

# Particle-based Fast Jet Simulation at the LHC with Variational Autoencoders

Mary Touranakou · Nadezda Chernyavskaya · Javier Duarte · Dimitrios Gunopulos · Raghav Kansal · Breno Orzari · Maurizio Pierini · Thiago Tomei · Jean-Roch Vlimant

Received: date / Accepted: date

**Abstract** We study how to use Deep Variational Autoencoders for a fast simulation of jets of particles at the LHC. We represent jets as a list of constituents, characterized by their momenta. Starting from a simulation of the jet before detector effects, we train a Deep Variational Autoencoder to return the corresponding list of constituents after detection. Doing so, we bypass both the time-consuming detector simulation and the collision reconstruction steps of a traditional processing chain, speeding up significantly the events generation workflow. Through model optimization and hyper-

parameter tuning, we achieve state-of-the-art precision on the jet four-momentum, while providing an accurate description of the constituents momenta, and an inference time comparable to that of a rule-based fast simulation.

## 1 Introduction

At particle colliders, collimated sprays of particles are produced as a consequence of the parton shower and hadronization processes typical of Quantum Chromodynamics (QCD). These sprays of particles, called *jets*, are reconstructed applying a recombination clustering algorithm, exploiting physics-inspired metrics such as the anti- $k_t$  distance [1]. Often, jets are clustered from energy deposits recorded in the electromagnetic and hadronic calorimeters of a particle detector. At the LHC, jets can be clustered from a list of reconstructed particles, the so-called particle-flow (PF) candidates [2, 3]. In this case, jets would be sparse sets of objects (the constituents), each represented by its momentum<sup>1</sup> and possibly a set of auxiliary features, such as the nature of the particle (electron, muon, etc.), its electromagnetic charge, etc.

A time- and resource-effective strategy to simulate jet production is a fundamental asset for physics studies

<sup>1</sup> As common for collider physics, we use a Cartesian coordinate system with the  $z$  axis oriented along the beam axis, the  $x$  axis on the horizontal plane, and the  $y$  axis oriented upward. The  $x$  and  $y$  axes define the transverse plane, while the  $z$  axis identifies the longitudinal direction. The azimuth angle  $\phi$  is computed with respect to the  $x$  axis. The polar angle  $\theta$  is used to compute the pseudorapidity  $\eta = -\log(\tan(\theta/2))$ . The transverse momentum ( $p_T$ ) is the projection of the particle momentum on the  $(x, y)$  plane. We fix units such that  $c = \hbar = 1$ .

Mary Touranakou  
European Organization for Nuclear Research (CERN), CH-1211 Geneva 23, Switzerland  
National and Kapodistrian University of Athens, Athens 157 72, Greece

Nadezda Chernyavskaya  
European Organization for Nuclear Research (CERN), CH-1211 Geneva 23, Switzerland

Javier Duarte  
University of California San Diego, La Jolla, CA 92093, USA

Dimitrios Gunopulos  
National and Kapodistrian University of Athens, Athens 157 72, Greece

Raghav Kansal  
University of California San Diego, La Jolla, CA 92093, USA

Breno Orzari  
Universidade Estadual Paulista, São Paulo/SP - CEP 01049-010, Brazil

Maurizio Pierini  
European Organization for Nuclear Research (CERN), CH-1211 Geneva 23, Switzerland

Thiago Tomei  
Universidade Estadual Paulista, São Paulo/SP - CEP 01049-010, Brazil

Jean-Roch Vlimant  
California Institute of Technology, Pasadena, CA 91125, USA

at the CERN Large Hadron Collider (LHC). Whether testing predictions of the Standard Model (SM), searching for evidence of physics beyond the SM, or assessing systematic uncertainties associated to a given measurement, physicists rely on an accurate simulation of the full collision process and of the detector response. This implies the need for a detailed simulation of a chain of very different steps, from the proton collision to the generation of the collected signal in the detector sensors. Typically, physicists simulate datasets at least 10 times larger than the amount of collected data, so that the precision on the final measurement is not limited by the amount of simulated data at hand.

A typical high-energy physics (HEP) simulation software relies on the GEANT4 [4] library to model the interaction of particles traversing the detector material. This approach, based on Monte Carlo (MC) techniques, provides a typical accuracy at the percentage level, but it comes at a high cost in terms of computing resource utilization. At the LHC, the simulation workflow consumes up to  $\sim 50\%$  of the total computing resources of an experiment. With the amount of collected data increasing, the need for MC simulation is going beyond what the available computing infrastructure could sustain. Projected to the planned High-Luminosity LHC upgrade, this trend will eventually become unsustainable [5]. Jet simulation is one of the most expensive tasks, since jets are very abundant in LHC collisions and are made of many particles. Since the simulation of each of these particles is a demanding operation, the possibility of simulating all particles in a jet at once would be a major improvement. The main difficulty of this task is to match state-of-the-art accuracy, which also depends on the specific use case, e.g., which quantity a specific analysis uses. For instance, an algorithm reproducing the jet kinematic but not describing the angular distribution of the particles in the jet might be suitable for an analysis needing a good description of the jet momentum (e.g., a dijet resonance search), but not for an analysis exploiting jet substructure techniques (e.g., an all-jet diboson resonance search). Certainly, an algorithm describing every aspect of jet physics would be an ideal solution.

As a first step, a typical LHC simulation makes use of an *event generator*, modeling a proton-proton collision, the consequent production of quarks and gluons (among other particles), and their hadronization into jets of particles. Since no detector response is involved, this step typically requires a relatively modest amount of computing resources<sup>2</sup>. In addition, its content is independent of experimental aspects (reconstruction soft-

ware, detector configuration, etc.), so that a dataset of these *generator-level* events can be stored for long term and used many times (as is the case for the CMS experiment). Computing requirements significantly increase when detector effects are to be taken into account. At first, one typically uses GEANT4 to simulate the detector response. Then, the reconstruction software runs on the event and produces the objects (e.g., the PF candidates for CMS), eventually clustered into jets. These two steps could be bypassed by a *jet response function*, taking as input the list of jet constituents at generator level and returning the list of constituents at reconstruction level. In this paper, we aim at approximating this jet response function with a Variational Autoencoder (VAE), trained using generator-level jets as input and the corresponding reconstruction-level jets as a target. We represent a jet as a list of particles' momenta. Doing so, the VAE returns a reconstructed jet in a format that is already compatible with a typical PF-based analysis software. A different approach to the problem of data sparsity consists of representing the jet as a point cloud, as proposed for many HEP-specific problems [7]. We investigated that approach when training a Generative Adversarial Network [8] to generate the list jet constituent momenta from random numbers. A similar approach was presented in Ref. [9], where a graph VAE was used to generate detector *hits* in a jet, from which jet constituents could be reconstructed using standard rule-based algorithms, e.g., PF reconstruction [2, 3]. This work has many common points with Ref. [9], with two main differences: we do not use graph architectures, and we aim at learning the detector response and bypassing the standard rule-based reconstruction algorithms (to offer further speed up of the simulation process). We do so by taking the reconstructed jet as a target. In this respect, our algorithm could be used to replace detector parametrization approaches now used in Fast Simulation tools [10–12], while the algorithm of Ref. [9] aims at speeding up a GEANT-based full simulation. In the future, both approaches will be useful to HEP experiment and, most likely, the ultimate generative model will emerge from a combination of the two.

This paper is organized as follows: Section 2 discusses related work. Sections 3 and 4 describe the benchmark dataset and the model architecture, respectively. A strategy to apply the model to realistic use case is discussed in Section 5. Training results are discussed in Section 6. Conclusions are given in Section 7.

<sup>2</sup> This picture could in principle change if next-to-leading order precision would be adopted as a default. On the other

hand, ongoing work on parallelizing event generation libraries on GPUs [6] may compensate for this precision increase.

## 2 Related Works

In the recent past, several studies explored the possibility of speeding up the data simulation process using generative models based on deep neural networks (NNs). In particular, convolutional neural networks (CNNs) have been proposed to generate single-particle showers in a calorimeter [13–19], full jets at the LHC [20–22], multi-dimensional functions of kinematic quantities [23, 24], event kinematics at colliders [25, 26], and cosmic ray showers [27]. Both generative adversarial networks (GANs) [28–30] and variational autoencoders (VAEs) [31] were considered.

These studies clearly demonstrate that integrating deep generative models in the data simulation workflows of HEP experiments could lead to an important saving in terms of computing resources. But there is an objective difficulty when scaling up these proof-of-concept solutions to production-ready simulation tools. The main problem lies in the complexity of a typical HEP detector, characterized by detector elements with different technology and geometry, partially overlapping with each other and with passive material (e.g., absorbers in calorimeters) in between. As a consequence of this, a typical HEP dataset consists of a sparse set of energy deposits, which often cannot be represented as a regular grid of pixels. Future detectors will be characterized by higher granularity, with small-size sensors designed to resolve hits from individual particles in dense environments, such as jet cores. This will make the sparsity of the event even more complicated. This is the main reason why most of the great ideas based on CNNs had so far a little impact on HEP experiments. Instead, other approaches were explored as alternatives to CNNs, e.g., a recurrent neural network (RNN) trained adversarially [32], graph neural networks [8, 9, 33], or normalizing flows [34]. Similar issues are present in other domains, e.g., galaxy simulation in cosmology [35].

In this paper, we investigate an alternative strategy to overcome difficulties with the peculiar nature of HEP data. In previous studies, we discussed how to sample jets as sparse data from a probability density function, modeled using deep generative models. To this purpose, we considered both GANs [8] and VAEs [36]. Here we take a different approach, in which the input is a generator-level jet (as opposed to a vector of random coordinates in some latent space) and the aim of the training is to learn a morphing function from generator to reconstruction level.

The strategy is similar to what is discussed in Ref. [37], where a similar approach is followed to morph a set of analysis-specific features from generator- to reconstruction-level precision.

## 3 Dataset

The reference dataset consists of jets generated in  $pp \rightarrow WW$  collisions at a center-of-mass energy  $\sqrt{s} = 13$  TeV. The  $W$  bosons are forced to decay to quarks, that then shower to jets. The event generation is performed using PYTHIA8 [38]. The generated list of particles is passed to DELPHES [10], which applies detector effects using the CMS DELPHES description. At this stage, additional collision events are superimposed to the generated collision, to mimic the effect of so-called *pileup*. The number of collisions is randomly sampled from a Poisson distribution with expectation value set to 50, in agreement with the expected LHC running conditions for Run 3. The DELPHES particle-flow reconstruction algorithm is applied to the event, returning the list of reconstructed particles. Reconstructed particles are required to have  $p_T > 250$  MeV and be within  $|\eta| < 3.2$ . These particles are then clustered into jets using the anti- $k_t$  [1] algorithm with jet-size parameter  $R = 0.5$ . Jets with  $p_T > 200$  GeV and within  $|\eta| < 2.5$  are retained. These jets represent the target dataset. With the same setting, generator-level jets are clustered from the stable and detectable particles produced in the collision, before detector effects are taken into account. These jets represent the input dataset.

Target jets within  $p_T$  and  $\eta$  acceptance are matched to input jets minimizing the angular distance  $\Delta R = \sqrt{\Delta\phi^2 + \Delta\eta^2}$ . An input-target pair is formed taking the closest input jet to each target jet. For both input and target jets, constituents are ordered by decreasing  $p_T$  and the first 50 particles are retained. When fewer particles are present, the list is zero-padded. The list contains the momentum of each constituent in Cartesian coordinates  $(p_x, p_y, p_z)$ . The constituent mass is implicitly assumed to be zero. The main advantages of this specific choice are: it retains information of the position of the jet in the detector (as opposed to local coordinate choice); the distribution of these quantities is unbounded and symmetric around 0, which makes the learning process easier. In particular, we avoid issues related to the periodicity of  $\phi$  and hard-threshold at boundaries (e.g., on  $p_T$ ).

We apply feature-dependent standardization by subtracting the mean and scaling the features to unit variance. During early stages of this work, we verified that these choices help the model training to converge to more accurate configurations of the network weights. After this pre-processing, each jet is represented as a 2D array of  $3 \times 50$  numbers. The whole dataset includes  $\sim 1.7$ M jets. We split these data in three parts: 60% for training, 20% for validation, and 20% for testing. The dataset is published on Zenodo [39].

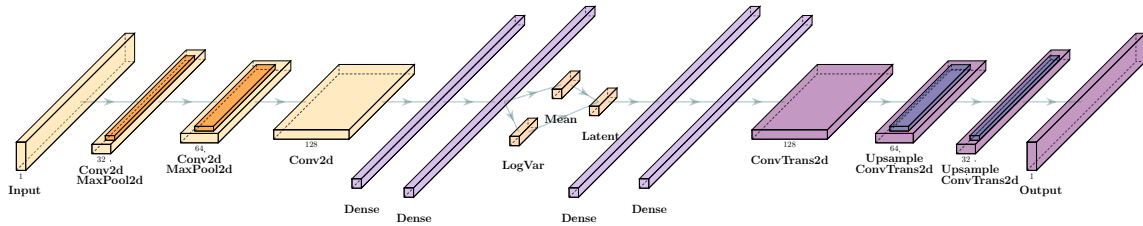


Fig. 1 Graphical representation of the VAE architecture.

#### 4 VAE Architecture

The VAE architecture is schematically shown in Fig. 1. The encoder receives a single-channel  $3 \times 50$  table, which is processed by three 2D convolutional layers, with 32  $3 \times 5$  kernels, 64  $1 \times 5$  kernels, and 128  $1 \times 5$  kernels, respectively. The stride is set to 1 and zero padding is used when the kernel arrives at the edge of the table. The output tensor is flattened and passed to two dense layers, with 640 and 150 nodes, respectively. From the second layer, two 20-dimension vectors are derived, corresponding to the mean  $\mu$  and log-variance values of the latent-space variables  $z$ . These values are used to define the Gaussian prior function from which a set of  $z$  values is sampled and passed to the decoder. The decoder architecture mirrors the encoder, with the Conv2D layers being replaced by ConvTrans2D layers. Leaky ReLU activation functions are used across the whole architecture with the coefficient of the negative slope set to 0.1, except for the encoder and decoder output layers, for which a linear activation function is used. Max pooling with the kernel size  $1 \times 2$  and stride of 2 is applied after each of the first two convolution operations in the encoder. Respectively, two upsampling operations are used in the decoder, each placed before the last two ConvTrans2D layers with a bilinear interpolation scheme. Dropout is added to the output of the first dense layer in the encoder and to the output of the two dense layers in the decoder with a dropout rate of 0.2. The deep learning (DL) model is implemented in PyTorch [40].

The model is trained using an input dataset, containing the list of particles at generator level, and a target dataset, containing the corresponding list after detector effects and event reconstruction. In this way, the model is trained to regress the detector response function starting from a generator-level jet, i.e. it corresponds to what a Fast Simulation software in HEP computing literature. For this kind of application, a typical state-of-the-art simulation has a 10% accuracy on jet kinematic properties.

The model training is performed minimizing a domain-specific loss function:

$$L^{\text{VAE}} = \frac{1}{N} \sum_{i=0}^N \left[ \beta D_{\text{KL}}^i + (1 - \beta) (L_{\text{R}}^i + \alpha_m (m_{\text{jet}}^i - \hat{m}_{\text{jet}}^i)^2 + \alpha_{p_{\text{T}}} (p_{\text{T}}^{\text{jet},i} - \hat{p}_{\text{T}}^{\text{jet},i})^2) \right] \quad (1)$$

where  $N$  is the dataset size,  $L_{\text{R}}$  is the reconstruction loss (i.e., a distance between the target and the output),  $D_{\text{KL}}$  is the Kullback–Leibler (KL) divergence regularizer usually employed to force the data distribution in the latent space to a multi-dimensional Gaussian with unitary covariance matrix [41], and  $\beta$  is a parameter that controls the relative importance of the two terms [42]. The reconstruction loss  $L_{\text{R}}$  is computed using the permutation-invariant Chamfer loss [43]:

$$L_{\text{R}} = \sum_i \min_j (p_i - \hat{p}_j)^2 + \sum_j \min_i (p_i - \hat{p}_j)^2, \quad (2)$$

where  $p_i$  is the feature for the  $i$ -th target particle, and  $\hat{p}_j$  is the corresponding quantity for the  $j$ -th output particle. By construction, this quantity is invariant under the permutation of the input or output particle lists. We also experimented with a mean squared error (MSE) loss, observing typically worse results.

In Eq. (1),  $p_{\text{T}}^{\text{jet}}$  and  $m^{\text{jet}}$  are the transverse momentum and mass of a target jet, respectively. Jet features are computed from the momenta of the target-jet constituents, while  $\hat{p}_{\text{T}}^{\text{jet}}$  and  $\hat{m}^{\text{jet}}$  are the corresponding quantities computed from the model output. The coefficients  $\alpha_m = 1.0$  and  $\alpha_{p_{\text{T}}} = 0.1$  were chosen such that the reconstruction loss  $L_{\text{R}}$  and the jet- $p_{\text{T}}$  and jet-mass MSE constraints in Eq. (1) have similar magnitudes. The expression in Eq. (1) is only one of the possible ways one could enforce kinematic constraints on the jet generator. Similar approaches have been followed in previous works, e.g., for particle energy in GAN-based single-particle generators [14, 44]. The main difference here is that the quantity on which the constraint is applied on is analytically computed from the output list, as opposed of being regressed from an image. We also

tried other combinations of kinematic constraints, e.g., the three momentum components in Cartesian coordinates, observing similar or worse results after training.

## 5 Target application

The aim of this work is to create a fast-simulation workflow for an analysis demanding a large sample of multijet events. As a reference, we consider the case of dijet resonance searches. Traditionally, these searches are carried on through a bump-hunt maximum-likelihood fit, in which the background is analytically modeled [45, 46]. Large samples of simulated multijet events are used to find an adequate model. In addition, a novel simulation-assisted strategy uses ratios of simulated distributions to avoid the need of a specific background analytical model [47]. Also in this case, having at hand a large simulated sample is crucial. Finally, the proposed strategy could be crucial to move the default event-generation precision to next-to-leading order, trading simulation computing time for generation computing time. This would be also relevant for analyses exploiting angular information about the dijet system [48]. Similar considerations hold for multijet searches.

The reference analysis requires an accurate model of jet kinematic properties for jets momenta larger than 200 GeV. As we will see, this can be achieved. But some care is required to model the sharp threshold at 200 GeV. As we experienced in the early stages of this study, an ML-based simulation struggles to model such a sharp threshold. As a solution, we extend the jet phase space in the training sample down to  $p_T > 130$  GeV and we apply the selection of  $p_T > 200$  GeV on the jets obtained as output of the VAE. A similar problem exists for the  $p_T$  threshold of the jet constituents. In this case, we also extend the  $p_T$  range of the predicted model down to  $p_T > 0$  MeV and apply the selection  $p_T > 250$  MeV afterwards. Similar considerations hold for  $\eta$ , where the acceptance requirements are imposed on the predicted jets after inference.

As discussed in Section 6, this setup provides an adequate description of the jet kinematic but it fails in providing an accurate description of jet substructure. In this respect, the proposed model cannot be extended to other dijet bump-hunt analyses, e.g., diboson resonance searches, where the accurate modelling of jet substructure is crucial. One could have then enforced a limited scope from the beginning and avoided generating jet constituents, working directly at the level of jet four momenta. However, we see two added values in working with jet constituents: on the one hand, we obtain a faithful description of the jet mass, the most crucial jet-substructure high-level features; on the other hand, we

establish a baseline model which could further improve to also model jet substructure. This will be the subject of future studies exploiting a permutation-equivariant graph architecture.

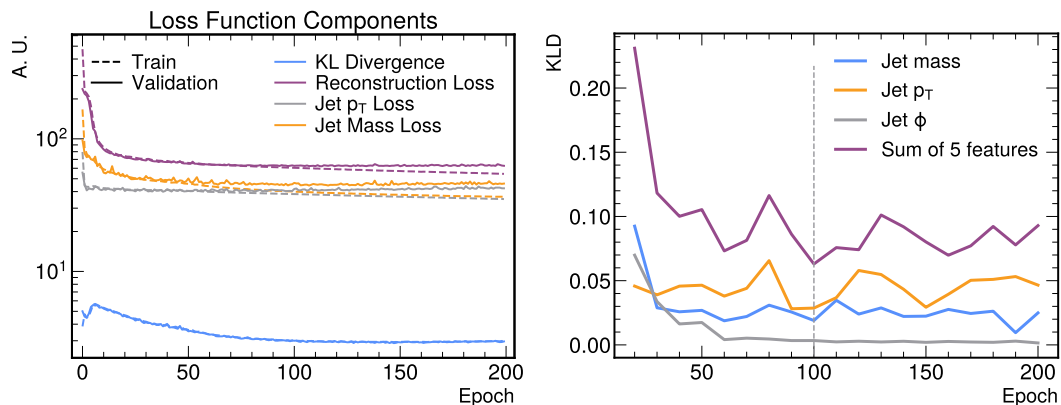
## 6 Results

We train all models using the Adam [49] optimizer with a learning rate of 0.0001 for 300 epochs. The training was repeated for several values of  $\beta$ , and the value corresponding to the best agreement between input and target ( $\beta = 1/9$ ) was chosen. During the training, we monitor the values of the total loss and its individual components evaluated on the training and validation datasets to check for over-training. To quantitatively evaluate the performance of different training settings, in addition to comparing the loss values, we also use the symmetrized version of the KL divergence (KLD) [50] between probabilities of the predicted and target jet-kinematic distributions (mass,  $p_T$ , energy,  $\eta$ ,  $\phi$ ). The KLD is computed every 10 epochs on the testing dataset, after rescaling the DL-predicted and target distributions so that the reconstructed distribution is contained in the  $[0, 1]$  range. The best model is selected based on the values of total and individual components of the loss and the KLD, while ensuring no over-training.

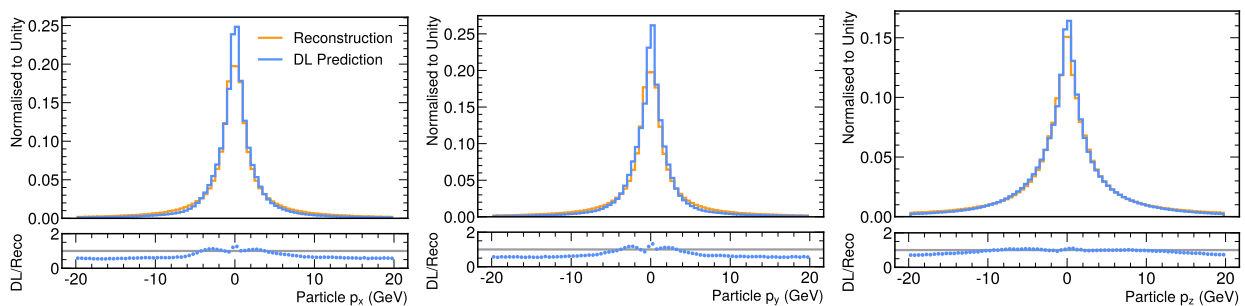
Figure 2 (left) shows the evolution of various contributions to the loss function (see Eq.(1)), evaluated on the training and validation datasets, as a function of the epoch. The monitored evolution of the KLD computed on the testing dataset for the three jet features (mass,  $p_T$ ,  $\phi$ ) and the total KLD sum of 5 features (mass,  $p_T$ , energy,  $\eta$ ,  $\phi$ ) is shown in Fig. 2 (right). The model from epoch number 100 is chosen as the best.

We compare the distributions of the DL predicted and the target  $p_x$ ,  $p_y$ , and  $p_z$  of the jet constituents in Fig. 3. These distributions are obtained applying the constituents acceptance thresholds (see Section 3) to the list of particles which is output by the VAE, as discussed in Section 5. We observe a good agreement between the model prediction and the target reconstruction.

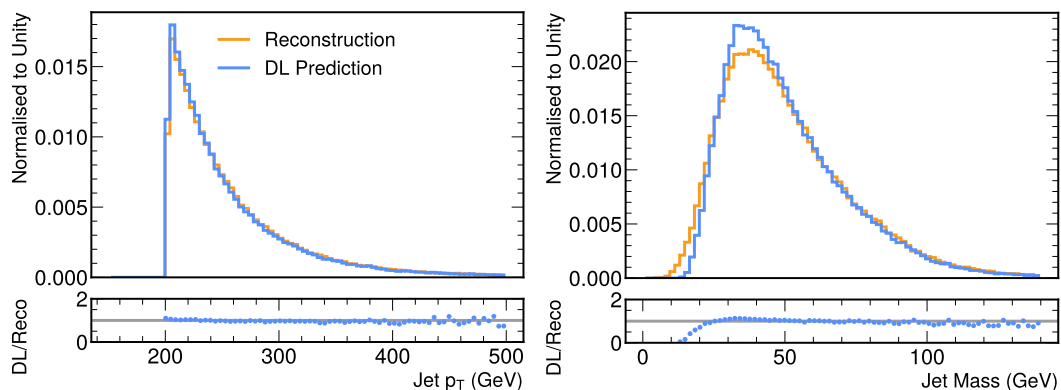
The output list of particles is then used to analytically compute the jet kinematic properties. Figure 4 (Fig. 5) shows the distribution of the jet kinematic properties explicitly used (not used) in the likelihood. The jet acceptance thresholds (see section 5) are imposed on the jet  $p_T$  and  $\eta$  for both the target and output jets. In general, a good agreement is observed. The residual discrepancies between the model prediction and the target reconstruction are smaller than the modelling differences typically observed between the



**Fig. 2** (Left) Evolution of various contributions to the loss, evaluated on the training (dashed lines) and validation (solid lines) datasets, as a function of the training epoch. (Right) Evolution of the KLD as a function of the training epoch. The KLD is computed on the testing dataset. The total KLD, a sum of the KLD for the 5 jet features (mass,  $p_T$ , energy,  $\eta$ ,  $\phi$ ), and three individual components are shown. The dashed vertical line indicates the chosen epoch.



**Fig. 3** Comparison of the jet-constituent  $p_x$  (left),  $p_y$  (center), and  $p_z$  (right) distributions, for the output (DL prediction) and target (Reconstruction) datasets. In the bottom panels, the ratio between the two distributions is shown. These distributions are obtained removing the zero-padding particles from the target list, and enforcing on the output of the DL model the same acceptance requirements that define the jet constituents (see section 3).

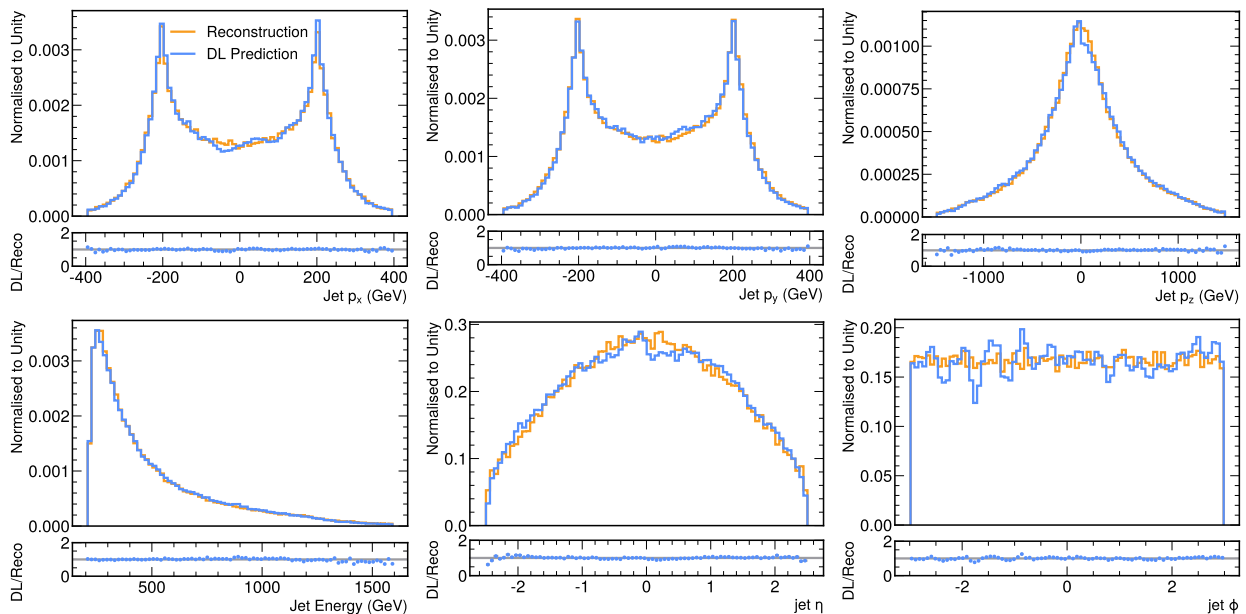


**Fig. 4** Comparison of the jet  $p_T$  (left) and mass (right) distributions, for the output (DL prediction) and target (Reconstruction) datasets. In the bottom panels, the ratio between the two distributions is shown. These distributions are obtained removing the zero-padding particles from the target list, and enforcing on the output of the DL model the same acceptance requirements that define the jet constituents (see section 3).

jets data and MC reconstruction. Remarkably, once forced to learn the jet mass and transverse momentum components, the model learns to model the entire jet kinematic, including non-linear functions of the three quantities above. This aspect proves that the training pro-

cess converges to a solution that preserves the main physics of the jet shower. At this stage, such a generator would be useful to generate events for most of the physics studies performed at the LHC.





**Fig. 5** Comparison of the jet  $p_x$  (top-left),  $p_y$  (top-center),  $p_z$  (top-right), energy (bottom-left),  $\eta$  (bottom-center), and  $\phi$  (bottom-right) distributions, for the output (DL prediction) and target (Reconstruction) datasets. In the bottom panels, the ratio between the two distributions is shown. These distributions are obtained removing the zero-padding particles from the target list, and enforcing on the output of the DL model the same acceptance requirements that define the jet constituents (see section 3).

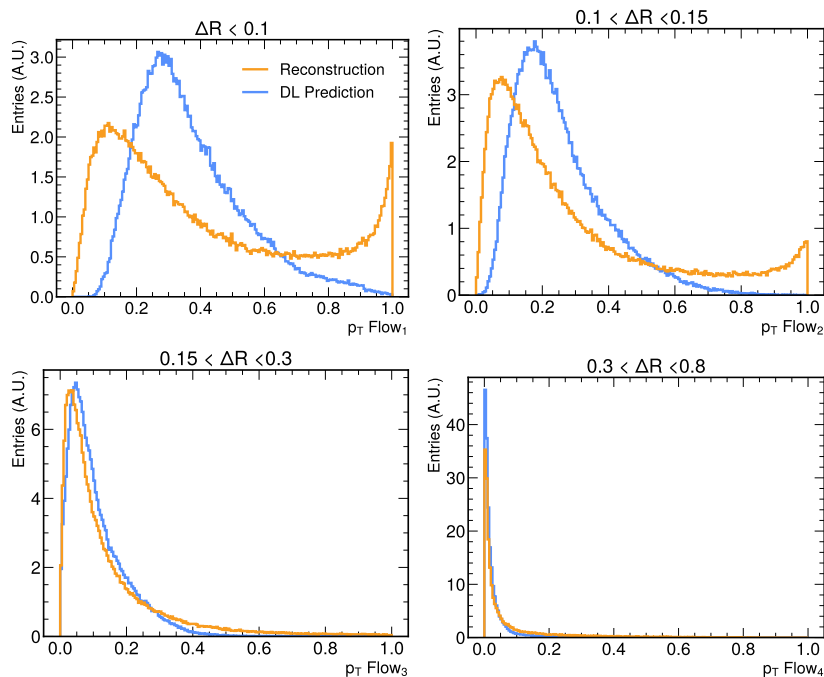
In Appendix A, we show the distribution of the jet features in the entire generation phase space, i.e., without enforcing the jet  $p_T > 200$  GeV requirement on the target and output jets. There, the problem of modeling the sharp  $p_T$  threshold is visible. Remarkably, this issue has little impact on the agreement observed in the jet mass distribution.

While our algorithm could serve the bulk of data analyses at the LHC, it still fails in faithfully describing the jet dynamic at constituents level. In fact, we verified that jet substructure quantities are not well reproduced. This is shown in Fig. 6, where the distribution of four momentum flows [51] are shown. The momentum flows are computed as  $Flow_n = \sum_p \frac{p_T^p}{p_T^{jet}}$ , where the sum runs over all particles with distance  $\Delta R = \sqrt{\Delta\phi^2 + \Delta\eta^2}$  from the jet axis falling within  $(n-1)/4 \times R$  and  $n/4 \times R$ , where  $R$  is the jet size parameter. We tracked the cause of the mismodeling to the noise induced by zero-momentum fake particles (both in input and target), resulting from zero-padding the jet representation to a fixed dimension. This problem could be solved moving to a graph-based VAE architecture, as in Ref. [9]. Through a PyTorch Geometric [52] implementation, for example, one could avoid the need of zero-padding the datasets, possibly leading to a better representation of the jet substructure. This approach will be investigated in future studies.

The inference speed of the algorithm was measured running it on 1000 generator-level jets, and measuring the execution time. The test was performed calling the Pytorch library from a python script and running the algorithm on different hardware platforms. We obtain an inference time of 0.007 (0.004) seconds per event when running on a Intel Xeon Silver 4216 CPU (NVIDIA T4), while the traditional approaches typically require  $O(100)$  seconds per event of CPU time. This demonstrates how the proposed strategy represents a major speedup with respect to currently employed simulation algorithms. One should also keep in mind that this is an overestimate of the actual inference time in real world C++ computing environment, where tools such as ONNX run time [53] typically offer a further speed up with respect to a python environment. When running on a GPU, one could further increase the throughput by running the difference on a batch of all jets in an event.

## 7 Conclusions

We present a jet fast-simulation algorithm based on a Variational Autoencoder, trained to learn the detector response function to a generator-level jet, represented as a list of particle momenta, and returning a list of reconstructed particle momenta. This algorithm correctly



**Fig. 6** Distribution of the four  $Flow_n$  quantities (see text) for the output and target jets.

captures the reconstructed jet kinematic with high accuracy.

By bypassing the detector simulation and particle reconstruction step, an algorithm of this kind could be important to make simulation on demand a concrete possibility at the High-Luminosity LHC.

The main strength of the current algorithm is in its speed and high accuracy when modeling jet kinematic quantities, which makes it applicable to the majority of LHC physics studies. Its main limitation stands with the poor description of the jet substructure, a consequence of the noise induced by zero-momentum ghost particles introduced to equalize the length of the input particle list. A possible solution to this problem could be the use of a graph architecture with variable-length input, which we aim at investigating in the future.

## Acknowledgement

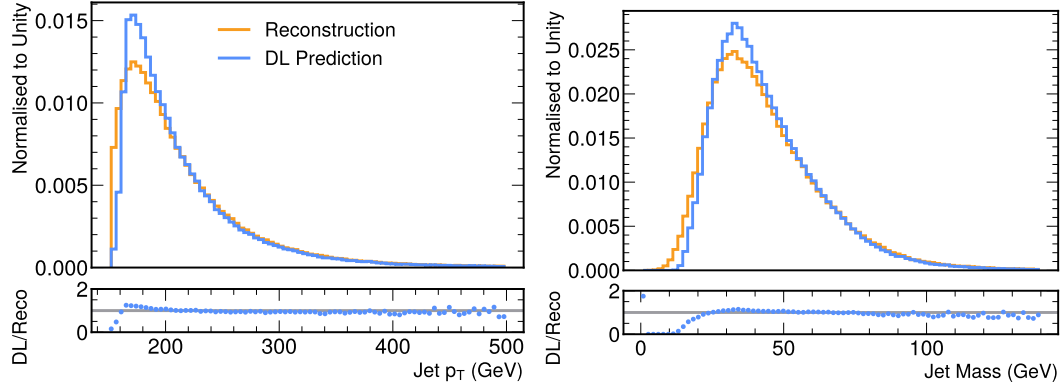
This work was supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreement No. 772369). R. K. was partially supported by an IRIS-HEP fellowship through the U.S. National Science Foundation under Cooperative Agreement OAC-1836650, and by the LHC Physics Center at Fermi National Accelerator Laboratory, managed and operated by Fermi Research Alliance, LLC under Contract No. DE-AC02-07CH11359 with the U.S. Department of En-

ergy (DOE). J. D. is supported by the U.S. Department of Energy (DOE), Office of Science, Office of High Energy Physics Early Career Research program under Award No. DE-SC0021187. D. G. is partially supported by the EU ICT-48 2020 project TAILOR (No. 952215). J-R. V. is partially supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (Grant Agreement No. 772369) and by the U.S. DOE, Office of Science, Office of High Energy Physics under Award No. DE-SC0011925, DE-SC0019227, and DE-AC02-07CH11359. B. O. and T. T. are supported by grant #2018/25225-9, São Paulo Research Foundation (FAPESP). B. O. is also supported by grant #2020/06600-3, São Paulo Research Foundation (FAPESP). This work was supported in part by NSF awards CNS-1730158, ACI-1540112, ACI-1541349, OAC-1826967, the University of California Office of the President, and the University of California San Diego’s California Institute for Telecommunications and Information Technology/Qualcomm Institute. Thanks to CENIC for the 100 Gpbs networks.

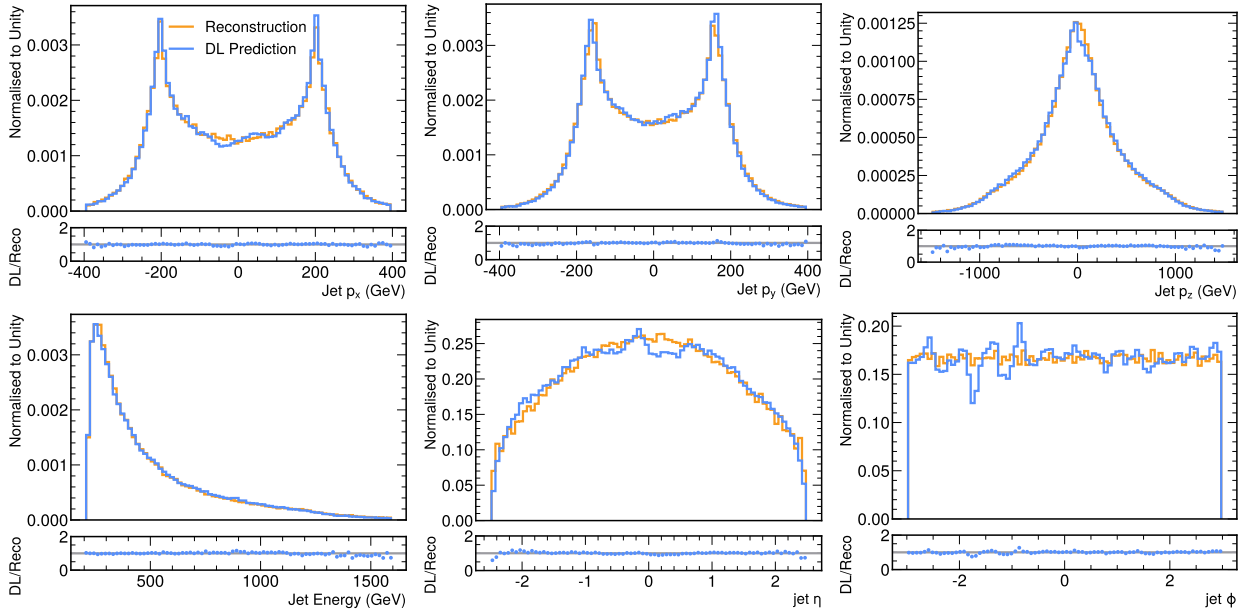


## A Appendix

In this appendix, we show the distribution of the jet features in the entire generation phase space before applying the jet  $p_T$  selection  $p_T > 200$  GeV. Figure 7 (Fig. 8) shows the distribution of the jet kinematic properties explicitly used (not used) in the likelihood within the extended jet phase space before any selections.



**Fig. 7** Comparison of the jet  $p_T$  (left) and mass (right) distributions, for the output (DL prediction) and target (Reconstruction) datasets before applying the jet  $p_T$  selection  $p_T > 200$  GeV. In the bottom panels, the ratio between the two distributions is shown. These distributions are obtained removing the zero-padding particles from the target list, and enforcing on the output of the DL model the same acceptance requirements that define the jet constituents (see section 3).



**Fig. 8** Comparison of the jet  $p_x$  (top-left),  $p_y$  (top-center),  $p_z$  (top-right), energy (bottom-left),  $\eta$  (bottom-center), and  $\phi$  (bottom-right) distributions, for the output (DL prediction) and target (Reconstruction) datasets before applying the jet  $p_T$  selection  $p_T > 200$  GeV. In the bottom panels, the ratio between the two distributions is shown. These distributions are obtained removing the zero-padding particles from the target list, and enforcing on the output of the DL model the same acceptance requirements that define the jet constituents (see section 3).

## References

1. Matteo Cacciari, Gavin P. Salam, and Gregory Soyez. The anti- $k_t$  jet clustering algorithm. *JHEP*, 04:063, 2008.
2. A. M. Sirunyan et al. Particle-flow reconstruction and global event description with the CMS detector. *JINST*, 12(10):P10003, 2017.
3. Morad Aaboud et al. Jet reconstruction and performance using particle flow with the ATLAS Detector. *Eur. Phys. J. C*, 77(7):466, 2017.
4. S. Agostinelli et al. GEANT4: A Simulation toolkit. *Nucl. Instrum. Meth.*, A506, 2003.
5. Johannes Albrecht et al. A Roadmap for HEP Software and Computing R&D for the 2020s. *Comput. Softw. Big Sci.*, 3(1):7, 2019.
6. K. Hagiwara, J. Kanzaki, Q. Li, N. Okamura, and T. Stelzer. Fast computation of MadGraph amplitudes on graphics processing unit (GPU). *Eur. Phys. J. C*, 73:2608, 2013.
7. J. Shlomi, P. Battaglia, and J. R. Vlimant. Graph Neural Networks in Particle Physics. *Machine Learning for Science and Technology*, 2(2):021001, 2021.
8. R. Kansal et al. Graph Generative Adversarial Networks for Sparse Data Generation in High Energy Physics. In *34th Conference on Neural Information Processing Systems*, 2020.
9. A. Hariri, D. Dyachkova, and S. Gleyzer. Graph generative models for fast detector simulations in high energy physics, 2021.
10. J. de Favereau et al. DELPHES 3, A modular framework for fast simulation of a generic collider experiment. *JHEP*, 02:057, 2014.
11. S. Sekmen. Recent Developments in CMS Fast Simulation, 2017.
12. Georges Aad et al. AtlFast3: the next generation of fast simulation in ATLAS, 9 2021.
13. M. Paganini, L. de Oliveira, and B. Nachman. Accelerating Science with Generative Adversarial Networks: An Application to 3D Particle Showers in Multilayer Calorimeters. *Phys. Rev. Lett.*, 120(4), 2018.
14. M. Paganini, L. de Oliveira, and B. Nachman. CaloGAN : Simulating 3D high energy particle showers in multilayer electromagnetic calorimeters with generative adversarial networks. *Phys. Rev. D*, 97(1), 2018.
15. M. Erdmann, J. Glombitza, and T. Quast. Precise simulation of electromagnetic calorimeter showers using a Wasserstein Generative Adversarial Network. *Comput. Softw. Big Sci.*, 3(1), 2019.
16. D. Salamani et al. Deep generative models for fast shower simulation in atlas. In *2018 IEEE 14th International Conference on e-Science (e-Science)*, page 348, 2018.
17. Dawit Belayneh et al. Calorimetry with deep learning: particle simulation and reconstruction for collider physics. *Eur. Phys. J. C*, 80(7):688, 2020.
18. E. Buhmann et al. Getting High: High Fidelity Simulation of High Granularity Calorimeters with High Speed. *Comput. Softw. Big Sci.*, 5(1):13, 2021.
19. E. Buhmann et al. Fast and Accurate Electromagnetic and Hadronic Showers from Generative Models. *EPJ Web Conf.*, 251:03049, 2021.
20. L. de Oliveira, M. Paganini, and B. Nachman. Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis. *Comput. Softw. Big Sci.*, 1(1), 2017.
21. P. Musella and F. Pandolfi. Fast and Accurate Simulation of Particle Detectors Using Generative Adversarial Networks. *Comput. Softw. Big Sci.*, 2(1), 2018.
22. S. Carrazza and F. A. Dreyer. Lund jet images from generative and cycle-consistent adversarial networks. *Eur. Phys. J. C*, 79(11), 2019.
23. Sydney Otten et al. Event Generation and Statistical Sampling for Physics with Deep Generative Models and a Density Information Buffer. *Nature Commun.*, 12(1):2985, 2021.
24. B. Hashemi et al. LHC analysis-specific datasets with Generative Adversarial Networks, 2019.
25. R. Di Sipio et al. DijetGAN: A Generative-Adversarial Network Approach for the Simulation of QCD Dijet Events at the LHC. *J. High Energy Phys.*, 08, 2020.
26. A. Butter, T. Plehn, and R. Winterhalder. How to GAN LHC Events. *SciPost Phys.*, 7, 2019.
27. M. Erdmann et al. Generating and refining particle detector simulations using the Wasserstein distance in adversarial networks. *Comput. Softw. Big Sci.*, 2(1), 2018.
28. I. J. Goodfellow et al. Generative Adversarial Networks, 2014.
29. M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein GAN, 2017.
30. I. Gulrajani et al. Improved training of Wasserstein GANs, 2017.
31. D. P Kingma and M. Welling. Auto-Encoding Variational Bayes, 2013.
32. J. Arjona Martínez et al. Particle Generative Adversarial Networks for full-event simulation at the LHC and their application to pileup description. *J. Phys. Conf. Ser.*, 1525(1):012081, 2020.
33. V. Belavin and A. Ustyuzhanin. Electromagnetic shower generation with graph neural networks. *Journal of Physics: Conference Series*, 1525:012105, 2020.
34. C. Krause and D. Shih. CaloFlow: Fast and Accurate Generation of Calorimeter Showers with Normalizing Flows, 2021.
35. F. Lanusse. Machine learning in cosmology, 2019. ACAT 2019, Saas Fee (CH).
36. B. Orzari et al. Sparse Data Generation for Particle-Based Simulation of Hadronic Jets in the LHC. In *38th International Conference on Machine Learning Conference*, 2021.
37. C. Chen et al. Data Augmentation at the LHC through Analysis-specific Fast Simulation with Deep Learning. *Comput Softw Big Sci*, 5(15), 2021.
38. T. Sjöstrand et al. An introduction to pythia 8.2. *Computer Physics Communications*, 191, 2015.
39. M. Touranakou et al. Particle-based Fast Jet Simulation at the LHC with Variational Autoencoders: generator-level and reconstruction-level jets dataset, 2022. <https://doi.org/10.5281/zenodo.6047873>.
40. A. Paszke et al. Pytorch: An imperative style, high-performance deep learning library, 2019. <https://arxiv.org/pdf/1912.01703.pdf>.
41. D. J. Rezende and S. Mohamed. Variational inference with normalizing flows, 2016.
42. I. Higgins et al.  $\beta$ -VAE: Learning basic visual concepts with a constrained variational framework. In *5th International Conference on Learning Representations*, 2017.
43. H. Fan, H. Su, and L. Guibas. A point set generation network for 3d object reconstruction from a single image, 2016.
44. D. Belayneh et al. Calorimetry with deep learning: particle simulation and reconstruction for collider physics. *The European Physical Journal C*, 80(7), 2020.
45. V. Khachatryan et al. Search for narrow resonances decaying to dijets in proton-proton collisions at  $\sqrt{s} = 13$  TeV. *Phys. Rev. Lett.*, 116(7):071801, 2016.

46. M. Aaboud et al. Search for new phenomena in dijet events using  $37 \text{ fb}^{-1}$  of  $pp$  collision data collected at  $\sqrt{s} = 13 \text{ TeV}$  with the ATLAS detector. *Phys. Rev. D*, 96(5):052004, 2017.
47. A. M. Sirunyan et al. Search for high mass dijet resonances with a new background prediction method in proton-proton collisions at  $\sqrt{s} = 13 \text{ TeV}$ . *JHEP*, 05:033, 2020.
48. A. M. Sirunyan et al. Search for new physics in dijet angular distributions using proton-proton collisions at  $\sqrt{s} = 13 \text{ TeV}$  and constraints on dark matter and other models. *Eur. Phys. J. C*, 78(9):789, 2018.
49. D. K. Kingma and J. Ba. Adam: A method for stochastic optimization, 2015.
50. S. Kullback and R. A. Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
51. M. L. Mangano et al. Physics at a 100 TeV pp Collider: Standard Model Processes, 2016.
52. M. Fey and J. E. Lenssen. Fast graph representation learning with pytorch geometric, 2019.
53. Onnx runtime. <https://www.onnxruntime.ai>, 2021.