# Moments Of Clarity in Machine Learning for Jet Physics

Rikab Gambhir

Email me questions at rikab@mit.edu!
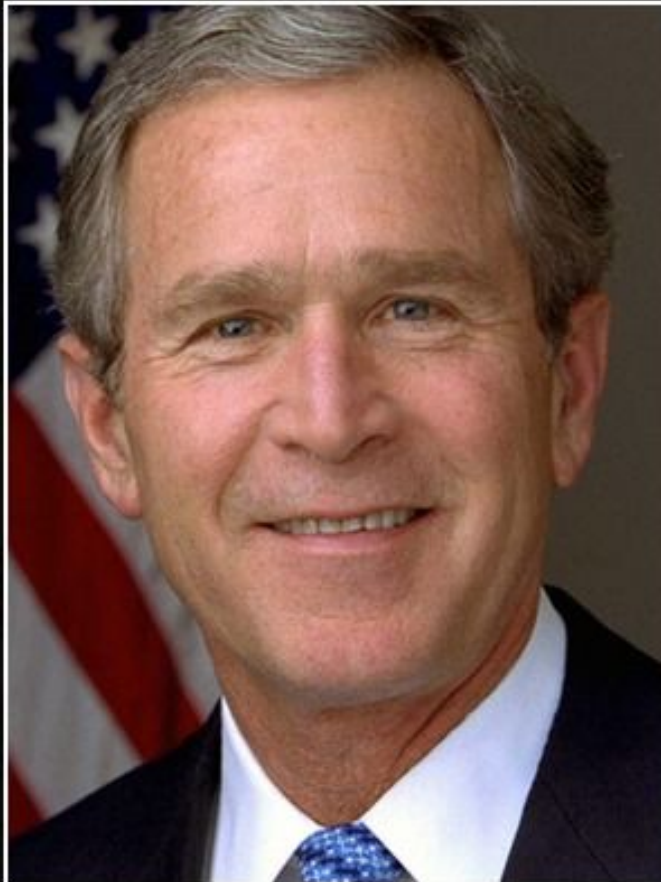
Based on [**RG**, Nachman, Thaler, 2205.03413]
[**RG**, Nachman, Thaler, 2205.05084]
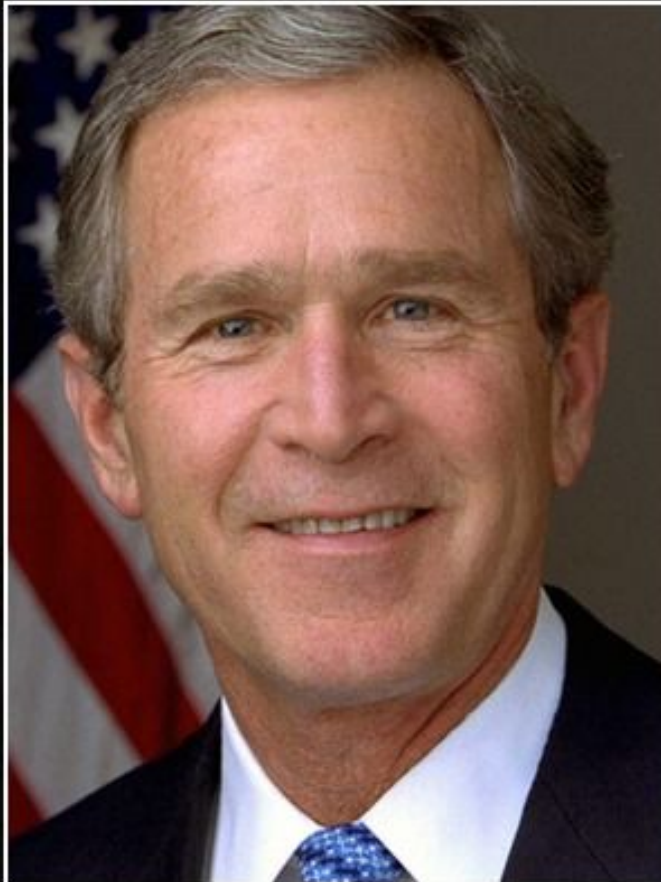[Ba, Dogra, **RG**, Tasissa, Thaler, 2302.12266]
[**RG**, Thaler, Wu, WIP]
[**RG**, Osathapan, Tasissa, Thaler, 2403.08854]

Rikab Gambhir – SLAC – 28 May 2024

Rarely is the question asked: Is our children learning?
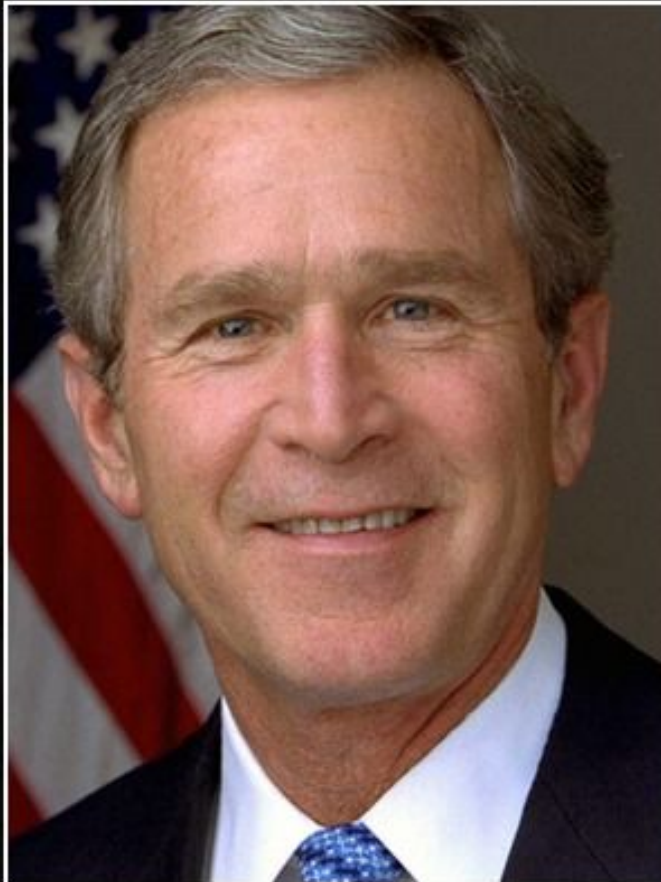
— George W. Bush —

AZ QUOTES

Rarely is the question asked: Is our ~~children~~ machines learning?

— George W. Bush —

AZ QUOTES

Rarely is the question asked: Is our ~~children~~ learning? ~~machines~~
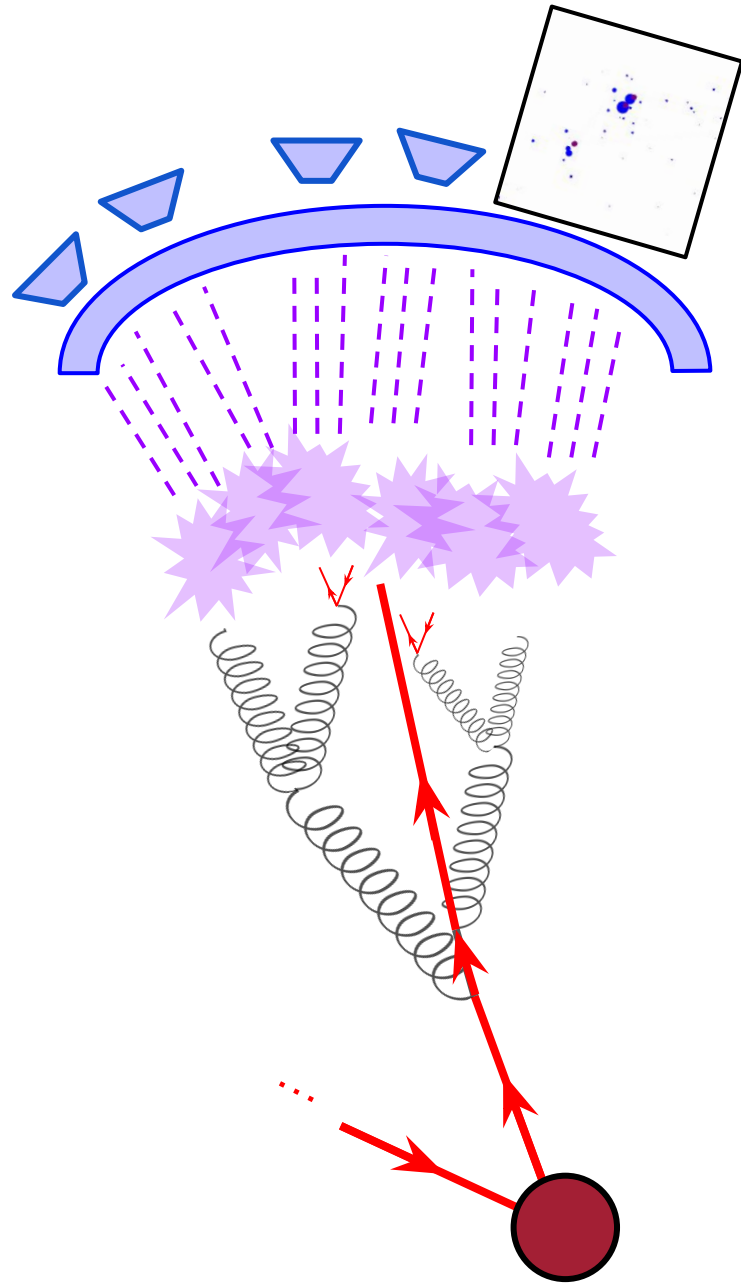
physicists using machines
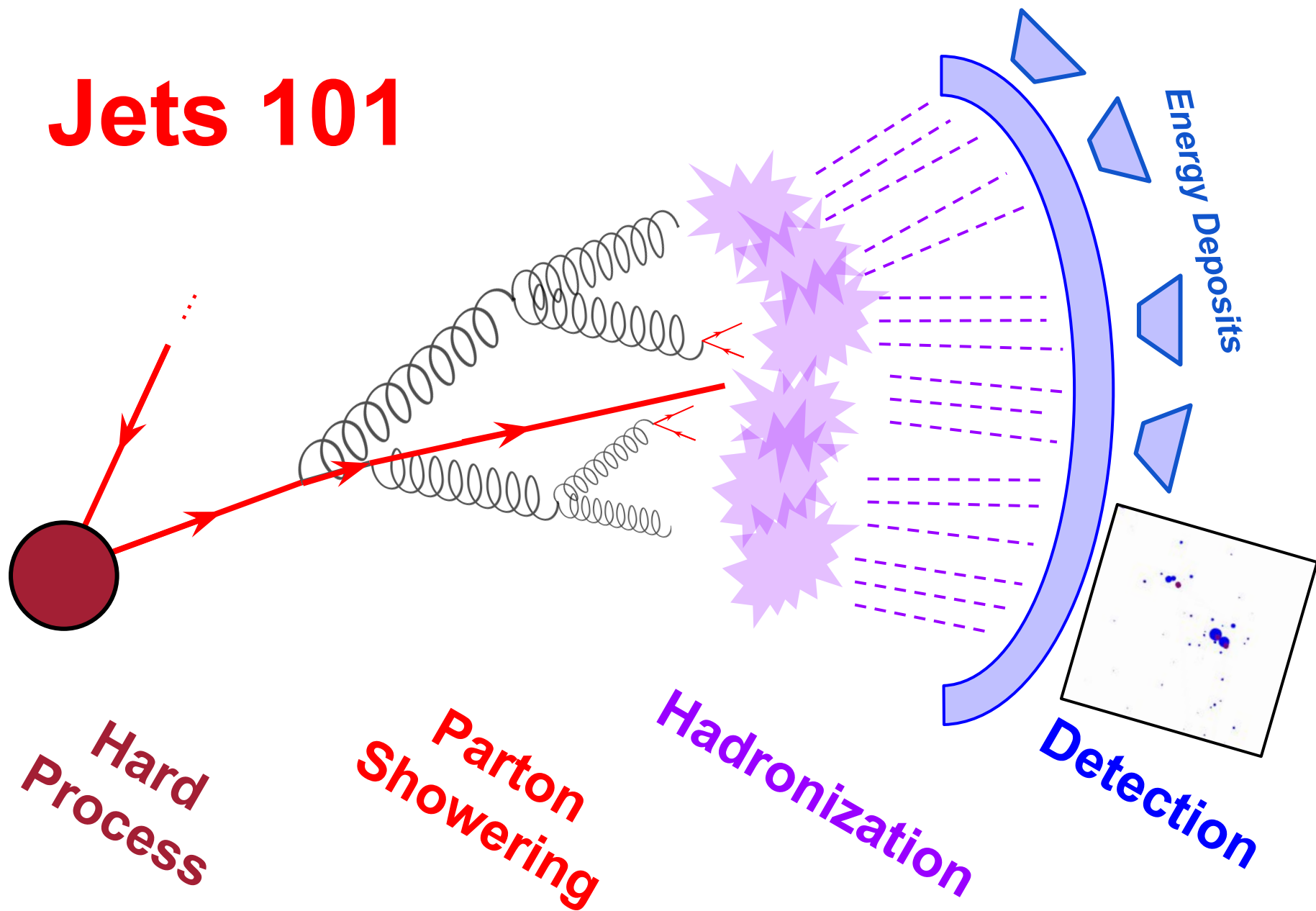
— George W. Bush —

AZ QUOTES

# Introduction

I want to study **jets** at the **Large Hadron Collider** (LHC). Jets are complicated!

I want to use **machine learning** (ML) to make my job easier, but with safeguards in place to make sure the physics I learn makes some sense!

**This talk**: 3 vignettes from my work in jet physics designing ML algorithms to give me exactly what I want.
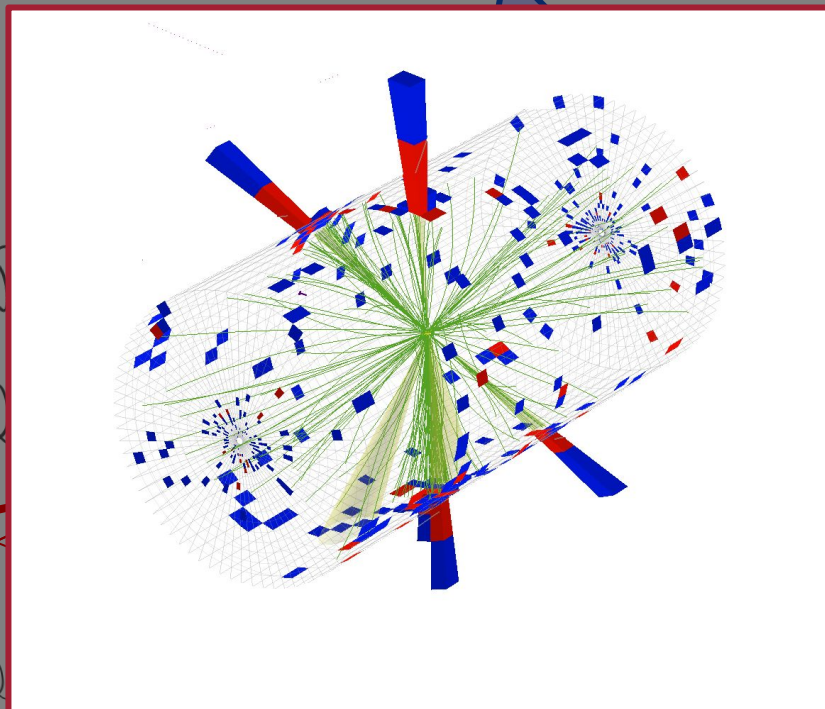
# Jets 101

Hard Process

Parton Showering

Hadronization

Detection

Energy Deposits

# Jets 101

What a jet actually looks like at the CMS detector

My theorist's picture of a jet

**Detector Geometry** $\mathscr{X}$

Detector Region $X$

$\mathscr{E}(X)$

$\phi$

$\eta$

PYTHIA8 AK8 Top Jet + Pileup

Azimuthal Angle

Rapidity

# Jets 101

Jets have **rich latent structure** and subtle correlations – amenable to machine learning!

Jets are naturally represented as **point clouds** – amenable to machine learning!
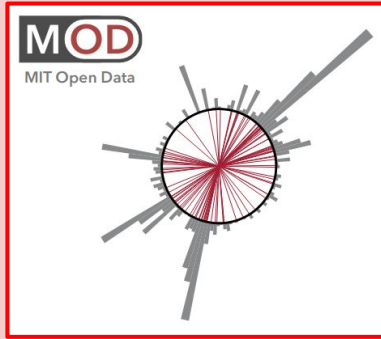
We have **sophisticated detector models** (Geant4) – amenable to machine learning!

Jet physics is an *excellent* playground for machine learning!
Lots of high dimensional structures for the machine to learn!
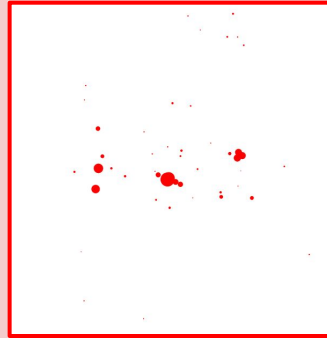
**But what can *I* learn? How can we extract this information?**
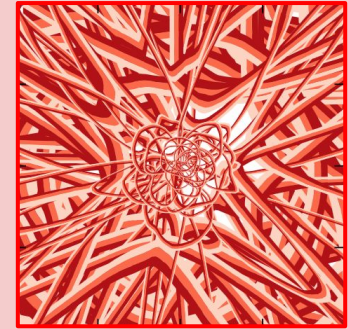
# Things a **machine** can understand

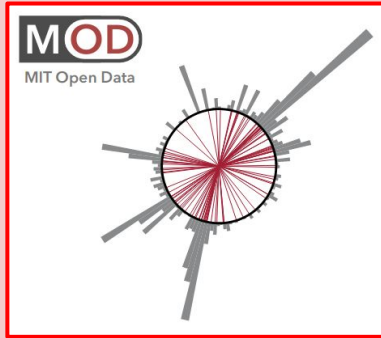Sophisticated Detector Models

Complicated Point Clouds

Huge Latent Spaces

Rikab Gambhir – SLAC – 28 May 2024

# Things a **machine** can understand

## Sophisticated Detector Models



## Complicated Point Clouds



## Huge Latent Spaces





Gaussians



1) Eff. Rad: 0.07, Ecc: 0.95, $z_\delta$, $z_\Theta$: (0.10, 0.13)
2) Eff. Rad: 0.37, Ecc: 0.78, $z_\delta$, $z_\Theta$: (0.23, 0.06)
3) Eff. Rad: 0.04, Ecc: 0.85, $z_\delta$, $z_\Theta$: (0.25, 0.24)
EMD: 0.038

Basic Kindergarten Shapes



Moment EFN Quark/Gluon
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4, p_T \in [500, 550]$ GeV
— Latent Space
- - Fit: $c_1 + c_2\log(c_3 + y)$
$k = 4, L = 1$
$c_1 = -3.584$
$c_2 = -0.847$
$c_3 = 0.005$

Addition, multiplication, and a few elementary functions
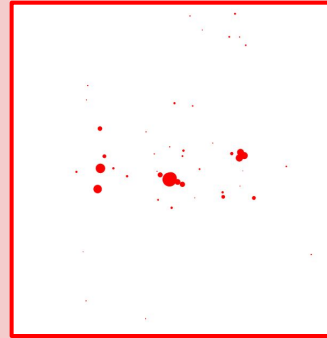
# Things my **mortal physicist brain** can understand
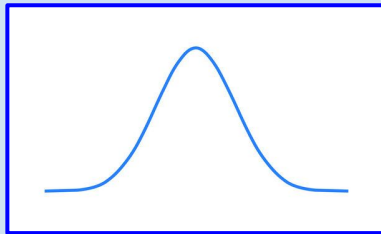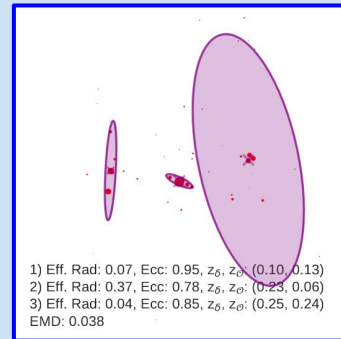
# Things a **machine** can understand

Sophisticated Detector Models



Complicated Point Clouds



Huge Latent Spaces



**Interface**: *Targeted* and *goal-motivated* network architectures and loss functions!

Using the Gaussian Ansatz and DV Loss

Using faithful optimal transport

Using Moment Pooling



Gaussians



1) Eff. Rad: 0.07, Ecc: 0.95, $z_\delta$, $z_{\mathcal{O}}$: (0.10, 0.13)
2) Eff. Rad: 0.37, Ecc: 0.78, $z_\delta$, $z_{\mathcal{O}}$: (0.23, 0.06)
3) Eff. Rad: 0.04, Ecc: 0.85, $z_\delta$, $z_{\mathcal{O}}$: (0.25, 0.24)
EMD: 0.038

Basic Kindergarten Shapes



Moment EFN Quark/Gluon
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4, p_T \in [500, 550]$ GeV
— Latent Space
- - Fit: $c_1 + c_2 \log(c_3 + y)$
$k = 4, L = 1$
$c_1 = -3.584$
$c_2 = -0.847$
$c_3 = 0.005$

Addition, multiplication, and a few elementary functions

## Things my **mortal physicist brain** can understand
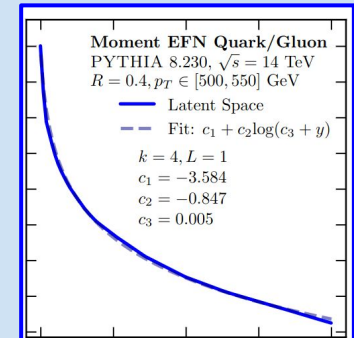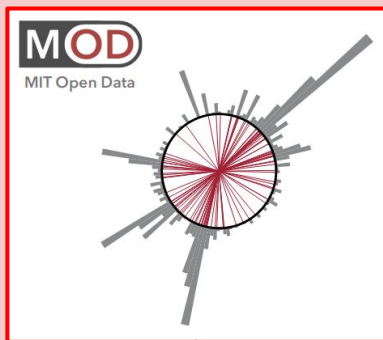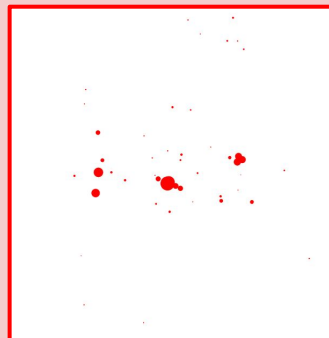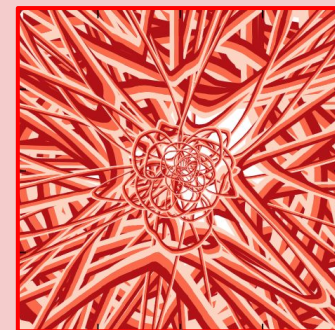
Sophisticated Detector Models     Complicated point clouds     Huge Latent Spaces



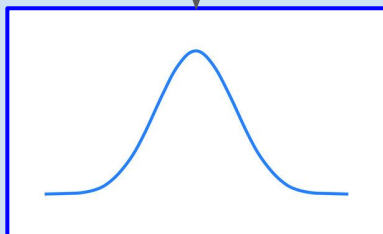Learning **Uncertainties** the **Frequentist** Way

[**RG**, Nachman, Thaler, 2205.05084
[**RG**, Nachman, Thaler, 2205.05084]

Can you Hear the **Shape** of a Jet?

[Ba, Dogra, **RG**, Tasissa, Thaler, 2302.12266]

The **Moments** of Clarity: Streamlining Latent Spaces

MPA Latent Space: k = 4, L = 16
MPA Latent Space: k = 4, L = 1

[**RG**, Osathapan, Thaler, 2403.08854]

Things my **mortal physicist brain** can understand

# Part I

# Learning **Uncertainties** the **Frequentist** Way: Calibration and Correlation

Download
our repo!

Try *pip install*
*GaussianAnsatz*

Rikab Gambhir – SLAC – 28 May 2024

# I want to **invert** this picture!

Energy Deposits

Parton with **true** energy *z* → Measured particles *x*

**Model**: Pythia + Geant4

Energy Deposits

Measured particles **$x$**

What function does this?
Machine learn it!

Parton with **estimated** energy **$\hat{z}$**

**Energy Deposits**

**Rich existing literature!**

Simulation based inference & Uncertainty Estimation:
[Cranmer, Brehmer, Louppe 1911.01429;
Alaa, van der Schaar 2006.13707;
Abdar et. al, 2011.06225;
Tagasovska, Lopez-Paz, 1811.00908;
And many more!]

Bayesian techniques:
[Jospit et. al, 2007.06823;
Wang, Yeung 1604.01662;
Izmailov et. al, 1907.07504;
Mitos, Mac Namee, 1912.1530;
And many more!]

Measured particles $x$

What function does this?
Machine learn it!

Parton with **estimated** energy $\hat{z}$

# Problem Statement

Given a training set of (*x, z*) pairs, can we machine learn a function *f* such that *f(x) = z*? With uncertainties?

**Yes. Extremely easily.** This is just bread-and-butter least-squares regression with a neural network *g*:

$$\underset{g}{\text{argmin}}\, \mathbb{E}_{\text{train}}[(g(X) - Z)^2]$$

It's the first thing you learn how to do !

Basic regression: Predict fuel efficiency 🔖 ▾

CO Run in Google Colab    ⬡ View source on GitHub    ⬇ Download notebook

In a *regression* problem, the aim is to predict the output of a continuous value, like a price or a probability. Contrast this with a *classification* problem, where the aim is to select a class from a list of classes (for example, where a picture

# Problem Statement

Given a training set of [**Not so fast!**] e learn a function *f* such that *f(x) = z*? With



**Not so fast!**

Same "detector" sim *p(x|z)*, only different priors *p(z)*!

This doesn't make a good, universal calibration.

It's the first thing

Upon closer inspection, MSE gives us prior dependent neural regressions!
Is this what we really want? Back to the drawing board!
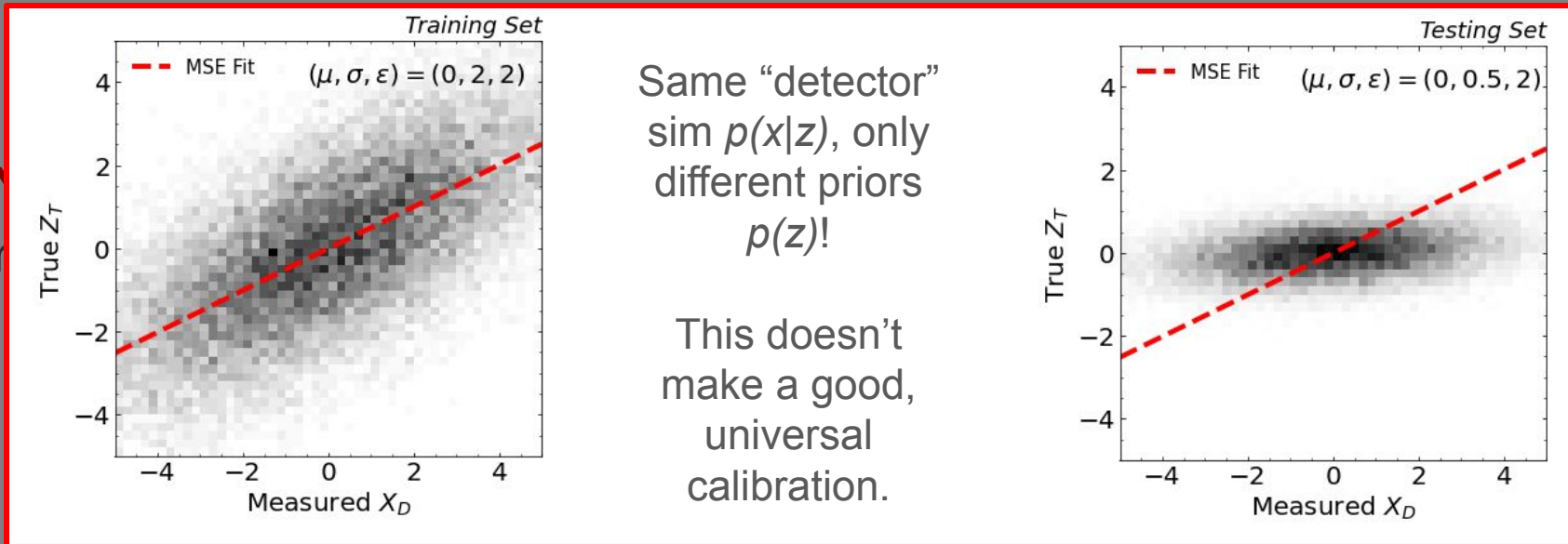
In a regression problem, the aim is to predict the output of a continuous value, like a price or a probability. Contrast this with a classification problem, where the aim is to select a class from a list of classes (for example, where a picture

# Calibration

Let's agree on what makes a good calibration, then design a loss for it!

1. **Closure**: On average, $f(x)$ should be correct for each $x$! That is, $f$ is **unbiased**.

$$b(z) = \mathbb{E}_{\text{test}}[f(X) - z | Z = z]$$
$$= 0$$

2. **Universality**: $f(x)$ should not depend on the choice of sampling for $z$. That is, $f$ is **prior-independent**.

$f$ depends only on $p(x|z)$, and not $p(z)$

**Likelihood**: Detector simulation, noise model, etc

What if the detector simulation is imperfect? Ask me later!

# Calibration

Our network should make it easy to extract unbiased [design a loss for it!] estimates! (e.g. via maximum likelihood)

1. **Closure**: On average, $f(x)$ should be correct for each $x$! That is, $f$ is **unbiased**.

$$b(z) = \mathbb{E}_{\text{test}}[f(X) - z | Z = z]$$
$$= 0$$

Our loss should give us this likelihood!

2. **Universality**: $f(x)$ should not depend on the choice of sampling for $z$. That is, $f$ is **prior-independent**.

$f$ depends only on $p(x|z)$, and not $p(z)$

**Likelihood**: Detector simulation, noise model, etc

What if the detector simulation is imperfect? Ask me later!

# Learning the likelihood

The **Donsker-Varadhan Representation (DVR)** of the KL divergence has been used in the statistics literature for mutual information estimation

$$\mathcal{L}_{\text{DVR}}[T] = -\left( \mathbb{E}_{P_{XZ}}[T] - \log\left( \mathbb{E}_{P_X \otimes P_Z}\left[e^T\right]\right)\right)$$

Interestingly, a nonlocal loss!

Strict bound on *I(X;Z)*

**What we want!**

Minimized when $\quad T(x,z) = \log \dfrac{p(x|z)}{p(x)} + c$

Unimportant

Lots of other losses also work, but DVR has very nice convergence properties - ask me later!

# Inference using *T*

We can use a neural net to parameterize *T(x,z)*, and use standard gradient descent techniques to minimize the DVR loss. Then we can do …

$$\hat{z}(x) = \underset{z}{\operatorname{argmax}}\ T(x,z)$$

$$\left[\hat{\sigma}_z^2(x)\right]_{ij} = -\left[\frac{\partial^2 T(x,z)}{\partial z_i\,\partial z_j}\right]^{-1}\Bigg|_{z=\hat{z}}$$

**Inference**                    **Gaussian Uncertainty Estimation**

**BUT!**

- Maxima are hard to estimate – even *more* gradient descent!
- Second derivatives are sensitive to the choice of activations in *T* – ReLU spoils everything!

We solve both problems with the **Gaussian Ansatz**

Rikab Gambhir – SLAC – 28 May 2024

# The Gaussian Ansatz

Parameterize *T(x,y)* in the following way (the **Gaussian Ansatz**):

$$T(x, z) = A(x)$$
$$+ (z - B(x)) \cdot D(x)$$
$$+ \frac{1}{2} (z - B(x))^T \cdot C(x, z) \cdot (z - B(x))$$

Where *A(x), B(x), C(x,z), and D(x)* are learned functions. Then, if *D→0*, our inference and (Gaussian) uncertainties are given by …
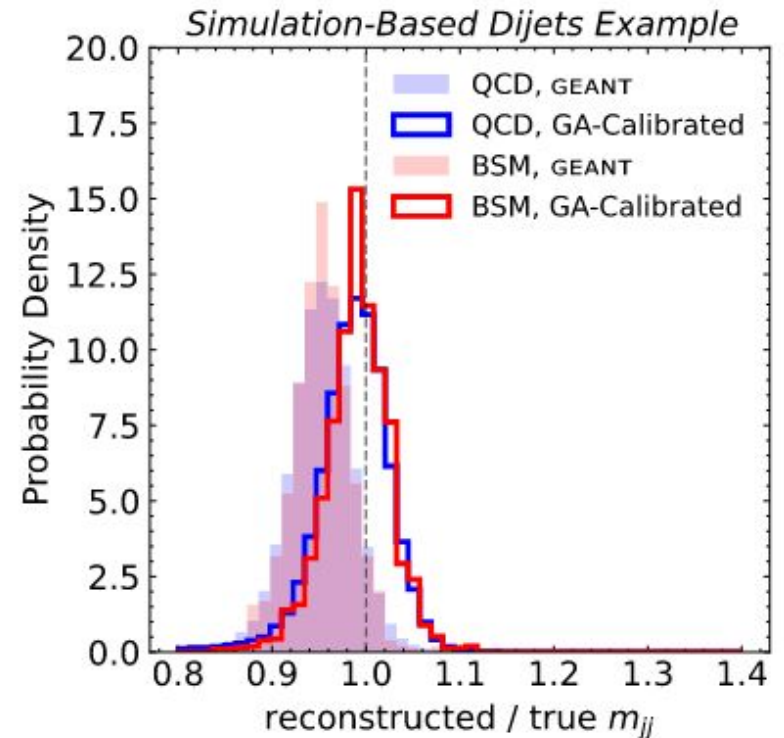
$$\hat{z}(x) = B(x) \qquad \hat{\sigma}_z^2(x) = -\left[ C(x, B(x)) \right]^{-1}$$

**No additional post processing or numerical estimates required!**

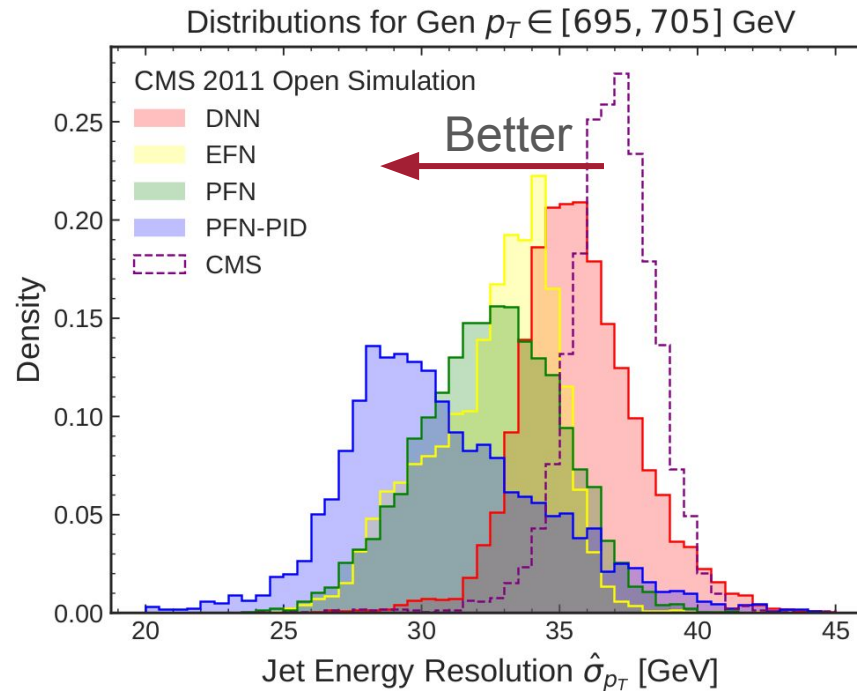# Example 1: QCD and BSM Dijets



(Left) MSE-fitted network.

(Right) Gaussian Ansatz-fitted network

Clever loss function choice: *Prior dependence is built-in!*

# Example 2: Jet Energy Resolution



Distributions for Gen $p_T \in [695, 705]$ GeV

CMS 2011 Open Simulation
- DNN
- EFN
- PFN
- PFN-PID
- CMS

Better

Density

Jet Energy Resolution $\hat{\sigma}_{p_T}$ [GeV]

**Best of both worlds**: Using ML to extract more information[*] than hand-crafted features, while still being able to extract resolutions in a prior-independent and unbiased way!

[*]I mean information literally. Ask me later about the cool information theory of the DV loss!
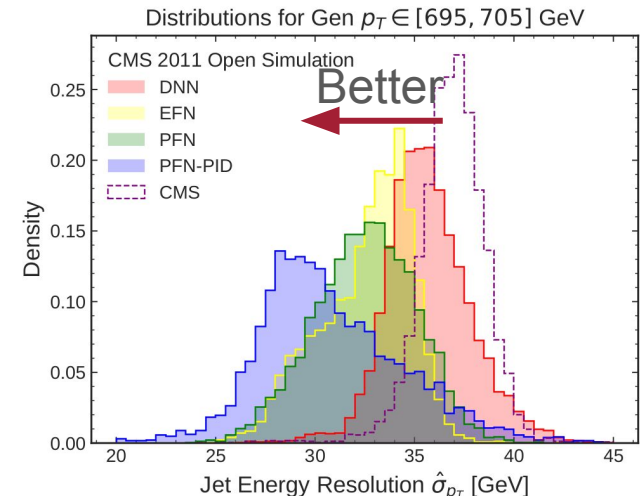
# Part I Summary

By choosing the following loss:

$$\mathcal{L}_{\text{DVR}}[T] = -\left( \mathbb{E}_{P_{XZ}}[T] - \log\left( \mathbb{E}_{P_X \otimes P_Z}[e^T] \right) \right)$$

With the following network parameterization:

$$\begin{aligned}
T(x,z) = A(x) \\
+ \left( z - B(x) \right) \cdot D(x) \\
+ \frac{1}{2} \left( z - B(x) \right)^T \cdot C(x,z) \cdot \left( z - B(x) \right)
\end{aligned}$$

We can get **unbiased**, **prior-independent** inference, and easily extract **maximum likelihood** and **resolution** estimates!
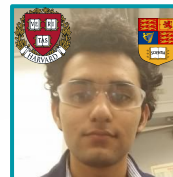


Distributions for Gen $p_T \in [695, 705]$ GeV

CMS 2011 Open Simulation
DNN
EFN
PFN
PFN-PID
CMS

Better

Density

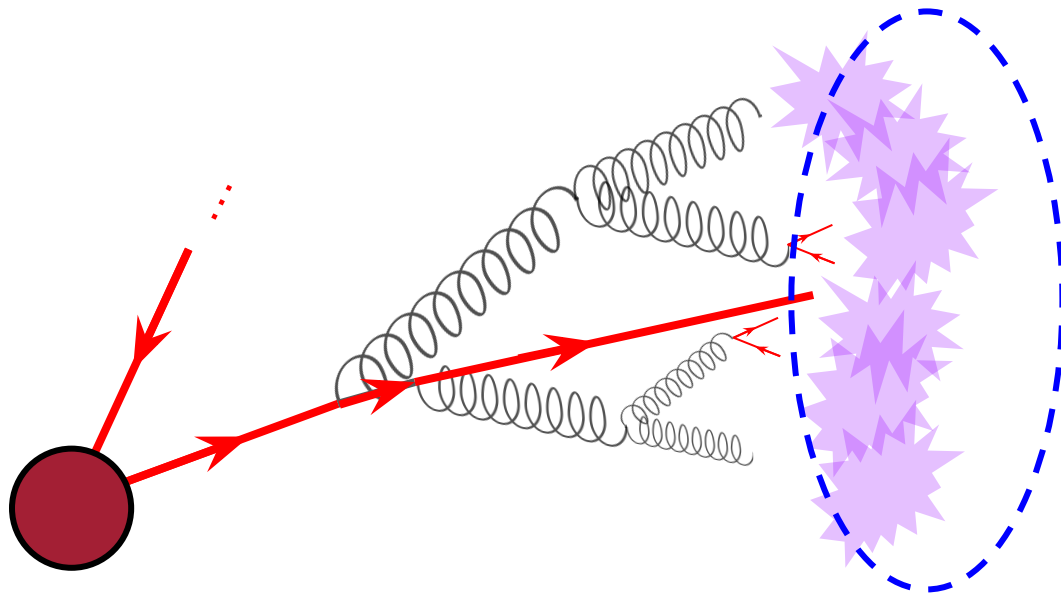Jet Energy Resolution $\hat{\sigma}_{p_T}$ [GeV]

# Part II
# Can you Hear the **Shape** of a Jet?

Download our repo!

Try *pip install pyshaper*

How "*big*" is this jet?

… What do we even mean by the "size" of a jet?

Azimuthal Angle $\phi$

Rapidity $\eta$

I don't mean the fixed size of a jet algorithm, e.g. AK4, but rather a dynamical notion of a jet size!

Rikab Gambhir – SLAC – 28 May 2024

How "*big*" is this jet?

… What do we even mean by the "size" of a jet?

By following this line of thinking, we'll be able to come up with a new class of QCD observables!

Azimuthal Angle $\phi$

Rapidity $\eta$
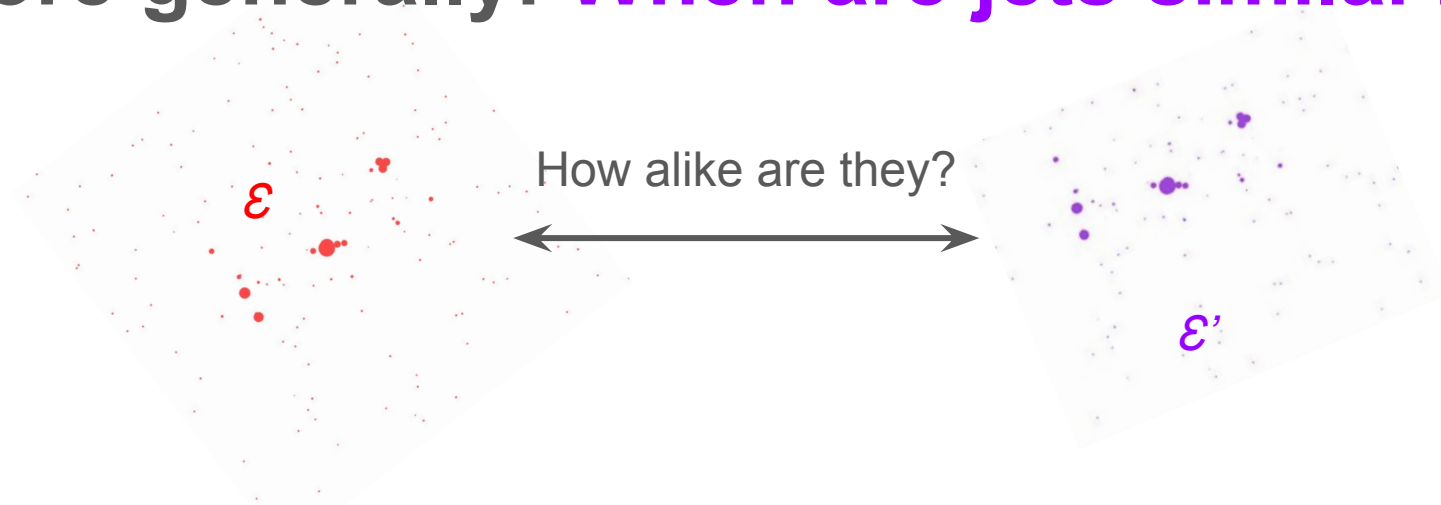
I don't mean the fixed size of a jet algorithm, e.g. AK4, but rather a dynamical notion of a jet size!

# More generally: **When are jets similar?**



How alike are they?

$\mathcal{E}$ ⟷ $\mathcal{E}'$
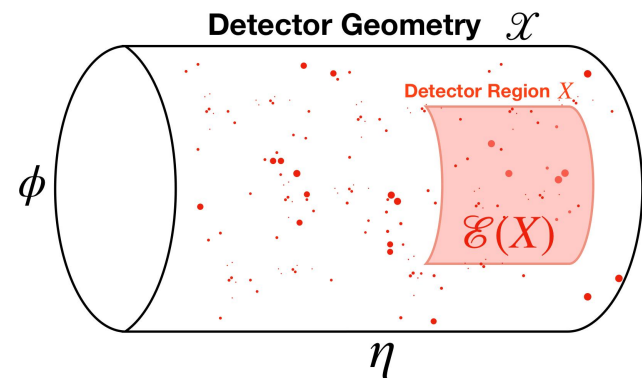
Jets can be represented as **point clouds** – let's scour through the ML and computer vision literature for a metric on point clouds we like![*]

Detector Coordinate $(\eta, \phi)$

$$\mathcal{E}(y) = \sum_i z_i \delta(y - y_i)$$

Energy Fraction $E_i / E_{Tot}$



**Detector Geometry** $\mathcal{X}$

**Detector Region** $X$

$\phi$

$\mathcal{E}(X)$

$\eta$

[*]Or better yet, construct one!

# The Wasserstein Metric

Let's demand the following reasonable properties of our metric on point clouds:

1.  … is nonnegative and finite
2.  … is **IRC-safe**[*] (Calculable and robust)
3.  … is translationally invariant
4.  … is invariant to particle labeling
5.  respects the detector metric ***faithfully***[**]

[*]Ask me for more details on this offline!
[**] Preserves distances between *extended* objects, not just points

Rikab Gambhir – SLAC – 28 May 2024

# The Wasserstein Metric

Let's demand the following reasonable properties of our metric on point clouds:

1.  … is nonnegative and finite
2.  … is **IRC-safe**[*] (Calculable and robust)
3.  … is translationally invariant
4.  … is invariant to particle labeling
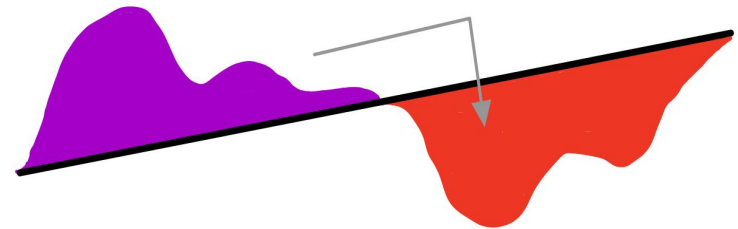5.  respects the detector metric ***faithfully***[**]

EMD = Work done to move "dirt" optimally

It turns out that the *only* metric that satisfies this is the **Wasserstein Metric / EMD**!

$$\text{EMD}^{(\beta,R)}(\mathcal{E}, \mathcal{E}') = \min_{\pi \in \mathcal{M}(\mathcal{X} \times \mathcal{X})} \left[ \frac{1}{\beta R^\beta} \left\langle \pi, d(x,y)^\beta \right\rangle \right] + |\Delta E_{\text{tot}}|$$

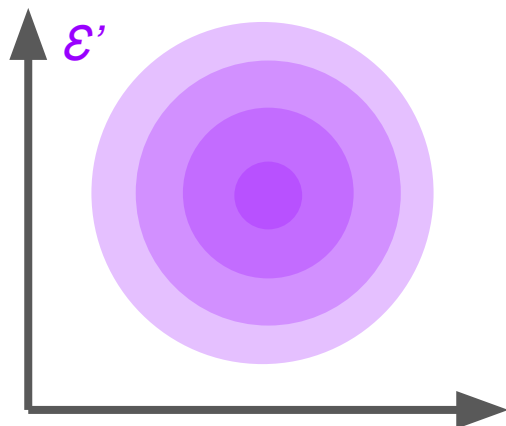$$\pi(\mathcal{X}, Y) \le \mathcal{E}'(Y), \quad \pi(X, \mathcal{X}) \le \mathcal{E}(X), \quad \pi(\mathcal{X}, \mathcal{X}) = \min(E_{\text{tot}}, E'_{\text{tot}})$$

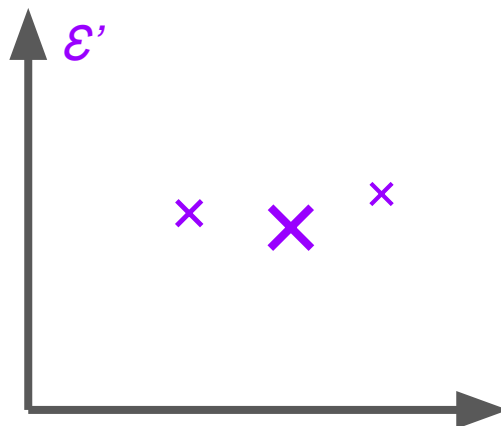A staple of computer vision ML is also useful for jets!

Rikab Gambhir – SLAC – 28 May 2024

# **Shapes** as Energy Flows

Energy flows don't have to be real events – they can be *any* energy distribution in detector space, or **shape**.

Can make anything you want! Even continuous or complicated shapes.
(Or, something you calculate in perturbative QCD)
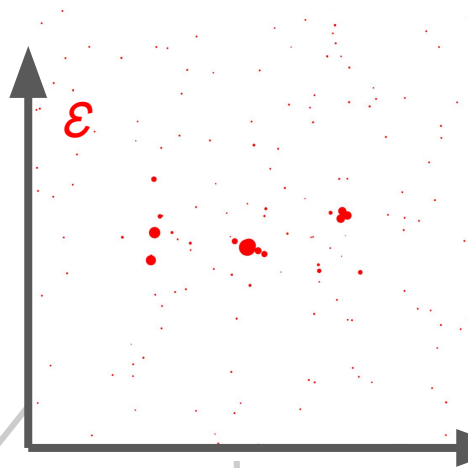


**Shape** = 2D Gaussian       **Shape** = 3 Points       **Shape** = Smile

Rikab Gambhir – SLAC – 28 May 2024

# Shapiness

The EMD between a real event or jet $\mathcal{E}$ and idealized shape $\mathcal{E}'$ is the [shape]iness of $\mathcal{E}$ – a well defined IRC-safe observable!
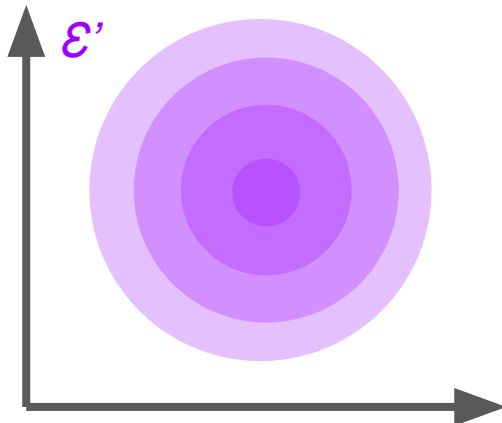
Answers the question: "How much like the shape $\mathcal{E}'$ is my event $\mathcal{E}$?"



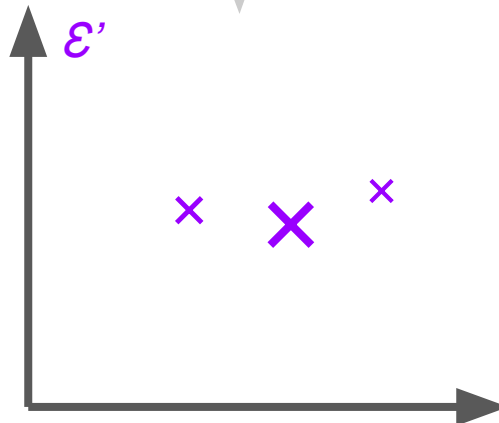Med EMD($\mathcal{E}$, $\mathcal{E}'$) "Gaussiness"

Low EMD($\mathcal{E}$, $\mathcal{E}'$) "3-Pointiness"
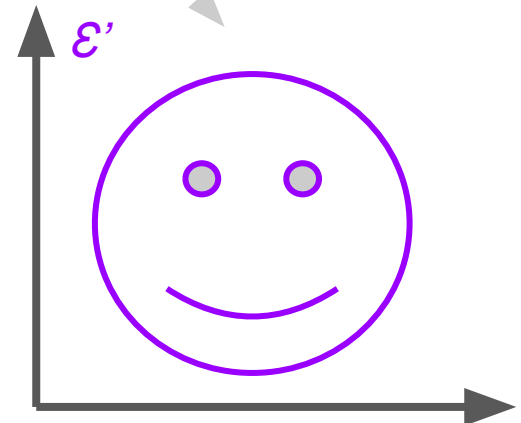AKA "3-Subjettiness"

High EMD($\mathcal{E}$, $\mathcal{E}'$) "Smileyness"

**Shape** = 2D Gaussian

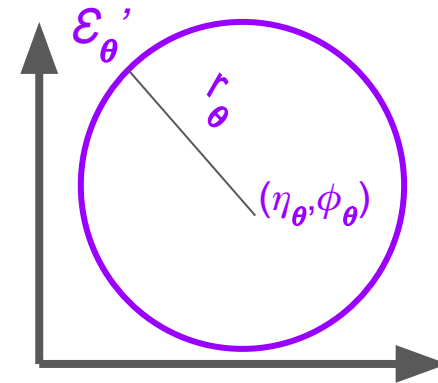**Shape** = 3 Points

**Shape** = Smile

# Mathematical Details - **Shapiness**

Rather than a single shape, consider a **parameterized manifold** $\mathcal{M}$ of **energy flows.**

e.g. The manifold of uniform circle energy flows:

$$\mathcal{E}_{\theta}'(y) = \begin{cases} \frac{1}{2\pi r_\theta} & |\vec{y} - \vec{y}_\theta| = r_\theta \\ 0 & |\vec{y} - \vec{y}_\theta| \neq r_\theta \end{cases}$$



Then, for an event $\mathcal{E}$, define the **shapiness** $\mathcal{O}_{\mathcal{M}}$ and **shape parameters** $\theta_{\mathcal{M}}$, given by:

$$\mathcal{O}_{\mathcal{M}}(\mathcal{E}) \equiv \min_{\mathcal{E}_\theta \in \mathcal{M}} \mathrm{EMD}^{(\beta,R)}(\mathcal{E}, \mathcal{E}_\theta)$$

$$\theta_{\mathcal{M}}(\mathcal{E}) \equiv \underset{\mathcal{E}_\theta \in \mathcal{M}}{\mathrm{argmin}}\ \mathrm{EMD}^{(\beta,R)}(\mathcal{E}, \mathcal{E}_\theta)$$

# This is basically just a **W-GAN**!

Fitting a 2D image distribution with a fully flexible generative neural network



Fitting a 2D point cloud with a small parameterized generator



$z_{PU}$: 0.00
1) Eff. Rad: 0.00, Ecc: nan, $z_\delta$, $z_\mathcal{O}$: (0.24, 0.00)
2) Eff. Rad: 0.00, Ecc: nan, $z_\delta$, $z_\mathcal{O}$: (0.30, 0.00)
3) Eff. Rad: 0.00, Ecc: nan, $z_\delta$, $z_\mathcal{O}$: (0.46, 0.00)
EMD: inf

Rather than using a fully-flexible neural network to fit our distributions with the Wasserstein metric, as in a W-GAN, we craft specific parameters of interest!

# Observables ⟺ Manifolds of Shapes

Observables can be specified solely by defining a **manifold of shapes**:

$$\mathcal{O}_{\mathcal{M}}(\mathcal{E}) \equiv \min_{\mathcal{E}_\theta \in \mathcal{M}} \mathrm{EMD}^{(\beta,R)}(\mathcal{E}, \mathcal{E}_\theta),$$
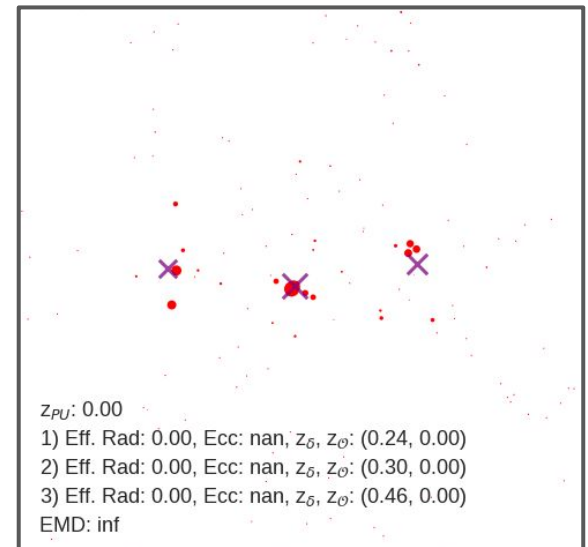
$$\theta_{\mathcal{M}}(\mathcal{E}) \equiv \operatorname*{argmin}_{\mathcal{E}_\theta \in \mathcal{M}} \mathrm{EMD}^{(\beta,R)}(\mathcal{E}, \mathcal{E}_\theta),$$

Many well-known observables[*] already have this form!

| Observable | Manifold of Shapes |
|---|---|
| $N$-Subjettiness | Manifold of $N$-point events |
| $N$-Jettiness | Manifold of $N$-point events with floating total energy |
| Thrust | Manifold of back-to-back point events |
| Event / Jet Isotropy | Manifold of the single uniform event |

… and more!

All of the form "How much like **[shape]** does my **event** look like?"
Generalize to *any* shape.

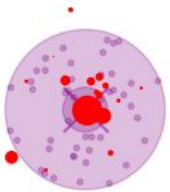[*]These observables are usually called event shapes or jet shapes in the literature – we are making this literal!

Rikab Gambhir – SLAC – 28 May 2024

# Hearing Jets **Ring**

(and Disk, and Ellipse)

$\mathcal{O}_{\mathcal{M}}(\mathcal{E})$ answers: "How much like a shape in $\mathcal{M}$ does my event **ε** look like?"

$\theta_{\mathcal{M}}(\mathcal{E})$ answers: "Which shape in $\mathcal{M}$ does my event **ε** look like?"

Can define complex manifolds to probe increasingly subtle geometric structure!

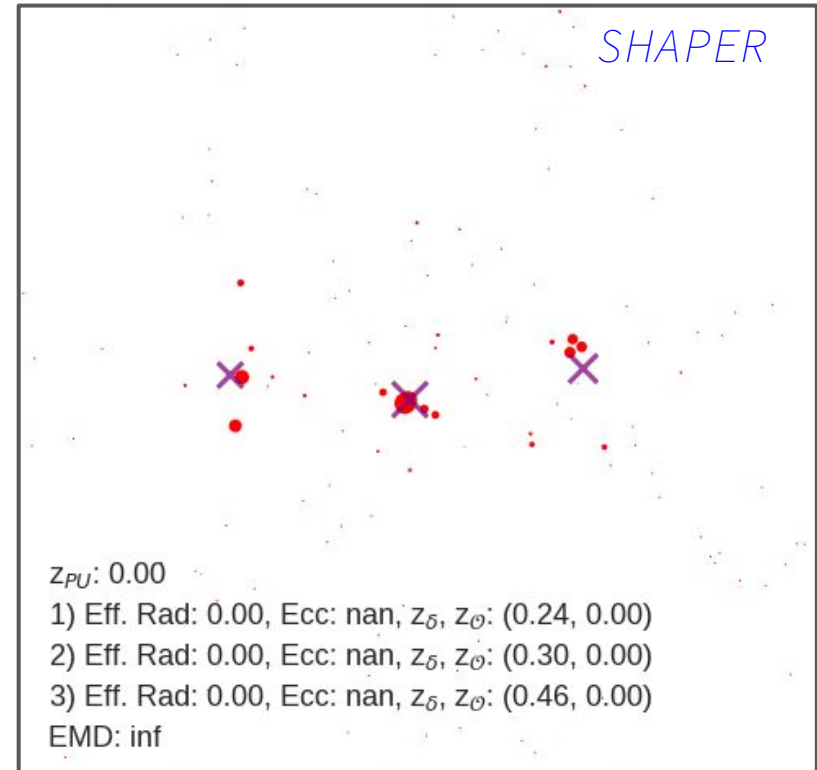| Shape | Specification | Illustration |
|---|---|---|
| **Ringiness** $\mathcal{O}_R$ | **Manifold of Rings** $\mathcal{E}_{x_0, R_0}(x) = \frac{1}{2\pi R_0}$ for $|x - x_0| = R_0$ <br> $x_0 = $ Center, $R_0 = $ Radius |  |
| **Diskiness** $\mathcal{O}_D$ | **Manifold of Disks** $\mathcal{E}_{x_0, R_0}(x) = \frac{1}{\pi R_0^2}$ for $|x - x_0| \leq R_0$ <br> $x_0 = $ Center, $R_0 = $ Radius |  |
| **Ellipsiness** $\mathcal{O}_E$ | **Manifold of Ellipses** $\mathcal{E}_{x_0, a, b, \varphi}(x) = \frac{1}{\pi ab}$ for $x \in \text{Ellipse}_{x_0, a, b, \varphi}$ <br> $x_0 = $ Center, $a, b = $ Semi-axes, $\varphi = $ Tilt |  |
| **(Ellipse Plus Point)iness** | **Composite Shape** $\mathcal{O}_E \oplus \tau_1$ <br> Fixed to same center $x_0$ |  |
| **N-(Ellipse Plus Point)iness Plus Pileup** | **Composite Shape** $N \times (\mathcal{O}_E \oplus \tau_1) \oplus \mathcal{I}$ |  |



**Disk** + δ-function     **Ring** + δ-function

Probing **collinear** (δ-function) and **soft** (shape) structure!

Some examples of new shapes you can define!

[see also L., B. Nachman, A. Schwartzman, C. Stansbury, 1509.02216;
see also B. Nachman, P. Nef, A. Schwartzman, M. Swiatlowski, C. Wanotayaroj, 1407.2922;
see also M. Cacciari, G. Salam, 0707.1378]

# New IRC-Safe Observables

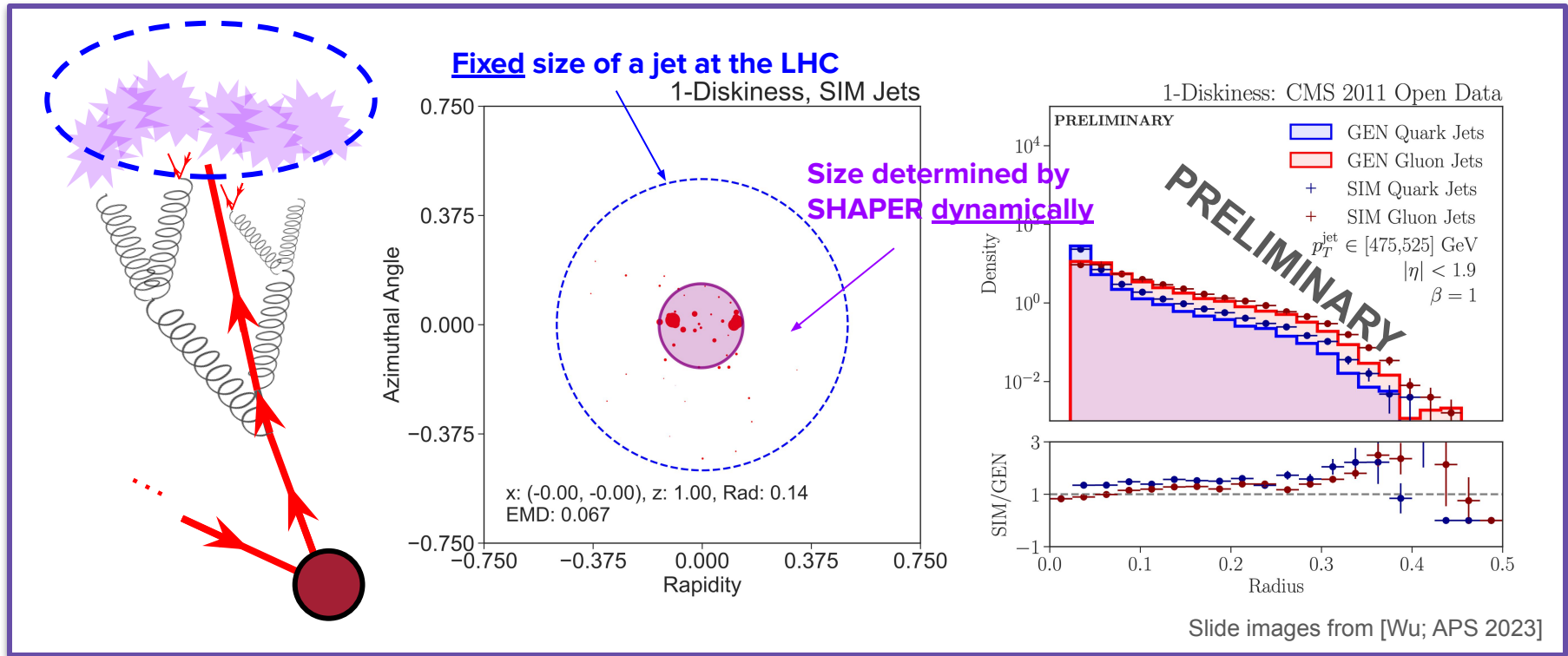The *SHAPER* framework makes it easy to algorithmically invent new jet observables!

e.g. **N-(Ellipse+Point)iness+Pileup** as a jet algorithm:

- Learn jet centers + collinear radiation
- Dynamic jet radii (no *R* parameter!)
- Dynamic eccentricities and angles
- Dynamic jet energies
- Dynamic pileup (no $z_{cut}$ parameter!!!)
- Learned parameters for discrimination

Can design custom specialized jet algorithms to learn jet substructure!



*SHAPER*

$z_{PU}$: 0.00
1) Eff. Rad: 0.00, Ecc: nan, $z_\delta$, $z_\mathcal{O}$: (0.24, 0.00)
2) Eff. Rad: 0.00, Ecc: nan, $z_\delta$, $z_\mathcal{O}$: (0.30, 0.00)
3) Eff. Rad: 0.00, Ecc: nan, $z_\delta$, $z_\mathcal{O}$: (0.46, 0.00)
EMD: inf

Rikab Gambhir – SLAC – 28 May 2024

# Back to our question: How Big are Jets?



**Fixed** size of a jet at the LHC

**Size determined by SHAPER dynamically**

1-Diskiness, SIM Jets

x: (-0.00, -0.00), z: 1.00, Rad: 0.14
EMD: 0.067

1-Diskiness: CMS 2011 Open Data

PRELIMINARY

GEN Quark Jets
GEN Gluon Jets
+ SIM Quark Jets
+ SIM Gluon Jets
$p_T^{\text{jet}} \in [475,525]$ GeV
$|\eta| < 1.9$
$\beta = 1$

PRELIMINARY

Slide images from [Wu; APS 2023]

We can now answer this question in a precise way!

# Part II Summary

By choosing the following loss function:

$$\text{EMD}^{(\beta,R)}(\mathcal{E}, \mathcal{E}') = \min_{\pi \in \mathcal{M}(\mathcal{X} \times \mathcal{X})} \left[ \frac{1}{\beta R^\beta} \left\langle \pi, d(x,y)^\beta \right\rangle \right] + |\Delta E_{\text{tot}}|$$

$$\pi(\mathcal{X}, Y) \leq \mathcal{E}'(Y), \quad \pi(X, \mathcal{X}) \leq \mathcal{E}(X), \quad \pi(\mathcal{X}, \mathcal{X}) = \min(E_{\text{tot}}, E'_{\text{tot}})$$

Based on IRC-safety and geometric faithfulness, we can build well-defined and robust observables that quantify *targeted* geometric properties of point clouds!

We can make jet shapes well defined!



3-Diskiness, SIM Jets, Epoch init

1) x: (-0.06, -0.14), z: 0.11, Rad: 0.00
2) x: (-0.03, 0.00), z: 0.34, Rad: 0.00
3) x: (0.03, 0.03), z: 0.54, Rad: 0.00
EMD: inf

Rikab Gambhir – SLAC – 28 May 2024

# Part III

# Moments of Clarity: Streamlining Latent Spaces

Rikab Gambhir – SLAC – 28 May 2024

# Jet Classification

Energy Deposits

Given the final measured **point cloud**, was the **initiating parton** a **quark** or a **gluon**?

A staple machine learning task!

# Typical **Machine Learning** Setup



**Collider Data** → **Latent Space** → **Observables**

~1000 Dimensional ~10-100 Dimensional ~1-10 Dimensional

Quark vs. Gluon Jets
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4$, $p_T \in [500, 550]$ GeV

PFN-ID
PFN-Ex
PFN-Ch    EFN
PFN       EFPs

Pictured: An Energy Flow Network (**EFN**):

$$\mathcal{O}(\{p_1, \ldots, p_M\}) = F\left(\sum_{i=1}^{M} z_i \Phi(\hat{p}_i)\right)$$

(How) can we understand and constrain this?

(How) can we be more efficient?

Rikab Gambhir – SLAC – 28 May 2024

# Deep Learning on Point Clouds

The **Deep Sets Theorem** tells us how to parameterize functions on point clouds[*]:

Set of momenta

a = 1 … L, the *Latent Dimension* index

$$\mathcal{O}(\mathcal{P}) = F\left(\langle \phi^a \rangle_{\mathcal{P}}\right)$$

The Deep Sets Theorem guarantees that any function on sets $\mathcal{P}$ can be written this way, for "sufficiently complex" $F$ and $\phi$ and large enough $L$.

$$\langle \phi \rangle_{\mathcal{P}} \equiv \sum_i z_i \phi(\hat{p}_i)$$

For this talk, I will focus primarily on Energy Flow Networks (EFNs) where the sums are energy-weighted.

[*]More sophisticated architectures still reduce to Deep Sets in some limit, e.g. graph networks

# Deep Learning on Point Clouds

The **Deep Sets Theorem** tells us how to parameterize functions on point clouds[*]:

Set of momenta        $a = 1 \dots L$, the *Latent Dimension* index

$$\mathcal{O}(\mathcal{P}) = F\left(\langle \phi^a \rangle_{\mathcal{P}}\right)$$

The Deep Sets Theorem guarantees that any function on sets $\mathcal{P}$ can be written this way, for "sufficiently complex" $F$ and $\phi$ and large enough $L$.

$$\langle \phi \rangle_{\mathcal{P}} \equiv \sum_i z_i \phi(\hat{p}_i)$$

But how complex do $F$ and $\phi$ need to be? Can I reduce the number of latent dimensions?

For this talk, I will focus primarily on Energy Flow Networks (EFNs) where the sums are energy-weighted.

[*]More sophisticated architectures still reduce to Deep Sets in some limit, e.g. graph networks

Rikab Gambhir – SLAC – 28 May 2024

# Moment Pooling

$$\mathcal{O}(\mathcal{P}) = F\left(\langle \phi^a \rangle_{\mathcal{P}}\right)$$

**Generalize** to more moments! "**Moment Pooling**"

k = Highest order moment considered

$$\mathcal{O}_k(\mathcal{P}) = F_k\left(\underbrace{\langle \phi^a \rangle_{\mathcal{P}}}_{\text{1}^{\text{st}} \text{ Moment}}, \underbrace{\langle \phi^{a_1}\phi^{a_2} \rangle_{\mathcal{P}}}_{\text{2}^{\text{nd}} \text{ Moment}}, ..., \underbrace{\langle \phi^{a_1}...\phi^{a_k} \rangle_{\mathcal{P}}}_{k^{\text{th}} \text{ Moment}}\right)$$

$$L_{\text{eff}} = \binom{L+k}{k}$$

Example: 2$^{\text{nd}}$ Moment
$$\langle \Phi^{a_1}\Phi^{a_2} \rangle_{\mathcal{P}} = \sum_{i \in \mathcal{P}} z_i \Phi^{a_1}(x_i)\Phi^{a_2}(x_i)$$

# **Moment Pooling** – Why?

k = Highest order moment considered

$$\mathcal{O}_k(\mathcal{P}) = F_k \left( \underbrace{\langle \phi^a \rangle_{\mathcal{P}}}_{\text{1st Moment}}, \underbrace{\langle \phi^{a_1} \phi^{a_2} \rangle_{\mathcal{P}}}_{\text{2nd Moment}}, ..., \underbrace{\langle \phi^{a_1} ... \phi^{a_k} \rangle_{\mathcal{P}}}_{k^{\text{th}} \text{ Moment}} \right)$$

$$L_{\text{eff}} = \binom{L + k}{k}$$

- **Explicit Multiplication**: Neural nets are mostly piecewise linear! But most functions we learn aren't. Moments are just multiplication, but for distributions!
- **Latent Space Compression**: For large $L$, there are effectively $L^k$ latent dimensions due to combinatorics, but still made of only $L$ functions! Fewer latent dimensions means easier analysis!

# Moment Pooling Results



A high order **Moment EFN** can achieve the same top performance on quark/gluon classification as an **ordinary EFN**, but with far fewer latent dimensions!*

A Moment EFN is a more efficient encoding of the same data!

*Interestingly, the performance is constant in the *effective* latent dimensions and number of parameters. Ask me about this!

# Latent Spaces



EFN ($k = 1$)
$L = 128$

$k = 2$
$L = 64$

$k = 3$
$L = 32$

$k = 4$
$L = 16$

The 45% - 55% contours of each of the $L$ different $\phi$ functions

*Interestingly, the performance is constant in the *effective* latent dimensions and number of parameters. Ask me about this!

# Moment Pooling Results (Cont'd)



Let's instead look at networks with just a *single* latent dimension.

A **k = 4 Moment EFN** with *just one* latent dimension does just as good as an **ordinary EFN** with *four* latent dimensions!

# Going down to *L = 1*

The more efficient encoding is extremely useful when we can get to *L = 1*!
Look how much simpler allowing neural networks to multiply makes things!

EFN (*k = 1*)
*L = 4*

Order *k = 4*
*L = 1*



Same information!



- 4 Different Functions
- Combined in a complicated way
- No symmetry
- ???

- Single function
- Expected to be, and is, radially symmetric
- We can interpret this!

# The **Moment** of Truth

We can extract physics from this!

$$\Phi_{\mathcal{L}}(r) = c_1 + c_2 \log(c_3 + r)$$

This defines a **log angularity** observable, related to the $\square \rightarrow 0$ limit of ordinary angularities

We can even see non perturbative physics!

$$c_3 \sim \frac{\Lambda_{\mathrm{QCD}}}{p_T R} \sim 0.001$$

Just this alone is enough to get an IRC-safe quark-gluon classifier with an AUC ~ 0.83!
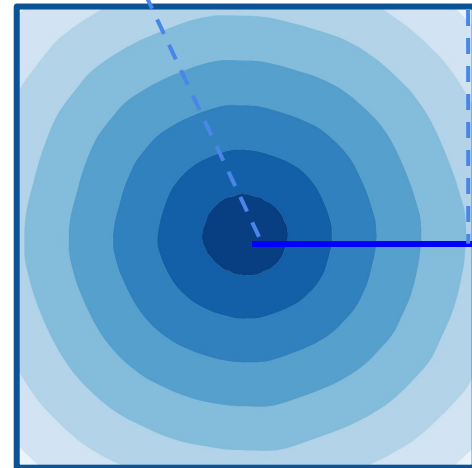
Moment EFN Quark/Gluon
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4$, $p_T \in [500, 550]$ GeV
— Latent Space
- - Fit: $c_1 + c_2 \log(c_3 + y)$
$k = 4, L = 1$
$c_1 = -3.584$
$c_2 = -0.847$
$c_3 = 0.005$

~log(r + c)

**Radial Profile**

Latent Embedding $\phi(y, 0)$ vs Rapidity $y$

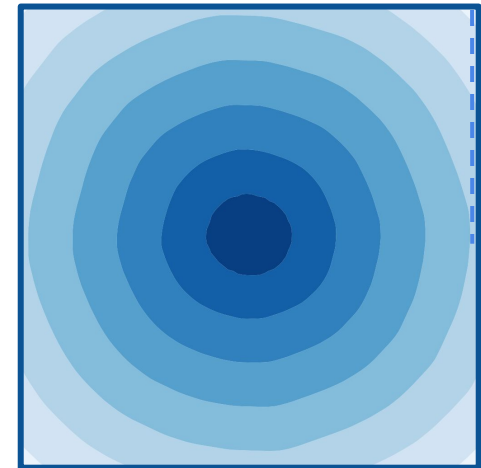Order $k = 4$ $L = 1$

# Part III Summary

By choosing the following new architecture for Deep Sets:

$$\mathcal{O}_k(\mathcal{P}) = F_k \left( \langle \phi^a \rangle_{\mathcal{P}}, \langle \phi^{a_1} \phi^{a_2} \rangle_{\mathcal{P}}, ..., \langle \phi^{a_1} ... \phi^{a_k} \rangle_{\mathcal{P}} \right)$$

We expand the basis of primitive operations to include multiplication, allowing for **streamlined latent spaces** that are easier to extract physics information from!
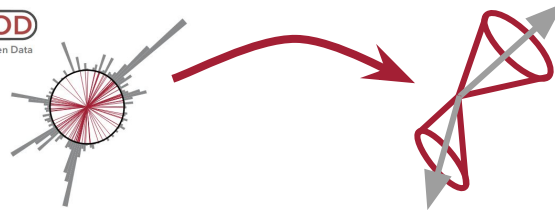
Moment EFN Quark/Gluon
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4, p_T \in [500, 550]$ GeV
— Latent Space
-- Fit: $c_1 + c_2 \log(c_3 + y)$
$k = 4, L = 1$
$c_1 = -3.584$
$c_2 = -0.847$
$c_3 = 0.005$

**~log(r + c)**

**Radial Profile**

Order $k = 4$ $L = 1$

Rikab Gambhir – SLAC – 28 May 2024

# Part IV: Conclusion

Learning **Uncertainties** the **Frequentist** Way
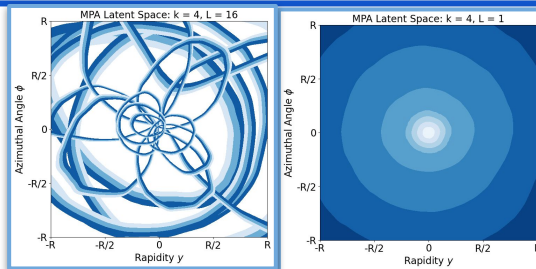
[**RG**, Nachman, Thaler, 2205.05084
**RG**, Nachman, Thaler, 2205.05084]

Can you Hear the **Shape** of a Jet?
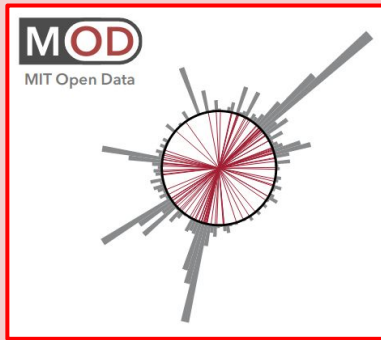
[Ba, Dogra, **RG**, Tasissa, Thaler, 2302.12266]
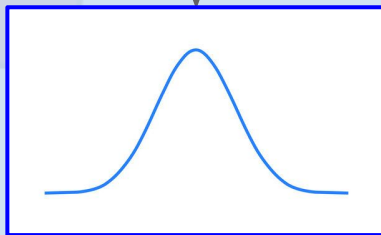
The **Moments** of Clarity: Streamlining Latent Spaces

[**RG**, Osathapan, Thaler, 2403.08854]

MPA Latent Space: k = 4, L = 16

MPA Latent Space: k = 4, L = 1

# Things a **machine** can understand

## Sophisticated Detector Models



## Complicated Point Clouds



## Huge Latent Spaces



Using the Gaussian Ansatz and DV Loss

Using faithful optimal transport

Using Moment Pooling





1) Eff. Rad: 0.07, Ecc: 0.95, $z_\delta$, $z_\oslash$: (0.10, 0.13)
2) Eff. Rad: 0.37, Ecc: 0.78, $z_\delta$, $z_\oslash$: (0.23, 0.06)
3) Eff. Rad: 0.04, Ecc: 0.85, $z_\delta$, $z_\oslash$: (0.25, 0.24)
EMD: 0.038



**Moment EFN Quark/Gluon**
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4, p_T \in [500, 550]$ GeV
— Latent Space
-- Fit: $c_1 + c_2\log(c_3 + y)$
$k = 4, L = 1$
$c_1 = -3.584$
$c_2 = -0.847$
$c_3 = 0.005$

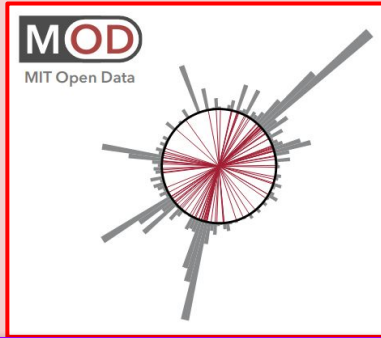**Gaussians**

**Basic Kindergarten Shapes**

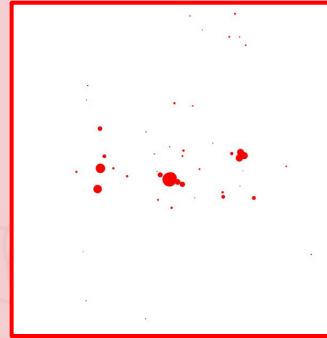**Addition, multiplication, and a few elementary functions**

## Things my **mortal physicist brain** can understand
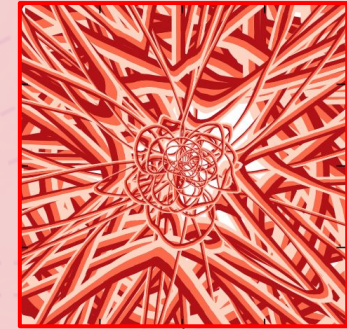
# Things a **machine** can understand

## Sophisticated Detector Models



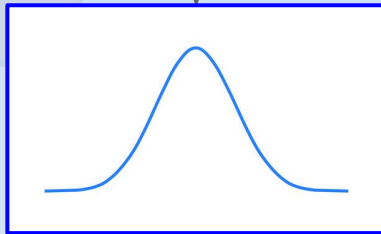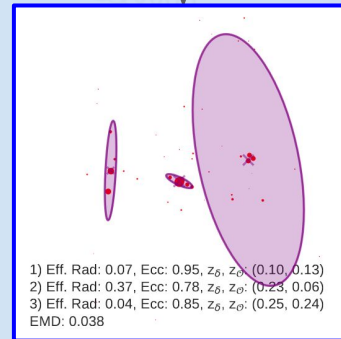## Complicated Point Clouds



## Huge Latent Spaces



**Main Takeaway**: By using purpose-built machine learning architectures and losses, we can make sure we extract the **physics we want** from our machines!

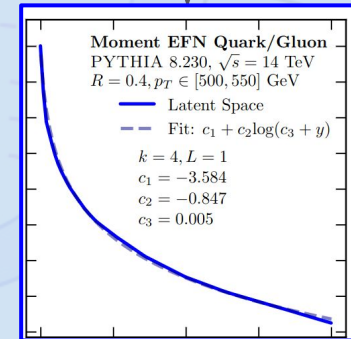Using the Gaussian Ansatz and DV Loss

Using faithful optimal transport

Using Moment Pooling





1) Eff. Rad: 0.07, Ecc: 0.95, $z_\delta$, $z_\phi$: (0.10, 0.13)
2) Eff. Rad: 0.37, Ecc: 0.78, $z_\delta$, $z_\phi$: (0.23, 0.06)
3) Eff. Rad: 0.04, Ecc: 0.85, $z_\delta$, $z_\phi$: (0.25, 0.24)
EMD: 0.038



Moment EFN Quark/Gluon
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4, p_T \in [500, 550]$ GeV
— Latent Space
-- Fit: $c_1 + c_2 \log(c_3 + y)$
$k = 4, L = 1$
$c_1 = -3.584$
$c_2 = -0.847$
$c_3 = 0.005$

**Gaussians**

**Basic Kindergarten Shapes**

Addition, multiplication, and a few elementary functions

## Things my **mortal physicist brain** can understand

Email me questions at rikab@mit.edu!
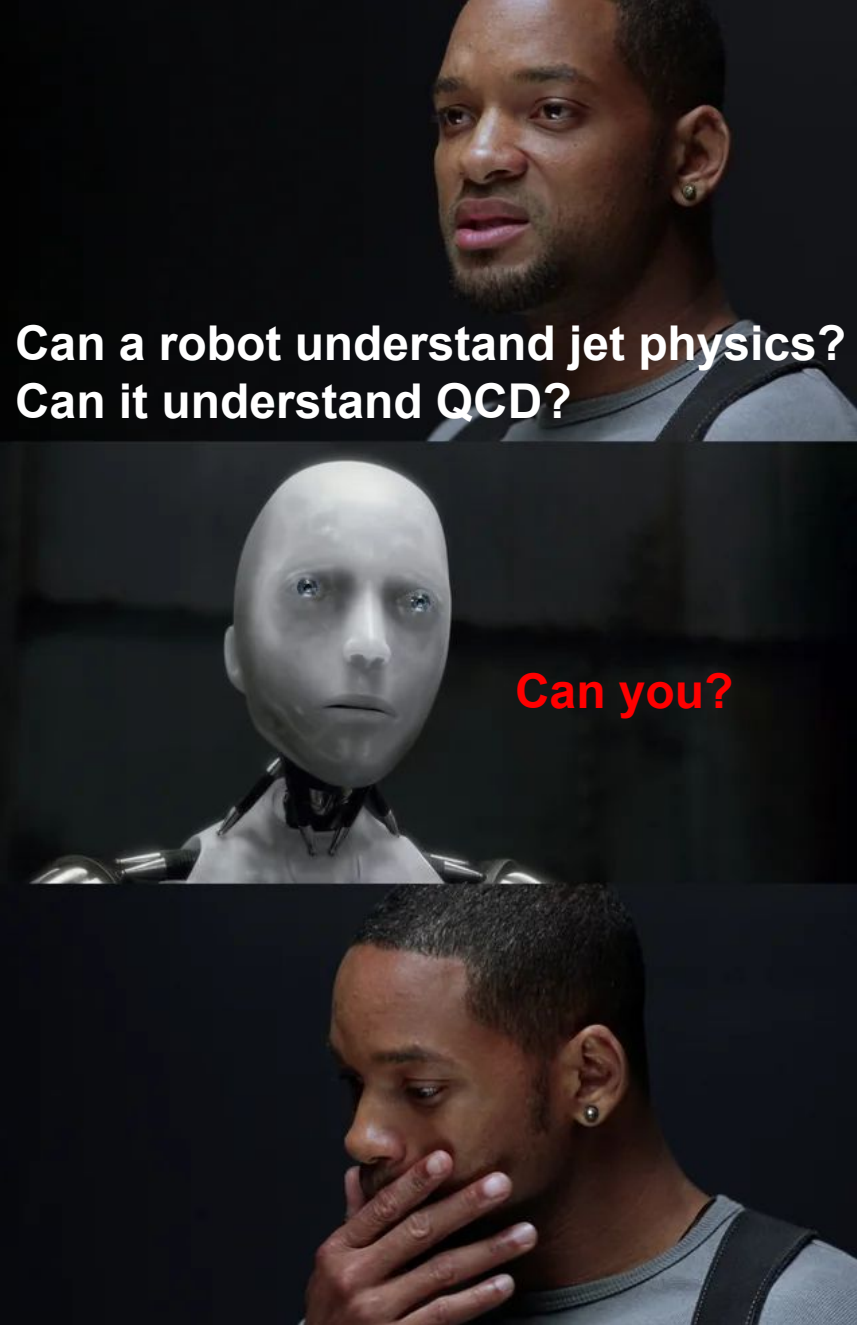
Based on [**RG**, Nachman, Thaler, 2205.03413]
[**RG**, Nachman, Thaler, 2205.05084]
[Ba, Dogra, **RG**, Tasissa, Thaler, 2302.12266]
[**RG**, Thaler, Wu, WIP]
[**RG**, Osathapan, Tasissa, Thaler, 2403.08854]

# Any Questions?



Can a robot understand jet physics?
Can it understand QCD?

Can you?

Rikab Gambhir – SLAC – 28 May 2024

# Backup

# Learning MLC

How do we calculate *f*?

$$f_{\text{MLC}}(x) = \underset{z}{\text{argmax}}\, p_{\text{train}}(x|z)$$

$$= \underset{z}{\text{argmax}}\, \log \underbrace{\frac{p_{\text{train}}(x,z)}{p_{\text{train}}(x)p_{\text{train}}(z)}}_{T(x,z)}$$

The function *T* is the likelihood ratio between *p(x,z)* and *p(x)p(z)*.
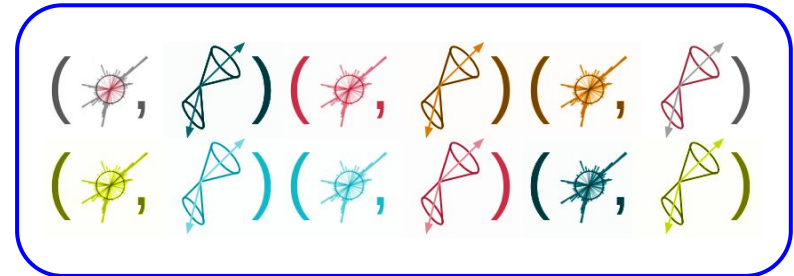
Neyman–Pearson

*T* is the **optimal classifier** between (*x,z*) pairs and shuffled (*x,z*) pairs!

**Class P**



**Class Q**



Classify between *P* and *Q*!

Rikab Gambhir – SLAC – 28 May 2024

# Aside: Mutual information

A measure for non-linear interdependence is the **Mutual Information**:

$$I(X; Z) = \int \mathrm{d}x\, \mathrm{d}z\, p(x, z) \log \frac{p(x, z)}{p(x)\, p(z)}$$
$$= \mathbb{E}_{\mathrm{train}} T(X, Z)$$

Answers the question: How much information, in terms of bits, do you learn about $Z$ when you measure $X$ (or vice versa)?

When doing calibration this way, we get a measure of the **correlation** between $X$ and $Z$, *for free*.
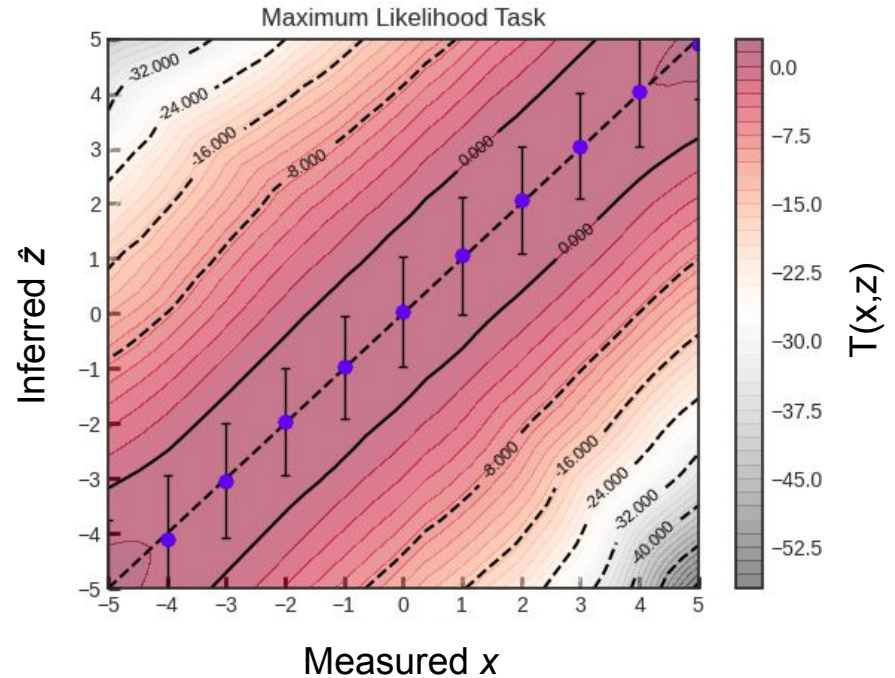
# Example 1: Gaussian Calibration Problem
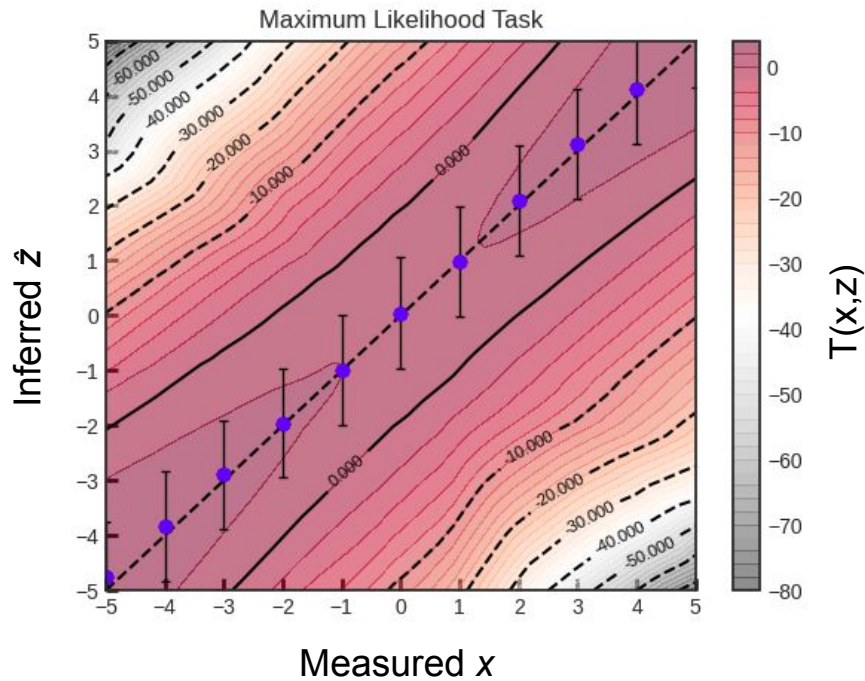
Gaussian noise model: $p(x|z) \sim N(z, 1)$

**Model**:

- The *A*, *B*, *C*, and *D* networks are each Dense networks with 4 layers of size 32
- ReLU activations
- All parameters have an L2 regularization ($\lambda$ = 1e-6)
- The D network regularization slowly increased to ($\lambda_D$ = 1e-4)

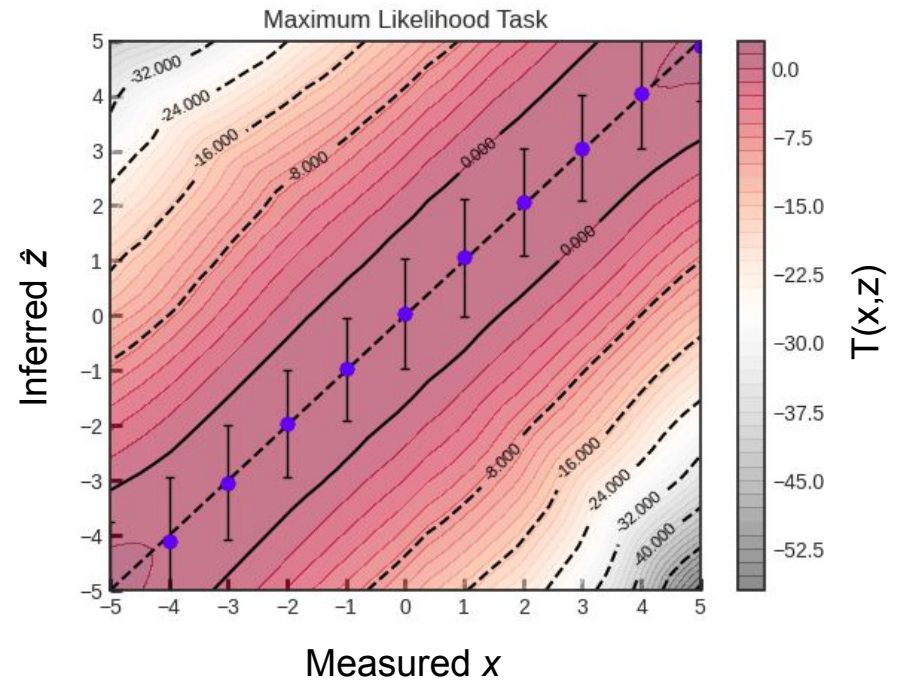Learned mutual information of 1.05 natural bits

Reproduces the expected maximum likelihood outcome and the correct resolution!



Maximum Likelihood Task

# Example 1 - Prior Independence

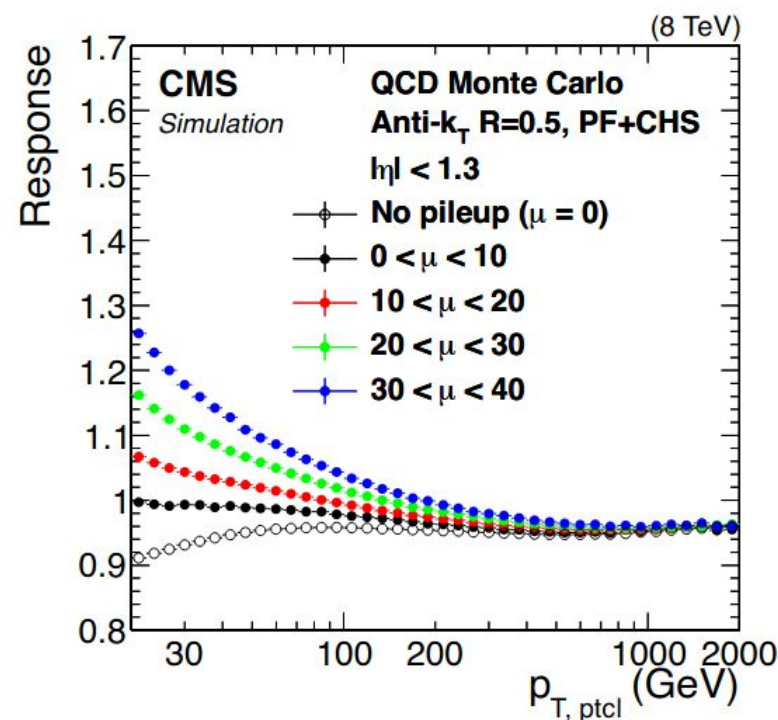

P(z) ~ N(0, 2.5)

P(z) ~ U(-5, 5)

# Example 3: Jet Energy Calibrations

Measure a set particle flow candidates *x* in the detector. What is the underlying jet $p_T$, *x*, and its uncertainty?

Define the **jet energy scale (JES)** and **jet energy resolution (JER)** as the ratio of the underlying (GEN) jet $p_T$ (resolution) to the measured total (SIM) jet $p_T$

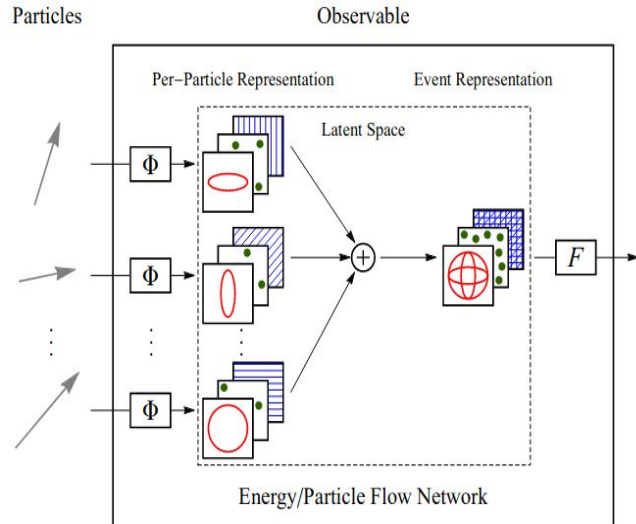$$\hat{p}_T \equiv \text{JEC} \times p_{T,\text{SIM}} \approx p_{T,\text{GEN}}$$

$$\hat{\sigma}_{p_T} = \text{JER} \times p_{T,\text{SIM}}$$

# Example 3: Models

- **DNN**: $X$ = (Jet $p_T$, Jet $\eta$, Jet φ), Dense Neural Network
- **EFN**: $X$ = {(PFC $p_T$, PFC $\eta$, PFC φ)}, Energy Flow Network
- **PFN**: $X$ = {(PFC $p_T$, PFC $\eta$, PFC φ)}, Particle Flow Network
- **PFN-PID**: $X$ = {(PFC $p_T$, PFC $\eta$, PFC φ, PFC PID)}, Particle Flow Network

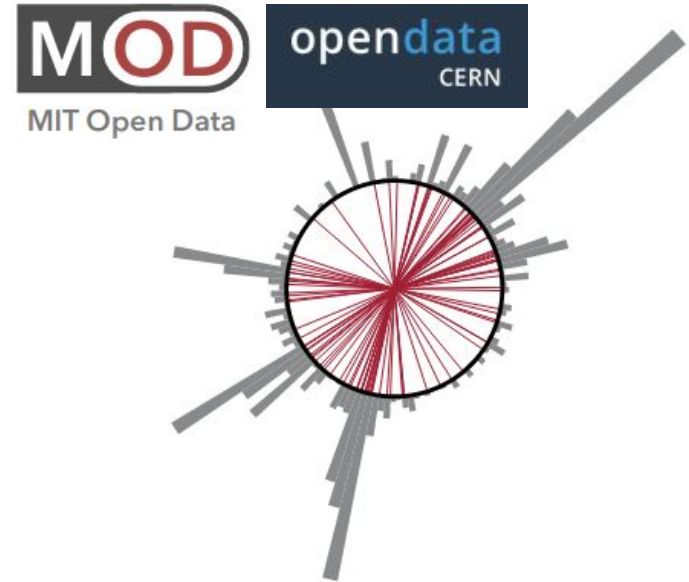For each model, $A(x)$, $B(x)$, $C(x,z)$, and $D(x)$ are all of the same type.



Permutation-invariant function of point clouds
For EFN's, manifest IRC Safety

Details on hyperparameters can be found in [**RG**, Nachman, Thaler, PRL 129 (2022) 082001]

Rikab Gambhir – SLAC – 28 May 2024

# Example 3: Jet Dataset

Using CMS Open Data:
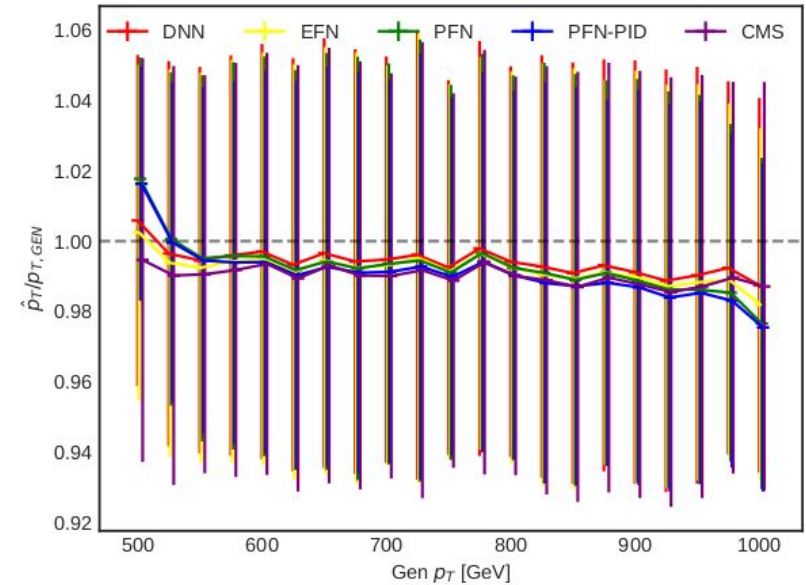
- *CMS2011AJets* Collection, SIM/GEN QCD Jets (AK 0.5)
- Select for jets with 500 GeV < Gen $p_T$ < 1000 GeV, $|\eta|$ < 2.4, quality ≥ 2
- Select for jets with ≤ 150 particles
- Jets are rotated such that jet axis is centered at (0,0)
- Train on 100k jets

Rikab Gambhir – SLAC – 28 May 2024

# Jet Energy Scales

For jets with a true $p_T$ between 695-705 GeV, we should expect well-trained models to predict 700 GeV on average!

| Model | Gaussian Fit [GeV] |
|---|---|
| DNN | 695 ± 38.2 |
| EFN | 692 ± 37.7 |
| PFN | 702 ± 37.4 |
| PFN-PID | 693 ± 35.9 |
| CMS Open Data | 695 ± 37.4 |



Close to 1.00 – unbiased estimates!

# Data Based Calibration

"What if my detector simulation *p(x|z)* is imperfect"?

Given a *bad* simulator $p_{SIM}(x|z)$, we can correct it by matching it to data:

$$\hat{p}(x_D|z_T) = p_{sim}(h(x_D)|z_T)|h'(x_D)|$$

Where

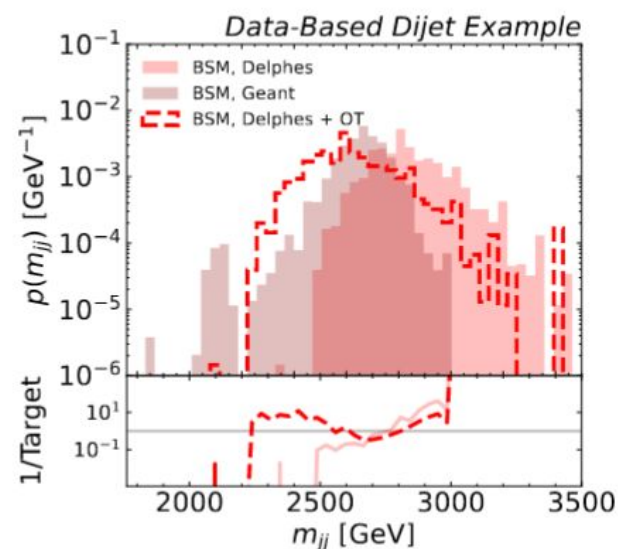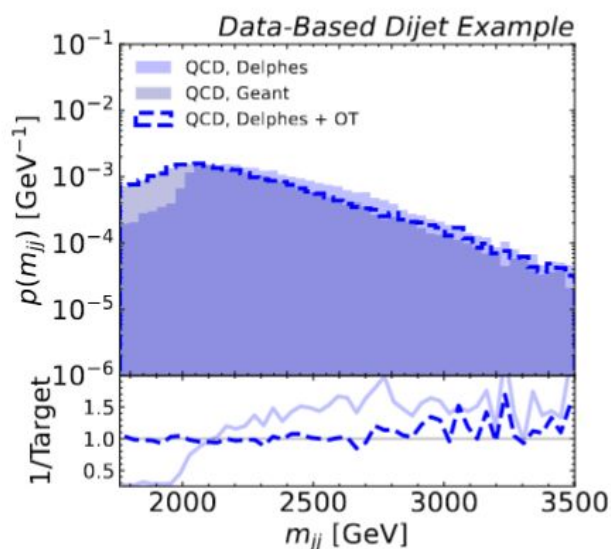$$h(x_D) = P_{data}^{-1}(P_{sim}(x_D))$$

The function *h* "optimally transports" points to where they belong and reweights them.

# Data Based Calibration

**BUT!** There is a cost. We have to give up prior independence.



"Fixing" the Delphes simulation to match Geant4 works when trained on **Prior 1 (QCD)**, but fails miserably when applied to **Prior 2 (BSM)**, despite being the same detector simulation!
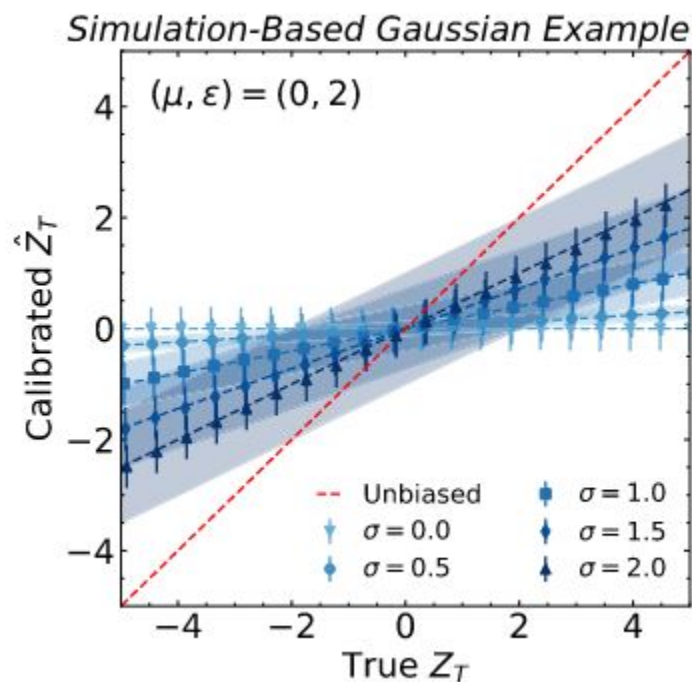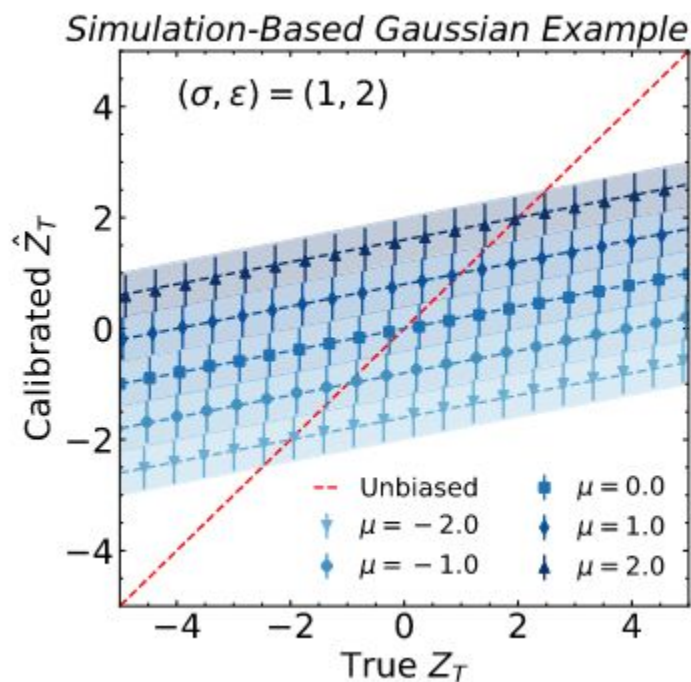
No (known) method of prior independent DBC, but no proof it is impossible!

# Prior dependence of MSE

MSE fits for a gaussian noise model, for different choices of $z$ prior.
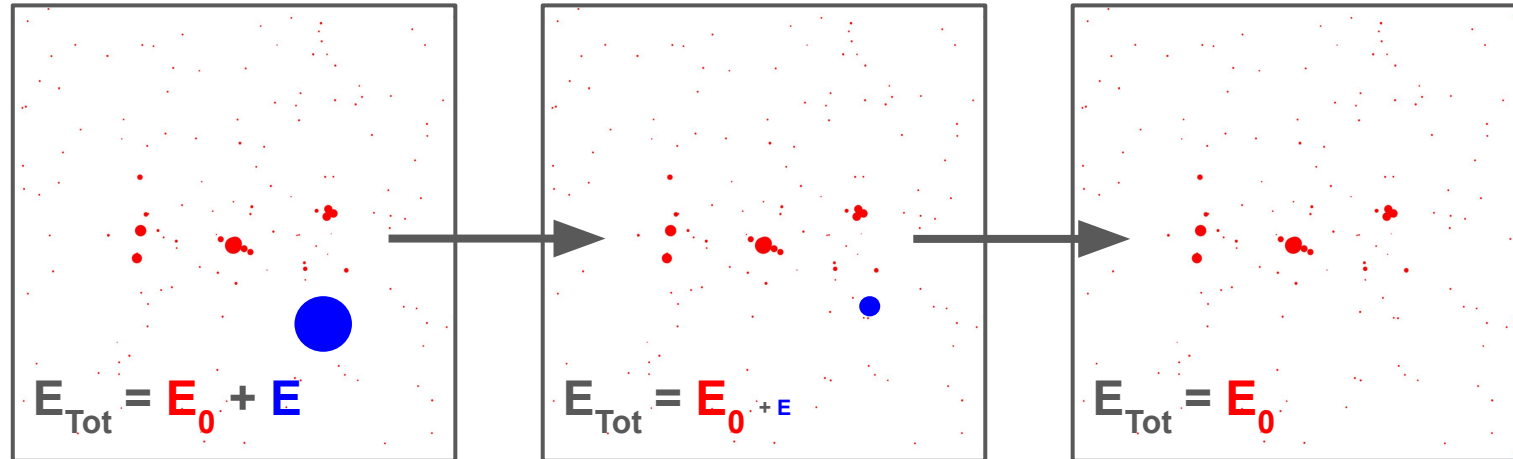
Left: Different choices of mean               Right: Different choices of width

# Mathematical Details – Topology

Definition (The **Weak\* Topology**): A sequence of measures converges if all of their expectation values converge, as real numbers.



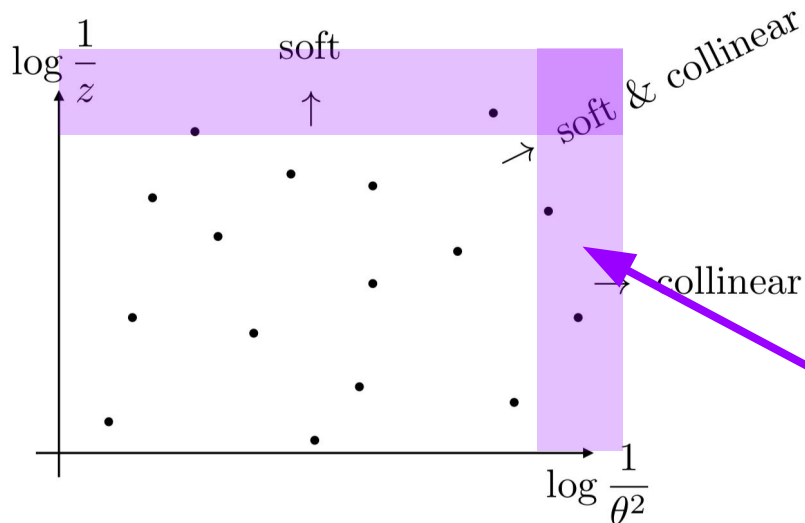$$E_{Tot} = E_0 + E \qquad E_{Tot} = E_{0+E} \qquad E_{Tot} = E_0$$

Definition (**Weak\* Continuity**): An observable $\mathcal{O}(\mathcal{E})$ is continuous with respect to energy flows if, for any sequence of measures $\mathcal{E}_n$ that converges to $\mathcal{E}$, the sequence of real numbers $\mathcal{O}(\mathcal{E}_n)$ converges to $\mathcal{O}(\mathcal{E})$.

# Topology ⟺ IRC-Safety

Two ways to change the expectation values of an energy flow:

1. Change a particle's energy slightly, or add a low-energy particle - **IR**
2. Move a particle's position slightly, or split particles in two - **C**



An observable $\mathcal{O}$ is continuous if it changes only slightly under the above perturbations.

The regions of phase space causing IRC divergences is suppressed — $\mathcal{O}$ is **IRC-Safe**!

# Mathematical Details - **Geometry**

When are two events similar? We need a metric to compare!

Properties we want:

1.  … is non-negative, non-degenerate, symmetric and finite
2.  … is weak* continuous (IRC-safe)
3.  … lifts the detector metric **faithfully**

Explained shortly!

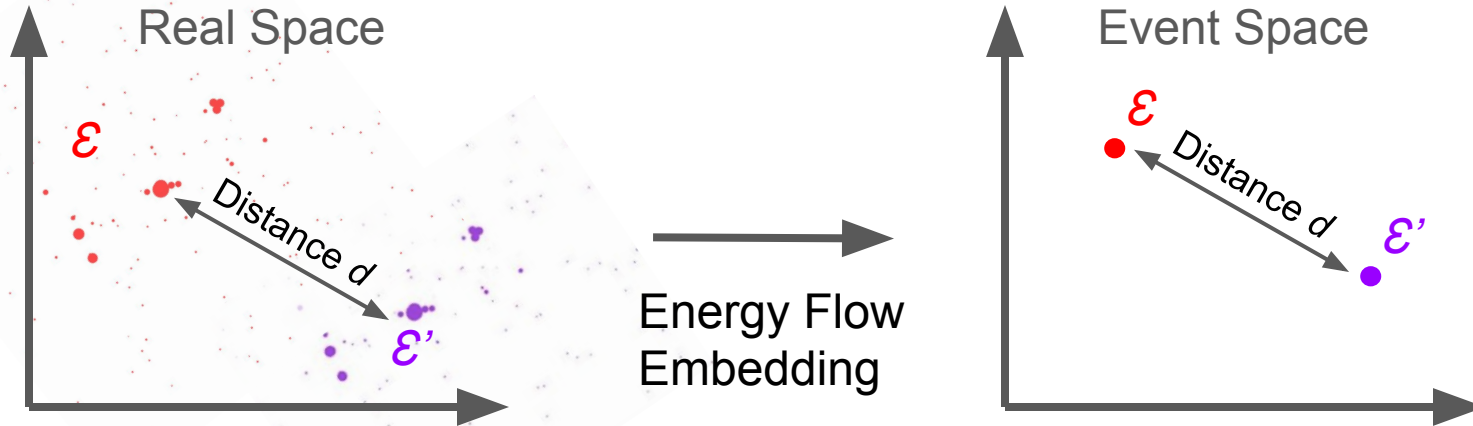The only* metric on distributions satisfying the above is the **Wasserstein Metric**:

$$\mathrm{EMD}^{(\beta,R)}(\mathcal{E},\mathcal{E}') = \min_{\pi \in \mathcal{M}(\mathcal{X}\times\mathcal{X})} \left[ \frac{1}{\beta R^\beta} \left\langle \pi, d(x,y)^\beta \right\rangle \right] + |\Delta E_{\mathrm{tot}}|$$

$$\pi(\mathcal{X}, Y) \le \mathcal{E}'(Y), \quad \pi(X, \mathcal{X}) \le \mathcal{E}(X), \quad \pi(\mathcal{X},\mathcal{X}) = \min(E_{\mathrm{tot}}, E'_{\mathrm{tot}})$$

*There exist other metrics on distributions that are faithful only for very specific real-space distance norms, but we want them all!
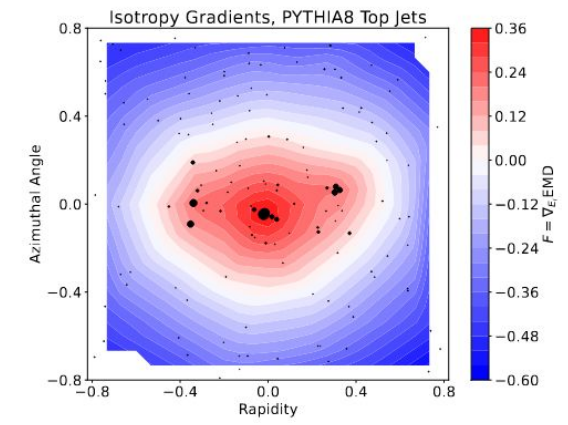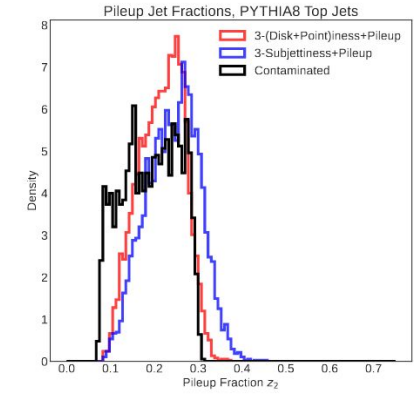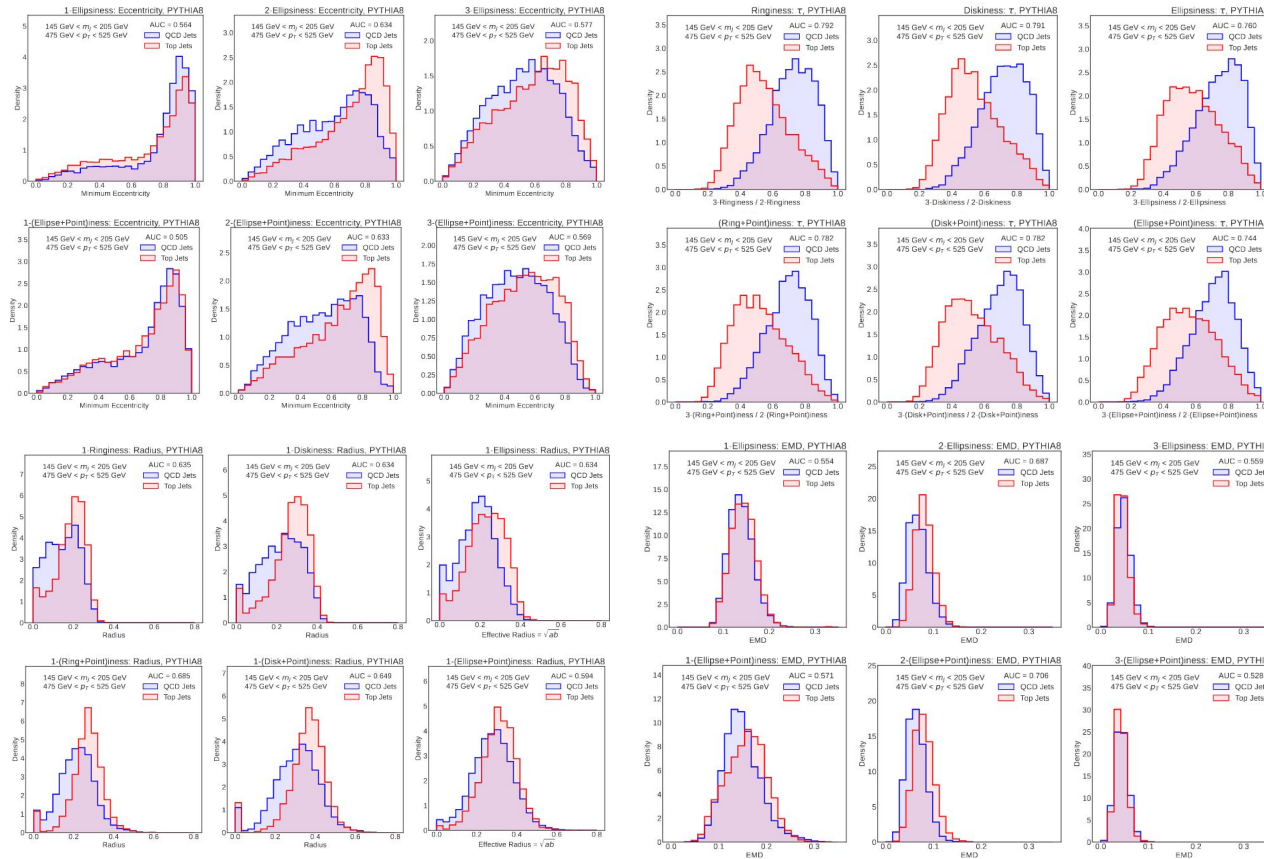
# The Importance of Being Faithful*

Real Space

$\mathcal{E}$

Distance $d$

$\mathcal{E}'$

Energy Flow
Embedding

Event Space

$\mathcal{E}$

Distance $d$

$\mathcal{E}'$

A metric on events is **faithful** if, whenever two otherwise identical events $\mathcal{E}$ and $\mathcal{E}'$ are separated in real space by a distance $d$, the distance between the events is also $d$. Or any predetermined invertible function of $d$

Only the **Wasserstein Metric** does this! → Can use to build **event** and **jet shapes** (old and new)!

Faithfulness also ensures very nice **numerical properties**, including no vanishing or exploding gradients.

# New IRC-Safe Observables



… **Lots** of extractable information!
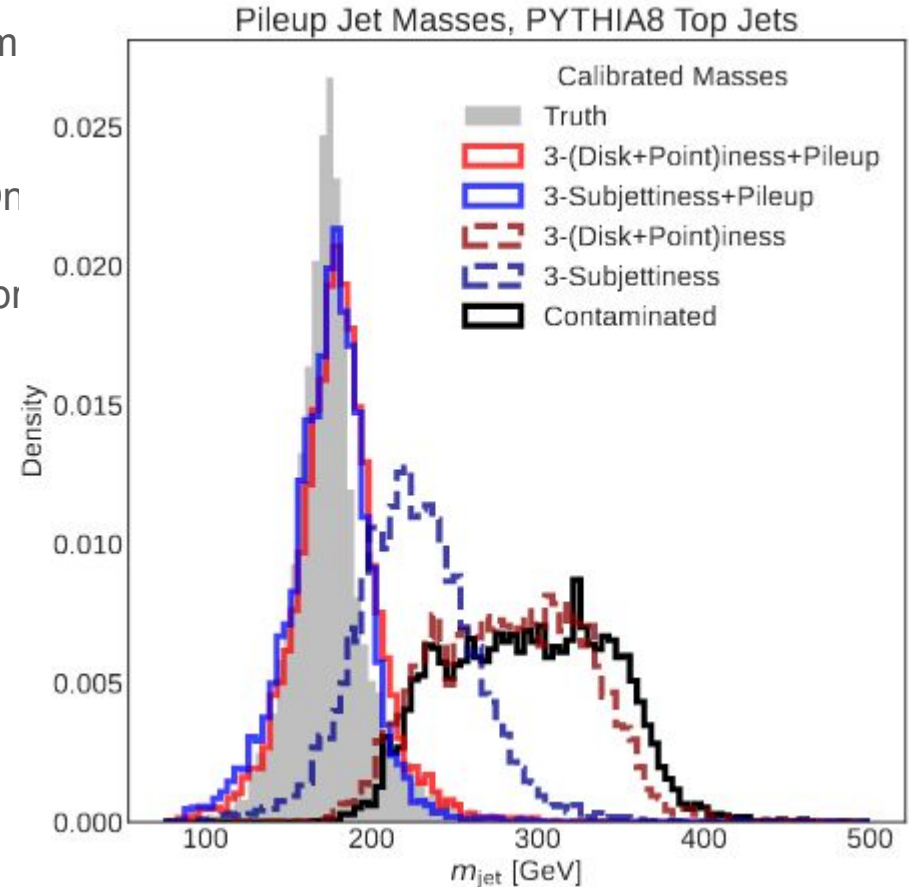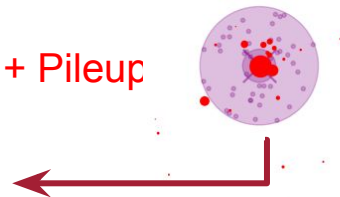
# **Automatic** Grooming with Shapes

Use shapes to approximate events and extract m
background with floating weight!

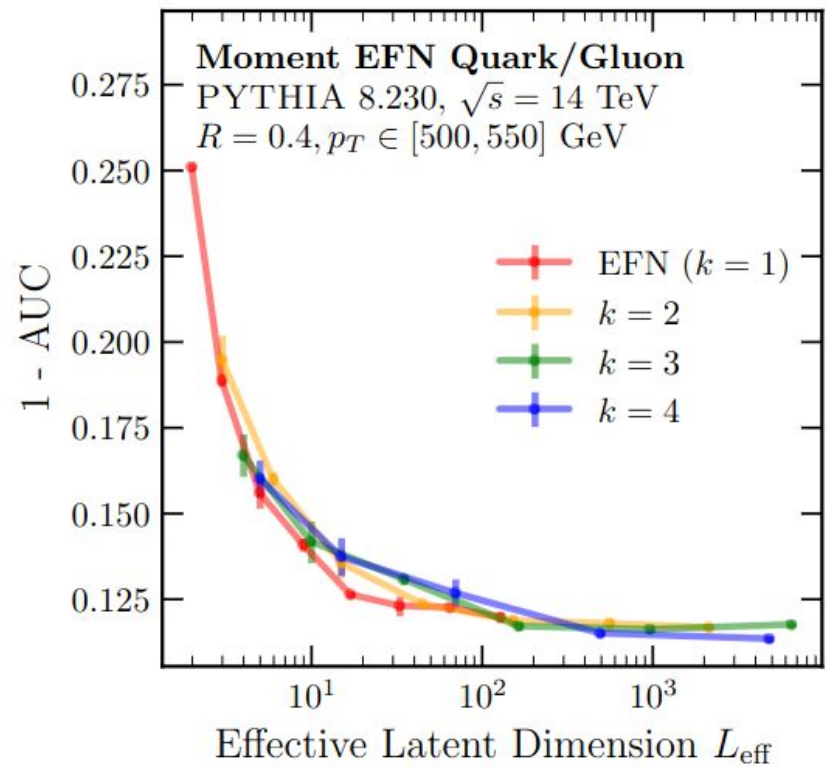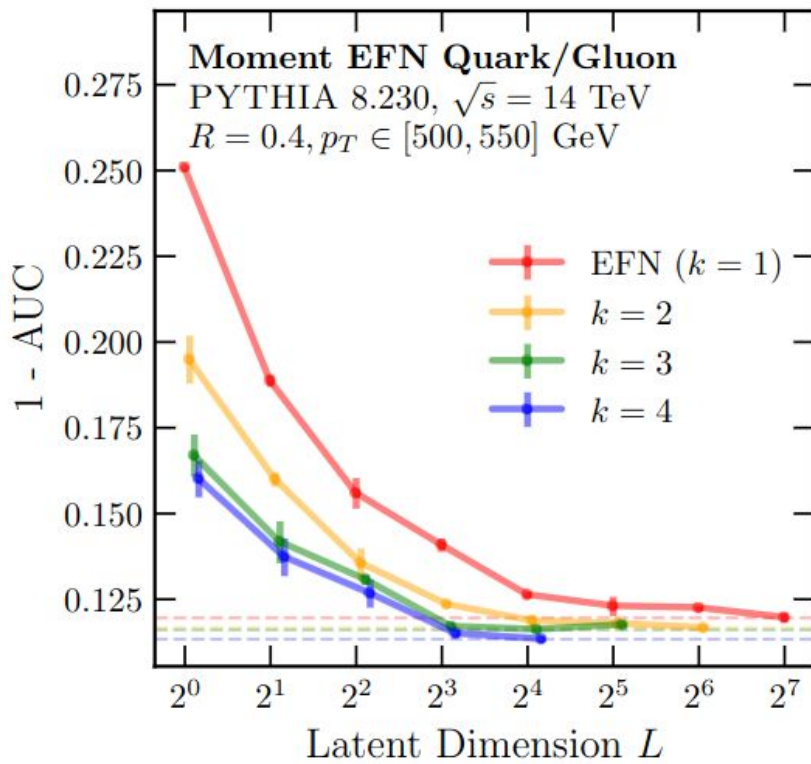*No external hyperparameters*, unlike softdrop. On

Contaminate top jets with 5-30% extra energy spr

Consider 4 shapes:

- 3-Subjettiness
- 3-Subjettiness + Pileup
- 3-(Disk+Point)iness
- 3-(Disk+Point)iness + Pileup



Pileup Jet Masses, PYTHIA8 Top Jets

Calibrated Masses
- Truth
- 3-(Disk+Point)iness+Pileup
- 3-Subjettiness+Pileup
- 3-(Disk+Point)iness
- 3-Subjettiness
- Contaminated

Density

$m_{jet}$ [GeV]

Can also consider ellipses instead of disks – only marginally better performance

Rikab Gambhir – SLAC – 28 May 2024

Quark/Gluon
Analytic Latent Space
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4, p_T \in [500, 550]$ GeV

$k = 1$, AUC $= 0.730$
$k = 2$, AUC $= 0.784$
$k = 3$, AUC $= 0.816$
$k = 4$, AUC $= 0.821$
$L = 1$ Mom. EFN, AUC $= 0.835$
$k = 4$ DNN, $c_3 \to 0$, AUC $= 0.799$

# Attention is all you need



[Vaswani et. al., 1706.03762]

$$\sim \quad \langle \phi^{a_1} \phi^{a_2} \rangle_{\mathcal{P}}$$

Can view moment pooling as a multi-headed self-attention-like mechanism
Each latent variable weights each other latent variable

Rikab Gambhir – SLAC – 28 May 2024

Energy Deposits

Rikab Gambhir – SLAC – 28 May 2024