# CMS at the frontier of Deep Learning applications in Particle Physics

#### Jennifer Ngadiuba (Fermilab) on behalf of the CMS Collaboration

Aspen Winter Conference on Particle Physics 26-31 March, 2023 Aspen Center for Physics, Aspen (USA)





#### Intro: Machine Learning @ CMS

#### • Growing number of applications of machine learning techniques in CMS

- increasingly **more sophisticated** exploiting advances from non-scientific domains
- include also co-existing effort of expanding central software and hardware infrastructures to enable inference and training at scale
- **the CMS ML group** at the border between computing and physics — coordinating and overseeing technical and innovation aspects across the experiment

CMS RELEASES OPEN DATA FOR MACHINE LEARNING DEEP LEARNING TECHNIQUE IDENTIFIES COMPLEX DECAYS OF TAU LEPTONS MACHINING JETS



USING MACHINE LEARNING TO IMPROVE THE DETECTION OF NEW PHYSICS IN THE INTERACTIONS OF THE TOP QUARK

USING ARTIFICIAL NTELLIGENCE TO SEARCH FOR NEW EXOTIC PARTICLES

USING MACHINE LEARNING IMPROVE THE DETECTION O NEW PHYSICS IN THE INTERACTIONS OF THE TOP QUARK

**Check CMS news page!** 

#### Intro: Machine Learning @ CMS

#### **ML for classification:**

heavy jet tagging heavy flavour jet tagging exotic jets tau leptons event-level (also unsupervised)

#### **ML beyond classification:**

mass and energy regression background estimation detector faults detection reconstruction & simulation triggering

#### **Computing software & hardware for ML:**

optimized inference in central software for CPU/GPU GPU hardware on-site for software trigger system & grid sites more powerful chips in hardware trigger system & development of portable tools\* ML-friendly central data format (NanoAOD) and scalable processing tools (Coffea)

#### **DISCLAIMERS:**

a lot ongoing that cannot be covered today for more info contact <u>CMS ML group conveners</u>!

\* see Javier's talk

## **ML for classification**

## Heavy jet tagging

- Identification of jets arising from hadronization of boosted W/Z/H/top is a key task in LHC physics:
  - new physics searches, standard model measurements, higgs sector
  - **unique signature** from hadrons merging in single jet with substructure
  - exploit to suppress **overwhelming background from multijet** processes in most sensitive all-hadronic and semi-leptonic channels
- A topic of interest in both theory and experiment communities since ~ 30 years
- Recent years advancement in ML enabled more powerful algorithms





## Heavy jet tagging @ CMS

- A variety of ML algorithms developed and deployed in CMS since beginning of Run 2
- SOA: <u>ParticleNet</u> first graph-based tagger at LHC !
  - jet raw representation based on point cloud formed by its constituents
  - a dynamic graph CNN architecture used to fully exploit inherent permutation invariance of the representation ("EdgeConv" block)





EdgeConv

H. Ou et al.: Phys. Rev. D 101, 056019 (2020)



Softmax

(a)

- Also heavy flavour (b/c) jet tagging allows to access a variety of rare physics processes
- Unique signature due to sizeable lifetime of B/C hadrons [ $c\tau = O(mm)$ ] resulting in displaced tracks and secondary vertices
- We have built classifiers based on ML for HF jets at the LHC since Run 1
  - successful in reducing time to discovery!





#### **H→bb** discovery in VH production

- Also heavy flavour (b/c) jet tagging allows to access a variety of rare physics processes
- Unique signature due to sizeable lifetime of B/C hadrons [ $c\tau = O(mm)$ ] resulting in displaced tracks and secondary vertices
- We have built classifiers based on ML for HF jets at the LHC since Run 1
  - successful in reducing time to discovery!

CMS

• Application of more powerful ML architectures allowed recently setting the most stringent constraints on HH production!

138 fb<sup>-1</sup> (13 TeV)





CMS

- Brand new development in CMS uses particle-based transformers *learn which neighbour particles are relevant through attention mechanism* 
  - input embedding of both single particle and pair-wise features information
  - the pair-wise features encode physics principles → modifiers of standard dot-product attention weights in **Particle Attention Block**



- Recent development in CMS uses particle-based transformers *learn which neighbour* particles are relevant through attention mechanism
  - input embedding of both single particle and pair-wise features information
  - the pair-wise features encode physics principles → modifiers of standard dot-product attention weights in **Particle Attention Block**
  - **Class Attention Block:** compute the attention between a global class and all the particles via standard Multi-Head Attention (MHA)



#### **Class Attention Block**

 $10^{-3}$ 

- Comparison with current default CMS algo <u>Deeplet</u> based on 1D CNN architecture on per-particle plus b-tagging features
- Known challenges in tagging high p<sub>T</sub> jets
   — worse reconstruction of key features like
   impact parameter and secondary vertices
- Introduction of pair-wise features mitigate performance loss thanks to increased expressivity

nisid. probability **CMS** Simulation Preliminary QCD multijet events jet p\_ > 300 GeV ParticleTransformerAK 10-DeepJet 10<sup>-2</sup> udsg С 10<sup>-3</sup> 0.2 0.3 0.4 0.5 0.9 0 0.1 0.6 0.70.8 efficiency √s=13 TeV nisid. probability CMS Simulation Preliminary tī events jet p\_ > 30 GeV ParticleTransformerAk 10- $10^{-2}$ udsg

Finalising integration in CMS software... Looking forward for deployment in Run 3 analyses — first attention-based tagger at LHC!

С

0.8

0.9

efficiency

0.6

0.7

0.5

0.4

0.3

0.2

√s=13 TeV

# A Control of the second second

 Recent exciting CMS results reporting observation of the rare 4 top quarks production process at 4σ for combination of multiple final states [CMS-TOP-21-005 — Submitted to PLB]



**OBSERVATION JUST REPORTED LAST WEEK AT MORIOND!** 

- Recent exciting CMS results reporting evidence for the rare 4 top quarks production process at 4σ for combination of multiple final states [CMS-TOP-21-005 — Submitted to PLB]
- Analysis of data in the challenging **all-hadronic final state performed for the first time**



- Start with defining 5 CRs as in "extended" ABCD method (i.e. number of CRs > 3)
- Main idea: for each CR learn transformation  $\tau: \mathscr{P} \to \mathscr{P}'$  using NAF as universal approximator of bijective transformation with sigmoidal function  $\sigma$

 $\mathscr{P}'(\overrightarrow{x'} \mid \overrightarrow{c'}) = \tau(\overrightarrow{x'}; \overrightarrow{x} \mid \overrightarrow{c'}; \overrightarrow{c}) \otimes \mathscr{P}(\overrightarrow{x} \mid \overrightarrow{c})$ 

- $\mathscr{P}': t\bar{t}$ +QCD background = data ( $t\bar{t}t\bar{t}$  + minor MC backgrounds)
- $\mathcal{P}: t\overline{t}$  simulation

•  $\vec{x}$  : input variables **BDT and HT shapes** 

- $\overrightarrow{y}$  : output variables
- $\vec{c}$  : conditions, including control variables N(Jets) and N(bJets) plus preceding inputs
- loss function: maximum-mean-discrepancy
- activation function: sigmoid





- After training in the 5 CRs, the transformation between MC and data is applied to *tt* simulation in the SR
   → morphed to predict the shape of the *tt*+QCD multijet background in the SR
- Extensive closure checks in validation region: identical to SR but with N(jets) = 8
- Uncertainties derived from discrepancies in the VR and applied to corresponding SR





[CMS-TOP-21-005 Submitted to PLB]

## ML for detector health monitoring

- The online data quality monitoring (DQM) of CMS aims at **detecting faults almost** in real-time O(s)
  - initiate reaction from detector experts to promptly identify and fix problems
  - minimize downtime, maximize time for physics!
- **Standard workflow:** set of histograms that are populated with a set of events seen by the detector
- Challenge: anomalies come in all shapes and sizes
   → *impossible to anticipate all possible failure modes*
- Why ML can help?
  - provides robust anomaly detection and localization (via semi/unsupervised learning)
  - eliminate the need for hand coded rules for every possible fault
  - automated adaptivity to changing running conditions and experimental setup

#### **Ex: CMS ECAL barrel**



GREEN = good
RED = bad
BROWN = known problem
YELLOW = no data (which may or may not be
problematic – depends on context).

#### ML for detector health monitoring

- First development using an autoencoder in production since 2017 for Drift Tubes subsystem in muon detectors → first example of ML for online DQM at LHC!
- More parallel efforts now ongoing mostly based on autoencoders **under or planned commissioning in Run 3 data taking**:
  - Resistive Plate Chambers subsystem of muon detectors [ACAT '22]
  - Electromagnetic and Hadronic Calorimeters subsystems [CMS-DP-2022-043, IML workshop]
  - Pixel Silicon Tracker subsystem [CMS-DP-2022-013]

#### **Electromagnetic Calorimeter example:**

- Training set of certified GOOD occupancy map images
- Test set of both synthetic and real Run 2 & 3 anomalies
- Autoencoder model: ResNet for both encoding and decoding
- Preprocessing to take into account variable spatial response and time-dependent nature of faults

#### Metric: False Discovery Rate (FDR) at 99% anomaly detection (i.e. what fractions of shifter calls will be false alarm?)

FDR at 99% anomaly detection								
Scenario	Missing Supermodule	Zero Occupancy Tower			Hot Tower (10% hot for EB, 20% hot for EE)			
	Barrel	Barrel	EE+	EE-	Barrel	EE+	EE-	
No correction	3.6 %	51 %	86 %	87 %	2.8 %	0.01 %	<1/30k	
After <i>spatial</i> correction	3.1 %	49 %	13 %	14 %	2.9 %	0.06 %	0.05 %	
After <i>spatial and time</i> correction	0.13 %	4.1 %	5.6 %	6.3 %	<1/10k	<1/30k	<1/30k	
	с	Mostly due t ontaminating (Se	to <b>actual and</b> certified GO ee Backup).	o <b>malies</b> OD data.	~10 <sup>4</sup> is the size of the validation set			

#### ML for detector health monitoring

- First development using an autoencoder in production since 2017 for Drift Tubes subsystem in muon detectors → first example of ML for online DQM at LHC!
- More parallel efforts now ongoing mostly based on autoencoders **under or planned commissioning in Run 3 data taking**:
  - Resistive Plate Chambers subsystem of muon detectors [ACAT '22]
  - Electromagnetic and Hadronic Calorimeters subsystems [CMS-DP-2022-043, IML workshop]
  - Pixel Silicon Tracker subsystem [CMS-DP-2022-013]

#### **Electromagnetic Calorimeter example:**

- Training set of certified GOOD occupancy map images
- Test set of both synthetic and real Run 2 & 3 anomalies
- Autoencoder model: ResNet for both encoding and decoding
- Preprocessing to take into account variable spatial response and time-dependent nature of faults

#### Successful commissioning tests from last year run!

#### ML Quality plot from ECAL Online DQM during a Run3 run



## **Computing aspects**

ompact Muon Solenoid

## ML for particle flow reconstruction

- General effort in the experimental community to **replace standard rule-based reconstruction algorithms with more computationally efficient and scalable ML models** in view of future increased challenges
- One example in CMS is the **Particle flow (PF) algorithm**: it combines information from all subdetector to reconstruct particles and thus improve the resolution
  - ex. track + hadronic energy= charged hadron
  - ex. no track + electromagnetic energy= photon
- It starts from calorimeter clusters & tracks and outputs particle candidates
   → replace with ML model



## ML for particle flow reconstruction

- Model based on **dynamic graph CNN** generating on the fly multiple internal kNN graphs based on embedded features
- **Per-particle loss function** to simultaneously perform multi-classification and regression tasks

$$||Y - Y'|| \equiv \sum_{j \in \text{event}} L(y_j, y'_j),$$
$$L(y_j, y'_j) \equiv \text{CLS}(c_j, c'_j) + \alpha \text{REG}(p_j, p'_j)$$



- target (predicted) particle, - no target (predicted) particle

- First version target baseline particle flow reconstruction → does not allow to do better than that
- Second version target generator-level information → can one improve response versus baseline?

#### ML for particle flow reconstruction

Run 3 (14 TeV)

**CMS** Simulation Preliminary





Eur. Phys. J. C (2021) 81: 381 ACAT '21 ACAT '22

- Overall, a good particle-level agreement is observed between PF and MLPF algos
- Missing transverse energy mismodelling under investigation
- Using gen-level information gives similar performance

#### **Computing infrastructure**

- A crucial metric when developing ML models is the computational efficiency
- Cannot exploit full ML power without extensive work to support and optimize ML inference in the CMS software
  - continuous improvements in supporting direct inference on local GPU/CPU
  - promising effort in supporting inference as a service on cloud resources



#### XGBoost TensorFlow+Keras PyTorch scikit-learn Python in your env.) Training engine KGB learn torch.onnx default default tf2onnx saved\_model\_cli sklearn-onnx Inference engine (CMSSW) XGB XLA/AOT ONNX XGBoost TensorFlow PyTorch CPU CPU GPU GPU = in development

**Current CMS software support for ML** 

First results from MLPF in standalone setup promising: TENSORFLOW implementation and exported to ONNX

> See <u>CHEP 2023 this May</u> for full results in CMS software stack!

## ML for triggering



Thanks to advances in computing tools/infrastructures **DL possible in CMS HLT since Run 2 & SOA models like graph-based tagger possible in Run 3 data taking!** 

## ML for triggering



Effort also ongoing for Phase 2 developments where increased L1T system capabilities will allow to do even more!

Thanks to advances in computing tools/infrastructures **DL possible in CMS HLT since Run 2 & SOA models like graph-based tagger possible in Run 3 data taking!** 

#### Conclusions

**Continuous innovation in CMS on both algorithms and system side** exploiting modern Deep Learning techniques

Push developments today and use acquired expertise in the future at HL-LHC

**ML** for classification **ML beyond classification Computing software & hardware for ML** The CMS ML group L2 Group: Machine Learning Goal: enable, support, guide and foster ML developments in CMS computing, physics objects and physics analyses groups Gregor Kasiecka Jennifer Ngadiuba

Three subgroups coordinating multiple efforts to achieve the goal!



<u>cms-conveners-ml@cern.ch</u>

# Thank you!

compact Muon Solenoid