

---

# CMS Draft Analysis Note

*The content of this note is intended for CMS internal use and distribution only*

---

2017/07/21

Head Id: 417226

Archive Id: -1:417226M

Archive Date: 2017/07/21

Archive Tag: trunk

## Deep learning for jet reconstruction

The CMS Collaboration  
CERN

### Abstract

Deep learning lead to several breakthroughs outside the field of high energy physics, yet in jet reconstruction at CMS it was not used so far. This report shows results of applying deep learning strategies to jet reconstruction at the stage of tagging and calibration. Jets with cone radius 0.4 and fat jets with radius 0.8 where investigated. We introduce deep neural network structures that were not yet proposed in this context and show that in all cases we studies significant gain in performance can be achieved by this approach with respect to the established CMS methods. We systematically also review other recently proposed network structures in context of sub-structure tagging without flavours. The proposed strategy is a multi-classification and for jets with radius of 0.4 and 0,8, as well as also a transverse momentum estimation for slim jets. The classes include heavy objects like, top, H, Z, W for fat jets and heavy flavour, light quark and gluon tagging for the slim jets.

This box is only visible in draft mode. Please make sure the values below make sense.

PDFAuthor: none  
PDFTitle: Deep learning for jet reconstruction  
PDFSubject: CMS  
PDFKeywords: CMS, physics, software, computing

Please also verify that the abstract does not use any user defined symbols



---

## 1 Introduction

The reconstruction of jets is a central element in high energy physics collider experiments. Recently several studies using simplified simulation made first studies on using deep neural networks (DNNs) to identify (tag) the particle that caused a jet. Some used the analogy of the calorimeter cells to pixels in photographs to apply convolutional or dense networks that are often used for photo labeling [1–4]. The results were mixed, ranging from some improvement to no improvements with respect to established methods. Also recurrent neural networks were proposed [5, 6]. CMS and ATLAS released public documents [7, 8] on applying DNNs in context of flavour tagging and in CMS the default flavour tagger is derived from a DNN, that, for the first time, showed the gain in performance in real data for a real detector.

In this note we present results of using new DNN structures in the context of jet tagging and regression for jets with radii of 0.4 and 0.8, which are the default jets in CMS. In section 2, we discuss the samples used for to train the different tagger, the input variables (we from now on use the machine learning term: features) used for the tagging (classification), the generator level truth of the different particles ID (in the following we call these labels), e.g. B-hadron, and finally the pre-processing applied to the raw features. Section 3 describes the DNN architectures chosen for AK4 and AK8. Finally, in section 4, we show the results compared to the standard tools in CMS.

## 2 The setup for DNN training

### 2.1 Training samples for slim jets

For the training of AK4 jets we use the QCD and  $t\bar{t}$  samples listed in table 1. The generator used for  $t\bar{t}$  is POWHEGv1.0 [9–13] generators were used. Showering and hadronization is done by the PYTHIA 8.2 package [14] and the detector simulation by the GEANT4 [15] package. QCD is done with PYTHIA only. All samples are with simulated using phase1 detector design. After a pres-election, which reduces the gluon jets we altogether have about 80M jets for training, testing and validation.

The samples used for the training of the AK8 jets are listed in table 2. We are currently using 2016 samples which have much larger statistics than the available PhaseI 2017 MC samples (training a DNN using samples with limited number of simulated events can impact the performance).

Table 1: Simulated phase1 samples for the training for slim jets.

Sample	Full name
QCD	/QCD_Pt_XtoY_TuneCUETP8M1.13TeV-pythia8/PhaseFall16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM
tt	/TT_TuneCUETP8M2T4.13TeV-powheg-pythia8/PhaseFall16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM

Table 2: Simulated MC samples used for the training of the AK8 jets. The samples are from the Summer16 campaign.

Sample	Full name
t	/TT_Mit*_*TuneCUETP8M2T4.13TeV-powheg-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /ZprimeToTTjet_M*_*TuneCUETP8M1.13TeV-ncaclino-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /RSGLuonToTT_M*_*TuneCUETP8M1.13TeV-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /BulkGravToWWlepWhad_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /BulkGravToWW_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /RadiationToWlepWhad_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /RadiationToWW_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /RadiationToWW_width0p*_M*_*TuneCUETP8M1.13TeV-madgraph-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /WprimeToWZtoWhadZinv_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /WprimeToWZtoWhadZlep_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /WprimeToWZtoWhadZlep_width0p*_M*_*TuneCUETP8M1.13TeV-madgraph-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /ZprimeToWWtoWlepWhad_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /ZprimeToWW_width0p*_M*_*TuneCUETP8M1.13TeV-madgraph-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /ZprimeToWW_width0p*_M*_*TuneCUETP8M1.13TeV-madgraph-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /BulkGravToZZtoZhadZinv_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /BulkGravToZZtoZhadZinv_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /RadiationToZZtoZlepZhad_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /WprimeToWZtoWlepZhad_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /WprimeToWZtoWlepZhad_width0p*_M*_*TuneCUETP8M1.13TeV-madgraph-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /BulkGravTohhTohhVhbb_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /BulkGravTohhTohhbbbb_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /WprimeToWhToWlepHbb_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /ZprimeToZhToZinvHbb_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /ZprimeToZhToZlepHbb_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /WprimeToWZtoWhadZhad_narrow_M*_*13TeV-madgraph-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /WprimeToWhToWhadHbb_narrow_M*_*13TeV-madgraph/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM /QCD_Pt_*_*TuneCUETP8M1.13TeV-pythia8/RuntlSummer16MiniAODv2-PUMoriond17.80X.mcRun2.asymptotic.2016.TrancheIV_v6-v1/MINIAODSIM
W	
Z	
H	
W/Z	
W/H	
QCD	

## 2.2 Labeling

We define the labeling in major and minor labels. The major labels define a broader category of jets, e.g. jets with at least one B hadron. The minor labels are a further sub-division of the major labels, e.g. jets with two or more B hadrons. All major labels with respect to all other major labels and all minor labels with respect to all other minor labels are orthogonal. The major label are the sum of their sub-labels.

### 2.2.1 Slim jet labels

Heavy flavor hadrons, scaled to negligible transverse momentum in order not to impact the final properties of the jet, are added to the list of stable particles to be clustered by the AK4 jet algorithm. Jets containing one or more heavy flavor objects in its constituents are assigned one of the major heavy flavor labels, which than are further sub-divided to minor labels to separate different decays or the number of heavy flavour hadrons in a jet. Jets not containing any clustered heavy-flavor hadron are labeled according to the flavor of the hardest (maximum transverse momentum) parton with PYTHIA 8 `status = 23`, assigning either the `light quark` or the `gluon` labels. This labeling of `light quarks` or `gluons` is following the “physics definition” as defined in [16]. Jets with no heavy-flavor hadron clustered among the constituents, but with a heavy-flavor quark as hardest parton, are considered undefined and excluded from the training and evaluation procedure. This is a limitation of the miniAOD data format used to extract the dataset. A summary of all major and minor label is shown in table 3

Table 3: The three major flavour label and the sub division of each of these flavour in even more detailed label. The sum of all minor (sub) labels is equivalent to the major label.

Major label	Minor (sub-)label
B, $\geq 1$ B hadron	bb, two ore more B hadrons $b_{lep}$ , exactly one B hadron with leptonic decay $b$ , exactly one B hadronic hadronic decay
C, $\geq 1$ C hadron and no B hadron	cc, two or more C hadrons c, exactly one C hadron
L, None of the above and parton matched	uds, physics definition [16] g, physics definition [16]

### 2.2.2 AK8 labels

The multi-classification (top, H, Z and W tagging) approach followed for the AK8 jets requires mutually exclusive labelling. Priority is given to the hadronically decaying heavy objects (i.e. top, H, Z, W). The generated heavy object,  $X$ , and its decay products,  $X_{decay}$  are matched to the AK8 jet following the conditions:  $\Delta R(X, AK8) < 0.6$  and  $\Delta R(X_{decay}, AK8) < 0.6$ . Then AK8 jets with heavy flavor content are identified following the BTV-style. The remaining jets are identified as light quarks/gluons. The proposed labelling is summarized in table 4 and is designed to have high granularity. The various labels can be easier combined. We therefore two differnt kind of labels, “major” and “minor”, to target the different analysis needs. We consider AK8 jets with transverse momentum,  $p_T$ , greater than 300 GeV.

The proposal described above is currently pending approval from the relevant CMS sub-groups. For the results shown in this version of the note we follow a simplified approach presented in 5.

Table 4: Summary of the various labels proposed for the classification of the AK8 jets.

Major label	minor label
t	bcq bqq bc bq
W	cq qq
Z	bb cc qq
H	bb cc qq qqqq
QCD	bb b q g

Table 5: Summary of the labelling used for the results presented in the current version of the note for the classification of the AK8 jets.  $q_X$  refers to the quarks from the hadronic decay of the boson  $X$ .

Major label	Requirement
t	$\Delta R(b, AK8) < 0.8$ and $\Delta R(q_W, AK8) < 0.8$
W	if W is from t decay: $\Delta R(b, AK8) > 0.8$ and $\Delta R(q_W, AK8) < 0.8$ else: $\Delta R(q_W, AK8) < 0.8$
Z	$\Delta R(q_Z, AK8) < 0.8$
H	$\Delta R(q_H, AK8) < 0.8$
QCD	anything else

## 66 2.3 Preprocessing

67 Before the input variables are feed into the neural networks for the training or the evaluation  
68 of a trained model, they are subject to a preprocessing. Goal of the preprocessing is to avoid  
69 inputs to the neural network with significantly different scales or unintended biases as this  
70 leads to easier convergence of the minimization.

### 71 2.3.1 Slim jet preprocessing

The scales are unified using the mean  $\langle x \rangle$  and the standard deviation  $\sigma_x$  of each feature  $x$  which is rescaled to

$$x^I = \frac{x - \langle x \rangle}{\sigma_x}. \quad (1)$$

72 The  $p_T$  and  $\eta$  of each jet are direct input to the neural network, such that the evaluation of  
73 the input features can be adjusted according to the jet kinematics. However, the kinematics of  
74 jets originating from a different parton flavours show partially significant deviations from each  
75 other. In consequence the neural network could learn to assign a jet flavour purely based on

76 the  $p_T$  and  $\eta$  of the jet. To avoid such biases, jets are removed from the training sample, such  
 77 that their  $p_T$  and  $\eta$  distributions agree for all jet flavours. As reference  $p_T$  and  $\eta$  shape we use  
 78 the shape of the b-jets, i.e. finally the other labels will have all the b-jet shape. In addition we  
 79 remove 50% of the gluons in order to avoid a gluon dominated training dataset. The probabil-  
 80 ities to remove a jet are calculated based on the entire training sample to minimise the impact  
 81 of statistical fluctuations.

82 In case a feature is missing we put in default values, that are not far from the normalized scale  
 83 and are not overlapping with the core distribution.

84

### 85 2.3.2 AK8 preprocessing

The input variables for AK8 tagging are preprocessed in a similar way as in AK4 tagging. Each  
 input variable  $x$  is transformed according to Eq. (2.3.2),

$$x^I = \frac{x - p_{50\%}}{p_{84\%} - p_{50\%}}, \quad (2)$$

86 where  $p_{50\%}$  and  $p_{84\%}$  are the 50th and 84th percentiles of the variable  $x$ . In the case when  
 87  $p_{84\%} = p_{50\%}$  (which happens for some discrete variables), the denominator is taken to be 1. The  
 88 use of percentiles instead of the mean and the standard deviation tends to be less sensitive to  
 89 outliers and distributions with long tails, leading to more unified scales for different variables.  
 90 The transformed values are further clipped to be in the range of  $[-5, 5]$  before feeding into the  
 91 neural networks, which are found to help improve the stability of neural network training.

92 To avoid biases from the difference in the jet  $p_T$  spectrum, jets in the training sample are  
 93 reweighted to have a flat distribution in  $p_T$ , and the contribution of each source (top, W, Z,  
 94 Higgs, and QCD) is equalized. However, from our studies we found that applying weights to  
 95 the neural network training often causes degradation of performance or even failure in con-  
 96 vergence. As a result, the reweighting is done “on-the-fly” by randomly sampling the training  
 97 dataset according to the “flattening” weights, thus effectively achieves the reweighting without  
 98 losing statistics.

99 For evaluating the performance, jets in the testing sample are reweighted such that different  
 100 signal processes (top, W, Z, Higgs) all have the same  $p_T$  spectrum as the background process  
 101 (QCD).

## 102 2.4 Input features

103 The basis for the taggers are the Particle Flow [17] jet constituents (particle candidates), namely  
 104 charged and neutral PF candidates as well as reconstructed secondary vertices within the jet.

### 105 2.4.1 slim jet input features

106 For the DeepFlavour tagger, several features of the jet constituents and of secondary vertices  
 107 within a cone of  $\Delta R = 0.4$  with respect to the jet axis are used. In some cases, their variation is  
 108 restricted to a reasonable range to avoid large outliers e.g. due to mis-measurements having a  
 109 strong effect on the training without providing any discrimination power. In addition, partic-  
 110 ular inputs are shifted by a constant offset, such that 0 corresponds to a value that is outside of  
 111 the bulk of the distribution and provides no handle on the flavour separation. For the charged  
 112 PF candidates, the majority of the input features are calculated following previous b-tagging  
 113 algorithms [18]. These are in the following indicated as BTV features and their exact definition  
 114 can be found in the reference [18]. The additional variables are described in the following.

- 115 •  $p_T(j)$ : jet  $p_T$
- 116 •  $\eta(j)$ : jet  $\eta$
- 117 •  $N_{cPF}$ : number of charged PF candidates within the jet
- 118 •  $N_{nPF}$ : number of neutral PF candidates within the jet
- 119 •  $N_{SV}$ : number of secondary vertices within the jet
- 120 •  $N_{PV}$ : number of primary vertices in the event
- 121 •  $p_T(cPF)/p_T(j)$ : relative  $p_T$  of a charged jet constituent with respect to the jet  $p_T$
- 122 •  $p_T(nPF)/p_T(j)$ : relative  $p_T$  of a neutral jet constituent with respect to the jet  $p_T$
- 123 •  $\Delta R_m(cPF, SV)$ :  $\Delta R$  of charged candidate and closest secondary vertex within the jet
- 124 •  $\Delta R_m(nPF, SV)$ :  $\Delta R$  of neutral candidate and closest secondary vertex within the jet
- 125 • VTAss: flags indicating whether the charged particle track is used in the primary
- 126 vertex fit, includes steps from low purity to high purity requirements.
- 127 • fromPV: similar to VTAss, but partially including information about the primary
- 128 vertex fit quality. Can indirectly include lepton information
- 129 •  $w_p(cPF)$ : weight assigned to the charged particle by the PUPPI [19] algorithm
- 130 •  $w_p(nPF)$ : weight assigned to the neutral particle by the PUPPI algorithm
- 131 •  $\chi^2$ : charged PF candidate track  $\chi^2$
- 132 • quality: flag that indicates the charged particle track reconstruction quality, from
- 133 passing low purity to high purity requirements
- 134 •  $\Delta R(cPF)$ :  $\Delta R$  to jet axis of a charged candidate
- 135 •  $\Delta R(nPF)$ :  $\Delta R$  to jet axis of a neutral candidate
- 136 • isGamma: flag whether a neutral candidate passes loose photon identification re-
- 137 quirements
- 138 • hadFrac: fraction of energy deposits in the hadronic calorimeter, only for neutral
- 139 candidates
- 140 •  $p_T(SV)$ : secondary vertex  $p_T$
- 141 •  $\Delta R(SV)$ :  $\Delta R$  between jet axis and secondary vertex flight direction
- 142 •  $m_{SV}$ : invariant mass of reconstructed secondary vertex
- 143 •  $N_{tracks}(SV)$ : number of tracks associated to the secondary vertex
- 144 •  $\chi^2(SV)$ : secondary vertex  $\chi^2$
- 145 •  $\chi_n^2(SV)$ : secondary vertex  $\chi^2$  normalised to degrees of freedom
- 146 •  $d_{xy}(SV)$ : transverse impact parameter of secondary vertex
- 147 •  $S_{xy}(SV)$ : transverse impact parameter significance of secondary vertex
- 148 •  $d_{3D}(SV)$ : 3D impact parameter of secondary vertex
- 149 •  $S_{3D}(SV)$ : 3D impact parameter significance of secondary vertex
- 150 •  $\cos \theta(SV)$ :  $\cos \theta$  of secondary vertex with respect to primary vertex
- 151 •  $E_{rel}(SV)$ : ratio of secondary vertex energy with respect to the jet

152 All global features with per-jet values that are considered are summarised in Table 6. No offsets,  
 153 upper or lower bounds are applied. These are applied to particular properties or charged and  
 154 neutral PF candidates, and secondary vertices as listed in Tables 7, 8 and 9.



Table 6: List of global input features for the AK4 DeepFlavour tagger

feature	comment
$p_T(j)$	
$\eta(j)$	
$N_{cPF}$	
$N_{nPF}$	
$N_{SV}$	
$N_{PV}$	
trackSumJetEtRatio	BTV
trackSumJetDeltaR	BTV
vertexCategory	BTV
trackSip2dValAboveCharm	BTV
trackSip2dSigAboveCharm	BTV
trackSip3dValAboveCharm	BTV
trackSip3dSigAboveCharm	BTV
jetNSelectedTracks	BTV
jetNTracksEtaRel	BTV

Table 7: Full list of charged PF candidate features used as input to the DeepFlavour network for AK4 jets

feature	offset	lower bound	upper bound	comment
trackEtaRel	-	-5	15	BTV
trackPtRel	-	-	4	BTV
trackPPar	-	$-10^5$	$10^5$	BTV
trackDeltaR	-	-5	5	BTV
trackPParRatio	-10	100	-	BTV
trackSip2dVal	-	-	70	BTV
trackSip2dSig	-	-	$4 \cdot 10^4$	BTV
trackSip3dVal	-	-	$10^5$	BTV
trackSip3dSig	-	-	$4 \cdot 10^4$	BTV
trackJetDistVal	-	-20	1	BTV
trackJetDistSig	-	-1	$10^5$	BTV
$p_T(cPF)/p_T(j)$	-1	-1	0	
$\Delta R_m(cPF, SV)$	-5	-5	0	
fromPV	-	-	-	
VTXass	-	-	-	
$w_p(cPF)$	-	-	-	
$\chi^2$	-	-	-	
quality	-	-	-	

Table 8: Full list of neutral PF candidate features used as input to the DeepFlavour network for AK4 jets

feature	offset	lower bound	upper bound
$p_T(nPF)/p_T(j)$	-1	-1	0
$\Delta R_m(nPF, SV)$	-5	-5	0
isGamma	-	-	-
hadFrac	-	-	-
$\Delta R(nPF)$	-0.6	-0.6	0
$w_p(cPF)$	-	-	-

Table 9: Full list of secondary vertex features used as input to the DeepFlavour network for AK4 jets

feature	offset	lower bound	upper bound
$p_T(SV)$			
$\Delta R(SV)$	-0.5	-2	0
$m_{SV}$	-	-	-
$N_{\text{tracks}}(SV)$	-	-	-
$\chi^2(SV)$			
$\chi_n^2(SV)$	0	-1000	1000
$d_{xy}(SV)$	-	-	-
$S_{xy}(SV)$	-	-	800
$d_{3D}(SV)$	-	-	-
$S_{3D}(SV)$	-2	-2	0
$\cos \theta(SV)$	-	-	-
$E_{rel}(SV)$	-	-	-

155 The particles and vertices are ordered using a hierarchical sorting algorithm. Charged can-  
 156 didates and secondary vertices are sorted by impact parameter significance. If the charged  
 157 candidate was used in the primary vertex fit, they are appended starting from the lowest  
 158  $\Delta R_m(cPF, SV)$  value. If no secondary vertex is present within the jet, the particle  $p_T$  is used  
 159 instead. The latter two sorting requirements are also applied to neutral PF candidates.

## 160 2.4.2 AK8 input features

161 The input features used by AK8 tagging are similar to those used in the AK4 DeepFlavour  
 162 tagger. They are organized into three groups: inclusive (charged and neutral) PF candidates,  
 163 charged PF candidates, and secondary vertices. We take up to 100 inclusive PF candidates,  
 164 sorted in descending  $p_T$  order, and up to 60 charged PF candidates and up to 5 secondary  
 165 vertices, ordered by impact parameter significance. The full lists of variables used in each  
 166 group are summarized in Table 10 to 12.

Table 10: Full list of charged PF candidate features used as input to the DeepAK8 network

feature	comment
trackEtaRel	BTV
trackPtRatio	BTV
trackPParRatio	BTV
trackSip2dVal	BTV
trackSip2dSig	BTV
trackSip3dVal	BTV
trackSip3dSig	BTV
trackJetDistVal	BTV
$p_T(cPF) / p_T(j)$	
$E_{rel}(cPF)$	
$\Delta\phi(cPF, j)$	
$\Delta\eta(cPF, j)$	
$\Delta R(cPF, j)$	
$\Delta R_m(cPF, SV)$	
$\Delta R(cPF, \text{subject 1})$	
$\Delta R(cPF, \text{subject 2})$	
$\chi_n^2$	
quality	
$d_z$	
$S_z$	
$d_{xy}$	
$S_{xy}$	
track_dptdpt	track covariance
track_detadeta	track covariance
track_dphidphi	track covariance
track_dxydxy	track covariance
track_dzdz	track covariance
track_dxydz	track covariance
track_dphidxy	track covariance
track_dlambdadz	track covariance

## 167 3 Deep neural network architectures

168 The neural network structure was designed to be able to make good use of the large input we  
 169 give to the neural network. In contrast to previous proposals we use more information per  
 170 particles candidate or vertex. This lead to the special challenge to digest the huge amount of  
 171 input features. In order to not expose the later layer to such a huge amount of features we  
 172 build a reduced set features per particle (or per few particles) candidate or vertex by so called

Table 11: Full list of inclusive PF candidate features used as input to the DeepAK8 network

feature
$p_T(PF)/p_T(j)$
$E_{rel}(PF)$
$\Delta\phi(PF, j)$
$\Delta\eta(PF, j)$
$\Delta R(PF, j)$
$\Delta R_m(PF, SV)$
$\Delta R(PF, \text{subject 1})$
$\Delta R(PF, \text{subject 2})$
$w_p(PF)$
$f_{HCAL}$

Table 12: Full list of secondary vertex features used as input to the DeepAK8 network

feature
$p_T(SV)/p_T(j)$
$E_{rel}(SV)$
$\Delta\phi(SV, j)$
$\Delta\eta(SV, j)$
$\Delta R(SV, j)$
$p_T(SV)$
$m_{SV}$
$N_{tracks}(SV)$
$\chi_n^2(SV)$
$d_{xy}(SV)$
$S_{xy}(SV)$
$d_{3D}(SV)$
$S_{3D}(SV)$
$\cos\theta(SV)$

convolutional layers. Convolutional layers learn a transformation from a typically higher dimensional representation to a lower representation of features, which in our physics jargon would be similar to building a few variables from a larger input. This is done simultaneously with the overall optimization, i.e. the transformation is trained to be ideal for the classification. Convolution networks are very spread in image recognition, where they effectively summarize small region of the image and build more useful features than the raw pixels, like edges or alike, which than are feed to the following layers. In our case a particle candidate or vertex takes the role of such a small region of an image. While slim and fat jets share this basic structure in the beginning, we currently use slightly different networks structures in the later layers.

### 3.1 Slim jet DNN architecture

The first layers are convolutional layers as explained in the previous paragraph. Figure 1 indicates the number of layer and nodes for these convolutional layers. To allow non-linearities we use up to four convolutional layers. The convolution are done  $1 \times 1$ , i.e. they are applied only to individual particle candidates and they only reduce the dimension of the feature per candidate or vertex, but are not a summary of several candidates. We use the rectified linear unit (ReLU) activation function.

From convolutional layer we get sequences of features of particle candidates. The sequence order is still defined from the input particle candidate (or vertex) sorting. They are sorted by displacement significance. The most displaced are the last in the list. In case the particles are not displaced and no secondary is in the jet, they are sorted with increasing  $p_T$ . Exact sorting details are in 2.4. These sequences are than feed into recurrent neural networks (LSTM) and by that compressed to a single vector per sequence, i.e. charged and neutral candidates and vertices. When using recurrent networks the ordering is important, thus our underlying assumption is that the most displaced (in case of displacement) or the highest  $p_T$  candidates matter the most.

The output of the recurrent layers is than combined with the global variables, like  $p_T$  and  $\eta$ . This is put into a fully connected neural network with 8 layers. The first layer has 200 nodes and the latter 100. Again we use ReLU activation.

In between the layer we use a dropout of 0.1 and do batch normalization apart for the input (layer 0) and output (last layer before loss). For the final layer we use the softmax function as activation and cross entropy as the loss to minimize. For the minimization we use the adam [20] optimizer and train for 50 epochs. The workflow was implemented using [21] that relies on [22, 23] for the neural network implementation. To check for over-training we use separate sample that is not used for training and no over-training was found. The final ROCs curves in the results section 4 using another third set of independent samples.

### 3.2 DNN AK8

The task of tagging heavy objects (top, W, Z, Higgs) with AK8 jets is more challenging than b-tagging in some aspects. With a larger jet radius, a typical AK8 jet has many more constituent particles than AK4 jets. And the interrelationship between these particles, like the spatial pattern and the energy correlation, is more crucial for heavy objects than for b-tagging. Thus, a more complex DNN architecture is adopted for AK8 tagging.

Similar to the DNN model for AK4 tagging, the DNN model for AK8 tagging, as illustrated in Fig. 2, first processes inclusive PF candidates, charge PF candidates and SVs separately with convolutional neural networks, and then combines outputs from these three networks in a fully-connected layer before yielding the final prediction. The network is trained as a whole to

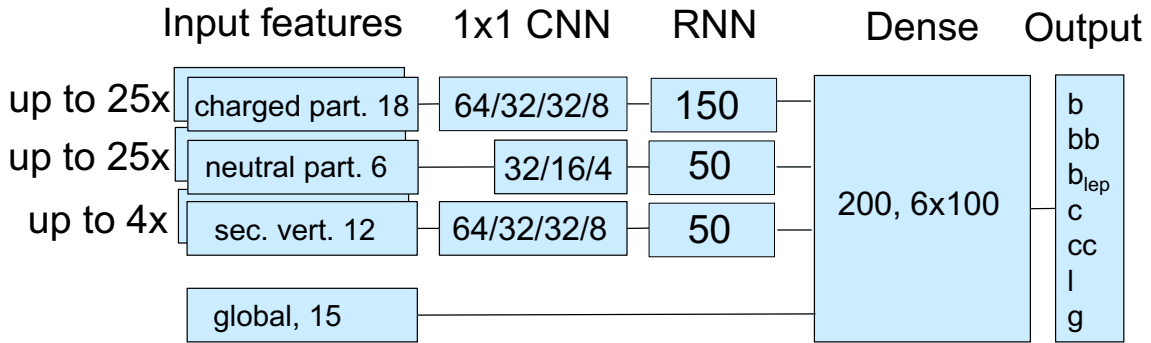


Figure 1: DNN architecture illustration. Dropout and batch normalizations are not indicated. The number in the boxes indicate the number of nodes per layer.

219 optimize all the components simultaneously.

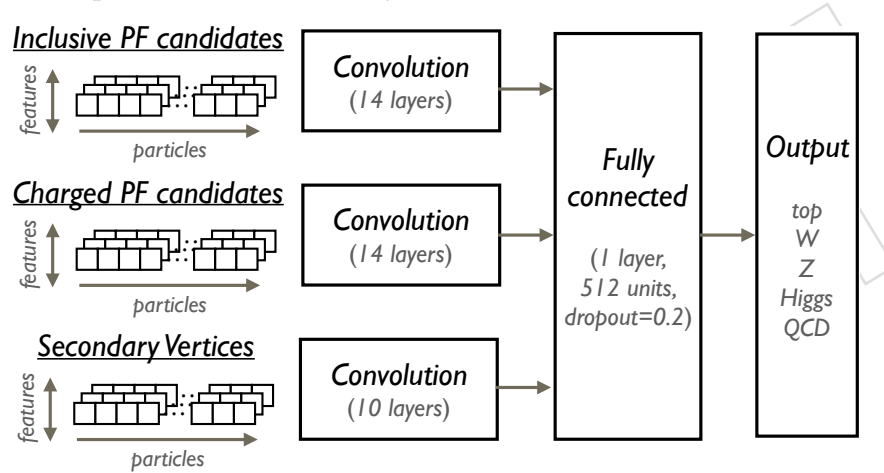


Figure 2: Illustration of the overall DNN architecture used by AK8 tagging.

220 Different from the convolutional neural networks in AK4 tagging, where the convolution is  
 221 performed for each individual particle (i.e., “1x1”), the convolution here is performed for each  
 222 adjacent particle triplet (i.e., “3x1”) with overlaps (in CNN jargon, we use a kernel size of  
 223 3 and a stride of 1). Such “3x1” convolutions are stacked on top of each other, thus allow  
 224 the DNN to see the correlation between nearby particles at earlier stages, and to have a more  
 225 global view of the particle correlations at later stages. The design of the convolutional neural  
 226 networks model is largely based on the ResNet model [24], which is one of the state-of-the-art  
 227 model for image recognition. We adapt it to work with one-dimensional particle list instead  
 228 of two-dimensional pictures, but adopt the main structure and all important ingredients such  
 229 as residual connection [25], batch normalization [26], and ReLU [27] activation function. The  
 230 depth of the convolutional network is 14 for inclusive PF candidates and charge PF candidates,  
 231 and 10 for SVs. The filter sizes (i.e., the number of output features) for each convolutional layer  
 232 ranges between 32 to 128.

233 The outputs from the three separate convolutional neural networks are combined in a fully-  
 234 connected layer with 512 units, followed a ReLU activation and a DropOut layer with a rate  
 235 of 0.2. We use the softmax function in the final layer to yield the final prediction, and cross  
 236 entropy as the loss function to minimize. The neural network is implemented with the MXNet

237 package [28] and trained with the Adam [20] optimizer with a learning rate of 0.001.

## 238 4 Results in simulation

239 We compare the results of the classification for AK4 and AK8 jets to references that are used in  
240 public analysis in CMS. We use physics sample with label composition and  $p_T$  and  $\eta$  shapes as  
241 they come naturally from the samples process in question. We reject jets with undefined labels.

### 242 4.1 AK4 jet results

243 For AK4 we use the CSVs2 and DeepCSV b-tagger as reference [18, 29] and for the quark-gluon  
244 discrimination we compare to [30] and alternative deep neural network structures. Figure 3  
245 compares by showing the ROC curves the DeepFlavour tagger results to the former default  
246 CMS tagger CSVv2 and DeepCSV for different processes and  $p_T$ . For both physics processes,  
247  $t\bar{t}$  and QCD we see significant gain in all region of the ROC curves. For very high b-jet  $p_T$  the  
248 b-jet efficiency is increased by 50% with respect to the DeepCSV for a light fake rate of 1%.

249 At higher  $p_T$  of jets gluon splitting leads to an increased amount of jets in QCD with two b-  
250 hadrons inside the jet. In a high  $p_T$  region we thus show in Figure 4 the efficiency for jets with  
251 single b-hadron using only the single b and leptonic b labels as discriminator. Identifying single  
252 bs is slightly more difficult than double bs and the performance is slightly less good. We also  
253 show the tagging performance for the bb label, using also the bb discriminator and it can be  
254 seen that separating bb from light jets is easier, as the performance is improved with respect to  
255 the single b case. The second curve in the double b ROC is the separation of b and bb, using the  
256 probabilities as binary classifier (binary means here that two estimated probabilities (double b  
257 and single b) used are renormalized to add up to one, before they are used as discriminator).  
258 It is interesting that about 1/3 of the jets can be identified as double b with only a fake-rate of  
259 1% for single bs. Figure 4 also shows the efficiency of leptonic decays of single bs vs the mistag  
260 rate for light jets and we see a decent separation, just in between the bb and single b case. The  
261 separation of hadronic and leptonic bs, even using the separation as binary classifier, does not  
262 lead to good leptonic decay separation. It should be noted that we did not explicitly add lepton  
263 information to the tagger.

264 The c-tagging of DeepJet is compared to DeepCSV in Figure 5. We do see a gain with respect  
265 to DeepCSV also for c-tagging.

266 We defined three working points, which lead to a light jet efficiency of 10%, 1%, and 0.1% respec-  
267 tively for jets of the QCD sample in a range from 80 to 120 GeV. Using these working points,  
268 called loose, medium and tight, we illustrate the dependence of the tagger performance as a  
269 function of  $p_T$ ,  $\eta$ , and number of primary vertices in Figure 6. We see the expected degradation  
270 of performance with higher  $p_T$  and at low  $\eta$ , for large  $\eta$  and large number of primary vertices.  
271 Note that the QCD sample has a relatively flat  $p_T$  distribution, thus the  $p_T$  integrated illustra-  
272 tions are dominated by high  $p_T$ .

273 An overview of the discriminator shapes us given in Figure 7. The discriminators shown  
274 in Figure 7 are according to labels that were present in CMS before and methods to estimate  
275 the data simulation agreement re present. The minor labels of the b-hadron jet major label  
276 are shown in the appendix Figure A. Especially the double b vs single b separation seems  
277 promising and it might motivate a dedicated effort to develop methods to also establish these  
278 discriminators in data in the future.

279 Figure 8 and 9 show the comparison of DeepJet quark gluon separation to the default like-  
280 likelihood method for different  $p_T$  in the central and forward region of the detector, respectively.  
281 The output of DeepJet was made to a binary classifier to compare to the quark-gluon likelihood

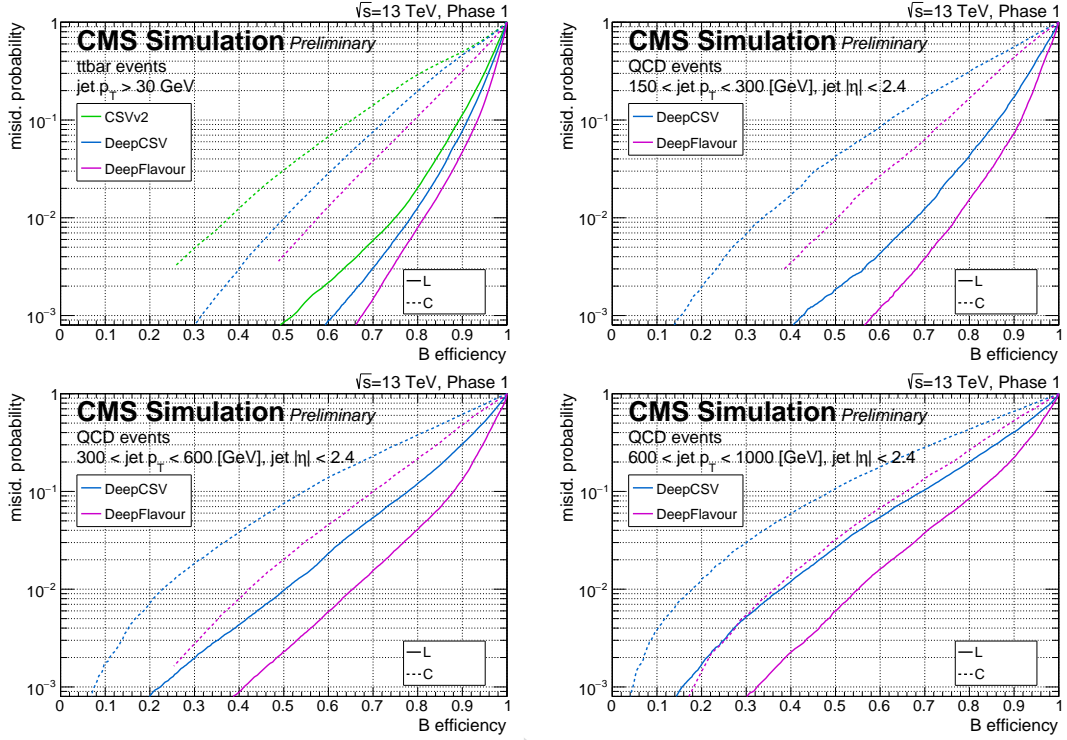


Figure 3: ROC curves for different processes and  $p_T$  ranges for b-tagging. Here all three categories,  $b$ ,  $b_{lep}$ , and  $bb$  are considered as b-jet. Light includes the gluon and uds-quark categories. The top left plots show the ROC evaluated with  $t\bar{t}$  events. The latter show the ROCs for QCD samples with increasing  $p_T$ .

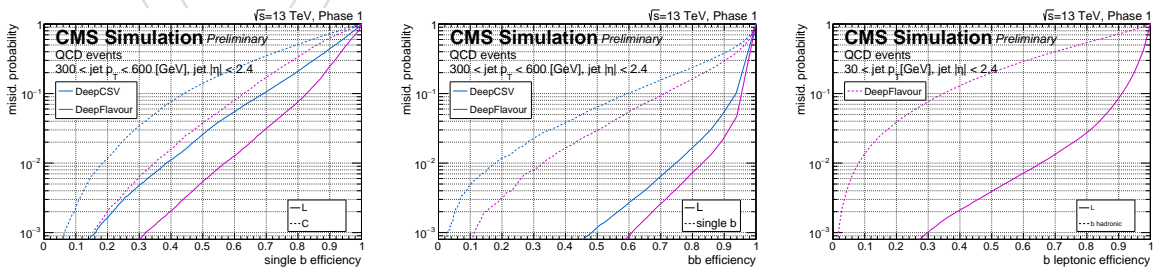


Figure 4: ROC curves for different b categories. For the efficiency vs light and c-jet always the unmodified multi-label probability of the b category is used. For the double b vs b and leptonic b decay vs. hadronic the binary classifiers between these categories are used.



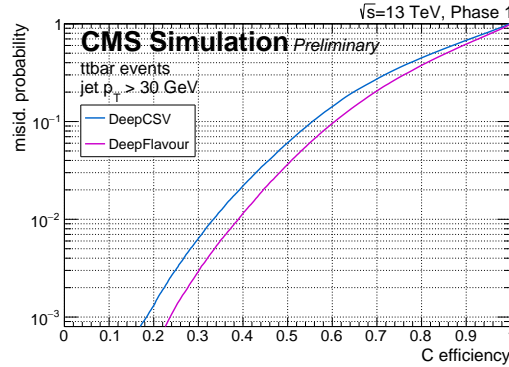


Figure 5: ROC curves for C-tagging in  $t\bar{t}$  events using the estimated probabilities as binary classifier for light jet and c-jet separation.

282 method, which technically was done by scaling the estimated probabilities for light quarks and  
 283 gluons such that they add up to one. We see constantly a significant improvement of about  
 284 5-10% absolute better efficiency for light quarks compared to likelihood method. Note that  
 285 for this comparison all jets of the QCD sample are taken into account that pass the kinematic  
 286 criteria and where we did find that the label was well defined, i.e. no balance selection of jets  
 287 was applied and instead it was only checked that a parton was found. We also conformed that  
 288 a balance selection as done by other quark gluon taggers did not effect the conclusions, but  
 289 moved all ROCS by a tiny amount towards better tagging.

290 We also added as comparison other DNN architectures that should also be able to find jet  
 291 structures, but would be blind to heavy flavour. For these DNNs we use only  $p_T$  of the candi-  
 292 dates relative to the jet, relative  $\phi$ , relative  $\eta$ , if the particles are charged or not, and the puppi  
 293 weight. We try two DNN structures, one according to [3], i.e. an image with 22 bins and a  
 294 convolutional DNN. We use relative eta from each particle candidate to define the bin which  
 295 the particle belongs to. For each bin we store the  $p_T$  sum of the particles with puppi weight  
 296 (we tried both, with and without puppi weights) and the multiplicities of charged and neutral  
 297 particle candidates. Also the exact details of the layer structure are used as in the reference with  
 298 the only difference that we could remove the regularization layers, as we use larger samples.  
 299 Alternatively we also took the list of particles candidates (charged and neutral), with the same  
 300 above mentioned information sorted in descending  $p_T$ . As for DeepFlavour we use a recurrent  
 301 (LSTM) network followed by several dense layer. We compared these different flavour blind  
 302 neural network structures to DeepFlavour. All three structure give similar results as seen in 8  
 303 and 9.

304 Another study we did is the impact of reduced input, i.e. using a few human made variables.  
 305 We gave as input the five variables currently used by the BDT quark gluon effort and added  
 306  $p_T$ ,  $\eta$  and rho. This made altogether 8 input variables. Again we used a DNN very similar to  
 307 the one of the recommended CMS flavour tagger DeepCSV. Only 15 nodes per layer are used,  
 308 which as for DeepCSV is between 1-2 times the inputs. We used 7 hidden layers. The compari-  
 309 son can be found in 10 and 11. We see a gain by using a full DNN with larger input than only  
 310 the human made variables. The effect varies depending on the  $p_T$  and  $\eta$  and can be sizable at  
 311 background rejection of around 10%.

312 To have an illustrative example on how the multi-classification simplifies real-life tasks, we  
 313 show how one can select light quarks in  $t\bar{t}$  events. traditionally this was quite difficult, as we  
 314 had a tagger to separate gluons and quarks, one to separate b-jets from "light" jets and yet  
 315 another to separate c-jets from "light" jets. Finally one would apply cuts on all three taggers or  
 316 even build a meta-tagger based on other taggers output. For DeepJet it is sufficient to just ask

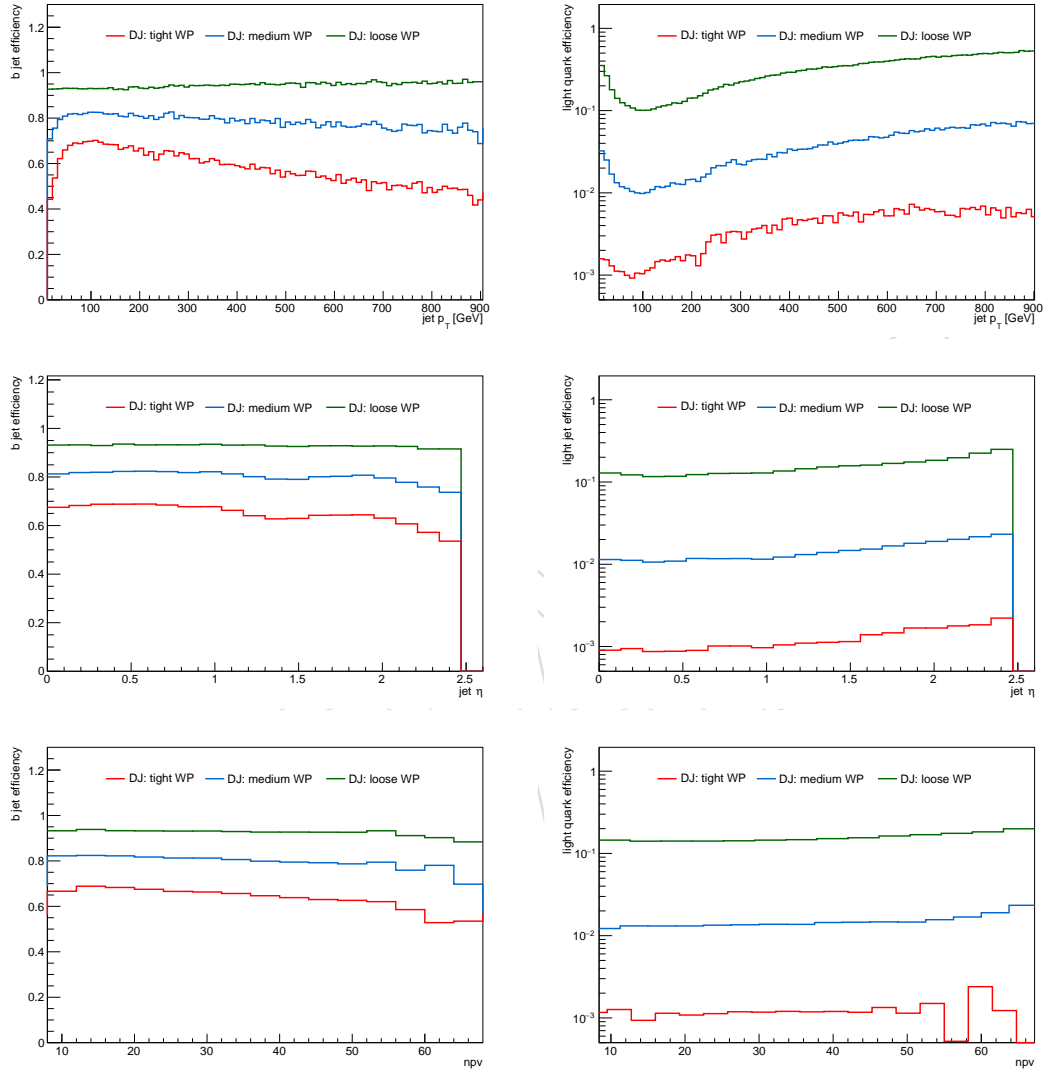


Figure 6: Efficiency for B and L in the QCD sample as functions of  $p_T$ ,  $\eta$  and number of primary vertices (npv). To not be dominated by high  $p_T$  jet in QCD, for the  $\eta$  and npv dependencies only jets with  $p_T > 30$  and  $< 150$  GeV were used. Efficiencies for three working points: loose, medium and tight are shown.

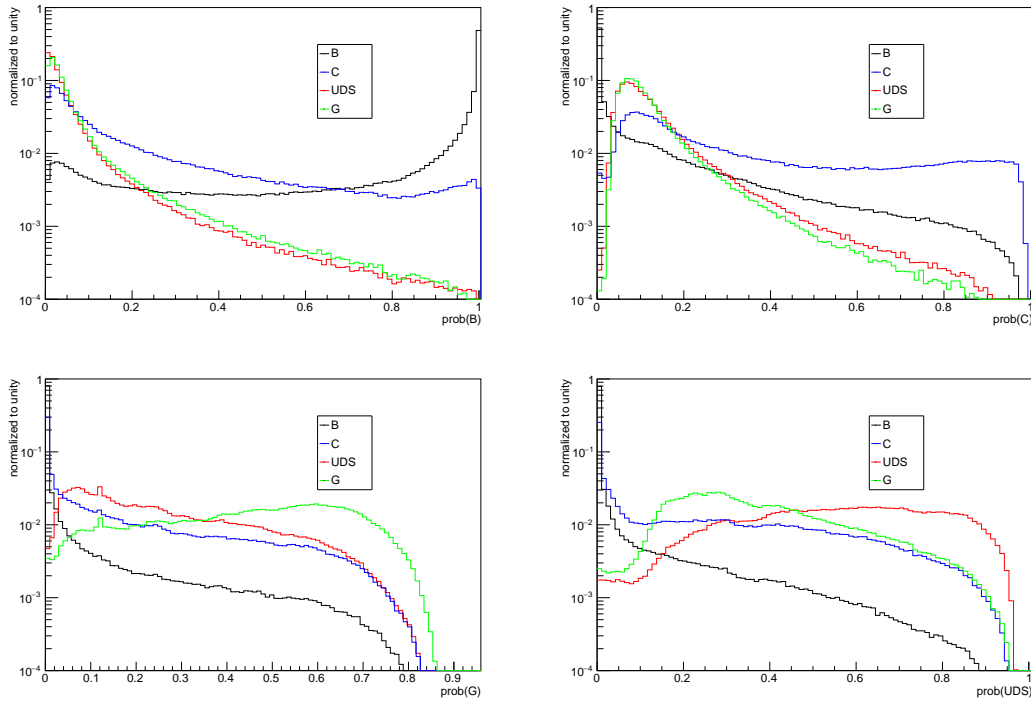


Figure 7: The different estimated probabilities for different jet labels as observed in the  $t\bar{t}$  sample.

317 for the uds probability. A comparison is shown in Figure 13, where we show how the quark-  
 318 gluon likelihood compares to DeepJet in answering the question if a jet is a uds-jet. While it  
 319 is difficult and often sub-optimal to try to extract a jet label from different standalone taggers,  
 320 it is straight forward to reduce a multiclass tagger to a tagger with only two label by rescaling  
 321 the tagger information as e.g. done in the previous paragraph for the quark-gluon separation.  
 322

323 For the results shown in Figures 8 and 9, the area under the curve is listed in Table 13. In  
 324 addition, the Table shows the efficiency  $\epsilon$  to select a light quark for two working points, defined  
 325 by a misidentification probability of 0.2 (loose), 0.1 (medium), or 0.01 (tight). The reference for  
 326 these working points is the QCD sample with  $\hat{p}_T = 80 - 120$  GeV for central  $\eta$ . The numbers are  
 327 extracted from the same sample and for the sample with  $\hat{p}_T = 300 - 470$  GeV. The conclusion  
 328 is similar for all  $\hat{p}_T$  ranges.

## 329 4.2 AK8 jet results

330 This section discusses the performance of the classification of the AK8 jets as originating from  
 331  $t$ ,  $H$ ,  $Z$ ,  $W$ , or QCD jet using the DNN structure detailed in Section 3.2. The distributions of  
 332 the individual probabilities as obtained from the DeepAK8 multi-tagger for different cases of  
 333 truth-matched jets are displayed in Fig 15.

334 The performance of the DeepAK8 multi-tagger is compared with the performance a boosted  
 335 decision tree (BDT) classifier heavily based on the  $t$  and  $W$  BDT developed in the all-hadronic  
 336 search for direct stop production [31]. Details about the selection of variables, the training,  
 337 as well as on the performance in MC and data can be found in Section 3.3 and Appendix B  
 338 in [32]. Very briefly the input variables exploit jet kinematics,  $N_{\text{subjettness}}$  ratios, soft-drop

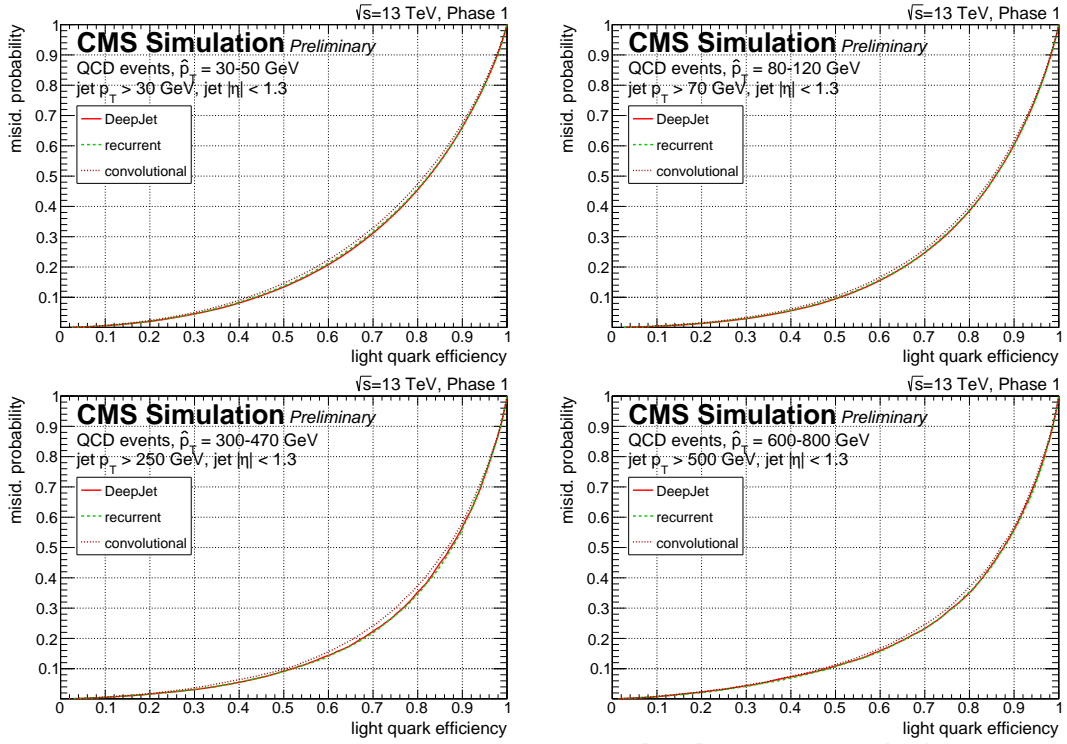


Figure 8: ROC curves in QCD simulation for different  $p_T$  ranges in the central region ( $|\eta| < 1.3$ ). Compared are the default DeepJet and the recurrent and convolutional approaches.

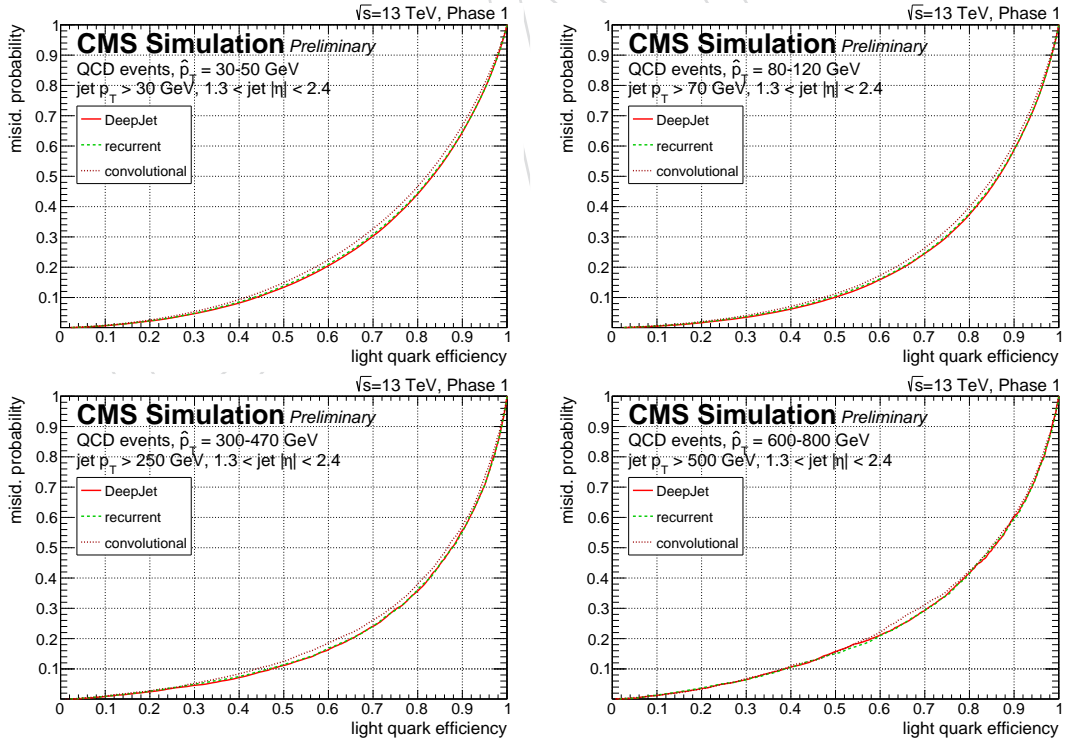


Figure 9: ROC curves in QCD simulation for different  $p_T$  ranges in the forward region ( $1.3 < |\eta| < 2.4$ ). Compared are the default DeepJet and the recurrent and convolutional approaches.

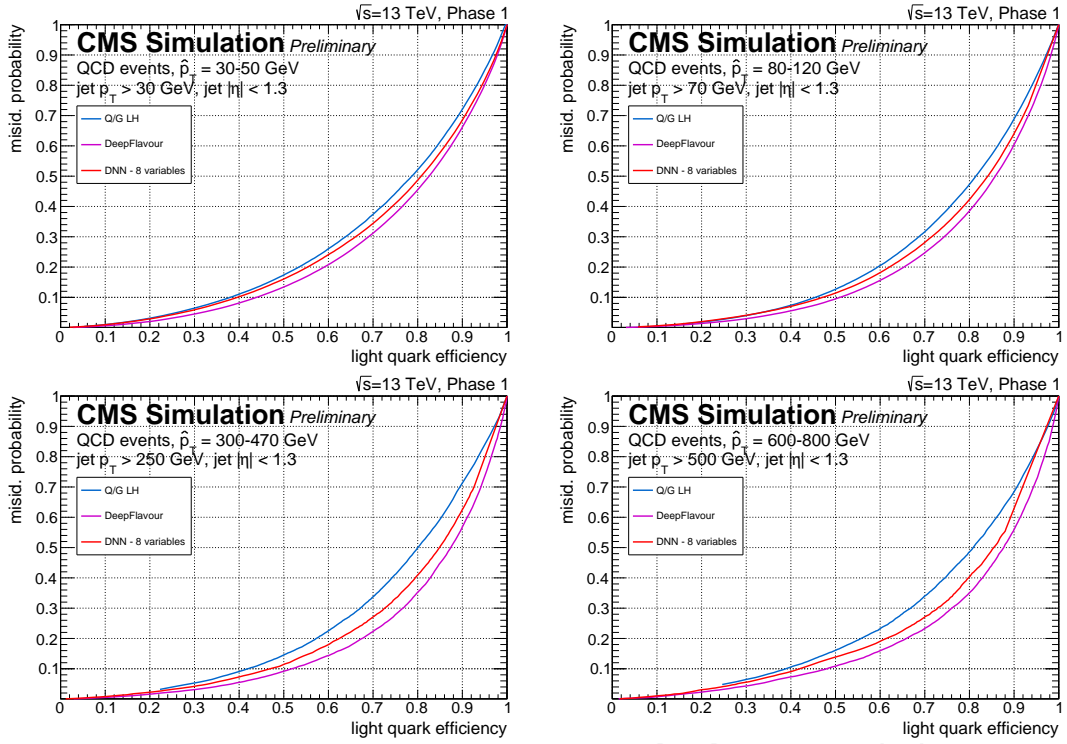


Figure 10: ROC curves in QCD simulation for different  $p_T$  ranges in the central region ( $|\eta| < 1.3$ ). Compared are the likelihood method, the 8 parameter DNN and the default DeepFlavour.

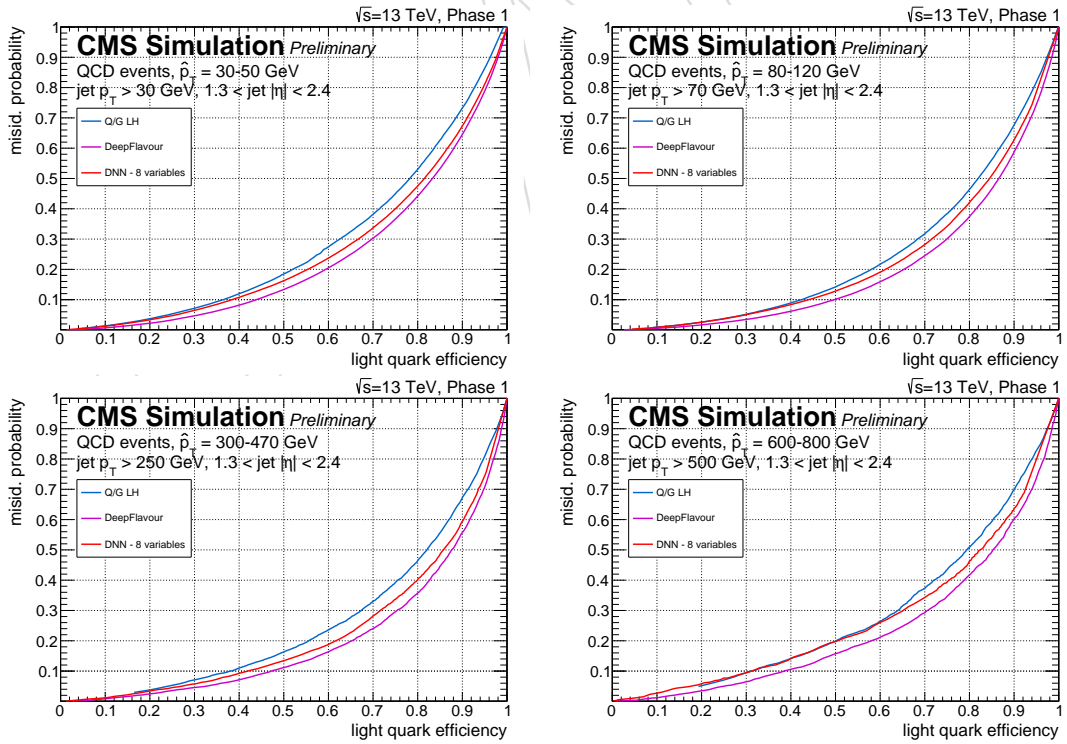


Figure 11: ROC curves in QCD simulation for different  $p_T$  ranges in the forward region ( $1.3 < |\eta| < 2.4$ ). Compared are the likelihood method, the 8 parameter DNN and the default DeepFlavour.

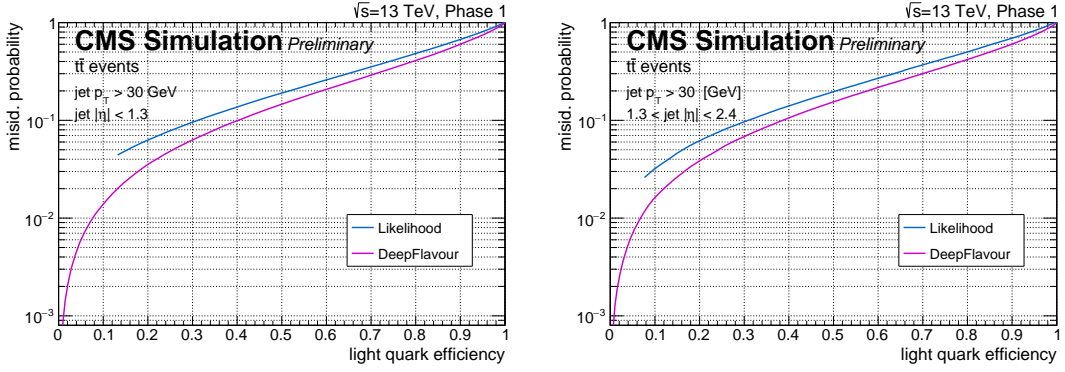


Figure 12: ROC curves in  $t\bar{t}$  simulation for  $p_T > 30$  GeV in the central ( $|\eta| < 1.3$ ) and forward region ( $1.3 < |\eta| < 2.4$ ) for light-quark and gluon separation.

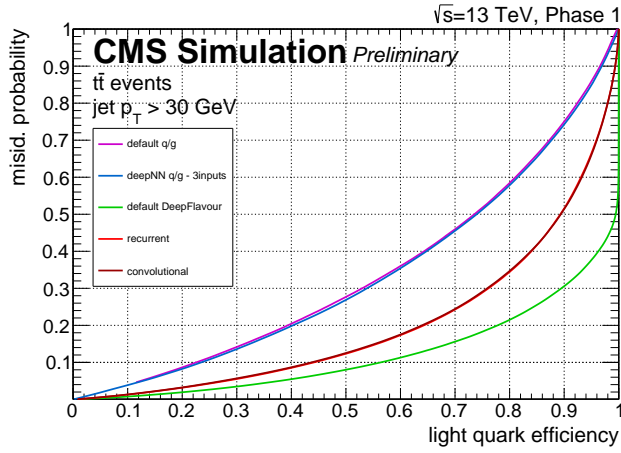


Figure 13: ROC curves in  $t\bar{t}$  simulation for  $p_T > 30$  GeV for light-quark-jets as signal and any other labeled jet as background.

Table 13: Area under the ROC curve and efficiencies for two selected working points for the different DNN-based approaches for quark-gluon tagging, evaluated for QCD samples with different  $\hat{p}_T$  and jet  $p_T$  thresholds.

	Area under ROC	$\epsilon(\text{tight})$	$\epsilon(\text{medium})$	$\epsilon(\text{loose})$
QCD $\hat{p}_T = 80 - 120$ GeV, jet $p_T > 70$ GeV				
DeepJet central	0.204	0.17	0.51	0.65
DeepJet forward	0.203	0.15	0.50	0.65
Convolution central	0.211	0.15	0.49	0.64
Convolution forward	0.215	0.13	0.47	0.63
Recurrent central	0.205	0.16	0.51	0.65
Recurrent forward	0.205	0.14	0.49	0.65
QCD $\hat{p}_T = 300 - 470$ GeV, jet $p_T > 250$ GeV				
DeepJet central	0.193	0.15	0.52	0.68
DeepJet forward	0.201	0.11	0.47	0.65
Convolution central	0.203	0.13	0.50	0.66
Convolution forward	0.214	0.10	0.44	0.62
Recurrent central	0.191	0.15	0.52	0.68
Recurrent forward	0.203	0.10	0.47	0.65

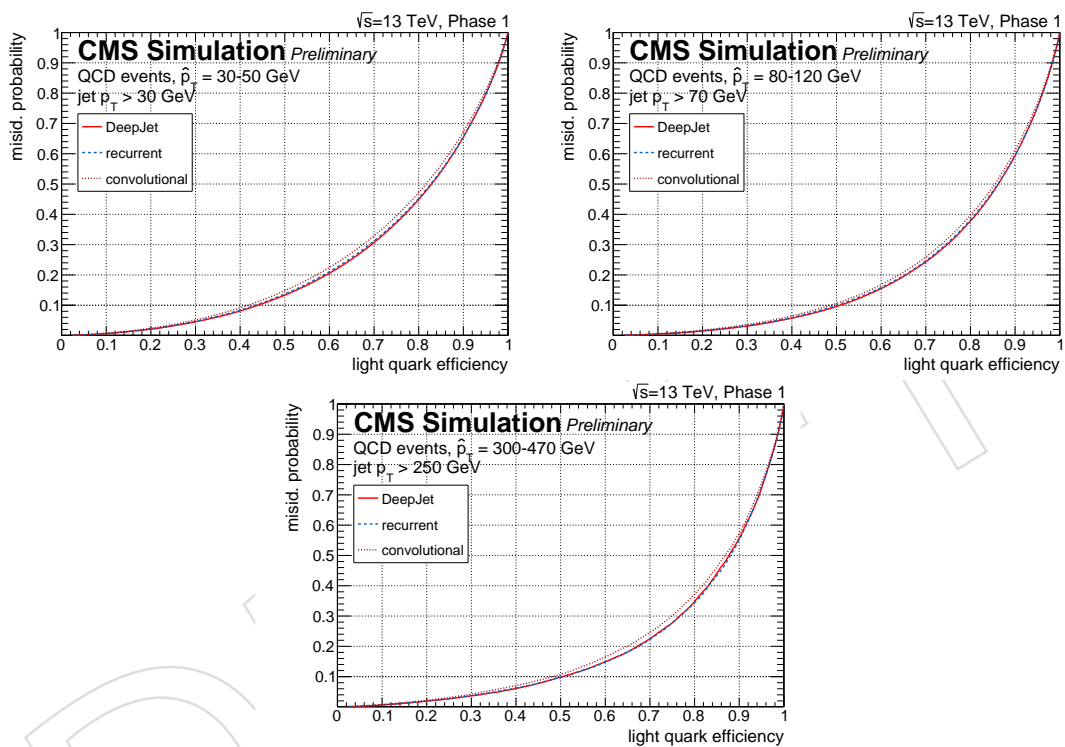


Figure 14: ROC curves in Pythia 8 QCD simulation for different  $p_T$  ranges in the full  $|\eta|$  region. Compared are the default DeepJet, the recurrent and convolutional approaches. Jets matched to uds quarks are considered light quark jets. Gluon jets are defined by a matching to gluons. Jets with heavy-flavour hadrons are excluded from the jets considered. This also applies to gluon splitting to BB and CC.

339 (SD) mass, Q/G variables and CSV discriminants of the SD subjets. Given that the DeepAK8  
340 multi-tagger targets also hadronic decays of the H and Z bosons, to allow for a fair comparison,  
341 the input variables used by the boosted double-b tagger [33] are also included in the t/W BDT.  
342 This results to a total of 45 input variables to the BDT. The BDT is retrained using the same  
343 samples as the DeepAK8 multi-tagger.

344 The performance of the DeepAK8 and BDT-based multi-taggers is evaluated in different re-  
345 gions of the  $p_T$  of the AK8 jet in terms of receiver operating characteristic (ROC) curves. An  
346 independent sample (i.e. not the sample used for train or validation) is used to produce the  
347 ROC curves. The results are displayed in Fig. 16-19. The efficiency of correctly identifying a  
348 t, H, Z, or W (signal efficiency, x-axis) is always compared against the QCD efficiency (back-  
349 ground efficiency, y-axis). The DeepAK8 multi-tagger outperforms the BDT multi-tagger in all  
350 cases. For example in a working point corresponding to a background efficiency of  $\sim 1\%$ , the  
351 DeepAK8 multi-tagger yields 10-25% larger signal efficiency in all classes.

352 One of the advantages of a multi-tagger is that allows separation between different objects. As  
353 an example, in Fig. 20 we compare the performance of the DeepAK8 and BDT multi-taggers to  
354 separate W and Z jets. This is a very challenging problem given the similar mass of the two  
355 bosons. The DeepAK8 multi-tagger shows significantly better performance over a wide range  
356 of  $p_T$ .

357 Appendix B includes comparison of the performance plots of the DeepAK8 and BDT multi-  
358 tagger using the samples and matching definitions described in [34]. In addition to the DeepAK8  
359 and the BDT multi-tagger ROC curves, we include the performance of two “cut-based” work-  
360 ing points (high and low purity) for each heavy object as described in [35] and [36].

## 361 5 Conclusion

362 For cases of our muti taggers we see significant gain with repect the the CMS reconstruction  
363 defaults taggers in the performance evaluated in simulation by ROC curves. This is true for  
364 slim jet taggin for heavy flaavours and quark gluon separation as well as for heavy resonances,  
365 H, top, W, and Z tagging for fat jet. The improvements range from a couple of % to even factors  
366 of 2 in efficiency gain at some mistag rate. The next step will be study of these gains in real  
367 data.

## 368 References

- 369 [1] J. Cogan, M. Kagan, E. Strauss, and A. Schwartzman, “Jet-Images: Computer Vision  
370 Inspired Techniques for Jet Tagging”, *JHEP* **02** (2015) 118,  
371 doi:10.1007/JHEP02(2015)118, arXiv:1407.5675.
- 372 [2] G. Kasieczka, T. Plehn, M. Russell, and T. Schell, “Deep-learning Top Taggers or The End  
373 of QCD?”, *JHEP* **05** (2017) 006, doi:10.1007/JHEP05(2017)006,  
374 arXiv:1701.08784.
- 375 [3] P. T. Komiske, E. M. Metodiev, and M. D. Schwartz, “Deep learning in color: towards  
376 automated quark/gluon jet discrimination”, *JHEP* **01** (2017) 110,  
377 doi:10.1007/JHEP01(2017)110, arXiv:1612.01551.
- 378 [4] P. Baldi et al., “Jet Substructure Classification in High-Energy Physics with Deep Neural  
379 Networks”, *Phys. Rev.* **D93** (2016), no. 9, 094034,  
380 doi:10.1103/PhysRevD.93.094034, arXiv:1603.09349.



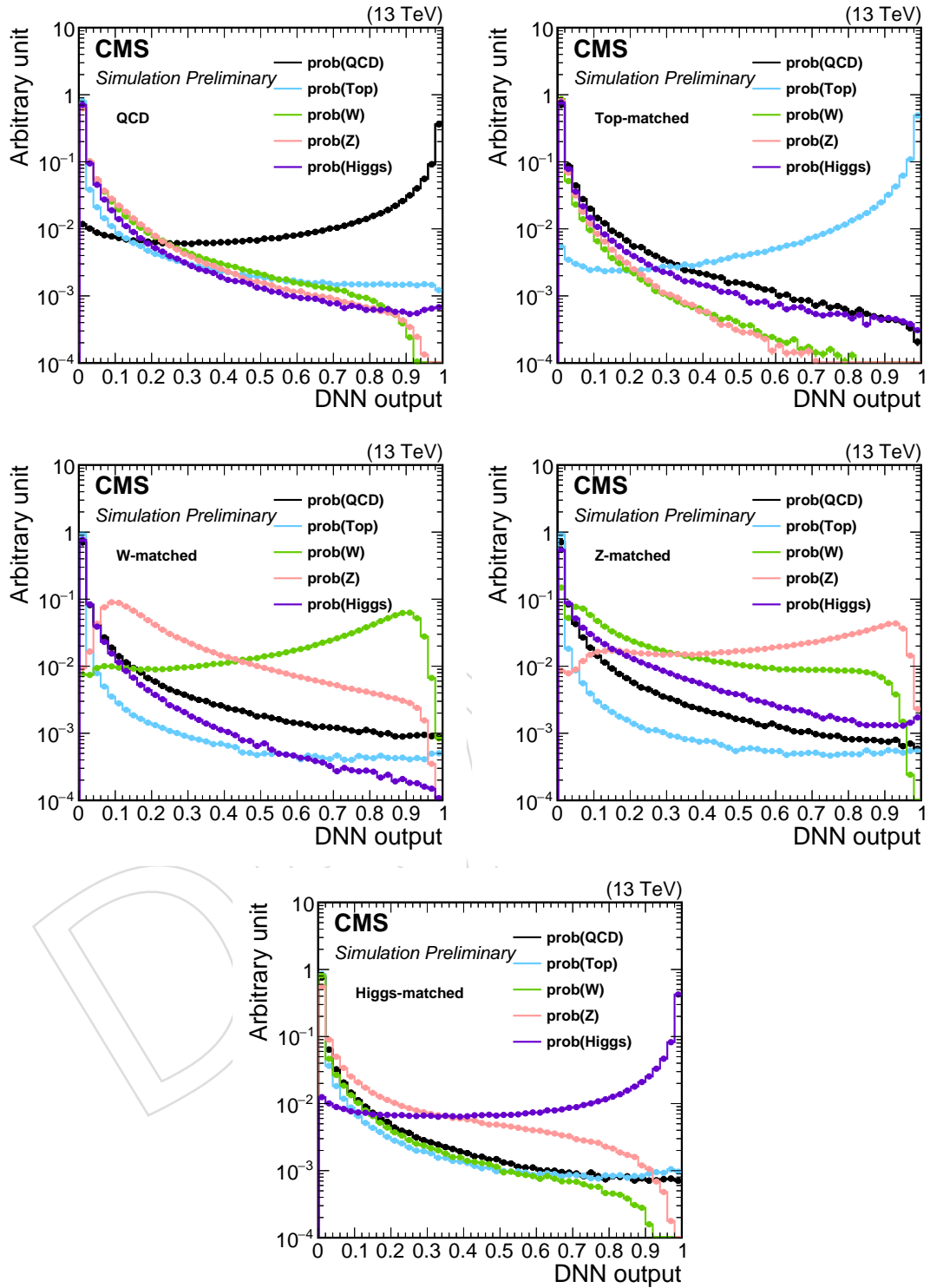


Figure 15: Distribution of the individual probabilities (DNN output) as obtained from the DeepAK8 tagger for different cases of truth-matched jets. Each truth-matched case is indicated on the plot. In this example AK8 jets with  $p_T > 400$  GeV are considered.

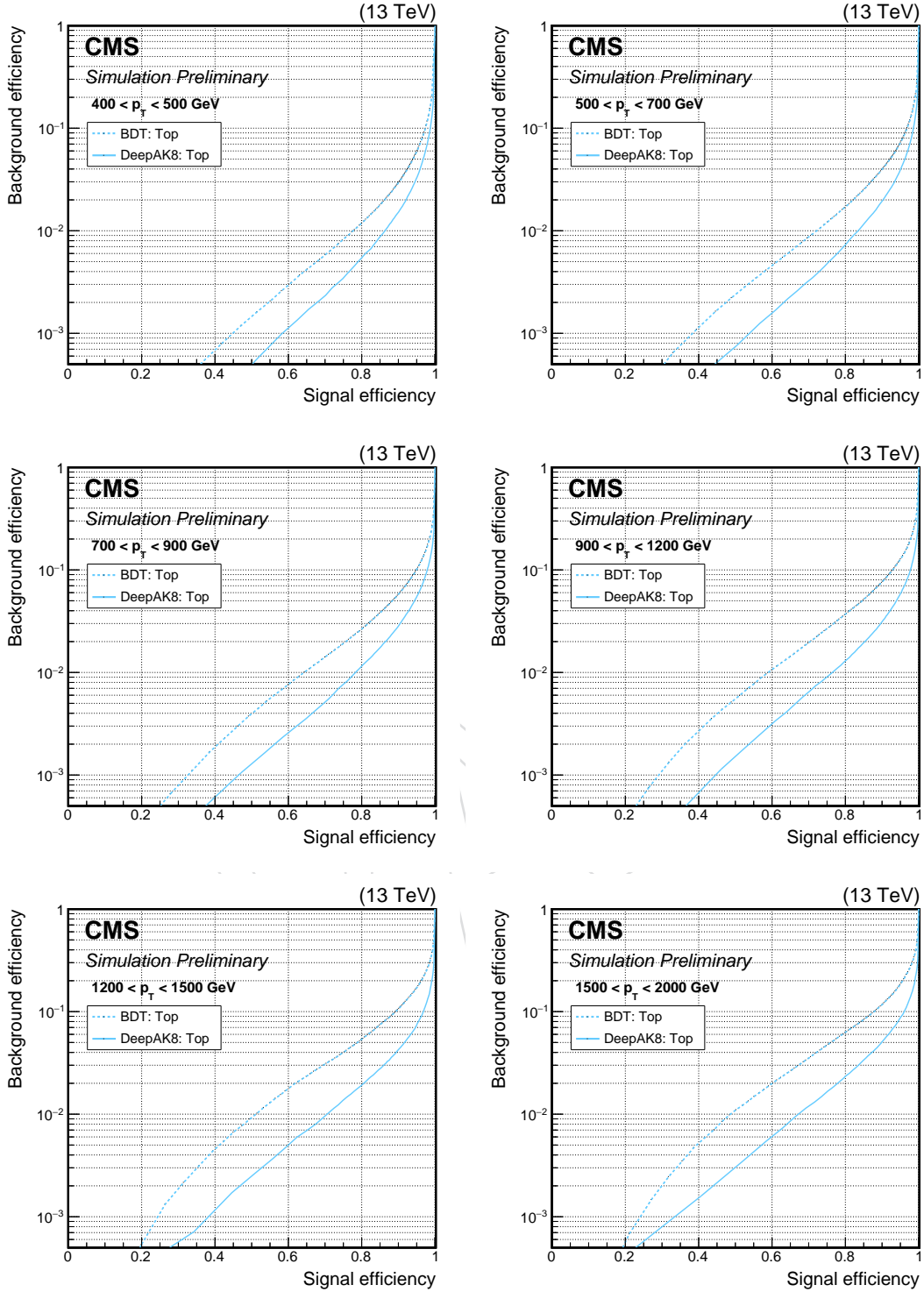


Figure 16: Comparison of the ROC curves obtained with the DeepAK8 multi-tagger (solid lines) and the BDT multi-tagger (dashed lines) in MC simulated events for  $t$  jets as signal and QCD jets as background. The plots correspond to different  $p_T$  ranges of the AK8 jet. upper:  $400 < p_T < 500$  GeV (left) and  $500 < p_T < 700$  GeV (right), middle:  $700 < p_T < 900$  GeV (left) and  $900 < p_T < 1200$  GeV (right), lower:  $1200 < p_T < 1500$  GeV (left) and  $p_T > 1500$  GeV (right).

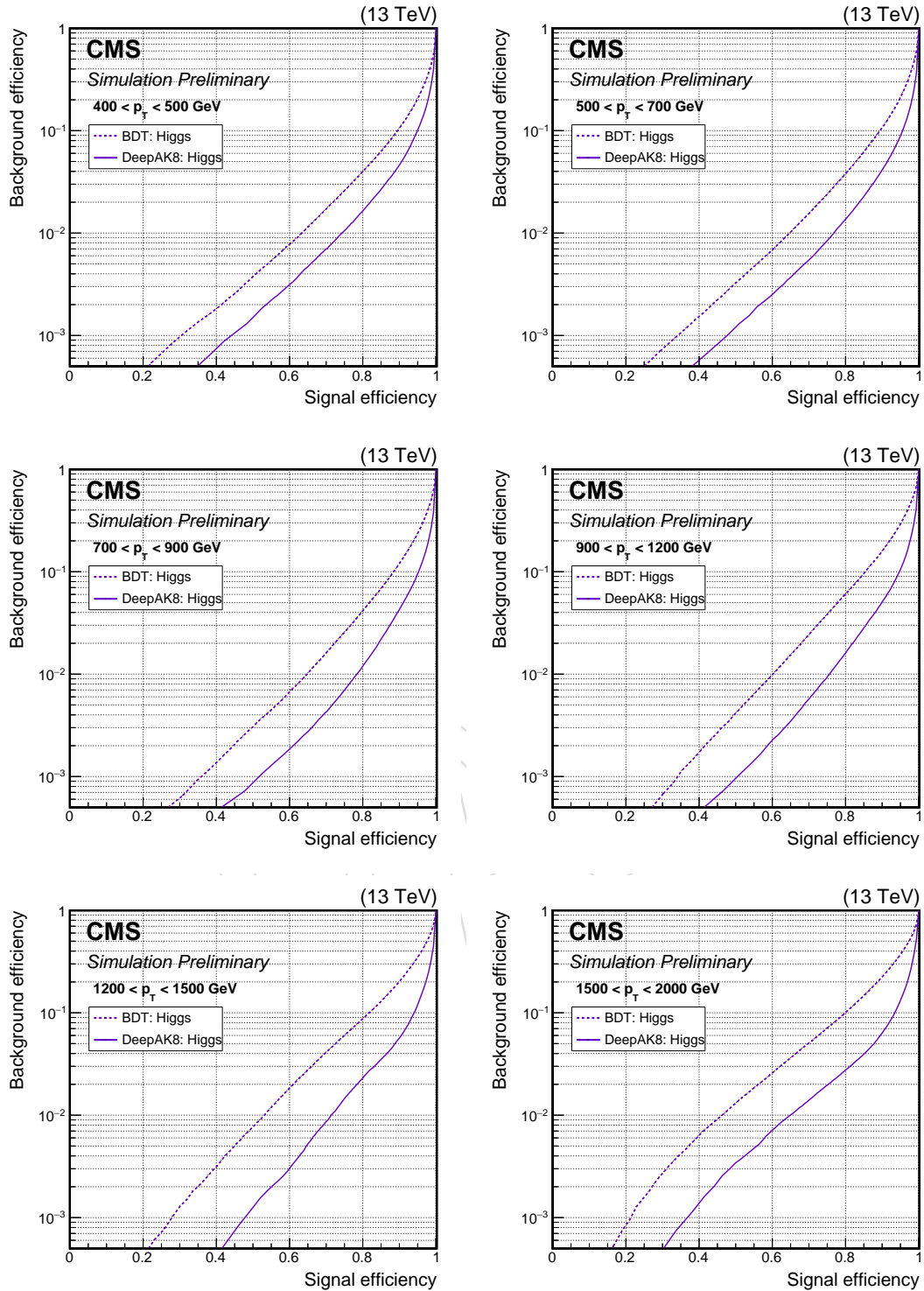


Figure 17: Comparison of the ROC curves obtained with the DeepAK8 multi-tagger (solid lines) and the BDT multi-tagger (dashed lines) in MC simulated events for H jets as signal and QCD jets as background. The plots correspond to different  $p_T$  ranges of the AK8 jet. upper:  $400 < p_T < 500$  GeV (left) and  $500 < p_T < 700$  GeV (right), middle:  $700 < p_T < 900$  GeV (left) and  $900 < p_T < 1200$  GeV (right), lower:  $1200 < p_T < 1500$  GeV (left) and  $p_T > 1500$  GeV (right).

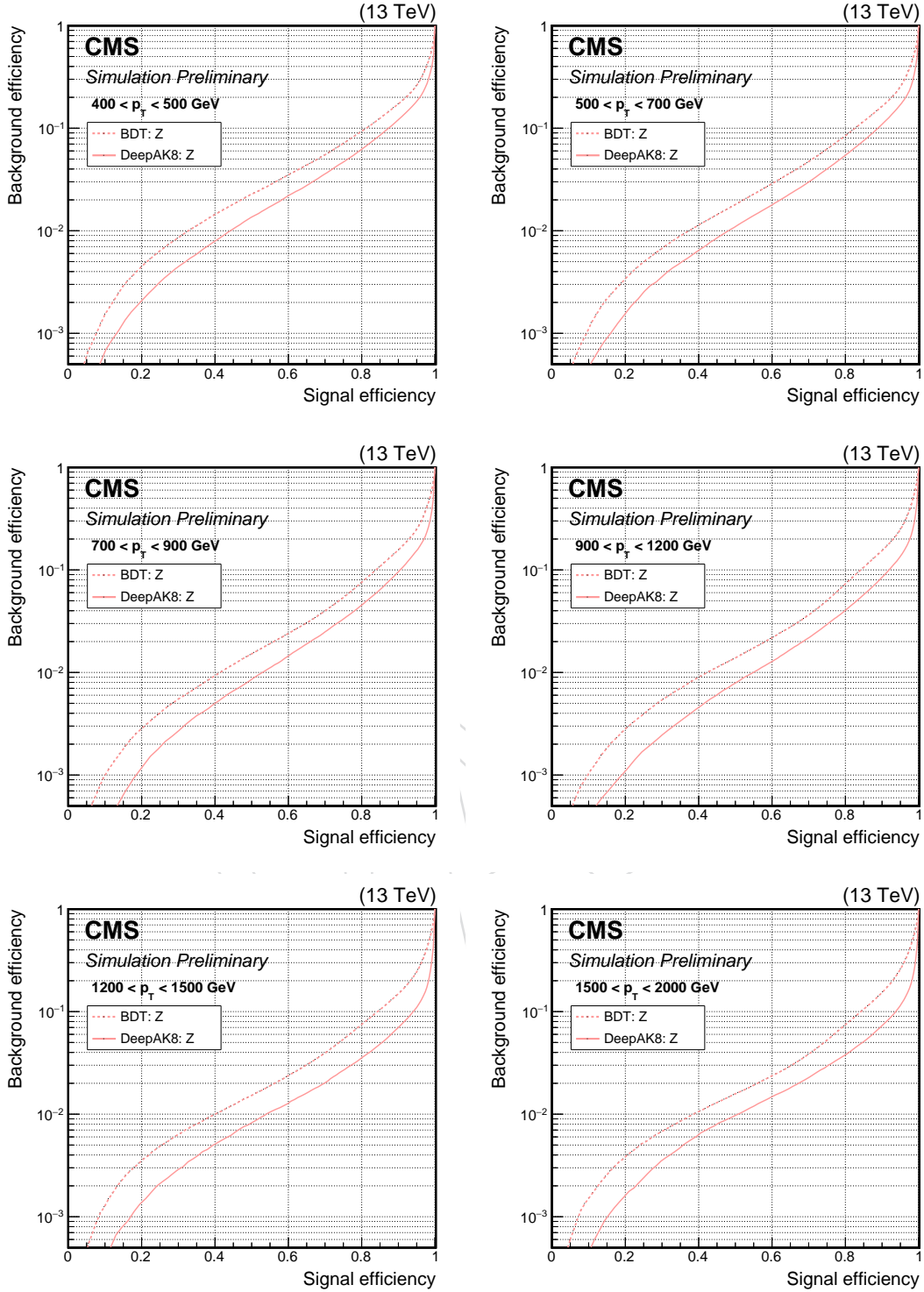


Figure 18: Comparison of the ROC curves obtained with the DeepAK8 multi-tagger (solid lines) and the BDT multi-tagger (dashed lines) in MC simulated events for Z jets as signal and QCD jets as background. The plots correspond to different  $p_T$  ranges of the AK8 jet. upper:  $400 < p_T < 500$  GeV (left) and  $500 < p_T < 700$  GeV (right), middle:  $700 < p_T < 900$  GeV (left) and  $900 < p_T < 1200$  GeV (right), lower:  $1200 < p_T < 1500$  GeV (left) and  $p_T > 1500$  GeV (right).

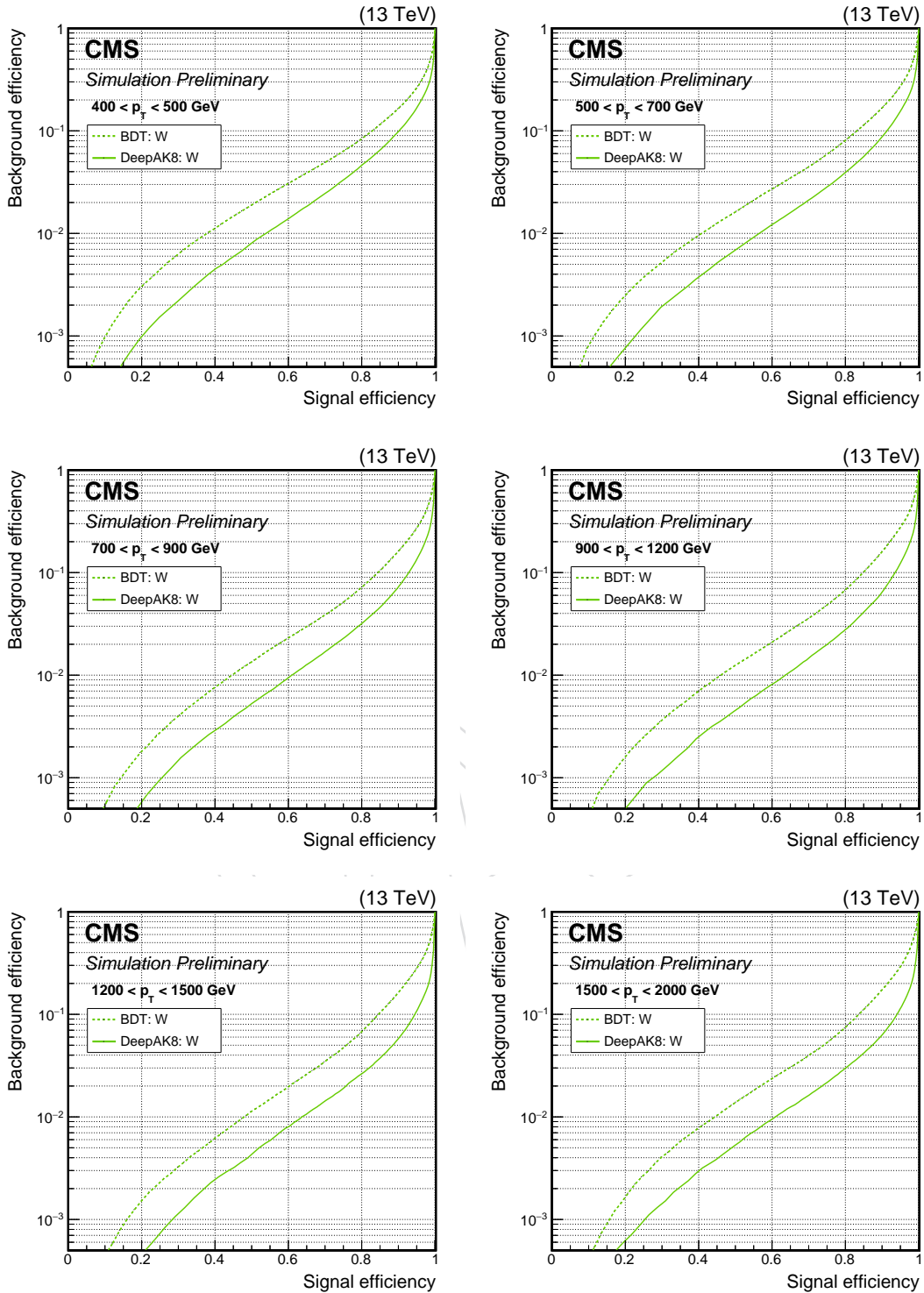


Figure 19: Comparison of the ROC curves obtained with the DeepAK8 multi-tagger (solid lines) and the BDT multi-tagger (dashed lines) in MC simulated events for W jets as signal and QCD jets as background. The plots correspond to different  $p_T$  ranges of the AK8 jet. upper:  $400 < p_T < 500$  GeV (left) and  $500 < p_T < 700$  GeV (right), middle:  $700 < p_T < 900$  GeV (left) and  $900 < p_T < 1200$  GeV (right), lower:  $1200 < p_T < 1500$  GeV (left) and  $p_T > 1500$  GeV (right).

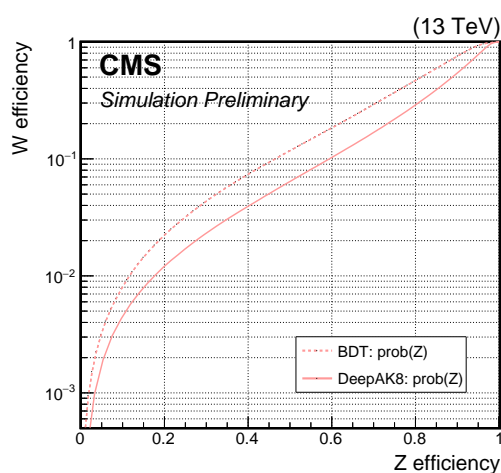


Figure 20: Comparison of the performance of the DeepAK8 (solid line) and BDT (dashed line) multi-taggers to separate W and Z jets in MC simulated events. In this example AK8 jets with  $p_T > 400$  GeV are considered.

- 381 [5] D. Guest et al., “Jet Flavor Classification in High-Energy Physics with Deep Neural  
382 Networks”, *Phys. Rev. D* **D94** (2016), no. 11, 112002,  
383 doi:10.1103/PhysRevD.94.112002, arXiv:1607.08633.
- 384 [6] G. Louppe, K. Cho, C. Becot, and K. Cranmer, “QCD-Aware Recursive Neural Networks  
385 for Jet Physics”, arXiv:1702.00748.
- 386 [7] CMS Collaboration, “CMS Phase 1 heavy flavour identification performance and  
387 developments”, CMS Detector Performance Summary CMS-DPS-17-013, 2017.
- 388 [8] ATLAS Collaboration, “Identification of Jets Containing b-Hadrons with Recurrent  
389 Neural Networks at the ATLAS Experiment”, ATLAS note ATL-PHYS-PUB-2017-003,  
390 2017.
- 391 [9] P. Nason, “A new method for combining NLO QCD with shower Monte Carlo  
392 algorithms”, *JHEP* **11** (2004) 040, doi:10.1088/1126-6708/2004/11/040,  
393 arXiv:hep-ph/0409146.
- 394 [10] S. Frixione, P. Nason, and C. Oleari, “Matching NLO QCD computations with parton  
395 shower simulations: the POWHEG method”, *JHEP* **11** (2007) 070,  
396 doi:10.1088/1126-6708/2007/11/070, arXiv:0709.2092.
- 397 [11] S. Alioli, P. Nason, C. Oleari, and E. Re, “A general framework for implementing NLO  
398 calculations in shower Monte Carlo programs: the POWHEG BOX”, *JHEP* **06** (2010) 043,  
399 doi:10.1007/JHEP06(2010)043, arXiv:1002.2581.
- 400 [12] S. Alioli, P. Nason, C. Oleari, and E. Re, “NLO single-top production matched with  
401 shower in POWHEG: s- and t-channel contributions”, *JHEP* **09** (2009) 111,  
402 doi:10.1007/JHEP02(2010)011, 10.1088/1126-6708/2009/09/111,  
403 arXiv:0907.4076. [Erratum: *JHEP* **02** (2010) 011].
- 404 [13] E. Re, “Single-top Wt-channel production matched with parton showers using the  
405 POWHEG method”, *Eur. Phys. J. C* **71** (2011) 1547,  
406 doi:10.1140/epjc/s10052-011-1547-z, arXiv:1009.2450.

- 407 [14] T. Sjöstrand et al., “An Introduction to PYTHIA 8.2”, *Comput. Phys. Commun.* **191** (2015)  
408 159, doi:10.1016/j.cpc.2015.01.024, arXiv:1410.3012.
- 409 [15] S. Agostinelli et al., “GEANT4 — a simulation toolkit”, *Nucl. Instr. and Meth. A* **506**  
410 (2003) 250, doi:10.1016/S0168-9002(03)01368-8.
- 411 [16] CMS Collaboration, “TWiki: Jet Flavour Identification (MC Truth), section on physics  
412 definition in PAT”, cms twiki, 2017.
- 413 [17] CMS Collaboration, “ParticleFlow Event Reconstruction in CMS and Performance for  
414 Jets, Taus, and  $E_T^{miss}$ ”, CMS Physics Analysis Summary CMS-PAS-PFT-09-001, 2009.
- 415 [18] CMS Collaboration, “Identification of b quark jets at the CMS Experiment in the LHC  
416 Run 2”, CMS Physics Analysis Summary CMS-PAS-BTV-15-001, 2015.
- 417 [19] D. Bertolini, P. Harris, M. Low, and N. Tran, “Pileup Per Particle Identification”, *JHEP* **10**  
418 (2014) 059, doi:10.1007/JHEP10(2014)059, arXiv:1407.6013.
- 419 [20] D. P. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization”, *CoRR*  
420 **abs/1412.6980** (2014).
- 421 [21] M. V. Jan Kieseler, Markus Stoye, “DeepJet: Framework for the development of  
422 deep-neural-network based reconstruction in high-energy-physics”, CMS Analysis Note  
423 AN 2017/126, 2017.
- 424 [22] F. Chollet, “Keras”. <https://github.com/fchollet/keras>, <https://keras.io/>,  
425 2015.
- 426 [23] M. Abadi et al., “TensorFlow: Large-Scale Machine Learning on Heterogeneous  
427 Systems”, 2015. Software available from tensorflow.org.
- 428 [24] K. He, X. Zhang, S. Ren, and J. Sun, “Identity mappings in deep residual networks”,  
429 *arXiv preprint arXiv:1603.05027* (2016).
- 430 [25] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition”, in  
431 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.  
432 2016.
- 433 [26] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by  
434 reducing internal covariate shift”, in *International Conference on Machine Learning*,  
435 pp. 448–456. 2015.
- 436 [27] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann  
437 machines”, in *Proceedings of the 27th international conference on machine learning (ICML-10)*,  
438 pp. 807–814. 2010.
- 439 [28] T. Chen et al., “Mxnet: A flexible and efficient machine learning library for heterogeneous  
440 distributed systems”, *arXiv preprint arXiv:1512.01274* (2015).
- 441 [29] CMS Collaboration, “Heavy flavor identification at CMS with deep neural networks”.
- 442 [30] CMS Collaboration, “Performance of quark/gluon discrimination in 8 TeV pp data”,  
443 CMS Physics Analysis Summary CMS-PAS-JME-13-002, 2013.
- 444 [31] CMS Collaboration, “Search for direct top squark pair production in the all-hadronic final  
445 state in proton-proton collisions at  $\sqrt{s} = 13$  TeV”, pending submission to JHEP.

- 446 [32] SUS-16-049 authors, “Search for direct production of top squark pairs in the  
 447 fully-hadronic final state with data collected in pp collisions at 13 TeV using the entire  
 448 2016 dataset”, *CMS-AN-16-437* (2016).
- 449 [33] CMS Collaboration, “Identification of double-b quark jets in boosted event topologies”.
- 450 [34] CMS Collaboration, “TWiki: Reference for JMAR Heavy Resonance Overview”, cms  
 451 twiki, 2017.
- 452 [35] CMS Collaboration, “TWiki: JMAR general page”, cms twiki, 2017.
- 453 [36] CMS Collaboration, “TWiki: Usage of b/c Tag Objects for 13 TeV Data in 2016 and 80X  
 454 MC”, cms twiki, 2017.

## 455 A Slim jet minor labels probabilities

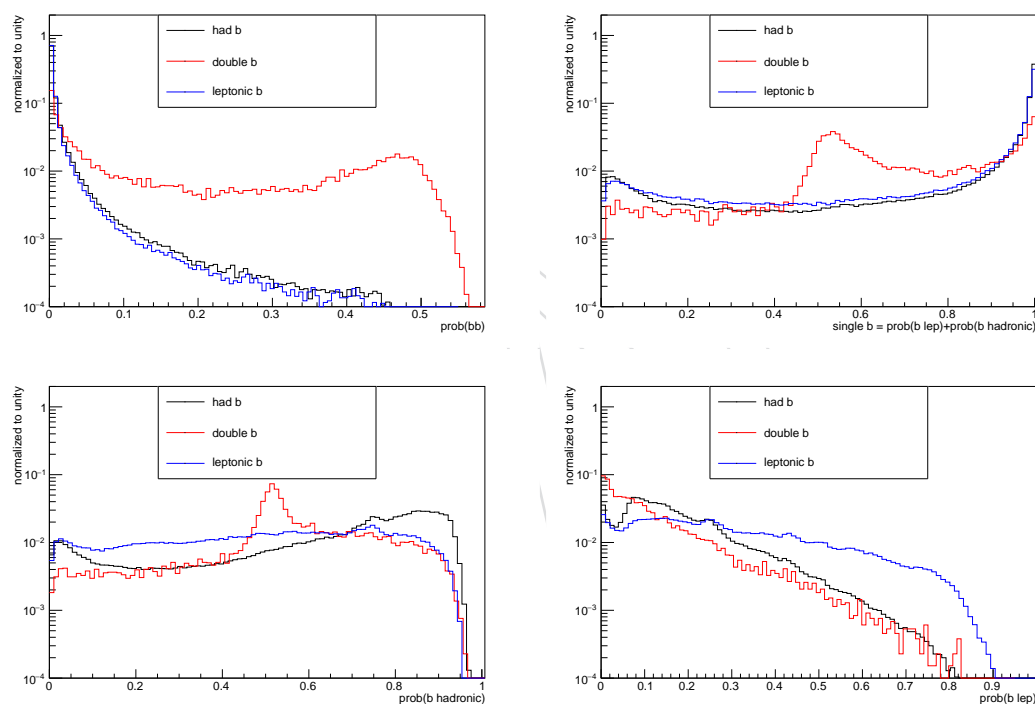


Figure 21: Minor b jet labels separately. Note that no direct lepton information is included in the input. We would not recommend using the separation of leptonic and hadronic single b-hadron jets at this point for analysis purposes.



---

456 **B Performance plots of DeepAK8 multi-tagger using the JMAR**  
457 **definition**

458 Performance plots of the DeepAK8 and BDT multi-tagger, as well as two “Cut based” working  
459 points (low and high purity) for each heavy objects as defined from the corresponding POGs  
460 [i.e. JMAR and BTV]. The ROC curves are estimated using the samples and matching definition  
461 suggested by JMAR. More details in the main text.

DRAFT

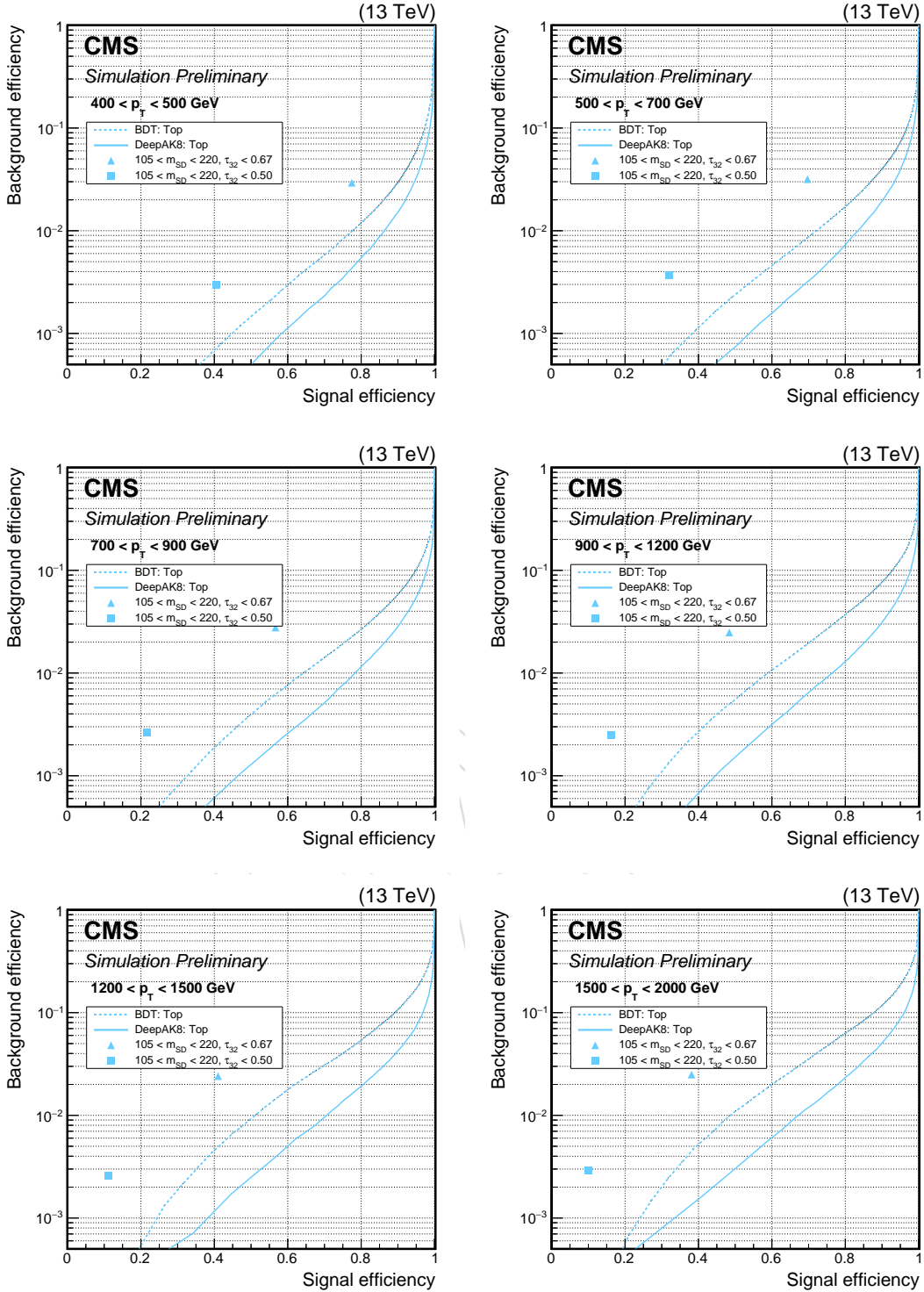


Figure 22: Comparison of the ROC curves obtained with the DeepAK8 multi-tagger (solid lines) and the BDT multi-tagger (dashed lines) in MC simulated events for  $t$  jets as signal and QCD jets as background. Two “cut-based” working points (low and high purity) are included in the plot. The working points are defined by the corresponding POG. The plots correspond to different  $p_T$  ranges of the AK8 jet. upper:  $400 < p_T < 500$  GeV (left) and  $500 < p_T < 700$  GeV (right), middle:  $700 < p_T < 900$  GeV (left) and  $900 < p_T < 1200$  GeV (right), lower:  $1200 < p_T < 1500$  GeV (left) and  $p_T > 1500$  GeV (right).

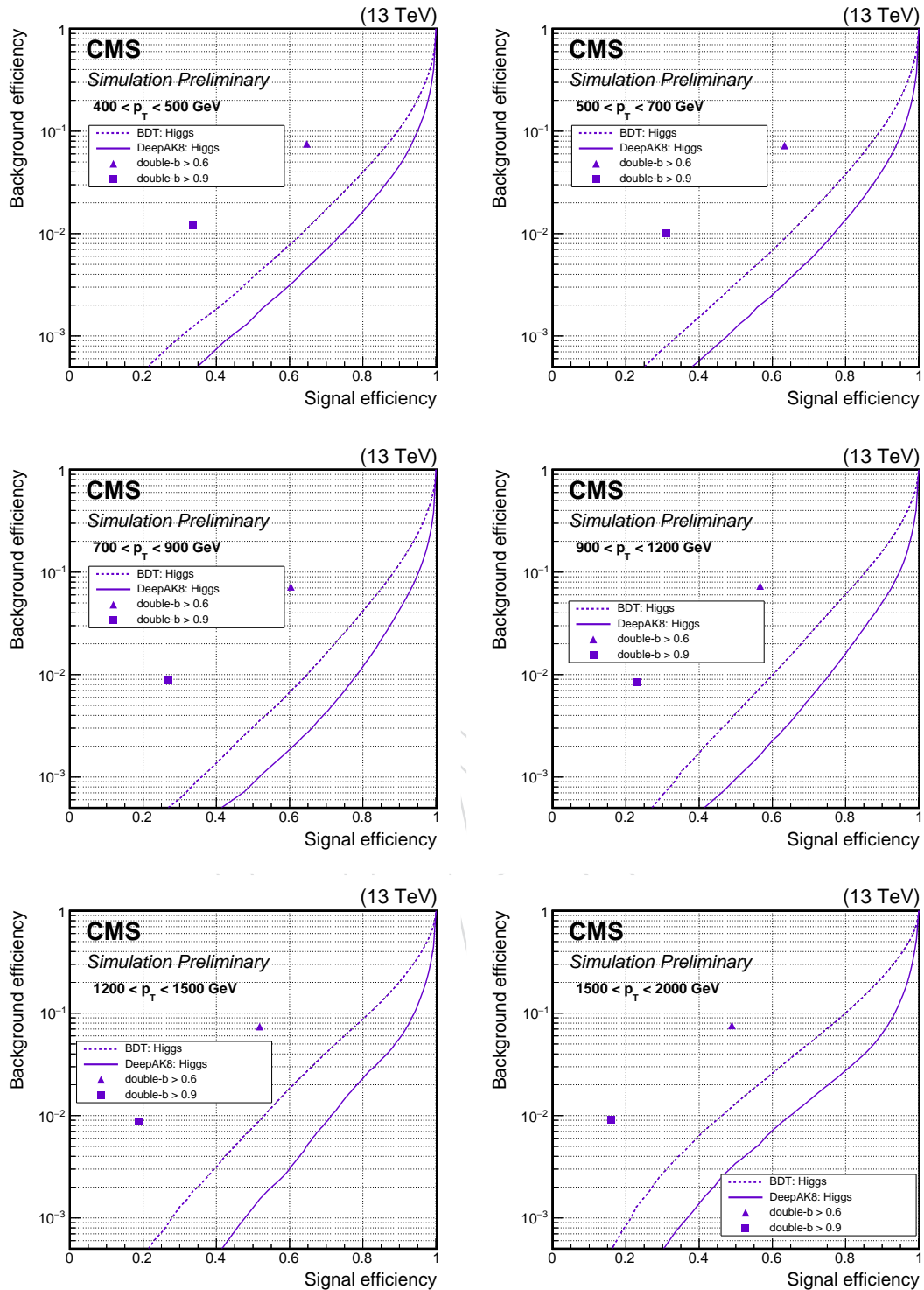


Figure 23: Comparison of the ROC curves obtained with the DeepAK8 multi-tagger (solid lines) and the BDT multi-tagger (dashed lines) in MC simulated events for H jets as signal and QCD jets as background. Two “cut-based” working points (low and high purity) are included in the plot. The working points are defined by the corresponding POG. The plots correspond to different  $p_T$  ranges of the AK8 jet. upper:  $400 < p_T < 500$  GeV (left) and  $500 < p_T < 700$  GeV (right), middle:  $700 < p_T < 900$  GeV (left) and  $900 < p_T < 1200$  GeV (right), lower:  $1200 < p_T < 1500$  GeV (left) and  $p_T > 1500$  GeV (right).

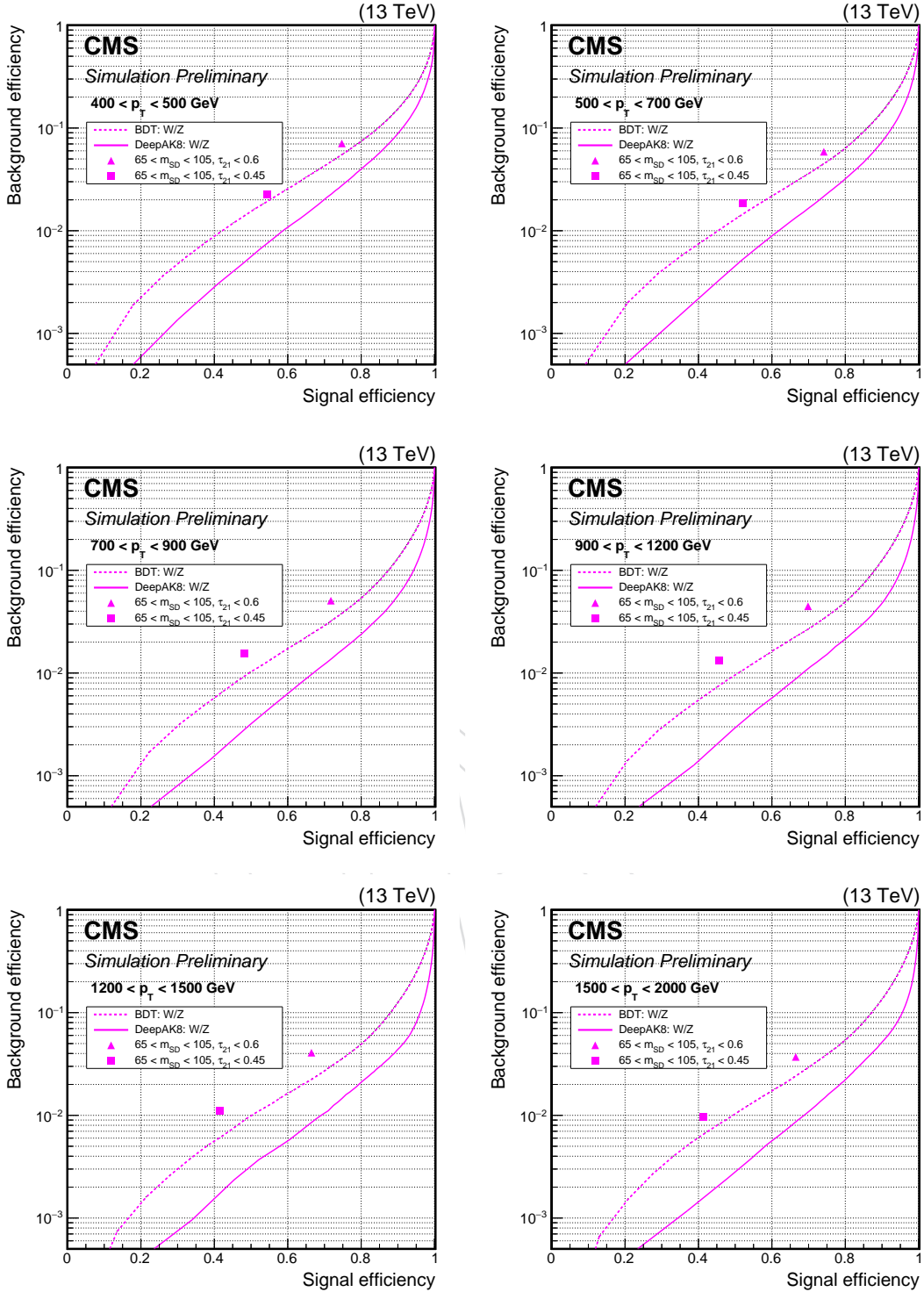


Figure 24: Comparison of the ROC curves obtained with the DeepAK8 multi-tagger (solid lines) and the BDT multi-tagger (dashed lines) in MC simulated events for Z or W jets as signal and QCD jets as background. Two “cut-based” working points (low and high purity) are included in the plot. The working points are defined by the corresponding POG. The plots correspond to different  $p_T$  ranges of the AK8 jet. upper:  $400 < p_T < 500$  GeV (left) and  $500 < p_T < 700$  GeV (right), middle:  $700 < p_T < 900$  GeV (left) and  $900 < p_T < 1200$  GeV (right), lower:  $1200 < p_T < 1500$  GeV (left) and  $p_T > 1500$  GeV (right).

## C Additional plots on the DeepAK8 multi-tagger

### C.1 Performance of the DeepAK8 multi-tagger as a function of jet $p_T$ , jet $\eta$ , and the number of primary vertices

Figure 25 to 28 show the performance of the DeepAK8 multi-tagger as a function of jet  $p_T$ , jet  $\eta$ , and the number of primary vertices in the event. We have chosen two working points based on the misidentification rate of 10% (loose) and 1% (tight). We focus on the top, W and QCD classes of the multi-classifier. The matching definition and the samples used to test the performance follow the recommendations from JMAR in [34].

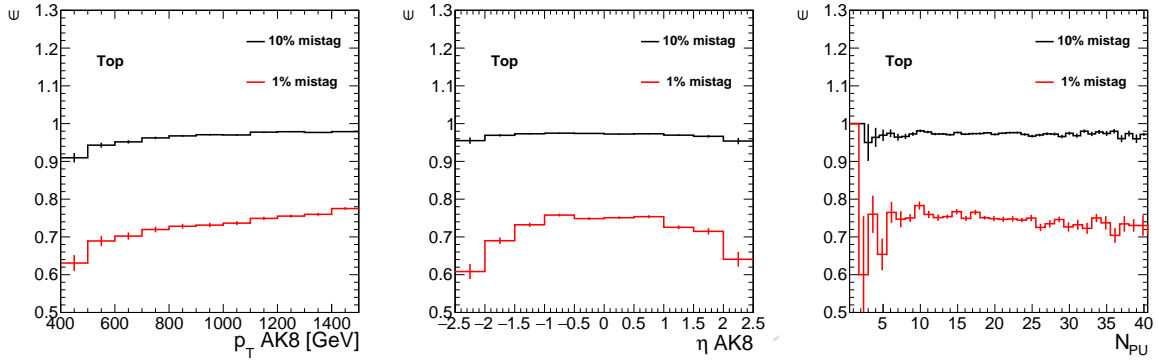


Figure 25: Efficiency of tagging truth-matched top quarks as a function of jet  $p_T$  (left), jet  $\eta$  (middle) and the number of primary vertices in the event (right).

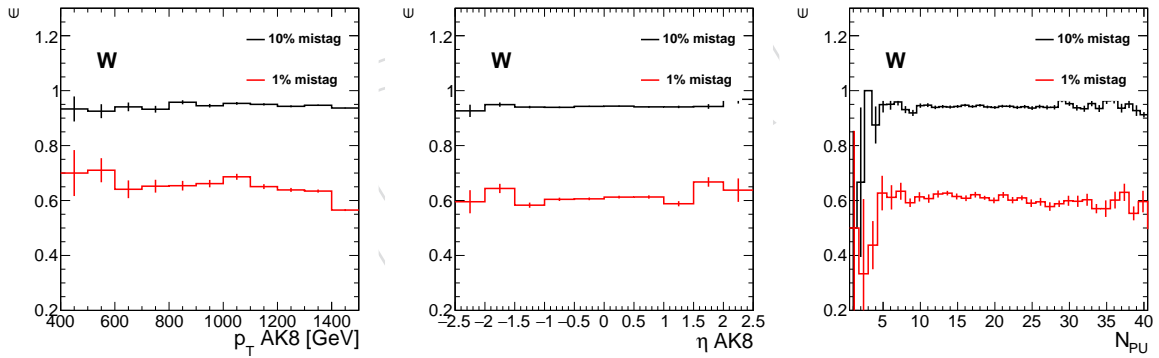


Figure 26: Efficiency of tagging truth-matched W bosons as a function of jet  $p_T$  (left), jet  $\eta$  (middle) and the number of primary vertices in the event (right).

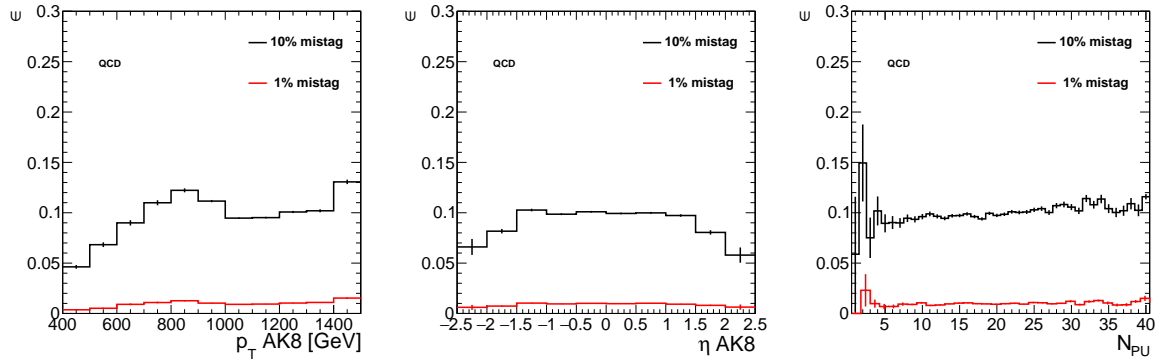


Figure 27: Rate of misidentifying QCD jets as top quarks as a function of jet  $p_T$  (left), jet  $\eta$  (middle) and the number of primary vertices in the event (right).

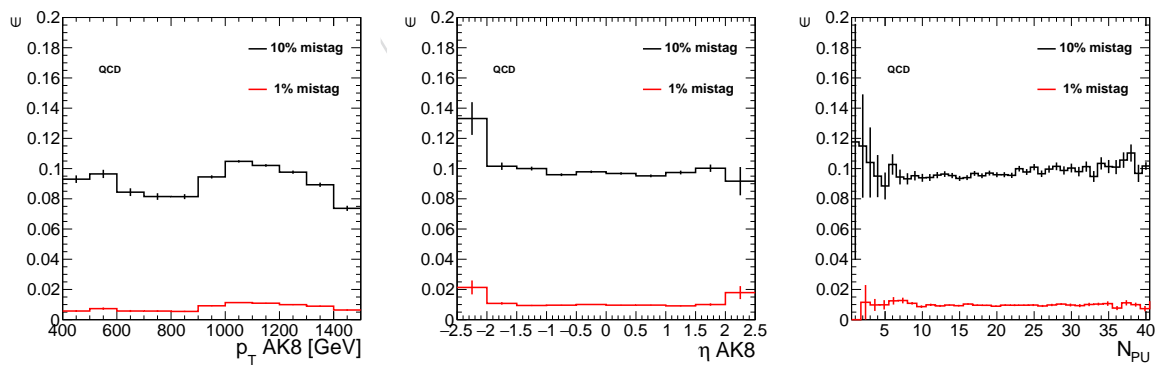


Figure 28: Rate of misidentifying QCD jets as W bosons as a function of jet  $p_T$  (left), jet  $\eta$  (middle) and the number of primary vertices in the event (right).

470 **C.2 Correlation between the DeepAK8 multi-tagger and traditional tagging vari-**  
 471 **ables**

472 Figure 29 to 32 show the correlation between the DeepAK8 multi-tagger output with traditional  
 473 jet tagging variables like the N-subjettiness ratios  $\tau_{21}$ ,  $\tau_{32}$  and the soft-drop mass. We focus on  
 474 the top,  $W$  and QCD classes of the multi-classifier. The matching definition and the samples  
 475 used to test the performance follow the recommendations from JMAR in [34].

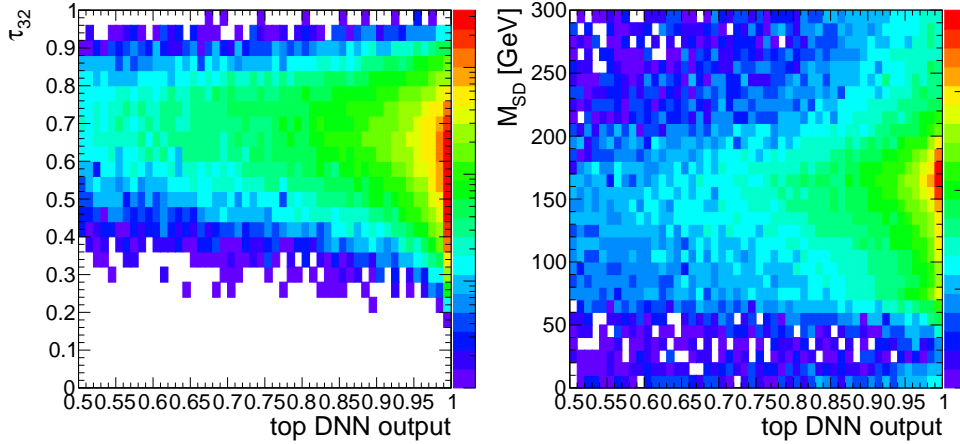


Figure 29: Correlation between the DNN output of the top class and the N-subjettiness ratio  $\tau_{32}$  (left) and the soft-drop mass  $M_{SD}$  (right) in truth-matched top jets.

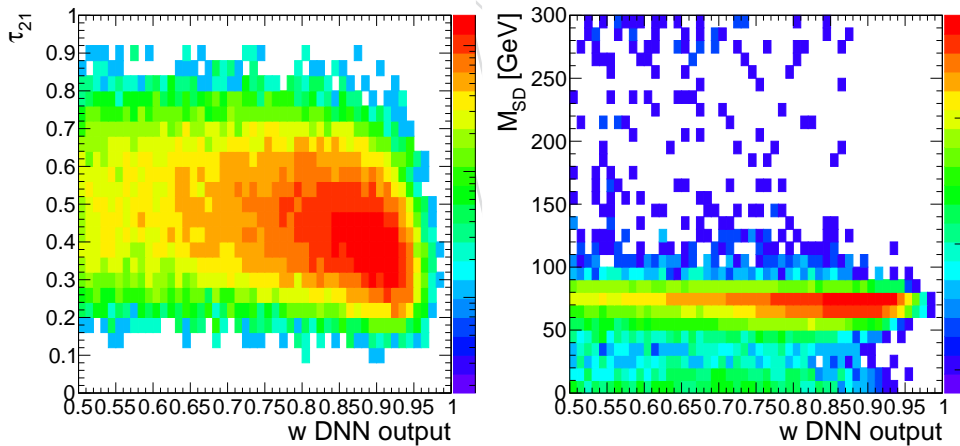


Figure 30: Correlation between the DNN output of the top class and the N-subjettiness ratio  $\tau_{21}$  (left) and the soft-drop mass  $M_{SD}$  (right) in truth-matched  $W$  jets.

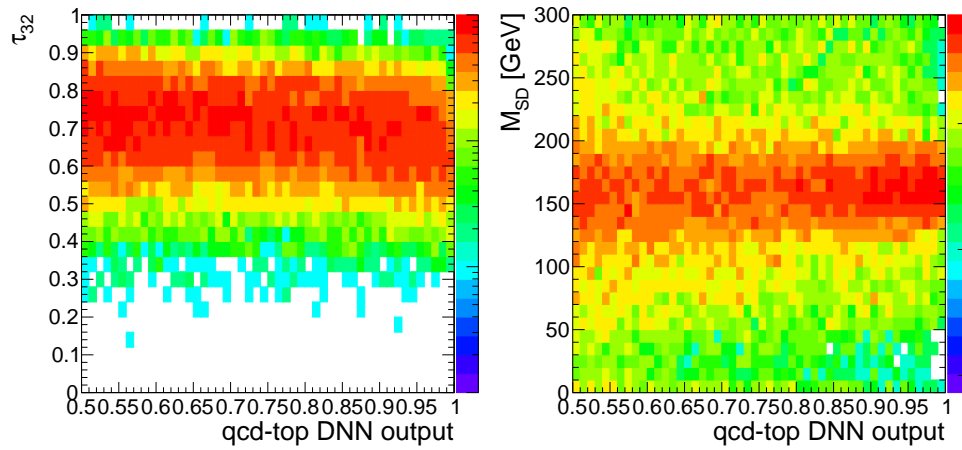


Figure 31: Correlation between the DNN output of the top class and the N-subjettiness ratio  $\tau_{32}$  (left) and the soft-drop mass  $M_{SD}$  (right) in QCD jets.

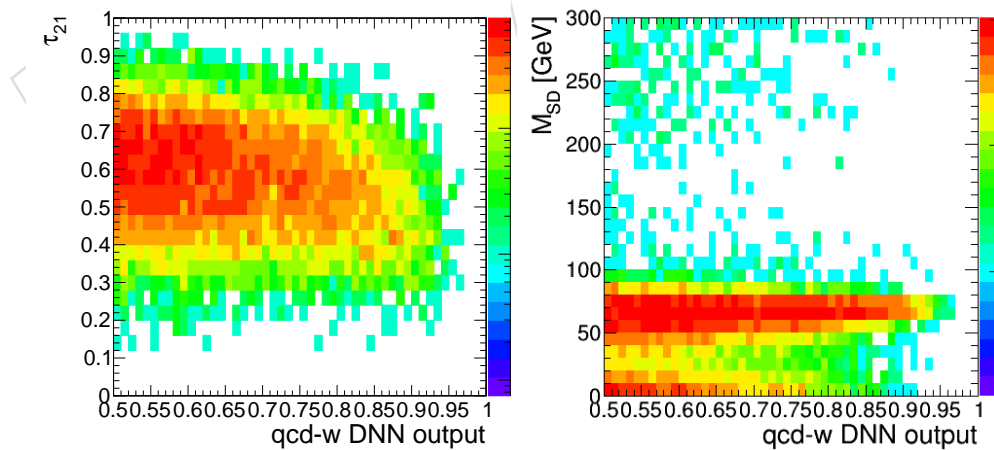


Figure 32: Correlation between the DNN output of the top class and the N-subjettiness ratio  $\tau_{21}$  (left) and the soft-drop mass  $M_{SD}$  (right) in QCD jets.



476 **C.3 Jet mass distribution**

477 Figure 33 and 34 show the soft-drop mass of the jets inclusively and after the loose and tight  
 478 working points in truth-matched jets and QCD jets, respectively.

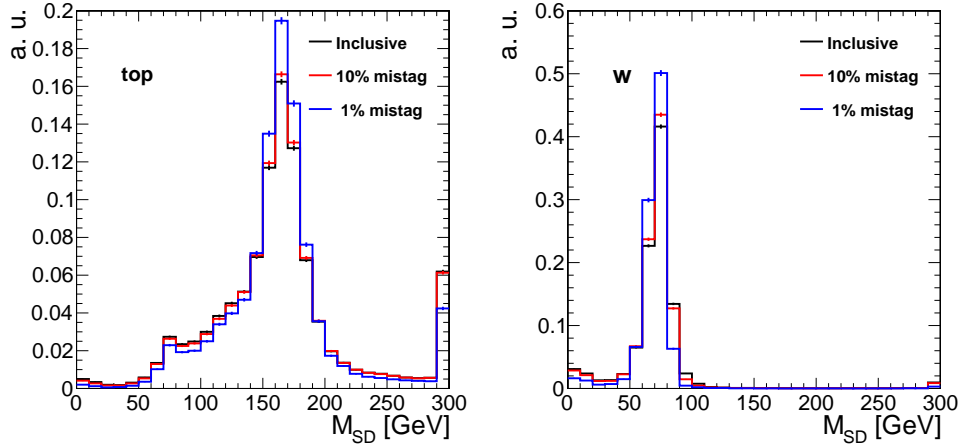


Figure 33: Soft-drop mass of the truth-matched jets inclusively (black line), after passing the loose working point (red line) and after passing the tight working point (blue line). Left: jets matched to top quarks; Right: jets matched to W bosons.

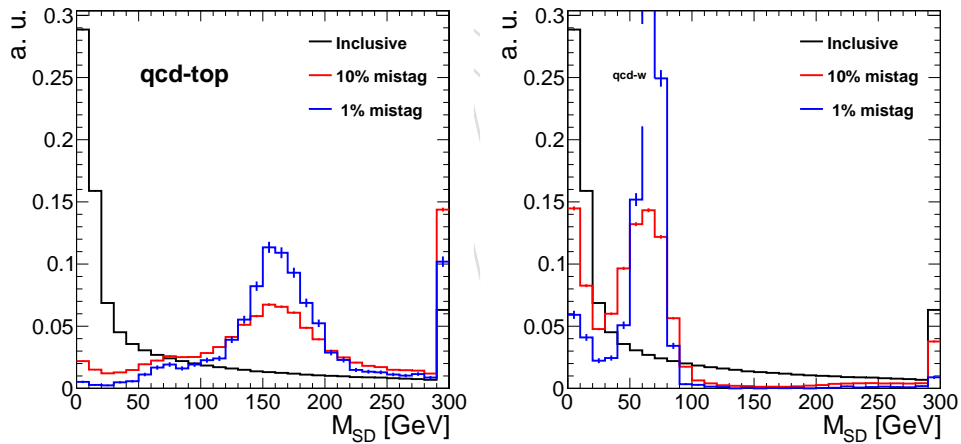


Figure 34: Soft-drop mass of the QCD jets inclusively (black line), after passing the loose working point (red line) and after passing the tight working point (blue line). Left: for top tagging; Right: for W tagging.