

# Towards an Understanding of the Correlations Within Jet Substructure

Report of BOOST2013, hosted by the University of Arizona, 12<sup>th</sup>-16<sup>th</sup> of August 2013.

D. Adams<sup>1</sup>, A. Arce<sup>2</sup>, L. Asquith<sup>3</sup>, M. Backovic<sup>4</sup>, T. Barillari<sup>5</sup>, P. Berta<sup>6</sup>,  
D. Bertolini<sup>2</sup>, A. Buckley<sup>8</sup>, J. Butterworth<sup>9</sup>, R. C. Camacho Toro<sup>10</sup>, J. Caudron<sup>9</sup>,  
Y.-T. Chien<sup>11</sup>, J. Cogan<sup>12</sup>, B. Cooper<sup>9</sup>, D. Curtin<sup>17</sup>, C. Debenedetti<sup>18</sup>,  
J. Dolen<sup>9</sup>, M. Eklund<sup>22</sup>, S. El Hedri<sup>22</sup>, S. D. Ellis<sup>22</sup>, T. Embry<sup>22</sup>, D. Ferencek<sup>23</sup>,  
J. Ferrando<sup>24</sup>, S. Fleischmann<sup>16</sup>, M. Freytsis<sup>25</sup>, M. Giuliani<sup>21</sup>, Z. Han<sup>27</sup>,  
D. Hare<sup>4</sup>, P. Harris<sup>4</sup>, A. Hinzmann<sup>4</sup>, R. Hoing<sup>4</sup>, A. Hornig<sup>22</sup>, M. Jankowiak<sup>4</sup>,  
K. Johns<sup>28</sup>, G. Kasieczka<sup>23</sup>, T. Knight<sup>24</sup>, G. Kasieczka<sup>29</sup>, R. Kogler<sup>30</sup>, W. Lampl<sup>4</sup>,  
A. J. Larkoski<sup>4</sup>, C. Lee<sup>31</sup>, R. Leone<sup>31</sup>, P. Loch<sup>31</sup>, D. Lopez Mateos<sup>27</sup>, H. K. Lou<sup>27</sup>,  
M. Low<sup>27</sup>, P. Maksimovic<sup>32</sup>, I. Marchesini<sup>32</sup>, S. Marzani<sup>32</sup>, L. Masetti<sup>33</sup>,  
R. McCarthy<sup>32</sup>, S. Menke<sup>32</sup>, D. W. Miller<sup>35</sup>, K. Mishra<sup>36</sup>, B. Nachman<sup>32</sup>, P. Nef<sup>4</sup>,  
F. T. O'Grady<sup>24</sup>, A. Ovcharova<sup>23</sup>, A. Picazio<sup>37</sup>, C. Pollard<sup>38</sup>, B. Potter Landua<sup>29</sup>,  
C. Potter<sup>29</sup>, S. Rappoccio<sup>39</sup>, J. Rutherford<sup>40</sup>, G. P. Salam<sup>10,11</sup>, J. Schabinger<sup>23</sup>,  
A. Schwartzman<sup>4</sup>, M. D. Schwartz<sup>27</sup>, B. Shuve<sup>43</sup>, P. Sinervo<sup>44</sup>, D. Soper<sup>45</sup>,  
D. E. Sosa Corral<sup>45</sup>, M. Spannowsky<sup>32</sup>, E. Strauss<sup>34</sup>, M. Swiatkowski<sup>4</sup>, J. Thaler<sup>34</sup>,  
C. Thomas<sup>34</sup>, E. Thompson<sup>1</sup>, N. V. Tran<sup>36</sup>, J. Tseng<sup>36</sup>, E. Usai<sup>36</sup>, L. Valery<sup>36</sup>,  
J. Veatch<sup>23</sup>, M. Vos<sup>23</sup>, W. Waalewijn<sup>4</sup>, and C. Young<sup>47</sup>

<sup>1</sup> Columbia University, Nevis Laboratory, Irvington, NY 10533, USA

<sup>2</sup> Duke University, Durham, NC 27708, USA

<sup>3</sup> Argonne National Laboratory, Lemont, IL 60439, USA

<sup>4</sup> SLAC National Accelerator Laboratory, Menlo Park, CA 94025, USA

<sup>5</sup> Deutsches Elektronen-Synchrotron, DESY, D-15738 Zeuthen, Germany

<sup>6</sup> Cornell University, Ithaca, NY 14853, USA

<sup>7</sup> Lund University, Lund, SE 22100, Sweden

<sup>8</sup> University of Edinburgh, EH9 3JZ, UK

<sup>9</sup> University College London, WC1E 6BT, UK

<sup>10</sup> LPTHE, UPMC Univ. Paris 6 and CNRS UMR 7589, Paris, France

<sup>11</sup> CERN, CH-1211 Geneva 23, Switzerland

<sup>12</sup> CAFPE and U. of Granada, Granada, E-18071, Spain

<sup>13</sup> McGill University, Montreal, Quebec H3A 2T8, Canada

<sup>14</sup> Iowa State University, Ames, Iowa 50011, USA

<sup>15</sup> Rutgers University, Piscataway, NJ 08854, USA

<sup>16</sup> Bergische Universitaet Wuppertal, Wuppertal, D-42097, Germany

<sup>17</sup> YITP, Stony Brook University, Stony Brook, NY 11794-3840, USA

<sup>18</sup> University of Manchester, Manchester, M13 9PL, UK

<sup>19</sup> UNESP - Universidade Estadual Paulista, Sao Paulo, 01140-070, Brazil

<sup>20</sup> INFN and University of Naples, IT80216, Italy

<sup>21</sup> University of Geneva, CH-1211 Geneva 4, Switzerland

<sup>22</sup> University of Washington, Seattle, WA 98195, USA

<sup>23</sup> Instituto de Física Corpuscular, IFIC/CSIC-UVEG, E-46071 Valencia, Spain

<sup>24</sup> University of Glasgow, Glasgow, G12 8QQ, UK

<sup>25</sup> Berkeley National Laboratory, University of California, Berkeley, CA 94720, USA

<sup>26</sup> Universidad de Buenos Aires, AR-1428, Argentina

<sup>27</sup> Harvard University, Cambridge, MA 02138, USA

<sup>28</sup> Weizmann Institute, 76100 Rehovot, Israel

<sup>29</sup> Universitaet Hamburg, DE-22761, Germany

<sup>30</sup> Universitaet Heidelberg, DE-69117, Germany

<sup>31</sup> University of Arizona, Tucson, AZ 85719, USA

<sup>32</sup> IPPP, University of Durham, Durham, DH1 3LE, UK

<sup>33</sup> Universitaet Mainz, DE 55099, Germany

<sup>34</sup> MIT, Cambridge, MA 02139, USA

<sup>35</sup> University of Chicago, IL 60637, USA

<sup>36</sup> Fermi National Accelerator Laboratory, Batavia, IL 60510, USA

<sup>37</sup> Indiana University, Bloomington, IN 47405, USA

<sup>38</sup> University of California, Davis, CA 95616, USA

<sup>39</sup> Johns Hopkins University, Baltimore, MD 21218, USA

<sup>40</sup> INFN and University of Pisa, Pisa, IT-56127, Italy

<sup>41</sup> Texas A & M University, College Station, TX 77843, USA

<sup>42</sup> INFN and University of Calabria, Rende, IT-87036, Italy

<sup>43</sup> Brown University, Richmond, RI 02912, USA

<sup>44</sup> Yale University, New Haven, CT 06511, USA

<sup>45</sup> CEA Saclay, Gif-sur-Yvette, FR-91191, France

<sup>46</sup> University of Illinois, Chicago, IL 60607, USA

<sup>47</sup> University of California, Berkeley, CA 94720, USA

**Abstract** Abstract for BOOST2013 report

**Keywords** boosted objects · jet substructure · beyond-the-Standard-Model physics searches · Large Hadron Collider

## 1 Introduction

Jet substructure has been around a while now, and it's time to study the correlations between the plethora of observables that have been developed and used. Previous BOOST reports [1, 2, 3] studied some of these things.

## 2 Monte Carlo Samples

### 2.1 Quark/gluon and $W$ tagging

Samples were generated at  $\sqrt{s} = 8$  TeV for QCD dijets and  $W^+W^-$  pairs decaying hadronically off a (pseudo) scalar resonance. The QCD events were split into sub-samples of  $gg$  and  $q\bar{q}$  events, allowing for tests of both  $W$  and quark-gluon discrimination.

Individual quark and gluon samples were produced at leading order (LO) using MADGRAPH5, while  $W^+W^-$  samples were generated using the JHU GENERATOR to allow for separation of longitudinal and transverse polarizations. Both were produced in exclusive  $p_T$  bins of 100 GeV and generated using CTEQ6L1 PDFs. The slicing parameter was chosen to be the  $p_T$  of any final state parton or  $W$ . Since no matching was performed, a cut on any parton was equivalent. These were then showered through PYTHIA8 (version 8.176) using the default tune 4C.

The showered events were clustered with FASTJET 3.03 using the anti- $k_t$  algorithm with jet radii of  $R = 0.4, 0.8, 1.2$ . A cut on the jet  $p_T$  is once again applied after showering/clustering, to ensure similar  $p_T$  spectra for signal and background in each bin.

### 2.2 Top tagging

Samples were generated at  $\sqrt{s} = 14$  TeV. Standard Model dijet and top pair samples were produced with SHERPA 2.0.0, with matrix elements with up to two extra partons matched to the shower. The top samples included only hadronic decays and were generated in exclusive  $p_T$  bins of width 100 GeV, taking as slicing parameter the maximum of the top/anti-top  $p_T$ . The QCD samples were generated with a cut on the leading

parton-level jet  $p_T$ , where parton-level jets are clustered with the anti- $k_t$  algorithm with jet radius  $R = 1.2$ . The matching scale is selected to be  $Q_{\text{cut}} = 40, 60, 80$  GeV for the  $p_{T\text{min}} = 600, 1000, \text{ and } 1500$  GeV bins, respectively.

The analysis again relies on FASTJET 3.0.3 for jet clustering and calculation of jet substructure observables, with the same cuts applied after showering and clustering as for  $\sqrt{s} = 8$  TeV data.

## 3 Jet Algorithms and Grooming Approaches

Describe the jet algorithms and grooming approaches that we will use in the report. Give the nomenclature that we will use to refer to e.g. the groomed mass in the rest of the report.

## 4 Substructure Variables/Taggers

Describe the specific substructure variables and tagging approaches that we will be using in this report e.g. n-subjettiness, Q-jets, HTT, JH tagger. Give the nomenclature that we will use to refer to these variables/taggers in the rest of the report.

## 5 Quark-Gluon Discrimination

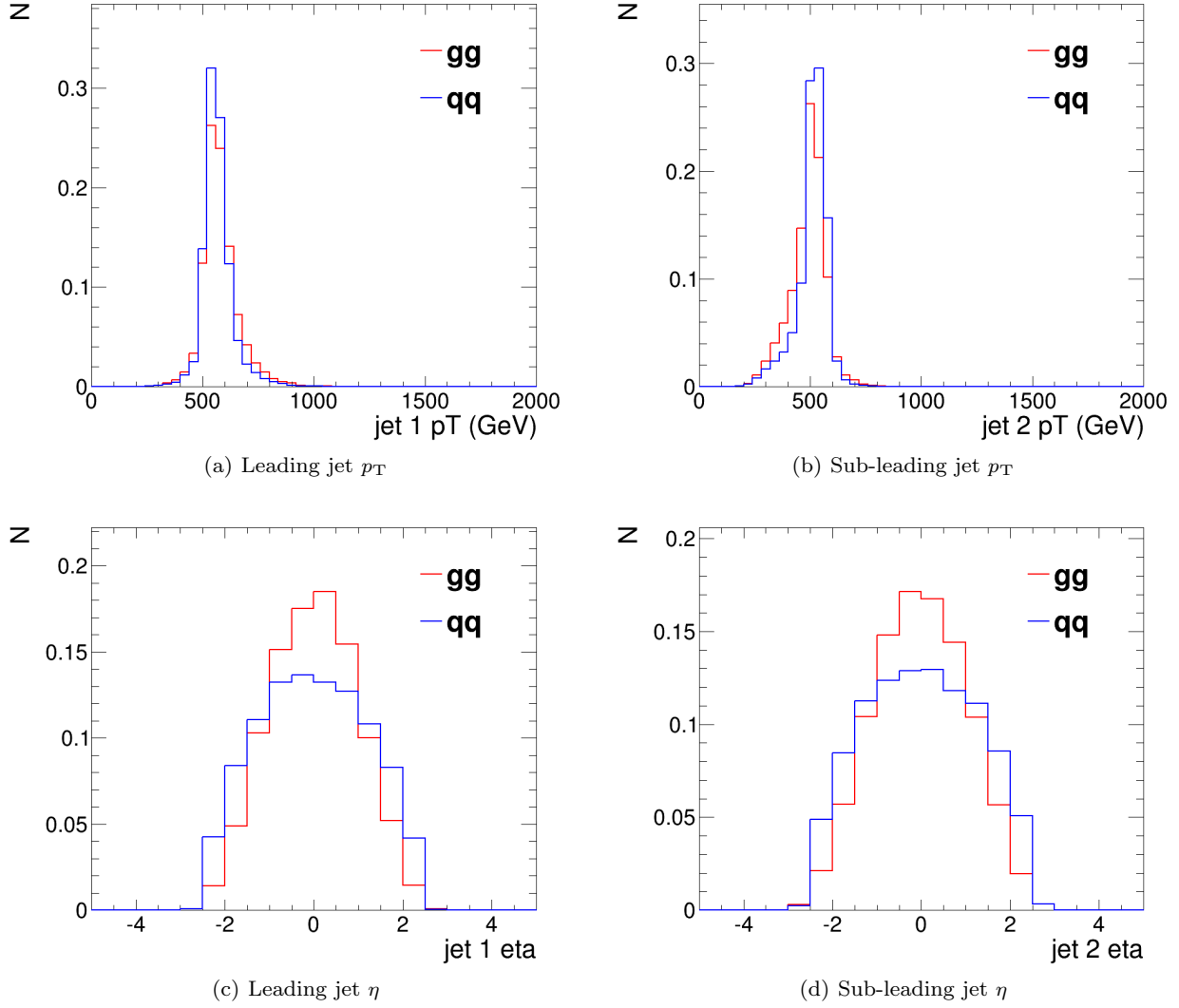
In this section we examine the differences between quark and gluon initiated jets in terms of the substructure variables, and to what extent these variables are correlated. Along the way, we attempt to provide some theoretical understanding of these observations. The motivation for these studies comes not only from the desire to “tag” a jet as being quark or gluon initiated, but also from the point of view of understanding the quark and gluon components to the QCD background to boosted boson and boosted top tagging.

### 5.1 Methodology

These studies use the  $qq$  and  $gg$  samples, described previously in Section 2.

Jets are reconstructed using the anti- $k_T$  algorithm, and have various jet grooming approaches applied, as described in Section 3. The following event selection is then applied to these samples....(presumably this will vary depending on which kinematic bin is used, as will the actual samples used - maybe summarize in a table).

Go on to explain how we produce the ROC curves, how the BDT training is done etc.



**Fig. 1** Comparisons of quark and gluon distributions in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm: basic kinematic distributions.

Figure 1 shows a comparison of the quark and gluon samples in some basic kinematic distributions.

- Dependence on  $R$ .
- Dependence on  $p_T$ .

## 5.2 Single Variable Discrimination

Figure 2 compares the quark and gluon samples in the mass distributions for the different groomers, and Figure 3 in the different substructure variables.

Figure 4 shows the single variable ROC curves in the  $p_T$  500 GeV bin for the anti- $k_T$   $R=0.8$  algorithm, compared to the ROC curve for a BDT combination of all the variables. Only the ungroomed mass is shown. One can see that the single most discriminant variables are  $n_{\text{constits}}$  and  $C_1^{\beta=0}$ .

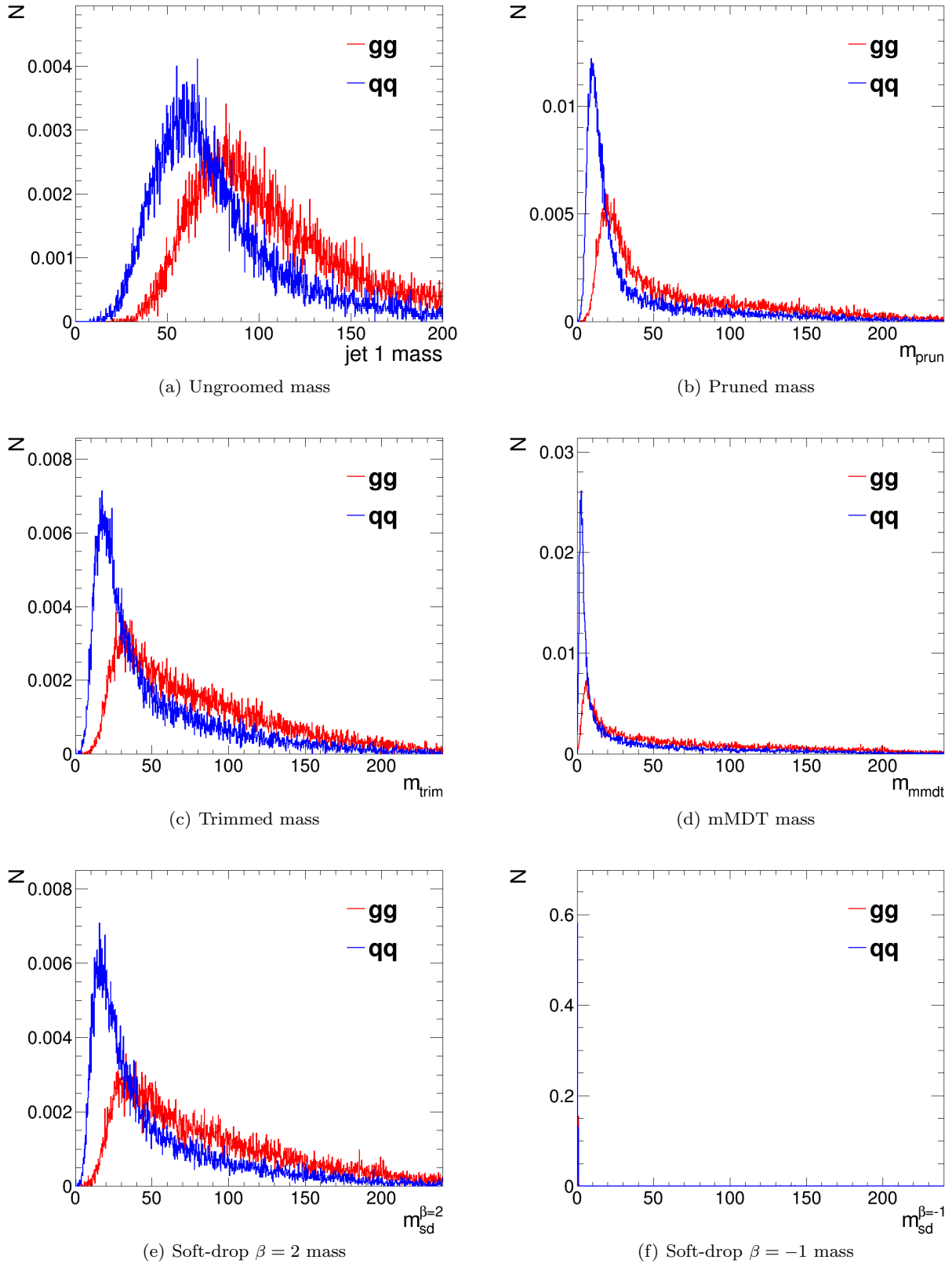
*We want to look also at:*

## 5.3 Correlations

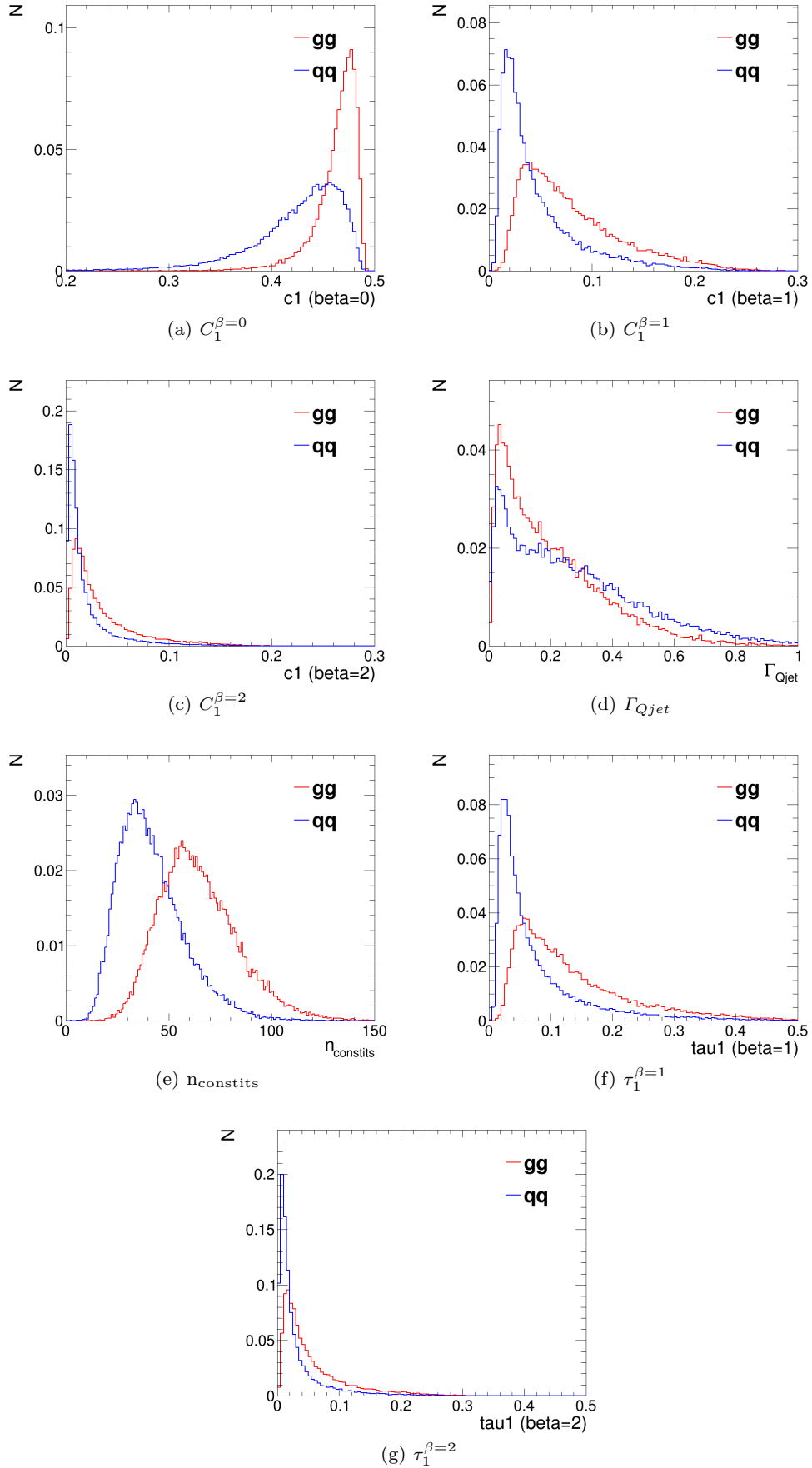
*Put in 2-D plots of correlations between variables (see theory discussions below)*

## 5.4 Combined Performance of Quark-Gluon Tagging

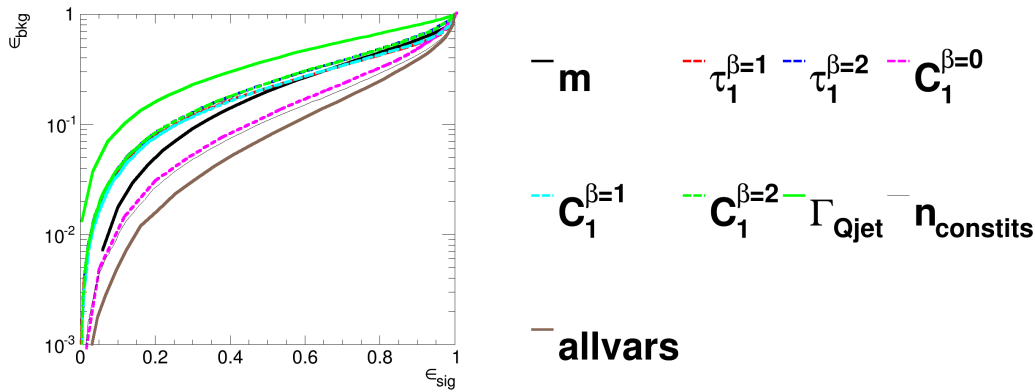
*Put in ROC curves of BDT combination of variables*



**Fig. 2** Comparisons of quark and gluon distributions in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm: leading jet mass distributions.



**Fig. 3** Comparisons of the quark and gluon distributions in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm: substructure variables.



**Fig. 4** The ROC curve for all single variables considered for quark-gluon discrimination in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm.

### 5.5 QJets Volatility and $p_T D$ ( $C_1^{(\beta=0)}$ )

Simple explanation of correlation, or why does combining volatility and  $p_T D$  improve quark versus gluon discrimination.  $p_T D$  ( $C_1^{(\beta=0)}$ ) takes small (large) values for a jet with near-democratic energy sharing between particles and large (small) values when the energy of the jet is contained in a few particles. Because we expect gluons to radiate more particles, we expect that  $p_T D_g < p_T D_q$  (or  $C_1^{(\beta=0)}_g > C_1^{(\beta=0)}_q$ ). Now, we expect the volatility of gluon jets to be in general smaller than that of quark jets because there is a greater probability (by a factor of about  $C_A/C_F = 9/4$ ) that there was a relatively hard emission in a jet that is not groomed away. By measuring both volatility and  $p_T D$ , we are sensitive to both regions of phase space: where a relatively hard emission dominates the mass of the jet as well as the region where many soft emissions set the jet mass.

*The following is Steve's discussion of volatility difference between quarks and gluons:*

Here is the (qualitative) thinking: typical QCD jet mass distributions look as illustrated on slide 17, although you should really be thinking in terms of plot versus  $m/p_T$ , since  $p_T$  is what sets the scale in the plot. Qualitatively there is a (very) large peak for  $m/p_T \lesssim 0.1$  and you should think of these jets as having masses that arise from multiple soft emissions, some of which are at substantial angles. It is these components of the jet that are operated on by pruning (reducing the mass dramatically) and that yield the large volatility tail for QCD jets. For larger  $m/p_T$  values there is typically a shoulder (my description is clearest on a semi-log plot) that runs out to about  $m/p_T \sim 0.40.5$  (where the distribution decreases rapidly). These are the QCD jets (a small fraction of the total in a given  $p_T$  bin) that contain

a hard, relatively large angle emission, which supplies the bulk of the jet mass. Such jets are effected only slightly by pruning and should exhibit much smaller volatility than the jets in the (smaller mass) peak region.

With that picture in mind and recalling that the size of the shoulder is given by low order perturbation theory (the probability of the one hard emission), we expect that the shoulder will be higher for gluons than for quarks (essentially by the usual  $C_A/C_F$  color charge factor), as suggested by the lower right plot on slide 17. Since the shoulder presumably plays a more important role for gluons (since it is larger), one would expect that the volatility distribution for gluons is narrower than quarks, as suggested in the upper left plot on slide 17. Am I making sense?

On the other hand, the volatility distribution plot indicates that the Q vs G distributions for your cuts are not really very different, which is presumably why it is not a very good discriminant by itself. But I expect this to depend in detail on where we are operating on the  $m/p_T$  distributions. This leads to my request above. Your  $p_T$  bin is pretty broad and I don't expect the q and g samples to have the same shape within the bin. Of course, this may not be an issue, but I would like to check.

### 5.6 Comparison of Groomed Jet Masses

## 6 Boosted $W$ -Tagging

In this section we study the performance of various groomed jet masses, substructure variables, and BDT combinations of groomed mass and substructure, in terms of the identification of a boosted hadronically decaying  $W$  signal against a gluon-gluon background. We produce

Receiver Operating Characteristic (ROC) curves that elucidate the performance of the various groomed mass and substructure variables that are capable of providing discrimination between signal and background. A range of different distance parameter settings for the anti- $k_T$  jet algorithm are explored, in a variety of kinematic regimes (lead jet  $p_T$  200-300 GeV, 500-600 GeV, 1.0-1.1 TeV), to explore the performance as a function of jet radius and jet boost, and to see where substructure approaches may break down. The groomed mass and substructure variables are then combined in a Boosted Decision Tree (BDT), and the performance of the resulting BDT discriminant explored through ROC curves to understand the degree to which variables are correlated and exploiting the same information, and how this changes with jet boost and jet radius.

## 6.1 Methodology

These studies use the  $X \rightarrow WW$  samples as signal and the  $gg$  samples to model the QCD background, described previously in Section 2. Whilst only gluonic backgrounds are explored here, the conclusions as to the dependence of the performance and correlations on the jet boost and radius have been verified to hold also for  $qq$  backgrounds. *To be checked!*

Jets are reconstructed using the anti- $k_T$  algorithm, and have various jet grooming approaches applied, as described in Section 3. The following event selection is then applied to these samples....(presumably this will vary depending on which kinematic bin is used, as will the actual samples used - maybe summarize in a table).

Figure 5 shows a comparison of the leading jet  $p_T$  for the signal and background in the  $p_T$  300-400 GeV bin, for the two different anti- $k_T$  jet algorithm distance parameters explored in this bin ( $R=0.8$  and  $R=1.2$ ). Figures 6 and 7 show the same for the  $p_T$  500-600 GeV bin and  $p_T$  1.0-1.1 TeV bin respectively, where for the  $p_T$  1.0-1.1 TeV bin the distance parameter  $R=0.4$  is also explored.

Go on to explain how we produce the ROC curves, how the BDT training is done etc.

## 6.2 Single Variable Performance

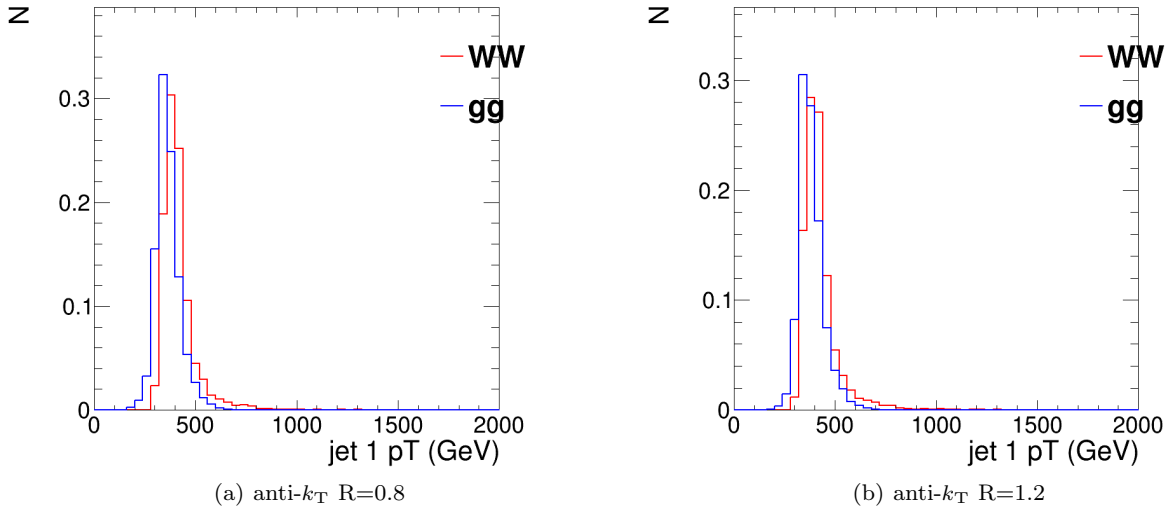
In this section we will explore the performance of the various groomed jet mass and substructure variables in terms of discriminating signal and background, and how this performance changes depending on the kinematic bin and jet radius considered.

Figure 8 compares the signal and background in terms of the different groomed masses explored for the

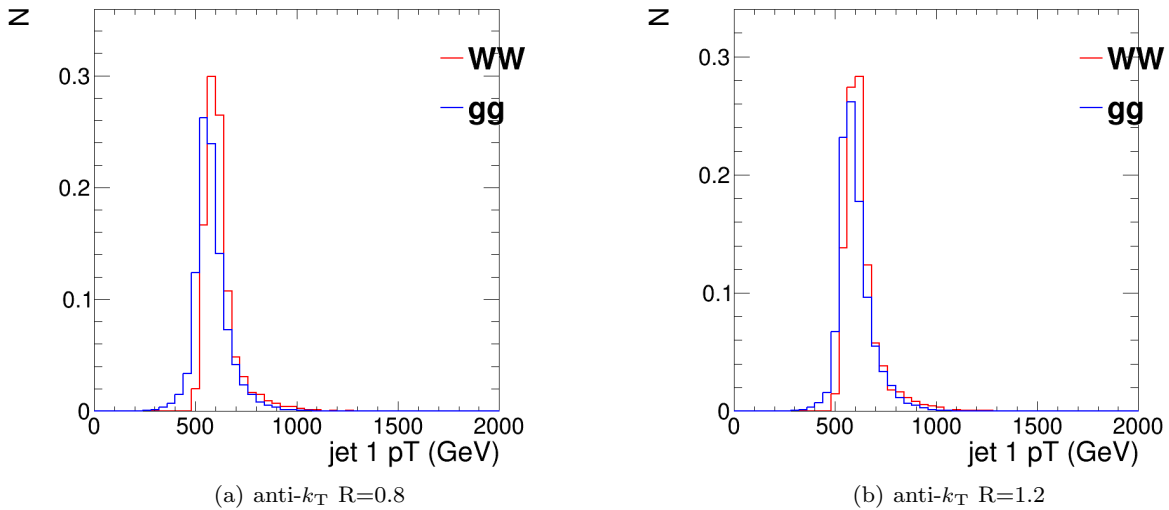
anti- $k_T$   $R=0.8$  algorithm in the  $p_T$  500-600 bin. One can clearly see that in terms of separating signal and background the groomed masses will be significantly more performant than the ungroomed anti- $k_T$   $R=0.8$  mass. *Need to comment on the soft drop  $B=-1$  mass here* Figure 9 compares signal and background in the different substructure variables explored for the same jet radius and kinematic bin.

Figures 10,11 and 12 show the single variable ROC curves compared to the ROC curve for a BDT combination of all the variables (labelled “allvars”), for each of the anti- $k_T$  distance parameters considered in each of the kinematic bins. One can see that, in all cases, the “allvars” option is considerably more performant than any of the individual single variables considered, indicating that there is considerable complementarity between the variables, that will be explored further in the next section. The best performant individual variables for a reasonable signal efficiency are the groomed masses, which all have a similar level of performance that is superior to that of any of the substructure variables considered.

Although the ROC curves give all the relevant information, it is hard to compare performance quantitatively. In Figures 13,14 and 15 matrices are shown which give the background rejection for a signal efficiency of 50% when two variables (that on the x-axis and that on the y-axis) are combined in a BDT. Thus, the diagonal of these plots can be examined to see quantitatively the individual single variable performance. Because we have not attempted to optimise the grooming parameter settings of each grooming algorithm, we do not want to place too much emphasis here on the relative performance of the groomed masses, but instead look at the trends versus  $p_T$  and  $R$ . One can see clearly that the background rejection power of the groomed mass variables increases as the  $p_T$  is increased. Within a  $p_T$  bin, one can also see that the groomed mass performance is rather invariant to changes in the jet radius. In contrast, the substructure variable performance varies considerably as the jet radius is changed. In general, the background rejection power of individual jet substructure variables gets worse as the jet radius is increased. The only exception to this is in the highest  $p_T$  bin, where the background rejection power of  $C_2^{\beta=1}$  improves when going from jet radius  $R=0.4$  to  $R=0.8$ , but then gets worse again as we go to  $R=1.2$ . *Insert some nice discussion/explanation of why jet substructure power generally gets worse as we go to large jet radius, but groomed mass performance does not*



**Fig. 5** Comparisons of the leading jet  $p_T$  spectrum of the  $gg$  background to the  $WW$  signal in the  $p_T$  300-400 GeV bin using the different anti- $k_T$  jet distance parameters explored.



**Fig. 6** Comparisons of the leading jet  $p_T$  spectrum of the  $gg$  background to the  $WW$  signal in the  $p_T$  500-600 GeV bin using the different anti- $k_T$  jet distance parameters explored.

### 6.3 Combined Performance

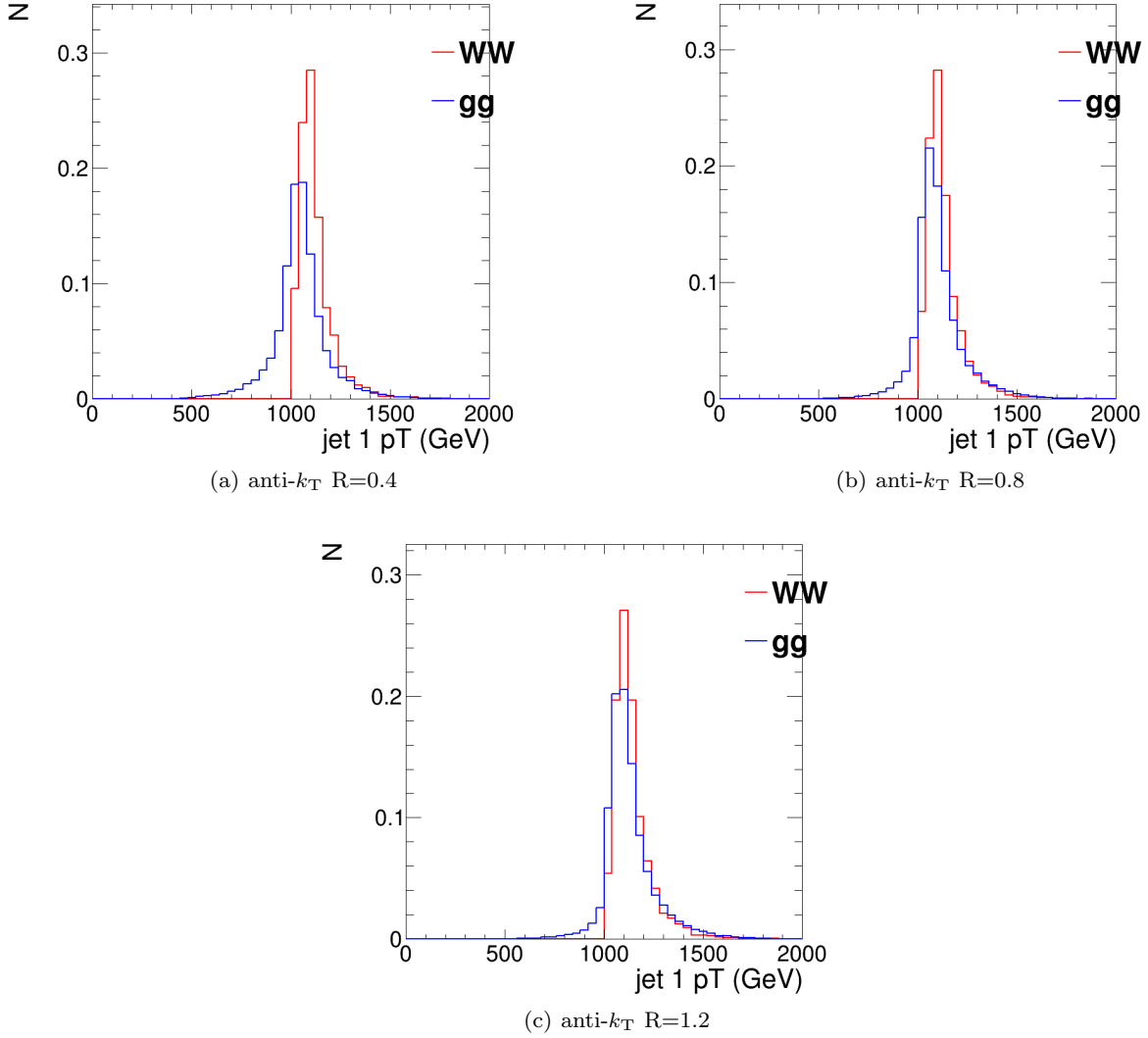
#### *Mass + X Performance*

Figure 16 shows the background efficiency for a fixed signal efficiency (50%) of each BDT combination of each pair of variables considered, in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm. One can see that the best background rejection is achieved using combinations of the groomed mass variables with other substructure variables (with the exception of the soft drop mass with  $\beta = -1$ ). Combinations of the mass vari-

ables themselves are not particularly powerful, but are interesting for understanding the correlations between the masses (see Section 6.3). Equally, combination of the substructure variables, without using a mass, are not powerful.

Figure 17 shows the actual ROC curves of the BDT combinations of each mass variable with every other variable considered in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm. *Can we drop the combinations of mass + mass from these plots to make them clearer? Also would be good to put the single variable mass curve on these plots, so you can see how much*





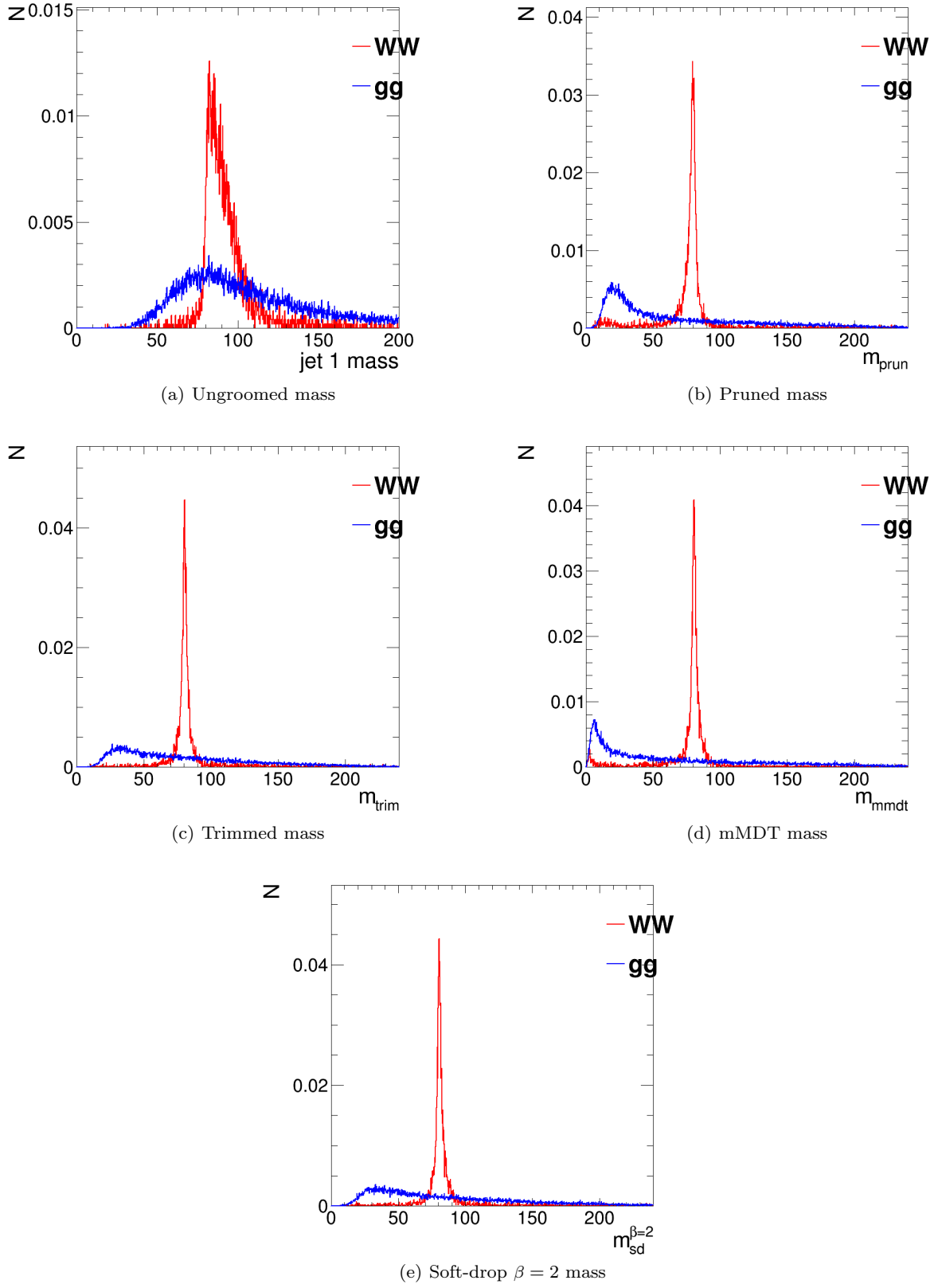
**Fig. 7** Comparisons of the leading jet  $p_T$  spectrum of the  $gg$  background to the  $WW$  signal in the  $p_T$  1.0-1.1 TeV bin using the different anti- $k_T$  jet distance parameters explored.

improvement the combination gives, and the “all variables” curve.

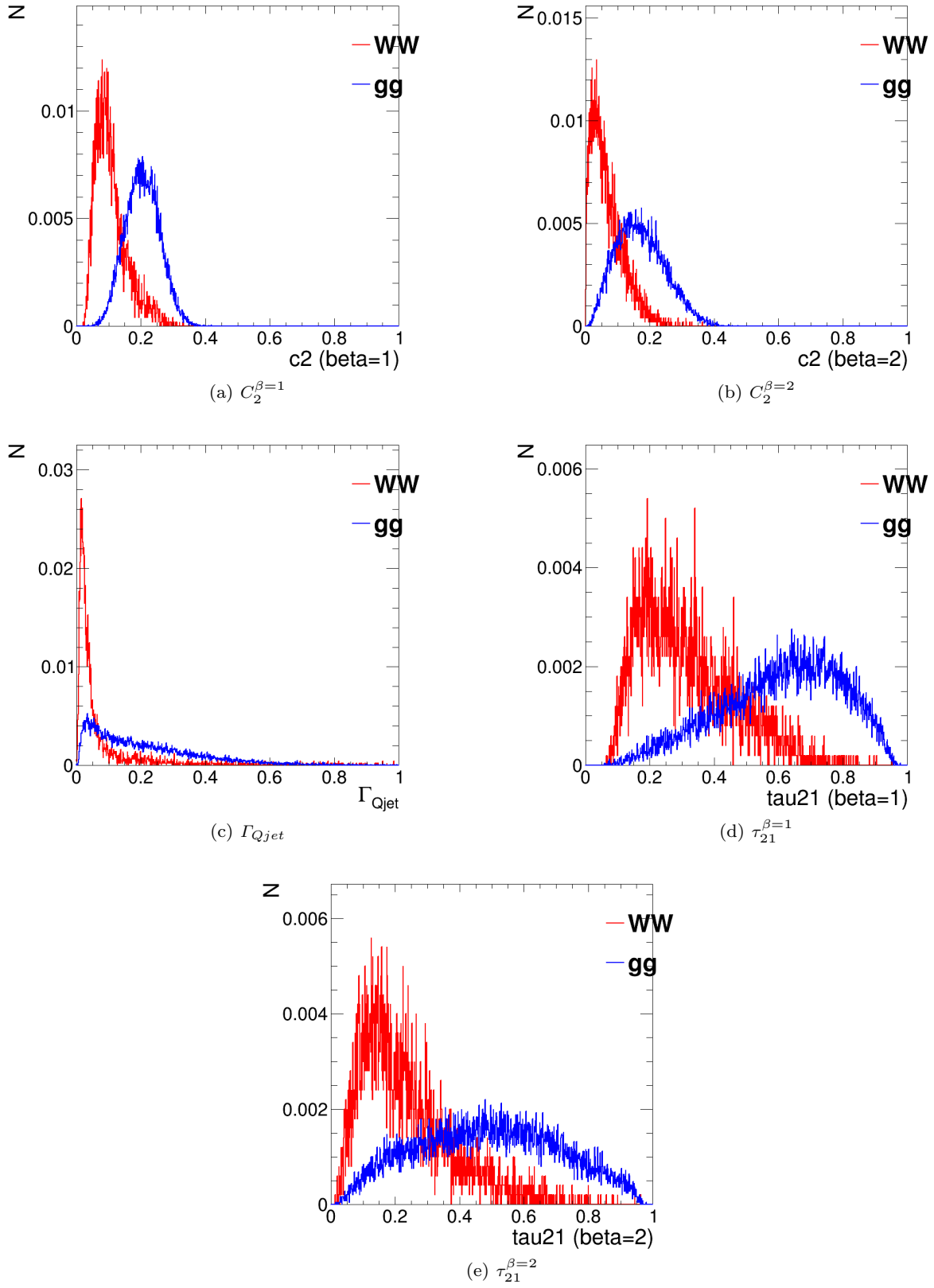
No combination with other variables can recover the poor performance of the ungroomed mass and the soft drop mass with  $\beta = -1$ . Figures 16 and 17 show that the other groomed/filtered masses are all most improved by combination with the  $C_2^{\beta=1}$  energy correlation function. Figure 18 shows the 2-D correlation plots between the mMDT mass and the  $C_2^{\beta=1}$ ,  $\Gamma_{Qjet}$  and  $\tau_{21}^{\beta=1}$  variables. One can clearly see that there is substantially less correlation between the mass and  $C_2^{\beta=1}$  than the other variables. Similar results are seen for the other groomed masses.

Figure 19 shows the background efficiency for a fixed signal efficiency (50%) of each BDT combination of each pair of variables considered, in the  $p_T$  500 GeV bin,

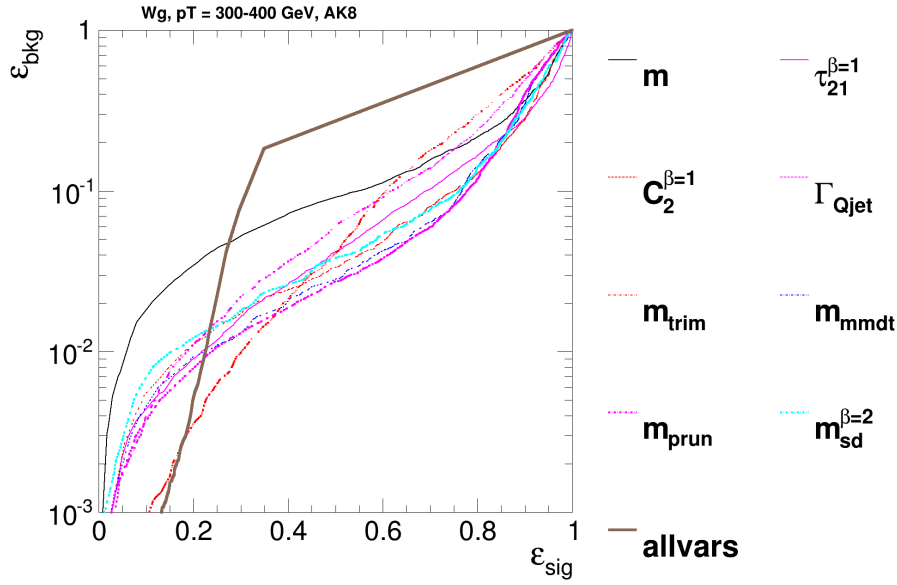
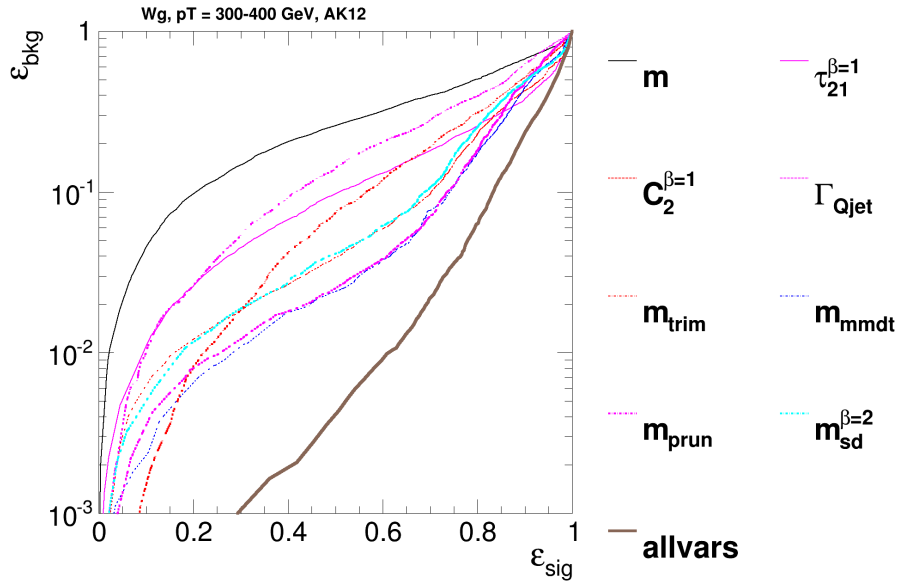
now using the anti- $k_T$   $R=1.2$  algorithm. Compared to Figure 16, the overall trends are similar, but there are clear differences in the relative power of the mass + X combinations. Interestingly, the groomed masses are now all most improved by combination with the  $\tau_{21}^{\beta=1}$  variable, in contrast with  $C_2^{\beta=1}$  which performed best for the smaller radius of  $R=0.8$ . Figure 20 shows the actual ROC curves for the BDT combinations of the best performing groomed masses with every other variable considered in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=1.2$  algorithm. One can see from Figure ?? that the single variable discrimination of  $\tau_{21}^{\beta=1}$  and  $C_2^{\beta=1}$  changes quite markedly when the distance parameter  $R$  is varied, although in both cases  $C_2^{\beta=1}$  is a better single variable discriminant (except for very high signal efficiencies). Figure 21 shows how the actual distribu-



**Fig. 8** Comparisons of the QCD background to the WW signal in the  $p_T$  500-600 GeV bin using the anti- $k_T$   $R=0.8$  algorithm: leading jet mass distributions.



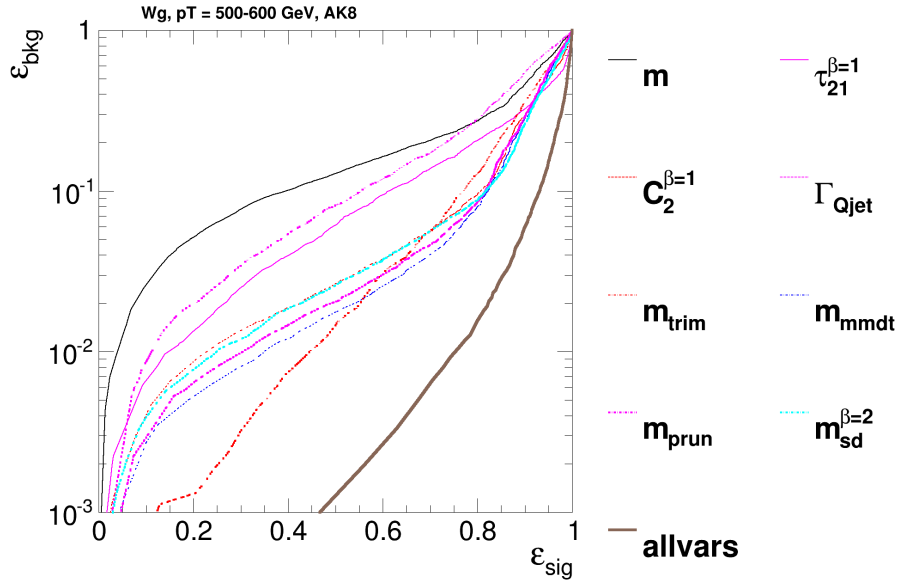
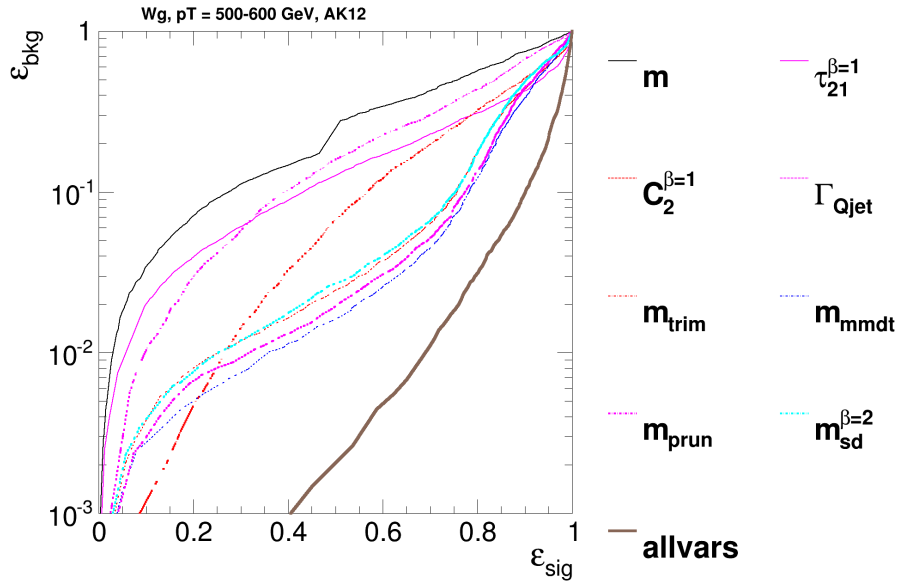
**Fig. 9** Comparisons of the QCD background to the WW signal in the  $p_T$  500-600 GeV bin using the anti- $k_T$  R=0.8 algorithm: substructure variables.

(a) anti- $k_T$   $R=0.8$ ,  $p_T$  300-400 GeV bin(b) anti- $k_T$   $R=1.2$ ,  $p_T$  300-400 GeV bin

**Fig. 10** The ROC curve for all single variables considered for  $W$  tagging in the  $p_T$  300-400 GeV bin using the anti- $k_T$   $R=0.8$  algorithm (top) and  $R=1.2$  algorithm (bottom).

tions of the  $C_2^{\beta=1}$  and  $\tau_{21}^{\beta=1}$  change when we change the distance parameter. Figure 22 shows the 2-D correlation plots between the mMDT mass and the  $C_2^{\beta=1}$ ,  $\Gamma_{Qjet}$  and  $\tau_{21}^{\beta=1}$  variables for the  $R=1.2$  case. It is hard to see a substantial difference in the correlations here versus Figure 18, but perhaps  $C_2^{\beta=1}$  is marginally more correlated with the mass for  $R=1.2$  compared to  $R=0.8$ . Andrew to add his explanation of why discrimination power of  $C_2$  versus  $\tau_{21}$  gets worse when we go to larger jet radii (email 0606/2014)

Now show a plot which compares on one plot the best combined performance for each groomed mass +  $X$  for both  $R=0.8$  and  $1.2$  cases e.g. mass +  $C_2^{\beta=1}$  for  $R=0.8$  and mass +  $\tau_{21}^{\beta=1}$  for  $R=1.2$ , and draw on also the all variables curve for both  $R=0.8, 1.2$ . Then we can see if there is much dependence on choice of mass once you combine with another variable, and compare directly the two distance parameters. This plot is just for one kinematic bin, we should make the same plot for others.


 (a) anti- $k_T$   $R=0.8$ ,  $p_T$  500-600 GeV bin

 (b) anti- $k_T$   $R=1.2$ ,  $p_T$  500-600 GeV bin

**Fig. 11** The ROC curve for all single variables considered for  $W$  tagging in the  $p_T$  500-600 GeV bin using the anti- $k_T$   $R=0.8$  algorithm (top) and  $R=1.2$  algorithm (bottom).

Repeat these studies for different  $R$  and different kinematic bins. Finally make plots which compare best combined performance for different  $R$  and kinematics.

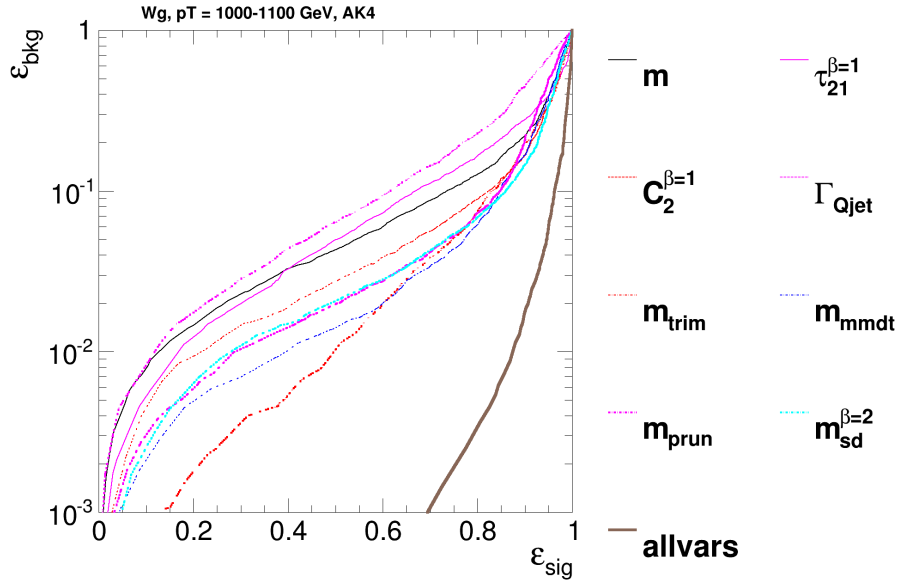
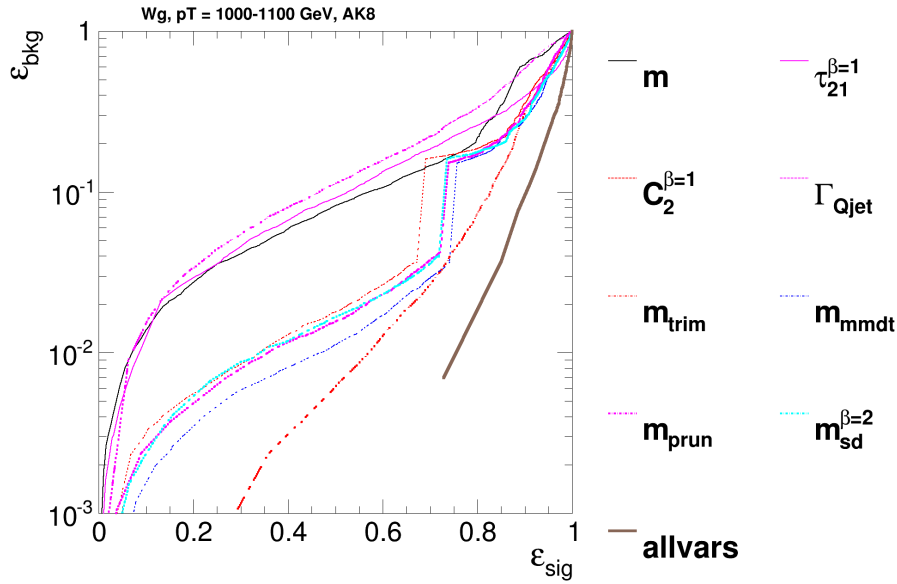
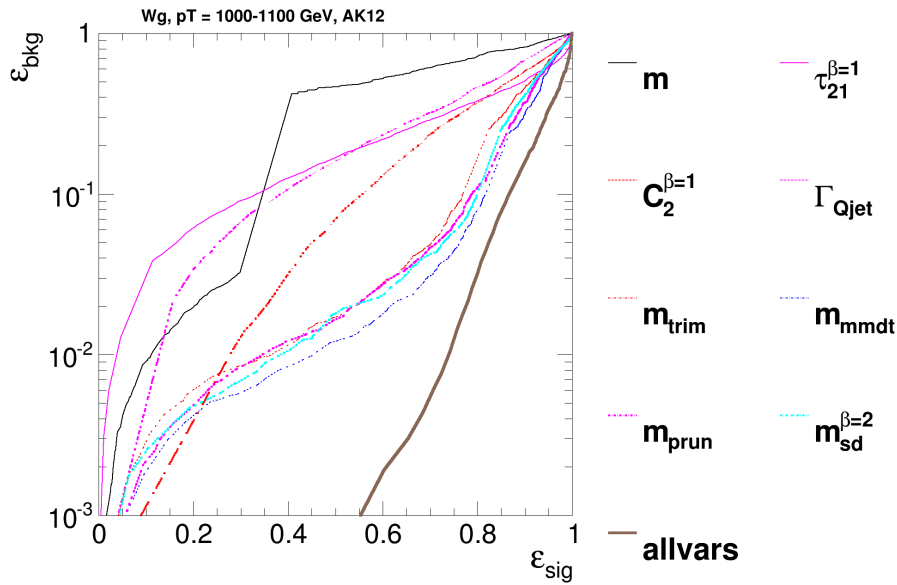
Do we want to look at other combinations of variables which don't involve mass? Practically I think we will always be making mass +  $X$  though.

#### Mass + Mass Performance

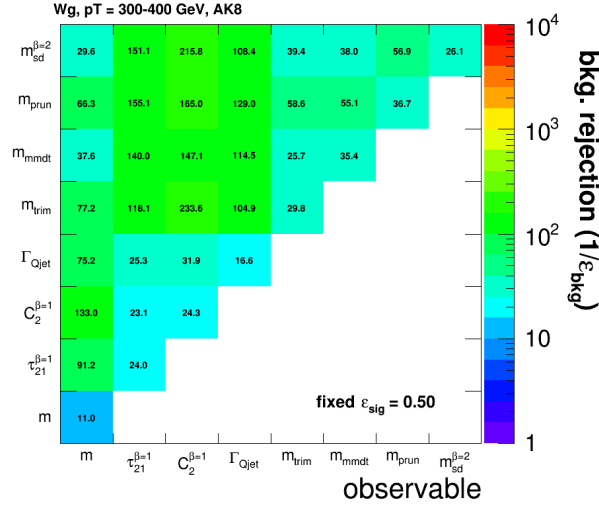
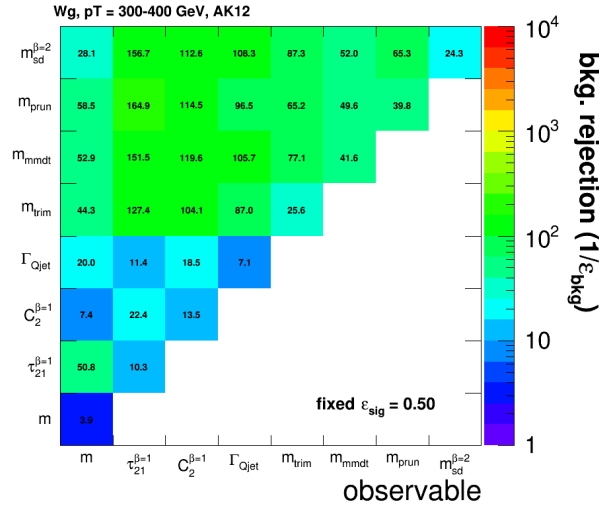
It's interesting also to study and understand how the different groomed masses relate to each other and how they are correlated.

Figures 23 and Figures 24 shows 2-D correlation plots of the different types of groomed mass in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm.

Worth also showing some ROC curves for mass + mass combinations?

(a) anti- $k_T$   $R=0.4$ ,  $p_T$  1.0-1.1 TeV bin(b) anti- $k_T$   $R=0.8$ ,  $p_T$  1.0-1.1 TeV bin(c) anti- $k_T$   $R=1.2$ ,  $p_T$  1.0-1.1 TeV bin

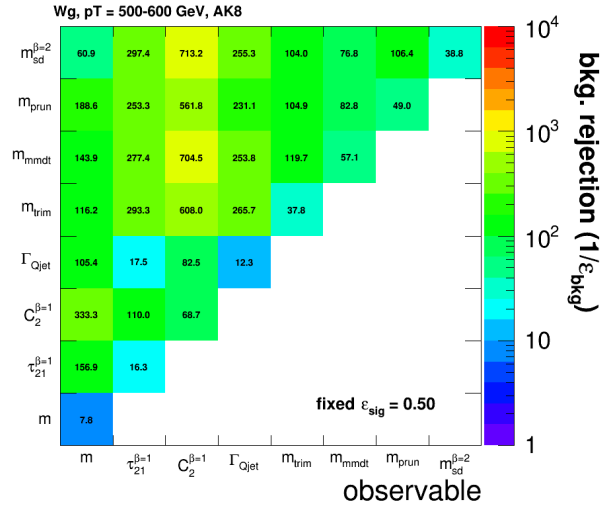
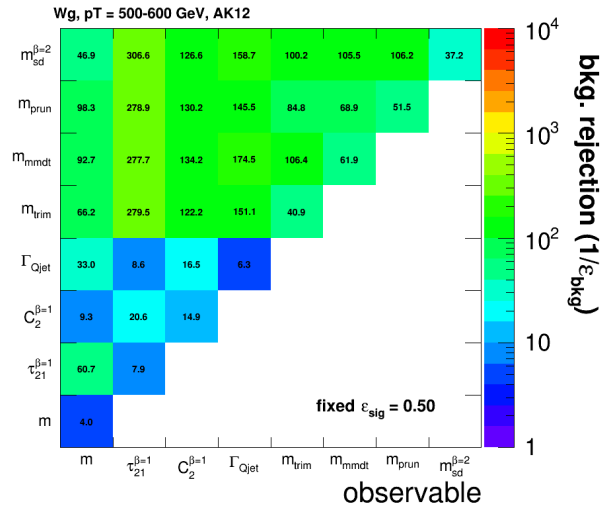
**Fig. 12** The ROC curve for all single variables considered for  $W$  tagging in the  $p_T$  1.0-1.1 TeV bin using the anti- $k_T$   $R=0.4$  algorithm (top), anti- $k_T$   $R=0.8$  algorithm (middle) and  $R=1.2$  algorithm (bottom).


 (a) anti- $k_T$   $R=0.8$ ,  $p_T$  300-400 GeV bin

 (b) anti- $k_T$   $R=1.2$ ,  $p_T$  300-400 GeV bin

**Fig. 13** The background rejection for a fixed signal efficiency (50%) of each BDT combination of each pair of variables considered, in the  $p_T$  300-400 GeV bin using the anti- $k_T$   $R=0.8$  algorithm (top) and  $R=1.2$  algorithm (bottom).

**Table 1** Action of various groomers on the jet mass distribution in the different phase space regions. For pruning,  $a_{\text{prune}} = z_{\text{cut}} R_0$  and for trimming  $a_{\text{trim}} = \sqrt{z_{\text{cut}}} R_{\text{sub}}$ .

Action	Pruning	Trimming	mMDT	SD ( $\beta > 0$ )
$m > \sqrt{z_{\text{cut}}} R_0 p_T$	—	—	—	—
$m < \sqrt{z_{\text{cut}}} R_0 p_T$ $m > a_x p_T$	cuts soft & soft-collinear	cuts soft & soft-collinear	cuts soft & soft-collinear	cuts soft & partially ( $\beta$ ) on soft-collinear
$m < a_x p_T$	cuts partially on both soft & soft-collinear	—	cuts soft & soft-collinear	cuts soft & partially ( $\beta$ ) on soft-collinear

(a) anti- $k_T$   $R=0.8$ ,  $p_T$  500-600 GeV bin(b) anti- $k_T$   $R=1.2$ ,  $p_T$  500-600 GeV bin

**Fig. 14** The background rejection for a fixed signal efficiency (50%) of each BDT combination of each pair of variables considered, in the  $p_T$  500-600 GeV bin using the anti- $k_T$   $R=0.8$  algorithm (top) and  $R=1.2$  algorithm (bottom).



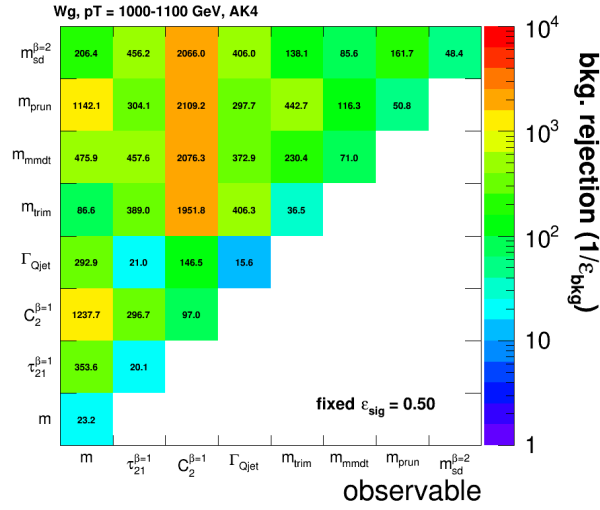
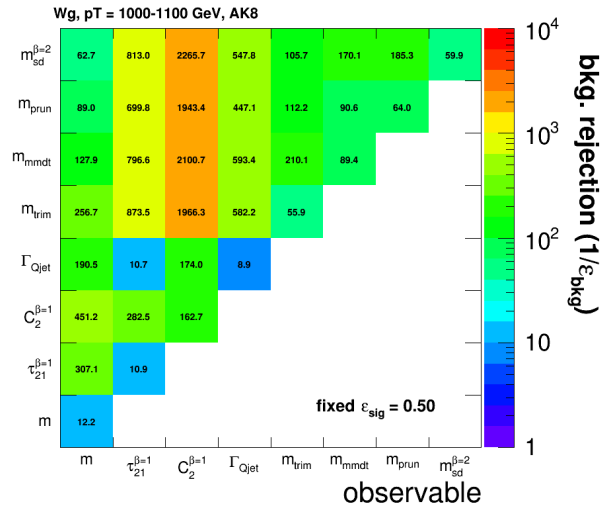
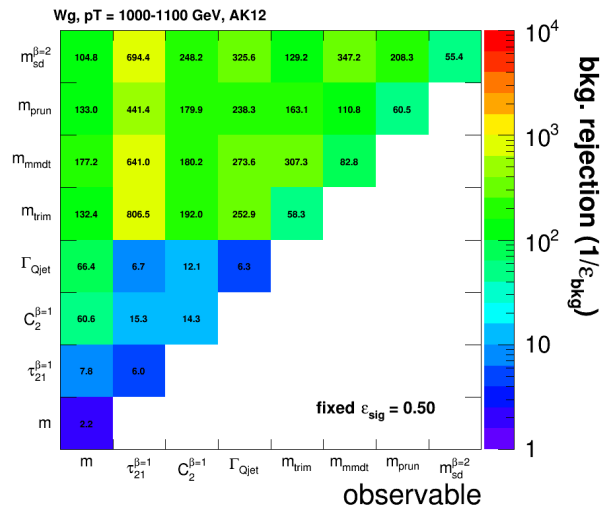
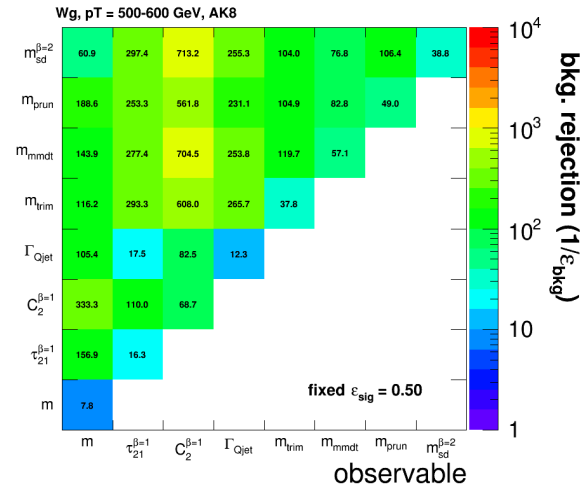
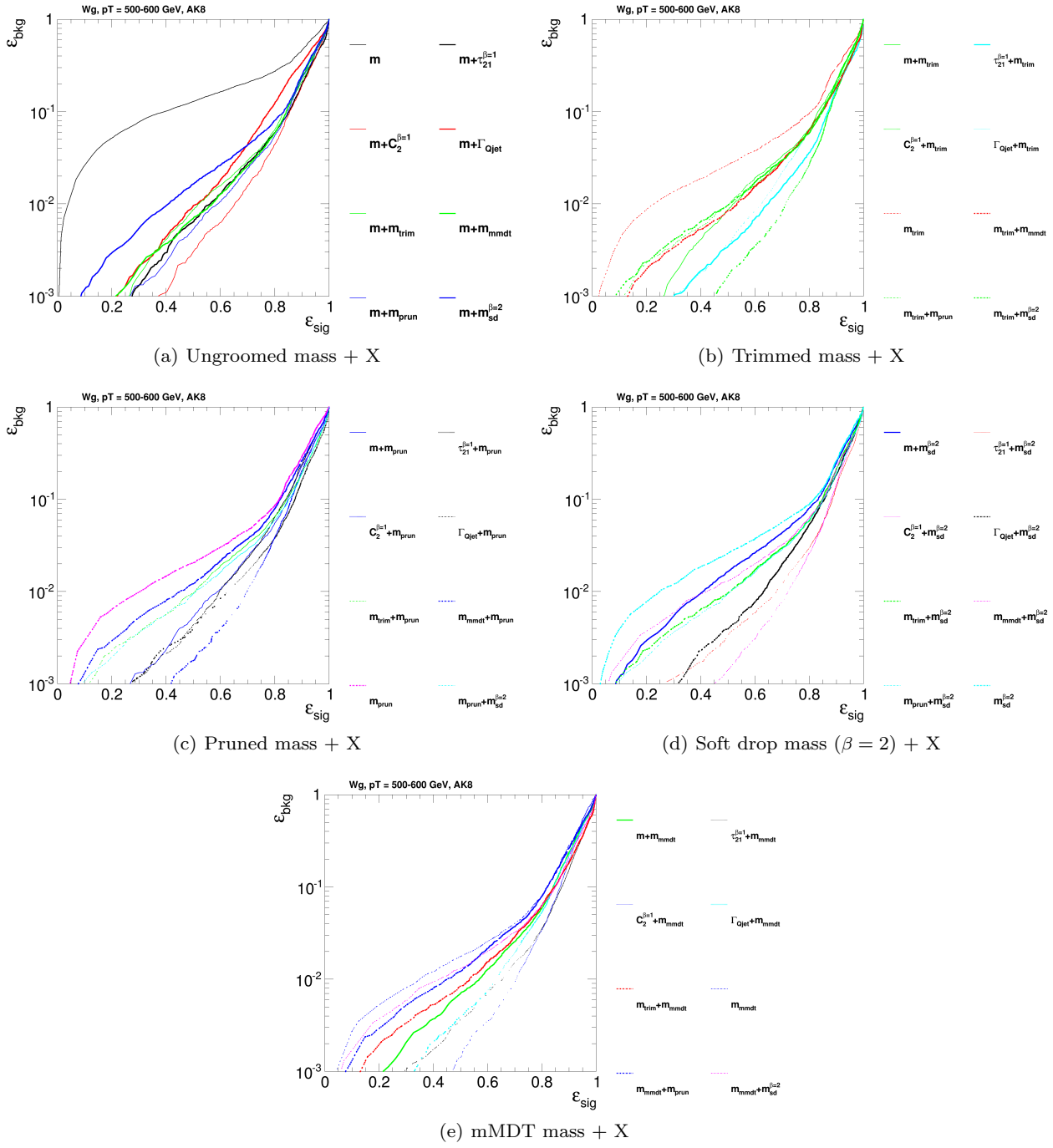

 (a) anti- $k_T$  R=0.4,  $p_T$  1.0-1.1 TeV bin

 (b) anti- $k_T$  R=0.8,  $p_T$  1.0-1.1 TeV bin

 (c) anti- $k_T$  R=1.2,  $p_T$  1.0-1.1 TeV bin

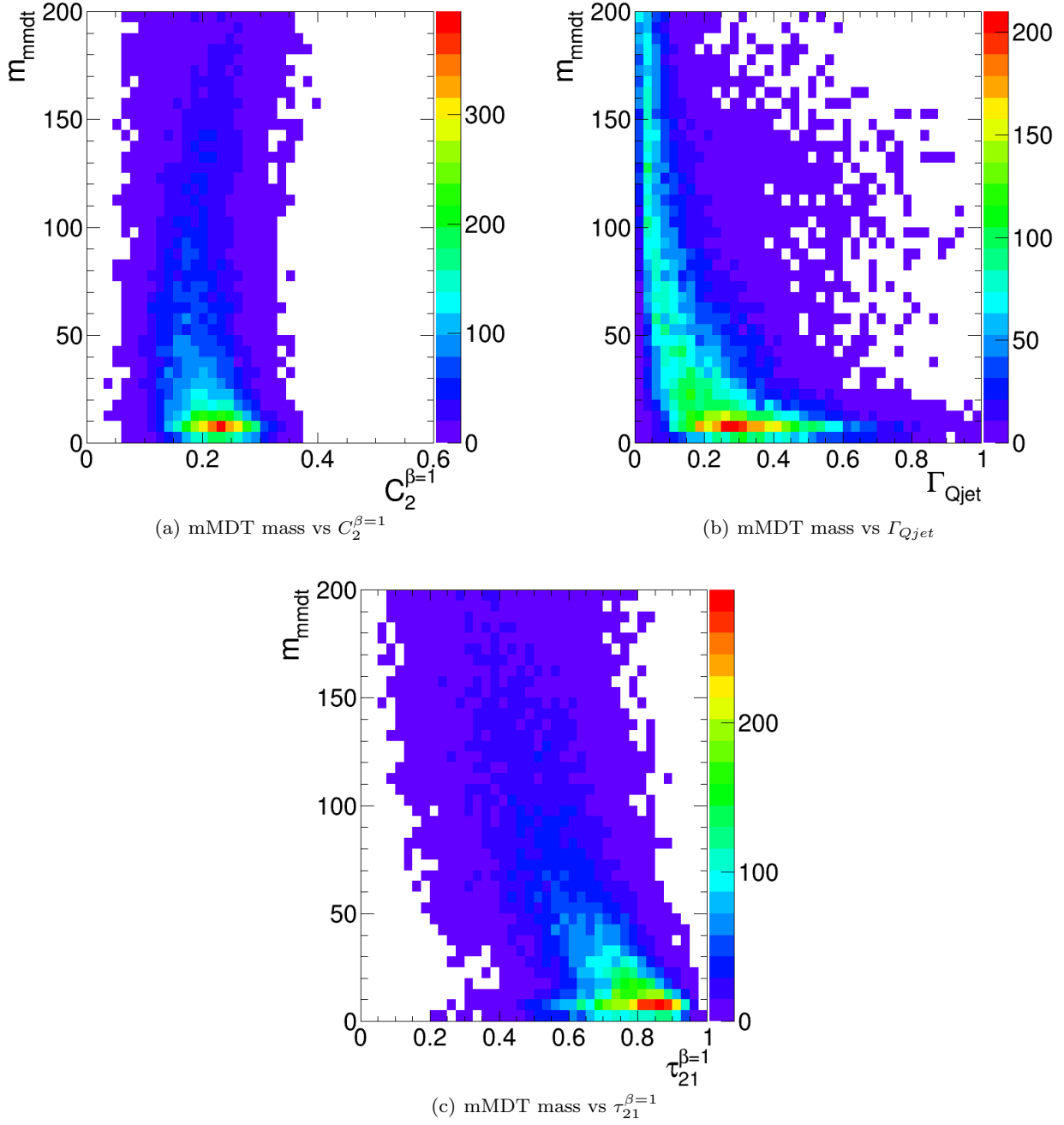
Fig. 15



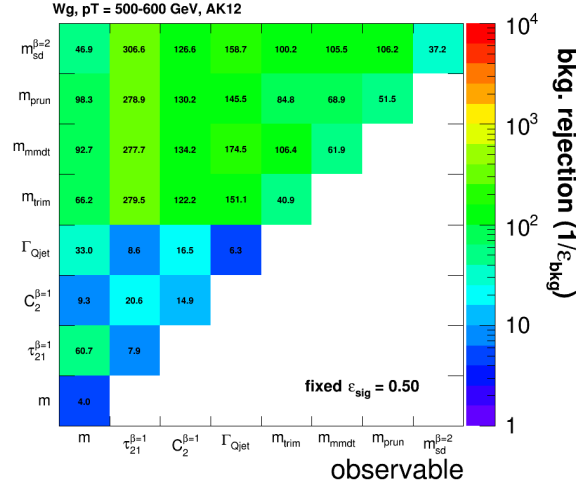
**Fig. 16** The background efficiency for a fixed signal efficiency (50%) of each BDT combination of each pair of variables considered, in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm.



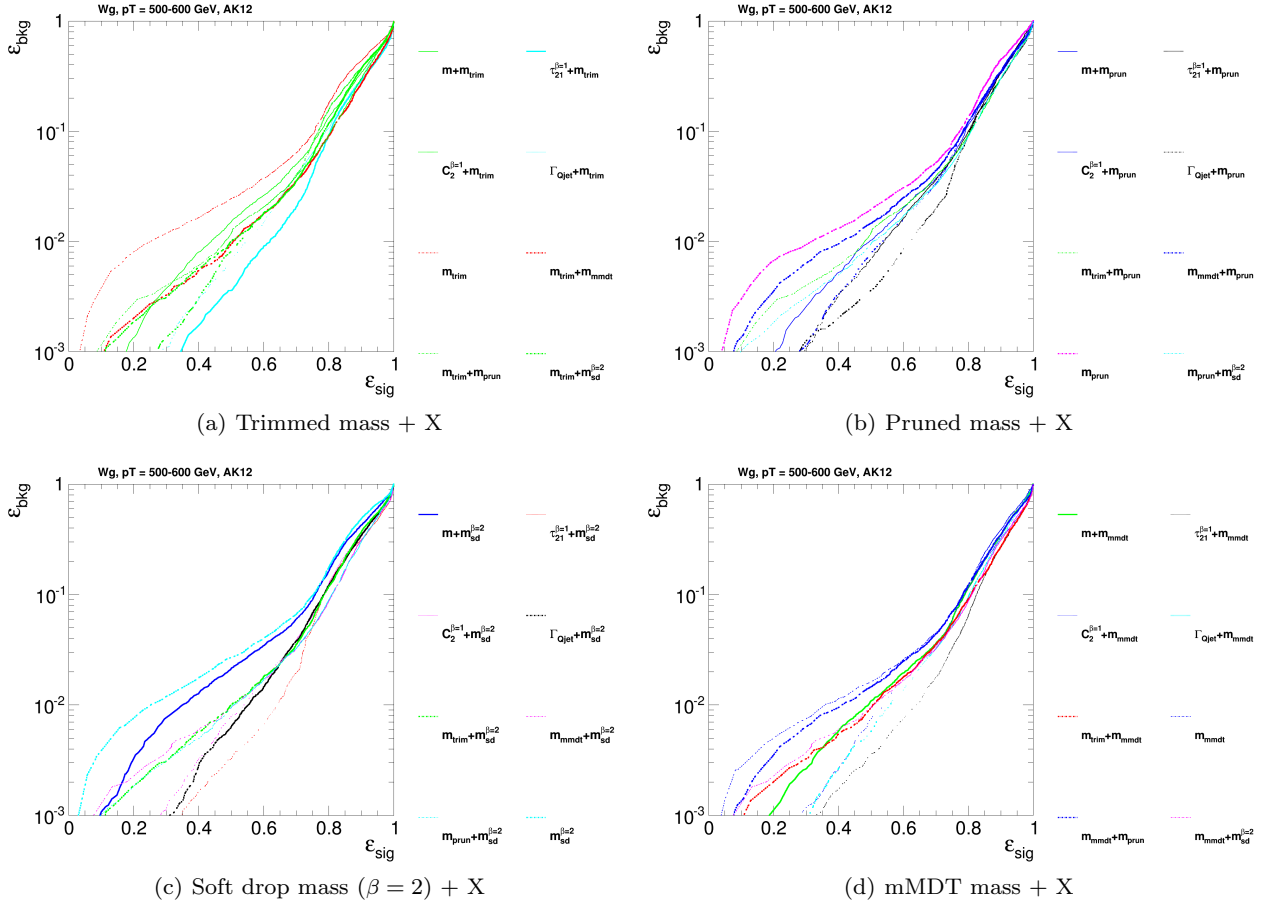
**Fig. 17** The BDT combinations of each mass variable with every other variable considered in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm.



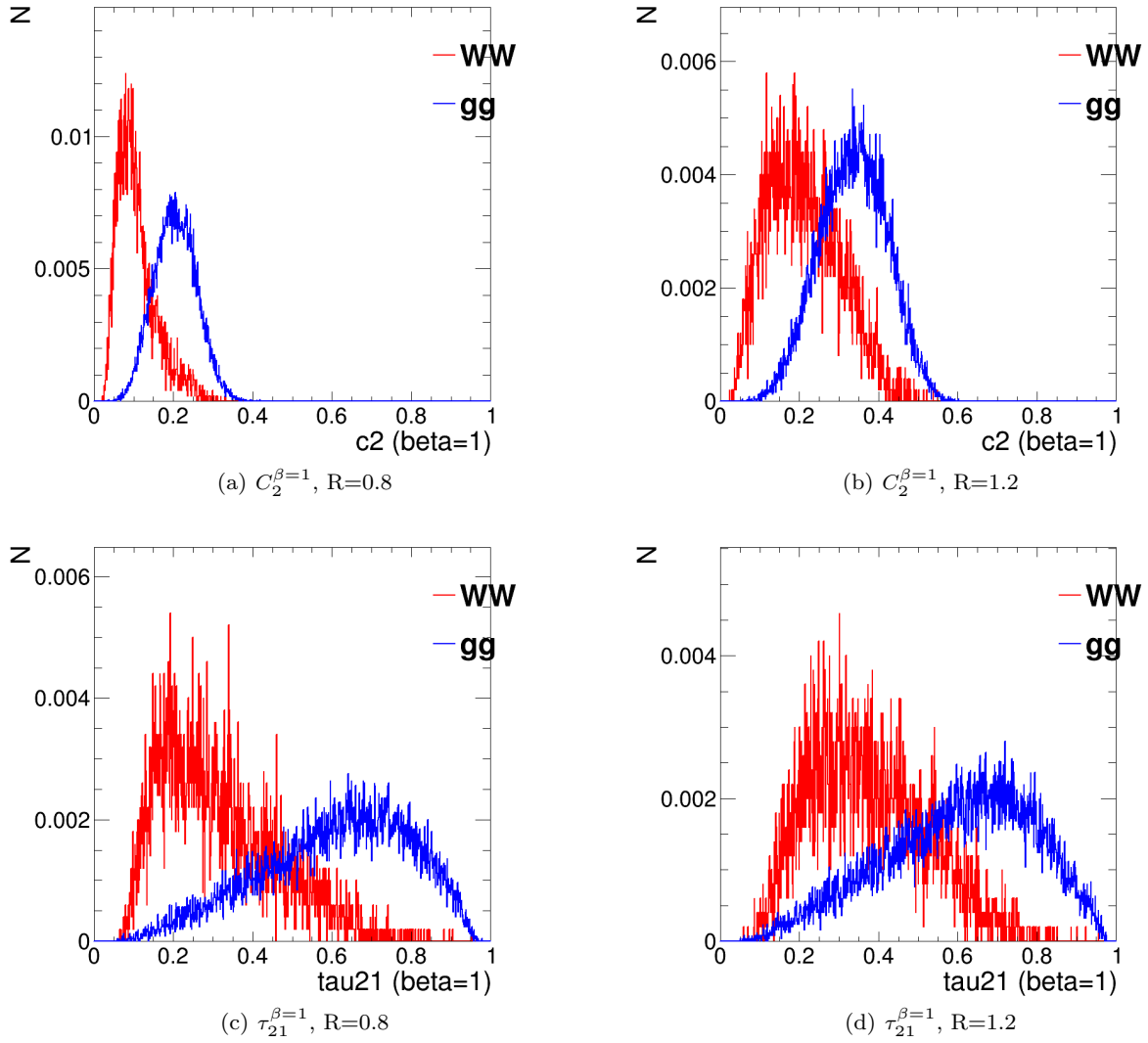
**Fig. 18** 2-D plots showing the correlation between mMDT mass and various substructure variables in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm in the gg sample.



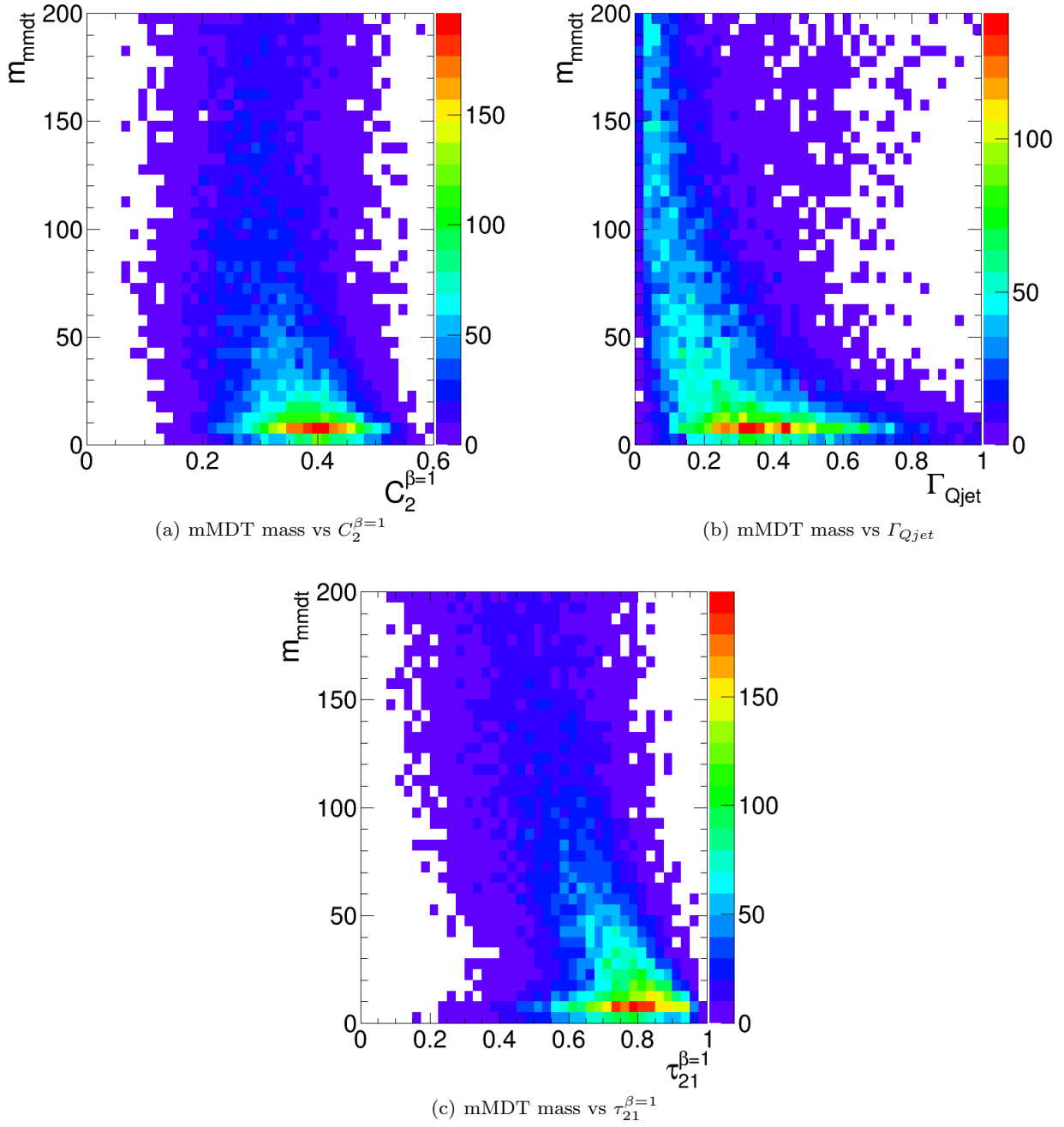
**Fig. 19** The background efficiency for a fixed signal efficiency (50%) of each BDT combination of each pair of variables considered, in the  $p_T$  500 GeV bin using the anti- $k_T$  R=1.2 algorithm.



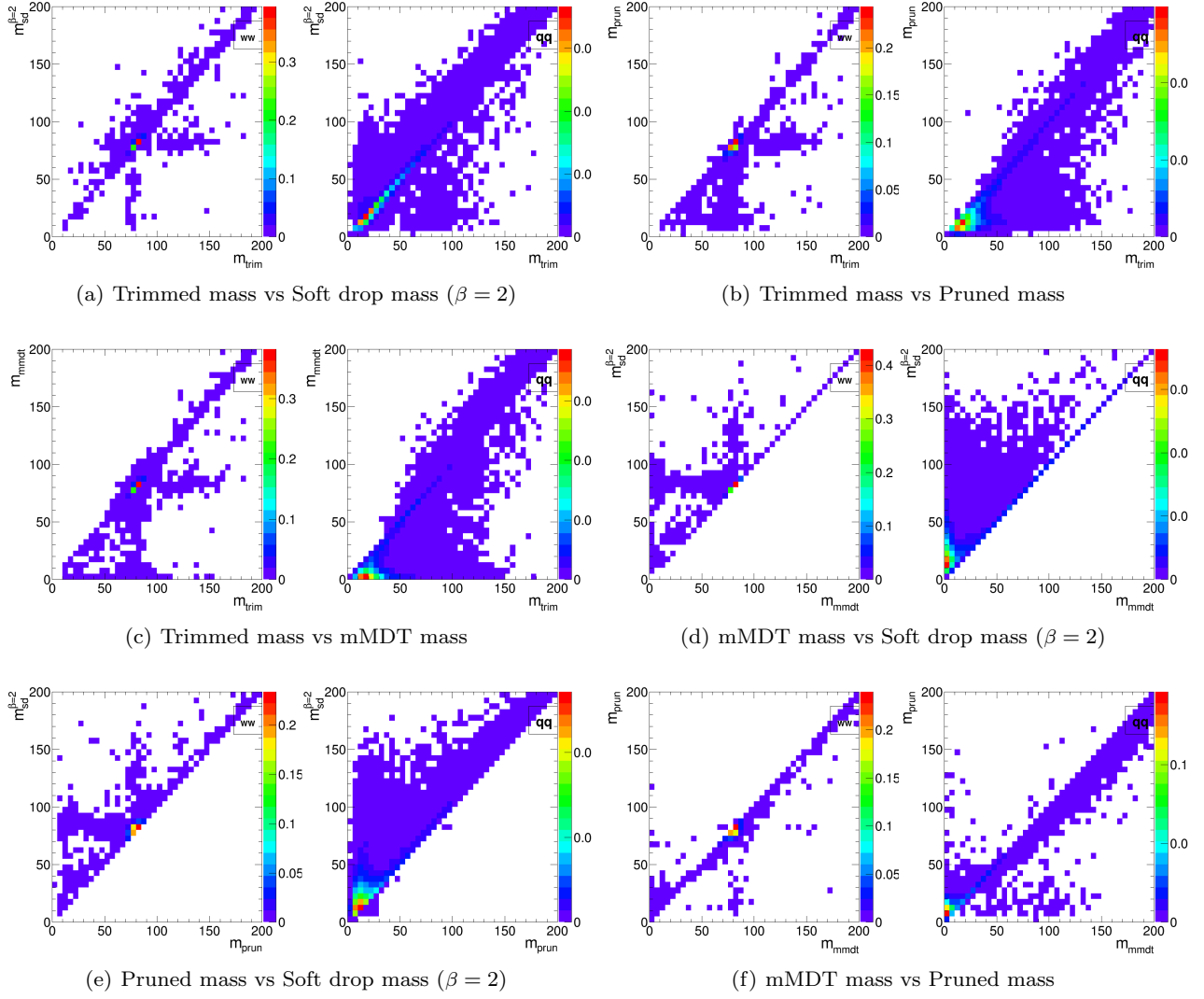
**Fig. 20** The BDT combinations of each mass variable with every other variable considered in the  $p_T$  500 GeV bin using the anti- $k_T$  R=1.2 algorithm.



**Fig. 21** Comparisons of the QCD background to the WW signal in the  $p_T$  500 GeV bin for  $C_2^{\beta=1}$  and  $\tau_{21}^{\beta=1}$  variables and using the R=0.8 and R=1.2 anti- $k_T$  distance parameters.

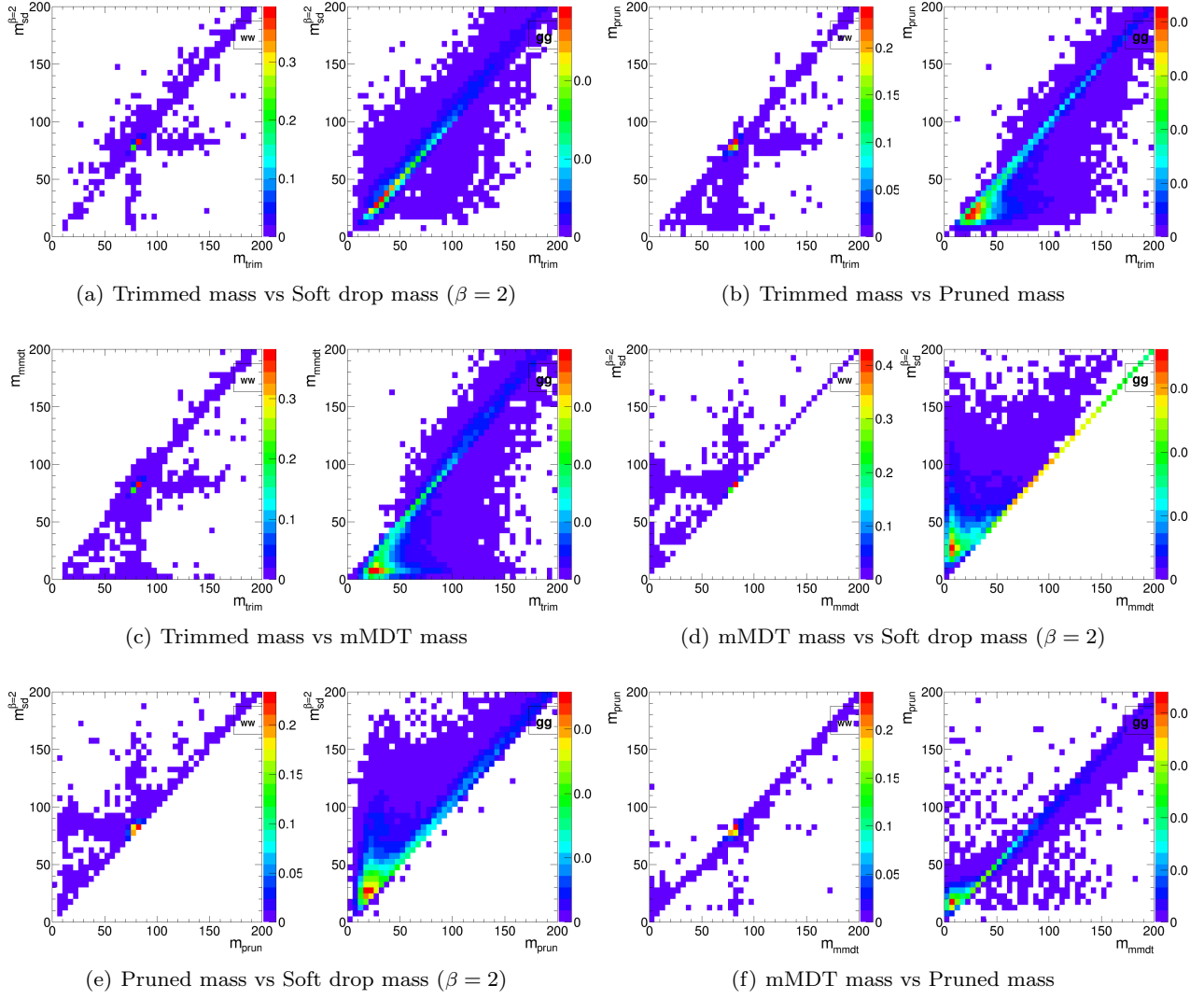


**Fig. 22** 2-D plots showing the correlation between mMDT mass and various substructure variables in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=1.2$  algorithm in the gg sample.



**Fig. 23** 2-D plots showing the correlation between different types of groomed mass in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm, separately for the jets in the  $X \rightarrow WW$  sample and the jets in the quark-quark sample.





**Fig. 24** 2-D plots showing the correlation between different types of groomed mass in the  $p_T$  500 GeV bin using the anti- $k_T$   $R=0.8$  algorithm, separately for the jets in the  $X \rightarrow WW$  sample and the jets in the gluon-gluon sample.

## 7 Top Tagging

In this section, we study the identification of boosted top quarks at Run II of the LHC. Boosted top quarks result in large-radius jets with complex substructure, containing a  $b$ -subjett and a boosted  $W$ . The additional kinematic handles coming from the reconstruction of the  $W$  mass and  $b$ -tagging allows a very high degree of discrimination of top quark jets from QCD backgrounds.

We consider top quarks with moderate boost (600-1000 GeV), and perhaps most interestingly, at high boost ( $\gtrsim 1500$  GeV). Top tagging faces several challenges in the high- $p_T$  regime. For such high- $p_T$  jets, the  $b$ -tagging efficiencies are no longer reliably known.

Also, the top jet can also be accompanied by additional radiation with  $p_T \sim m_t$ , leading to combinatoric ambiguities of reconstructing the top and  $W$ , and the possibility that existing taggers or observables shape the background by looking for subjet combinations that reconstruct  $m_t/m_W$ . To study this, we examine the performance of both mass-reconstruction variables, as well as shape observables that probe the three-pronged nature of the top jet and the accompanying radiation pattern.

### 7.1 Methodology

We study a number of top-tagging strategies, in particular:

1. HEPTopTagger
2. Johns Hopkins Tagger (JH)
3. Trimming
4. Pruning

The top taggers have criteria for reconstructing a top and  $W$  candidate, while the grooming algorithms (trimming and pruning) do not incorporate a  $W$ -identification step. For a level playing field, we construct a  $W$  candidate from the three leading subjets by taking the pair of subjets with the smallest invariant mass; in the case that only two subjets are reconstructed, we take the mass of the leading subjet. All of the above taggers and groomers incorporate a step to remove pile-up and other soft radiation.

We also consider the performance of jet shape observables. In particular, we consider the  $N$ -subjettiness ratios  $\tau_{32}^{\beta=1}$  and  $\tau_{21}^{\beta=1}$ , energy correlation function ratios  $C_3^{\beta=1}$  and  $C_2^{\beta=1}$ , and the Qjet mass volatility  $\Gamma$ . In addition to the jet shape performance, we combine the jet shapes with the mass-reconstruction methods listed above to determine the optimal combined performance.

To quantify the performance of each set of variables, we combine the relevant tagger output observables and/or jet shapes into a boosted decision tree (BDT), which determines the optimal multivariable cut. Additionally, because each tagger has two inputs (list, or maybe refer back to Section 3), we scan over reasonable values of the inputs to determine the optimal value for each top tagging signal efficiency. This allows a direct comparison of the optimized version of each tagger.

### 7.2 Single-observable performance

We start by investigating the behavior of individual jet substructure observables. Because of the rich, three-pronged structure of the top decay, it is expected that

combinations of masses and jet shapes will far outperform single observables in identifying boosted tops. However, a study of the top-tagging performance of single variables facilitates a direct comparison with the  $W$  tagging results in Section 6, and also allows a straightforward examination of the performance of each observable for different  $p_T$  and jet radius.

Fig. 25 shows the ROC curves for each of the top-tagging observables, with the bare jet mass also plotted for comparison. Unlike  $W$  tagging, the jet shape observables perform more poorly than jet mass. (*Check reasoning: this argument due to Andrew Larkoski*). As an example illustrating why this is the case, consider  $N$ -subjettiness. The  $W$  is two-pronged and the top is three-pronged; therefore, we expect  $\tau_{21}$  and  $\tau_{32}$  to be the best-performant  $N$ -subjettiness ratio, respectively. However,  $\tau_{21}$  also contains an implicit cut on the denominator,  $\tau_1$ , which is strongly correlated with jet mass. Therefore,  $\tau_{21}$  combines both mass and shape information to some extent. By contrast, and as is clear in Fig. 25(a), the best shape for top tagging is  $\tau_{32}$ , which contains no information on the mass. Therefore, it is unsurprising that the shapes most useful for top tagging are less sensitive to the jet mass, and under-perform relative to the corresponding observables for  $W$  tagging.

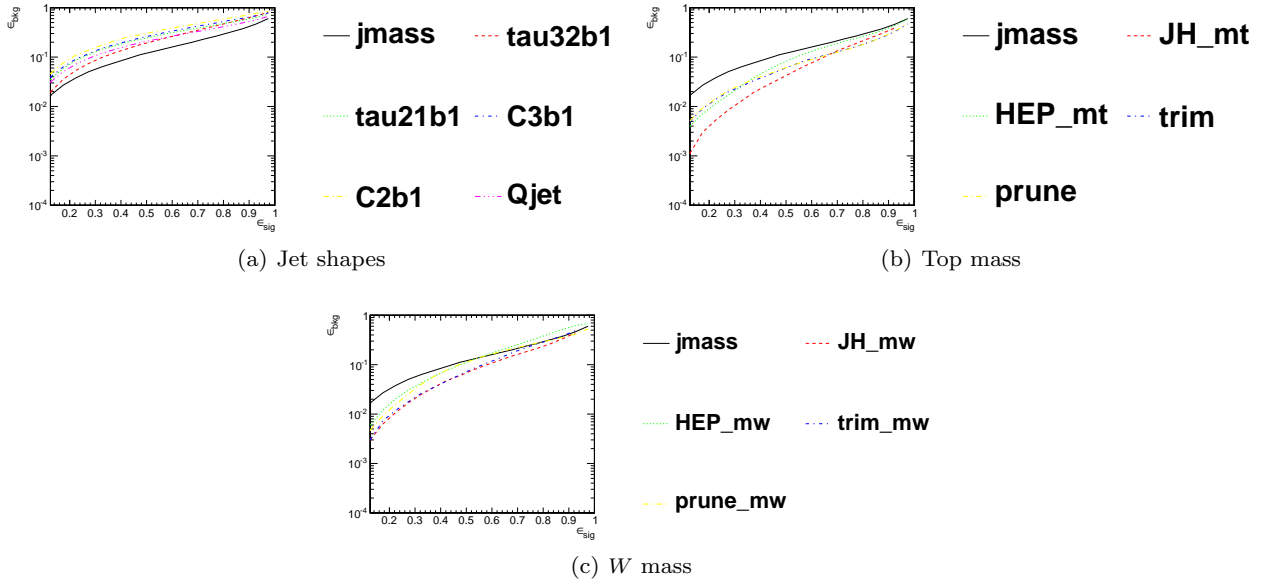
Of the two top tagging algorithms, the Johns Hopkins (JH) tagger out-performs the HEPTopTagger in its signal-to-background separation of both the top and  $W$  candidate masses, with larger discrepancy at higher  $p_T$  and larger jet radius. In Fig. 26, we show the histograms for the top mass output from the JH and HEPTopTagger for different  $p_T$  and  $R$ , optimized at a signal efficiency of 30%. The likely reason for this behavior is that, in the HEPTopTagger algorithm, the jet is filtered to select the five hardest subjets, and then three subjets are chosen which reconstruct the top mass. This requirement tends to shape a peak in the QCD background around  $m_t$  for the HEPTopTagger, while the JH tagger has no such requirement. It has been suggested by Anders *et al.* [4] that performance in the HEPTopTagger may be improved by selecting the three subjets reconstructing the top only among those that pass the  $W$  mass constraints, which somewhat reduces the shaping of the background. *Maybe try this out with my code to see if it helps?*

We also directly compare each variable's performance for different jet  $p_T$  and radius. The results are shown in Figs. 27-29 for different  $p_T$  bins and Figs. 30-32 for different  $R$  values. The input parameters of the taggers, groomers, and shape variables are separately optimized for each  $p_T$  and radius. If we only optimize the tagger inputs for one value of  $p_T$  and  $R$ , the ROC curve behavior does not change substantially from one where the

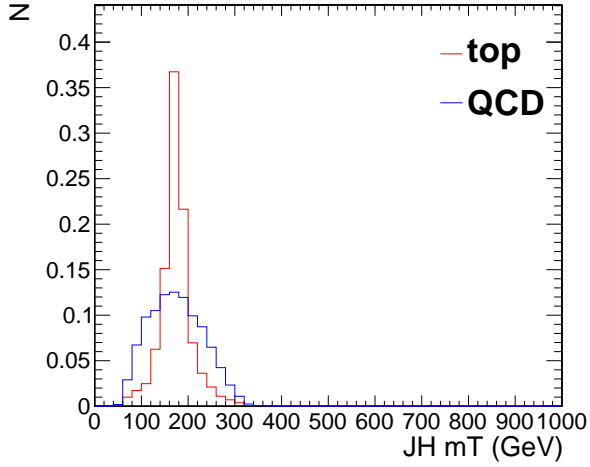
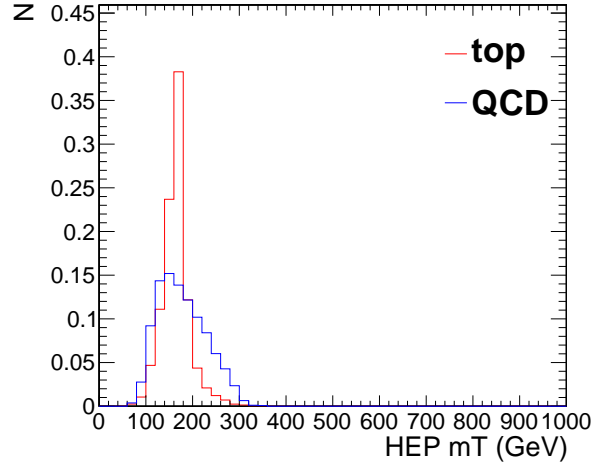
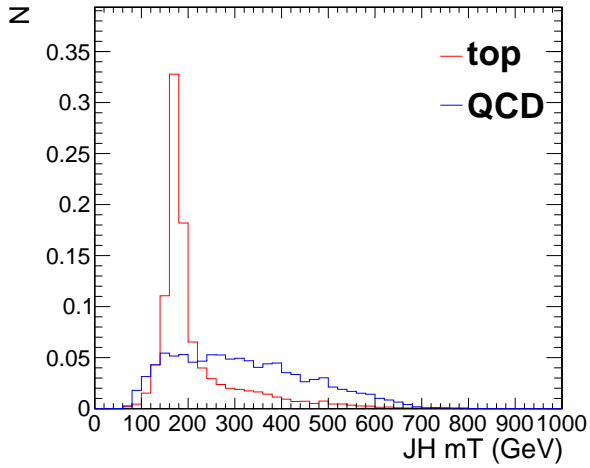
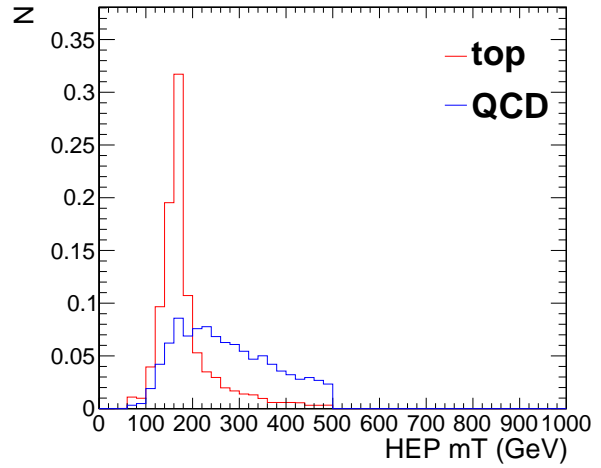
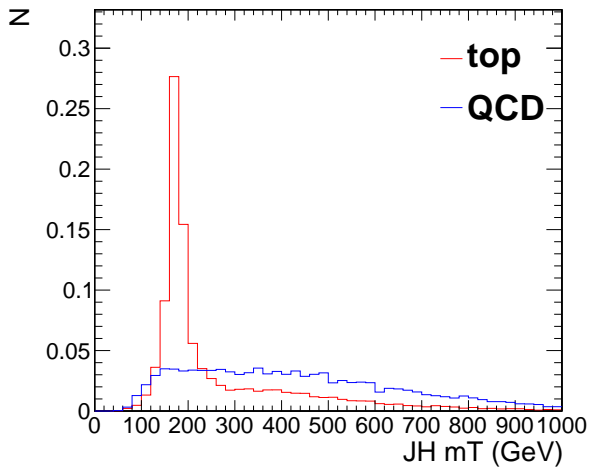
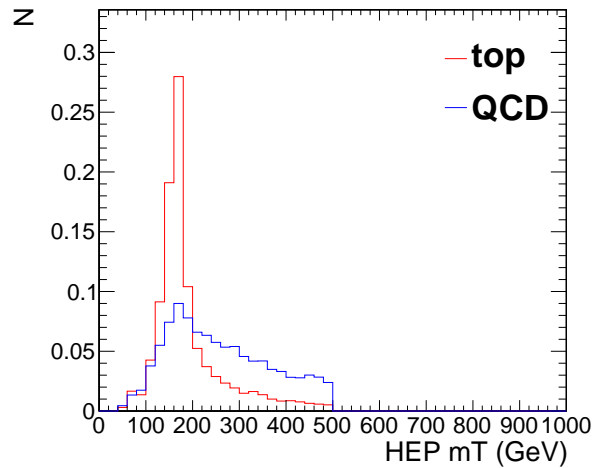
inputs are optimized at each  $p_T$  and  $R$  value; however, not all signal efficiencies are possible for every choice of tagger input, since the baseline selection efficiency might be too low.

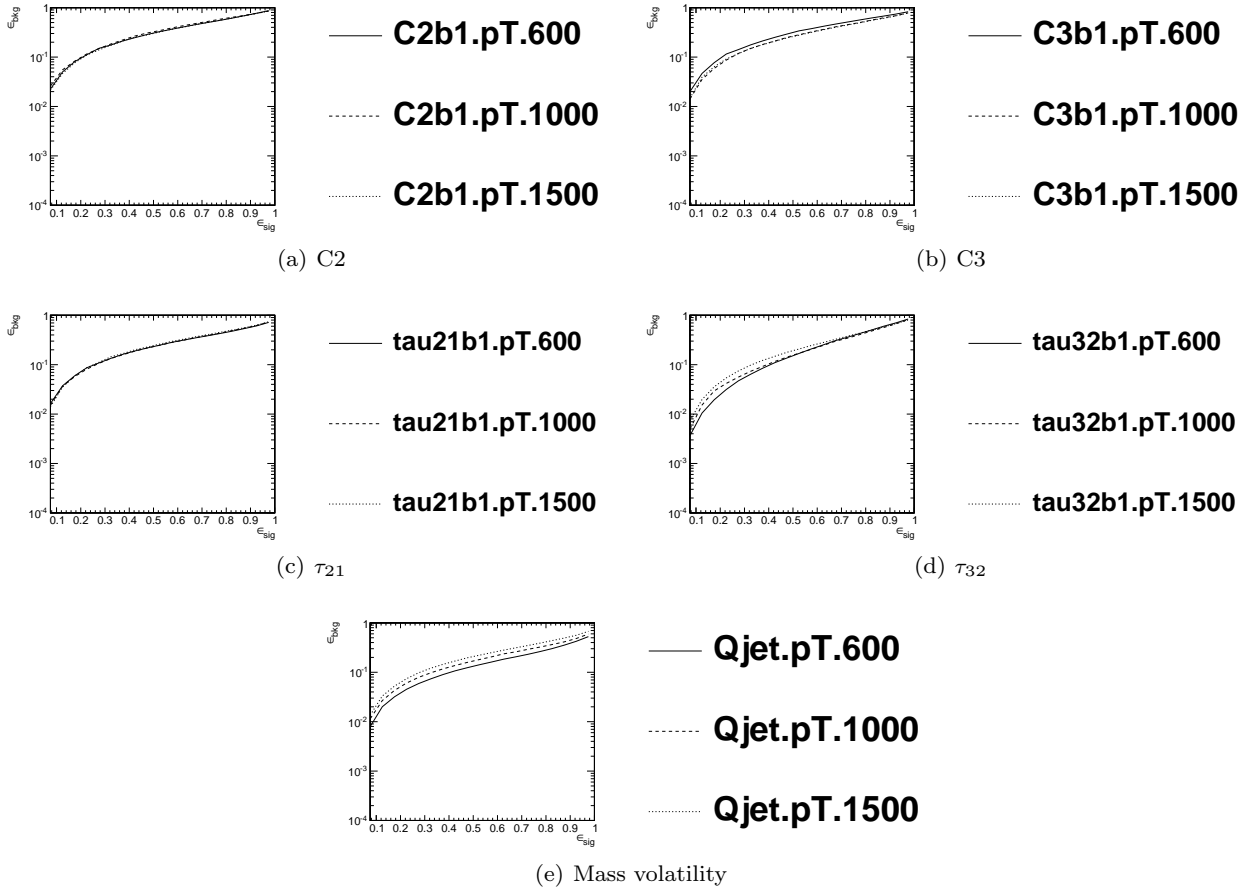
### 7.3 Performance of multivariable combinations

### 7.4 Performance at Sub-Optimal Working Points

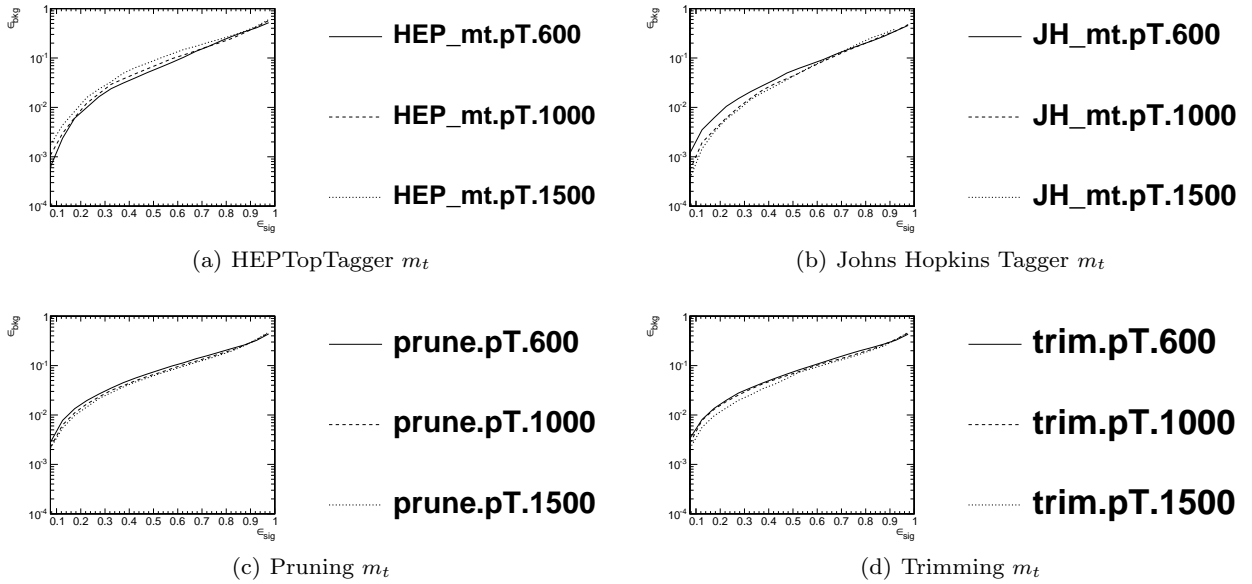


**Fig. 25** Comparison of single-variable top-tagging performance in the  $p_T$  1000-1100 GeV bin using the anti- $k_T$ ,  $R=0.8$  algorithm.

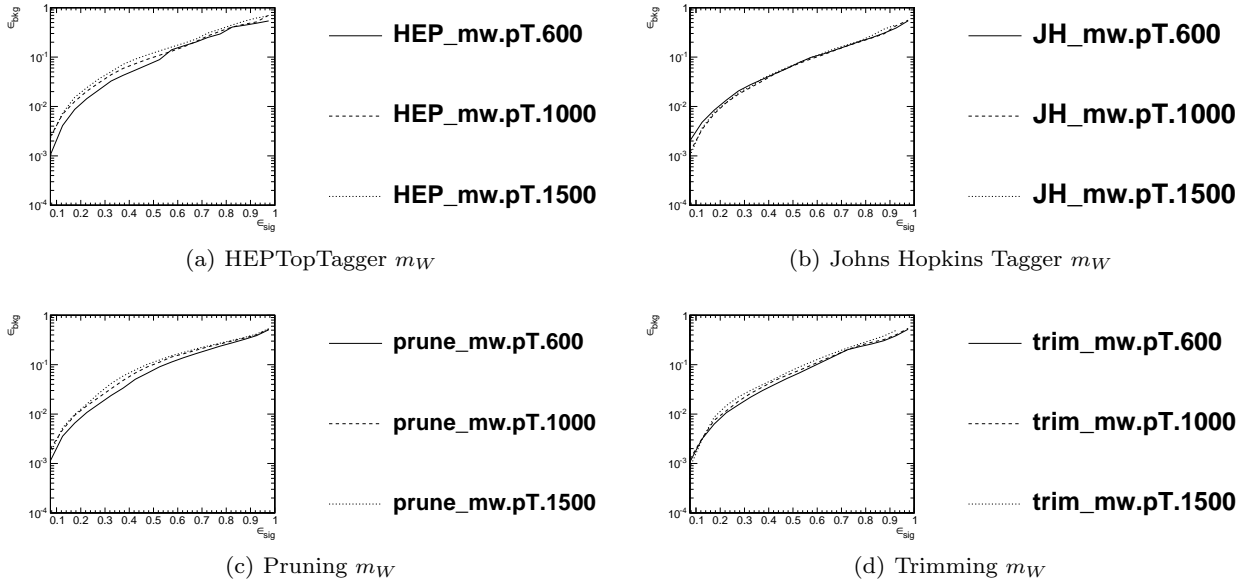
(a) Johns Hopkins Tagger,  $p_T = 600 - 700$  GeV,  $R = 0.8$ (b) HEPTopTagger,  $p_T = 600 - 700$  GeV,  $R = 0.8$ (c) Johns Hopkins Tagger,  $p_T = 1500 - 1600$  GeV,  $R = 0.8$ (d) HEPTopTagger,  $p_T = 1500 - 1600$  GeV,  $R = 0.8$ (e) Johns Hopkins Tagger,  $p_T = 1500 - 1600$  GeV,  $R = 1.2$ (f) HEPTopTagger,  $p_T = 1500 - 1600$  GeV,  $R = 1.2$ **Fig. 26** Comparison of individual jet shape performance at different  $p_T$  using the anti- $k_T$   $R=0.8$  algorithm.



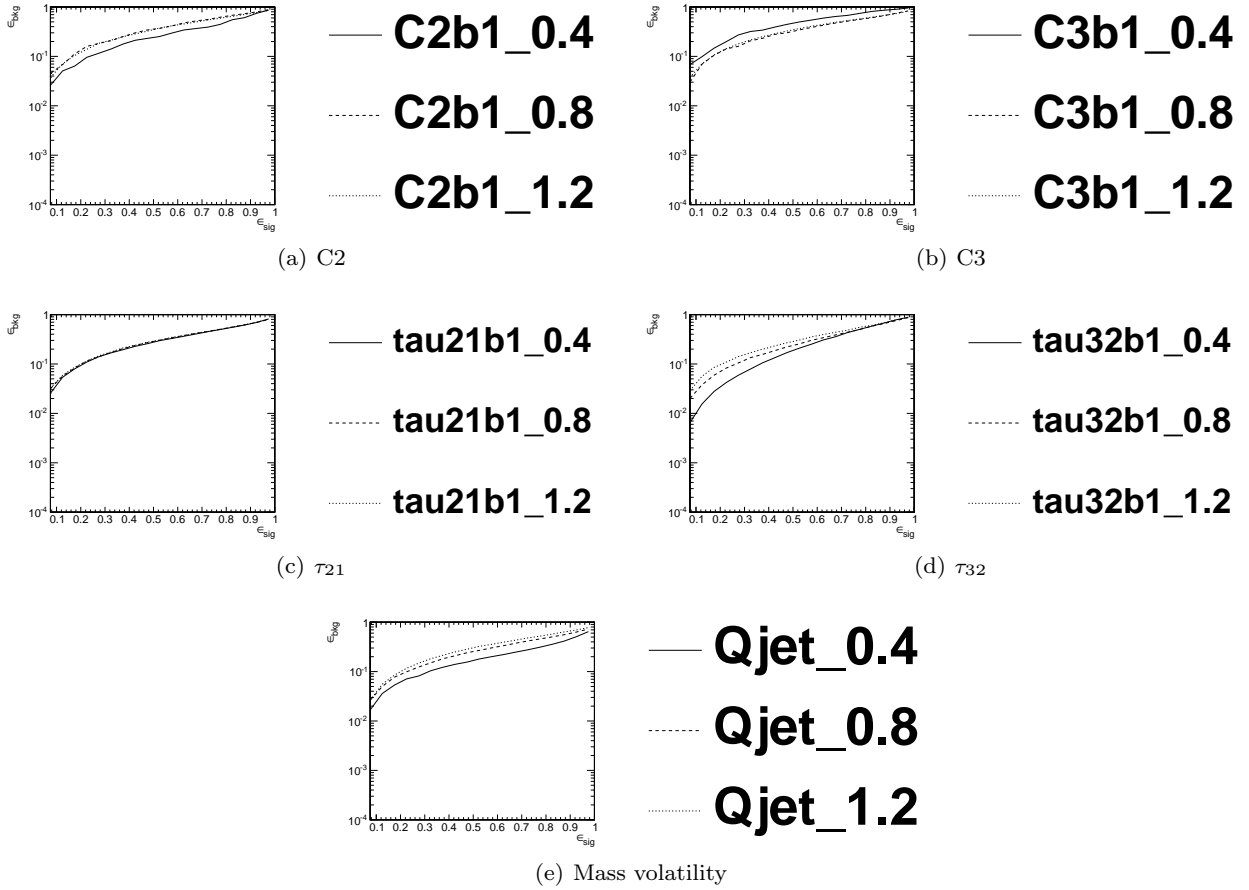
**Fig. 27** Comparison of individual jet shape performance at different  $p_T$  using the anti- $k_T$   $R=0.8$  algorithm.



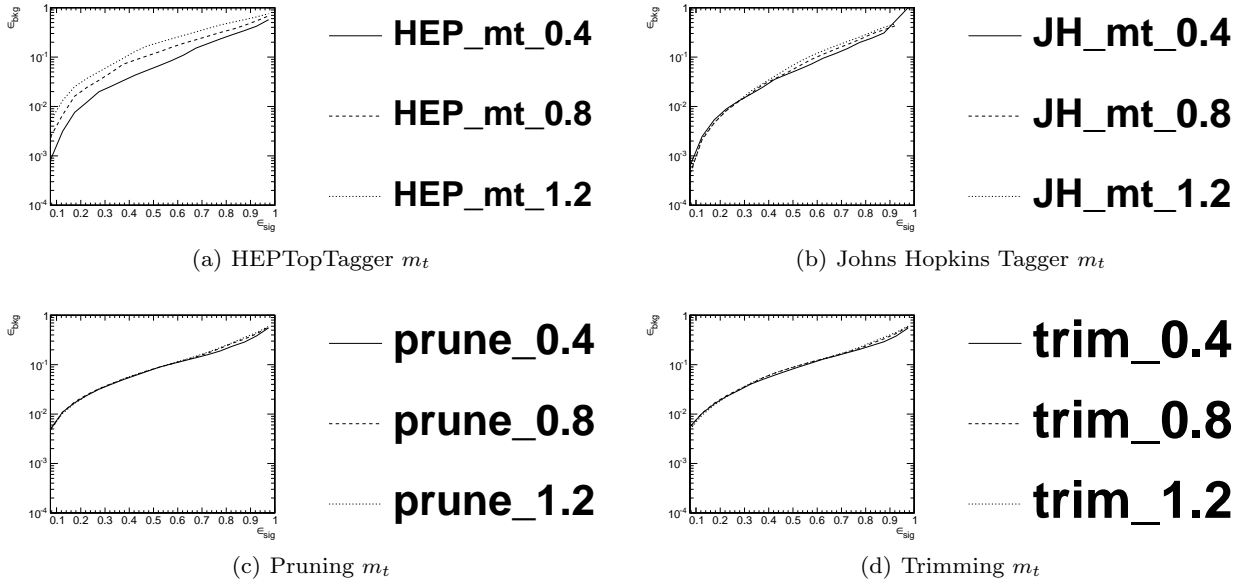
**Fig. 28** Comparison of top mass performance of different taggers at different  $p_T$  using the anti- $k_T$   $R=0.8$  algorithm.



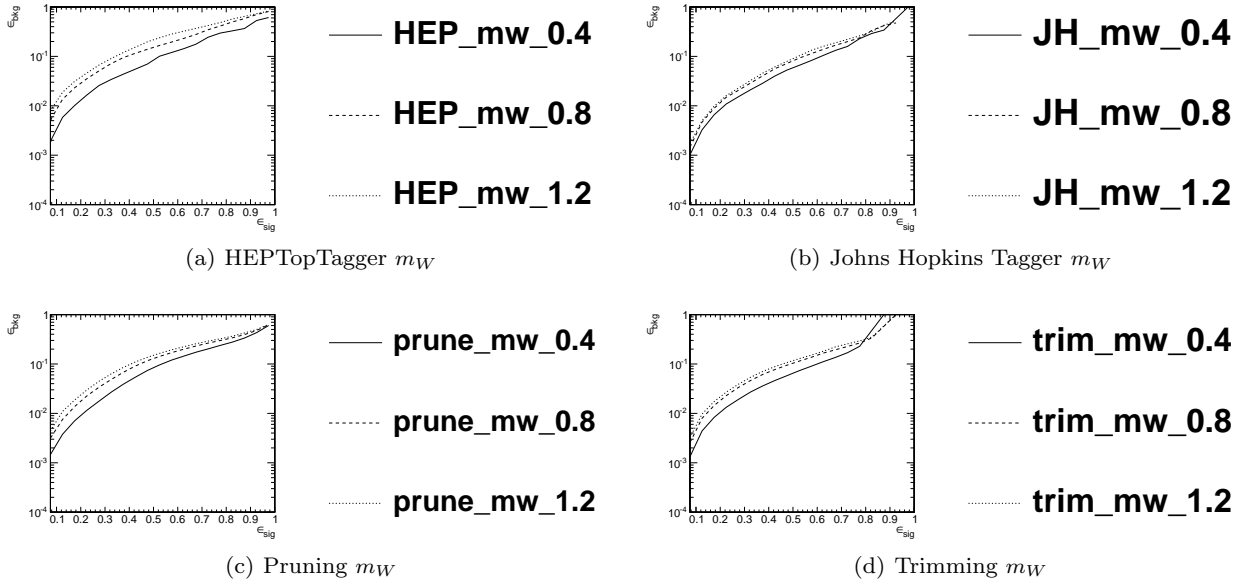
**Fig. 29** Comparison of  $W$  mass performance of different taggers at different  $p_T$  using the anti- $k_T$   $R=0.8$  algorithm.



**Fig. 30** Comparison of individual jet shape performance at different  $R$  in the  $p_T = 1500 - 1600$  GeV bin.

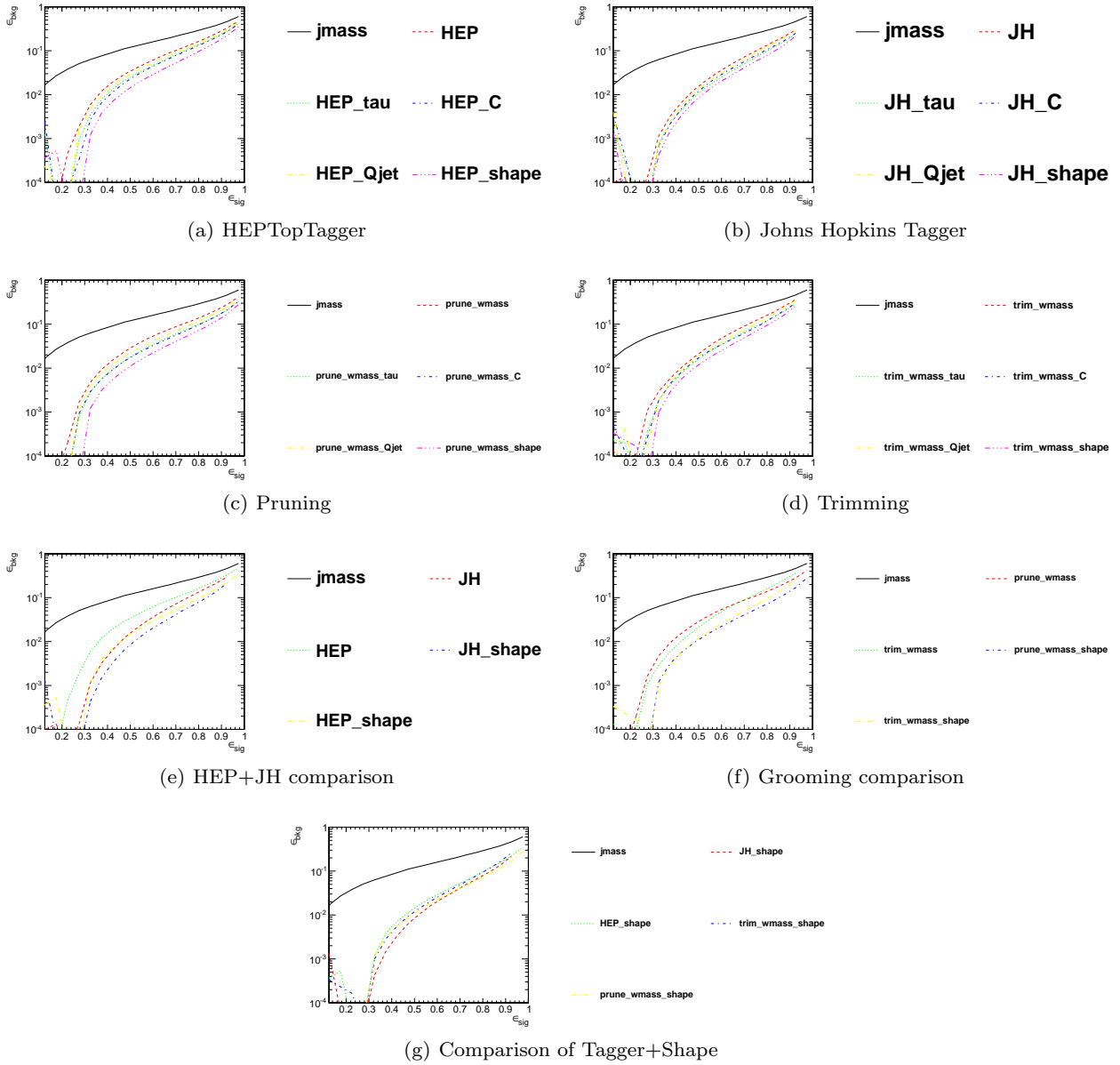


**Fig. 31** Comparison of top mass performance of different taggers at different  $R$  in the  $p_T = 1500 - 1600$  GeV bin.

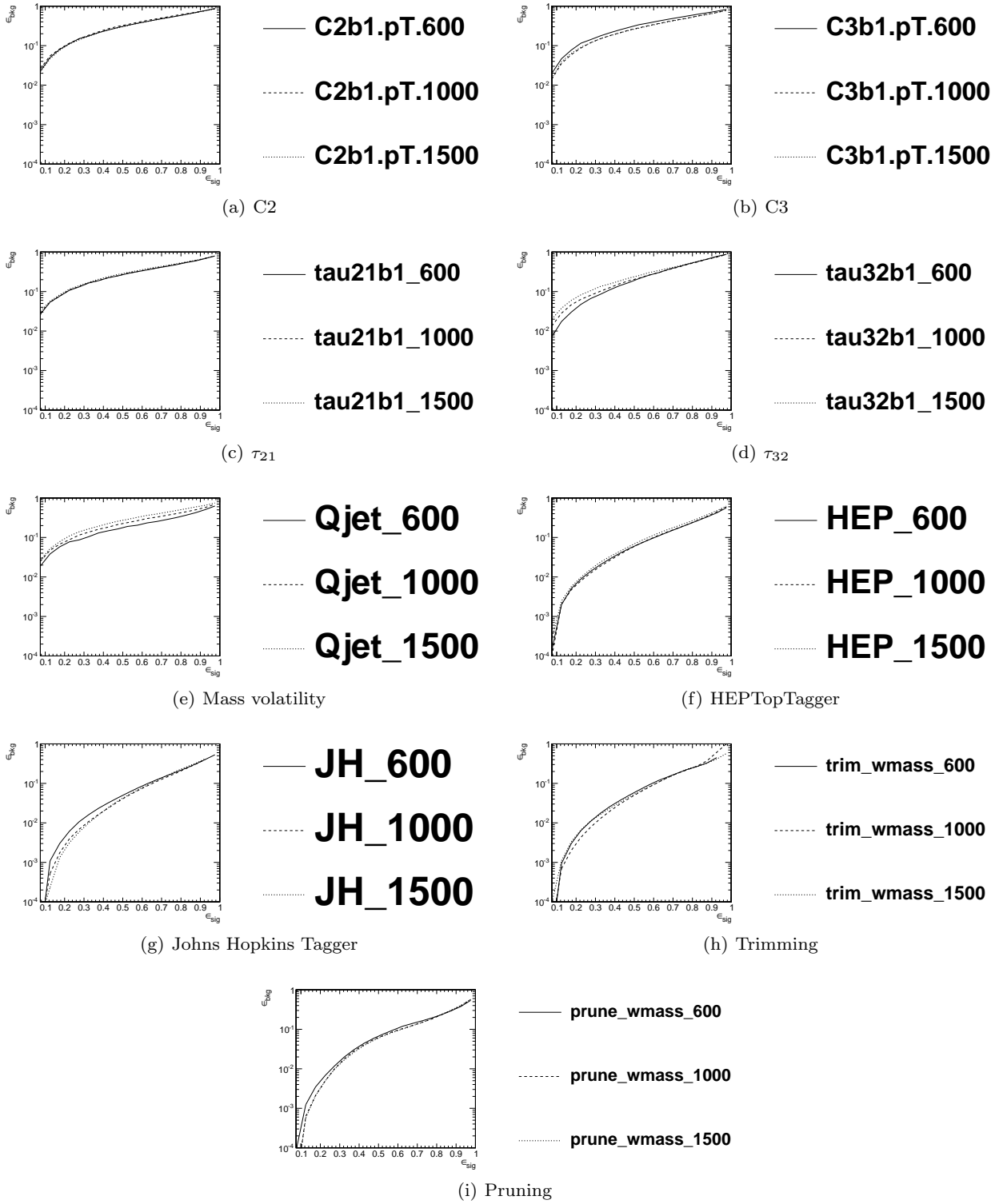


**Fig. 32** Comparison of  $W$  mass performance of different taggers at different  $R$  in the  $p_T = 1500 - 1600$  GeV bin.

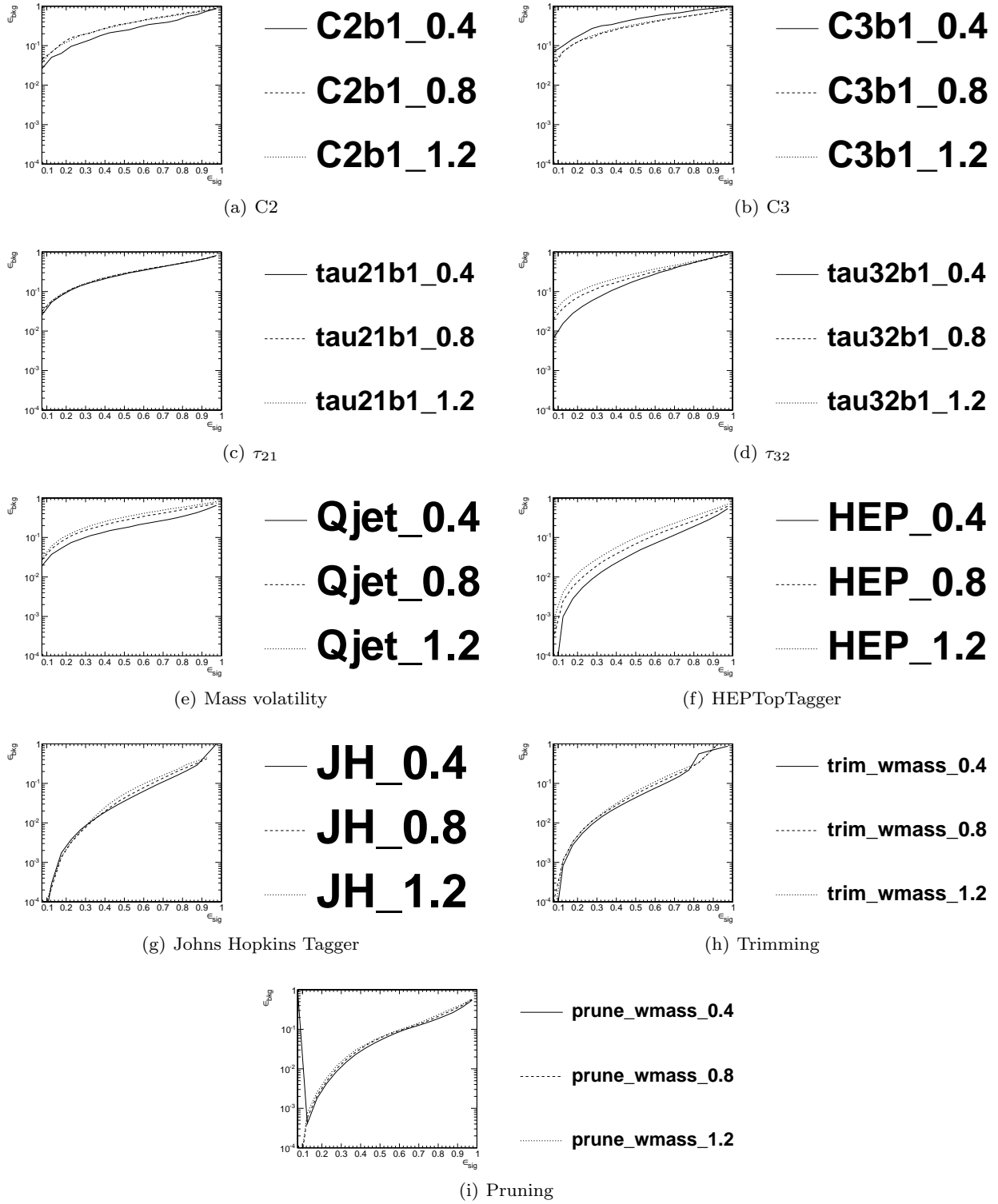




**Fig. 33** The BDT combinations in the  $p_T$  1000-1100 GeV bin using the anti- $k_T$   $R=0.8$  algorithm.



**Fig. 34** Comparison of tagger and jet shape performance at different  $p_T$  using the anti- $k_T$   $R=0.8$  algorithm.



**Fig. 35** Comparison of tagger and jet shape performance at different radius at  $p_T = 1.5\text{--}1.6$  TeV.

## 8 Summary & Conclusions

This report discussed the correlations between observables and looked forward to jet substructure at Run II of the LHC at 14 TeV center-of-mass collisions energies.

4. C. Anders, C. Bernaciak, G. Kasieczka, T. Plehn, and T. Schell, *Benchmarking an Even Better HEPTopTagger*, *Phys.Rev.* **D89** (2014) 074047, [[arXiv:1312.1504](#)].

## Acknowledgements

We thank the Department of Physics at the University of Arizona and for hosting the conference at the Little America Hotel. We also thank Harvard University for hosting the event samples used in this report. We also thank Hallie Bolonkin for the BOOST2013 poster design and Jackson Boelts' ART465 class (fall 2012) at the University of Arizona School of Arts VisCom program. (NEED TO ASK PETER LOCH FOR MORE ACKNOWLEDGEMENTS)

## References

1. A. Abdesselam, E. B. Kuutmann, U. Bitenc, G. Brooijmans, J. Butterworth, et al., *Boosted objects: A Probe of beyond the Standard Model physics*, *Eur.Phys.J.* **C71** (2011) 1661, [[arXiv:1012.5412](#)].
2. A. Altheimer, S. Arora, L. Asquith, G. Brooijmans, J. Butterworth, et al., *Jet Substructure at the Tevatron and LHC: New results, new tools, new benchmarks*, *J.Phys.* **G39** (2012) 063001, [[arXiv:1201.0008](#)].
3. A. Altheimer, A. Arce, L. Asquith, J. Backus Mayes, E. Bergeaas Kuutmann, et al., *Boosted objects and jet substructure at the LHC*, [arXiv:1311.2708](#).