# *Graph Neural Networks*
# *for (Experimental) Particle Physics*

Huilin Qu

# ABOUT ME

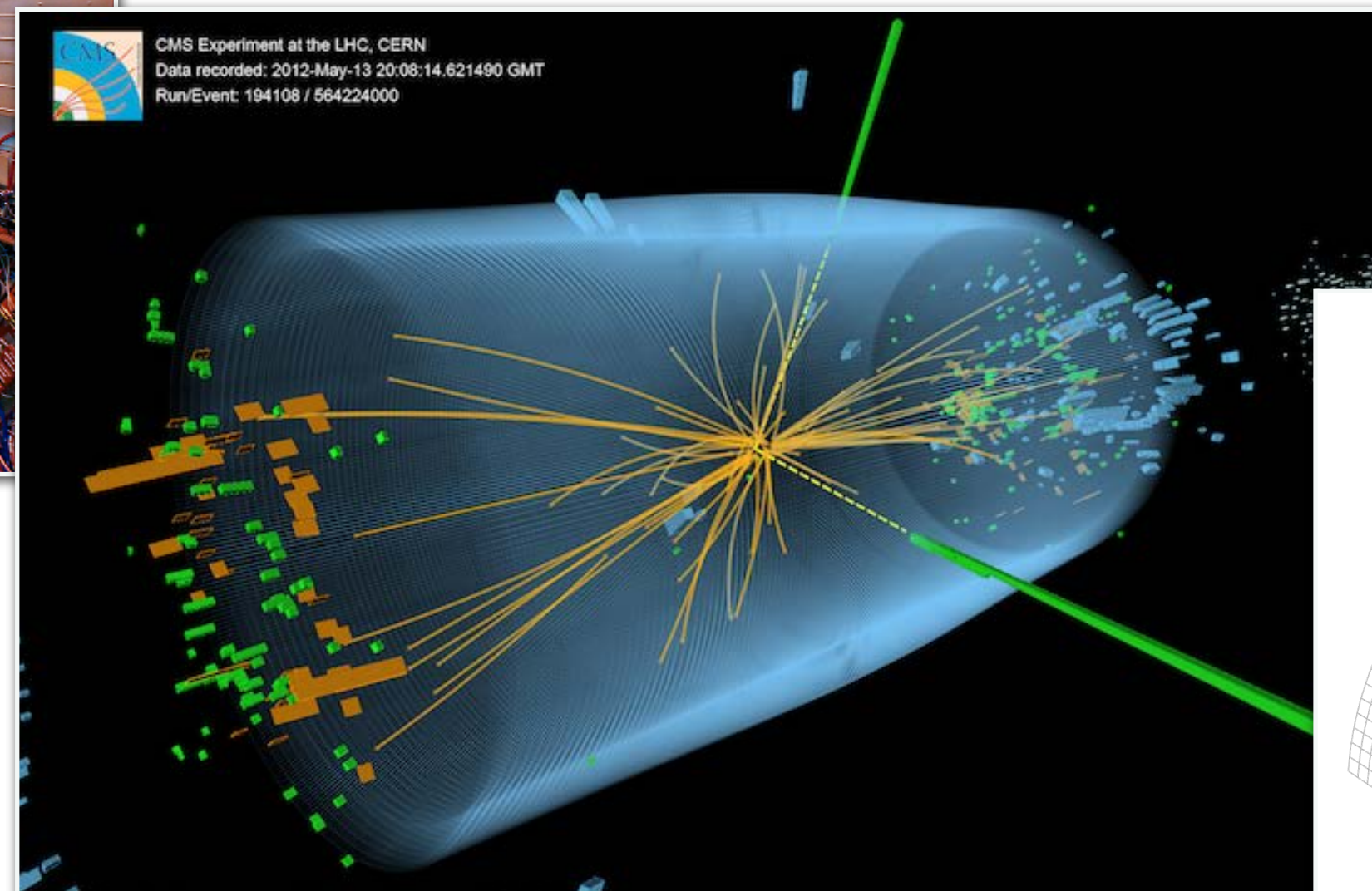- Experimental particle physicist at CERN

  - interests: Higgs and New Physics @ LHC, jet physics, machine learning, ...

LARGE HADRON COLLIDER

LHCb

ATLAS

CERN Meyrin

CERN Prévessin

SPS 7 km

SUISSE
FRANCE

ALICE

CMS

LHC 27 km

# PARTICLE COLLISION

*proton beam*  ⟹

⟸  *proton beam*

# PARTICLE COLLISION

# PARTICLE COLLISION

# PARTICLE COLLISION

# PARTICLE COLLISION



*What we observe…*

# PARTICLE COLLISION

*What we want to know…*

*What we observe…*

# EXAMPLE: WHAT IS THE COLLISION EVENT?

*Electron*

*Neutrino
(momentum imbalance)*

**Jet:** *a collimated spray of particles*

*Electron*

*Neutrino*
*(momentum imbalance)*

*Key question:*
**What particle initiates the jet?**



*Electron*

*Neutrino*
*(momentum imbalance)*

# THE DATA CHALLENGE IN HIGH ENERGY PHYSICS

*HEP*



*Large volume of data, complex topology, …*

# AI + HEP: At the Collision Point



*HEP*

*AI*

*Collimate HEP and AI to make them collide!*

*Large volume of data, complex topology, …*

# *Data Representations*

# DATA REPRESENTATION

HEP

ML



×



*Collision events, detector hits, sensor arrays, …*

*First and foremost:*
**How to represent the data?**

# DATA REPRESENTATION: IMAGE

*HEP*

*Image*



*Collision events, detector hits, sensor arrays, …*

e.g., de Oliveira, Kagan, Mackey, Nachman and Schwartzman, arXiv:1511.05190

- Convert to 2D/3D image => **Computer vision**

  - then use convolutional neural networks (CNNs)

  - but:

    - inhomogeneous geometry, high sparsity, …

# DATA REPRESENTATION: SEQUENCE

*HEP*

*Sequence*



*Collision events, detector hits, sensor arrays, …*

e.g., Guest, Collado, Baldi, Hsu, Urban, Whiteson
arXiv: 1607.08633

- Convert to a sequence => **Natural language processing (NLP)**

  - recurrent neural network (RNN), e.g., GRU/LSTM; 1D CNNs; etc.

# DATA REPRESENTATION: SEQUENCE?

*HEP*

*Sequence*



= | Permutation symmetry

≠

- Convert to a sequence => **Natural language processing (NLP)**

  - recurrent neural network (RNN), e.g., GRU/LSTM; 1D CNNs; etc.

  - but:

    - must impose an ***ordering*** on the particles/hits, which can limit the learning performance

# POINT CLOUD

*An unordered set of points in space
(e.g., produced by a LiDAR on self-driving cars)*

# DATA REPRESENTATION: POINT CLOUD

**HEP**

**Point cloud**



$\phi$

$\eta$

*Collision events, detector hits, sensor arrays, …*

- HEP data as a point cloud
  - each particle / detector cell is a point in the cloud
    - for each point: (spatial) coordinates + any additional properties (energy/momentum, detector response, …)
  - key feature: ***permutation symmetry***

# LEARNING ON POINT CLOUDS

*HEP*

*Point cloud*

$\phi$

$\eta$

*Collision events, detector hits, sensor arrays, …*

- Desired algorithms for learning on point cloud data

  - symmetry-preserving: the outputs should be invariant under permutation of the points

  - high expressiveness: capable of fully exploiting the correlations between points

  - low computational cost: scalable from O(10) to O(1000) points, and even up to O(1M) points in some cases

# LEARNING ON POINT CLOUDS

HEP

Point cloud

Graph neural network - A unified framework



Edge block   Node block   Global block

Review in Shlomi, Battaglia, Vlimant, arXiv:2007.13681

Unshared, deep GN stack

$G_0 \rightarrow \boxed{GN_1} \rightarrow G_1 \rightarrow \boxed{GN_2} \rightarrow \cdots \rightarrow \boxed{GN_M} \rightarrow G_M$

Shared, recurrent GN stack

Collision events, detector hi

$\eta$

- Desired algorithms for learning on point cloud data
  - symmetry-preserving: the outputs should be invariant under permutation of the points
  - high expressiveness: capable of fully exploiting the correlations between points
  - low computational cost: scalable from O(10) to O(1000) points, and even up to O(1M) points in some cases

GNNs for Particle Physics  - July 2023 - Huilin Qu (CERN)

# A Journey Through Graph Neural Networks

# WHAT IS A GRAPH?

*graph level attributes*

*Graph:* $\quad G = (\mathbf{u}, V, E) \quad$ with $N_v$ vertices and $N_e$ edges

*Vertices (nodes)*

$$V = \{\mathbf{v}_i\}_{i=1:N_v}$$

*attributes of the i-th node*

*Edges (links)*

$$E = \{(\mathbf{e}_k, r_k, s_k)\}_{k=1:N_e}$$

*attributes of the k-th edge*

*indices of the two nodes (receiver and sender) connected by the k-th edge*

$v_1$ $\quad$ $v_2$

$e_1$

2

1

$e_2$

3

$v_3$

27

# HOW TO BUILD THE GRAPH?

- From point clouds to graphs:

  - points (particles/hits/sensors) naturally become the **nodes** of the graph

  - but how to define the ***edges***?

**Set: no edges**

**Fully connected graph**
- *i.e., connect each node to all other nodes*

**Hierarchical trees:**
- *decay chain*
- *jet clustering history*

**Locally connected graph**
- *i.e., connect each node only to neighbor nodes*
  - *k-nearest neighbors*
  - *fixed radius*

static

(dynamically) learned

# GRAPH NETWORK FORMALISM

- Typical graph neural networks (GNNs) can be described in the "Message Passing" framework

**Graph:** $G = (\mathbf{u}, V, E)$  with $N_v$ vertices and $N_e$ edges

*Vertices (nodes)*  $s$*Edges (links)*

$$V = \{\mathbf{v}_i\}_{i=1:N_v}$$  $$E = \{(\mathbf{e}_k, r_k, s_k)\}_{k=1:N_e}$$

$v_1$  $e_1$  $2$  $v_2$

$1$  $(\boldsymbol{u}, V, E)$

$e_2$  $3$  $v_3$

$e'_k$: message computed for edge $k$ connecting nodes $r_k$, $s_k$

Shlomi, Battaglia, Vlimant, arXiv:2007.13681

$$e'_k = \phi^e(\mathbf{e}_k, \boldsymbol{v}_{r_k}, \boldsymbol{v}_{s_k}, \mathbf{u})$$

$\mathbf{u}$  $\rightarrow \mathbf{u}'$

$V$  $\rightarrow V'$

$E$ — $\phi^e$  $\rightarrow E'$

Edge block

# GRAPH NETWORK FORMALISM

- Typical graph neural networks (GNNs) can be described in the "Message Passing" framework

*Graph:* $G = (\mathbf{u}, V, E)$    with $N_v$ vertices and $N_e$ edges

*Vertices (nodes)*      *sEdges (links)*

$$V = \{\mathbf{v}_i\}_{i=1:N_v} \qquad E = \{(\mathbf{e}_k, r_k, s_k)\}_{k=1:N_e}$$

$(\boldsymbol{u}, V, E)$

$e_k'$: message computed for edge $k$ connecting nodes $r_k, s_k$

$v_i'$: node feature update based on aggregated messages and previous features

Shlomi, Battaglia, Vlimant, arXiv:2007.13681

$$e_k' = \phi^e(\mathbf{e}_k, v_{r_k}, v_{s_k}, \mathbf{u}) \qquad \bar{e}_i' = \rho^{e \to v}(E_i')$$

$$v_i' = \phi^v\left(\bar{e}_i', v_i, \boldsymbol{u}\right)$$

- Typical graph neural networks (GNNs) can be described in the "Message Passing" framework

**Graph:** $G = (\mathbf{u}, V, E)$  with $N_v$ vertices and $N_e$ edges

*Vertices (nodes)*  *Edges (links)*

$$V = \{\mathbf{v}_i\}_{i=1:N_v}$$  $$E = \{(\mathbf{e}_k, r_k, s_k)\}_{k=1:N_e}$$

$(\boldsymbol{u}, V, E)$

$e'_k$: message computed for edge $k$ connecting nodes $r_k$, $s_k$

$v'_i$: node feature update based on aggregated messages and previous features

$u'$: global feature update based on aggregated, updated node and edge features

$$e'_k = \phi^e(\mathbf{e}_k, v_{r_k}, v_{s_k}, \mathbf{u}) \qquad \bar{e}'_i = \rho^{e \to v}(E'_i)$$

$$v'_i = \phi^v\left(\bar{e}'_i, v_i, u\right) \qquad \bar{e}' = \rho^{e \to u}(E')$$

$$u' = \phi^u(\bar{e}', \bar{v}', u) \qquad \bar{v}' = \rho^{v \to u}(V')$$

Shlomi, Battaglia, Vlimant, arXiv:2007.13681



Edge block  Node block  Global block

# GRAPH NETWORK FORMALISM

- Typical graph neural networks (GNNs) can be described in the "Message Passing" framework

**Graph:** $G = (\mathbf{u}, V, E)$   with $N_v$ vertices and $N_e$ edges

$(\mathbf{u}, V, E)$

*Vertices (nodes)*

$V = \{\mathbf{v}_i\}_{i=1:N_v}$

*Edges (links)*

$E = \{(\mathbf{e}_k, r_k, s_k)\}_{k=1:N_e}$

$e_k'$: message computed for edge $k$ connecting nodes $r_k$, $s_k$

$v_i'$: node feature update based on aggregated messages and previous features

$u'$: global feature update based on aggregated, updated node and edge features

$$e_k' = \phi^e(\mathbf{e}_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}, \mathbf{u})$$

$$v_i' = \phi^v\left(\bar{\mathbf{e}}_i', \mathbf{v}_i, \mathbf{u}\right)$$

$$u' = \phi^u(\bar{\mathbf{e}}', \bar{\mathbf{v}}', \mathbf{u})$$

$$\bar{\mathbf{e}}_i' = \rho^{e \to v}(E_i')$$

$$\bar{\mathbf{e}}' = \rho^{e \to u}(E')$$

$$\bar{\mathbf{v}}' = \rho^{v \to u}(V')$$

Shlomi, Battaglia, Vlimant, arXiv:2007.13681



Edge block    Node block    Global block

*Shared-weight NN*    *Symmetric functions (e.g., sum, mean, max, etc.)*    → *Permutation invariance*

# GRAPH NETWORK FORMALISM

- Typical graph neural networks (GNNs) can be described in the "Message Passing" framework

*GN layer*

*Full GNN*

# EXAMPLE: DEEP SETS

---

**Deep Sets**

[1703.06114]

---

**Manzil Zaheer**[1,2]**, Satwik Kottur**[1]**, Siamak Ravanbhakhsh**[1]**,**
**Barnabás Póczos**[1]**, Ruslan Salakhutdinov**[1]**, Alexander J Smola**[1,2]
[1] Carnegie Mellon University      [2] Amazon Web Services

**Deep Sets Theorem [63].** *Let* $\mathfrak{X} \subset \mathbb{R}^d$ *be compact,* $X \subset 2^{\mathfrak{X}}$ *be the space of sets with bounded cardinality of elements in* $\mathfrak{X}$*, and* $Y \subset \mathbb{R}$ *be a bounded interval. Consider a continuous function* $f : X \to Y$ *that is invariant under permutations of its inputs, i.e.* $f(x_1, \ldots, x_M) = f(x_{\pi(1)}, \ldots, x_{\pi(M)})$ *for all* $x_i \in \mathfrak{X}$ *and* $\pi \in S_M$*. Then there exists a sufficiently large integer* $\ell$ *and continuous functions* $\Phi : \mathfrak{X} \to \mathbb{R}^{\ell}$*,* $F : \mathbb{R}^{\ell} \to Y$ *such that the following holds to an arbitrarily good approximation:*[1]

*Set: no edges*

$$f(\{x_1, \ldots, x_M\}) = F\left( \sum_{i=1}^{M} \Phi(x_i) \right)$$



Node block | Global block

# EXAMPLE: DYNAMIC GRAPH CNN (DGCNN)

**Dynamic** *locally connected graph*
- *k-nearest neighbors*

*Key building block: EdgeConv*

For the 1st layer:
kNN in input coordinates (xyz/η-φ)



Wang, Sun, Liu, Sarma, Bronstein, Solomon, arXiv:1801.07829

For subsequent layers:
kNN in learned latent space



Edge block      Node block      Global block

$$e'_{ij} = \phi^e(v_i, v_j) = MLP(v_i, v_j)$$

$$v'_i = \rho^{e \to v}(E'_i) = \square_j e'_{ij}$$

$$\square = sum, mean, max, etc.$$

# EXAMPLE: TRANSFORMER

- **Transformers**: the new state-of-the-art architecture in ML — foundation of LLM like BERT/GPT

  - core concept: self-attention mechanism

Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser, Polosukhin, arXiv:1706.03762 ["*Attention Is All You Need*"]

# EXAMPLE: TRANSFORMER

- The transformer architecture is also permutation-invariant as long as positional encoding is not used



*Fully connected graph*

$$e'_{ij} = \phi^e(v_i, v_j) = (W_Q v_i)^T (W_K v_j)$$

$$w_{ij} = \text{softmax}_j \frac{e'_{ij}}{\sqrt{d_k}}$$

$$v'_i = \sum_j w_{ij}(W_V v_j)$$

Edge block    Node block    Global block

# GRAPH NEURAL NETWORKS IN ACTION

# GRAPH ML TASKS

# GRAPH ML TASKS

CMS Experiment at LHC, CERN
Data recorded: Sat Aug  5 15:32:22 2017 CEST
Run/Event: 300515 / 205888132

*Jet:* a collimated spray of particles

p

p

H   ?

W   ?

Z   ?

t   ?

......

Key question:

**What type of particle initiates the jet?**

*The answer — Jet tagging!*

42

# JET TAGGING

- Jet tagging: identifying the origin of a jet, i.e., what kind of particle initiates the jet
  - essentially a classification task from the machine learning perspective



*Jet flavor tagging*

*Boosted jet tagging*

*Hadronic τ tagging*

GNNs for Particle Physics - July 2023 - Huilin Qu (CERN)

# BOOSTED JET TAGGING

- Hadronic decays of highly Lorentz-boosted heavy particles (Higgs/W/Z/top) lead to large-radius jets with distinctive characteristics:

  - different radiation patterns ("**substructure**")

    - 3-prong (top), 2-prong (W/Z/H) vs 1-prong (gluon/light quark jet)

  - different **flavor** content: existence of one or more b-/c-quarks

- Boosted jet tagging:

  - simultaneously exploiting both **substructure** and **flavor** to maximize the performance

  - significant performance leap thanks to deep learning techniques

*3-prong*

top
b
q
q

*2-prong*

W/Z/H
q
q

*1-prong*

q/g
q/g

jet
displaced tracks
charged lepton
heavy-flavour jet
SV
flight distance
IP
PV
jet
jet

# PARTICLENET

- ParticleNet: jet tagging via particle clouds

  - treating a jet as an **unordered set of particles**, distributed in the η − φ space

  - **graph neural network architecture**, adapted from Dynamic Graph CNN [arXiv:1801.07829]

    - treating a point cloud as a graph: each point is a vertex

      - for each point, a local patch is defined by finding its k-nearest neighbors

    - designing a permutation-invariant "convolution" function

      - define "edge feature" for each center-neighbor pair: $e_{ij} = MLP(x_i, x_j)$

      - aggregate the ... way:

*ParticleNet architecture*



coordinates   features

k-NN

k-NN indices → edge features

Linear
BatchNorm
ReLU
Linear
BatchNorm
ReLU
Linear
BatchNorm
ReLU
Aggregation
⊕
ReLU

coordinates   features

EdgeConv Block
k = 16, C = (64, 64, 64)

EdgeConv Block
k = 16, C = (128, 128, 128)

EdgeConv Block
k = 16, C = (256, 256, 256)

Global Average Pooling

Fully Connected
256, ReLU, Dropout = 0.1

Fully Connected
2

Softmax

*GNNs for Particle Physics – July 2023 – Huilin Qu (CERN)*

*cf.* P.T. Komiske, E. M. Metodiev and J. Thaler, *JHEP 01 (2019) 121*;
V. Mikuni and F. Canelli, *Eur. Phys. J. Plus 135, 463 (2020); Mach.Learn.Sci.Tech. 2 (2021) 3, 035027*.

# PARTICLENET: PERFORMANCE

- Top performance among a variety of deep learning taggers on the community-wide top tagging benchmark

| | 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 [30] | 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]   *(Preliminary ver.)* | 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 |
| 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 |
| *ParticleNet-Lite* | 0.984 | 0.937 | 1262±49 | | | 26k |
| *ParticleNet* | **0.986** | **0.940** | **1615±93** | | | 366k |

*Ensemble of all taggers*

# PARTICLENET: BEYOND JET TAGGING

**BES III**

$\Lambda_c^+ \to ne^+\nu$ search

*Yunxuan Song, Yangu Li et al.*



**FCC**

Particle identification

*Eur.Phys.J.Plus 137 (2022) 1, 39*
*Eur.Phys.J.C 82 (2022) 7, 646*



**KM3NeT**

Muon bundle reconstruction

*JINST 16 (2021) 10, C10011,*
*PoS ICRC2021 (2021) 1048*



**(b)** events with two or more muons



Cosmic ray pattern identification

*Astropart.Phys. 126 (2021) 102527*

# GNNs for Event Reconstruction

# EVENT RECONSTRUCTION

- Event reconstruction: deciphering the detector signals

  - what are the outgoing particles?

  - what are their momenta, energy, …?



Key:
- Muon
- Electron
- Charged Hadron (e.g. Pion)
- Neutral Hadron (e.g. Neutron)
- Photon

Transverse slice through CMS

3.8T

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

2T

Iron return yoke interspersed with Muon chambers

0m  1m  2m  3m  4m  5m  6m  7m

# CALORIMETER RECONSTRUCTION

Community Detection

- GNNs also powerful tools for event reconstruction, particularly for non-uniform detector geometry

- Distance-weighted GNNs: GarNet/GravNet
  - much lower computational cost than DGCNN
  - GarNet: lightweight, can be <u>implemented on FPGA</u> for e.g., event triggering

- Object condensation: one-stage multi-object reconstruction
  - simultaneously predict the number of showers and their properties
  - in addition: cluster hits belonging to shower in a clustering space by using attractive/repulsive potentials in the loss



GarNet/GravNet

Truth

Reconstructed

Time and memory usage

S. R. Qasim, J. Kieseler, Y. Iiyama and M. Pierini [*EPJC 79 (2019) 7, 608*]; J. Kieseler [*EPJC 80 (2020) 9, 886*]; S. R. Qasim et. al., [*EPJC 82, 753 (2022)*]

# PARTICLE FLOW

Node Classification

- Use GNNs to directly perform end-to-end particle flow reconstruction
  - comparable/better performance than rule-based PF on Delphes dataset
  - runtime scales linearly with input size, no quartic explosion

$t\bar{t}$, 14 TeV, 200 PU
— Tracks
■ ECAL clusters
■ HCAL clusters
× Truth particles

*Delphes simulation*

**Event as input set**
$X = \{x_i\}$

**Event as graph**
$X = \{x_i\}, A = A_{ij}$

**Transformed inputs**
$H = \{h_i\}$

**Graph building**
LSH+kNN
$\mathscr{F}(X \mid w) = A$

**Message passing**
GCN
$\mathscr{G}(X, A \mid w) = H$

**Target set** $Y = \{y_j\}$

**Output set** $Y' = \{y'_j\}$

Elementwise loss $L(y_j, y'_j)$
classification & regression

**Decoding**
elementwise FFN
$\mathscr{D}(x_j, h_j \mid w) = y'_j$

$x_i = [\text{type}, p_\text{T}, E_\text{ECAL}, E_\text{HCAL}, \eta, \phi, \eta_\text{outer}, \phi_\text{outer}, q, \ldots], \quad \text{type} \in \{\text{track, cluster}\}$
$y_j = [\text{PID}, p_\text{T}, E, \eta, \phi, q, \ldots], \quad \text{PID} \in \{\text{none, charged hadron, neutral hadron}, \gamma, e^\pm, \mu^\pm\}$
$h_i \in \mathbb{R}^{256}$
Trainable neural networks: $\mathscr{F}, \mathscr{G}, \mathscr{D}$
● - track, ■ - calorimeter cluster, ■ - encoded element
■ - target (predicted) particle, ■ - no target (predicted) particle

*Resolution*

Particles
QCD, 14 TeV, PU200
Charged hadrons
Rule-based PF
$\mu = -0.01, \sigma = 0.21$
MLPF
$\mu = 0.03, \sigma = 0.14$
$p_\text{T}$ resolution, $(p'_\text{T} - p_\text{T})/p_\text{T}$

*Inference time*

Average runtime / event [ms]
$t\bar{t}$, 14 TeV
▼ 40 PU
▲ 80 PU
● 200 PU
● MLPF scaling
Average event size [elements]

51

# GNNS FOR TRACKING

- Charged particle tracking as an edge prediction task within the GNN framework

  - each hit is a node of the graph

  - edges constructed between pairs of hits with geometrically plausible relations

  - classify whether each edge connects hits belonging to the same track or not

G. DeZoort et al.
[Comput. Softw. Big Sci. 5, 26 (2021)]



$$a_{ij}^{(1)} = \phi_{R,1}(x_i^{(0)}, x_j^{(0)}, a_{ij}^{(0)}) \quad w_{ij}^{(1)} = \phi_{R,2}(x_i^{(1)}, x_j^{(1)}, a_{ij}^{(1)})$$

$$x_i^{(1)} = \phi_O(x_i^{(0)}, \sum_{j \in N(i)} a_{ij}^{(1)})$$

GNNs for Particle Physics – July 2023 – Huilin Qu (CERN)

See also: S. Farrell et al. [1810.06111]; X. Ju et al. [2003.11603];
C. Biscarat, S. Caillou, C. Rougier, J. Stark and J. Zahreddine [2103.00916]; X. Ju et al. [2103.06995]; etc.

52

# GRAPH GENERATIVE MODELS

*https://towardsdatascience.com/graph-convolutional-networks-deep-99d7fee5706f*

53

# PARTICLE CLOUD GENERATION

- Exploit GNNs for "particle cloud" generation

  - enables fast detector simulation

R. Kansal, J. Duarte, H. Su, B. Orzari, T. Tomei, M. Pierini,
M. Touranakou, J. R. Vlimant and D. Gunopulos
[*NeurIPS 2021*]

# ANOMALY DETECTION

Graph Embedding

- GNN based autoencoders for anomaly detection
  - enables automated and model-agnostic new physics search



Recluster R=1.5 jets to R=0.1 microjets and use as inputs

$$L_{node} = \sqrt{\sum_{ia} \frac{(\hat{x}_i^a - x_i^a)^2}{N \times 5}}, \qquad L_{edge} = \sum_a \sqrt{\sum_{ij} \frac{(\hat{A}_{ij}^a - A_{ij}^a)^2}{N \times N}}$$

$$L_{auto} = \lambda_{node} L_{node} + \lambda_{edge} L_{edge}$$

# THE ROAD AHEAD

# THE ROAD AHEAD

- Can we better incorporate physics knowledge into the network design?

  - physics aware data representation, symmetry group equivariant architecture, …

# JETS IN THE LUND PLANE

- The Lund jet plane provides an efficient description of the radiation patterns within a jet



- each emission (splitting) is mapped to a point in the 2D (angle, transverse momentum) plane

  - further emissions (of the secondary particles) are represented in additional leaf planes

- different kinematic regimes are clearly separated in the Lund plane

- a natural input for ML algorithms on jets since it essentially encodes the full radiation patterns of a jet

# LundNet

- LundNet: a graph neural network based on the Lund jet plane

  - technically, the input is a binary tree (from Cambridge/Aachen clustering)

    - equivalent to the **full** Lund plane

  - each node corresponds to an emission

    - a set of variables are be defined for the current splitting

$$\Delta^2 = (y_a - y_b)^2 + (\phi_a - \phi_b)^2, \quad k_t \equiv p_{tb}\Delta_{ab}, \quad m^2 \equiv (p_a + p_b)^2,$$

$$z \equiv \frac{p_{tb}}{p_{ta} + p_{tb}}, \quad \kappa \equiv z\Delta, \quad \psi \equiv \tan^{-1}\frac{y_b - y_a}{\phi_b - \phi_a},$$

- Similar network architecture as ParticleNet

  - but the graph structure is fixed by the Lund tree

    - instead of the (dynamic) k-nearest neighbors

- Two variants of LundNet studied

  - LundNet-5: using all five Lund variables, $\quad (\ln k_t, \ln \Delta, \ln z, \ln m, \psi)$

  - LundNet-3: using only three Lund variables, $\quad (\ln k_t, \ln \Delta, \ln z)$

# LUNDNET: PERFORMANCE

- LundNet achieves very high performance at significant lower computational cost than ParticleNet

  - due to fewer number of neighbors in a binary tree & static graph structure

- Moreover, LundNet provides a systematic way to control the robustness of the tagger

  - the non-perturbative region can be effectively rejected by applying a $k_t$ cut on the Lund plane



QCD rejection v. Top tagging efficiency

Pythia 8.223 simulation
signal: $pp \to t\bar{t}$, background: $pp \to jj$
anti-$k_t$ $R = 1$ jets, $p_t > 500$ GeV

*better*

*Top tagging*

— LundNet-5
— LundNet-3
— RecNN (LCBC '17)
— Lund+LSTM (DSS '18)
— ParticleNet (QG '19)

| | Number of parameters | Training time [ms/sample/epoch] | Inference time [ms/sample] |
|---|---|---|---|
| LundNet | 395k | 0.472 | 0.117 |
| ParticleNet | 369k | 3.488 | 1.036 |
| Lund+LSTM | 67k | 0.424 | 0.131 |

performance v. resilience

LundNet-3
LundNet-5
RecNN (LCBC '17)
Lund+LSTM (DSS '18)
ParticleNet (QG '19)

Pythia 8.223 simulation
signal: pp->WW, background: pp->jj
anti $k_t$ R=1 jets, $p_t$>2 TeV, $\varepsilon_W$=0.7

*increasing $k_t$*

ln $k_t$ cuts: ∅, -2, -1.5, -1, -0.5, 0, 0.5, 1, 1.5, 2

**Better performance**

**Better resilience to non-pert. effects** →

Primary Lund-plane regions

$\ln(k_t/\text{GeV})$

ISR (large $\Delta$)

hard-collinear (large $z$)

soft-collinear

MPI/UE

non-pert. (small $k_t$)

*increasing $k_t$*

$\ln(1/\Delta)$

*\* Resilience assessed by applying
the model trained on hadron-level
samples to parton-level samples
and compare the difference*

# LORENTZNET

- Incorporating Lorentz symmetry into graph neural network architecture

Coordinate input: $x^0$

Feature input: $h_i^0$

**Lorentz 4-vector**
**Lorentz scalar**

Message: $m_{ij}^l = \phi_e\left(h_i^l, h_j^l, \psi(\|x_i^l - x_j^l\|^2), \psi(\langle x_i^l, x_j^l \rangle)\right)$

**Scalars**     **Pairwise Lorentz invariants**

Coordinate update: $x_i^{l+1} = x_i^l + c \sum_{j\in[N]} \phi_x(m_{ij}^l) \cdot x_j^l$

Feature update: $h_i^{l+1} = h_i^l + \phi_h(h_i^l, \sum_{j\in[N]} w_{ij} m_{ij}^l)$



MLP     Sum Pooling     Minkowski Norm & Inner Product

**Lorentz Group Equivariant Block (LGEB)**     **LorentzNet**

*cf. A. Bogatskiy, B. Anderson, J. Offermann, M. Roussi, D. Miller and R. Kondor, arXiv: 2006.04780 ["LGN"];*
*A. Bogatskiy, T. Hoffman, D. Miller and J. Offermann, arXiv: 2211.00454 ["PELICAN"];*
*I. Batatia, M. Geiger, J. Munoz, T. Smidt, L. Silberman and C. Ortner, arXiv: 2306.00091 ["lie-nn"];*

*GNNs for Particle Physics – July 2023 – Huilin Qu (CERN)*

61

# LORENTZNET: BENEFITS FROM SYMMETRY

- Benefits from the symmetry preservation

  - model response invariant under Lorentz transformation

  - sample efficiency: incorporation of Lorentz symmetry allows to train with very few samples



Model stability under Lorentz boost

Performance when trained on a fraction of the top-tagging dataset

| Training Fraction | Model | Accuracy | AUC | $1/\varepsilon_B$ $(\varepsilon_S = 0.5)$ | $1/\varepsilon_B$ $(\varepsilon_S = 0.3)$ |
|---|---|---|---|---|---|
| 0.5% (~6k jets) | ParticleNet | 0.913 | 0.9687 | $77 \pm 4$ | $199 \pm 14$ |
| | LorentzNet | **0.929** | **0.9793** | **$176 \pm 14$** | **$562 \pm 72$** |
| 1% | ParticleNet | 0.919 | 0.9734 | $103 \pm 5$ | $287 \pm 19$ |
| | LorentzNet | **0.932** | **0.9812** | **$209 \pm 5$** | **$697 \pm 58$** |
| 5% | ParticleNet | 0.931 | 0.9807 | $195 \pm 4$ | $609 \pm 35$ |
| | LorentzNet | **0.937** | **0.9839** | **$293 \pm 12$** | **$1108 \pm 84$** |

# THE ROAD AHEAD

- Can we better incorporate physics knowledge into the network design?
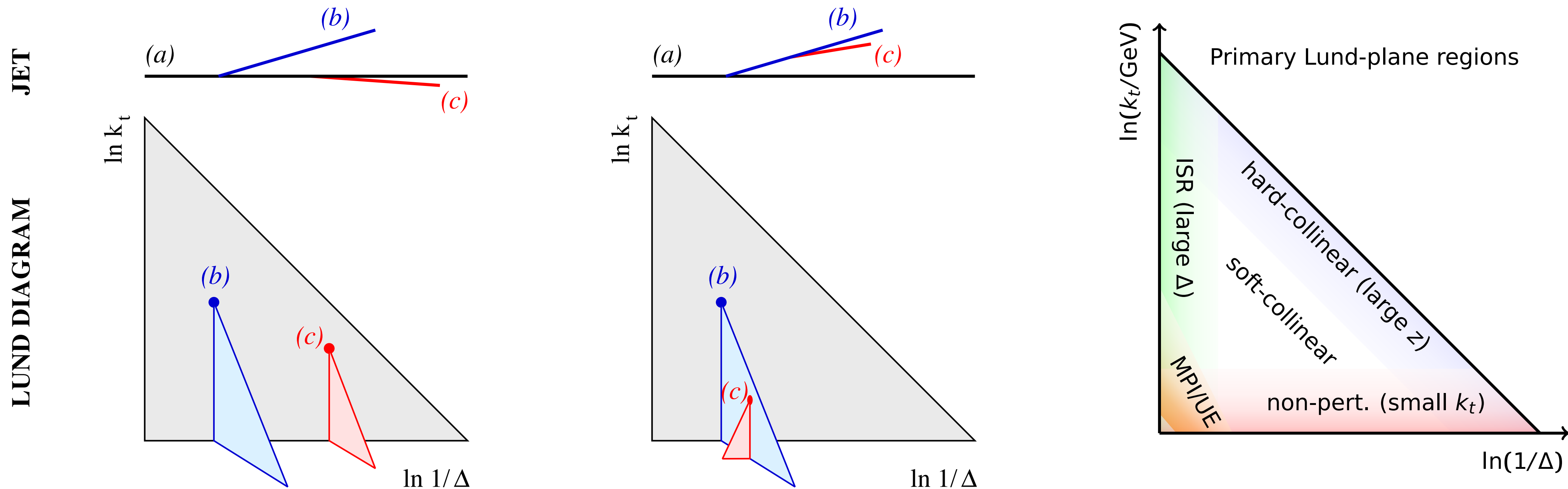
  - physics aware data representation, symmetry group equivariant architecture, …

- Can we scale up to a large model for HEP?

  - large datasets, pre-training, multi-modal learning, …

# LARGE PHYSICS MODEL?

- Large Language Models (like GPT) has transformed NLP. How about a Large Physics Model?



*R. Das, G. Kasieczka and D. Shih, arXiv: 2212.00046*

# A First Step: New Dataset

- **JetClass**: a new large and comprehensive jet simulation dataset
  - 100M jets in 10 classes: ~two orders of magnitude larger than existing public datasets

Embe

Linear

LN

MHA

LN

concat

$x_{\text{class}}$    $\mathbf{x}^L$

**(c) Class Attention Block**

*L* blocks

*Class token* ⊙

rticles → Embedding → $\mathbf{x}^0$ → Particle Attention Block → $\mathbf{x}^1$ → Particle Attention Block → $\mathbf{x}^{L-1}$ → Particle Attention Block → $\mathbf{x}^L$ → Class Attention Block → Class Attention Block → MLP → SoftMax →

ctions → Embedding → $\mathbf{U}$

**(a) Particle Transformer**

MatMul

$V$

SoftMax

$\mathbf{U}$ → ⊕

Mask

**P-MHA**

MatMul

$V$

SoftMax

⊕

Mask

Scale

MatMul

$Q$ ↑    $K$ ↑

Linear   Linear   Linear

$\mathbf{x}$

$$\Delta = \sqrt{(y_a - y_b)^2 + (\phi_a - \phi_b)^2},$$
$$k_{\text{T}} = \min(p_{\text{T},a}, p_{\text{T},b})\Delta,$$
$$z = \min(p_{\text{T},a}, p_{\text{T},b})/(p_{\text{T},a} + p_{\text{T},b}),$$
$$m^2 = (E_a + E_b)^2 - \|\mathbf{p}_a + \mathbf{p}_b\|^2,$$

*and many other possible
pairwise features…*

Linear          Linear

LN    Mask   LN

⊕            ⊕

**P-MHA**           **MHA**

$$\text{P-MHA}(Q, K, V) = \text{SoftMax}(QK^T/\sqrt{d_k} + \mathbf{U})V,$$

Scale

LN     MatMul    LN

$\mathbf{x}^{l-1}$    $Q$ ↑    concat    $\mathbf{U}$

$x_{\text{class}}$    $K$

*Injection of (physics-inspired) pairwise features to
"bias" the dot-product self-attention*

**(b) Particle Attention Block**    **(c) Class Attention Block**

Linear      Linear      Linear

# PARTICLE TRANSFORMER: PERFORMANCE

| | All classes | | $H \to b\bar{b}$ | $H \to c\bar{c}$ | $H \to gg$ | $H \to 4q$ | $H \to \ell\nu qq'$ | $t \to bqq'$ | $t \to b\ell\nu$ | $W \to qq'$ | $Z \to q\bar{q}$ |
| | Accuracy | AUC | $\text{Rej}_{50\%}$ | $\text{Rej}_{50\%}$ | $\text{Rej}_{50\%}$ | $\text{Rej}_{50\%}$ | $\text{Rej}_{99\%}$ | $\text{Rej}_{50\%}$ | $\text{Rej}_{99.5\%}$ | $\text{Rej}_{50\%}$ | $\text{Rej}_{50\%}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| PFN | 0.772 | 0.9714 | 2924 | 841 | 75 | 198 | 265 | 797 | 721 | 189 | 159 |
| P-CNN | 0.809 | 0.9789 | 4890 | 1276 | 88 | 474 | 947 | 2907 | 2304 | 241 | 204 |
| ParticleNet | 0.844 | 0.9849 | 7634 | 2475 | 104 | 954 | 3339 | 10526 | 11173 | 347 | 283 |
| **ParT** | **0.861** | **0.9877** | **10638** | **4149** | **123** | **1864** | **5479** | **32787** | **15873** | **543** | **402** |
| ParT (plain) | 0.849 | 0.9859 | 9569 | 2911 | 112 | 1185 | 3868 | 17699 | 12987 | 384 | 311 |

JETCLASS dataset (100M jets)

- Particle Transformer (ParT): significant performance improvement!

  - compared to the existing state-of-the-art, ParticleNet

    - 1.7% increase in accuracy

    - **up to 3x increase in background rejection** ($\text{Rej}_{X\%}$) ┄┄┄➤ $$\text{Rej}_{X\%} \equiv 1/\text{FPR at TPR} = X\%,$$

# PARTICLE TRANSFORMER: PERFORMANCE

| | All classes | | $H \to b\bar{b}$ | $H \to c\bar{c}$ | $H \to gg$ | $H \to 4q$ | $H \to \ell\nu qq'$ | $t \to bqq'$ | $t \to b\ell\nu$ | $W \to qq'$ | $Z \to q\bar{q}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | AUC | $\mathrm{Rej}_{50\%}$ | $\mathrm{Rej}_{50\%}$ | $\mathrm{Rej}_{50\%}$ | $\mathrm{Rej}_{50\%}$ | $\mathrm{Rej}_{99\%}$ | $\mathrm{Rej}_{50\%}$ | $\mathrm{Rej}_{99.5\%}$ | $\mathrm{Rej}_{50\%}$ | $\mathrm{Rej}_{50\%}$ |
| PFN | 0.772 | 0.9714 | 2924 | 841 | 75 | 198 | 265 | 797 | 721 | 189 | 159 |
| P-CNN | 0.809 | 0.9789 | 4890 | 1276 | 88 | 474 | 947 | 2907 | 2304 | 241 | 204 |
| ParticleNet | 0.844 | 0.9849 | 7634 | 2475 | 104 | 954 | 3339 | 10526 | 11173 | 347 | 283 |
| **ParT** | **0.861** | **0.9877** | **10638** | **4149** | **123** | **1864** | **5479** | **32787** | **15873** | **543** | **402** |
| ParT (plain) | 0.849 | 0.9859 | 9569 | 2911 | 112 | 1185 | 3868 | 17699 | 12987 | 384 | 311 |

- Particle Transformer (ParT): significant performance improvement!

  - compared to the existing state-of-the-art, ParticleNet

    - 1.7% increase in accuracy

    - **up to 3x increase in background rejection** ($\mathrm{Rej}_{X\%}$)

- ParT (plain): plain Transformer w/o interaction features

  - 1.2% drop in accuracy compared to full ParT

  - **Physics-driven modification of self-attention plays a key role!**

JETCLASS dataset (100M jets)

Model complexity

| | Accuracy | # params | FLOPs |
|---|---|---|---|
| PFN | 0.772 | 86.1 k | 4.62 M |
| P-CNN | 0.809 | 354 k | 15.5 M |
| ParticleNet | 0.844 | 370 k | 540 M |
| **ParT** | **0.861** | 2.14 M | 340 M |
| ParT (plain) | 0.849 | 2.13 M | 260 M |

# PARTICLE TRANSFORMER: PRE-TRAINING + FINE-TUNING

- The large Transformer-based model enables new training paradigm

  - (supervised) pre-training on a large dataset (e.g., JETCLASS) & fine-tuning to downstream tasks

  - significantly outperforms existing models

*Top quark tagging benchmark (~2M jets)* [SciPost Phys. 7 (2019) 014]

| | Accuracy | AUC | $\text{Rej}_{50\%}$ | $\text{Rej}_{30\%}$ |
|---|---|---|---|---|
| P-CNN | 0.930 | 0.9803 | $201 \pm 4$ | $759 \pm 24$ |
| PFN | — | 0.9819 | $247 \pm 3$ | $888 \pm 17$ |
| ParticleNet | 0.940 | 0.9858 | $397 \pm 7$ | $1615 \pm 93$ |
| JEDI-net (w/ $\sum O$) | 0.930 | 0.9807 | — | 774.6 |
| PCT | 0.940 | 0.9855 | $392 \pm 7$ | $1533 \pm 101$ |
| LGN | 0.929 | 0.964 | — | $435 \pm 95$ |
| rPCN | — | 0.9845 | $364 \pm 9$ | $1642 \pm 93$ |
| LorentzNet | 0.942 | 0.9868 | $498 \pm 18$ | $2195 \pm 173$ |
| ParT | 0.940 | 0.9858 | $413 \pm 16$ | $1602 \pm 81$ |
| ParticleNet-f.t. | 0.942 | 0.9866 | $487 \pm 9$ | $1771 \pm 80$ |
| **ParT-f.t.** | **0.944** | **0.9877** | **$691 \pm 15$** | **$2766 \pm 130$** |

*Quark-gluon tagging benchmark (~2M jets)* [JHEP 01 (2019) 121]

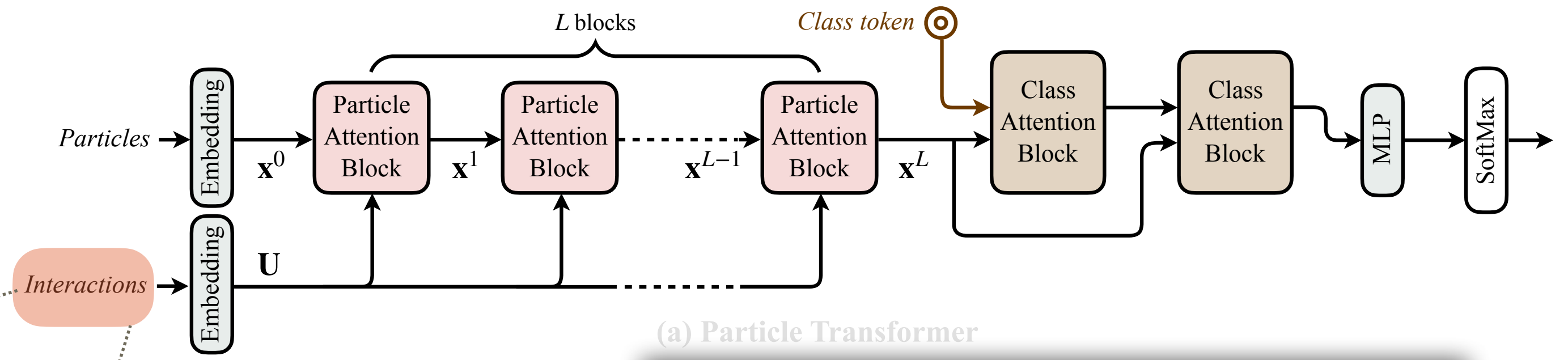| | Accuracy | AUC | $\text{Rej}_{50\%}$ | $\text{Rej}_{30\%}$ |
|---|---|---|---|---|
| $\text{P-CNN}_{\text{exp}}$ | 0.827 | 0.9002 | 34.7 | 91.0 |
| $\text{PFN}_{\text{exp}}$ | — | 0.9005 | $34.7 \pm 0.4$ | — |
| $\text{ParticleNet}_{\text{exp}}$ | 0.840 | 0.9116 | $39.8 \pm 0.2$ | $98.6 \pm 1.3$ |
| $\text{rPCN}_{\text{exp}}$ | — | 0.9081 | $38.6 \pm 0.5$ | — |
| $\text{ParT}_{\text{exp}}$ | 0.840 | 0.9121 | $41.3 \pm 0.3$ | $101.2 \pm 1.1$ |
| $\text{ParticleNet-f.t.}_{\text{exp}}$ | 0.839 | 0.9115 | $40.1 \pm 0.2$ | $100.3 \pm 1.0$ |
| **$\text{ParT-f.t.}_{\text{exp}}$** | **0.843** | **0.9151** | **$42.4 \pm 0.2$** | **$107.9 \pm 0.5$** |
| $\text{PFN}_{\text{full}}$ | — | 0.9052 | $37.4 \pm 0.7$ | — |
| $\text{ABCNet}_{\text{full}}$ | 0.840 | 0.9126 | $42.6 \pm 0.4$ | $118.4 \pm 1.5$ |
| $\text{PCT}_{\text{full}}$ | 0.841 | 0.9140 | $43.2 \pm 0.7$ | $118.0 \pm 2.2$ |
| $\text{LorentzNet}_{\text{full}}$ | 0.844 | 0.9156 | $42.4 \pm 0.4$ | $110.2 \pm 1.3$ |
| $\text{ParT}_{\text{full}}$ | 0.849 | 0.9203 | $47.9 \pm 0.5$ | $129.5 \pm 0.9$ |
| **$\text{ParT-f.t.}_{\text{full}}$** | **0.852** | **0.9230** | **$50.6 \pm 0.2$** | **$138.7 \pm 1.3$** |

# Going Beyond?

- **JetClass**: a new large and comprehensive jet simulation dataset
  - 100M jets in 10 classes: ~two orders of magnitude larger than existing public datasets

HQ, C. Li, S. Qian,
*ICML 2022*



*We invite the community to explore and experiment with this dataset and extend the boundary of deep learning and HEP even further.*

# THE ROAD AHEAD

- Can we better incorporate physics knowledge into the network design?

  - physics aware data representation, symmetry group equivariant architecture, …

- Can we scale up to a large model for HEP?

  - large datasets, pre-training, multi-modal learning, …

- Can we improve the computational efficiency of GNNs?

  - emerging specialized libraries for GNN training and inference (PyG, DGL, TF-GNN, …)

  - accelerated inference on specialized ASICs / FPGAs (e.g., for triggering), software hardware co-design, …

- Can we improve the robustness of GNNs (e.g., data/simulation difference)?

  - domain adaption? calibration? uncertainty aware training? …

- Can we improve the interpretability and explainability of GNNs?

# THE ROAD AHEAD

- Can we better incorporate physics knowledge into the network design?

  - physics aware data representation, symmetry group equivariant architecture, …

- Can we scale up to a large model for HEP?

  - large datasets, pre-training, multi-modal learning, …

- Can we impro…

  - emerging s…                 …nes for GNN training and inference (PyG, DGL, TF-GNN, …)

  - accelerated inference on specialized ASICs / FPGAs (e.g., for triggering), software hardware co-design, …

- Can we improve the robustness of GNNs (e.g., data/simulation difference)?

  - domain adaption? calibration? uncertainty aware training? …

- Can we improve the interpretability and explainability of GNNs?

*Your innovation and creativity can make a big difference!*

# *Extra: Partical Jet Tagging (in CMS)*

# BOOSTED JET TAGGING

- Hadronic decays of highly Lorentz-boosted heavy particles (Higgs/W/Z/top) lead to large-radius jets with distinctive characteristics:

  - different radiation patterns ("**substructure**")

    - 3-prong (top), 2-prong (W/Z/H) vs 1-prong (gluon/light quark jet)

  - different **flavor** content: existence of one or more b-/c-quarks

- Boosted jet tagging:

  - simultaneously exploiting both **substructure** and **flavor** to maximize the performance

  - significant performance leap thanks to deep learning techniques

# DᴇᴇᴘAK8

- Advanced deep learning-based algorithm for boosted jet tagging, using AK8 (anti-$k_T$ R=0.8) jets

  - **multi-class classifier** for top quark and W, Z, Higgs boson tagging

  - **directly uses jet constituents** (particle-flow candidates / secondary vertices)

  - **1D convolutional neural network (CNN)**, based on the ResNet [arXiv: 1512.03385] architecture

## Inputs

**Substructure**  **Flavour**

### Particles
- Up to 100 PF candidates(*)
- Sorted in descending $p_T$ order
- Uses basic kinematic variables, Puppi weights, and track properties (quality, covariance, displacement, etc.)

### Secondary vertices
- Up to 7 SVs(*) (inside jet cone)
- Sorted in descending $S_{IP2D}$ order
- Uses SV kinematics and properties (quality, displacement, etc.)

*(*) Number chosen to include all candidates for ≥ 90% of the events*

## Output

| Category | Label |
|----------|-------|
| Higgs | H (bb) |
| | H (cc) |
| | H (VV*→qqqq) |
| Top | top (bcq) |
| | top (bqq) |
| | top (bc) |
| | top (bq) |
| W | W (cq) |
| | W (qq) |
| Z | Z (bb) |
| | Z (cc) |
| | Z (qq) |
| QCD | QCD (bb) |
| | QCD (cc) |
| | QCD (b) |
| | QCD (c) |
| | QCD (others) |

## Architecture

### Particles
*features* — *filter* — *particles, ordered by $p_T$* → 1D CNN (14 layers)

### Secondary Vertices
*features* — *filter* — *SVs, ordered by $S_{IP2D}$* → 1D CNN (10 layers)

→ Fully connected (1 layer) → Output

## Top quark tagging
(13 TeV)

**CMS**
*Simulation*
**Top quark vs. QCD multijet**
$1000 < p_T^{gen} < 1500$ GeV, $|\eta^{gen}| < 2.4$
$105 < m_{SD}^{AK8} < 210$ GeV
$110 < m_{SD}^{CA15} < 210$ GeV
$140 < m_{HOTVR} < 220$ GeV

JINST 15 (2020) P06005

*better*

- DeepAK8
- DeepAK8-MD
- ImageTop
- ImageTop-MD
- $m_{SD} + \tau_{32}$
- $m_{SD} + \tau_{32} + b$
- BEST
- HOTVR
- $N_3$-BDT (CA15)

Background efficiency — Signal efficiency

75

# PARTICLENET

- ParticleNet: jet tagging via particle clouds

  - treating a jet as an **unordered set of particles**, distributed in the η − φ space

  - **graph neural network architecture**, adapted from Dynamic Graph CNN [arXiv:1801.07829]

    - treating a point cloud as a graph: each point is a vertex

      - for each point, a local patch is defined by finding its k-nearest neighbors

    - designing a permutation-invariant "convolution" function

      - define "edge feature" for each center-neighbor pair: $e_{ij}$ = MLP($x_i$, $x_j$)

      - aggregate the                                way:

*ParticleNet architecture*

*cf.* P.T. Komiske, E. M. Metodiev and J. Thaler, *JHEP 01 (2019) 121*;
V. Mikuni and F. Canelli, *Eur. Phys. J. Plus 135, 463 (2020)*; *Mach.Learn.Sci.Tech. 2 (2021) 3, 035027*.

GNNs for Particle Physics  – July 2023 – Huilin Qu (CERN)

# CORRELATION WITH THE JET MASS

Plain training:
no mass decorrelation



Background jet mass



- One feature of these taggers is the correlation with the jet mass

  - jet mass shape of the background becomes similar to that of the signal after selection with the tagger: **"mass sculpting"**

  - not necessarily a problem, but a mass-independent tagger is often more desirable:

    - allows to use the mass variable to further separate signal and background

    - enables tagging signal jets with an unknown mass

    - ...

77

Plain training:
no mass decorrelation

Background sample
Signal sample (fixed mass)

A.U.

Jet mass

Mass-decorrelated DeepAK8:
"adversarial training"

Nominal DeepAK8

Feature extractor

ID CNN

Classifier

Fully Connected

back propagation

Classification output

Joint loss
$L = L_C - \lambda L_{MP}$

Mass predictor

Fully Connected

Mass prediction

Loss
$L_{MP}$

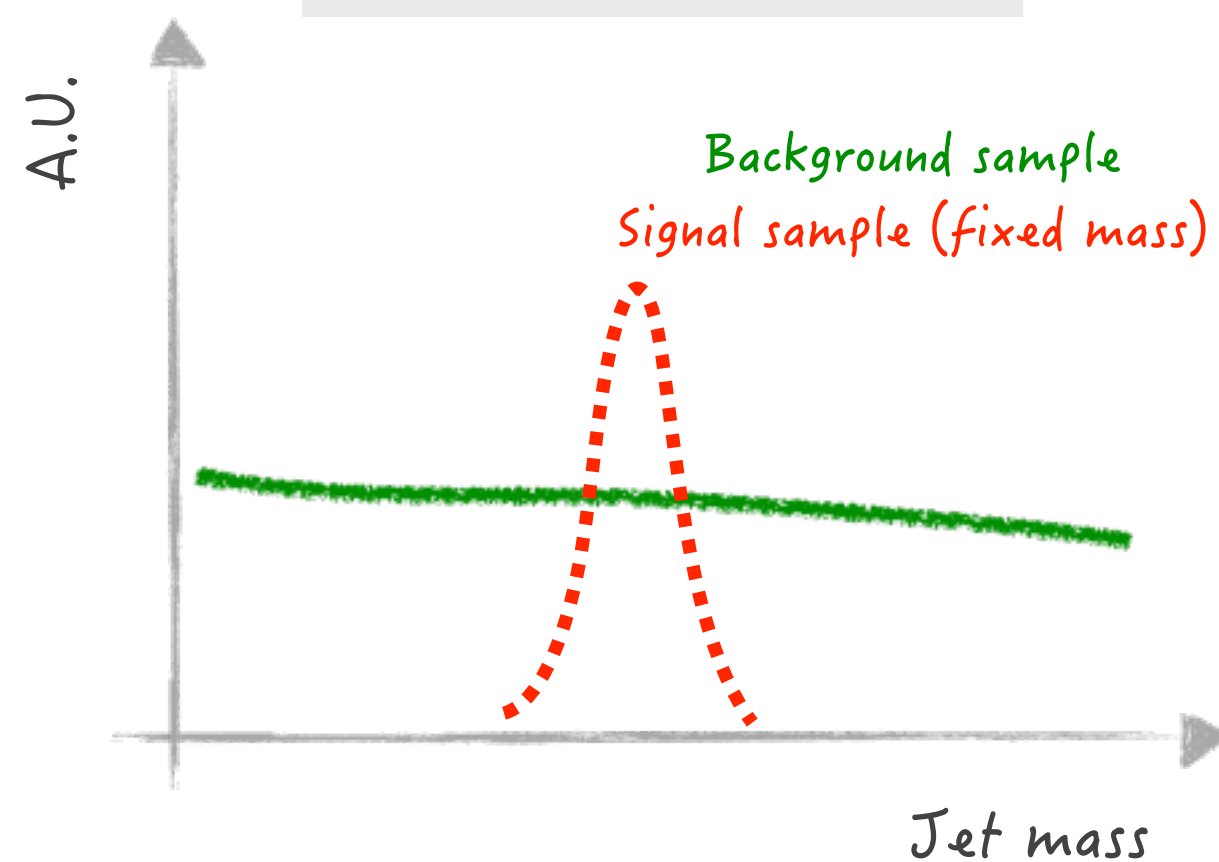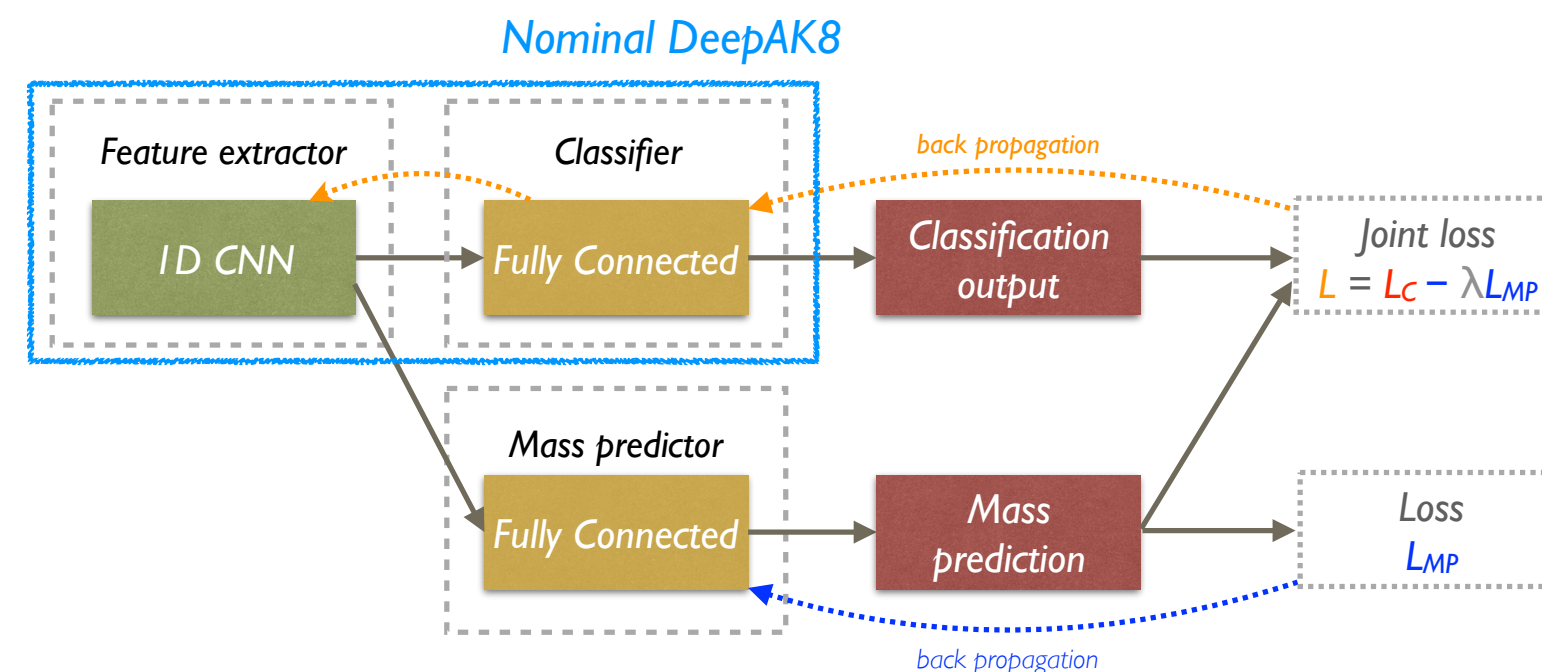back propagation

Mass-decorrelated ParticleNet:
training with variable-mass signal

A.U.

Background sample
Signal sample (variable mass)

Jet mass

Background jet mass

(13 TeV)

CMS
Simulation Preliminary

Dijet sample

H→cc̄ tagging: DeepAK8

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

Inclusive
$\varepsilon_B = 5\%$
$\varepsilon_B = 1\%$
$\varepsilon_B = 0.5\%$

$m_{SD}$ [GeV]

Background jet mass

(13 TeV)

CMS
Simulation Preliminary

Dijet sample

H→cc̄ tagging: DeepAK8-MD

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

Inclusive
$\varepsilon_B = 5\%$
$\varepsilon_B = 1\%$
$\varepsilon_B = 0.5\%$

$m_{SD}$ [GeV]

Background jet mass

(13 TeV)

CMS
Simulation Preliminary

Dijet sample

H→cc̄ tagging: ParticleNet-MD

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

Inclusive
$\varepsilon_B = 5\%$
$\varepsilon_B = 1\%$
$\varepsilon_B = 0.5\%$

$m_{SD}$ [GeV]

# PERFORMANCE COMPARISON

H→bb tagging performance

- ParticleNet-MD

  - using a special signal sample for training

    - hadronic decays of a spin-0 particle X

      - $X \to bb, X \to cc, X \to qq$

    - not a fixed mass, but a flat mass spectrum

      - $m(X) \in [15, 250]$ GeV

    - allows to easily reweight both signal and background to a ~flat 2D distribution in $(p_T, mass)$ for the training

- ParticleNet-MD shows the best performance

  - ~3-4x better background rejection compared to DeepAK8-MD (based on "adversarial training")

  - only slight performance loss compared to the nominal version w/o mass decorrelation

# MASS REGRESSION

- Jet mass: one of the most powerful observables for boosted jet tagging

  - characteristic mass peak for top/W/Z/H jets v.s. continuum for QCD jets

- Mass regression:

  - exploit deep learning to reconstruct jet mass with the highest possible resolution

  - training setup similar to the ParticleNet tagger

    - but: predict the jet mass directly from the jet consitituents

- Regression target:

  - signal (X → bb/cc/qq): generated particle mass of X [flat spectrum in 15 – 250 GeV]

  - background (QCD) jets: soft drop mass of the generated particle-level jet

- Loss function

  - LogCosh:   $L(y, y^p) = \sum_{i=1}^{n} \log(\cosh(y_i^p - y_i))$



CMS, JINST 15 (2020) P06005

# MASS REGRESSION: PERFORMANCE

**CMS Simulation Preliminary**

*anti-$k_T$ jets*

······· H -> cc (soft drop)

——— H -> cc (regression)

*R = 1.5*

*$p_T$ > 200 GeV*

Jet mass response: H→cc jets

~50% improvement in resolution

event fraction

$M_{reco}$ / $M_{target}$

**CMS Simulation Preliminary**

A.U.

*QCD sample (jets R = 1.5)*

—— inclusive

—— $\varepsilon_B$ = 5 %

—— $\varepsilon_B$ = 1 %

—— $\varepsilon_B$ = 0.5 %

$H \to c\bar{c}$ ParticleNet tagger

$p_T$ > 200 GeV

Regressed mass vs Tagger WP

Minimal mass sculpting on background QCD jets

$M_{reg}$ [GeV/c$^2$]

*Jet Tagging in the Era of Deep Learning - June 9, 2022 - Huilin Qu (CERN)*

81

# TAGGER CALIBRATION IN DATA

- Crucial to calibrate these taggers in real data for them to be used in analyses

- Top/W tagging efficiency



- measured using the single-μ sample enriched in semi-leptonic ttbar events

- fit jet mass templates in the "pass" and "fail" categories simultaneously to extract efficiency in data

  - simulation-to-data scale factors SF := eff(data) / eff(MC) derived to correct the simulation

- jet mass scale and resolution scale factors can also be extracted

- Mistag rates of background jet typically derived directly from analysis-specific control regions

GNNs for Particle Physics – July 2023 – Huilin Qu (CERN)

# Calibration of the cc-tagger

❑ Need to measure ParticleNet cc-tagging efficiency in data

  ▪ no pure sample of H → cc jets (or even Z → cc) in data

  ▪ using g → cc in QCD multi-jet events as a proxy

❑ Difficulty: select a phase-space in g → cc that resembles H → cc
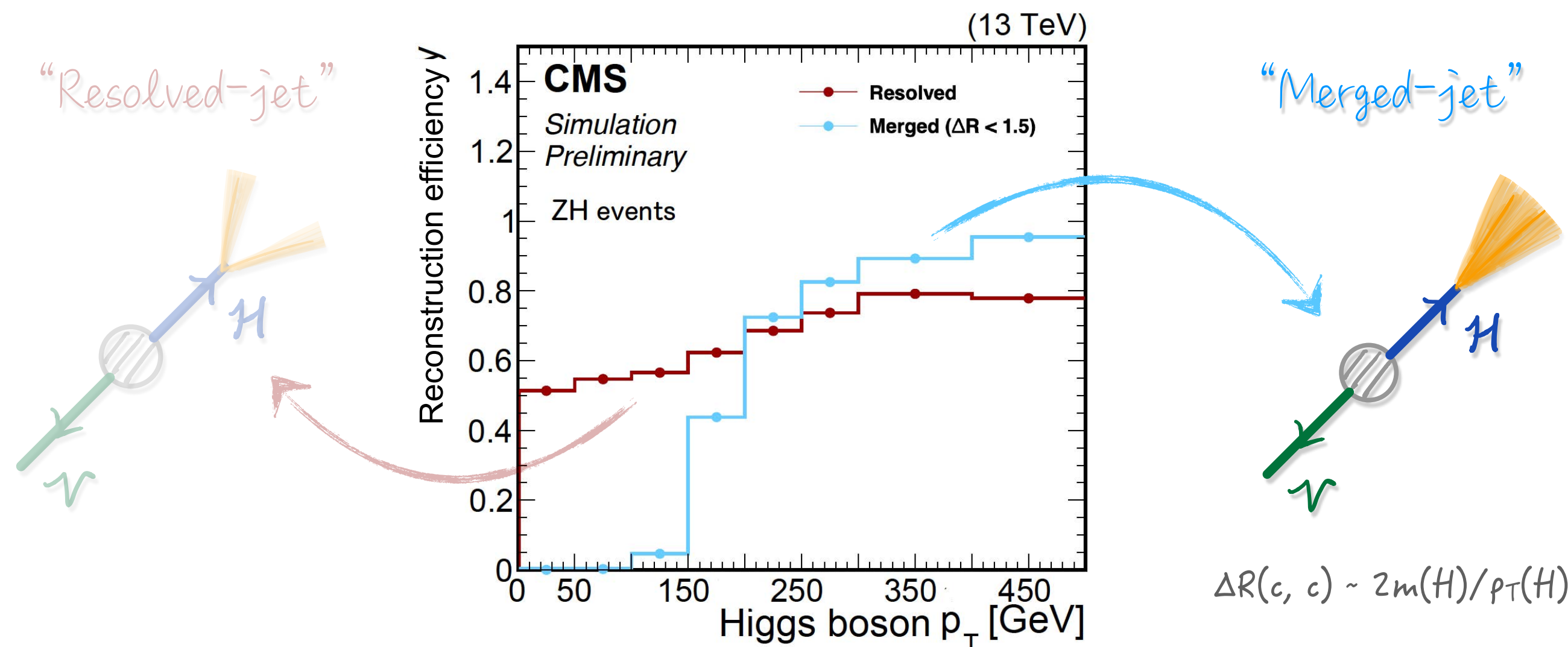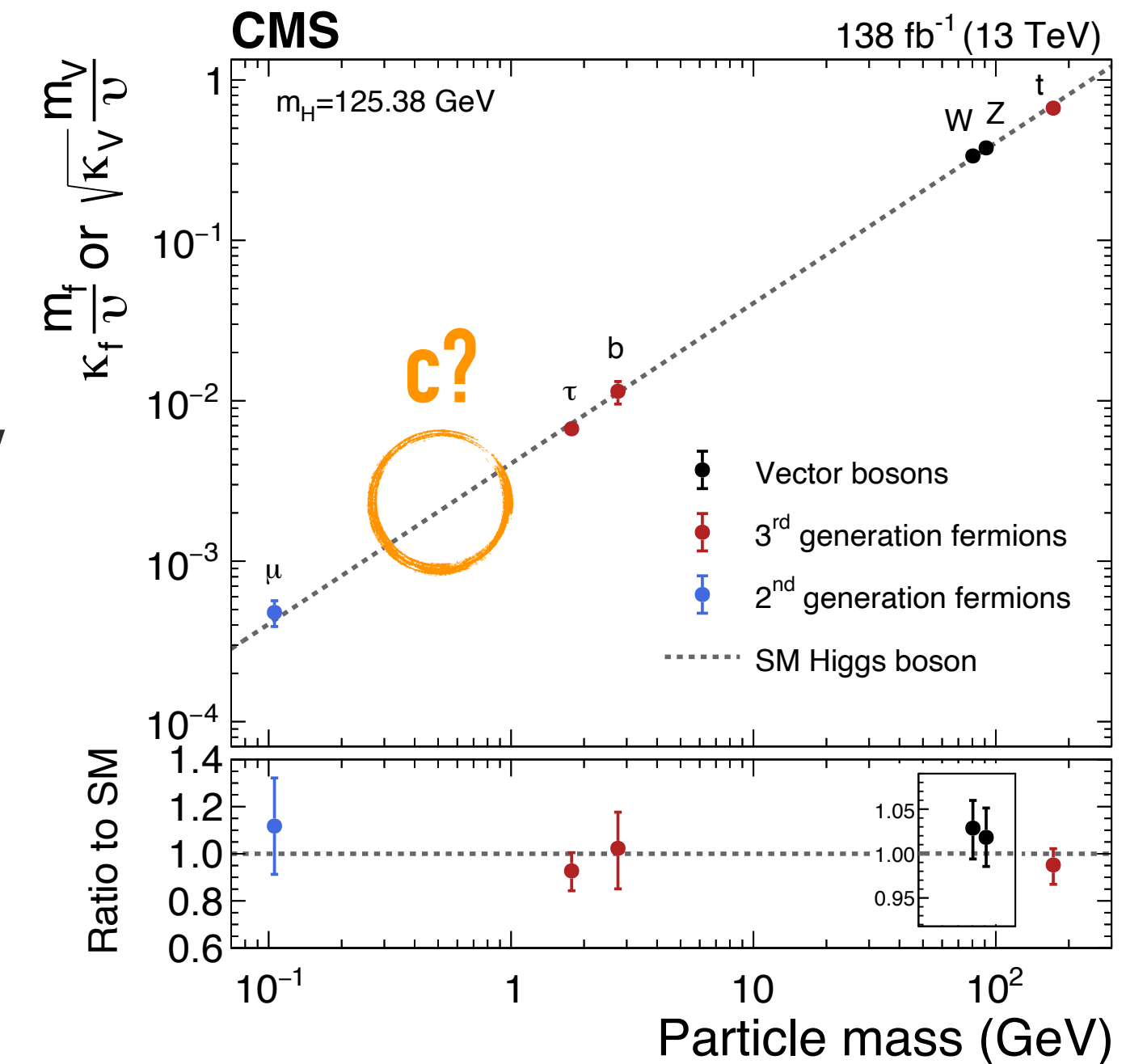
  ▪ solution: a **dedicated BDT** developed to distinguish **hard 2-prong splittings** (*i.e., high quark contribution to the jet momentum*) from **soft cc radiations** (*i.e., high gluon contribution to the jet momentum*)

  ▪ also allows to adjust the similarity between proxy and signal jets

    ▪ by varying the sfBDT cut — treated as a systematic uncertainty

❑ Perform a fit to the secondary vertex mass shapes in the "passing" and "failing" regions simultaneously to extract the scale factors

  ▪ three templates: cc (+ single c), bb (+ single b), light flavor jets

❑ Derived cc-tagging scale factors typically 0.9—1.3

  ▪ corresponding uncertainties are 20—30%



*H → cc like*

*Soft radiations: Dominant contribution!*

*Effects of the BDT*

A.U.

g→cc (all)     H→cc
g→cc (sfBDT>0.85)
g→cc (sfBDT>0.90)
g→cc (sfBDT>0.95)

**ParticleNet cc discriminant**

**Higgs-charm coupling: next milestone in Higgs physics**

- a crucial test of fermion mass generation mechanism in SM

- H→cc: extremely challenging search at the LHC

  - small branching fraction (~3%) vs enormous backgrounds — **charm tagging** is the key

**Innovative approach**: search for VH(H→cc) in the **"merged-jet" topology**

- reconstructs H→cc decay with one large-R jet (R=1.5)

- then: exploits advanced ML for H→cc identification





*"Resolved-jet"*          *"Merged-jet"*

$\Delta R(c, c) \sim 2m(H)/p_T(H)$

*Why merged-jet topology?*
- *better signal purity at higher $p_T$*
- *higher reconstruction efficiency with large-R jets*
- *better exploiting correlations between the two charm quarks — especially with deep learning*

# PARTICLENET IN ACTION: H→CC SEARCH

- **ParticleNet for H→cc jet tagging and mass reconstruction: substantial improvements**



ParticleNet tagger for H→cc tagging

*>2x improvement in final sensitivity*

ParticleNet-based jet mass regression

*~20-25% improvement in final sensitivity*

85

- **ParticleNet for H→cc jet tagging and mass reconstruction: substantial improvements**



*Most stringent limit on H→cc to date.*
- *~4x higher sensitivity than the ATLAS search*
- *Comparable to previous HL-LHC projection, but with only 5% of the data.*

*First observation of Z→cc at a hadron collider!*

# *Extra: More about practicalities*

# DISCLAIMER

- These are based on my very personal experiences in using ML to solved HEP problems

  - and highly biased to collider experiments / jet tagging

  - so please take them with a large grain of salt


- My take is that ML is 50% science and 50% engineering

  - and probably another 20% alchemy…

  - so things that should work may not necessarily work in reality…

# DATA MATTERS

- Always inspect your training data first

  - check the distributions for different classes / in different phase space ($p_T$, energy scale, vs time, …)

    - do they make sense?

    - are the trends expected?

    - do you see expected / unexpected separation power between different classes?

  - check for significant outliers / NaN / Inf / etc.


- Think carefully about how to choose your training data, how to define training target (truth labels, etc.)

  - highly case dependent, but this can have significant impact on the performance, generalization power, etc.

# DATA MATTERS (II)

- Mindful preprocessing

  - neural networks work best with "Gaussian-like" inputs

    - transform the inputs if needed, e.g., log(...) or tanh(const x ...) for long-tail distributions (energy, $p_T$, mass, $d_0/dz$, ...)

    - shift/scale the inputs, and then truncate (if needed) — extreme outliers can destabilize training and affect performance

    - use normalization layers (BatchNorm, LayerNorm, ...)

  - dealing with phase space difference between classes => reweighting (or better, sampling) if needed

  - decorrelation (e.g., mass decorrelation in jet tagging)

- Get more data whenever you can

  - if can not: consider data augmentation (rotation, reflection, smearing, dropout, ...)

# BASELINE FIRST, THEN ITERATE

- A good practice is to always establish a baseline algorithm first before developing more advanced approaches

  - with a baseline ready, then one can easily evaluate if the new algorithm is too good (to be true), or too bad (so probably missing something obvious), or just promising :)

  - if the problem is not new and a baseline already exists — just use/adapt it

  - the baseline can be a cut-based / rule-based algo, or a shallow model (e.g., BDT)

    - consider trying newer BDT libraries — even XGBoost is no longer the state-of-the-art

      - e.g., TensorFlow Decision Forests (TF-DF), LightGBM, CatBoost, ...

# SOME USEFUL LINKS

- Tutorial / hands-on textbook on Deep Learning:

  - **Dive into Deep Learning**: https://d2l.ai/

- More on Graph Neural Networks:

  - https://distill.pub/2021/gnn-intro/

  - https://distill.pub/2021/understanding-gnns/

  - and many more interactive ML posts on https://distill.pub

- HEP x ML:

  - **A Living Review of Machine Learning for Particle Physics**: https://github.com/iml-wg/HEPML-LivingReview

- My little framework:

  - **weaver** (data loading and classification/regression training pipelines): https://github.com/hqucms/weaver-core

  - weaver-examples (under construction): https://github.com/hqucms/weaver-examples