





# Boosted jet tagging at the CMS experiment

Donato Troiano<sup>1,2</sup>

October 8, 2024

CMS Italia 2024

1 INFN of Bari

2 University of Bari

### Boosted jets tagging

Hadronic decay products of highly boosted particles are collimated.

- Particles merged in one anti- $k_T$  jet (radius R = 0.8  $\rightarrow$  AK8).

Identifying the particle generating the jet is a suitable task for ML:

- large number of particles;
- diverse input variables for each particle.



### **Evolution of boosted jets tagging at CMS**

early Run 2		sta	te of art for Run 3 CMS
double-b	DeepAK8-MD	DeepDoubleX	ParticleNet-MD
			time
Model: Boosted decision Tree	1D Convolutional Neural Network (CNN)	1D CNN + Recursive NN	Dynamic Graph CNN (DGCNN)
Inputs: tracks and secondary vertices (SVs)	particle flow (PF) candidates and SVs	PF candidates and SVs	PF candidates, SVs and tracks
Outputs: • $H \rightarrow b\overline{b} vs QCD$	<ul> <li>X → bb̄ vs QCD</li> <li>X → cc̄ vs QCD</li> </ul>	• $X \rightarrow b\overline{b} vs QCD$ • $X \rightarrow c\overline{c} vs QCD$ • $X \rightarrow b\overline{b} vs X \rightarrow c\overline{c}$	<ul> <li>p(X→bb̄), p(X→cc̄),</li> <li>p(X→light quarks),</li> <li>p(X→ττ̄), p(QCD)</li> </ul>

### ParticleNet-MD

### Architecture

- Graph based architecture describing the jet as a particle cloud (unordered sample).

### **EdgeConv block:**

- NN module part of the ParticleNet architecture;
- New features vector associated to each jet constituent and based on the features of the k-nearest neighbors.





### ParticleNet-MD new training performance



5

### My work: First Run3 ParticleNet-MD validation

- Validation done on 2022 data events containing highly boosted  $Z \rightarrow b\overline{b}$  + jets (<u>CMS-DP-2024-055</u>).
- Likelihood fit on the soft-drop mass  $(m_{SD})$  to match data and SM prediction  $(Z \rightarrow q\overline{q}, W \rightarrow qq')$ and QCD).
  - m<sub>SD</sub>: jet mass after removing soft and wide-angle radiation.
- Data-driven estimation of QCD multijet background:
  - average of the fits of the Z-candidate m<sub>SD</sub> distributions in nine mass sidebands.



### ParticleNet-MD categorization

- Events categorized in five region of PNet-MD<sub>bbvsQCD</sub>.
  - 20% of signal events in each PNet-MD<sub>bbvsQCD</sub> region.
  - Lowest score region (<0.641) almost completely populated by QCD events and not included in the fit. \_\_\_\_\_\_CMS simulation Preliminary \_\_\_\_\_\_(13.6 TeV)
  - Independent Likelihood fit performed in each of the four highest score regions.



### Boosted $Z \rightarrow b\overline{b}$ validation

- Good data-prediction agreement.
- Clear Z-peak in the data distribution.
  - Good signal discrimination from background.

#### CMS-DP-2024-055



### Future plans for development of PNet at HLT

- ParticleNet-MD in HLT paths not optimized since Run 2.
- My plan: ParticleNet-MD optimization at HLT
  - ☑ ParticleNet-MD input variables implemented in the <u>BTV</u> <u>ntuplizer</u>.
  - Produce ntuples with flat p<sub>t</sub> and mass distributions for mass decorrelation (MC production ongoing).
  - ☑ Train on HLT features using BTV <u>framework b-hive</u> (targeting December for first HLT optimization).





### Conclusions

- Over the years, CMS has developed algorithms for boosted objects with increasing efficiency and performance.

- ParticleNet-MD best performing tagger for boosted jets, now optimized for Run 3.

- Validation on  $Z \rightarrow b\overline{b}$ -like events in 2022 data.
  - Stable performance, with some improvement with respect to Run 2.
  - Good data-prediction agreement.
  - Clear Z-peak in data distribution at high score.
- Currently working on boosted jet tagging at HLT.



# Backup

### Jets as particle clouds: ParticleNet-MD

# **ParticleNet-MD** state-of-art for CMS boosted jet tagging.

- Graph based architecture describing the jet as a particle cloud (unordered sample).

### EdgeConv block:

- NN module part of the ParticleNet architecture;
- New features vector associated to each jet constituent and based on the features of the *k*-nearest neighbors.

### Mass decorrelation:

 Trained on Monte Carlo (MC) simulations containing boosted resonances (X), decaying in taus and quarks, with a flat distributions in both of p<sub>t</sub> and mass, as the signal sample, and the QCD multijet sample (reweighted to yield flat distributions) as the background sample.



### ParticleNet architecture



### ParticleNet architecture



# **ROC** curve $b\overline{b}$ tagging performances (Run 2)



### ROC curve $c\overline{c}$ tagging performances (Run 2)



# Boosted $Z \rightarrow b\overline{b}$ event selection

- **High-Level-Trigger paths :** PFHT1050, PFJet500, AK8PFJet500, AK8PFJet400\_TrimMass30, AK8PFJet420\_TrimMass30, AK8PFHT800\_TrimMass50

- Leading-p<sub>T</sub> AK8 jet (Z-boson candidate):  $p_T > 450$  GeV and  $|\eta| < 2.4$
- Sub-leading-p<sub>T</sub> AK8 jet:  $p_T > 200$  GeV and  $|\eta| < 2.4$
- Veto events with at least one electron or muon with  $p_T$  > 20 GeV,  $|\eta|$  < 2.4, and satisfying the loosest identification and isolation working point
- Veto events with a b-tagged AK4 jet having  $p_T > 30$  GeV,  $|\eta| < 2.4$  and a distance  $\Delta R$  from the leading AK8 jet greater than 0.8
  - The DeepJet medium working point is used to tag AK4 jets as originating from bquark

# **PNet-MD**<sub>bbvsQCD</sub> categorization



# **PNet-MD**<sub>bbvsQCD</sub> validation: Likelihood fit

- The likelihood fit is performed within the signal mass window in the four highest score regions defined in slide 7.

- The parameters of interest of the fit (the MC  $Z \rightarrow q\bar{q}$  and  $W \rightarrow q\bar{q}$ ' normalization factors) are obtained independently in each score region.
- The background from QCD multijet events is estimated using the average of the fits of the Z-candidate  $m_{\rm SD}$  distributions outside one of nine alternative mass windows.
- The following uncertainties are considered:
  - uncertainty on QCD estimate due to the Z-candidate m<sub>SD</sub> distribution fit functions;
  - uncertainty on QCD estimate due to the use of the nine mass windows;
  - statistical uncertainties for MC Z  $\rightarrow q\bar{q}$  and W  $\rightarrow q\bar{q'}$ ;
  - jet energy scale corrections for MC Z  $\rightarrow q\bar{q}$  and W  $\rightarrow q\bar{q'}$ ;
  - the luminosity uncertainty.

- All the uncertainties, with the exception of the luminosity one, are assumed uncorrelated in the different score regions.

### **QCD** evaluation

- QCD estimation Data-driven technique
  - QCD estimated by fitting data m<sub>sD</sub> distribution ([50-150] GeV) outside [70;126], [62;134], [78;118], [54;142], [86,110], [62;126], [78;126], [70,134] and [70;118] GeV windows.
    - QCD obtained as the average of the fits.
  - 2. Chebyshev polynomial functions have been used to fit Data.
    - Order established by means of the Fisher test (CL at 5%).
- QCD distributions assumed with the following uncertainties:
  - fit: symmetric error in each m<sub>SD</sub> bin equal to the average of the fit function errors associated to each bin.
  - window: symmetric error in each  $m_{SD}$  bin equal to the root mean square deviation of the  $m_{SD}$  bin values obtained from the different fit functions.

### **QCD** evaluation



# **PNet-MD**<sub>bbvsQCD</sub> score post-fit distribution



- PNet-MD<sub>bbvsQCD</sub> distribution of the stack of  $W \rightarrow q\overline{q'}$  and  $Z \rightarrow q\overline{q}$  (split on the basis of its decay products) yield, and the data, once the QCD contribution is subtracted (Data - QCD), of the leading-p<sub>t</sub> AK8.

- Bottom pad: ratio between Data-QCD and the stack.

- Good data-prediction agreement.

- The amount of  $W \rightarrow q\overline{q'}$ and  $Z \rightarrow q\overline{q}$  events, with q (light) or c quark, decreases increasing the score value.

### Signal fraction

	Z→qq (q=b,c,u,d,s) percentage		
	Z→bb	Z→cc	$Z \rightarrow uu/dd/ss$
4 <sup>th</sup> highest score region	48.8%	35.5%	15.7%
3 <sup>rd</sup> highest score region	73.9%	20.5%	5.6%
2 <sup>nd</sup> highest score region	87.7%	10.9%	1.4%
highest score region	96.6%	3.18%	0.26%

- MC mis-identified  $Z \rightarrow b\overline{b}$  percentage less than 4% in highest score region.
- 20% of  $Z \rightarrow b\overline{b}$  signal events lost moving from a score region to the next.

### **ParticleNet-MD validation results**

	r <sub>Z</sub>
$0.641 < PNet-MD_{bbvsQCD} \le 0.875$	$0.9 \pm 0.5$ (stat) $\pm 0.3$ (syst)
$0.875 < PNet-MD_{bbvsQCD} \le 0.957$	1.37 ± 0.26 (stat) ± 0.21 (syst)
$0.957 < PNet-MD_{bbvsQCD} \le 0.988$	1.25 ± 0.15 (stat) ± 0.12 (syst)
0.988 < PNet-MD <sub>bbvsQCD</sub> ≤ 1	$1.01 \pm 0.07$ (stat) $\pm 0.07$ (syst)

- $Z \rightarrow q\overline{q}$  (q = u, d, s, c, b) normalization factors (r<sub>Z</sub>) with the corresponding error (split in statistical and systematic errors) in the four highest score regions.
- Normalization factors compatible with unity within uncertainties.

### 4<sup>th</sup> highest score region



### 3<sup>rd</sup> highest score region



### 2<sup>nd</sup> highest score region



### Highest score region



### ParticleTranfrormer

- ParticleTransformer architecture shows superior performance than ParticleNet [arXiv:2202.03772].

- Next goal: extend tagging to a large set of final states including top-tagging (3 prong),  $X \rightarrow VV$  with hadronic and leptonic decays, taus etc.
- ParticleTransformer used already in HH $\rightarrow$ VV $\rightarrow$ 4b analysis [CMS-PAS-HIG-23-012]





### New Jet tagging developments (Jet charge tagger)

- DGCNN based on the ParticleNet architecture predicting the charge of the AK8 jet.
  - Discriminate W<sup>+</sup> and W<sup>-</sup> bosons from Z boson.
- Observed good data-MC agreement selecting semileptonic tt events.



# Jet charge tagger

- DGCNN based on the ParticleNet architecture predicting the charge of the AK8 jet.
  - Discriminate W<sup>+</sup> and W<sup>-</sup> bosons from Z boson.
- Training samples:
  - semileptonic tt MC simulation is used to get a sample with W<sup>+</sup> and W<sup>-</sup> jets;
  - Z+jets MC simulation.
- Validation done on a region enriched of semileptonic  $t\bar{t}$  events.
  - Good data-MC agreement.



### Jet charge tagger

semileptonic tt enriched region output score: W<sup>+</sup>

#### semileptonic tt enriched region output score: Z (not expected Z bosons)



# New Jet tagging developments (HOTVR)

- Heavy Object Tagger with Variable Radius.
  - $R = 600 / p_T$  (R min 0.1).
  - Useful for 4 top final states in the intermediate region where the top quark is neither resolved ( $p_t < 200 \text{ GeV}$ ) nor boosted ( $p_t > 800 \text{ GeV}$ ).
  - For Run 2 top tagging was performed using a cut-based approach.
  - Developed a BDT, using as inputs the cut variables of the Run 2 analysis, to distinguish top quarks from QCD.
    - Observed improved with respect to the top tagging cut-based.



# Top tagging with variable sized jets

#### Heavy Object Tagger with Variable Radius

- $R = 600 / p_t$  (R min 0.1, R max 1.5).
- Useful for 4 top final states where the top quark is not completely boosted (200 < p<sub>t</sub> < 800 GeV).</li>
- Efficiency as the ratio between the generated top quarks matching a reconstructed jet within ΔR and all the generated top quarks.

# Developed a BDT to distinguish top quarks from QCD:

- Training on QCD multijet and the ttZ to simulate the background and the signal, respectively;
- Tested on a Z+jets enriched selection
  - Two opposite sign leptons (80 < m<sub>ℓℓ</sub> < 101 GeV) + ≥1 HOTVR

#### <u>CMS-DP-2024-038</u>



### Top tagging with variable sized jets

