



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

France-Berkeley PHYSTAT Conference on Unfolding
10.06.2024

Generative Unfolding

SPONSORED BY THE

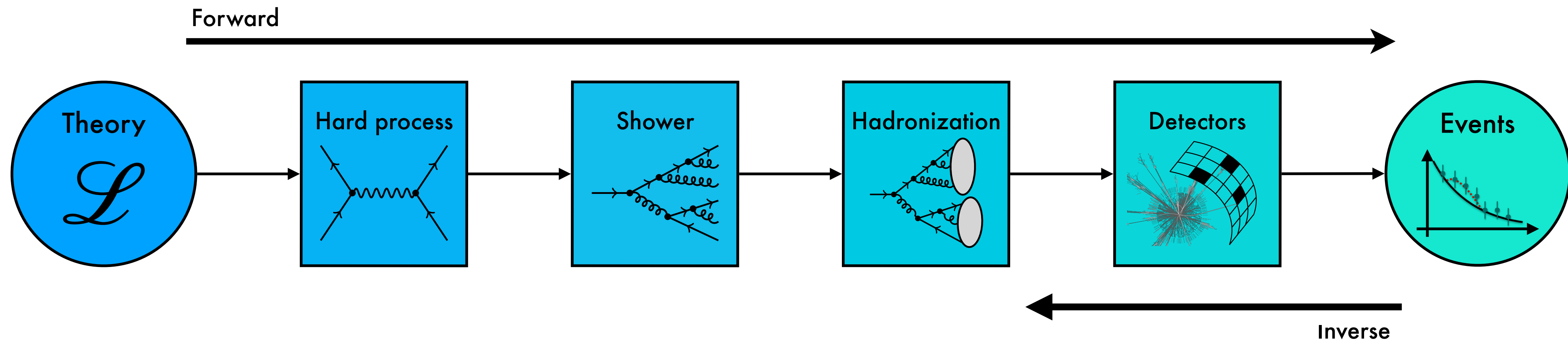


Federal Ministry
of Education
and Research

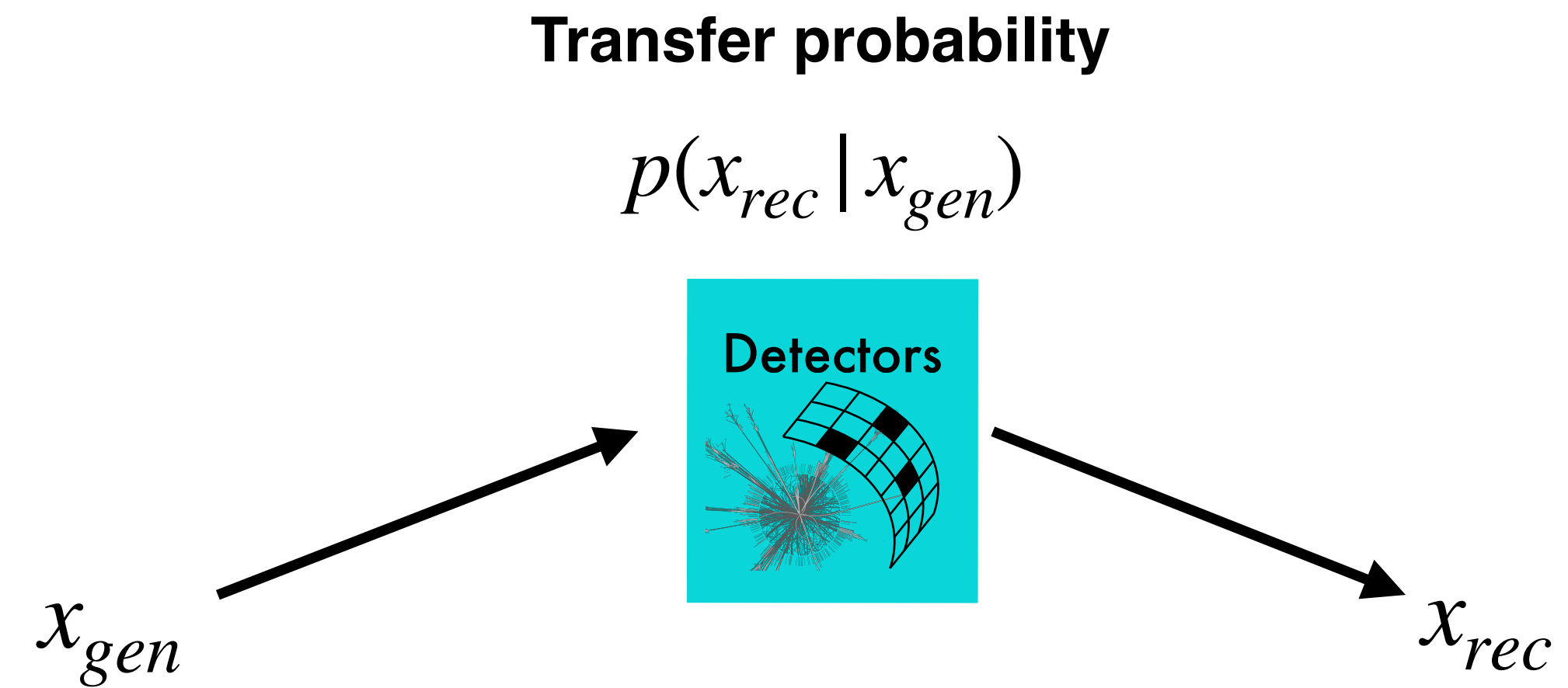
Nathan Huetsch

Huetsch et al. 2404.18807
The Landscape of Unfolding with Machine Learning

Simulation chain — Inversion

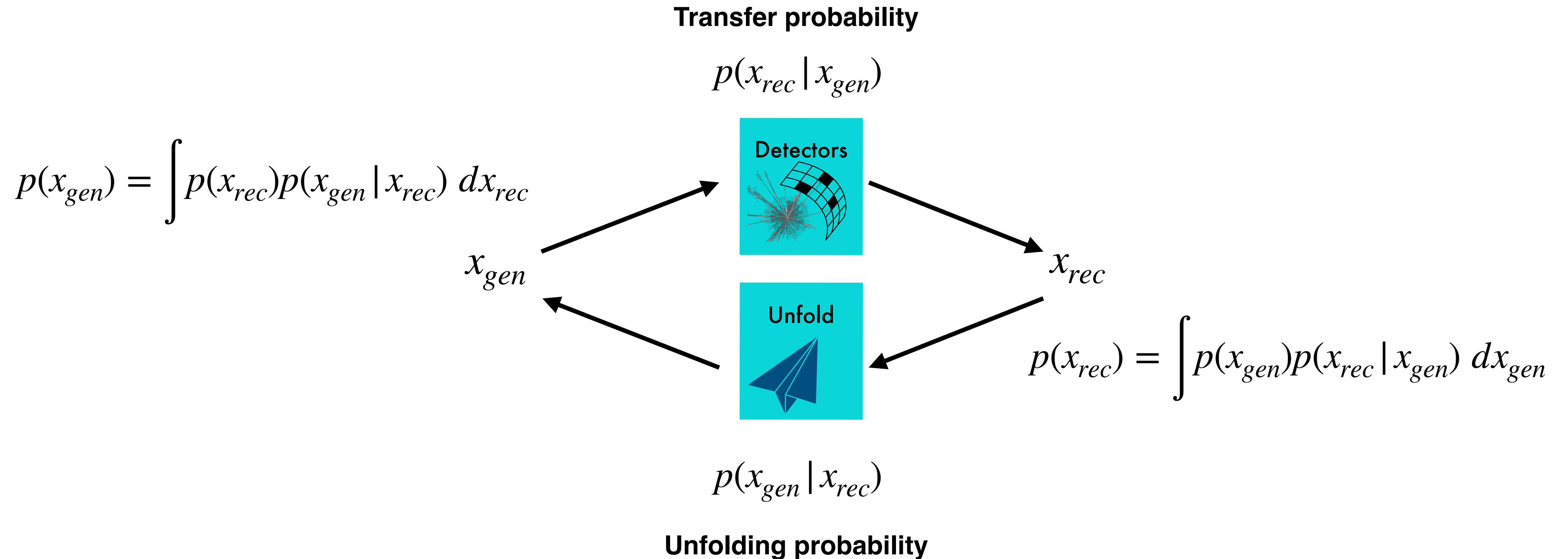


Probabilistic transfer

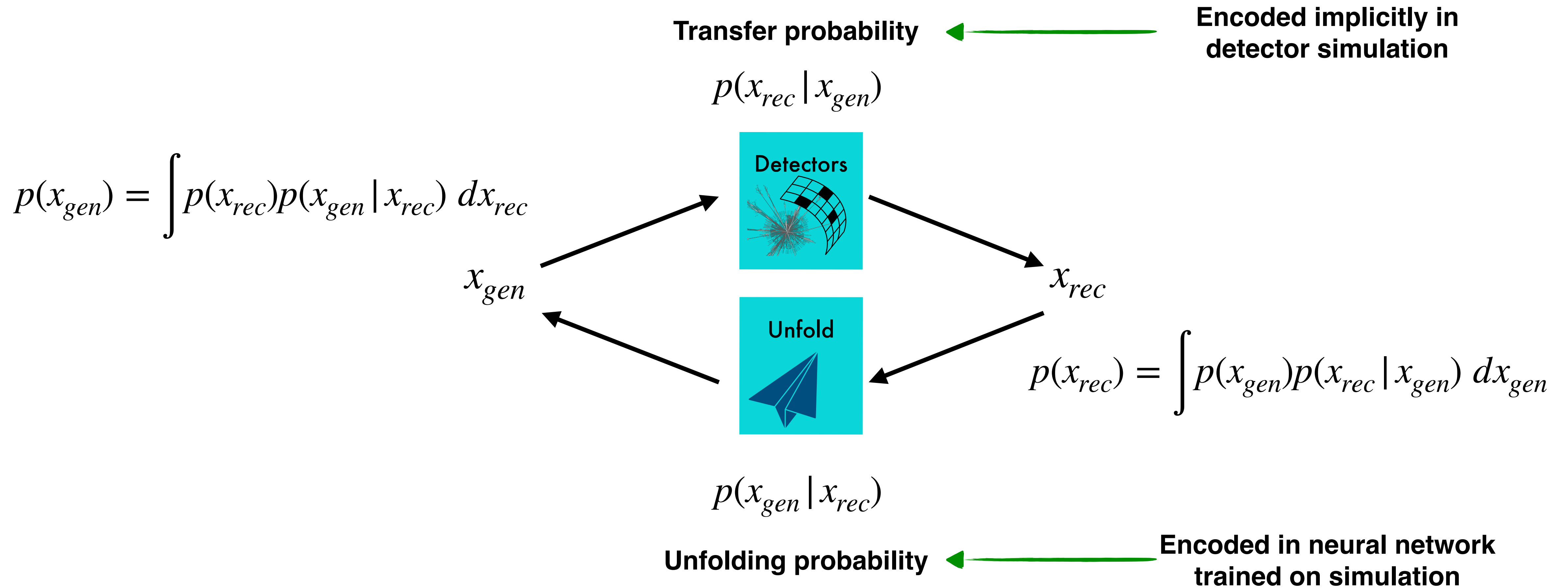


$$p(x_{rec}) = \int p(x_{gen})p(x_{rec} | x_{gen}) dx_{gen}$$

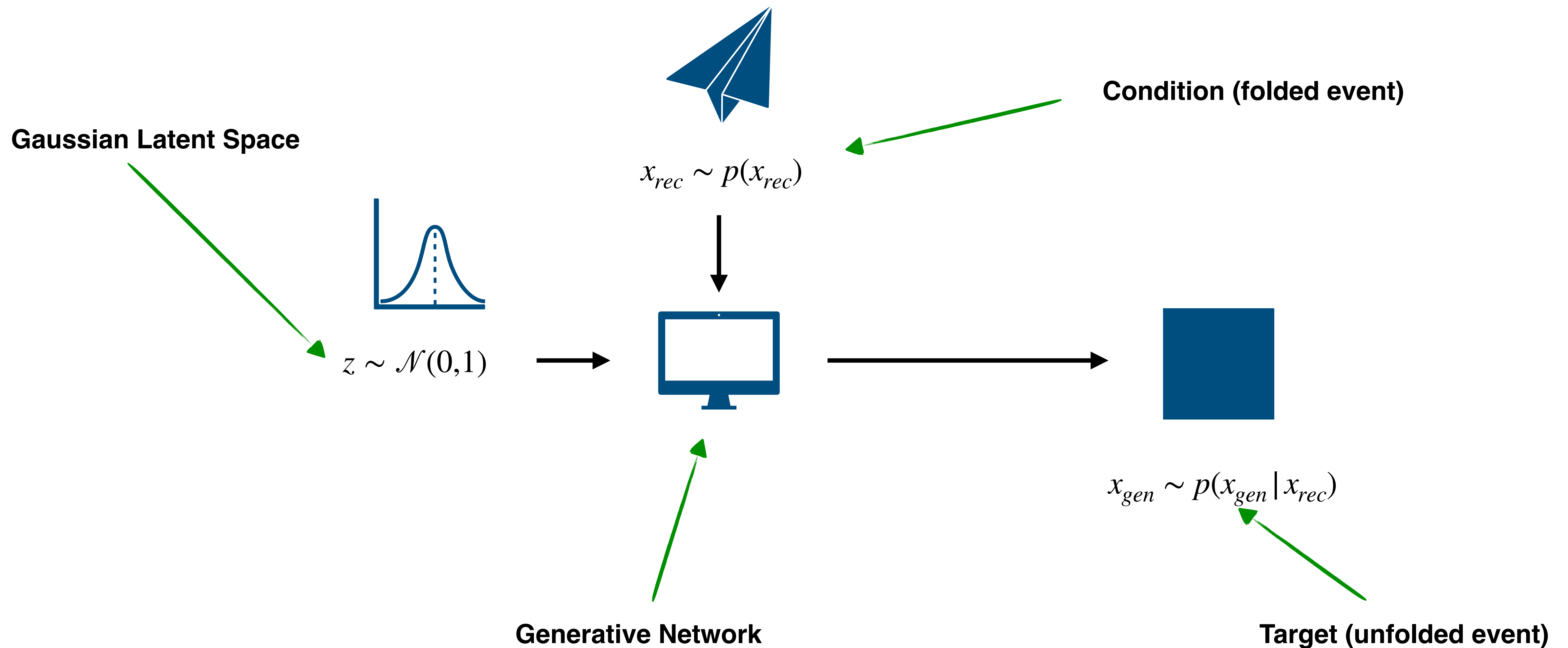
Generative unfolding



Generative unfolding



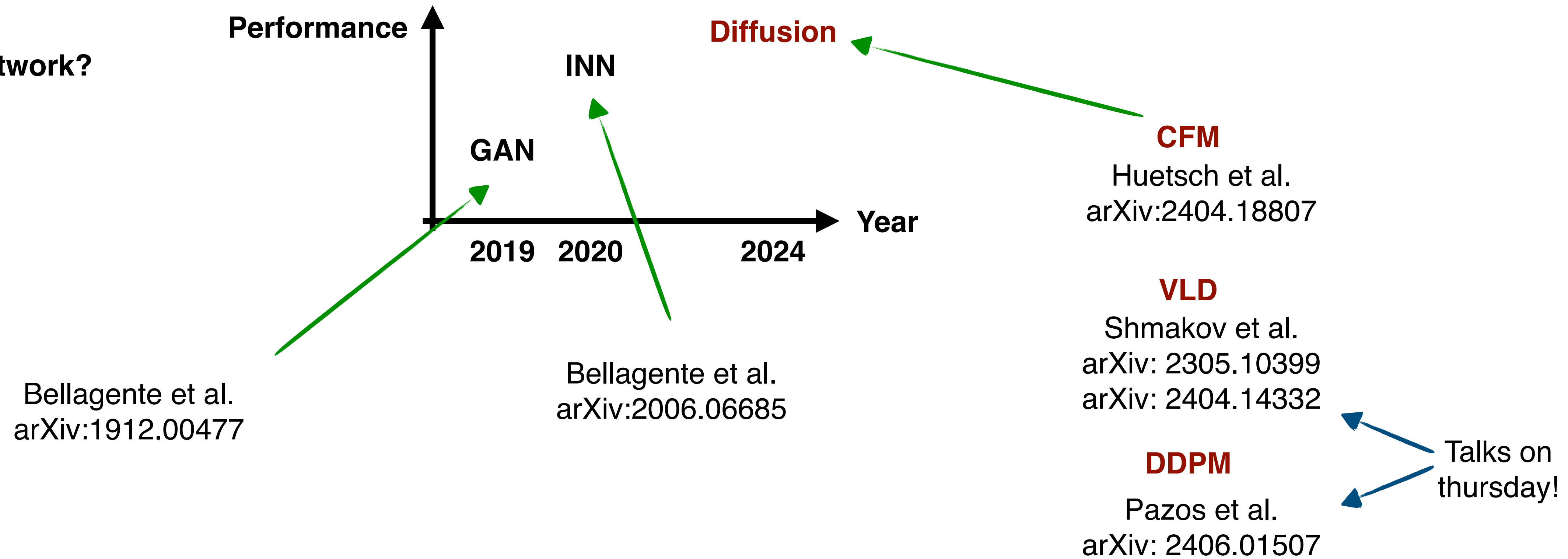
Conditional generative networks



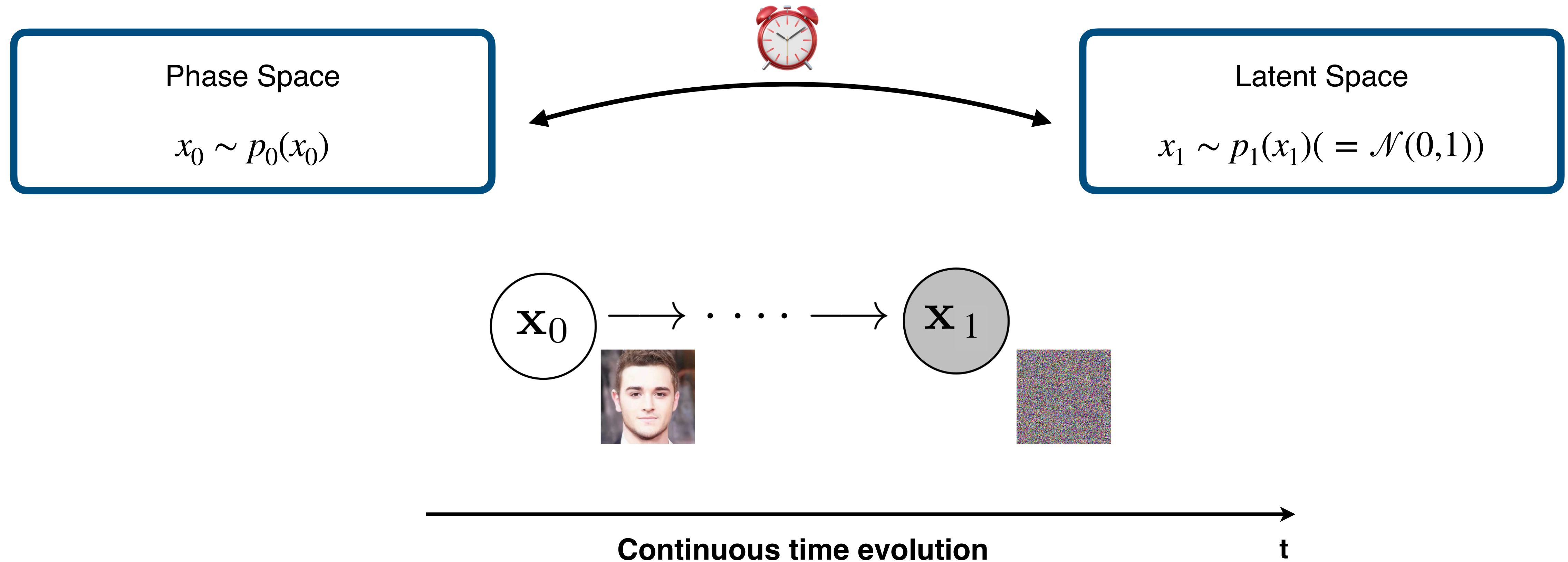
Conditional generative networks



Which generative Network?



Flow Matching (Lipman et al. 2210.02747)



Flow Matching (Lipman et al. 2210.02747)



Phase Space
 $t = 0$

Latent Space
 $t = 1$

Individual Samples
 $x_0 = x_{gen} \sim p_0(x_0)$

Individual Samples
 $x_1 = \epsilon \sim p_1(x_1)$

$$\frac{dx_t}{dt} = v_\theta(x_t, t)$$



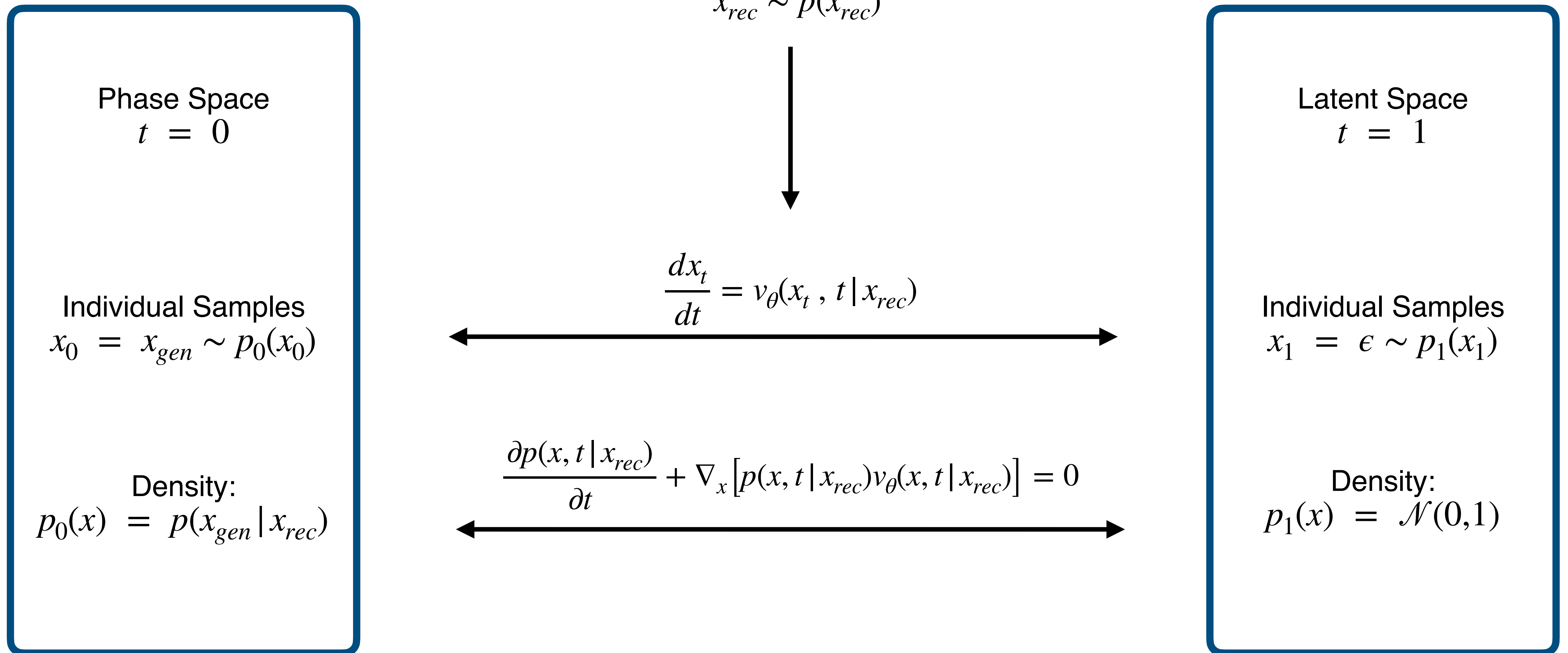
Density:
 $p_0(x) = p(x_{gen})$

Density:
 $p_1(x) = \mathcal{N}(0,1)$

$$\frac{\partial p(x, t)}{\partial t} + \nabla_x [p(x, t)v_\theta(x, t)] = 0$$

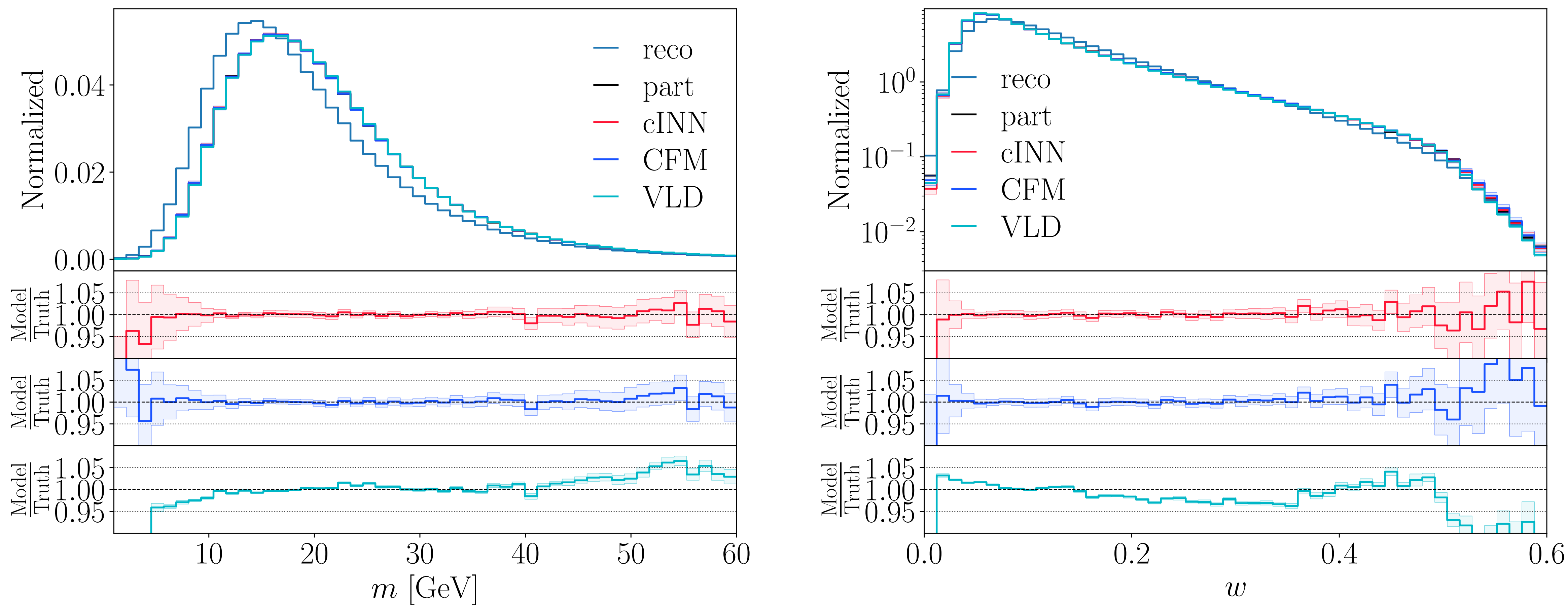


Flow Matching (Lipman et al. 2210.02747)

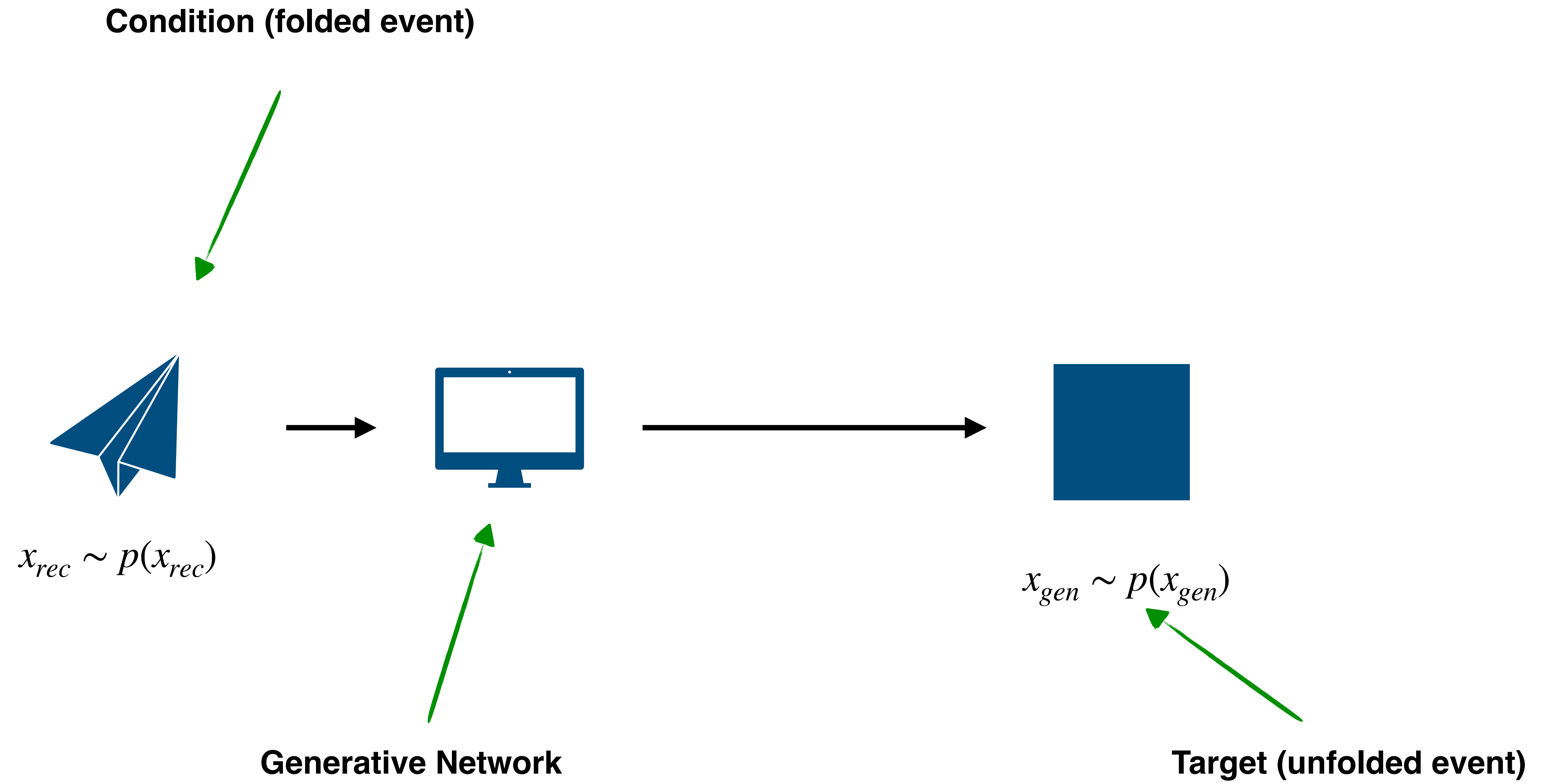


Z + jet results

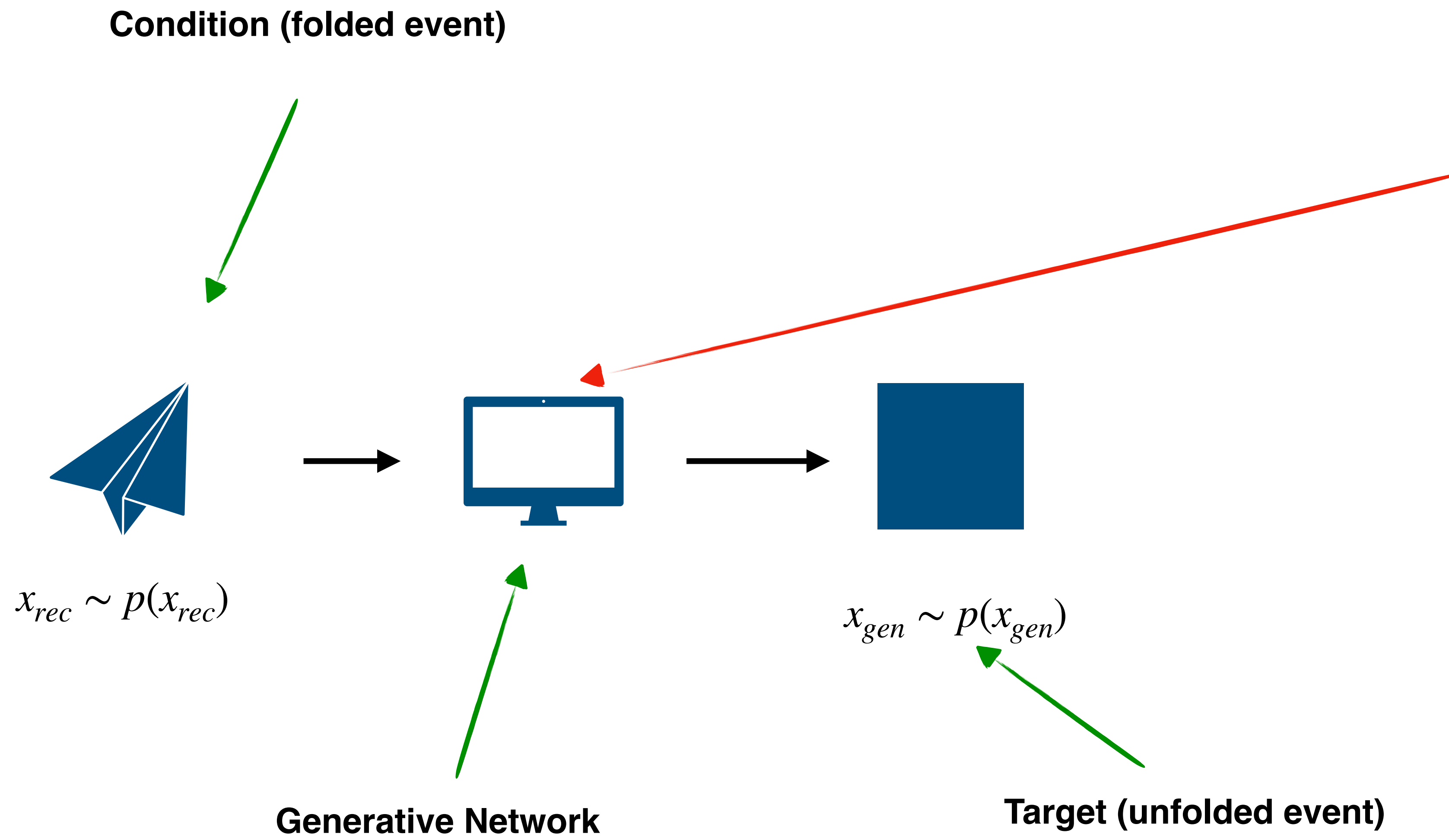
$pp \rightarrow Z + \text{jets}$ events following Andreassen et al. arXiv: 1911.09107



Distribution mapping



Distribution mapping



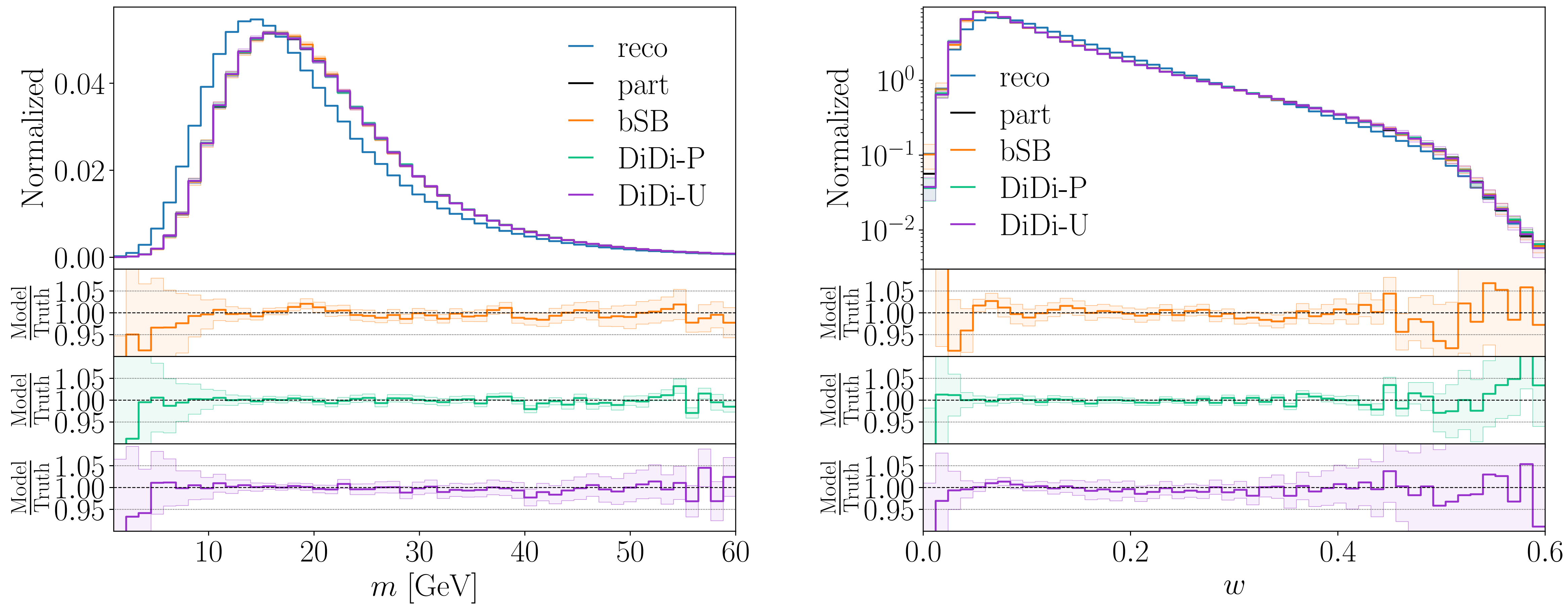
Schroedinger Bridge

Diefenbacher et al.
arXiv:2308.12351

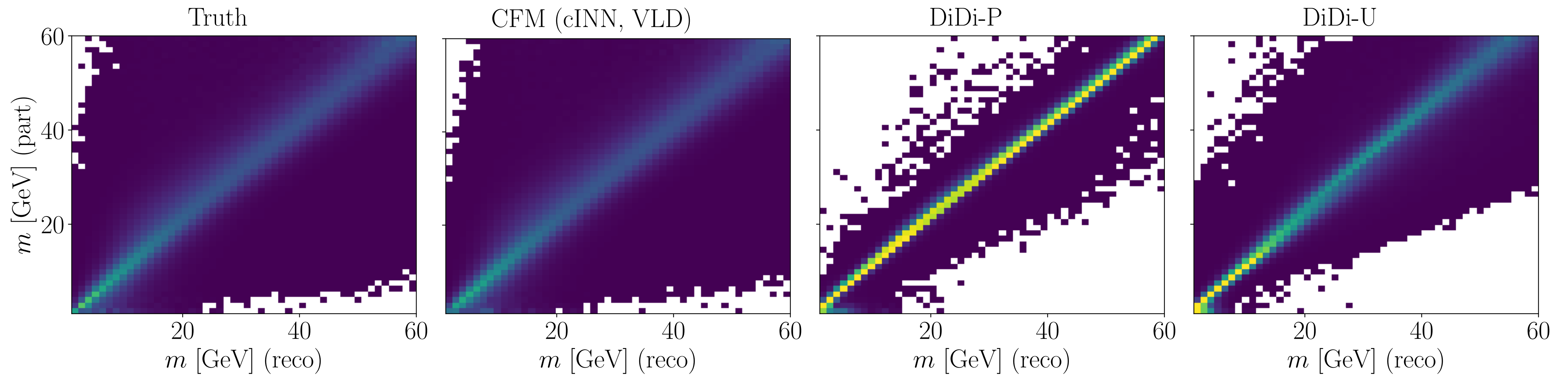
Direct Diffusion

Butter et al.
arXiv:2311.17175
Huetsch et al.
arXiv:2404.18807

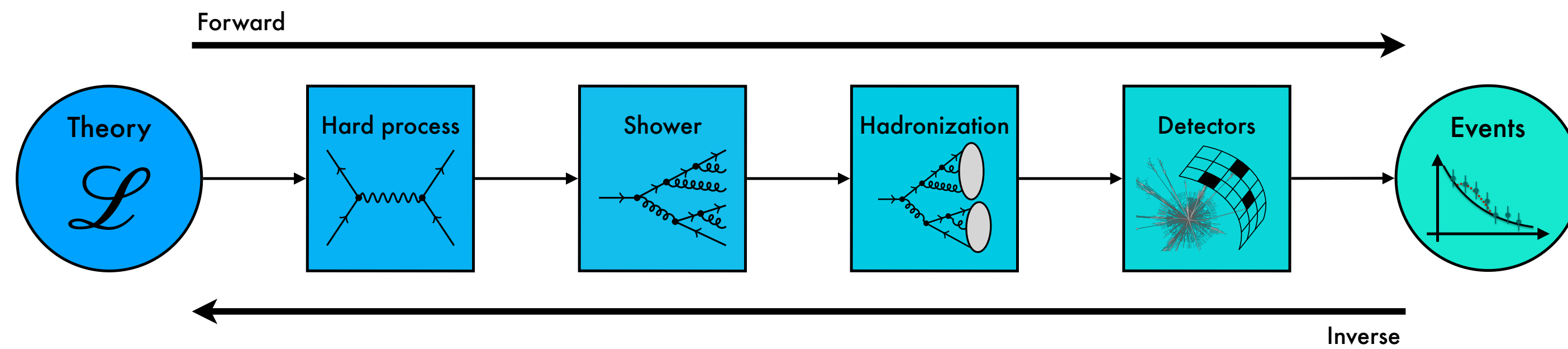
Z + jet results



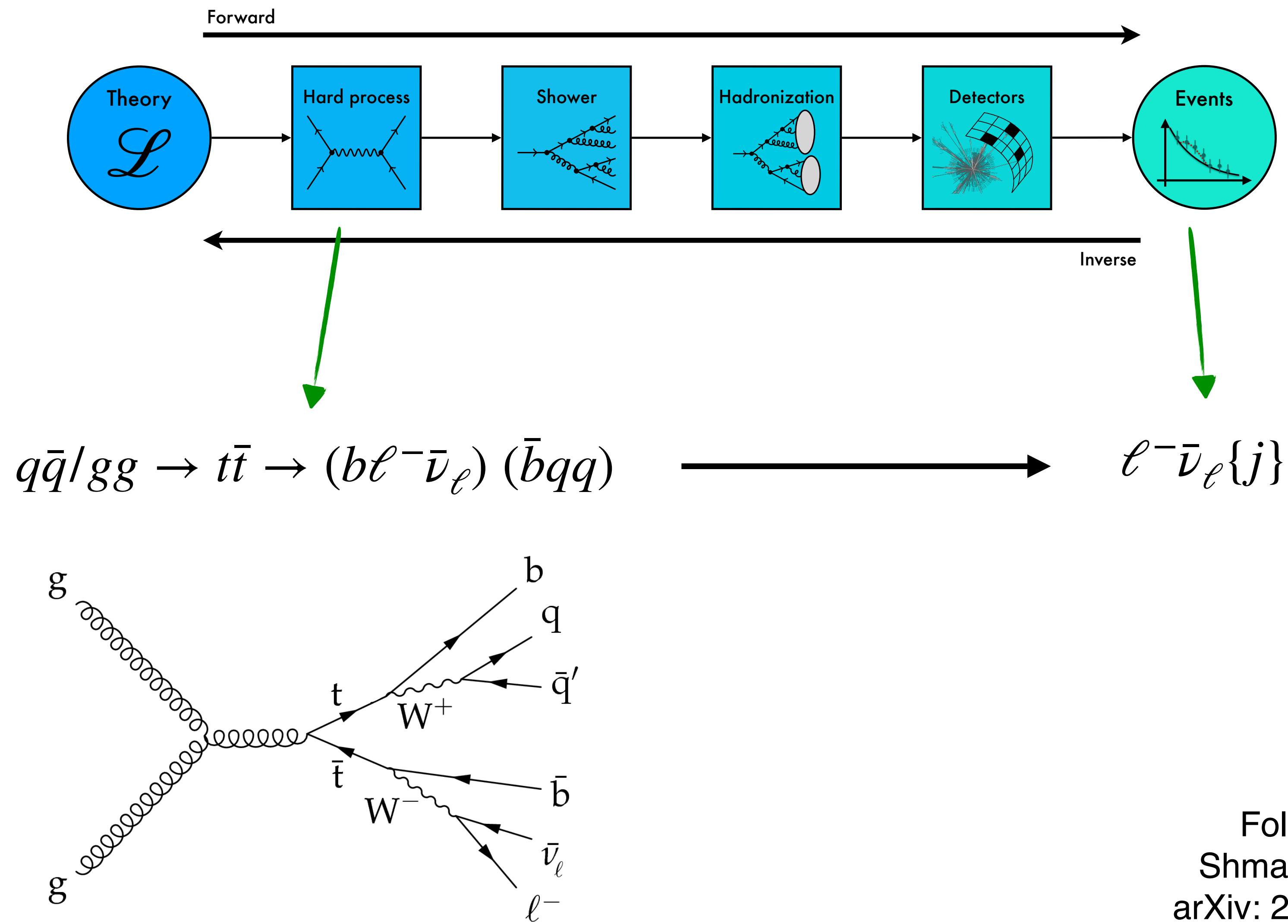
Z + jet migration



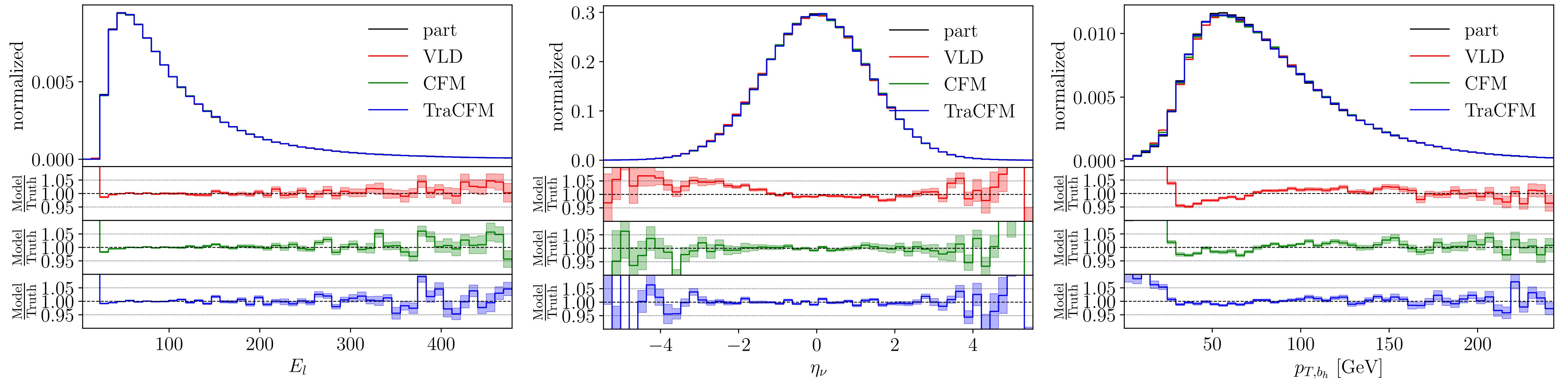
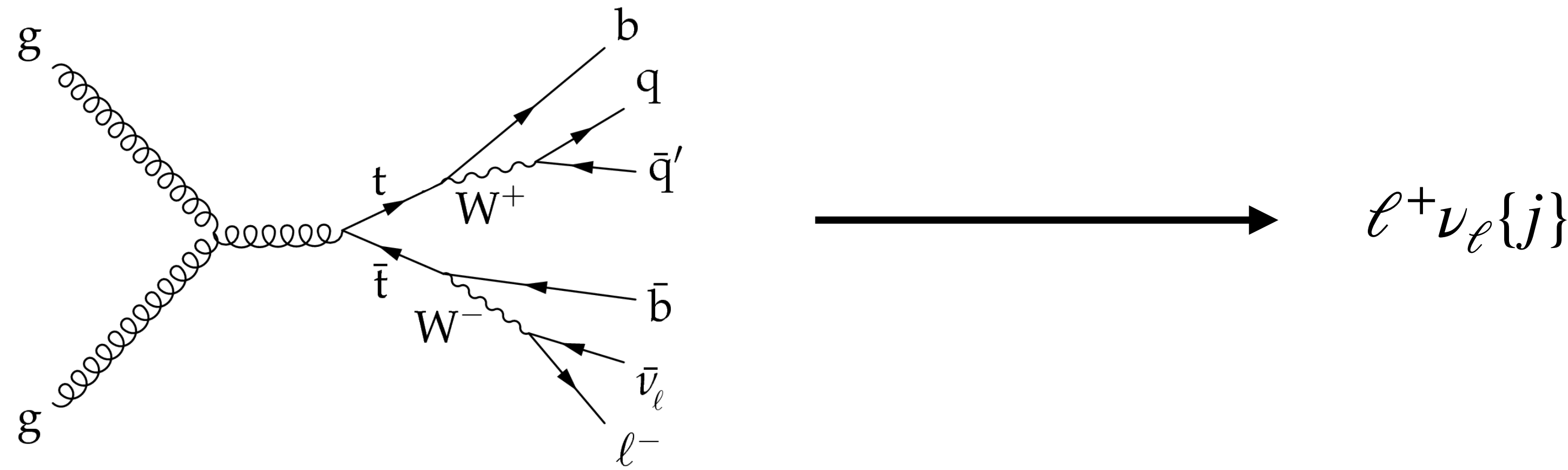
Parton-level unfolding — $t\bar{t}$ decay



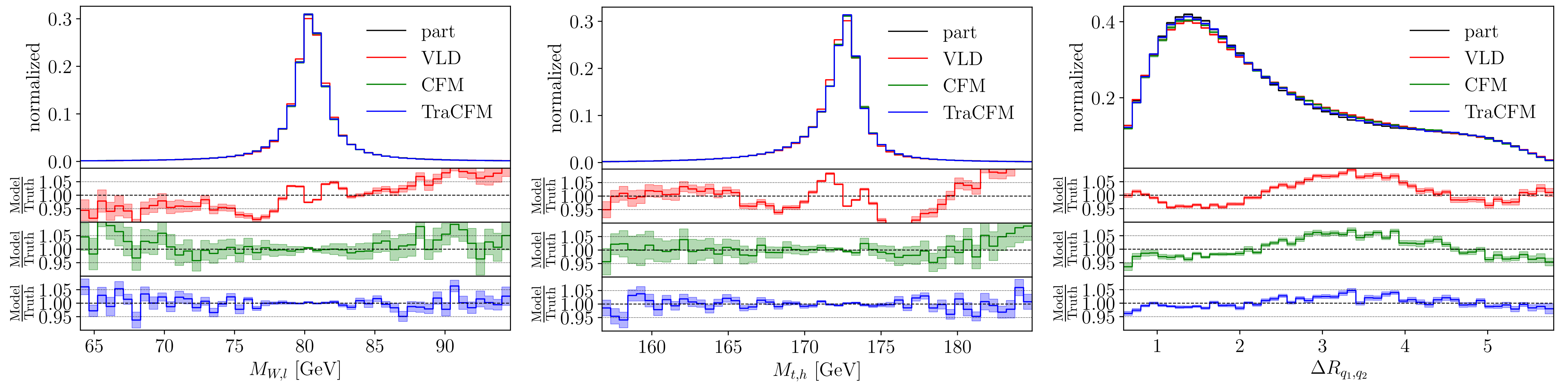
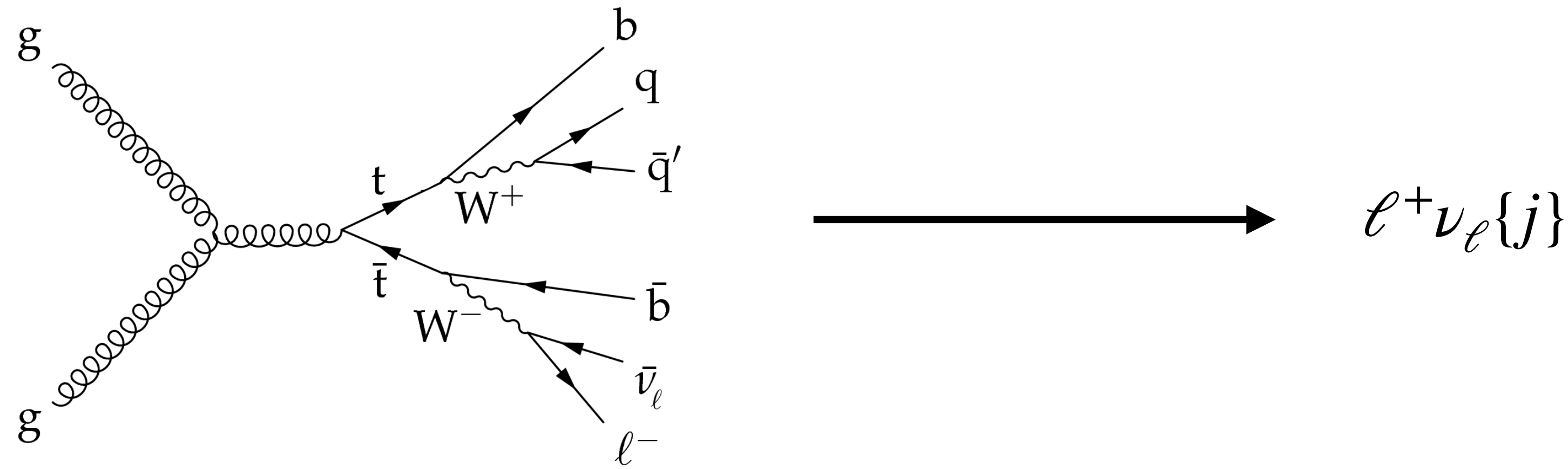
Parton-level unfolding — $t\bar{t}$ decay



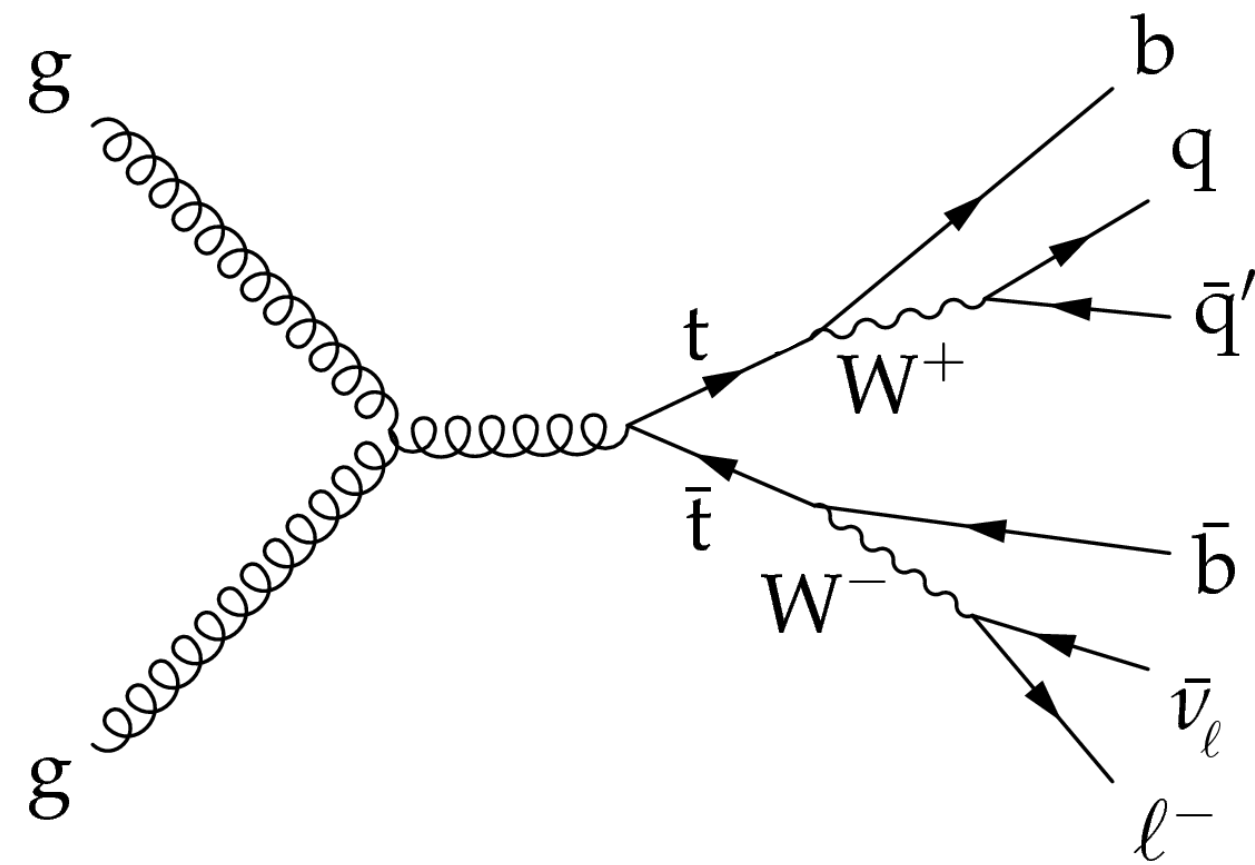
Parton-level unfolding — $t\bar{t}$ decay



Parton-level unfolding — $t\bar{t}$ decay



Parton-level unfolding — $t\bar{t}$ decay

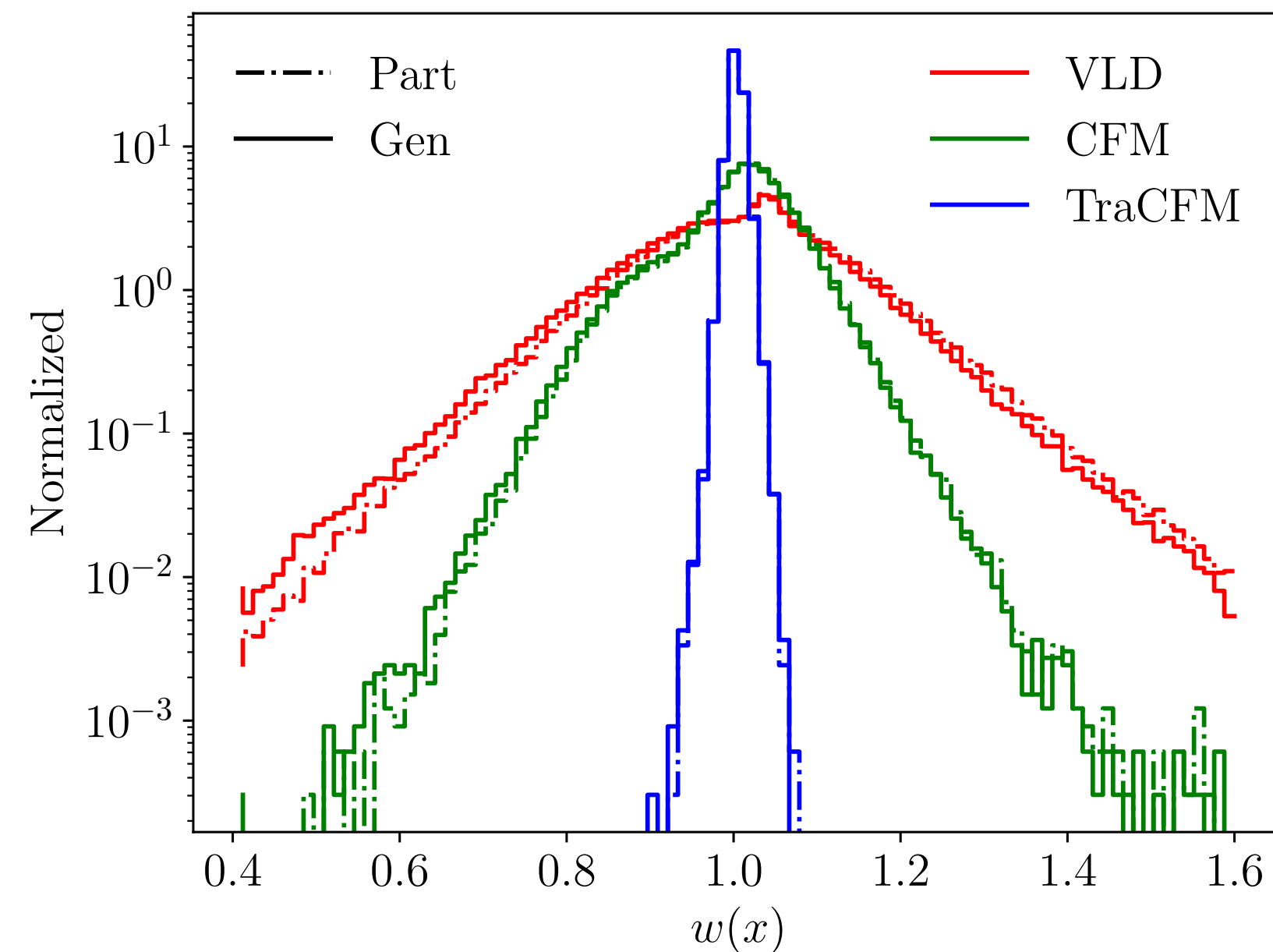


→ $\ell^+ \nu_\ell \{j\}$

Train a classifier classifier
between $p_{gen}(x)$ and $p_{unfold}(x)$

It learns the likelihood ratio

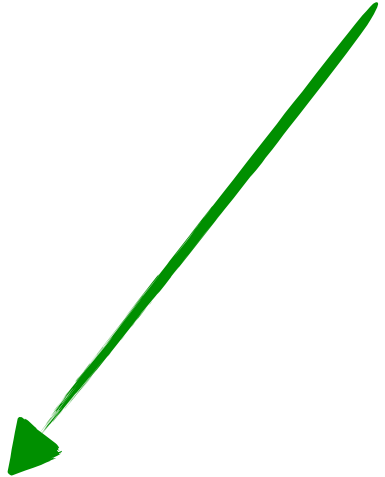
$$w(x) = \frac{p_{gen}(x)}{p_{unfold}(x)}$$



What about model dependence?

$$p(x_{gen} | x_{rec}) = \frac{p(x_{rec} | x_{gen})p(x_{gen})}{p(x_{rec})}$$

Prior



What about model dependence?

This problem is common to a long list of unfolding methods, with and without ML

Solution: Follow an iterative approach where we update our prior after each iteration

The same is done in Iterative Bayesian Unfolding, RooUnfold

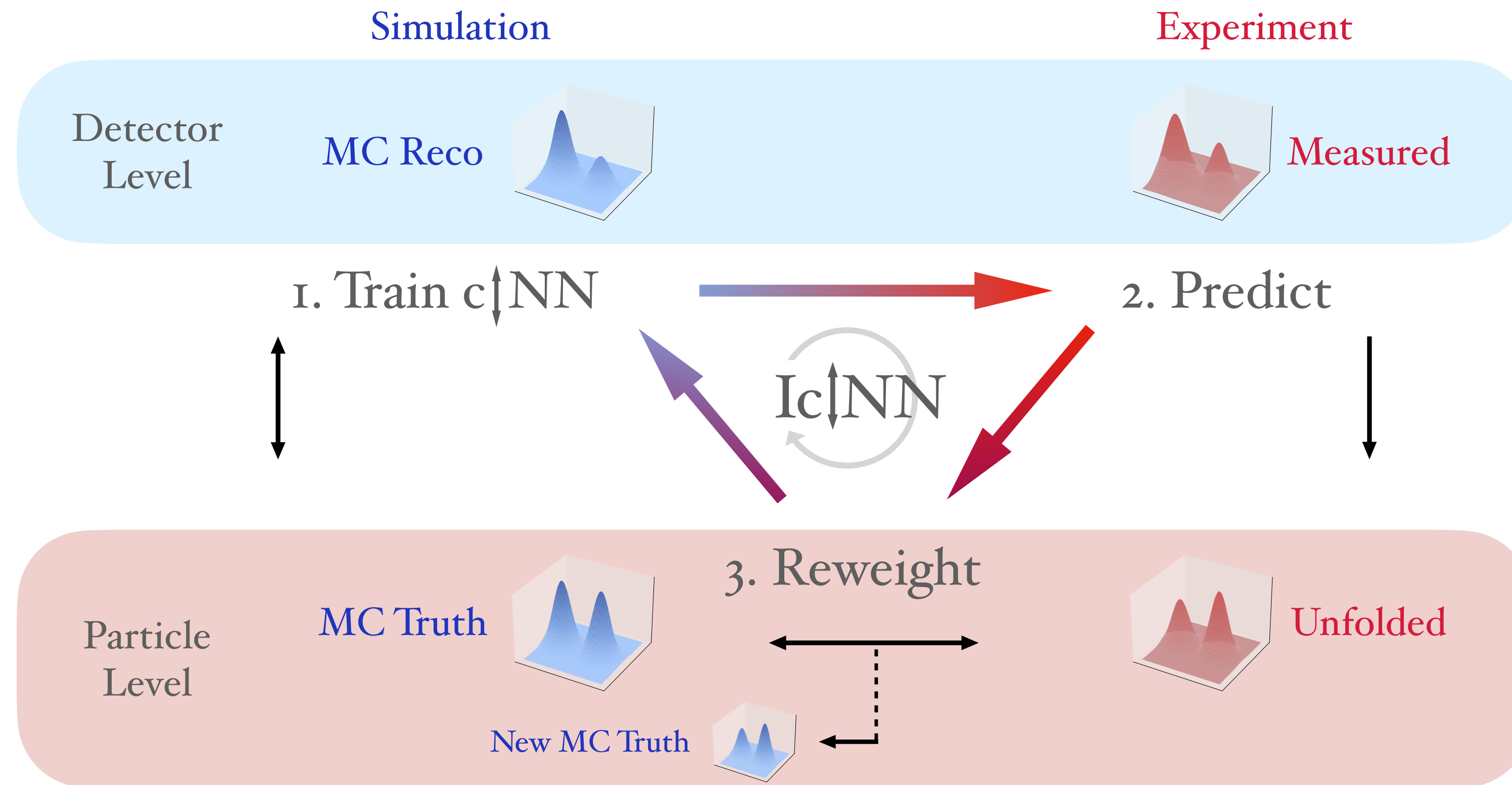
$$p(x_{gen} | x_{rec}) = \frac{p(x_{rec} | x_{gen})p(x_{gen})}{p(x_{rec})}$$

Prior

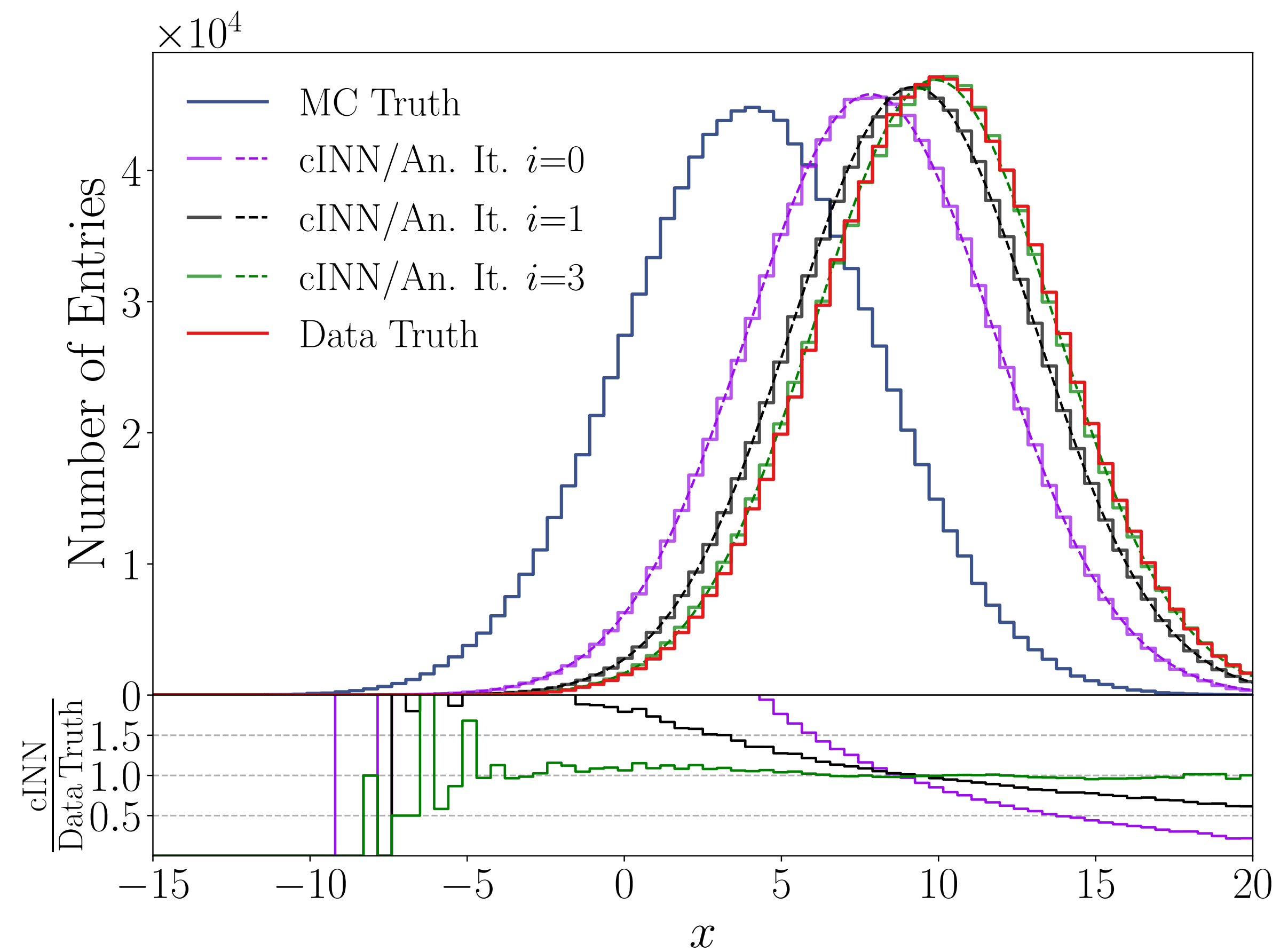
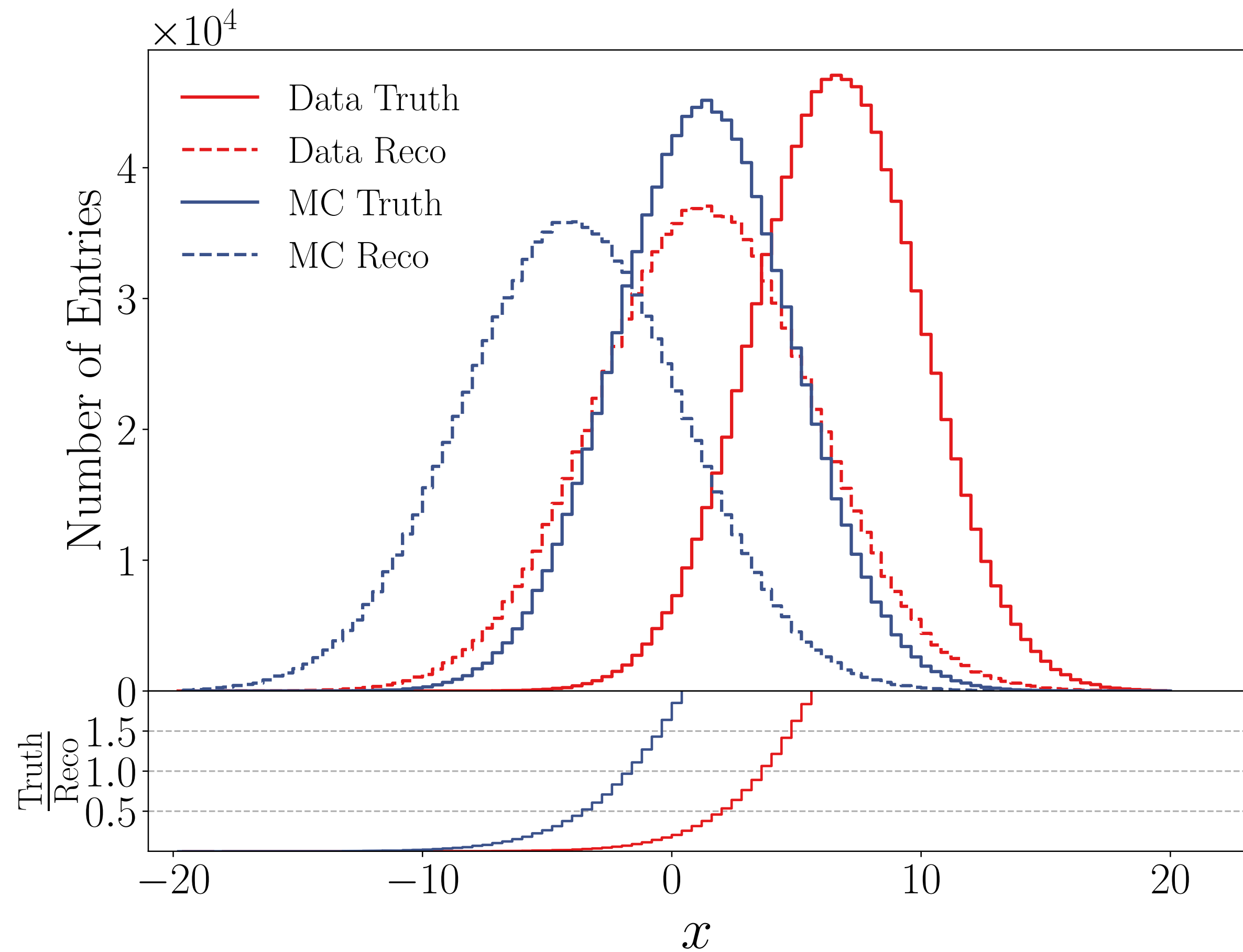
$$p_{unfold}(x_{gen}) = \int p_{data}(x_{rec})p(x_{gen} | x_{rec}) dx_{rec}$$

Use as new prior and start over

Iterative generative unfolding

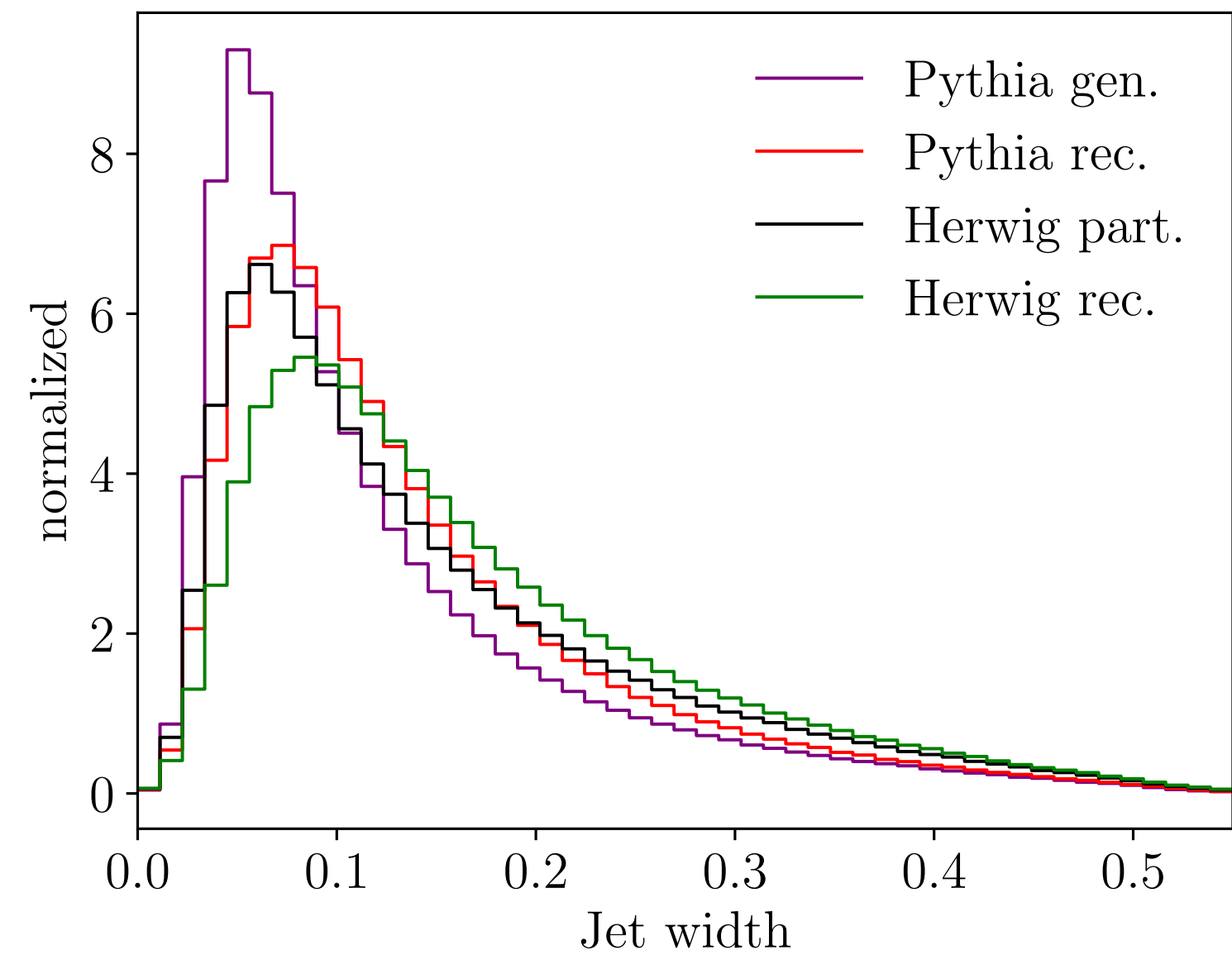


Iterative generative unfolding

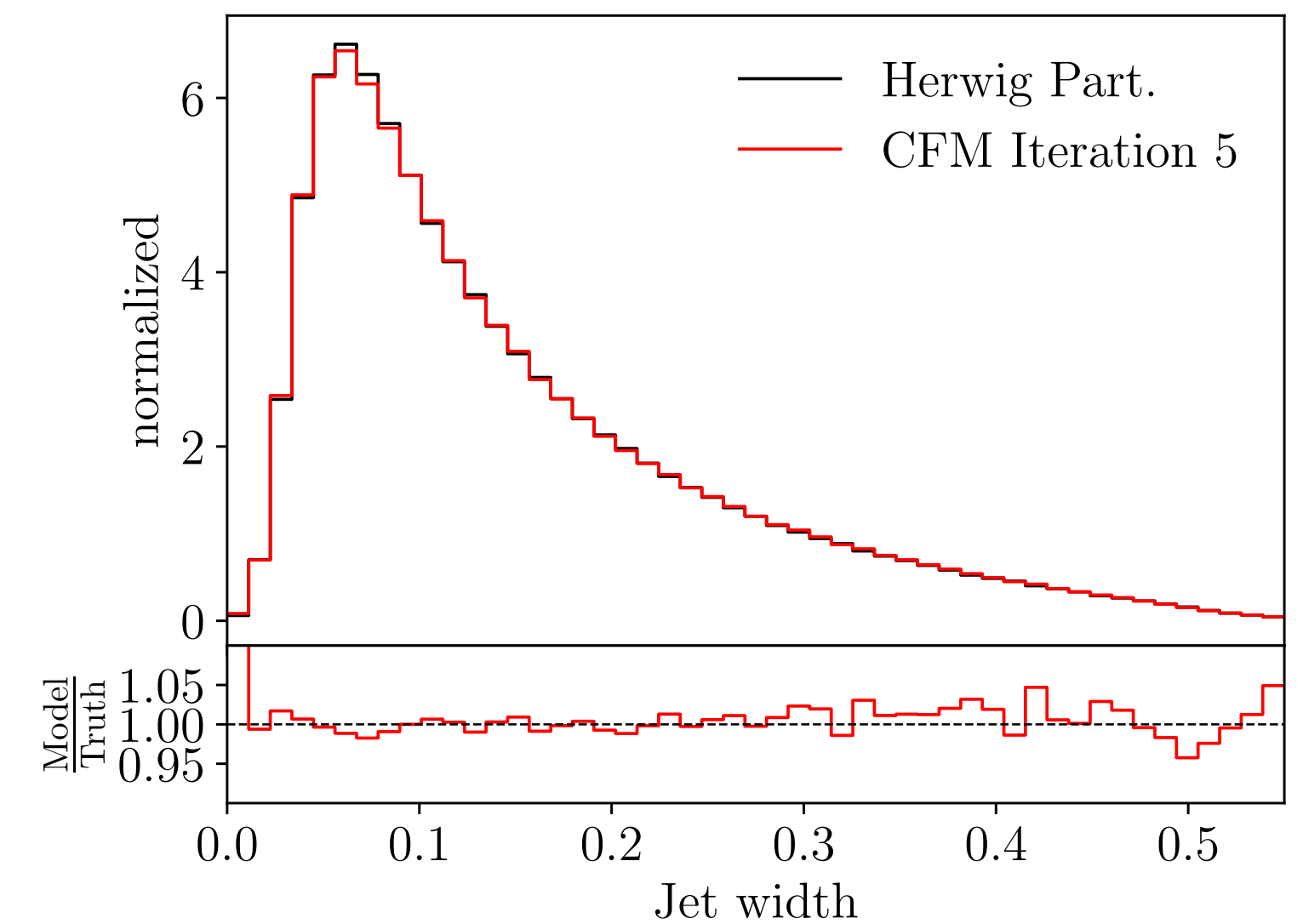
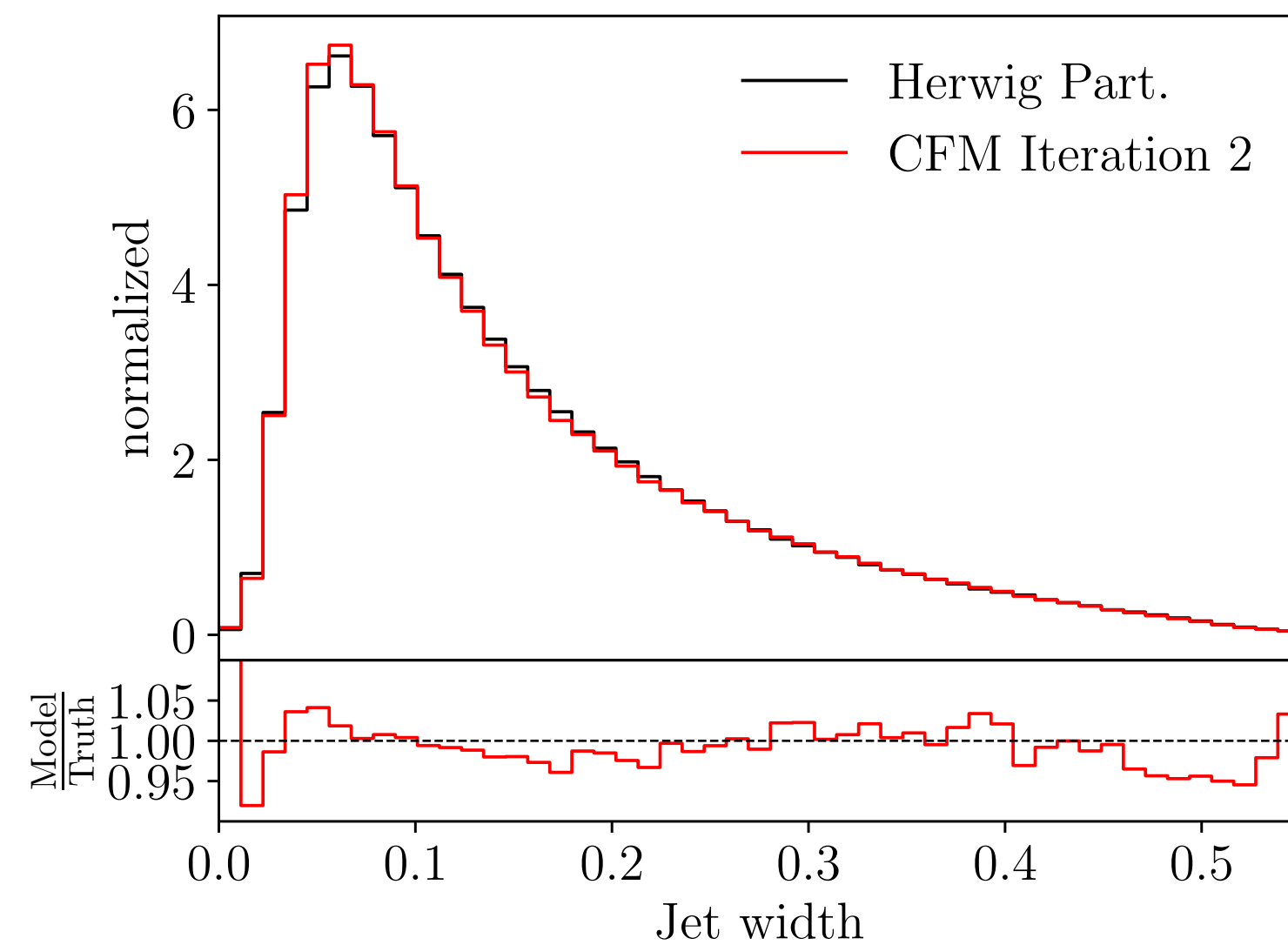
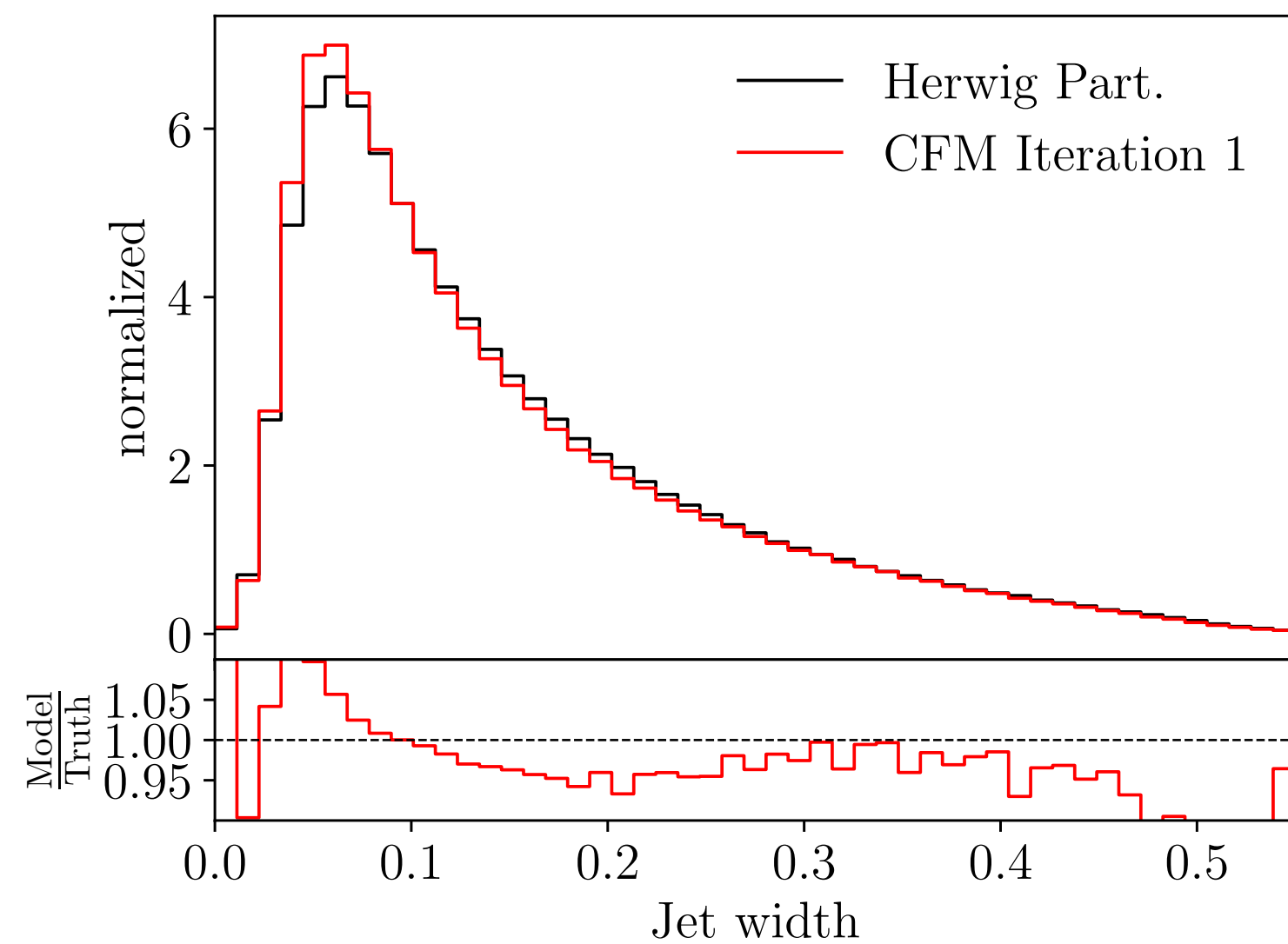


Z+jets: Pythia vs Herwig simulation

Use Pythia simulation as MC
Use Herwig simulation as Data



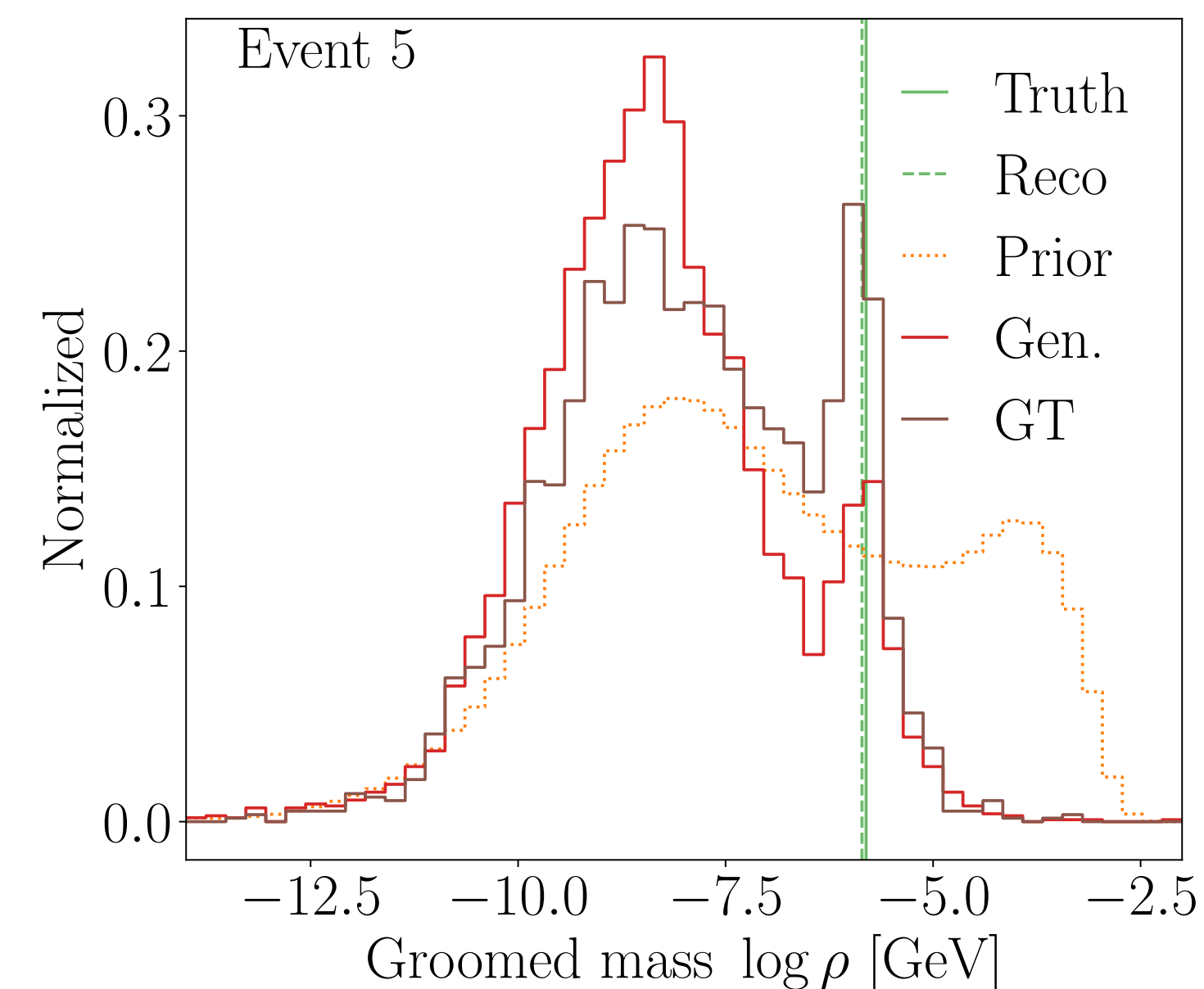
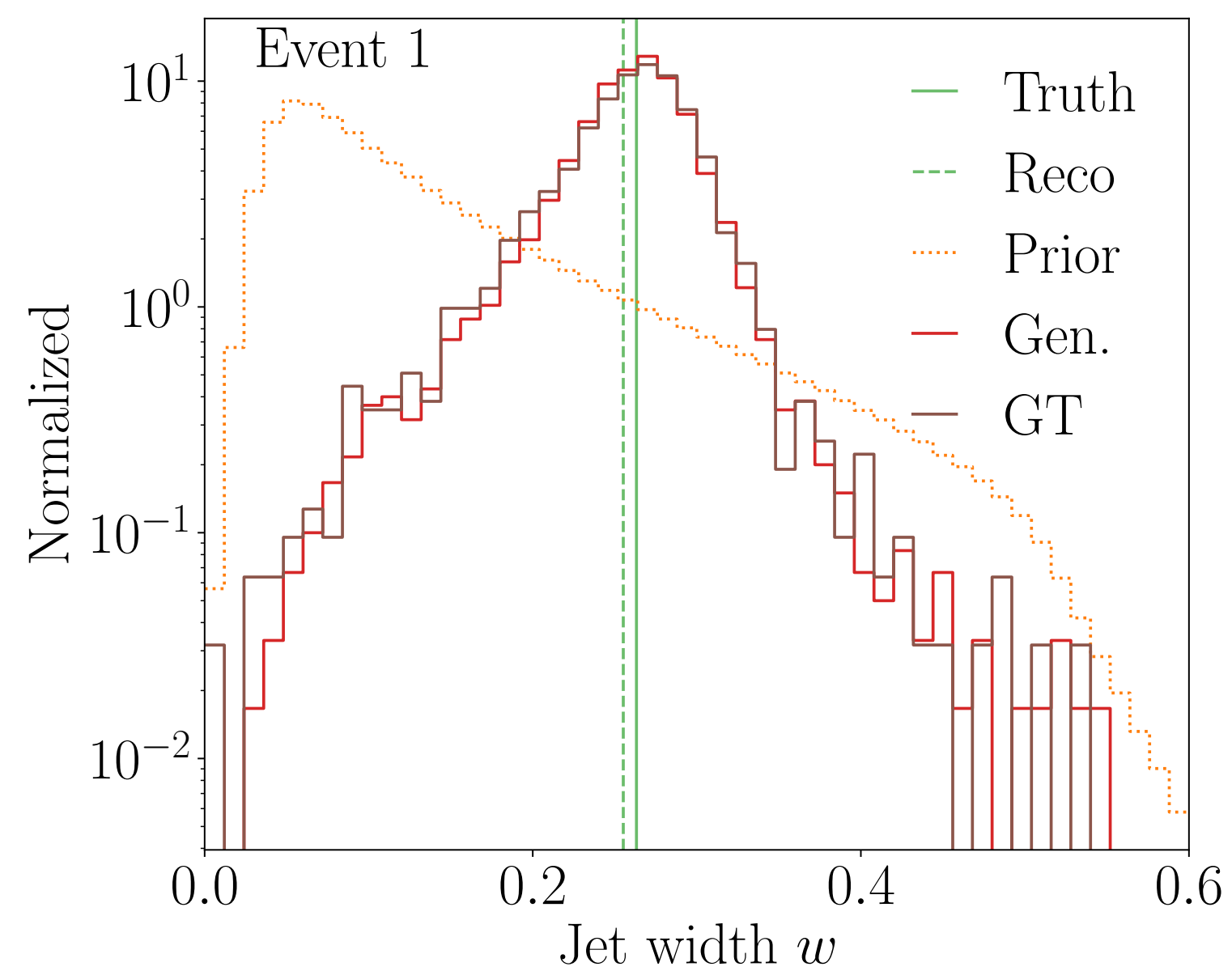
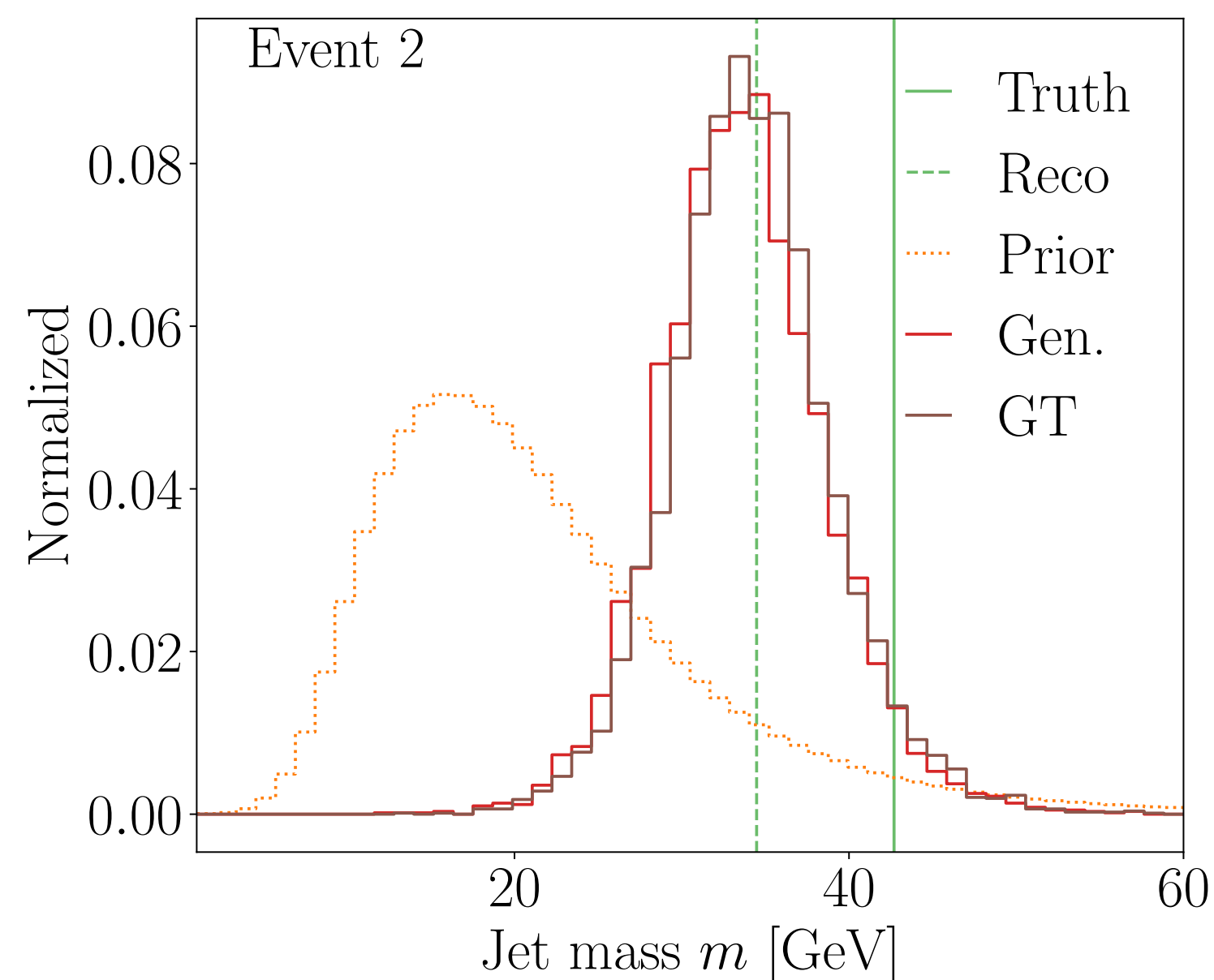
Following
Andreassen et al.
arXiv: 1911.09107



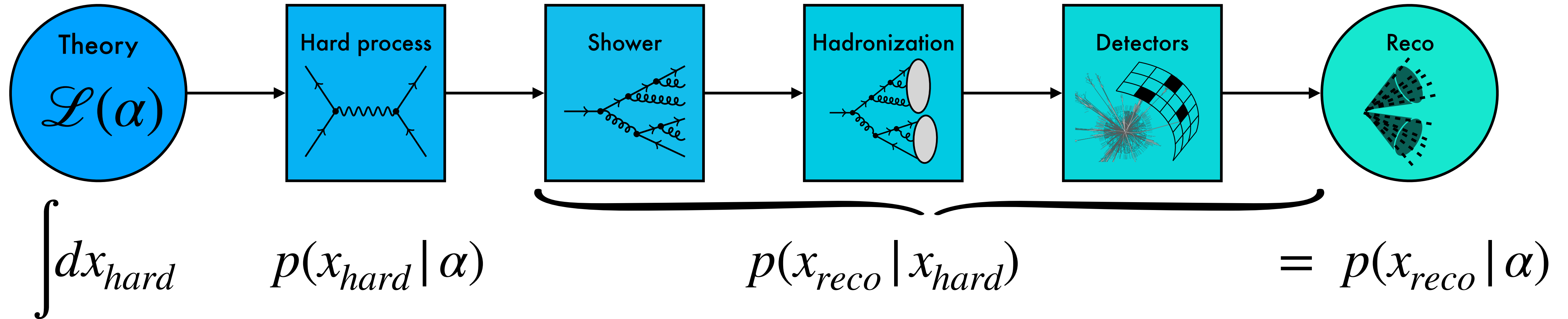
Single-event posteriors

Our network learned the posterior distribution $p(x_{gen} | x_{rec})$

We can sample from this for a single reco event to map out the posterior



Making use of single-event posteriors



We want to explicitly calculate reco-level likelihoods for our observed data given some theory hypothesis $\mathcal{L}(\alpha)$

Let us say we somehow now the transfer probability $p(x_{reco} | x_{hard})$ and can efficiently calculate it

How can we solve the difficult high-dimensional integral precisely and efficiently?

Making use of single-event posteriors

Challenging

$$\int dx_{hard} p(x_{hard} | \alpha) p(x_{reco} | x_{hard}) = p(x_{reco} | \alpha)$$

$$= \left\langle \frac{1}{q(x_{hard})} p(x_{hard} | \alpha) p(x_{reco} | x_{hard}) \right\rangle_{x_{hard} \sim q}$$

Integral becomes trivial if : $q(x_{hard} | x_{reco}, \alpha) = \underbrace{p(x_{hard} | x_{reco}, \alpha)}$

This is exactly our unfolding network !

For details see Butter et al. 2210.00019, Heimgel et al. 2310.07752

Conclusion

SciPost Physics

2404.18807

Submission

Generative Unfolding works !

Conditional Networks and Distribution Mapping both achieve percent-level precision

Conditional Networks learn the true migration encoded in the detector simulation.

Distribution Mapping Networks learn a simple optimal-transport based mapping

Conditional Generative Unfolding enables probabilistic inversion of simulation chain

Prior dependence is a solvable problem

Conditional Generative Models map out the single-event posterior distributions

The Landscape of Unfolding with Machine Learning

Nathan Huetsch¹, Javier Mariño Villadamigo¹, Alexander Shmakov², Sascha Diefenbacher³,
Vinicius Mikuni³, Theo Heimel¹, Michael Fenton², Kevin Greif²,
Benjamin Nachman^{3,4}, Daniel Whiteson², Anja Butter^{1,5}, and Tilman Plehn^{1,6}



Further investigation of failure modes

Further investigation of uncertainties

Further investigation of background, cuts, efficiencies

Flow Matching (Lipman et al. 2210.02747)

Training

1. Sample paired data from our simulation

$$(x_0, c) = (x_{gen}, x_{rec}) \sim p(x_{gen}, x_{rec})$$

2. Sample noise and a timestep

$$x_1 = \epsilon \sim \mathcal{N}(0,1), t \sim \mathcal{U}([0,1])$$

3. Calculate the trajectory

$$x_t = (1 - t)x_0 + tx_1$$
$$v_t = \frac{dx_t}{dt} = -x_0 + x_1$$

4. Predict the velocity field

$$\mathcal{L} = \left| v_\theta(x_t, t, c) - v_t \right|^2$$

Generation

1. Sample a reco event from our measured data

$$c = x_{rec} \sim p(x_{rec})$$

2. Sample noise as initial condition

$$x_1 = \epsilon \sim \mathcal{N}(0,1)$$

3. Solve the ODE numerically

$$x_0 = x_{gen} = x_1 + \int_1^0 v_\theta(x_t, t, c) dt$$