

CMS Draft Analysis Note

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

2013/04/10
Head Id: 180541
Archive Id: 180539:180543M
Archive Date: 2013/04/10
Archive Tag: trunk

New gluon/quark jet discriminator

A.C. Marini, F. Pandolfi¹, S. Bansal, T. Cornelis, N. van Remortel³, P. Azzurri⁴, and D. del Re⁵

¹ Institute for Particle Physics, ETH Zürich, CH-8093 Zürich, Switzerland

² European Organization for Nuclear Research, CERN

³ University of Antwerp, Antwerpen, Belgium

⁴ INFN Sezione di Pisa, Italy

⁵ University and INFN Rome, ???, Italy

Abstract

A new discriminator has been developed to distinguish jets originating from quarks and gluons, that extends the separation capabilities to jets in the forward region, up to $|\eta| = 5$, and increases the overall performances, in particular for jets of low transverse momentum, down to $p_T = 20\text{ GeV}$. Discriminating observables have been optimized and chosen based on performances established with Monte Carlo simulated QCD events. Two alternative multivariate quark-gluon jet taggers are presented with slightly different performances. Comparisons of 2012 dijet and Z plus jets data and Monte Carlo simulations for the input discriminating variables and the output tagger distributions have been performed. Finally a recipe to evaluate systematic uncertainties related to the use of either tagger is given, based on the observed data versus Monte Carlo differences.

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

PDFAuthor: CMS Collaboration

PDFTitle: New gluon/quark jet discriminator

PDFSubject: CMS

PDFKeywords: CMS, physics, jet, quark, gluon, QG, discrimination, tagger

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

1 Contents

2	1	Introduction	2
3	2	Motivations	2
4	3	Object Definition	3
5	4	Improving the Previous Tagger	5
6	4.1	The Considered Variables	5
7	4.2	The Choice of the Variables	9
8	4.3	The New Taggers	13
9	5	Pythia6/Herwig++ Comparisons	24
10	6	Validation on data	30
11	6.1	Validation on Z+jet Events	30
12	6.2	Validation on Dijet Events	47
13	7	Systematics evaluation	59
14	7.1	Usage: data vs data	59
15	7.2	Usage: data vs mc	59
16	8	Conclusions	64
17	A	Bjets	66

DRAFT

18 1 Introduction

19 It is known, both from first principles and a large collection of experimental measurements [1–
 20 4], that hadronic jets initiated from gluons exhibit differences with respect to jets from light
 21 flavor quarks. In jet fragmentation, final state hadrons are produced in a cascade process of sub-
 22 sequent bremsstrahlungs, followed by hadronisation setting in at a scale of $\Lambda_{QCD} \sim 200$ MeV.
 23 Bremsstrahlung is directly proportional to the coupling of the radiated gluon to the radiator,
 24 in this case a gluon or light quark. The probabilities are therefore directly proportional to the
 25 color factor $C_A = 3$ for gluons and $C_F = 4/3$ for quarks. Neglecting gluon splitting to quark-
 26 antiquark pairs, the ratio of the number of gluons radiated from a gluon to that from quarks
 27 is therefore expected the approach 9/4 at leading order and for infinite phase-space. Correc-
 28 tions to this ratio can can be calculated analytically up to NLO with the inclusion of energy
 29 (but not momentum) conservation. Numerical approaches exist as well, which are typically
 30 implemented in Monte Carlo generators, and have the advantage of taking the phase space
 31 limits more accurately into account, but by introducing ambiguities in the scale at which the
 32 evolution of the gluon multiplicities radiated from a parton are being evaluated. Regardless of
 33 the approach, a larger number of softer particles are expected to be produced in a gluon jet and
 34 the gluon jet is expected to be wider than a quark jet of the same energy.

35 Many experimental measurements of these fundamental properties have been made. Albeit
 36 with ambiguities in the gluon-jet separation, the measurements are in qualitative agreement
 37 with the analytical expectations and are described quantitatively by most recent Monte Carlo
 38 models and have been used to extract the ratio C_A/C_F which is compatible with the SU(3)
 39 expectation. The most accurate experimental results originate from LEP, supplemented with
 40 more data from hadron colliders such as the Tevatron. The following main experimental obser-
 41 vations were made:

- 42 • The charge multiplicity is higher in gluon jets than in light quark jets, independent
 of the particle species;
- 44 • the fragmentation function of gluon jets is considerably softer than that of a quark
 jet, which can be explained by the higher multiplicity of soft gluons radiated off a
 gluon and double string formation necessary before hadronization in case of $g \rightarrow q\bar{q}$
 splitting;
- 48 • gluon jets are less collimated than quark jets, an observation that can be quantified
 in the jet broadness or the major and minor axis of the jet cone.

50 In this study we will exploit these differences to construct a probability tagger capable of dis-
 51 criminating jets initiated by light quark partons from those initiated by gluons. The use of
 52 Particle Flow jet reconstruction allows us to access particle-level information during jet recon-
 53 struction, exploiting the full granularity of the CMS detector. We shall show how this advan-
 54 tage allows us to define a broad range of discriminating variables.

55 2 Motivations

56 Some analyses rely on the identification of an exclusive final state that includes a fixed num-
 57 ber of hadronic jets that generally originate from light or heavy quarks. Some examples are
 58 hadronic W, Z and top decays, the production of vector bosons or Higgs via vector boson fu-
 59 sion, associated vector boson and Higgs production and cascade decays of squarks. In many
 60 cases these hadronic final states suffer from overwhelming backgrounds from multi-jet QCD
 61 production or from electroweak backgrounds with hard initial or final state gluon radiation.

62 The tagging of quark jets is also useful in the mass reconstruction of hadronically decaying
 63 objects, where the resolution is generally degraded due to combinatorial backgrounds.

64 In addition, there is also the increasing amount of pile-up that can result in the reconstruction
 65 of artificially created jets by clustering energy deposits and/or tracks originating from different
 66 proton-proton interactions in teh same bunch crossing. These fake jets will have characteristics
 67 that are different from showering gluons or quarks and can therefore also be discriminated.

68 Developing a general tool that suits all purposes is difficult, as many analyses use different
 69 kinds of jet objects, work in different kinematic regimes and/or use loose or stringent detector
 70 acceptances and quality requirements.

71 The need for this tool arose from an analysis dedicated to the measurement of a Higgs boson,
 72 produced via the vector-boson-fusion process and subsequently decaying to b-quarks. This
 73 process is characterized by a four jet final state, of which two are b-tagged and two forward
 74 going light quark jets with a large rapidity separation. The multi-jet QCD backgrounds for this
 75 process are overwhelming, resulting in signal-to-noise ratios of 10^{-5} , even after using dedicated
 76 triggers. A gluon/quark tagging tool will undoubtedly help in suppressing these types of
 77 background. When developing the tool, we will try to cover the most generally used jet objects
 78 and cover a large pseudorapidity and transverse momentum range. The approach can also
 79 serve as a guideline for those who want to develop their own discriminators when using more
 80 specialized jet definitions or quality cuts.

81 **3 Object Definition**

82 The CMS software framework supports a susbtantial variety of jet objects, some dedicated to
 83 QCD physics, some to electroweak analyses and searches. In addition, there are several levels
 84 of jet energy scale corrections and procedures to subtract pile-up effects. In this study we use
 85 jets clustered with the Anti- k_t algorithm with a distance parameter $R = 0.5$ using charged and
 86 neutral particles that are reconstructed with the Particle Flow algorithm. Residual jet energy
 87 corrections are applied in order to make the response as linear as possible with unit slope. In
 88 this study we have used the 53X(V7) jet energy corrections. The same clustering algorithm is
 89 applied to stable generator particles in simulated events, thus defining generator jets.

90 All jets are further subjected to a set of basic quality requirements:

- 91 • they should originate from an event with at least one good quality vertex with a min-
 92 imum of four degrees of freedom and a maximal distance to the nominal interaction
 93 point along the beam axis of 24 cm.
- 94 • a tight jet identification, according to the criteria in [5].

95 The acceptance in pseudorapidity is kept as large as possible, but the jets will be categorized
 96 in three distinct acceptance regions: the central region corresponding to $|\eta| < 2$, the transition
 97 region corresponding to $2 < |\eta| < 3$, and the forward region corresponding to $3 < |\eta| < 5$. The
 98 jet transverse momenta are required to exceed 20 GeV/c and can be as large as kinematically
 99 allowed. The training of the tool will split up the sample in several narrow jet p_T bins, and will
 100 be conducted in the central region ($|\eta| < 2$) and the forward region ($3 < |\eta| < 5$). The results
 101 will then be extended to cover the transition region, so that the taggers defined in the central
 102 region will be commissioned up to $|\eta| = 2.5$, and the forward tagger will be used also in the
 103 $2.5 < |\eta| < 3$ pseudorapidity region.

104 The studies carried on the Monte Carlo simulation will assign a flavor to each reconstructed

105 jet. In QCD, flavor is not an infrared-safe observable, and is therefore not well defined if not
 106 at tree level. We define the jet flavor with an empirical approach, by accessing the generator
 107 information (partons and generator jets) in the following way:

- 108 • if no generator jet is found with a $\Delta R = 0.3$ cone around the reconstructed jet direc-
 109 tion, the jet is considered to be a pile-up (PU) jet;
- 110 • if a generator jet is present, the reconstructed jet is matched to the closest (in terms
 111 of ΔR) parton in the event (partons are particles in the event with PYTHIA status
 112 = 3 and PDG identification number compatible with either a quark or a gluon). If a
 113 parton is found within $\Delta R = 0.3$ of the reconstructed jet direction, the flavor of the
 114 parton is assigned to the jet; if not, the jet is considered ‘undefined’.

115 In the above (as in the following), ΔR represents the distance in the $\eta - \phi$ plane, and is defined,
 116 for pairs of objects whose rapidities (azimuths) differ by $\Delta\eta$ ($\Delta\phi$), as $\Delta R = \sqrt{\Delta\eta^2 + \Delta\phi^2}$.

117 As is shown in more detail in Appendix A, jets initiated by b quarks present hadronization
 118 properties which are closer to gluon-initiated jets than to light quark jets. On the other hand, b-
 119 jets can be efficiently be identified with the use of b-tagging algorithms, such as the Combined
 120 Secondary Vertex (CSV) tagger, which will be used in this note, therefore this tagger will be
 121 specialized in the identification of light quark jets. In order to do so, all PFJets considered
 122 in this note will be required to be anti-b-tagged, i.e. if a jet is found to have a positive loose
 123 Combined Secondary Vertex (CSV) tag, it is discarded.

In addition to being complementary to b-tagging, the quark-gluon discriminator is also meant to be complementary to PU-jet identification. For this reason we require, for jets reconstructed in the tracker-covered region, that they are not identified as PU jets by using the standard requirement on the β^* variable, which is defined as the energy fraction of tracks of the jet which are not compatible with the primary vertex of the event. We therefore require for all jets considered in this study to satisfy:

$$\beta^* < 0.2 * \ln(N_{PV} - 0.67)$$

124 where N_{PV} is the number of reconstructed primary vertexes in the event. Standard PU identi-
 125 fication also makes use of another variable, the RMS of the spread of the jet candidates in the
 126 $\eta - \phi$ plane. But as this variable is also strongly influenced by hadronization effects, and is
 127 discriminating between quarks and gluons, we do not apply any requirement on this variable,
 128 in order to avoid to introduce a bias in the sample.

129 The amount of PU (and underlying event) activity in a given event is measured by the diffuse
 130 particle flow energy density (ρ), computed with the median of the transverse momenta of jets
 131 clustered with the k_T clustering algorithm with radius parameter $R = 0.6$, after injecting in the
 132 event a very large number of infinitely soft particles (ghosts). For the studies presented in this
 133 note, two definitions of the ρ variable have been employed:

- 134 • the ‘standard’ ρ , which is computed with all PFCandidates in the event (up to $|\eta| =$
 135 5);
- 136 • an energy density variable computed in the same manner, but restricting the median
 137 computation to the candidates reconstructed in the tracker acceptance ($|\eta| < 2.5$).
 138 This variable will be from here onwards called ρ_{iso} .

139 The two above measurements of the amount of PU in an event are equivalent, and the two
 140 variables have very high degrees of correlation. The use of two different variables is mainly
 141 due to historical reasons.

142 4 Improving the Previous Tagger

143 Previous studies, conducted in 2011, have already led to the definition of a quark-gluon dis-
 144 criminator [6], which has been used in published results, such as [7]. The tagger was a simple
 145 multiplicative likelihood discriminator, based on three variables:

- 146 • **Charged Multiplicity:** the number of reconstructed charged hadron candidates in
 147 the PFJet;
- 148 • **Neutral Multiplicity:** the number of reconstructed photon and neutral hadron can-
 149 didates in the PFJet;
- **$p_T D$:** defined as:

$$p_T D = \frac{\sqrt{\sum p_T^2}}{\sum p_T}$$

150 where the sums are extended over all candidates within the PFJet.

151 This tagger has been developed focusing on the discrimination of hard jets ($p_T \gtrsim 100\text{GeV}$), and
 152 was defined only for jets reconstructed within the tracker acceptance. Of the three used vari-
 153 ables, the discrimination power was almost exclusively left to $p_T D$ up to transverse momenta
 154 of 100-200 GeV, and from then on the role of the multiplicities became increasingly important.

155 The studies reported in this note describe how we have improved the performance of quark-
 156 gluon discrimination at CMS. The improvements to the discriminator have been lead by two
 157 main directives:

- 158 • the extension of the discriminator to the forward rapidity region;
- 159 • the optimization of the discrimination performance for soft jets ($p_T < 50\text{ GeV}$).

160 We have furthermore decided to define two distinct taggers: one will be a simple likelihood
 161 product, similar to what was done in 2011; the other will be a multivariate discriminant (a
 162 multi-layer perceptron, or MLP), implemented within the TMVA [8] framework.

163 The re-optimization of the tagger performance starts from an improved set of discriminating
 164 variables. We will now show how the choice of the new variable set was made.

165 4.1 The Considered Variables

166 A large number of variables can be considered in order to discriminate between quark and
 167 gluon hadronization. As we have seen, we are looking for variables which are able to distin-
 168 guish quark and gluon hadronization based on the number, energy and spread of the stable
 169 particles produced during such a process. We have therefore tried to define the largest set of
 170 variables which could be able of doing so. The complete list of variables we have considered is
 171 the following:

- 172 • charged multiplicity ($nChg$);
- 173 • neutral multiplicity ($nNeutral$);
- 174 • total multiplicity ($nPFCand$);
- 175 • candidate RMS of in the $\eta - \phi$ plane ($RMSCan$);
- 176 • candidate second moment minor axis in the $\eta - \phi$ plane ($axis_2$);
- 177 • candidate second moment major axis in the $\eta - \phi$ plane ($axis_1$);
- 178 • candidate asymmetry (or $pull$);

- 179 • fractional leading candidate transverse momentum (R);
 180 • candidate p_T distribution (p_TD).

181 Each variable has been computed in two varieties: the default computation uses all PFCandidates
 182 within the jet as input. An additional value of the variable is obtained when restricting
 183 the input candidate list: out of all charged constituents, only those whose tracks satisfy qual-
 184 ity control (QC) criteria are considered. These criteria require the tracks to be assigned to the
 185 hardest reconstructed primary vertex of the event, in a similar way as what is done by the
 186 charged hadron subtraction in PFNoPU. Only ‘highPurity’ tracks assigned to the hardest PV
 187 are kept, and they are further required to satisfy $|d_z/\sigma(d_z)| < 5$ and $|d_0/\sigma(d_0)| < 3$. Quality
 188 control criteria are intended to mitigate the effect of pile up in the tracker covered region of the
 189 detector.

190 One additional variation was considered, only for the multiplicity and R variables: the set of
 191 PFCandidates was restricted to the candidates with a p_T greater than 1 GeV. This requirement
 192 is also intended to limit the effects of PU, especially in the forward, tracker-less, region of the
 193 detector.

194 What follows is the complete list of variables we have considered, with their definition.

195 **Multiplicities**

196 The simplest and most studied variable that we can construct is the multiplicity, i.e. the total
 197 number of candidates reconstructed within the jet. Its properties are summarized in different
 198 articles (e.g. [9]), where it is shown that the ratio of the multiplicity from quark and gluon jets
 199 should converge to the color factor ratio $\frac{C_A}{C_F} = \frac{9}{4}$ and provide a good discrimination, especially
 200 at high p_T . We construct three multiplicity variables:

- 201 • **nPFCand**: the total PFCandidate multiplicity of the PFJet;
- 202 • **nChg**: the multiplicity of charged PFCandidates of the PFJet (charged hadrons, elec-
 203 trons, muons);
- 204 • **nNeutral**: the neutral PFCandidate multiplicity of the PFJet (photons, neutral hadrons,
 205 HF hadrons, HF EM particles).

206 Quark jets are expected to have, on average, lower values of multiplicities with respect to
 207 gluon jets, for a given transverse momentum. Figures 4.1(a,b,c) show the expected normal-
 208 ized distributions for nPFCand, nChg and nNeutral for quark (blue) and gluon (red) jets with
 209 $80 < p_T < 120$ GeV reconstructed in the central part of the detector ($|\eta| < 2$).

210 **RMS_cand**

211 Jets have a conical structure that can be projected in the $\eta - \phi$ plane. From first principles
 212 we expect gluon hadronization to produce jets which are ‘wider’ than jets induced by quark
 213 hadronization. The jet width may be quantified by looking at the second moments of the con-
 214 stituent distribution in the $\eta - \phi$ plane.

The simplest way of doing this is to treat the η and ϕ directions democratically, and define a global RMS of the candidate spread:

$$\text{RMS}_{\text{c}}\text{and} = \sqrt{\frac{\sum_i w_i \Delta R_i^2}{\sum_i w_i}}$$

215 where the sums run on all PFCandidates within the jet, the ΔR_i are distances in the $\eta - \phi$ plane

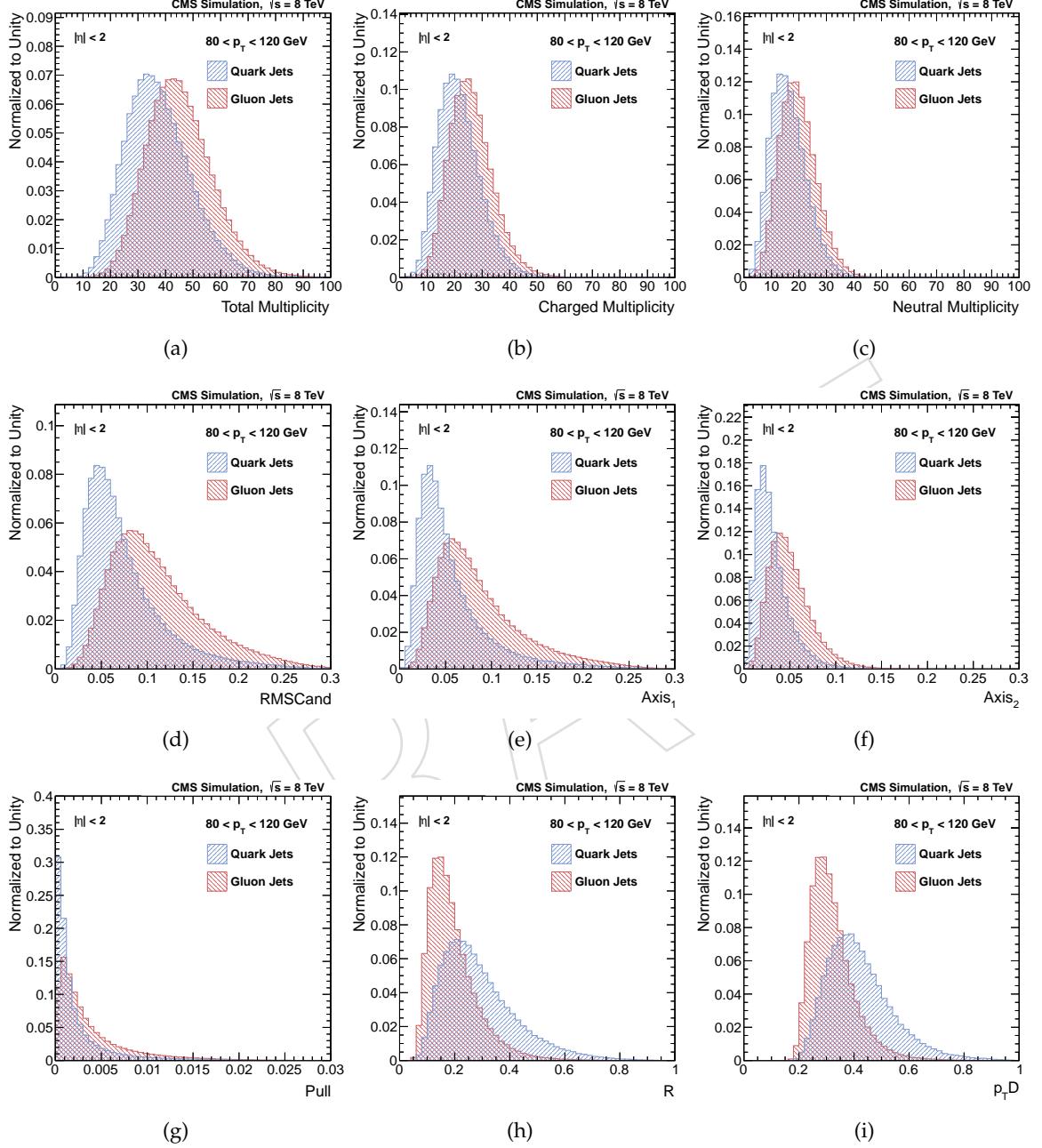


Figure 4.1: Normalized distributions of the considered discriminating variables, for quark (blue) and gluon (red) jets with $80 < p_T < 120 \text{ GeV}$ reconstructed in the central part of the detector ($|\eta| < 2$). (a): total multiplicity (nPFcand); (b): charged multiplicity (nChg); (c): neutral multiplicity (nNeutral); (d): RMScand; (e): axis₁; (f): axis₂; (g): pull; (h): R; (i): $p_T D$.

($\Delta R_i = \sqrt{\Delta\eta_i^2 + \Delta\phi_i^2}$, where $\Delta\eta_i$ and $\Delta\phi_i$ are, respectively, the rapidity and azimuthal differences between the given candidate and the mean of the candidates), and the w_i are weights. Since jets consist of particles with a power law p_t spectrum, it is natural to compute the jet axes by means of p_t based weights. It was found in [?] that the weights, w_i give an optimal separation between gluon and quark jets when they depend on the squared transverse momentum of each constituent particle: $w_i = p_{T,i}^2$. This definition of weights will be used in all variables.

Quark jets are expected to produce, on average, narrower jets and therefore lower values of RMS and with respect to gluon jets, for a given transverse momentum. Figure 4.1(d) shows the expected normalized distribution for RMS and for quark (blue) and gluon (red) jets with $80 < p_T < 120$ GeV reconstructed in the central part of the detector ($|\eta| < 2$).

Axes

A slightly more sophisticated approach is to compute, instead of the simple RMS of the candidate $\eta - \phi$ spread, the two principal components of the second moments. The shape of a jet can be approximated by an ellipse which is characterized by its two principle axes, the major and the minor axis, and the orientation of the major axis in the plane. The shape of the jet can thus be specified in terms of a 2x2 symmetric matrix, M , with the following elements:

$$M_{11} = \sum_i w_i \Delta\eta_i^2 \quad (4.1)$$

$$M_{22} = \sum_i w_i \Delta\phi_i^2 \quad (4.2)$$

$$M_{12} = M_{21} = \sum_i w_i \Delta\eta_i^2 \Delta\phi_i \quad (4.3)$$

where $\Delta\eta_i$ and $\Delta\phi_i$ are taken as the pseudorapidity and azimuthal distance between each constituent particle, i , and the jet direction.

The major and minor axes of the jet can then be computed as the eigenvalues $\lambda_{1,2}$ of the matrix M by:

$$\text{axis}_1 = (\lambda_1 / \sum_i w_i)^{1/2} \quad (4.4)$$

$$\text{axis}_2 = (\lambda_2 / \sum_i w_i)^{1/2}. \quad (4.5)$$

By defining them in this way, the two axes are statistically independent.

Quark jets are expected to produce, on average, narrower jets and therefore lower values of the axes variables with respect to gluon jets, for a given transverse momentum. Figures 4.1(e,f) show the expected normalized distribution for axis_1 and axis_2 for quark (blue) and gluon (red) jets with $80 < p_T < 120$ GeV reconstructed in the central part of the detector ($|\eta| < 2$).

Pull

Another discriminator is the jet asymmetry or pull, that expresses how much the jet direction is determined by one (or more) particularly high p_t particle that are not symmetrically distributed around the jet direction. It is computed as the vector sum of p_T^2 weighted distances of each jet constituent with respect to the jet direction (r_i):

$$\text{pull} = \left| \frac{\sum_i w_i |r_i| \vec{r}_i}{\sum_i w_i} \right| \quad (4.6)$$

242 Quark jets are expected to have, on average, larger values of the pull variable with respect to
243 gluon jets, for a given transverse momentum. Figure 4.1(g) shows the expected normalized
244 distribution for the pull variable for quark (blue) and gluon (red) jets with $80 < p_T < 120$ GeV
245 reconstructed in the central part of the detector ($|\eta| < 2$).

246 **R**

The R variable is defined as the fractional transverse momentum of the leading particle inside the jet:

$$R = \frac{\max(p_{T,i})}{\sum_i p_{T,i}}. \quad (4.7)$$

247 Quark jets are expected to hadronize more asymmetrically with respect to gluon jets, and there-
248 fore have R values closer to unity. Figure 4.1(h) shows the expected normalized distribution
249 for the R variable for quark (blue) and gluon (red) jets with $80 < p_T < 120$ GeV reconstructed
250 in the central part of the detector ($|\eta| < 2$).

251 **$p_T D$**

A slightly more complex way of assessing how asymmetrically the transverse momentum of a given jet is distributed throughout its constituents is given by the $p_T D$ variable, defined by:

$$p_T D = \frac{\sqrt{\sum_i p_{T,i}^2}}{\sum_i p_{T,i}}$$

252 From its definition, it stems that $p_T D \rightarrow 1$ for jets made of only one particle, which carries all
253 of its momentum, and $p_T D \rightarrow 0$ for a jet made of an infinite number of particles, similarly to
254 the R variable.

255 It is interesting to notice how $p_T D$ and the axes are linked together by an elegant description
256 of the candidates in terms of the introduction of momenta of the jets, and so they can be seen as
257 term of the momenta expansion as shown in [10].

258 Quark jets are expected to hadronize more asymmetrically with respect to gluon jets, and there-
259 fore have $p_T D$ values closer to unity. Figure 4.1(i) shows the expected normalized distribution
260 for the $p_T D$ variable for quark (blue) and gluon (red) jets with $80 < p_T < 120$ GeV recon-
261 structed in the central part of the detector ($|\eta| < 2$).

262 **4.2 The Choice of the Variables**

263 The choice of the set of variables on which to base the new discriminator has been done by
264 studying their discriminating performance on the simulation. The discriminating power was
265 assessed by studying their receiver operating characteristics (ROC) in terms of quark efficiency
266 and gluon rejection. These studies have been performed in a specific ρ_{iso} interval ($12 < \rho_{iso} <$
267 14 GeV), which roughly corresponds to the position of the mode of the data. No dramatic
268 difference in terms of variable performance is expected as a function of ρ_{iso} , nor any has been
269 observed.

270 The ROC plots are shown in Figures 4.2 and 4.3, respectively for jet reconstructed in the central
271 part of the detector ($|\eta| < 2$) and in the forward ($3 < |\eta| < 5$). In each figure, different
272 transverse momentum intervals are considered: for central jets $30 < p_T < 40$ GeV (top left),
273 $80 < p_T < 120$ GeV (top right) and $300 < p_T < 400$ GeV (bottom); for forward jets $20 < p_T <$
274 40 GeV (left) and $50 < p_T < 100$ GeV (right).

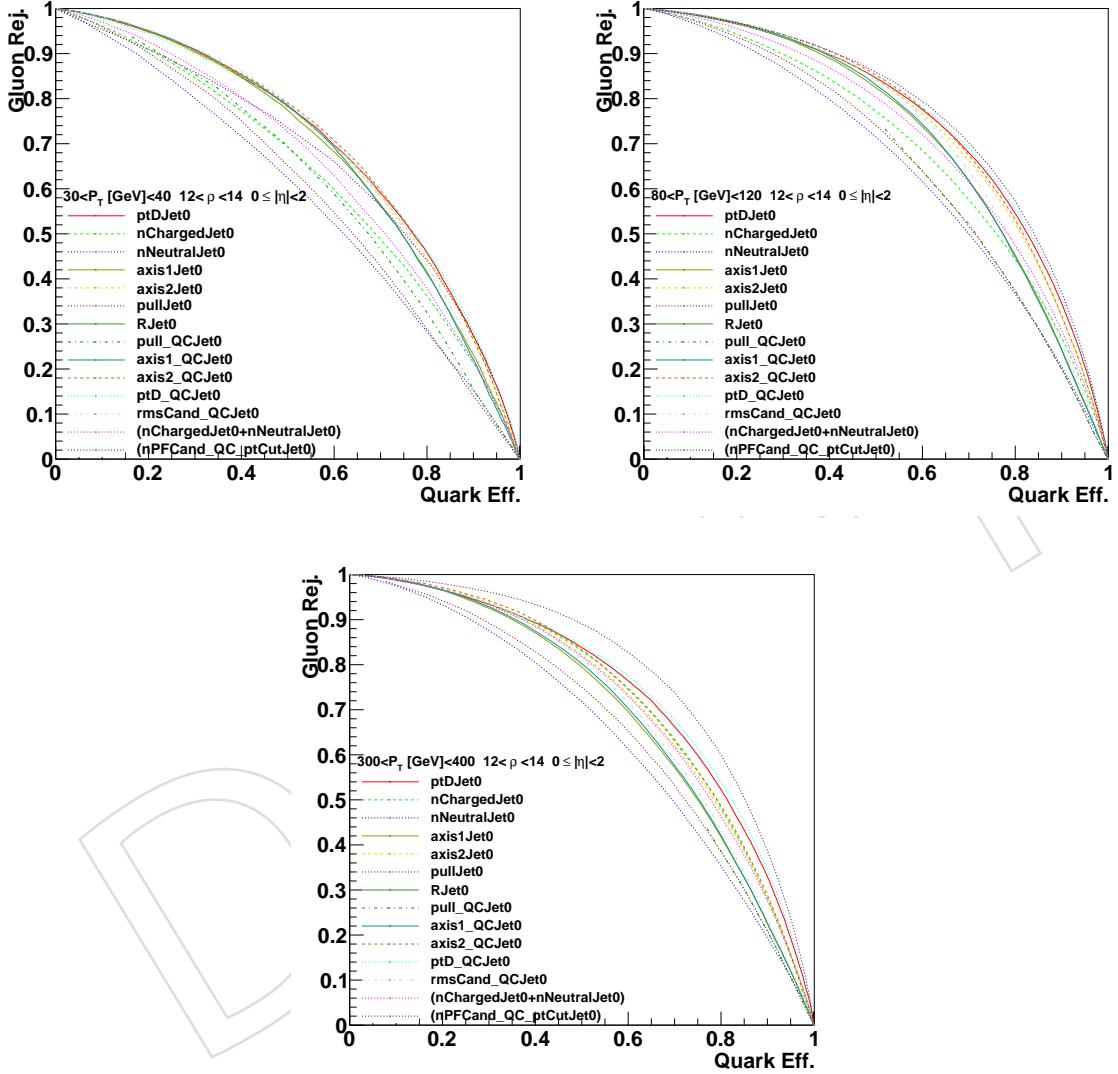


Figure 4.2: Single variable ROC curves for central jets ($|\eta| < 2$) in three transverse momentum bins: $30 < p_T < 40$ GeV (top left), $80 < p_T < 120$ GeV (top right) and $300 < p_T < 400$ GeV (bottom). The figures show the single variable performance in the quark efficiency (x axis) - gluon rejection (y axis) plane.

275 We start by analyzing the central jet ROC figures. The figures are very busy, but (with the help
276 of some zooming) some results may be assessed:

- 277 • in general, quality control (QC) variables perform (slightly) better than their stan-
278 dard counterparts;
- 279 • of the multiplicity variables, the charged multiplicity significantly outperforms the
280 neutral multiplicity, but the total multiplicity is found to have even better discrim-
281 ination; the total multiplicity with quality control on tracks and minimal 1 GeV p_T
282 requirement on the neutral (nPFcand_QC_ptCut) shows very good discrimination
283 starting from jets of about 100 GeV onwards;
- 284 • out of the ‘spread’ variables, the second moment minor axis (axis_2) has the best per-
285 formance, the second moment major axis (axis_1) has the worst, and the averaged
286 RMS (RMScand) lies approximately in between;
- 287 • the pull variable does not provide much discrimination power;
- 288 • out of the two asymmetrical hadronization variables (R and p_TD), p_TD has the best
289 performance, and actually seems to be the variable with the best discrimination
290 power across all the p_T range.

291 Most of these observations are valid also in the forward region (see Figure 4.3). Here no tracker
292 multiplicity is present, so the corresponding entries in the legend for the total multiplicity cor-
293 respond to the total neutral multiplicity. As you can see, the introduction of the minimal 1 GeV
294 transverse momentum threshold increases significantly the discrimination.

295 We therefore proceed to choose the set of variables on which to build the discriminator. The
296 choice criteria may be summarized as follows: we identify the minimal set of uncorrelated
297 variables which present the highest single-variable discrimination power. From what has been
298 observed above, it is clear that such a set is constituted by the following variables:

- 299 • total multiplicity
- 300 • axis_2
- 301 • p_TD

302 To these three variables, a fourth one can be added:

- 303 • axis_1

304 Although this variable doesn’t have very high discrimination power, its definition makes it
305 independent from axis_2 , therefore it brings additional, uncorrelated, information, which will
306 improve the discrimination power of a multivariate discriminator.

307 All these variables are taken in their quality control (QC) variation, and only the neutral can-
308 didates which have transverse momentum greater than 1 GeV are used in the computation of
309 the total multiplicity.

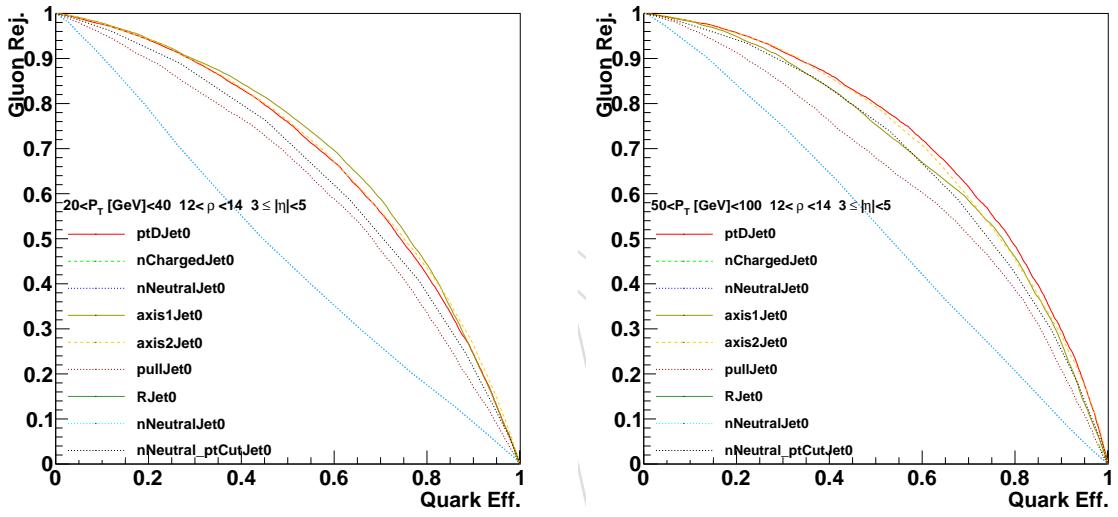


Figure 4.3: Single variable ROC curves for forward jets ($3 < |\eta| < 5$) in two transverse momentum bins: $20 < p_T < 40$ GeV (left) and $50 < p_T < 100$ GeV (right). The figures show the single variable performance in the quark efficiency (x axis) - gluon rejection (y axis) plane.

310 **4.3 The New Taggers**

311 Once the set of variables is found, we have to choose the discriminator. We define two discrim-
312 inators:

- 313 • a simple 3-variable likelihood discriminator, built on the product of nPFCand, axis₂
314 and $p_T D$;
- 315 • a multivariate 4-variable discriminator, based on nPFCand, axis₂, $p_T D$, and axis₁,
316 which uses the Multi-Layer Perceptron (MLP) Artificial Neural Network provided
317 within the TMVA framework.

318 We will now detail how these two discriminators are built, as they differ not only in the algo-
319 rithmical computation of the output, but also in the treatment of the input variables' depen-
320 dence on the jet transverse momentum and on the amount of PU activity.

321 **Likelihood Discriminant**

322 The likelihood discriminant is built as a simple product of the three variables nPFCand, axis₂
323 and $p_T D$. This approach ensures simplicity, transparency and robustness. As it explicitly ne-
324 glects the treatment of variable correlations, though, it might be suboptimal in phase space
325 regions in which these are complex, as for instance the high p_T limit.

326 The probability density functions (PDFs) on which the likelihood is based are built from jets
327 in simulated events, which have been successfully tagged as light quark (u, d, s) or gluon jets.
328 The PDFs are computed separately in two rapidity regions: a central region ($|\eta| < 2$) and a
329 forward region ($3 < |\eta| < 5$). The PDFs obtained in the central region are then used for jets up
330 to $|\eta| = 2.5$, and jets with $2.5 < |\eta| < 3$ use the forward region PDFs.

331 In order to take into account the strong dependance of the means and shapes of the variables
332 both as a function of the jet transverse momentum and the amount of PU activity in the event,
333 the PDFs are computed double-differentially in bins of jet transverse momentum and ρ_{iso} . In
334 the central region, a total of 21 bin in transverse momentum are defined in the central region,
335 which span in logarithmic spacing between 20 and 4000 GeV: 20, 26, 32, 40, 51, 64, 80, 101, 127,
336 159, 201, 252, 317, 400, 503, 633, 797, 1003, 1262, 1589, 2000, 4000 GeV. In the forward region, all
337 jets with transverse momentum greater than 127 GeV are put in one big 127-4000 GeV bin, in
338 order to take into account the obvious energetic limitations at high rapidity.

339 As for the binning in ρ_{iso} , it is linearly spaced in 1 GeV stpdf, between 0 and 45 GeV. Then a
340 final bin integrates all values between 45 and 100 GeV.

Once the three variables' probability density functions (in the form of histograms derived from the simulation) are obtained in all of the $p_T \times \rho$ bins, a likelihood discriminant can be defined. This has been accomplished as a simple product of the three variables' distributions, in each $p_T \times \rho$ bin. A given PFJet identifies a vector \vec{x} in the three-dimensional space of the structure variables. Probability density functions for gluons (G) and quarks (Q) can then be defined as the product of each variable's probability density function (f^i) computed at the given variable's value ($x[i]$):

$$G(\vec{x}) = \prod_i f_G^i(x[i]) \quad Q(\vec{x}) = \prod_i f_Q^i(x[i])$$

so that a likelihood estimator can be defined as:

$$L(\vec{x}) = \frac{Q(\vec{x})}{Q(\vec{x}) + G(\vec{x})}$$

341 and interpreted as the probability of a given PFJet to be originated from a quark parton.

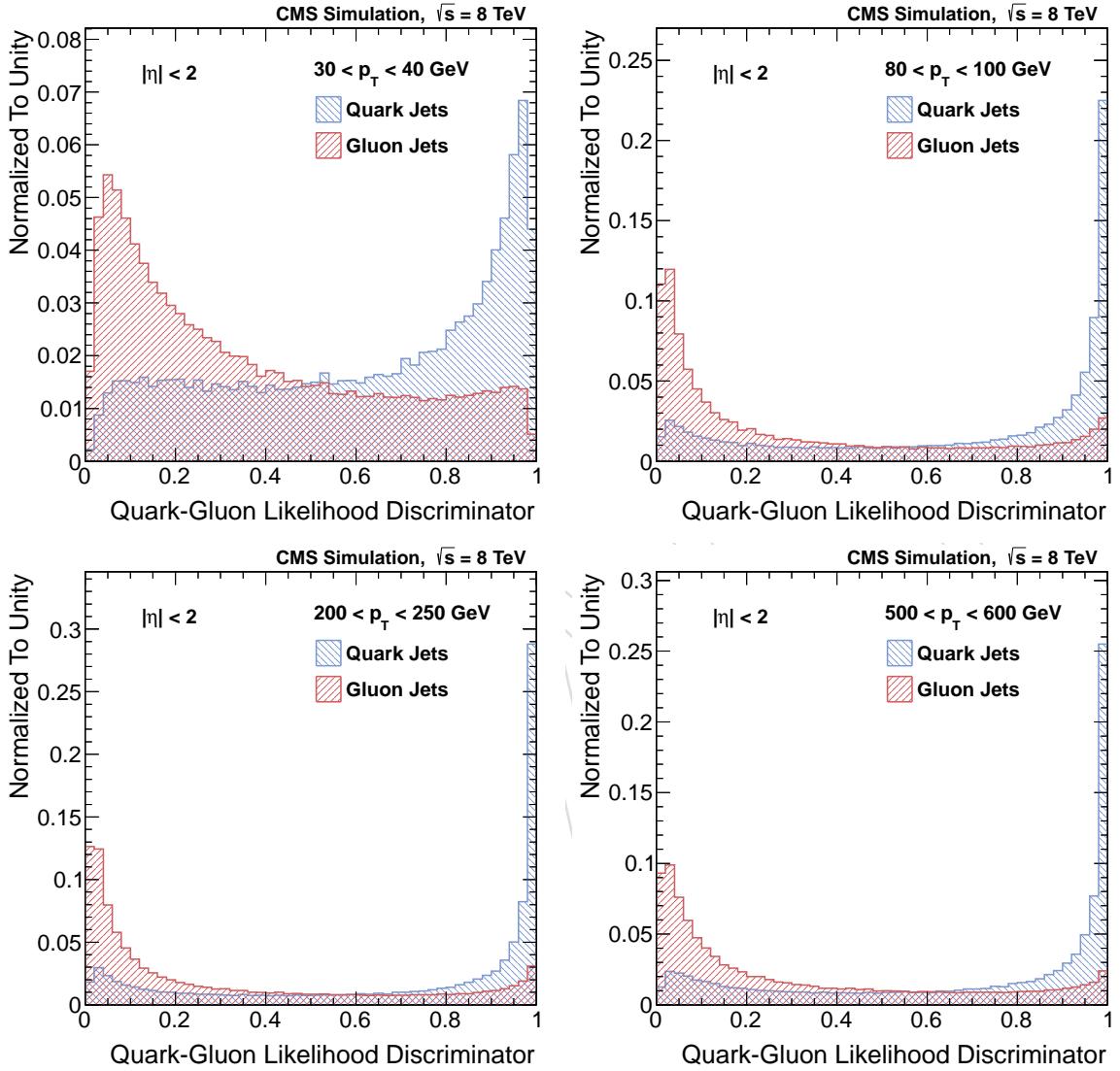


Figure 4.4: Normalized distributions of the Likelihood Quark-Gluon discriminant for quark (blue) and gluon (red) jets reconstructed in the central part of the detector ($|\eta| < 2$). Four representative transverse momentum bins are shown: $30 < p_T < 40 \text{ GeV}$ (top left), $80 < p_T < 100 \text{ GeV}$ (top right), $200 < p_T < 250 \text{ GeV}$ (bottom left) and $500 < p_T < 600 \text{ GeV}$ (bottom right).

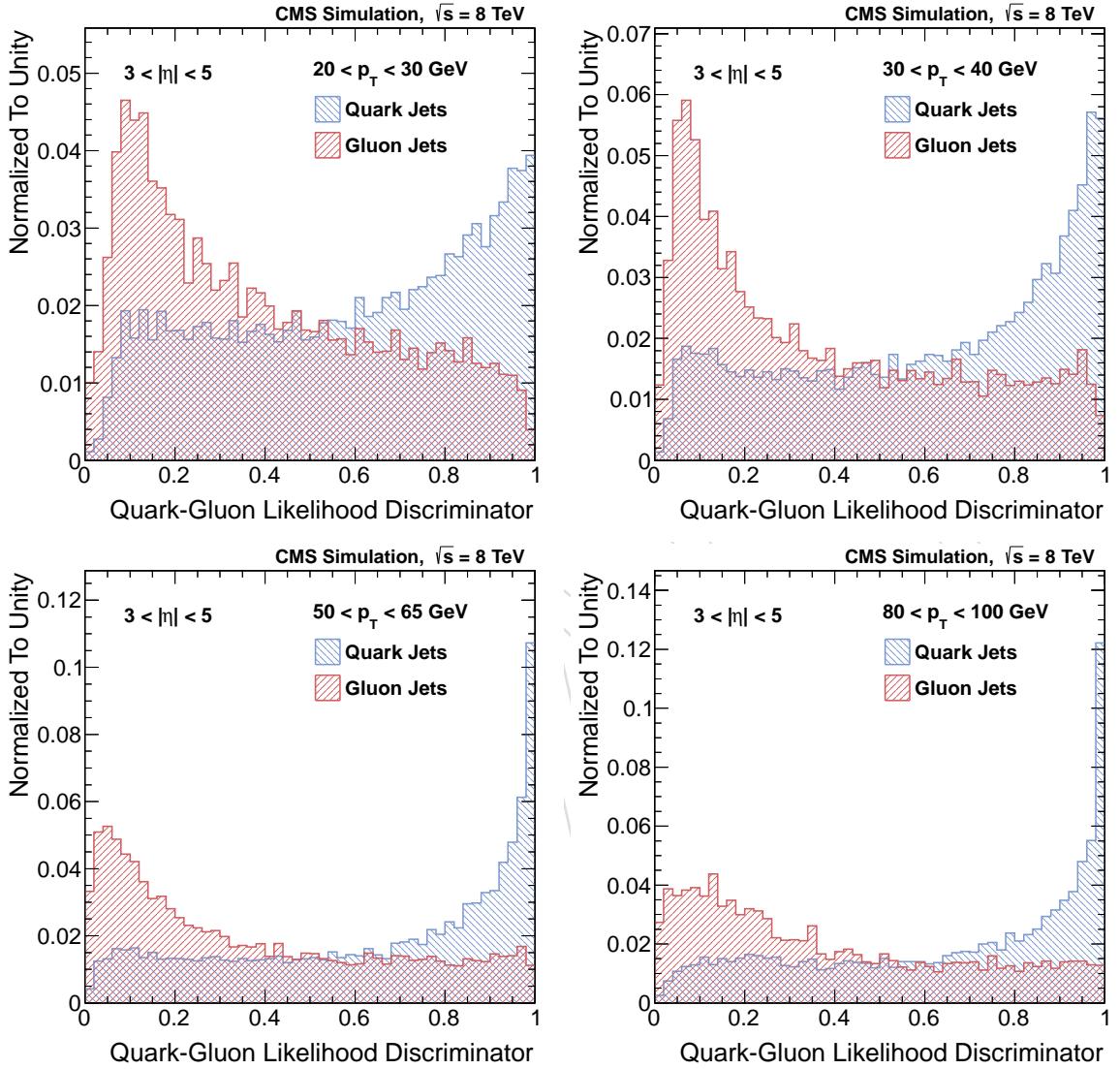


Figure 4.5: Normalized distributions of the Likelihood Quark-Gluon discriminant for quark (blue) and gluon (red) jets reconstructed in the forward part of the detector ($3 < |\eta| < 5$). Four representative transverse momentum bins are shown: $20 < p_T < 30 \text{ GeV}$ (top left), $30 < p_T < 40 \text{ GeV}$ (top right), $50 < p_T < 65 \text{ GeV}$ (bottom left) and $80 < p_T < 100 \text{ GeV}$ (bottom right).

342 Figures 4.4 and 4.5 show, respectively for the central and forward rapidity regions, normalized
 343 Likelihood distributions for quark and gluon jets in the simulation, in four representative
 344 transverse momentum bins: $30 < p_T < 40$ GeV, $80 < p_T < 100$ GeV, $200 < p_T < 250$ GeV
 345 and $500 < p_T < 600$ GeV in the central region, $20 < p_T < 30$ GeV, $30 < p_T < 40$ GeV,
 346 $50 < p_T < 65$ GeV and $80 < p_T < 80$ GeV in the forward region.

347 MLP Discriminant

348 The MLP discriminant is built from the MLP method, the recommended neural network in the
 349 TMVA framework, using four variables: nPFCand, axis₁, axis₂ and $p_T D$. This method is able to
 350 treat correlations and has good robustness against overtraining, but it has little transparency.

351 The neural network is trained on generated dijet events, which were divided in two rapidity
 352 regions: a central region ($|\eta| < 2.5$) and a forward region ($2.5 < |\eta| < 4.7$). To account for the
 353 strong dependence of the variables as a function of the jet transverse momentum, a separate
 354 MLP training was done for each of 20 p_T bins, logarithmically spaced between 20 GeV and
 355 600 GeV. In the forward region, only 18 bins were used covering a range between 20 GeV and
 356 427 GeV. For every p_T bin, the dependence of the variables as a function of the amount of PU
 357 activity in the event is measured by a linear fit of the 4 variables versus ρ . This fit is done in the
 358 2 GeV $< \rho < 30$ GeV range, as shown in figure 4.6. The four variables are then corrected for
 359 this dependence before passing them on to the MLP training phase.

360 A drawback with the MLP approach are the different shapes and ranges of the response distri-
 361 butions in the various p_T bins and rapidity regions. Therefore the MLP Quark-Gluon discrimi-
 362 nant uses the classification probability, also provided within the TMVA framework, in order to
 363 have a discriminating value which can be interpreted in the same way for all p_T bins.

364 In the application phase, the four variables are again corrected for the ρ dependence. The MLP
 365 classification probability is retrieved from the training bins which have their average jet p_T just
 366 above and below the tagging jet p_T . The final MLP Quark-Gluon discriminating value is then
 367 the linear interpolation of these two outputs.

368 Figures 4.7 and 4.8 show, respectively for the central and forward rapidity regions, normalized
 369 distributions of the MLP Quark-Gluon discriminant for quark and gluon jets in the simulation,
 370 in four representative transverse momentum bins: $30 < p_T < 40$ GeV, $80 < p_T < 100$ GeV,
 371 $200 < p_T < 250$ GeV and $500 < p_T < 600$ GeV in the central region, $20 < p_T < 30$ GeV,
 372 $30 < p_T < 40$ GeV, $50 < p_T < 65$ GeV and $80 < p_T < 80$ GeV in the forward region.

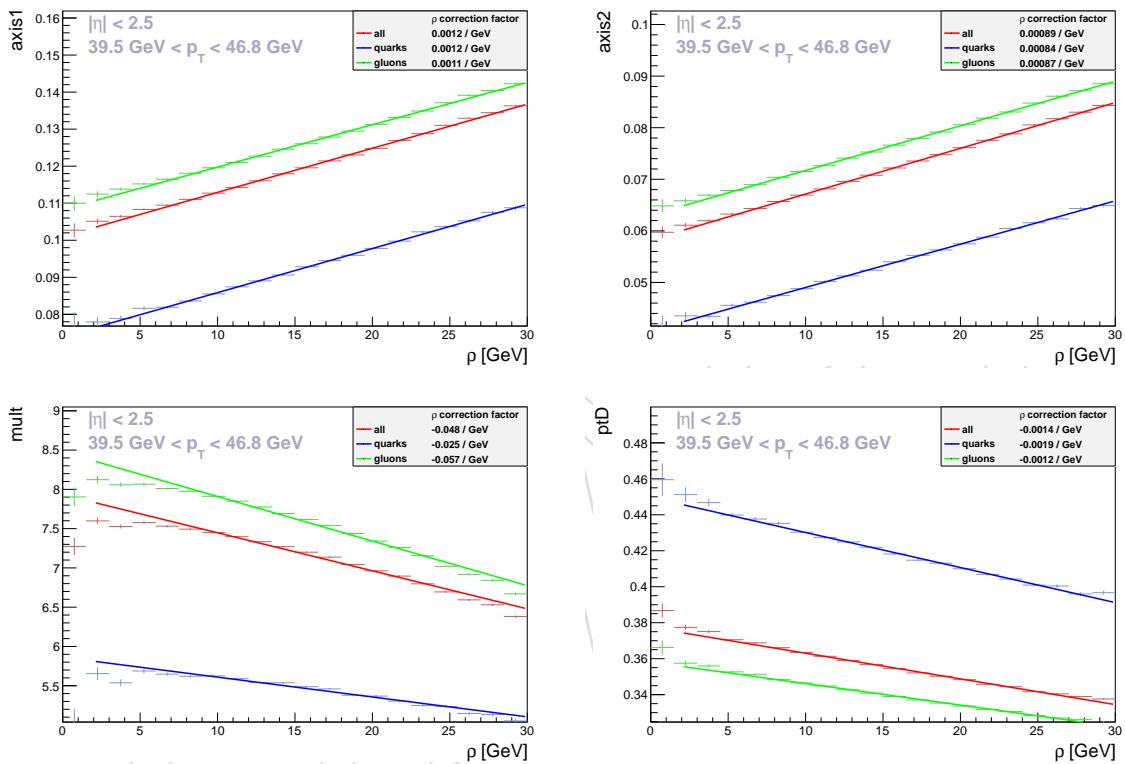


Figure 4.6: Linear fit of the ρ dependence for the axis_1 (top left), axis_2 (top right), $\text{nPF}C_{\text{and}}$ (bottom left) and $p_T D$ variables (bottom right) in one selected p_T bin. The correction factor is extracted from the fit to the complete sample, the separate gluon and quark fits are shown as reference.

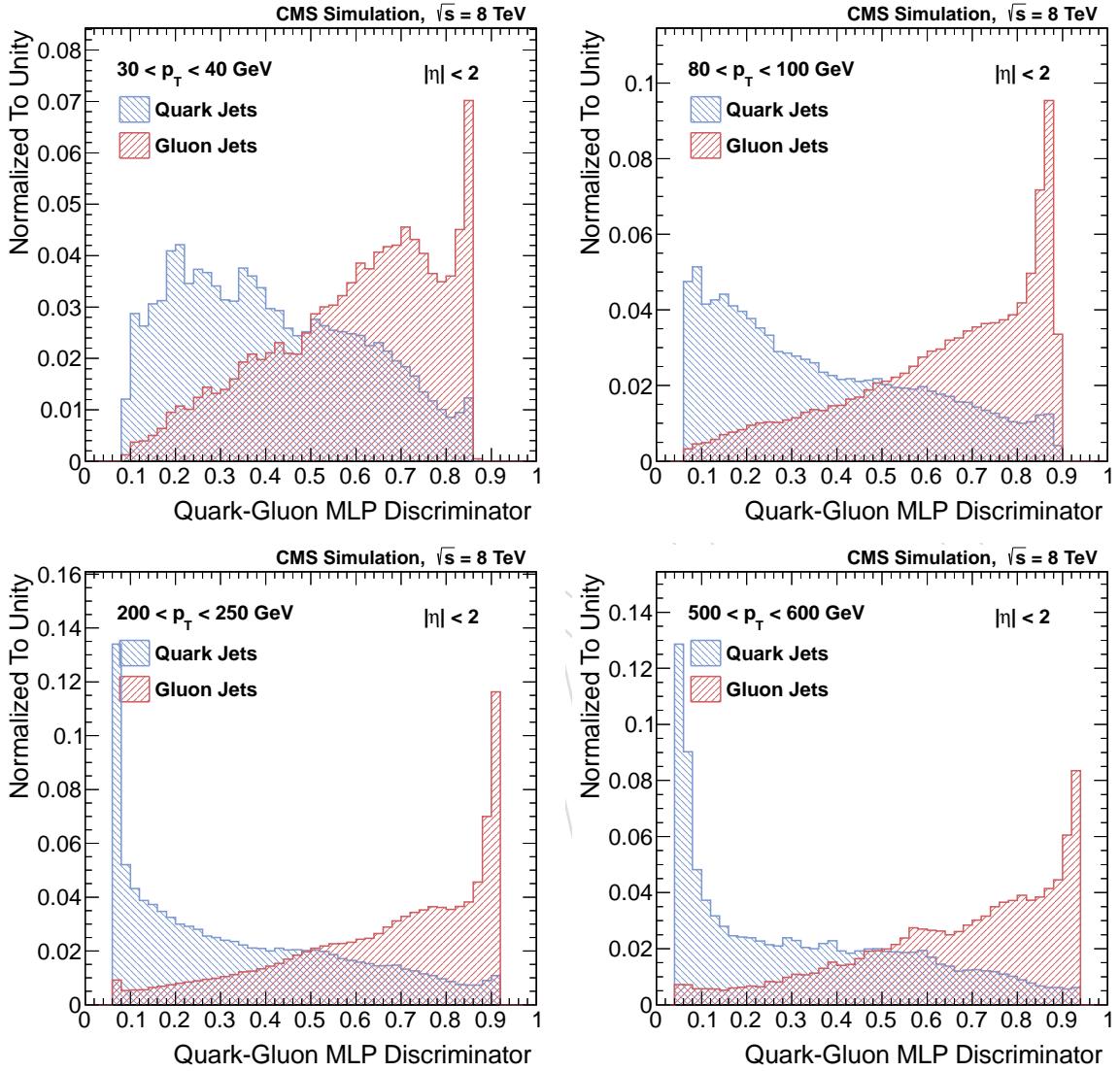


Figure 4.7: Normalized distributions of the MLP Quark-Gluon discriminant for quark (blue) and gluon (red) jets reconstructed in the central part of the detector ($|\eta| < 2$). Four representative transverse momentum bins are shown: $30 < p_T < 40 \text{ GeV}$ (top left), $80 < p_T < 100 \text{ GeV}$ (top right), $200 < p_T < 250 \text{ GeV}$ (bottom left) and $500 < p_T < 600 \text{ GeV}$ (bottom right).

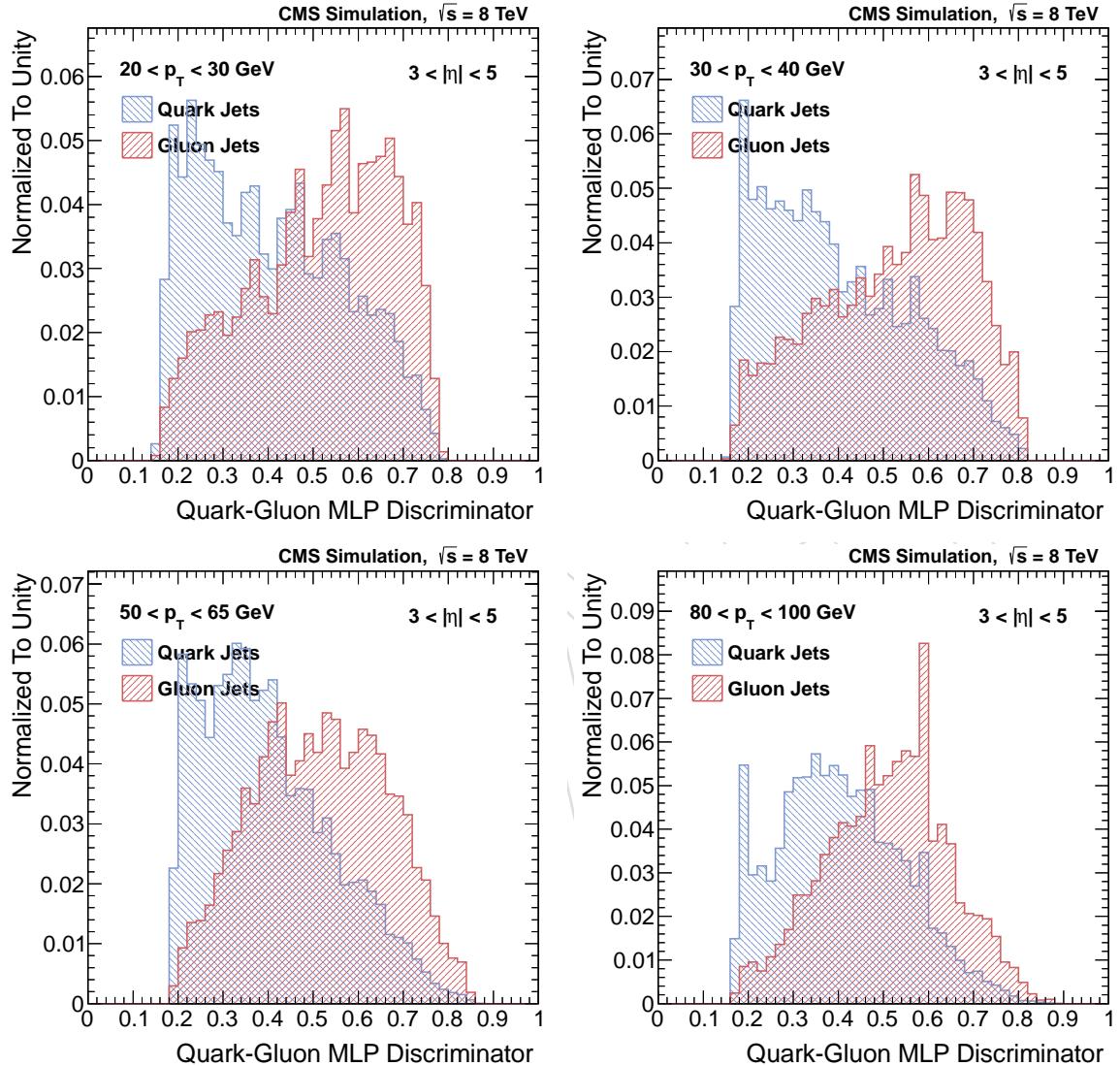


Figure 4.8: Normalized distributions of the MLP Quark-Gluon discriminant for quark (blue) and gluon (red) jets reconstructed in the forward part of the detector ($3 < |\eta| < 5$). Four representative transverse momentum bins are shown: $20 < p_T < 30 \text{ GeV}$ (top left), $30 < p_T < 40 \text{ GeV}$ (top right), $50 < p_T < 65 \text{ GeV}$ (bottom left) and $80 < p_T < 100 \text{ GeV}$ (bottom right).

373 **4.3.1 Performance Comparison**

374 A comparison of the performance of the two taggers is presented, in terms of ROC curves built
375 on the simulation, in Figures 4.9, 4.10 and 4.11, respectively for central rapidity ($|\eta| < 2$),
376 transition ($2 < |\eta| < 2.5$) and the forward ($3 < |\eta| < 5$), in four representative transverse
377 momentum intervals. The new likelihood discriminant is shown with hollow brown markers
378 and it is compared to the MLP, which is shown with solid green markers. In the tracker covered
379 region, the two are further compared to the old likelihood discriminant, shown in yellow.

380 The following conclusions can be made:

- 381 • both the new likelihood and the MLP increase the discrimination performance of the
382 old tagger, across the phase space;
- 383 • in the central part of the detector, the two new taggers have similar performance
384 at low p_T , but the MLP's ability to handle correlations gives it an edge for high p_T
385 (> 200 GeV) jets;
- 386 • discrimination has been successfully extended to the forward region of the detector;
- 387 • in the forward, the two taggers again have similar performance, but this time the
388 likelihood discriminant seems to lead to better discrimination.

DRAFT

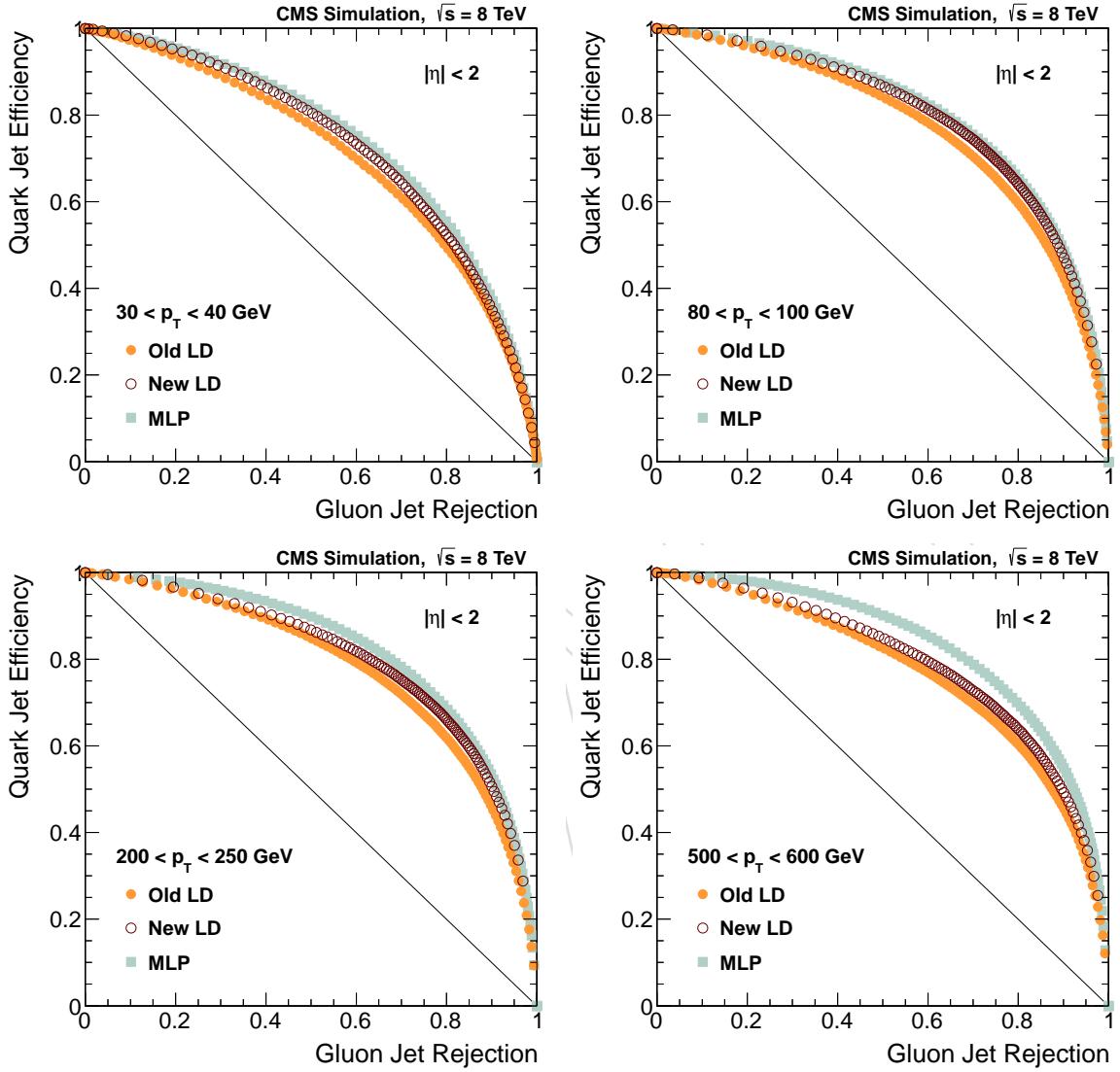


Figure 4.9: Receiver operating characteristic (ROC) curves for jets reconstructed in the central part of the detector ($|\eta| < 2$), in four representative transverse momentum intervals: $30 < p_T < 40 \text{ GeV}$ (top left), $80 < p_T < 100 \text{ GeV}$ (top right), $200 < p_T < 250 \text{ GeV}$ (bottom left) and $500 < p_T < 600 \text{ GeV}$ (bottom right). The likelihood discriminator (open brown markers) is compared to the MLP discriminant (green markers). The two are further compared to the old (2011) likelihood discriminant (yellow markers).

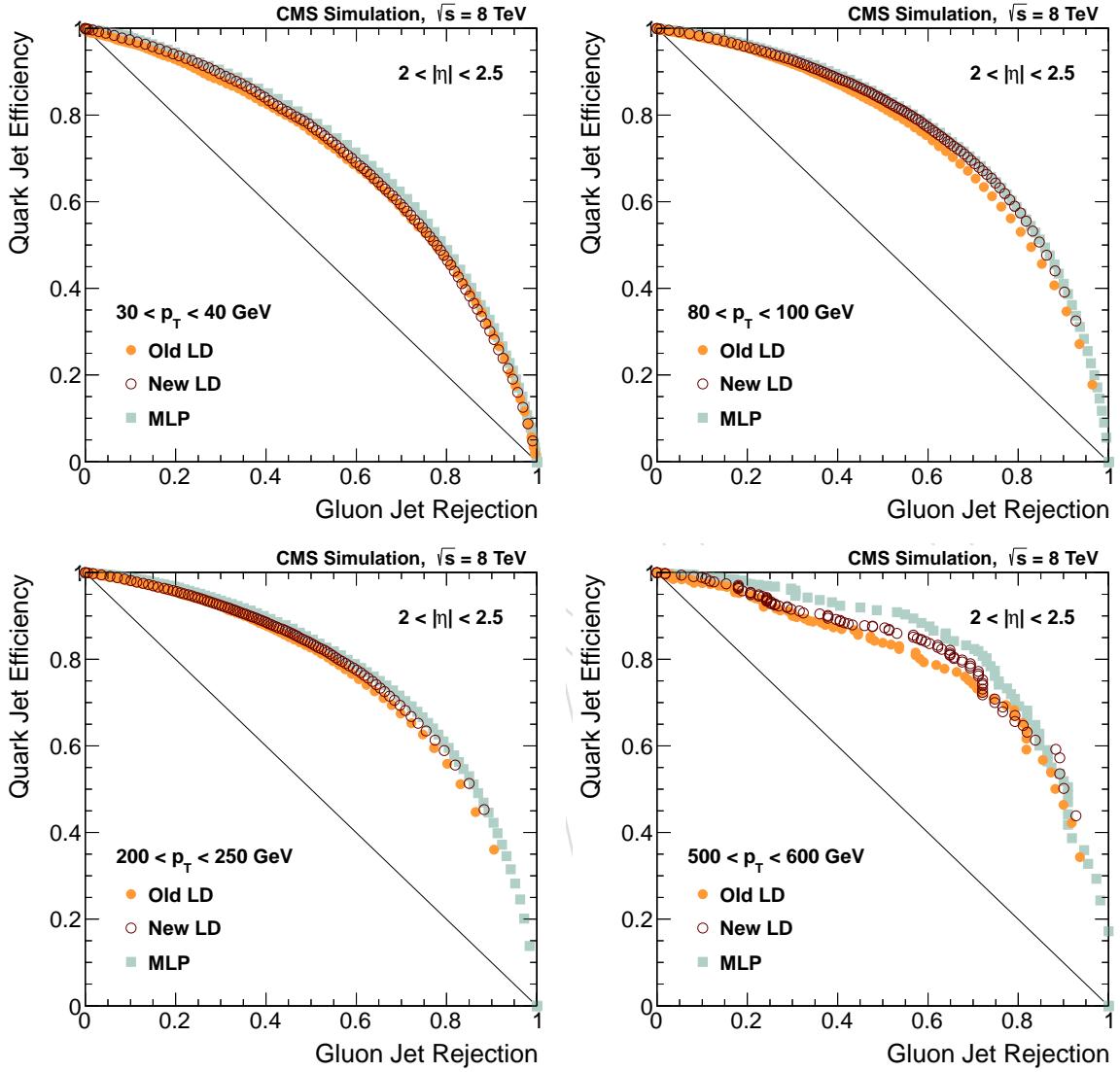


Figure 4.10: Receiver operating characteristic (ROC) curves for jets reconstructed in the detector transition region ($2 < |\eta| < 2.5$), in four representative transverse momentum intervals: $30 < p_T < 40 \text{ GeV}$ (top left), $80 < p_T < 100 \text{ GeV}$ (top right), $200 < p_T < 250 \text{ GeV}$ (bottom left) and $500 < p_T < 600 \text{ GeV}$ (bottom right). The likelihood discriminator (open brown markers) is compared to the MLP discriminant (green markers). The two are further compared to the old (2011) likelihood discriminant (yellow markers).

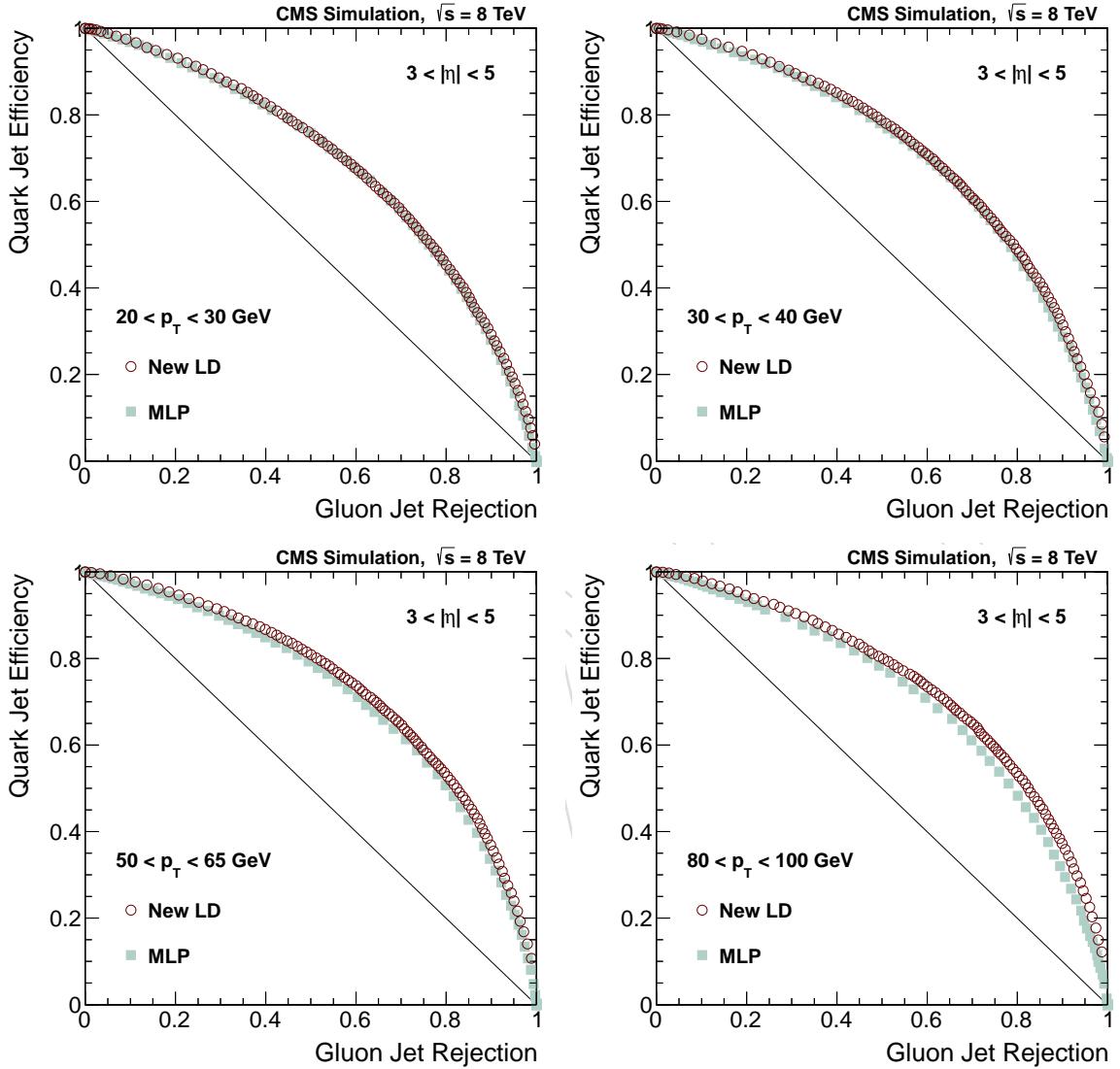


Figure 4.11: Receiver operating characteristic (ROC) curves for jets reconstructed in the detector forward ($3 < |\eta| < 5$), in four representative transverse momentum intervals: $20 < p_T < 30 \text{ GeV}$ (top left), $30 < p_T < 40 \text{ GeV}$ (top right), $50 < p_T < 65 \text{ GeV}$ (bottom left) and $80 < p_T < 100 \text{ GeV}$ (bottom right). The likelihood discriminator (open brown markers) is compared to the MLP discriminant (green markers).

389 5 Pythia6/Herwig++ Comparisons

390 As most of the results presented in this note are based on events simulated in Pythia 6, we
 391 here want to investigate how different things would look like by using a different parton
 392 shower/hadronization model. In order to do so we compare Pythia6 with Herwig++. This
 393 will become even more relevant in the following validation Sections, as it is observed than
 394 under some kinematical regimes Herwig++ seems to describe the data in a more satisfactory
 395 manner.

396 All of the results presented in this Section are based on the QCD ‘flat’ samples, namely:

397 Pythia6: /QCD_Pt-15to3000_TuneZ2star_Flat_8TeV_pythia6

398 Herwig++: /QCD_Pt-15to3000_TuneEE3C_Flat_8TeV_herwigpp

399 Both samples have been reconstructed in the same release (CMSSW_5_3_X) with the same pileup
 400 scenario (S10).

401 Comparisons between the two generators, for the four input variables, in three representative
 402 generator jet transverse momentum / rapidity bins, can be found in Figures 5.1, 5.2, 5.3 and
 403 5.4. As can be seen, for all variables, the observed trend is that quarks have very similar had-
 404 ronization properties in the two generators (at least for what can be probed by the considered
 405 variables), whereas there is a significant difference in gluons. Gluon jets in Herwig++ are much
 406 more similar to quarks, compared to Pythia. This is compatible with what has been reported
 407 in [11]. This is true for all intervals of transverse momentum and rapidity, we show here some
 408 bins just for the sake of conciseness.

409 This means that, if nature hadronizes like Herwig++, the discriminators presented in this note
 410 would have worse discrimination. This is shown in Figure 5.5, where the quark-gluon likeli-
 411 hood receiver operating characteristic obtained on Pythia QCD (open brown markers) is com-
 412 pared to the one obtained on Herwig++ (solid yellow markers), for three transverse momenta
 413 in the centerand one in the forward. As can be seen, for a quark efficiency, Herwig++ obtaines
 414 15-20% worse gluon rejection.

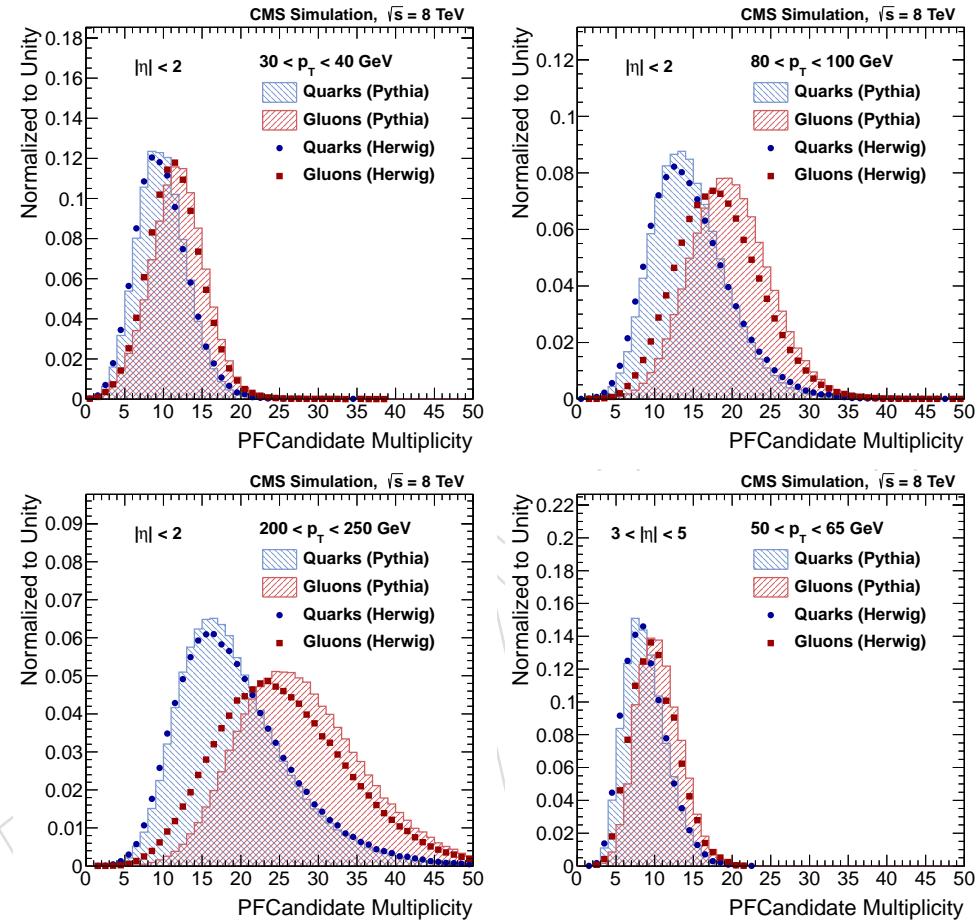


Figure 5.1: Comparison between Pythia and Herwig++ for the total PFCandidate multiplicity: Pythia distributions are shown with hashed histograms, Herwig++ ones with markers. Quark jets are shown in blue, gluon jets in red. Four phase space regions are shown: central jets with transverse momentum between 30 and 40 GeV (top left), 80 and 100 GeV (top right), 200 and 250 GeV (bottom left) and forward jets with transverse momentum between 50 and 65 GeV (bottom right).

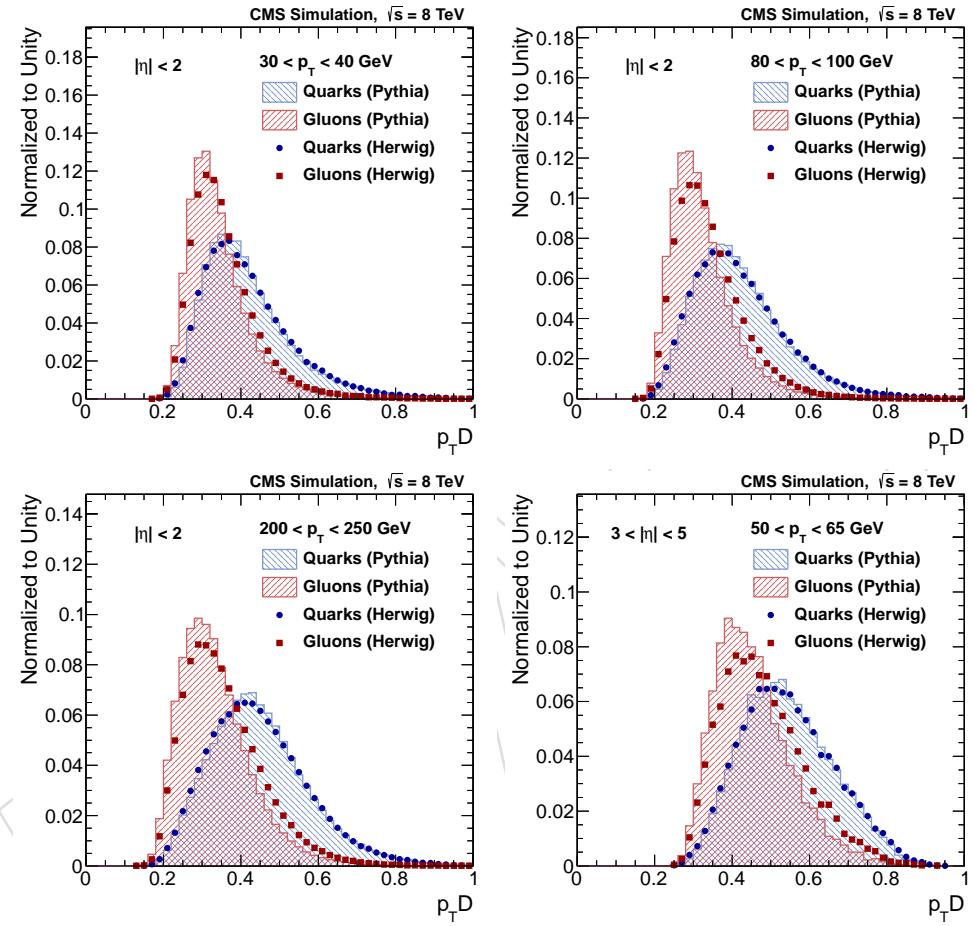


Figure 5.2: Comparison between Pythia and Herwig++ for p_{TD} : Pythia distributions are shown with hashed histograms, Herwig++ ones with markers. Quark jets are shown in blue, gluon jets in red. Four phase space regions are shown: central jets with transverse momentum between 30 and 40 GeV (top left), 80 and 100 GeV (top right), 200 and 250 GeV (bottom left) and forward jets with transverse momentum between 50 and 65 GeV (bottom right).

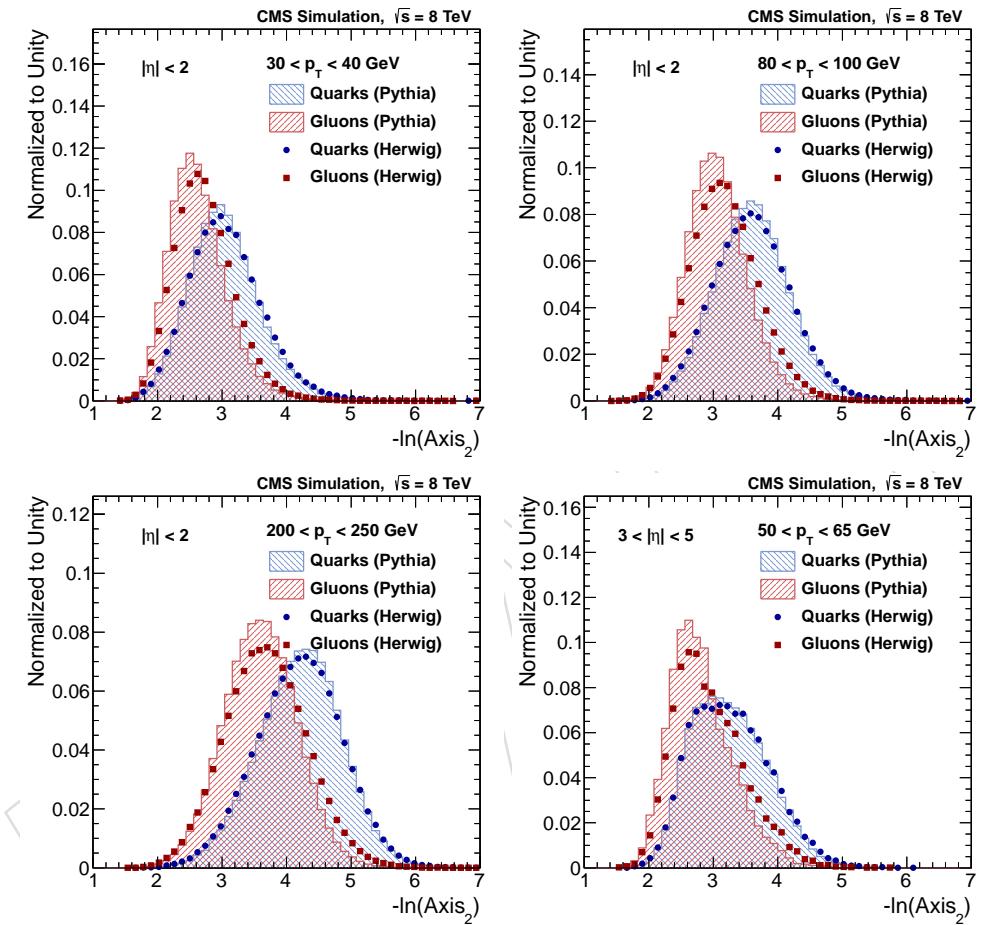


Figure 5.3: Comparison between Pythia and Herwig++ for Axis₂: Pythia distributions are shown with hashed histograms, Herwig++ ones with markers. Quark jets are shown in blue, gluon jets in red. Four phase space regions are shown: central jets with transverse momentum between 30 and 40 GeV (top left), 80 and 100 GeV (top right), 200 and 250 GeV (bottom left) and forward jets with transverse momentum between 50 and 65 GeV (bottom right).

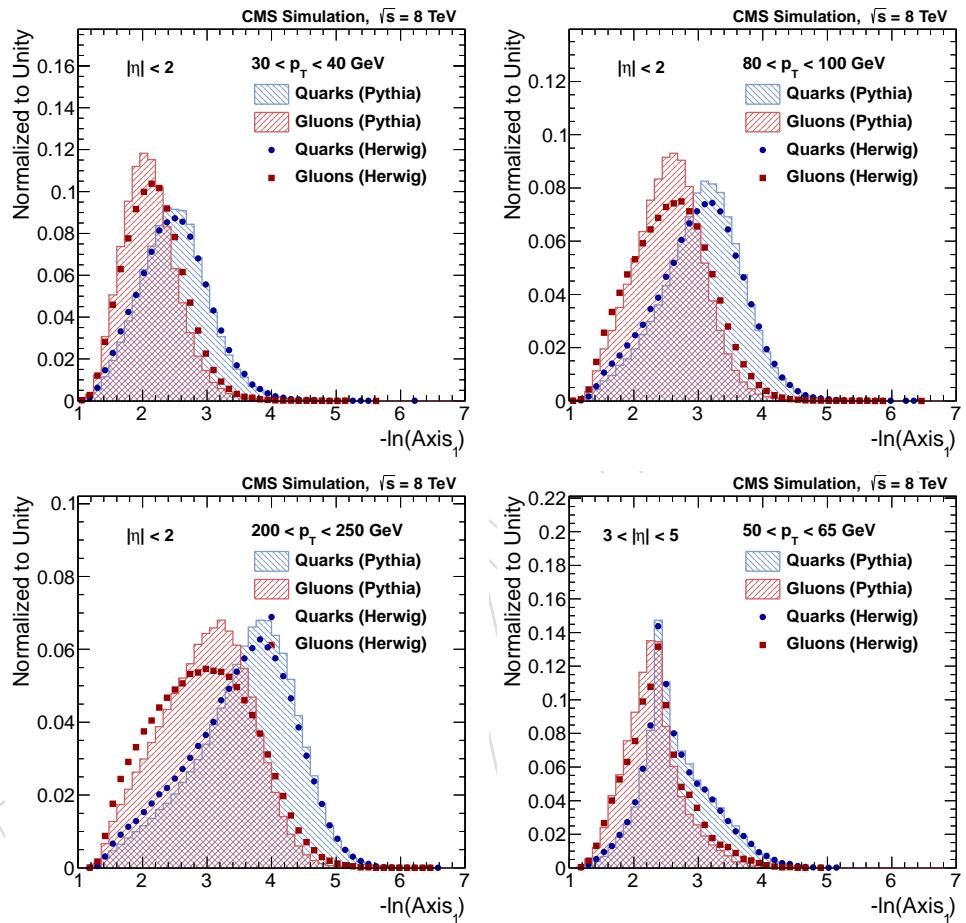


Figure 5.4: Comparison between Pythia and Herwig++ for Axis_1 : Pythia distributions are shown with hashed histograms, Herwig++ ones with markers. Quark jets are shown in blue, gluon jets in red. Four phase space regions are shown: central jets with transverse momentum between 30 and 40 GeV (top left), 80 and 100 GeV (top right), 200 and 250 GeV (bottom left) and forward jets with transverse momentum between 50 and 65 GeV (bottom right).

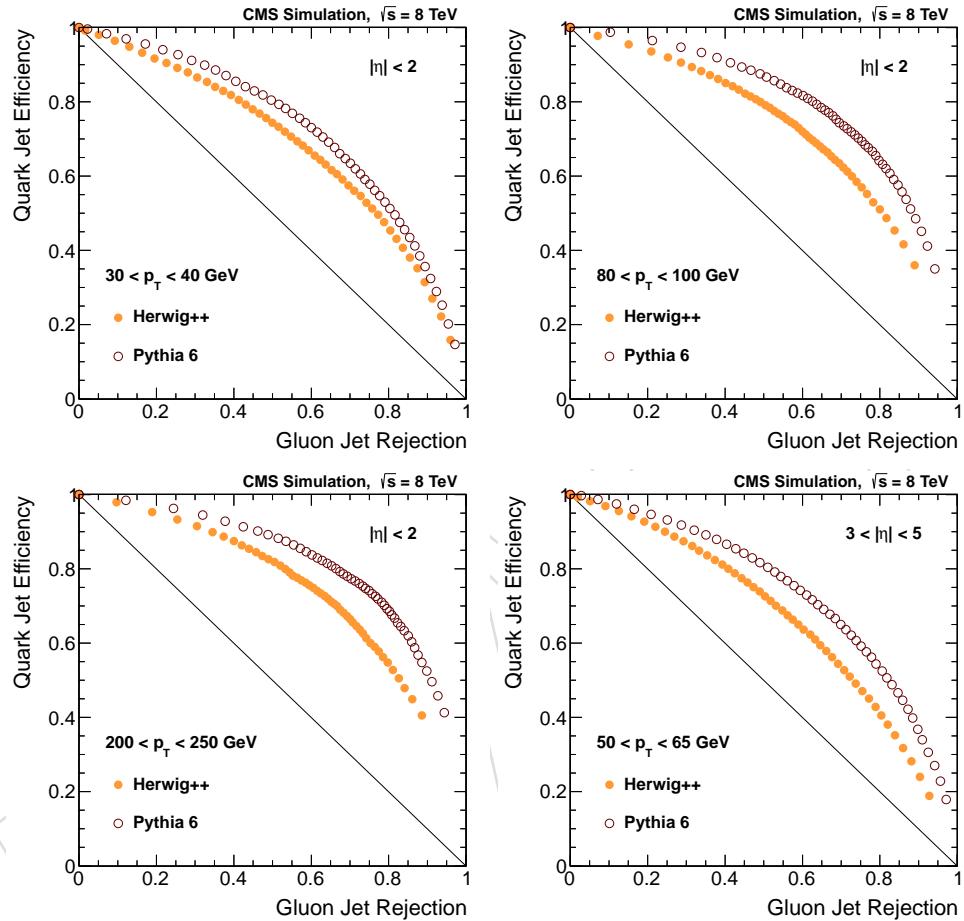


Figure 5.5: Performance, in terms of receiver operating characteristics on the quark efficiency / gluon rejection plane, of the likelihood quark-gluon discriminator for Pythia6 (hollow brown markers) and Herwig++ (yellow markers) jets. Four phase space regions are shown: central jets with transverse momentum between 30 and 40 GeV (top left), 80 and 100 GeV (top right), 200 and 250 GeV (bottom left) and forward jets with transverse momentum between 50 and 65 GeV (bottom right).

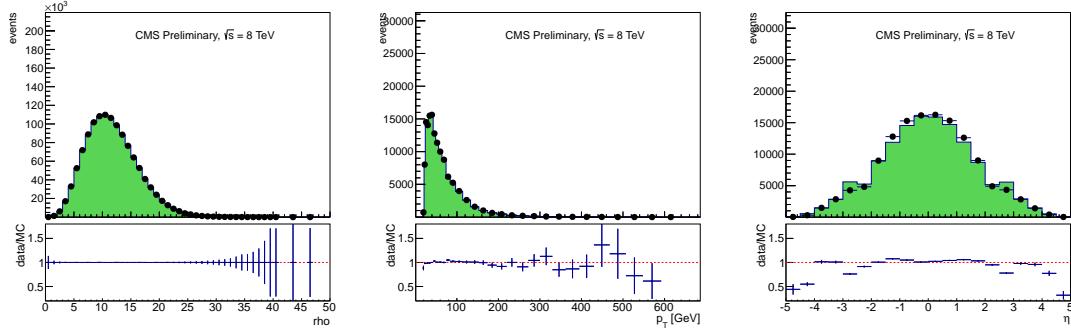


Figure 6.1: Data-MC comparisons for the particle-flow energy density ρ (left), the leading jet transverse momentum (center) and its pseudorapidity (right). The reweighting procedure is defined by the ρ distribution, so it is expected to find perfect agreement there.

6 Validation on data

We will now proceed to validate the discriminators on 8 TeV collision data. In order to do this, two main samples have been identified:

- **Z+jet events**, which are quark-enriched;
- **dijet events**, which are gluon-enriched.

By the simultaneous use of these two control samples, we can verify the functioning of the discriminators on both parton flavors, and across the whole phase space. The following subsections detail the event selection and the obtained results in these two control samples. In what follows, all Monte Carlo distributions are normalized to the integral of the data, as we are interested only in shape comparisons.

6.1 Validation on Z+jet Events

The Z+jet control sample offers a relatively pure quark set of jets, as can be seen in Figure ??, where the fraction of quark (blue), gluon (red), pile up (yellow) and undefined (grey) leading jets are shown as a function of the hard scattering transverse momentum (\hat{p}_T). As can be seen, 80% or more of leading jets in Z+jet events originate from light quark hadronization, across the transverse momentum spectrum.

We choose to use only events in which the Z boson has decayed to muons, as they are less sensitive to pile up effects and therefore provide the cleanest sample. We plan to include also electronic decays of the Z soon. In order to select Z+jet events in data, we run on the full 2012 dataset recorded with the dimuon HLT path HLT_Mu17_Mu8, and compare it to the Madgraph Drell-Yan sample (/DYJetsToLL_M-50_TuneZ2Star_8TeV-madgraph-tarball) produced in the Summer12 production with pile up profile S10. The event selection further requires:

- the presence of two tight muons of opposite charge;
- the transverse momenta of the muons are required to be greater than 20 and 10 GeV respectively;
- the dimuon invariant mass is required to fall in the 70-110 GeV range;
- the subleading jet in the event is required to have a transverse momentum smaller than 30% that of the dimuon system.

The leading jet is considered, and as previously mentioned it is required to pass jet ID, anti-b-tagging, pile up ID. In order to compare the data and the simulation under similar pile up

445 conditions, a pile up reweighting procedure is enforced. This is done by simply reweighting
446 the MC distribution of the particle flow energy density ρ to match the one observed in the data.
447 This leads to a satisfactory comparison also for the leading jet's transverse momentum and
448 pseudorapidity distributions, as is shown in Figure 6.1.

449 Figures 6.2, 6.3, 6.4 and 6.5 show, respectively, the data-MC comparisons of the four input
450 variables (total multiplicity, $p_T D$, Axis₂ and Axis₁), in four representative p_T bins, for jets re-
451 constructed in the center of the detector. The corresponding plots for the forward are shown
452 in Figures 6.6, 6.8, 6.9 and 6.10. As can be seen in Figure 6.6, a significant discrepancy between
453 data and the simulation is observed in the multiplicity of forward jets. As will be seen in the
454 next Section, a similar discrepancy is observed in an orthogonal dataset (dijets). We therefore
455 assume it is caused by a mismodeling of pile-up effects in the forward region. The fact that it
456 is seen only in the multiplicity and is not so striking in the other variables is explained by the
457 fact that the multiplicity variables is sensitive to first-order effects, whereas all other variables
458 are sensitive only to higher-order corrections.

459 We therefore decide to correct the data multiplicity distributions before they are fed to the
460 discriminators. This is similar, conceptually to what is done in L1 jet energy corrections, where
461 different Offset corrections are employed in data and MC, in order to take into account possible
462 discrepancies in pile up modeling. Data multiplicities for forward jets are therefore decreased
463 by one unity. This leads to better compatibility between the data and the simulation, as can be
464 seen by comparing the distributions before (Figure 6.6) and after (Figure 6.7) the correction.

465 The corresponding data-MC comparisons for the output discriminants are shown in Figures 6.11
466 for the likelihood discriminant and in Figure 6.14 for the MLP. The comparisons for the forward
467 discriminators are shown respectively in Figures 6.12 and 6.15 before the multiplicity correc-
468 tion, and in Figures 6.13 and 6.16 after the multiplicity correction. As can be seen significantly
469 better agreement is observed after the application of the correction.

470 Overall, the level of agreement between the data and the simulation is satisfactory.

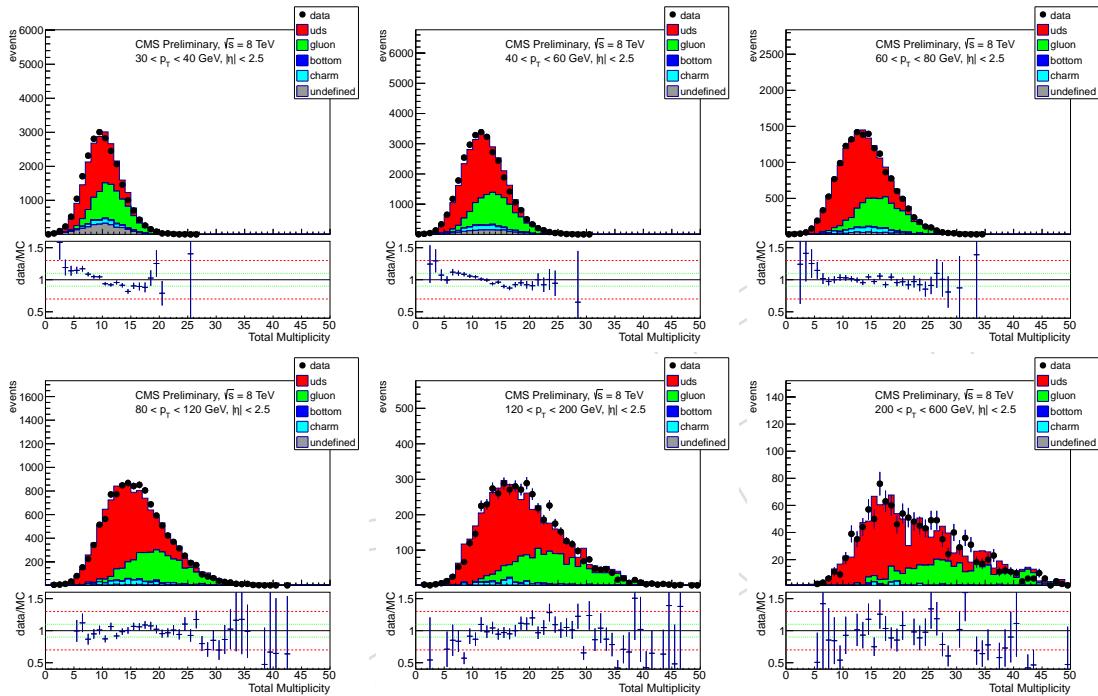


Figure 6.2: Data-MC comparisons for the total PFCandidate multiplicity in central jets in Z+jet events, for six p_T bins: 30-40, 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

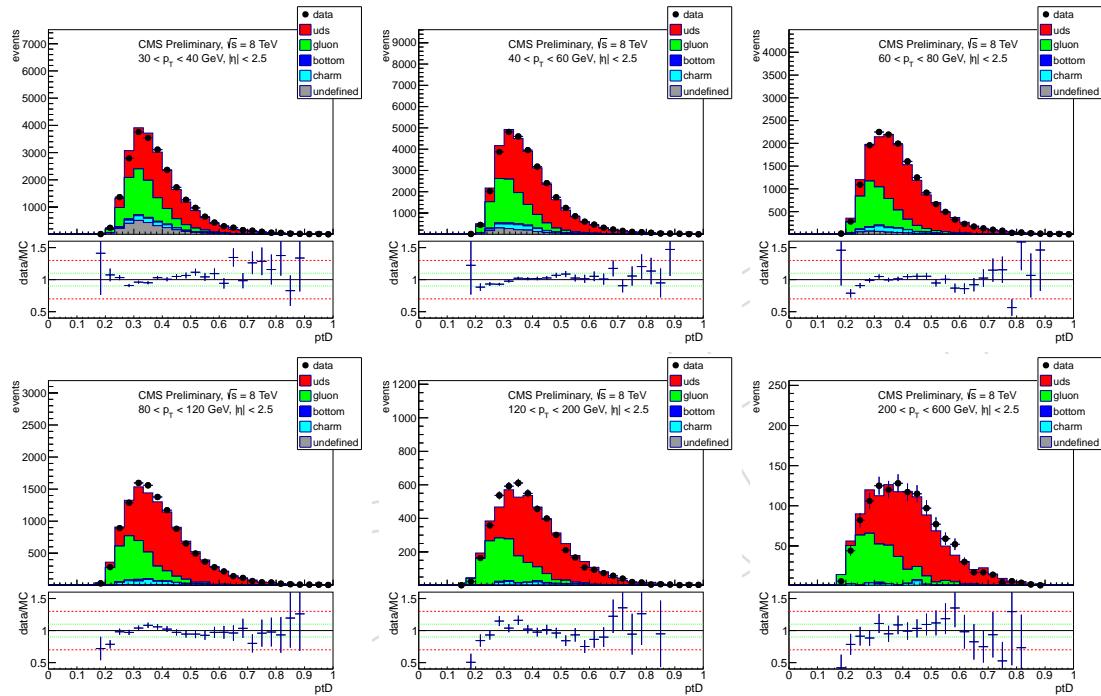


Figure 6.3: Data-MC comparisons for $p_T D$ in central jets in Z+jet events, for six p_T bins: 30-40, 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

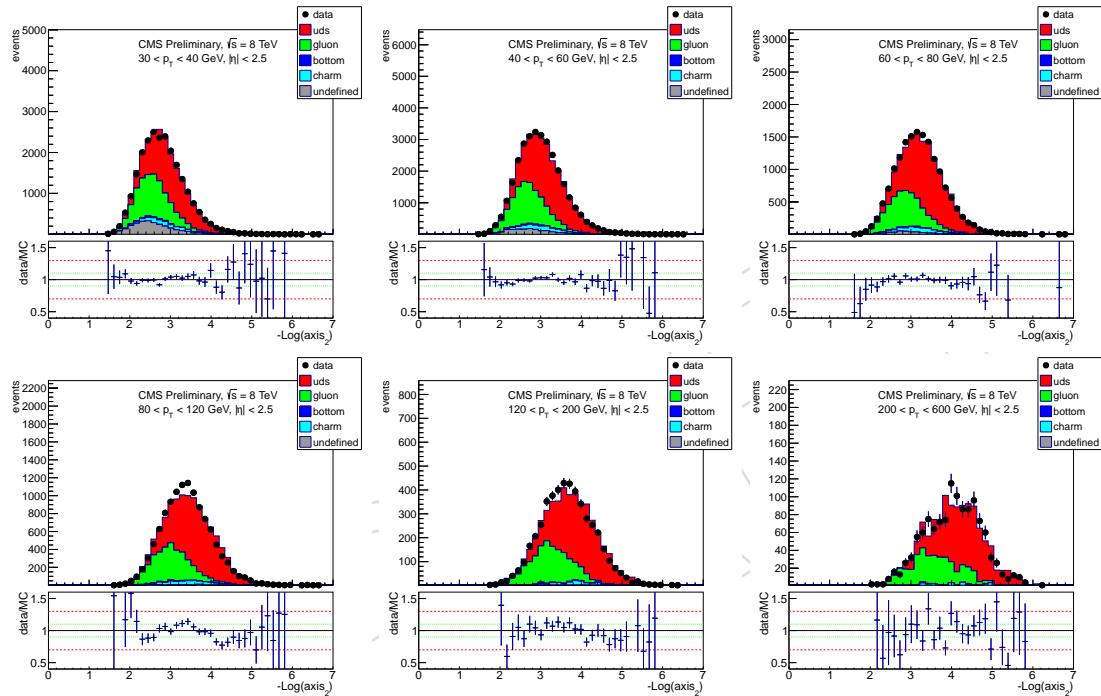


Figure 6.4: Data-MC comparisons for Axis₂ in central jets in Z+jet events, for six p_T bins: 30-40, 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

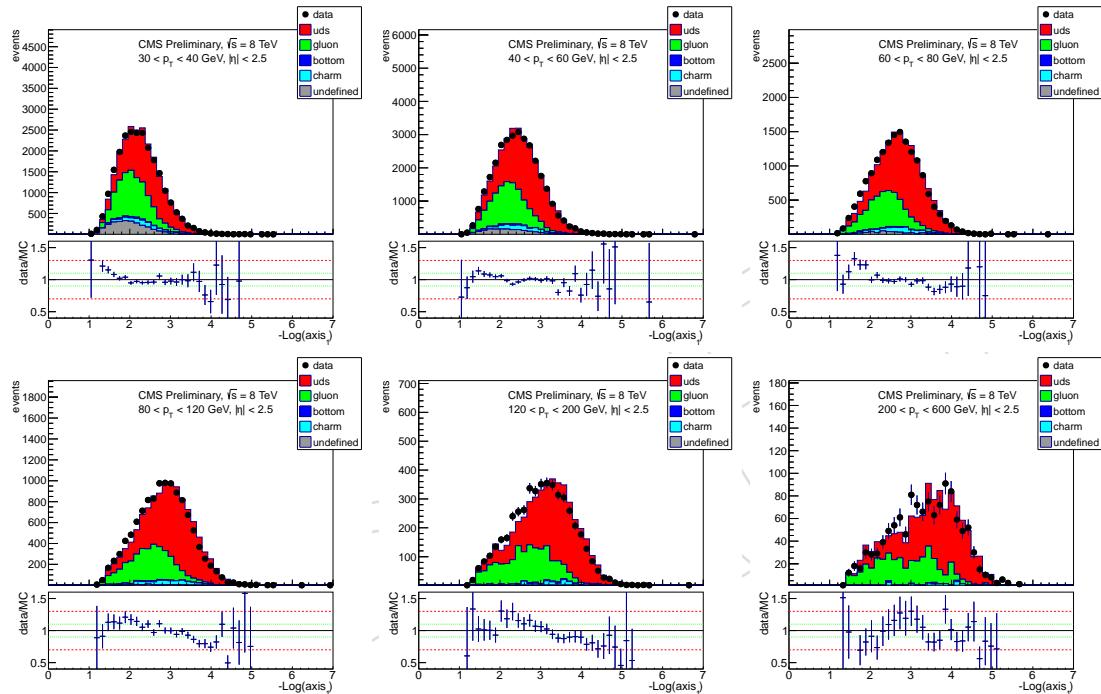


Figure 6.5: Data-MC comparisons for Axis_1 in central jets in Z+jet events, for six p_T bins: 30-40, 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

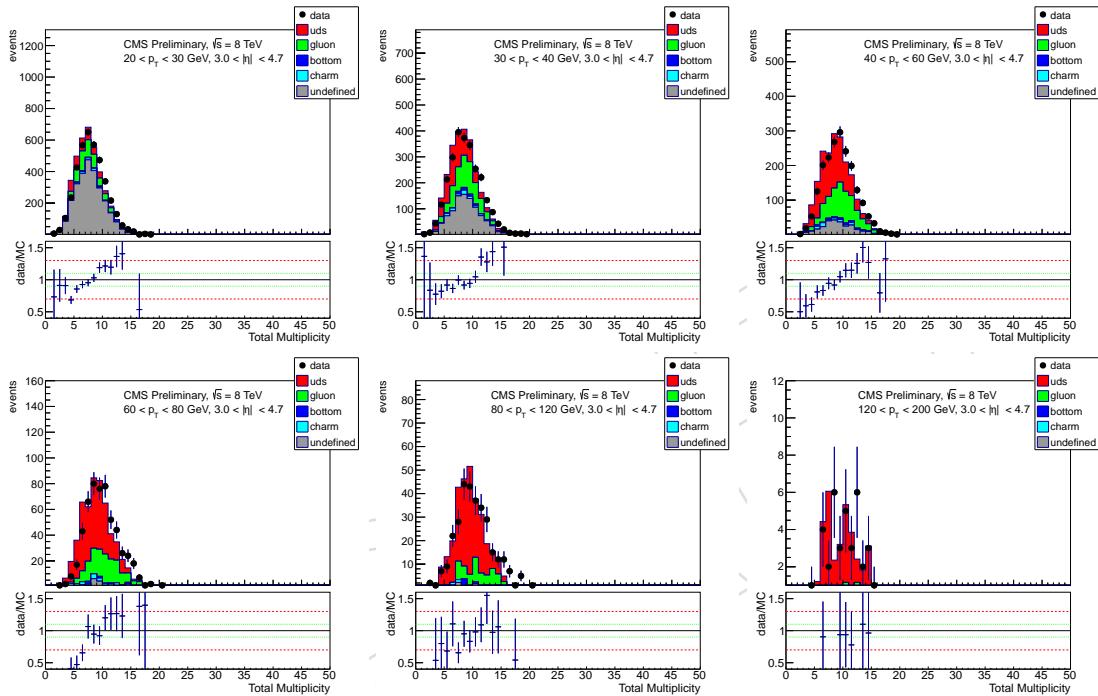


Figure 6.6: Data-MC comparisons for the total PFCandidate multiplicity in forward jets in Z+jet events **before the correction**, for six p_T bins: 20-30, 30-40, 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

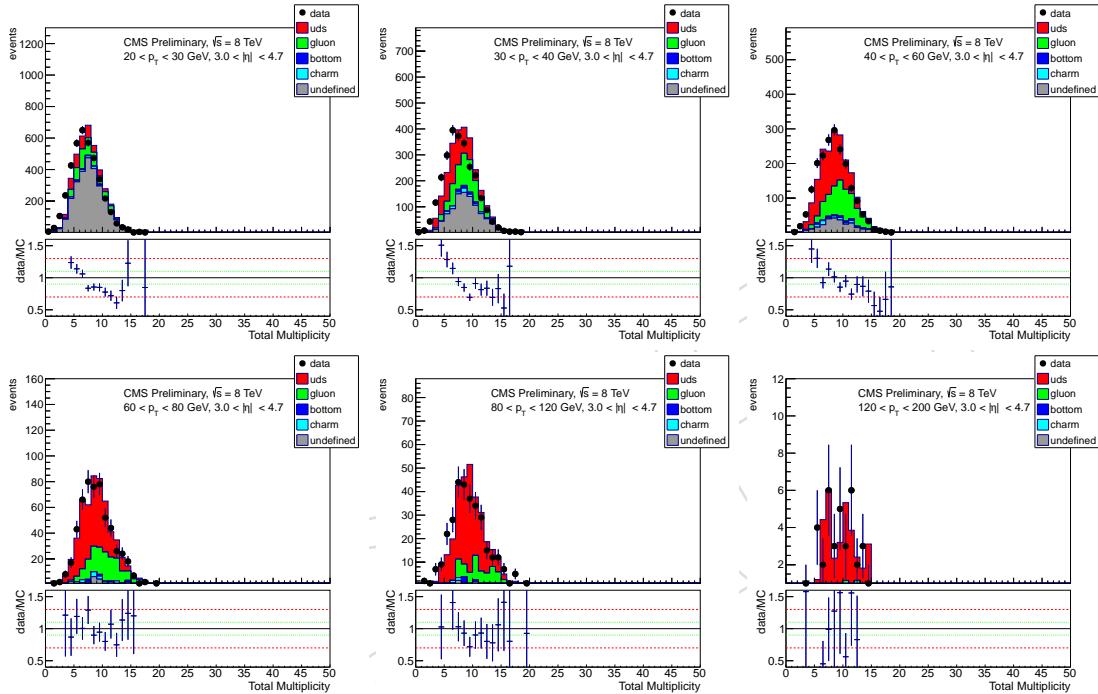


Figure 6.7: Data-MC comparisons for the total PFCandidate multiplicity in forward jets in Z+jet events **after the correction**, for six p_T bins: 20-30, 30-40, 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

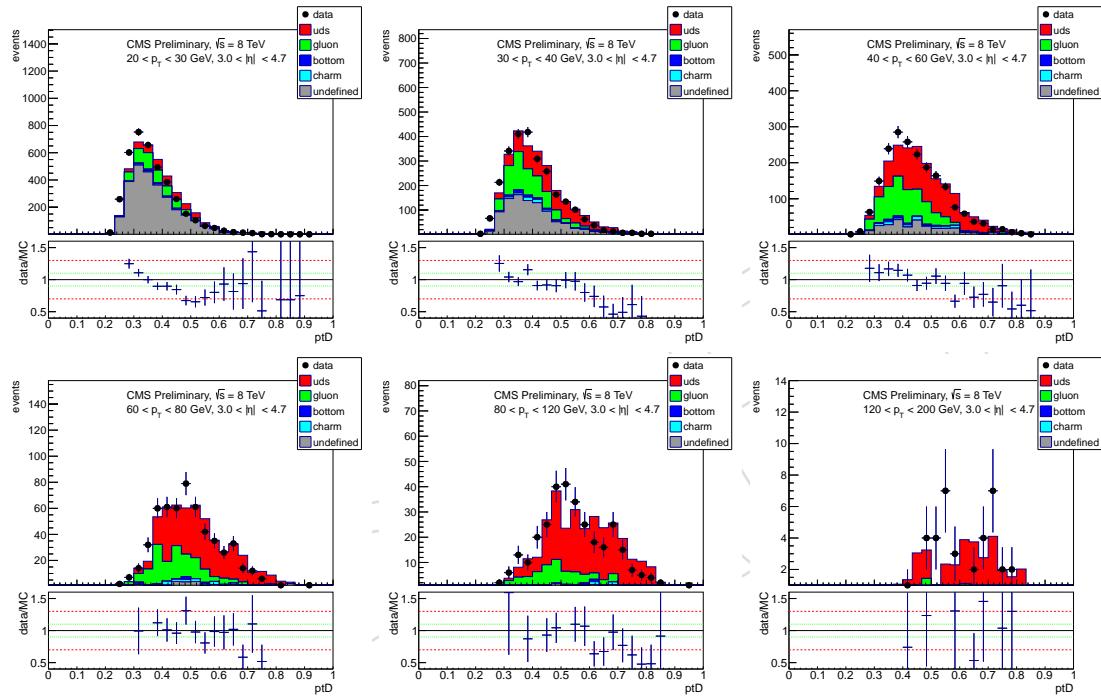


Figure 6.8: Data-MC comparisons for $p_T D$ in forward jets in Z+jet events, for six p_T bins: 20-30, 30-40, 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

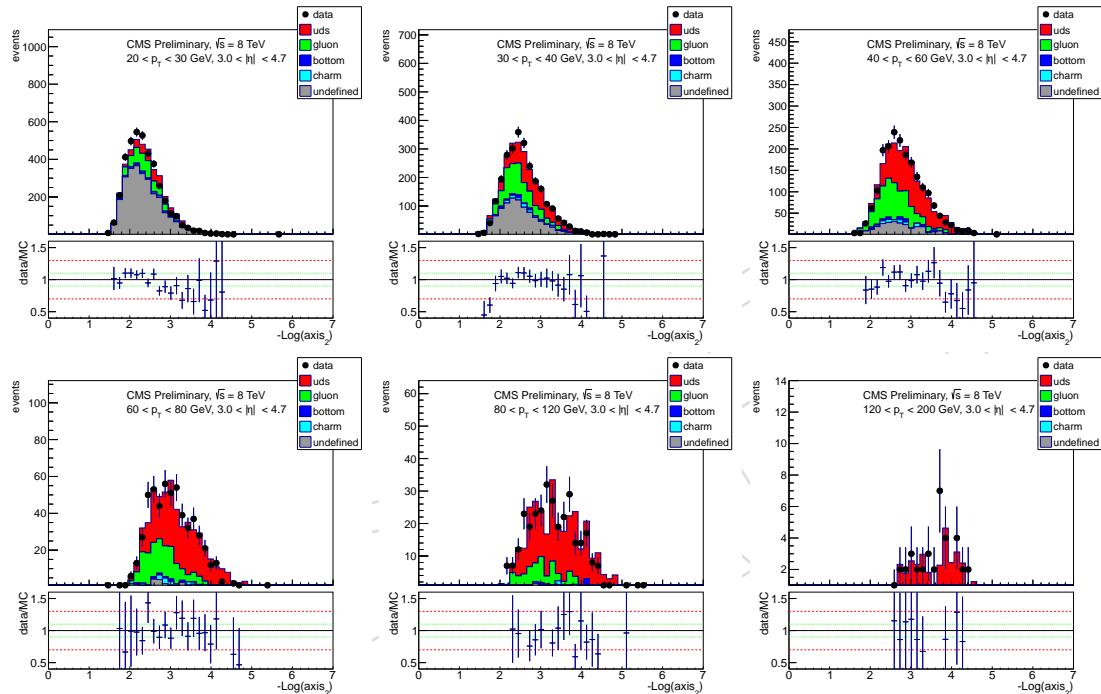


Figure 6.9: Data-MC comparisons for Axis₂ in forward jets in Z+jet events, for six p_T bins: 20-30, 30-40, 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

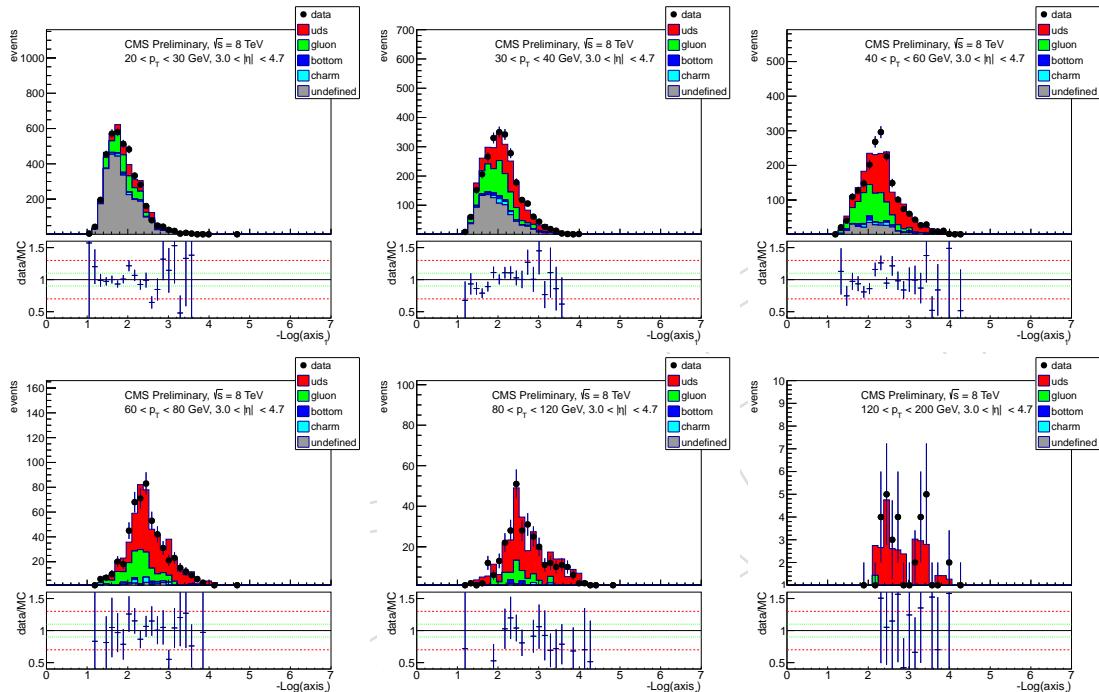


Figure 6.10: Data-MC comparisons for Axis_1 in forward jets in $Z+\text{jet}$ events, for six p_T bins: 20-30, 30-40, 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

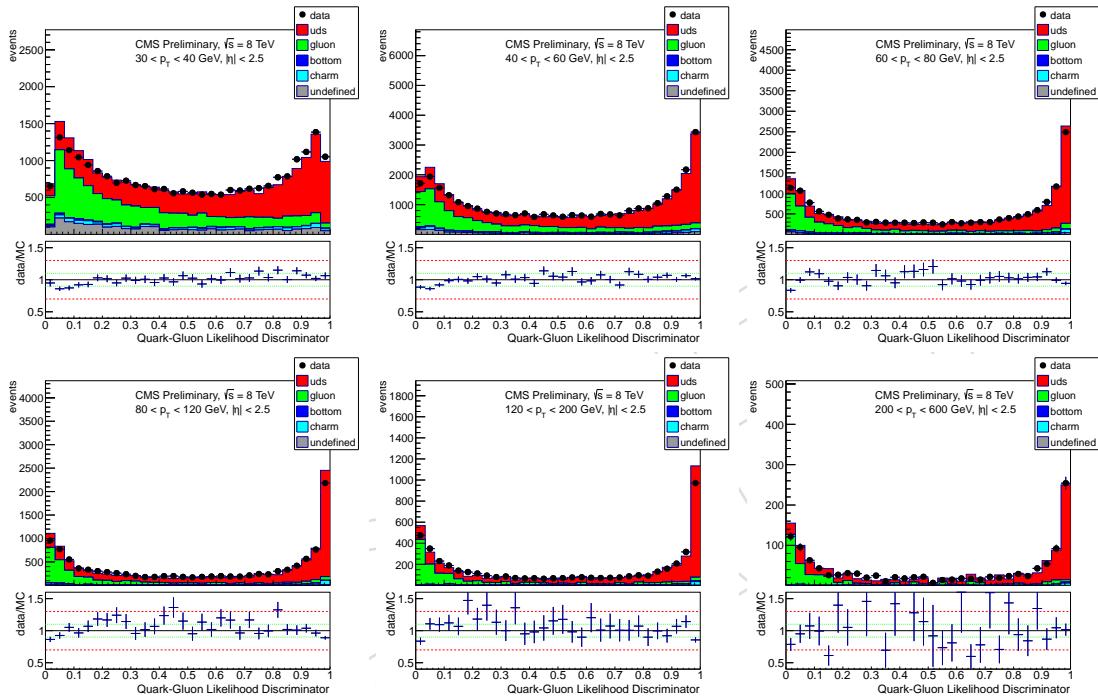


Figure 6.11: Data-MC comparisons for the quark-gluon likelihood discriminant in central jets in Z+jet events, for six p_T bins: 30-40, 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

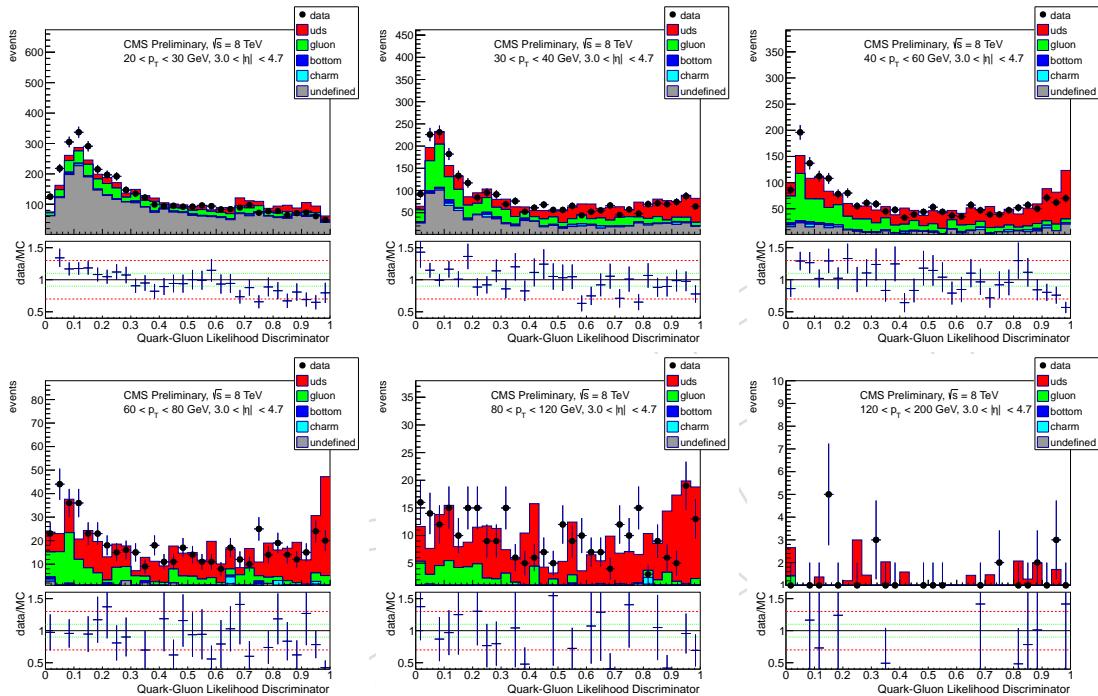


Figure 6.12: Data-MC comparisons for the quark-gluon likelihood discriminant in forward jets in Z+jet events **before the multiplicity correction**, for six p_T bins: 20-30, 30-40, 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

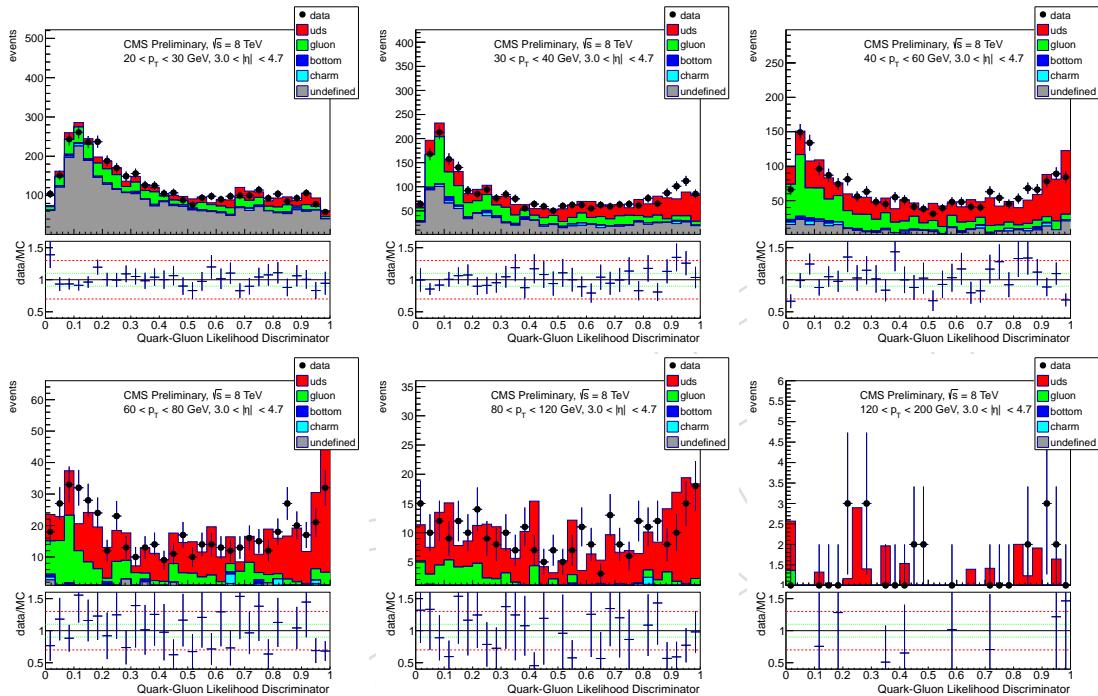


Figure 6.13: Data-MC comparisons for the quark-gluon likelihood discriminant in forward jets in Z+jet events **after the multiplicity correction**, for six p_T bins: 20-30, 30-40, 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

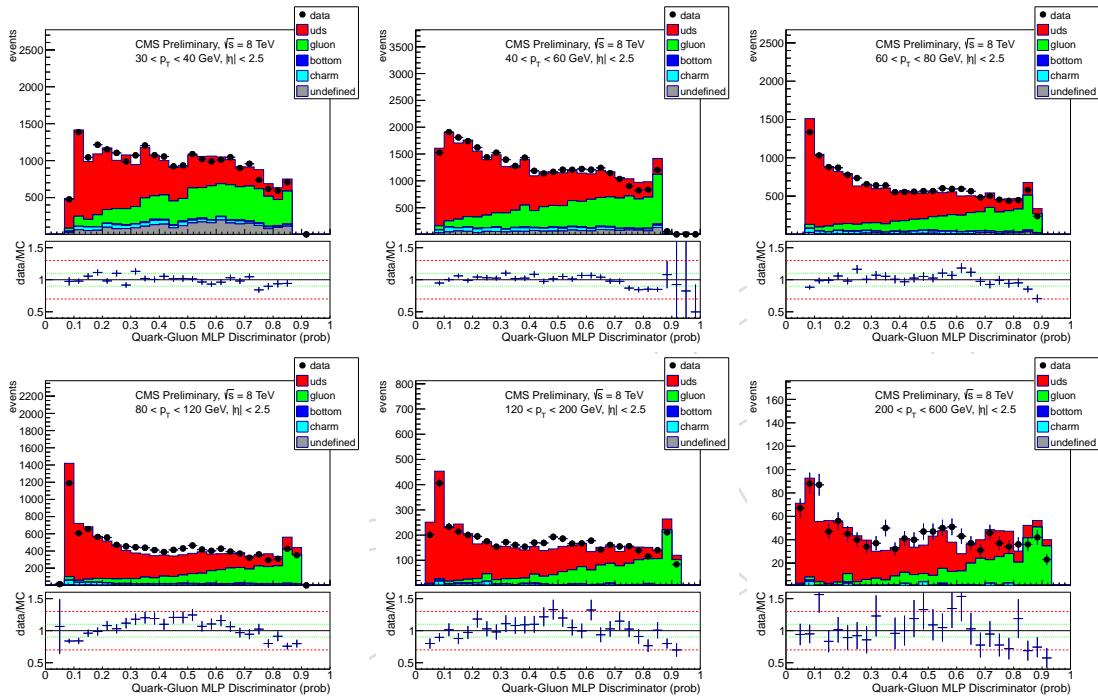


Figure 6.14: Data-MC comparisons for the quark-gluon MLP discriminant in central jets in Z+jet events, for six p_T bins: 30-40, 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

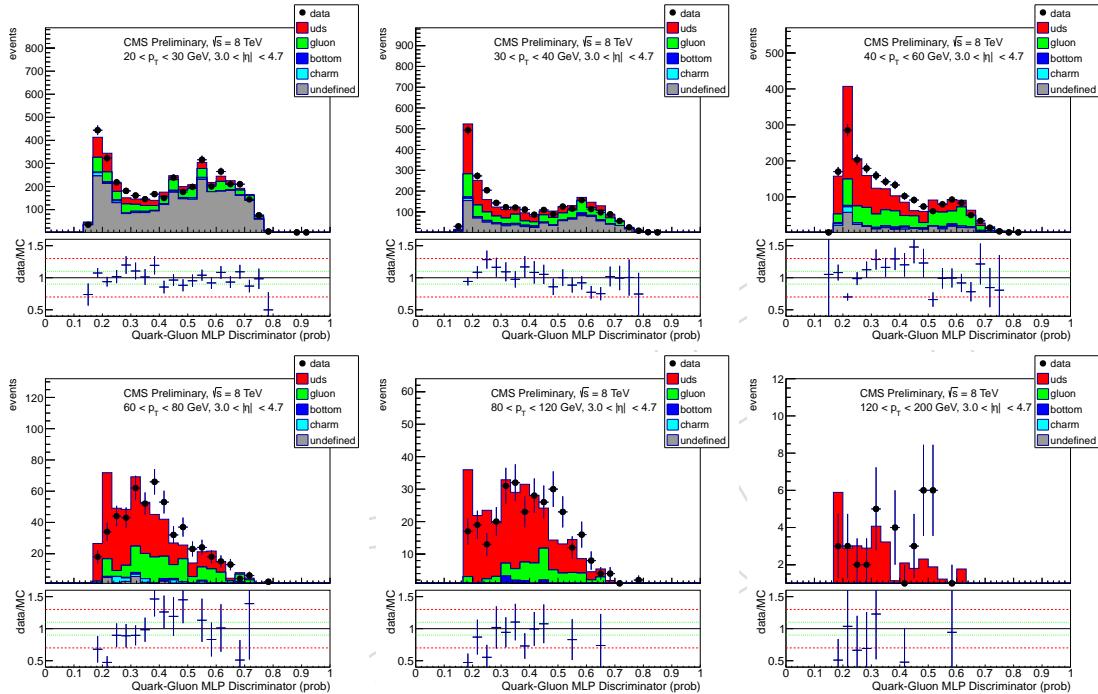


Figure 6.15: Data-MC comparisons for the quark-gluon MLP discriminant in forward jets in Z+jet events **before the multiplicity correction**, for six p_T bins: 20-30, 30-40, 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

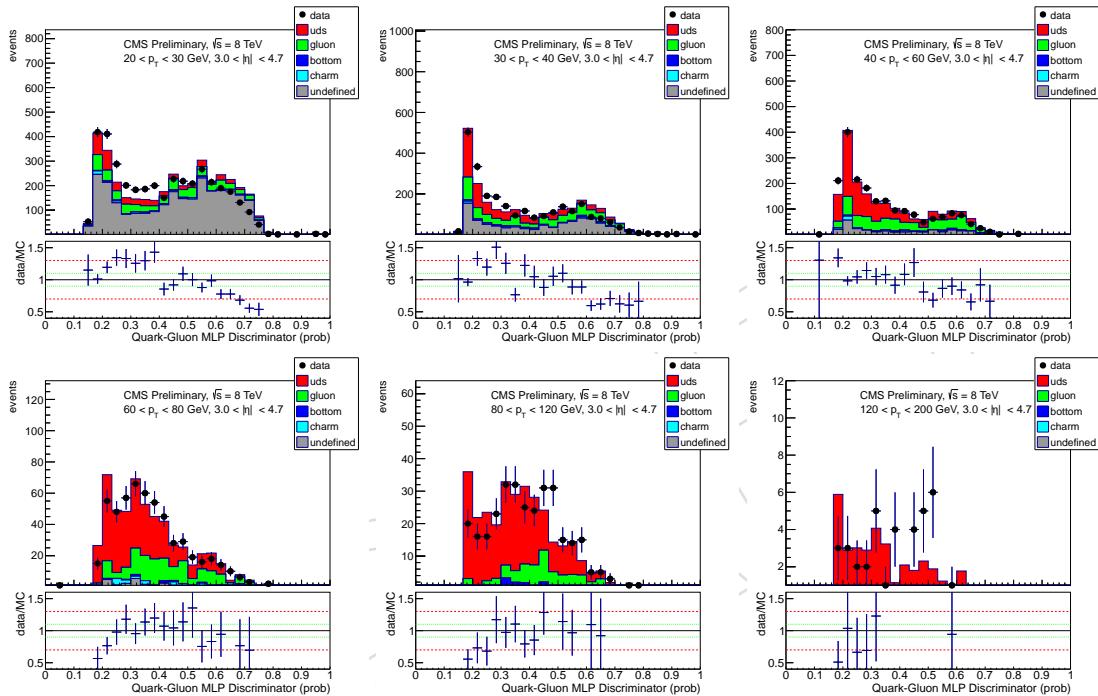


Figure 6.16: Data-MC comparisons for the quark-gluon MLP discriminant in forward jets in Z+jet events **after the multiplicity correction**, for six p_T bins: 20-30, 30-40, 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey). The data/MC ratio is showed on the bottom pad, where a solid black line marks unity, and the dotted green and red lines mark, respectively, $\pm 10\%$ and $\pm 30\%$.

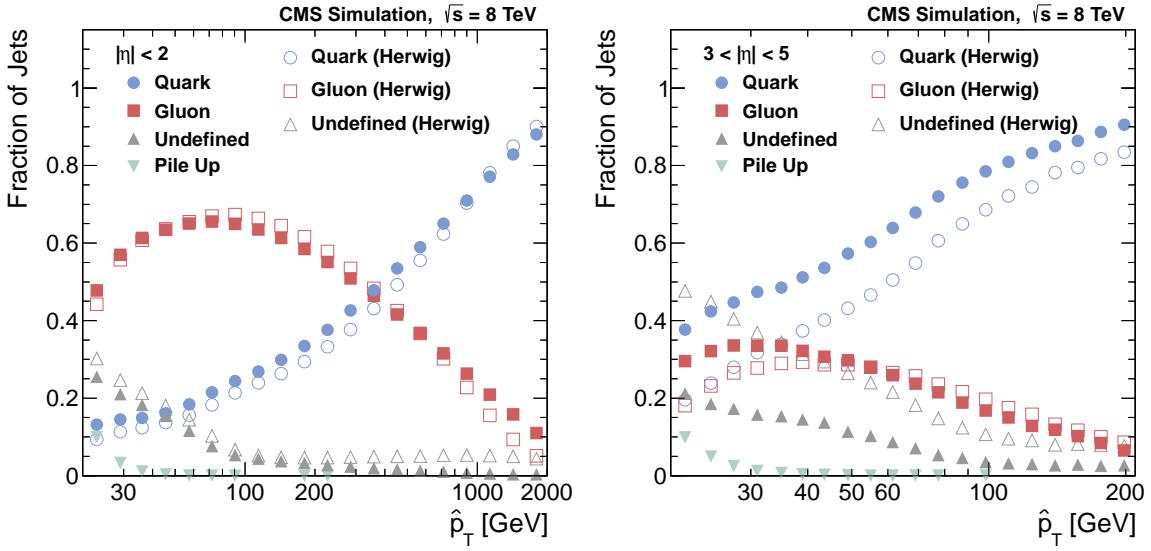


Figure 6.17: Jet flavour composition in dijet events, for central (left) and forward (right) jets, as a function of the event hard scattering transverse momentum (\hat{p}_T). Jets originating from light quark partons are shown with blue circles, gluon jets with red squares, undefined jets with grey triangles and pile up jets with yellow triangles. Solid markers indicate events simulated in Pythia 6, whereas hollow ones with Herwig++. As can be seen, whereas for $|\eta| < 2$ the majority of jets originate from gluons (at least up to 300-400 GeV), in the forward quark jets are dominant. The dijet sample is therefore a useful gluon enriched sample for central jets, and offers an important cross check to Z+jets for forward quark jets, as it has higher statistics.

6.2 Validation on Dijet Events

The flavour composition of the dijet dataset is somewhat more complex than the Z+jet case, as can be seen in Figure 6.17, where the flavor composition of jets in dijet events is shown as a function of the hard scattering transverse momentum (\hat{p}_T) of the event. The left plot shows central jets, the right forward jets. Jets originating from light quark partons are shown with blue circles, gluon jets with red squares, undefined jets with grey triangles and pile up jets with yellow triangles. Solid markers indicate events simulated in Pythia 6, whereas hollow ones with Herwig++. As can be seen, whereas for $|\eta| < 2$ the majority of jets originate from gluons (at least up to 300-400 GeV), in the forward quark jets are dominant. The dijet sample is therefore a useful gluon enriched sample for central jets, and offers an important cross check to Z+jets for forward quark jets, as it has higher statistics.

Dijet events are selected in data by running on the full 2012 dataset analyzing events recorded with the single jet trigger HLT_PFJet40. They are compared to dijet events simulated in Pythia, from the ‘Flat’ Summer12 dataset reconstructed in CMSSW_5_3_X:

/QCD_Pt-15to3000_TuneZ2star_Flat_8TeV_pythia6.

Further comparisons have been carried out on the corresponding Herwig++ sample:

/QCD_Pt-15to3000_TuneEE3C_Flat_8TeV_herwigpp.

The event selection further requires:

- two jets with transverse momentum greater than 40 and 20 GeV, respectively;
- the two leading jets are required to be back-to-back in the transverse plane by requiring their azimuthal difference to be greater than 2.5 rad;
- the third jet in the event is required to have a transverse momentum inferior to 30% than the average p_T of the first two jets.

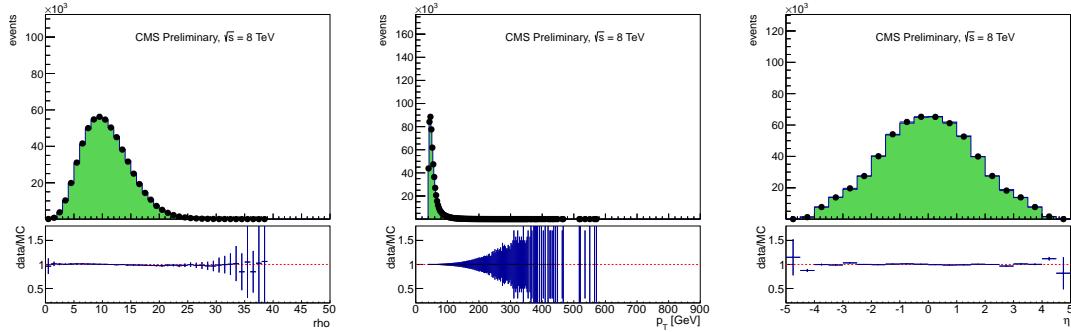


Figure 6.18: Data-MC comparisons for the particle-flow energy density ρ (left), the leading jet transverse momentum (center) and its pseudorapidity (right). The reweighting procedure is defined by the ρ distribution, so it is expected to find perfect agreement there.

494 Only the leading jet is used to fill the histograms (each in their own p_T bin), and it is required to
 495 pass the usual jet ID, pile up ID and anti-btagging requirements. In order to match the pile up
 496 distribution observed in the data, the Monte Carlo is reweighted based on the *a posteriori*
 497 particle flow energy density (ρ) distribution, in the same manner as is done in the Z+jets validation.
 498 As the simulated samples have a flat transverse momentum spectrum, though, an additional
 499 reweighting is necessary: we reweight simulated events so that the leading jet transverse mo-
 500 mentum and pseudorapidity spectra match the ones observed in the data. This is done with a
 501 2-dimensional weight matrix, so that the fact that transverse momentum and pseudorapidity
 502 are not completely independent is taken into account. The results of this reweighting is shown
 503 in Figure 6.18: the right plot shows the ρ distribution, the center plot shows the leading jet
 504 transverse momentum and the right plot its pseudorapidity.

505 Figures 6.19, 6.20, 6.21 and 6.22 show, respectively, the data-MC comparisons of the four input
 506 variables (total multiplicity, $p_T D$, Axis₂ and Axis₁), in four representative p_T bins, for jets re-
 507 constructed in the center of the detector. The corresponding plots for the forward are shown
 508 in Figures 6.23, 6.25, 6.26 and 6.27. Similarly to the validation performed on Z+jet events, also
 509 in dijet events we observe a significant discrepancy between the data and the simulation for
 510 what concerns the candidate multiplicity forward jets. We apply the same correction to data
 511 multiplicities before feeding them to the discriminators, i.e. we subtract one (unity) from the
 512 multiplicity. The data-MC comparisons for the forward multiplicity after this correction are
 513 shown in Figure 6.24.

514 The corresponding data-MC comparisons for the output discriminants are shown in Figures 6.28
 515 for the likelihood discriminant and in Figure 6.31 for the MLP. The comparisons for the forward
 516 discriminators are shown respectively in Figures 6.29 and 6.32 before the multiplicity correc-
 517 tion, and in Figures 6.30 and 6.33 after the multiplicity correction. As can be seen significantly
 518 better agreement is observed after the application of the correction.

519 As was seen in Section 5, Herwig++ predicts gluon jets to be more similar to quark jets than
 520 Pythia. This will have significant implications in the dijet sample, which has a large contamina-
 521 tion of gluon jets. We have therefore performed a second validation on Herwig++ dijet events.
 522 While the single-variable data/MC comparisons are moved to Appendix ??, we show the data-
 523 Herwig++ comparisons for central jets in Figures 6.34 and 6.36, respectively for the likelihood
 524 and the MLP discriminants. Similar comparisons for forward jets are shown in Figures 6.35 and
 525 ???. For forward jets the same correction to the jet candidate multiplicity that was used in the
 526 Pythia comparisons is used when comparing to Herwig++. It must be noted that whereas for
 527 the MLP tagger the same weight files have been used both for Pythia and Herwig++ samples,

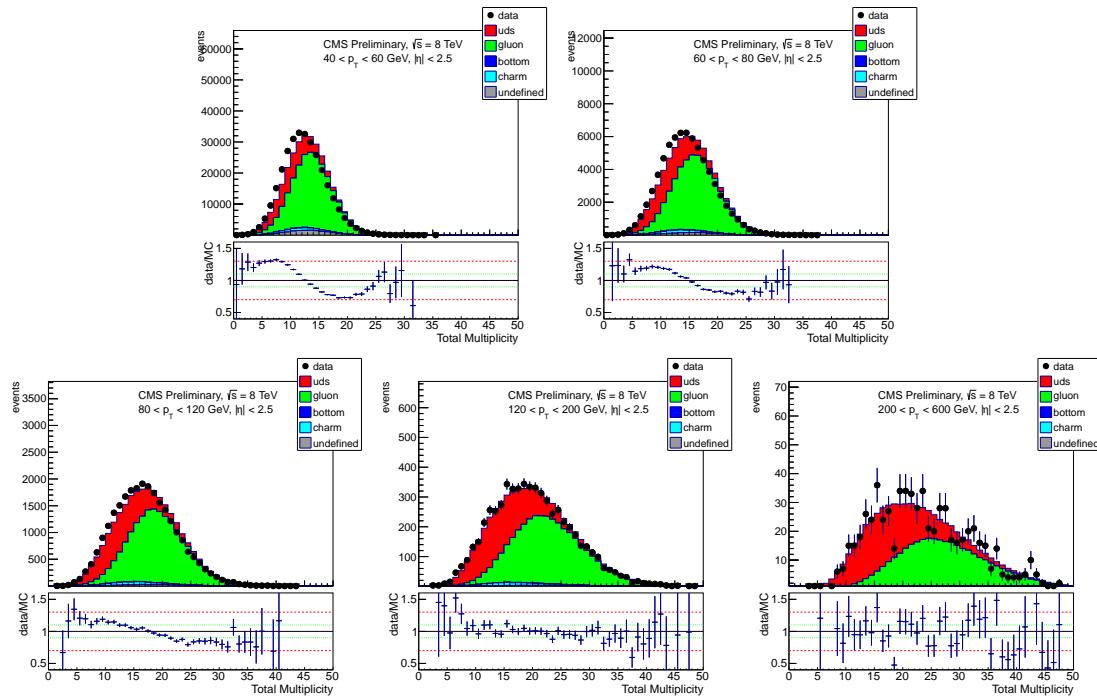


Figure 6.19: Data-MC comparisons for the total PFCandidate multiplicity in central jets in dijet events, for five p_T bins: 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

528 the likelihood has been re-defined for the Herwig++ comparison, by constructing the single
 529 variable PDFs on the Herwig++ sample.

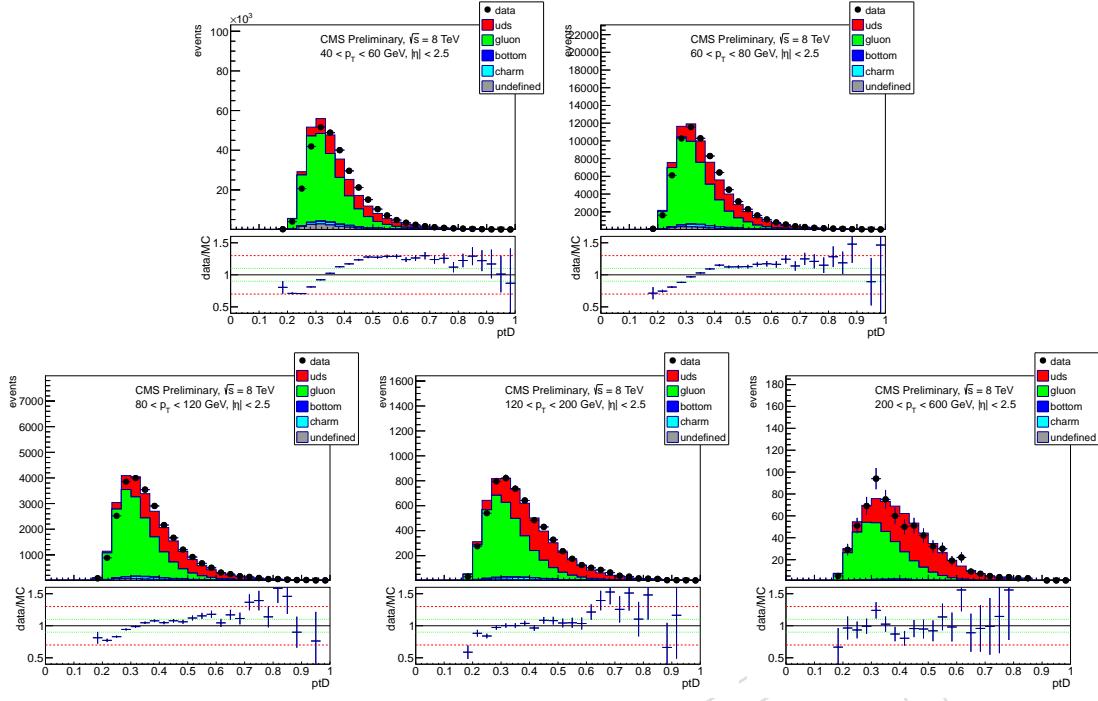


Figure 6.20: Data-MC comparisons for $p_T D$ in central jets in dijet events, for five p_T bins: 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

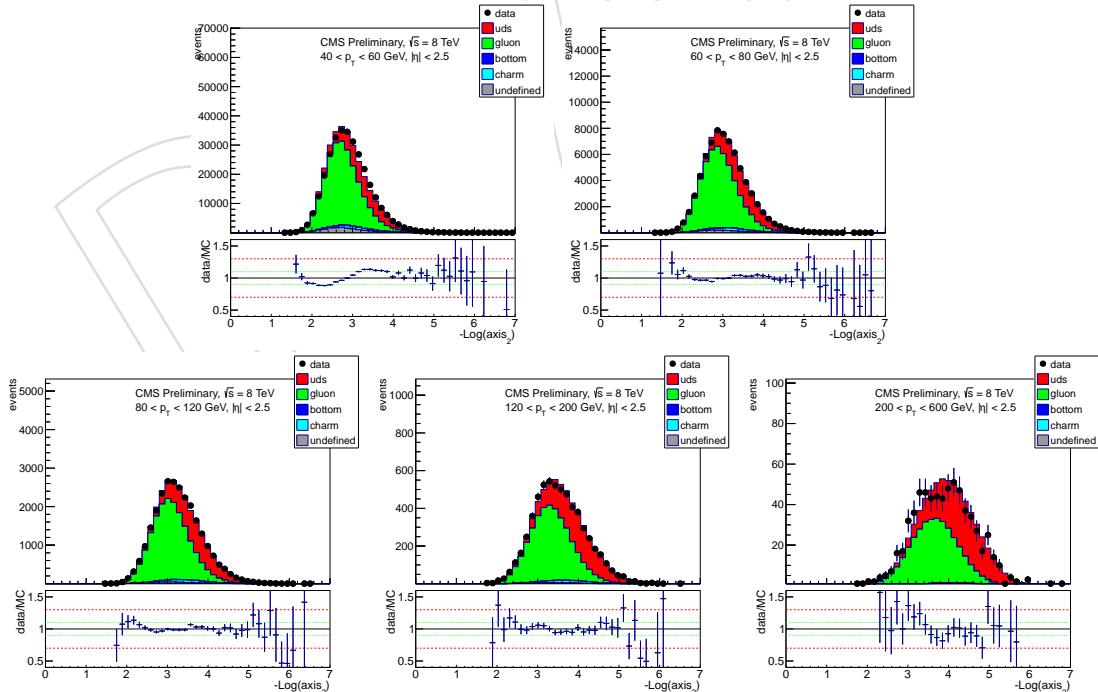


Figure 6.21: Data-MC comparisons for Axis₂ in central jets in dijet events, for five p_T bins: 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

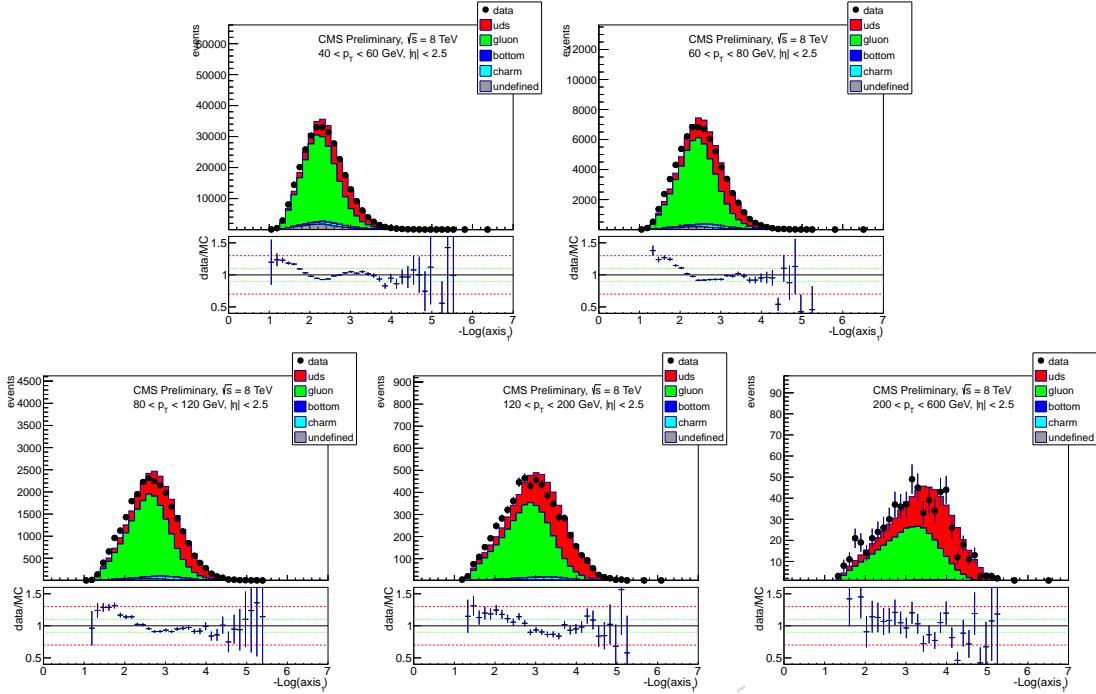


Figure 6.22: Data-MC comparisons for Axis₁ in central jets in dijet events, for five p_T bins: 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

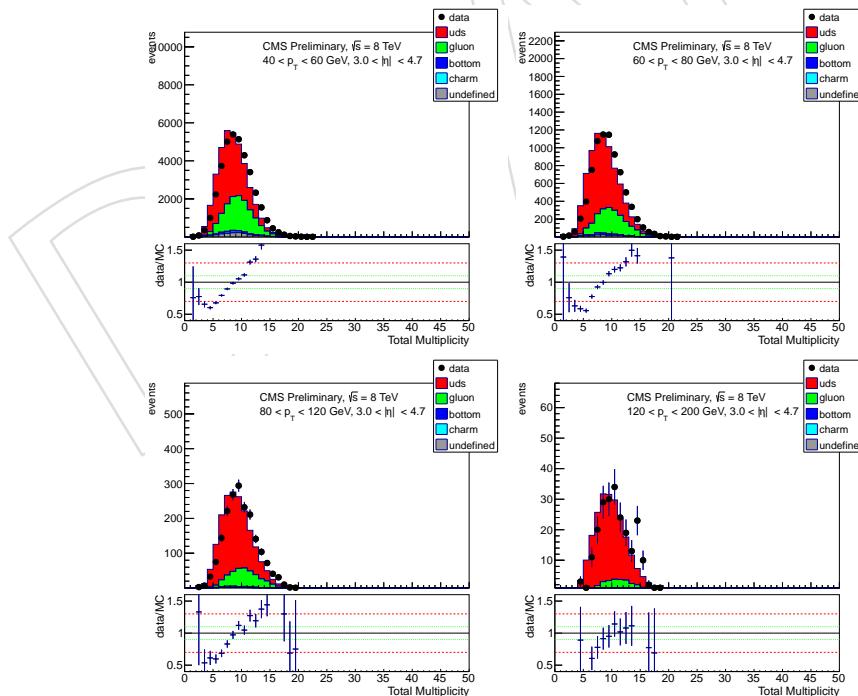


Figure 6.23: Data-MC comparisons for the total PFCandidate multiplicity in forward jets in dijet events **before the correction**, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

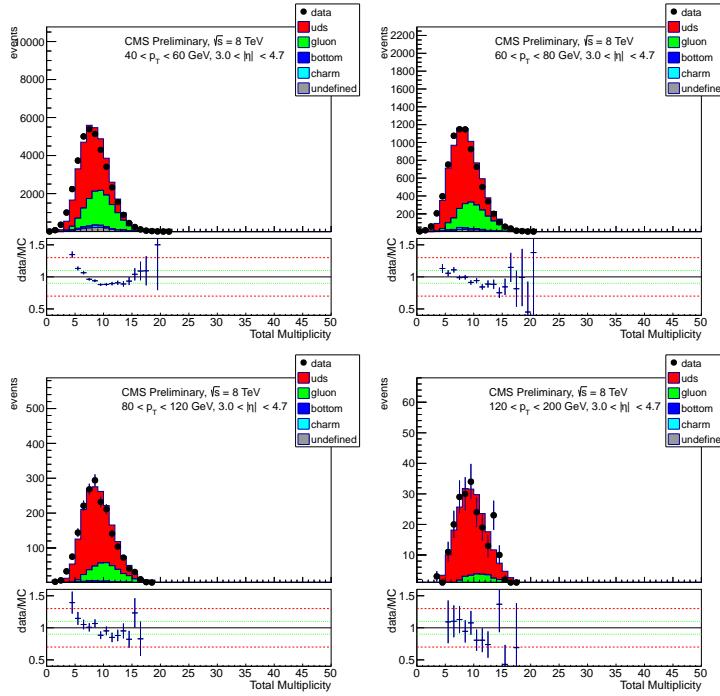


Figure 6.24: Data-MC comparisons for the total PFCandidate multiplicity in forward jets in dijet events **after the correction**, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

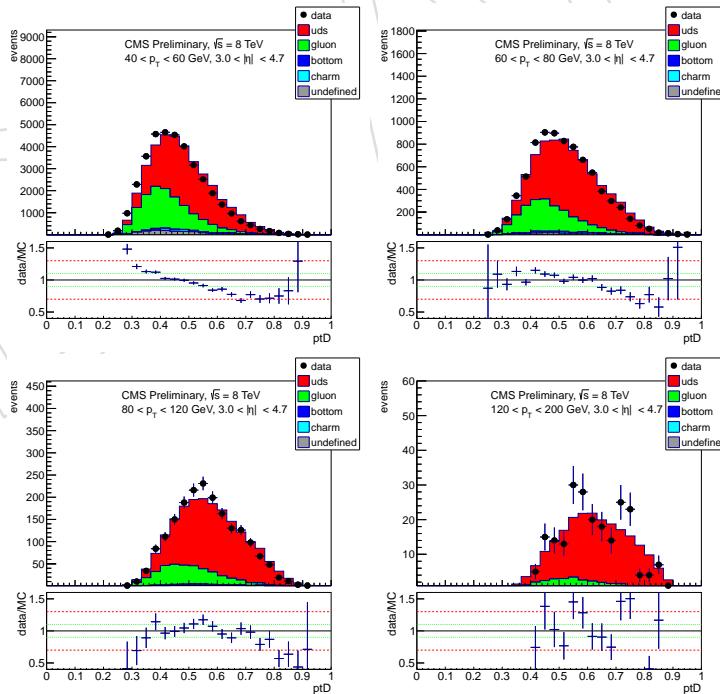


Figure 6.25: Data-MC comparisons for $p_T D$ in forward jets in dijet events, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

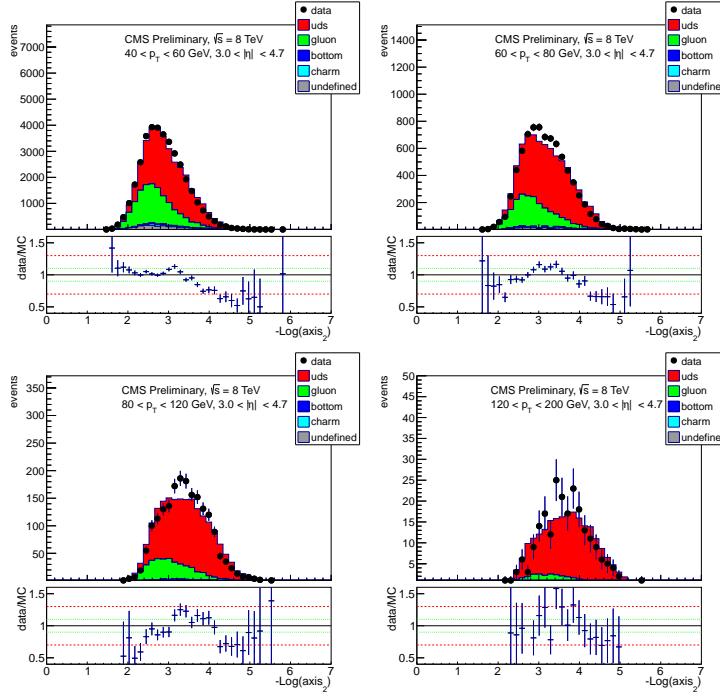


Figure 6.26: Data-MC comparisons for Axis₂ in forward jets in dijet events, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

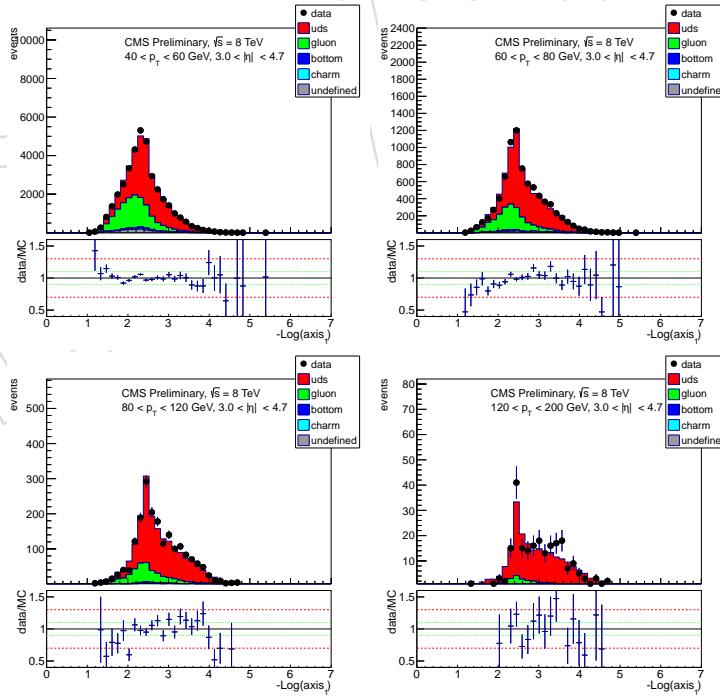


Figure 6.27: Data-MC comparisons for Axis₁ in forward jets in dijet events, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

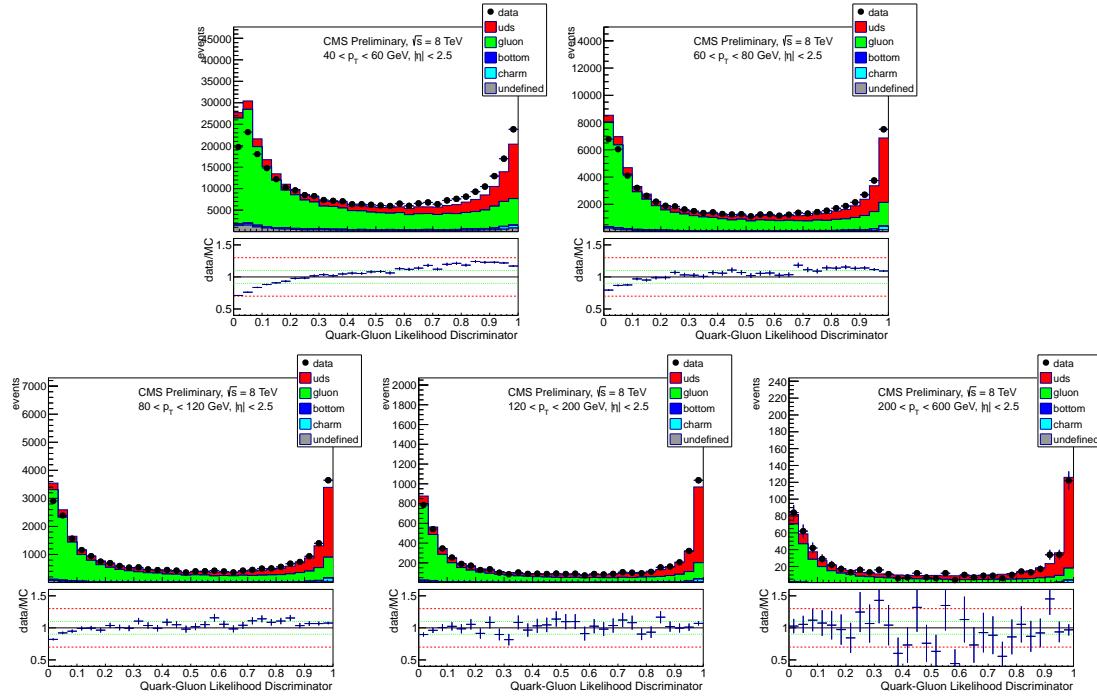


Figure 6.28: Data-MC comparisons for the quark-gluon likelihood discriminant in central jets, for five p_T bins: 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

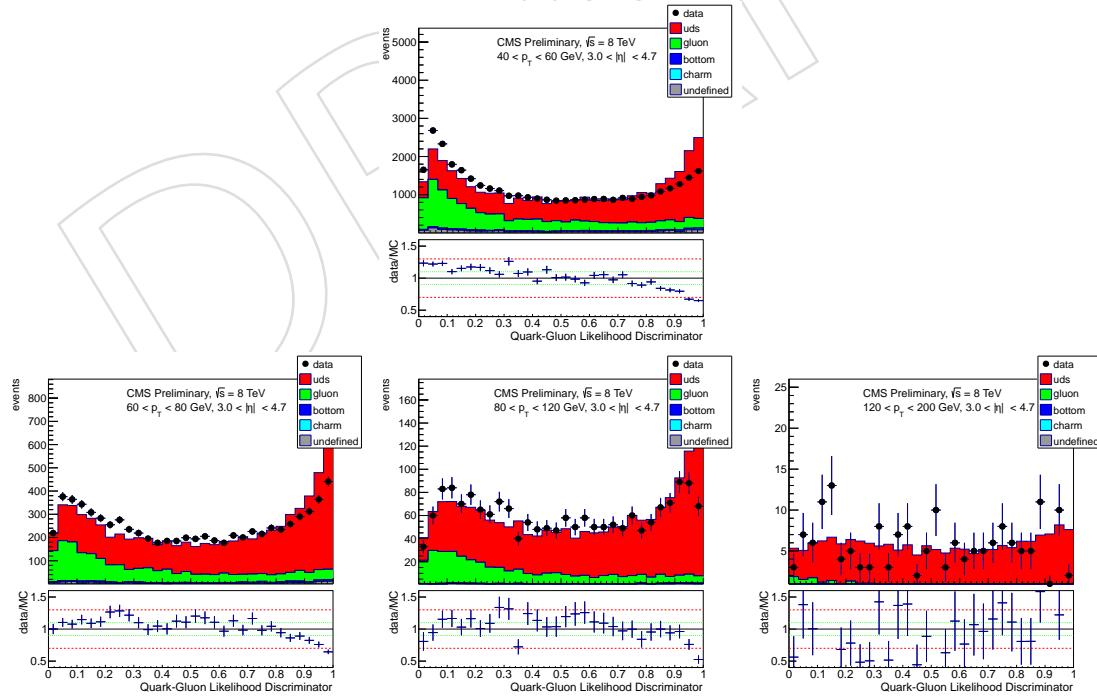


Figure 6.29: Data-MC comparisons for the quark-gluon likelihood discriminant in forward jets in dijet events **before the multiplicity correction**, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

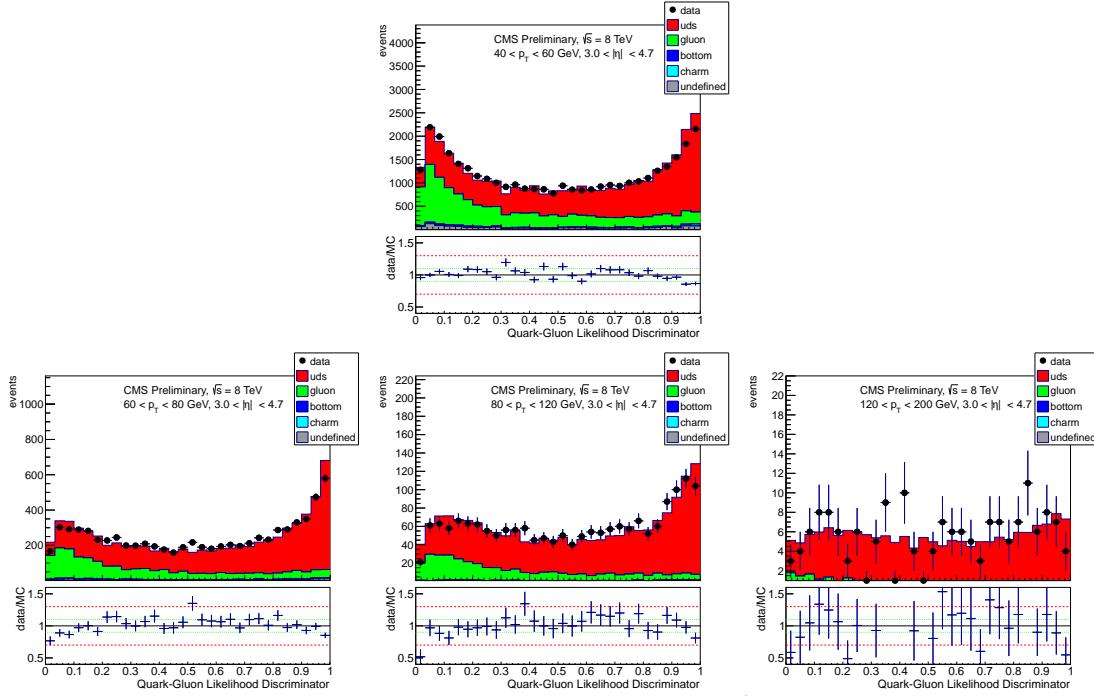


Figure 6.30: Data-MC comparisons for the quark-gluon likelihood discriminant in forward jets in dijet events **after the multiplicity correction**, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

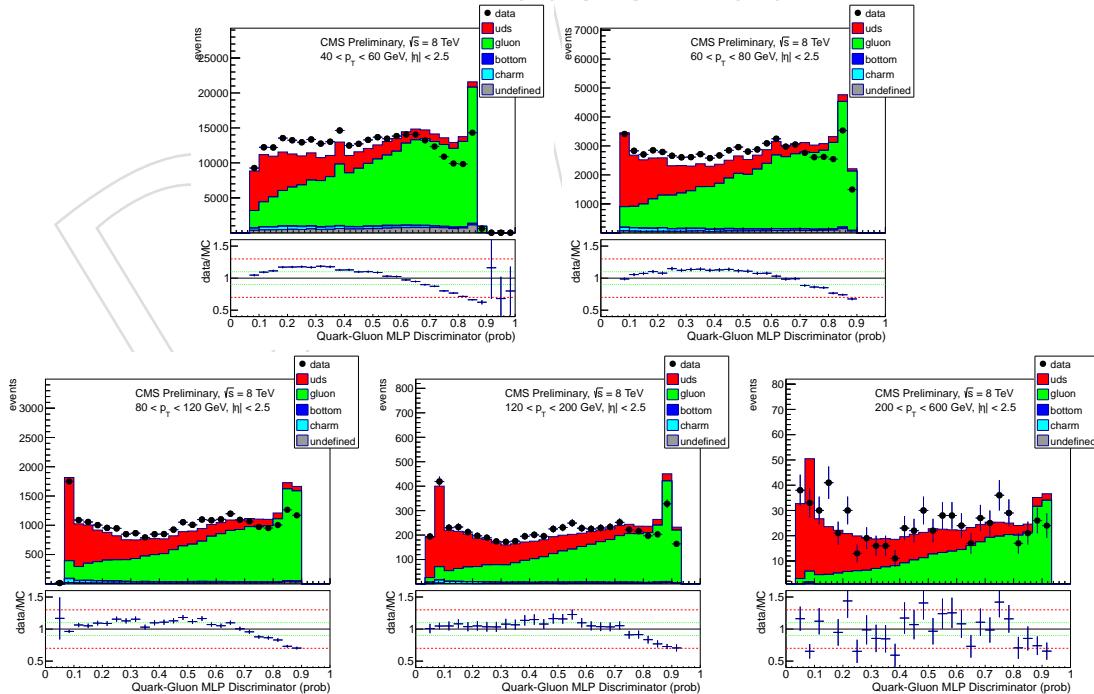


Figure 6.31: Data-MC comparisons for the quark-gluon MLP discriminant in central jets in dijet events, for five p_T bins: 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

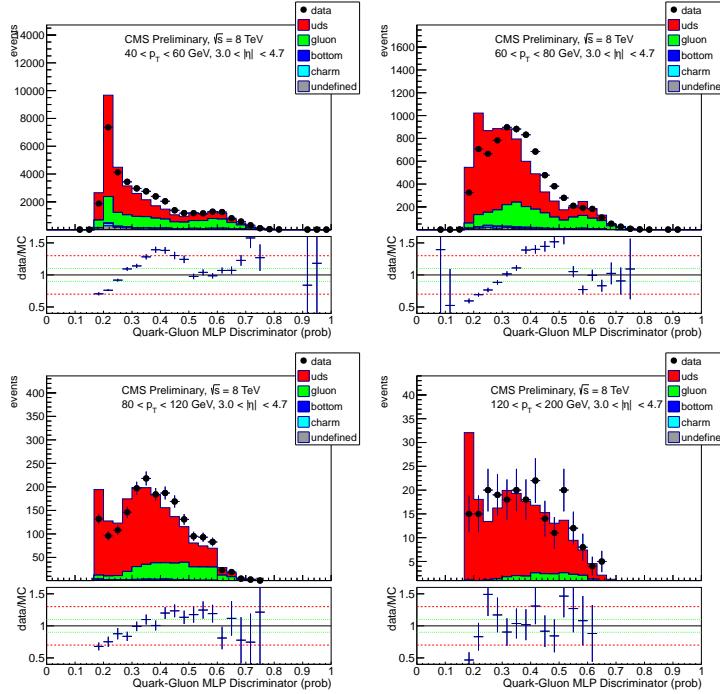


Figure 6.32: Data-MC comparisons for the quark-gluon MLP discriminant in forward jets in dijet events **before the multiplicity correction**, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

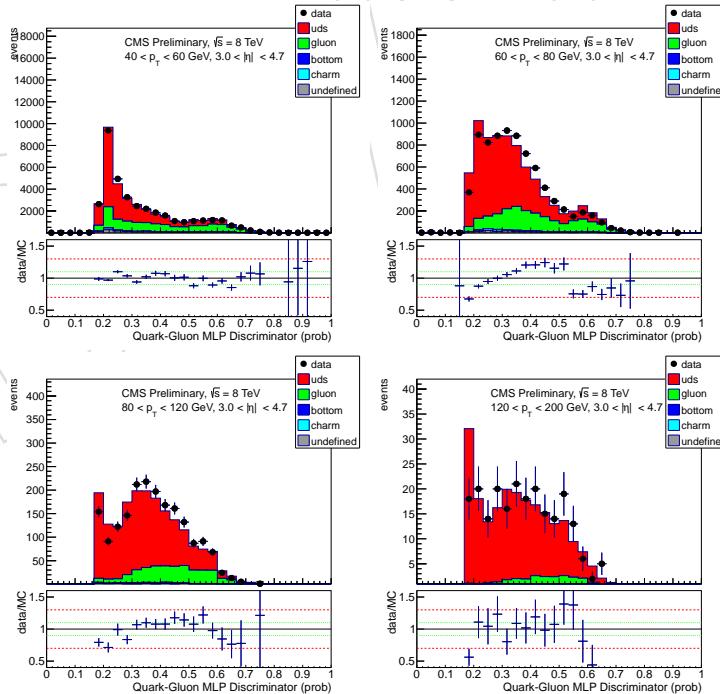


Figure 6.33: Data-MC comparisons for the quark-gluon MLP discriminant in forward jets in dijet events **after the multiplicity correction**, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

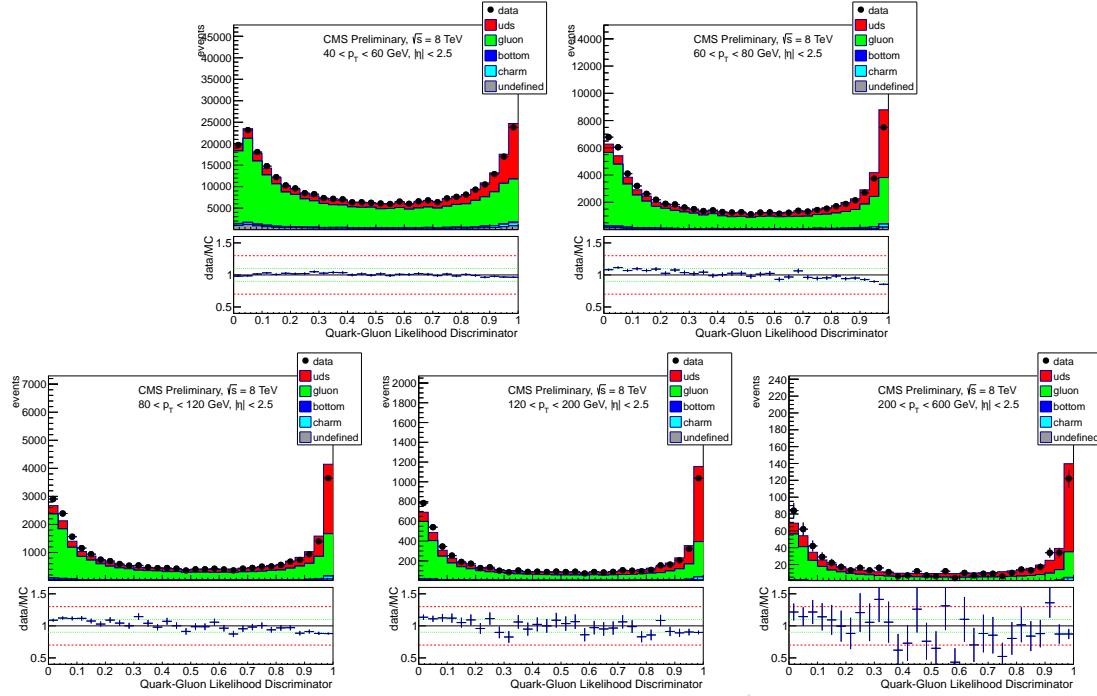


Figure 6.34: Data-Herwig++ comparisons for the quark-gluon likelihood discriminant in central jets in dijet events, for five p_T bins: 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

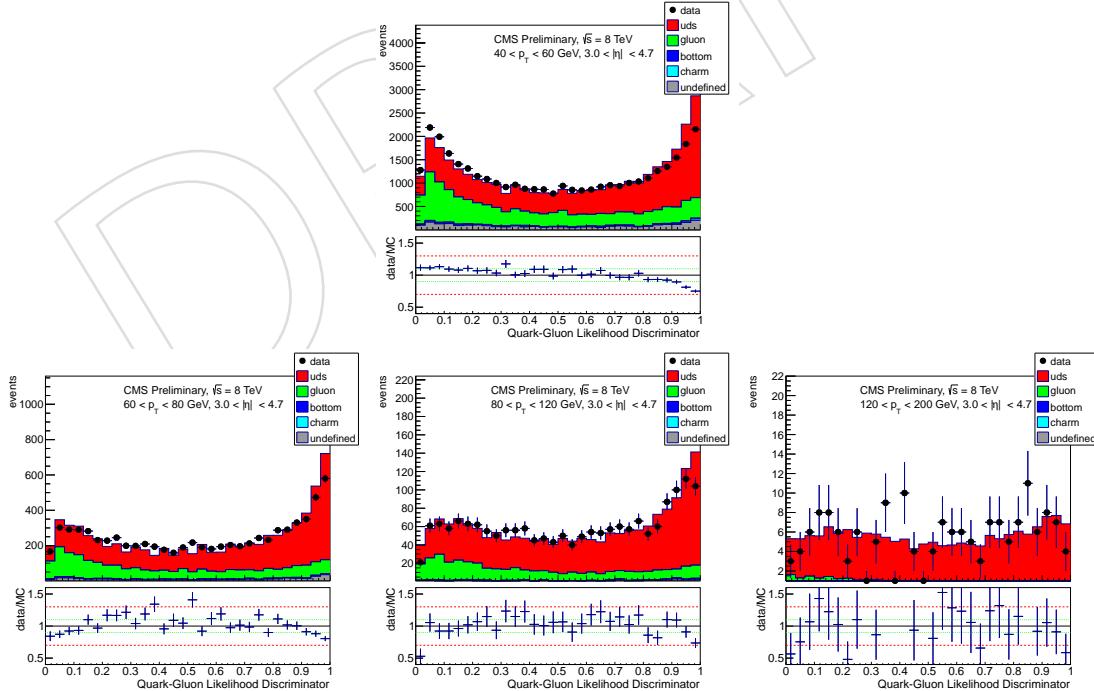


Figure 6.35: Data-Herwig++ comparisons for the quark-gluon likelihood discriminant in forward jets in dijet events **after the multiplicity correction**, for four p_T bins: 40-60, 60-80, 80-120 and 120-200 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

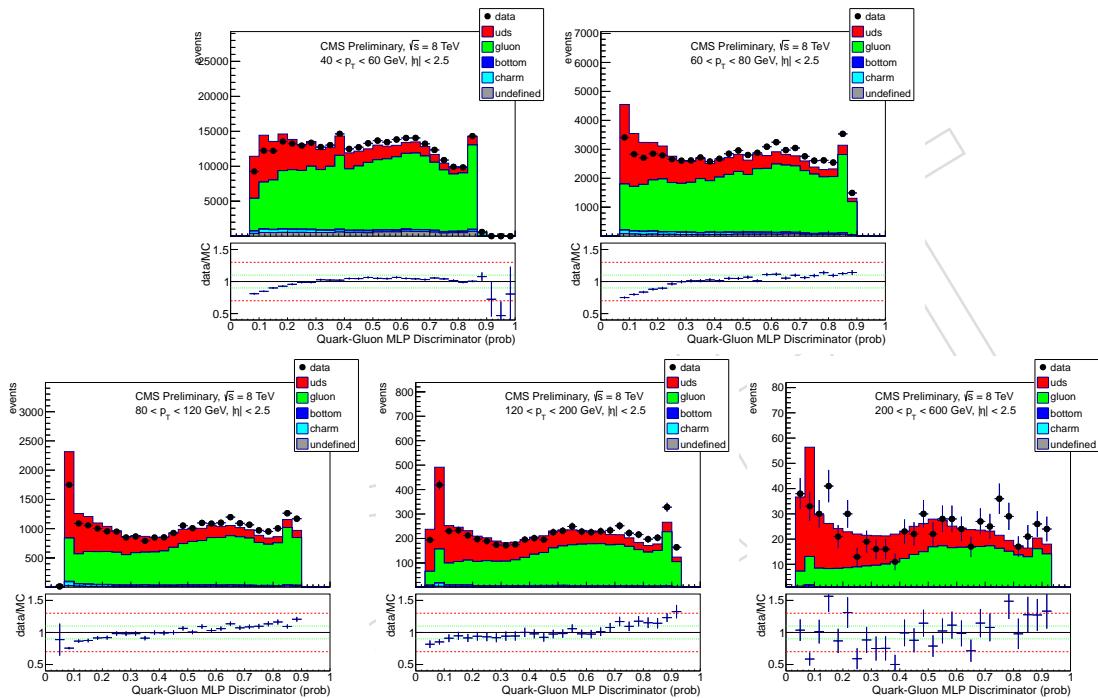


Figure 6.36: Data-Herwig++ comparisons for the quark-gluon MLP discriminant in central jets in dijet events, for five p_T bins: 40-60, 60-80, 80-120, 120-200 and 200-600 GeV. The different jet flavours of simulated jets are shown: light quarks (red), gluon (green), bottom (blue), charm (cyan), and undefined (grey).

530 7 Systematics evaluation

531 7.1 Usage: data vs data

532 The comparison between data and data it is far more less sensitive to any model used to derive
 533 the selections. The main goal in this case is to show that the analysis is not particularly sensitive
 534 to where the cut is placed in the discriminator (or in any other) and it is left to the analysts: a
 535 stability under a variation of the cut value of some percent should be expected.

536 7.2 Usage: data vs mc

537 This is the first and most delicate systematics that it is need to evaluate, and that allow to make
 538 a fair comparison between the results obtained with data and the one with Monte Carlo. A
 539 general comparison between data and Monte Carlo can show discrepancies in the discriminator
 540 for many reasons. Among them it worth to notice the hadronization model, the detector model,
 541 the pile-up model, the hard scattering, the pdf ...

542 We are therefore interest to assset systematics in order to cover the events not due to physics.
 543 In order to reduce the impact due to physics, the validation samples used in the validation are
 544 used to extract and verify the systematic error.

545 A genaral usage of the discriminator, being ething the likelihood or the MLP, implies one of the
 546 following two scenarios:

- 547 • use with cuts
- 548 • use with multivariate analysis

549 In order to allow these two shenario a systematic method which takes into account the whole
 550 shape is developed.

551 Among a set of function that maps the interval $[0, 1]$ into itself, the following were considered
 552 because have a reasonable number of parameters (≤ 2) that can be used to shift the population
 553 through the center or move the population to one of the two extremities.

The functions studied are:

$$f(x, a, b) = \arctan \left(a \tan \left(x\pi - \frac{\pi}{2} \right) + b \right) / \pi + \frac{1}{2} \quad (7.1)$$

$$g(x, a, b) = \tanh \left(a \operatorname{arctanh} (2x - 1) + b \right) / 2 + \frac{1}{2} \quad (7.2)$$

554 The behaviour of these two function is shown in Fig. 7.1. It worths to be noticed the difference
 555 behaviour of the two set of functions around the extreme points of the interval. This difference
 556 allows bigger changes for the function g (Eqn. 7.2) at the extreme with less modification in the
 557 middle range. For this reason it perform better on the likelihood. For this reason this function
 558 was chosen.

559 7.2.1 Fit

560 The Z+Jets data and Monte Carlo are used in order to perform the systematics evaluation. The
 561 Z +Jet selection is performed very tight in order to have the less as possible contriubition from
 562 background or pile-up with a β^* cut on the jets and a balance in p_T and ϕ between the leading
 563 jet and the Z, and a jet veto on the second jet.

564 The minimization is perform w.r.t the χ^2 : the Monte Carlo entries of the discriminator are var-
 565 ied according to the chosen function (Eqn. 7.2) and the two parameters are therefore evaluated

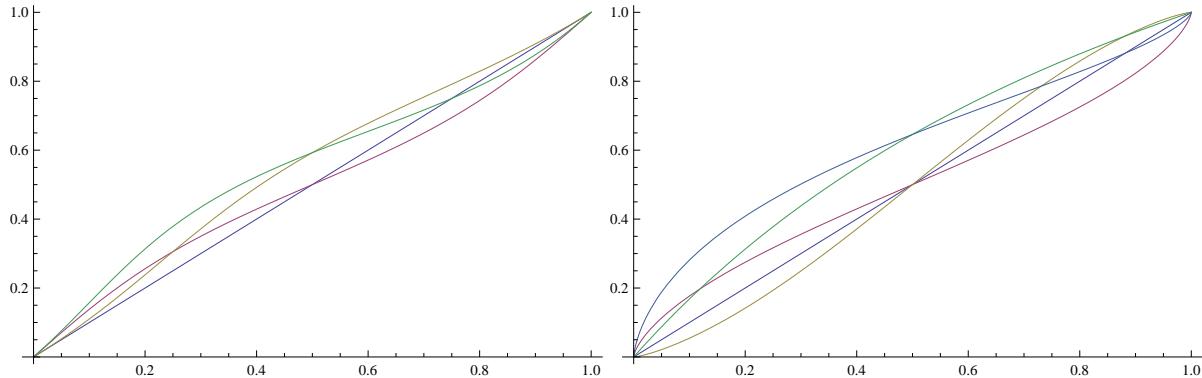


Figure 7.1: Representation of the functions in Equation 7.1 on the left, and in Equation 7.2 on the right for different values of the parameters.

- 566 in large bin of the jet $p_{T,\eta}$, and ρ in order to exploit possible differences in any of these direc-
 567 tions.

For the MLP an additional step is added to the minimization process: on the given bin and selection the minimum and maximum of the tagger are recorded from the Monte Carlo (t_{\min}, t_{\max}) and a linear transformation (l) is appended to transform $[t_{\min}, t_{\max}]$ in $[0, 1]$ and anti-transformed after the function is computed, i.e., the following modified functions are used:

$$\begin{aligned}\tilde{f}(x, a, b, t_{\min}, t_{\max}) &= l^{-1}(x, t_{\min}, t_{\max}) \circ f(x, a, b) \circ l(x, t_{\min}, t_{\max}) \\ &= f\left(\frac{x - t_{\min}}{t_{\max} - t_{\min}}, a, b\right) \cdot (t_{\max} - t_{\min}) + t_{\min}\end{aligned}\quad (7.3)$$

$$\begin{aligned}\tilde{g}(x, a, b, t_{\min}, t_{\max}) &= l^{-1}(x, t_{\min}, t_{\max}) \circ g(x, a, b) \circ l(x, t_{\min}, t_{\max}) \\ &= g\left(\frac{x - t_{\min}}{t_{\max} - t_{\min}}, a, b\right) \cdot (t_{\max} - t_{\min}) + t_{\min}\end{aligned}\quad (7.4)$$

- 568 In this way the minimization is still on two parameters (a and b), while $t_{\min/\max}$ are evaluated
 569 in advance. If a value will follow out of that interval in the future it will be mapped in the
 570 systematic to the closest point in that interval (i.e., aut t_{\min} aut t_{\max}).
 571 The parameters, result of the minimization process, are reported in Figure 7.2 for the different
 572 bins. No trends can be assessed from these results.
 573 On Figures 7.3, 7.4 are shown the data with the Monte Carlo before and after the shift for the
 574 likelihood tagger and the MLP tagger respectively.

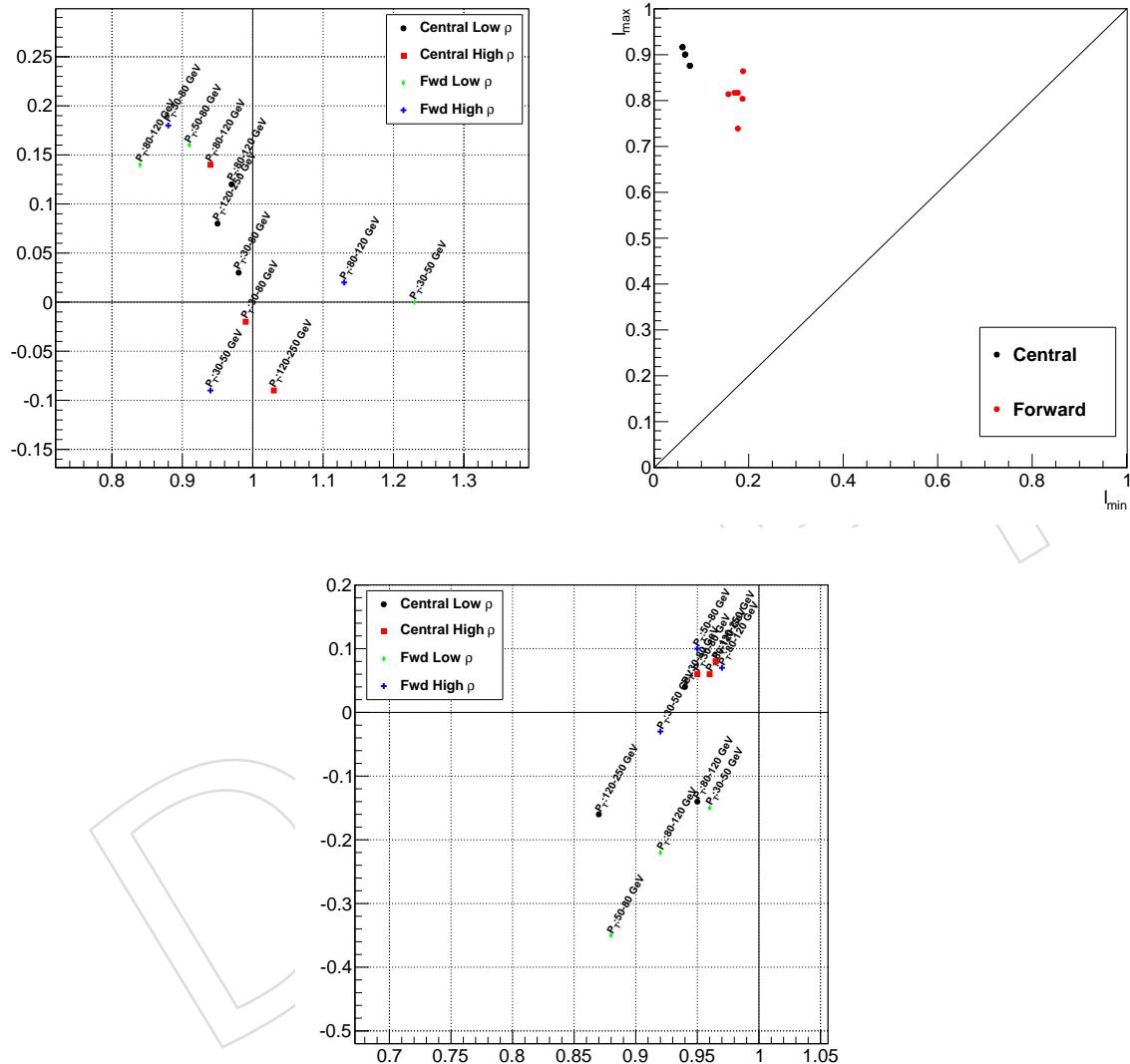


Figure 7.2: Parameters fitted for the chosen function (Eqn. 7.2) with respect to the different bins: on the top the MLP parameters (with the ranges on the right), on the bottom the Likelihood ones.

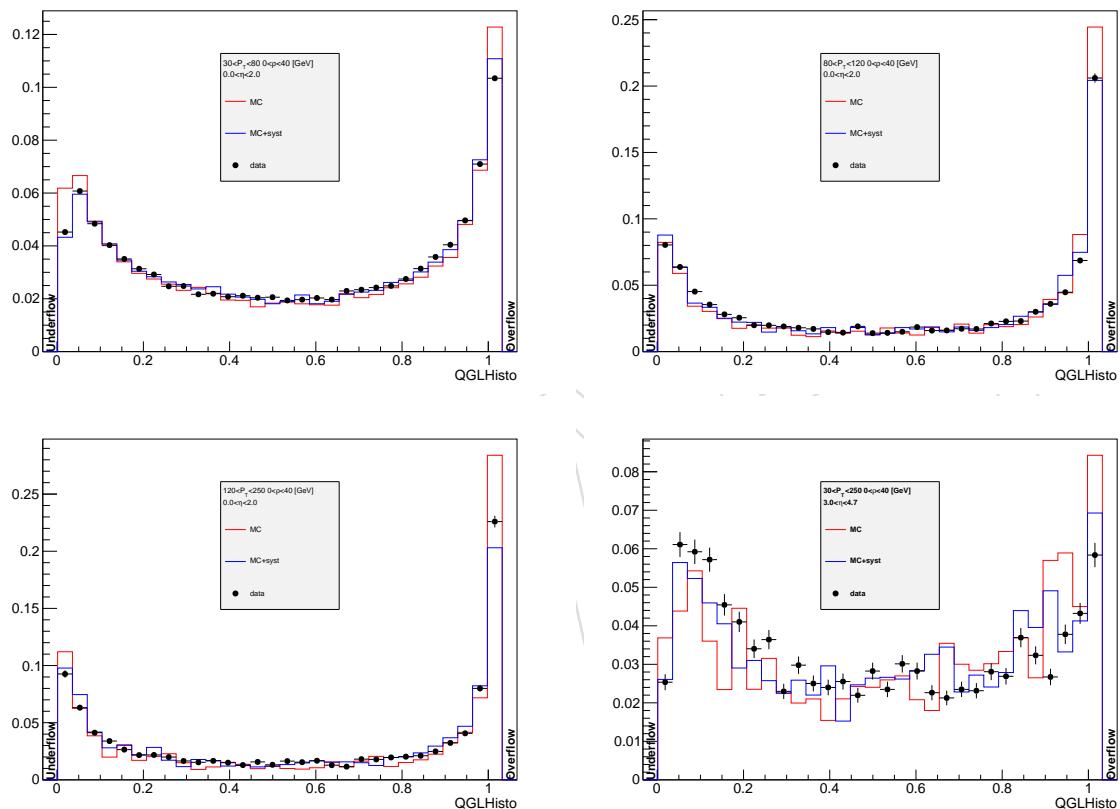


Figure 7.3: Data/MC agreement before and after fit for systematics evaluation in bins of jet p_T for the likelihood tagger.

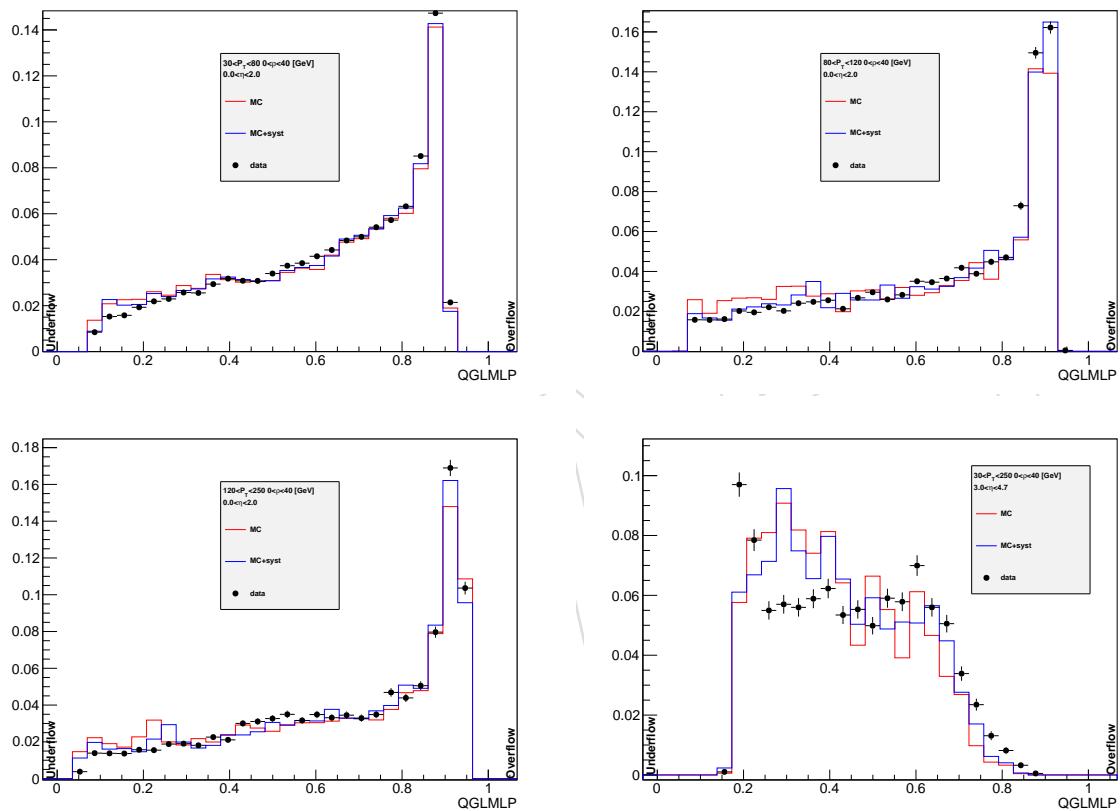


Figure 7.4: Data/MC agreement before and after fit for systematics evaluation in bins of jet p_T for the MLP Tagger.

575 8 Conclusions

576 A new tool to separate CMS particle flow jets originating from the showering and hadronization
577 of final states gluons or light (uds) quarks has been presented. The new discriminator
578 makes use of internal jet composition properties and exploits known and well established dif-
579 ferences in the showering of quarks and gluons. With respect to the previously developed dis-
580 criminator, the presented ones achieve improved performances and extend the discrimination
581 power to the CMS forward region, up to $|\eta| = 5$. The new discriminator has been validated
582 using 2012 8 TeV collision data, both with dijet and Z+jet event samples. Finally a recipe to eval-
583 uate the systematic uncertainties associated to the use of the new discriminator is given, based
584 on the observed data versus Monte Carlo differences in the validation samples. The recipe
585 makes use of a two-parameter analytical transformation function that distorts the Monte Carlo
586 output distributions in order to reproduce better the observed data outputs: it can therefore be
587 used to determine uncertainties on Monte Carlo efficiencies and rejections both for a cut-based
588 selection or when using the whole discriminator output for another multivariate analysis step.



589 References

- 590 [1] OPAL Collaboration Collaboration, "A Study of differences between quark and gluon jets
591 using vertex tagging of quark jets", *Z.Phys.* **C58** (1993) 387–404,
592 doi:10.1007/BF01557696.
- 593 [2] OPAL Collaboration Collaboration, "A Model independent measurement of quark and
594 gluon jet properties and differences", *Z.Phys.* **C68** (1995) 179–202.
- 595 [3] DELPHI Collaboration Collaboration, "Energy dependence of the differences between
596 the quark and gluon jet fragmentation", *Z.Phys.* **C70** (1996) 179–196,
597 doi:10.1007/s002880050095.
- 598 [4] ALEPH Collaboration Collaboration, "Quark and gluon jet properties in symmetric three
599 jet events", *Phys.Lett.* **B384** (1996) 353–364, doi:10.1016/0370-2693(96)00849-0.
- 600 [5] N. Saoulidou, "Particle Flow Jet Identification Criteria", CMS Note 2010/03, (2010).
- 601 [6] A. C. Marini and F. Pandolfi, "Quark-Gluon Discrimination in CMS", CMS Analysis Note
602 AN-2012/074, (2012).
- 603 [7] CMS Collaboration, "Search for a Higgs boson in the decay channel
604 $H \rightarrow ZZ(*) \rightarrow q\bar{q}\ell^-\ell^+$ in pp collisions at $\sqrt{s} = 7$ TeV", *Journal of High Energy Physics*
605 **2012** (2012) 1–36, doi:10.1007/JHEP04(2012)036.
- 606 [8] A. Hoecker et al., "TMVA: Toolkit for Multivariate Data Analysis", *PoS ACAT* (2007)
607 040, arXiv:physics/0703039.
- 608 [9] E. D. Malaza, "Multiplicity distributions in quark and gluon jets", *Zeitschrift fur Physik C
609 Particles and Fields* **31** (1986) 143–150. 10.1007/BF01559605.
- 610 [10] A. C. Marini et al., "Quark-Gluon Jet Discrimination through Particle Flow Jet Structure",
611 CMS Note 2011/215, (2011).
- 612 [11] J. Gallicchio and M. D. Schwartz, "Quark and Gluon Jet Substructure",
613 arXiv:1211.7038.
- 614 [12] A. C. Marini and F. Pandolfi, "Quark-Gluon Discrimination", CMS Note 2012/74, (2012).
- 615 [13] G. P. Salam, "Towards Jetography", *Eur.Phys.J.* **C67** (2010) 637–686,
616 doi:10.1140/epjc/s10052-010-1314-6, arXiv:0906.1833.

617 A Bjets

618 It was already observed that b-jets have similar shape as gluon jets **FIXME**: Ref. We find similar
619 results, in particular we observed that at low p_T b-jets are closer to gluon rather than to light
620 quark jets in all variables studied, while at high p_T they come closer and closer to the quark
621 jets. The explanation to this behaviour is laid in two opposite facts: from being massive the
622 b-partons are less coupled to the QCD field as measured by **FIXME**: Ref. However, the chain
623 of weak decays which will form the b-jets as we observe them,i.e., after the decay of the heavy
624 B and D mesons produced during hadronization, contribute to enlarge the jet and enrich it
625 in multiplicity. This behaviour is shown in Figure A.1 where the Charged Multiplicity (all)
626 compared to the Charged Multiplicity (QC) show that the extra-number of particles comes from
627 secondary vertexes (and so do not match QC). Also the other variables will be affected by these
628 weaks decays. The c-jets should have the same but much more less pronounced behaviour. In
629 particular at very high p_T c and light (uds) quark are undistinguishable (in our variables).

630 The goal of this product is not to identify b-jets (or c-jets) rather than to be able to discriminate
631 between light quark jets and gluon jets. Special efforts are already developed in the Collabora-
632 tion for this kind of jets, and therefore these tools should be used. In particular, since this tool
633 can be used in addition to b-jets tagger, an explicit b-jets vetos was required in the training.

DRAFT

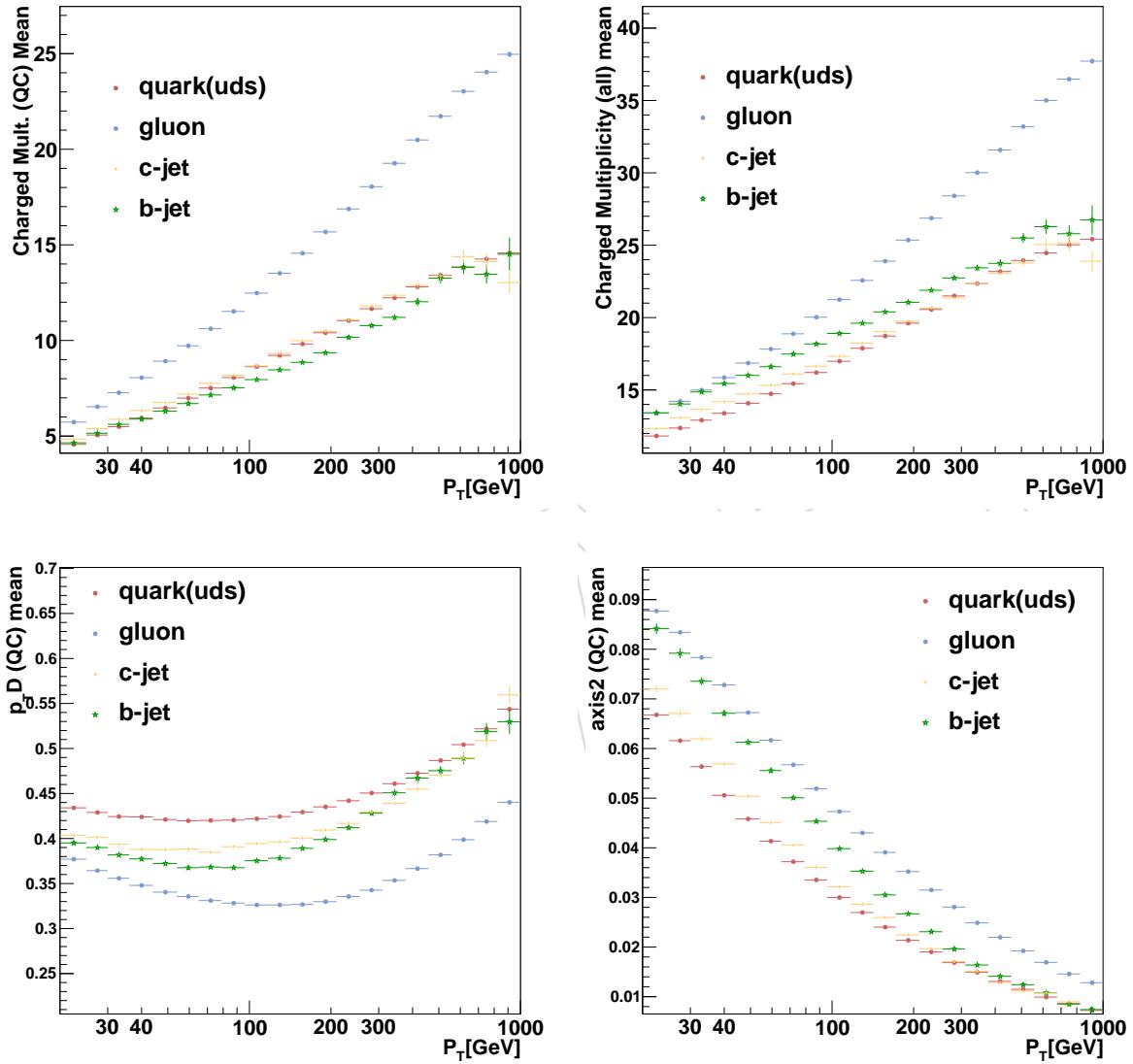


Figure A.1: In the figures is shown the mean of Charged Multiplicity (QC), the Charge Multiplicity, $p_T D(QC)$ and $\text{axis2}(QC)$ vs p_T of the jet.