

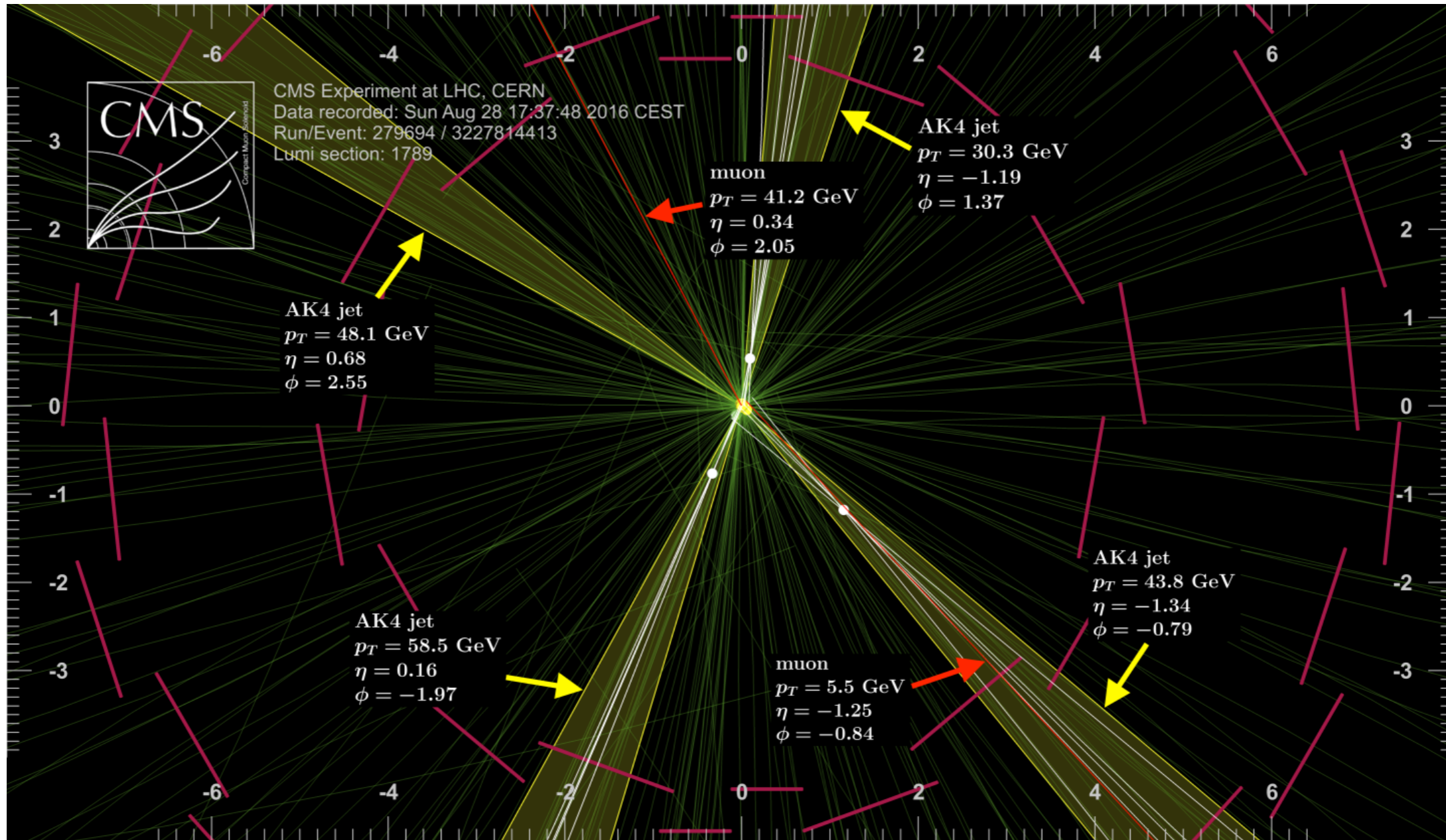
Transformer models for Heavy Flavor Jet Identification in CMS



Alexandre De Moor, Congqiao Li, Denise Müller, Huilin Qu, Sitian Qian
on behalf of CMS Collaboration

ML4Jets 2022 @ Rutgers University, 2022/11/01

INTRODUCTION: JET TAGGING



Jet tagging:

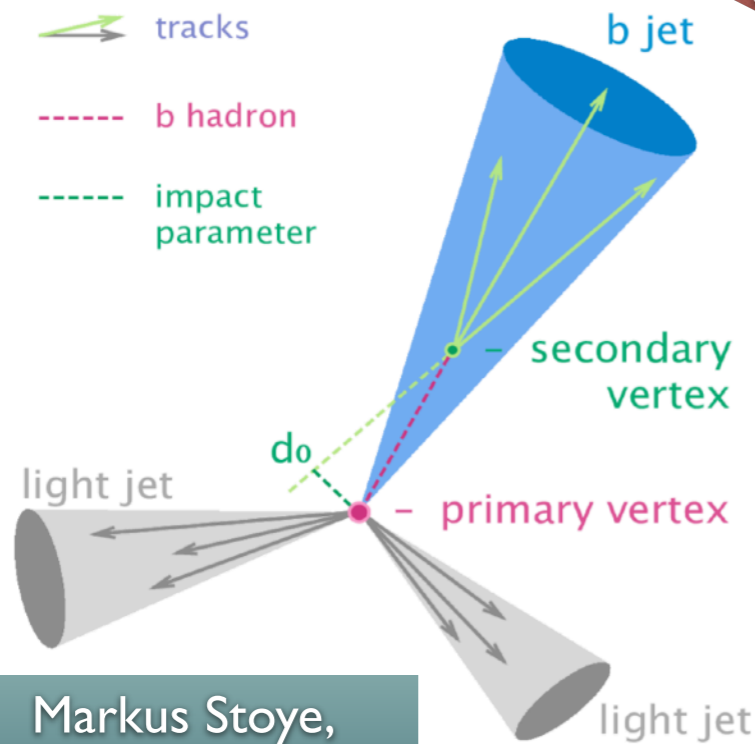
Identifying the particle that initiates the jet with experimentally observed quantities

INTRODUCTION: JET TAGGING



Different types of jet tagging

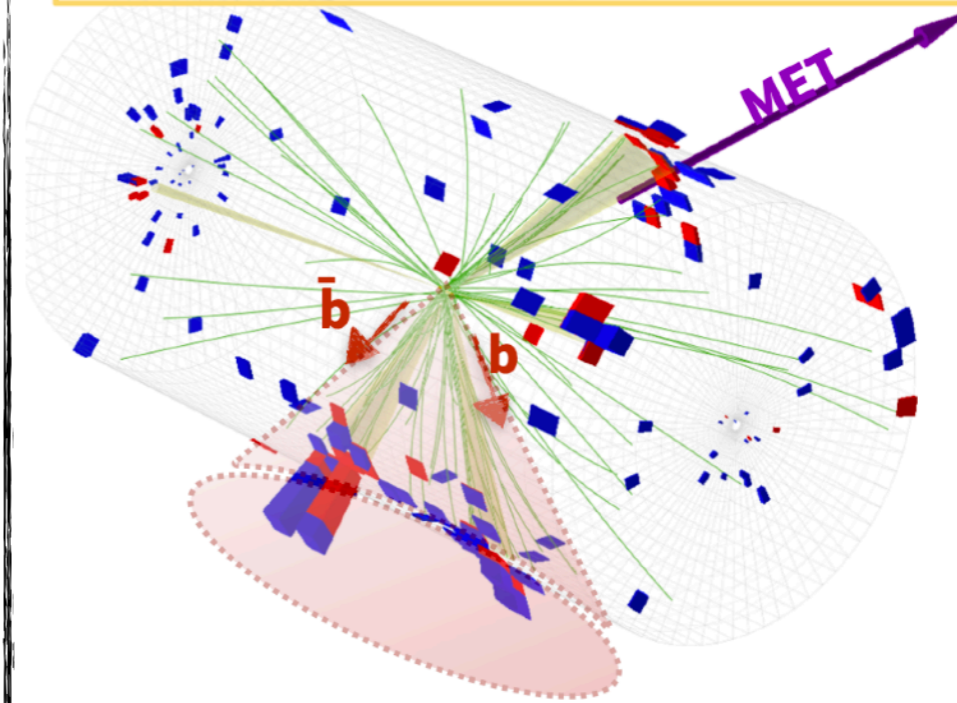
Today's focus



Markus Stoye,
HAP workshop 2018

Small radius tagging:
Mainly on the flavor
information (light quark/
gluon/charm/beauty/tau
lepton...)

a simulated $Z(\nu\nu)H(bb)$ event from CMS: the boosted $H(bb)$ object can be tagged with a large-R jet



Congqiao Li, ML4jets 2021

Large radius tagging:
Targeting on boosted heavy
resonance (W/Z/H/Top...)

HEAVY FLAVOR TAGGING



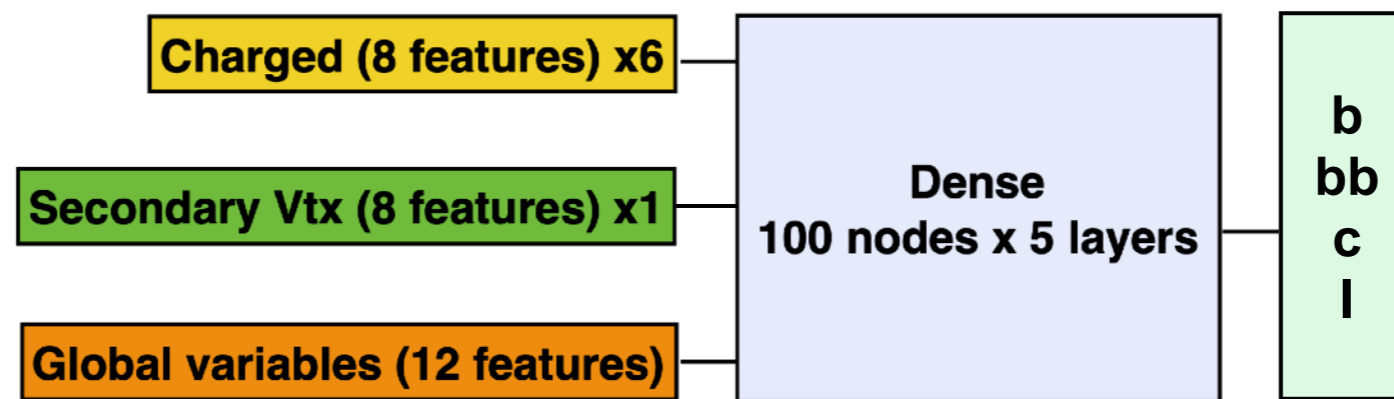
- Problem setup of heavy flavor (HF) tagging:
 - Focus on small radius jet (in CMS: anti-kT with $R=0.4$)
 - Usually 3 cases:
 - Tagged as “b”
 - Tagged as “c”
 - Tagged as light flavor (“udsg”)
 - Can involve more categories including hadronic tau, “uds” vs “gluon” etc.
- Key features of HF tagging:
 - Track displacement is of vital importance
 - Due to longer life of heavy flavor hadrons!

HEAVY FLAVOR TAGGING @ CMS

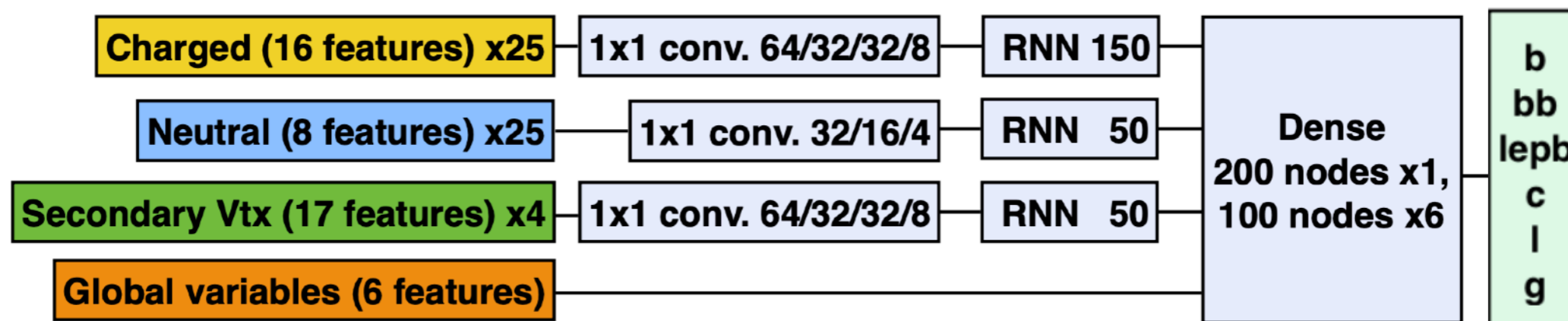


- Brief review of CMS HF tagging evolution:
 - ~ 2016: **CSVv2(NN)** and cMVA_{v2}(BDT)
 - CSVv2 is based on a very simple 1 layer NN
 - 2017: deepCSV

ArXiv: 1712.07158
CSVv2, cMVA_{v2} & deepCSV



- Mid 2017: deepJet (deepFlavour) ArXiv: 2008.10519



HEAVY FLAVOR TAGGING @ CMS



- Brief review of CMS HF tagging evolution:

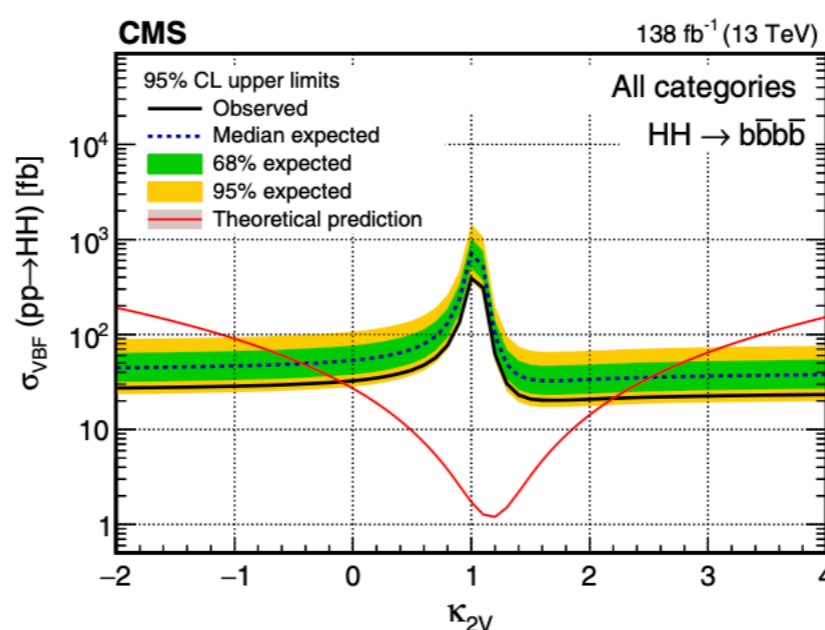
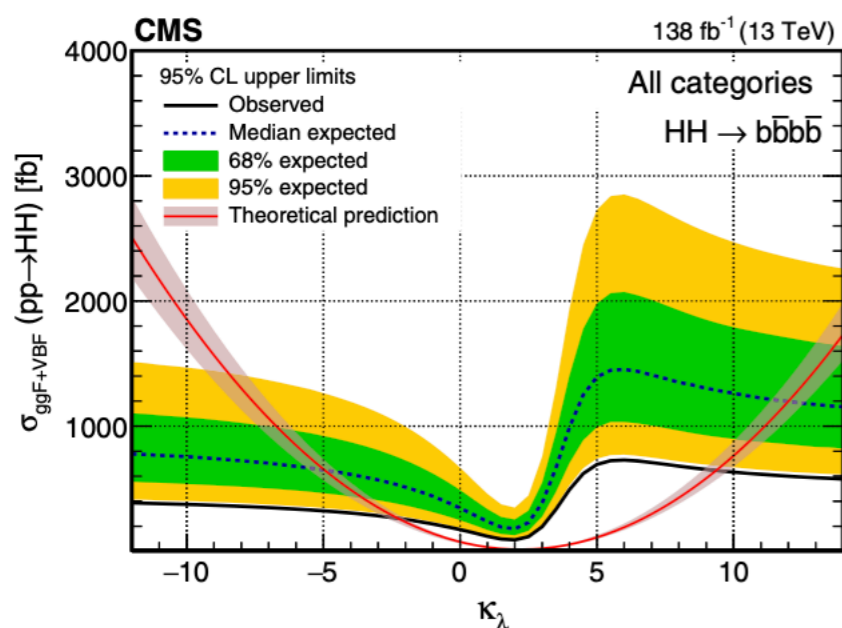
- ~ 2016: **CSVv2(NN)** and cMVA_{v2}(BDT)

ArXiv: 1712.07158
CSVv2, cMVA_{v2} & deepCSV

- 2017: deepCSV

- Mid 2017: deepJet (deepFlavour) ArXiv: 2008.10519

- Rich physical results have continuously come out with the help of CMS HF tagging algorithms!



PhysRevLett. 129.081802

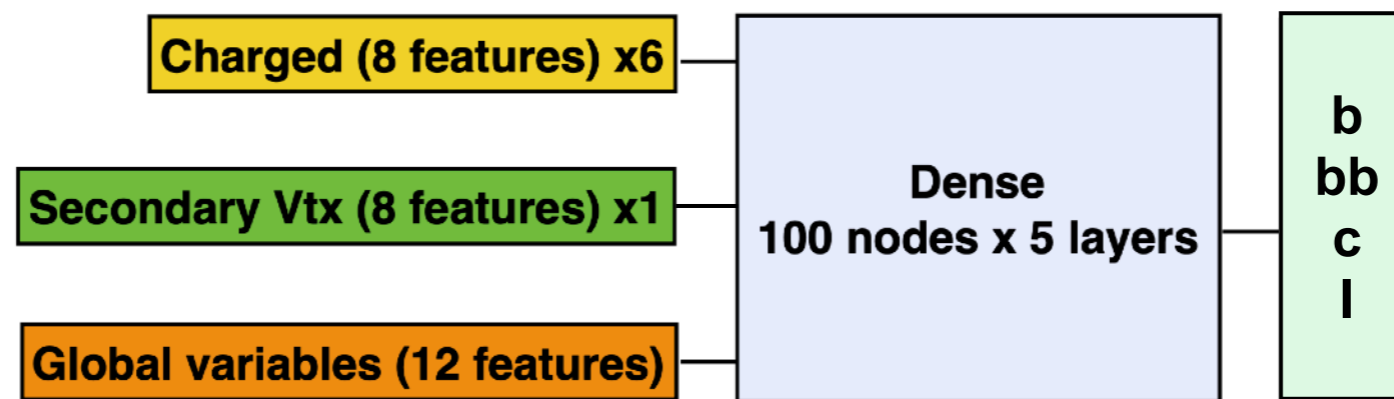
Exciting result from
CMS Full Run2
nonresonant HH4b
measurement, b
tagging is based on
DeepJet algorithm

HEAVY FLAVOR TAGGING @ CMS

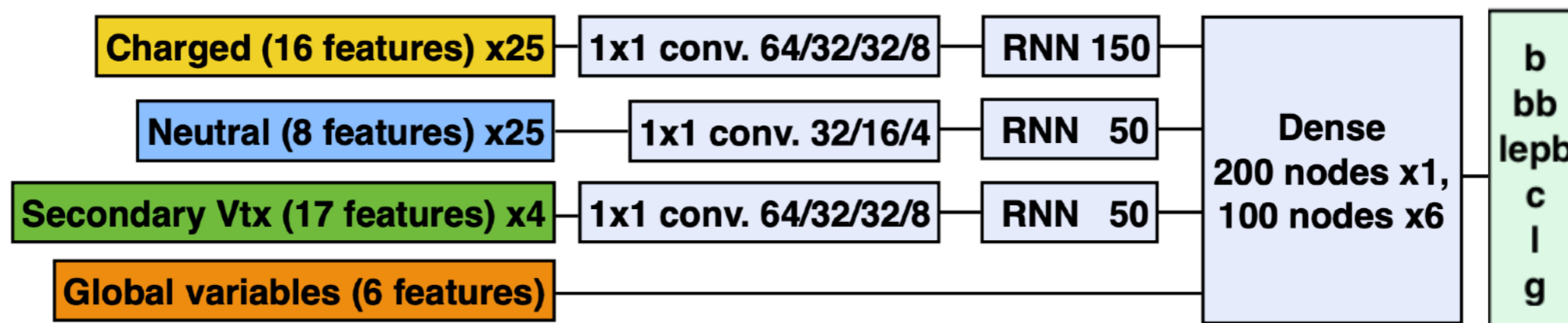


- Brief review of CMS HF tagging evolution:
 - ~ 2016: **CSVv2(NN)** and cMVA_{v2}(BDT)
 - CSVv2 is based on a very simple 1 layer NN
 - 2017: deepCSV

ArXiv: 1712.07158
CSVv2, cMVA_{v2} & deepCSV



- Mid 2017: deepJet (deepFlavour) ArXiv: 2008.10519

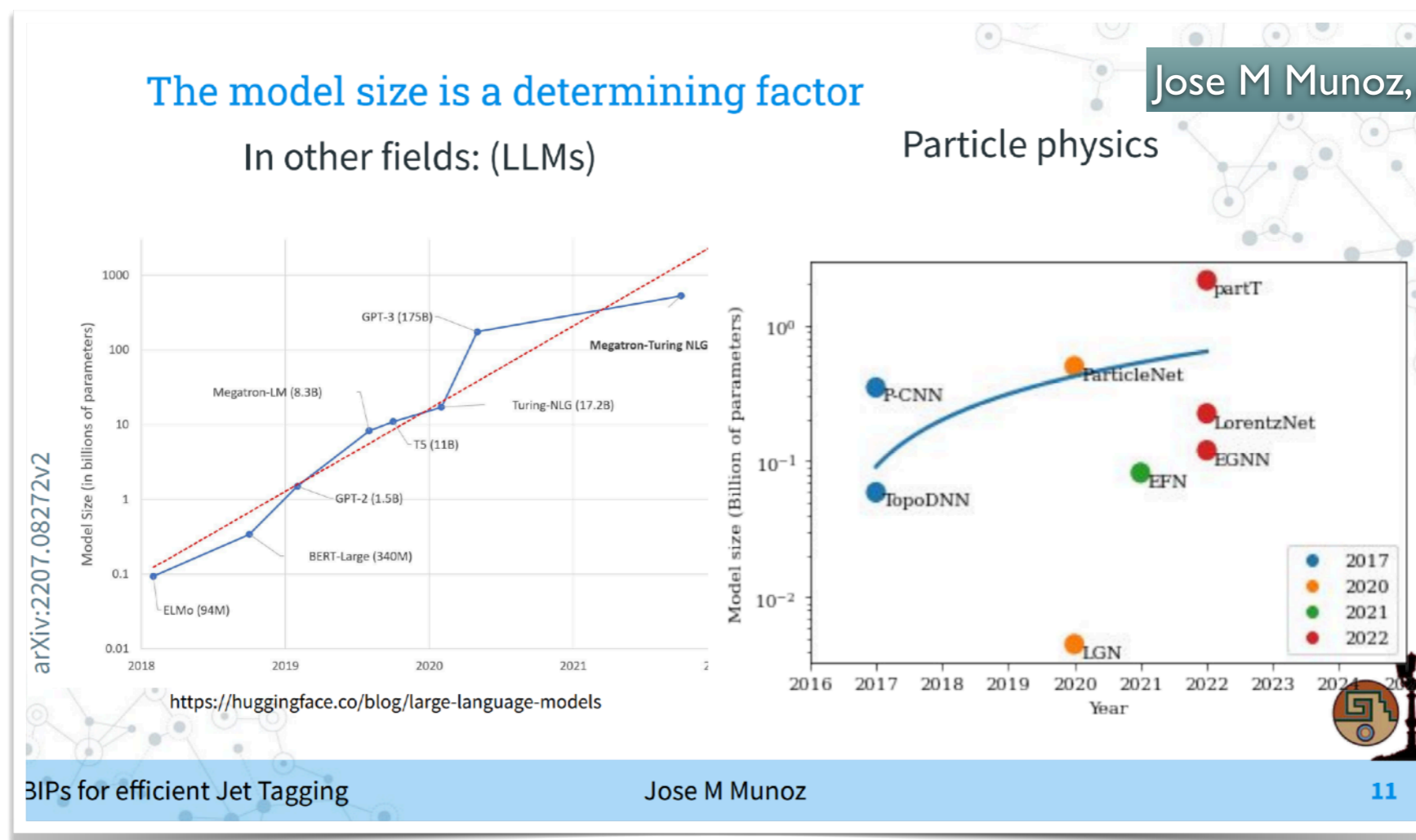


- What about Run3?

EVOLUTION OF JET TAGGING ALGO.



- Observation from CMS HF tagging evolution:
 - Models become larger and larger
- CMS HF tagging is not alone!



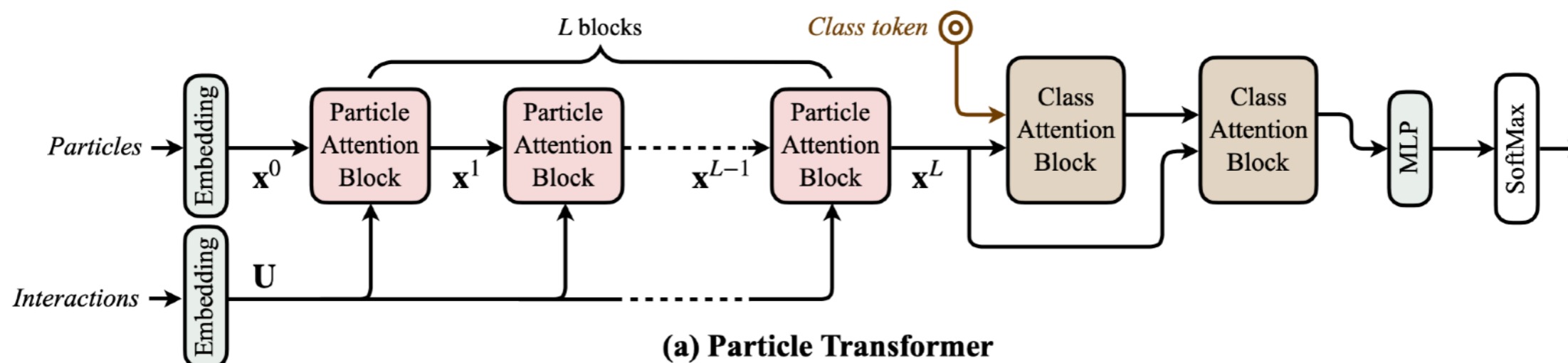
- Larger model is now a trend in not only ML community, but also in HEP! -> Motivates the usage of largest model :) -> Particle Transformer (ParT)

PART 101



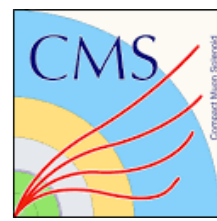
ArXiv: 2202.03772

- Particle Transformer:
 - the transformer designed for particle physics



- Input embedding: Not only inject single particle information, but also include pair-wise features
- Particle Attention Block: Multi-Head Attention (MHA)
 - Pair-wise feature are introduced as the attention mask (P-MHA)
- Class Attention Block: Multi-Head Attention
 - Class token is used for the MHA calculation

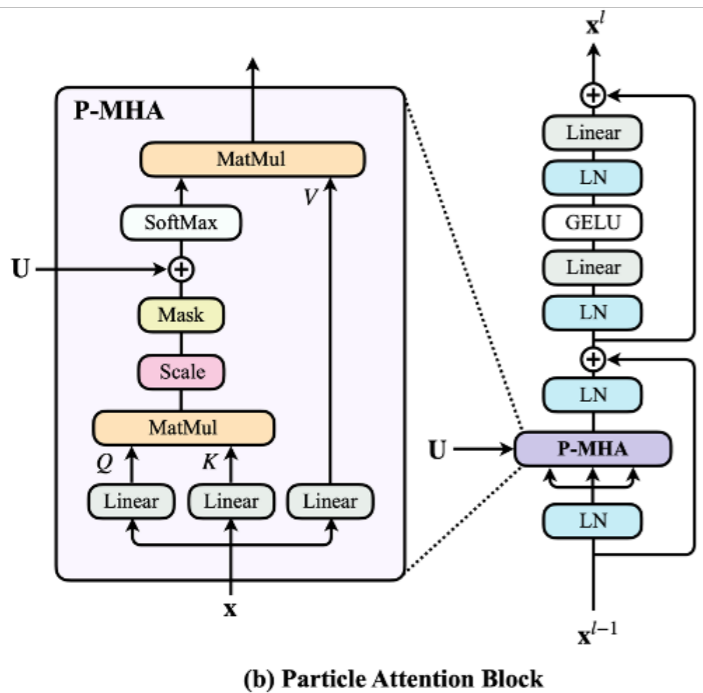
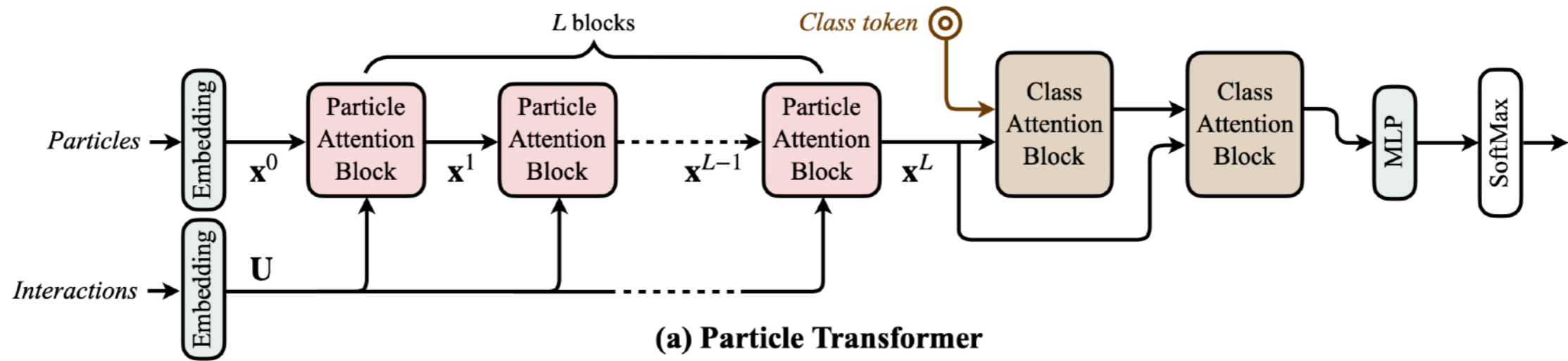
PART 101



Particle Transformer:

ArXiv: 2202.03772

- the transformer designed for particle physics



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

d_k : dimension of K

Choice of the pair-wise features: from LundNet

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

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

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

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

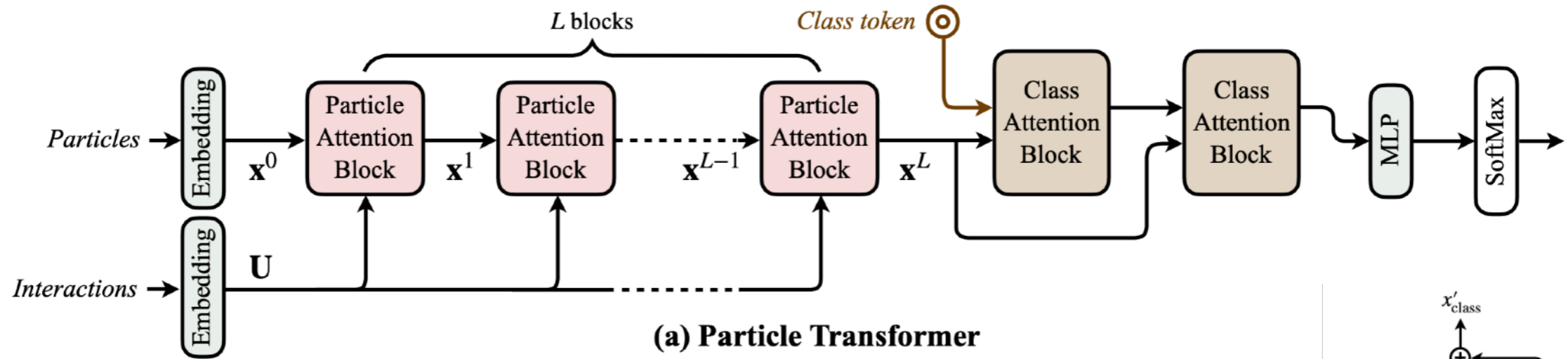
PART 101



- Particle Transformer:

ArXiv: 2202.03772

- the transformer designed for particle physics

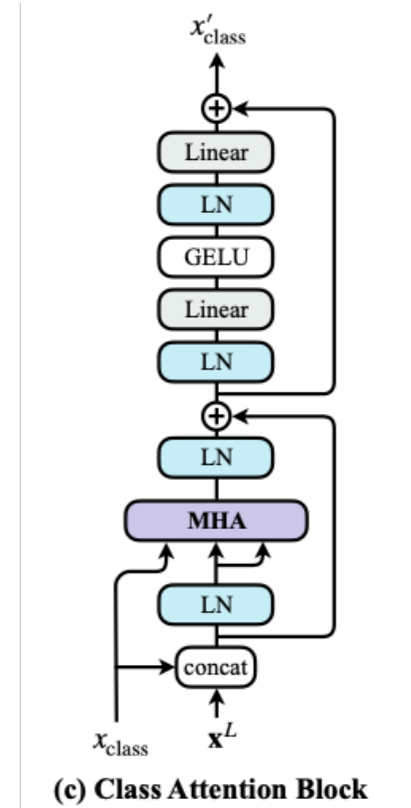


$$\text{MHA}_C(Q_C, K_C, V_C) = \text{SoftMax}(Q_C K_C^T / \sqrt{d_{kC}}) V_C$$

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

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

Concatenate class information and particle embedding

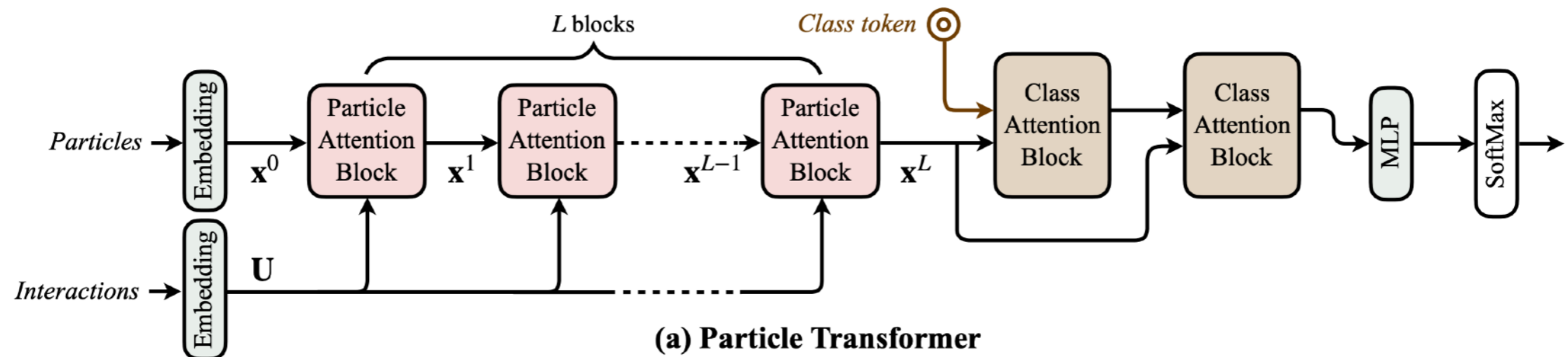


PART @ CMS FOR HF TAGGING



- The debut of ParT in CMS: Run3 HF tagging

ArXiv: 2202.03772



- ParT Architect for CMS HF tagging:

- 3 Particle Attention Blocks + 1 Class Attention Block
- GELU as activation for attention blocks, ReLU as activation for initial embedding + final MLP
- Number of “Heads” in MHA: 8
- Number of feature: 128 (for \mathbf{x}^l 's)
 - For initial embedding the MLPs are (128-512-128) and (64-64-64-8) for *Particles* and *Interactions* respectively.

PART @ CMS FOR HF TAGGING



- ParT @ CMS:
 - The input features are almost identical with DeepJet inputs
 - (+): 4 momenta of the constituents; (-): global features.
 - FLOPs: ~117M, #Params: ~1.4M
 - Inference time is similar with ParticleNet's
- Training details:
 - Dataset: ~65M jet from $t\bar{t}$ and QCD multi-jets
 - Ranger optimizer
 - LR initially set to $1e-3$ and will decay linearly after 70% of training to reach $1e-5$ in the end
 - Batch size: 512
 - Loss function: cross entropy
 - Model with best loss performance on validation set will be kept

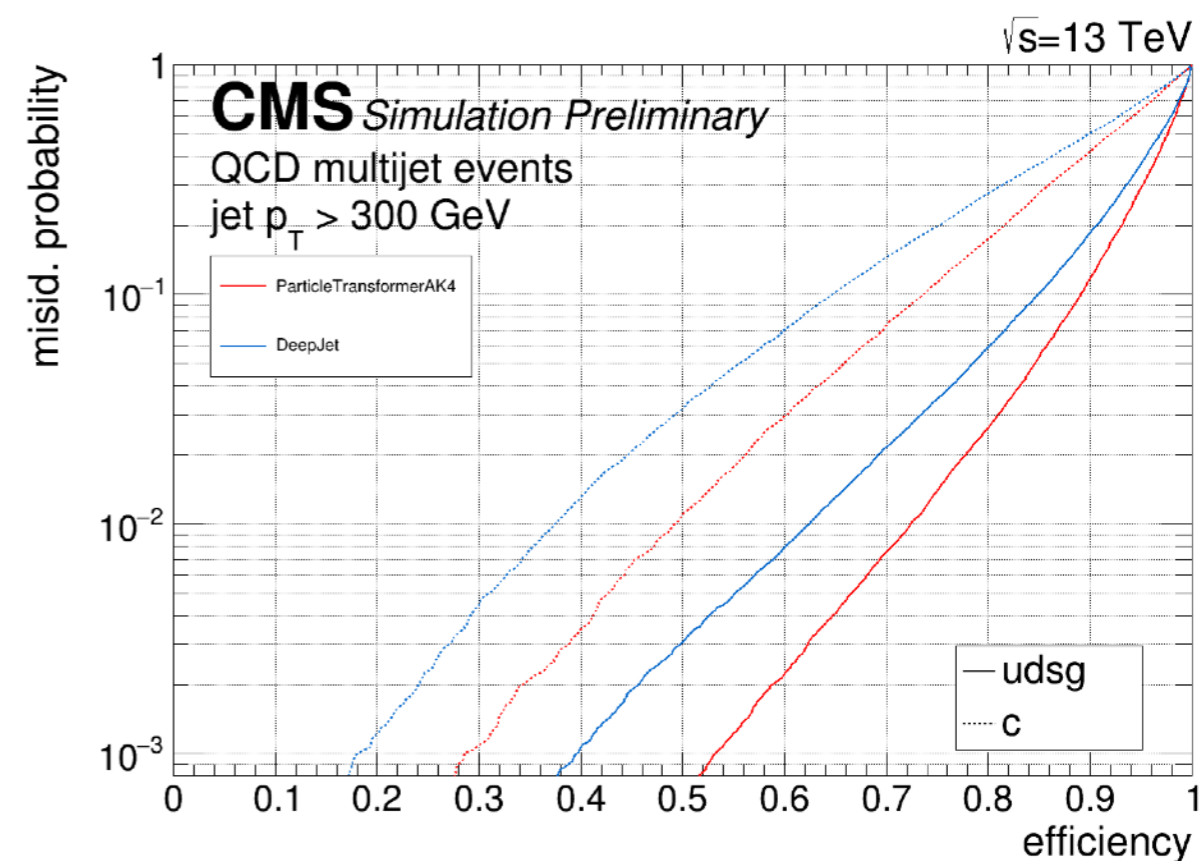
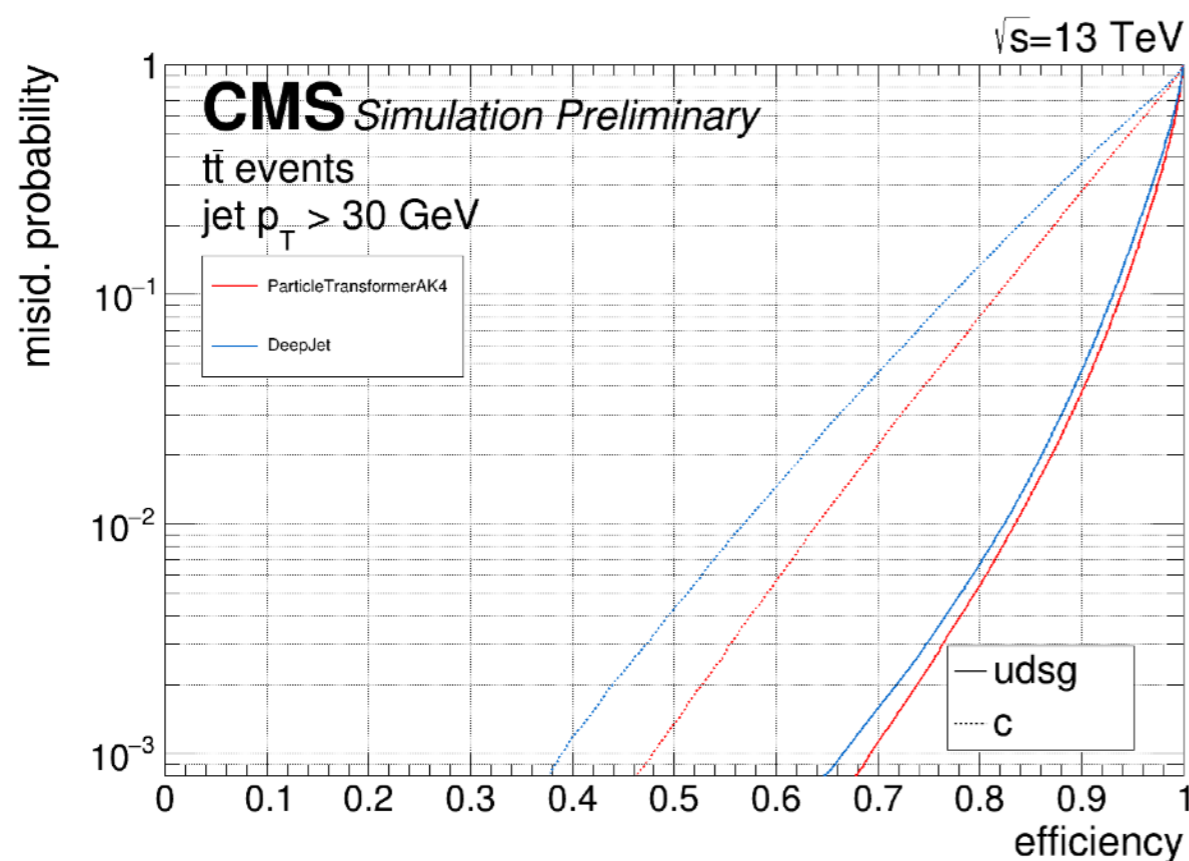
PART @ CMS FOR HF TAGGING



■ Performance: b tagging

DP-2022/050

Left: $t\bar{t}$ events, Right: QCD multijets



Performance on b tagging of ParT @ CMS:
Solid(dashed) ROC curves indicate the mistagging rates for
“udsg”(“c”) jet at given b tagging signal efficiency

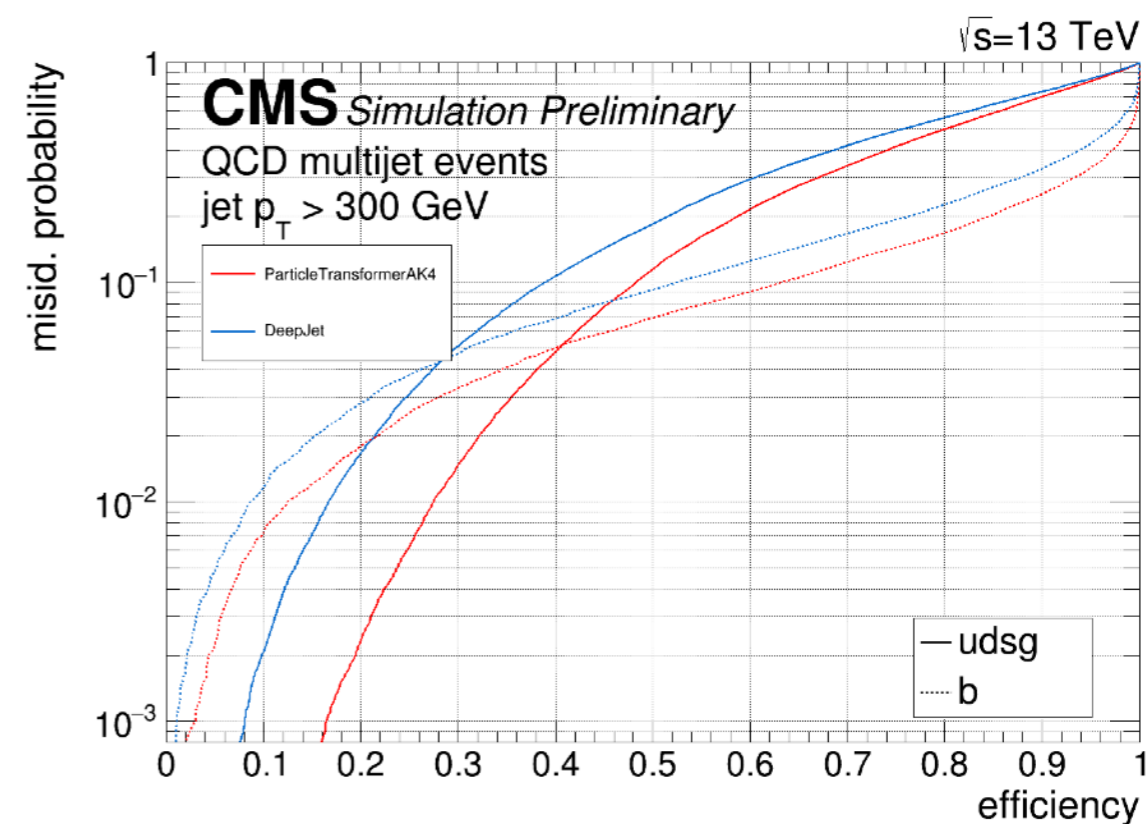
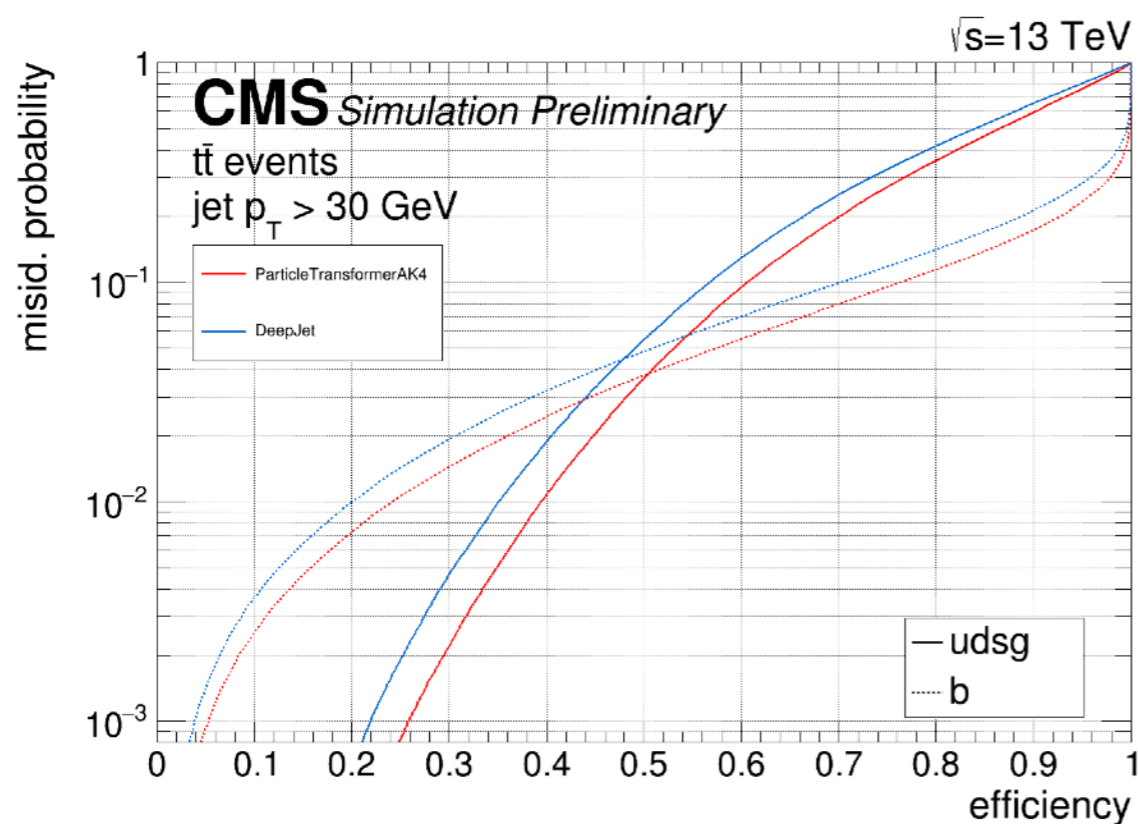
PART @ CMS FOR HF TAGGING



■ Performance: c tagging

DP-2022/050

Left: $t\bar{t}$ events, Right: QCD multijets



Performance on c tagging of ParT @ CMS:
Solid(dashed) ROC curves indicate the mistagging rates for
“udsg”(“b”) jet at given c tagging signal efficiency

PART @ CMS FOR HF TAGGING



- Summary:
 - Jet tagging is an important yet challenging task in HEP, among which flavor tagging is of vital importance
 - CMS has kept investigating HF tagging with more and more advanced models
 - For Run3, Particle Transformer (ParT) is considered as the official HF tagging algorithm
 - Significant improvement on HF tagging performance is observed from the comparison between ParT and previous CMS State-Of-The-Art model: DeepJet

BACK UP

PART VS PARTICLENET ON AK8 JET



Table 4. 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
ParT	0.861	2.14 M	340 M
ParT (plain)	0.849	2.13 M	260 M

Table 1. Jet tagging performance on the JETCLASS dataset. ParT is compared to PFN (Komiske et al., 2019b), P-CNN (CMS Collaboration, 2020b) and the state-of-the-art ParticleNet (Qu & Gouskos, 2020). For all the metrics, a higher value indicates better performance. The ParT architecture using plain MHAs instead of P-MHAs, labelled as ParT (plain), is also shown for comparison.

	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 _{50%}	Rej _{50%}	Rej _{50%}	Rej _{50%}	Rej _{99%}	Rej _{50%}	Rej _{99.5%}	Rej _{50%}	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

PART VS PARTICLENET ON AK8 JET



Table 3. Impacts of the training dataset size. Entries in bold correspond to the training using the full 100 M training dataset.

	All classes		$H \rightarrow b\bar{b}$	$H \rightarrow c\bar{c}$	$H \rightarrow gg$	$H \rightarrow 4q$	$H \rightarrow \ell\nu qq'$	$t \rightarrow bqq'$	$t \rightarrow b\ell\nu$	$W \rightarrow qq'$	$Z \rightarrow q\bar{q}$
	Accuracy	AUC	Rej _{50%}	Rej _{50%}	Rej _{50%}	Rej _{50%}	Rej _{99%}	Rej _{50%}	Rej _{99.5%}	Rej _{50%}	Rej _{50%}
ParticleNet (2 M)	0.828	0.9820	5540	1681	90	662	1654	4049	4673	260	215
ParticleNet (10 M)	0.837	0.9837	5848	2070	96	770	2350	5495	6803	307	253
ParticleNet (100 M)	0.844	0.9849	7634	2475	104	954	3339	10526	11173	347	283
ParT (2 M)	0.836	0.9834	5587	1982	93	761	1609	6061	4474	307	236
ParT (10 M)	0.850	0.9860	8734	3040	110	1274	3257	12579	8969	431	324
ParT (100 M)	0.861	0.9877	10638	4149	123	1864	5479	32787	15873	543	402

Table 6. Comparison between ParT and existing models on the quark-gluon tagging dataset. ParT refers to the model trained from scratch on this dataset. ParticleNet-f.t. and ParT-f.t. denote the corresponding models pre-trained on JETCLASS and fine-tuned 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), rPCN (Shimmin, 2021), and LorentzNet (Gong et al., 2022). The subscript “exp” and “full” distinguish models using partial or full particle identification information.

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

Table 5. Comparison between ParT and existing models on the top quark tagging dataset. ParT refers to the model trained from scratch on this dataset. ParticleNet-f.t. and ParT-f.t. denote the corresponding models pre-trained on JETCLASS and fine-tuned 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), rPCN (Shimmin, 2021), and LorentzNet (Gong et al., 2022).

	Accuracy	AUC	Rej _{50%}	Rej _{30%}
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
LorentzNet	0.942	0.9868	498 ± 18	2195 ± 173
ParT	0.940	0.9858	413 ± 16	1602 ± 81
ParticleNet-f.t.	0.942	0.9866	487 ± 9	1771 ± 80
ParT-f.t.	0.944	0.9877	691 ± 15	2766 ± 130