# TRANSFORMER FOR PARTICLE PHYSICS **MIHOKO NOJIRI** (KEK)

with Ahmed Hammad and Sung Hak Lim

### MLPHYSICS GRANT IN JAPAN"MACHINE LEARNING PHYSICS "

MLPhys Foundation of "Machine Learning Physics" Grant-in-Aid for Transformative Research Areas (A) CONTACT

Members only





#### message

Head Investigator

#### Koji Hashimoto

Professor Particle Physics Theory Group Department of physics, Kyoto University



The research area "Machine Learning Physics" will begin with the aim of discovering new laws and pioneering new materials

Hello. My name is Koji Hashimoto, Professor of Graduate School of Science, Kyoto University. Let me explain about the "Learning Physics Domain" that we are just now trying to create. This new transformative research area aims to revolutionize fundamental physics by combining machine learning and physics.

- B01 Math and application of DL
- B02 Statistical data and ML
- B03 Topology and Geometry of ML
- A01 Lattice

A02 Mihoko Nojiri HEP Junichi Tanaka (ICEPP Tokyo, ATLAS) Masako lawasaki (Osaka Metropolitan Belle II ) Noriko Takemura and Hajime Nagahara (Data Science)
A03 Condensed Matter
A04 Quantum and Gravity

> MLPhysics PD. Ahmed Hammad 2017-2020: Ph.D Basel University, Basel Switzerland 2020-2023: SeoulTech, Korea 2023- KEK



# Deep Learning is changing particle physics Flavor Tagging Performance

#### b/c-tagging performance

- Promising performance compared to previous taggers
  - ×3 better light jet rejection (at b-jet eff 70%) than DeepJet
  - ×2 better light rejection + ×2 better b-jet rejection (at c-jet eff 35 %) than DeepJet





#### epJet t c-iet eff 35 %) than DeepJe



# **Basic ingredient of ML**

$$x_i^{(K+1)} = \varphi(w_{ij}^{(K)}x_j^{(K)} + b_i^{(K)})$$
  
\varphi: activation function —-for r



#### non lineality



Rectified Linear Unit(ReLU)

 $\varphi_{\mathrm{ReLU}}(x) = x\theta(x)$ 

# $\varphi(\cdot) \rightarrow \varphi(w_{ij}x_j + b_i)$

#### activation output

# **Training for classification**

#### data pool of images

2222222222222222222





(28x28)の画像データをn個

✓表現力 expressive power ✓データを学習 learn from data ✓微分可能Simple linear algebra + activation minimization of loss function

 $L(y, \hat{y})$ 

output  $\hat{y} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_{10}), \sum \hat{y}_i = 1$  $\exp(x_i)$ Ŷ;  $\sum_{i} \exp(x_{i})$ 

 $\hat{y}$ 's represent likeliness to be y







## TRANSFORMER

- Successful language model— good at long correlation(context)
- good at long term memory (keeping information)
  - Scaling low of LLM (Large model and large data →improve forever

**Recent Apple paper** (They just know patterns very well)

https://chatgpt.com/share/6854a7be-2768-800c-a9a4-34dfd0431c4a

Apple LLM 推論限界



ChatGPT 40 ~

大規模言語モデルの限界について最近の Apple の論文が話題だと聞いたのですが、 どの論文のことでしょうか。

Appleの研究者による注目の論文は、大規模言語モデル(LLMs)の 「推論能力」における根源的な限界を示唆するものです。特に、 「Large Reasoning Models (LRMs)」が高度に複雑な問題に対し て"完全な精度崩壊(complete accuracy collapse)"を起こすことを

報告しています。具体的には以下の内容が取り上げられています。









## "PARTICLE TRANSFORMER" :SELF ATTENTION LAYERS

$$\operatorname{Attention}(Q, K, V) = \operatorname{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

#### SETUP

- input X: n( particles in the jet) x d(features)
- <u>n x n Attention matrix</u>  $A = QK^T$  $K = W_K X, Q = W_Q X$  (W trainable parameters)
- Multiply  $V = W_V X$  to get n x d output  $\delta X$
- stack attention layers for  $X \rightarrow X'=X + \delta X \rightarrow X''$ skip connection help to keep original information
- · completely built from Low level feature
- e)  $\cdot$  using all correlation inclusing long range to short range.



### Features in the context of jet classification

	Category	Variable
Particle momentum	Kinematics	$\begin{array}{l} \Delta\eta\\ \Delta\phi\\ \log p_{\mathrm{T}}\\ \log p_{\mathrm{T}}\\ \log E\\ \log \frac{p_{\mathrm{T}}}{p_{\mathrm{T}}(\mathrm{jet})}\\ \log \frac{E}{E(\mathrm{jet})}\\ \Delta R \end{array}$
charge,particle ID	Particle identification	charge Electron Muon Photon CH NH
displaced vertex	Trajectory displacement	$ anh d_0 \  anh d_z \ \sigma_{d_0} \ \sigma_{d_z}$

#### Definition

difference in pseudorapidity  $\eta$  between the particle and the jet axis difference in azimuthal angle  $\phi$  between the particle and the jet axis logarithm of the particle's transverse momentum  $p_{\rm T}$ 

logarithm of the particle's energy

logarithm of the particle's  $p_{\rm T}$  relative to the jet  $p_{\rm T}$ 

logarithm of the particle's energy relative to the jet energy

angular separation between the particle and the jet axis  $(\sqrt{(\Delta \eta)^2 + (\Delta \phi)^2})$ 

electric charge of the particle
if the particle is an electron (|pid|==11)
if the particle is an muon (|pid|==13)
if the particle is an photon (pid==22)
if the particle is an charged hadron (|pid|==211 or 321 or 2212)
if the particle is an neutral hadron (|pid|==130 or 2112 or 0)

hyperbolic tangent of the transverse impact parameter value hyperbolic tangent of the longitudinal impact parameter value error of the measured transverse impact parameter error of the measured longitudinal impact parameter



## PHYSICS APPLICATION OF MACHINE LEARNING: CONCERNS

1.bias and valiance trade off on the training results

ML can approximate "any function" in infinite paramter limit, therefore unstable against the fluctuation of the training sample.

#### 2.Interpretability

It gives you the results, but the reasoning is deep inside the network and hard to extract. This is problematic when you do not understand "true" distirbution. This is often the case for collider physics, where MC do not reproduce experimental data in detail.

#### 3. scaling vs symmetry

Large data and Large network wins(it says), but actually nature has symmetry, Permutation Invariance, boost invariance, Lorentz symmetry.. which is not easy to reconstruct from scratch.



## THE PERFORMANCE FOR TOP VS QCD CLASSIFICATION

	Accuracy	AUC	$1/\epsilon_B(\epsilon_s=0.5)$	$1/\epsilon_B(\epsilon_s=0.3)$	Parameters
	Lorentz	z invaria	nce based net	works	
PELICAN[35]	0.9426	0.987		$2250\pm75$	208K
LorentzNet[70]	0.942	0.9868	$498 \pm 18$	$2195 \pm 173$	224K
L-GATr[71]	0.942	0.9870	$540 \pm 20$	$2240\pm70$	
	At	tention	based network	S	
ParT[49]	0.940	0.9858	$413 \pm 6$	$1602\pm81$	$2.14\mathrm{M}$
MIParT[50]	0.942	0.9868	$505\pm8$	$2010\pm97$	$720.9 \mathrm{K}$
Mixer[21]	0.940	0.9859	$416 \pm 5$		$86.03\mathrm{K}$
OmniLearn[72]	0.942	0.9872	$568 \pm 9$	$2647 \pm 192$	$1.6\mathrm{M}$
Plain Transformer*	0.927	0.979	$362\pm7$	$780\pm73$	$1.7\mathrm{M}$
<b>IAFormer</b> *	0.942	0.987	$510\pm6$	$2012\pm30$	$\mathbf{211K}$

#### Yellow bands highlight our works!



## TRANSFORMER VS PHYSICS

#### With Ahmed Hammad

1. fast, lightweight, while keeping performance Reduction of network parameter 2M→100k

scaling behavior MIXER network arXiv 2404 14677 JHEP 06 (2024) 176

- 2. Jet analysis  $\rightarrow$  event analysis. analysis of  $H \rightarrow hh$  arXiv 2401.00452 JHEP 03(2024)
- **3. Respect symmetry** Replacing "attention from generic features" → "pairwise boost invariant information " (IAFormer)

**Reduce valiance in training**  $\rightarrow$  via differential attention arXiv 2505.03258 With Sung Hak Lim

- 4. Identify the key parameters for classifications → via comparison with the simple MLP model using High leavel features
- 5. calibration of MC in future. arXiv 2312.11760 JHEP 07 (2024) 146 arXiv 2503.01452 JHEP(accepted)

Incorporate physics picture, Cross attention between subjet vs hadron inspired by QCD



### 1. PARTON SHOWER AND NETWORK STRUCTURE

"Ahmed Hammad, MN "Streamlined jet tagging network assisted by jet prong structure" arXiv 2404 14677 JHEP 06 (2024) 176

- Hard Process = Partons(quarks and gluons) {y}
- a jet: P(hadrons in jets | parton ~ jet) =  $P({x_i} | {y})$
- · jet with substructure
- · Maybe several fatjets in an event (factorization)

Why don't we construct the network forcusing on parton(subjet) vs hadrons

 $P(\{x_i\} | \{y_{\alpha}\})$ 

 $P(\{x_i\}, \{x_i'\}, \{y_{\alpha}\}, \{y_{\beta}\}) \sim P(\{x_i\} | \{y_{\alpha}\}) P(\{x_i'\} | \{y_{\beta}\}) P(\{y_{\alpha}, y_{\beta}'\})$ 







## CROSS ATTENTION LAYERS

• restrict network to cross attention (subjet) x (jet constituent)  $A = QK^T$ 

· jet constituent Q

subjet ~ parton, shower K V

Structure: High scale feature (subjet) gives extra weight to jets constituents

The performace is not reduced significantly from Transformer while networks size is very small becuase we respect QCD

subjet kinematics



## 2. CROSS ATTENTION AND GLOBAL EVENT KINEMATICS

## $X \rightarrow HH$



H(bb)H(bb) most sensitive channel for  $m_X > 400/500 \text{ GeV}$ H(yy)H(bb) complement in the low mass



**SLAC** Caterina Vernieri · LCWS 2024 · Tokyo

#### Phys. Rev. Lett. 132 (2024) 231801



17



## cross attention motention for 2 fatjet events

jet constituent information relevant gives extra weight to the corresponding jets though backward propagation

We can replace transformer to "mixer+subjet" network

•



### IMPROVEMENT USING CROSS ATTENTION COMBINIE TO NEXT



## NETWORK BUDGET BEHIND CROSS ATTENTION

Isn't that good enough to add subjet information to particle and do transformer?



1. optimization of main term

2. Lost in minor term optimization

our network kill this term and keep off-diagonal part only





### 3. (IAFormer =InterAction transFormer)

Transformer Input:  $\eta, \phi, p_T$ ...

BUT we only want to use lorentz covariant or boost invariant for LHC) information

$$\Delta R = \sqrt{(\eta - \eta')^2 + (\phi - \phi')^2} \quad z = mi$$

so using  $K = W_K V, Q = W_O V$  is too much...

Lorentz covariant network is proposed and performance is high. but boost invariance is also known to be useful at LHC

 $in(z_1, z_2)/(z_1 + z_2) \dots$ 

 $\alpha = \operatorname{softmax}\left([Q \times K^T]^{n \times n}\right) \to \alpha = \operatorname{softmax}(\mathscr{I}^{n \times n})$ 

**I**: function of pairwise variable



## STABILITY OF THE TRAINING-SPARSE ATTENTION

- Sparse attention: a rule to use only part of attention matrix for quick convergence and reasonings
- static attention—use "fixed patterns" to filter attention
- This is probably very important for Language model but does not looks right for particle physics.
  - **band attention** "大きなりんご big apple" "赤い りんご red apple"
  - Dilated attention 赤い<u>りんご</u>が <u>落ちた</u>のを <u>みた</u> (I saw a red apple falling )





https://developers.agirobots.com/jp/sparse-attention/

### **DIFFERENTIAL ATTENTION** AN EXAMPLE OF DYNAMIC ATTENTION

- Dynamic attention "ask network to learn viable sparse attention pattern" •
- arXiv:2410.05258

### • $\alpha = \operatorname{softmax}(\mathscr{I}) \to \alpha^{(i)} = \operatorname{softmax}(\mathscr{I}_1^{(i)}) - \beta^{(i)} \operatorname{softmax}(\mathscr{I}_2^{(i)})$

 $\mathcal{J}^{(1)}, \beta^{(1)}$ 

Each layer built different filters dynamically

• We introduce a new dynamic attention called "differential attention" to the network.





### **EFFICIENCY AND INTERPRETATION** Compare hidden layers using CKA(centered kernel alignment) similaity



 $X(d \times h_1)$  $\rightarrow M = XX^{\dagger} (d \times d)$ 

 $Y(d \times h_2) \to N = YY^{\dagger} \quad (d \times d)$ 

$$CKA(M,N) = \frac{HSIC(M,N)}{\sqrt{HSIC(M,M)HSIC(N,N)}}$$

$$HSIC(M,N) = \frac{1}{(d-1)^2} Tr(MHNH) \qquad H = \delta_{ij} - \frac{1}{d}$$

#### CKA <<1 if new information is captured



Figure 3. CKA reveals when depth becomes pathological. Top: Linear CKA between layers of individual networks of different depths on the CIFAR-10 test set. Titles show accuracy of each network. Later layers of the 8x depth network are similar to the last layer. Bottom: Accuracy of a logistic regression classifier trained on layers of the same networks is consistent with CKA.







### IAFormar constantly improve classification performance

IAFormer( $\epsilon_b(50\%) = 510$ )

	12	0.6	0.5	0.7	0.8	0.8	0.8	0.8	0.8	0.9	0.9	0.9	1.0
	<u>-1</u>	0.6	0.5	0.7	0.8	0.8	0.8	0.9	0.9	0.9	0.9	1.0-	0.9
	) 1	0.0	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.0	1.0	1.0	0.0
	1(	0.7	0.5	0.7	0.8	0.8	0.9	0.9	0.9	0.9	1.0	0.9	0.9
	6 -	0.7	0.5	0.7	0.8	0.9	0.9	0.9	0.9	1.0	0.9	0.9	0.9
yer	~ ~	0.7	0.5	0.7	0.8	0.9	0.9	0.9	1.0	0.9	0.9	0.9	0.8
on la	<b>⊳</b> -	0.7	0.6	0.8	0.9	0.9	0.9	1.0	0.9	0.9	0.9	0.9	0.8
entic	9 -	0.7	0.6	0.8	0.9	0.9	1.0	0.9	0.9	0.9	0.9	0.8	0.8
Att	ۍ ۲0 -	0.6	0.6	0.8	0.9	1.0	0.9	0.9	0.9	0.9	0.8	0.8	0.8
	4-	0.6	0.6	0.8	1.0	0.9	0.9	0.9	0.8	0.8	0.8	0.8	0.8
	က -	0.6	0.8	1.0	0.8	0.8	0.8	0.8	0.7	0.7	0.7	0.7	0.7
	- 7	0.6	1.0	0.8	0.6	0.6	0.6	0.6	0.5	0.5	0.5	0.5	0.5
		1.0	0.6	0.6	0.6	0.6	0.7	0.7	0.7	0.7	0.7	0.6	0.6
		1	2	3	4	5	6	7	8	9	10	11	12
					1	Atte	entic	on l	aye	ſ			

CKA - IAFormer

1.0

-0.9

-0.8

-0.7

-0.6

-0.5

### ParT $\epsilon_b(50\%) = 413$

### **Plain Transformer** $\epsilon_b(50\%) = 360$

			Cł	KA	- T	ran	sfoi	rme	er+	$I_{i,j}$							CK	A -	- Pl	ain	Tr	ans	form	ner		
12	0.5	0.5	0.6	0.7	0.9	0.9	1.0	1.0	1.0	1.0	1.0	1.0		12	0.9	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.(
11	0.5	0.6	0.6	0.8	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	11	0.9	0.9	1.0	1.0			<b>V</b> .0	<b>P</b> 1.0	<b>9</b> .0	<b>D</b> <sub>1.0</sub>	1.0	1.(
10	0.5	0.6	0.6	0.8	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-0.0	10	0.9	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
o -	0.5	0.6	0.6	0.8	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.9	o -	0.9	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
yer 8	0.5	0.6	0.7	0.8	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	-0.8	yer 8	0.9	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.(
n la	0.6	0.6	0.7	0.8	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.0	n la	0.9	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.(
entic	0.6	0.7	0.7	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.9	-07	entic 6	0.9	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.(
Att.	0.7	0.8	0.8	1.0	1.0	1.0	0.9	0.9	0.9	0.9	0.9	0.9	0.1	Att.	0.9	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.(
4.	0.8	0.9	1.0	1.0	1.0	0.9	0.8	0.8	0.8	0.8	0.8	0.7	-0.6	4 -	0.9	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.(
က -	1.0	1.0	1.0	1.0	0.8	0.7	0.7	0.7	0.6	0.6	0.6	0.6	0.0	က -	0.9	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.(
0	1.0	1.0	1.0	0.9	0.8	0.7	0.6	0.6	0.6	0.6	0.6	0.5	-0.5	- 17	1.0	1.0	1.0	0.9	0.9	0.9	0.9	0.9	0.9	0.9	0.9	0.9
	1.0	1.0	1.0	0.8	0.7	0.6	0.6	0.5	0.5	0.5	0.5	0.5	0.0		1.0	1.0	0.9	0.9	0.9	0.9	0.9	0.9	0.9	0.9	0.9	0.9
	1	2	3	4	5 Atte	6 enti	7 on l	8 aye	9 r	10	11	12			1	2	3	4	$\dot{5}$	6 entio	7 on l	8 ayei	9 r	10	11	12





#### **4. Important feature for jet classification** Amon Furuichi, Sung Hak Lim, Mihoko M. Nojiri JHEP 07 (2024) 146 2312.11760

#### pt distribution of constituents

Jet spectrum two point Energy correlation (unlocalized sampling )

 $S_{2,ab}(R) \stackrel{\text{def}}{=} \sum_{i \in a} \sum_{j \in b} p_{T,i} p_{T,j} \delta(R - R_{ij}).$ 

Energy flow polinomials

 $i_1 \in J$   $i_N \in J$ 

 $\mathrm{EFP}_{G}^{(\kappa,\beta)} = \sum \cdots \sum z_{i_{1}}^{(\kappa)} \cdots z_{i_{N}}^{(\kappa)} \prod \theta_{i_{m}i_{\ell}}^{(\beta)}.$ 

 $(m,\ell) \in G$ 



Subjet Localized sampling momentum and counting for various angular sccale R=0.1, 0.2, 0.3

Minkowski Functionals geometry of jet cosntituent distribution



## NETWORK USING HL INPUTS VS PARTICLE TRANSFORMER





arXiv 2312.11760 JHEP 07 (2024) 146 arXiv 2503.01452 JHEP(accepted)

	Pythi	a	Herv	vig	PY vs QCI	HW D	PY vs TO
AM	85.7	343	61.3	<b>244</b>	2.78	5.43	2.82
[ <mark>] [33</mark> ]	90.5	372	62.6	242	2.77	5.38	3.07

for Systematical error

correction

The network using HL feature is more stable and useful to fine difference among Monte Carlo and reduce systematical errors (currenty more than 20%) because variance is low.







#### **QCD** correction

Matching

SOMETHING

Parton shower

Hadronization

#### 5. Calibration of MC data

- QCD connecting BSM and events
- Madgraph: Automatic Amplitude calculation in NLO level
- 2001 Matrix element and Parton shower matching MLM CKKW → 2007 Madgraph
- Pythia, Herwig, Sherpa evolving toward Dipole shower, Higher order collection of PS (Panscale)
- 2006 QCD aware definition of jets(fastjet)

## Should we give up Theory? No.



## It is time to invest the parton shower and hadronization algorithm in the level relevant to DeepLearning Era

#### Comparison to LEP data



Colour is handled using the NODS scheme which gives full colour accuracy at NLL for global observables (includes those shown)

\*This should be taken as an average  $\alpha_s^{\mathrm{eff}}$  not an  $\alpha_s^{\overline{MS}}$ 

J.Helliwell (U.O.O)

NNLL Parton Showers

Inclusion of NNLL potentially resolves the issue of needing an anomalously large value of  $\alpha_s(m_Z)$  to achieve good agreement with LEP data.  $(\alpha_s(m_Z) = 0.137 \text{ in Pythia's}$ Monash 13 tune \* arxiv:1404.5630, Skands, Carrazza,

Some caution needed as no 3-jet NLO matching, which is known to be relevant away from the 2-jet region.

 A comprehensive study of shower uncertainties is still to be done.

# PanScale shower reproducing $\alpha_s(m_Z)$ at last!

https://gsalam.web.cern.ch/ panscales/



# TAKE AWAY MESSAGE

- 1. fast, lightweight, while keeping performance
- 2. Incorporate physics picture
- 3. Jet analysis  $\rightarrow$  event analysis.(H $\rightarrow$ hh)
- 4. Respect symmetry Replacing "attention from generic features" → "pairwise boost invariant information " (IAFormer)
- 5. Reduce valiance in training
- 6. Identify the key parameters for classifications

**RESPCETING QCD** 

**Cross attention is important** 

Symmetry

Improved stability within DL

Identify Important variables in DL era Improving MC simulation

