

A Measurement of the Underlying Event Distributions in Proton-Proton Collisions at \sqrt{s}
= 7 TeV in Events containing Charged Particle Jets using the ATLAS Detector at the
Large Hadron Collider

by

Joseph Salvatore Virzi

A dissertation submitted in partial satisfaction of the

requirements for the degree of

Doctor of Philosophy

in

Physics

in the

Graduate Division

of the

University of California, Berkeley

Committee in charge:

Professor M. D. Shapiro, Chair

Professor Robert Jacobsen

Professor Jenny Harrison

Spring 2012



A Measurement of the Underlying Event Distributions in Proton-Proton Collisions at \sqrt{s}
= 7 TeV in Events containing Charged Particle Jets using the ATLAS Detector at the
Large Hadron Collider

Copyright 2012

by

Joseph Salvatore Virzi

Abstract

A Measurement of the Underlying Event Distributions in Proton-Proton Collisions at $\sqrt{s} = 7$ TeV in Events containing Charged Particle Jets using the ATLAS Detector at the Large Hadron Collider

by

Joseph Salvatore Virzi

Doctor of Philosophy in Physics

University of California, Berkeley

Professor M. D. Shapiro, Chair

Underlying Event distributions are studied in events containing at least one charged-particle jet produced in proton-proton collisions at $\sqrt{s} = 7$ TeV. Jets are reconstructed from charged particles using the anti- k_r algorithm with radius parameter $R = 0.6$. The jet with the largest transverse momentum p_T^{jet} and $|\eta^{\text{jet}}| \leq 1.5$ defines the azimuthal ϕ^{jet} direction. Distributions of the charged particle multiplicity, the scalar sum of the transverse momenta (p_T) of charged particles, and the average charged particle p_T are measured as functions of p_T^{jet} in the transverse region ($\frac{\pi}{3} \leq |\phi - \phi^{\text{jet}}| \leq \frac{2\pi}{3}$) for $4 \text{ GeV} \leq p_T^{\text{jet}} \leq 100 \text{ GeV}$. The data are compared to predictions from the Monte Carlo generators which have been tuned to data from the Large Hadron Collider, and are found in to be good agreement.

Dedicado a la memoria de mi madre

Acknowledgments

I certainly did not do this work alone. Maria Elena, *mi alma gemela y amor de mi vida*, and my children Sofia Elena, Maria Elena and Valentina Elena, *las adoro con toda mi alma*. Everyone in my life has played a role in the completion of my dissertation. Marjorie Shapiro, my advisor, was there for me through every step. Many ATLAS physicists provided extensive peer review of this work, including the SM groups and editorial boards for the supporting note of a publication based on this analysis. I would like to thank the members of my committee, Bob Jacobsen and Jenny Harrison. Special thanks to Kevin Einsweiler, Paul Laycock, Emily Nurse, Liza Mijovic, Frank Paige, Claudia Glasman, Andy Pilkington, Jon Butterworth and Deepak Kar. Gabe Hare jump-started me. The eye-opening discussions with Andrei Gaponenko and Ayana Holloway-Arce. Johannes Muelmenstaedt was the voice of insanity, and always there for the kind of support no one else could provide, day or night. Mostly night. Aathi always a great sounding board. I am grateful to Hung-Chung Fang for thesis-altering physics and analysis discussions and invaluable critique, especially right at the end when I could no longer read my own writing. My father, mother, and brothers, Nicolas and Vincent, shaped my life throughout this journey. I can't believe how helpful Patrick McGuire was in helping me work out technical problems. I would acknowledge more people, but I've already submitted this dissertation.

Contents

List of Figures	v
List of Tables	vii
1 Introduction and Theoretical Overview	1
2 Monte Carlo	9
2.1 Generation	9
2.2 Simulation	10
2.3 Remarks	11
3 The ATLAS Inner Detector at the Large Hadron Collider	13
3.1 The Large Hadron Collider at CERN	13
3.2 The ATLAS detector	13
3.3 The ATLAS Inner Detector	14
3.3.1 The PIXEL detector	15
3.3.2 The Semiconductor Tracker	16
3.3.3 The Transition Radiation Tracker	16
3.4 ATLAS Trigger Overview	17
3.4.1 Minimum Bias Trigger	18
3.5 Event Reconstruction	18
3.5.1 Track Reconstruction	18
3.5.2 Vertex Reconstruction	19
3.5.3 Jet Reconstruction	21
4 Measurement of Raw UE Distributions	23
4.1 Event and Data Selection	23
4.1.1 Track Selection	23
4.1.2 Jet Selection	24
4.2 Measuring the Distributions from Data	25

5	Correcting the UE Distributions for Detector Effects	30
5.1	The Response Matrix	30
5.1.1	Construction of the Response Matrix	34
5.1.2	Purity and Stability	34
5.2	Bayesian Iterative Unfolding with RooUnfold	36
5.3	Validation of the Unfolding Procedure	39
5.4	Corrected Distributions	39
5.4.1	Mean Values of Corrected Distributions	44
6	Uncertainty Analysis	46
6.1	Track Reconstruction	47
6.1.1	Track Momentum Resolution	47
6.1.2	Tracking Efficiency	47
6.2	Uncertainty in the Unfolding Procedure	49
6.2.1	Unfolding Uncertainty for Individual Bins of p_T^{jet} and \mathcal{O}	49
6.2.2	Unfolding Uncertainty in the Mean Values of \mathcal{O}	54
6.3	Sensitivity to the Response Matrix	54
6.4	Statistical Uncertainty in the Response Matrix	55
6.5	Misidentification of the Leading Jet	56
6.6	Discretization Effects	57
6.7	Dependence on Number of Iterations	59
6.8	Statistical Uncertainties	59
6.9	Summary of Total Systematic Uncertainties	60
6.10	Consistency Checks - Refolding the Distributions	60
7	Conclusions	71
	Bibliography	72
A	Track Quality	75
B	Discretization Effects	77
B.1	Discretization Effects in Σp_T	77
B.2	Discretization Effects in N_{ch}	80
B.3	Discretization Effects in \bar{p}_T	81
B.4	Validation of the Spline-based Methods	81

List of Figures

1.1	Example diagram for $pp \rightarrow pp + jj$	5
1.2	Example diagram for $pp \rightarrow pp + jj + \pi^+ \pi^-$	5
1.3	Regions of the Underlying Event	6
3.1	The CERN Accelerator Complex	14
3.2	Diagram of the ATLAS detector.	15
3.3	An $r - z$ View of the Inner Detector	16
3.4	Cutaway View of the ATLAS Inner Detector	17
3.5	Schematic of the different stages of track reconstruction.	20
3.6	Jet reconstruction efficiency	21
3.7	Measured jet p_T	22
4.1	Reconstructed track p_T and η distributions	25
4.2	Reconstructed track p_T and η distributions in the TRANSVERSE region	26
4.3	Measured Σp_T distributions	27
4.4	Measured N_{ch} distributions	28
4.5	Measured \bar{p}_T distributions	29
5.1	A slice of an example response matrix	32
5.2	Example response matrix	32
5.3	Stability for PYTHIA 6 (AMBT1)	37
5.4	Purity for PYTHIA 6 (AMBT1)	38
5.5	Closure tests	39
5.6	Corrected Σp_T distributions	41
5.7	Corrected N_{ch} distributions	42
5.8	Corrected \bar{p}_T distributions	43
5.9	Mean values of the corrected Σp_T distributions	44
5.10	Mean values of the corrected N_{ch} distributions	45
5.11	Mean values of the corrected \bar{p}_T distributions	45
6.1	Uncertainties in mean value due to tracking efficiency	48
6.2	Distribution of Σp_T corresponding to $14 \text{ GeV} \leq p_T^{\text{jet}} \leq 19 \text{ GeV}$	49
6.3	Slice in Σp_T comparing different Monte Carlo	51

6.4	Example of unfolding uncertainty for slice in Σp_T	53
6.5	Comparison of unfolding and statistical uncertainties	53
6.6	Mean values of UE distributions using reweighted response matrix	55
6.7	Comparing mean values using binned and spline-based methods	58
6.8	Mean values of Σp_T unfolding using different number of iterations	60
6.9	Uncertainties in the Σp_T distributions	62
6.10	Uncertainties in the Σp_T distributions (cont.)	63
6.11	Uncertainties in the N_{ch} distributions	64
6.12	Uncertainties in the N_{ch} distributions (cont.)	65
6.13	Uncertainties in the \bar{p}_T distributions	66
6.14	Uncertainties in the \bar{p}_T distributions (cont.)	67
6.15	Refolded Σp_T distributions	68
6.16	Refolded N_{ch} distributions	69
6.17	Refolded \bar{p}_T distributions	70
A.1	Track selection for ID hit requirements	75
A.2	Track selection for primary vertex requirements	76
B.1	Comparison of the Σp_T mean values using different methods	78
B.2	Comparison of the N_{ch} mean values using different methods	81
B.3	Spline-based predictions of the UE distributions	82

List of Tables

2.1	GEANT4-simulated Monte Carlo generation	12
4.1	Track selection criteria	24
5.1	Example unnormalized response matrix	31
5.2	Example response matrix	33
5.3	Example purity and stability calculations	36
6.1	Systematic uncertainties in the mean values of UE distributions	46
6.2	Example calculation of unfolding uncertainty	52

Chapter 1

Introduction and Theoretical Overview

Performing precision physics measurements at a hadron collider, such as the Large Hadron Collider (LHC) [1] where opposing beams of protons collide at unprecedented energies, requires that we be able to model not only the relevant energetic (hard scattering) processes, but also the softer (long-distance) components of the interactions.

Within the Standard Model (SM) theory of physics, Quantum Chromodynamics (QCD) is the theory of the *strong* nuclear force, the dominant force governing interactions in hadron colliders. The principle of local gauge symmetry states that the laws of physics exhibit an invariance under certain classes of transformations, and that these transformations can vary from point to point in spacetime. Using $\mathbf{M}(\mathbf{x}) \equiv \exp(i\lambda \mathbf{G}(\mathbf{x}) \tau)$ to denote a local unitary gauge transformation, terms containing the particle fields ϕ in the Lagrangian formulation of the SM remain invariant under the transformation $\phi \rightarrow \mathbf{M}\phi$. The spacetime dependence of the transformation is absorbed in the 8 fields $\mathbf{G}(\mathbf{x})$, one for each of the SU(3) group generators denoted by τ , and λ is the QCD coupling constant that determines the strength of the fields and their interactions. Three generations of *quark* pairs (arranged as (up, down), (charm, strange), and (top, bottom)) are described by the *fermionic* fields ϕ , and we refer to the \mathbf{G} fields as gluons, the *gauge bosons* of QCD. The terms fermion and boson are used to describe the representations of the fields under Lorentz (O(3,1)) transformations of spacetime coordinates, and can be characterized by degrees of freedom more commonly known as *spin*. Fermions have spin 1/2, 3/2, and so on. Bosons have integer values of spin.

The Lagrangian formalism contains not only the fields ϕ , but also derivatives of ϕ adding dynamics to the theory. Without these “kinetic” terms, there is no spacetime evolution of the fields, and the theory would be sterile. These kinetic terms of the form $\partial_x \phi$ invoke derivatives of the transformation when applying local gauge invariance, resulting in a huge number of terms in the Lagrangian^{1 2}.

¹Quantum field theories require the presence of all terms compatible with a proposed symmetry, or have good reason to exclude it.

²For the reader familiar with differential geometry, the gluons are analogous to the Christoffel symbols. The principle of local gauge symmetry gives rise to a covariant derivative, analogous to the covariant deriva-

Analogous to the theory of Laurent expansion of complex-valued functions, the theory of perturbative Quantum Chromodynamics (pQCD) attempts to predict cross sections (rates of interaction) by expanding the theory using λ . Just as $e^{\alpha x} \approx 1 + \alpha x$ for small αx , pQCD posits that if $\lambda \ll 1$, the effects of the transformation $\exp(i\lambda \mathbf{G}(\mathbf{x}) \tau)$ can be approximated by an expression of the form $1 + i\lambda \mathbf{G}(\mathbf{x}) \tau$. The interactions between particles and fields are then described by the expression of the form $\langle \phi(\mathbf{x}_0) | \phi(\mathbf{x}_1) \rangle = \langle \phi(\mathbf{x}_0) | e^{i(t_1 - t_0) \mathbf{H}} \phi(\mathbf{x}_0) \rangle$, where \mathbf{H} is the Hamiltonian operator derived from the Lagrangian. This compact expression contains a huge number of terms, due to the large number of terms in the Lagrangian. All these terms appear in the Hamiltonian which in turn "multiply" (couple to) each other, when calculating the expectation values $\langle \phi(\mathbf{x}_0) | \phi(\mathbf{x}_1) \rangle$. The complexity of the calculations was reduced significantly to a "topological" art form by Feynman, when he introduced the famous diagrams bearing his name. Each interaction term can be denoted graphically, and the diagrams carry rules for their computation.

We introduced λ , the QCD coupling "constant". λ is a constant with respect to the spacetime coordinates, but exhibits a dependence on the energy scales in question. This phenomenon is referred to as running of the coupling constant, $\lambda = \lambda(E)$ [2]. The behavior of λ is determined by the β -function [2] which is in turn determined by the $SU(3)$ group structure of QCD. λ is large (technically speaking, divergent) for low energies and small $\mathcal{O}(0.1)$ above $\Lambda_{QCD} \approx 200$ MeV, the so-called QCD scale. This is the phenomenon known as asymptotic freedom [3].

We use pQCD to calculate cross sections for interactions of *partons* (gluons and quarks) at high energy scales, where the QCD coupling constant is sufficiently small to make an expansion of the QCD interaction terms meaningful. At lower energy scales, however, the coupling constant becomes large and pQCD can no longer be used. We must therefore resort to an empirical approach to model the interaction terms.

The SM has a much richer structure than we have thus far described, $SU(3) \times SU(2) \times U(1)$ to be exact, incorporating the *electroweak* interactions. The same general principles we discussed about QCD apply to the electroweak interactions. The group structure is $SU(2) \times U(1)$, which has four generators and, therefore, four bosons. The electroweak bosons are the W^\pm , Z^0 and the photons. There are two coupling constants, much smaller than λ , and the theory is always perturbative. The fermion fields are the 3 generations of *lepton* pairs (arranged as (e^-, ν_e) , (μ^-, ν_μ) , and (τ^-, ν_τ) .) The quark fields also couple to the electroweak bosons. The left-handed components of the above quark and lepton pairs are *doublets* of the electroweak $SU(2)$ group.

We have discussed partons, but these are not the particles we actually observe in our detector. The principle of confinement states that free quarks do not exist in nature [4]. We can only directly observe leptons and final state *hadrons* which are bound states of 2 or 3 quarks. For two quarks, these are the QCD color singlets (color-neutral combinations) of the $SU(3) \times SU(3)$ product group (mesons). The observable bound states of 3 quarks are the color singlets of the $SU(3) \times SU(3) \times SU(3)$ product group (baryons). The momenta

of the hard scatter partons are highly correlated with the momenta of the hadrons.

As an example of perturbation theory applied to the electroweak and strong interactions, we can calculate the cross section for a top quark decaying to a W and a bottom quark ($t \rightarrow Wb$). We observe the W experimentally through either its leptonic ($W \rightarrow \ell + \nu$) or hadronic ($W \rightarrow q\bar{q}$) signature. For the reasons already discussed, we cannot detect an isolated b quark. The b quark will interact with the rest of the proton (beam remnants), possibly producing pairs of quarks from the vacuum, eventually *hadronizing* (forming hadrons) into, say, a B^0 or B^\pm . In our example, the direction of the B -meson is highly correlated to the direction of the outgoing b quark.

Hadronization of partons are extremely complex interactions, and occur at much lower energy scales (longer time scales) than does the hard scatter. These interactions elude calculation, and we must adopt alternate models to account for their effects. We feed back knowledge obtained from experiments to adjust parameters these models will invariably have.

We can use pQCD to calculate the high energy component of interactions between gluons and quarks, but the subsequent evolution of the interaction is impossible to calculate, as final state particles are produced through complex intermediate QCD parton exchanges (described by Feynman diagrams), Essentially, we cannot calculate the distributions of final state particles from first principles. The motivation for our measurement begins to nucleate. In order to relate physics measurements back to theoretical predictions, we must understand how the observed data feeds back to the theoretical model. We must determine how does the distribution of particles in our detector affects our interpretation of the hard scatter between two colliding protons.

One of the most indispensable tools we have are Monte Carlo (MC) generators, such as PYTHIA [5], SHERPA [6] and HERWIG++ [7]. These are computer programs that implement phenomenological physics models, and we shall discuss this topic in more detail in Chapter 2. These generators model the hard scatter between colliding protons, evolve the interactions through a series of models (e.g. - fragmentation, hadronization), and provide distributions of final state particles. The extent to which we believe the MC predictions depends on how well they can model the distributions of final state particles. Our measurement focuses on characterizing the agreement in the "tails" of these distributions, the lesser populated regions of phase space. Specifically, we measure the distributions of particles far away from the primary regions of interest and compare to the MC predictions.

Historically, the *Underlying Event* (UE) has been a catch-all term, relating to the distributions of particles away from the directions defined by the more interesting hard scatter. The CDF and D0 experiments at the Fermilab Tevatron performed related measurements [8, 9]. The concept of the UE is important to hadron collider physics because it enters the uncertainty analysis for many precision measurements. Every UE analysis is required to precisely define the concept and parameters for itself; the current analysis is no exception. As a result, there are many definitions of the UE in the literature describing the same concept.

We motivate our definition and subsequent measurement of the Underlying Event with

a highly contrived toy example. We use the uncertainty analysis of a *different*, hypothetical measurement to motivate our UE measurement.

Suppose we were measuring the differential cross section of parton scattering as a function of Q^2 in diffractive proton-proton collisions, using an unrealistic detector capable of infinite track momentum and spatial resolution with perfect reconstruction efficiency. This interaction is modelled at leading order (in λ) by a parton (gluon or quark) from one proton scattering with a quark or gluon in the other proton. An example is shown in Fig. 1.1. We refer to these partons as *incoming*. The leading order terms in pQCD contain two *outgoing* partons, meaning two partons are produced in the hard scatter. In Fig. 1.1, a u-quark from one proton scatters against a d-quark from the other proton, each quark radiating a gluon. Those gluons interact, producing a "s-channel" gluon, which then produces a quark-antiquark pair. In our detector, the momenta of the outgoing quark-antiquark pair manifest themselves as a collimated collection of final state hadrons, known as *jets*. This situation is visualized in Fig. 1.3(a). The extent of collimation is characterized by the jet radius parameter (R). Low momentum transfers require a larger radius to capture the hadronic shrapnel from the interaction; high momentum transfers will collimate strongly with less dispersion, thus requiring a smaller radius. For our analysis, we settled on the value $R = 0.6$ as the appropriate jet radius to account for dynamic $4 \text{ GeV} \leq p_{\text{T}}^{\text{jet}} \leq 100 \text{ GeV}$ range of our jets. As an example of the difference between a quark and a jet, an outgoing parton with transverse momentum $p_{\text{T}} = 10 \text{ GeV}$ may hadronize into 5 more more final state particles moving approximately in the same direction, each having approximately $p_{\text{T}} = 2 \text{ GeV}$.

In practice, jets are defined by the algorithm used to construct them. There are various algorithms available such as the anti- k_t [10] and SiSCone [11], etc., each having their advantages and drawbacks. This analysis uses the anti- k_t algorithm to construct *track jets* from charged particles, using only the inner detector (ID). The efficiency and resolution of the tracking detector allows us to probe very low energy jets. The anti- k_t algorithm is an IR- and collinear-safe version of the geometrically-intuitive cone algorithm, which uses tracks inside a cone in $\phi - \eta$ in to characterize its energy. The terms track jets and charged particle jets are used synonymously.

In our example analysis, the most natural strategy would be to reconstruct track jets, compute and histogram the Q^2 . and compare the results to the predictions of Monte Carlo generators. One relevant diagram of a contributing process is $pp \rightarrow pp + jets$, which can proceed through $pp \rightarrow pp + q\bar{q}$ (See Fig. 1.1). For simplicity, let us assume the outgoing quarks each have $p_{\text{T}} = 20 \text{ GeV}$ and the direction of the quark momentum vectors are $\eta = 0$ and $\phi = 0$ for the quark, and $\phi = \pi$ for the antiquark³. This scenario is just a quark pair produced back to back in the lab rest frame. We further assume the outgoing protons travel directly down the beampipe ($\eta = \pm\infty$). The Q^2 of this interaction is 40^2 GeV^2 . We will not discuss the plethora of other leading order diagrams. In our example, detailed numeric information is provided for definiteness.

³The ATLAS experiment uses a right-handed coordinate system. The x -direction \hat{x} points radially inward from the interaction point to the center of the LHC ring, \hat{y} points upward, and \hat{z} points along the beampipe.

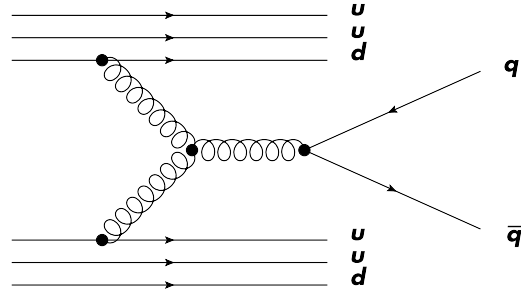


Figure 1.1: An example diagram for $pp \rightarrow pp + jj$

In our example Feynman diagram, we have taken the liberty of conflating quarks and jets. In general, the picture is far from straightforward once we start to probe the hadronization process. Again, we focus on one simplified scenario out of an infinite number of possibilities. The outgoing quarks interact with gluons and spectator quarks from the protons, an example of which can be seen in Fig. 1.2, where the outgoing anti-quark interacts with one of the spectator quarks from the proton via a gluon. The outgoing quark radiates a gluon which (a) splits into a $d\bar{d}$ pair, and (b) pulls a $u\bar{u}$ pair from the vacuum to form a final state π^+ and π^- . The quarks continue through the hadronization process to form a collection of final state particles. The charged particles leave tracks in the Inner Detector from which we reconstruct track jets. In our example, the leading track jet (having the largest p_T) has $p_T = 21$ GeV, $\phi = 0$ and $\eta = 0.1$, close to but not coincidental with the direction of the outgoing quarks from the hard scatter. The subleading track jet (track jet with the next highest p_T) has $p_T = 18$ GeV, $\phi = \pi$ and $\eta = 0.2$. The directions of the π mesons need not follow the direction of the parent quark due to the interactions with the beam remnants. In our example, the π^+ is produced with $p_T = 3$ GeV, $\eta = 0.5$ and $\phi = \pi/2$; the π^- is produced with $p_T = 2$ GeV, $\eta = -0.8$ and $\phi = -\pi/2$. This situation is depicted in Fig. 1.3.

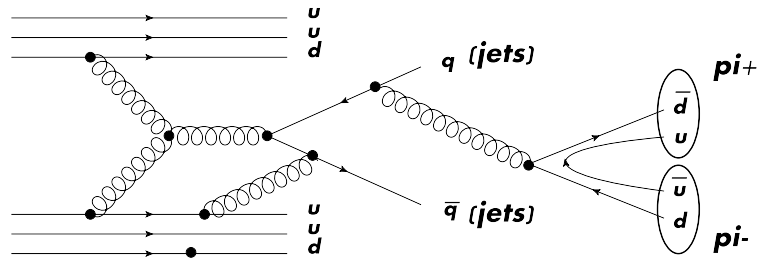


Figure 1.2: A example diagram for $pp \rightarrow pp + jj + \pi^+ \pi^-$

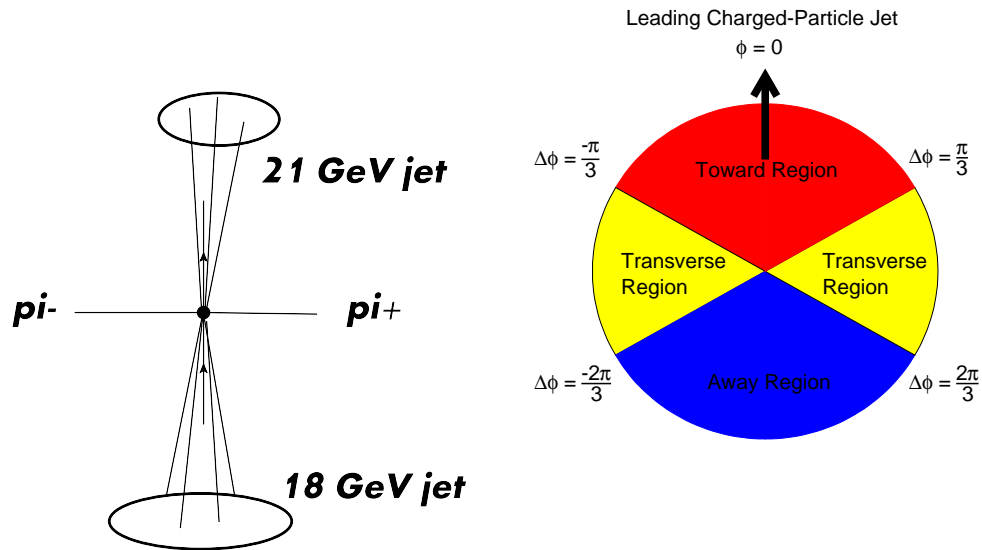


Figure 1.3: An example of particles in the TRANSVERSE region, with jets in the TOWARD and AWAY side.

We now proceed to do something useful with all of this detailed information. This event enters our hypothetical cross section measurement because it contains at least two energetic track jets. The measured Q^2 is 39^2 GeV^2 , according to the sum of the track jet p_T . The question arises as to how we account for the pions that were produced away from the direction of the hard scatter. Their contribution is neglected because the reconstructed jets did not account for their energy. If the Monte Carlo generator models the pion pair production as ultimately originating from hard scatter partons, then the pions are accounted for by jet energy resolution. If MC produces the pions from a secondary interaction, not tracing its “ancestrage” to the hard scatter, we must understand how the MC models these effects. The question we are asking is, “How are meaningful comparisons to Monte Carlo generator predictions to be made?”

MC generators tend to simplify the hadronization process in such a way that the final state particles tend to follow the direction of the hard scatter. The population of phase space away from this direction needs a correction in order to match observed data. One current model that performs this correction is known as multiple parton interactions (MPI) [12], where a secondary, softer interaction occurs generating particles much more isotropically, thus filling in phase space. The MPI will also produce particles overlapping in phase space with the hard scatter, contributing particles to the jet and overestimating its p_T . The MC produces the hard scatter with a cross section calculated unaware of this secondary interaction⁴. MC must account for the additional energy in order to remain consistent with the

⁴Some newer models attempt to reconcile QCD color differences between the hard scatter, beam remnants

data.

We are now ready to motivate the definition of the *Underlying Event* (UE) for track jets. The leading track jet is used to define a direction in phase space. The track jet can be described by four variables (p_T , η , ϕ , mass). All variables except η are relativistically invariant under boosts along the z axis. We use ϕ^{jet} to define the direction ϕ_0 , and categorize the jet using its p_T . Jet mass does not enter the analysis. We define the TRANSVERSE region to be the area of phase space $\pi/3 \leq |\phi_0 - \phi| \leq 2\pi/3$.

We tentatively define the Underlying Event as the charged particle activity in the TRANSVERSE region, using the multiplicity and scalar sum p_T of the tracks as relativistically invariant measures of particle activity. The caveat is that since our tracker offers excellent tracking to $|\eta| \leq 2.5$, we impose a selection criteria on the tracks. We settle on the final definition of the Underlying Event, as it pertains to this analysis:

The Underlying Event is the charged particle activity in the TRANSVERSE region ($\frac{\pi}{3} \leq |\phi - \phi_0| \leq \frac{2\pi}{3}$), for tracks having $|\eta| \leq 1.5$ and $p_T \geq 500$ MeV. Charged particle jets are constructed using the anti- k_t algorithm, with the clustering radius R -parameter value fixed at 0.6, from tracks having $|\eta| \leq 2.5$ and $p_T \geq 500$ MeV. The leading jet with $|\eta| \leq 1.5$ and $p_T \geq 1$ GeV is used to define ϕ_0 on a per event basis. The charged particle activity is characterized by the scalar sum of the individual track p_T (Σp_T), the number of tracks (N_{ch}), and the average p_T per track ($\bar{p}_T \equiv \Sigma p_T / N_{ch}$).

We now summarize and discuss the Underlying Event. Experimentally, this definition of the UE is an unambiguous and a well-defined quantity. Theoretically, one encounters problems in the interpretation, typically arising from too literal an interpretation of Feynman diagrams. Feynman diagrams are invaluable tools that help calculate cross sections, but we must be cautious and remember that any diagram is one of an infinite number of diagrams that must be consistently summed, *before* interpretation. One popular working definition is that the UE is "everything" in the event except the hard scatter, which already presupposes we can unambiguously classify a hadron as having originated from the hard scatter.

The different available MC use different models to populate areas of phase space complementary to those defined by the hard scatter. Such models as Initial-State Radiation (ISR) and Final-State Radiation (FSR), MPI, beam remnant interactions, bremsstrahlung, etc. all work differently and complement the hard scatter model. In the end, the goal is to describe the physics correctly. We refrain here from describing the various MC models; it is not germane to the current analysis and detailed information can be found in the literature. We focus on characterizing their performance in reproducing the relevant physics distributions, instead of analyzing the success or failure modes. The results of this analysis will provide the authors of MC generators another test of their models, by providing truth-level distributions which can be compared without reference to any detector.

A similar measurement was performed simultaneously by the CMS experiment [13]. The CMS measurement of the track jet-based UE used charged particle jets reconstructed

and secondary interactions. This is known as *color reconnection*.

using the SISCone algorithm [11] with $R=0.5$. The CMS analysis used the same $p_T \geq 0.5$ GeV acceptance for tracks, but the $|\eta| \leq 2$ was different than the acceptance in this analysis. ATLAS also performed a similar measurement [14] to ours, except the direction of the UE was determined by the track with the largest p_T .

This analysis measures three different quantities simultaneously, most often with identical methods for each. We will usually outline the method for one of the observables, noting the similar or identical approach for the remaining observables. We mention differences, if any, in the treatment of the different observables. Alternately, we refer to the measured quantities generically as \mathcal{O} . For example, the distributions of \mathcal{O} vs p_T^{jet} refer to each of Σp_T , N_{ch} and \bar{p}_T vs p_T^{jet} .

UE analyses have been performed at different experiments. Although they may use different objects (e.g. - leading track, Z^0) to define the directions, these analyses consistently label the different UE regions. The phase space selection criteria may vary, but the concepts of TOWARD, TRANSVERSE and AWAY regions are used consistently. The different regions of the UE, as defined for the current analysis, are shown in Fig. 1.3.

Chapter 2

Monte Carlo

2.1 Generation

Monte Carlo generators are algorithms, implemented as computer programs, employing phenomenological models to simulate physics processes. The output of these generators is typically a list of partons or final state particles, and their properties, whose origin reflects the physics process being modelled. For example, if we are modelling a top quark decaying to a W -boson and a b -quark, the list of final state particles would usually include a B -meson (containing the b -quark). In high energy physics, the most ubiquitous generator is PYTHIA. Examples of other generators include HERWIG++, SHERPA, AcerMC [15], ALPGEN [16], etc.

Monte Carlo generators are extremely important in a wide class of analyses, because it helps us relate what we observe in our detectors to the fundamental physics processes. The distribution of particles in our (imperfect) detectors does not uniquely point to the responsible physics process. It must usually be inferred by examining its consistency with different scenarios. As we search for potential new physics in our experiments, we have to ask many "what if?" type questions. For example, let us assume we are looking for a spin-0 particle with a mass of 120 GeV. We would model the relevant physics process and compare the output to the observed data. What if the particle had a mass of 115 GeV? or 125 GeV? How would the output of our simulation change, and could our analyses resolve the differences? What if the particle were a vector boson (force mediator described by a field with vectorial transformation properties) instead? The distributions of the decay products differ depending on whether the particle is described by a scalar or vector field. By modelling the different scenarios, comparing and analyzing the data, we determine the ability of our analysis to resolve new physics from the predictions of the Standard Model. Although we used a sleek example to motivate the importance of Monte Carlo generators, the same line of reasoning applies to precision measurements of the Standard Model.

In the context of this analysis, we ask what the distributions of particles are, how many are there, and their energy content. This is one of the first UE measurements performed at the LHC made with early 2010 data. Physicists fed back knowledge from previous experi-

ments into existing MC generators and made an educated guess at the physics distributions at the LHC. The generators are configured to give distributions that agree with data from previous experiments. The output of these generators are then to be compared to new data.

Monte Carlo generators usually have a set of parameters that can be adjusted to modify its behavior, to make it agree better with experimental data. A particular configuration of parameters is known as a *tune*. When referring to MC with a particular configuration, we use the convention of specifying the tune in parentheses after the name of the generator. For example, PYTHIA 6 (AMBT1) refers to the AMBT1 tune of PYTHIA6.

This analysis uses PYTHIA 6 (MC09) [17] to validate the analysis techniques, and in the evaluation of systematic uncertainties. The available sample had 20M events, approximately half the size of the data sample. Other MC were used in the evaluation of systematic uncertainties. The exhaustive list of Monte Carlo generators considered in this analysis is

- MC09 tune of PYTHIA 6 [17]
- AMBT1 tune of PYTHIA 6 [18]
- UE7-2 tune of HERWIG++ [7]
- Perugia2010 tune of PYTHIA6 [19]
- Perugia2011 tune of PYTHIA6 [19]
- Perugia2011 (without color reconnection) tune of PYTHIA6 [19]
- 4C tune of PYTHIA8 [20]
- Z1 tune of PYTHIA6 [21]
- AUET2B tune of PYTHIA6 [22]

2.2 Simulation

We described Monte Carlo generators in Sec. 2.1. We now discuss simulation of the detector. Any measurement requires knowledge of how the apparatus responds to input stimuli.

The output of a Monte Carlo generator is usually a list of particles and properties that we might observe with a perfect detector. In order to compare the predictions of MC to the data, we need a more realistic description of its output. By simulating the detector response, and applying it to the MC output, we obtain a modified list of particles with which we can make a meaningful comparison to the data.

This analysis makes extensive use of PYTHIA 6 (AMBT1) [18] that has undergone full detector simulation using the GEANT4 package [23]. As will be discussed in Sec. 5, this tune is used in the correction procedure to account for detector effects. At the time this

analysis was performed, this Monte Carlo sample had the best available statistics. The AMBT1 configuration was not tuned to LHC data.

We require high statistics in all of the relevant phase space, in order to obtain an accurate description of the detector. Unfortunately, modelling the detector is *very* time-consuming, roughly 15 minutes per event. This means we cannot afford to fully simulate every generated event; we must carefully choose the events we wish to simulate. One of the variables in our analysis is the transverse momentum of charged particle jet (p_T^{jet}). The cross section for generating MC samples drops rapidly as a function of p_T^{jet} (See Fig.3.7). Without making specific cuts during the generation process, obtaining full-simulation samples at high p_T would be a very inefficient process. By making specific cuts on the transverse momentum of truth jets, we are able to efficiently populate all relevant phase space. Samples generated with such cuts are referred to as *slices*, and the events in these slices are properly weighted to form consistent distributions when histogrammed.

The details of GEANT4-simulated MC sample generation are tabulated in Table 2.1. PYTHIA 6 (AMBT1) is listed twice; there are two statistically independent samples.

2.3 Remarks

Because PYTHIA 6 (AMBT1) plays a major role in the derivation of the final measurement, we compare most of its physics distributions to the data, before application of any corrections. For some variables, the differences in the data and MC distributions are fairly large. Binning in these variables removes the leading order effects of such differences, and we are still able to use the Monte Carlo to obtain good results. We account for residual differences as a source of systematic uncertainty in Sec. 6.

We use the terms *baseline measurement* and *central value* to refer to the final corrected measurements.

We compare the central values to the predictions of PYTHIA 6 (Z1) and PYTHIA 6 (AUET2B), which were tuned to CMS and ATLAS data, respectively. These are the best configurations CMS and ATLAS had at the time our measurements were completed. These tunes were not used to obtain the central values, but the comparisons are interesting because these Monte Carlo were tuned using LHC results from other measurements.

Table 2.1: GEANT4-simulated Monte Carlo generation

GEANT4 simulation samples		
generator	generator cut	# of events
Pythia 6 (MC09)	-	19,693,365
Pythia 6 (AMBT1)	$4\text{GeV} \leq p_T^{\text{jet}} \leq 15\text{GeV}$	19,823,155
	$15\text{GeV} \leq p_T^{\text{jet}} \leq 30\text{GeV}$	19,660,690
	$30\text{GeV} \leq p_T^{\text{jet}} \leq 60\text{GeV}$	19,618,890
	$60\text{GeV} \leq p_T^{\text{jet}}$	9,745,804
Pythia 6 (AMBT1)	-	4,907,480
Pythia 6 (Perugia2010)	$4\text{GeV} \leq p_T^{\text{jet}} \leq 15\text{GeV}$	2,477,628
	$15\text{GeV} \leq p_T^{\text{jet}} \leq 30\text{GeV}$	2,445,198
	$30\text{GeV} \leq p_T^{\text{jet}} \leq 60\text{GeV}$	2,424,625
	$60\text{GeV} \leq p_T^{\text{jet}}$	1,224,549
Pythia 8.145 (4C)	-	4,004,064
Data	-	42,617,085

Chapter 3

The ATLAS Inner Detector at the Large Hadron Collider

3.1 The Large Hadron Collider at CERN

Located in the environs of Geneva, Switzerland, the Large Hadron Collider (LHC) [1] collides proton beams in opposing directions. Designed to operate at 14 GeV center-of-mass energy with a luminosity $\mathcal{L} = 10^{34} \text{cm}^{-2} \text{s}^{-1}$, the LHC was operating at 7 GeV center-of-mass energy and peak luminosity $\mathcal{L} = 6.6 \times 10^{28} \text{cm}^{-2} \text{s}^{-1}$ in early 2010, when our measurement was made. The LHC is 27 kilometers in circumference, and located up to 175 meters underground.

The beams are produced as hydrogen ions and a chain of accelerators successively boost these protons to increasing energies. A linear accelerator (LINAC2) brings the protons to 50 MeV and feeds the (PSB) Proton Synchrotron Booster. The PSB boosts the beam to 1.4 GeV and feeds the SP (Proton Synchrotron). The SP boosts the proton beam to 25 GeV and feeds the Super Proton Synchrotron (SPS). The SPS boosts the protons to 450 GeV, which then feeds the last accelerator in the chain, the LHC. After proton bunches in each opposing beam are accelerated to 3.5 GeV, they are collided. The collisions occur at 4 locations, referred to by the name of the experiment located at the points - ATLAS, CMS, LHCb and ALICE. The LHC uses a circular array of more than 1600 superconducting magnets keep the protons in their trajectory. The CERN accelerator complex is shown in Fig. 3.1.

3.2 The ATLAS detector

The ATLAS detector is actually an ensemble of many detectors, forming a general-purpose hermetic detector. The hadronic calorimetry system measures the energy deposited by strongly-interacting particles. The electromagnetic calorimetry measures the energy deposited by electrons and photons. Both calorimeters have fine segmentation (*granularity*)

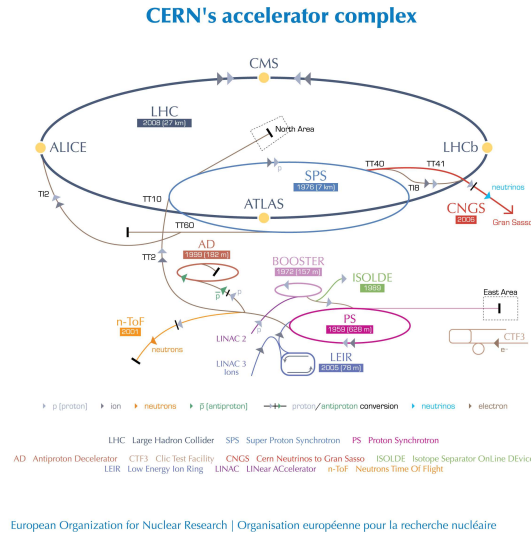


Figure 3.1: The CERN accelerator complex [24]

allowing for good spatial resolution and longitudinal energy deposit resolution (dE/dx). As the outermost detector, the muon system tracks electrically charged particles that "punch through" all the other detectors, including the calorimetry. We identify these minimum ionizing particles (MIP) as muons. Neutrinos do not interact with the detector; their signature is missing energy. The innermost detector, the ID (inner detector) resolves charged particle trajectories and is described in more detail next, as this is the relevant subsystem for our analysis. The ATLAS detector is depicted in Fig. 3.2.

3.3 The ATLAS Inner Detector

The ATLAS detector is a system of complex and complementary subdetectors. This analysis, however, makes use of a very small part of the whole detector. The entire analysis is based on *tracks*, which only require the inner detector (ID). We focus our discussion on this relevant subdetector. Detailed information about the different components of the ATLAS detector can be found in the literature. [25]

The ATLAS Inner Detector (ID) consists of 3 separate and complementary subdetectors. These are the PIXEL detector [26], the SemiConductor Tracker (SCT) [27] and the Transition Radiation Tracker (TRT) [28, 29], that comprise the innermost component of the ATLAS detector. The envelope of the ID is located just outside the beampipe, extending 1.2m radially and $-3.5\text{m} \leq z \leq 3.5\text{m}$. Geometric details of the ID layout can be seen in Fig. 3.3. The entire ID is immersed in a 2T axial magnetic field, causing charged particles to trace out curved trajectories according to their momentum. The ID is responsible for finding charged particles and resolving their momenta.

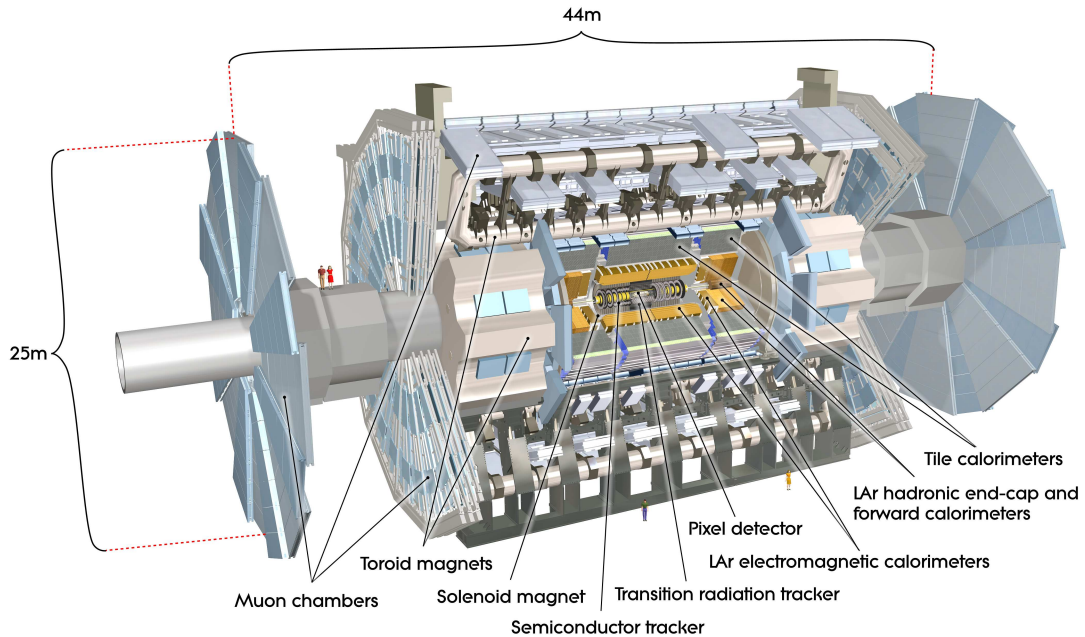


Figure 3.2: A diagram of the ATLAS detector indicating its sub-components [25].

3.3.1 The PIXEL detector

The PIXEL detector is an array of approximately 80 million $50\mu\text{m} \times 400\mu\text{m}$ silicon-based charge detectors, known as pixels. The pixels detect charged particles traversing their geometry, providing 3-dimensional spacepoint information about the track. The geometrical layout of the pixel detector was designed to provide ≥ 3 spacepoints (See Fig. 3.5) for tracks with $|\eta| \leq 2.5$. The pixels are arranged in 3 cylinders concentric to the beampipe, and 6 parallel disks (3 on each side). Fig. 3.3 shows geometric detail about the PIXEL detector. The 3 cylinders, collectively referred to as the pixel barrel, are located at $r = 50.5\text{mm}$ (B-layer), 88.5mm (Layer 2) and 122.5mm (Layer 3) from the beamline, and span $|z| \leq 400.5\text{mm}$. On either side of the barrel, 3 disks are arranged at $|z| = 495\text{mm}$, 580mm and 650mm , and span $88.8\text{mm} \leq r \leq 149.6\text{mm}$. The pixels disks are commonly referred to as the end caps. The pixels are arranged into *modules* with 46080 pixels each. The modules are tiled onto the disk surfaces in the end caps, and onto long carbon fiber strips (staves) in the barrel. To provide full quality coverage in ϕ , the PIXEL detector provides some module overlap between modules in the azimuthal direction. As a result, a track may have more than one hit in on the same layer, if it crosses the the region of overlap. Detailed information about the pixels can be found in [26]. Fig. 3.4 shows the rendered image of the PIXEL detector.

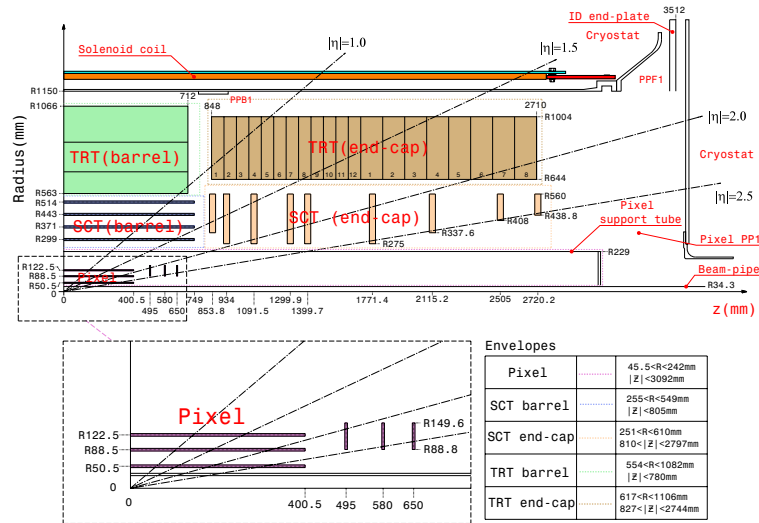


Figure 3.3: An $r - z$ view of the ATLAS Inner Detector with details of the geometric layout. Reproduced from [25].

3.3.2 The Semiconductor Tracker

Located just outside the PIXEL detector, the SCT detector is also an array of silicon-based charge detectors. The $80\mu\text{m} \times 1260\text{mm}$ geometry of the SCT detector elements is much larger length-wise than the pixels, and are referred to as strips. The SCT has a barrel and endcaps, one on each side of the barrel along the beampipe. In the barrel region, there are 4 double-sided layers, with the strips arranged stereographically. The layout in the endcap is more intricate; Fig. 3.3 shows the detail of the geometrical layout of the strips, which cover tracks with $|\eta| \leq 2.5$. The strips do not offer the same resolution in z (1260mm vs $400\mu\text{m}$) as the pixels, but compensate by providing many more channels and covering more surface area. Unless lost to efficiency, or interaction with material in the SCT infrastructure, tracks normally register 2 hits per layer. Detailed information about the SCT can be found in [27]. Fig. 3.4 shows the rendered image of the SCT detector.

3.3.3 The Transition Radiation Tracker

The outermost component of the ID, located outside the SCT, the Transition Radiation Tracker (TRT) is an array of straw drift tubes. The TRT provides continuous tracking information between the SCT and the outer envelope of the TRT (approximately 1m). In the barrel, $4\text{mm} \times 370\text{mm}$ straws, running axially to the beampipe, are filled with $\text{Xe}/\text{CO}_2/\text{O}_2$. In the endcap, the $4\text{mm} \times 1440\text{mm}$ straws run radially. A fine tungsten wire anode at the center of the straw is held at ground; the tube (cathode) is held at -1.5kV . Charged particles traversing the straw tube ionize the gas; the potential difference between cathode and anode

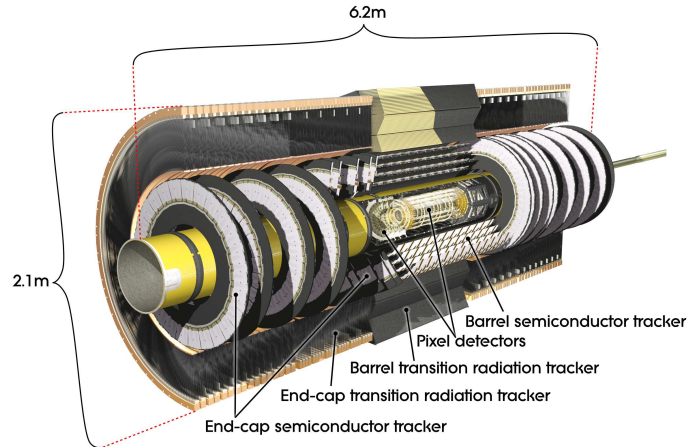


Figure 3.4: An $r - z$ view of the ATLAS Inner Detector. Reproduced from [25].

causes electrons to drift to the anode. Drift time provides further spatial resolution. A particle traversing the TRT barrel ($|\eta| \leq 0.8$) is expected to have 36 TRT hits. Geometric considerations must be taken into account to determine the number of expected hits in the TRT end cap ($1 \leq |\eta| \leq 1.9$), and in the transition region between the barrel and endcap. Polypropylene fibers (in the barrel) and foils (in the endcap) between the straws act as a transition radiation generator to help discriminate between electrons and pions. Different voltage thresholds are used to detect the difference between transition radiation from electrons and minimum ionization from pions. Detailed information about the TRT can be found in [28, 29]. Fig. 3.4 shows the rendered image of the TRT detector.

3.4 ATLAS Trigger Overview

The ATLAS detector has a 3 level trigger system, known as the Level 1 (L1), Level 2 (L2) and Event Filter (EF).

The L1 trigger system is a sum of hardware triggers from different components of the detector. Detector components are designed to trigger within $2.5 \mu\text{seconds}$ of a significant signal in that subdetector. Due to various constraints, including full detector readout time, this trigger can operate at a maximum rate of 75 kHz.

The L2 software trigger is performed outside the detector, using dedicated computers to perform optimized reconstruction in the various regions of interest. The L2 trigger system reduces the event rate down to approximately 3.5 kHz.

The last component in the trigger chain is the Event Filter (EF). The decision to commit the event to permanent storage is based on full event reconstruction. The EF passes approximately 200 events per second.

This analysis uses events selected with the L1 Minimum Bias trigger.

3.4.1 Minimum Bias Trigger

The different components of the ATLAS detector (calorimeters, muons, etc.) are each capable of firing the L1 trigger in response to signal detection. The minimum bias trigger consists of Beam Pickup Timing (BPTX) devices and the Minimum Bias Trigger Scintillator (MBTS).

The BPTX Trigger

Formally considered part of the LHC machine, even though they are operated by ATLAS, there are two BPTX stations on either side of the ATLAS detector, located at $\pm 175\text{m}$ from the nominal interaction point. Each BPTX station consists of 4 electrostatic button pickup devices, arranged symmetrically around and attached to the beam pipe. The BPTX devices pick up the signal from passing proton bunches. A coincidental trigger from both sides of the detector indicates that two proton bunches have collided. Detailed description and performance of the BPTX can be found in [30]

The MBTS Trigger

On either side of the ATLAS detector, the MBTS system consists of a disk (its face perpendicular to the beamline), with scintillator counters mounted on two radial rings. Each ring is divided into 8 equal segments in ϕ , for a total of 16 segments on each side of the ATLAS detector. The two rings span $2.09 \leq |\eta| \leq 3.84$, and are located at $\pm 3.56\text{ m}$ from the nominal interaction point. Particles traversing any segment deposit energy into the scintillator, and the light is guided to a photomultiplier tube (PMT). After signal shaping, a hit is defined as a signal over the discriminator threshold.

The MBTS trigger efficiency is ≥ 0.97 for events with 2 selected tracks, rising to > 0.99 for ≥ 3 tracks. The trigger, by construction, does not introduce a significant selection bias, and the efficiency does not affect the measurements in our analysis. Our event selection is based on forming charged particle jets, naturally selecting events with larger track multiplicities, and therefore the MBTS efficiency is essentially 100% [31].

Events in this analysis are selected with the Minimum Bias trigger, with at least one MBTS hit and a coincidence in both sides of the BPTX.

3.5 Event Reconstruction

3.5.1 Track Reconstruction

Because the Inner Detector (ID) is immersed in a 2T magnetic field, a charged particle will trace out a helical trajectory, the parameters of which depend on the particle's momentum and production vertex. As the charged particle traverses the various components of the ID, the registered hits are recorded and used to reconstruct its trajectory. This

trajectory is also referred to as a *track*. We summarize the salient features of the reconstruction algorithms, those which are relevant to this analysis. Details about the reconstruction algorithms and their performance can be found in [32].

A helix can be described by the following five parameters, assuming knowledge of the event's primary vertex (PV), the point where the protons collided. The particle's production vertex may differ from the PV.

- p_T - the transverse momentum of the particle. We measure the track curvature ρ and use it to determine the particle's p_T using its charge q and the magnetic field strength B , via the relation $p_T = qB/\rho$.
- $\eta = -\log\left(\tan\left(\frac{\theta}{2}\right)\right)$, the pseudorapidity of the particle's production vertex
- d_0 - the transverse impact parameter, is the distance of closest approach, in the $r - \phi$ plane, to the PV
- z_0 - the longitudinal impact parameter, is the z coordinate of the closest point to the PV
- ϕ_0 - the azimuthal coordinate of the point of closest approach to the PV

Hits from adjacent pixels or strips are gathered into clusters (contiguous combinations of hits), which are used as seeds in the tracking algorithm. All combinations of three clusters (from any pixel layer or the innermost SCT strip) are used to define a *road*, which is essentially is a track candidate. Hits from the ID are then associated with the track, which is refitted after every hit association using a Kalman filter [33] and a simplified model of the detector geometry. After the hit association is complete, the track undergoes further quality checks. The track is refit using a more sophisticated description of the detector, and scored accounting for the quality of the fit, the number of hits, the number of holes ("missing" hits) and the χ^2/N_{dof} is used to select good tracks. Hits which are shared between tracks are reassigned to the highest quality track. Tracks are subsequently extrapolated to the TRT, and the analogous procedure of attaching TRT hits is repeated.

3.5.2 Vertex Reconstruction

Primary vertex reconstruction begins after track reconstruction is complete, requiring at least two tracks having

- $p_T^{\text{track}} \geq 100 \text{ MeV}$
- $|d_0^{\text{BS}}| \leq 4 \text{ mm}$, where d_0^{BS} is the transverse distance of closest approach to the beamspot
- the uncertainty on d_0^{BS} , $\sigma(d_0^{\text{BS}}) \leq 5 \text{ mm}$
- the uncertainty on z_0^{BS} , $\sigma(z_0^{\text{BS}}) \leq 10 \text{ mm}$

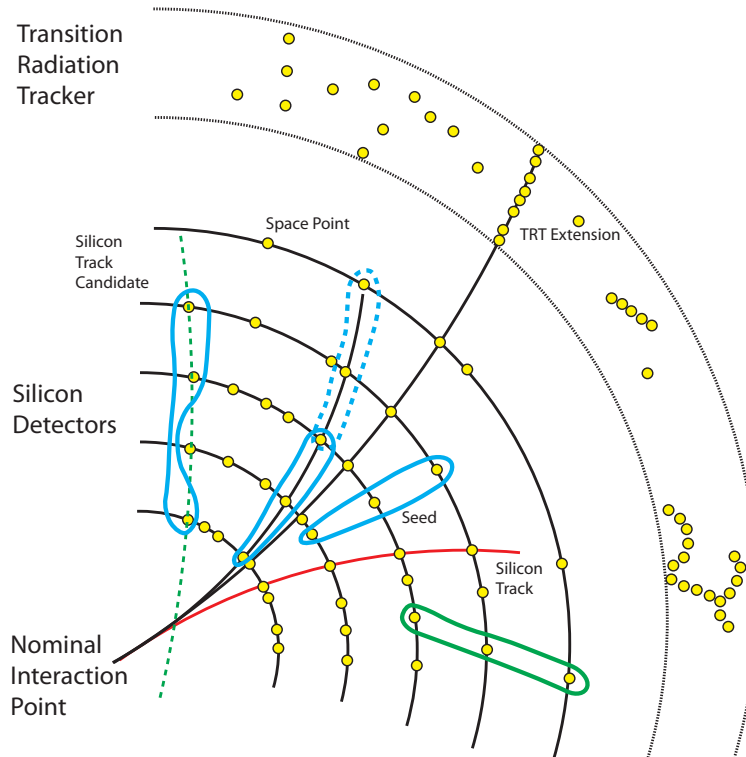


Figure 3.5: Schematic of the different stages of track reconstruction. Spacepoints are represented by the yellow dots and blue lines indicate track seeds. The dashed blue line is a seed which shares a single hit and the green line illustrates a seed which was rejected prior to hit association. The green dashed line indicates a track candidate which failed the impact parameter cuts. The red line represents a silicon-only track. The black line indicates a track including TRT hits. Figure reproduced from [34] with permission of the author.

- at least one pixel hit ($N_{\text{pix}} \geq 1$)
- at least 4 SCT hits ($N_{\text{sct}} \geq 4$)
- at least 6 silicon hits total ($N_{\text{pix}} + N_{\text{sct}} \geq 6$)

The vertex fitter is seeded with the maximum of the z_0 distribution of the tracks. Tracks are tested for consistency with the candidate vertex. The adaptive vertex fitter [35] uses a χ^2 -based algorithm to iteratively reduce the contribution from outlying tracks, which can become candidates for another vertex. The algorithm is complete when the track collection is exhausted or no further vertices are found. If the beamspot is known, it is also used to constrain the fit. Vertex reconstruction is described in detail in [36].

3.5.3 Jet Reconstruction

Selected tracks with $|\eta| \leq 2.5$ (c.f. Sec. 4.1.1) are clustered into charged particle jets (*track jets*) using the anti- k_t algorithm, using $R = 0.6$ for the clustering radius R-parameter. The R-parameter is often referred to as the *jet radius*. Charged truth jets are formed from Monte Carlo, applying the anti- k_t algorithm to primary particles in the event HepMC collection. Fig. 3.6 shows the reconstruction efficiency for charged truth jets that have been matched to a charged particle jet, with a matching criteria $\Delta R \equiv \sqrt{(\Delta\eta)^2 + (\Delta\phi)^2} \leq 0.3$. Fig. 3.7 shows the p_T spectrum of reconstructed jets, comparing data to PYTHIA 6 (AMBT1). The MC p_T spectrum is *harder* (having larger p_T) than the data.

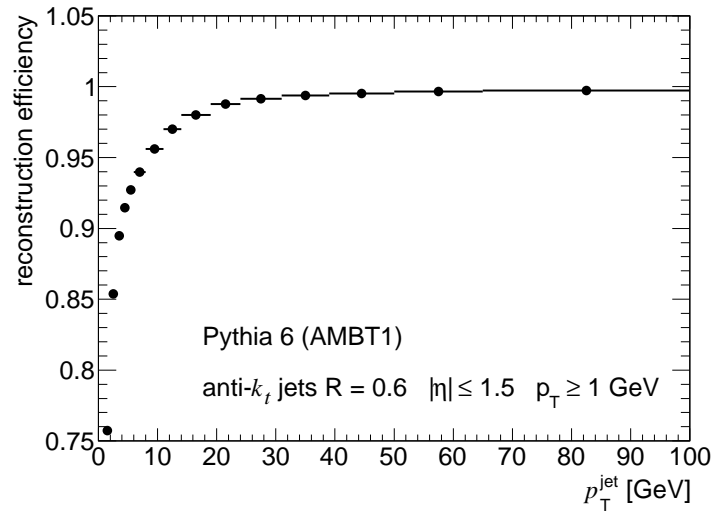


Figure 3.6: The reconstruction efficiency for charged truth jets. Reconstructed charged particle jets have $p_T^{\text{jet}} \geq 1$ GeV.

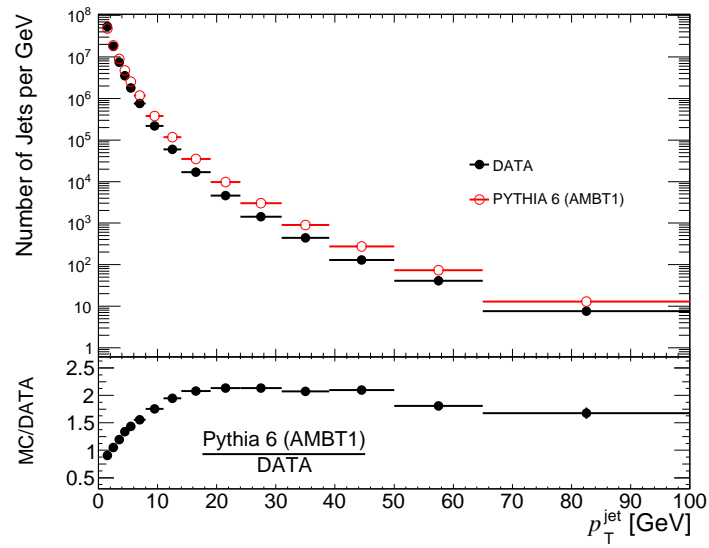


Figure 3.7: The measured charged particle jet p_T spectrum compared to PYTHIA 6 (AMBT1)

Chapter 4

Measurement of Raw UE Distributions

We describe the measurement of the Underlying Event distributions directly from the data, before any corrections are applied.

4.1 Event and Data Selection

The data used in this analysis were taken early 2010, accumulating a total integrated luminosity of $800 \mu\text{b}^{-1}$ (after highly prescaled triggers). More than half the data were taken with $\mu \leq 0.01$, where μ is the average number of collisions per bunch crossing, and never exceeded $\mu = 0.14$ throughout the relevant data-taking period. Therefore, the effects of *pile-up* (more than one collision per bunch crossing) are minimal. The relevant triggers (Sec. 3.4.1) and detectors (Sec. 3.3) were fully functional. An event was selected for analysis if it had exactly one primary vertex (PV) and a charged particle jet with $p_{\text{T}}^{\text{jet}} > 1\text{GeV}$ and $|\eta^{\text{jet}}| \leq 1.5$.

4.1.1 Track Selection

In this section, we outline the criteria for selecting primary tracks. Primary tracks have been studied comprehensively in Minimum Bias studies in ATLAS. We adopt the same selection criteria as the ATLAS Minimum Bias analysis [37]. A complementary UE measurement based on the leading track to define the direction of the UE [14] also makes use of these selection criteria. Using the same track selection criteria allows us to make comparisons, and use the same infrastructure for tracking efficiency uncertainty analysis. The same tracks are used to make track jets and calculate the UE observables, with one important exception. Whereas the tracks used in jet reconstruction are allowed to have $|\eta| \leq 2.5$, *analysis tracks* (those used for calculating the UE observables), are restricted to $|\eta| \leq 1.5$.

Some history is required to explain the choice of η acceptance used in our measurement. This analysis considers anti- k_r track jets with a jet radius $R = 0.6$. We performed

variable	cut
p_T	0.5 GeV
SCT hits	≥ 6
Pixel hits including B-Layer hit if expected	≥ 1
$ z_0 \sin \theta $	$\leq 1.5\text{mm}$
$ d_0 $	$\leq 1.5\text{mm}$
total tracks with $ \eta \leq 2.5$	404,137,798
total tracks with $ \eta \leq 1.5$	259,643,051

Table 4.1: Track selection criteria

companion analyses [38, 39], considering jet radii ranging from $0.2 \leq R \leq 1.0$, using the same analysis techniques outlined in this thesis. To obtain the best p_T^{jet} resolution, all tracks should be within the acceptance region $|\eta^{\text{trk}}| \leq 2.5$. Restricting jets to $|\eta^{\text{jet}}| \leq 1.5$ ensures that all constituent tracks are within the track acceptance. Analysis tracks, however, should remain within the jet acceptance to avoid cases where $p_T^{\text{track}} \geq p_T^{\text{jet}}$, and suppress the effects of mismeasured high p_T tracks.

The track selection criteria are listed in Table 4.1. Tracks must register at least six (6) hits in the SCT, and at least one (1) hit in the pixel detector. To reduce the contamination of secondary tracks, if the corresponding module in the B-layer in the pixel detector is operational, the track must register a hit in the B-layer, automatically satisfying the ≥ 1 pixel hit requirement. To select tracks from the primary vertex, tracks must have a transverse impact parameter $|d_0| \leq 1.5\text{mm}$ and longitudinal impact parameter $|z_0 \sin(\theta)| \leq 1.5\text{mm}$.

A total of 404,137,798 tracks passed the selection criteria for track jet construction. A total of 259,643,051 tracks passed the selection criteria for calculation of the UE observables. The distributions of the variables used to select tracks are shown in Appendix A. Fig. [4.1] shows the p_T and η distributions of the tracks used to construct charged particle jets. Fig. [4.2] shows the p_T and η distributions of the tracks that enter the calculation of the UE observables.

4.1.2 Jet Selection

Charged particle jets are accepted if they have $p_T^{\text{jet}} \geq 1\text{ GeV}$ and $|\eta| \leq 1.5$. This analysis reports results for $p_T^{\text{jet}} \geq 4\text{ GeV}$. The expanded acceptance is used for the purposes of analyzing the uncertainties associated with jets "smearing" in from outside the acceptance.

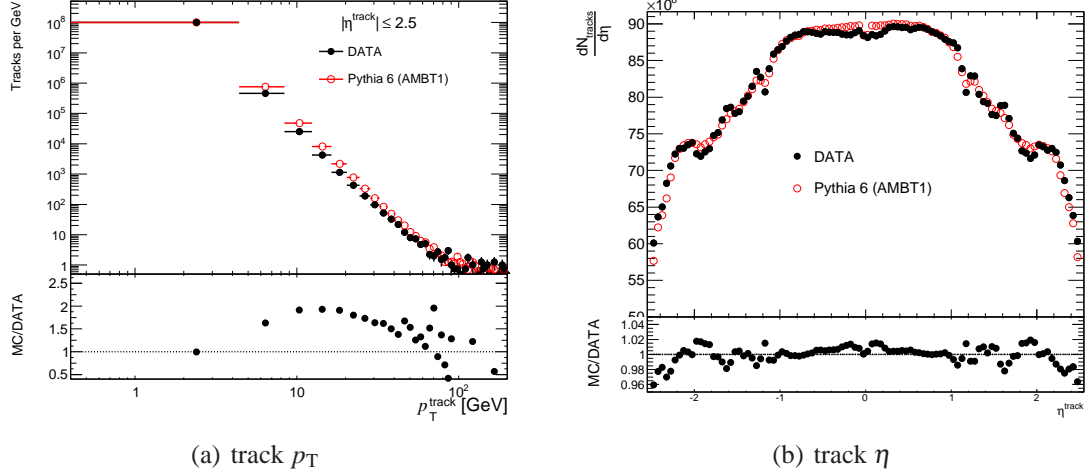


Figure 4.1: The p_T and η distributions for tracks used to construct charged particle jets. The Monte Carlo is normalized to the data.

4.2 Measuring the Distributions from Data

Tracks in selected events are clustered into jets using the anti- k_r algorithm, using $R=0.6$ as the clustering radius R -parameter. For each event, the leading jet with $|\eta| \leq 1.5$ is selected to define the $\phi = \phi_0$ direction. Tracks in the event are selected for the UE calculations if the relative azimuth to the leading jet satisfies $\frac{\pi}{3} \leq |\phi - \phi_0| \leq \frac{2\pi}{3}$. For each event, p_T^{ext} is defined as the transverse momentum of the hardest jet satisfying $1.5 \leq |\eta| \leq 2.5$. If there are no jets satisfying this condition, $p_T^{ext} \equiv 0$. For each observable \mathcal{O} , we record p_T^{jet} , p_T^{ext} and \mathcal{O} in a 3-dimensional histogram.

The UE observables are calculated for each event as follows:

- N_{ch} = the number of tracks in the TRANSVERSE region.
- $\Sigma p_T \equiv \sum_{k=1}^{N_{ch}} p_{T,k}^{track}$ = scalar sum of the track p_T
- \bar{p}_T = the average p_T per track $\equiv \frac{\Sigma p_T}{N_{ch}}$

Figs. [4.3-4.5] show slices in each observable \mathcal{O} , holding p_T^{jet} fixed and integrating over p_T^{ext} . As we will discuss later, p_T^{ext} enters the correction procedure when we adjust for detector effects; its measurement is important for determining the effects of jets that *smear* in from outside the acceptance. PYTHIA 6 (AMBT1) is compared to the measured data and is *not* in good agreement, potentially leading to biases. These differences are accounted for in Sec. 6, when we discuss systematic uncertainties.

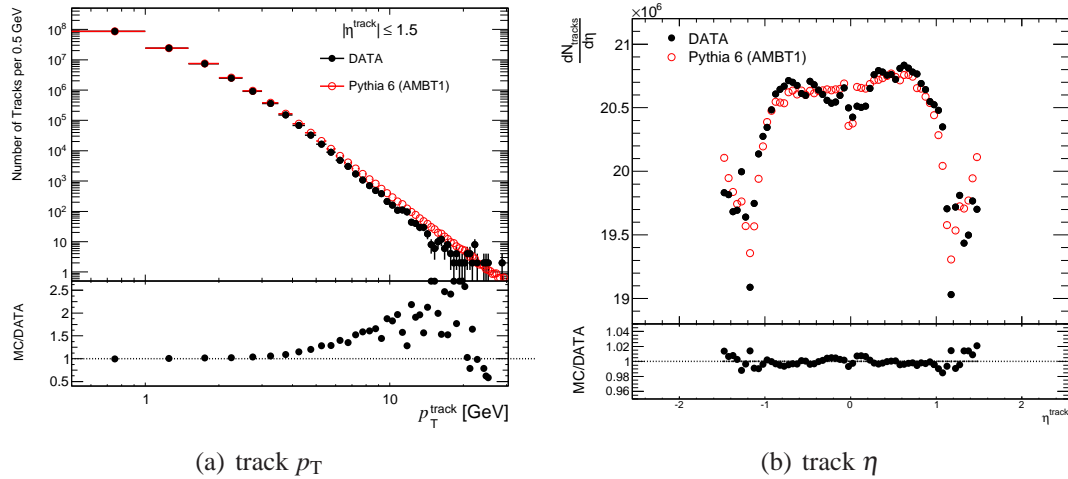


Figure 4.2: The p_T and η distributions for tracks used to calculate the UE observables. The Monte Carlo is normalized to the data.

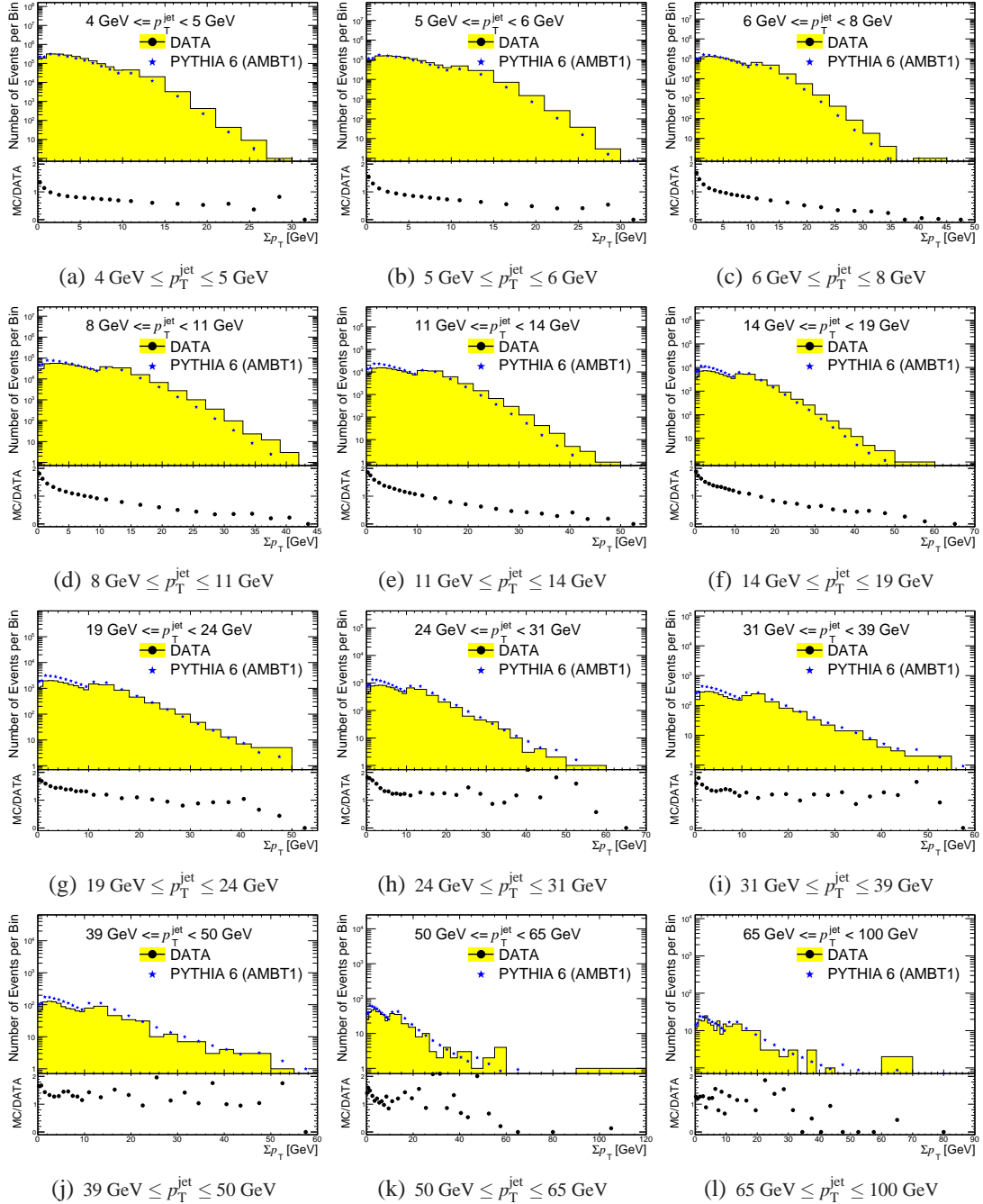


Figure 4.3: The measured Σp_T distributions before any corrections are applied. The data are compared to PYTHIA 6 (AMBT1). The ratio is shown in the bottom plot.

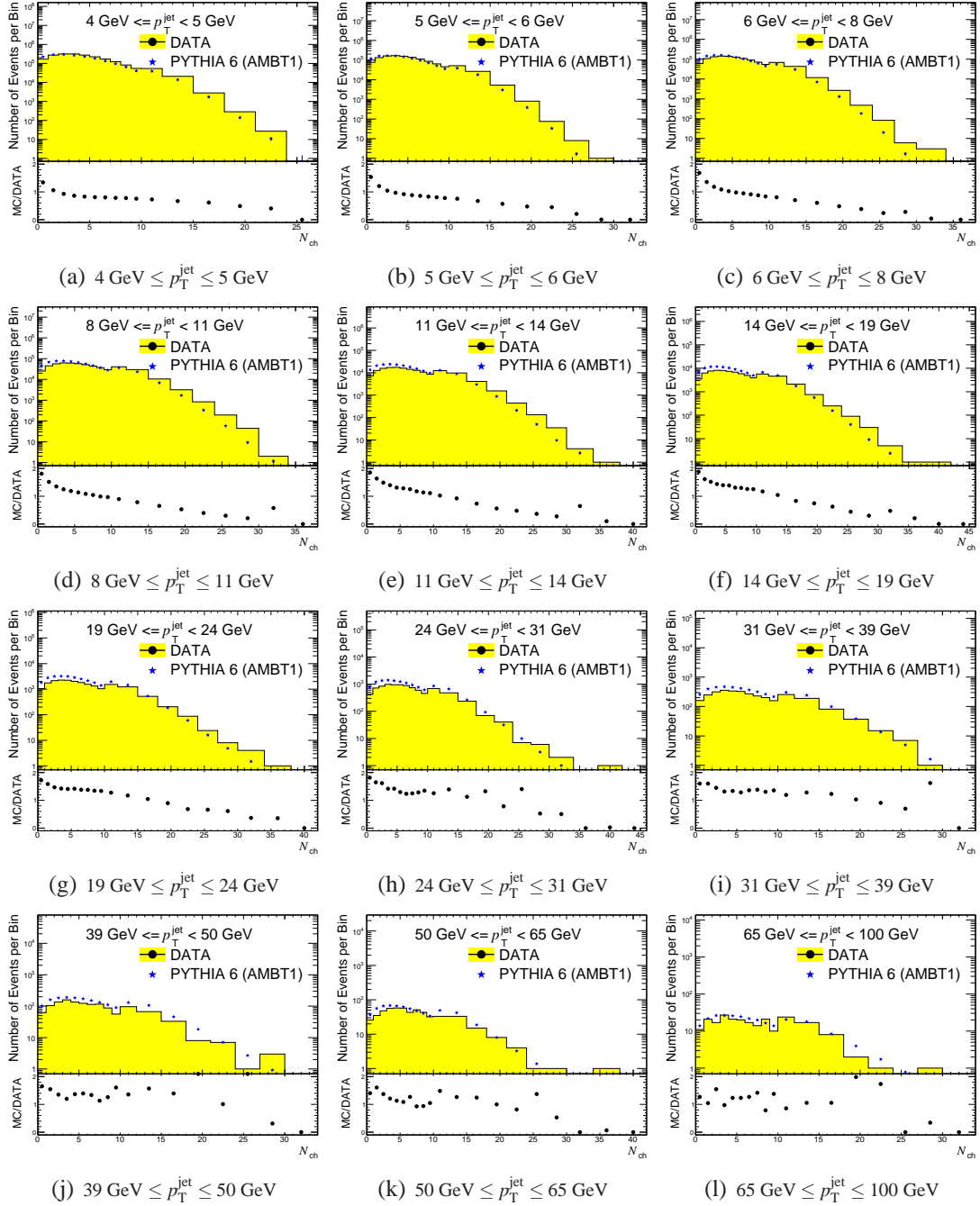


Figure 4.4: The measured N_{ch} distributions before any corrections are applied. The data are compared to PYTHIA 6 (AMBT1). The ratio is shown in the bottom plot.

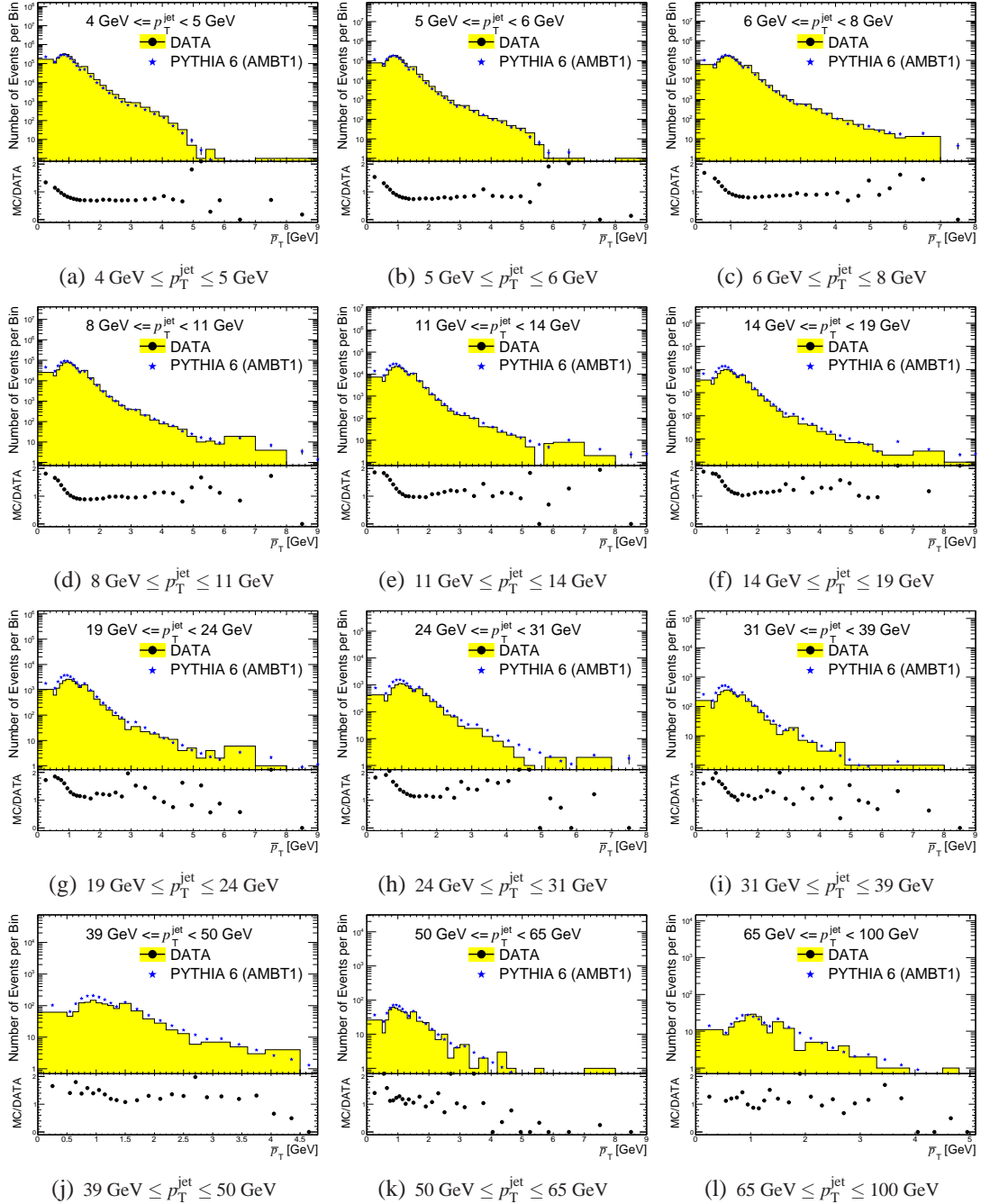


Figure 4.5: The measured \bar{p}_T distributions before any corrections are applied. The data are compared to PYTHIA 6 (AMBT1). The ratio is shown in the bottom plot.

Chapter 5

Correcting the UE Distributions for Detector Effects

As discussed in Sec.3.5, the ATLAS detector has finite momentum resolution and reconstruction inefficiencies for tracks and jets. Our goal is to provide distributions that can be compared to predictions of Monte Carlo (MC) generators, without reference to any detector. In this section, we outline the procedure for making corrections to the observed distributions that will account for the detector effects. We label the desired distributions generically as $f^{\text{true}}(\vec{\mathbf{x}}^{\text{true}})$. We need to relate them to the measured quantities $f^{\text{reco}}(\vec{\mathbf{x}}^{\text{reco}})$. The equation governing the procedure is

$$f^{\text{reco}}(\vec{\mathbf{x}}^{\text{reco}}) = \int \mathbf{R}(\vec{\mathbf{x}}^{\text{true}}, \vec{\mathbf{x}}^{\text{reco}}) f^{\text{true}}(\vec{\mathbf{x}}^{\text{true}}) d\vec{\mathbf{x}}^{\text{true}} \quad (5.1)$$

where we have introduced the concept of the detector *response matrix* $\mathbf{R}(\vec{\mathbf{x}}^{\text{true}}, \vec{\mathbf{x}}^{\text{reco}})$. At first glance, Eq. 5.1 defines a matrix equation that might be invertible. For reasons discussed below, matrix inversion is not the appropriate procedure to use.

The detector response is built using Monte Carlo that has undergone full detector simulation using the GEANT4 framework [23]. We have access to the true values of the distributions $f_{\text{MC}}^{\text{true}}(\vec{\mathbf{x}}^{\text{true}})$, and the reconstructed distributions $f_{\text{MC}}^{\text{reco}}(\vec{\mathbf{x}}^{\text{reco}})$. This analysis uses PYTHIA 6 (AMBT1) and the ATLAS simulation to determine the detector response matrix. Sixty-eight million events were generated in *slices* having specific cuts on truth jet p_{T} for efficient population of the high p_{T} regions of phase space.

5.1 The Response Matrix

In the previous section, we introduced the response matrix \mathbf{R} , encapsulating the detector response to charged particles. To develop the concept and define the response matrix, it is instructive to start with small, but real, examples. The first example starts with the *transfer function* for N_{ch} (the track multiplicity in the transverse region.)

We count the number of events (766387) where the transverse region has a single track at the generator-level. In Table 5.1, we tabulate the number of reconstructed events with 1 track in the transverse region (572545 events), 2 tracks (44749 events), and 3 tracks (6122 events). This corresponds to having 3 bins along the $N_{\text{ch}}^{\text{reco}}$ axis. This situation is also depicted as a histogram in Fig. 5.1. If we divide this histogram by the total number of events with a single track, as we have done in Fig. 5.2(a), we obtain the transfer function for $N_{\text{ch}}^{\text{true}} = 1$. The total visible area in the histogram is the efficiency for an event with $N_{\text{ch}}^{\text{true}} = 1$ to be reconstructed with $1 \leq N_{\text{ch}}^{\text{reco}} \leq 3$. We interpret the individual bin contents of Fig. 5.2(a) as the probabilities that an event with $N_{\text{ch}}^{\text{true}} = 1$ will be reconstructed as an event with $N_{\text{ch}}^{\text{reco}} = 1$, $N_{\text{ch}}^{\text{reco}} = 2$ and $N_{\text{ch}}^{\text{reco}} = 3$, respectively.

We have defined the concept of the transfer function; it is a probability distribution function of the $N_{\text{ch}}^{\text{reco}}$ spectrum corresponding to charged truth jets with specific cuts on $N_{\text{ch}}^{\text{true}}$. Turning our focus back to Table 5.1, we look at the other rows corresponding to different values of $N_{\text{ch}}^{\text{true}}$. By dividing each row by the total number of corresponding events (in the column labelled TOTAL), we have constructed a set of 4 transfer functions, each having 3 bins. The results of this operation are tabulated in Table 5.2, and depicted as histograms in Fig. 5.2. We use the nomenclature \mathbf{R}_J to denote the J^{th} transfer function. For each transfer function \mathbf{R}_J , we define \mathbf{R}_{JK} to be the contents of its K^{th} bin. Note that we have chosen 4 transfer functions with 3 bins; the number of transfer functions need not match the number of bins. \mathbf{R}_{JK} is one of the simplest examples of the response matrix, with elements enclosed within the double lines in Table 5.2. The response matrix is often referred to as the smearing matrix in the literature.

	TOTAL	$N_{\text{ch}}^{\text{reco}} = 1$	$N_{\text{ch}}^{\text{reco}} = 2$	$N_{\text{ch}}^{\text{reco}} = 3$
$N_{\text{ch}}^{\text{true}} = 1$	766387	572545	44749	6122
$N_{\text{ch}}^{\text{true}} = 2$	843600	236531	508282	53453
$N_{\text{ch}}^{\text{true}} = 3$	845925	69646	278381	417639
$N_{\text{ch}}^{\text{true}} \geq 4$	4131439	29187	152576	468731

Table 5.1: Tabulated values of the example response matrix

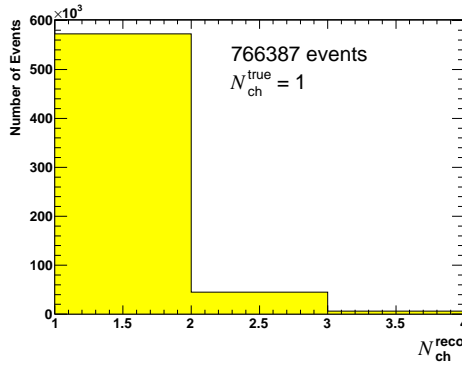


Figure 5.1: A slice of an example response matrix

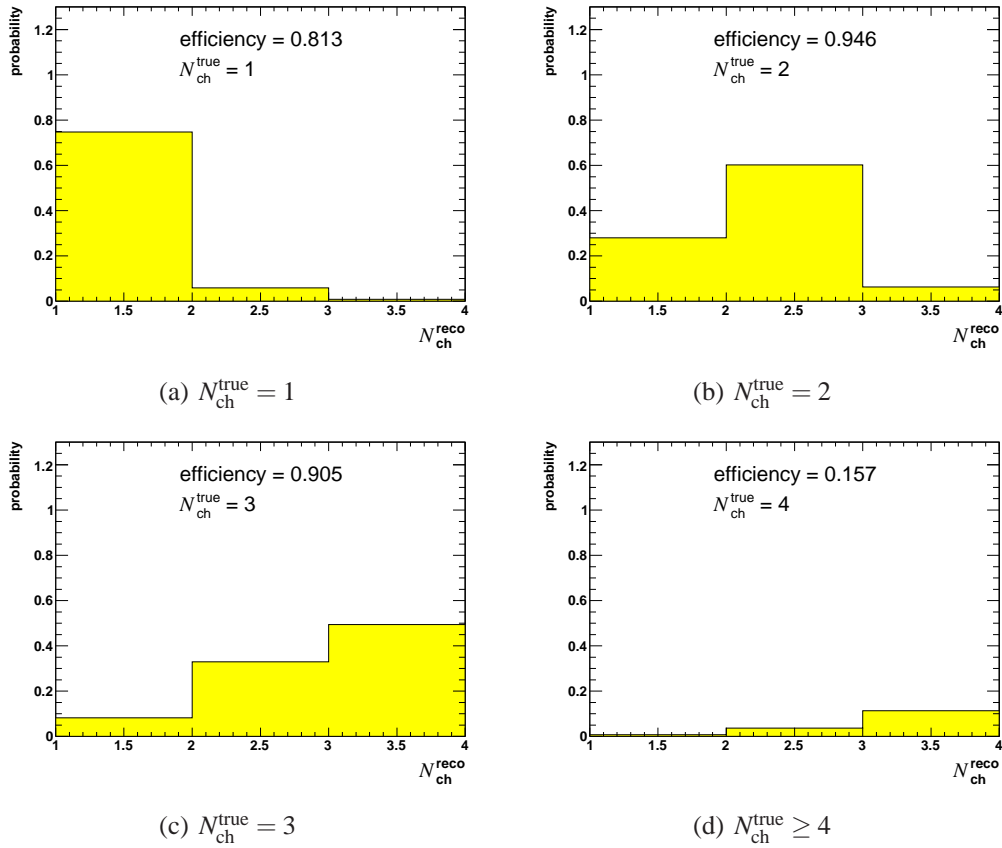


Figure 5.2: The example response matrix for N_{ch} depicted in slices

	EFFICIENCY	$N_{\text{ch}}^{\text{reco}} = 1$	$N_{\text{ch}}^{\text{reco}} = 2$	$N_{\text{ch}}^{\text{reco}} = 3$
$N_{\text{ch}}^{\text{true}} = 1$	0.813	0.747	0.058	0.008
$N_{\text{ch}}^{\text{true}} = 2$	0.946	0.280	0.603	0.063
$N_{\text{ch}}^{\text{true}} = 3$	0.905	0.082	0.329	0.494
$N_{\text{ch}}^{\text{true}} \geq 4$	0.157	0.007	0.037	0.113

Table 5.2: Tabulated values of the example response matrix and efficiency calculations

We have shown a simple example of the response matrix. We are in a position to make a more general definition. The example above used bins, which almost every analysis would do, but we define the "continuous" version of the response matrix $\mathbf{R}(\vec{\mathbf{x}}^{\text{true}}, \vec{\mathbf{x}}^{\text{reco}})$ to be the probability that an object with property $\vec{\mathbf{x}}^{\text{true}}$ is reconstructed with property $\vec{\mathbf{x}}^{\text{reco}}$. The reconstruction efficiency $\varepsilon(\vec{\mathbf{x}}^{\text{true}})$ is given by

$$\varepsilon(\vec{\mathbf{x}}^{\text{true}}) \equiv \int \mathbf{R}(\vec{\mathbf{x}}^{\text{true}}, \vec{\mathbf{x}}^{\text{reco}}) d\vec{\mathbf{x}}^{\text{reco}} \quad (5.2)$$

The entry corresponding to $N_{\text{ch}}^{\text{true}} \geq 4$ in Table 5.2 has a very low efficiency. This points out a subtlety in the definition of efficiency. In our example, the efficiency accounts for events reconstructed $1 \leq N_{\text{ch}}^{\text{reco}} \leq 3$. If we had expanded the range for $N_{\text{ch}}^{\text{reco}}$, we would have a substantially larger efficiency for the $N_{\text{ch}}^{\text{true}} \geq 4$ bin. The definition of efficiency must also specify a range of reconstruction values.

We have denoted the response matrix by $\mathbf{R}(\vec{\mathbf{x}}^{\text{true}}, \vec{\mathbf{x}}^{\text{reco}})$, where the composite nature of the quantity(ies) in question is emphasized using vector-like notation. The variables being corrected are binned, leading to a modified nomenclature $\mathbf{R}_{\mathbf{J}\mathbf{K}}$, where \mathbf{J} refers to the bin(s) containing the truth-level quantity, and \mathbf{K} refers to the bin(s) containing the reconstructed quantities. The response matrix is neither required to have (a) the same binning for truth and reconstructed quantities, nor (b) an equal number of bins for truth and reconstructed quantities. The unfolding algorithm constructs one large, flat vector out of the multiple binned variables, for both truth and reconstructed levels, thereby rendering the response matrix a 2-dimensional matrix. We relax the nomenclature to R_{jk} , without any loss of generality. The efficiency assumes the form:

$$\varepsilon_j \equiv \sum_k R_{jk} \quad (5.3)$$

To make contact with some of the definitions used in the literature [40], which characterize the procedure using "Bayesian terminology", the response matrix is $\mathbf{P}(\mathbf{E}_{\mathbf{K}}|\mathbf{C}_{\mathbf{J}})$. The response matrix is the probability that the J^{th} cause (generator-level object) gave rise to the K^{th} effect (reconstructed object).

5.1.1 Construction of the Response Matrix

For each of the UE observables ($\mathcal{O} = \Sigma p_T, N_{ch}$ and \bar{p}_T), a separate response matrix is created using full GEANT4-simulated PYTHIA 6 (AMBT1). The " Σp_T " response matrix maps three variables ($p_T^{\text{jet}}, p_T^{\text{ext}}, \Sigma p_T$) against the same set of corresponding reconstructed values, where p_T^{jet} is the transverse momentum of the hardest jet with $|\eta| \leq 1.5$ and p_T^{ext} is the transverse momentum of the hardest jet with $|\eta| \geq 1.5$. The construction procedure of the response matrices process for N_{ch} and \bar{p}_T is analogous to that of Σp_T .

The response matrices are constructed for each observable \mathcal{O} , as follows.

- Events are accepted if there exists at least one charged truth jet (with $p_T \geq 1$ GeV and $|\eta| \leq 1.5$), and at least one accepted charged particle jet with the same kinematic acceptance.
 - If the event satisfies the truth-level acceptance criteria, but has no charged particle jet inside the acceptance, the event is recorded as lost due to efficiency. Dedicated bins in the response matrix retain efficiency information.
 - Events with at least one accepted charged particle jet, but no charged truth jets, are not accounted for in the response matrix. This situation is treated as a systematic uncertainty in Section 6.
- The truth level UE observable $\mathcal{O}^{\text{true}}$ in the TRANSVERSE region is calculated using charged primary particles, with $p_T \geq 0.5$, $|\eta| \leq 1.5$ and $\pi/3 \leq |\phi - \phi_0| \leq 2\pi/3$, where ϕ_0 is the azimuth of the leading truth jet. If there are no particles for calculating \mathcal{O} , $\Sigma p_T \equiv N_{ch} \equiv \bar{p}_T \equiv 0$.
- The measured UE observable $\mathcal{O}^{\text{reco}}$ in the TRANSVERSE region is calculated using reconstructed tracks, with $p_T \geq 0.5$, $|\eta| \leq 1.5$ and $\pi/3 \leq |\phi - \phi_0| \leq 2\pi/3$, where ϕ_0 is the azimuth of the leading track jet. If there are no tracks for calculating \mathcal{O} , $\Sigma p_T \equiv N_{ch} \equiv \bar{p}_T \equiv 0$.
- $p_T^{\text{ext,true}}$ is the transverse momentum of the hardest charged truth jet with $|\eta| \geq 1.5$. If no such jet exists, $p_T^{\text{ext}} \equiv 0$.
- $p_T^{\text{ext,reco}}$ is the transverse momentum of the hardest track jets with $|\eta| \geq 1.5$. If no such jet exists, $p_T^{\text{ext}} \equiv 0$.
- The values $(p_T^{\text{jet}}, p_T^{\text{ext}}, \mathcal{O})^{\text{true}}$ are recorded with $(p_T^{\text{jet}}, p_T^{\text{ext}}, \mathcal{O})^{\text{reco}}$.

5.1.2 Purity and Stability

The response matrix is often characterized with two figures of merit, *purity* and *stability*. To develop these concepts in the context of our current analysis, we turn our

focus back to the simple example response matrix from the previous section. Looking at the 766387 events with a single track (at generator level) in the TRANSVERSE region, we see that 572545 are reconstructed with a single track, 44749 are reconstructed with 2 tracks, 6122 are reconstructed with 3 tracks, and so on. The probabilities of reconstruction into the $N_{\text{ch}}^{\text{reco}} = \{1, 2, 3\}$ bins are $\{0.747, 0.058, 0.008\}$, respectively, with an efficiency equal to 0.813. We see that events with a single truth particle have a high probability to be reconstructed into a single bin, in this case it is the bin corresponding to $N_{\text{ch}}^{\text{reco}} = 1$. This is an example of high stability. Stability is the maximum probability that events originating in one bin are reconstructed into a *single* bin. The working definition of stability only concerns itself with reconstructed events, so we divide by the efficiency. The stability of the j -th truth bin is

$$\begin{aligned} \text{stability}_j (j = 1, 2, \dots, N_{\text{bins}}^{\text{true}}) &\equiv \frac{\max \{R_{jk}\} (k = 1, 2, \dots, N_{\text{bins}}^{\text{reco}})}{\sum_{k=1}^{N_{\text{bins}}^{\text{reco}}} R_{jk}} \quad (5.4) \\ &= \frac{\max \{R_{jk}\} (k = 1, 2, \dots, N_{\text{bins}}^{\text{reco}})}{\varepsilon_j} \end{aligned}$$

In our example, the " $N_{\text{ch}}^{\text{true}} = 1$ " bin has a stability equal to $0.747 / 0.813 = 0.918$.

The concept of purity complements stability. Purity is the maximum probability that objects reconstructed in a bin originated in a *single* bin at truth level. In our example (see Table 5.3), we see that of the 946690 events that were reconstructed with 3 tracks in the TRANSVERSE region, 745 originated as events with no truth particles, 6122 originated as events with a single truth particle, and so on. Scanning down the column corresponding to $N_{\text{ch}}^{\text{reco}} = 3$ we have an exhaustive list of possible origins for any event. An important subtlety in the construction of the response matrix is that every reconstructed object must have originated in exactly one truth bin.

The definition for the purity of the k -th reconstructed bin, as used in this analysis, is

$$\text{purity}_k \{k = 1, 2, \dots, N_{\text{bins}}^{\text{reco}}\} \equiv \frac{\max (R_{jk}) (j = 1, 2, \dots, N_{\text{bins}}^{\text{true}})}{\sum_{j=1}^{N_{\text{bins}}^{\text{true}}} R_{jk}} \quad (5.5)$$

Applying Eqn. 5.5, the purity for the bin corresponding to $N_{\text{ch}}^{\text{reco}} = 3$ is 0.495.

High stability is desirable, usually indicating a sharp detector response. High stability can also be achieved by making larger bins. Low stability indicates that the resolution of the truth objects is wide, smeared out across more than one bin. In practice, as it relates to the Bayesian Iterative Unfolding algorithm, higher stabilities are desirable in the regions of low statistics. Lower stability values in regions with high statistics will not significantly affect the results, but one must exercise caution to quantify the size of the effect. We have verified this is the case for our analysis. In regions of

	$N_{ch}^{reco} = 1$	$N_{ch}^{reco} = 2$	$N_{ch}^{reco} = 3$	STABILITY
$N_{ch}^{true} = 0$	26666	2833	745	-
$N_{ch}^{true} = 1$	572545	44749	6122	0.918
$N_{ch}^{true} = 2$	236531	508282	53453	0.637
$N_{ch}^{true} = 3$	69646	278381	417639	0.545
$N_{ch}^{true} \geq 4$	29187	152576	468731	0.721
purity	0.613	0.515	0.495	

Table 5.3: Purity and stability for the example response matrix

low statistics, large statistical fluctuations tend to couple into neighboring bins. The extent of the coupling depends on the stability - high stability keeps the correlations between neighboring bins low. We chose the bin sizes to balance between reasonable stability and good resolution; the emphasis is on keeping the bin widths small for better resolution. Fig. 5.3 (5.4) shows the stability (purity) for p_T^{jet} , Σp_T , N_{ch} , and \bar{p}_T .

5.2 Bayesian Iterative Unfolding with RooUnfold

The RooUnfold framework [41] implements an iterative algorithm proposed in [40], based on Bayes' Theorem in the following form:

$$P(C_j|E_k) = \frac{P(E_k|C_j)P(C_j)}{\sum_{j=1}^{n_c} P(E_k|C_j)P(C_j)} \quad (5.6)$$

where $P(C_j)$ is the probability of the j^{th} cause, $P(E_k|C_j)$ is the conditional probability of the j^{th} cause to produce the k^{th} effect, and $P(C_j|E_k)$ is the probability that the k^{th} effect was due to the j^{th} cause. Translating into our formalism:

$$P(C_j|E_k) = \frac{y_j R_{jk}}{\sum_k y_j R_{jk}} \quad (5.7)$$

where y_k denotes the content of the k -th truth bin. The algorithm starts by using an initial distribution (*prior*) as an estimate of the final distribution (*posterior*). We use the relevant PYTHIA 6 (AMBT1) distribution as the prior for the first iteration, from which the algorithm produces an improved estimate of the final distribution. Each subsequent iteration processes the output distribution of the previous iteration, to produce yet another improved estimate. Since each iteration uses Bayes' Theorem in an intermediate step, the algorithm is often referred to as Bayesian Iterative Unfolding. We continue the process of iteration until the output has stabilized, at which point

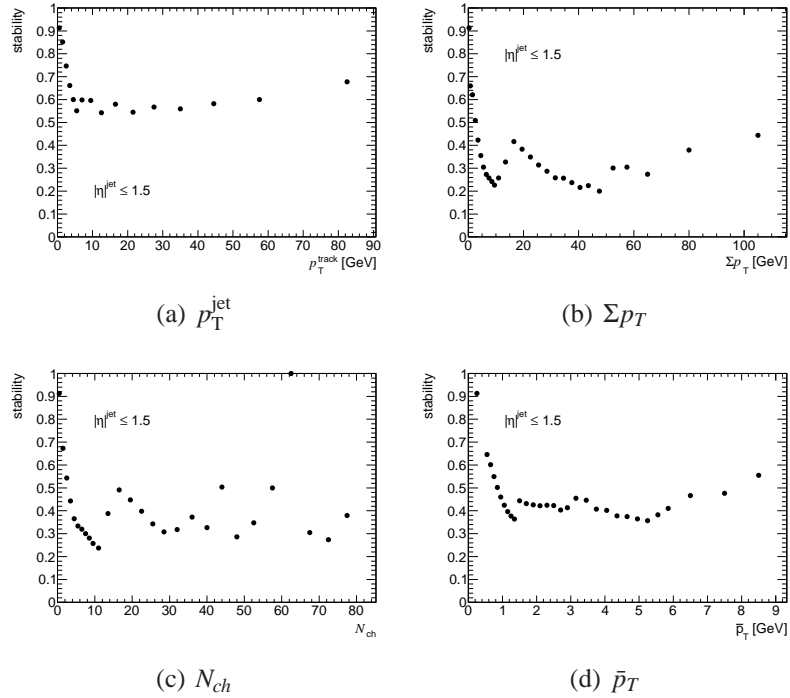


Figure 5.3: The stability for PYTHIA 6 (AMBT1)

we stop. It is important not to perform excessive iterations after the output has stabilized, for reasons we will discuss shortly. We characterize stabilization by visually examining the χ^2 between the prior and posterior distributions of each iteration. Our decision to use 4 iterations for each of the UE observables was based on observing the behavior on Monte Carlo distributions.

The difference between stabilization and convergence is a very important distinction, and a topic in the theory of *regularization* in unfolding methods [42, 43, 44]. The salient point is that if we allow the solution to converge, it starts to track the statistical fluctuations in the input distribution, producing unphysical ripples in the output. Stability plays a role in controlling these ripples. The problem with these ripples is that, from a mathematical viewpoint, they are the correct solution. We want to iterate the solution until the large scale structure is resolved (the basic form of the curves), and stop before we track the fine scale structure (the statistical fluctuations). Different number of iterations give different results, leading to an uncertainty in the final answers. We account for this uncertainty in Section 6.2.

It is outside the scope of this work to delve into the theoretical details of the algorithm. We describe the mechanical aspects of the algorithm as it pertains to the current analysis, and how it uses the response matrix R to perform the unfolding.

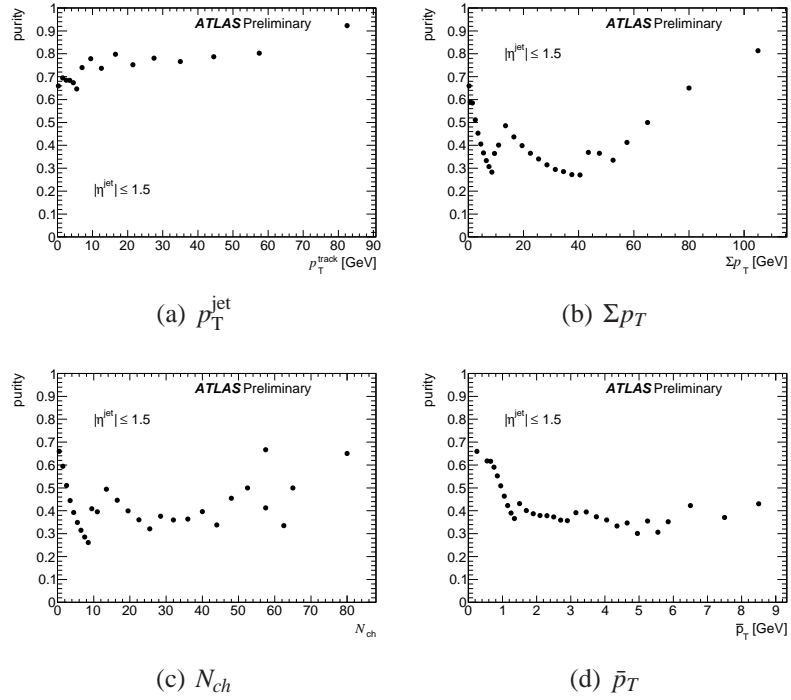


Figure 5.4: The purity for PYTHIA 6 (AMBT1)

1. Let $\mathbf{p}_0 \equiv (p_1, p_2, \dots, p_M)$ be an initial set of probabilities (derived from the input spectrum, or even a constant value) for an event to be found in each bin, and $n_{tot} = \sum_i n_i$ be the total number of entries.
2. Define

$$\hat{\mu}_0 \equiv n_{tot} \mathbf{p}_0 \quad (5.8)$$

3. Update to a new value μ , using the response matrix R in indexed form R_{jk} . This step is motivated by Eqn. 5.6.

$$\hat{\mu}_i = \frac{1}{\epsilon_i} \sum_{j=1}^N \left(\frac{R_{ij} p_i}{\sum_k R_{kj} p_k} \right) n_j \quad (5.9)$$

4. Form new probabilities \mathbf{p}

$$\mathbf{p}_k = \frac{\hat{\mu}_k}{n_{tot}} \quad (5.10)$$

5. Iterate steps 3 and 4 on Monte Carlo, until the change in χ^2 between iterations indicates the distribution has stabilized, and before the unfolded distributions

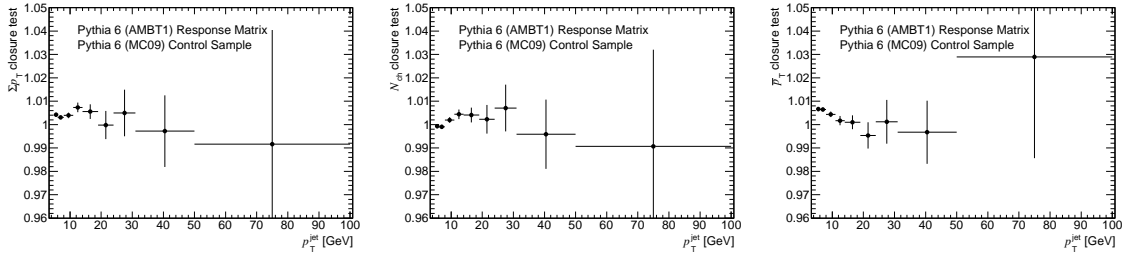
start to track the statistical fluctuations from the input. The same number of iterations N_{iter} is used when unfolding data. This analysis uses $N_{iter} = 4$, as described earlier in this section.

5.3 Validation of the Unfolding Procedure

To characterize the unfolding process, we perform *closure tests* using full GEANT4 simulation PYTHIA 6 (MC09) [17] as a control sample. We apply the unfolding techniques described in Sec. 5.2 to unfold the PYTHIA 6 (MC09) control sample and compare the results to the known truth values. The response matrix is used to unfold the control sample distributions and calculate the mean values as functions of p_T^{jet} .

We quantify the closure tests by taking the ratio of the mean values of the corrected distributions to the true mean values. A value of 1.0 indicates total closure - indicating the corrected and true values agree perfectly. The closure tests for the mean values of the UE observables are shown in Fig. [5.5], indicating a 1% performance level for the mean values of the UE observables for $p_T^{\text{jet}} \leq 50\text{GeV}$. The \bar{p}_T closure tests degrade to the 3% level above $p_T^{\text{jet}} \geq 50\text{GeV}$, albeit with a large statistical uncertainty.

These tests confirm that the correction procedure works as intended. We return to the closure test performance in Chapter 6, when we analyze systematic uncertainties.



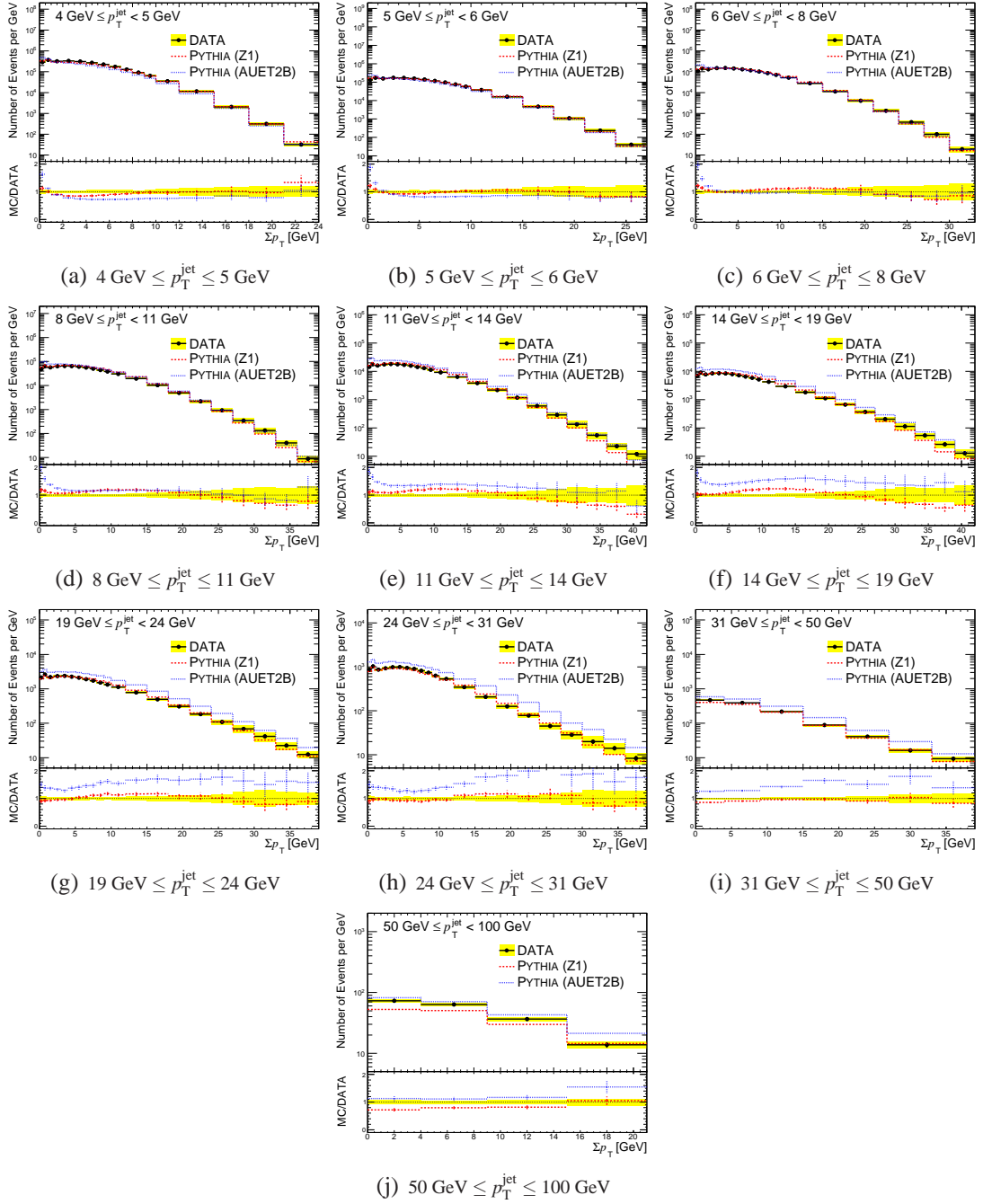
(a) Σp_T in the TRANSVERSE re- (b) N_{ch} in the TRANSVERSE region (c) \bar{p}_T in the TRANSVERSE region
 gion

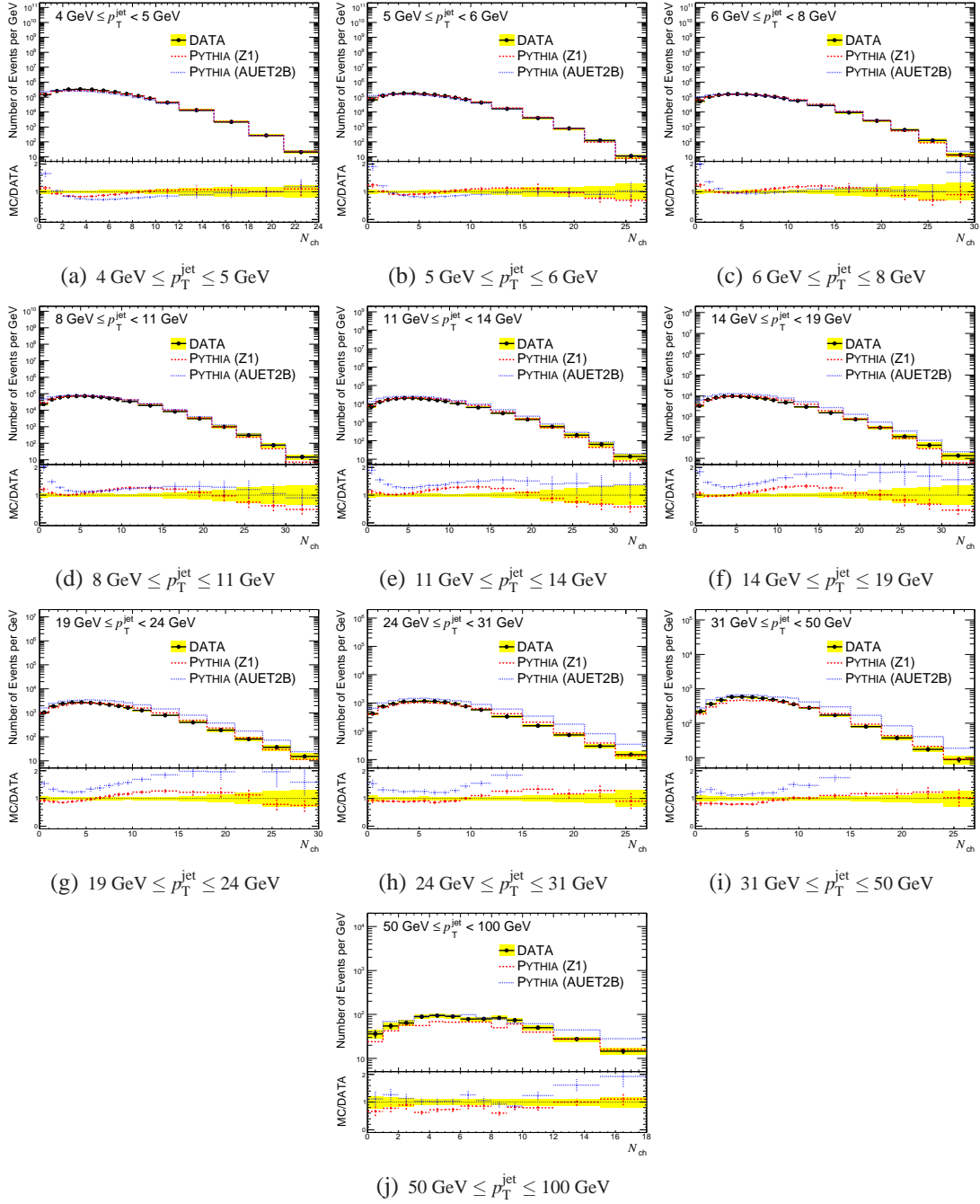
Figure 5.5: The closure tests for the mean values of the UE observables, as a function of p_T^{jet} . Error bars reflect statistical uncertainties.

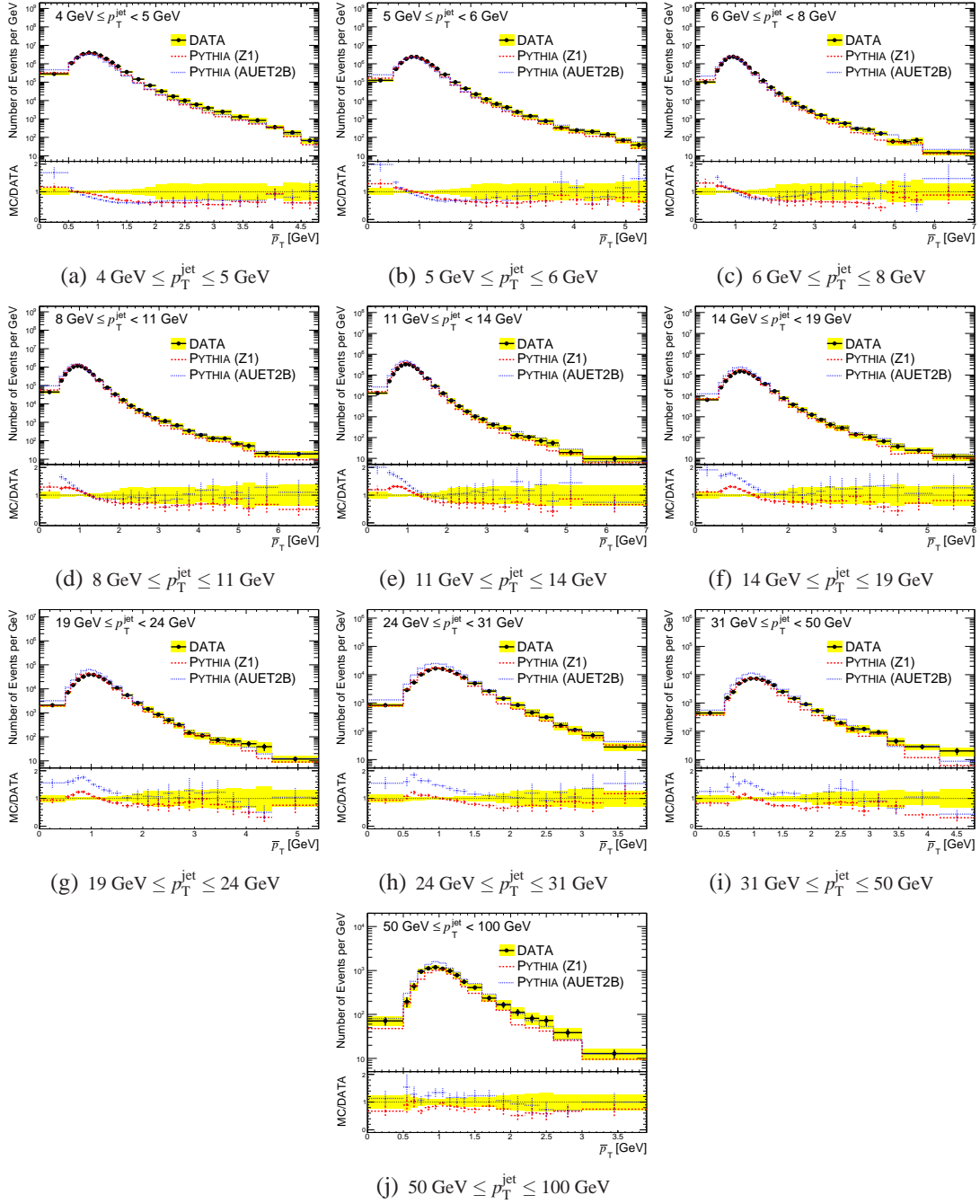
5.4 Corrected Distributions

The output of the correction process is a set of 3D histograms containing the final distributions of p_T^{jet} , p_T^{ext} and \mathcal{O} . The final form of the p_T^{ext} spectrum does not interest us, so we integrate over it, essentially projecting it onto the remaining axes. The results of

this projection are 2D histograms (\mathcal{O} vs p_T^{jet}). The information in these histograms is best rendered as slices in the \mathcal{O} variable (projections along y-axis), holding p_T^{jet} fixed along the x-axis (Figs.[5.6-5.8]). The data are compared to PYTHIA (Z1) and PYTHIA (AUET2B). The agreement is good, but PYTHIA (Z1) reproduces the data distributions better than the AUET2B tune.

Figure 5.6: The corrected Σp_T distributions

Figure 5.7: The corrected N_{ch} distributions

Figure 5.8: The corrected \bar{p}_T distributions

5.4.1 Mean Values of Corrected Distributions

As pointed out in the previous section, the output of the correction process is a histogram, with binned contents. It is straightforward to taking the binned mean value, defined in Eqn. 5.11,

$$\mu_{\text{binned}} = \frac{\sum_{k=1}^N n_k x_k}{\sum_{k=1}^N n_k} \quad (5.11)$$

but potentially misleading if interpreted as the true mean value of the distribution defined in Eqn. 5.12.

$$\mu_{\text{true}} = \frac{\int_0^{\infty} x n(x) dx}{\int_0^{\infty} n(x) dx} \quad (5.12)$$

If the bins are sufficiently small, the difference between the binned mean and true mean values is small, and the uncertainty in the (true) mean value associated with finite bin widths is negligible. The available statistics for this analysis preclude small bins at high p_T and/or large values of the UE observables; we compensate with larger bins. Fortunately, we can use cubic splines to correct for the bias due to large bin widths. Fitting a cubic spline to the cumulative distributions, not the actual distributions, is a well-defined process. In calculating the mean value of the distribution, the integral of the distribution is more useful than the actual distribution itself. This is outlined in detail in Appendix B, Only Σp_T and N_{ch} require correction; \bar{p}_T is sufficiently finely binned and doesn't require the spline-based corrections.

Figs. [5.9-5.11] show the mean values of the corrected distributions of Σp_T , N_{ch} , and \bar{p}_T , obtained using the spline-based approach for removing the bias due to large bin widths.

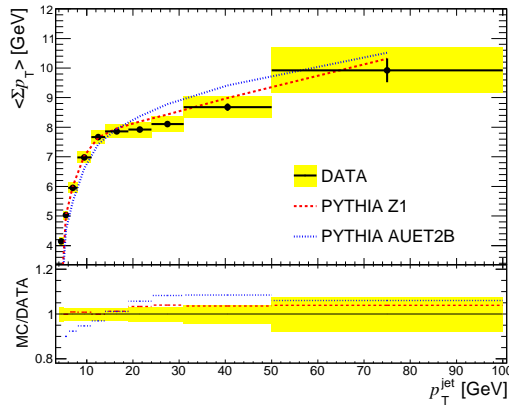


Figure 5.9: The mean values of the corrected Σp_T distributions, as functions of p_T^{jet} , are compared to Monte Carlo. The error bars indicate the statistical uncertainty; the shaded area shows the combined statistical and systematic uncertainties.

The data are compared to PYTHIA (Z1) and PYTHIA (AUET2B). The agreement is good, but PYTHIA (Z1) reproduces the data distributions better than the AUET2B tune.

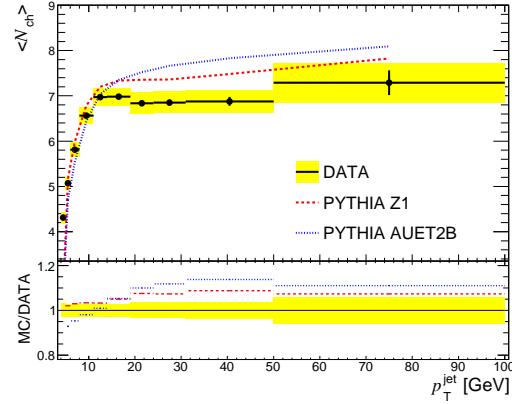


Figure 5.10: The mean values of the corrected N_{ch} distributions, as functions of p_T^{jet} , are compared to Monte Carlo. The error bars indicate the statistical uncertainty; the shaded area shows the combined statistical and systematic uncertainties.

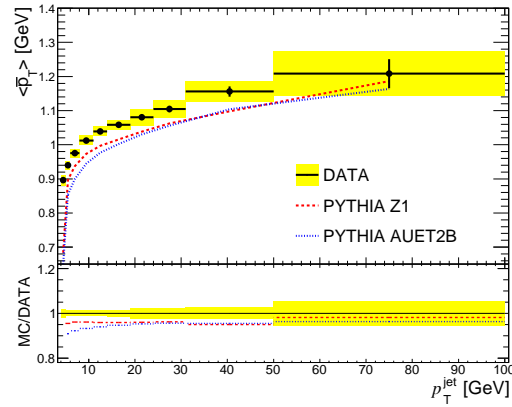


Figure 5.11: The mean values of the corrected \bar{p}_T distributions, as functions of p_T^{jet} , are compared to Monte Carlo. The error bars indicate the statistical uncertainty; the shaded area shows the combined statistical and systematic uncertainties.

Chapter 6

Uncertainty Analysis

Table 6.1: The systematic uncertainties associated with measurement of the mean values of Σp_T , N_{ch} and \bar{p}_T .

Relative Systematic Uncertainties			
source	Σp_T	N_{ch}	\bar{p}_T
Track Reconstruction	2.3%	2.1%	0.2%
Unfolding	1.5%-6%	1.5%-4%	1%-4%
Response Matrix	0.5%-1%	0.5%-1%	0.5%-1%
Lead Jet Misidentification	$\leq 1\%$	$\leq 1\%$	$\leq 1\%$
Discretization Effects	$\leq 0.5\%$	$\leq 0.5\%$	$\leq 0.5\%$
Dependence on Number of Iterations	$\leq 0.5\%$	$\leq 0.5\%$	$\leq 0.5\%$
Total	2.9%-6.5%	2.7%-4.6%	1.3%-4.1%

Table 6.1 summarizes the systematic uncertainties associated with the measurement of the UE distributions. In this section, we discuss each of these sources of uncertainty, which are

1. Track Reconstruction - the effects of imperfect efficiency and momentum resolution
2. Uncertainty in the Unfolding Procedure - potential bias from the unfolding procedure
3. Sensitivity to the Response Matrix - potential bias due to differences in distributions between the data and MC used to build the response matrix
4. Misidentification of the Leading Jet - the leading jet corresponds to a subleading jet
5. Discretization Effects - large bin widths introduce a potential bias in the mean values
6. Dependence on Number of Iterations - the optimal number of iterations used in the unfolding procedure

6.1 Track Reconstruction

6.1.1 Track Momentum Resolution

The momentum resolution uncertainty was studied in detail [37]. We assess the induced uncertainty in the baseline measurements, due to the uncertainty in track momentum resolution, by smearing the track momentum [45]. The momentum resolution for the tracks in our sample is excellent; the track momentum resolution uncertainty induces a negligible uncertainty in our measurement ($\leq 0.1\%$ for $p_T^{\text{jet}} \leq 20$ GeV and $\leq 0.5\%$ for $p_T^{\text{jet}} \geq 20$ GeV.)

6.1.2 Tracking Efficiency

Because the ID has substantial material, charged particles can be lost due to hadronic interactions. Uncertainties in the ID material budget [37] result in an uncertainty in the track reconstruction efficiency, which propagate into our measurements. The uncertainties in the tracking efficiency are approximately 2% for $|\eta| \leq 1.3$, 3% for $1.3 \leq |\eta| \leq 1.9$, 4% for $1.9 \leq |\eta| \leq 2.3$ and rises to 7% for $2.3 \leq |\eta| \leq 2.5$ [46] for tracks with $p_T \geq 0.5$ GeV [37, 47].

We propagate the uncertainties in the track reconstruction efficiency into an uncertainty in the measurement of the underlying event as follows:

1. For each track in the TRANSVERSE regions, generate a uniform random number x between 0 and 1.
2. If the track has $|\eta| \leq 1.3$, retain it if $0.98 \leq x$. Otherwise, the track is discarded.
3. If the track has $1.3 \leq |\eta| \leq 1.5$, retain it if $0.97 \leq x$. Otherwise, the track is discarded.
4. Perform a measurement of the UE observables using the retained reconstructed tracks from the previous steps.
5. Unfold the measurement in (4) using the baseline response matrix.
6. Compare the results of the unfolding procedure in the previous step to the baseline measurement.
7. The relative deviation from the baseline is taken as the uncertainty in the measurement

The method of discarding tracks to simulate a different tracking efficiency only works for a lower efficiency; it will not work for a higher efficiency. The propagated uncertainties are relatively small; we symmetrize the uncertainty due to tracking efficiency.

Fig. [6.1] shows the uncertainties in Σp_T , N_{ch} and \bar{p}_T due to uncertainties in the tracking efficiency. The uncertainties in Σp_T and N_{ch} are generally between 2% to 3%; for

\bar{p}_T , the uncertainties are much lower ($< 0.5\%$), as the uncertainties factor out in the ratio $\Sigma p_T / N_{ch}$. We assess the final systematic uncertainty by performing a fit to a *constant* value throughout the entire p_T range, to compensate for the loss in statistical power at high p_T^{jet} . The final values assessed for the uncertainties are denoted by the dotted lines.

Figs. [6.9 - 6.14] show the uncertainty for the individual bins of p_T^{jet} and \mathcal{O} . We point out the contribution from the uncertainty in the tracking efficiency to the uncertainty in the calculation of the mean values of the UE, may be substantially smaller than those of the individual bins. This effect is due to the high correlation in the uncertainties in the individual bins; high upward fluctuations in some bins guarantee a downward fluctuation in others. Fig. [6.2] illustrates the correlation in the uncertainties in the individual bins.

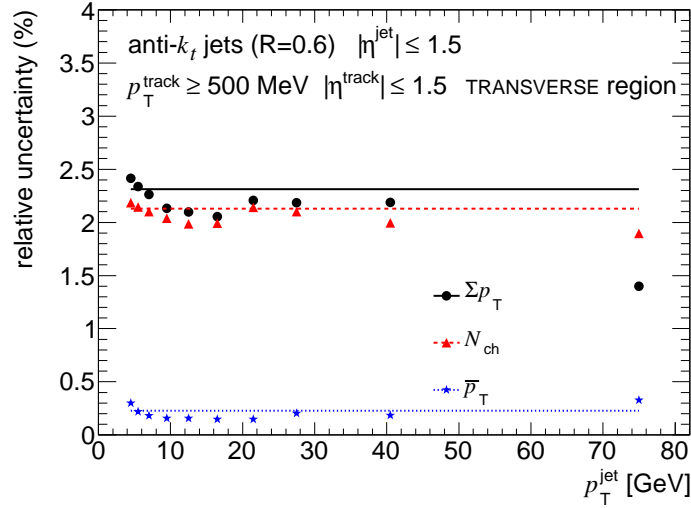


Figure 6.1: The relative uncertainties in the measurement of mean values of the UE distributions, due to uncertainties in the tracking efficiency, as functions of p_T^{jet} . The horizontal lines indicate the assessed uncertainty.

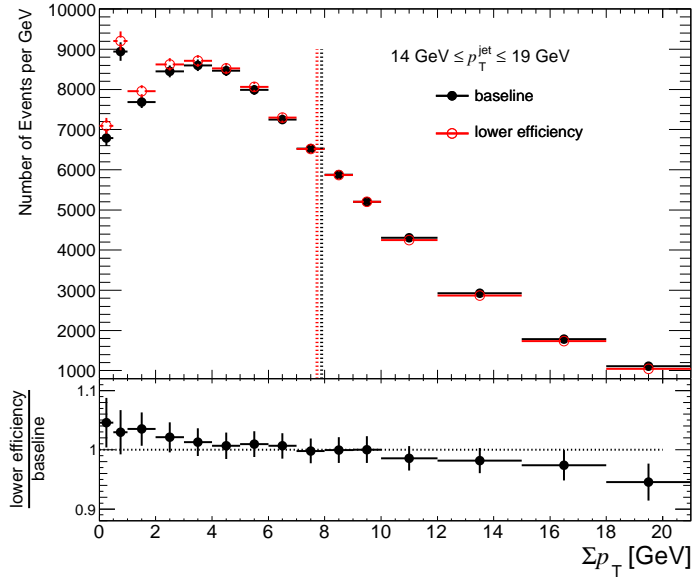


Figure 6.2: The *slice* of Σp_T corresponding to $14 \text{ GeV} \leq p_T^{\text{jet}} \leq 19 \text{ GeV}$. The black points represent the baseline measurement. The red points represent the measurement, propagating the uncertainty in the tracking efficiency. The vertical lines depict the mean values of the respective distributions.

6.2 Uncertainty in the Unfolding Procedure

6.2.1 Unfolding Uncertainty for Individual Bins of p_T^{jet} and \mathcal{O}

The closure test for a single bin in p_T and \mathcal{O} , where the value N_{truth} is expected, is the value

$$\frac{N_{\text{corrected}}}{N_{\text{truth}}} \quad (6.1)$$

The unfolding process provides an excellent, but not perfect, description of the true physics distributions. In this section, we describe how accurately we can expect the unfolding to predict the true p_T^{jet} and \mathcal{O} distributions. We saw in Section 5.3 that the unfolding process reproduces the mean values of the UE observables within a few percent, with the uncertainties on individual bins of p_T^{jet} and \mathcal{O} somewhat higher. To estimate how well the unfolding procedure works on the data, we should use Monte Carlo control samples that resemble the data, in both the physics distributions and statistics. Since the MC and data are not in agreement with respect to the jet p_T and UE observables, we reweight the MC to reproduce the data. We compensate for the lack of statistics using the bootstrap method [48] to increase the statistical power of the MC. The definitions of the variables and details of this process follow:

1. N_{MC} and N_{DATA} are the numbers of events in the MC and data, respectively.
2. $H_{DATA}(p_T^{\text{jet}}, \mathcal{O})$ is the unfolded data distribution of track p_T^{jet} and \mathcal{O} , represented by a 2D-histogram.
3. $H_{MC}(p_T^{\text{jet}}, \mathcal{O})$ is the truth MC distribution of charged truth p_T^{jet} and \mathcal{O} , represented by a 2D-histogram.
4. For each event, ω_0 is a random number drawn from a Poisson distribution with mean $\mu = \frac{N_{DATA}}{N_{MC}}$.
5. $H_0(p_T^{\text{jet}}, \mathcal{O}) \equiv \frac{H_{DATA}}{H_{MC}} \frac{N_{MC}}{N_{DATA}}$ is the weighting 2D-histogram.
6. For each MC event, the leading charged truth p_T^{jet} and \mathcal{O} are used to index the weight ω_1 from H_0 .
7. The same weight $\omega \equiv \omega_0 \times \omega_1$ is applied to the event at generator and reconstructed levels, when constructing new \mathcal{O} vs p_T^{jet} histograms.

After reweighting the Monte Carlo control sample to reproduce the p_T^{jet} and \mathcal{O} distributions, there may exist residual differences between the data and MC. For example, the topological distribution of tracks inside the TRANSVERSE region will differ for the Monte Carlo and the data. The tracks in the data may be more uniformly distributed in the TRANSVERSE region, whereas the Monte Carlo may exhibit "clumpiness". Another example is the subleading p_T^{jet} distribution. A control sample with a harder subleading p_T^{jet} distribution will have a higher likelihood of erroneously promoting the subleading jet to the leading jet, due to track jet p_T resolution. These differences introduce a potential bias in the unfolding procedure.

In general, as we repeat the closure tests using different MC, with closure test x , we would obtain a collection F of closure tests. The mean value (μ_F) of F is the bias in the unfolding procedure. The RMS (σ_F) of F is the dispersion in the closure tests. The total uncertainty in the unfolding procedure is given by

$$\sigma_{\text{tot}} \equiv \mu_F \oplus \sigma_F \equiv \sqrt{\mu_F^2 + \sigma_F^2} \quad (6.2)$$

The uncertainty is calculated as follows, using Σp_T as an example. The analogous procedures for N_{ch} and \bar{p}_T are otherwise identical.

- For each MC and ($p_T^{\text{jet}}, \Sigma p_T$) bin, N_0 is the expected (generator-level) value and N is the corrected value.
- Form the sum $\bar{\omega} \equiv \sum_{MC} 1/(\sigma_S^{\text{MC}})^2$, where σ_S^{MC} is the statistical uncertainty in the unfolding procedure of each MC.

- ω is a weight formed from the statistical uncertainty (σ_S^{MC}) in the unfolding procedure. $\omega \equiv 1 / \left((\sigma_S^{\text{MC}})^2 \bar{\omega} \right)$.
- $\rho \equiv N/N_0$ is the closure test.
- μ_C is the weighted mean of the closure tests. $\mu_C \equiv \sum_i \omega_i \rho_i / \sum_i \omega_i$
- β_C is the bias $\equiv |\mu_C - 1|$.
- σ_C is the weighted RMS of the closure tests. $\sigma_C^2 \equiv \frac{\sum_i \omega_i}{(\sum_i \omega_i)^2 - \sum_i \omega_i^2} \sum_i \omega_i (\rho_i - \mu_C)^2$.

We have four fully simulated Monte Carlo control samples available - PYTHIA (6) with MC09, AMBT1 and Perugia2010 tunes, and PYTHIA (8.145) with the 4C tune. To illustrate the procedure for evaluation of this uncertainty, Fig. 6.3 compares the unfolded and generator-level spectra of Σp_T for $11 \text{ GeV} \leq p_T^{\text{jet}} \leq 14 \text{ GeV}$, obtained by reweighting the control samples to reproduce the data.

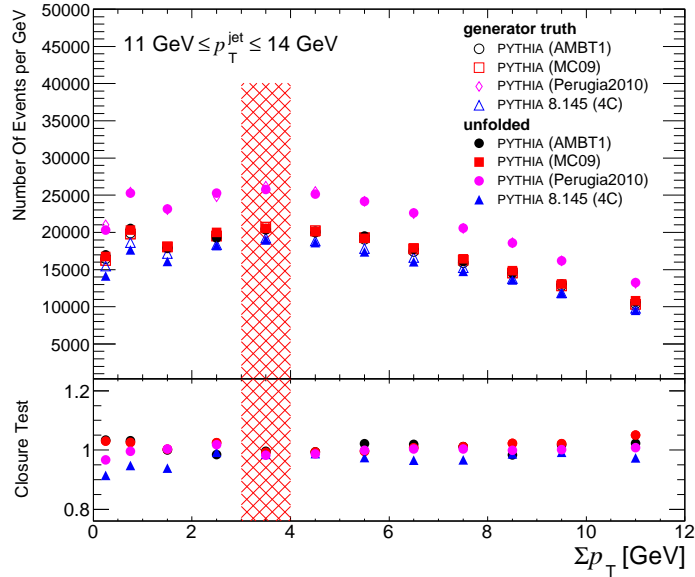


Figure 6.3: The spectra of Σp_T corresponding to $11 \text{ GeV} \leq p_T^{\text{jet}} \leq 14 \text{ GeV}$, for different Monte Carlo control samples. The bottom plot shows the ratio of the unfolded values to the generator-level (truth) values.

As an example, we will work out the detailed calculation of the uncertainty for one bin, corresponding to $3 \text{ GeV} \leq \Sigma p_T \leq 4 \text{ GeV}$ (hatched region in Fig. 6.3). The expected (generator-level) and corrected (unfolded) values are tabulated in Table [6.2], along with the numbers of merit used to calculate the uncertainty. The uncertainty for this bin can be seen as the large red star in Fig. 6.4, along with the uncertainties for the other bins in Fig. 6.3.

Control Sample	Expected N_0	Unfolded $N \pm \sigma_S$	Closure Test x	Weight ω
PYTHIA 6 (AMBT1)	20438	20349 ± 332	0.996	0.151
PYTHIA 6 (MC09)	20701	20548 ± 180	0.993	0.515
PYTHIA 8	19267	18980 ± 382	0.985	0.114
PYTHIA 6 (Perugia 2010)	26227	25752 ± 276	0.982	0.220

Property	Value
bias	-0.010
RMS	0.006
bias \oplus RMS	0.012

Table 6.2: Example calculation of unfolding uncertainty

We must discuss another important effect before finalizing the calculation of the unfolding uncertainty. The unfolding uncertainty has a statistical component which can be large, especially in regions with low statistics. The effects of a large statistical component of the unfolding uncertainty can be seen in Fig. 6.5, where we plot the unfolding uncertainty for $3\text{GeV} \leq \Sigma p_T \leq 4\text{GeV}$, as a function of p_T^{jet} . The black points show the uncertainties (with error bars), as just discussed in the text. The red points are the statistical uncertainty taken directly from the unfolding algorithm. We can see that the statistical uncertainty strongly influences the calculations of the unfolding uncertainty. To avoid overestimation of the statistical uncertainty in our final measurements, we must properly remove this component. Assuming the resolution of Σp_T in the TRANSVERSE region is independent of p_T^{jet} ¹, we harness the power of high statistics in the lower p_T^{jet} regions by fitting to a constant value for the unfolding uncertainty for each bin in Σp_T . The results can be seen in Fig. 6.5. The numbers from the yellow band comprise the curves labelled "unfolding" in Figs.[6.9 - 6.14].

¹Whereas the mean values of the Σp_T and p_T^{jet} distributions are correlated, the resolution of each variable is not because they are in different regions of the detector.

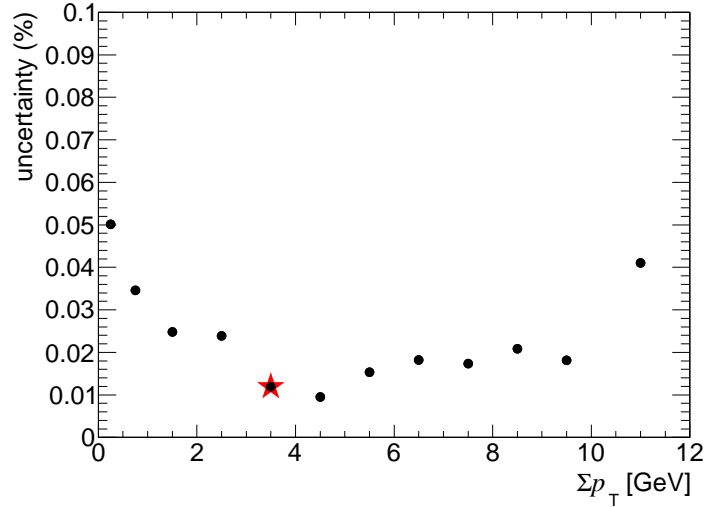


Figure 6.4: The uncertainty in the Σp_T spectra corresponding to $24 \text{ GeV} \leq p_T^{\text{jet}} \leq 31 \text{ GeV}$. The red star indicates the uncertainty discussed and calculated in the text, as an example.

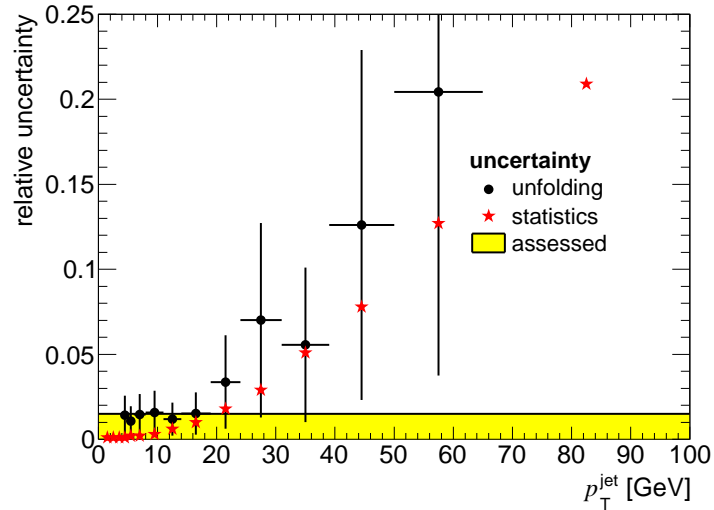


Figure 6.5: The unfolding uncertainty in Σp_T as a function of p_T^{jet} , compared to the statistical uncertainty. The yellow band denotes the final assessed value of the uncertainty, obtained by performing a fit to a constant value across the p_T^{jet} axis.

We have just prescribed a method for evaluating the uncertainty in the number of events in any bin for p_T^{jet} and Σp_T . For the mean values of the Σp_T N_{ch} and \bar{p}_T distributions, as functions of p_T^{jet} , we take an analogous approach to assessing the uncertainty in the unfolding procedure. The number of events in a bin is replaced with the mean value of the spectrum.

6.2.2 Unfolding Uncertainty in the Mean Values of \mathcal{O}

We calculate the unfolding uncertainty in the mean value calculations much in the same manner as for the individual bins. Due to a high degree of correlation in the uncertainties between the individual bins of p_T^{jet} and \mathcal{O} , we obtain a better estimate in the uncertainty in the mean value by examining it directly, not simply propagating the individual bin uncertainties through Eqn. 5.11.

The uncertainty is calculated as follows, using Σp_T as an example. The analogous procedures for N_{ch} and \bar{p}_T are otherwise identical.

- For each control sample and $(p_T^{\text{jet}}, \Sigma p_T)$ bin, μ_0 is the expected (generator-level) value and μ is the corrected value.
- Form the sum $\bar{\omega} \equiv \sum_{\text{MC}} 1/(\sigma_S^{\text{MC}})^2$, where σ_S^{MC} is the statistical uncertainty in the unfolding procedure of each MC.
- ω is a weight formed from the statistical uncertainty (σ_S^{MC}) in the unfolding procedure. $\omega \equiv 1/((\sigma_S^{\text{MC}})^2 \bar{\omega})$.
- $\rho \equiv \mu/\mu_0$ is the closure test.
- μ_C is the weighted mean of the closure tests. $\mu_C \equiv \sum_i \omega_i \rho_i / \sum_i \omega_i$
- β_C is the bias $\equiv |\mu_C - 1|$.
- σ_C is the weighted RMS of the closure tests. $\sigma_C^2 \equiv \frac{\sum_i \omega_i}{(\sum_i \omega_i)^2 - \sum_i \omega_i^2} \sum_i \omega_i (\rho_i - \mu_C)^2$.

6.3 Sensitivity to the Response Matrix

The baseline response matrix was constructed using PYTHIA 6 (AMBT1), which has different p_T^{jet} and \mathcal{O} distributions than the data. These differences lead to a potential bias in the measurement. We estimate the size of this bias by constructing an alternate response matrix, formed from PYTHIA 6 (AMBT1) which has been reweighted to reproduce the corrected distributions from data.

Reweighting the Monte Carlo used to construct the response matrix closes the differences between it and the data. Residual differences contribute second order effects. Using

the reweighted response matrix to perform the correction procedure, we calculate the final distributions we would obtain were the Monte Carlo were in excellent agreement with the data.

The results from this reweighted unfolding is compared to the baseline measurement, and the difference is interpreted as the bias. Fig. [6.6] compares the baseline measurements (black circles) to those made by unfolding the \mathcal{O} distributions with the reweighted response matrix (red circles.) The ratio is shown at the bottom of each plot; the yellow band denotes the small assessed uncertainties ($\mathcal{O}(0.5\%)$ for $5 \text{ GeV} \leq p_T^{\text{jet}}$.) The lowest $p_T^{\text{jet}} = 4 \text{ GeV}$ bin is slightly higher - 1%.

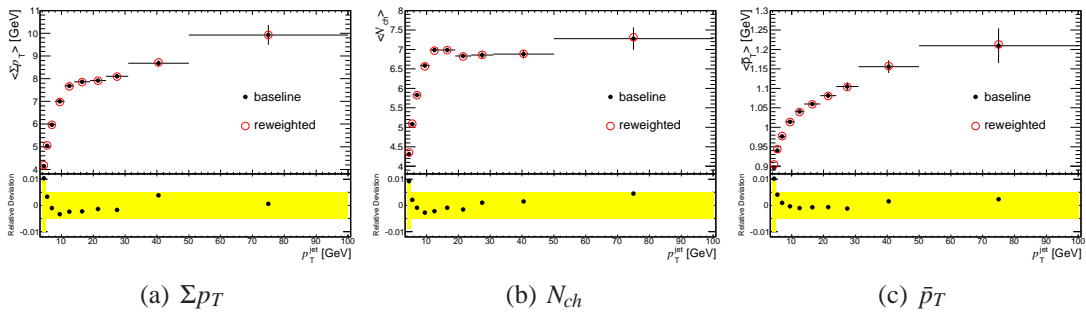


Figure 6.6: The mean values of the unfolded UE distributions. The baseline measurements (black points) are compared to the measurements (red circles) made using a response matrix constructed from PYTHIA 6 (AMBT1) that has been reweighted to reproduce the unfolded p_T^{jet} and \mathcal{O} . The ratio is shown in the bottom plot. The yellow band indicates the assessed uncertainty.

6.4 Statistical Uncertainty in the Response Matrix

The response matrix was constructed using PYTHIA 6 (AMBT1) with high statistics. Since the MC was generated with specific cuts on truth jet p_T , the high regime p_T ($20 \text{ GeV} \leq p_T^{\text{jet}}$) has a very high population, much higher than the data. For truth jet $p_T \simeq 20 \text{ GeV}$, however, the situation is reversed; the data has a higher population (approximately $2\times$) than the MC. We explore the effects of statistical uncertainty in the response matrix by using the *bootstrap* [48] method to form statistical perturbations of the baseline response matrix. We compare the baseline measurements to the results of the unfolding procedure using these alternate response matrices. The RMS width of the resulting spectrum of measurements is the associated uncertainty, and is negligibly less than 0.1%.

6.5 Misidentification of the Leading Jet

Due to track jet reconstruction efficiency and p_T resolution, the leading track jet may be matched to a non-leading charged truth jet. As a result, the direction of the UE will be incorrectly specified, leading to an uncertainty in the measurement of the UE observables.

The rate at which the leading charged truth jet fails reconstruction as the leading track jet is encoded in the p_T distributions of the leading (p_T^{jet}) and subleading jets (p_T^{sub}). As an example, assume the leading jet spectrum between two Monte Carlo (MC_A and MC_B) samples were identical, but MC_A has a harder p_T^{sub} distribution. Because of the p_T^{jet} transfer function, the subleading jet from MC_A is more likely to reconstruct as the leading jet, potentially confusing the true direction of the underlying event. However, if the leading and subleading jets were perfectly back-to-back ($\Delta\phi = \pi$), then the misidentification would have no effect on the TRANSVERSE region because of symmetry. The $\Delta\phi$ spectrum between the leading and subleading jet captures the effects of leading jet misidentification (*jet swap*).

In principle, we would consider the p_T and $\Delta\phi$ spectrum of all the jets. Due to the available statistics, we only consider the first subleading jet. The steeply falling jet multiplicity curve indicates consideration of subsubleading jets would provide 2^{nd} order corrections.

Due to differences in the p_T^{jet} , p_T^{sub} and $\Delta\phi$ distributions between the data and the PYTHIA 6 (AMBT1) used to construct the response matrix, the correction procedure may introduce a potential bias. The strategy taken to evaluate the bias due to jet swap is to create an alternate response matrix, where PYTHIA 6 (AMBT1) has been reweighted to reproduce the data p_T^{jet} , p_T^{sub} and $\Delta\phi$ distributions. The baseline measurements are compared to those obtained using the reweighted response matrix. The relative deviation is taken as the bias due to jet swap.

The procedure for Σp_T is as follows; the treatment of N_{ch} and \bar{p}_T is identical.

1. Create a new response matrix to unfold in four variables - p_T of the leading and subleading jets, $\Delta\phi$ and Σp_T . Events with only one jet are assigned $p_T \equiv 0$ for the subleading jet and $\Delta\phi \equiv 0$.
2. Unfold data using the new response matrix to obtain truth level distribution of the unfolded variables $\equiv F_{DATA}^0(p_T^{\text{jet}}, p_T^{\text{sub}}, \Delta\phi, \Sigma p_T)$.
3. Derive the equivalent distribution $F_{MC}^0(p_T^{\text{jet}}, p_T^{\text{sub}}, \Delta\phi, \Sigma p_T)$ from Pythia 6 (AMBT1) truth.
4. Normalize F_{DATA}^0 and F_{MC}^0 to unity when integrated over all variables, and form the weight $\omega \equiv F_{DATA}^0 / F_{MC}^0$.
5. Reweight PYTHIA 6 (AMBT1) by ω .
6. Create another (reweighted) response matrix using the reweighted PYTHIA 6 (AMBT1).
7. Use the reweighted response matrix to unfold data to obtain the next approximation to the correct distribution $\equiv F_{DATA}^1(p_T^{\text{jet}}, p_T^{\text{sub}}, \Delta\phi, \Sigma p_T)$.

8. Use $F_{DATA}^1(p_T^{\text{jet}}, p_T^{\text{sub}}, \Delta\phi, \Sigma p_T)$ to plot the distributions of Σp_T as functions of p_T^{jet} .
9. Plot the mean values of Σp_T as a function of p_T^{jet} .
10. Interpret the ratio of the baseline measurement to the mean values in the previous step, as the induced bias.

Steps (1) - (2), in principle, provide a central value that can be used to estimate the bias. Steps (3) - (8) attempt to correct for (a) insufficient MC statistics and (b) technical complications (memory limitations) using the RooUnfold package and ROOT, forcing us to use different binning than the baseline measurement in this analysis.

6.6 Discretization Effects

Sec. 5.4.1 discussed the issues of obtaining an unbiased mean value of a distribution, when calculating the mean values using a histogram with binned contents. For Σp_T and N_{ch} , a correction for the bias was made using cubic splines. Fig. [6.7] indicates the level of performance of these methods. By sampling various Monte Carlo samples (Pythia 6 with Z1, AMBT1, MC09, Perugia2011 tunes, Pythia 8, and Herwig ++ with UE7-2 tune), making the spline corrections (to the binned distributions) and comparing the results to the true (unbinned) mean values, we obtain a distribution of the performance index (ratio of corrected to truth) of the spline-based methods. The solid red lines indicate the mean values of the distribution (*not* the baseline measurement of this analysis); the yellow bands indicate the RMS. We take the RMS value as the uncertainty due to discretization effects for Σp_T and N_{ch} . The uncertainty is negligible for low track jet p_T and rises to 0.3% - 0.5% at high p_T . The splines definitely help to correct the bias due to discretization effects, as can be seen by comparing Fig. 6.7(a) to Fig. 6.7(b), and comparing Fig. 6.7(c) to Fig. 6.7(d).

For \bar{p}_T , no spline-based corrections were made as the binning bias was seen to be negligible. The binning allows an accurate \mathcal{O} (0.1%) calculation of the mean value (See Fig. [6.7(e)].)

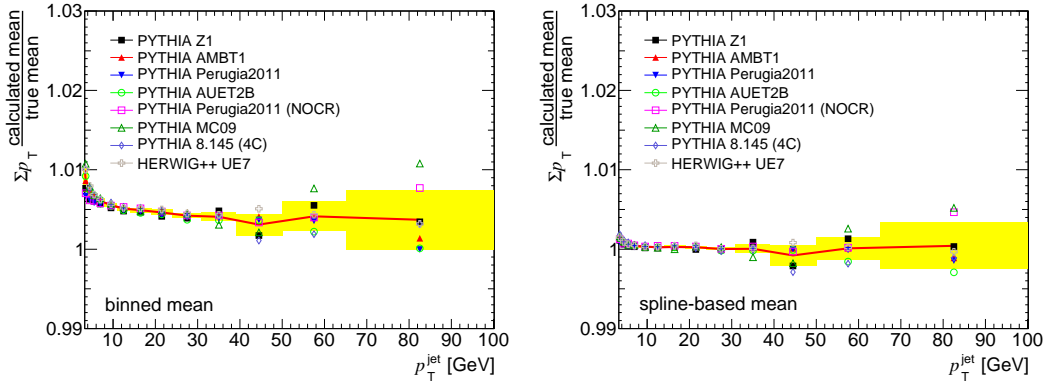
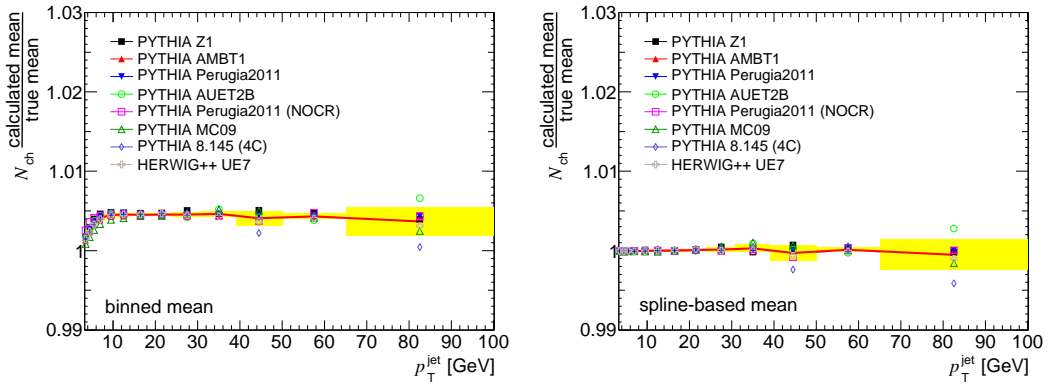
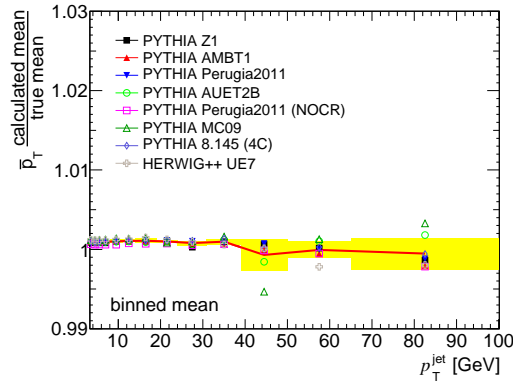
(a) binned mean Σp_T (b) spline mean Σp_T (c) binned mean N_{ch} (d) spline mean N_{ch} (e) binned mean \bar{p}_T

Figure 6.7: The ratios of calculated (binned or spline-based) mean value to the true mean value in the TRANSVERSE region. The results for various Monte Carlo samples and all jet radii are plotted together to form a scatter plot of the ratios. The red line represents the mean value of the dispersion of points; the yellow band represents 1 standard deviation.

6.7 Dependence on Number of Iterations

We use the Bayesian Iterative Unfolding algorithm to calculate the true p_T^{jet} and UE distributions, from the observed data and the response matrix. Sec. 5.2 discussed using 4 iterations in the algorithm to obtain the central values, a choice motivated by the performance of the closure tests and χ^2 values between iterations. We have no reliable method of determining the optimal number(s) of iterations. To determine the size of the uncertainty associated with this ambiguity, we repeat the analysis using different numbers of iterations, yielding a spectrum of measurements. We interpret the unweighted ² RMS (cf. Eqn. 6.4) of these measurements as the uncertainty.

$$\bar{x} = \frac{1}{N} \sum_{iter=4}^{iter=8} x_{iter} \quad (6.3)$$

$$\sigma^2 = \frac{1}{N-1} \sum_{iter=4}^{iter=8} (x_{iter} - \bar{x})^2 \quad (6.4)$$

where $N = 5$ (iterations = 4, 5, 6, 7, 8).

Fig. [6.8] shows the mean values of the UE distributions using different numbers of iterations. The bottom plot shows the (relative) difference between the mean value for each iteration number and the baseline mean value. the yellow band denotes the RMS width. We assess the uncertainty due to the choice of number of iterations as 0.5%, a small contribution.

6.8 Statistical Uncertainties

The statistical uncertainty in the measurement is provided by the RooUnfold package. This uncertainty is propagated through the unfolding procedure as outlined in [40, 49]. The statistical uncertainties in the measurements of the individual bins of \mathcal{O} vs p_T^{jet} are shown in Figs.[6.9 - 6.14].

As a cross check of the reported statistical uncertainty in the measurement, we used the bootstrap method to derive statistical perturbations of the (measured) data distributions, and unfolded them to give a spectrum of corrected distributions. We compared the dispersion (width) of the spectrum to the statistical uncertainties reported by RooUnfold. On average, RooUnfold gave uncertainties approximately 2% higher than our results. We report the errors given by RooUnfold.

²Since each point has approximately equal weight, namely its statistical uncertainty, the difference between weighted and unweighted RMS is not important.

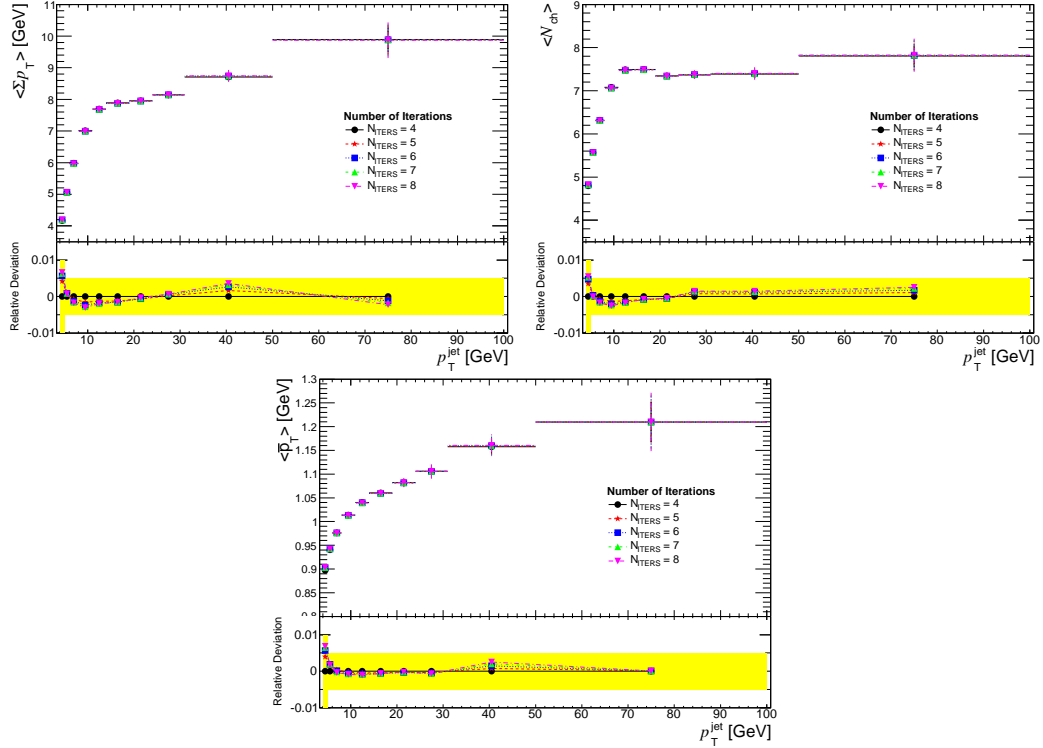


Figure 6.8: The mean values of Σp_T as a function of track jet p_T , using different number of iterations for the unfolding algorithm. The bottom plot shows the ratio of the measurements to the baseline measurement. The width of the yellow band is the dispersion (RMS) of the different measurements relative to the baseline.

6.9 Summary of Total Systematic Uncertainties

We have discussed the different sources of the systematic uncertainties and estimated their sizes. We add each of the sources in quadrature (cf. Eqn. 6.5).

$$\sigma_{\text{tot}}^2 \equiv \sum_k \sigma_k^2 \quad (6.5)$$

The uncertainties for the individual bins of \mathcal{O} and p_T^{jet} are shown in Figs.[6.9-6.12].

6.10 Consistency Checks - Refolding the Distributions

Closure tests on fully simulated Monte Carlo control samples are extremely important. Knowledge of the truth distributions allows us to calibrate our expectations of the unfolding procedure. We do not have the luxury of this knowledge for the data, but we can perform

some consistency checks that increase confidence in our results. The concept is simple; we *refold* (See Eqn. 6.6) the unfolded data with the response matrix, and compare to the measured distributions.

The refolding procedure is defined as follows:

$$y_k^{\text{reco}} = \sum_j^{N_{\text{bins}}^{\text{true}}} R_{jk} y_j^{\text{true}} \quad (6.6)$$

where \mathbf{y}^{reco} and \mathbf{y}^{true} are the observed and corrected data, respectively. Since there exists a high degree of correlation between the refolded data and observed data, it would be difficult to gauge the performance of such tests as KS and χ^2 . Comparison of the refolded data to the measured data still retain power. These consistency checks cannot tell us that we have the correct answers; but they do indicate that our results are consistent. To the extent that *all* bins in p_T^{jet} and \mathcal{O} simultaneously agree, then the corrections we have performed on the measured distributions are a feasible approximation of the true physics distributions.

The refolded data is compared with the observed data (before corrections) in Figs. [6.15 - 6.17]. The bottom plots show the ratio of the refolded data to the measured data. A value of 1 indicates perfect agreement. We see very good agreement between the observed and corrected data in the regions of high statistics. Significant deviations from unity occur in regions of low statistics, where the uncertainties are higher.

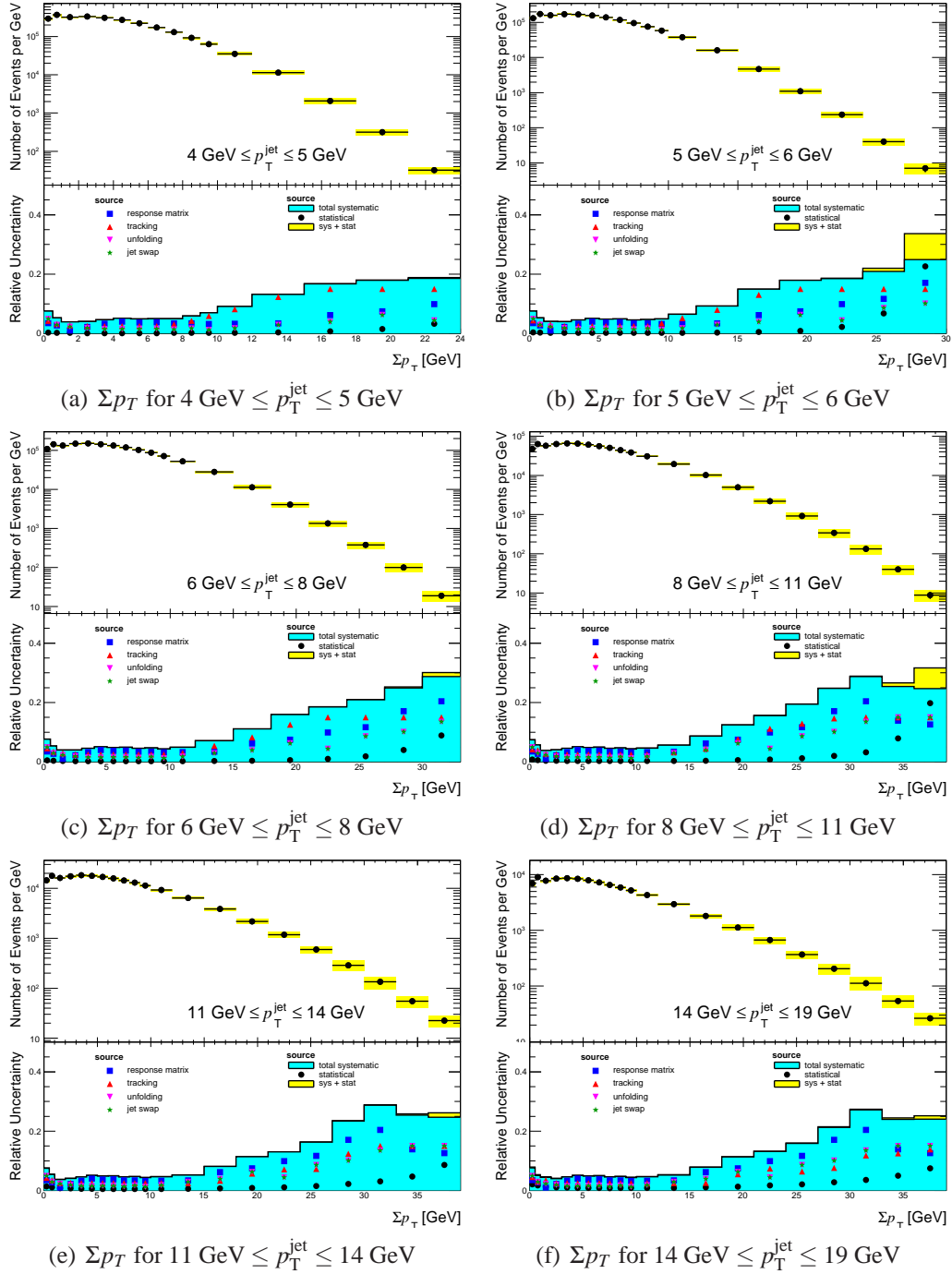


Figure 6.9: The corrected Σp_T data and uncertainties.

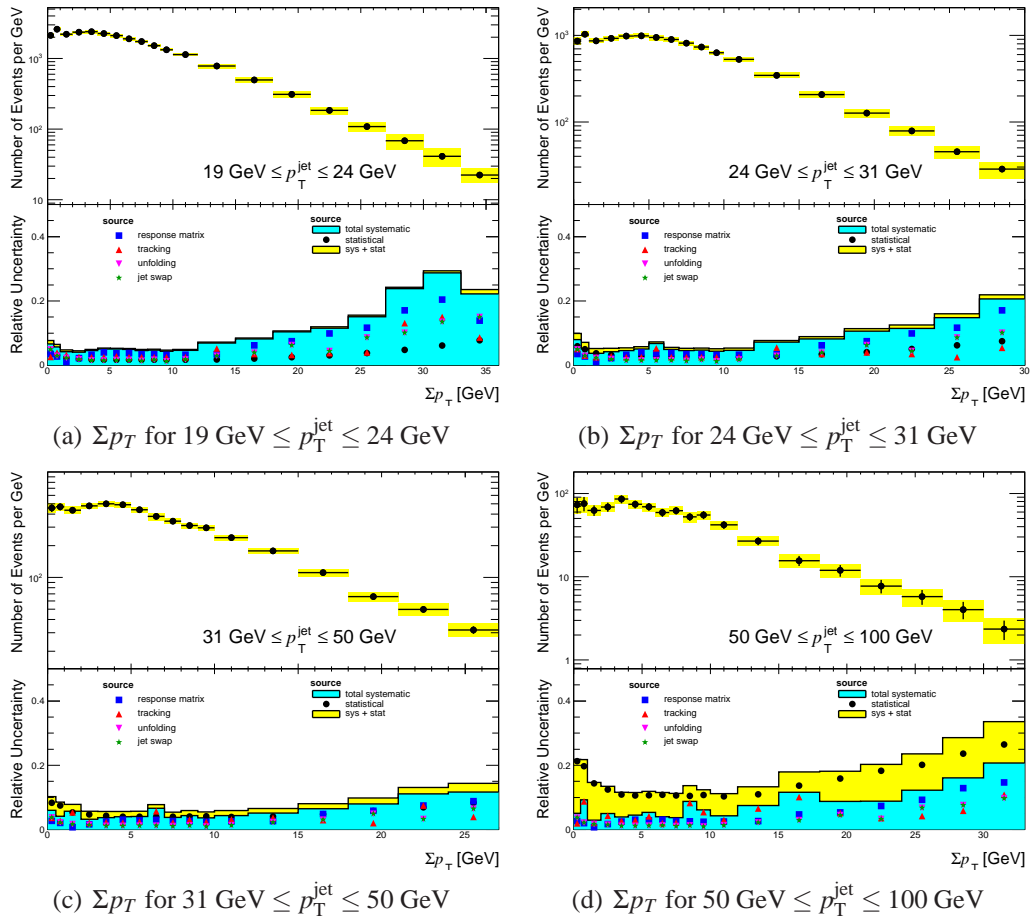


Figure 6.10: The corrected Σp_T data and uncertainties (cont.)

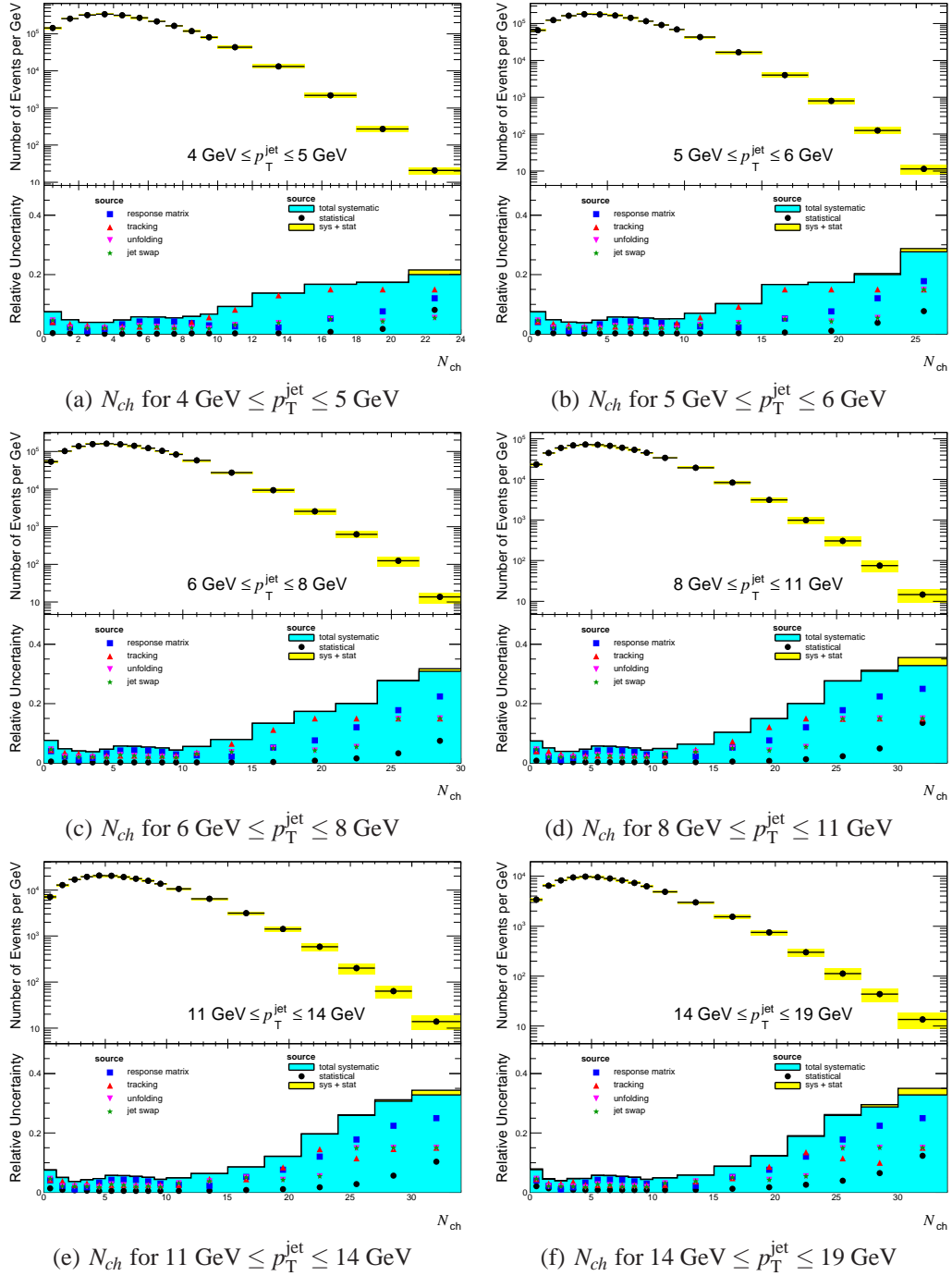


Figure 6.11: The corrected N_{ch} data and uncertainties

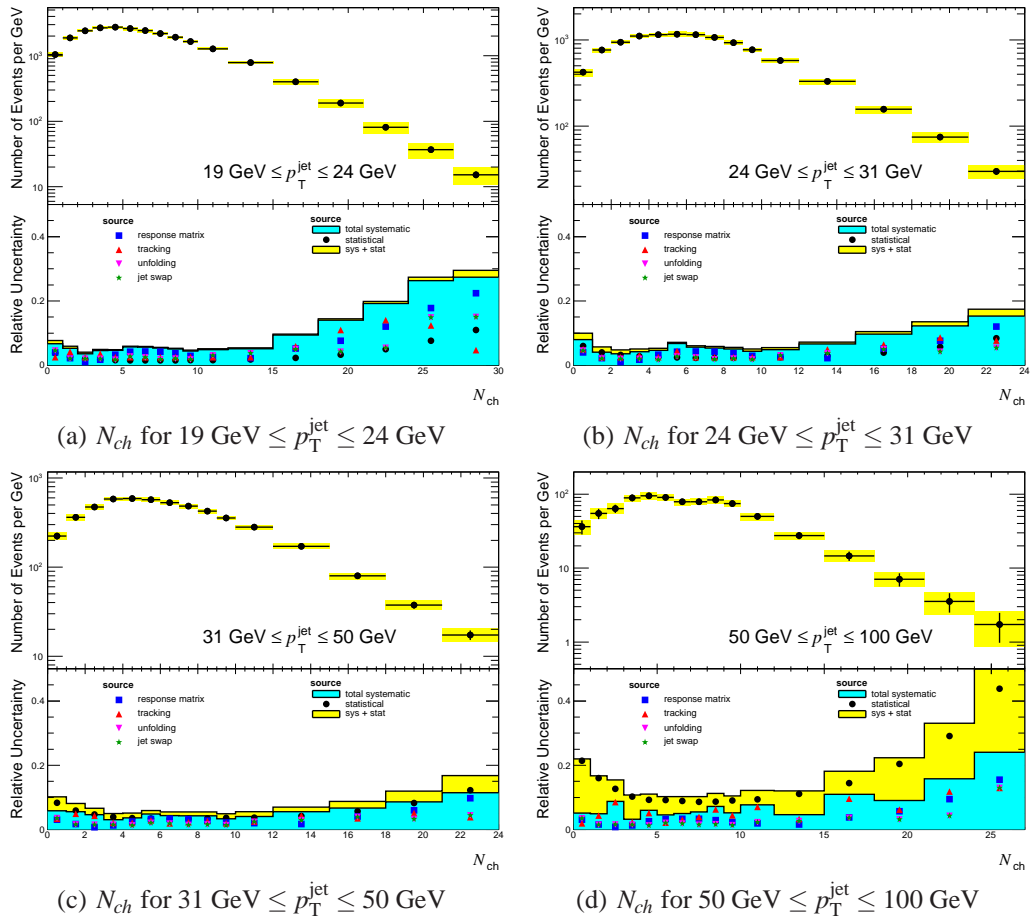
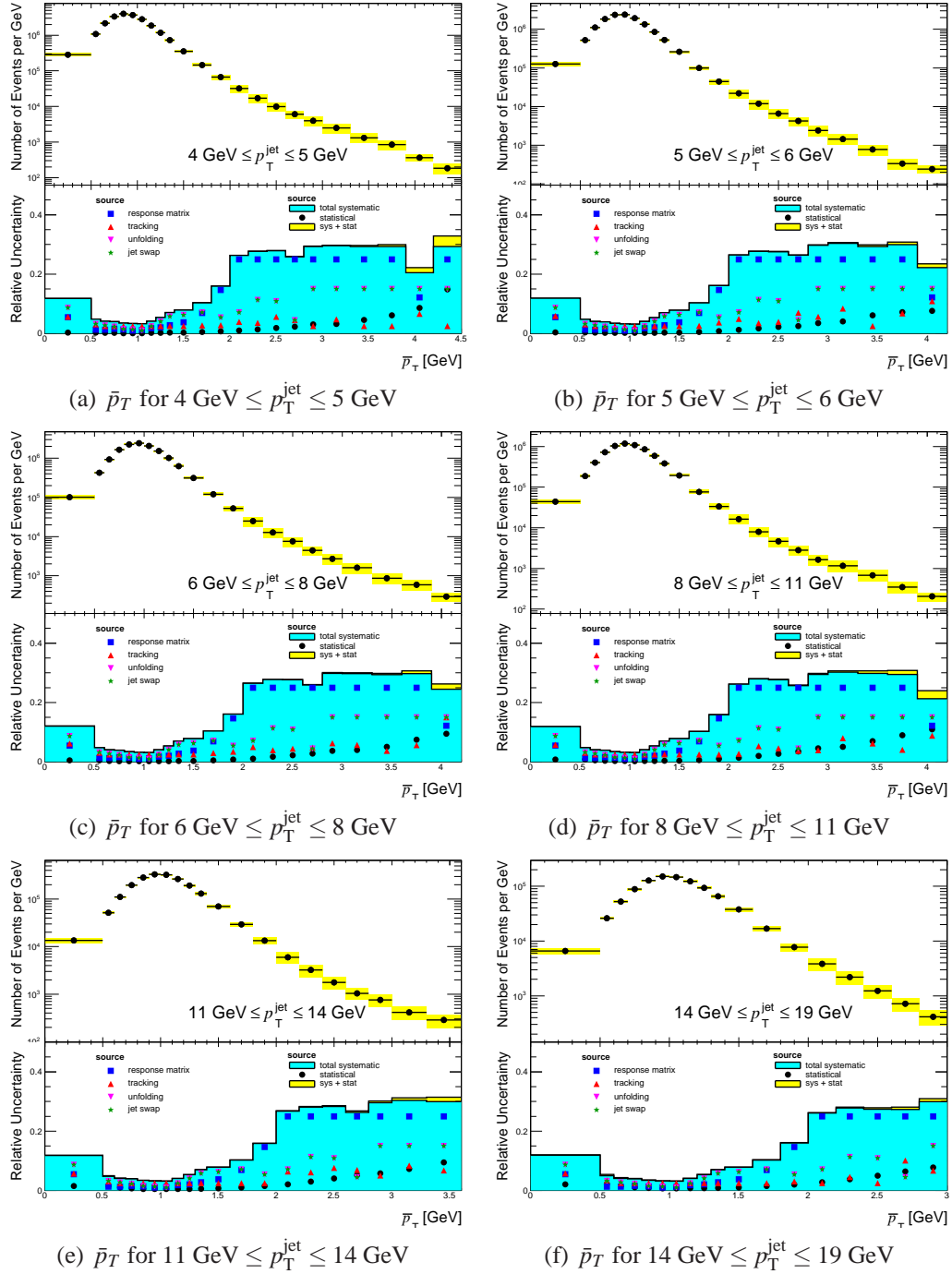


Figure 6.12: The corrected N_{ch} data and uncertainties (cont.)

Figure 6.13: The corrected \bar{p}_T data and uncertainties

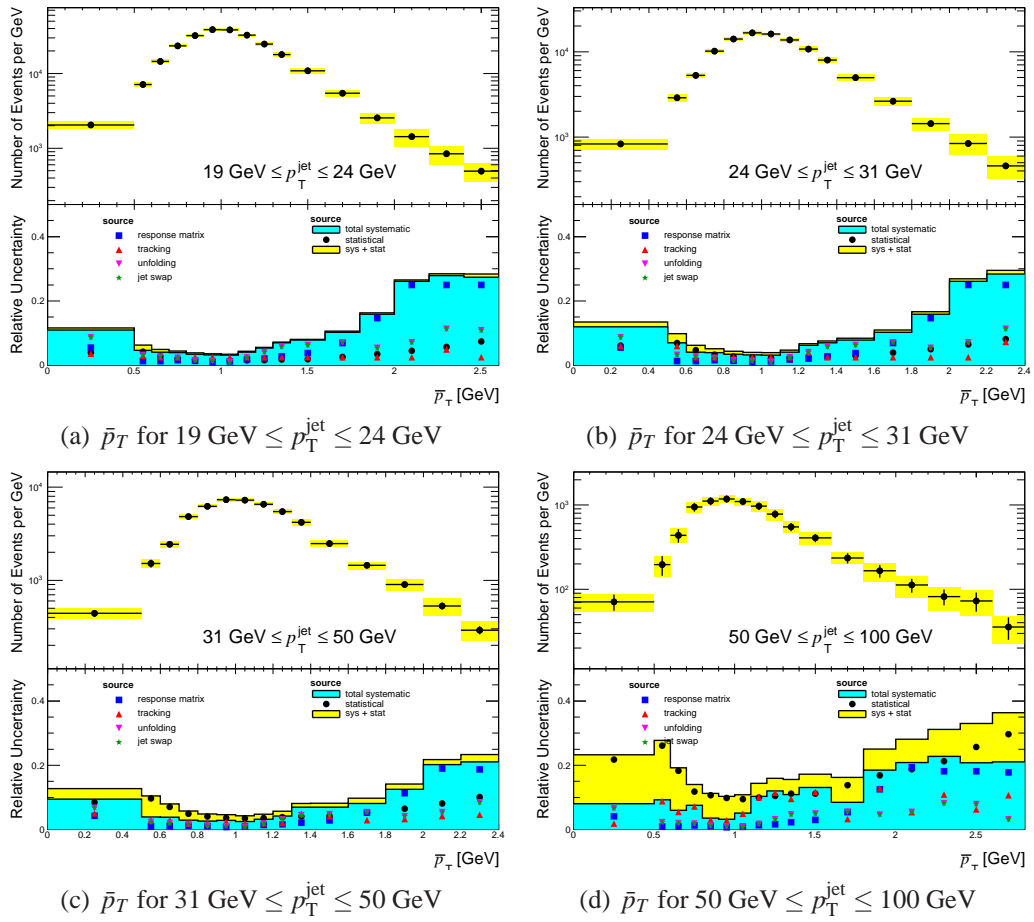
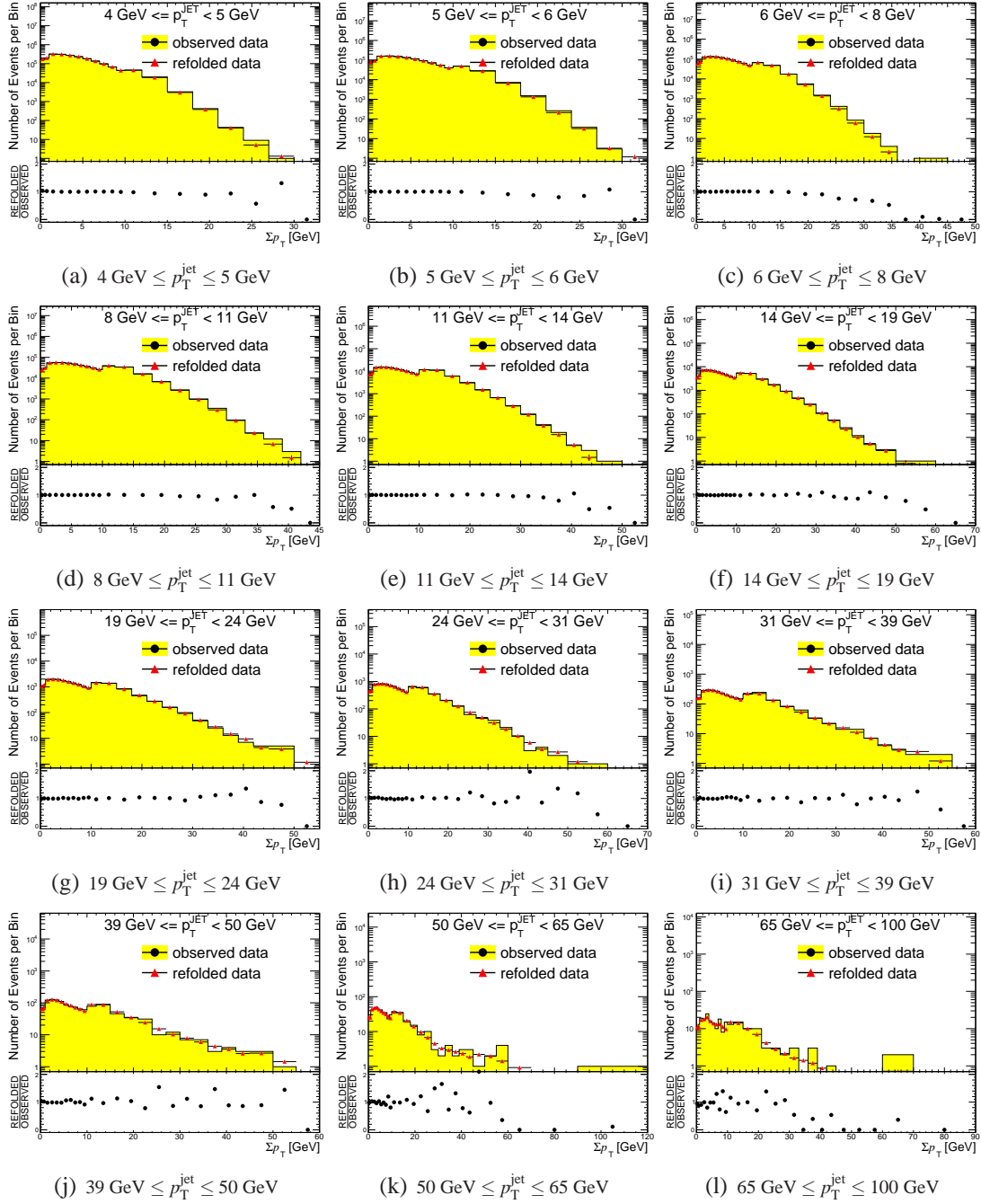
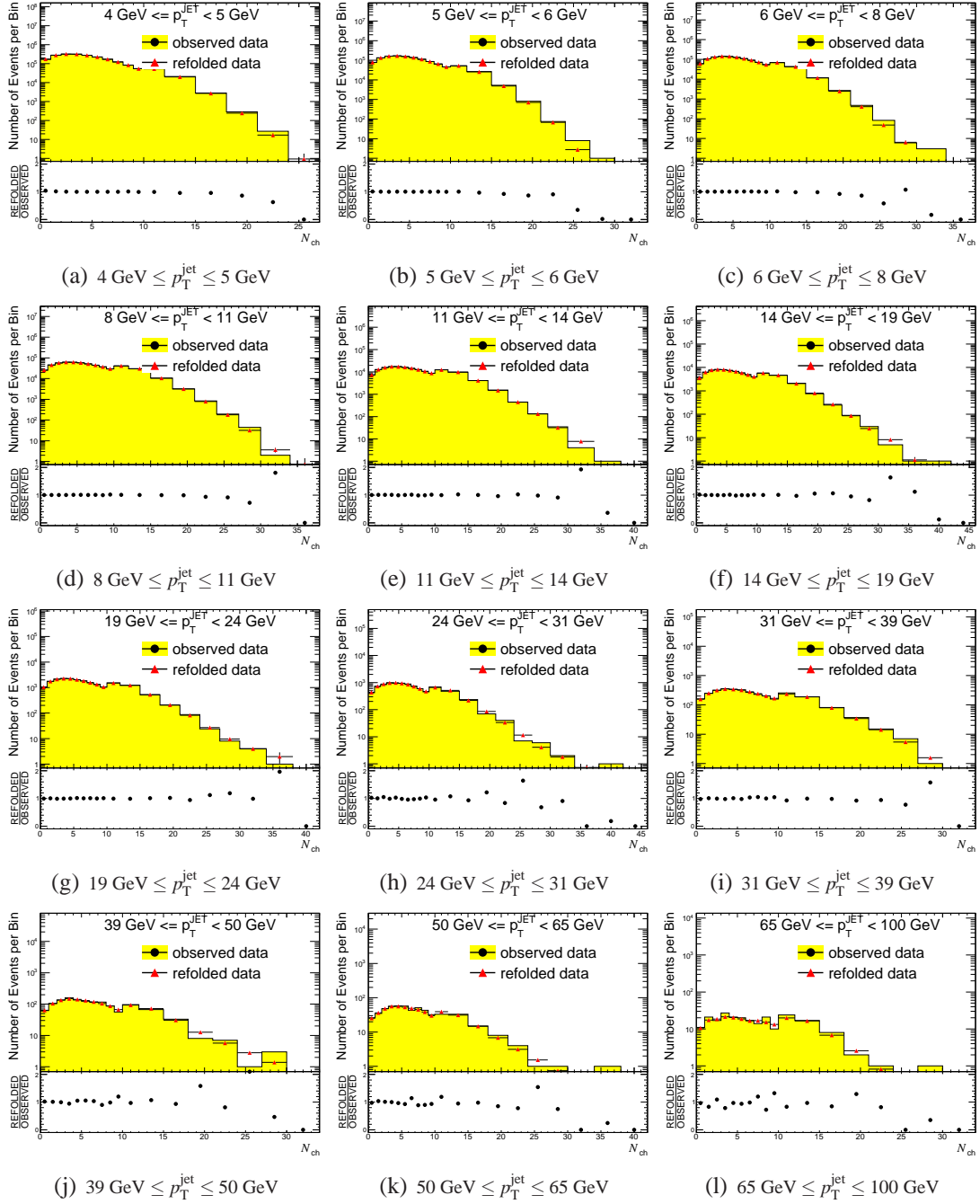
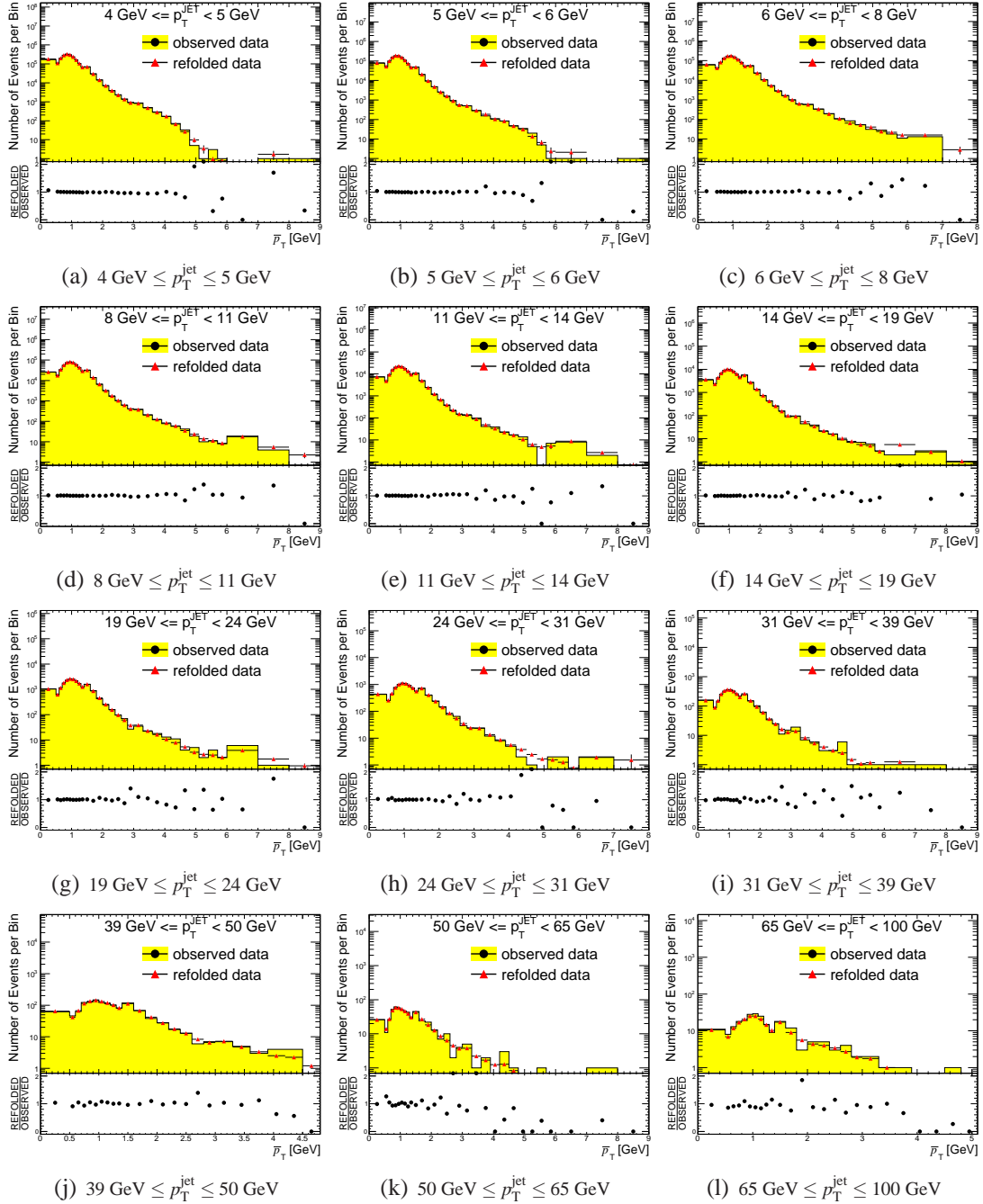


Figure 6.14: The corrected \bar{p}_T data and uncertainties (cont.)

Figure 6.15: The refolded Σp_T distributions are compared to the raw (uncorrected) data.

Figure 6.16: The refolded Σp_T distributions are compared to the raw (uncorrected) data.

Figure 6.17: The refolded \bar{p}_T distributions are compared to the raw (uncorrected) data.

Chapter 7

Conclusions

We have measured the full distributions of the UE observables (Σp_T , N_{ch} and \bar{p}_T) in the TRANSVERSE region, in slices of p_T^{jet} . Jets were constructed using the anti- k_t algorithm, using a value of 0.6 as the R-parameter, from reconstructed tracks in the ATLAS Inner Detector. The excellent tracking performance of the Inner Detector allowed us to probe very low $p_T^{\text{jet}} \geq 4$ GeV. The relative systematic uncertainties in the mean values were 2.9% - 6.5% (Σp_T), 2.7%-4.6% (N_{ch}) and 1.3%-4.1% (\bar{p}_T). PYTHIA 6 Z1 performed better than PYTHIA 6 AUET2B, although both gave good agreement with the data. The measurements presented provide another testing ground for further tuning of Monte Carlo generators. Future analyses would benefit from larger data samples, which would allow us to probe the high $p_T^{\text{jet}} \geq 50$ GeV range.

Bibliography

- [1] L. Evans (ed.), P. Bryant (ed.), JINST **3**, S08001 (2008).
- [2] M. Peskin and D. Schroeder, *An Introduction to Quantum Field Theory* (Addison-Wesley, 1995).
- [3] D. J. Gross and F. Wilczek, Phys. Rev. Lett. **30**, 1343 (1973).
- [4] V. Barger and R. Phillips, *Collider Physics* (Addison-Wesley, 1997).
- [5] T. Sjöstrand, S. Mrenna and P. Skands, JHEP **05**, 026 (2006).
- [6] T. Gleisberg *et al.*, JHEP **0902**, 007 (2009), 0811.4622.
- [7] S. Gieseke *et al.*, Herwig++ 2.5 Release Note, arXiv.1102.1672, 2011.
- [8] D. Acosta *et al.*, Phys. Rev. **D70**, 072002 (2004), hep-ex/0404004.
- [9] T. Affolder *et al.*, Phys. Rev. **D65**, 092002 (2002).
- [10] M. Cacciari, G. P. Salam and G. Soyez, JHEP **0804**, 063 (2008), 0802.1189.
- [11] G. P. Salam and G. Soyez, JHEP **0705**, 086 (2007), 0704.0292.
- [12] T. Sjöstrand and M. van Zijl, Phys. Rev. **D36**, 2019 (1987).
- [13] S. Chatrchyan *et al.*, JHEP **09**, 109 (2011), 1107.0330.
- [14] ATLAS Collaboration, Phys. Rev. D **83**, 112001 (2011).
- [15] B. P. Kersevan and E. Richter-Was, (2004), hep-ph/0405247.
- [16] M. L. Mangano, M. Moretti, F. Piccinini, R. Pittau and A. D. Polosa, JHEP **0307**, 001 (2003), hep-ph/0206293.
- [17] ATLAS Collaboration, ATLAS Monte Carlo Tunes for MC09, ATLAS-PHYS-PUB-2010-002.
- [18] ATLAS Collaboration, ATL-COM-PHYS-2010-267.

- [19] P. Z. Skands, Phys. Rev. **D82**, 074018 (2010), 1005.3457.
- [20] T. Sjöstrand, S. Mrenna and P. Z. Skands, Comput. Phys. Commun. **178**, 852 (2008).
- [21] R. Field, (2010), 1010.3558, Invited talk at HCP2010, Toronto, August 23, 2010.
- [22] ATLAS Collaboration, ATLAS-CONF-2011-009.
- [23] S. Agostinelli *et al.*, Nucl. Instrum. Meth. **A506**, 250 (2003).
- [24] CERN-DI-0812015.
- [25] ATLAS Collaboration, JINST **3**, S08003 (2008).
- [26] G. Aad *et al.*, JINST **3**, P07007 (2008).
- [27] A. Abdesselam *et al.*, Nucl.Instrum.Meth. **A568**, 642 (2006).
- [28] E. Abat *et al.*, JINST **3**, P02014 (2008).
- [29] E. Abat *et al.*, JINST **3**, P10003 (2008).
- [30] C. Ohm and T. Pauly, Nucl. Instrum. Meth. **A623**, 558 (2010), 0905.3648.
- [31] ATLAS Collaboration, Phys. Rev. D **84** (2011).
- [32] T. Cornelissen *et al.*, J. Phys. Conf. Ser. **119**, 032014 (2008).
- [33] R. Frühwirth *et al.*, Nucl. Instr. Meth. **A262** (1987).
- [34] H. Gray, *The Charged Particle Multiplicity at Center of Mass Energies from 900 GeV to 7 TeV measured with the ATLAS Experiment at the Large Hadron Collider*, PhD thesis, 2010.
- [35] R. Frühwirth, W. Waltenberger and P. Vanlaer, J. Phys. **G34**, N343 (2007).
- [36] G. Piacquadio, K. Prokofiev and A. Wildauer, J. Phys. Conf. Ser. **119**, 032033 (2008).
- [37] ATLAS, Charged particle multiplicities in pp interactions at $\sqrt{s} = 7$ TeV measured with the ATLAS detector at the LHC, ATLAS-CONF-2010-024.
- [38] ATLAS Collaboration, ATLAS-CONF-2011-1678.
- [39] ATLAS Collaboration, ATLAS-CONF-2011-408.
- [40] G. D'Agostini, Nucl.Instrum.Meth. **A362**, 487 (1995).
- [41] T. Adye, <http://hepunix.rl.ac.uk/~adye/software/unfold/RooUnfold.html>.

-
- [42] A. Hocker and V. Kartvelishvili, Nucl.Instrum.Meth. **A372**, 469 (1996), hep-ph/9509307.
- [43] V. Anikeev, A. Spiridonov and V. Zhigunov, Nucl.Instrum.Meth. **A322**, 280 (1992).
- [44] N. D. Gagunashvili, (2010), 1004.2006.
- [45] M. Shapiro and I. Hinchliffe, Measurement of the jet fragmentation function and transverse profile in proton-proton collisions at a center-of-mass energy of 7 TeV with the ATLAS detector at the LHC, ATL-COM-PHYS-2011-923.
- [46] ATLAS Collaboration, ATLAS-CONF-2010-069.
- [47] G. Aad *et al.*, Track Reconstruction Efficiency in $\sqrt{s} = 7$ TeV Data for Tracks with $p_T > 100$ MeV, ATL-COM-PHYS-2010-682.
- [48] B. Efron, Annals of Statistics **7**, 1 (1979).
- [49] G. D'Agostini, Improved iterative Bayesian unfolding,
<http://roma1.infn.it/~dagos/unf2.pdf>.

Appendix A

Track Quality

The track selection criteria for tracks used in this analysis are plotted in Fig. [A.1-A.2]. The black curves indicate the tracks which have passed all other selection criteria; the yellow filled areas depict tracks that also pass the corresponding selection criteria. Out of a total of 407.9 million tracks that pass all other selection criteria, 3.8 million fail the B-Layer requirement. The plots that follow are so-called “N-1” plots, where the indicated variable is shown for tracks that have passed all other selection criteria. The tracks passing the selection criteria for that variable is shown in yellow.

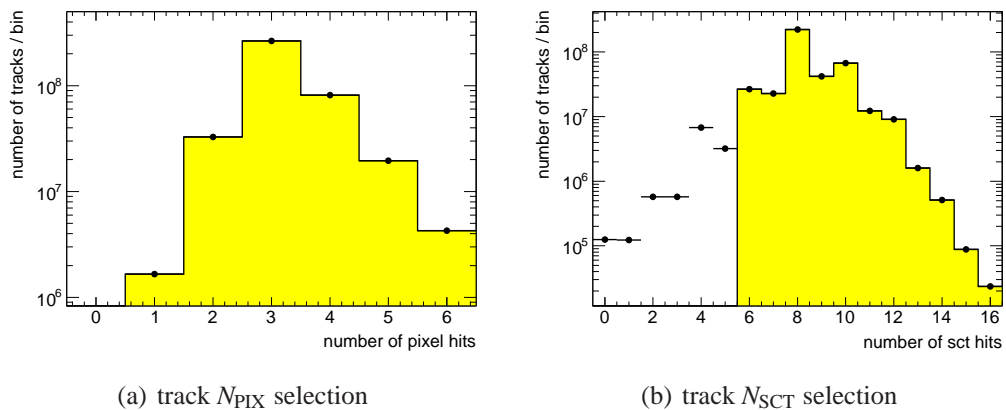


Figure A.1: The “N-1” distributions for the ID track selection criteria. Tracks passing all other selection criteria are plotted in black. The yellow area corresponds to the tracks that also pass the indicated cut.

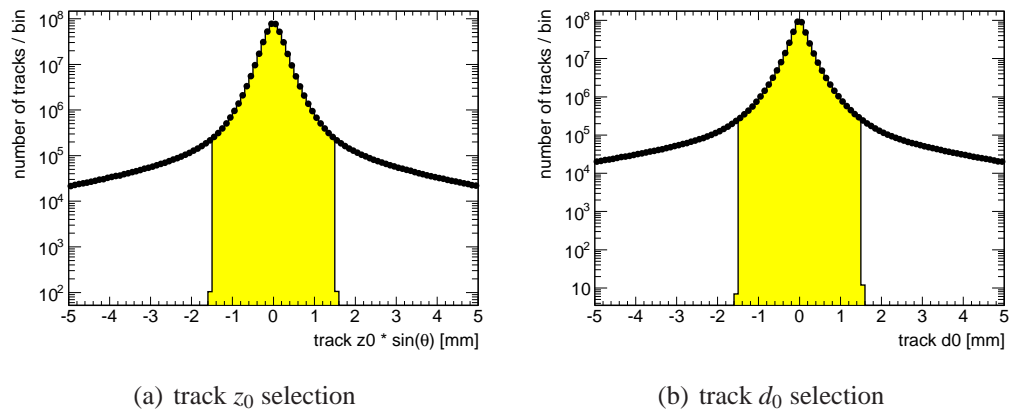


Figure A.2: The “N-1” distributions for the z_0 and d_0 track selection criteria. Tracks passing all other selection criteria are plotted in black. The yellow area corresponds to the tracks that also pass the indicated cut.

Appendix B

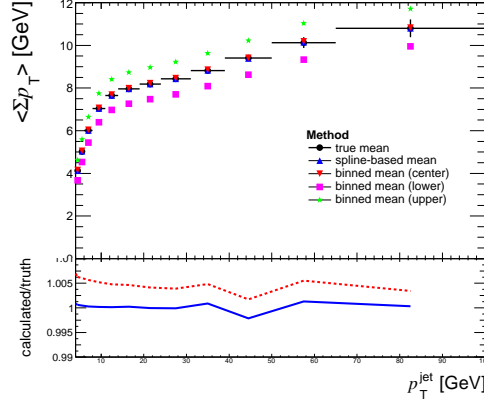
Discretization Effects

For measured (uncorrected) data, it is straightforward to calculate the true mean values of the UE distributions because we have access to the variables on an event-by-event basis. The correction process, described in Sec. 5.2, requires *discretization* (binning) of the variables in question. After the correction procedure, we no longer have access to the variables on a per-event basis, but rather only to the number of events in bins of the given variables. This situation prevents a straightforward calculation of the true mean value of the corrected variables. One can decide to use the bin center as the representative abscissa of the bin. This choice may be appropriate for distributions that populate the bin uniformly, but this is not our situation. Most of the distributions in this analysis have an exponential dependence on the abscissa, which could lead to significant differences between the mean values calculated using different representative abscissae. For example, in Fig. [B.1], we calculate the mean value using the lower edge of the bin, the bin center, and the upper edge. We see differences $\mathcal{O}(10\%)$ between the different calculations. Using the center of the bin appears to be the best choice, but we have no guarantee this condition will persist with different distributions. We propose a method to overcome this obstacle using numerical methods based on cubic splines. Each UE observable (Σp_T , N_{ch} and \bar{p}_T) is treated differently in the following sections.

B.1 Discretization Effects in Σp_T

The binning used to analyze Σp_T , being of finite width, introduces an uncertainty in the x-value (abscissa) when calculating the mean value. For example, a bin with coordinates $4\text{GeV} \leq \Sigma p_T \leq 5\text{GeV}$ lumps events having $\Sigma p_T = 4.1\text{GeV}$ together with events having $\Sigma p_T = 4.9\text{GeV}$. When calculating a *binned* mean value, as in Eq. B.1, all the events in that bin contribute equally to the mean value, although there are many more events with $\Sigma p_T = 4.1\text{GeV}$ than there are with $\Sigma p_T = 4.9\text{GeV}$. We have lost information, i.e. - introduced an uncertainty, pertaining to the abscissa of the bin. The uncertainty due to discretization diminishes as the bins become smaller. Due to lack of statistics, we are not able to make

sufficiently small bins in the high p_T region. Fig. [B.1] indicates the size of the uncertainty due to discretization, comparing the binned mean value (red graph) (cf. Eq. B.1) to the true mean value (black graph). Fig. [B.1] compares the binned and true mean values for uncorrected data and truth-level Monte Carlo (Pythia 6 with AMBT1 tune) distributions. The failure (binned vs. unbinned) is $\mathcal{O}(1\%)$ throughout the track jet p_T spectrum.



(a) Pythia 6 (Z1) truth Σp_T in the XVERSE region

Figure B.1: Comparison of the Σp_T mean values to the mean values obtained via binned calculations and a spline-based calculation, as a function of track jet p_T . The binned calculations are performed using the lower edge, center and upper edge of the bins. The plots on the bottom are the ratios of the predicted mean values to the known mean values.

$$\mu_{\text{binned}} = \frac{\sum_{k=2}^N n_k x_k}{\sum_{k=1}^N n_k} \quad (\text{B.1})$$

$$\sigma_{\text{binned}}^2 = \frac{\sum_k n_k (x_k - \mu)^2}{\sum_k n_k} \quad (\text{B.2})$$

where $\{n_k\}$ and $\{x_k\}$ are the histogram contents and bin centers, respectively. The reader will notice that the summations in the numerators in Eqs. [B.1-B.2] start at 2, ignoring the lowest bin. The first bin, spanning $0 \leq \Sigma p_T \leq 0.5 \text{ GeV}$, is populated by events having no tracks in the TRANSVERSE region, for which Σp_T is identically 0. There is no uncertainty in the abscissa of this bin. The summations which include a factor of x_k can safely omit the first bin, as it identically contributes 0. We will have more to say on this subject shortly, referring to this as the "0-bin" effect.

We now propose a method to compensate for the loss of resolution due to discretization. The procedure involves fitting a cubic spline to the *integral* of Σp_T . Fitting directly to the variable is not a well-defined procedure because, as was mentioned above, the abscissa of the bin is not well-defined. Knowledge of the number of events are in a bin, however, is

equivalent to knowledge of the integral of the function, over the bin coordinates. Furthermore, the integral of the function may be a more suitable quantity for calculating the mean value. To see this, let $F(s)$ be the antiderivative of $N(s)$, where $s \equiv \Sigma p_T$ and $N(s)$ is the number of events with $\Sigma p_T = s$.

$$\frac{dF(s)}{ds} = N(s) \quad (\text{B.3})$$

The rule of integration by parts gives the following formulae for the true mean and standard deviation:

$$\mu_{TRUE} \equiv \frac{\int_a^b s N(s) ds}{\int_a^b N(s) ds} = \frac{\int_a^b s F'(s) ds}{\int_a^b F'(s) ds} = \frac{b F(b) - a F(a) - \int_a^b F(s) ds}{F(b) - F(a)} \quad (\text{B.4})$$

$$\begin{aligned} \sigma_{TRUE}^2 &\equiv \frac{\int_a^b (s - \mu)^2 N(s) ds}{\int_a^b N(s) ds} = \frac{\int_a^b (s - \mu)^2 F'(s) ds}{\int_a^b F'(s) ds} \\ &= \frac{b^2 F(b) - a^2 F(a) - 2 \int_a^b s F(s) ds}{F(b) - F(a)} - \mu^2 \quad (\text{B.5}) \end{aligned}$$

The final results on the right hand side (RHS) depend on $F(s)$ and its integral, but not on $F'(s) = N(s)$ itself.

The exact process for the construction of the cubic spline is outlined next:

1. Given the Σp_T histogram with N bins with contents n_k , for $k = 1, 2, 3, \dots, N$, define N pairs of (x, y) -values. As was discussed earlier in this section, the "0-bin" (covering $\Sigma p_T = 0$), is omitted in these calculations. We define:
 - $x[k] =$ upper edge of k^{th} bin ($k = 1, \dots, N$)
 - $y[1] = 0$
 - $y[k] = \sum_{i=2}^k n_i$ ($k = 2, \dots, N$)
2. Fit a standard cubic spline, with the N (x, y) -pairs as *knots*.¹
3. The derivative of the cubic spline at x_N is specified as the contents of the overflow bin.

The total number of events n_{TOT} must include the contents of the "0-bin", retaining consistency with the definition of the mean value of a distribution.

$$n_{TOT} \equiv \int_0^\infty \Sigma p_T dp_T = \int_0^\infty N_{ch} dp_T = \int_0^\infty \bar{p}_T dp_T \quad (\text{B.6})$$

¹In the context of splines, a knot is a point where the value of the function is specified.

In summary, the mean values of the Σp_T and \bar{p}_T are calculated in Eqs. [B.7 - B.8], invoking Eq. B.4. Note the different limits of integration for each variable. The lower limits are chosen to exclude the "0-bin" discussed above. The upper limits are practical limits set by the available data and Monte Carlo statistics.

$$\langle \Sigma p_T \rangle \equiv \frac{1}{n_{TOT}} \int_{0.5\text{GeV}}^{120\text{GeV}} s N(s) ds \quad (\text{B.7})$$

$$\langle \bar{p}_T \rangle \equiv \frac{1}{n_{TOT}} \int_{0.5\text{GeV}}^{9\text{GeV}} \bar{p}_T N(\bar{p}_T) d\bar{p}_T \quad (\text{B.8})$$

Fig. [B.1] compares the true mean values (black graph) to the mean values obtained via the spline-based methods (blue graph) described in this section, summarized in Eqs.[B.7 - B.8]. The uncertainties associated with these cubic spline-based methods are almost negligible, as will be further discussed in Sec. B.4.

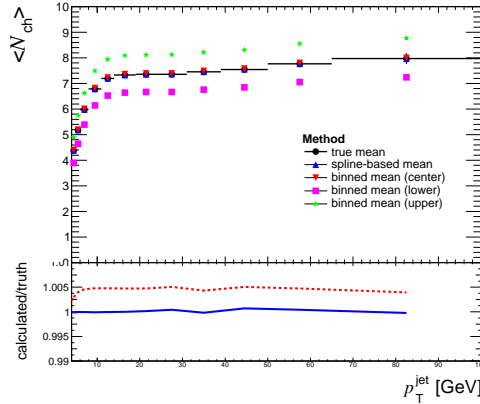
B.2 Discretization Effects in N_{ch}

The distributions of N_{ch} are treated differently than those of Σp_T and \bar{p}_T described in the previous section. N_{ch} , the number of tracks, is a discrete variable (a non-negative integer), whereas Σp_T and \bar{p}_T are continuous variables. For the bins numbered sequentially from 0, incrementing by 1, the abscissae of the bins are exactly known. After a certain point, it becomes statistically unfeasible to continue to increment the bins by 1. However, the abscissae of bins spanning two or more integers (such as $8 \leq N_{ch} \leq 11$) is not well-defined. Fig. [B.2] compares the mean values to the binned mean values. The difference is a bias of $\mathcal{O}(0.5\%)$. We correct this bias using the spline to evaluate $N(n)$ at each of the integers, as follows. Just as for Σp_T and \bar{p}_T , described in Sec. B.1, we fit a cubic spline $F(n)$ to the integral of $N(n)$, where $N(n)$ is the number of events with $N_{ch} = n$. The discrete "derivative" of the spline $F(n+1) - F(n)$ approximates $N(N_{ch})$ at each of the integers, $N_{ch} = n = 0, 1, \dots$.² The mean value of N_{ch} is calculated as

$$\langle N_{ch} \rangle \equiv \frac{1}{n_{TOT}} \sum_{k=1}^{60} k \times (F(k+1) - F(k)) \quad (\text{B.9})$$

The upper limit of $N_{ch} = 60$ is determined by the available statistics.

²The derivative is identically $N(N_{ch})$ at the knots.



(a) N_{ch} for Pythia 6 (Z1) in the TRANSVERSE region

Figure B.2: Comparison of the N_{ch} mean values to the mean values obtained via binned calculations and a spline-based calculation, as a function of track jet p_T . The binned calculations are performed using the lower edge, center and upper edge of the bins. The plots on the bottom are the ratios of the predicted mean values to the known mean values.

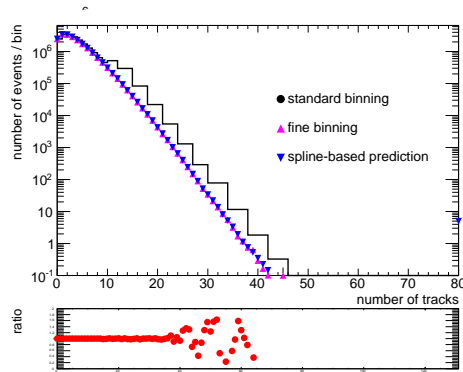
B.3 Discretization Effects in \bar{p}_T

For \bar{p}_T , the mean value of the track p_T , has sufficiently fine binning that the difference between the true and binned mean values is acceptable. See Fig. [6.7(e)]. No further corrections are performed.

B.4 Validation of the Spline-based Methods

Splines are extremely powerful tools, but must be used with caution to avoid fluctuations in the solutions, associated with the rapidly falling spectra³ we are attempting to describe. In Fig. [B.3], the distributions obtained using the standard (baseline) binning are compared to the distributions obtained using a fine binning. For the purpose of evaluating the performance of the spline-based predictions, we examine both measured (uncorrected) data distributions and truth-level Monte Carlo distributions. The spline-based predictions for the contents of the fine bins are compared to the known contents. The predictions for Σp_T are extremely good (within 0.01%) for $\Sigma p_T \leq 40\text{GeV}$, after which fluctuations cause the predictions to fail. Similarly, for $N_{ch} \leq 25$, the predictions are extremely good. The performance diminishes substantially for $N_{ch} \geq 25$. These bins contribute negligibly to the mean and standard deviations.

³The errors in polynomial interpolations tend to be proportional to the n^{th} derivative of the approximated functions.



(a) Pythia 6 (AMBT1) N_{ch} in the TRANSVERSE region

Figure B.3: Comparison of the distributions of the UE observables to the predictions made by the spline-based methods. The standard binning refers to the bins used to obtain the central values in this analysis. The plots on the bottom are the ratios of the predicted values to the known values in the fine bins.