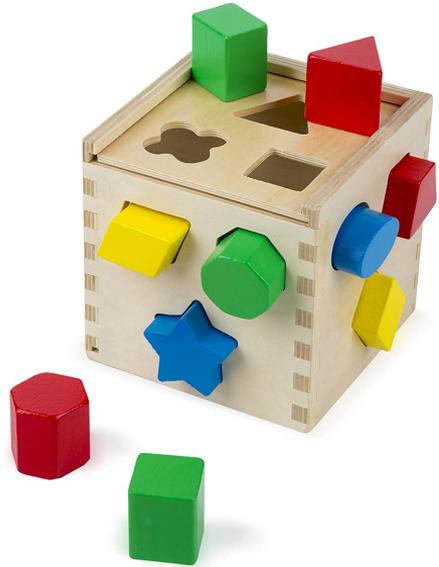


# Can You Hear the Shape of A Jet?



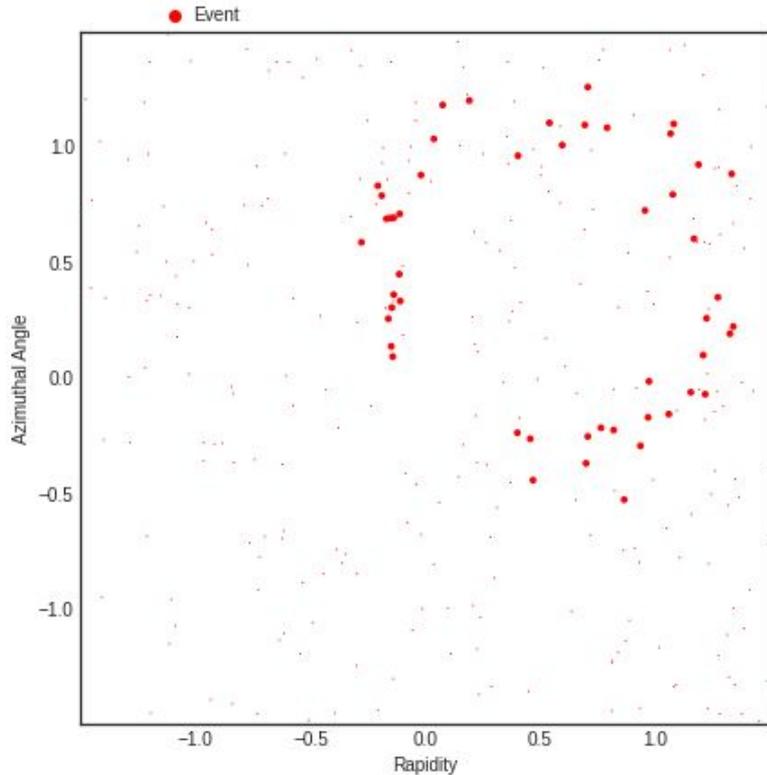
*An IAIFI Story*

Rikab Gambhir

With Akshunna S. Dogra (FI  ) ,  
Demba Ba (FI  ) ,  
& Jesse Thaler (FI )



# Fundamental Question: What shape is this?

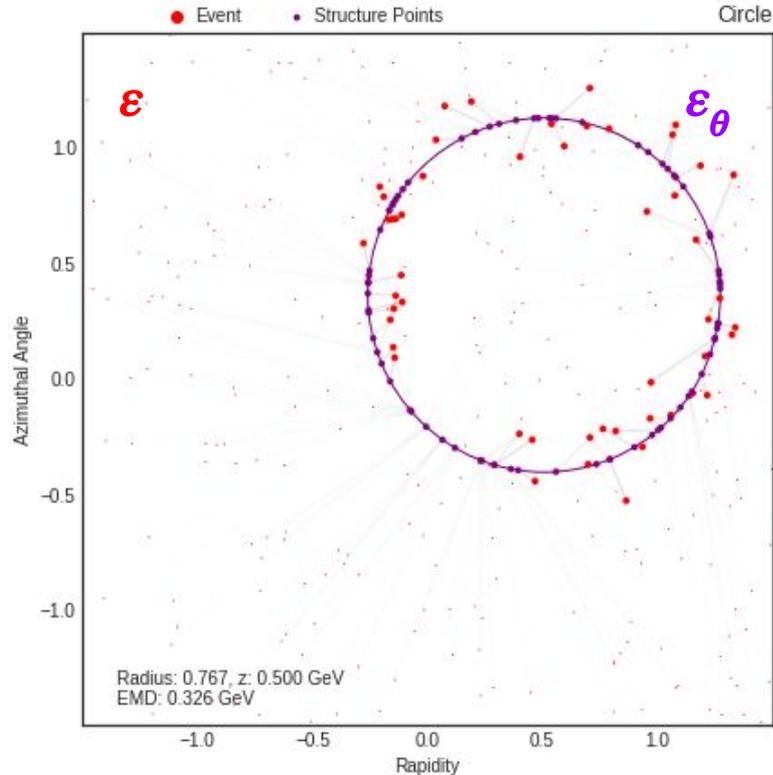


Pictured: (Fake) event that you might have measured at the LHC

Red dots are detector hits on a patch of the LHC cylinder, weighted by energy

**Goal:** Construct an observable  $\mathcal{O}$  that generically answers this question!

# Fundamental Question: What shape is this?

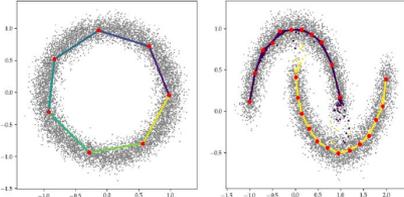


Using the **SHAPER** framework and optimal transport .....

$$\mathcal{O}_{\mathcal{M}}(\mathcal{E}) = \min_{\mathcal{E}'_\theta \in \mathcal{M}} \text{EMD}(\mathcal{E}, \mathcal{E}'_\theta)$$
$$\theta = \operatorname{argmin}_{\mathcal{E}'_\theta \in \mathcal{M}} \text{EMD}(\mathcal{E}, \mathcal{E}'_\theta)$$

**Circle** with radius 0.767, center (0.50, 0.36) and a “circle-ness” value of 0.32.

**Yes, you CAN hear the shape of a jet!**

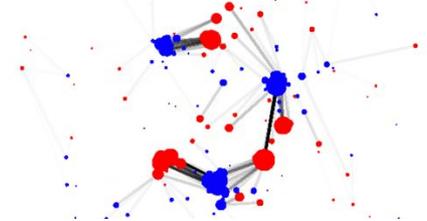


Piecewise-Linear Manifold Approximation with K-Deep Simplicies (KDS, [2012.02134](#))

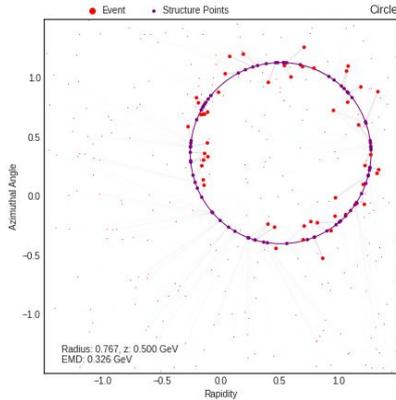
Ai

FI

fi



Well-Defined Metric on Particle Collisions using Energy Mover's Distance (EMD, [2004.04159](#))



## SHAPER: Learning the Shape of Collider Events

$$\mathcal{O}_{\mathcal{M}}(\mathcal{E}) = \min_{\mathcal{E}'_{\theta} \in \mathcal{M}} \text{EMD}(\mathcal{E}, \mathcal{E}'_{\theta})$$

$$\theta = \operatorname{argmin}_{\mathcal{E}'_{\theta} \in \mathcal{M}} \text{EMD}(\mathcal{E}, \mathcal{E}'_{\theta})$$

Framework for defining and calculating useful observables for collider physics!

# Building SHAPER

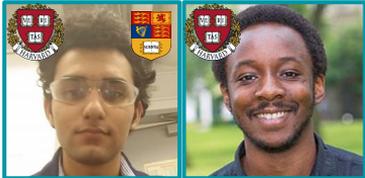
Key Component: The Loss function! Step 1: Manifold Learning

$$\mathcal{L}_R(\mathcal{E}, \mathcal{E}') = \min_{\pi_{ij} \geq 0} \left[ \sum_{i=1}^M \sum_{j=1}^{M'} \pi_{ij} \frac{|x_i - x'_j|}{R} \right],$$

where  $\sum_{i=1}^M \pi_{ij} = 1$

Dogra

Ba



# Ai

K-Deep Simplices,  
Dictionary Learning, &  
Manifold Learning

# Building SHAPER

Key Component: The Loss function! Step 2: Physical Principles

$$\mathcal{L}_R(\mathcal{E}, \mathcal{E}') = \min_{\pi_{ij} \geq 0} \left[ \sum_{i=1}^M \sum_{j=1}^{M'} \pi_{ij} \frac{|x_i - x'_j|}{R} \right] + \left| \sum_{i=1}^M z_i - \sum_{j=1}^{M'} z'_j \right|,$$

where  $\sum_{i=1}^M \pi_{ij} \leq z'_j, \sum_{j=1}^{M'} \pi_{ij} \leq z_i, \sum_{i,j}^{M,M'} \pi_{ij} = \min \left( \sum_{i=1}^M z_i, \sum_{j=1}^{M'} z'_j \right)$

Dogra

Ba



Ai

K-Deep Simplices,  
 Dictionary Learning, &  
 Manifold Learning

fi

IRC Safety,  
 Unclustered Radiation, &  
 Wasserstein Geometry

Gambhir Thaler



# Building SHAPER

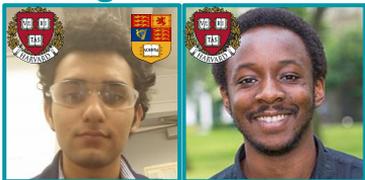
Key Component: The Loss function! Step 3: Synthesis

$$\mathcal{L}_R(\mathcal{E}, \mathcal{E}') = \min_{\pi_{ij} \geq 0} \left[ \sum_{i=1}^M \sum_{j=1}^{M'} \pi_{ij} \frac{|x_i - x'_j|}{R} \right] + \left| \sum_{i=1}^M z_i - \sum_{j=1}^{M'} z'_j \right|,$$

where  $\sum_{i=1}^M \pi_{ij} \leq z'_j$ ,  $\sum_{j=1}^{M'} \pi_{ij} \leq z_i$ ,  $\sum_{i,j}^{M,M'} \pi_{ij} = \min \left( \sum_{i=1}^M z_i, \sum_{j=1}^{M'} z'_j \right)$

Dogra

Ba



# Ai

K-Deep Simplices,  
Dictionary Learning, &  
Manifold Learning

+

# fi

IRC Safety,  
Unclustered Radiation, &  
Wasserstein Geometry

Gambhir Thaler



# Building SHAPER



function! Step 3: Synthesis

Nolte

Williams

Kitouni

Conversations with:


$$\sum_{i=1}^M \sum_{j=1}^M \pi_{ij} \frac{|x_i - x_j|}{R}$$

- Connections to **LHCb Trigger** development
- Potential applications to **future colliders**
- Discussion of **implementation details**

Dogra

Ba

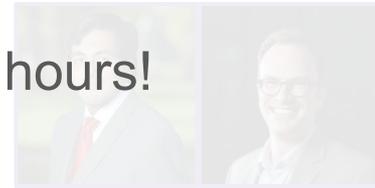


Connections made at the **IAIFI** penthouse and coffee hours!

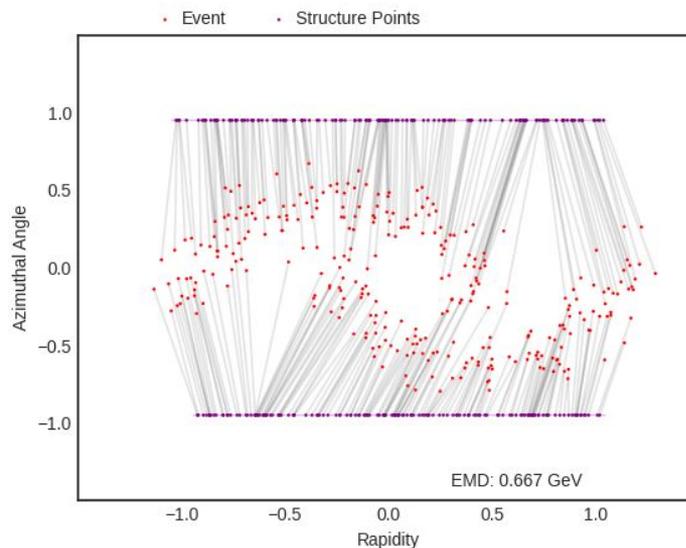
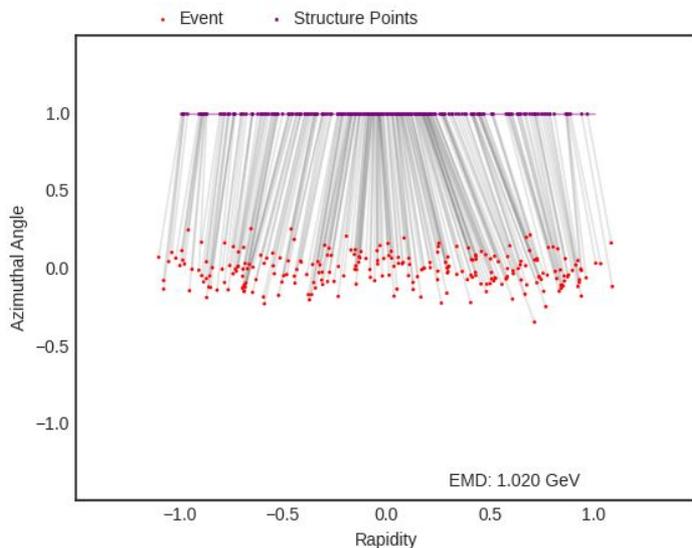
Deep Simplicial  
Dictionary Learning, &  
Manifold Learning

IRG Society  
Unclustered Radiation, &  
Wasserstein Geometry

Gambhir Thaler



# Fun Animations

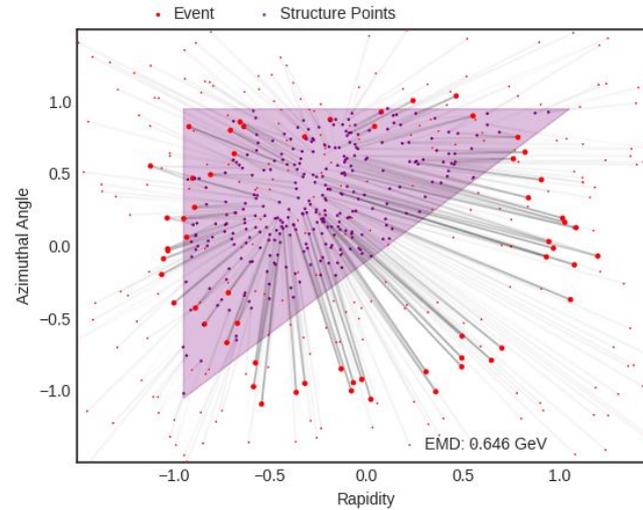
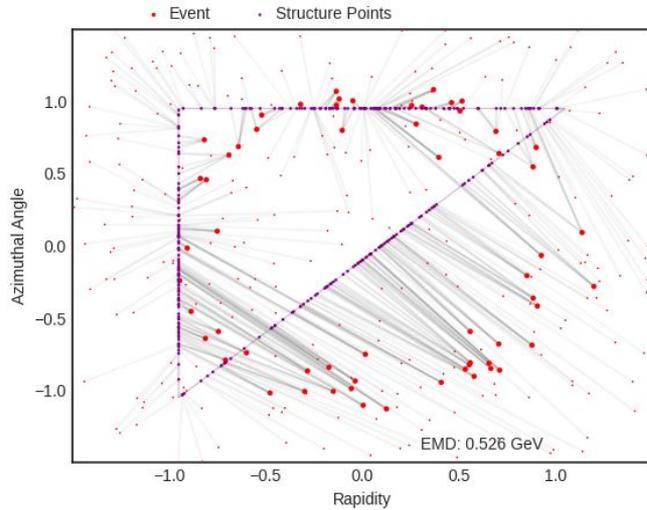


**Red:** Event  $Y$

**Purple:** Shape  $\varepsilon_\theta$  with structure points  $a_i$

**Grey:** Matrix  $x_{ij}$  connecting  $y$ 's and  $a_i$ 's

# Fun Animations Cont'd



**Red:** Event  $Y$

**Purple:** Shape  $\varepsilon_\theta$  with structure points  $a_i$

**Grey:** Matrix  $f_{ij}$  connecting  $y$ 's and  $a_i$ 's

**Left:**  $\varepsilon_\theta =$   EMD = 0.245

**Right:**  $\varepsilon_\theta =$   EMD = 0.279

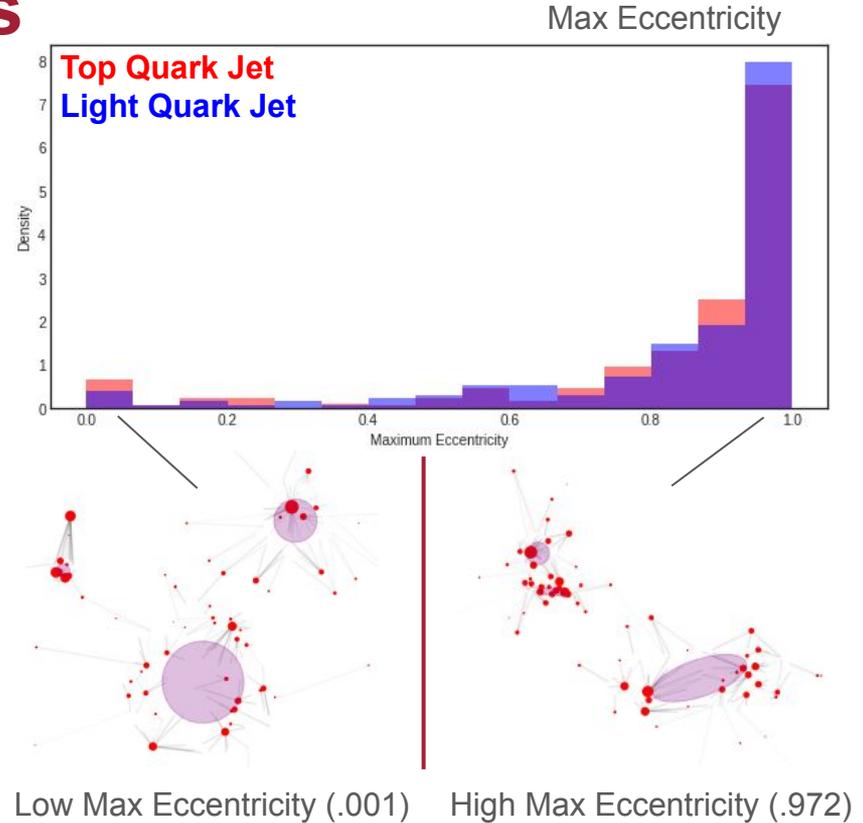
# New IRC-Safe Observables

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

e.g. ***N*-Ellipsiness+Pileup** as a jet algorithm.

- Learn jet centers
- Dynamic jet radii (no  $R$  hyperparameter)
- Dynamic **eccentricities** and angles
- Dynamic jet energies
- Uniform Pileup Subtraction
- Learned parameters for discrimination

Can design custom specialized jet algorithms to learn jet substructure!



# Other Developments: Statistics in Physics

Thaler



Gambhir



Nachman



## Machine Learning Calibrations (2205.05084)

Bias and Priors in Machine Learning Calibrations for High Energy Physics

Rikab Gambhir,<sup>1,2,\*</sup> Benjamin Nachman,<sup>3,4,†</sup> and Jesse Thaler<sup>1,2,‡</sup>

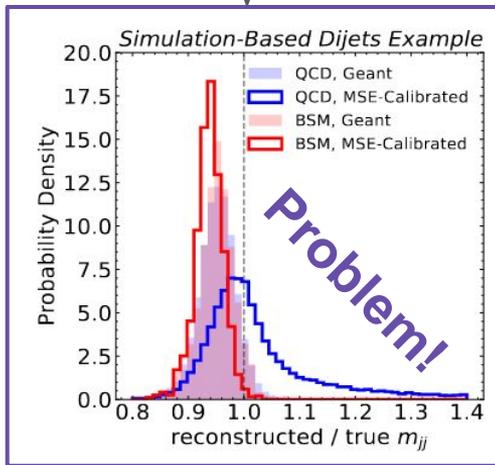
<sup>1</sup>Center for Theoretical Physics, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

<sup>2</sup>The NSF AI Institute for Artificial Intelligence and Fundamental Interactions

<sup>3</sup>Physics Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

<sup>4</sup>Berkeley Institute for Data Science, University of California, Berkeley, CA 94720, USA

Machine learning offers an exciting opportunity to improve the calibration of nearly all reconstructed objects in high-energy physics detectors. However, machine learning approaches often depend on the spectra of examples used during training, an issue known as prior dependence. This is an undesirable property of a calibration, which needs to be applicable in a variety of environments. The purpose of this paper is to explicitly highlight the prior dependence of some machine learning-based calibration strategies. We demonstrate how some recent proposals for both simulation-based and data-based calibrations inherit properties of the sample used for training, which can result in biases for downstream analyses. In the case of simulation-based calibration, we argue that our recently proposed Gaussian Ansatz approach can avoid some of the pitfalls of prior dependence, whereas prior-independent data-based calibration remains an open problem.



Ai

## Gaussian Ansatz Statistical Framework (2205.03413)

Learning Uncertainties the Frequentist Way:  
Calibration and Correlation in High Energy Physics

Rikab Gambhir,<sup>1,2,\*</sup> Benjamin Nachman,<sup>3,4,†</sup> and Jesse Thaler<sup>1,2,‡</sup>

<sup>1</sup>Center for Theoretical Physics, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

<sup>2</sup>The NSF AI Institute for Artificial Intelligence and Fundamental Interactions

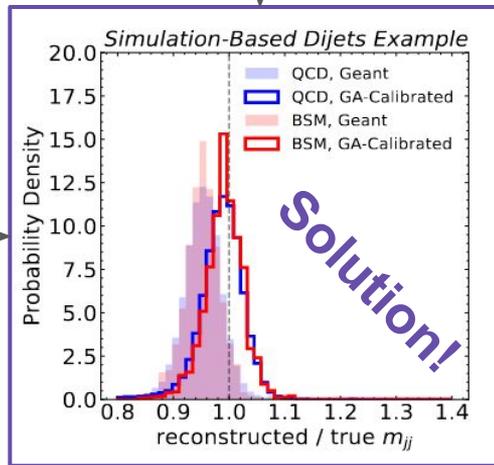
<sup>3</sup>Physics Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

<sup>4</sup>Berkeley Institute for Data Science, University of California, Berkeley, CA 94720, USA

Calibration is a common experimental physics problem, whose goal is to infer the value and uncertainty of an unobservable quantity  $Z$  given a measured quantity  $X$ . Additionally, one would like to quantify the extent to which  $X$  and  $Z$  are correlated. In this paper, we present a machine learning framework for performing frequentist maximum likelihood inference with Gaussian uncertainty estimation, which also quantifies the mutual information between the unobservable and measured quantities. This framework uses the Donsker-Varadhan representation of the Kullback-Leibler divergence – parametrized with a novel Gaussian Ansatz – to enable a simultaneous extraction of the maximum likelihood values, uncertainties, and mutual information in a single training. We demonstrate our framework by extracting jet energy corrections and resolution factors from a simulation of the CMS detector at the Large Hadron Collider. By leveraging the high-dimensional feature space inside jets, we improve upon the nominal CMS jet resolution by upwards of 15%.

continues!

fi



# Other Developments: MSRP



## Gaussian Ansatz Statistical Framework (2205.03413)

### Learning Uncertainties the Frequentist Way: Calibration and Correlation in High Energy Physics

Rikab Gambhir,<sup>1,2,\*</sup> Benjamin Nachman,<sup>3,4,†</sup> and Jesse Thaler<sup>1,2,‡</sup>

<sup>1</sup>Center for Theoretical Physics, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

<sup>2</sup>The NSF AI Institute for Artificial Intelligence and Fundamental Interactions

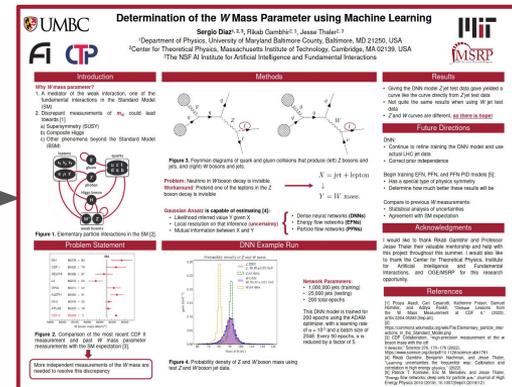
<sup>3</sup>Physics Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

<sup>4</sup>Berkeley Institute for Data Science, University of California, Berkeley, CA 94720, USA

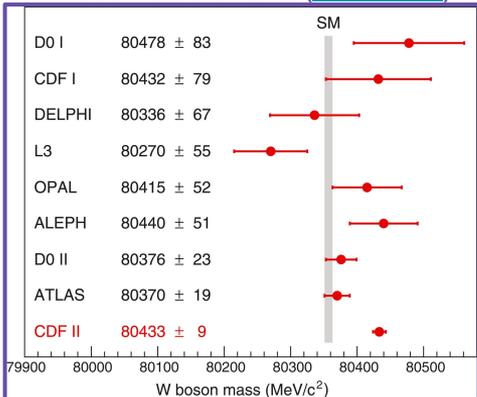
Calibration is a common experimental physics problem, whose goal is to infer the value and uncertainty of an unobservable quantity  $Z$  given a measured quantity  $X$ . Additionally, one would like to quantify the extent to which  $X$  and  $Z$  are correlated. In this paper, we present a machine learning framework for performing frequentist maximum likelihood inference with Gaussian uncertainty estimation, which also quantifies the mutual information between the unobservable and measured quantities. This framework uses the Donsker-Varadhan representation of the Kullback-Leibler divergence—parametrized with a novel Gaussian Ansatz—to enable a simultaneous extraction of the maximum likelihood values, uncertainties, and mutual information in a single training. We demonstrate our framework by extracting jet energy corrections and resolution factors from a simulation of the CMS detector at the Large Hadron Collider. By leveraging the high-dimensional feature space inside jets, we improve upon the nominal CMS jet resolution by upwards of 15%.



Sergio Diaz



### W Mass Measurements (DOI: 10.112)

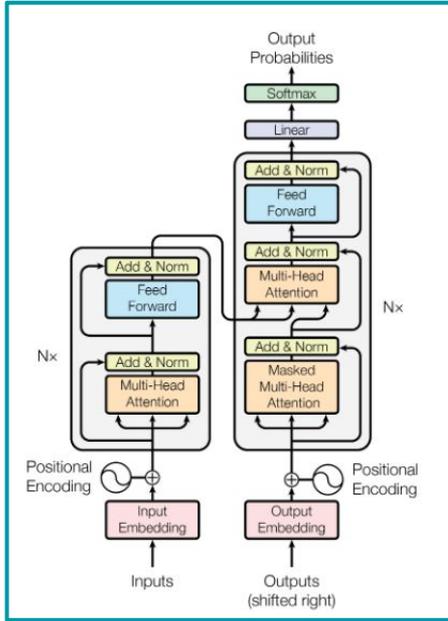


Exposing students to *both* **particle physics** and **machine learning** to explore new ways to **synthesize** the two!

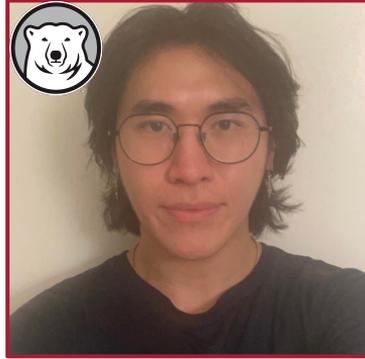


# Other Developments: Summer Students

Attention Is All You Need (1706.03762)

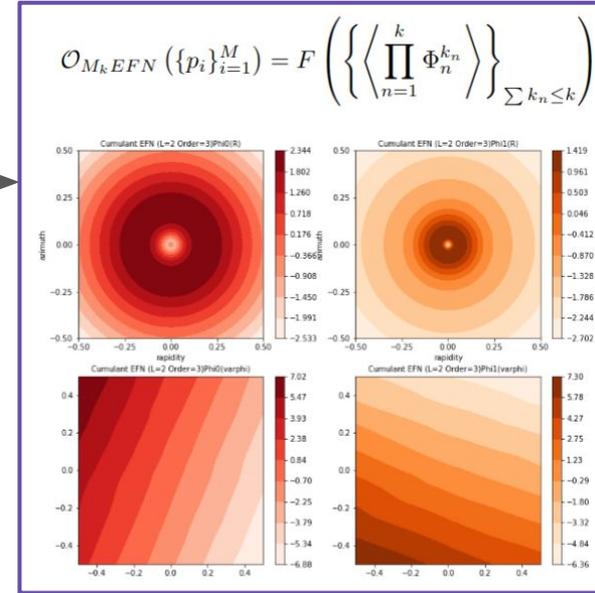


$A_i$



$f_i$

Moment Pooling (WIP)



Translating machine learning language into physics language:  
What does the attention mechanism look like for a physicist?

# Outlook

Exciting research in **physics** and **machine learning** enabled by **IAIFI!**

- Ideas from dictionary and manifold learning to analyze jet data
- Statistical frameworks for precision electroweak measurements
- Efficient machine learning architectures translated to physics language

Made possible by collaborations across fields and institutions!

