

# Simulation & Generative Models

Gregor Kasieczka

Email: [gregor.kasieczka@uni-hamburg.de](mailto:gregor.kasieczka@uni-hamburg.de)

Twitter/X: [@GregorKasieczka](https://twitter.com/GregorKasieczka)

COFI Winter School 2023

**CLUSTER OF EXCELLENCE**  
QUANTUM UNIVERSE

**U+H**  
  
Universität Hamburg  
DER FORSCHUNG | DER LEHRE | DER BILDUNG

  
**KISS**  
CDCS  
CENTER FOR DATA AND COMPUTING  
IN NATURAL SCIENCES

  
**FSP**  
CMS

  
**PUNCH**  
4 NFDI

**DASHH**

  
**PIER**  
Partnership of  
Universität Hamburg and DESY

GEFÖRDERT VOM

  
Bundesministerium  
für Bildung  
und Forschung

**Emmy  
Noether-  
Programm**  
Deutsche  
Forschungsgemeinschaft  
**DFG**  


























$$\begin{matrix} \frac{FA}{16} \\ \frac{E}{\Delta} \\ C) \rho \Delta = \rho a \end{matrix}$$

$$= E \frac{F}{E} \left( \frac{E}{C} H = \frac{\hat{v} \tau^7 c = = u f^p}{j^{1p} d^7 h c}$$

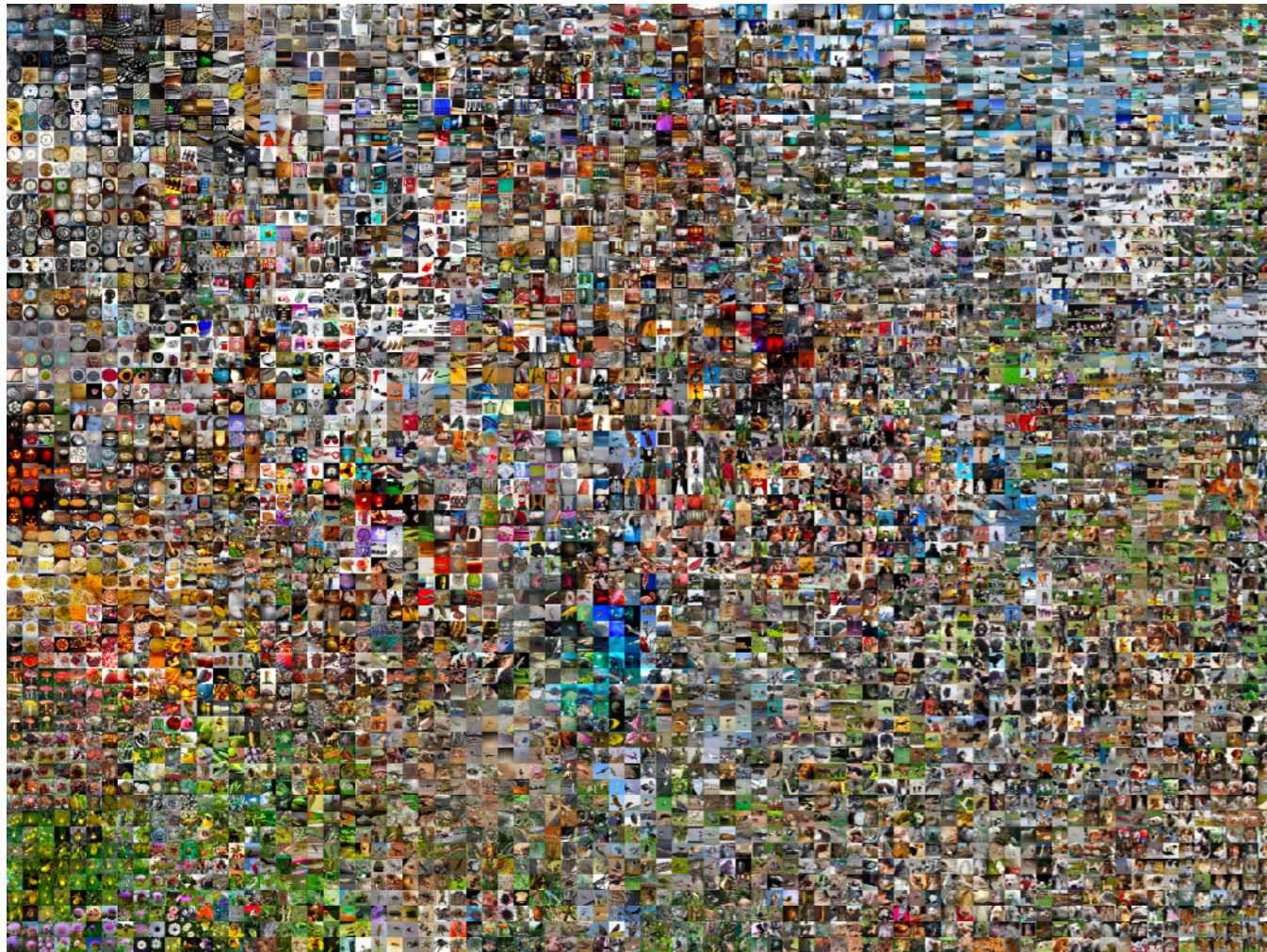
$$\begin{aligned} &= \frac{ns}{z} - B - \Pi \frac{0 C - )}{\frac{z R}{1}} \frac{3 C - )}{\frac{F A \sim}{15}} C X = -2 \\ &C - D \rho_2 (1 - \rho_2) (=) C_F = ) B \cup \rho_2 (y - N = -u = -v = \frac{n D}{B} D \frac{r i}{2 T} \frac{u}{A L c} D \\ &-^A = \frac{20}{10} | C \frac{\Delta R_{\sim}}{\lambda L} P - ) - D - L - B \frac{25}{\lambda D} C s = \frac{z N}{-1} \\ &z C = C F \cdot j^7 E = \rho = n e \end{aligned}$$

$$\begin{aligned} C P = P ) &= A \cdot X \\ &= A \cdot K \cdot P - X \end{aligned}$$



# Motivation

**Have:** input examples  
(collision events,  
detector readouts, ...)



**Want:** **more data**

Specifically: new data similar to  
the input, but not exact copies

How to encode in neural net?

Uses:

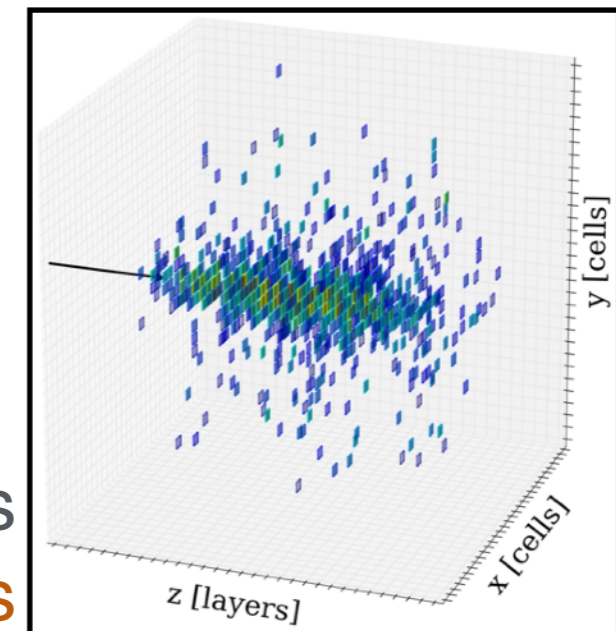
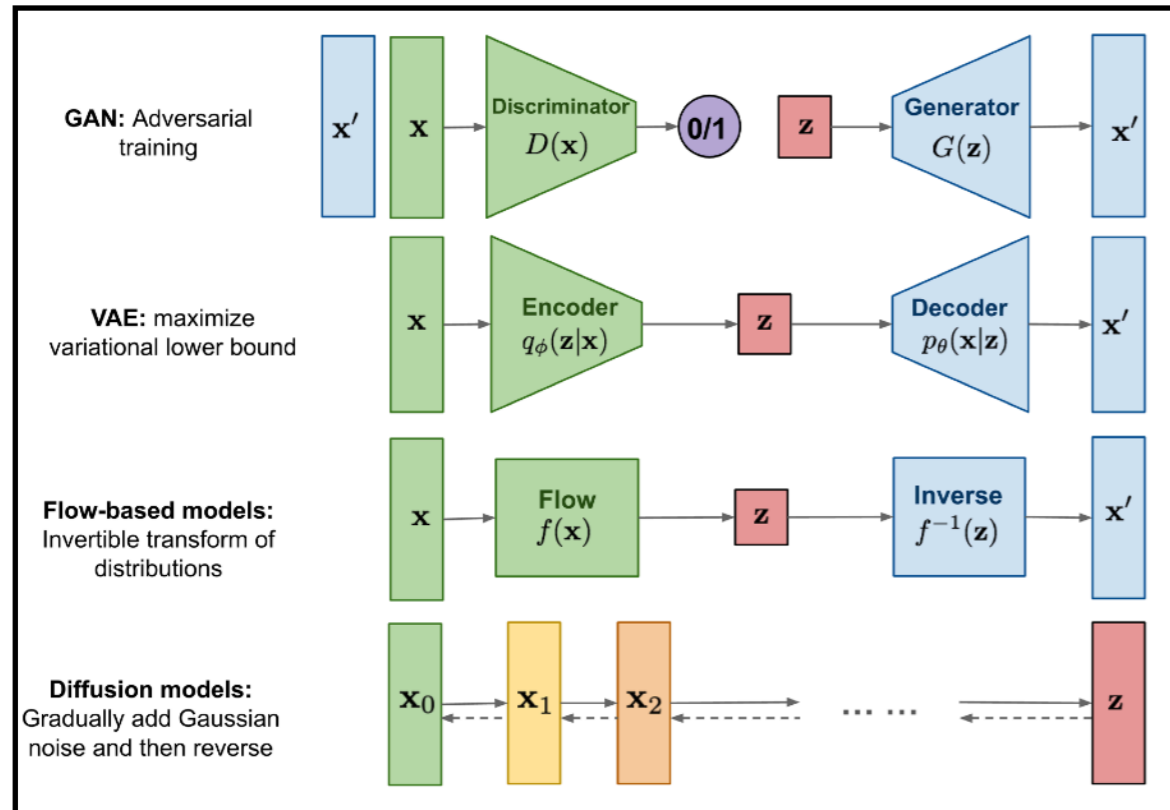
- Detector Simulation
- In-situ background estimation
- Surrogate models
- ...



# Overview

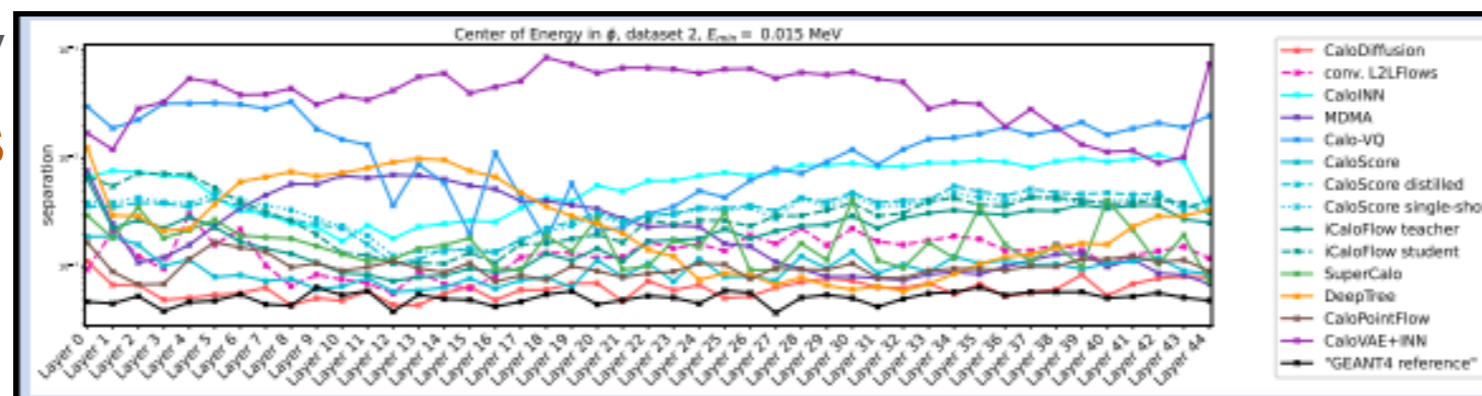
## 1. Common architectures\*

- > GANs, VAEs, NF today
- > Diffusion & CNF tomorrow



## 2. Physics applications

## 3. Quality metrics

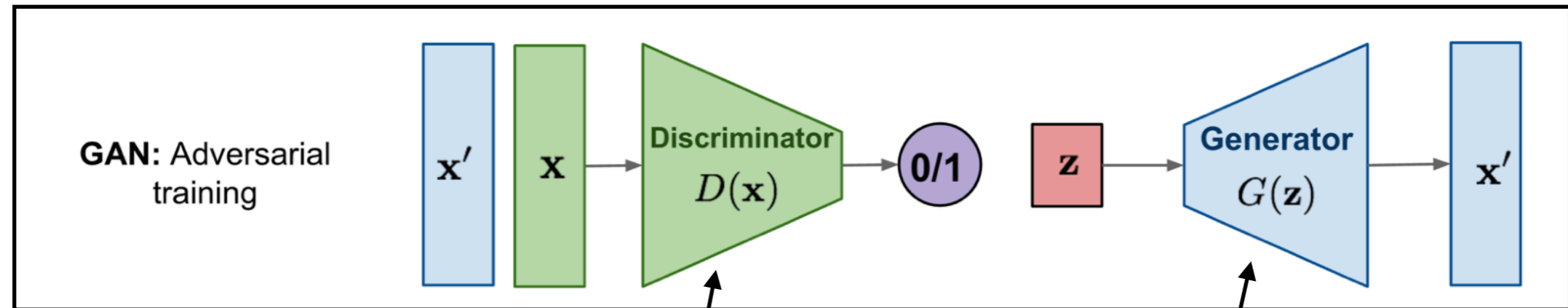




# **Generative Adversarial Networks**



# Generative Adversarial Networks



Generative Adversarial Networks (GANs) consist of **2 networks**

Provides feedback on quality of examples

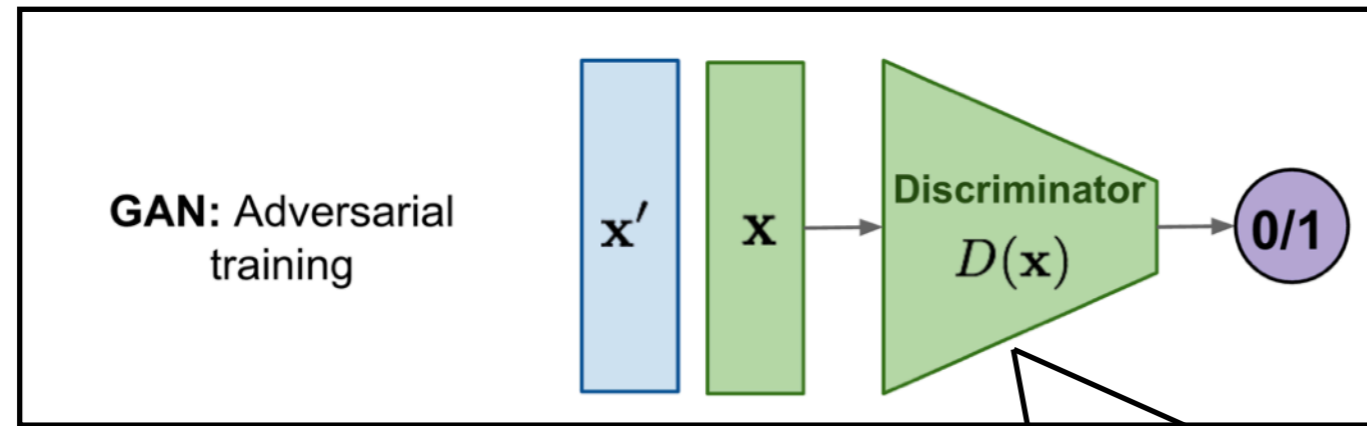
Maps random noise to realistic examples







# Generative Adversarial Networks



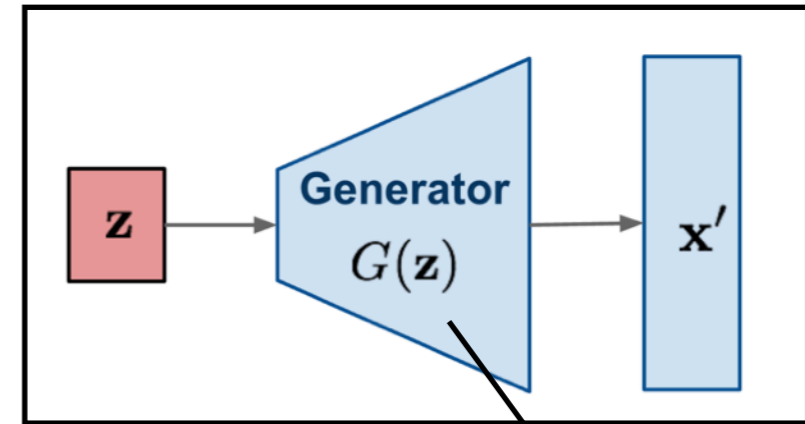
Training objective:  
Binary cross entropy

Maximise for  
discriminator

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$



# Generative Adversarial Networks



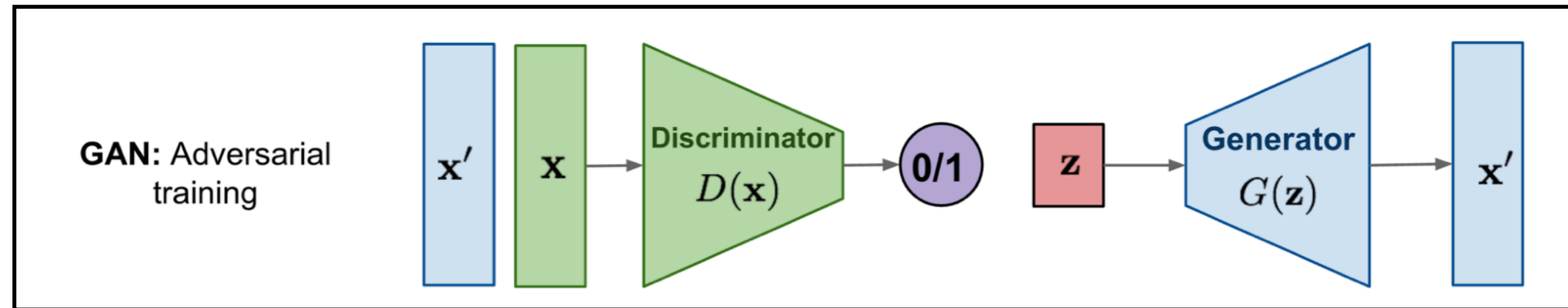
Training objective:  
Binary cross entropy

Minimise for generator

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$



# Generative Adversarial Networks



Training objective:

Binary cross entropy

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

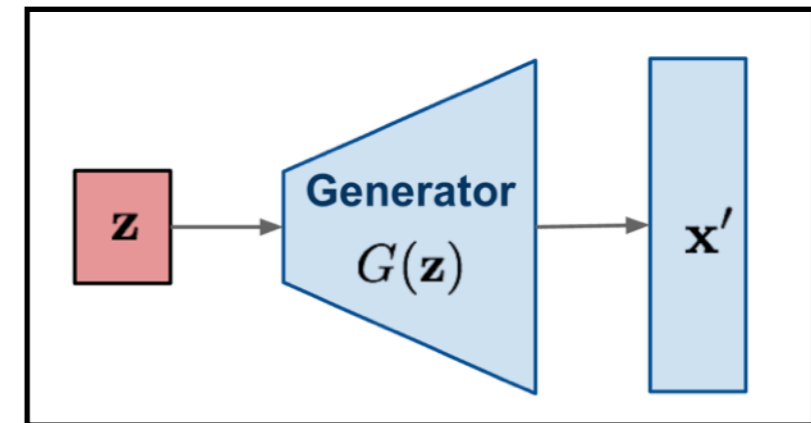
At (Nash) equilibrium:

Generator produces realistic examples

Discriminator is maximally confused



# Generative Adversarial Networks



Training objective:

Binary cross entropy

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

For generation:

Sample from **Generator**

Discard **Discriminator**



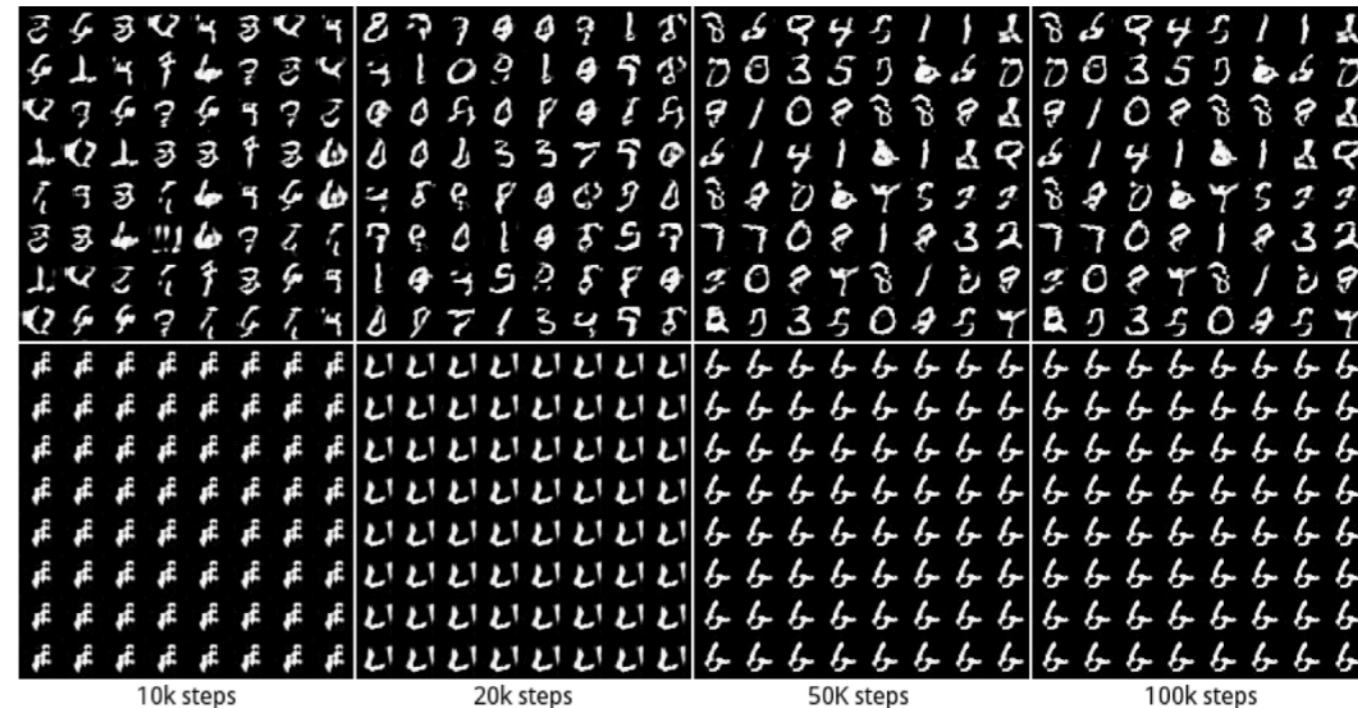
# Comments on GANs

## Architecture:

- Low complexity, **fast** and adaptable

## Learning:

- **Unstable** training
- Matching of generator/discriminator (**vanishing gradients**)
- Mode collapse
- Loss function not interpretable



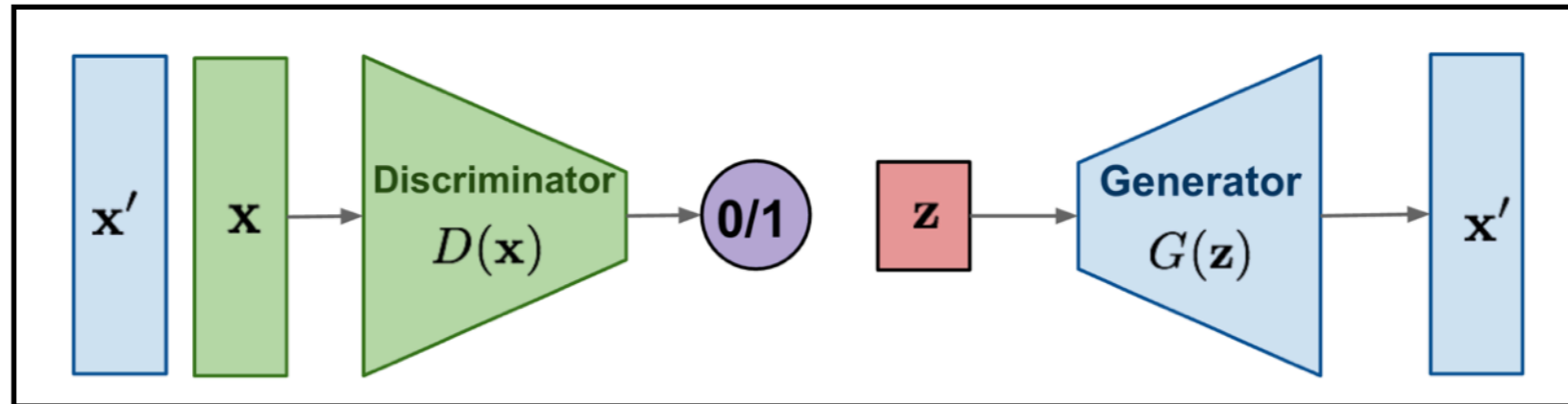
Mode collapse

## Maturity:

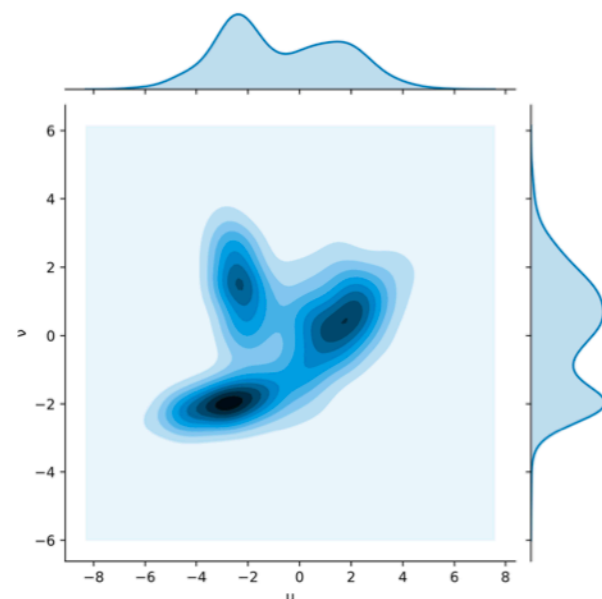
- Well established, many variants and extensions



# Wasserstein GAN



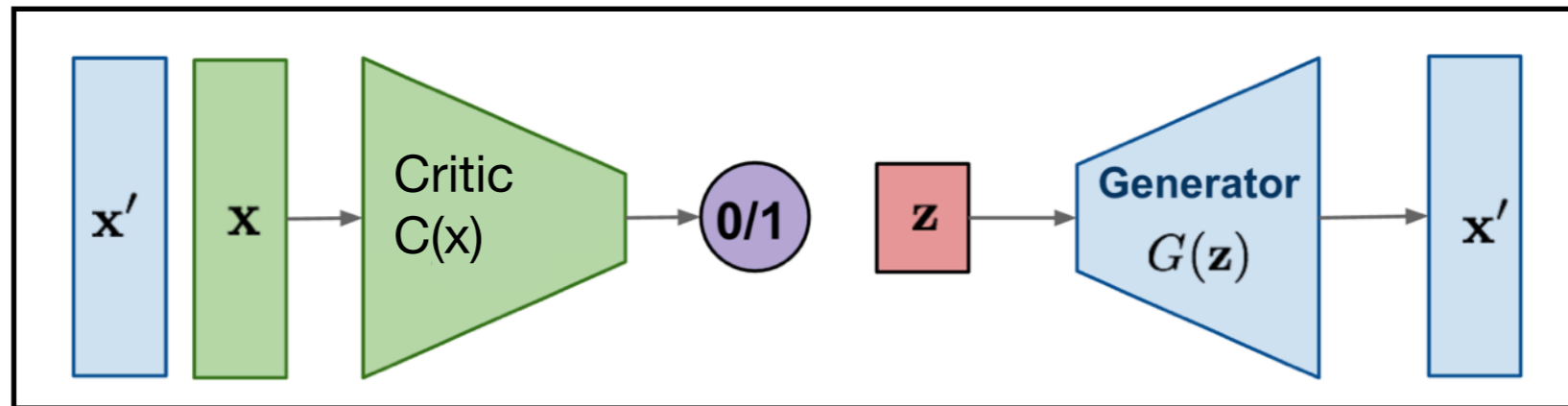
- Standard GANs minimise Jensen-Shannon divergence of generator output and true data
  - Not best measure, e.g. for non-overlapping distributions
- Replace with Wasserstein / Earth-Mover-Distance



$$W_p(\mu, \nu) = \left( \inf_{\gamma \in \Gamma(\mu, \nu)} \mathbf{E}_{(x, y) \sim \gamma} d(x, y)^p \right)^{1/p}$$



# Wasserstein GAN



GAN loss: 
$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [\log(D(\mathbf{x}))] + \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [\log(1 - D(\tilde{\mathbf{x}}))]$$

Wasserstein GAN

loss\*:

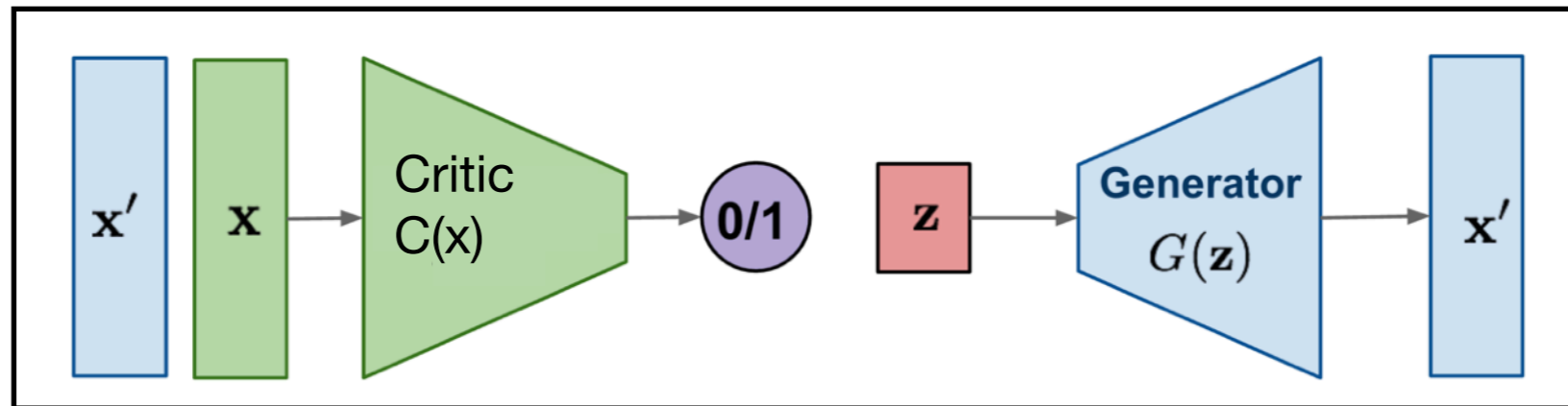
$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})]$$

Requires bounded Lipschitz norm,  
e.g. via term in loss

\* Some mathematics  
involved from earth  
mover distance to here



# Wasserstein GAN



GAN loss: 
$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [\log(D(\mathbf{x}))] + \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [\log(1 - D(\tilde{\mathbf{x}}))]$$

Wasserstein GAN

loss: 
$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})]$$

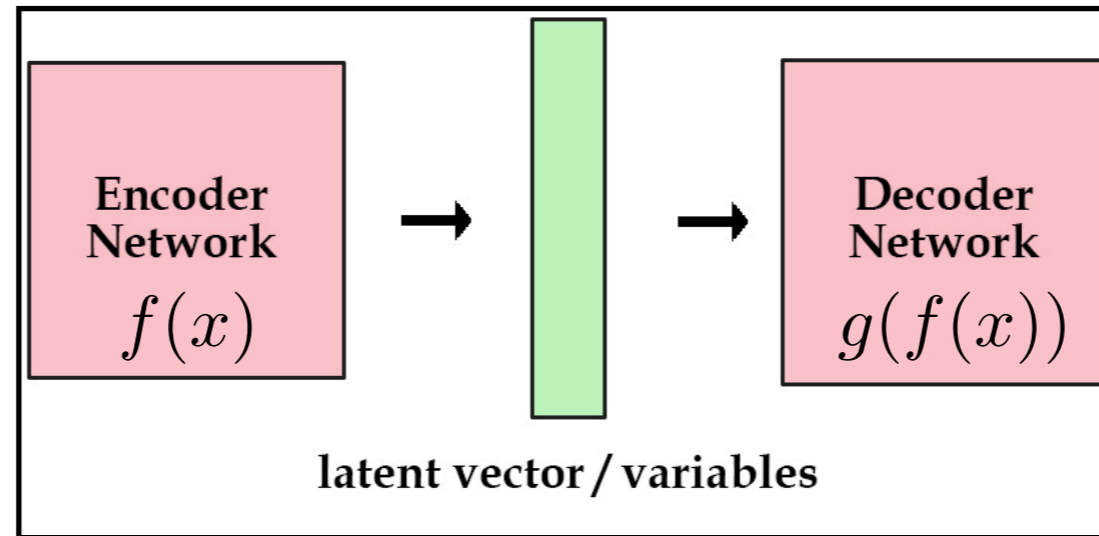
Improves training stability and sample quality (e.g. mode collapse)



# **Variational Autoencoders**



# Autoencoder



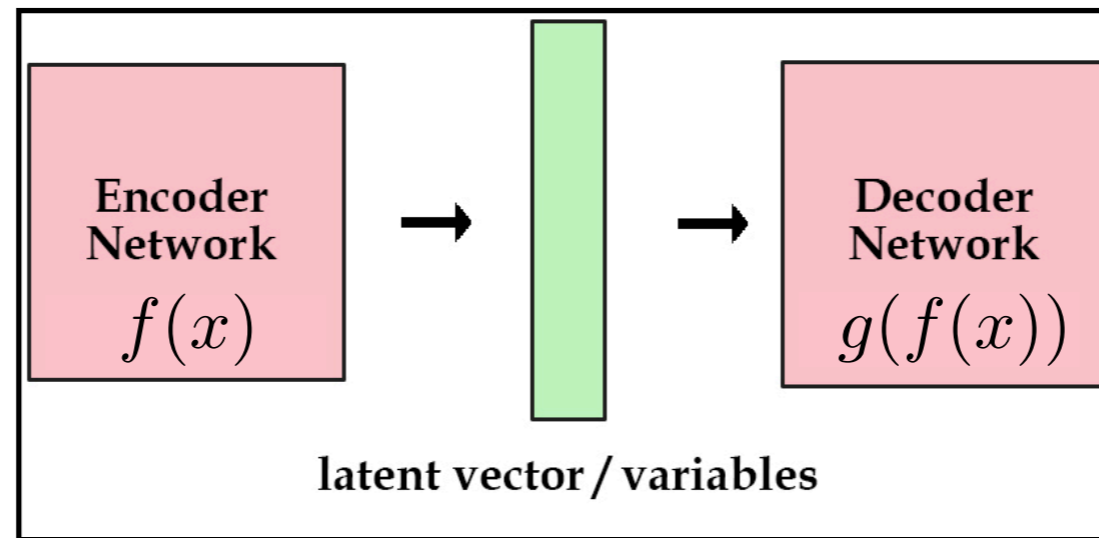
Two networks

**Encoder:** data  $\rightarrow$  latent space

**Decoder:** latent space  $\rightarrow$  data



# Autoencoder



Two networks

**Encoder:** data  $\rightarrow$  latent space

**Decoder:** latent space  $\rightarrow$  data

Training objective:

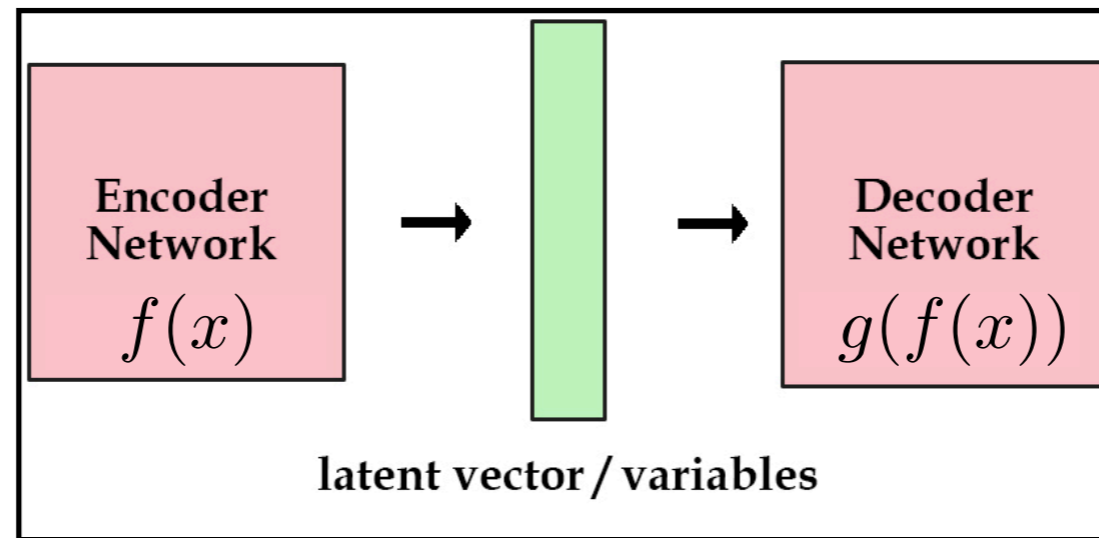
**Minimise input/output difference**

$$L = (x - f(g(x)))^2$$

Decoder Encoder



# Autoencoder



Two networks

Encoder: data  $\rightarrow$  latent space

Decoder: latent space  $\rightarrow$  data

Training objective:

Minimise input/output difference

$$L = (x - f(g(x)))^2$$

Decoder Encoder

**Uses:**

Dimension reduction

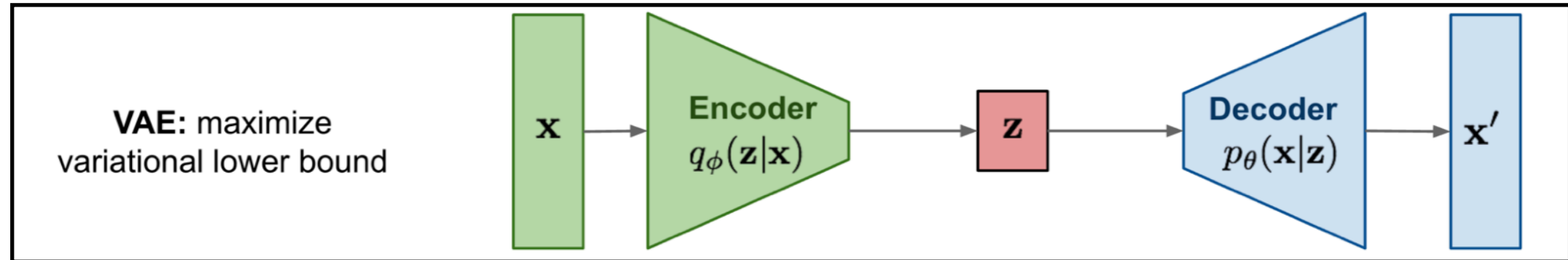
Denoising

Anomaly detection

Generation?



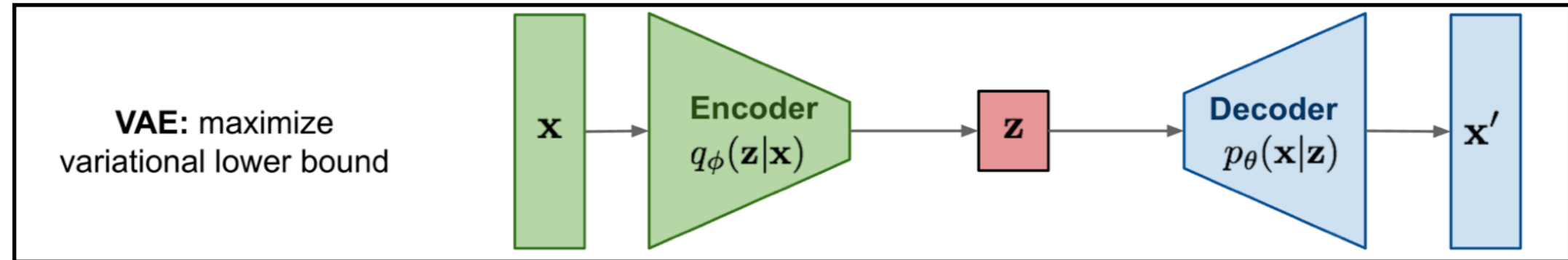
# Variational Autoencoder



Variational Autoencoder (VAE):  
Split latent space

$$f(x) = (\mu, \sigma)$$

# Variational Autoencoder



Variational Autoencoder (VAE):  
Split latent space  
Sample before decoder

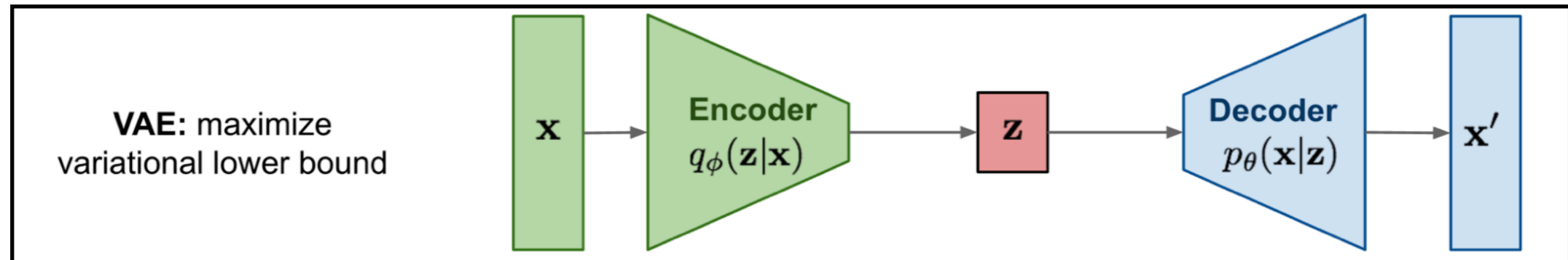
$$f(x) = (\mu, \sigma)$$

$$z = \text{Gaussian}(\mu, \sigma)$$

$$x' = g(z)$$



# Variational Autoencoder



Variational Autoencoder (VAE):

Split latent space

Sample before decoder

Penalty so mean/std are close to unit Gaussian

$$f(x) = (\mu, \sigma)$$

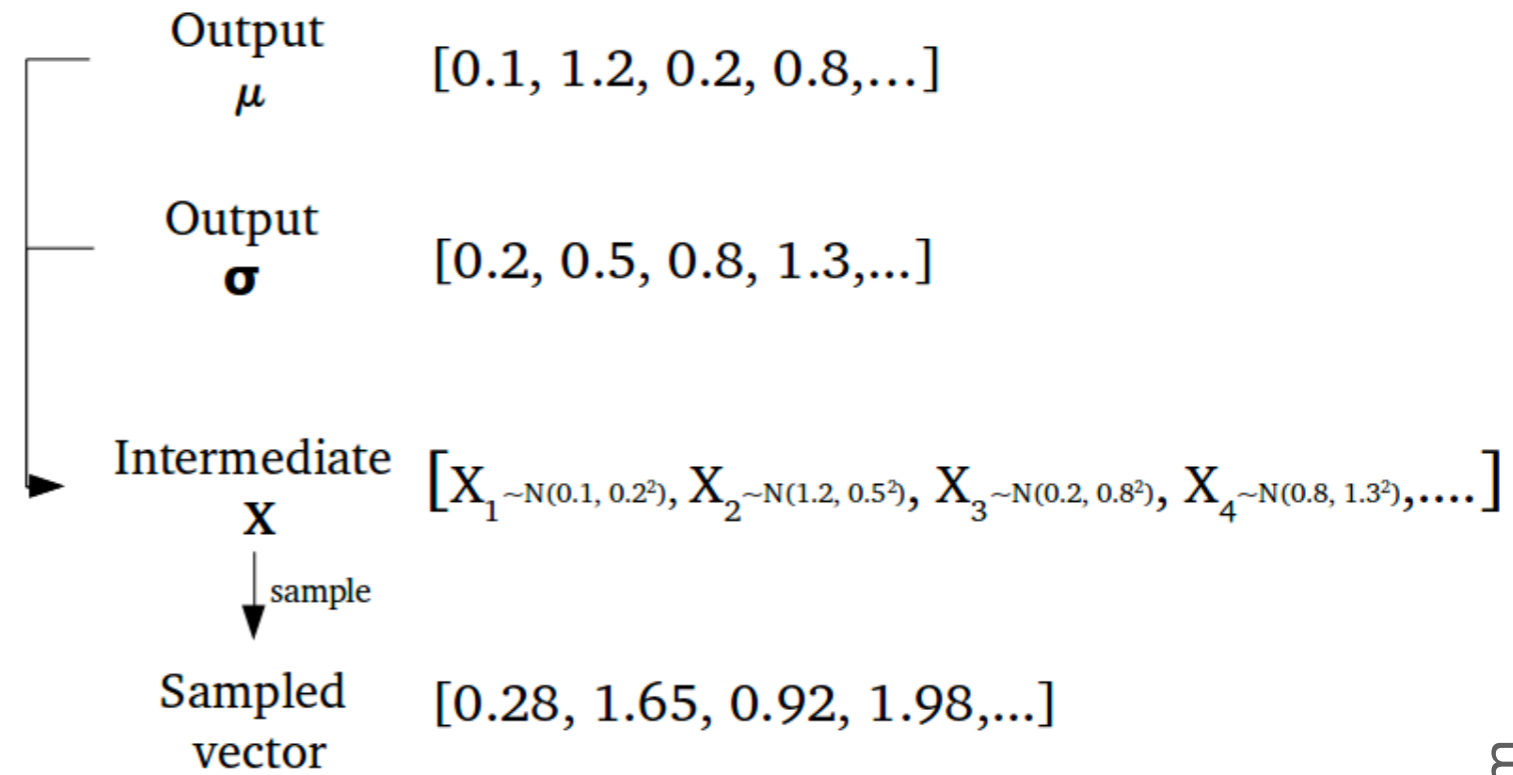
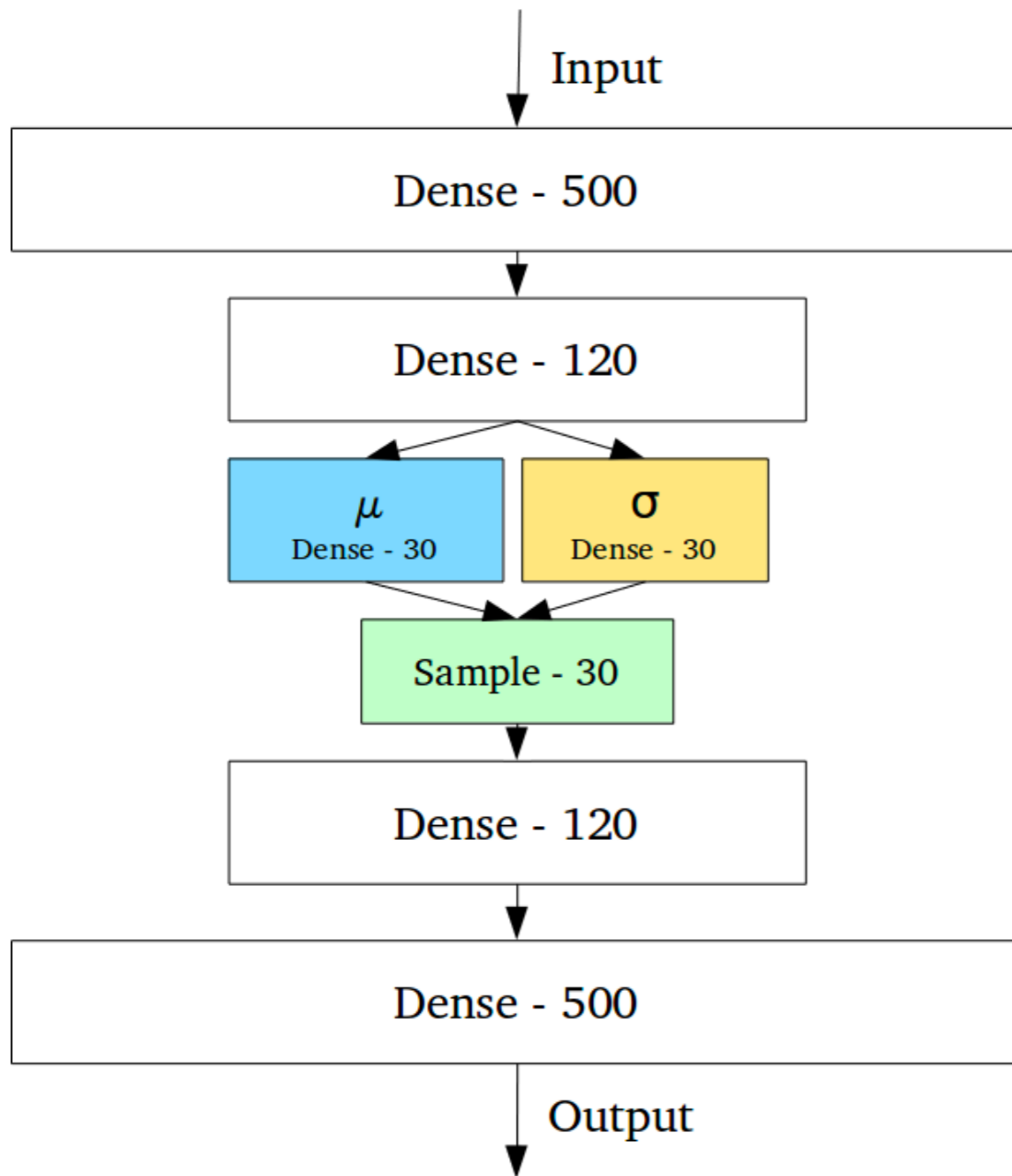
$$z = \text{Gaussian}(\mu, \sigma)$$

$$x' = g(z)$$

$$L = (x - g(z))^2 + \sigma^2 + \mu^2 - \log(\sigma) - 1$$

(Calculate KL-divergence between Gaussians)

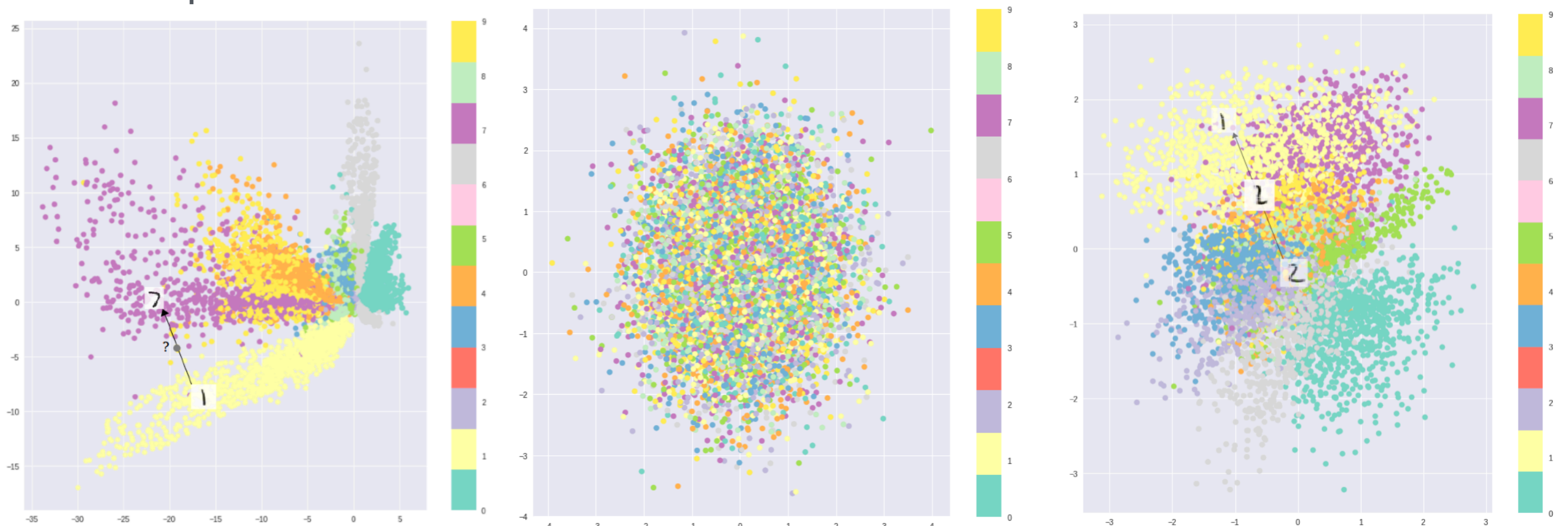
# VAE Example





# Loss terms

## Latent space of MNIST VAE



$$(x - g(z))^2$$

Reconstruction

$$\sigma^2 + \mu^2 - \log(\sigma) - 1$$

Regularisation

Both terms

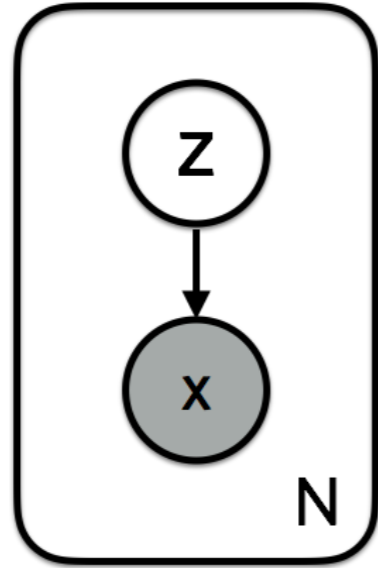
# Loss terms

$$L = (x - g(z))^2 + \sigma^2 + \mu^2 - \log(\sigma) - 1$$

How did we get here?



# Loss terms



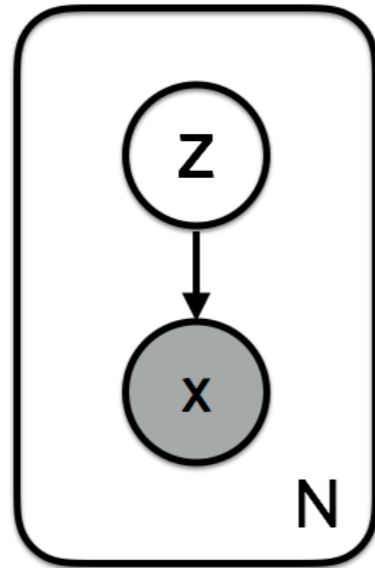
Sample from latent variables  $z$

$$z_i \sim p(z)$$

Produce data points  $x$

$$x_i \sim p(x | z)$$

# Loss terms



Sample from latent variables  $z$

$$z_i \sim p(z)$$

Produce data points  $x$

$$x_i \sim p(x | z)$$

To choose correct latent distribution given data, could use **Bayes theorem**:

$$p(z | x) = \frac{p(x | z)p(z)}{p(x)}$$

Conditional      Prior

Evidence

Difficult due to  $p(x)$



# Loss terms

To choose correct latent distribution given data, could use Bayes theorem:

$$p(z | x) = \frac{p(x | z)p(z)}{p(x)}$$

Instead, approximate with family of posterior distributions (variational inference):

$$\begin{aligned} \text{KL}(q_\lambda(z | x) || p(z | x)) = \\ \mathbf{E}_q[\log q_\lambda(z | x)] - \mathbf{E}_q[\log p(x, z)] + \log p(x) \end{aligned}$$

And find optimal approximation:

$$q_\lambda^*(z | x) = \arg \min_\lambda \text{KL}(q_\lambda(z | x) || p(z | x))$$

Still difficult due to (hidden)  $p(x)$  term!

# Loss terms

To choose correct latent distribution given data, could use Bayes theorem:

$$p(z | x) = \frac{p(x | z)p(z)}{p(x)}$$

$$\mathbb{KL}(q(x) || p(x)) = - \int dx q(x) \log \frac{p(x)}{q(x)}$$

Kullback-Leibler definition

Instead, approximate with family of posterior distributions (variational inference):

$$p(x, z) = p(z|x)p(x)$$

$$\mathbb{KL}(q_\lambda(z | x) || p(z | x)) = \text{Reminder}$$

$$\mathbf{E}_q[\log q_\lambda(z | x)] - \mathbf{E}_q[\log p(x, z)] + \log p(x)$$

And find optimal approximation:

$$q_\lambda^*(z | x) = \arg \min_\lambda \mathbb{KL}(q_\lambda(z | x) || p(z | x))$$

Still difficult due to  $p(x)$  term!



# Loss terms

$$\mathbb{KL}(q_\lambda(z | x) || p(z | x)) = \mathbf{E}_q[\log q_\lambda(z | x)] - \mathbf{E}_q[\log p(x, z)] + \log p(x)$$

Introduce

$$ELBO(\lambda) = \mathbf{E}_q[\log p(x, z)] - \mathbf{E}_q[\log q_\lambda(z | x)]$$

Rewrite

$$\log p(x) = ELBO(\lambda) + \mathbb{KL}(q_\lambda(z | x) || p(z | x))$$

As KL is  $\geq 0$ , ELBO is a lower limit for  $p(X)$   
ELBO: Evidence Lower Bound

# Loss terms

Maximise

$$ELBO(\lambda) = \mathbf{E}_q[\log p(x, z)] - \mathbf{E}_q[\log q_\lambda(z | x)]$$

Rewrite for samples, using neural networks:

$$ELBO_i(\theta, \phi) = \mathbb{E}_{q_\theta(z | x_i)}[\log p_\phi(x_i | z)] - \mathbb{KL}(q_\theta(z | x_i) || p(z))$$

Reconstruction term

Regularisation term

Assume normal distribution

Difference between normal and standard normal

$$L = (x - g(z))^2 + \sigma^2 + \mu^2 - \log(\sigma) - 1$$



# Loss terms

Maximise

$$ELBO(\lambda) = \mathbf{E}_q[\log p(x, z)] - \mathbf{E}_q[\log q_\lambda(z | x)]$$

Rewrite for samples, using neural networks:

$$ELBO_i(\theta, \phi) = \mathbb{E}_{q_\theta(z | x_i)}[\log p_\phi(x_i | z)] - \mathbb{KL}(q_\theta(z | x_i) || p(z))$$

Reconstruction term

Regularisation term

Assume normal distribution

Difference between normal and standard normal

$$L = (x - g(z))^2 + \sigma^2 + \mu^2 - \log(\sigma) - 1$$

# Comments on VAEs

## Architecture:

- Low complexity, fast and adaptable
- Target: Maximise lower bound on likelihood

## Learning:

- **Stable** training
- Average prediction → blurrier output
- Interpretable latent space



## Maturity:

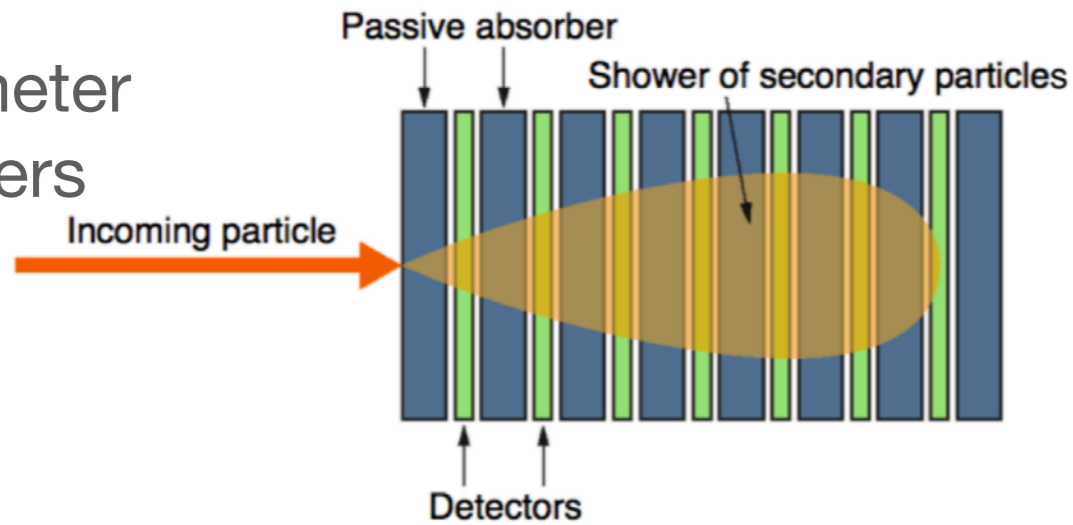
- Well established,  
many variants and extensions



# Applications I

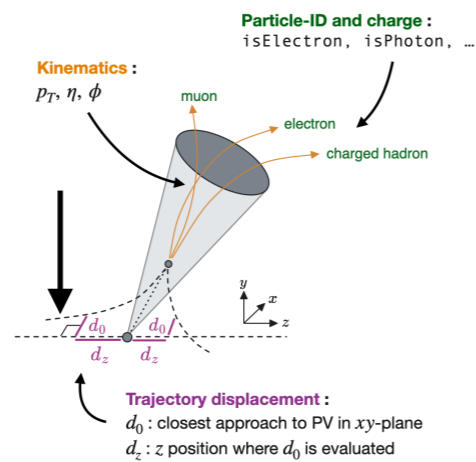
# (Some) Simulation targets

Calorimeter  
Showers



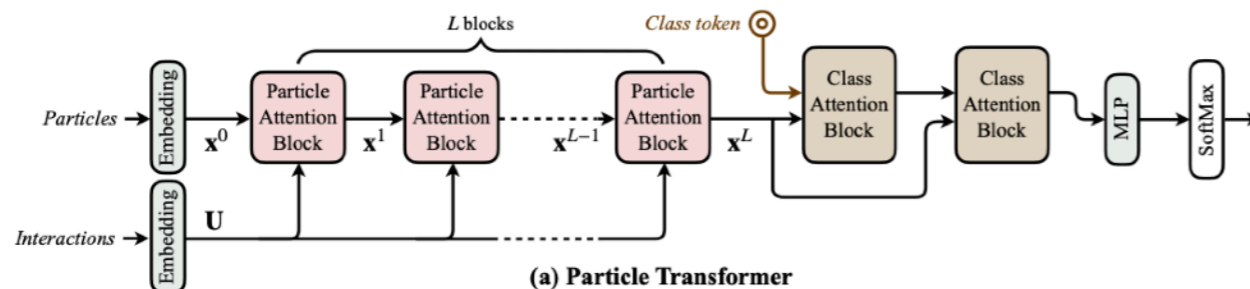
Reduce computational bottleneck

Jet Constituents



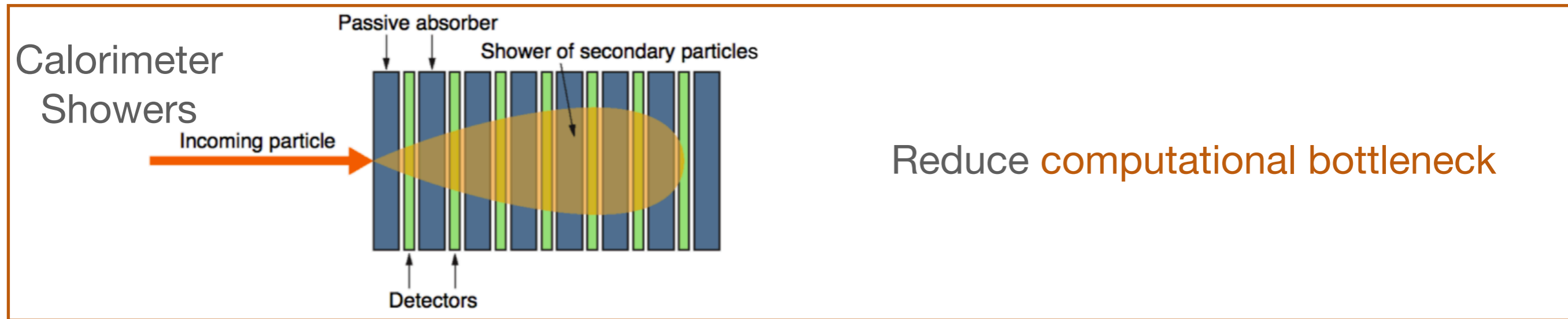
Predict background from data

Classification and  
Reconstruction tasks

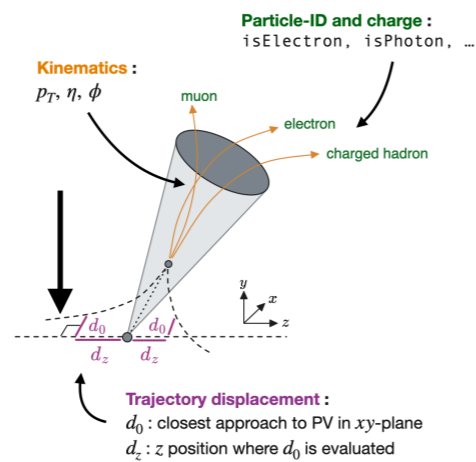


Act as surrogate models

# (Some) Simulation targets

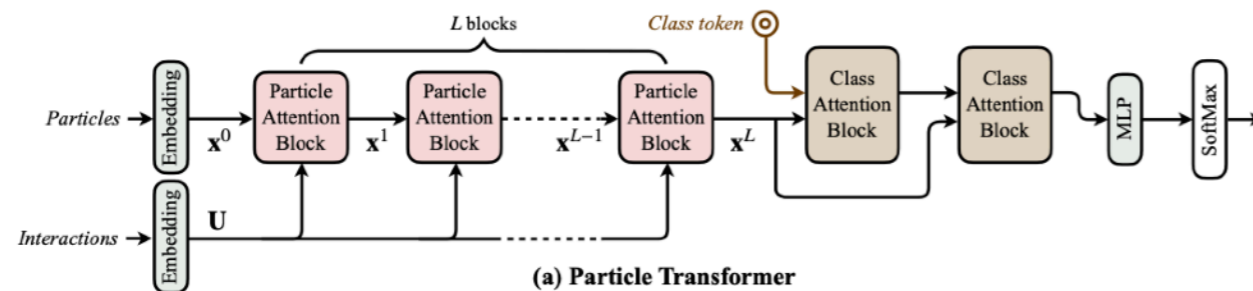


## Jet Constituents



Predict background from data

## Classification and Reconstruction tasks



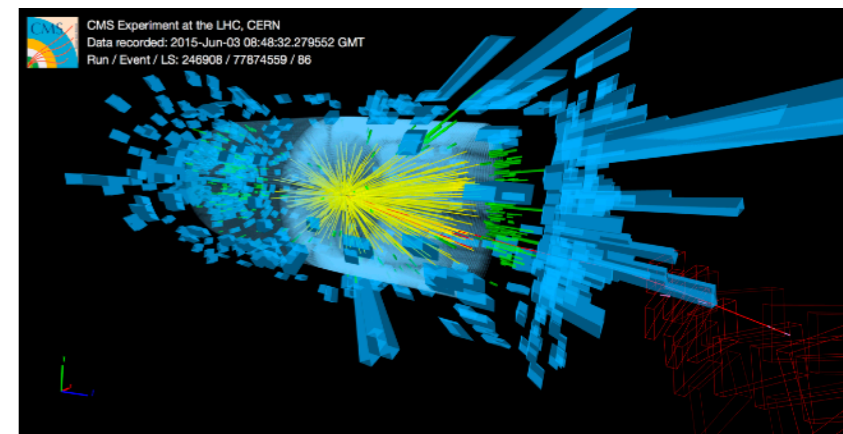
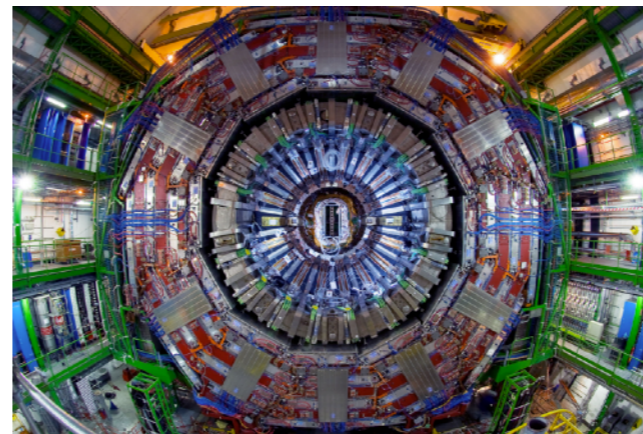
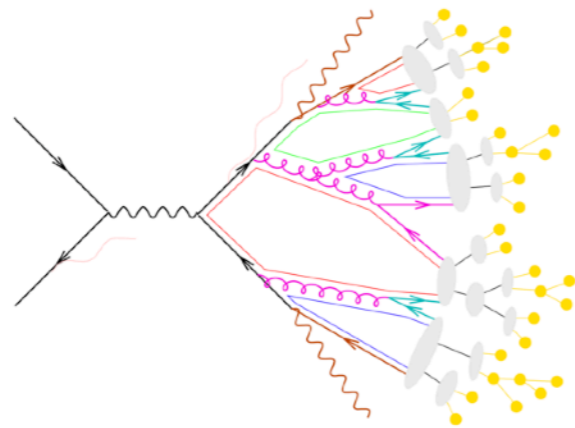
Act as surrogate models



# Generative Models

This happens in the experiment

$$\begin{aligned} \mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\ & + i\bar{\psi}\not{D}\psi + h.c. \\ & + \chi_i Y_{ij} \chi_j \phi + h.c. \\ & + |D_\mu\phi|^2 - V(\phi) \end{aligned}$$



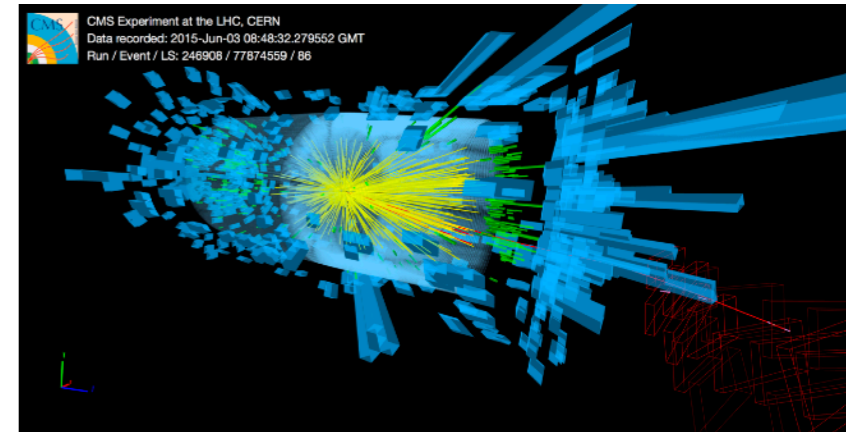
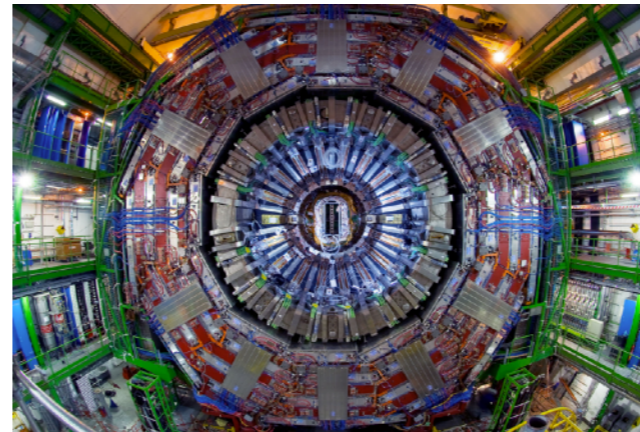
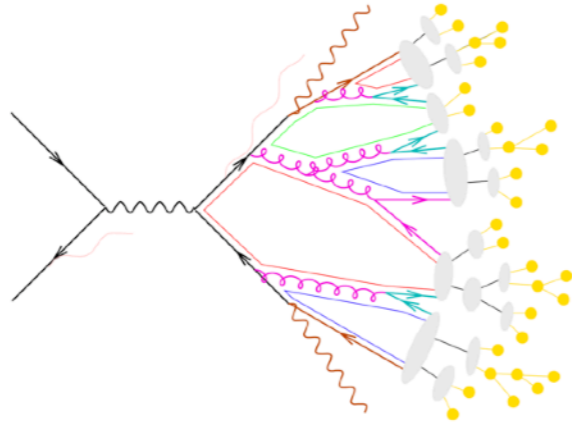
This is what we want to know

Simulation is crucial to connect  
experimental data with theory  
predictions

# Generative Models

This happens in the experiment

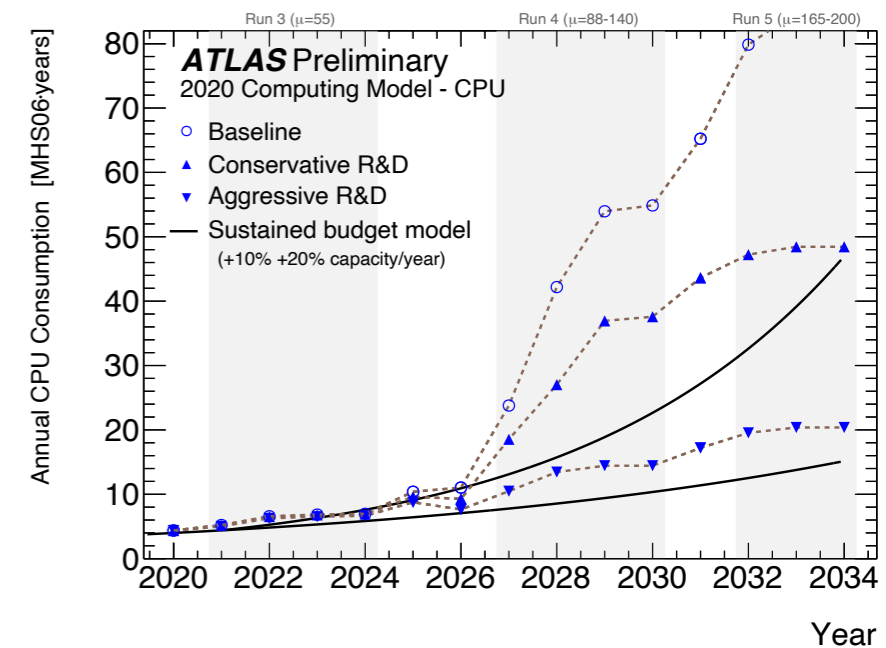
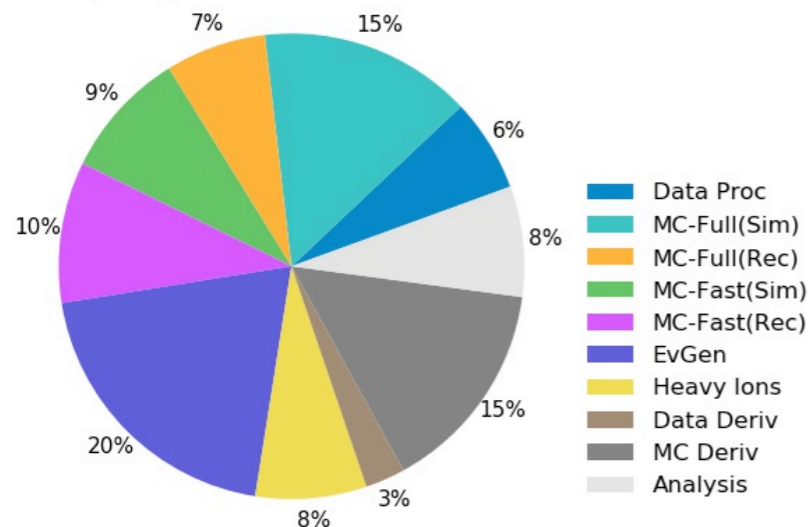
$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} + i\bar{\psi}\not{D}\psi + h.c. + \chi_i Y_{ij} \chi_j \phi + h.c. + |D_\mu \phi|^2 - V(\phi)$$



This is what we want to know

Simulation is crucial to connect experimental data with theory predictions, **but computationally very costly**

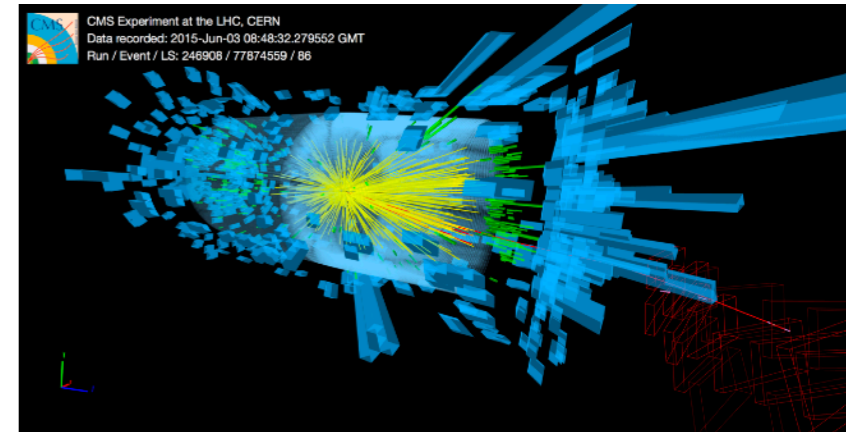
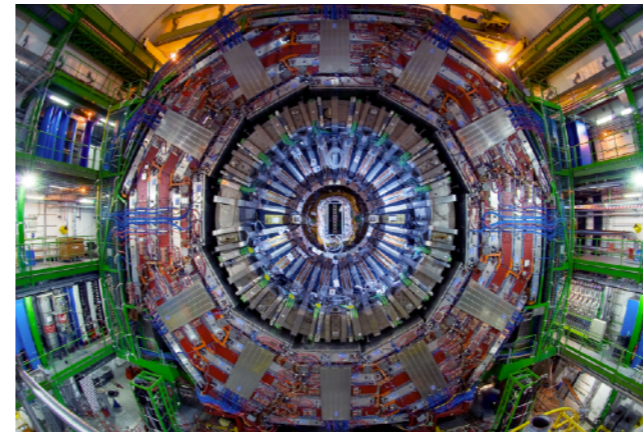
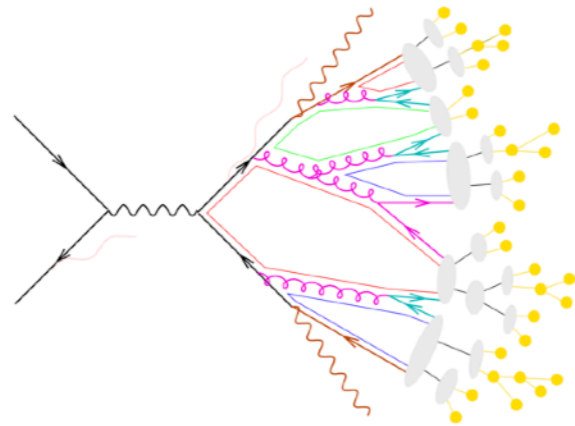
**ATLAS Preliminary**  
2020 Computing Model - CPU: 2030: Baseline



# Generative Models

This happens in the experiment

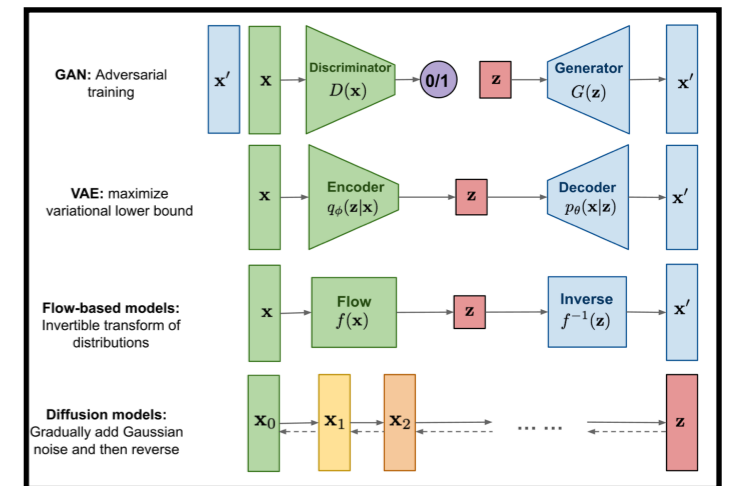
$$\begin{aligned} \mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\ & + i\bar{\psi} \not{D} \psi + h.c. \\ & + \chi_i Y_{ij} \chi_j \phi + h.c. \\ & + |D_m \phi|^2 - V(\phi) \end{aligned}$$



This is what we want to know

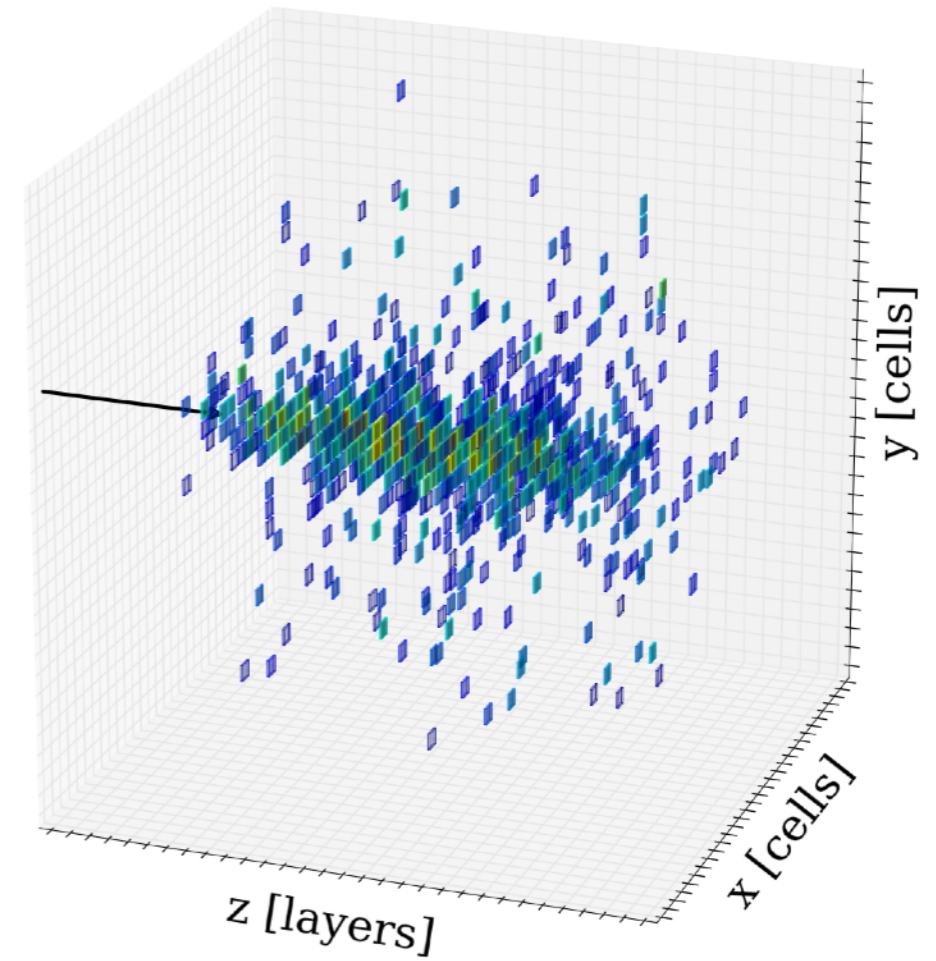
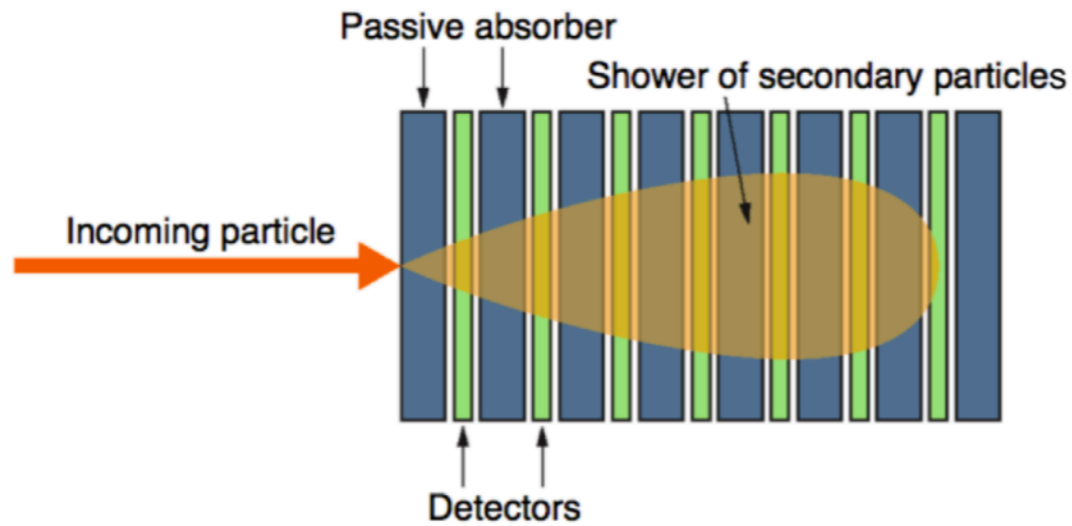
Simulation is crucial to connect experimental data with theory predictions, but computationally very costly

→ Use generative models trained on simulation or data to augment simulations



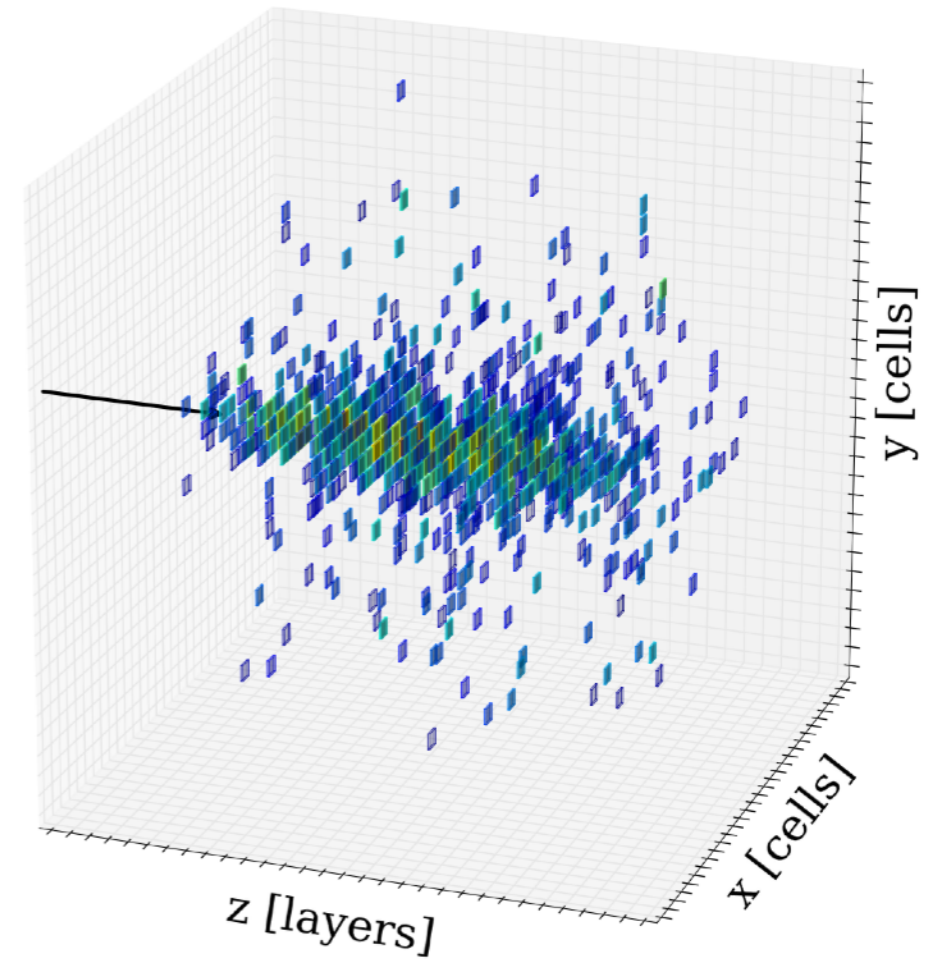
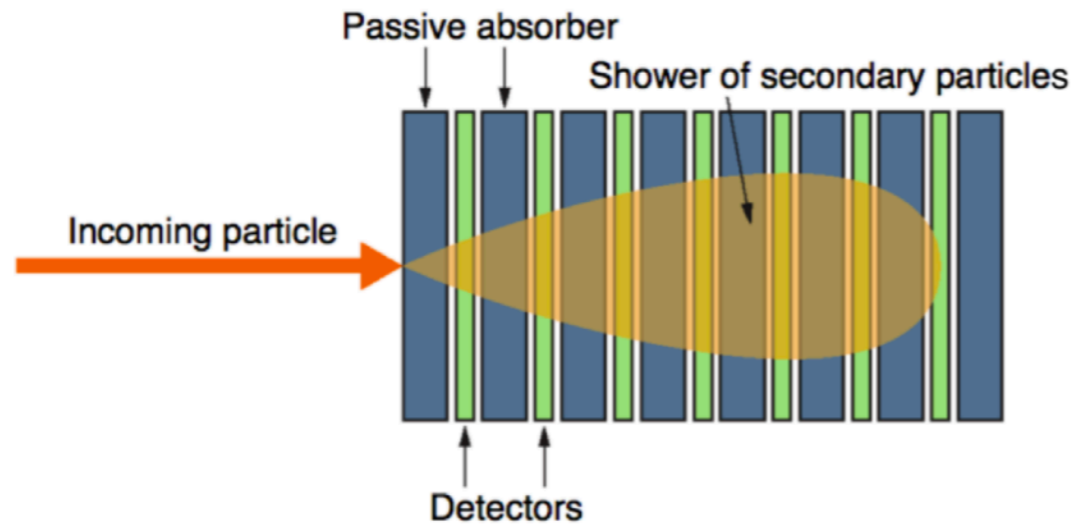


# Simulation targets



How to represent?

# Simulation targets

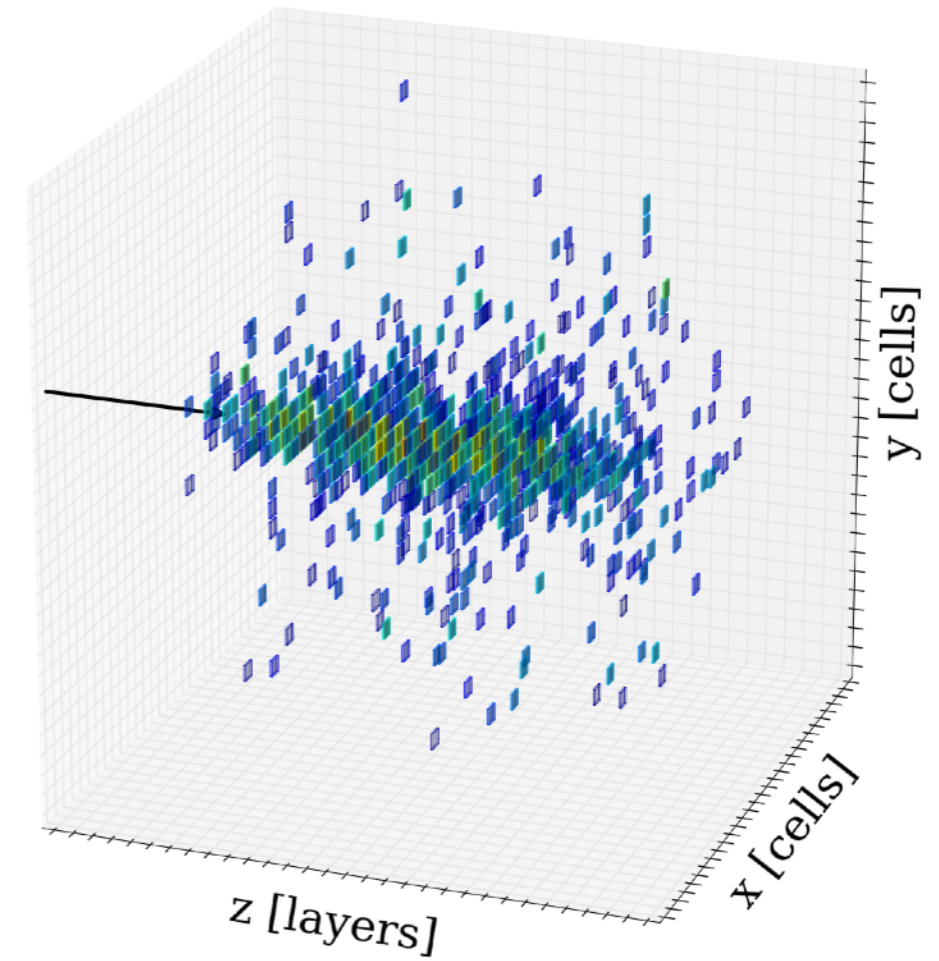
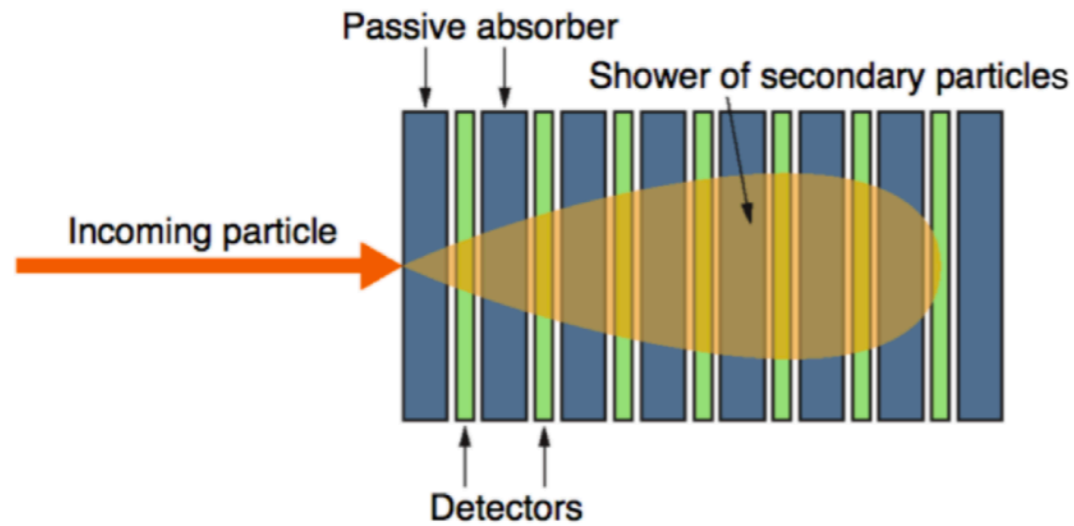


How to represent?

**Tabular data:**

Easy, insufficient for high-dimensions

# Simulation targets



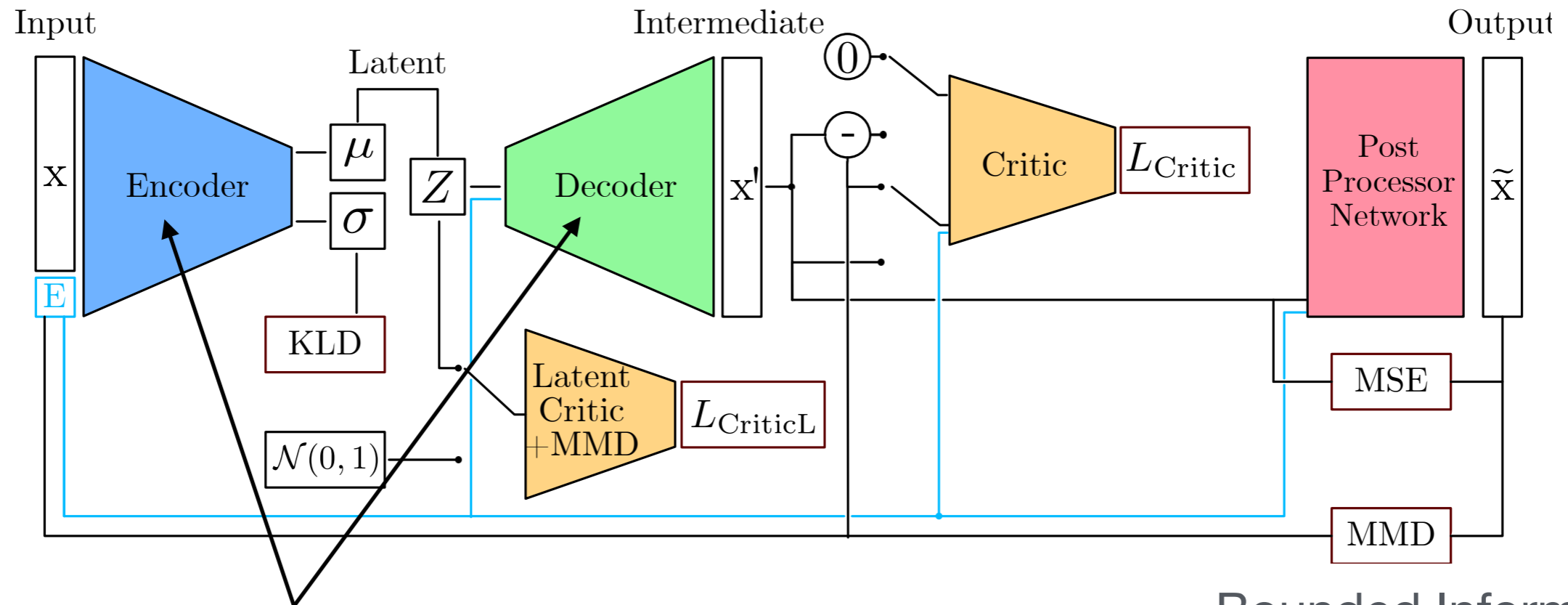
How to represent?

Tabular data

Fixed grid (voxels)



# Generative results



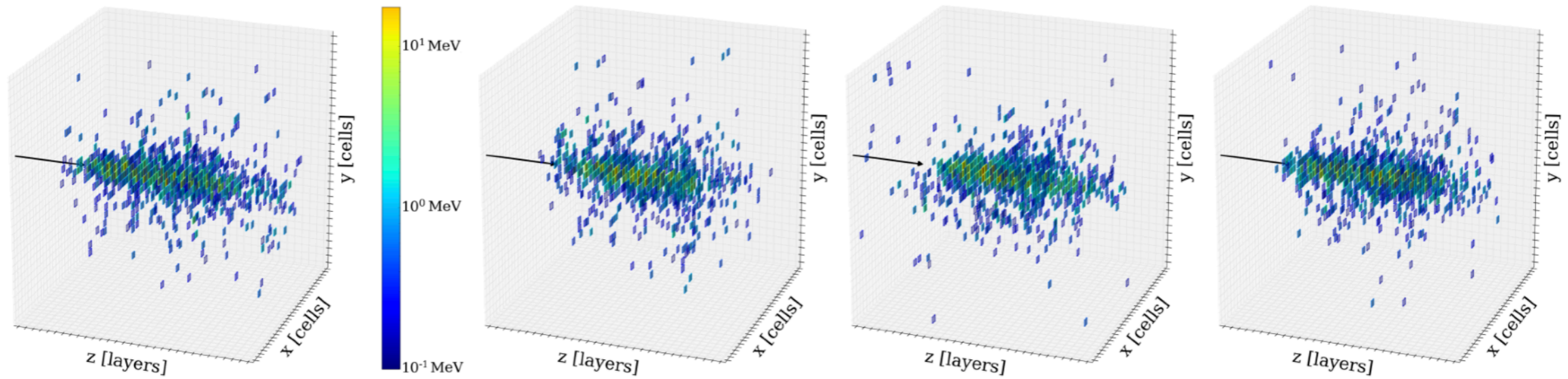
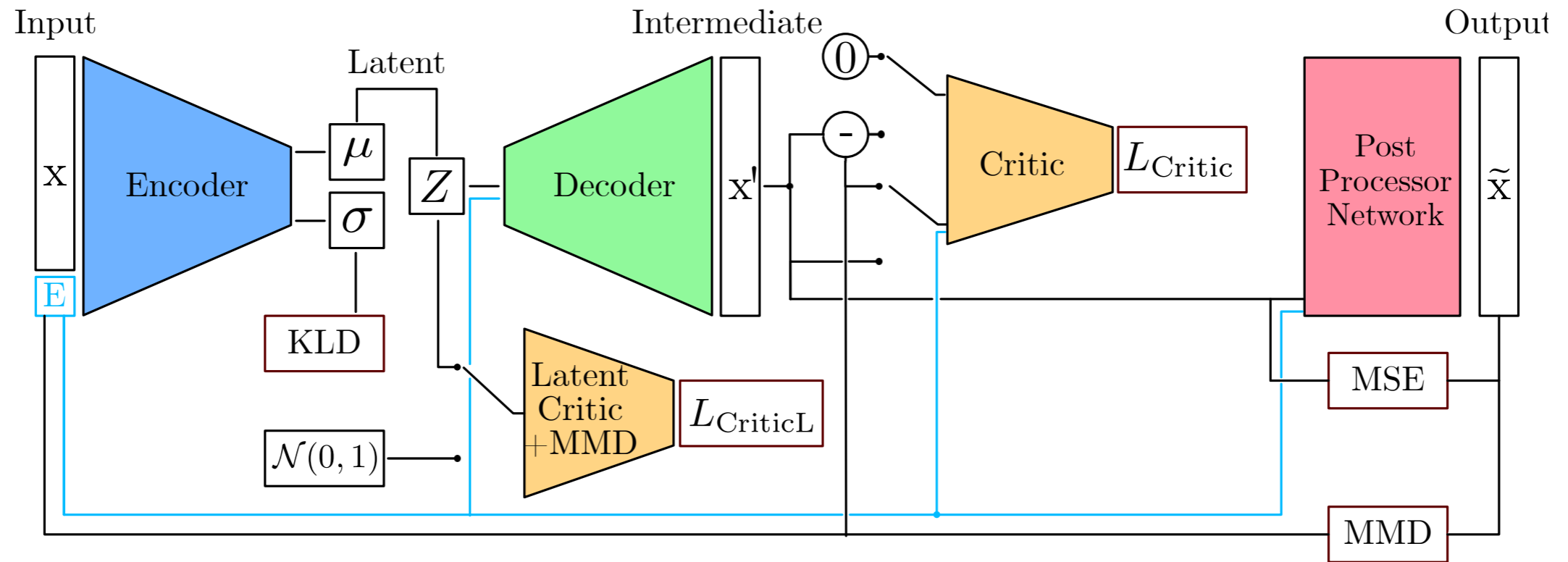
(Transposed) Convolution

Bounded Information Bottleneck AE

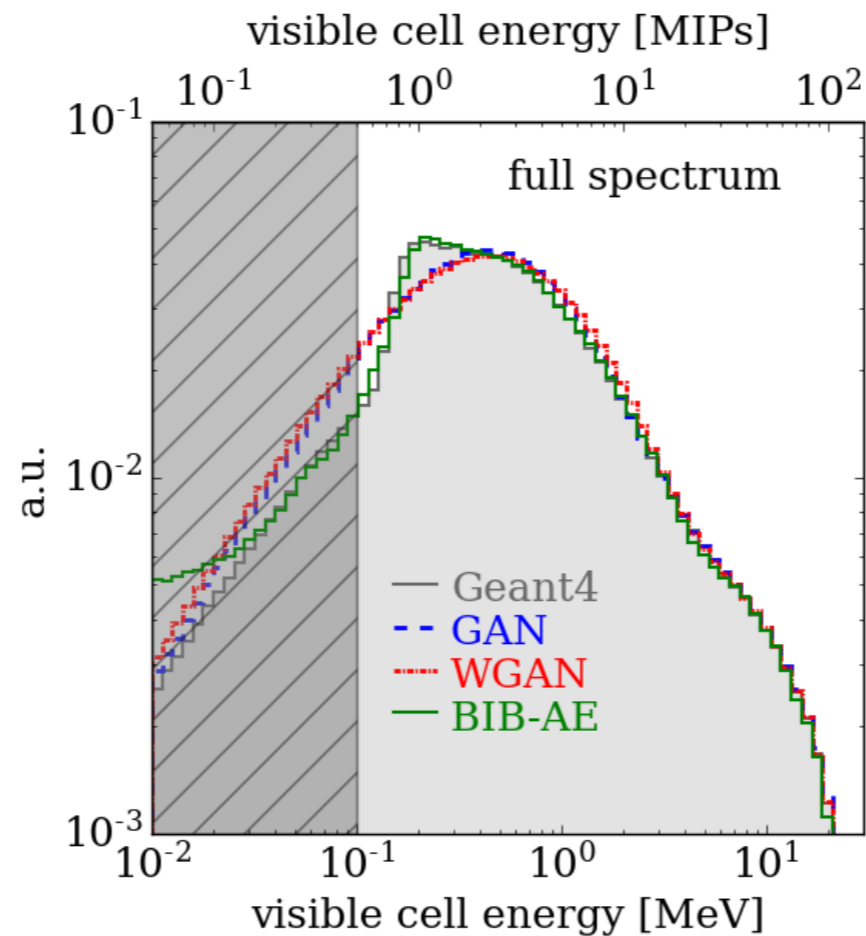
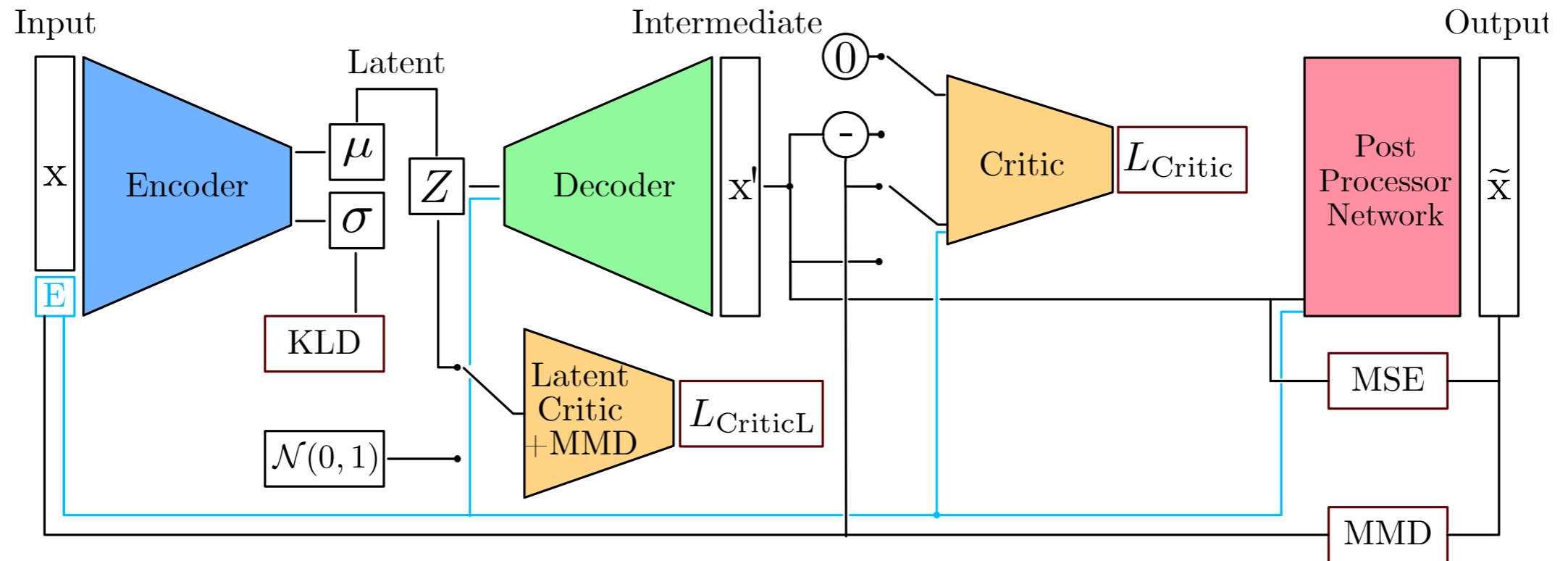
**BIB-AE (GAN + VAE)**

$$\begin{aligned}
 L_{\text{BIB-AE}} = & -\beta_{C_L} \cdot \mathbb{E}[C_L(N_E(x))] && \text{Latent Critic} \\
 & -\beta_C \cdot \mathbb{E}[C_E(D_E(N_E(x)))] && \text{Critic} \\
 & -\beta_{C_D} \cdot \mathbb{E}[C_{D,E}(D_E(N_E(x)) - x)] && \text{Difference Critic} \\
 & + \beta_{\text{KLD}} \cdot \text{KLD}(N_E(x)) && \text{Latent Regularisation} \\
 & + \beta_{\text{MMD}} \cdot \text{MMD}(N_E(x), \mathcal{N}(0, 1)).
 \end{aligned}$$

# Generative results



# Generative results



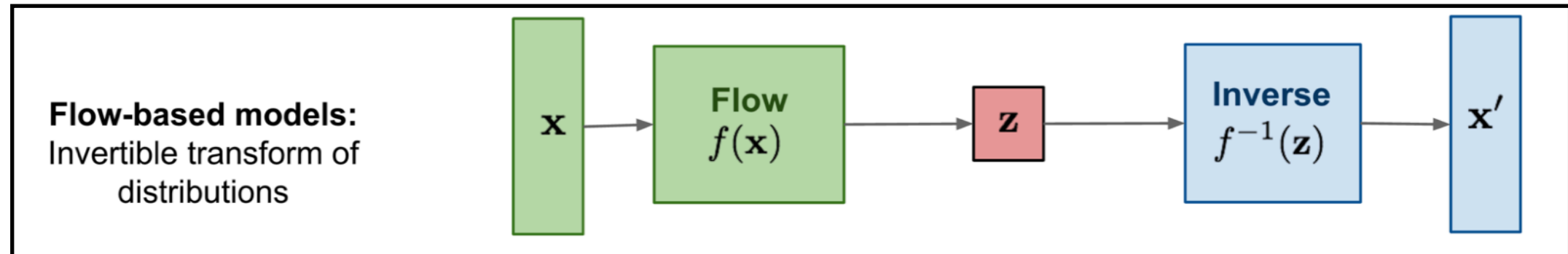


**Go with the...**

**(Normalising) Flows**



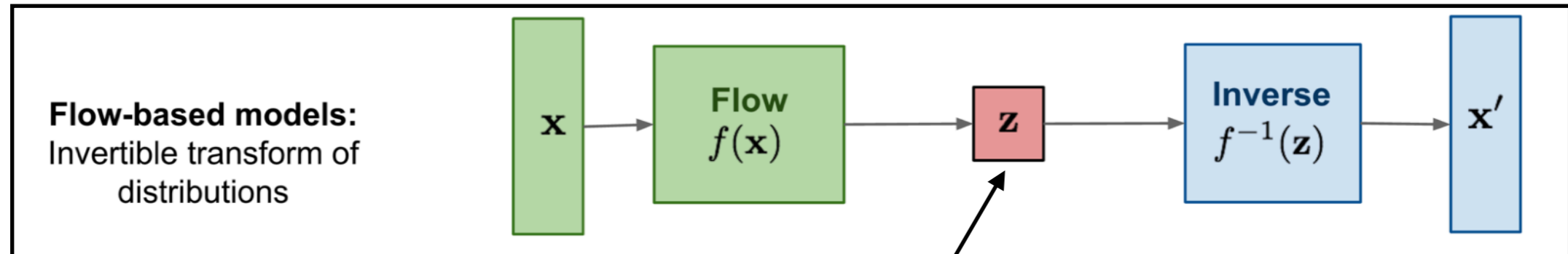
# Generative models



In auto-encoders, the decoder learns to ‘undo’ the encoder

Can we make this exact?

# Generative models

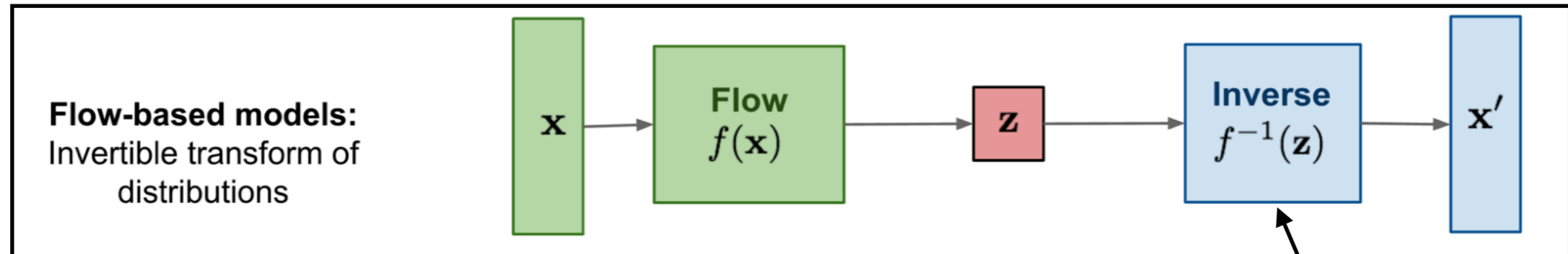


Choose latent space, e.g. standard normal distribution (normalising flow!)  
Same dimension as data!

Learn a diffeomorphism between data and latent-space



# Generative models

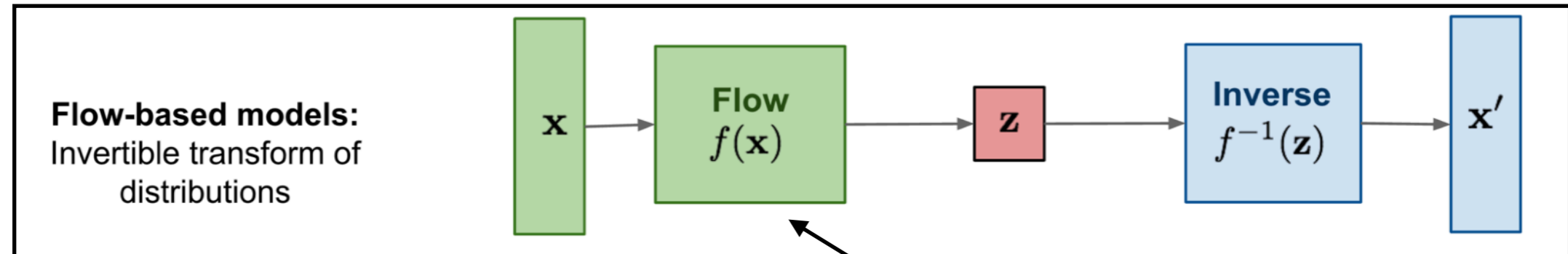


$f^{-1}$  is not a learned inversion, but exact inverse by construction

Learn a diffeomorphism between data and latent-space

Bijjective, invertable

# Generative models



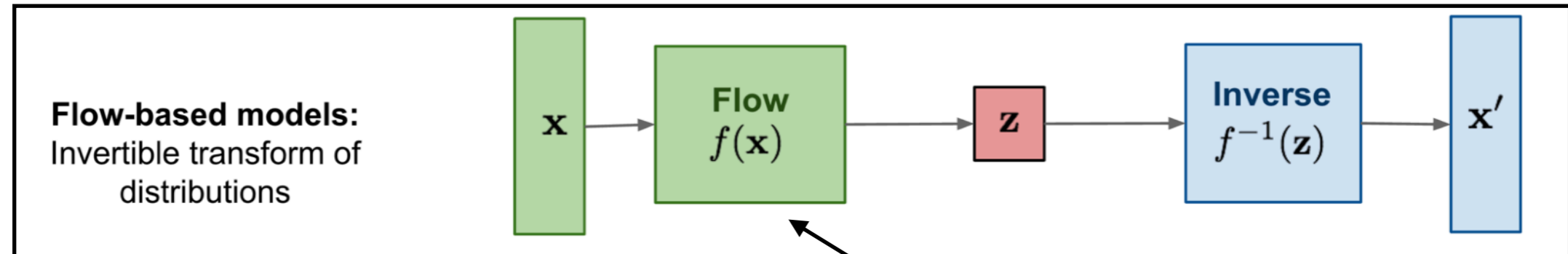
Learn a diffeomorphism between data and latent-space

Bijjective, invertable

Learn likelihood of data

Take into account Jacobian determinant to evaluate probability density

# Generative models



2 challenges:

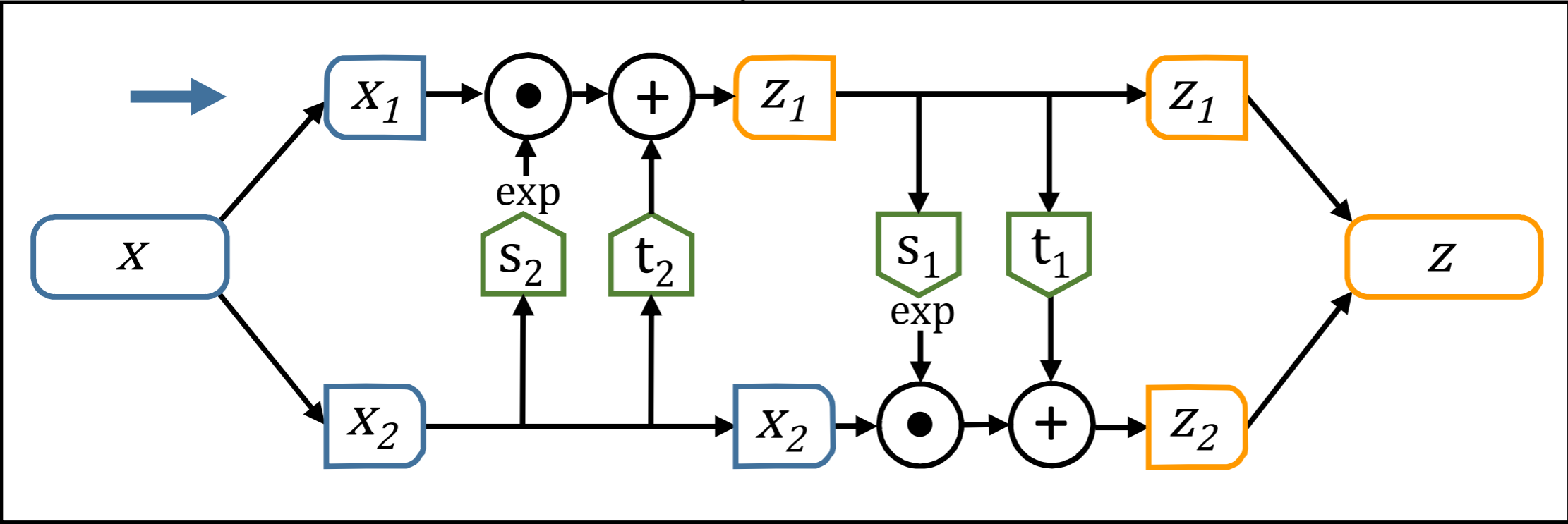
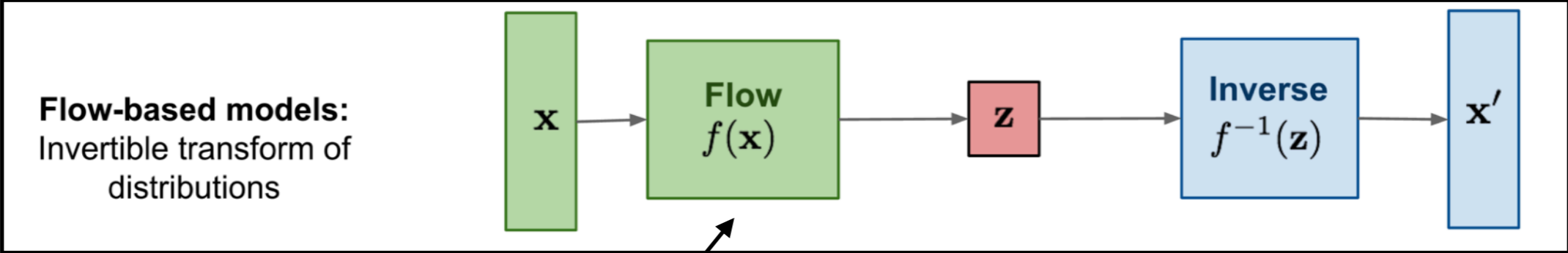
**Invertible**

Easy-to-calculate Jacobean

Take into account Jacobian determinant to evaluate probability density

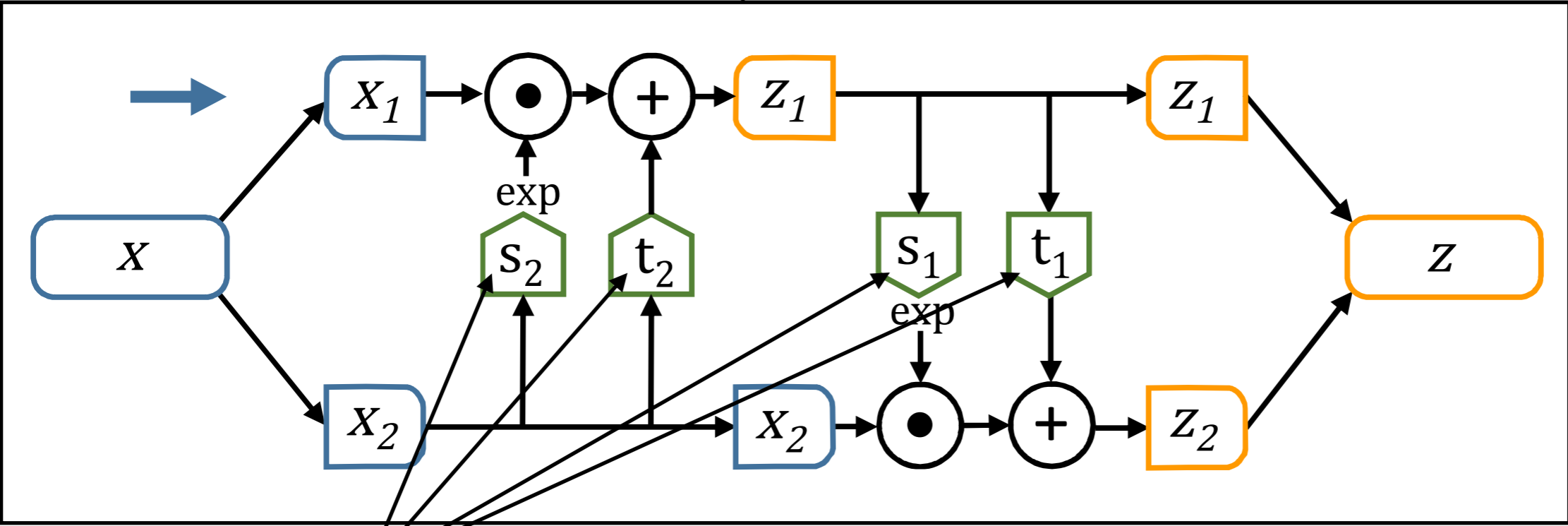
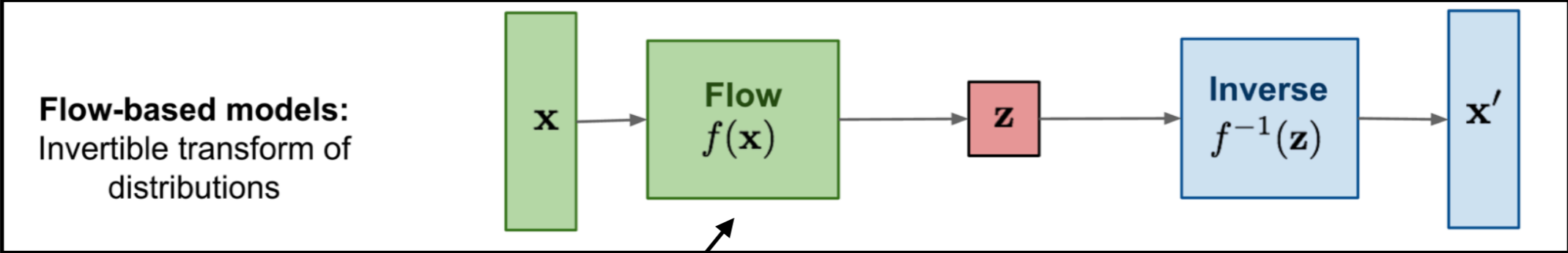


# Coupling flows



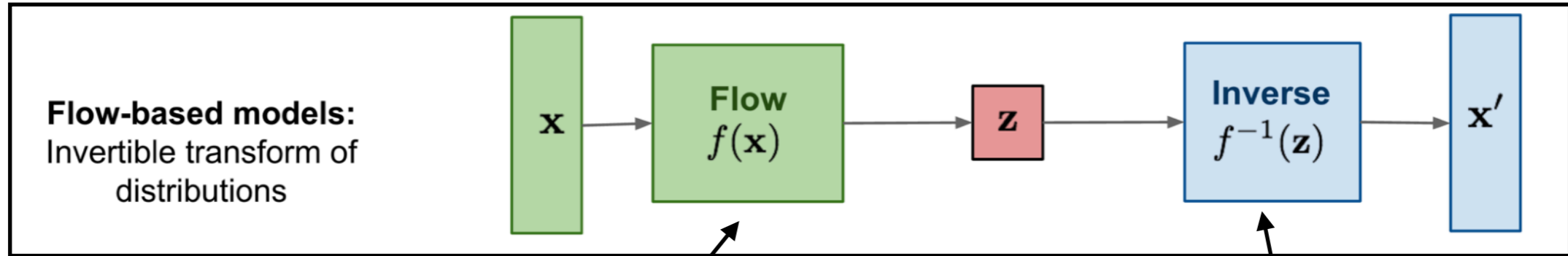
Coupling layers: Not the most expressive, but useful for illustration/understanding

# Coupling flows

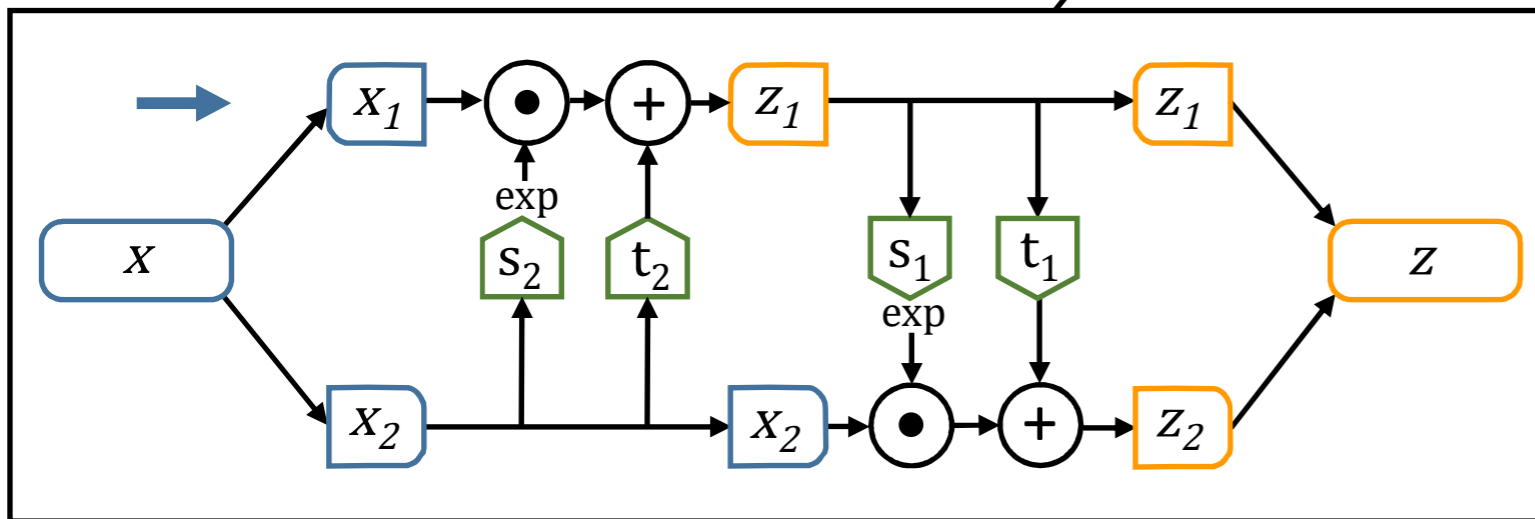


Simple (e.g. dense) neural networks

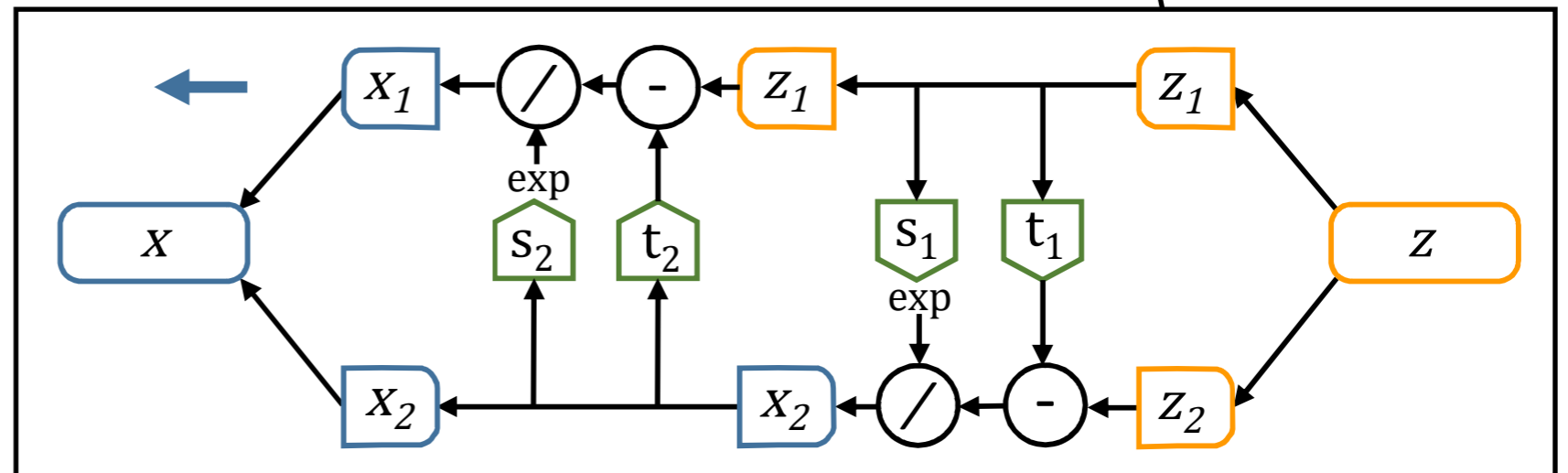
# Coupling flows



Forward direction

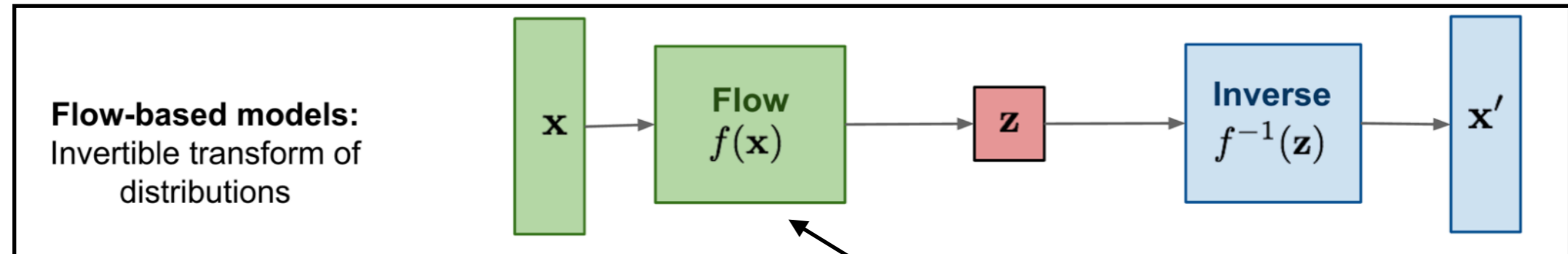


Inverse direction





# Generative models



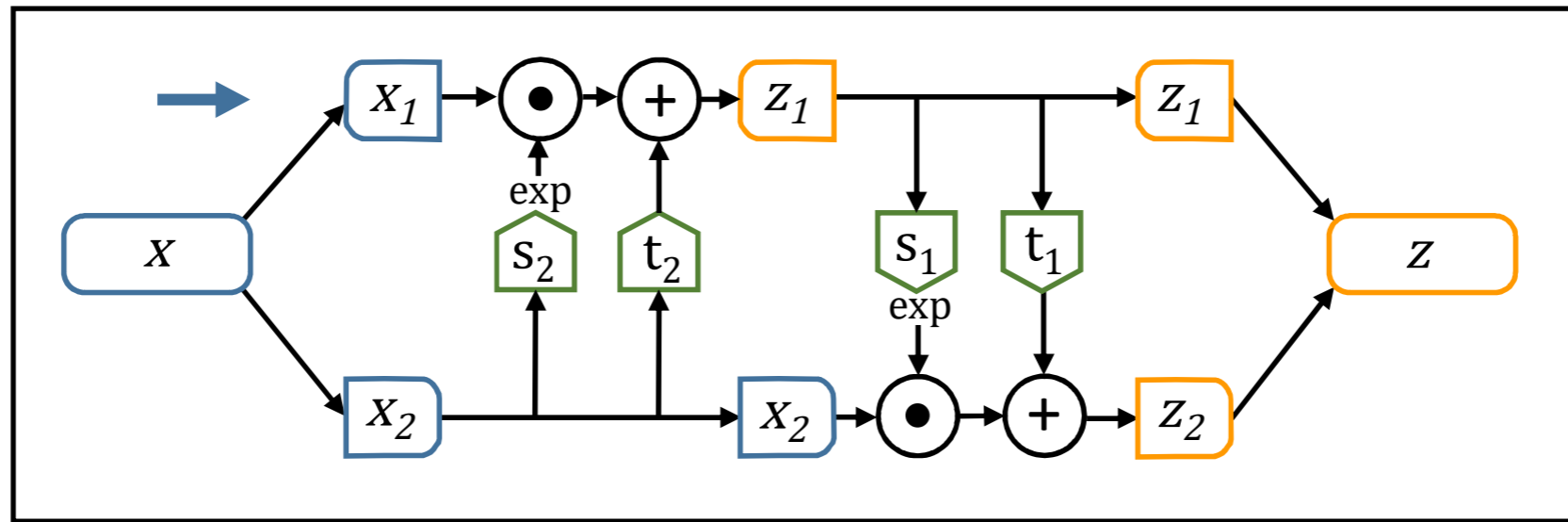
2 challenges:

Invertible

Easy-to-calculate **Jacobian**

Take into account Jacobian determinant to evaluate probability density

# Calculating Jacobian determinant



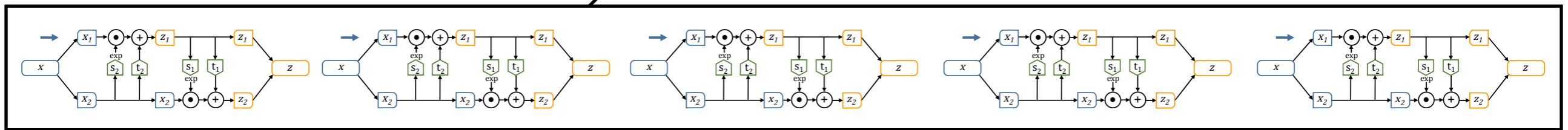
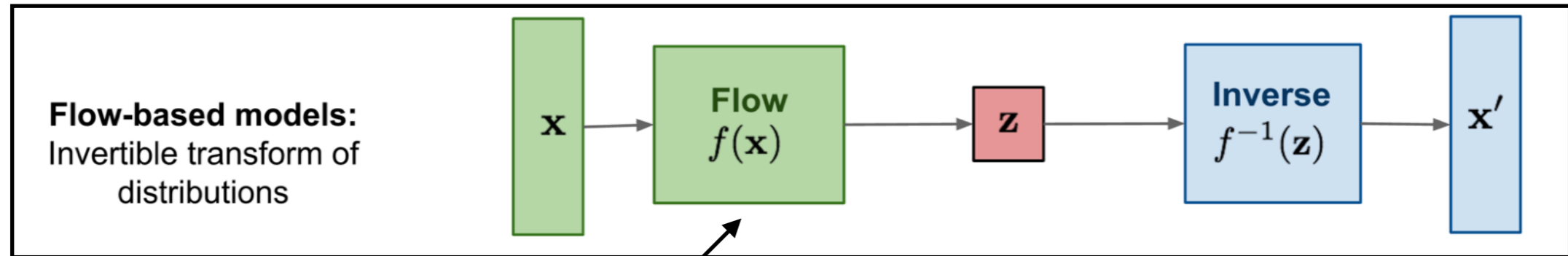
$$\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \xrightarrow{f_1} \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{x}_2 \end{pmatrix} \xrightarrow{f_2} \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{pmatrix} \text{ with } \begin{array}{l} \mathbf{x}_1 \xrightarrow{f_1} \mathbf{z}_1 = \mathbf{x}_1 \odot \exp(s_2(\mathbf{x}_2)) + t_2(\mathbf{x}_2) \\ \mathbf{x}_2 \xrightarrow{f_1} \mathbf{x}_2. \end{array}$$

$$\mathbf{J}_1 = \begin{pmatrix} \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}_2} \\ \frac{\partial \mathbf{z}_2}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{z}_2}{\partial \mathbf{x}_2} \end{pmatrix} = \begin{pmatrix} \text{diag}(\exp(s_2(\mathbf{x}_2))) & \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}_2} \\ 0 & \mathbb{1} \end{pmatrix}$$

Triangular by construction

$$\det \mathbf{J}_1 = \prod \exp(s_2(\mathbf{x}_2)) = \exp \left( \sum s_2(\mathbf{x}_2) \right)$$

# Composition



Composition of bijective functions remains bijective

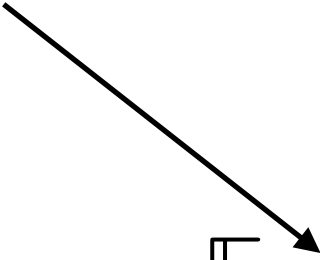
Chain rule: Jacobian determinant of composition is product of determinants



# How to train NF?

Training objective: Minimise negative log likelihood of data

Sample points from training data


$$\mathcal{L} = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[ -\frac{1}{2} \|f(\mathbf{x})\|_2^2 + \sum s(\mathbf{x}) \right]$$

# How to train NF?

Training objective: Minimise negative log likelihood of data

$$\mathcal{L} = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[ -\frac{1}{2} \|f(\mathbf{x})\|_2^2 + \sum s(\mathbf{x}) \right]$$

Transform into latent space and evaluate probability there

$$\varphi(\mathbf{z}) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\mathbf{z}^2}$$


$$-\log(\varphi(f(\mathbf{x}))) = -\log\left(\frac{1}{\sqrt{2\pi}}\right) - \log\left(e^{-\frac{1}{2}f(\mathbf{x})^2}\right) = -\log\left(\frac{1}{\sqrt{2\pi}}\right) + \frac{1}{2}f(\mathbf{x})^2$$

# How to train NF?

Training objective: Minimise negative log likelihood of data

$$\mathcal{L} = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[ -\frac{1}{2} \|f(\mathbf{x})\|_2^2 + \sum s(\mathbf{x}) \right]$$

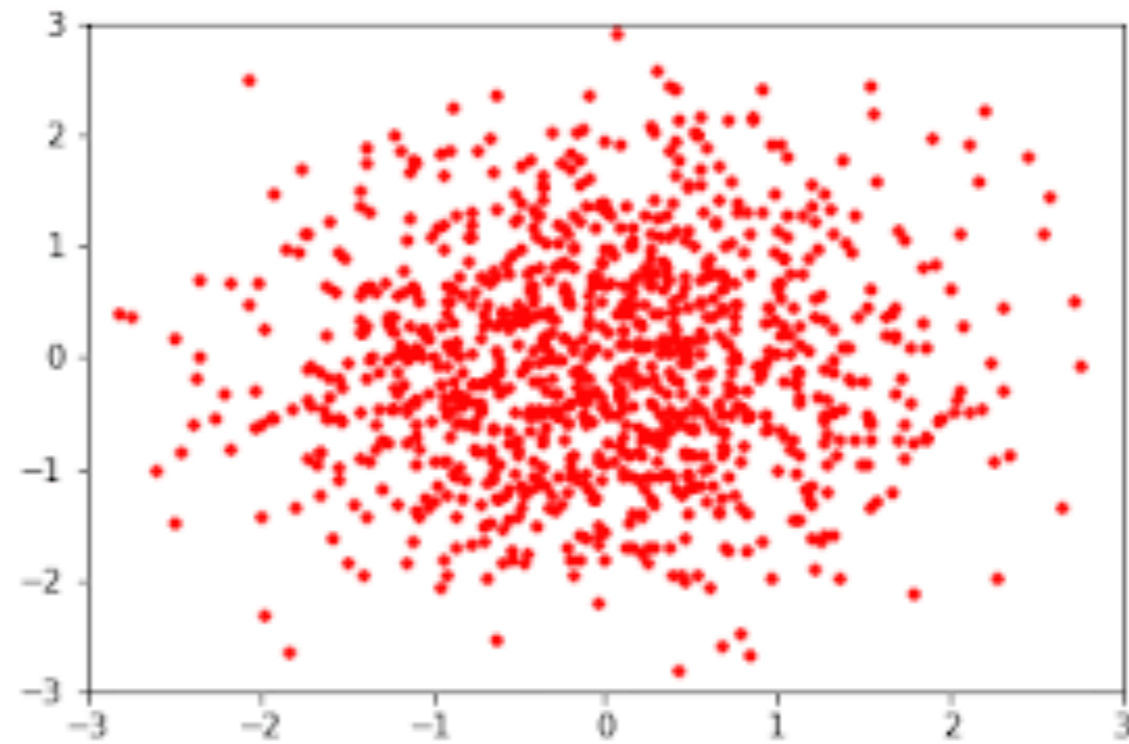
Contribution from Jacobian determinant



$$\det \mathbf{J} = \exp \left( \sum s(\mathbf{x}) \right)$$

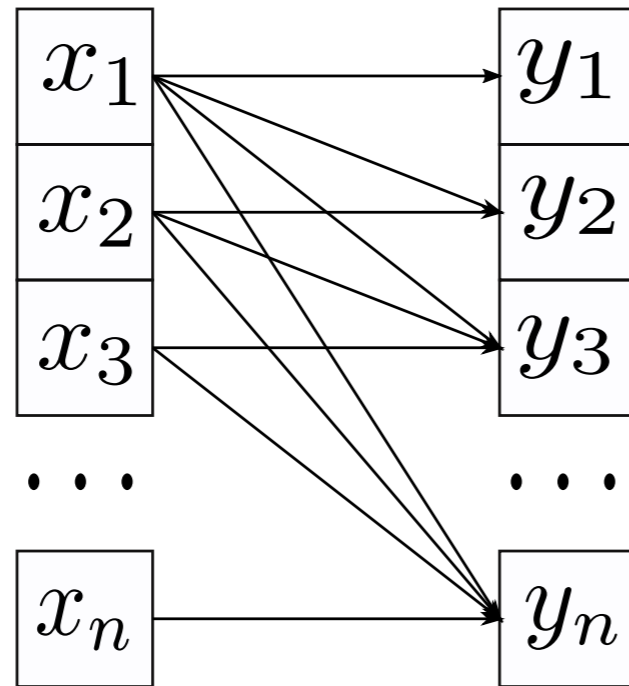
$$-\log(\det \mathbf{J}) = -\sum s(\mathbf{x})$$

# Animation



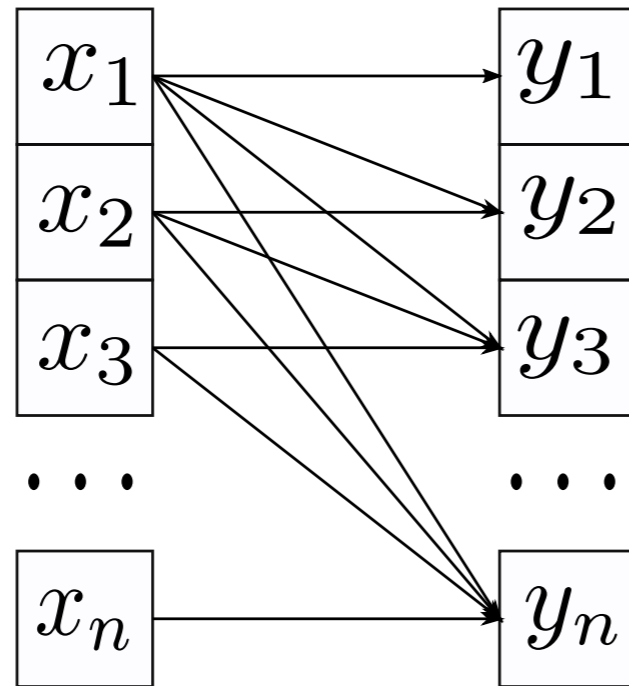


# Autoregressive Flows



Alternative to coupling flows:  
Outputs conditioned on previous inputs

# Autoregressive Flows



Bijjective function

Parameters

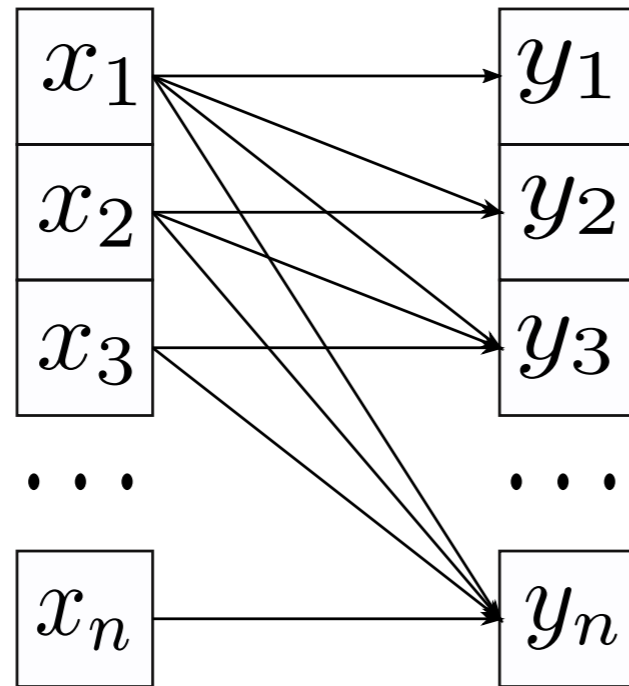
$$y_t = h(x_t; \Theta_t(\mathbf{x}_{1:t-1}))$$

Alternative to coupling flows:

Outputs conditioned on previous inputs

Again: simple Jacobian and invertible functions

# Autoregressive Flows



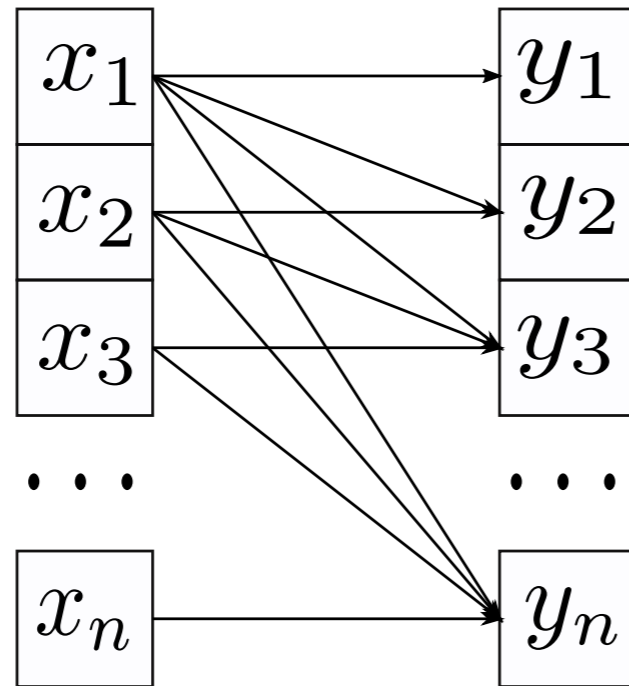
$$y_t = h(x_t; \Theta_t(\mathbf{x}_{1:t-1}))$$

Masked autoregressive flow (MAF):

Fast: Data  $\rightarrow$  latent space

Slow: Latent space  $\rightarrow$  data

# Autoregressive Flows

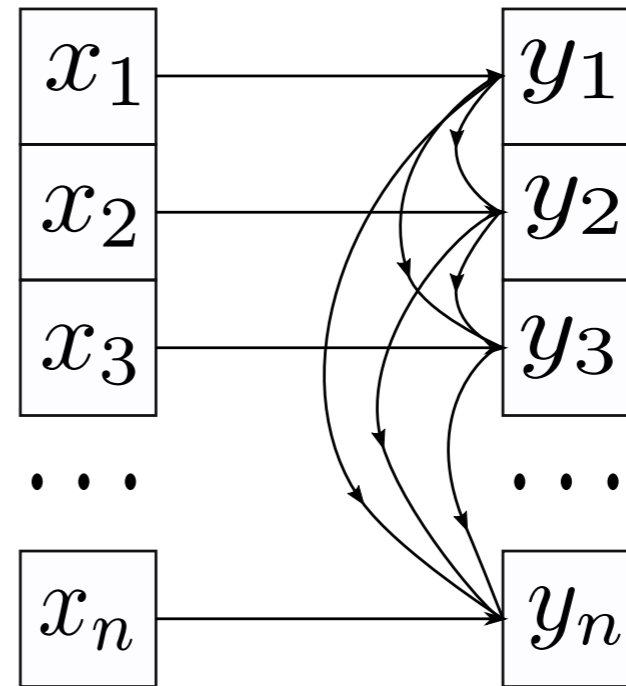


$$y_t = h(x_t; \Theta_t(\mathbf{x}_{1:t-1}))$$

Masked autoregressive flow (MAF):

Fast: Data  $\rightarrow$  latent space

Slow: Latent space  $\rightarrow$  data



$$y_t = h(x_t; \theta_t(\mathbf{y}_{1:t-1}))$$

Inverse autoregressive flow (IAF):

Slow: Data  $\rightarrow$  latent space

Fast: Latent space  $\rightarrow$  data



# Comments on Flows

Only scratched the surface:  
more constructions available

§3.1 Elementwise bijections Non-linear elementwise transform	→	Problem: no mixing of variables
§3.2 Linear flows Affine combination of variables	→	Problem: limited representational power
§3.3 Planar and radial flows Non-linear transforms	→	Problem: hard to compute inverse
§3.4.1 Coupling flows	→	Depend on coupling functions
§3.4.2 Autoregressive flows Architectures that allow invertible non-linear transformations.		
§3.5 Residual flows Invertible residual networks		
§3.6 Infinitesimal flows Continuous flows depending on ODEs or SDEs		

§3.4.4 Coupling functions

- Affine
- Non-linear squared
- Continuous mixture CDFs
- Splines
- Neural autoregressive
- Sum-of-squares polynomial
- Piecewise-bijective

# Comments on Flows

Only scratched the surface:  
more constructions available

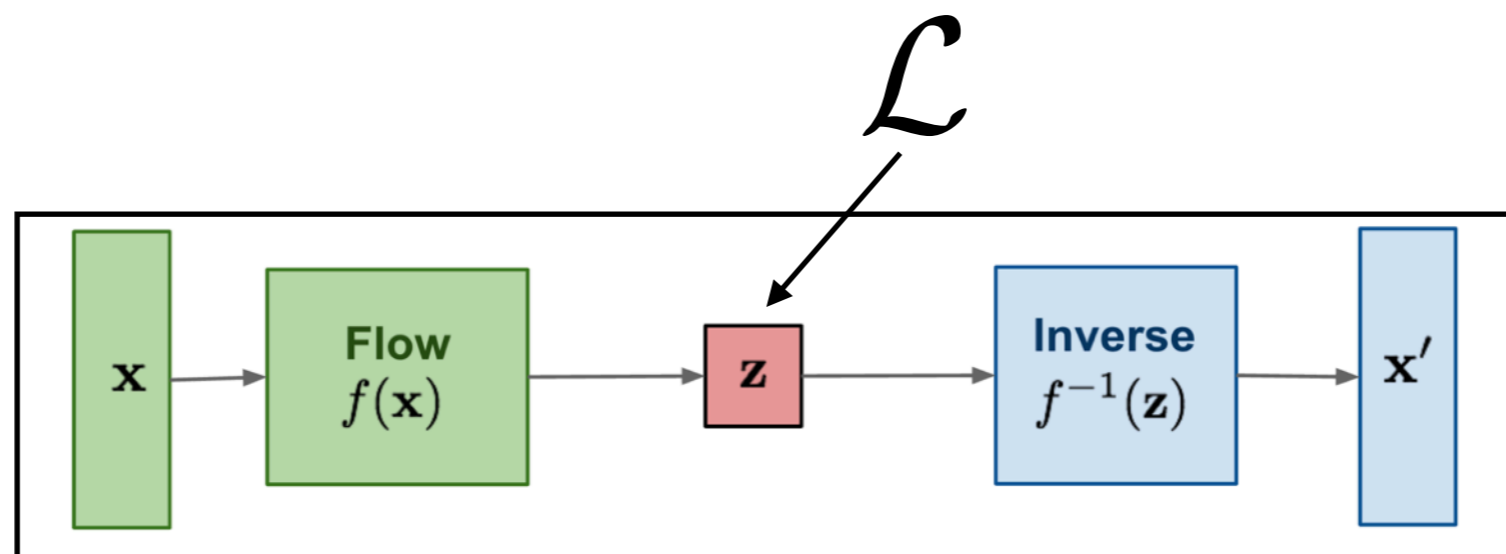
Exact learning of likelihood  
→ Better generative fidelity  
→ Can evaluate likelihood of  
data

More complex  
→ Slower, choice of fast direction

§3.1 Elementwise bijections Non-linear elementwise transform	→ Problem: no mixing of variables
§3.2 Linear flows Affine combination of variables	→ Problem: limited representational power
§3.3 Planar and radial flows Non-linear transforms	→ Problem: hard to compute inverse
§3.4.1 Coupling flows §3.4.2 Autoregressive flows Architectures that allow invertible non-linear transformations.	→ Depend on coupling functions
§3.5 Residual flows Invertible residual networks	
§3.6 Infinitesimal flows Continuous flows depending on ODEs or SDEs	

§3.4.4 Coupling functions

- Affine
- Non-linear squared
- Continuous mixture CDFs
- Splines
- Neural autoregressive
- Sum-of-squares polynomial
- Piecewise-bijective





# Simulation & Generative Models Part II

Gregor Kasieczka

Email: [gregor.kasieczka@uni-hamburg.de](mailto:gregor.kasieczka@uni-hamburg.de)

Twitter/X: [@GregorKasieczka](https://twitter.com/GregorKasieczka)

COFI Winter School 2023

**CLUSTER OF EXCELLENCE**  
QUANTUM UNIVERSE



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG



**KISS**



**FSP**  
CMS

CDCS

CENTER FOR DATA AND COMPUTING  
IN NATURAL SCIENCES



**DASHH**

**PIER**

Partnership of  
Universität Hamburg and DESY

GEFÖRDERT VOM



Bundesministerium  
für Bildung  
und Forschung

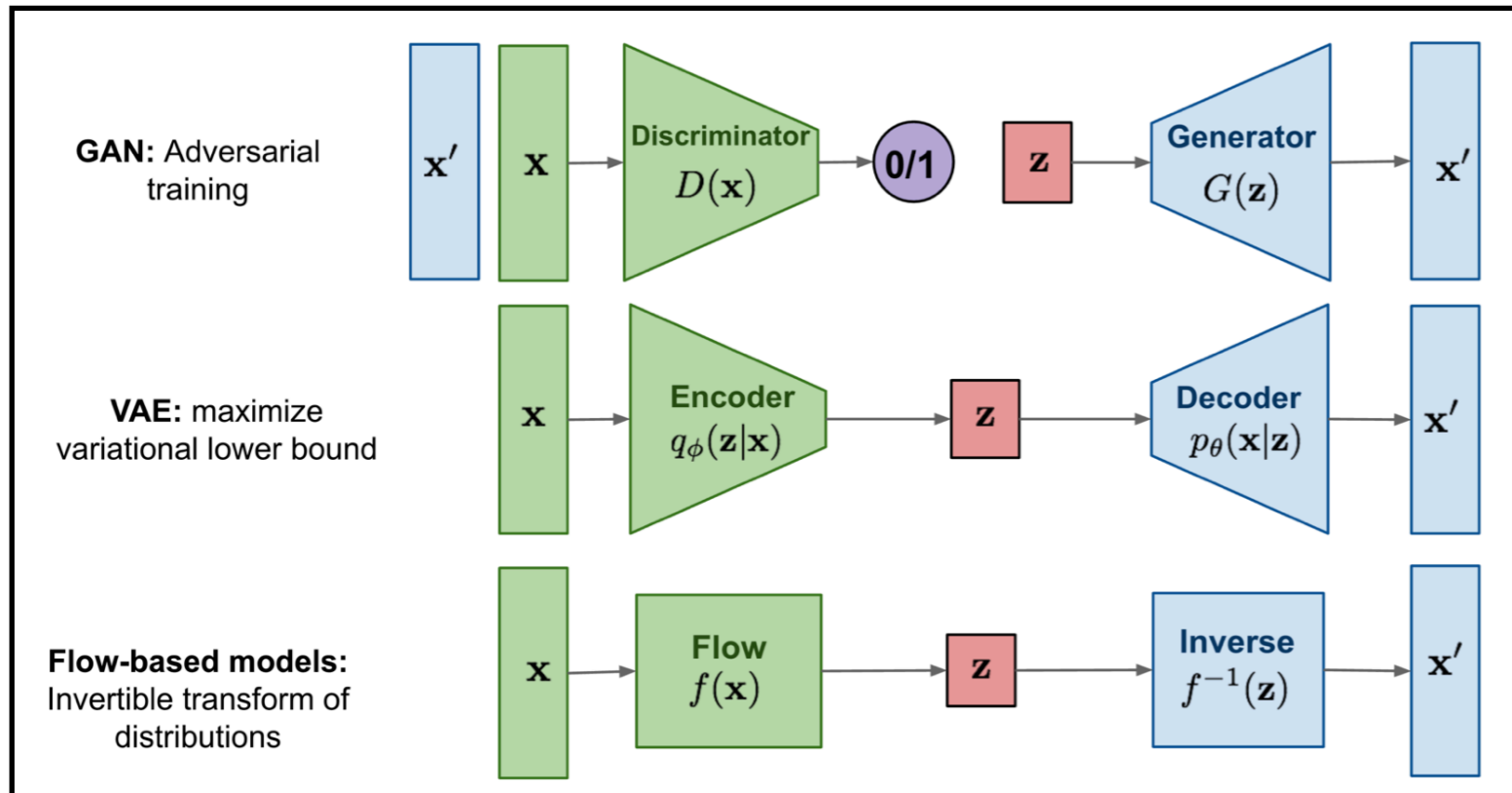
Emmy  
Noether-  
Programm

Deutsche  
Forschungsgemeinschaft  
**DFG**



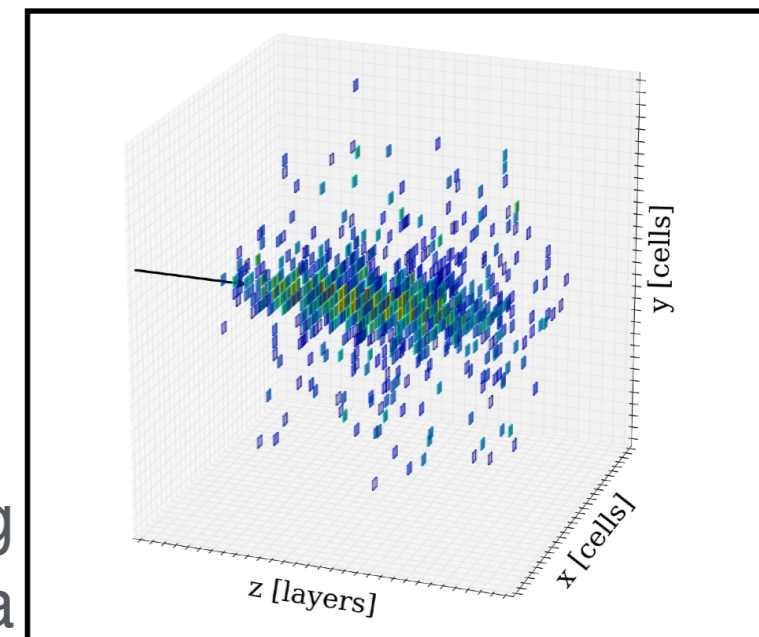


# Yesterday



Understand three basic generative architectures

First look at simulating fixed grid data





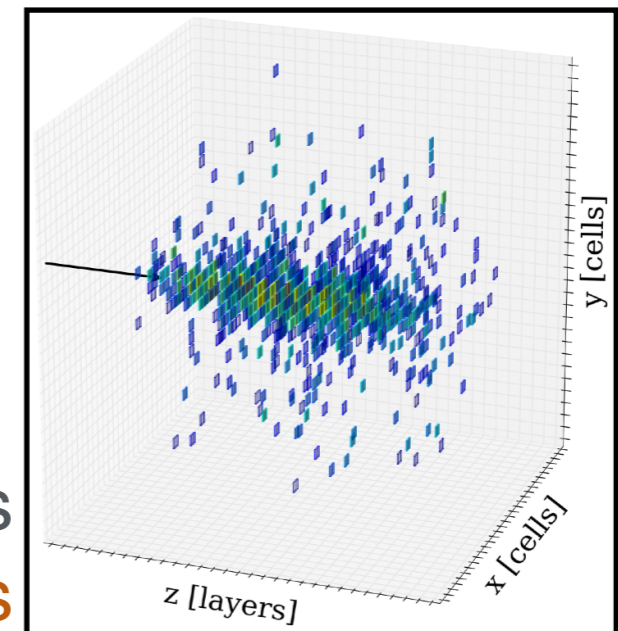
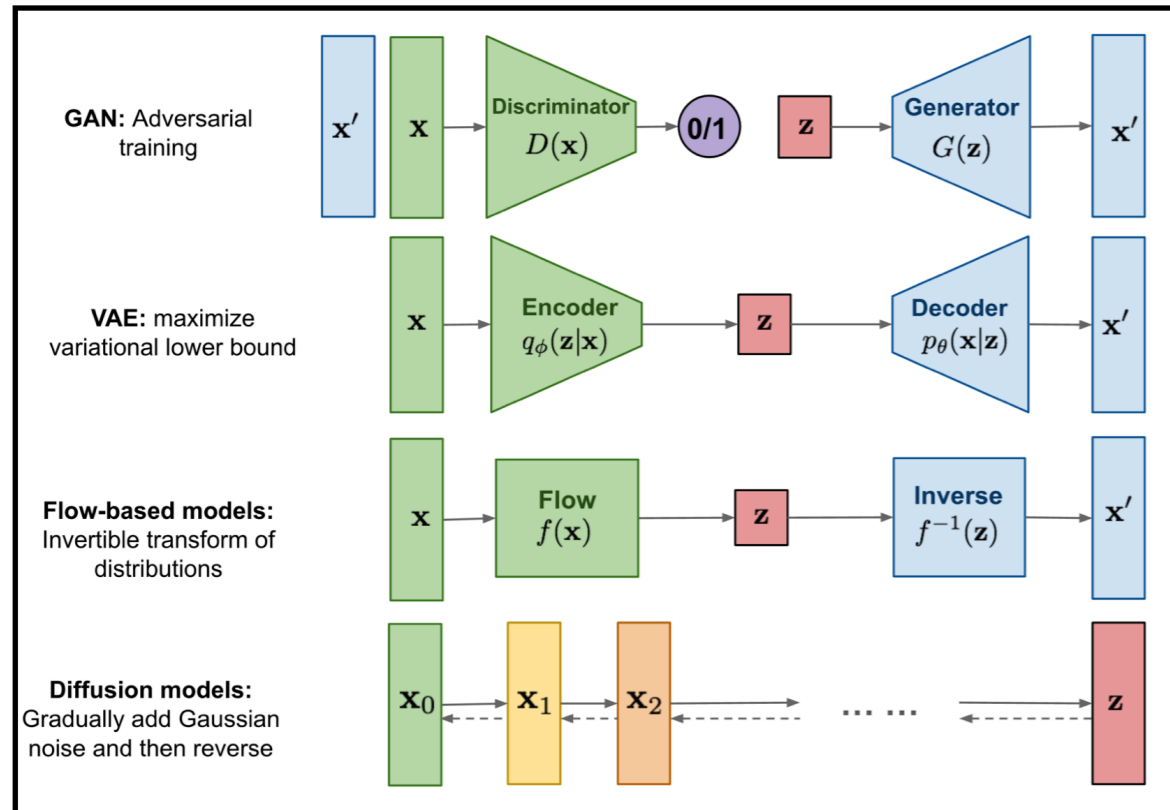
# Overview

## 1. Common architectures

architectures

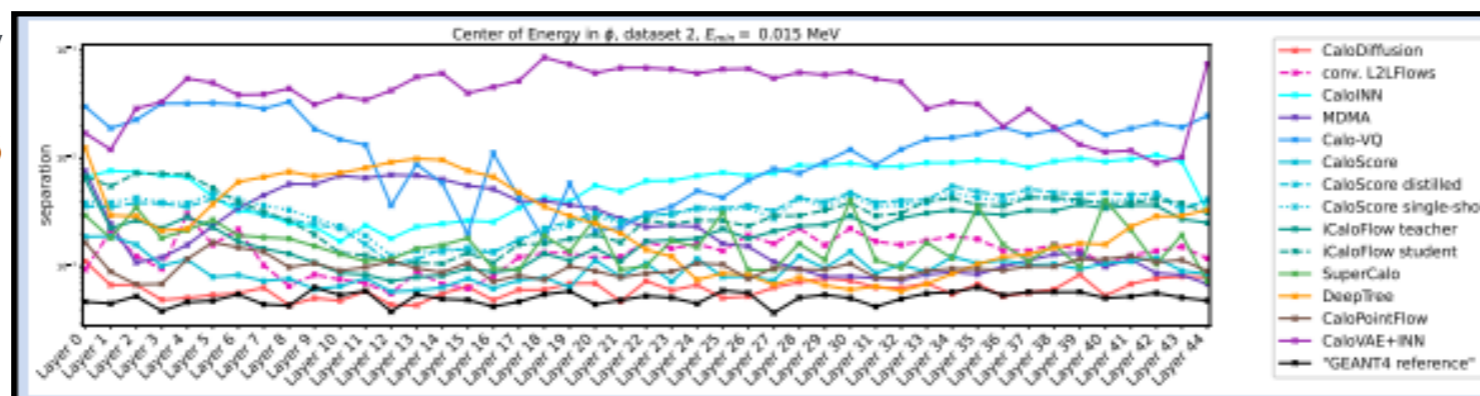
-> GANs, VAEs, NF yesterday

-> Diffusion & CNF today



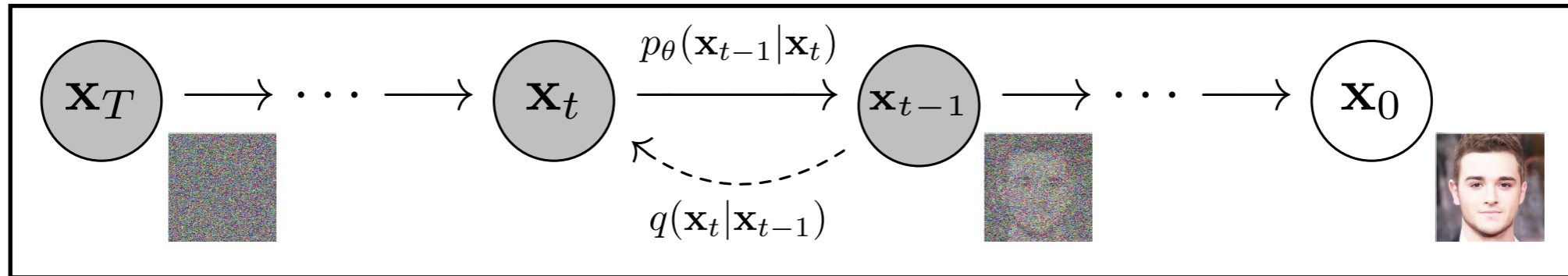
## 2. Physics applications

## 3. Quality metrics



# Diffusion Models

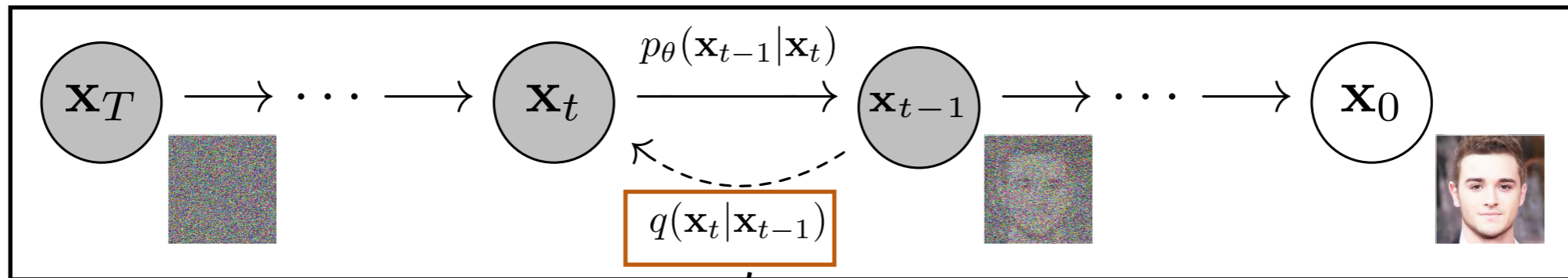
# Diffusion Model



Core idea: Stepwise transition  
from pure noise to data

Markov chain

# Diffusion Model



Core idea: Stepwise transition from pure noise to data

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

Individual step

Forward  
(Data  $\rightarrow$  Noise)

Noise schedule  
(hyper-parameter)

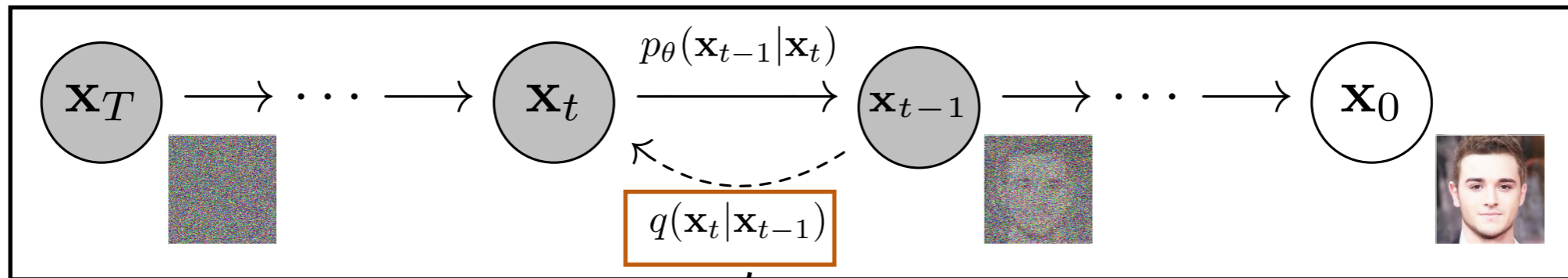
$$q(\mathbf{x}_{1:T} | \mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

Full chain

This is added in the training process



# Diffusion Model



Core idea: Stepwise transition from pure noise to data

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

Forward  
(Data  $\rightarrow$  Noise)

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

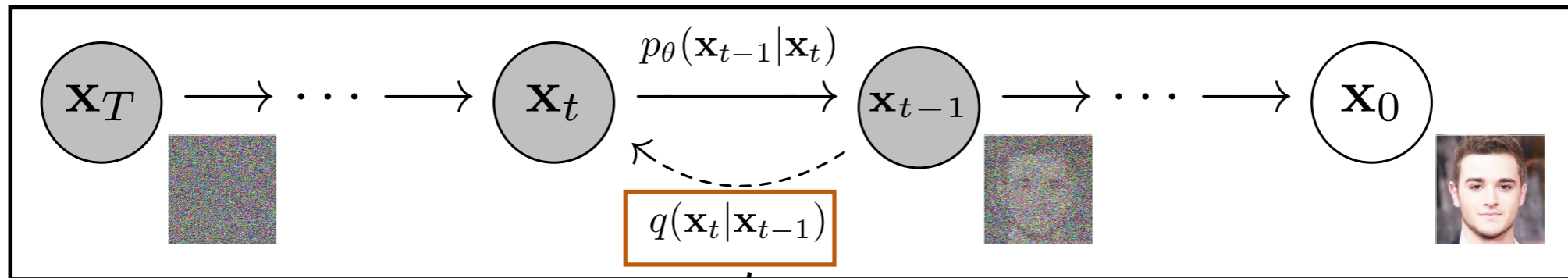
Shortcut:  
transition to any time

$$\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$$

$$\alpha_t := 1 - \beta_t$$

This is added in the training process

# Diffusion Model



Core idea: Stepwise transition from pure noise to data

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

Forward  
(Data  $\rightarrow$  Noise)

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

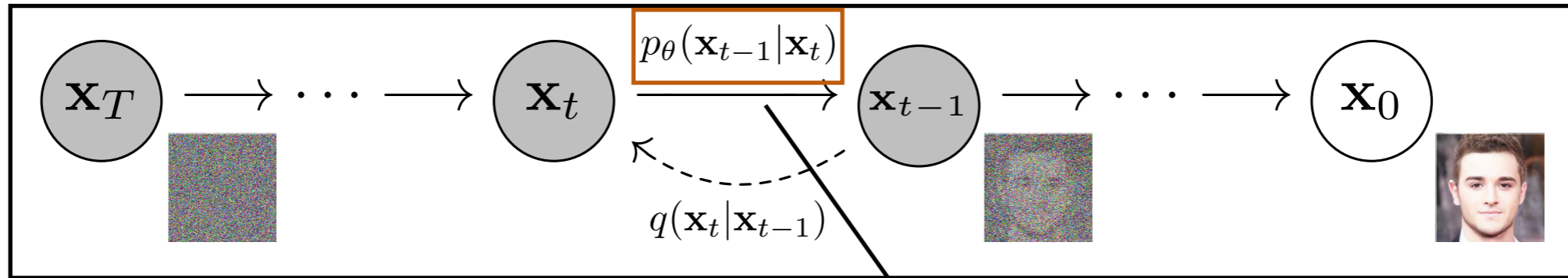
$$\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon} \text{ for } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Rewrite: State at **any time**

Will try to **predict**

This is added in the training process

# Diffusion Model



Core idea: Stepwise transition  
from pure noise to data

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

To be **learned**  
(in principle)

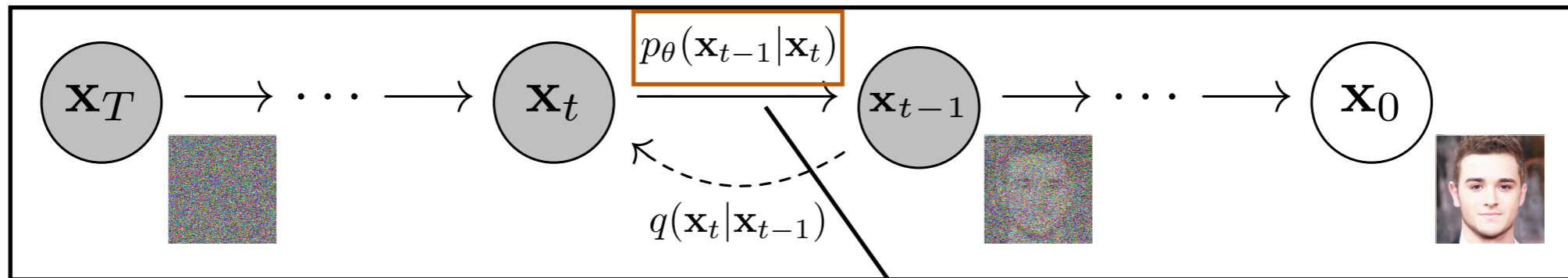
$\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$   
(hyper-parameter)

Reverse  
(Noise  $\rightarrow$  data)

$$p_\theta(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$

# Diffusion Model



Core idea: Stepwise transition  
from pure noise to data

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

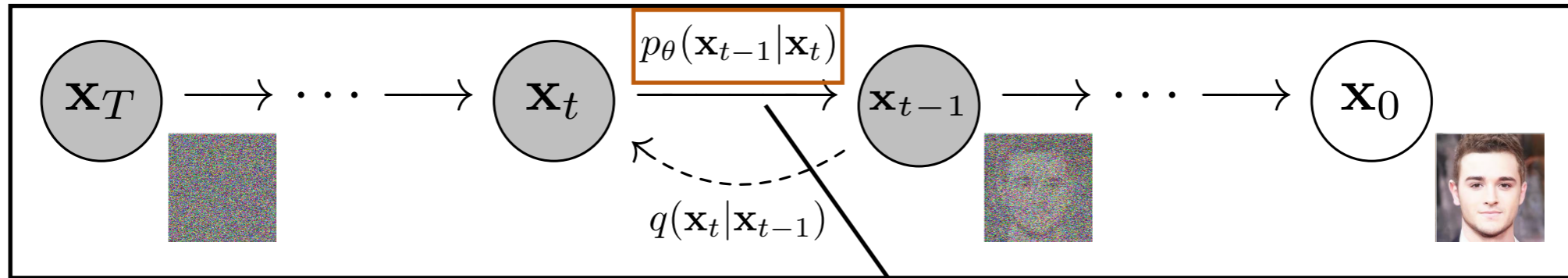
Reverse  
(Noise  $\rightarrow$  data)

$$\frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

Optimal **expected mean** (take into  
account known noise schedule)



# Diffusion Model



Core idea: Stepwise transition  
from pure noise to data

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

Reverse  
(Noise  $\rightarrow$  data)

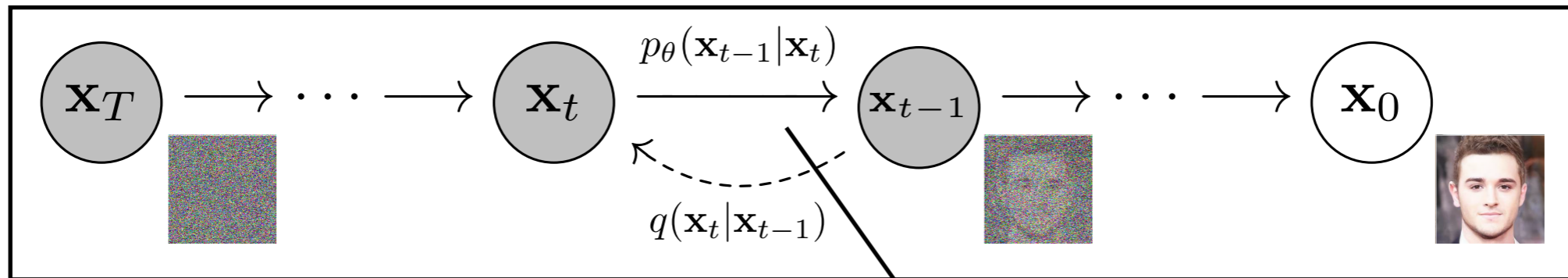
$$\frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

Use to parametrise function:

$$\boldsymbol{\mu}_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(x_t, t) \right)$$

and learn to **predict noise**

# Diffusion Model



Core idea: Stepwise transition  
from pure noise to data

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

Reverse  
(Noise  $\rightarrow$  data)

Resulting learning objective

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}} \left[ \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2 \right]$$

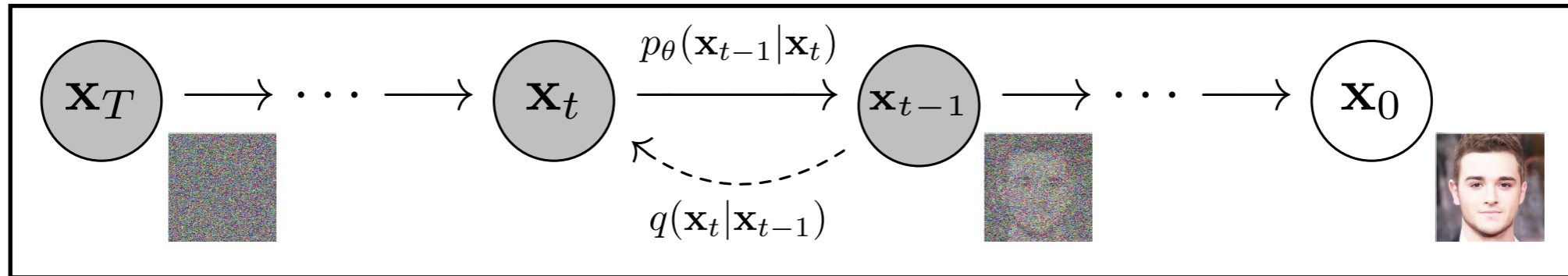
Noisy image

Reminder: Forward  
diffusion to time t

$$\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$

Timestep

# Diffusion Model



Core idea: Stepwise transition  
from pure noise to data

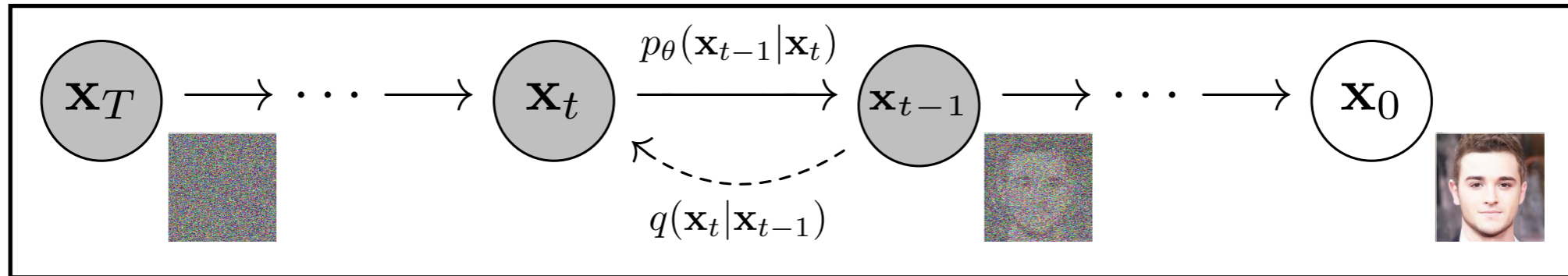
---

## Algorithm 1 Training

---

- 1: **repeat**
  - 2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
  - 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
  - 4:  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 5: Take gradient descent step on  
$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$
  - 6: **until** converged
-

# Diffusion Model



Core idea: Stepwise transition  
from pure noise to data

---

## Algorithm 1 Training

---

- 1: **repeat**
  - 2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
  - 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
  - 4:  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 5: Take gradient descent step on  
$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$
  - 6: **until** converged
- 

---

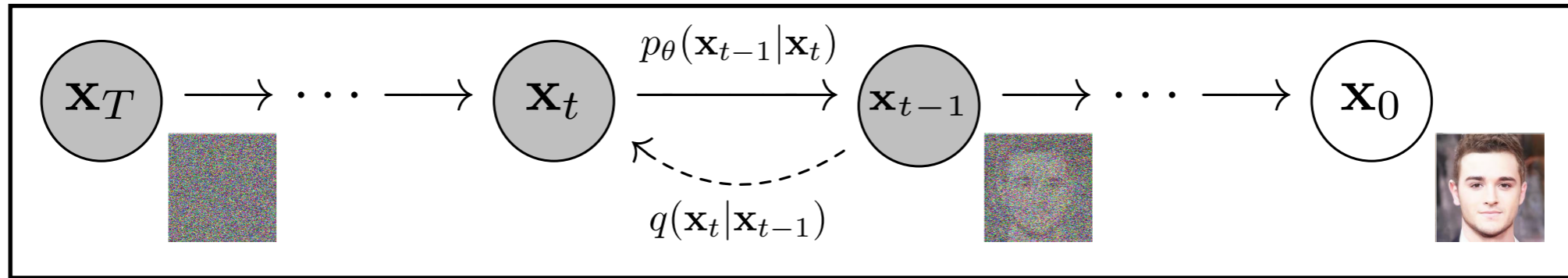
## Algorithm 2 Sampling

---

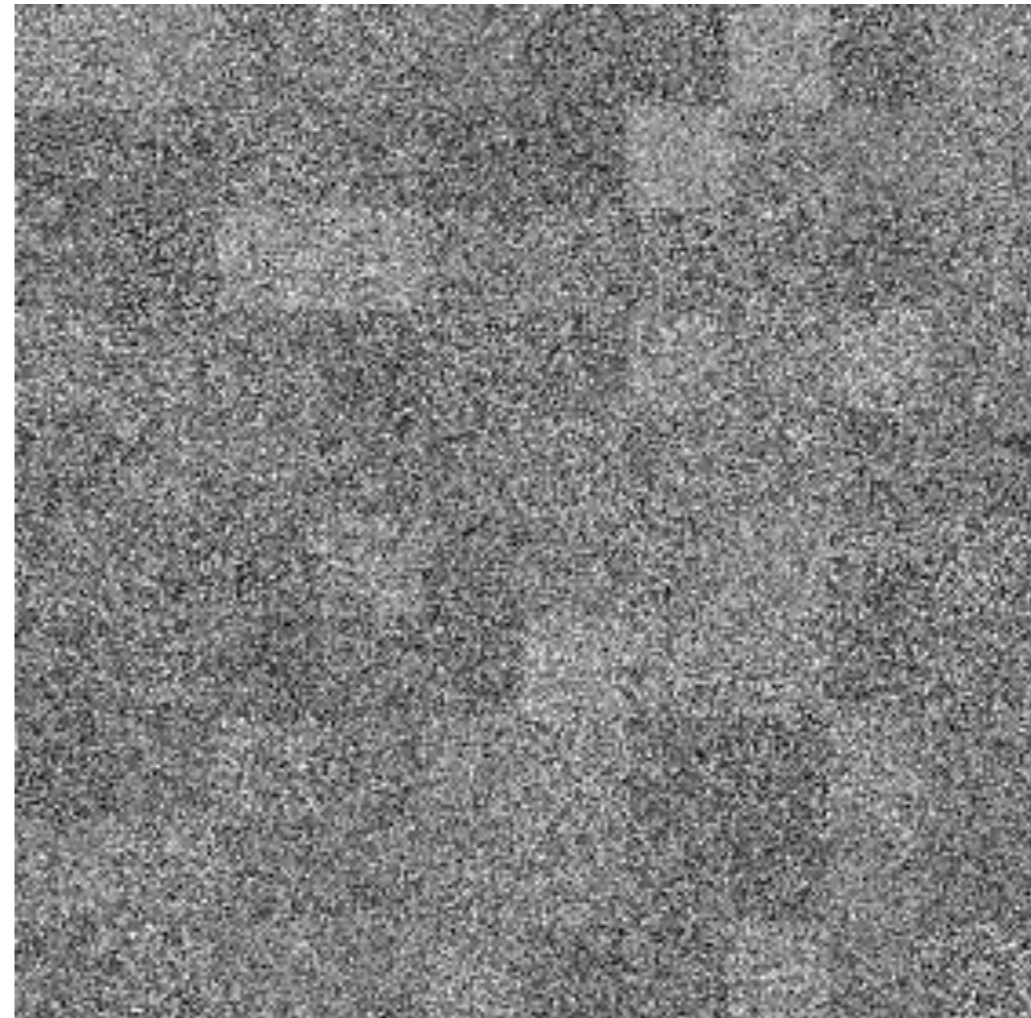
- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 2: **for**  $t = T, \dots, 1$  **do**
  - 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
  - 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
  - 5: **end for**
  - 6: **return**  $\mathbf{x}_0$
-



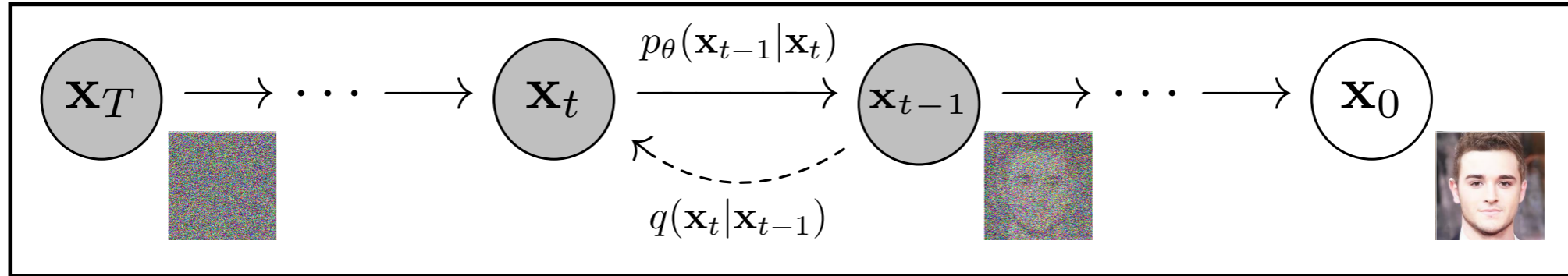
# Diffusion Model



Core idea: Stepwise transition from pure noise to data

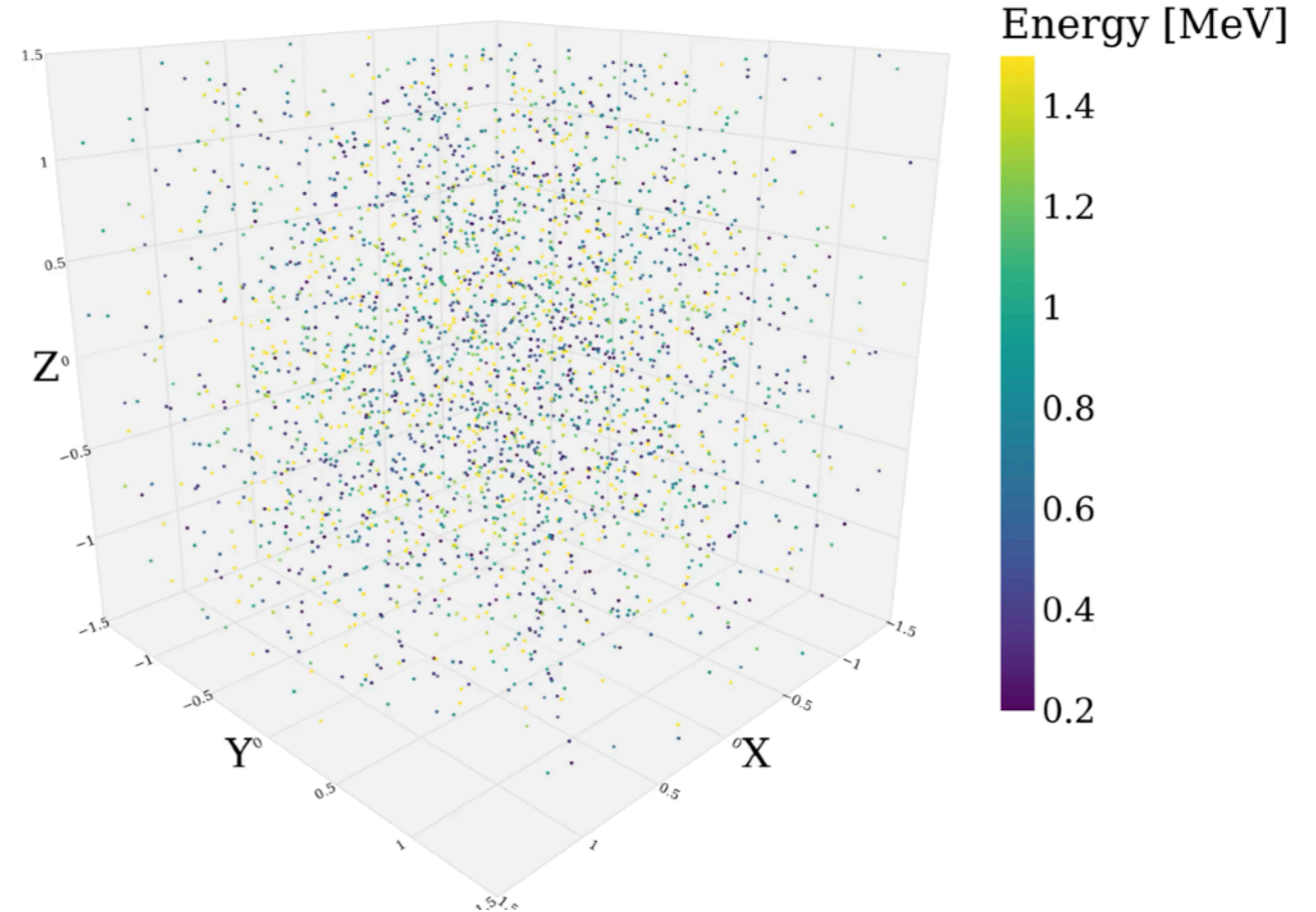


# Diffusion Model

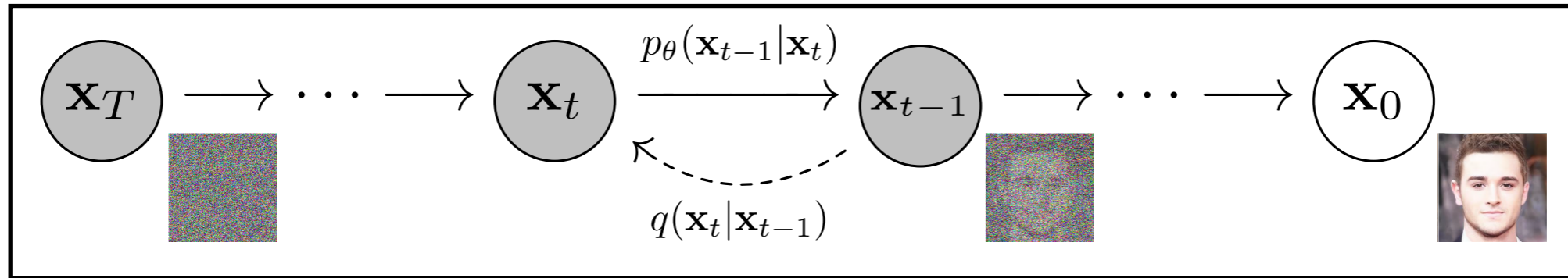


Core idea: Stepwise transition from pure noise to data

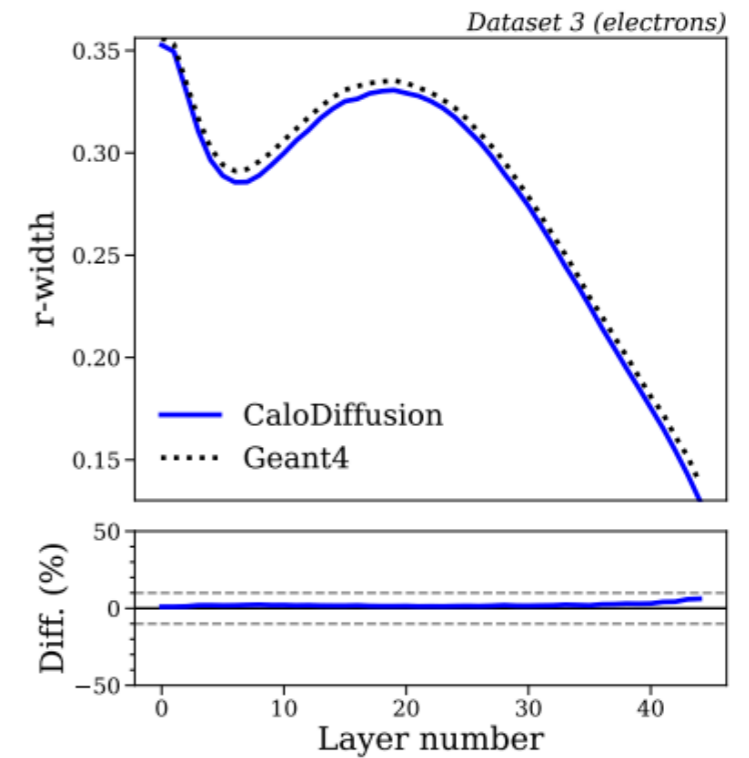
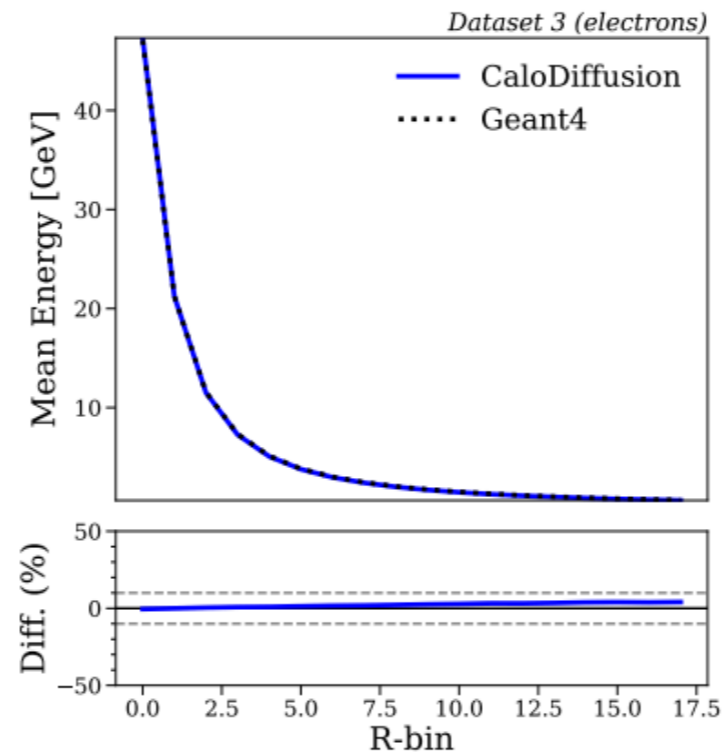
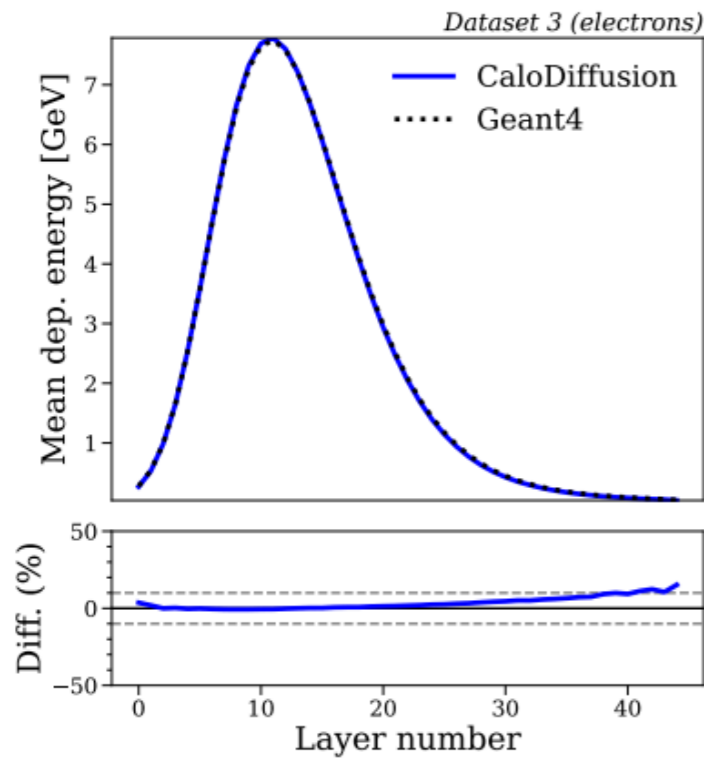
CaloCloud, time stamp:  $t_{99}$



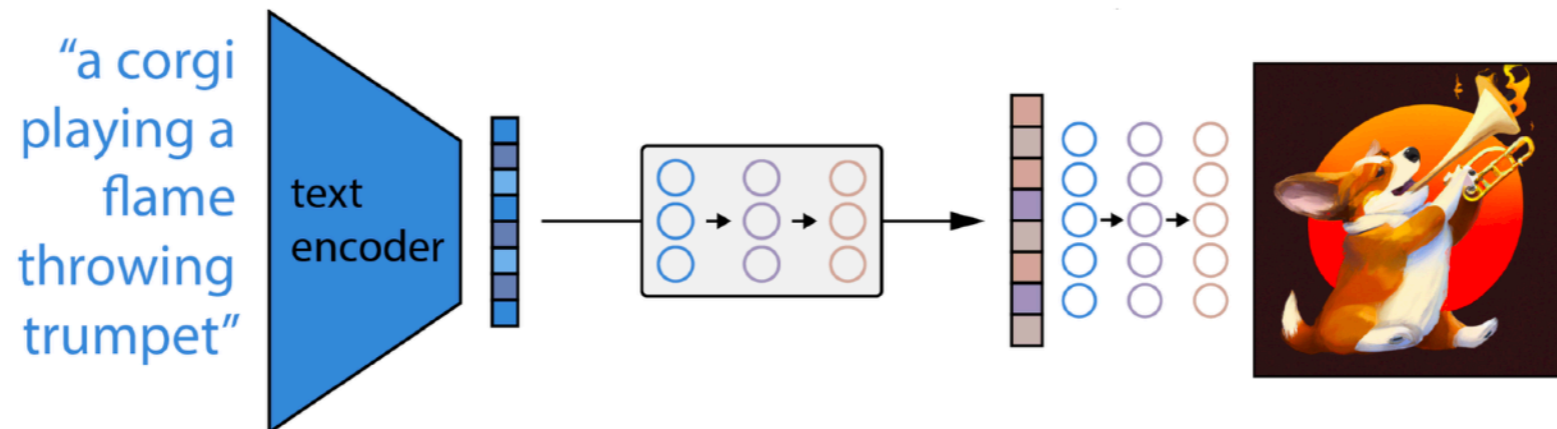
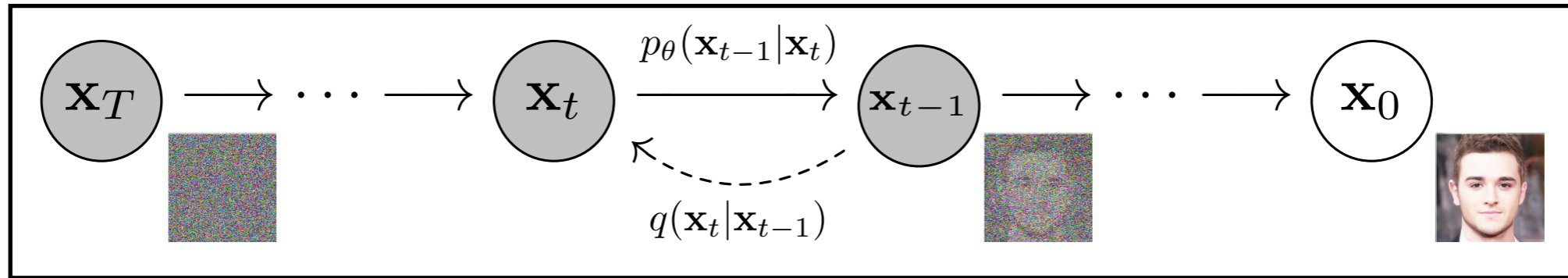
# Diffusion Model



Core idea: Stepwise transition from pure noise to data



# Diffusion Model

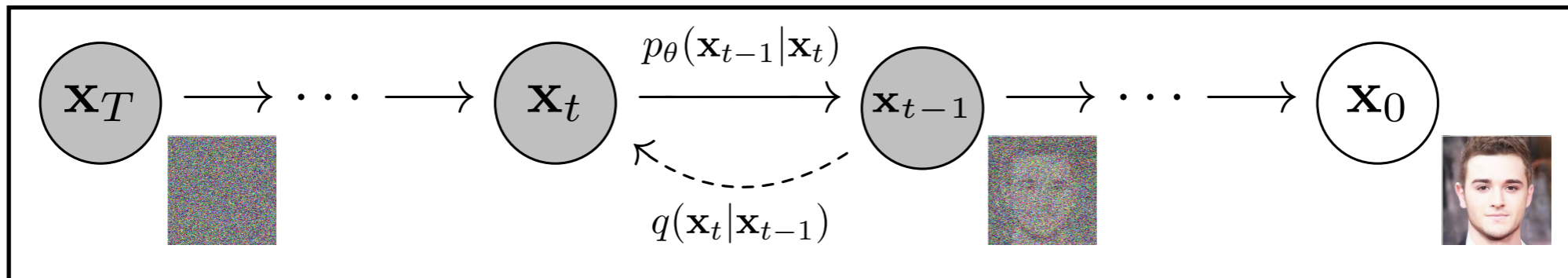


*DALLE-2 (Roughly)*

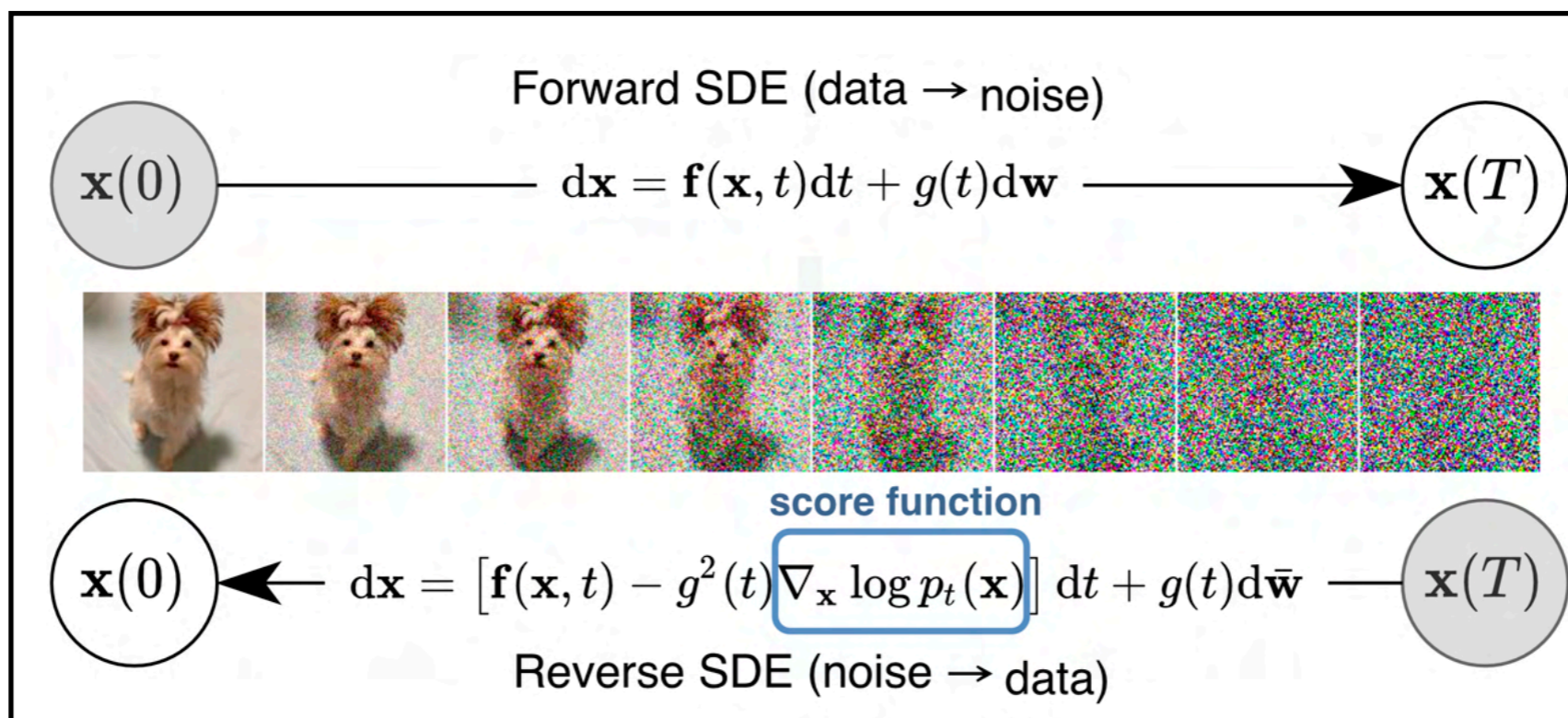
*Contrastive learning of joint image/text embedding +  
translation of embedding +  
conditioned diffusion*



# Continuous Diffusion Model



Replace **discrete** noise steps



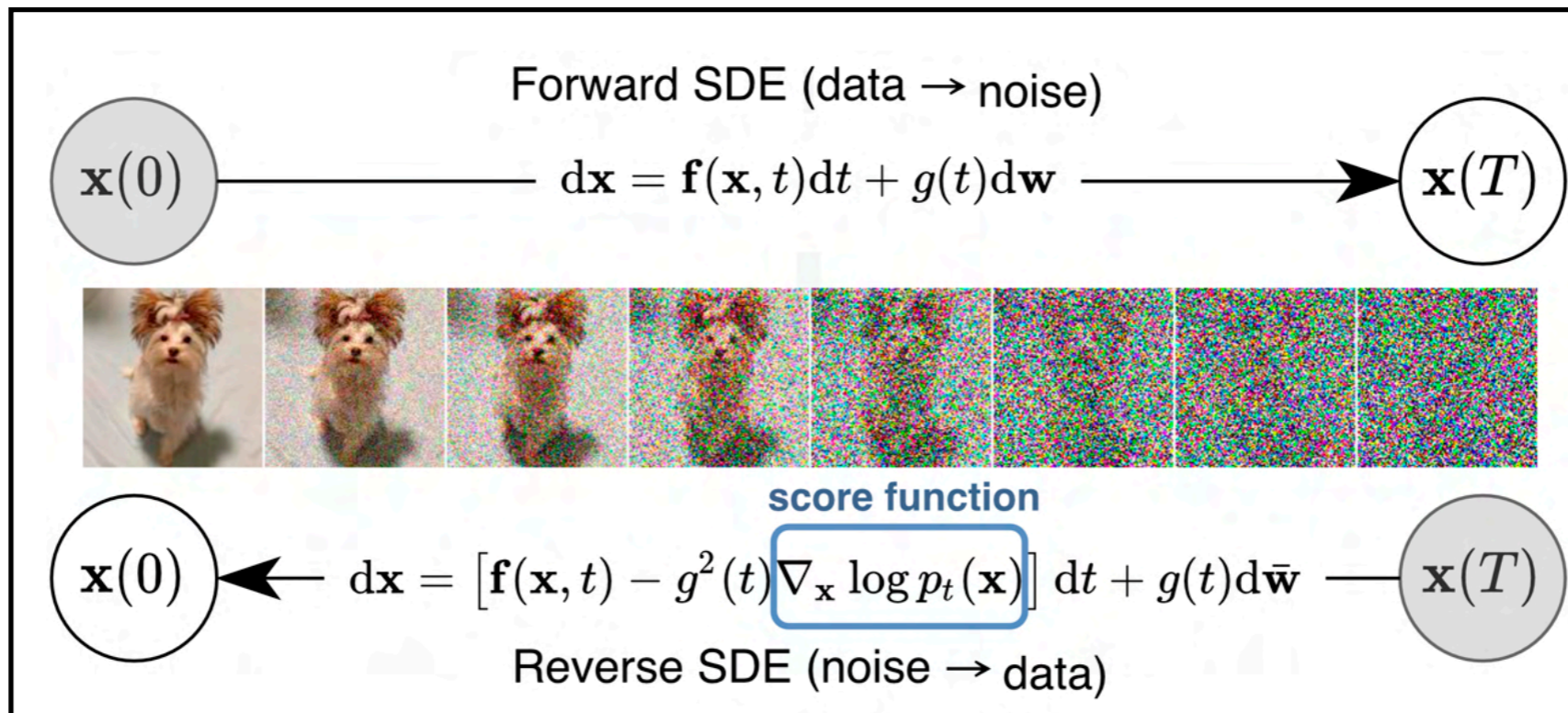
with **stochastic differential equations** (SDEs)

# Continuous Diffusion Model

Forward SDE: 
$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w}$$

↑ Drift term      ↑ Diffusion term      ← Wiener process/  
 Brownian motion  
 (independent  
 Gaussian  
 increments)

(Correspond to the noise schedule in discrete case)



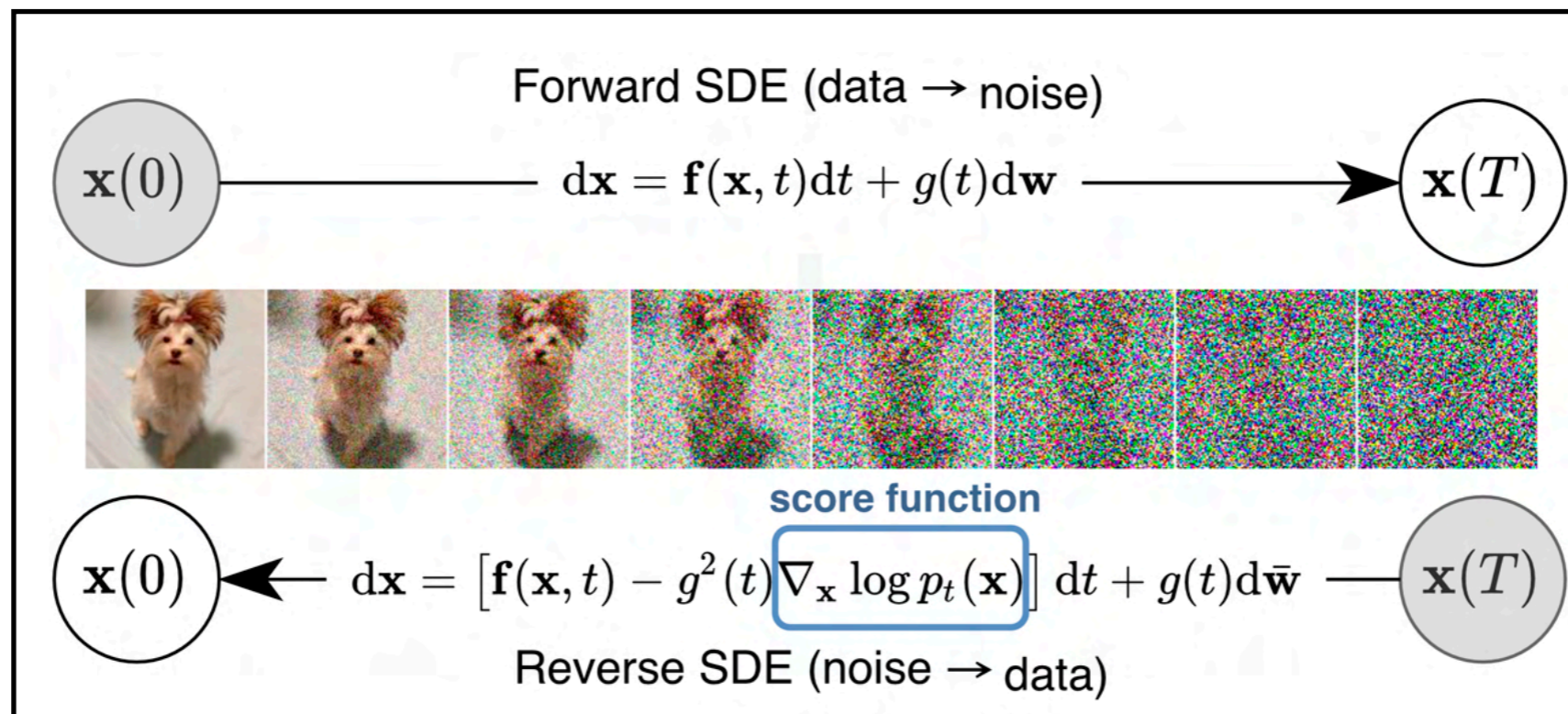
# Continuous Diffusion Model

Probability density of  $x(t)$

Reverse SDE: 
$$dx = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t) d\bar{w}$$

Score function

Reverse of a diffusion process is also a diffusion



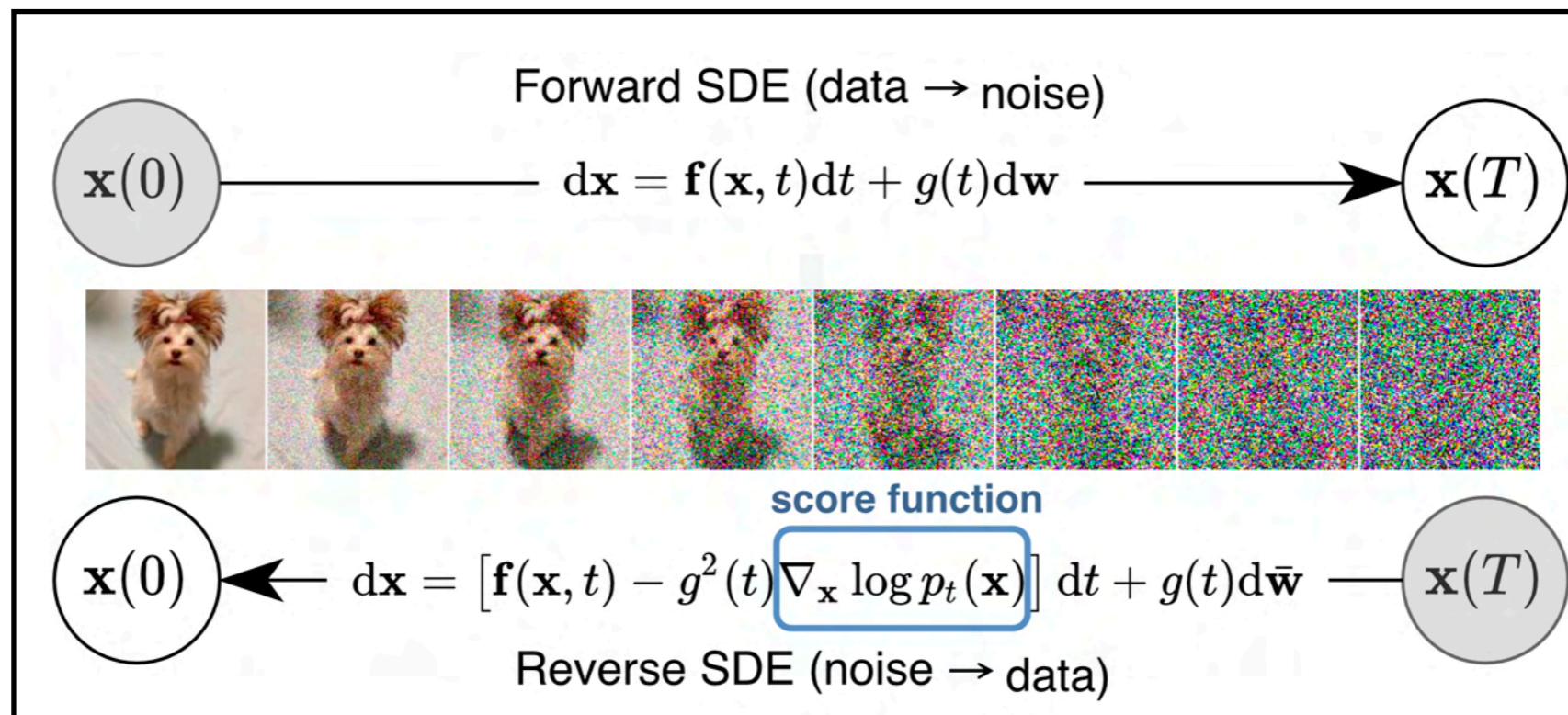


# Continuous Diffusion Model

Reverse SDE: 
$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t) d\bar{\mathbf{w}}$$

$$\theta^* = \arg \min_{\theta} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{\mathbf{x}(0)} \mathbb{E}_{\mathbf{x}(t)|\mathbf{x}(0)} \left[ \left\| \mathbf{s}_{\theta}(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) \right\|_2^2 \right] \right\}$$

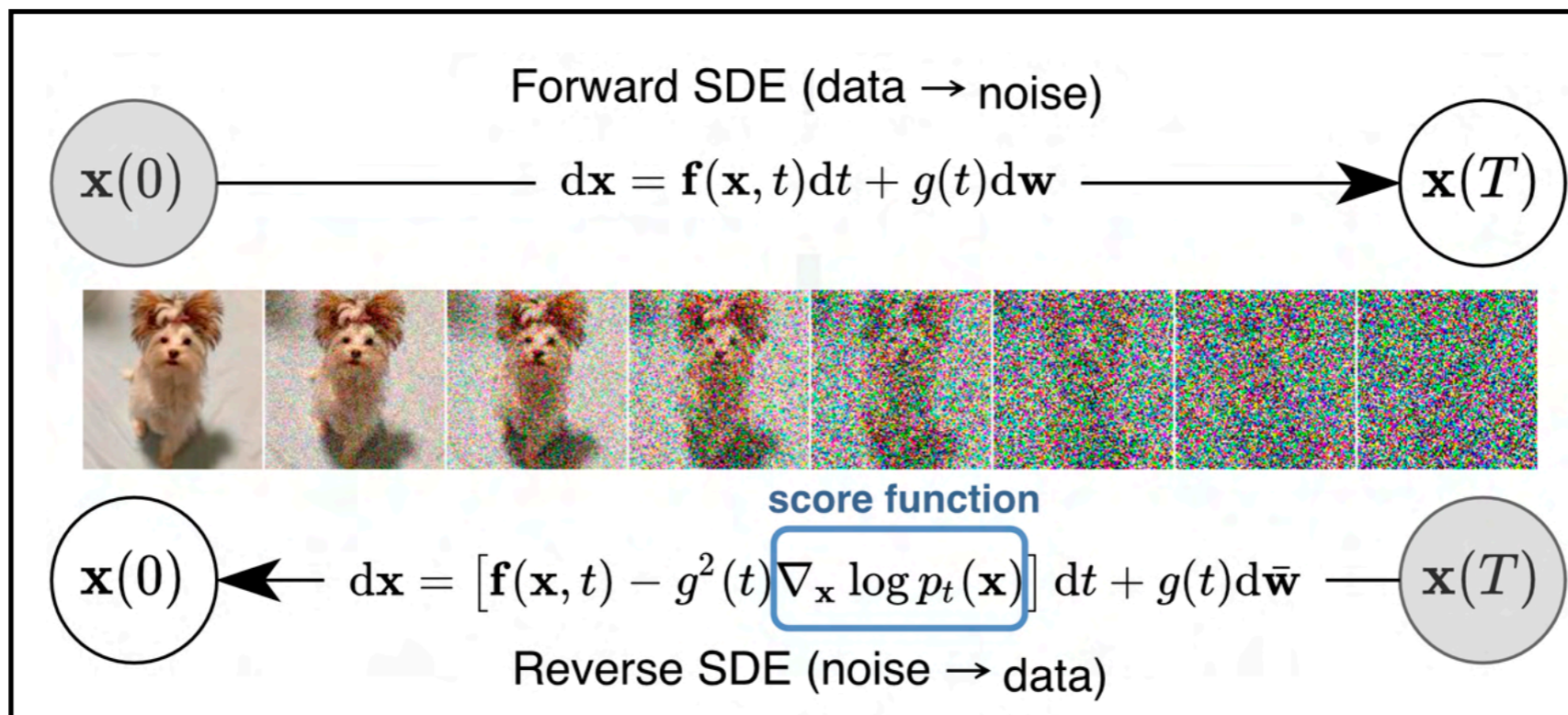
Learn to approximate score function with neural network





# Continuous Diffusion Model

Once trained: Sample latent space and numerically solve SDE to transport to data space



# Aside

The image shows two presentation slides. The left slide is titled "Forward SDE (data → noise)" and shows a sequence of images of a dog's face being corrupted by noise over time, with the equation  $dx = f(x, t)dt + g(t)dw$ . Below it, the "Reverse SDE (noise → data)" is shown with the equation  $dx = [f(x, t) - g^2(t) \nabla_x \log p_t(x)] dt + g(t)dw$ , where the term  $\nabla_x \log p_t(x)$  is labeled as the "score function". The right slide is a flow diagram showing a green box labeled "x" pointing to a green box labeled "Flow f(x)", which points to a red box labeled "z", which then points to a blue box labeled "Inverse f^{-1}(z)", which finally points to a blue box labeled "x'".

**Corporate needs you to find the differences between this picture and this picture.**

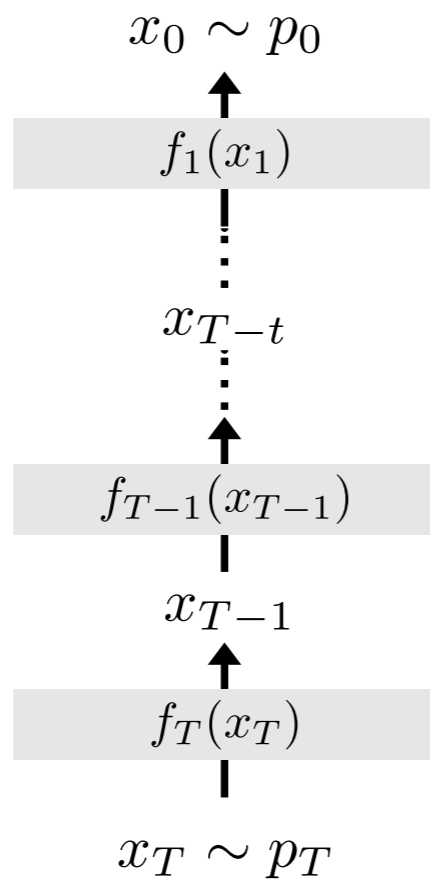
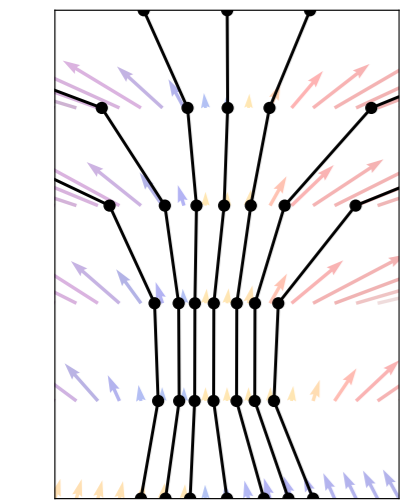
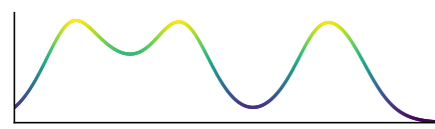


**They're the same picture.**

Well, almost

# Flow Matching

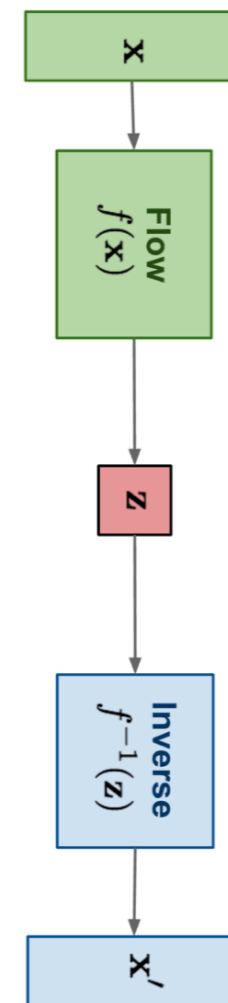
# Remember: Normalising Flows



↑ generate

↓ train

Can view each layer in the flow as a **discrete time step**



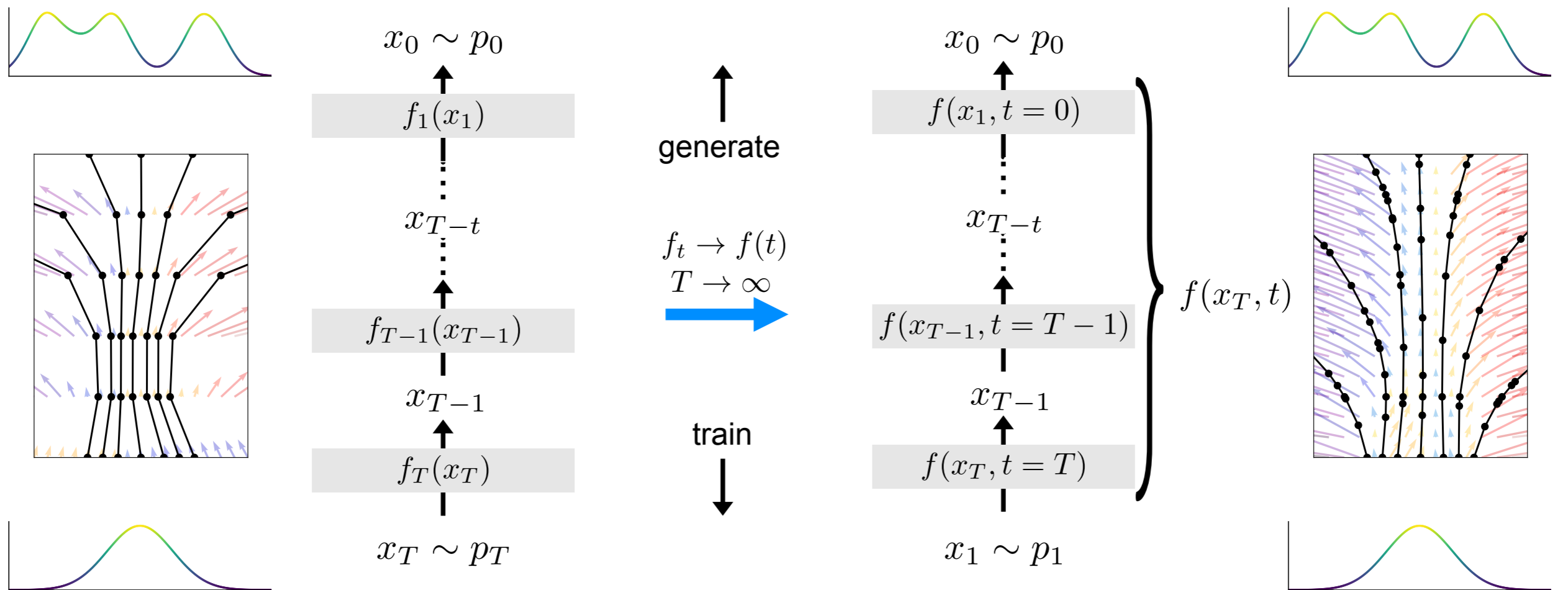
discrete normalizing flow

$$x_0 = f_1^{\theta_1} \dots \circ f_{T-1}^{\theta_{T-1}} \circ f_T^{\theta_T}(x_T)$$

$$\log p_0(x_0) = \log p_T(x_T) - \sum_t \log |\det J_t|$$



# Remember: Normalising Flows



discrete normalizing flow

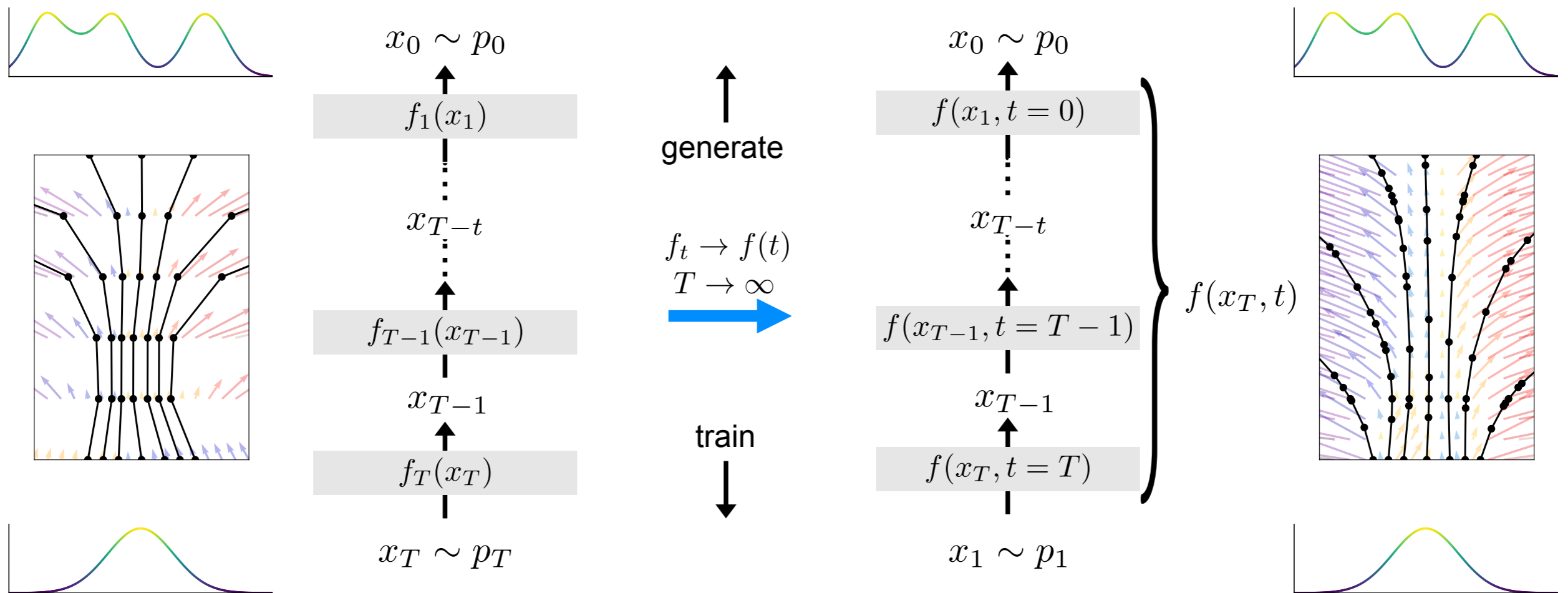
$$x_0 = f_1^{\theta_1} \dots \circ f_{T-1}^{\theta_{T-1}} \circ f_T^{\theta_T}(x_T)$$

continuous normalizing flow

$$\frac{dx_t}{dt} = v_t^\theta(x_t), \quad x_t := f(x_T, t)$$

Take continuous limit:  
 Approximate vector field with network

# Remember: Normalising Flows



discrete normalizing flow

$$x_0 = f_1^{\theta_1} \dots \circ f_{T-1}^{\theta_{T-1}} \circ f_T^{\theta_T}(x_T)$$

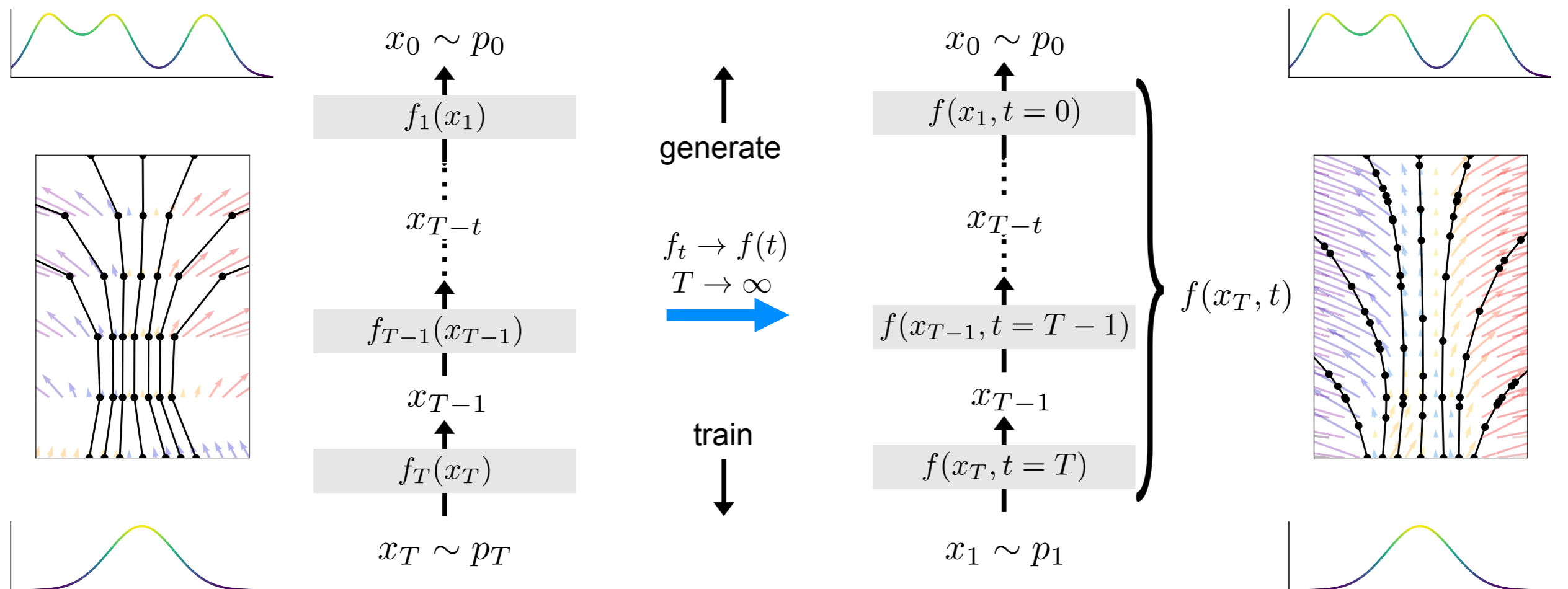
continuous normalizing flow

$$\frac{dx_t}{dt} = v_t^\theta(x_t), \quad x_t := f(x_T, t)$$

Change in loss:

$$\log p_0(x_0) = \log p_T(x_T) - \sum_t \log |\det J_t| \longrightarrow \log p_0(x_0) = \log p_1(x_1) - \int_{t_1}^{t_0} \text{tr} \left( \frac{\partial u_t}{\partial x_t} \right) dt$$

# Remember: Normalising Flows



discrete normalizing flow

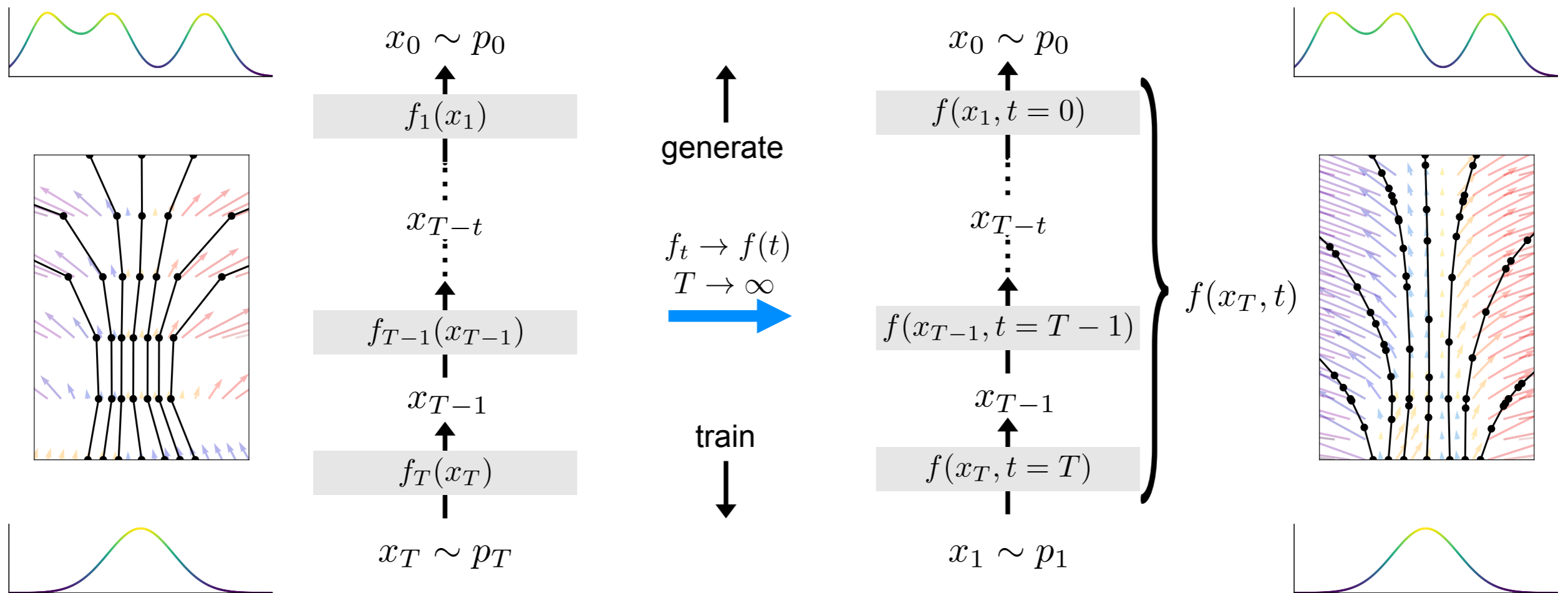
$$x_0 = f_1^{\theta_1} \dots \circ f_{T-1}^{\theta_{T-1}} \circ f_T^{\theta_T}(x_T)$$

continuous normalizing flow

$$\frac{dx_t}{dt} = v_t^\theta(x_t), \quad x_t := f(x_T, t)$$

For sampling:  
Solve differential equation (ODE)

# Remember: Normalising Flows



discrete normalizing flow

$$x_0 = f_1^{\theta_1} \dots \circ f_{T-1}^{\theta_{T-1}} \circ f_T^{\theta_T}(x_T)$$

continuous normalizing flow

$$\frac{dx_t}{dt} = v_t^\theta(x_t), \quad x_t := f(x_T, t)$$

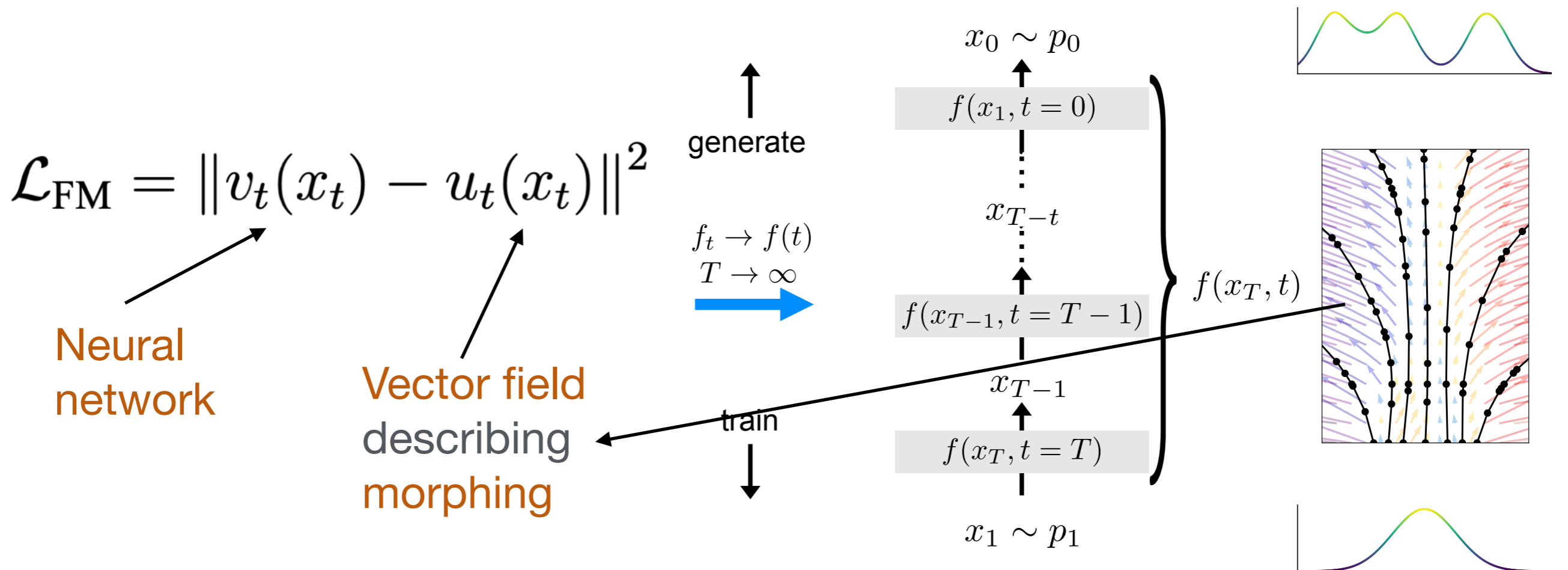
Training costly:

Solve ODE to evaluate  $p(x)$

Calculate trace for change in probability volume



# Flow Matching



Training costly:

Solve ODE to evaluate  $p(x)$

Calculate trace for change in probability volume

Instead, learn to approximate target vector field

# Flow Matching

How to get the vector field describing morphing?

Constrains:

At  $t=0$ : Input data

At  $t=1$ : Target latent distribution

Marginal path

conditional path

$$p_t(x) = \int dx_0 p_t(x|x_0)p_0(x_0)$$

Use mixture of conditional morphings

$$u_t(x) = \int dx_0 u_t(x|x_0) \frac{p_t(x|x_0)p_0(x_0)}{p_t(x)}$$

Would allow calculating corresponding vector field

$$\mathcal{L}_{\text{FM}} = \|v_\theta(x_t|t) - u_t(x_t|x_0)\|^2$$

Still expensive, let network only approximate conditional vector field

# Flow Matching

How to get the vector field describing morphing?

Constrains:

At  $t=0$ : Input data

At  $t=1$ : Target latent distribution

Use mixture of conditional morphings

But what are the **actual paths?**

Can use **Gaussians** with **linear interpolation**

$$p_t(x|x_0) = \mathcal{N}(x|\gamma_t, \sigma_t)$$

$$\gamma_t^{\text{FM}} = (1 - t)x_0,$$

$$\sigma_t^{\text{FM}} = \sigma_{\min} + (1 - \sigma_{\min})t$$

Vector field:

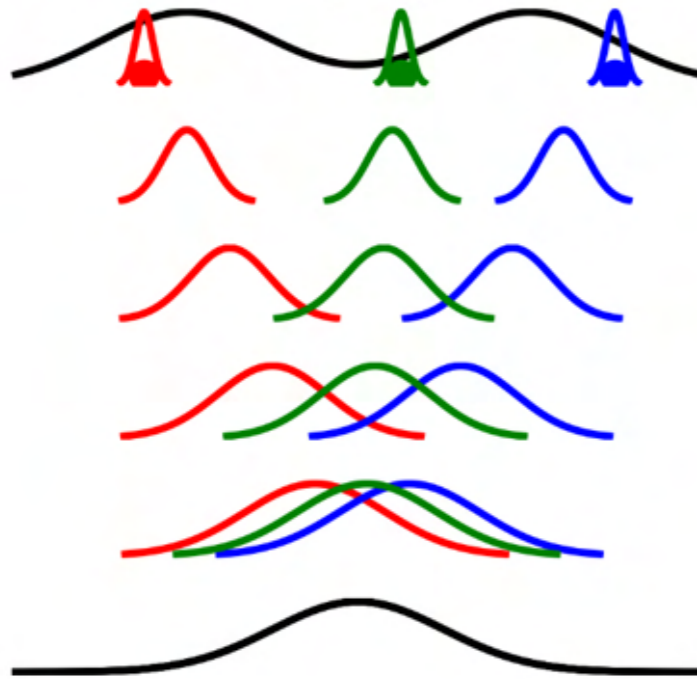
$$u_t(x_t|x_0) = \dot{\gamma}_t + \dot{\sigma}_t \epsilon = (1 - \sigma_{\min})\epsilon - x_0$$

Loss:

$$\mathcal{L}_{\text{FM}} = \|(v_{\theta}(x_t, t) - (1 - \sigma_{\min})\epsilon - x_0)\|^2$$

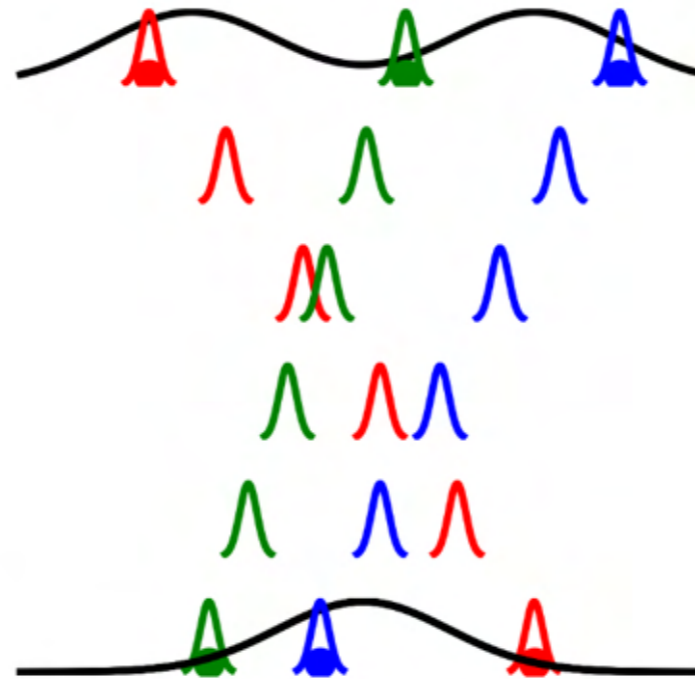
# Optimal Transport

Flow Matching (Lipman et al.)



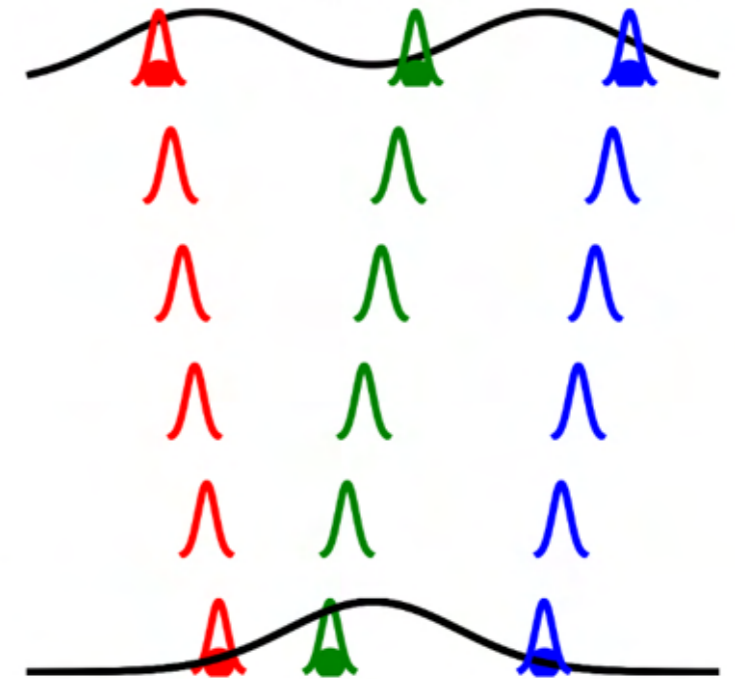
Can further improve paths by:

Conditional Flow Matching



Conditioning on start- and endpoint

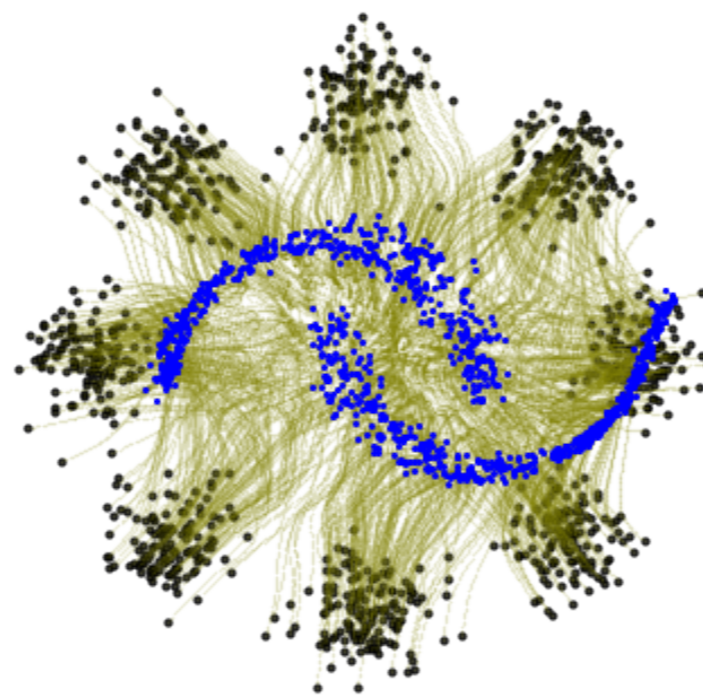
OT Conditional Flow Matching



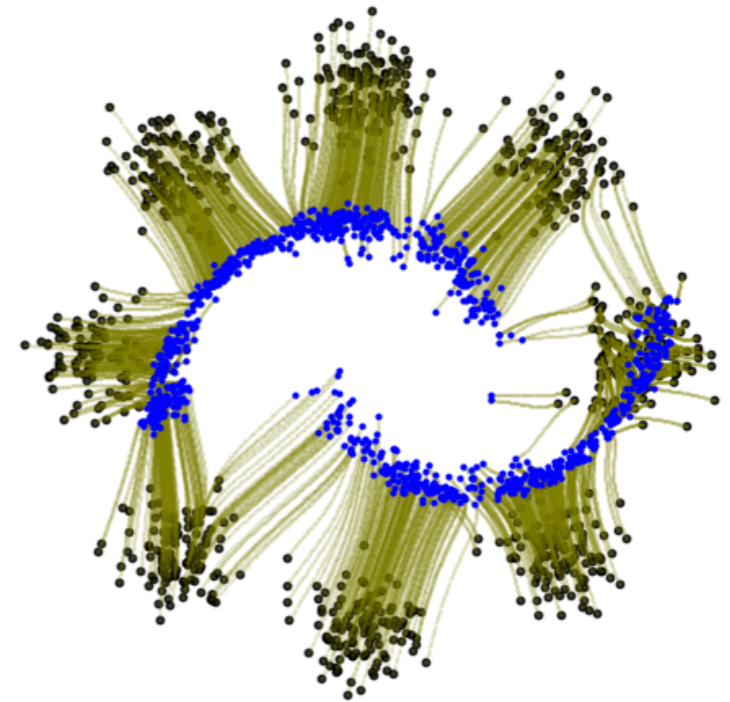
Adding optimal transport condition



# Optimal Transport

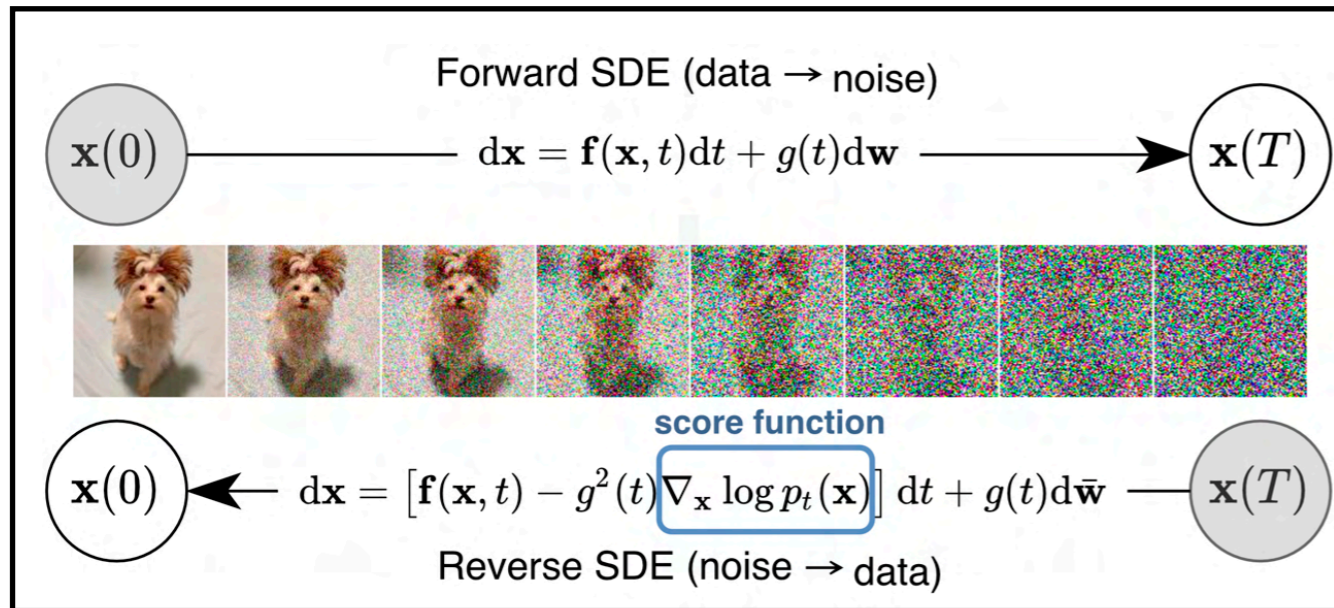


Conditioning on  
start- and endpoint

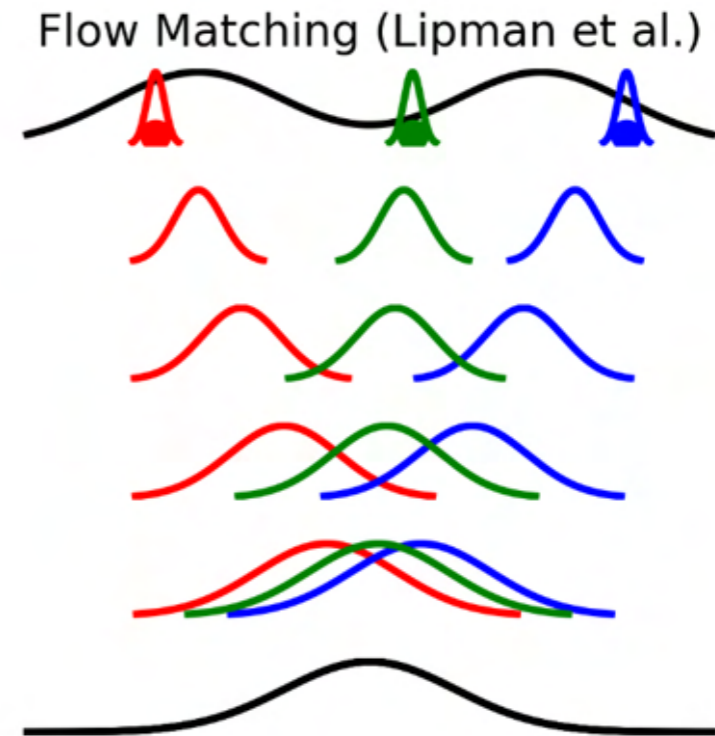


Adding optimal  
transport condition

# Comments



Score matching objective



Flow matching objective

Can show **equivalence for Gaussian probability paths**

Continuous Normalising Flow  
more general framework (other definition of paths)

# Comments

Close relation of CNF and diffusion models

Maturity:

Very recent, much in flux

Sample quality:

Very high

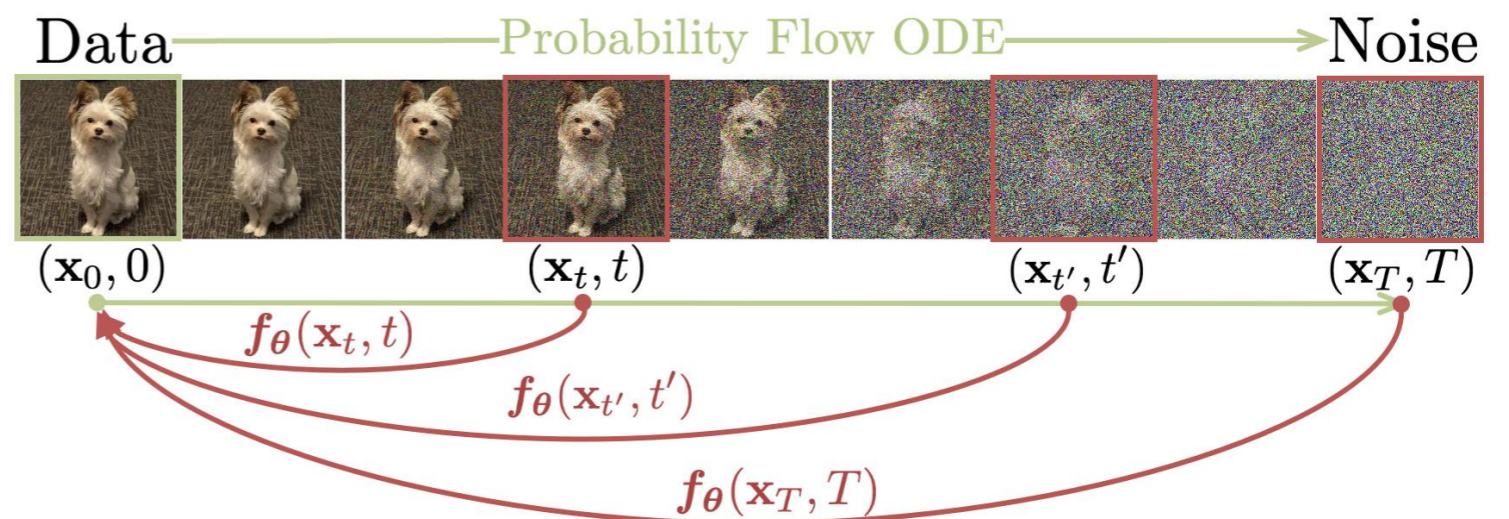
Training:

Stable

Sampling:

Expensive

(however: consistency distillation)



# Applications II



# Comments on Flows

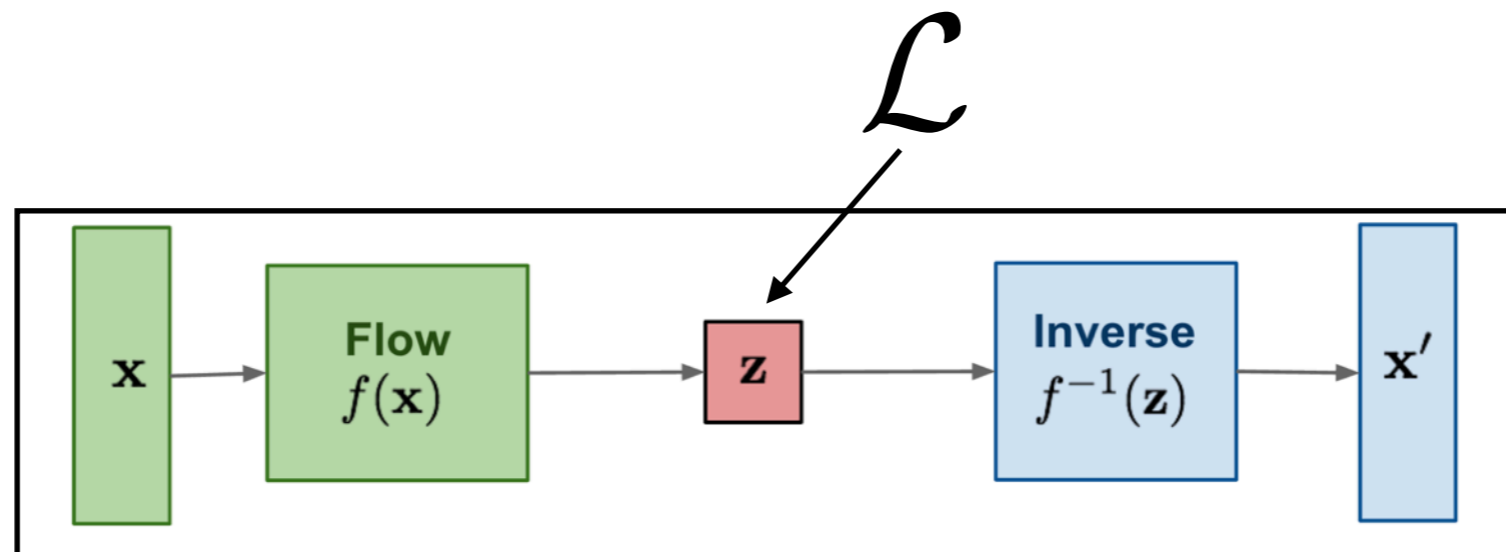
Exact learning of likelihood

→ Better generative fidelity

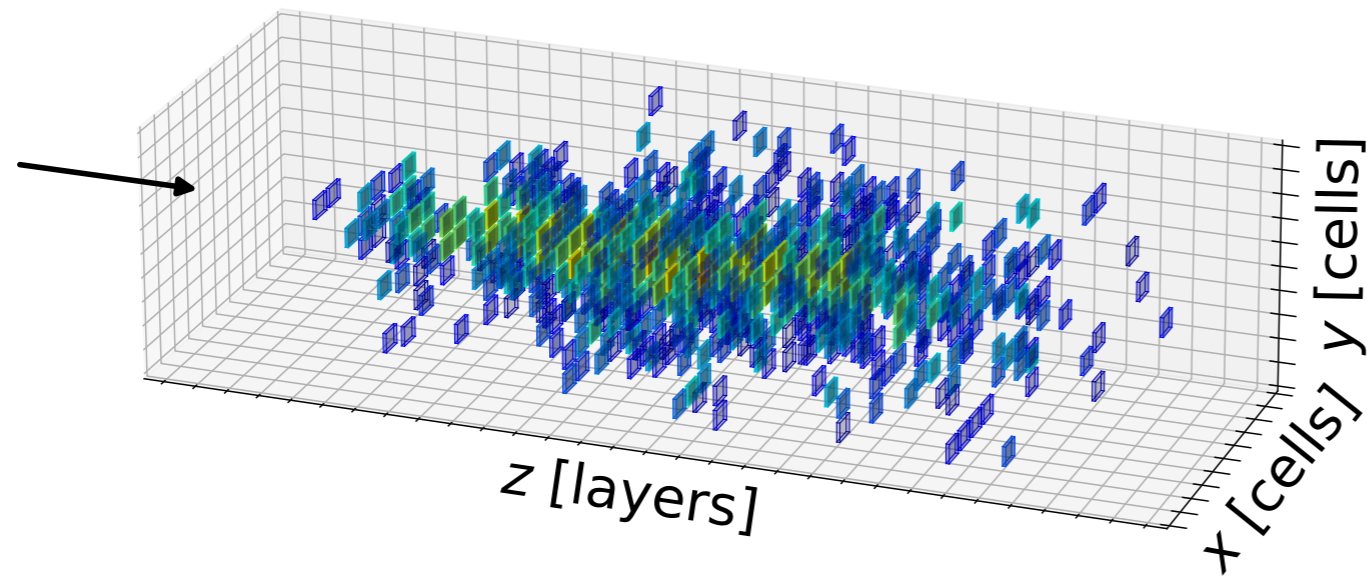
→ Can evaluate likelihood of data

More complex

→ Slower, choice of fast direction



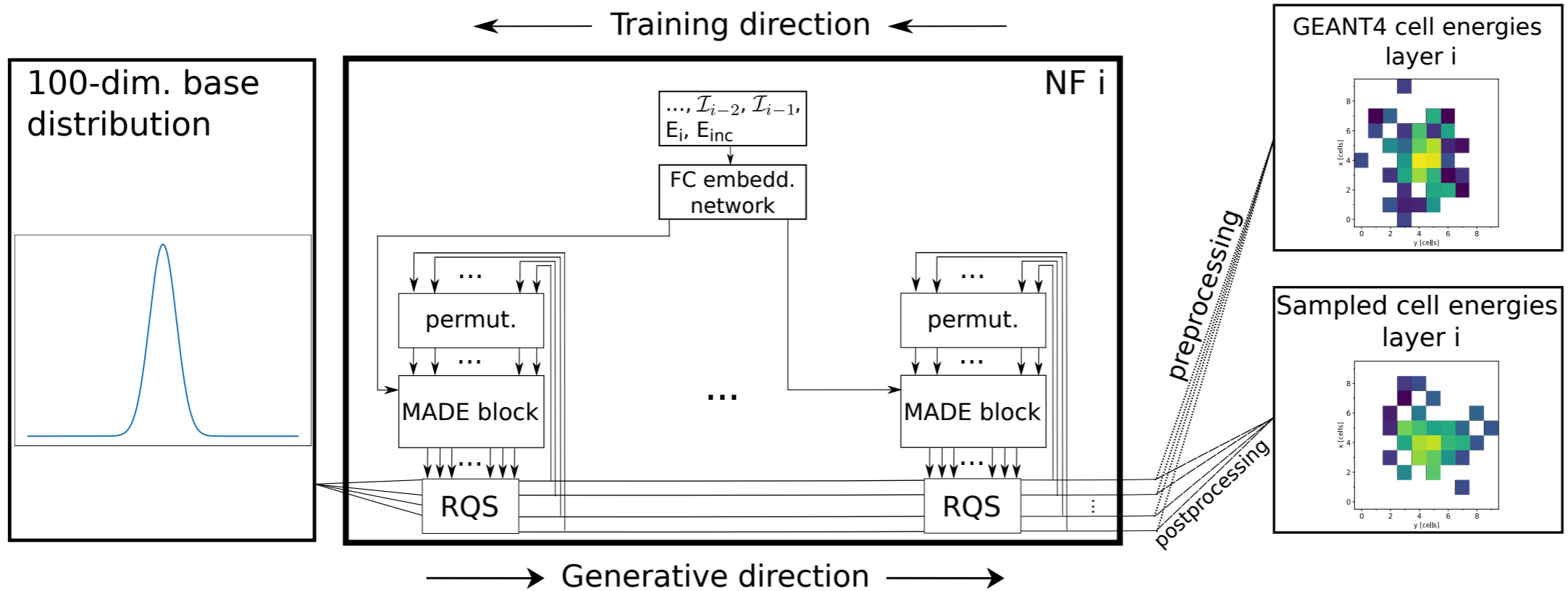
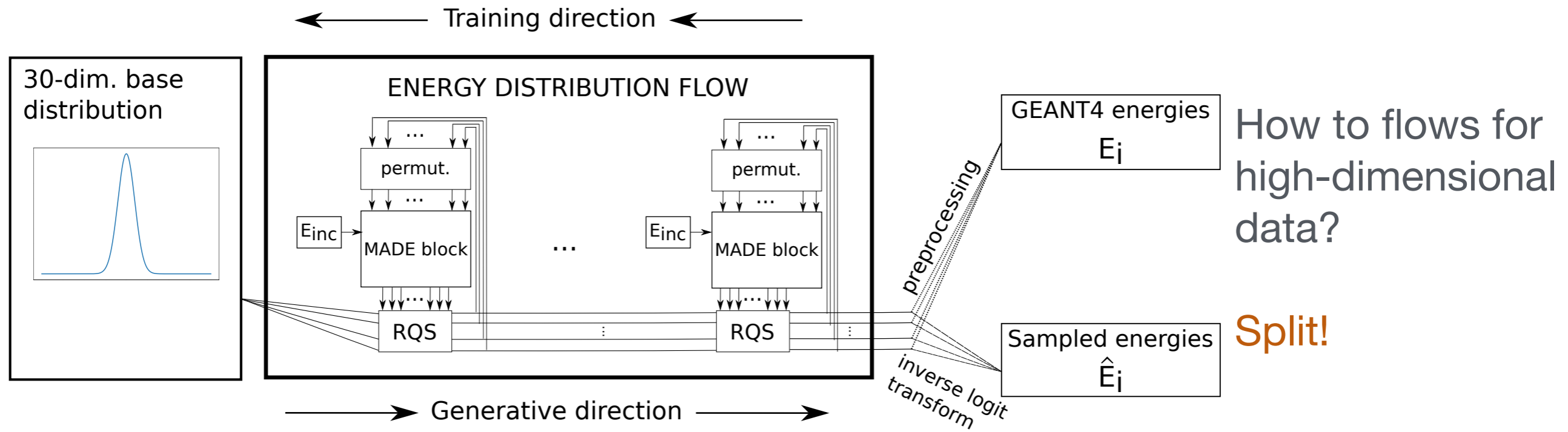
# Generative results II



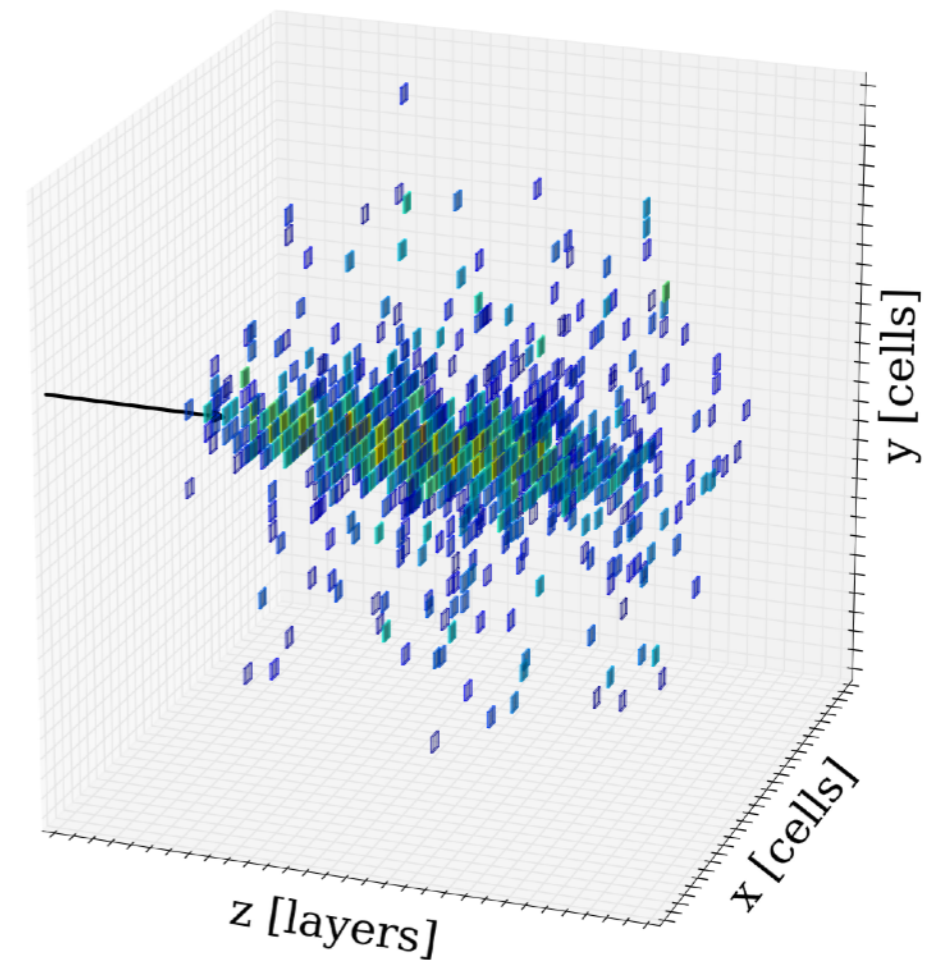
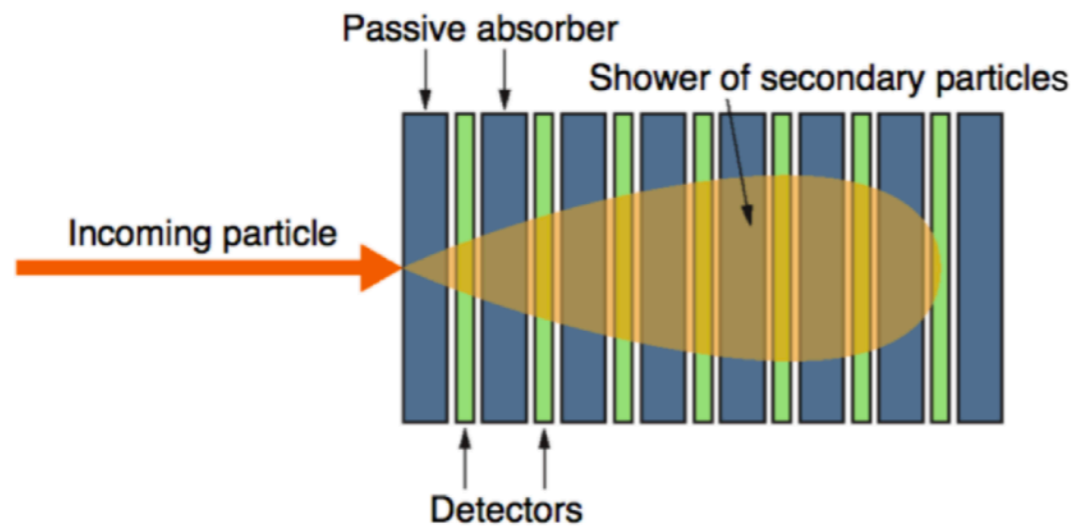
10x10 cells / layer  
30 layers

How to flows for  
high-dimensional  
data?

# Generative results II



# Simulation targets



How to represent?

Tabular data

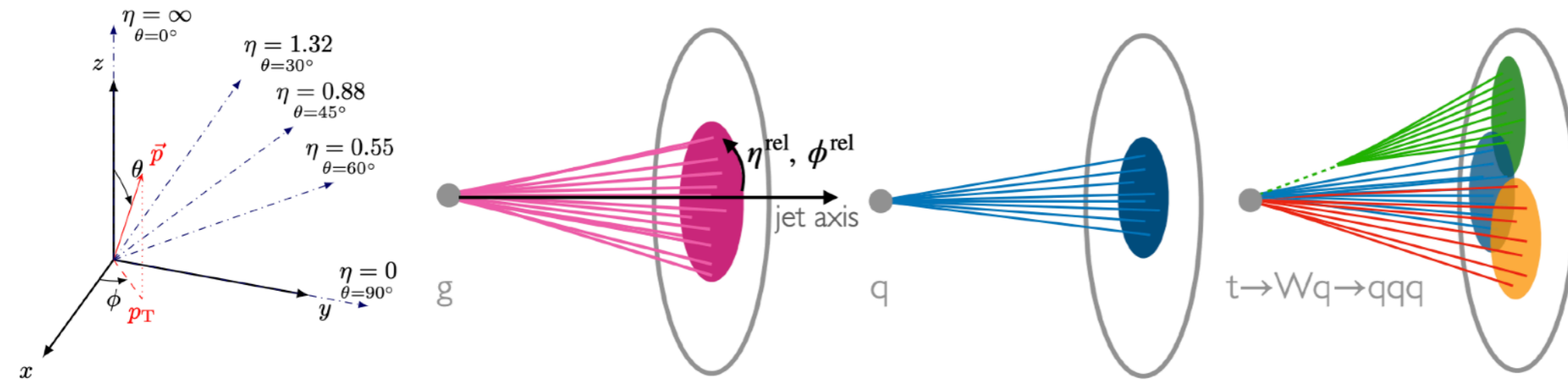
Fixed grid (voxels)

Limiting for high-dimensions (sparse data)

Point clouds / graphs



# Simulation targets



Before tackling showers in calorimeters:  
Look at jet constituents (JetNet data):  
3 features per constituents  
up to 30/150 constituents/jet

How to represent?

Tabular data

Fixed grid (voxels)

Limiting for high-dimensions (sparse data)

Point clouds / graphs

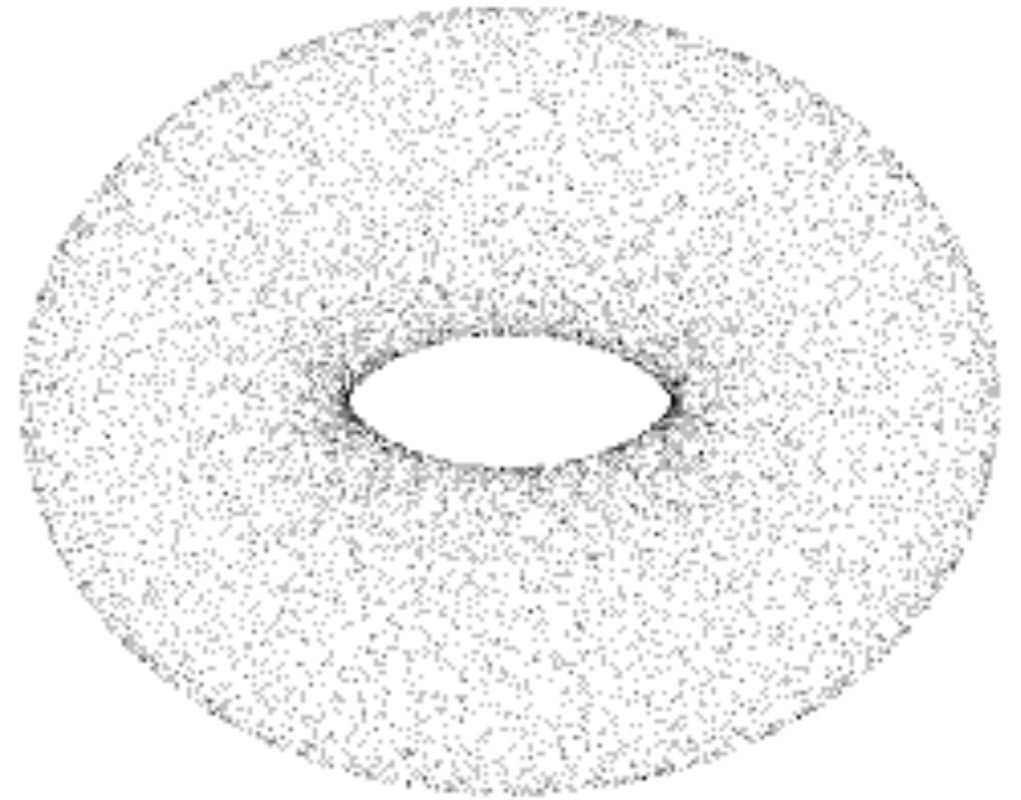
Why?

Useful stepping stone

In-situ background

# Point Clouds

- Example:  
*Sensors in a space*
  - Fixed grid vs arbitrary positions
  - Potential sparsity of data
- Permutation symmetry
- Can view as trivial graph



Total data  $\{x^j\}_{j=1 \dots N}$  <sup>Examples</sup>

with  $x^j = \{ \vec{p}_{11}^j, \dots, \vec{p}_{L(j)}^j \}$

and  $\vec{p}_i \in \mathbb{R}^D$

# Deep Sets

**Theorem 7** Let  $f : [0, 1]^M \rightarrow \mathbb{R}$  be a permutation invariant continuous function iff it has the representation

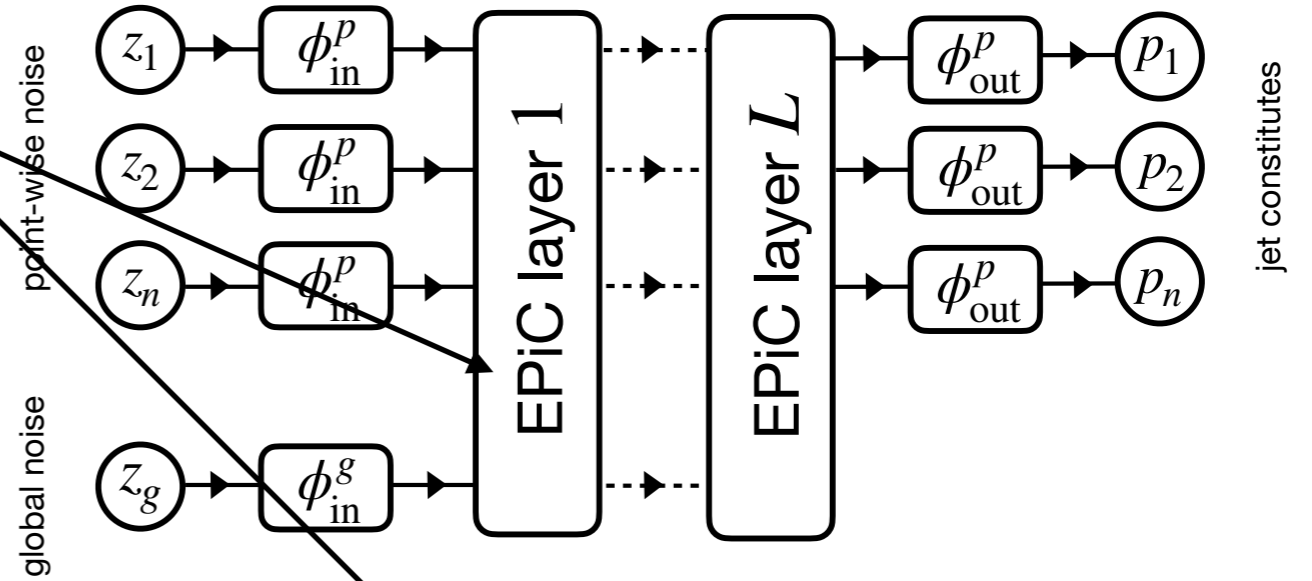
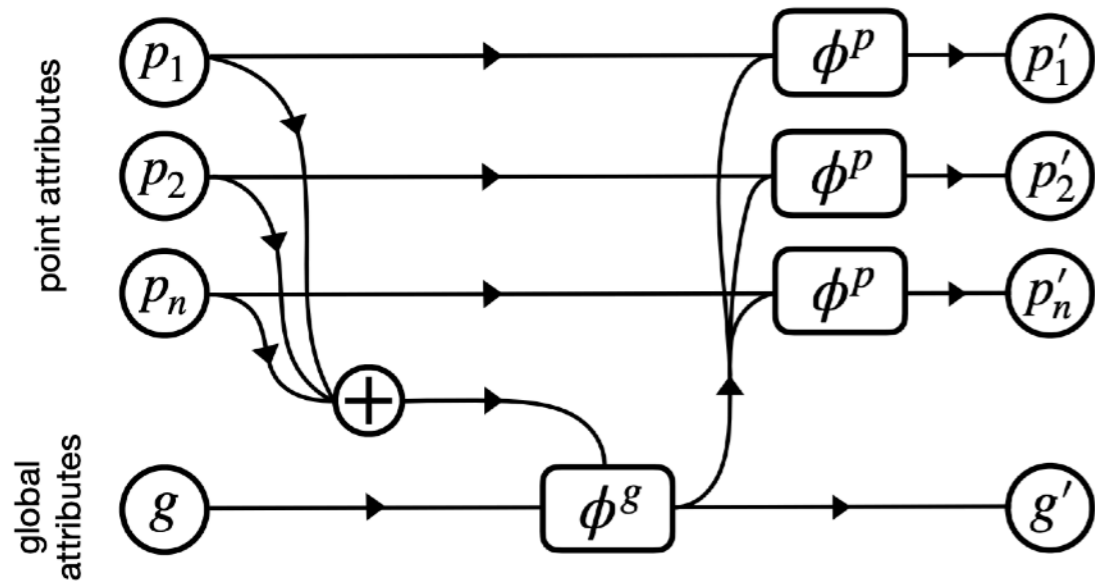
$$f(x_1, \dots, x_M) = \rho \left( \sum_{m=1}^M \phi(x_m) \right) \quad (18)$$

for some continuous outer and inner function  $\rho : \mathbb{R}^{M+1} \rightarrow \mathbb{R}$  and  $\phi : \mathbb{R} \rightarrow \mathbb{R}^{M+1}$  respectively. The inner function  $\phi$  is independent of the function  $f$ .

$$\begin{array}{l}
 x = \left. \begin{array}{l} \vec{p}_1 = \{r_1, \theta_1, \varphi_1, T_1\} \\ \vdots \\ \vec{p}_L = \{r_L, \theta_L, \varphi_L, T_L\} \end{array} \right\} \\
 f(x) = \rho \left( \sum_i \phi(\vec{p}_i) \right) \\
 \begin{array}{l}
 \rho: \mathbb{R}^{4 \times L} \rightarrow \mathbb{R}^1 \\
 \phi: \mathbb{R}^4 \rightarrow \mathbb{R}^s \\
 \Sigma: \mathbb{R}^{s \times L} \rightarrow \mathbb{R}^s \\
 \rho: \mathbb{R}^s \rightarrow \mathbb{R}^1
 \end{array}
 \end{array}$$

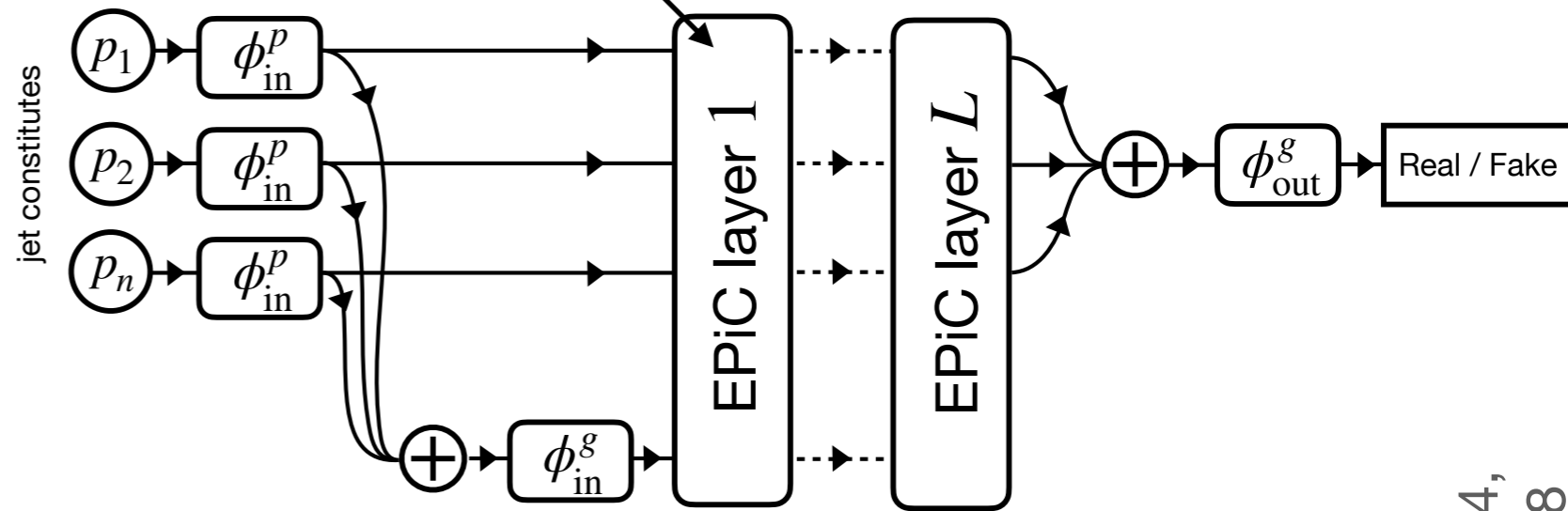
Latent Space Dimension

# How to GAN with it



(a) Generator

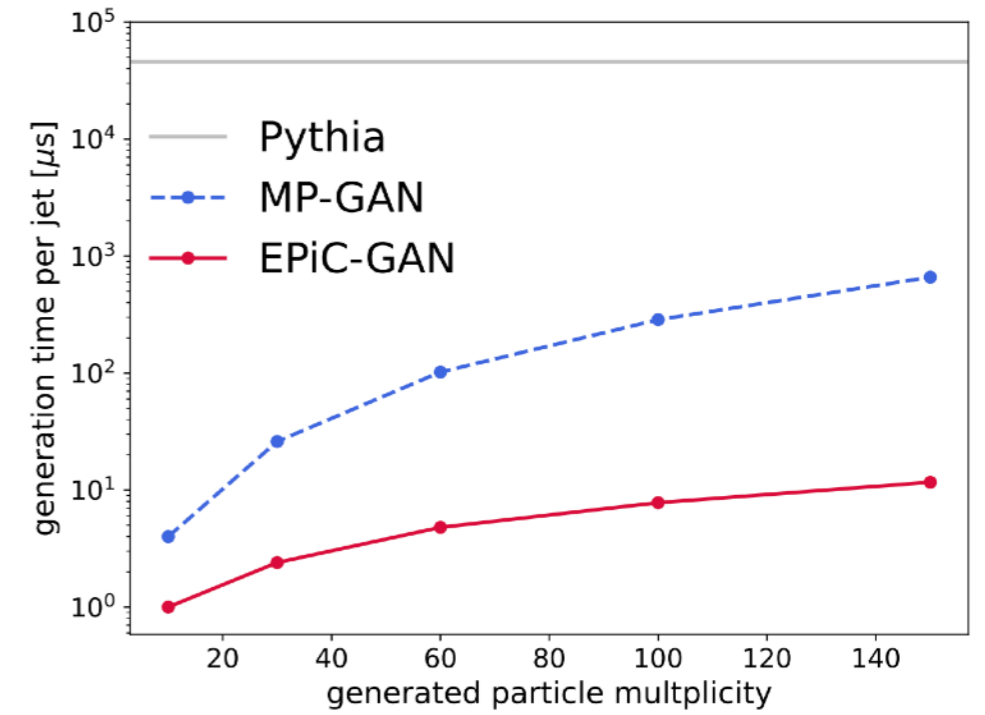
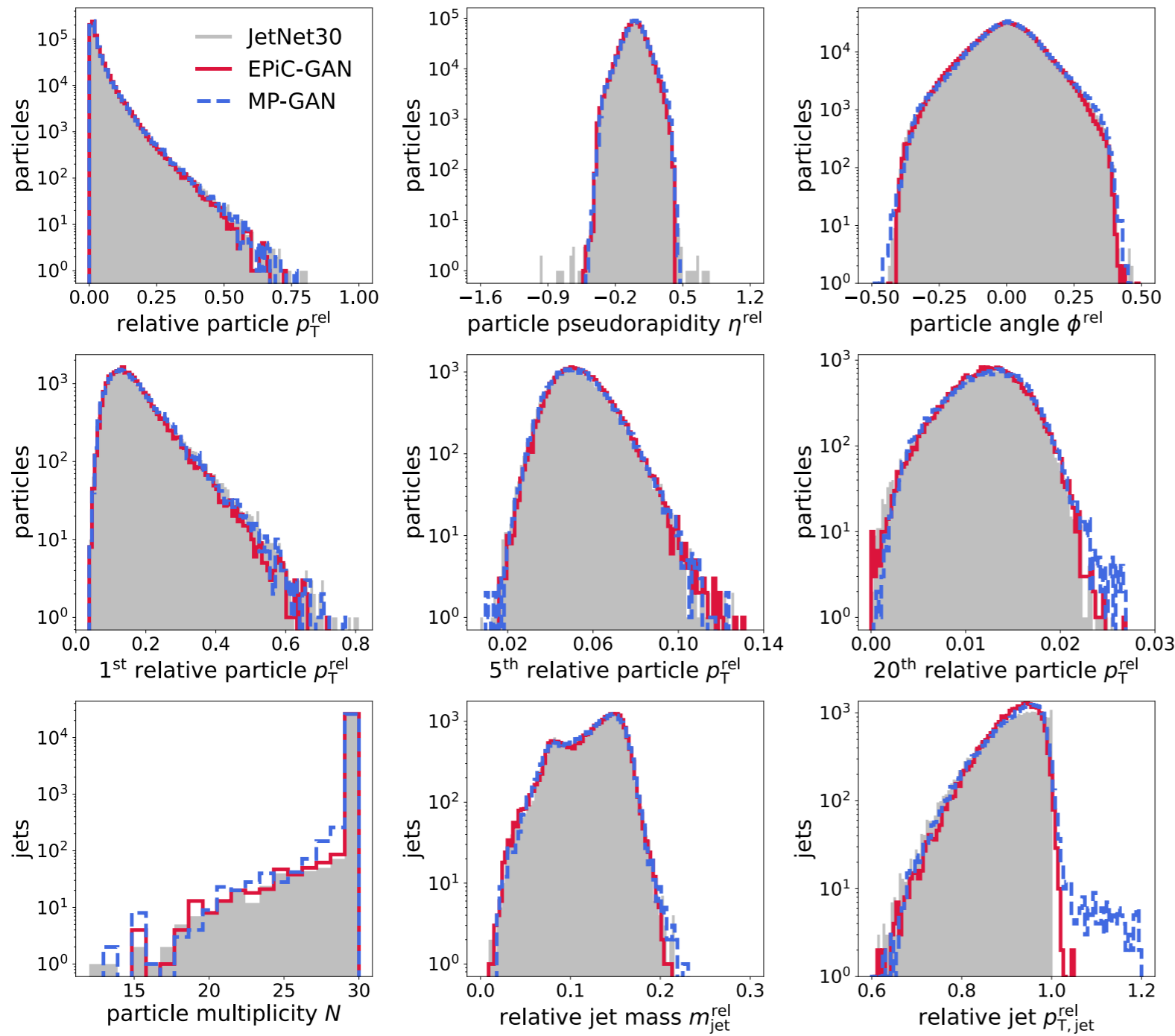
Add **permutation symmetry** to GAN architecture



(b) Discriminator

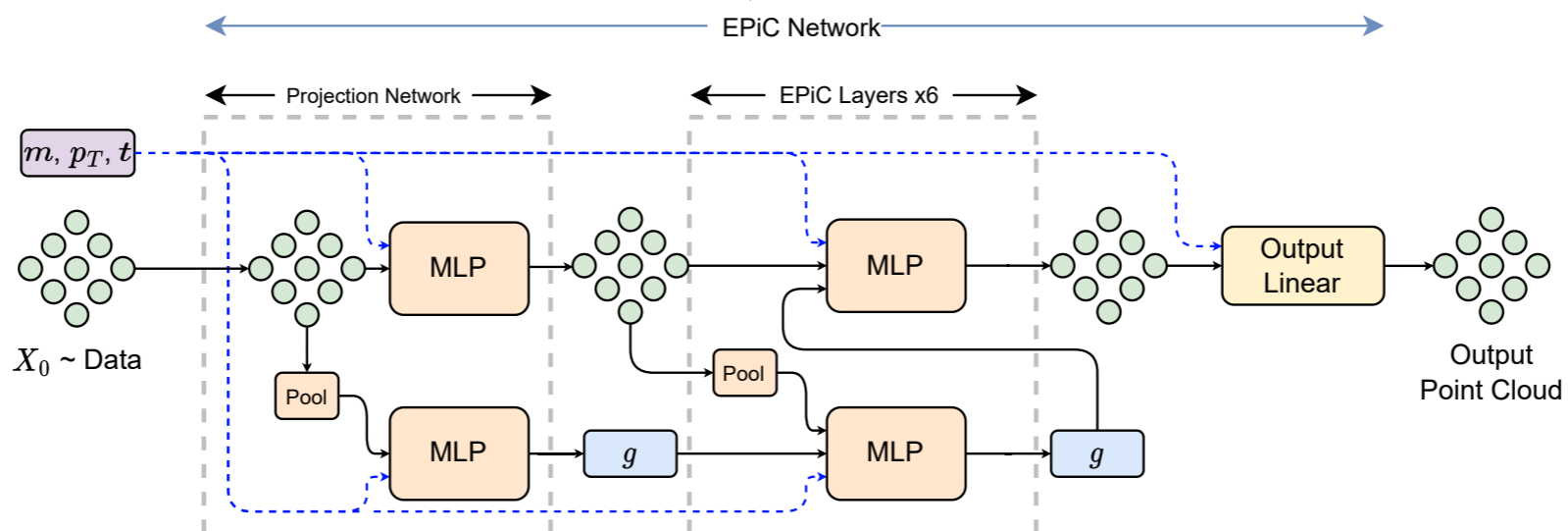
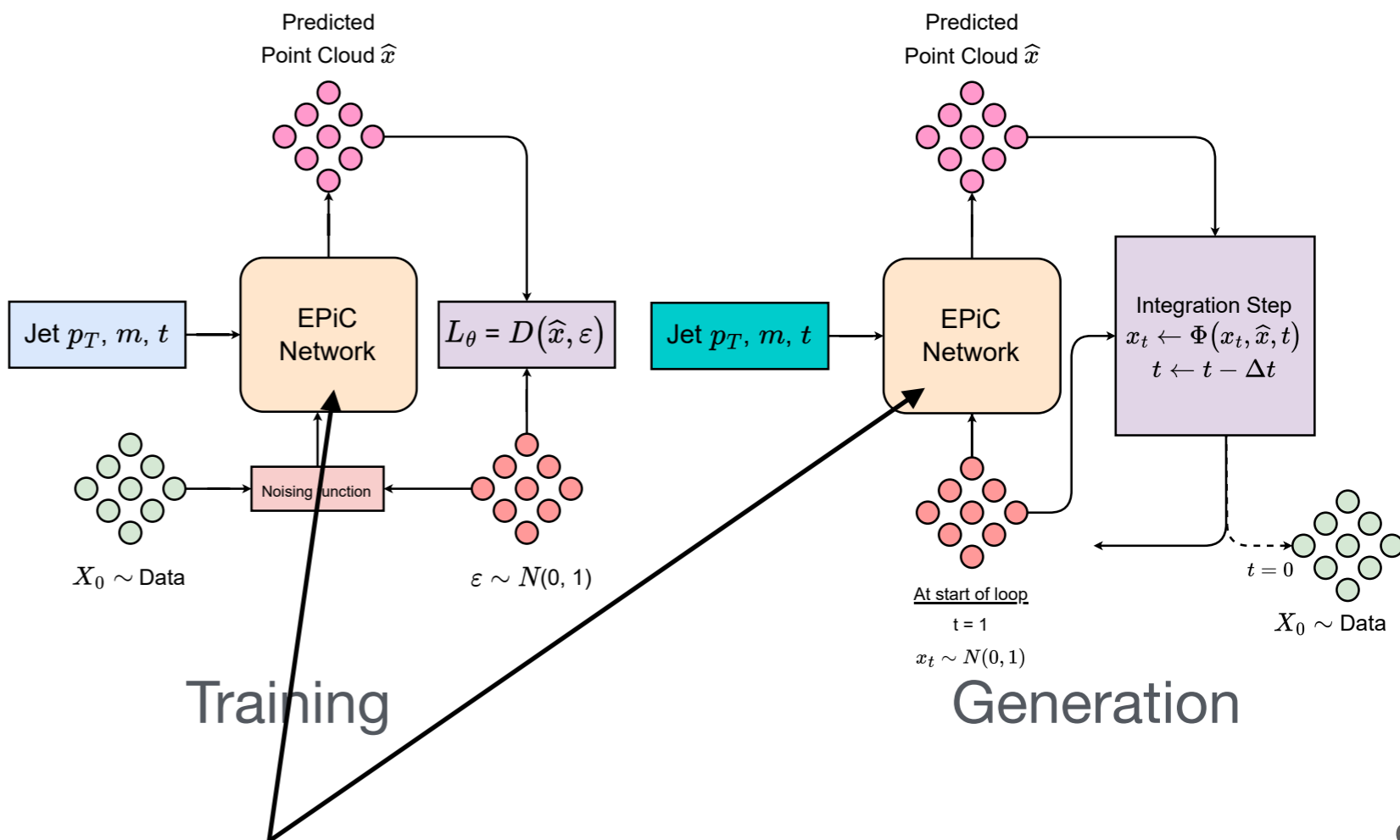
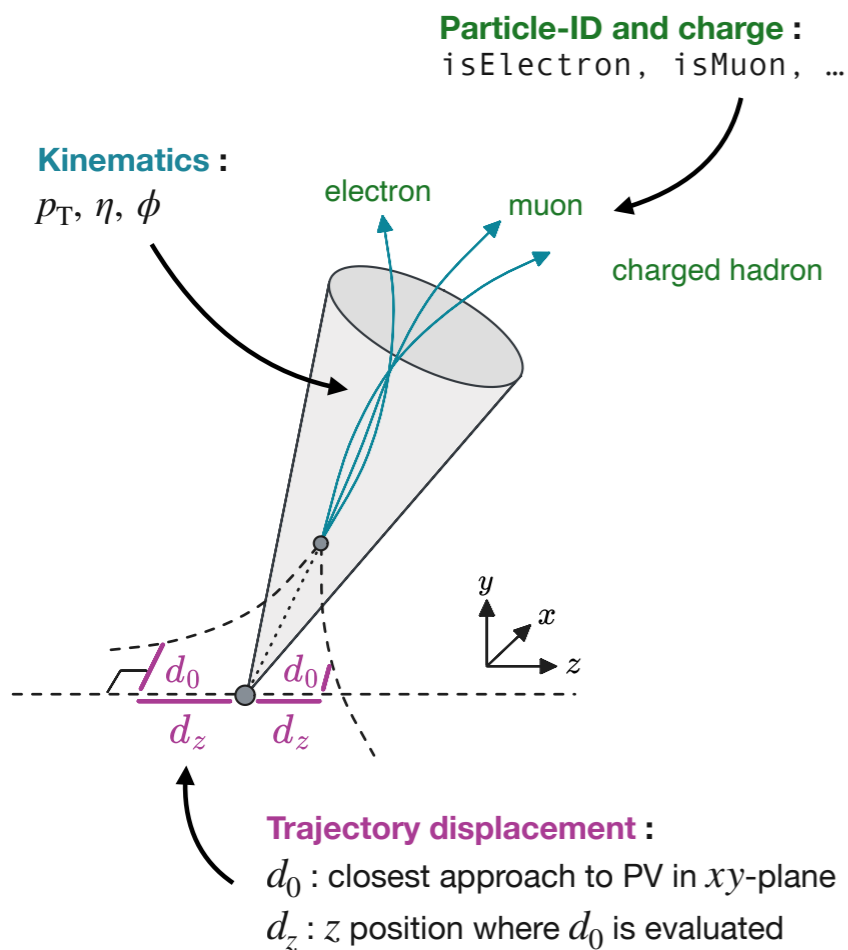


# Generative results III

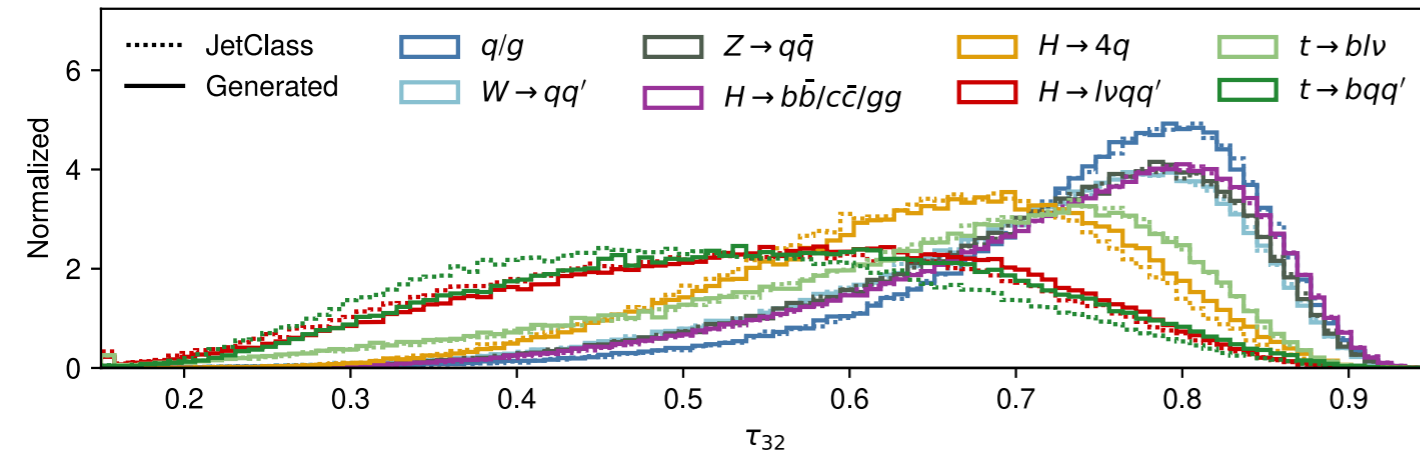


Point clouds/graph GANs  
successfully generate jet  
constituents

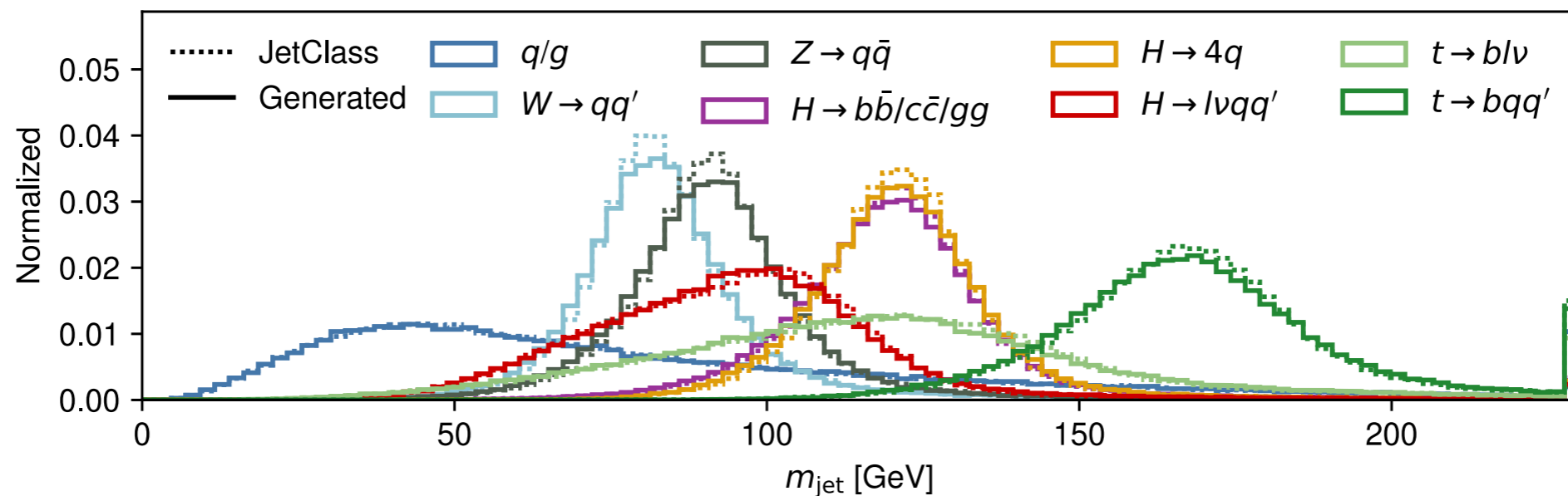
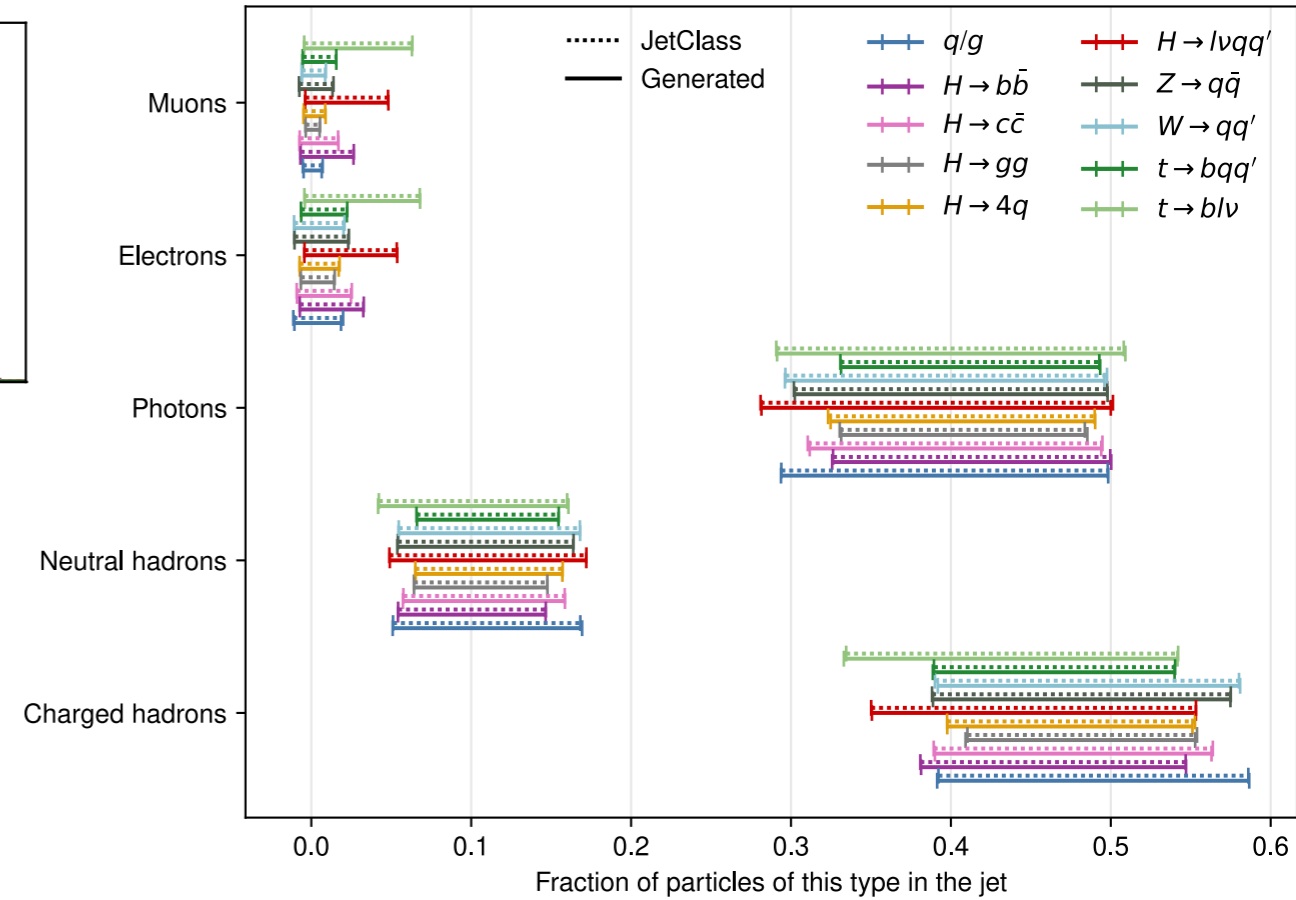
# Beyond kinematics



# Beyond kinematics



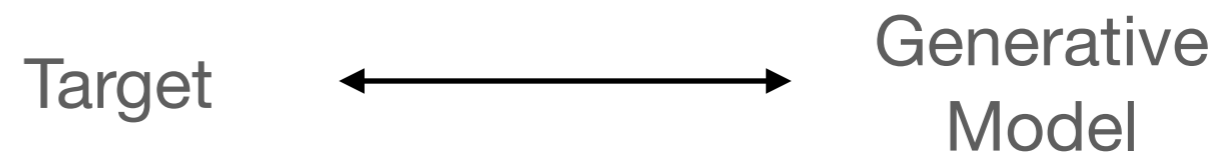
Good agreement across distributions



# Quality Metrics

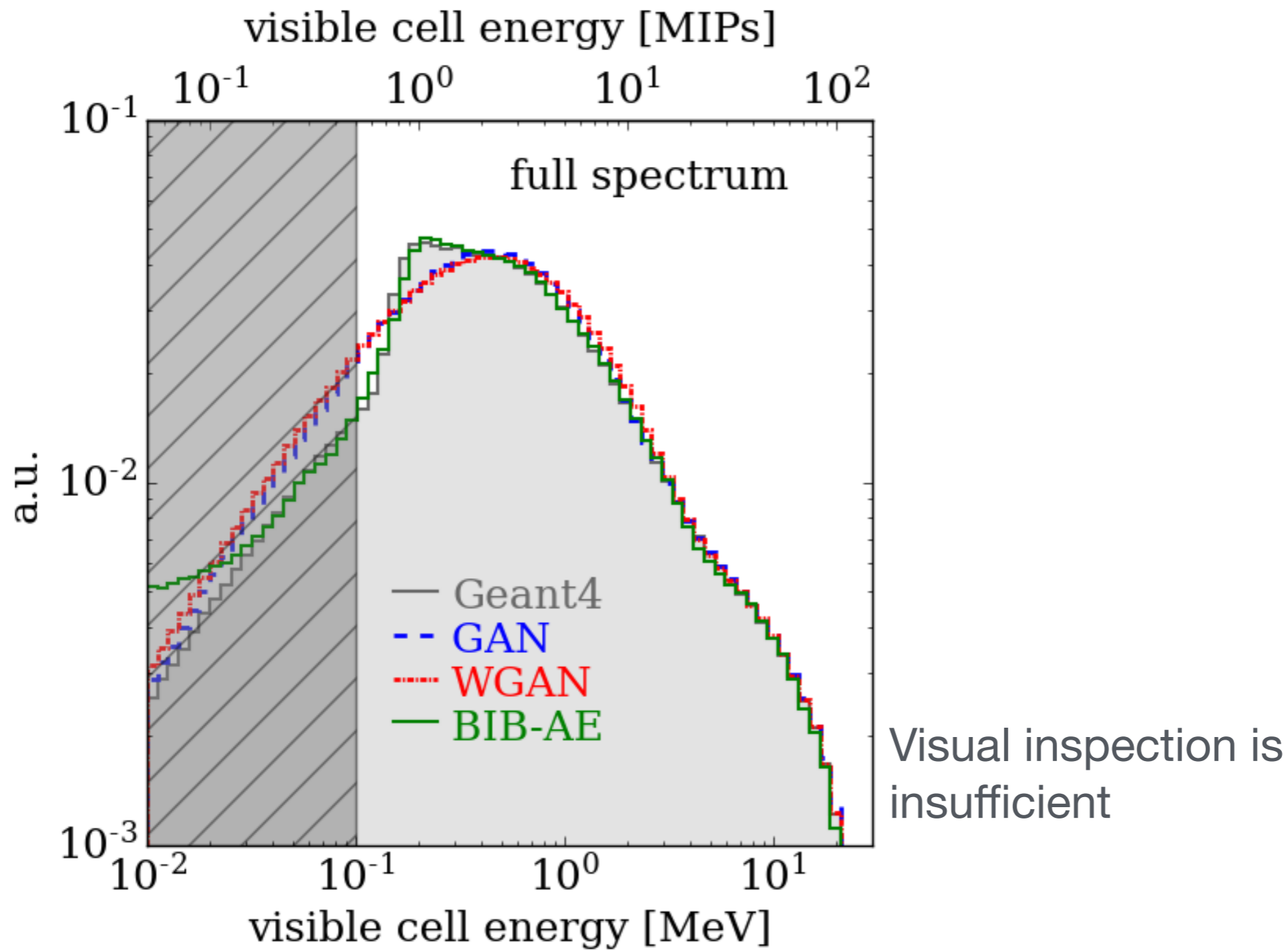


# Quality of simulation

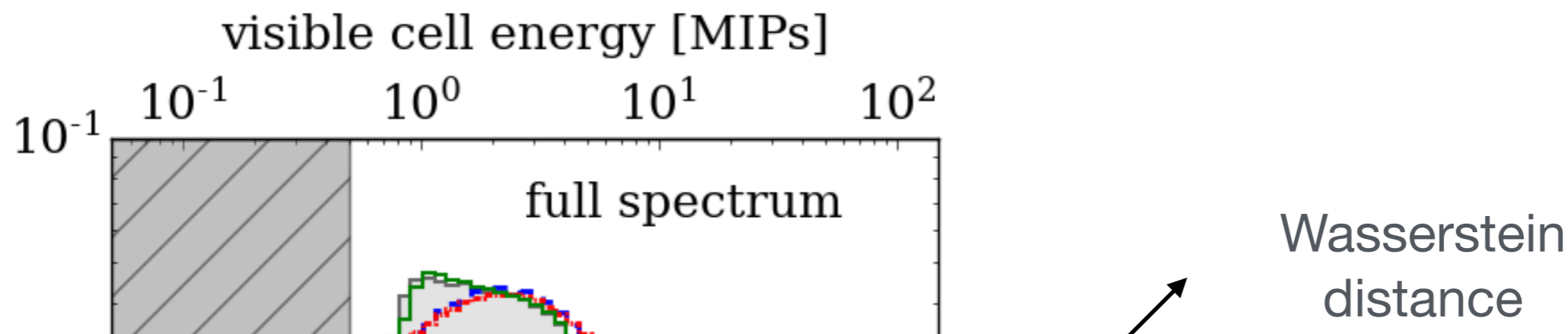


How well does the  
generative model  
describe the training  
data?

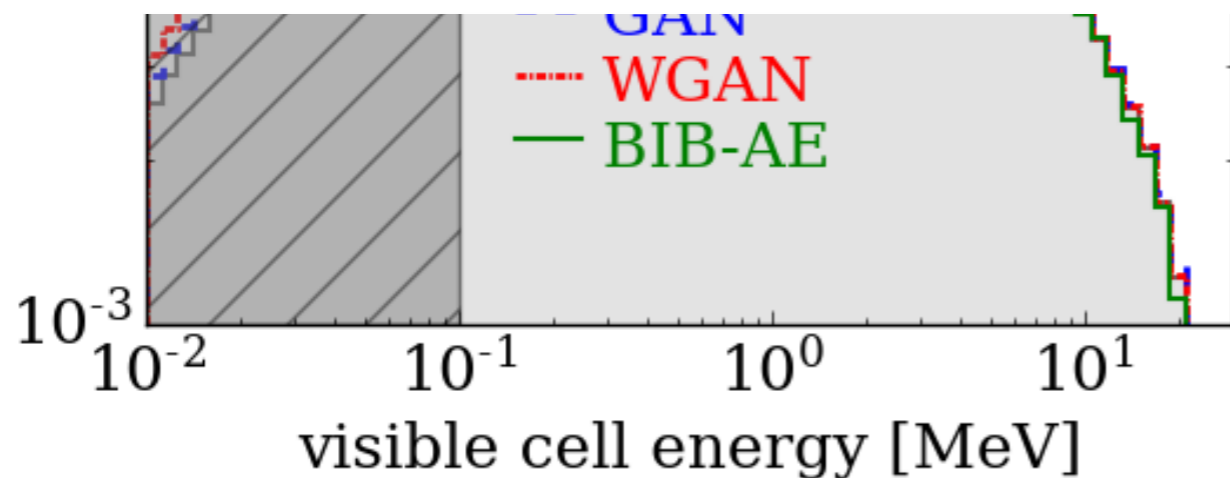
# One-dimensional metrics



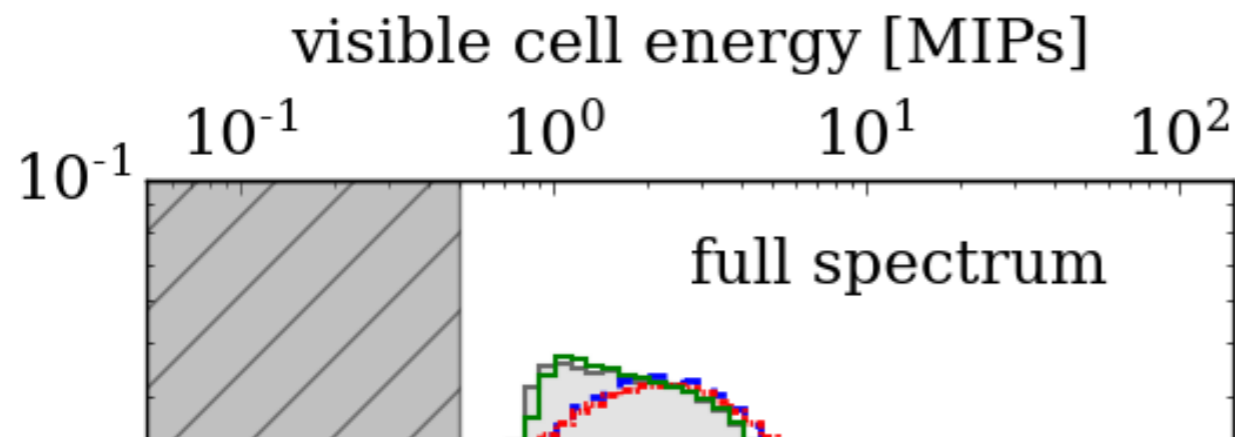
# One-dimensional metrics



Simulator	$W_1^{N_{\text{hits}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{vis}}/E_{\text{inc}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{cell}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{long}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{radial}}}$ ( $\times 10^{-3}$ )	$W_1^{m_{1,X}}$ ( $\times 10^{-3}$ )	$W_1^{m_{1,Y}}$ ( $\times 10^{-3}$ )	$W_1^{m_{1,Z}}$ ( $\times 10^{-3}$ )
GEANT4	0.7 ± 0.2	0.8 ± 0.2	0.9 ± 0.4	0.7 ± 0.8	0.7 ± 0.1	0.9 ± 0.1	1.1 ± 0.3	0.9 ± 0.3
CALOCLOUDS	<b>2.5 ± 0.3</b>	11.4 ± 0.4	15.9 ± 0.7	<b>2.0 ± 1.3</b>	38.8 ± 1.4	4.0 ± 0.4	8.7 ± 0.3	1.4 ± 0.5
CALOCLOUDS II	3.6 ± 0.5	26.4 ± 0.4	<b>15.3 ± 0.6</b>	3.7 ± 1.6	11.6 ± 1.5	<b>2.4 ± 0.4</b>	<b>7.6 ± 0.2</b>	3.9 ± 0.4
CALOCLOUDS II (CM)	6.1 ± 0.7	<b>9.8 ± 0.5</b>	16.0 ± 0.7	<b>2.0 ± 1.4</b>	<b>8.3 ± 1.9</b>	3.0 ± 0.4	9.5 ± 0.6	<b>1.2 ± 0.5</b>



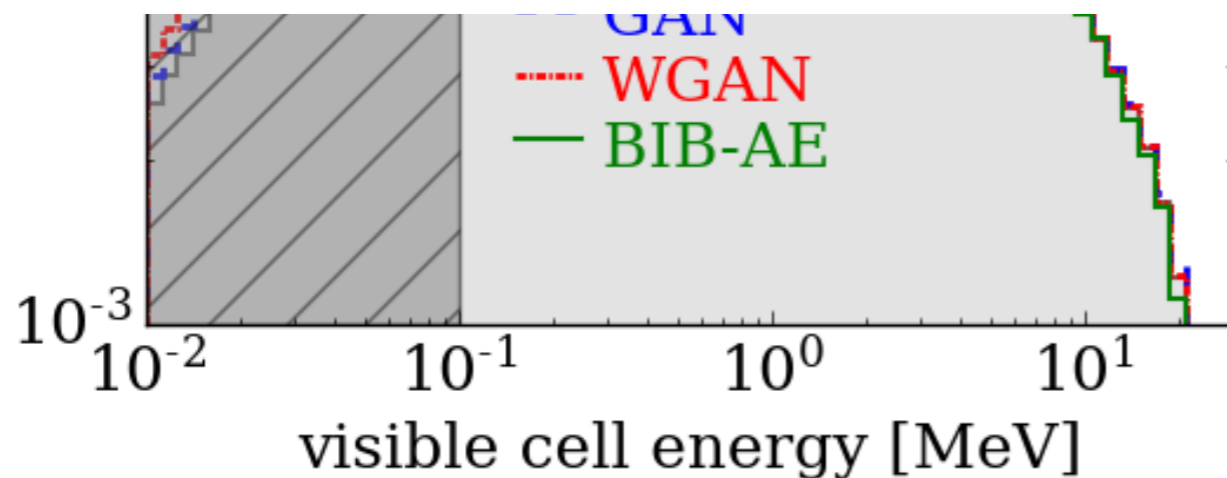
# One-dimensional metrics



Wasserstein distance

Simulator	$W_1^{N_{\text{hits}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{vis}}/E_{\text{inc}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{cell}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{long}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{radial}}}$ ( $\times 10^{-3}$ )	$W_1^{m_1, X}$ ( $\times 10^{-3}$ )	$W_1^{m_1, Y}$ ( $\times 10^{-3}$ )	$W_1^{m_1, Z}$ ( $\times 10^{-3}$ )
GEANT4	0.7 ± 0.2	0.8 ± 0.2	0.9 ± 0.4	0.7 ± 0.8	0.7 ± 0.1	0.9 ± 0.1	1.1 ± 0.3	0.9 ± 0.3
CALOCLOUDS	2.5 ± 0.3	11.4 ± 0.4	15.9 ± 0.7	2.0 ± 1.3	38.8 ± 1.4	4.0 ± 0.4	8.7 ± 0.3	1.4 ± 0.5
CALOCLOUDS II	3.6 ± 0.5	26.4 ± 0.4	15.3 ± 0.6	3.7 ± 1.6	11.6 ± 1.5	2.4 ± 0.4	7.6 ± 0.2	3.9 ± 0.4
CALOCLOUDS II (CM)	6.1 ± 0.7	9.8 ± 0.5	16.0 ± 0.7	2.0 ± 1.4	8.3 ± 1.9	3.0 ± 0.4	9.5 ± 0.6	1.2 ± 0.5

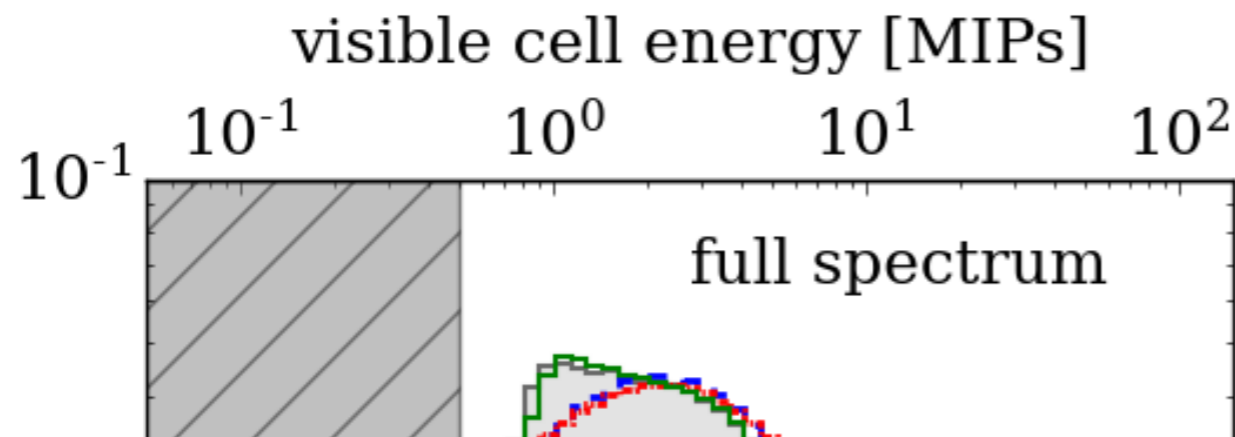
a.u.



Floor is given by sample-size effects of GEANT4 vs itself

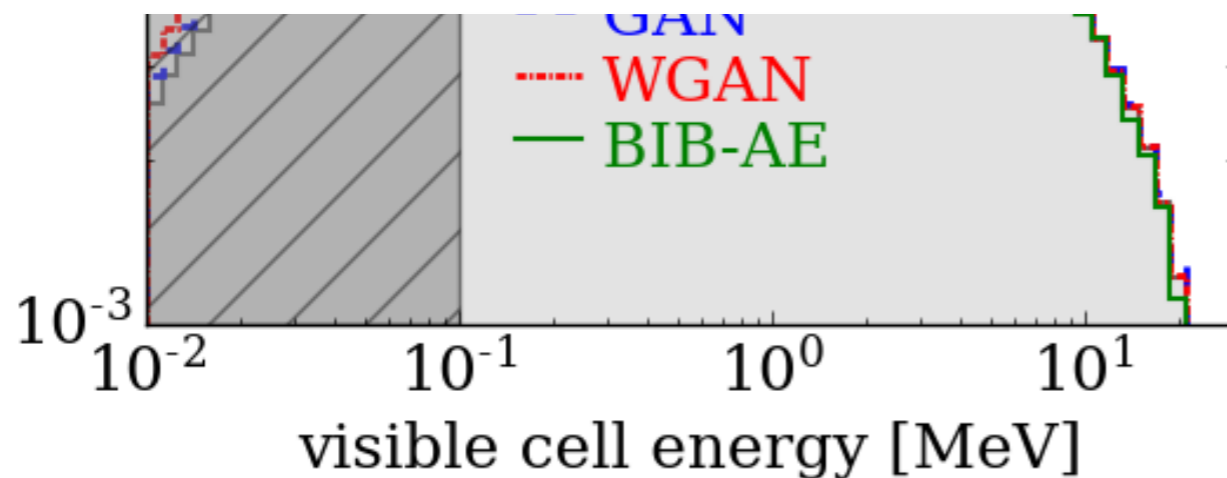


# One-dimensional metrics



Wasserstein distance

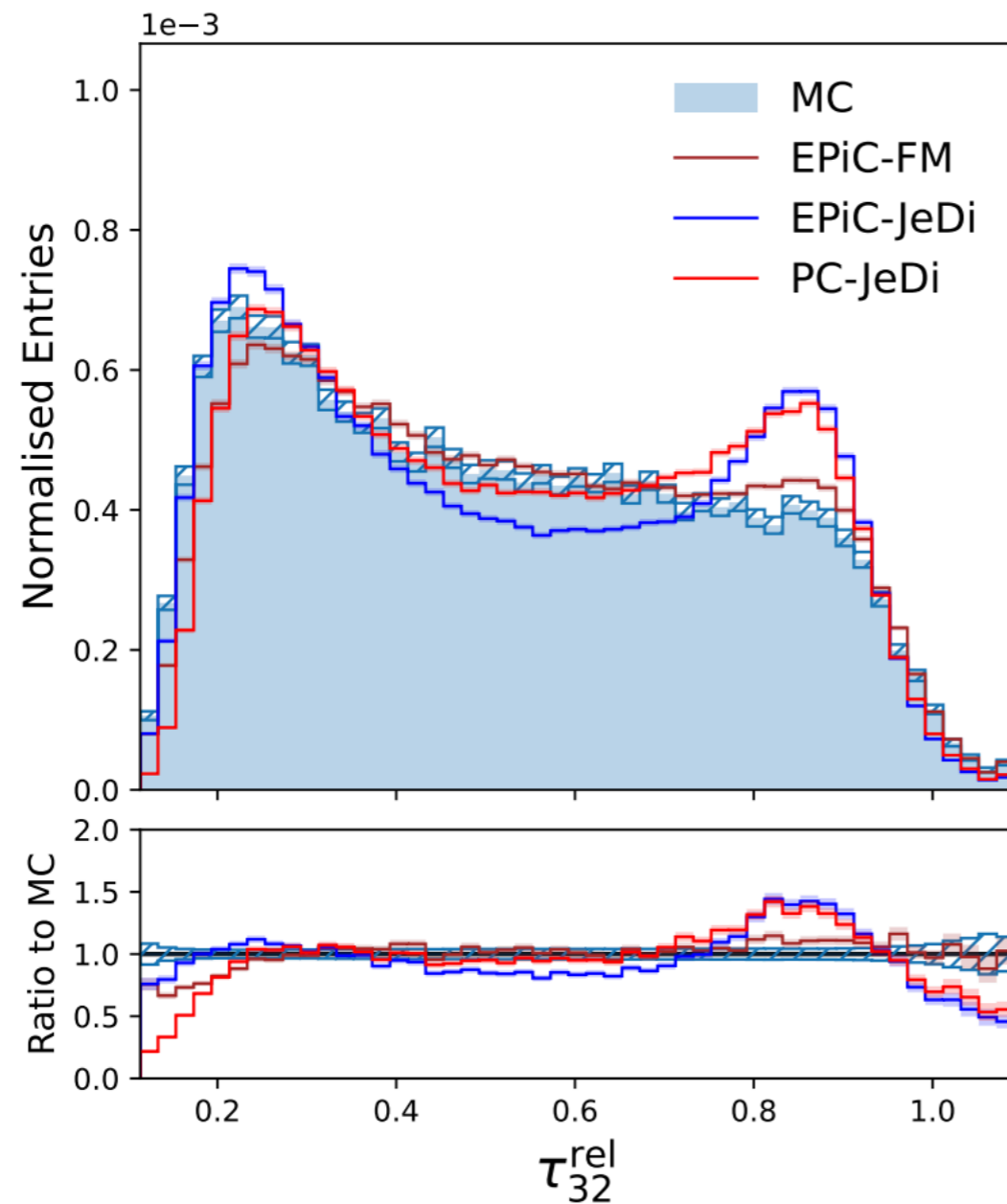
Simulator	$W_1^{N_{\text{hits}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{vis}}/E_{\text{inc}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{cell}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{long}}}$ ( $\times 10^{-3}$ )	$W_1^{E_{\text{radial}}}$ ( $\times 10^{-3}$ )	$W_1^{m_1, X}$ ( $\times 10^{-3}$ )	$W_1^{m_1, Y}$ ( $\times 10^{-3}$ )	$W_1^{m_1, Z}$ ( $\times 10^{-3}$ )
GEANT4	0.7 ± 0.2	0.8 ± 0.2	0.9 ± 0.4	0.7 ± 0.8	0.7 ± 0.1	0.9 ± 0.1	1.1 ± 0.3	0.9 ± 0.3
CALOCLOUDS	<b>2.5 ± 0.3</b>	11.4 ± 0.4	15.9 ± 0.7	<b>2.0 ± 1.3</b>	38.8 ± 1.4	4.0 ± 0.4	8.7 ± 0.3	1.4 ± 0.5
CALOCLOUDS II	3.6 ± 0.5	26.4 ± 0.4	<b>15.3 ± 0.6</b>	3.7 ± 1.6	11.6 ± 1.5	<b>2.4 ± 0.4</b>	<b>7.6 ± 0.2</b>	3.9 ± 0.4
CALOCLOUDS II (CM)	6.1 ± 0.7	<b>9.8 ± 0.5</b>	16.0 ± 0.7	<b>2.0 ± 1.4</b>	<b>8.3 ± 1.9</b>	3.0 ± 0.4	9.5 ± 0.6	<b>1.2 ± 0.5</b>



Floor is given by sample-size effects of GEANT4 vs itself

Uncertainty is standard deviation of 10 independent samples

# One-dimensional metrics



Wasserstein distance

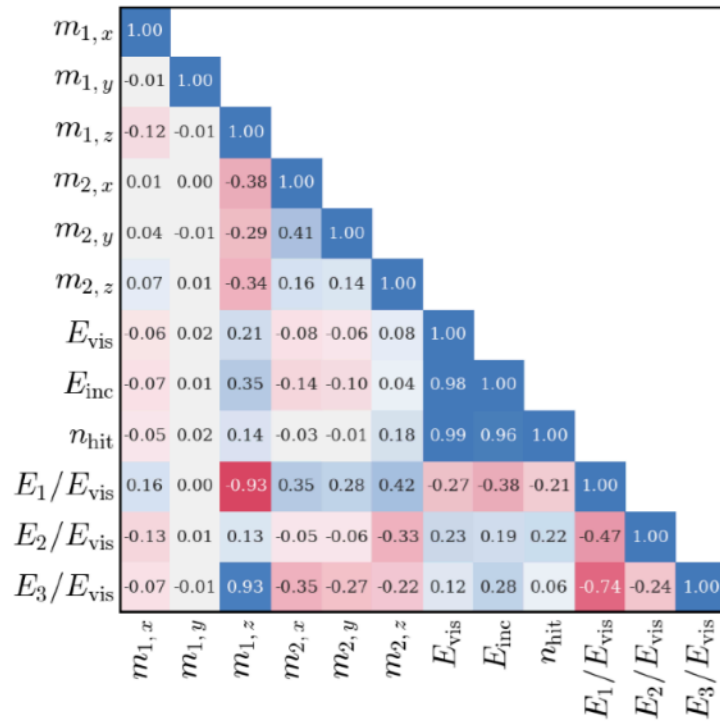
More robust, well defined also for non-overlapping distributions

Kullback-Leibler divergence

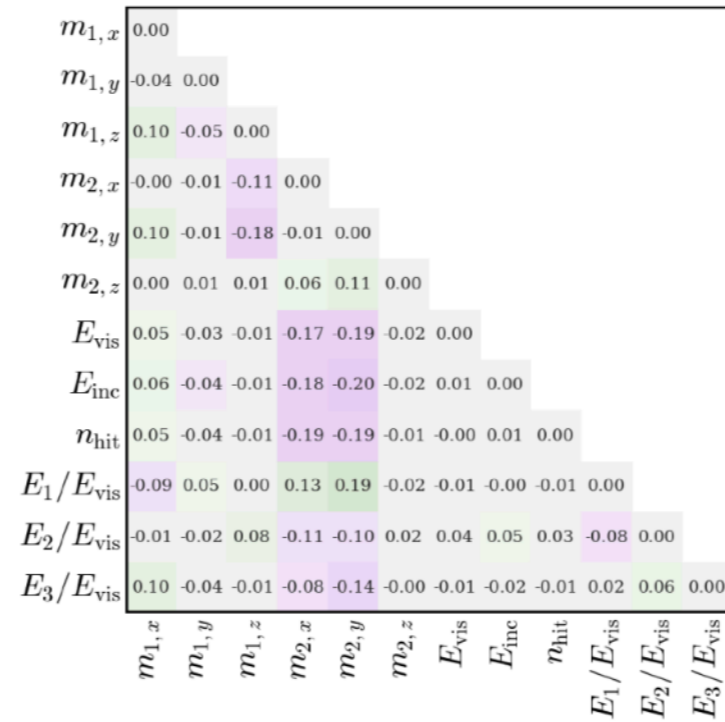
More intuitive for shape-mismatching

# Two-dimensional metrics

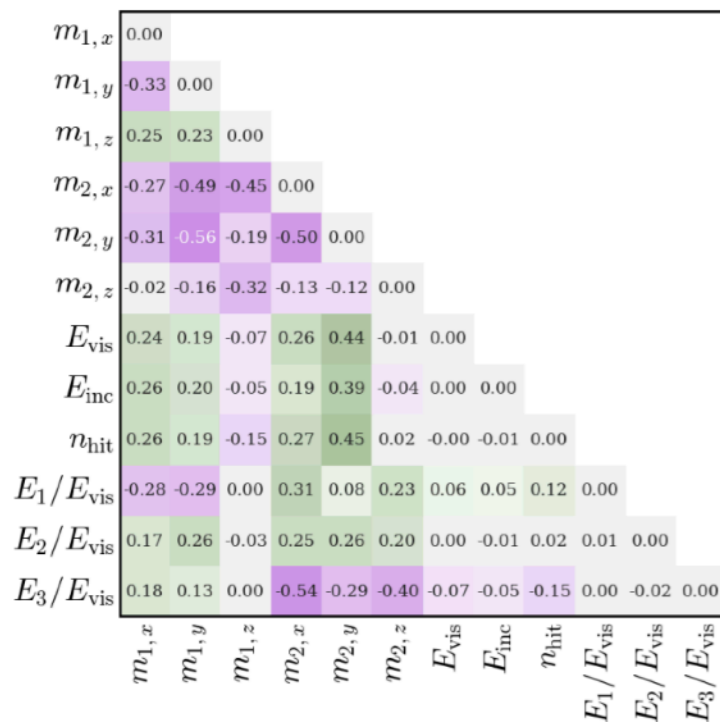
Geant4



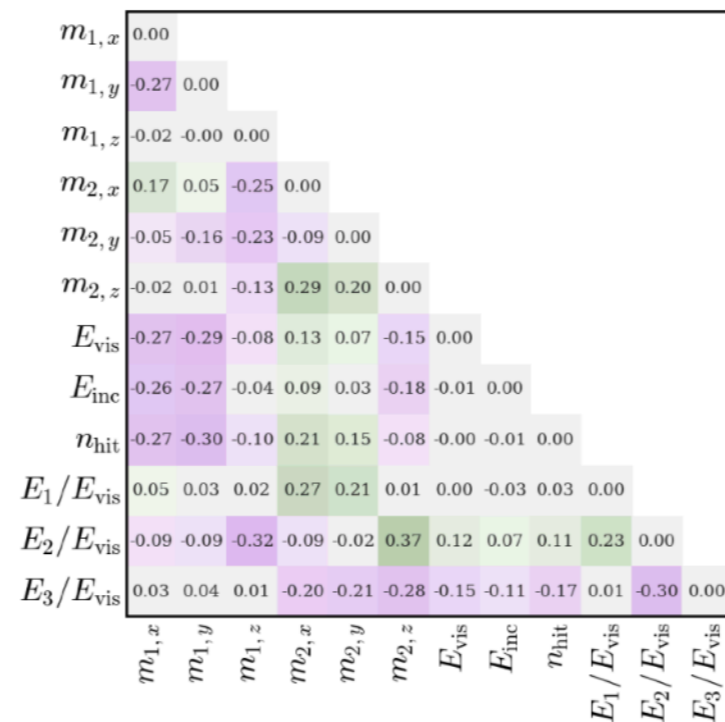
Geant4 - GAN



Geant4 - WGAN



Geant4 - BIB-AE PP



Pair-wise correlations contain more information

# Multi-dimensional metrics

# Showers per simulator	AUC GEANT4 vs L2LFlows	AUC GEANT4 vs BIB-AE
95k	$0.8518 \pm 0.0042$	$0.9947 \pm 0.0025$
190k	$0.8768 \pm 0.0029$	–
380k	$0.8962 \pm 0.0024$	–
760k	$0.9402 \pm 0.0011$	–

(Can also compare multiple models with multi-class classifier and then evaluate on data)

(Or use latent space of pre-trained classifier: Frechet Inception/Particle Distance)

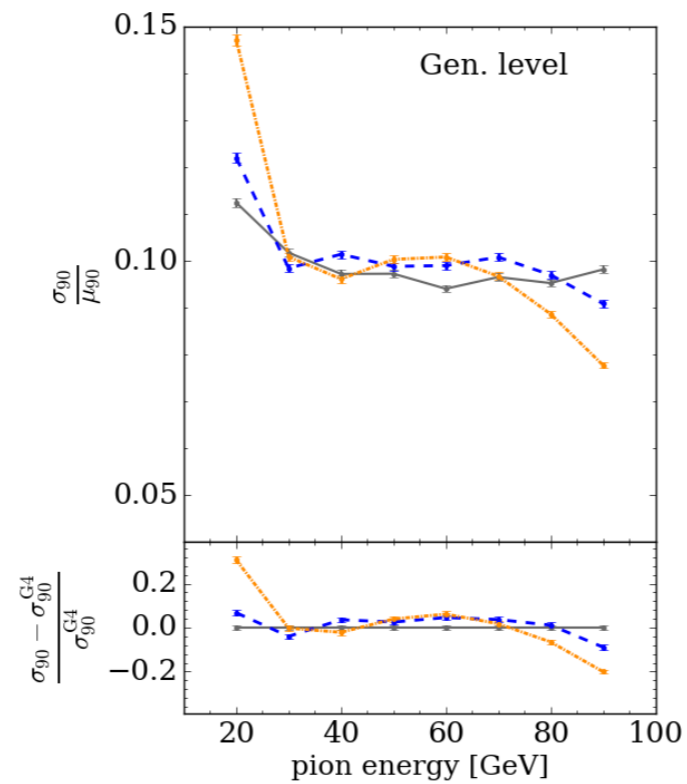
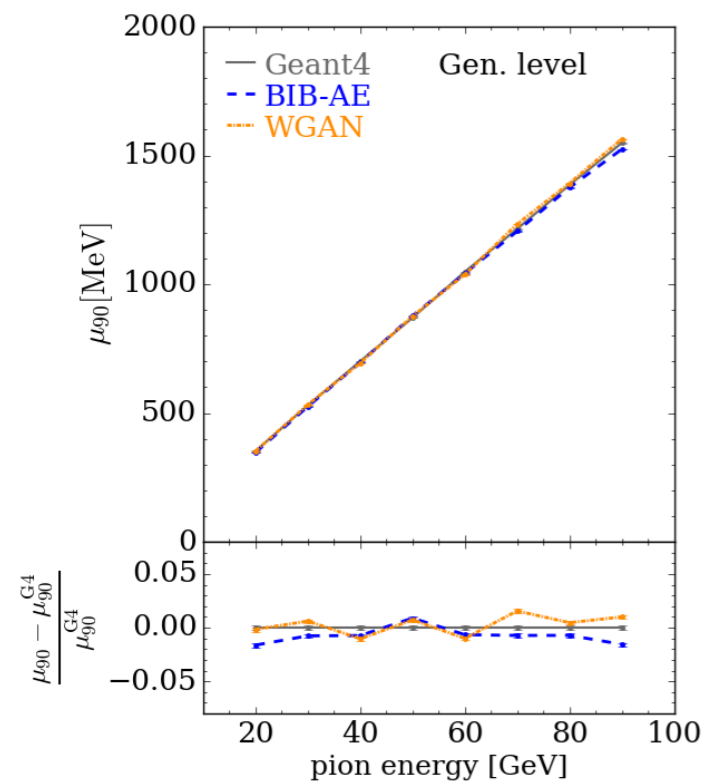
Capture full phase space information with classifiers

Still depends on training data

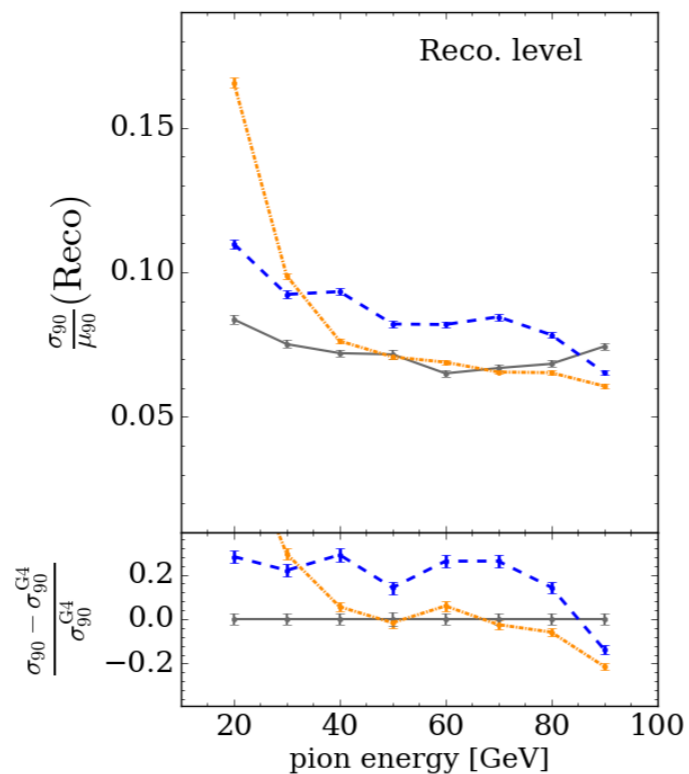
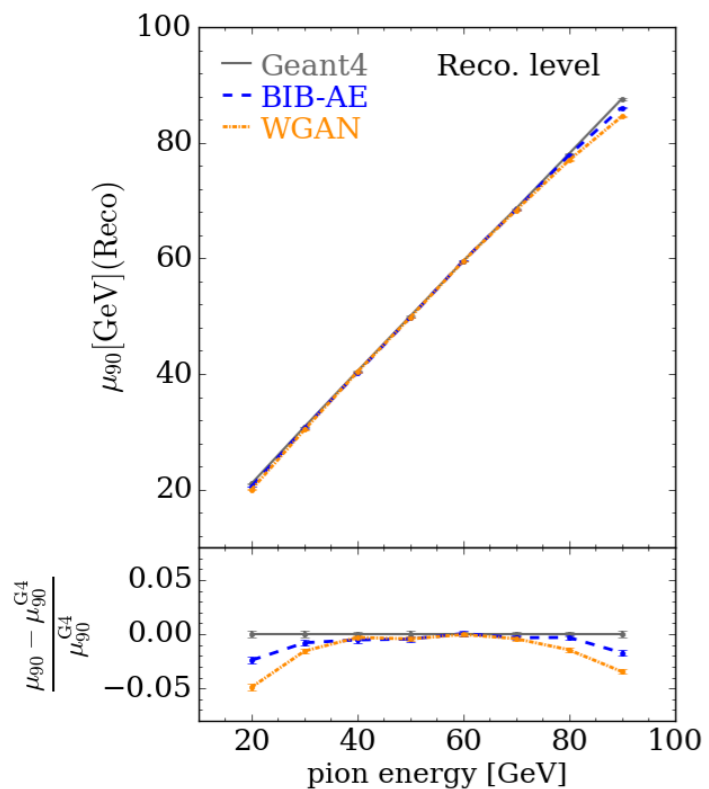
Choice of classifier

How good, is good enough really?

# Adding reconstruction



Without reconstruction



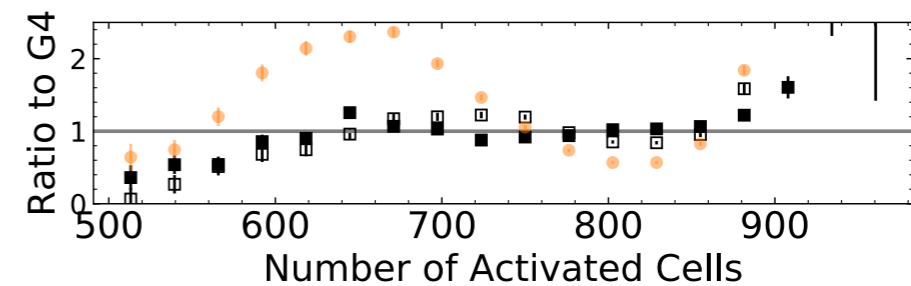
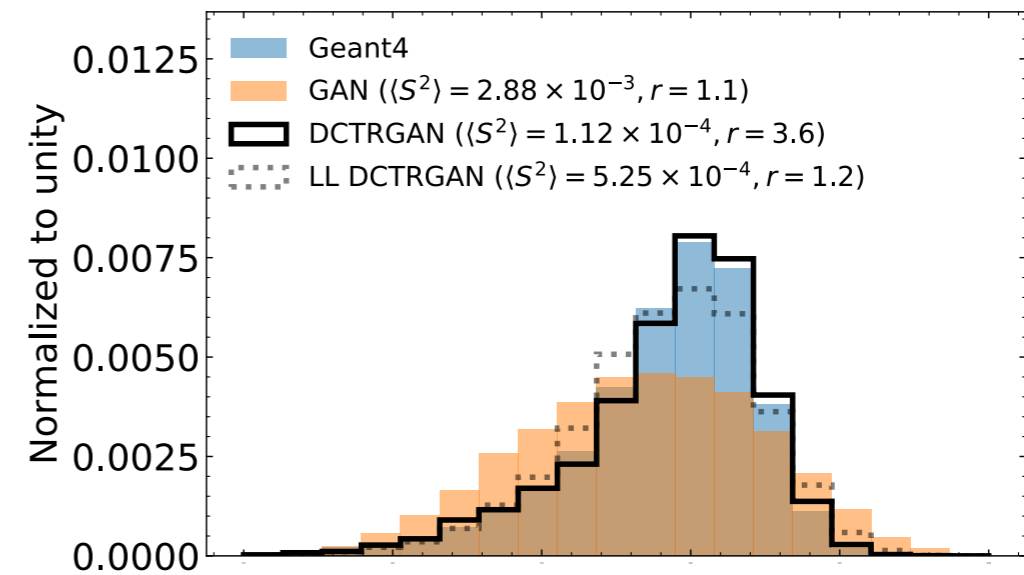
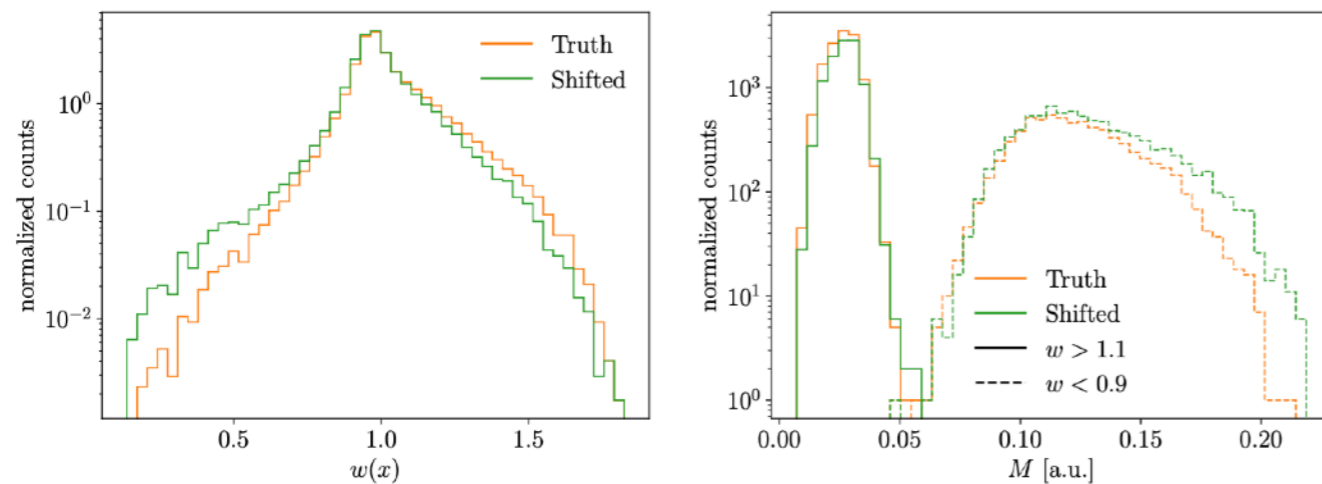
With reconstruction

→ Non-linear effects of reconstruction can change relative performance



# Weights

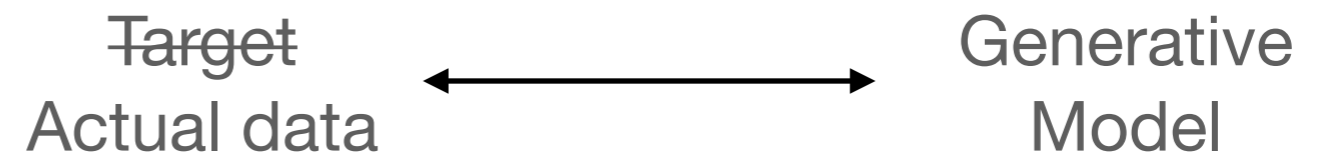
$$w(x) = \frac{p_{\text{data}}(x)}{p_{\text{model}}(x)} = \frac{C(x)}{1 - C(x)} \quad \text{with} \quad C(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_{\text{model}}(x)}$$



Train classifiers to find mismodelling

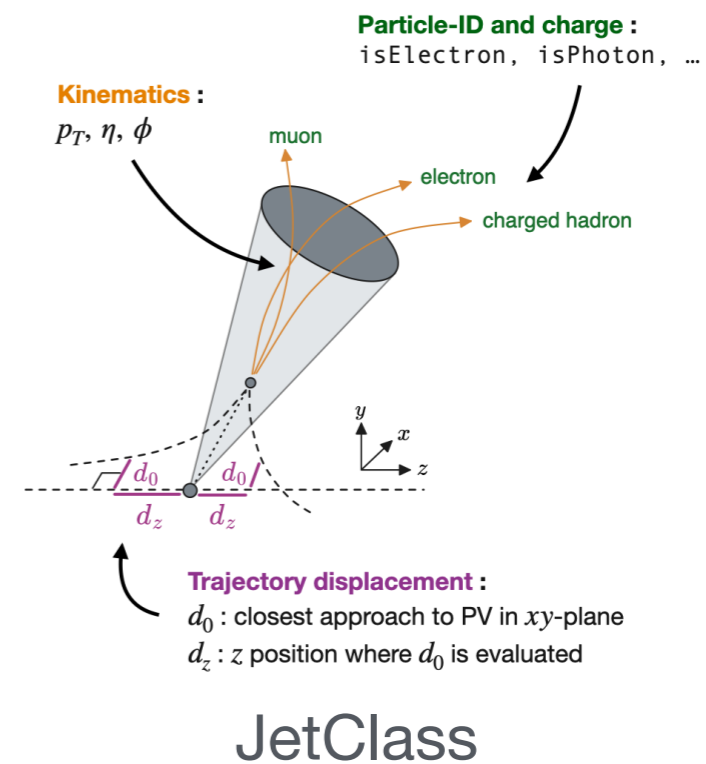
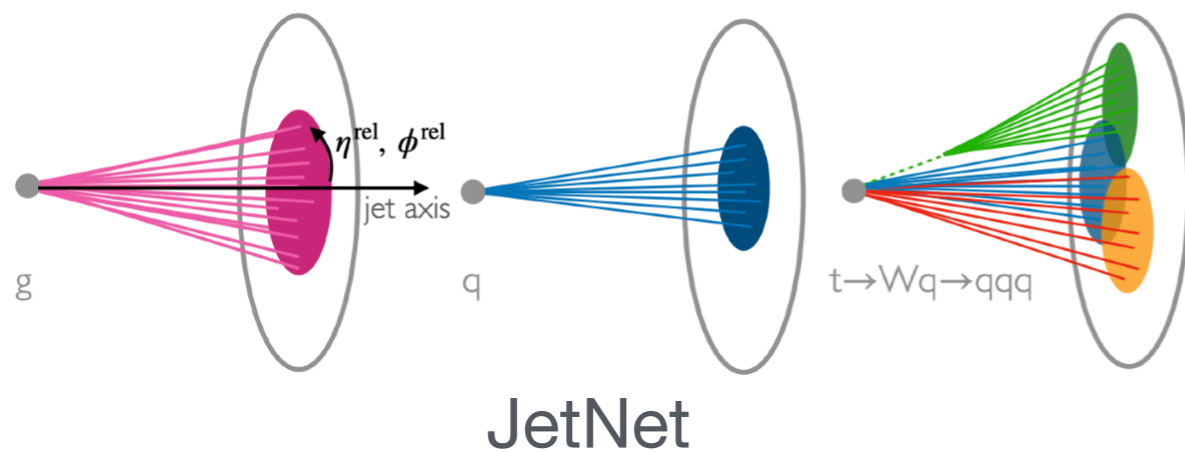
or to reweight distributions

# Quality of simulation



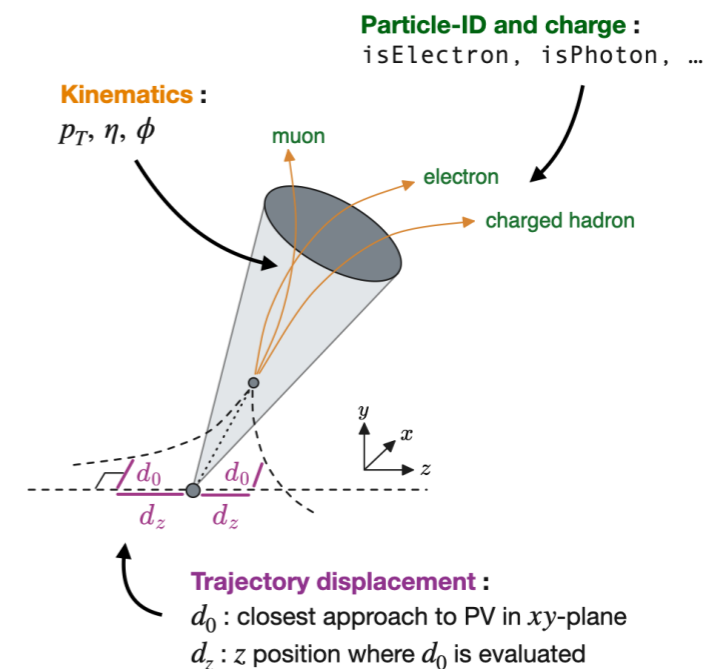
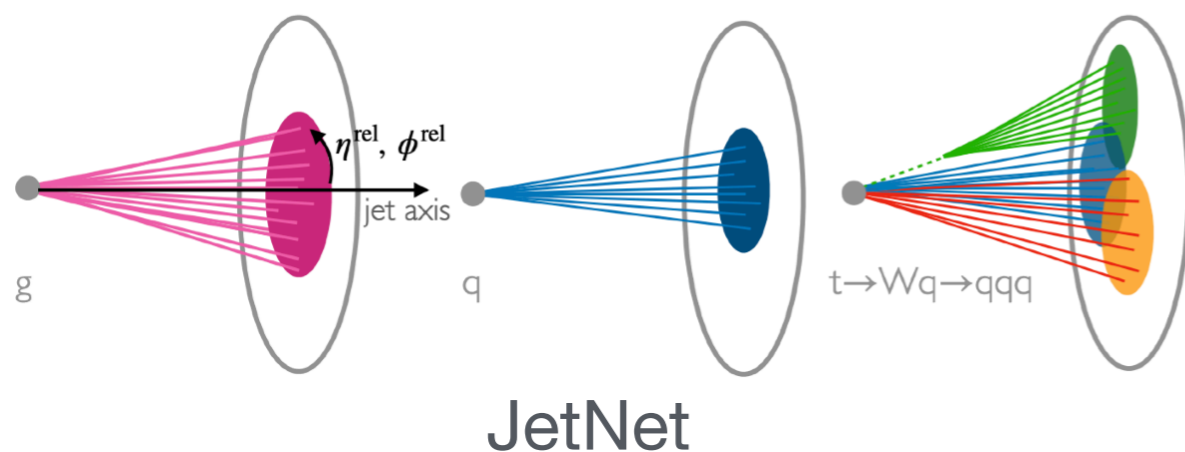
# Closing II

# Common Datasets



	JetNet [3]	JetClass [1]
Jet types	5 types	10 types ( several decay channels for top and H jets )
Dataset size	180 thousand jets per class	12.5 million jets per class ( 70x more than JetNet )
Features	Kinematics	Kinematics, Particle-ID and charge, trajectory displacement

# Common Datasets



**Fast Calorimeter Simulation Challenge 2022**

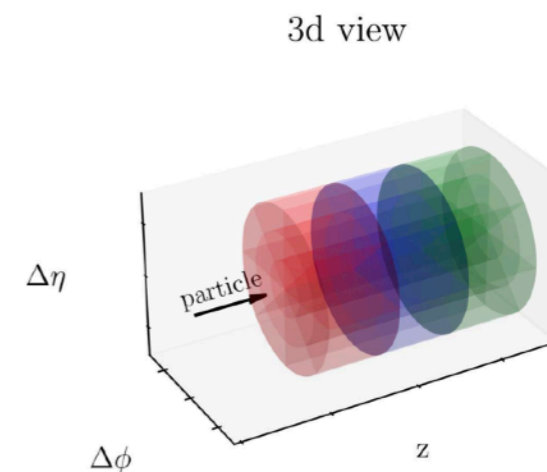
[View on GitHub](#)

Welcome to the home of the first-ever Fast Calorimeter Simulation Challenge!

The purpose of this challenge is to spur the development and benchmarking of fast and high-fidelity calorimeter shower generation using deep learning methods. Currently, generating calorimeter showers of interacting particles (electrons, photons, pions, ...) using GEANT4 is a major computational bottleneck at the LHC, and it is forecast to overwhelm the computing budget of the LHC experiments in the near future. Therefore there is an urgent need to develop GEANT4 emulators that are both fast (computationally lightweight) and accurate. The LHC collaborations have been developing fast simulation methods for some time, and the hope of this challenge is to directly compare new deep learning approaches on common benchmarks. It is expected that participants will make use of cutting-edge techniques in generative modeling with deep learning, e.g. GANs, VAEs and normalizing flows.

This challenge is modeled after two previous, highly successful data challenges in HEP - the [top tagging community challenge](#) and the [LHC Olympics 2020 anomaly detection challenge](#).

3 datasets of increasing complexity:  
 DS1: Up to 533 voxels (ATLAS)  
 DS2: 6480 voxels  
 DS3: 40500 voxels





# Generative Cheat Sheet

