



*Proceedings of the 39th International Conference  
on Machine Learning, PMLR 162:18281-18292*

# Part

# *Particle Transformer for Jet Tagging*

ArXiv: 2202.03772, Huilin Qu (CERN), Congqiao Li & Sitian Qian (PKU)

The 21st International Workshop on Advanced Computing and Analysis Techniques

in Physics Research (**ACAT2022**)

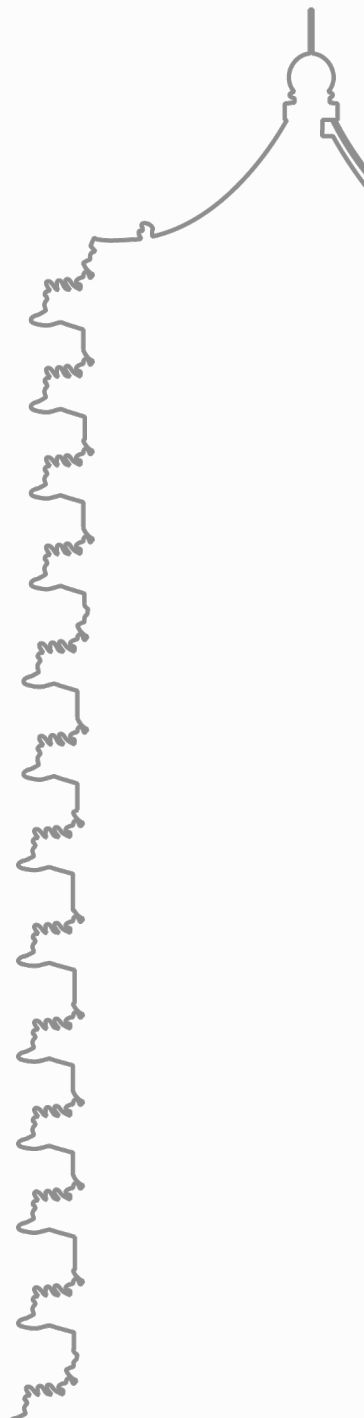
2022/10/24~2022/10/28

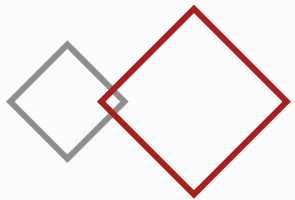
Sitian Qian (PKU)

ACAT 2022



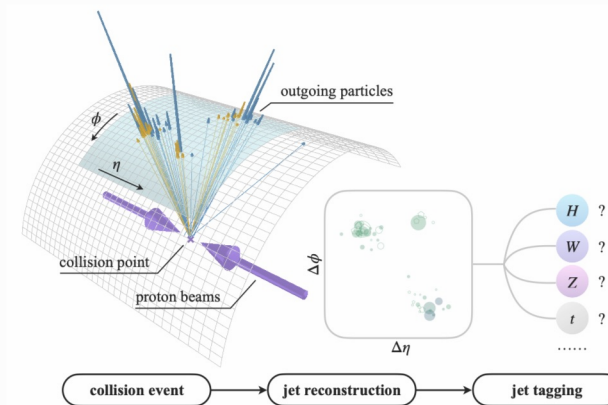
BARI

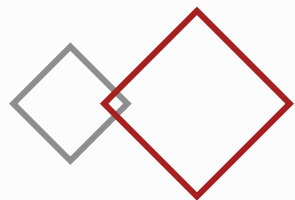




# Introduction

- Jet tagging: an ideal bridge linking HEP & ML
  - Important tool for HEP community
  - Well-defined task for ML enthusiasts
- Graph Neural Networks (GNN) + Point Cloud (PC) Representation = State-Of-The-Art (SOTA)
  - First show up in ParticleNet: previous SOTA for top tagging benchmark
  - Since then, various models are proposed under GNN+PC framework
    - ABCNet
    - ParticleNeXt
    - HEP application of Point Cloud Transformer
  - Can we add another term to the right-hand side?





# Attention Mechanism And Transformer

- Attention mechanism has drawn the attention of ML community
  - Success in various ML communities: Natural Language Processing (NLP), Computer Vision (CV)...

*BERT (1810.04805):*

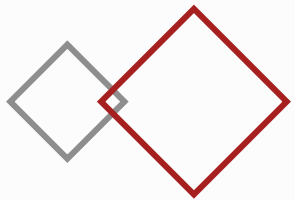
*the very first transformer, overperformed every other model even human in all kinds of NLP task!*

| System                | MNLI-(m/mm)<br>392k | QQP<br>363k | QNLI<br>108k | SST-2<br>67k | CoLA<br>8.5k | STS-B<br>5.7k | MRPC<br>3.5k | RTE<br>2.5k | Average     |
|-----------------------|---------------------|-------------|--------------|--------------|--------------|---------------|--------------|-------------|-------------|
| Pre-OpenAI SOTA       | 80.6/80.1           | 66.1        | 82.3         | 93.2         | 35.0         | 81.0          | 86.0         | 61.7        | 74.0        |
| BiLSTM+ELMo+Attn      | 76.4/76.1           | 64.8        | 79.8         | 90.4         | 36.0         | 73.3          | 84.9         | 56.8        | 71.0        |
| OpenAI GPT            | 82.1/81.4           | 70.3        | 87.4         | 91.3         | 45.4         | 80.0          | 82.3         | 56.0        | 75.1        |
| BERT <sub>BASE</sub>  | 84.6/83.4           | 71.2        | 90.5         | 93.5         | 52.1         | 85.8          | 88.9         | 66.4        | 79.6        |
| BERT <sub>LARGE</sub> | <b>86.7/85.9</b>    | <b>72.1</b> | <b>92.7</b>  | <b>94.9</b>  | <b>60.5</b>  | <b>86.5</b>   | <b>89.3</b>  | <b>70.1</b> | <b>82.1</b> |

| System                             | Dev         | Test        |
|------------------------------------|-------------|-------------|
| ESIM+GloVe                         | 51.9        | 52.7        |
| ESIM+ELMo                          | 59.1        | 59.2        |
| OpenAI GPT                         | -           | 78.0        |
| BERT <sub>BASE</sub>               | 81.6        | -           |
| BERT <sub>LARGE</sub>              | <b>86.6</b> | <b>86.3</b> |
| Human (expert) <sup>†</sup>        | -           | 85.0        |
| Human (5 annotations) <sup>†</sup> | -           | 88.0        |

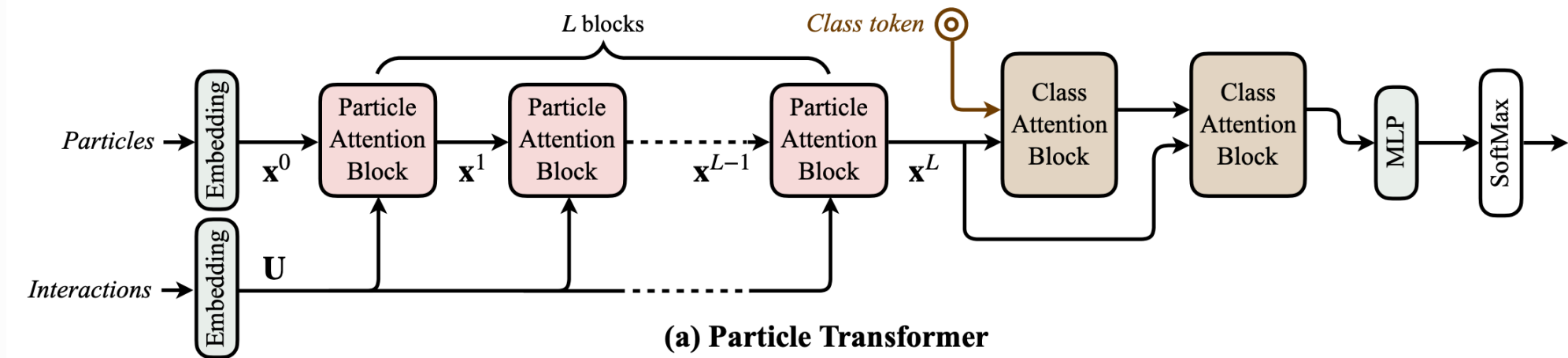
- Already several attempts in HEP context:
  - ABCNet, ParticleNeXt, Point Cloud Transformer

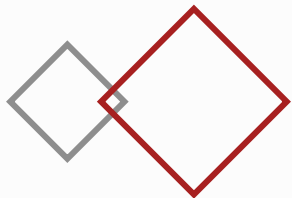




# Particle Transformer

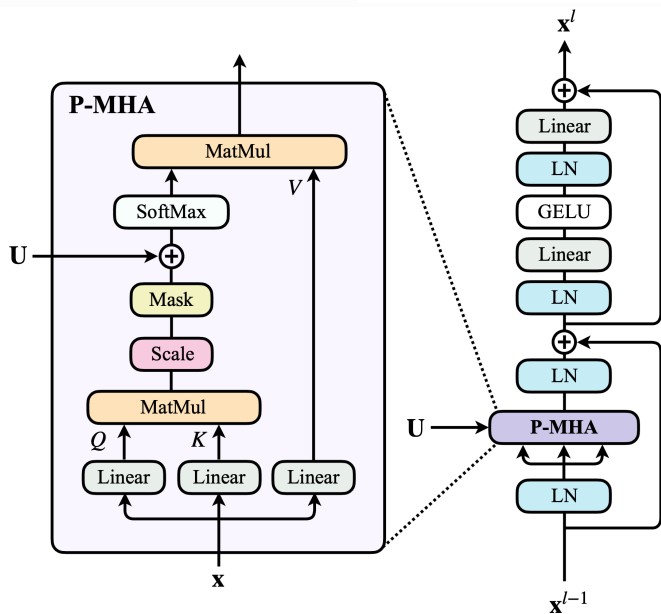
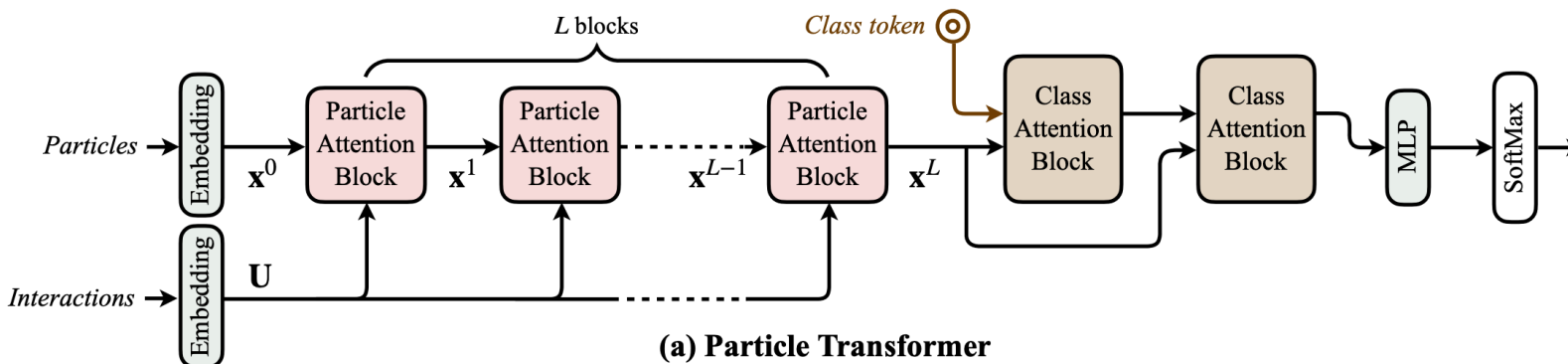
- Particle Transformer (ParT): transformer designed for particle physics
  - Input embedding: Not only inject single particle information, but also include pair-wise feature





# Particle Transformer

## Particle Attention Block



Multi-head Attention (MHA)  
powered feature extraction (embedding)

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

Kinematic Variables:  
based on LundNet

$d_k$ : dimension of  $K$

$$\Delta = \sqrt{(y_a - y_b)^2 + (\phi_a - \phi_b)^2}$$

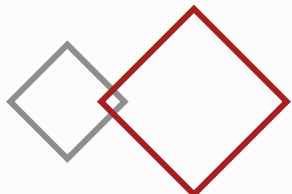
$$k_T = \min(p_{T,a}, p_{T,b}) \cdot \Delta$$

$$z = \min(p_{T,a}, p_{T,b}) / (p_{T,a} + p_{T,b})$$

$$m^2 = (E_a + E_b)^2 - \|\mathbf{p}_a + \mathbf{p}_b\|^2$$

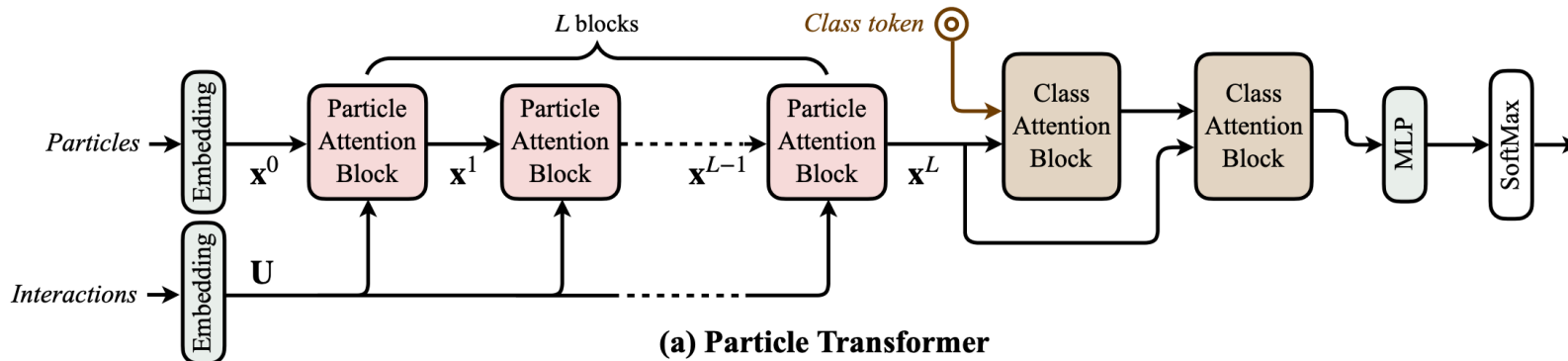
Particle interactions come in as  
pair-wise features





# Particle Transformer Class Attention Block

Output of class attention blocks will be imported to MLP + softmax for final classification scores



(a) Particle Transformer

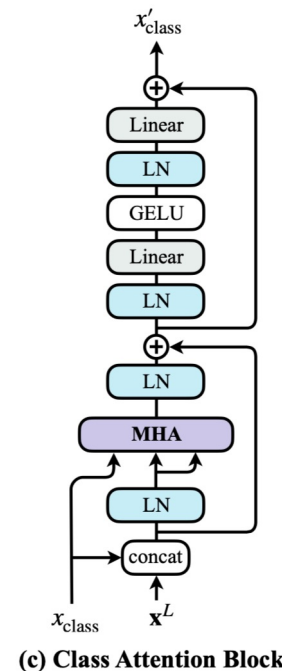
Multi-head Attention (MHA) powered, class information comes in for classification

$$\text{MHA}_C(Q_C, K_C, V_C) = \text{SoftMax}(Q_C K_C^T / \sqrt{d_{kC}}) V_C$$

$$Q_C = W_{qC} x_{\text{class}} + b_{qC} \quad K_C = W_{kC} \mathbf{z} + b_{kC} \quad V_C = W_{vC} \mathbf{z} + b_{vC} \quad d_{kC}: \text{dimension of } K_C$$

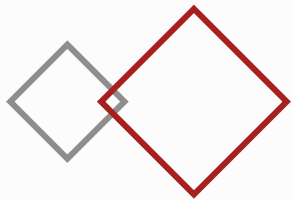
$$\mathbf{z} = [x_{\text{class}}, \mathbf{x}^L]$$

Concatenate class information and particle embedding



(c) Class Attention Block





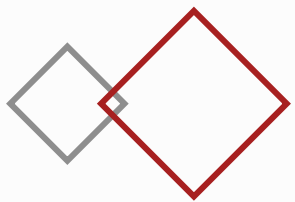
# Larger Model Calls for Larger Dataset

- ParT is one of the largest ML model for jet tagging
  - Require enough statistics to avoid overtraining!
- Existing datasets however are not large
  - Top tagging: 2M jets
  - Quark-gluon tagging: 2M jets
  - JetNet: 500k jets
  - Jedi-Net: 880k jets
  - Higgs boson tagging: 3.9M signal jets, 1.9M background jets
- A large dataset is needed
  - JetClass dataset

Table 2. Number of trainable parameters and FLOPs.

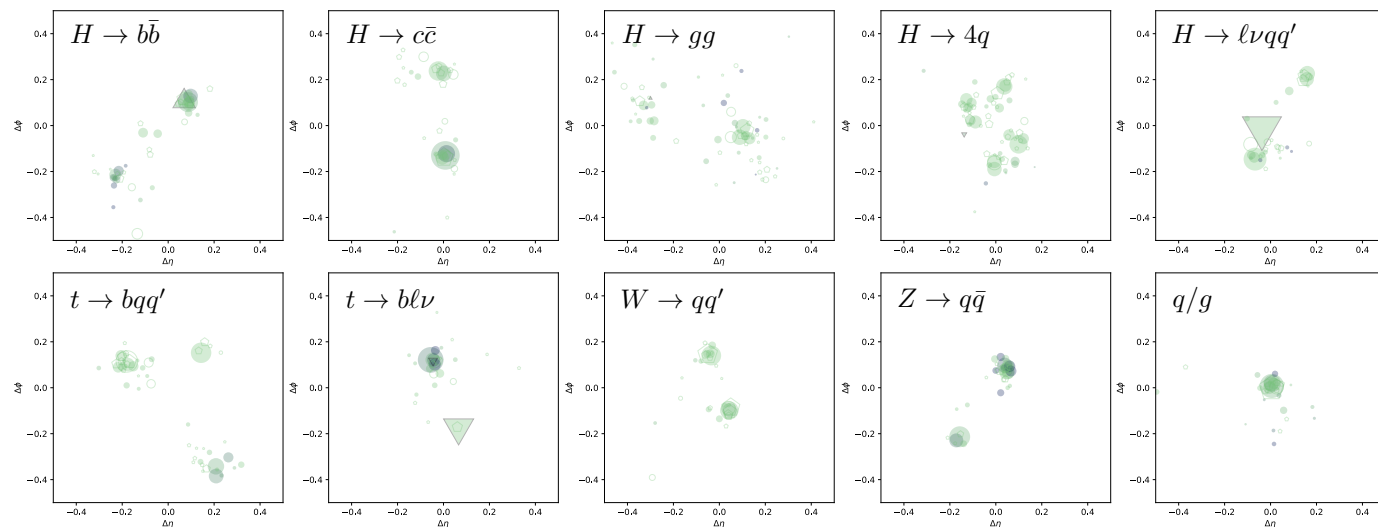
|              | Accuracy     | # params | FLOPs  |
|--------------|--------------|----------|--------|
| PFN          | 0.772        | 86.1 k   | 4.62 M |
| P-CNN        | 0.809        | 354 k    | 15.5 M |
| ParticleNet  | 0.844        | 370 k    | 540 M  |
| <b>ParT</b>  | <b>0.861</b> | 2.14 M  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



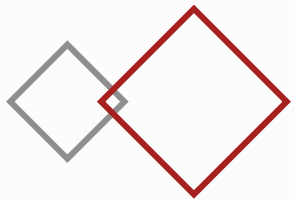


# The JetClass Dataset

- JetClass is inclusive:
  - 10 types of jets
  - Kinematics,
  - PID,
  - trajectory displacement
- JetClass is large:
  - 100M jets for training  $\rightarrow$  10M each class
  - 5M for validation
  - 20M for test  $\rightarrow$  2M each class



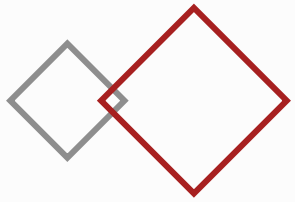




# The JetClass Dataset

- Simulation details:
  - MadGraph5\_aMC@NLO: ME level, production & decay of top, W/Z & Higgs boson
  - Pythia8: Parton showering and hadronization
  - Delphes: fast simulation of detector response, CMS configuration
  - Jets: anti- $k_T$  algorithm,  $R=0.8$  on Delphes E-Flow objects,
    - $|\eta| < 2$ ,  $500 \text{ GeV} < p_T < 1000 \text{ GeV}$
    - "high quality" jets only: jets fully containing decay products.
- Proposals of evaluation metrics for classification with JetClass:
  - Common metrics: Accuracy (Acc.) and Area Under (ROC) Curve (AUC)
  - HEP interest: Background rejection at given signal efficiency





# Numerical Experiments With JetClass

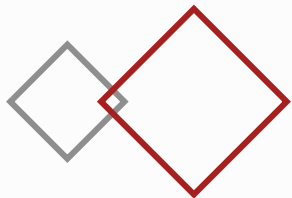


|              | All classes  |               | $H \rightarrow b\bar{b}$ | $H \rightarrow c\bar{c}$ | $H \rightarrow gg$ | $H \rightarrow 4q$ | $H \rightarrow l\nu qq'$ | $t \rightarrow bqq'$ | $t \rightarrow bl\nu$ | $W \rightarrow qq'$ | $Z \rightarrow q\bar{q}$ |
|--------------|--------------|---------------|--------------------------|--------------------------|--------------------|--------------------|--------------------------|----------------------|-----------------------|---------------------|--------------------------|
|              | Accuracy     | AUC           | Rej <sub>50%</sub>       | Rej <sub>50%</sub>       | Rej <sub>50%</sub> | Rej <sub>50%</sub> | Rej <sub>99%</sub>       | Rej <sub>50%</sub>   | Rej <sub>99.5%</sub>  | Rej <sub>50%</sub>  | Rej <sub>50%</sub>       |
| 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                      |
| <b>ParT</b>  | <b>0.861</b> | <b>0.9877</b> | <b>10638</b>             | <b>4149</b>              | <b>123</b>         | <b>1864</b>        | <b>5479</b>              | <b>32787</b>         | <b>15873</b>          | <b>543</b>          | <b>402</b>               |
| ParT (plain) | 0.849        | 0.9859        | 9569                     | 2911                     | 112                | 1185               | 3868                     | 17699                | 12987                 | 384                 | 311                      |

ParT (plain): no pair-wise feature mask applied

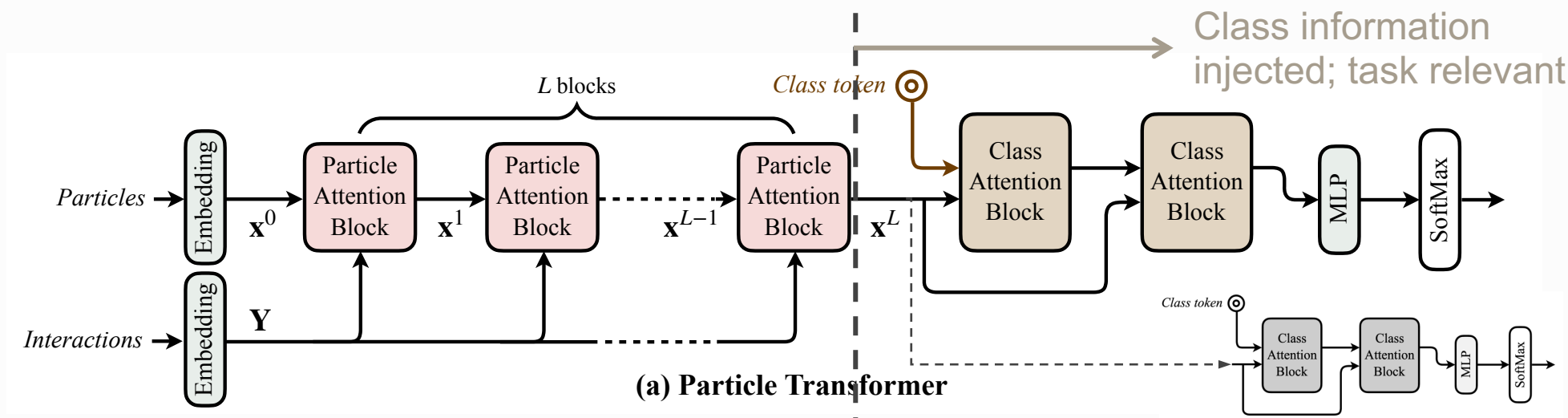
- ✓ ParT achieves SOTA performance in every classification task!
- ✓ Adding pair-wise feature enhances performance even further!
- ✓ Take  $H \rightarrow c\bar{c}b\bar{b}$  as an example: doubled background rejection  $\rightarrow \sim 1.4x$  significance!
  - ✓ Same significance reach with half data!





# Pre-training: The Lesson From JetClass

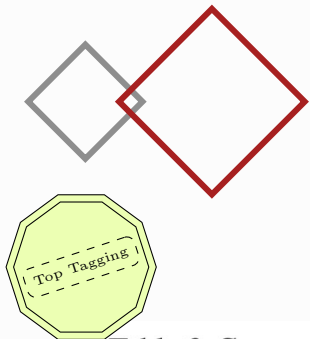
- Pre-train + Fine-tune becomes the trend in ML community
  - Self-attention from transformer → task irrelevant embedding
  - Large dataset → embedding captures generic information



No class information;  
task-irrelevant

Task-irrelevant information  
can be used for different tasks:  
Top tagging,  
quark/gluon discrimination





# Fine-tuning With different datasets

Table 3. Comparison between ParT and existing models on the top quark tagging dataset. ParT-f.t. denotes the model pre-trained on JETCLASS and fine-tuned on this dataset. ParT refers to the model trained from scratch on this dataset. Results for other models are quoted from their published results: P-CNN and ParticleNet (Qu & Gouskos, 2020), PFN (Komiske et al., 2019b), JEDI-net (Moreno et al., 2020), PCT (Mikuni & Canelli, 2021), LGN (Bogatskiy et al., 2020), and rPCN (Shimmin, 2021).

|                         | Accuracy     | AUC           | Rej <sub>50%</sub> | Rej <sub>30%</sub> |
|-------------------------|--------------|---------------|--------------------|--------------------|
| P-CNN                   | 0.930        | 0.9803        | 201 ± 4            | 759 ± 24           |
| PFN                     | —            | 0.9819        | 247 ± 3            | 888 ± 17           |
| ParticleNet             | 0.940        | 0.9858        | 397 ± 7            | 1615 ± 93          |
| JEDI-net (w/ $\sum O$ ) | 0.930        | 0.9807        | —                  | 774.6              |
| PCT                     | 0.940        | 0.9855        | 392 ± 7            | 1533 ± 101         |
| LGN                     | 0.929        | 0.964         | —                  | 435 ± 95           |
| rPCN                    | —            | 0.9845        | 364 ± 9            | 1642 ± 93          |
| ParT                    | 0.940        | 0.9858        | 413 ± 16           | 1602 ± 81          |
| <b>ParT-f.t.</b>        | <b>0.944</b> | <b>0.9877</b> | <b>691 ± 15</b>    | <b>2766 ± 130</b>  |

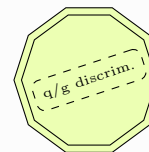


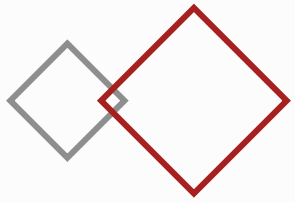
Table 4. Comparison between ParT and existing models on the quark-gluon tagging dataset. ParT-f.t. denotes the model pre-trained on JETCLASS and fine-tuned on this dataset. ParT refers to the model trained from scratch on this dataset. Results for other models are quoted from their published results: P-CNN and ParticleNet (Qu & Gouskos, 2020), PFN (Komiske et al., 2019b), ABCNet (Mikuni & Canelli, 2020), PCT (Mikuni & Canelli, 2021), and rPCN (Shimmin, 2021). The subscript “exp” and “full” distinguish models using partial or full particle identification information.

|                                 | Accuracy     | AUC           | Rej <sub>50%</sub> | Rej <sub>30%</sub> |
|---------------------------------|--------------|---------------|--------------------|--------------------|
| P-CNN <sub>exp</sub>            | 0.827        | 0.9002        | 34.7               | 91.0               |
| PFN <sub>exp</sub>              | —            | 0.9005        | 34.7 ± 0.4         | —                  |
| ParticleNet <sub>exp</sub>      | 0.840        | 0.9116        | 39.8 ± 0.2         | 98.6 ± 1.3         |
| rPCN <sub>exp</sub>             | —            | 0.9081        | 38.6 ± 0.5         | —                  |
| ParT <sub>exp</sub>             | 0.840        | 0.9121        | 41.3 ± 0.3         | 101.2 ± 1.1        |
| <b>ParT-f.t.<sub>exp</sub></b>  | <b>0.843</b> | <b>0.9151</b> | <b>42.4 ± 0.2</b>  | <b>107.9 ± 0.5</b> |
| PFN <sub>full</sub>             | —            | 0.9052        | 37.4 ± 0.7         | —                  |
| ABCNet <sub>full</sub>          | 0.840        | 0.9126        | 42.6 ± 0.4         | 118.4 ± 1.5        |
| PCT <sub>full</sub>             | 0.841        | 0.9140        | 43.2 ± 0.7         | 118.0 ± 2.2        |
| ParT <sub>full</sub>            | 0.849        | 0.9203        | 47.9 ± 0.5         | 129.5 ± 0.9        |
| <b>ParT-f.t.<sub>full</sub></b> | <b>0.852</b> | <b>0.9230</b> | <b>50.6 ± 0.2</b>  | <b>138.7 ± 1.3</b> |

Take home:

Pre-train with JetClass helps ParT to reach SOTA performance!





# Welcome to Jet-Universe!



- We are more than glad to share our work to the whole community:
  - Both source code of ParT and JetClass dataset are now public at Jet Universe
  - Hope to see more enthusiasts onboard!

jet-universe

Overview Repositories 1 Projects Packages Teams People 3

Popular repositories

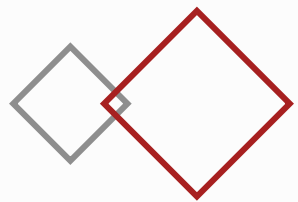
**particle\_transformer** Public  
Official implementation of "Particle Transformer for Jet Tagging".  
Python ☆ 8 🍴 3

Repositories

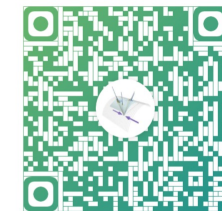
Find a repository... Type Language Sort New

**particle\_transformer** Public  
Official implementation of "Particle Transformer for Jet Tagging".  
Python ☆ 8 MIT 🍴 3 🗨 0 📄 0 Updated 12 days ago



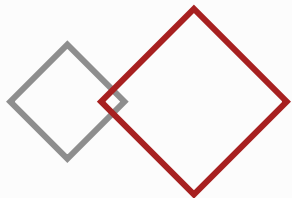


## Summary



- Particle Transformer (ParT):
  - Dedicated transformer architecture for jet tagging
  - SOTA in various benchmarks
  - Particle interactions help ParT to perform better
- JetClass dataset:
  - Large and inclusive: order of magnitude higher in statistics and classes of jets
  - Pretrain with JetClass enhance ParT's performance on other datasets
- ParT and JetClass are publicly available on Github!
  - Welcome to [Jet Universe](#)



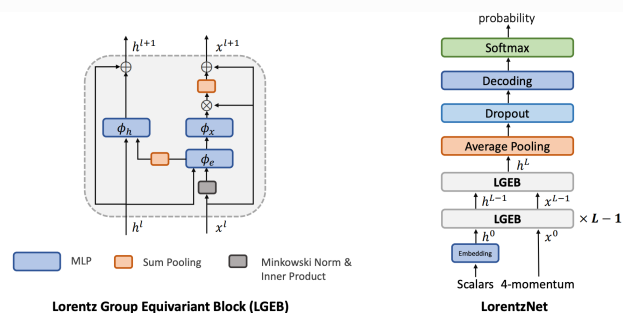


# End of the Journey?

- There is still room for further improvement:
  - “Physics” should be a term added to the “SOTA equation”
    - Pair-wise features from particle interactions enhance ParT’s performance
    - Improvements has been seen in physics inspired models:
      - LundNet, LorentzNet (Lorentz symmetry applied, competitive performance),...

## An Efficient Lorentz Equivariant Graph Neural Network for Jet Tagging

Shiqi Gong<sup>a,c,1</sup> Qi Meng<sup>b</sup> Jue Zhang<sup>b</sup> Huilin Qu<sup>c</sup> Congqiao Li<sup>d</sup> Sitian Qian<sup>d</sup> Weitao Du<sup>a</sup> Zhi-Ming Ma<sup>a</sup> Tie-Yan Liu<sup>b</sup>



JHEP07(2022)030

| Model       | Accuracy     | AUC           | $1/\epsilon_B$<br>( $\epsilon_S = 0.5$ ) | $1/\epsilon_B$<br>( $\epsilon_S = 0.3$ ) |
|-------------|--------------|---------------|--|--|
| ResNeXt     | 0.936        | 0.9837        | 302 ± 5                                  | 1147 ± 58                                |
| P-CNN       | 0.930        | 0.9803        | 201 ± 4                                  | 759 ± 24                                 |
| PFN         | 0.932        | 0.9819        | 247 ± 3                                  | 888 ± 17                                 |
| ParticleNet | 0.940        | 0.9858        | 397 ± 7                                  | 1615 ± 93                                |
| EGNN        | 0.922        | 0.9760        | 148 ± 8                                  | 540 ± 49                                 |
| LGN         | 0.929        | 0.9640        | 124 ± 20                                 | 435 ± 95                                 |
| LorentzNet  | <b>0.942</b> | <b>0.9868</b> | <b>498 ± 18</b>                          | <b>2195 ± 173</b>                        |
| ParT        | 0.940        | 0.9858        | 413 ± 16                                 | 1602 ± 81                                |
| ParT-ft.    | <b>0.944</b> | <b>0.9877</b> | <b>691 ± 15</b>                          | <b>2766 ± 130</b>                        |

| Model                   | Accuracy     | AUC           | $1/\epsilon_B$<br>( $\epsilon_S = 0.5$ ) | $1/\epsilon_B$<br>( $\epsilon_S = 0.3$ ) |
|-------------------------|--------------|---------------|--|--|
| ResNeXt                 | 0.821        | 0.8960        | 30.9                                     | 80.8                                     |
| P-CNN                   | 0.827        | 0.9002        | 34.7                                     | 91.0                                     |
| PFN                     | –            | 0.9005        | 34.7 ± 0.4                               | –  |
| ParticleNet             | 0.840        | 0.9116        | 39.8 ± 0.2                               | 98.6 ± 1.3                               |
| EGNN                    | 0.803        | 0.8806        | 26.3 ± 0.3                               | 76.6 ± 0.5                               |
| LGN                     | 0.803        | 0.8324        | 16.0                                     | 44.3                                     |
| LorentzNet              | <b>0.844</b> | <b>0.9156</b> | <b>42.4 ± 0.4</b>                        | <b>110.2 ± 1.3</b>                       |
| ParT <sub>exp</sub>     | 0.840        | 0.9121        | 41.3 ± 0.3                               | 101.2 ± 1.1                              |
| ParT-ft. <sub>exp</sub> | <b>0.843</b> | <b>0.9151</b> | <b>42.4 ± 0.2</b>                        | <b>107.9 ± 0.5</b>                       |

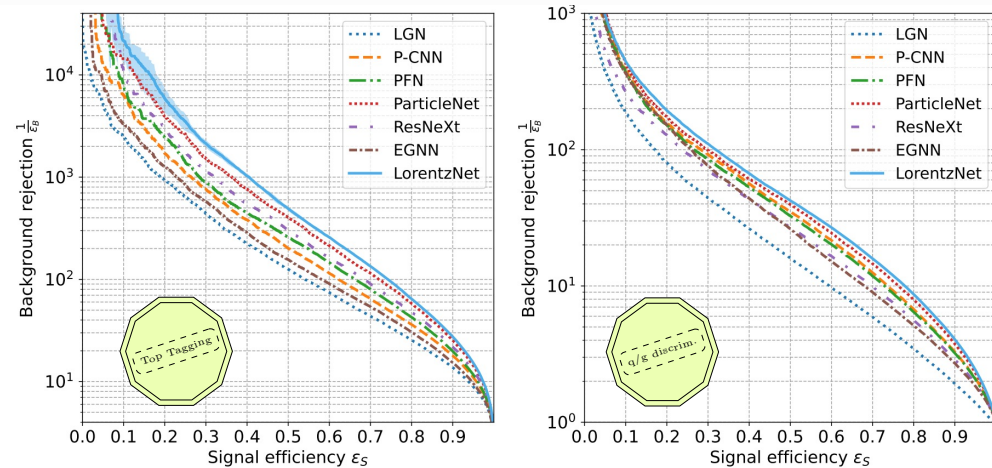
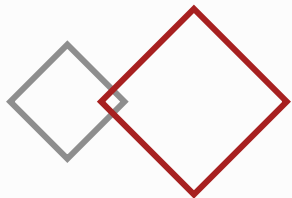


Figure 2. A comparison of ROC curves between LorentzNet and other algorithms on top tagging dataset (left) and quark-gluon dataset (right).





# End of the Journey?

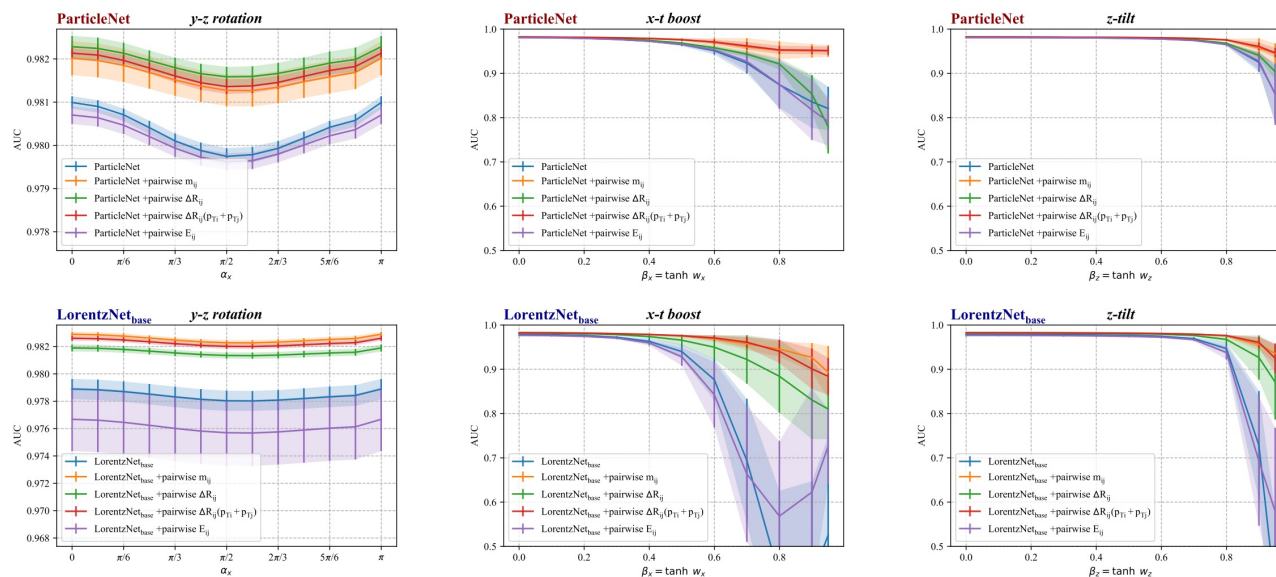
- There is still room for further improvement:
  - “Physics” should be a term added to the “SOTA equation”
    - Pair-wise features from particle interactions enhance ParT’s performance
    - Improvements has been seen in physics inspired models:
      - LundNet, LorentzNet (Lorentz symmetry applied, competitive performance), ...
      - Lorentz symmetry is indeed important!

Does Lorentz-symmetric design boost network performance in jet physics?

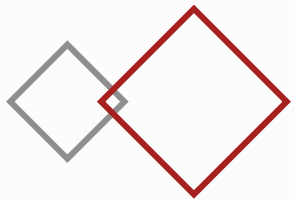
Congqiao Li,<sup>1,\*</sup> Huilin Qu,<sup>2</sup> Sitian Qian,<sup>1</sup> Qi Meng,<sup>3</sup> Shiqi Gong,<sup>4</sup> Jue Zhang,<sup>3</sup> Tie-Yan Liu,<sup>3</sup> and Qiang Li<sup>1</sup>

Our answer: **YES!**

Adding Lorentz symmetry inspired features (either pairwise or elementwise) will enhance performance of NN



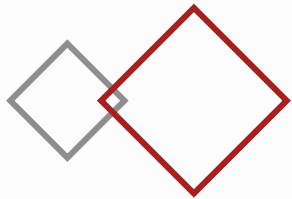




## End of the Journey?

- There is still room for further improvement:
  - “Physics” should be a term added to the “SOTA equation”
    - Pair-wise features from particle interactions enhance ParT’s performance
    - Improvements has been seen in physics inspired models:
      - LundNet, LorentzNet (Lorentz symmetry applied, competitive performance),...
      - Lorentz symmetry is indeed important!
  - Education from ML community:
    - We have already benefited a lot from ML community (transformer, pair-wise features, pretrain, etc.)
    - Attempts on novel techniques (Multi-modal transformers, neighbor embedding...) are promising to explore!
- A practical perspective: model complexity vs efficient computation
  - Larger and Larger ML models call for model compression techniques.





# Back Up Input Features

Table 5. Particle input features used for jet tagging on the JETCLASS, the top quark tagging (TOP) and the quark gluon tagging (QG) datasets. For QG, we consider two scenarios: QG<sub>exp</sub> is restricted to use only the 5-class experimentally realistic particle identification information, while QG<sub>full</sub> uses the full set of particle identification information in the dataset and further distinguish between different types of charged hadrons and neutral hadrons.

| Category                | Variable                             | Definition  | JETCLASS | TOP | QG <sub>exp</sub> | QG <sub>full</sub> |
|-------------------------|--------------------------------------|---|----------|-----|-------------------|--------------------|
| Kinematics              | $\Delta\eta$                         | difference in pseudorapidity $\eta$ between the particle and the jet axis                                 | ✓        | ✓   | ✓                 | ✓                  |
|                         | $\Delta\phi$                         | difference in azimuthal angle $\phi$ between the particle and the jet axis                                | ✓        | ✓   | ✓                 | ✓                  |
|                         | $\log p_T$                           | logarithm of the particle's transverse momentum $p_T$   | ✓        | ✓   | ✓                 | ✓                  |
|                         | $\log E$                             | logarithm of the particle's energy  | ✓        | ✓   | ✓                 | ✓                  |
|                         | $\log \frac{p_T}{p_{T(\text{jet})}}$ | logarithm of the particle's $p_T$ relative to the jet $p_T$   | ✓        | ✓   | ✓                 | ✓                  |
|                         | $\log \frac{E}{E(\text{jet})}$       | logarithm of the particle's energy relative to the jet energy   | ✓        | ✓   | ✓                 | ✓                  |
|                         | $\Delta R$                           | angular separation between the particle and the jet axis ( $\sqrt{(\Delta\eta)^2 + (\Delta\phi)^2}$ )     | ✓        | ✓   | ✓                 | ✓                  |
| Particle identification | charge                               | electric charge of the particle   | ✓        | —   | ✓                 | ✓                  |
|                         | Electron                             | if the particle is an electron ( <code> pid ==11</code> )   | ✓        | —   | ✓                 | ✓                  |
|                         | Muon                                 | if the particle is an muon ( <code> pid ==13</code> )   | ✓        | —   | ✓                 | ✓                  |
|                         | Photon                               | if the particle is an photon ( <code>pid==22</code> )   | ✓        | —   | ✓                 | ✓                  |
|                         | CH                                   | if the particle is an charged hadron ( <code> pid ==211</code> or <code>321</code> or <code>2212</code> ) | ✓        | —   | ✓                 | ✓ <sup>a</sup>     |
|                         | NH                                   | if the particle is an neutral hadron ( <code> pid ==130</code> or <code>2112</code> or <code>0</code> )   | ✓        | —   | ✓                 | ✓ <sup>b</sup>     |
| Trajectory displacement | $\tanh d_0$                          | hyperbolic tangent of the transverse impact parameter value   | ✓        | —   | —                 | —                  |
|                         | $\tanh d_z$                          | hyperbolic tangent of the longitudinal impact parameter value   | ✓        | —   | —                 | —                  |
|                         | $\sigma_{d_0}$                       | error of the measured transverse impact parameter   | ✓        | —   | —                 | —                  |
|                         | $\sigma_{d_z}$                       | error of the measured longitudinal impact parameter   | ✓        | —   | —                 | —                  |

<sup>a</sup> (`|pid|==211`) + (`|pid|==321`)\*0.5 + (`|pid|==2212`)\*0.2

<sup>b</sup> (`|pid|==130`) + (`|pid|==2112`)\*0.2.

