



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

IMPRS-Seminar 21.09.2023

Jet-Diffusion vs Jet-GPT

Modern Networks for the LHC

Anja Butter, Nathan Huetsch, Sofia Palacios Schweitzer, Tilman Plehn, Peter Sorrenson, Jonas Spinner
arXiv: 2305.10475



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

IMPRS-Seminar 21.09.2023

Jet-Diffusion vs Jet-GPT

Modern Networks for the LHC

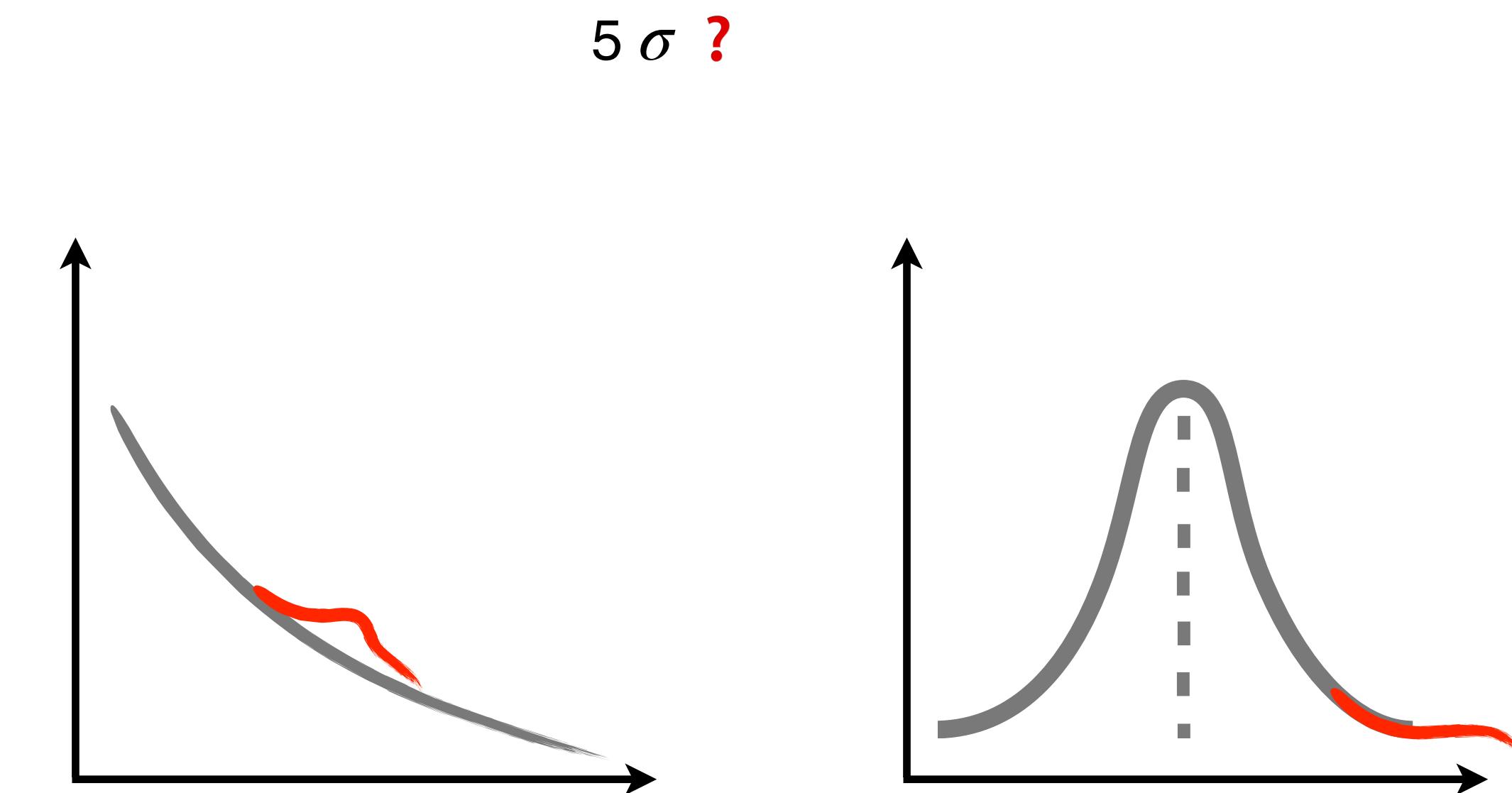
Anja Butter, Nathan Huetsch, Sofia Palacios Schweitzer, Tilman Plehn, Peter Sorrenson, Jonas Spinner
arXiv: 2305.10475

Why Event Generation?

Vast amount of data collected by collider experiments

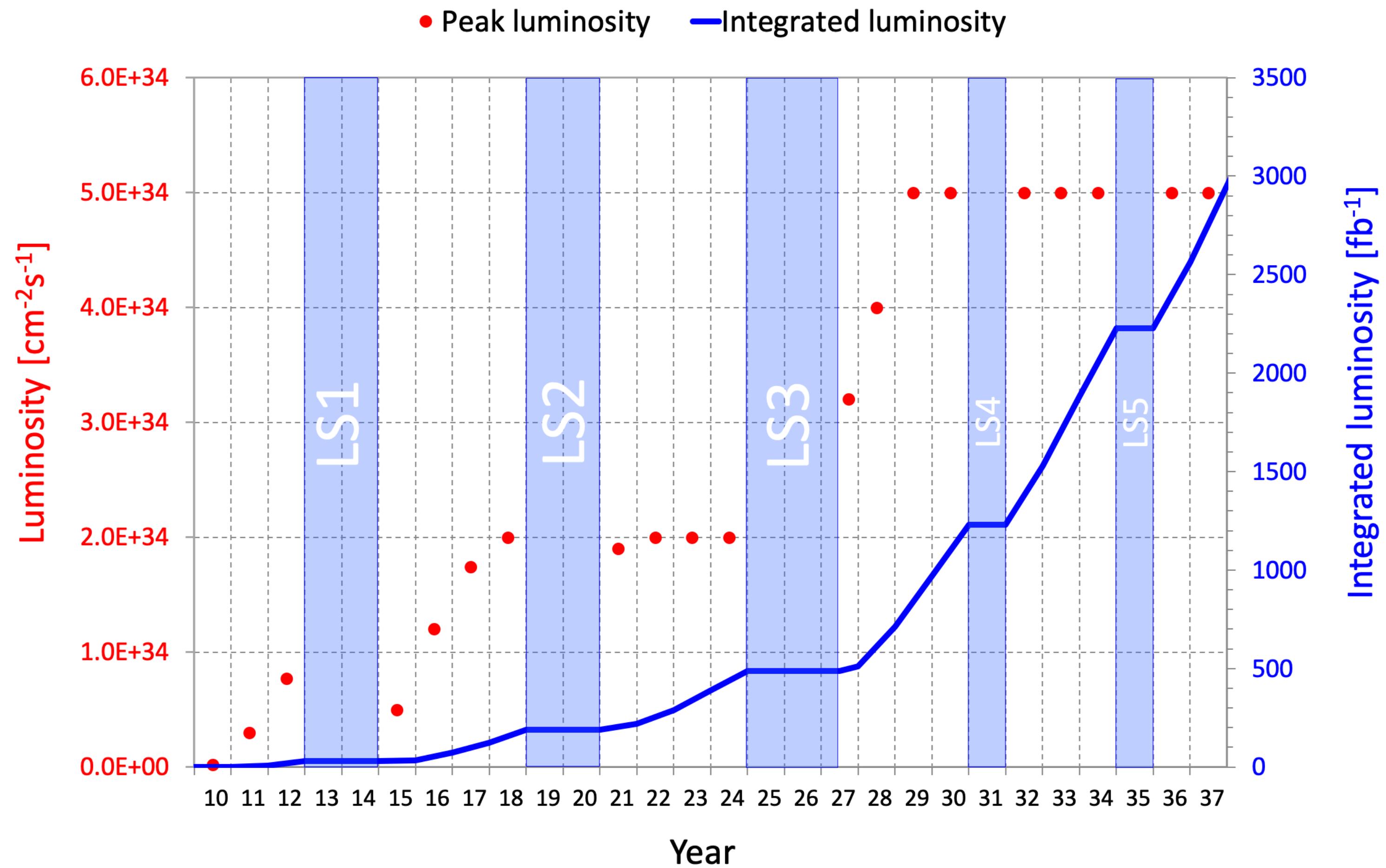
Standard Model is probed

Theoretical predictions (simulation) needs to match experimental statistics



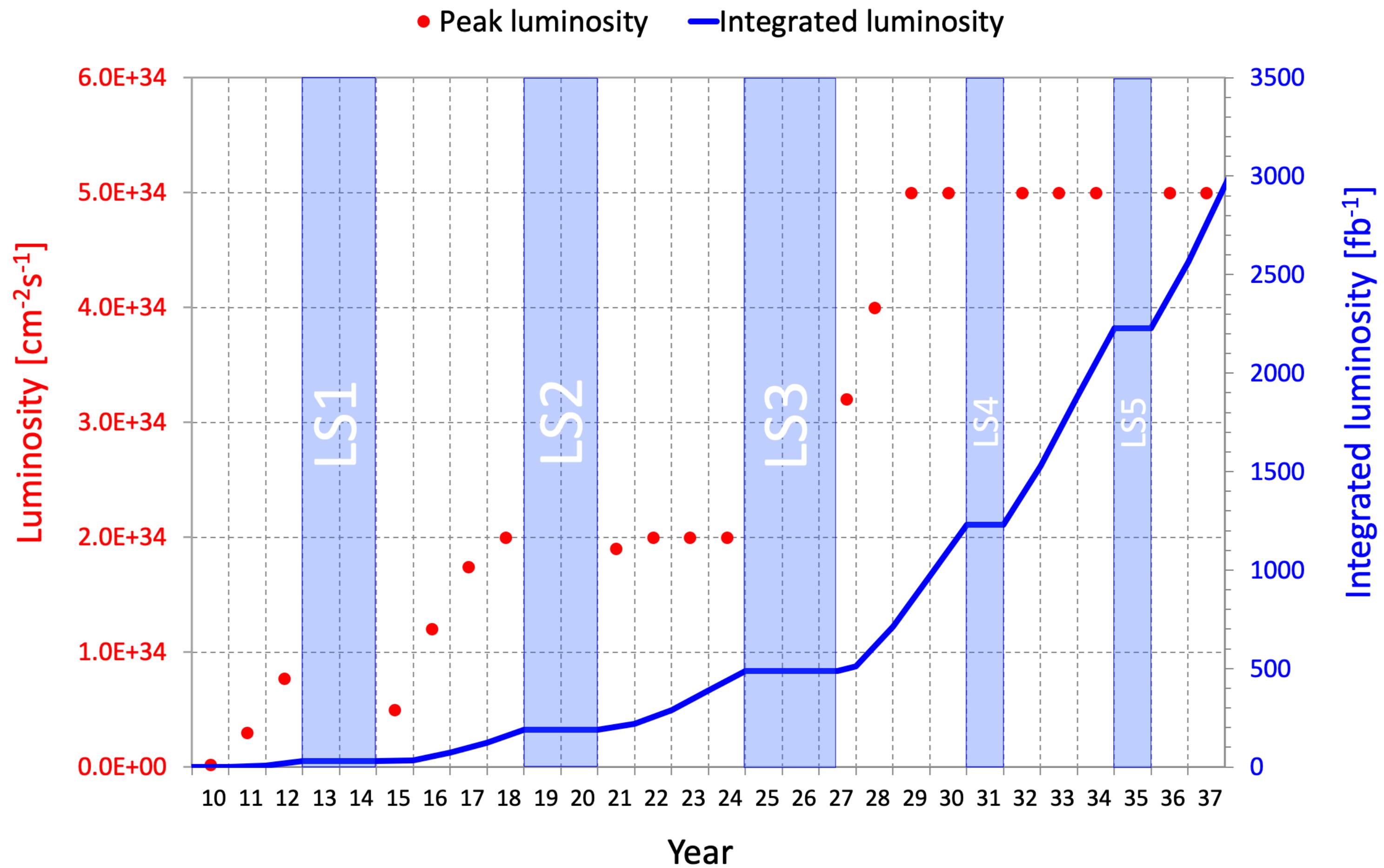
Why ML Event Generation?

After high luminosity runs $\rightarrow \sim 5$ times as much data
Theoretical predictions needs to be even more precise (include higher correction terms)

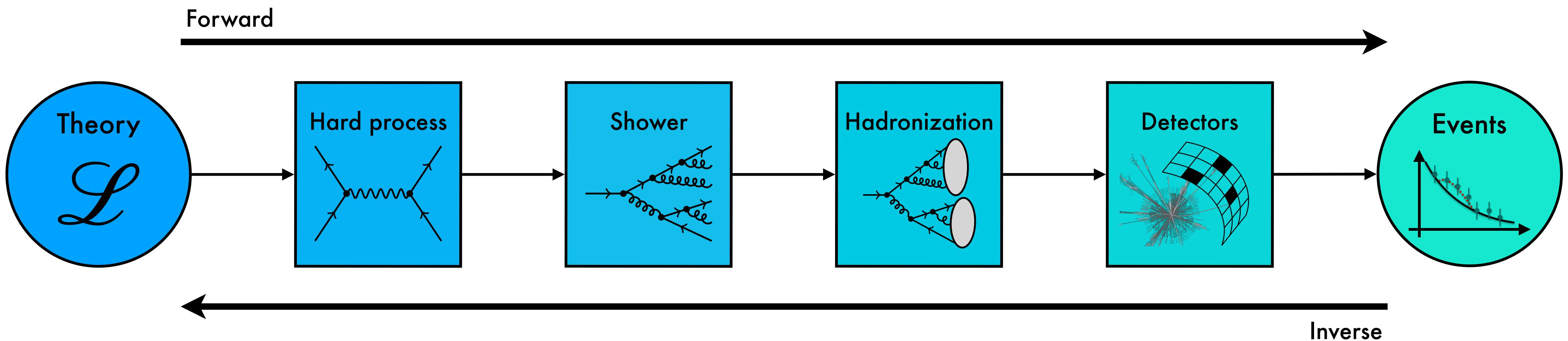


Why ML Event Generation?

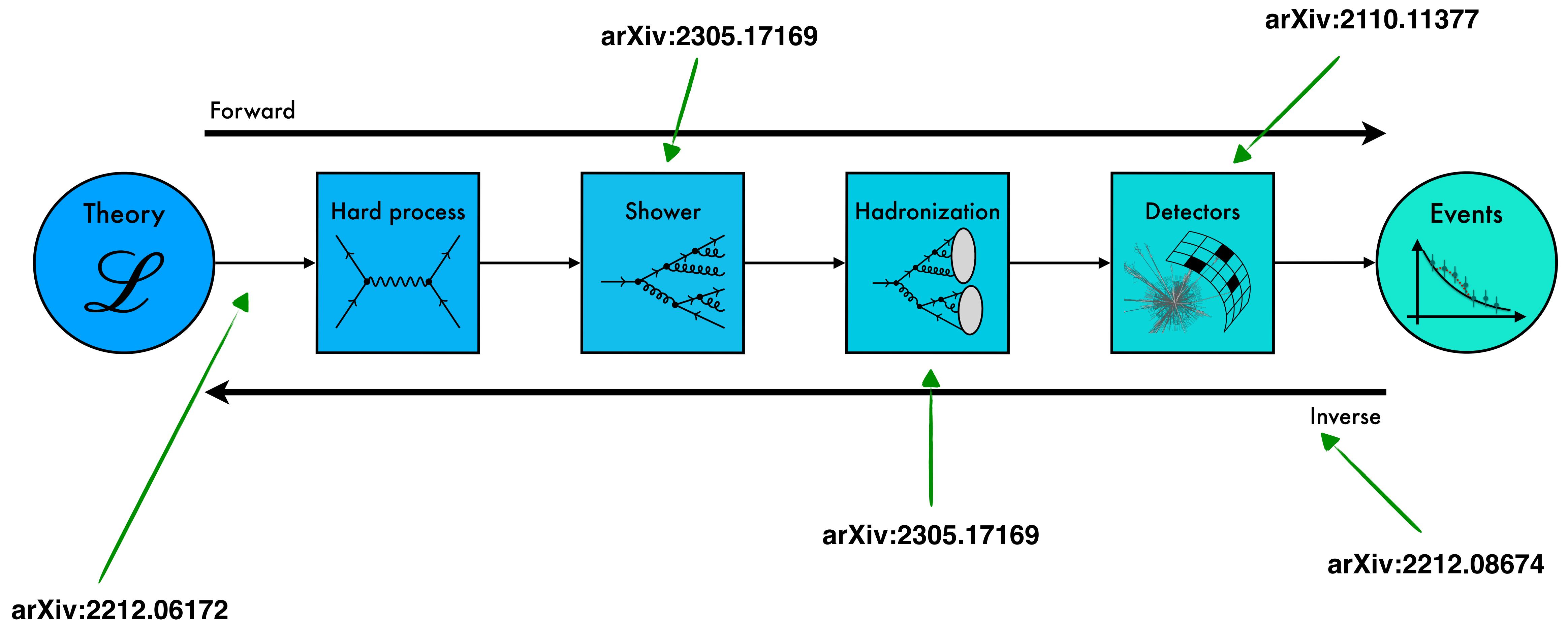
- After high luminosity runs → ~ 5 times as much data
- Theoretical predictions needs to be even more precise (include higher correction terms)
- But: Currently computational expensive



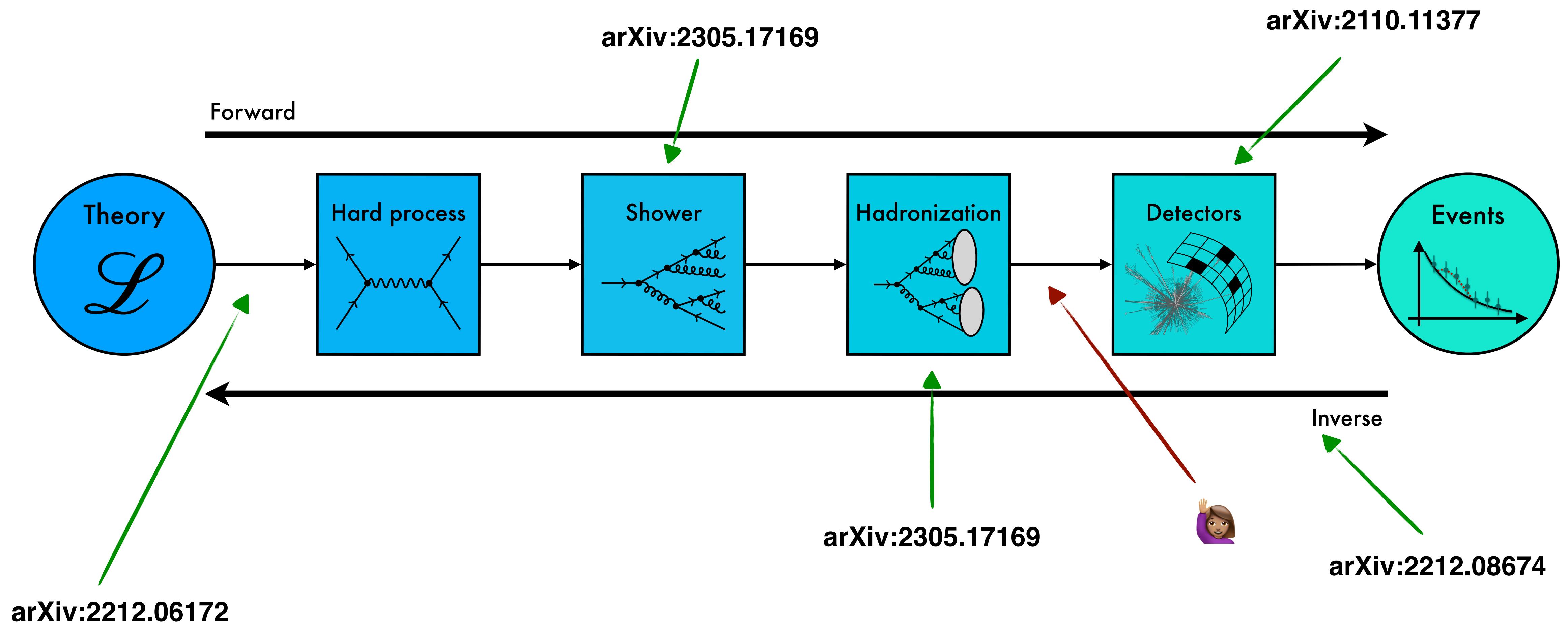
Where ML Event Generation?



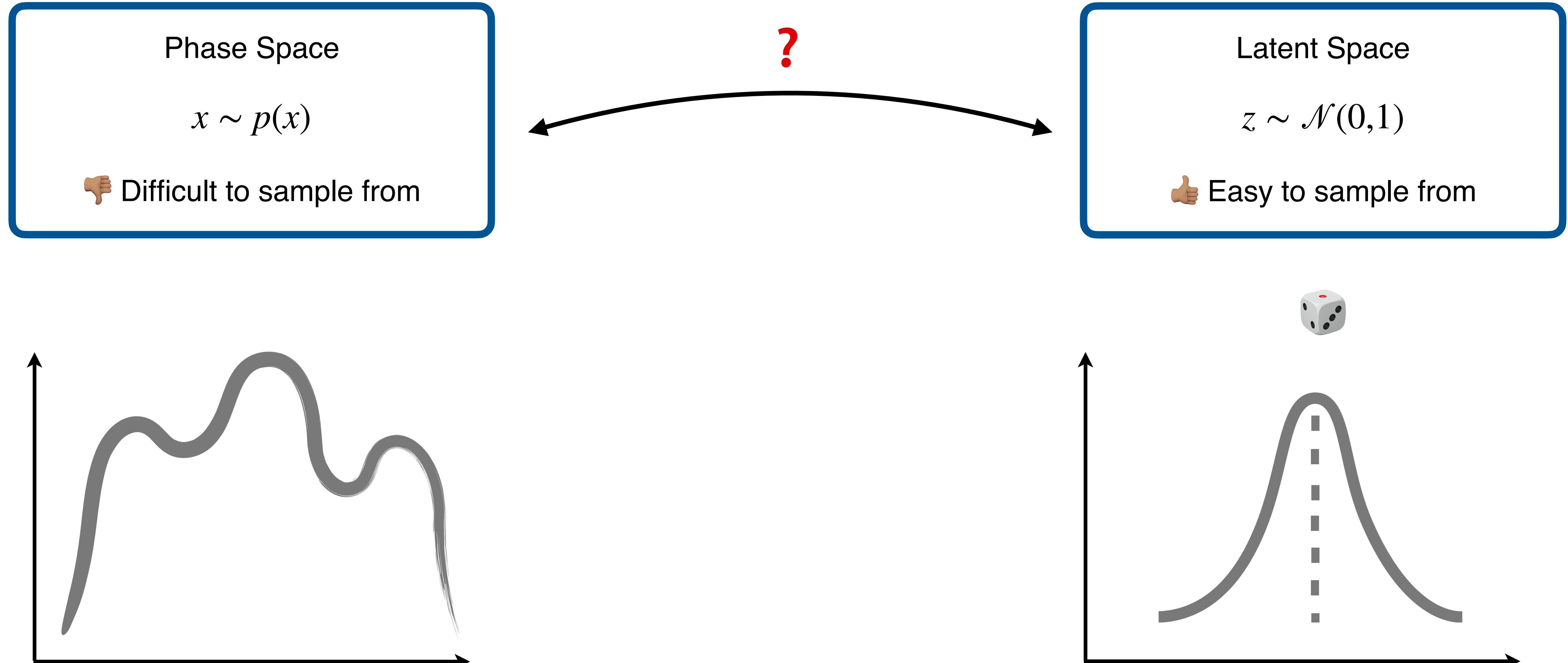
Where ML Event Generation?



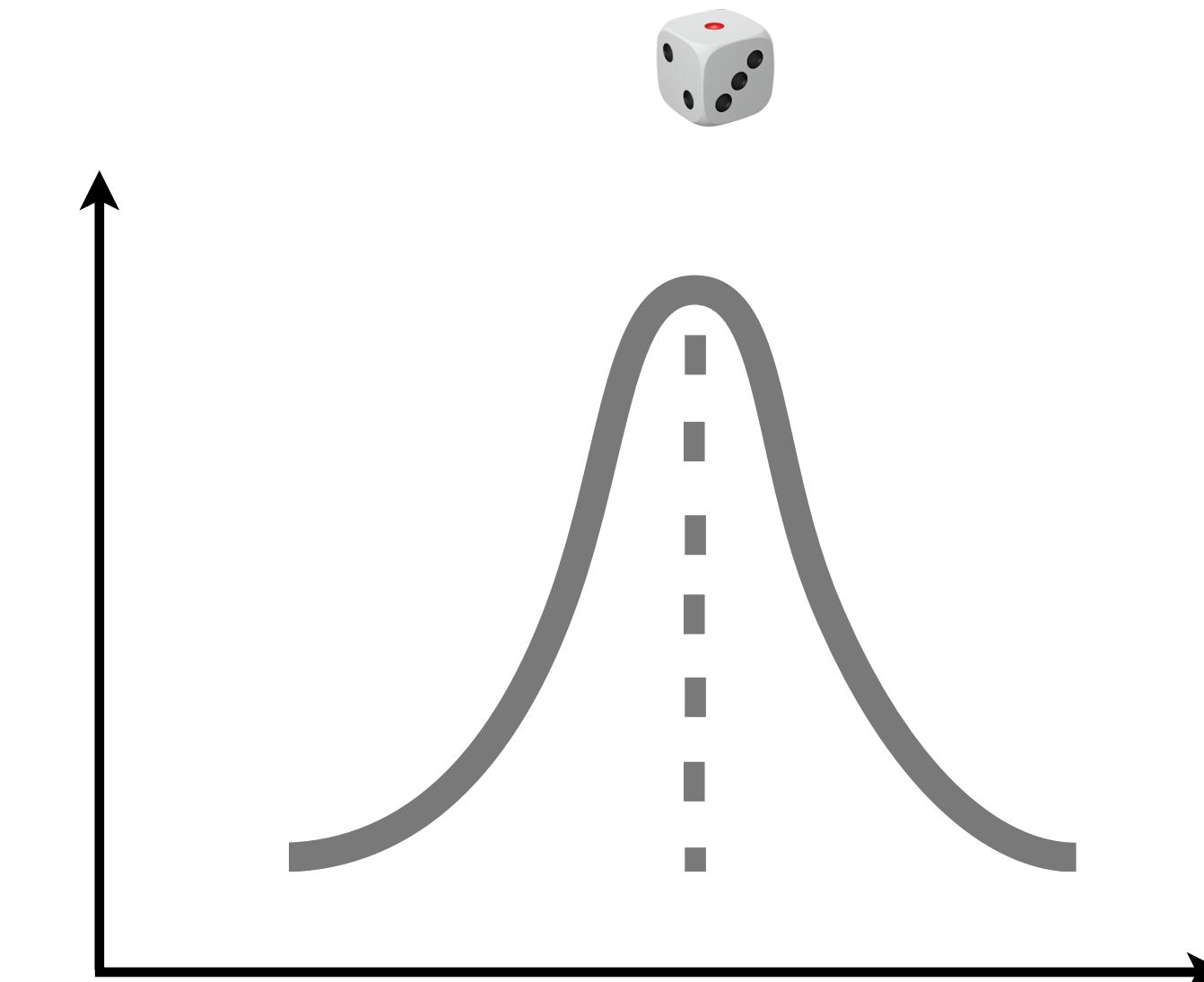
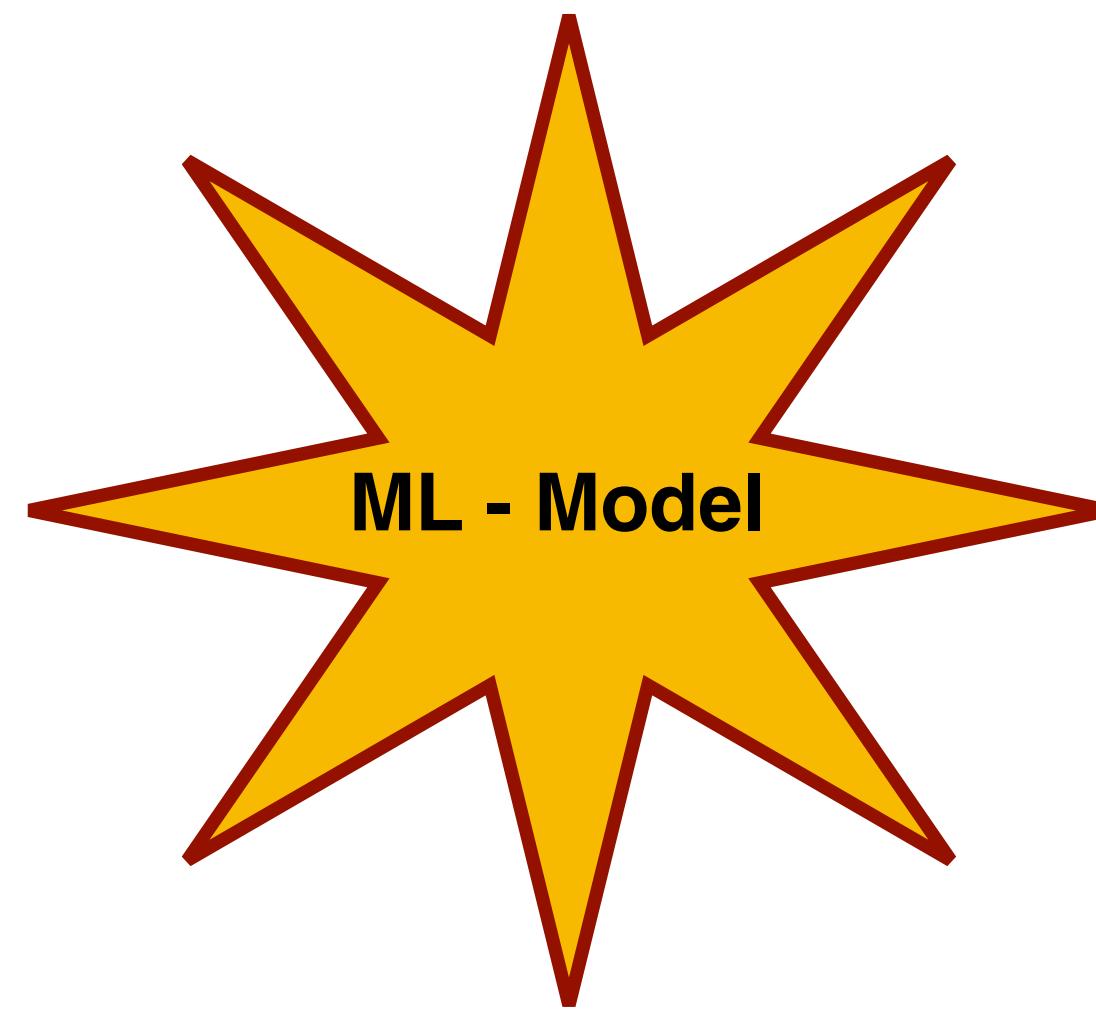
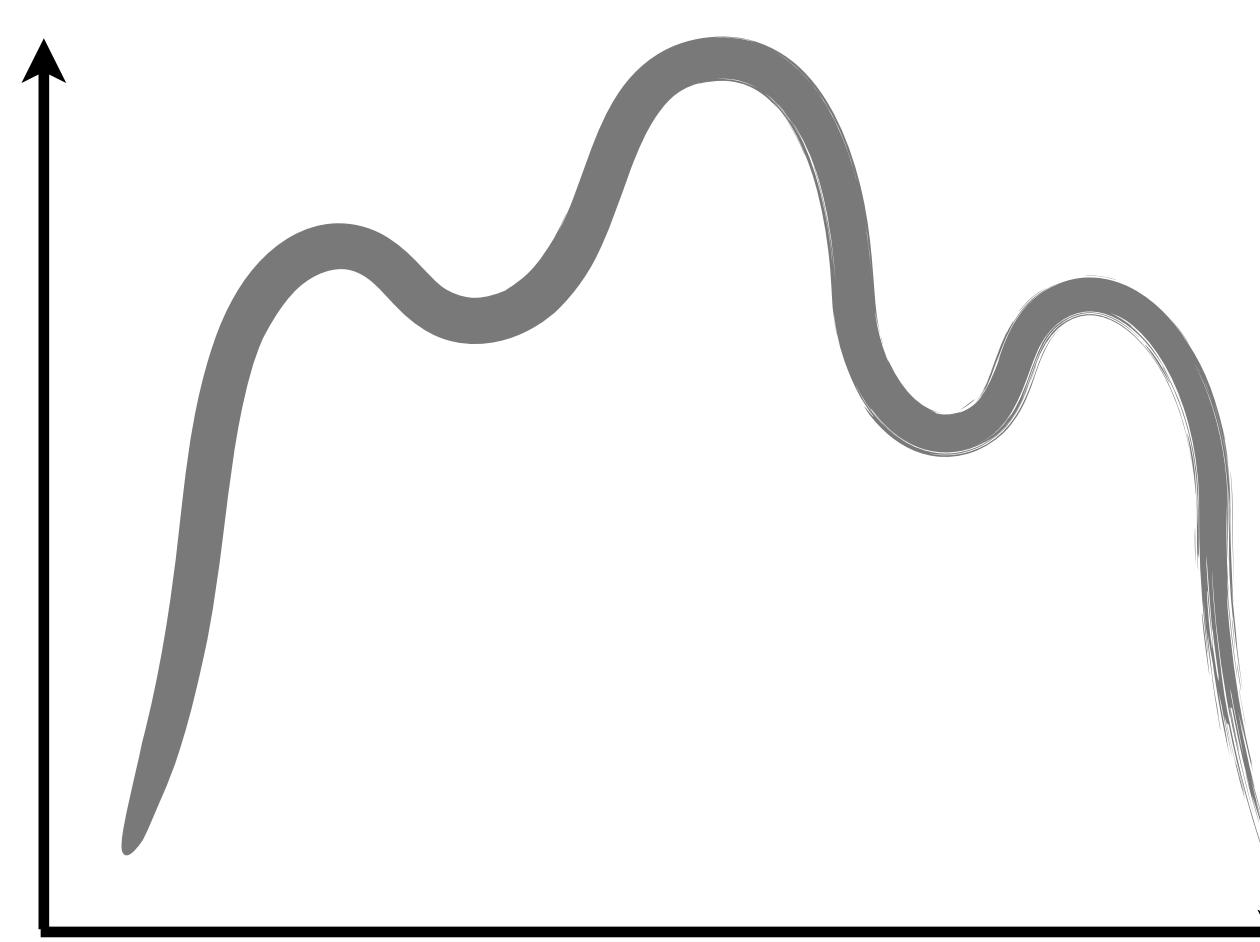
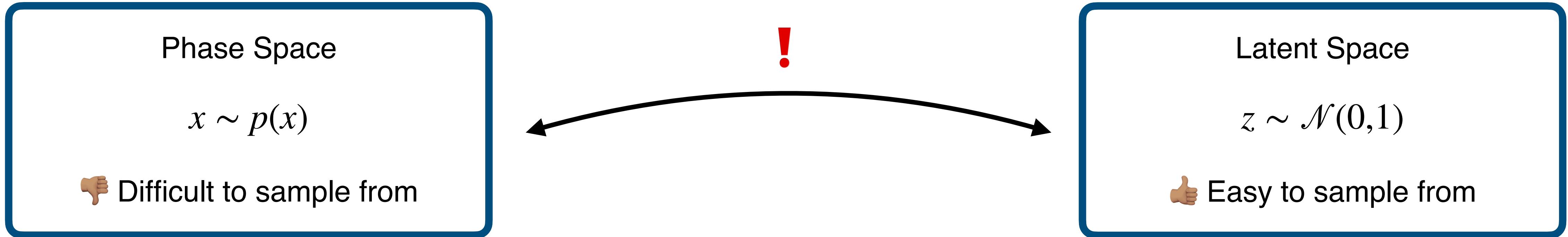
Where ML Event Generation?



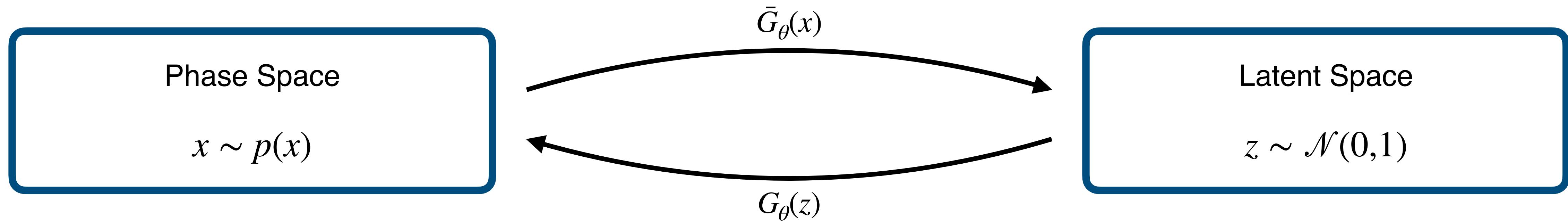
How to be generative



How to be generative



Invertible Neural Networks (INNs)

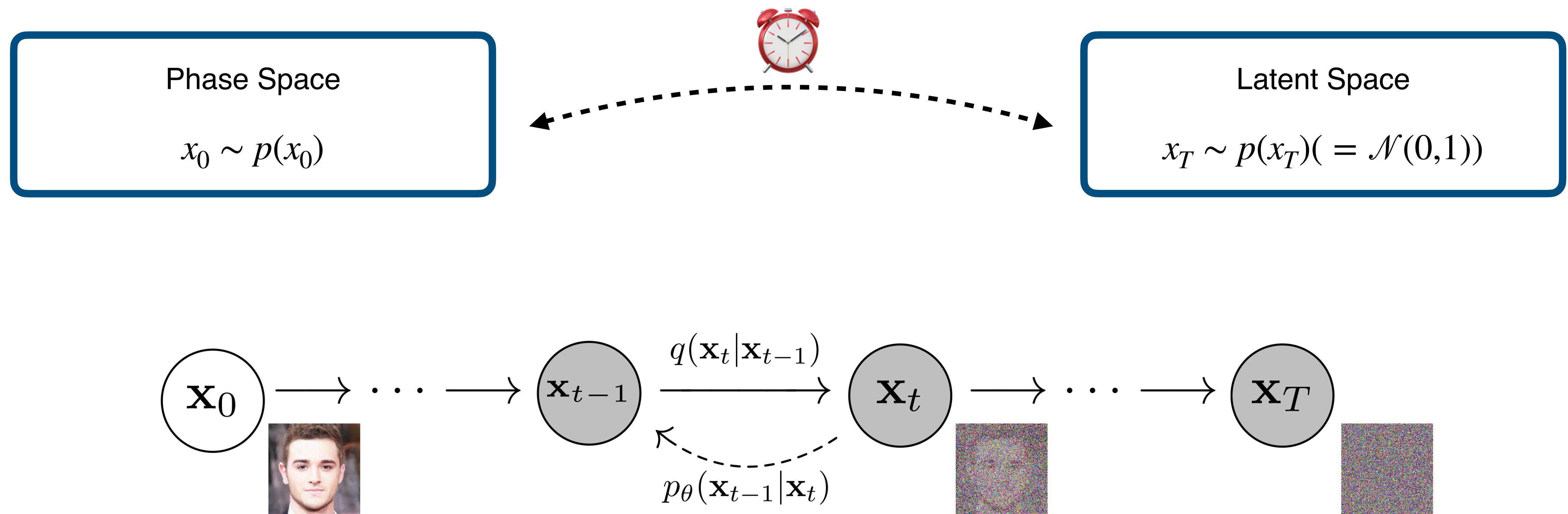


Bijective mapping $G_\theta(x)$

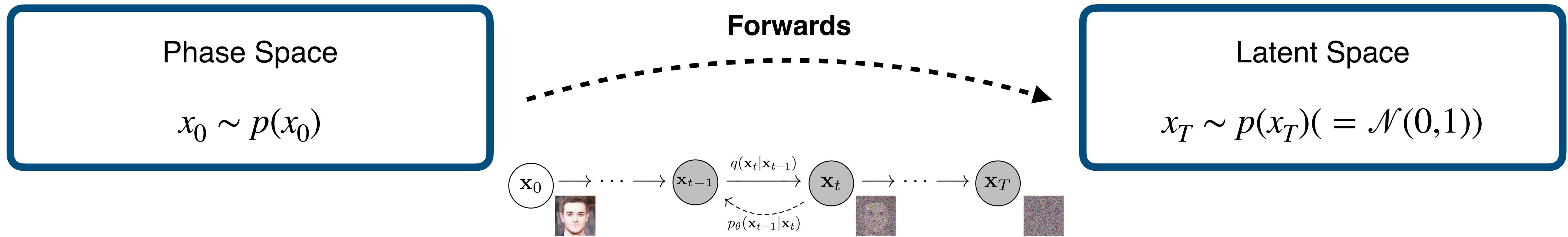
$$p_\theta(x) = p(z) \frac{dz}{dx} = p(z) \left| \frac{\partial \bar{G}_\theta(x)}{\partial x} \right|$$

$$\mathcal{L}_{INN} = -\log p_\theta(x) = -\log p(\bar{G}_\theta(x)) - \log \left| \frac{\partial \bar{G}_\theta(x)}{\partial x} \right|$$

Diffusion Models



Diffusion Models (DDPM)

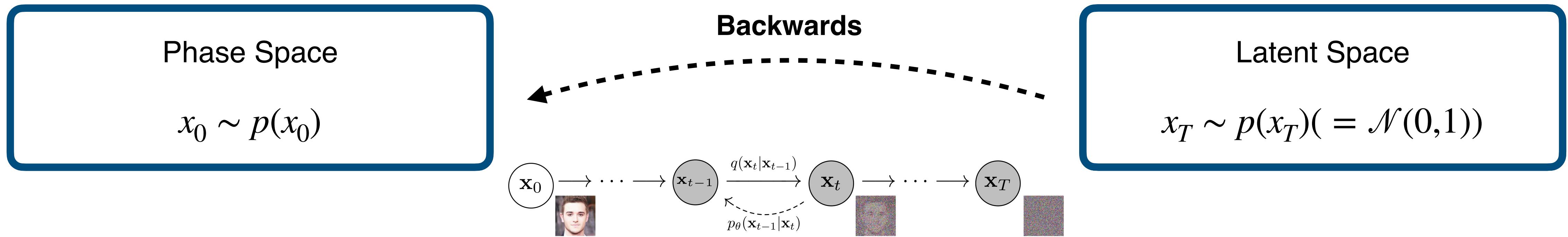


$$q(x_1, \dots, x_T | x_0) = \prod_{t=1}^T q(x_t | x_{t-1})$$

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t)$$

where β_t follows noise scheduler

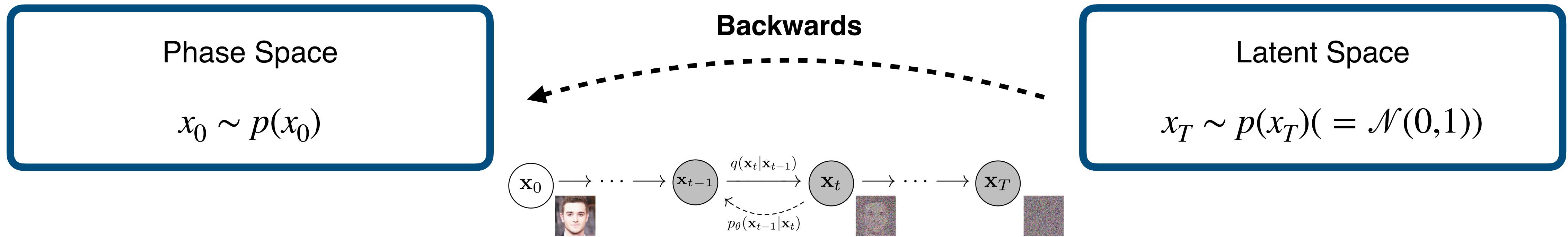
Diffusion Models (DDPM)



To reverse: $q(x_{t-1} | x_t) = \frac{q(x_t | x_{t-1}) q(x_{t-1})}{q(x_t)}$

...but we don't know $q(x_t)$ & $q(x_{t-1})$ 🙃 🙃

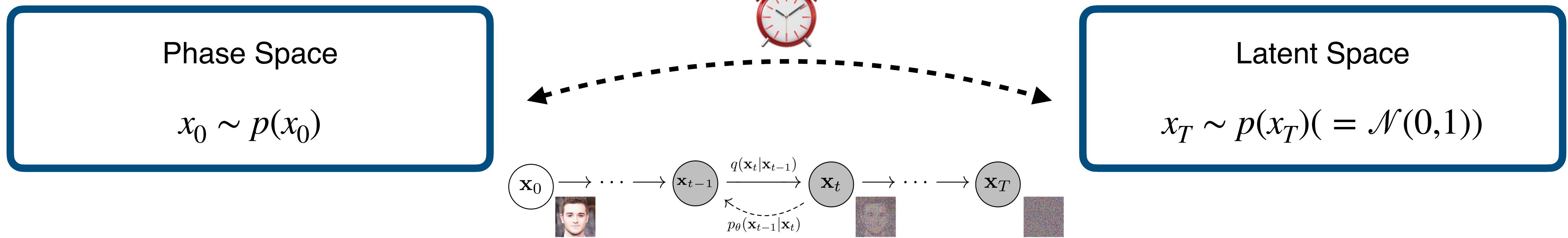
Diffusion Models (DDPM)



Instead, learn: $p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_\theta^2(x_t, t))$

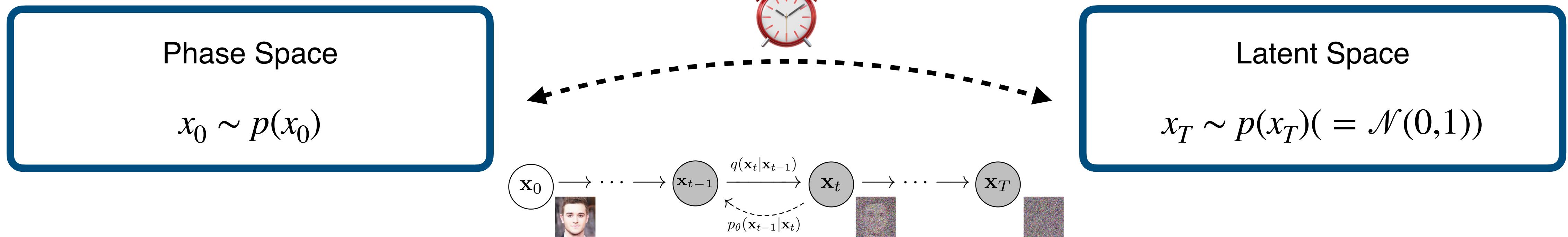
$$p(x_0, \dots, x_T | \theta) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t)$$

Diffusion Models (DDPM)



$$\mathcal{L}_{DDPM} = -\log p_\theta(x_0) \approx \frac{1}{2\sigma_t^2} \frac{\beta_t^2}{(1 - \beta_t)\bar{\beta}_t} |\epsilon(t) - \epsilon_\theta(t)|^2$$

Diffusion Models (DDPM)

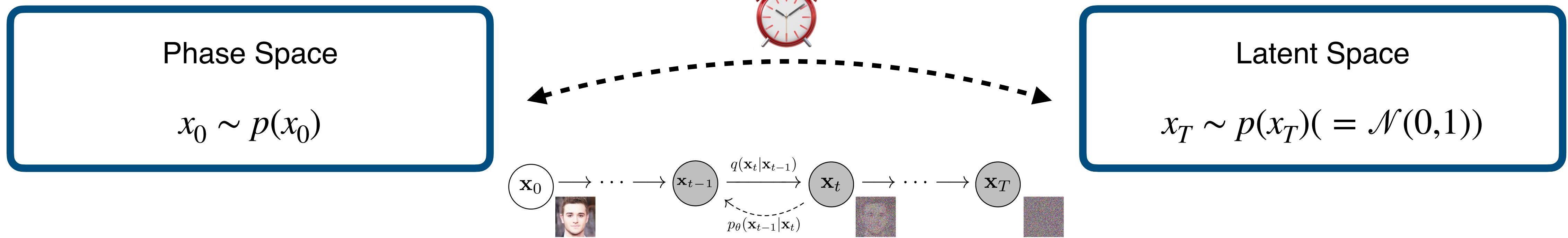


$$\mathcal{L}_{DDPM} = -\log p_\theta(x_0) \approx \frac{1}{2\sigma_t^2} \frac{\beta_t^2}{(1 - \beta_t)\bar{\beta}_t} |\epsilon(t) - \epsilon_\theta(t)|^2$$

- Reparametrization
• Estimation
• ...

Predicted and actual
noise added at time t

Diffusion Models (DDPM)

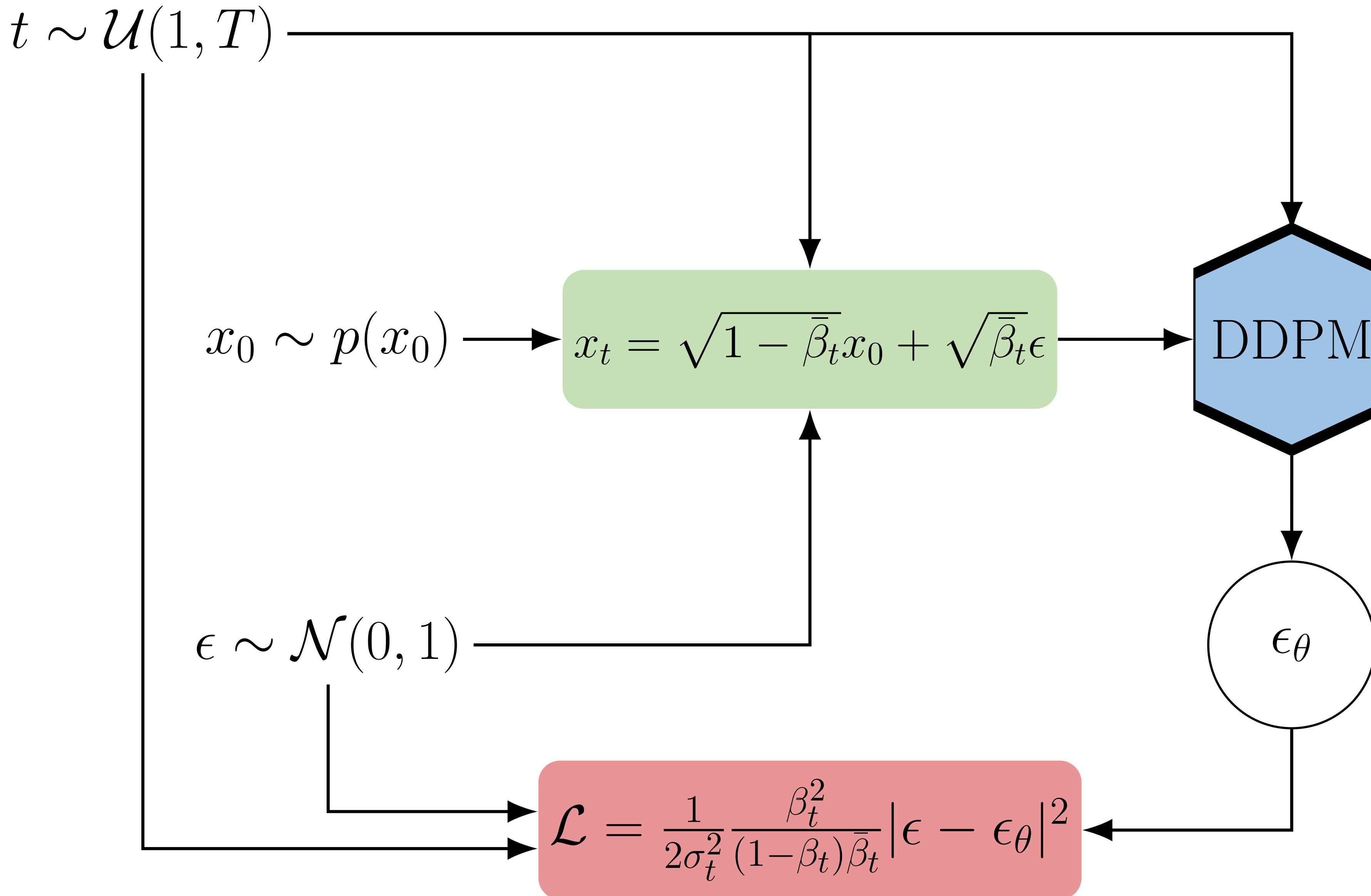


$$\mathcal{L}_{DDPM} = -\log p_\theta(x_0) \approx \frac{1}{2\sigma_t^2} \frac{\beta_t^2}{(1 - \beta_t)\bar{\beta}_t} |\epsilon(t) - \epsilon_\theta(t)|^2$$

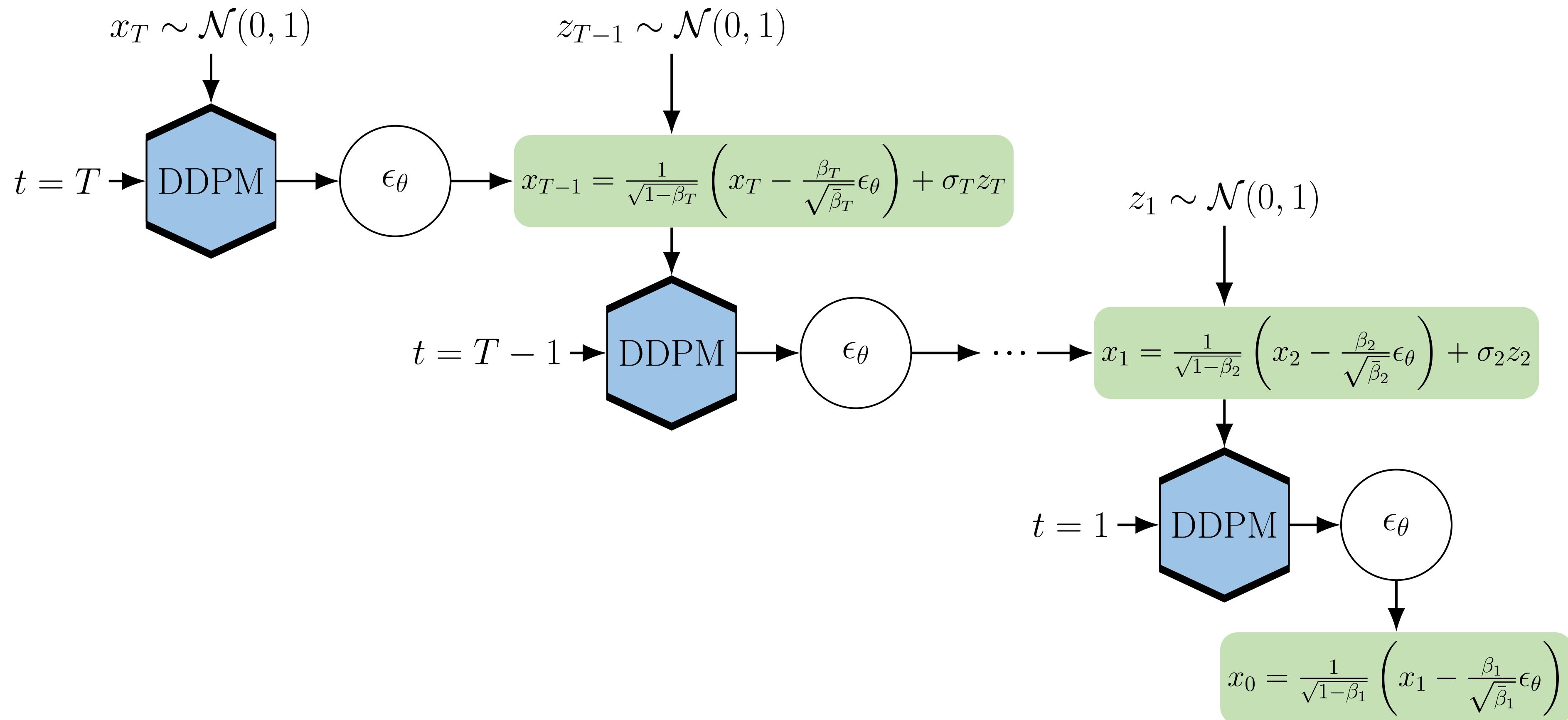


Denoising

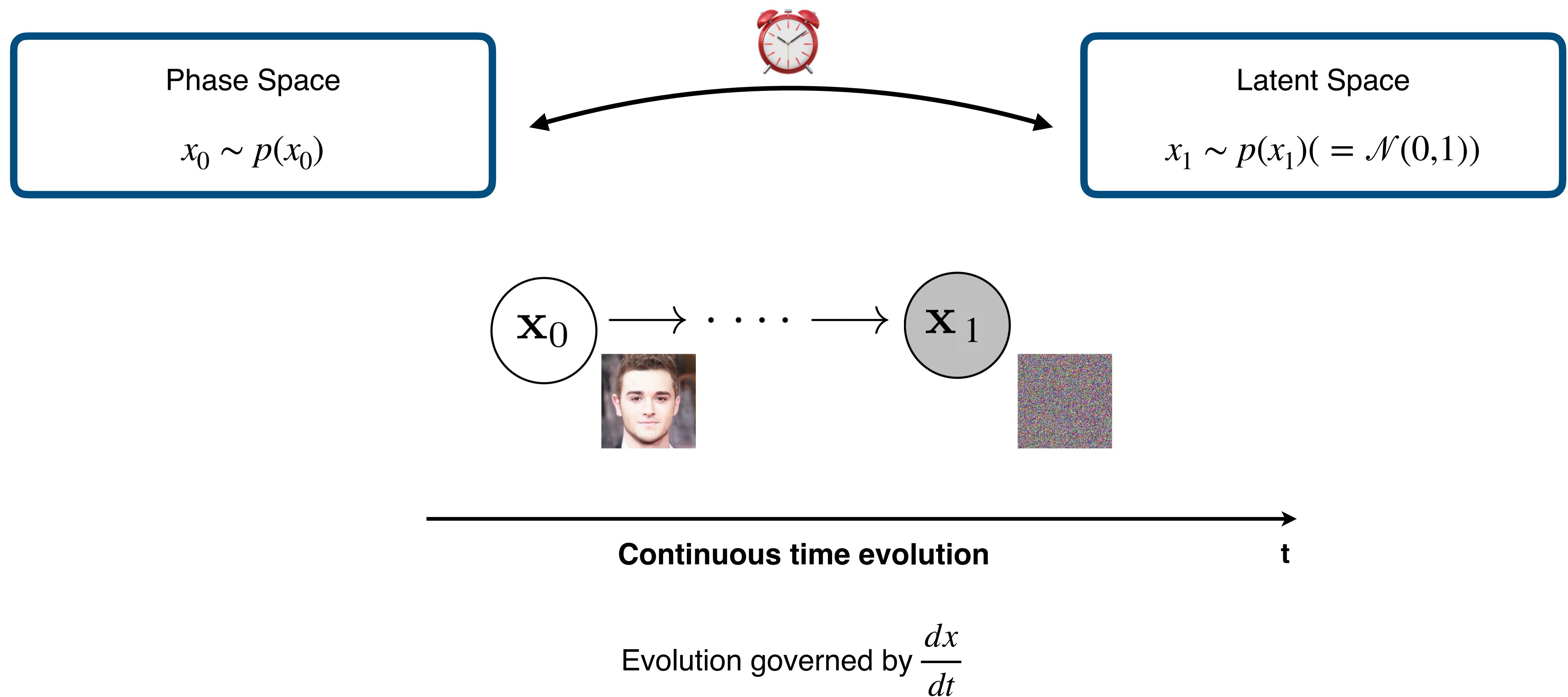
Diffusion Models (DDPM)



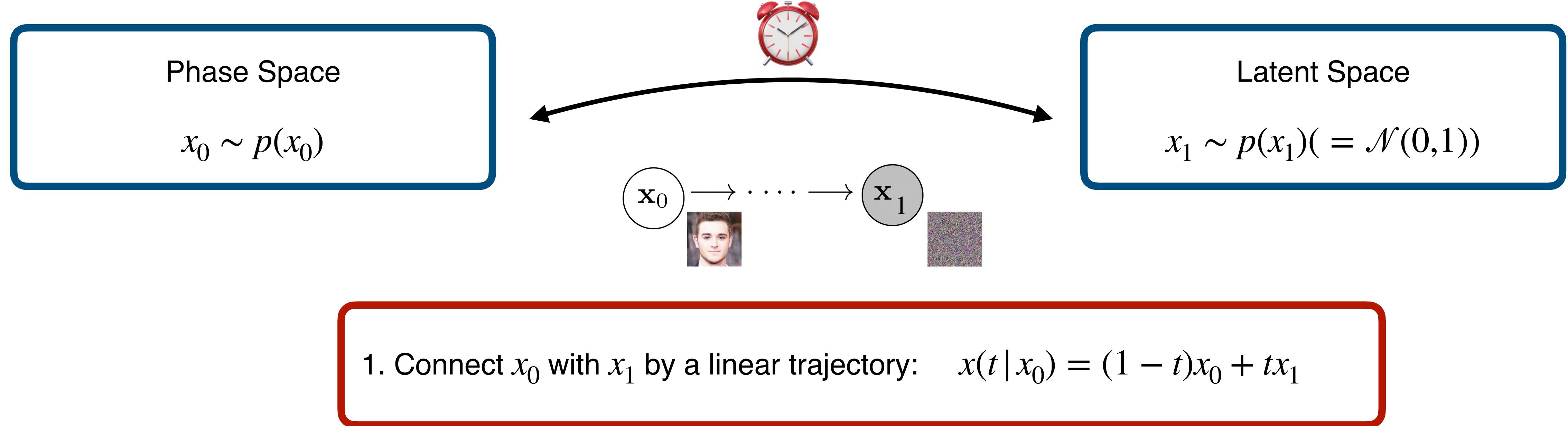
Diffusion Models (DDPM)



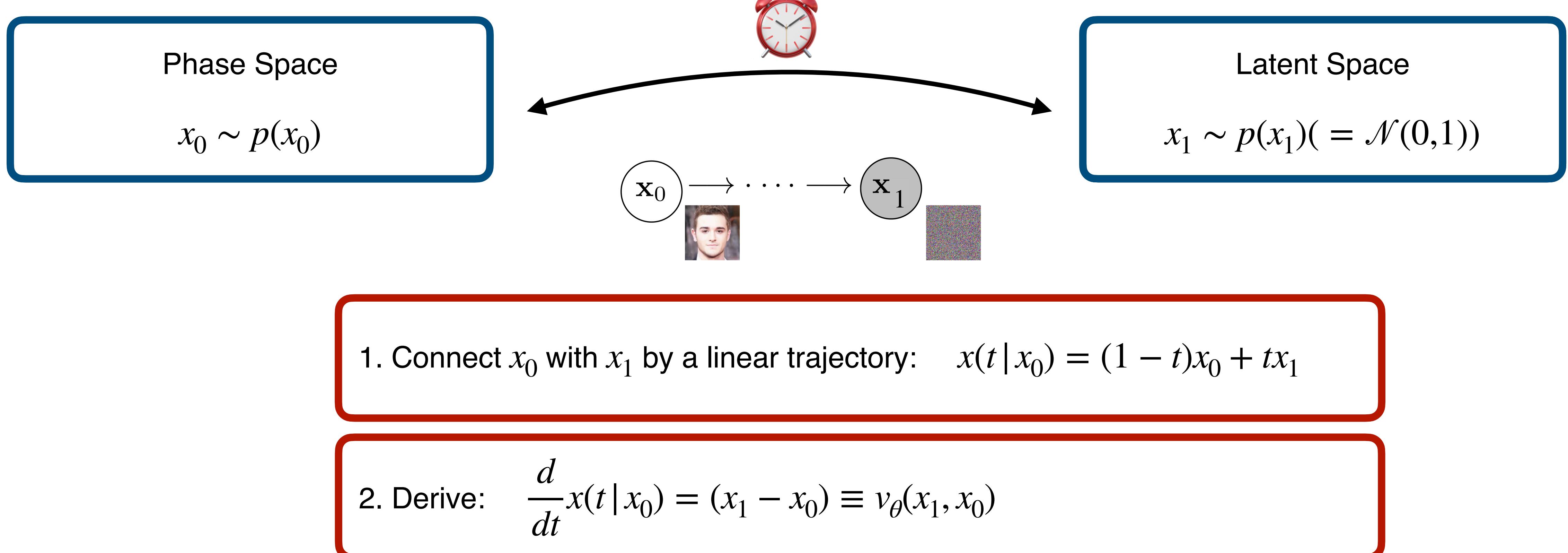
Diffusion Models (CFM)



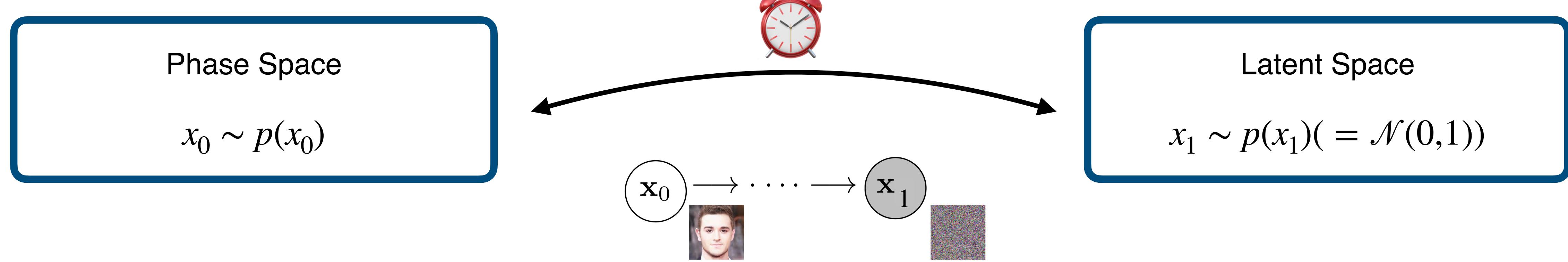
Diffusion Models (CFM)



Diffusion Models (CFM)



Diffusion Models (CFM)

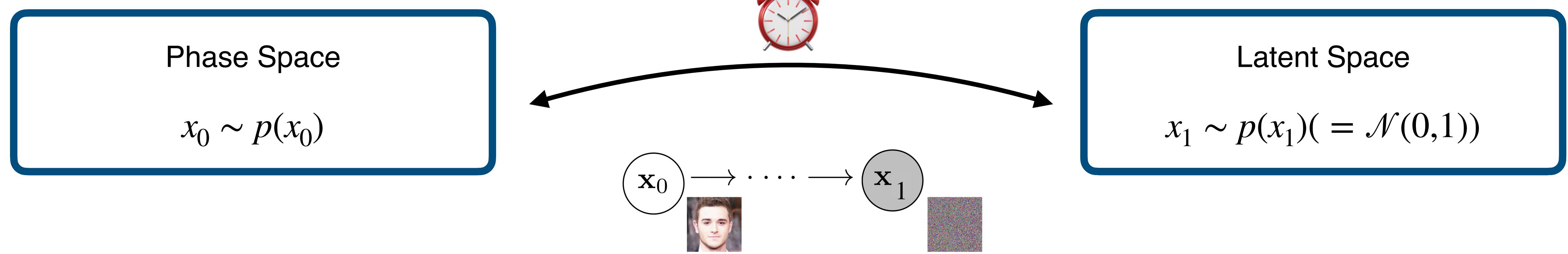


1. Connect x_0 with x_1 by a linear trajectory: $x(t | x_0) = (1 - t)x_0 + tx_1$

2. Derive: $\frac{d}{dt}x(t | x_0) = (x_1 - x_0) \equiv v_\theta(x_1, x_0)$

Velocity Field (can
be learned by simple
MSE loss)

Diffusion Models (CFM)

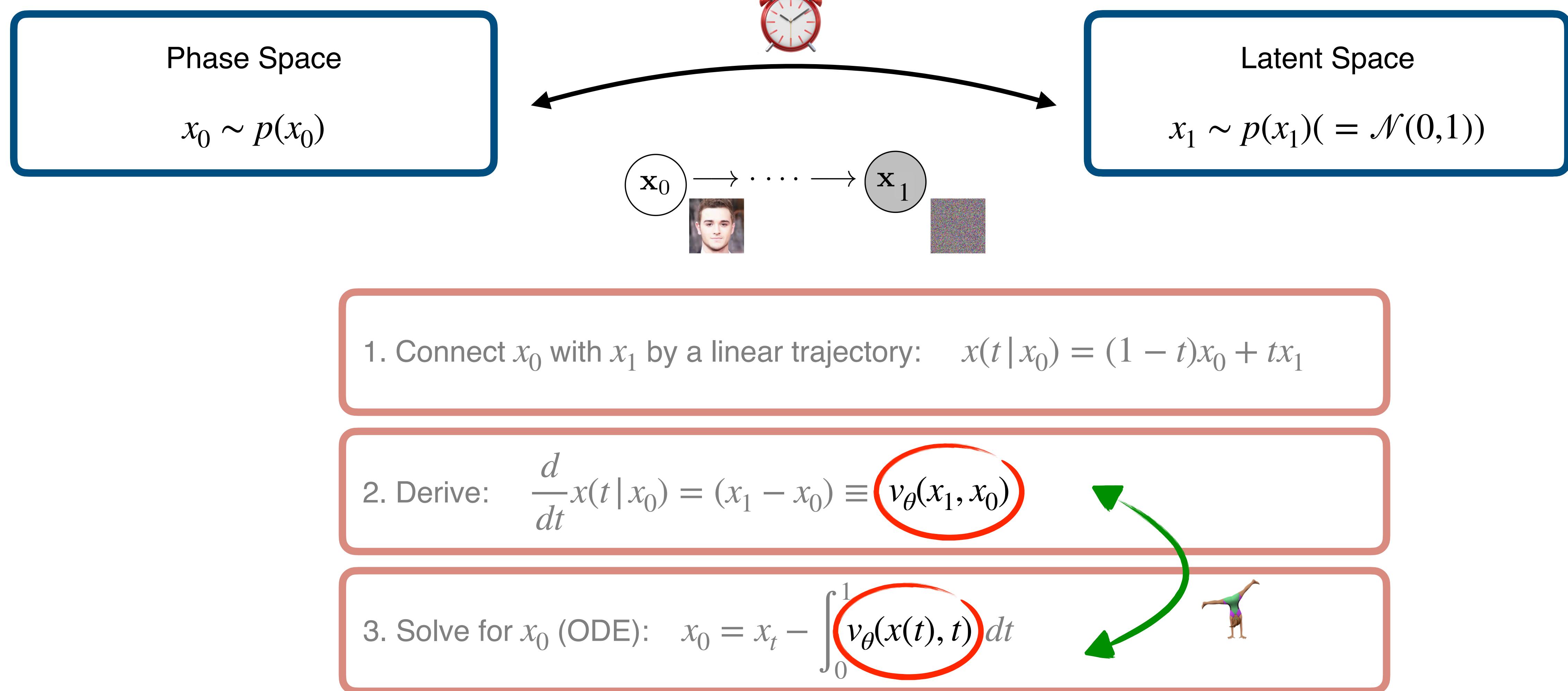


1. Connect x_0 with x_1 by a linear trajectory: $x(t | x_0) = (1 - t)x_0 + tx_1$

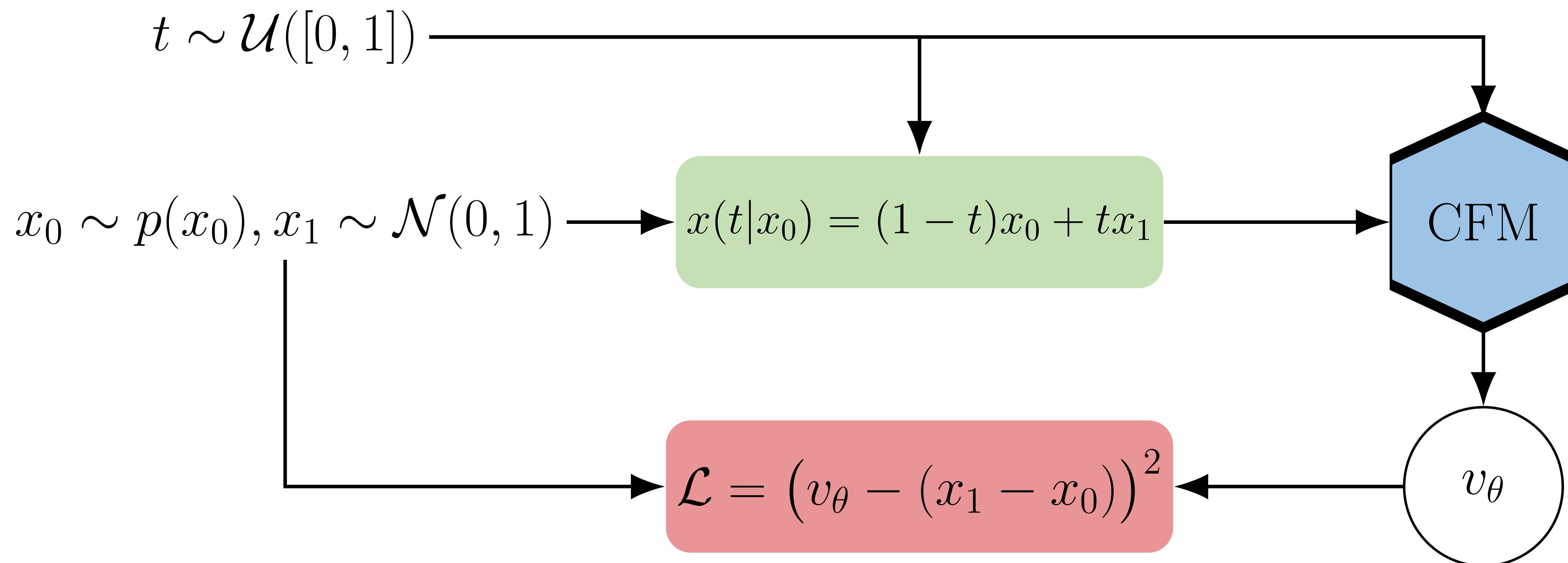
2. Derive: $\frac{d}{dt}x(t | x_0) = (x_1 - x_0) \equiv v_\theta(x_1, x_0)$

3. Solve for x_0 (ODE): $x_0 = x_t - \int_0^1 v_\theta(x(t), t) dt$

Diffusion Models (CFM)



Diffusion Models (CFM)



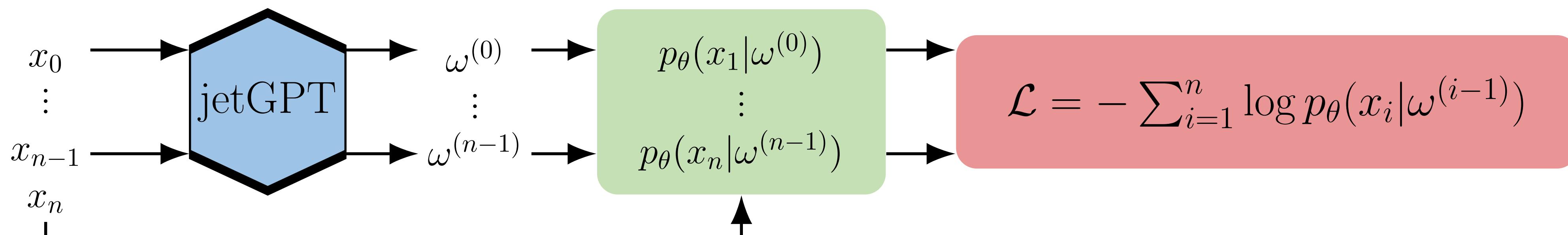
Autoregressiv Transformer

Estimate the density autoregressively

$$p_{\theta}(x) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) \approx p(x)$$



i goes over phase space dimension



Autoregressiv Transformer

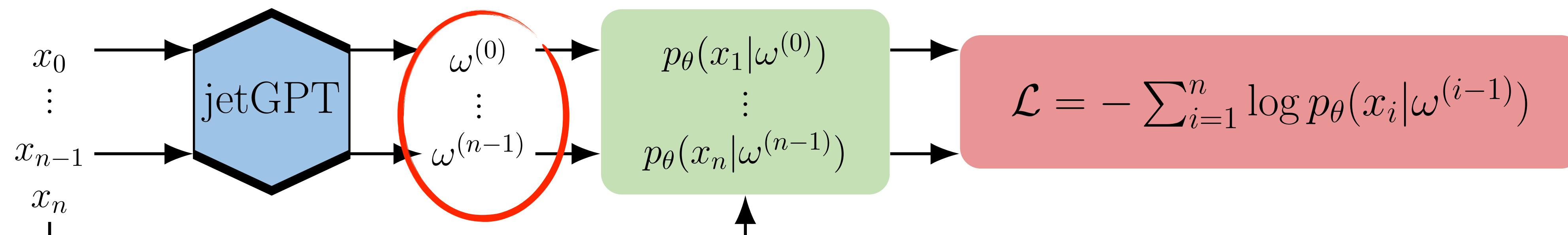
Parametrization

Estimate the density autoregressively

$$p_{\theta}(x) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) \approx p(x)$$

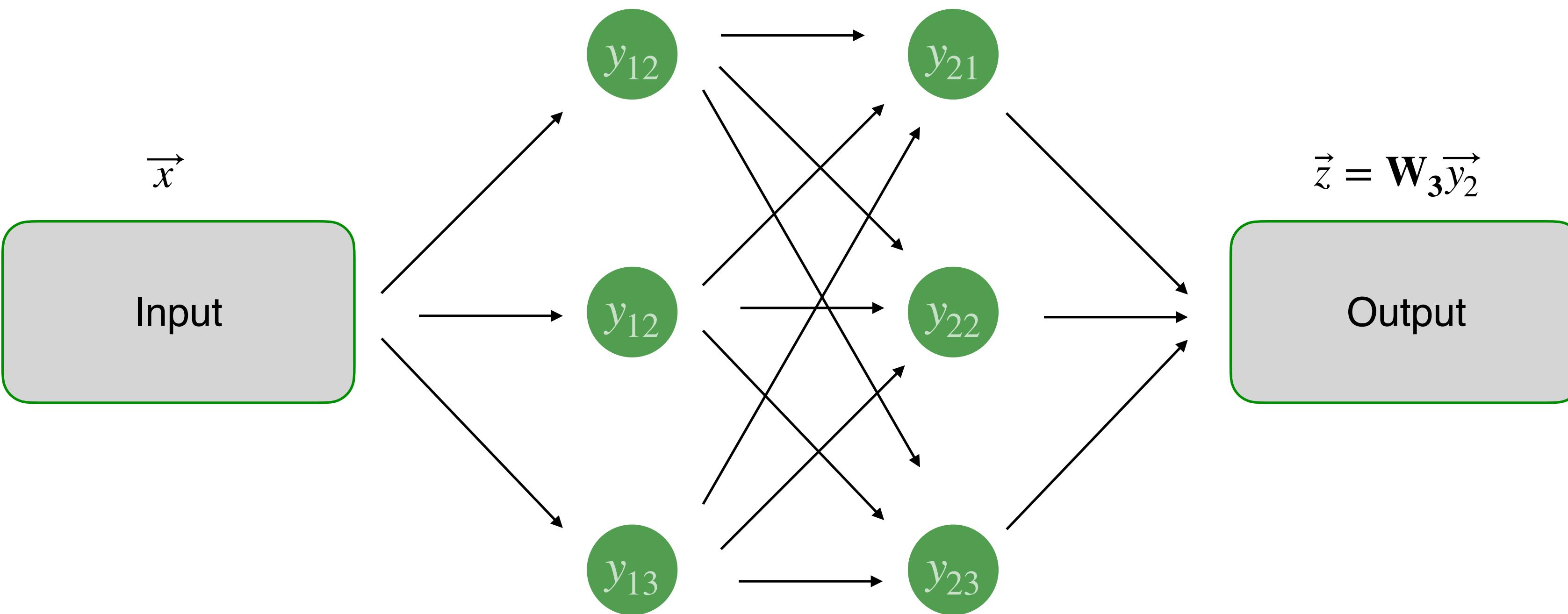


i goes over phase space dimension



What about uncertainties?

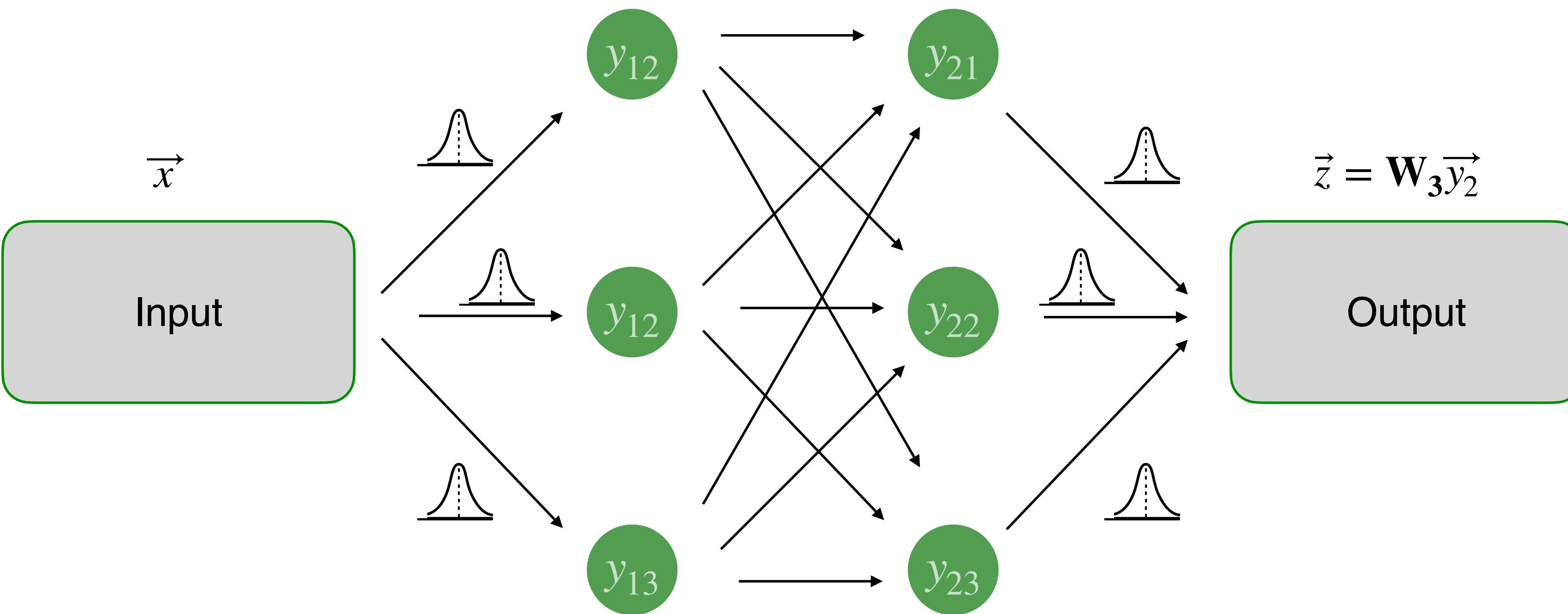
$$\vec{y}_1 = \mathbf{W}_1 \vec{x} \quad \vec{y}_2 = \mathbf{W}_2 \vec{y}_1$$



Once training is done: $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3$ fixed ("Network output is deterministic")

What about uncertainties?

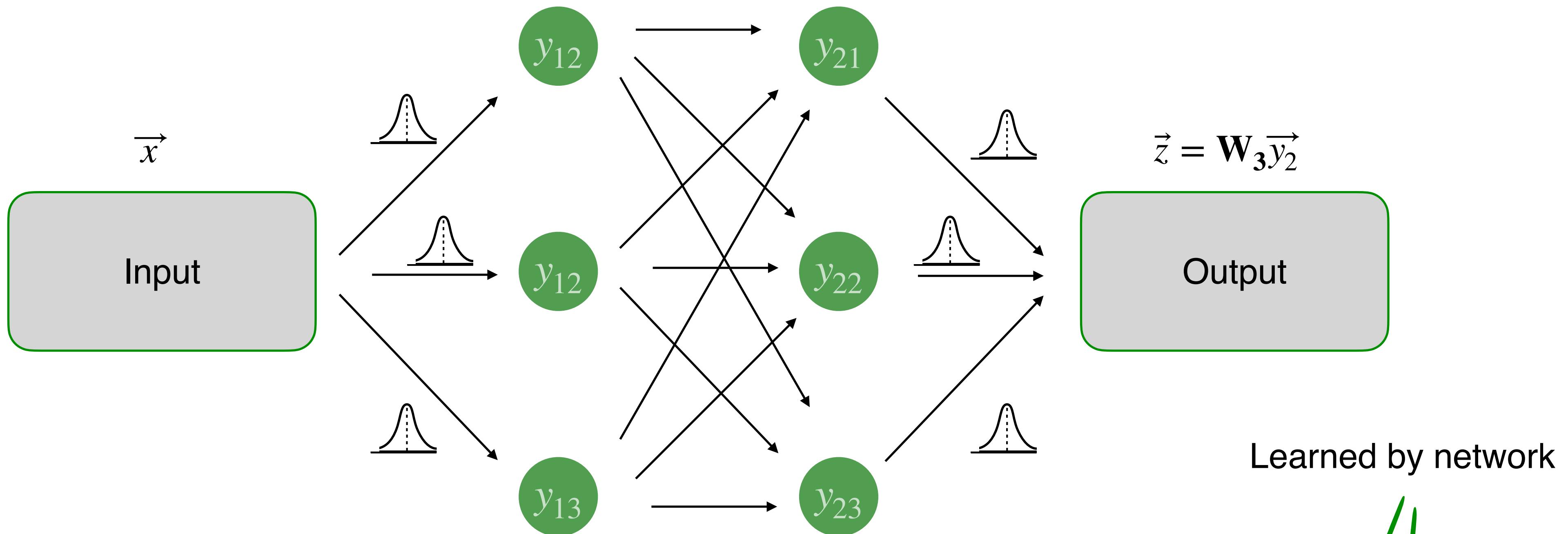
$$\vec{y}_1 = \mathbf{W}_1 \vec{x} \quad \vec{y}_2 = \mathbf{W}_2 \vec{y}_1$$



Bayesianization: We draw each entry from $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3$ from distribution $q(w | \mu_\phi, \sigma_\phi)$

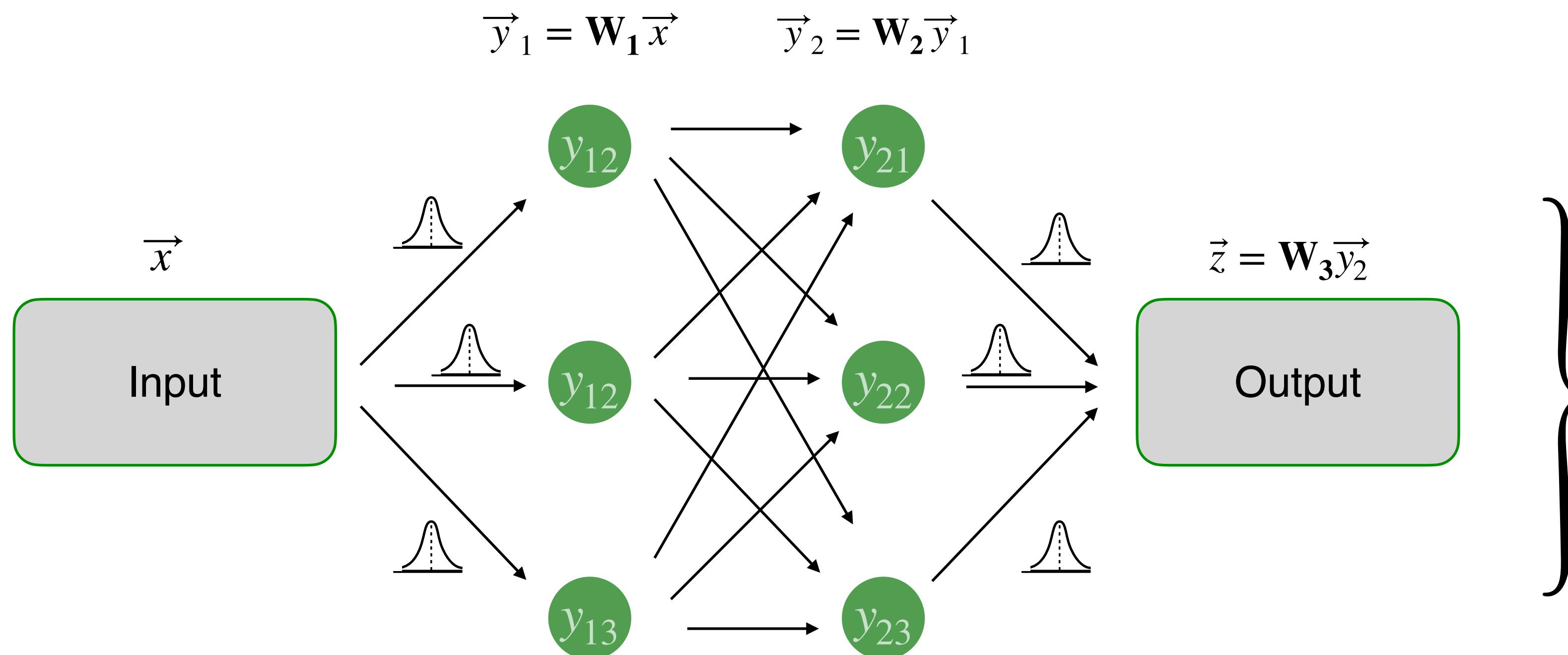
What about uncertainties?

$$\vec{y}_1 = \mathbf{W}_1 \vec{x} \quad \vec{y}_2 = \mathbf{W}_2 \vec{y}_1$$



Bayesianization: We draw each entry from $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3$ from distribution $q(w | \mu_\phi, \sigma_\phi)$

What about uncertainties?

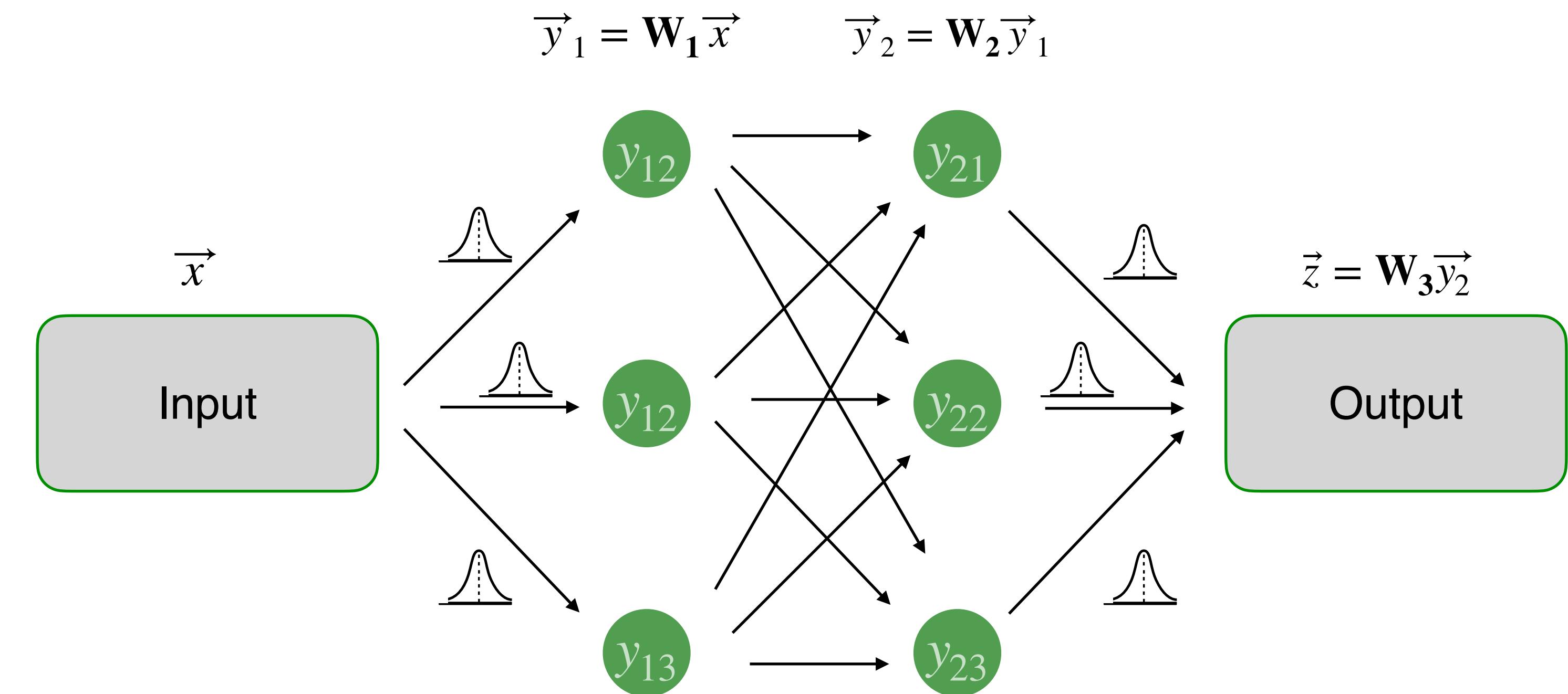


$$\langle \vec{z} \rangle = \frac{1}{N} \sum_i \vec{z}_i$$

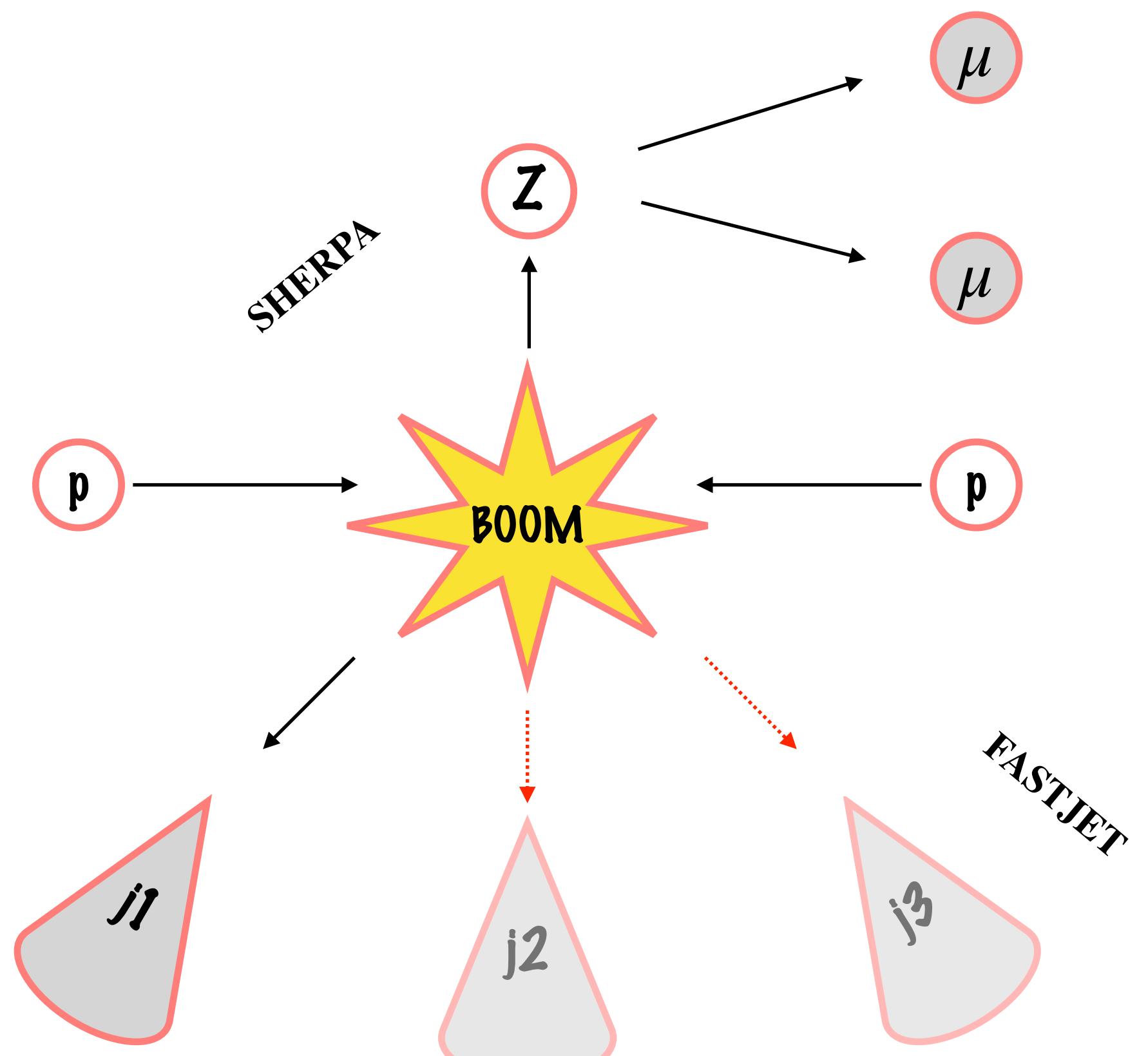
$$\sigma_{pred}^2 = \frac{1}{N} \sum_{i=1}^N (\langle \vec{z} \rangle - \vec{z}_i)^2$$

How to Bayesianize

1. Replace each linear layer with a *Bayesian* layer
2. Add additional regularisation term to loss



Concrete Application – LHC



$Z \rightarrow \mu\mu + \text{jets}$:

3 - 5 final state particles (including jets)

12 - 20 dimensional phase space

Smart preprocessing:

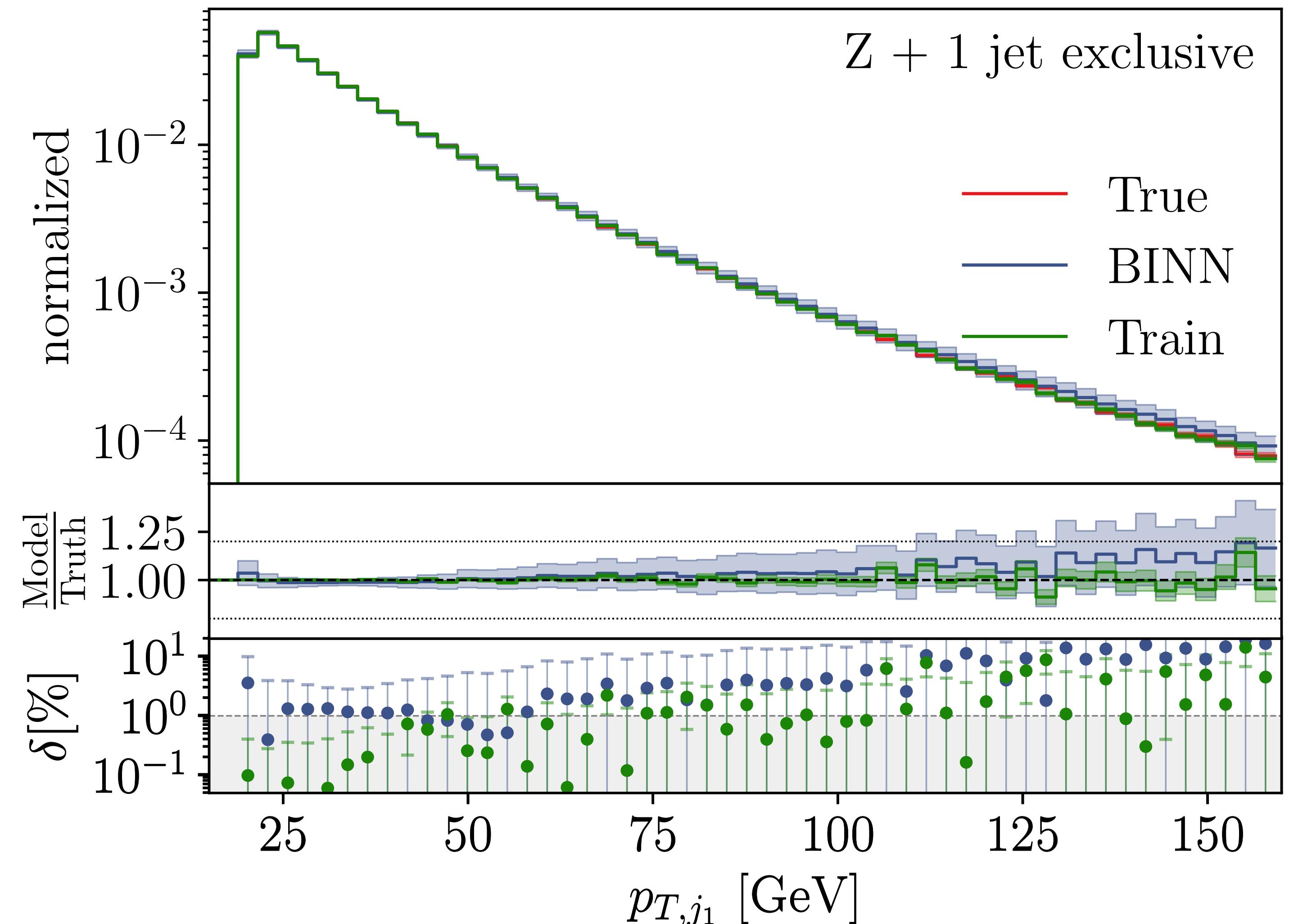
Global Phase Shift

Drop muon masses

→ reduces phase space to
9 - 17 dimensions

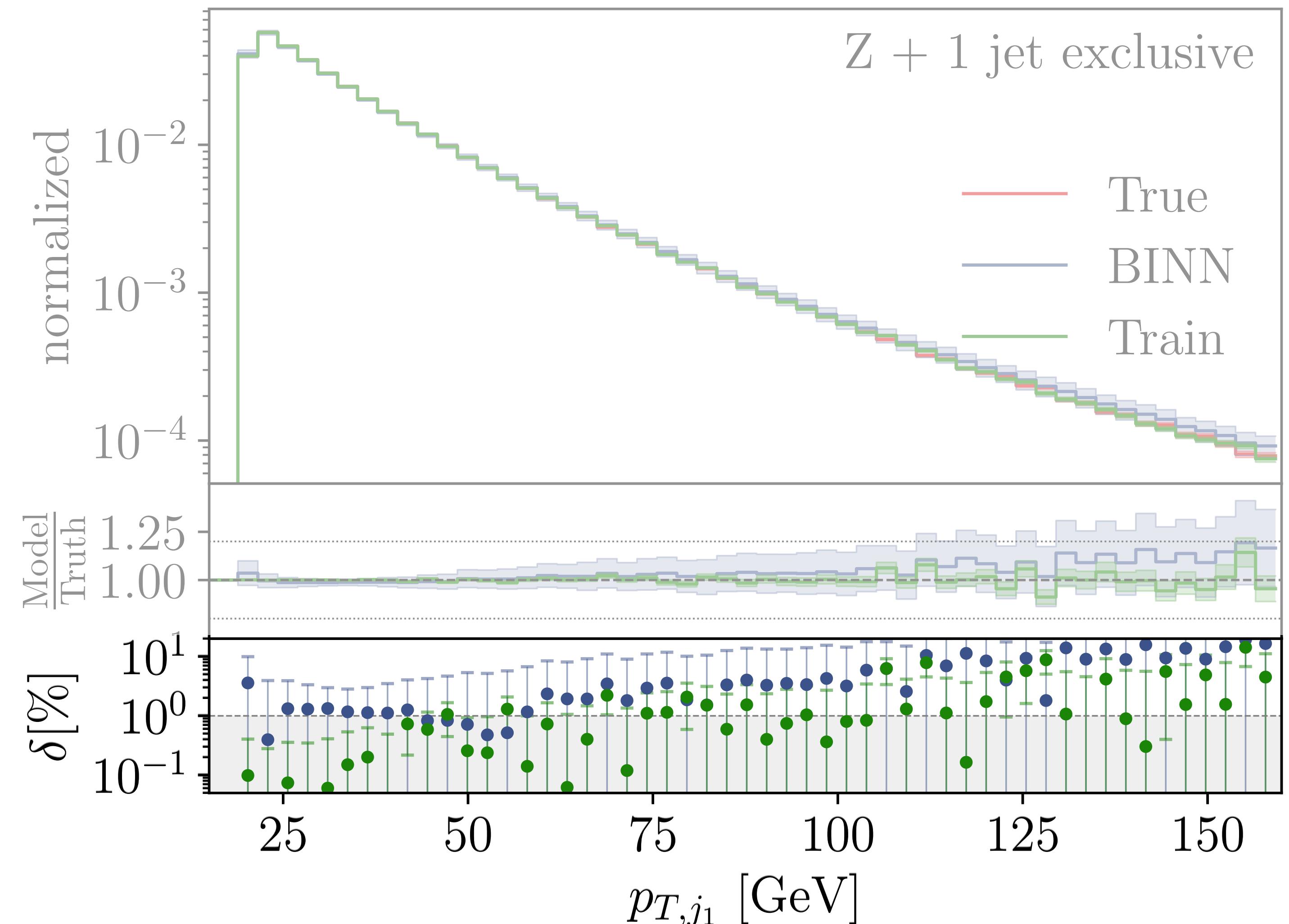
To be precise

- Previous studies showed: INNs can reach precision benchmark
- Percent level precision (comparable to statistical uncertainty)
- Uncertainty well defined



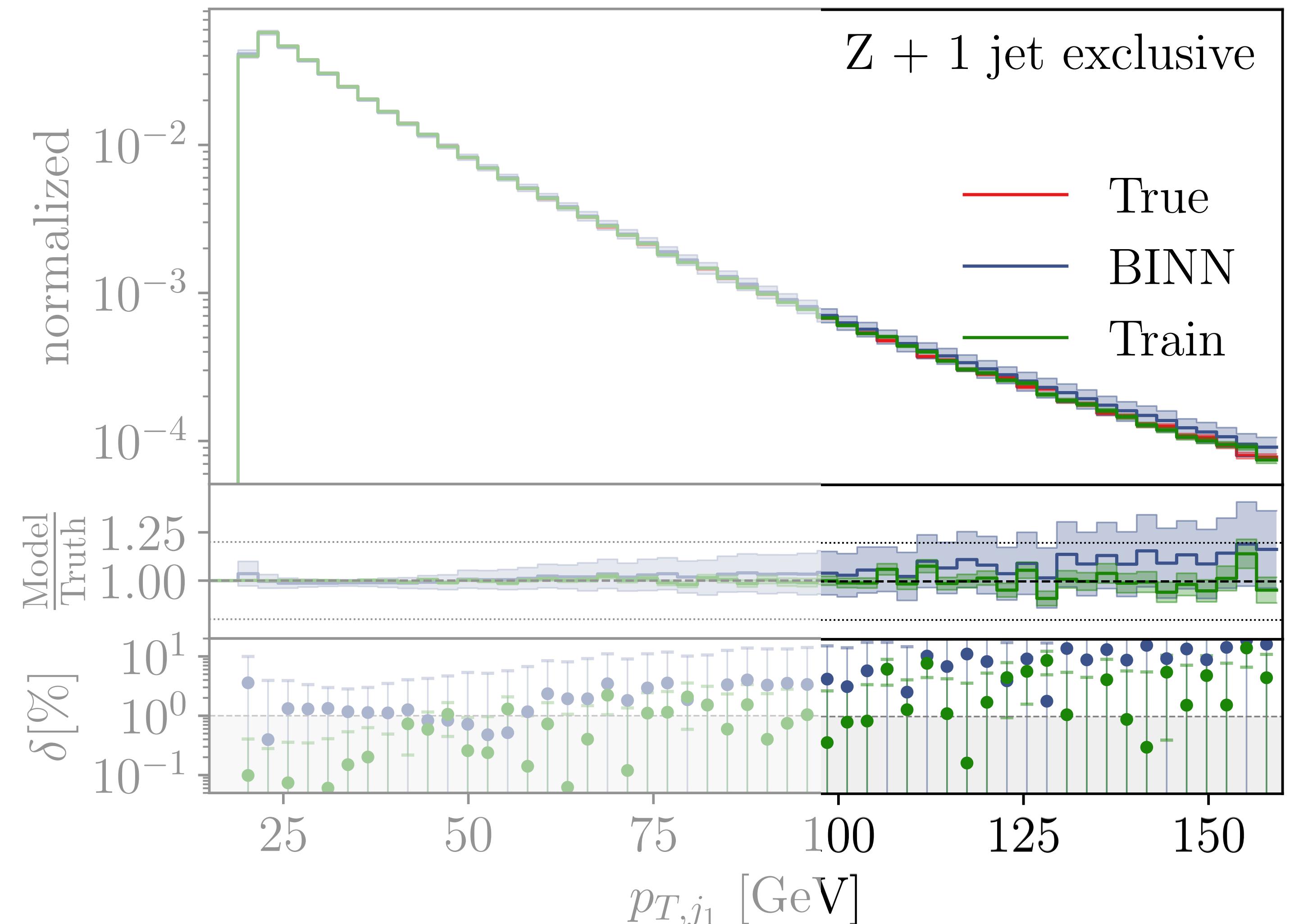
To be precise

- Previous studies showed: INNs can reach precision benchmark
- Percent level precision (comparable to statistical uncertainty)
- Uncertainty well defined



To be precise

- Previous studies showed: INNs can reach precision benchmark
- Percent level precision (comparable to statistical uncertainty)
- Uncertainty well defined



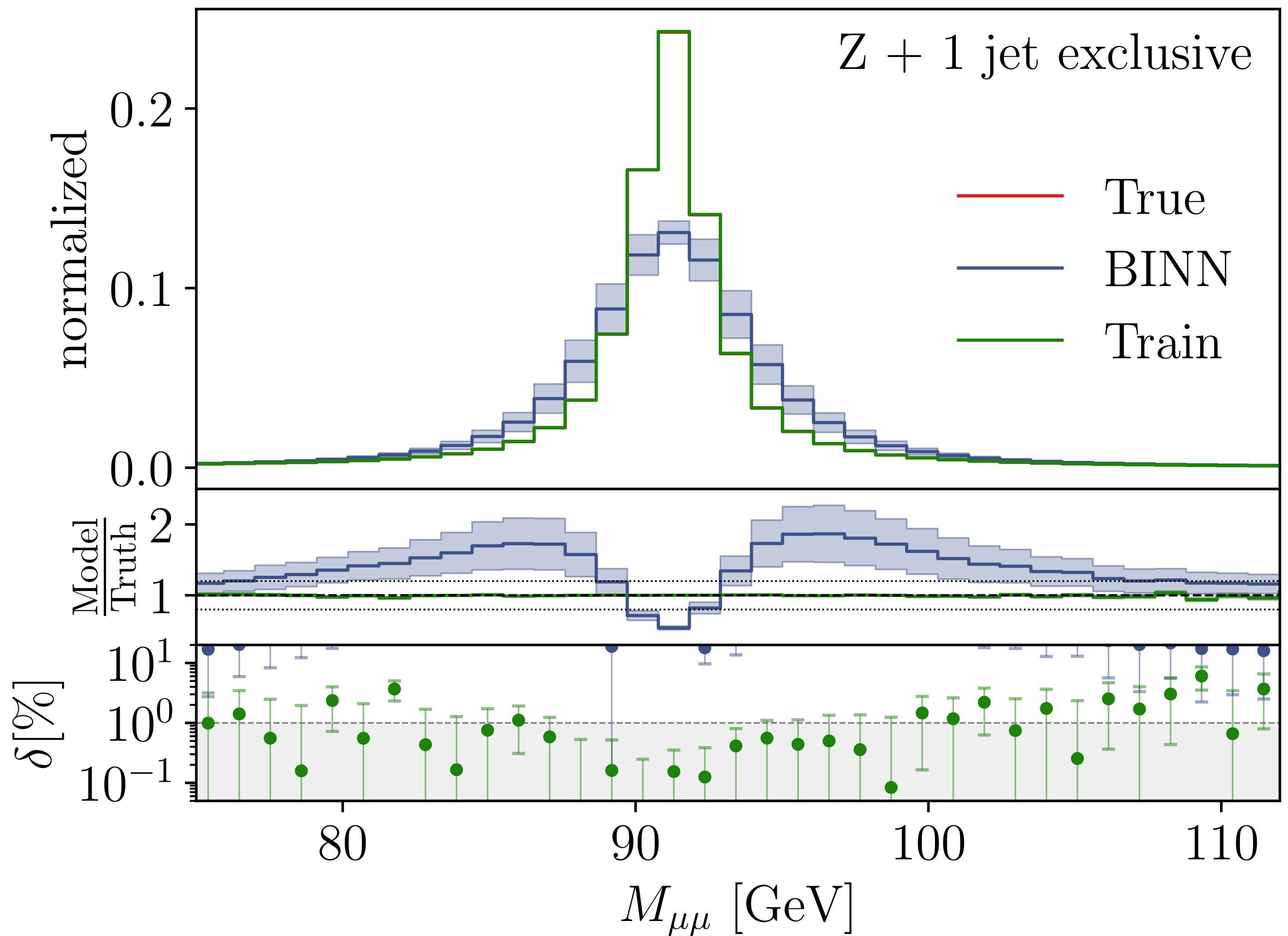
To be precise

Previous studies showed: INNs can reach precision benchmark

Percent level precision (comparable to statistical uncertainty)

Uncertainty well defined

BUT: only in one dimensional marginal distributions



To be precise

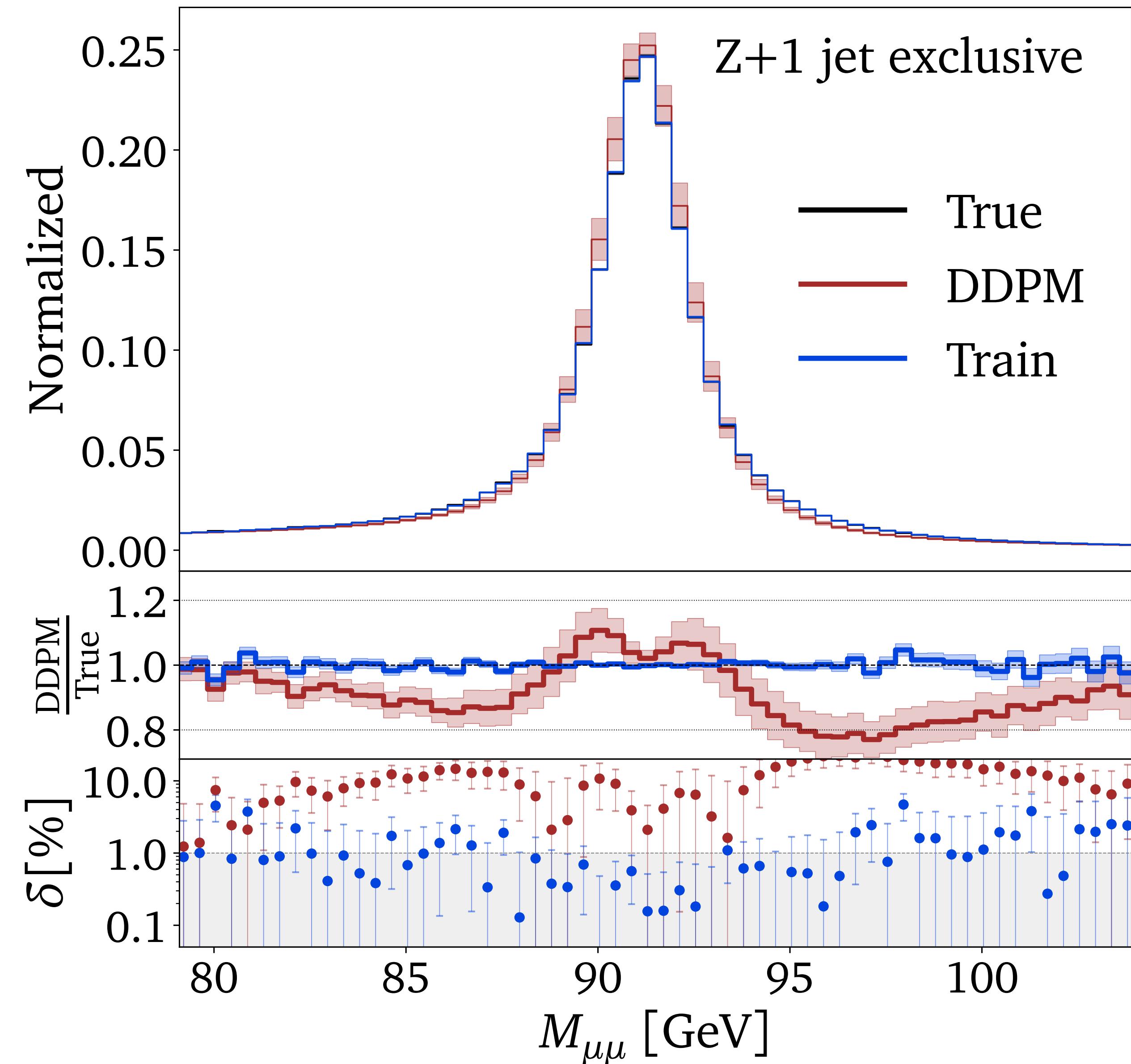
Previous studies showed: INNs can reach precision benchmark

Percent level precision (comparable to statistical uncertainty)

Uncertainty well defined

BUT: only in one dimensional marginal distributions

→ Diffusion Models surpass precision



To be precise

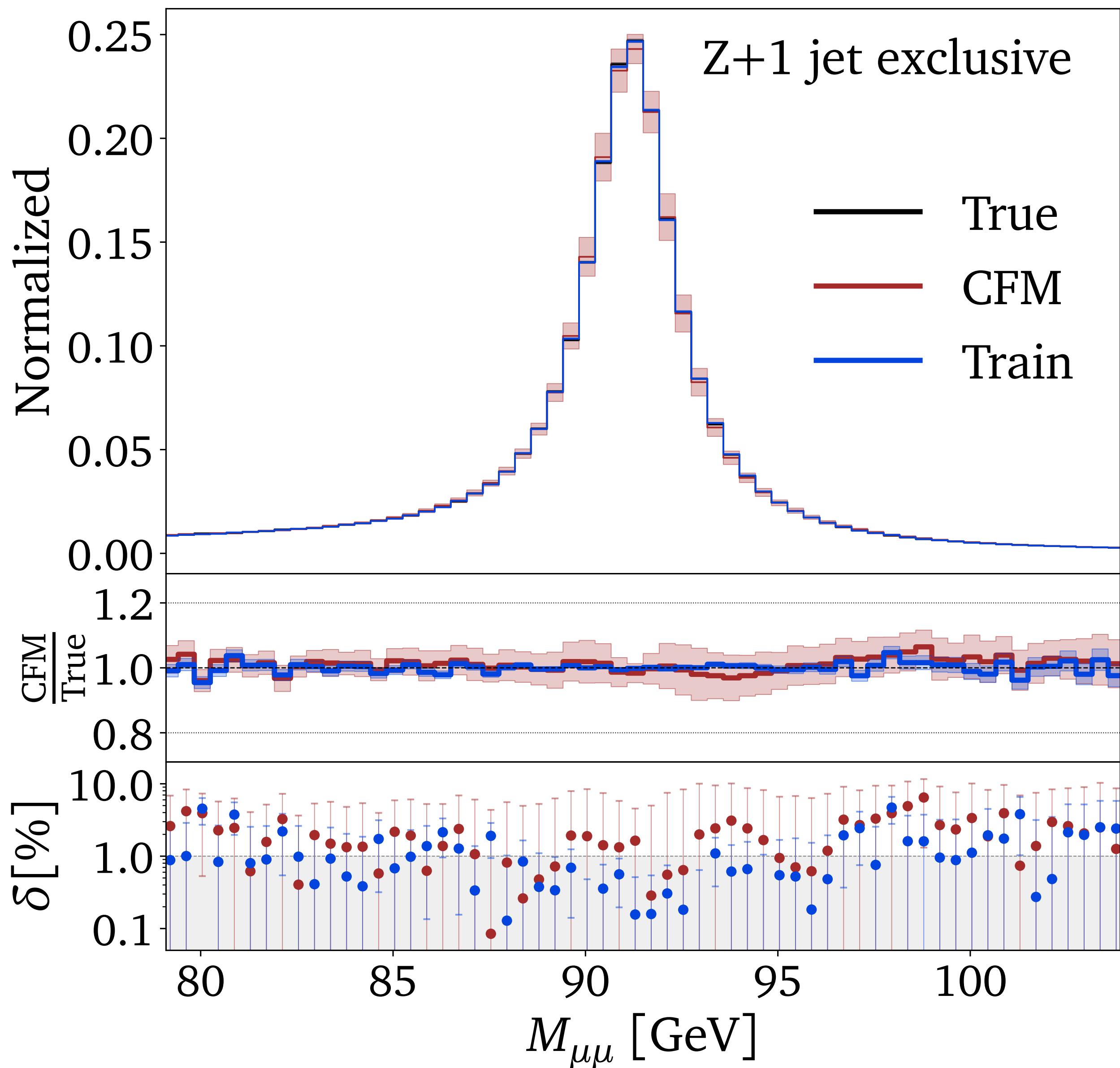
Previous studies showed: INNs can reach precision benchmark

Percent level precision (comparable to statistical uncertainty)

Uncertainty well defined

BUT: only in one dimensional marginal distributions

→ Diffusion Models surpass precision

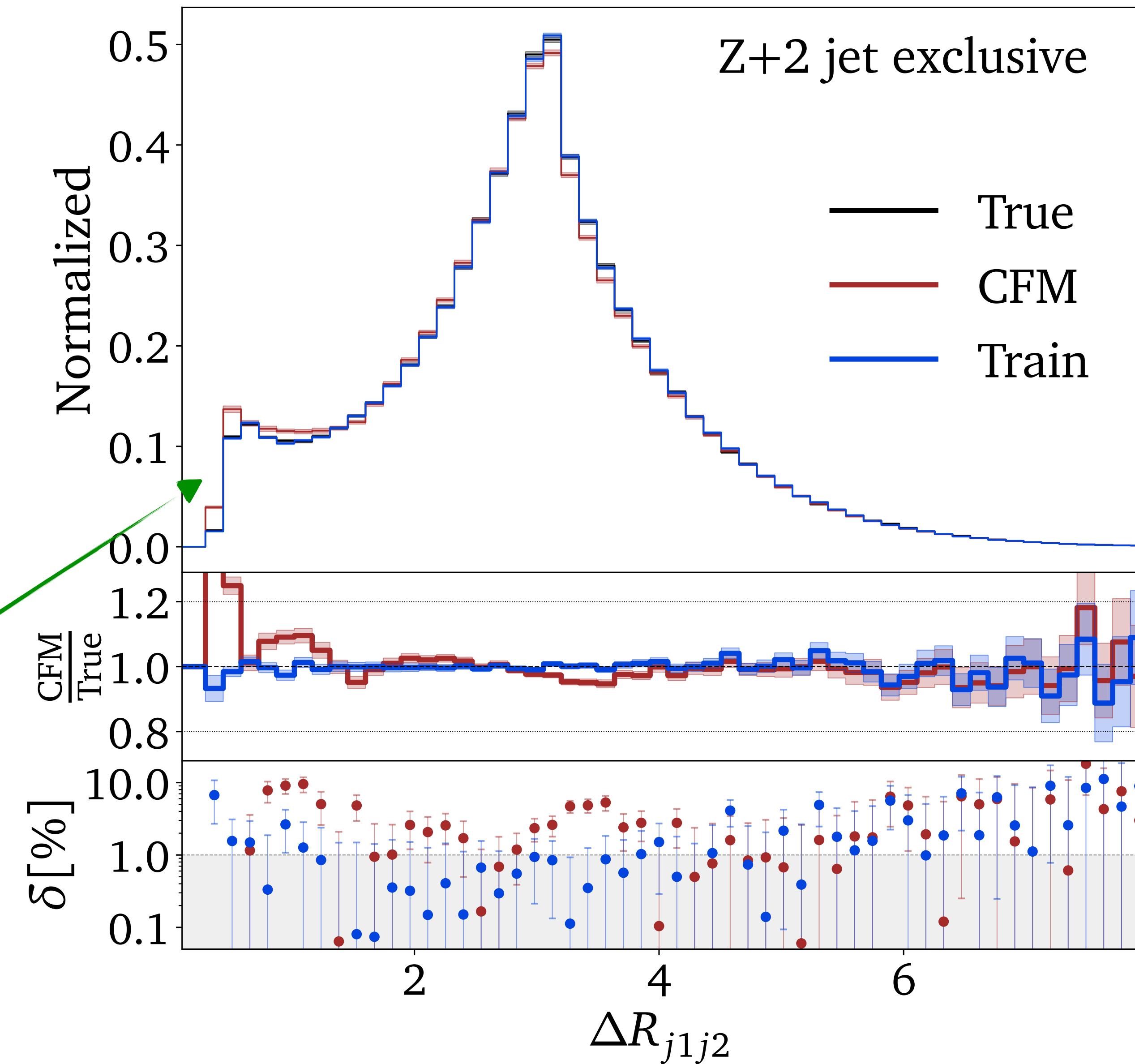


To be precise

For all non-autoregressiv models:
Difficult is to learn **sharp cuts** in correlations



Some tricks already applied



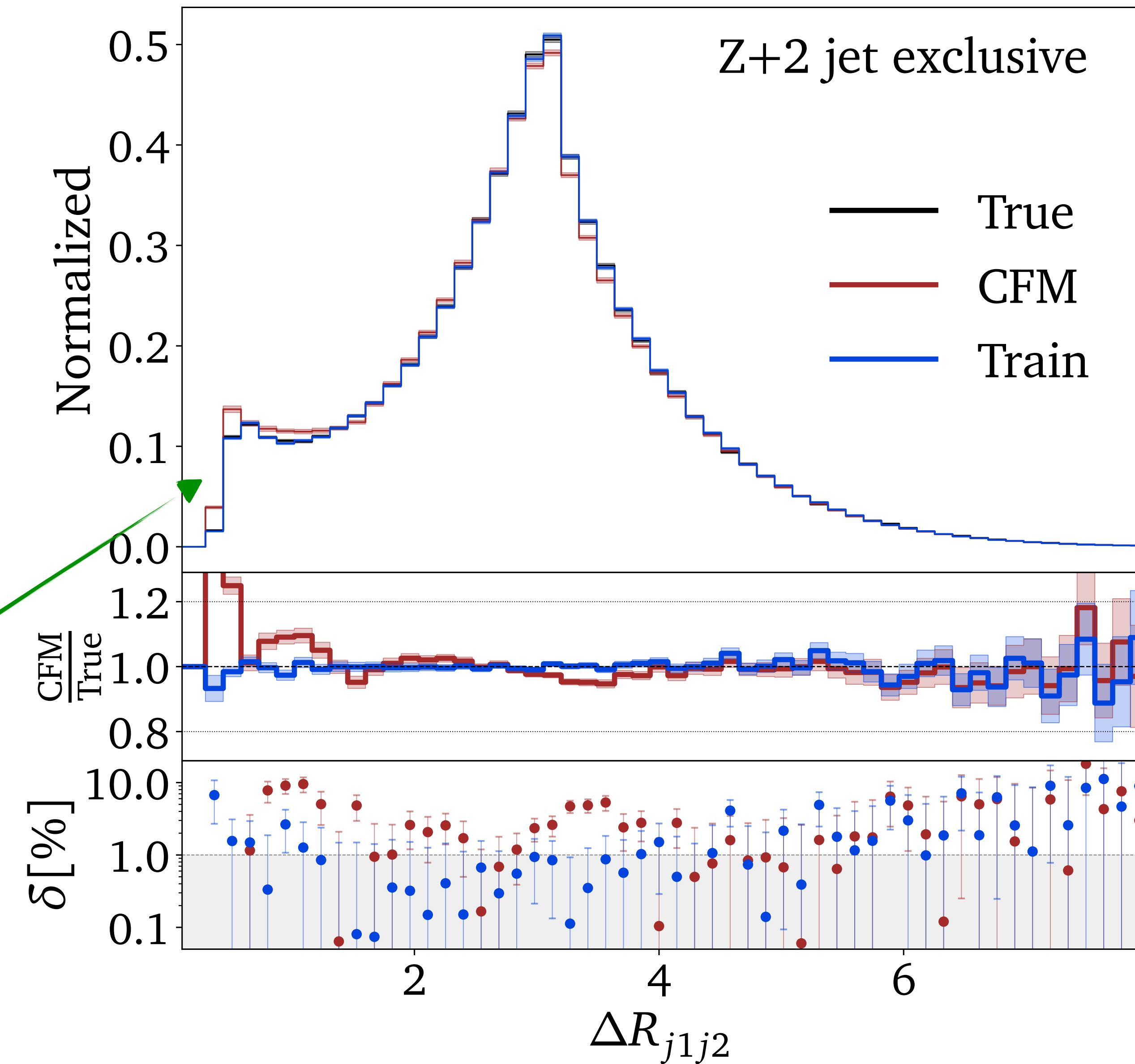
To be precise

For all non-autoregressiv models:
Difficult is to learn **sharp cuts** in correlations



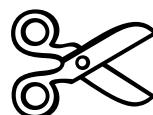
Could an autoregressiv transformer help?

Some tricks already applied



To be precise

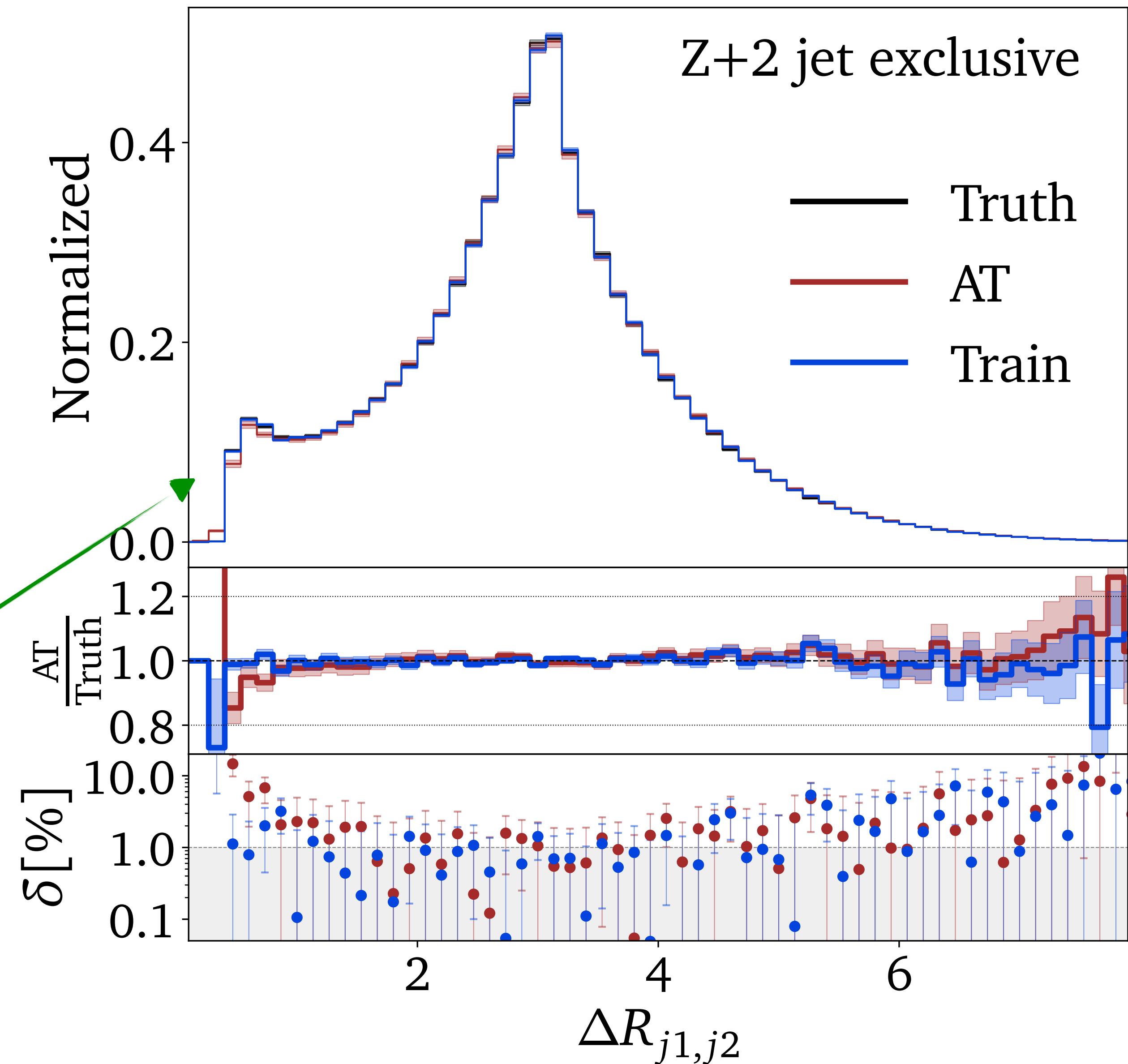
For all non-autoregressiv models:
Difficult is to learn **sharp cuts** in correlations



Could an autoregressiv transformer help?

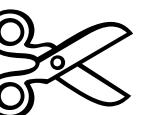
Yes

No tricks applied



To be precise

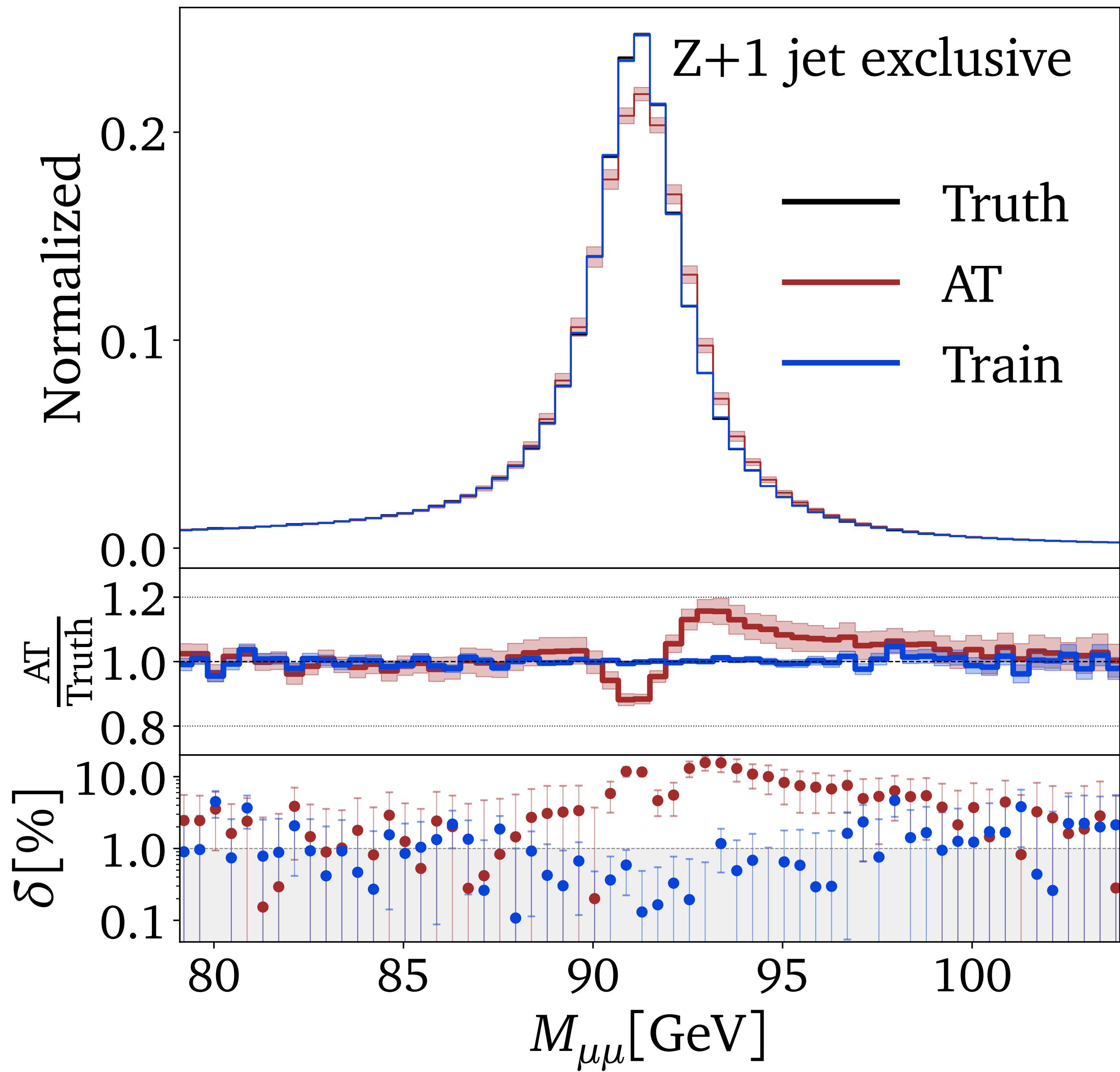
For all non-autoregressiv models:
Difficult is to learn **sharp cuts** in correlations



Could an autoregressiv transformer help?

Yes

BUT: to which cost?

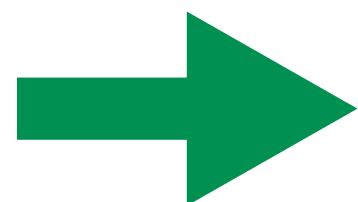


And now what?

New, “hyped” ML-Models can compete with current benchmark

All of them come with their own advantages and disadvantages

A lot of on-going research (generation speed up, precision, etc.)



Modern networks show potential to be applied to particle physics tasks