

Replication of detector simulations using supervised 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: October 29, 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 (η), can be viewed as a multivariate transform F of the
 44 original variable ξ_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 η , flavor of jets and other effects that can be inferred
 49 from the truth level. Similarly if particular detector modules in the azimuthal angle (ϕ) are
 50 not active, this would introduce an additional dependence of this transform on ϕ .

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 η). 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, η , ϕ 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 η of the object.

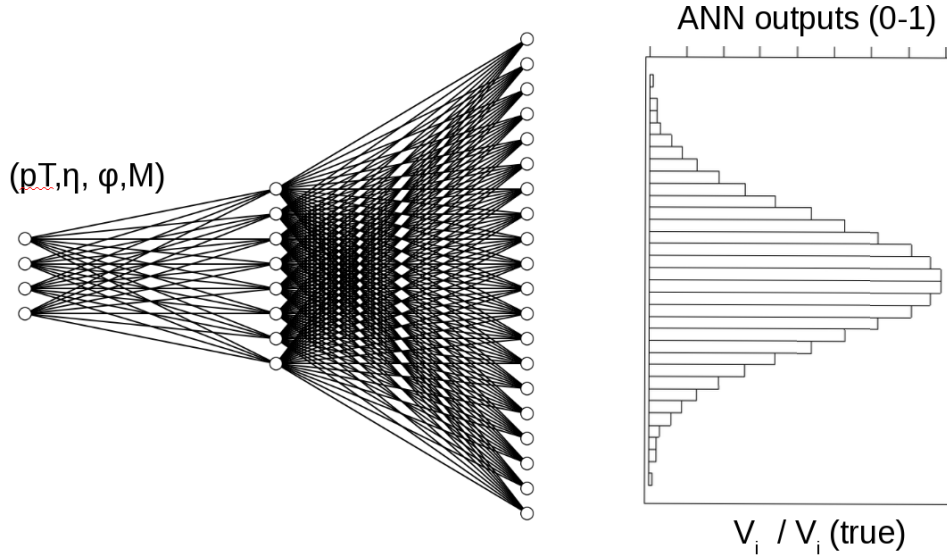


FIG. 1. A schematic representation of the NN architecture for modelling the detector response to truth-level input variables. The output of this NN is PDF for a single variable, e.g. $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$, which modified by the detector.

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 γ +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 ΔR 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 , η , ϕ , 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 , η , ϕ , 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 1st and 99th
 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 , η , ϕ , and m , and consist of five
 110 layers with 100 nodes each and with each node having a rectifier linear unit (ReLU) activation

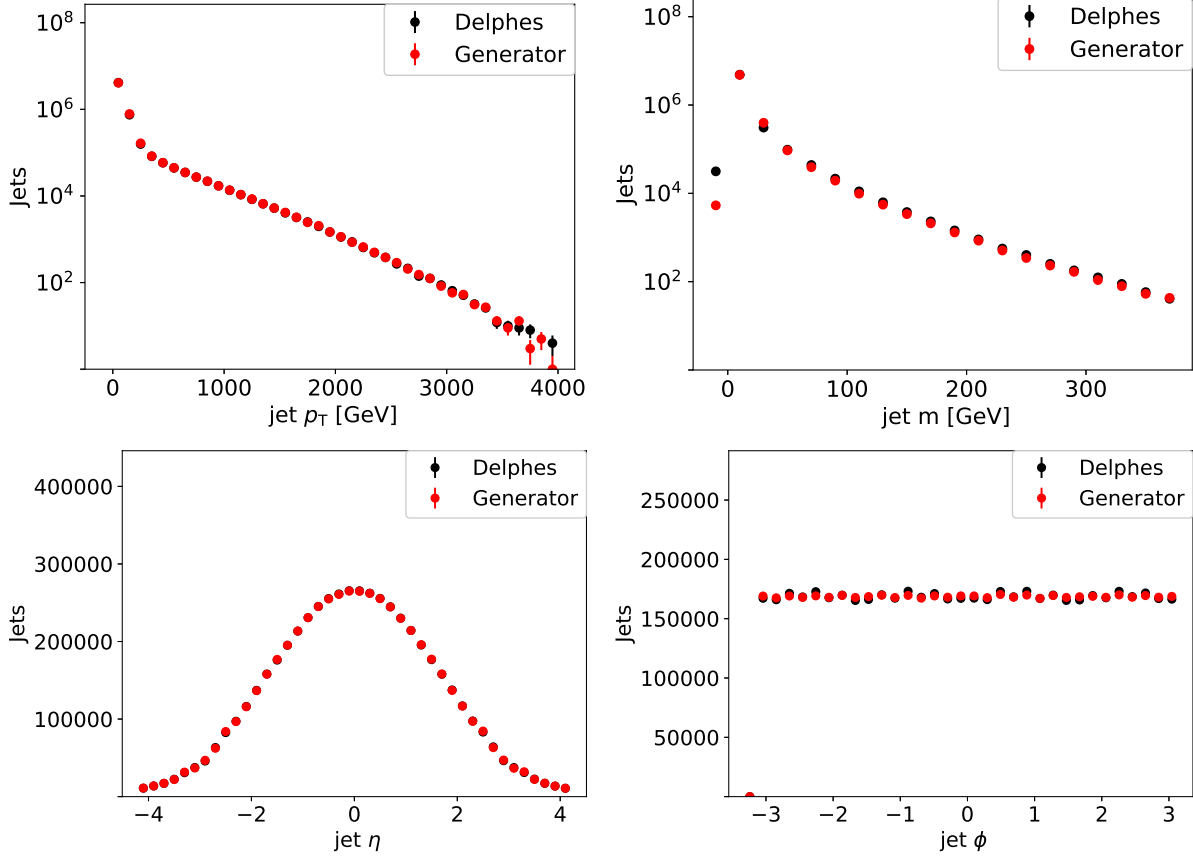


FIG. 2. Input variable shapes for truth-level (red) and detector-level quantities (black).

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 1×10^{-4} . The NN is implemented using Keras [8] with a TensorFlow [9] 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 Fig. ?? . Good agreement is observed between the Delphes and NN PDFs, showing that the NN has learned the bulk distribution.

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 Fig. 4. These PDFs are then randomly sampled to produce a NN jet that mimic the detector-level jet. A comparison

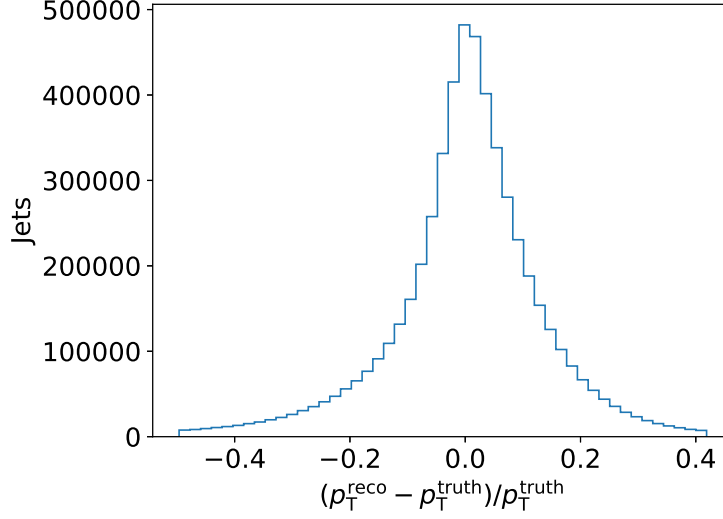


FIG. 3. Differences between truth-level and detector-level p_T .

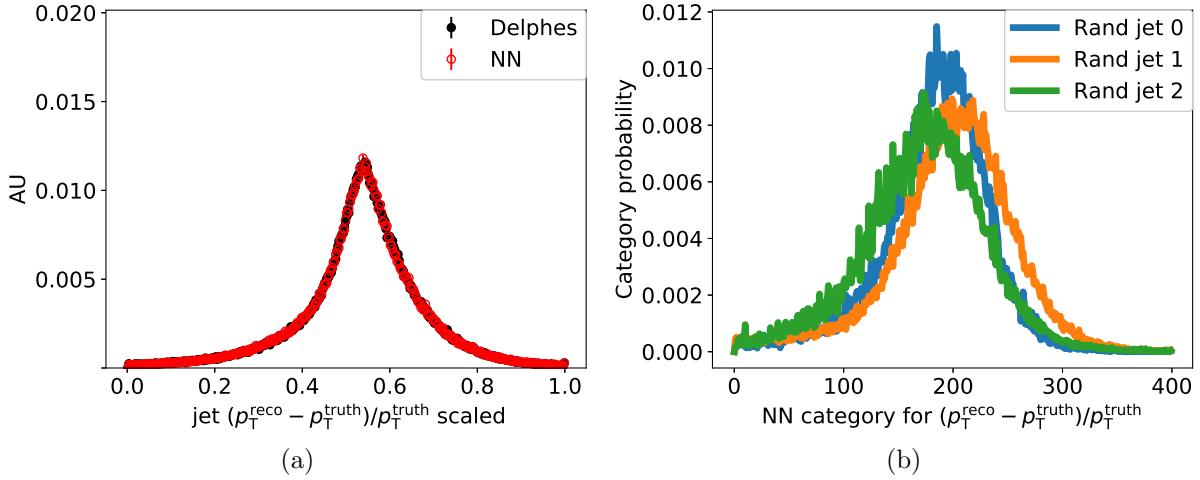


FIG. 4. NN-generated jet $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$ compared to detector-level jet $(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}$ (a). NN-generated jet PDFs for three randomly selected truth-level jets (b).

of the NN-generated and Delphes jet p_T distributions for the testing sample and for NNs trained with batch sizes 10000, 1000, 10, and 5 are shown in Fig. ?? while the distributions at lower p_T are shown in Fig. ?. The NN reproduces the jet p_T distribution of Delphes within 5% for reconstructed jets with $p_T > 20$ GeV for all batch size less than 100.

To test whether the NN learned correlations between input parameters and the p_T resolution, defined as $\frac{(p_T^{\text{reco}} - p_T^{\text{truth}})/p_T^{\text{truth}}}{p_T^{\text{truth}}}$, the jets were divided into central ($|\eta| < 3.2$) and forward ($|\eta| > 3.2$) jets, then the p_T resolution is compared between the two regions for both the

Delphes jets as well as the NN-generated jets. These two regions in the detector simulation have different calorimeter resolutions which results in different jet p_T resolutions and thus the p_T resolution is correlated with $|\eta|$. The resulting resolutions for both regions and for the training samples and jets are shown in Fig. 5. To quantify how the training batch size affects the NNs ability to reproduce the resolution of the forward region, which has significantly lower statistics and is a small subsample of the total training sample, the two-sample Kolmogorov–Smirnov statistic was calculated for NN and Delphes p_T resolution distribution for both the central and forward regions. The training sample was chosen for these tests because it had significantly more jets and due the low number of jets that have $|\eta| > 3.2$, as can be seen in Fig. 2. The NN clearly performs better for the rare subsample (the forward region) when trained on smaller batch sizes but the NN performs worse for batch size of five.

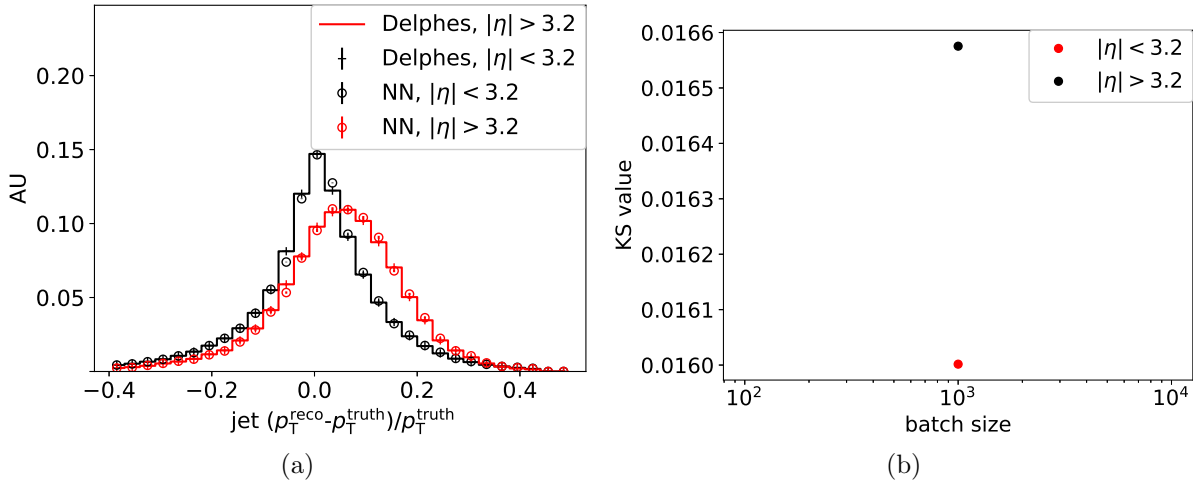
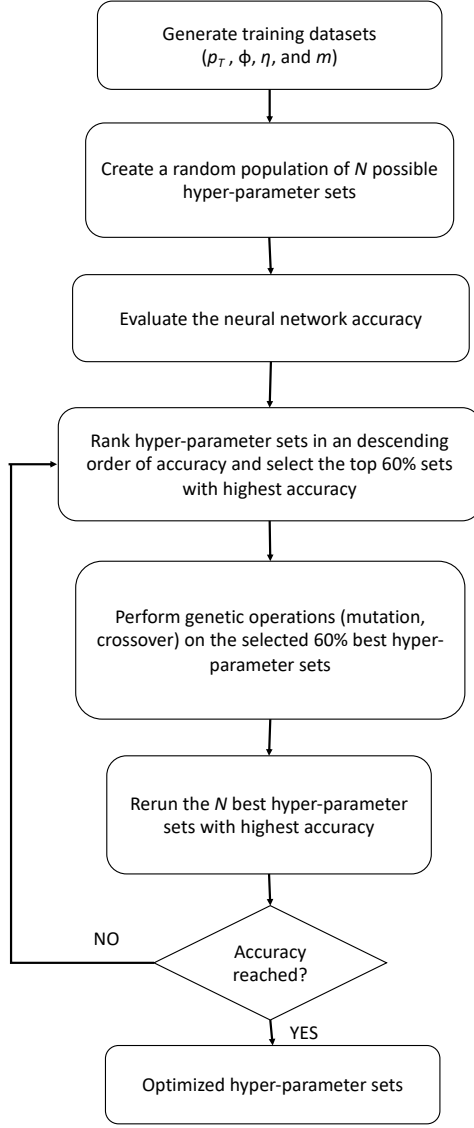


FIG. 5. Resolutions for the jet p_T (a) for the training sample for both the central and forward region using a batch size of 10 during the training of the model. The KS statistic for NN and Delphes resolutions as a function of batch size (b) is also shown. The red dots show the KS statistic for the central region while the black dots show the KS statistic for the forward region.

VII. HYPERPARAMETER OPTIMIZATION

In order to acquire the optimal performance with an NN with a minimal amount of parameters, a genetic algorithm (GA) was used to optimize the NN hyperparameters. Hyperparametrization, that is, the determination of the optimal parameters associated with a complex NN (number of layers, number of neurons on each layer, choice of the learning

rate, mini-batch size), is a challenging optimization problem in a high dimensionality space,
 especially for nonlinear NNs. Ultimately, the range of hyperparameter values that must be
 explored grows dramatically as the complexity of the deep NN increases. To address this
 challenge, we have used a (GA), which is an evolutionary algorithm that mimics the process
 of natural selection. Previously successful applications of the GA in similar contexts include
 the determination of the parameters of a complex force field [10, 11]. A flowchart of the
 ML/GA optimization protocol is depicted in Fig. 6. The optimization starts with a set of
 parent parameters, which is defined as the population. For hyperparameter optimization
 of the NN, the parents sets are number of layers, number of neurons on each layer, choice
 of the learning rate or mini-batch size. Some of the individuals in this population present
 a better fit, which in the context of NNs means a higher value for the accuracy of the
 network. Fit parents survive and are allowed to mate, which is accomplished by crossing
 patterns with other fit individuals. During crossover, random mutations in the genes are
 also allowed, to a certain degree, to avoid a stagnant gene pool and a better sampling of
 the parameters space. The offspring individuals form the next generation of parents and
 this process continues until some predefined criteria are met. Sets of parameters are then
 ranked in ascending order. After the ranking, a nonlinear roulette wheel selection[?] was
 performed to select the best 60% members, that is, the ones with lowest values of accuracy,
 which were then subjected to genetic operations: mutation and crossover with a 3% crossover
 rate. These mutations introduce sufficient diversity into the population, and the nonlinear
 selection scheme helps to avoid premature convergence of the ML/GA run. After the genetic
 operations, both the parent and offspring sets of parameters are ranked by their value of
 accuracy. The best hyperparameter sets are then chosen to constitute the next generation.
 Such an optimization routine ensures that only satisfactory hyperparameter sets survive
 after each generation; upon repeating this workflow for sufficient generations and sampling
 viable regions in the parameter space, we performed three separate ML/GA runs starting
 with different random populations. From each of the converged ML/GA run, we chose the
 final hyperparameter set corresponding the highest value of accuracy.



(a)

FIG. 6. Flow chart of the hyperparameter optimization used to optimize the resolution NN.

VIII. 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 changes in resolutions of different regions of the detector based on the NN inputs during training. For the NN to learn these correlations between variables for small subsets of objects, such as for forward jets, small batch size training was necessary. The NN learned the truth-to-reconstruction transformation without requiring manual binning to capture the differences in resolutions of particular

181 subsamples. This method should be easily extendable to additional reconstructed quantities
182 and could be used to model the ATLAS and CMS detector. The method described in this
183 paper thus allows for automated detector parameterization which can facilitate phenom-
184 logical, efficient Geant4 simulation, and upgrade studies.

185 ACKNOWLEDGMENTS

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

-
- 196 [1] S. Agostinelli et al. GEANT4: A Simulation toolkit. *Nucl. Instrum. Meth. A*, 506:250, 2003.
197 [2] J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaître, A. Mertens, and M. Sel-
198 vaggi. DELPHES 3, A modular framework for fast simulation of a generic collider experiment.
199 *JHEP*, 02:057, 2014.
200 [3] J. Alwall, R. Frederix, S. Frixione, V. Hirschi, F. Maltoni, et al. The automated computation
201 of tree-level and next-to-leading order differential cross sections, and their matching to parton
202 shower simulations. *JHEP*, 07:079, 2014.
203 [4] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. FastJet User Manual. *Eur. Phys. J. C*,
204 72:1896, 2012.
205 [5] Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. The anti- k_t jet clustering algorithm.
206 *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.