



Performance of heavy-flavour jet identification in boosted topologies in CMS 13 TeV data

Congqiao Li (Peking University)

on behalf of the CMS Collaboration, based on CMS-PAS-BTV-22-001

BOOST 2023 · Berkeley 31 July, 2023

Boosted heavy-flavour jet tagging

- ➔ Boosted topologies are crucial in LHC physics program
 - ★ highly Lorentz-boosted resonance
 → decay products are collimated



- → Tagging the heavy resonance X to bb/cc̄ decay is an important technique for Higgs and BSM resonance search
 - algorithms developed in CMS during Run 2:
 - ParticleNet-MD, DeepDoubleX, DeepAK8-MD, double-b
 - calibrating the algorithm from data is a necessary step for physics measurements

A simulated $Z(v\bar{v})H(b\bar{b})$ events in CMS



H→bb

A simulated $Z(v\bar{v})H(b\bar{b})$ events in CMS

Boosted heavy-flavour jet tagging

- ➔ Boosted topologies are crucial in LHC physics program
 - ♦ highly Lorentz-boosted resonance
 → decay products are collimated



- → Tagging the heavy resonance X to bb/cc̄ decay is an important technique for Wigge and DCM resonance
 - Aim of this talk
 - Introduce the X→bb/cc̄ taggers developed in CMS Run 2
 - Highlight their performance on simulated jets
 - Summarise and benchmark three calibration methods



algorith

calibrati

step for

Parti

doub

 $\mathbf{\mathbf{x}}$

X→bb/cc̄ taggers: ParticleNet-MD



See full architecture in backup

X→bb/cc̄ taggers: ParticleNet-MD



$$bbvsQCD = \frac{p(X \to b\overline{b})}{p(X \to b\overline{b}) + p(QCD)}$$
$$ccvsQCD = \frac{p(X \to c\overline{c})}{p(X \to c\overline{c}) + p(QCD)}$$

X→bb/cc̄ taggers: ParticleNet-MD



X→bb/cc̄ taggers: DeepDoubleX



- → DeepDoubleX:
 - architecture: 1D CNN + RNN (with gated recurrent units, GRU)
 - inputs: PF candidates + SV + jet global features
 - irrelevant features pruned using layer-wise relevance propagation

→ Train three binary networks:

- X→bb vs. QCD (BvL), X→cc vs. QCD (CvL), X→bb vs. X→cc (CvB)
- studied tagger: BvL and CvL



Mass decorrelation achieved by training on signal and background jets with uniform distribution



CMS-DP-2022-041



X→bb/cc̄ taggers: DeepAK8-MD & double-b



- → DeepAK8-MD:
 - architecture: 1D CNN
 - inputs: same with ParticleNet-MD
 - classes: t/H/Z/W with dedicated final states
 - mass decorrelation: achieved by adversarial training



Define discriminants:

$$bbvsQCD = \frac{p(H \to bb) + p(Z \to bb)}{p(H \to b\overline{b}) + p(Z \to b\overline{b}) + p(QCD)}$$
$$ccvsQCD = \frac{p(H \to c\overline{c}) + p(Z \to c\overline{c})}{p(H \to c\overline{c}) + p(Z \to c\overline{c}) + p(QCD)}$$

X→bb/cc̄ taggers: DeepAK8-MD & double-b



- → DeepAK8-MD:
 - architecture: 1D CNN
 - inputs: same with ParticleNet-MD
 - classes: t/H/Z/W with dedicated final states
 - mass decorrelation: achieved by adversarial training



Define discriminants:

$$bbvsQCD = \frac{p(H \to b\overline{b}) + p(Z \to b\overline{b})}{p(H \to b\overline{b}) + p(Z \to b\overline{b}) + p(QCD)}$$
$$ccvsQCD = \frac{p(H \to c\overline{c}) + p(Z \to c\overline{c})}{p(H \to c\overline{c}) + p(Z \to c\overline{c}) + p(QCD)}$$

→ double-b:

- a BDT for distinguishing $H \rightarrow b\bar{b}$ and QCD jets
- inputs: jet variables constructed from tracks and SVs
- ★ mass decorrelation: by choosing input with weak correlation on p_T and mass

Tagger discriminants spectra



CMS-PAS-BTV-22-001

- Signal: H→bb,
 H→cc̄ jets from ggH
- BKG: QCD jets
 - QCD bb (cc): number of ghostmatched b and c hadrons:

N_b=2 (N_b=0 & N_c=2)

More taggers in backup

CMS-PAS-BTV-22-001

Tagger discriminants spectra

Conggiao Li (Peking University)



BOOST 2023

LP

low

purity

80%

50%

Performance on simulation

CMS-PAS-BTV-22-001

Low p_T: (450, 600) GeV, similar performance for high p_T (see <u>backup</u>)



ROC comparison on $H \rightarrow b\bar{b}$ (c \bar{c}) jets vs. inclusive QCD

Performance on simulation

CMS-PAS-BTV-22-001

Low p_T: (450, 600) GeV, similar performance for high p_T (see <u>backup</u>)



7

WEW

Three calibration methods

- sfBDT method
- μ-tagged method
- boosted Z method

sfBDT method: key concept



The general idea is to select **g**→**bb**/**cc**̄ jets from QCD events as the proxy of H→b**b**/**cc**̄

sfBDT method: key concept



 $\overline{b}(\overline{c})$

data

(proxy jets)

selected jets from QCD multijet events (requiring 2 SVs)

then used as SF for $H \rightarrow b\bar{b}/c\bar{c}$ jets

The general idea is to select $g \rightarrow bb/c\bar{c}$ jets from QCD events as the proxy of $H \rightarrow b\bar{b}/c\bar{c}$

Problem:

Preselected jets from QCD events, enriched in $\mathbf{g} \rightarrow \mathbf{b} \mathbf{b} / \mathbf{c} \mathbf{\bar{c}}$, still do not closely resemble **H→bb/cc** jets.

 \rightarrow Further selection needed

g

sfBDT method: key concept



The general idea is to select **g→bb/cc** jets from QCD events as the proxy of **H→bb/cc**

Problem:

Preselected jets from QCD events, enriched in g→bb/cc̄, still do not closely resemble H→bb/cc̄ jets.

→ Further selection needed

Solution:

Train a BDT variable which selects more signal-like jets from QCD



Select two regions using a generator-level variable

 one resembles more signal H→bb/cc jets, one does not



Select two regions using a generator-level variable

 one resembles more signal H→bb/cc̄ jets, one does not

Train a BDT (named sfBDT) against each other, using reconstructed-level variables as input (jet N-subjettiness, track/SV features)





Design 1: Published in <u>CMS-DP-2022-005</u>



- QCD jets are likely to be contaminated with extra gluons
- Define variables from partons

$$\kappa_g = \frac{\sum_{i \in \{g\}} p_{\mathrm{T},i}}{\sum_{i \in \{g,q\}} p_{\mathrm{T},i}}$$

Select two regions using a generator-level variable

 one resembles more signal H→bb/cc̄ jets, one does not

Train a BDT (named sfBDT) against each other, using reconstructed-level variables as input (jet N-subjettiness, track/SV features)



to evaluate the multi-prongness of a QCD jet smaller $\tau_{31} \rightarrow$ less complexity in prongness \rightarrow signal-like



Key logic:

- sfBDT as a handle to "tune" proxy—signal similarity
- evaluate the SF's dependence on sfBDT selection extract an extra uncertainty from it



Key logic:

- sfBDT as a handle to "tune" proxy—signal similarity
- evaluate the SF's dependence on sfBDT selection extract an extra uncertainty from it

Design *a systematic* workflow that applies to all taggers

- step1: transform the tagger to a uniform shape
- step2: define a set of sfBDT curves on the sfBDT-tagger plane, ٠ that match the proxy tagger shape to the signal (but with different selection efficiency)

CMS Simulation Preliminary

 MC (b) proxy (1st sfBDT curve) MC (b) proxy (2nd sfBDT curve)

MC (b) proxy (3rd sfBDT curve) MC (b) proxy (4th sfBDT curve)

MC (b) proxy (5th sfBDT curve) MC (b) proxy (6th sfBDT curve)

MC (b) proxy (7th sfBDT curve) MC (b) proxy (8th sfBDT curve) MC (b) proxy (9th sfBDT curve)

0.4

0.6

Signal and proxy jets

sfBDT method

---- H → bb signal

0.2

A.U.

0.5

0.4

0.3

0.2

0.1

0.0





Key logic:

- sfBDT as a handle to "tune" proxy—signal similarity
- evaluate the SF's dependence on sfBDT selection extract an extra uncertainty from it

Design *a systematic* workflow that applies to all taggers

- step1: transform the tagger to a uniform shape
- step2: define a set of sfBDT curves on the sfBDT-tagger plane, that match the proxy tagger shape to the signal (but with different selection efficiency)

CMS Simulation Preliminary

MC (b) proxy (1st sfBDT curve)

MC (b) proxy (2nd sfBDT curve)

MC (b) proxy (3rd sfBDT curve) MC (b) proxy (4th sfBDT curve)

MC (b) proxy (5th sfBDT curve)

MC (b) proxy (6th sfBDT curve)

MC (b) proxy (7th sfBDT curve) MC (b) proxy (8th sfBDT curve)

MC (b) proxy (9th sfBDT curve)

0.4

0.6

Signal and proxy jets

sfBDT method

---- H → bb signal

0.2

A.U.

0.5

0.3

0.2

0.1

0.85

- step3: select on those curves to derive nominal SFs; select different curves for the "pass" and "fail" tagger region to modify the proxy—signal tagger shape similarity
- step4: from the deviation of SFs, we obtain an external uncertainty assigned to the final SF







Data and MC distribution on transformed tagger discriminant, applying the central sfBDT selection

More taggers in **backup**

sfBDT method: fit



µ-tagged method: concept

- → Select proxy jet:
 - - b/c flavour content enriched
 - incorporate online trigger selection: requiring a soft muon in AK4/AK8 jets
 - largely orthogonal to the phase space explored in the sfBDT method
 - * τ₂₁ variable as a tune to modify signal—proxy similarity
 - for main fit: apply $\tau_{21} < 0.3$





µ-tagged method: concept

- → Select proxy jet:
 - - b/c flavour content enriched
 - incorporate online trigger selection: requiring a soft muon in AK4/AK8 jets
 - largely orthogonal to the phase space explored in the sfBDT method
 - - for main fit: apply $\tau_{21} < 0.3$





µ-tagged method: concept

- → Select proxy jet:
 - - b/c flavour content enriched
 - incorporate online trigger selection: requiring a soft muon in AK4/AK8 jets
 - largely orthogonal to the phase space explored in the sfBDT method
 - * τ₂₁ variable as a tune to modify signal—proxy similarity
 - for main fit: apply $\tau_{21} < 0.3$
- → Feasibility of the proxy
 - ★ taggers not trained with particle ID input
 → suitable to calibrate using the µ-tagged proxy





µ-tagged method: selected jets



Data and MC distribution applying the μ -tagged requirement and the central $\tau_{21} < 0.3$ cut

More taggers in backup

BOOST 2023

μ-tagged method: fit



Boosted Z method: concept

- → Proxy jets: Z→bb jets
 - ✤ $Z \rightarrow bb$ jet is a better proxy to $H \rightarrow bb$ jet
 - measure the SF at Z peak on top of the overwhelming QCD BKG



Boosted Z method: concept

- → Proxy jets: $Z \rightarrow bb$ jets
 - ✤ Z→bb jet is a better proxy to H→bb jet
 - measure the SF at Z peak on top of the overwhelming QCD BKG
- → Signal and background estimation
 - QCD yield in "fail" tagger region estimated from data
 - transferred to the QCD yield in "pass" region by a fitted polynomial on jet p_T and "mass"
 - "mass": a regressed mass via ParticleNet DNN architecture [See <u>CMS-DP-2021-017</u>]
 - signal: rely on accurate modelling of the Z+jet process (NLO correction applied)



Boosted Z method: fit



Systematic uncertainties (sfBDT & µ-tagged)

	sfBDT method	relative contribution to overall uncertainty
Uncertainty	source	$\Delta SF / (\Delta SF)_{tot}$
Statistical		49%
Theory		
Fractio	n of jet flavours	30%
ISR and FSR in parton shower		10%
Experiment	al	
Effect of varying sfBDT thresholds		38%
Effect of applying "reweighting schemes"		64%
Jet energy scale and resolution		6.4%
Integrated luminosity		1.0%
Pileup	reweighting	7.4%

µ-tagged method

Uncertainty source	$\Delta SF / (\Delta SF)_{tot}$
Statistical	34%
Theory	
Fraction of jet flavours	42%
ISR and FSR in parton shower	14%
QCD jet modelling	6.5%
Experimental	
Effect of varying τ_{21} thresholds	69%
Effect of "simulation-to-data reweighting"	28%
Jet energy scale and resolution	5.0%
Integrated luminosity	3.5%
Pileup reweighting	3.9%

Varying the yields of b/c/light flavour jets by 20%, separately

→ compensate for the mismodelling of flavour proportion in QCD simulation

Systematic uncertainties (sfBDT & µ-tagged)

sfBDT method	elative contribution to overall uncertainty		Varying the yields of b/c/light flavour jots
Uncertainty source Statistical Theory Fraction of jet flavours ISR and ESR in parton shower	$\frac{\Delta SF / (\Delta SF)_{tot}}{49\%}$	and the second	 by 20%, separately → compensate for the mismodelling of flavour proportion in QCD simulation
Experimental Effect of varying sfBDT thresholds Effect of applying "reweighting schemes" Jet energy scale and resolution Integrated luminosity Pileup reweighting	38% 64% 6.4% 1.0% 7.4%	× × × × × × × × × × × × × × × × × × ×	An external uncertainty is assigned to
µ-tagged method			selection (sfBDT or τ_{21}) varies

 $\Delta SF / (\Delta SF)_{tot}$

34%

42%

14%

6.5%

69%

28%

5.0%

3.5%

3.9%

→ a substantial uncertainty is obtained from this source

Integrated luminosity

Pileup reweighting

Fraction of jet flavours

QCD jet modelling

ISR and FSR in parton shower

Effect of varying τ_{21} thresholds

Jet energy scale and resolution

Effect of "simulation-to-data reweighting"

Uncertainty source

Statistical

Experimental

Theory

Systematic uncertainties (sfBDT & µ-tagged)

elative contribution to overall uncertainty	Varying the yields of b/c/light flavour jets
$\Delta SF / (\Delta SF)_{tot}$	by 20%, separately
49%	→ compensate for the mismodelling of
30%	flavour proportion in QCD simulation
10%	
38%	
64%	
6.4%	
1.0%	An externel un certainty is assigned to
7.4%	Solution in the span of derived SFs, when the
	elative contribution to overall uncertainty $\frac{\Delta SF/(\Delta SF)_{tot}}{49\%}$ $\frac{30\%}{10\%}$ $\frac{38\%}{64\%}$ 6.4% 1.0% 7.4%

µ-tagged method

Uncertainty source	$\Delta SF / (\Delta SF)_{tot}$	
Statistical	34%	
Theory		
Fraction of jet flavours	42%	Í I
ISR and FSR in parton shower	14%	
QCD jet modelling	6.5%	
Experimental		
Effect of varying τ_{21} thresholds	69%	
Effect of "simulation-to-data reweighting"	28%	*****
Jet energy scale and resolution	5.0%	
Integrated luminosity	3.5%	
Pileup reweighting	3.9%	

he selection (sfBDT or τ_{21}) varies

 \rightarrow a substantial uncertainty is obtained from this source

A residue mismodelling on the SV mass is

observed, especially in the sfBDT method, which is difficult to fix by using alternative simulation

→ adopt a conservative approach: reweight directly on the fitted mass variable in "pass+fail" region to check the change of SF

Systematic uncertainties (boosted Z)

boosted Z method	relative contribution to overall uncertainty	
Uncertainty sourceStatisticalTheoryISR and FSR in parton showerNLO correctionsPDF uncertaintiesExperimentalJet mass scale and resolutionJet energy scale and resolutionTrigger effiencyIntegrated luminosityPileup reweighting	$\frac{\Delta SF/(\Delta SF)_{tot}}{67\%}$ $\frac{45\%}{43\%}$ 5.5% 6.9% 25% 5.4% 11% 2.5%	Large statistical uncertainty for determining $R_{\rm P/F}$ is a major limiting power Theoretical uncertainties on Z+jets affect the Z peak yield, directly resulting to a larger SF uncertainty

Results and combination

ParticleNet-MD bbvsQCD (HP)





- Combination among methods via BLUE (best linear unbiased estimate)
- treating common uncertainty sources as correlated, and others as uncorrelated

Full SF results in backup

Results and combination

ParticleNet-MD bbvsQCD (HP)





- Combination among • methods via **BLUE** (best linear unbiased estimate)
- treating common uncertainty sources as correlated, and others as uncorrelated

Full SF results in backup

- **Remarks:**
 - All methods yield consistent results
 - Results also consistent with previous calibrated SFs, on the central SF, yeardependence, and level of uncertainties
 - The methods summarised in our works offer a systematic approach to handle all available X→bb/cc̄ taggers in CMS

Conggiao Li (Peking University)

BOOST 2023

Summary

CMS-PAS-BTV-22-001

- → Present a comprehensive summary of $X \rightarrow b\bar{b}/c\bar{c}$ taggers:
 - a review of taggers developed in CMS Run 2
 - performance comparison in ROC curves
 - calibration methods
 - sfBDT methods / μ-tagged methods / boosted Z method
- → Showcase individual and combined SF measurements
 - consistency found among methods, and with previous results
- → Benchmark the calibration methods for CMS Run 2 and Run 3
 - developed to handle a variety of tagger discriminants
 - marks the maturity of boosted-jet tagging and calibration technique at this stage

Backup

ParticleNet architecture



Congqiao Li (Peking University)

BOOST 2023

1.0

More spectra



Performance on simulation

CMS-PAS-BTV-22-001

high p_T: (600, +∞) GeV



ROC comparison on $H \rightarrow b\bar{b}$ (c \bar{c}) jets vs. inclusive QCD



Data and MC distribution on transformed tagger discriminant, applying the central sfBDT selection



Data and MC distribution on transformed tagger discriminant, applying the central sfBDT selection

sfBDT method: fit



- Fit on the SV mass (log scale) to distinguish the b/c/light templates
- each flavour template assigned a free-floating SF in fit
- simultaneous fit in the "pass" and "fail" region

sfBDT method: full workflow

See also in CMS-PAS-BTV-22-001

→ Target proxy jets: g→bb/cc̄ jets from QCD multijet events, selected by a BDT trained on QCD jets

→ Workflow:

- Select events with logical OR of H_T triggers with multiple thresholds
- select both leading and subleading jet (if it exists), passing p_T > 200 GeV, |η| < 2.4, and 500 < m_{SD} < 200 GeV, and requires N(matched-SV) ≥ 2
- ✤ reweight total simulated events to data on (jetIndex, H_T, p_T) bins
- all simulated jets are classified to b, c, and light flavour (based on the number of ghost-associated b and c hadrons, N_b and N_c)
- ★ [individual step] train the sfBDT using a special QCD sample (enriched in bbb/ccc by production), with jets selected in the same criteria and requiring N_b ≥ 2 or (N_b = 0 & N_c ≥ 2), and the gen-level τ₃₁ on hadrons applied to define signal and BKG jets for BDT training. With the sfBDT, for each target discriminant, determine the sfBDT curves on the 2D sfBDT—transformed tagger plane
- ◆ apply jet p_T selection for a dedicated fit point
- ✤ apply selection via the sfBDT curve, with 9×9 combinations in total
- ◆ perform fit on $\log(M_{SV_1}^{corr} / GeV)$: fit simultaneously in the "pass" and "fail" tagger region, with three SFs assigned to the b/c/light categories in the "pass" region. Target SF is SF_b (SF_c) for calibrating X→bb (cc̄) tagger
- assign two sources of external uncertainties:
 - Merge the target SFs (with all varied sfBDT selections) into one: the uncertainty in the final SF is larger than that of the individual fit result, and its central value is positioned in the middle of the entire SF series.
 - consider the deviation of SF as an external uncertainty, from the nominal SF to the SF derived in the schemes to reweight on the fit variable, or the sfBDT discriminant

µ-tagged method: selected jets



Data and MC distribution applying the μ -tagged requirement and the central τ_{21} < 0.3 cut

µ-tagged method: selected jets



Data and MC distribution applying the central τ_{21} < 0.3 cut

µ-tagged method: proxy—signal similarity



Signal and proxy similarity when different τ_{21} selections are applied

Congqiao Li	(Peking	University)
-------------	---------	-------------

μ-tagged method: fit



- Fit on the invariant mass of SV 4-vectors (log scale)
- each flavour template assigned a free-floating SF in fit
- simultaneous fit in the "pass" and "fail" region

μ-tagged method: full workflow

See also in CMS-PAS-BTV-22-001

→ Target proxy jets: g→bb/cc̄ jets from QCD multijet events (triggered by the existence of a soft muon in jets), with the additional offline requirements on the soft muon and N-subjettiness ratio τ₂₁

→ Workflow:

- select events with online triggers: requiring AK4 or AK8 jet with p_T > 300, and including a muon with p_T > 5 GeV
- ★ select both leading and subleading jets (if exists) with offline p_T > 350 GeV, |η| < 2.4, m_{SD} > 40 GeV
- apply a jet-based reweighting from QCD MC to data subtracting the tt and single top and V+jets MC, on the (p_T, η, τ₂₁)

(note: reweighting on τ₂₁ mitigate the mismodelling of gluon splitting jets by PYTHIA-based QCD MC, and reduce its discrepancy with the MG-based QCD MC on tau21 and all tagger discriminants)

- all simulated jets are classified to b, c, and light flavour (based on the number of ghost-associated b and c hadrons, N_b and N_c)
- ◆ apply jet p_T selection for a dedicated fit point
- * apply the selection via τ_{21} : the main scheme has $\tau_{21} < 0.3$, and it is varied from 0.2 to 0.4 in the auxiliary fit
- ◆ perform fit on $\log(M(\sum \vec{p}_{SV}^{corr}) / \text{GeV})$: fit simultaneously in the "pass" and "fail" tagger region, with three SFs assigned to the b/c/light categories in the "pass" region. Target SF is SF_b (SF_c) for calibrating X→bb̄ (cc̄) tagger
- assign two sources of external uncertainties:
 - merge the target SFs (with different τ₂₁ selections) into one, by expanding the uncertainty of the nominal fitted SF to cover the maximum deviation of the SFs.
 - consider the deviation of SF as an external uncertainty, from the nominal SF to the SF derived in the scheme to reweight on the fit variable

Boosted Z method: ROC from data



Boosted Z method: full workflow

See also in <u>CMS-PAS-BTV-22-001</u>

→ Target proxy jets: Z→bb jets from Z+jets events

→ Workflow:

- select events with online triggers: a combination of requirements on the jet p_T, trimmed mass, and event H_T
- ★ select the leading AK8 jet with p_T > 450 GeV, |η| < 2.4, M_{PNet} > 40 GeV as the target jet; presence of a subleading AK8 jet with p_T > 200 GeV and |η| < 2.4 as the recoil jet is required
- veto events with loose electrons/muons, and events with b-tagged AK4 jets (to suppress tt
 events)
- Perform the fit on the 2D binned histogram of (M_{PNet}, p_T)
 - Z+jets and W+jets modelled by MC
 - QCD in the "fail" region estimated directly from data, after subtracting other MC processes
 - QCD in the "pass" region estimated from the "fail" region through a fitted polynomial
 - b three SFs assigned to Z+jets (3 p_T bins) and one SF to W+jets, in the "pass" region
 - the target SF is that for the Z+jets since $Z \rightarrow b\bar{b}$ is the predominated contribution in the "pass region"
- → Note: the boosted Z method is only used to derive SFs for:
 - ParticleNet-MD bbvsQCD in HP, MP, LP; DeepDoubleBvL in HP, MP; double-b in HP
 - DeepAK8-MD bbvsQCD has a residue mass correlation which sculpts the QCD background in the "pass" region, making the method invalid

Full results: ParticleNet-MD bbvsQCD



31 July, 2023 **38**

BOOST 2023

Full results: DeepDoubleBvL

Congqiao Li (Peking University)



31 July, 2023 **39**

BOOST 2023

Full results: DeepAK8-MD bbvsQCD



31 July, 2023 **40**

BOOST 2023

Full results: double-b

Congqiao Li (Peking University)



BOOST 2023

Full results: ParticleNet-MD ccvsQCD



31 July, 2023 **42**

BOOST 2023

Full results: DeepDoubleCvL



31 July, 2023 **43**

BOOST 2023

Full results: DeepAK8-MD ccvsQCD



BOOST 2023