# Towards a data-driven model of Hadronization using Normalizing Flows

## Pheno 2024

Based on **arXiv:2311.09296**, **SciPost Phys. 14, 027 (2023)**, and 2407.XXXXX

**Ahmed Youssef**

**Ph.D. Candidate, University of Cincinnati**
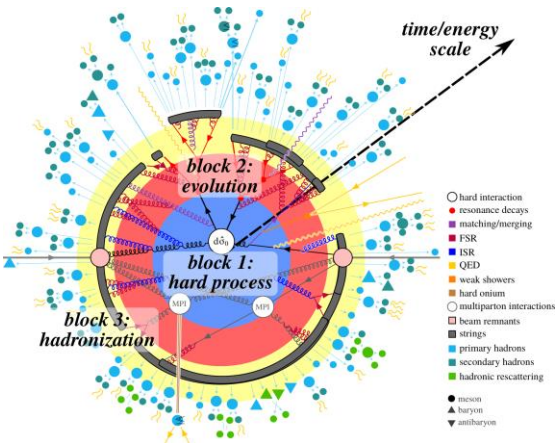
**youssead@ucmail.uc.edu**

May 15th, 2024

**In collaboration with:**

**C. Bierlich, P. Ilten, T. Menzo, S. Mrenna, M. Szewc, M.K. Wilkinson, and J. Zupan**

## Simulating Collision



*time/energy scale*

block 2:
*evolution*

$d\hat{\sigma}_0$

block 1:
*hard process*

MPI          MPI

block 3:
*hadronization*

○ hard interaction
● resonance decays
■ matching/merging
■ FSR
■ ISR
■ QED
■ weak showers
■ hard onium
○ multiparton interactions
■ beam remnants
■ strings
■ primary hadrons
■ secondary hadrons
■ hadronic rescattering

● meson
▲ baryon
▼ antibaryon

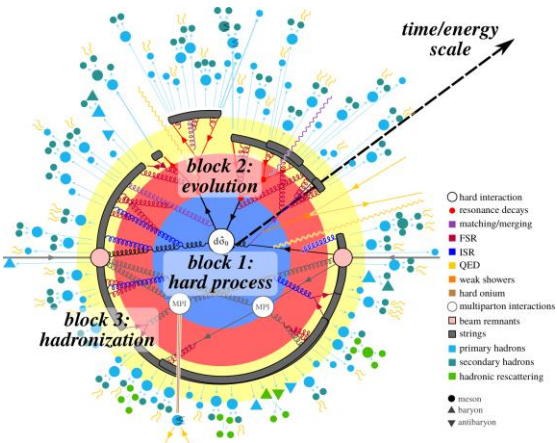➡ **Hard process:**
in itial high-energy interaction

➡ **Evolution:**
parton shower

*perturbative*

➡ **Hadronization:**
com bine quarks and gluons

*non-perturbative*

## Simulating Collision



➤ **Hard process:**
initial high-energy interaction

➤ **Evolution:**
parton shower

*perturbative*

➤ **Hadronization:**
combine quarks and gluons

*non- -ve*

Use ML!

MLHAD

University of CINCINNATI

MLHAD

**A series of progressive steps needs to be done before practically useful in Pythia simulations**

SciPost Phys. 14, 027 (2023)

**Train on truth level Pythia output (not obs. In exp)**

We are here

**Train on real data (i.e., just already measured information)**

arXiv:2311.09296

**Develop a framework to propagate errors**

arXiv:2311.09296, 2407.XXXXX

**Train on mock data (i.e., just observable information)**

Partial results

**Replace/Complement Pythia string model**

**A series of progressive steps needs to be done before practically useful in Pythia simulations**

SciPost Phys. 14, 027 (2023)

**Train on truth level Pythia output (not obs. In exp)**

arXiv:2311.09296

**Develop a framework to propagate errors**

arXiv:2311.09296, 2407.XXXXX

**Train on mock data (i.e., just observable information)**

**We are here**

**Partial results**

**Train on real data (i.e., just already measured information)**

**Replace/Complement Pythia string model**

**Pythia**   **String model**



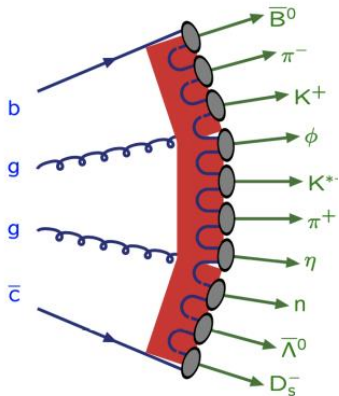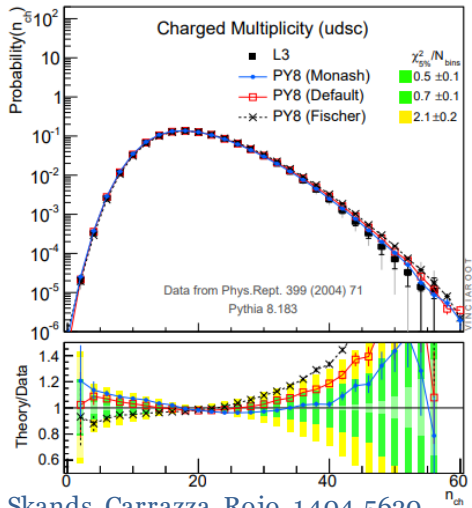Fig from Vitev, IV YR workshop, 2020

**Pythia**     **String model**



Fig from Vitev, IV YR workshop, 2020



Observables

Skands, Carrazza, Rojo, 1404.5630

University of
CINCINNATI

**Pythia**    **String model**



Not Observed!

$\overline{B}^0$

$\pi^-$

$K^+$

$\phi$

$K^{*-}$

$\pi^+$

$\eta$

$n$

$\overline{\Lambda}^0$

$D_s^-$

b

g

g

$\overline{c}$

Observables



Skands, Carrazza, Rojo, 1404.5630

Fig from Vitev, IV YR workshop, 2020

# MAGIC

## (Microscopic Alterations Generated from IR Collections)

| Step 1 | Step 2 |
|---|---|
| Train a Base (B) Model to reproduce Pythia | Fine Tune (FT) the B Model on Observables |

➔ Only access to hadron level information

➔ Only access to observables

# MAGIC
## (Microscopic Alterations Generated from IR Collections)
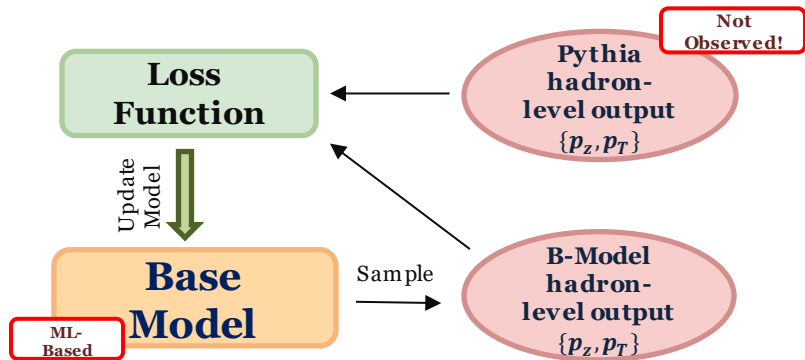
**Step 1**
Train a Base (B) Model to reproduce Pythia

➡ Only access to hadron level information

**Step 2**
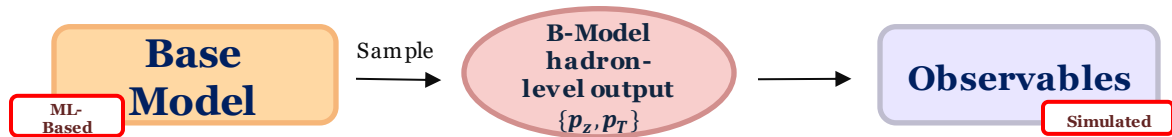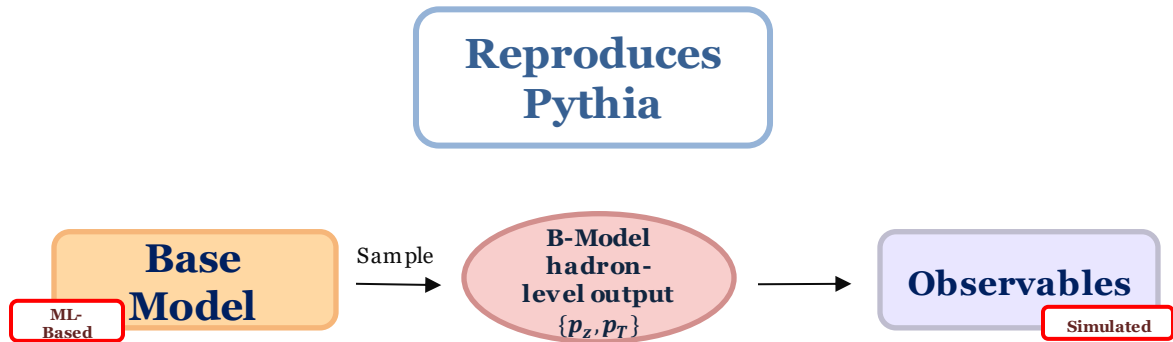Fine Tune (FT) the B Model on Observables

➡ Only access to observables

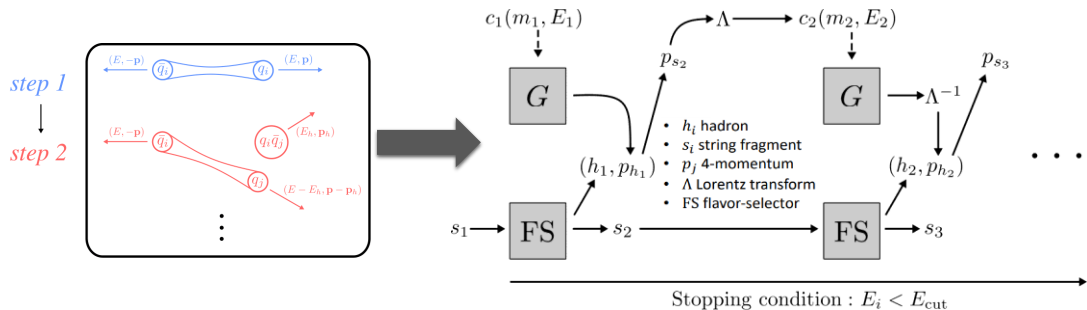# Step 1: Train Base (B) - Model on Pythia generated hadron-level output

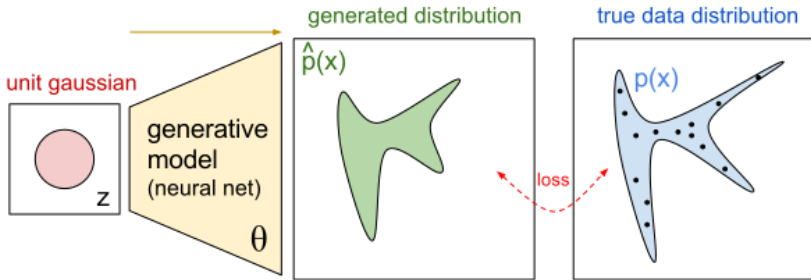**Step 1: Train Base (B) - Model on Pythia generated hadron-level output**



Base Model
*ML-Based*

→ Sample →

B-Model hadron-level output $\{p_z, p_T\}$

→

Observables
*Simulated*

**Step 1: Train Base (B) - Model on Pythia generated hadron-level output**



Reproduces Pythia

Base Model — ML-Based

Sample →

B-Model hadron-level output $\{p_z, p_T\}$

→ Observables — Simulated

$c_1(m_1, E_1)$ $\quad\quad\quad\quad \Lambda \longrightarrow c_2(m_2, E_2)$

- $h_i$ hadron
- $s_i$ string fragment
- $p_j$ 4-momentum
- $\Lambda$ Lorentz transform
- FS flavor-selector

Stopping condition : $E_i < E_{cut}$

**We need a generative model!**

| Sample hadron kinematics:<br>Train on $\{p_z, p_T\}$ | Emission of different Mesons:<br>Condition on mass ($m$) and energy ($E$) |

University of CINCINNATI
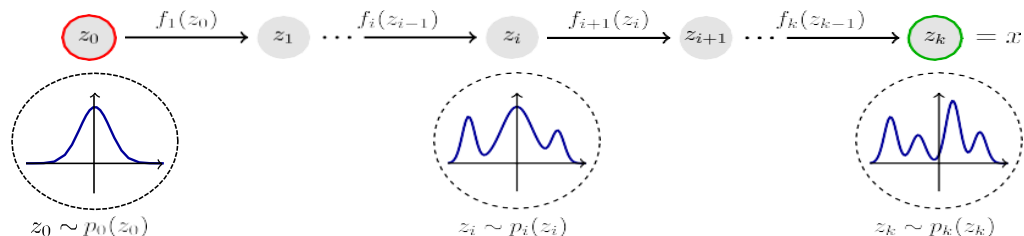
MLHAD

https://openai.com/research/generative-models



Source: generative models

$\Rightarrow$ Task: Learn the probability distribution $p(x)$ of the data

## Which generative model should we choose?

| Is it able to learn **complex distributions**? | Do we have access to the **exact probability distribution**? |

A. Youssef, Towards data-driven models of hadronization    youssead@ucmail.uc.edu

$$z_0 \xrightarrow{f_1(z_0)} z_1 \cdots \xrightarrow{f_i(z_{i-1})} z_i \xrightarrow{f_{i+1}(z_i)} z_{i+1} \cdots \xrightarrow{f_k(z_{k-1})} z_k = x$$

$z_0 \sim p_0(z_0)$        $z_i \sim p_i(z_i)$        $z_k \sim p_k(z_k)$

Zo - random vector
sampled from a
Gaussian $p_0(z_0)$

Fi − invertible NN that
transforms $p_0(z_0)$ to $p_i(z_i)$
by change of variables
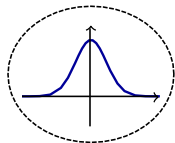
Complex target distribution
$p_k(z_k)$ is learned

⇒ **Can learn complex distributions!**

**Exact probability distribution is
obtained by change of variables**
$$p_k(z_k) = p_0(z_0) \prod_{i=1}^{K} \left| \det\left( \frac{\partial f_i(z_{i-1})}{\partial z_{i-1}} \right) \right|^{-1}$$

hHps://github.com/janosh/awe
some-normalizing-flows

⇒ **Access to the exact probability distribution**

University of CINCINNATI

MLHAD



$z_0$ ──── $f_1(z_0)$ ──── $f_i(z_{i-1})$ ──── $f_{i+1}(z_i)$ ──── $f_k(z_{k-1})$ ──── $z_k = x$

$z_0 \sim p_0(z_0)$

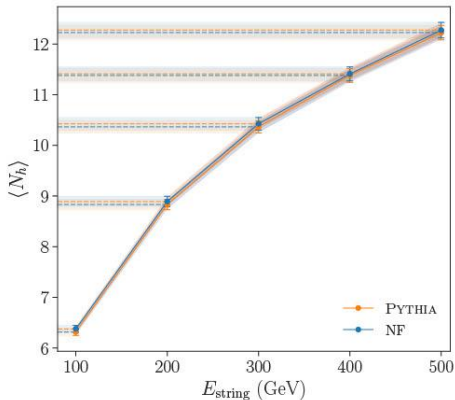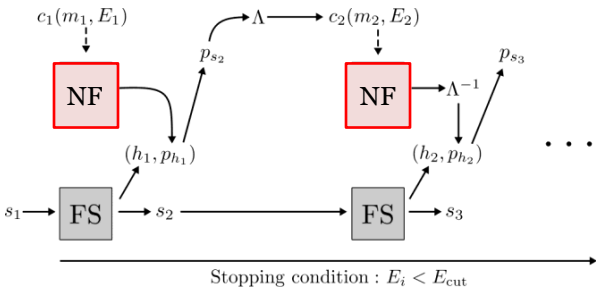$z_k \sim p_k(z_k)$

Zo- random vector sampled from a Gaussian $p_0(z_0)$

omplex target distribution $p_k(z_k)$ is learned

**learn complex distributions!**

Exact p
obtain

$\left. \dfrac{(z_{i-1})}{z_{i-1}} \right)^{-1}$

— NF
PYTHIA

Density

0.025
0.020
0.015
0.010
0.005
0.000

$p_z$ (GeV): 0.0 0.2 0.4 0.6 0.8 1.0

$m_\perp$ (GeV): 0.172 0.27 0.39 0.511

hHps://github.com/janosh/awe some-normalizing-flows

⇒**Access to the exact probability distribution**

A. Youssef, Towards data-driven models of hadronization

youssead@ucmail.uc.edu

**Implement NF in the fragmentation chain to obtain physical observables**



$\Rightarrow$ *Multiplicity obtained by MLHad agrees with Pythia!*

# MAGIC
**(M**icroscopic **A**lterations **G**enerated from **IR C**ollections)


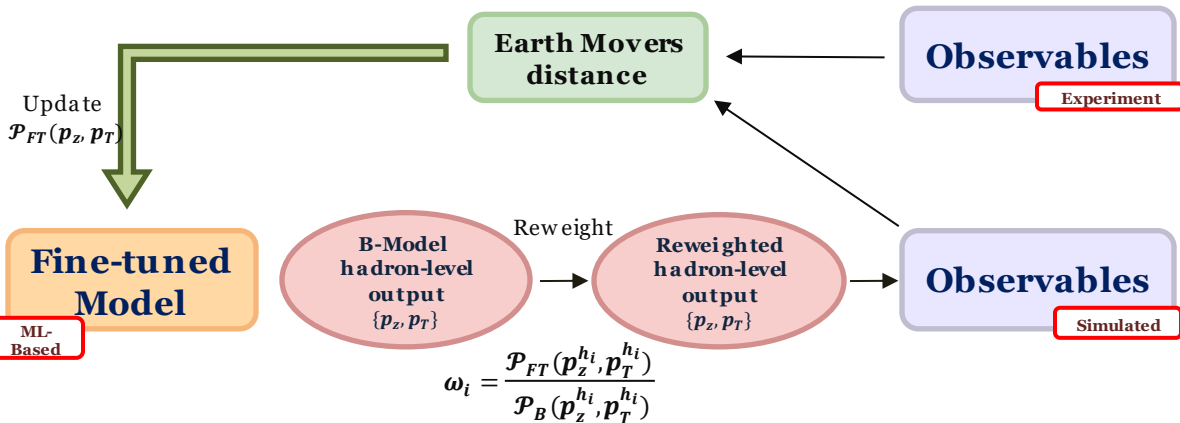
**Step 1**
Train a Base (B) Model
to reproduce Pythia
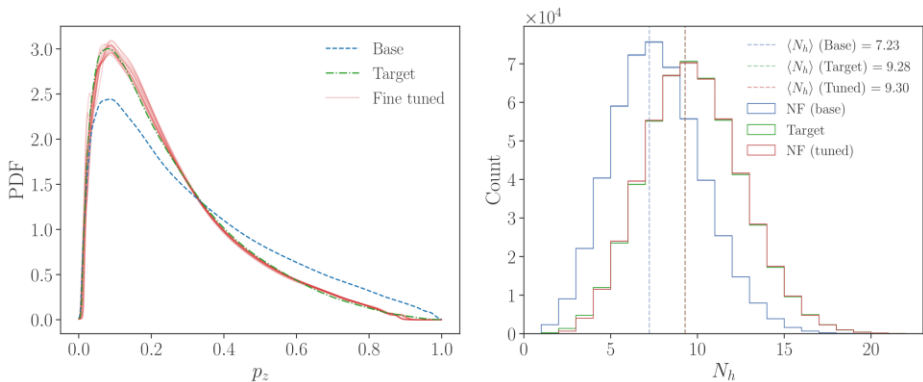
➡ Only access to hadron
level information

**Step 2**
Fine Tune (FT) the B
Model on Observables

➡ Only access
to observables

## Step 2: Fine-tune B-Model on physical observables



$$\omega_i = \frac{\mathcal{P}_{FT}(p_z^{h_i}, p_T^{h_i})}{\mathcal{P}_B(p_z^{h_i}, p_T^{h_i})}$$

**Base: Pythia default parameters**
**Target: Pythia perturbed; aLund=1.5**

A. Youssef, Towards data-driven models of hadronization

youssead@ucmail.uc.edu

# MAGIC is a very promising methods for data-driven hadronization models!

➜ Excellent results by training on only one observable (multiplicity)!

➜ More details on MAGIC and uncertainty quantification in **arXiv: 2311.09296**

## More MLHAD work

- Reweighting Monte Carlo Predictions and Automated Fragmentation Variations in Pyhia 8 (arXiv:2308.13459)

- Pythia Flavor Reweigthing (arXiv:24NN.NNNNN)

- Collective Reweighting Method - two part reweighter (arXiv:2407.XXXXX)

- Tuning Hadronization Models

**Project Homepage:**
**https://uchep.gitlab.io/mlhad-docs/**

# Backup

**Uncertainty estimation is crucial for event generator predictions!**

**„Classical" Neural Networks**



Weights have a fixed value
→ Weight values are updated in each epoch

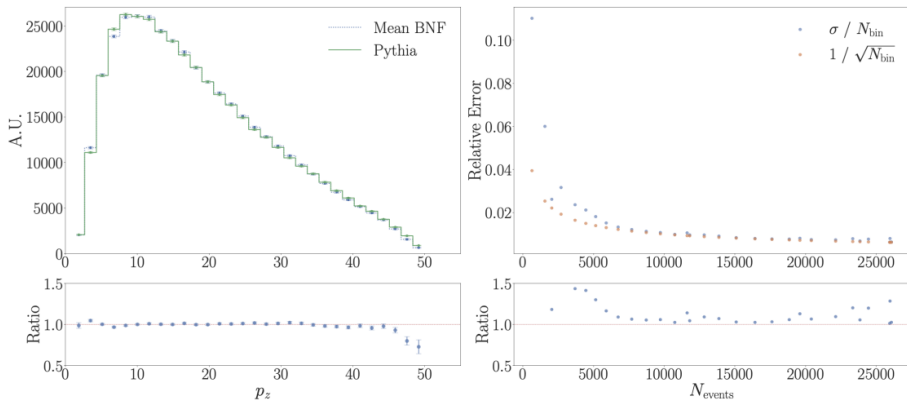University of CINCINNATI

MLHAD

## „Classical" Neural Networks



Weights have a fixed value
→ Weight values are updated in each epoch

## Bayesian Neural Networks (BNN)



Weights are sampled from a distribution
→ Distribution parameter are updated in each epoch

→ **BNN are easy to implement: Add additional loss function for weight distribution**

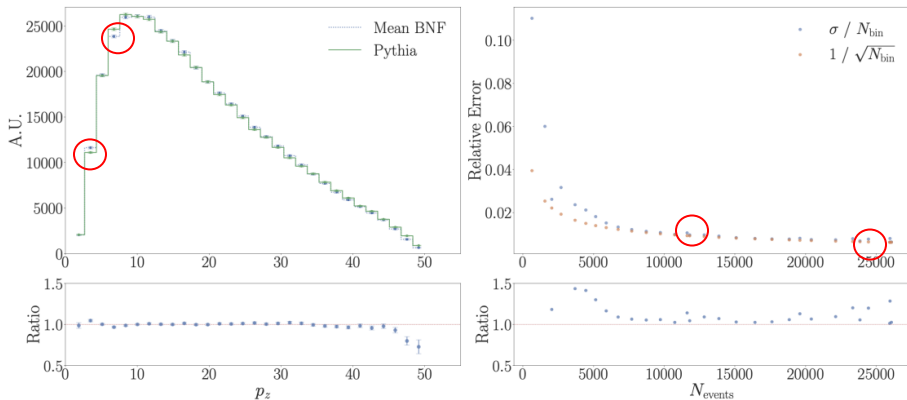→ **Capture statistical and training uncertainties**

(Image source: The very Basics of Bayesian Neural Networks )

**Pythia Sample**:
One sample with errors corresponding to $\sqrt{N_{bin}}$

**Mean BNF**:
$5 \times 10^5$ samples with errors corresponding to the standard deviation

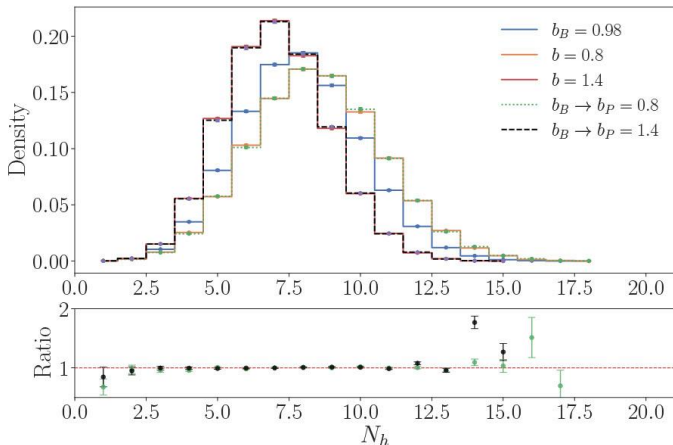**BNF capture the statistical and training uncertainties**

b is a free parameter in the Lund function used in Pythia: StringZ:bLund

Train nominal NF conditioned on different b → Get likelihood

→ Reweight nominal output using ratio of likelihoods:

$$w = \prod_i \frac{p_{nom}^{(i)}(z)}{p_{pert}^{(i)}(z)}$$
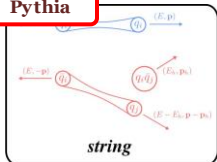
**When is a hadronization model successful?**

## When is a hadronization model successful?

➡ **The performance is judged by their description of experimental measurements!**

## When is a hadronization model successful?

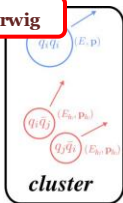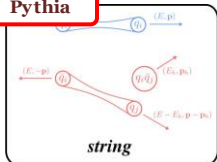➡ **The performance is judged by their description of experimental measurements!**

**Phenomenological Models (String, Cluster) are currently state of art and are overall very successful, however:**

➡ **comparison of data from proton-proton and ion-ion collision with Pythia**
   ➡ **discrepancies at the level of $O(20\%)$ to $O(50\%)$** N. Fischer and T. Sjöstrand, JHEP 01, 140 (2017), 1610.09818.

➡ **recovering collective effects can be challenging, for instance, heavy baryon production at high event multiplicities** Alice Collaboration, arXiv: 1807.11321

➡ **no efficient estimation of Uncertainties**

## When is a hadronization model successful?

➡ **The performance is judged by their description of experimental measurements!**
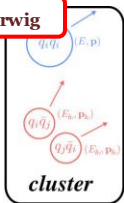

*string*
Pythia


*cluster*
Herwig

**Phenomenological Models (String, Cluster) are currently state of art and are overall very successful, however:**

➡ **comparison of data from proton-proton and ion-ion collision with Pythia**
  ➡ **discrepancies at the level of O(20%) to O(50%)** N. Fischer and T. Sjöstrand, JHEP 01, 140 (2017), 1610.09818.

➡ **recovering collective effects can be challenging, for instance, heavy baryon production at high event multiplicities** Alice Collaboration, arXiv: 1807.11321
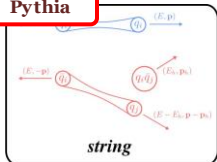
➡ **no efficient estimation of Uncertainties**

➡ **Both models have a discrepancy in describing experimental measurements!**

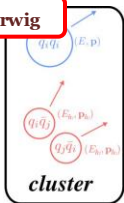University of
CINCINNATI

MLHAD

## When is a hadronization model successful?

➡ **The performance is judged by their description of experimental measurements!**

**Pythia**


*string*

**Phenomenological Models (String, Cluster) are currently state of art and are overall very successful, however:**

➡ comparison of data from proton-proton and ion-ion ~~collision~~ with Pythia
  ➡ discrepancies at the level of O(20%) to O~~...~~ and T. Sj̈ostrand, [**~~01, 140 (2017), 1610.09818~~**]

➡ recovering collective effects ca~~...~~ng, for instance, heavy baryon production at high event~~...~~  **Alice Collaboration, [arXiv: 1807.11321]**

➡ no efficient estim~~...~~ncertainties
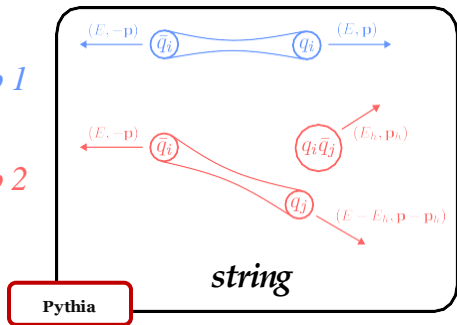
**Herwig**


*cluster*

*We need an innovative approach!*

➡ **Both models have a discrepancy in describing experimental measurements!**
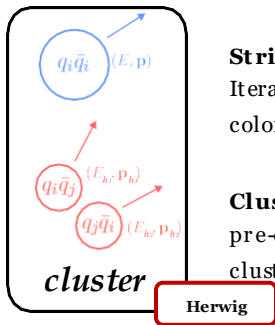
**Two primary hadronization models are used**



**String model:**
Iteratively split parton connected by QCD color strings with linear potential

**Cluster model:**
pre-confine partons into proto-clusters, then split by two-body decays

MLhad: Ilten, Menzo, Youssef, Zupan, 2203.04983, https://gitlab.com/uchep/mlhad

HadML: (Chan, Ghosh,) Ju, (Kania), Nachman, (Sangli,) Siodmok, 2203.12660, 2305.17169

University of
CINCINNATI

MLHAD

Uncertainty estimation is crucial for event generator predictions!

A. Youssef, Towards data-driven models of hadronization

youssead@ucmail.uc.edu

**Uncertainty estimation is crucial for event generator predictions!**

→ **Hard matrix element**

→ **Parton shower**

**Efficient solutions exist!**

**perturbative calculations depend on choices of scale, values of gauge and other couplings, particle masses, and nonperturbative inputs**

Giele et al, Phys. Rev. D84, 054003 (2011)

S. Mrenna and P. Skands, Phys. Rev. D94(7), 074005 (2016)

University of
CINCINNATI

MLH▵D

> **Uncertainty estimation is crucial for event generator predictions!**

➡ **Hard matrix element**

➡ **Parton shower**

**Efficient solutions exist!**

**perturbative calculations depend on choices of scale, values of gauge and other couplings, particle masses, and nonperturbative inputs**

Giele et al, Phys. Rev. D84, 054003 (2011)
S. Mrenna and P. Skands, Phys. Rev. D94(7), 074005 (2016)

➡ **Hadronization**

**Efficient solution has remained elusive!**

**Standard procedure: perform repeated simulations with different sets of values for the model parameters**

➡ **Computationally very expensive!**

University of
CINCINNATI

MLHaD

> **Uncertainty estimation is crucial for event generator predictions!**

➡ **Hard matrix element**

➡ **Parton shower**

**Efficient solutions exist!**

**perturbative calculations depend on choices of scale, values of gauge and other couplings, particle masses, and nonperturbative inputs**

Giele et al, Phys. Rev. D84, 054003 (2011)
S. Mrenna and P. Skands, Phys. Rev. D94(7), 074005 (2016)

➡ **Hadronization**

**Efficient solution has rem̶a̶i̶n̶ed̶ ̶e̶l̶u̶si̶ve!**

**Standard procedure̶s̶ ̶r̶e̶peated simulations with different s̶e̶t̶s̶ ̶f̶or the model parameters**

̶c̶o̶mputationally very expensive!

*Need a more efficient way!*

**Small Detour: No ML, only Had**

### Reweighting Monte Carlo Predictions and Automated Fragmentation Variations in PYTHIA 8

Christian Bierlich[1♠], Phil Ilten[2†], Tony Menzo[2★], Stephen Mrenna[2,3✦], Manuel Szewc[2∥], Michael K. Wilkinson[2⊥], Ahmed Youssef[2‡], and Jure Zupan[2§]

[1] Department of Physics, Lund University, Box 118, SE-221 00 Lund, Sweden
[2] Department of Physics, University of Cincinnati, Cincinnati, Ohio 45221, USA
[3] Scientific Computing Division, Fermilab, Batavia, Illinois, USA

♠ christian.bierlich@hep.lu.se, † philten@cern.ch, ★ menzoad@mail.uc.edu, ✦ mrenna@fnal.gov, ∥ szewcml@ucmail.uc.edu, ⊥ michael.wilkinson@uc.edu, ‡ youssead@ucmail.uc.edu, § zupanje@ucmail.uc.edu
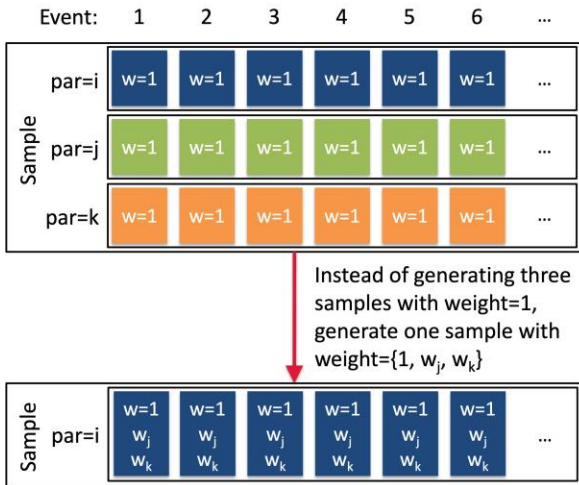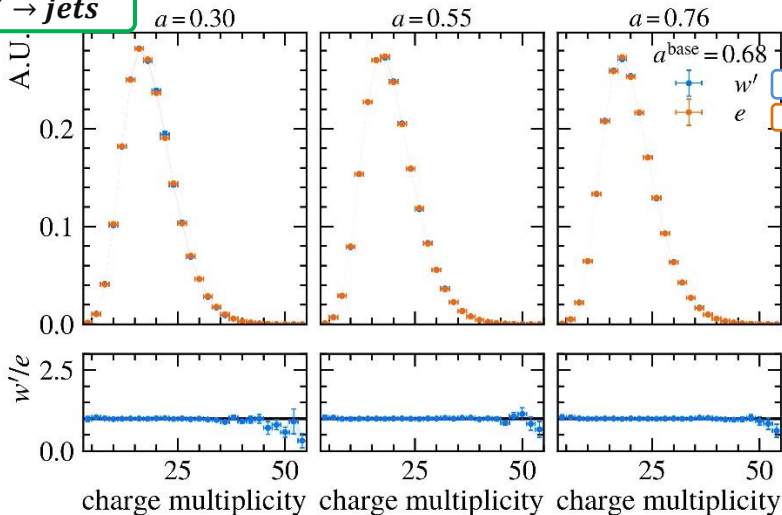
MLHAD

### Abstract

This work reports on a method for uncertainty estimation in simulated collider-event predictions. The method is based on a Monte Carlo-veto algorithm, and extends previous work on uncertainty estimates in parton showers by including uncertainty estimates for the Lund string-fragmentation model. This method is advantageous from the perspective of simulation costs: a single ensemble of generated events can be reinterpreted as though it was obtained using a different set of input parameters, where each event now is accompanied with a corresponding weight. This allows for a robust exploration of the uncertainties arising from the choice of input model parameters, without the need to rerun full simulation pipelines for each input parameter choice. Such explorations are important when determining the sensitivities of precision physics measurements. Accompanying code is available at gitlab.com/uchep/mlhad-weights-validation.
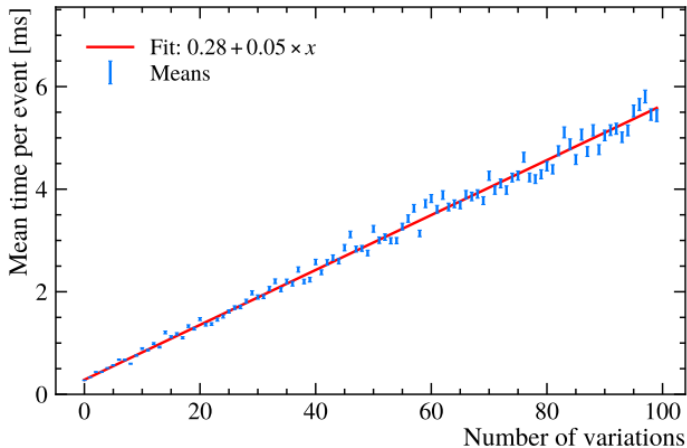
University of
CINCINNATI



- **Event generation is time consuming**
  - **We want to reweight events without regenerating**

- **Use a modified veto algorithm**
  - **New event weights for different hadronization param are book kept**

- **We calculate event weights for different hadronization options in a single event generation!**

Instead of generating three samples with weight=1, generate one sample with weight={1, $w_j$, $w_k$}

$e^+e^- \rightarrow Z \rightarrow jets$

- **Generate 100 samples with different variations of aLund**

- **Each sample has 1000 events**

- **Cost per additional parameter variation is around 0.05 ms**

- **We have a speed up by a factor ~3**

## Variational Autoencoder (VAE)
Kingma et al, [arXiv:1312.6114](#)



KL-divergence limits the latent space to a simple analytic distribution

**Vanilla VAE**

**VAE latent space**
[arXiv: 1804.01947](#)

# Generative Models
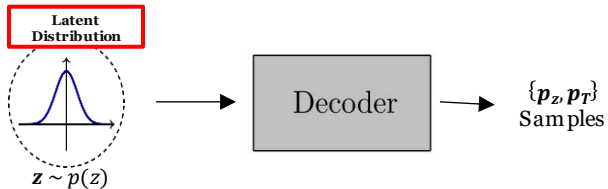
## Variational Autoencoder (VAE)
### Kingma et al, arXiv:1312.6114

Vanilla VAE

KL-divergence limits the latent space to a simple analytic distribution
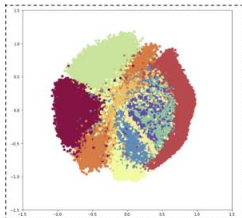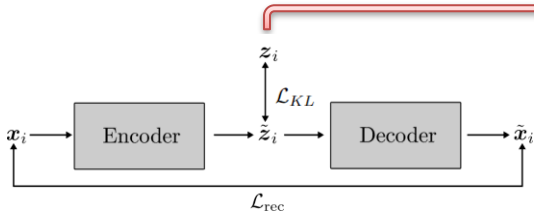
VAE latent space
arXiv: 1804.01947

### Inference

Latent Distribution

$z \sim p(z)$

Decoder

$\{p_z, p_T\}$ Samples

University of CINCINNATI

MLHAD

**Variational Autoencoder (VAE)**
Kingma et al, arXiv:1312.6114



**Vanilla VAE**

KL-divergence limits the latent space to a simple analytic distribution

**VAE latent space**
arXiv: 1804.01947

**Complex input data**

**Simple latent space**

Encoder

Decoder

⇒ **Complex distribution are hard to learn!**

# Generative Models

## Variational Autoencoder (VAE)
Kingma et al, arXiv:1312.6114

**Vanilla VAE**

**Complex input data**

KL-divergence limits the latent space to a simple analytic distribution

VAE latent space
arXiv: 1804.01947

**How can we make VAEs learn more complex distribution?**
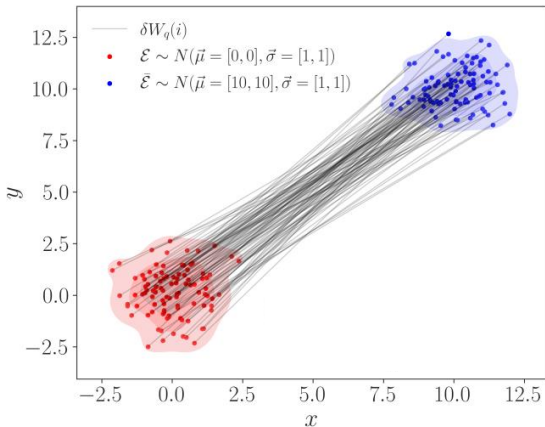
⇒ Complex distribution are hard to learn!

**Use Sliced Wasserstein Distance as latent loss function!**

**Use Sliced Wasserstein Distance as latent loss function!**

**Wasserstein distance (WD)**

$$W_q(\mathcal{E}, \bar{\mathcal{E}}) = \left[ \min_{\{f_{ij} \geq 0\}} \sum_{i=1}^{N} \sum_{j=1}^{\bar{N}} f_{ij} \left( \hat{d}_{ij} \right)^q \right]^{1/q}$$
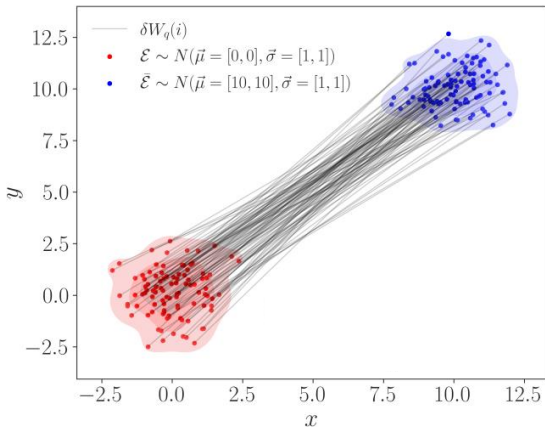
## Use Sliced Wasserstein Distance as latent loss function!

### Wasserstein distance (WD)

$$W_q(\mathcal{E}, \bar{\mathcal{E}}) = \left[ \min_{\{f_{ij} \geq 0\}} \sum_{i=1}^{N} \sum_{j=1}^{\bar{N}} f_{ij} \left(\hat{d}_{ij}\right)^q \right]^{1/q}$$
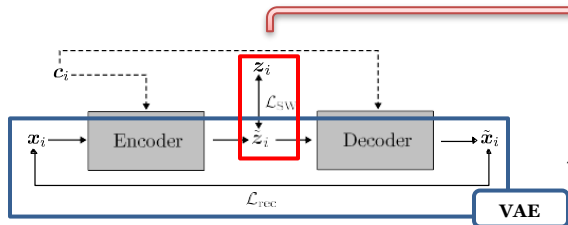
### Sliced Wasserstein distance

➡ **Projects high dimensional data into one dimensional "slices"**

➡ **WD in 1D has a closed form solution**

  ➡ **Sorted Difference of the two samples**

## Conditional Sliced Wasserstein (SW) Autoencoder (cSWAE)
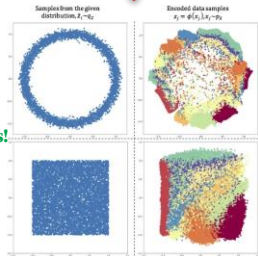
**Restricted to Pion emissions**



**cSWAE architecture**
(Architecture used in SciPost Phys. 14, 027 (2023) )

**SW distance enables learning any sampleable latent distribution**
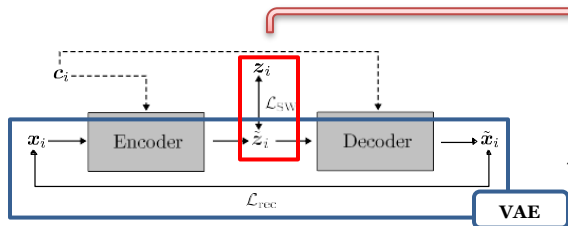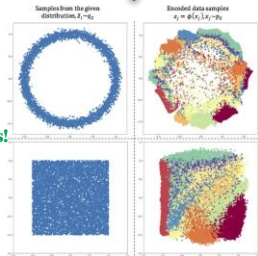
⇒ **Can learn complex distributions!**



**SWAE latent space**
(arXiv: 1804.01947 )

**Conditional Sliced Wasserstein (SW) Autoencoder (cSWAE)**

**Restricted to Pion emissions**



**SW distance enables learning any sampleable latent distribution**

⇒ **Can learn complex distributions!**

**cSWAE architecture**
(Architecture used in SciPost Phys. 14, 027 (2023) )

**SWAE latent space**
(arXiv: 1804.01947 )



$\{p_z, p_T\}$ Samples

$z \sim p(z)$

**Latent Distribution**

**Decoder "just" generates samples**

⇒ **No access to the probability distribution**

$z_0 \xrightarrow{f_1(z_0)} z_1 \cdots \xrightarrow{f_i(z_{i-1})} z_i \xrightarrow{f_{i+1}(z_i)} z_{i+1} \cdots \xrightarrow{f_k(z_{k-1})} z_k = x$

$z_0 \sim p_0(z_0)$

$z_i \sim p_i(z_i)$

$z_k \sim p_k(z_k)$

Zo- r andom vector
sampled from a
Gaussian $p_0(z_0)$

Fi − invertible NN that
tr ansforms $p_0(z_0)$ to $p_i(z_i)$
by change of variables

Complex target distribution
$p_k(z_k)$ is l earned

⇒ **Can learn complex distributions!**

**Exa ct probability distribution is
obtained by change of variables**

$$p_k(z_k) = p_0(z_0) \prod_{i=1}^{K} \left| \det\left( \frac{\partial f_i(z_{i-1})}{\partial z_{i-1}} \right) \right|^{-1}$$

⇒ **Access to the exact probability distribution**

**Removed pion
emission restriction**

University of
CINCINNATI

MLHAD

�ड **Propagation of errors**

 ➔ **ML architecture with Bayesian Normalizing Flows (presented in part)**

➔ **Train on observables only**

 ➔ **Two part reweighter (not part of the talk)**

 ➔ **Train on global observables with Fine tuning (results not shown in this talk)**

➔ **To train on experimental data**

 ➔ **Want fast evaluation of parameter dependency**

 ➔ **Use reweighting method**

 ➔ **First implementation in Pythia for Lund string model (to be released soon in Pythia)**