## Modeling Hadronization using Machine Learning

Phil Ilten <sup>1†</sup>, Tony Menzo <sup>1\*</sup>, Ahmed Youssef<sup>1‡</sup>, and Jure Zupan<sup>1§</sup>

Department of Physics, University of Cincinnati, Cincinnati, Ohio 45221,USA
 <sup>†</sup> philten@cern.ch, \* menzoad@mail.uc.edu, <sup>‡</sup> youssead@ucmail.uc.edu, <sup>§</sup> zupanje@ucmail.uc.edu,

March 1, 2023

## Abstract

We present the first steps in the development of a new class of hadronization models utilizing machine learning techniques. We successfully implement, validate, and train a conditional sliced-Wasserstein autoencoder to replicate the PYTHIA generated kinematic distributions of first-hadron emissions, when the Lund string model of hadronization implemented in PYTHIA is restricted to the emissions of pions only. The trained models are then used to generate the full hadronization chains, with an IR cutoff energy imposed externally. The hadron multiplicities and cumulative kinematic distributions are shown to match the PYTHIA generated ones. We also discuss possible future generalizations of our results.

# Contents

1	Introduction			
<b>2</b>	Conditional SWAEs and hadronization			
	2.1 The simplified Lund string hadronization model	3		
	2.2 The cSWAE architecture	6		
	2.3 Training	9		
3	Reproducing the simplified Pythia fragmentation model	10		
	3.1 First emission trained models	11		
	3.2 Labels and $E$ dependent distributions	13		
	3.3 Hadronization chain	15		
4	Conclusion and Outlook	16		
A	A Public code MLhad_v0.1			
в	B Sliced Wasserstein distance			
С	Latent distributions	20		
References				

### 1 Introduction

A typical particle physics Monte Carlo event generator factorizes into three distinct steps or blocks of code: (i) the generation of the hard process, (ii) parton shower, and (iii) hadronization (including color reconnections). The first two steps are perturbative in their nature, and thus under good theoretical control, with significant efforts devoted to improving the precision even further [1–4]]. The algorithmic challenges are efficient sampling of final state particle configurations, and taming the factorial growth of the calculations with the increasing number of particles. The simulation of the hard matrix element is performed either by a specialized code, e.g., MADGRAPH [5], which only calculates the hard process, or is directly included in complete event generators, such as PYTHIA [6], HERWIG [7], or SHERPA [8], that also perform the parton showering.

In contradistinction, the hadronization step is inherently non-perturbative. One is therefore forced to resort to phenomenological models inspired by non-perturbative discriptions such as lattice QCD. The two main models used in simulating hadronization are the Lund string model [9–11] and cluster model [12–14]. In the string model, quark-antiquark pairs are thought of being connected by a string, a flux tube of the strong force confined in the lateral direction. As the quark-anti-quark pair moves apart, the string breaks, creating new quark-anti-quark pairs in the process, resulting in the emission of hadrons. These emissions are performed iteratively, breaking the string either from the left or the right side, with the final step modified *post hoc* in order to provide an emission similar to the previous steps. This model requires extra parameters to describe the hadrons' transverse momenta and heavy particle suppression, and has some challenges describing baryon production. Over  $\mathcal{O}(20)$  parameters are required by the string model to describe the hadronization.

In the cluster model, gluons are forced to split into quark–anti-quark pairs at longer distances (lower energy). All quark–anti-quark pairs are grouped into color singlet combinations with a distance scale that depends only on the evolution step, and not the hard process step of the Monte Carlo even generation. Hadrons are emitted from these universally pre-confined clusters via a series of two-body decays until only physical hadrons remain. The model has fewer parameters and naturally generates hadron transverse momenta. However, the decays of massive clusters lead to phenomenological problems such as predicting heavy baryon distributions which do not match data well.

Machine Learning (ML) techniques offer the possibility to build alternatives to the above two models of hadronization. Such ML models could be directly built from data and provide insights into the current phenomenological models. While ML techniques have recently entered into the development of event generators, through adaptive integration [15–20], [parton showers [21–30], ML based fast detector or event simulations [31–57]], and model parameter tuning [58,59], the application of ML to the problem of hadronization as the final step in the Monte Carlo pipeline is entirely new, to the best of our knowledge. The present manuscript represents [the first step toward building a full-fledged ML based hadronization framework.]

In principle, Generative Adversarial Networks (GANs) [60], Variational Auto-Encoders (VAEs) [61] [ and Normalizing Flows (NF) [62]] have demonstrated the ability for ML to generate convincing physical observables. [In addition, conditional generative models give more flexibility and control of the output [63, 64]] [Extending the ML techniques for hadronization faces] [three challenges]: (i) producing sets of physical observables that vary in size (unlike a fixed number of pixels), ranging from  $\mathcal{O}(1)$  to  $\mathcal{O}(10^4)$ ; (ii) strictly conserving certain physical quantities, e.g., momentum and energy; and (iii) learning from limited training sets which only provide coarse-grain detail. In this paper we present

an architecture based on conditional sliced-Wasserstein autoencoders (cSWAE) [65, 66], that overcomes the above challenges. The resulting code, MLHAD, is publicly available, see Appendix A. We demonstrate the capabilities of MLHAD by training it on specially prepared PYTHIA hadronization outputs with an explicit IR cut-off. To speed up the training we perform a transformation that captures the bulk of the energy dependence of the PYTHIA hadronization output. However, we also show that, if this transformation is not performed, the cSWAE can still reproduce the energy dependence and thus should be able to reproduce any additional energy dependence that may be present in the hadronization process realized in nature. We expect that the first version of the cSWAE architecture presented here can be upgraded to eventually be trained directly on data [ (details about further steps to achieve this can be found in section 4)].

The paper is structured as follows. In Section 2 we introduce conditional sliced-Wasserstein autoencoders and describe how these can be used to reproduce the Lund string model of hadronization. In Section 3 we then compare the trained MLHAD models to the results of a simplified PYTHIA hadronization model. Section 4 contains our conclusions and a brief discussion of future directions. Appendix A contains details about the publicly accessible MLHAD code, while Appendix B gives further details on the sliced-Wasserstein distance.

## 2 Conditional SWAEs and hadronization

### 2.1 The simplified Lund string hadronization model

As the first step toward building a machine learning (ML) based simulator of hadronization, we create a ML architecture that is able to reproduce a somewhat simplified Lund string model for hadronization. [Hadronization is the last step in the Monte Carlo simulation of the particle collision, and describes the creation of hadrons from quarks and gluons, a process that occurs at the nonperturbative scale of a few 100 MeV. The distributions of quarks and gluons at low scales is obtained using a parton shower simulation, which describes the emission of particles between the hard scale of the collisions, typically a few 100 GeV, down to low energies. In a Lund string model the quarks and gluons are thought of being connected by QCD [color flux tubes, or strings,] that carry significant amounts of energy, and shed it in the process of hadron creation. While there were already attempts to use ML to improve parton shower simulations [28, 67-73], this manuscript represents the first attempt to use ML for hadronization. In both cases the physics is described by a Markov chain, however, for different reasons. The semi-classical evolution of a parton shower, where gluons and quarks are radiated in a Markov chain, can be justified in the small angle emission limit. The hadronization, on the other hand, can be represented as a Markov chain process because string fragmentations occur at causally disconnected points.

The physical process we want to describe is depicted in Fig. 1. It shows a  $q_i \bar{q}_i$  fragmentation event in the center-of-mass frame in which the individual partons, each with flavor index *i* and initial energy *E*, travel with equal and opposite momenta and are connected via a QCD string. String breaking produces a composite hadron  $h \sim q_i \bar{q}_j$  and a new  $q_j \bar{q}_i$ -string system depicted in the lower part of Fig. 1.<sup>1</sup> The hadron *h* is ejected with some energy and momentum  $(E_h, \vec{p}_h)$ , while the new string system has the energy and momentum  $(2E - E_h, -\vec{p}_h)$ , so that the total energy and momentum are conserved.

<sup>&</sup>lt;sup>1</sup>The depiction in Fig. 1 is for a string breaking occurring on the quark side. The string breaking on the anti-quark side produces similarly a hadron with quark composition  $h \sim q_j \bar{q}_i$ , and the new  $q_i \bar{q}_j$ -string.



**Figure 1:** [Schematic of a single fragmentation event, for an initial quark–anti-quark pair,  $q_i \bar{q}_i$ , into a hadron with quarks  $q_i \bar{q}_j$  and new endpoints  $\bar{q}_i q_j$ .]

[The goal of our ML framework will be do properly describe the probabilities of emitting a hadron of given energy and momentum.]

After boosting to the center-of-mass frame of the new string, one has essentially the same initial state, a quark-anti-quark pair going back to back connected by a string, but with reduced energy E' and a different quark flavor composition. Such fragmentation events stack one after the other and form a fragmentation chain, one hadron emission at a time, until the entire energy of the initial two-parton system (2E) is converted into hadrons. The end of the string used for each splitting is chosen at random. Until relatively low string energies of a few GeV, the selection of flavor and the kinematics of the hadron emission are taken to be independent processes. In the final stages of hadronization, when the string energy is close to the nonperturbative scale, the two processes, on the other hand, become intertwined. To simplify the problem, we therefore terminate fragmentation events at a center-of-mass string energy  $E_{\rm cut} = 5$  GeV. We also consider a simplified string system which allows for u and d quarks as string ends, as well as their respective anti-quarks, and pions as final states.

Note that each step in the above hadronization chain is independent from the previous one. A successful hadronization simulator therefore takes as the input the string energy E (i.e., the energy of one of the endpoint quarks in the center-of-mass frame) as well as its flavor composition, and gives the flavor and kinematics of the hadron after first emission,  $(E_h, \vec{p}_h)$ . Repeating the first emission generates the full hadronization chain. Since  $E_h^2 = \vec{p}_h^2 + m_h^2$ , where  $m_h$  is the hadron mass, the kinematics of the emission are fully described by specifying  $\vec{p}_h$  and flavor of the created hadron h. We orient the coordinate system such that the z axis is along the direction of the initial string, while the x and ycoordinates are perpendicular to it. The transverse components of the  $\vec{p}_h$  vector are given by

$$p_x = p_T \cos \varphi, \quad p_y = p_T \sin \varphi,$$
 (1)

where  $p_T \equiv \sqrt{p_x^2 + p_y^2}$  and  $\varphi$  is the polar angle. The string breaking and hadron emission are assumed to be axially symmetric in PYTHIA, i.e., independent of  $\varphi$ , and thus the problem of simulating the hadronization event reduces to a two variable problem of generating the  $p_z$  and  $p_T$  distributions for the first emission.

A special feature of the hadronization event and the chosen kinematic variables is the ability to render the  $p_z$  kinematic distributions independent of the initial parton energy, E, through a simple rescaling transformation

$$p'_z \equiv E_{\rm ref} \frac{p}{E},\tag{2}$$

where E is the energy of the quark in the center of mass for the initial string, and  $E_{\rm ref}$ 



Figure 2: The  $p_z$  distributions (left) and the rescaled  $p'_z$ , Eq. (2), distributions (right) from PYTHIA hadronization events for the first-hadron emission with initial parton energies E = 10,100,1000 GeV shown with blue, red, and green solid lines, respectively.

is a conveniently chosen reference energy that renders p' dimensionful. In the rest of the paper we set  $E_{\rm ref} = 50$  GeV. The transformation of the  $p_z$  distribution with respect to the initial parton energy E can be seen in Fig. 2.

The fragmentation process implemented in PYTHIA is constructed in momentum space as an iterative walk through production vertices. To do so a stochastic variable termed the longitudinal momentum fraction z is defined, describing the fraction of longitudinal momentum taken away by the emitted hadron.<sup>2</sup> Given the longitudinal momentum fraction,  $p_z$  can straight-forwardly be obtained via the relation  $z = (p_z + E_h)/2E$  where 2E is the total energy of the initial fragmenting system. The probability distribution f(z) from which z is sampled is called the Lund left-right symmetric scaling function (also Lund sampling or fragmentation function) and is given by

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(-b\frac{m_{h,T}^2}{z}\right),\tag{3}$$

where  $m_{h,T}^2 \equiv m_h^2 + p_T^2$  is the transverse mass, and the normalization prefactor is omitted for clarity. The phenomenological parameters a, b are chosen to match experimental data. The  $p_T^2$  term in the transverse mass squared,  $m_{h,T}^2$ , captures the tunneling probability for a string breaking to occur away from the classical position of the string end, such that the additional energy required for the transverse momentum kick can be released from the string. It leads to a correlation between transverse and longitudinal distributions of hadron momenta (in the center-of-mass frame of the initial string), i.e., the average value of z increases with increasing  $p_T$ . In the default implementation of the Lund model in PYTHIA, the hadron  $p_T$  distribution is assumed to be Gaussian distributed, with average  $\langle \vec{p}_T \rangle = 0$ , and a width  $\sigma_0 \sim \mathcal{O}(300 \text{ MeV})$ , reflecting that its origin is an inherently quantum process occurring at the nonperturbative QCD scale.<sup>3</sup>

The above basic setup of the Lund model becomes more involved when full complexity of the experimental data needs to be explained. Most of the  $\mathcal{O}(20)$  parameters that give more flexibility to the PYTHIA implementation of the Lund string model are related to

<sup>&</sup>lt;sup>2</sup>In Section 2.1,  $z_i$  denote the latent-space variables. Despite similarity in notation there is no relation between the two variables.

<sup>&</sup>lt;sup>3</sup>The configurable PYTHIA parameter name is StringPT:sigma.

the differences in hadronizations of the light quarks compared to the strange, c and bquarks. For instance, each quark flavor can in principle have a different a; in PYTHIA strange quarks are allowed to have different values of a than for u and d quarks, while for heavier c and b quarks the Lund fragmentation is also allowed to be multiplied by an extra z-dependent factor with new flavor-dependent parameters. Similarly, the  $p_T$ distributions can deviate from the Gaussian form. While this gives quite some flexibility to the hadronization model, it does have its own drawbacks. On one hand, the number of parameters to be tuned to data is already quite large. On the other hand, one may worry that the analytic form of the scaling function in Eq. (3), while well motivated, is not flexible enough, with higher order corrections in z potentially becoming important, e.g., at low string energies. Generative ML models, such as the architecture that we introduce in the next section, can be used as effective tools to address both of these issues. For the purposes of this paper we will not yet train our ML architecture on the physics data, but rather on the synthetic data generated by PYTHIA. However, we anticipate that the expressibility of the ML framework, which we demonstrate below, will allow for a better description of the physics data sensitive to hadronization than the Lund left-right symmetric scaling function in Eq. (3) does right now.]

### 2.2 The cSWAE architecture

The ML model of hadronization used here is based on the conditional sliced-Wasserstein Autoencoder (cSWAE) [65, 66] (for an example of a use of SWAE architecture in particle physics simulations see [40]). The motivation for using cSWAE is two-fold, i) the flexibility of being able to use a wide variety of latent-space distributions and thus optimize the performance of the hadronization model, and ii) the ability to incorporate the energy dependence of hadronization through a two dimensional condition vector c. We expect the second feature to become most relevant once MLHAD is trained on experimental data, for which small breakings of the energy independence exhibited by the Monte Carlo generated  $p'_z$  data, Fig. 2, may be anticipated. [The main advantage of SWAEs over VAEs is the flexibility in the choice of the latent space distribution, which allows the user to chose any sampleable distribution as latent space distribution.] [This is achieved by introducing a sliced Wasserstein distance (i.e. an approximate of the real Wasserstein distance between the desired and the obtained latent space distributions) in the cost function, see Eq. (6) below. This is then added to the usual reconstruction loss estimate term in the cost function, see Eq. (5) below.]

The schematic of the cSWAE architecture is given in Fig. 3. It has two parts, the encoder and the decoder:

The encoder  $\phi$  takes as inputs the data vectors  $\mathbf{x}_i$  and labels  $\mathbf{c}_i$  and returns a latentspace vector  $\tilde{\mathbf{z}}_i = \phi(\mathbf{x}_i, \mathbf{c}_i)$ . Depending on the value of  $\mathbf{c}_i$  the encoder will transform  $\mathbf{x}_i$ to different regions in the latent space, as shown in the graphical representation of Fig. 4. The dimension of the latent space,  $d_z$ , needed for the application to hadronization is anywhere from  $d_z = 2$  to  $d_z = 30$ , see also Table 1. The latent-space vectors  $\tilde{\mathbf{z}}_i$  are trained to be distributed according to the target latent-space distribution,  $\tilde{\mathbf{z}}_i \sim I(\tilde{\mathbf{z}}_i, \mathbf{c}_i)$ , which is ensured through the use of sliced-Wasserstein distance,  $SW_p$ , in the loss function. In particular, the latent-space variable  $\tilde{\mathbf{z}}_i$  need not be normally distributed. We found that this feature translated to significant improvements in the performance of MLHAD. With cSWAE one can choose a custom probability distribution such that the encoding of the information about the first emission hadron kinematics leads to optimal results. This is the main practical difference between cSWAE and the conditional Variational Autoencoder (cVAE). The cVAE use KL-divergence in the loss function, which typically require



Figure 3: [The cSWAE architecture for simulating hadronization. Inputs  $x_i$  have condition  $c_i$ , which parametrizes the string energy. The decoder takes  $\tilde{z}_i$  as inputs and generates the predicted hadron kinematics  $\tilde{x}_i = \{\tilde{p}_{z,k}^{(i)}\}$ . The sliced-Wasserstein-distance loss function,  $\mathcal{L}_{SW}$ , constraints the latent-space vectors  $\tilde{z}_i$  to the target distribution  $\tilde{z}_i \sim I(\tilde{z}_i, c_i)$ . The reconstruction loss function,  $\mathcal{L}_{rec}$ , minimizes the difference between  $x_i$  and  $\tilde{x}_i$ .]



Figure 4: [Illustration of the conditional vector  $c_i = c(E_i)$  mapping the input data  $x_i$  into different regions of the latent space,  $\tilde{z}$ .]

that the latent-space variables follow simple distributions, such as a normal distribution. The cSWAE uses instead the sliced-Wasserstein distance,  $SW_p$ , see Appendix B for more details. This gives the architecture significantly more flexibility, as one can choose the latent-space distributions to follow almost any distribution, as long as it is sampleable (in particular, the analytic form of  $I(z, c_i)$  is not required to exist).

The decoder  $\psi$  takes as inputs the condition vector  $\mathbf{c}_i$  and the latent-space vector  $\tilde{\mathbf{z}}_i$ . It returns the reconstructed hadron kinematics  $\tilde{\mathbf{x}}_i = \psi(\phi(\mathbf{x}_i, \mathbf{c}_i))$ , where  $\tilde{\mathbf{x}}_i$  is the  $N_e$  dimensional vector consisting of sorted kinematic variables, either  $p'_{z,k}^{(i)}$  or  $p_{T,k}^{(i)}$ . Through the minimization of the loss function [65]

$$\mathcal{L}(\psi, \phi) = \mathcal{L}_{\rm rec} + \mathcal{L}_{\rm SW},\tag{4}$$

where

$$\mathcal{L}_{\text{rec}} = \frac{1}{N_{\text{tr}}} \sum_{i=1}^{N_{\text{tr}}} \left[ \frac{1}{Q} d_2^2(\boldsymbol{x}_i, \boldsymbol{\psi}(\boldsymbol{\phi}(\boldsymbol{x}_i, \boldsymbol{c}_i))) + d_1(\boldsymbol{x}_i, \boldsymbol{\psi}(\boldsymbol{\phi}(\boldsymbol{x}_i, \boldsymbol{c}_i))) \right],$$
(5)

$$\mathcal{L}_{SW} = \frac{\lambda}{LN_{tr}} \sum_{\ell=1}^{L} \sum_{i=1}^{N_{tr}} d_{SW}(\boldsymbol{\theta}_{\ell} \cdot \boldsymbol{z}_{[i]_{\ell}}, \boldsymbol{\theta}_{\ell} \cdot \boldsymbol{\phi}(\boldsymbol{x}_{[i]_{\ell}}, \boldsymbol{c}_{i})),$$
(6)

with  $z_i \sim I(z_i, c_i)$ , the training attempts to reproduce the training data distribution  $x_i$  with the generated data distribution  $\tilde{x}_i$ , while the latent-space vectors  $\tilde{z}_i$  follow the desired target distribution  $\tilde{z}_i \sim I(\tilde{z}_i, c_i)$ . The reconstruction loss  $\mathcal{L}_{\text{rec}}$  is a measure of the differences between the input,  $x_i$ , and generated kinematic vectors,  $\tilde{x}_i$ . It is the sum of



Stopping condition :  $E_i < E_{cut}$ 

Figure 5: [Illustration of MLHAD generating hadronization chains. Random variables  $z_i$  are passed through the decoder D with condition vector  $c_i$  to generate the hadron momentum, given the string energy  $E_i$ . A modified PYTHIA flavor selector FS, generates the new string flavor,  $s_{i+1}$ , and emitted hadron species,  $h_i$ . Before each emission, the string is boosted to its center-of-mass frame using a Lorentz transformation  $\Lambda$ .]

two terms for each of the 1D distributions that we consider,

$$d_2^2(\boldsymbol{x}_i, \boldsymbol{\psi}(\boldsymbol{\phi}(\boldsymbol{x}_i, \boldsymbol{c}_i))) = \begin{cases} \sum_k \left( p_{z,k}^{\prime(i)} - \tilde{p}_{z,k}^{\prime(i)} \right)^2, & \text{for } p_z' \text{ distributions,} \\ \sum_k \left( p_{T,k}^{(i)} - \tilde{p}_{T,k}^{(i)} \right)^2, & \text{for } p_T \text{ distributions,} \end{cases}$$
(7)

$$d_1(\boldsymbol{x}_i, \boldsymbol{\psi}(\boldsymbol{\phi}(\boldsymbol{x}_i, \boldsymbol{c}_i))) = \begin{cases} \sum_k |p_{z,k}^{\prime(i)} - \tilde{p}_{z,k}^{\prime(i)}|, & \text{for } p_z' \text{ distributions,} \\ \sum_k |p_{T,k}^{(i)} - \tilde{p}_{T,k}^{(i)}|, & \text{for } p_T \text{ distributions,} \end{cases}$$
(8)

where  $p_{z,k}^{\prime(i)}$  and  $p_{T,k}^{(i)}$  are the components of the training-dataset vectors  $\boldsymbol{x}_i$ , while  $\tilde{p}_{z,k}^{\prime(i)}$  and  $\tilde{p}_{T,k}^{(i)}$  are the components of the output vectors  $\tilde{\boldsymbol{x}}_i$ . For the relative weight between the two terms in  $\mathcal{L}_{\text{rec}}$  we take Q = 1 GeV. The two contributions of  $\mathcal{L}_{\text{rec}}$  are sensitive to distinct scales allowing for fast convergence  $(d_1)$  and continual improvement  $(d_2)$  throughout training while also heavily penalizing outliers.

The second term in Eq. (4),  $\mathcal{L}_{SW}$ , is the implementation of the sliced-Wasserstein distance  $SW_1$  between the distribution of latent-space vectors  $\tilde{z}_i$  created by the encoder, and the target latent-space distribution  $I(z_i, c_i)$ . [The sliced-Wasserstein distance is the approximation of the true Wasserstein distance between the two distributions, and is smaller the closer the latent space distribution is to the desired one. The sliced-Wasserstein distance approximation becomes better and better the higher the number of 1D slices (or probes) of the distributions one uses. The advantage is that the computation of Wasserstein distances for 1D slices can be done very efficiently, leading to a significant speed up of the algorithm.]

[The computation of  $SW_1$  is done as follows.] The vectors  $z_i$  in Eq. (6) are randomly drawn from this target distribution,  $z_i \sim I(z, c_i)$ . The scalar products with the unit vectors  $\theta_l$ , defining the L slices, give the one dimensional projections of the latent-space distributions, for which the Wasserstein distances,  $W_1$ , are straightforward to compute. They are given simply by the average sum of the distances between the sorted data points, see Appendix B for further details. Note that for one dimensional latent space  $SW_1 = W_1$ , and in the sum in Eq. (4) one can set L = 1.

### 2.3 Training

The input data to the encoder are  $N_e$  PYTHIA generated first-hadron emissions for a fixed initial string energy  $E_i = 50$  GeV. In all of the numerical examples below we take  $N_e = 100$ , so that the input is an  $N_e$  dimensional vector  $\boldsymbol{x}_i$  of either  $p_{z,k}^{(i)}$  or  $p_{T,k}^{(i)}$ ,  $k = 1, \ldots, N_e$ . That is, in this manuscript we apply cSWAE to the case where the  $p_z$ and  $p_T$  distributions are uncorrelated and treat each of them separately. However, the architecture is flexible enough that correlated 2D or higher dimensional distributions could also be used as inputs.

The elements of the input vectors  $\boldsymbol{x}_i$  are sorted, i.e.,  $p'_{z,1}^{(i)} \leq p'_{z,2}^{(i)} \leq \cdots \leq p'_{z,N_e}^{(i)}$ (and similarly for  $p_{T,k}^{(i)}$ ).<sup>4</sup> The training dataset consists of  $N_{\text{tr}}$  such  $\boldsymbol{x}_i$  input vectors,  $i = 1, \ldots, N_{\text{tr}}$ , and  $N_{\text{val}} \boldsymbol{y}_j$  validation vectors,  $j = 1, \ldots, N_{\text{val}}$ , where typically  $N_{\text{tr}}$  is taken to be  $N_{\text{tr}} = \mathcal{O}(4000)$  and  $N_{\text{val}}$  an order of magnitude smaller. To summarize, the training and validation datasets are created by generating  $N \equiv N_e(N_{\text{tr}} + N_{\text{val}}) = 4 \times 10^5$  PYTHIA first hadron emission events. The emission data ( $p_z$  or  $p_T$ ) is then partitioned randomly into  $N_{\text{tr}} + N_{\text{val}}$  vectors of length  $N_e = 100$ . Finally, the elements in each vector are sorted from least to greatest.

The string energy  $E_i$ , or equivalently mass in the center-of-mass frame, is converted to a unit condition vector  $\mathbf{c}_i = (\bar{c}_i, 1 - \bar{c}_i)$  with  $\bar{c}_i \in [0, 1]$  a floating point number such that

$$E_i = E_{\min}\bar{c}_i + E_{\max}\left(1 - \bar{c}_i\right), \quad \text{and thus} \quad \bar{c}_i = \frac{E_{\max} - E_i}{E_{\max} - E_{\min}}, \quad (9)$$

where  $E_{\min}$  and  $E_{\max}$  are the reference minimal and maximal energies. A good choice for  $E_{\max}$  is the maximal partonic collision energy in the simulation, while  $E_{\min}$  can be taken to be the IR cutoff  $E_{\text{cut}}$ .

In general, the cSWAE allows for the initial string energy  $E_i$  of each  $x_i$  to be different (but the same for all the  $N_e$  components of  $x_i$ ). For the PYTHIA generated events the kinematic variable  $p_z$  can be made E independent through the transformation in Eq. (2) and thus  $E_i$  can be set to a constant value,  $E_i = 50$  GeV. As a proof of principle we also show in Section 3.2 that cSWAE models can be trained on E-dependent  $x_i$ .

The algorithm for training the cSWAE is as follows. Applying the encoder to the input data sample  $\{\boldsymbol{x}_1, ..., \boldsymbol{x}_{N_{tr}}\}$  gives the latent-space vectors  $\{\tilde{\boldsymbol{z}}_1, ..., \tilde{\boldsymbol{z}}_{N_{tr}}\}$ . To compute the sliced-Wasserstein distance term, Eq. (6), the unit vectors  $\{\boldsymbol{\theta}_1, ..., \boldsymbol{\theta}_L\}$  are randomly sampled from the  $(d_z-1)$ -dimensional unit sphere  $S^{d_z-1}$ , while the  $N_{tr}$  latent-space vectors  $\{\boldsymbol{z}_1, ..., \boldsymbol{z}_{N_{tr}}\}$  are sampled from the target distribution,  $\boldsymbol{z}_i \sim I(\boldsymbol{z}_i, \boldsymbol{c}_i)$ . For each  $\boldsymbol{\theta}_\ell$ , the scalar products  $\boldsymbol{\theta}_\ell \cdot \tilde{\boldsymbol{z}}_i = \boldsymbol{\theta}_l \cdot \boldsymbol{\phi}(\boldsymbol{x}_i)$  and  $\boldsymbol{\theta}_\ell \cdot \boldsymbol{z}_i$  are sorted in the following way. First the energy labels  $c_i$  (and the corresponding  $\tilde{z}_i, z_i$ ) are sorted into  $N_c$  bins of increasing  $c_i$  intervals with boundaries  $\bar{c}_{[1]} < \bar{c}_{[2]} < \cdots < \bar{c}_{[N_c]}$ . That is, the latent-space data are binned according to their energies,  $E_i$ , where the bins are chosen such that the distributions  $I(\boldsymbol{z}_i, \boldsymbol{c}_i)$  do not have large dependence on  $c_i$  within the bin. The generated and target  $I(\boldsymbol{z}_i, \boldsymbol{c}_i)$  distributions are then compared within each energy bin. This is achieved by first sorting the scalar products of  $\tilde{\boldsymbol{z}}_i$  and  $\boldsymbol{z}_i$  with  $\boldsymbol{\theta}_\ell$  within each  $c_i$  bin, and then combined into the lists  $\{\boldsymbol{\theta}_\ell \cdot \tilde{\boldsymbol{z}}_{[1]_\ell}, \ldots, \boldsymbol{\theta}_\ell \cdot \tilde{\boldsymbol{z}}_{[N_{tr}]_\ell}\}$  and  $\{\boldsymbol{\theta}_\ell \cdot \boldsymbol{z}_{[1]_\ell}, \ldots, \boldsymbol{\theta}_\ell \cdot \boldsymbol{z}_{[N_{tr}]_\ell}\}$ , respectively. The SW loss function  $\mathcal{L}_{SW}$  in Eq. (6) is then the average over the latent space distances between the two sorted lists,

$$d_{\rm SW}(\boldsymbol{\theta}_{\ell} \cdot \boldsymbol{z}_{[i]_{\ell}}, \boldsymbol{\theta}_{\ell} \cdot \boldsymbol{\phi}(\boldsymbol{x}_{[i]_{\ell}}) = \left| \boldsymbol{\theta}_{\ell} \cdot \boldsymbol{z}_{[i]_{\ell}} - \boldsymbol{\theta}_{\ell} \cdot \boldsymbol{\phi}(\boldsymbol{x}_{[i]_{\ell}}) \right|,\tag{10}$$

 $<sup>{}^{4}</sup>$ For 2D or higher dimensional problems the data would first be clustered in predefined 1D bins and then sorted within each bin.



Figure 6: [Illustration of Lorentz boosting ( $\Lambda$ ) from the lab frame to the string center-of-mass frame. Red and blue lines are the string system's longitudinal momentum with the total area equal string system's longitudinal momentum  $E + p_z$ . Each box is a new string.]

averaged also over all the L slices and multiplied by the relative weight prefactor  $\lambda$ . The final step in the algorithm is applying the decoder to  $\tilde{z}_i$ , which gives  $\{\tilde{x}_1, \ldots, \tilde{x}_{N_{\text{tr}}}\}$ . The distances between input dataset,  $\{x_1, \ldots, x_{N_{\text{tr}}}\}$ , and the generated sets  $\{\tilde{x}_1, \ldots, \tilde{x}_{N_{\text{tr}}}\}$ are then calculated using Eqs. (7) and (8), giving the reconstruction loss function  $\mathcal{L}_{\text{rec}}$ . Eq. (5). The decoder and encoder are updated in steps, trying to minimize the combined loss function, Eq. (4). Overfitting is avoided by monitoring the value of loss function when applied to the validation dataset, i.e., the loss function (4) with  $x_i \to y_i$ ,  $N_{\text{tr}} \to N_{\text{val}}$ .

Fig. 5 illustrates how the trained MLHAD decoder is used, along with the PYTHIA flavor selector, to generate the hadronization chain. Note, the full PYTHIA flavor selector is not needed here, but included to allow for subsequent development. The flavor selector takes as input the initial string flavor ID,  $s_i$ , and gives as the output the flavor ID of the emitted hadron,  $h_i$ , which also defines the flavor of the new string fragment,  $s_{i+1}$ . The MLHAD decoder takes as input the latent-space vector  $\mathbf{z}_i \sim I(\mathbf{z}_i, \mathbf{c}_i)$  sampled from the target distribution  $I(\mathbf{z}_i, \mathbf{c}_i)$ , where  $\mathbf{c}_i$  is the label encoding the center-of-mass energy of the string  $s_i$ , see Eq. (9). The MLHAD decoder returns the  $N_e$ -dimensional vector with a list of possible momenta for the emitted hadron,  $\tilde{p}_{z,k}^{\prime(i)}$  (or  $\tilde{p}_{T,k}^{(i)}$ ). We randomly choose one of these as the actual hadron kinematics, and modify accordingly the kinematics of the remaining string fragment,  $s_{i+1}$ , such that the energy and momentum are conserved. The emitted hadron is boosted to its rest frame, see Fig. 6. Its center-of-mass energy defines the label  $\mathbf{c}_{i+1}$  used as the input in the decoder for the next hadron emission. These steps are repeated until the string energy in its rest frame reaches the IR cutoff energy  $E_{\text{cut}}$ .

We have implemented the cSWAE architecture described above using PyTORCH [74]. The code can be accessed via a public repository, see Appendix A for details.

### 3 Reproducing the simplified Pythia fragmentation model

To demonstrate the viability and capability of the cSWAE based machine learning algorithm implemented in MLHAD, we reproduce the PYTHIA hadronization outputs. We analyze a  $q_i \bar{q}_i$  hadronization event in the center-of-mass frame in which the individual partons, each with flavor index *i* and initial energy *E*, travel with equal and opposite momenta producing a string between them. After the string breaks this produces a new string and the first emission hadron, see Section 2.1 for more details.

While MLHAD treats all the hadron emissions on an equal footing, PYTHIA treats the first emission slightly differently; in the first emission  $m_{T,h}$  in Eq. (3) is set to  $m_h$ 

Variable $\boldsymbol{x}$	Target $\boldsymbol{z}$	t (epochs)	$d_z$	λ	L
	Pythia	150	35	35	15
$p'_z$	Trapezoidal	300	2	20	30
	Triangular	150	2	30	25
	Pythia	100	20	30	30
$p_T$	Skew-norm	120	4	20	25
	Triangular	120	4	15	25

**Table 1:** The cSWAE training configurations, [where x is the input data, z the target latent-space distribution, t the number of epochs,  $d_z$  the dimension of the latent space,  $\lambda$  the regularization parameter of the sliced-Wasserstein loss, and L the number of latent space projections (slices).]

(i.e.,  $p_T = 0$ ), while for all subsequent emissions  $p_x$  and  $p_y$  are sampled from a normal distribution with a width  $\sigma_0$  (we set this tunable PYTHIA parameter to  $\sigma_0 = 0.335 \text{ GeV}$ ). Therefore, in training MLHAD we only aim to reproduce the PYTHIA output on average, which is in line with the physical limitations of the problem, since one cannot trace in nature each individual emission in the hadronization event.

Our model is trained on kinematic distributions for transformed variables,  $p'_z$ ,  $p_T$ , Eq. (2), obtained from the PYTHIA first emission events. With a uniformly sampled polar angle  $\varphi$  in the transverse plane, these kinematic variables then completely define the phase space of the system through Eqs. (1), (2). The MLHAD decoder is then used with a fixed shifted value transverse mass  $m_{T,h}^2 = m_h^2 + \sigma^2$ , with  $\sigma = \sigma_0/\sqrt{2}$ . This accounts for using only PYTHIA produced first emission data where  $p_T = 0$  GeV. For flavor selection we rely on PYTHIA's probabilistic model, and limit ourselves to light quarks, u, d and only pions as the final state hadrons.

The independence of the distributions from the initial parton energy, see Fig. 2, allows the cSWAE model to be trained on a dataset using an arbitrary initial parton energy,  $E_{\rm ref}$ , while the outputs of cSWAE hadronization generator can be transformed accordingly to obtain the distributions for any desired initial energy, E, using Eq. 2. While in the PYTHIA output the complete energy dependence is already captured with the simple rescaling in Eq. (2) we do not expect this to be entirely true for actual physical hadronization events realized in nature, for which subleading deviations from the scaling law in Eq. (2) may be anticipated. In Section 3.2 we demonstrate that such corrections to the scaling law can be captured by the cSWAE architecture.

#### 3.1 First emission trained models

The cSWAE trained models differ according to the target latent-space distribution, I(z, c), the dimension of the latent space  $d_z$ , training time t (epochs), the value of the sliced-Wasserstein regularization parameter  $\lambda$ , and the number of slices L, as shown in Table 1. In all the cases we fix the string energy to be E = 50 GeV. The first emissions for other string energies can be obtained by inverting the rescaling of the  $p'_z$  distributions in Eq. (2), while  $p_T$  distributions do not scale with E, although this is an assumption of the PYTHIA model. For PYTHIA generated  $p'_z$  data we use the transverse pion mass  $m_{T,\pi}^2 = m_{\pi}^2 + \sigma^2$ , instead of the actual pion mass. Because of the different treatment of first and subsequent hadron emissions in PYTHIA, this choice for a pion mass will then reproduce the average PYTHIA hadronization results for full hadronization chains, as discussed in the beginning of Section 3 and shown explicitly in Section 3.3 below.



Figure 7: Three choices for latent-space target distributions I(z, c) for  $p'_z$  inputs (left) and for  $p_T$  inputs (right). See Appendix C for more details.

A key feature of the SWAE algorithm and the sliced-Wasserstein loss is the ability to 'push' the encoded latent space towards a target latent-space distribution. The choice of target distribution affects the total training time and the speed of kinematic data generation. Choosing a target latent-space distribution which is similar to the training data set distribution generally requires a fewer number of epochs to train the model to a specified accuracy compared to a target latent space which is dissimilar. This may come at a cost during the generation of kinematic data for hadronization events due to the generation of a large number of random variables obeying potentially complex probability distributions.

We demonstrate this flexibility by training with multiple target latent-space distributions, see Fig. 7. A total of six models are trained, three for each kinematic variable  $p'_z$ and  $p_T$ , with the results shown in Figs. 8 and 9. Of the three models in each kinematic variable, one model is trained using a target latent-space distribution equivalent to the training set distribution, i.e., the PYTHIA generated distribution of  $p'_z$  or  $p_T$ . The other two trained models have target latent-space distributions that are distinctly different from the training set distributions. For  $p'_z$  we choose trapezoidal and triangular target latent distributions and for  $p_T$  we choose a skewed normal and triangular target latent-space distributions. The latent-space distributions are shown in Fig. 7, while their analytic forms can be found in Appendix C. Regardless of the choice of the latent-space distribution, the trained and the target (prior) data distributions are in good agreement.

The dimension of the latent space is a tunable discrete hyperparameter, taking values  $d_z \in [2, 35]$ , see the fourth column in Table 1. The regularization parameter  $\lambda$  controls the magnitude of the sliced-Wasserstein loss and determines its relative weight in the total loss function, see Eq. (4). In practice, the regularization parameter determines how closely the encoded latent-space distribution will agree with the chosen target latent-space distribution Eq. (4) takes values  $\lambda \in [15, 35]$ , as listed in the fifth column in Table 1. Larger values are chosen in models where the target latent-space distribution is similar to the training distribution. Large values of  $\lambda$  effectively reduce the size of the explored manifold which maps decoder weight-configurations to values of the loss function (if we think of the decoder as a partition function and the loss function as a functional, large values of  $\lambda$  place the decoder near a saddle-point configuration). This improves the convergence to the minimum of  $\mathcal{L}_{rec}$ , resulting in shorter training times. This can also be explained by



Figure 8: Top: MLHAD generated  $p_z$  distributions for first-hadron emission from a string with an energy E = 50 GeV, using three different latent-space distributions, PYTHIA (blue), trapezoidal (red), and triangular (green), compared to the PYTHIA generated target distribution (purple), as well as the ratios of MLHAD generated to PYTHIA generated distributions. Bottom: comparison of the trained and target latent-space distributions for the three cases.

describing the correlation between the minimization of  $\mathcal{L}_{SW}$  and  $\mathcal{L}_{rec}$ .

The number of slices or projections used in the sliced-Wasserstein loss is also a tunable hyperparameter taking values  $L \in [15, 30]$ , as listed in the last column in Table 1. Each model uses the kinematic data generated from  $N = 4 \times 10^5$  first emission events partitioned into  $N/N_e = 4000 N_e$ -dimensional vectors, where 80% of the data is used as the training and 20% as the validation set. We use an initial learning rate value of  $10^{-3}$  and utilize PYTORCH's dynamic learning-rate scheduler to reduce the learning rate according to plateaus of the loss function during training.

#### **3.2** Labels and *E* dependent distributions

The trained models for the first-hadron emission presented in the previous section were all obtained for a fixed initial string energy, E. To reproduce the PYTHIA model for the first-hadron emissions (for string fragments with energies above  $E_{\rm cut}$ ) this is all that is required. The  $p'_z$  distributions for any string energy can be obtained from the reference value of E = 50 GeV that we used in the training by performing the rescaling, cf. Eq. (2) and Fig. 2. The  $p_T$  distributions for first emissions, on the other hand, are independent of the initial string energy.

However, the above scaling behaviors are not expected to be exact in nature. For one,



Figure 9: Top: MLHAD generated  $p_T$  distributions for first-hadron emission using three different latent-space distributions, PYTHIA (blue), skewed-normal (red), and triangular (green), compared to the PYTHIA generated target distribution (purple), as well as the ratios of MLHAD generated to PYTHIA generated distributions. Bottom: comparison of the trained and target latent-space distributions for the three cases.

at lower string energies the approximations in deriving the string Lund model are likely to fail - the quarks are not massless, and there may be couplings between  $p_T$  and  $m_h$  that are not captured by the simple transverse mass tunneling ansatz, Eq. (3). Furthermore, the origin of  $p_T$  distributions for first emissions is purely non-perturbative in nature, and thus the *E* independence of  $p_T$  distribution assumed in PYTHIA is not rooted in first principles.

The MLHAD architecture is flexible enough to allow for the dependence of first emissions on the string energy, E. This is achieved by training the conditional SWAE on label-dependent datasets, which we demonstrate next. The training proceeds in a similar way as in the previous section, but now on a dataset comprising of first-hadron emissions for four distinct string energies,  $E = \{5, 30, 700, 1000\}$  GeV.<sup>5</sup> Each  $x_i$  input vector is therefore accompanied by one of the four discrete values for the two-dimensional vectors  $c_i = (1 - c_i, c_i)$  encoding the string energy through the label  $c_i$  as defined in Eq. (9), taking  $E_{\min} = 5$  GeV and  $E_{\max} = 1000$  GeV.

The decoder in the trained cSWAE was then used to generate the first-hadron emissions at a different set of string energies,  $E = \{100, 200, 300, 400, 500\}$  GeV. Importantly, because the conditional vector is not discrete but rather depends on a continuous pa-

<sup>&</sup>lt;sup>5</sup>One could also have used emission data for continuous values of E, but binned finely enough in string energy values. We choose discrete string energies to demonstrate clearly that the cSWAE decoder can interpolate between the input labels.



Figure 10: [MLHAD generated  $p_z$  distributions using the cSWAE model trained on data with string energies different from training and compared with PYTHIA (black).]

rameter defined between the minimum and maximum energies  $(E_{\min}, E_{\max})$  the trained decoder is able to interpolate between labels (ones which the decoder has not trained on explicitly, see Fig. 4) and rescale the kinematic distributions accordingly. This considerably increases the flexibility of generating training datasets as the user is able to choose the number of interpolation points which the model can use as anchors in generating data with a unique energy label. The comparison of MLHAD and PYTHIA generated  $p_z$  distributions for the first-hadron emissions is shown in Fig. 10, demonstrating that MLHAD reproduces faithfully the PYTHIA results.

### 3.3 Hadronization chain

As shown in the previous subsections the cSWAE trained models in MLHAD are able to accurately reproduce PYTHIA's first emission kinematics for a hadronized  $q\bar{q}$  system in the center-of-mass frame of the string. In this section we show how well the MLHAD decoder reproduces the full PYTHIA hadronization event. The implementation can be summarized as follows: from the initial string system, one string end is chosen randomly, while PYTHIA flavor selector is used to determine the flavor ID of the emitted hadron. Given the energy of the initial string end in the center-of-mass frame,  $p'_z$  and  $p_T$  are sampled using the corresponding cSWAE models. The  $p'_z$  and  $p_T$  of the emitted hadron are transformed to  $p_x, p_y, p_z$  variables using Eqs. (1) and (2), and boosted to the lab frame. The string fragment is boosted to its center-of-mass frame, see Fig. 6, after which one repeats the hadron emission process until the string energy in the center of mass of the remaining string fragment falls below the IR cutoff,  $E_{cut}$ . The implemented fragmentation chain architecture is illustrated in Fig. 5.

Fig. 11 shows a comparison between the hadronization chain multiplicities obtained by PYTHIA (blue) and by the MLHAD model trained on first emission data (red). In both cases, starting from the initial string energy of E = 50 GeV, on average 9.1 hadron emissions occur before the string fragment energy drops below the cutoff energy,  $E_{\text{cut}} = 5$ 



Figure 11: Comparison of the number of hadrons produced in the fragmentation chain of a single string for a sample of  $10^4$  strings, compared between PYTHIA (blue) and MLHAD (red) generated hadronization events.

GeV. The MLHAD decoder also reproduces well the distribution of hadronization chain multiplicities. Only a few hadronization events result in just a few hadrons, a bulk of hadronization events contain between 7 to 13 hadrons, and both hadronization chain generators feature a tail of rather long hadronization chains. The differences between the PYTHIA and MLHAD hadron multiplicity distributions are in most cases at the level of 5-10%, where the largest deviations occur for hadronization events with just a few hadron emissions. This is expected, given that PYTHIA and MLHAD models of hadronization differ in the treatment of the very first emission, see the discussion at the beginning of Section 3.

In Fig. 12 we also show the comparison of the average multiplicity of the hadronization chain as a function of the initial parton energy, obtained either with PYTHIA (blue solid line) or with MLHAD (red). We see that MLHAD is able to reproduce the PYTHIA fragmentation chain length averages, and in particular also give the expected log E dependence of the average number of produced hadrons. For each energy the multiplicity distributions also match well, which we checked explicitly, while in the figure we only show the result for MLHAD to guide the eye (red density plot). The density plot scan was performed by randomly choosing an initial parton energy E between 20 GeV-1000 GeV and binning each fragmentation chain length with a parton energy resolution of 22 GeV and chain length resolution of 1.7 hadrons for a total of  $2 \times 10^4$  fragmentation events. The minimal initial string energy was chosen to be 20 GeV such that it is still well above the imposed hadron emission cut  $E_{\rm cut} = 5$  GeV.

## 4 Conclusion and Outlook

The cSWAE architecture that was developed in this work appears to be well suited for modeling the nonperturbative process of hadronization – the creation of hadrons from the energy stored in the string connecting a  $q\bar{q}$  pair. We have demonstrated this by training the MLHAD hadronization models to a simplified version of PYTHIA hadronization, limited to only light quark flavor endings of the string, and allowing only for pions to be the final-state



Figure 12: Comparison of the average number of hadrons produced in the fragmentation chain of a single string as a function of the initial parton energy E $(E_{\text{string}} = 2E)$ , produced using PYTHIA (blue) and MLHAD (red). The density plot shows the multiplicity distributions obtained with MLHAD for  $2 \times 10^4$ fragmentation chains.

hadrons. Furthermore, we utilized the scaling properties of the PYTHIA hadronization model that simplified the cSWAE training, requiring training at just a single string energy. Even so, the results shown in Figs. 8, 9 and 11 are very encouraging. The PYTHIA first-hadron emission distributions at a fixed string energy, Fig. 8, 9, are faithfully reproduced by the MLHAD decoder, as are the hadron multiplicities for full hadronization chains, Fig. 11.

The cSWAE architecture also has enough built in flexibility that it should be possible to extend the MLHAD model to handle all possible string flavors and kinematics. We have already shown that the inclusion of a label allows for an interpolation of the hadronization models to different string energies, see Fig. 10. This should then also allow to extend the MLHAD models below the string energy cut of 5 GeV that we imposed in this preliminary exploration. Similarly, the conditional label could be used for MLHAD to handle the generation of hadron flavors, including possible kinematic dependencies. The MLHAD architecture should also allow us to model any correlations between  $p_z$  and  $p_T$ distributions of the emitted hadrons, if these are present in data, even though currently we used the absence of such correlations in PYTHIA generated data to simplify the training of MLHAD models. Another important feature that we anticipate to be particularly important once MLHAD is trained directly on experimental data, is the flexibility in the choice of the latent-space distributions, making it easier to adapt to any possible features not captured by the rather constrained form of the Lund fragmentation function underlying the hadronization implementation in PYTHIA. Finally, some of the planned extensions of the MLHAD hadronization framework may require more thought, most notably how to best model the hadronization of baryons and include gluons.

While in this paper the training of MLHAD was performed on the first hadron emissions in the Pythia output, such training will not be possible when using real experimental data, since such information is physically not possible to extract directly from data. Instead, the training will need to be performed on the physically accessible observables constructed from particle flows measured either in  $e^+e^-$  or pp collisions with two, three or more jets in the final state. We anticipate that this is where the machine learning approach to hadronization will prove most useful — capturing the many observables in principle available in the data, such as hadron multiplicities, angular separations and momentum distributions for various hadrons [(see [75–80] for a selection of potentially useful observables)]. [While many of these observables are not currently available in the literature, open-data efforts by a number of collaborations have or will make access possible.] This data-collection is tedious when performed through human intervention and is a problem that calls for a machine learning based optimization. We believe that the presented MLHAD cSWAE architecture is well suited to achieve this next step, [and we are in the process of building a pipeline to perform training of MLHAD on actual data]. [In addition different generative models like Normalizing Flows will be explored, which provide a tractable probability distribution function.]

## Acknowledgments

We thank Jared Evans for collaboration in the initial stages of this work, and Stephen Mrenna, Manuel Szewc, and Mike Williams for useful comments on the manuscript.

**Funding information.** AY, JZ, and TM acknowledge support in part by the DOE grant de-sc0011784 and NSF OAC-2103889. PI is supported in part by NSF OAC-2103889.

# A Public code MLhad\_v0.1

The public code may be accessed through https://gitlab.com/uchep/mlhad. The public directory includes example files allowing the user to train and implement cSWAE models in full fragmentation chains. The programs are written in Python and extensively use the PYTHIA, PYTORCH and SCIKIT-LEARN libraries. Installation instructions can be found on the respective installation pages for each library.

The provided programs can be split into two categories: training cSWAE models and generating hadronization events. The latter relies on the former. However, we have also provided pre-trained models such that the user can generate hadronization events without explicitly training a model.

Training a unique model configuration can be done by modifying the files pT\_SWAE.py, pz\_SWAE.py, or pz\_cSWAE.py. The SWAE programs contain examples of label-independent training, while the cSWAE program provides an example of label-dependent training. The model hyperparameters and target latent distribution described in Section 2 have been set to default values to provide a reasonable starting configuration but may be modified. Label independent kinematic training datasets for  $p_z$  and  $p_T$  have been provided as well as a label-dependent  $p_z$  dataset.

Full hadronization events use the trained model decoder to generate hadronic kinematics. An example of generating this kinematic data from SWAE trained model decoders can be found in model\_pxpypz.py. The setup of our modified fragmentation chain which utilizes these kinematics can be seen in frag\_chain.py.

### **B** Sliced Wasserstein distance

In this appendix we give a short overview of the Wasserstein distance and the sliced-Wasserstein distance.

**The Wasserstein distance.** The Earth mover's distance or the Wasserstein distance gives a measure of how different two distributions are, given a metric space  $\Omega$  and a space of Borel probability measures  $\mathcal{P}(\Omega)$  on  $\Omega$ . The *p*-Wasserstein distance  $W_p(\mu, \nu)$  between any two probability measures  $\mu \in \mathcal{P}(X)$  and  $\nu \in \mathcal{P}(Y)$  is [81]

$$W_p(\mu,\nu) := \left(\inf_{\pi \in \Pi(\mu,\nu)} \int_X c(x,y) d\pi(x,y)\right)^{\frac{1}{p}},\tag{11}$$

where c(x, y) is the cost function,  $\Pi(\mu, \nu)$  is the set of all transportation plans, with  $\pi \in \Pi(\mu, \nu)$ , while  $p \in [1, \infty)$ . The distance  $W_1$  is also commonly called the Kantorovich-Rubinstein distance.

If  $\mu$  and  $\nu$  are one-dimensional measures, the Wasserstein distance has a closed-form expression

$$W_p(\mu,\nu) = \left(\int_0^1 |F_{\mu}^{-1}(z) - F_{\nu}^{-1}(z)|^p dz\right)^{1/p},\tag{12}$$

where  $F_{\mu(\nu)}(x) = \int_{-\infty}^{x} I_{\mu(\nu)}(\tau) d\tau$  are the cumulative distribution functions, with  $I_{\mu}$  and  $I_{\nu}$ the probability density functions for the measures  $\mu$  and  $\nu$ , respectively. The  $W_p(\mu, \nu)$  for the one dimensional case can therefore be calculated by simply sorting the samples from the two distributions and calculating the average cost.

**Radon transform and the sliced-Wasserstein distance.** An approximate value for the Wasserstein distance  $W_p$  between two higher dimensional distributions on  $X = \mathcal{R}^d$ can be obtained efficiently from a set of projections to one-dimensional distributions, since for each of these one can use the closed form of Eq. (12). The projection from the higher dimensional distribution to the one-dimensional representation is done by the Radon transform.

The *d*-dimensional Radon transform R maps a function  $I \in L^1(\mathcal{R}^d)$  to [82]

$$RI(t,\theta) = \int_{\mathcal{R}^d} |I(x)| \delta(t - \langle x, \theta \rangle) dx, \qquad (13)$$

with  $(t,\theta) \in \mathcal{R} \times S^{d-1}$ , where  $S^{d-1}$  is the unit sphere in  $\mathcal{R}^d$ ,  $\delta(\cdot)$  is the delta function and  $\langle , \rangle$  is the Euclidean scalar product. For a fixed direction  $\theta$  the Radon transform  $RI_{\mu}(\cdot,\theta)$  therefore gives a one dimensional marginal distribution of  $I_{\mu}$  that is obtained by integrating  $I_{\mu}$  over the hyperplane orthogonal to  $\theta$ .

The sliced-Wasserstein distance  $SW_p(I_\mu, I_\nu)$  between  $I_\mu$  and  $I_\nu$  is defined as

$$SW_p(I_{\mu}, I_{\nu}) = \left(\int_{\mathcal{S}^{d-1}} W_p(RI_{\mu}(\cdot, \theta), RI_{\nu}(\cdot, \theta)d\theta)\right)^{\frac{1}{p}}.$$
(14)

The Wasserstein distance between each of the one dimensional projections (slicings)  $RI_{\mu}(\cdot, \theta)$ and  $RI_{\nu}(\cdot, \theta)$  is obtained straightforwardly using the closed form result of Eq. (12). The integral over the unit sphere vectors  $\theta$  probes all the possible slicings. Furthermore,  $SW_p(I_{\mu}, I_{\nu})$  approximates  $W_p(I_{\mu}, I_{\nu})$  "well enough" [83]. The integration in Eq. (14) over the unit sphere in  $\mathcal{R}^d$  can be estimated using a Monte Carlo integration that draws samples  $\{\theta_l\}$  from the uniform distribution on  $\mathcal{S}^{d-1}$ , which substitutes a finite sample average for the integral [84],

$$SW_p(I_{\mu}, I_{\nu}) \approx \left(\frac{1}{L} \sum_{l=1}^{L} W_p(RI_{\mu}(\cdot, \theta_l), RI_{\nu}(\cdot, \theta_l))\right)^{\frac{1}{p}},\tag{15}$$

where L is the number of projections (slicings). With this result, the sliced-Wasserstein distance is obtained by solving a finite number of one-dimensional optimal transport problems, each of which has a closed-form solution. Furthermore, the sliced-Wasserstein distance approximates well the Wasserstein distance and thus can be used as a useful discriminator for the similarity of distributions. More details can be found in [84] and [65].

### C Latent distributions

The analytic forms of the latent target distributions used in the training of cSWAE in Section 3.1 are  $\left(2\left(2-2\right)\right)$ 

$$I_{\text{tri.}}(z; a, b, c) = \begin{cases} \frac{2(z-a)}{(b-a)(c-a)}, & a \le z \le c, \\ \frac{2(b-z)}{(b-a)(b-c)}, & c < z \le b, \end{cases}$$
(16)

for the triangular distribution, and

$$I_{\text{trap.}}(z; a, b, c, d) = \begin{cases} \frac{2}{d + c - a - b} \frac{z - a}{b - a}, & a \le z < b, \\ \frac{2}{d + c - a - b}, & b \le z < c, \\ \frac{2}{d + c - a - b} \frac{d - z}{d - c}, & c \le z \le d, \end{cases}$$
(17)

for the trapezoidal distribution. For a given initial parton energy E the choices of parameters a, b, c, d can be seen in Table 2. The target latent-space distributions are then given by

$$I_{\text{tri.}}(\boldsymbol{z}, \boldsymbol{c}) = \prod_{k=1}^{N_e} I_{\text{tri.}}(z_k; a, b, c), \qquad I_{\text{trap.}}(\boldsymbol{z}, \boldsymbol{c}) = \prod_{k=1}^{N_e} I_{\text{trap}}(z_k; a, b, c, d), \tag{18}$$

that is we take the same values of a, b, c, d parameters for all  $d_z$  latent dimensions.

The normal and skewed-normal distributions are given by

$$I_{\text{Gauss}}(z;\mu,\sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(z-\mu)^2}{2\sigma^2}\right),\tag{19}$$

$$I_{\text{Skew-Gauss}}(z;\mu,\sigma,\alpha) = 2I_{\text{Gauss}}(z;\mu,\sigma)\Phi\left(\frac{\alpha(z-\mu)}{\sigma}\right),\tag{20}$$

respectively, where

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-t^2/2} dt.$$
 (21)

The  $\mu$ ,  $\sigma$ , and  $\alpha$  are the fit parameters corresponding to the mean, standard deviation, and skewness of the distribution, respectively. As in Eq. (18) the  $d_z$  dimensional latent-space distributions are products of one dimensional ones with the same  $\mu$ ,  $\sigma$ ,  $\alpha$  parameters. For  $p_T$  we have  $\mu = 0.099$ ,  $\sigma = 0.257$ , and  $\alpha = 4.259$ .

Variable $\boldsymbol{x}$	Target $\boldsymbol{z}$	a	b	c	d
~/	Trapezoidal	0.04E	0.16E	0.24E	E
$p_z$	Triangular	0.04E	0.2E	E	_
$p_T$	Triangular	0.0	0.3	1.0	_

**Table 2:** The  $p'_z$  and  $p_T$  latent-space distribution parameters.

## References

- R. Frederix, S. Frixione, V. Hirschi, D. Pagani, H. S. Shao and M. Zaro, *The au-tomation of next-to-leading order electroweak calculations*, JHEP 07, 185 (2018), doi:10.1007/JHEP11(2021)085, [Erratum: JHEP 11, 085 (2021)], e-print:1804.10017.
- [2] J. Bellm, S. Gieseke and S. Plätzer, Merging NLO Multi-jet Calculations with Improved Unitarization, Eur. Phys. J. C 78(3), 244 (2018), doi:10.1140/epjc/s10052-018-5723-2, e-print:1705.06700.
- [3] J. M. Campbell, S. Höche, H. T. Li, C. T. Preuss and P. Skands, Towards NNLO+PS Matching with Sector Showers (2021), e-print:2108.07133.
- [4] S. Höche and S. Prestel, The midpoint between dipole and parton showers, Eur. Phys. J. C 75(9), 461 (2015), doi:10.1140/epjc/s10052-015-3684-2, e-print:1506.05057.
- [5] J. Alwall, R. Frederix, S. Frixione, V. Hirschi, F. Maltoni, O. Mattelaer, H. S. Shao, T. Stelzer, P. Torrielli and M. Zaro, *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), doi:10.1007/JHEP07(2014)079, e-print:1405.0301.
- [6] T. Sjöstrand, S. Ask, J. R. Christiansen, R. Corke, N. Desai, P. Ilten, S. Mrenna, S. Prestel, C. O. Rasmussen and P. Z. Skands, An introduction to PYTHIA 8.2, Comput. Phys. Commun. 191, 159 (2015), doi:10.1016/j.cpc.2015.01.024, eprint:1410.3012.
- [7] J. Bellm et al., Herwig 7.0/Herwig++ 3.0 release note, Eur. Phys. J. C 76(4), 196 (2016), doi:10.1140/epjc/s10052-016-4018-8, e-print:1512.01178.
- [8] E. Bothmann et al., Event Generation with Sherpa 2.2, SciPost Phys. 7(3), 034 (2019), doi:10.21468/SciPostPhys.7.3.034, e-print:1905.09127.
- B. Andersson, G. Gustafson, G. Ingelman and T. Sjostrand, Parton Fragmentation and String Dynamics, Phys. Rept. 97, 31 (1983), doi:10.1016/0370-1573(83)90080-7.
- [10] B. Andersson, The Lund model, Camb. Monogr. Part. Phys. Nucl. Phys. Cosmol. 7, 1 (1997).
- [11] S. Ferreres-Solé and T. Sjöstrand, The space-time structure of hadronization in the Lund model, Eur. Phys. J. C 78(11), 983 (2018), doi:10.1140/epjc/s10052-018-6459-8, e-print:1808.04619.
- [12] R. D. Field and S. Wolfram, A QCD Model for e+ e- Annihilation, Nucl. Phys. B 213, 65 (1983), doi:10.1016/0550-3213(83)90175-X.

- [13] T. D. Gottschalk, An Improved Description of Hadronization in the {QCD} Cluster Model for e<sup>+</sup>e<sup>-</sup> Annihilation, Nucl. Phys. B 239, 349 (1984), doi:10.1016/0550-3213(84)90253-0.
- [14] B. Webber, A QCD Model for Jet Fragmentation Including Soft Gluon Interference, Nucl. Phys. B 238, 492 (1984), doi:10.1016/0550-3213(84)90333-X.
- [15] F. Bishara and M. Montull, (Machine) Learning amplitudes for faster event generation (2019), e-print:1912.11055.
- [16] S. Badger and J. Bullock, Using neural networks for efficient evaluation of high multiplicity scattering amplitudes, JHEP 06, 114 (2020), doi:10.1007/JHEP06(2020)114, e-print:2002.07516.
- [17] C. Gao, J. Isaacson and C. Krause, *i-flow: High-dimensional Integration and Sampling with Normalizing Flows*, Mach. Learn. Sci. Tech. 1(4), 045023 (2020), doi:10.1088/2632-2153/abab62, e-print:2001.05486.
- [18] C. Gao, S. Höche, J. Isaacson, C. Krause and H. Schulz, Event Generation with Normalizing Flows, Phys. Rev. D 101(7), 076002 (2020), doi:10.1103/PhysRevD.101.076002, e-print:2001.10028.
- [19] I. Chahrour and J. D. Wells, Function Approximation for High-Energy Physics: Comparing Machine Learning and Interpolation Methods (2021), e-print:2111.14788.
- [20] R. Winterhalder, V. Magerya, E. Villa, S. P. Jones, M. Kerner, A. Butter, G. Heinrich and T. Plehn, *Targeting Multi-Loop Integrals with Neural Networks* (2021), e-print:2112.09145.
- [21] R. Kansal, J. M. Duarte, B. Orzari, T. Tomei, M. Pierini, M. Touranakou, J. R. Vlimant and D. Gunopoulos, *Graph generative adversarial networks for sparse data generation in high energy physics*, ArXiv abs/2012.00173 (2020).
- [22] J. W. Monk, Deep learning as a parton shower, Journal of High Energy Physics 2018(12) (2018), doi:10.1007/jhep12(2018)021.
- [23] R. Kansal, J. Duarte, H. Su, B. Orzari, T. Tomei, M. Pierini, M. Touranakou, J. Vilmant and D. Gunopulos, *Particle cloud generation with message passing generative* adversarial networks, ArXivorg.
- [24] Y. S. Lai, D. M. Neill, M. Plosko'n and F. M. Ringer, Explainable machine learning of the underlying physics of high-energy particle collisions, Physics Letters B (2022).
- [25] B. Orzari, T. Tomei, M. Pierini, M. Touranakou, J. Duarte, R. Kansal, J.-R. Vlimant and D. Gunopulos, Sparse data generation for particle-based simulation of hadronic jets in the lhc, doi:10.48550/ARXIV.2109.15197 (2021).
- [26] M. Touranakou, N. Chernyavskaya, J. Duarte, D. Gunopulos, R. Kansal, B. Orzari, M. Pierini, T. Tomei and J.-R. Vlimant, *Particle-based fast jet simulation at the lhc* with variational autoencoders, doi:10.48550/ARXIV.2203.00520 (2022).
- [27] E. Bothmann and L. Debbio, Reweighting a parton shower using a neural network: the final-state case, JHEP 01, 033 (2019), doi:10.1007/JHEP01(2019)033, e-print:1808.07802.

- [28] S. Tsan, R. Kansal, A. Aportela, D. Diaz, J. Duarte, S. Krishna, F. Mokhtar, J.-R. Vlimant and M. Pierini, Particle Graph Autoencoders and Differentiable, Learned Energy Mover's Distance, In 35th Conference on Neural Information Processing Systems (2021), e-print:2111.12849.
- [29] C. Krause and D. Shih, Caloflow: Fast and accurate generation of calorimeter showers with normalizing flows, ArXiv abs/2106.05285 (2021).
- [30] C. Krause and D. Shih, Caloflow ii: Even faster and still accurate generation of calorimeter showers with normalizing flows, ArXiv abs/2110.11377 (2021).
- [31] K. T. Matchev, A. Roman and P. Shyamsundar, Uncertainties associated with GANgenerated datasets in high energy physics (2020), e-print:2002.06307.
- [32] Y. Alanazi et al., Simulation of electron-proton scattering events by a Feature-Augmented and Transformed Generative Adversarial Network (FAT-GAN) (2020), doi:10.24963/ijcai.2021/293, e-print:2001.11103.
- [33] B. Nachman and J. Thaler, Neural resampler for Monte Carlo reweighting with preserved uncertainties, Phys. Rev. D 102(7), 076004 (2020), doi:10.1103/PhysRevD.102.076004, e-print:2007.11586.
- [34] B. Stienen and R. Verheyen, Phase space sampling and inference from weighted events with autoregressive flows, SciPost Phys. 10(2), 038 (2021), doi:10.21468/SciPostPhys.10.2.038, e-print:2011.13445.
- [35] A. Butter, S. Diefenbacher, G. Kasieczka, B. Nachman and T. Plehn, GANplifying event samples, SciPost Phys. 10(6), 139 (2021), doi:10.21468/SciPostPhys.10.6.139, e-print:2008.06545.
- [36] M. Backes, A. Butter, T. Plehn and R. Winterhalder, How to GAN Event Unweighting, SciPost Phys. 10(4), 089 (2021), doi:10.21468/SciPostPhys.10.4.089, eprint:2012.07873.
- [37] K. Danziger, T. Janßen, S. Schumann and F. Siegert, Accelerating Monte Carlo event generation – rejection sampling using neural network event-weight estimates (2021), e-print:2109.11964.
- [38] A. Butter, T. Heimel, S. Hummerich, T. Krebs, T. Plehn, A. Rousselot and S. Vent, Generative Networks for Precision Enthusiasts (2021), e-print:2110.13632.
- [39] G. Bíró, B. Tankó-Bartalis and G. G. Barnaföldi, Studying Hadronization by Machine Learning Techniques (2021), e-print:2111.15655.
- [40] J. N. Howard, S. Mandt, D. Whiteson and Y. Yang, Foundations of a Fast, Data-Driven, Machine-Learned Simulator (2021), e-print:2101.08944.
- [41] G. Quétant, M. Drozdova, V. Kinakh, T. Golling and S. Voloshynovskiy, Turbo-Sim: a generalised generative model with a physical latent space (2021), e-print:2112.10629.
- [42] S. Bieringer, A. Butter, S. Diefenbacher, E. Eren, F. Gaede, D. Hundhausen, G. Kasieczka, B. Nachman, T. Plehn and M. Trabs, *Calomplification – The Power* of Generative Calorimeter Models (2022), e-print:2202.07352.

- [43] E. Buhmann, S. Diefenbacher, E. Eren, F. Gaede, G. Kasieczka, A. Korol and K. Krüger, *Getting High: High Fidelity Simulation of High Granularity Calorimeters* with High Speed, Comput. Softw. Big Sci. 5(1), 13 (2021), doi:10.1007/s41781-021-00056-0, e-print:2005.05334.
- [44] E. Bothmann, T. Janßen, M. Knobbe, T. Schmale and S. Schumann, Exploring phase space with Neural Importance Sampling, SciPost Phys. 8(4), 069 (2020), doi:10.21468/SciPostPhys.8.4.069, e-print:2001.05478.
- [45] L. de Oliveira, M. Paganini and B. P. Nachman, Learning particle physics by example: Location-aware generative adversarial networks for physics synthesis, Computing and Software for Big Science 1, 1 (2017).
- [46] 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, 014021 (2018), doi:10.1103/PhysRevD.97.014021.
- [47] M. Paganini, L. de Oliveira and B. Nachman, Accelerating science with generative adversarial networks: An application to 3d particle showers in multilayer calorimeters, Physical Review Letters 120(4) (2018), doi:10.1103/PhysRevLett.120.042003.
- [48] P. Musella and F. Pandolfi, Fast and accurate simulation of particle detectors using generative adversarial networks, Computing and Software for Big Science 2, 1 (2018), doi:10.1007/s41781-018-0015-y.
- [49] M. Erdmann, L. Geiger, J. Glombitza and D. Schmidt, Generating and refining particle detector simulations using the Wasserstein distance in adversarial networks, Comput. Softw. Big Sci. 2(1), 4 (2018), doi:10.1007/s41781-018-0008-x, eprint:1802.03325.
- [50] M. Erdmann, J. Glombitza and T. Quast, Precise simulation of electromagnetic calorimeter showers using a Wasserstein Generative Adversarial Network, Comput. Softw. Big Sci. 3, 4 (2019), doi:10.1007/s41781-018-0019-7, e-print:1807.01954.
- [51] ATLAS Collaboration, Energy resolution with aGANfor Shower ATLAS-SIM-2019-004, Fast Simulation inATLAS, Https://atlas.web.cern.ch/Atlas/GROUPS/PHYSICS/PLOTS/SIM-2019-004/ (2019).
- [52] D. Belayneh et al., Calorimetry with Deep Learning: Particle Simulation and Reconstruction for Collider Physics, Eur. Phys. J. C 80(7), 688 (2020), doi:10.1140/epjc/s10052-020-8251-9, e-print:1912.06794.
- [53] E. Buhmann, S. Diefenbacher, E. Eren, F. Gaede, G. Kasieczka, A. Korol and K. Krüger, *Decoding Photons: Physics in the Latent Space of a BIB-AE Generative Network*, EPJ Web Conf. **251**, 03003 (2021), doi:10.1051/epjconf/202125103003, e-print:2102.12491.
- [54] C. Chen, O. Cerri, T. Q. Nguyen, J. R. Vlimant and M. Pierini, Analysis-Specific Fast Simulation at the LHC with Deep Learning, Comput. Softw. Big Sci. 5(1), 15 (2021), doi:10.1007/s41781-021-00060-4.
- [55] C. Krause and D. Shih, CaloFlow: Fast and Accurate Generation of Calorimeter Showers with Normalizing Flows (2021), e-print:2106.05285.

- [56] C. Krause and D. Shih, CaloFlow II: Even Faster and Still Accurate Generation of Calorimeter Showers with Normalizing Flows (2021), e-print:2110.11377.
- [57] M. Paganini, L. de Oliveira and B. Nachman, CaloGAN: Simulating 3d high energy particle showers in multilayer electromagnetic calorimeters with generative adversarial networks, Physical Review D 97(1) (2018), doi:10.1103/physrevd.97.014021.
- [58] P. Ilten, M. Williams and Y. Yang, Event generator tuning using Bayesian optimization, JINST 12(04), P04028 (2017), doi:10.1088/1748-0221/12/04/P04028, eprint:1610.08328.
- [59] A. Andreassen and B. Nachman, Neural Networks for Full Phase-space Reweighting and Parameter Tuning, Phys. Rev. D 101(9), 091901 (2020), doi:10.1103/PhysRevD.101.091901, e-print:1907.08209.
- [60] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville and Y. Bengio, *Generative adversarial nets*, In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence and K. Weinberger, eds., *Advances in Neural Information Processing Systems*, vol. 27. Curran Associates, Inc. (2014).
- [61] D. P. Kingma and M. Welling, Auto-encoding variational bayes (2014), eprint:1312.6114.
- [62] D. J. Rezende and S. Mohamed, Variational inference with normalizing flows, In Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37, ICML'15, p. 1530–1538. JMLR.org (2015).
- [63] M. Bellagente, A. Butter, G. Kasieczka, T. Plehn, A. Rousselot, R. Winterhalder, L. Ardizzone and U. Köthe, *Invertible Networks or Partons to Detector and Back Again*, SciPost Phys. 9, 74 (2020), doi:10.21468/SciPostPhys.9.5.074.
- [64] M. Bellagente, A. Butter, G. Kasieczka, T. Plehn and R. Winterhalder, How to GAN away Detector Effects, SciPost Phys. 8, 70 (2020), doi:10.21468/SciPostPhys.8.4.070.
- [65] S. Kolouri, C. E. Martin and G. K. Rohde, Sliced-wasserstein autoencoder: An embarrassingly simple generative model, CoRR abs/1804.01947 (2018), eprint:1804.01947.
- [66] I. O. Tolstikhin, O. Bousquet, S. Gelly and B. Schölkopf, Wasserstein auto-encoders, CoRR abs/1711.01558 (2017), e-print:1711.01558.
- [67] A. Andreassen, I. Feige, C. Frye and M. D. Schwartz, Binary JUNIPR: an interpretable probabilistic model for discrimination, Phys. Rev. Lett. 123(18), 182001 (2019), doi:10.1103/PhysRevLett.123.182001, e-print:1906.10137.
- [68] A. Andreassen, I. Feige, C. Frye and M. D. Schwartz, JUNIPR: a Framework for Unsupervised Machine Learning in Particle Physics, Eur. Phys. J. C 79(2), 102 (2019), doi:10.1140/epjc/s10052-019-6607-9, e-print:1804.09720.
- [69] E. Bothmann and L. Debbio, Reweighting a parton shower using a neural network: the final-state case, JHEP 01, 033 (2019), doi:10.1007/JHEP01(2019)033, e-print:1808.07802.
- [70] 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), 4 (2017), doi:10.1007/s41781-017-0004-6, e-print:1701.05927.

- [71] J. W. Monk, Deep Learning as a Parton Shower, JHEP 12, 021 (2018), doi:10.1007/JHEP12(2018)021, e-print:1807.03685.
- [72] K. Dohi, Variational Autoencoders for Jet Simulation (2020), e-print:2009.04842.
- [73] B. Orzari, T. Tomei, M. Pierini, M. Touranakou, J. Duarte, R. Kansal, J.-R. Vlimant and D. Gunopulos, Sparse Data Generation for Particle-Based Simulation of Hadronic Jets in the LHC, In 38th International Conference on Machine Learning Conference (2021), e-print:2109.15197.
- [74] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf et al., Pytorch: An imperative style, high-performance deep learning library, In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox and R. Garnett, eds., Advances in Neural Information Processing Systems 32, pp. 8024–8035. Curran Associates, Inc. (2019).
- [75] Y.-T. Chien, A. Deshpande, M. M. Mondal and G. Sterman, Probing hadronization with flavor correlations of leading particles in jets, Phys. Rev. D 105(5), L051502 (2022), doi:10.1103/PhysRevD.105.L051502, e-print:2109.15318.
- [76] R. A. Khalek, V. Bertone and E. R. Nocera, Determination of unpolarized pion fragmentation functions using semi-inclusive deep-inelastic-scattering data, Phys. Rev. D 104(3), 034007 (2021), doi:10.1103/PhysRevD.104.034007, e-print:2105.08725.
- [77] V. Bertone, N. P. Hartland, E. R. Nocera, J. Rojo and L. Rottoli, *Charged hadron fragmentation functions from collider data*, Eur. Phys. J. C 78(8), 651 (2018), doi:10.1140/epjc/s10052-018-6130-4, e-print:1807.03310.
- M. Soleymaninia, H. Hashamipour and H. Khanpour, Neural network QCD analysis of charged hadron fragmentation functions in the presence of SIDIS data, Phys. Rev. D 105(11), 114018 (2022), doi:10.1103/PhysRevD.105.114018, e-print:2202.10779.
- [79] H. Chen, I. Moult, J. Thaler and H. X. Zhu, Non-Gaussianities in collider energy flux, JHEP 07, 146 (2022), doi:10.1007/JHEP07(2022)146, e-print:2205.02857.
- [80] P. T. Komiske, I. Moult, J. Thaler and H. X. Zhu, Analyzing N-point Energy Correlators Inside Jets with CMS Open Data (2022), e-print:2201.07800.
- [81] C. e. Villani, Optimal transport, old and new, Springer, Berlin (2008).
- [82] S. Helgason, Integral Geometry and Radon Transforms, Springer, New York (2015).
- [83] F. Santambrogio, Optimal Transport for Applied Mathematicians, Springer, Switzerland (2015).
- [84] S. Kolouri, K. Nadjahi, U. Simsekli, R. Badeau and G. K. Rohde, Generalized sliced wasserstein distances, CoRR abs/1902.00434 (2019), e-print:1902.00434.