

Automated detector simulation and reconstruction parametrization using machine learning

D. Benjamin,¹ S. Chekanov,¹ W. Hopkins,¹ Y. Li,² and J. R. Love¹

¹*High Energy Physics Division, Argonne National Laboratory,
9700 S. Cass Avenue, Argonne, IL 60439, USA*

²*Computational Science Division, Argonne National Laboratory,
9700 S. Cass Avenue, Argonne, IL 60439, USA*

(Dated: January 29, 2020)

Abstract

Rapidly applying the effects of detector resolutions and experimental reconstruction algorithms to physics objects (e.g. electrons, muons, showers of particles) is essential in high energy physics. Currently available tools for the transformation from truth-level physics objects to reconstructed detector-level physics objects involve manually defining resolution functions. These resolution function are typically derived in bins of variables that are correlated with the resolution (e.g. pseudorapidity and transverse momentum). This process is time consuming, requires manual updates when detector conditions change, and can miss important correlations. Machine learning offers a way to automate the process of building these truth-to-reconstructed object transformation functions and can capture complex correlation of these functions for any given set of input variables. Such machine learning algorithms, with sufficient optimization, could have a wide range of applications: improving phenomenological studies by using a better detector representation, speeding up fast simulations based on parametric description of LHC detector responses, and allowing for more efficient production of Geant4 simulation by only simulating events within an interesting part of phase space.

9 I. INTRODUCTION

10 A cornerstone of particle collision experiments is Monte Carlo (MC) simulations of physics
11 processes which generates particles produced by a collision of particles followed by simula-
12 tions of detector responses and object reconstruction. The MC simulation produces objects
13 (jets, electrons, muons, etc) with properties (four momenta, particle types) which entirely
14 depend on the physics processes occurring. These objects commonly referred to as “truth”
15 objects. These objects are altered by interactions with the detector and are reconstructed
16 with experimental algorithms. Such objects, that have undergone a transformation due to
17 detector interaction and reconstruction will be referred to as “reco” objects. With increased
18 complexity of such experiments, such as those at the Large Hadron Collider (LHC), the de-
19 tector simulations become increasing complex and time consuming. Parameterized detector
20 simulations, such as Delphes [1], have been proven to be a vital tools for physics performance
21 and phenomenological studies (i.e. to estimate the sensitivity of an experiment to a new physics
22 model). An approximation of the detector responses and experimental object reconstruction
23 can, however, also be performed by neural networks (NN) trained using the Geant4-based
24 simulations that have gone through an experiment’s reconstruction algorithm. This NN
25 could then computationally rapidly transform truth MC objects (jets and other identified
26 particles) to objects modified by a detector and experimental reconstruction algorithms.

27 The main advantage of detector parametrization based on machine learning (ML) is that
28 a neural network can automatically learn the features introduced by detailed full simulations,
29 therefore, handcrafting parameters to represent resolutions and inefficiencies, as it was done
30 in Delphes and for upgrade studies, is not required. An NN trained using realistic detector
31 simulation could memorize the transformation from truth to the reco quantities without
32 manual binning of quantities by analyzers. Another advantage is that the NN approach can
33 introduce a complex interdependence of variables which is currently difficult to implement in
34 parameterized simulations. Finally, since the underlying libraries used for ML (e.g. Keras,
35 pyTorch, etc) are optimized for a wide range of hardware, an NN-based truth-to-reco trans-
36 formation would be able to run efficiently on heterogeneous hardware resources (resources
37 that use a varied set of processors such as GPUs and CPUs).

38 As a first step towards parameterized detector simulations with ML, it is instructive to
39 investigate how a transformation from the truth to reco objects can be performed, leaving

40 aside the question of introducing objects that are created by misreconstructions or objects
 41 that are lost due to inefficiencies.

42 II. TRADITIONAL PARAMETERIZED FAST SIMULATIONS

43 In abstract terms, a typical variable ξ_i^{reco} that characterizes a particle/jet, such as trans-
 44 verse momentum (p_T) or pseudorapidity (η), can be viewed as a multivariate transform, F ,
 45 of the original variable ξ_1^{truth} at truth-level:

$$\xi_1^{\text{reco}} = F(\xi_1^{\text{truth}}, \xi_2^{\text{truth}}, \xi_3^{\text{truth}}, \dots, \xi_N^{\text{truth}}).$$

46 Generally, such a transform depends on several other variables $\xi_2^{\text{truth}} \dots \xi_N^{\text{truth}}$ characterizing
 47 this (or other) objects at the truth level. For example, the extent at which jet transverse
 48 momentum, p_T is modified by a detector depends on the original truth-level transverse
 49 momentum ($\xi_1^{\text{truth}} = p_T^{\text{truth}}$), pseudorapidity ($\xi_1^{\text{truth}} = \eta^{\text{truth}}$), flavor of jets and other effects
 50 that can be inferred from truth quantities. Similarly, if particular detector modules in the
 51 azimuthal angle (ϕ) are not active, this would introduce an additional dependence of this
 52 transform on ϕ .

53 Typical parameterized simulations ignore the full range of correlations between the vari-
 54 ables. In most cases, the above transform is reduced to a single variable, or two (as in the
 55 case of Delphes simulations where the energy resolution of clusters depends on the original
 56 energies of particles and their positions in η). In order to take into account correlations be-
 57 tween multiple parameters characterizing transformations to reconstruction objects a grid
 58 in the hypercube with the dimension N_b^N , where N_b is the number of histogram bins for the
 59 distributions $(\xi^{\text{reco}} - \xi^{\text{truth}})/\xi^{\text{truth}}$ representing “resolution” must be created. This method-
 60 ology results in a large number of histograms when there are many correlated variables that
 61 affect the resolution.

62 It should be pointed out that the calculation speed for parameterized simulations of one
 63 variable that depends on N other variables at the truth level depends as N_b^N since each
 64 object at the truth level should be placed inside the grid defined by N_b bins. Therefore,
 65 complex parameterisations of resolutions and efficiencies for $N > 2$ becomes CPU intensive.

66 III. JET TRUTH-TO-RECO TRANSFORMATION WITH ML

67 To test the viability of using ML to transform truth objects to reco objects, we studied
68 the truth-to-reco transformation for jets. Jet truth-level quantities, such as jet η , p_T , ϕ
69 and jet mass (m) are used as training inputs to an NN while the output is an array of
70 nodes that represent the binned probability density function (PDF) of the resolution for a
71 single variable (such as jet p_T). Additional input variables could be any variable that can
72 influence the resolution of a jet, such as jet flavor at the truth level, jet radius, etc. Figure 1
73 shows a schematic representation of the NN architecture for modelling detector response for
74 a single output variable. The aim is to have the NN learn the shape of the resolution PDF,
75 for example for the p_T , depending on other input variables such as the η of the object. A
76 binned output (multi-categorization) was used to so that the precision of resolution PDF
77 modelling can be selected.

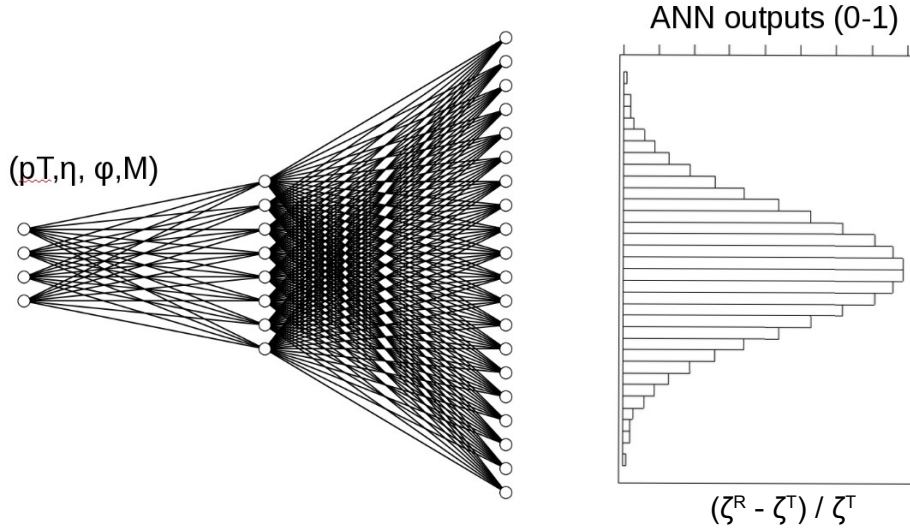


FIG. 1. A schematic representation of the NN architecture for modelling the detector response and affect of reconstruction algorithms on truth-level input variables. The output of this NN is a PDF for the resolution of single variable, e.g. $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$.

78 IV. MONTE CARLO SIMULATED EVENT SAMPLES

79 Monte Carlo events used for this analysis were produced using the Madgraph generator [2].
80 The simulated processes are a combination of equal parts top pair production ($t\bar{t}$ +jets)
81 and photons produced in association with jets (γ +jets), which give a high rate of jets
82 in different environments. Hadronic jets were reconstructed with the FASTJET package [3]
83 using the anti- k_t algorithm [4] with a distance parameter of 0.4. The detector simulation was
84 performed with the Delphes package with a detector geometry which is similar to the ATLAS
85 geometry. The event samples used in this paper, before and after the fast simulation, are
86 available from the HepSim database [5]. In this paper only the transformation from truth
87 jets (which have truth particle constituents) to reconstructed jets (which calorimeter cell
88 constituent) and only for p_T was performed, however the methodology should be object and
89 parameter agnostic. Truth jets which were matched to a reconstructed Delphes jet are used.
90 For the matching criteria the reconstructed jet that has the smallest $\Delta R = \sqrt{\Delta\phi^2 + \Delta\eta^2}$,
91 where $\Delta\phi = \phi^{\text{truth}} - \phi^{\text{reco}}$ and $\Delta\eta = \eta^{\text{truth}} - \eta^{\text{reco}}$, with respect to the truth jet is chosen.
92 If this minimum ΔR is greater than 0.2, the truth jet is discarded. No other requirements
93 are made on truth and reconstructed Delphes jets other than the $p_T > 15$ GeV requirement
94 made by Delphes. Only matched jets are used for this study since the aim of the study is to
95 test whether an NN can changes in detector resolution as a function of kinematic properties
96 of the jet (e.g. p_T , η , ϕ , m). The final number of training jets used is two million while
97 500,000 jets were used as a testing sample. The distributions of quantities used as the input
98 for the NN, p_T , η , ϕ , m , are shown in Figure 2.

99 To facilitate gradient descent in all direction of the input variables, the input variables are
100 scaled to be in the range [0,1]. This avoids the p_T and the mass from having a disproportional
101 affect on the training of the NNs. The output variable, $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$, is also scaled
102 to have values between 0 and 1. Only objects that are within the 1st and 99th percentile of
103 the $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$ distribution are considered in this study since objects outside this
104 range are typically not used in physics analyses.

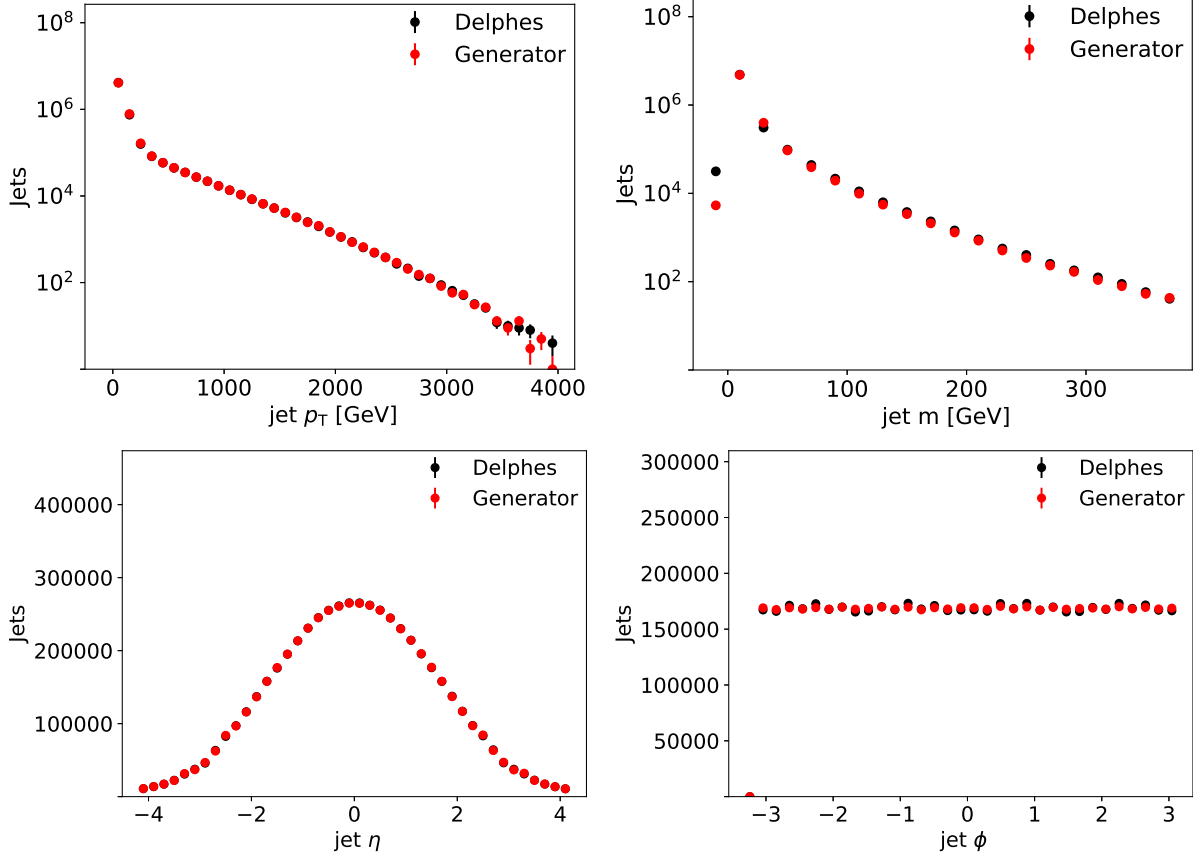


FIG. 2. Distributions for input variables for truth (red) and reco quantities (black).

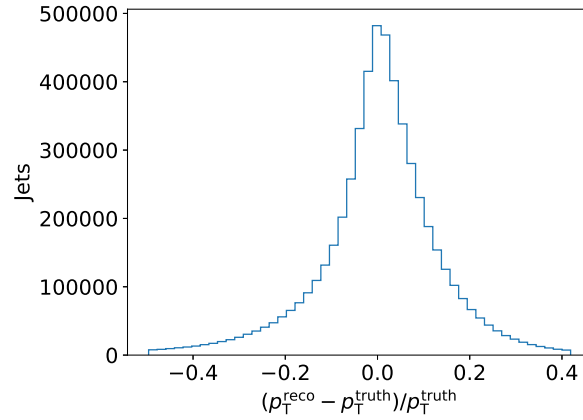


FIG. 3. Relative differences between truth and reco p_T .

105 V. NEURAL NETWORK STRUCTURES

106 An NN is trained with four input parameters, the scaled p_T , η , ϕ , and m , and consist
 107 of five layers with 100 nodes each and with each node having a rectifier linear unit (ReLU)

activation function. The output layer has 400 nodes with a softmax activation function. Finally, the NN is trained over 1000 epochs with batch size 1000 using the Adam [6] optimizer with a learning rate of 10^{-4} . The NN is implemented using Keras [7] with a TensorFlow [8] backend.

VI. RESULTS

After the NN has been trained to learn the PDF of $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$, the resulting learned PDF is compared to the Delphes PDF using the testing sample in Figure 4a. Good agreement is observed between the Delphes and NN PDFs, showing that the NN has learned the bulk distribution.

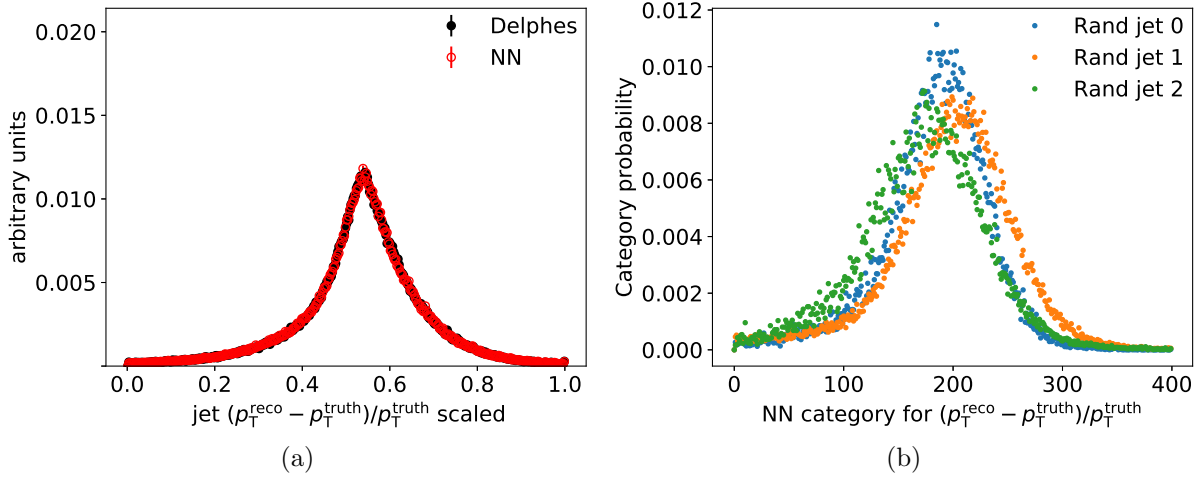


FIG. 4. NN-generated jet $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$ compared to reco Delphes jet $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$ (a). NN-generated jet PDFs for three randomly selected truth jets (b).

The NN predicts a PDF for each jet based on its input parameters (i.e. p_T , ϕ , η , and m). The PDFs for a set of randomly selected jets are shown in Figure 4b. These PDFs are then randomly sampled to produce an NN jet that mimics the reco jet. A comparison of the NN-generated and Delphes jet p_T distribution for the testing sample is shown in Figure 5. The NN reproduces the jet p_T distribution of Delphes within 5% for reconstructed jets with $p_T > 20$ GeV.

To test whether the NN learned correlations between input parameters and the p_T resolution, the jets were divided into central ($|\eta| < 3.2$) and forward ($|\eta| > 3.2$) jets. The p_T resolution is then compared between the two regions for both the Delphes jets as well as the

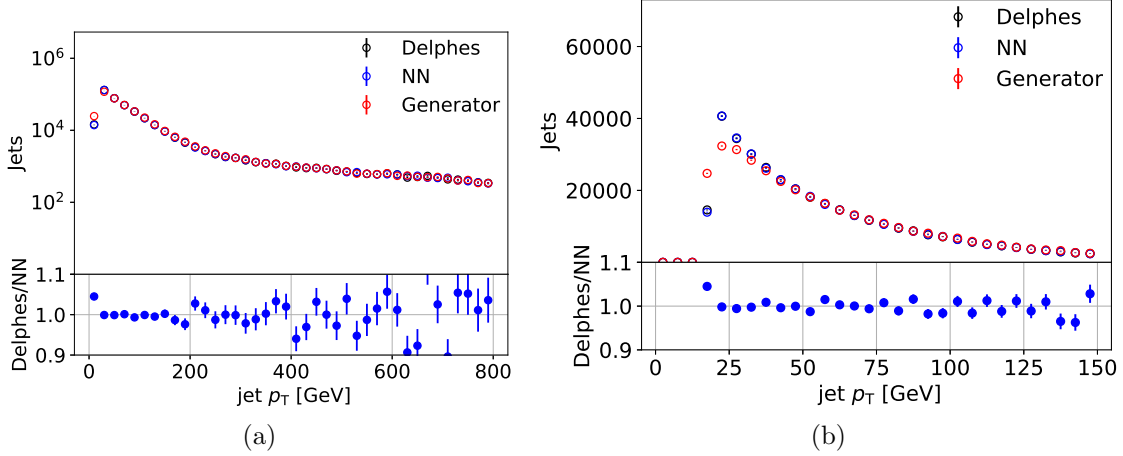


FIG. 5. Delphes and NN-generated jet p_T disitributions for a wide (a) and narrow (b) p_T range.

126 NN-generated jets. These two regions in the detector simulation have different calorimeter
 127 resolutions which results in different jet p_T resolutions and thus the p_T resolution is corre-
 128 lated with $|\eta|$. The resulting resolutions for both regions are shown in Figure 6 using the
 129 training sample. The training sample was chosen for this comparison because forward jets
 130 make up a small subsample of all jets, as can be seen in Figure 2. Several batch-size and
 131 number-of-epoch combinations were used in an attempt to optimize the sensitivity to a small
 132 subsample (the forward jets) of the training sample. The number of backpropagations (N_{bp})
 133 were held constant by keeping the ratio of the number of epochs (N_e) and batch size (N_b)
 134 constant since $N_{bp} = \frac{N_t}{N_b} N_e$ where N_t is the number of training jets. Batch size and number
 135 of epochs of 5, 10, 20, 100, 200, 1000 were tested resulting in similar performance of the NN
 136 in both the central and forward region.

137 To test whether the NN learned the jet resolution for a large range of p_T s, the mean
 138 and standard deviation of the resolution (shown, inclusively, in Figure 4) was plotted as a
 139 function of p_T in Figure 7. The mean of the resolution for the NN is systematically higher
 140 than the resolution for Delphes but this effect is small when considering the width of the
 141 resolution. The standard deviation of the resolution, however, are the same for the NN and
 142 Delphes across the p_T range.

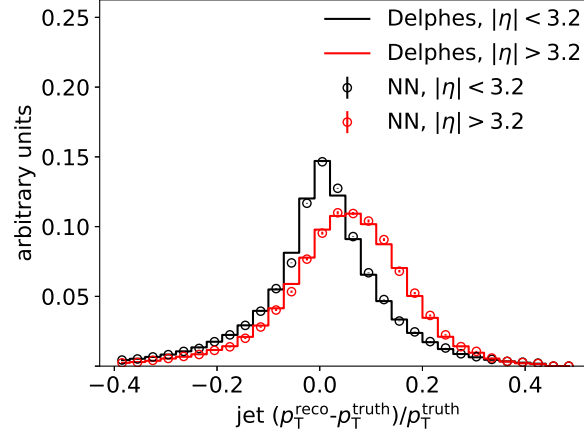


FIG. 6. Jet p_T resolution for the training sample for both the central and forward region.

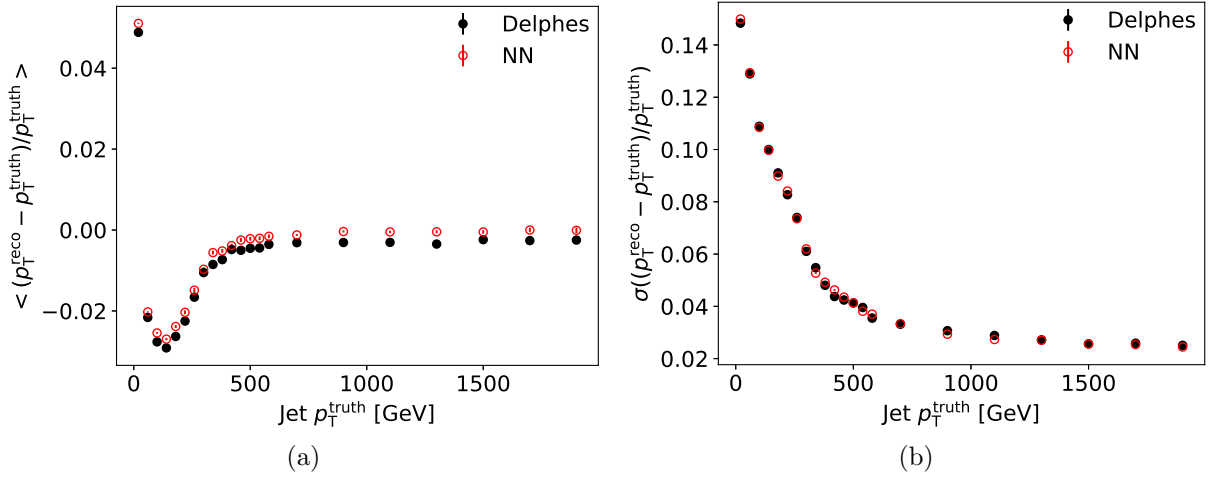


FIG. 7. The mean (a) and standard deviation of jet p_T resolution for Delphes and NN-generated as a function truth jet p_T .

VII. CONCLUSION

We have shown that a truth-level quantity can be transformed to a reconstruction-level quantity using a multi-categorizing NN. The NN learned the truth-to-reconstruction transformation without requiring manual binning to capture the differences in resolutions of particular subsamples (central and forward jets). Additional improvements could be made by including more information about the objects (e.g. whether a b -quark is present in a jet, kinematic information from other objects in the event) making this method more robust. This method should be easily extendable to additional reconstructed quantities and could

151 be used to model the ATLAS and CMS detector. The method described in this paper thus
152 allows for automated detector parametrization which can facilitate phenomenological studies,
153 efficient truth event selection, and upgrade studies.

154 ACKNOWLEDGMENTS

155 The submitted manuscript has been created by UChicago Argonne, LLC, Operator of
156 Argonne National Laboratory (“Argonne”). Argonne, a U.S. Department of Energy Office
157 of Science laboratory, is operated under Contract No. DE-AC02-06CH11357. The U.S. Gov-
158 ernment retains for itself, and others acting on its behalf, a paid-up nonexclusive, irrevocable
159 worldwide license in said article to reproduce, prepare derivative works, distribute copies to
160 the public, and perform publicly and display publicly, by or on behalf of the Government.
161 The Department of Energy will provide public access to these results of federally sponsored
162 research in accordance with the DOE Public Access Plan. [http://energy.gov/downloads/](http://energy.gov/downloads/doe-public-access-plan)
163 [doe-public-access-plan](http://energy.gov/downloads/doe-public-access-plan). Argonne National Laboratory’s work was funded by the U.S.
164 Department of Energy, Office of High Energy Physics under contract DE-AC02-06CH11357.

-
- 165 [1] J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaître, A. Mertens, and M. Sel-
166 vaggi. DELPHES 3, A modular framework for fast simulation of a generic collider experiment.
167 *JHEP*, 02:057, 2014.
- 168 [2] J. Alwall, R. Frederix, S. Frixione, V. Hirschi, F. Maltoni, et al. The automated computation
169 of tree-level and next-to-leading order differential cross sections, and their matching to parton
170 shower simulations. *JHEP*, 07:079, 2014.
- 171 [3] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. FastJet User Manual. *Eur. Phys. J. C*,
172 72:1896, 2012.
- 173 [4] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. The anti- k_t jet clustering algorithm.
174 *JHEP*, 04:063, 2008.
- 175 [5] S.V. Chekanov. HepSim: a repository with predictions for high-energy physics experiments.
176 *Advances in High Energy Physics*, 2015:136093, 2015. Available as [http://atlaswww.hep.](http://atlaswww.hep.anl.gov/hepsim/)
177 [anl.gov/hepsim/](http://atlaswww.hep.anl.gov/hepsim/).

- 178 [6] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *arXiv*
179 *e-prints*, page arXiv:1412.6980, Dec 2014.
- 180 [7] François Chollet et al. Keras. <https://keras.io>, 2015.
- 181 [8] Martín Abadi et al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015.
182 Software available from tensorflow.org.