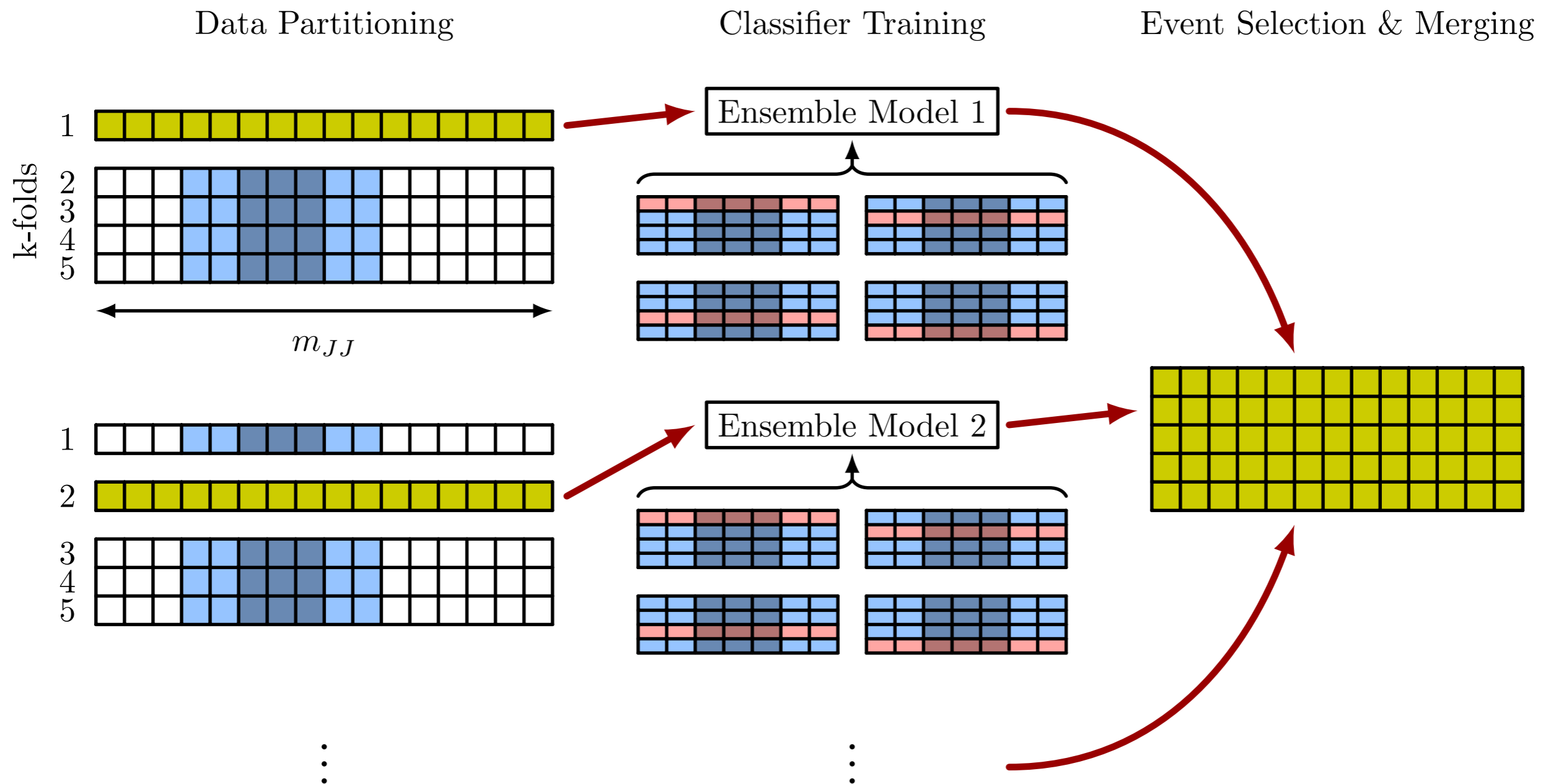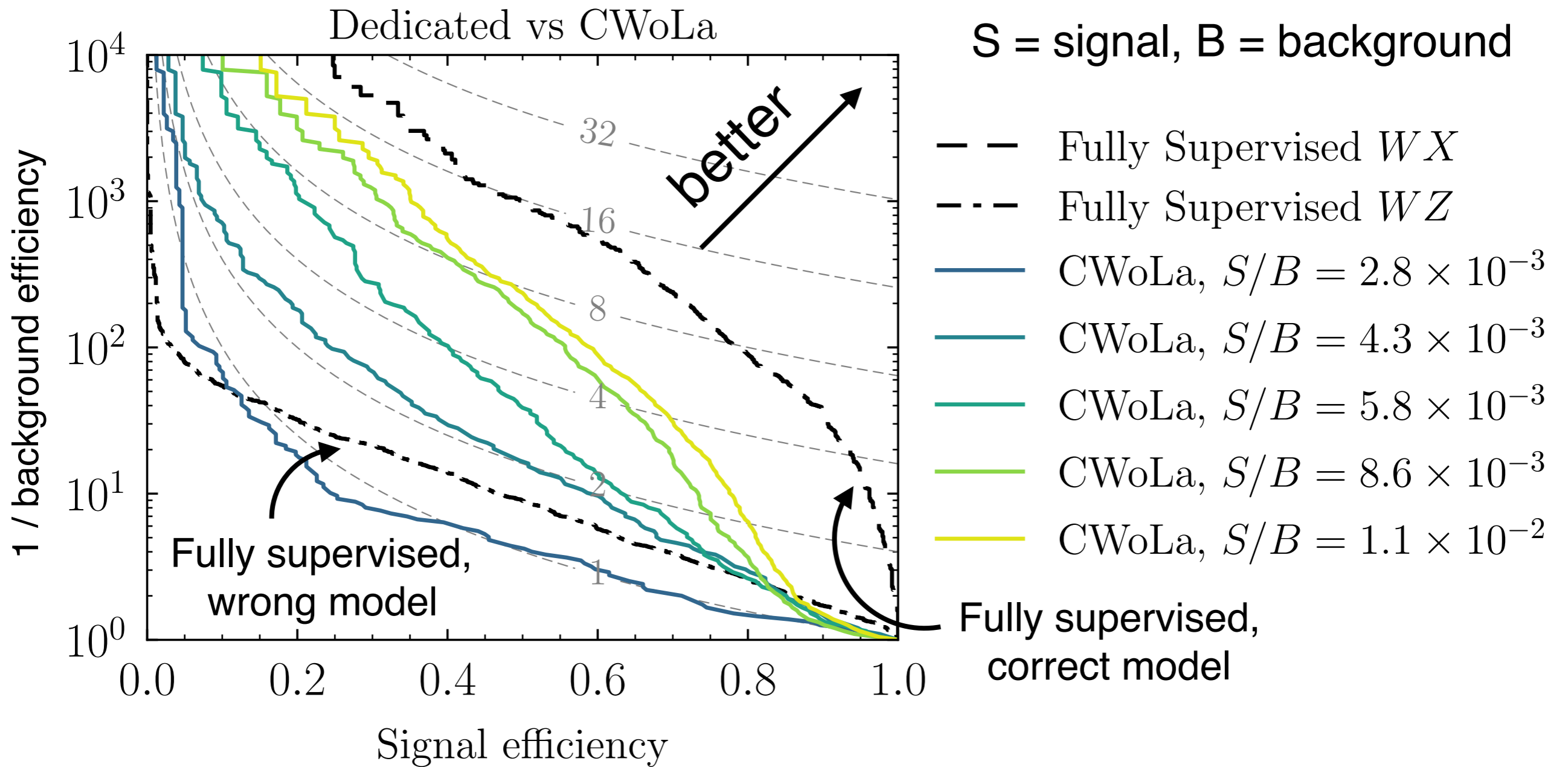
Fin.

Slightly better closure for the simultaneous calibration.

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

Dedicated vs CWoLa

S = signal, B = background

- – – Fully Supervised $WX$
- – · – Fully Supervised $WZ$
- —— CWoLa, $S/B = 2.8 \times 10^{-3}$
- —— CWoLa, $S/B = 4.3 \times 10^{-3}$
- —— CWoLa, $S/B = 5.8 \times 10^{-3}$
- —— CWoLa, $S/B = 8.6 \times 10^{-3}$
- —— CWoLa, $S/B = 1.1 \times 10^{-2}$

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!