

MODELING HADRONIZATION USING MACHINE LEARNING

Technion HEP Journal Club

Tony Menzo

University of Cincinnati

May 10, 2022

Based on [2203.04983](#) with Phil Ilten, Ahmed Youssef, and Jure Zupan

Monte Carlo Event Generators

- Currently there is no phenomenological model of hadronization which correctly reproduces all experimental data.
- No ‘new’ models of fragmentation in ~ 30 years

Event generators are ubiquitous in HEP:

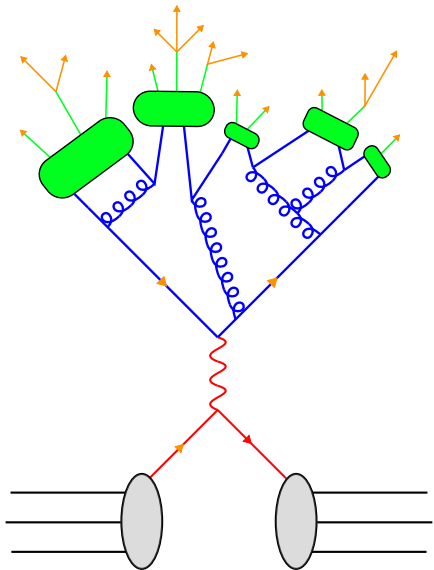
Citation data according to INSPIRE

Event generator	hep-ex (Citations)	hep-ph (Citations)
PYTHIA 6	8,314	4,459
PYTHIA 8	4,528	2,100
HERWIG 6	2,499	1,103
HERWIG++	1,906	971
SHERPA	2,569	1,073

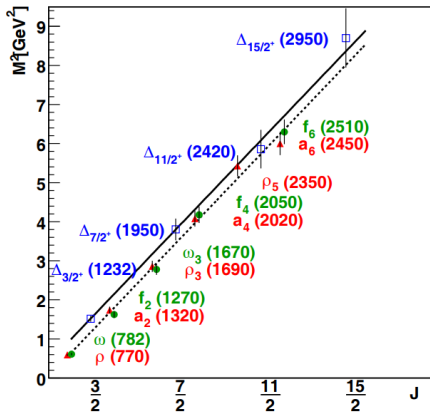
Simulation of hard process

The simulation is ordered according to the magnitude of momentum transfer:

1. **Hard interaction**
2. **Parton Showers**
3. **Hadronization**
4. **Unstable particle decay**



QCD, Confinement, and Strings



Leading Regge trajectories $J = L + S$ [1].

Regge Trajectories:

$$J = \alpha_0 + \alpha' M^2$$

Nambu suggested a stringy origin: consider a $q\bar{q}$ pair connected by a rigid rod with linear energy density κ

$$J = \frac{M^2}{2\pi\kappa} \quad (1)$$

Define $\alpha' \equiv 1/2\pi\kappa \sim 0.9$
GeV⁻²

QCD, Confinement, and Strings

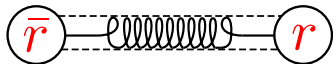
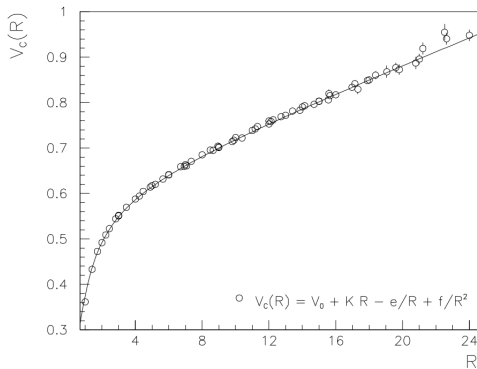
Around Λ_{QCD} particles charged under QCD begin to confine into color neutral hadrons. Separations > 1 fm generate thin $\mathcal{O}(1/\Lambda_{\text{QCD}})$ color flux tubes between charged particles.

The $q\bar{q}$ potential is given by

$$V_{\text{QCD}}(r) \approx -\frac{4}{3} \frac{\alpha_s}{r} + \kappa r$$

with $\kappa \approx 1$ GeV/fm

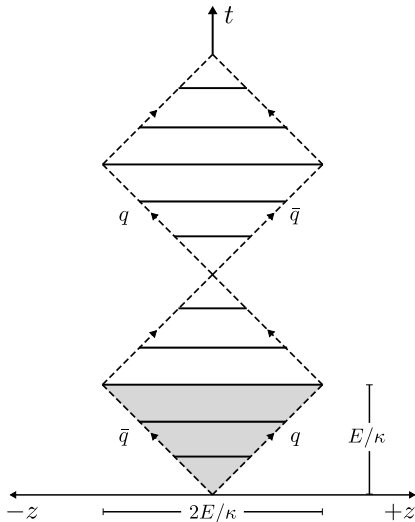
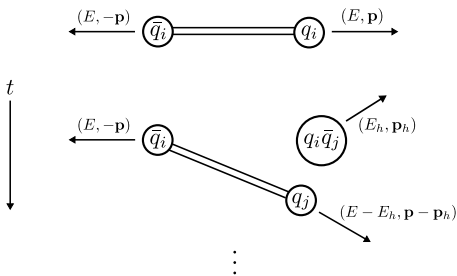
- Motivates a classical description at large separations



Confirmation from quenched lattice QCD simulation. arXiv:hep-lat/9210003

Lund Model: *Classical* model of hadronization utilizing string picture. Assume $r \gg 1 \rightarrow$

$$V_{\text{Lund}} = \kappa r$$

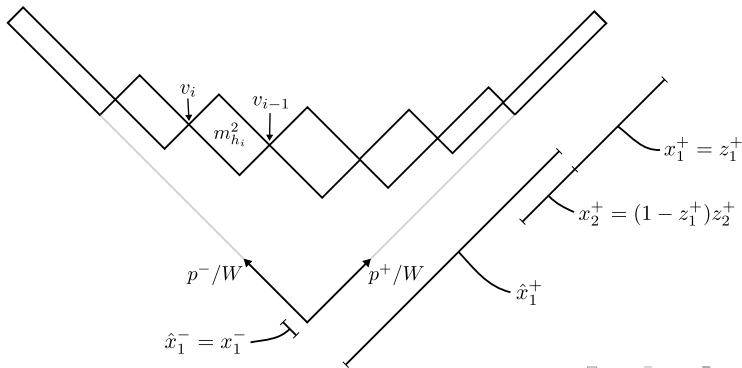


Yo-yo string motion in the absence of string breaks as seen in the center of mass frame of the string system.

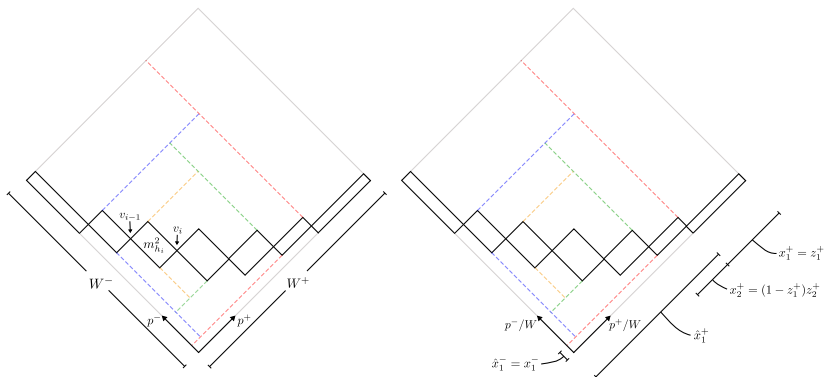
- Fragmentation is implemented as a stochastic process through production vertices in momentum space.
- Two inputs: mass of hadron m_h and longitudinal momentum fraction z . Left-right symmetric Lund fragmentation function:

$$f(z)dz \propto \frac{(1-z)^a}{z} \exp\left(-b\frac{m_h^2}{z}\right) dz$$

where a and b are fit parameters.

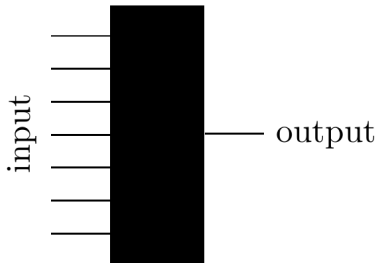


- Each iteration is a ‘scaled down’ replica of the initial system.



'Primer' to Machine Learning

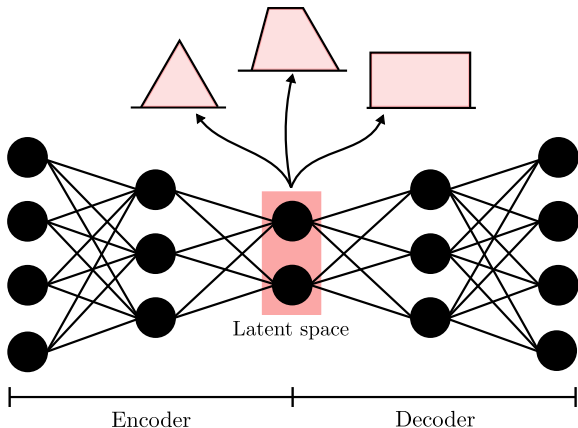
Gross oversimplification: A black box which takes input, performs manipulations, and returns output.



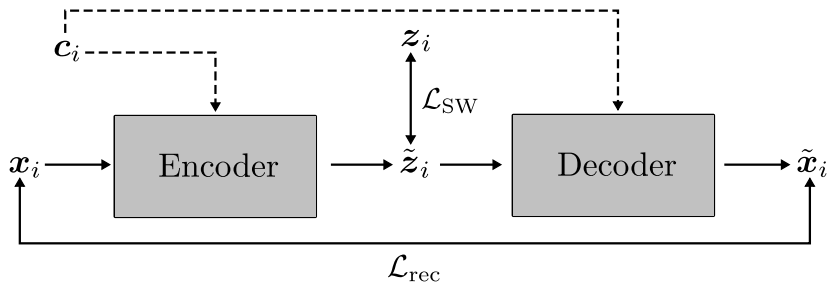
For theorists machine learning algorithms should be thought of as a powerful tool!

Sliced-Wasserstein Autoencoder (SWAE)

Goal: Input number sampled from a 'simple' probability distribution and return a number distributed according to the hadronic kinematic distribution



Conditional SWAE

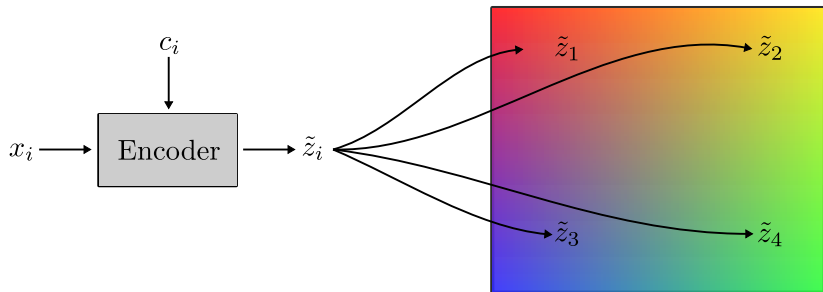


Program architecture for our implementation of the SWAE.

$$\mathcal{L}_{\text{Total}} = \mathcal{L}_{\text{SW}} + \mathcal{L}_{\text{rec}}.$$

Tunable pParameters: learning rate, latent dimension, regularization parameter, number of slices, number of epochs.

Conditional SWAE

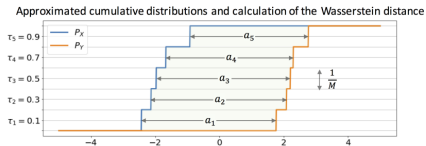
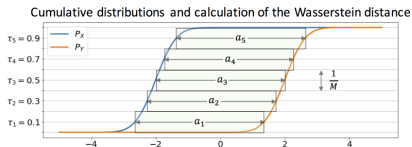
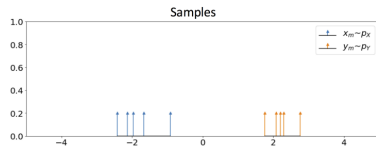
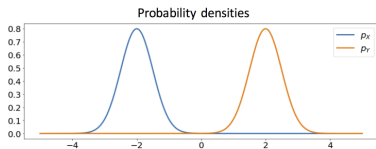


Condition dependent latent space.

$$c_i = \frac{E_{\max} - E_i}{E_{\max} - E_{\min}} \quad (2)$$

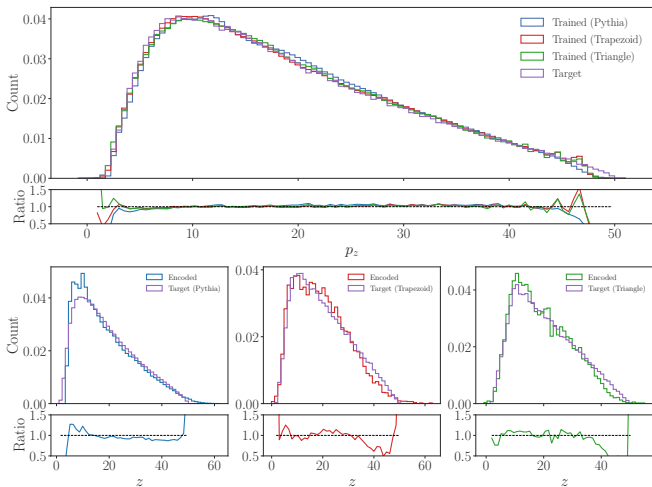
Wasserstein or 'earth movers' distance

$$f_{\text{WD}}(p_1, p_2) \rightarrow \# \in \mathbb{R} \quad (3)$$



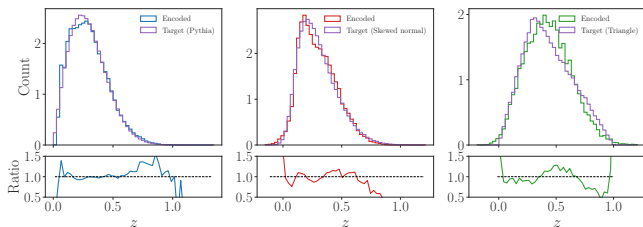
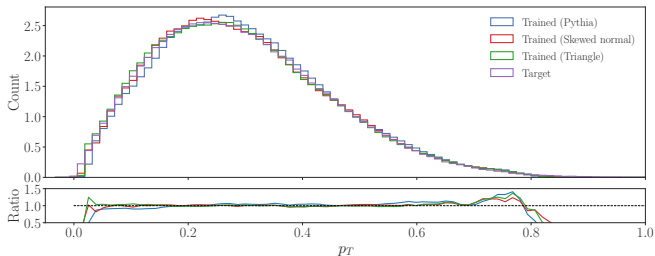
Wasserstein distance for 1D PDF [2].

SWAE-trained models: p_z



Pion first emission z -component of momentum distribution compared with encoder/decoder results from three SWAE-trained models.

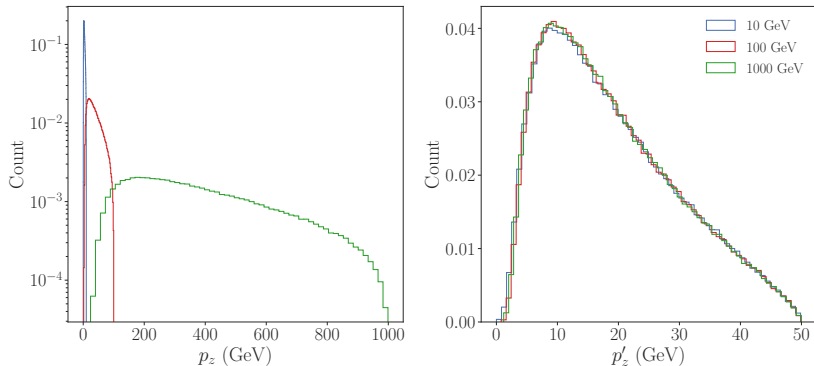
SWAE-trained models: p_T



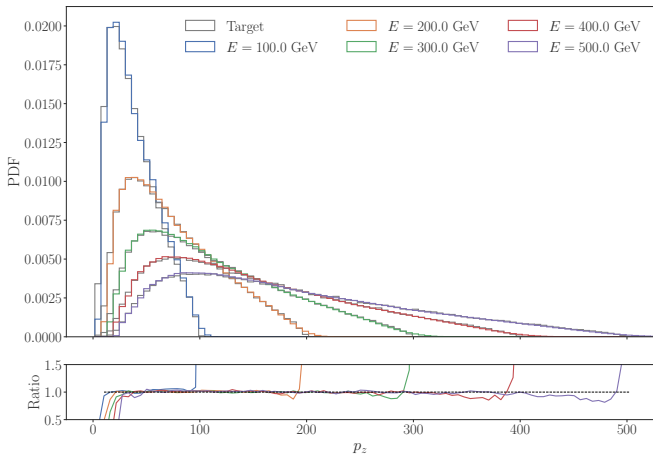
Pion first emission transverse momentum distribution compared with encoder/decoder results from three SWAE-trained models.

SWAE-trained models: Labels

Dependence on initial parton energy can be resolved - linear rescaling. Pythia output:



SWAE-trained models: Label-dependent p_z



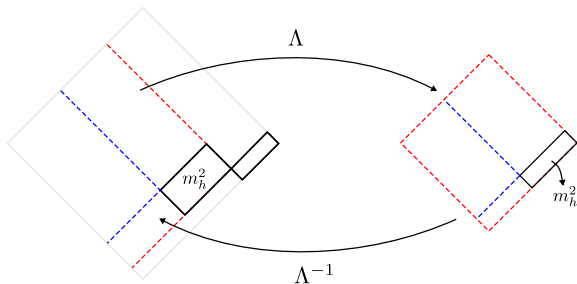
Inclusive first emission magnitude of momentum distribution compared with encoder/decoder results from three SWAE-trained models.

Modified Lund model

We require a change of inputs $(m_h, z) \rightarrow (E_{\text{CM}}, \tilde{z})$ and rely on:

1. Causal disconnection of events
2. Simple rescaling of the kinematic distributions with respect to energy. For example,

$$p'_z = p_{z,\text{ref}} \frac{p_z}{E} \quad (4)$$



Issues

- Production pairs have transverse momenta distributed according to the semi-classical tunneling probability

$$f(p_{\perp}) \propto \exp(\pi p_{\perp}^2 / \kappa)$$

- Replace $m_h^2 \rightarrow m_{\perp}^2 = m_h^2 + p_{\perp}^2$

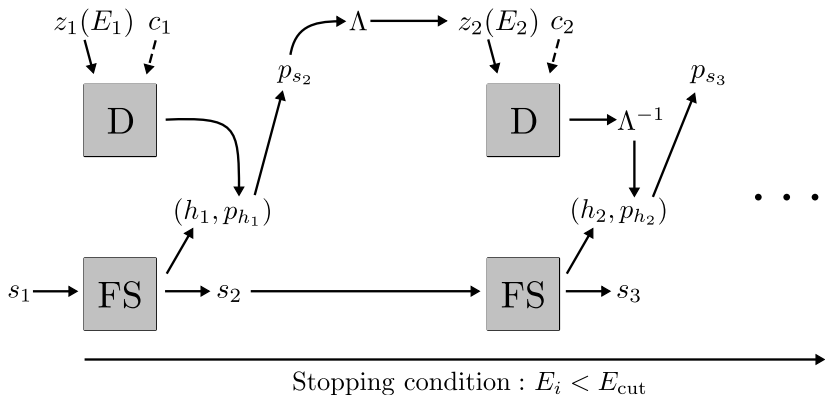
This creates two issues:

1. Our fragmentation chain always samples from the CM frame whereas the transverse momenta are generated in the lab frame.
2. The dependence on $f(z)$ on m_{\perp} creates a correlation between z and p_{\perp} through each iteration.

Approximate solution

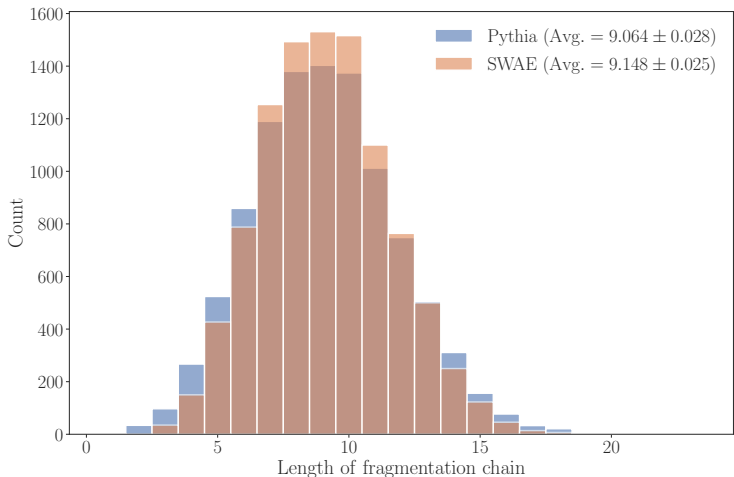
- Sample p_{\perp} from distribution above
- Train on p_z distributions where the hadron mass is increased by the variance of the p_{\perp} distribution
- Ideally train on m_T -labeled dataset

Fragmentation chain



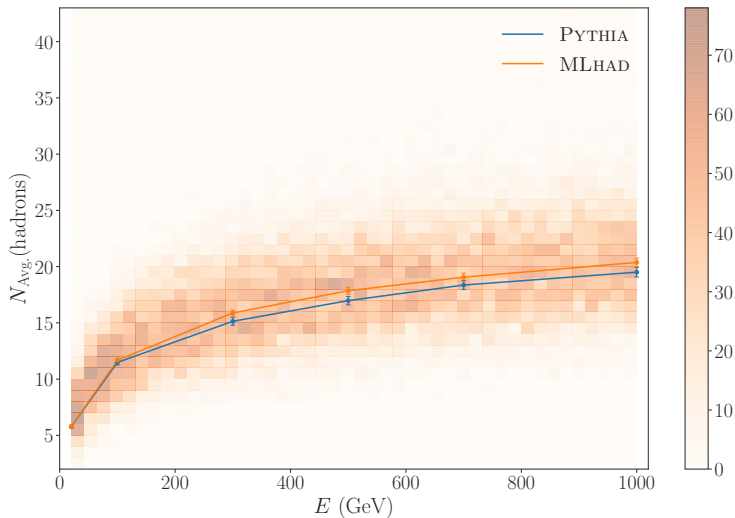
Program architecture for our implementation of the fragmentation chain.

Fragmentation Multiplicities



Simplified hadronization chain, only includes $u(\bar{u})$ and $d(\bar{d})$ quarks as string ends and pion (π^0, π^\pm) final states

Scaling with energy



Consistency check: the total number of hadrons should scale $\sim \ln E$.

Public code MLHAD!



Open-source!

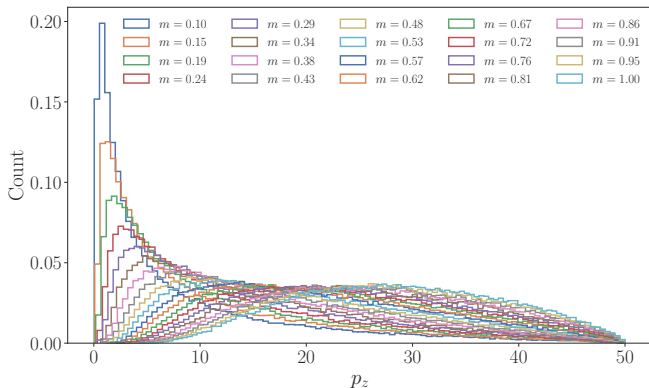
All the code written for this project is publically available in the following repo:

<https://gitlab.com/uchep/mlhad>

What's next? Immediate tasks

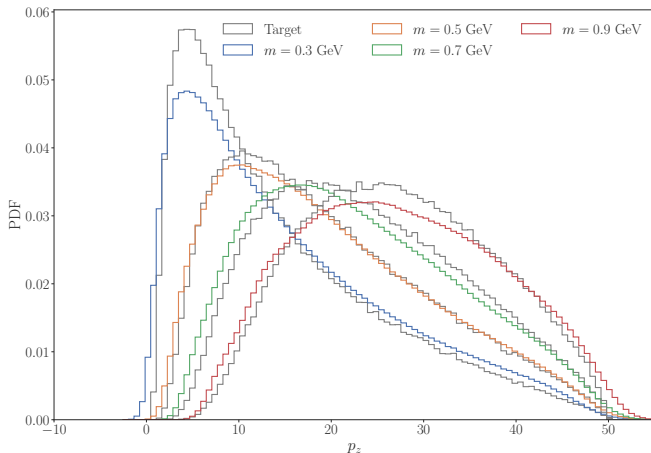
The model must be generalized to fully accommodate the Pythia hadronization model including all mesons and baryon production. Switch from one-dimensional training to multidimensional training.

Include mass-label!



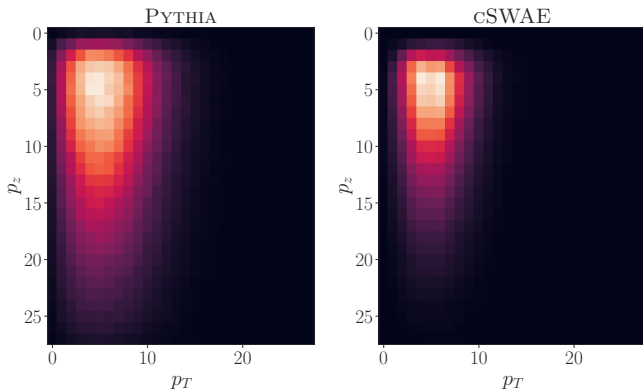
What's next? Immediate tasks

Super-preliminary results: ~ 100 epochs of training.



What's next? Immediate tasks

Super-preliminary results: 2D training, ~ 100 of training.



What else?

Our main goal is to train on data. To achieve this we must further develop the architecture and collect/develop ‘clean’ datasets.

Requires a slight shift in mindset: training on particle flow data and looking at global observables such as particle multiplicity and other jet observables.

Other future directions

1. A new model of flavor selection via machine learning
2. Additional event topologies: $gg \rightarrow ggh$, $qq\bar{q} \rightarrow qq\bar{q}h$, ...
3. Integration with FASTJET to test model consequences in jet observables
4. Feasibility of collecting training datasets from jet data (CMS Open Data)
5. Lattice QFT applications?

Conclusion

- Machine learning can be a valuable tool in modeling hadronization
- Ultimate goal: new generative hadronization paradigm - modeling directly from data rather than fitting to it to produce a complete phenomenological model of hadronization.
- Public code available! <https://gitlab.com/uchep/mlhad>
- See also [2203.04983](#) for more details ☺



References

1. Klempt et al “Multiplet classification of light-quark baryons.” The European Physical Journal A 48.9 (2012): 1-15.
2. Kolouri, Soheil, et al. “Sliced Wasserstein auto-encoders.” International Conference on Learning Representations. 2018.