

Data-driven hadronization models

Rutgers NHEC Seminar
November 19th, 2024

Tony Menzo

PhD candidate, University of Cincinnati

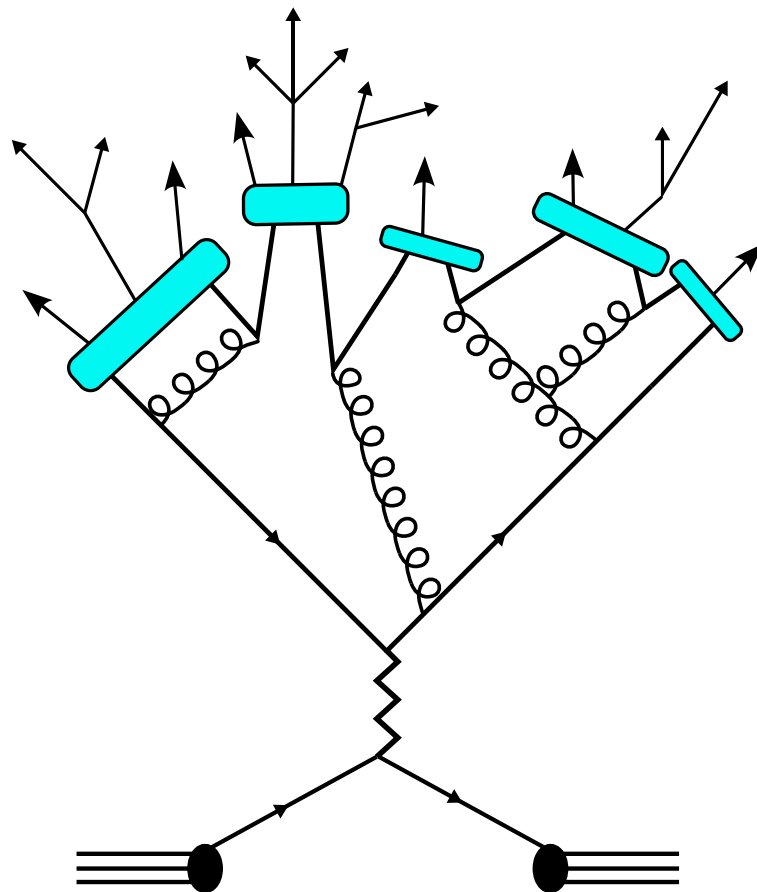
In collaboration with:

MLHAD

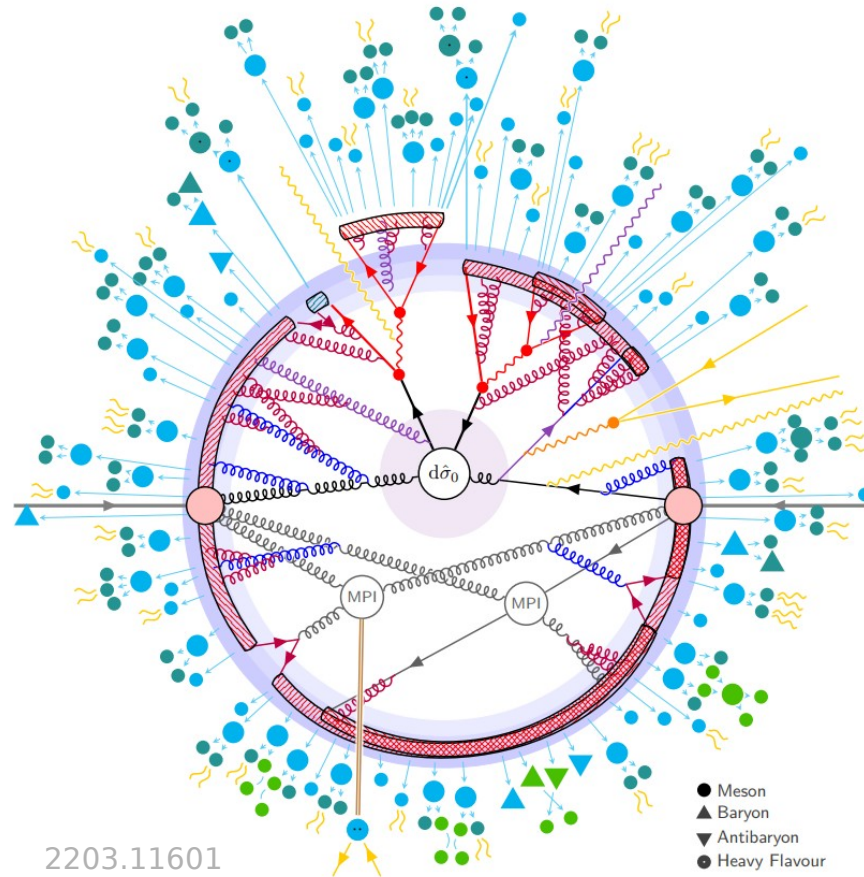
Christian Bierlich, Phil Ilten, Stephen Mrenna, Manuel Szwec,
Michael Wilkinson, Ahmed Youssef, Jure Zupan,

Nick Heller, Ben Nachman, Andrzej Siodmok

Based upon work in [2203.04983](#), [2308.13459](#), [2311.09296](#), [2410.06342](#), [2411.02194](#)



Event generators

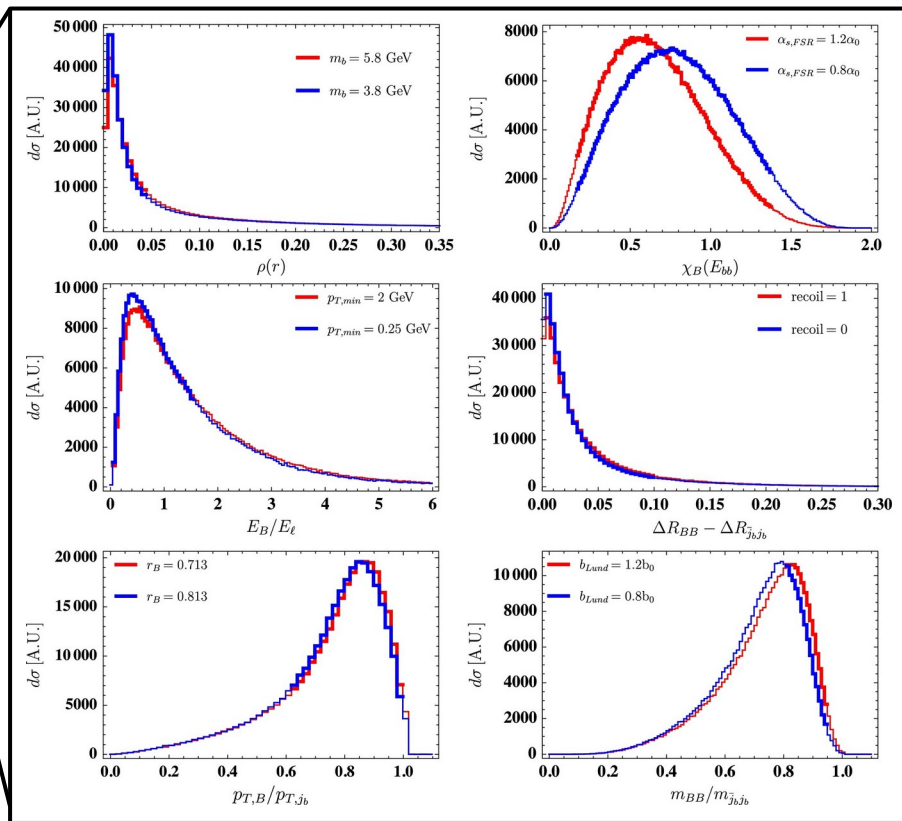


2203.11601

Motivation

Precision (exclusive) measurements dependent on hadronization!

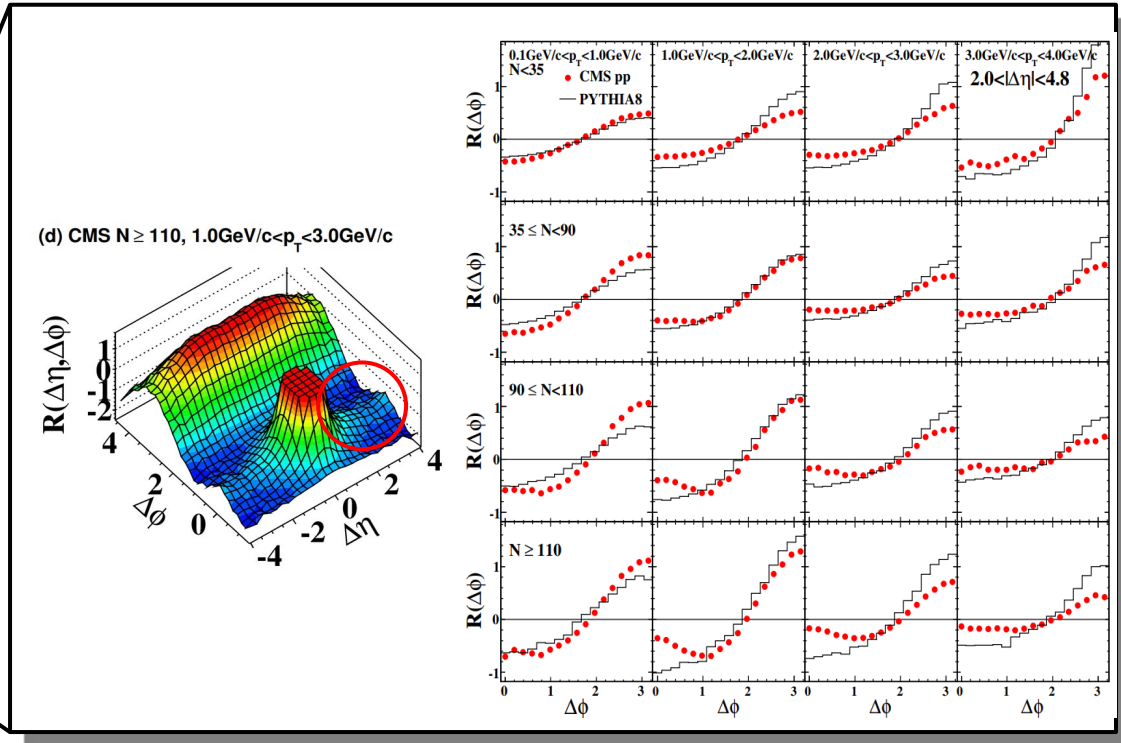
- Uncertainty reduction
 - Top quark mass measurement (r_b)
 - e^+e^- determination of α_s
- Mis-modeling
 - High-multiplicity events
 - Tuning discrepancies



Motivation

Precision (exclusive) measurements dependent on hadronization!

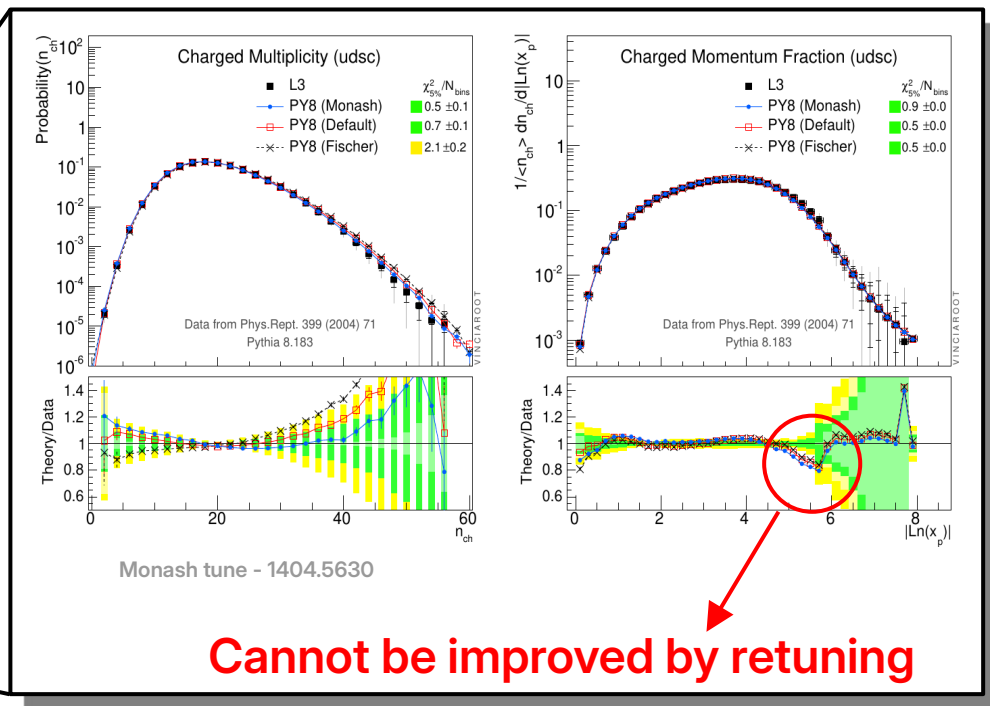
- Uncertainty reduction
 - Top quark mass measurement (r_b)
 - e^+e^- determination of α_s
- Mis-modeling
 - High-multiplicity events
 - Tuning discrepancies



Motivation

Precision (exclusive) measurements dependent on hadronization!

- Uncertainty reduction
 - Top quark mass measurement (r_b)
 - e^+e^- determination of α_s
- Mis-modeling
 - High-multiplicity events
 - Tuning discrepancies



Three roads to improvement

- **Improve model**

- MPIs, rope hadronization, transverse mass suppression, flavor asymmetries, hadronic rescattering, multiscale models (string → hydrodynamical), flavor selector, etc.

Hard to come up with mathematically precise model without established calculational techniques

- **Tend towards model independence**

- Sample directly from global distributions

Non-universal and extremely difficult to convert into **representative** particle flow data (uninterpretable)

- **Improve observables**

- Come up with better, hadronization sensitive (IR unsafe), observables

Non-trivial to come up with new observables given current data resolution

Three* roads to improvement

- Improve model

- MPIs, rope hadronization, transverse mass suppression, flavor asymmetries, hadronic rescattering, multiscale models (string → hydrodynamical), flavor selector, etc.

Hard to come up with mathematically precise model without established calculational techniques

- Tend towards model independence

- Sample directly from global distributions

Non-universal and extremely difficult to convert into **representative** particle flow data (uninterpretable)

- Improve observables

- Come up with better, hadronization sensitive (IR unsafe), observables

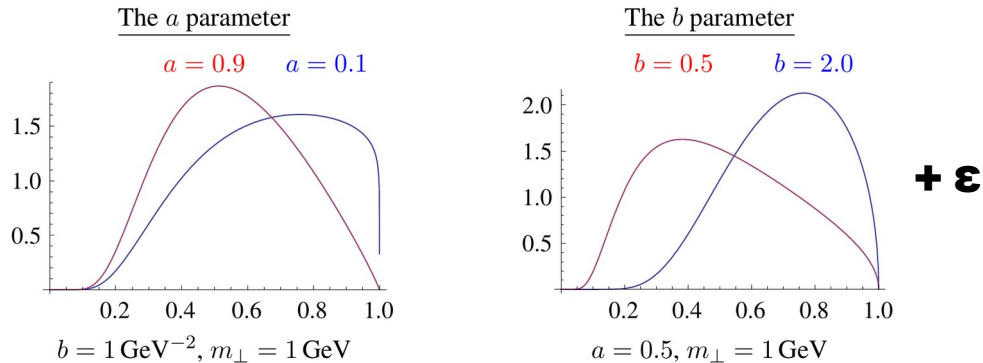
Non-trivial to come up with new observables given current data resolution

*or some combination of all three

Hybrid data-driven approach

The phenomenological models of hadronization already give an acceptable description of a large amount of data.

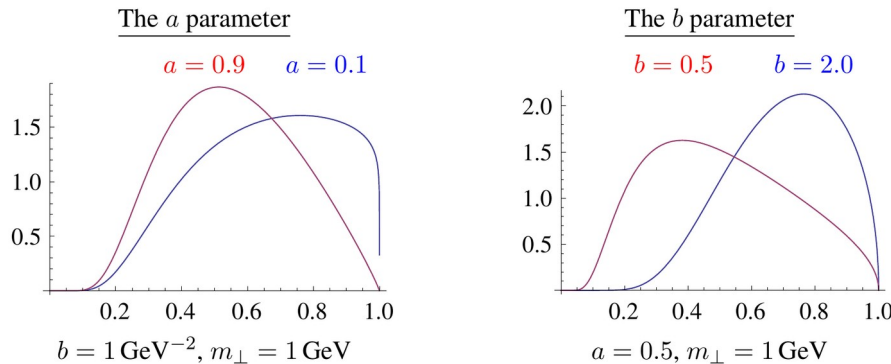
Hybrid approach: Keep the underlying paradigm e.g. strings but modify the microscopic kinematics to accommodate the discrepant global experimental observables.



Hybrid data-driven approach

The phenomenological models of hadronization already give an acceptable description of a large amount of data.

Hybrid approach: Keep the underlying paradigm e.g. strings but modify the microscopic kinematics to accommodate the discrepant global experimental observables.



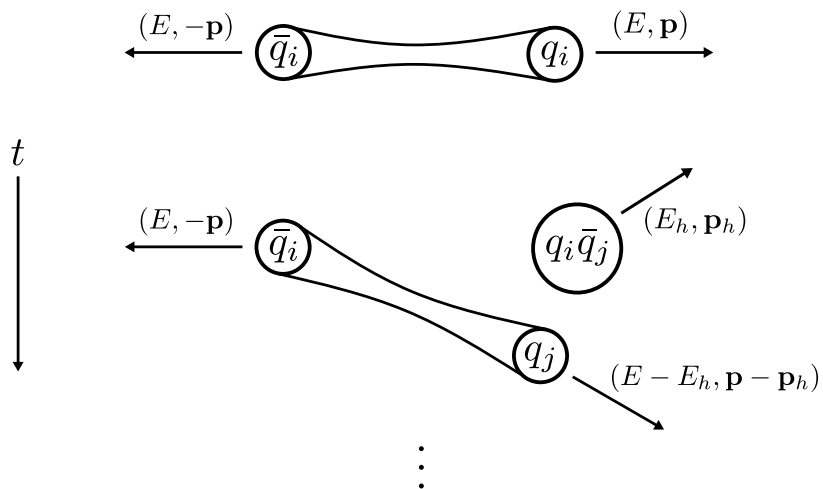
+ ϵ

Machine learning offers a nice framework to tackle this problem

Phenomenological models

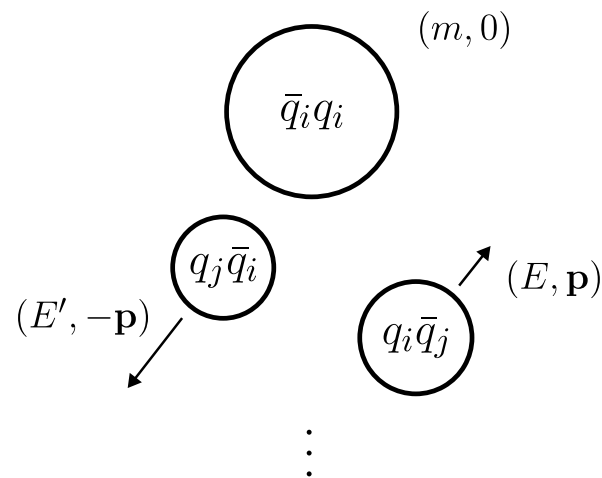
Lund string model

(used in Pythia)



Cluster model

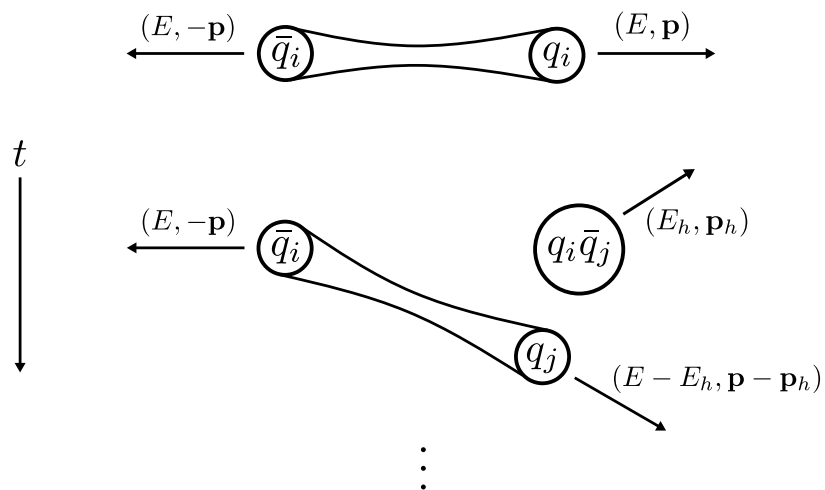
(used in Herwig)



Phenomenological models

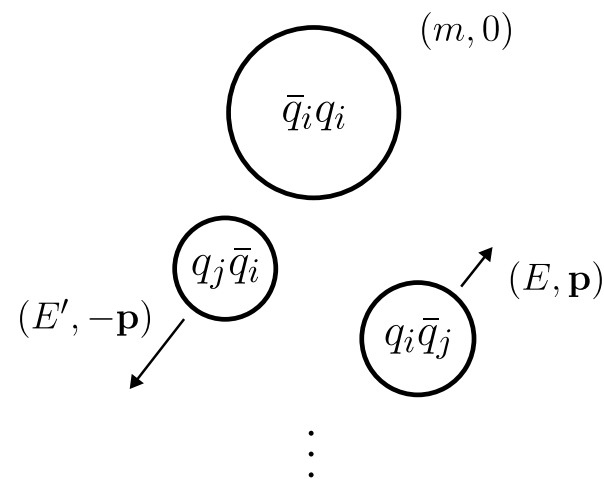
Lund string model

(used in Pythia)



Cluster model

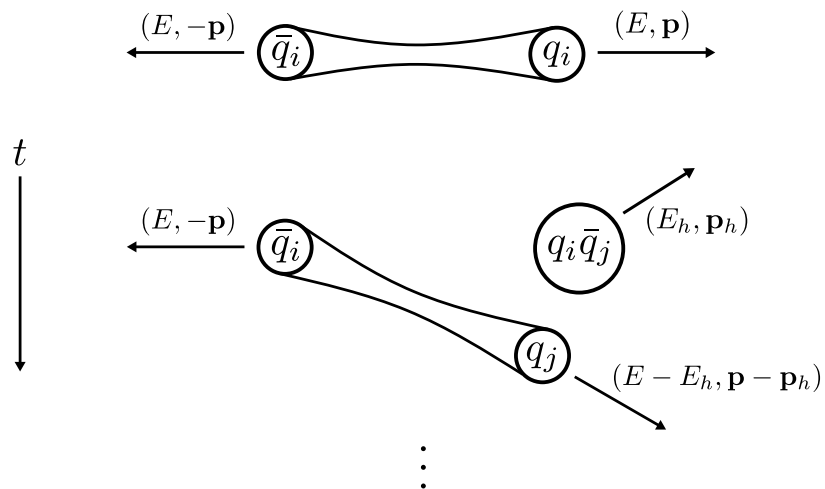
(used in Herwig)



Phenomenological models

Lund string model

(used in Pythia)



Cluster model

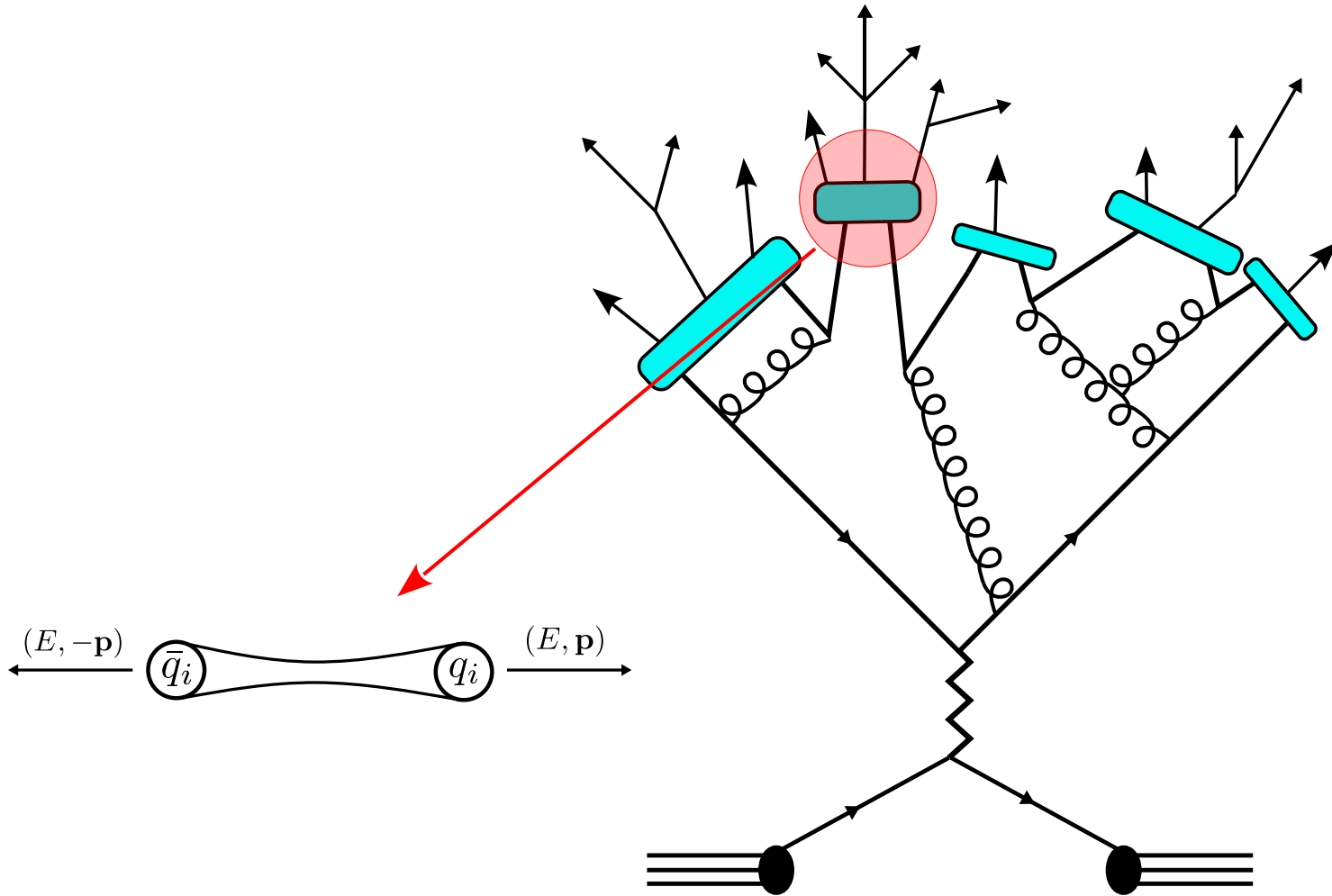
(used in Herwig)

HadML collaboration

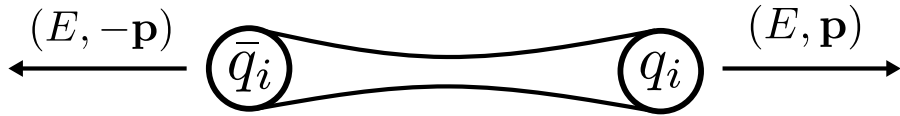
For ML methods applied to the cluster model:

2203.12660, 2305.17169,
2312.08453

\vdots



The algorithm ($q\bar{q}$)



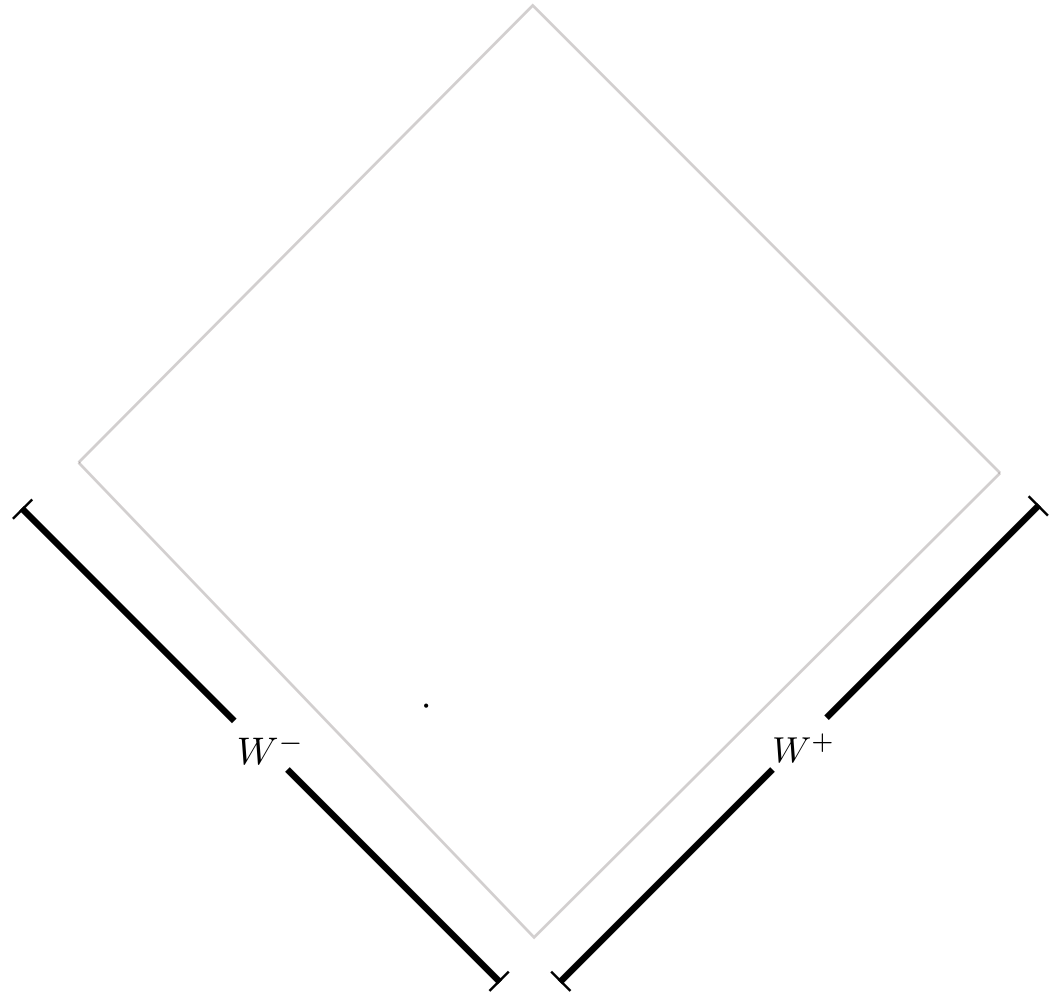
- 1) Randomly select one of the string ends
- 2) Sample new quark flavor
- 3) Sample transverse momentum of new quarks

$$\mathcal{P}(p_x, p_y; \sigma_{p_T}) = \frac{1}{\pi\sigma_{p_T}^2} \exp\left(-\frac{p_x^2 + p_y^2}{\sigma_{p_T}^2}\right)$$

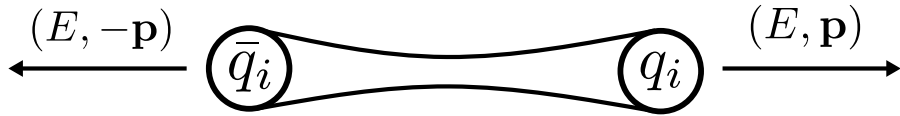
- 4) Sample longitudinal momentum fraction of new hadron

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(\frac{-bm_T^2}{z}\right), \quad z = \frac{p_z + E_h}{2E}$$

- 5) Repeat steps 1-4



The algorithm ($q\bar{q}$)



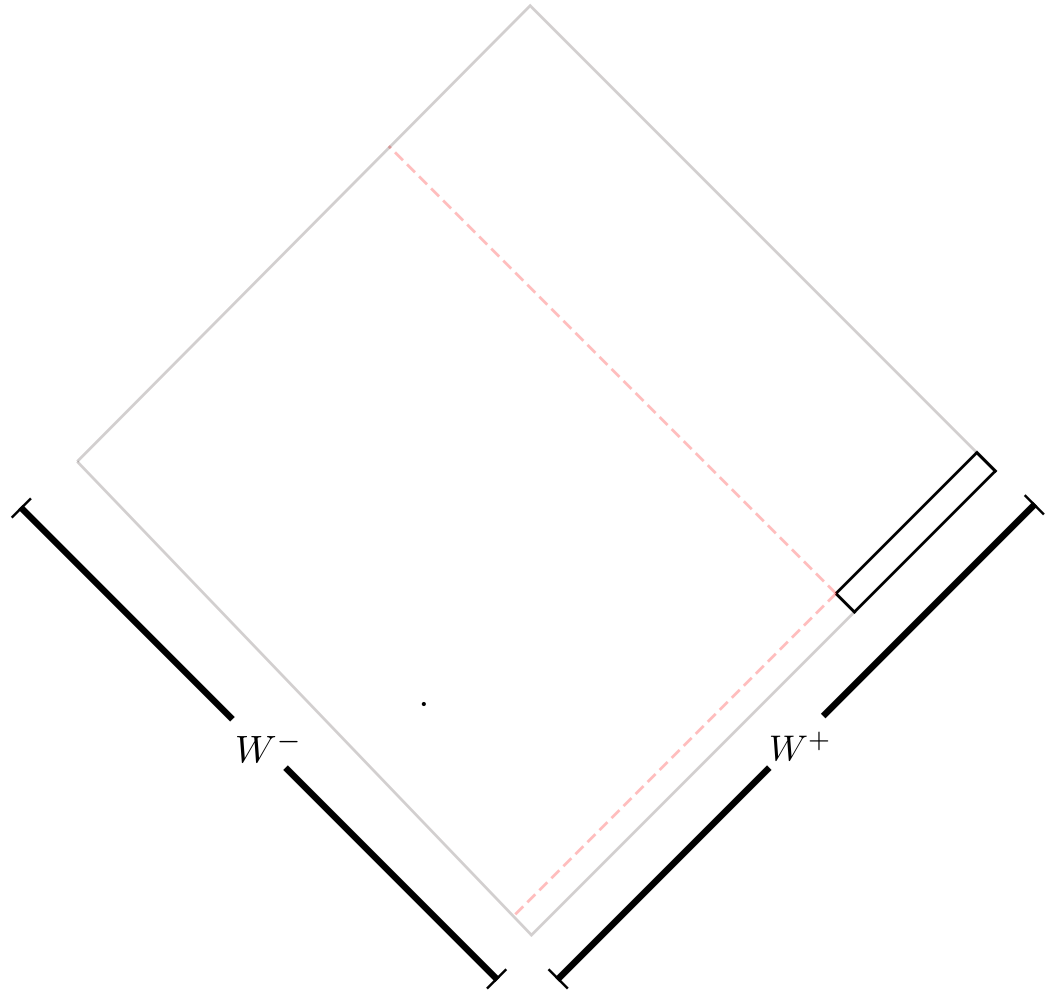
- 1) Randomly select one of the string ends
- 2) Sample new quark flavor
- 3) Sample transverse momentum of new quarks

$$\mathcal{P}(p_x, p_y; \sigma_{p_T}) = \frac{1}{\pi \sigma_{p_T}^2} \exp\left(-\frac{p_x^2 + p_y^2}{\sigma_{p_T}^2}\right)$$

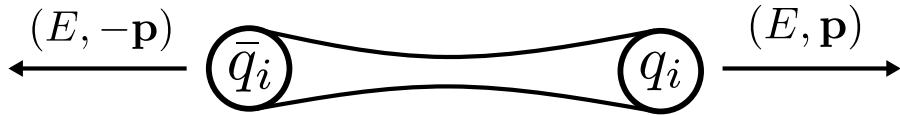
- 4) Sample longitudinal momentum fraction of new hadron

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(-\frac{bm_T^2}{z}\right), \quad z = \frac{p_z + E_h}{2E}$$

- 5) Repeat steps 1-4



The algorithm ($q\bar{q}$)



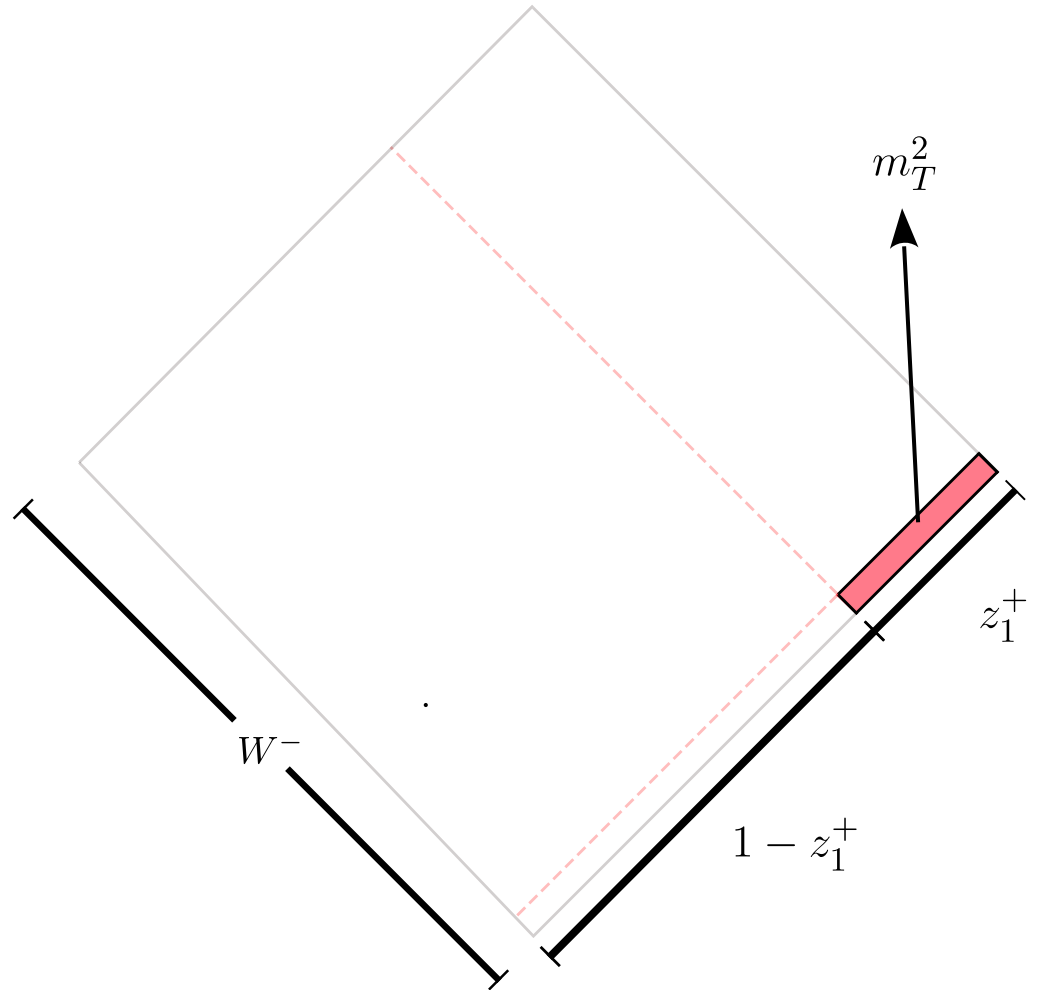
- 1) Randomly select one of the string ends
- 2) Sample new quark flavor
- 3) Sample transverse momentum of new quarks

$$\mathcal{P}(p_x, p_y; \sigma_{p_T}) = \frac{1}{\pi \sigma_{p_T}^2} \exp\left(-\frac{p_x^2 + p_y^2}{\sigma_{p_T}^2}\right)$$

- 4) Sample longitudinal momentum fraction of new hadron

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(\frac{-bm_T^2}{z}\right), \quad z = \frac{p_z + E_h}{2E}$$

- 5) Repeat steps 1-4



The algorithm ($q\bar{q}$)

$(E, -\mathbf{p})$ \leftarrow (\bar{q})

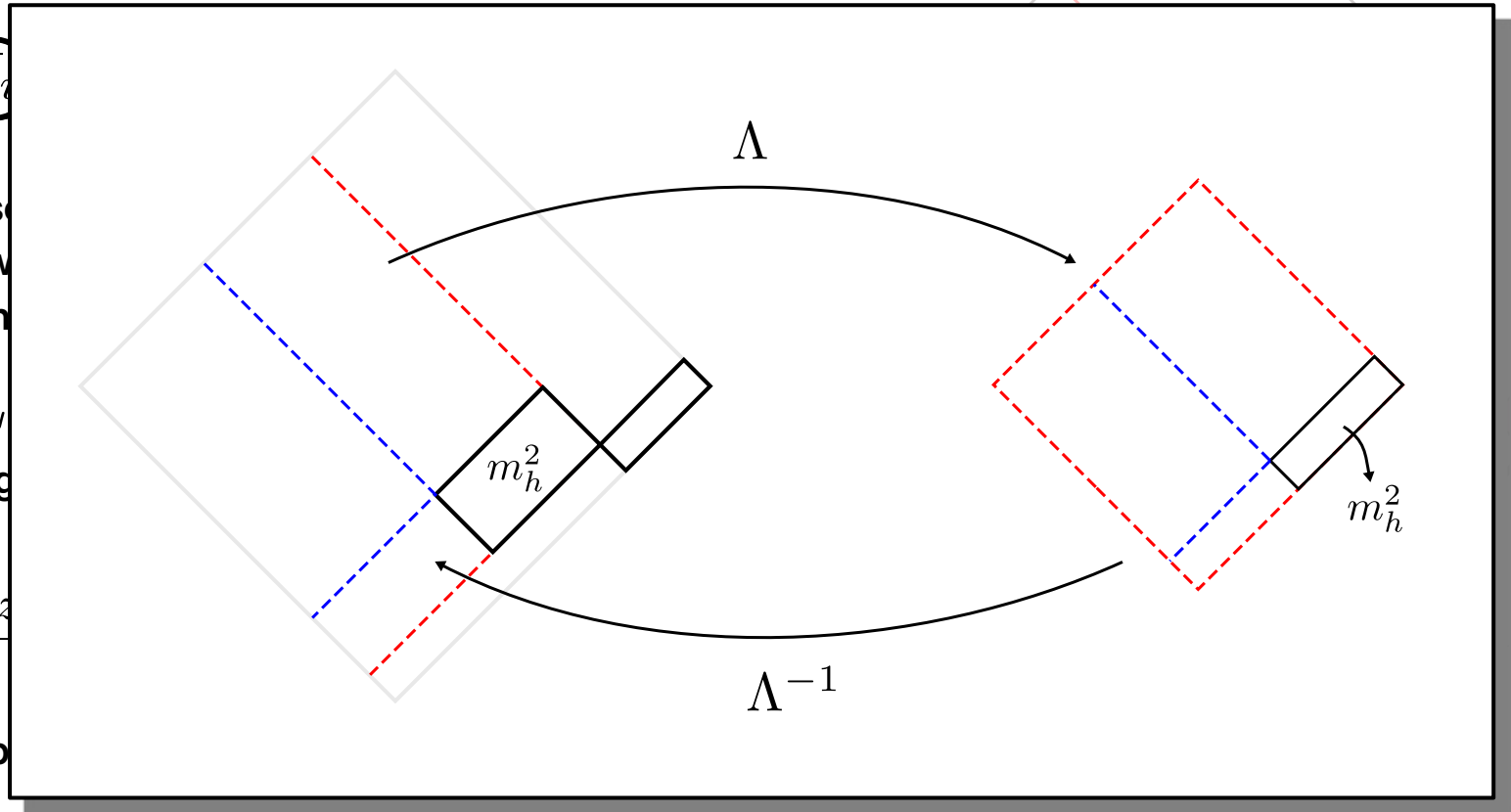
- 1) Randomly select
- 2) Sample new
- 3) Sample trans

$$\mathcal{P}(p_x, p_y)$$

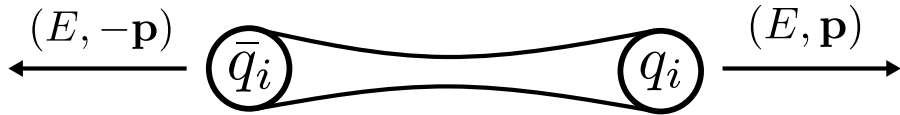
- 4) Sample longitudinal
new hadron

$$f(z) \propto \frac{(1-z)}{z}$$

- 5) Repeat step



The algorithm ($q\bar{q}$)



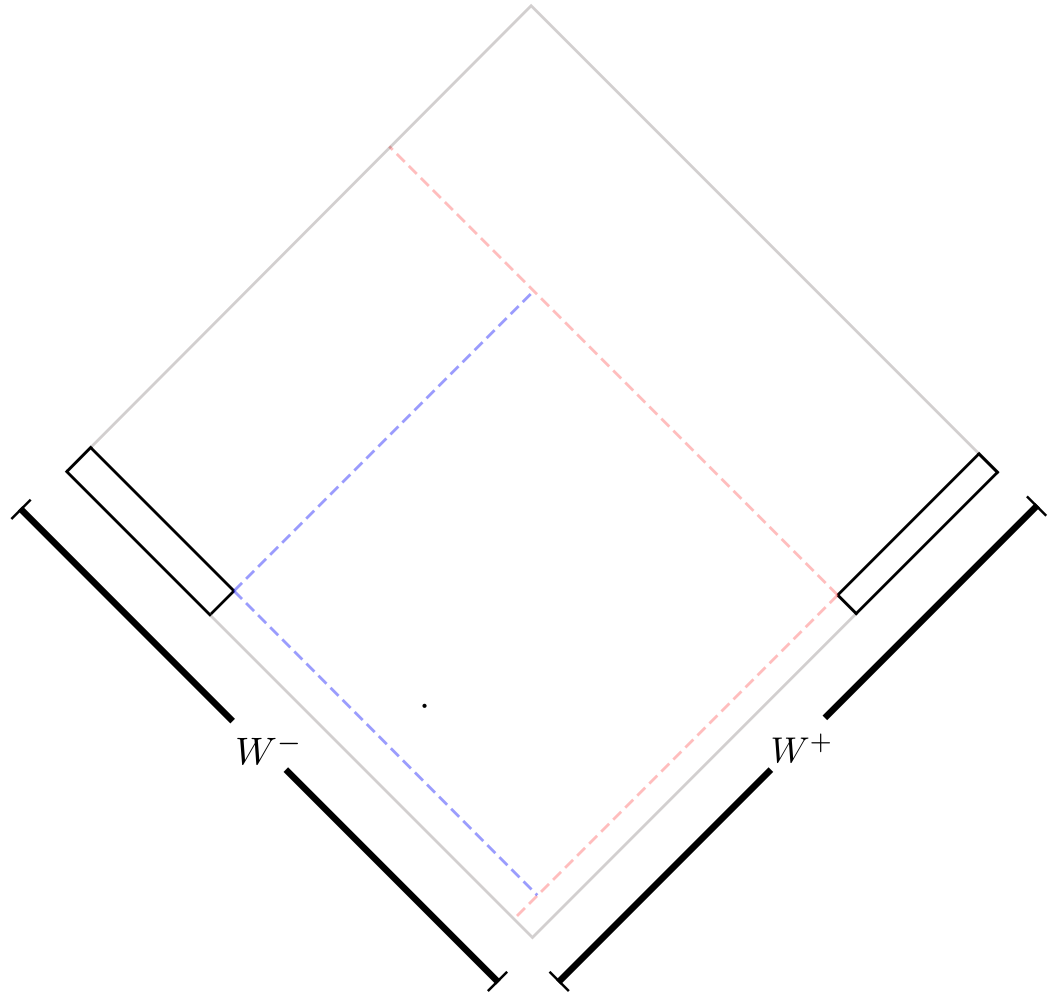
- 1) Randomly select one of the string ends
- 2) Sample new quark flavor
- 3) Sample transverse momentum of new quarks

$$\mathcal{P}(p_x, p_y; \sigma_{p_T}) = \frac{1}{\pi\sigma_{p_T}^2} \exp\left(-\frac{p_x^2 + p_y^2}{\sigma_{p_T}^2}\right)$$

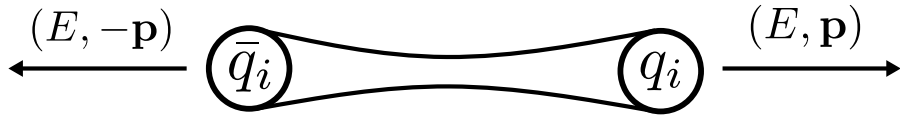
- 4) Sample longitudinal momentum fraction of new hadron

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(\frac{-bm_T^2}{z}\right), \quad z = \frac{p_z + E_h}{2E}$$

- 5) Repeat steps 1-4



The algorithm ($q\bar{q}$)



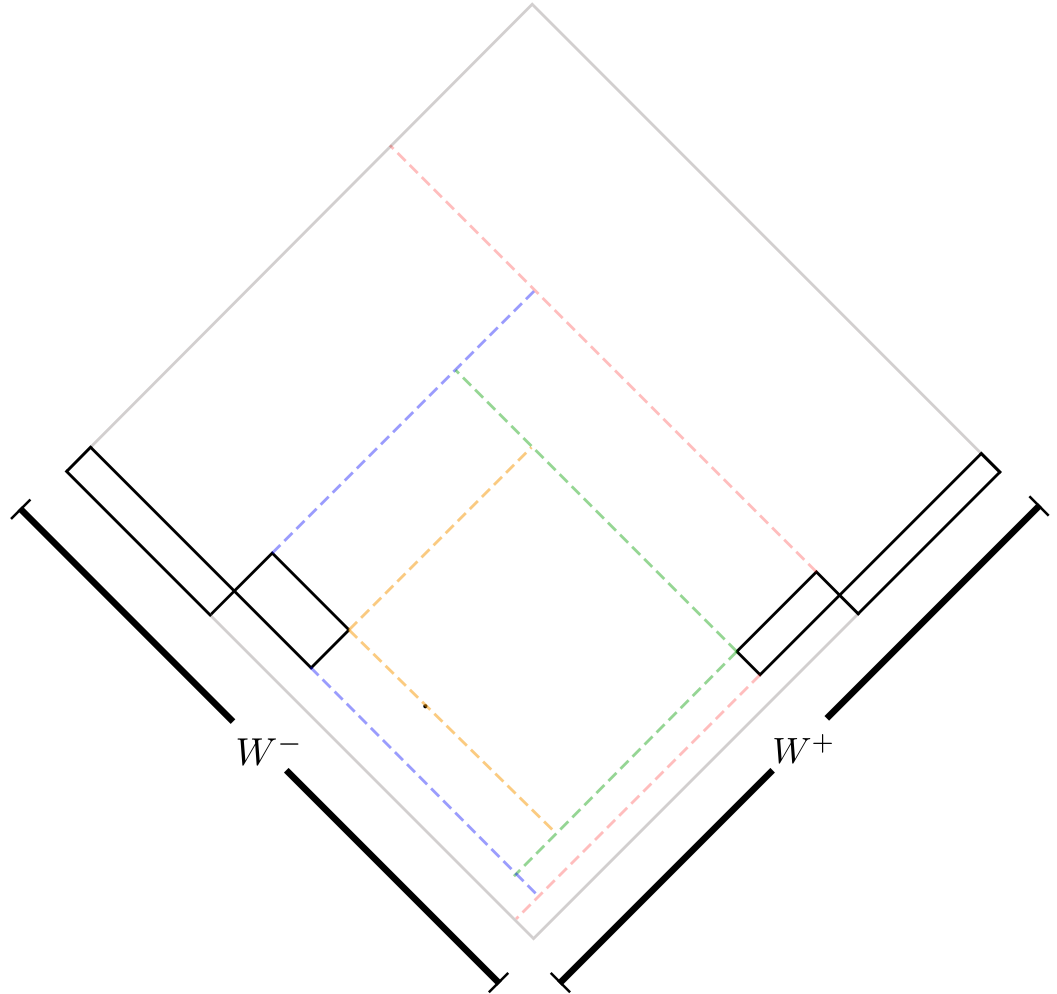
- 1) Randomly select one of the string ends
- 2) Sample new quark flavor
- 3) Sample transverse momentum of new quarks

$$\mathcal{P}(p_x, p_y; \sigma_{p_T}) = \frac{1}{\pi\sigma_{p_T}^2} \exp\left(-\frac{p_x^2 + p_y^2}{\sigma_{p_T}^2}\right)$$

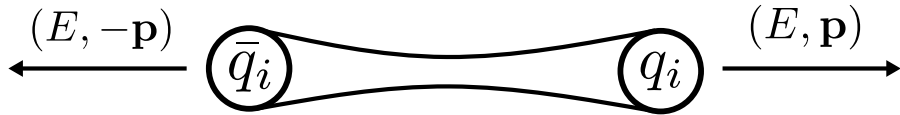
- 4) Sample longitudinal momentum fraction of new hadron

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(-\frac{bm_T^2}{z}\right), \quad z = \frac{p_z + E_h}{2E}$$

- 5) Repeat steps 1-4



The algorithm ($q\bar{q}$)



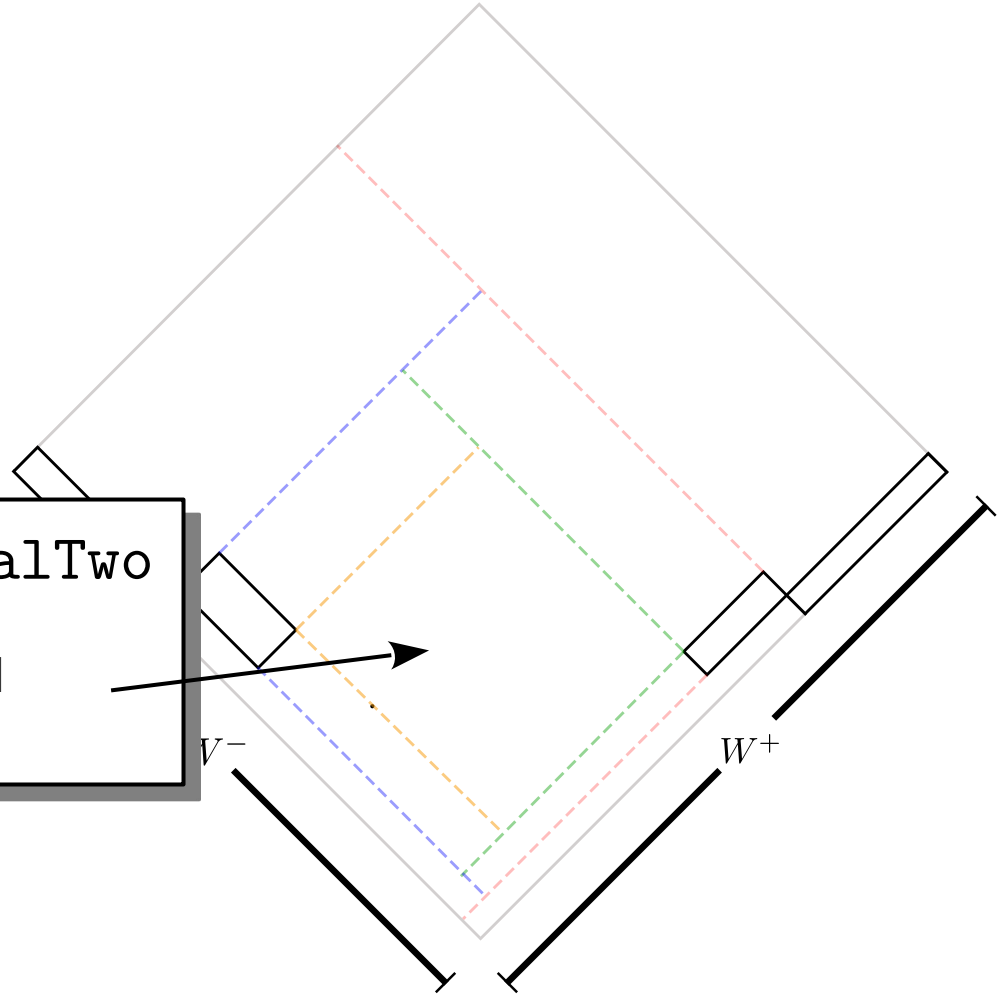
- 1) Randomly select one of the string ends
- 2) Sample new quark flavor
- 3) Sample transverse momentum of new quarks

$\mathcal{P}(z)$
4) Sample new hadron

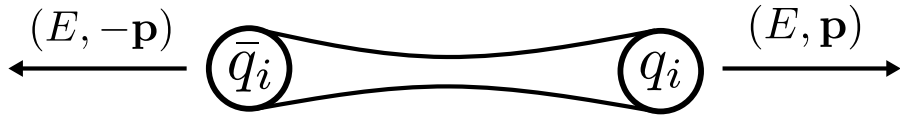
When E_{CM} goes below ~ 2 GeV, finalTwo is called. This can fail, when it does, the full string system is re-simulated from the beginning.

$$f(z) \propto \frac{1}{z} \exp\left(-\frac{1}{z}\right), \quad z = \frac{p_{\perp}}{2E}$$

- 5) Repeat steps 1-4



The algorithm ($q\bar{q}$)



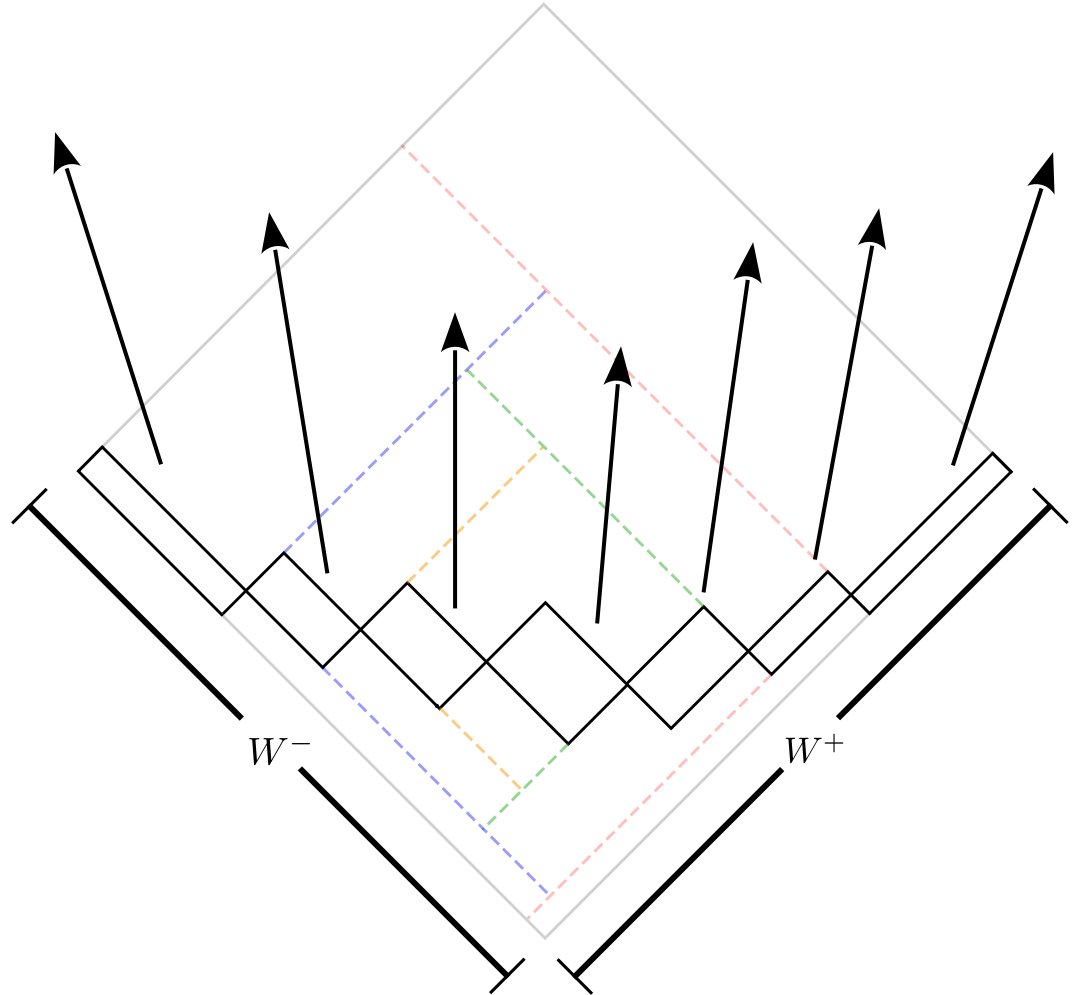
- 1) Randomly select one of the string ends
- 2) Sample new quark flavor
- 3) Sample transverse momentum of new quarks

$$\mathcal{P}(p_x, p_y; \sigma_{p_T}) = \frac{1}{\pi \sigma_{p_T}^2} \exp\left(-\frac{p_x^2 + p_y^2}{\sigma_{p_T}^2}\right)$$

- 4) Sample longitudinal momentum fraction of new hadron

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(-\frac{bm_T^2}{z}\right), \quad z = \frac{p_z + E_h}{2E}$$

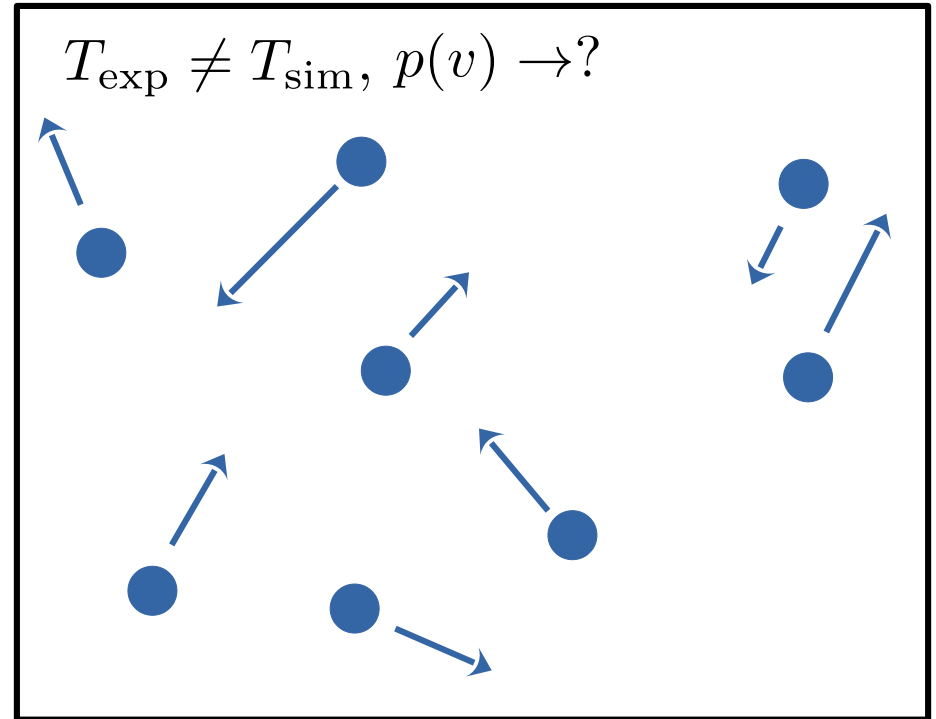
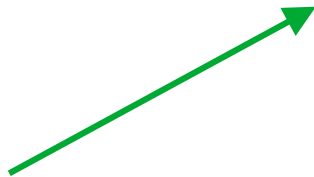
- 5) Repeat steps 1-4



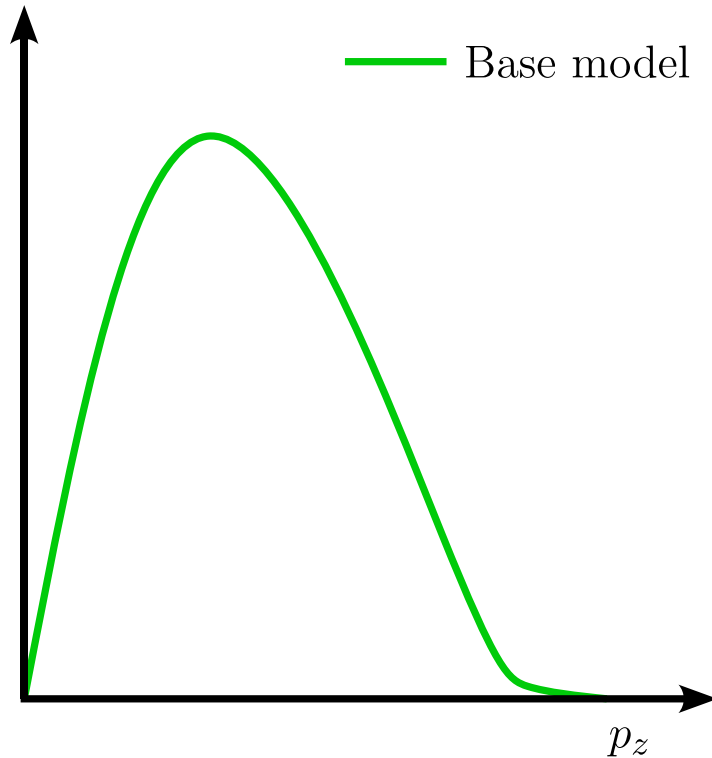
The inverse problem

Task: Given a discrepant macroscopic observable between simulation and experiment, find a way to modify the microscopic dynamics of the simulation to accommodate the observable.

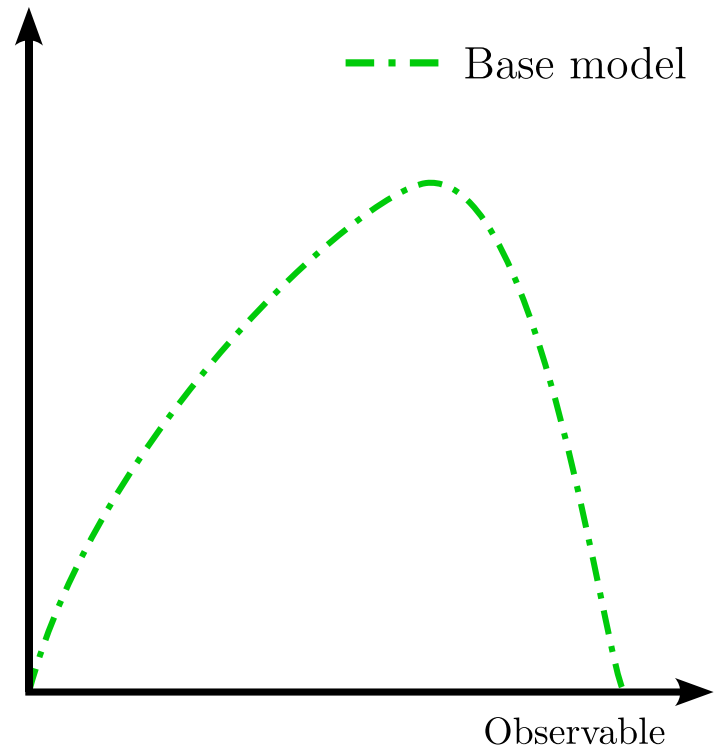
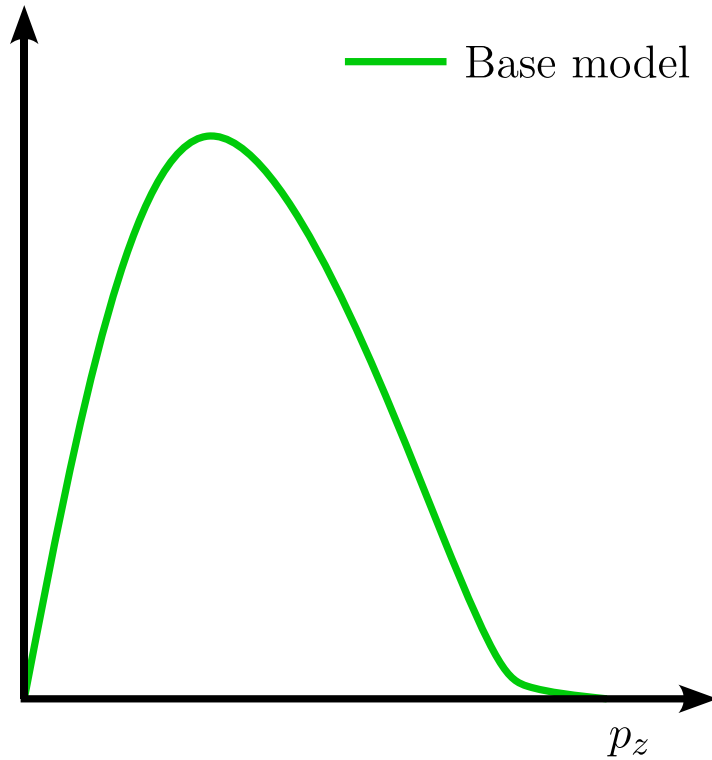
Akin to...



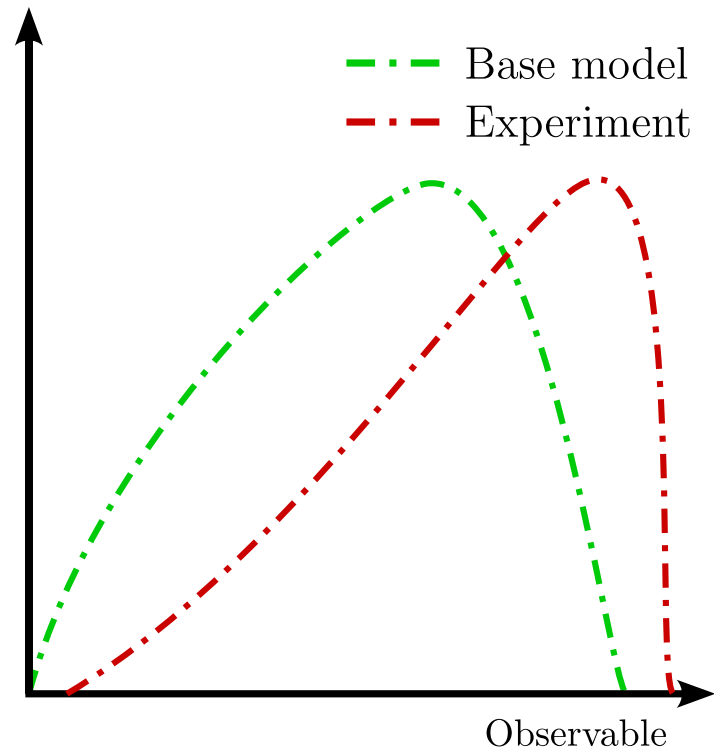
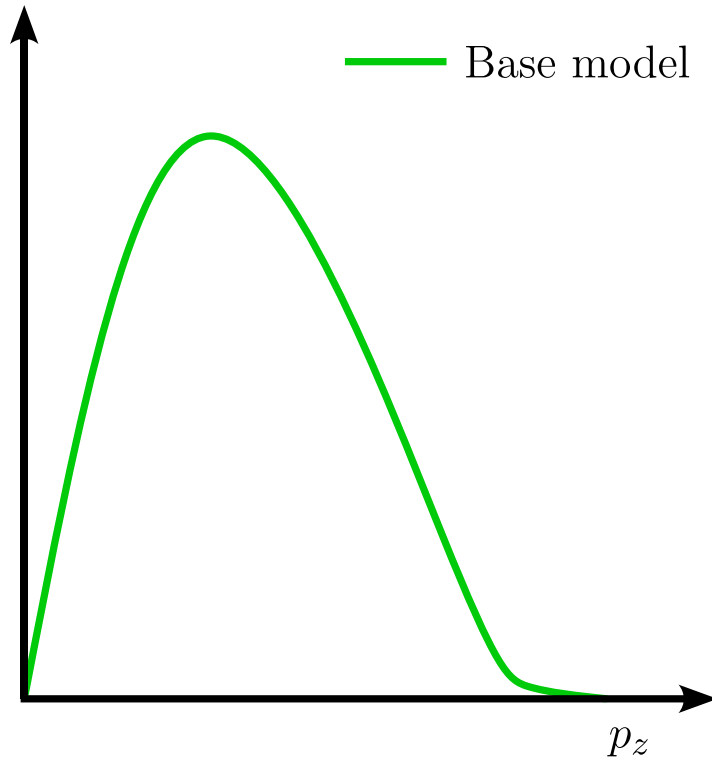
Micro from macro: big picture



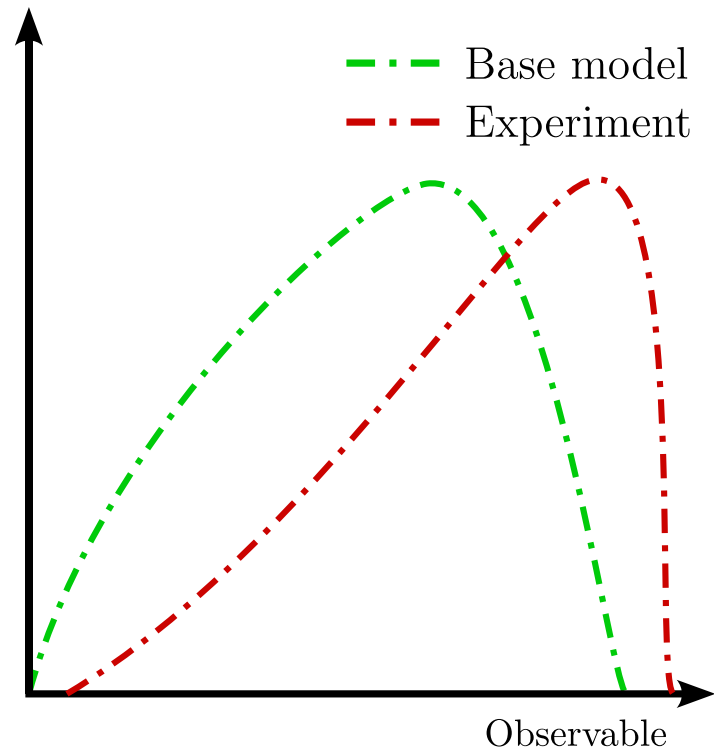
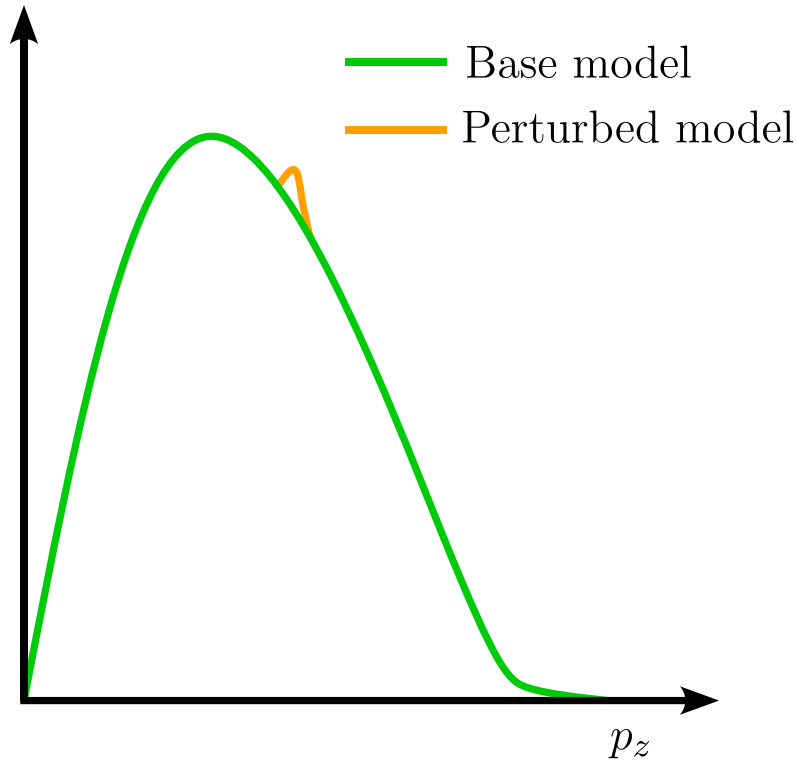
Micro from macro: big picture



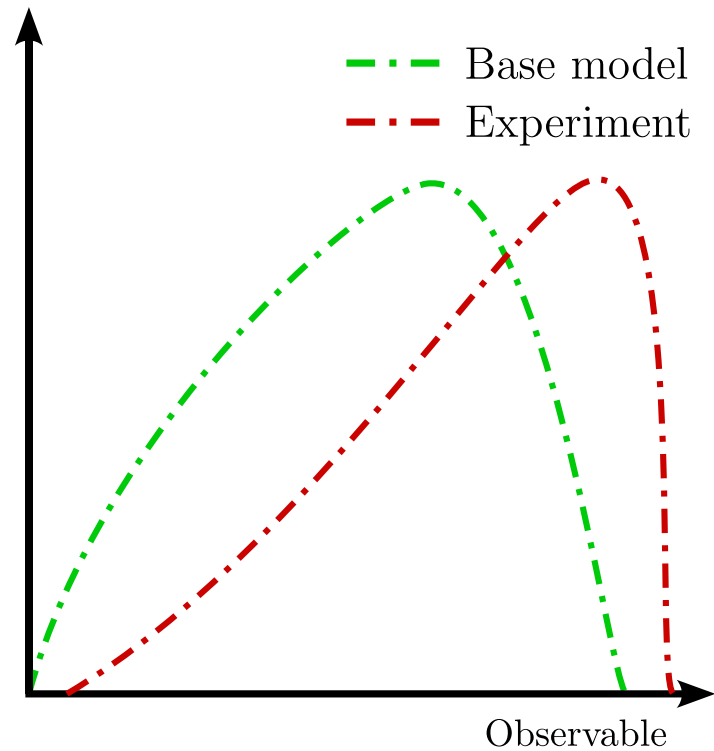
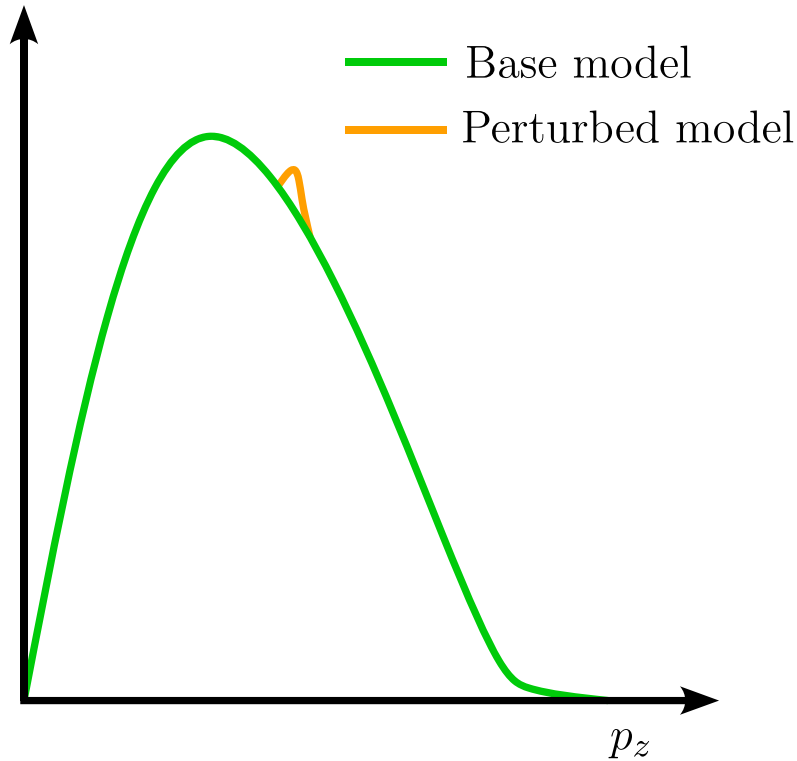
Micro from macro: big picture



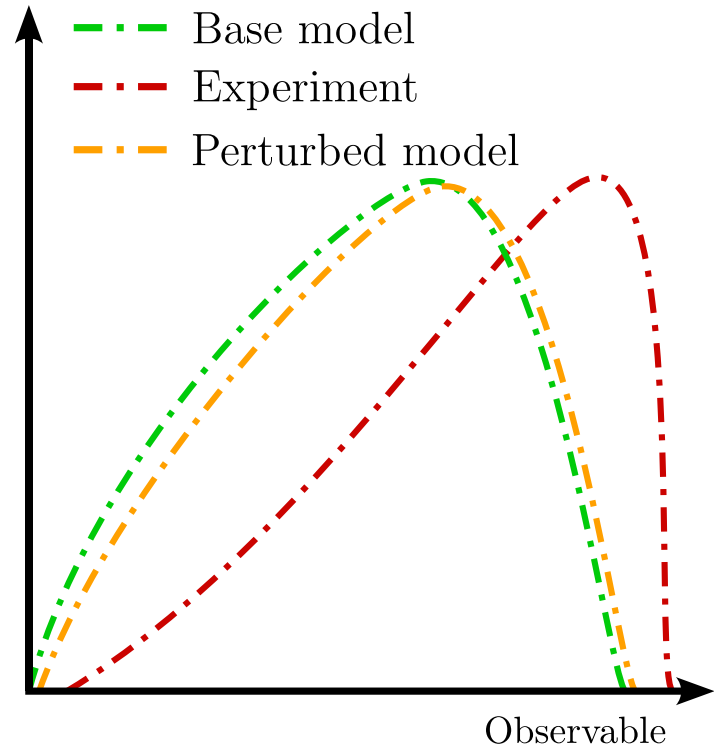
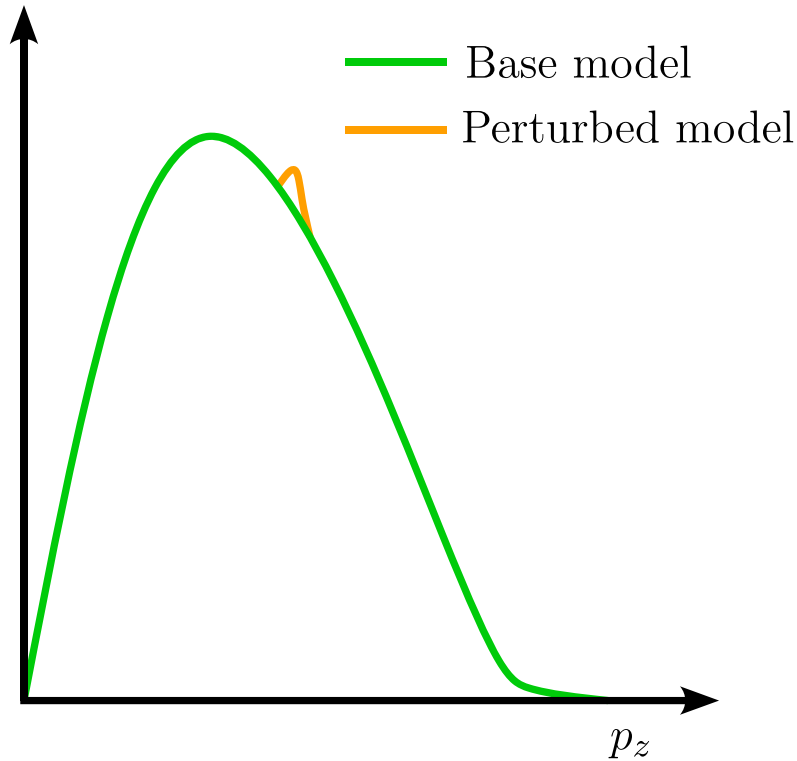
Micro from macro: big picture



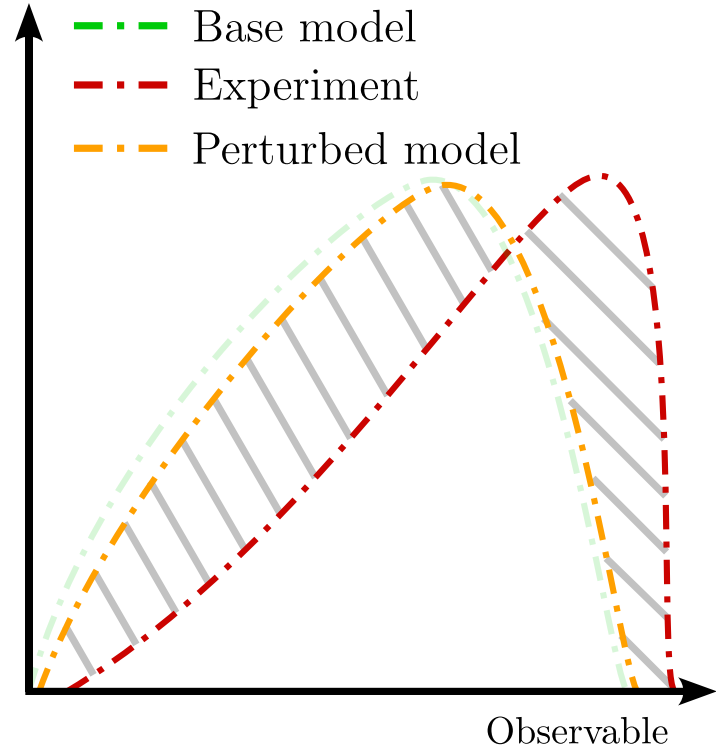
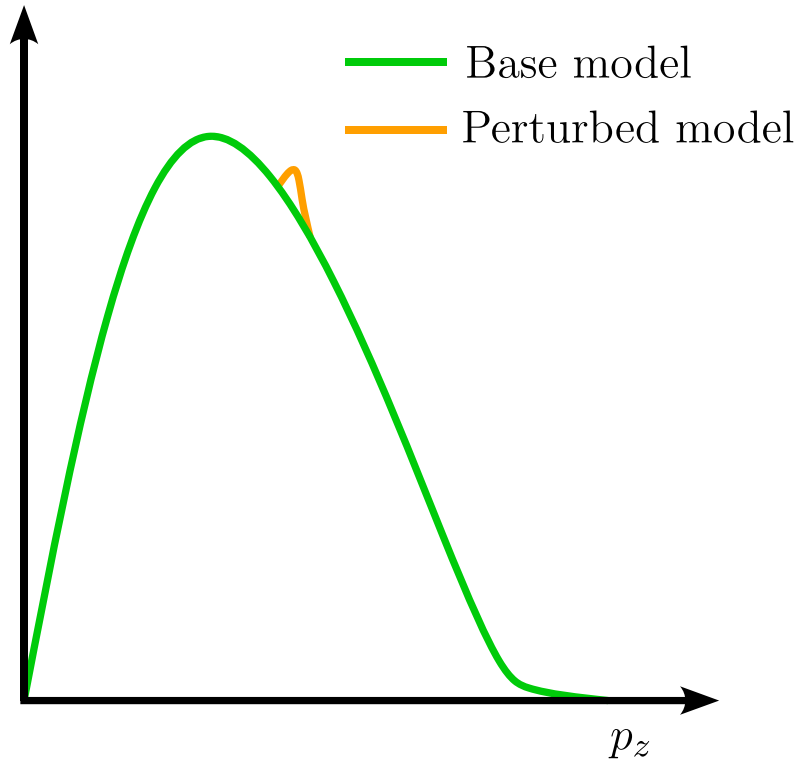
Micro from macro: big picture



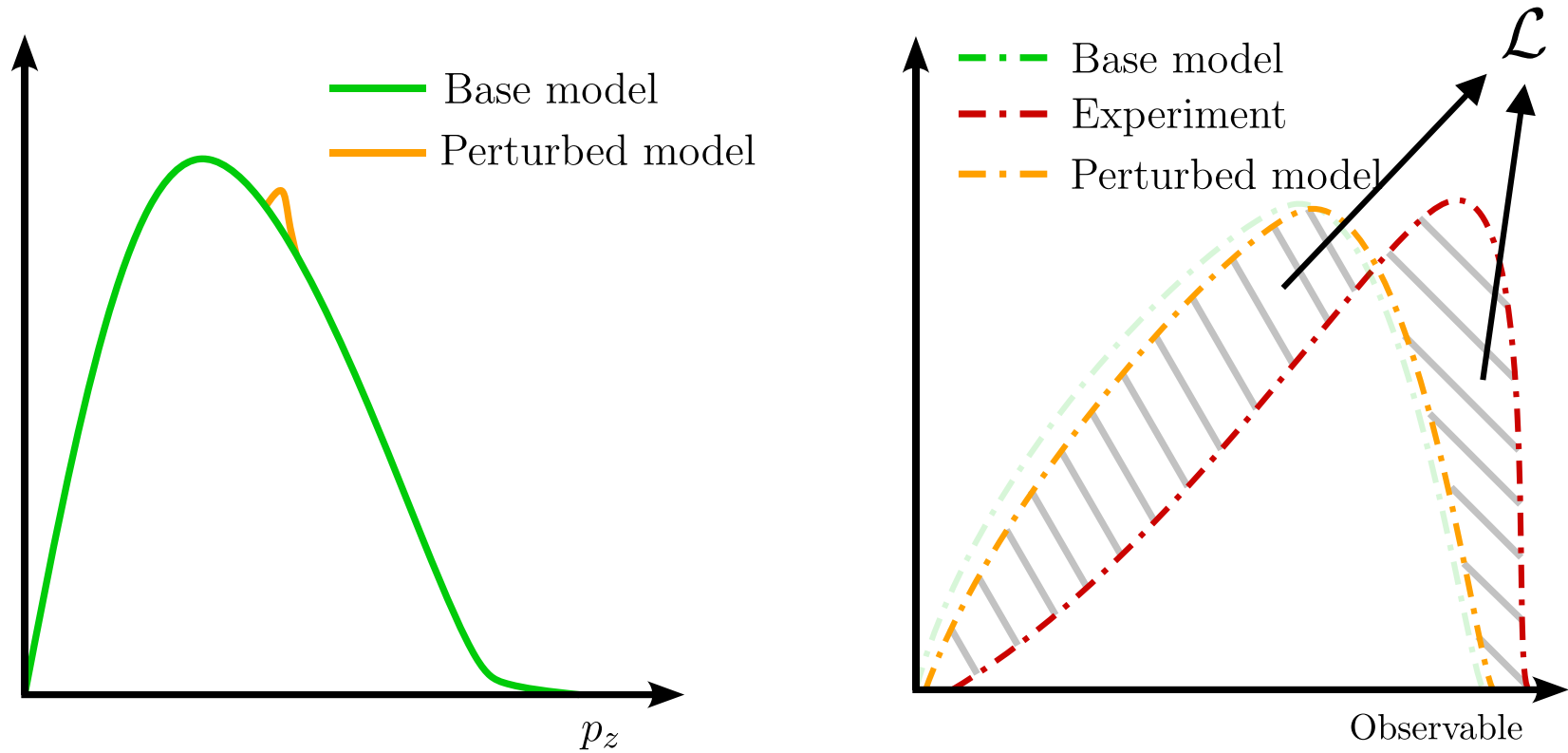
Micro from macro: big picture



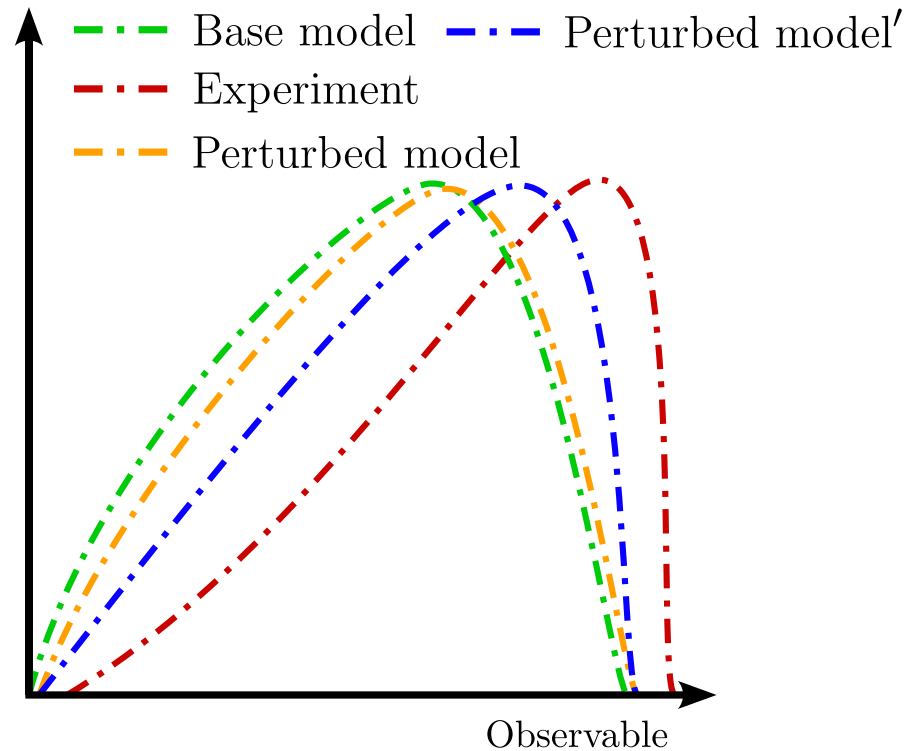
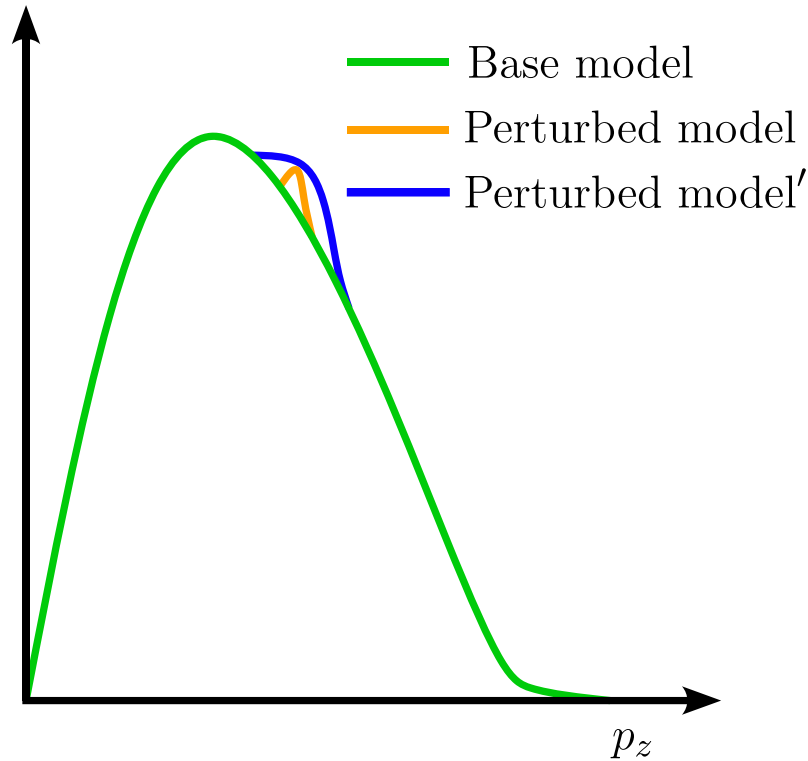
Micro from macro: big picture



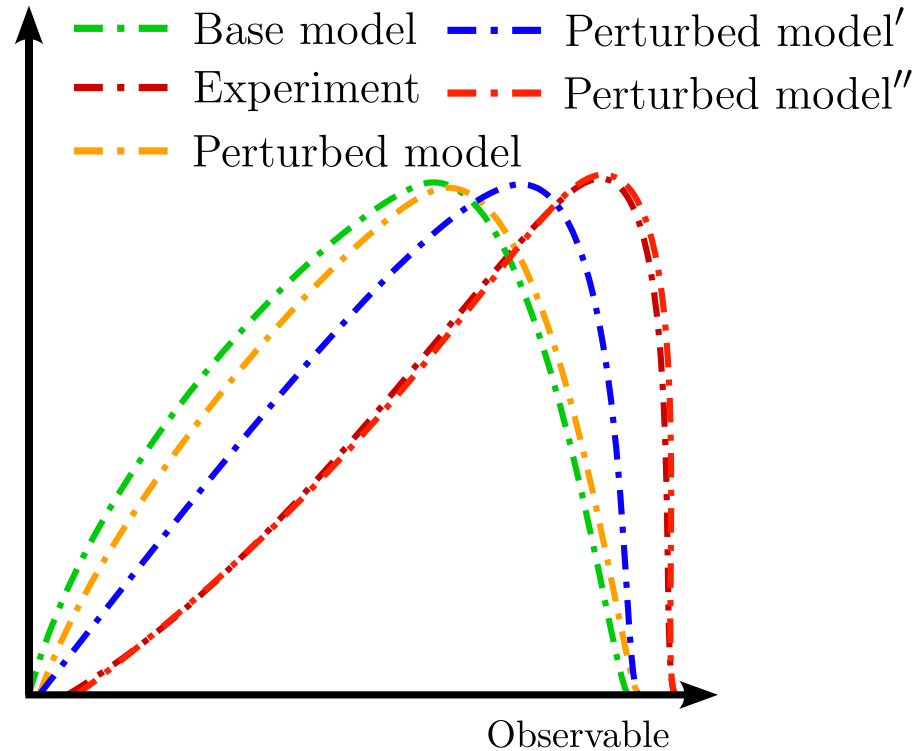
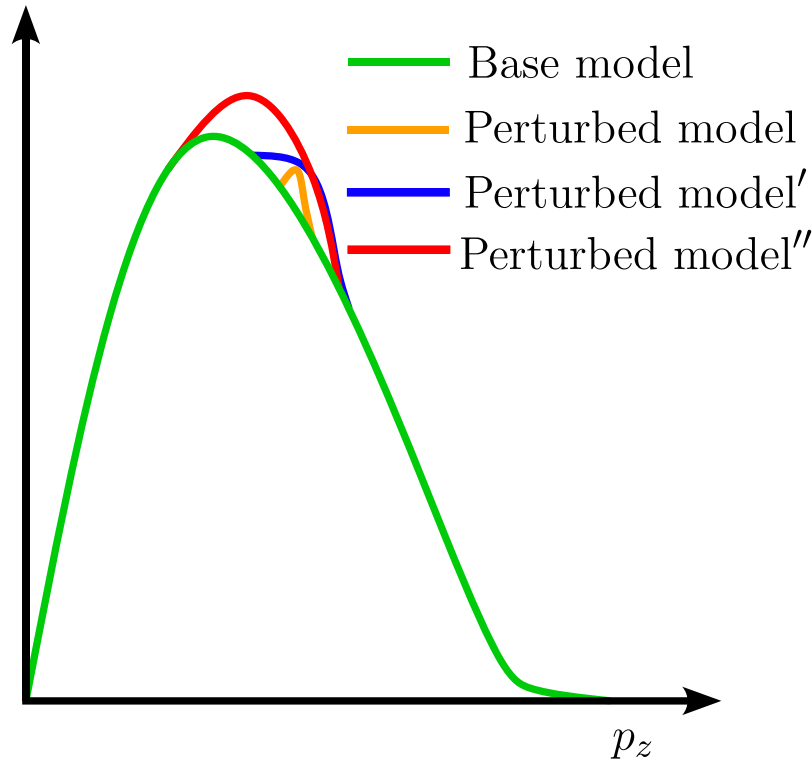
Micro from macro: big picture



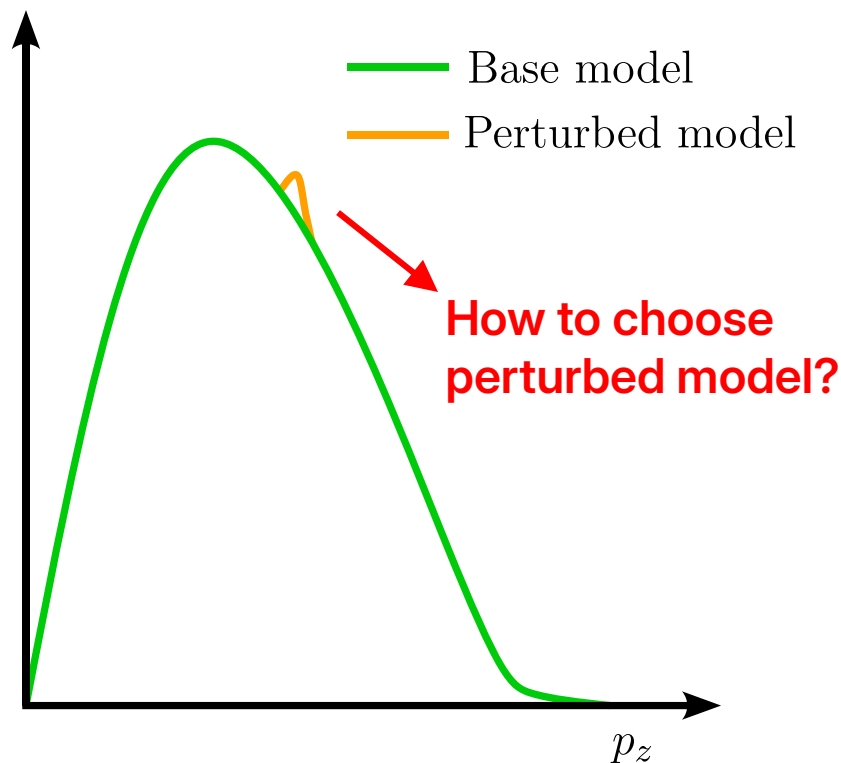
Micro from macro: big picture



Micro from macro: big picture

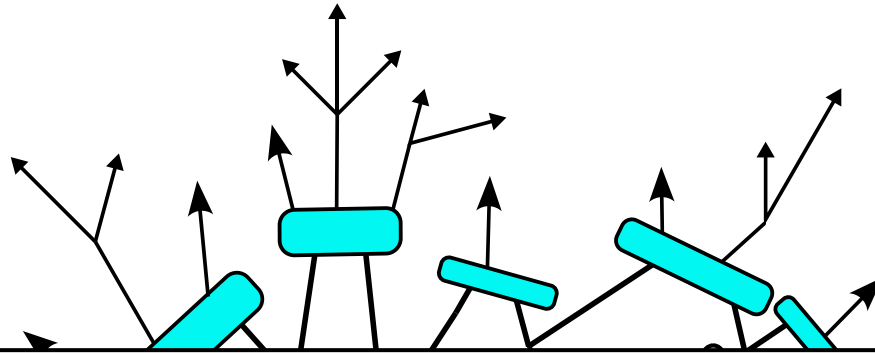


Micro from macro: details



1) Perform an analytic fit of the Lund parameters (traditional fit)

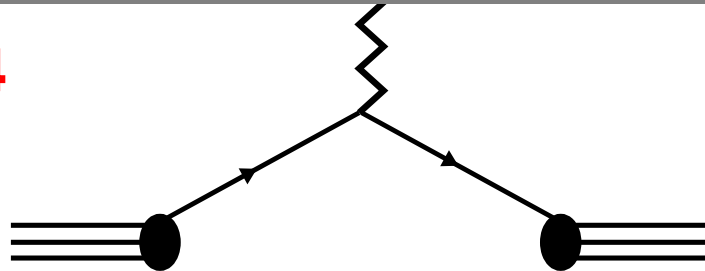
2) Perform an over-parameterized fit (ML fit)



1.

Rejection Sampling with Autodifferentiation (RSA)

2411.02194

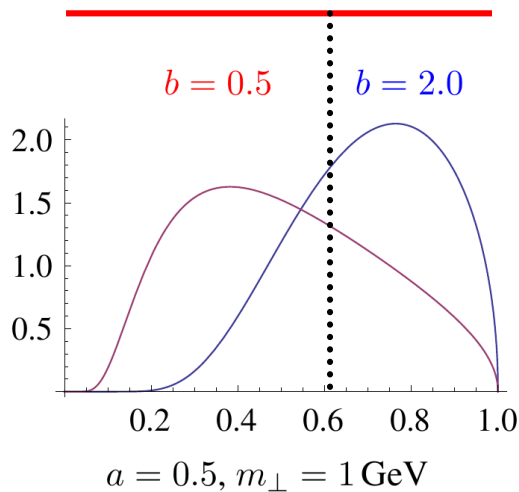


Kinematic reweighting (2308.13459)

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(\frac{-bm_{\perp}^2}{z}\right)$$

Can I understand the change in observable distributions for a different values of a and b WITHOUT re-simulating?

Compute event weights i.e. ratios of probabilities



Rejection sampling:

Acceptance probability:

$$P_{\text{accept}} = \frac{p(z, \theta)}{\hat{P}}$$

Rejection probability

$$P_{\text{reject}} = 1 - P_{\text{accept}}$$

$$w_{\text{accept}} = \frac{P_{\text{accept}}(z; \theta')}{P_{\text{accept}}(z; \theta)}$$

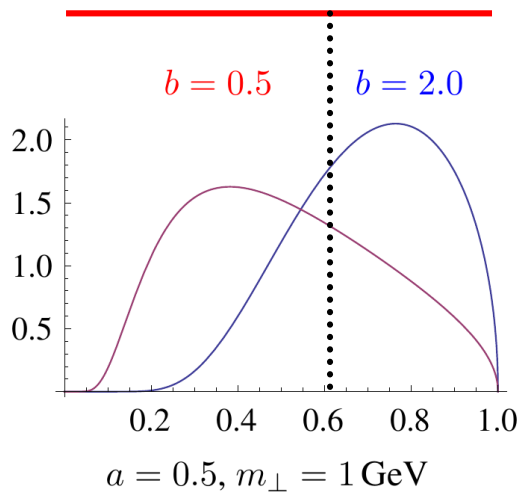
$$w_{\text{reject}} = \frac{1 - P_{\text{accept}}(z; \theta')}{1 - P_{\text{accept}}(z; \theta)}$$

Kinematic reweighting (2308.13459)

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(\frac{-bm_{\perp}^2}{z}\right)$$

Can I understand the change in observable distributions for a different values of a and b WITHOUT re-simulating?

Compute event weights i.e. ratios of probabilities



Rejection sampling:

Acceptance probability:

$$P_{\text{accept}} = \frac{p(z, \theta)}{\hat{P}}$$

Rejection probability

$$P_{\text{reject}} = 1 - P_{\text{accept}}$$

$$w_{\text{accept}} = \frac{p(z; \theta')}{p(z; \theta)}$$

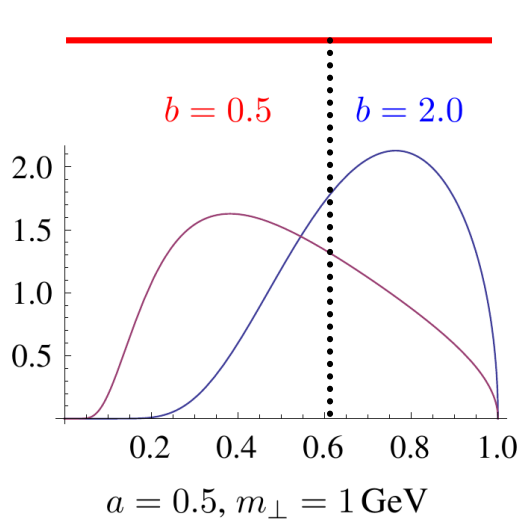
$$w_{\text{reject}} = \frac{\hat{P} - p(z; \theta')}{\hat{P} - p(z; \theta)}$$

Kinematic reweighting (2308.13459)

$$f(z) \propto \frac{(1-z)^a}{z} \exp\left(\frac{-bm_{\perp}^2}{z}\right)$$

Can I understand the change in observable distributions for a different values of a and b WITHOUT re-simulating?

Compute event weights i.e. ratios of probabilities



Rejection sampling:

Acceptance probability:

$$P_{\text{accept}} = \frac{p(z, \theta)}{\hat{P}}$$

Rejection probability

$$P_{\text{reject}} = 1 - P_{\text{accept}}$$

$$w_{\text{accept}} = \frac{p(z; \theta')}{p(z; \theta)}$$

$$w_{\text{reject}} = \frac{\hat{P} - p(z; \theta')}{\hat{P} - p(z; \theta)}$$

$$w = w_{\text{accept}} \prod_{i=1}^{n_{\text{rej.}}} w_{\text{reject}}^i$$

Kinematic reweighting (2308.13459)

$$(1 - z)^a \quad (-bm_{\perp}^2)$$

Can I understand the change in observable

Hadronization data-structure:

$$z = \begin{pmatrix} z_1 = \left(\begin{array}{l} \{m_T^{h_1}, z_{\text{accept}}^{h_1}, z_{\text{reject}}^{1,h_1}, \dots, z_{\text{reject}}^{n_{h_1},h_1}\} \\ \{m_T^{h_2}, z_{\text{accept}}^{h_2}, z_{\text{reject}}^{1,h_2}, \dots, z_{\text{reject}}^{n_{h_2},h_2}\} \\ \{m_T^{h_3}, z_{\text{accept}}^{h_3}, z_{\text{reject}}^{1,h_3}, \dots, z_{\text{reject}}^{n_{h_3},h_3}\} \end{array} \right)_1 \\ z_2 = \left(\begin{array}{l} \{m_T^{h_1}, z_{\text{accept}}^{h_1}, z_{\text{reject}}^{1,h_1}, \dots, z_{\text{reject}}^{n_{h_1},h_1}\} \\ \vdots \\ \{m_T^{h_4}, z_{\text{accept}}^{h_4}, z_{\text{reject}}^{1,h_4}, \dots, z_{\text{reject}}^{n_{h_4},h_4}\} \\ \vdots \\ z_N = \dots \end{array} \right)_2 \end{pmatrix}$$

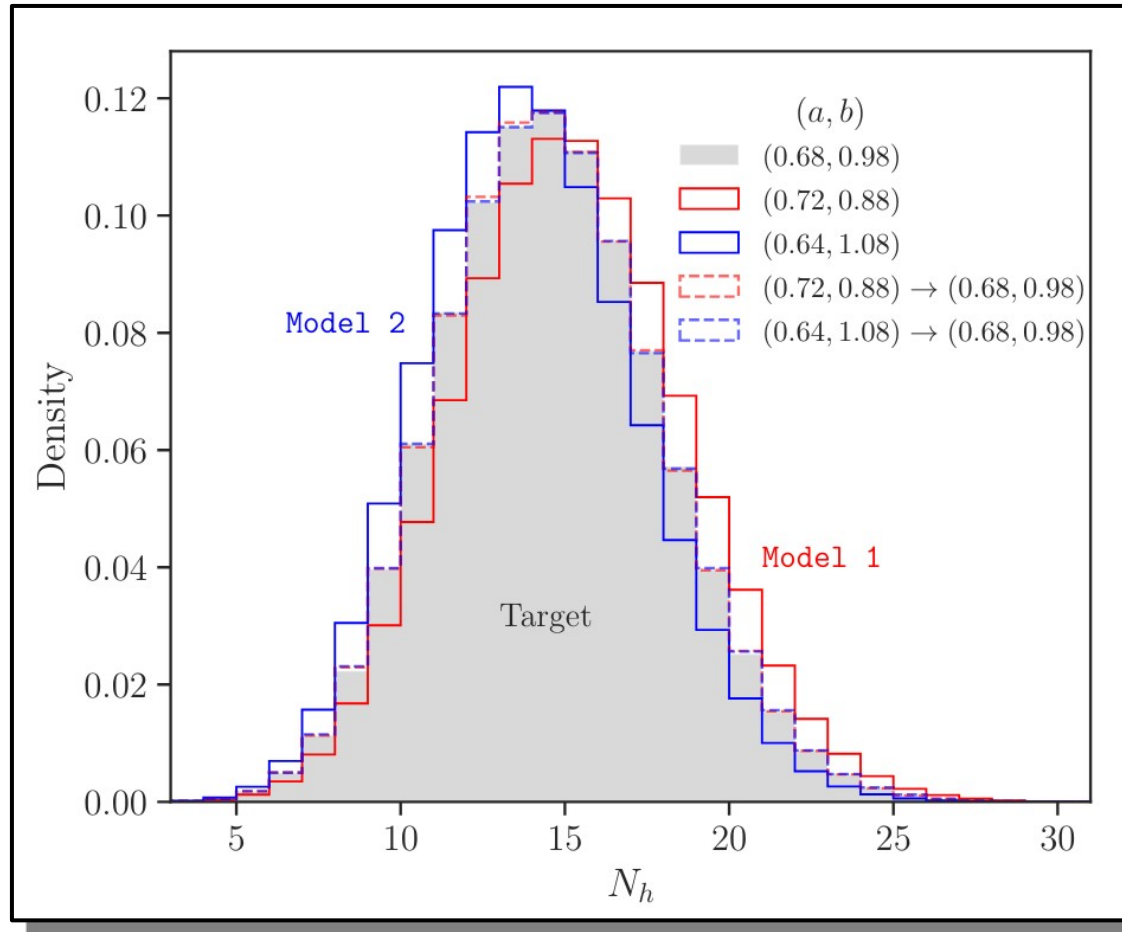
$$w_n = \prod_{i=1}^{\tilde{N}_{h,n}} \left(\frac{f(z_{\text{accept}}^{h_i}; \{a, b\}_P)}{f(z_{\text{accept}}^{h_i}; \{a, b\}_B)} \right) \times \prod_{j=1}^{n_{h_i}} \left(\frac{\hat{f} - f(z_{\text{reject}}^{j,h_i}; \{a, b\}_P)}{\hat{f} - f(z_{\text{reject}}^{j,h_i}; \{a, b\}_B)} \right)$$

$$a = 0.5, m_{\perp} = 1 \text{ GeV}$$

$$P_{\text{reject}} = 1 - P_{\text{accept}}$$

$$\omega_{\text{accept}} \prod_{i=1}^{\tilde{N}_{h,n}} \omega_{\text{reject}}$$

Kinematic reweighting (2308.13459)



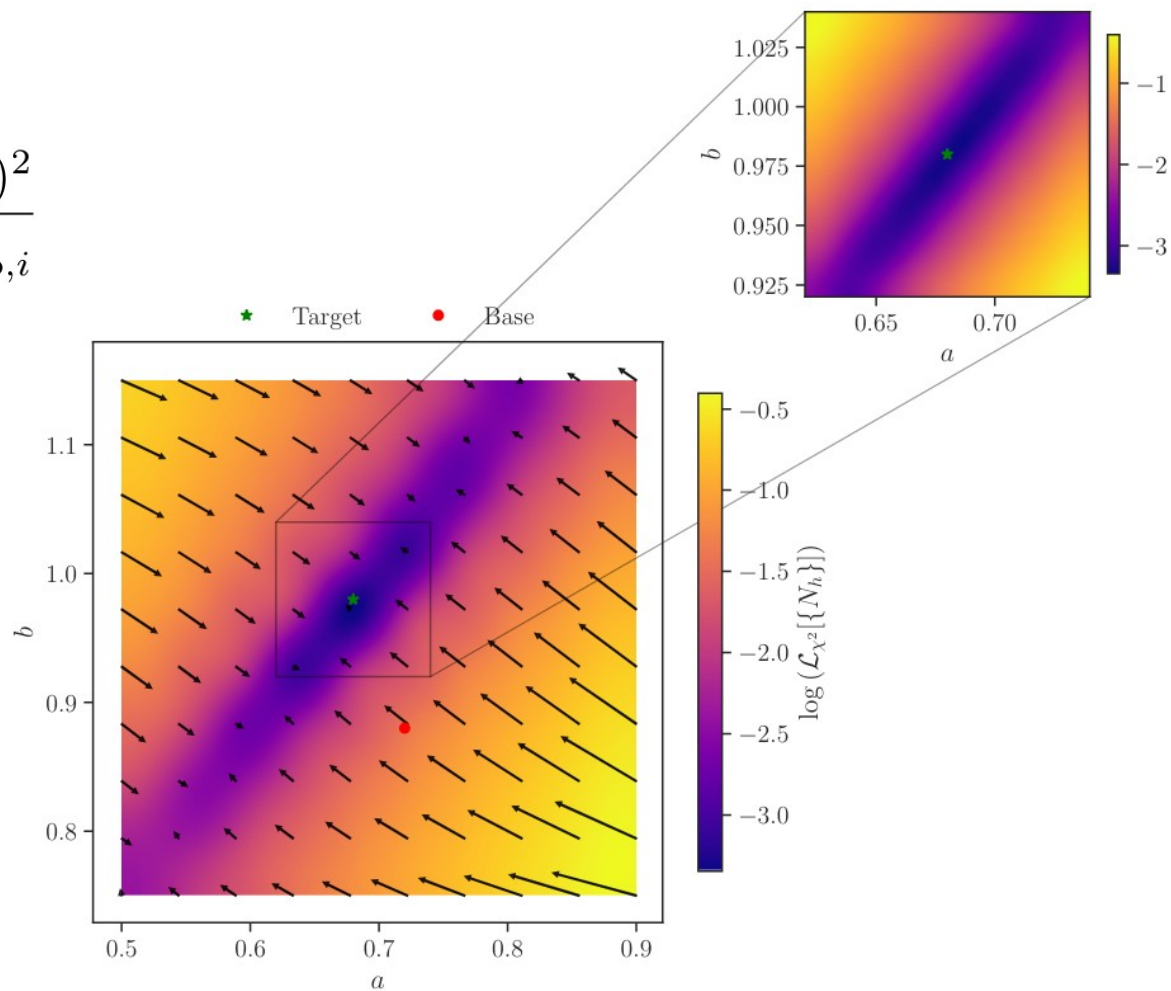
Applications

- Efficient means of exploring parameter space
- Inherently differentiable
- **Very useful for tuning!**
 - Embed the computation of weights into numerical autodifferentiation engine
 - Picking new parameters in the update step facilitated by well-developed optimizers (SGD, Adam, etc.)

Rejection sampling with Autodifferentiation (RSA)

χ^2 -loss landscape

$$\mathcal{L}_{\chi^2}(\mathbf{y}_{\text{sim}}, \mathbf{y}_{\text{exp}}; \mathbf{w}) = \sum_{i=1}^{n_{\text{bins}}} \frac{(y_{\text{sim}}^{(i)} - y_{\text{exp}}^{(i)})^2}{\sigma_{\text{sim},i}^2 + \sigma_{\text{exp},i}^2}$$



The "score" observable

- Train a deepsets classifier to distinguish simulation from experiment

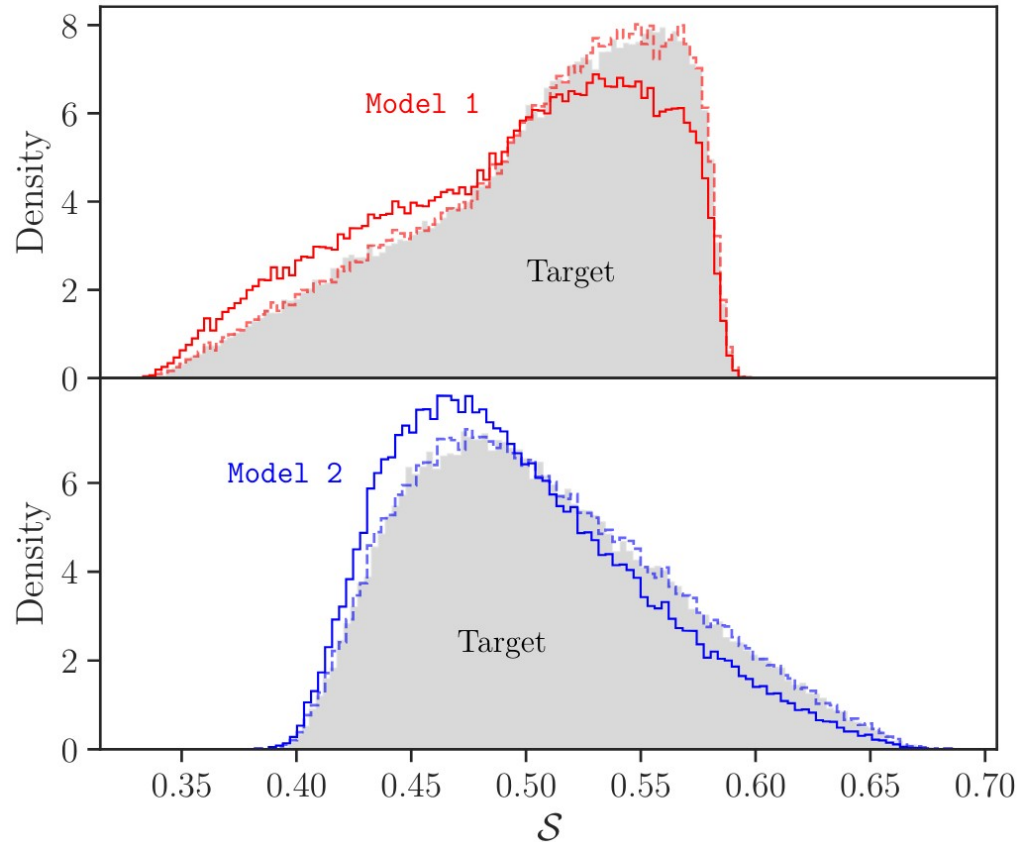
- Takes full event information as input

$(E, p_x, p_y, p_z)_1$

$(E, p_x, p_y, p_z)_2$

...

Use trained classifier output
(score) as an observable

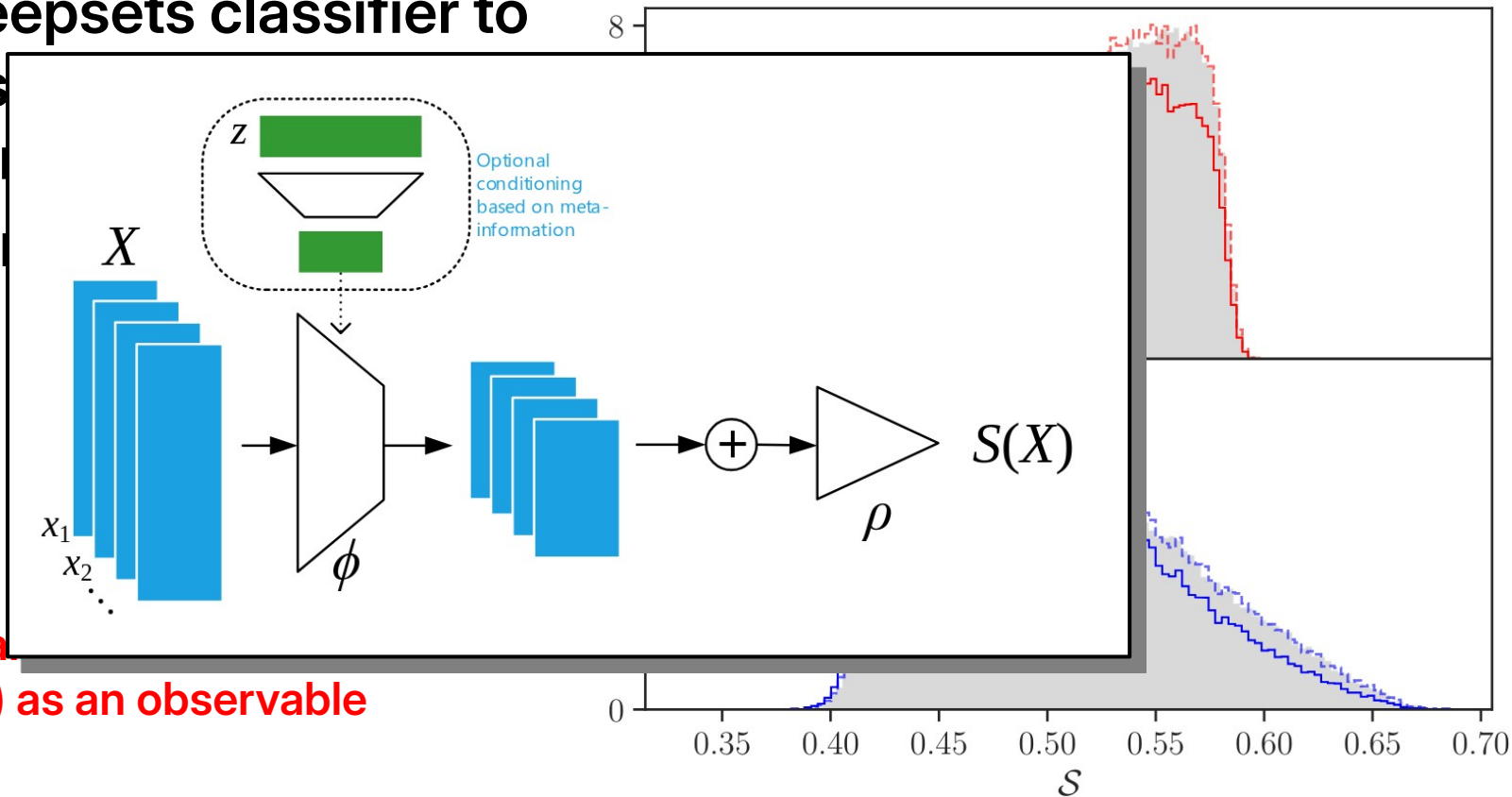


The "score" observable

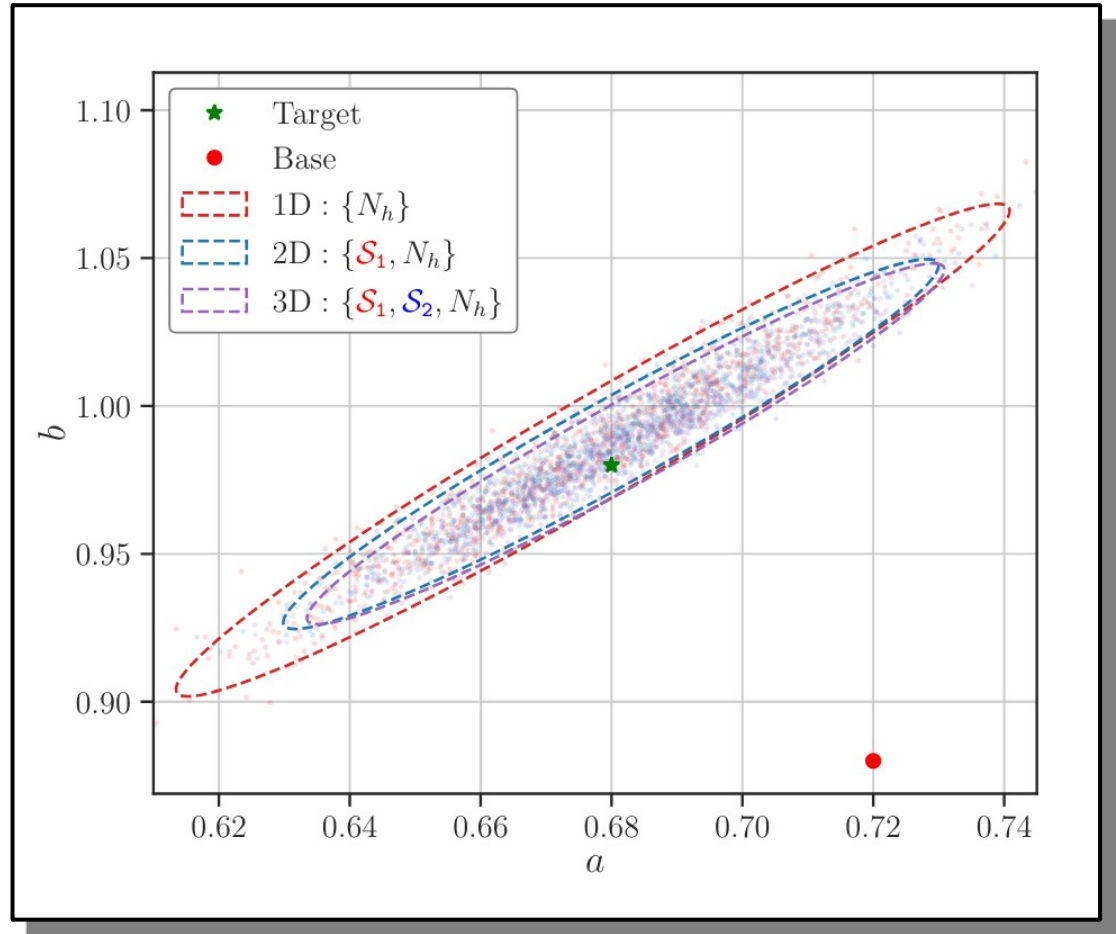
- Train a deepsets classifier to distinguish experimental

- Takes full input

Use training data (score) as an observable



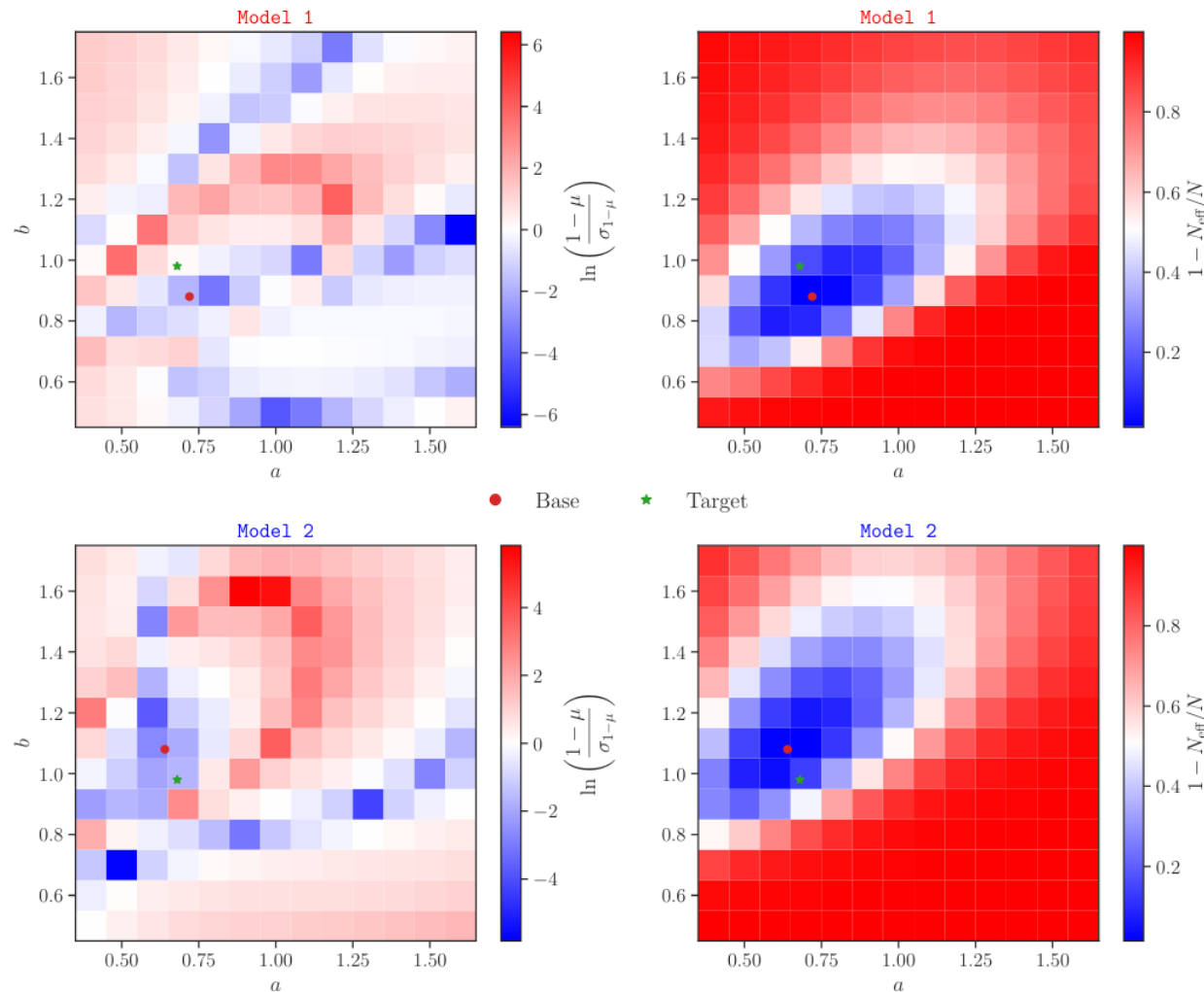
Classifier score with full event
info improves tuning
convergence



No free lunch

Statistical power drops off quickly as you move away from the base parameterization

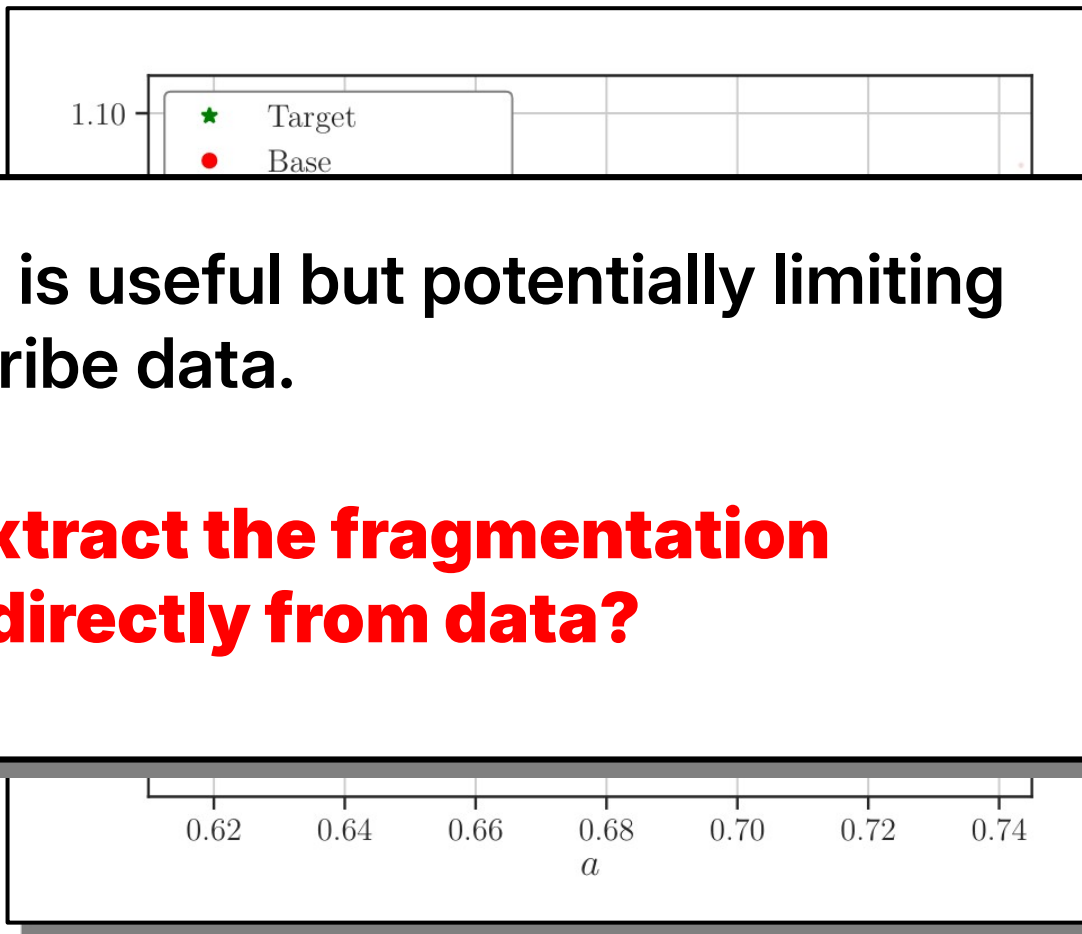
$$\mu \equiv \sum_{i=1}^N \frac{w_i}{N}, \quad N_{\text{eff}} = \frac{\left(\sum_{i=1}^N w_i\right)^2}{\sum_{i=1}^N w_i^2}$$

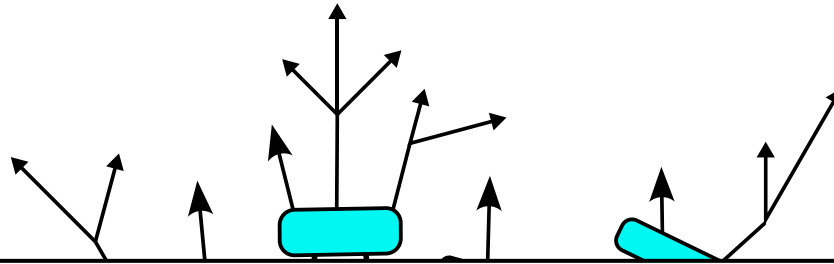


Clas
info
con

Analytic form of $f(z)$ is useful but potentially limiting when trying to describe data.

Can we extract the fragmentation function directly from data?



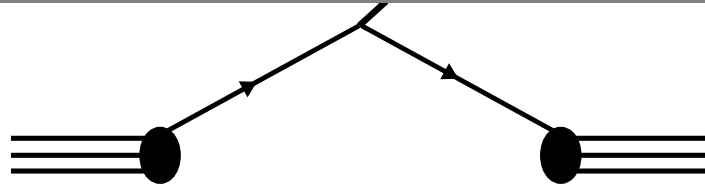


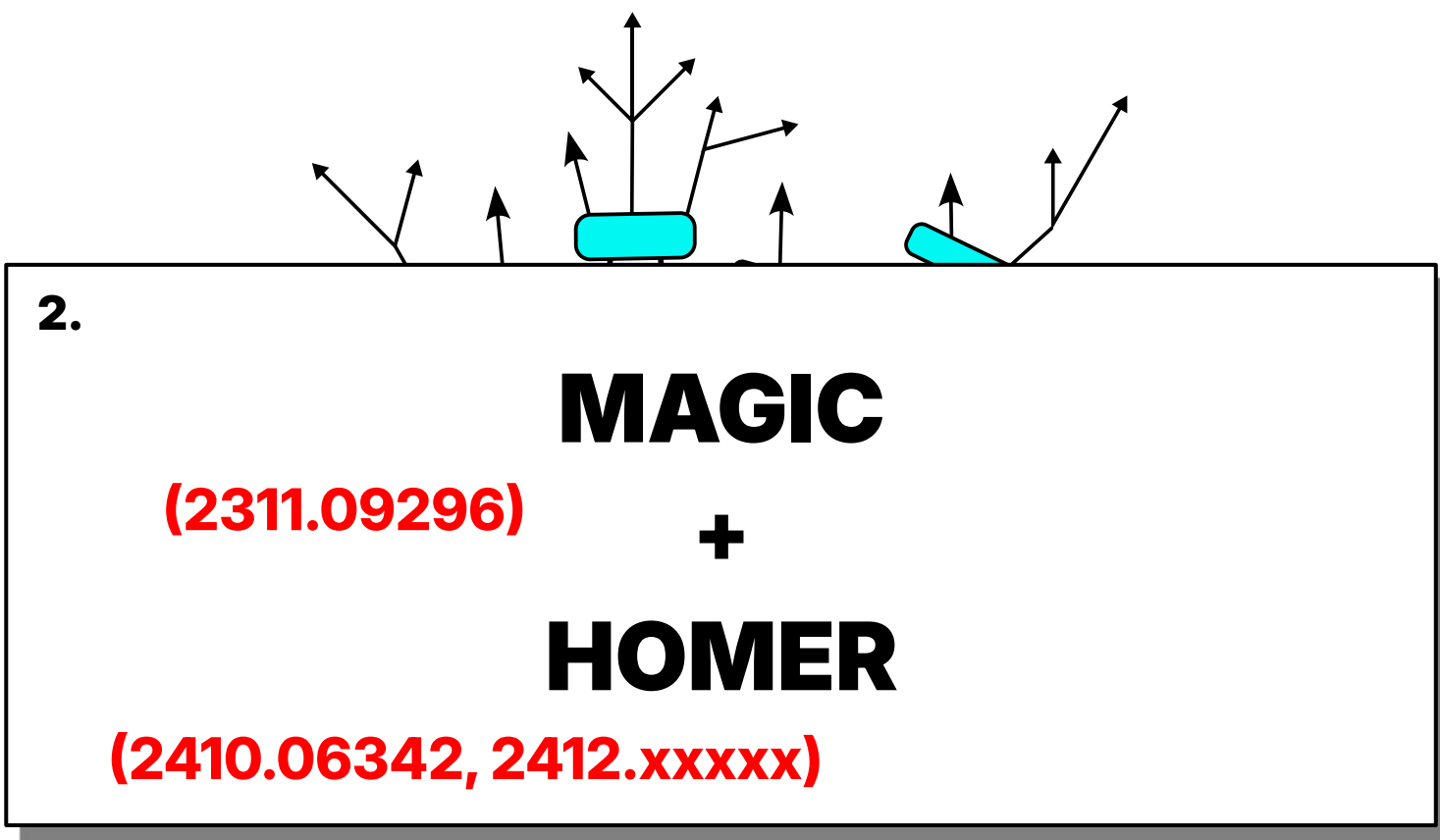
2. Microscopic Alterations Generated from IR Collections

(2311.09296) +

Histories and Observables for Monte Carlo Event Reweighting

(2410.06342, 2412.xxxxx)





Two methods

MAGIC

- "Top-down"
- Relies on the construction of an explicit likelihood (normalizing flow)
- Must be embedded into the simulation pipeline

HOMER

- "Bottom-up"
- Relies on the reconstruction of a classifier weight (derived from the score) from a product of string-break weights
- No modifications to the simulation pipeline required

MAGIC: Invertible neural networks (INN)

a.k.a normalizing flow

Invertible Real NVP transformations:

$$z'_1 = z_1 \odot \exp(s_1(z_2)) + t_1(z_2),$$

$$z'_2 = z_2 \odot \exp(s_2(z'_1)) + t_2(z'_1),$$

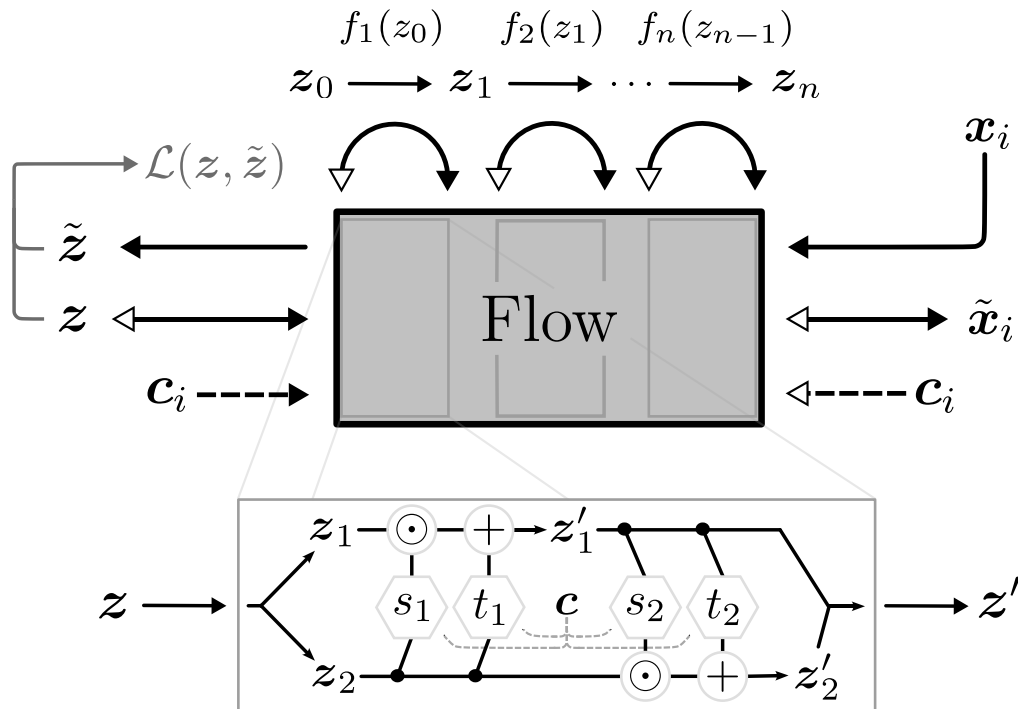
Scale transform

Translation transform

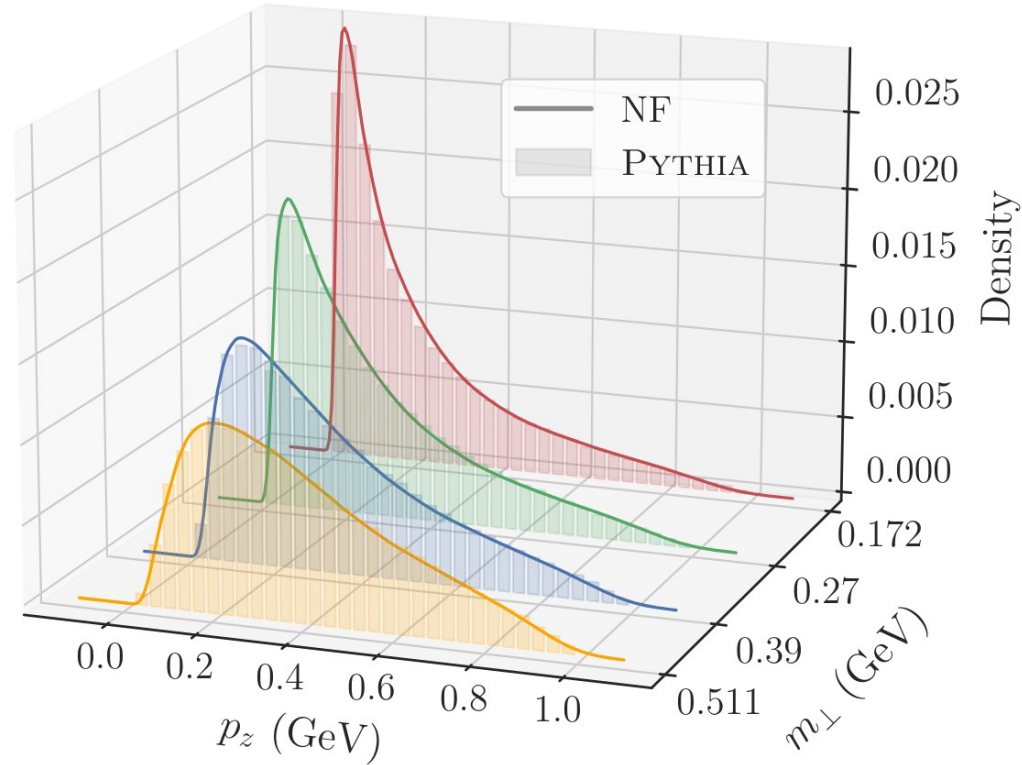
Inverse:

$$z_2 = (z'_2 - t(z'_1)) \odot \exp(-s(z'_1)),$$

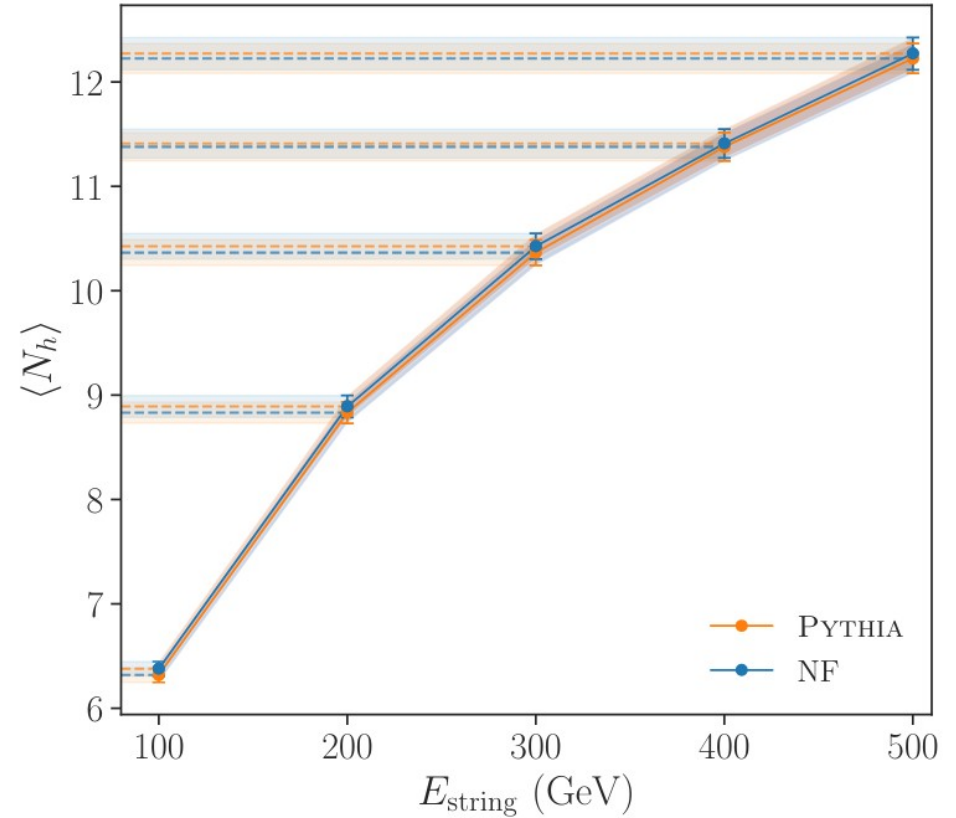
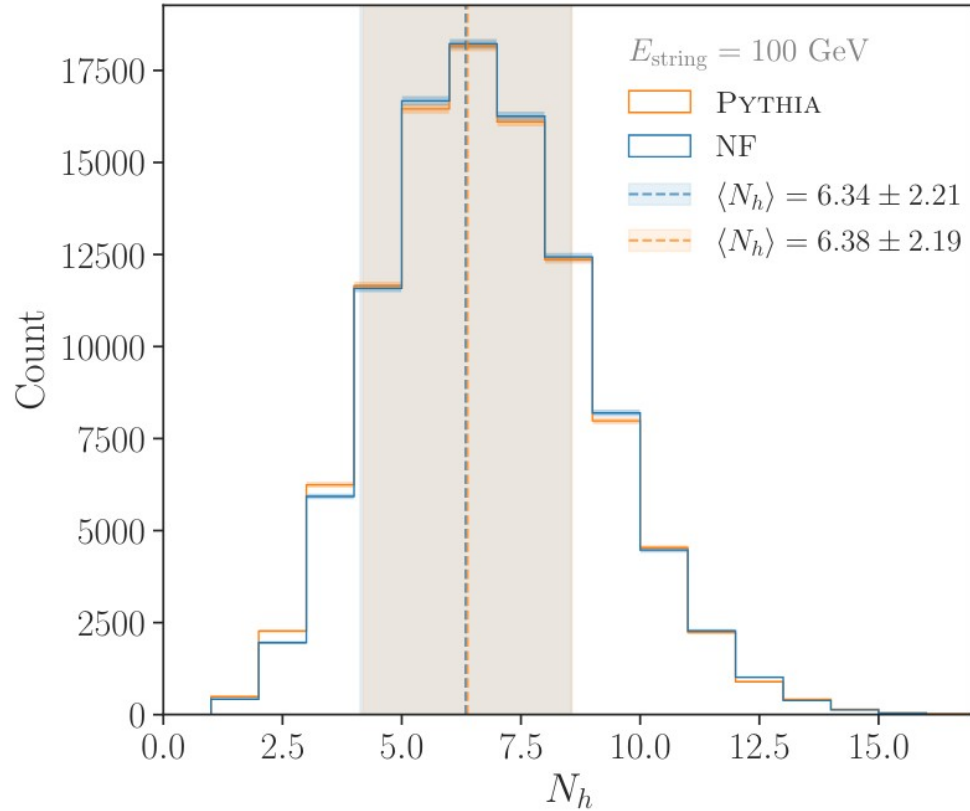
$$z_1 = (z'_1 - t(z_2)) \odot \exp(-s(z_2)),$$



MAGIC: Base model



MAGIC: Base model



MAGIC: $q\bar{q}$

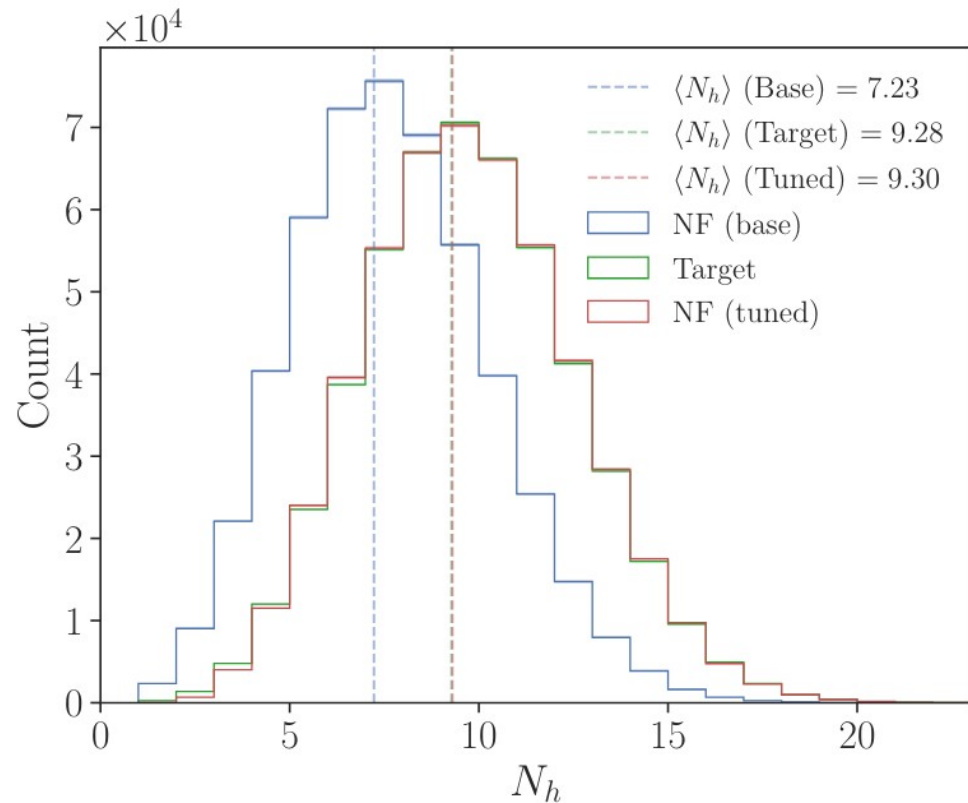
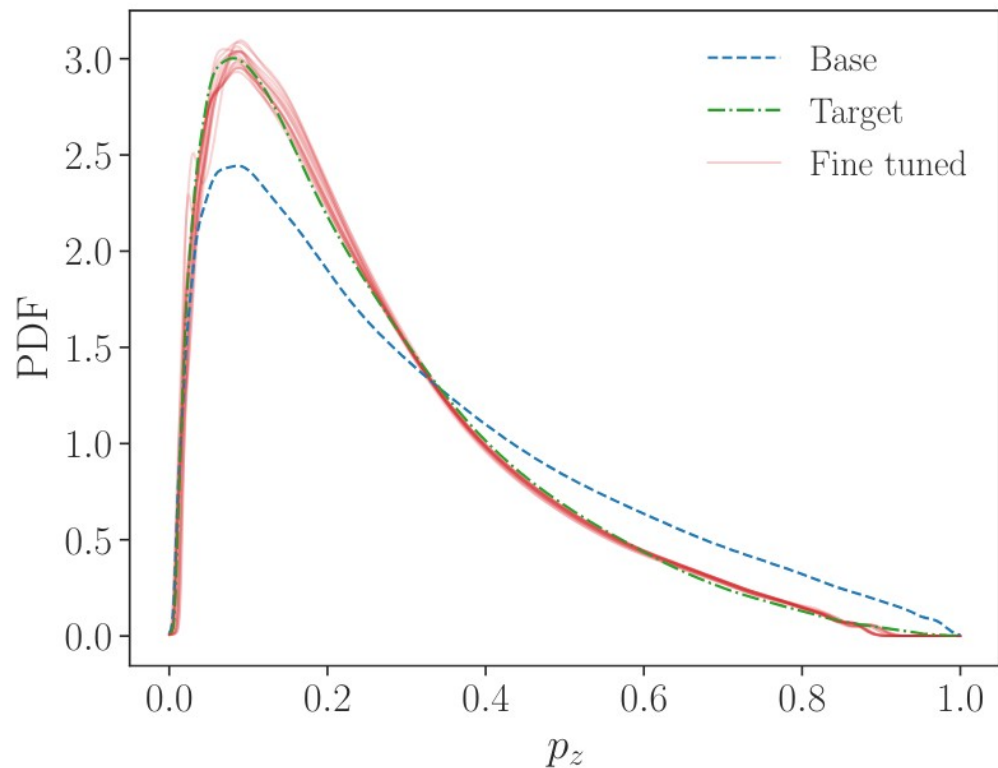
Use hadron multiplicity N as global event-level observable

$$\mathbf{X} = \begin{pmatrix} [\{p_z^{h_1}, p_T^{h_1}\}, \{p_z^{h_2}, p_T^{h_2}\}, \{p_z^{h_3}, p_T^{h_3}\}]_1 \\ [\{p_z^{h_1}, p_T^{h_1}\}, \{p_z^{h_2}, p_T^{h_2}\}, \{p_z^{h_3}, p_T^{h_3}\}, \{p_z^{h_4}, p_T^{h_4}\}, \{p_z^{h_5}, p_T^{h_5}\}]_2 \\ \vdots \\ [\{p_z^{h_1}, p_T^{h_1}\}, \{p_z^{h_2}, p_T^{h_2}\}]_n \end{pmatrix}, \quad \mathbf{Y}^{\text{sim}} = \begin{pmatrix} N_1 = 3 \\ N_2 = 5 \\ \vdots \\ N_n = 2 \end{pmatrix}, \quad \mathbf{Y}^{\text{exp}} = \begin{pmatrix} N_1 = 7 \\ N_2 = 3 \\ \vdots \\ N_n = 2 \end{pmatrix}$$

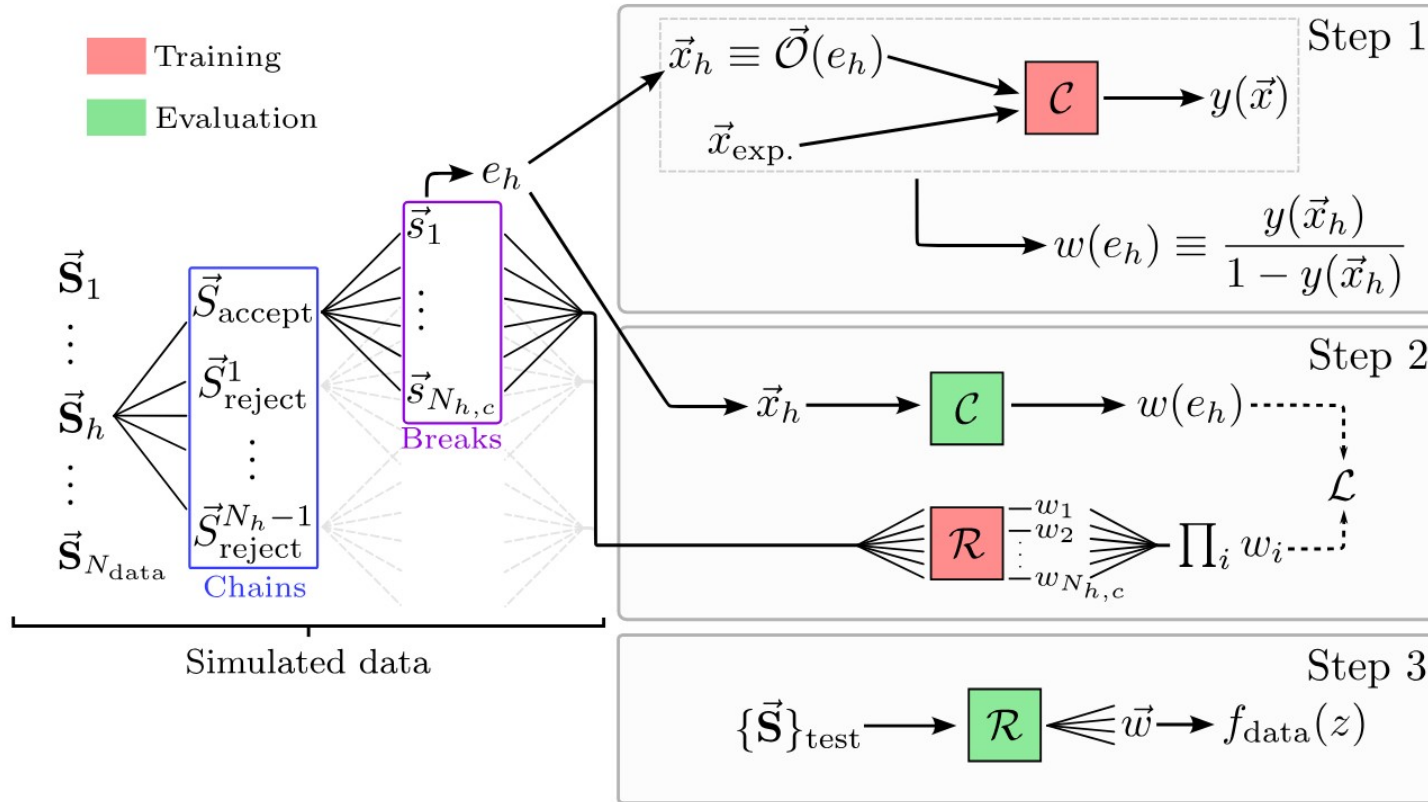
Event weights: $w = \begin{pmatrix} \prod_{i=1}^{N_1} w_i \\ \prod_{j=1}^{N_2} w_j \\ \vdots \\ \prod_{k=1}^{N_n} w_k \end{pmatrix}$ where $w_i = \frac{p_X^{F'}(p_z^{h_i}, p_T^{h_i})}{p_X^F(p_z^{h_i}, p_T^{h_i})}$,

Loss (Earth mover's distance): $\mathcal{L}_{\text{EMD}}(\mathbf{Y}^{\text{sim}}, \mathbf{w}, \mathbf{Y}^{\text{exp}})$

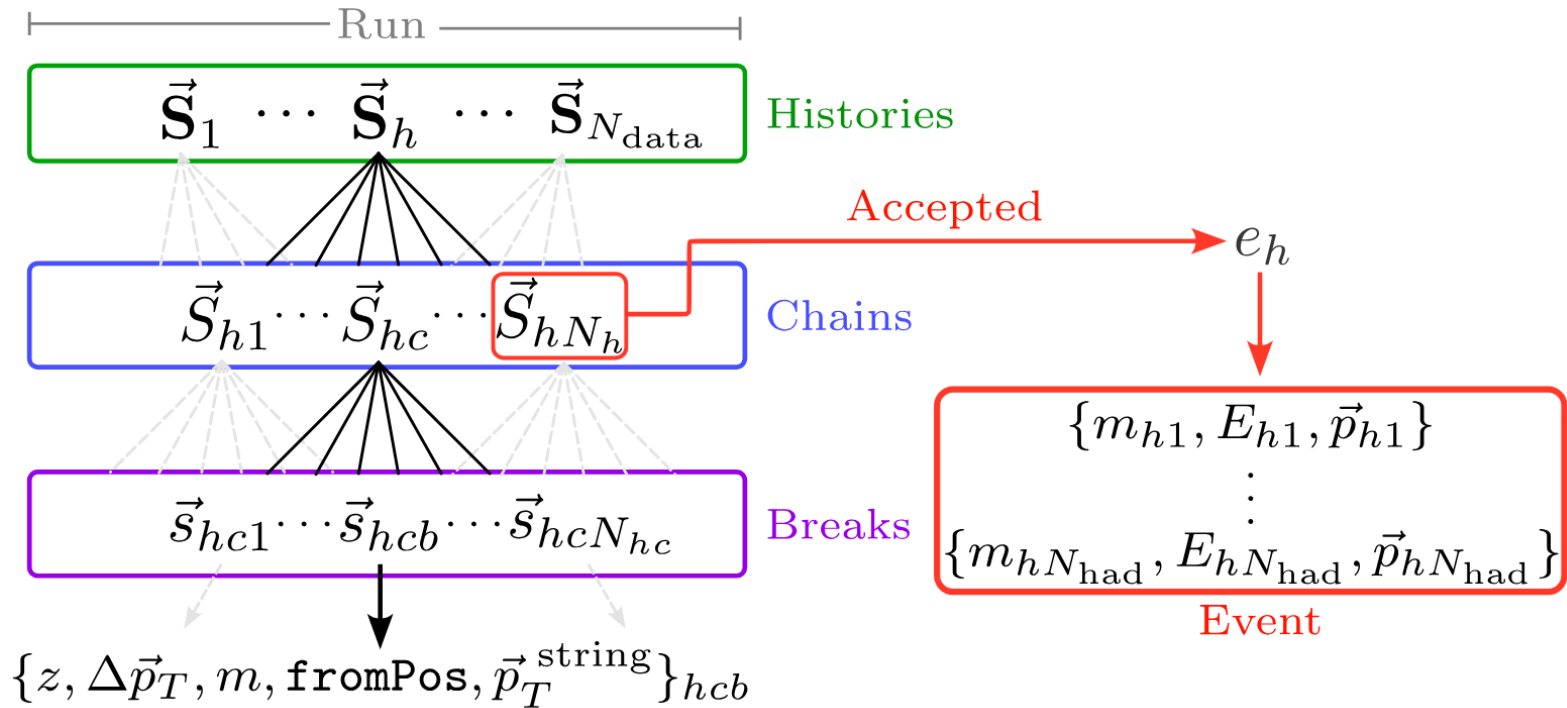
MAGIC: $q\bar{q}$



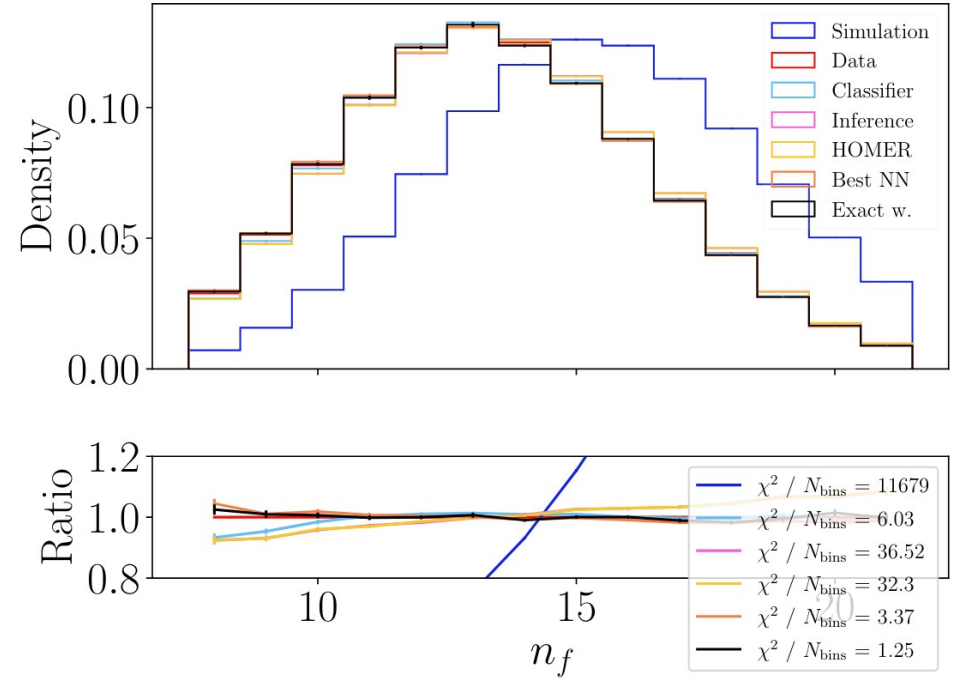
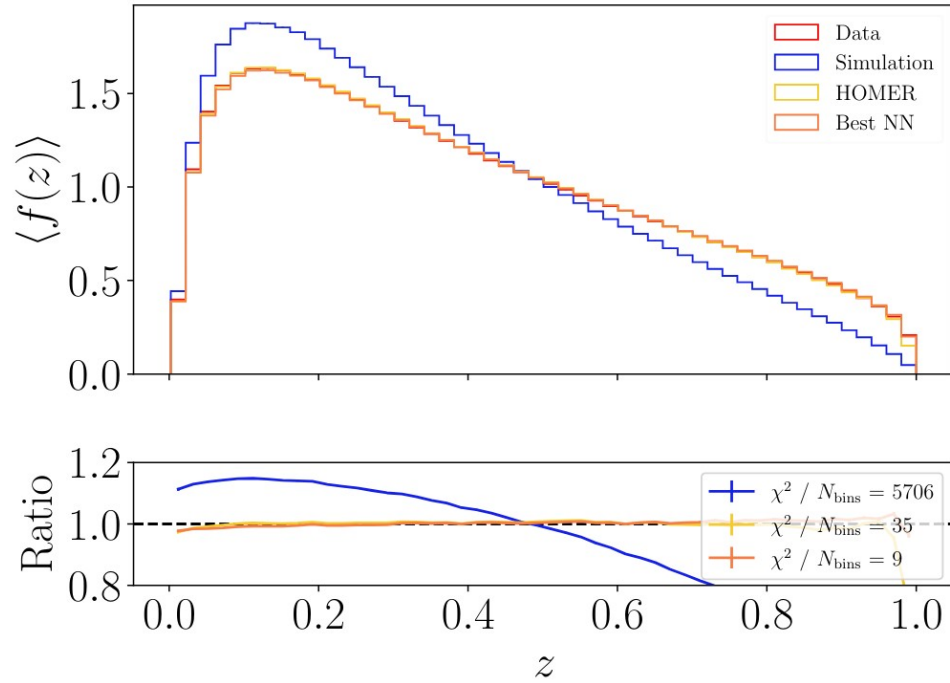
HOMER



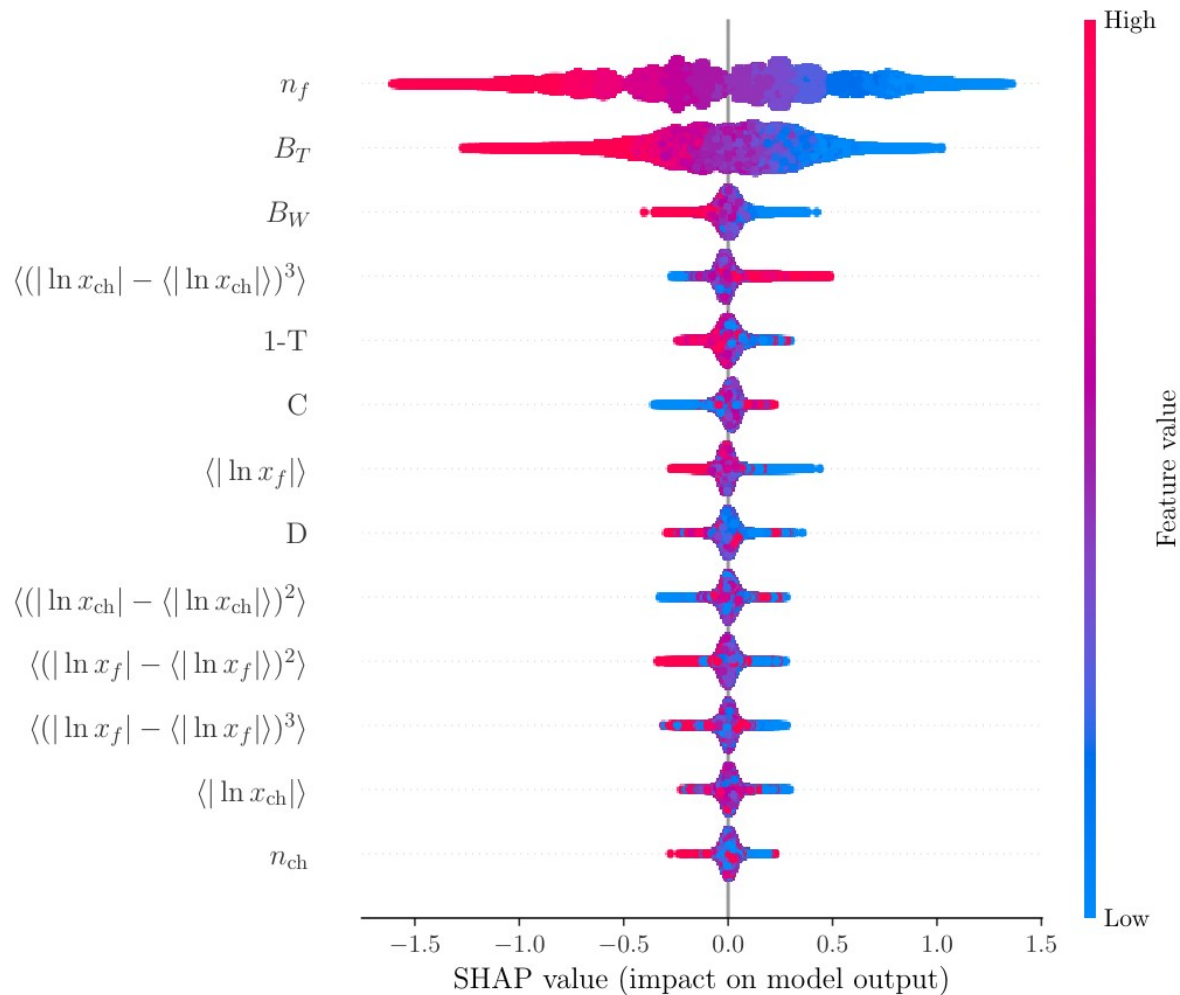
HOMER: $q\bar{q}$



HOMER: $q\bar{q}$



SHAP values



Conclusions

- **RSA: flexible and scalable approach to parameter estimation**
 - Fully differentiable simulations
 - Uncertainty quantification, flavor tune
- **Data-driven fragmentation function extraction:**
 - HOMER and MAGIC offer as complementary methods
 - HOMER + gluons (2412.xxxxx), flavor, tackle real data soon?
 - Better observables needed
 - HOMER x MAGIC
 - Better likelihood models? (MAGIC+)
 - Topology and simulation informed weights



<https://uclep.gitlab.io/mlhad-docs/>

Conclusions

- **RSA: flexible and scalable approach to parameter estimation**
 - Fully differentiable simulations
 - Uncertainty quantification, flavor tune
- **Data-driven fragmentation function extraction:**
 - HOMER and MAGIC offer as complementary methods
 - HOMER + gluons (2412.xxxxx), flavor, tackle real data soon?
 - Better observables needed
 - HOMER x MAGIC
 - Better likelihood models? (MAGIC+)
 - Topology and simulation informed weights



<https://uclep.gitlab.io/mlhad-docs/>

Thank you :)

Back-up

Stringy hadronization: overview

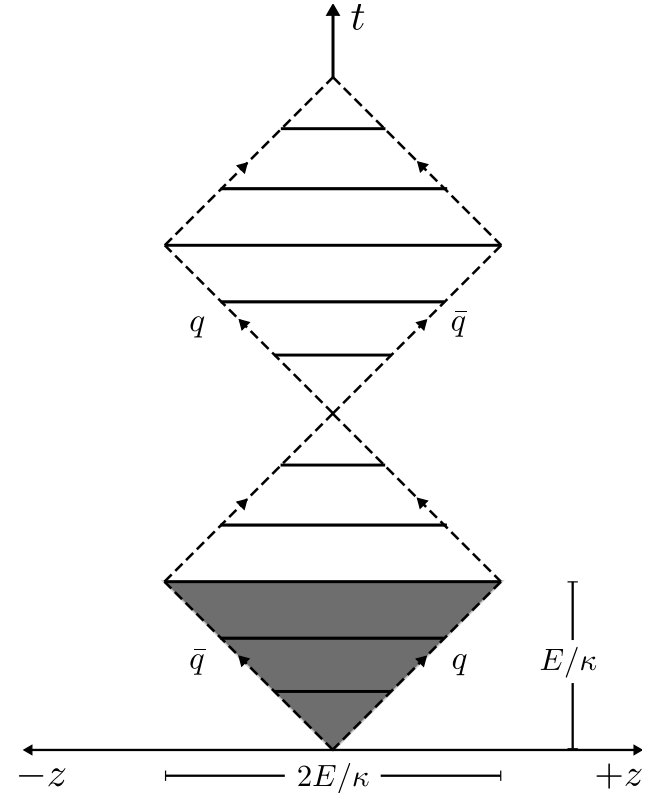
Consider the simplest hadronizing system:

A $q\bar{q}$ pair oriented along the z -axis, with equal and opposite momentum.

Treat this as a semi-classical system with potential

$$V(r) = \kappa r$$

A color flux tube forms between the $q\bar{q}$ pair and in the absence of string breaks the system follows a 'yo-yo' motion



Invertible neural networks (INN)

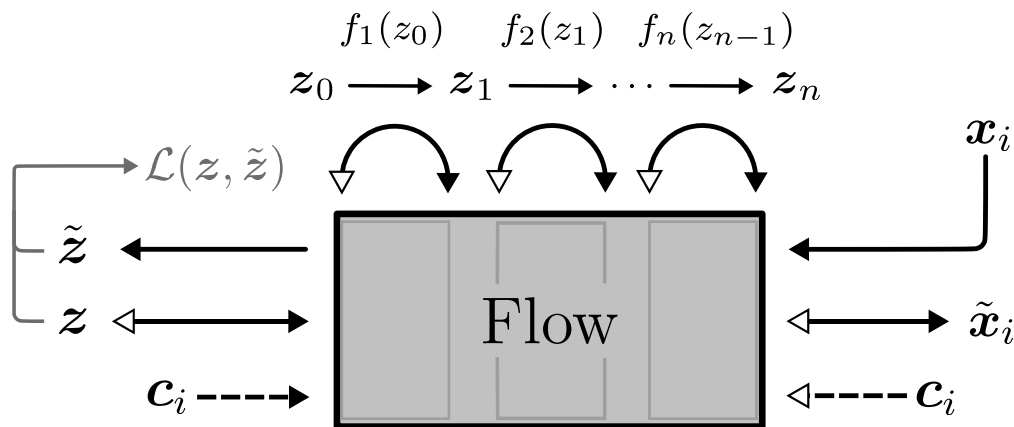
a.k.a normalizing flow

Learn a composition of n independent bijective transformations that relate a probability distribution $p_Z(\mathbf{z})$ on latent space Z to the target distribution $p_X(\mathbf{x})$ on target space X .

The probability distribution for the random variable $\mathbf{x} = f(\mathbf{z})$ is given by

$$p_X^f(\mathbf{x}) = p_Z(\mathbf{z}) |\det J_f(\mathbf{z})|^{-1}$$

$$J_f = \partial f / \partial \mathbf{x}$$



For n iterative transformations:

$$p_X^F(\mathbf{x}) = p_Z(\mathbf{z}_0) \prod_{i=1}^n |\det J_{f_i}(\mathbf{z}_{i-1})|^{-1}$$

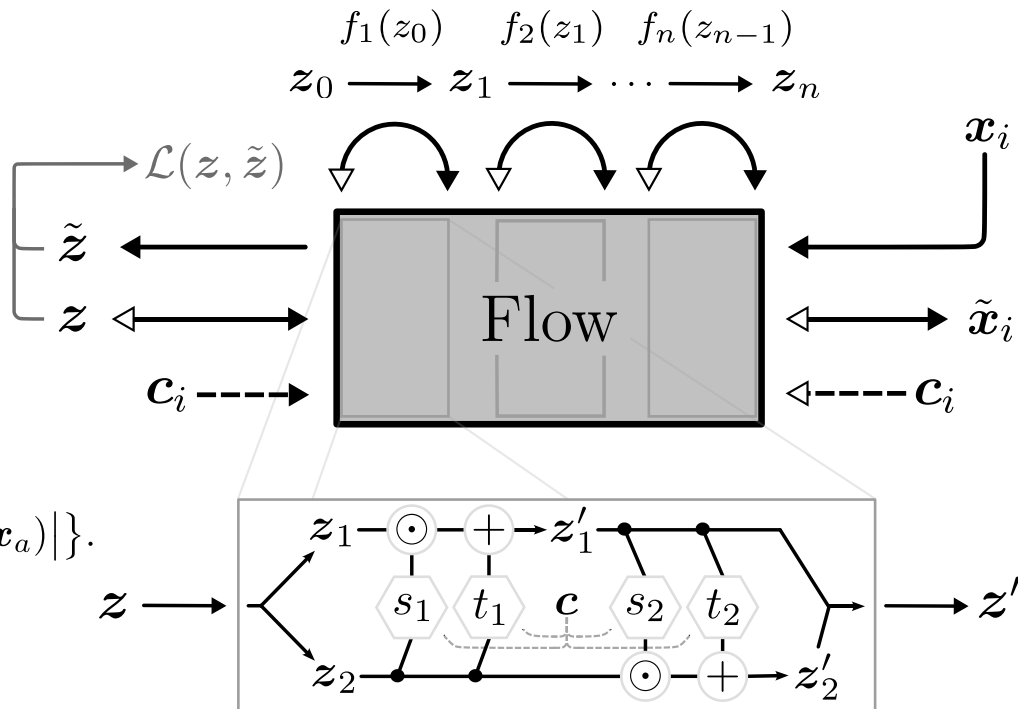
Invertible neural networks (INN)

a.k.a normalizing flow

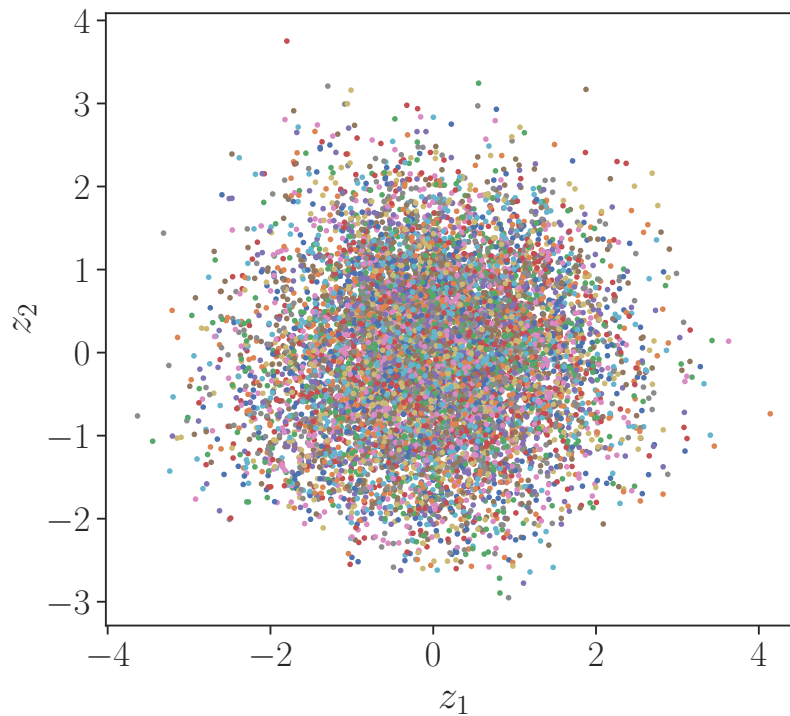
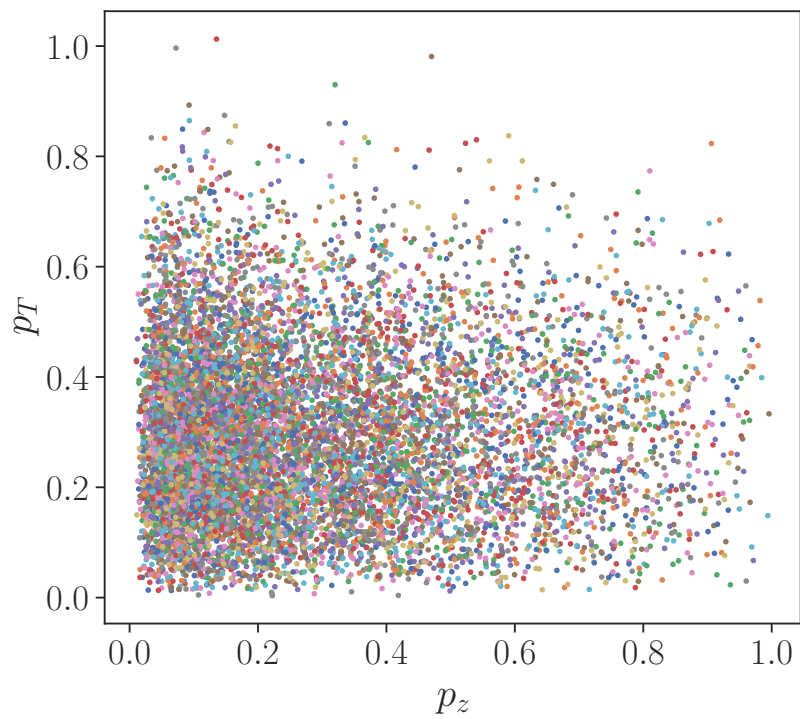
$$p_X^F(\mathbf{x}) = p_Z(\mathbf{z}_0) \prod_{i=1}^n |\det J_{f_i}(\mathbf{z}_{i-1})|^{-1}$$

Train with the negative log likelihood:

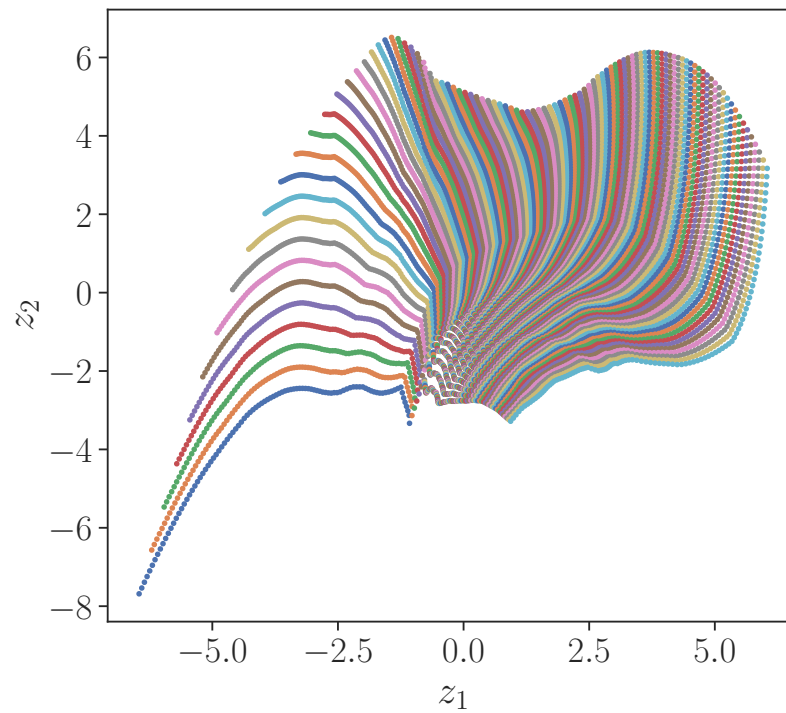
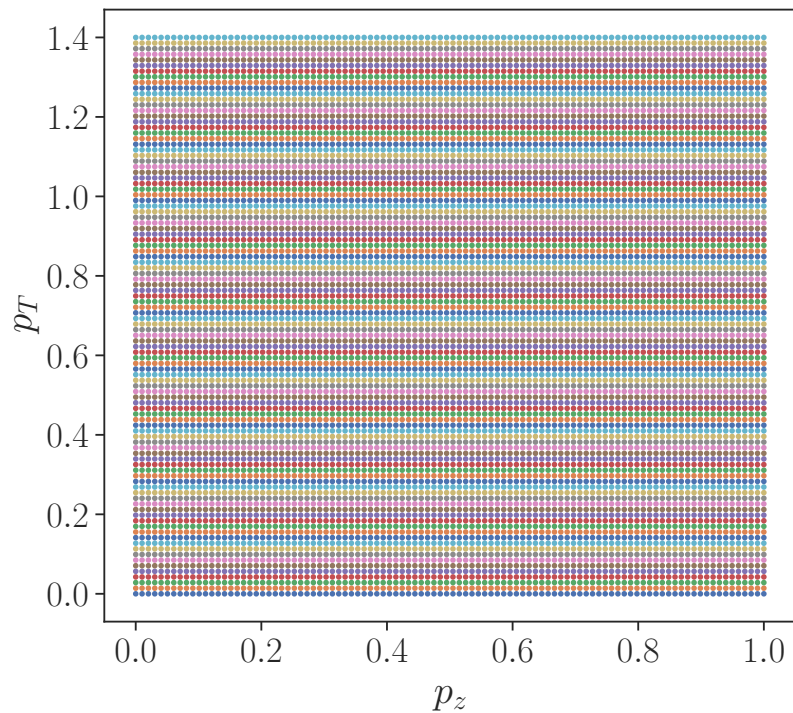
$$\begin{aligned} \mathcal{L}_{\text{NF}} &= - \sum_{a=1}^N \log p_X^F(\mathbf{x}_a; \boldsymbol{\theta}, \mathbf{c}_a) \\ &= \sum_{a=1}^N \{ - \log p_Z(F^{-1}(\mathbf{x}_a; \boldsymbol{\theta}, \mathbf{c}_a)) + \log |\det J_{F^{-1}}(\mathbf{x}_a)| \}. \end{aligned}$$



INN learned mapping

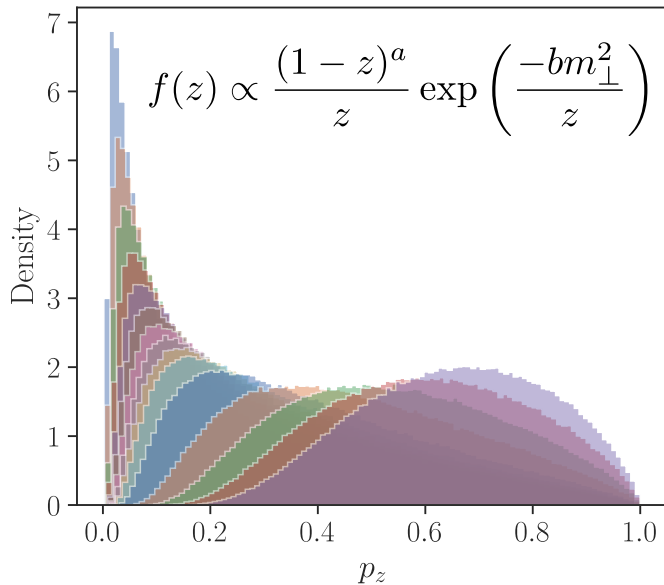


INN learned mapping

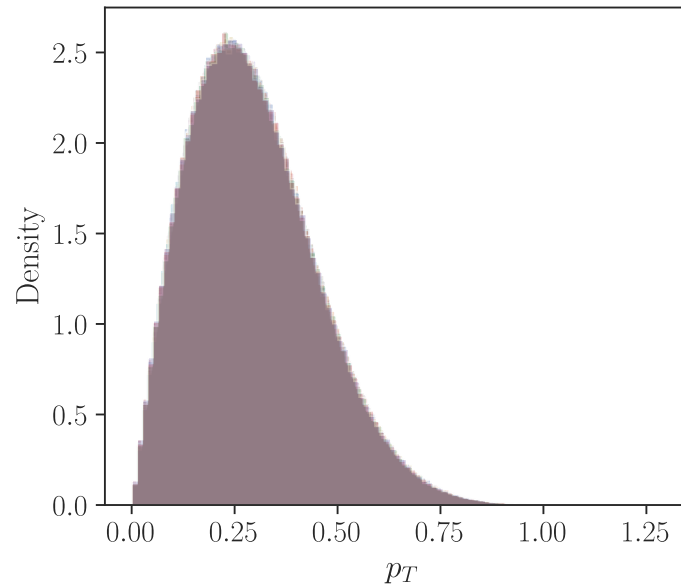


Training data

The implementation of full hadronization event using single emission kinematics requires an independent p_T sampling followed by a p_T -dependent sampling of p_z sampling (due to the dependence of $f(z)$ of transverse mass).



Tony Menzo (U. Cincinnati)



Data-driven hadronization models

Generate Pythia $q\bar{q} \rightarrow h$ events, first with no p_T kicks at different values of the hadron mass, record p_z . Generate events again, with kicks turned on, record p_T .