# Overview of Machine Learning for Particle Physics

Convolution    Max-Pool

Jet Image

## Benjamin Nachman

*Lawrence Berkeley National Laboratory*

bpnachman.com    @bpnachman    bnachman
bpnachman@lbl.gov

**US ATLAS ML Training**
July 26, 2023

# Particle Physics + Machine Learning

**3**

**Theory of everything**

**Physics simulators**

Detector-level observables

↓

Pattern recognition

Fast simulation / phase space

Parameter estimation / unfolding

**Nature**

↓

Experiment

↓

Detector-level observables

↓

Pattern recogn

Experimental Design

Online processing & quality control

Data curation

Inference

Dimensionality reduction; "signal" versus "background"

calibration
clustering
tracking
noise mitigation
particle identification
...

Theory of everything

↓

**Physics simulators**

↓

Detector-level observables

↓

Pattern recognition

←→

**Nature**

↓

Experiment

↓

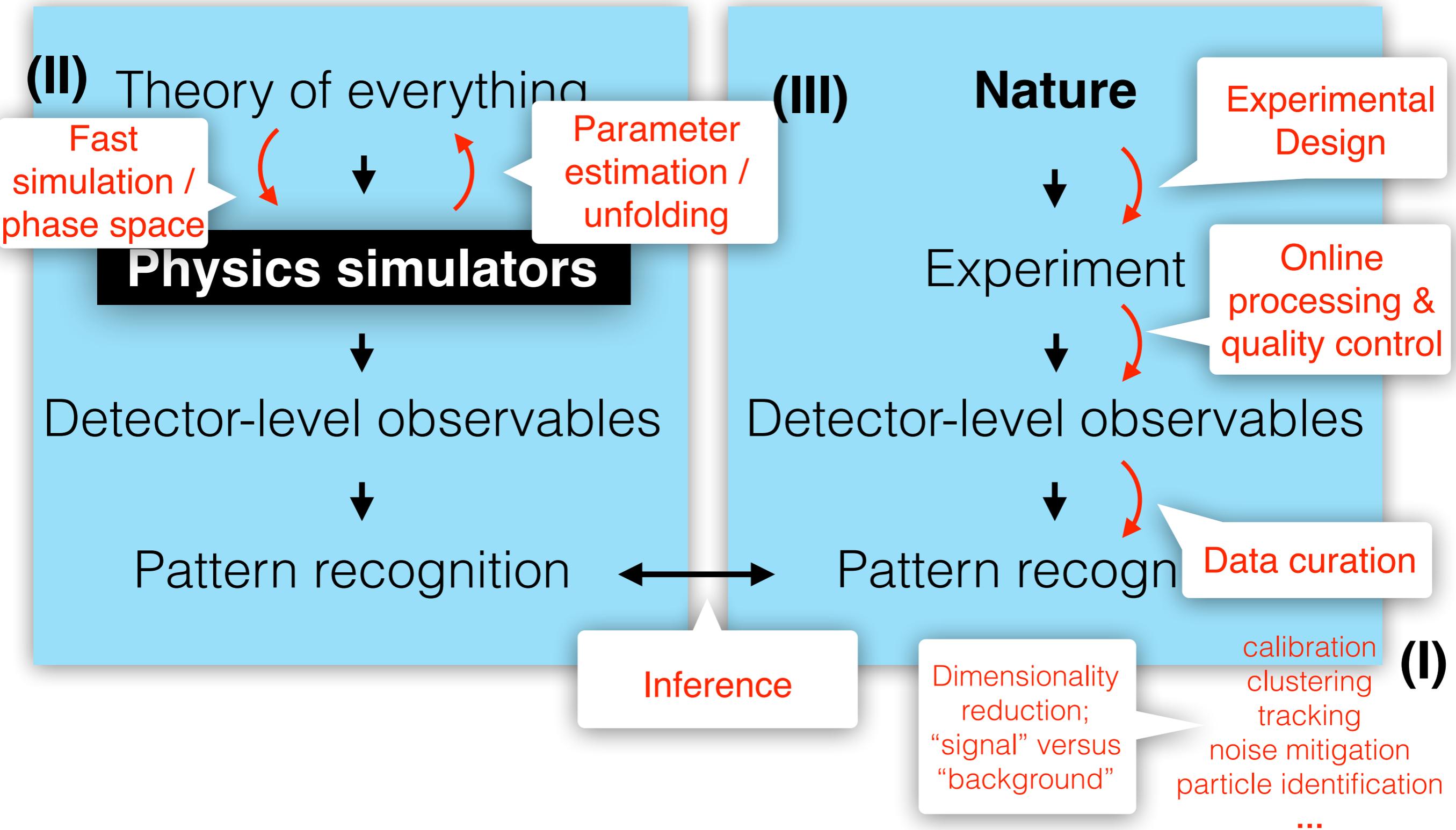Detector-level observables

↓

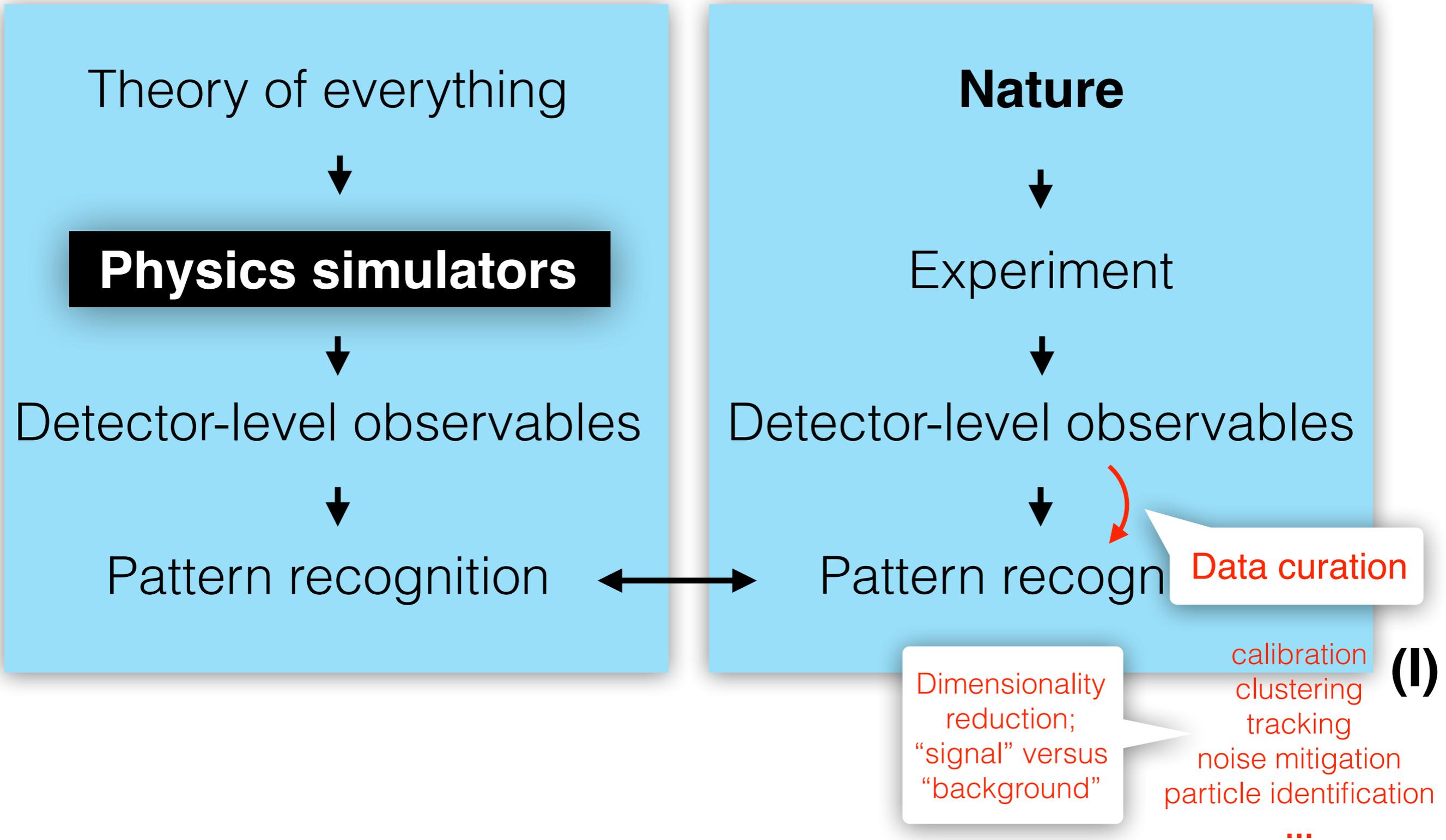Pattern recogn

Data curation

**(I)**

Dimensionality reduction; "signal" versus "background"

calibration
clustering
tracking
noise mitigation
particle identification
...

# Step 1: how to represent our data

**6**

**Fixed sets**

$$J = \{p_1^\mu, p_2^\mu, \ldots, p_n^\mu\}$$

*Dense networks*

**Variable sets**

*Deep sets*

**Images**

*Convolutional NNs (CNNs)*

**Sequences**

*Recurrent NNs*

**Trees**

*Recursive NNs*

**Graphs**

*Graph CNNs*

# Step 1: how to represent our data

Fixed
sets

$J = \{p_1^\mu, p_2^\mu, \ldots, p_n^\mu\}$

Dense networks

Variable
sets

Deep sets

**Images**

Convolutional NNs (CNNs)

Sequences

Recurrent NNs

Graph CNNs

Trees

Recursive NNs

Graphs

*p*

*p*

*p*

*p*

Full Event Information

Neutral $p_T$

Charge Multiplicity

Charge $p_T$

Charge Multiplicity

Repeat

Dense

Full Event Information

Neutral $p_T$   Charge Multiplicity

Neutral $p_T$   Charge Multiplicity

Charge Multiplicity

Padding Layer

Convolutional Layer

Pooling

Higgs or not?

Convolutional Layer   Pooling   Dense
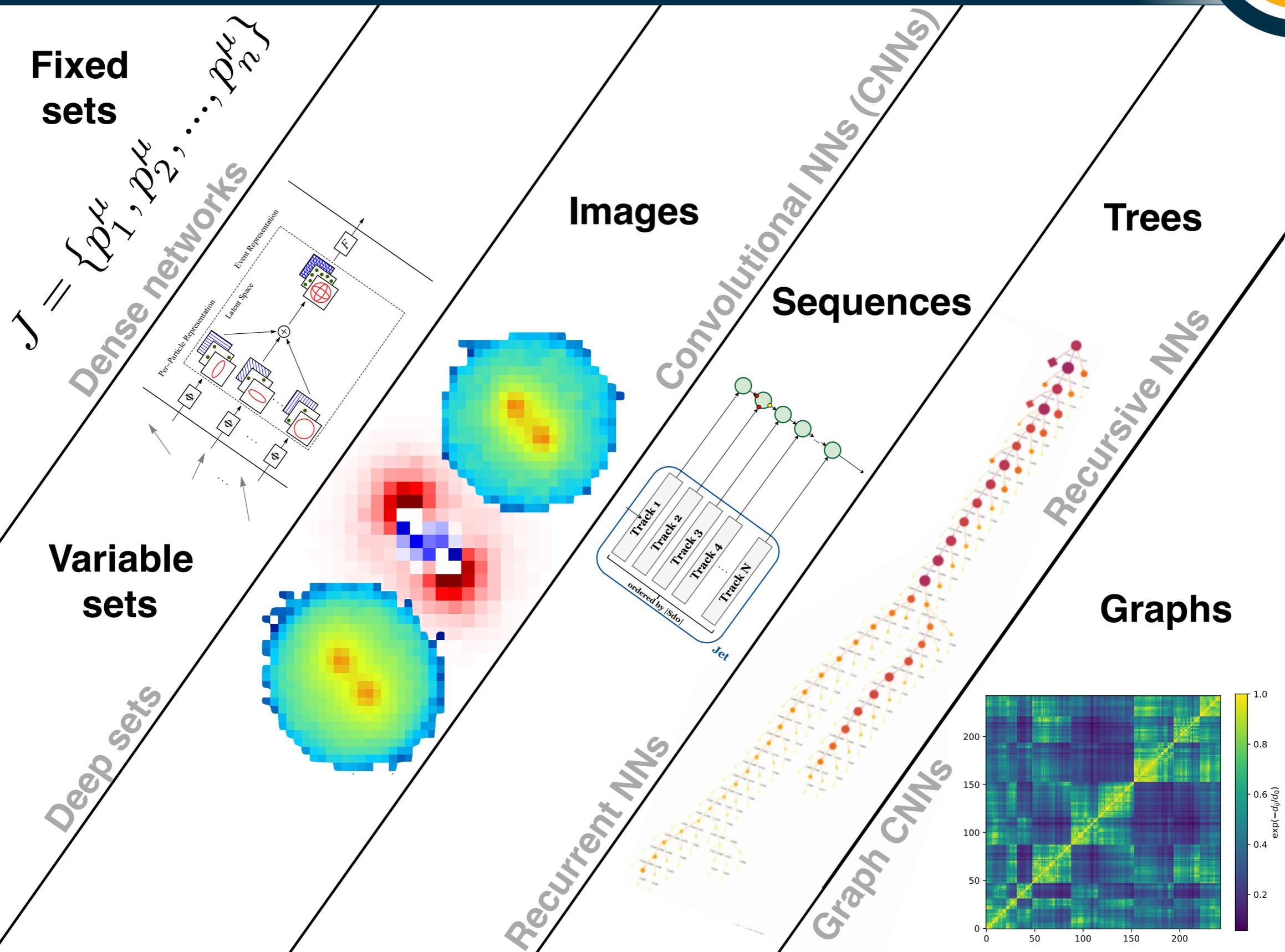
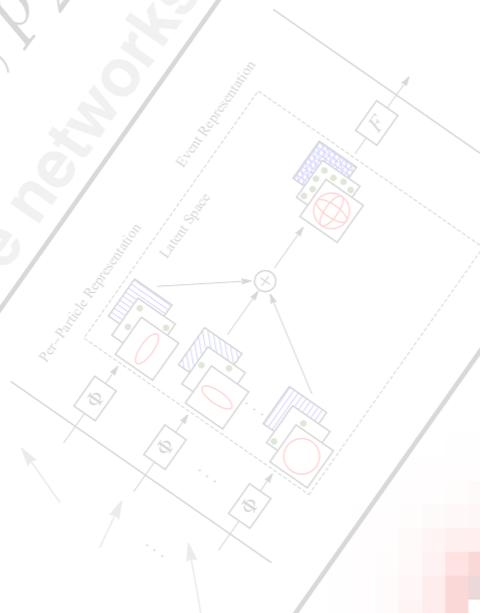Can combine local and global information from jet images and "event" images.

# Step 1: how to represent our data

**Fixed sets**

$$J = \{p_1^\mu, p_2^\mu, \ldots, p_n^\mu\}$$

*Dense networks*

**Images**

*Convolutional NNs (CNNs)*

**Trees**

*Recursive NNs*

**Sequences**

**Variable sets**

**Graphs**

*Deep sets*

*Recurrent NNs*

*Graph CNNs*

# Step 1: how to represent our data

Fixed sets

$J = \{p_1^\mu, p_2^\mu, \ldots, p_n^\mu\}$

Dense networks

Variable sets

Deep sets

Images

Convolutional NNs (CNNs)

**Sequences**

Track 1
Track 2
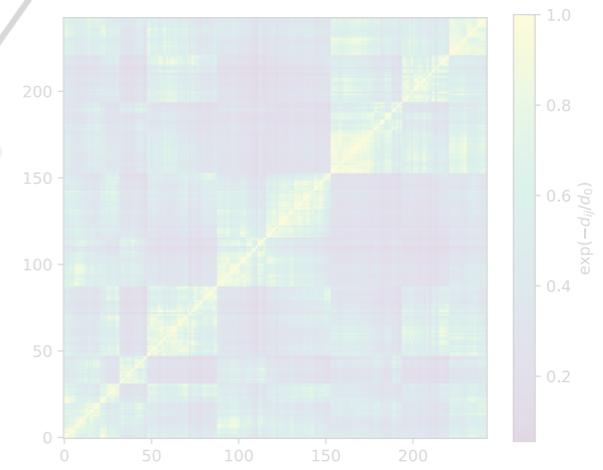Track 3
Track 4
...
Track N

ordered by |Sd0|

Jet

Recurrent NNs

Graph CNNs

Trees

Recursive NNs

Graphs

One key challenge with images is that they have a fixed size.

*In many contexts, this is ideal, because the data also have a fixed size. However, this is not always the case.*

For example, events / jets have a variable number of particles.

One can represent these particles as a sequence in order to apply variable-length approaches that can access the full feature granularity.

Flavor tagging (classify jets from b-quark or not) has a long history of ML.  Use features of the charged-particle tracks inside jets.

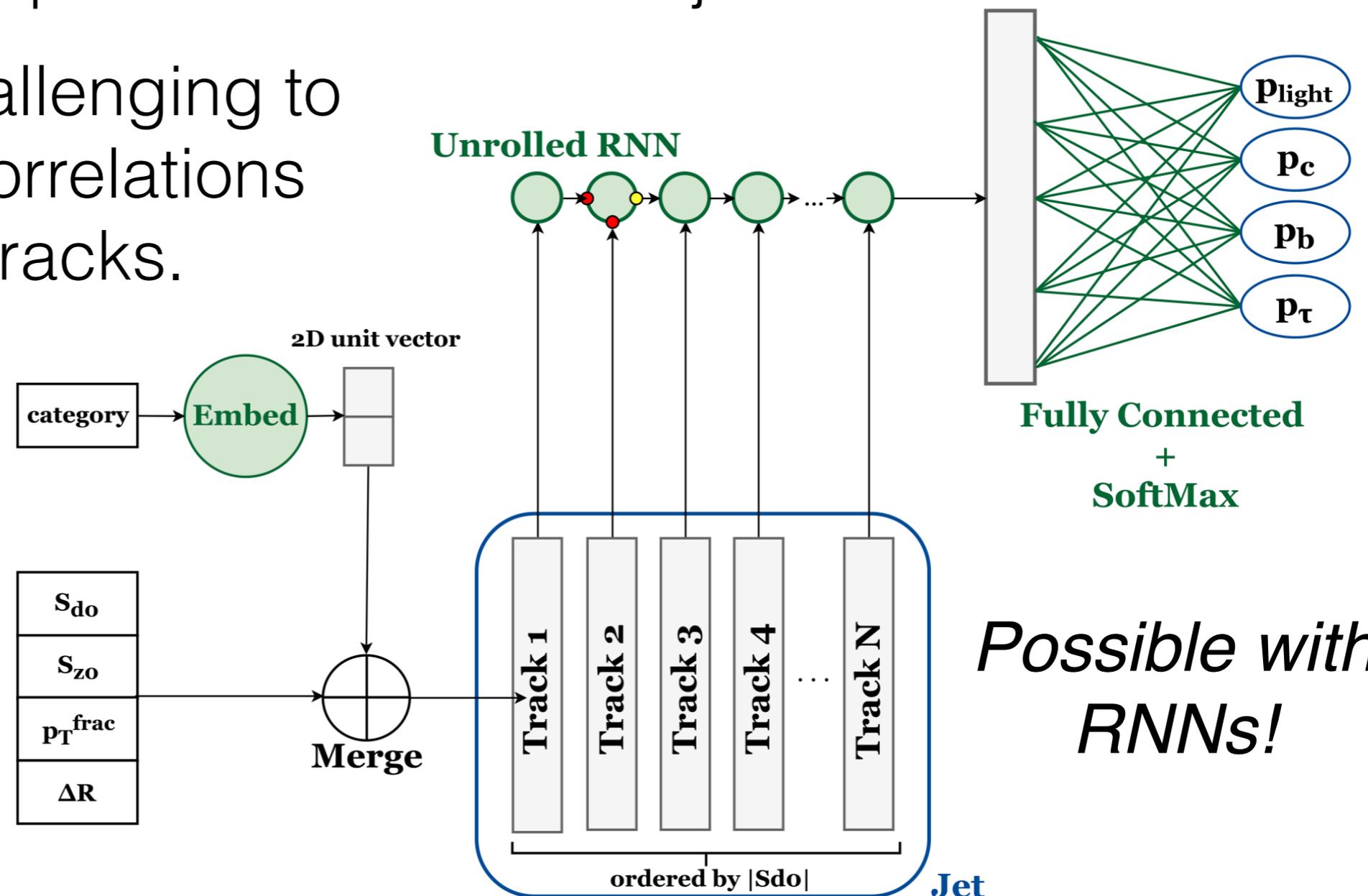In the past, challenging to incorporate correlations between tracks.

Flavor tagging (classify jets from b-quark or not) has a long history of ML.  Use features of the charged-particle tracks inside jets.

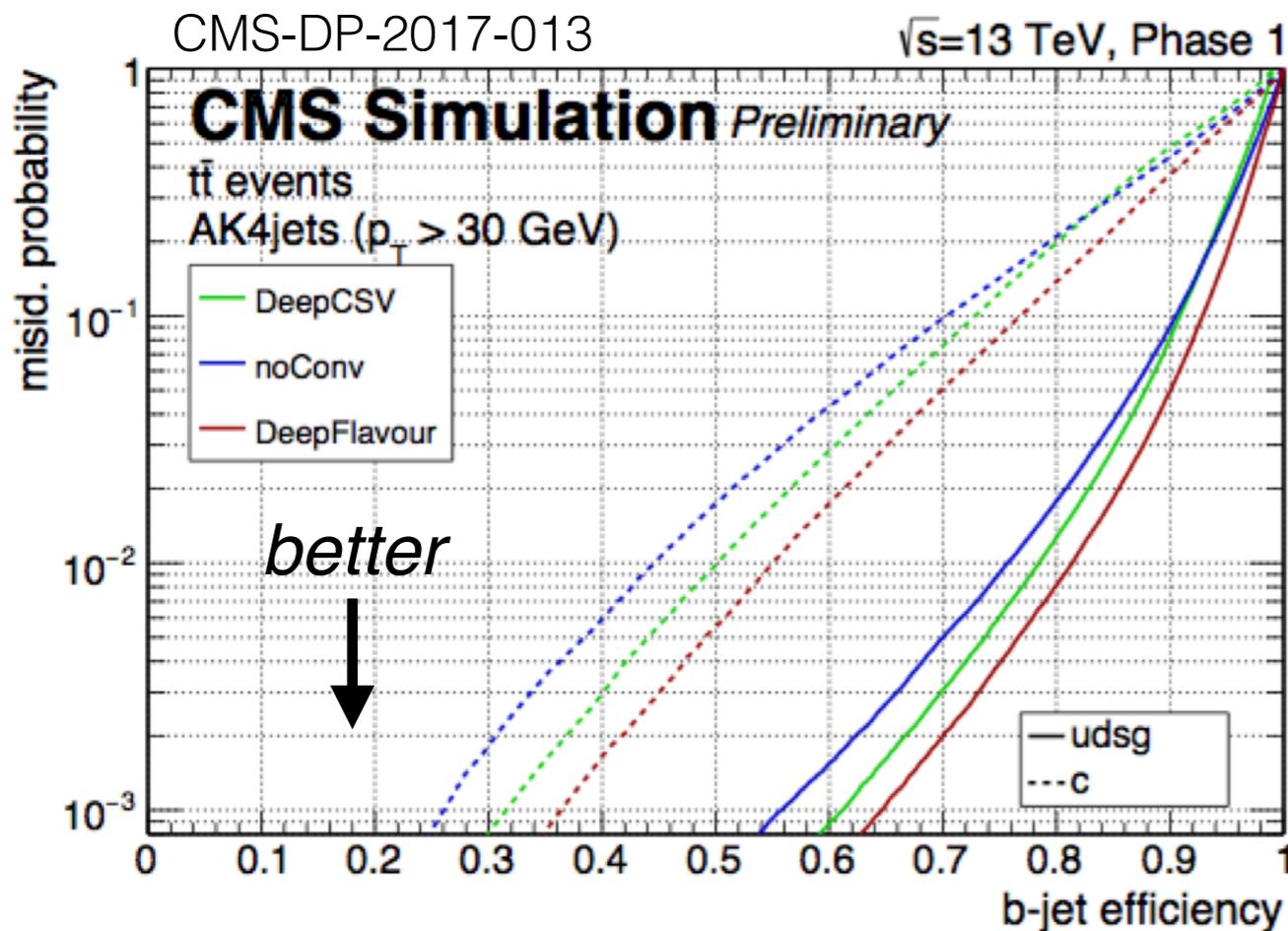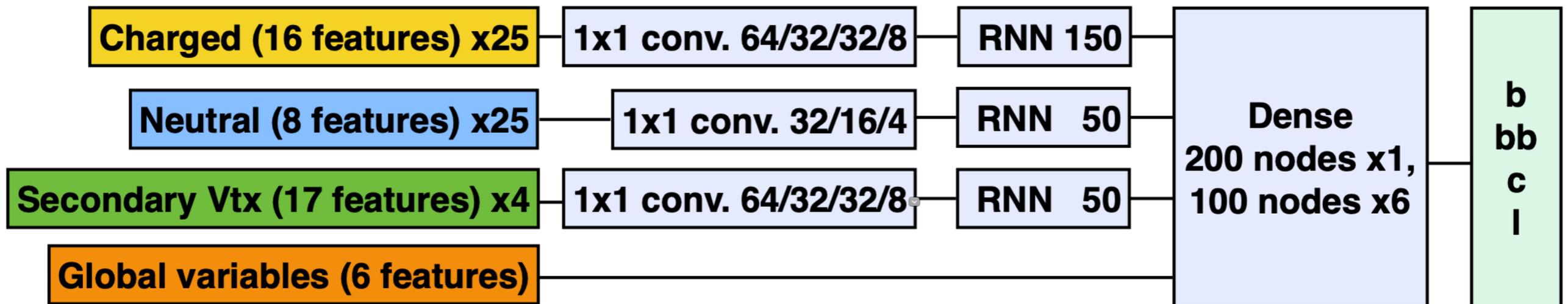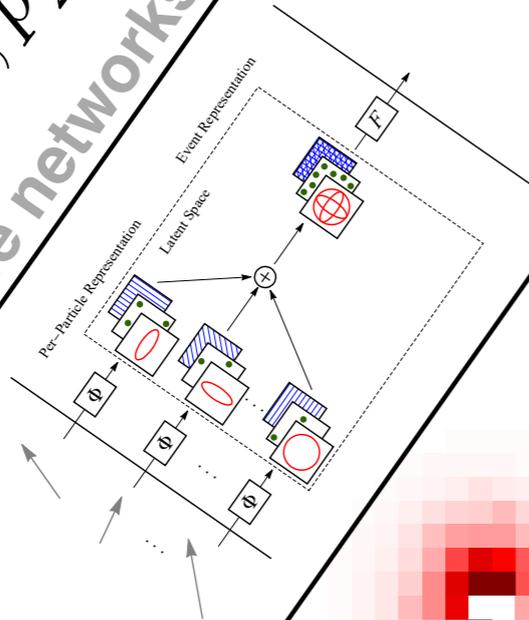In the past, challenging to incorporate correlations between tracks.

**Unrolled RNN**

**2D unit vector**

category → **Embed** →

| $S_{do}$ |
| $S_{zo}$ |
| $p_T^{frac}$ |
| $\Delta R$ |

**Merge**

Track 1 · Track 2 · Track 3 · Track 4 · … · Track N

ordered by |Sdo|

**Jet**

$p_{light}$

$p_c$

$p_b$

$p_\tau$

**Fully Connected + SoftMax**

*Possible with RNNs!*

# Hybrid methods

| Charged (16 features) x25 | 1x1 conv. 64/32/32/8 | RNN 150 | | b |
| Neutral (8 features) x25 | 1x1 conv. 32/16/4 | RNN 50 | Dense 200 nodes x1, 100 nodes x6 | bb c l |
| Secondary Vtx (17 features) x4 | 1x1 conv. 64/32/32/8 | RNN 50 | | |
| Global variables (6 features) | | | | |



CMS-DP-2017-013

RNN + 1x1 CNNs for dimensionality reduction.

This reduction improved the performance of the overall classifier.

**Fixed sets**

$$J = \{p_1^\mu, p_2^\mu, \ldots, p_n^\mu\}$$

**Dense networks**

**Images**

**Convolutional NNs (CNNs)**

**Sequences**

**Trees**

**Recursive NNs**

**Variable sets**

Track 1
Track 2
Track 3
Track 4
Track N
ordered by |Sd0|
Jet

**Deep sets**

**Recurrent NNs**

**Graph CNNs**

**Graphs**

Fixed
sets

$J = \{p_1^\mu, p_2^\mu, \ldots, p_n^\mu\}$

Dense networks

Variable
sets
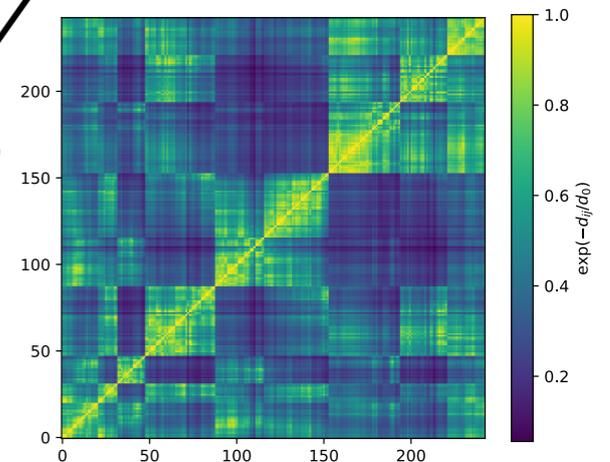
Deep sets
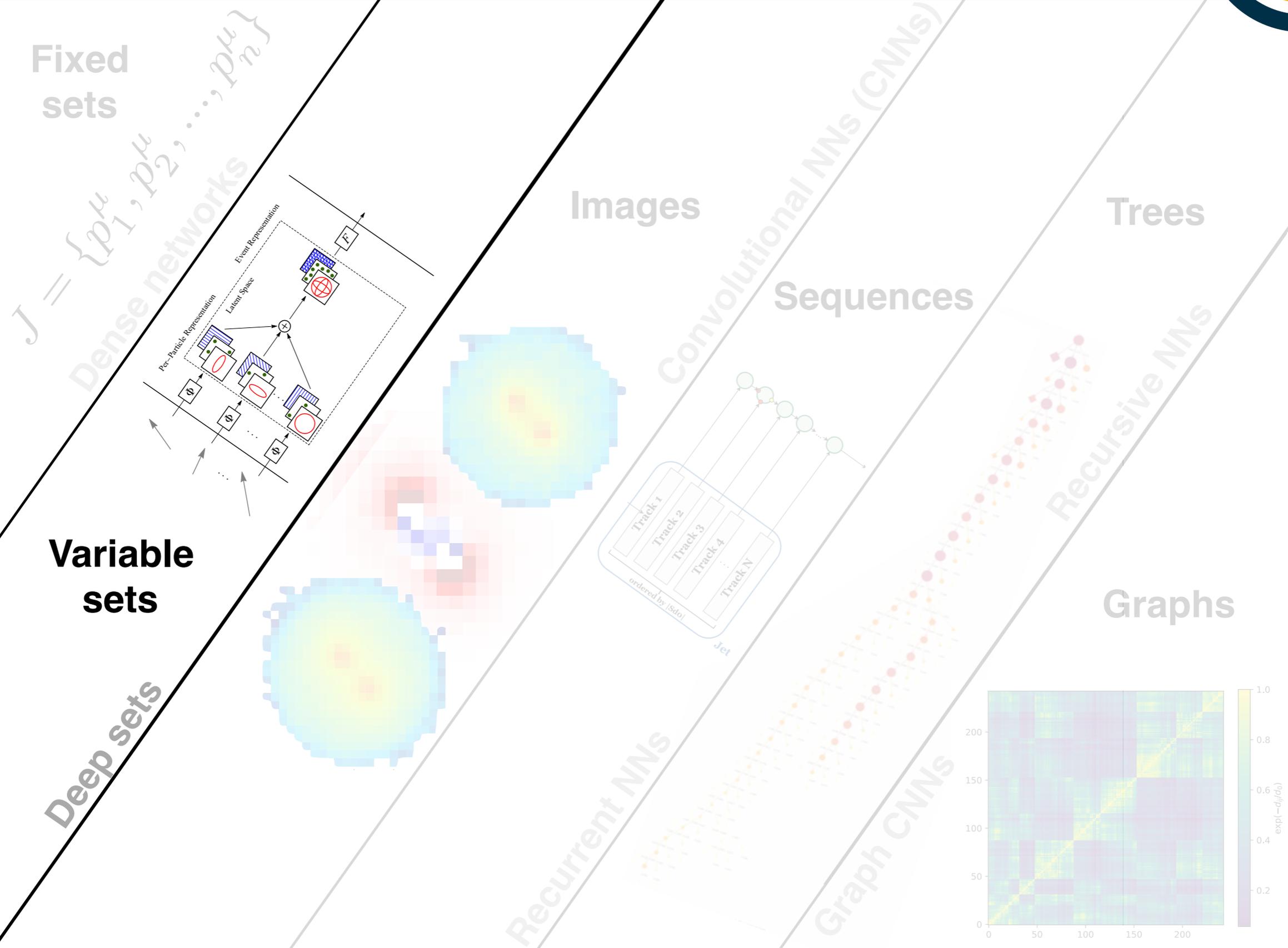
Images

Convolutional NNs (CNNs)

Sequences

Track 1
Track 2
Track 3
Track 4
Track N
ordered by |Sd0|
Jet

Recurrent NNs

Trees

Recursive NNs

Graphs

Graph CNNs

A challenge with sequence learning is that thanks to quantum mechanics, there is often no unique order.

A common scenario is that we have a variable-length **SET** of particles and we would like to learn from them directly.

Solution: set learning / point cloud approaches

Factorize the problem into two networks: one that embeds the set into a fixed-length latent space and one that acts on a permutation invariant operation on that latent space:
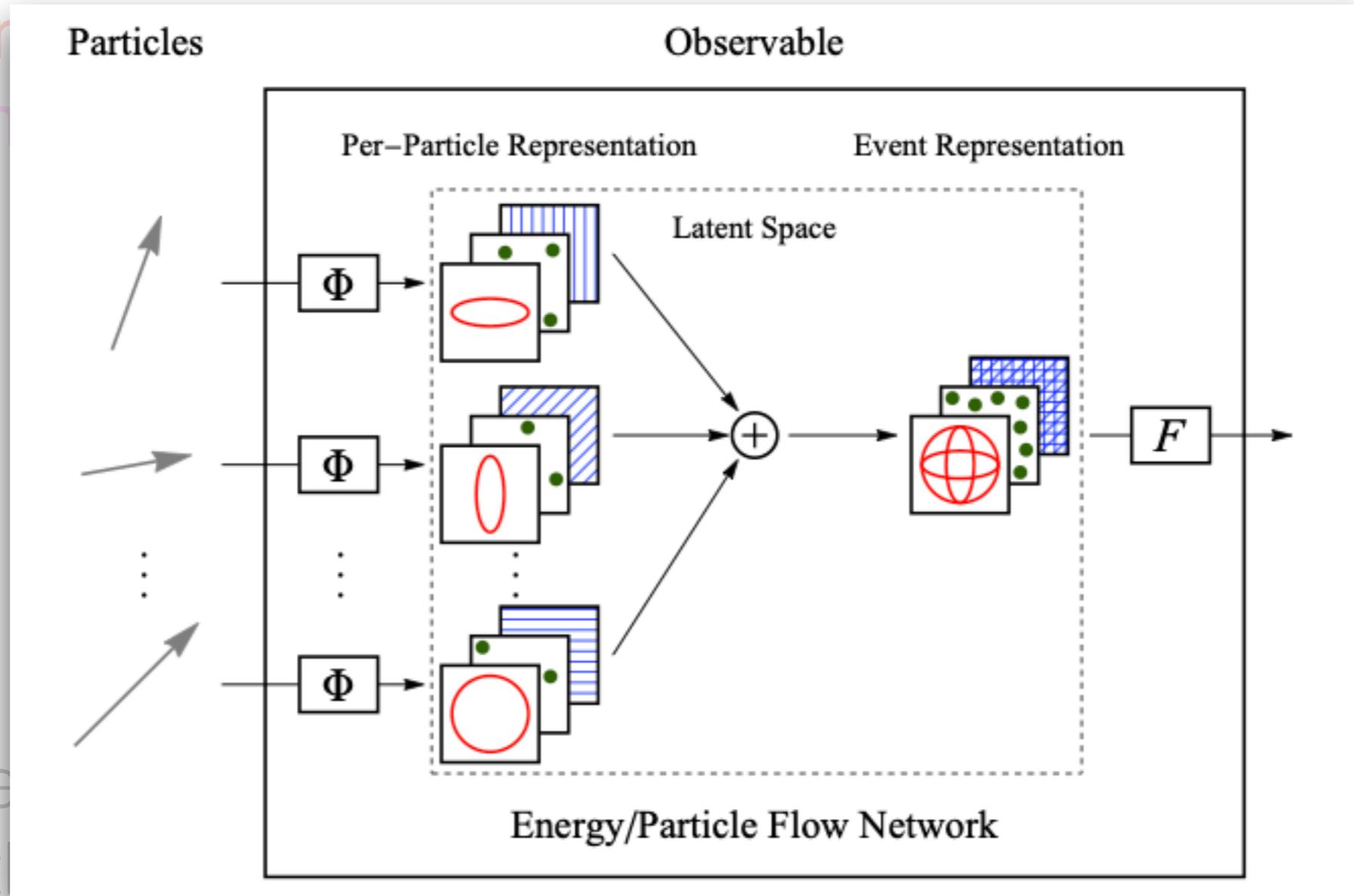
$$f(\{x_1, \ldots, x_M\}) = F\left(\sum_{i=1}^{M} \Phi(x_i)\right)$$

Due to the sum, this structure can operate on any length set and the order of the inputs doesn't matter.

Factorize the problem into two networks: one that embeds the set in [...] acts on a perm[...] space:



Particles

Observable

Per−Particle Representation

Event Representation

Latent Space

$\Phi$

$\Phi$

$\Phi$

$+$

$F$

Energy/Particle Flow Network

Due [...] any
lengt[...] atter.

- Can readily incorporate per-particle features
- Can be made infrared and collinear safe (EFN) safe

Latent space in IRC safe case is interpretable (and predictable!)

ATL-PHYS-PUB-2020-014

Faster to train than RNN so can do R&D on input features to improve overall performance.

Latent space in IRC safe case is interpretable (and predictable!)

M. Zaheer et al. https://arxiv.org/abs/1703.06114; P. Komiske, E. Metodiev, & J. Thaler, JHEP 01 (2019) 12

Classic CNNs operate on a fixed grid and are not invariant under the permutation of points

Can generalize CNNs to act on graphs



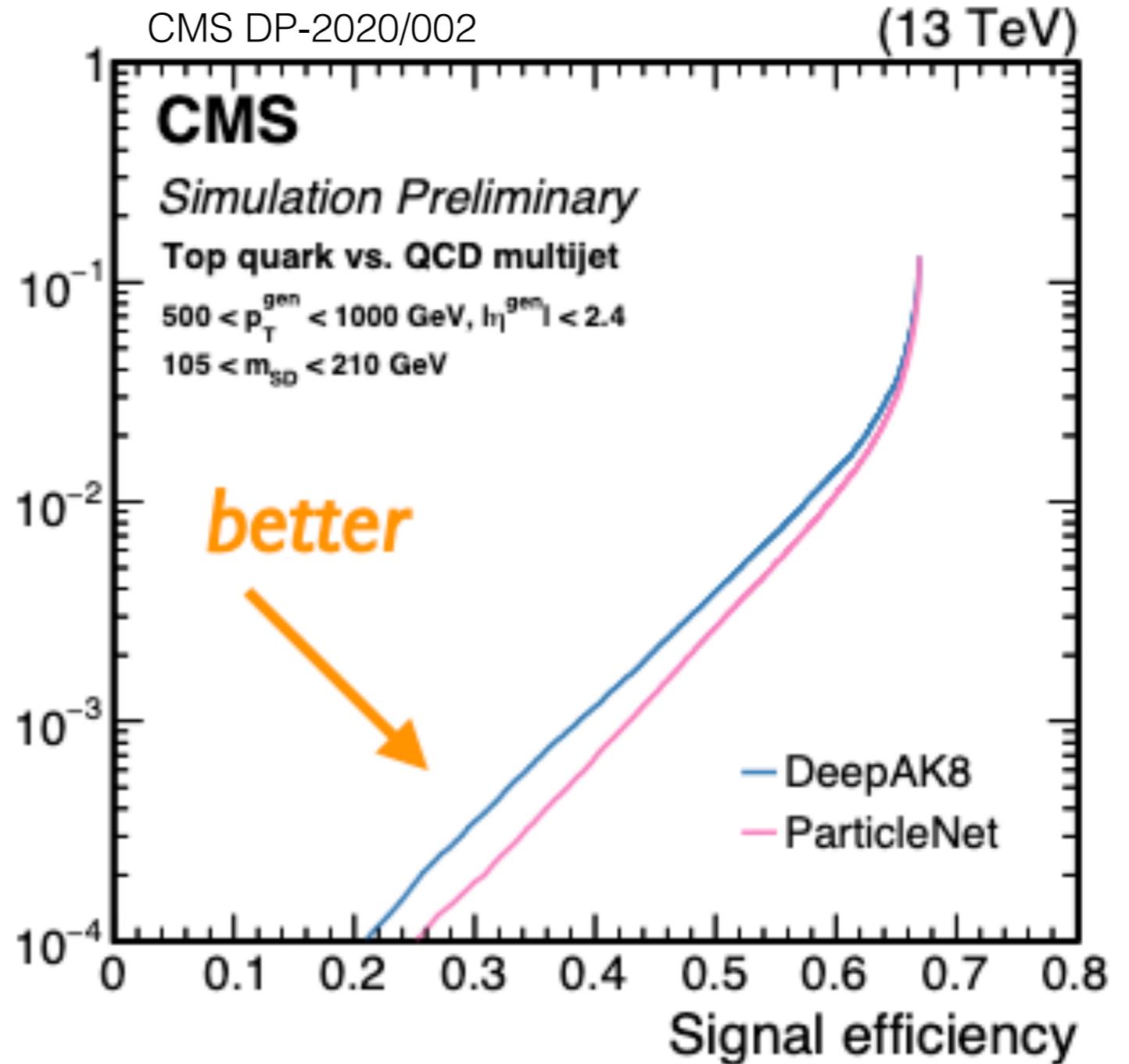Need to define distances using particle properties

**See also Javier's talk!**

Classic CNNs operate on a fixed grid and are not invariant un

Can generaliz

**Competitive performance to other state-of-the-art methods**

Need to define dis

CMS DP-2020/002

(13 TeV)

**CMS**

*Simulation Preliminary*

**Top quark vs. QCD multijet**

$500 < p_T^{gen} < 1000$ GeV, $|\eta^{gen}| < 2.4$

$105 < m_{SD} < 210$ GeV

Background efficiency

Signal efficiency

**better**

— DeepAK8

— ParticleNet

1801.07829 , 1902.08570

# Step 1: how to represent our data

**Fixed sets**

$$J = \{p_1^\mu, p_2^\mu, \ldots, p_n^\mu\}$$

Dense networks

**Images**

Convolutional NNs (CNNs)

**Trees**

Recursive NNs

**Sequences**

**Variable sets**

Track 1
Track 2
Track 3
Track 4
Track N

ordered by |Sd0|

Jet

**Graphs**

Deep sets

Recurrent NNs

Graph CNNs

One way to categorize methods is based on their level of *supervision*

**Unsupervised** = no labels
**Weakly-supervised** = noisy labels
**Semi-supervised** = partial labels
**Supervised** = full label information

This is 99% of the ML.   We have labeled examples and we train a model to predict the labels from the examples.

Need to be careful about what loss function to pick (more on that in a little bit…)

**Unsupervised** = no labels

Typically, the goal of these methods is to implicitly or explicitly estimate $p(x)$.



One strategy (autoencoders) is to try to compress events and then uncompress them. When $x$ is far from uncompres(compress($x$)), then $x$ probably has low $p(x)$.

Talking point: anomaly detection!

**Weakly-supervised** = noisy labels

Typically, the goal of these methods is to estimate
*p(possibly signal-enriched)/p(possibly signal-depleted)*



Signal enriched

Signal depleted

**Semi-supervised** = partial labels

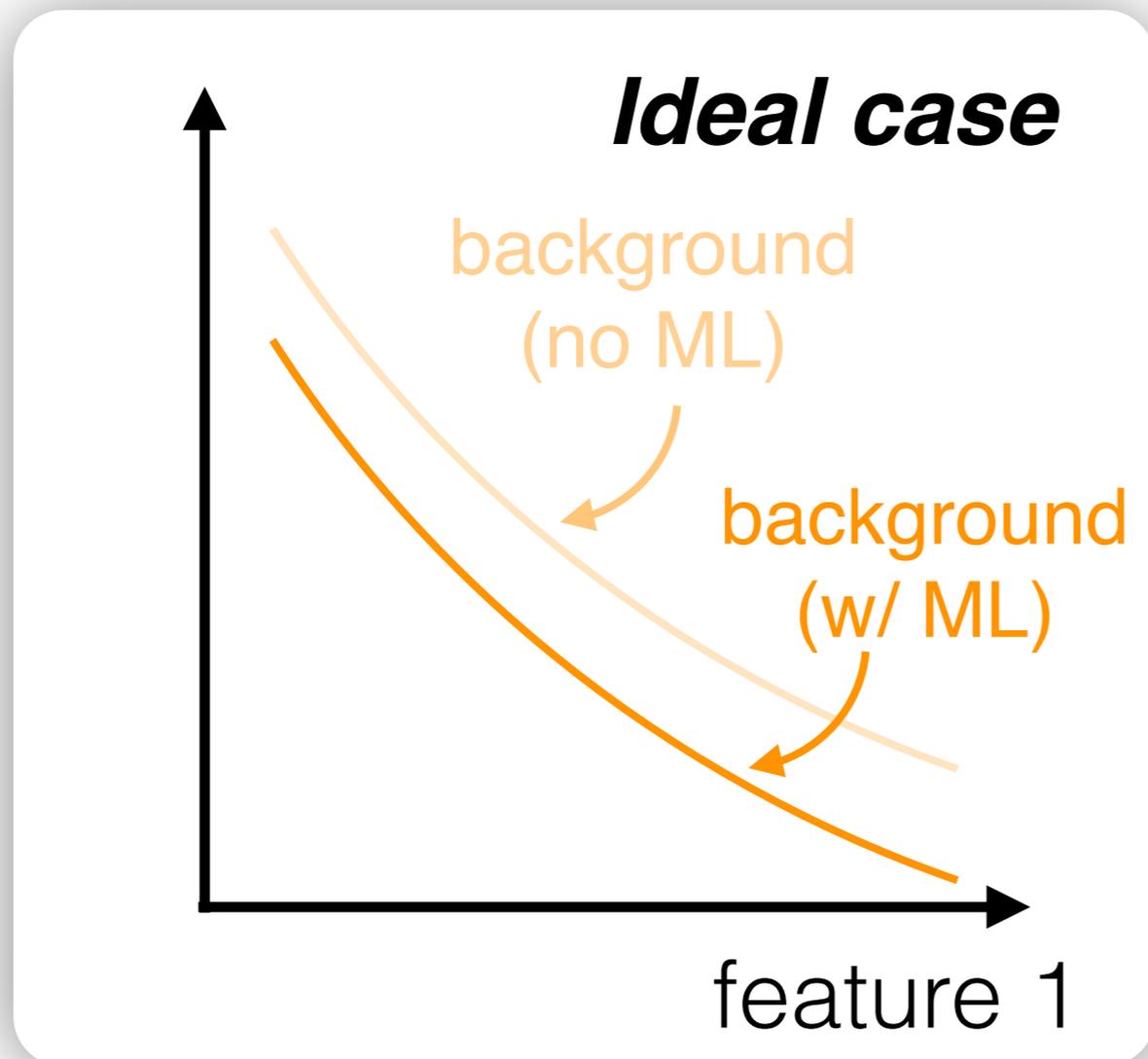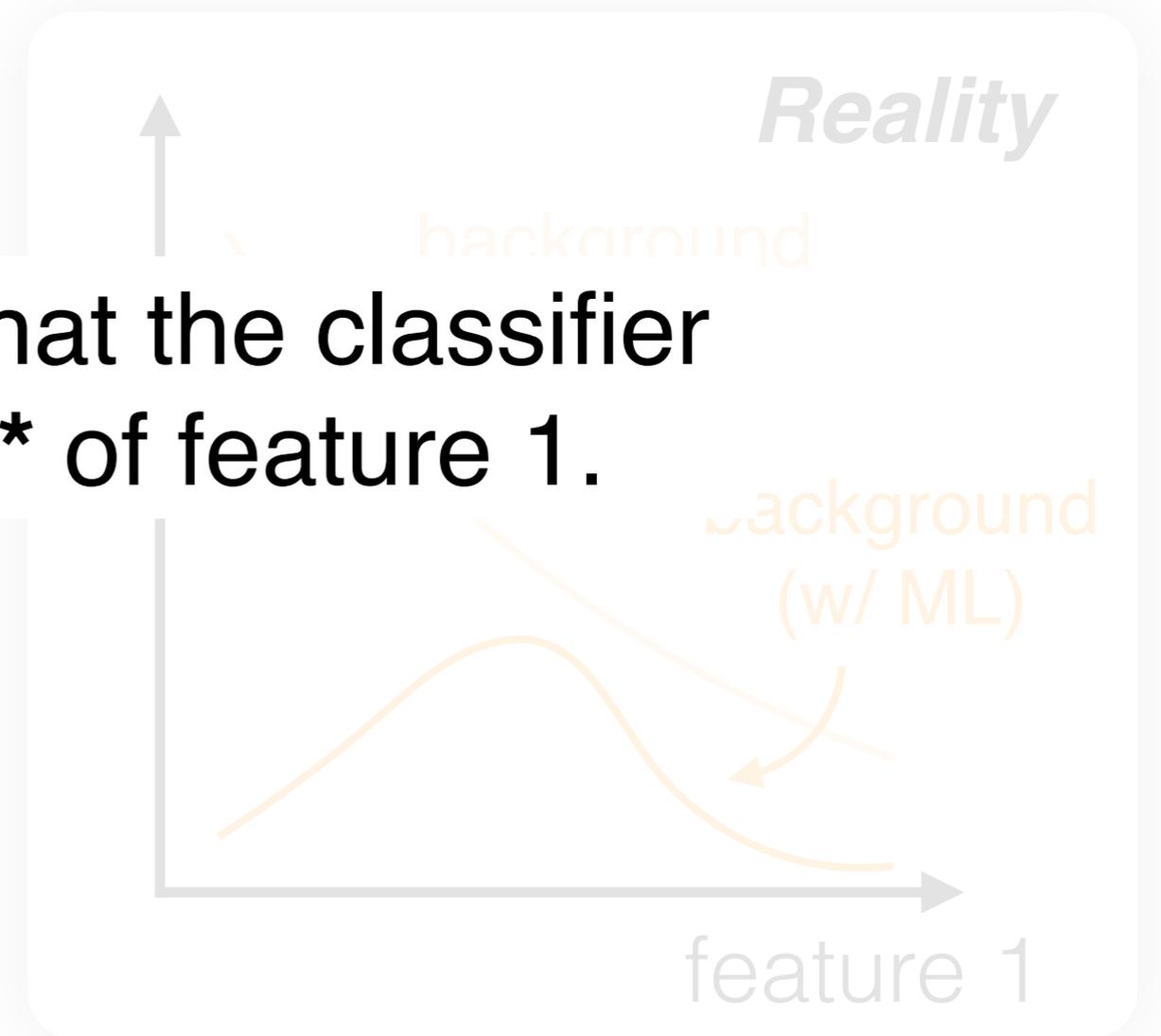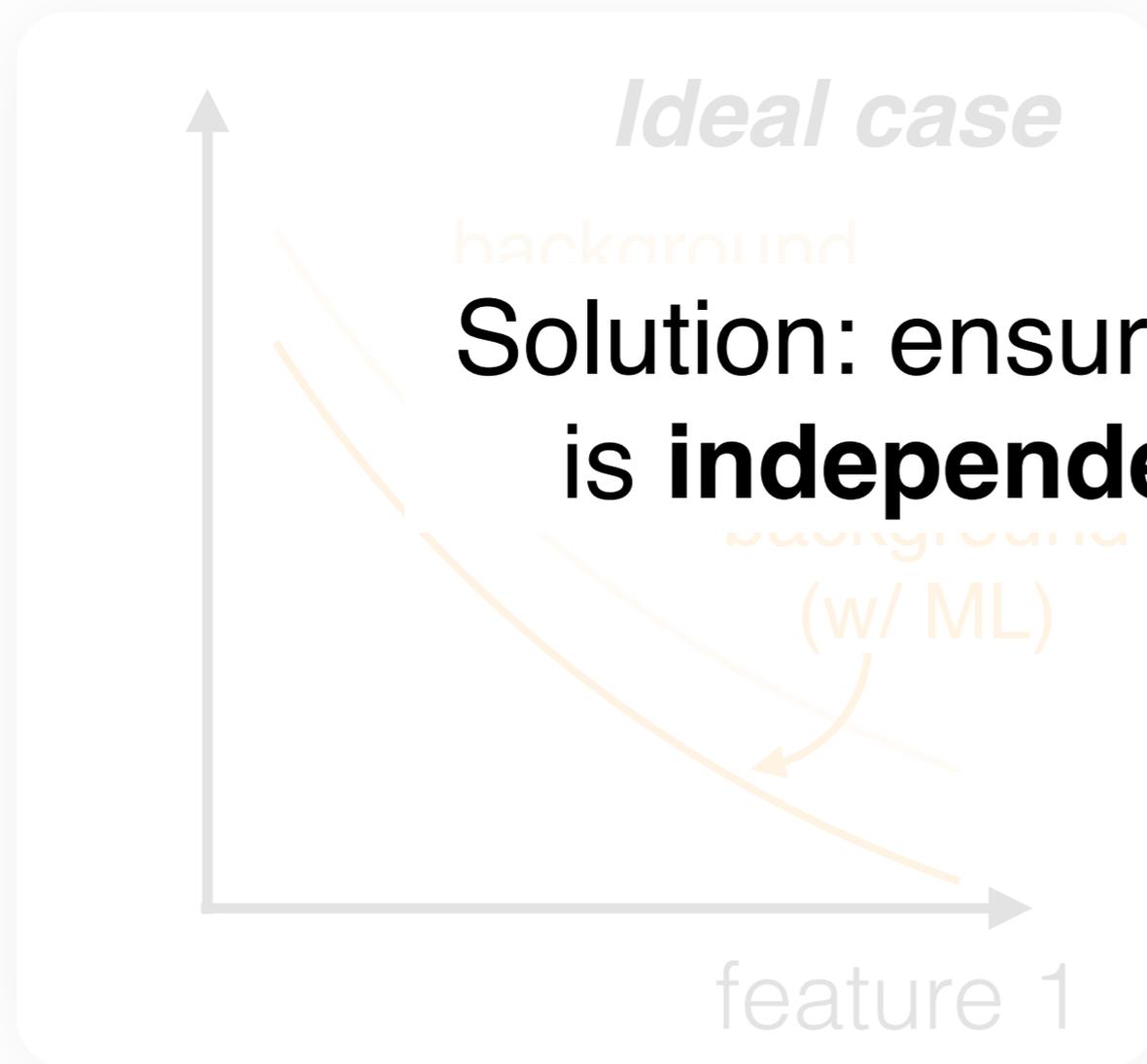Typically, these methods use some signal simulations to build signal sensitivity



Image credit: https://www.particlezoo.net

**vs**

e.g. SM background versus many signals

*How can we learn a classifier that does not sculpt a bump in the background?*

**Ideal case**

background
(no ML)

background
(w/ ML)

feature 1

**Reality**

background
(no ML)

background
(w/ ML)

feature 1

*How can we learn a classifier that does not sculpt a bump in the background?*

*Ideal case*

*Reality*

Solution: ensure that the classifier is **independent\*** of feature 1.

feature 1

feature 1

*\*This is actually sufficient but unnecessary.  There are many dependencies (e.g. linear) that would not sculpt bumps.*
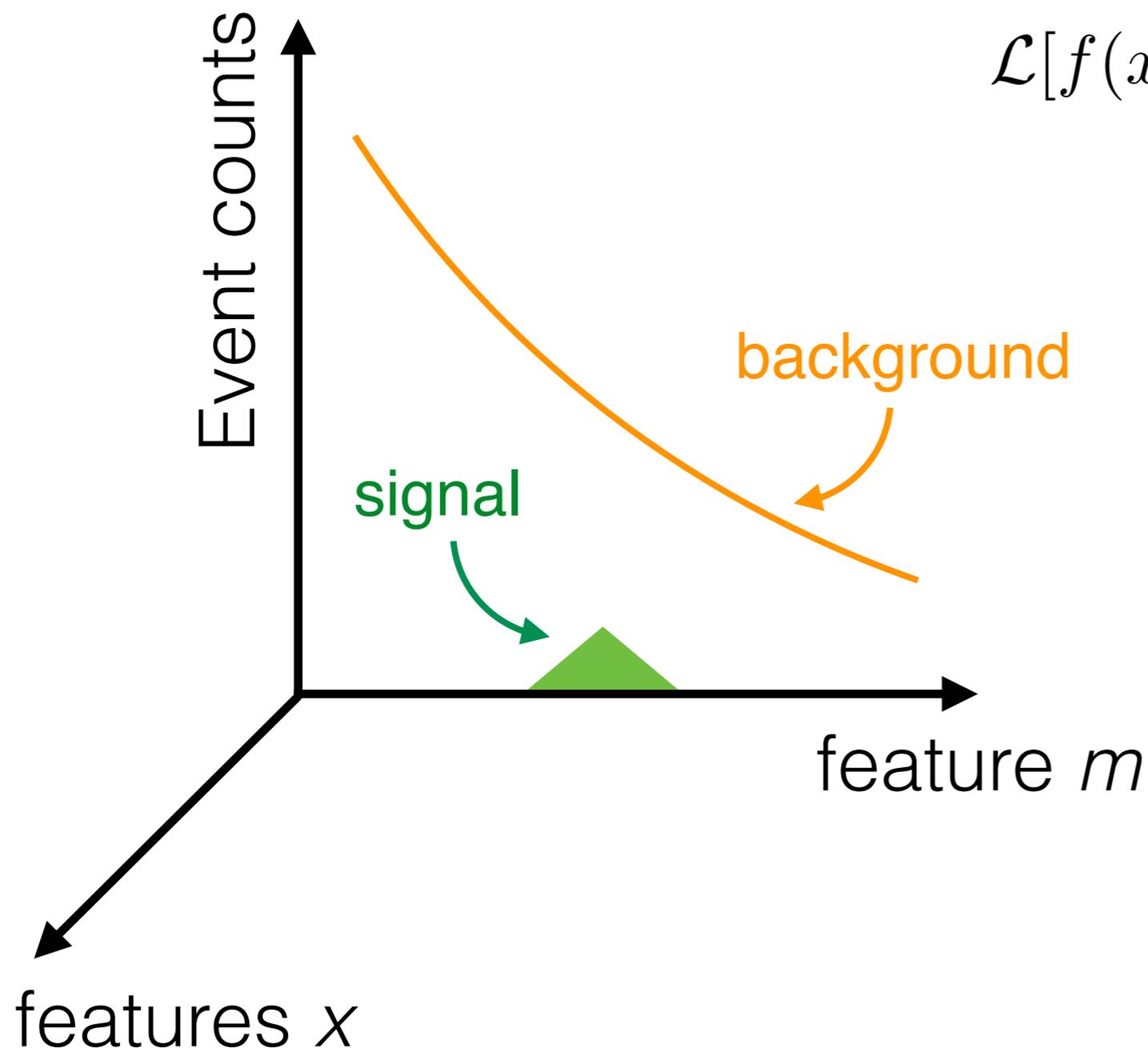
Train e.g. a neural network

$$\mathcal{L}[f(x)] = \textcolor{green}{\sum_{i \in s} L_{\text{classifier}}(f(x_i), 1)}$$
$$\textcolor{orange}{+ \sum_{i \in b} L_{\text{classifier}}(f(x_i), 0)}$$

Event counts

background

signal

feature *m*

features *x*

$L_{classifier}$ is the usual classifier loss, e.g. cross entropy or mean squared error.

Train e.g. a neural network with a **custom loss functional**

$$\mathcal{L}[f(x)] = \textcolor{green}{\sum_{i \in s} L_{\text{classifier}}(f(x_i), 1)}$$
$$+ \textcolor{orange}{\sum_{i \in b} L_{\text{classifier}}(f(x_i), 0)}$$
$$+ \lambda \sum_{i \in b} L_{\text{decor}}(f(x_i), m_i)$$

$L_{classifier}$ is the usual classifier loss, e.g. cross entropy or mean squared error.

$L_{decor}$ is large when $f(x)$ and $m$ are "correlated"

Event counts

background

signal

feature $m$

features $x$

Train e.g. a neural network with a **custom loss functional**

$$\mathcal{L}[f(x)] = \textcolor{green}{\sum_{i \in s} L_{\text{classifier}}(f(x_i), 1)} + \textcolor{orange}{\sum_{i \in b} L_{\text{classifier}}(f(x_i), 0)}$$
$$+ \lambda \sum_{i \in b} L_{\text{decor}}(f(x_i), m_i)$$

*Recent proposals:*

**Adversaries**: $L_{decor}$ is the loss of **a 2nd NN** (adversary) that tries to learn *m* from *f(x)*.
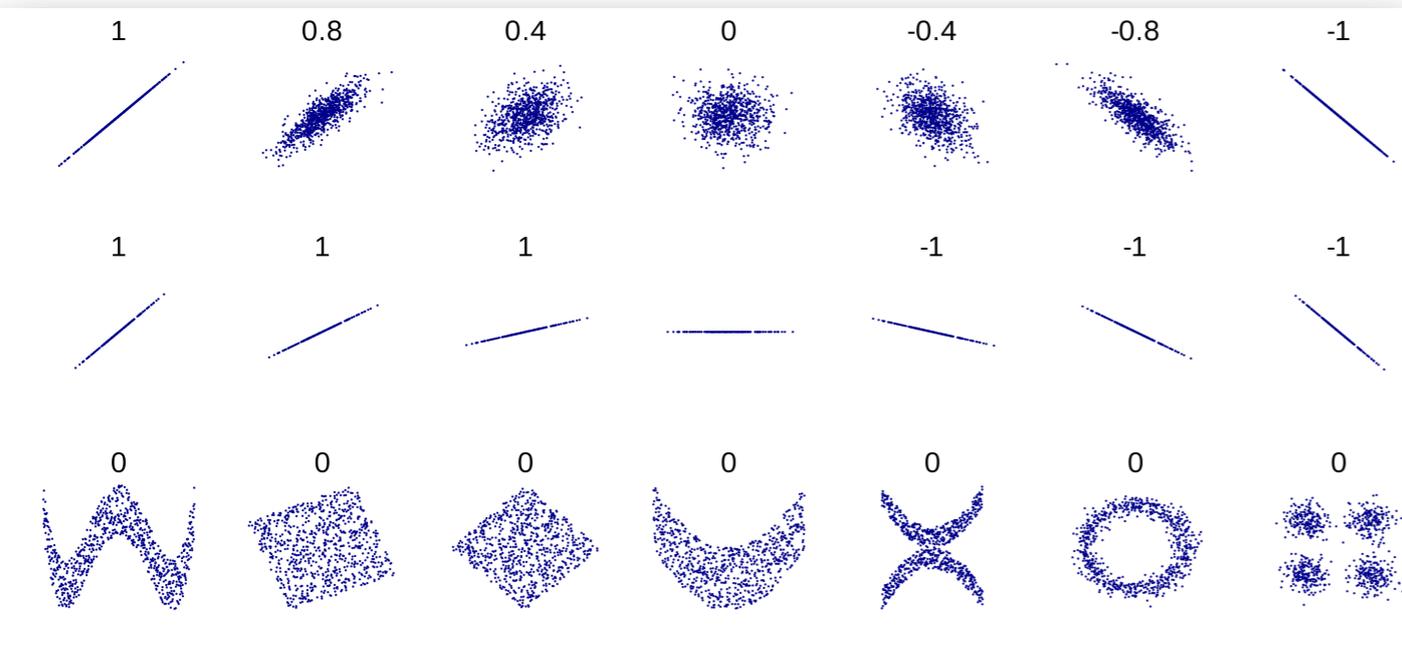
**Distance Correlation**: $L_{decor}$ is **distance correlation** (generalizes Pearson correlation) between *m* and *f(x)*.

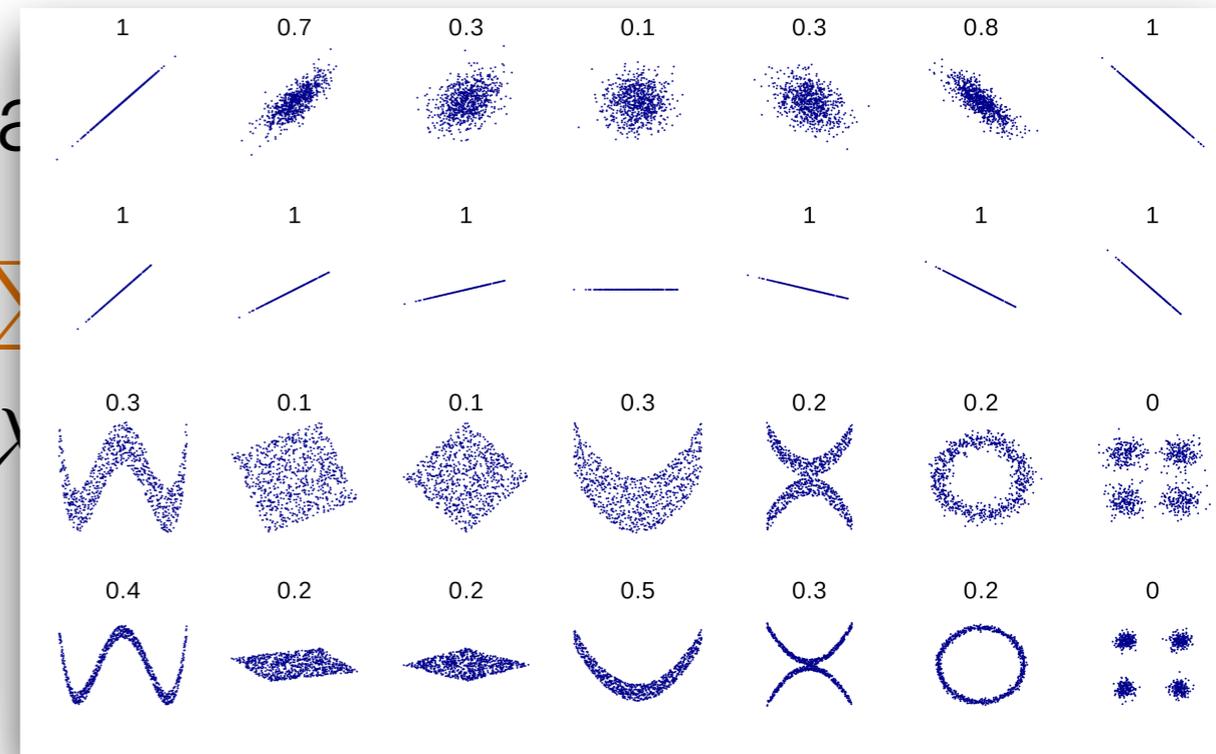**Mode Decorrelation**: $L_{decor}$ is small when the **CDF** of *f(x)* is the same across different values of *m*.

Train e.g. a neural network with a **custom loss functional**

$$\mathcal{L}[f(x)] = \sum_{i \in s} L_{\text{classifier}}(f(x_i), 1) + \sum_{i \in b} L_{\text{classifier}}(f(x_i), 0)$$
$$+ \lambda \sum_{i \in b} L_{\text{decor}}(f(x_i), m_i)$$

*Recent proposals:*

**Adversaries**: $L_{decor}$ is the loss of **a 2nd NN** (adversary) that tries to learn *m* from *f(x)*.

**Distance Correlation**: $L_{decor}$ is **distance correlation** (generalizes Pearson correlation) between *m* and *f(x)*.

**Mode Decorrelation**: $L_{decor}$ is small when the **CDF** of *f(x)* is the same across different values of *m*.

Train e.g. a neural network with a **custom loss functional**

$$\mathcal{L}[f(x)] = \textcolor{green}{\sum_{i \in s} L_{\text{classifier}}(f(x_i), 1)} + \textcolor{orange}{\sum_{i \in b} L_{\text{classifier}}(f(x_i), 0)}$$
$$+ \lambda \sum_{i \in b} L_{\text{decor}}(f(x_i), m_i)$$

*Recent proposals:*

**Adversaries**: $L_{decor}$ is the loss of **a 2nd NN** (adversary) that tries to learn *m* from *f(x)*.

**Distance Correlation**: $L_{decor}$ is **distance correlation** (generalizes Pearson correlation) between *m* and *f(x)*.

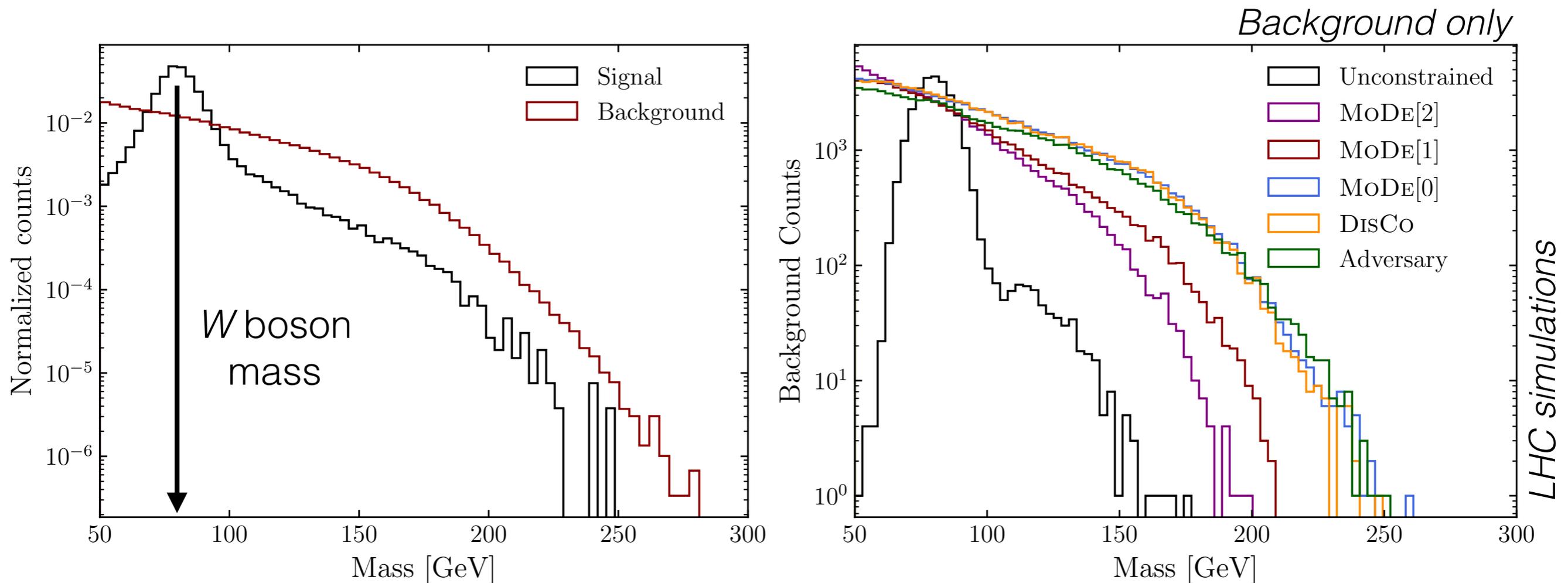**Mode Decorrelation**: $L_{decor}$ is small when the **CDF** of *f(x)* is the same across different values of *m*.

**Pearson Correlation**

**Distance Correlation**



Image credit: Denis Boigelot

**Adversaries**: $L_{decor}$ is the loss of **a 2ⁿᵈ NN** (adversary) that tries to learn *m* from *f(x)*.

**Distance Correlation**: $L_{decor}$ is **distance correlation** (generalizes Pearson correlation) between *m* and *f(x)*.

**Mode Decorrelation**: $L_{decor}$ is small when the **CDF** of *f(x)* is the same across different values of *m*.

Train e.g. a neural network with a **custom loss functional**

$$\mathcal{L}[f(x)] = \sum_{i \in s} L_{\text{classifier}}(f(x_i), 1) + \sum_{i \in b} L_{\text{classifier}}(f(x_i), 0)$$
$$+ \lambda \sum_{i \in b} L_{\text{decor}}(f(x_i), m_i)$$

*Recent proposals:*

**Adversaries**: $L_{decor}$ is the loss of **a 2nd NN** (adversary) that tries to learn *m* from *f(x).*

**Distance Correlation**: $L_{decor}$ is **distance correlation** (generalizes Pearson correlation) between *m* and *f(x).*

**Mode Decorrelation**: $L_{decor}$ is small when the **CDF** of *f(x)* is the same across different values of *m.*

**Adversaries**: $L_{decor}$ is the loss of **a 2nd NN** (adversary) that tries to learn $m$ from $f(x)$.

**Pros:** Very flexible and $m$ can be multidimensional

**Cons:** Hard to train (minimax problem) & many parameters

**Distance Correlation**: $L_{decor}$ is **distance correlation** (generalizes Pearson correlation) between $m$ and $f(x)$.

**Pros:** Convex (easier to train) and no free parameters

**Cons:** Memory intensive to compute distance correlation

**Mode Decorrelation (MoDE)**: $L_{decor}$ is small when the **CDF** of *f(x)* is the same across different values of *m.*

**Pros:**  Readily generalizes beyond independence (can require linear, quadratic (+monotonic), …

No free parameters and small memory footprint

**Cons:**  In its simplest form, need discrete bins in *m* (does not seem to be fundamental)

## Real world example: the search for Lorentz-boosted *W* bosons at the Large Hadron Collider



*W* boson mass

Signal
Background

*Background only*

Unconstrained
MoDe[2]
MoDe[1]
MoDe[0]
DisCo
Adversary

*LHC simulations*

MoDe[0] enforces independence, [1] is linear, [2] is monotonic quadratic, …

*N.B. think twice about using decorrelation for uncertainties! See 2109.08159.*

Sometimes, we need a model (often for calibration) that does not depend on the training sample properties.

For example, a particle of a given energy hits our detector and registers measurements in a number of sensors

e.g. the particle energy is uniform during training, but exponential for certain running conditions.

(usually not an issue for classification)

Sometimes, we need a model (often for calibration) that does not depend on the training sample properties.

For example, a particle of a given energy hits our detector and registers measurements in a number of sensors

e. ng,
s.

*Your first instinct here might have been to train a classifier to estimate the true value given measured values using simulated data.*

**Claim: this is prior dependent !**

For example, a particle of a given energy hits our detector and registers measurements in a number of sensors

e.                                    ng,
                                      s.

*Your first instinct here might have been to train a classifier to estimate the true value given measured values using simulated data.*

Suppose you have some features <span style="color:blue">x</span> and you want to predict <span style="color:red">y</span>.

*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].

*If you have not seen this before, please let me know if you need help with the proof!

Suppose you have some features <span style="color:blue">x</span> and you want to predict <span style="color:red">y</span>.

*detector energy*　　　　*true energy*

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

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

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

*Simulation-Based Gaussian Example*

$(\mu, \sigma, \varepsilon) = (0, 1, 2)$

MSE Fit

*Simulation-Based Gaussian Example*

$(\mu, \sigma, \varepsilon) = (0, 1, 2)$

Analytic MSE
Analytic MLC
Fit MSE
Gaussian Ansatz

2205.05084

*Dijet Distributions*

# Physics Example



Simulation-Based Dijets Example

Simulation-Based Dijets Example

2205.05084; see also ATL-PHYS-PUB-2018-013

Can we train a neural network to emulate the detector simulation?

Grayscale images: Pixel intensity = energy deposited

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



Deep generative models: the map is a deep neural network.

**See also Sascha's talk!**

**GANs**
*Generative Adversarial Networks*

**Score-based**

**NFs**
*Normalizing Flows*

**VAEs**
*Variational Autoencoders*

Deep generative models: the map is a deep neural network.

Generative Adversarial Networks (GANs):
*A two-network game where one **maps noise to structure** 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 or data

Variational Autoencoders (VAEs):
*A pair of networks that embed the data into a latent space with a given prior and decode back to the data space.*

latent space

E

D

$p(z|x)$

$p(x|z)$

Physics-based simulator or data

*Probabilistic* **encoder**

*Probabilistic* **decoder**

Normalizing Flows (NFs):
*A series of invertible transformations mapping a known density into the data density.*

Optimize via
maximum likelihood



latent
space

*Invertible transformations
with tractable Jacobians*

p(z)

p(x) = p(z) |dF⁻¹/dx|

p(x)

Score-based

*Learn the gradient of the density instead of the probability density itself.*

**Forward diffusion (training)**



**Reverse-time diffusion (data generation)**

From 2206.11898

1705.02355

**GAN**

2102.12491

**VAE**

2110.11377

**NF**

Many papers on this subject - see the living review for all



**Score**

2206.11898

See also https://calochallenge.github.io/homepage/ and http...

One image per calo layer

One network per particle type; input particle energy
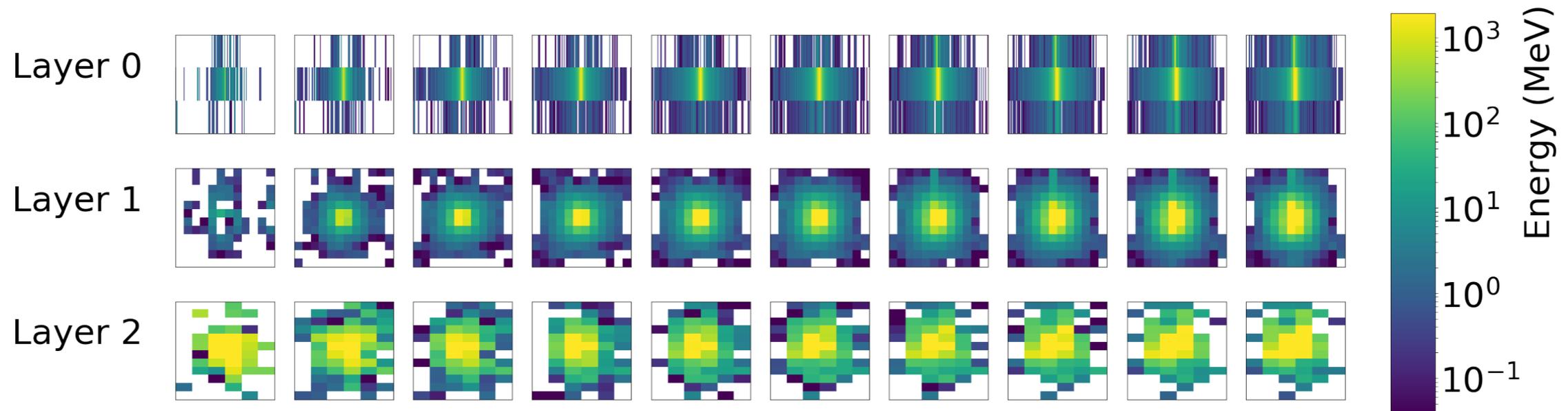
**Generator network**

INPUTS

OUTPUTS

1705.02355

LA = Locally Aware, like a CNN

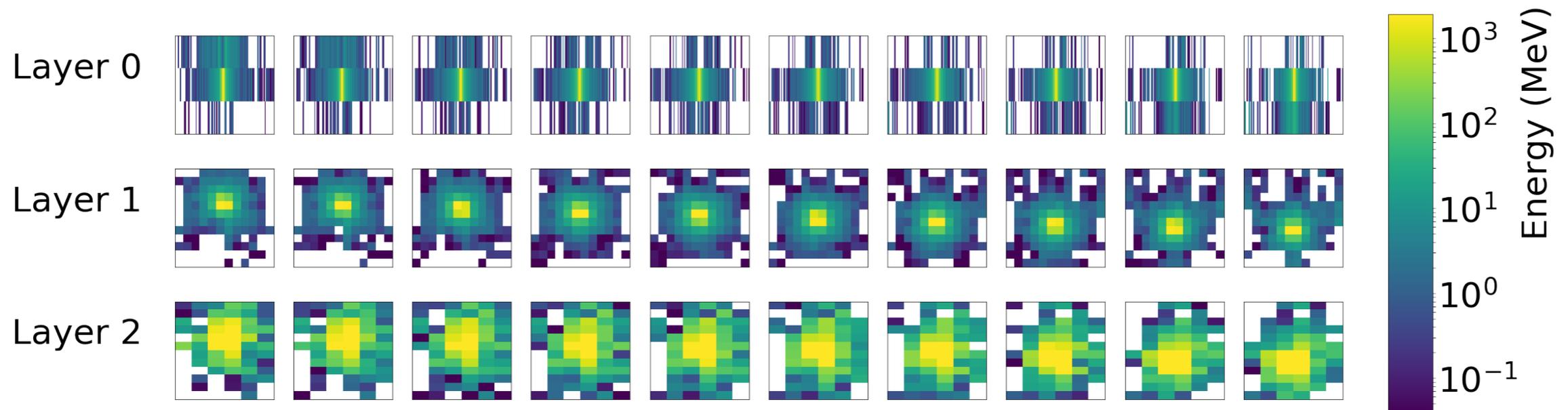use layer i as input to layer i+1

ReLU to encourage sparsity

# Conditioning

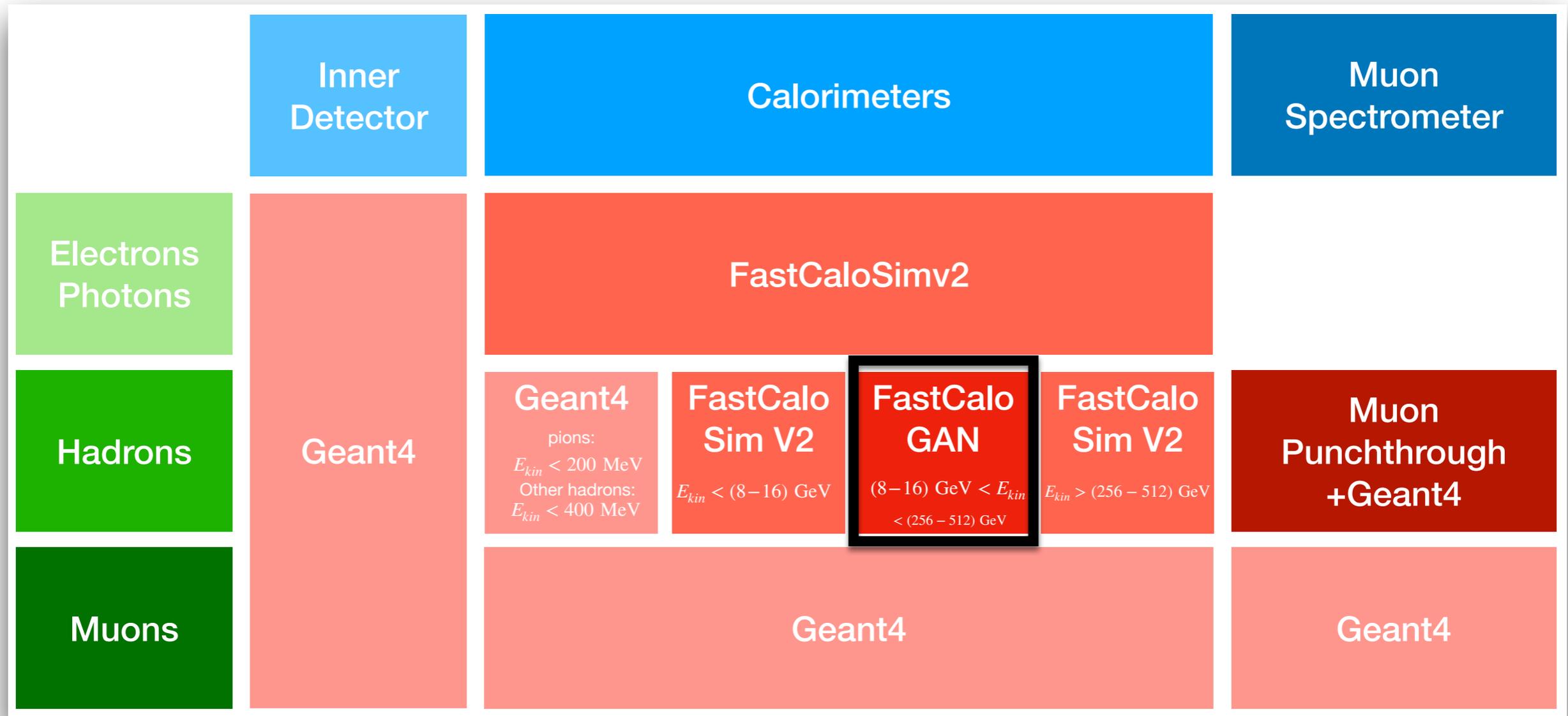Fix noise, scan latent variable corresponding to energy



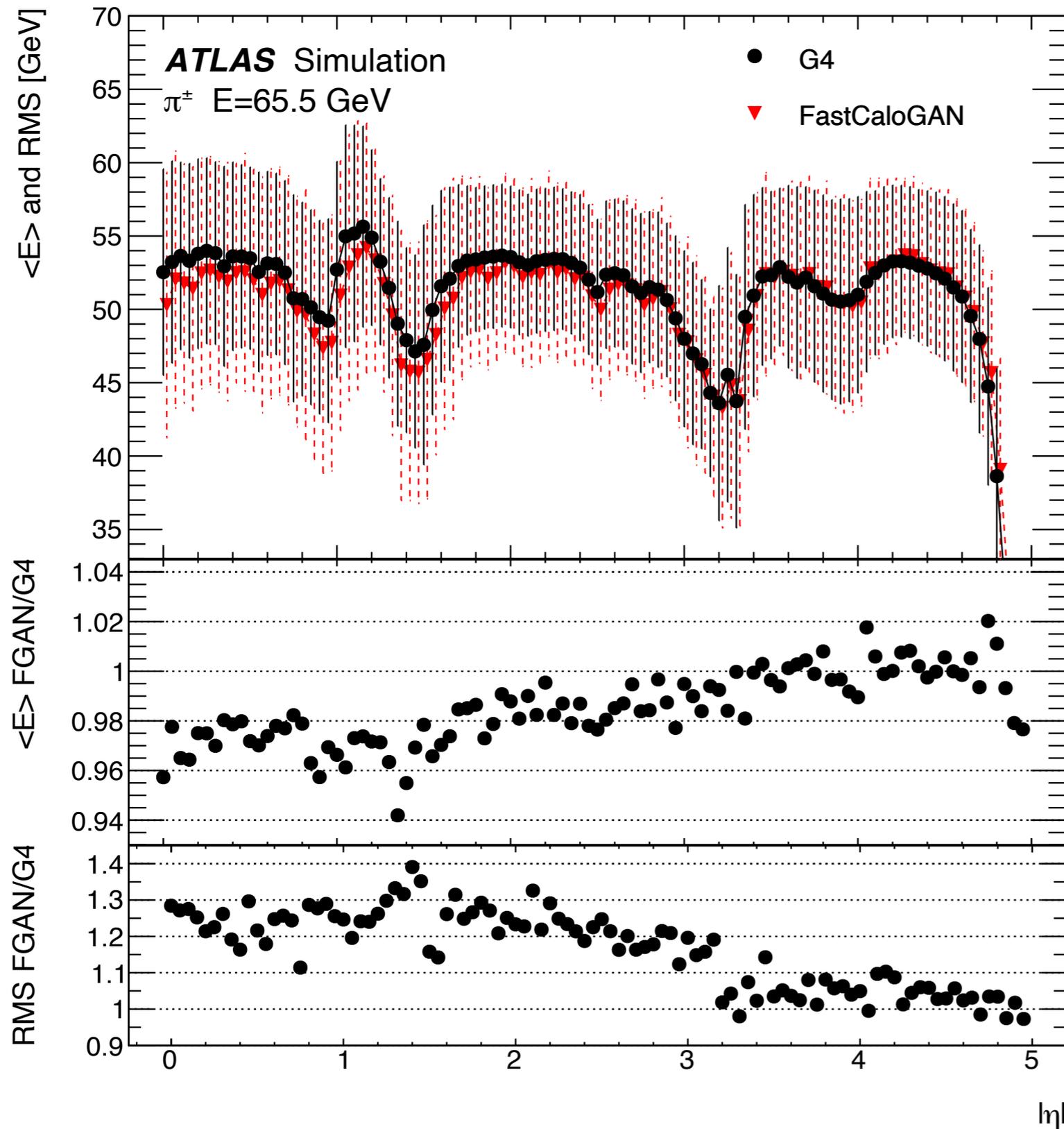Fix noise, scan latent variable corresponding to x-position

# Integration into real detector sim.

Our (ATLAS Collaboration) fast simulation (AF3) now includes a GAN at intermediate energies for pions
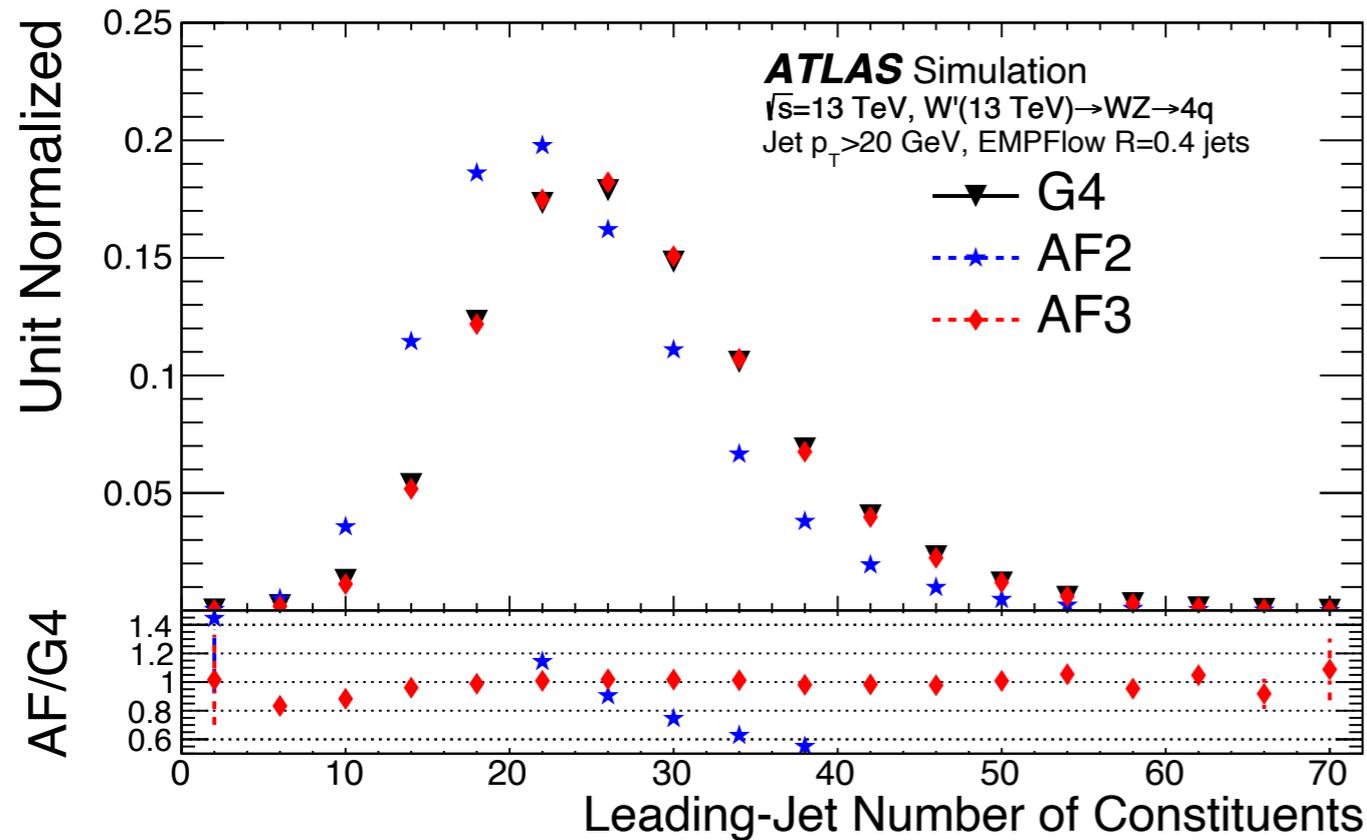
# Integration into real detector sim.

The GAN architecture is relatively simple, but it is able to match the energy scale and resolution well.
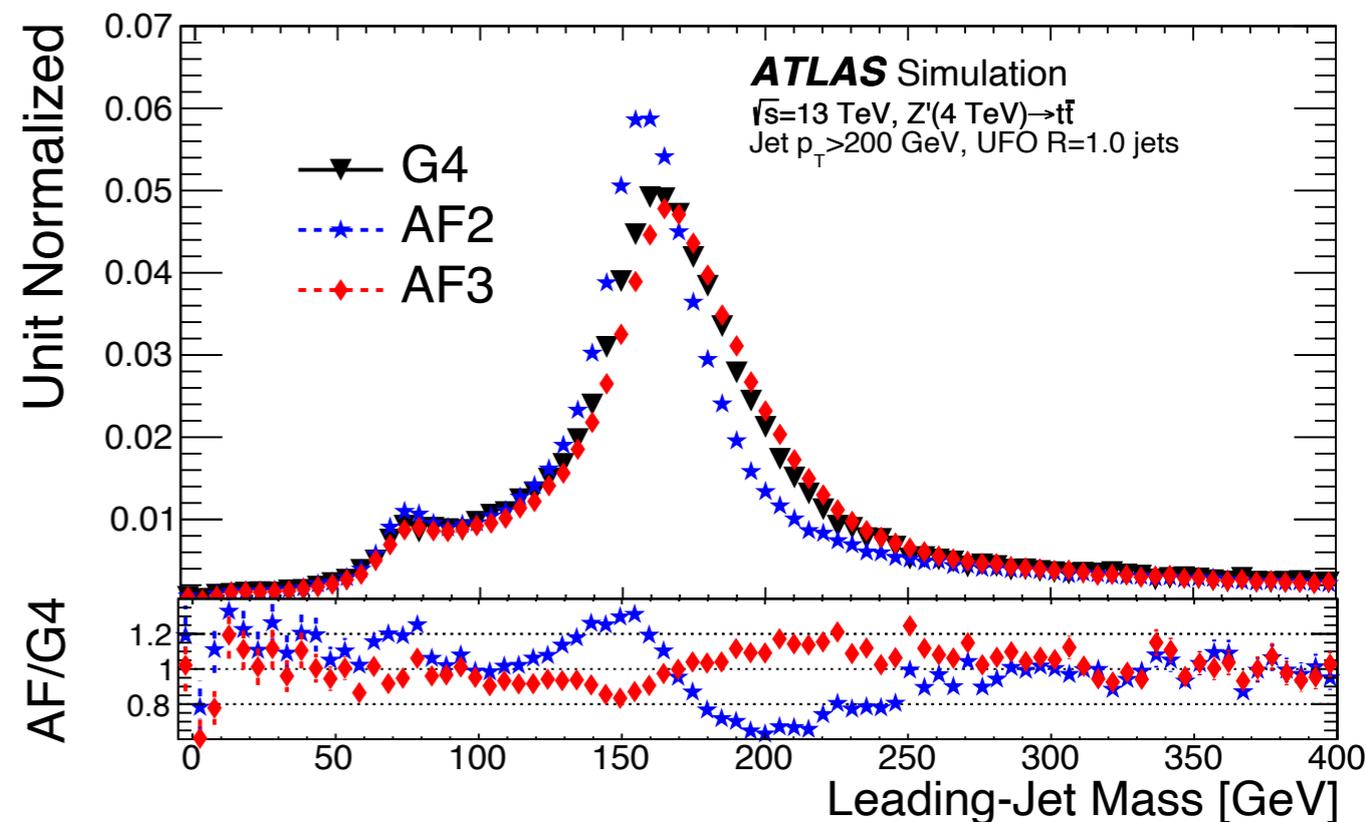
There is one GAN per η slice
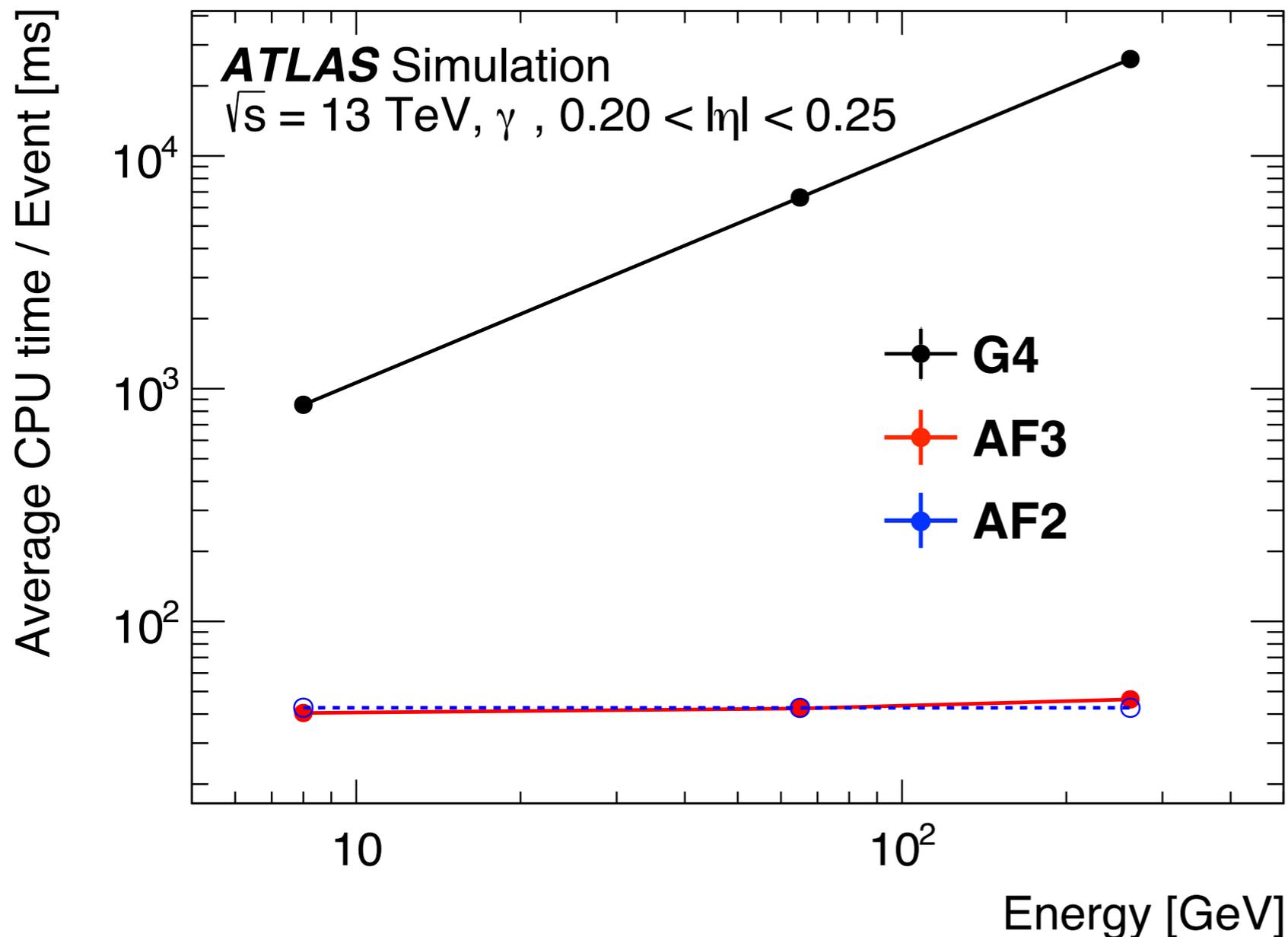
# Integration into real detector sim.

The new fast simulation (**AF3**) significantly improves jet substructure with respect to the older one (**AF2**)

Ideally, the same calibrations derived for full sim. (Geant4-based) can be applied to the fast sim.

# Integration into real detector sim.

As expected, the fast sim. timing is independent of energy, while Geant4 requires more time for higher energy.

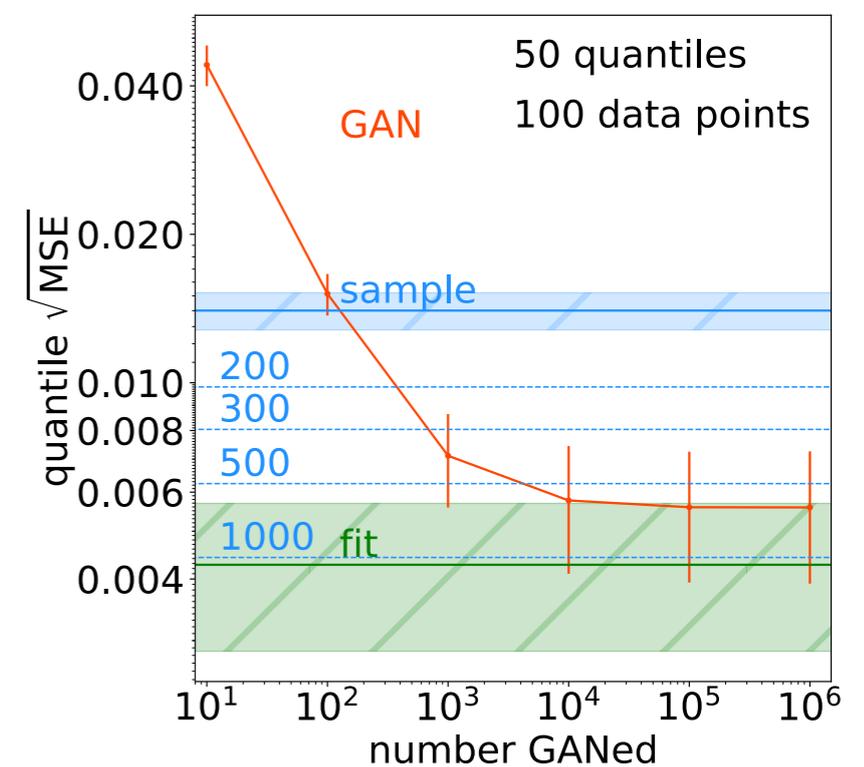Common question: if we train on N events and sample M >> N events, do we have the statical power of M or N?
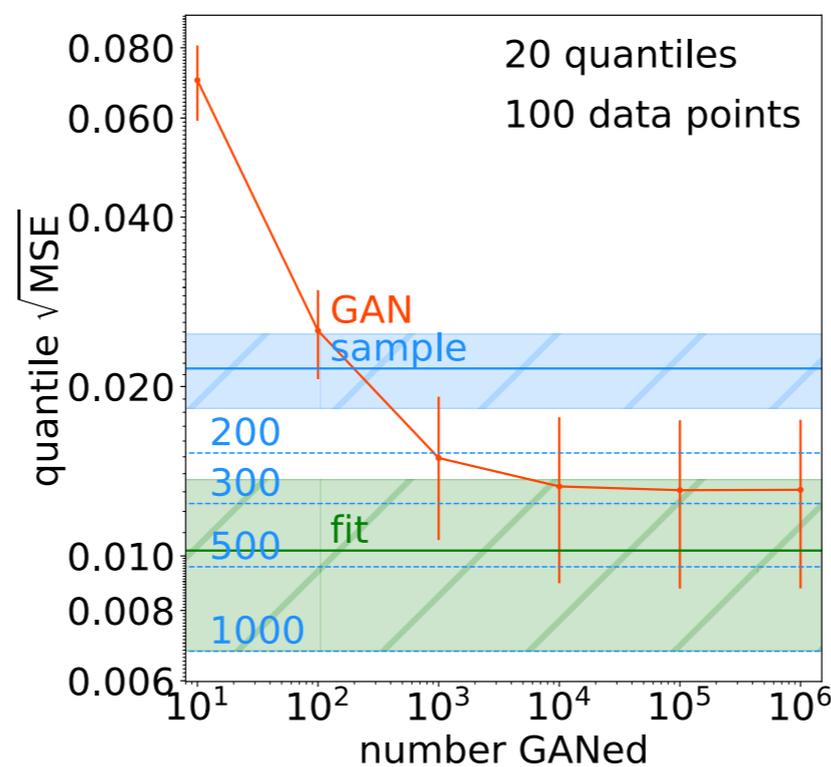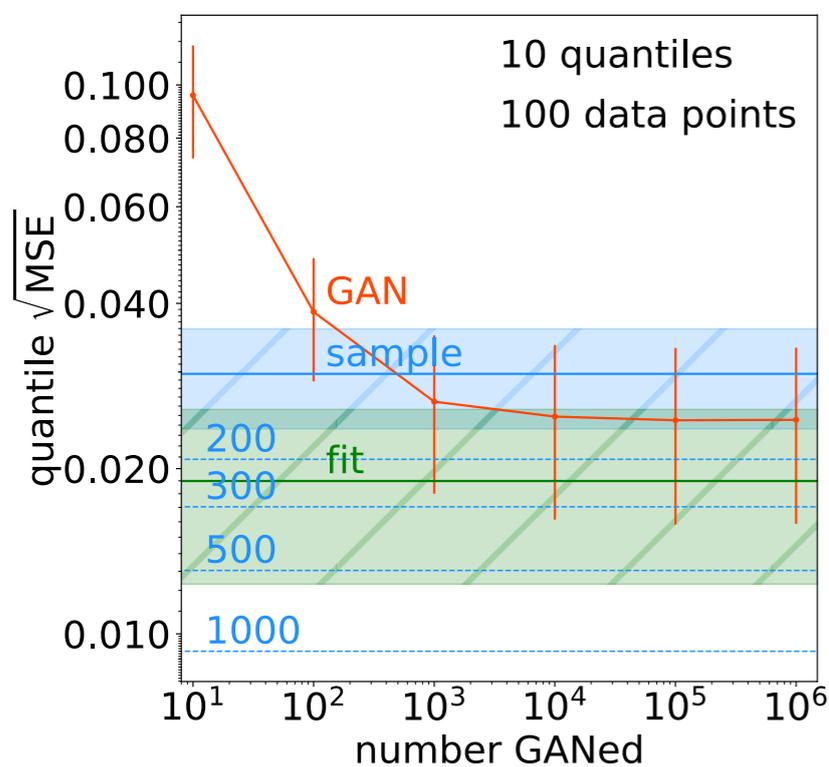
No free lunch - only win with **inductive bias**.  Examples: factorization, symmetries, smoothness, …

Common question: if we train on N events and sample M >> N events, do we have the statical power of M or N?

No free lunch - only win with **inductive bias**.  Examples: factorization, symmetries, smoothness, …



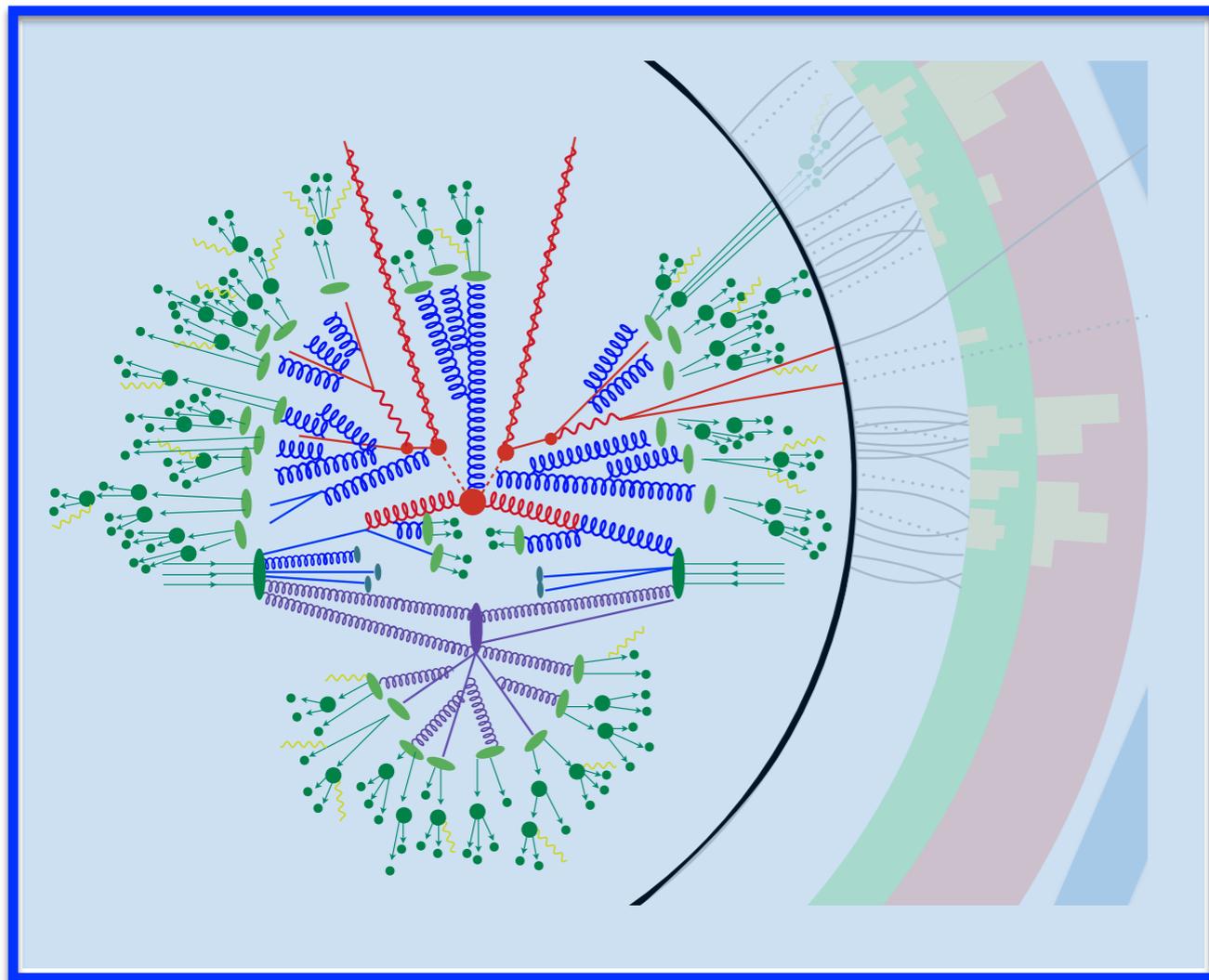2006.06685; 2202.07352

Theory of everything

**Physics simulators**

Parameter estimation / unfolding

Detector-level observables

Pattern recognition

**(III)**

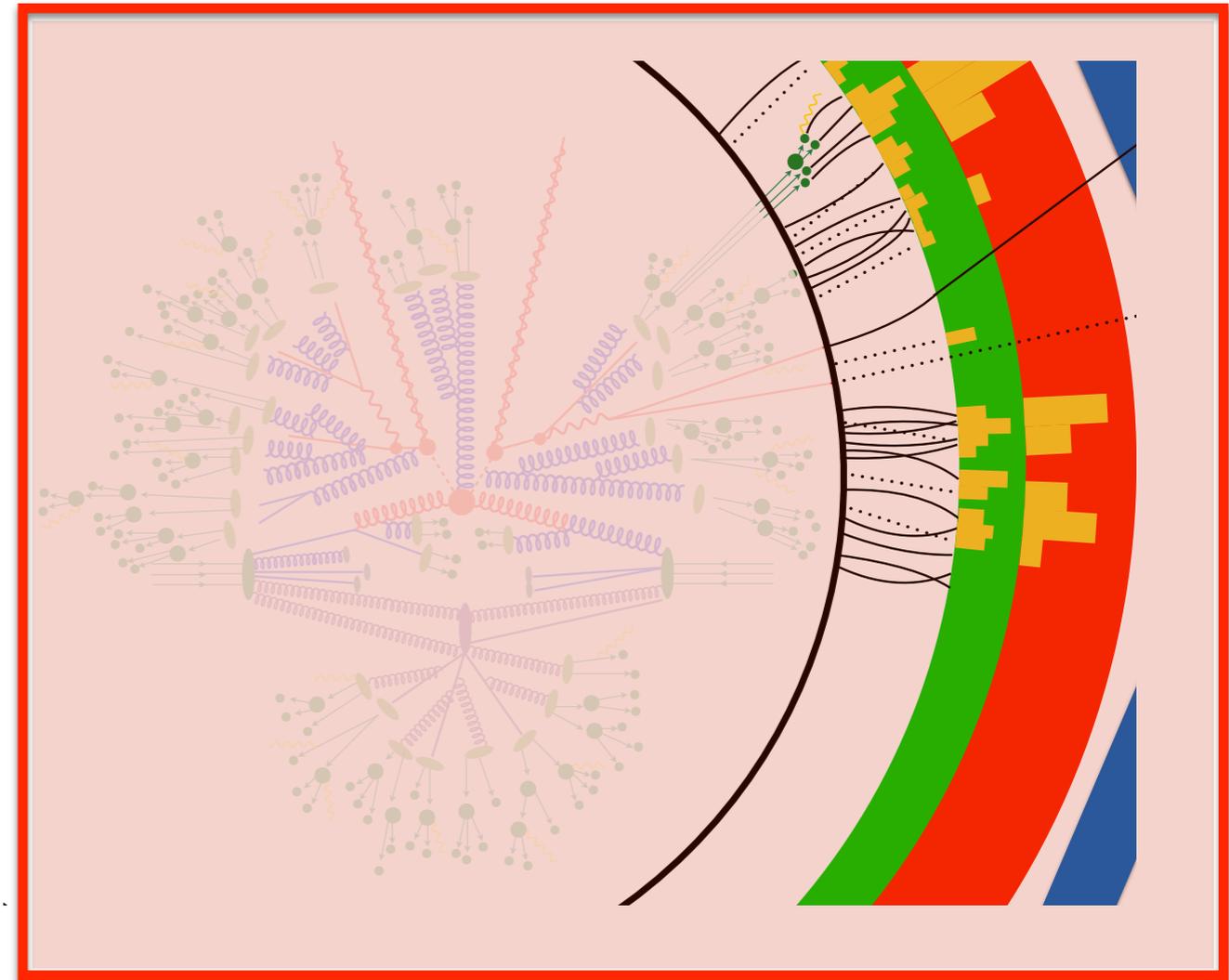**Nature**

Experiment

Detector-level observables

Pattern recognition

**Want this**

(or the parameters of the generative model)

**Measure this**



remove detector distortions (unfolding) or parameter estimation

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

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



**Want this**

**Measure this**

*For parameter estimation, replace **true** with θ*

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** !*

*For parameter estimation, replace **true** with $\theta$*

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

$$unfolded = \underset{true}{argmax}\ p(measured\ |\ true)$$

⚠ 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**

*For parameter estimation, replace **true** with $\theta$*

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 $= q(x)/p(x)$

(or something monotonically related to it in a known way)

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

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

*This is easy to prove. If you have not seen it before, please ask!*

$$L[f] = \sum (f(x_i) - c)^2$$

*Try yourself with BCE!*

$$\approx \int dx\, p(x, c)\, (f(x) - c)^2$$

$$\frac{\delta L[f, f']}{\delta f} = \frac{\partial L}{\partial f} - \frac{d}{dx}\frac{\partial L}{\partial f'} = 0$$

*Euler-Lagrange Equation*

$$L[f] = \sum (f(x_i) - c)^2$$

*Try yourself with BCE!*

$$\approx \int dx \, p(x, c) \, (f(x) - c)^2$$

$$\frac{\delta L[f, f']}{\delta f} = \frac{\partial L}{\partial f} - \frac{d}{dx} \frac{\partial L}{\partial f'} = 0$$

*Euler-Lagrange Equation*

Basically just a regular derivative:

$$\int dc \, p(x, c)(f(x) - c) = 0 \implies f(x) = E[c \,|\, x]$$

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

(or something monotonically related to it in a known way)

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

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

*This is easy to prove. If you have not seen it before, please ask!*

# Example

87

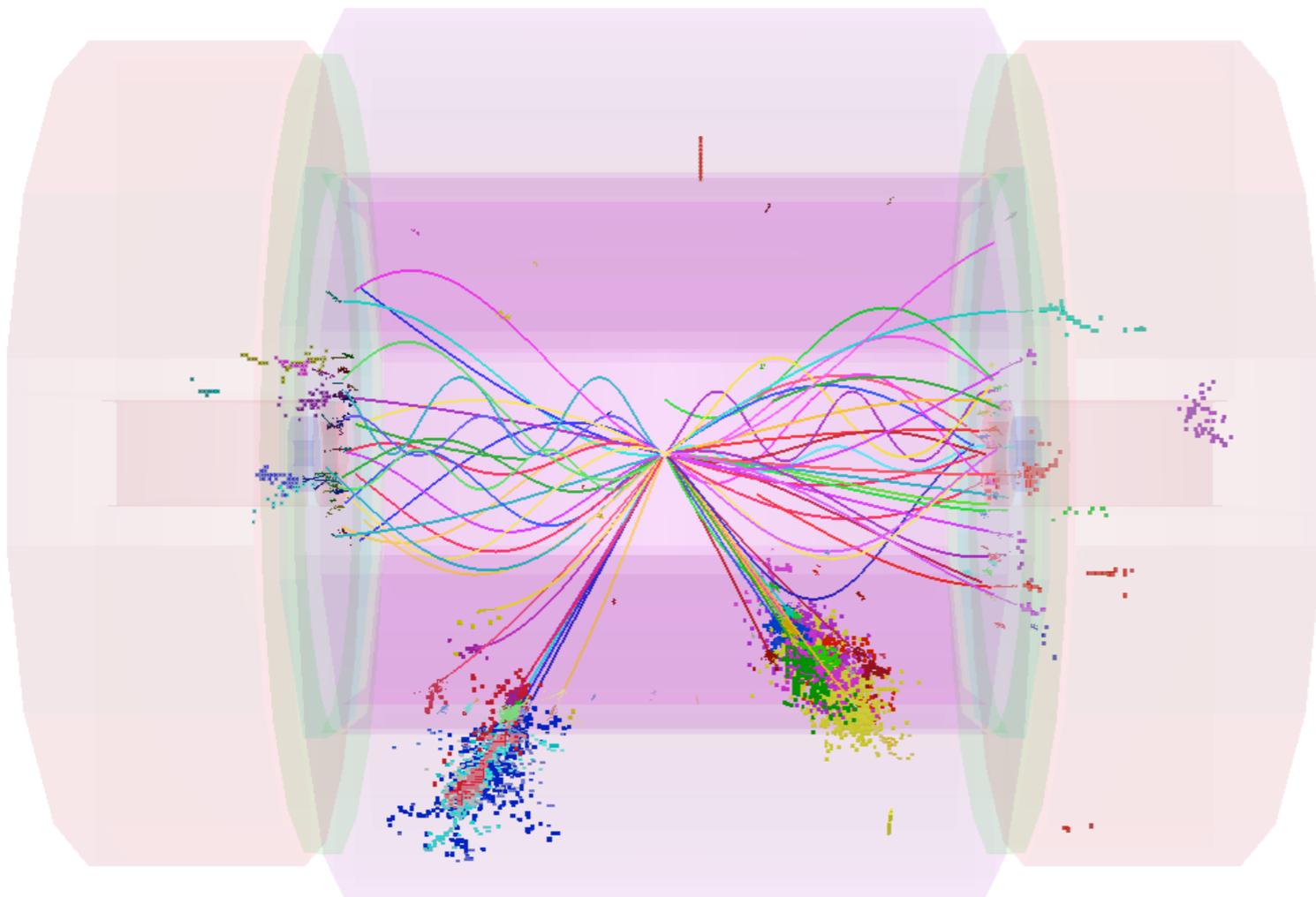Here, instead of emulating $p(x|\theta)$ directly, we learn $\dfrac{p(x|\theta)}{p(x|\theta_0)}$

(turns the problem of generation into classification)

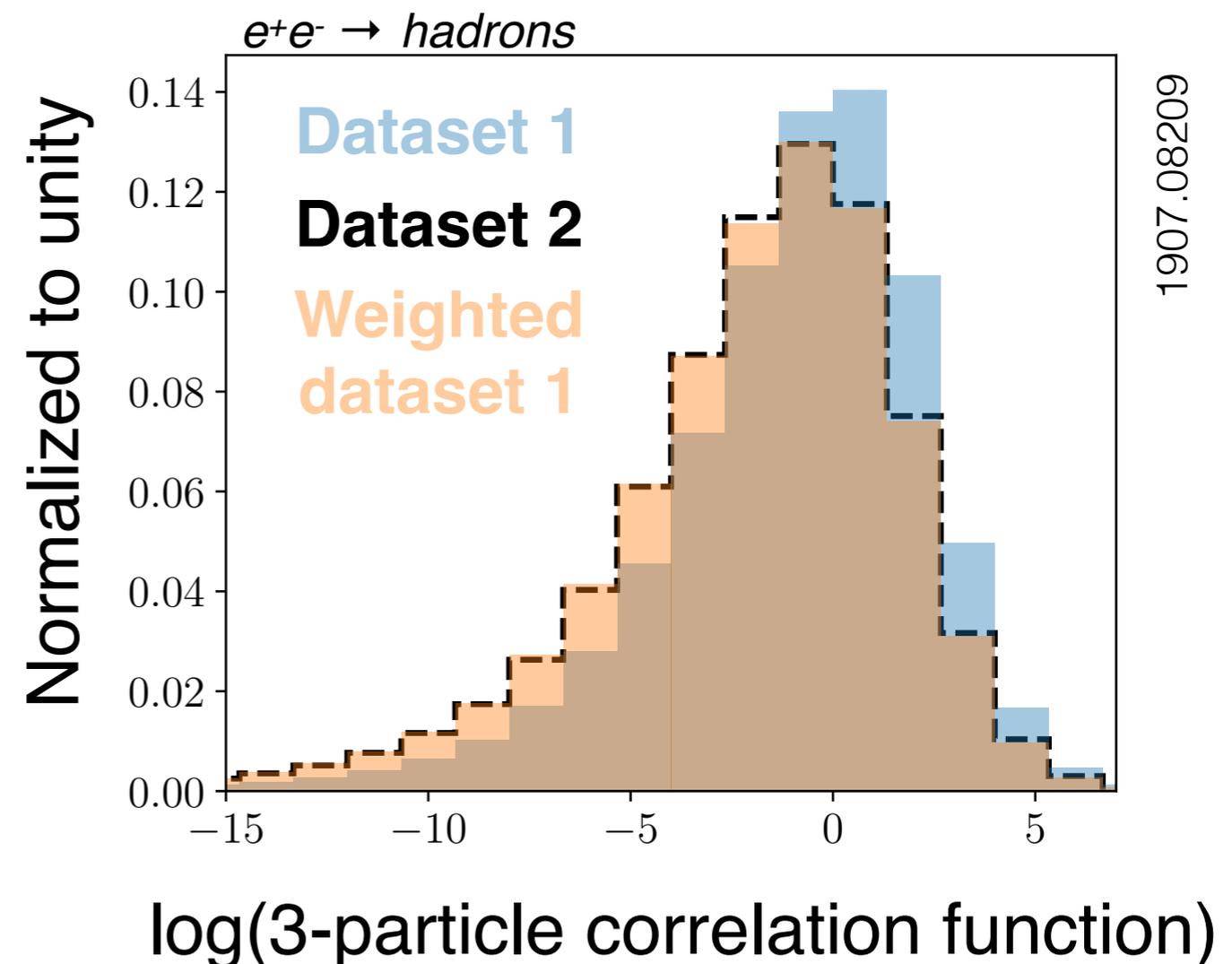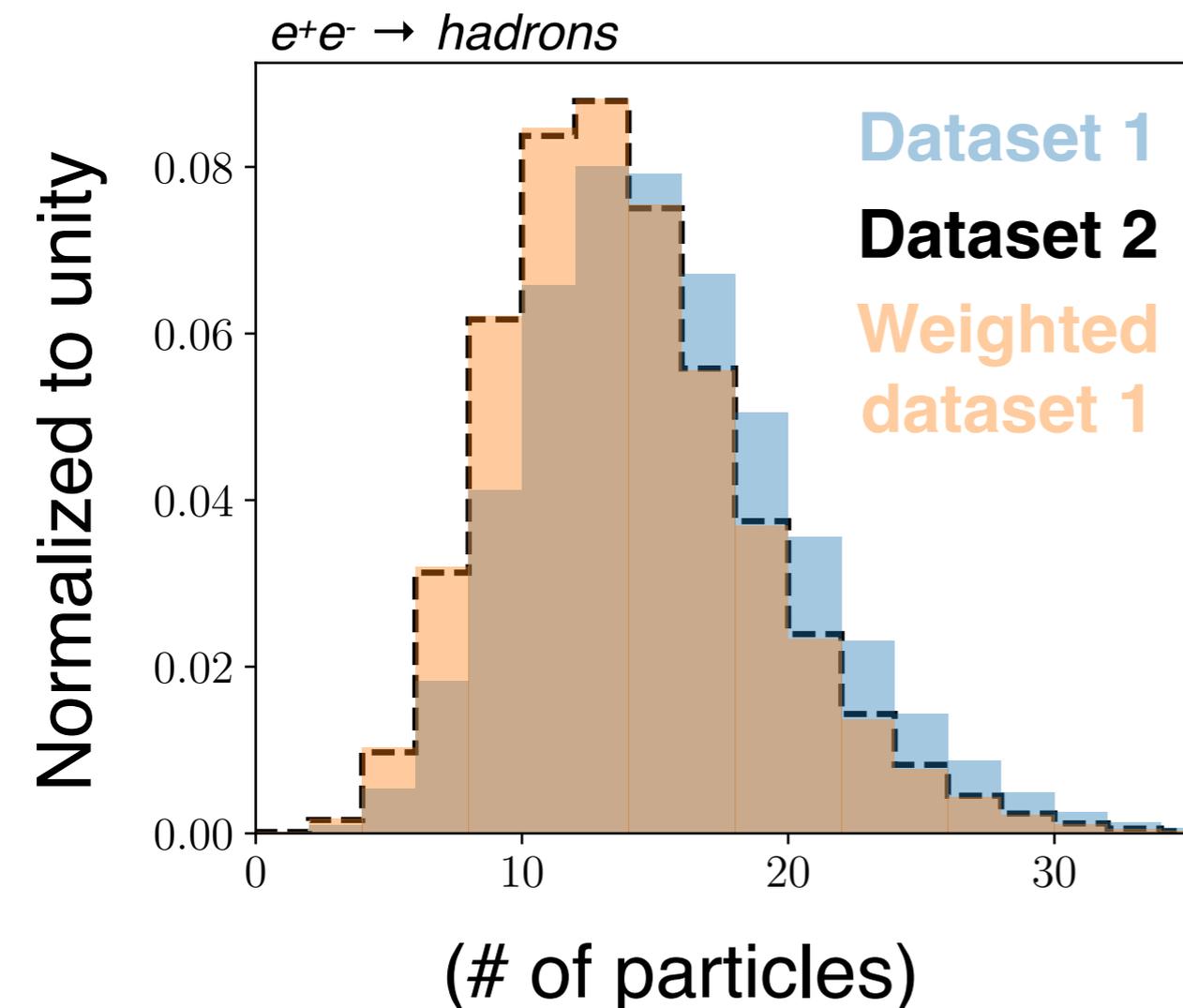Here, instead of emulating $p(x|\theta)$ directly, we learn $\dfrac{p(x|\theta)}{p(x|\theta_0)}$

(turns the problem of generation into classification)

Benefit: easy to integrate complex data structure (symmetries, etc.)

Downside: large weights when $\theta$ is far from $\theta_0$

Image: Linear Collider Detector Project

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

*e⁺e⁻ → hadrons*

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

Normalized to unity

(# of particles)

*e⁺e⁻ → hadrons*

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

Normalized to unity

log(3-particle correlation function)

1907.08209

Works also when the differences between the two simulations are **small** (left) or **localized** (right).

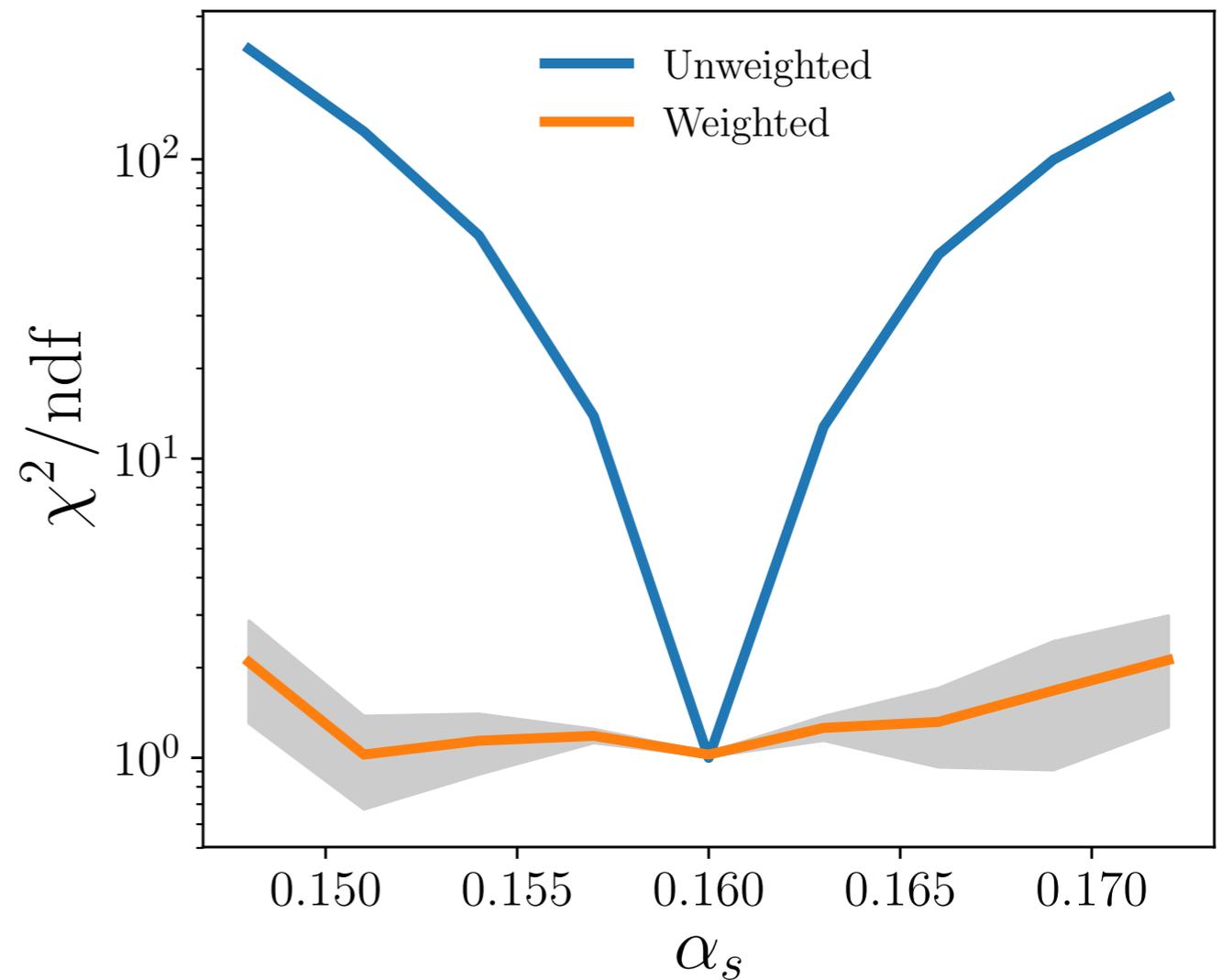*These are histogram ratios for a series of one-dimensional observables*

1907.08209

What if we have a new simulation with multiple continuous parameters θ?

Easy - learn a parameterized classifier* !

…simply add the parameter as a feature to the network during training and let it learn to interpolate.



1907.08209

*1506.02169

Step 1: Differentiable Surrogate Model

$$f(x, \theta) = \operatorname*{argmax}_{f'} \sum_{i \in \boldsymbol{\theta_o}} \log f'(x_i, \theta) + \sum_{i \in \boldsymbol{\theta}} \log(1 - f'(x_i, \theta))$$

See also 1805.00020

## Step 1: Differentiable Surrogate Model

$$f(x, \theta) = \operatorname*{argmax}_{f'} \sum_{i \in \boldsymbol{\theta_0}} \log f'(x_i, \theta) + \sum_{i \in \boldsymbol{\theta}} \log(1 - f'(x_i, \theta))$$

## Step 2: Gradient-based optimization

$$\theta^* = \operatorname*{argmax}_{\theta'} \sum_{i \in \boldsymbol{\theta_0}} \log f(x_i, \theta') + \sum_{i \in \boldsymbol{\theta_1}} \log(1 - f(x_i, \theta'))$$

See also 1805.00020

Obs log10(x), e

log10(x)

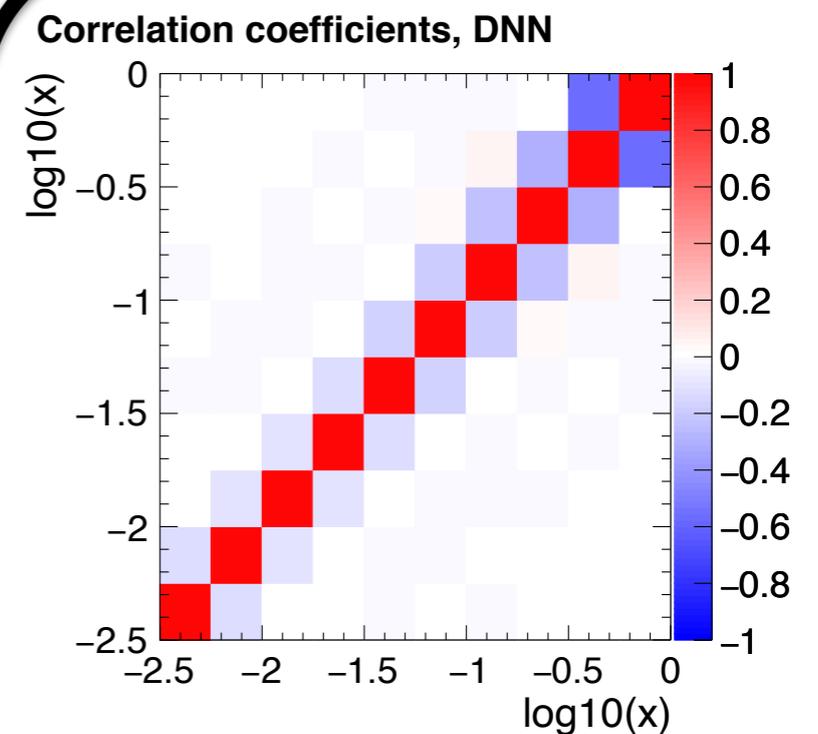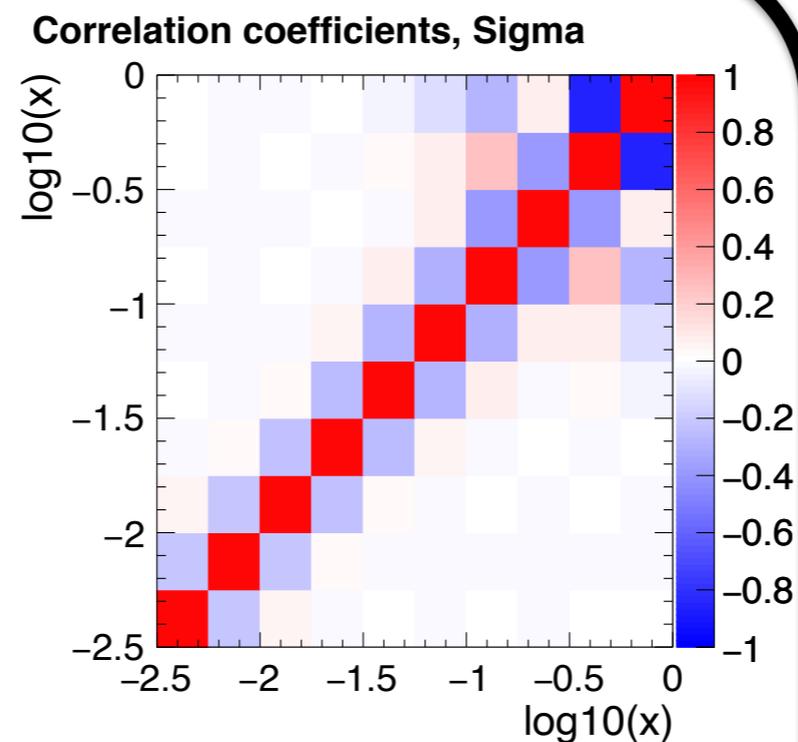log10(x)

**Unfolded distribution, electron**

**Correlation coefficients, electron**

Learn tailored observables; no reason detector level
needs to be same observable as particle level!

**Correlation coefficients, electron**

**Correlation coefficients, Sigma**

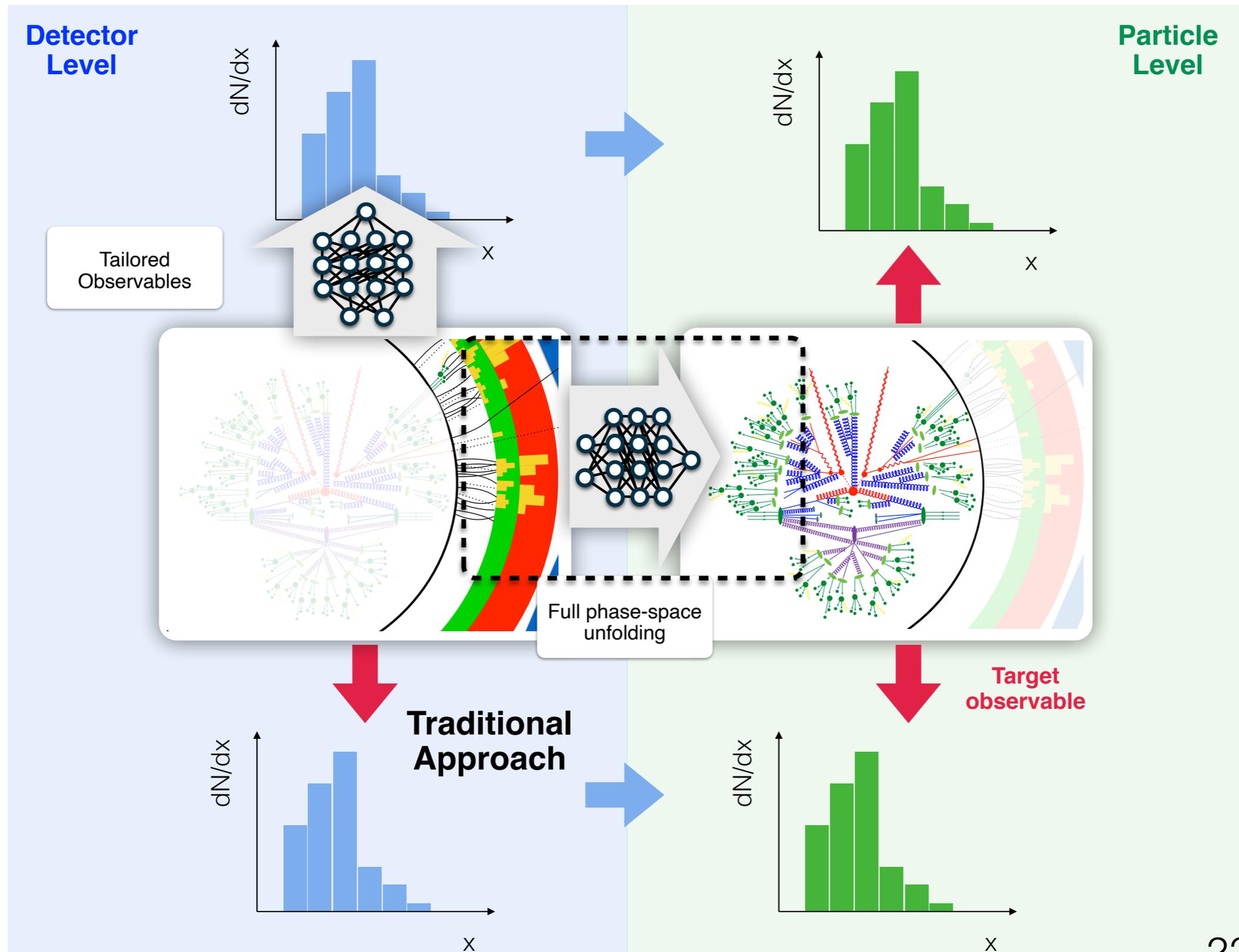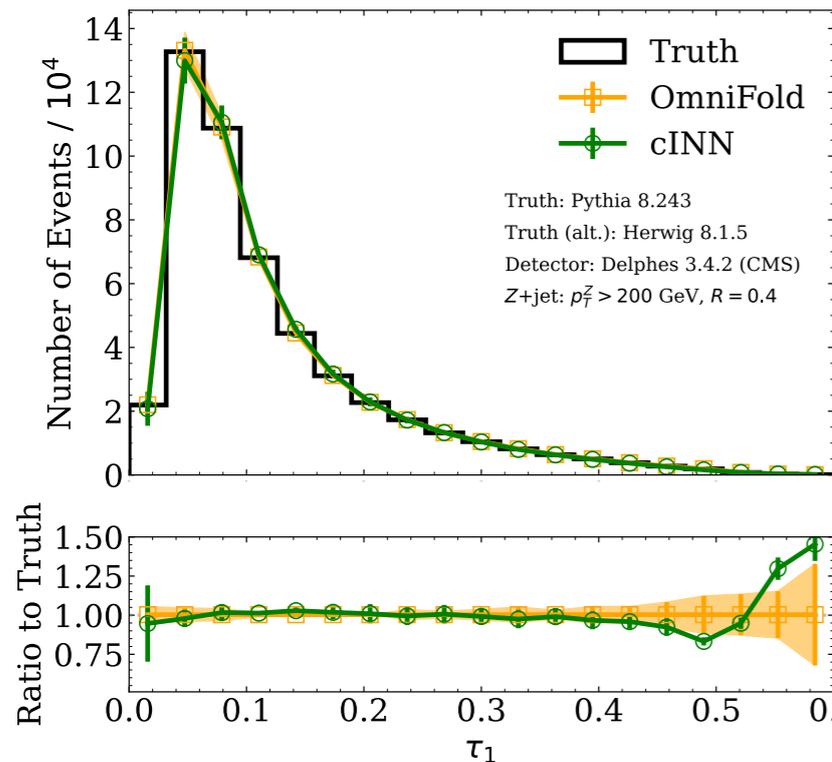**Correlation coefficients, DNN**

**Correlation coefficients, Sigma**

**Correlation coefficients, DNN**

Classical Observables

Neural Network

2203.16722

**Detector Level**

**Particle Level**

Tailored Observables

Full phase-space unfolding

**Traditional Approach**

**Target observable**

2203.16722

ML allows us to do unfolding **unbinned** and in **high dimensions**!
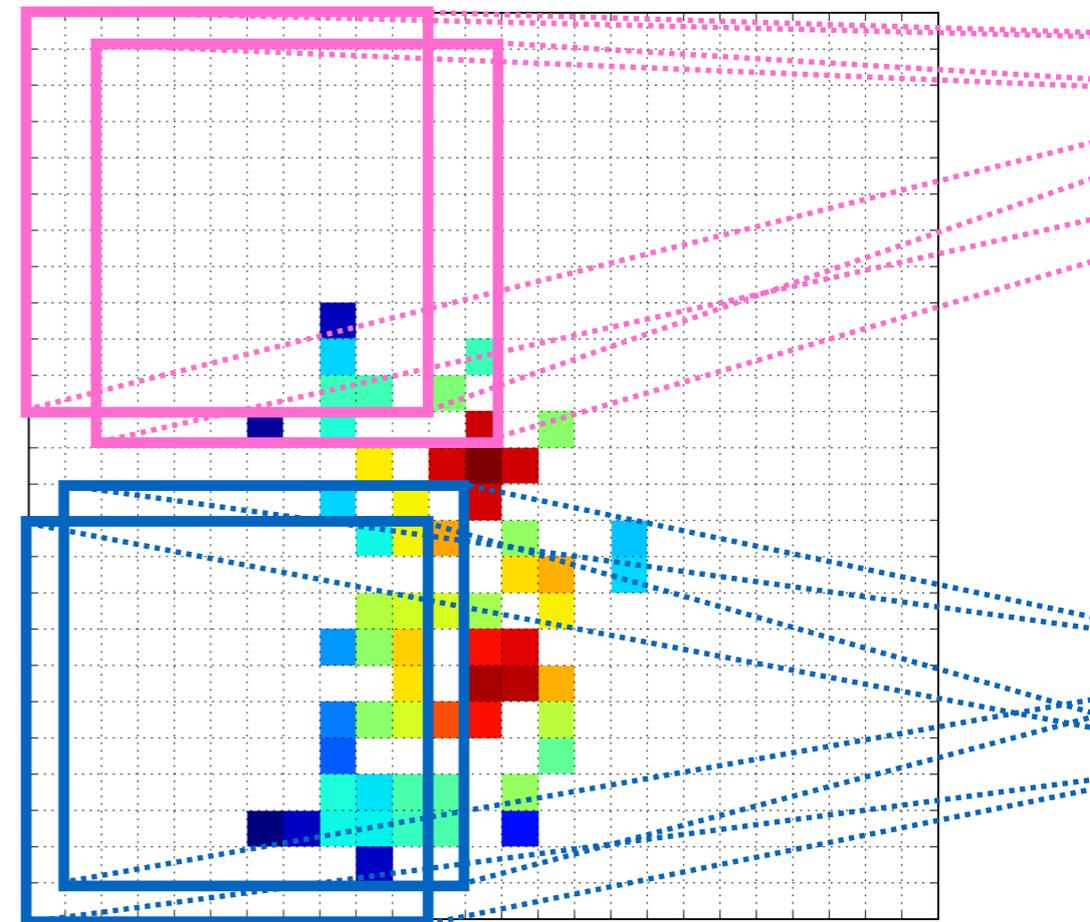
**See Vinny's talk for more!**

# Conclusions and Out

AI/ML has a great potential to **enhance**, **accelerate**, and **empower** all areas of HEP
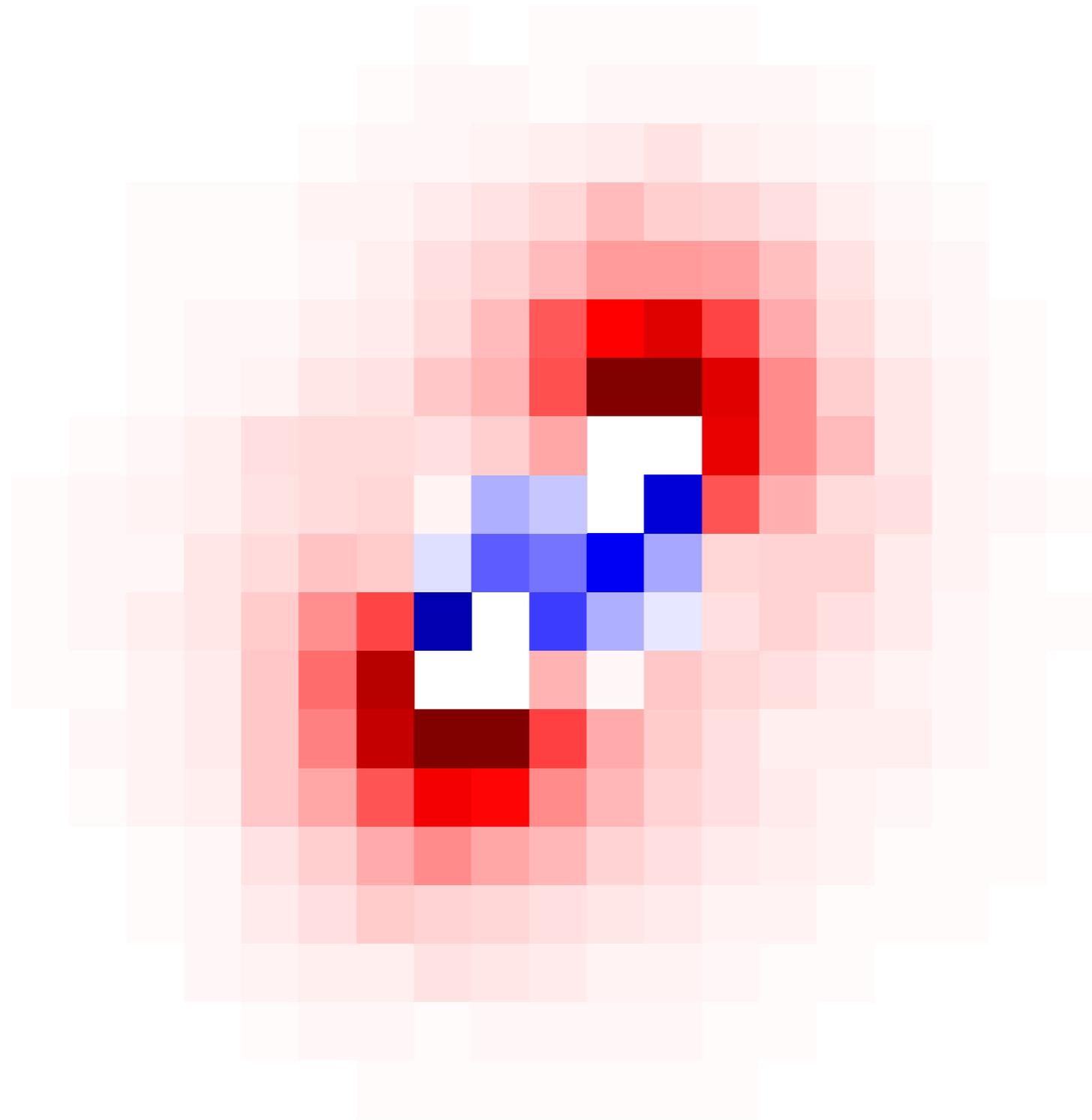
There are applications now that were unthinkable before ML and new ideas are incoming!

*We need you to help develop, adapt, and deploy new methods*

I've provided some specific examples today, but see the Living Review, 2102.02770, for more!

Note that I could not cover everything!  e.g. equivariance



Fin.