

t-SNE (7.3e+02 sec)

Machine Learning in HEP data processing

李 刚

中国科学院高能物理研究所

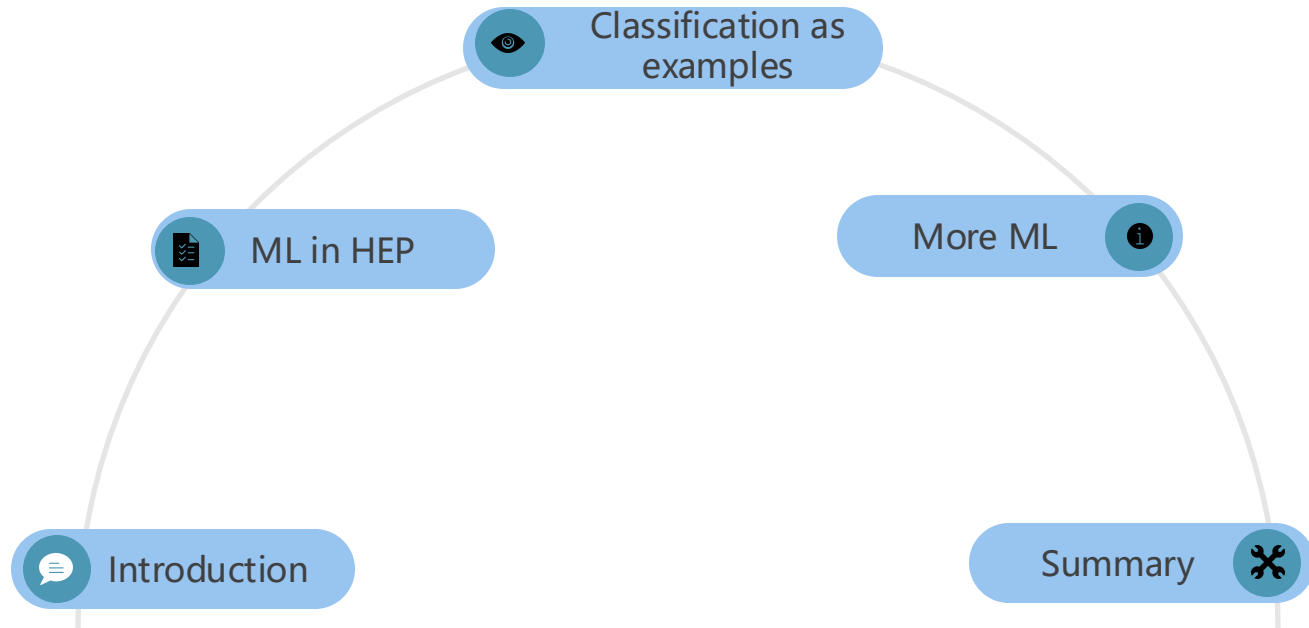
IHEP School of Computing 2024, 2024.08.21-23

北京延庆

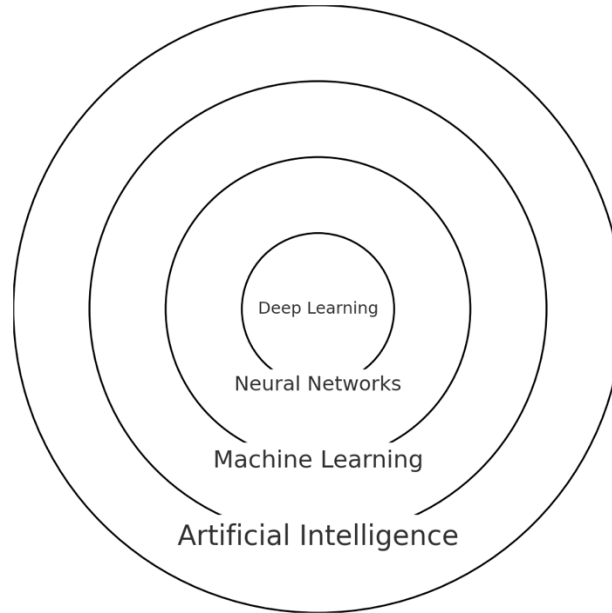
Disclaimers

- This is a very personal review, highly biased
- And mainly focusing on classification problems in offline data processing

Outline

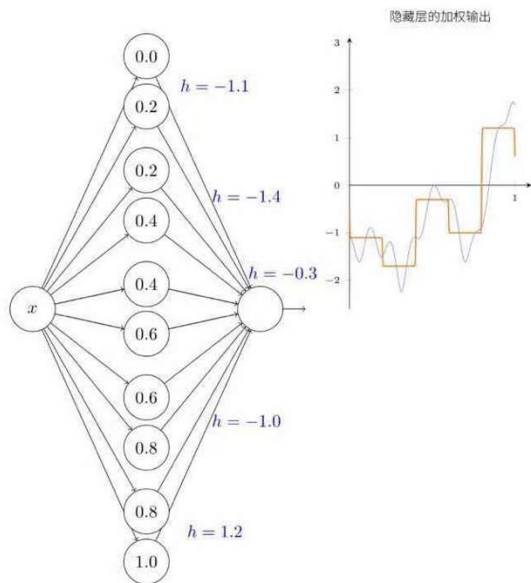


What is Machine Learning ?

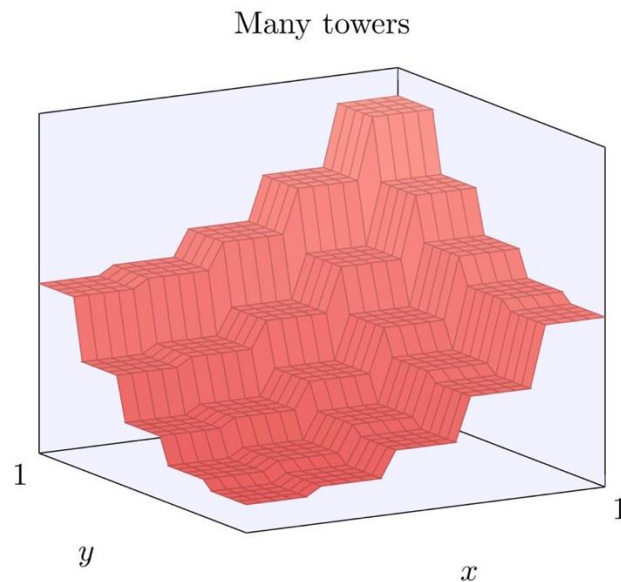


- ✓ Field of study that gives computers the ability to learn without being explicitly programmed
- ✓ A set of rules that allows systems to learn directly from examples, data and experience
- ✓ A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E
- ✓ Machine learning is a set of methods that can automatically detect patterns in data, and then use the uncovered patterns to predict future data or other outcomes of interest
- ✓

事实 1: Neural network as universal function approximator



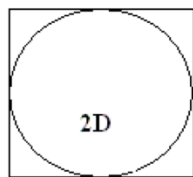
1D



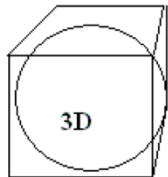
2D

A notable fact about neural networks is that they can approximate a continuous function to any desired level of precision, provided that there are enough neurons in the hidden layers.

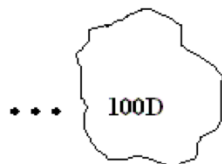
事实 2 : Curse of dimensionality



ratio: $4/\pi = 1.27$



ratio: $6/\pi = 1.91$

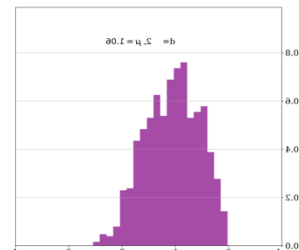
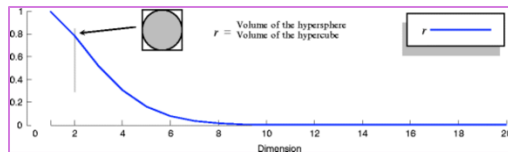


ratio: $4.2 \cdot 10^{39}$

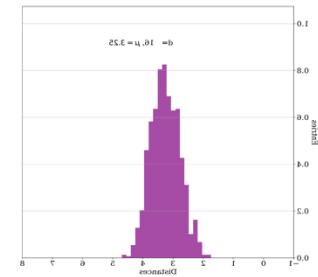
$$\frac{A_{\text{circle}}}{A_{\text{square}}} = \frac{\pi}{4} \text{ for } d=2$$

$$\frac{V_{\text{sphere}}}{V_{\text{cube}}} = \frac{\pi}{6} \text{ for } d=3$$

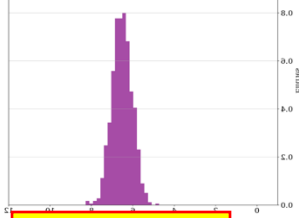
$$\frac{V_{\text{hypersphere}}}{V_{\text{hypercube}}} = \frac{\pi^{d/2}}{d2^{d-1}\Gamma(d/2)} \rightarrow 0 \text{ as } d \rightarrow \infty$$



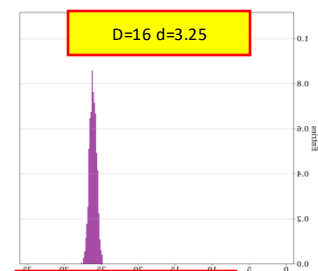
D=2, d=1.06



D=16, d=3.25



D=64, d=6.5



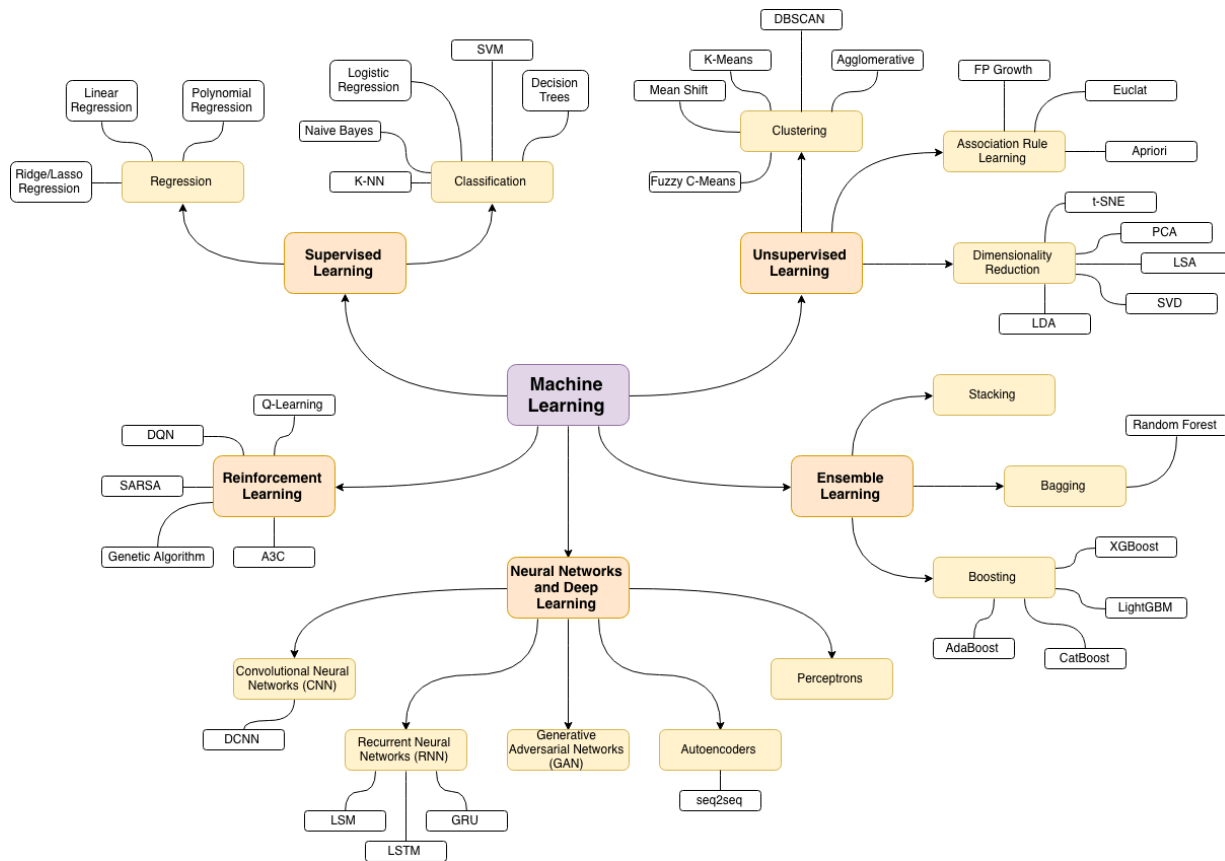
D=1024, d=26.12

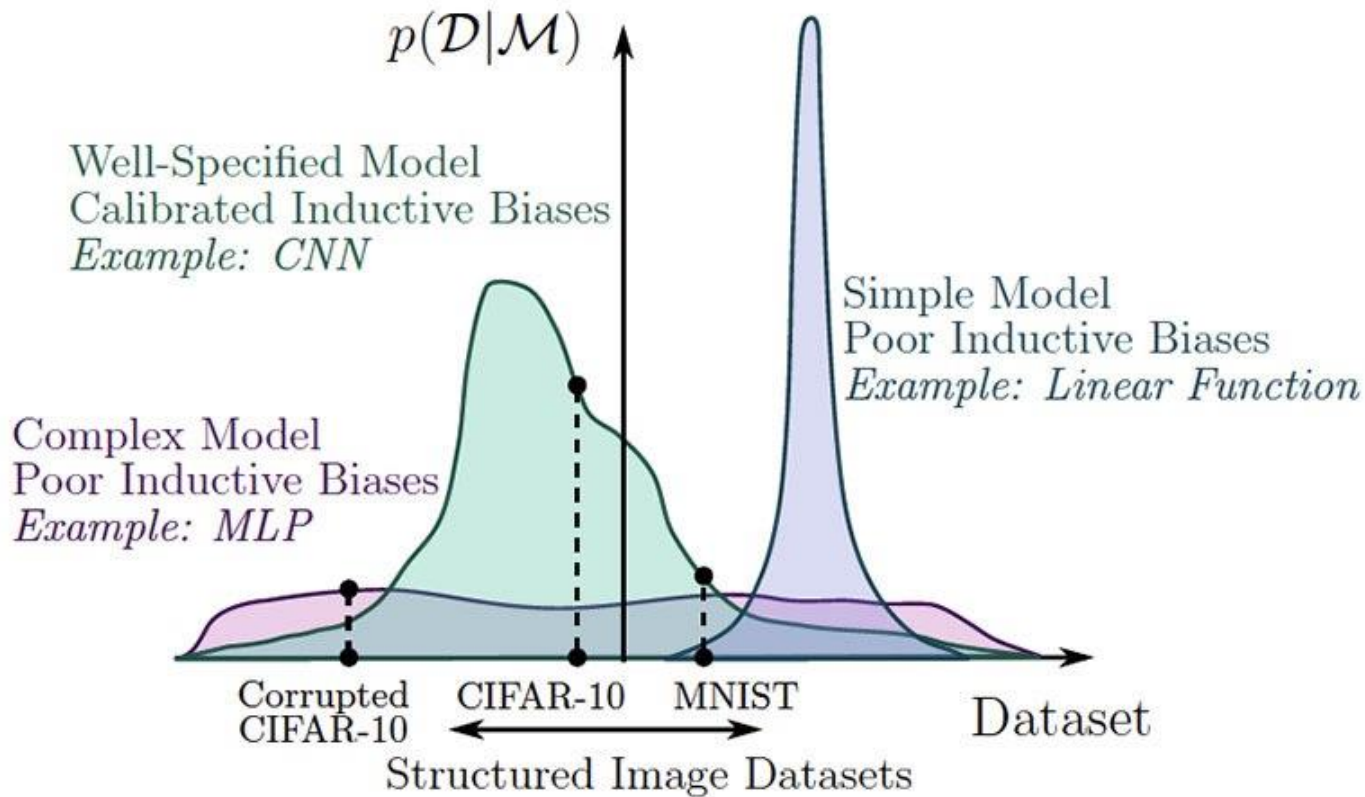
- When D=1: 100 evenly distributed points can sample a unit interval with a distance no greater than 0.01;
- When D=10: it requires 10^{20} sampling points to achieve the same sampling rate.
- Almost all points in high-D are isolated

On one hand: fortunately most specific problems can be reduced in dimensionality!

Neural networks have demonstrated their ability to effectively address the dimension problem!

事实 3: No free lunch theorem (<http://www.no-free-lunch.org>)
There is no single algorithm that is universally the best for all problems
Performance of a learning algorithm is problem-specific





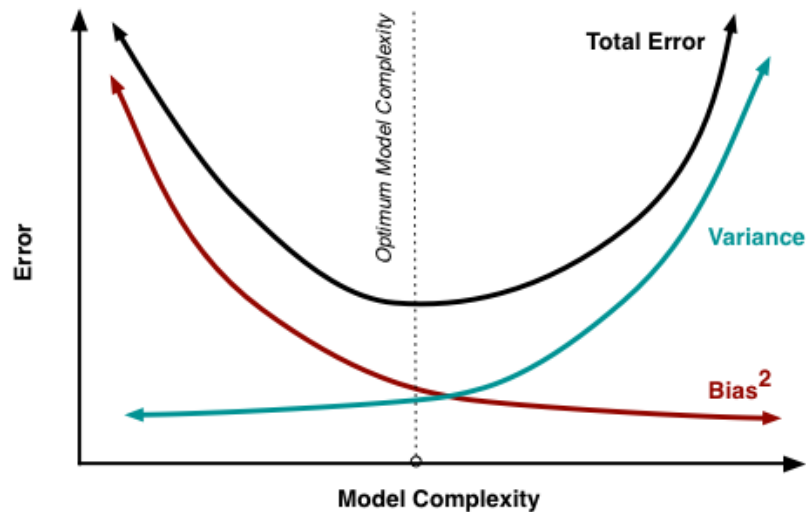
Why do some models perform well on certain datasets? Inductive bias

事实 4: bias-variance tradeoff

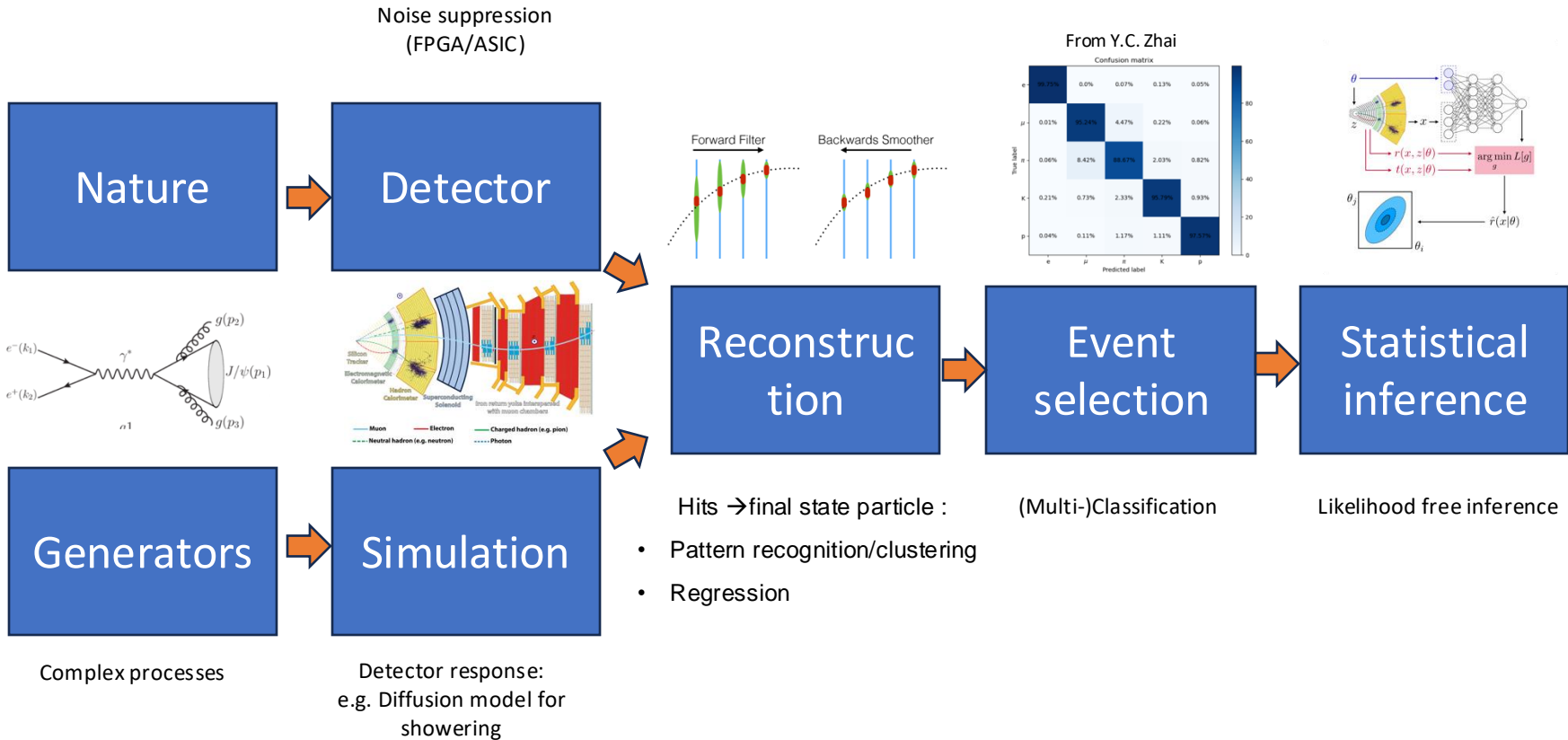
- 机器学习模型的目标是最小化预测误差。
- 预测误差可以分解为偏差、方差和不可减少的误差。
- Bias: 模型预测值与真实值之间的差异。高偏差导致模型在训练集上的表现不佳。
- Variance :模型在不同数据集上的预测结果的波动。高方差导致模型在新数据上的泛化能力差

$$\text{Total Error}^2 = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}^2$$

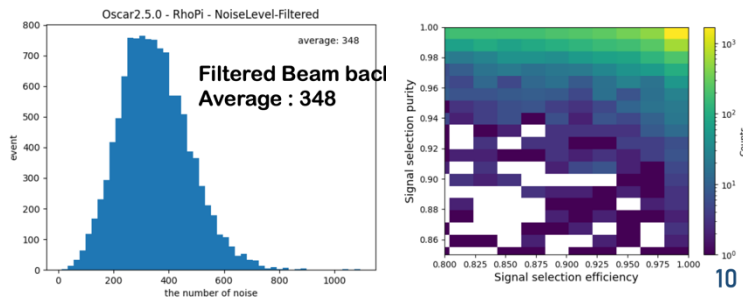
- ✓ 理解偏差和方差的关系对于构建有效模型至关重要。
- ✓ 持续调整模型以找到最佳平衡点。



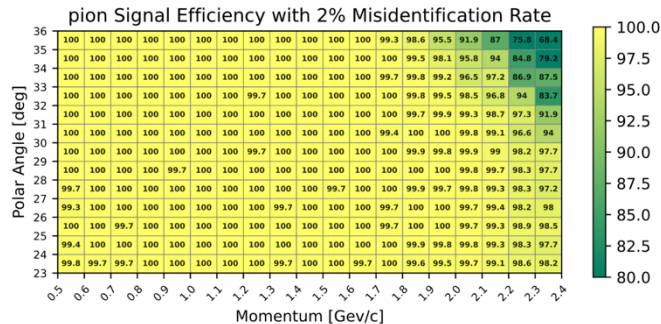
ML in HEP experiments



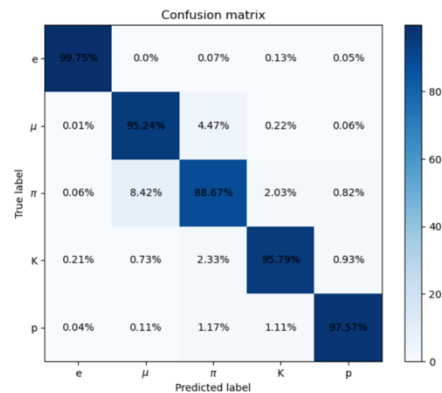
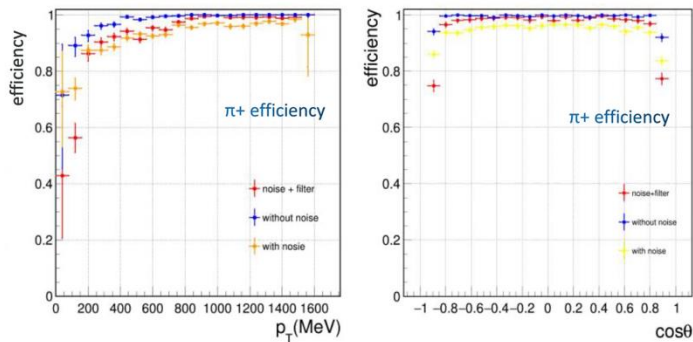
X. Jia: CNN for tracking



Z. Yao: CNN for PID



Y. Zhai: BDT for (global)PID

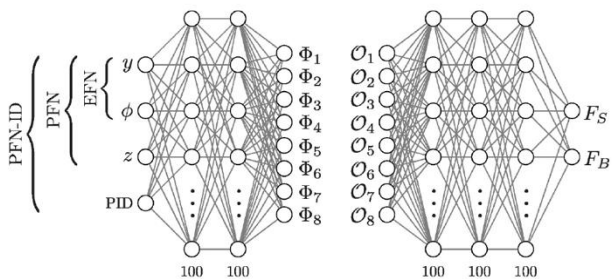


(Multi-)Classification problem

- Jet tagging/W tagger
- Event classification

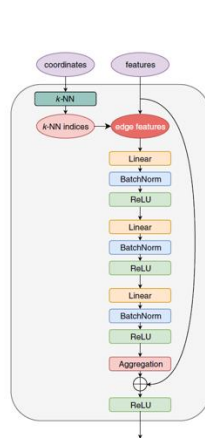
Algorithms used

Energy Flow Network(EFN) / Particle Flow Network(PFN)



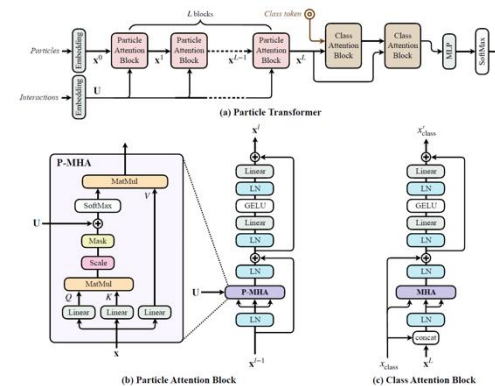
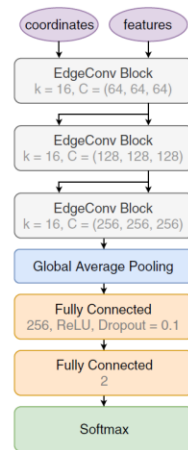
P. T. Komiske, E. M. Metodiev and J. Thaler
[\[JHEP01\(2019\)121\]](#)

ParticleNet



H. Qu and L. Gouskos [\[Phys.Rev.D 101 \(2020\) 5, 056019\]](#)

ParticleTransformers (ParT)

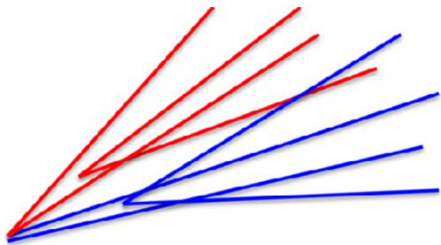


H. Qu, C. Li, S. Qian [\[2202.03772\]](#)

Jet (flavor) tagging(单个喷注)

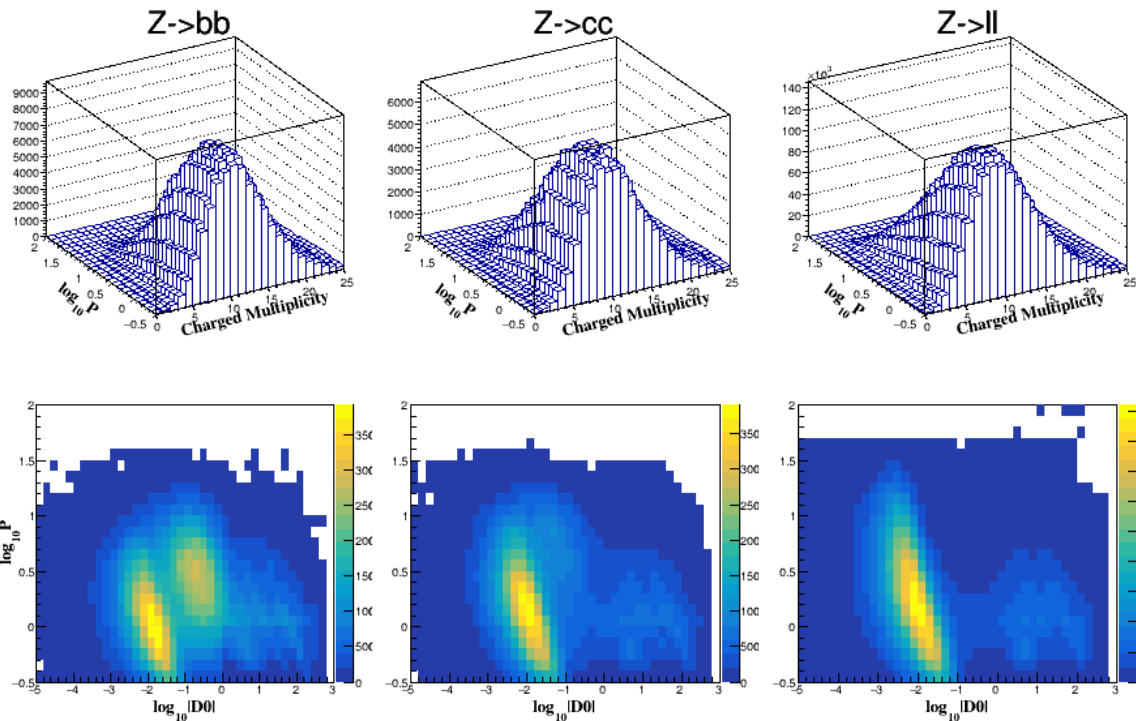
- 91 GeV
- $Z \rightarrow bb, cc, ll$ (uu,dd,ss)
- 450k events (900k jets) for each class

- Take particle level information as input

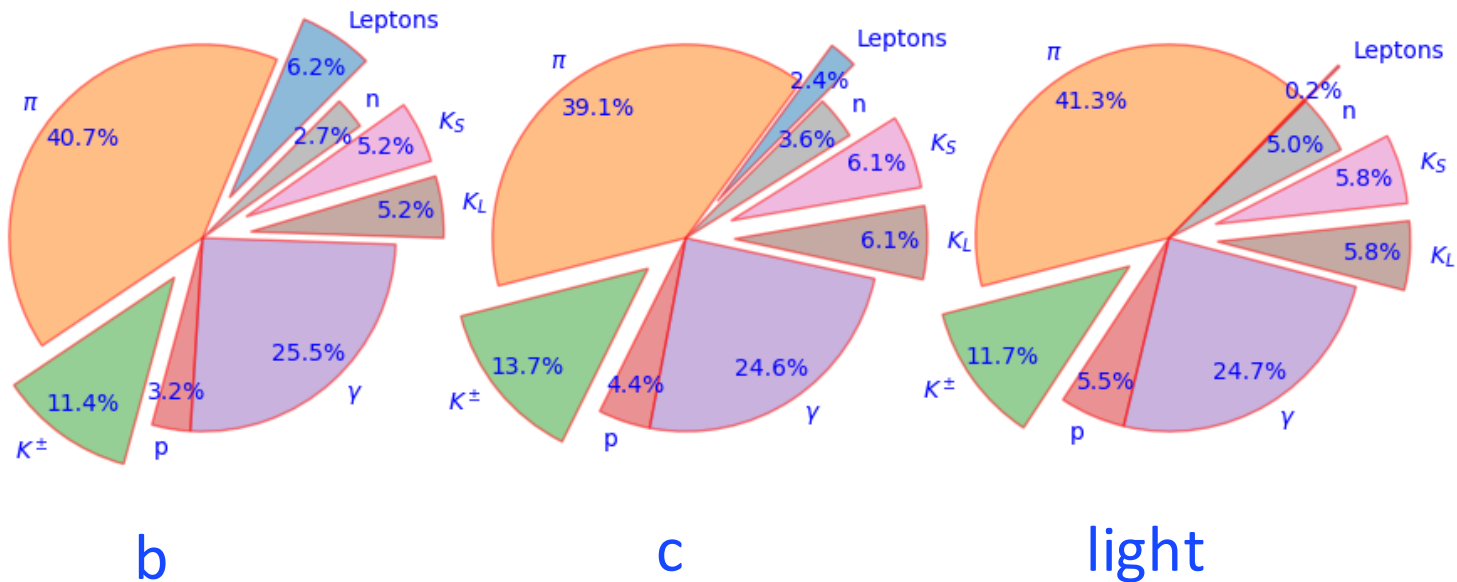


- 4-momenta
- d_0/z_0
- PID
-

Multiplicity, impact parameters



PID information



Weighted by momenta

Accuracy \longrightarrow

Algorithm	ParticleNet	PFN	DNN	BDT	GBDT	gforest	XGBoost
Accuracy	0.872	0.850	0.788	0.776	0.794	0.785	0.801

Purity \times efficiency \longrightarrow

tag	$\epsilon_S(\%)$	$\epsilon \times \rho$			
		LCFIPlus	XGBoost	ParticleNet	PFN
b	60	-	-	0.589	0.596
	70	-	-	0.694	0.689
	80	-	0.747	0.780	0.763
	90	0.72	0.713	0.810	0.752
	95	-	0.609	0.721	0.645
c	60	0.36	-	0.548	0.485
	70	-	-	0.589	0.497
	80	-	0.345	0.584	0.467
	90	-	0.292	0.516	0.402
	95	-	0.251	0.451	0.348

Take c-tagging as example

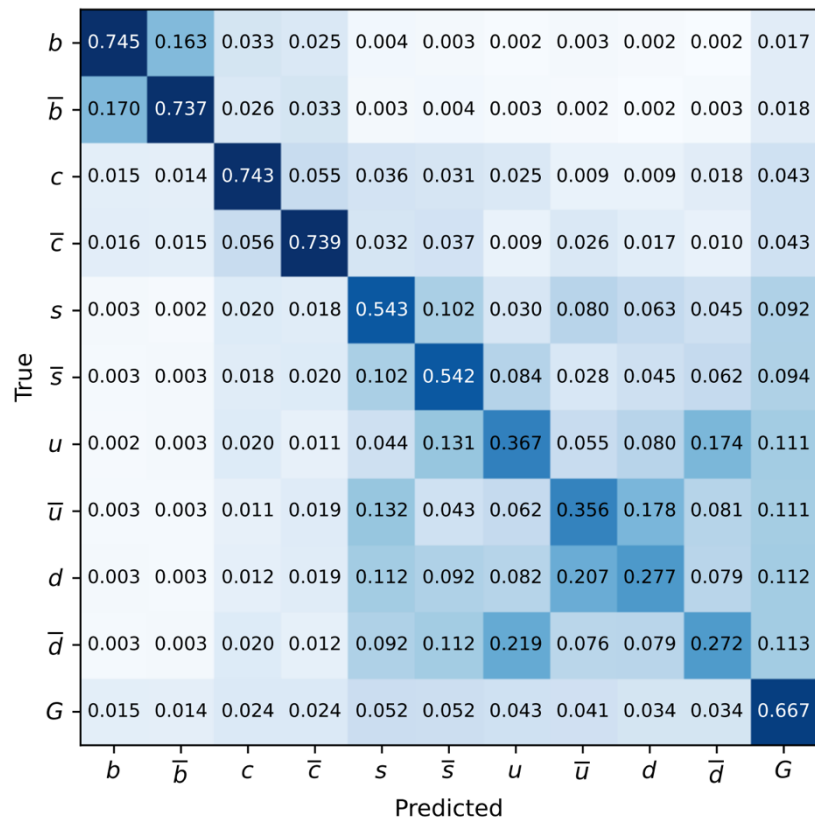
$$\text{sqrt}(0.584/0.345)=1.3$$

Statistical uncertainty: 30% \downarrow

$$\frac{1}{(\Delta\sigma_s)^2} = \frac{1}{\sigma_s} \mathcal{L}^{\epsilon_s \rho} = \frac{1}{\sigma_s^2} S_{\text{tot}} \epsilon_s \rho$$

11 classes

Ambitious test by M. Ruan

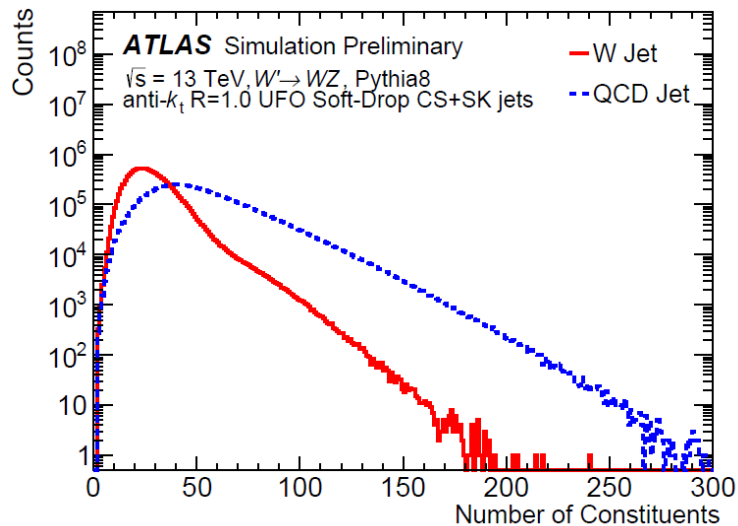


Phys. Rev. Lett. 132, 221802 (2024)

W Jet Taggers (ATLAS, by Shudong Wang)

(一对喷注)

- In this study, a maximum of 200 constituents are considered by all constituent-based taggers. Only a small portion of jets in the dataset have more than 200 constituents (less than 0.04%). As jet constituents are sorted by decreasing p_T , truncation eliminates the softest constituents of the jet.



Distributions of the number of constituents in a large- R jet.

W Jet Taggers

- Particle Flow Network(PFN)/Energy Flow Network(EFN)

- Based on Deep Sets Theorem
- [JHEP01\(2019\)121](#)

- ParticleNet

- Customized graph neural network architecture for jet tagging with the point cloud approach
- [Phys.Rev.D 101 \(2020\) 5, 056019](#)

- ParticleTransformer

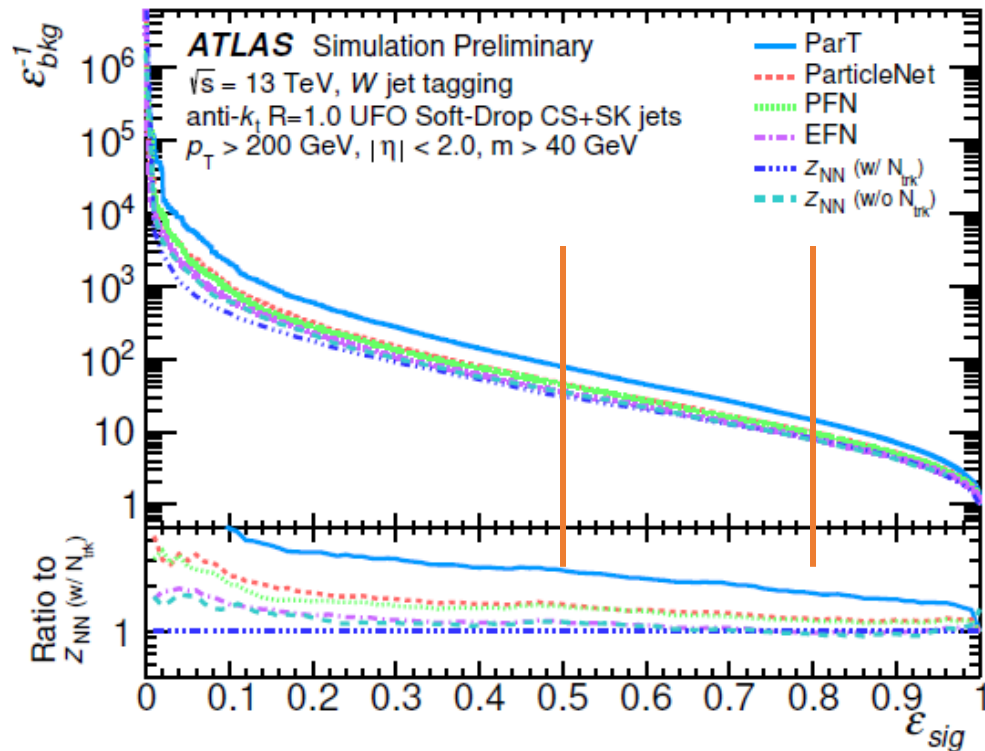
- Transformer designed for particle physics
- [arxiv: 2202.03772](#)

Models	Input variables
EFN	$\Delta\eta, \Delta\phi, \ln p_T$
PFN	$\Delta\eta, \Delta\phi, \ln p_T, \ln E, \ln \frac{p_T}{\sum_{jet} p_T}, \ln \frac{E}{\sum_{jet} E}, \Delta R$
ParticleNet	$\Delta\eta, \Delta\phi, \ln p_T, \ln E, \ln \frac{p_T}{\sum_{jet} p_T}, \ln \frac{E}{\sum_{jet} E}, \Delta R$
ParticleTransformer	$\Delta\eta, \Delta\phi, \ln p_T, \ln E, \ln \frac{p_T}{\sum_{jet} p_T}, \ln \frac{E}{\sum_{jet} E}, \Delta R$ (E, p_x, p_y, p_z)

Tagger Performance

Calculated using samples with steeply falling pT spectra, i.e. both sig & bkg are weighted to have falling pT spectra.

For a signal efficiency of 0.5 (0.8) case, the background rejection of ParticleTransformer is about 1.8-2.8 (1.6-2.7) times better than the baseline tagger.



Tagger Performance

Model	AUC	ACC	ε_{bkg}^{-1} @ $\varepsilon_{sig} = 0.5$	ε_{bkg}^{-1} @ $\varepsilon_{sig} = 0.8$	# Params	Inference Time
EFN	0.920	0.835	35.1	7.95	56.73k	0.065 ms
PFN	0.931	0.853	44.7	9.50	57.13k	0.11 ms
ParticleNet	0.933	0.826	46.2	9.76	366.16k	0.36 ms
ParticleTransformer	0.951	0.880	77.9	14.6	2.14M	0.28 ms

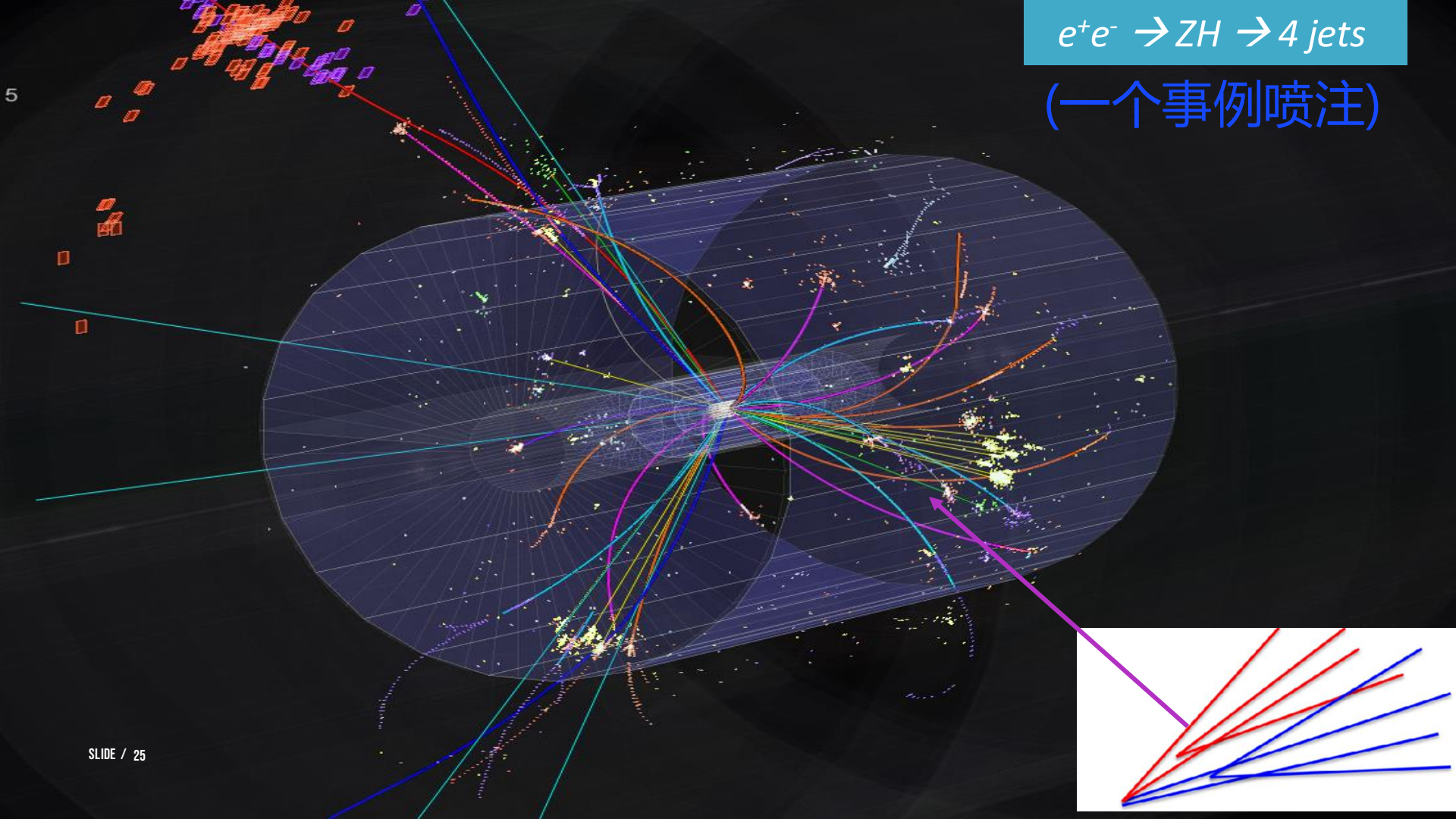
Table 3: The performance of each W jet tagger is measured with several metrics evaluated on the testing set.

Transformers the best

But the # of parameters is almost one order of magnitude larger

$e^+e^- \rightarrow ZH \rightarrow 4 \text{ jets}$

(一个事例喷注)



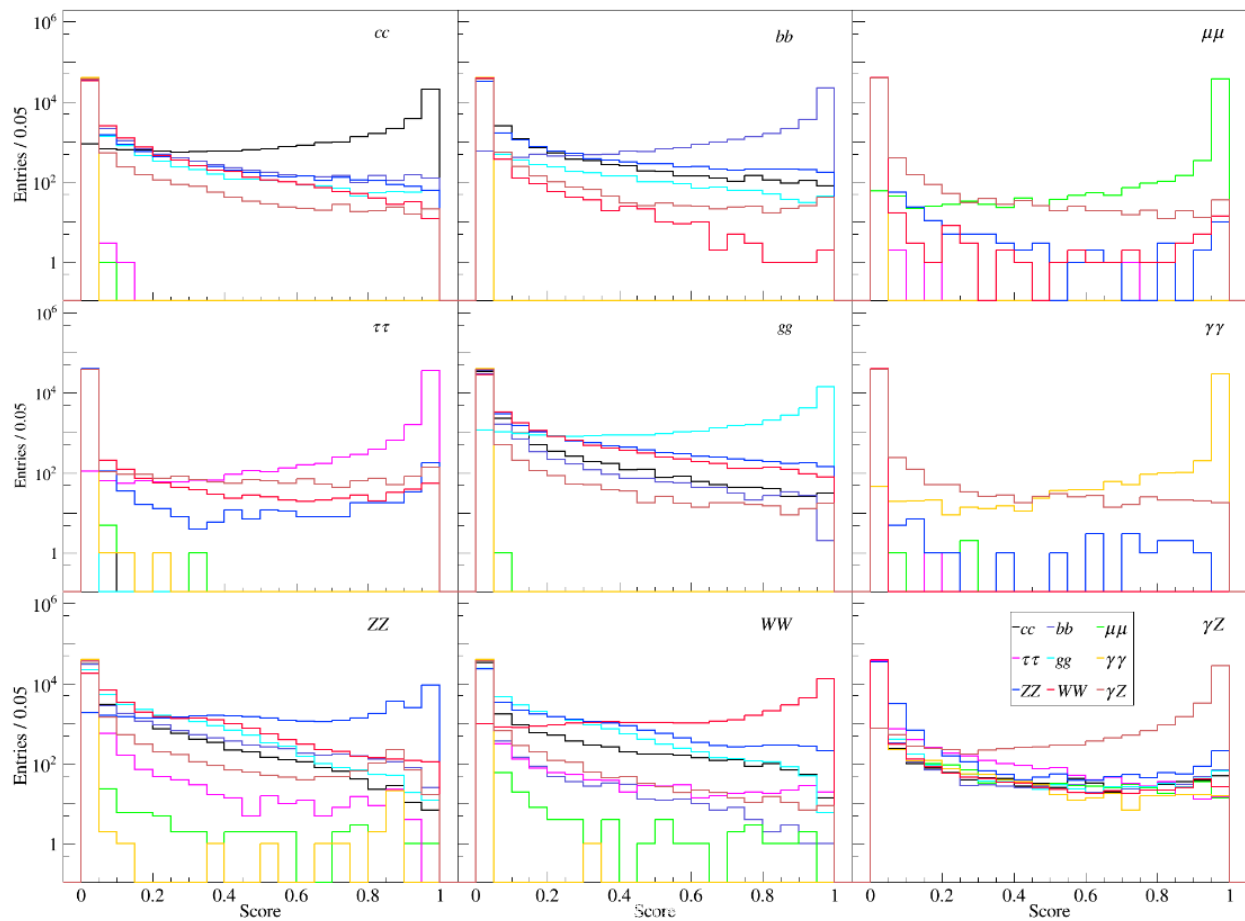
Many processes are selected simultaneously

Prod/decay	cc	bb	mm	$\tau\tau$	gg	gg	WW	ZZ	aZ	ee, uu,dd,ss
eeH	3	1	5	2	4	1	2	3	5	Not covered yet
mmH	3	1	5	2	4	1	2	3	5	
$\tau\tau$ H	3	1	5	2	4	1	2	3	5	
qqH	4	1	2	1	2	5	5	5	3	
nnH	5	1	3	2	3	5	4	2	4	

Consider: $\psi(2S) \rightarrow \pi^+ \pi^- J/\psi$, $J/\psi \rightarrow$ various processes

Try eeH first

Probability distributions of each class

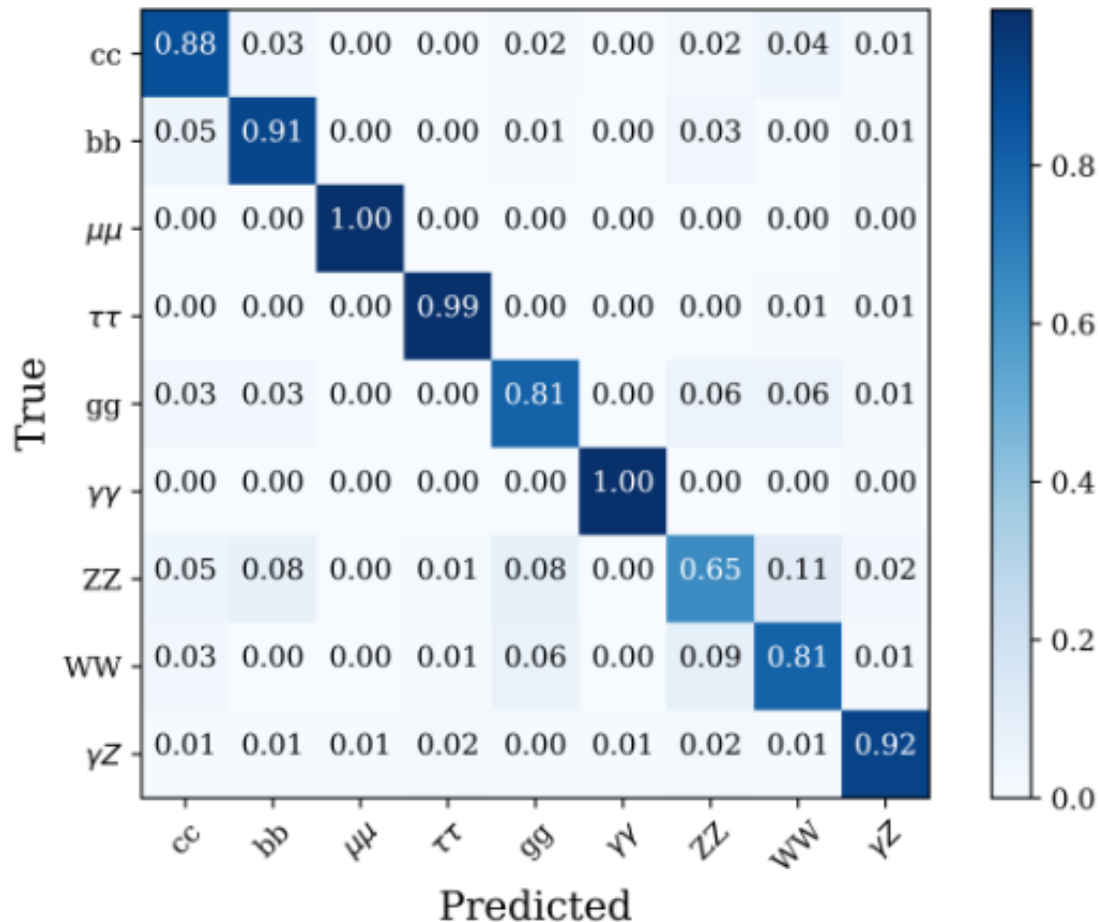


Try eeH first

Sufficiently good performance

Average Accuracy ~ 87%

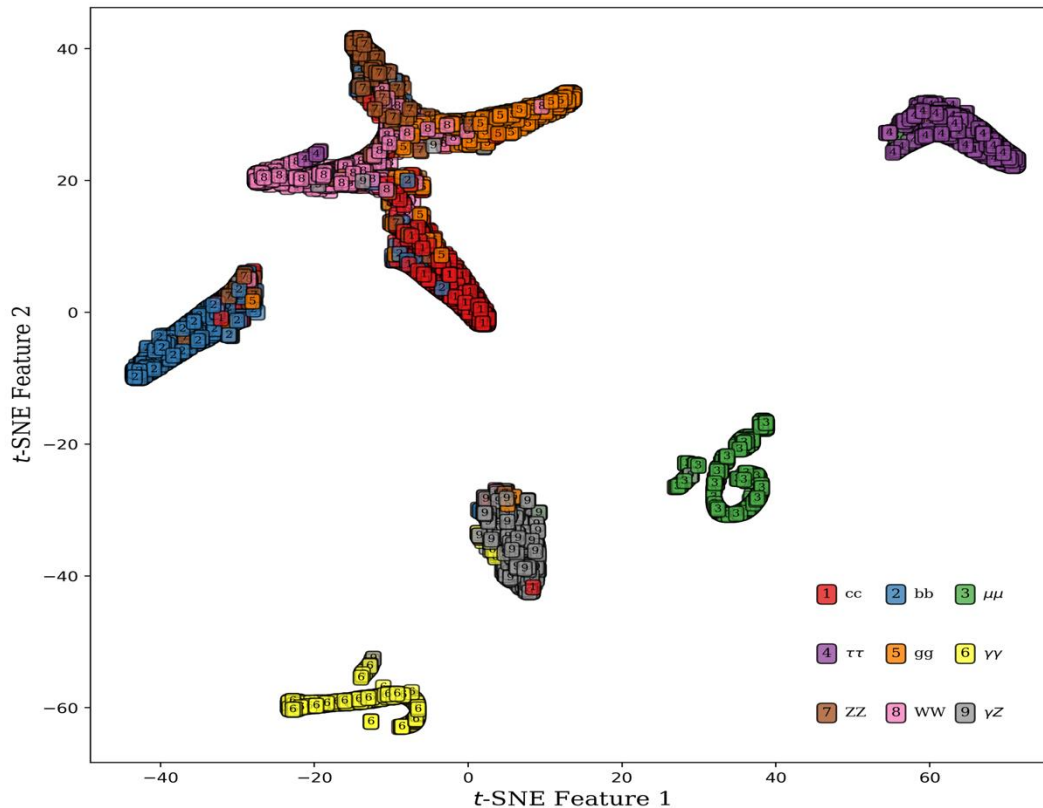
(11% for random guess)



Taking the one has largest probability (ArgMax)

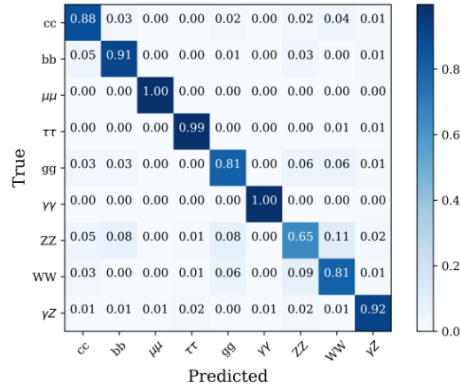
Dimension reduction tells us more

- ✓ $\mu\mu$, $\gamma\gamma$, $\tau\tau$ well classified as expected
- ✓ bb and γZ also good
- ✓ cc , gg , WW , and ZZ fake each other, but under control

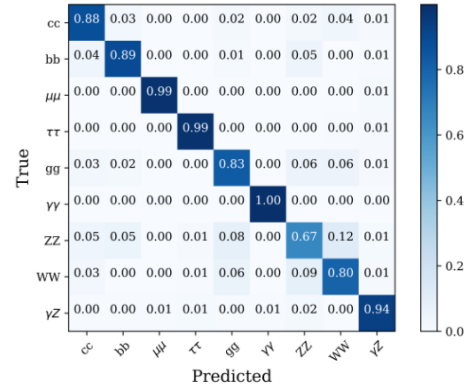


Dimensional reduction (t-SNE)

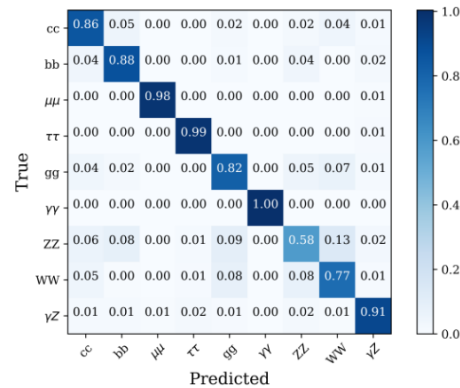
All 4 production modes



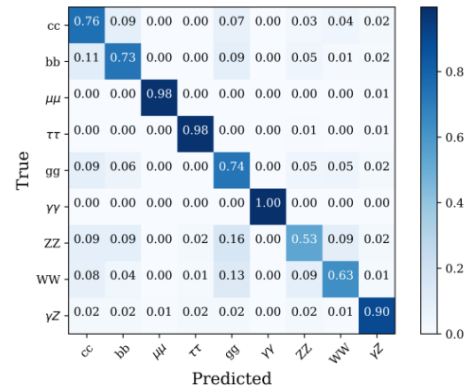
eeH



$\mu\mu H$



$\tau\tau H$



qqH



UC San Diego



Fermilab (LHC 上的 “事例”)

Based on [arXiv:2405.12972](https://arxiv.org/abs/2405.12972)

[\[Github\]](#) [\[Google Colab\]](#)

Sophon meets LHC: Accelerating resonance discovery via signature- oriented pre-training

Congqiao Li (李聪乔), *Peking University*

based on the work with our colleagues in the CMS Collaboration:

Antonios Agapitos¹, Jovin Drews², Javier Duarte³, Dawei Fu¹, Leyun Gao¹, Raghav Kansal³, Gregor Kasieczka²,
Louis Moureaux², Huilin Qu⁴, Cristina Mantilla Suarez⁵, Qiang Li¹

1) Peking U. 2) Hamburg U. 3) UC San Diego 4) CERN 5) FNAL

also thanks Yuzhe Zhao¹ for his contribution

Quantum Computing and Machine Learning Workshop 2024, Changchun

6 August, 2024

More?

Diffusion models

Large Language Models

... ..

Used as copilot-like assistant: Dr. Sai

Used to HEP data directly: tokenization is a key

How to represent a HEP events? Tokenization

Feature engineering

Some mathematical methods? Such as fox-wolfram moments

$$H_l \equiv \left(\frac{4\pi}{2l+1} \right) \sum_{m=-l}^{+l} \left| \sum_{\mathbf{i}} Y_l^m(\Omega_{\mathbf{i}}) \frac{|\vec{\mathbf{p}}_{\mathbf{i}}|}{\sqrt{s}} \right|^2$$
$$= \sum_{\mathbf{i}, \mathbf{j}} \frac{|\vec{\mathbf{p}}_{\mathbf{i}}| |\vec{\mathbf{p}}_{\mathbf{j}}|}{s} P_l(\cos\varphi_{\mathbf{i}\mathbf{j}}),$$

Autoencoder?

作业：实战内容（借用李聪乔的劳动）

- Particle Transformer 用作分类
- <https://gitee.com/colizz/qcml-2024-tutorial>
- `lxslcaaa.ihep.ac.cn` → `lxloginaaa.ihep.ac.cn`（需要修改的地方）
 - ✓ chap1 是视频教学：看看就行
 - ✓ chap2 是 general 的展示：需要科学上网，了解下就行
 - ✓ chap3 我实测可以跑（CPU）：复现结果

Summary

- Machine learning is statistical learning (NFL)
- Machine learning is useful (~~CoD~~): high dimensional HEP data
- **Machine learning method with proper bias** is helpful and easy to explain
- Machine learning methods can be applied to almost all aspects of HEP experiments.
- LLM demonstrated astonishing capabilities, which are worth exploring from two aspects:
 - Use LLMs as language-based assistants
 - Employing LLMs to directly to process data: how to represent HEP events is the key