

# Top quark physics at the LHC with the CMS detector



Sergey Senkin  
School of Physics  
University of Bristol

*A dissertation submitted to the University of Bristol  
in accordance with the requirements of the degree of  
Doctor of Philosophy in the Faculty of Science*

March 16, 2014

# Contents

<b>Contents</b>	<b>i</b>
<b>1 Theoretical Background</b>	<b>1</b>
1.1 The Standard Model of Particle Physics . . . . .	1
1.1.1 Electroweak theory . . . . .	1
1.1.2 Electroweak symmetry breaking . . . . .	1
1.2 Top Quark Physics within the Standard Model . . . . .	1
1.2.1 Top quark production at the LHC . . . . .	1
1.2.2 Top quark decay . . . . .	1
1.2.3 Background processes . . . . .	1
1.2.4 Top quark mass . . . . .	1
1.2.5 Importance of top quark physics . . . . .	1
1.3 Physics Beyond The Standard Model . . . . .	1
1.3.1 Supersymmetry . . . . .	1
1.3.2 Extra dimensions . . . . .	1
<b>2 The LHC and the CMS detector</b>	<b>3</b>
2.1 The Large Hadron Collider . . . . .	3
2.2 The CMS Detector . . . . .	5
2.2.1 Inner Tracking System . . . . .	7
2.2.2 Electromagnetic Calorimeter . . . . .	8
2.2.3 Hadron Calorimeter . . . . .	12
2.2.4 Superconducting Magnet . . . . .	13
2.2.5 Muon System . . . . .	14
2.2.6 Trigger and Data Acquisition . . . . .	15
2.3 Computing . . . . .	18
2.3.1 Event Data Model . . . . .	18
2.3.2 Analysis Software . . . . .	19
2.4 Object Reconstruction . . . . .	20

---

2.4.1	Electron Reconstruction . . . . .	20
2.4.1.1	Electron Identification . . . . .	23
2.4.1.2	Electron Isolation . . . . .	27
2.4.1.3	Identification of photon conversions . . . . .	27
2.4.2	Muon Reconstruction . . . . .	28
2.4.3	Jet Reconstruction . . . . .	30
2.4.3.1	Jet Energy Corrections . . . . .	30
2.4.3.2	Particle Flow Jet Identification . . . . .	33
2.4.3.3	b-tagging . . . . .	33
2.4.4	Missing Transverse Energy . . . . .	34
2.5	Summary . . . . .	35
<b>3</b>	<b>High level triggers for Top Physics</b>	<b>37</b>
3.1	Level-1 triggers . . . . .	38
3.2	High-level triggers for top physics . . . . .	39
3.3	Trigger efficiency measurement . . . . .	42
3.3.1	Validation of jet energy corrections and charged hadron subtraction for top triggers . . . . .	44
3.4	Summary . . . . .	50
<b>4</b>	<b>Top Quark Mass Measurement</b>	<b>51</b>
4.1	Data and Simulation . . . . .	51
4.1.1	Data . . . . .	51
4.1.2	Simulation of signal and background processes . . . . .	51
4.1.2.1	Monte Carlo event generators . . . . .	52
4.1.2.2	Matching Threshold and Factorisation Scale . . . . .	53
4.1.2.3	CMS detector simulation . . . . .	53
4.1.2.4	Monte Carlo samples . . . . .	54
4.2	Event Selection . . . . .	57
4.3	Kinematic Fit . . . . .	60
4.4	The Ideogram Method . . . . .	63
4.4.1	Event likelihood . . . . .	63
4.4.2	Combined likelihood fit . . . . .	68
4.5	Mass Calibration . . . . .	69
4.6	Systematic Uncertainties . . . . .	70
4.7	Results . . . . .	74
4.8	Summary . . . . .	75
<b>5</b>	<b>Differential cross section measurement</b>	<b>77</b>
5.1	Data and Simulation . . . . .	77
5.1.1	Data . . . . .	77
5.1.2	Monte Carlo samples . . . . .	77
5.1.3	Pile-up reweighting . . . . .	78
5.1.4	b-tagging corrections . . . . .	82

5.2	Event Selection . . . . .	84
5.2.1	Electron plus jets channel . . . . .	84
5.2.2	Muon plus jets channel . . . . .	87
5.2.3	Additional filters for missing transverse energy . . . . .	89
5.3	Data-driven QCD estimation . . . . .	90
5.3.1	QCD multi-jet events with electrons . . . . .	91
5.3.2	QCD multi-jet events with muons . . . . .	91
5.4	Data/MC comparison . . . . .	93
5.5	Differential cross section measurement . . . . .	98
5.5.1	Primary variables . . . . .	98
5.5.2	Choice of binning . . . . .	99
5.5.3	Fitting procedure . . . . .	101
5.5.4	Unfolding . . . . .	107
5.5.5	Differential cross section calculation . . . . .	109
5.6	Systematic Uncertainties . . . . .	111
5.7	Results . . . . .	112
5.8	Summary . . . . .	115
<b>A</b>	<b>High level triggers</b>	<b>117</b>
<b>B</b>	<b>Choice of binning</b>	<b>119</b>
<b>C</b>	<b>Fitting templates</b>	<b>123</b>



# 1. Theoretical Background

## 1.1 The Standard Model of Particle Physics

1.1.1 Electroweak theory

1.1.2 Electroweak symmetry breaking

## 1.2 Top Quark Physics within the Standard Model

1.2.1 Top quark production at the LHC

1.2.2 Top quark decay

1.2.3 Background processes

1.2.4 Top quark mass

1.2.5 Importance of top quark physics

## 1.3 Physics Beyond The Standard Model

1.3.1 Supersymmetry

1.3.2 Extra dimensions



## 2. The LHC and the CMS detector

### 2.1 The Large Hadron Collider

The LHC [1] is currently the largest and the most powerful particle accelerator ever built. It is installed in the 26.7 km tunnel that was originally constructed for the LEP accelerator in the 1980s. The tunnel lies at a depth of 45 m to 170 m underground between the Jura mountain and Lake Geneva, being the main part of the CERN accelerator complex.

The machine is designed to accelerate proton beams and provide collisions at a centre of mass energy of  $\sqrt{s} = 14$  TeV. Unlike particle-antiparticle colliders, the LHC requires two rings with opposite magnetic dipole fields in order to maintain and collide two counter-rotating proton beams. Since the tunnel was originally designed for the electron-positron LEP, it has an internal diameter of 3.7 m which is not enough to install two separate independent rings. Therefore, a twin-bore magnet design was adopted [2], which resulted in substantial cost savings.

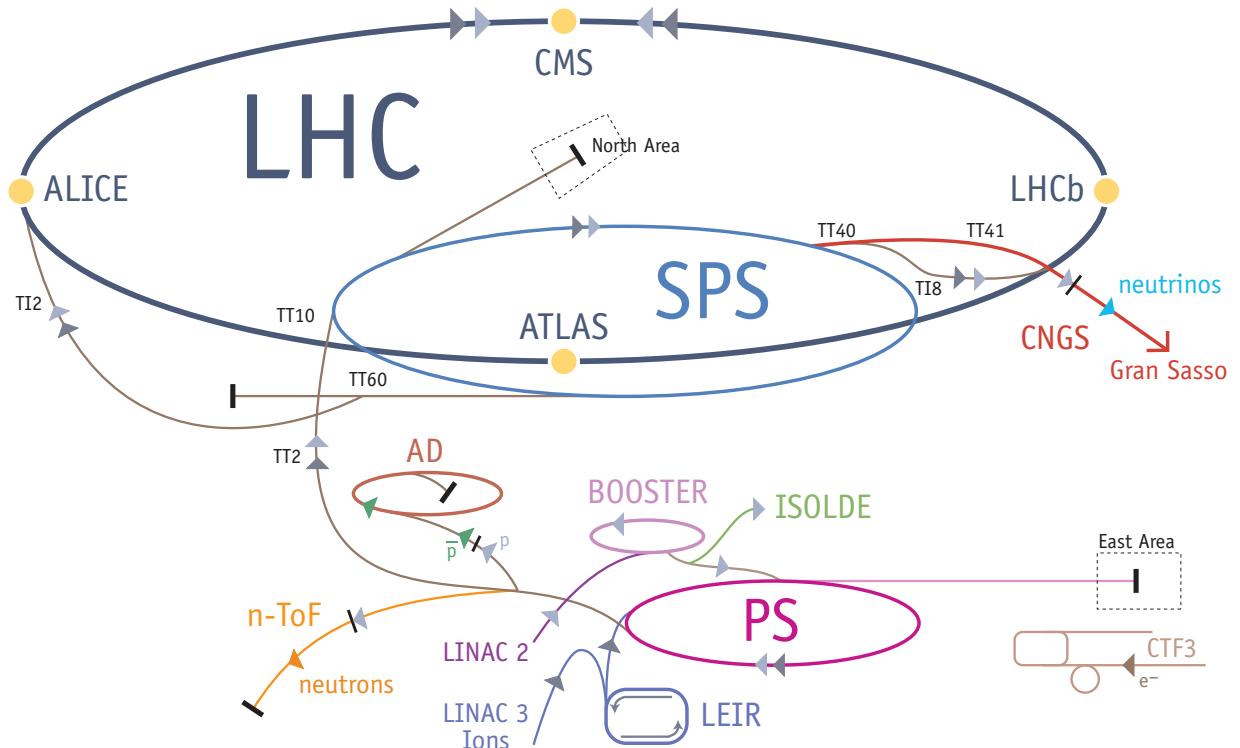


Figure 2.1: CERN accelerator complex.

A schematic view of the LHC accelerator chain is shown in Figure 2.1. Initially, the protons are obtained by stripping orbiting electrons from hydrogen atoms. Then they are injected into the linear accelerator LINAC2 to reach the energy of 50 MeV and enter the Proton Synchrotron Booster (PSB). The booster accelerates them to 1.4 GeV and passes the beam to the Proton Synchrotron (PS) where the energy rises to 25 GeV. In the next step, protons enter the Super Proton Synchrotron (SPS) where they are accelerated to 450 GeV. Finally, the beam is transferred to the LHC in both clockwise and anti-clockwise directions where it takes about 20 minutes to reach the design 7 TeV energy (per beam).

The LHC has four interaction points, providing collisions to four major experiments. Two of them, CMS and ATLAS, are multi-purpose high-luminosity experiments with a peak luminosity of  $L = 10^{34} \text{ cm}^{-2}\text{s}^{-1}$ . The other two experiments operate at low luminosities and have more specific physics goals: LHCb studies b-meson decays, and Alice is a dedicated heavy ion experiment.

The instantaneous luminosity of a collider can be calculated as

$$L = \frac{n_1 n_2 f}{4\pi\sigma_x\sigma_y}, \quad (2.1)$$

where  $n_1$  and  $n_2$  are the numbers of particles in each of the colliding bunches,  $f$  is the revolution frequency,  $\sigma_x$  and  $\sigma_y$  are the horizontal and vertical beam sizes, assuming the two beams have the same size.

The number of events generated in the collisions per second is given by

$$N_{events} = L \times \sigma, \quad (2.2)$$

where  $\sigma$  is the cross section of the process under study.

The LHC started operating on the 10th of September 2008, with the first beams fully circulating in both rings. However, only 9 days later a magnet quench occurred in two sectors of the tunnel, which was caused by an electrical fault due to a bad connection between two magnets. A consequent liquid helium explosion damaged a total of 53 superconducting magnets. Over a year was spent on repairs and tests, and the first collisions were recorded on the 23rd of November 2009 at a centre of mass energy of 0.9 TeV. The following few months showed the continuous ramp up of the beam energies up to 3.5 TeV per beam which was achieved on the 30rd of March 2010 when the LHC physics programme started.

Throughout the rest of 2010, the two general-purpose LHC experiments (CMS and ATLAS) recorded approximately  $40 \text{ pb}^{-1}$  of data, which resulted in the first measurements of

various physics processes at the LHC. The following year became the main 7 TeV data-taking period, with about  $5 \text{ fb}^{-1}$  of data recorded by ATLAS and CMS. On the 5th of April 2012 the centre of mass energy was increased to 8 TeV, and July of 2012 marked the first major discovery of a new boson which was later shown to be consistent with the Standard Model Higgs boson, according to approximately  $21.8 \text{ fb}^{-1}$  of data recorded until early 2013. A long shut-down is planned for the following two years with various upgrades scheduled. The next physics run is expected in 2015 with the beam energy increased up to 6 or 7 TeV.

## 2.2 The CMS Detector

The Compact Muon Solenoid [3] is a general-purpose detector designed to carry out precise measurements of the Standard Model and searches for physics beyond it. The primary design requirement was the ability to discover the nature of electroweak symmetry breaking, and the first observation of a Higgs boson was obtained in the Summer of 2012 [4].

The detector is installed at one of the LHC interaction points (Point 5) at about 100 m underground near the French village of Cessy, between the Jura mountains and Lake Geneva. The overall dimensions of the CMS detector are a length of 21.6 m, a diameter of 14.6 m and a total weight of 12 500 t.

The sectional view of CMS is shown in Figure 2.2. In the centre of the detector, tracking and calorimetry systems are surrounded by the superconducting solenoid. On the outermost part of it the magnetic flux is returned through the iron yoke in which the muon system is also integrated. All the sub-systems are discussed in the following sections in more detail.

The cylindrical shape of the CMS detector dictates using a cylindrical coordinate system, with the origin centred at the interaction point, the  $x$ -axis pointing towards the centre of the LHC ring, the  $y$ -axis pointing upwards and the  $z$ -axis pointing along the beamline in the anti-clockwise direction. The azimuthal angle  $\phi$  is measured from the  $x$ -axis in the transverse ( $x - y$ ) plane and the polar angle  $\theta$  is measured from the  $z$ -axis. The radial distance to the beamline is denoted by  $r$ . Pseudorapidity is defined as:

$$\eta = -\ln \tan \frac{\theta}{2}. \quad (2.3)$$

This implies that the particles moving in the transverse plane (perpendicular to the beamline) have a pseudorapidity of 0, whereas the beam direction has an infinite pseudorapidity. Considering the cylindrical shape of the detector, it has barrel and endcap regions, with the transition occurring at  $\eta \sim 1.4$ . The momentum and energy transverse to the beamline are denoted by  $p_T$  and  $E_T$  respectively; the imbalance of the energy measured in the transverse

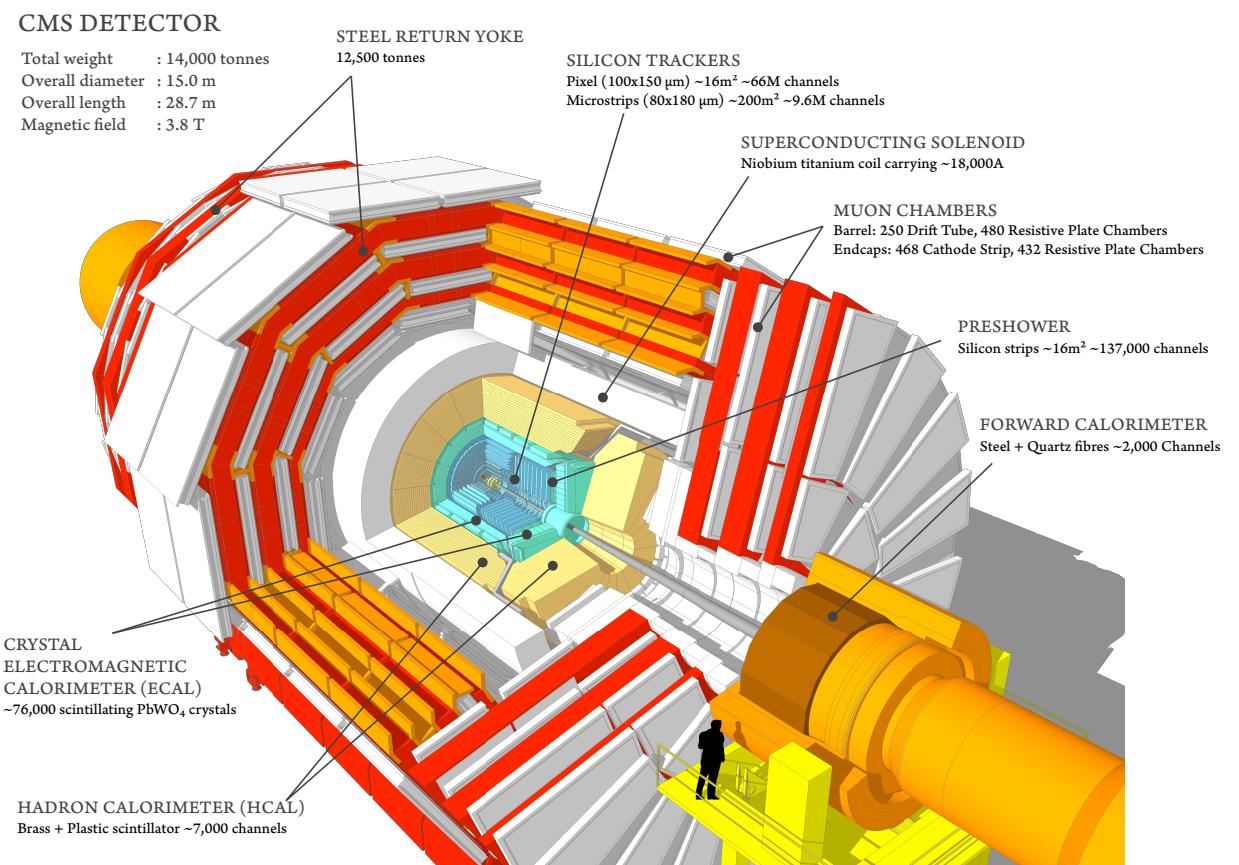


Figure 2.2: Sectional view of the CMS detector.

plane, called missing transverse energy, is denoted by  $E_T^{\text{miss}}$ .

### 2.2.1 Inner Tracking System

The tracking system lies in the heart of the CMS detector and is the closest to the interaction point where the particle flux has the highest value. This imposes demanding requirements on the configuration of the system. At design luminosity of  $L = 10^{34} \text{ cm}^{-2}\text{s}^{-1}$  with the bunch spacing of 25 ns, an average of 1000 particles from about 25 proton-proton interactions (pile-up vertices) is expected to traverse the tracker for each bunch crossing. However, up until the long shutdown a bunch spacing of 50 ns was used, which meant a higher number of protons in each bunch leading to approximately twice the number of pile-up vertices. Therefore, in order for the particle tracks to be identified reliably and separately for each bunch crossing, the tracker requires very fine granularity and fast response parameters. Another complication caused by the intense particle flux is the severe radiation damage, so the tracker has to be highly resilient in operating in the harsh environment for a reasonable lifetime.

To meet these requirements on granularity, response time and radiation resilience, the tracker design was chosen to be based on silicon detector technology. Although capable of meeting such conditions, this technology has a disadvantage of a high power density of on-detector electronics. This implies the necessity of an efficient cooling system. Moreover, a large amount of dense material interacting with the particles leads to higher multiple scattering, bremsstrahlung, photon conversions and nuclear interactions. Therefore, there are complications in the reconstruction of the tracks, meaning some loss of efficiency and precision. This will be discussed in detail later on in the object reconstruction section.

Figure 2.3 shows the overall layout of the tracking system. It consists of the inner pixel detector, located in the vicinity of the interaction point, and silicon strip tracker detectors: inner barrel and disks (TIB and TID), outer barrel (TOB) and endcaps (TEC). The geometrical acceptance of the tracker system goes up to  $|\eta| < 2.5$ . The outer radius of the CMS tracker reaches approximately 110 cm, and its total length is about 540 cm.

The pixel detector consists of three layers of pixel sensors at radii of 4.4 cm, 7.3 cm and 10.2 cm from the beamline in the barrel region. In addition there are two endcap disks on each side at  $|z| = 34.5$  cm and 46.5 cm. The pixel size equals  $100 \times 150 \mu\text{m}^2$  in  $r\phi \times z$  coordinates. The pixel detector has 66 million pixels and the total area of about 1 m<sup>2</sup>.

The silicon strip tracker consists of several layers of silicon microstrip detectors. It covers the region between 20 cm to 110 cm in radius and extends up to  $\pm 280$  cm in the z direction. The Tracker Inner Barrel (TIB) is made out of 4 layers and the Tracker Outer Barrel (TOB) has 6 layers in it. The tracker endcaps (TEC) comprise 9 disks, and there are also the tracker



Figure 2.3: Cross section of the CMS tracker system [3].

inner disks (TID) that consist of 3 disks filling the gap between TIB and TEC as shown in Figure 2.3. There are 9.3 million silicon strips covering the area of about  $200 \text{ m}^2$ . The silicon sensors' thickness varies between 320 and 500  $\mu\text{m}$  and the strip pitch varies from 80  $\mu\text{m}$  in the TIB to 180  $\mu\text{m}$  in TOB and TEC.

The silicon detectors of the tracker, the readout electronics and support structure form a considerable amount of material for the particles traversing from the interaction point. Figure 2.4 [3] shows the material budget of the CMS tracker in units of radiation lengths<sup>1</sup> ( $X_0$ ). It grows from about  $0.4 X_0$  to  $1.8 X_0$  in the barrel region, and then decreases to about  $1 X_0$  in the endcaps. This causes a substantial conversion rate for photons and electrons in the tracker material; it also will be discussed in more detail in the electron reconstruction section.

## 2.2.2 Electromagnetic Calorimeter

The next detector subsystem which is surrounding the tracker is the electromagnetic calorimeter, or ECAL. It is of a primary importance for the analyses described in this thesis, as it provides information for the electron and positron reconstruction. Combination of this information with that from the tracking system must ensure a precise measurement of electron position and momentum, and also sufficient background removal. It has to effectively

<sup>1</sup>A material's radiation length is the mean distance over which a high-energy electron loses all but  $1/e$  of its energy by bremsstrahlung; this is equal to  $7/9$  of the mean free path for pair production by a high-energy photon.

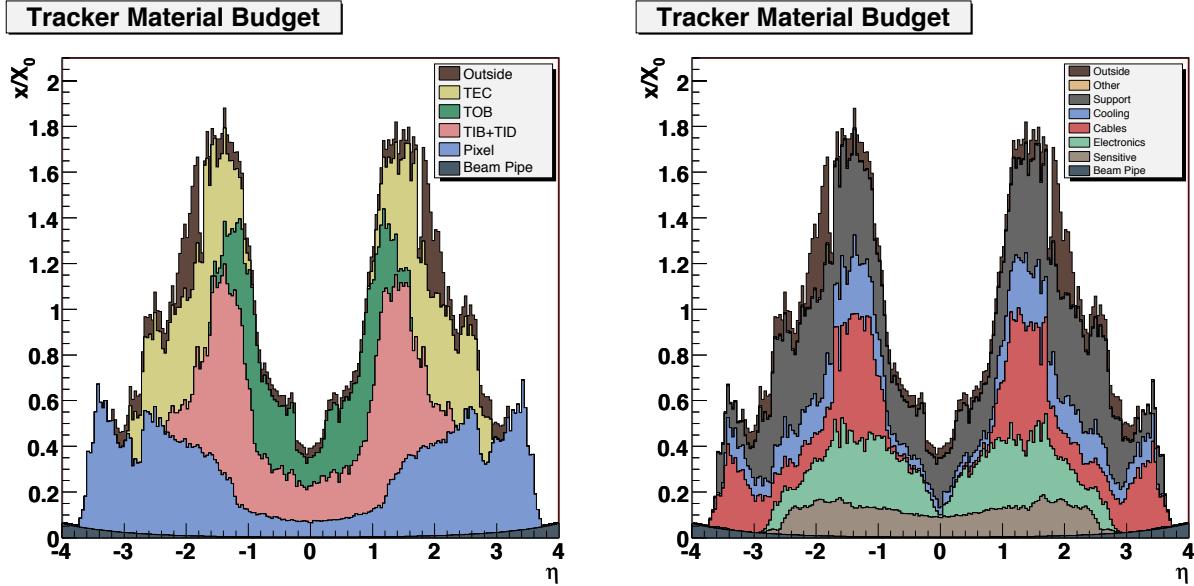


Figure 2.4: Material budget as a function of pseudorapidity  $\eta$  for the different sub-detectors of the tracker (left) and broken down into the functional contributions (right), in units of radiation length [3].

distinguish the energy deposit shape of an electromagnetic particle from the one of a hadronic particle, which requires good segmentation and high resolution.

ECAL is a hermetic, high-granularity, high-resolution scintillating crystal calorimeter consisting of 61 200 lead tungstate ( $\text{PbWO}_4$ ) crystals located in the central barrel region ( $|\eta| < 1.479$ ), and 7324 crystals in each of the two endcaps ( $1.479 < |\eta| < 3.0$ ). All crystals are followed by photodetectors reading and amplifying their scintillation: avalanche photodiodes (APD) are used in the barrel, and vacuum phototriodes (VPTs) are used in the endcaps. These different choices were caused by the configuration of the magnetic field and the expected level of radiation.

The layout of the ECAL sub-detector is shown in Figure 2.5. An additional preshower detector is used in the endcap region to lower the required detector depth. Its principal aim is to identify neutral pions in the endcaps, but it also helps to distinguish neutral pions and electrons from minimum ionising particles and improves the position determination of electrons and photons with high granularity.

The main geometrical characteristics of the ECAL crystals are shown in Table 2.1. The choice of lead tungstate was driven by the constraints of the CMS design. It is a very dense material ( $8.28 \text{ g/cm}^3$ ) with a short radiation length of  $X_0 = 0.89 \text{ cm}$ , which allows the calorimeter to fit inside the compact magnet. Lead tungstate also has a small Molière

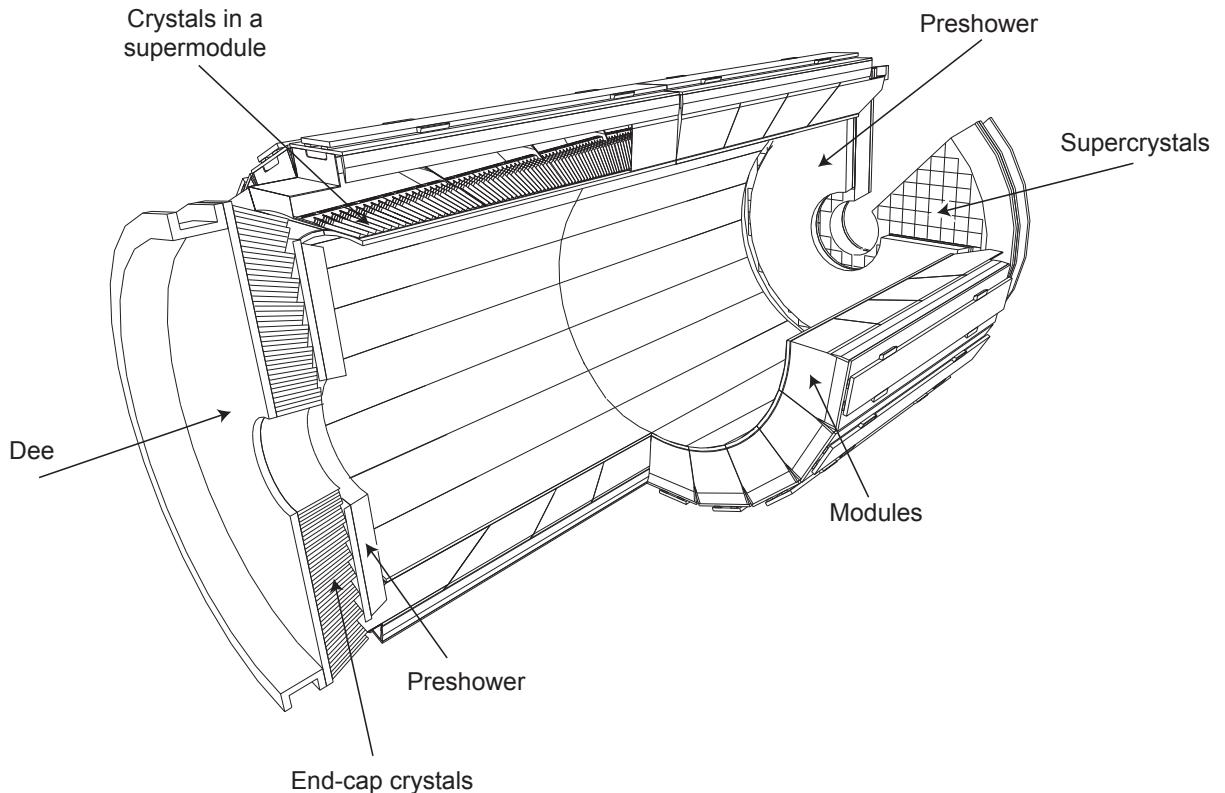


Figure 2.5: Layout of the CMS electromagnetic calorimeter [3].

Table 2.1: ECAL crystal characteristics

	Barrel	Endcaps
number of crystals	61 200	14 648
crystal cross section in $(\eta, \phi)$	$0.0174 \times 0.0174$	varies
crystal cross section at the front	$22 \times 22 \text{ mm}^2$	$28.62 \times 28.62 \text{ mm}^2$
crystal cross section at the rear	$26 \times 26 \text{ mm}^2$	$30 \times 30 \text{ mm}^2$
crystal length	230 mm ( $25.8X_0$ )	220 mm ( $24.7X_0$ )

radius<sup>1</sup> of 2.2 cm, which allows a calorimeter with fine granularity. Finally, the crystals emit 80 % of their scintillation light in just 25 ns, however the light yield is relatively low. At 18 °C, about 4.5 photoelectrons per MeV are collected. The dependence of the light yield on temperature requires a cooling system capable of keeping the crystal temperature stable within ±0.05 °C to preserve energy resolution [5].

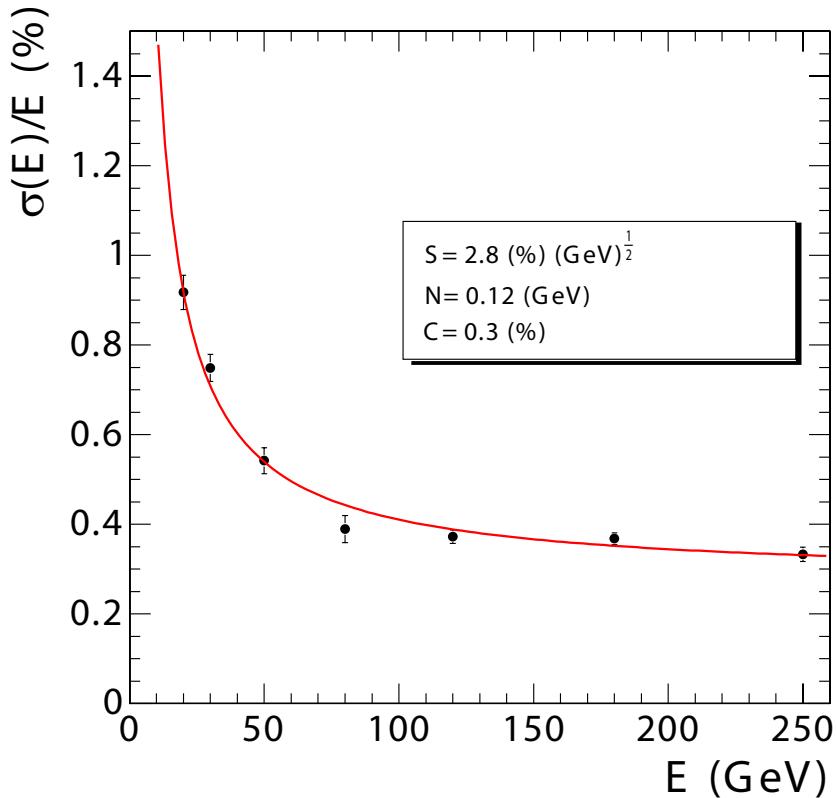


Figure 2.6: ECAL energy resolution derived from the test beam measurements as a function of deposited energy. The stochastic, noise, and constant contributions are shown [3].

The energy-dependent resolution of the calorimeter can be parameterised as follows [3]:

$$\left(\frac{\sigma}{E}\right)^2 = \left(\frac{S}{\sqrt{E}}\right)^2 + \left(\frac{N}{E}\right)^2 + C^2. \quad (2.4)$$

where  $S$  is the stochastic term,  $N$  is the noise term, and  $C$  is the constant term. Figure 2.6

---

<sup>1</sup>The Molière radius  $R_\mu$  is a characteristic constant of a material giving the scale of the transverse dimension of the fully contained electromagnetic showers initiated by an incident high energy electron or photon. It is defined as the radius of a cylinder containing an average of 90 % of the shower's energy deposition.

shows the energy resolution measured using incident electrons, during the beam tests in 2004.

### 2.2.3 Hadron Calorimeter

The hadron calorimeter (HCAL) is the next sub-detector located mostly inside the solenoid and completing the CMS calorimetry system. It is essential for the measurement of hadron jets and missing transverse energy.

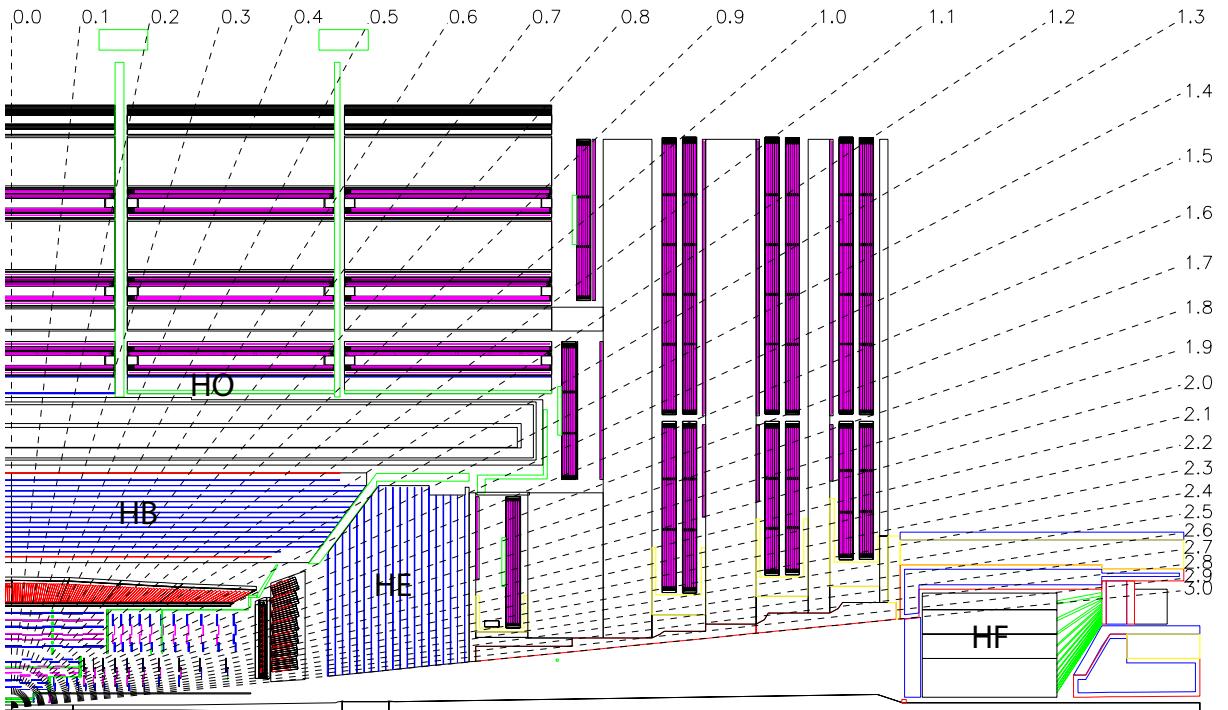


Figure 2.7: Longitudinal view of the CMS detector showing the locations of the hadron barrel (HB), endcap (HE), outer (HO) and forward (HF) calorimeters [3].

As shown in Figure 2.7, HCAL consists of four subsystems: the hadron barrel calorimeter (HB), the hadron endcap calorimeter (HE), the hadron outer calorimeter (HO) and the hadron forward calorimeter (HF). The barrel and endcap parts (HB, HE) cover the pseudorapidity range up to  $|\eta| < 3.0$ , and the forward part (HF) extends it to a total coverage of  $|\eta| < 5.0$ . HCAL surrounds ECAL from its outer limit of 1.77 m from the beamline, to the inner limit of the magnet coil at 2.95 m from the beamline. However, due to space limitations the barrel calorimeters do not contain complete hadronic showers, therefore an outer calorimeter (HO) was designed to measure the energy leakage. It is placed in the muon system just outside of the solenoid in the barrel region.

HCAL is a sampling calorimeter consisting of alternating layers of brass and stainless steel absorbers, and plastic scintillators as active elements. The choice of the absorber material was caused by its short hadronic interaction length and its property of being non-magnetic, which is crucial in the strong magnetic field of the CMS magnet. The scintillation light is guided by embedded wavelength-shifting (WLS) fibres. The light from the WLS is then transmitted via a network of clear fibres, arranged in read-out towers, to hybrid photodiodes (HPDs) [3].

Both HB and HE scintillators have a granularity of  $\Delta\eta \times \Delta\phi = 0.087 \times 0.087$  for  $|\eta| < 1.6$ , and  $\Delta\eta \times \Delta\phi = 0.17 \times 0.17$  for  $|\eta| \geq 1.6$ . The tower segmentation of the forward calorimeter (HF) varies from  $\Delta\eta \times \Delta\phi = 0.175 \times 0.175$  at  $|\eta| = 3.0$  to  $\Delta\eta \times \Delta\phi = 0.3 \times 0.35$  at  $|\eta| = 5.0$ . The HF is placed at about 11 m from the interaction point, and is essential to reconstruct very forward hadron jets. Together with HO, it provides the hermeticity of the calorimetry system, making it possible to measure the transverse missing energy to a reasonable precision.

#### 2.2.4 Superconducting Magnet

The superconducting solenoid is a central feature of the CMS apparatus, essentially giving it its name. The magnet has a length of 12.5 m, diameter of 6.3 m and mass of 220 t. Although it was initially designed to sustain a uniform magnetic field of 4 T within the 5.9 m diameter free bore, operation at 3.8 T was chosen in order to increase the lifetime. The magnetic field is returned by a massive iron yoke. The main parameters of the CMS magnet are shown in Table 2.2.

Table 2.2: Parameters of the CMS superconducting solenoid [5, 6].

Field	3.8 T
Inner Bore	5.9 m
Length	12.5 m
Number of Turns	2168
Current	18 160 kA
Stored energy	2.3 GJ

The large bending power of the solenoid is required to bend the tracks of high energy charged particles to an extent where good momentum resolution is achieved. The design requirement for the strength of the magnetic field was the ability to unambiguously determine the sign of the electric charge for muons with a momentum of  $\approx 1$  TeV/c [5].

The solenoid coil is constructed from four layers of superconducting high-purity niobium-titanium cable co-extruded with pure aluminium, which acts as a thermal stabiliser. The cold mass is cooled down to 4.5 K by liquid helium. If a fast discharge happens (e.g. caused by a magnet quench), about 3 days are necessary to re-cool the coil.

### 2.2.5 Muon System

The last sub-detector placed on the outermost part of CMS is the muon system. Since the muons are the most penetrating particles detectable by CMS, they have the cleanest signature and play an important role in many physics analyses. Due to their ability to travel through the many layers of the calorimeters, muons are relatively easy to identify and separate from the background.

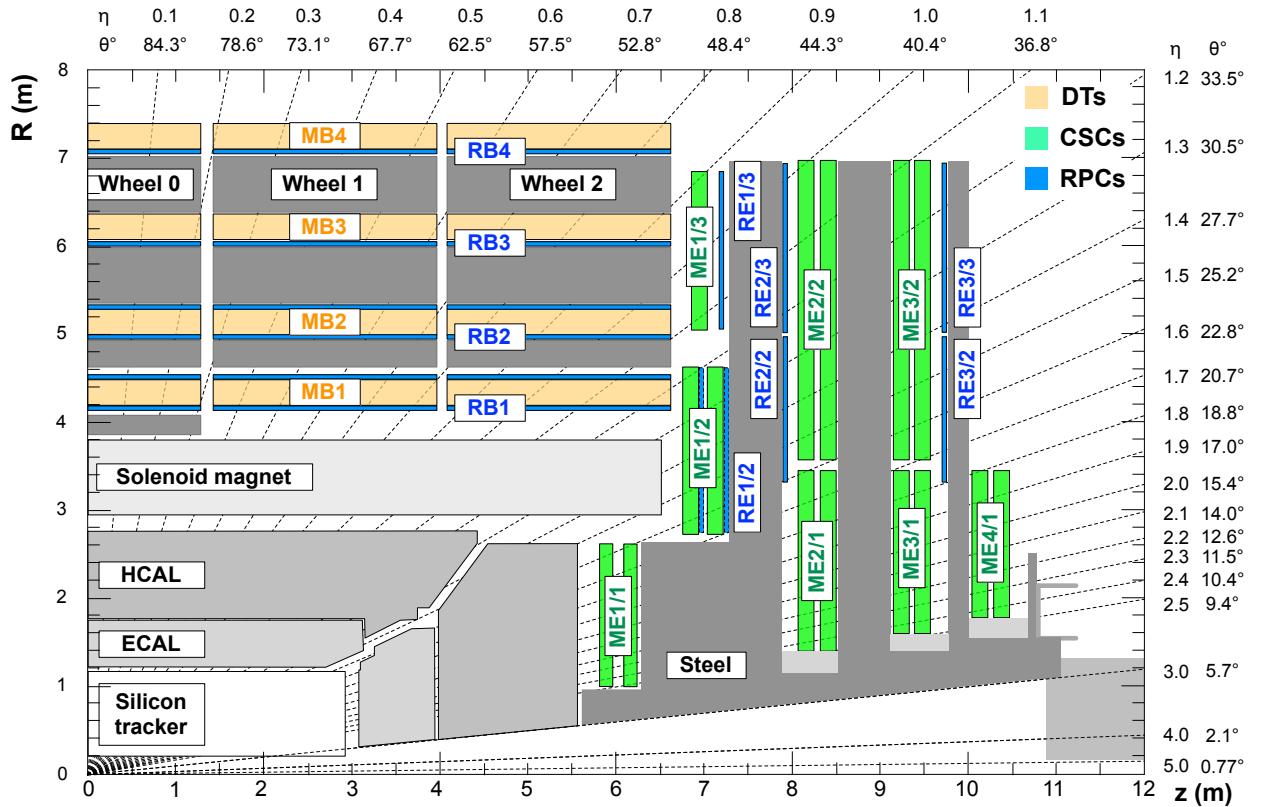


Figure 2.8: Layout of one quarter of the CMS muon system. Four drift tube (DT, in light orange) stations are labeled MB (âĂĲmuon barrelâĂ İl) and the cathode strip chambers (CSC, in green) are labeled ME (âĂĲmuon endcapâĂ İl). Resistive plate chambers (RPC, in blue) are in both the barrel and the endcaps of CMS, where they are labeled RB and RE, respectively.

The layout of the CMS muon system is shown in Figure 2.8. It consists of the drift tubes

(DT), cathode strip chambers (CSC) and resistive plate chambers (RPC). The entire system surrounds the solenoid and covers the pseudorapidity region of  $|\eta| < 2.4$ .

The drift tubes are located in the barrel region ( $|\eta| < 1.2$ ). Consisting of four stations, they form concentric cylinders around the beam line; there are 250 drift chambers with about 172 000 sensitive wires in total. When a muon passes through the volume, it knocks electrons off the atoms of the gas, which then follow the electric field and reach the positively-charged wires, providing information on the muon's position. The chambers are filled with the gas mixture of 85 % Ar and 15 % CO<sub>2</sub>, where the muon drift time does not exceed 380 ns. Although this value is bigger than the typical bunch crossing time (25 or 50 ns), it is sufficient because of the small muon rate in this region.

In the endcaps, the cathode strip chambers cover the pseudorapidity region of  $0.9 < |\eta| < 2.4$ . Each of 468 CSCs is a trapezoidal multi-wire proportional chamber consisting of 6 gas gaps with a plane of radial cathode strips and a plane of anode wires which are roughly perpendicular. A charged muon traversing each plane of a chamber causes gas ionisation and a subsequent electron avalanche which produces a charge on the anode wire and an image charge on the cathode strips. The gas used in CSCs is a mixture of Ar, CO<sub>2</sub> and CF<sub>4</sub>.

The resistive plate chambers system is complementary to both DT and CSC systems, and is located in both barrel and endcap regions ( $|\eta| < 2.1$ ). RPCs also operate in avalanche mode with a gas mixture of C<sub>2</sub>H<sub>2</sub>F<sub>4</sub>, C<sub>4</sub>H<sub>10</sub> and SF<sub>6</sub>, and due to an excellent time resolution of about 1 ns they provide fast information for triggering. The spacial resolution is, however, quite limited ( $\approx 1$  cm, compared to  $\approx 100$   $\mu$ m for DTs and CSCs).

The muon momentum is measured in both the tracker and the muon system. As it can be seen on Figure 2.9, both sub-systems contribute to the momentum resolution at different  $p_T$  values. This happens due to the difference in the magnetic field and detector technology. For low- $p_T$  muons, the best momentum resolution is obtained in the tracker, whereas in the high- $p_T$  region the muon system provides a significant improvement. Therefore, by using information from both the silicon tracker and the muon chambers (i.e. reconstructing the “global muon”), the momentum resolution is improved in the whole  $p_T$  region up to a  $\approx 1$  TeV/c level.

## 2.2.6 Trigger and Data Acquisition

At design LHC luminosity of  $L = 10^{34}$  cm<sup>-2</sup>s<sup>-1</sup>, approximately 25 collisions are expected to occur at each crossing of the proton bunches. The bunch spacing of 25 ns corresponds to a crossing rate of 40 MHz. Since every event produces  $\sim 1$  MB of raw data, it corresponds to a total data production of 40 TB s<sup>-1</sup>. Attempting to store all of this data is clearly beyond the

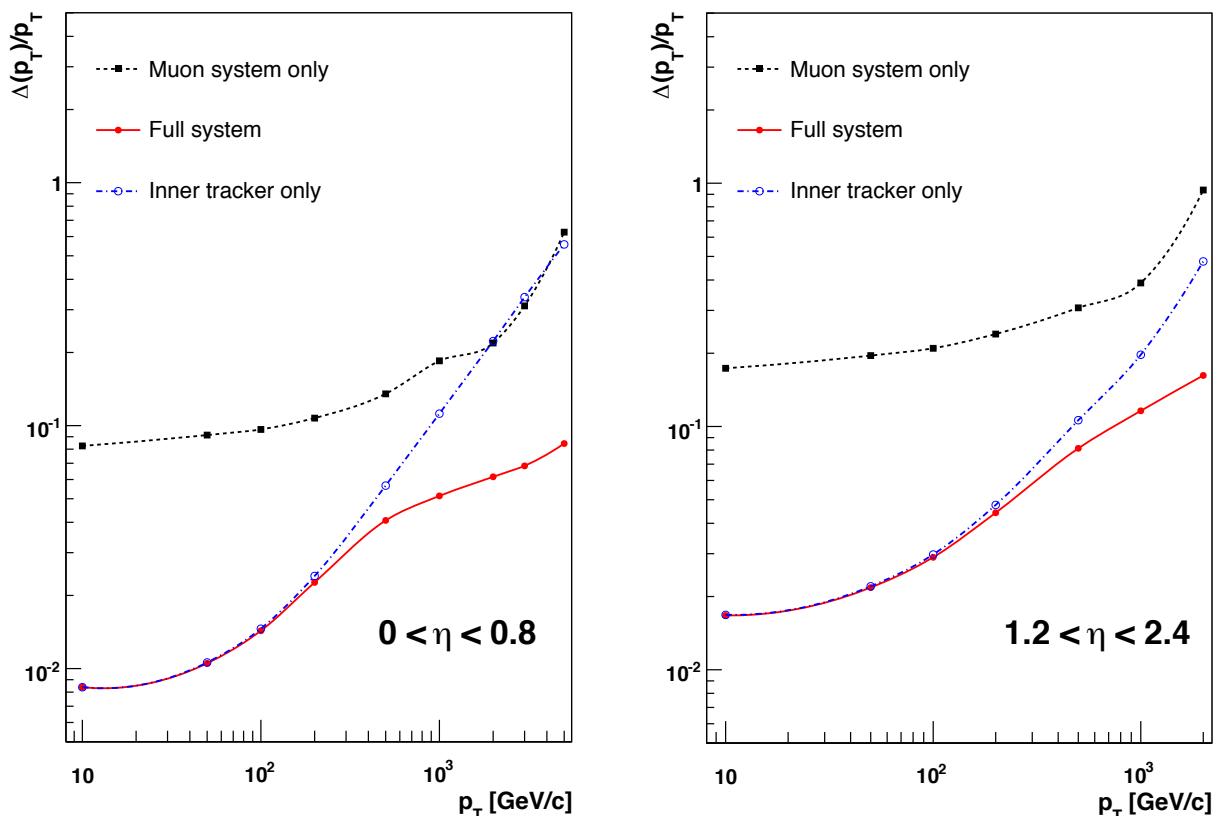


Figure 2.9: The muon transverse momentum resolution as a function of the transverse momentum ( $p_T$ ) using the muon system only (black), the inner tracking only (blue), and both (red), in regions of  $|\eta| < 0.8$  (left) and  $1.2 < |\eta| < 2.4$  [3].

available technology. Moreover, only a fraction of events contain hard scattering processes that are of interest, therefore an effective trigger system had to be implemented.

The CMS trigger is a two-level system, consisting of two independent parts: the Level-1 (L1) trigger and the High-Level Trigger (HLT). The L1 trigger is a hardware system implemented in programmable electronics residing partly on detector, and partly in the underground control room located at approximately 90 m from the experimental cavern. The maximum latency between the collision and the L1 accept decision received by front-end electronics is 3.2  $\mu$ s. During this amount of time, the complete event information is buffered in pipelined memories on the detector. The only information used for the L1 trigger decision is that from the muon system and the calorimetry. Since the reconstruction of tracks exceeds the time scale required for the L1 decision, the tracker information can't be used. The L1 trigger reduces the event rate from  $\sim 40$  MHz to  $\sim 100$  kHz, corresponding to a data flow of about 100  $\text{GB s}^{-1}$ . These events are fed into the HLT system.

The High-Level Trigger is a software system implemented in a single CPU farm, sometimes referred to as the “Event Filter Farm”. Having access to the full event information, customised algorithms of increasing complexity are used which results in a highly flexible trigger system. The event rate is reduced down to  $\sim 300$  Hz, with the final data rate of approximately 300  $\text{MB s}^{-1}$  being stored on a large disk cache at the experimental site (the Storage Manager) and later on transferred to CERN Tier 0 for further processing (see Section 2.3).

Since the start of the LHC running, the operating conditions have been changing drastically. During the start-up year of 2010, the instantaneous luminosity went up from about  $10^{27} \text{ cm}^{-2}\text{s}^{-1}$  to approximately  $0.2 \times 10^{32} \text{ cm}^{-2}\text{s}^{-1}$ . In 2011 the luminosity ramped up to a factor of 20 above that of 2010, reaching approximately  $4 \times 10^{33} \text{ cm}^{-2}\text{s}^{-1}$ . This required a lot of continuous effort to control the trigger rates at reasonable level, whilst also keeping its efficiency acceptable. In 2012 the luminosity was more stable, peaking at  $\approx 7.6 \times 10^{33} \text{ cm}^{-2}\text{s}^{-1}$  which is just a factor of 2 above the 2011 values. However, it still came as a challenge because of the impact of pile-up. At a bunch spacing of 50 ns and increased centre-of-mass energy of 8 TeV, the average number of pile-up vertices nearly doubled comparing to that in 2011, which required a major CPU extension and implementation of sophisticated PU mitigation techniques at the HLT level. The author's contribution to the HLT development of the trigger paths important for top physics is described in Chapter 3.

## 2.3 Computing

The vast amounts of data delivered by the CMS detector impose high requirements on the offline computing system. During 2010–2012 operation, CMS collected  $\sim 10$  PB of raw data per year. Including Monte Carlo simulations, reconstructed data and analysis skims, the total annual amount of data essentially doubles. To handle the distributed storage and processing of this data, not just for CMS but for the entire high energy physics community using the LHC, a worldwide LHC computing grid (WLCG) has been put in place.

WLCG is a global collaboration of more than 150 computing centres in about 40 countries. The grid has a tiered architecture, comprising 4 tiers with different resources and services. The first one, Tier 0, is based at CERN and is responsible for data-taking. It accepts raw data from the data acquisition system and repacks it into primary datasets according to the trigger information. The raw data is archived to tape, and is also prompt-reconstructed (within 48 hours) before being distributed to the Tier 1 (T1) centres around the world. There are 8 T1 sites based at large national laboratories in collaborating countries (e.g. RAL in the UK and FNAL in the US). Each of the T1 centres is used for large-scale centrally organised data-processing activities. The data is then distributed in the reduced format (see Section 2.3.1) to a more numerous set of Tier 2 centres, typically located at collaborating universities. Each of these centres is used for the grid-based analysis and Monte Carlo simulation for the whole experiment, as well as local services for groups maintaining them. The last stage of computing system, Tier 3, is meant solely for the local institution’s user analysis.

### 2.3.1 Event Data Model

In the basis of the CMS Event Data Model lies the concept of an event, which is physically a result of a single collision in the LHC. From a software point of view, the event is a C++ object container storing raw data from a single readout of detector electronics (e.g. hits in various sub-detectors), as well as reconstructed data which is based on this information, such as tracks, clusters and physics objects. All these C++ objects are stored in ROOT format [7].

The EDM makes use of three main data formats, based on different levels of detail and precision:

- RAW format, containing full information from the detector as well as L1 and HLT trigger decisions, with the event size of  $\sim 1.5$  MB.
- RECO (reconstructed data) format, which is obtained from raw data by application of pattern recognition and compression algorithms. This data includes reconstructed

detector hits, clusters and physics objects (electrons, muons, etc.). The typical event size is  $\sim 250$  kB.

- AOD (Analysis Oriented Data) format, produced by filtering the RECO data from the reconstructed detector objects, leaving just the high-level physics objects required for analysis. The event size is reduced down to  $\sim 50$  kB.

The RECO and AOD data are analysis-ready data formats, produced centrally and used by many physics analysis groups. However, further simplification of the data is also a common practice. By transforming the C++ objects produced by CMS software into plain basic types or vectors of them, only including the analysis-specific content, the event size can be reduced down to  $\sim 3$  kB level depending on the needs of a particular analysis. This data format is often referred to as private “ntuples”, and it requires specific analysis software capable of restructuring the data into user-defined classes. By following this approach, the analysis can be run locally and generally much faster than processing the RECO or AOD data. However, it requires “ntuplising” this data every time when new centrally-recommended physics objects or corrections are produced.

### 2.3.2 Analysis Software

Both of the analyses described in this thesis use the CMS software framework (CMSSW [8]), as well as Bristol Analysis Tools (BAT [9]). The differential cross section analysis also uses an additional level of python scripts for post-processing [10].

CMSSW is the key CMS software framework built around the Event Data Model (see Section 2.3.1). The framework is essential for purposes of Monte Carlo simulation, detector calibration and alignment, as well as data reconstruction and analysis. CMSSW has a modular architecture, consisting of one configurable executable (`cmsRun`) and a large set of plug-in modules that contain all the code needed for event processing (reconstruction algorithms, calibration, etc.). Different versions of CMSSW were used for different analyses:

- `CMSSW_4_2_8` for the top mass analysis on 2011 data;
- `CMSSW_4_4_4` for the missing transverse energy analysis on 2011 data;
- `CMSSW_5_3_9` for the top cross pair cross section analysis on 2012 data.

Corresponding versions were used to produce ntuples for processing by BAT, which was used to read the data, apply selections, calculate high-level variables and to create various histograms of distributions. BAT was originally started in 2010 by Dr. Lukasz Kreczko

for the needs of the Bristol top group, later on also developed by the author and other researchers from Bristol and affiliated top groups. Like CMSSW, this framework has a modular structure, with its classes falling in four main categories:

- readers, for translating plain data types from ROOT files into C++ objects;
- RECO objects, i.e. output of the readers (physical objects like leptons, jets and its collections);
- selections, for application of event selections;
- analysers for creating histograms, applying selections, algorithms, and filling histograms.

All analysers are independent from each other, making the analysis chain stable and reliable. The final set of python scripts is used to prepare the histograms, perform fitting and unfolding procedures (in case of cross section analysis), and producing final tables and plots. Rootpy package [11] was used to access ROOT libraries in python interface, and matplotlib [12] was used to create plots.

## 2.4 Object Reconstruction

Most CMS analyses, including the ones described in this thesis, adopt a reconstruction technique called Particle Flow (PF) [13]. This algorithm is used to obtain a global event description at level of individually reconstructed particles by means of combining information coming from all sub-detector systems. The ultimate goal is to determine type, energy and momentum of all the particles in the event with highest possible precision and in the most optimal way. The types of these particles include electrons, muons, charged hadrons, neutral hadrons and photons. All these particles are then used to reconstruct jets (Section 2.4.3), missing transverse energy (Section 2.4.1) and tau leptons from their decay products.

### 2.4.1 Electron Reconstruction

The reconstruction of the  $t\bar{t}$  pair with an electron in the final state imposes high requirements on the electron identification and its energy-momentum measurement, precision of which is of major importance for both top mass and  $t\bar{t}$  cross section measurements.

Although the CMS detector is equipped with highly accurate ECAL and tracker systems, electron identification and reconstruction is still a challenging task due to the large amount of tracker material (see Section 2.2.1). This results in a significant Bremsstrahlung photon emission, which often causes an ECAL energy deposit to be widely spread in azimuthal

direction because of the high magnetic field. Therefore, dedicated algorithms were developed in order to collect all Bremsstrahlung energy deposits in the calorimeter (Bremsstrahlung recovery), and also to take into account the kinks in the electron trajectory caused by photon emissions.

Electron reconstruction in CMS has following distinct stages: seeding, track finding, pre-identification, Bremsstrahlung recovery, track-cluster linking and final identification. Historically, the original seeding algorithm was designed and optimised for isolated high- $p_T$  electrons. This approach starts from ECAL clusters, and therefore is called the “ECAL-driven” seeding. It is based on the property of the ECAL energy deposits to have narrow width in the  $\eta$  coordinate, and to be widely spread in  $\phi$  (azimuthal direction) like it was mentioned above. The electron and all the associated Bremsstrahlung energy deposits form a single “super-cluster”, and the ability to correctly determine it affects the overall performance of this method. Only super-clusters with transverse energy above 4 GeV are taken into account. Super-clusters are then matched to pairs or triplets of hits in the inner tracker layers, forming the track seeds on which the electron tracks are built upon.

The performance of the ECAL-driven method is not very well suited for non-isolated and low- $p_T$  electrons. This occurs mainly due to the fact that the super-cluster position and energy can be highly biased by the impact of overlapping particles, especially if the electron happens to be within a jet and therefore non-isolated. Also, high track multiplicity complicates the backward propagation from a super-cluster, because it can be consistent with a number of track seeds corresponding to other particles. To minimise the number of these fake seeds, the ratio between the HCAL and ECAL energy deposits ( $H/E$ ) is required to be smaller than 0.15. The HCAL towers used in the calculation of this ratio are taken within a cone of  $\Delta R = 0.3$  behind the super-cluster position. Although this helps to keep the fake seed rate under control, the efficiency for non-isolated electrons becomes rather limited. As for the low- $p_T$  electrons, the wider azimuthal spread of Bremsstrahlung photons leads to poorer reconstruction of the super-cluster, biasing its position and therefore preventing the efficient matching with a track seed.

Within the particle flow method, efficient reconstruction of non-isolated and low- $p_T$  electrons is particularly important since it affects the reconstruction of jets and missing transverse energy. Therefore, a different (‘tracker-driven’) seeding algorithm is used, which starts from reconstruction of tracks. The baseline of the CMS track reconstruction is the Kalman filter (KF) [14], which is a linear least-squares estimator based solely on Gaussian probability density functions. It is particularly suitable for muon reconstruction since it is dominated by multiple Coulomb scattering and its impact is well modelled by Gaussian fluctuations.

However, this approach usually fails for electrons because Bremsstrahlung photon emission is highly non-Gaussian. To accommodate for the resulting kinks in the electron trajectory, the Gaussian-Sum Filter (GSF) [15] is used, which is essentially a non-linear generalisation of the Kalman Filter. In this method, Bremsstrahlung energy loss is modelled by a Gaussian mixture, therefore GSF track fit provides a better estimate for the inner and outer track momentum comparing to the KF algorithm. The downside of this approach is its high CPU usage, which means it can be run on a limited number of seeds.

The GSF tracks are reconstructed upon all ECAL-driven seeds. In case of the “tracker-driven” seeds, a pre-identification based on high-purity KF tracks has been adopted. This procedure starts with the tracks reconstructed with very tight criteria, thus decreasing the fake rate yet compromising on tracking efficiency. Then an iterative-tracking strategy is carried out by means of removing hits unambiguously assigned to tracks from the previous iteration, and also progressively relaxing track seeding criteria. This approach leads to both high efficiency and low fake rate, which is crucial for low- $p_T$  and non-isolated electrons.

In the next step, track-cluster matching has to be performed. In case if electron has negligible Bremsstrahlung emission, the track is well reconstructed with the KF algorithm all the way to the ECAL internal surface, where the closest cluster is matched to the track. The corresponding cluster energy is compared with the track momentum, and if the ratio ( $E/p$ ) is close to unity, the track is selected. On the contrary, if the electron experiences a significant Bremsstrahlung emission, other track characteristics have to be exploited. In this case a selection based on the number of hits in the tracker and the  $\chi^2_{\text{KF}}$  of the KF fit is applied before running a GSF refit. Finally, the number of hits, the GSF refit  $\chi^2_{\text{GSF}}$ ,  $\chi^2_{\text{KF}}/\chi^2_{\text{GSF}}$  ratio, the energy loss measured by the track and the quality of the ECAL cluster-track matching are fed into a multivariate analysis using a Boosted Decision Trees (BDT) estimator.

Both tracker-driven and ECAL-driven seeds are used to obtain the GSF track collection of electron candidates. In the particle flow algorithm, it is necessary to link both electron and Bremsstrahlung energy deposits to the GSF track. A super-cluster is linked to a track if the extrapolated position from the outermost tracker measurement is within the boundaries of one of the ECAL cells at the expected depth of the electron shower maximum. The preshower-ECAL and ECAL-HCAL links are made in a similar way.

Another important particle flow procedure, also driven by GSF tracks, is Bremsstrahlung recovery. In order to reconstruct an electron with correctly assigned energy and momentum, it is crucial to identify all energy deposits from Bremsstrahlung photons, thus forming a super-cluster. This procedure is carried out for each tracker layer by computing a straight-line extrapolation tangent to the track, up to the calorimeter. To determine a Bremsstrahlung

photon, track-cluster linking is performed as described above. To limit the charged hadron contamination, clusters already assigned to KF tracks are not included in the calculation. Also, the distance in  $\eta$  coordinate between the extrapolation and the cluster is required to be smaller than 0.015, which helps to reduce the neutral particles background.

#### 2.4.1.1 Electron Identification

Electron reconstruction in CMS is based on a characteristic signature that electrons leave in the tracker and calorimetry systems. However, other objects like charged hadrons, jets or photon conversions can produce very similar signatures and therefore may be reconstructed as electrons. Therefore, in order to distinguish these “fake” electrons from “real” ones, a further selection has to be applied. This procedure is referred to as electron identification (or electron ID).

Initially the electron identification is performed at the final stage of the electron reconstruction process. The working points of the cuts applied are selected to be loose enough in order to satisfy most CMS analyses requirements. Afterwards, more specific (tighter) ID cuts are applied for each individual analysis, defining the working point in the trade-off between selection efficiency and fakes contamination. This will be discussed separately in the selection description for each analysis in corresponding chapters.

There are several electron identification algorithms used by various CMS analyses, and four of them are used in the analyses described in this thesis: simple cut-based (SCB ID), cuts in categories (CiC ID), particle flow (PF ID) and multi-variate analysis (MVA ID).

Simple cut-based identification is used in the High-Level Trigger, and therefore has to be as simple, fast and robust as possible. Cuts are applied on the following variables:

- $H/E$ , i.e. the ratio of hadronic energy of the HCAL towers centred at the super-cluster position of the electron in a cone of radius  $\Delta R = 0.15$ , and the super-cluster electromagnetic energy;
- $\Delta\eta_{\text{in}}$ , i.e. the difference in  $\eta$  between the extrapolated track position and  $\eta$  of the super-cluster;
- $\Delta\phi_{\text{in}}$ , i.e. the difference in  $\phi$  between the extrapolated track position and  $\phi$  of the super-cluster;
- $\sigma_{i\eta,i\eta}$ , i.e. cluster shape of the electron energy in a  $5 \times 5$  block of crystals around the seed crystal (the one with the highest energy) [16]:

$$\sigma_{i\eta,i\eta} = \sum_{5 \times 5 \text{ crystals}} (\eta_i - \eta_{\text{seed cluster}})^2 \frac{E_i}{E_{\text{seed cluster}}} \quad (2.5)$$

Although this identification method has an advantage of its simplicity, it does not show the best signal efficiency and background rejection. One of the more complex methods is the CiC ID, which exploits the categorisation of electrons. It is optimised to select electrons from different sources ( $W$ ,  $Z$  and  $J/\psi$  decays) and reject fakes from jets or conversions. In order to achieve higher efficiency, electron candidates are split into categories with different signal to background ratio, allowing to better tune the working points of the cuts. The first step in categorisation is done by division between barrel and endcap regions of the detector. Since the properties of both the tracker and calorimetry systems differ significantly in these two regions, different cuts are applied. The second step is based on the following observables:

- $f_{\text{brem}}$ , or measured bremsstrahlung fraction, defined as:

$$f_{\text{brem}} = \frac{p_{\text{in}} - p_{\text{out}}}{p_{\text{in}}} \quad (2.6)$$

where  $p_{\text{in}}$  is the initial track momentum at the vertex and  $p_{\text{out}}$  is the track momentum at the last hit;

- $E/p$ , which is the ratio of the super-cluster energy and the initial track momentum.

Three categories of electron candidates are distinguished [17]:

- “Low-Brem”:  $0.9 < E/p < 1.2 - f_{\text{brem}} < 0.12$  (barrel),  $0.82 < E/p < 1.22 - f_{\text{brem}} < 0.2$  (endcap), the fake-like region with high number of both real and fake electrons;
- “Bremming”:  $0.9 < E/p < 1.2 - f_{\text{brem}} > 0.12$  (barrel),  $0.82 < E/p < 1.22 - f_{\text{brem}} > 0.2$  (endcap), the electrons-like region with a little contamination from fakes;
- “Bad-Track”: remaining regions with a low number of real electrons.

The third and the final step in categorisation is based simply on the electron transverse energy, with the lower threshold taken as 10 GeV. In each category, the selection is performed on the basic ID variables already mentioned above for the simple cut-based identification. On top of these variables, isolation and conversion rejection variables are also used, they are discussed in Section 2.4.1.2 and Section 2.4.1.3, respectively.

The cuts within CiC ID are applied in order to maximise the signal to background ratio. Depending on the needs of different analyses, nine levels of cut severity are implemented: *VeryLoose*, *Loose*, *Medium*, *Tight*, *SuperTight*, *HyperTight(1-4)*. Each step decreases the fake rate by about a factor of two for electrons with  $E_T > 20$  GeV [17]. The CiC ID is used as a primary electron identification method in the Top Mass analysis (Chapter 4).

Particle flow ID is the final step of the particle flow electron reconstruction. Following the particle flow concept, it uses information from all the CMS sub-detectors obtained in previous reconstruction steps to build new observables for electron identification. The complete list of PF ID observables is given below:

- $p_T$  and  $\eta$  of the GSF track;
- GSF  $\sigma_{p_T}/p_T$ , transverse momentum resolution of the GSF track;
- $\#hits_{KF}$ , number of reconstructed KF track hits;
- $\chi^2_{GSF}$  and  $\chi^2_{KF}$ , GSF and KF goodness-of-fits;
- $\Delta\eta$ : distance in  $\eta$  between the position of the cluster and the extrapolated position of the GSF track;
- $\sigma_{\eta\eta}$ , cluster shape (Equation 2.5) of the ECAL cluster linked to the GSF track;
- $H/(H+E_e)$ , hadron fraction of the shower, where  $H$  is the energy of the hadron cluster linked to the GSF track;
- $(E_e + \sum E_\gamma)/p_{in}$ , ratio between the super-cluster energy and the inner track momentum;
- $E_e/p_{out}$ , ratio between the electron cluster energy and the track outer-momentum;
- $\sum E_\gamma/(p_{in} - p_{out})$ , the ratio between the Bremsstrahlung photon energy as measured by ECAL and by the tracker;
- *EarlyBrem*, flag of  $(E_e + \sum E_\gamma) > p_{in}$  inequality, corresponding to an electron emitting an “early” Bremsstrahlung photon, i.e. before it has crossed at least three tracker layers;
- *LateBrem*, flag of  $E_e > p_{out}$  inequality, corresponding to an electron emitting a “late” Bremsstrahlung electron, when the ECAL clustering is not able to disentangle the overlapping electron and photon showers;
- $f_{brem} = (p_{in} - p_{out})/p_{in}$ , Bremsstrahlung fraction (Equation 2.6).

All these variables are combined into a single discriminator by a multivariate analysis technique (BDT method), which has been trained on signal and background Monte Carlo samples. PF ID is a relatively loose identification method, since it has to satisfy the needs of all analyses using particle flow collections.

Finally, MVA electron identification is used in top cross sections analysis on 2012 data. It is another multivariate analysis technique, optimised to select isolated electrons from W and Z decays. The variables used in MVA ID are also combined into a single discriminator, they are largely similar to the ones used in PF ID:

- $p_T$  and  $\eta$  of the GSF track;
- $\#_{\text{hits}_{\text{KF}}}$ , number of reconstructed KF track hits;
- $\chi^2_{\text{GSF}}$  and  $\chi^2_{\text{KF}}$ , GSF and KF goodness-of-fits;
- $\Delta\eta$ : distance in  $\eta$  between the position of the cluster and the extrapolated position of the GSF track;
- $\Delta\phi$ : distance in  $\phi$  between the position of the cluster and the extrapolated position of the track;
- $\Delta\eta_{vtx}$ : distance in  $\eta$  between the position of the cluster and the position of the GSF track at vertex;
- $\Delta\phi_{vtx}$ : distance in  $\phi$  between the position of the cluster and the position of the GSF track at vertex;
- $\sigma_{\eta\eta}$ , cluster shape in  $\eta$ ;
- $\sigma_{\phi\phi}$ , cluster shape in  $\phi$ ;
- $\eta$  width of the super-cluster;
- $\phi$  width of the super-cluster;
- $H/E$ , ratio of HCAL and ECAL cluster energy;
- $E_{\text{super-cluster}}/p_T$ , ratio between the super-cluster energy and the track momentum;
- $1/E_{\text{super-cluster}} - 1/p_T$ , difference between inverse super-cluster energy and inverse track momentum;
- $E_e/p_{\text{out}}$ , ratio between the electron cluster energy and the track outer-momentum;
- $1 - E_{1\times 5}/E_{5\times 5}$ , where  $E_{1\times 5}$  is the energy in the central  $1 \times 5$  strip of the  $5 \times 5$  electron cluster and  $E_{5\times 5}$  its total energy;

- $E_{3\times 3}/E_{\text{super-cluster, raw}}$ , ratio of the energy of a cluster of  $3 \times 3$  and the uncorrected (raw) energy of the super-cluster;
- $E_{\text{PS}}/E_{\text{super-cluster, raw}}$ , ratio of the energy in the preshower detector and the raw super-cluster energy (only in the endcap region).

#### 2.4.1.2 Electron Isolation

Isolation is an observable that allows to distinguish prompt electrons (i.e. the ones from W and Z decays) from jets faking electrons and electrons within jets. It is essentially a measure of activity around the particle. In CMS, there are two different ways of quantifying isolation: detector-based and particle-based. The detector-based isolation is defined separately for each detector sub-system (tracker, ECAL and HCAL):

- Tracker isolation, calculated as the sum of transverse momenta of all track within a cone of  $\Delta R = 0.3$  around the electron, excluding the electron momentum itself;
- ECAL isolation, i.e. the sum of transverse energy of all ECAL clusters within a cone of  $\Delta R = 0.3$  around the super-cluster position. The footprint of the original electron is also removed;
- HCAL isolation, defined as the sum of transverse energy of all HCAL towers within a cone of  $\Delta R = 0.3$  centred at the super-cluster position.

Particle-based isolation exploits the particle flow information: it is calculated as the sum of transverse energy of all PF particles in the cone of  $\Delta R = 0.3$  around the electron. Both detector-based and particle-based isolation definitions are often normalised to the electron transverse momentum (or energy in case of calorimeter isolation) in order to improve signal efficiency. The normalised sum of the tracker, ECAL and HCAL isolation variables is referred to as detector-based relative isolation (or reliso), similarly normalised particle-based isolation is called PF reliso. These quantities are crucial in various methods to estimate the QCD background contribution for many analyses, and will be used for both electrons and muons throughout this thesis.

#### 2.4.1.3 Identification of photon conversions

Interacting with detector material, photons can convert into electron-positron pairs. Due to the large amount of material budget in the tracker (Figure 2.4), especially in the endcap region, there is a high chance of conversions to happen. The resulting electrons can

successfully fake prompt signal electrons, passing all the identification criteria and appearing isolated if the initial photon was isolated, too. Therefore conversions constitute a large proportion of the QCD background to top signal. The electron-positron pair may not be symmetrical in transverse momenta, therefore a veto on a second electron is not sufficient to reject such electrons and other conversion identification criteria are necessary.

The simplest method is based on counting the number of missing hits in the tracker. Conversions are most likely to happen at some distance from the interaction point, essentially anywhere between the point of photon production and the end of the tracker. Therefore electrons produced in such conversions are likely not to traverse through all the pixel layers. To separate these electrons from the prompt ones, a cut on the number of missing layers can be used. However, this approach fails if conversion occurs in the beam pipe, or if the reconstructed electron track is paired with unrelated hits in the tracker, making it look like a prompt electron.

A slightly more sophisticated method is a partner track method. It is based on geometrical cuts on *dist* and *dcot* variables, which refer to the distances between tangent points of any two tracks in the  $r - \phi$ -plane and  $r - z$ -plane, respectively:

$$dist = ||\vec{r}_1 - \vec{r}_2| - r_1 - r_2| \quad (2.7)$$

$$dcot = \left| \frac{1}{\tan \theta_1} - \frac{1}{\tan \theta_2} \right| \quad (2.8)$$

Here  $\vec{r}_{1,2}$  are radial vectors of the two tracks,  $r_{1,2}$  are the track radii and  $\theta_{1,2}$  – angles between the tracks and the beam pipe. Apart from the cuts on *dist* and *dcot* variables, tracks should have opposite charge in order to be considered to come from a photon conversion.

These two methods for conversion identification were used in the top mass analysis on 2011 data. For the top cross sections analysis on 2012 data, along with the number of hits method, a more advanced technique called vertex fit was used. It is essentially a full vertex fit of all pairs of tracks, with the selection being made on the fit probability. The method benefits from the full use of track uncertainties and covariances, combining all information into a single discriminator. However, increased complexity of this technique makes it more CPU-intensive than the partner track method.

### 2.4.2 Muon Reconstruction

A good muon reconstruction and identification was one of the main CMS design requirements. Initially, muons are reconstructed independently in the tracker (“tracker tracks”) and in the muon system (“standalone muon tracks”) using the Kalman Filter technique [14].

Based on these objects, two different reconstruction methods are used: [18]

- Global muon reconstruction. Each standalone muon track reconstructed in the muon system is matched with the tracker track by propagating onto a common surface. A global fit is performed on the combined collection of hits from both tracks, and the resulting muon is referred to as a *global muon*.
- Tracker muon reconstruction. All tracker tracks above certain threshold ( $p_T > 0.5$  GeV,  $p > 0.5$  GeV) are extrapolated to the muon system, taking into account possible energy losses, magnetic field and multiple Coulomb scattering. If the track matches at least one muon segment in the muon system, i.e. a short track stub made of DT or CSC hits, it is considered a *tracker muon*.

Global muon reconstruction has a high efficiency for high- $p_T$  muons, penetrating through more than one muon station. On the contrary, tracker muon reconstruction is more efficient for low- $p_T$  muons ( $p_T < 5$  GeV), as it requires just a single muon segment. Since the  $t\bar{t}$  analyses described in this thesis have a semileptonic signature with exactly one energetic muon in the final state (in case of the muon channel), only the global muon reconstruction method is used, as its momentum resolution at high transverse momenta benefits from both the tracker and the muon system (Figure 2.9).

Following the reconstruction, the quality of the muon objects is verified by applying identification criteria. This is done in order to suppress hadronic punch-through, muons from decays in flight and cosmic muons. Selection is applied on the following observables:

- normalised  $\chi^2$  ( $\chi^2/\text{number of degrees of freedom}$ ) of the global muon fit;
- number of muon chamber hits in the global muon fit;
- number of muon stations with muon segments;
- transverse impact parameter  $d_{xy}$  (closest approach of the track to the primary vertex);
- longitudinal distance  $d_z$  of the tracker track w.r.t. the primary vertex;
- number of hits in the pixel detector;
- number of hits in the tracker layers.

The muon identification also exploits the CMS particle flow event reconstruction, making use of information coming from all sub-detectors. Particle flow muons (PF muons) are identified by imposing selection on all muon candidates reconstructed with the global muon method. This selection was optimised for identification of muons in jets, minimising the fake rate from misidentified charged hadrons, which is crucial for correct reconstruction of jets and missing transverse energy (described in more detail in the following sections). Depending on the analysis specifics, isolation requirements can also be applied; isolation definitions closely follow the ones for electrons (Section 2.4.1.2).

### 2.4.3 Jet Reconstruction

Hadronisation of quarks and gluons leads to production of narrow cones of particles moving in approximately one direction, called jets. This happens due to colour confinement, as particles carrying a colour charge cannot exist in free form, they have to fragment into hadrons before they can be detected directly. Therefore, to measure the initial parton's momentum and energy, all these particles must be combined into jets.

The CMS particle flow algorithm implies reconstruction of individual particles (charged and neutral hadrons, electrons, muons, photons) before combining (or clustering) them into jets. A few different jet clustering techniques exist, but the one used predominantly in CMS and exclusively in this work is anti- $k_t$  algorithm [19], which defines the distance between constituent particles as:

$$d_{ij} = \min \left( \frac{1}{k_{t,i}^2}, \frac{1}{k_{t,j}^2} \right) \frac{\Delta_{i,j}^2}{R^2} \quad (2.9)$$

where  $\Delta_{i,j}^2 = (y_i - y_j)^2 + (\phi_i - \phi_j)^2$ ,  $k_{t,i/j}$ ,  $y_{i/j}$  and  $\phi_{i/j}$  are respectively transverse momenta, rapidities and azimuth angles of particles  $i/j$ , and  $R$  is the radius parameter.

The clustering proceeds by identifying the smallest distances between particles and recombining them until all jets are formed and no particles are left. An event typically has a few well-separated hard (high- $p_T$ ) particles and a large amount of soft (low- $p_T$ ) particles. The distance  $d_{1i}$  between a hard particle 1 and a soft particle  $i$  will be fully determined by the transverse momentum of the hard particle and the  $\Delta_{1i}$  separation. On the other hand, distance between soft particles with similar separation will be substantially larger. Therefore, soft particles tend to cluster around the hard ones before they cluster amongst themselves. On the output, the anti- $k_t$  algorithm forms conical jets with boundaries resilient to soft radiation.

#### 2.4.3.1 Jet Energy Corrections

Due to non-linear and non-uniform response of the calorimetry systems, jets reconstructed using detector inputs typically have energies different to the ones of corresponding Monte

Carlo particle jets (or generator jets), reconstructed by clustering the four-momenta of all stable particles generated in Monte Carlo simulation. Therefore some mapping procedure is necessary. Corrections applied to reconstructed jets in order to translate the measured energy to the true particle (or parton) energy are referred to as jet energy corrections, or JEC.

CMS has adopted a factorised approach of applying jet energy corrections, meaning that each correction level takes care of a different effect. The set of corrections is applied sequentially, i.e. the output of each step is the input to the next one. Essentially, each level of correction is a scaling of a jet four-momentum, with a scale factor depending on various parameters of the jet, typically pseudorapidity and transverse momentum. Currently, the correction levels go as follows:

- L1 Offset correction;
- L2 Relative ( $\eta$ ) correction;
- L3 Absolute ( $p_T$ ) correction;
- L4 EMF (electromagnetic energy fraction) correction;
- L5 Flavour correction;
- L6 UE (underlying event) correction;
- L7 Parton correction.

The goal of the L1 correction is subtraction of pile-up and electronic noise contributions from the jet energy. The energy from pile-up vertices can be deposited in calorimeters since the additional proton-proton collisions occur close enough in time to the hard scattering process. Electronic noise in calorimeter readouts also creates additional energy offset which needs to be corrected for. The scale factors for L1 correction are derived using data-driven methods (zero-bias collisions).

The relative L2 correction is designed to flatten the jet response in pseudorapidity. These corrections are extracted with respect to the barrel region in bins of  $p_T$ , and therefore are uncorrelated with the following L3 corrections. Both Monte Carlo and data-driven (dijet balance) methods are used to derive the L2 scale factors. Once a jet is corrected for  $\eta$  dependence, it is corrected back to the particle level by applying absolute L3 correction. The goal of this correction is to flatten the jet response in transverse momentum. Derivation

of L3 scale factors is done by using Monte Carlo truth information, or data-driven  $Z/\gamma$ +jet balance techniques.

The L4 EMF correction takes care of variations in jet response with respect to electromagnetic energy fraction, i.e. the fraction of energy deposited in ECAL by hadrons. Although it has been shown that EMF-dependent correction in three parameters ( $p_T$ ,  $\eta$ , EMF) can improve the observed jet resolution, this correction is optional for most CMS analyses and hasn't been used in this work.

Another optional correction is the L5 flavour correction, which is intended to correct the jets depending on the flavour of initiating partons, i.e. light quarks, heavy ( $b$  and  $c$  quarks), or gluons. Jets originating from heavy quarks are different from light quark jets in a few aspects. The average number of charged hadrons within  $b$ -jets is higher than the one for the light jets, although the average momentum of charge hadrons is smaller. Moreover,  $b$ -hadrons have a larger branching fraction into semileptonic decays with neutrinos in the final state that cannot be detected. All these factors may lead to a substantially smaller charged hadron energy fraction comparing to the light jets, therefore the average calorimeter energy response is lower for the  $b$ -jets. These effects are meant to be taken care of by the L5 correction, which is derived either from Monte Carlo or data-driven  $t\bar{t}$  events.

Underlying event [20] refers to all the activity in the proton-proton collision apart from the process of interest, i.e. the hard scattering process. It includes particles coming from additional parton interactions (or multiple scatterings) and beam remnants. Depending on the analysis goals, underlying event may also include initial and final state radiation (ISR and FSR) which represent the soft gluon radiation before and after the hard scattering process, respectively. Underlying event activity results in production of additional soft jets which can bias the jet energy measurement. Optional L6 UE correction attempts to mitigate such bias. However, it was not used in the analyses described in this thesis, since underlying event is an intrinsic part of proton-proton interactions effectively modelled by Monte Carlo simulation. In particular, the top cross sections analysis is designed to be very sensitive to such higher-order processes, and not taking them into account may result in compromising the discriminating power between Monte Carlo generators.

Finally, L7 parton correction attempts to correct the jet energies back to the parton level. It is also optional and was only used in the top mass analysis. The correction factors were derived from Monte Carlo simulations by comparing the generator jets momenta to their matched partons.

Due to the fact that CMS simulation is not perfectly tuned to the data yet, additional residual L2L3 corrections are applied in order to achieve better agreement between data and

simulation. It is essentially a small residual  $\eta$ - and  $p_T$ -dependent calibration applied exclusively to data, which will remain in place until CMS develops a perfectly tuned simulation reproducing the data features out of the box.

#### 2.4.3.2 Particle Flow Jet Identification

To ensure the quality of the jets, a final identification criteria are applied to all jet objects. Particle flow jet identification (or PF jet ID) is used to reduce the noise and rate of electrons reconstructed as jets. The cuts are applied on the following observables :

- number of constituent particles;
- NHF (neutral hadron energy fraction);
- CHF (charged hadron energy fraction);
- NEF (neutral electromagnetic energy fraction);
- CEF (charged electromagnetic energy fraction);
- NCH (number of charged particles).

A rather loose identification cuts (referred to as “loose PF jet ID”) were used in this work: a requirement of more than one constituent particle in the jet,  $\text{NHF} < 0.99$ ,  $\text{NEF} < 0.99$ . In addition, for a pseudorapidity region of  $\eta < 2.4$  the following requirements are imposed:  $\text{CEF} < 0.99$ ,  $\text{CHF} > 0$  and  $\text{NCH} > 0$ .

#### 2.4.3.3 b-tagging

The identification of jets originating from b-quarks, or b-tagging, is one of the most important tools in top quark physics, capable of significantly decreasing the background contamination of signal processes. A variety of b-tagging algorithms has been developed by CMS [21], and the one that was used in this work is called combined secondary vertex (CSV) algorithm. It is based on the fact that b-hadrons have a significant lifetime ( $\sim 10^{-12}$  s) and can travel a distance of a few centimetres before decaying. Therefore jets originating from b-quarks are likely to have a secondary vertex located at a considerable distance from the primary vertex, which can be used as an efficient discriminator between light jets and b-jets. In order to maximise its efficiency, apart from the secondary vertex information the CSV algorithm also exploits the track-based lifetime information. The following variables are used in the algorithm:

- number of tracks in the jet;

- number of tracks at the secondary vertex;
- secondary vertex category;
- secondary vertex invariant mass;
- ratio of the total track energy at secondary vertex with respect to all tracks in the jet;
- pseudorapidities of tracks at secondary vertex with respect to the jet axis;
- impact parameter significance (i.e. ratio of IP to its uncertainty) of the first track that increases the invariant mass above the charm threshold of 1.5 GeV (tracks are ordered by IP significance; the mass of the system is recalculated after adding each track);
- impact parameter significances of each track in the jet.

All these observables are combined into a single discriminator, the “medium” working point of which is used in this work. It provides  $\sim 70\%$  b-tagging efficiency for mis-tag rate of approximately 1 % [21].

#### 2.4.4 Missing Transverse Energy

Due to the energy and momentum conservation, the sum of transverse momenta and energy of all particles in the final state of proton-proton collisions is expected to be zero. However, some particles can escape the detector without being reconstructed, therefore creating the imbalance in transverse momentum which is referred to as missing transverse energy, defined as:

$$\vec{E}_T^{\text{miss}} = - \sum_i \vec{p}_T^i \quad (2.10)$$

where  $\vec{p}_T^i$  are the transverse momentum vectors of all reconstructed particles. The modulus of  $\vec{E}_T^{\text{miss}}$  vector is denoted by  $E_T^{\text{miss}}$ .

Accurate reconstruction of missing transverse energy is crucial for precise measurements of Standard Model processes with neutrinos in the final state. Top quark pair semileptonic decay is one of such processes, therefore analyses covered in this thesis implicitly (top mass) or explicitly ( $t\bar{t}$  cross section with respect to  $E_T^{\text{miss}}$ -related variables) rely on efficient reconstruction of  $E_T^{\text{miss}}$ . Misidentification and misreconstruction of any visible particles in the event contribute to  $E_T^{\text{miss}}$  measurement, therefore it is a rather demanding task.

Just like in the case of leptons and jets, particle flow was used for  $E_T^{\text{miss}}$  reconstruction: in equation 2.10 the transverse momentum vectors are of the particles reconstructed using particle flow algorithm. This procedure gives so-called raw  $E_T^{\text{miss}}$  on the output. The raw

$E_T^{\text{miss}}$  is systematically different from true  $E_T^{\text{miss}}$ , which denotes the transverse momentum carried by invisible particles. This happens mainly due to non-compensating nature of the calorimeters, effects of pile-up, noise, etc. Therefore, a set of corrections is applied:

- Type-0, which corrects  $E_T^{\text{miss}}$  for pile-up;
- Type-I, a propagation of jet energy corrections (Section 2.4.3.1) to  $E_T^{\text{miss}}$ ;
- $xy$ -shift correction, reducing the  $E_T^{\text{miss}} \phi$  modulation.

The causes of systematic  $E_T^{\text{miss}} \phi$  modulation include detector misalignment, beam spot displacement, inactive calorimeter cells and anisotropic detector response.  $xy$ -shift correction mitigates these effects, making the measured  $E_T^{\text{miss}}$  distribution closer to true  $E_T^{\text{miss}}$  distribution which is flat in  $\phi$  because of the rotational symmetry of the collisions.

All these corrections were applied to missing transverse energy in the analyses described in this work. However, the  $xy$ -shift correction was not applied in the top mass analysis, since it was not available at the time. As this analysis is not particularly sensitive to  $E_T^{\text{miss}}$ , the  $\phi$  modulation is not expected to affect the top mass measurement and its resolution.

## 2.5 Summary

In this chapter, the Large Hadron Collider (LHC) has been introduced to the reader. The CMS experiment including all its subsystems has been described in detail, the overview of the CMS computing model and analysis software have been shown. Reconstruction and identification methods of various analysis objects including particle flow algorithm have been discussed in detail.



## 3. High level triggers for Top Physics

The LHC is often referred to as a top quark factory, producing a  $t\bar{t}$  pair nearly each second of its nominal operation. While the production rate of  $\approx 1$  Hz seems manageable in terms of recording the data, it is significantly complicated by background processes with similar signatures occurring at much higher rates.

The trigger is the starting point of any physics event selection process, and therefore is clearly important for any physics analysis. As it was mentioned in Section 2.2.6, the CMS L1 trigger rate is limited to  $\sim 100$  kHz. In order to meet the data recording constraints of approximately  $300 \text{ MB s}^{-1}$ , this rate is further reduced down to  $\sim 300$  Hz, which is done by the HLT system. The total rate budget has to be shared between various physics analysis groups (e.g. Top, Higgs, Exotica, etc). Corresponding allocations are determined by CMS trigger coordination according to CMS physics goals, and can be a matter of serious debate.

During the LHC operation in 2011 and 2012 under conditions of gradual increase of instantaneous luminosity and pile-up, but very limited rate budget, trigger developers constantly tackled the challenge of finding the best compromise between growing rates and maintaining reasonable signal acceptance. While the simplest approach is tightening the cuts on physical quantities like lepton or jet transverse momenta, it is not favourable since it lowers the number of stored signal events and decreases the phase space which is crucial for new physics searches as well as Standard Model precision measurements. Therefore, development of more efficient algorithms allowing to keep high level of acceptance for signal events, whilst effectively rejecting background events, is the most preferable solution. This can often be achieved by increasing the level of approximation of the online (HLT) object reconstruction, making the algorithms closer to their sophisticated offline counterparts. However, it leads to a higher execution time, which is limited by computing resources available for HLT reconstruction. Hence, the CPU timing is another major constraint faced by the HLT developers.

This chapter covers the author’s contribution to development and verification of High-Level Triggers for top physics with semileptonic signature, where one of the W bosons decays into an electron and a neutrino. The description of top triggers, efficiency measurement, validation of jet energy corrections and pile-up subtraction applied at the HLT level are discussed in relevant sections of this chapter.

### 3.1 Level-1 triggers

The CMS L1 trigger [22] is built of custom electronics processing the data from the calorimeters and the muon system. It is the first filtering stage of any physics selection, therefore it has high requirements on efficiency and phase space. The input rate of events, i.e. the beam crossing frequency of  $\sim 5 \times 10^8$  Hz must be reduced down to  $\sim 100$  kHz, being a major challenge considering a gradual increase of instantaneous luminosity throughout the LHC operation. The rate budget of  $\sim 100$  kHz is shared between several L1 triggers corresponding to different physics objects being present in the event. For top physics signatures with a single electron in the final state, the following triggers have been used:

- L1\_SingleEG\_18
- L1\_SingleEG\_20
- L1\_SingleEG\_22

These triggers are referred to as electron/photon ( $e/\gamma$ , or EG) triggers. They scan the ECAL for groups of crystal cells ( $4 \times 4$ ) with a cumulative energy above the threshold that is suggested by the name of the trigger (18 GeV, 20 GeV and 22 GeV). The choice of threshold is determined by sustainability of the trigger rate at a given instantaneous luminosity.

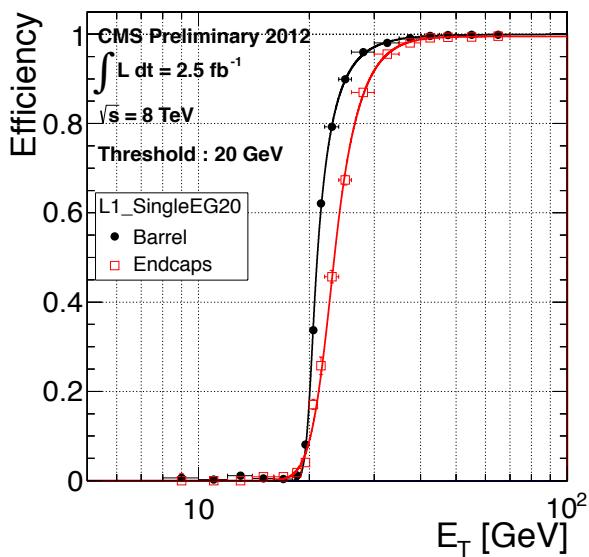


Figure 3.1: Efficiency of the 20 GeV  $e/\gamma$  trigger as a function of offline  $E_T$ , shown separately for barrel and endcap regions [23].

The efficiency of the L1\_SingleEG\_20 trigger as a function of the offline-reconstructed electron transverse momentum is shown on Figure 3.1. The L1 trigger decisions are then used as an input “seeds” to the HLT system, which is described in the following section.

## 3.2 High-level triggers for top physics

The High-Level Trigger [24] is the crucial part of CMS event selection process. As it was mentioned in Section 2.2.6, it is based on software algorithms running on the Event Filter Farm, i.e. a large cluster of commercial CPUs. The HLT reconstruction, often referred to as online reconstruction, is implemented in the same software framework (CMSSW) which is used for offline reconstruction, and the algorithms can be very similar. However, the key difference between online and offline reconstruction is the running time. Since the HLT selection has to be performed in real time, it imposes a significant constraint on computing resources, enforcing to compromise on robustness and efficiency of online algorithms. With an exception of small samples for performance monitoring, data rejected by the HLT is lost irrevocably. Therefore, correct and efficient operation of the HLT is of major importance for CMS physics programme.

Modular structure of CMSSW provides high flexibility in use of selection and reconstruction algorithms, allowing their continuous optimisation according to changes in physics needs and data-taking conditions. Various modules account for reconstruction of different physics objects and their matching with L1 objects, filtering, logging, monitoring, etc. All these modules are grouped into so-called trigger “paths”, ultimately giving trigger decisions on whether to accept an event or not. A typical example of a trigger path used for top physics is an electron-plus-jets trigger, which requires the presence of isolated electron and at least three energetic jets in the event.

There are two types of trigger paths: prescaled and unprescaled. Prescaling is a method of reducing the trigger rate by only recording a fraction of events that pass the trigger selection. For example, if the prescale value of a particular trigger is 10, only one out of ten events that fire the trigger will be actually recorded. Usually prescaled triggers are used for performance monitoring and background estimation as they can impose a much looser selection criteria yet still have manageable rate. However, such triggers are unfavourable for signal selection, as they greatly reduce signal acceptance.

A set of trigger paths and their prescale factors is combined in the HLT configuration, referred to as the HLT menu or table [25]. Since the start of data-taking in 2010, the CMS trigger coordination adopted a “trigger train” model, implying the schedule of regular deadlines for updates in the trigger menu. These deadlines are usually imposed fortnightly,

according to changes in instantaneous luminosity or other alterations in data-taking conditions. In order for a trigger path to be included in the trigger menu, it has to be implemented in the HLT configurations database, validated by measuring the trigger rate and signal efficiency, tested for CPU timing and finally approved by the trigger studies group.

As it was mentioned before, the top quark pair decay covered in this work has a semileptonic signature, containing an electron, at least four jets and a neutrino in form of  $E_T^{\text{miss}}$ . Due to the nature of the decay, all these objects are highly energetic and can be triggered on with rather high thresholds on transverse energy. During the start-up year of 2010 when the instantaneous luminosity went up from  $\sim 10^{27} \text{ cm}^{-2}\text{s}^{-1}$  to  $\sim 10^{31} \text{ cm}^{-2}\text{s}^{-1}$ , top analyses with the signature of interest used the single electron trigger, requiring just one electron with certain isolation and transverse momentum criteria. However, by extrapolating the trigger rates on higher luminosities foreseen in the following years it became obvious that the single electron trigger rate would quickly become too high for the rate budget restrictions at the time. Therefore, alternative ways of adding other objects from the  $t\bar{t}$  signature were explored, naturally leading to electron-plus-jets trigger solution.

A typical electron-plus-jets trigger consists of a standard electron module, a jet module, a cleaning module to remove overlap between jet and electron collections, and a jet multiplicity filter. The electron module is constructed from the following sub-modules:

- L1 object and ECAL super-cluster matching;
- ECAL transverse energy filter;
- ECAL electron ID and isolation filter;
- HCAL electron ID and isolation filter;
- tracker electron ID and isolation filter.

To minimise the total running time, all the sub-modules are ordered by speed, starting from the fastest. The matching module finds the ECAL super-cluster closest to the L1 object in  $\eta$  and  $\phi$  dimensions. If the super-cluster is located within a certain  $\eta$  and  $\phi$  range, the event is passed to the next module, otherwise it is rejected.

The following (ECAL) module calculates the super-cluster transverse energy ( $E_T$ ), and in case if it is above 25 GeV, the electron ID and isolation criteria are imposed for both ECAL and HCAL (see Section 2.4.1). The working points of identification and isolation and their naming conventions are shown in Table 3.1.

In the main trigger path used for signal rather than background estimation, Very Tight (VT) working point for calorimeter identification (CaloID) and Tight (T) working point for calorimeter isolation (CaloIso) were used. This was motivated by consistency with similar criteria for the single electron trigger used in 2010.

For the events that pass all described criteria in the electron module, a simplified algorithm for offline iterative tracking is applied. This algorithm has a faster running time at the expense of worse resolution. In order to decrease the CPU usage, the Combinatorial Track Finder (CTF) [26] is used instead of the GSF one, and the tracking is confined to a small region around the electron. The events are required to pass the tight criteria of the tracker electron ID (TrkId) which is calculated at this stage. Finally, the tight working point of the tracker isolation (TrkIso) is applied, where the track  $p_T$  is summed over all tracks within  $\Delta R = 0.3$  cone around the electron track, excluding the electron itself. Provided that event passes all the electron sub-modules, the jet module is then executed.

The jet reconstruction is performed with the anti- $k_t$  algorithm (see Section 2.4.3) with a radius parameter of  $R = 0.5$ . Historically, at the early stages of LHC operation, CMS analyses mostly used calorimeter jets, i.e. jets reconstructed using calorimeter information exclusively. Therefore, initially the trigger jet module used calorimeter jet collection, which provided fast reconstruction yet worse resolution hence lower trigger efficiency comparing to PF jets which were introduced later on during 2011 data-taking. In order to work online, PF jet reconstruction was modified to be compatible with the simplified tracking algorithm. Despite the fact that running the jet module with PF jet reconstruction is highly CPU-intensive and therefore can not be used as an unprescaled stand-alone trigger, it works effectively together with the electron module since it reduces the input event rate considerably. Both calorimeter and PF jet versions of the module impose the 30 GeV cut on transverse momentum on the jets as well as the pseudorapidity cut of  $|\eta| < 2.6$ . Subsequently, the jet

Working point	CaloID	CaloIso	TrkId	TrkIso
VeryLoose (VL)	$H/E < 0.15$ (0.10) $\sigma_{inj} < 0.024$ (0.040)	ECAL iso/ $E_T < 0.2$ (0.2) HCAL iso/ $E_T < 0.2$ (0.2)	$\Delta\eta < 0.01$ (0.01) $\Delta\phi < 0.15$ (0.10)	track iso/ $p_T < 0.2$ (0.2)
Loose (L)	$H/E < 0.15$ (0.10) $\sigma_{inj} < 0.014$ (0.035)			
Tight (T)	$H/E < 0.10$ (0.075) $\sigma_{inj} < 0.011$ (0.031)	ECAL iso/ $E_T < 0.125$ (0.075) HCAL iso/ $E_T < 0.125$ (0.125)	$\Delta\eta < 0.008$ (0.008) $\Delta\phi < 0.07$ (0.05)	track iso/ $p_T < 0.125$ (0.125)
VeryTight (VT)	$H/E < 0.05$ (0.05) $\sigma_{inj} < 0.011$ (0.031)			
WP80	$H/E < 0.10$ (0.05) $\sigma_{inj} < 0.01$ (0.03)	ECAL iso/ $E_T < 0.15$ (0.1) HCAL iso/ $E_T < 0.1$ (0.1)	$\Delta\eta < 0.007$ (0.007) $\Delta\phi < 0.06$ (0.03)	track iso/ $p_T < 0.05$ (0.05)

Table 3.1: Naming conventions for electron trigger working points. The given values are for the barrel region of the detector and the values in brackets are for the endcap region.

collection is passed to the cleaning module.

Within CMS, electron and jet collections are reconstructed independently, and since they often leave similar footprints in the detector, these collections can overlap, causing potential double-counting. The electron-jet cleaning module removes electron candidates found by the electron module from the collection of jet candidates found by the jet module. A jet is removed from the jet collection if there is an electron within a cone of  $\Delta R = 0.3$  around the jet axis. However, occasionally a genuine jet can be incorrectly identified as an electron. In this case the fake electron can be removed from the jet collection, biasing the performance of the jet multiplicity filter. To tackle this issue, separate jet collections are created for all electrons in the event, where jets are only cleaned from additional (non-signal) electrons. All these jet collections are passed to the following jet multiplicity filter, which accepts an event if at least one collection contains a desired number of cleaned jets.

Electron-plus-jets triggers successfully functioned throughout 2011 and 2012 years of data taking. During 2011, they were used in the top mass analysis and differential cross section analysis with respect to missing transverse energy, described in this thesis. The 2012 cross section analysis used a single electron trigger with a lepton  $p_T$  threshold of 27 GeV, reasonably tight (WP80) lepton isolation requirements, but no specific jet requirements. In contrast to 2011 running, this trigger was decided to be unprescaled in the trigger menu for the whole period of data-taking in 2012, regardless of having a significant rate. It is favourable for many analyses as it is much more straightforward in terms of calculating efficiency and acceptance scale factors, also reducing the complexity of estimating the systematic errors associated with the triggering.

### 3.3 Trigger efficiency measurement

One of the most important characteristics of a trigger path is its efficiency. In a broad sense, selection efficiency can be defined as a conditional probability that a single event passes the selection, given all other conditions (detector configuration, preselection, etc). In the context of the trigger, efficiency highly depends on the offline selection it is measured with respect to. As a matter of fact, different studies can apply different offline selections, therefore the same trigger may have different efficiencies for each of them.

The trigger efficiency can be measured in both data and Monte Carlo simulation using various methods. For all of them, the study has to be performed on some initial preselected dataset. Ideally, this preselection needs to be unbiased with respect to selection imposed by the trigger and offline selection. However, in reality it implies running on vast amounts of data, with very little number of events satisfying the trigger selection, leading to insufficient

statistics or non-feasible computing resources needed for the study. Therefore, some other trigger is used for preselecting the initial dataset, somewhat correlated to the trigger under study in order to obtain reasonable statistics. To correctly measure the efficiency of a trigger path, this correlation has to be taken into account.

The efficiency of a trigger  $A$  can be written as the following conditional probability [27]:

$$\epsilon_A = P(A|\vec{x}, T, D) \quad (3.1)$$

where  $\vec{x}$  are reference quantities used for triggering (e.g. lepton  $p_T$ ), and  $T$  is an offline selection, and  $D$  is a combination of all other factors like detector effects. Assuming that these factors, along with the offline selection and quantities used for triggering are fixed, we can denote  $\epsilon_A = P(A)$  for brevity.

If  $B$  is the trigger used for preselection, then

$$P(A, B) = P(A|B) \cdot P(B) = P(B|A) \cdot P(A) \quad (3.2)$$

Therefore, as it immediately follows from Bayes' theorem,

$$P(A) = \frac{P(A|B) \cdot P(B)}{P(B|A)} \quad (3.3)$$

Here the conditional probability  $P(A|B)$  can be estimated by taking the ratio of the numbers of events that pass the triggers  $A$  and  $B$ , given that they all pass offline selection. Efficiency of the auxiliary trigger  $\epsilon_B = P(B)$  can be estimated from data by using a dataset obtained with looser (“minimum bias”) selection. Unfortunately, conditional probability  $P(B|A)$  can not be measured using real data, as by the nature of the study the trigger  $A$  only considers events preselected by the trigger  $B$ . However, the trigger  $B$  is usually chosen in such way that it is safe to assume that  $P(B|A) = 1$  with negligible uncertainty. This is the case if the preselection trigger criteria are considerably looser than those of the trigger under study. To summarise, the trigger efficiency can be measured as:

$$\epsilon_A = P(A|B) \cdot P(B) = \frac{N_A}{N_B} \cdot \epsilon_B \quad (3.4)$$

where  $N_{A(B)}$  is the number of events that pass the offline selection and fire the trigger  $A$  ( $B$ ).

In case of electron plus jets triggers, the efficiency can be factorised in contributions of

leptonic and hadronic parts of the trigger:

$$\epsilon_{ele+jets} = \epsilon_{ele} \cdot \epsilon_{jets} \quad (3.5)$$

This method can be used under assumption that the probability of finding an electron in the event is independent of the presence of the jets in the event, i.e. leptonic and hadronic “legs” of the trigger are uncorrelated. Although such correlations do exist, it has been shown that most of them are negligible and factorisation method provides a meaningful estimate of the trigger efficiency [28].

Trigger efficiency can be parametrised by various physical quantities of the triggered objects, e.g. reconstructed jet pseudorapidity and transverse momentum. As the background spectrum and trigger rate fall down as a function of jet  $p_T$ , trigger efficiency can be described by a turn-on curve approaching a plateau region, where the efficiency is maximum and constant within the statistical uncertainty. It is important for the trigger to have a sharp turn-on, with the nominal 95 % efficiency point being below or approximately at the cut value of offline reconstructed jet  $p_T$ . This can be achieved by adjusting the online cuts to be lower than offline ones: if these thresholds are close, then the trigger would bias the analysis distributions. In this case, non-flat  $p_T$ -dependent scale factors have to be applied to the total number of signal events, which also complicates estimation of the systematic uncertainty attributed to the trigger. As for the efficiency dependency on the jet angular distribution, it is expected to be flat – however, as it can be seen from the following section, this is not always the case.

### 3.3.1 Validation of jet energy corrections and charged hadron subtraction for top triggers

In the end of 2011, to match the jet reconstruction method used by most CMS analyses, the calorimeter-based jet module in electron-plus-jets was replaced with PF-based one. This was done in an attempt to improve the  $p_T$  resolution of online jet reconstruction, reduce the rate and obtain sharper turn-on curves of the trigger efficiency. In the first iteration of PF-based trigger paths, no jet energy corrections were applied to the jets online, which resulted in observation of a non-flat trigger efficiency with respect to jet pseudorapidity. After a few months of work, responsible physics object group (known as “JetMET”) provided a series of jet energy corrections to be applied at the HLT level, which were tested by the author.

Pile-up turned out to be another major issue particularly specific to 2012 data-taking period. Increased centre-of-mass energy and instantaneous luminosity at a bunch spacing of 50 ns effectively doubled the average number of pile-up vertices comparing to 2011 conditions.

This led to a significant increase in rate of PF-jets triggers, requiring implementation of pile-up mitigation techniques at the HLT level. As a result, charged hadron subtraction (CHS) for PF jets, also known as “PFnoPU”, was introduced online. As the name suggests, it essentially removes the charged hadrons associated with secondary vertices from PF candidates, which helped to reduce the trigger rate dramatically. As top triggers operate with PF jets, they were greatly affected by pile-up, therefore this correction also had to be validated for electron-plus-jets triggers.

In this work, an impact of both jet energy corrections and pile-up subtraction on electron-plus-jets trigger efficiency was investigated. A rather complicated name of the main signal trigger under study (HLT\_Ele25\_CaloIdVT\_CaloIsoT\_TrkIdT\_TrkIsoT\_TriCentralPF-Jet30) suggests working points of tracker and calorimeter ID and isolation criteria for the signal electron (see Table 3.1), as well as transverse momenta requirements for both electron (25 GeV) and each of the three particle flow jets (30 GeV). Since only hadronic part of the trigger was of interest for this particular study, the efficiency was measured by using a “Single Electron” primary dataset of raw 2012 data. This primary dataset was constructed by accepting events passing at least one of the several single electron triggers with various  $p_T$  and isolation requirements. One of the loosest triggers, referred to as HLT\_Ele27\_WP80, operates with working point “WP80”, also shown in Table 3.1, and electron  $p_T$  threshold of 27 GeV.

Offline selection used to obtain a clean  $t\bar{t}$  sample is outlined in Section 4.2, with an exception of the jet multiplicity requirement: here at least three jets were required. A tighter requirement of at least four jets was also of interest, as described below. The trigger efficiency is then calculated as:

$$\epsilon = \frac{N_{\text{ fired}}}{N_{\text{ selected}}} \quad (3.6)$$

where  $N_{\text{ selected}}$  is the number of events from the primary dataset passing the described offline selection, and  $N_{\text{ fired}}$  is the number of events which in addition fired the electron plus three jet trigger.

All the efficiencies were measured in bins of jet  $p_T$  and  $\eta$ . For these variables, different fit functions were used. The efficiency in  $p_T$  was fitted with the default function describing the turn-on curve:

$$\epsilon(p_T) = A \times e^{(B \times e^{C \times p_T})} \quad (3.7)$$

while the efficiency in  $\eta$  was fitted with

$$\epsilon(\eta) = A \times \eta^2 + B \times \eta + C \quad (3.8)$$

where  $A$ ,  $B$  and  $C$  are the fit parameters.

Although the  $\eta$  dependency was expected to be flat, the observed structure in the  $\eta$  distribution of the PF jet triggers suggested the use of a parabolic fit function. In fact, this drop of the efficiency in the endcap region was the main motivation for the study, as it was believed to be due to the fact that the  $p_T$ - and  $\eta$ -dependent jet energy corrections were not applied online, causing inconsistency between HLT and offline reconstructed jets. Efficiency of the main electron-plus-jets trigger with no jet energy corrections applied is shown on Figure 3.2. Clearly, parabolic shape is seen on the efficiency plot with respect to pseudorapidity of a third reconstructed jet (all the jets were sorted by the value of transverse momentum).

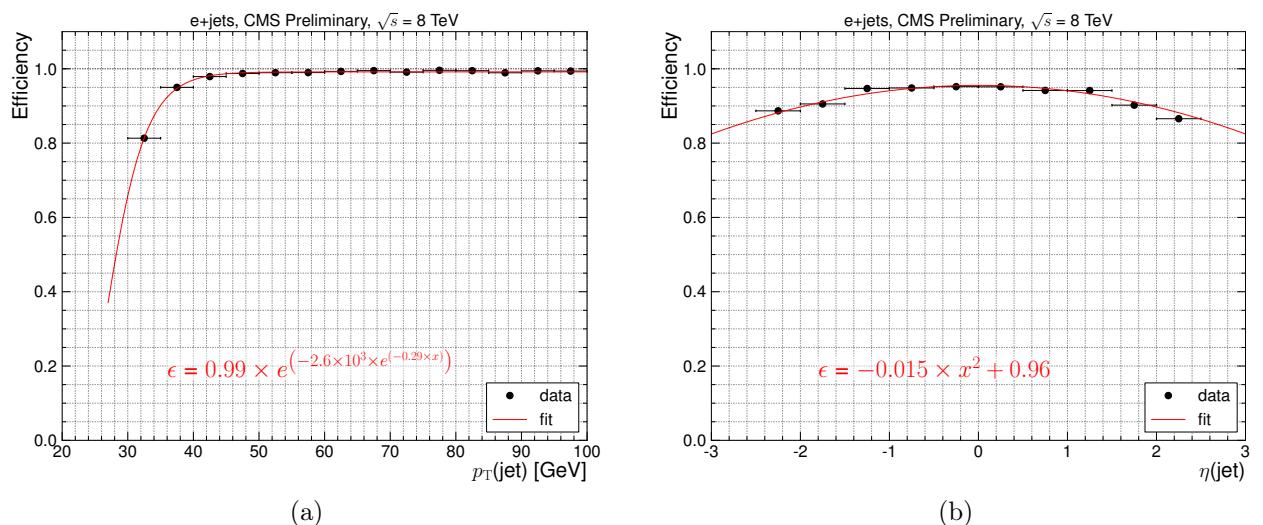


Figure 3.2: Trigger efficiency for HLT\_Ele25\_CaloIdVT\_CaloIsoT\_TrkIdT\_TrkIsoT\_Tri-CentralPFJet30 (uncorrected) as a function of jet  $p_T$  (a) and  $\eta$  (b) for events with at least three jets

It has to be emphasised that only efficiencies with respect to offline selection with at least three jets were affected. Requiring a fourth jet in the event increases the trigger efficiency to nearly 100 %, as it can be seen on Figure 3.3. This implies that observed inefficiency in the endcap region was not critical for analyses requiring at least four jets in the event (including the ones in this thesis), however, it was still important to investigate this effect in order to better understand any possible systematic errors due to the triggering.

For the purpose of the study, a modified HLT menu was created, containing new versions of the main signal trigger with jet energy corrections and charged hadron subtraction applied. This menu was run on a  $t\bar{t}$  skim of a RAW dataset mentioned above, and resulting trigger

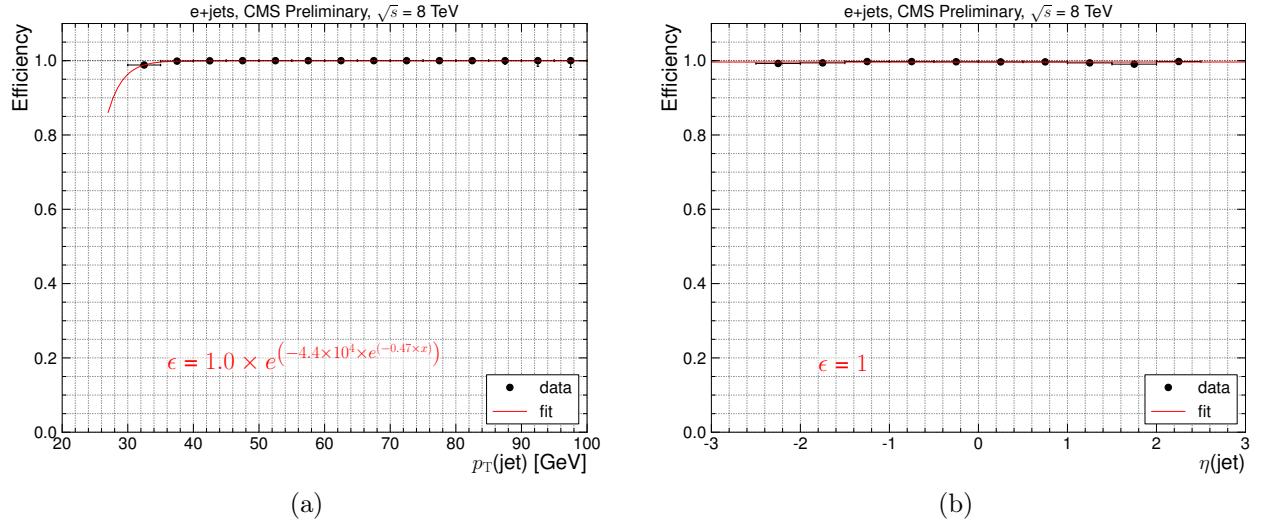


Figure 3.3: Trigger efficiency for HLT\_Ele25\_CaloIdVT\_CaloIsoT\_TrkIdT\_TrkIsoT\_Tri-CentralPFJet30 (uncorrected) as a function of jet  $p_T$  (a) and  $\eta$  (b) for events with at least four jets

flags as well as jet collections were used to produce the final plots.

Figure 3.4 presents the results for the aforementioned trigger path but with jet energy corrections applied, and Figure 3.5 shows it for the same path with both charged hadron subtraction and jet corrections applied. It is obvious that all these corrections haven't solved the inefficiency problem in the endcaps, which caused some confusion in the group that provided them and triggered additional scrutiny from the trigger and jet/ $E_T^{\text{miss}}$  experts.

In order to gain further understanding, jet response distributions were investigated. Jet response is a ratio of the transverse momenta of offline-reconstructed and HLT jets, ideally produced from the same objects in the detector. To find this ratio, all the HLT jets were matched with offline jets found within a cone of  $\Delta R = 0.3$  for each online jet. Matching efficiency proved to be close to unity, and following distributions were produced only for the successfully matched jets. Figure 3.6 shows the response plots for trigger paths with different corrections: uncorrected, only jet energy corrected (JEC) and with both jet energy corrections and charged hadron subtraction (JEC+CHS) applied. It can be seen that corrections work reasonably well in the barrel region, but substantially underestimate the energy of online jets in the endcap region. The response plot shown in slices of offline jet transverse momenta (c) suggests that corrections work better for jets with higher energy, which can also be seen from the distribution with respect to jet  $p_T$ , however, the overall response in the endcap region is still far from being perfect.



Figure 3.4: Trigger efficiency for HLT\_Ele25\_CaloIdVT\_CaloIsoT\_TrkIdT\_TrkIsoT\_Tri-CentralPFJet30 with jet energy corrections applied online, as a function of jet  $p_T$  (a) and  $\eta$  (b) for events with at least three jets



Figure 3.5: Trigger efficiency for HLT\_Ele25\_CaloIdVT\_CaloIsoT\_TrkIdT\_TrkIsoT\_Tri-CentralPFJet30 with jet energy corrections and charged hadron subtraction applied online, as a function of jet  $p_T$  (a) and  $\eta$  (b) for events with at least three jets

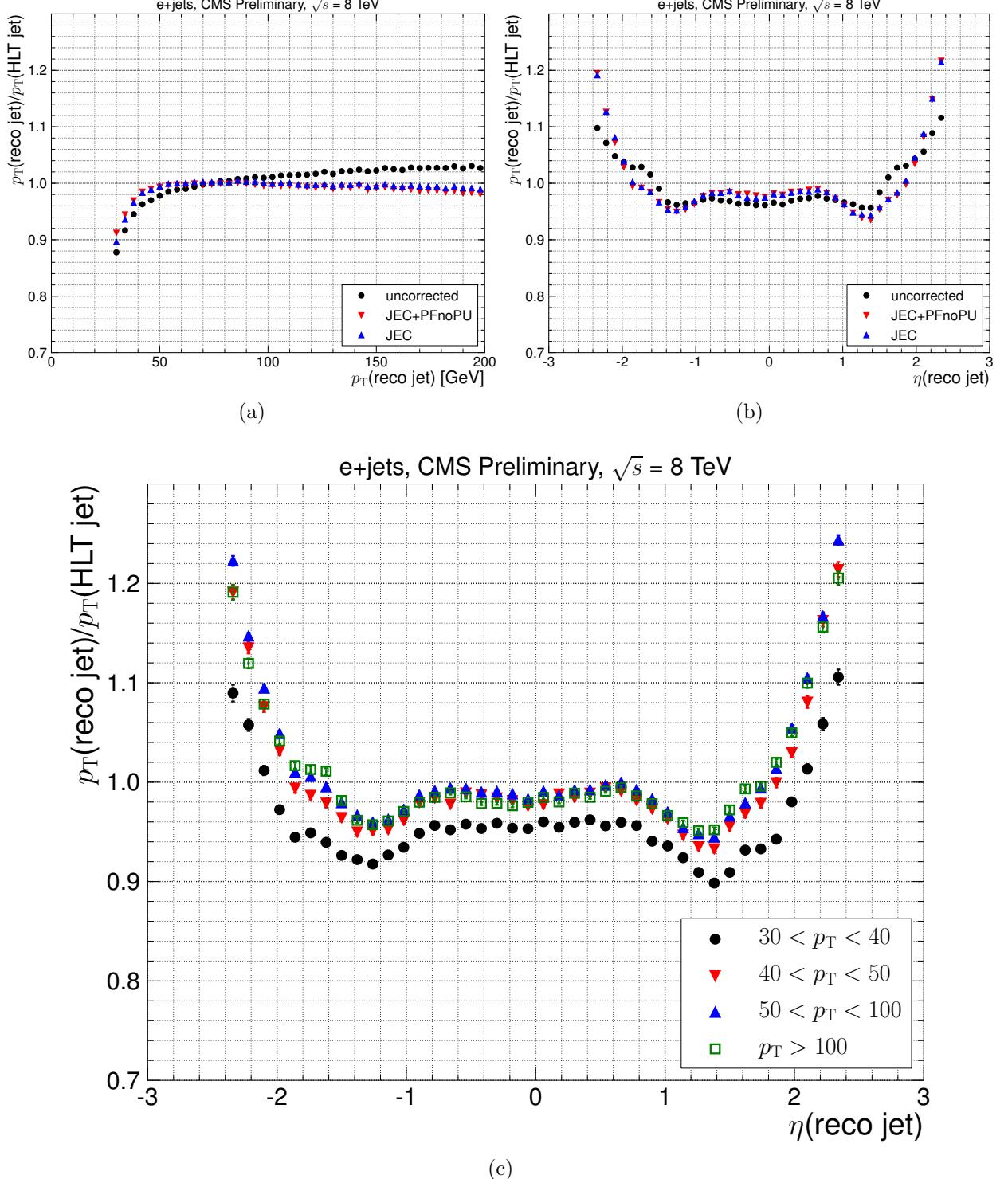


Figure 3.6: Ratio of matched offline and online reconstructed jet transverse momenta as a function of offline jet  $p_T$  (a) and  $\eta$  (b) for different jet corrections and in slices of offline jet  $p_T$  (c)

Shortly after these results were presented to the trigger coordination, they were confirmed independently by other experts within the JetMET group. Because of the time constraints associated with the luminosity increase schedule, a following workaround solution was suggested. In order to re-gain the trigger efficiency, the threshold for the third jet was dropped from 30 GeV to 20 GeV, which allowed to apply jet energy corrections and charged hadron subtraction with a minimum loss in the efficiency in the endcap region. However, this was a subject of a different study.

### **3.4 Summary**

This chapter covered the service work performed by the author for CMS collaboration, namely development, validation and maintenance of high-level triggers for top quark physics, specifically for the electron channel of semileptonic  $t\bar{t}$  decay signature. The importance and typical challenges of trigger development has been discussed. The modular structure of high level triggers has been described in detail for electron-plus-jets paths. A core of author's contribution: the investigation of the effect of jet energy corrections applied on the HLT level to top trigger paths has been presented, showing how it affects trigger efficiency and jet response distributions. An impact of charged hadron subtraction applied online has also been studied and demonstrated.

# 4. Top Quark Mass Measurement

The top quark mass measurement using 2011 LHC data recorded by the CMS detector at a centre of mass energy of  $\sqrt{s} = 7$  TeV is presented in this chapter. The analysis was a cross-check to the official CMS top quark mass measurement at 7 TeV in the lepton plus jets channel published in 2012 [29], which currently remains the most precise single measurement of the top quark mass. Only the electron plus jets channel is described in this thesis, as the muon side of the analysis was performed by a different group at CERN, although in close collaboration.

The mass extraction technique used in this analysis, referred to as the ideogram method, is essentially the same as one used in the CMS measurement of the mass difference between top and antitop quarks [30] and the top quark mass measurement in 2010 [31]. In this chapter, this technique is described in detail. Additionally, data and Monte Carlo samples, event selection process, the kinematic fit used for  $t\bar{t}$  reconstruction, as well as the calibration of the ideogram method, systematic uncertainties and the final result are presented in relevant sections of the chapter.

## 4.1 Data and Simulation

### 4.1.1 Data

The data used in this work is the full 2011 dataset recorded by the CMS detector with a total integrated luminosity of  $5.0 \pm 0.1 \text{ fb}^{-1}$ . As only the electron channel is covered by this particular analysis, all data were pre-selected by single electron plus jets high level triggers, described in detail in Chapter 3.

### 4.1.2 Simulation of signal and background processes

To develop and test any analysis technique in particle physics, simulated events from Monte Carlo (MC) generators as well as a working simulation model of the detector are inevitably needed. In this analysis,  $t\bar{t}$  signal decays and background processes were generated with a set of different MC generators including MADGRAPH [32], PYTHIA [33, 34] and POWHEG [35]. A supplementary package (TAUOLA [36]) was used for the tau lepton decay description. In this section, all the event generators as well as the different MC tunes used for evaluation of systematic uncertainties are briefly described. Additionally, the GEANT4-based [37] CMS detector simulation is mentioned, and finally Monte Carlo samples for all processes are presented.

#### 4.1.2.1 Monte Carlo event generators

Simulation of  $t\bar{t}$  signal and background processes, described in Section 1.2, was performed using several Monte Carlo event generators. The process of simulating a hadron collision event is rather complicated, and usually includes generation of the following steps:

- incoming hadrons (protons);
- hard scattering of partons;
- boson radiation ( $Z, \gamma, g$ );
- parton shower;
- hadronisation of partons.

The performance of MC generators is usually optimised for one or more of these steps, therefore they two or more generators often used in combination, e.g.  $t\bar{t}$  signal is modelled with **MADGRAPH** and **PYTHIA**. Standalone generators can also be used for modelling theory systematic effects, which is also exploited in this thesis.

**MADGRAPH** [32] is a matrix-element Monte Carlo event generator. For any renormalisable Lagrangian-based model it produces all possible Feynman diagrams and automatically generates its matrix elements at the tree level, performing the integration over all phase-space. This calculation is then used to produce the cross section of various processes and subprocesses as well as partons and kinematics of an event, including decays that are described using spin correlations. Partons from matrix element calculations are then matched to parton showers from hadronisation of quarks and gluons, which is simulated in **PYTHIA** (see below) along with the fragmentation of initial protons and the soft scattering of underlying event (mentioned in Section 2.4.3.1). The matching is done according to the so-called MLM prescription [38] if a parton-jet pair satisfies a certain  $\Delta R$  separation criterion. If no or more than one matched jets are found, events are rejected. There is also a certain transverse energy threshold requirement for the partons to be considered in the matching. Default values as well as variations used to estimate the systematic effect of the matching threshold are shown in Table 4.1.

**PYTHIA** [33, 34] is one of the most widely used MC generators in high energy physics. It is a standard tool used for simulation of quark and gluon hadronisation, multi-particle production, beam remnants, initial protons fragmentation, etc. In this work, **PYTHIA** is

used to generate QCD multi-jet production, and to simulate the underlying event on top of other generators providing partons from hard processes.

**MC@NLO** [39] is a parton shower MC generator providing next-to-leading order (NLO) corrections for a simulated process. It represents a major improvement in precision and accuracy of physics simulations comparing to leading order (LO) generators like PYTHIA. NLO implementation can handle additional partons in the final state coming from the hard scattering, which is not possible with LO calculations. Therefore it includes additional dynamic features, giving a major advantage for heavy flavour physics.

**POWHEG** (Positive Weight Hardest Emission Generator) [35] also works at NLO precision. Its main difference to MC@NLO generator lies in the order of process generation: POWHEG initially generates the hardest radiation. This technique allows to avoid double-counting coming from subsequent radiation which may happen with MC@NLO method when matching QCD NLO calculations with parton showers.

#### 4.1.2.2 Matching Threshold and Factorisation Scale

As described above with the example of MADGRAPH generator, there is a certain threshold for partons transverse momenta from matrix element calculations to be matched with parton showers, often referred to as ME-PS threshold. To estimate the systematic uncertainty arising from it, additional MC samples were generated with the value of the threshold either halved or doubled. This procedure was performed for signal  $t\bar{t} + \text{jets}$  and background  $W + \text{jets}$  and  $Z + \text{jets}$  samples. Similarly, the systematic effect of the choice of factorisation scale  $Q^2$  (described in Section 1.2.1) is determined. A summary of central values and variations used to estimate these systematic uncertainties is shown in Table 4.1, whereas the list of systematic MC samples is shown in Table 4.4.

#### 4.1.2.3 CMS detector simulation

Detector simulation is an essential final step on top of the full physics simulation of proton-proton collisions. All generated particles resulting from these collisions are propagated to the simulated layers of the detector so that their interaction with the detector material and detector response can also be modelled. The CMS simulation is based on the GEANT4 (GEometry ANd Tracking [37]) package, which is implemented in CMS software CMSSW [8]. This package fully recreates the geometry of the detector, including all its subsystems as well as a detailed map of the magnetic field.

In general, the simulation procedure includes modelling of the following steps:

- interaction region;

parameter	factorisation scale	matching threshold
$t\bar{t}$ events		
central value	$Q^2 = m_t^2 + \sum_{\text{jets}} p_T^2$	20 GeV
variations	$(0.5 \cdot Q)^2, (2 \cdot Q)^2$	10 GeV and 40 GeV
W + jets and Z + jets events		
central value	$Q^2 = m_{W/Z}^2 + \sum_{\text{jets}} p_T^2$	10 GeV
variations	$(0.5 \cdot Q)^2, (2 \cdot Q)^2$	5 GeV and 20 GeV

Table 4.1: Central values and variations used for systematic uncertainties of factorisation/renormalisation scale ( $Q^2$ ) and matrix element – parton shower (ME-PS) matching in MADGRAPH MC generator.

- particles traversing all the layers of the detector and the interaction processes;
- multiple interactions per beam crossing (pile-up) and event overlay effects;
- digital readout of the detector (digitisation).

The data stream generated with this process is very similar to the output of the real detector, although naturally in addition to reconstructed objects, generated events also include collections of generated objects, allowing users to access the Monte Carlo truth information, which is essential for any particle physics analysis.

#### 4.1.2.4 Monte Carlo samples

The nominal signal  $t\bar{t}$  MC sample used in this analysis was generated using the MADGRAPH generator with a top quark mass of  $m_t = 172.5$  GeV. To calibrate the mass extraction technique, as explained in Section 4.5, 8 additional  $t\bar{t}$  samples were used with different top quark masses ranging between 161.5 GeV and 184.5 GeV. The W + jets and Z + jets background processes (see Section 1.2.3) were also generated with MADGRAPH, whereas the single top background was generated with POWHEG generator. A full list of signal  $t\bar{t}$  and background MC samples with values of cross section, numbers of generated events and corresponding integrated luminosities is shown in Table 4.2. All samples were produced using the CTEQ6L Parton Distribution Functions (PDF) [40].

The QCD background samples are available in two different sets, one of which is pre-filtered to contain jets with high electromagnetic content ( $e/\gamma$  enriched), and the second one is enriched with electrons coming from decays of heavy flavour quarks ( $b/c \rightarrow qe\nu$ ).

This pre-selection is performed at generator level, to ensure reasonable statistics after the  $t\bar{t}$  selection imposed in top quark analyses with electrons in the final state. Both sets of QCD samples are generated with PYTHIA generator in different exclusive bins of  $\hat{p}_T$  (i.e. transverse momentum of the outgoing partons defined in the rest frame of the hard process). The  $\gamma+$ jets MC samples were produced with MADGRAPH in different bins of  $H_T$  at generator level, which denotes the sum of transverse momenta of all high- $p_T$  objects. Table 4.3 gives a list of QCD and  $\gamma+$ jets samples with values of production cross sections, pre-selection efficiencies, numbers of generated events and corresponding integrated luminosities.

Finally, Table 4.4 lists the additional samples generated to allow estimation of the systematic uncertainties due to factorisation scale and matching threshold choices, as discussed earlier.

Process	Generator	$\sigma$ (pb)	# events	$\int \mathcal{L} dt$ (fb $^{-1}$ )
$t\bar{t} + \text{jets}$	MADGRAPH			
$m_t = 161.5 \text{ GeV}$		157.5	1620072	10.3
$m_t = 163.5 \text{ GeV}$		157.5	1633197	10.4
$m_t = 166.5 \text{ GeV}$		157.5	1669034	10.6
$m_t = 169.5 \text{ GeV}$		157.5	1606570	10.2
$m_t = 172.5 \text{ GeV}$		157.5	7490162	47.6
$m_t = 175.5 \text{ GeV}$		157.5	1538301	9.8
$m_t = 178.5 \text{ GeV}$		157.5	1648519	10.5
$m_t = 181.5 \text{ GeV}$		157.5	1665350	10.6
$m_t = 184.5 \text{ GeV}$		157.5	1671859	10.6
$W + \text{jets } (W \rightarrow l\nu)$	MADGRAPH	31314	81345381	2.6
$Z/\gamma^* \rightarrow l^+l^- + \text{jets}, m(ll) > 50 \text{ GeV}$	MADGRAPH	3048	36222153	11.9
Single top	POWHEG			
top $t$ -channel		42.6	3814228	89.5
anti-top $t$ -channel		22.0	1944822	88.4
top $s$ -channel		2.72	259971	92.6
anti-top $s$ -channel		1.49	137980	92.6
top $tW$ -channel		5.3	814390	153.7
anti-top $tW$ -channel		5.3	809984	152.8

Table 4.2: Signal and background Monte Carlo samples with cross sections at  $\sqrt{s} = 7 \text{ TeV}$ , numbers of generated events and corresponding integrated luminosities.

Process	Generator	$\sigma$ (pb)	filter efficiency	# events	$\int \mathcal{L} dt$ (fb $^{-1}$ )
QCD ( $e/\gamma$ enriched)	PYTHIA	$2.355 \times 10^8$	$7.3 \times 10^{-3}$	34720808	$2.0 \times 10^{-2}$
20 GeV < $\hat{p}_T$ < 30 GeV		$5.93 \times 10^7$	0.059	70375915	$2.0 \times 10^{-2}$
30 GeV < $\hat{p}_T$ < 80 GeV		$9.06 \times 10^5$	0.148	8150669	$6.1 \times 10^{-2}$
QCD (b/c $\rightarrow e\nu$ )	PYTHIA	$2.355 \times 10^8$	$4.6 \times 10^{-4}$	2002588	$1.8 \times 10^{-2}$
20 GeV < $\hat{p}_T$ < 30 GeV		$5.93 \times 10^7$	$2.34 \times 10^{-3}$	2030030	$1.5 \times 10^{-2}$
30 GeV < $\hat{p}_T$ < 170 GeV		$9.06 \times 10^5$	0.0104	1082690	0.1
$\gamma + \text{jets}$	MADGRAPH	23 620		12730863	0.5
40 GeV < $H_T$ < 100 GeV		3476		1536287	0.4
100 GeV < $H_T$ < 200 GeV		485		9377168	19.3
$H_T > 200$ GeV					

Table 4.3: QCD multi-jet background and  $\gamma + \text{jets}$  MC samples with cross sections at  $\sqrt{s} = 7$  TeV, numbers of generated events and corresponding integrated luminosities.

Process	Generator	$\sigma$ (pb)	# events	$\int \mathcal{L} dt$ (fb $^{-1}$ )
$t\bar{t} + \text{jets}$	MADGRAPH	157.5	1607808	10.2
0.5 $\times$ matching threshold		157.5	4029823	25.5
2 $\times$ matching threshold		157.5	4004587	25.4
0.5 $\times Q$		157.5	3696269	23.4
2 $\times Q$				
W + jets (W $\rightarrow l\nu$ )	MADGRAPH	29690	9956679	0.3
0.5 $\times$ matching threshold		30290	10461655	0.3
2 $\times$ matching threshold		33300	10092532	0.3
0.5 $\times Q$		32000	9784907	0.3
2 $\times Q$				
Z + jets (Z $\rightarrow ll$ )	MADGRAPH	2888	1614808	0.6
0.5 $\times$ matching threshold		2915	1641121	0.6
2 $\times$ matching threshold		3312	1658855	0.5
0.5 $\times Q$		2954	1592742	0.5
2 $\times Q$				

Table 4.4: Systematic MC samples with cross sections at  $\sqrt{s} = 7$  TeV, numbers of generated events and corresponding integrated luminosities. Factorisation scale  $Q$  and matching threshold systematic uncertainties (see Section 4.1.2.2) are estimated with variations of  $t\bar{t} + \text{jets}$ , W + jets and Z + jets samples.

## 4.2 Event Selection

The event selection applied in this analysis essentially follows the so-called reference selection, i.e. the one recommended by the CMS Top Quark Physics Analysis Group for SM top quark measurements. It is designed to select the semileptonic signature of a  $t\bar{t}$  decay with an isolated lepton, four jets (including two jets from b-quarks) and a neutrino in the form of missing transverse energy. All objects are reconstructed using the particle flow method as described in Section 2.4.

As this analysis focuses on the electron channel, the event selection proceeds as follows:

1. preselection;
2. trigger;
3. electron candidate selection;
4. dilepton veto;
5. conversion veto;
6. muon veto;
7. jet selection;
8. b-tagging.

### Preselection

Preselection is performed in order to reduce the number of events stored for local analysis. All events are required to have at least one electron with  $E_T$  above 30 GeV within pseudorapidity region of  $|\eta| < 2.5$ , or at least one muon with  $p_T$  above 20 GeV and  $|\eta| < 2.1$ . Additional event cleaning is performed at this stage to decrease the number of events with substantial detector noise. A small fraction of events with anomalous HCAL noise is removed by an HCAL noise filter. Furthermore, beam scraping is mitigated by the appropriate filter requiring at least 25 % tracks with high purity if there are more than 10 tracks in the event. Finally, the preselection stage includes the requirement of at least one good primary vertex in the event which has to be located within 24 cm in the  $z$ -direction from the centre of CMS and within a radial distance of 2 cm. The primary vertex is also required to have at least four degrees of freedom which is determined from the vertex fit and is strongly correlated with the number of tracks compatible with the primary interaction region [41].

### Trigger

The trigger requirement is imposed as described in Chapter 3. Essentially, all the events are required to fire the electron-plus-three-jets trigger in the version current at the time of data-taking. A full list of triggers can be found in Appendix A.

### Electron candidate selection

Exactly one electron candidate satisfying all the following criteria is required to be present in the event. As the electron-plus-jets trigger has a lepton  $p_T$  threshold of 25 GeV, the electron candidate transverse momentum is required to be above 30 GeV, in order to be in the trigger efficiency plateau region. Furthermore, the electron has to be within the tracker region of  $|\eta| < 2.5$  excluding the ECAL barrel-endcap transition regions of  $1.4442 < |\eta| < 1.566$ . The CiC electron ID (see Section 2.4.1) with the tightest working point (“HyperTight”) is used for electron identification purposes. This working point provides a misidentification rate of less than 1 %, and overall identification efficiency of 75 % [17].

In order to only select electrons originating at the interaction vertex, the  $x$ - $y$  distance ( $d_{xy}$ ) between the electron track and the interaction point (2D impact parameter) is required to be less than 0.02 cm. The contribution of electrons inside jets from QCD background events, as well as fake electrons, is reduced by imposing the PF-based relative isolation cut on the variable  $\text{reliso}$  defined in Section 2.4.1.2. A  $\Delta R$ -cone with size of 0.3 is used in calculation of the  $\text{reliso}$  variable, and events with  $\text{reliso} < 0.1$  are accepted.

### Dilepton veto

A veto requirement on a second electron candidate (or dilepton veto) is used to reject events with any additional electrons. Looser criteria are used to identify the second electron in an event. Namely, the event is rejected if there is a second electron satisfying a lower  $E_T$  threshold of 20 GeV and a looser cut on  $\text{reliso} (< 0.2)$ .

Additionally, in order to reduce  $Z + \text{jets}$  background, the following veto is applied. If the event contains a second electron with  $E_T > 30$  GeV and  $\text{reliso} < 1.0$  that forms such an invariant mass with the signal electron candidate that it is close to the  $Z$  mass peak, it is also rejected. The invariant mass window around the  $Z$  peak is chosen to be  $76 \text{ GeV} < m_{ee} < 106 \text{ GeV}$ .

### Conversion veto

As mentioned in Section 2.4.1.3, two methods are used in this analysis to remove electrons coming from photon conversions: missing pixel layers method and partner track matching method. If missing pixel layer hits or a matching partner track are present in the event, it is rejected.

### Muon veto

Other  $t\bar{t}$  channels, including muon and dilepton decay modes, can contaminate events passing the electron plus jets selection. To reduce this contamination, events containing an isolated global muon (see Section 2.4.2) are rejected. The muon is required to have a  $p_T$  above 10 GeV,  $|\eta| < 2.5$  and  $\text{reliso} < 0.2$  with a  $\Delta R$  cone of 0.4.

### Jet selection and b-tagging

Final steps in the selection help to further reduce the background by imposing constraints on the number of jets. This is particularly effective to mitigate the  $W + \text{jets}$  contamination, as the number of  $W + \text{jets}$  events decreases exponentially with increasing number of jets. In this analysis, at least four jets with  $p_T > 30 \text{ GeV}$  and  $|\eta| < 2.4$  are required to be present in the event. These jets have to pass the loose PF jet ID (see Section 2.4.3). Finally, the CSV b-tagging algorithm with medium working point (Section 2.4.3.3) is used to identify the two b-quarks from the  $t\bar{t}$  decay.



Figure 4.1: Kinematic variable distributions after all selection cuts. (a) electron  $p_T$ , (b) electron  $\eta$ , (c) jet  $p_T$  for all jets passing the selection, (d)  $E_T^{\text{miss}}$ .

### 4.3 Kinematic Fit

In order to precisely measure the top quark mass, objects in the final state need to be associated with the originating partons of the top quark pair decay. In this analysis, a kinematic fit is employed to fully reconstruct the event kinematics under the  $t\bar{t}$  hypothesis, thus improving the resolution of measured quantities by exploiting the knowledge of decay process.

The HITFIT fitting package [42] used in this work originates from the DØ collaboration. Based on the SQUAW algorithm [43] developed in Lawrence Berkeley National Laboratory, the original HITFIT code was written by Scott Stuart Snyder in Fortran for Run I of the Tevatron. It was successfully used in the first direct measurement of the top quark mass by DØ for lepton plus jets  $t\bar{t}$  events [44]. The package was then migrated to C++ for Run II of the Tevatron, and exploited in a series of  $t\bar{t}$  analyses including the updated top quark mass measurement using the ideogram method [45], first measurement of the forward-backward charge asymmetry in  $t\bar{t}$  production [46] and some others.

One of the key features of the HITFIT package is its independence of specific experiments and detector geometries. The only external dependence is on CLHEP [47], a C++ library that provides utility classes for linear algebra, geometry, vector arithmetic, etc., which was specifically developed for high energy physics simulation and analysis software and is notably used by GEANT4 package. In 2011 HITFIT was ported to CMSSW and later on successfully used in a number of CMS top quark analyses, including the top mass measurement in the lepton plus jets channel [29].

HITFIT is designed to fit lepton plus jets events under the  $t\bar{t}$  hypothesis. Specifically, it strives to constrain the event to a hypothesis of production of two heavy particles, each decaying into a W boson and b quark. According to the semileptonic signature, one of the W boson decays into an electron-neutrino pair, while the other decays into a light quark-antiquark pair.

The input to the fit includes the four-momenta of the lepton, four leading jets and missing transverse energy, as well as their respective resolutions that were calculated using Monte Carlo samples. Since only four leading jets are used, there are  $4! = 24$  ways to associate these four reconstructed jets with the partons from the  $t\bar{t}$  decay. The number of permutations is reduced to  $4!/2 = 12$  since the light jets from the hadronic W decay are interchangeable. However, for each hypothesis the  $z$  component of neutrino four-momentum needs to be determined. This is done by requiring the two heavy particles to have equal mass ( $m_t = m_{\bar{t}}$ ), which yields a quadratic equation which has either two real solutions, or two complex solutions that share the real part. In the case of the two real solutions, the fit

is performed for each of them. If there are two complex solutions, the fit is performed once for the common real part of them, but the fit result is included in the event content twice. Therefore, each hypothesis has two possible values of the neutrino four-momentum, which doubles the number of permutations back to  $4! = 24$ .

The kinematic fit is performed by minimising the  $\chi^2$  function defined as

$$\chi^2 = (\mathbf{x} - \mathbf{x}^m)^T \mathbf{G} (\mathbf{x} - \mathbf{x}^m) \quad (4.1)$$

where  $\mathbf{x}^m$  is the vector of measured observables,  $\mathbf{x}$  is the vector of fitted variables and  $\mathbf{G}$  is the inverse of the error matrix given by the resolutions of the observables.

The following kinematic constraints are imposed:

1. Reconstructed top quarks from both hadronic and leptonic legs are required to have equal masses:

$$m(t \rightarrow \ell\nu b) = m(\bar{t} \rightarrow q\bar{q}\bar{b}), \quad (4.2)$$

2. The W boson masses reconstructed in both hadronic and leptonic legs are fixed:

$$m(\ell\nu) = m(q\bar{q}) = m_W = 80.4 \text{ GeV}. \quad (4.3)$$

Apart from exploiting the knowledge of the W mass, the b quark mass of 4.7 GeV is also used. Light quarks and leptons are considered to be massless.

The full description of the fitting algorithm can be found in the appendix of the HITFit author's PhD thesis [48]. Essentially, since the constraints are non-linear, an iterative technique is used. Starting with the measured observables, the constraint equations are linearised by expansion in a power series around the starting point. The minimisation is then solved with these linearised constraints, providing the new starting value for the next minimisation step which is repeated until the  $\chi^2$  converges or a maximum number of iterations is reached. Events are rejected if there is no permutation resulting in a converging fit. Some kinematic variables for permutations giving the best (smallest)  $\chi^2$  values are shown in Figure 4.2.

Typically, an event has a few permutations resulting in a converged fit. The simplest approach of picking a permutation with the best  $\chi^2$  value is disfavoured, since Monte Carlo studies have shown that the correct permutation does not necessarily have the best  $\chi^2$  [48]. This value can change drastically with relatively small variations of the input parameters, as different permutations may obtain the smallest value of  $\chi^2$ . On the other hand, some solutions yield very large  $\chi^2$  value, which typically happens for badly reconstructed or background events. Therefore, a loose cut of  $\chi^2 = 20$  is applied for all permutations on the



Figure 4.2: Distribution of (a) fitted top quark mass, (b) its standard deviation, (c)  $\chi^2$  and (d) fit probability for the best fitted hypotheses.

output of the converged fit. This is the final selection criterion on top of the reference selection mentioned in Section 4.2.

The compatibility of a permutation with the  $t\bar{t}$  hypothesis is quantified by the following fit probability:

$$P_{\text{fit}} = \exp\left(-\frac{1}{2}\chi^2\right). \quad (4.4)$$

This probability is used in the next step of the analysis, referred to as the ideogram method.

## 4.4 The Ideogram Method

The ideogram method is a powerful mass extraction technique. Its main idea lies in finding the likelihood distribution of the particle mass given the particular data sample, which is performed by evaluating the likelihood of measured observables given a generated particle mass for each event in the sample. The signal likelihood is evaluated from the theoretically expected Breit-Wigner distribution convoluted with the Gaussian experimental resolution on event-by-event basis. The whole procedure is calibrated using Monte Carlo simulation of a series of particle masses in the expected region.

One of the first applications of the ideogram method tracks back to the W boson mass measurement by the DELPHI Collaboration at CERN LEP collider [49, 50]. It was later used in a series of top quark mass measurements, including the DØ measurement in the lepton plus jets channel [45] and CDF measurement in the all-hadronic channel [51]. The main CMS top quark mass measurement in lepton plus jets channel at 7 TeV [29] which is cross-checked by this analysis is also based on the ideogram method with, however, *in situ* JES (jet energy scale) measurement in a joint likelihood fit. As JES is a dominant systematic in this analysis, this approach proved to be more precise, but let us not neglect the importance of a cross-check.

### 4.4.1 Event likelihood

In the basis of the ideogram method lies the event likelihood, also referred to as the ideogram. It is calculated for each event in the sample as a function of the hypothesised top quark mass  $m_t$  in the following way:

$$\mathcal{L}_{\text{event}}(\mathbf{x}|m_t, f_{t\bar{t}}) = f_{t\bar{t}} P_{t\bar{t}}(\mathbf{x}|m_t) + (1 - f_{t\bar{t}}) P_{\text{bkg}}(\mathbf{x}). \quad (4.5)$$

Here  $\mathbf{x}$  is the set of event observables with the top mass information from the kinematic fit,  $f_{t\bar{t}}$  is the fraction of  $t\bar{t}$  events in the data sample, and  $P_{t\bar{t}}$  and  $P_{\text{bkg}}$  are probability densities for signal and background events, respectively. Essentially, by construction of the

likelihood, the first term corresponds to the probability for the event to be signal, whilst the second term is its probability to be background. The kinematic information  $\mathbf{x}$  includes the fitted mass  $m_i$ , the estimated standard deviation for each fitted mass  $\sigma(m_i)$  and the goodness-of-fit  $\chi^2_i$  for all permutations, i.e. all possible jet-parton assignments and neutrino solutions surviving the selection criteria on the output of the fit.

### Jet-parton assignment weights

The signal and background probabilities forming the event likelihood are calculated as a weighted sum over jet-parton assignments and neutrino solutions coming from the kinematic fit. Each permutation is weighted by the probability that it is correct, which is essentially the fit probability  $P_{\text{fit}}$  introduced in Equation 4.4. However, this particular implementation of the ideogram method also exploits the b-tagging information, thus the fit probability is multiplied by an additional term  $p_{\text{btags}}$ , which is the probability that the particular jet-parton assignment is compatible with the observed b-tags:

$$p_{\text{btags}} = \prod_{j=1}^4 p^j \quad (4.6)$$

Here the index  $j$  runs over all four jets in the kinematic fit, and the probability  $p^j$  can be either  $\varepsilon_l$ ,  $(1 - \varepsilon_l)$ ,  $\varepsilon_b$ , or  $(1 - \varepsilon_b)$ , depending on the flavour of the jet as derived by the fit (light jet or b-jet), and whether the jet is b-tagged or not. As a recap from Section 2.4.3.3, for the medium working point of CSV b-tagging algorithm used in this analysis, b-tagging efficiency is approximately  $\sim 70\%$ , whilst the mis-tag rate (or  $(1 - \varepsilon_l)$ ) is approximately  $1\%$ .

Therefore, for each such permutation the relative weight is  $w_i$  is calculated as:

$$w_i = p_{\text{fit}} \cdot p_{\text{btags}} = \exp\left(-\frac{1}{2}\chi^2\right) \cdot \prod_{j=1}^4 p^j \quad (4.7)$$

These weights are used in the calculation of signal and background probability densities, described hereafter.

### Signal probability

The  $t\bar{t}$  signal probability in the event likelihood shown in Equation 4.5 is calculated as a weighted sum over all permutations:

$$P_{t\bar{t}}(\mathbf{x}|m_t) = \sum_i^{24} w_i \left( f_{cp} \cdot \int_{m_{min}}^{m_{max}} G(m'|m_i, \sigma_i) RBW(m'|m_t, \Gamma_t) dm' + (1 - f_{cp}) WP(m_i|m_t) \right) \quad (4.8)$$

Here  $f_{cp}$  is the fraction of events in which the maximum weight is assigned to the correct permutation, estimated from simulated  $t\bar{t}$  events. For each permutation, the first term expresses the probability that it has a correct jet-parton assignment (correct permutation). It is calculated as a convolution of a relativistic Breit-Wigner distribution  $RBW(m'|m_t, \Gamma_t)$  and a Gaussian  $G(m'|m_i, \sigma_i)$ .

The relativistic Breit-Wigner distribution describes the theoretical top quark mass distribution for a given hypothesised top quark mass value  $m_t$  with a top quark width  $\Gamma_t$  set to 2 GeV, according to the latest experimental results [52]. As suggested by Standard Model electroweak corrections and  $s$ -dependence of the phase space [53], the following form of  $RBW$  is used:

$$RBW(m'|m_t, \Gamma_t) \propto \frac{m'^2}{(m'^2 - m_t^2)^2 + m'^4 \frac{\Gamma_t^2}{m_t^2}} \quad (4.9)$$

The Gaussian in the convolution represents the detector resolution. It is centred at the fitted mass value  $m_i$  with a width of the fitted mass uncertainty  $\sigma_i$  as estimated by the fit:

$$G(m'|m_i, \sigma_i) = \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(\frac{(m' - m_i)^2}{2\sigma_i^2}\right) \quad (4.10)$$

The resulting probability density for  $t\bar{t}$  events with correct permutations is shown in Figure 4.3. The Monte Carlo sample with top mass of 166.5 GeV was used, for illustration.

The second term shows the probability of a wrong permutation in the signal event. This probability density, denoted  $WP(m_i|m_t)$ , is calculated by fitting the Crystal Ball analytical function to the fitted top quark mass distribution for permutations with wrong jet-parton assignment, estimated purely from Monte Carlo  $t\bar{t}$  sample whilst taking the permutation weights (Equation 4.7) into account. The wrong permutations fall in two categories: the first one corresponds to events where the four leading jets come from  $t\bar{t}$  decay, but are not assigned to the partons correctly; the second one are events with one or more of the four leading jets coming from soft gluon radiation (ISR or FSR). It was found that the fitted top quark mass distributions have similar shapes for both categories, therefore they were

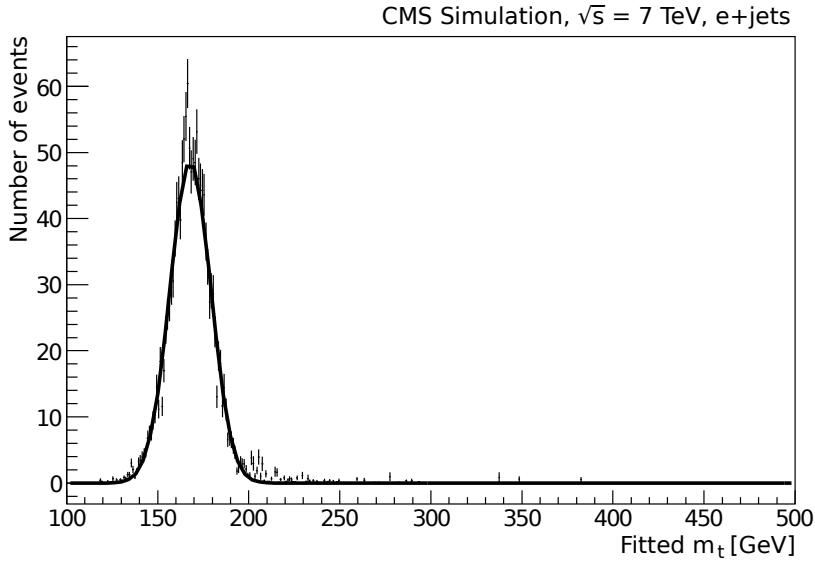


Figure 4.3: The fitted top quark mass distribution for  $t\bar{t}$  events with correct permutations (MC truth,  $m_t = 166.5$  GeV).

combined for the purpose of fitting. The Crystal Ball function has the form of:

$$\text{CB}(x; \mu, \sigma, \alpha, n) = N \times \begin{cases} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) & \text{if } \frac{x-\mu}{\sigma} > -\alpha, \\ \frac{1}{\sqrt{2\pi}\sigma} \left(\frac{n}{|\alpha|}\right)^n \exp\left(-\frac{|\alpha|^2}{2}\right) \left(\frac{n}{|\alpha|} - |\alpha| - \frac{x-\mu}{\sigma}\right)^{-n} & \text{if } \frac{x-\mu}{\sigma} \leq -\alpha. \end{cases} \quad (4.11)$$

The exponential order of  $n = 5$  was used, and parameters  $\mu$ ,  $\sigma$ , and  $\alpha$  were parametrised as functions of the top quark mass. The fitted distributions of wrong permutations spectra for three different top masses (166.5 GeV, 172.5 GeV and 178.5 GeV) are shown in Figure 4.4.

### Background probability

The shape of the background probability density in Equation 4.5 is assumed not to depend on the top quark mass, and is estimated from the background Monte Carlo samples. The procedure is done in a similar way to the wrong permutations probability, and the jet-parton assignment weights (Equation 4.7) are also taken into account. Since the dominant background is  $W + \text{jets}$ , it was used exclusively for extracting the fitted top quark mass distribution by running the kinematic fit. The contribution from single top,  $Z + \text{jets}$  and QCD was expected to be small, and it was confirmed that their probability density distributions have similar shapes to that of the  $W + \text{jets}$  sample. Moreover, the Monte Carlo calibration described later on in Section 4.5 takes all mentioned background sources into account,

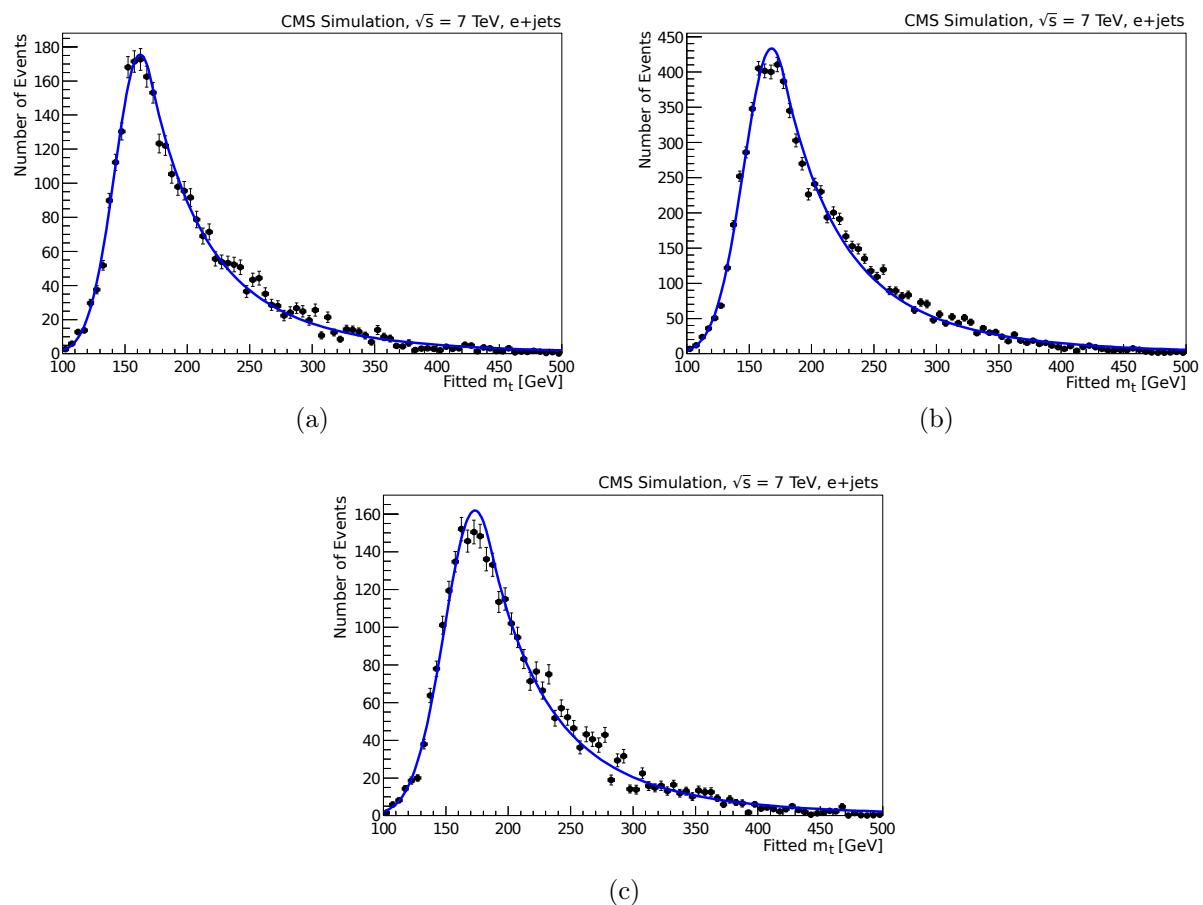


Figure 4.4: The fitted top quark mass distribution for  $t\bar{t}$  events with wrong permutations (MC truth) for (a)  $m_t = 166.5$  GeV, (b)  $m_t = 172.5$  GeV, (c)  $m_t = 178.5$  GeV.

therefore any possible bias due to this approximation is corrected for.

The fitted top quark mass spectrum for background events as estimated from the W+jets MC sample is shown in Figure 4.5. This distribution is described by the Landau function, fitted by the ROOT package [54].

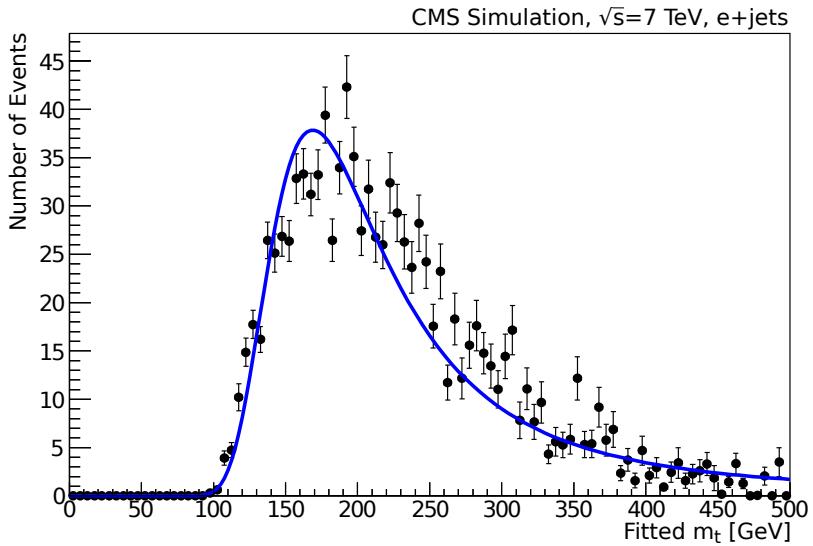


Figure 4.5: The fitted top quark mass distribution for  $t\bar{t}$  events in W + jets Monte Carlo sample.

#### 4.4.2 Combined likelihood fit

Finally, for the extraction of the top quark mass the total likelihood for the entire data sample is calculated as the product of all the event likelihoods:

$$\mathcal{L}_{\text{sample}} = \prod_{\text{all events}} \mathcal{L}_{\text{event}}(\mathbf{x}|m_t) \quad (4.12)$$

To extract the most probable top quark mass, this likelihood needs to be maximised. In practice, it is simplified by performing the minimisation of the summed negative log-likelihoods:

$$-2 \log \mathcal{L}_{\text{sample}} = -2 \sum_{\text{all events}} \log \mathcal{L}_{\text{event}}(\mathbf{x}|m_t) \quad (4.13)$$

The result with the full dataset is shown in Figure 4.6. The top quark mass and its statistical uncertainty are calculated by using the parabolic interpolation with three lowest points in the log-likelihood curve. However, this result can not be regarded as the final

measurement, as it does not take into account the possible mass bias of the ideogram method. The mass calibration technique is described in the next section.

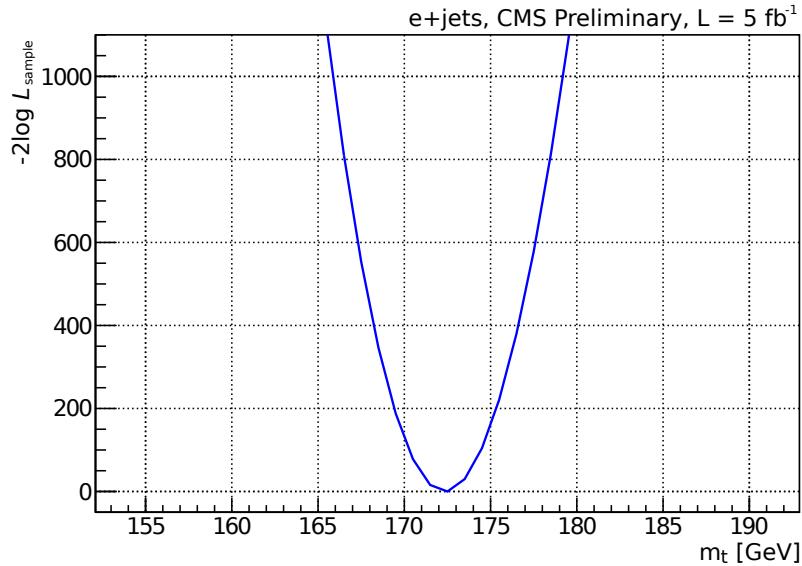


Figure 4.6: Log-likelihood curve as a function of top quark mass for e+jets dataset before calibration.

## 4.5 Mass Calibration

By construction of the ideogram method, it has approximations which can lead to possible bias in the top quark mass measurement. To account for this bias, a calibration is performed by running the measurement on Monte Carlo simulation with known true top quark mass. Pseudo-experiments are used to measure the bias in the fitted mass, as well as to calibrate the statistical uncertainty of the measurement.

For each of the nine generated  $t\bar{t}$  samples with mass points between 161.5 GeV and 184.5 GeV, 5000 pseudo-experiments were performed with events picked from Monte Carlo samples representing different sources of events, so that the numbers of events follow a Poisson distribution centred around the corresponding numbers observed in real data. The observed value of 23 727 electron plus jets events with an average signal fraction of 76 % was used in each pseudo-experiment. The background was estimated by using events from the  $W + \text{jets}$  MC sample, whereas the impact of other background sources was accounted for separately as a systematic uncertainty (Section 4.6). For every pseudo-experiment all the event likelihoods are calculated, the total log-likelihood is formed and the top quark mass ( $m_i$ ) with its statistical uncertainty ( $\sigma_i$ ) are extracted in the same way as in the real measurement described above. These results are used to derive the bias and pull distributions,

defined as follows:

$$\text{bias} = m_t^{\text{meas}} - m_t^{\text{gen}} \quad \text{and} \quad \text{pull}_i = \frac{m_i - m_t^{\text{meas}}}{\sigma_i}, \quad (4.14)$$

where  $m_t^{\text{meas}}$  is the mean top quark mass over all pseudo-experiments, for a given generated mass point. The width of the pull is defined as the standard deviation of a Gaussian fitted to the pull distribution.

The mass bias and the width of the pull as a function of the generated top quark mass are shown in Figure 4.7. Clearly, using the ideogram method out of the box reveals a significant bias which needs to be corrected for. It can also be seen that the width of the pull distribution is approximately 11 % larger than unity, suggesting that the ideogram method underestimates the statistical uncertainty.

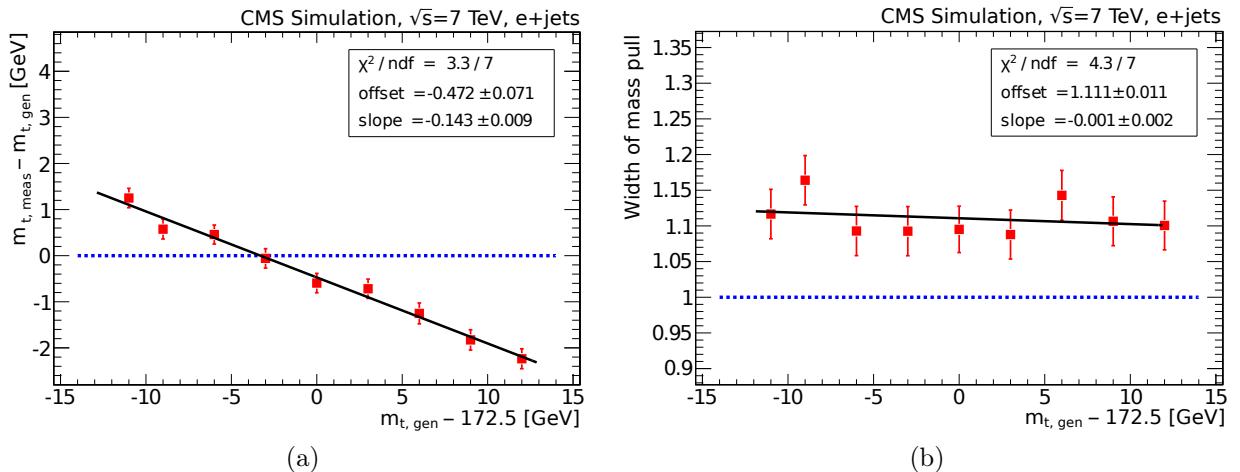


Figure 4.7: Average mass bias (a) and width of the pull (b) as a function of true top quark mass.

Due to the apparent bias, the fitted linear calibration curves in both distributions are used to correct the final measurement of the top quark mass. As a closure test, the mass bias and pull width as a function of true top quark mass are calculated after applying the full calibration. The distributions are shown in Figure 4.8. Evidently, the closure tests confirm that the top quark mass and its statistical uncertainty are unbiased after calibration.

## 4.6 Systematic Uncertainties

Like any measurement, the top quark mass measurement is prone to various systematic errors. In attempt to evaluate them, several sources of systematic effects are considered in

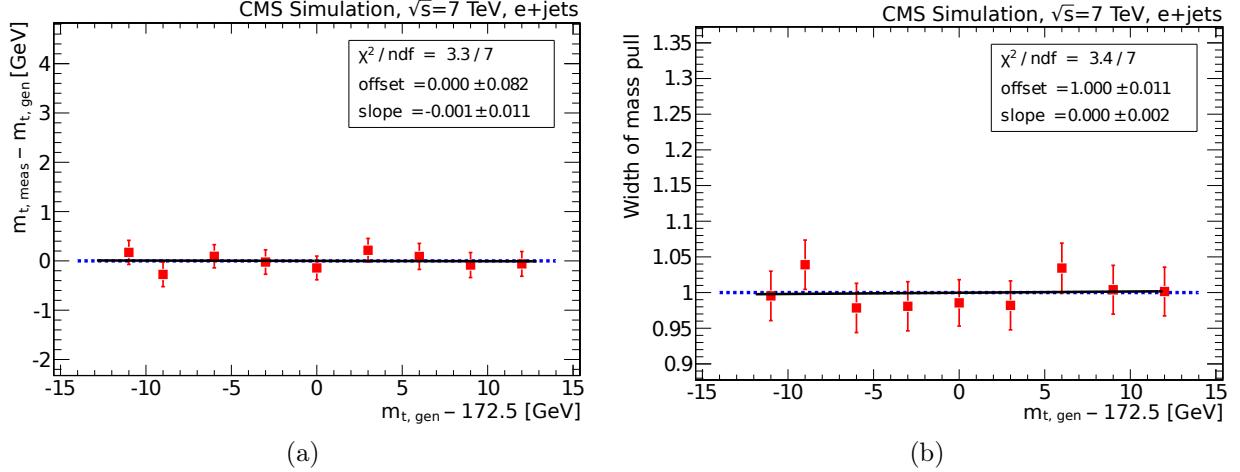


Figure 4.8: Closure test, showing the average mass bias (a) and width of the pull (b) as a function of true top quark mass after calibration.

this analysis. These include theoretical uncertainties on Monte Carlo modelling, assumptions made within the ideogram method, as well as imperfect understanding of detector effects and reconstruction algorithms.

Systematic uncertainties are calculated by running the analysis with tuned MC samples where systematic effects are varied by  $\pm 1$  standard deviation, and comparing the result with the nominal one. The obtained difference in fitted  $\Delta m_t$  is regarded as an uncertainty due to a given systematic effect. The total systematic error is calculated as a sum of all components in quadrature, under the assumption that all the effects are uncorrelated. A summary of all systematic uncertainties is given in Table 4.5, whereas all corresponding sources are described below.

**Fit calibration statistics.** To estimate the systematic effect of the calibration with pseudo-experiments described in Section 4.5, the statistical uncertainty of this calibration is propagated to the final measurement by regarding it as a systematic error. This uncertainty was evaluated with the central  $t\bar{t}$  MC sample with the top quark mass of 172.5 GeV.

**b-tagging.** The nominal (“medium”) working point of the CSVM algorithm is varied in order to reflect the uncertainty of the b-tagging efficiency of  $\sim 4\%$  [21]. The average effect of this change on the mass measurement (approximately 0.1 GeV) is quoted as a systematic uncertainty.

**Jet energy scale.** Since the top quark mass measurement is based on knowledge about at least four energetic jets in semileptonic  $t\bar{t}$  decay, it is very sensitive to any uncertainties in

Source of the systematic uncertainty	$\Delta m_t$ (GeV)
Fit calibration statistics	0.12
b-tagging	0.10
Jet Energy Scale (overall data/MC)	1.21
Factorisation scale	1.38
ME-PS matching threshold	0.93
Colour Reconnection	0.03
Underlying event	0.08
Non-tt background	0.61
Pile-up	0.22
PDF	0.15
Total	2.17

Table 4.5: Overview of the systematic uncertainties and their effects on the top quark mass measurement.

jet energies. Hence, the jet energy scale is a dominant systematic in this analysis, particularly depending on our understanding of the CMS detector. To quantify this systematic effect, the jet energies are varied within  $\pm 1\sigma$  uncertainty of the overall jet energy correction, obtained from a dedicated jet energy calibration and resolution study [55]. This  $p_T$ - and  $\eta$ -dependent uncertainty incorporates various contributions, including an offset due to noise and pile-up, absolute and relative discrepancies between data and MC, flavour-dependent corrections, etc. After applying asymmetric  $\pm 1\sigma$  uncertainties, an average systematic effect of 1.21 GeV on the fitted top mass was observed.

**Factorisation scale.** As mentioned in Section 4.1.2.2, the factorisation scale is varied by factors of 0.5 and 2 with respect to the nominal value used in the central generated samples. The procedure is performed for  $t\bar{t} + \text{jets}$ ,  $W + \text{jets}$  and  $Z + \text{jets}$  Monte Carlo samples independently as the factorisation scales are considered uncorrelated. To combine the systematic uncertainty, the average absolute size of the up and down effects are added in quadrature. Due to the lack of statistics available in variations of  $W + \text{jets}$  and  $Z + \text{jets}$  samples, this systematic uncertainty is one of the largest in the top quark mass measurement.

**ME-PS matching threshold.** The systematic effect of choice of ME-PS threshold, also mentioned in Section 4.1.2.2, is estimated by varying the default value of 20 GeV by factors of 0.5 and 2. Similarly to factorisation scale systematic uncertainty, the matching thresholds are varied independently for  $t\bar{t} + \text{jets}$ ,  $W + \text{jets}$  and  $Z + \text{jets}$  Monte Carlo samples and then added in quadrature. The same issue with the lack of statistics appears here, although to a

lesser extent.

**Colour reconnection.** Different modelling of colour reconnections [56] result in a systematic uncertainty that is estimated by varying the underlying event tune in simulation. The observed shifts in top quark mass are quoted as a systematic uncertainty.

**Underlying event.** The underlying event [20] modelling uncertainty was estimated by using different PYTHIA tunes with increased and decreased underlying event activity. The central tune, referred to as Perugia 2011, was compared to Perugia 2011 mpiHi and Perugia 2011 Tevatron tunes [57] with more and less activity, respectively. The largest mass shift ( $\approx 0.08$  GeV) represents the systematic uncertainty due to underlying event modelling.

**Background modelling.** As mentioned in the description of mass calibration (Section 4.5), the background in pseudo-experiments was estimated by using only  $W + \text{jets}$  MC sample. The impact of other background sources was studied separately and accounted for as a systematic uncertainty. Implementing the  $Z + \text{jets}$ , single top and QCD backgrounds in the calibration and separately varying their contributions by 30 % yielded the shifts in the measurement, the quadrature sum of which ( $\approx 0.61$  GeV) was quoted as a systematic uncertainty due to the background composition.

**Pile-up.** Since the number of pile-up vertices in data can only be known after data-taking process, it does not usually match the according number in MC simulation. To account for this discrepancy, the simulated events are re-weighted to match the observed pile-up distribution in data. The systematic uncertainty due to pile-up re-weighting was estimated by scaling the pile-up weights up and down, leading to shifts in top quark mass measurement. The largest of these shifts (approximately 0.22 GeV) was quoted as a systematic uncertainty.

**Parton distribution functions.** The CTEQ parton distribution functions (PDFs) [40] used in the simulation of Monte Carlo samples can be described by 22 orthogonal parameters. To estimate the systematic uncertainty due to the choice of PDFs, the up and down variations of these parameters were employed to calculate the impact on the measurement. Events in nominal  $t\bar{t} + \text{jets}$  sample were re-weighted according to the deviation of each PDF parameter from the central value. The average absolute values of the up and down shifts for each variation are then summed in quadrature, yielding the value of 0.15 GeV which we quote as a PDF systematic uncertainty.

## 4.7 Results

After extracting the top quark mass using combined likelihood fit, performing mass calibration and estimating systematic uncertainties as shown in previous sections, the following result was obtained in the electron plus jets channel:

$$m_t = 172.87 \pm 0.27 \text{ (stat.)} \pm 2.17 \text{ (syst.) GeV.} \quad (4.15)$$

Being a cross-check result, this measurement is consistent with the official CMS top quark mass measurement at 7 TeV in the lepton plus jets channel [29]:

$$m_t = 173.49 \pm 0.43 \text{ (stat.+JES)} \pm 0.98 \text{ (syst.) GeV.} \quad (4.16)$$

This result expectedly has a significantly lower uncertainty due to *in situ* measurement of the jet energy scale (JES) in a joint likelihood fit. It was obtained using the combined fit in the electrons plus jets and muon plus jets channels, whereas separate fits in each channel yield following results:

$$\begin{aligned} e+jets: m_t &= 173.72 \pm 0.66 \text{ (stat.+JES)} \pm 1.00 \text{ (syst.) GeV,} \\ \mu+jets: m_t &= 173.22 \pm 0.56 \text{ (stat.+JES)} \pm 1.06 \text{ (syst.) GeV.} \end{aligned} \quad (4.17)$$

All results are consistent with the latest Tevatron combination result of  $m_t = 173.18 \pm 0.56 \text{ (stat.)} \pm 0.75 \text{ (syst.) GeV}$  [58] and the latest LHC combination result of  $m_t = 173.29 \pm 0.23 \text{ (stat.)} \pm 0.92 \text{ (syst.) GeV}$  [59].

## 4.8 Summary

In this chapter, a top quark mass measurement in the electron plus jets channel was presented using the full 2011 dataset recorded by the CMS detector with a total integrated luminosity of  $5.0 \text{ fb}^{-1}$ . The mass extraction technique called ideogram method was described in detail, as well other technicalities of the analysis. The result is consistent with the official CMS top quark mass measurement in the lepton plus jets channel at 7 TeV, as well as latest combination results from the Tevatron and the LHC experiments.



# 5. Differential cross section measurement

Measurement of the differential cross section of the top quark pairs with respect to different variables is an important precision measurement, and provides a sensitive probe for new physics. For instance, deviations in the tail of the missing transverse energy distribution could signal new resonances with invisible particles produced in association with top quarks.

The analysis described in this chapter represents the top quark differential cross section measurement with respect to event-level (or global) distributions, including the missing transverse momentum ( $E_T^{\text{miss}}$ ), the scalar sum of jet transverse momenta ( $H_T$ ), the scalar sum of the transverse momenta of all objects in the event ( $S_T$ ), and both the transverse mass ( $M_T^W$ ) and transverse momentum ( $p_T^W$ ) of the leptonically decaying W boson produced in the top quark decay [60, 61]. The analysis is focused on semileptonic  $t\bar{t}$  decays, where a lepton is either an electron or a muon. The cross section measurement is performed to validate different Monte Carlo generators and theoretical predictions of the effects due to variations of various modelling parameters.

## 5.1 Data and Simulation

### 5.1.1 Data

This analysis uses the full 2012 dataset collected by the CMS detector at a centre of mass energy of 8 TeV, with a total integrated luminosity of  $19.7 \text{ fb}^{-1}$ . Only certified events were used in the data, i.e. from such periods of data-taking when all of the detector subsystems were functioning with no errors. Depending on the channel, the data were preselected with the single electron or single muon trigger. The preselection procedure as well as full event selection will be described in Section 5.2.

### 5.1.2 Monte Carlo samples

The Monte Carlo generators used in this work were presented in Section 4.1.2.1. The list of signal and background MC samples is shown in Table 5.1, and is largely similar to that from the top quark mass analysis. In addition, the signal  $t\bar{t} + \text{jets}$  sample is also available with POWHEG and MC@NLO generators in order to be able to differentiate between them. Moreover, to extend the statistics of the  $W/Z + \text{jets}$  samples, they were generated in four exclusive jet multiplicity bins:  $W/Z$  boson plus one/two/three and at least four jets.

Table 5.2 presents the list of QCD multi-jet background and  $\gamma + \text{jets}$  samples used in the estimation of the QCD background in the electron plus jets channel. The analogous set

of muon-enriched QCD samples used in the muon plus jets channel is shown in Table 5.3. Although the QCD background is estimated using data-driven techniques in both channels (see Section 5.3), the MC simulation is still used for normalisation purposes.

Process	Generator	$\sigma$ (pb)	# events	$\int \mathcal{L} dt$ (fb $^{-1}$ )
$t\bar{t} + \text{jets}, m_t = 172.5 \text{ GeV}$	MADGRAPH	245.8	6854416	27.9
	POWHEG	245.8	21675970	88.2
	MC@NLO	245.8	32706581	133.1
W + jets ( $W \rightarrow l\nu$ )	MADGRAPH	5400.0	23141598	4.3
			1750.0	34044921
			519.0	15539503
			214.0	13349346
				62.4
$Z/\gamma^* \rightarrow l^+l^- + \text{jets}, m(ll) > 50 \text{ GeV}$	MADGRAPH	561.0	24045248	42.9
			181.0	21852156
			51.1	11015445
			23.04	6402827
				277.9
Single top	POWHEG	55.5	3758227	67.7
			30.0	1935072
			3.89	259961
			1.76	139974
			11.18	497658
			11.18	493460

Table 5.1: Signal and background Monte Carlo samples with cross sections at  $\sqrt{s} = 8 \text{ TeV}$ , numbers of generated events and corresponding integrated luminosities.

### 5.1.3 Pile-up reweighting

Since the number of simulated pile-up interactions does not necessarily represent the same distribution in data, a reweighting procedure for MC events is required. To produce the weights, measured instantaneous luminosity is used to obtain a distribution of true pile-up vertices in data, i.e. before any vertex reconstruction inefficiencies. This distribution is then divided by the simulated one, and for each number of pile-up interactions, the weight is calculated. Figure 5.1 shows the number of interactions before and after reweighting for both electron and muon channels. Evidently, the reweighting procedure helps to achieve a good agreement between data and simulation.

Process	Generator	$\sigma$ (pb)	filter efficiency	# events	$\int \mathcal{L} dt$ (fb $^{-1}$ )
QCD ( $e/\gamma$ enriched)	PYTHIA				
$20 \text{ GeV} < \hat{p}_T < 30 \text{ GeV}$		$2.886 \times 10^8$	$1.01 \times 10^{-2}$	34339883	$1.2 \times 10^{-2}$
$30 \text{ GeV} < \hat{p}_T < 80 \text{ GeV}$		$7.433 \times 10^7$	$6.21 \times 10^{-2}$	32537408	$7.0 \times 10^{-3}$
$80 \text{ GeV} < \hat{p}_T < 170 \text{ GeV}$		$1.191 \times 10^6$	0.154	34542763	0.19
$170 \text{ GeV} < \hat{p}_T < 250 \text{ GeV}$		30 990.0	0.148	22862259	5.0
$250 \text{ GeV} < \hat{p}_T < 350 \text{ GeV}$		4250.0	0.131	32505856	58.4
$\hat{p}_T > 350 \text{ GeV}$		810.0	0.11	33981105	381.4
QCD ( $b/c \rightarrow e\nu$ )	PYTHIA				
$20 \text{ GeV} < \hat{p}_T < 30 \text{ GeV}$		$2.886 \times 10^8$	$5.8 \times 10^{-4}$	1740229	$1.0 \times 10^{-2}$
$30 \text{ GeV} < \hat{p}_T < 80 \text{ GeV}$		$7.433 \times 10^7$	$2.25 \times 10^{-3}$	2048152	$1.2 \times 10^{-2}$
$80 \text{ GeV} < \hat{p}_T < 170 \text{ GeV}$		$1.191 \times 10^6$	$1.09 \times 10^{-2}$	1945525	0.15
$170 \text{ GeV} < \hat{p}_T < 250 \text{ GeV}$		30 990.0	$2.04 \times 10^{-2}$	1948112	3.1
$250 \text{ GeV} < \hat{p}_T < 350 \text{ GeV}$		4250.0	$2.43 \times 10^{-2}$	2026521	19.6
$\hat{p}_T > 350 \text{ GeV}$		810.0	$2.95 \times 10^{-2}$	1948532	81.5
$\gamma + \text{jets}$	MADGRAPH				
$200 \text{ GeV} < H_T < 400 \text{ GeV}$		960.5		1	10479625
$H_T > 400 \text{ GeV}$		107.5		1	1611963

Table 5.2: QCD multi-jet background and  $\gamma + \text{jets}$  MC samples used in the electron plus jets channel with cross sections at  $\sqrt{s} = 8 \text{ TeV}$ , numbers of generated events and corresponding integrated luminosities.

Process	Generator	$\sigma$ (pb)	filter efficiency	# events	$\int \mathcal{L} dt$ (fb $^{-1}$ )
QCD ( $\mu$ enriched)	PYTHIA				
$15 \text{ GeV} < \hat{p}_T < 20 \text{ GeV}$		$7.022 \times 10^8$	$3.9 \times 10^{-3}$	1722681	$6.3 \times 10^{-4}$
$20 \text{ GeV} < \hat{p}_T < 30 \text{ GeV}$		$2.87 \times 10^8$	$6.5 \times 10^{-3}$	8486904	$4.5 \times 10^{-3}$
$30 \text{ GeV} < \hat{p}_T < 50 \text{ GeV}$		$6.609 \times 10^7$	$1.22 \times 10^{-2}$	9560265	$1.2 \times 10^{-2}$
$50 \text{ GeV} < \hat{p}_T < 80 \text{ GeV}$		$8.082 \times 10^6$	$2.18 \times 10^{-2}$	10365230	$5.9 \times 10^{-2}$
$80 \text{ GeV} < \hat{p}_T < 120 \text{ GeV}$		$1.024 \times 10^6$	$3.95 \times 10^{-2}$	9238642	0.23
$120 \text{ GeV} < \hat{p}_T < 170 \text{ GeV}$		$1.578 \times 10^5$	$4.73 \times 10^{-2}$	8501935	1.1
$170 \text{ GeV} < \hat{p}_T < 300 \text{ GeV}$		34 020.0	$6.76 \times 10^{-2}$	7669947	3.3
$300 \text{ GeV} < \hat{p}_T < 470 \text{ GeV}$		1757.0	$8.64 \times 10^{-2}$	7832261	51.6
$470 \text{ GeV} < \hat{p}_T < 600 \text{ GeV}$		115.2	0.102	3783069	322.0
$600 \text{ GeV} < \hat{p}_T < 800 \text{ GeV}$		27.01	0.0996	4119000	1531.1
$800 \text{ GeV} < \hat{p}_T < 1000 \text{ GeV}$		3.57	0.1033	4107853	11 139.0
$\hat{p}_T > 1000 \text{ GeV}$		0.774	0.1097	3873970	45 625.6

Table 5.3: QCD multi-jet background MC samples used in the muon plus jets channel with cross sections at  $\sqrt{s} = 8 \text{ TeV}$ , numbers of generated events and corresponding integrated luminosities.

Process	Generator	$\sigma$ (pb)	# events	$\int \mathcal{L} dt$ (fb $^{-1}$ )
$t\bar{t} + \text{jets}$	MADGRAPH	245.8	5476728	10.2
			5306710	25.5
			5387181	25.4
			5009488	23.4
$W + \text{jets } (W \rightarrow l\nu)$	MADGRAPH	29690	21364637	0.7
			20976082	0.7
			20719363	0.6
			20784770	0.6
$Z + \text{jets } (Z \rightarrow ll)$	MADGRAPH	2888	2112387	0.6
			1985529	0.7
			1934901	0.6
			2170270	0.7

Table 5.4: Systematic MC samples with cross sections at  $\sqrt{s} = 8$  TeV, numbers of generated events and corresponding integrated luminosities. Factorisation scale  $Q$  and matching threshold systematic uncertainties are estimated with variations of  $t\bar{t} + \text{jets}$ ,  $W + \text{jets}$  and  $Z + \text{jets}$  samples.

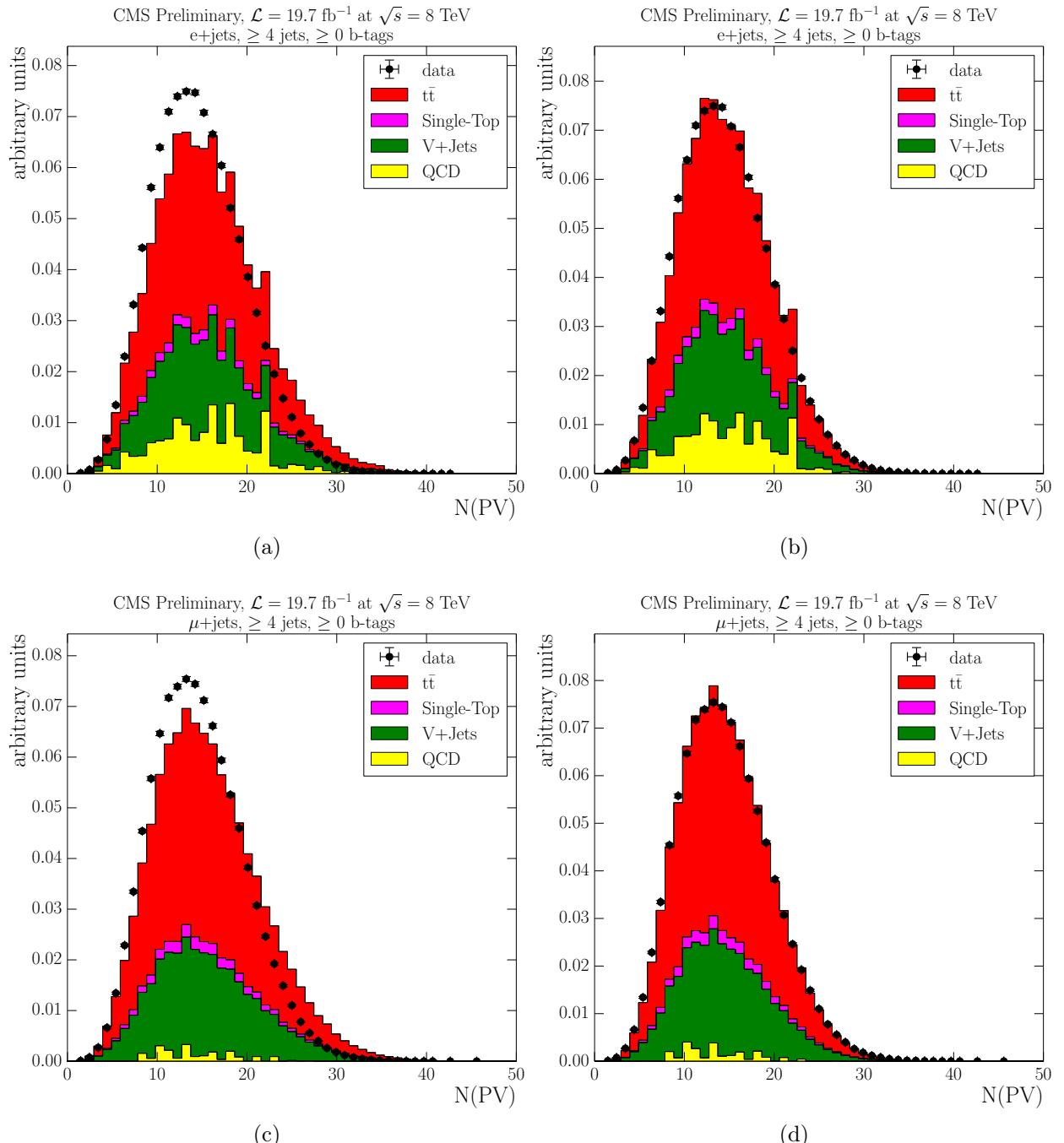


Figure 5.1: Number of reconstructed vertices per event before (left) and after pile-up reweighting (right) in the electron channel (top) and in the muon channel (bottom). Both data and sum of the MC samples are normalised to unit area.

To estimate the number of pile-up vertices in data, the measured instantaneous luminosity in each bunch crossing is multiplied by an average total inelastic proton-proton cross section. Therefore, two sources of uncertainty arise from these factors: the luminosity uncertainty, measured to be 4.4 % by a dedicated study [62], and the uncertainty on the total inelastic cross section. To obtain the total cross section value, the CMS measurement of  $68.0 \pm 4.5$  mb based on 7 TeV data [63] was extrapolated to the 8 TeV value of 69.3 mb. The recommended uncertainty on this value, set to be 5 %, covers the modelling and physics aspects of pile-up simulation that have not been fully studied up to date.

The effect of the  $\pm 5$  % variations in total inelastic cross section on the estimate of the true number of vertices in data is shown in Figure 5.2, whereas its effect on the simulated distributions is shown in Figure 5.3 for both electron and muon channels.

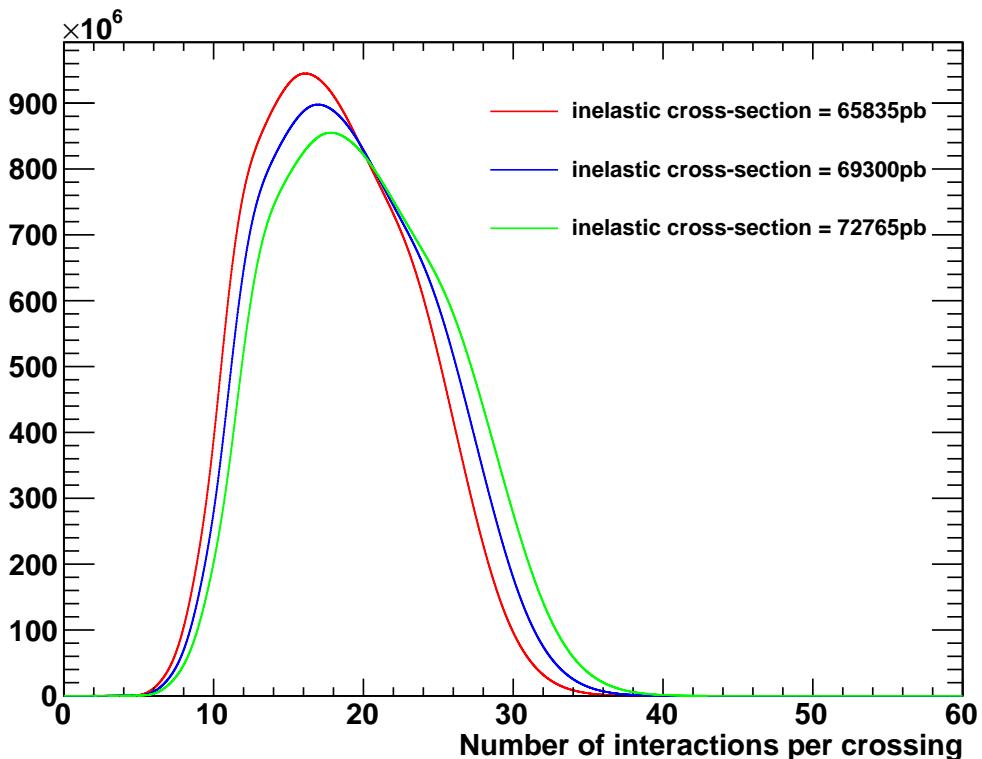


Figure 5.2: Number of expected vertices per event for central and  $\pm\sigma$  variations of the total inelastic cross section for the 2012 data.

### 5.1.4 b-tagging corrections

In order to mitigate the background (QCD in particular), the b-tagging is applied in the event selection of this analysis. Similarly to the top quark mass analysis, the CSV

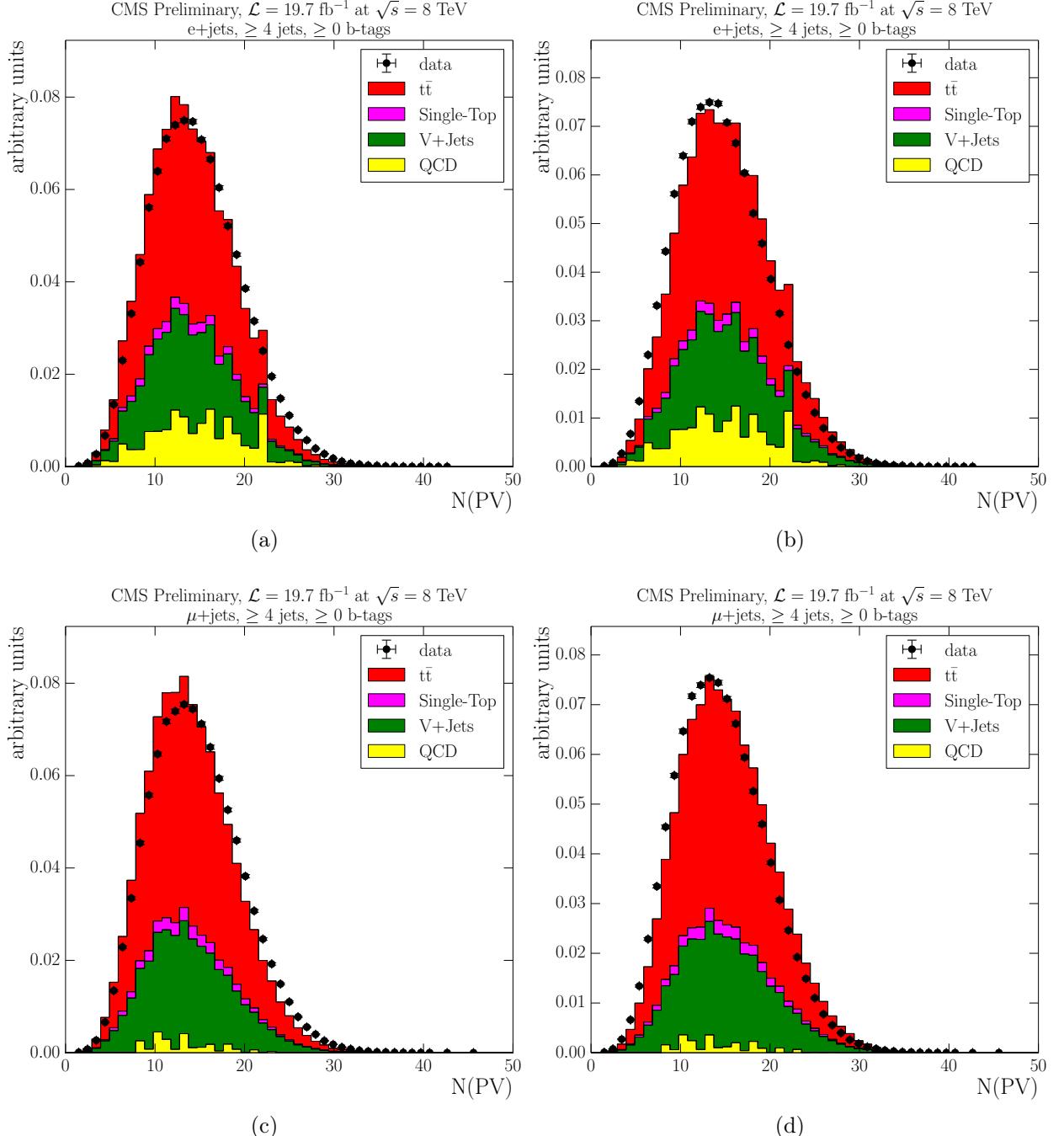


Figure 5.3: Number of reconstructed vertices per event for the  $-\sigma$  (left) and the  $+\sigma$  variations (right) of the pile-up reweighting procedure for the electron channel (top) and the muon channel (bottom).

algorithm with the medium working point providing  $\sim 70\%$  efficiency and  $\sim 1\%$  mis-tag rate (Section 2.4.3.3) was used. However, in contrast with the 2011 analysis where the b-tagging corrections were implemented in the construction of the event likelihood, in this work the scale factors are applied to Monte Carlo events. This is done in order to correct for the mismatch between the algorithm efficiency seen in data and MC [64].  $p_T$  and  $\eta$ -dependent scale factors were derived by the CMS b-tag Physics Object Group [65] and are used to reweight the simulated events. The overall jet content of the event is taken into account, such as the number of jets originating from light and heavy quarks.

Figure 5.4 shows the results of the reweighting procedure for electron and muon channels. It can be seen that the numbers of simulated events in b-tag multiplicity bins of 0 and 1 are scaled up, whereas for other bins the number of MC events is scaled down. The systematic uncertainty due to b-tagging corrections is calculated by applying  $\pm 1\sigma$  variations of the scale factors in the reweighting procedure.

## 5.2 Event Selection

As in the top quark mass analysis, the event selection in this work is based on the standard CMS selection optimised for the Standard Model  $t\bar{t}$  production. It also exploits the semileptonic signature of a  $t\bar{t}$  decay with exactly one isolated lepton and at least four jets, two of which are b-tagged. One of the main goals of the analysis is the measurement of the missing transverse energy distribution represented by a neutrino from the leptonic decay of the W boson, however, there is no specific requirement on the  $E_T^{\text{miss}}$  in the selection.

Comparing to the 2011 event selection described in Section 4.2, enhanced pile-up subtraction and lepton identification techniques are used. Improved understanding of the detector in form of updated alignment and calibration constants resulted in better resolution of jets and  $E_T^{\text{miss}}$ . The preselection (or skimming) step used in order to reduce the number of events for local analysis processing is identical to that of 2011 analysis. However, additional filters are used to reject the events with artificial  $E_T^{\text{miss}}$  caused by noise or known problems with one of the detector subsystems. These filters are discussed in Section 5.2.3.

Two semileptonic  $t\bar{t}$  decay modes are explored in this analysis: the electron plus jets and the muon plus jets channels. The event selection for each of these decay modes is described below.

### 5.2.1 Electron plus jets channel

Similarly to the electron plus jets channel in the top quark mass analysis, the selection consists of the following steps:

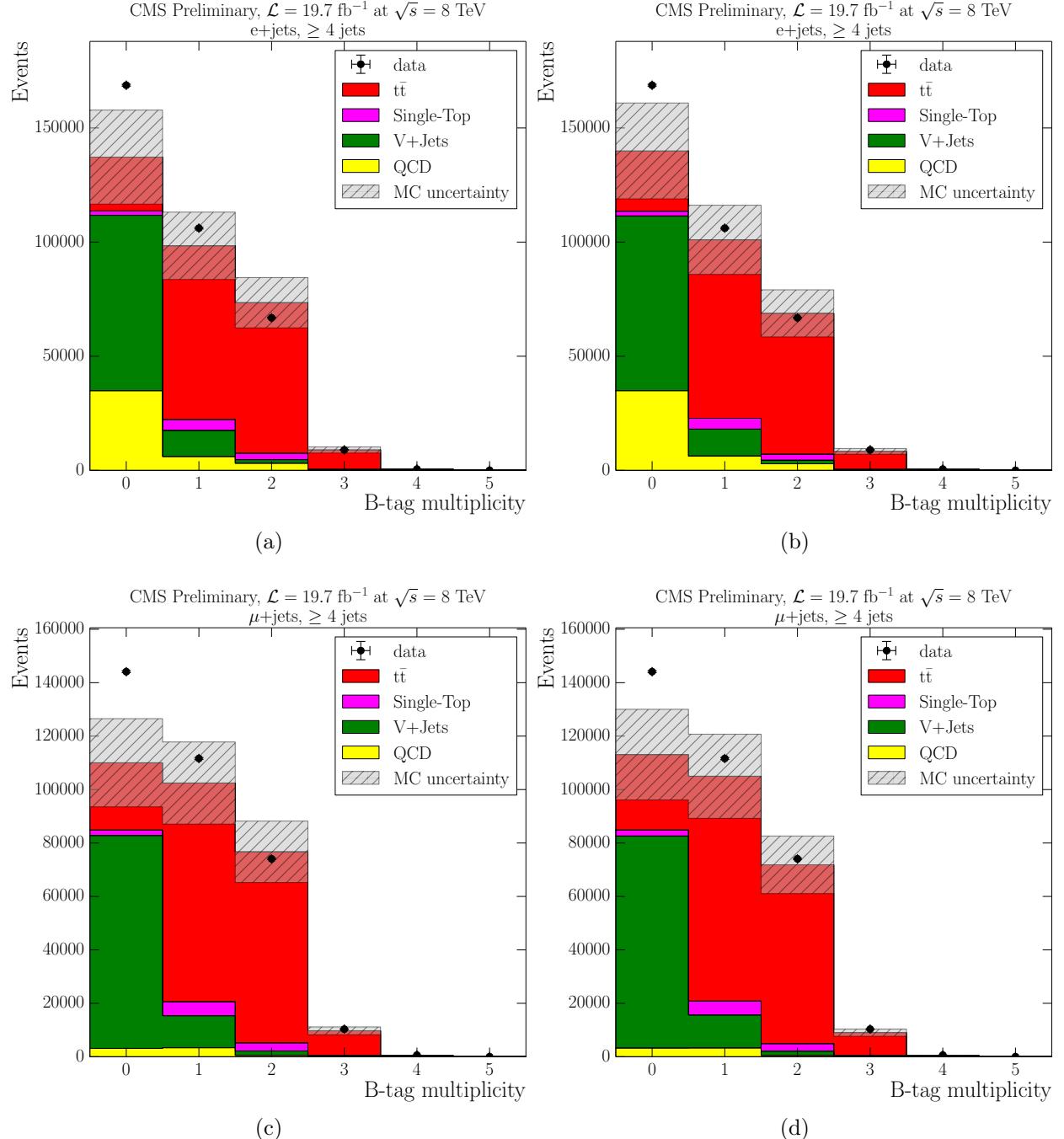


Figure 5.4: b-tag multiplicity before (left) and after applying b-tag scale factors (right) in the electron channel (top) and the muon channel (bottom).

1. preselection;
2. trigger;
3. electron candidate selection;
4. muon veto;
5. dilepton veto;
6. conversion veto;
7. jet selection;
8. b-tagging.

## Trigger

As opposed to 2011 analysis using cross triggers, a single electron trigger is used in this analysis. Referred to as HLT\_Ele27\_WP80, this trigger has a lepton  $p_T$  threshold of 27 GeV and “WP80” lepton isolation requirement (see Table 3.1), but no specific jet requirements. In contrast to 2011 data-taking, the single electron trigger was unprescaled for the whole period of 2012 running, despite a comparatively large rate. This decision was made due to simplicity of the trigger selection, requiring flat efficiency corrections.

## Electron candidate selection

Similarly to 2011 selection, exactly one electron candidate passing is required in the event, passing the following criteria:

- $E_T > 30$  GeV;
- $|\eta| < 2.5$ , excluding the ECAL barrel-endcap transition regions of  $1.4442 < |\eta| < 1.566$ ;
- MVA electron ID (Section 2.4.1) with a discriminator cut of  $> 0.5$ ;
- transverse impact parameter with respect to the primary vertex  $d_{xy} < 0.02$  cm;
- PF-based relative isolation (Section 2.4.1.2)  $\text{reliso} < 0.1$  with a  $\Delta R$  cone of 0.3.

## Muon veto

To reduce contamination from other  $t\bar{t}$  decay modes (i.e. muon and dilepton channels), events containing a global or a tracker muon (see Section 2.4.2) are rejected. The requirements on the muon are identical to those of 2011 selection:  $p_T > 10$  GeV,  $|\eta| < 2.5$  and  $\text{reliso} < 0.2$  with a  $\Delta R$  cone of 0.4.

## Dilepton veto

To reject events with additional electrons, dilepton veto is applied. Events are rejected if they contain a second electron candidate with the same ID and  $\eta$  requirements, but looser  $E_T$  and  $\text{reliso}$  criteria ( $> 20$  GeV and  $< 0.15$ , respectively).

### Conversion veto

As described in Section 2.4.1.3, a vertex fit technique combined with the number of missing hits in the tracker is used for photon conversion rejection.

### Jet selection and b-tagging

To help further reduce the background, at least four jets with  $p_T > 30$  GeV and  $|\eta| < 2.5$  are required in the event. These jets are required to pass the loose PF jet ID (see Section 2.4.3). Before the multiplicity and identification requirements, the collection of jets is cleaned against the signal electron, i.e. any jet within a  $\Delta R$  cone of 0.3 around the signal lepton is excluded from the analysis. This is done due to the fact that an electron may be reconstructed as a jet by leaving a characteristic signature in calorimeters. Any additional jets with  $p_T > 20$  GeV and the same cleaning and ID requirements are also used for calculation of  $H_T$  and  $S_T$  variables (Section 5.5.1). Finally, the combined secondary vertex (CSV) b-tagging algorithm with the medium working point (Section 2.4.3.3) is used to identify the jets originating from b-quarks in the  $t\bar{t}$  decay. At least two b-tagged jets are required in this analysis.

## 5.2.2 Muon plus jets channel

The muon channel event selection follows similar strategy as that of the electron channel:

1. preselection;
2. trigger;
3. muon candidate selection;
4. second muon veto;
5. electron veto;
6. jet selection;
7. b-tagging.

### Trigger

An isolated single muon trigger, referred to as HLT\_IsoMu24\_eta2p1, is used in this channel. Partially suggested by the name of the trigger path, the following criteria on the muon are imposed at the HLT level:

- $p_T > 24$  GeV;
- $|\eta| < 2.1$ ;
- Combined tracker and calorimeter-based relative isolation (Section 2.4.1.2)  $\text{reliso} < 0.15$  with a  $\Delta R$  cone of 0.3.

### Muon candidate selection

For selection of a muon candidate, exactly one muon satisfying the following requirements is required:

- $p_T > 26$  GeV;
- $|\eta| < 2.1$ ;
- PF-based relative isolation (Section 2.4.1.2)  $\text{reliso} < 0.12$  with a  $\Delta R$  cone of 0.4;
- PF-muon ID criteria (Section 2.4.2):
  - identified as a particle flow muon and as a global muon;
  - a normalised of the global muon fit  $\chi^2 < 10$ ;
  - number of muon chamber hits  $> 0$ ;
  - number of muon stations with muon segments  $> 1$ ;
  - transverse impact parameter w.r.t. the primary vertex  $d_{xy} < 0.02$  cm;
  - longitudinal distance of the tracker track w.r.t. the primary vertex  $d_z < 0.5$  cm;
  - number of hits in the pixel detector  $> 0$ ;
  - number of hits in the tracker layers  $> 5$ .

### Second muon veto

Events containing an additional loose muon are rejected, with loose ID requirements as follows:

- $p_T > 10$  GeV;
- $|\eta| < 2.5$ ;
- PF-based relative isolation (Section 2.4.1.2)  $\text{reliso} < 0.2$  with a  $\Delta R$  cone of 0.4;
- identified as a particle flow muon, and either global or tracker muon.

These criteria are identical to the muon veto requirements in the electron channel.

### Electron veto

To further reduce background contamination, events with electrons are rejected. The requirements on additional loose electron are the same to the dilepton veto in the electron channel.

### Jet selection and b-tagging

Finally, the jet requirements are applied. These are also identical to those of the electron channel, mentioned above.

#### 5.2.3 Additional filters for missing transverse energy

There is no specific requirement on missing transverse energy in the selection, however, a set of corrections is applied to  $E_T^{\text{miss}}$ . As described in Section 2.4.4, jet energy corrections are propagated to  $E_T^{\text{miss}}$ , and it is also corrected for pile-up and various detector effects causing  $\phi$  modulation.

Additionally, a set of optional filters, prescribed for analyses with high sensitivity to  $E_T^{\text{miss}}$ , was applied in the selection for both electron and muon channels. Most of these filters were put in place after discovering events with anomalously high missing transverse energy. A short description of each filter is given below.

**CSC beam halo filter.** Protons in the LHC beams can interact with residual particles of beam collimators, which produces secondary particle showers referred to as beam halo [66]. This machine-induced background can significantly contaminate collision events, particularly affecting the  $E_T^{\text{miss}}$  observable due to specific trajectories of halo muons. A dedicated algorithm based on exploiting the halo signatures in the Cathode Strip Chambers (CSCs) helps to mitigate this background.

**HCAL laser filter.** The HCAL laser system is used for calibration and monitoring the detector response [5]. During the 2012 data-taking, laser pulses were accidentally fired and consequently polluted a small fraction of the recorded physics dataset. Using the fact that unlike physics events, laser pulses produce hits in calibration channels of the HCAL, the dataset is filtered from this contamination.

**ECAL dead cell filter.** Both ECAL endcap and barrel regions are known to have faulty crystals producing too much noise in detector readouts. There are also crystals corresponding to front-end cards with defective data link. All these channels are masked in reconstruction and constitute to about 1 % of the total amount. However, they can still cause a significant amount of energy deposits to be lost, therefore contributing to fake  $E_T^{\text{miss}}$ . One of the approaches to tackle this issue uses the so-called trigger primitive information, i.e. digital quantities produced by the Level 1 trigger electronics [22]. Another approach uses the energy deposits surrounding the masked channels (boundary energy). These methods allow to filter events with the estimated energy leakage above certain threshold ( $\sim 10$  GeV).

**Tracking failure filter.** Some events have been observed to have a lack of tracks corresponding to standard or large calorimeter deposits. There are two understood sources of this misreconstruction: in a first type of such events, the tracking algorithm (explained in Section 2.4.1) halts for some of its iterations due to an exceeding number of calorimeter clusters. In a second type the hard collision happens at a distance of approximately 75 cm from the nominal interaction point. A cut on the summed  $p_T$  of the tracks originating from the primary vertices passing the preselection criteria, divided by the summed  $p_T$  of all jets in the event, was found to clearly distinguish between pathological and unaffected events. A cut value of 10 % is imposed.

**Bad ECAL endcap super-cluster filter.** In 2012, two  $5 \times 5$  super-clusters in the ECAL endcap regions were observed to occasionally produce anomalous pulses, leading to artificially high  $E_T^{\text{miss}}$ . The source of this issue can be related to the High-Voltage system, but is subject to further investigation. A filter removing the problematic events is used, configured to cut on the total energy of reconstructed hits not passing the nominal quality criteria in the two super-clusters, with the cut value of 1 TeV.

**ECAL laser correction filter.** Lead tungstate crystals of the ECAL naturally lose their transparency with radiation exposure. To maintain high precision, this effect is corrected for by frequently injecting laser pulses and reading out the response in order to calibrate each crystal [5]. During 2012 data-taking, a handful of crystals unphysically large values of laser corrections, resulting in anomalously high  $E_T^{\text{miss}}$  in recorded events. The filter removes such events from the dataset.

**Strip tracker noise filter.** A few events were affected by a large coherent noise in the silicon strip tracker, which resulted in a much larger number of the strip clusters compared to according pixel cluster multiplicity in the tracking algorithm. A dedicated filter rejects events with this anomaly.

### 5.3 Data-driven QCD estimation

The QCD multi-jet events constitute an important background to semileptonic  $t\bar{t}$  decays, and one of the most difficult to model (Section 1.2.3). The two b-tagged jets requirement in the event selection is mainly motivated by the necessity to mitigate the effect of this particular background. In this analysis, the final QCD yield is obtained from the fitting process as explained in Section 5.5. However, the initial estimate is performed using the data-driven techniques. Only the shape of the distribution is of importance here as the

normalisation is determined in the fitting process. The QCD estimation methods for both electron and muon channel are described in detail hereafter.

### 5.3.1 QCD multi-jet events with electrons

Two control regions have been investigated in the electron channel:

- conversion control region: conversion veto is inverted in the full event selection;
- non-isolated electron control region: the electron isolation requirement is changed from  $\text{reliso} < 0.1$  to  $\text{reliso} > 0.2$ .

Additionally, only events with no b-tagged jets are taken into account to increase the statistics and QCD purity in both control regions. Both control regions are presented on Figure 5.5, and the QCD shapes are compared for the electron absolute pseudorapidity distribution as this observable is used in the fitting procedure. It can be seen that the distributions are not well described by the Monte Carlo, e.g. occasional high peaks in the QCD MC represent the lack of statistics in simulation. Moreover, the two control regions have considerably different shapes in the electron  $|\eta|$  distributions. For the conversion control region, the shape is mostly flat in the barrel region ( $|\eta| < 1.5$ ) and peaks in the endcap region ( $|\eta| > 1.5$ ), which is expected as the conversions are more likely to happen in the region with higher material budget. In contrast, the non-isolated control region peaks in the pseudorapidity region of  $1.0 < |\eta| < 1.5$ . Such behaviour is correlated with the jet  $|\eta|$  distribution, since this control region favours non-isolated electrons coming from heavy quark decays.

Due to the fact that the non-isolated control region has a significant contamination from signal and W/Z+jets events, the conversion region was chosen to produce the QCD template for the electron channel. However, since the relative contributions of the QCD events from both control regions in the final selected sample are not known, the real QCD shape can differ considerably from the prediction. Therefore, the non-isolated region is used to estimate the systematic uncertainty due to the QCD shape choice.

### 5.3.2 QCD multi-jet events with muons

For the QCD shape extraction in the muon plus jets channel, the inverted isolation cut is applied in the event selection: at least one muon with  $\text{reliso} > 0.3$  is required, and events with isolated muons with  $\text{reliso} < 0.3$  are vetoed. Similarly to the electron channel, only events with no b-tagged jets are taken into account. Additionally, to increase the available statistics, events with only three jets are also considered ( $\geq 3$  jet bin).

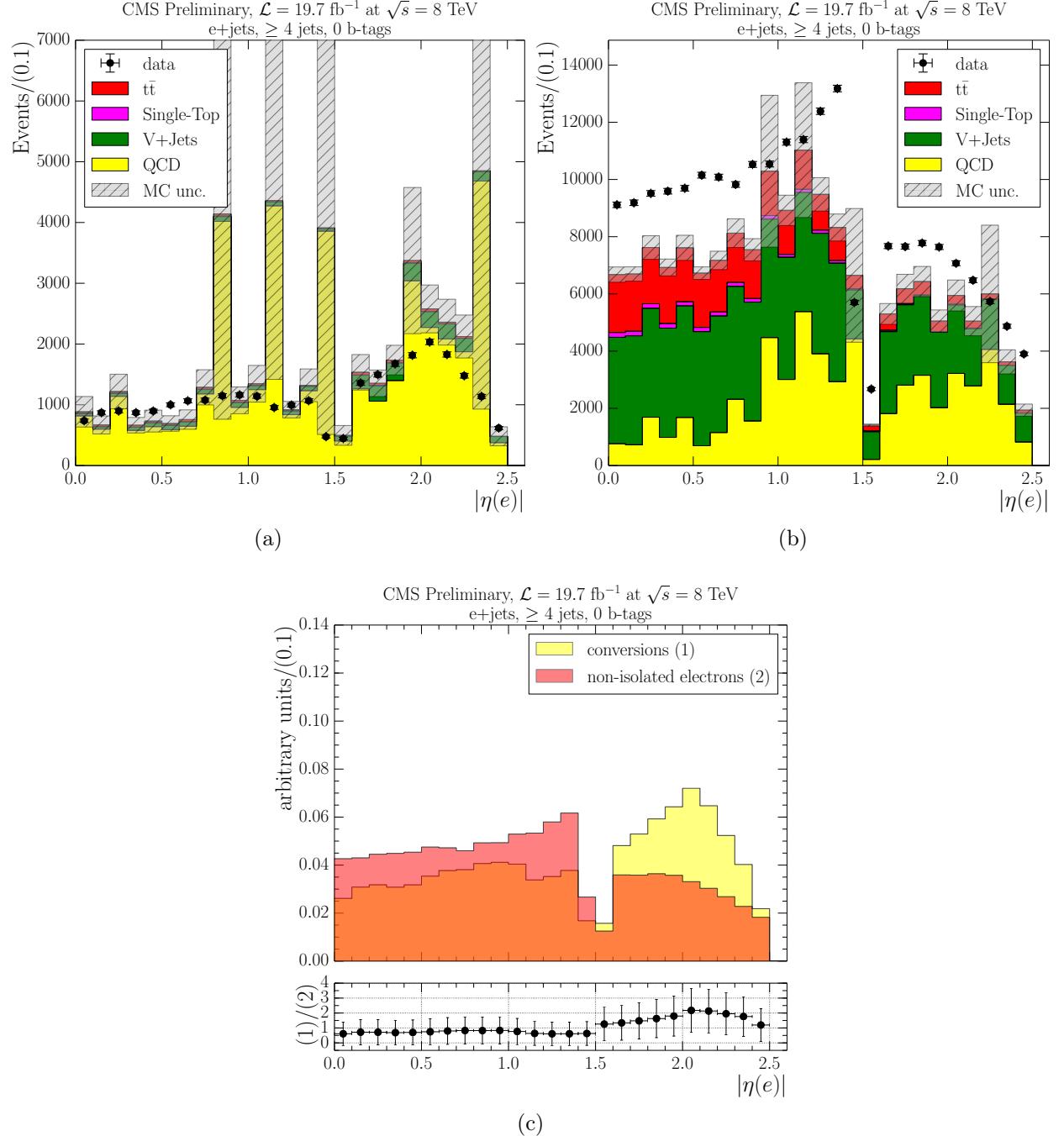


Figure 5.5: Distribution of the electron  $|\eta|$  for conversion selection (a), non-isolated electron selection (b) and the shape comparison of both QCD control region in data (c) in electron plus jets events with at least four jets and 0 b-tagged jets.

The relative isolation distribution with the cut value and the muon  $|\eta|$  for the non-isolated region are shown on Figure 5.6. Clearly, the control region is dominated by QCD events and generally there is a reasonable agreement between data and simulation. The contamination from other processes is subtracted using Monte Carlo.

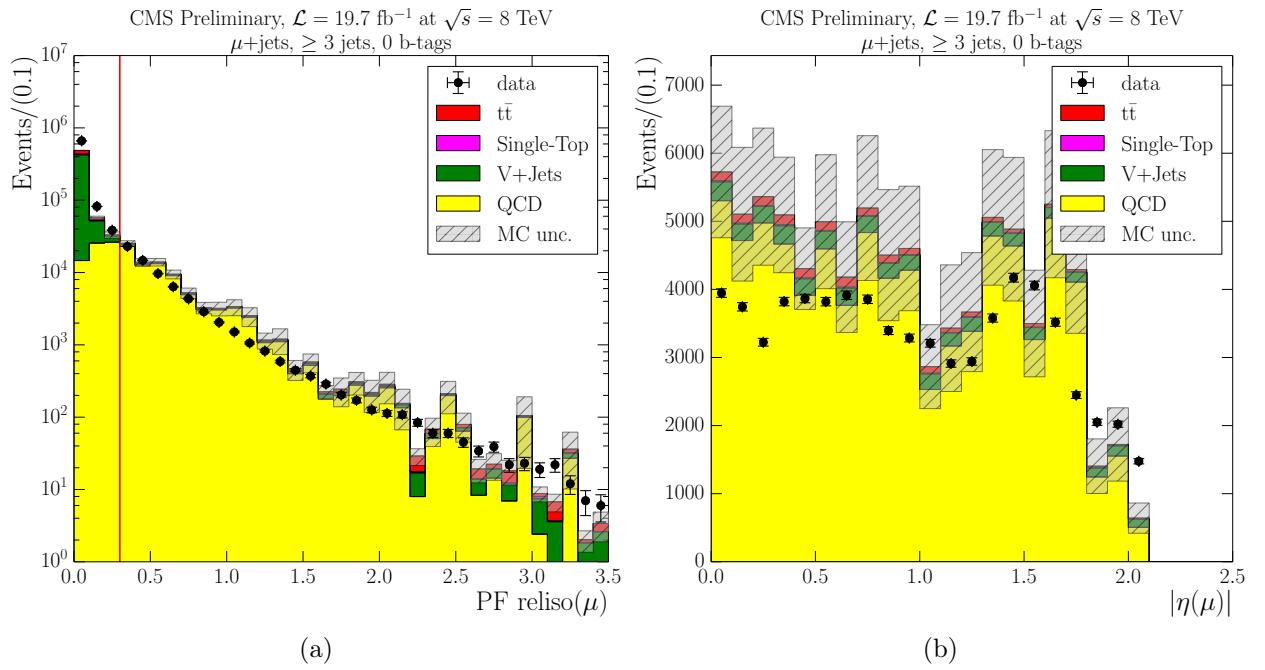


Figure 5.6: Muon relative isolation with no isolation cut (a) and muon  $|\eta|$  (b) for non-isolated ( $\text{reliso} > 0.3$ ) selection in muon plus jets events with at least three jets and 0 b-tagged jets. The vertical red line on plot (a) shows the  $\text{reliso}$  cut value of 0.3.

## 5.4 Data/MC comparison

In this section, the agreement between simulation and data for various kinematic distributions is reviewed after the final event selection. The event yields for all selection steps described in Section 5.2 are shown in Tables 5.5 and 5.6 for the electron and muon channels, respectively. All background events are normalised to the number of expected events for the total integrated luminosity of  $19.7 \text{ fb}^{-1}$ .

The QCD multi-jet background events shown in plots are derived using the data-driven methods explained previously. Figure 5.7 shows the lepton  $p_T$  and  $|\eta|$  distributions for both electrons and muons. In Figure 5.8, missing transverse energy distributions are shown, also including the logarithmic scale for higher values of  $E_T^{\text{miss}}$ . Figure 5.9 shows the azimuthal angle of  $E_T^{\text{miss}}$  ( $\phi(E_T^{\text{miss}})$ ) and jet multiplicity distributions for both channels. The  $t\bar{t}$  uncer-

Table 5.5: Numbers of expected and observed events from simulation and data for the electron channel out of the box, i.e. before the fitting process.

Selection step	t̄ + jets	W + jets	Z + jets	Single top	QCD	Sum MC	Data
Preselection	1346476 ± 1028	1727900 ± 1138	380944 ± 234	169689 ± 262	130308513 ± 475118	133933524 ± 475120	13042702
Event cleaning/HLT	385009 ± 553	516544 ± 600	147515 ± 143	37308 ± 127	3854386 ± 83441	4940764 ± 83445	5846672
One isolated electron	342394 ± 522	446691 ± 548	100243 ± 117	33000 ± 120	578895 ± 29802	1501225 ± 29812	1688811
Muon veto	323759 ± 508	446579 ± 548	99805 ± 117	32301 ± 118	578860 ± 29802	1481306 ± 29812	1668851
Dilepton veto	319019 ± 504	446469 ± 548	73876 ± 101	32117 ± 118	578821 ± 29802	1450305 ± 29812	1628009
Conversion veto	310863 ± 498	430191 ± 538	70929 ± 99	31287 ± 116	321292 ± 21826	1164563 ± 21839	1396638
≥ 1 jets	310863 ± 498	430179 ± 538	70928 ± 99	31287 ± 116	321292 ± 21826	1164550 ± 21839	1396638
≥ 2 jets	310838 ± 498	429195 ± 536	70723 ± 99	31279 ± 116	320506 ± 21819	1162543 ± 21831	1396506
≥ 3 jets	306406 ± 494	386892 ± 493	63653 ± 90	29953 ± 114	226938 ± 16321	1013843 ± 16337	1215535
≥ 4 jets	174701 ± 373	76646 ± 179	13323 ± 35	9765 ± 67	44203 ± 4340	318640 ± 4361	351194
≥ 1 b-tagged jets	148287 ± 341	11280 ± 69	2152 ± 14	7675 ± 58	9304 ± 2055	178700 ± 2085	182481
≥ 2 b-tagged jets	70103 ± 228	1320 ± 23	337 ± 5	2900 ± 35	3035 ± 1789	77697 ± 1804	76379

Table 5.6: Numbers of expected and observed events from simulation and data for the muon channel out of the box, i.e. before the fitting process.

Selection step	t̄ + jets	W + jets	Z + jets	Single top	QCD	Sum MC	Data
Preselection	1346476 ± 1028	1727900 ± 1138	380944 ± 234	169689 ± 262	104079124 ± 236784	107704135 ± 236789	20284215
Event cleaning/HLT	443056 ± 581	644505 ± 688	148143 ± 142	44727 ± 135	1664036 ± 33747	2944468 ± 33759	3063569
One isolated muon	360897 ± 521	473825 ± 556	88283 ± 106	35546 ± 121	83535 ± 5833	1042088 ± 5885	1327738
Second muon veto	353646 ± 516	473776 ± 556	53263 ± 82	35260 ± 120	82863 ± 5823	998811 ± 5874	1254896
Electron veto	336474 ± 503	473631 ± 556	52935 ± 82	34633 ± 119	82841 ± 5823	980516 ± 5873	1237495
≥ 1 jets	336474 ± 503	473622 ± 556	52934 ± 82	34633 ± 119	82841 ± 5823	980507 ± 5873	1237495
≥ 2 jets	336457 ± 503	472377 ± 554	52768 ± 81	34620 ± 119	81211 ± 5777	977436 ± 5827	1237428
≥ 3 jets	332863 ± 501	425334 ± 507	47493 ± 74	33146 ± 117	33505 ± 2726	872343 ± 2822	1108272
≥ 4 jets	188466 ± 377	83248 ± 184	10144 ± 29	10556 ± 67	7006 ± 1155	299422 ± 1231	340786
≥ 1 b-tagged jets	160223 ± 344	12341 ± 71	1690 ± 12	8322 ± 59	3763 ± 910	186341 ± 977	196667
≥ 2 b-tagged jets	76068 ± 231	1465 ± 24	248 ± 4	3096 ± 35	481 ± 413	81361 ± 476	85028

tainty shown in the plots refers to the theoretical uncertainty on the NNLO calculation of the total t̄t cross section at 8 TeV [67] used in this analysis.

Overall, good agreement is found between data and Monte Carlo simulation in all kinematic distributions. One should note that the comparison shown is made before the fitting process, which is described in Section 5.5.

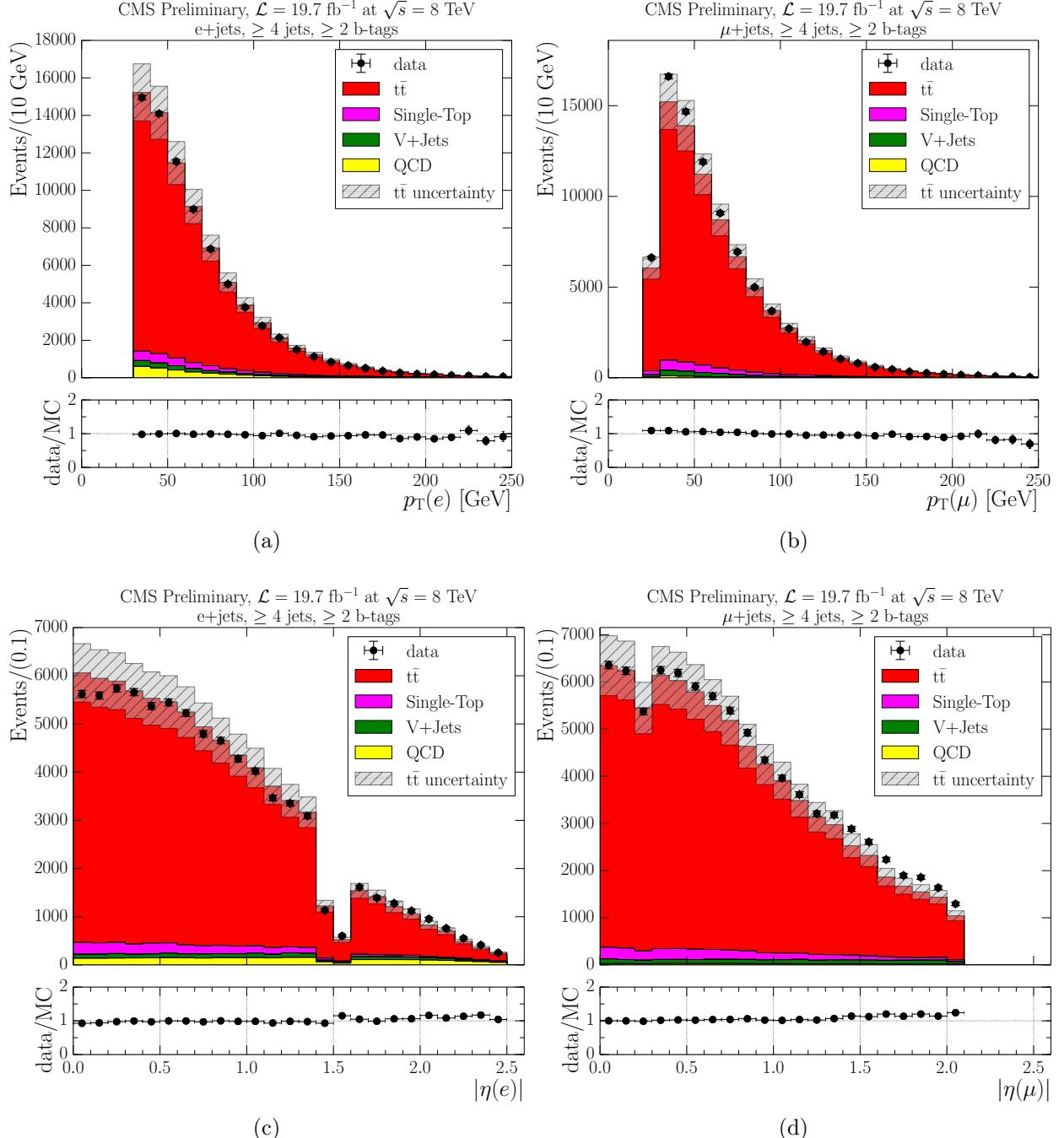


Figure 5.7: Data/MC comparison plots after the final event selection: electron  $p_T$  (a), muon  $p_T$  (b), electron  $|\eta|$  (c) and muon  $|\eta|$  (d). Left-hand plots: electron plus jets selection, right-hand plots: muon plus jets selection.

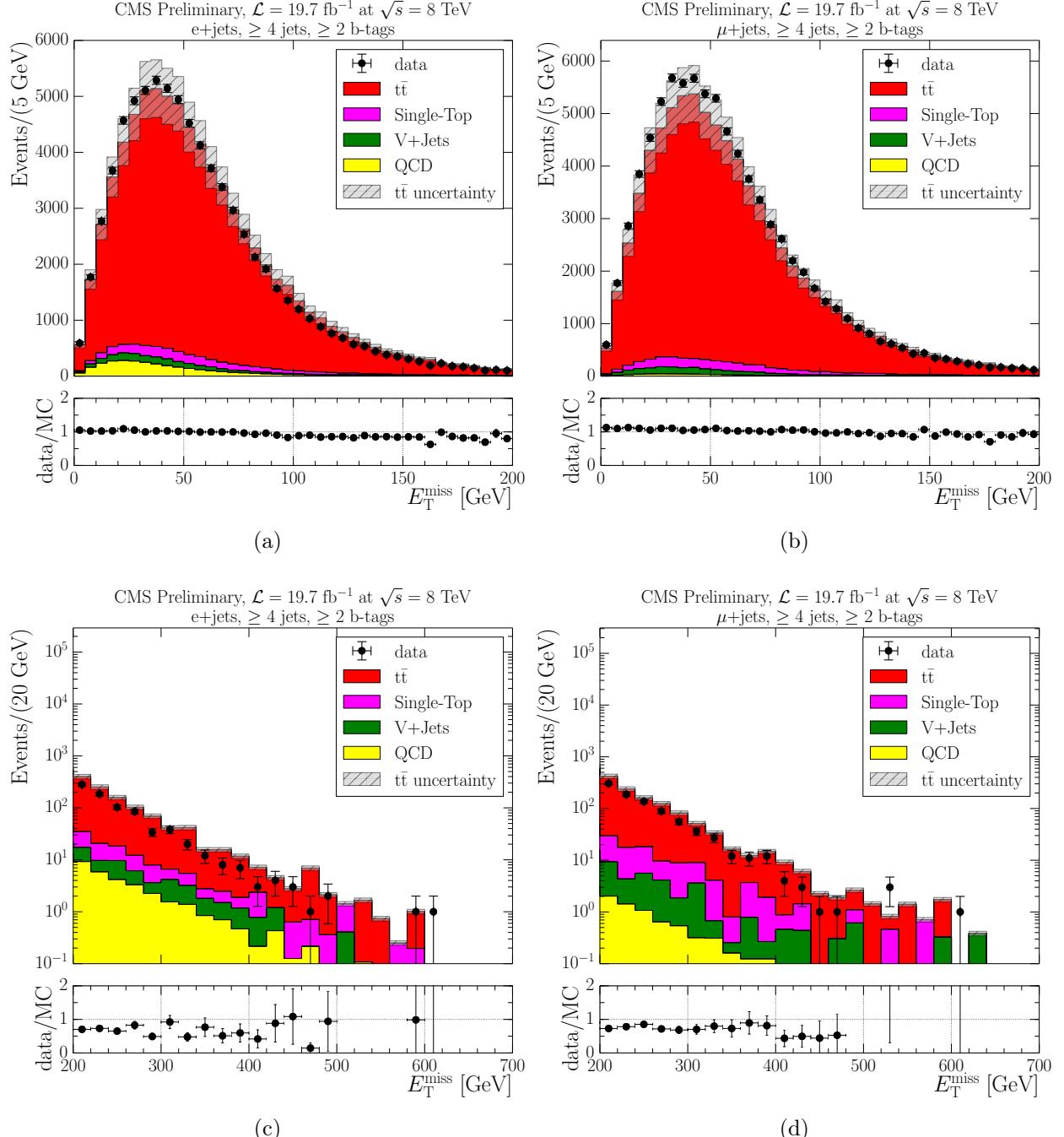


Figure 5.8: Data/MC comparison plots of  $E_T^{\text{miss}}$  (a, b) and logarithmic view of  $E_T^{\text{miss}}$  (c, d) after the final event selection. Left-hand plots: electron plus jets selection, right-hand plots: muon plus jets selection.

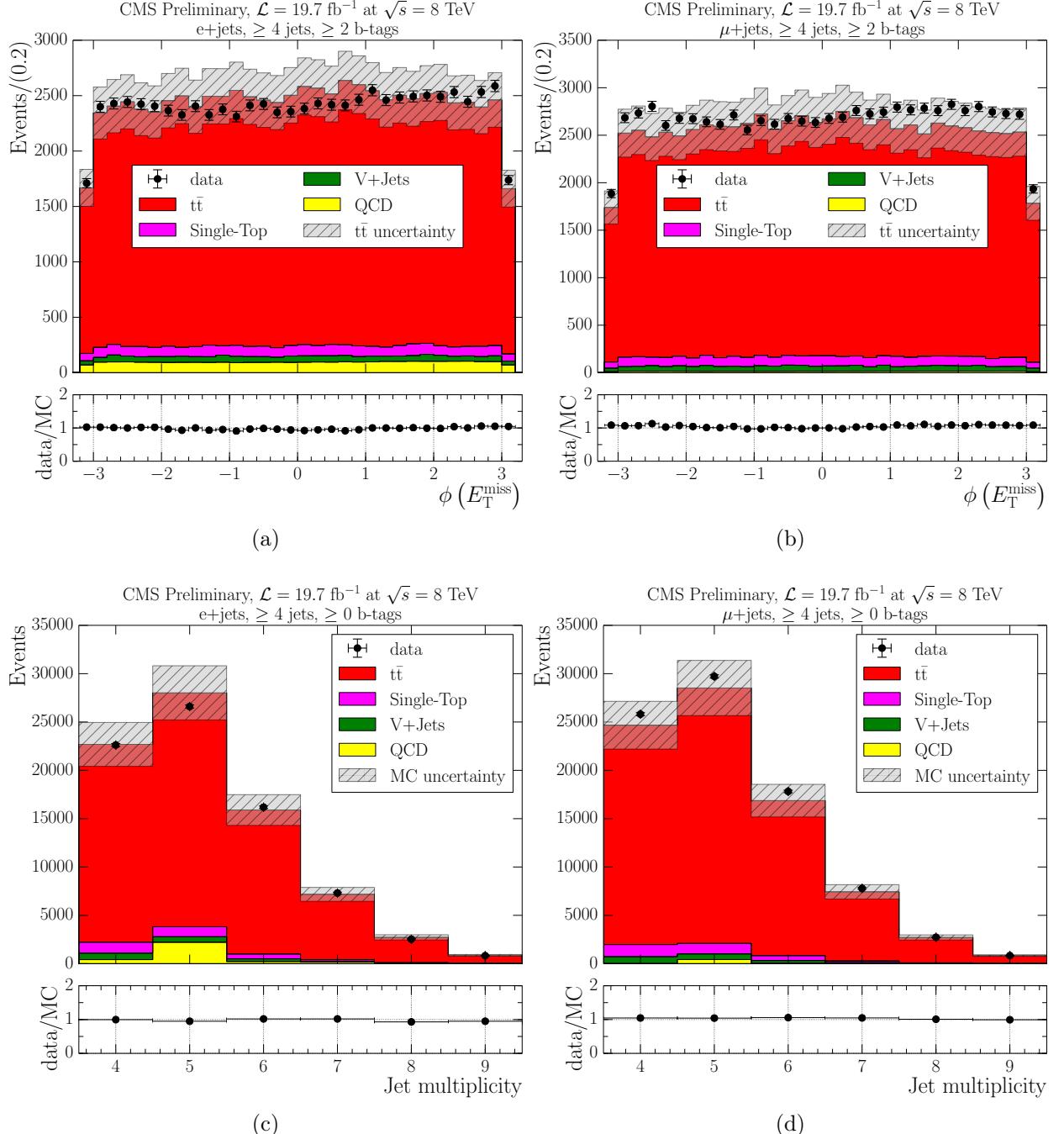


Figure 5.9: Data/MC comparison plots of  $\phi(E_T^{\text{miss}})$  (a, b) and jet multiplicity (c, d) after the final event selection. Left-hand plots: electron plus jets selection, right-hand plots: muon plus jets selection.

## 5.5 Differential cross section measurement

### 5.5.1 Primary variables

In this analysis, the normalised differential cross section of  $t\bar{t}$  production is studied with respect to five primary variables:

- $E_T^{\text{miss}}$ , or missing transverse momentum;
- $H_T$ , or scalar sum of jet transverse momenta;
- $S_T$ , or scalar sum of the transverse momenta of all objects in the event;
- $M_T^W$ , or transverse mass of the leptonically decaying W boson;
- $p_T^W$ , or transverse momentum of the leptonically decaying W boson.

These observables are often called event-level (or global) variables due to the fact that no high-level kinematic reconstruction (e.g. that of a  $t\bar{t}$  pair) is required to calculate them. A lower level of complexity makes these variables less prone to inefficiencies, systematic errors or biases arising from various kinematic reconstruction techniques. Hence, global variables can provide better sensitivity in comparison of different theoretical models or generator tunes.

$E_T^{\text{miss}}$  is commonly referred to as the missing transverse energy, which is equivalent to missing transverse momentum under assumption that the missing particles contributing to  $E_T^{\text{miss}}$  are massless. It is defined as the negative of the vector sum of the transverse momenta of all reconstructed particles in an event. Therefore, its direction is opposite to the observed sum of transverse momentum, with the magnitude of

$$E_T^{\text{miss}} = - \left[ \left( \sum_i p_x^i \right)^2 + \left( \sum_i p_y^i \right)^2 \right]^{\frac{1}{2}}$$

where the sums are calculated over all measured particles in the event.

$H_T$  is defined as the scalar sum of the transverse momenta of all reconstructed jets in the event:

$$H_T = \sum_{\text{all jets}} p_T^{\text{jet}}$$

Here all jets are required to have a  $p_T > 20$  GeV as described in Section 5.2.

$S_T$  represents the sum of  $H_T$ ,  $E_T^{\text{miss}}$  and the  $p_T$  of the single isolated lepton:

$$S_T = H_T + E_T^{\text{miss}} + p_T^{\text{lepton}}$$

$p_T^W$  refers to the transverse momentum of the leptonically decaying W boson, which is calculated from the single isolated lepton and  $E_T^{\text{miss}}$  assumed to originate from the  $t\bar{t}$  decay:

$$p_T^W = \sqrt{(p_x^{\text{lepton}} + p_x^{\text{miss}})^2 + (p_y^{\text{lepton}} + p_y^{\text{miss}})^2}$$

Finally,  $M_T^W$  is defined as the transverse mass of the leptonically decaying W boson, also using the single isolated lepton and  $E_T^{\text{miss}}$ :

$$M_T^W = \sqrt{(E_T^{\text{lepton}} + E_T^{\text{miss}})^2 - p_T^W{}^2}$$

where  $E_T^{\text{lepton}}$  denotes the transverse energy of the lepton.

### 5.5.2 Choice of binning

The choice of binning for primary variables is an important part of differential cross section measurement. It is directly influenced by available statistics and detector resolution. An ideal measurement would benefit from the cross section calculation in as many bins as possible. However, insufficient statistics leads to high statistical uncertainty in a particular bin, whereas limited detector resolution causes substantial migration effects between bins. Bin migration happens if an event originating from one bin appears in a different bin after reconstruction. Unreasonable binning choice would lead to a high number of such events, which can potentially compromise the cross section measurement.

To quantify the bin migration, purity and stability variables are introduced, defined as:

$$p_i = \frac{N_i^{\text{rec\&gen}}}{N_i^{\text{rec}}} \tag{5.1}$$

$$s_i = \frac{N_i^{\text{rec\&gen}}}{N_i^{\text{gen}}} \tag{5.2}$$

where  $N_i^{\text{rec\&gen}}$  is the number of events both generated and reconstructed in the same bin, and  $N_k^{\text{rec}}$  ( $N_i^{\text{gen}}$ ) is the number of reconstructed (generated) events within a bin.

The purity  $p_i$  quantifies migration into a bin  $i$ , whereas the stability  $s_i$  is sensitive to migration out of a bin. The binning is chosen in such way that both purity and stability are above 0.5, or do not fall far below this value. This process is performed for both electron

and muon channels, and common binning is derived in order to be able to combine both channels in the final cross section measurement.

Figure 5.10 shows the two-dimensional distributions of reconstructed versus generated  $E_T^{\text{miss}}$ , which were used to calculate purity and stability. The distributions were obtained with the central signal  $t\bar{t}$  sample. Red lines represent the bin edges for the chosen binning. Tables 5.7 and 5.8 show the numerical values of purity and stability, as well as the selected binning for  $E_T^{\text{miss}}$  variable. Results for other primary variables can be found in Appendix B.

Table 5.7: Stability and purity of chosen  $E_T^{\text{miss}}$  bins in the electron channel for  $t\bar{t}$  MC events.

bin, GeV	$0 < E_T^{\text{miss}} < 25$	$25 \leq E_T^{\text{miss}} < 45$	$45 \leq E_T^{\text{miss}} < 75$	$75 \leq E_T^{\text{miss}} < 100$	$100 \leq E_T^{\text{miss}} < 150$	$E_T^{\text{miss}} \geq 150$
events	10337	17982	19503	12594	7220	2832
purity	0.52	0.48	0.44	0.42	0.54	0.69
stability	0.42	0.41	0.47	0.52	0.66	0.86

Table 5.8: Stability and purity of chosen  $E_T^{\text{miss}}$  bins in the muon channel for  $t\bar{t}$  MC events.

bin, GeV	$0 < E_T^{\text{miss}} < 25$	$25 \leq E_T^{\text{miss}} < 45$	$45 \leq E_T^{\text{miss}} < 75$	$75 \leq E_T^{\text{miss}} < 100$	$100 \leq E_T^{\text{miss}} < 150$	$E_T^{\text{miss}} \geq 150$
events	11506	20348	22594	14413	8417	3274
purity	0.52	0.49	0.45	0.42	0.52	0.69
stability	0.43	0.42	0.47	0.5	0.66	0.87

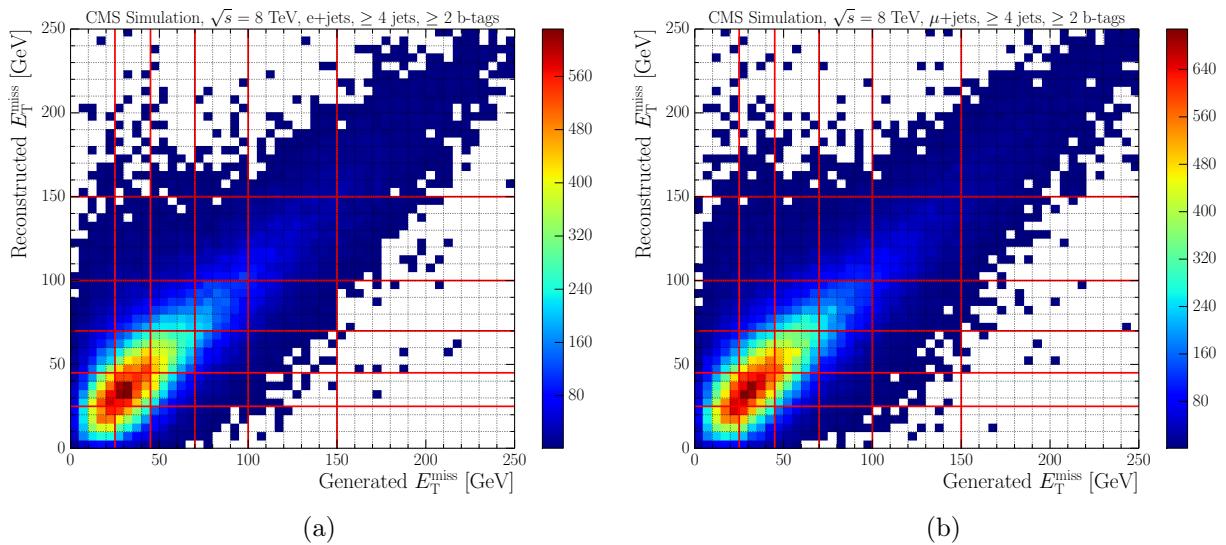


Figure 5.10: Reconstructed versus generated  $E_T^{\text{miss}}$  for electron plus jets (a) and muon plus jets events (b).

### 5.5.3 Fitting procedure

To calculate the differential cross section as a function of a primary variable, the number of signal  $t\bar{t}$  event has to be measured in each primary variable bin. For this purpose, a template fit of the observed distribution of lepton pseudorapidity ( $|\eta_l|$ ) is performed in each bin in order to estimate the number of  $t\bar{t}$  events by subtracting the primary backgrounds. Three separate normalised templates are produced for the input of the template fit:

- $|\eta_l|$  distribution of top-like events, i.e. combined template of  $t\bar{t}$  and single top MC;
- $|\eta_l|$  distribution of  $V + \text{jets}$  events, i.e. combined  $W + \text{jets}$  and  $Z + \text{jets}$  templates;
- $|\eta_l|$  shape of QCD events, estimated from data (Section 5.3).

The combination of  $W + \text{jets}$  and  $Z + \text{jets}$  templates is justified by the similarity of the shapes of distributions, as well as by limited MC statistics. Figures 5.11 and 5.12 show templates for the  $E_T^{\text{miss}}$  variable for electron and muon channels, respectively. Templates for other primary variables are attached in Appendix C.

The template fit is performed independently in each primary variable bin by minimising the log-likelihood defined as:

$$LL(\lambda_i, d_i) = -2 \log \left( \prod_i \frac{\lambda_i^{d_i} \cdot e^{-\lambda_i}}{d_i!} \right) = -2 \sum_i \log \left( \frac{\lambda_i^{d_i} \cdot e^{-\lambda_i}}{d_i!} \right), \quad (5.3)$$

where  $\lambda_i$  is the sum of expected events from all the templates  $d_i$  is the number of data events in each bin  $i$ , and summation is performed over all bins  $i$  of the  $|\eta_l|$  distribution. The expected lepton pseudorapidity (array of  $\lambda_i$ ) is modelled by normalising the templates  $\theta_{ij}$  by normalisation factors  $N_j$ :

$$\lambda_i = \sum_j N_j \theta_{ij}, \text{ where } \sum_i \theta_{ij} = 1 \text{ for every process } j. \quad (5.4)$$

The minimisation of the log-likelihood is achieved by varying the fit parameters, i.e. normalisation factors for each template. Therefore, the template fit strives for equality of  $\lambda_i$  and  $d_i$  in all lepton pseudorapidity bins. The initial normalisation values for the templates are taken from the expected number of events for  $t\bar{t}$ , single top,  $W + \text{jets}$  and  $Z + \text{jets}$  as per Table 5.1. Although the QCD template shape is data-driven, the initial normalisation for it is also taken from Monte Carlo simulation (Tables 5.2 and 5.3).

The similarity of QCD and  $V + \text{jets}$  templates, particularly for muon plus jets events (Figure 5.12), can lead to unphysical bin-to-bin fluctuations of these backgrounds. To mitigate

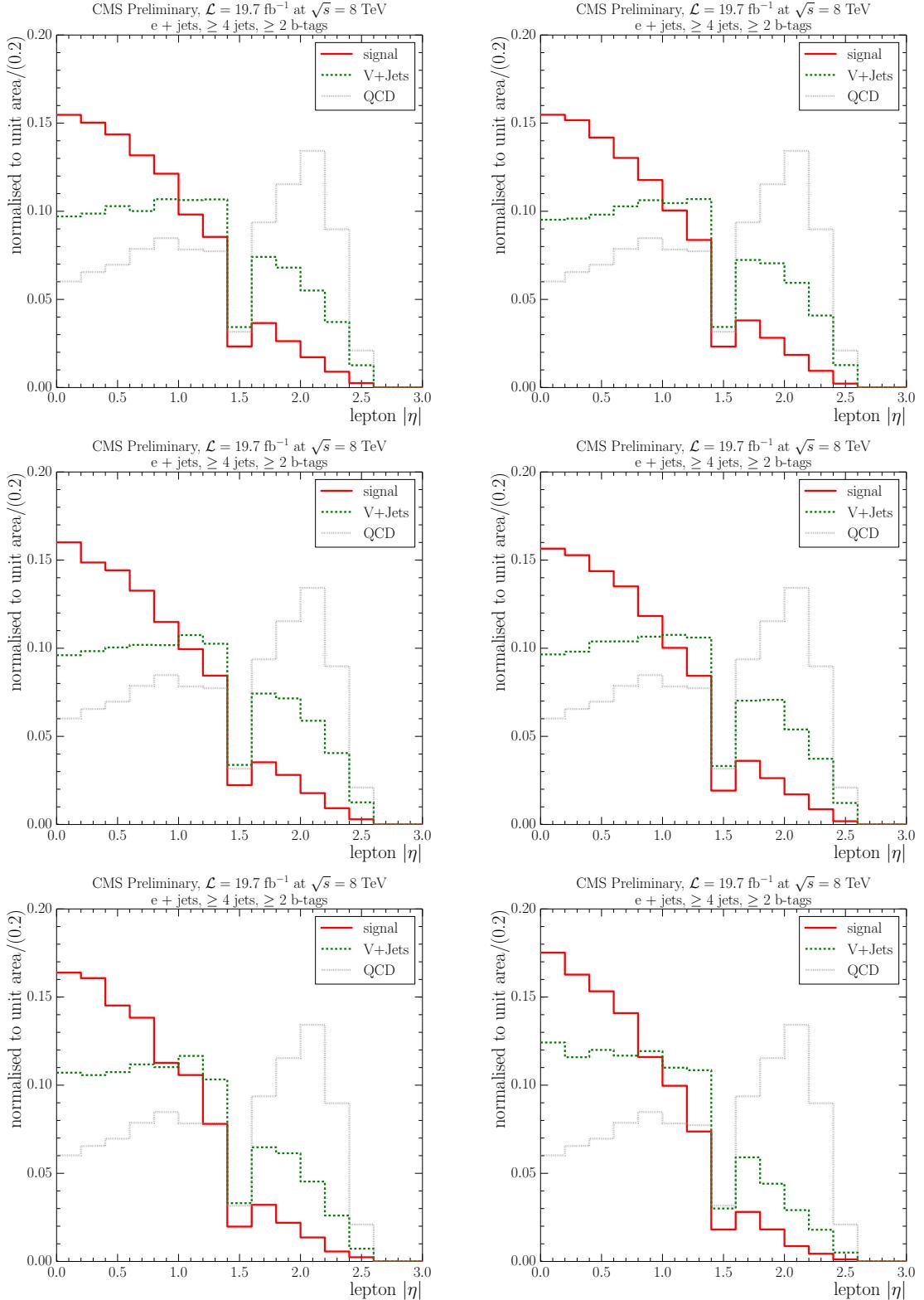


Figure 5.11: Electron  $|\eta|$  templates for the fit in different bins of  $E_T^{\text{miss}}$ , from top left to bottom right: 0–25 GeV, 25–45 GeV, 45–70 GeV, 70–100 GeV, 100–150 GeV and  $\geq 150$  GeV.

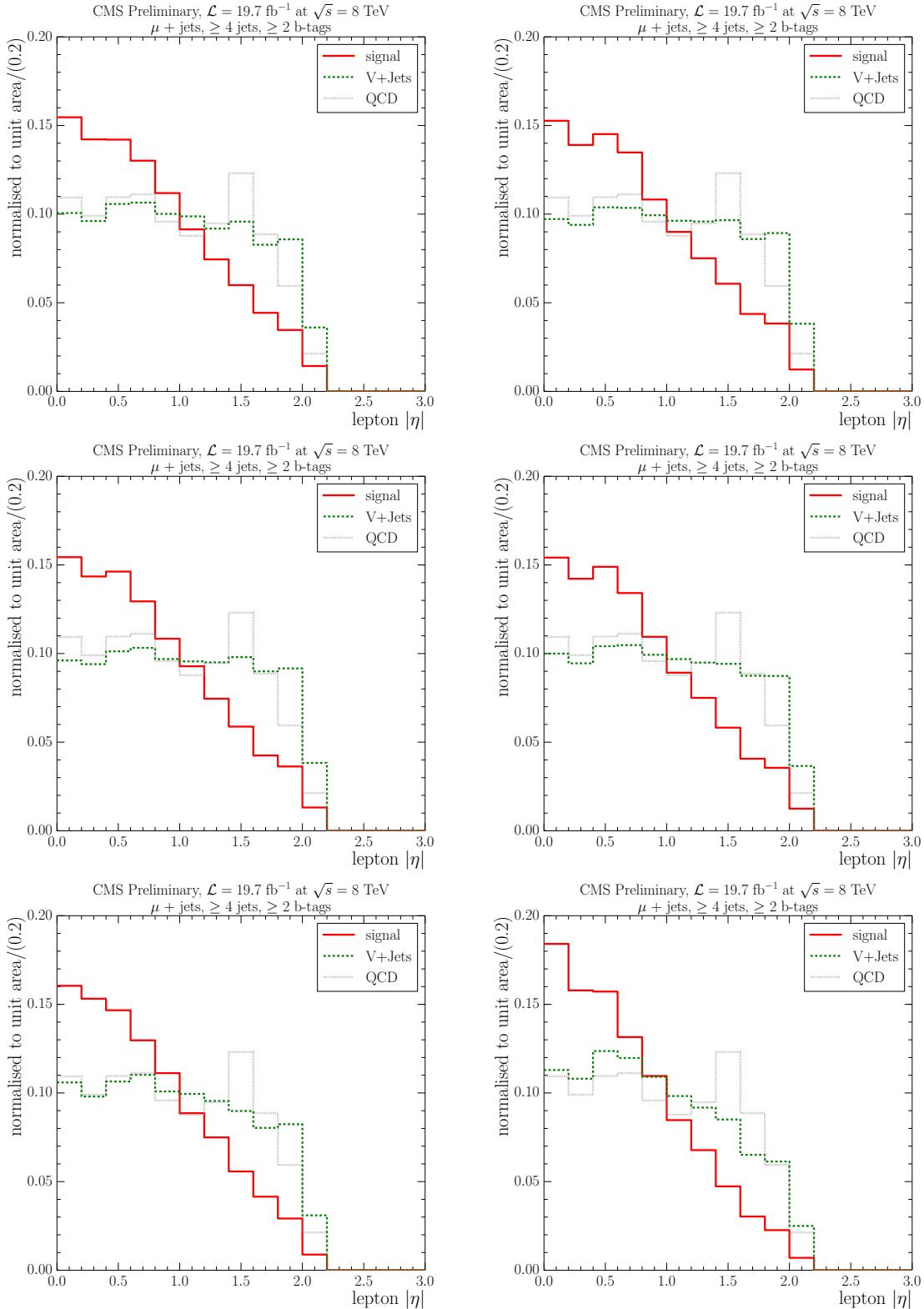


Figure 5.12: Muon  $|\eta|$  templates for the fit in different bins of  $E_T^{\text{miss}}$ , from top left to bottom right: 0–25 GeV, 25–45 GeV, 45–70 GeV, 70–100 GeV, 100–150 GeV and  $\geq 150$  GeV.

such behaviour, the following conservative Gaussian constraints were added to the likelihood function of the fit:

$$(N_{V+jets}^{\text{fit}} - N_{V+jets}^{\text{MC}})^2 / (0.5 \times N_{V+jets}^{\text{MC}})^2 \quad (5.5)$$

$$(N_{\text{QCD}}^{\text{fit}} - N_{\text{QCD}}^{\text{MC}})^2 / (2 \times N_{\text{QCD}}^{\text{MC}})^2 \quad (5.6)$$

The first term represents the constraint of the number of V+jets events to within 50 % of that predicted by theory, whereas the second term is the 200 % constraint for QCD events.

Figure 5.13 shows the control plots of data/MC comparison for  $E_{\text{T}}^{\text{miss}}$  variable in both electron and muon channels. Here the normalisation for signal and backgrounds is taken from the fitted results. Similar plots for  $H_{\text{T}}$  and  $S_{\text{T}}$  are shown in Figure 5.14, whereas Figure 5.15 shows control plots for  $p_{\text{T}}^{\text{W}}$  and  $M_{\text{T}}^{\text{W}}$ . Overall, the data and simulation agree within the fit uncertainty. The observed discrepancies can be attributed to the known issue of the event generators incorrectly modelling the  $p_{\text{T}}$  spectrum of the top quarks [68, 69]. This effect is accounted for as a source of systematic uncertainty, as explained in Section 5.6.

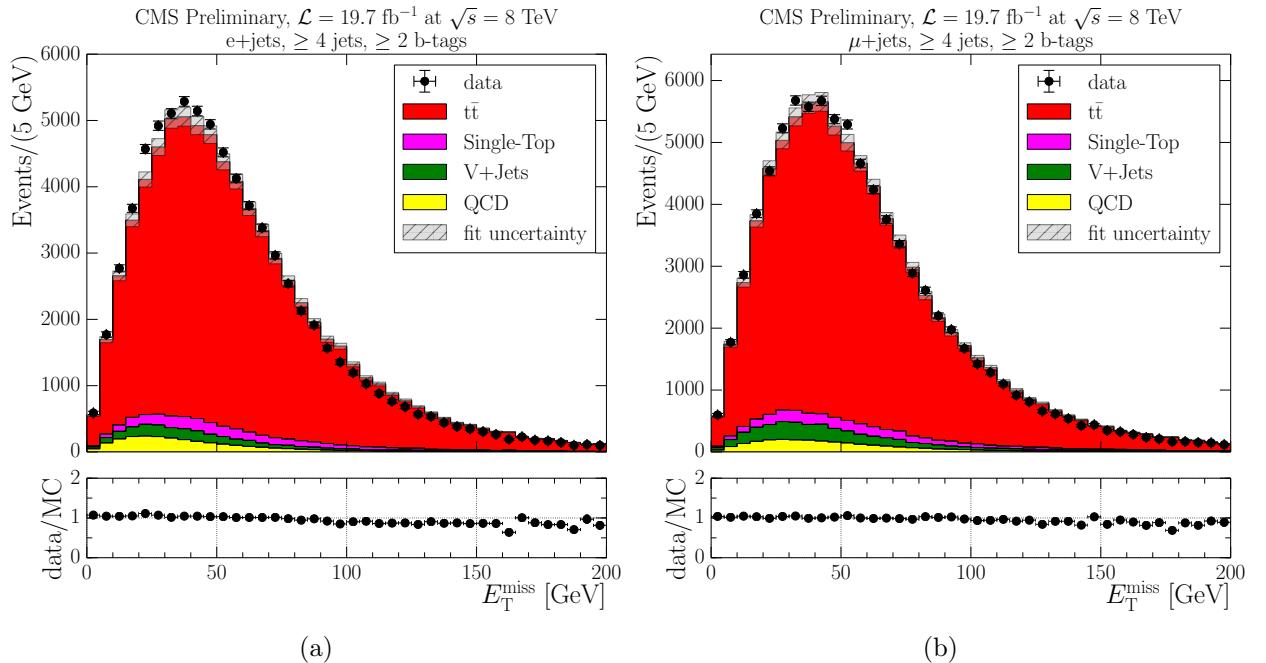


Figure 5.13: Data/MC comparison plots of  $E_{\text{T}}^{\text{miss}}$  for electron plus jets (a) and muon plus jets (b) events after the final event selection using the normalisation from the fitted results.

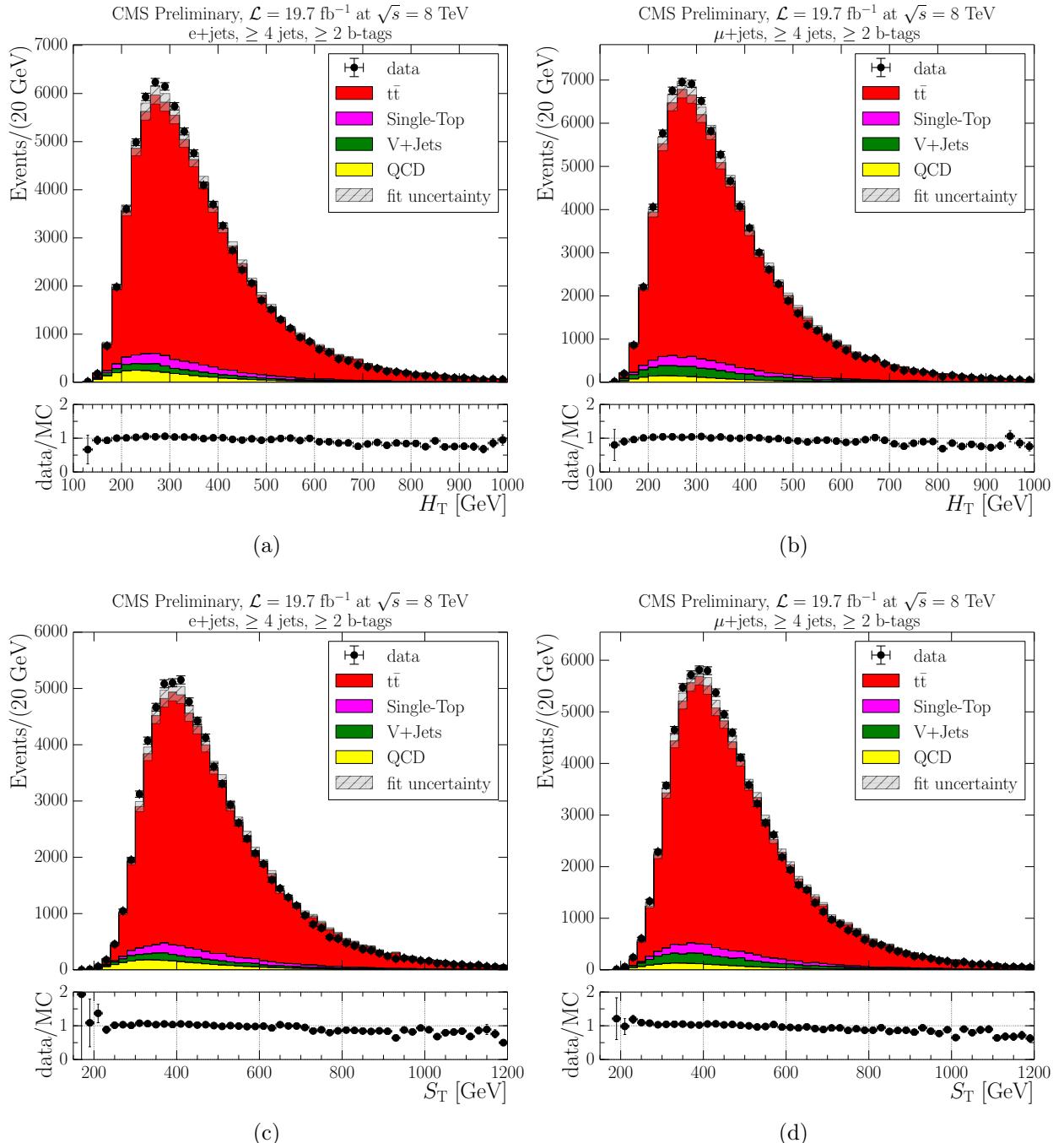


Figure 5.14: Data/MC comparison plots of  $H_T$  (a, b) and  $S_T$  (c, d) after the final event selection using the normalisation from the fitted results. Left-hand plots: electron plus jets selection, right-hand plots: muon plus jets selection.

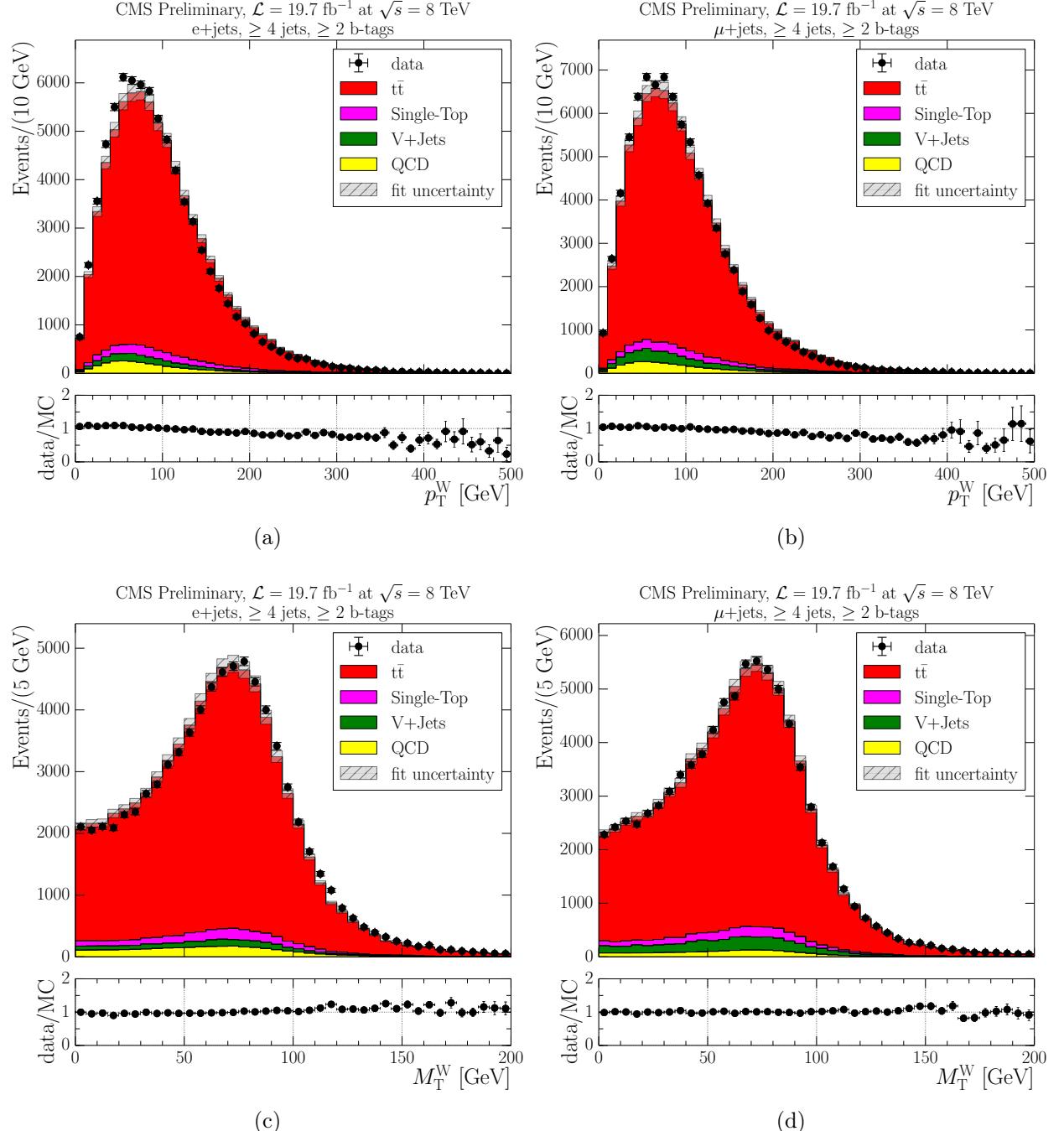


Figure 5.15: Data/MC comparison plots of  $p_T^W$  (a, b) and  $M_T^W$  (c, d) after the final event selection using the normalisation from the fitted results. Left-hand plots: electron plus jets selection, right-hand plots: muon plus jets selection.

The output of the fit yields the number of top-like events in each bin, which includes both  $t\bar{t}$  and single top processes since the lepton pseudorapidity distribution is very similar between them. The number of  $t\bar{t}$  events is extracted from the signal fit result by using the expected ratio of  $t\bar{t}$  and single top events before the fitting process:

$$N_{t\bar{t}}^{\text{fit}} = N_{\text{signal}}^{\text{fit}} \times \frac{N_{t\bar{t}}^{\text{exp}}}{N_{\text{signal}}^{\text{exp}}}, \quad (5.7)$$

where  $N_{t\bar{t}}^{\text{fit}}$  ( $N_{\text{signal}}^{\text{fit}}$ ) is the number of  $t\bar{t}$  ( $t\bar{t} + \text{single top}$ ) events as found by the fit, and  $N_{t\bar{t}}^{\text{exp}}$  ( $N_{\text{signal}}^{\text{exp}}$ ) is the number of  $t\bar{t}$  ( $t\bar{t} + \text{single top}$ ) events as predicted by theory.

#### 5.5.4 Unfolding

Measurements of physical observables are inevitably distorted by various detector effects, including limited detector resolution, experimental acceptance and selection inefficiency. Therefore, it is often very challenging to compare the results with theory predictions, as well as other experiments. To tackle this issue, unfolding can be used to estimate true underlying distributions of measured physical observables.

In order to describe the concept of unfolding, let us consider the distribution of a measured observable as a vector  $b$ . The measurement of this observable can be usually simulated by using Monte Carlo techniques according to some theoretical model and detector simulation. Let  $x_0$  be the “true” generated distribution, i.e. unaffected by detector effects, and  $b_0$  the corresponding measured distribution obtained by performing reconstruction in simulation:

$$\hat{A}x_0 = b_0. \quad (5.8)$$

Here  $\hat{A}$  is the response matrix of the detector, representing all aforementioned detector effects. Therefore, the given measured observable  $b$  and corresponding true distribution  $x$  are similarly related:

$$\hat{A}x = b. \quad (5.9)$$

Attempting to solve this system of equations by direct inversion of the matrix  $\hat{A}$  usually results in rapidly oscillating, unacceptable solutions. Therefore, a more sophisticated approach is needed. The process of finding the true distribution  $x$  is referred to as an unfolding, and in this particular analysis, Singular Value Decomposition (SVD) method of unfolding [70] is used.

The SVD unfolding incorporates the following factorisation of the response matrix:

$$\hat{A} = U S V^T, \quad (5.10)$$

where  $S$  is a diagonal matrix with non-negative elements (called singular values of matrix  $\hat{A}$ ), and  $U$  and  $V$  are orthogonal matrices (called the left and right singular vectors). The inverse of the response matrix then takes the form:

$$\hat{A}^{-1} = V S^{-1} U^T. \quad (5.11)$$

Such factorisation allows simple manipulation of the response matrix and its inverse, which is particularly useful for ill-determined cases of near-degenerate matrices when singular values of  $\hat{A}$  are close to zero. Additionally, if statistical errors of the measurement are substantial, it can result in amplification of essentially random components in the solution. These issues can be resolved by implementing regularisation in the unfolding procedure. The regularisation parameter  $k$  is used to select the number of statistically significant equations [70]. The value of  $k$  is determined by analysing the vector  $d$  obtained by rotating the measured distribution:

$$d = U^T \times b \quad (5.12)$$

For a reasonably smooth measured distribution (an a priori knowledge used for regularisation), only the first few terms of the decomposition are expected to be significant, with the rest of the terms corresponding to quickly oscillating vectors. This implies that the distribution  $d_i$  of the vector  $d$  components is expected to gradually fall towards a Gaussian-distributed random values for large  $i$ . The number of first few significant components ( $|d_i| \geq 1$ ) is therefore chosen as a regularisation parameter.

In this work, RooUnfold package [71] is used to implement the SVD method of data unfolding. The unfolded  $E_T^{\text{miss}}$  is defined as the generated  $p_T$  of the neutrino produced in the semileptonic  $t\bar{t}$  decay. The underlying  $H_T$  (and therefore  $S_T$ ) distribution is calculated by applying the full jet clustering reconstruction described in Section 2.4.3 to the generated particles, with the  $p_T > 20$  GeV requirement for the generated jets.  $M_T^W$  and  $p_T^W$  variables are constructed from generated  $E_T^{\text{miss}}$  and leptons, as shown in Section 5.5.1. The response matrices are built from true and measured distributions of primary variables, based on  $t\bar{t}$  Monte Carlo samples.

To test and validate the performance of the unfolding, a set of 300 pseudo-experiments (also known as toy MC) was created by generating a random number of events in each bin of a

primary variable according to a Poisson distribution with the mean around the initial central value. This was performed for both true and measured distributions in the unfolding, with the response matrix recalculated accordingly. Then the unfolding procedure is performed for the measured distribution in each pseudo-experiment by using the SVD decomposition of a response matrix from every other pseudo-experiment, making a total of  $300 \times 299 = 89700$  different measurements. For each combination, a pull is calculated:

$$p_i = \frac{N_i^{\text{unfolded}} - N_i^{\text{truth}}}{\sigma_i}, \quad (5.13)$$

where  $\sigma_i$  is the combined statistical and unfolding uncertainty in bin  $i$ . The pull distributions for all bins of the  $E_T^{\text{miss}}$  variable are shown in Figure 5.16. Clearly, there is no significant bias in the unfolding procedure, as all the distributions are centred around zero. The width of the distributions is close to unity, meaning that the errors are reliably calculated. The unfolding performance was similarly validated for other primary variables.

### 5.5.5 Differential cross section calculation

The unfolded number of  $t\bar{t}$  events  $N_{t\bar{t}}^k$  can be converted to a partial cross section of  $t\bar{t}$  production in each bin  $k$  of a given variable:

$$\Delta\sigma_{t\bar{t}}^k = \frac{N_{t\bar{t}}^k}{\text{BR} \times \mathcal{L}}, \quad (5.14)$$

where  $\mathcal{L}$  is the total integrated luminosity and BR is the theoretical branching ratio of semileptonic  $t\bar{t}$  decay. These values cancel in the normalisation of the cross section, which is performed by division by the bin width  $\Delta X^k$  of a given variable  $X$ , yielding the differential cross section:

$$\frac{d\sigma_{t\bar{t}}^k}{dX} = \frac{\Delta\sigma_{t\bar{t}}^k}{\Delta X^k}. \quad (5.15)$$

Finally, the normalised differential cross section is obtained by using the sum of the partial differential cross sections:

$$\frac{1}{\sigma_{t\bar{t}}^{\text{tot}}} \frac{d\sigma_{t\bar{t}}^k}{dX} = \frac{1}{\sum_k \frac{d\sigma_{t\bar{t}}^k}{dX}} \frac{d\sigma_{t\bar{t}}^k}{dX} = \frac{1}{\Delta X^k} \frac{N_{t\bar{t}}^k}{\sum_k N_{t\bar{t}}^k}. \quad (5.16)$$

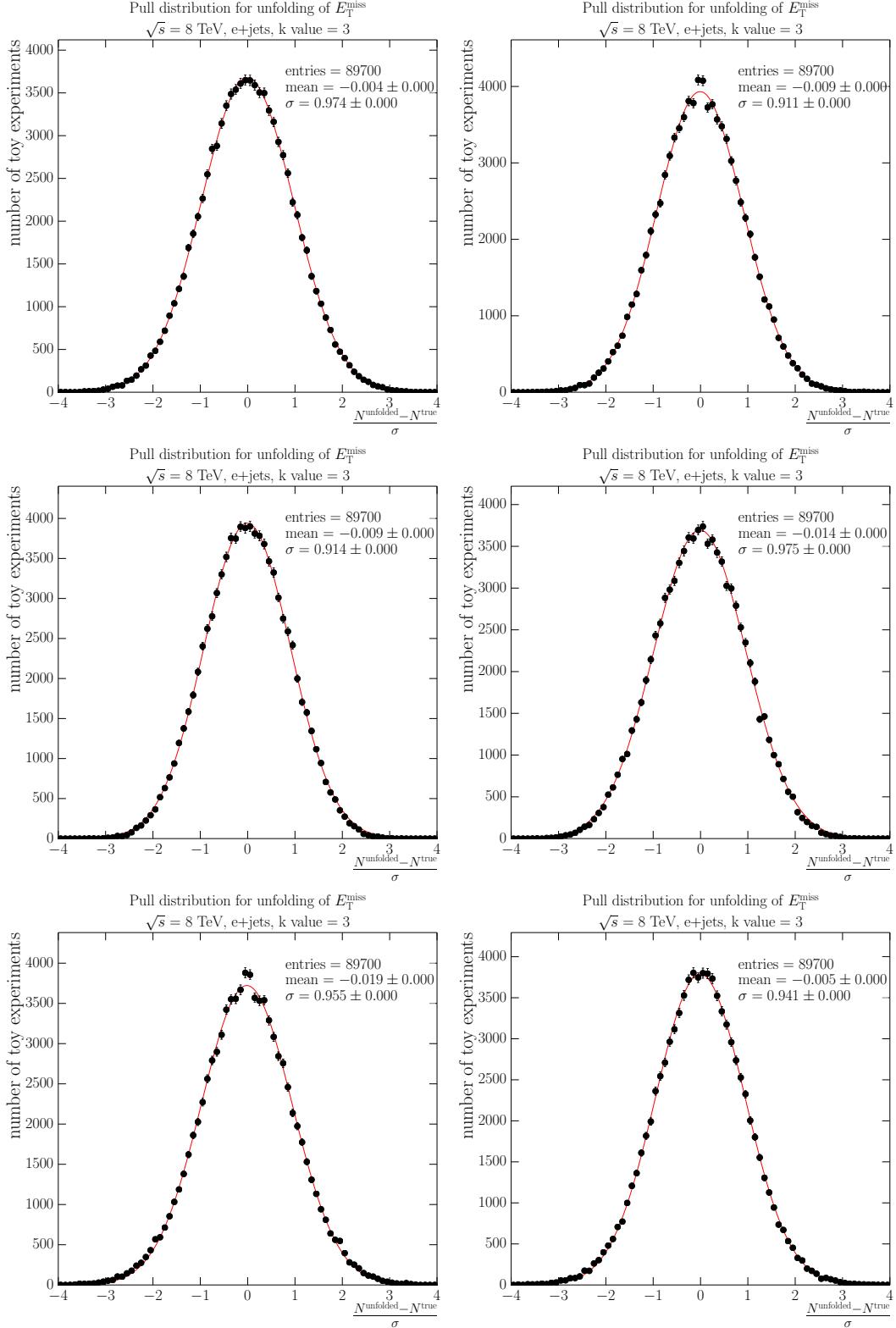


Figure 5.16: Pull distributions for the SVD unfolding method for all  $E_T^{\text{miss}}$  bins, from top left to bottom right: 0 GeV to 25 GeV, 25 GeV to 45 GeV, 45 GeV to 70 GeV, 70 GeV to 100 GeV, 100 GeV to 150 GeV and  $\geq 150$  GeV.

## 5.6 Systematic Uncertainties

The sources of systematic uncertainties considered in this analysis are largely the same to those of the top quark mass analysis (Section 4.6). Similarly, the errors are estimated by varying the input quantities according to their theoretical or experimental variations, and measuring the change in the final result. All systematic errors are added to the fitting and unfolding errors in quadrature.

The full list of systematic uncertainties for the normalised cross section measurement with respect to  $E_T^{\text{miss}}$  variable is shown in Table 5.9.

Table 5.9: Systematic uncertainties for the normalised  $t\bar{t}$  cross section measurement with respect to  $E_T^{\text{miss}}$  variable (combination of electron and muon channels).

Uncertainty source	0–25 GeV	25–45 GeV	45–70 GeV	70–100 GeV	100–150 GeV	$\geq 150$ GeV
b-jets - (%)	0.13	0.06	-0.05	-0.09	-0.07	-0.03
b-jets + (%)	0.01	0.02	0.01	-0.02	-0.04	-0.05
JER - (%)	0.08	0.04	-0.04	-0.06	-0.04	0.01
JER + (%)	0.05	0.04	-0.01	-0.06	-0.07	-0.05
JES - (%)	1.75	1.11	-0.02	-1.36	-2.57	-3.47
JES + (%)	-1.72	-1.02	0.11	1.29	2.26	3.01
Light jet - (%)	0.09	0.04	-0.03	-0.06	-0.04	-0.01
Light jet + (%)	0.04	0.03	-0.00	-0.04	-0.06	-0.06
Pile-up - (%)	0.07	0.04	-0.02	-0.06	-0.07	-0.08
Pile-up + (%)	0.09	0.03	-0.03	-0.06	-0.03	0.02
QCD shape uncertainty (%)	-0.85	-0.34	0.24	0.53	0.57	0.53
hadronisation uncertainty (%)	-1.15	0.33	0.24	-0.26	0.48	0.85
$p_T(t, \bar{t})$ reweighting (%)	0.44	0.26	0.00	-0.27	-0.66	-1.12
$t\bar{t}$ (matching down) (%)	-0.57	0.01	-0.03	0.16	0.79	-0.18
$t\bar{t}$ (matching up) (%)	-0.41	-0.48	-0.56	0.31	2.65	2.06
$t\bar{t}$ ( $Q^2$ down) (%)	1.14	0.58	-0.90	-0.80	0.57	0.18
$t\bar{t}$ ( $Q^2$ up) (%)	-1.20	-0.48	-0.27	1.06	1.61	2.24
V+jets (matching down) (%)	0.30	0.18	-0.17	-0.44	0.03	0.86
V+jets (matching up) (%)	0.80	0.65	0.08	-0.88	-1.45	-1.45
V+jets ( $Q^2$ down) (%)	0.12	-0.19	-0.35	-0.03	0.85	2.02
V+jets ( $Q^2$ up) (%)	1.89	0.82	-0.56	-1.37	-1.25	-0.32
Electron energy $-1\sigma$ (%)	0.09	0.03	-0.08	-0.08	0.09	0.25
Electron energy $+1\sigma$ (%)	-0.12	0.02	0.12	0.02	-0.17	-0.33
Muon energy $-1\sigma$ (%)	0.01	-0.02	-0.04	0.01	0.09	0.18
Muon energy $+1\sigma$ (%)	-0.02	0.02	0.04	-0.01	-0.07	-0.12
Tau energy $-1\sigma$ (%)	-0.02	0.00	0.01	0.00	-0.01	-0.01
Tau energy $+1\sigma$ (%)	0.04	-0.00	-0.03	-0.01	0.03	0.05
Unclustered energy $-1\sigma$ (%)	1.00	0.67	-0.00	-0.85	-1.50	-1.86
Unclustered energy $+1\sigma$ (%)	-1.10	-0.70	0.05	0.88	1.54	1.96
Total (%)	3.81	2.20	1.80	2.96	5.08	5.82

## 5.7 Results

The template fit explained in Section 5.5.3 was performed for both electron and muon channels, yielding two distributions of signal  $t\bar{t}$  events. These distributions are then separately unfolded using the SVD unfolding method described in Section 5.5.4. The two unfolded distributions are summed together to obtain the total distribution of  $t\bar{t}$  events in semileptonic production mode. Finally, the normalised cross section is calculated as outlined in Section 5.5.5.

The final combined results are shown in Figures 5.17–5.21 for all primary variables. The unfolded data are compared with predictions by MADGRAPH, POWHEG and MC@NLO event generators (Section 4.1.2.1). Results are also compared with theoretical predictions showing the effects due to variations of modelling parameters, including matching threshold and normalisation scale (Section 4.1.2.2). In all figures, the inner error bars represent the combined fitting and unfolding uncertainties, whereas the outer error bars show the uncertainty due to systematic variations. Naturally, the systematic errors due to matching threshold and normalisation scale choice are excluded from comparison with variations due to these parameters, and the hadronisation systematic is excluded from the comparison with different generators. Overall, the results are consistent with predictions of the Standard Model.

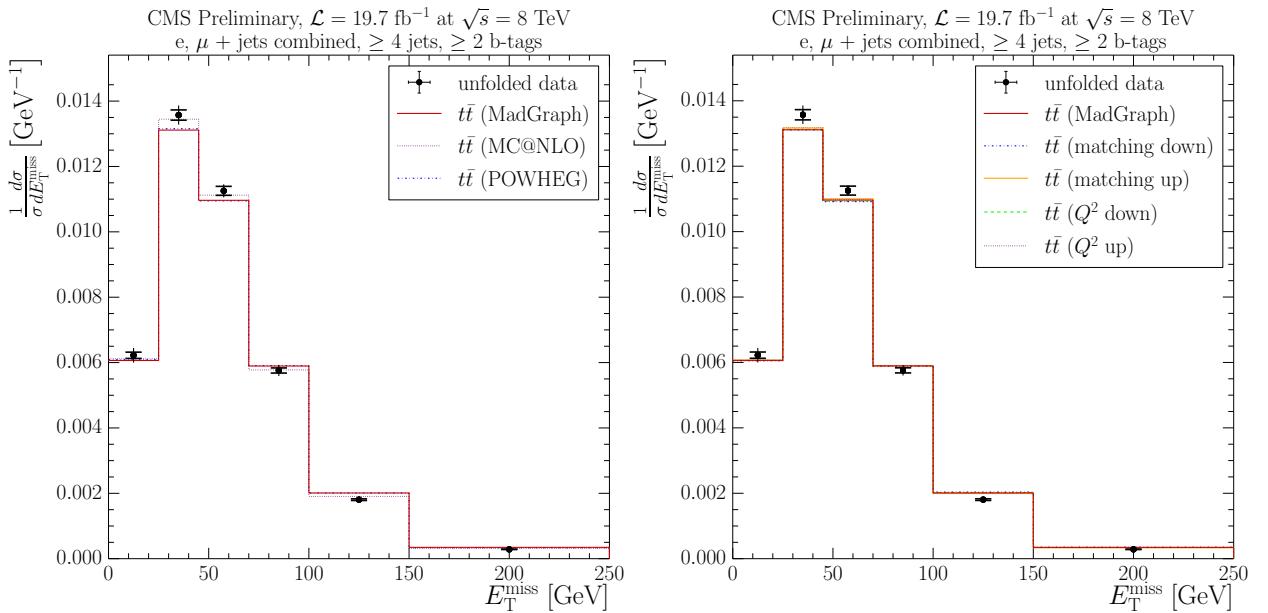


Figure 5.17: The normalised semileptonic  $t\bar{t}$  differential cross section with respect to  $E_T^{\text{miss}}$ . The data are compared to predictions of three different MC generators (left) and theoretical predictions showing the variation due to modelling uncertainties (right).

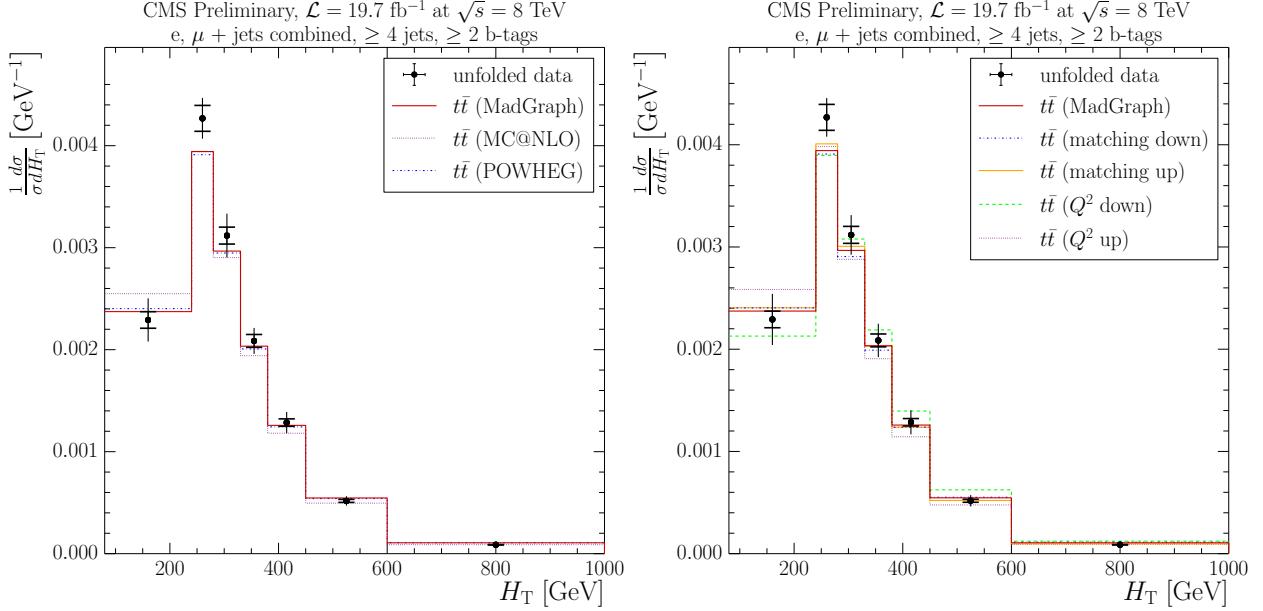


Figure 5.18: The normalised semileptonic  $t\bar{t}$  differential cross section with respect to  $H_T$ . The data are compared to predictions of three different MC generators (left) and theoretical predictions showing the variation due to modelling uncertainties (right).

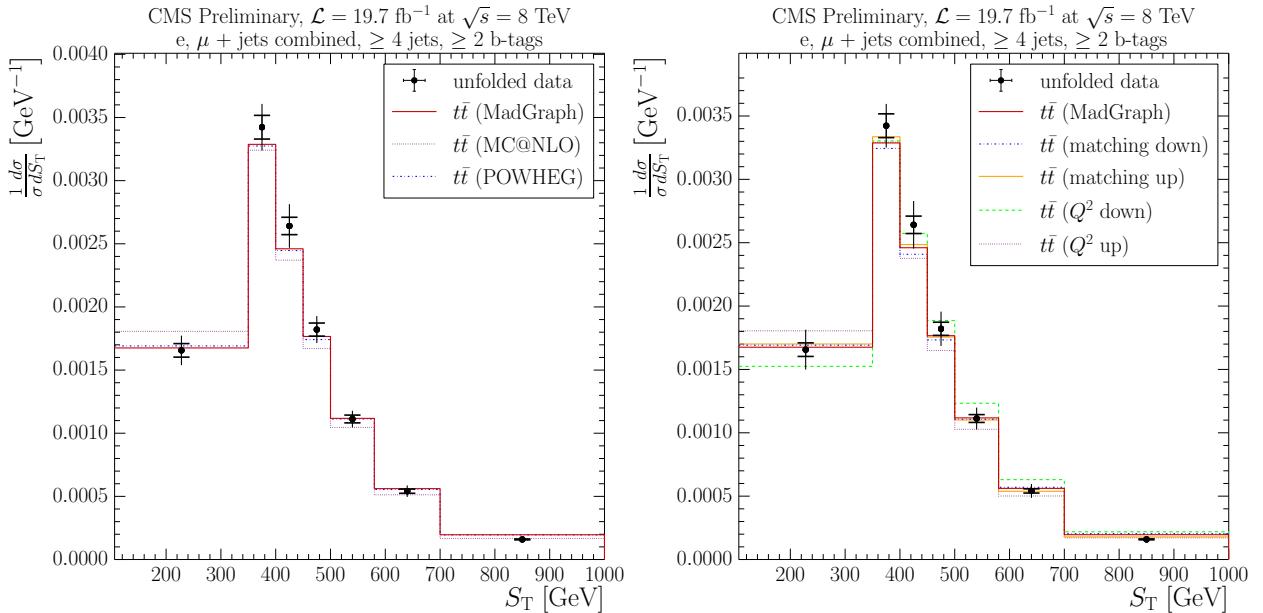


Figure 5.19: The normalised semileptonic  $t\bar{t}$  differential cross section with respect to  $S_T$ . The data are compared to predictions of three different MC generators (left) and theoretical predictions showing the variation due to modelling uncertainties (right).

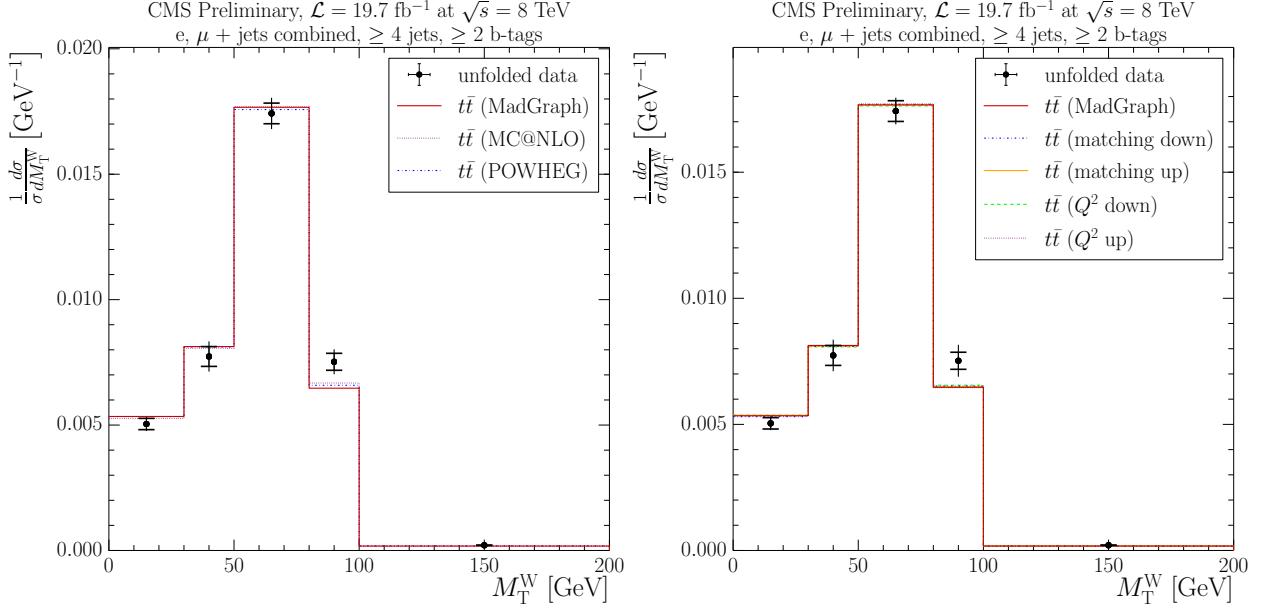


Figure 5.20: The normalised semileptonic  $t\bar{t}$  differential cross section with respect to  $M_T^W$ . The data are compared to predictions of three different MC generators (left) and theoretical predictions showing the variation due to modelling uncertainties (right).

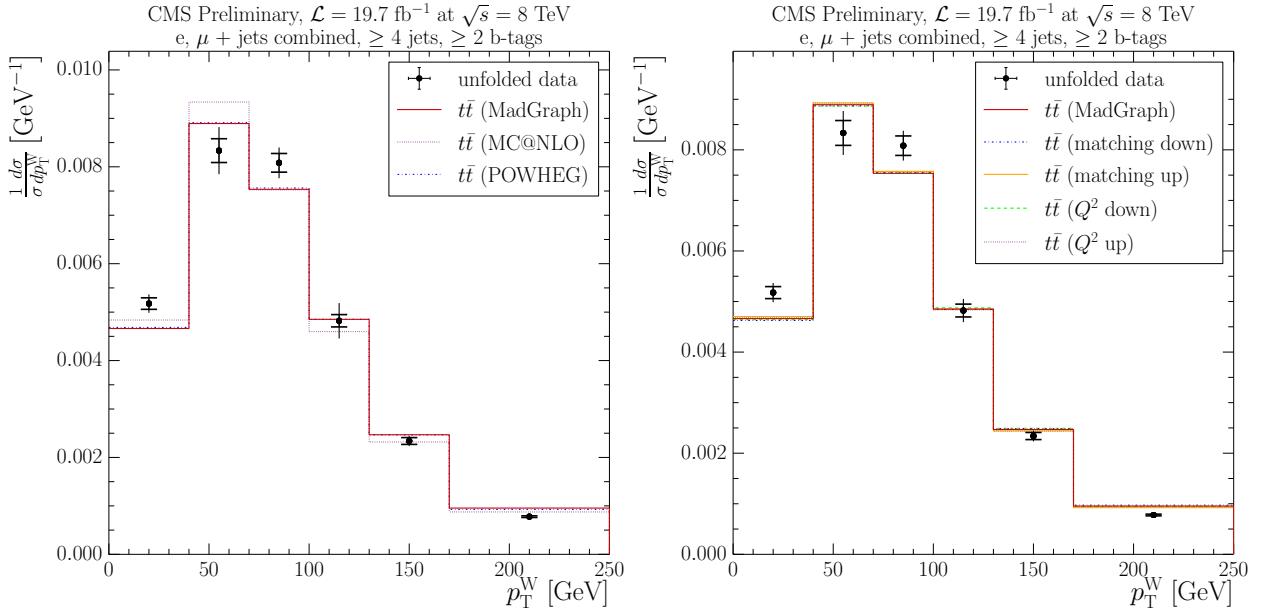


Figure 5.21: The normalised semileptonic  $t\bar{t}$  differential cross section with respect to  $p_T^W$ . The data are compared to predictions of three different MC generators (left) and theoretical predictions showing the variation due to modelling uncertainties (right).

## 5.8 Summary

A differential cross section measurement of top quark pair production with respect to  $E_T^{\text{miss}}$  and other global variables was performed in proton-proton collisions at a centre of mass energy of  $\sqrt{s} = 8$  TeV, using the full 2012 LHC dataset corresponding to integrated luminosity of  $19.7 \text{ fb}^{-1}$ , collected by the CMS experiment. The analysis selected events with a single isolated highly-energetic electron or muon, which is assumed to come from one of the W bosons in top quark and anti-quark decay. The shape of QCD multi-jet events was estimated using data-driven techniques. A template fit method was used to determine the sample composition and the number of  $t\bar{t}$  events, which was corrected for misidentification, detector resolution and acceptance effects using the SVD unfolding technique. The results from both semileptonic decay channels were combined and compared with different Monte Carlo generators, as well as with variations in theoretical predictions due to modelling uncertainties. No significant deviations from the Standard Model predictions were observed.



## A. High level triggers

A full list of electron triggers including their version numbers used in the top quark mass analysis analysis (Chapter 4):

- HLT-Ele25-CaloidVT-TrkIdT-CentralTriJet30-v1  
for run number  $\leq 161216$
- HLT-Ele25-CaloidVT-TrkIdT-CentralTriJet30-v2  
for run number  $\leq 163269$
- HLT-Ele25-CaloidVT-TrkIdT-TriCentralJet30-v3  
for run number  $\leq 165969$
- HLT-Ele25-CaloidVT-CaloiSoT-TrkIdT-TrkIsoT-TriCentralJet30-v1  
for run number  $\leq 166967$
- HLT-Ele25-CaloidVT-CaloiSoT-TrkIdT-TrkIsoT-TriCentralJet30-v2  
for run number  $\leq 167913$
- HLT-Ele25-CaloidVT-CaloiSoT-TrkIdT-TrkIsoT-TriCentralJet30-v4  
for run number  $\leq 173235$
- HLT-Ele25-CaloidVT-CaloiSoT-TrkIdT-TrkIsoT-TriCentralJet30-v5  
for run number  $\leq 178380$
- HLT-Ele25-CaloidVT-CaloiSoT-TrkIdT-TrkIsoT-TriCentralPFJet30-v2  
for run number  $\leq 179889$
- HLT-Ele25-CaloidVT-CaloiSoT-TrkIdT-TrkIsoT-TriCentralPFJet30-v3  
for run number  $\leq 180252$



## B. Choice of binning

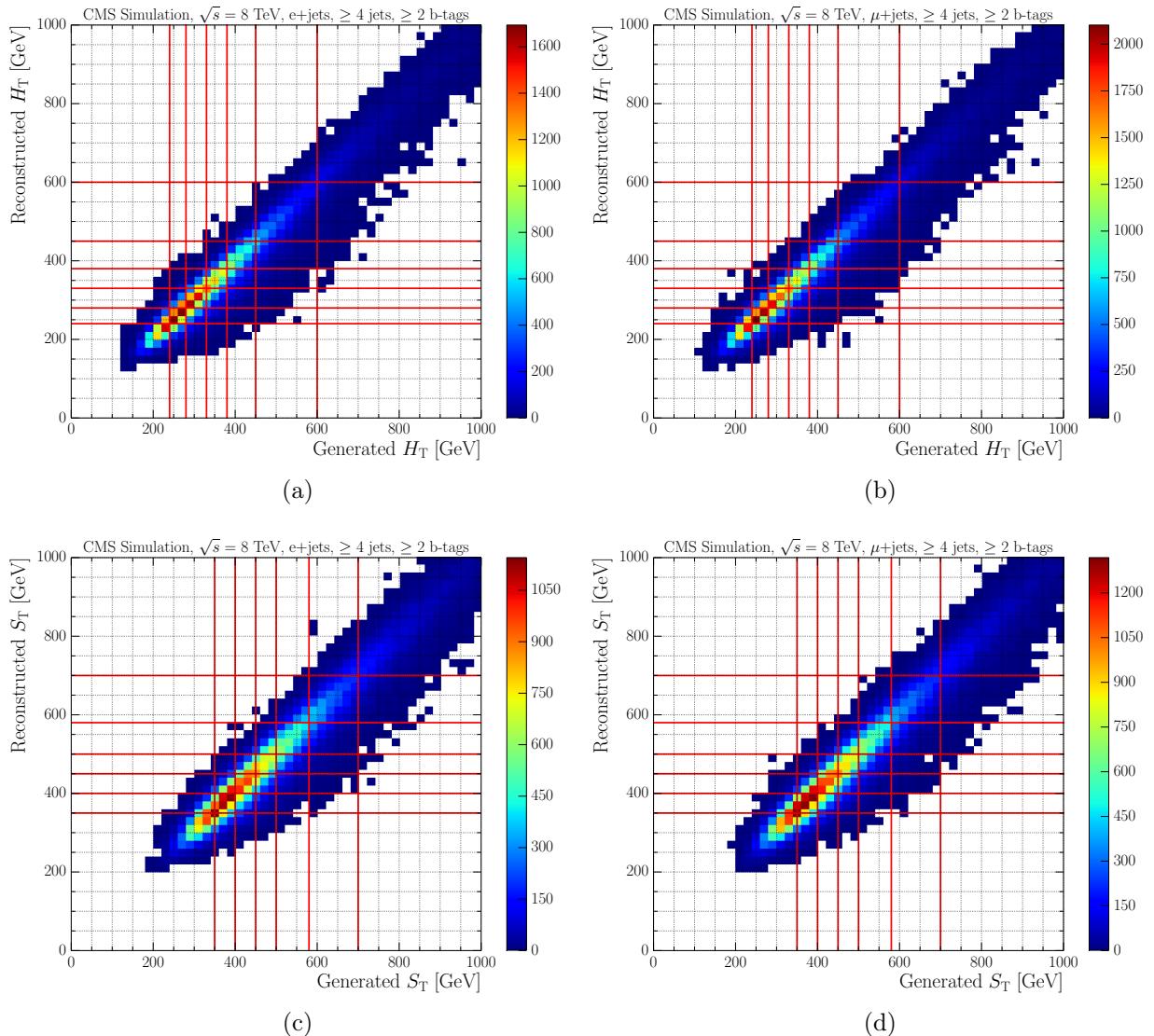


Figure B.1: Reconstructed versus generated  $H_T$  (a, b) and  $S_T$  (c, d) for electron plus jets (left) and muon plus jets events (right).

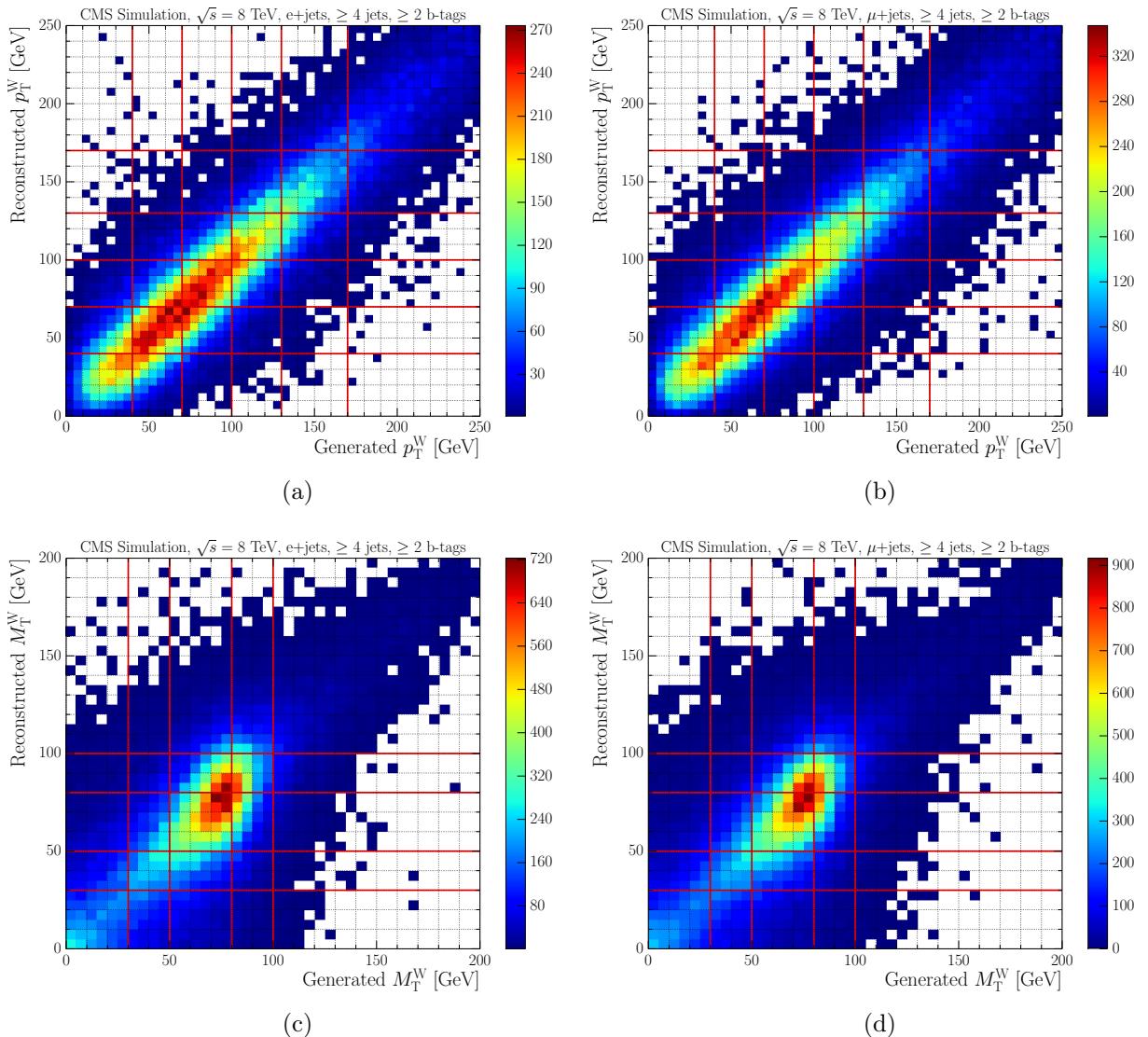


Figure B.2: Reconstructed versus generated  $p_T^W$  (a, b) and  $M_T^W$  (c, d) for electron plus jets (left) and muon plus jets events (right).

Table B.1: Stability and purity of chosen  $H_T$  bins in the electron channel for  $t\bar{t}$  MC events.

bin, GeV	$80 < H_T < 240$	$240 \leq H_T < 280$	$280 \leq H_T < 330$	$330 \leq H_T < 380$	$380 \leq H_T < 450$	$450 \leq H_T < 600$	$H_T \geq 600$
events	10601	12569	14372	12397	10893	10721	5609
purity	0.82	0.62	0.63	0.62	0.68	0.81	0.9
stability	0.76	0.63	0.65	0.64	0.68	0.8	0.9

Table B.2: Stability and purity of chosen  $H_T$  bins in the muon channel for  $t\bar{t}$  MC events.

bin, GeV	$80 < H_T < 240$	$240 \leq H_T < 280$	$280 \leq H_T < 330$	$330 \leq H_T < 380$	$380 \leq H_T < 450$	$450 \leq H_T < 600$	$H_T \geq 600$
events	12337	14821	16747	14268	12258	12119	5857
purity	0.82	0.63	0.64	0.63	0.67	0.81	0.9
stability	0.76	0.64	0.65	0.65	0.68	0.81	0.89

Table B.3: Stability and purity of chosen  $S_T$  bins in the electron channel for  $t\bar{t}$  MC events.

bin, GeV	$106 < S_T < 350$	$350 \leq S_T < 400$	$400 \leq S_T < 450$	$450 \leq S_T < 500$	$500 \leq S_T < 580$	$580 \leq S_T < 700$	$S_T \geq 700$
events	11183	13246	12073	10709	11317	9166	8373
purity	0.83	0.6	0.53	0.53	0.63	0.71	0.87
stability	0.73	0.6	0.55	0.55	0.66	0.73	0.91

Table B.4: Stability and purity of chosen  $S_T$  bins in the muon channel for  $t\bar{t}$  MC events.

bin, GeV	$106 < S_T < 350$	$350 \leq S_T < 400$	$400 \leq S_T < 450$	$450 \leq S_T < 500$	$500 \leq S_T < 580$	$580 \leq S_T < 700$	$S_T \geq 700$
events	13593	15717	13878	12307	12732	10135	8875
purity	0.83	0.61	0.54	0.54	0.64	0.71	0.87
stability	0.74	0.61	0.55	0.56	0.66	0.74	0.9

Table B.5: Stability and purity of chosen  $p_T^W$  bins in the electron channel for  $t\bar{t}$  MC events.

bin, GeV	$0 < p_T^W < 40$	$40 \leq p_T^W < 70$	$70 \leq p_T^W < 100$	$100 \leq p_T^W < 130$	$130 \leq p_T^W < 170$	$p_T^W \geq 170$
events	10005	15849	16053	12203	9532	7717
purity	0.64	0.54	0.52	0.5	0.56	0.77
stability	0.63	0.54	0.52	0.51	0.57	0.76

Table B.6: Stability and purity of chosen  $p_T^W$  bins in the muon channel for  $t\bar{t}$  MC events.

bin, GeV	$0 < p_T^W < 40$	$40 \leq p_T^W < 70$	$70 \leq p_T^W < 100$	$100 \leq p_T^W < 130$	$130 \leq p_T^W < 170$	$p_T^W \geq 170$
events	12379	18821	18330	13495	10260	8394
purity	0.67	0.55	0.52	0.5	0.55	0.76
stability	0.63	0.55	0.53	0.51	0.56	0.76

Table B.7: Stability and purity of chosen  $M_T^W$  bins in the electron channel for  $t\bar{t}$  MC events.

bin, GeV	$0 < M_T^W < 30$	$30 \leq M_T^W < 50$	$50 \leq M_T^W < 80$	$80 \leq M_T^W < 100$	$MT \geq 100$
events	11509	11633	26288	14824	8921
purity	0.57	0.35	0.65	0.4	0.36
stability	0.62	0.37	0.53	0.41	0.67

Table B.8: Stability and purity of chosen  $M_T^W$  bins in the muon channel for  $t\bar{t}$  MC events.

bin, GeV	$0 < M_T^W < 30$	$30 \leq M_T^W < 50$	$50 \leq M_T^W < 80$	$80 \leq M_T^W < 100$	$MT \geq 100$
events	13209	13636	30669	16974	9179.3
purity	0.56	0.36	0.66	0.42	0.37
stability	0.64	0.38	0.54	0.42	0.66

## C. Fitting templates

### $H_T$ variable

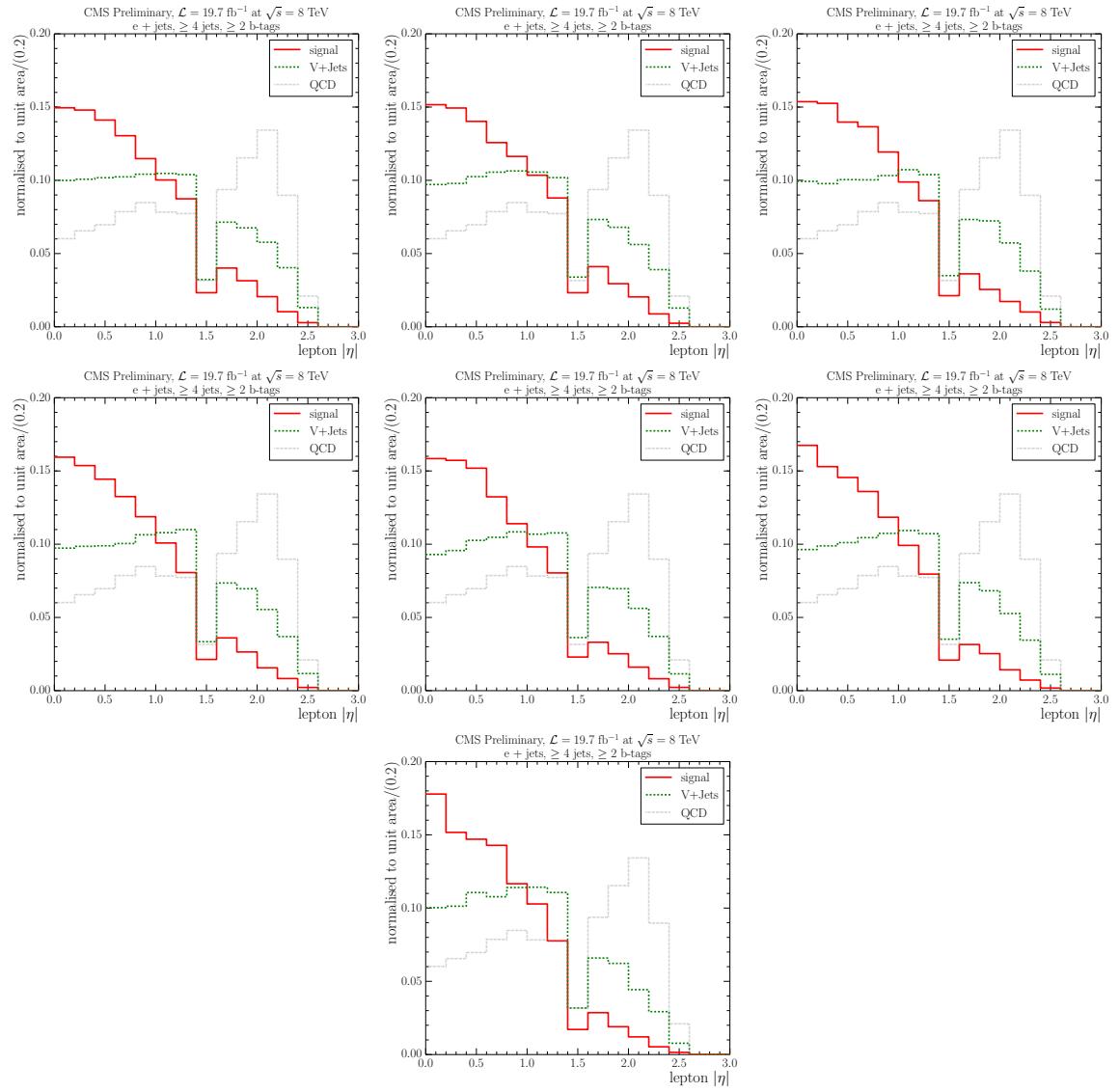


Figure C.1: Electron  $|\eta|$  templates for the fit in different bins of  $H_T$ , from top left to bottom right: 0–240 GeV, 240–280 GeV, 280–330 GeV, 330–380 GeV, 380–450 GeV, 450–600 GeV and  $\geq 600$  GeV.

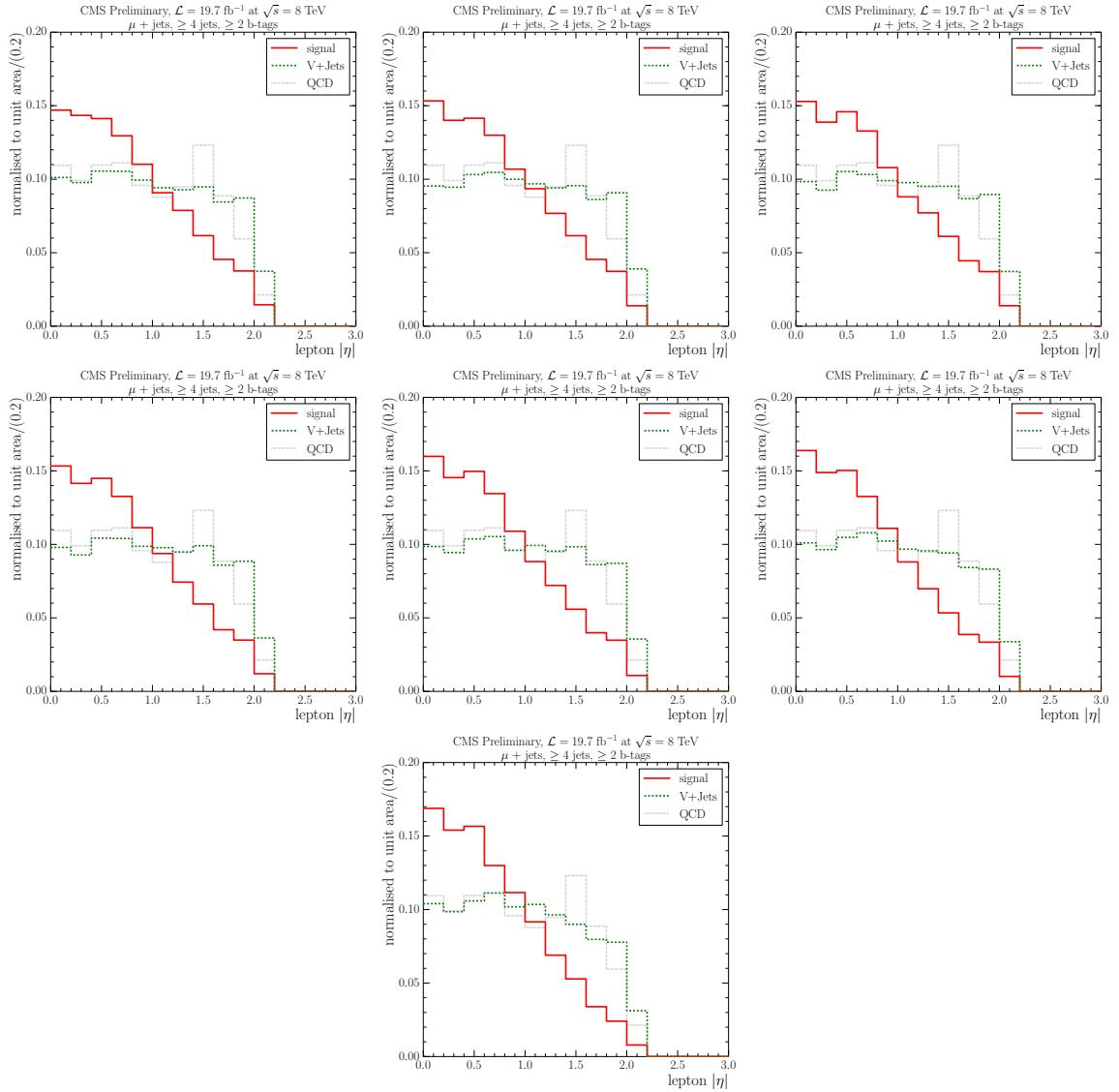


Figure C.2: Electron  $|\eta|$  templates for the fit in different bins of  $H_T$ , from top left to bottom right: 0–240 GeV, 240–280 GeV, 280–330 GeV, 330–380 GeV, 380–450 GeV, 450–600 GeV and  $\geq 600$  GeV.

## $S_T$ variable

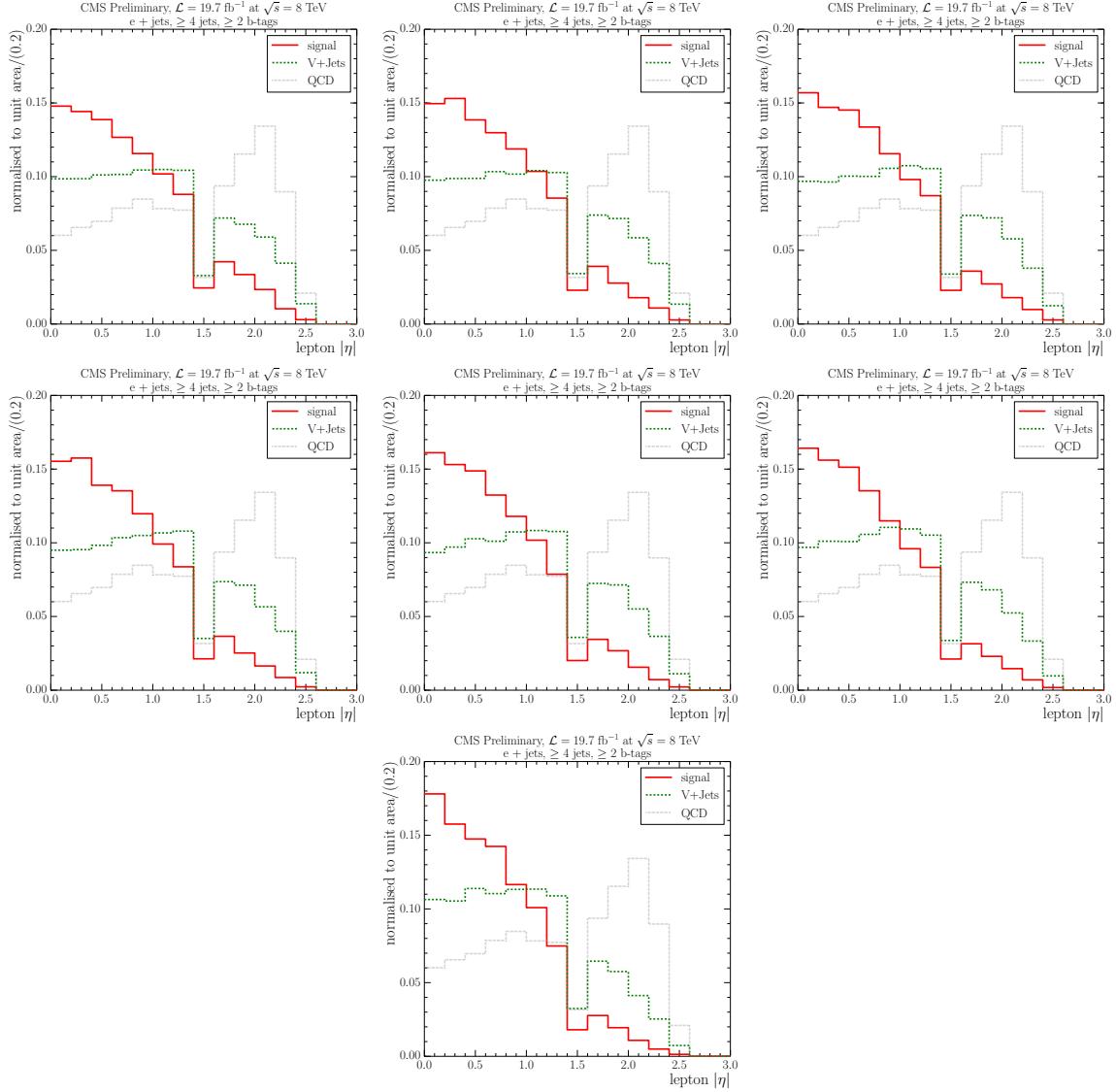


Figure C.3: Electron  $|\eta|$  templates for the fit in different bins of  $S_T$ , from top left to bottom right: 0–350 GeV, 350–400 GeV, 400–450 GeV, 450–500 GeV, 500–580 GeV, 580–700 GeV and  $\geq 700$  GeV.

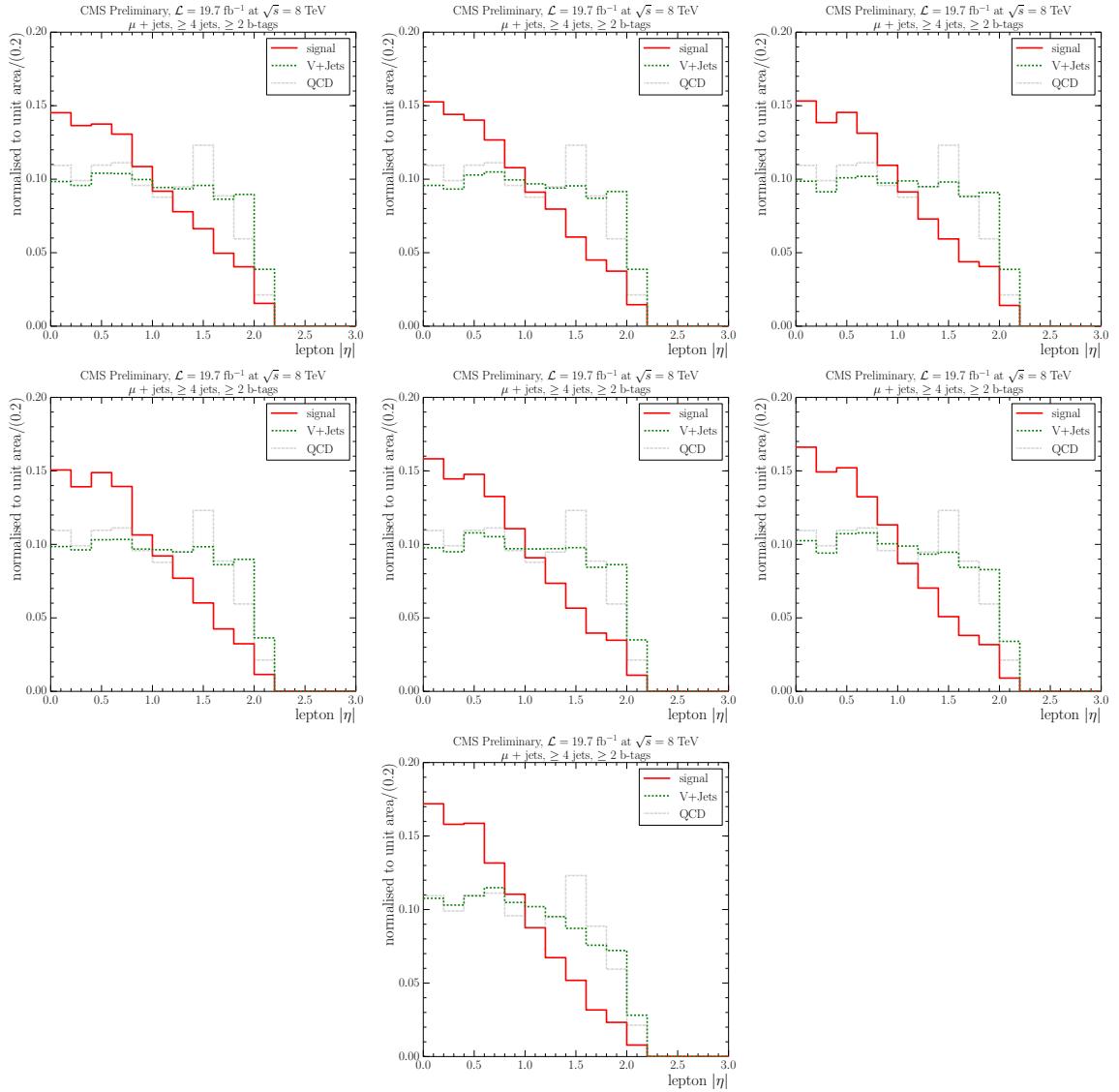


Figure C.4: Electron  $|\eta|$  templates for the fit in different bins of  $S_T$ , from top left to bottom right: 0–350 GeV, 350–400 GeV, 400–450 GeV, 450–500 GeV, 500–580 GeV, 580–700 GeV and  $\geq 700$  GeV.

## $p_T^W$ variable

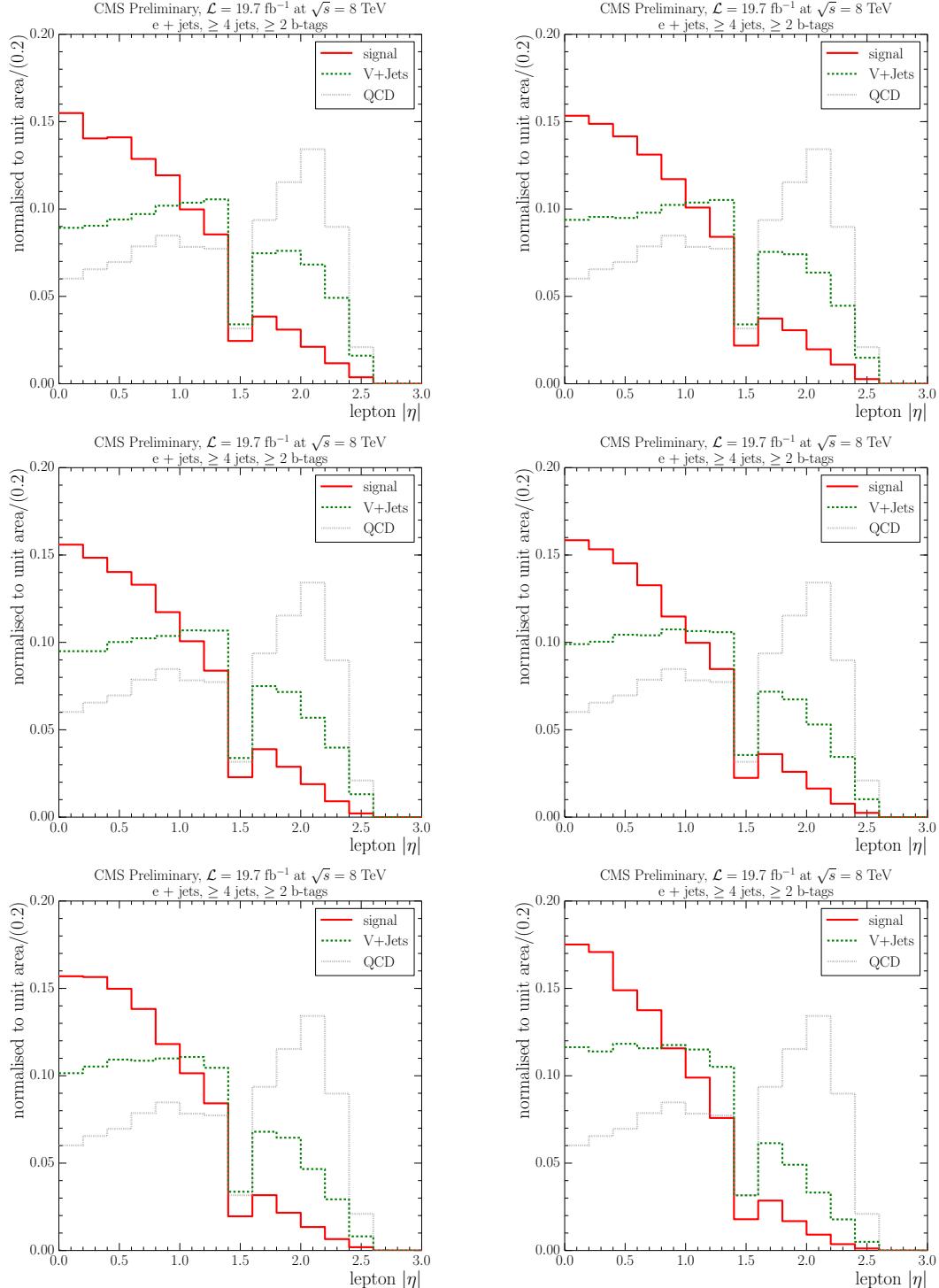


Figure C.5: Electron  $|\eta|$  templates for the fit in different bins of  $p_T^W$ , from top left to bottom right: 0–40 GeV, 40–70 GeV, 70–100 GeV, 100–130 GeV, 130–170 GeV and  $\geq 170$  GeV.

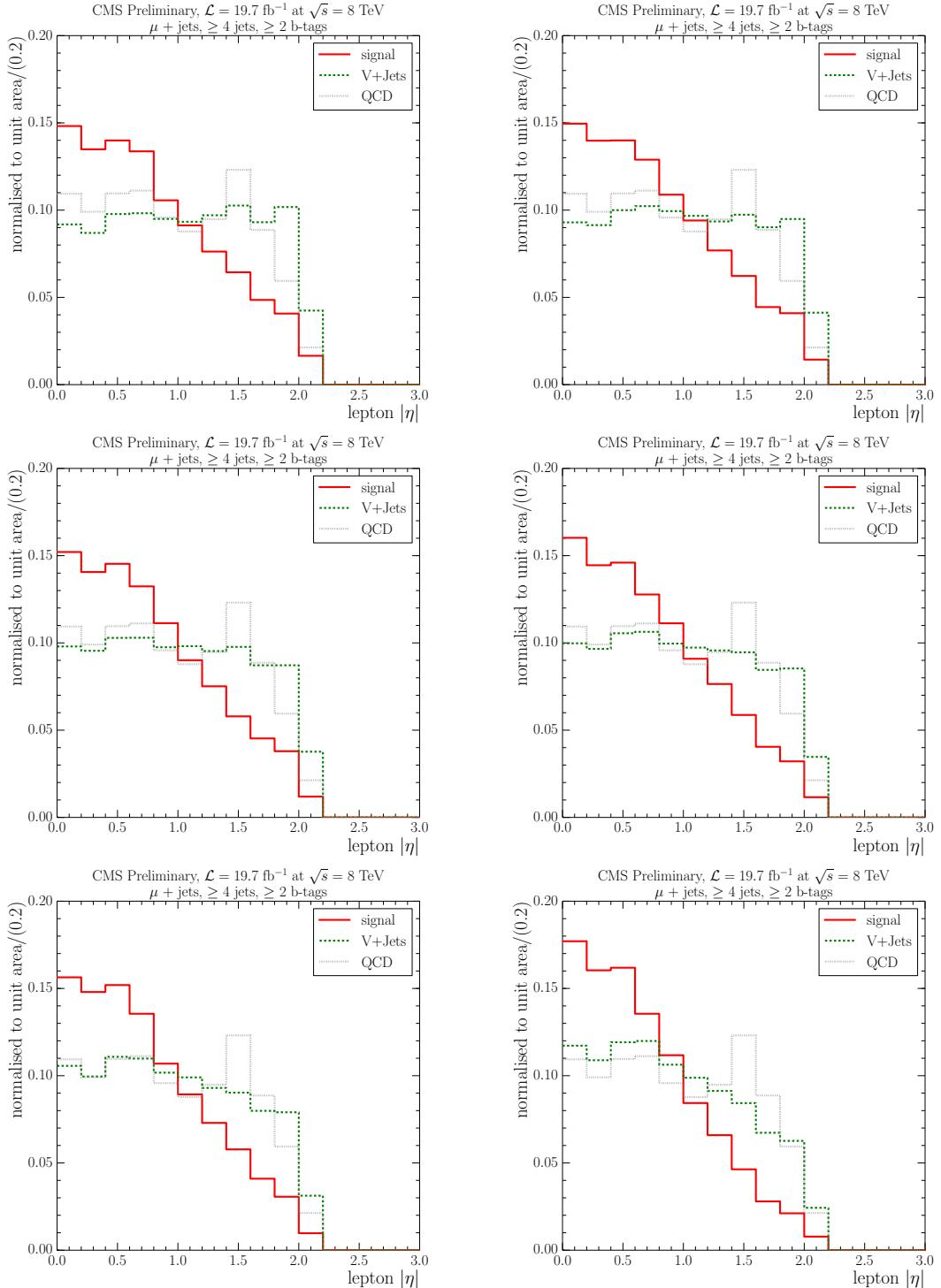


Figure C.6: Muon  $|\eta|$  templates for the fit in different bins of  $p_T^W$ , from top left to bottom right: 0–40 GeV, 40–70 GeV, 70–100 GeV, 100–130 GeV, 130–170 GeV and  $\geq 170$  GeV.

## $M_T^W$ variable

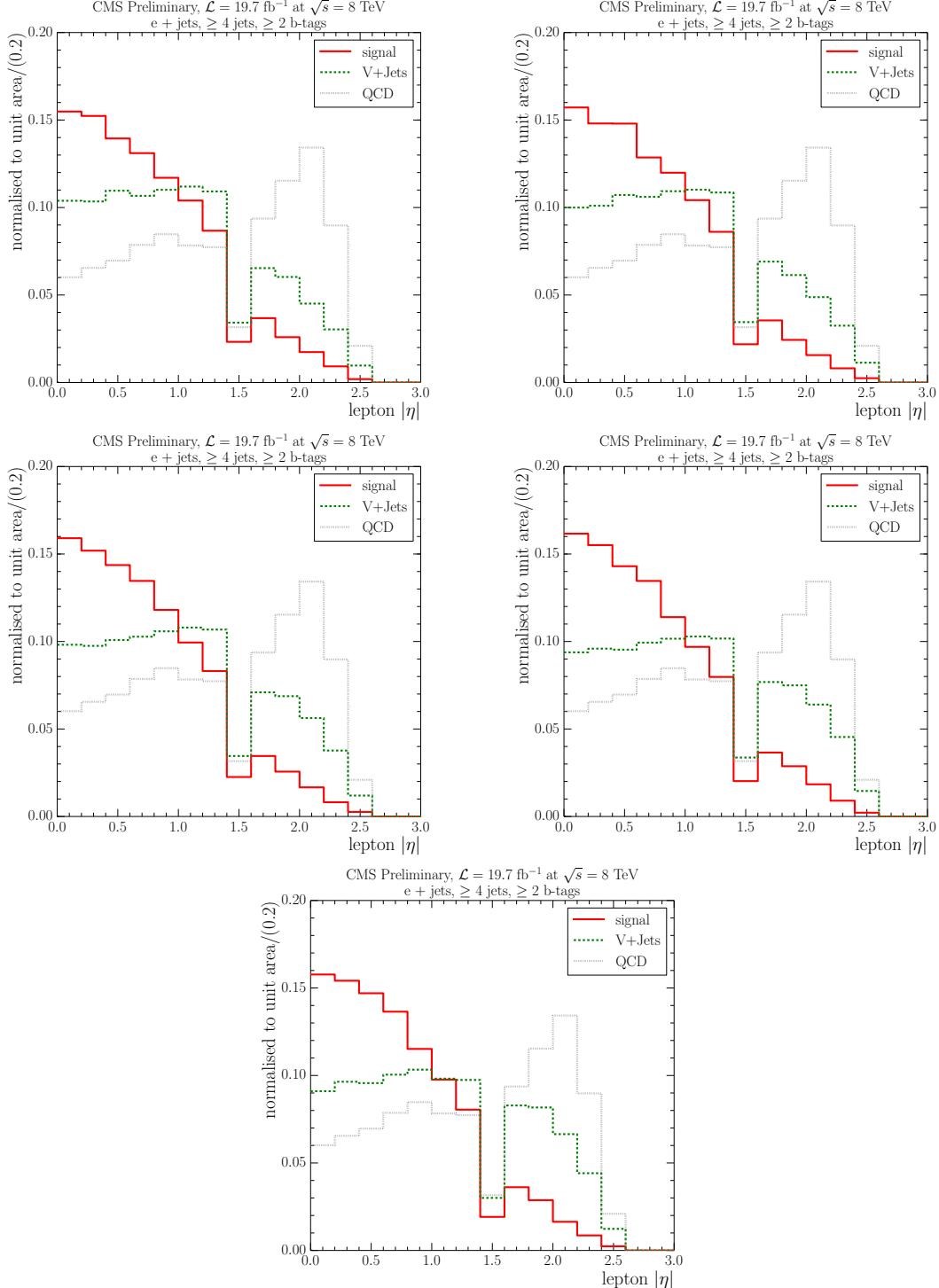


Figure C.7: Electron  $|\eta|$  templates for the fit in different bins of  $M_T^W$ , from top left to bottom right: 0–30 GeV, 30–50 GeV, 50–80 GeV, 80–100 GeV and  $\geq 100$  GeV.

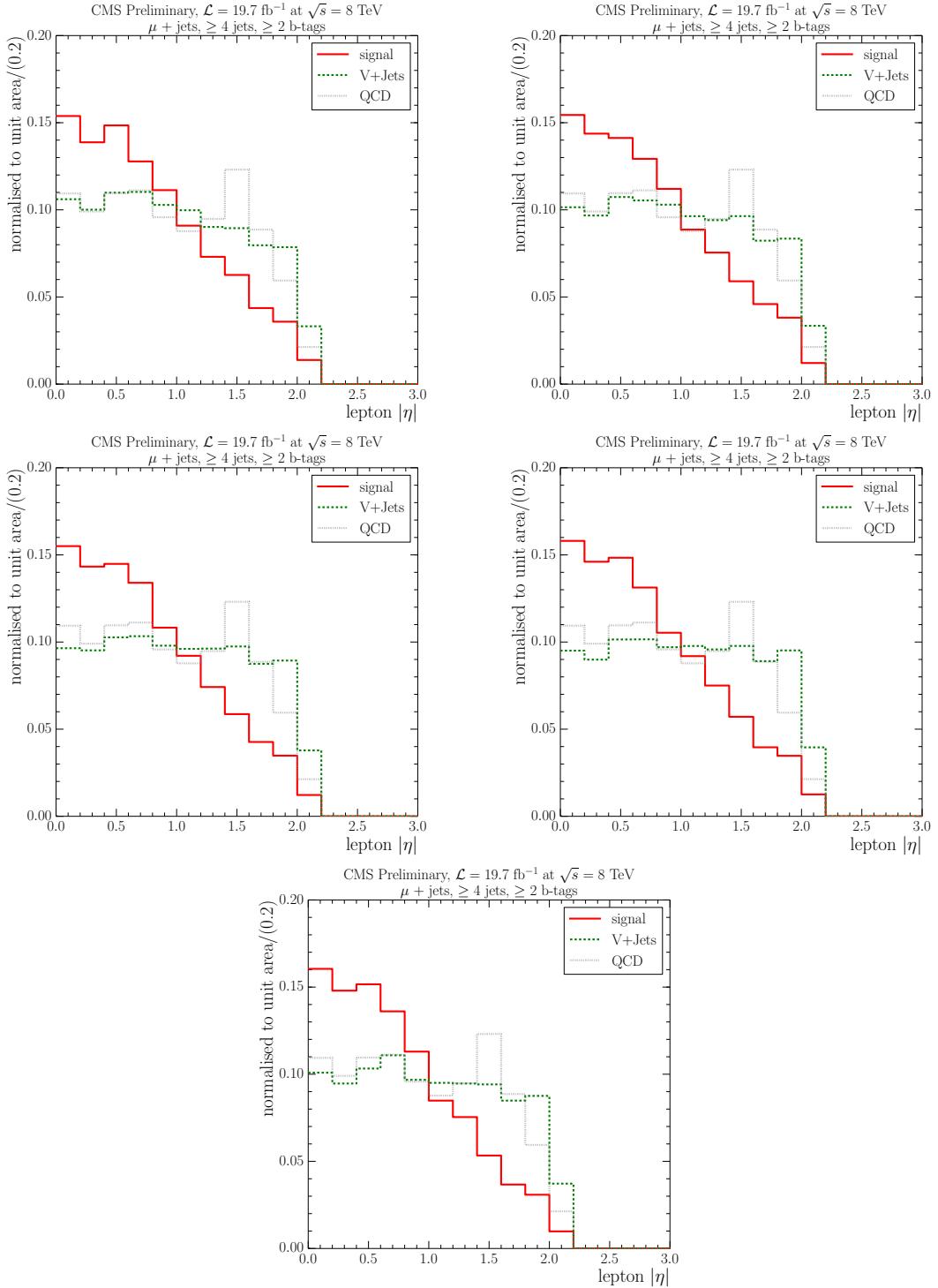


Figure C.8: Muon  $|\eta|$  templates for the fit in different bins of  $M_T^W$ , from top left to bottom right: 0–30 GeV, 30–50 GeV, 50–80 GeV, 80–100 GeV and  $\geq 100$  GeV.

# References

- [1] L. Evans and P. Bryant. “LHC Machine”. *JINST* 3.08 (2008), S08001. DOI: 10.1088/1748-0221/3/08/S08001.
- [2] J.P. Blewett. “200 GeV intersecting storage accelerators”. *Proceedings of The 8th International Conference on High-Energy Accelerators*. CERN. Geneva, Switzerland, 1971, p. 501.
- [3] The CMS Collaboration. “The CMS experiment at the CERN LHC”. *JINST* 3.08 (2008), S08004. DOI: 10.1088/1748-0221/3/08/S08004.
- [4] The CMS Collaboration. “Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC”. *Physics Letters B* 716.1 (2012), pp. 30–61. DOI: 10.1016/j.physletb.2012.08.021.
- [5] The CMS Collaboration. *CMS Physics: Technical Design Report Volume 1: Detector Performance and Software*. Technical Design Report CMS. Geneva, Switzerland: CERN, 2006.
- [6] The CMS Collaboration. “Precise mapping of the magnetic field in the CMS barrel yoke using cosmic rays”. *Journal of Instrumentation* 5.03 (2010). DOI: 10.1088/1748-0221/5/03/T03021.
- [7] Rene Brun and Fons Rademakers. “ROOT – An object oriented data analysis framework”. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 389.1–2 (1997). New Computing Techniques in Physics Research V, pp. 81 –86. DOI: 10.1016/S0168-9002(97)00048-X.

- [8] *CMSSW software framework*. URL: <http://cms-sw.github.io/cmssw>.
- [9] *Bristol Analysis Tools*. URL: <https://github.com/BristolTopGroup/AnalysisSoftware>.
- [10] *Daily Python Scripts*. URL: <https://github.com/BristolTopGroup/DailyPythonScripts>.
- [11] *rootpy: Pythonic ROOT package*. URL: <http://rootpy.org>.
- [12] *matplotlib plotting library*. URL: <http://matplotlib.org>.
- [13] The CMS Collaboration. *Particle-Flow Event Reconstruction in CMS and Performance for Jets, Taus, and MET*. Tech. rep. CMS-PAS-PFT-09-001. Geneva, Switzerland: CERN, 2009.
- [14] R. Frühwirth. “Application of Kalman filtering to track and vertex fitting”. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 262.2–3 (1987), pp. 444 –450. DOI: 10.1016/0168-9002(87)90887-4.
- [15] W Adam et al. “Reconstruction of electrons with the Gaussian-sum filter in the CMS tracker at the LHC”. *Journal of Physics G: Nuclear and Particle Physics* 31.9 (2005), N9. DOI: 10.1088/0954-3899/31/9/N01.
- [16] S. Baffioni et al. “Electron reconstruction in CMS”. *The European Physical Journal C* 49.4 (2007), pp. 1099–1116. DOI: 10.1140/epjc/s10052-006-0175-5.
- [17] The CMS Collaboration. *Cuts in Categories (CiC) Electron Identification*. Online; accessed 04-September-2013. URL: <https://twiki.cern.ch/twiki/bin/view/CMSPublic/SWGuideCategoryBasedElectronID>.
- [18] The CMS collaboration. “Performance of CMS muon reconstruction in pp collision events at  $\sqrt{s} = 7$  TeV”. *Journal of Instrumentation* 7.10 (2012), P10002. DOI: 10.1088/1748-0221/7/10/P10002.
- [19] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. “The anti- $k_t$  jet clustering algorithm”. *Journal of High Energy Physics* 2008.04 (2008), p. 063. DOI: 10.1088/1126-6708/2008/04/063.
- [20] Rick D. Field. “The underlying event in hard scattering processes”. *eConf* C010630 (2001), P501. arXiv: hep-ph/0201192.
- [21] The CMS collaboration. “Identification of b-quark jets with the CMS experiment”. *Journal of Instrumentation* 8.04 (2013), P04013. DOI: 10.1088/1748-0221/8/04/P04013.

- [22] The CMS Collaboration. *CMS TriDAS project: Technical Design Report, Volume 1: The Trigger Systems*. CERN-LHCC-2000-038; CMS-TDR-6-1. CERN/LHCC, 2000.
- [23] J. Brooke. “Performance of the CMS Level-1 Trigger”. *Proceedings of Science* ICHEP2012 (2013), p. 508. arXiv: 1302.2469 [[hep-ex](#)].
- [24] The CMS collaboration. “The CMS High Level Trigger System”. *2007 15th IEEE-NPSS Real-Time Conference*. Institute of Electrical and Electronics Engineers, 2007. DOI: 10.1109/RTC.2007.4382773.
- [25] The CMS collaboration. “Commissioning of the CMS High Level Trigger”. *Journal of Instrumentation* 4.10 (2009), P10005. DOI: 10.1088/1748-0221/4/10/P10005.
- [26] Wolfgang Adam et al. *Track Reconstruction in the CMS tracker*. Tech. rep. CMS-NOTE-2006-041. Geneva, Switzerland: CERN, 2006.
- [27] Diego Casadei. “Estimating the selection efficiency”. *Journal of Instrumentation* 7.08 (2012), P08021. DOI: 10.1088/1748-0221/7/08/P08021.
- [28] M. Agelou et al. *Top Trigger Efficiency Measurements and the top trigger package*. Tech. rep. DØ Note 4512. 2004.
- [29] The CMS Collaboration. “Measurement of the top-quark mass in  $t\bar{t}$  events with lepton+jets final states in pp collisions at  $\sqrt{s} = 7$  TeV”. *Journal of High Energy Physics* 2012.12 (2012), pp. 1–37. DOI: 10.1007/JHEP12(2012)105.
- [30] The CMS Collaboration. “Measurement of the mass difference between top and anti-top quarks”. *Journal of High Energy Physics* 2012.6 (2012), pp. 1–36. DOI: 10.1007/JHEP06(2012)109.
- [31] The CMS Collaboration. *Measurement of the top quark mass in the l+jets channel*. Tech. rep. CMS-PAS-TOP-10-009. Geneva, Switzerland: CERN, 2011.
- [32] Johan Alwall et al. “MadGraph 5: going beyond”. *Journal of High Energy Physics* 2011.6 (2011), pp. 1–40. DOI: 10.1007/JHEP06(2011)128.
- [33] Torbjörn Sjöstrand et al. “High-energy-physics event generation with Pythia 6.1”. *Computer Physics Communications* 135.2 (2001), pp. 238–259. DOI: 10.1016/S0010-4655(00)00236-8.
- [34] Torbjörn Sjöstrand, Stephen Mrenna, and Peter Z. Skands. “PYTHIA 6.4 Physics and Manual”. *JHEP* 05 (2006). DOI: 10.1088/1126-6708/2006/05/026. arXiv: hep-ph/0603175.

- [35] Stefano Frixione, Paolo Nason, and Carlo Oleari. "Matching NLO QCD computations with parton shower simulations: the POWHEG method". *Journal of High Energy Physics* 11 (2007). DOI: 10.1088/1126-6708/2007/11/070.
- [36] Z. Was. "TAUOLA the library for Ď lepton decay, and KKMC/KORALB/KORALZ/â€ status report". *Nuclear Physics B - Proceedings Supplements* 98.1â€ 3 (2001), pp. 96 –102. DOI: 10.1016/S0920-5632(01)01200-2.
- [37] GEANT4 Collaboration. "GEANT4 - a simulation toolkit". *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 506.3 (2003), pp. 250 –303. DOI: 10.1016/S0168-9002(03)01368-8.
- [38] Stefan Höche et al. "Matching parton showers and matrix elements" (2006). Proceedings of the "HERA and the LHC" workshop, CERN/DESY 2004/2005". arXiv: hep-ph/0602031.
- [39] Stefano Frixione and Bryan R. Webber. "Matching NLO QCD computations and parton shower simulations". *Journal of High Energy Physics* 06 (2002). DOI: 10.1088/1126-6708/2002/06/029.
- [40] Daniel R. Stump. "A new generation of CTEQ parton distribution functions with uncertainty analysis". *Proceedings to ICHEP 2002* 117 (2003), pp. 265 –267. DOI: 10.1016/S0920-5632(03)90541-X.
- [41] The CMS Collaboration. *Tracking and Primary Vertex Results in First 7 TeV Collisions*. Tech. rep. CMS-PAS-TRK-10-005. Geneva, Switzerland: CERN, 2010.
- [42] Suharyo Sumowidagdo. *Programming documentation of HitFit: A Kinematic Fitter for top quark-antiquark pair lepton+jets event*. CMS AN-2011/171. 2011.
- [43] Orin I. Dahl, T.B. Day, and Frank T. Solmitz. *SQUAW Kinematic Fitting Program*. Tech. rep. P-126. Berkeley, California: Lawrence Radiation Laboratory, 1968.
- [44] DØ Collaboration. "Direct measurement of the top quark mass by the DØ Collaboration". *Phys. Rev. D* 58 (5 1998). DOI: 10.1103/PhysRevD.58.052001.
- [45] DØ Collaboration. "Measurement of the top quark mass in the lepton+jets channel using the ideogram method". *Phys. Rev. D* 75 (9 2007). DOI: 10.1103/PhysRevD.75.092001.

- [46] DØ Collaboration. “Measurement of the Forward-Backward Charge Asymmetry in Top-Quark Pair Production”. *Phys. Rev. Lett.* 100 (14 2008). DOI: 10.1103/PhysRevLett.100.142002.
- [47] *CLHEP - A Class Library for High Energy Physics*. URL: <http://cern.ch/CLHEP>.
- [48] Scott Stuart Snyder. “Measurement of the Top Quark Mass at DØ”. FERMILAB-THESIS-1995-27. PhD thesis. Stony Brook University, 1995.
- [49] DELPHI Collaboration. “Measurement of the W-pair cross-section and of the W mass in  $e^+e^-$  interactions at 172 GeV”. *The European Physical Journal C - Particles and Fields* 2.4 (1998), pp. 581–595. DOI: 10.1007/s100529800857.
- [50] DELPHI Collaboration. “Measurement of the mass and width of the W boson in  $e^+e^-$  collisions at  $\sqrt{s} = 161\text{--}209$  GeV”. *The European Physical Journal C* 55.1 (2008), pp. 1–38. DOI: 10.1140/epjc/s10052-008-0585-7.
- [51] CDF Collaboration. “Measurement of the Top-Quark Mass in All-Hadronic Decays in  $p\bar{p}$  Collisions at CDF II”. *Phys. Rev. Lett.* 98 (14 2007). DOI: 10.1103/PhysRevLett.98.142001.
- [52] DØ Collaboration. “Improved determination of the width of the top quark”. *Phys. Rev. D* 85 (9 2012), p. 091104. DOI: 10.1103/PhysRevD.85.091104.
- [53] The ALEPH Collaboration et al. “Precision electroweak measurements on the Z resonance”. *Physics Reports* 427.5–6 (2006), pp. 257 –454. DOI: 10.1016/j.physrep.2005.12.006.
- [54] K.S. Kölbig and B. Schorr. “A program package for the Landau distribution”. *Computer Physics Communications* 31.1 (1984), pp. 97 –111. DOI: 10.1016/0010-4655(84)90085-7.
- [55] The CMS collaboration. “Determination of jet energy calibration and transverse momentum resolution in CMS”. *Journal of Instrumentation* 6.11 (2011). DOI: 10.1088/1748-0221/6/11/P11002.
- [56] Daniel Wicke and Peter Z. Skands. “Non-perturbative QCD Effects and the Top Mass at the Tevatron”. *Nuovo Cim.* B123 (2008), S1. arXiv: 0807.3248.
- [57] Peter Z. Skands. “Tuning Monte Carlo generators: The Perugia tunes”. *Phys. Rev. D* 82 (7 2010), p. 074018. DOI: 10.1103/PhysRevD.82.074018.

- [58] CDF and DØ Collaborations. “Combination of the top-quark mass measurements from the Tevatron collider”. *Phys. Rev. D* 86 (9 2012), p. 092003. doi: 10.1103/PhysRevD.86.092003.
- [59] ATLAS and CMS Collaborations. *Combination of ATLAS and CMS results on the mass of the top-quark using up to 4.9 fb<sup>-1</sup> of  $\sqrt{s} = 7$  TeV LHC data*. Tech. rep. ATLAS-CONF-2013-102. Geneva: CERN, 2013.
- [60] The CMS Collaboration. *Measurement of missing transverse energy in top pair events at  $\sqrt{s} = 7$  TeV*. Tech. rep. CMS-PAS-TOP-12-019. Geneva: CERN, 2012.
- [61] The CMS Collaboration. *Measurement of missing transverse energy and other event-level observables in top pair events at  $\sqrt{s} = 8$  TeV*. Tech. rep. CMS-PAS-TOP-12-042. Geneva: CERN, 2013.
- [62] The CMS Collaboration. *CMS Luminosity Based on Pixel Cluster Counting - Summer 2012 Update*. Tech. rep. CMS-PAS-LUM-12-001. Geneva: CERN, 2012.
- [63] The CMS Collaboration. *Inelastic pp cross section at 7 TeV*. Tech. rep. CMS-PAS-FWD-11-001. Geneva: CERN, 2011.
- [64] The CMS Collaboration. *Results on b-tagging identification in 8 TeV pp collisions*. Tech. rep. CMS-DP-2013-005. 2013.
- [65] Andrea Rizzi. “Events weights from b-tagging efficiency SF”. URL: <https://twiki.cern.ch/twiki/pub/CMS/BTagWeight/BTagWeight.pdf>.
- [66] Pushpalatha Bhat et al. *Simulation Studies of the Radiation Environment in the CMS Detector for pp Collisions at the LHC*. Tech. rep. CERN-CMS-NOTE-2013-001. Geneva: CERN, 2010.
- [67] Michał Czakon, Paul Fiedler, and Alexander Mitov. “The total top quark pair production cross-section at hadron colliders through O( $\alpha_S^4$ )”. *Phys. Rev. Lett.* 110 (25 2013). doi: 10.1103/PhysRevLett.110.252004.
- [68] The CMS Collaboration. “Measurement of differential top-quark-pair production cross sections in pp collisions at  $\sqrt{s} = 7$  TeV”. *The European Physical Journal C* 73.3 (2013). doi: 10.1140/epjc/s10052-013-2339-4.
- [69] The CMS Collaboration. *Measurement of differential top-quark pair production cross sections in the lepton+jets channel in pp collisions at 8 TeV*. Tech. rep. CMS-PAS-TOP-12-027. Geneva: CERN, 2013.

- [70] Andreas Höcker and Vakhtang Kartvelishvili. “SVD approach to data unfolding”. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 372.3 (1996), pp. 469–481. DOI: 10.1016/0168-9002(95)01478-0.
- [71] T. Adye. “Unfolding algorithms and tests using RooUnfold”. *Proceedings of the PHYSTAT 2011 Workshop*. CERN. Geneva, Switzerland, 2011, pp. 313–318.