

## On statistical aspects of Qjets

---

Stephen D. Ellis,<sup>a</sup> Andrew Hornig,<sup>b</sup> David Krohn<sup>c</sup> and Tuhin S. Roy,<sup>b,d</sup>

<sup>a</sup>*Department of Physics, University of Washington,  
Seattle WA, 98195, U.S.A.*

<sup>b</sup>*Theoretical Division T-2, Los Alamos National Laboratory,  
Los Alamos, NM, 87544, U.S.A.*

<sup>c</sup>*Department of Physics, Harvard University,  
Cambridge MA, 02138, U.S.A.*

<sup>d</sup>*Department of Theoretical Physics, Tata Institute of Fundamental Research,  
Mumbai 400005, India*

*E-mail:* [sdellis@u.washington.edu](mailto:sdellis@u.washington.edu), [ahornig@lanl.gov](mailto:ahornig@lanl.gov),  
[dkrohn@physics.harvard.edu](mailto:dkrohn@physics.harvard.edu), [tuhinsroy@lanl.gov](mailto:tuhinsroy@lanl.gov)

**ABSTRACT:** The process by which jet algorithms construct jets and subjects is inherently ambiguous and equally well motivated algorithms often return very different answers. The Qjets procedure was introduced by the authors to account for this ambiguity by considering many reconstructions of a jet at once, allowing one to assign a weight to each interpretation of the jet. Employing these weighted interpretations leads to an improvement in the statistical stability of many measurements. Here we explore in detail the statistical properties of these sets of weighted measurements and demonstrate how they can be used to improve the reach of jet-based studies.

**KEYWORDS:** QCD Phenomenology, Jets

**ARXIV EPRINT:** [1409.6785v1](https://arxiv.org/abs/1409.6785v1)

---

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Statistical uncertainties</b>	<b>4</b>
2.1	Cross-section measurement	6
2.2	Mass measurement	9
<b>3</b>	<b>Review of Qjets</b>	<b>12</b>
<b>4</b>	<b>Results from phenomenological studies</b>	<b>15</b>
<b>5</b>	<b>Understanding the statistical effects</b>	<b>18</b>
<b>6</b>	<b>Understanding the physics effects</b>	<b>22</b>
<b>7</b>	<b>Conclusions</b>	<b>25</b>
<b>A</b>	<b>Validation of section 2 with pseudo-experiments</b>	<b>27</b>
<b>B</b>	<b>Jet mass</b>	<b>27</b>

---

## 1 Introduction

Jets arise in many important scattering processes encountered at the LHC. Therefore, much experimental and theoretical effort has recently gone into creating better tools for handling them. Techniques now exist to identify jets arising from the decay of boosted heavy particles [1–21], to remove unwanted radiation from jets [4, 22–34], and to measure properties of the partons initiating jets [35–40]. See refs. [41–43] for an overview. To this toolkit the authors recently added Qjets [44]. In the discussion that follows we are specifically thinking in terms of “tagging” jets as either containing (the decay products of) a heavy boosted object (the signal), or as being an ordinary QCD-jet (the background).

The motivation behind Qjets comes from the observation that as jets are produced through a stochastic process there is an inherent ambiguity in their reconstruction. That is, even with a perfect algorithm one could never hope to unambiguously associate each hadron to an individual jet — instead one typically makes a best-guess assignment using a well motivated procedure. This is not ideal as it removes all information about the ambiguity in jet processing and tagging; any two jets that pass a set of selection cuts are assigned the same weight, even if one is unambiguously signal-like and the other is only marginally so. To address this concern the Qjets procedure processes and tags a jet using a range of plausible algorithms and grooming procedures, assigning a *distribution* of possible

properties to each jet. The initial Qjets description [44] presented two central ideas: (i) a new observable volatility that characterizes the *width* of the mass distribution generated by the Qjets procedure and can help distinguish jets arising from boosted heavy objects from QCD jets; and (ii) the use of the Qjets distributions to improve the statistical stability of the measurements of jet observables. The former is an intuitively reasonable result in the sense that one expects that a jet with an underlying mass scale (i.e., the mass of the heavy object) will exhibit a jet mass that is more robust under changes in the details of the jet algorithm and grooming procedure compared to a background QCD jet. Volatility as a discriminating variable has recently been validated [45, 46] by both the ATLAS and the CMS collaborations of the LHC. The Qjets improvement in the statistical behavior of jet measurements is less intuitive and the current work has the goal of explaining the how and why of this statistical improvement.

In order to explain why the Qjets procedure is associated with non-standard statistical analyses, let us first distinguish it from a conventional, or “classical” approach, in which a jet is first groomed and then tagged to be a signal jet if its groomed mass falls within a pre-defined signal-mass window. Such a conventional approach therefore assigns both a groomed mass  $\mu_j^C$  and a tagging efficiency  $\tau_j^C$  to each jet  $j$ . The conventional tagging efficiency is a *binary* tagging variable, which takes the value 1, if the mass of the jet is within the mass-window ( $\Omega$ ),  $\mu_j^C \in \Omega$ , and takes the value 0 if the mass of the jet is not in the window,  $\mu_j^C \notin \Omega$ . For Qjets there is a well defined procedure (reviewed in more detail in section 3) to groom an individual jet in a variety of ways leading a *distribution* of groomed masses. The corresponding Qjets tagging efficiency  $\tau_j^Q$  is the fraction of those masses that fall within the mass-window, while the Qjets measure of the jet mass  $\mu_j^Q$  is the average of the masses that fall within the mass-window. The fundamental difference in the statistical analysis of the Qjets case arises from the fact that  $\tau_j^Q$  exhibits a *continuous* range of values in the interval  $[0, 1]$ , in contrast to the binary values of  $\tau_j^C$ .

To illustrate the unconventional features of a continuous weight  $\tau_j^Q$  more specifically, consider the goal of identifying boosted  $W$ -jets. A binary  $\tau_j^C$  implies a jet is *either*  $W$ -like or QCD-like, whereas a continuous  $\tau_j^Q$  allows a jet to be treated as partially  $W$ -like *and* partially QCD-like. Now consider an example where in an experiment the conventional approach identifies two jets with masses  $\mu_1^C = 80$  GeV and  $\mu_2^C = 85$  GeV with  $\tau_1^C = \tau_2^C = 1$ . One therefore reports that the experiment sees 2 tagged  $W$ -jets and measures the masses of the tagged jets to be  $(80 + 85) / 2$  GeV = 82.5 GeV. Contrast that result with the Qjets procedure that might assign these two jets the same masses as the conventional approach (i.e.,  $\mu_j^Q = \mu_j^C$ ), but finds one jet to be more  $W$ -like than the other (say,  $\tau_1^Q = 0.9$  and  $\tau_2^Q = 0.2$ ). So, using the Qjets procedure, the experiment instead finds  $(0.9 + 0.2) = 1.1$   $W$ -jets, and measures the  $W$ -mass to be  $(0.9 \times 80 + 0.2 \times 85) / (0.9 + 0.2)$  GeV = 80.9 GeV. Furthermore, as we explain below, both of these observables (the number of tagged jets and the measured mass from the tagged jets) are statistically more robust in the case of the Qjets procedure than in the conventional approach. In fact, one can make a definite statement:

$$\left(\frac{\delta N_T}{N_T}\right)^C = \frac{1}{\sqrt{\epsilon N}} \quad \text{and} \quad \frac{1}{\sqrt{N}} \leq \left(\frac{\delta N_T}{N_T}\right)^Q \leq \frac{1}{\sqrt{\epsilon N}}, \quad (1.1)$$

where  $N_T$  represents the number of tagged jets that arise from a physical process expected to yield  $N$  total jets and  $\epsilon$  represents the efficiency of the conventional tagging procedure,  $\epsilon = N_T/N$ . So, if a process is expected to yield  $N = 100$  jets reconstructed at  $\epsilon = 50\%$  efficiency, one expects unweighted measurements of the cross section to have a statistical uncertainty of  $14\%$  ( $= 1/\sqrt{50}$ ). On the other hand, if one employs the Qjets procedure with the average tagging efficiency  $\epsilon$  still at  $50\%$ , one can achieve an uncertainty somewhere between  $10\%$  and  $14\%$ . Thus, with Qjets one can hope to obtain more precise results using the same data.

More specifically, the claims in ref. [44] regarding the uncertainties of various measurements can be stated as

$$\frac{S^Q/\delta B^Q}{S^C/\delta B^C} > 1 \quad \text{and} \quad \frac{\delta m^Q/m^Q}{\delta m^C/m^C} < 1. \quad (1.2)$$

These expressions use the definitions (to be explained in more detail later):

$$S^{Q/C} = \sum_{j \in \text{signal}} \tau_j^{Q/C}: \text{total number of signal jets correctly tagged in an experiment.}$$

$$B^{Q/C} = \sum_{j \in \text{bkg}} \tau_j^{Q/C}: \text{total number of QCD jets incorrectly tagged in an experiment.}$$

$$m^{Q/C} = \frac{\sum_j \mu_j^{Q/C} \tau_j^{Q/C}}{\sum_j \tau_j^{Q/C}}: \text{the (average) mass of the tagged jets as measured in an experiment.}$$

$$\delta B^{Q/C}, \delta m^{Q/C}: \text{the fluctuations in the corresponding measurements.}$$

In the phenomenological studies presented below we will confirm the inequalities in eq. (1.2) and attempt to provide intuitive explanations of why they hold. Note that the explanation is not as straightforward as for eq. (1.1).

It is helpful to think in terms of *two* types of effects contributing to the fact that the left-hand sides in eq. (1.2) are different from 1. As described above, an essential difference of the Qjets procedure is the shift from the binary tagging efficiency of the conventional approach,  $\tau_j^C$  ( $= 0$  or  $1$ ), to the continuously valued  $\tau_j^Q$  ( $0 \leq \tau_j^Q \leq 1$ ). Thus jets with  $\tau_j^C = 1$  can have  $\tau_j^Q < 1$ , while jets with  $\tau_j^C = 0$ , which make no contribution to the conventional analysis, can have  $\tau_j^Q > 0$  and contribute to the Qjets analysis. These changes impact both the counting of jets and the values of weighted averages, as in the weighted average mass defined just above. One of the important results of the Qjets analysis described below is that the distribution of jet-masses assigned by the Qjets procedure ( $\mu_j^Q$ ) for  $W$ -jets is found to be more sharply peaked around  $M_W$  than the  $\mu_j^C$  distribution. The Qjets procedure, since it samples a variety of pruning scenarios, can include scenarios that remove unwanted radiation from a  $W$ -jet more effectively than the single conventional pruning scenario [24, 25]. Since it is exactly these more effective scenarios that lead to larger weights in the Qjets analysis, the resulting weighted average mass tends to be closer to the physical  $W$ -mass. Thus the Qjets procedure can provide a better “groomer” than the classical pruning [24, 25]. In summary, the improvement indicated in eq. (1.2) stems

from both the “purely statistical” enhancement inherent in the shift from the binomial distribution of  $\tau_j^C$  to the continuous distribution of  $\tau_j^Q$ , which we label the “statistical” effect, and from the possible improvement in the measured signal mass distribution inherent in the shift of mass observable from  $\mu_j^C$  to  $\mu_j^Q$ , which we label the “physics” effect.

Of course, the “statistical” and “physics” effects are not explicitly independent. In an effort to provide a quantitative separation of these two effects, we define a third, hybrid pair of variables,  $(\mu_j^Q, \tilde{\tau}_j^Q)$ , where the mass variable remains the same as for the usual Qjets procedure, but the tagging probability variable  $\tilde{\tau}_j^Q$  follows a binomial distribution (similar to  $\tau_j^C$ ) defined by

$$\tilde{\tau}_j^Q = \begin{cases} 0 & \text{for } \tau_j^Q = 0 \\ 1 & \text{otherwise.} \end{cases} \quad (1.3)$$

With our definition of  $\mu_j^Q$  in Qjets,  $\tilde{\tau}_j^Q$  corresponds to tagging a jet based on whether  $\mu_j^Q$  is in the bin or not — i.e. tagging efficiency is derived just like in the conventional case, but using  $\mu_j^Q$  instead of  $\mu_j^C$ . Further, we define  $\mu_j^Q$  such that its value is in the bin if *any* of the Qjet masses for a given jet are in the bin, which is why all nonzero values for  $\tau_j^Q$  yield a  $\tilde{\tau}_j^Q$  value of 1.

The left-hand sides in eq. (1.2) can then be represented as *products* of statistical pieces and physics pieces:

- *statistical quantities:*  $\frac{S^Q/\delta B^Q}{\bar{S}^Q/\delta \bar{B}^Q}$  and  $\frac{\delta m^Q/m^Q}{\delta \tilde{m}^Q/\tilde{m}^Q}$ , exhibiting the impact of using a continuous versus binary variable,  $\tau_j^Q$  versus  $\tilde{\tau}_j^Q$ ;
- *physics quantities:*  $\frac{\bar{S}^Q/\delta \bar{B}^Q}{\bar{S}^C/\delta \bar{B}^C}$  and  $\frac{\delta \tilde{m}^Q/\tilde{m}^Q}{\delta \tilde{m}^C/\tilde{m}^C}$ , primarily exhibiting the impact of the differing distributions for the mass variables  $\mu_j^C$  versus  $\mu_j^Q$ .

The present article aims to clarify these points by presenting an explicit framework for calculating the statistics of jets obtained from the Qjets procedure, as applied to a jet-tagging analysis. The paper is structured as follows: in section 2 we introduce a statistical formalism for evaluating the uncertainties associated with the measurement of cross-section and mass for a tagging efficiency described by a continuous variable, in section 3 we review the Qjets procedure and discuss, in particular, how it leads to a mass and a tagging efficiency for a given jet, in section 4 we apply the formalism derived in section 2 to estimate the statistical and physics quantities outlined above, in section 5 and section 6 we present simple phenomenological pictures to assist in the understanding of the results for the statistical (section 5) and physics (section 6) effects presented in section 4, and in section 7 we provide concluding remarks. A validation of our analytical results, derived in section 2, using Monte Carlo pseudo-experiments is provided in appendix A, and more mathematical details are included in appendix B.

## 2 Statistical uncertainties

In this section we lay out the mathematical foundation needed to understand the statistical fluctuations of measurements when using the Qjets procedure (i.e., non-binary tagging). This analysis applies to both signal and background measurements.

One can think of the statistical uncertainties in jet-based measurements as arising from two sources: (1) Poisson uncertainty, and (2) sampling uncertainty:

- *Poisson uncertainty* refers to the uncertainty in the number of events (or jets) of a certain variety produced by a process yielding discrete counts at some continuous rate. For example, if a collider is expected to yield on average  $N$  events (of the given variety) with a given luminosity ( $N = \mathcal{L}\sigma$ , where  $\sigma$  is the production cross section for this kind of event or jet) then the probability of it producing  $n$  events is given by the Poisson distribution:

$$\text{Pois}(n|N) \equiv \frac{e^{-N}N^n}{n!}, \quad \langle n \rangle_{\text{Pois}} = N, \quad \sigma_{\text{Pois}}^2 = N. \quad (2.1)$$

Thus the variance ( $\sigma_{\text{Pois}}^2$ ) of this distribution is  $N$  as indicated, which tells us that the characteristic size of the variation in the number of events (of the given variety) produced with a given luminosity from one experimental run to the next is  $\sqrt{N}$ .

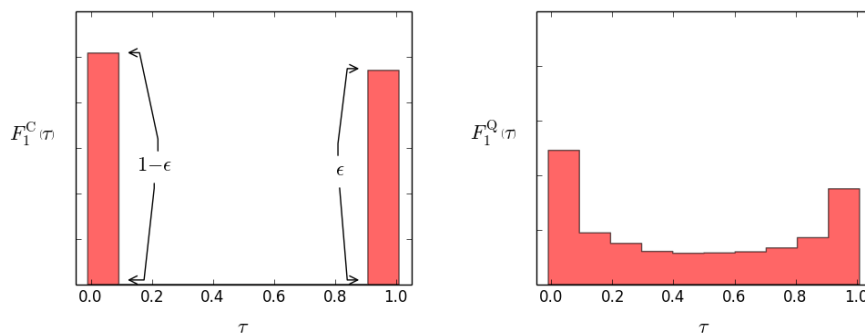
- *Sampling uncertainty* refers to the uncertainty in the way the events will be reconstructed by the analysis procedure, leading to fluctuations in the tagging rate sample-to-sample. Let us illustrate this point with an explicit example. Consider that we are trying to identify jets containing  $W$  decays with an algorithm characterized by a given tagging efficiency (say 70%). By sampling uncertainty we refer to the fact that for one sample of 100 signal jets the procedure might tag 75 jets as  $W$ -like, while for another sample of 100 signal jets it might only tag 65.

The next step is to explain why the probability distribution describing the tagging of jets in the Qjets procedure is fundamentally different from the conventional procedure, resulting in qualitatively (and quantitatively) different expressions for the sampling uncertainty, as well as the total statistical uncertainty. Recall that a conventional tagging procedure assigns a binary valued weight  $\tau$  of either 1 or 0 (i.e., tagged or not-tagged) to a jet. Such a procedure is usually characterized by a tagging efficiency  $\epsilon$ , which means that, on average, a fraction  $\epsilon$  of jets selected at random from a sample of  $W$ -jets will be tagged. Thus the explicit probability distribution function (or pdf) for tagging 1-jet, picked at random from a set of  $W$ -jets, by a conventional (C) procedure can be simply represented as:

$$F_1^{\text{C}}(\tau) = (1 - \epsilon)\delta(\tau) + \epsilon\delta(\tau - 1), \quad (2.2)$$

(where  $\delta(\tau)$  is the usual delta function that vanishes everywhere except at  $\tau = 0$ , but is sufficiently singular at  $\tau = 0$  to satisfy  $\int d\tau f(\tau)\delta(\tau) = f(0)$  for any range of integration that includes the origin). This form is illustrated in the left-hand plot in figure 1. In contrast, the weight  $\tau$  assigned by a Qjets procedure can have any value in the interval  $[0, 1]$ . We label the pdf for tagging 1-jet (picked at random from a set of  $W$ -jets) with probability  $\tau$  by the Qjets procedure as  $F_1^{\text{Q}}(\tau)$ . Note that, unlike eq. (2.2),  $F_1^{\text{Q}}(\tau)$  is a continuous function of  $\tau$ , as illustrated in the right-hand plot in figure 1.

These 1-jet tagging probability distribution functions (the  $F_1(\tau)$ 's illustrated in figure 1) are central to our analysis. As we will show later, the statistical uncertainties



**Figure 1.** Illustration of how  $F_1^C(\tau)$  (left) and  $F_1^Q(\tau)$  (right) may look for a sample of  $W$ -jets. Note the binomial nature of  $F_1^C(\tau)$  as opposed to the continuous distribution of  $\tau$  in  $F_1^Q(\tau)$ .

associated with various tagged 1-jet measurements are given entirely in terms of the first few moments of the  $F_1(\tau)$ . In particular, we define the average and variance (and the normalization) of  $F_1(\tau)$  to be:

$$\langle \tau \rangle = \int_0^1 \tau F_1(\tau) d\tau, \quad \text{and} \quad \sigma_\tau^2 = \int_0^1 (\tau - \langle \tau \rangle)^2 F_1(\tau) d\tau, \quad \left( \int_0^1 F_1(\tau) d\tau = 1 \right). \quad (2.3)$$

Note that, in the special case of the conventional procedure as in eq. (2.2),

$$\langle \tau^C \rangle = \int_0^1 \tau F_1^C(\tau) d\tau = \epsilon, \quad \text{and} \quad (\sigma_\tau^C)^2 = \int_0^1 (\tau - \langle \tau \rangle)^2 F_1^C(\tau) d\tau = \epsilon(1 - \epsilon), \quad (2.4)$$

where, as above,  $\epsilon$  is the average tagging efficiency in the conventional procedure. The results of eq. (2.4) are just what we expect from the binomial distribution corresponding to a *binary* valued weight. Note that the differences between the two distributions in figure 1 (binary versus continuous, where the latter has more support near the average value) already suggest that  $\sigma_\tau^Q < \sigma_\tau^C$ , even for cases where  $\langle \tau^Q \rangle \approx \epsilon$ .

For the corresponding hybrid analysis of eq. (1.3) we have a distribution similar to eq. (2.2),

$$\tilde{F}_1^Q(\tilde{\tau}) = (1 - \tilde{\epsilon})\delta(\tilde{\tau}) + \tilde{\epsilon}\delta(\tilde{\tau} - 1), \quad (2.5)$$

with moments

$$\langle \tilde{\tau}^Q \rangle = \tilde{\epsilon}, \quad \text{and} \quad (\sigma_{\tilde{\tau}}^Q)^2 = \tilde{\epsilon}(1 - \tilde{\epsilon}). \quad (2.6)$$

Note that, since the jets that were tagged in the conventional analysis are typically still tagged, and the Qjets procedure allows more jets to be tagged, with  $\tilde{\tau}_j^Q = 1$  in the hybrid analysis, we expect that  $\tilde{\epsilon} > \epsilon$ .

## 2.1 Cross-section measurement

As a first detailed example consider the statistical uncertainties inherent in the measurement of the production cross-section of jets containing the desired heavy particle. First,

imagine that  $N_S$  jets are selected at random from a set of  $W$ -jets. The total number of (correctly) tagged  $W$ -jets (or  $N_T$ ) is then given as

$$N_T = \sum_{j=1}^{N_S} \tau_j. \quad (2.7)$$

Since  $N_T$  is a sum of weights, it can exhibit non-integral values for the Qjets procedure. The probability distribution describing  $N_T$ , for a given sample size  $N_S$ , can be constructed in terms of  $F_1$ ,

$$F_{N_S}(N_T) = \left[ \prod_{k=1}^{N_S} \int F_1(\tau_k) d\tau_k \right] \delta \left( N_T - \sum_{k=1}^{N_S} \tau_k \right). \quad (2.8)$$

For future reference the first two moments of this general distribution are

$$\langle N_T \rangle_{N_S} = \int N_T dN_T F_{N_S}(N_T) = \left[ \prod_{k=1}^{N_S} \int F_1(\tau_k) d\tau_k \right] \sum_{k=1}^{N_S} \tau_k = N_S \langle \tau \rangle, \quad (2.9)$$

and

$$\begin{aligned} \langle N_T^2 \rangle_{N_S} &= \int N_T^2 dN_T F_{N_S}(N_T) = \left[ \prod_{k=1}^{N_S} \int F_1(\tau_k) d\tau_k \right] \left( \sum_{k=1}^{N_S} \tau_k \right)^2 \\ &= \left[ \prod_{k=1}^{N_S} \int F_1(\tau_k) d\tau_k \right] \left( \sum_{k=1}^{N_S} \tau_k^2 + \sum_{k \neq l}^{N_S} \tau_k \tau_l \right) \\ &= N_S \langle \tau^2 \rangle + N_S(N_S - 1) \langle \tau \rangle^2 = N_S^2 \langle \tau \rangle^2 + N_S (\langle \tau^2 \rangle - \langle \tau \rangle^2) \equiv N_S^2 \langle \tau \rangle^2 + N_S \sigma_\tau^2. \end{aligned} \quad (2.10)$$

For the conventional procedure (with a binary valued  $\tau$ )  $F_{N_S}^C(N_T)$  is given by the probability of selecting  $N_T$  objects from a set of  $N_S$  objects and the pdf is given by a *Binomial* distribution of mean  $\epsilon$ :

$$\begin{aligned} F_{N_S}^C(N_T) &= \left[ \prod_{k=1}^{N_S} \int F_1^C(\tau_k) d\tau_k \right] \delta \left( N_T - \sum_{k=1}^{N_S} \tau_k \right) \\ &= \frac{N_S!}{N_T!(N_S - N_T)!} \epsilon^{N_T} (1 - \epsilon)^{N_S - N_T} \equiv B(N_T | N_S, \epsilon), \end{aligned} \quad (2.11)$$

with moments

$$\begin{aligned} \langle N_T^C \rangle_{N_S} &= N_S \epsilon \\ \langle (N_T^C)^2 \rangle_{N_S} &= N_S^2 \epsilon^2 + N_S \epsilon (1 - \epsilon). \end{aligned} \quad (2.12)$$

Next we consider measuring the production cross section for the tagged jets. As noted above, the total statistical uncertainty depends on both the Poisson uncertainty and the sampling uncertainty. If the expected number of jets (for a given luminosity  $\mathcal{L}$ ) is  $N$ , on average the probability  $P$  of tagging  $N_T$  jets is given by:

$$P(N_T | N) = \sum_{N_S=N_T}^{\infty} \text{Pois}(N_S | N) \times F_{N_S}(N_T). \quad (2.13)$$



Evaluating eq. (2.13) in the conventional case is easier than one might expect, since the combination of a Poisson process and a Binomial process is still a Poisson process. We have

$$\begin{aligned} P^C(N_T|N) &= \sum_{N_S=N_T}^{\infty} \text{Pois}(N_S|N) \times F_{N_S}^C(N_T) \\ &= \sum_{N_S=N_T}^{\infty} \text{Pois}(N_S|N) \times B(N_T|N_S, \epsilon) = \text{Pois}(N_T|\epsilon N), \end{aligned} \quad (2.14)$$

i.e., it is a Poisson distribution with mean  $\epsilon N$ . Thus we can still apply our “ $\sqrt{N}$ ” intuition. Using eq. (2.14) (and eq. (2.1)) we find that the fractional uncertainty in the number of conventionally tagged jets is

$$\frac{\delta N_T^C}{N_T^C} = \frac{\sqrt{\sigma_{\text{Pois}}^2(N_T)}}{\langle N_T \rangle_{\text{Pois}}} = \frac{\sqrt{\epsilon N}}{\epsilon N} = \frac{1}{\sqrt{\epsilon N}}, \quad (2.15)$$

as already noted in eq. (1.1).

Thus, if we observe 100 events with tagged signal jets in  $\mathcal{L} = 1 \text{ fb}^{-1}$  with  $\epsilon = 50\%$ , we would report a cross section for signal jets of  $200 \pm 20 \text{ fb}$  (i.e.,  $\sigma = N_T/\epsilon/\mathcal{L} = 100/0.5 \text{ fb}$ , and  $\delta\sigma/\sigma = \delta N_T/N_T = 1/\sqrt{100} = 1/10$ ).

Evaluating statistical uncertainties for a general  $F_1(\tau)$ , e.g., a Qjets  $F_1^Q(\tau)$ , is slightly more complicated. In particular, for the Qjets case  $N_T$  is a sum of non-integer weights and so can exhibit non-integer values. For example, consider a sample of 5 jets/events. If, at the non-integer value 4.5,  $F_5^Q(4.5) = 0.1$ , then we interpret this to mean that the probability of measuring one jet/event in the bin  $4.5 \pm \rho/2$  is  $0.1 \times \rho$ , for infinitesimal  $\rho$ . In the following manipulations we treat  $N_T$  as a continuous variable. The mean of the distribution  $P(N_T|N)$  is obtained from (recall eq. (2.9))

$$\begin{aligned} \langle N_T \rangle &= \int N_T dN_T P(N_T|N) = \sum_{N_S=0}^{\infty} \text{Pois}(N_S|N) \int_0^{N_S} N_T dN_T F_{N_S}(N_T) \\ &= \sum_{N_S=0}^{\infty} \text{Pois}(N_S|N) N_S \langle \tau \rangle = \langle \tau \rangle N. \end{aligned} \quad (2.16)$$

The second moment of  $P(N_T|N)$  is found from (recall eq. (2.10))

$$\begin{aligned} \langle N_T^2 \rangle &= \int N_T^2 dN_T P(N_T|N) = \sum_{N_S=0}^{\infty} \text{Pois}(N_S|N) \int_0^{N_S} N_T^2 dN_T F_{N_S}(N_T) \\ &= \sum_{N_S=0}^{\infty} \text{Pois}(N_S|N) (N_S \sigma_\tau^2 + N_S^2 \langle \tau \rangle^2) \\ &= N \sigma_\tau^2 + N(N+1) \langle \tau \rangle^2. \end{aligned}$$

So the desired variance is

$$\begin{aligned} (\delta N_T)^2 &\equiv \langle N_T^2 \rangle - \langle N_T \rangle^2 = N \sigma_\tau^2 + N(N+1) \langle \tau \rangle^2 - N^2 \langle \tau \rangle^2 \\ &= N (\sigma_\tau^2 + \langle \tau \rangle^2). \end{aligned} \quad (2.17)$$

This is the general result including the analysis above for the conventional case in eq. (2.15), when we recall that in the conventional scenario (as in eq. (2.4))  $\langle \tau^C \rangle = \epsilon$ ,  $(\sigma_\tau^C)^2 = \epsilon(1 - \epsilon)$  so that  $(\sigma_\tau^C)^2 + \langle \tau^C \rangle^2 = \epsilon$ . In the Qjets analysis the distribution  $F_1(\tau)$  becomes non-zero at intermediate  $\tau$  values ( $\tau \neq 0, 1$ ), which, as already suggested, serves to *reduce*  $\sigma_\tau$  from its “conventional” value, as we will see explicitly shortly.

So it follows that for a general probability distribution  $F_1(\tau)$  we have

$$\frac{\delta N_T}{N_T} = \frac{1}{\sqrt{N}} \times \sqrt{1 + \frac{\sigma_\tau^2}{\langle \tau \rangle^2}}. \tag{2.18}$$

Since in the general case,  $\tau_k \leq 1.0$  and thus  $\tau_k^2 \leq \tau_k$ , the averages must obey

$$\langle \tau^2 \rangle \leq \langle \tau \rangle \quad \Rightarrow \quad \sigma_\tau^2 \equiv \langle \tau^2 \rangle - \langle \tau \rangle^2 \leq \langle \tau \rangle(1 - \langle \tau \rangle). \tag{2.19}$$

Thus we obtain (essentially as claimed in the Introduction) that

$$\frac{1}{\sqrt{N}} \leq \frac{\delta N_T}{N_T} \leq \frac{1}{\sqrt{\langle \tau \rangle N}}. \tag{2.20}$$

Comparing this to eq. (2.15) we see that the *upper* limit is saturated for the conventional procedure with binary valued tagging. This allows for the the fractional uncertainty in the cross-section measurement to be reduced by up to a factor of  $\sqrt{\langle \tau^C \rangle}$  ( $= \sqrt{\epsilon}$ ) if weighted jets are used in the measurement. *This is the advantage of using weighted jets — while we are still subject to the Poisson uncertainties in eq. (2.14), the sampling uncertainties, encoded in  $B(N_T|N_S, \epsilon)$  for a conventional tagging procedure, are reduced.*

## 2.2 Mass measurement

The statistical uncertainty of a cross section measurement is straightforward to compute with Qjets because the probability distribution for the number of tagged events factorizes nicely into one factor capturing the effects of Poisson uncertainties and one capturing the effects of sampling uncertainties (see eq. (2.13)). This is not generally true for other quantities that involve a weighted average rather than a simple sum, e.g., the average tagged jet mass is defined by

$$m_T = \frac{\sum_{j=1}^{N_S} \mu_j \tau_j}{\sum_{j=1}^{N_S} \tau_j} = \frac{1}{N_T} \sum_{j=1}^{N_S} \mu_j \tau_j. \tag{2.21}$$

The corresponding expression relevant to the hybrid analysis of eq. (1.3) is

$$\tilde{m}_T = \frac{\sum_{j=1}^{N_S} \mu_j \tilde{\tau}_j}{\sum_{j=1}^{N_S} \tilde{\tau}_j} = \frac{1}{\tilde{N}_T} \sum_{j=1}^{N_S} \mu_j \tilde{\tau}_j. \tag{2.22}$$

One can still relate the relevant uncertainties to the underlying probability distribution functions; however, the resulting expressions are more complicated. In particular,  $F_1(\tau)$  is no longer enough. We now need to know the probability distribution as a function of *both*

$\tau$  and  $\mu$ . We label this distribution  $F_1(\mu, \tau)$ , which denotes the probability distribution in the  $(\mu, \tau)$  plane. Note that  $F_1(\tau)$  is simply related to  $F_1(\mu, \tau)$  by

$$F_1(\tau) = \int d\mu F_1(\mu, \tau). \quad (2.23)$$

For illustration we show the  $F_1^C(\mu, \tau)$  and  $F_1^Q(\mu, \tau)$  distributions in figure 2 derived from a sample of  $W$ -jets. In the conventional procedure, a jet with jet mass inside a pre-defined mass window, for example,  $\Omega = (70 - 90)$  GeV for  $W$ -tagging, is tagged (with  $\tau = 1$ ). This fact is demonstrated by  $F_1^C(\mu, \tau)$ , where all non-zero entries are in the bin at  $\tau = 1$  and the jet mass distribution peaks around the  $W$ -mass. On the other hand,  $F_1^Q(\mu, \tau)$  shows that there are non-zero probabilities for tagging jets with efficiency  $\tau^Q$  in the full range  $[0, 1]$  for jet masses in the tagging window  $\Omega$ . Note that the contributions that lead to the strictly  $\tau = 0$  part of the distribution (see, for example, eq. (2.2)) all arise from  $\mu$  values *outside* of  $\Omega$ .

In this section, we simply define moments of the two-dimensional distribution, and leave all technical details to appendix B. The moments of interest, the single averages, the two-dimensional mean, variance and covariance are given by

$$\begin{aligned} \langle \tau \rangle &\equiv \int d\mu \int_0^1 d\tau \tau F_1(\mu, \tau) = \int_0^1 d\tau \tau F_1(\tau), \\ \langle \mu\tau \rangle &\equiv \int d\mu \int_0^1 d\tau \mu\tau F_1(\mu, \tau), \\ \sigma_\tau^2 &\equiv \int d\mu \int_0^1 d\tau (\tau - \langle \tau \rangle)^2 F_1(\mu, \tau), \\ \sigma_{\mu\tau}^2 &\equiv \int d\mu \int_0^1 d\tau (\mu\tau - \langle \mu\tau \rangle)^2 F_1(\mu, \tau), \\ \sigma(\tau, \mu\tau) &\equiv \int d\mu \int_0^1 d\tau (\mu\tau - \langle \mu\tau \rangle)(\tau - \langle \tau \rangle) F_1(\mu, \tau). \end{aligned} \quad (2.24)$$

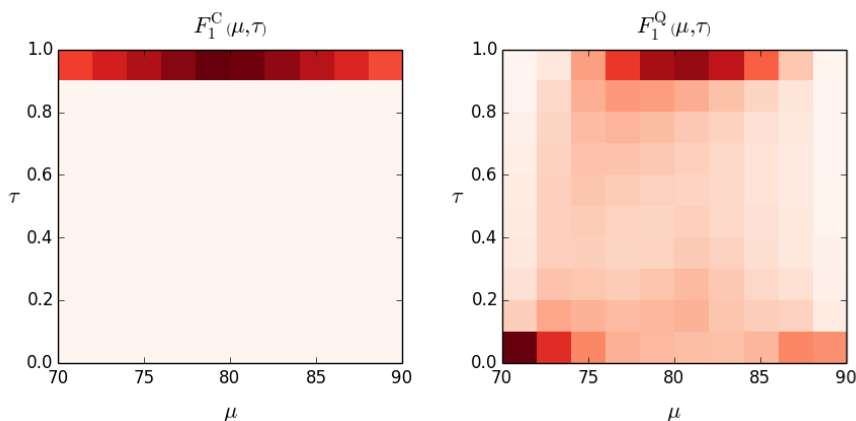
Note especially that, since  $\tau$  and  $\mu$  are correlated by  $F_1(\mu, \tau)$ , we are now interested in both the variance of the parameters  $\tau$  and  $\mu\tau$  and in the covariance  $\sigma(\tau, \mu\tau)$ .

So we are now ready to consider the measurement of the average (weighted) jet mass as defined in eq. (2.21), where we want to understand the expected improvement in precision from using the Qjets technique. Proceeding essentially as we did in the cross section case, the expected average value of  $m_T$  in a sample of  $N_S$  jets, is given by (recall eq. (2.21))

$$\langle m_T \rangle_{N_S} \simeq \frac{\langle \mu\tau \rangle}{\langle \tau \rangle} \left[ 1 + \frac{\sigma_\tau^2}{N_S \langle \tau \rangle^2} - \frac{\sigma(\tau, \mu\tau)}{N_S \langle \mu\tau \rangle \langle \tau \rangle} \right]. \quad (2.25)$$

As explained in the appendix, we are expanding in the fluctuations around the average values and assuming that the higher order fluctuations are negligible. The corresponding variance in this quantity is given by

$$(\delta m_T)_{N_S}^2 = \langle (m_T - \langle m_T \rangle_{N_S})^2 \rangle_{N_S} \simeq \frac{\langle \mu\tau \rangle^2}{N_S \langle \tau \rangle^2} \left[ \frac{\sigma_{\mu\tau}^2}{\langle \mu\tau \rangle^2} + \frac{\sigma_\tau^2}{\langle \tau \rangle^2} - 2 \frac{\sigma(\tau, \mu\tau)}{\langle \mu\tau \rangle \langle \tau \rangle} \right]. \quad (2.26)$$



**Figure 2.** Illustration of how  $F_1^C(\mu, \tau)$  and  $F_1^Q(\mu, \tau)$  may behave for a sample of  $W$ -jets. Since these plots are for illustration purposes only, we do not provide the numerical values associated with different shades of red. Qualitatively, the lightest shade in these plots represents a vanishing  $F_1(\mu, \tau)$  value, and a darker shade represents a larger value of  $F_1$ . Note that all jets in the jet mass window (70 – 90) GeV are tagged (with  $\tau = 1$ ) in a conventional procedure and so only the  $\tau = 1$  boxes will be non-zero for  $F_1^C(\mu, \tau)$ . On the other hand,  $F_1^Q(\mu, \tau)$  can be non-zero in the entire  $(\mu, \tau)$  plane.

If we average over samples (to take into account the Poisson uncertainties) within an experiment with a given luminosity, then we have  $N_S \rightarrow N = \sigma \mathcal{L}$  in the denominator of both eq. (2.25) and eq. (2.26), plus corrections of order  $1/N^2$ . Combining the above results, the ratio of the fluctuations to the average value can be written as:

$$\left(\frac{\delta m_T}{\langle m_T \rangle}\right)^2 = \frac{1}{N} \left( \frac{\sigma_{\mu\tau}^2}{\langle \mu\tau \rangle^2} + \frac{\sigma_\tau^2}{\langle \tau \rangle^2} - 2 \frac{\sigma(\tau, \mu\tau)}{\langle \mu\tau \rangle \langle \tau \rangle} \right) + \mathcal{O}\left(\frac{1}{N^2}\right). \quad (2.27)$$

We can easily evaluate this quantity for the conventional binary tagging procedure. By definition, and as illustrated in figure 2,  $\tau = 1$  for  $\mu \in \Omega$  when we consider the pdf  $F_1^C(\mu, \tau)$ . The fractional mass uncertainty of eq. (2.27) for the conventional tagging procedure with average tagging efficiency  $\epsilon$  is then

$$\left(\frac{\delta m_T^C}{\langle m_T^C \rangle}\right)^2 = \frac{1}{N} \times \frac{(\sigma_\mu^C)^2}{\epsilon \langle \mu^C \rangle^2} + \mathcal{O}\left(\frac{1}{N^2}\right), \quad (2.28)$$

where we define the properly normalized mass distribution moments *in* the mass window  $\Omega$  as

$$\begin{aligned} \langle \mu \rangle &\equiv \frac{1}{N_\Omega} \int_\Omega d\mu \int_0^1 d\tau \mu F_1(\mu, \tau), \\ \sigma_\mu^2 &\equiv \frac{1}{N_\Omega} \int_\Omega d\mu \int_0^1 d\tau (\mu - \langle \mu \rangle)^2 F_1(\mu, \tau), \\ \text{with } N_\Omega &= \int_\Omega d\mu \int_0^1 d\tau F_1(\mu, \tau). \end{aligned} \quad (2.29)$$

Here  $N_\Omega$  fixes the normalization of the pdf  $F_1$  in the mass window  $\Omega$ . In eq. (2.28) we follow the convention in eq. (2.4) to denote that the moments  $\langle \mu^C \rangle$  and  $\sigma_\mu^C$  are calculated from eq. (2.29) using the conventional pdf  $F_1^C(\mu, \tau)$ . Note that in the conventional case the normalization is

$$(N_\Omega)^C = \int_\Omega d\mu \int_0^1 d\tau F_1^C(\mu, \tau) = \int_0^1 d\tau \left( \int_\Omega d\mu F_1^C(\mu, \tau) \right) = \int_0^1 d\tau \epsilon \delta(1-\tau) = \epsilon, \quad (2.30)$$

where we have used the fact that in the conventional or classical analysis all jets in the tagging window have  $\tau = 1$ . Once again, the reader is directed to appendix B for details.

### 3 Review of Qjets

The purpose of this section is to demonstrate how the Qjets procedure assigns a jet mass ( $\mu_j^Q$ ) and tagging efficiency ( $\tau_j^Q$ ) to a given jet  $j$ . Before describing the details, let us first review the general idea of the procedure. As suggested in ref. [44], we start with jets identified using a standard algorithm like Anti- $k_T$  [47]. We recluster the constituents of the given jet using a sequential and probabilistic recombination algorithm, such as  $k_T$  [48, 49] or Cambridge/Aachen (C/A) [50–52]. During clustering, pruning [24, 25] is performed in order to remove unwanted elements in the jet, i.e., those elements not arising from the decay of the desired heavy object. Through pruning we map a jet to its pruned version. If the above set of steps is repeated on the same jet using a slightly different recombination metric as explained below (the Qjets procedure), we obtain a different four-vector after pruning due to the probabilistic nature of the Qjets clustering algorithm. We iterate the procedure a number of times (say  $N_{\text{iter}}$ ) to map a jet to a set of pruned four-vectors. The quantities  $\mu_j^Q$  and  $\tau_j^Q$  are then calculated from the invariant masses of these pruned four-vectors.

In more detail, sequential recombination algorithms build up jets by merging four-momenta in pairs over many steps. The behavior of the algorithms is determined by the metric for measuring the “distance” between four-momenta. At each stage in the jet clustering, one identifies the pair of four-momenta with the smallest distance and merges them together (i.e., adds the corresponding 4-momenta and replaces the merged pair with this sum in the updated list of 4-momenta). This merging step is repeated on the list of 4-momenta until all remaining 4-momenta are separated by more than a predefined cutoff. See ref. [53] for a more comprehensive discussion. For instance, the  $k_T$  [48, 49] and C/A [50–52] algorithms correspond to the following metrics:

$$d_{ij}^{k_T} \equiv \min\{p_{T_i}^2, p_{T_j}^2\} \Delta R_{ij}^2 \quad \text{and} \quad d_{ij}^{C/A} \equiv \Delta R_{ij}^2, \quad (3.1)$$

where  $\Delta R_{ij}^2 = \Delta y_{ij}^2 + \Delta \phi_{ij}^2$  is the squared angular distance between a pair of four-momenta  $i$  and  $j$  (with  $y$  the usual rapidity and  $\phi$  the azimuthal angle). Thus the C/A algorithm merges the 4-momenta in strict order of their angular separation with closest merged first. The  $k_T$  algorithm, in contrast, gives some emphasis to merging the smallest  $p_T$  elements first and so the two algorithms will tend to identify jets with slightly different constituents.

As implemented in ref. [44], the Qjets procedure also processes jets via pairwise mergings with pruning applied at each merging step. However, unlike traditional clustering which works deterministically, Qjets uses a probabilistic clustering procedure:

1. At every stage of clustering, for each pair of four-vectors (say,  $i$  and  $j$ ), the conventional distance metric  $d_{ij}$  from eq. (3.1) (for  $k_T$  or C/A) is evaluated for all such pairs. This is translated into a weight  $\omega_{ij}^{(\alpha)}$  via

$$\omega_{ij}^{(\alpha)} \equiv \exp \left\{ -\alpha \frac{(d_{ij} - d^{\min})}{d^{\min}} \right\}, \quad (3.2)$$

where  $d^{\min}$  is the smallest  $d_{ij}$  at this stage in the clustering process and  $\alpha$  (termed *rigidity*) is a continuous real parameter. This weight is then used to assign a probability  $\Omega_{ij}$  to each pair via

$$\Omega_{ij} = \omega_{ij} / N, \text{ where } N = \sum_{\langle ij \rangle} \omega_{ij}. \quad (3.3)$$

2. A random number is generated and used to select a pair  $\langle ij \rangle$  with probability  $\Omega_{ij}$ . Note that the conventional clustering process will always choose the pair with the minimum  $d_{ij}$  at this point and corresponds to the limit  $\alpha \rightarrow +\infty$ .
3. Having chosen the pair  $\langle ij \rangle$ , the standard pruning procedure is applied. The softer of the two selected four-momentum pair  $\langle ij \rangle$  is discarded, if *both* of the following criteria are satisfied for a given set of parameters ( $z_{\text{cut}}, D_{\text{cut}}$ ).

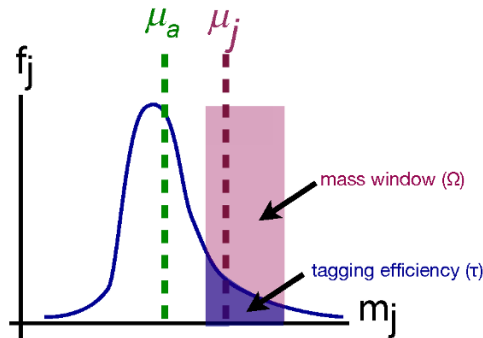
$$z \equiv \frac{\min(p_{T_i}, p_{T_j})}{p_{T_p}} < z_{\text{cut}} \quad \text{and} \quad \Delta R_{ij} > D_{\text{cut}}. \quad (3.4)$$

Otherwise, the pair is merged.

4. Steps (1-3) are repeated until all constituents are clustered. The invariant mass of the resultant pruned four-vector is stored for further analysis.
5. Steps (1-4) are repeated  $N_{\text{iter}}$  times. This procedure yields a set of  $N_{\text{iter}}$  masses for every jet it operates on. Due to the random numbers in step 2 these masses are generally not the same, but instead define a distribution of masses.

In summary, the Qjets procedure maps the initial jet  $j$  to a set of masses,  $\{m_{j,k}\}$ , where  $k$  takes integer values in  $[1, N_{\text{iter}}]$ . For each jet  $j$  we can construct a probability distribution  $f_j(m_j)$  as suggested in figure 3, with normalization  $\int f_j(m_j) dm_j = 1$ . For  $N_{\text{iter}} \gg 1$  this distribution will be relatively smooth and we will treat it as a continuous function,

$$f_j(m_j) \equiv \lim_{N_{\text{iter}} \gg 1} \frac{1}{N_{\text{iter}}} \sum_{k=1}^{N_{\text{iter}}} \delta(m_j - m_{k,j}). \quad (3.5)$$



**Figure 3.** Sketch of the pruned jet mass distribution for a jet processed many times with Qjets. The red area represents the mass window ( $\Omega$ ), the fraction of the jetmass distribution within the mass window (blue) is the tagging efficiency of the jet  $\tau_j^Q$ , and  $\mu_j^Q$  is the mean jetmass-in-the-window.  $\mu_a$  is the average jet mass for the entire distribution.

As described above, to define a tagging process, for example for W-jets, we focus on the W-like mass window  $\Omega$  illustrated in figure 3. For a given jet  $j$ , the tagging probability  $\tau_j^Q$  is the fraction of the  $N_{\text{iter}}$  clustering sequences yielding a pruned mass within the  $W$ -window,

$$\tau_j^Q = \frac{1}{N_{\text{iter}}} \sum_{k \ni m_{j,k} \in \Omega} 1 = \int_{\Omega} f_j(m_j) dm_j. \tag{3.6}$$

Similarly we define  $\mu_j^Q$  as the mean value of the pruned jet mass for these  $W$ -like interpretations for the same jet. Thus we have

$$\mu_j^Q = \frac{1}{\tau_j N_{\text{iter}}} \sum_{k \ni m_{j,k} \in \Omega} m_{j,k} = \frac{\int_{\Omega} f_j(m_j) m_j dm_j}{\int_{\Omega} f_j(m_j) dm_j}. \tag{3.7}$$

For comparison,  $\mu_a$  in figure 3 indicates the average jet mass for the full distribution, not just in the signal window. For a background (QCD) jet, this full-average mass value is generally quite different from  $\mu_j^Q$ .

Let us quickly review the Qjets procedure up to this point. We begin with a choice of the jet finding algorithm and kinematic cuts, e.g., the anti- $k_T$  jet algorithm with  $R = 1.0$  and kinematic cuts on the jet,  $p_T > 200$  GeV and rapidity  $|y| \leq 1.0$ . Then we subject the jets identified in this fashion to the Qjets procedure with specific choices of the Qjets parameters  $\alpha$  and  $N_{\text{iter}}$  to produce the single-jet pruned mass distribution in figure 3. With a specific signal jet in mind, say boosted W-jets, we define the mass window  $\Omega$  in figure 3. This procedure results in values for  $\tau_j^Q$  and  $\mu_j^Q$  from eqs. (3.6) and (3.7), which provide a measure of the likelihood that the given jet is a signal jet along with an estimate of the “true” mass of that signal jet.

## 4 Results from phenomenological studies

As an introduction to the following discussion of the results of our phenomenological studies, recall that the goal of the current work is to provide a more detailed explanation of the claim made in ref. [44] that the Qjets procedure improves the statistical stability of jet observables. The fundamental point is that, unlike a conventional binary tagging algorithm that identifies a jet as either tagged or not, the Qjets procedure yields a continuously valued tagging probability for a jet (as detailed in section 3). If observables are constructed using these tagging probabilities, it is non-trivial to estimate the statistical uncertainties associated with these observables as the tagging probabilities exhibits a continuous distribution on the interval  $[0, 1]$ . For example, the well known result that the statistical uncertainty associated with the measurement of the number of tagged jets is given by  $\delta N_T = \sqrt{N_T}$ , is no longer true. In section 2, we gave analytic expressions for these uncertainties. In this section, we report the results of our phenomenological studies, where we analyze carefully prepared event samples (generated by standard Monte Carlo event generators) and use the formulas from section 2 to demonstrate that indeed the Qjets procedure improves the statistical uncertainties associated with cross-section and mass measurements.

To be specific, we study the problem of tagging jets containing the decay products of  $W$ -particles. We treat a set of  $WW$  diboson events, where both  $W$ s decay hadronically, as signal events. We also consider QCD dijet events that provide the primary background to  $W$ -tagging. We generate both signal and background events for a 14 TeV LHC, using Pythia 8 [54]. Additionally, we use the “ATLAS UE Tune AU2-CTEQ6L1” [55] provided by Pythia 8 to give these events a realistically busy environment corresponding to actual proton-proton collisions. The detector simulation is provided by Delphes [56]. In particular, we use the default parameters provided by Delphes to simulate the ATLAS detector. Delphes output consists of energy flow four-vectors that are constructed out of the calorimeter cells, tracks, and muon elements of the detector. We do not impose any additional cut on rapidity or  $p_T$  on the Delphes output. We cluster the Delphes outputs into anti- $k_T$  jets with  $R = 0.7$  and  $p_T > 500$  GeV using Fastjet [57]. Only the leading jet from each event is selected for further analysis.

We perform the Qjets procedure using the publicly available Qjets plugin.<sup>1</sup> The constituents of the selected jets are reclustered for various values of the rigidity parameter (listed in table 1) using the C/A definition of the separation metric (see eq. (3.1)). For the pruning parameter  $D_{\text{cut}}$  we use  $D_{\text{cut}} = m/p_T$ , where  $m$  and  $p_T$  are the mass and the transverse momentum of the unpruned jet respectively. We perform our analysis for two  $z_{\text{cut}}$  parameter values, 0.1 and 0.15, where the smaller value corresponds to the default or optimized pruning case and the larger value should lead to a bit of “over”-pruning. Results for both of these values are listed in table 1. Finally, we set the  $W$ -mass window to (70 – 90) GeV for the purpose of tagging.

At this point, we reiterate that in this work we are interested in *both* the effects of a switching from a binary to a continuous tagging variable and from the corresponding change in the weighted average mass. In order to define a separation of these effects we introduced

<sup>1</sup><http://jets.physics.harvard.edu/Qjets>.



a new binary tagging efficiency  $\tilde{\tau}^Q$  for every  $\tau^Q$  obtained after the Qjets procedure (as defined in the Introduction). The ratio of the statistical uncertainty estimated using  $\tau^Q$  to that using  $\tilde{\tau}^Q$  therefore provides an estimate of the statistical improvement arising primarily from the differences between binary and continuous tagging variables (with the identical mass distribution). We also consider the differences between an analysis using  $(\mu^Q, \tilde{\tau}^Q)$  versus one using  $(\mu^C, \tau^C)$  to try to isolate the effects primarily due to the changes in the mass distribution (which we label “physics” effects).

To introduce our explicit numerical results it will be useful to make a few more comments to define the notation used:

- As noted above, we are studying both a sample of  $W$ -jets, or signal jets, and QCD-jets, or background jets. The corresponding results will be labeled by  $S$  and  $B$ .
- We also include results for the hybrid analysis of eq. (1.3) that is intended to separate statistical from physics effects. In particular, since this analysis uses a binary  $\tilde{\tau}^Q$  (with values only 0 and 1), the corresponding average tagging efficiency and fluctuation are given by  $\langle \tilde{\tau}^Q \rangle = \tilde{\epsilon}$  and  $\sigma_{\tilde{\tau}}^Q = \tilde{\epsilon}(1 - \tilde{\epsilon})$  respectively (see eq. (2.6)). The uncertainties associated with the measurement of the cross-section and mass in this hybrid analysis can be estimated from the corresponding formulas for conventional analysis in eq. (2.15) and eq. (2.28), respectively, using the substitutions  $\epsilon \rightarrow \tilde{\epsilon}$ ,  $\langle \mu^C \rangle \rightarrow \langle \tilde{\mu}^Q \rangle$ , and  $\sigma_{\mu}^C \rightarrow \sigma_{\tilde{\mu}}^Q$ . Once again we follow the convention that the appearance of  $\tilde{\tau}$  and  $\tilde{\mu}$  in these moments reflects the fact that these moments are calculated from their definitions in eqs. (2.3), (2.24), (2.29) using the hybrid pdf  $\tilde{F}_1^Q$ , which we discuss in more detail below.

The statistical quantities we look at are given by the following equations:

$$\begin{aligned} \frac{\delta S^Q / \sqrt{S^Q}}{\delta \tilde{S}^Q / \sqrt{\tilde{S}^Q}} &= \frac{\delta S^Q}{\sqrt{S^Q}} = \sqrt{\langle \tau_S \rangle + \frac{\sigma_{\tau_S}^2}{\langle \tau_S \rangle}}, & \frac{\delta B^Q / \sqrt{B^Q}}{\delta \tilde{B}^Q / \sqrt{\tilde{B}^Q}} &= \frac{\delta B^Q}{\sqrt{B^Q}} = \sqrt{\langle \tau_B \rangle + \frac{\sigma_{\tau_B}^2}{\langle \tau_B \rangle}} \\ \frac{S^Q / \delta B^Q}{\tilde{S}^Q / \delta \tilde{B}^Q} &= \left( \frac{\langle \tau_S \rangle}{\tilde{\epsilon}_S} \right) \times \sqrt{\frac{\tilde{\epsilon}_B}{\langle \tau_B \rangle}} \times \frac{1}{\sqrt{\langle \tau_B \rangle + \frac{\sigma_{\tau_B}^2}{\langle \tau_B \rangle}}}, & & (4.1) \\ \frac{\delta m_T^Q / m_T^Q}{\delta \tilde{m}_T^Q / \tilde{m}_T^Q} &= \sqrt{\left( \frac{\sigma_{\mu_S \tau_S}^2}{\langle \mu_S \tau_S \rangle^2} + \frac{\sigma_{\tau_S}^2}{\langle \tau_S \rangle^2} - 2 \frac{\sigma(\tau_S, \mu_S \tau_S)}{\langle \tau_S \rangle \langle \mu_S \tau_S \rangle} \right) / \left( \frac{\sigma_{\tilde{\mu}_S}^2}{\tilde{\epsilon}_S \langle \tilde{\mu}_S \rangle^2} \right)}, \end{aligned}$$

where we have used the equations derived in section 2.

In table 1, we tabulate the numerical estimations of the various observables for different values of  $z_{\text{cut}}$  and  $\alpha$ . In the remaining part of this section we provide a brief description of the patterns observed in table 1. Detailed explanations of these observations will be provided in the following two sections.

The first four observables in the table capture what we have labeled the statistical improvements seen in the Qjets procedure for the signal and the background samples. The quantity  $\delta N_T^Q / \sqrt{N_T^Q}$  for both signal and background represents the improvement in the

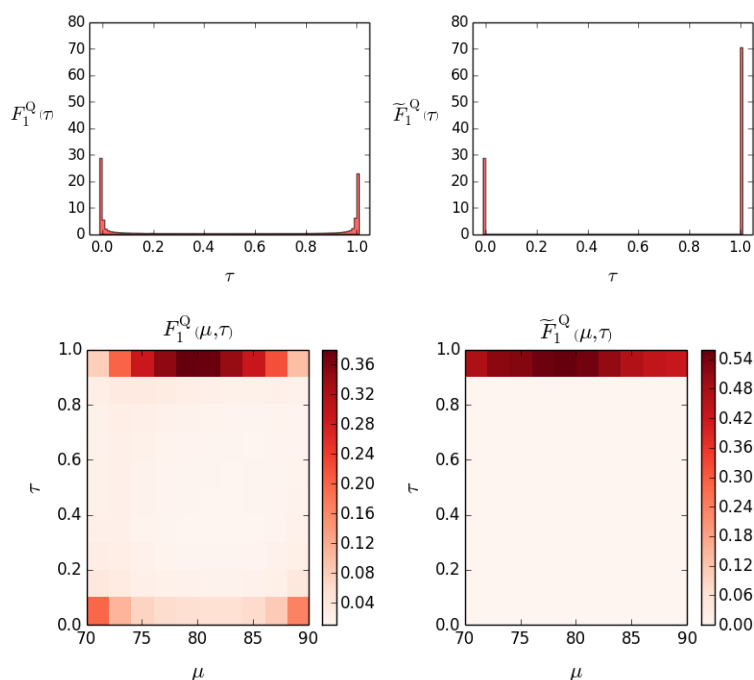
$\alpha$	Statistical Effects								Total uncertainty			
	$\frac{\delta S^Q}{\sqrt{S^Q}}$		$\frac{\delta B^Q}{\sqrt{B^Q}}$		$\frac{S^Q/\delta B^Q}{\tilde{S}^Q/\delta \tilde{B}^Q}$		$\frac{\delta m_T^Q/m_T^Q}{\delta \tilde{m}_T^Q/\tilde{m}_T^Q}$		$\frac{S^Q/\delta B^Q}{S^C/\delta B^C}$		$\frac{\delta m_T^Q/m_T^Q}{\delta m_T^C/m_T^C}$	
	$z_{\text{cut}}$		$z_{\text{cut}}$		$z_{\text{cut}}$		$z_{\text{cut}}$		$z_{\text{cut}}$		$z_{\text{cut}}$	
	0.10	0.15	0.10	0.15	0.10	0.15	0.10	0.15	0.10	0.15	0.10	0.15
10.0	0.99	0.98	0.94	0.94	1.17	1.15	1.00	1.01	1.05	1.05	0.96	0.96
1.00	0.95	0.94	0.85	0.85	1.42	1.38	1.00	1.05	1.16	1.17	0.86	0.86
0.10	0.90	0.88	0.74	0.72	1.63	1.57	1.00	1.08	1.26	1.29	0.73	0.71
0.01	0.86	0.82	0.69	0.66	1.61	1.54	0.98	1.00	1.22	1.25	0.65	0.56
0.00	0.87	0.82	0.77	0.72	1.28	1.24	0.88	0.92	1.00	1.03	0.60	0.52

**Table 1.** Statistical uncertainties associated with various measurements of cross-section and mass. Formulas used to estimate these quantities as listed in eq. (4.1).

uncertainty of the measured cross-section due to what we have labeled statistical effects. Note that this quantity is unity for a binary tagging variable, which is why the denominator becomes unity in the first line of eq. (4.1). For large  $\alpha$  these quantities are close to 1 for both values of the  $z_{\text{cut}}$  parameter. This situation reflects the fact that at high rigidity the individual mass distribution for each jet is quite narrow even after applying the Qjets procedure and  $\tau$  mostly has the values 0 or 1 (i.e., the Qjets procedure approaches the “classical” limit as  $\alpha \rightarrow \infty$ ). As indicated in table 1, the uncertainties decrease as  $\alpha$  is decreased and we include an increasing range of different clustering/pruning scenarios until a plateau is reached at  $\alpha \sim 0.01$  (for the background the uncertainty actually turns over and starts to increase again as  $\alpha \rightarrow 0$ ). It is interesting also to note that the improvement with decreasing  $\alpha$  is slightly better (i.e., smaller values of the ratio) for the less optimal  $z_{\text{cut}}$  value (0.15). This feature presumably arises from the fact that we start, in the classical limit, with less than optimal pruning, which allows the Qjets procedure more opportunity to include different clustering/pruning scenarios that improve the situation. The background case is somewhat less  $z_{\text{cut}}$  dependent as expected, as there is less of a clear definition of optimal pruning.

The statistical improvement in the discovery potential is captured by the third quantity, the ratio  $(S^Q/\delta B^Q) / (\tilde{S}^Q/\delta \tilde{B}^Q)$ . The larger this number becomes, the better is the chance that a precise measurement of the signal can be performed with a given luminosity. Once again we see that this observable is maximized for a small  $\alpha \sim 0.01$ . The small  $z_{\text{cut}}$  dependence in this case makes the not unexpected suggestion that it is best for the Qjets procedure to perturb around an optimal classical choice of parameters.

Finally, the fourth observable in table 1 provides an estimate of the uncertainty associated with the measurement of the jet mass arising from what we have labeled as statistical effects. We interpret the fact that this ratio remains near unity (except for very small values of  $\alpha \sim 0$ ) as confirmation that we have largely succeeded in separating the effects of binary versus continuous tagging variables, which we see are small for this variable, from the effects of changing the mass distribution itself, which will be important for this quantity. We refer the reader to section 5 for further explanation of these observations.



**Figure 4.** The probability distributions  $F_1^Q(\tau)$ ,  $\tilde{F}_1^Q(\tau)$ ,  $F_1^Q(\mu, \tau)$ ,  $\tilde{F}_1^Q(\mu, \tau)$  derived from a sample of  $W$ -jets. These particular distribution are produced using  $\alpha = 1.0$  and  $z_{\text{cut}} = 0.1$ . For the rest of the parameters, see section. 4.

For completeness we include our estimate of the total improvements provided by the Qjets procedure using the last two observables in table 1. These observables compare the uncertainties in the Qjets procedure to those in the conventional or classical procedure. As explained earlier, these quantities can be calculated from eq. (4.1) by the replacements  $\tilde{\epsilon} \rightarrow \epsilon$ ,  $\langle \tilde{\mu}^Q \rangle \rightarrow \langle \mu^C \rangle$ , and  $\sigma_{\tilde{\mu}}^Q \rightarrow \sigma_{\mu}^C$ . Overall we find that the behavior of the statistical uncertainties associated with the cross-section and mass is similar to what was described in ref. [44]. The cross-section measurement is most stable in the range  $0.1 \geq \alpha > 0.01$ , whereas the mass uncertainty prefers even smaller rigidity ( $0.01 > \alpha \geq 0.0$ ).

Note that the contribution to the uncertainties from what we have labeled physics effects can be found by simply dividing the total uncertainty by the corresponding statistical contribution. These results will be discussed in more detail in section 6. It is worthwhile noting that this exercise already tells us that the effects we labeled physics will be more important than the statistical effects for the mass measurement uncertainties, as we just suggested.

## 5 Understanding the statistical effects

In order to understand the uncertainties listed in table 1 it is essential to study the probability distributions  $F_1(\tau)$  and  $F_1(\mu, \tau)$ . In figure 4, we display these distributions as derived from a sample of  $W$ -jets. On the left are the distributions  $F_1^Q(\tau)$  and  $F_1^Q(\mu, \tau)$  arising from the full Qjets analysis, while the plots on the right illustrate the distributions

$\tilde{F}_1^Q(\tau)$  and  $\tilde{F}_1^Q(\mu, \tau)$  from the hybrid analysis. Recall that the latter analysis uses the binary tagging probability  $\tilde{\tau}^Q$  derived from the standard Qjets probability  $\tau^Q$  as defined in eq. (1.3), i.e., all nonzero  $\tau^Q$  values ( $\tau^Q > 0$ ) correspond to  $\tilde{\tau}^Q = 1$ . By construction,  $F_1^Q(\tau = 0) = \tilde{F}_1^Q(\tau = 0)$  as illustrated by the equal heights of the zero bins of  $F_1^Q(\tau)$  and  $\tilde{F}_1^Q(\tau)$  in figure 4. The difference between the two distributions arises from the fact that the rest of the probability in  $\tilde{F}_1^Q(\tau)$  all lies in the  $\tau = 1$  bin, whereas  $F_1^Q(\tau)$  exhibits nonzero probability at values of  $\tau$  between 0 and 1 (although it is still strongly peaked in the  $\tau = 1$  bin). In other words in moving from the  $\tilde{F}_1^Q(\tau)$  distribution to the  $F_1^Q(\tau)$  distribution (i.e., moving from a binary tagging probability to a continuous one), probability “leaks out” of the  $\tau = 1$  bin into the  $1 > \tau > 0$  bins.

The lower plots of  $F_1^Q(\mu, \tau)$  and  $\tilde{F}_1^Q(\mu, \tau)$  provide additional information. In particular, almost all jets that leak-out of the  $\tau = 1$  bin, as one moves from  $\tilde{F}_1^Q(\mu, \tau)$  to  $F_1^Q(\mu, \tau)$ , lie near or at one of the boundaries of the window in  $\mu$ . Also note that the distribution in  $\mu$  corresponding to  $\tau = 1$  is peaked near the  $W$  mass (as expected for an underlying  $W$ -jet sample), and that the  $\tau = 0$  bin does not actually appear in the lower plots as all of the corresponding  $\mu$  values are *outside* of the  $\mu$  window (by definition). Lastly, but perhaps most importantly, if we sum over  $\tau$  but with no explicit  $\tau$  weighting, the resulting mass distributions are identical,  $\int d\tau \tilde{F}_1^Q(\mu, \tau) = \int d\tau F_1^Q(\mu, \tau)$ . To make this last point explicit we note the following results for the moments of these two distributions,

$$\begin{aligned}
 \tilde{N}_\Omega^Q &= \int_\Omega d\mu \int_0^1 d\tau \tilde{F}_1^Q(\mu, \tau) = \tilde{\epsilon} = \langle \tilde{\tau}^Q \rangle = \int_\Omega d\mu \int_0^1 d\tau F_1^Q(\mu, \tau) = N_\Omega^Q, \\
 \langle \tilde{\mu}^Q \rangle &= \frac{1}{\tilde{N}_\Omega^Q} \int_\Omega d\mu \int_0^1 d\tau \mu \tilde{F}_1^Q(\mu, \tau) = \frac{1}{N_\Omega^Q} \int_\Omega d\mu \int_0^1 d\tau \mu F_1^Q(\mu, \tau) = \langle \mu^Q \rangle, \\
 \left(\sigma_{\tilde{\mu}^Q}^Q\right)^2 &= \frac{1}{\tilde{N}_\Omega^Q} \int_\Omega d\mu \int_0^1 d\tau (\mu - \langle \tilde{\mu}^Q \rangle)^2 \tilde{F}_1^Q(\mu, \tau) = \frac{1}{N_\Omega^Q} \int_\Omega d\mu \int_0^1 d\tau (\mu - \langle \mu^Q \rangle)^2 F_1^Q(\mu, \tau) \\
 &= \left(\sigma_{\mu^Q}^Q\right)^2.
 \end{aligned} \tag{5.1}$$

These equalities should help to confirm that comparing the Qjet and  $\tilde{Q}$ jet analyses, as in table 1, focuses on the statistical effects, while comparing the  $\tilde{Q}$ jet analysis with the conventional analysis focuses primarily on the physics effects caused by the changes in the mass distributions, as we will discuss in section 6.

With these insights, we can construct an explicit toy model that helps to illuminate the connection between the two distributions  $F_1^Q$  and  $\tilde{F}_1^Q$ . We can approximate the filled bins (closest to the boundaries) as being described by delta functions (recall the description of the conventional result in eq. (2.2)). Considering first adding just a single extra bin near the upper boundary, we have

$$\begin{aligned}
 \tilde{F}_1^Q(\tau) &= (1 - \tilde{\epsilon})\delta(\tau) + \tilde{\epsilon} \delta(\tau - 1) \\
 F_1^Q(\tau) &\simeq \tilde{F}_1^Q(\tau) - \Delta[\delta(1 - \tau) - \delta(1 - \eta - \tau)],
 \end{aligned} \tag{5.2}$$

where (as shown in eq. (2.5))  $\tilde{F}_1^Q(\tau)$  is represented by a binomial representation with mean  $\tilde{\epsilon}$  and variance  $\sigma_\tau^2 = \tilde{\epsilon}(1 - \tilde{\epsilon})$ . The extra term in the expression for  $F_1^Q(\tau)$  is intended to

present the fact that a small fraction of the jets,  $\Delta$ , have migrated from the  $\tau = 1$  bin to the  $\tau = (1 - \eta)$  bin ( $0 < \eta < 1$ ). It is straightforward to evaluate the corresponding approximate mean and variance of  $F_1^Q(\tau)$  in the limit  $\Delta \ll \tilde{\epsilon}$  in terms of the mean and variance of  $\tilde{F}_1^Q(\tau)$ . To first order in  $\Delta/\tilde{\epsilon}$  we find

$$\begin{aligned} \langle \tau \rangle &\simeq \tilde{\epsilon} - \Delta\eta, \\ \sigma_\tau^2 &\simeq \tilde{\epsilon}(1 - \tilde{\epsilon}) + \Delta\eta(\eta - 2(1 - \tilde{\epsilon})) = \sigma_{\tilde{\tau}}^2 + \Delta\eta(\eta - 2(1 - \tilde{\epsilon})). \end{aligned} \tag{5.3}$$

Applying this result to the first few column of table 1 we obtain

$$\frac{\delta S^Q/\sqrt{S^Q}}{\delta S^{\tilde{Q}}/\sqrt{S^{\tilde{Q}}}} = \frac{\delta S^Q}{\sqrt{S^Q}} = \sqrt{\langle \tau_s \rangle + \frac{\sigma_{\tau_s}^2}{\langle \tau_s \rangle}} \simeq 1 - \frac{\Delta}{2\tilde{\epsilon}}\eta(1 - \eta) \leq 1. \tag{5.4}$$

Noting that this expression is symmetric in  $\eta \rightarrow 1 - \eta$ , we see that the bins at both ends of the  $\tau$  distribution will contribute in a similar fashion, decreasing the scaled fluctuations in this observable. So, if we define a more accurate approximate expression for  $F_1(\tau)$ , including all of the filled in bins ( $\eta_k$  near 0 and near 1),

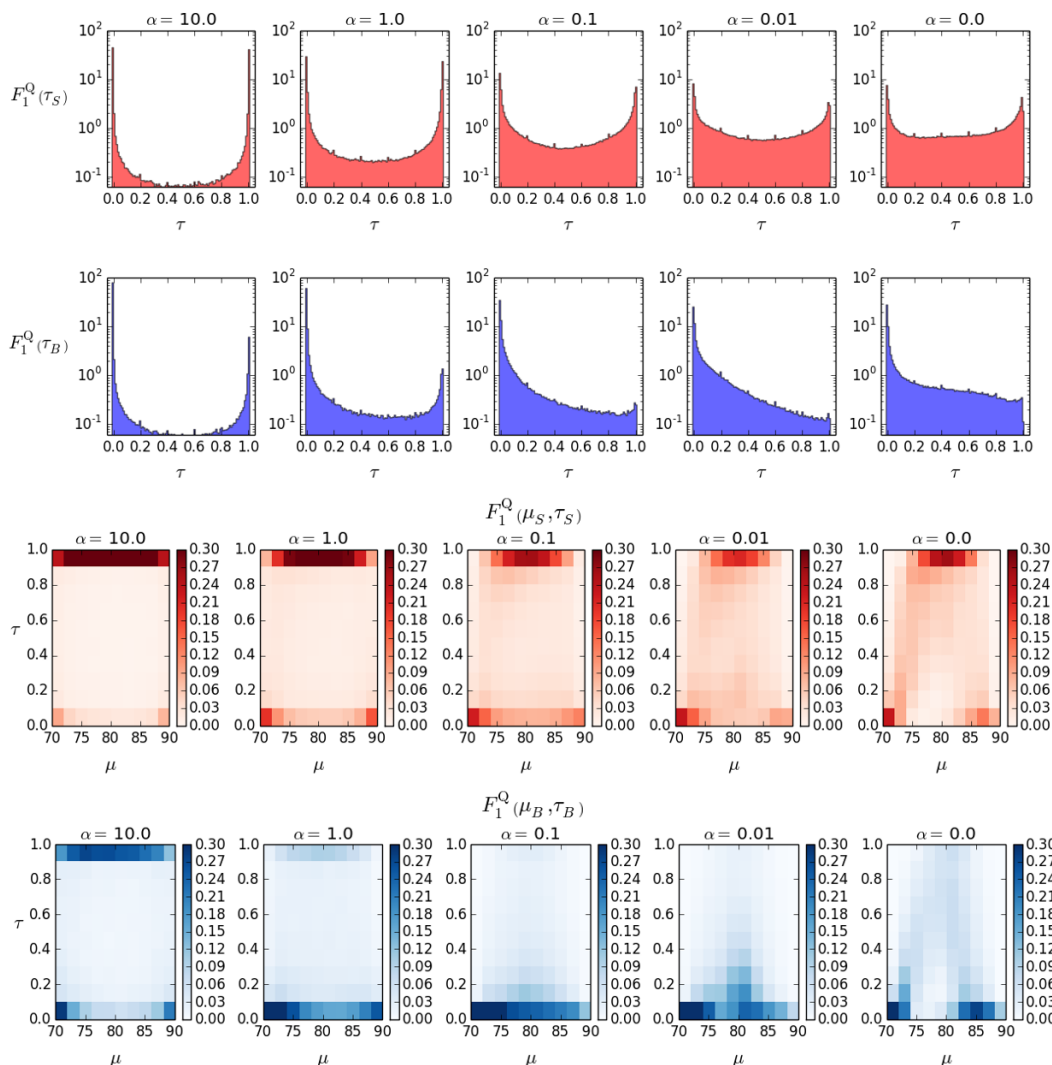
$$F_1^Q(\tau) \simeq \tilde{F}_1^Q(\tau) - \sum_k \Delta_k [\delta(1 - \tau) - \delta(1 - \eta_k - \tau)], \tag{5.5}$$

we find

$$\frac{\delta S^Q/\sqrt{S^Q}}{\delta S^{\tilde{Q}}/\sqrt{S^{\tilde{Q}}}} \simeq 1 - \sum_k \frac{\Delta_k}{2\tilde{\epsilon}} \eta_k (1 - \eta_k) \leq 1. \tag{5.6}$$

Since  $0 \leq \eta_k, \tilde{\epsilon} \leq 1$  and  $\Delta_k > 0$ , *all* of the terms in the sum serve to decrease these fluctuations (at least to leading order in  $\Delta_k/\tilde{\epsilon}$ ). As rigidity is decreased and the analysis moves further from the “classical” limit, we expect more bins away from the edges to be filled-in, which explains, at least qualitatively, the systematic decrease with decreasing rigidity (at least until we reach zero rigidity) in the first two columns, both signal and background, in table 1. Note that the deviation of the l.h.s. from 1 in eq. (5.6) is essentially proportional to the factor  $\sum_k \Delta_k/\tilde{\epsilon}$ . In our toy example, this represents the fraction of jets that occupied the  $\tau = 1$  bin in  $\tilde{F}_1^Q(\tau)$ , but correspond to a smaller  $\tau$  value in  $F_1^Q(\tau)$ . As we mentioned earlier, these jets have masses near the boundary of the mass window and we can say that  $\sum_k \Delta_k/\tilde{\epsilon}$  represents the fraction of jets that reside near the mass boundary (at least in our toy example) and that, after the full Qjets procedure, exhibit less than unit tagging probability.

To provide a more detailed picture of the rigidity ( $\alpha$ ) dependence exhibited in table 1, figure 5 shows plots of  $F_1^Q(\tau)$  and  $F_1^Q(\mu, \tau)$  for various values of the rigidity with  $z_{\text{cut}} = 0.1$ . To appreciate these plots it is important to note a couple of relevant features. In order to make visible the values at intermediate  $\tau$  values, the plots of  $F_1^Q(\tau)$  are semi-log plots, while the 2-D  $F_1^Q(\mu, \tau)$  plots use a linear scale for the color scale. Further, the  $\tau = 0$  bin at the extreme left of the  $F_1^Q(\tau)$  plots corresponds to  $\mu$  values not shown in the 2-D  $F_1^Q(\mu, \tau)$  plots, which show only the  $\mu$  values in the window  $\Omega$ . To understand the general structure of the plots in figure 5 recall that, as we lower the value of the rigidity parameter  $\alpha$ , the Qjets procedure explores an ever broader spectrum of clustering histories. This



**Figure 5.** The distributions  $F_1(\tau)$  and  $F_1(\mu, \tau)$  for signal (in red) and background (in blue) jets as functions of  $\alpha$ .

leads to changes in the plots arising from two distinct underlying effects, where both are associated with  $\mu$  values near the boundaries of the window  $\Omega$ . Individual jets, whose Qjet mass distributions are entirely within  $\Omega$  for large  $\alpha$  values, i.e., appear in the  $\tau = 1$  bin, may eventually exhibit Qjets mass distributions with tails that extend outside of  $\Omega$  for sufficiently small  $\alpha$  values. So, as the value of  $\alpha$  decreases, such jets will gradually populate the bins with  $\tau < 1$ , but still with  $\mu$  near the boundaries, at least initially. Likewise there are jets whose large  $\alpha$  Qjets mass distributions are entirely outside of  $\Omega$ , i.e., appear in the  $\tau = 0$  bin, but then develop tails inside of  $\Omega$  for sufficiently small  $\alpha$  values. These jets will gradually populate the bins at small  $\tau$  values, always moving inward from the boundaries in  $\mu$ . This simple picture is generally correct for both the signal and background samples as illustrated in figure 5. As  $\alpha$  decreases the distributions in  $\tau$ ,  $F_1^Q(\tau)$ , for both signal and

$\alpha$	$\frac{\tilde{S}^Q/\delta\tilde{B}^Q}{S^C/\delta B^C}$		$\tilde{S}^Q/S^C$		$\sqrt{B^C/\tilde{B}^Q}$	
	$z_{\text{cut}}$		$z_{\text{cut}}$		$z_{\text{cut}}$	
	0.10	0.15	0.10	0.15	0.10	0.15
10.0	0.90	0.91	1.16	1.18	0.776	0.774
1.00	0.82	0.85	1.48	1.61	0.554	0.530
0.10	0.77	0.82	1.81	2.12	0.427	0.385
0.01	0.76	0.81	1.91	2.32	0.399	0.351
0.00	0.78	0.83	1.93	2.37	0.406	0.350

**Table 2.** The physics component of the cross-section uncertainties as functions of  $z_{\text{cut}}$  and rigidity  $\alpha$ , found by dividing the total uncertainty by the purely statistical component (from table. 1). We also provide numerical values of different components in eq. (6.1) responsible for the physics part of the cross-section uncertainties.

background, gradually fill-in at intermediate  $\tau$  values leading to decreasing values of the fluctuation  $\sigma_\tau^Q$ . In turn this means that  $\delta N_T/\sqrt{N_T}$  decreases with decreasing  $\alpha$  values as indicated by the numbers in table 1, and the toy model analysis of eq. (5.6).

In the limit  $\alpha \rightarrow 0$  the  $F_1^Q(\mu, \tau)$  plots for both signal (red) and background (blue) suggest an “arc” structure, i.e., a ridge of enhanced probability connecting the  $\mu$  boundaries at small  $\tau$  via bins at intermediate  $\mu$  values at larger  $\tau$  values. The primary distinction between signal and background is the fact that the bins near  $\tau = 1$  for the background are rapidly depopulated as  $\alpha$  decreases, while for the signal the bins near  $\tau = 1$  maintain a population similar to those near  $\tau = 0$  and the bins near  $\tau = 1$  always exhibit a  $\mu$  distribution with a peak near the signal mass ( $M_W$ ).

Figure 5 also indicates that for large rigidity (say,  $\alpha \geq 1.0$ ), the probabilities for finding jets with intermediate  $\tau$  values ( $0 < \tau < 1$ ), are tiny. In this range of  $\alpha$ , the approximations made in our toy model above are quite appropriate. For smaller rigidity ( $\alpha < 1.0$ ), more and more jets occupy the intermediate  $\tau$  values, new patterns emerge in the pdf  $F_1^Q(\mu, \tau)$ , and the deviations of  $F_1^Q$  from  $\tilde{F}_1^Q$  are not necessarily small.

## 6 Understanding the physics effects

We list the components of the cross-section and mass uncertainties due to what we have labeled physics effects in table 2 and table 3 respectively. The numerical values of these components can be evaluated by dividing the total uncertainty in table 1 by its statistical part (also listed in table 1).

Understanding these physics quantities is relatively easier since one does not need to think of fractional tagging efficiencies. In particular, since both the conventional and hybrid analyses use binary tagging efficiencies, the fluctuations in the number of jets goes like  $1/\sqrt{N}$  (recall the discussion in section 2). For example, the quantity in the table 2, which measures the improvement in the cross-section measurement significance, simplifies to

$$\frac{\tilde{S}^Q/\delta\tilde{B}^Q}{S^C/\delta B^C} = \left(\frac{\tilde{S}^Q}{S^C}\right) \times \sqrt{\frac{B^C}{\tilde{B}^Q}}. \tag{6.1}$$



The improvement in statistical stability, therefore, depends on two independent ratios, the relative signal efficiency ( $\tilde{S}^Q/S^C$ ) and (1 over) the square root of the relative background efficiency ( $\tilde{B}^Q/B^C$ ), for the hybrid Qjets analysis compared to the conventional analysis. Table 2 separately exhibits the variation of these two components of eq. (6.1) with the rigidity parameter  $\alpha$ . To understand the exhibited behavior, we must recall our previous discussion. As we decrease  $\alpha$ , we include new clustering histories, and, as a result, find that jets, which were previously not tagged (for larger  $\alpha$  values), are now tagged. By construction  $\tilde{\tau}^Q = 1$  for these jets (even though they may have small  $\tau^Q$ ). This is why the ratio  $\tilde{N}_T^Q/N_T^C$  increases with decreasing  $\alpha$  for both signal and background. In the case of the (signal)  $W$ -jets, almost all jets are tagged even for large  $\alpha$  and so  $\tilde{S}^Q/S^C$  increases relatively slowly (but monotonically) as  $\alpha$  decreases. In the language of the simple model in the previous section (see eq. (5.6)), the behavior of the ratio  $\tilde{S}^Q/S^C$  is telling us about the magnitude of the leak-in effect,  $\sum_k \Delta_k/\tilde{\epsilon}$ , at least quantitatively (note that the effect is no longer small as  $\alpha$  approaches zero).

In the case of background jets, there are always more untagged jets than tagged ones, some of which can be tagged when we allow a broader range of clustering histories as  $\alpha$  decreases. Thus the ratio  $\tilde{B}^Q/B^C$  increases quite rapidly with decreasing  $\alpha$ , resulting in the somewhat slower but still rapid decrease of the factor  $\sqrt{B^C/\tilde{B}^Q}$ . By eq. (6.1), the physics component of the cross-section uncertainty in table 2 is the product of the corresponding values in the two right-hand columns in the table. Numerically the decrease of the background ratio is dominant, leading to a slowly decreasing cross-section uncertainty in the hybrid Qjets analysis compared to the conventional analysis with decreasing  $\alpha$  until  $\alpha$  reaches 0.01. For even smaller  $\alpha$  values the sampling of clustering histories is so broad that the qualitative behavior of the background ratio changes and the relative fluctuations begin to grow.

Note also that the variation with  $\alpha$  of the individual ratios, and the product, is somewhat stronger for the non-optimal  $z_{\text{cut}}$  value (0.15). This is to be expected as the non-optimal conventional result implies that more of the added clustering histories in the Qjets analysis will correspond to an improvement. Note that this is a statement about the *improvement* in the statistical stability. Overall one is better off starting with an optimal choice of the conventional pruning parameters to perform the Qjets procedure around. However, the results in table 2 do suggest that the Qjets procedure can help to moderate the impact of any initial poor choice of parameters.

Finally we turn to the mass measurement uncertainties as described by the results in table 3. The general expressions from eqs. (2.27) and (2.28) for the signal sample yield

$$\frac{\delta\tilde{m}_T^Q/\tilde{m}_T^Q}{\delta m_T^C/m_T^C} = \sqrt{\frac{\langle\tau_S^C\rangle}{\langle\tilde{\tau}_S^Q\rangle}} \times \frac{\sigma_{\tilde{\mu}_S}^Q}{\sigma_{\mu_S}^C} \times \frac{\langle\mu_S^C\rangle}{\langle\tilde{\mu}_S^Q\rangle} = \sqrt{\frac{S^C}{\tilde{S}^Q}} \times \frac{\sigma_{\tilde{\mu}_S}^Q}{\sigma_{\mu_S}^C} \times \frac{\langle\mu_S^C\rangle}{\langle\tilde{\mu}_S^Q\rangle} \quad (6.2)$$

The relative stability in the mass measurement depends on three important ratios, the relative signal efficiency ( $S^C/\tilde{S}^Q$ ), the relative fluctuation in the mass spectra ( $\sigma_{\tilde{\mu}_S}^Q/\sigma_{\mu_S}^C$ ), and the relative average mass ( $\langle\mu_S^C\rangle/\langle\tilde{\mu}_S^Q\rangle$ ). Table 3 exhibits the variation of these quantities with  $\alpha$ . As discussed in the previous paragraphs (and indicated in also table 2)  $\sqrt{S^C/\tilde{S}^Q}$



$\alpha$	$\frac{\delta\tilde{m}_T^Q/\tilde{m}_T^Q}{\delta m_T^C/m_T^C}$		$\sqrt{S^C/\tilde{S}^Q}$		$\sigma_{\tilde{\mu}_S}^Q/\sigma_{\mu_S}^C$		$\langle\mu_S^C\rangle/\langle\tilde{\mu}_S^Q\rangle$	
	$z_{\text{cut}}$		$z_{\text{cut}}$		$z_{\text{cut}}$		$z_{\text{cut}}$	
	0.10	0.15	0.10	0.15	0.10	0.15	0.10	0.15
10.0	0.96	0.95	0.93	0.92	1.03	1.03	1.00	1.00
1.00	0.86	0.82	0.82	0.79	1.05	1.04	1.00	1.00
0.10	0.73	0.66	0.74	0.69	0.98	0.95	1.01	1.01
0.01	0.66	0.56	0.72	0.66	0.91	0.86	1.01	1.02
0.00	0.69	0.57	0.72	0.65	0.95	0.86	1.01	1.02

**Table 3.** The physics component of the mass uncertainties as functions of  $z_{\text{cut}}$  and rigidity  $\alpha$ , found by dividing the total uncertainty by the purely statistical component (from table. 1). We also provide numerical values of different components in eq. (6.2) responsible for the physics part of the mass uncertainties.

is a slowly but monotonically decreasing function as  $\alpha$  decreases due to the increasing set of tagged jets in the hybrid analysis, i.e., the jets leaking-in at the edge of the window  $\Omega$  as measured by the quantity  $\sum_k \Delta_k/\tilde{\epsilon}$  in our simple model.

As shown in table 3, the average jet mass remains relatively constant ( $\langle\tilde{\mu}_S\rangle \simeq \langle\mu_S^C\rangle \simeq 80\text{ GeV}$ ) for all values of  $\alpha$ . In terms of the simple model presented in the previous section, the shift in the average jet mass (in the window  $\Omega$ ) in going from the conventional analysis to the hybrid analysis is proportional to the *difference* between the number of jets leaking-in from the upper edge of the window and the number leaking-in at the lower edge,  $(\sum_{k+} \Delta_{k+} - \sum_{k-} \Delta_{k-})/\tilde{\epsilon}$  (recall that the counting analysis, see eq. (5.6), involved the simple sum of these contributions). Since this leaking-in process is quite symmetrical (i.e., the signal sample itself is quite symmetrical about  $M_W$  with nearly identical numbers of jets just outside the window at both ends), any shift in the average jet mass is expected to be quite small, i.e., much smaller than the shift seen in the quantity  $\sqrt{S^C/\tilde{S}^Q}$ , in agreement with the results in table 3.

The last quantity (namely, the the relative fluctuation in the mass spectra,  $\sigma_{\tilde{\mu}_S}^Q/\sigma_{\mu_S}^C$ ) is especially interesting. Table 3 shows that this ratio first increases with decreasing  $\alpha$ , and then decreases. The simple model of the previous section suggests that the size of the deviation from unity for the ratio is again set by the fraction of tagged jets that are leaking in,  $\sum_k \Delta_k/\tilde{\epsilon}$ , but now with a coefficient that, not surprisingly, depends on the *shapes* of the jet mass distributions. The changes in the mass distribution can be qualitatively understood as follows. As  $\alpha$  is decreased and we move away from the conventional analysis, the initial change in the mass distribution is the leaking-in of jets just outside the mass window  $\Omega$  into mass bins just inside the window (as is evident in figure 5). Thus initially the mass distribution in the hybrid analysis is *broader* than in the conventional analysis and  $\sigma_{\tilde{\mu}_S}^Q/\sigma_{\mu_S}^C$  increases above unity with decreasing  $\alpha$ . However, eventually, as the mass distribution fills in the central region of the window (again see figure 5), the Qjets mass distribution again has a width similar to the conventional case and  $\sigma_{\tilde{\mu}_S}^Q/\sigma_{\mu_S}^C$  goes back to unity (for  $\alpha$  just above 0.1 in table 3). With a further decrease of  $\alpha$ ,  $0.1 > \alpha \geq 0$ , the

results in table 3 indicate that the jet mass distribution found by the Qjets procedure is *narrower* than the one found by pruning alone, i.e., the Qjets procedure provides a more efficient groomer than conventional or classical pruning.

Overall the relative uncertainty in the tagged mass measurement for the hybrid analysis versus the conventional analysis *decreases* with decreasing  $\alpha$  and the hybrid result becomes approximately 30% smaller than the fluctuations in the conventional analysis, i.e., it is the  $\sqrt{S^C/\tilde{S}^Q}$  factor that effectively controls the  $\alpha$  dependence shown in table 3. What we have labeled the physics part of the mass measurement uncertainty is minimized for  $0.01 \geq \alpha \geq 0$ .

## 7 Conclusions

The Qjets procedure is intuitively motivated by the idea that analyses of jet observables that depend on clustering histories can be improved by considering multiple clustering histories of a jet. On the other hand, the statistical treatment of the results can be unintuitive and opaque. Much of the confusion lies in the fact that, while all observables in the Qjets procedure are weighted with weights following a continuous distribution in the interval  $[0, 1]$ , the conventional approach applies no weight as long as jets are tagged, i.e., applies a simple binary weight. Even in sophisticated multivariate analyses, where many variables are combined in a likelihood and each jet/event is assigned a likelihood (a continuous distribution in the interval  $[0, 1]$ ) for being a signal, the likelihood variable only provides a discriminatory variable to separate signal from background. The measurements are subsequently estimated from the tagged jet/event sample (i.e., the jets/events that pass the cut on likelihood to be signal) with only a binary (0 or 1) weight.

The purpose of this paper is to address this issue, namely, to provide a platform in which the uncertainties associated with the measurements in the Qjets procedure can be evaluated. We also propose an alternative way to calculate the uncertainties of measurements. Uncertainties are traditionally estimated using Monte Carlo pseudo-experiments, in which jets/events are picked at random from a given *master-sample* of jets/events (either carefully prepared using a Monte Carlo event generator, or control-samples from collider events), and then repeating pseudo-experiments several times. Variations of observables over pseudo-experiments then provide an estimate of statistical uncertainties. While this method is straightforward, it is time consuming (since pseudo-experiments need to be repeated many times), and still does not provide any insights regarding these measurements. In this work we choose a different framework — we provide analytic formulas in section 2, which relate these uncertainties with various moments of the given jets/events sample. On the one hand, these expressions provide much faster ways to measure uncertainties; while on the other, they help explain the physics of the uncertainties. We have also presented a simple model of how the Qjets procedure impacts the probability distributions in both the tagging efficiency  $\tau$  and the jet mass  $\mu$ , which provides further insight into the observed numerical results.

We find that, while Poisson uncertainties associated with measurements are unavoidable, sampling uncertainties can be reduced by using weighted jets such as those returned

by Qjets. We show that this additional stability in measurements provided by the Qjets procedure can arise from two qualitatively different sources — from the transition from unweighted to weighted measurements (which we label the statistics effects), and from the Qjets generated changes in the distributions of jet-observables themselves, e.g., jet masses, (which we label physics effects). Our explicit numerical results indicate how these two kinds of effects often compete with each other, and how they vary as various Qjets parameters, especially the rigidity  $\alpha$ , are altered. Overall, however, the Qjets procedure acts to improve both the statistical stability of counting experiments and the precision of the measurement of jet observables like the jet mass. Further, we have seen that the Qjets procedure can largely remove the negative impact of a less-than-optimum choice of jet grooming parameters on a conventional analysis.

Before we conclude, let us note that the results in this work can be easily generalized. We obtained the expressions for uncertainties only for cross-section and mass measurements. Uncertainties for any other weighted measurements in the Qjets procedure can be performed by following the treatment for the mass measurement. Also note that, in deriving these formulas, we explicitly talked about jets. However, we can easily use the same formalism when we need to talk about events. In fact, we choose one jet per event in our calculations. Therefore, the expressions for uncertainties associated with the number of jets observed (for example), is identical to the uncertainties associated with the number of events observed. It is straightforward to apply the framework introduced in this work to explain the statistical improvements claimed by the recent proposals such as “Telescoping Jets” [58] and “Jet Sampling” [59]. Finally, we also expect that sophisticated, state-of-the-art multivariate techniques can be made more robust by estimating measurements using weighted events with the likelihood variable as the weight. Such an analysis could presumably follow the framework laid out in this paper.

## Acknowledgments

The authors would like to thank Matthew D. Schwartz for his collaboration at an earlier stage of this work. This work was supported in part by the US Department of Energy under contracts DE-FGO2-96ER40956 and DE-SC003916, by a Simons postdoctoral fellowship, Director’s fellowships from LANL, an LHC-TI travel grant and by the KITP, under NSF grant PHY05-51164. The bulk for the computations were performed on the Mapache cluster in the HCP facility at LANL. Some preliminary computations were also performed on the Odyssey cluster at Harvard University and the TEV cluster in the University of Washington.

**Note added.** While this manuscript was being finalized ref. [60] appeared on the arXiv. Ref. [60] studies the statistical effects in counting experiments (i.e., cross-sections) for Qjet-like observables using pseudo-experiments in the context of ref. [59]. In contrast, in this manuscript we explore both cross-sections and more general measurements such as jet mass and provide an analytical framework for calculating their statistical properties in terms of probability density functions (the results are then validated using pseudo-experiments in appendix A).

## A Validation of section 2 with pseudo-experiments

Traditionally, statistical uncertainties of complicated observables are estimated by using Monte-Carlo pseudo-experiments. In this procedure, one generates many sets of events, where the number of events is chosen according to a Poisson distribution with a given mean (see eq. (2.1)). One then measures the quantity of interest on each set of events, and, by considering the variation of the quantity across many pseudo-experiments, one can estimate the statistical uncertainty of the measurement considered. This procedure simultaneously accounts for both Poisson and sampling uncertainties.

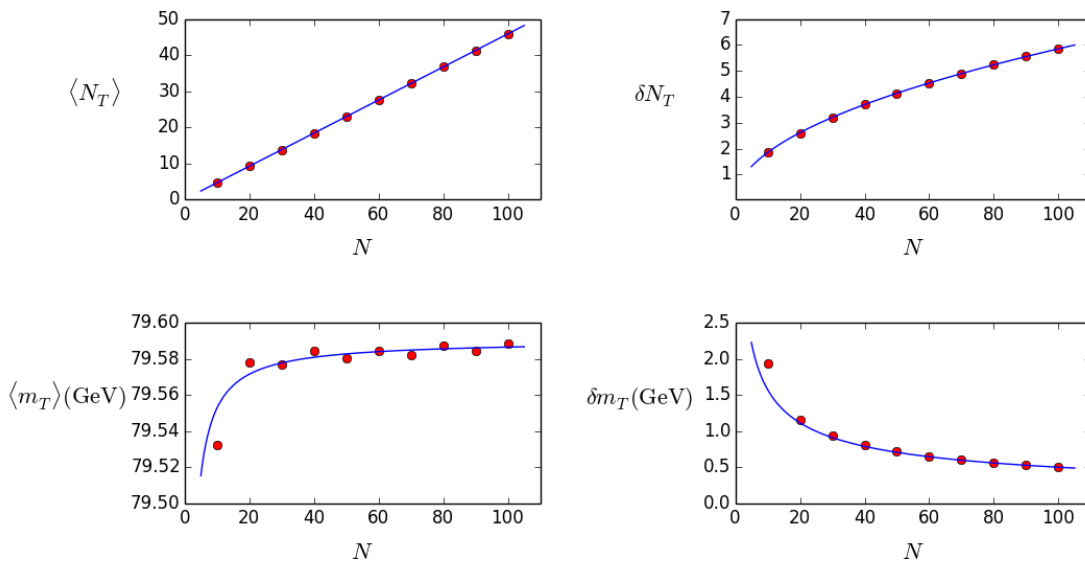
In this work, we advocate for a different method of calculating statistical uncertainties. As shown in section 2, analytical expressions may be derived, which relate these uncertainties to various moments of a probability distribution constructed from a sample of events. These analytical formulas carry more information than just performing Monte-Carlo pseudo-experiments, since they (like all analytical derivations) also explain “why the numbers are what they are.” One can use this improved understanding to devise ways to attempt to reduce uncertainties further.

The purpose of this section of the appendix is to validate the formulas derived in section 2 using pseudo-experiments. In order to do this, we choose a sample of  $W$ -jets (in fact, we choose the same set of hadronic  $WW$ -events as outlined in section 4, and use the same procedure and parameters to construct  $W$ -jets out of these events). We perform  $10^5$  pseudo-experiments, in each of which  $n$   $W$ -jets are chosen at random. As explained above (see especially section 2),  $n$  follows a Poisson distribution with mean  $N$ . Jets chosen in a pseudo-experiment, are then subjected to the Qjets procedure for a particular set of parameters ( $\alpha = 0.01, z_{\text{cut}} = 0.1, D_{\text{cut}} = m/p_T, \Omega = (70 - 90)$  GeV). Using the outputs of the Qjets procedure, we calculate observables  $N_T$  (number of jets tagged in an experiment) and  $m_T$  (tagged mass in the experiment) for each pseudo-experiment. The variations of these observables over the set of pseudo-experiments provides estimates of the statistical uncertainties  $\delta N_T$  and  $\delta m_T$ . We also estimate the same uncertainties using the analytic expressions derived in section 2 and compare them.

In figure 6 we compare the analytic estimates of the mean values and the uncertainties, represented by the blue lines, and the numerical values from the pseudo-experiments (the red points) as a function of the mean number of  $W$ -jets  $N$ . We begin with a measurement of the cross section in the top row of figure 6. Here we see that the average and the uncertainty in the number of tagged events follow essentially exactly the distribution in eq. (2.18). Next, consider a measurement of the average jet mass, as in the bottom row of figure 6. Here we see that the uncertainty falls as  $1/\sqrt{N}$  up to corrections whose effects are captured by terms of  $\mathcal{O}(1/N^2)$  in the formulas in section 2.

## B Jet mass

In this section of the appendix we derive the analytical expressions relevant for estimating the uncertainties associated with mass. In particular, we derive eqs. (2.25)–(2.28). For the sake of completeness, we repeat a few of the definitions introduced in section 2.2.



**Figure 6.** Left: variation of  $\langle N_T \rangle$ ,  $\delta N_T$ ,  $\langle m_T \rangle$ , and  $\delta m_T$  as functions of  $N$  for a sample of  $W$ -jets. See the text for details of the Qjets parameters used. The analytical results (calculated using formulas derived in section 2) are represented by blue lines. The red points denote the same quantities evaluated using Monte Carlo pseudo-experiments.

In an experiment where  $N_S$  jets with masses  $\{\mu_j\}$  and tagging probabilities  $\{\tau_j\}$  are chosen at random, the measured tagged mass is given by

$$m_T = \frac{\sum_{j=1}^{N_S} \mu_j \tau_j}{\sum_{j=1}^{N_S} \tau_j}.$$

We are interested in the average and the variance of  $m_T$ . Given the fact that the experiment started with  $N_S$  jets, we have

$$\begin{aligned} \langle m_T \rangle_{N_S} &= \left\langle \frac{\sum_{j=1}^{N_S} \mu_j \tau_j}{\sum_{j=1}^{N_S} \tau_j} \right\rangle_{N_S} \equiv \left\langle \frac{M_T}{N_T} \right\rangle_{N_S}, \quad \text{and} \\ (\delta m_T)_{N_S}^2 &= \langle (m_T - \langle m_T \rangle_{N_S})^2 \rangle_{N_S} = \langle m_T^2 \rangle_{N_S} - \langle m_T \rangle_{N_S}^2. \end{aligned}$$

In these expressions the notation  $M_T \equiv m_T N_T$  is used in order to simplify the results. We note that the probability distribution for  $M_T$  and  $N_T$ , for a given sample of size  $N_S$ , can be constructed in terms of  $F_1(\mu, \tau)$ ,

$$F_{N_S}(M_T, N_T) = \left[ \prod_{k=1}^{N_S} \int F_1(\mu_k, \tau_k) d\mu_k d\tau_k \right] \delta \left( N_T - \sum_{k=1}^{N_S} \tau_k \right) \delta \left( M_T - \sum_{k=1}^{N_S} \mu_k \tau_k \right). \quad (\text{B.1})$$

The relevant moments of this general distribution can be derived in terms of the moments

of  $F_1$  by repeating the manipulations in eqs. (2.9)–(2.10). We have

$$\langle N_T \rangle_{N_S} = \int dM_T dN_T N_T F_{N_S}(M_T, N_T) = N_S \langle \tau \rangle \quad (\text{B.2})$$

$$\langle M_T \rangle_{N_S} = \int dM_T dN_T M_T F_{N_S}(M_T, N_T) = N_S \langle \mu \tau \rangle \quad (\text{B.3})$$

$$\langle N_T^2 \rangle_{N_S} = \int dM_T dN_T N_T^2 F_{N_S}(M_T, N_T) = N_S^2 \langle \tau \rangle^2 + N_S \sigma_\tau^2 \quad (\text{B.4})$$

$$\langle M_T^2 \rangle_{N_S} = \int dM_T dN_T M_T^2 F_{N_S}(M_T, N_T) = N_S^2 \langle \mu \tau \rangle^2 + N_S \sigma_{\mu\tau}^2 \quad (\text{B.5})$$

$$\langle M_T N_T \rangle_{N_S} = \int dM_T dN_T M_T N_T F_{N_S}(M_T, N_T) = N_S^2 \langle \mu \tau \rangle \langle \tau \rangle + N_S \sigma(\tau, \mu\tau) \quad (\text{B.6})$$

Now we are ready to estimate the mean and variance of the tagged mass,  $m_T$ , distribution. These calculations are slightly non-trivial since  $m_T$  is a ratio of two independent variables. We use a Taylor series expansion to simplify the results. In particular, note that a generic bivariate function  $f(x, y)$  can be expanded using

$$f(x, y) \simeq f(x_0, y_0) + \left. \frac{\partial f}{\partial x} \right|_{x_0, y_0} (x - x_0) + \left. \frac{\partial f}{\partial y} \right|_{x_0, y_0} (y - y_0) \quad (\text{B.7})$$

$$+ \frac{1}{2} \left[ \left. \frac{\partial^2 f}{\partial x^2} \right|_{x_0, y_0} (x - x_0)^2 + \left. \frac{\partial^2 f}{\partial y^2} \right|_{x_0, y_0} (y - y_0)^2 + 2 \left. \frac{\partial^2 f}{\partial x \partial y} \right|_{x_0, y_0} (x - x_0)(y - y_0) \right] + \dots$$

Therefore, treating  $m_T$  as a function of  $M_T$  and  $N_T$ , we can expand  $m_T$  around  $M_T = \langle M_T \rangle_{N_S}$  and  $N_T = \langle N_T \rangle_{N_S}$ . We find that

$$m_T \simeq \frac{\langle M_T \rangle_{N_S}}{\langle N_T \rangle_{N_S}} + \frac{\langle M_T \rangle_{N_S}}{\langle N_T \rangle_{N_S}^3} (N_T - \langle N_T \rangle_{N_S})^2 \quad (\text{B.8})$$

$$- \frac{1}{\langle N_T \rangle_{N_S}^2} (N_T - \langle N_T \rangle_{N_S}) (M_T - \langle M_T \rangle_{N_S}) + \dots$$

It is now straightforward to find the average

$$\langle m_T \rangle_{N_S} \simeq \frac{\langle \mu \tau \rangle}{\langle \tau \rangle} \left[ 1 + \frac{\sigma_\tau^2}{N_S \langle \tau \rangle^2} - \frac{\sigma(\tau, \mu\tau)}{N_S \langle \mu \tau \rangle \langle \tau \rangle} \right] + \dots \quad (\text{B.9})$$

A similar expression can be derived for  $m_T^2$ ,

$$\langle m_T^2 \rangle_{N_S} \simeq \frac{\langle \mu \tau \rangle^2}{\langle \tau \rangle^2} \left[ 1 + \frac{\sigma_{\mu\tau}^2}{N_S \langle \mu \tau \rangle^2} + 3 \frac{\sigma_\tau^2}{N_S \langle \tau \rangle^2} - 4 \frac{\sigma(\tau, \mu\tau)}{N_S \langle \mu \tau \rangle \langle \tau \rangle} \right]. \quad (\text{B.10})$$

The final step in our calculation involves convolving with the Poisson distributions. This

yields

$$\langle m_T \rangle = \sum_{N_S=0}^{\infty} \text{Pois}(N_S|N) \langle m_T \rangle_{N_S} \simeq \frac{\langle \mu\tau \rangle}{\langle \tau \rangle} \left[ 1 + \frac{\sigma_\tau^2}{N\langle \tau \rangle^2} - \frac{\sigma(\tau, \mu\tau)}{N\langle \mu\tau \rangle \langle \tau \rangle} \right], \quad (\text{B.11})$$

$$\langle m_T^2 \rangle = \sum_{N_S=0}^{\infty} \text{Pois}(N_S|N) \langle m_T^2 \rangle_{N_S} \simeq \frac{\langle \mu\tau \rangle^2}{\langle \tau \rangle^2} \left[ 1 + \frac{\sigma_{\mu\tau}^2}{N\langle \mu\tau \rangle^2} + 3\frac{\sigma_\tau^2}{N\langle \tau \rangle^2} - 4\frac{\sigma(\tau, \mu\tau)}{N\langle \mu\tau \rangle \langle \tau \rangle} \right], \quad (\text{B.12})$$

$$(\delta m_T)^2 = \langle m_T^2 \rangle - \langle m_T \rangle^2 \simeq \frac{\langle \mu\tau \rangle^2}{N\langle \tau \rangle^2} \left[ \frac{\sigma_{\mu\tau}^2}{\langle \mu\tau \rangle^2} + \frac{\sigma_\tau^2}{\langle \tau \rangle^2} - 2\frac{\sigma(\tau, \mu\tau)}{\langle \mu\tau \rangle \langle \tau \rangle} \right]. \quad (\text{B.13})$$

In these expressions we have neglected terms of order  $1/N^2$  and higher.

Some simplifications arise for the case of the conventional tagging procedure. Since  $\tau$  is non-zero (and equal to one) only in the range  $\Omega$ , we find that ( $q > 0$ )

$$\langle (\mu^p \tau^q) \rangle^C = \int d\mu \int_0^1 d\tau \mu^p \tau^q F_1^C = \int_\Omega d\mu \int_0^1 d\tau \mu^p F_1^C = (N_\Omega \langle \mu^p \rangle)^C = \epsilon \langle (\mu^C)^p \rangle, \quad (\text{B.14})$$

where we use eq. (2.29) to derive the final expressions and borrow the notation  $\mu^C$  from eq. (2.28), to denote that the moment is to be calculated from eq. (2.29) using the conventional pdf  $F_1^C(\mu, \tau)$ .

Therefore we find the following identities (recall  $\langle \tau^C \rangle = \epsilon$ )

$$\left( \frac{\sigma_{\mu\tau}^2}{\langle \mu\tau \rangle^2} \right)^C = \frac{\langle (\mu^C)^2 \rangle - \epsilon \langle \mu^C \rangle^2}{\epsilon \langle \mu^C \rangle^2} \quad \text{and} \quad \left( \frac{\sigma(\tau, \mu\tau)}{\langle \mu\tau \rangle \langle \tau \rangle} \right)^C = \frac{\langle \mu^C \rangle (1 - \langle \tau^C \rangle)}{\langle \mu^C \rangle \langle \tau^C \rangle} = \frac{1 - \epsilon}{\epsilon} = \frac{(\sigma_\tau^C)^2}{\langle \tau^C \rangle^2}. \quad (\text{B.15})$$

The final expressions for the conventional average mass and its uncertainty then simplify to

$$\begin{aligned} \langle m_T^C \rangle &= \langle \mu^C \rangle, & (\delta m_T^C)^2 &= \frac{1}{N} \times \frac{1}{\epsilon} (\sigma_\mu^C)^2, \\ \left( \frac{\delta m_T^C}{\langle m_T^C \rangle} \right)^2 &= \frac{1}{N} \times \frac{1}{\epsilon} \frac{(\sigma_\mu^C)^2}{\langle \mu^C \rangle^2} \end{aligned} \quad (\text{B.16})$$

**Open Access.** This article is distributed under the terms of the Creative Commons Attribution License ([CC-BY 4.0](https://creativecommons.org/licenses/by/4.0/)), which permits any use, distribution and reproduction in any medium, provided the original author(s) and source are credited.

## References

- [1] M.H. Seymour, *Searches for new particles using cone and cluster jet algorithms: A Comparative study*, *Z. Phys. C* **62** (1994) 127 [[INSPIRE](#)].
- [2] J.M. Butterworth, J.R. Ellis and A.R. Raklev, *Reconstructing sparticle mass spectra using hadronic decays*, *JHEP* **05** (2007) 033 [[hep-ph/0702150](#)] [[INSPIRE](#)].
- [3] G. Brooijmans, *High  $p_T$  Hadronic Top Quark Identification*, [ATL-PHYS-CONF-2008-008](#) (2008).



- [4] J.M. Butterworth, A.R. Davison, M. Rubin and G.P. Salam, *Jet substructure as a new Higgs search channel at the LHC*, *Phys. Rev. Lett.* **100** (2008) 242001 [[arXiv:0802.2470](#)] [[INSPIRE](#)].
- [5] J.M. Butterworth, J.R. Ellis, A.R. Raklev and G.P. Salam, *Discovering baryon-number violating neutralino decays at the LHC*, *Phys. Rev. Lett.* **103** (2009) 241803 [[arXiv:0906.0728](#)] [[INSPIRE](#)].
- [6] D.E. Kaplan, K. Rehermann, M.D. Schwartz and B. Tweedie, *Top Tagging: A Method for Identifying Boosted Hadronically Decaying Top Quarks*, *Phys. Rev. Lett.* **101** (2008) 142001 [[arXiv:0806.0848](#)] [[INSPIRE](#)].
- [7] J. Thaler and L.-T. Wang, *Strategies to Identify Boosted Tops*, *JHEP* **07** (2008) 092 [[arXiv:0806.0023](#)] [[INSPIRE](#)].
- [8] L.G. Almeida et al., *Substructure of high- $p_T$  Jets at the LHC*, *Phys. Rev. D* **79** (2009) 074017 [[arXiv:0807.0234](#)] [[INSPIRE](#)].
- [9] T. Plehn, G.P. Salam and M. Spannowsky, *Fat Jets for a Light Higgs*, *Phys. Rev. Lett.* **104** (2010) 111801 [[arXiv:0910.5472](#)] [[INSPIRE](#)].
- [10] J. Thaler and K. Van Tilburg, *Identifying Boosted Objects with  $N$ -subjettiness*, *JHEP* **03** (2011) 015 [[arXiv:1011.2268](#)] [[INSPIRE](#)].
- [11] D.E. Soper and M. Spannowsky, *Finding physics signals with shower deconstruction*, *Phys. Rev. D* **84** (2011) 074002 [[arXiv:1102.3480](#)] [[INSPIRE](#)].
- [12] C. Englert, T.S. Roy and M. Spannowsky, *Ditau jets in Higgs searches*, *Phys. Rev. D* **84** (2011) 075026 [[arXiv:1106.4545](#)] [[INSPIRE](#)].
- [13] J. Thaler and K. Van Tilburg, *Maximizing Boosted Top Identification by Minimizing  $N$ -subjettiness*, *JHEP* **02** (2012) 093 [[arXiv:1108.2701](#)] [[INSPIRE](#)].
- [14] T. Plehn and M. Spannowsky, *Top Tagging*, *J. Phys. G* **39** (2012) 083001 [[arXiv:1112.4441](#)] [[INSPIRE](#)].
- [15] M. Jankowiak and A.J. Larkoski, *Jet Substructure Without Trees*, *JHEP* **06** (2011) 057 [[arXiv:1104.1646](#)] [[INSPIRE](#)].
- [16] S.D. Ellis, T.S. Roy and J. Scholtz, *Jets and Photons*, *Phys. Rev. Lett.* **110** (2013) 122003 [[arXiv:1210.1855](#)] [[INSPIRE](#)].
- [17] S.D. Ellis, T.S. Roy and J. Scholtz, *Phenomenology of Photon-Jets*, *Phys. Rev. D* **87** (2013) 014015 [[arXiv:1210.3657](#)] [[INSPIRE](#)].
- [18] D.E. Soper and M. Spannowsky, *Finding top quarks with shower deconstruction*, *Phys. Rev. D* **87** (2013) 054012 [[arXiv:1211.3140](#)] [[INSPIRE](#)].
- [19] A.J. Larkoski, G.P. Salam and J. Thaler, *Energy Correlation Functions for Jet Substructure*, *JHEP* **06** (2013) 108 [[arXiv:1305.0007](#)] [[INSPIRE](#)].
- [20] A.J. Larkoski, I. Moult and D. Neill, *Power Counting to Better Jet Observables*, *JHEP* **12** (2014) 009 [[arXiv:1409.6298](#)] [[INSPIRE](#)].
- [21] N.G. Ortiz, J. Ferrando, D. Kar and M. Spannowsky, *Reconstructing singly produced top partners in decays to  $Wb$* , *Phys. Rev. D* **90** (2014) 075009 [[arXiv:1403.7490](#)] [[INSPIRE](#)].
- [22] J.M. Butterworth, A.R. Davison, M. Rubin and G.P. Salam, *Jet substructure as a new Higgs search channel at the LHC*, *AIP Conf. Proc.* **1078** (2009) 189 [[arXiv:0809.2530](#)] [[INSPIRE](#)].



- [23] J.M. Butterworth, A.R. Davison, M. Rubin and G.P. Salam, *Jet substructure as a new Higgs search channel at the LHC*, [arXiv:0810.0409](#) [INSPIRE].
- [24] S.D. Ellis, C.K. Vermilion and J.R. Walsh, *Techniques for improved heavy particle searches with jet substructure*, *Phys. Rev. D* **80** (2009) 051501 [[arXiv:0903.5081](#)] [INSPIRE].
- [25] S.D. Ellis, C.K. Vermilion and J.R. Walsh, *Recombination Algorithms and Jet Substructure: Pruning as a Tool for Heavy Particle Searches*, *Phys. Rev. D* **81** (2010) 094023 [[arXiv:0912.0033](#)] [INSPIRE].
- [26] D. Krohn, J. Thaler and L.-T. Wang, *Jet Trimming*, *JHEP* **02** (2010) 084 [[arXiv:0912.1342](#)] [INSPIRE].
- [27] G. Soyez, G.P. Salam, J. Kim, S. Dutta and M. Cacciari, *Pileup subtraction for jet shapes*, *Phys. Rev. Lett.* **110** (2013) 162001 [[arXiv:1211.2811](#)] [INSPIRE].
- [28] D. Krohn, M.D. Schwartz, M. Low and L.-T. Wang, *Jet cleansing: Separating data from secondary collision induced radiation at high luminosity*, *Phys. Rev. D* **90** (2014) 065020 [[arXiv:1309.4777](#)] [INSPIRE].
- [29] M. Dasgupta, A. Fregoso, S. Marzani and G.P. Salam, *Towards an understanding of jet substructure*, *JHEP* **09** (2013) 029 [[arXiv:1307.0007](#)] [INSPIRE].
- [30] M. Dasgupta, A. Fregoso, S. Marzani and A. Powling, *Jet substructure with analytical methods*, *Eur. Phys. J. C* **73** (2013) 2623 [[arXiv:1307.0013](#)] [INSPIRE].
- [31] M. Cacciari, G.P. Salam and G. Soyez, *On the use of charged-track information to subtract neutral pileup*, [arXiv:1404.7353](#) [INSPIRE].
- [32] M. Cacciari, G.P. Salam and G. Soyez, *SoftKiller, a particle-level pileup removal method*, [arXiv:1407.0408](#) [INSPIRE].
- [33] D. Bertolini, P. Harris, M. Low and N. Tran, *Pileup Per Particle Identification*, *JHEP* **1410** (2014) 59 [[arXiv:1407.6013](#)] [INSPIRE].
- [34] A.J. Larkoski, S. Marzani, G. Soyez and J. Thaler, *Soft Drop*, *JHEP* **05** (2014) 146 [[arXiv:1402.2657](#)] [INSPIRE].
- [35] J. Gallicchio and M.D. Schwartz, *Seeing in Color: Jet Superstructure*, *Phys. Rev. Lett.* **105** (2010) 022001 [[arXiv:1001.5027](#)] [INSPIRE].
- [36] J. Gallicchio and M.D. Schwartz, *Pure Samples of Quark and Gluon Jets at the LHC*, *JHEP* **10** (2011) 103 [[arXiv:1104.1175](#)] [INSPIRE].
- [37] J. Gallicchio and M.D. Schwartz, *Quark and Gluon Tagging at the LHC*, *Phys. Rev. Lett.* **107** (2011) 172001 [[arXiv:1106.3076](#)] [INSPIRE].
- [38] J. Gallicchio and M.D. Schwartz, *Quark and Gluon Jet Substructure*, *JHEP* **04** (2013) 090 [[arXiv:1211.7038](#)] [INSPIRE].
- [39] D. Krohn, M.D. Schwartz, T. Lin and W.J. Waalewijn, *Jet Charge at the LHC*, *Phys. Rev. Lett.* **110** (2013) 212001 [[arXiv:1209.2421](#)] [INSPIRE].
- [40] A.J. Larkoski, J. Thaler and W.J. Waalewijn, *Gaining (Mutual) Information about Quark/Gluon Discrimination*, *JHEP* **11** (2014) 129 [[arXiv:1408.3122](#)] [INSPIRE].
- [41] A. Abdesselam et al., *Boosted objects: A Probe of beyond the Standard Model physics*, *Eur. Phys. J. C* **71** (2011) 1661 [[arXiv:1012.5412](#)] [INSPIRE].

- [42] A. Altheimer et al., *Jet Substructure at the Tevatron and LHC: New results, new tools, new benchmarks*, *J. Phys. G* **39** (2012) 063001 [[arXiv:1201.0008](#)] [[INSPIRE](#)].
- [43] A. Altheimer et al., *Boosted objects and jet substructure at the LHC. Report of BOOST2012, held at IFIC Valencia, 23rd-27th of July 2012*, *Eur. Phys. J. C* **74** (2014) 2792 [[arXiv:1311.2708](#)] [[INSPIRE](#)].
- [44] S.D. Ellis, A. Hornig, T.S. Roy, D. Krohn and M.D. Schwartz, *Qjets: A Non-Deterministic Approach to Tree-Based Jet Substructure*, *Phys. Rev. Lett.* **108** (2012) 182003 [[arXiv:1201.1914](#)] [[INSPIRE](#)].
- [45] ATLAS collaboration, *Performance and Validation of Q-Jets at the ATLAS Detector in pp Collisions at  $\sqrt{s} = 8$  TeV in 2012*, [ATLAS-CONF-2013-087](#) (2013).
- [46] CMS collaboration, *Identifying Hadronically Decaying Vector Bosons Merged into a Single Jet*, [CMS-PAS-JME-13-006](#) (2013).
- [47] M. Cacciari, G.P. Salam and G. Soyez, *The Anti- $k(t)$  jet clustering algorithm*, *JHEP* **04** (2008) 063 [[arXiv:0802.1189](#)] [[INSPIRE](#)].
- [48] S. Catani, Y.L. Dokshitzer, M.H. Seymour and B.R. Webber, *Longitudinally invariant  $K_t$  clustering algorithms for hadron hadron collisions*, *Nucl. Phys. B* **406** (1993) 187 [[INSPIRE](#)].
- [49] S.D. Ellis and D.E. Soper, *Successive combination jet algorithm for hadron collisions*, *Phys. Rev. D* **48** (1993) 3160 [[hep-ph/9305266](#)] [[INSPIRE](#)].
- [50] Y.L. Dokshitzer, G.D. Leder, S. Moretti and B.R. Webber, *Better jet clustering algorithms*, *JHEP* **08** (1997) 001 [[hep-ph/9707323](#)] [[INSPIRE](#)].
- [51] M. Wobisch and T. Wengler, *Hadronization corrections to jet cross-sections in deep inelastic scattering*, [hep-ph/9907280](#) [[INSPIRE](#)].
- [52] M. Wobisch, *Measurement and QCD analysis of jet cross-sections in deep inelastic positron proton collisions at  $\sqrt{s} = 300$ -GeV*, DESY-THESIS-2000-049, Hamburg germany (2014).
- [53] G.P. Salam, *Towards Jetography*, *Eur. Phys. J. C* **67** (2010) 637 [[arXiv:0906.1833](#)] [[INSPIRE](#)].
- [54] T. Sjöstrand, S. Mrenna and P.Z. Skands, *A Brief Introduction to PYTHIA 8.1*, *Comput. Phys. Commun.* **178** (2008) 852 [[arXiv:0710.3820](#)] [[INSPIRE](#)].
- [55] ATLAS collaboration, *ATLAS tunes of PYTHIA 6 and Pythia 8 for MC11*, [ATL-PHYS-PUB-2011-009](#) (2011).
- [56] DELPHES 3 collaboration, J. de Favereau et al., *DELPHES 3, A modular framework for fast simulation of a generic collider experiment*, *JHEP* **02** (2014) 057 [[arXiv:1307.6346](#)] [[INSPIRE](#)].
- [57] M. Cacciari, G.P. Salam and G. Soyez, *FastJet User Manual*, *Eur. Phys. J. C* **72** (2012) 1896 [[arXiv:1111.6097](#)] [[INSPIRE](#)].
- [58] Y.-T. Chien, *Telescoping jets: Probing hadronic event structure with multiple  $R_s$* , *Phys. Rev. D* **90** (2014) 054008 [[arXiv:1304.5240](#)] [[INSPIRE](#)].
- [59] D. Kahawala, D. Krohn and M.D. Schwartz, *Jet Sampling: Improving Event Reconstruction through Multiple Interpretations*, *JHEP* **06** (2013) 006 [[arXiv:1304.2394](#)] [[INSPIRE](#)].
- [60] Y.-T. Chien et al., *Quantifying the power of multiple event interpretations*, [arXiv:1407.2892](#) [[INSPIRE](#)].