

# Replication of detector simulations using supervised machine learning

D. Benjamin,<sup>1</sup> S. Chekanov,<sup>1</sup> W. Hopkins,<sup>1</sup> Y. Li,<sup>2</sup> and J. R. Love<sup>1</sup>

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

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

(Dated: October 31, 2019)

## Abstract

Accurately and computationally rapidly modeling stochastic detector response for complex LHC experiments involving many particles from multiple interaction points requires the development of novel techniques. A study aimed at finding a fast transformation from truth-level physics objects to reconstructed detector-level physics objects is presented. This study used Delphes fast simulation based on an LHC-like detector geometry for inputs for a multi-categorizing machine learning (ML) algorithms. This ML transfer 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 followed by simulations of detector responses. With increased complexity of such  
12 experiments, such as those at the Large Hadron Collider (LHC), the detector simulations  
13 become increasing complex and time consuming. For example, the time required to simulate  
14 Geant4 [1] hits and to reconstruct from such hits physics objects (electrons, muons, taus, jets)  
15 needs a factor 100-1000 more CPU time than the creation of typical Monte Carlo events  
16 that represent physics processes according to theoretical models (“truth level” MC event  
17 generation). A possible method to speed up simulations of detector responses is to apply  
18 neural networks (NN) trained using the Geant4-based simulations, and use such supervised  
19 NN for transforming truth-level MC objects (jets and other identified particles) to objects  
20 modified by detectors (“detector-level”).

21 A typical simulation of detector responses stochastically modifies positions and energies  
22 of particles and jets created by MC generators at the truth-level. Another important compo-  
23 nent of such simulations is to introduce additional particles due to misreconstructed energy  
24 deposits in active detector volumes (examples include misreconstructed electrons or photons  
25 which are, in fact, hadronic jets). The latter effects represent a significant complication for  
26 the so-called “fast” or “parameterized” detector simulations, such as Delphes [2]. Never-  
27 theless, parameterized detector simulations have been proven to be a vital tools for physics  
28 performance and phenomenological studies.

29 The main advantage of detector parameterization based on machine learning is that a  
30 neural networks can automatically learn the features introduced by detailed full simulations,  
31 therefore, handcrafting parameters to represent resolutions and inefficiencies, as it was done  
32 in Delphes and for upgrade studies, is not required. A neural network trained using realistic  
33 detector simulation should memorize the transformation from truth-level to the detector-  
34 level quantities without manual binning of quantities by analyzers. Another advantage  
35 is that the NN approach can introduce a complex interdependence of variables which is  
36 currently difficult to implements in parameterized simulations.

37 As a first step towards parameterized detector simulations using machine learning tech-  
38 niques, it is instructive to investigate how a transformation from the truth-level MC to  
39 detector-level objects can be performed, leaving aside the question of introducing objects

40 that are created by misreconstructions.

## 41 II. TRADITIONAL PARAMETERIZED FAST SIMULATIONS

42 In abstract terms, a typical variable  $f_i$  that characterizes a particle/jet, such as transverse  
 43 momentum ( $p_T$ ), pseudorapidity ( $\eta$ ), can be viewed as a multivariate transform  $F$  of the  
 44 original variable  $\xi_1^T$  at truth-level:

$$\xi_1 = F(\xi_1^T, \xi_2^T, \xi_3^T, \dots, \xi_N^T).$$

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

51 Typical parameterized simulations ignore the full range of correlations between the vari-  
 52 ables. In most cases, the above transform is reduced to a single variable, or two (as in the  
 53 case of Delphes simulations where energy resolution of clusters depend on the original ener-  
 54 gies of particles and their positions in  $\eta$ ). In order to take into account correlations between  
 55 multiple parameters characterizing transformations to the detector level, the following steps  
 56 have to be undertaken:

- 57 • create a grid in the hypercube with the dimension  $N_b^N$ , where  $N_b$  is the number of  
 58 histogram bins for the distributions  $f_1 - f_i^N$  representing “resolution” smearing. This  
 59 can be done numerically, using frequencies, or using analytically using “resolution  
 60 functions”.
- 61 • calculate “efficiencies” that model losses of particles/jets for each variable.

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 efficiency’s for  $N > 2$  becomes CPU intensive.

### 66 III. MACHINE LEARNING APPROACH FOR FAST SIMULATION

67 Unlike the traditional approach for fast simulation using parameterized density functions  
 68 for resolution variables and probability values for efficiency, a neural network approach offers  
 69 an opportunity to formulate this problem in terms of NN nodes and their connections that  
 70 scale as  $N_b^N \cdot N$ , which can speed up the fast simulations and, at the same time, can be used  
 71 for learning more complex full simulations in an automated way.

72 In the case of objects, such as jets, a typical truth-level input are jet transverse momen-  
 73 tum,  $\eta$ ,  $\phi$  and jet mass  $m$ , while the output is an array of output nodes that represent the  
 74 binned probability density function (PDF) of the resolution for a single variable (such as jet  
 75  $p_T$ ). Additional input variables can be jet flavor at the truth level, jet radius etc., i.e. any  
 76 variable that can influence the output of such neural network. Figure 1 shows a schematic  
 77 representation of the NN architecture for modelling detector response for a single output  
 78 variable. The inputs of the NN are variables that can affect the object resolutions while the  
 79 hidden layers are used to capture correlations between these variables and the the output  
 80 which consist of a binned PDF. The aim is to have the NN learn the shape of this PDF, for  
 81 example for the  $p_T$ , depending on other variables such as the  $\eta$  of the object.

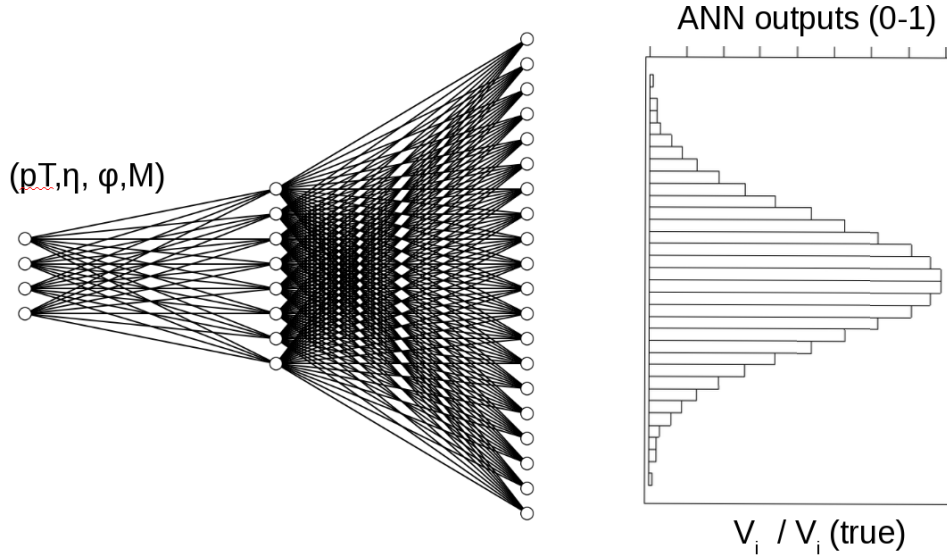


FIG. 1. A schematic representation of the NN architecture for modelling the detector response to truth-level input variables. The output of this NN is PDF for a single variable, e.g.  $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$ , which modified by the detector.

## 82 IV. MONTE CARLO SIMULATED EVENT SAMPLES

83 Monte Carlo events used for this analysis were generated using the Madgraph genera-  
 84 tor [3]. The simulated processes were a combination of equal parts  $t\bar{t}$ +jets and  $\gamma$ +jets, which  
 85 give a high rate of jets in different environments. Hadronic jets were reconstructed with the  
 86 FASTJET package [4] using the anti- $k_T$  algorithm [5] with a distance parameter of 0.4. The  
 87 detector simulation was performed with the Delphes package [2] with an ATLAS-like detec-  
 88 tor geometry. The event samples used in this paper, before and after the fast simulation, are  
 89 available from the HepSim database [6]. In this paper only the transformation from truth-  
 90 level jets to detector-level jets and only for  $p_T$  was performed, however the methodology  
 91 should be object and parameter agnostic. Only truth jets which have been matched to a  
 92 reconstructed Delphes jet are used. For the matching criteria the reconstructed jet that has  
 93 the smallest  $\Delta R = \sqrt{\Delta\phi^2 + \Delta\eta^2}$ , where  $\Delta\phi = \phi^{\text{truth}} - \phi^{\text{reco}}$  and  $\Delta\eta = \eta^{\text{truth}} - \eta^{\text{reco}}$ , with  
 94 respect to the truth jet is chosen. If this minimum  $\Delta R$  is greater than 0.2, the truth-level jet  
 95 is discarded. No other requirements are made on truth on reconstructed Delphes jets other  
 96 than the  $p_T > 15$  GeV requirement made by Delphes. Only matched jets are used for this  
 97 study since the aim of the study is to test whether an NN can changes in detector resolution  
 98 as a function of kinematic properties of the jet (e.g.  $p_T$ ,  $\eta$ ,  $\phi$ ,  $m$ ). The final number of  
 99 training jets used is two million while 500,000 jets were used as a testing sample.

100 The distributions of quantities used as the input for the NN,  $p_T$ ,  $\eta$ ,  $\phi$ ,  $m$ , are shown in  
 101 Figure 2.

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

## 108 V. NEURAL NETWORK STRUCTURES

109 An NN is trained with four input parameters, the scaled  $p_T$ ,  $\eta$ ,  $\phi$ , and  $m$ , and consist of five  
 110 layers with 100 nodes each and with each node having a rectifier linear unit (ReLU) activation

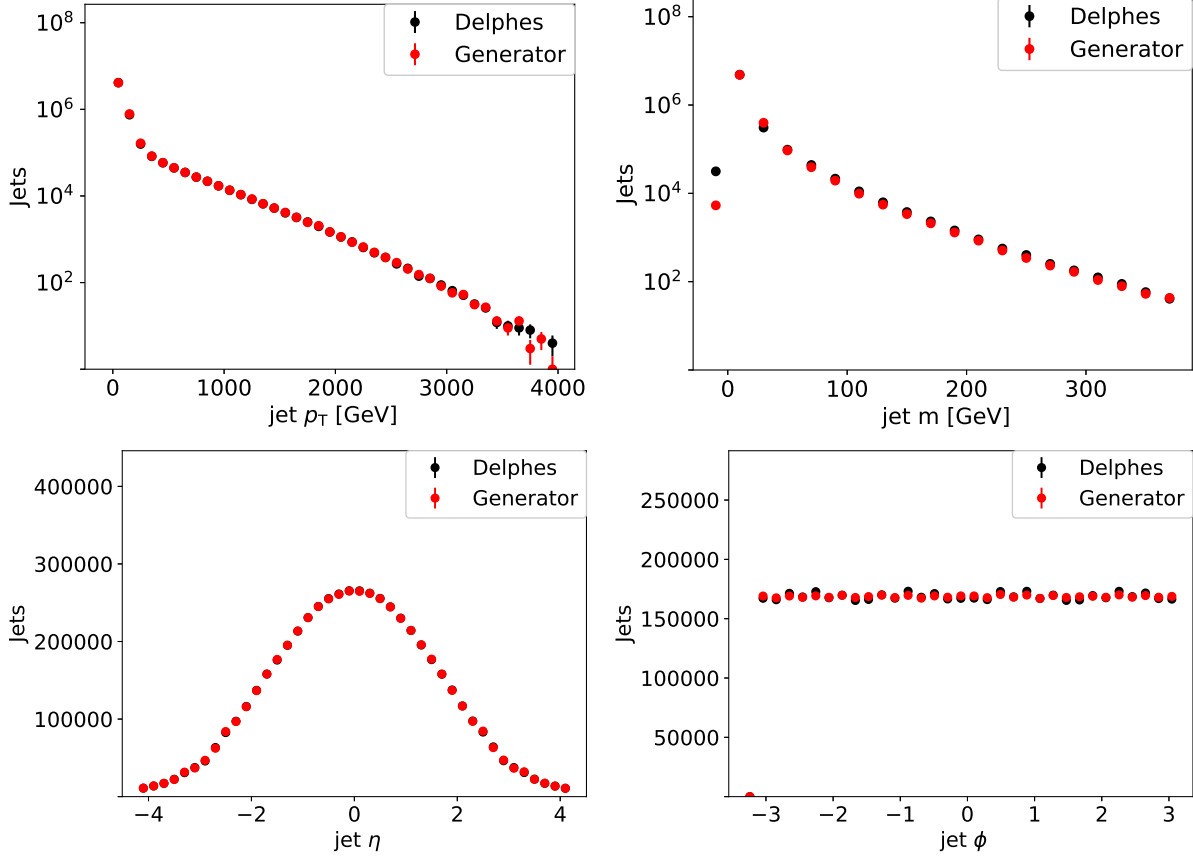


FIG. 2. Input variable shapes for truth-level (red) and detector-level quantities (black).

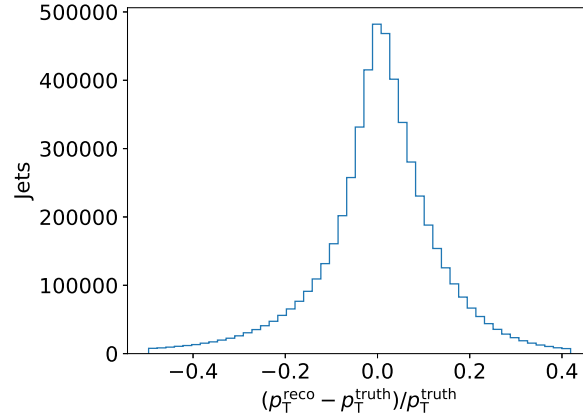


FIG. 3. Relative differences between truth-level and detector-level  $p_T$ .

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 [7] optimizer with a learning rate of  $1 \times 10^{-4}$ . The NN is implemented using Keras [8] with a TensorFlow [9]

114 backend.

## 115 VI. RESULTS

116 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  
 117 learned PDF is compared to the Delphes PDF using the testing sample in Fig. 4. Good  
 118 agreement is observed between the Delphes and NN PDFs, showing that the NN has learned  
 119 the bulk distribution.

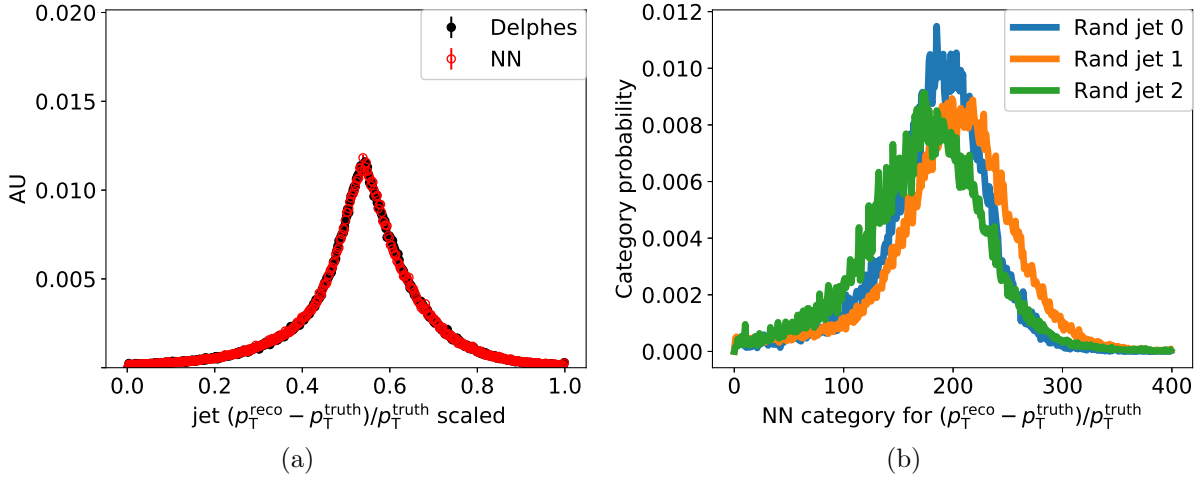


FIG. 4. NN-generated jet  $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$  compared to detector-level jet  $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$   
 (a). NN-generated jet PDFs for three randomly selected truth-level jets (b).

120 The NN predicts a PDF for each jet based on its input parameters (i.e.  $p_T$ ,  $\phi$ ,  $\eta$ , and  
 121  $m$ ). The PDFs for a set of randomly selected jets are shown in Fig. 4. These PDFs are then  
 122 randomly sampled to produce an NN jet that mimics the detector-level jet. A comparison of  
 123 the NN-generated and Delphes jet  $p_T$  distribution for the testing sample is shown in Fig. 5.  
 124 The NN reproduces the jet  $p_T$  distribution of Delphes within 5% for reconstructed jets with  
 125  $p_T > 20$  GeV.

126 To test whether the NN learned correlations between input parameters and the  $p_T$  res-  
 127 olution, the jets were divided into central ( $|\eta| < 3.2$ ) and forward ( $|\eta| > 3.2$ ) jets. The  $p_T$   
 128 resolution is then compared between the two regions for both the Delphes jets as well as the  
 129 NN-generated jets. These two regions in the detector simulation have different calorimeter  
 130 resolutions which results in different jet  $p_T$  resolutions and thus the  $p_T$  resolution is corre-  
 131 lated with  $|\eta|$ . The resulting resolutions for both regions and jets are shown in Fig. 6 using

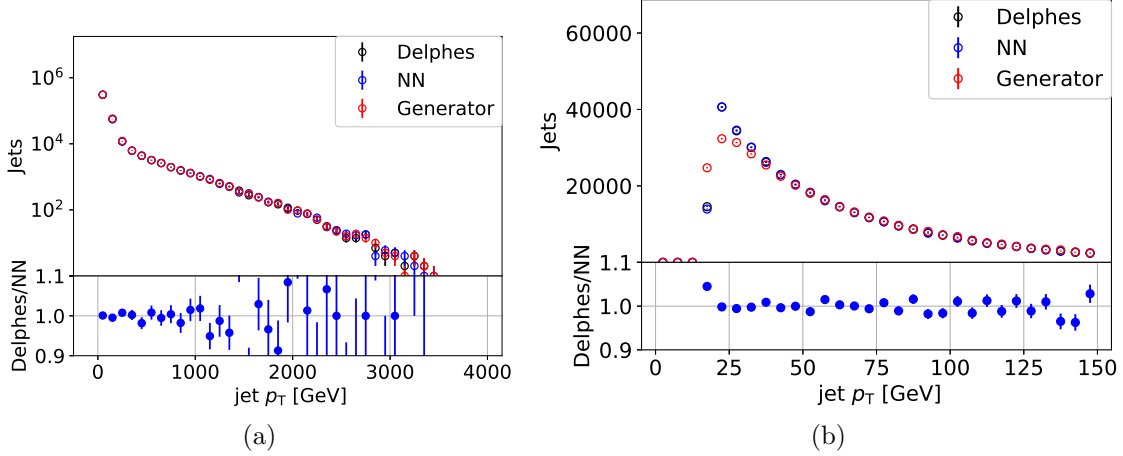


FIG. 5. Delphes and NN-generated jet  $p_T$  distributions for a wide (a) and narrow (b)  $p_T$  range.

the training sample. The training sample was chosen for this comparison because forward jets make up a small subsample of all jets, as can be seen in Fig. 2. Several batch-size and number-of-epoch combinations were used in an attempt to optimize the sensitivity to a small subsample (the forward jets) of the training sample. The number of backpropagations were held constant by keeping the ratio of the number of epochs ( $N_e$ ) and batch size ( $N_b$ ) constant since the number of backpropagations ( $N_{bp}$ ) is given by  $N_{bp} = \frac{N_t}{N_b} N_e$  where  $N_t$  is the number of training jets. Batch size and number of epochs of 5, 10, 20, 100, 200, 1000 were tested resulting in similar performance of the NN in both the central and forward region.

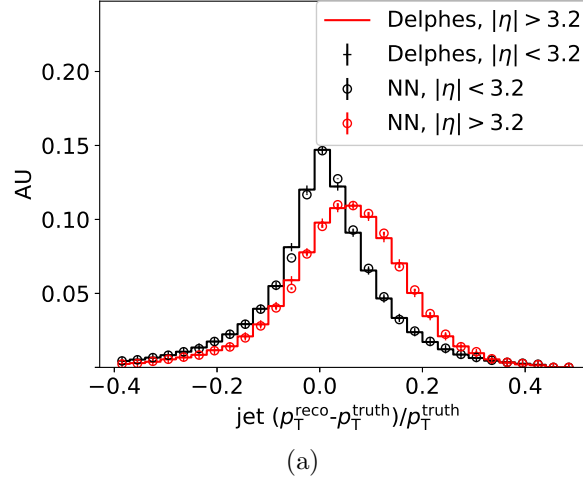


FIG. 6. Jet  $p_T$  resolution for the training sample for both the central and forward region.

To test whether what the NN learned is sample independent we applied the NN to a  $t\bar{t} + W$  sample which contained 900,000 jets. This sample is expected to have lower  $p_T$



jets in general but a higher fraction of high  $p_T$  jets that contain  $b$ -quarks ( $b$ -jets). The  $p_T$  distributions over two  $p_T$  ranges are shown in Fig. 7 while the comparison between the central and forward jet  $p_T$  resolutions can be seen in Fig. 8. The agreement in Fig. 7 is good for jets with  $p_T < 200$  GeV but a mismodeling trend appears at  $p_T$ s greater than 200 GeV. We expect this to be due high  $p_T$   $b$ -jets which are not as prevalent in the training set where the high  $p_T$  jets mainly come from the  $\gamma$ +jets sample which is produced to have a flattened  $p_T$  spectrum. For future studies and refinements one could add truth  $b$ -quark information to help improve the modelling of these types of jets.

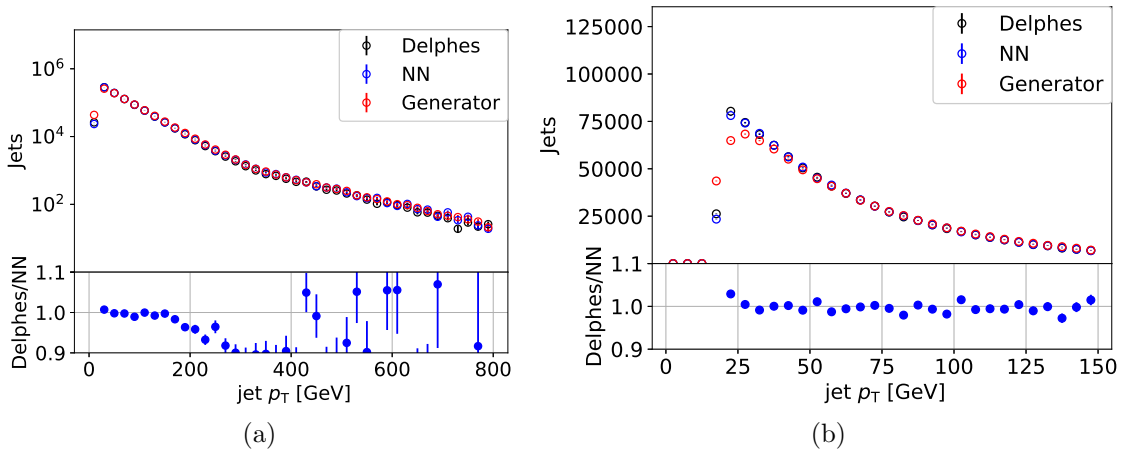


FIG. 7. Delphes and NN-generated jet  $p_T$  distributions for a wide (a) and narrow (b)  $p_T$  range for jets originating from  $t\bar{t} + W$  production.

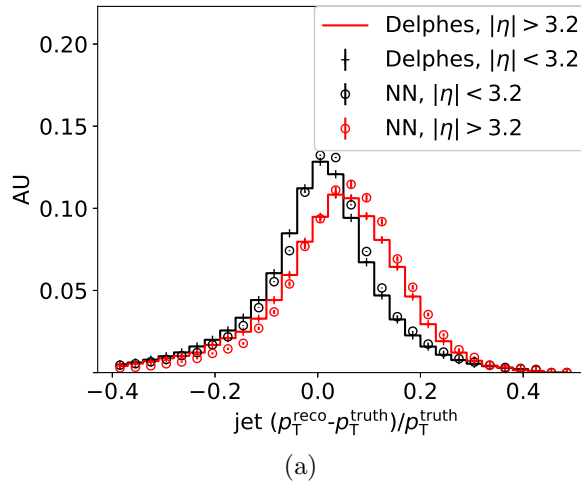
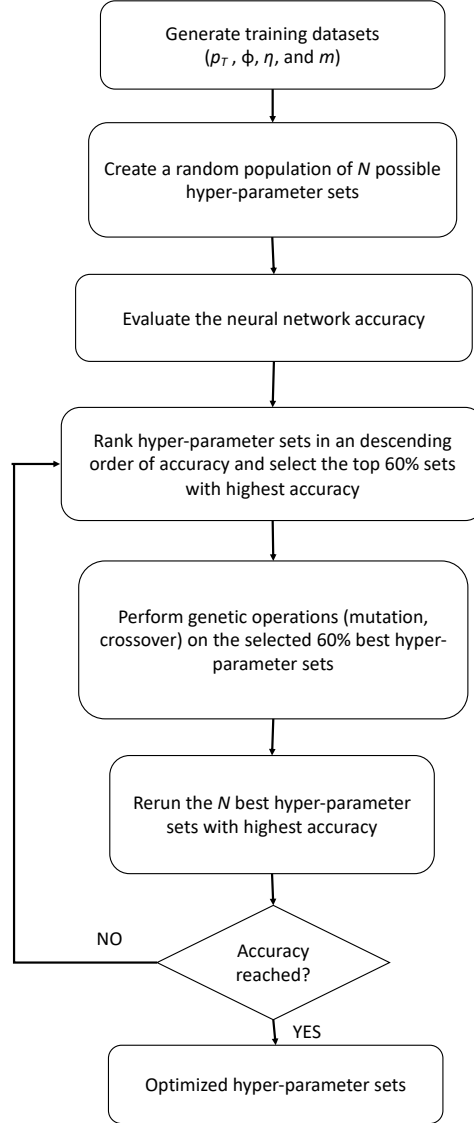


FIG. 8. Jet  $p_T$  resolution for the training sample for both the central and forward region for jets originating from  $t\bar{t} + W$  production.

## 150 VII. HYPERPARAMETER OPTIMIZATION

151 In order to acquire the optimal performance with an NN with a minimal amount of  
152 parameters, a genetic algorithm (GA) was used to optimize the NN hyperparameters. Hyper  
153 parametrization, that is, the determination of the optimal parameters associated with a  
154 complex NN (number of layers, number of neurons on each layer, choice of the learning  
155 rate, mini-batch size), is a challenging optimization problem in a high dimensionality space,  
156 especially for nonlinear NNs. Ultimately, the range of hyperparameter values that must be  
157 explored grows dramatically as the complexity of the deep NN increases. To address this  
158 challenge, we have used a GA, which is an evolutionary algorithm that mimics the process  
159 of natural selection. Previously successful applications of the GA in similar contexts include  
160 the determination of the parameters of a complex force field [10, 11]. A flowchart of the  
161 ML/GA optimization protocol is depicted in Fig. 9. The optimization starts with a set of  
162 parent parameters, which is defined as the population. For hyperparameter optimization  
163 of the NN, the parents sets are number of layers, number of neurons on each layer, choice  
164 of the learning rate or mini-batch size. Some of the individuals in this population present  
165 a better fit, which in the context of NNs means a higher value for the accuracy of the  
166 network. Fit parents survive and are allowed to mate, which is accomplished by crossing  
167 patterns with other fit individuals. During crossover, random mutations in the genes are  
168 also allowed, to a certain degree, to avoid a stagnant gene pool and a better sampling of  
169 the parameters space. The offspring individuals form the next generation of parents and  
170 this process continues until some predefined criteria are met. Sets of parameters are then  
171 ranked in ascending order. After the ranking, a nonlinear roulette wheel selection [12] was  
172 performed to select the best 60% members, that is, the ones with lowest values of accuracy,  
173 which were then subjected to genetic operations: mutation and crossover with a 3% crossover  
174 rate. These mutations introduce sufficient diversity into the population, and the nonlinear  
175 selection scheme helps to avoid premature convergence of the ML/GA run. After the genetic  
176 operations, both the parent and offspring sets of parameters are ranked by their value of  
177 accuracy. The best hyperparameter sets are then chosen to constitute the next generation.  
178 Such an optimization routine ensures that only satisfactory hyperparameter sets survive  
179 after each generation; upon repeating this workflow for sufficient generations and sampling  
180 viable regions in the parameter space, we performed three separate ML/GA runs starting

181 with different random populations. From each of the converged ML/GA run, we chose the  
 182 final hyperparameter set corresponding the highest value of accuracy.



(a)

FIG. 9. Flow chart of the hyperparameter optimization used to optimize the resolution NN.

## 183 VIII. CONCLUSION

184 We have shown that a truth-level quantity can be transformed to a reconstruction-level  
 185 quantity using a multi-categorizing NN. The NN learned the changes in resolutions of dif-  
 186 ferent regions of the detector based on the NN inputs during training. The NN learned

the truth-to-reconstruction transformation without requiring manual binning to capture the differences in resolutions of particular subsamples. This method should be easily extendable to additional reconstructed quantities and could be used to model the ATLAS and CMS detector. The method described in this paper thus allows for automated detector parameterization which can facilitate phenomenological, efficient truth event selection, and upgrade studies. Additional improvements could be made by including more information about the objects (e.g. whether a  $b$ -quark is present in a jet) could make this method more robust.

## ACKNOWLEDGMENTS

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

- 
- [1] S. Agostinelli et al. GEANT4: A Simulation toolkit. *Nucl. Instrum. Meth. A*, 506:250, 2003.
  - [2] J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaître, A. Mertens, and M. Selvaggi. DELPHES 3, A modular framework for fast simulation of a generic collider experiment. *JHEP*, 02:057, 2014.
  - [3] J. Alwall, R. Frederix, S. Frixione, V. Hirschi, F. Maltoni, et al. The automated computation of tree-level and next-to-leading order differential cross sections, and their matching to parton shower simulations. *JHEP*, 07:079, 2014.
  - [4] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. FastJet User Manual. *Eur. Phys. J. C*, 72:1896, 2012.

- [5] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. The anti- $k_t$  jet clustering algorithm. *JHEP*, 04:063, 2008.
- [6] S.V. Chekanov. HepSim: a repository with predictions for high-energy physics experiments. *Advances in High Energy Physics*, 2015:136093, 2015. Available as <http://atlaswww.hep.anl.gov/hepsim/>.
- [7] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *arXiv e-prints*, page arXiv:1412.6980, Dec 2014.
- [8] François Chollet et al. Keras. <https://keras.io>, 2015.
- [9] Martín Abadi et al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [10] Ying Li, Hui Li, Frank C. Pickard, Badri Narayanan, Fatih G. Sen, Maria K. Y. Chan, Subramanian K. R. S. Sankaranarayanan, Bernard R. Brooks, and Benoît Roux. Machine learning force field parameters from ab initio data. *Journal of Chemical Theory and Computation*, 13(9):4492–4503, 2017. PMID: 28800233.
- [11] Md Mahbubul Islam, Grigory Kolesov, Toon Verstraelen, Efthimios Kaxiras, and Adri C. T. van Duin. ereaxff: A pseudoclassical treatment of explicit electrons within reactive force field simulations. *Journal of Chemical Theory and Computation*, 12(8):3463–3472, 2016. PMID: 27399177.
- [12] Adam Lipowski and Dorota Lipowska. Roulette-wheel selection via stochastic acceptance. *CoRR*, abs/1109.3627, 2011.