

A Generative-Adversarial Network Approach for the Simulation of QCD Dijet Events at the LHC

Riccardo Di Sipio[†]

University of Toronto

E-mail: riccardo.disipio@utoronto.ca

Michele Fucci Giannelli

University of Edinburgh

E-mail: michele.fucci.giannelli@cern.ch

Sana Ketabchi Haghghat

University of Toronto

E-mail: sana.ketabchihaghghat@mail.utoronto.ca

Serena Palazzo

University of Edinburgh

E-mail: serena.palazzo@cern.ch

A Generative-Adversarial Network (GAN) based on convolutional neural networks is used to simulate the production of pairs of jets at the LHC. The GAN is trained on events generated using MADGRAPH5, PYTHIA8, and DELPHES3 fast detector simulation. We demonstrate that a number of kinematic distributions both at Monte Carlo truth level and after the detector simulation can be reproduced by the generator network.

XXIX International Symposium on Lepton Photon Interactions at High Energies - LeptonPhoton2019

August 5-10, 2019

Toronto, Canada

**This work was supported by the Natural Sciences and Engineering Research Council of Canada (NSERC) and of the Science, Technology and Facility Council (STFC). This project has received funding from the European Union Horizon 2020 research and innovation programme under grant agreement No 765710. We received the donation of two P6000 GPU cards (one per group) in support of our work by the NVIDIA Corporation GPU grant programme.

[†]Speaker.

1. Introduction

In the forthcoming years, experiments at the Large Hadron Collider (LHC) are expected to cope with a deluge of data. At the same time, the strategy to produce reliable and statistically large samples of simulated proton–proton (pp) inelastic collisions will be facing both technological limitations and new opportunities.

In the context of Machine Learning (ML) and specifically Deep Neural Networks (DNNs), an unsupervised learning technique called Generative-Adversarial Networks (GANs) has been proposed recently [1] to approximate a probability density function (pdf) from a set of unlabelled training examples. In practice, GANs have been widely used to generate photorealistic portraits [2], music [3], and recently also to simulate the response of a calorimeter to the passage of particles [4, 5, 6]. Other deep learning methodologies such as variational autoencoders (VAE) [7] have been applied to reproduce kinematic distributions learned from examples taken from Monte Carlo simulations. In this work, we successfully trained a GAN to reproduce kinematic distributions, improving on such previous attempts, thanks to a careful consideration of the implicit symmetries of the physics process under study and the employment of convolutional layers [8, 9, 10].

In a GAN, a generative network G transforms a vector of random numbers (input noise) $z \sim p_z$ into a sample carrying some physical meaning, which in this case are the four-momenta of the two jets. In practical applications, p_z is usually a N -dimensional uniform distribution in the range $[0, 1]^N$. Subsequently, a discriminative network D estimates the probability that a given sample comes either from the training data or the generator. The two samples are distributed with probability density functions p_{data} and p_{fake} respectively, with p_{data} fixed and usually estimated using a Monte Carlo method. The Nash equilibrium (min-max game) is reached when D is unable to distinguish fake examples from real data, hence the generator has been trained to be a good approximator of the data pdf , i.e. $p_{fake} \sim p_{data}$.

Currently, both the ATLAS [11] and CMS [12] experiments of the LHC at CERN deploy Monte Carlo (MC) event generators such as MADGRAPH5 [13], POWHEG-BOX [14] and MC@NLO [15] to simulate the hard-scattering (HS) process, PYTHIA8 [16] and HERWIG7 [17] for the parton-shower (PS) process, and a GEANT4 [18] simulation of the actual detector for the response of the experimental apparatus. Best estimates suggest that the simulation of a single event takes already several minutes [19], with $O(10^9)$ events to be generated for each simulation campaign, leading to a huge computational footprint both in terms of CPU usage and disk space. Detailed simulation based on GEANT4 will not be an affordable solution due to the large time required to simulate an event [19]. Both experiments are already using fast simulation and are developing new tools exploiting ML and other advanced statistical techniques. We will demonstrate that the same GAN used for reproducing the generator output can be also used to reproduce a simulation of a detector response with a significant time gain with respect to full simulation.

The code is publicly accessible on the online repository <https://gitlab.cern.ch/disipio/DiJetGAN>.

2. Physics of QCD dijet events

At hadron colliders such as the LHC, the most abundant kind of interaction between the two

colliding protons is the scattering between quarks and gluons (collectively referred to as *partons*). According to calculations based on the SM, these parton-scattering processes via strong interactions described by quantum chromodynamics (QCD) result in the overwhelming majority of cases in two outgoing partons which carry a net color charge and evolve from high to low virtuality producing parton showers, which eventually hadronize into collimated highly-energetic clusters of particles called *jets*. Hence, $2 \rightarrow 2$ parton scattering processes with a pair of jets in the final state are called *dijet events*. The relationship between the clusters and the original partons is revealed by the execution of a clustering algorithm [20]. One can think of a jet approximately as a cone of radius R whose axis correspond to the direction of flight of the initial parton. The size of the radius can be controlled by setting a distance parameter in the clustering algorithm.

For most analyses, the most relevant jets are produced with a large transverse momentum (p_T) and large angle with respect to the incoming partons. The jet mass, defined as the norm of the four-momentum sum of constituents inside a jet, is only loosely related to the mass of the originating parton, and comes mostly from the dynamics of strong interactions. Programs such as PYTHIA8 [16] and HERWIG7 [17] implement such calculations with beyond the leading-logarithm (LL) accuracy in what are called a Parton Shower algorithms. The jet mass also plays a key role in the identification of Lorentz-boosted hadronically decaying massive particles such as top quarks [21, 22], vector (W and Z) and Higgs bosons [23].

In the following sections, the agreement between MC calculations and the output of the GAN is evaluated by comparing the individual jets' and dijet system's transverse momentum, pseudo-rapidity¹ (η) and mass distributions. The χ^2 between the MC and the GAN distributions is used as figure of merit.

3. Monte Carlo Sample

A sample of 10 million dijet events has been generated using MADGRAPH5 and PYTHIA8, corresponding to an integrated luminosity of about 0.5 fb^{-1} . The response of the detector was simulated by a DELPHES3 [24] fast simulation, using settings that resemble the ATLAS detector. An average of 25 additional soft-QCD pp collisions (pile-up) were overlaid to reproduce more realistic data-taking conditions.

Electrons, muons, jets and missing transverse energy are reconstructed by DELPHES3 algorithms before and after the detector simulation. These two levels of reconstruction are referred to in the following sections as particle- and reco-level respectively. At particle-level, only stable final-state particles, *i.e.* particles that are not decayed further by the generator, and unstable particles² that are to be decayed later by the detector simulation, are considered. Jets were reconstructed using the anti- k_T algorithm [25] as implemented in FastJet [26], with a distance parameter $R = 1.0$.

To increase the number of events with both jets with $p_T > 250 \text{ GeV}$, a cut on the scalar sum of the transverse momenta of the outgoing partons $H_T > 500 \text{ GeV}$ was applied to the hard-scattering

¹Pseudorapidity is a commonly spatial coordinate describing the angle of a particle relative to the beam axis defined as $\eta = \frac{1}{2} \ln \frac{E+p_L}{E-p_L}$, where E is the energy and p_L is the longitudinal component of the momentum. It is related to the other components of the momentum via the relationship $|p| = p_T \cosh \eta$.

²Particles with a mean lifetime $\tau > 300 \text{ ps}$

process. Approximately 7.5 million events passed this selection at particle level, and about 4 million at reco level. The difference in efficiency between the two levels of the simulation can be understood in terms of distortions introduced by the detector, which smear the jets' transverse momentum distributions, hence lowering the number of events passing the final selection. These events were used to train the network in the subsequent steps.

4. Network Architecture

The overall architecture of the network, summarized in Fig. 1, is composed of two main blocks: a generator (G) and a discriminator (D), both based on convolutional layers. All layers have *LeakyReLU* activation functions [27] except the last layers that have either *tanh* or *sigmoid* for the generator and the discriminator respectively. The generator transforms a vector of 128 random numbers drawn from a uniform distribution in the $[0, 1]$ range into a vector of 7 elements representing the p_T , η and mass of the leading jet, and the p_T , η , ϕ and mass of the second-leading jet. The discriminator takes as input the array of 7 elements described above and gives as output a number d between 0 and 1 that is interpreted as the likelihood of the event being drawn from “real” MADGRAPH5 events ($d=1$) or from “fakes” generated by G ($d=0$).

The network is implemented and trained using KERAS v2.2.4 [28] with TENSORFLOW v1.12 [29] back-end. Input features are scaled in the range $[-1, 1]$ and pre- and post-processed using the SCIKIT-LEARN [30] and PANDAS [31] libraries. The loss function of the generator is mean squared error, while that of the discriminator is the binary cross-entropy. The optimizer is in both cases Adam [32] with learning rate $lr = 10^{-5}$, $\beta_1 = 0.5$ and $\beta_2 = 0.9$. The parameters described above are those that provide the best results among many values and configurations tested. Having reached a satisfactory performance, no further parameter optimisation was carried out.

5. Training

For the purpose of the training, all MC events were rotated so that the azimuthal ϕ angle of the leading jet is always zero. A significant performance improvement was achieved by exploiting the intrinsic ϕ symmetry in di-jet events; the ϕ of the leading jet is set to zero while the ϕ of the other jet is set to the absolute value of the difference in ϕ between the two generated jets. This transformation is reversed when events are generated. In order to further deploy the symmetries of dijet kinematics, every event is used twice: first in its original configuration, and then with the sign of the pseudorapidity of each jet reversed (η -flip). During event generation, the η of the jets is randomly flipped to remove any nonphysical effects that could be introduced by the GAN.

The network was then trained for 500,000 iterations with mini-batches of 128 events each, drawn from the original distribution and from the noise-generated fakes. It took about five hours to complete the training on a GPU NVIDIA Quadro P6000. For each iteration, we first trained the discriminator to distinguish between real and fake events. Then, the discriminator weights are fixed and the generator is trained. At the end of the training, the discriminator will be unable to distinguish between the two sets. In such balanced set-up, equilibrium is reached when the generator can fool the discriminator in 50% of cases.

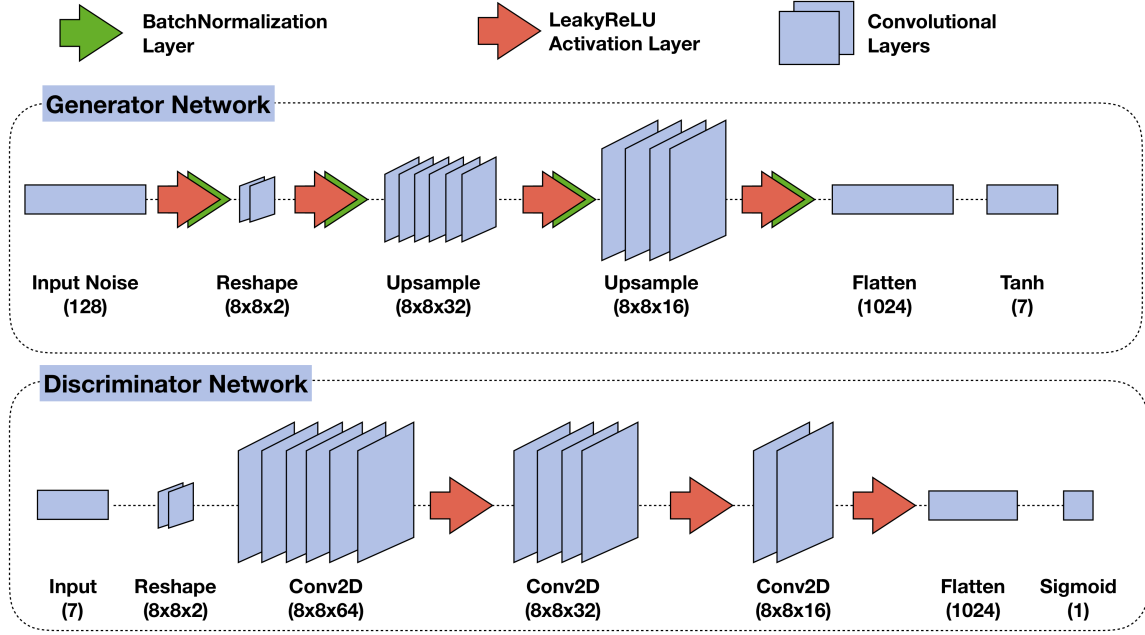


Figure 1: Network architecture: generator (top), discriminator (bottom). The GAN is composed by connecting the output of the generator to the input of the discriminator.

6. Event Generation and Final Results

During the training, the weights of the generator model are saved into a file every 5000 epochs and used subsequently to generate an arbitrary number of events. To match the size of the MADGRAPH5 samples, 10 million events were generated with the GAN. On average, it takes about 80 seconds to generate 1 million events on a GPU NVIDIA Quadro P6000. After the generation, events are filtered by applying the same kinematic cuts we applied to the real MC events, *i.e.* both jets with $p_T > 250$ GeV, ordered by decreasing p_T . Approximately 90% fulfill these requirements and are used to fill the histograms.

Figs. 2 shows the comparison of the two leading jets and dijet system kinematics, as they appear at the iteration that yields the best agreement in terms of overall χ^2 over degrees of freedom. Overall, the level of agreement is satisfactory over a large range of the kinematic regime.

To complement the investigations described above, we also trained our GAN on a sample of top-quark pairs decaying in the all-hadronic channel, *i.e.* $t\bar{t} \rightarrow WbW\bar{b} \rightarrow bq\bar{q}'\bar{b}q\bar{q}'$. A cut on the scalar sum of the transverse momenta of the outgoing partons $H_T > 700$ GeV was applied. Also, both jets at particle level are required to have a transverse momentum $p_T > 350$ GeV and mass < 500 GeV. In this region of the phase space, the jet mass is expected to have a peak around the top mass, which is set to 172.5 GeV in the MC simulation. In some cases the b -quarks are produced at an angle such that only the W boson is actually found within $\Delta R < 1.0$ from the jet axis. Thus, a secondary peak appears around the W boson mass, set to 80.4 GeV in the MC simulation. The total jet mass distribution has two peaks. As can be seen from Fig. 3, the agreement between the MADGRAPH5+PYTHIA8MC and the GAN output is in fact satisfactory.

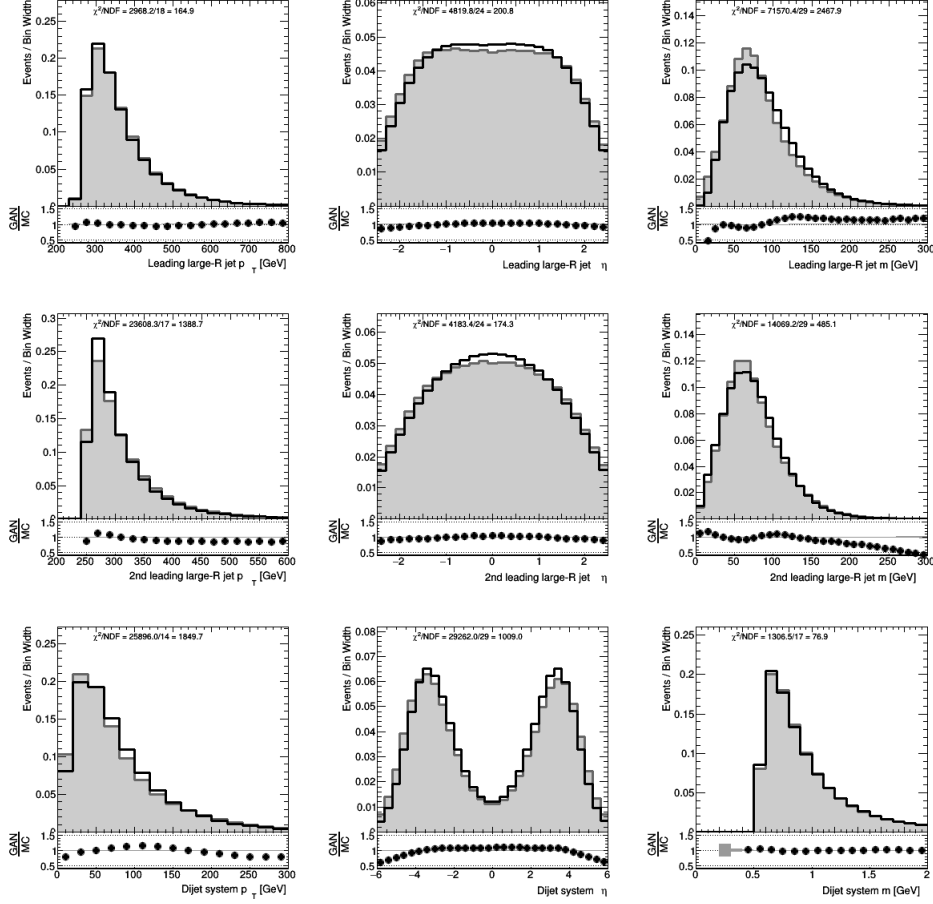


Figure 2: Comparison of kinematic observables with respect to reco-level (MADGRAPH5+PYTHIA8+DELPHES3) Monte Carlo simulation. The gray area represents the MC prediction, and the black line indicates the GAN output.

We further investigated the agreement in regions of the phase space with low cross-section, in particular where the dijet invariant mass is in the multi-TeV regime. This kinematic region is of particular interest for searches of physics beyond the SM. A very common approach is to fit the MC sample with the following four-parameters ($4p$) analytic function:

$$f(x) = \frac{p_0(1-x)^{p_1}}{x^{(p_2+p_3 \log x)}} \quad (6.1)$$

where $x = m_{jj}/\sqrt{s}$ and p_0 , p_1 , p_2 , p_3 are the free parameters of the fit. Such function is motivated by the structure of parton distribution functions and has been widely used by Tevatron and LHC experiments [33].

We trained the GAN using only a small fraction of the available events, about 15% corresponding to about 150,000 events with $m_{jj} > 1.5$ TeV. Then, we used the trained model to generate a sample of about 11 million events, a number much larger than that of events in the MADGRAPH5+PYTHIA8 training sample. A difference between the two approaches can be seen for

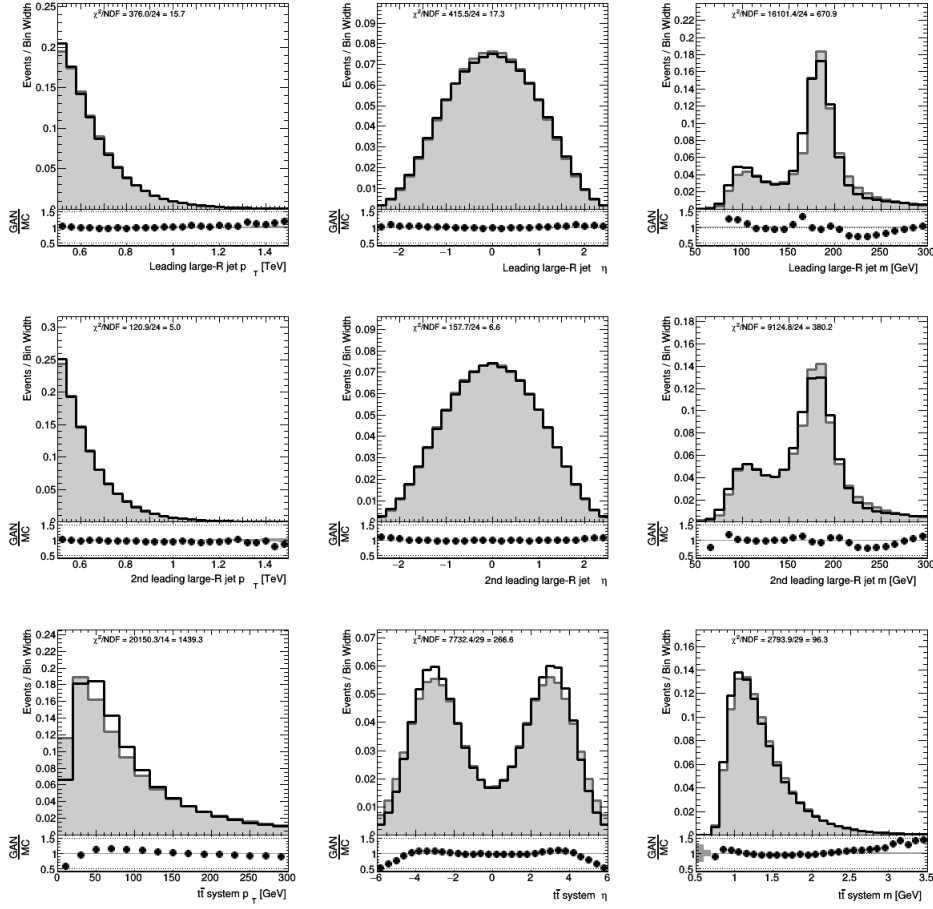


Figure 3: Comparison of kinematic observables for the all-hadronic $t\bar{t}$ production with respect to particle-level (MADGRAPH5+PYTHIA8) Monte Carlo simulation. The gray area represents the MC prediction, and the black line indicates the GAN output.

$m_{jj} > 8$ TeV, where the number of MC events is very small. The difference between the two methods can be interpreted as a source of systematic uncertainty.

As shown in Fig.4, limiting the fit in the region between 2.5 and 10 TeV, the $4p$ analytic function can predict the shape of the MC distribution with a $\chi^2/\text{NDF} = 0.98$. In the same kinematic region, the sample generated with the GAN shows an agreement with comparable χ^2/NDF . Besides the agreement in a single variable, one has to take in mind that the $4p$ fit does not allow the user to generate an event, but only to make an estimate of the background due to multijet production in that particular kinematic region and only for that specific observable. Therefore, events produced with our GAN can significantly expand the techniques used by analysis teams in determining their background.

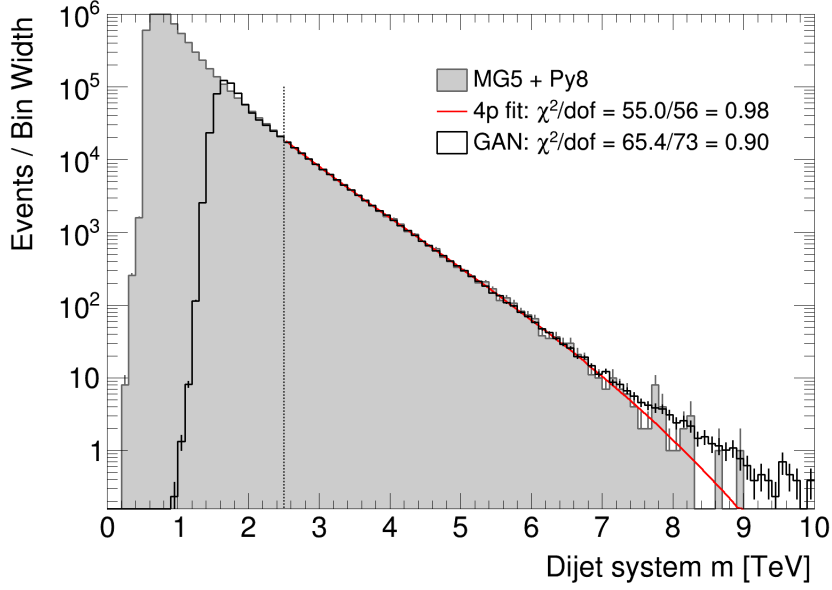


Figure 4: Comparison between a (MADGRAPH5 + PYTHIA8) Monte Carlo simulation sample and the GAN extrapolation to high dijet invariant mass. The gray area represents the MC prediction, the black line indicates the GAN output, and the red line is the fitted four-parameters analytic function.

7. Conclusions and Outlook

The Generative-Adversarial Network presented in this paper provides a novel and attractive solution to reduce the usage of CPU and possibly disk space to generate and simulate events at the LHC experiments. While still in its infancy, this method provides a unique opportunity to improve the quality of the MC used by the LHC collaborations as they will be able to use generators that are currently too time consuming to use. In the future, it should be possible to generalize this approach to more complicated processes such as top-quark pair or vector boson production in association with jets; the best MC predictions of these processes are also limited by high CPU requirements. Our results comparing simulated events show that our GAN can reproduce simulated events with high accuracy. This proof-of-concept shows the potential of these tools to provide an efficient solution to the large number of simulated events required by the ambitious physics programme of the LHC experiments. Future work will also focus on more advanced methods to further stabilize the training and avoid model collapse, while still being able to fit the relevant kinematic distributions in regions of the phase-space with low cross-sections.

References

- [1] Generative Adversarial Networks, I. J. Goodfellow, et al., Proceedings of NIPS 2014, arXiv:1406.2661
- [2] A Style-Based Generator Architecture for Generative Adversarial Networks, T. Karras et al., arXiv:1812.04948

- [3] C-RNN-GAN: Continuous recurrent neural networks with adversarial training, O. Mogren, arXiv:1611.09904, Accepted to Constructive Machine Learning Workshop (CML) at NIPS 2016
- [4] Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis, L. de Oliveira, M. Paganini, B. Nachman, *Comput Softw Big Sci* (2017) 1: 4
- [5] Accelerating Science with Generative Adversarial Networks: An Application to 3D Particle Showers in Multi-Layer Calorimeters, Michela Paganini, Luke de Oliveira, and Benjamin Nachman, *Phys. Rev. Lett.* 120, 042003 (2018)
- [6] CaloGAN: Simulating 3D high energy particle showers in multilayer electromagnetic calorimeters with generative adversarial networks, Michela Paganini, Luke de Oliveira, and Benjamin Nachman, *Phys. Rev. D* 97, 014021
- [7] Auto-Encoding Variational Bayes, D. P. Kingma and M. Welling, arXiv:1312.6114
- [8] Jet Constituents for Deep Neural Network Based Top Quark Tagging, J. Pearkes et al., arXiv:1704.02124
- [9] Jet-Images – Deep Learning Edition, L. de Oliveira et al., *JHEP* 07 (2016) 069
- [10] Boosted Jet Tagging with Jet-Images and Deep Neural Networks, M. Kagan et al., *EPJ Web of Conferences* 127, 00009 (2016)
- [11] The ATLAS experiment at the CERN Large Hadron Collider, ATLAS Collaboration, *J Instrum*, 2008, 3: S08003
- [12] The CMS experiment at the CERN LHC, The CMS Collaboration, *J Instrum*, 2008, 3: S08004
- [13] MadGraph 5 : Going Beyond, J. Allwall et al., *JHEP* 06 (2011) 128
- [14] A general framework for implementing NLO calculations in shower Monte Carlo programs: the POWHEG BOX, S. Alioli, P. Nason, C. Olearic, E. Re, *JHEP* 06 (2010) 043
- [15] Matching NLO QCD computations and parton shower simulations, S. Frixione, B.R. Webber, *JHEP* 06 (2002) 029
- [16] A Brief Introduction to PYTHIA 8.1, T. Sjöstrand, S. Mrenna and P. Z. Skands, *Comput. Phys. Commun.* 178 (2008) 852
- [17] Herwig 7.0/Herwig++ 3.0 release note, J. Bellm et al., *Eur. Phys. J. C* 76 (2016) 196
- [18] GEANT4 Collaboration, *Nuclear Instruments and Methods in Physics Research A* 506, 250 (2003).
- [19] A Roadmap for HEP Software and Computing R&D for the 2020s, The HEP software foundation, arXiv:1712.06982v5
- [20] Towards Jetography, G. Salam, *Eur. Phys. J. C* 67 (2010) 637
- [21] Measurements of $t\bar{t}$ differential cross-sections of highly boosted top quarks decaying to all-hadronic final states in pp collisions at $\sqrt{s} = 13$ TeV using the ATLAS detector, The ATLAS Collaboration, *Phys. Rev. D* 98, 012003 (2018)
- [22] Search for supersymmetry in the all-hadronic final state using top quark tagging in pp collisions at $\sqrt{s} = 13$ TeV, The CMS Collaboration, *Phys. Rev. D* 96, 012004 (2017)
- [23] Search for new resonances decaying into boosted W, Z and H bosons at CMS, M. Krohn and C. Vernieri, FERMILAB-CONF-17-429-PPD, arXiv:1710.02217

- [24] DELPHES 3: a modular framework for fast simulation of a generic collider experiment, J. de Favereau et al., J. High Energ. Phys. (2014) 2014: 57
- [25] The anti- k_T jet clustering algorithm, Matteo Cacciari, Gavin P. Salam and Gregory Soyez, Journal of High Energy Physics, Volume 2008, JHEP04(2008)
- [26] FastJet: a code for fast k_T clustering, and more, M. Cacciari, arxiv:hep-ph/0607071
- [27] Empirical Evaluation of Rectified Activations in Convolutional Network, Bing Xu et al., arXiv:1505.00853
- [28] F. Chollet, (2015) Keras, GitHub. <https://github.com/fchollet/keras>
- [29] TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems, M. Abadi et al. (2015), <https://www.tensorflow.org>
- [30] Scikit-learn: Machine Learning in Python, F. Pedregosa et al., Journal of Machine Learning Research vol. 12 (2011)
- [31] Data Structures for Statistical Computing in Python, Wes McKinney, Proceedings of the 9th Python in Science Conference (2010)
- [32] Adam: A Method for Stochastic Optimization, D. P. Kingma and J. Ba, arXiv:1412.6980, conference paper at the 3rd International Conference for Learning Representations, San Diego, 2015
- [33] Searches for Dijet Resonances at Hadron Colliders, R. H. Harris and K. Kousouris, Int. J. Mod. Phys. A26:5005-5055, 2011