# Double $B$-hadron Jet Tagging and Identification of Gluon to $b\bar{b}$ jets with the ATLAS Detector

Lic. María Laura González Silva

Tesis Doctoral en Ciencias Físicas

Facultad de Ciencias Exactas y Naturales

Universidad de Buenos Aires

Noviembre 2012

**UNIVERSIDAD DE BUENOS AIRES**

Facultad de Ciencias Exactas y Naturales

Departamento de Física

# Double $B$-hadron Jet Tagging and Identification of Gluon to $b\bar{b}$ jets with the ATLAS Detector

Trabajo de Tesis para optar por el título de

Doctor de la Universidad de Buenos Aires en el área Ciencias Físicas

por  **María Laura González Silva**

Director de Tesis: Dr. Ricardo Piegaia

Consejero de estudios: Dr. Daniel Deflorian

Lugar de Trabajo: Departamento de Física (CONICET-UBA)

Buenos Aires, 2012

## AGRADECIMIENTOS

Agradezco a...

**Abstract**

This thesis describes a method that allows the identification of double $B$-hadron jets originating from gluon-splitting. The technique exploits the kinematic differences between the so called "merged" jets and single $B$-hadron jets using track-based jet shape and jet substructure variables combined in a multivariate likelihood analysis. The ability to reject $b$-jets from gluon splitting is important to reduce and to improve the estimation of the $b$-tag background in Standard Model analyses and in new physics searches involving $b$-jets in the final state. In the simulation, the algorithm rejects 95% (50%) of merged $B$-hadron jets while retaining 50% (90%) of the tagged $b$-jets, although the exact values depend on the jet $p_T$.

# Contents

# Chapter 1

# Double $B$-hadron jet identification

## 1.1  Data and Monte Carlo samples

The tagging technique presented in this thesis relies on Monte Carlo predictions for the signal (single $b$) or background (merged $b$) hypotheses. The accuracy of the simulation is validated with data by comparing the distributions of the different variables explored.

Samples of jet events from proton-proton collision processes are simulated with PYTHIA8 [1] event generator using a $2 \rightarrow 2$ matrix element at leading order in the strong coupling to model the hard subprocess, and $p_T$-ordered parton showers are utilized to model additional radiation in the leading-logarithmic approximation [2]. Multiple parton interactions [3], as well as fragmentation and hadronisation based on the Lund string model [4] are also simulated. The ATLAS MC11 tune of the soft model parameters was used [5]. In order to have sufficient statistics over the entire $p_T$ spectrum, eight samples were generated with different thresholds of the hard-scattering

partonic transverse momentum $\hat{p}_T$. Events from different samples were mixed taking into account their respective production cross sections.

The GEANT4 [6] software toolkit within the ATLAS simulation framework [7] propagates the generated particles through the ATLAS detector and simulates their interactions with the detector material. The energy deposited by particles in the active detector material is converted into detector signals in the same format as the detector read-out. Finally the Monte Carlo generated events are processed through the trigger simulation package of the experiment, and are reconstructed and analyzed with the same software as for the real data. The simulated data sample used for the analysis gives an accurate description of the pile-up content and detector conditions for the full 2011 data-taking period.

The data samples employed correspond to proton-proton collisions at $\sqrt{s} = 7$ TeV delivered by the LHC and recorded by ATLAS between May and November 2011, with the LHC running with 50 ns bunch spacing, and bunches organized in bunch trains. Only data collected during stable beam periods in which all sub-detectors were fully operational are used. After the application of the data quality selection, the surviving data corresponds to an integrated luminosity of 4.7 fb$^{-1}$. The LHC performance steadily improved during 2011. In particular the average number of minimum-bias pile-up events, originating from collisions of additional protons in the same bunch as the signal collision, grew from from 3 to 20. This fact will be of importance when discussing the selection of discriminating variables.

For the study of systematic effects and for result comparison, other Monte Carlo samples were utilised. Results were produced with the HERWIG++ generator [8] and with PYTHIA8 using the Perugia tune [9]. The former is based on the event generator HERWIG, but redesigned in the $C++$ programming

language. The generator contains a few modlleing improvements. It also uses angular-ordered parton showers, but with an updated evolution variable and a better phase space treatment. Hadronisation is perfromed using the cluster model. The underlying event and soft inclusive interactions are described using a hard and soft multiple partonic interactions model [10]. The Perugia tune is an independent tune of PYTHIA with increased final state radiation to better reproduce the jet shapes and hadronic event shapes using LEP and TEVATRON data. In addition, paramteres sensitive to the production of particles with strangeness and related to jet fragmentation have been adjusted.

## 1.2    Data analysis

Jets are reconstructed using the anti-$k_t$ jet algorithm [11] with a distance parameter $R = 0.4$, using calorimeter topological clusters [12] as input. Several quality criteria are applied to eliminate "fake" jets that are caused by noise bursts in the calorimeters and energy depositions belonging to a previous bunch crossing [13].

The jet energies are corrected for inhomogeneities and for the non-compensating nature of the calorimeter by using $p_T - and\eta$-dependent calibration factors determined from Monte Carlo simulation [14]. This calibration is referred to as the EM+JES scale. Using test beam results, in-situ track and calorimeter measurements, estimations of pile-up energy depositions, and detailed Monte Carlo comparisons, an uncertainty on the absolute jet energy scale was estabished. This uncertainty is smaller than $\pm 10\%$ for $\eta < 2.8$ and $p_T > 20$ GeV. More sophisticated techniques undergoing commissioning, such as local cluster weighting, are expected to considerably improve the jet energy uncertainty and resolution [15].

4

The event sample for this analysis was collected using a logical OR of single jet triggers which select events with at least one jet with transverse energy above a given threshold at the highest trigger level. The ATLAS Trigger system uses three consecutive trigger levels. At the hardware Level 1 and local software Level 2, cluster-based jet triggers are used to select events. The Level 3, the so-called Event Filter, runs the offline anti-$k_t$ jet finding algorithm with $R = 0.4$ on topological clusters over the complete calorimeter. At this stage, the transverse energy thresholds, expressed in GeV, are: 20, 30, 40, 55, 75, 100, 135, 180. These triggers reach an efficiency of 99% for events having the leading jet with an offline energy higher than the corresponding trigger thresholds by a factor ranging between 1.5 and 2. The triggers with the lowest $p_T$ thresholds were prescaled by up to five orders of magnitude, and typically the same jet trigger is prescaled ten times more in the later data taking periods compared to the early ones.

The offline event selection requires at least one primary vertex candidate with 5 or more tracks. All jets, with transverse momentum between 40 and 480 GeV, were required to be in a region with full tracking coverage, $|\eta_{jet}| < 2.1$, and they were classified in eight $p_T$ bins chosen such as to match the jet trigger 99% efficiency thresholds (in GeV): 40, 60, 80, 110, 150, 200, 270, 360. Only jets tagged as $b$-jets using the MV1 $b$-tagging algorithm at the 60% efficiency working point were considered. $b$-tagged jets with close-by jets ($\Delta R < 0.8$) with $p_T$ higher than 7 GeV at electromagnetic scale were not included in the analysis. Unless otherwise indicated, performance plots are shown for one medium $p_T$ bin (80 to 110 GeV) and one high $p_T$ bin (200 to 270 GeV).

In the case of MC the reconstructed $b$-tagged jets were further classified into single and merged $b$-jets based on truth Monte Carlo information. A

$B$-hadron is considered to be associated to a jet if the $\Delta R$ distance in $\eta - \phi$ space between the direction of the hadron and the jet axis is smaller than 0.4. Jets were labeled as merged (single) $b$-jets if they contain two (only one) $B$-hadron.

It is important to select genuine tracks belonging to jets. Only tracks located within a cone of radius $\Delta R(jet^{\mathrm{reco}}, \mathrm{track}) \leq 0.4$ around the jet axis were considered. Cuts on $p_{\mathrm{T}}^{\mathrm{trk}} > 1.0$ GeV and the $\chi^2$ of the track fit, $\chi^2/ndf < 3$, are applied. In addition, tracks are required to have a total of at least seven precision hits (pixel or micro-strip) in order to guarantee at least 3 $z$-measurements. Tracks are also required to fulfill cuts on the transverse and longitudinal impact parameters at the perigee to ensure that they arise from the primary vertex. As cutting on impact parameter (IP) significance might be detrimental for $b$-jets, where large IP values are expected, the relaxed cuts were used, $|IP_{xy}| < 2$ mm, and $|IP_z \sin\theta| < 2$ mm, with $\theta$ being the polar angle measured with respect to the beam axis.

## 1.3 Preliminary studies in standalone PYTHIA generator

A small set of dijet events was generated using sc pythia. With the help of fastjet [16], anti-$k_t$ jets with distance parameters $R$ of 0.4 and 0.6 were built using as inputs:

- all stable particles
- charged stable particles
- 0.1 $k_t$ jets, using all particles
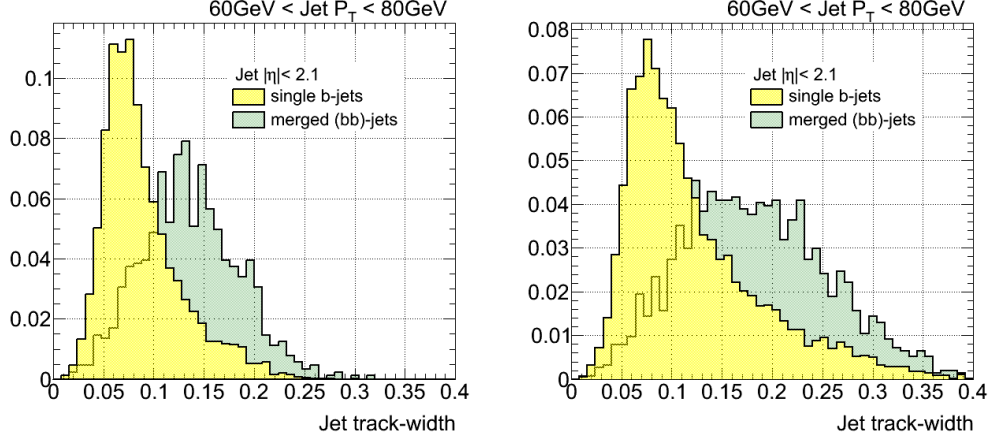- 0.1 $k_t$ jets, using charged particles

Figure 1.1: Distribution of track-jet width in anti-$k_T$ 0.4 (left) and 0.6 (right) jets, for single and merged $b$-jets between 60 GeV to 80 GeV. Jets were built using all stable particles in the simulation.

The labeling was done in the same procedure as in the Full ATLAS Monte Carlo analysis.

Figures 1.1 to 1.5 show the distributions and correlations of some of the tracking variables, for single and merged $b$-jets, in bins of the jet $p_T$.

For the single $b$-jets, $\tau_1$, $\tau_2$, and $\Delta R$ between the $k_T$ axes in the jet are all small which is expected for a pencil-like jet. For the $b\bar{b}$-jets, these variables are all large, which is typical of a gluon jet. But the correlations are really fascinating in merged $b$-jets. $\tau_1$ and $\Delta R$ between the $k_T$ axes are nearly linearly related, which is expected if there are two hard lobes of energy. But $\tau_2$ is almost independent of $\Delta R$ between the $k_T$ axes, meaning that regardless of where the axes are, the energy is uniformly distributed around them.
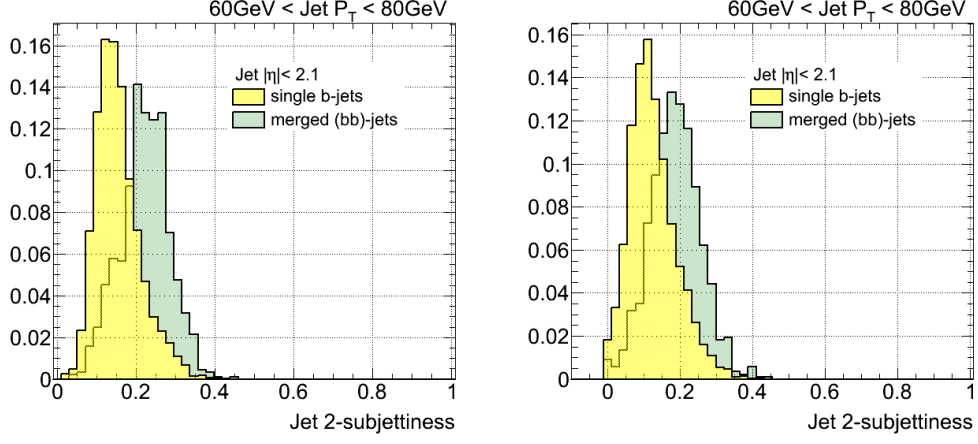
7

Figure 1.2: Distribution of $\tau_2$ in anti-$k_T$ 0.4 jets, for single and merged $b$-jets between 60 GeV to 80 GeV. Jets were built using all stable particles (left) and charged particles only (right).
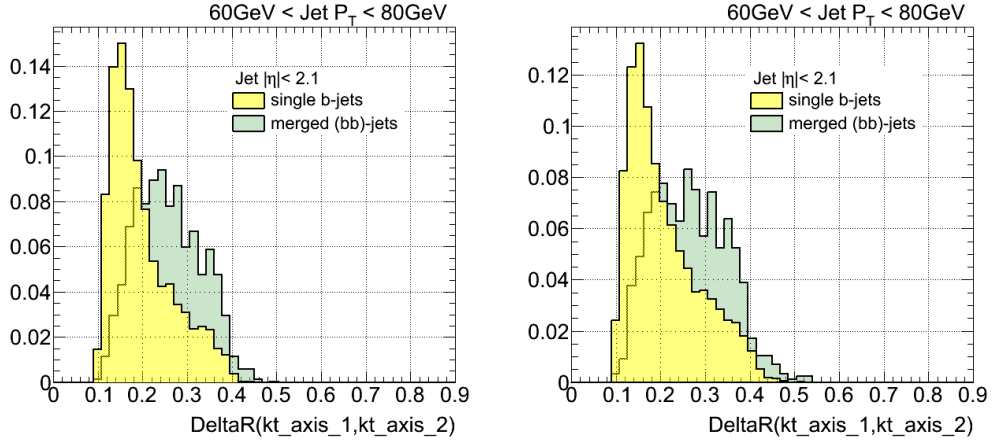


Figure 1.3: Distribution of $\Delta R$ between the axes of two $k_T$ subjets in anti-$k_T$ 0.4 jets, for single and merged $b$-jets between 60 GeV to 80 GeV. Jets were built using 0.1 $k_t$ jets from all stable particles (left) and charged particles only (right).
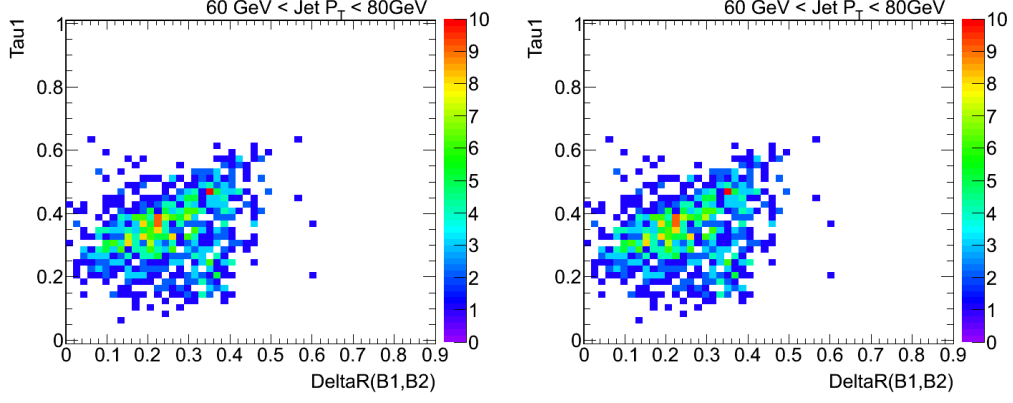
Figure 1.4: Correlation between $\tau_1$ (left) and $\tau_2$ (right) and the $\Delta R$ between the $B$-hadrons in merged anti-$k_T$ 0.4 jets between 60 GeV to 80 GeV. Jets were built using all stable particles.
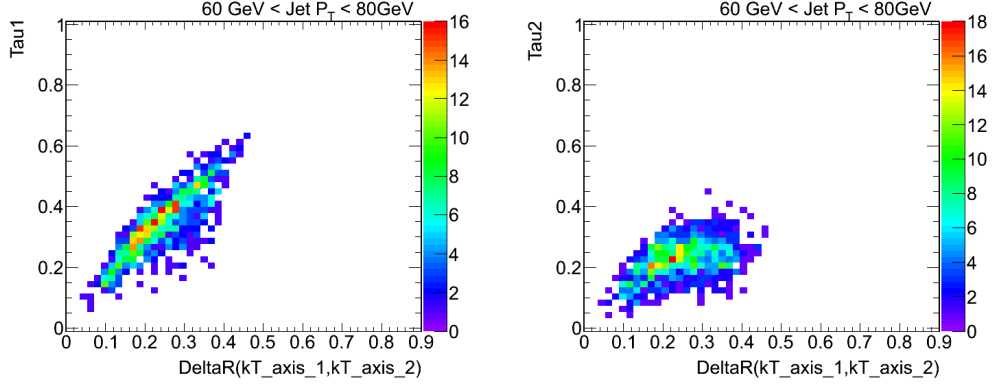


Figure 1.5: Correlation between $\tau_1$ (left) and $\tau_2$ (right) and the $\Delta R$ between the $k_T$ subjets in merged anti-$k_T$ 0.4 jets between 60 GeV to 80 GeV. Jets were built using all stable particles.

# 1.4 Kinematic differences between single and double $B$-hadron jets

The differences between genuine $b$-quark jets and $b\bar{b}$ jets are expected to arise from the two-subjet (two $B$-hadrons) substructure of merged jets. They are thus expected, for the same jet $p_T$, to have higher track-multiplicity and be wider than single $b$-jets. Based on these characteristics simulated QCD samples of $b$-tagged jets were used to study the following properties, discussed in the next paragraphs, built from jet constituents either at calorimeter level (topological clusters) or tracks associated to the jet:

- Jet multiplicity (number of constituents)
- Jet width, $p_T$ weighted
- Jet Mass
- Nr. of $k_t$ subjets
- Maximum $\Delta R$ between pairs of constituents
- $\Delta R$ between 2 $k_t$ subjets within the $b$-jet
- $\tau_2$: 2-subjettiness
- $\tau_2/\tau_1$
- $\Delta R$ of leading constituents
- Eccentricity

*I. Jet track multiplicity*

This variable is defined as the number of tracks associated to the jet, it is simple to calculate and carries important information of the jet inner structure. Figure 1.6 shows the distribution of the observable for single and merged $b$-jets. It was observed that merged $b$-jets contain on average around two
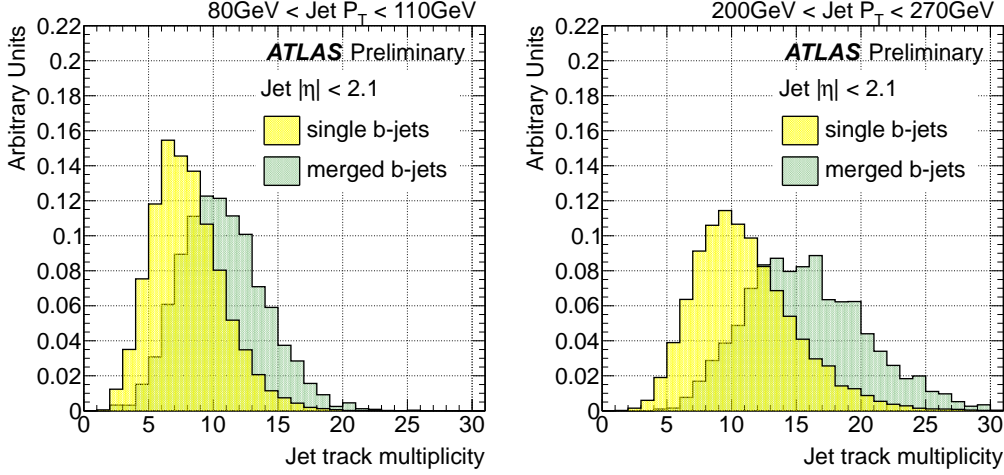
Figure 1.6: Distribution of the track multiplicity in jets for single and merged *b*-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).

more tracks than single *b*-jets at low jet $p_T$, with a larger difference at higher $p_T$ values. The jet track multiplicity corresponds to tracks with $p_T$ above 1 GeV, satisfying the quality cuts described in section 1.2. The effect of using a minimum track $p_T$ of 0.5 GeV was also examined. This was motivated by the fact that it could lead to an improvement in discrimination if it captured more information about the fragmentation process. On ther other hand, a lower minimum track $p_T$ can make the method more sensitive to pile-up with the addition of soft tracks incorrectly associated to the jets. What it was observed is that reducing the $p_T$ cut only widens the distributions without increasing the separation between single and merged jets.

*II. Jet width*

The jet width was computed as the $p_T$ weighted average of the $\Delta R$ distance between the associated constituent ("*const*") and the jet axis:

$$Jet\ width = \frac{\sum_{i=1}^{N} p_T{}^{const_i}\, \Delta R(const_i, jet)}{\sum_{i=1}^{N} p_T{}^{const_i}} \qquad (1.1)$$

11

where $N$ is the total number of calorimeter or track constituents.

Figure 1.7 shows the distribution for the Track-jet width. As expected, merged $b$-jets are wider than single $b$-jets. In Fig. 1.8 the correlation between the track-jet width and the jet track multiplicity is shown for single and merged $b$-jets. These two variables alone provide a good discrimination for tagging $b\bar{b}$ jets.

The calorimeter jet width ( using topological clusters) gives also good separation. However, this variable is more sensitive to the amount of pile-up in the event than its track-based counterpart. In Fig. 1.9 the distributions of calorimeter width for single and merged $b$-jets can be seen for events with low and high Number of Primary Vertices (NPV), in a low $p_T$ region where the effect of pile-up is more important. In Fig. 1.10 the same distributions are shown for the track-jet width. Calorimeter jet width varies with NPV and due to this behavior the track-based version is more suitable as a more robust discriminator. For similar reasons, the jet topological cluster multiplicity and the jet mass were discarded as discriminating variables.

### III. Maximum $\Delta R$ between track pairs

Figure 1.11 shows the distribution of the maximum $\Delta R$ between track pairs in the jets ($\mathrm{Max}\{\Delta R(trk, trk)\}$). Merged $b$-jets show significantly higher values for this variable over a broad range of jet $p_T$. The distinct characteristic of this variable is that the separation between single $b$-jets and merged does not depend on jet $p_T$. In spite of its good discrimination power, we have looked for alternatives to $\mathrm{Max}\{\Delta R(trk, trk)\}$ as it is not an infrared safe observable and is sensitive to soft tracks originating from pile-up.
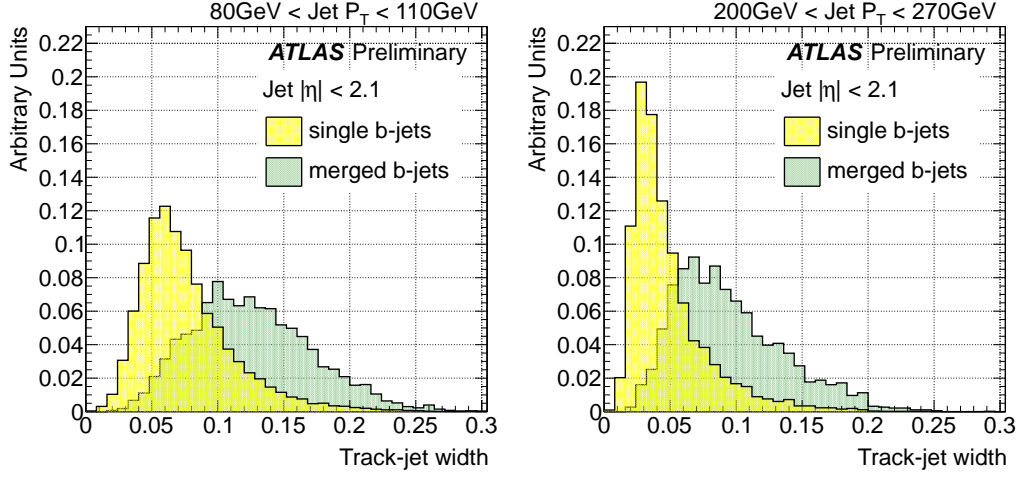
Figure 1.7: Distribution of track-jet width in jets for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).
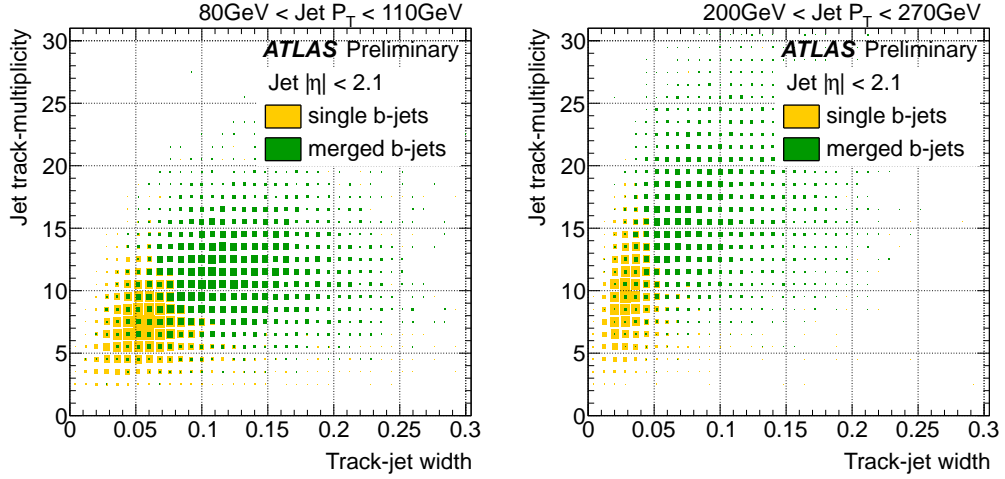


Figure 1.8: Correlation between jet track multiplicity and track-jet width for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).
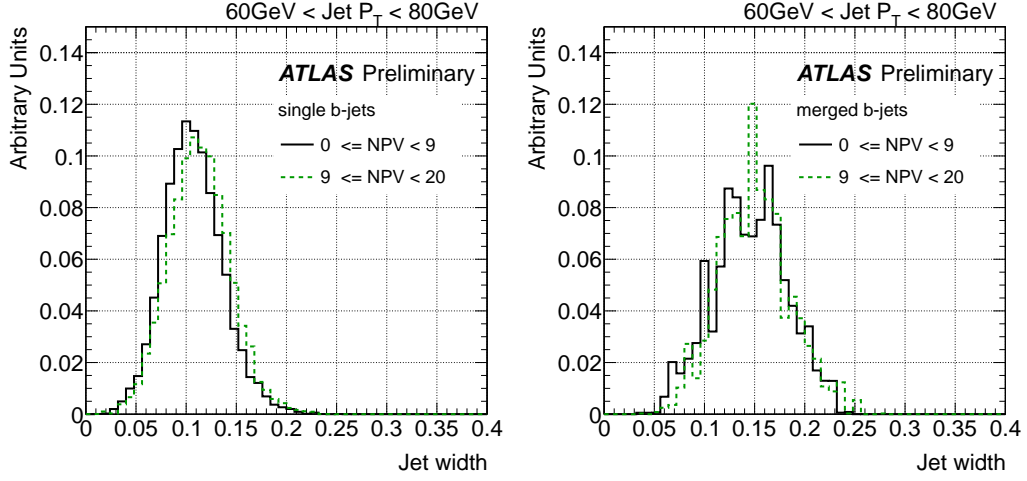
Figure 1.9: Distribution of calorimeter jet width (using topological clusters) for single (left) and merged (right) $b$-jets in two bins of Number of Primary Vertices for jets between 60 GeV to 80 GeV.



Figure 1.10: Distribution of track-jet width for single (left) and merged (right) $b$-jets in two bins of Number of Primary Vertices for jets between 60 GeV to 80 GeV.
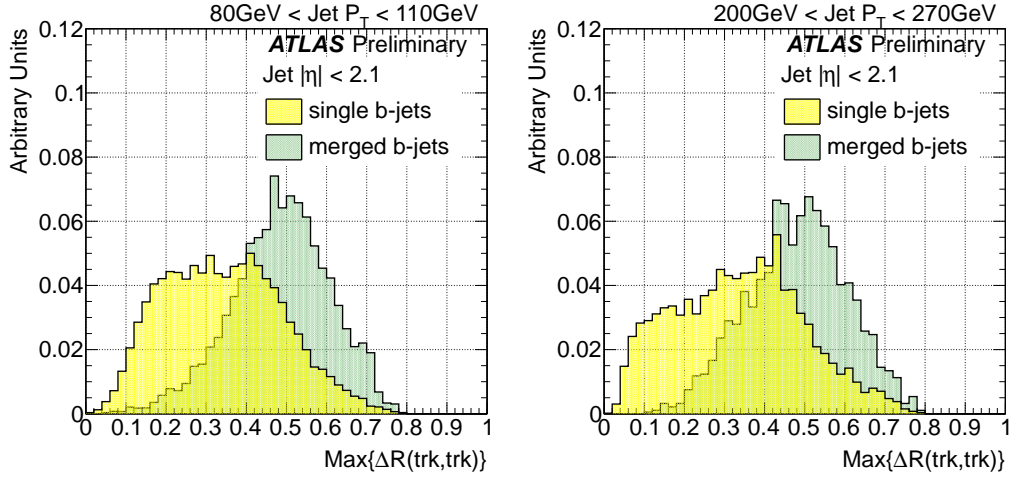
Figure 1.11: Distribution of the maximum $\Delta R$ between pairs of tracks in jets for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).

*IV. $\Delta R$ between the axes of two $k_t$ subjets*

The distribution of the $\Delta R$ between the axes of the two exclusive $k_t$ subjets in the jet is shown in Fig. 1.12 for single and merged $b$-jets. In order to build this variable the $k_t$ algorithm [17] is applied to all the tracks associated to the jet using a large $k_t$ distance parameter to ensure that all of them get clustered. The clustering is stopped once it reaches exactly two jets. We observe that this variable also provides good separation, with the advantage of infrared safeness and insensitivity to pile-up.

*V. N-subjettiness variables*

$N$-subjettiness variables, as described in Ref. [18], were originally designed to identify boosted objects, like electroweak bosons and top quarks,
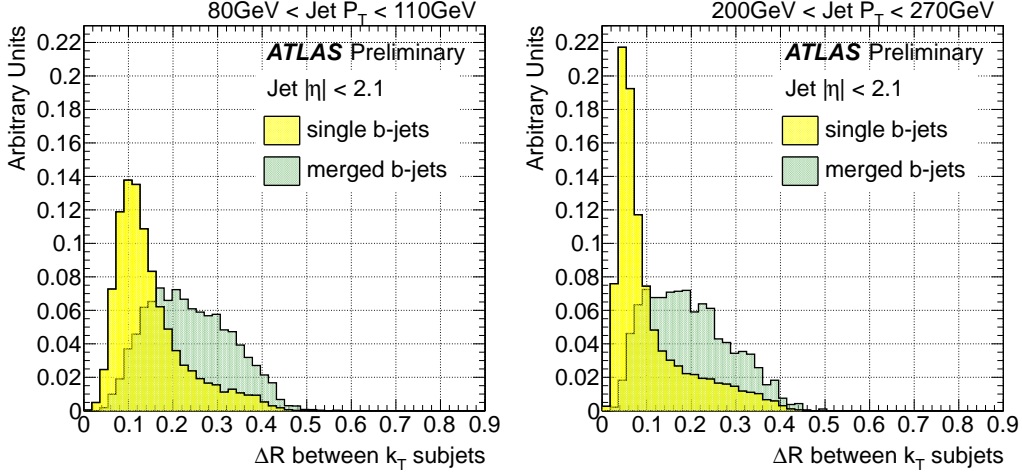
Figure 1.12: Distribution of the $\Delta R$ between the axes of the two $k_t$ subjets in the jet for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).

decaying into collimated shower of hadrons which a standard jet algorithm would reconstruct as single jets. It is defined as:

$$\tau_N = \frac{1}{\sum_k p_{Tk} R_0} \sum_k p_{Tk} \min\{\Delta R_{S_1,k}, \Delta R_{S_2,k}, ..., \Delta R_{S_N,k}\} \qquad (1.2)$$

where $R_0$ is the jet radius used in the jet clustering algorithm and the sum runs over the constituents of the jet. To avoid dependence on pile-up we consider the track-based $n$-subjettiness, where the sum is over the tracks in the $b$-tagged jet. $\Delta R_{S_j,k}$ is the distance in the rapidity-azimuth plane between the axis of subjet $j$ and constituent track $k$. This jet shape variable quantifies to what degree a jet can be regarded as composed of $N$ subjets. For instance, a jet with a two pronged structure, with all tracks clustered along two directions, is expected to have a smaller $\tau_2$ value than a jet with tracks uniformly distributed in $\eta - \phi$ space.

Plots of $\tau_2$ are shown in Fig. 1.13. In spite of its expected 2-prong substructure, merged $b$-jets have higher values of $\tau_2$ than single $b$-jets. The
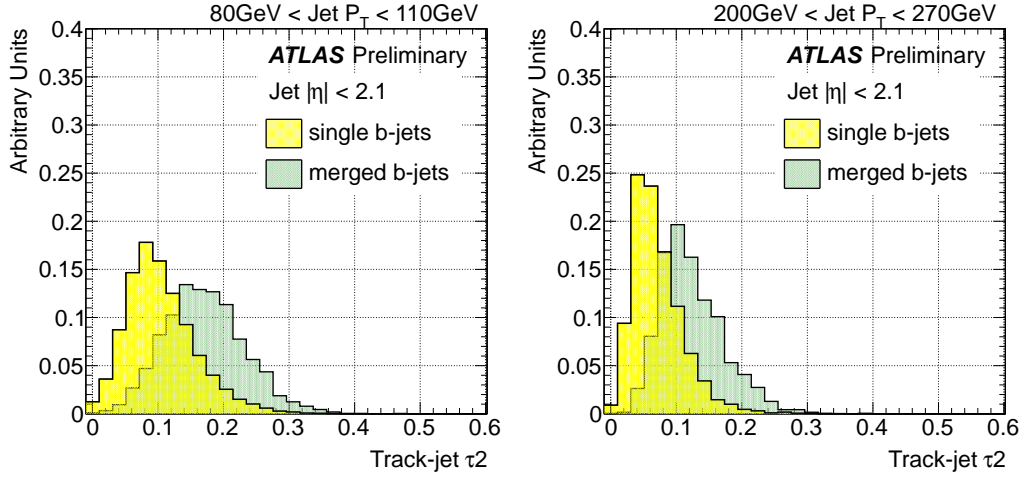
Figure 1.13: Distribution of $\tau_2$ in jets for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).

explanation of this behavior can be found in Fig. 1.14, where its correlation with track-jet width ($\sim \tau_1$) is shown for single and merged $b$-jets. The two variables are highly correlated and for this reason wider jets have a larger $\tau_2$. This suggests to switch from an absolute to a width-normalized $\tau_2$. Fig. 1.15 thus shows the distributions of $\tau_2/\tau_1$. This ratio is often used but, although as expected somewhat larger values are obtained for single than for merged $b$-jets, specially at high $p_T$, we decided not to use this variable as it offers only marginal discrimination.

*VI. Jet Mass*

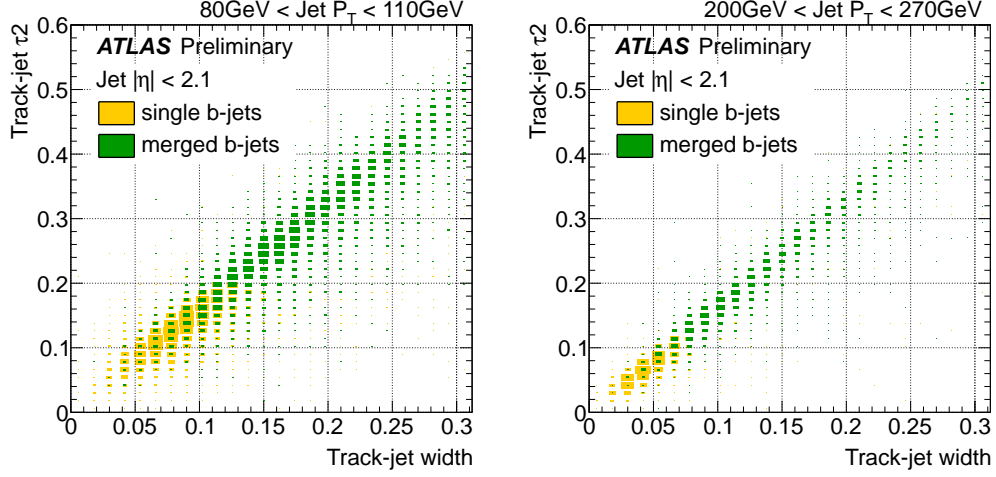Figure 1.16 shows the distribution of the jet mass for single and merged $b$-jets.

Figure 1.14: Correlation between $\tau_2$ and track-jet width for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).
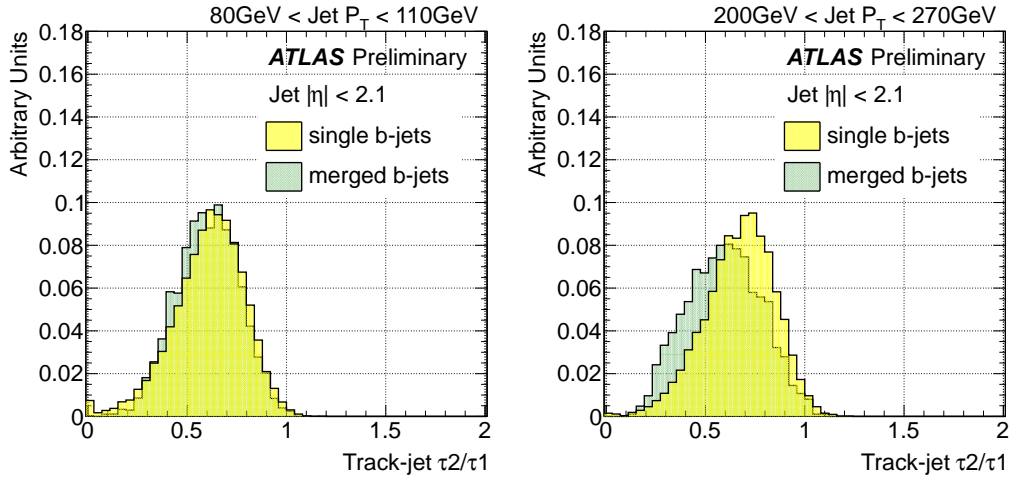


Figure 1.15: Distribution of $\tau_2/\tau_1$ in jets for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).
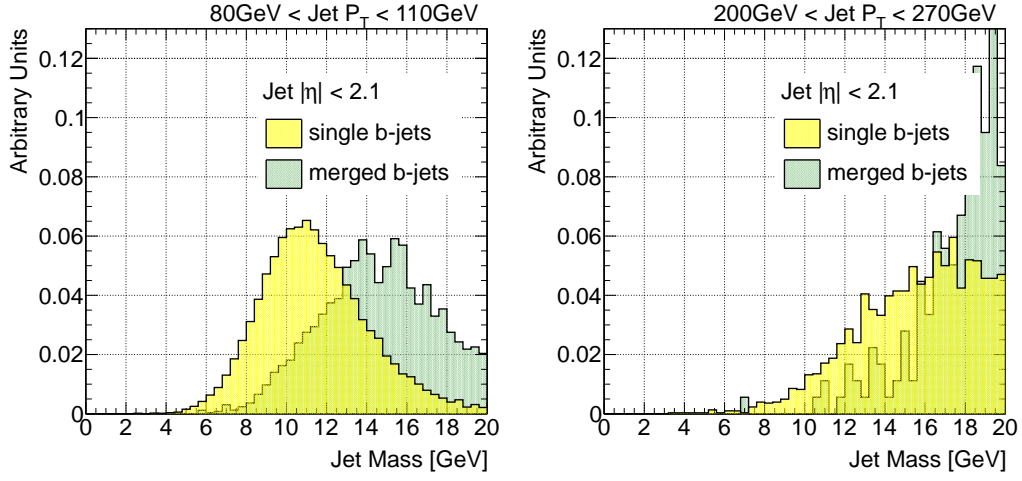
Figure 1.16: Distribution of jet mass in GeV for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).

*VII. Number of $k_t$ subjets*

Figure 1.17 shows the distribution of the number of sub-track-jets single and merged $b$-jets.

*VIII. $\Delta R$ between leading constituents*

Figure 1.18 shows the distribution of the number $\Delta R$ between leading tracks in the jet for single and merged $b$-jets.
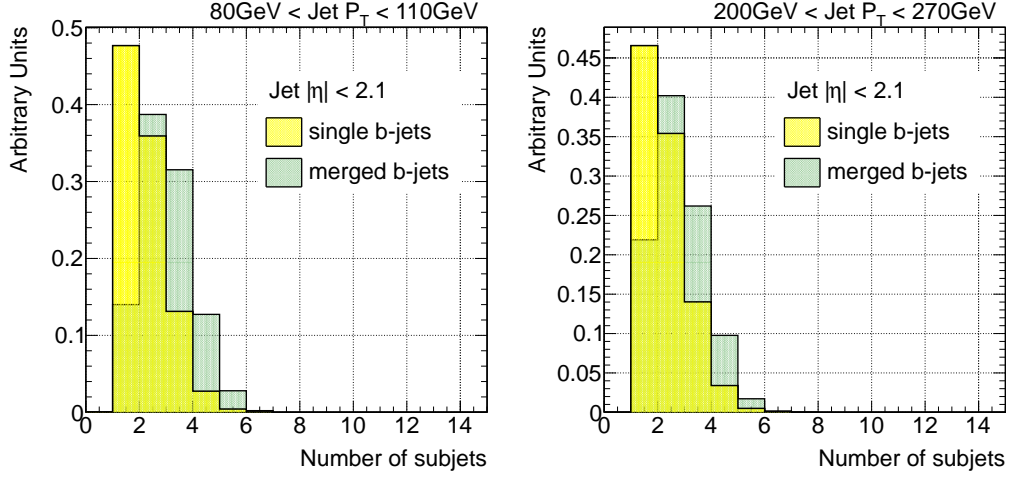
*IX. Jet eccentricity*

19

Figure 1.17: Distribution of the number of $k_t$ sub-track-jets for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).
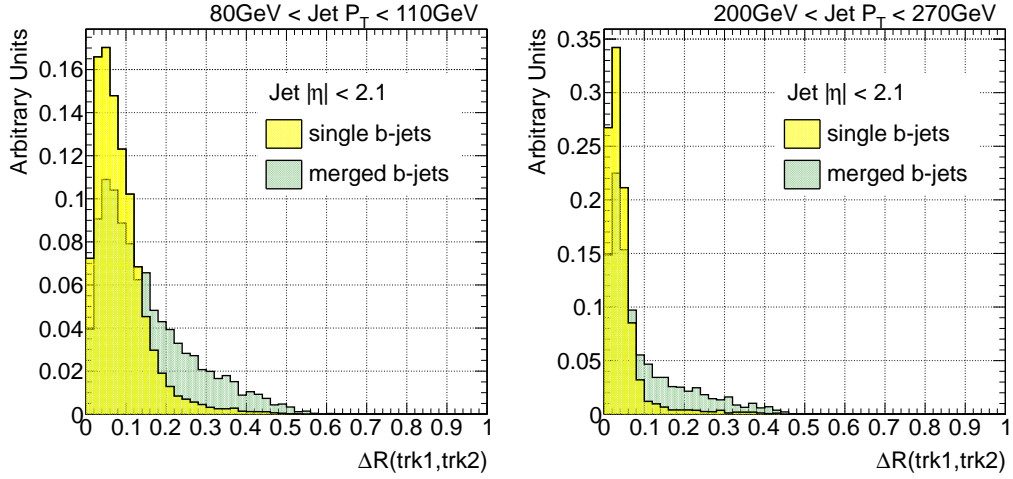


Figure 1.18: Distribution of $\Delta R$ between leading tracks for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).
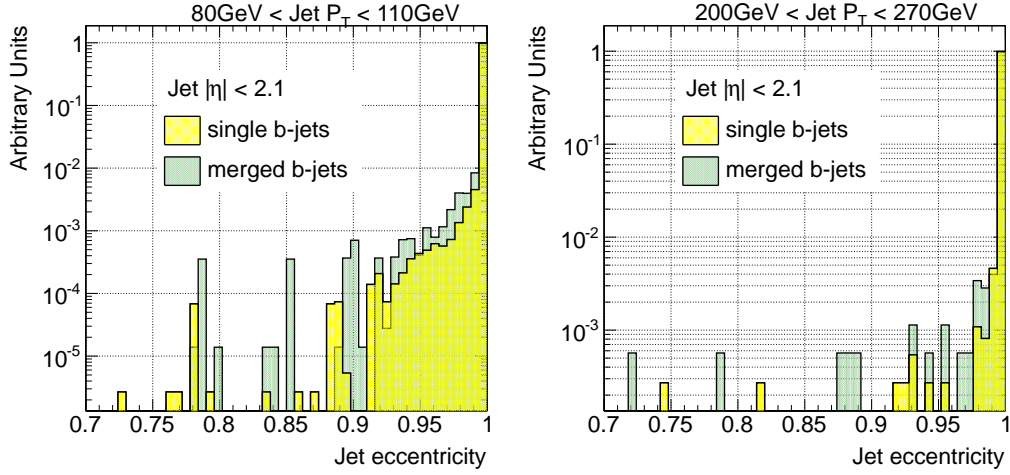
Figure 1.19: Distribution of the jet eccentricity for single and merged $b$-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right).

Figure 1.19 shows the distribution of the jet track-eccentricity for single and merged $b$-jets.

We also explored the potential improvement of constructing kinematic variables with only displaced tracks, as these are the ones expected to arise from the decay of B-hadrons. Cuts of 2, 2.5 and 3 on the track transverse impact parameter significance were investigated leading however to no gain in discrimation power.

In Figures 1.20 and 1.21 two examples are shown.

## 1.4.1 Further studies using "ghost-association" and bigger cone jets

In order to better understand the behavior observed for $\tau_2$, $\Delta R$ between the axes of $k_T$ subjets and jet eccentricity in anti-$k_T$ 0.4 jets, these variables were
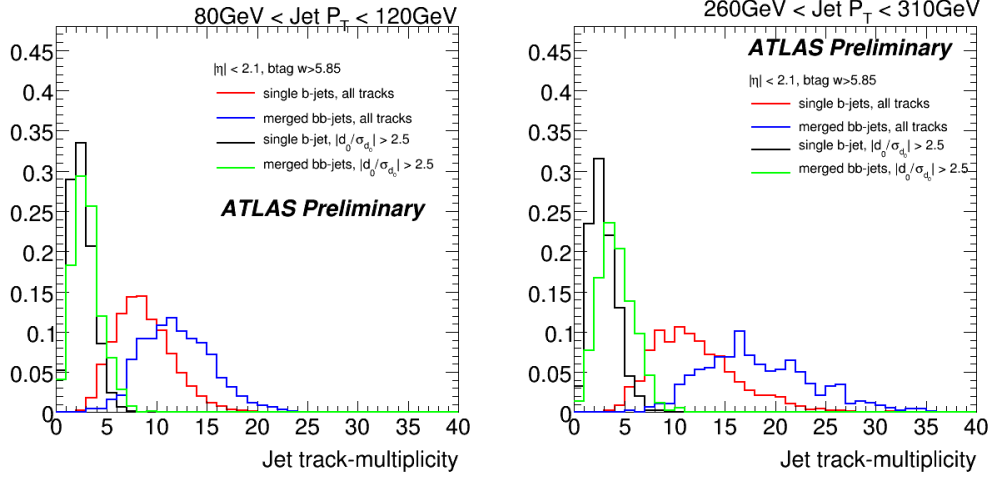
Figure 1.20: Distribution of the jet track multiplicity single and merged *b*-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right), for all and displaced tracks only.
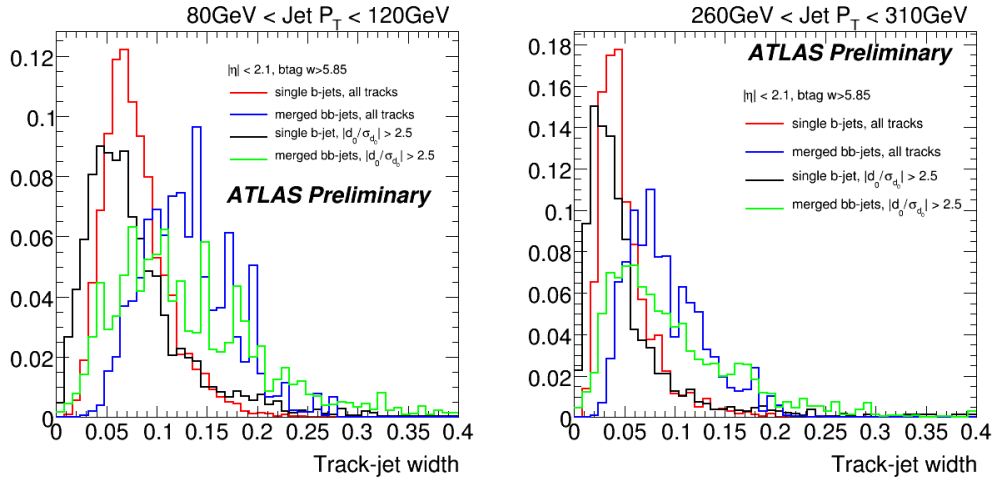


Figure 1.21: Distribution of the track-jet width for single and merged *b*-jets between 80 GeV to 110 GeV (left) and 200 GeV to 270 GeV (right), for all and displaced tracks only.
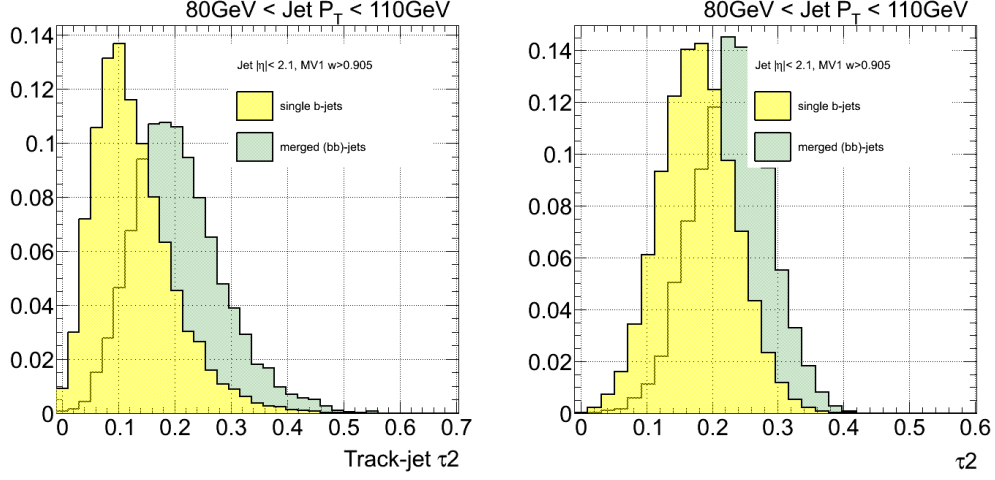
Figure 1.22: Distribution of $\tau_2$ for single and merged $b$-jets between 80 GeV to 110 GeV in anti-$k_T$ 0.6 jets using track constituents (left) and anti-$k_T$ 0.4 jets using the active area of the jet, with calorimeter topoclusters as input.

studied for other two different scenarios,

- using the active area of jets (with clusters used as input to jet reconstruction).
- using bigger 0.6 anti-$k_T$ jets

in order to enhance the efficiency to capture the decay products in gluon to $b\bar{b}$-jets.

Figures 1.22 to 1.24 show distributions of variables mentioned above for single and merged $b$-jets between 80 GeV to 110 GeV.

## 1.5 Validation of the jet variables in data

In order to study the extent to which the simulation reproduces the distributions observed in data for the different variables explored a set of comparison
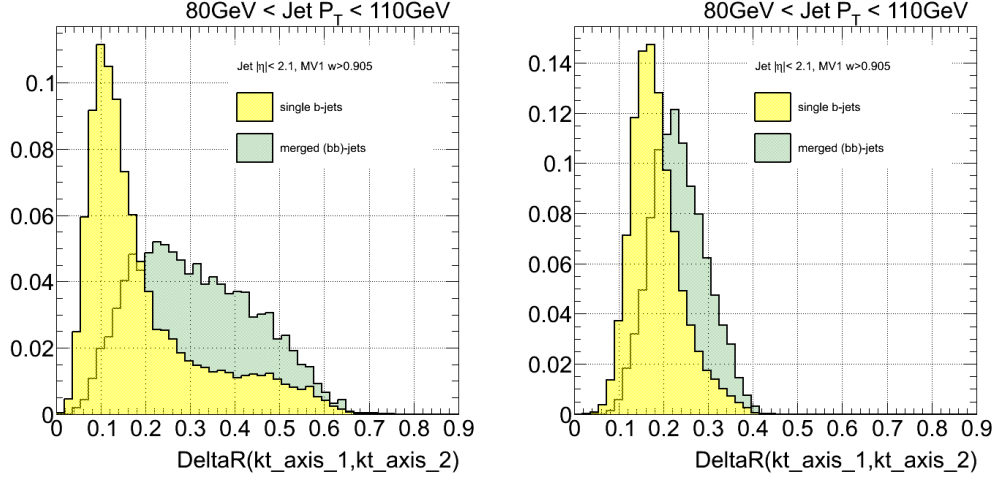
Figure 1.23: Distribution of $\Delta R$ between $k_T$ subjets for single and merged $b$-jets between 80 GeV to 110 GeV in anti-$k_T$ 0.6 jets using track constituents (left) and anti-$k_T$ 0.4 jets using the active area of the jet, with calorimeter topoclusters as input.
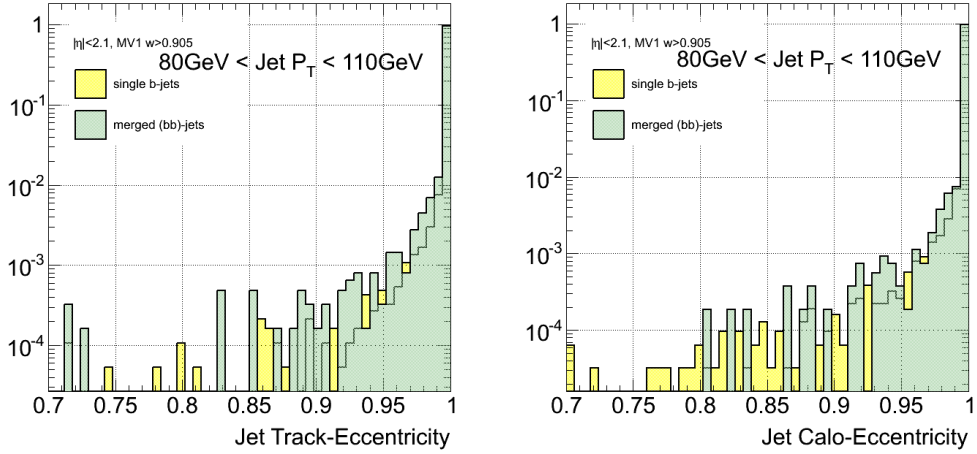


Figure 1.24: Distribution of the jet eccentricity for single and merged $b$-jets between 80 GeV to 110 GeV in anti-$k_T$ 0.6 jets using track constituents (left) and anti-$k_T$ 0.4 jets using the active area of the jet, with calorimeter topoclusters as input.

plots is presented. Fig. 1.25 shows the distributions of jet track multiplicity, track-jet width and $\Delta R$ between the axes of the two $k_t$ subjets, in two different jet $p_T$ bins in dijet Monte Carlo and data events collected by ATLAS during 2011. The distributions are normalized to unit area to allow for shape comparisons. There is a good agreement between data and simulation. It should be remarked that the observed agreement is actually not a direct validation of the description in the MC of the relevant variables, but its convolution with the simulated relative fractions of light-, $c$-, $b$- and $bb$-jets in the $b$-tagged generated jet sample. To some extent, some level of compensation can take place between these two effects.
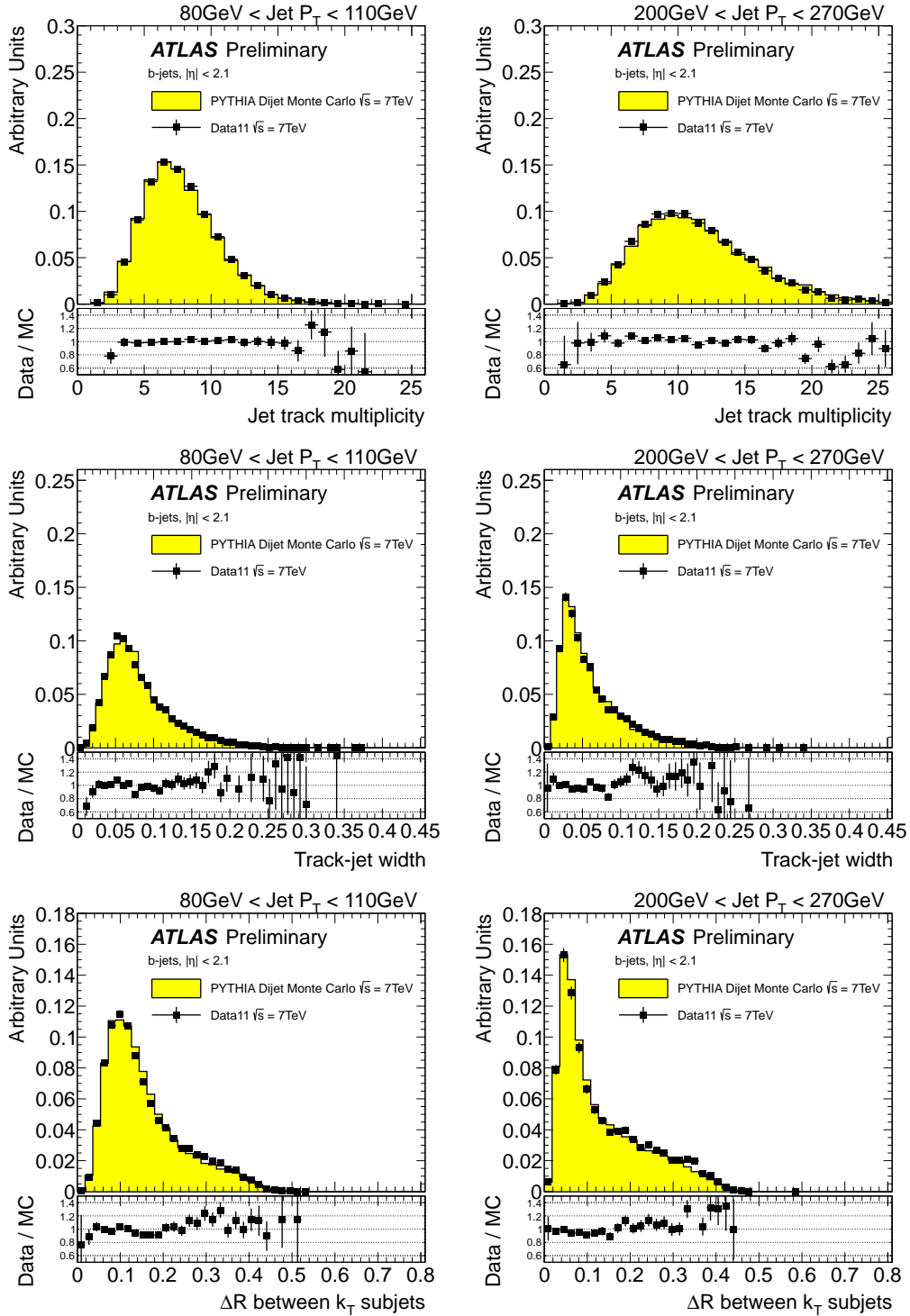
Figure 1.25: Distribution of three tracking variables in 2 different jet $p_T$ bins, for experimental data collected by ATLAS during 2011 (solid black points), and simulated data (filled histograms). The ratio data over simulation is shown at the bottom of each plot.

26

# Chapter 2

# Multivariate Analysis

## 2.1 The multivariate classifiers

The following multivariate methods were explored:

- Likelihood ratio estimators
- Neural Networks (NN)
- Boosted decision Trees (BDTs)

And different trainings were tested:

- Inclusive, with $p_T$-weighting
- In bins of jet $p_T$

Signal and background jets were not weighted by the dijet samples cross-sections to allow the contribution of subleading lower $p_T$ jets from high $p_T$ events, and thus increase the statistics of merged jets in the low $p_T$ bins.

Figure 2.1 and 2.2 show distributions of the MVA outputs in different bins of jet $p_T$ for the two proposed trainings. In figures 2.3 and 2.4 a comparison of the performance of all methods, for inclusive and "in-bins", training is illustrated.
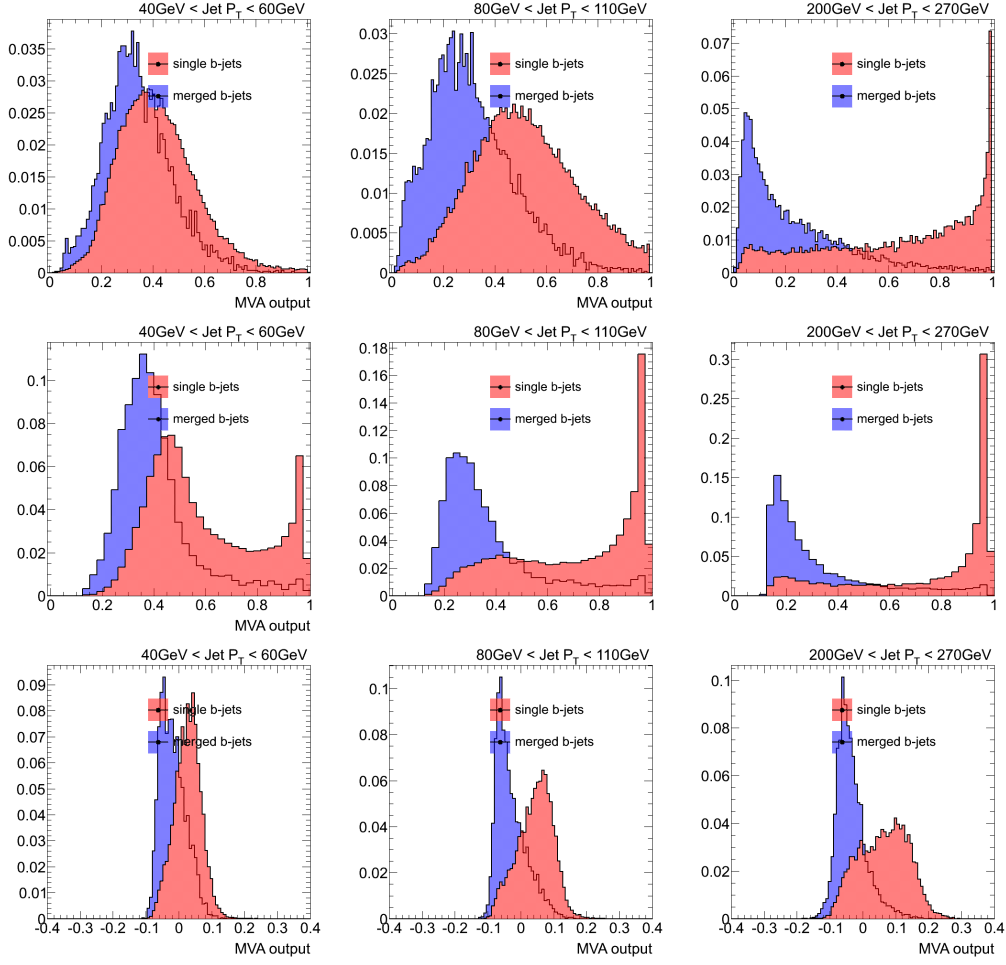
Figure 2.1: Distribution of the MVA discriminant outputs, for inclusive training, in single and merged $b$-jets, for low, medium and high jet $p_T$.

Figure 2.2: Distribution of the MVA discriminant outputs, for training in bins of jet $p_T$, in single and merged $b$-jets, for low, medium and high jet $p_T$.

Figure 2.3: Distribution of the MVA discriminant performance for inclusive training, in single and merged $b$-jets, for low and high jet $p_T$.
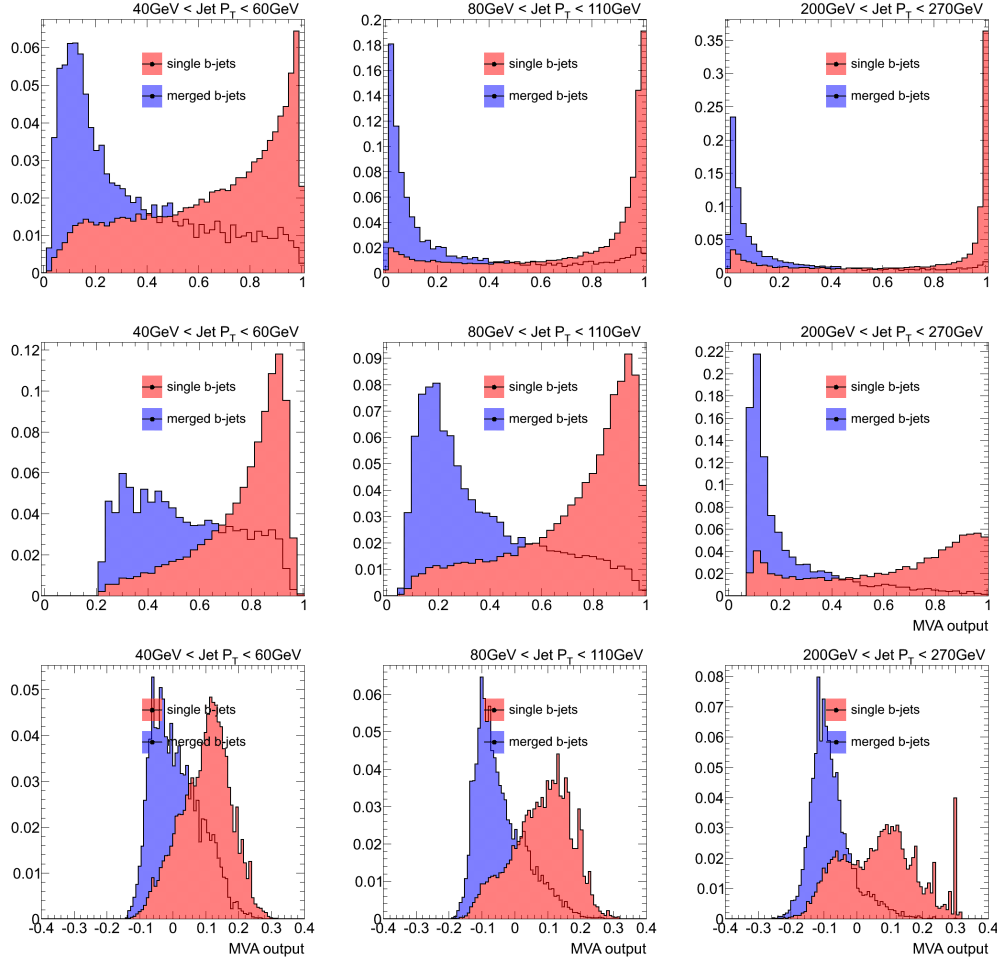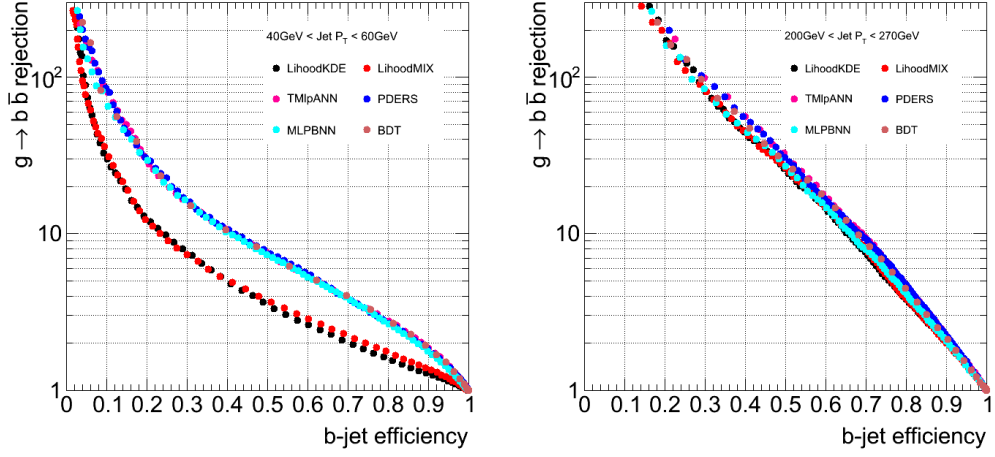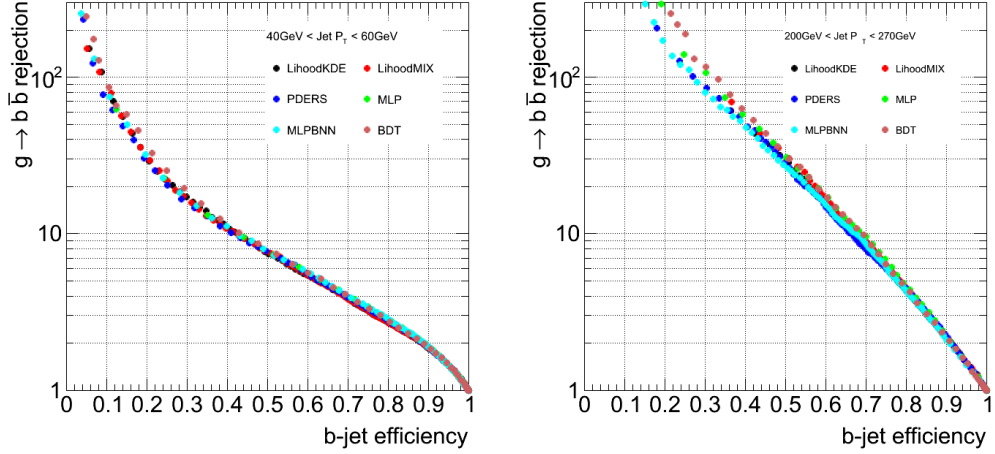


Figure 2.4: Distribution of the MVA discriminant performance for training in bins of jet $p_T$, in single and merged $b$-jets, for low and high jet $p_T$.
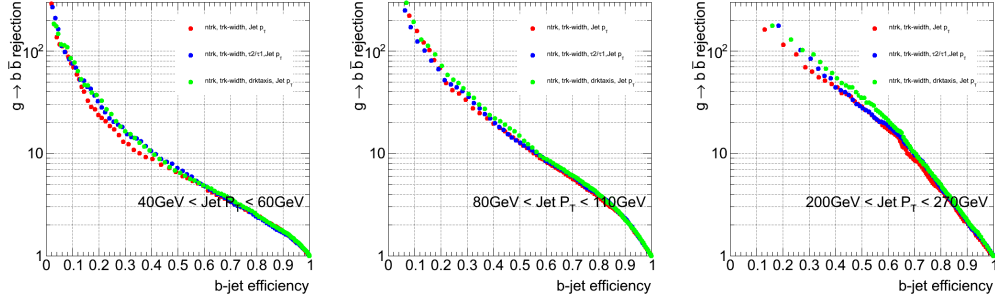
Figure 2.5: Distribution of the MVA discriminant performance for three sets of input variables, in single and merged $b$-jets, for low, medium and high jet $p_T$.

## 2.2 The input variables

Different groups of input variables were tested. Figure **??** shows the performance for three sets of variables for MVA classifier.

## 2.3 $g \to b\bar{b}$ likelihood training and performance

A discriminant between single $b$-jets and merged $b$-jets was built by training a simple likelihood estimator in the context of the Toolkit for Multivariate Data Analysis, TMVA [19].

A sub-set of the dijet Monte Carlo sample was used for training. After the event and jet selections were performed, the $b$-tagged jets with $|\eta| < 2.1$ were classified as signal (single $b$-jets) or background (merged $b$). The likelihood training was done in bins of calorimeter jet $p_T$. Signal and background jets were not weighted by the dijet samples cross-sections to allow the contribution of subleading lower $p_T$ jets from high $p_T$ events, and thus increase the statistics of merged jets in the low $p_T$ bins. For the evaluation of the method the same procedure was followed.
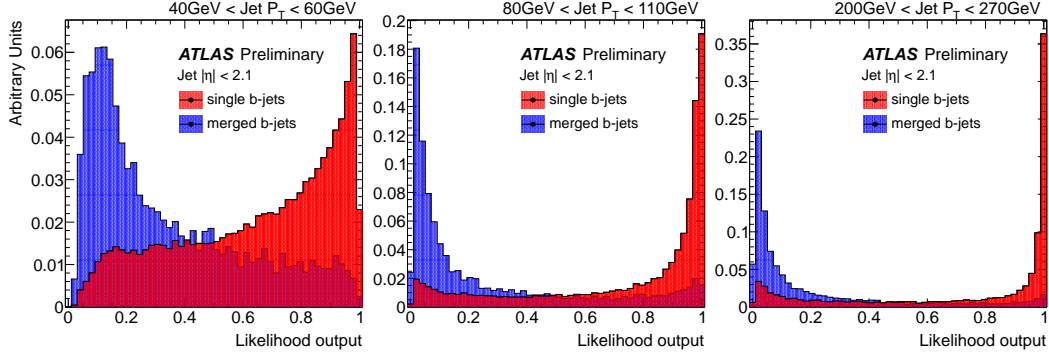
Figure 2.6: Distribution of the $g \to b\bar{b}$ likelihood output for single and merged $b$-jets for low, medium and high $p_T$ jets.

Several combinations of the tracking and jet shape variables studied in the previous section were tested as input variables. We found that the following three offer the best performance:

1. Jet track multiplicity
2. Track-jet width
3. $\Delta R$ between the axes of 2 $k_t$ subjets within the jet

A requirement of at least two matching tracks was imposed to all $b$-tagged jets in order to build the third variable listed. This cut was applied in both training and testing samples.

The distribution of the likelihood output for single and merged $b$-jets is shown in Fig. 2.6 for low, medium and high transverse momentum jets.

The performance of the $g \to b\bar{b}$ tagger in the simulation can be displayed in a plot of rejection $(1/\epsilon_{bkg})$ of merged $b$-jets as a function of single $b$-jet efficiency, where $\epsilon_{bkg}$ is the probability that a $b\bar{b}$-jet passes the tagger. This is shown in Fig. 2.7 for the eight bins of jet $p_T$ mentioned in section 1.2. The performance improves with $p_T$:

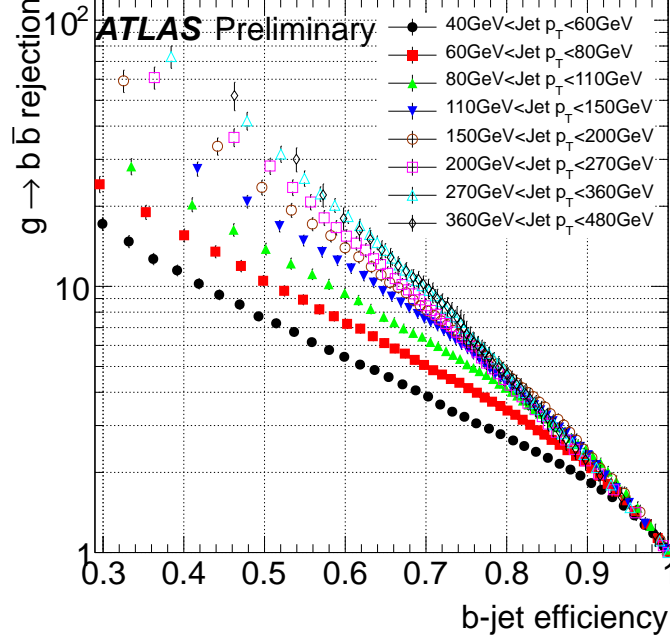- $p_T > 40$ GeV: rejection above 8 at 50% eff.

Figure 2.7: Rejection of $g \to b\bar{b}$ merged $b$-jets as a function of $b$-jet efficiency for dijet events in 8 jet $p_T$ bins.

- $p_T > 60$ GeV: rejection above 10 at 50% eff.
- $p_T > 200$ GeV: rejection above 30 at 50% eff.

The likelihood was trained with jets that had been first tagged by the MV1 algorithm. In order to use the $g \to b\bar{b}$ classifier for jets tagged by another tagger a new training is required.

The rejection of merged jets attained as a function of $p_T$ for the 50% and 60% efficiency working points are summarized in Table 2.1, together with their relative statistical error. These are propagated from the Poisson fluctuations of the number of events in the merged and single $b\bar{b}$ distributions. The error is slightly lower for the 60% efficiency working point because a higher efficiency allows for a greater number of Monte Carlo events to measure

the performance.

| Jet $p_T$ | single $b$-jet efficiency 50% | | single $b$-jet efficiency 60% | |
|---|---|---|---|---|
| (GeV ) | Rejection | stat.err. | Rejection | stat.err. |
| 40 - 60 | 8 | 4% | 5 | 3% |
| 60 - 80 | 10 | 4% | 7 | 4% |
| 80 - 110 | 14 | 5% | 9 | 4% |
| 110 - 150 | 19 | 5% | 12 | 4% |
| 150 - 200 | 23 | 5% | 14 | 5% |
| 200 - 270 | 30 | 7% | 16 | 6% |
| 270 - 360 | 36 | 7% | 19 | 6% |
| 360 - 480 | 41 | 8% | 18 | 8% |

Table 2.1: The merged $b$-jet rejection for the 50% and 60% efficiency working points in bins of $p_T$.

## 2.4   Systematic uncertainties

The development, training and performance determination of the tagger is based on simulated events. Although the agreement between simulation and data explored in section **??** is a necessary validation condition, it is also important to investigate how the tagger performance depends on systematics relevant in the data. In particular we have considered:

- presence of additional interactions (pile-up)
- uncertainty in the $b$-jet tagging efficiency
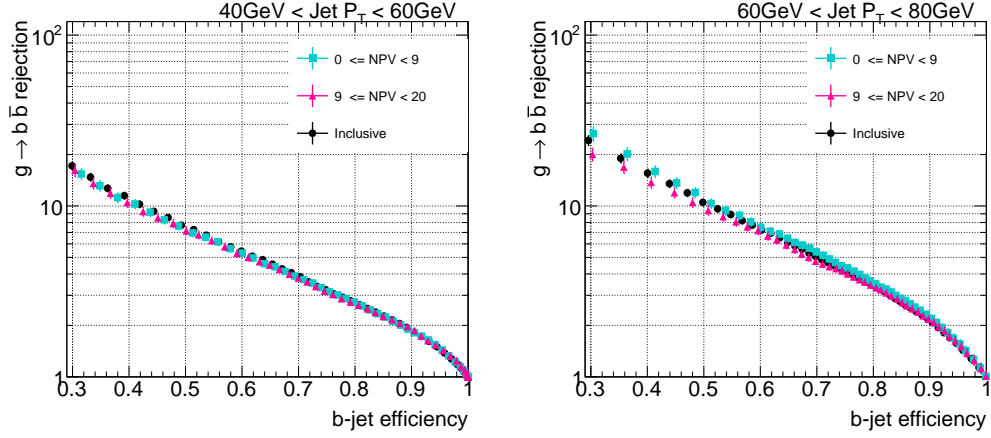- uncertainty in the track reconstruction efficiency

Figure 2.8: Rejection of $g \to b\bar{b}$ merged b-jets as a function of $b$-jet efficiency in bins of $N_{\mathrm{vtx}}$ for two low jet $p_T$ bins.

- uncertainty in the track transverse momentum resolution
- uncertainty in the jet transverse momentum resolution

*I. Pile-up*

The size of this effect was studied by comparing the performance of the likelihood discriminant with b-jets in events with small (1-9) and large (9-20) number of primary vertices. The comparison of the performance in these two sub-samples can be seen in Fig. 2.8. As expected from the use of tracking (as opposed to calorimeter) variables no significant dependence with pile-up is observed within statistics. Of the 16 determinations (2 working points with 8 $p_T$ bins each) of performance differences between high and low number of primary vertices events, it is observed that 6 of them are positive and 10 negative, with a global mean of 0.3%. We conclude that the effect is negligible compared to other source of uncertainties.

*II. b-tagging efficiency*

The performance of heavy-flavor tagging in Monte Carlo events is calibrated

to experimental data by means of the scale factors (SFs) measured by the $b$-tagging group. Such a measurement carries a systematic uncertainty, and in order to estimate its effect a conservative approach is followed: the SFs are varied in all the $p_T$ bins simultaneously by one standard deviation both in the up and down directions. The result of this procedure for the distribution of two of the tracking variables used in our discriminant is illustrated in Fig. 2.9.

The effect of the $b$-tagging calibration uncertainty on the likelihood peformance is $< 1\%$, negligible with respect to the statistical uncertainty as it can be seen in Fig. 2.10. This was indeed expected. The scale factors depend on the true flavor of the jet and on its $p_T$, but these are basically constant in the performance determination, which is based on single flavor (true $b$-) jets classified in $p_T$-bins.

*III. Track reconstruction efficiency*

This uncertainty arises from the limit in the understanding of the material layout of the Inner Detector. To test its impact a fraction of tracks determined from the track efficiency uncertainty was randomly removed following the method in Ref. [20].

The tracking efficiency systematics are given in bins of track $\eta$. For tracks with $p_T^{\text{track}} > 500$ MeV the uncertainties are independent of $p_T$: 2% for $|\eta^{\text{track}}| < 1.3$, 3% for $1.3 < |\eta^{\text{track}}| < 1.9$, 4% for $1.9 < |\eta^{\text{track}}| < 2.1$, 4% for $2.1 < |\eta^{\text{track}}| < 2.3$ and 7% for $2.3 < |\eta^{\text{track}}| < 2.5$ [21]. All numbers are relative to the corresponding tracking efficiencies.

The tracking variables were re-calculated and the performance of the nominal likelihood was evaluated in the new sample with worse tracking efficiency. The rejection-efficiency plots, shown in Fig. 2.11, show a small degradation
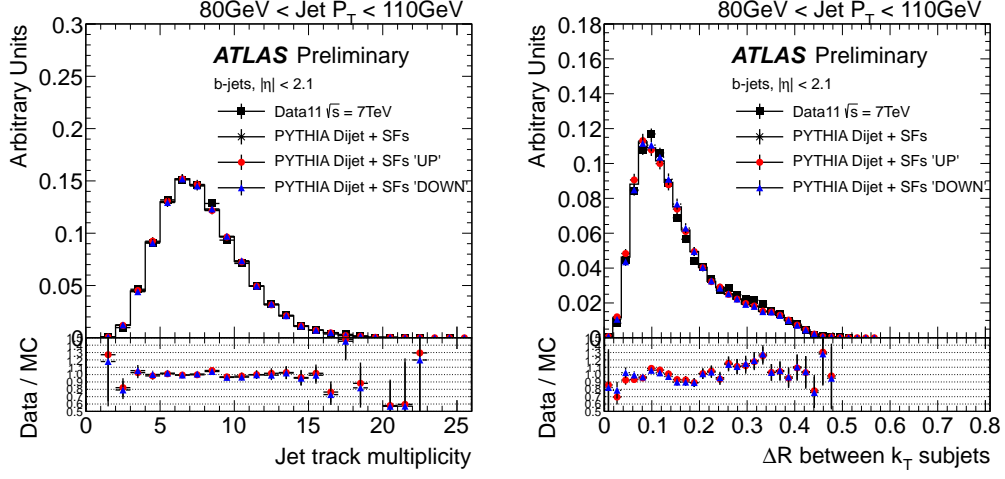
36

Figure 2.9: The effect of a variation in the *b*-tagging Scale Factors on the tracking variables distributions. Scale Factors were varied up (down) by 1-sigma to evaluate the systematic uncertainty from this source. The ratio data over MC is shown for MC PYTHIA with SFs varied up (circles) and down (triangles).
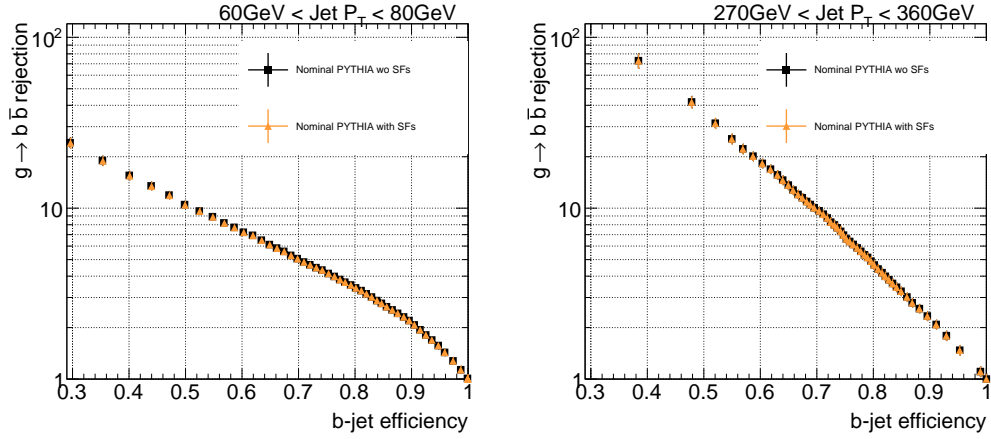


Figure 2.10: Rejection of $g \rightarrow b\bar{b}$ merged b-jets as a function of *b*-jet efficiency with and without scale factors.
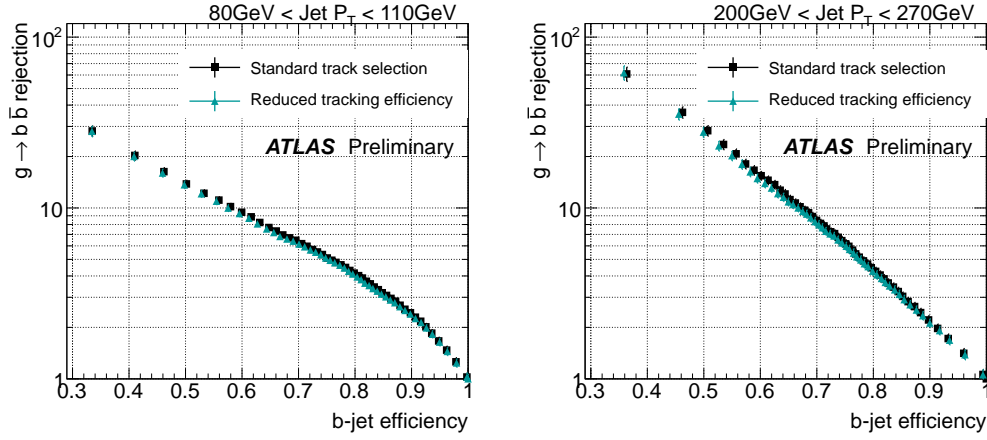
Figure 2.11: Rejection of $g \rightarrow b\bar{b}$ merged $b$-jets as a function of $b$-jet efficiency showing shift in likelihood performance caused by a reduction in the tracking efficiency .

of the performance which is comparable to the statistical uncertainty. The effect is however systematically present over all 16 $p_T$ bin/working points, without a clear $p_T$ dependence. We have thus taken the average over $p_T$, and obtained a global systematic uncertainty of 4% both for the 50% and 60% efficiency working points.

*IV. Track momentum resolution*

The knowledge of the track momentum resolution is limited by the precision both in the material description of the Inner Detector and in the mapping of the magnetic field. Its uncertainty propagates to the kinematic variables used in the $g \rightarrow b\bar{b}$ tagger. In order to study this effect, track momenta are over-smeared according to the measured resolution uncertainties before computing the rejection. The actual smearing is done in $1/p_T$, with an upper bound to the resolution uncertainty given by $\sigma(1/p_T)=0.02/p_T$ [22]. The effect is found to be negligible.

38

*V. Jet transverse momentum resolution*

The jet momentum resolution was measured for 2011 data and found to be in agreement with the predictions from the PYTHIA8-based simulation [23]. The precision of this measurement, determined in $p_T$ and $\eta$ bins,is typically 10%. The systematic uncertainty due to the calorimeter jet $p_T$ resolution was estimated by over-smearing the jet 4-momentum in the simulated data, without changing jet $\eta$ or $\phi$ angles. The performance is found to globally decrease by 6%, without a particular $p_T$ dependence.

The different contributions to the systematic uncertainty on the $g \to b\bar{b}$ rejection are summarized in Table 2.2.

| Systematic source | Uncertainty |
|---|---|
| pile-up | neglible |
| $b$-tagging efficiency | neglible |
| track reconstruction efficiency | 4% |
| track $p_T$ resolution | neglible |
| jet $p_T$ resolution | 6% |

Table 2.2: Systematic uncertainities in the merged $b$-jet rejection (common to both the 50% and the 60% efficiency working points).

## 2.5 Isolation studies

Although the tagger was derived with isolated jets it can also be applied to non-isolated jets. Studies were performed to evaluate the likelihood rejection in $b$-jets with close-by jet with $p_T$ between 7 GeV at electromagnetic scale
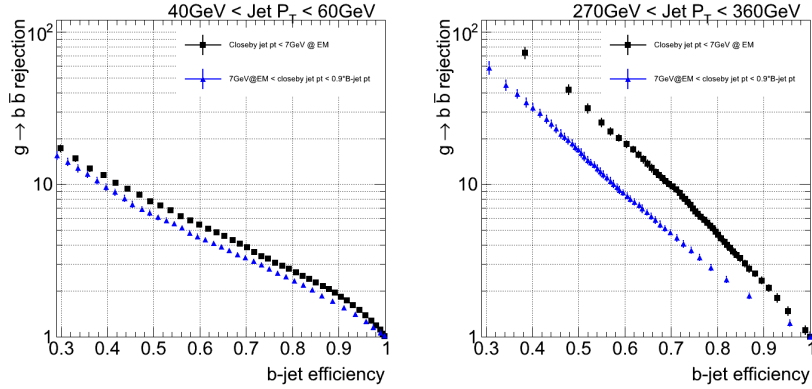
39

Figure 2.12: Rejection of $g \to b\bar{b}$ merged b-jets as a function of $b$-jet efficiency for two different isolation cuts.

scale and 90% of the $b$-jet $p_T$. The results can be seen in Fig. 2.12. The presence of close-by jets with a susbtancial fraction of the $b$-jet pt worsens the performance in more than 50% at very high $p_T$.

## 2.6    Other Monte Carlo generators

The development, training and performance determination of the tagger has been done using Monte Carlo events generated with the PYTHIA8 event simulator, interfaced to the GEANT4 based simulation of the ATLAS detector. An immediate question is what the performance would be if studied with a different simulation. In this section we investigate this question for the Perugia tune of PYTHIA8 and the HERWIG++ event generators.

Fig. 2.13 shows a comparison of the likelihood rejection, at the 50% efficiency working point, between nominal PYTHIA and the alternative simulations as a function of the jet $p_T$ . The larger errors are due to the reduced statistics available, which are even lower for the Perugia case than for HERWIG.
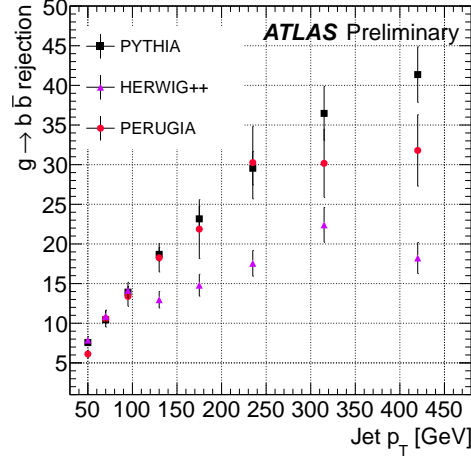
Figure 2.13: Rejection of $g \to b\bar{b}$ merged $b$-jets as a function of jet $p_T$ for different Monte Carlo generators, at the 50% efficiency working point.

The performance in HERWIG shows a systematic trend, with agreement at low $p_T$ and increasingly poor performances compared to PYTHIA as $p_T$ grows. For the Perugia tune, on the other hand, there is no definite behavior, with the performance fluctuating above or below the nominal simulation for different $p_T$ bins consistently with the statistical uncertainties.

The reason for the systematic difference observed between the performances of PYTHIA and HERWIG can be traced to the extent with which jets are accurately modelled. Fig. 2.14 compares the measured jet track multiplicity distributions in $b$-tagged jets and the prediction from both simulations, for low and high $p_T$ jets. It is observed that indeed HERWIG++ does not correctly reproduce the data, particularly at high $p_T$. The level of agreement is found to be better for track-jet width and the $\Delta R$ between the axes of the two $k_t$ subjets in the jet, the two other variables used for discrimination.
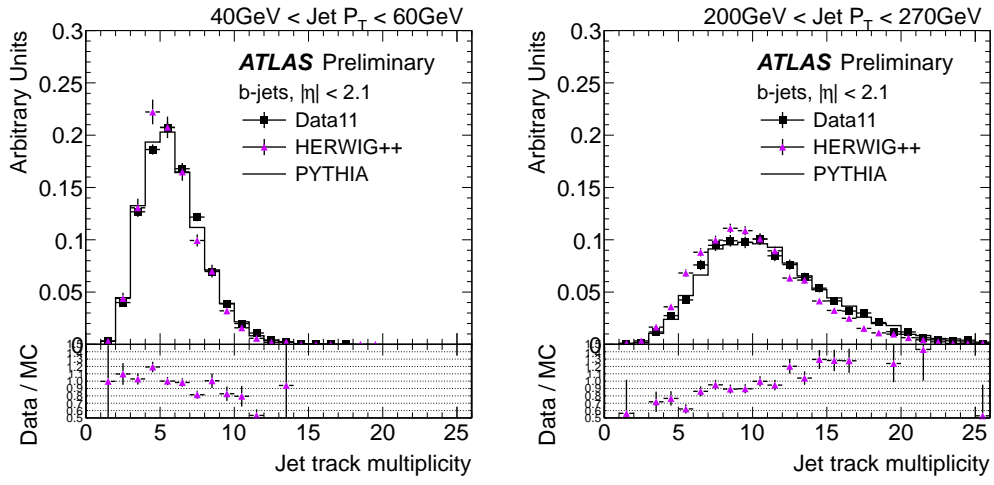
Figure 2.14: Distribution of the jet track multiplicity in 2 different jet $p_T$ bins, for experimental data collected during 2011 (solid black points) and HERWIG++ events (solid violet triangules). The ratio data over HERWIG++ simulation is shown at the bottom of the plot. PYTHIA distribution is also shown for reference.

# Chapter 3

# ??Fraction of gluon-splitting jets in data??

## 3.1 ??Template fits??

# Bibliography

[1] Torbjorn Sjostrand, Stephen Mrenna, and Peter Skands. A brief introduction to pythia 8.1. *Comput. Phys. Commun.*, 178:852, 2008.

[2] R. Corke and T. Sjöstrand. Improved parton showers at large transverse momenta. *European Physical Journal C*, 69:1, 2010.

[3] T. Sjöstrand and P. Z. Skands. Transverse-momentum-ordered showers and interleaved multiple interactions. *European Physical Journal C*, 39:129, 2005.

[4] B. Andersson et al. Parton fragmentation and string dynamics. *Phys. Rep.*, 97:31, 1983.

[5] Atlas tunes of pythia 6 and pythia 8 for mc11. Technical Report ATL-PHYS-PUB-2011-009, CERN, Geneva, Jul 2011.

[6] S. Agostinelli et al. Geant4 a simulation toolkit. *Nucl. Inst. Meth. Section A*, 506(3):250 – 303, 2003.

[7] ATLAS Collaboration. The ATLAS Simulation Infrastructure. *Eur.Phys.J.C*, 70:051, 2010.

[8] G. Corcella, I.G. Knowles, G. Marchesini, S. Moretti, K. Odagiri, et al. HERWIG 6: An Event generator for hadron emission reactions with in-

terfering gluons (including supersymmetric processes). *JHEP*, 0101:010, 2001.

[9] Peter Z. Skands. The Perugia Tunes. 2009.

[10] Manuel Bahr, Stefan Gieseke, and Michael H. Seymour. Simulation of multiple partonic interactions in Herwig++. *JHEP*, 0807:076, 2008.

[11] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. The anti-$k_t$ jet clustering algorithm. *JHEP*, 04:063, 2008.

[12] W Lampl et al. Calorimeter Clustering Algorithms: Description and Performance. (ATL-LARG-PUB-2008-002. ATL-COM-LARG-2008-003), Apr 2008.

[13] ATLAS Collaboration. Selection of jets produced in proton-proton collisions with the ATLAS detector using 2011 data. *ATLAS-CONF-2012-020*, 2012.

[14] ATLAS Collaboration. Jet energy scale and its systematic uncertainty for jets produced in proton-proton collisions at $\sqrt{s} = 7$ TeV and measured with the ATLAS detector. *ATLAS-CONF-2010-056*, 2010.

[15] G. Aad et al. Expected Performance of the ATLAS Experiment - Detector, Trigger and Physics. 2009.

[16] M Cacciari and G.P. Salam. Dispelling the $N^3$ myth for the $k_t$ jet-finder. *Phys. Lett. B}, 661:057, 2006.*

[17] S. Catani, Y.L. Dokshitzer, H. Seymour, and B.R. Webber. Longitudinally invariant K(t) clustering algorithms for hadron hadron collisions. Nucl. Phys., B406:187, 1993.

[18] Jesse Thaler and Ken Van Tilburg. Identifying Boosted Objects with N-subjettiness. JHEP, 1103:015:026, 2011.

[19] Andreas Hoecker, Peter Speckmayer, Joerg Stelzer, Jan Therhaag, Eckhard von Toerne, and Helge Voss. TMVA: Toolkit for Multivariate Data Analysis. PoS, ACAT:040, 2007.

[20] R Alon et al. Backup Note for Measurement of Jet Mass and Substructure in QCD with the ATLAS Experiment. ATL-COM-PHYS-2011-401, 2011.

[21] G. Aad et al. Charged-particle multiplicities in pp interactions measured with the ATLAS detector at the LHC. New J.Phys., 13:053033, 2011.

[22] ATLAS Collaboration. Estimating Track Momentum Resolution in Minimum Bias Events using Simulation and $K_s$ in $\sqrt{s}$ = 900 GeV collision data. ATLAS-CONF-2010-009, 2010.

[23] G. Romeo, A. Schwartzman, R. Piegaia, T. Carli, and R. Teuscher. Jet energy resolution from in-situ techniques with the atlas detector using proton-proton collisions at a center of mass energy $\sqrt{s}$ = 7 tev. ATL-COM-PHYS-2011-240, 2011.