48,000 UC academic workers on strike for fair compensation, support for working parents, international scholar rights, and more.
To find out how you can support: fairucnow.org

PROFESSIONAL RESEARCHERS

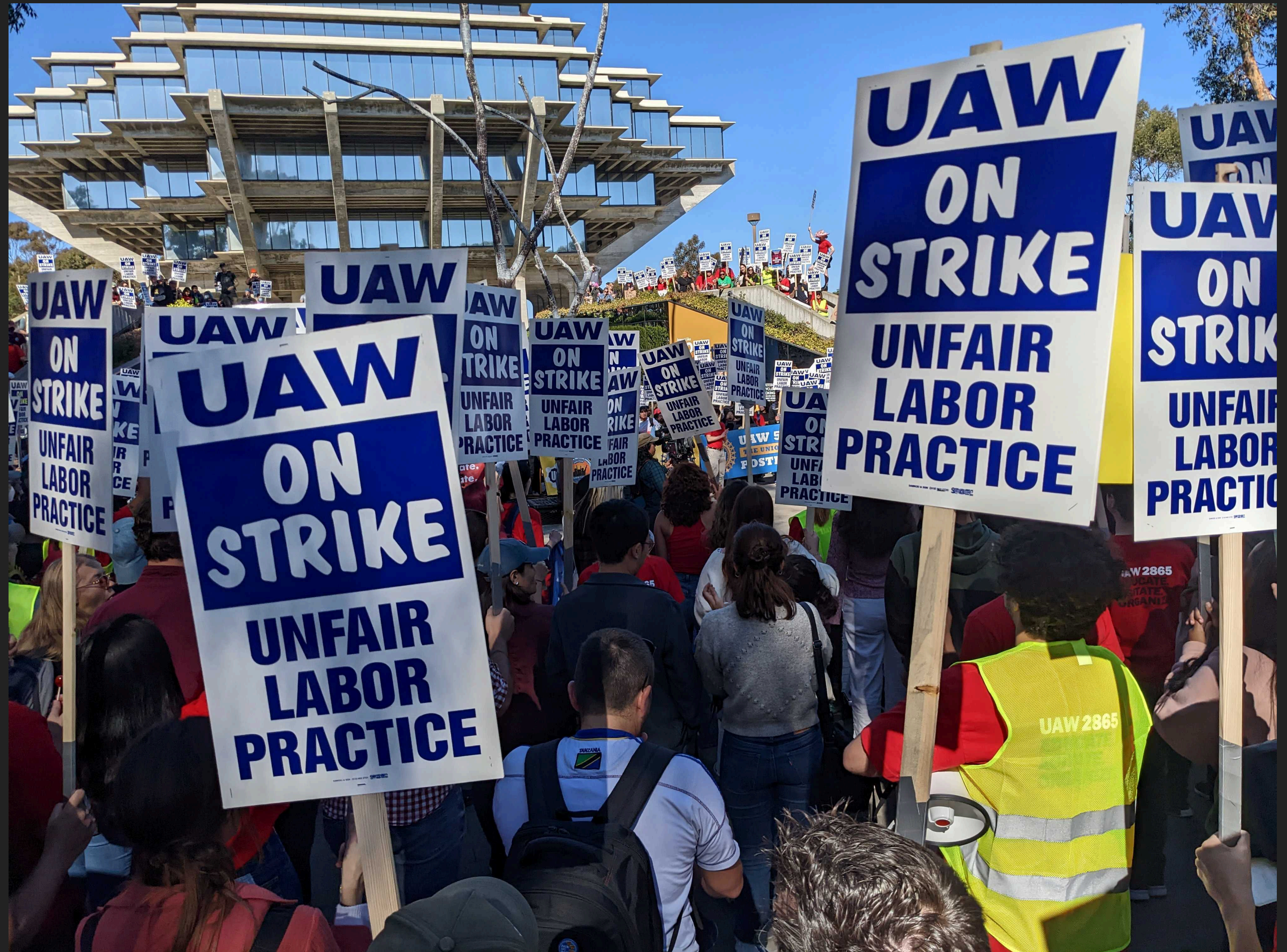STUDENT RESEARCHERS

PROJECT SCIENTISTS

SPECIALISTS

POSTDOCS

READERS

TUTORS

CPPs

TAs

48,000 STRONG

iml-wg.github.io/HEPML-LivingReview
github.com/jmduarte/Nomological_Net_ML_Particle_Physics

- High-level (expert) variables

- Ordered list of particles

- Images

- Set of particles

- Graph of particles

- Lorentz scalars/vectors

- Shallow neural network, boosted decision tree, …

- 1D convolutional neural network, recurrent neural network

- 2D convolutional neural network

- Deep set (energy flow network)

- Graph neural network

- Lorentz-equivariant network

▸ After "particle-flow reconstruction,"  can think of event as a collection of points in momentum space

▸ For jets (localized clusters of particles), dimensionality ($N_{\text{particles}} \sim 100, \ 4 + M$)

▸ Variable jet length requires:

  ▸ Preprocessing into another rep. (tab. data, jet images, …)

  ▸ Truncation to fixed size

  ▸ Graph NN

u,d or s jet

c or b jet

gluon jet

pileup jet

W or Z jet

Higgs jet

top jet

?

▸ Tabular data: use physics knowledge to preprocess jet information into a set of high-level features

▸ Substructure variable:

$$ {}_1e_3^\beta = \sum_{1 \le i < j < k \le n_J} z_i z_j z_k \min\{\Delta R_{ij}^\beta, \Delta R_{ik}^\beta, \Delta R_{jk}^\beta\} $$

$$ {}_2e_3^\beta = \sum_{1 \le i < j < k \le n_J} z_i z_j z_k \min\{\Delta R_{ij}^\beta \Delta R_{ik}^\beta, \Delta R_{ij}^\beta \Delta R_{jk}^\beta, \Delta R_{ik}^\beta \Delta R_{jk}^\beta\} $$

   ▸ jet mass

   ▸ energy correlation functions, e.g. $N_2^{\beta=1} = {}_2e_3^{\beta=1}/({}_1e_3^{\beta=1})^2$

▸ Jet images = pixelated versions of calorimeter hits in 2D (η, φ)

▸ Much lower level



top jet vs. u,d or s jet



top jet (on average)



QCD jet (on average)

# Boosted Boson Type Tagging

## *Jet ETmiss*



Convolutions

Convolved
Feature Layers

Max-Pooling

*W'→ WZ* event

Repeat

▸ CNNs among the best performing algorithms

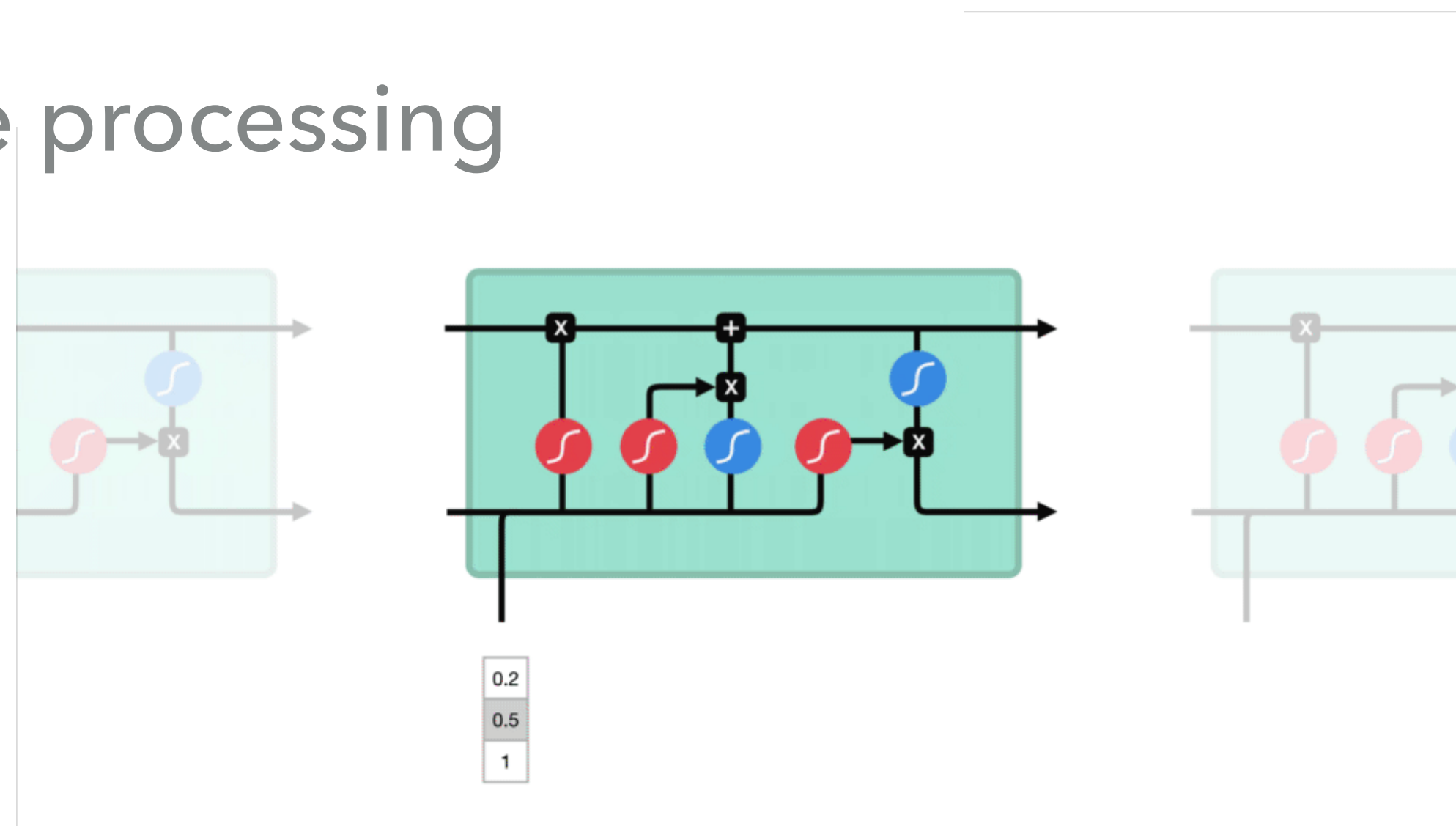| | AUC | Acc | $1/\epsilon_B$ ($\epsilon_S = 0.3$) | | | #Param |
|---|---|---|---|---|---|---|
| | | | single | mean | median | |
| CNN [16] | 0.981 | 0.930 | 914±14 | 995±15 | 975±18 | 610k |
| ResNeXt [31] | 0.984 | 0.936 | 1122±47 | 1270±28 | 1286±31 | 1.46M |
| TopoDNN [18] | 0.972 | 0.916 | 295±5 | 382± 5 | 378 ± 8 | 59k |
| Multi-body $N$-subjettiness 6 [24] | 0.979 | 0.922 | 792±18 | 798±12 | 808±13 | 57k |
| Multi-body $N$-subjettiness 8 [24] | 0.981 | 0.929 | 867±15 | 918±20 | 926±18 | 58k |
| TreeNiN [43] | 0.982 | 0.933 | 1025±11 | 1202±23 | 1188±24 | 34k |
| P-CNN | 0.980 | 0.930 | 732±24 | 845±13 | 834±14 | 348k |
| ParticleNet [47] | 0.985 | 0.938 | 1298±46 | 1412±45 | 1393±41 | 498k |
| LBN [19] | 0.981 | 0.931 | 836±17 | 859±67 | 966±20 | 705k |
| LoLa [22] | 0.980 | 0.929 | 722±17 | 768±11 | 765±11 | 127k |
| LDA [54] | 0.955 | 0.892 | 151±0.4 | 151.5±0.5 | 151.7±0.4 | 184k |
| Energy Flow Polynomials [21] | 0.980 | 0.932 | 384 | | | 1k |
| Energy Flow Network [23] | 0.979 | 0.927 | 633±31 | 729±13 | 726±11 | 82k |
| Particle Flow Network [23] | 0.982 | 0.932 | 891±18 | 1063±21 | 1052±29 | 82k |
| GoaT | 0.985 | 0.939 | 1368±140 | | 1549±208 | 35k |

▸ In deep learning, tailoring algorithms to the structure (and symmetries) of the data has led to groundbreaking performance

  ▸ CNNs for images

  ▸ RNNs for language processing

▸ Distributed unevenly in space
▸ Sparse
▸ Variable size
▸ No defined order
▸ Interconnections
  → Graphs

▸ What about high energy physics data like jets?

▸ Node features $\mathbf{v}_i$: particle 4-momentum

$$p = [E, p_x, p_y, p_z] \equiv [p_\mathrm{T}, \eta, \phi, m]$$

▸ Edge features $\mathbf{e}_k$: pseudoangular distance between particles

$$\sqrt{\Delta\eta^2 + \Delta\phi^2}$$

▸ Graph (global) features $\mathbf{u}$: jet mass

$$m = \sqrt{\sum_{i \in \mathrm{jet}} E_i^2 - p_{x,i}^2 - p_{y,i}^2 - p_{z,i}^2}$$

▸ Node-level tasks

  ▸ Identify "pileup" particles
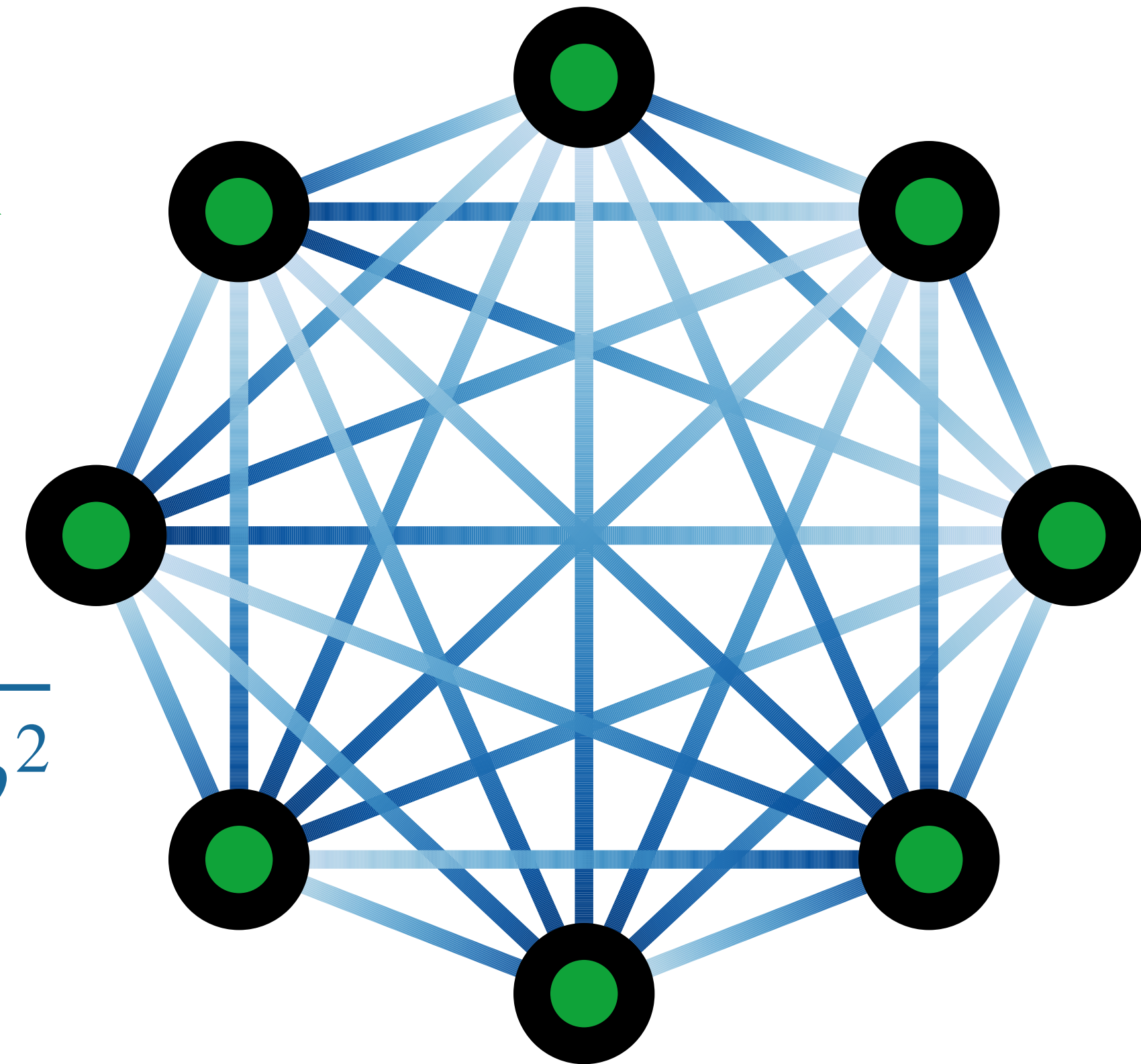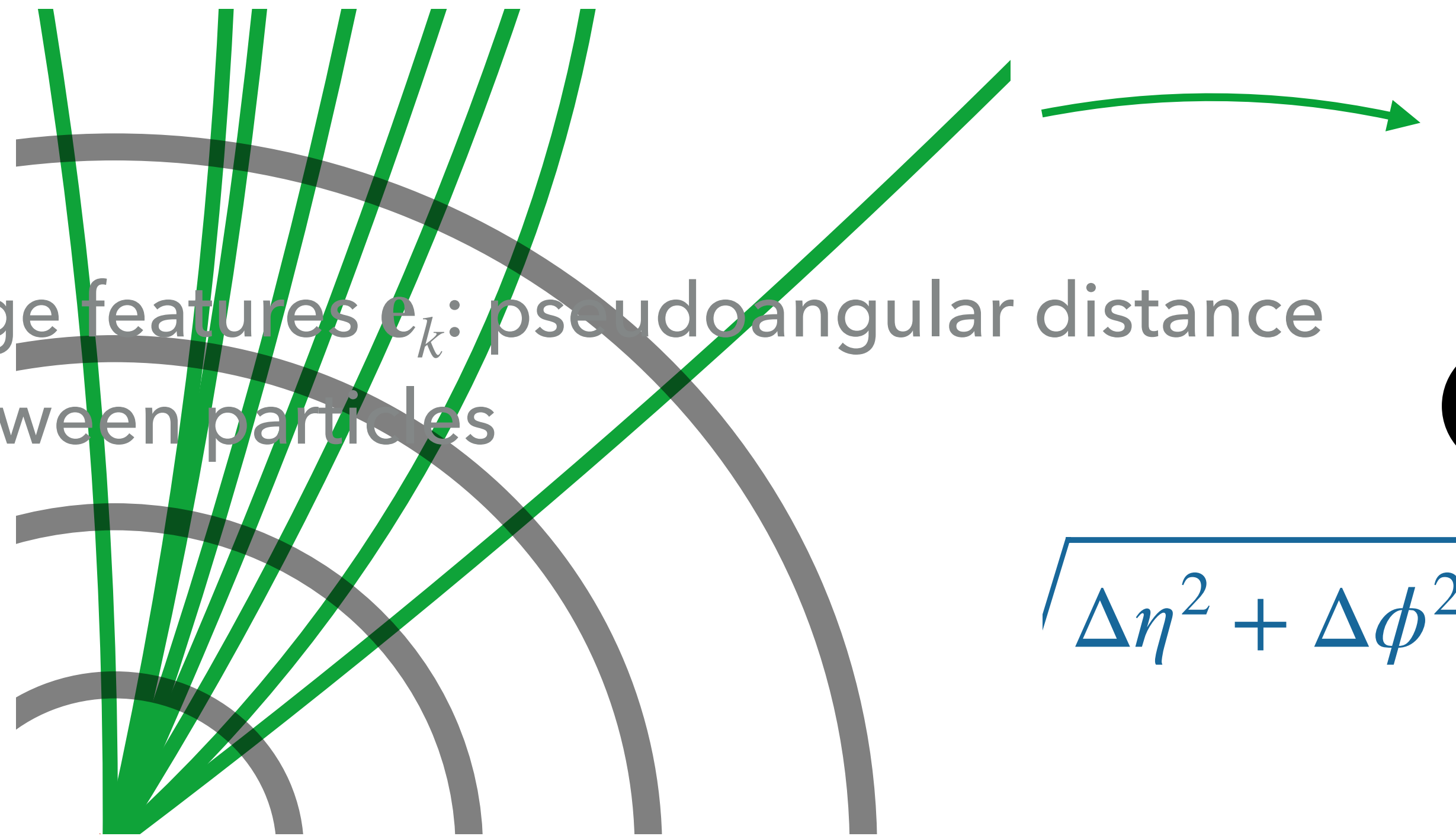
▸ Graph-level tasks

  ▸ **Jet tagging**



$\phi(x_i, x_j)$

$\mathbf{x}_i$

GNN

**Inputs**
$(\mathbf{X}, \mathbf{A})$

$\mathbf{h}_i$

**Latents**
$(\mathbf{H}, \mathbf{A})$

$\mathbf{z}_i$  **Node** classification
$\mathbf{z}_i = f(\mathbf{h}_i)$

$\mathbf{z}_G$  **Graph** classification
$\mathbf{z}_G = f\left(\bigoplus_{i \in \mathcal{V}} \mathbf{h}_i\right)$

$\mathbf{z}_{ij}$  **Link** prediction
$\mathbf{z}_{ij} = f(\mathbf{h}_i, \mathbf{h}_j, \mathbf{e}_{ij})$

  ▸ Edge-level tasks

    ▸ Identify good track doublets

▸ ParticleNet, using "dynamic edge convolutions:" graph is constructed based on "closeness" in an abstract "latent" space

▸ Identifies H(bb) with an efficiency of ~50% while rejecting 99.9% of background

▸ Symmetry-equivariant networks

    ▸ More economical (fewer, but more expressive parameters), interpretable, and trainable

**Invariance**

$$f(\rho_g(x)) = f(x)$$

**Equivariance**

$$f(\rho_g(x)) = \rho'_g(f(x))$$

▸ Lorentz-invariant networks:

  ▸ Boosting all particles into a new frame should give the same result

▸ Lorentz-equivariant networks:

  ▸ Boosting all particles into a new frame should give an output that transforms the same way

$$\begin{pmatrix} \gamma & -\gamma\beta & 0 & 0 \\ -\gamma\beta & \gamma & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Lorentz boost

arXiv:2201.08187

▸ State-of-the-art performance for top quark tagging

▸ Lorentz group invariance confirmed

# I. DATA REPRESENTATIONS & SYMMETRIES
# II. ANOMALY DETECTION
# III. GENERATIVE MODELING
# III. FAST INFERENCE
# VI. SUMMARY & OUTLOOK

▸ Supervised = full label information

▸ Semi-supervised = partial labels

▸ Weakly-supervised = noisy labels

▸ Unsupervised = no labels

- ▸ Example: autoencoders compress data and then uncompress it

- ▸ Assumption: if $x$ is far from $\text{Decoder}(\text{Encoder}(x))$, then $x$ has low $p_{\text{bkgd}}(x)$

background model independence

Some searches (train signal versus data)

**many new ideas!**

Most searches ("train" with simulations)

Train data versus background simulation

signal model independence

Encoder

Decoder

‣ Challenge with "black box" signals run in 2020–2021

‣ Plethora of new techniques



**3 Unsupervised**

3.1   Anomalous Jet Identification via Variational Recurrent Neural Network

3.2   Anomaly Detection with Density Estimation

3.3   BuHuLaSpa: Bump Hunting in Latent Space

3.4   GAN-AE and BumpHunter

3.5   Gaussianizing Iterative Slicing (GIS): Unsupervised In-distribution Anomaly Detection through Conditional Density Estimation

3.6   Latent Dirichlet Allocation

3.7   Particle Graph Autoencoders

3.8   Regularized Likelihoods

3.9   UCluster: Unsupervised Clustering

**4 Weakly Supervised**

4.1   CWoLa Hunting

4.2   CWoLa and Autoencoders: Comparing Weak- and Unsupervised methods for Resonant Anomaly Detection

4.3   Tag N' Train

4.4   Simulation Assisted Likelihood-free Anomaly Detection

4.5   Simulation-Assisted Decorrelation for Resonant Anomaly Detection

**5 (Semi)-Supervised**

5.1   Deep Ensemble Anomaly Detection

5.2   Factorized Topic Modeling

5.3   QUAK: Quasi-Anomalous Knowledge for Anomaly Detection

5.4   Simple Supervised learning with LSTM layers

▸ Several different strategies:

    ▸ Replace (part of) FullSim: increase speed, preserve accuracy

    ▸ Replace (part of) FastSim: maintain speed, increase accuracy

    ▸ Conditional: map generated → reconstructed events

    ▸ **End-to-end: map random noise → reconstructed events directly**

| | END-TO-END | | | | |
|---|---|---|---|---|---|

| | | CONDITIONAL | | | |
|---|---|---|---|---|---|

| | Full Simulation | FULL DETECTOR SIM W/ ML | DIGITIZATION EMULATION | RECONSTRUCTION | |
| HARD PROCESS GENERATION | SHOWERING/ HADRONIZATION/ UNDERLYING EVT | Fast Simulation — APPROXIMATE DETECTOR SIM W/ ML | PARTIAL DIGITIZATION EMULATION | SIMPLIFIED RECONSTRUCTION | ANALYSIS/ NTUPLING |
| | | Delphes — PARAMETRIZED SMEARING | | | |

# GENERATIVE ADVERSARIAL NETWORKS

Note: failure modes!

thispersond

▸ Train two neural networks in tandem:

  ▸ one to generate realistic "fake" data

  ▸ the other to discriminate "real" from "fake" data

▸ Evaluati

▸ We want

   ▸ the **qu**

   ▸ the **di**

   ▸ ultima

▸ To do so ... etrics

**On the Evaluation of Generative Models in High Energy Physics**

Raghav Kansal,* Anni Li, and Javier Duarte
*University of California, San Diego*

Nadezda Chernyavskaya, Maurizio Pierini
*European Center for Nuclear Research (CERN)*

Breno Orzari, Thiago Tomei
*Universidade Estadual Paulista, São Paulo/SP*
(Dated: November 16, 2022)

There has been a recent explosion in research into machine-learning- (ML-) based generative modeling to tackle computational challenges for simulations in high energy physics (HEP). In order to use such alternative simulators in practice, we need a well defined metrics to compare different generative models and evaluate their discrepancy from the true distributions. We present the first systematic review and investigation into evaluation metrics and their sensitivity to failure models of generative models, using the framework of two-sample goodness-of-fit testing, and their relevance and viability for HEP. Inspired by previous work in both physics and computer vision, we propose two new metrics, the Fréchet and Kernel Physics Distances (FPD and KPD), and perform a variety of experiments measuring their performance on simple Gaussian-distributed, and simulated high energy jet datasets. We find FPD, in particular, to be the most sensitive metric to all alternative jet distributions tested and recommend its adoption, along with KPD and Wasserstein distances between individual feature distributions, for evaluating generative models in HEP. We finally demonstrate the efficacy of these proposed metrics in evaluating and comparing a novel attention-based generative model, GAPT, to the state-of-the-art MPGAN jet simulation model.

assersstein
ance (W₁)

Qualit ✅

Diversi ✅

Physics F ✅

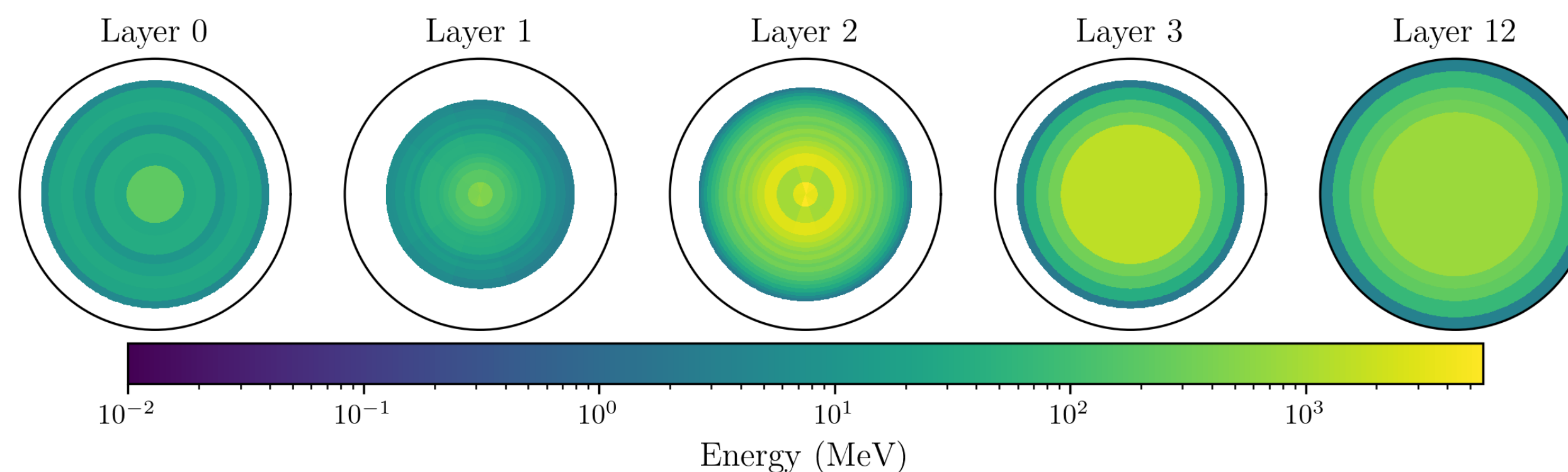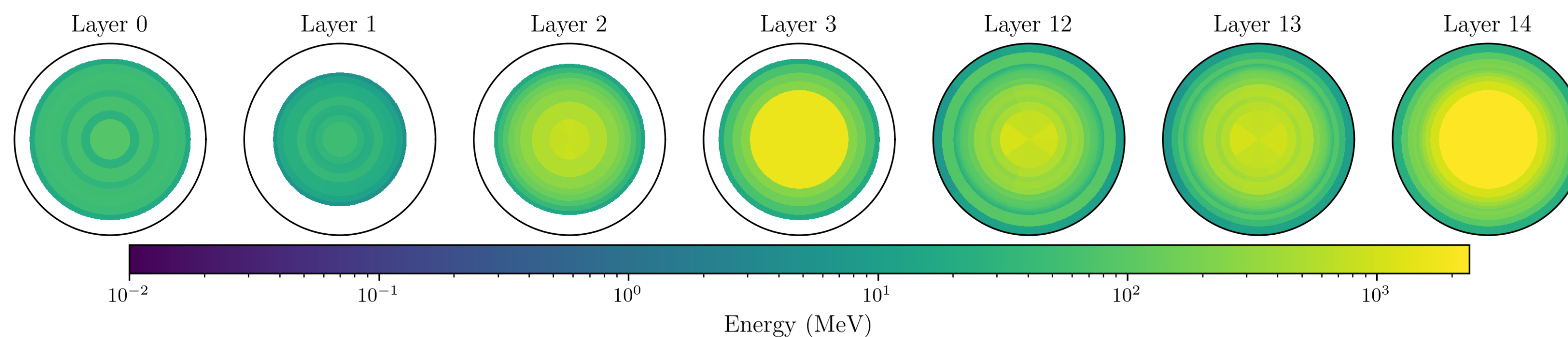▸ Ongoing challenge for generative modeling of calorimeter showers in HEP!

▸ Many new approaches presented at ML4Jets 2022: https://indico.cern.ch/event/1159913/



Shower average GEANT4 photon reference dataset

Shower average GEANT4 pion reference dataset

# DIFFUSION MODELS IN HEP

▸ Diffusion models have very recently dethroned GANs for natural images

▸ Generative model is trained using a diffusion process that slowly perturbs the data by adding noise – model learns to **denoise**

**Forward diffusion (training)**



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

▸ Generation of new samples by reversing the diffusion process

▸ Distribution of deposited energies for generated particle energies (top) and the energy deposition in a single layer of a calorimeter (bottom) vs time step

I.   DATA REPRESENTATIONS
     & SYMMETRIES
II.  ANOMALY DETECTION
III. GENERATIVE MODELING
III. FAST INFERENCE
VI.  SUMMARY & OUTLOOK

▸ **Codesign**: intrinsic development loop between algorithm design, training, and implementation

▸ Compression

  ▸ Maintain high performance while removing redundant operations

▸ Quantization

  ▸ Reduce precision from 32-bit floating point to 16-bit, 8-bit, …

▸ Parallelization

  ▸ Balance parallelization (how fast) with resources needed (how costly)

▸ hls4ml for scientists or ML experts to translate ML algorithms into RTL firmware

**Model**

**Keras
TensorFlow
PyTorch
…**

**Compressed
model**

**Machine learning model
optimization, compression**

AMD

XILINX®

intel®

Mentor®
A Siemens Business

▸ Challenge: if new physics has an unexpected signature that doesn't align with existing triggers, precious BSM events may be discarded at trigger level

▸ Can we use unsupervised algorithms to detect non-SM-like anomalies?

  ▸ Autoencoders (AEs): compress input to a smaller dimensional latent space then decompress and calculate difference

  ▸ Variational autoencoders (VAEs): model the latent space as a probability distribution; possible to detect anomalies purely with latent space variables

Encoder

μ

σ

Key observation: Can build an anomaly score from the latent space of VAE directly! No need to run decoder!

$$R_z = \sum_i \frac{\mu_i^2}{\sigma_i^2}$$

▸ CNNs as the basis for (V)AEs for anomaly detection

▸ Good anomaly detection performance for unseen signals
($LQ \to b\tau$, $A \to 4l$, **$h^{\pm} \to \tau\nu$**, $h^0 \to \tau\tau$)

▸ **VAE** fits in latency and resource requirements for HL-LHC!



Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Block 0:
Input 19x3x1
ZeroPadding (1,0)
BatchNorm

Output:
Conv2d 4 (1,(3,3))

CNN ROC $h^{\pm} \to \tau\nu$
— IO VAE (AUC = 95%)
— VAE $D_{KL}$ (AUC = 86%)
— VAE $R_z$ (AUC = 86%)
— IO AE (AUC = 96%)

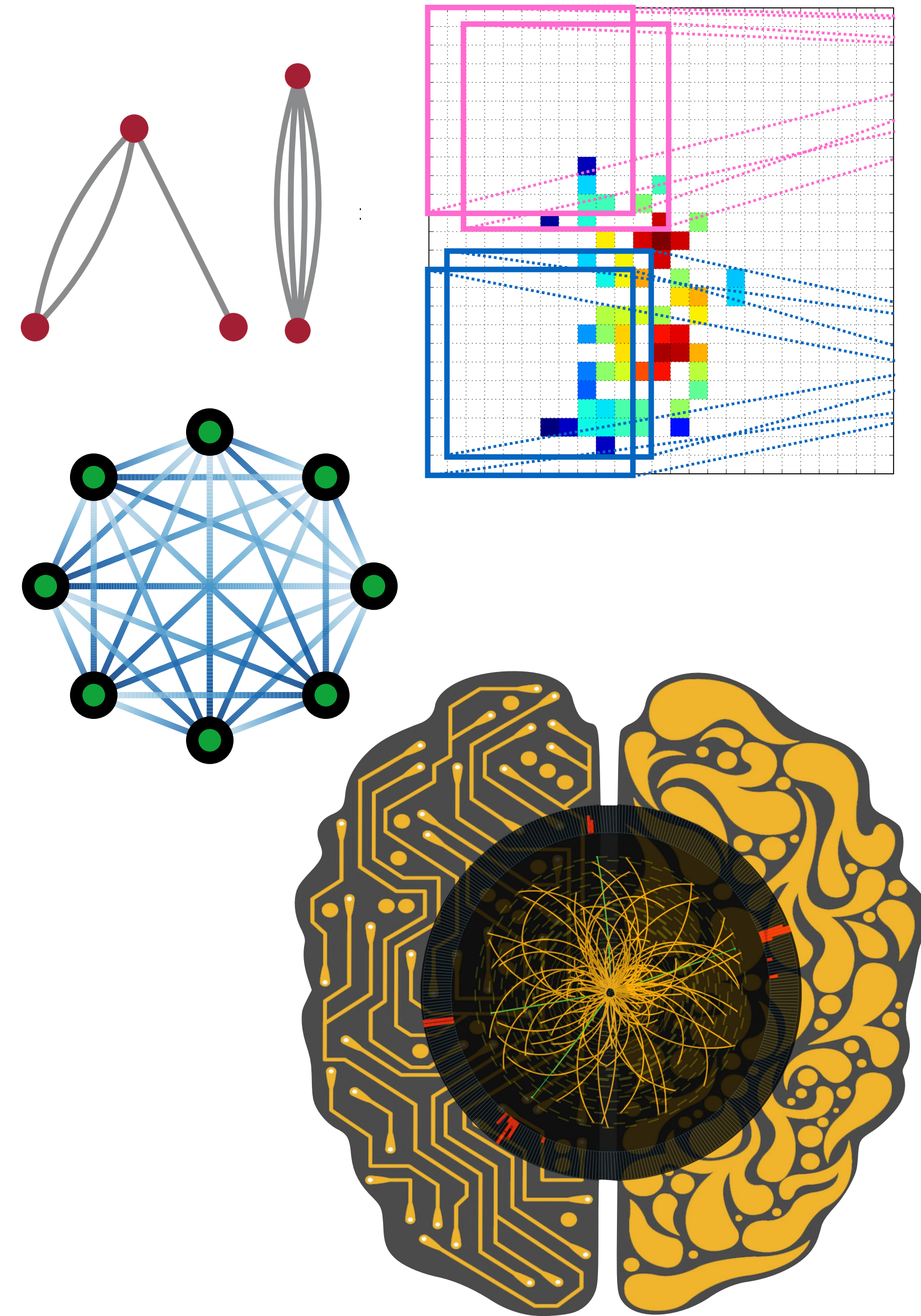| Model | DSP [%] | LUT [%] | FF [%] | BRAM [%] | Latency [ns] | II [ns] | AUC [%] | TPR @ FPR=$10^{-5}$ |
|---|---|---|---|---|---|---|---|---|
| CNN VAE $R_z$ | 10 | **12** | **4** | **2** | **365** | **115** | 86 | 0.06% |
| CNN AE | **7** | 47 | 5 | 6 | 1480 | 895 | **96** | **0.10%** |

I.   DATA REPRESENTATIONS
     & SYMMETRIES
II.  ANOMALY DETECTION
III. GENERATIVE MODELING
III. FAST INFERENCE
VI. SUMMARY & OUTLOOK

▸ Different representations of HEP data, from tabular data, image data, set data, graph data, paired with corresponding algorithms can achieve excellent performance

▸ Plethora of ML techniques in HEP from anomaly detection to generative modeling have exploded in recent years

# ts as Images

▸ Availability of public datasets and challenges have advanced the state of the art

φ) to a rectangular grid that allows for an image-ergy from particles are deposited in pixels in (η, φ)
em as the pixel intensities in a greyscale analogue.
st introduced by our group [JHEP 02 (2015) 118],
s event reconstruction and computer vision. We
he jet-axis, and normalize each image, as is often
scriminative difference in pixel intensities.

▸ Fast ML can accelerate science allowing us to test hypotheses faster, enhance performance of detectors/accelerators, and use potentially overlooked data

B. Nachman (SLAC)

Below, we have
signal and backg
difference-visualiza