

Simulation & Generative Models

Gregor Kasieczka

Email: 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














$$\begin{matrix} \frac{F \Delta}{\Delta x} \\ \frac{1}{\Delta x} \int_{x_0}^{x_0+\Delta x} E dx \\ C) \rho \Delta = \rho \Delta x \end{matrix}$$

$$= \epsilon \frac{E}{E} \left(\frac{\epsilon}{c} H = \frac{\nabla \times \mathbf{c} = u \mathbf{r}}{j \mathbf{r} \cdot d \mathbf{r} \cdot hc}$$

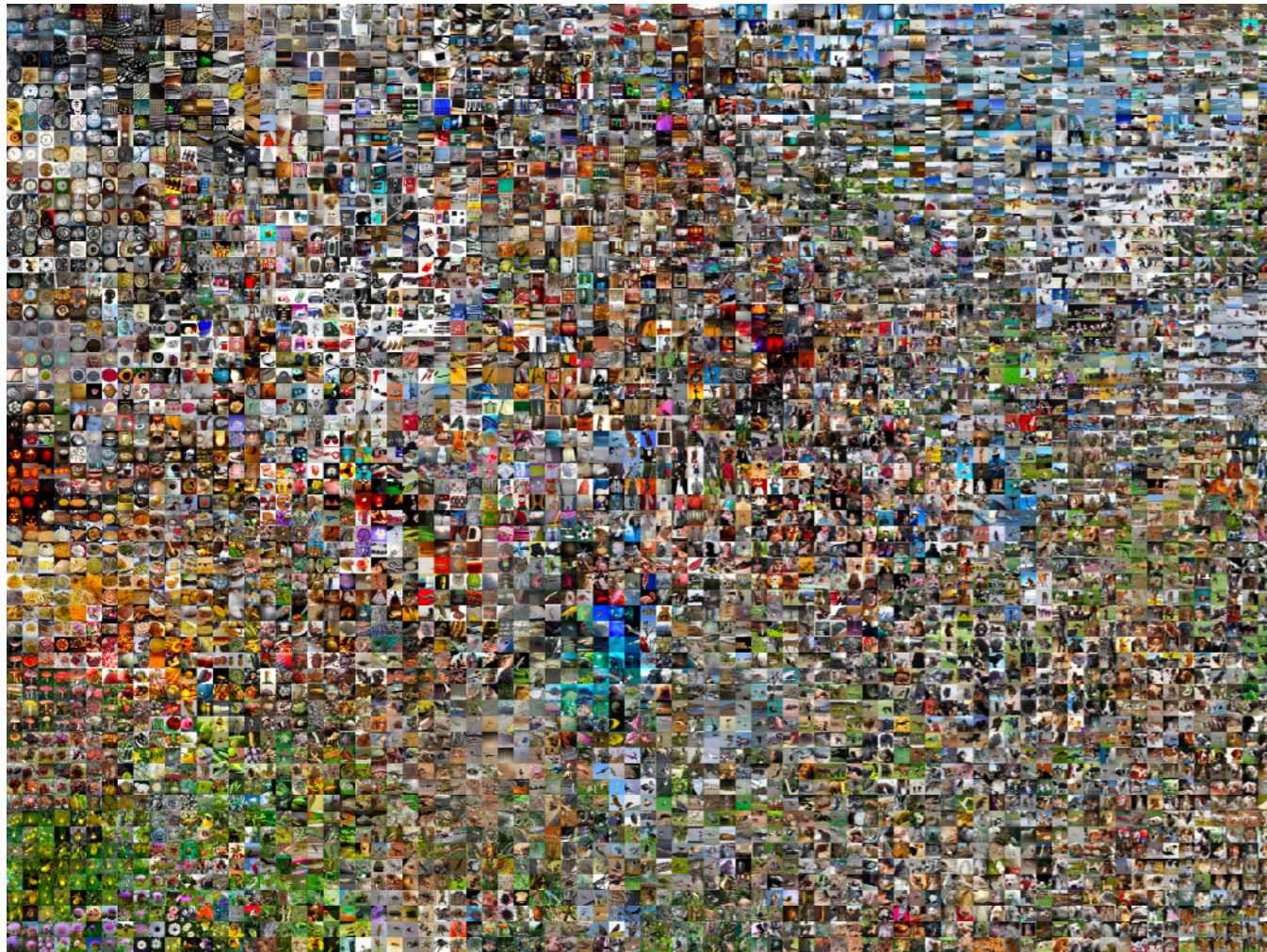
$$= \frac{ns}{z} - B - \Pi \frac{0 C - 1}{\Delta R} \frac{3 C - 1}{\Delta R} \frac{4 A}{\Delta R} IS = D - [5 - H \left(C \frac{a}{1} = \frac{D \mathbf{r} = A C I}{\tau \Delta U} = \mathbf{r} \right)$$
$$C - D \rho_2 (1 - \rho_2) (=) C_F =) B \cup \rho_2 (y - N = u \mathbf{c} = -v = \frac{n D}{B} D \frac{r \cdot i}{2 T} \frac{u}{A \epsilon} D)$$
$$- \frac{x}{A} = \frac{20}{10} \frac{\Delta R_{max}}{\Delta L} P - \frac{D}{x} - \frac{L}{x} B \frac{1.5}{\Delta D} C_s = \frac{c N}{-1}$$

$$\Delta C = C_F \cdot j \cdot E = \rho = n \epsilon$$

$$C \rho = P) = A \cdot \Delta x$$
$$= \Delta \cdot K \cdot \rho \cdot \Delta x$$

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?

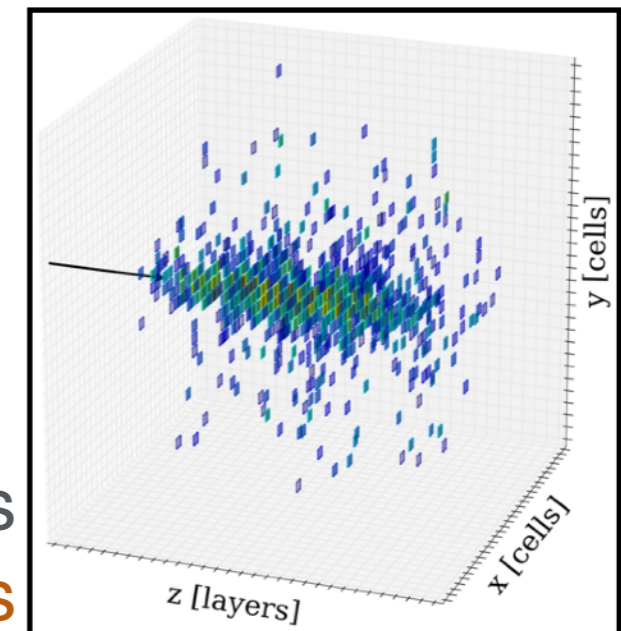
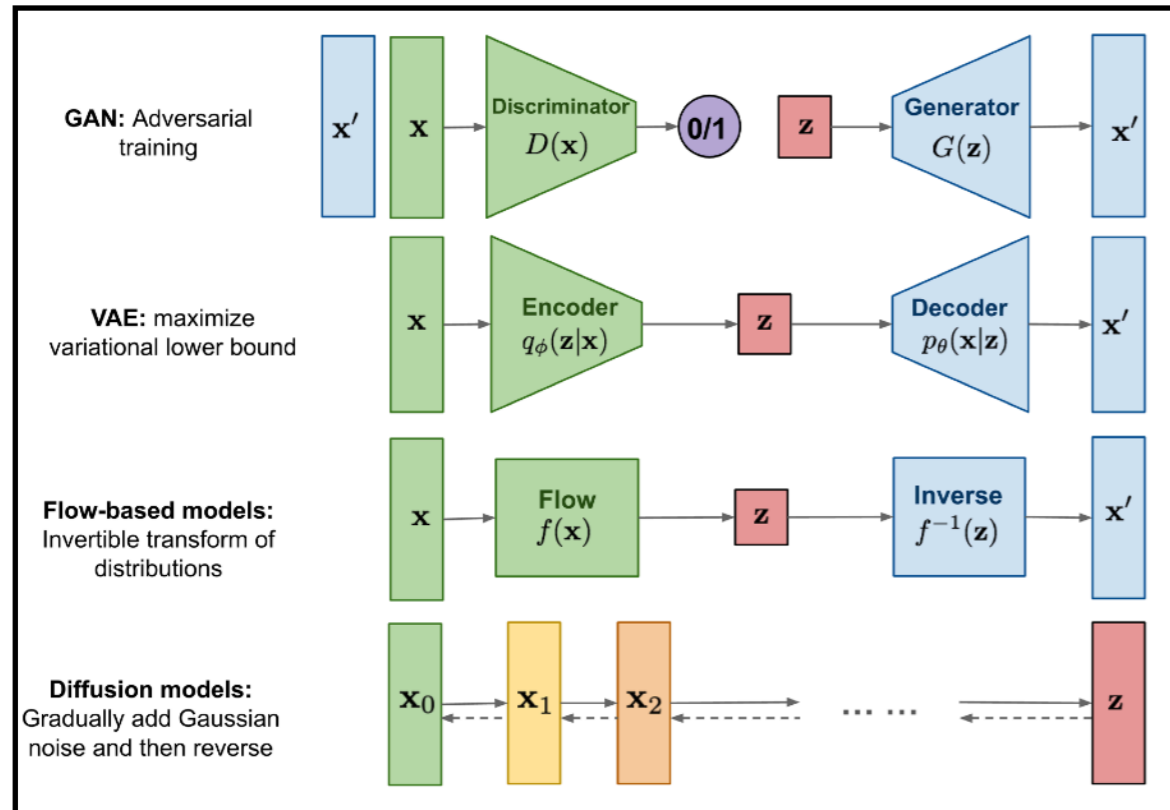
Uses:

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

Overview

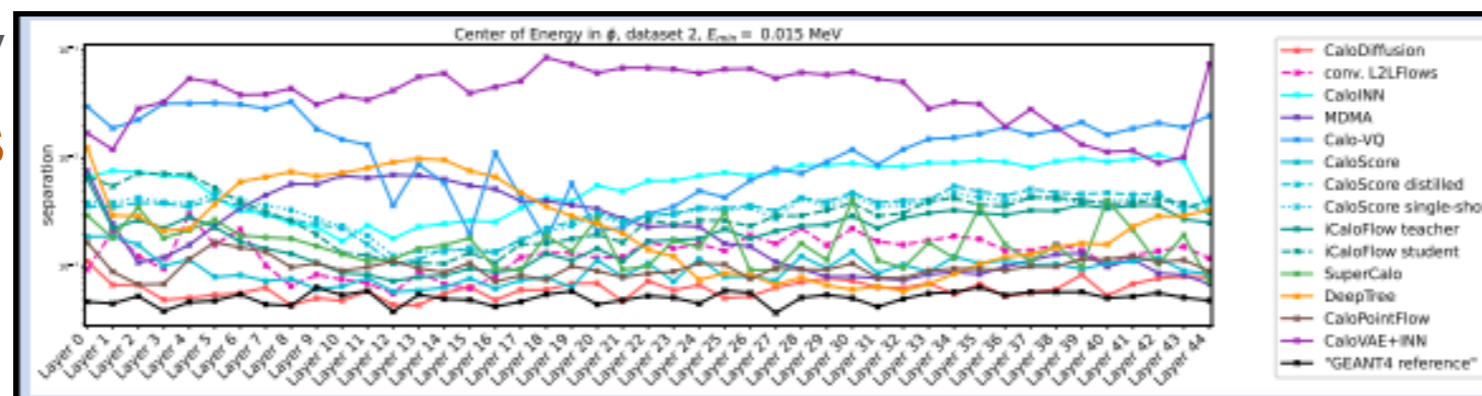
1. Common architectures*

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



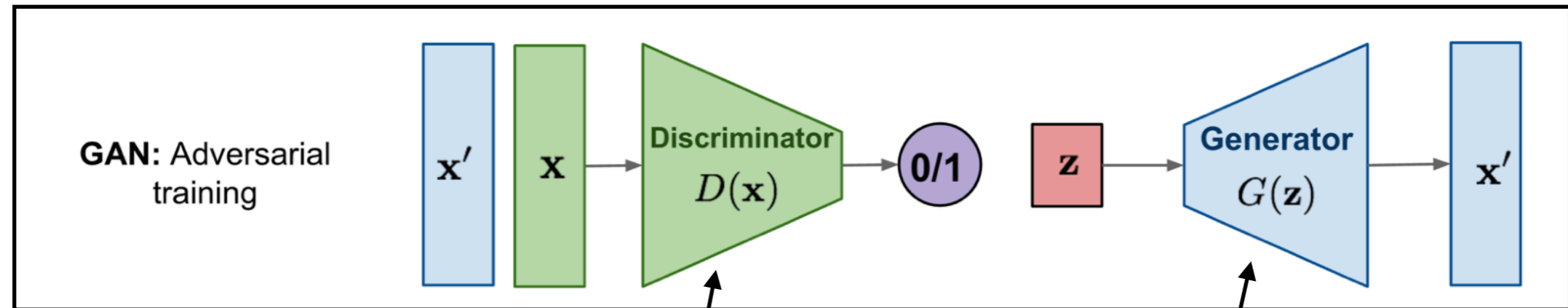
2. Physics applications

3. Quality metrics



Generative Adversarial Networks

Generative Adversarial Networks

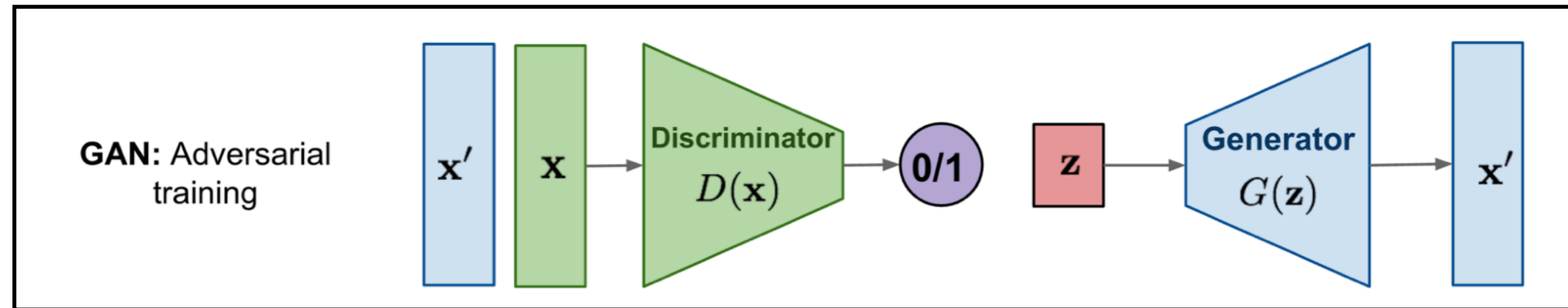


Generative Adversarial Networks (GANs) consist of **2 networks**

Provides feedback on quality of examples

Maps random noise to realistic examples

Generative Adversarial Networks

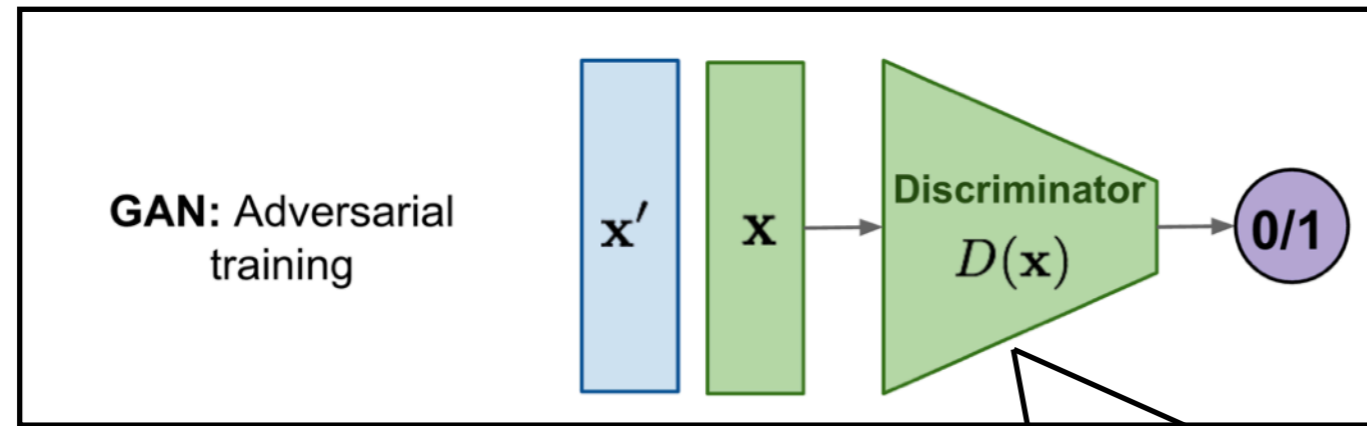


Training objective:
Binary cross entropy

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

True examples Fake examples

Generative Adversarial Networks

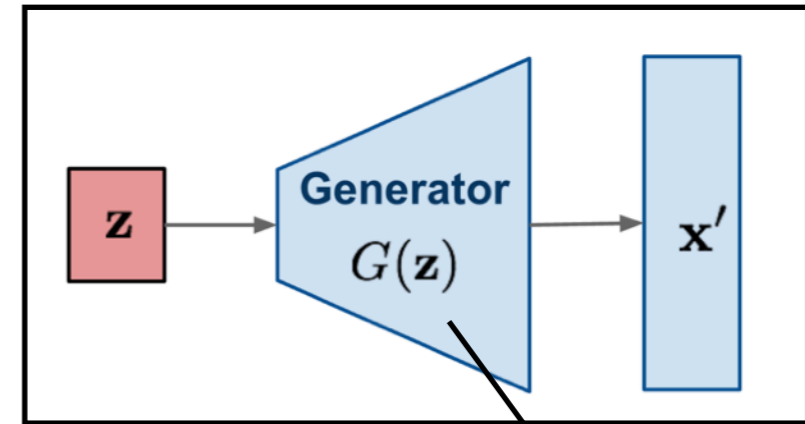


Training objective:
Binary cross entropy

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

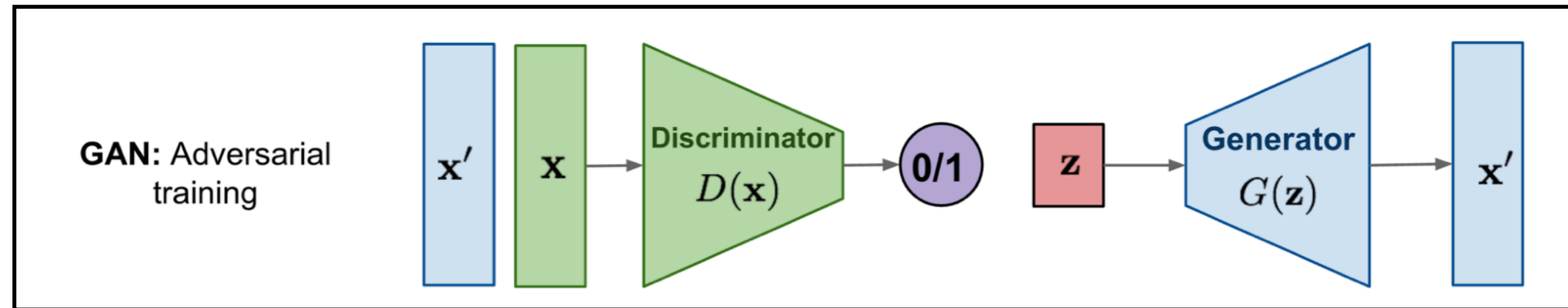


Training objective:
Binary cross entropy

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

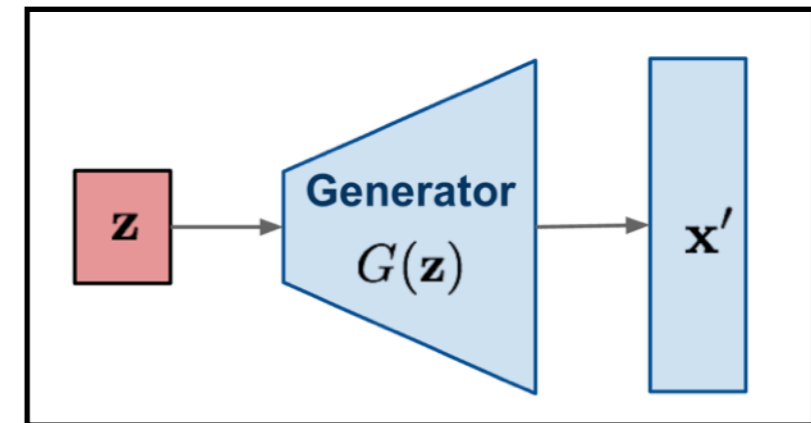
$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

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**

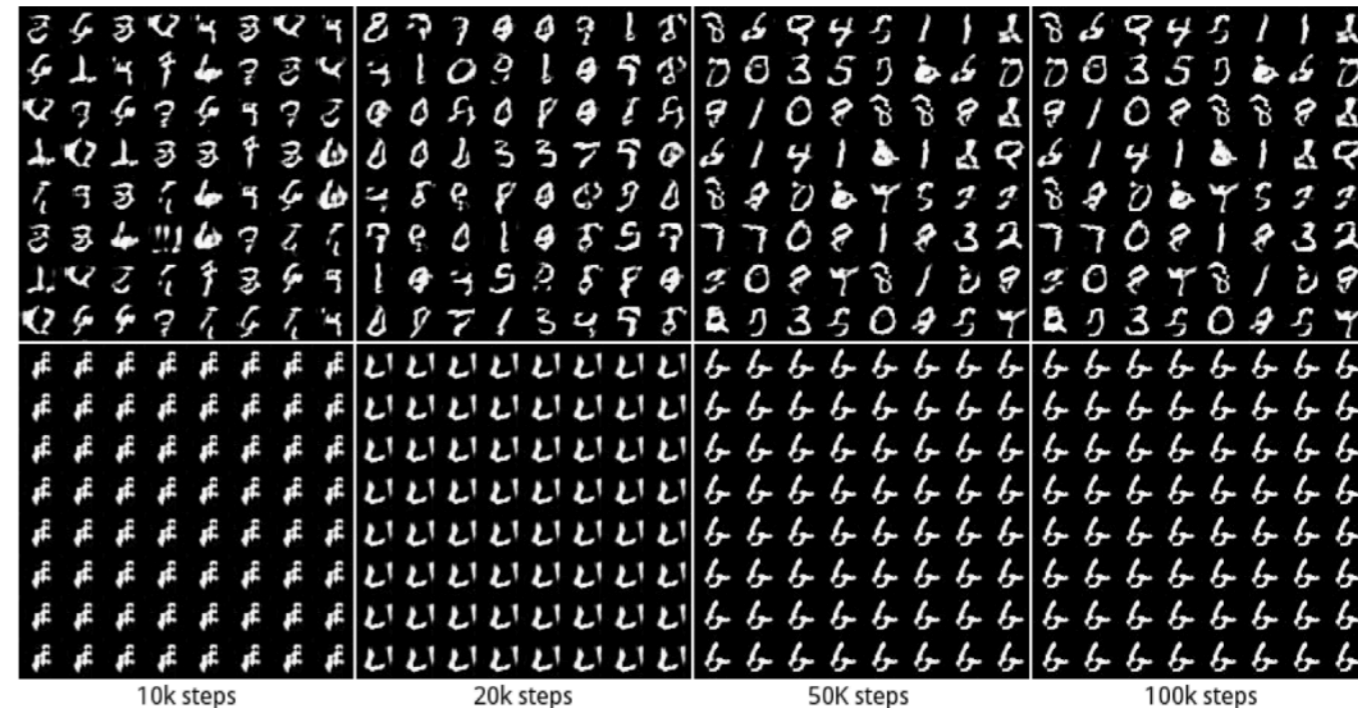
Comments on GANs

Architecture:

- Low complexity, **fast** and adaptable

Learning:

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

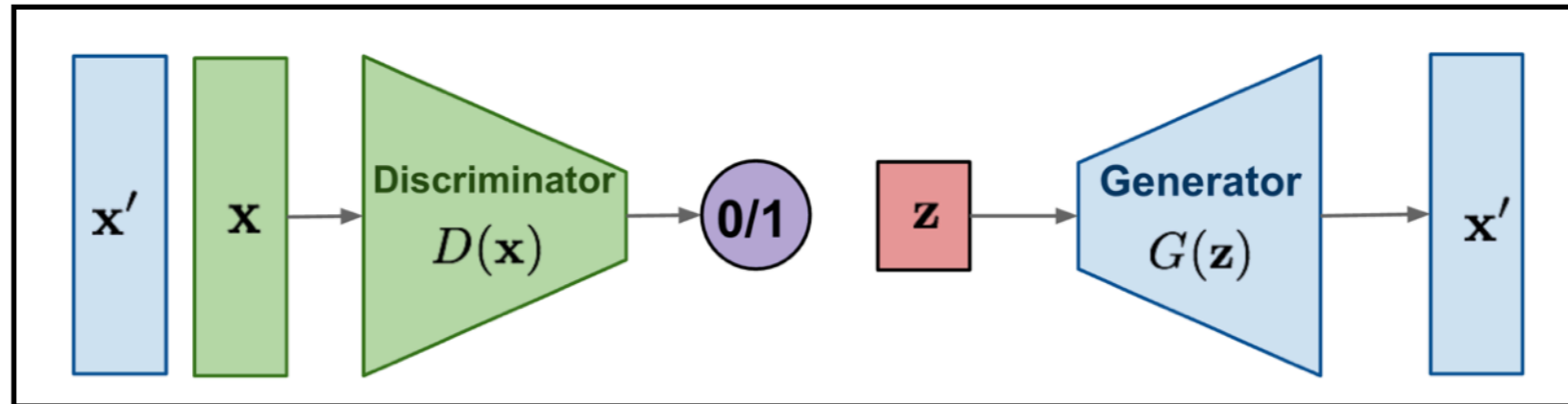


Mode collapse

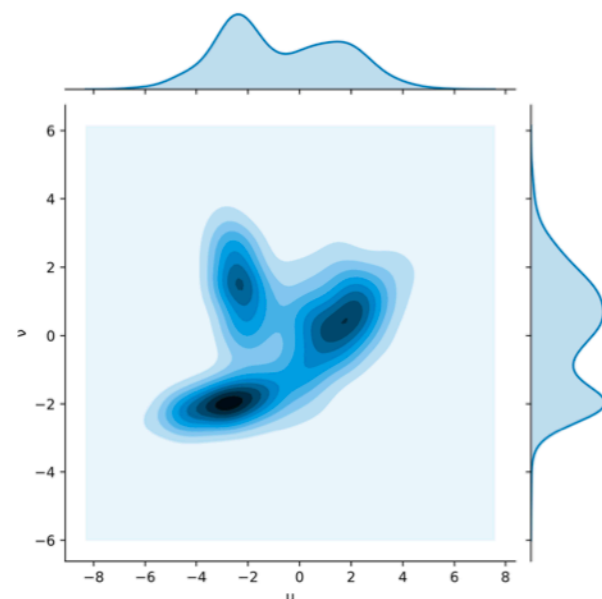
Maturity:

- Well established, many variants and extensions

Wasserstein GAN

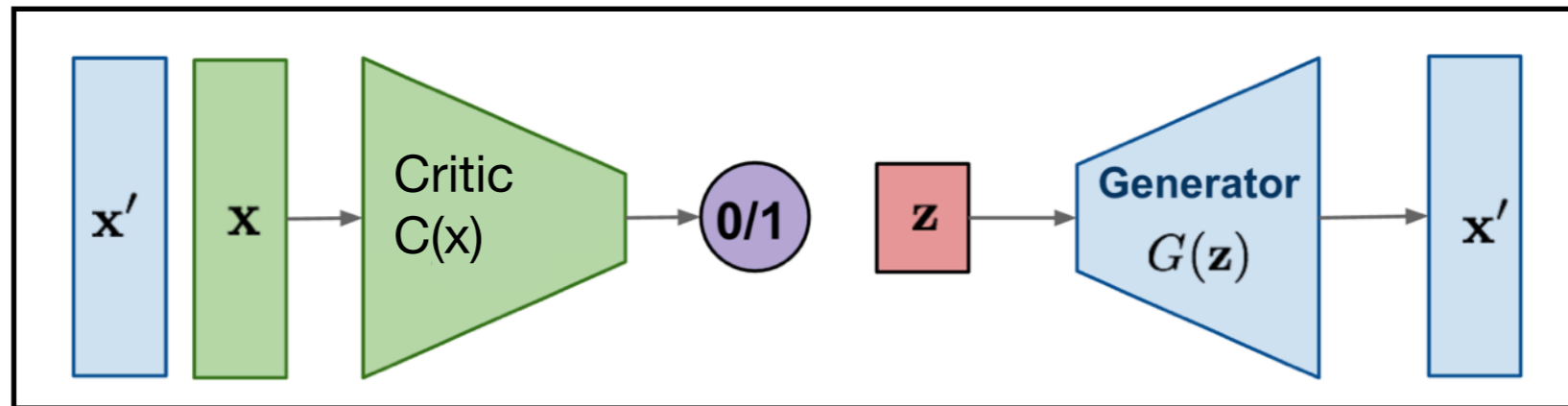


- Standard GANs minimise Jensen-Shannon divergence of generator output and true data
 - Not best measure, e.g. for non-overlapping distributions
- Replace with Wasserstein / Earth-Mover-Distance



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

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

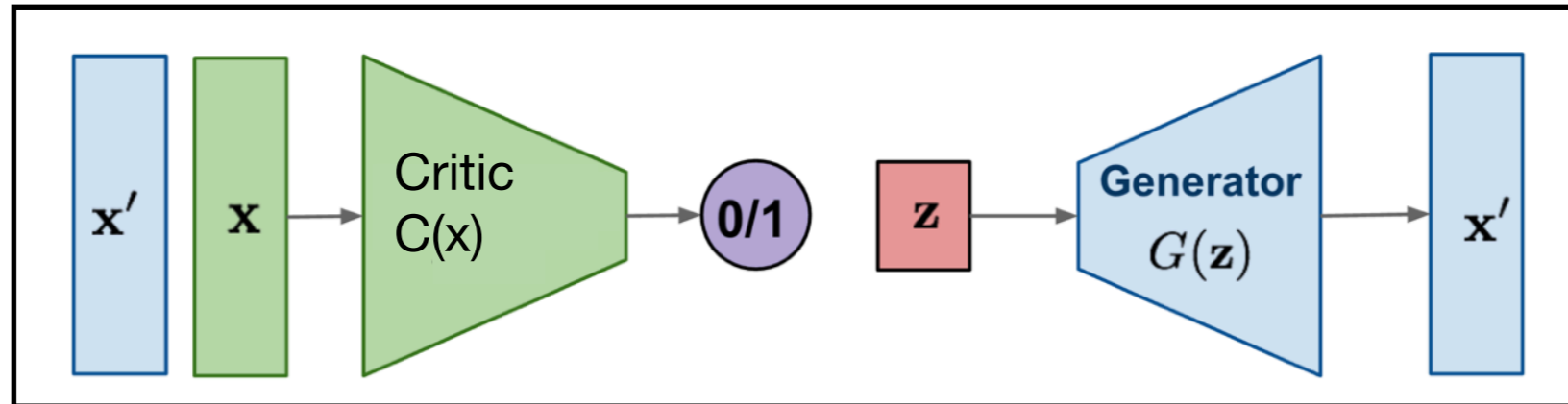
loss*:

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

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

* Some mathematics involved from earth mover distance to here

Wasserstein GAN



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

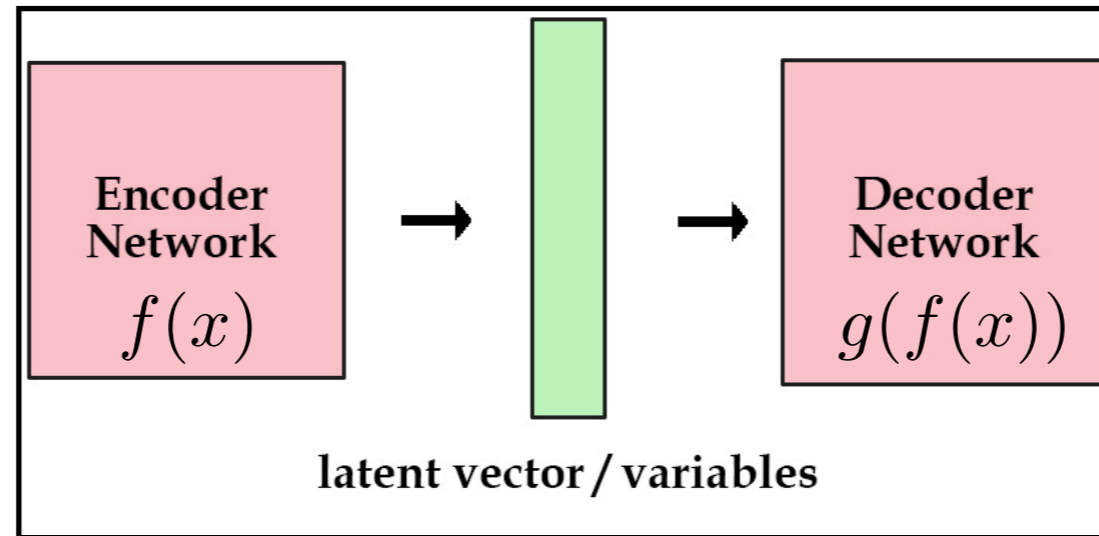
Wasserstein GAN

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

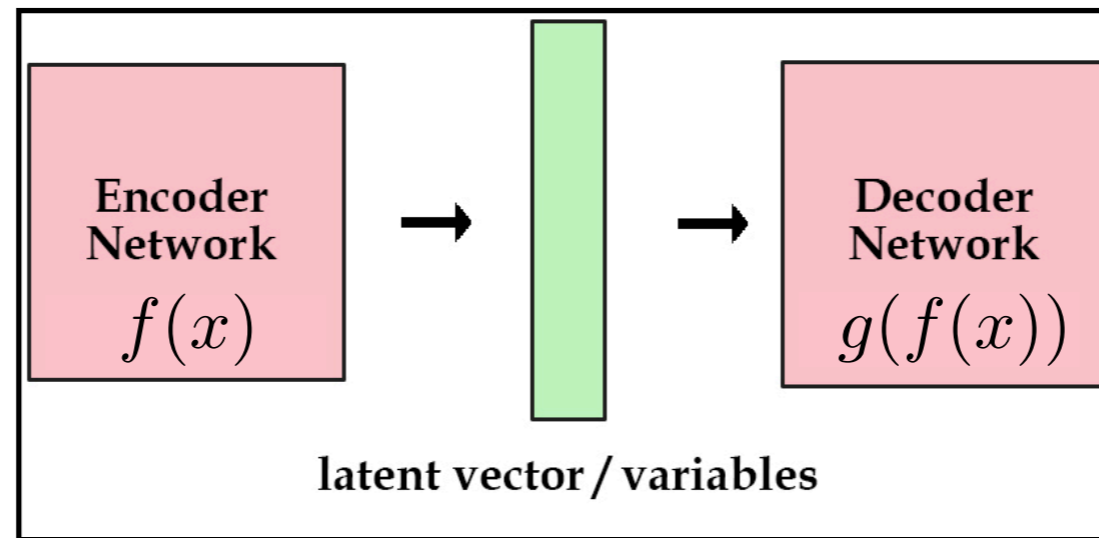


Two networks

Encoder: data \rightarrow latent space

Decoder: latent space \rightarrow data

Autoencoder



Two networks

Encoder: data \rightarrow latent space

Decoder: latent space \rightarrow data

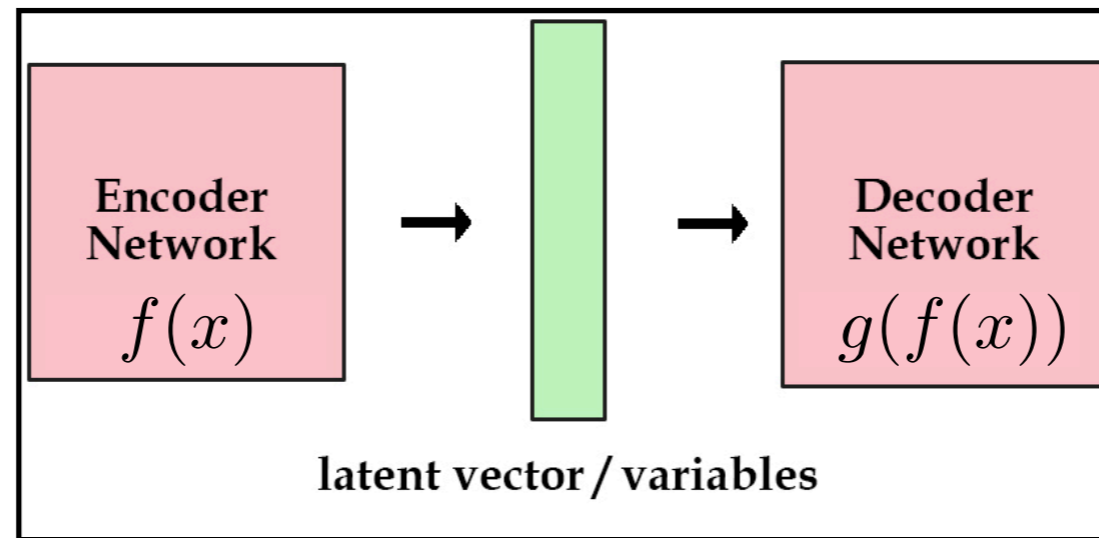
Training objective:

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:

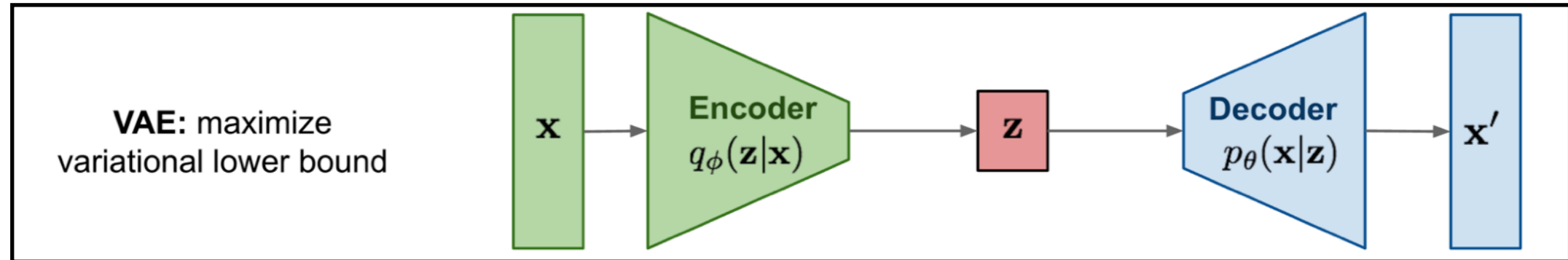
Dimension reduction

Denoising

Anomaly detection

Generation?

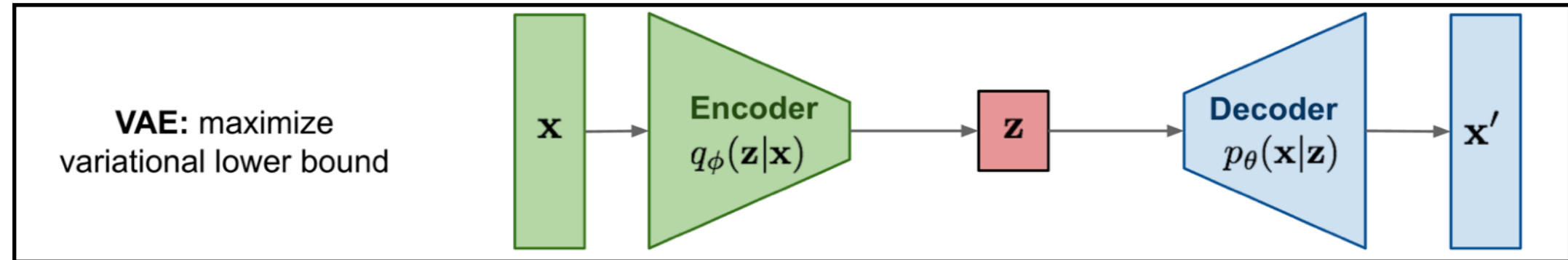
Variational Autoencoder



Variational Autoencoder (VAE):
Split latent space

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

Variational Autoencoder



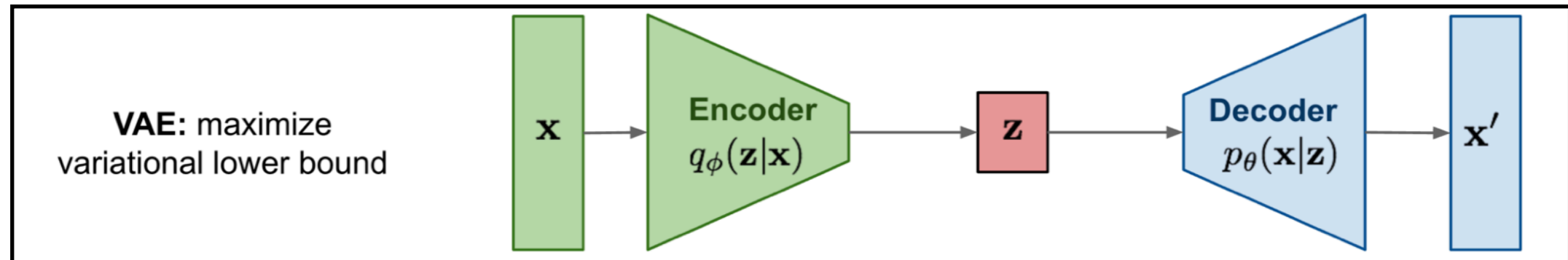
Variational Autoencoder (VAE):
Split latent space
Sample before decoder

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

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

$$x' = g(z)$$

Variational Autoencoder



Variational Autoencoder (VAE):

Split latent space

Sample before decoder

Penalty so mean/std are close to unit Gaussian

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

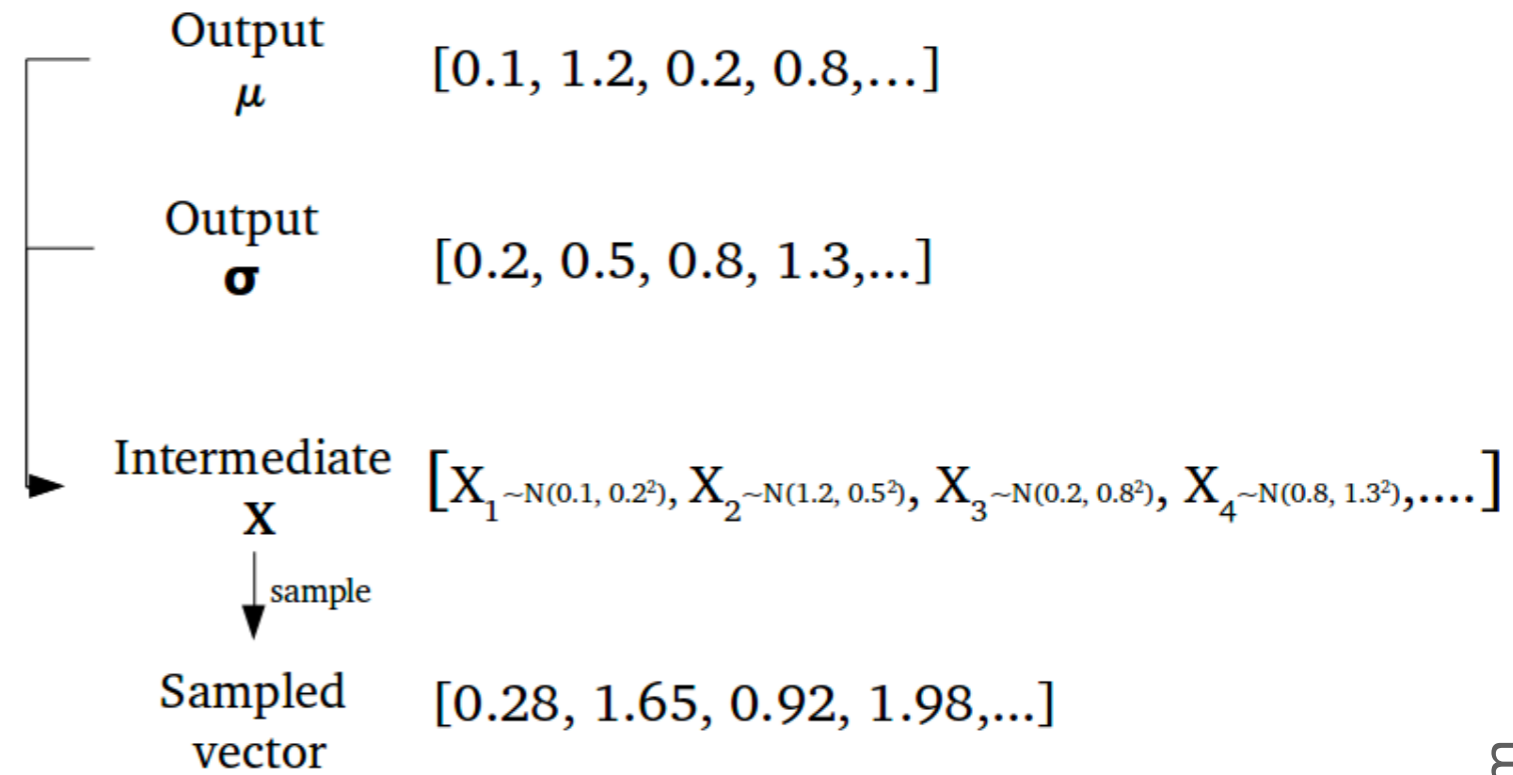
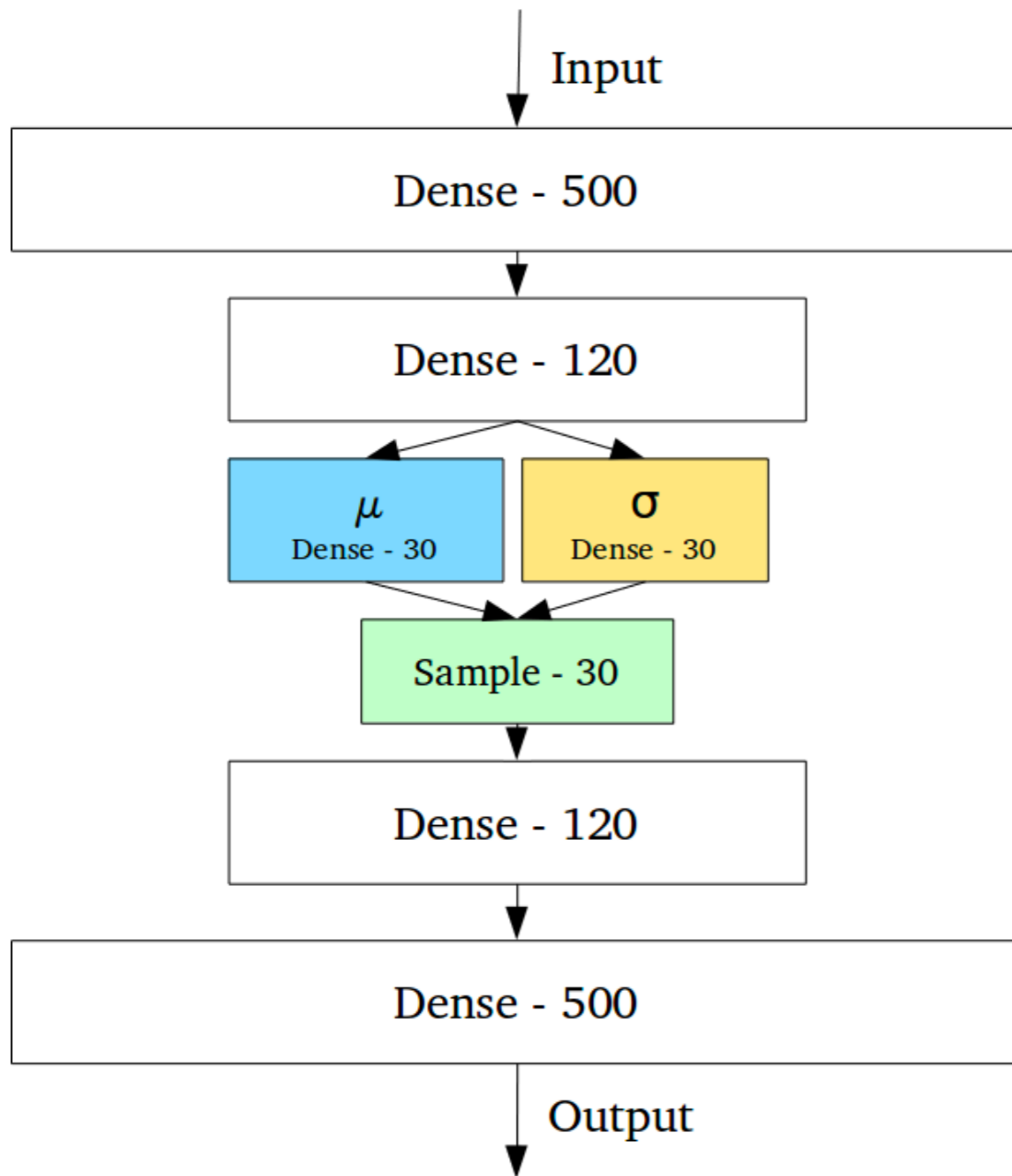
$$z = \text{Gaussian}(\mu, \sigma)$$

$$x' = g(z)$$

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

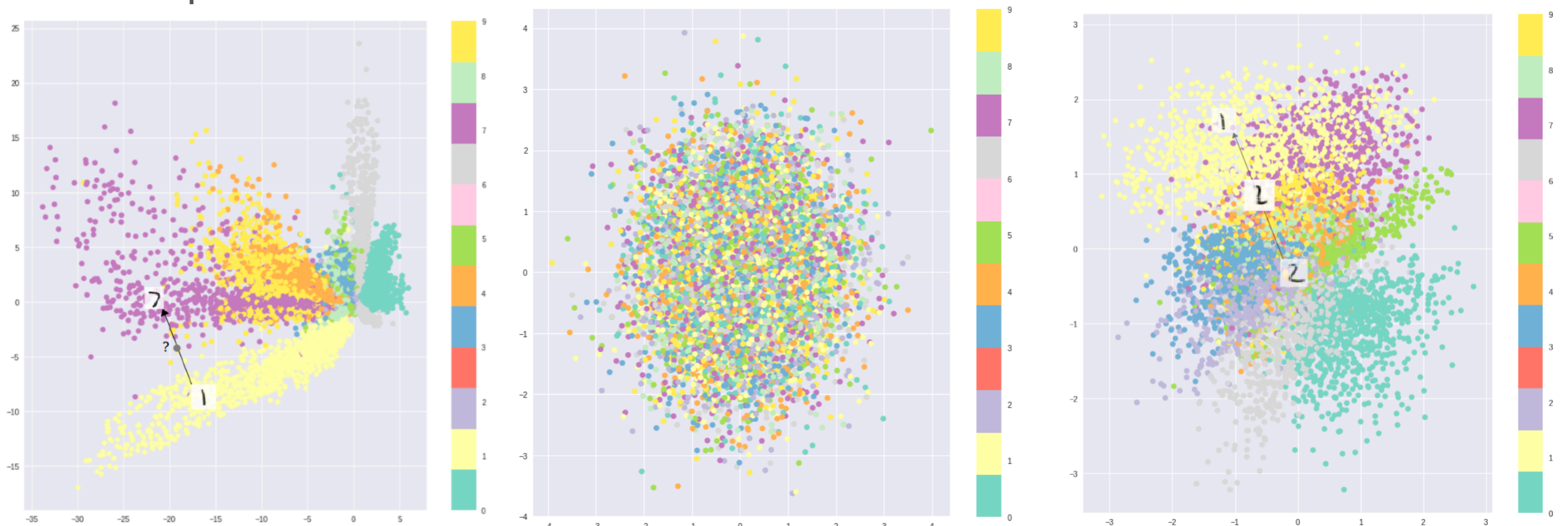
(Calculate KL-divergence between Gaussians)

VAE Example



Loss terms

Latent space of MNIST VAE



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

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

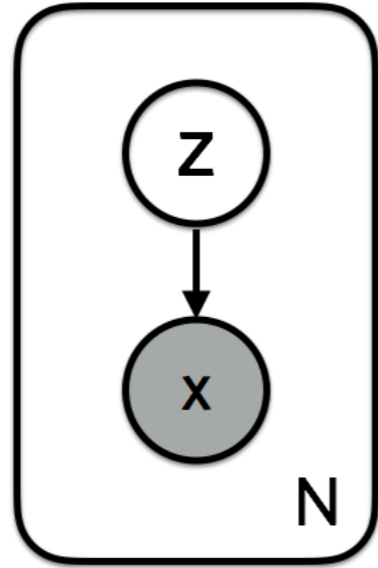
Both terms

Loss terms

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

How did we get here?

Loss terms



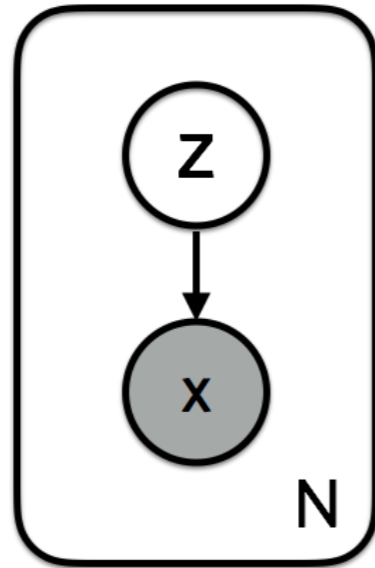
Sample from latent variables z

$$z_i \sim p(z)$$

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 ≥ 0 , ELBO is a lower limit for $\log p(x)$
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



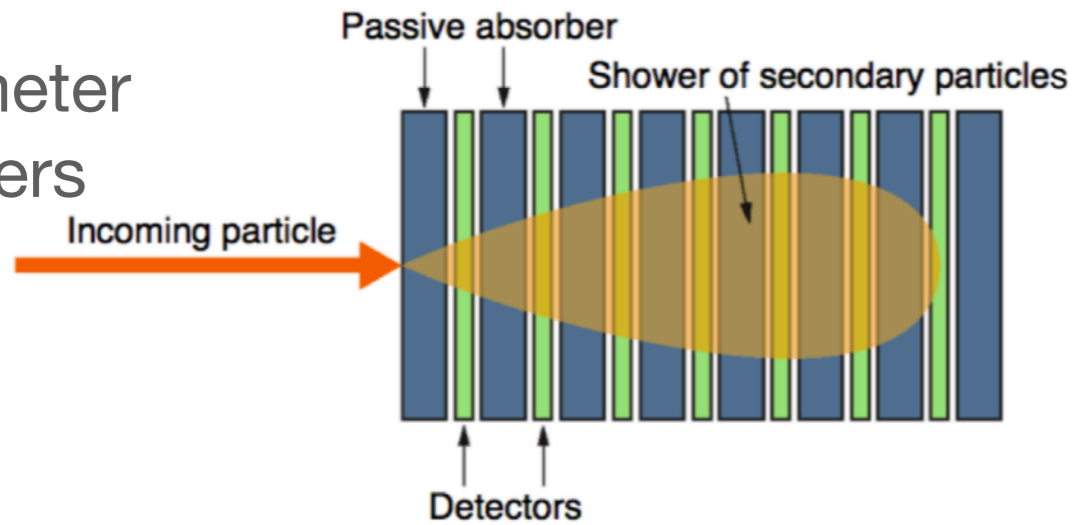
Maturity:

- Well established,
many variants and extensions

Applications I

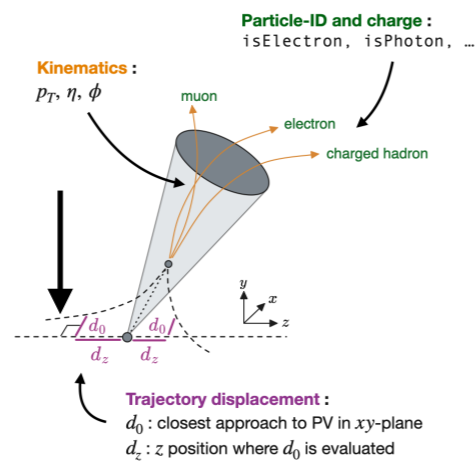
(Some) Simulation targets

Calorimeter
Showers



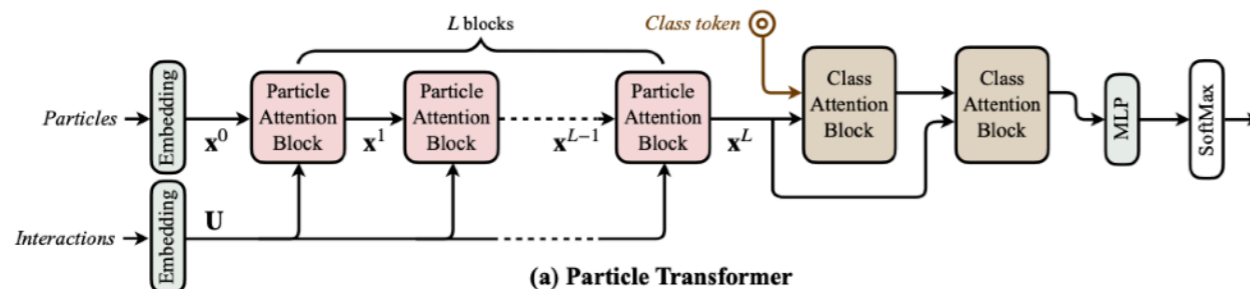
Reduce computational bottleneck

Jet Constituents



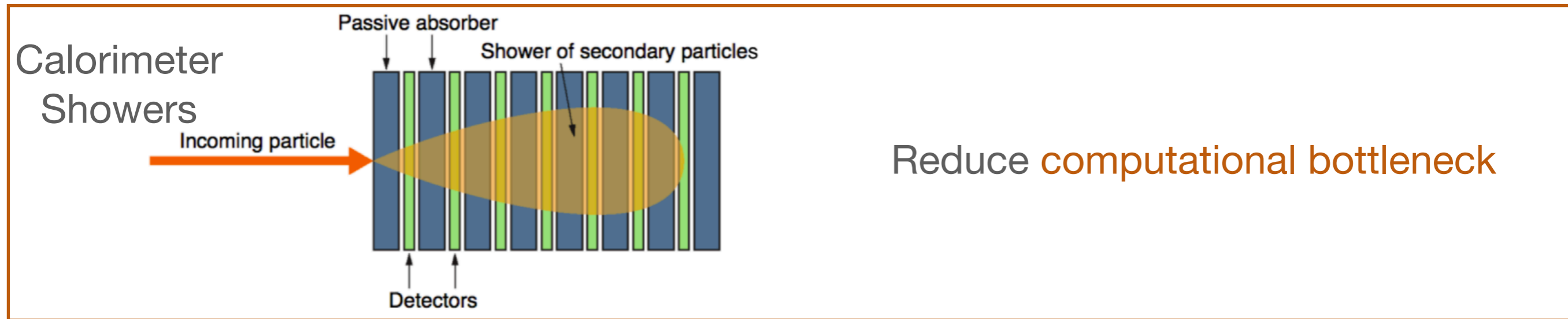
Predict background from data

Classification and
Reconstruction tasks

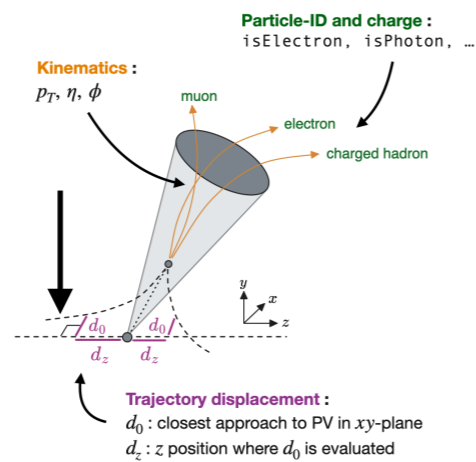


Act as surrogate models

(Some) Simulation targets

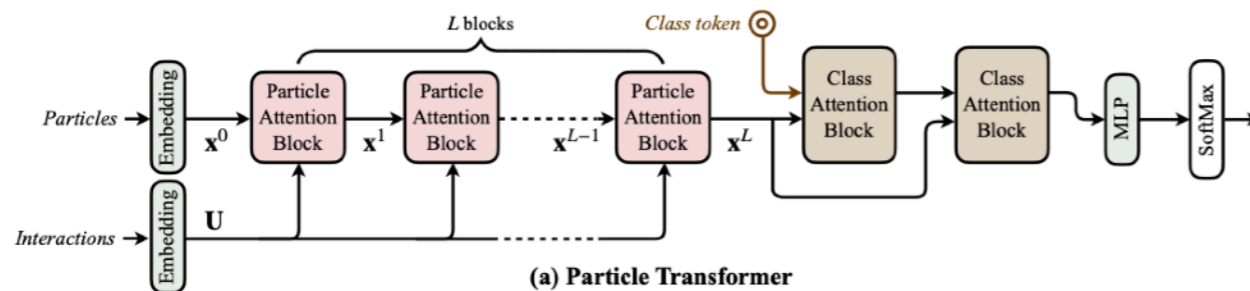


Jet Constituents



Predict background from data

Classification and Reconstruction tasks

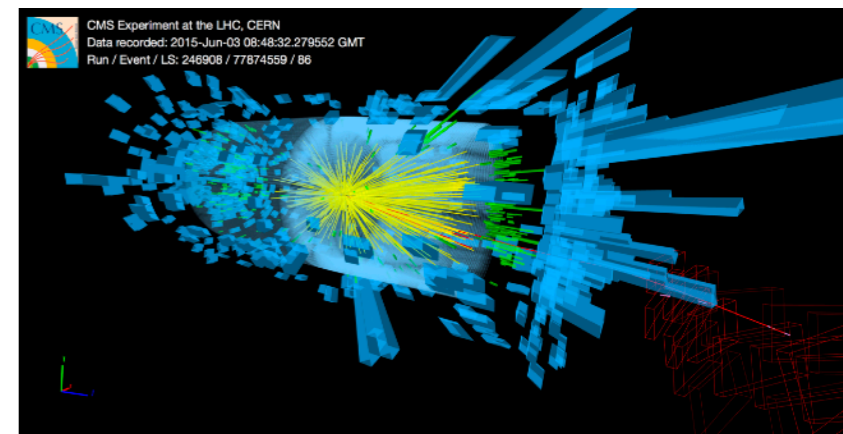
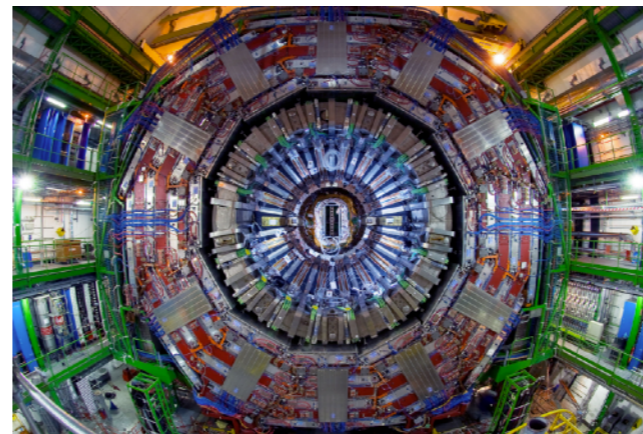
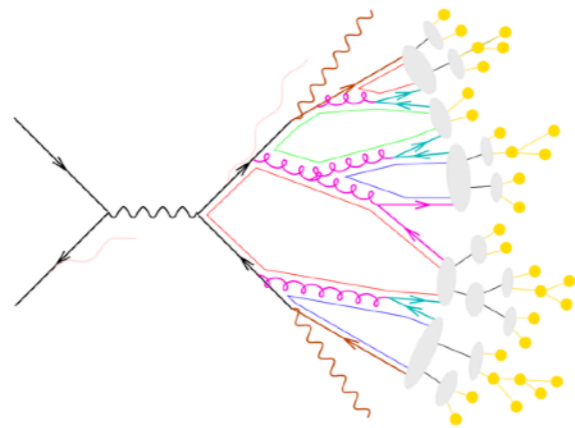


Act as surrogate models

Generative Models

This happens in the experiment

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



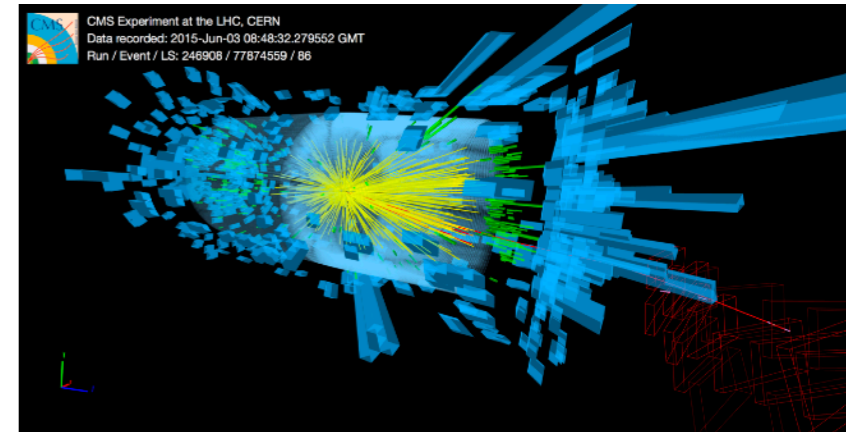
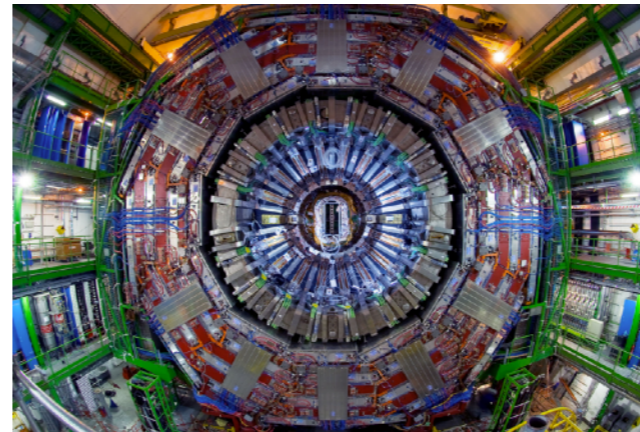
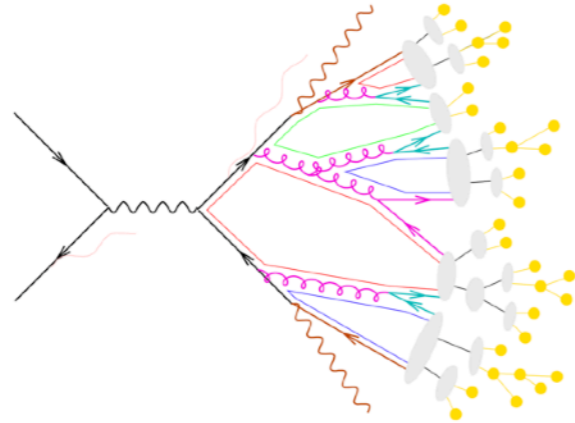
This is what we want to know

Simulation is crucial to connect
experimental data with theory
predictions

Generative Models

This happens in the experiment

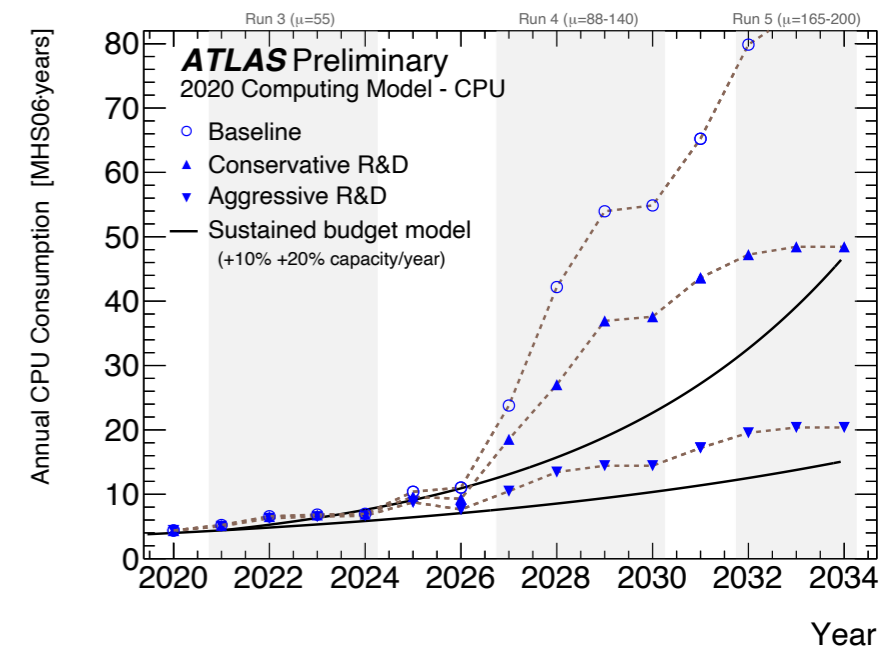
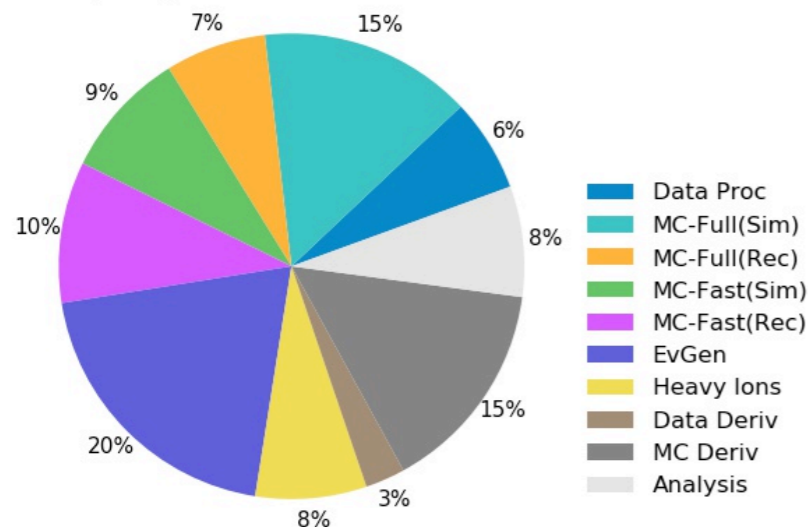
$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} + i\bar{\psi}\not{D}\psi + h.c. + \chi_i Y_{ij} \chi_j \phi + h.c. + |D_\mu \phi|^2 - V(\phi)$$



This is what we want to know

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

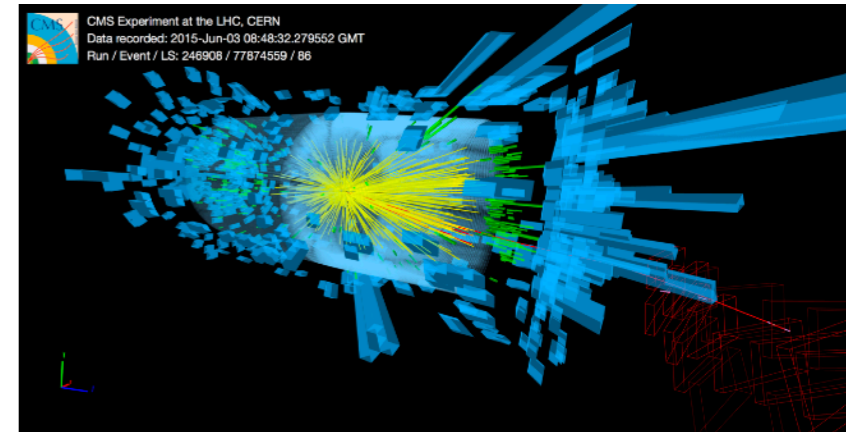
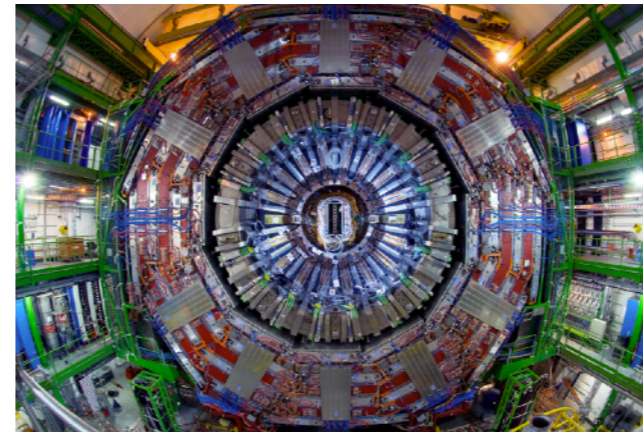
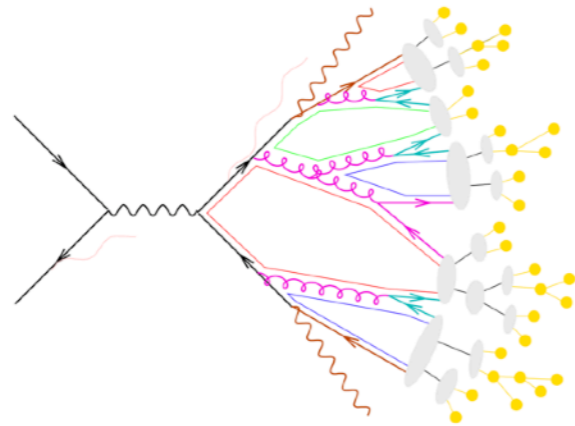
ATLAS Preliminary
2020 Computing Model - CPU: 2030: Baseline



Generative Models

This happens in the experiment

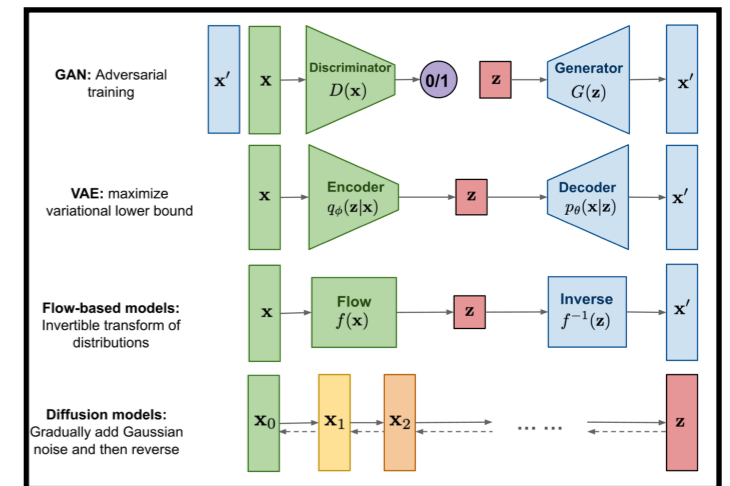
$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} + i\bar{\psi}\not{D}\psi + h.c. + \chi_i Y_{ij} \chi_j \phi + h.c. + |D_m \phi|^2 - V(\phi)$$



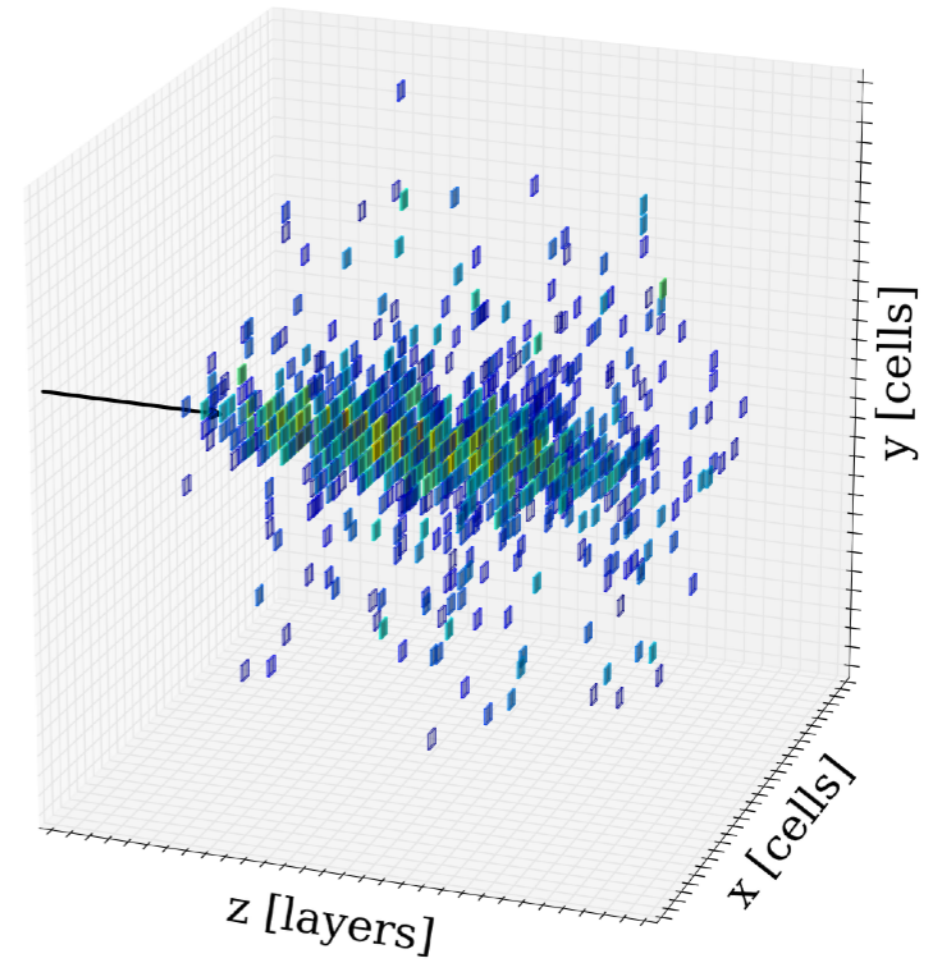
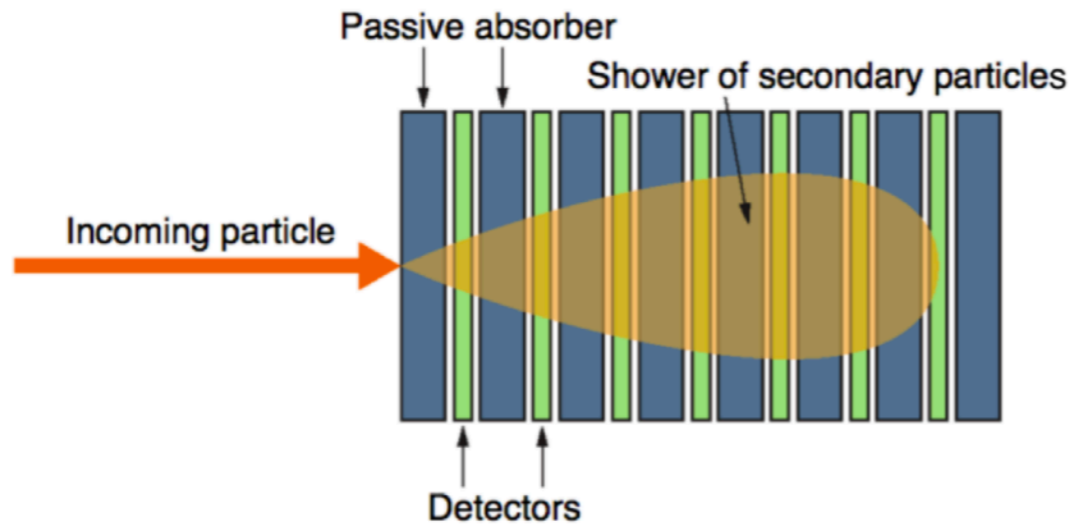
This is what we want to know

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

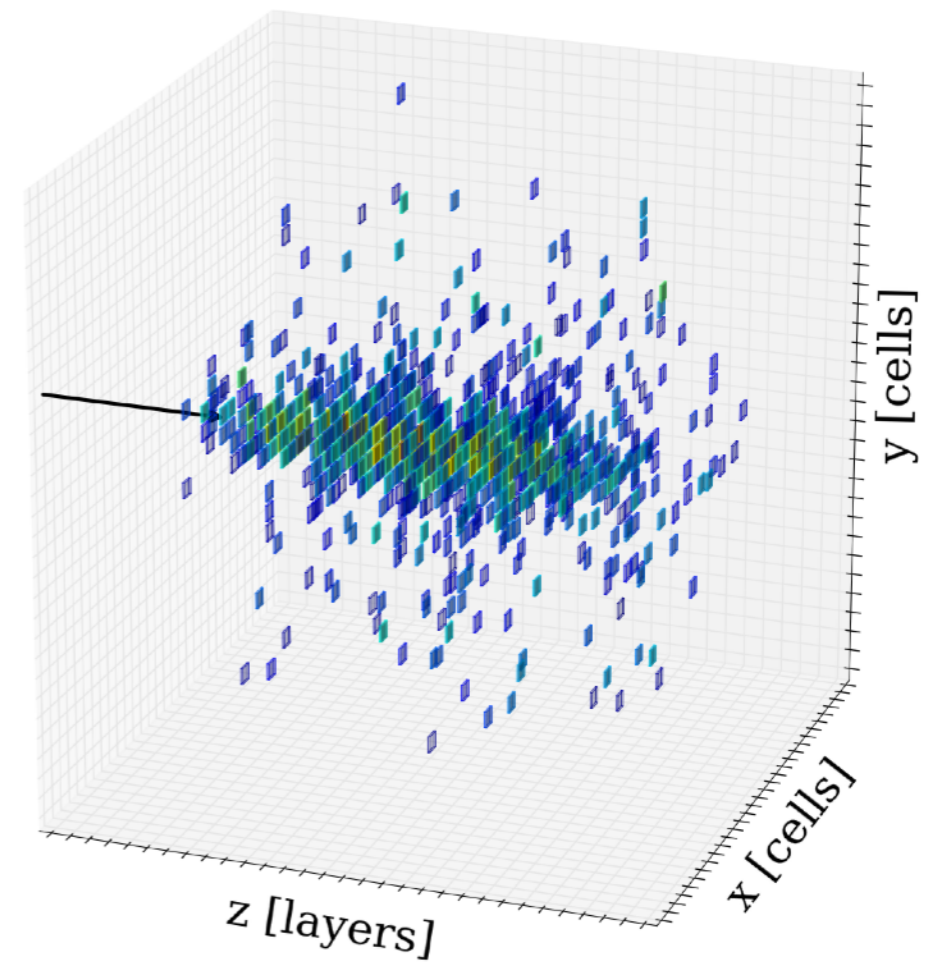
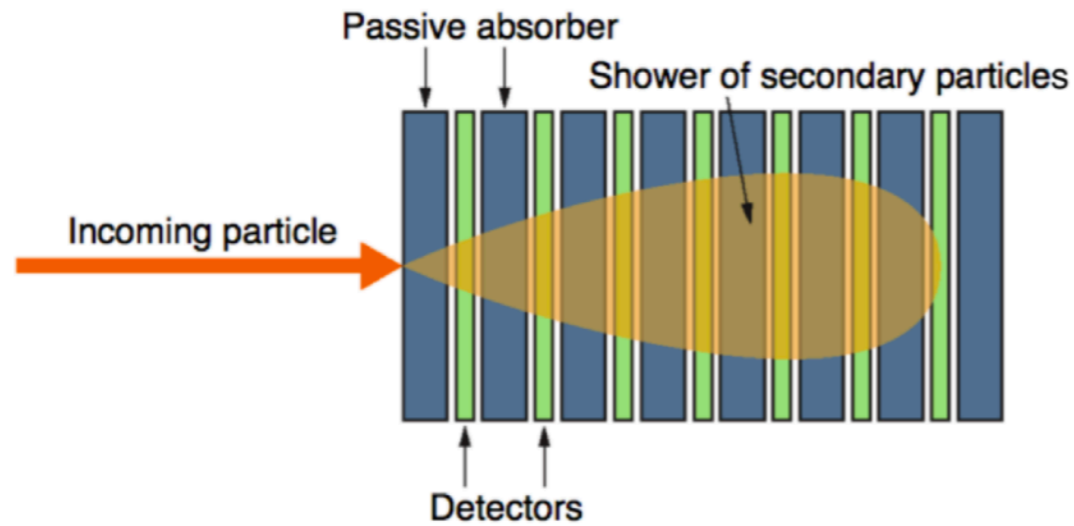


Simulation targets



How to represent?

Simulation targets

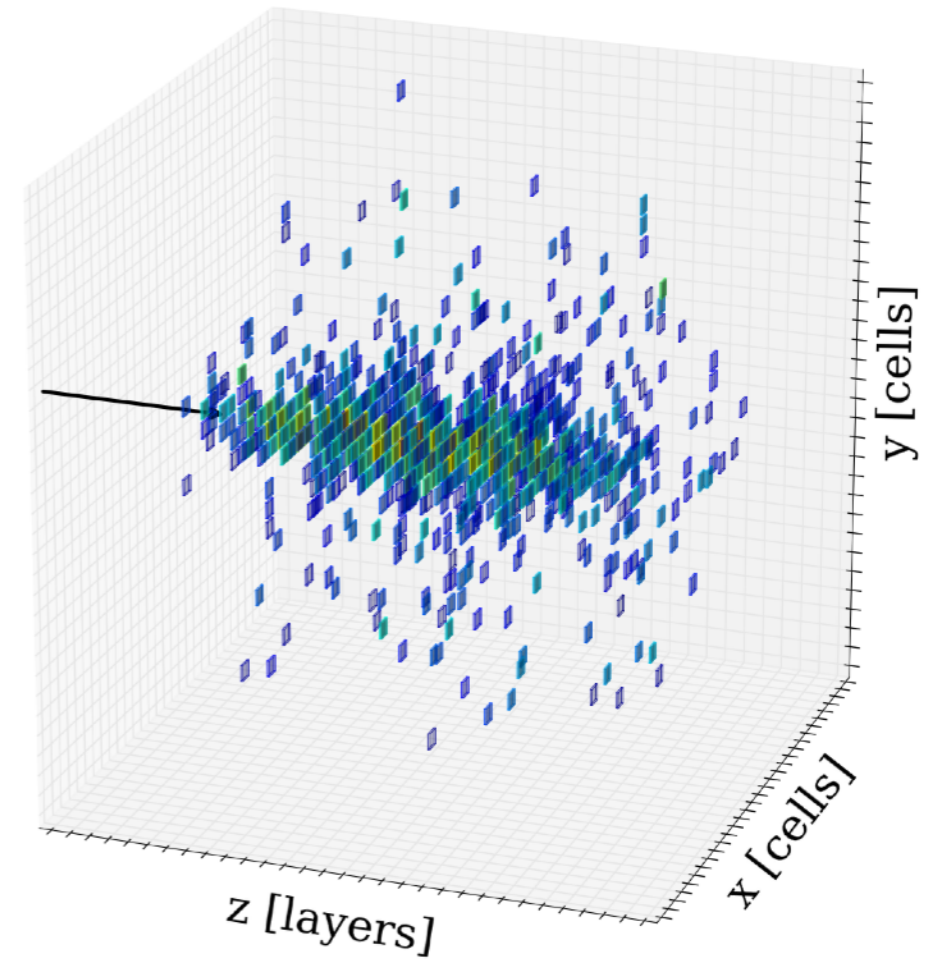
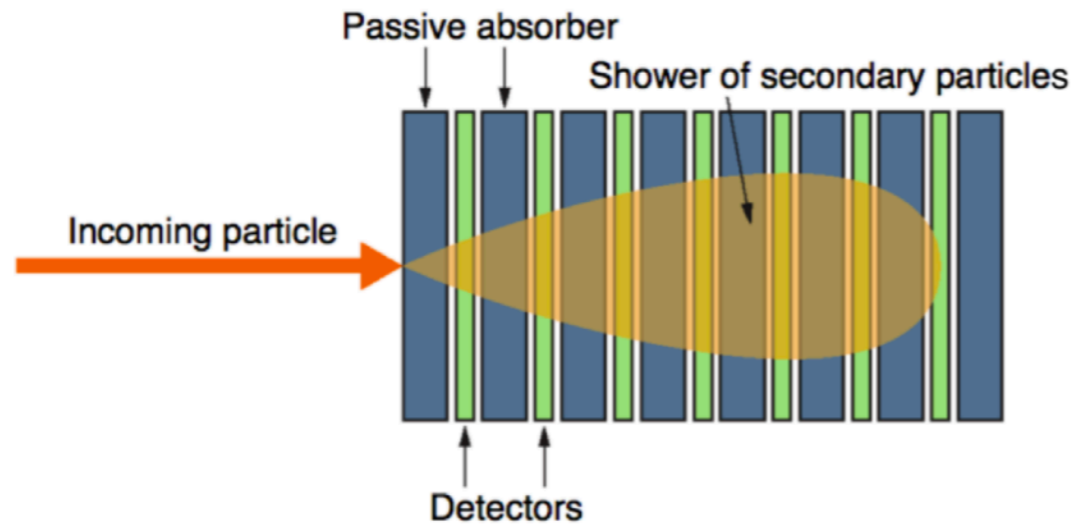


How to represent?

Tabular data:

Easy, insufficient for high-dimensions

Simulation targets

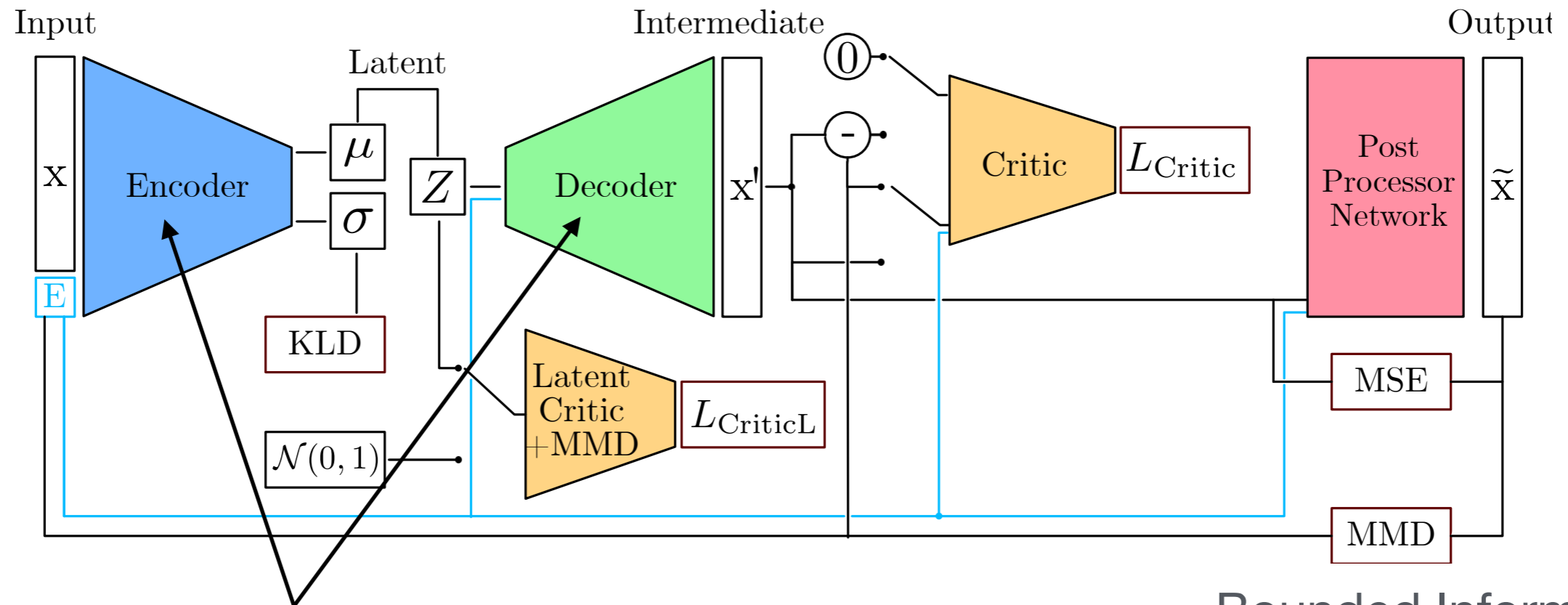


How to represent?

Tabular data

Fixed grid (voxels)

Generative results



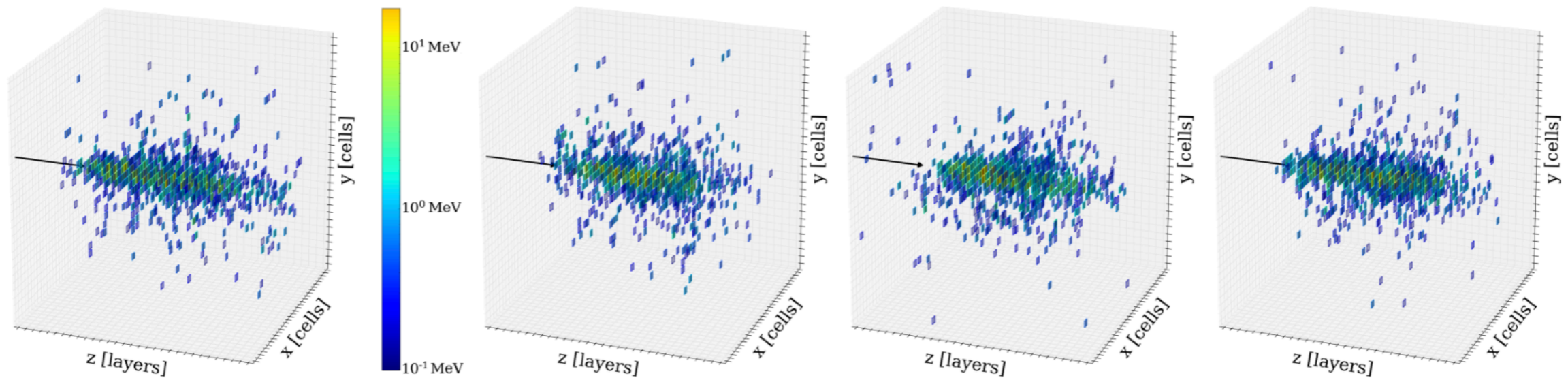
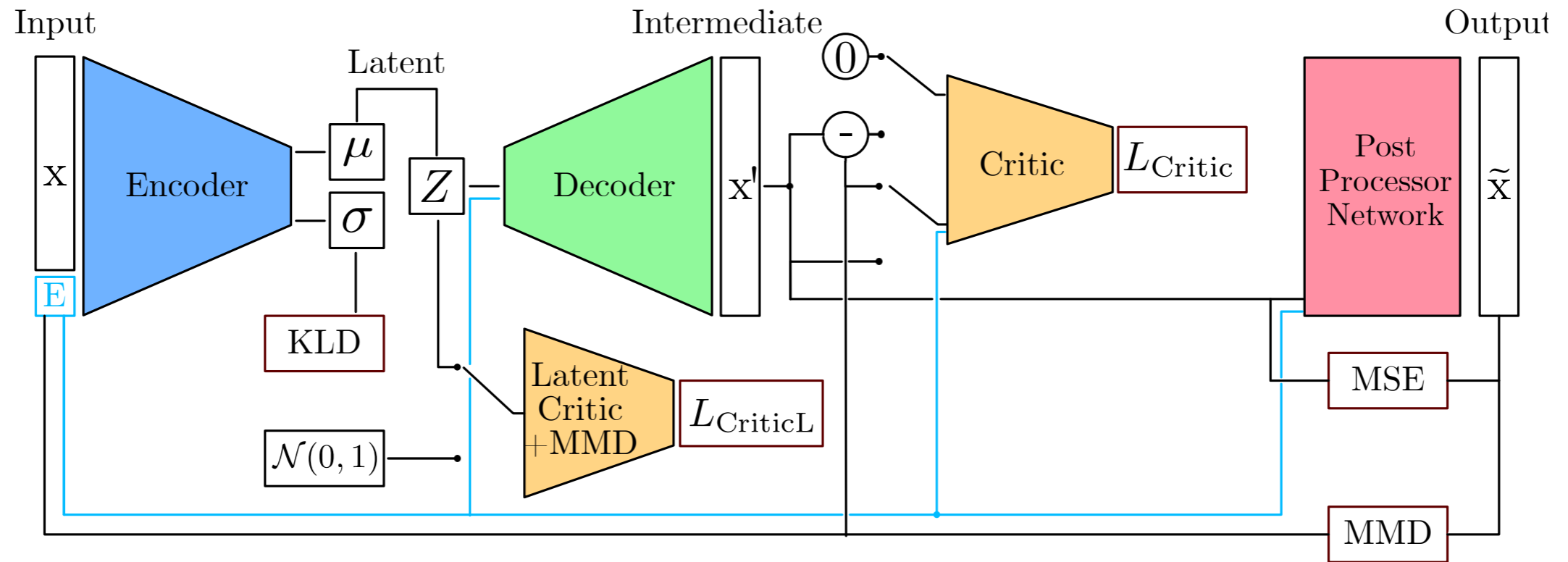
(Transposed) Convolution

Bounded Information Bottleneck AE

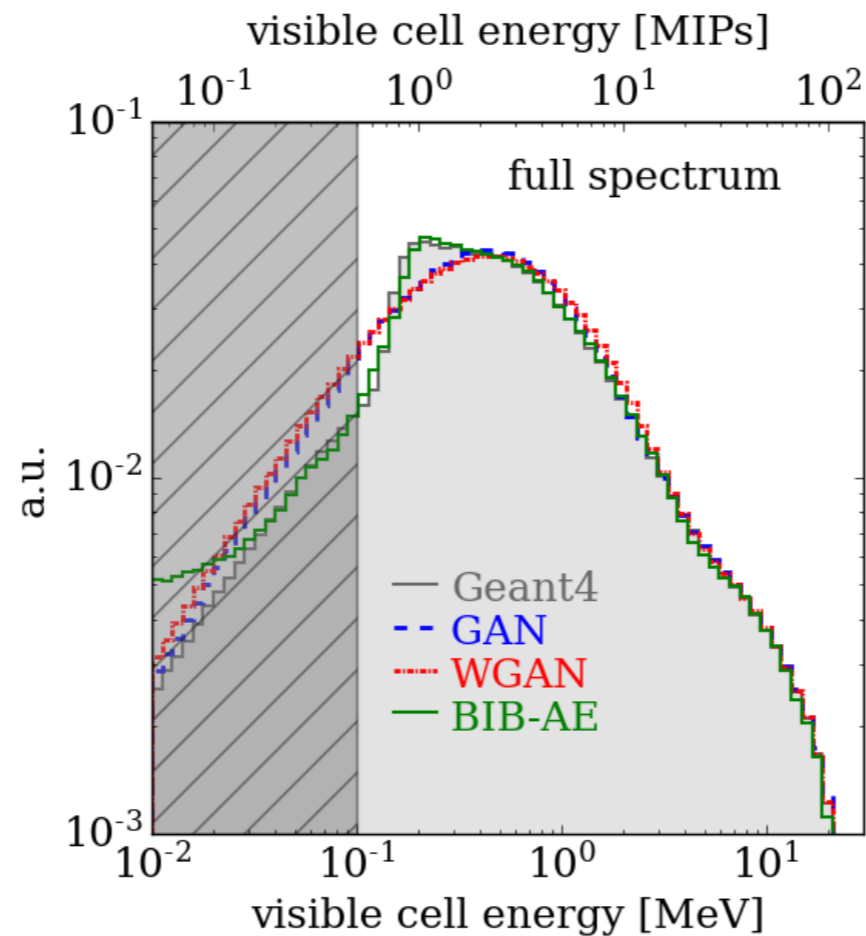
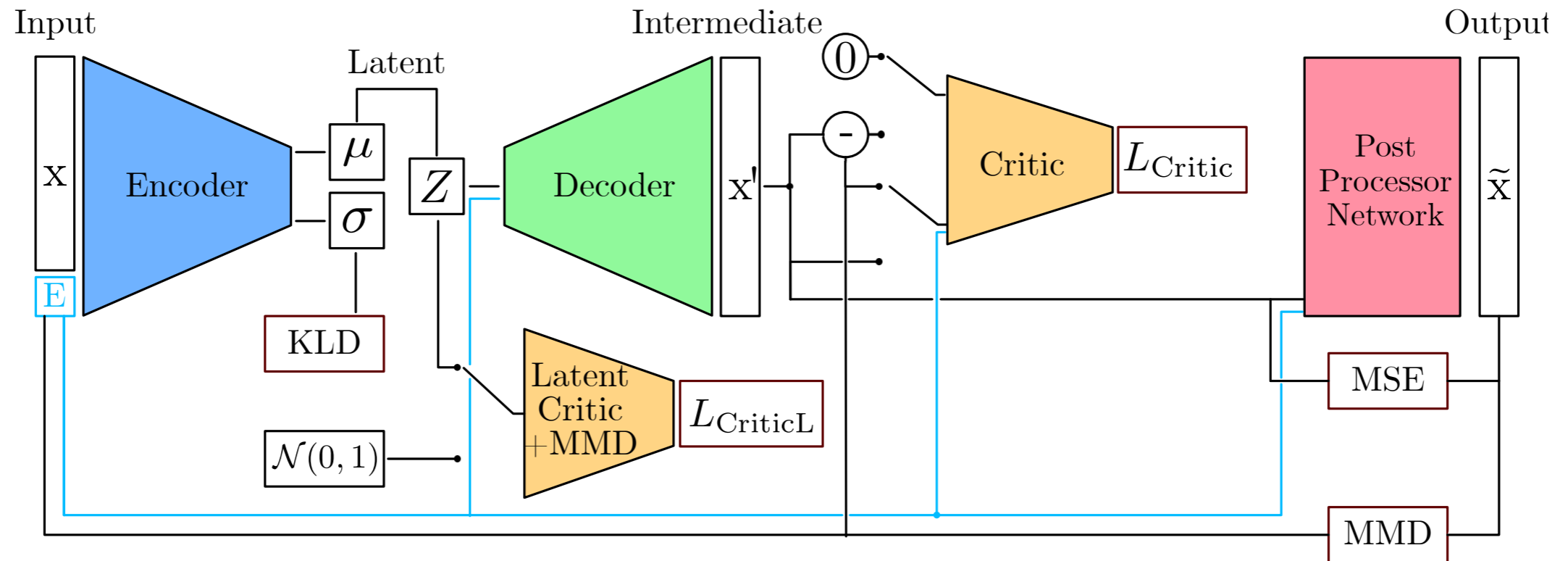
BIB-AE (GAN + VAE)

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

Generative results



Generative results

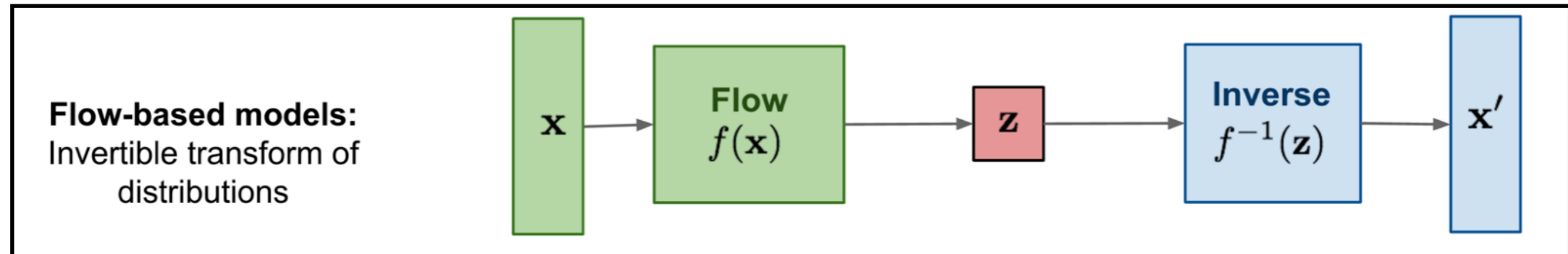


Go with the...

(Normalising) Flows



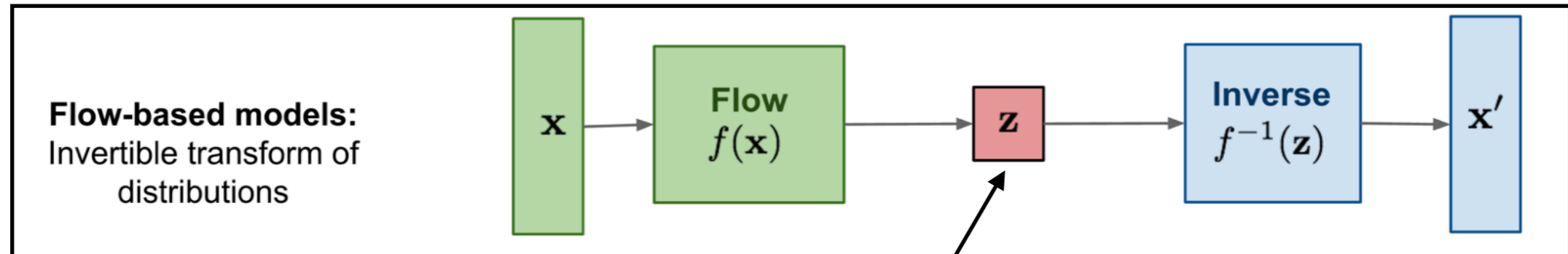
Generative models



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

Can we make this exact?

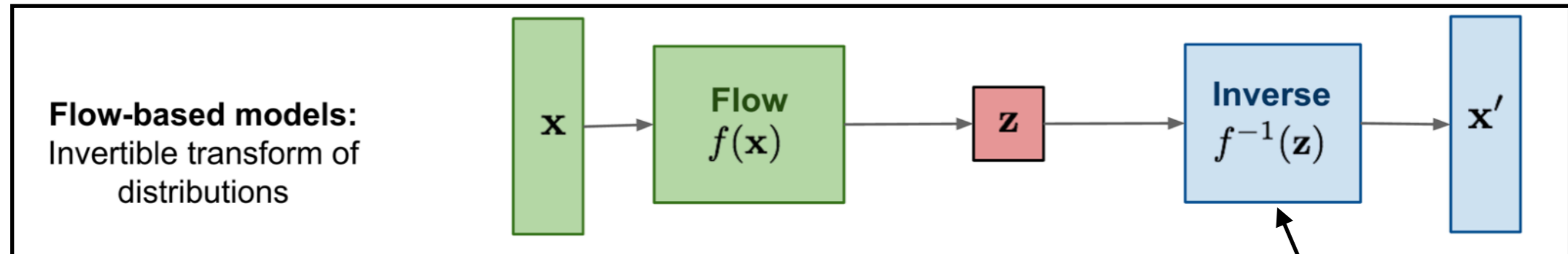
Generative models



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

Learn a diffeomorphism between data and latent-space

Generative models

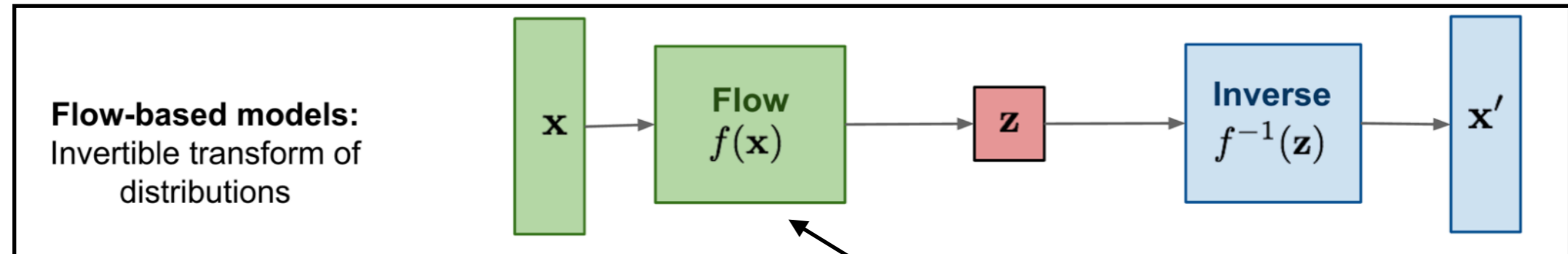


f^{-1} is not a learned inversion, but exact inverse by construction

Learn a diffeomorphism between data and latent-space

Bijjective, invertable

Generative models



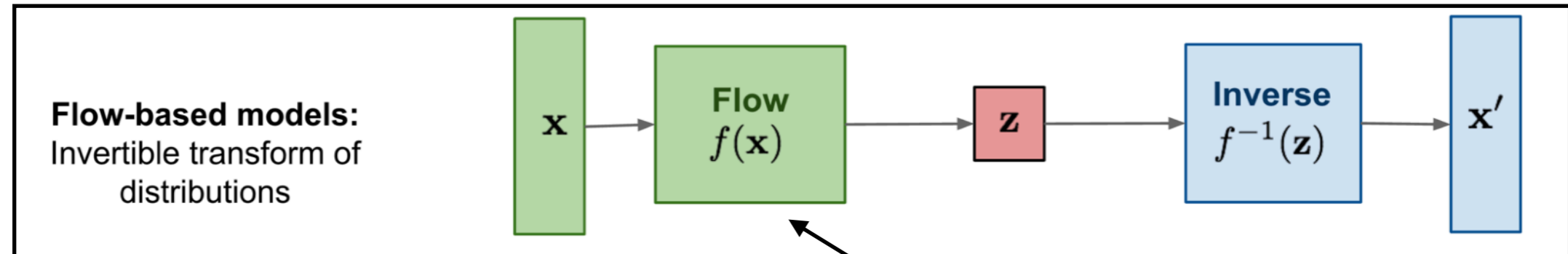
Learn a diffeomorphism between data and latent-space

Take into account Jacobian determinant to evaluate probability density

Bijjective, invertable

Learn likelihood of data

Generative models



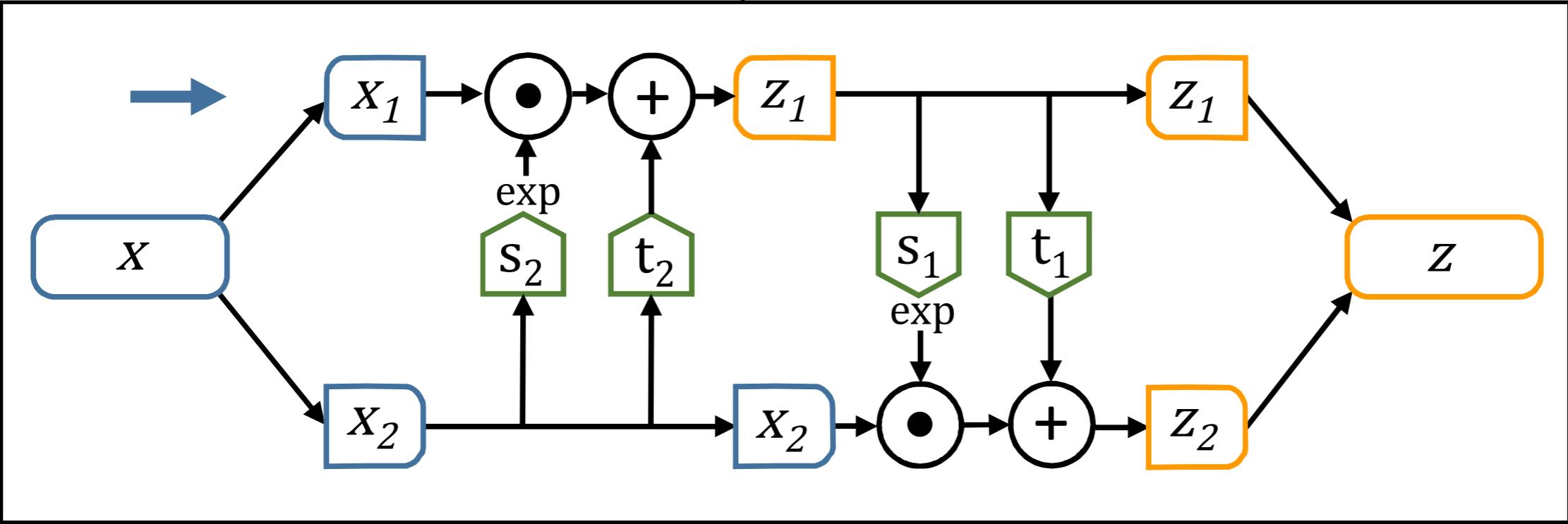
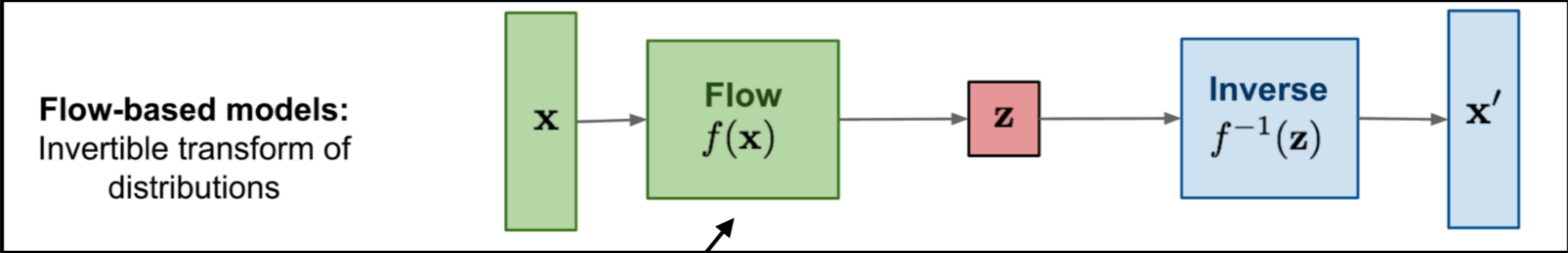
2 challenges:

Invertible

Easy-to-calculate Jacobean

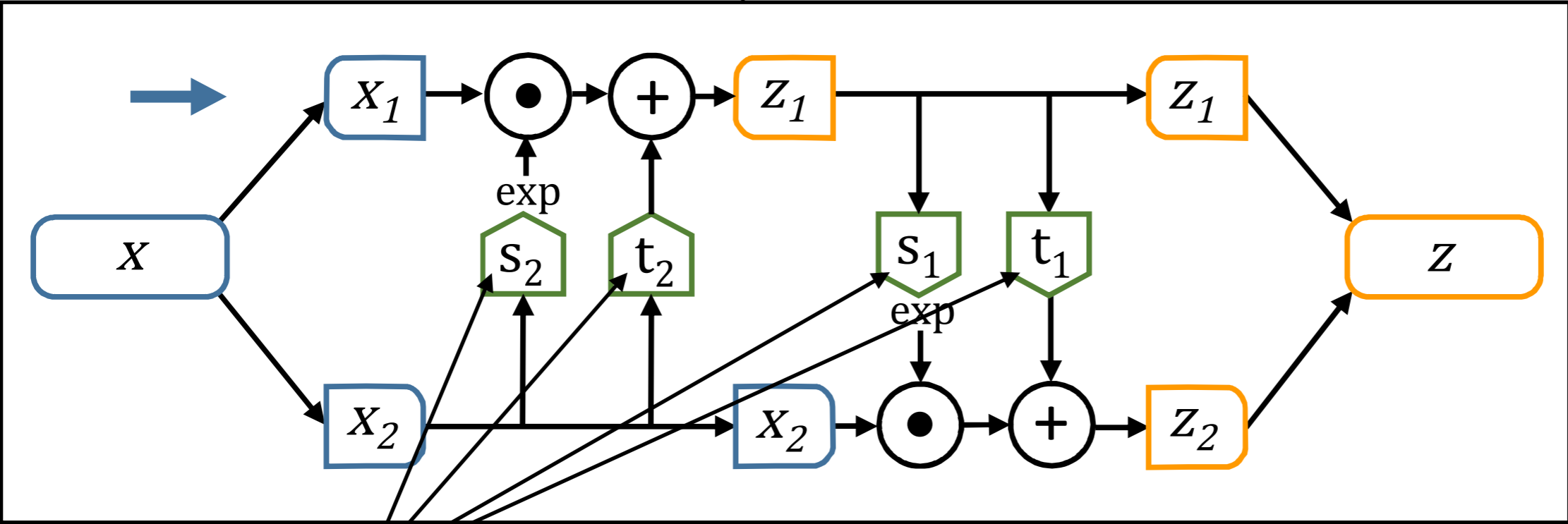
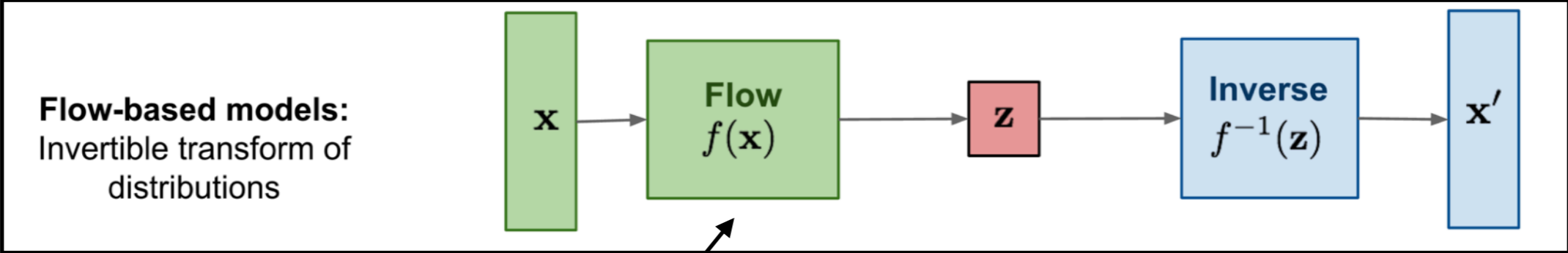
Take into account Jacobian determinant to evaluate probability density

Coupling flows



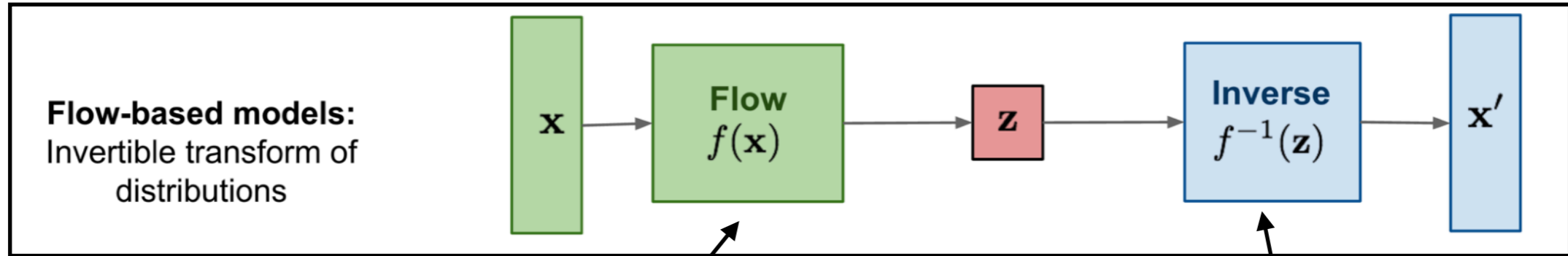
Coupling layers: Not the most expressive, but useful for illustration/understanding

Coupling flows

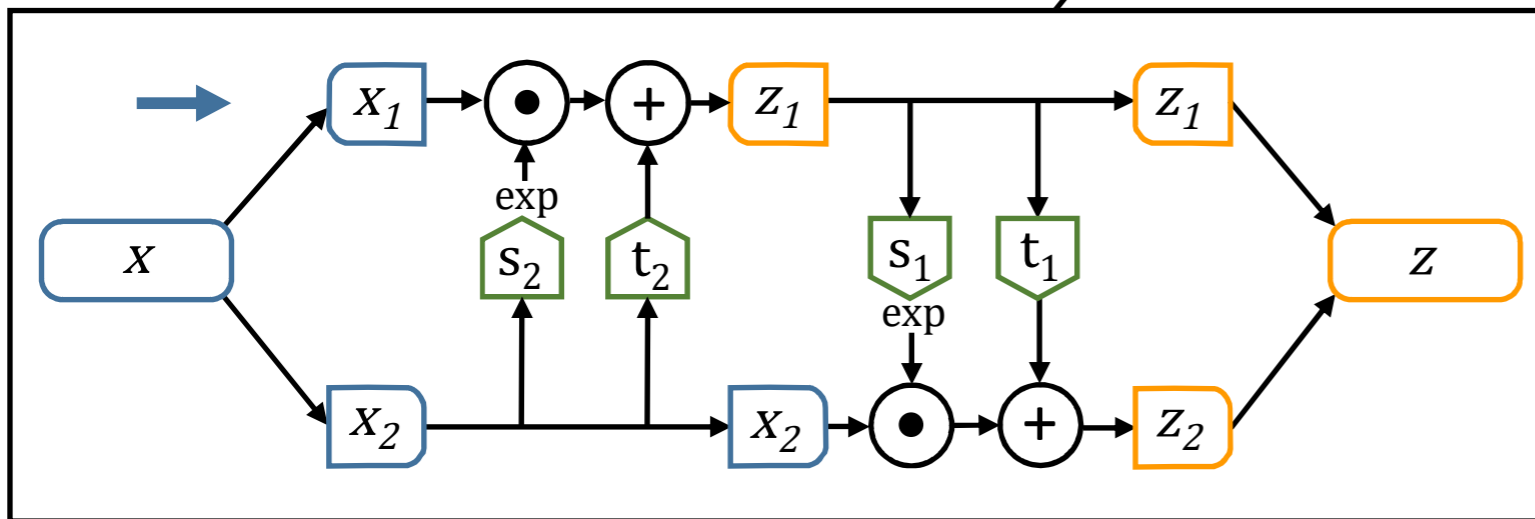


Simple (e.g. dense)
neural networks

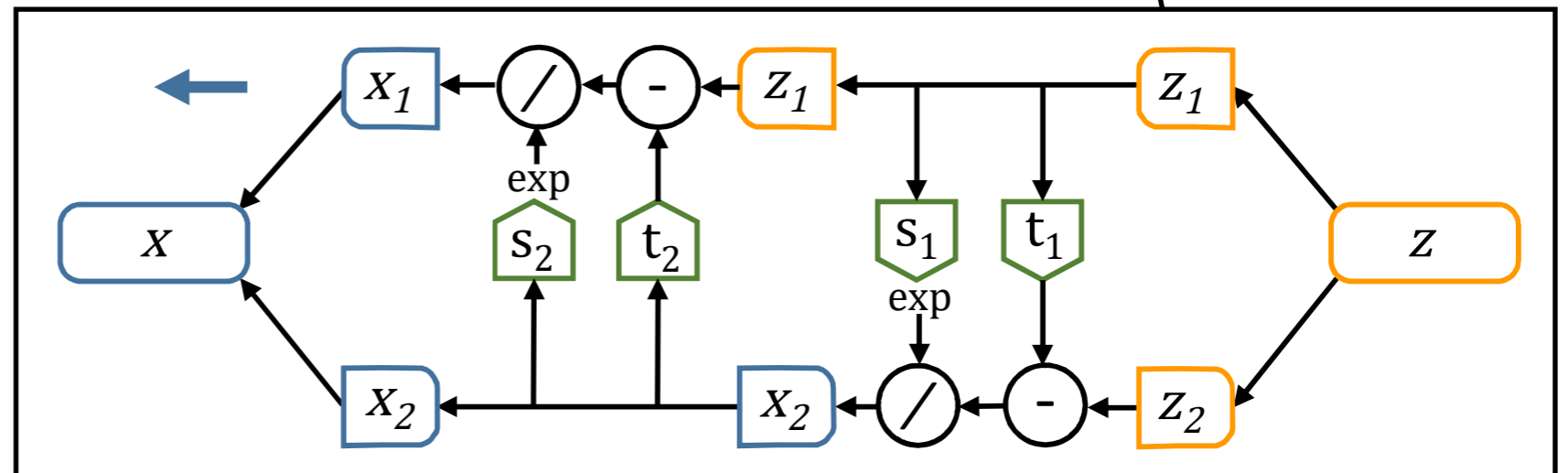
Coupling flows



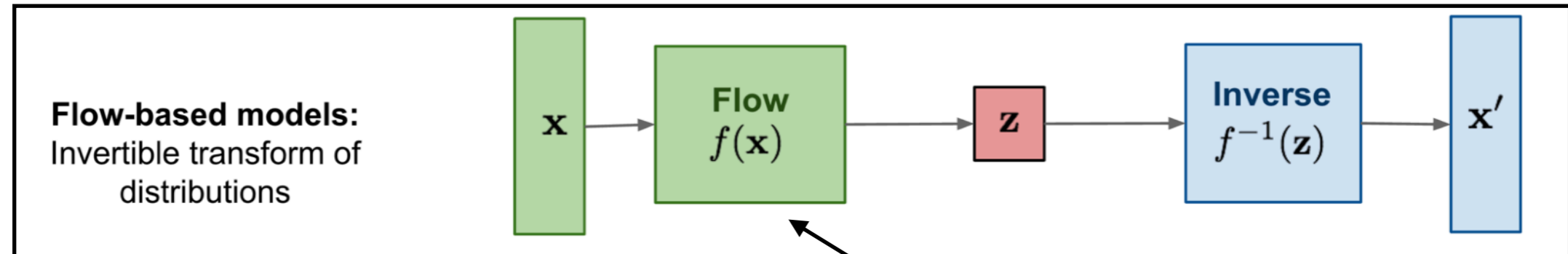
Forward direction



Inverse direction



Generative models



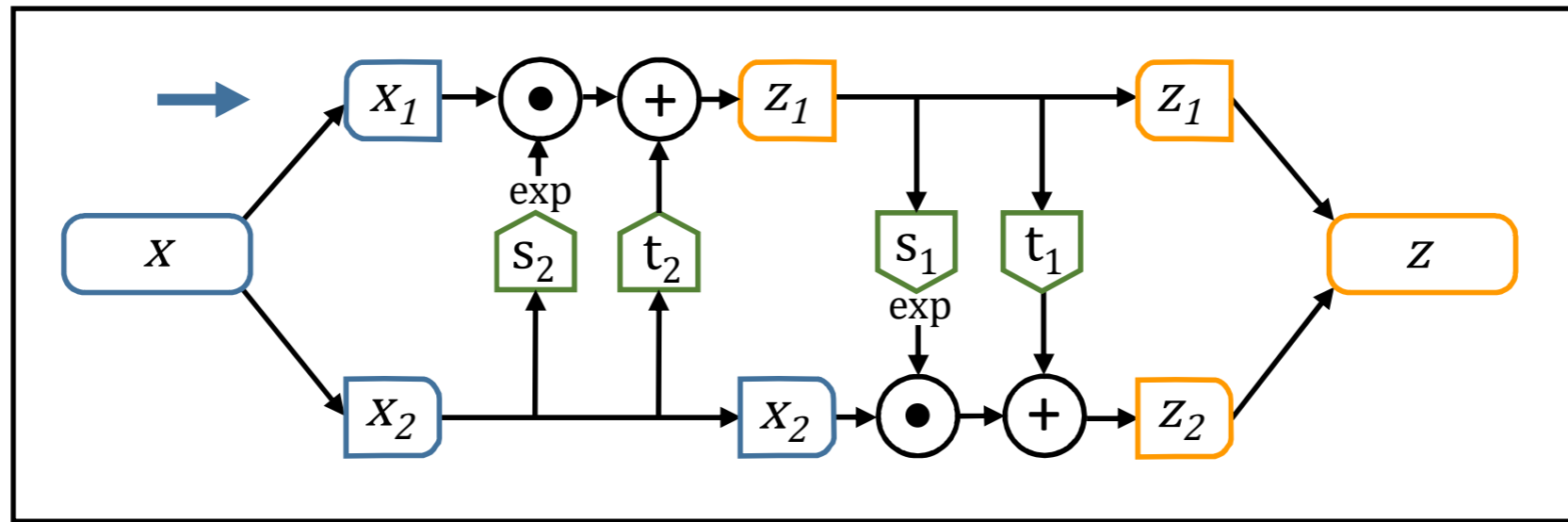
2 challenges:

Invertible

Easy-to-calculate **Jacobian**

Take into account Jacobian determinant to evaluate probability density

Calculating Jacobian determinant



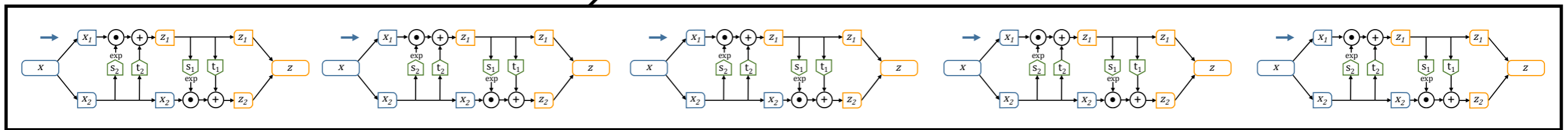
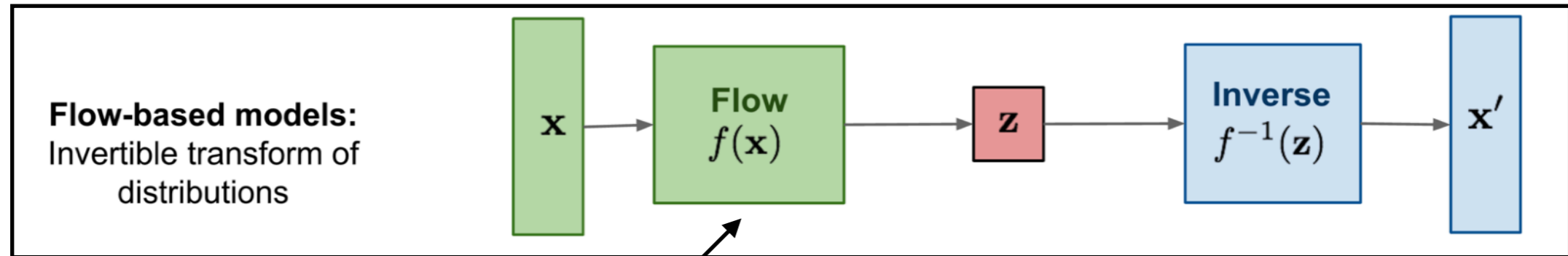
$$\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \xrightarrow{f_1} \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{x}_2 \end{pmatrix} \xrightarrow{f_2} \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{pmatrix} \text{ with } \begin{matrix} \mathbf{x}_1 \xrightarrow{f_1} \mathbf{z}_1 = \mathbf{x}_1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

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

Triangular by construction

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

Composition



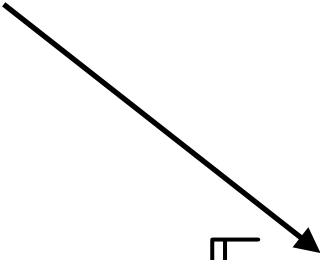
Composition of bijective functions
remains bijective

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

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

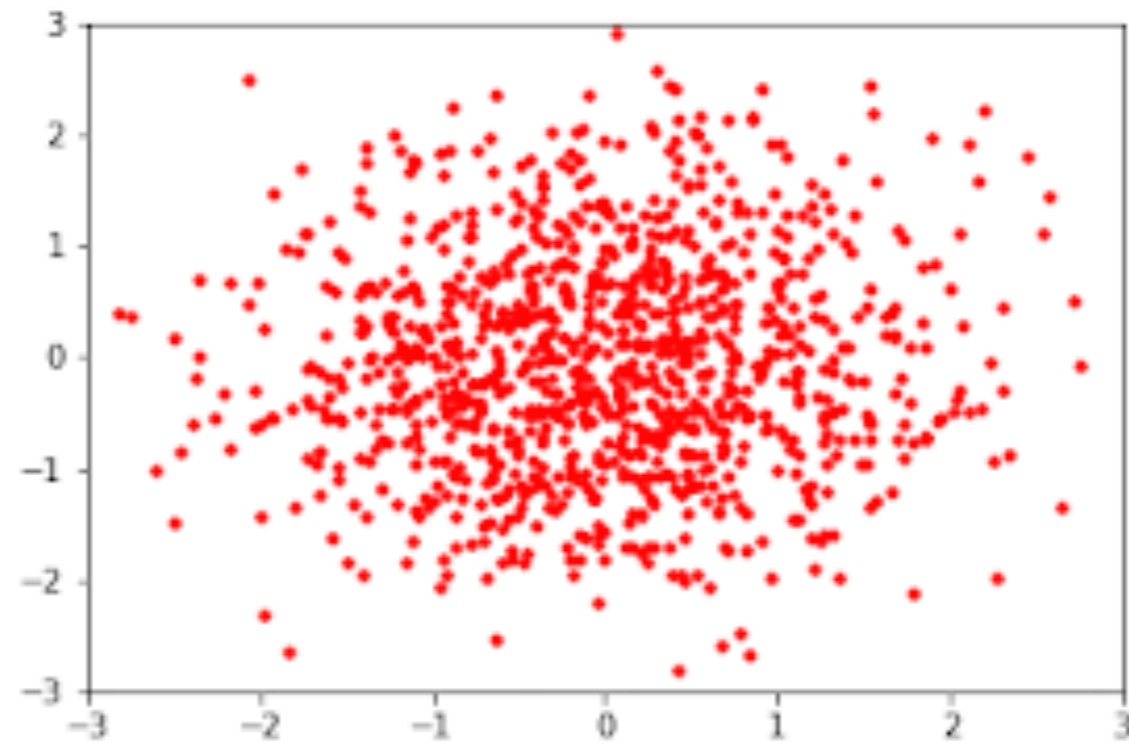
Contribution from Jacobian determinant



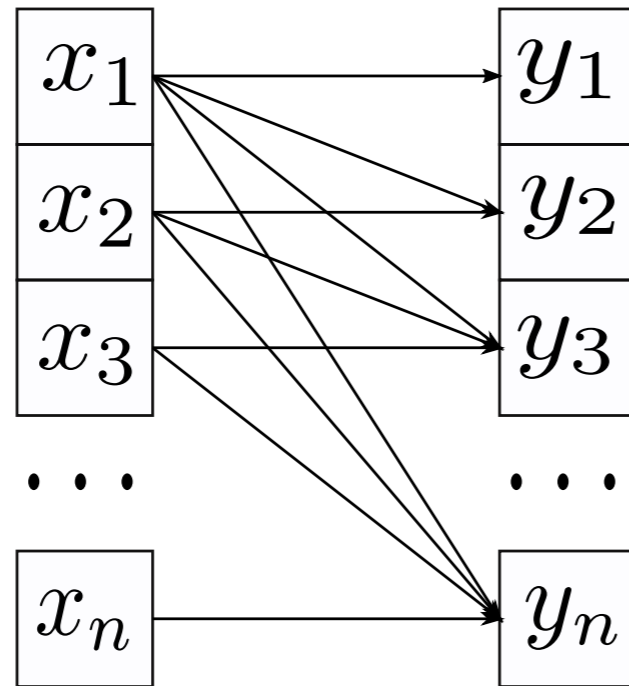
$$\det \mathbf{J} = \exp \left(\sum s(\mathbf{x}) \right)$$

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

Animation

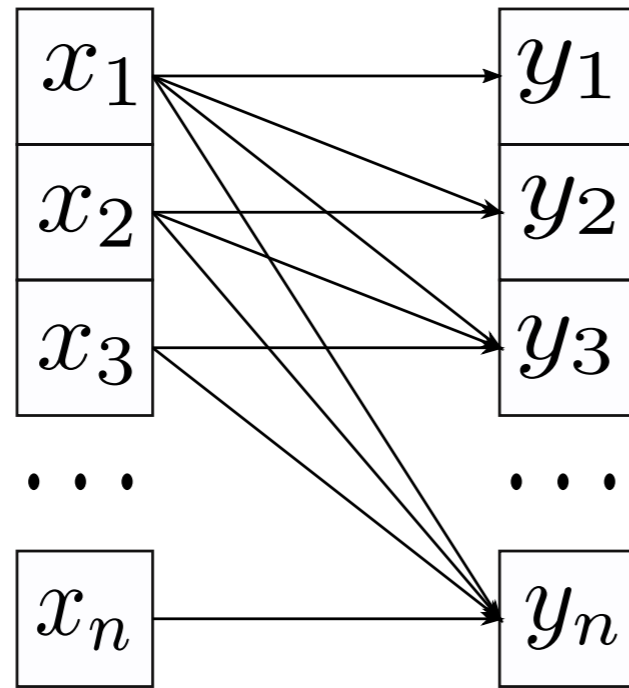


Autoregressive Flows



Alternative to coupling flows:
Outputs conditioned on previous inputs

Autoregressive Flows



Bijjective function

Parameters

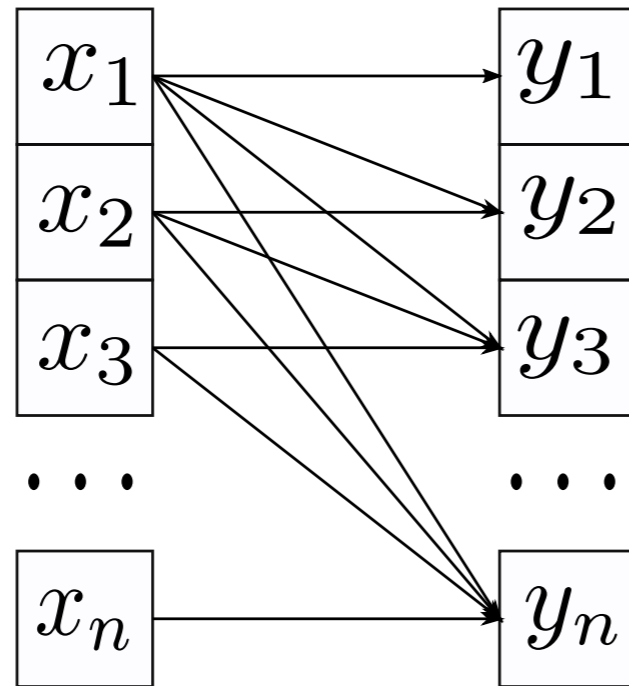
$$y_t = h(x_t; \Theta_t(\mathbf{x}_{1:t-1}))$$

Alternative to coupling flows:

Outputs conditioned on previous inputs

Again: simple Jacobian and invertible functions

Autoregressive Flows



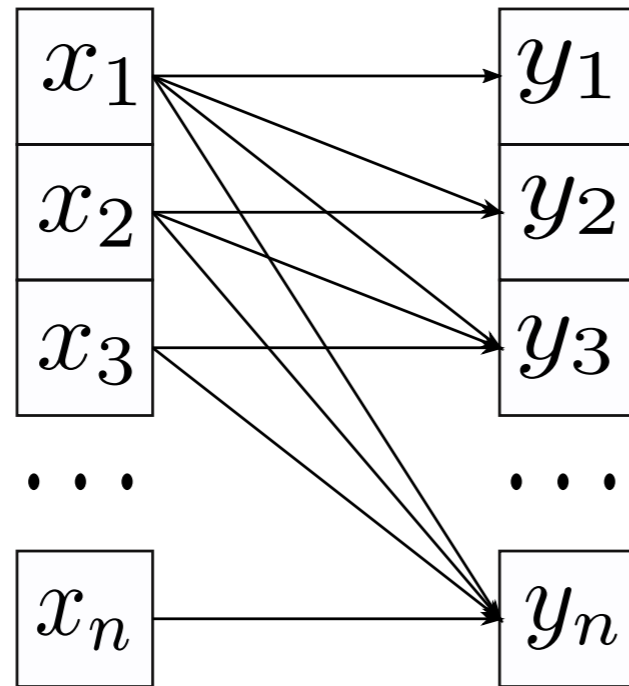
$$y_t = h(x_t; \Theta_t(\mathbf{x}_{1:t-1}))$$

Masked autoregressive flow (MAF):

Fast: Data \rightarrow latent space

Slow: Latent space \rightarrow data

Autoregressive Flows

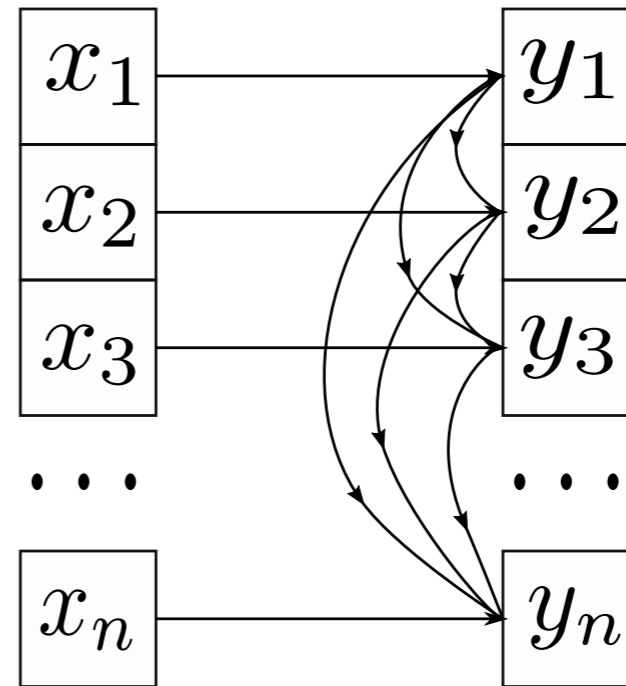


$$y_t = h(x_t; \Theta_t(\mathbf{x}_{1:t-1}))$$

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 §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

Comments on Flows

Only scratched the surface:
more constructions available

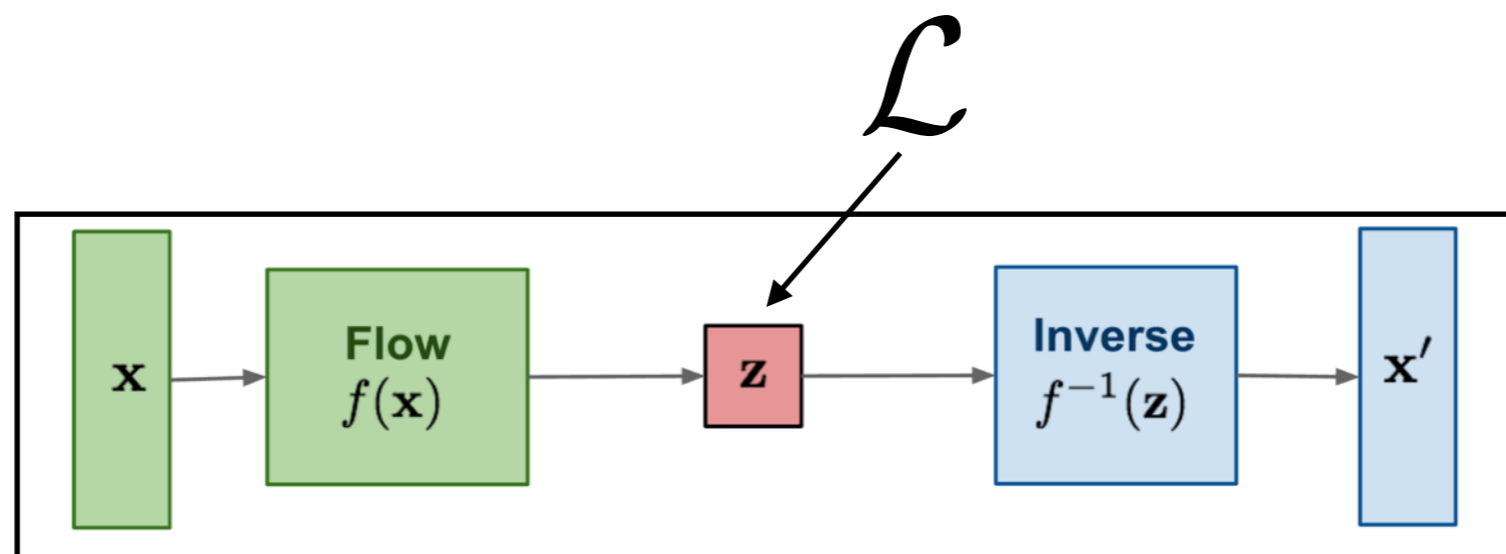
Exact learning of likelihood
→ Better generative fidelity
→ Can evaluate likelihood of
data

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 | |

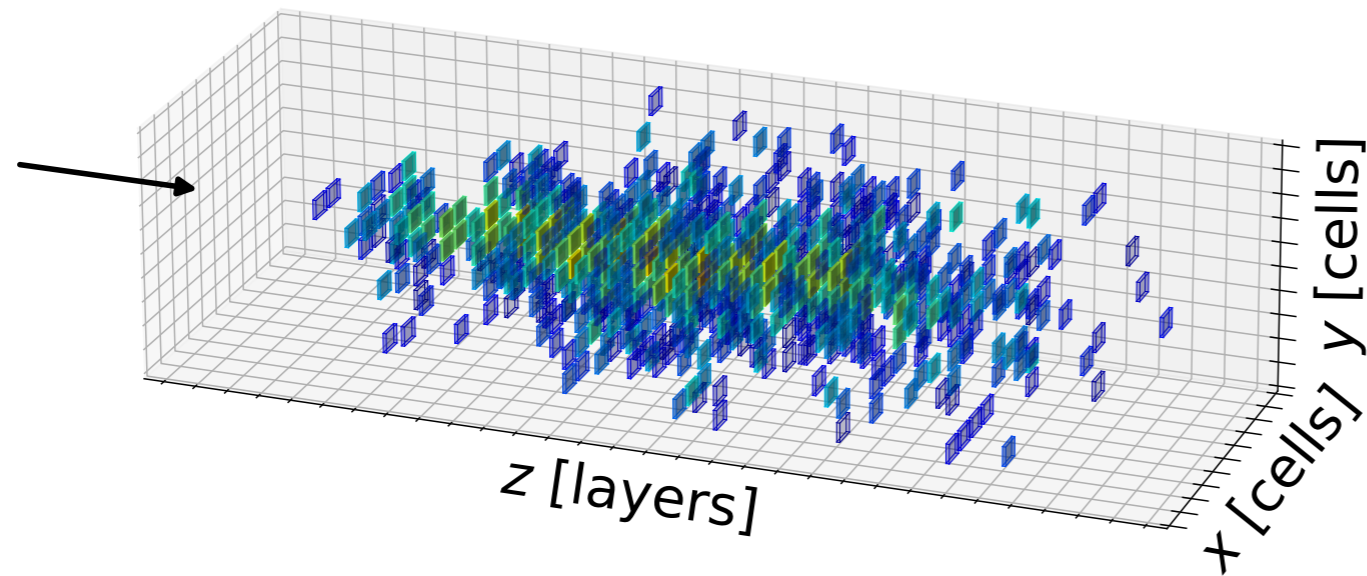
§3.4.4 Coupling functions

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



Applications II

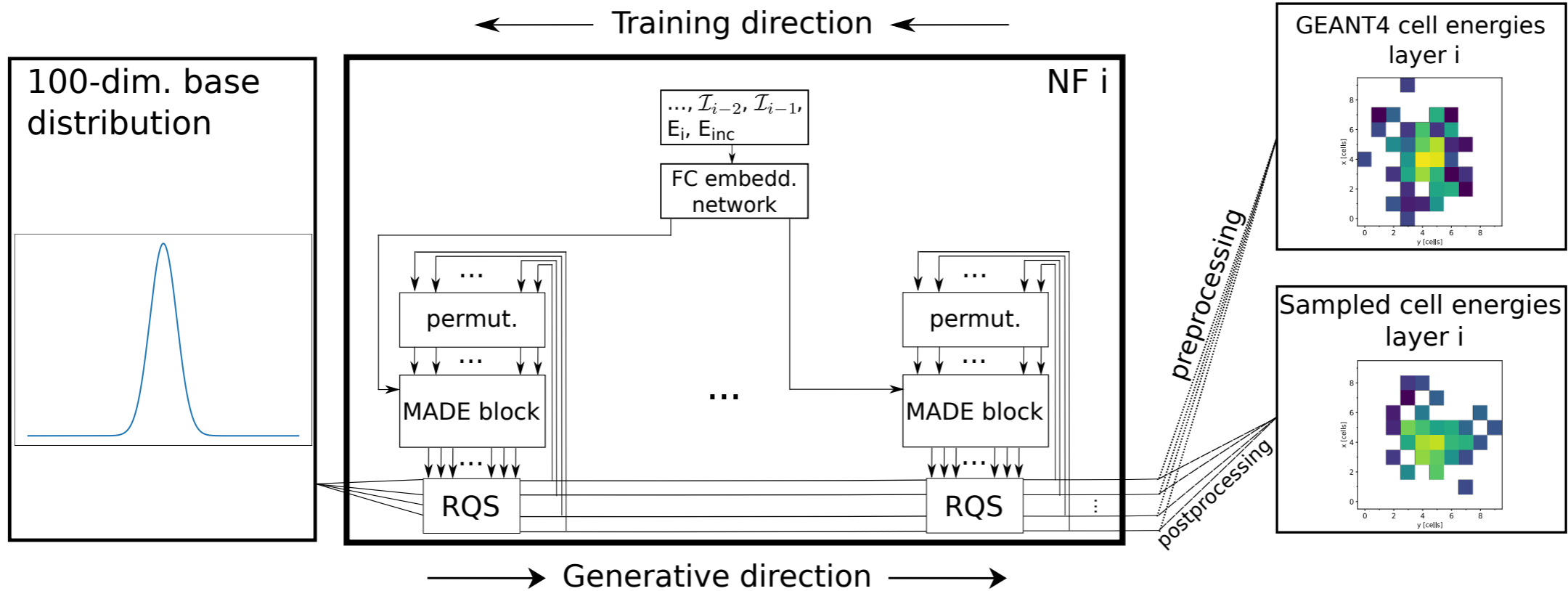
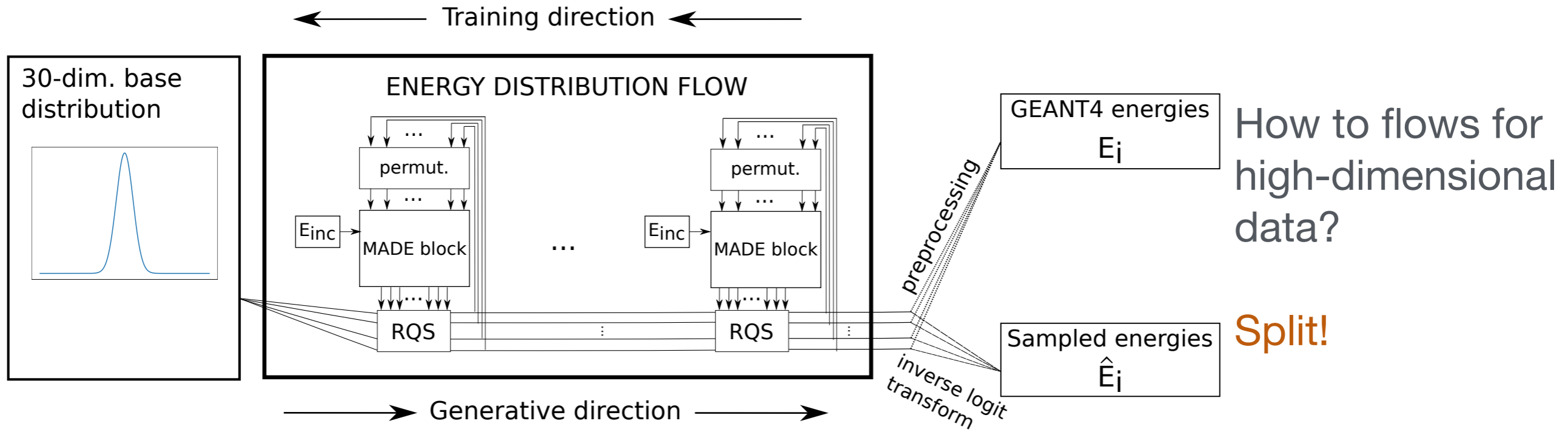
Generative results II



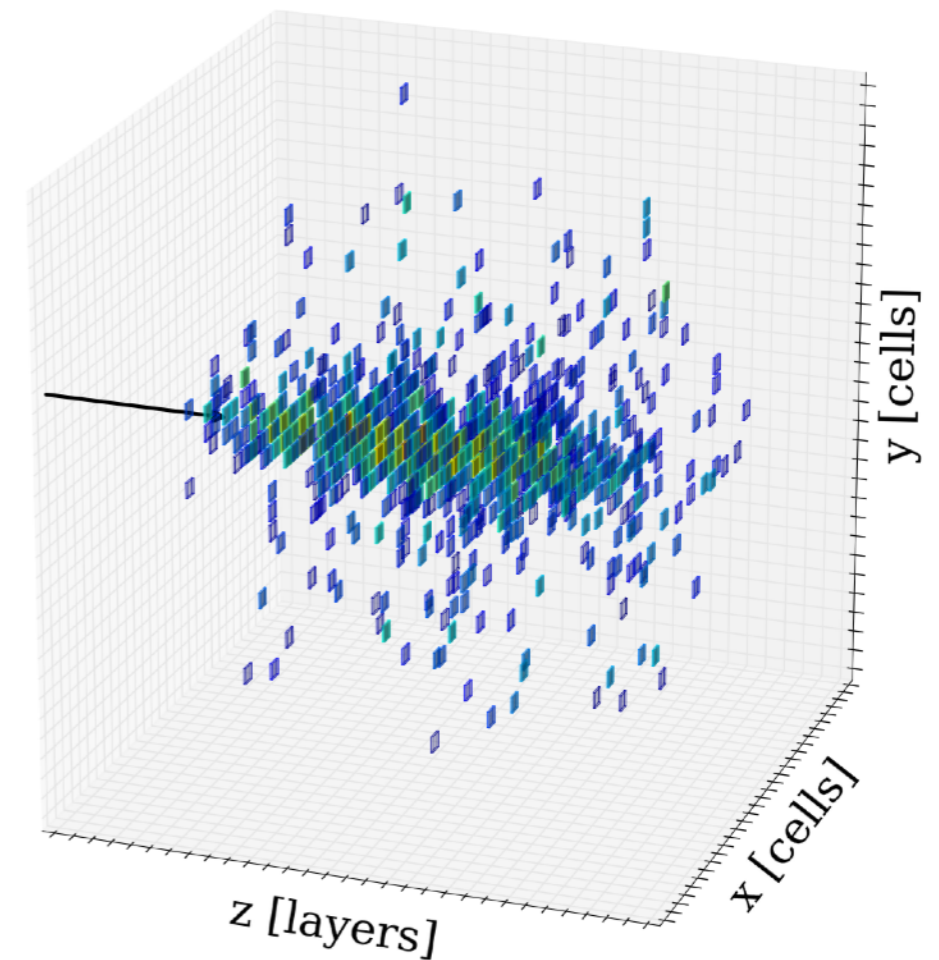
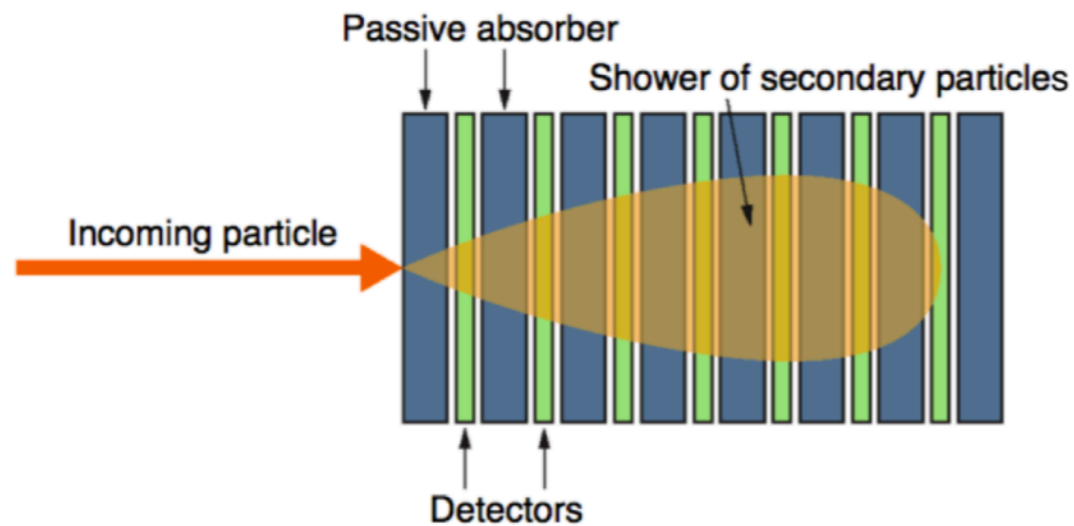
10x10 cells / layer
30 layers

How to flows for
high-dimensional
data?

Generative results II



Simulation targets



How to represent?

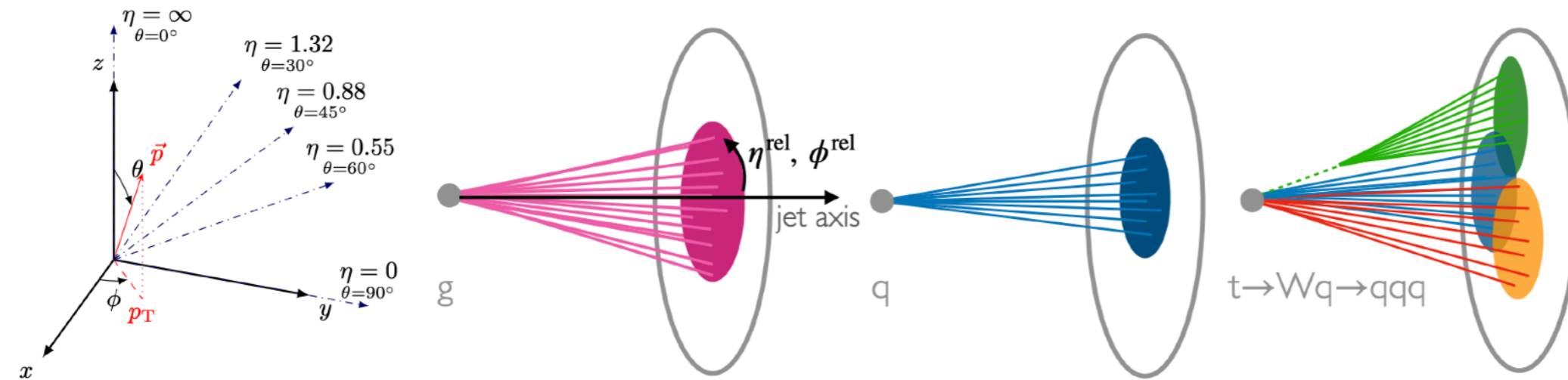
Tabular data

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

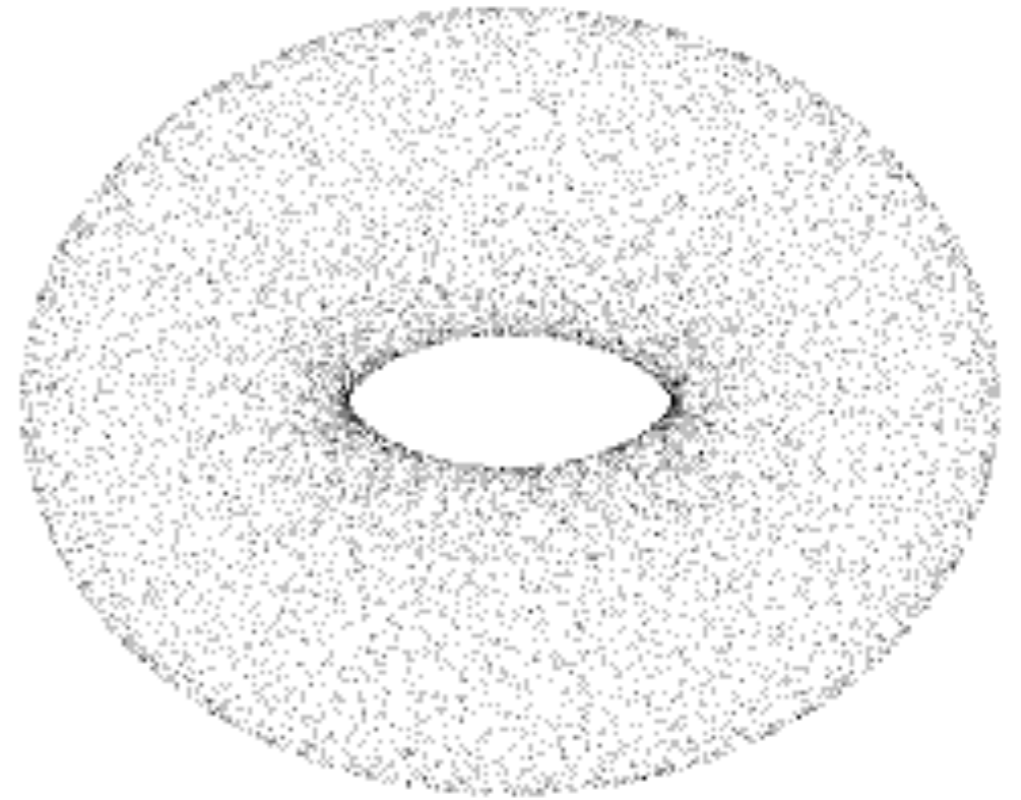
Why?

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 $\{x^j\}_{j=1 \dots N}$ ^{Examples}

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

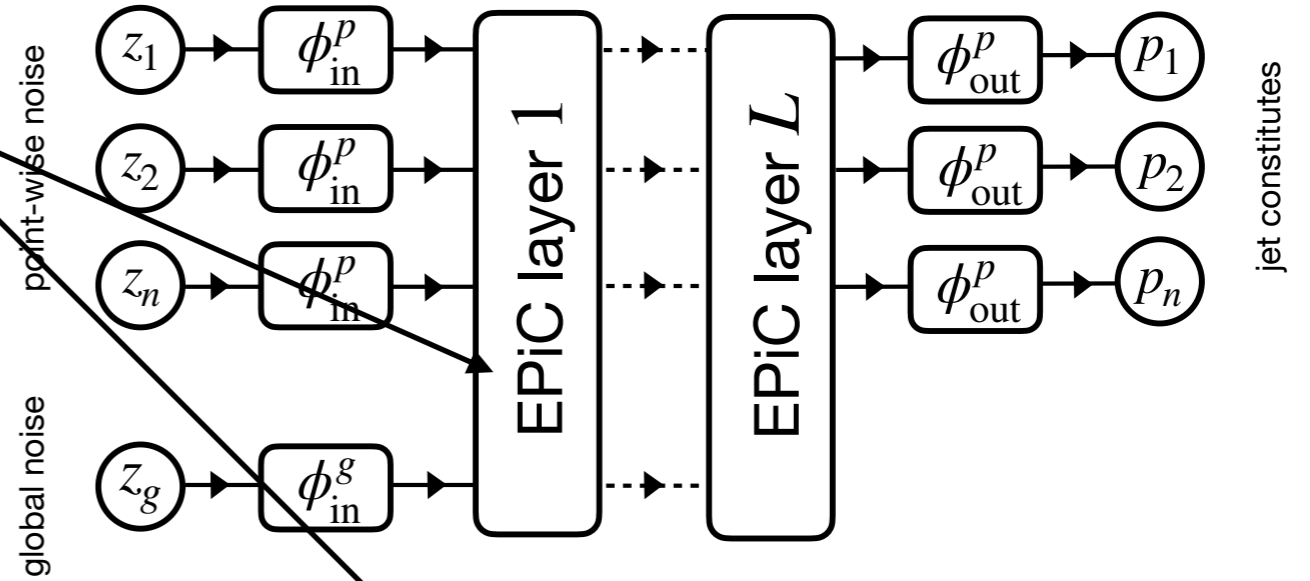
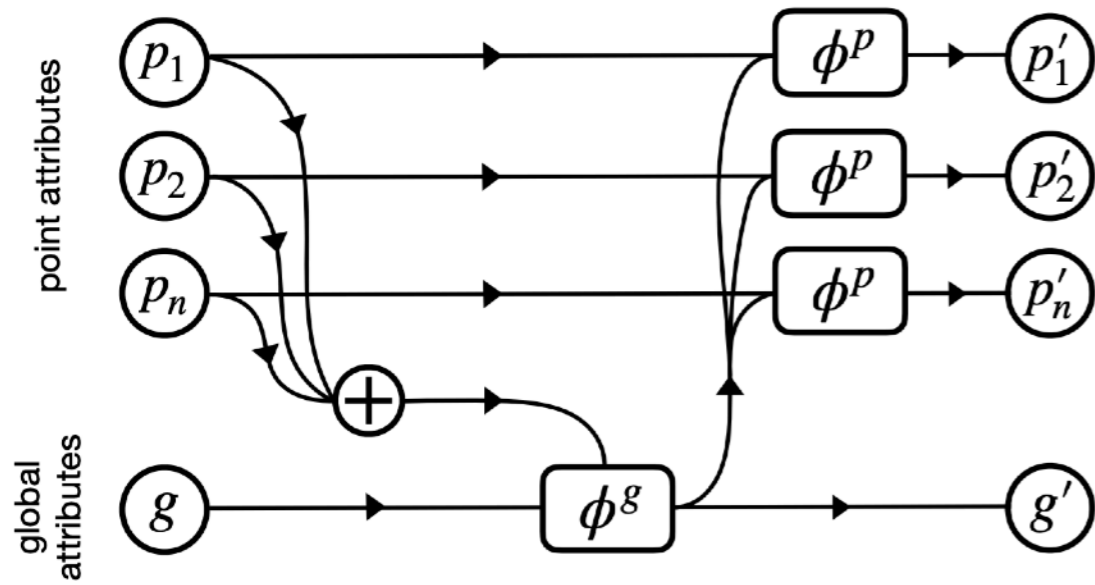
$$f(x_1, \dots, x_M) = \rho \left(\sum_{m=1}^M \phi(x_m) \right) \quad (18)$$

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 ϕ is independent of the function f .

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

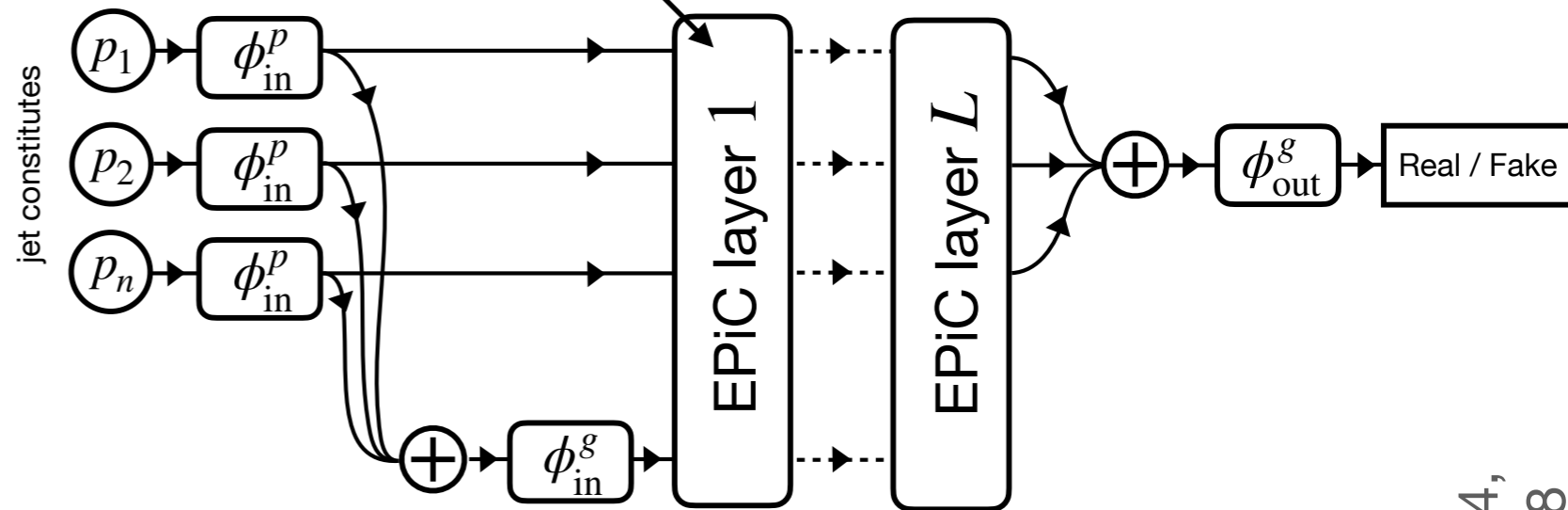
Latent Space Dimension

How to GAN with it



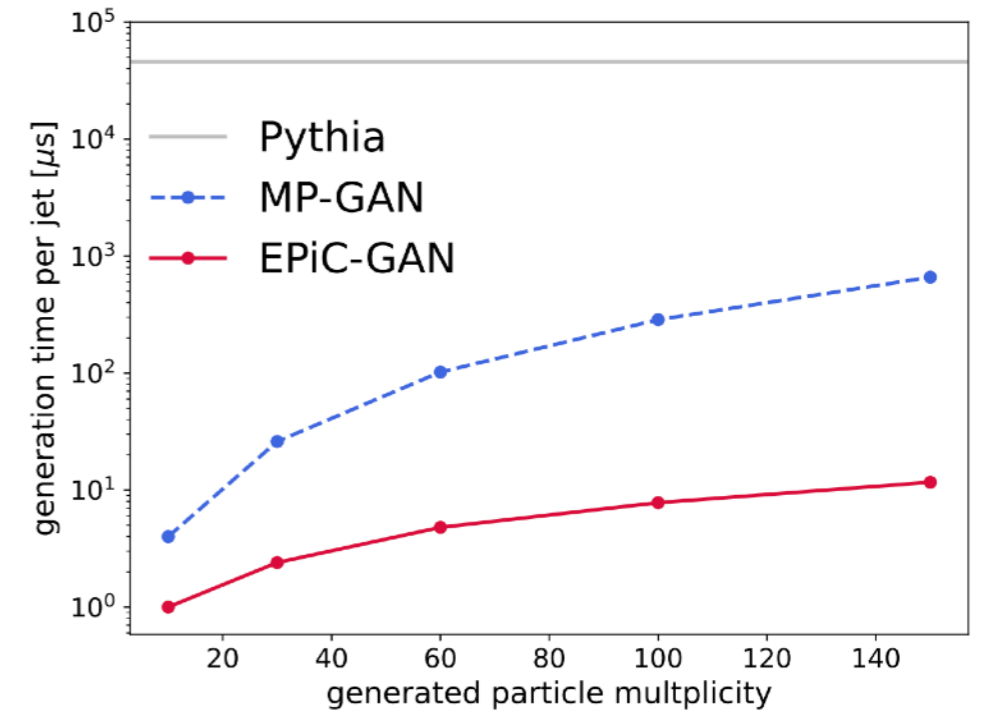
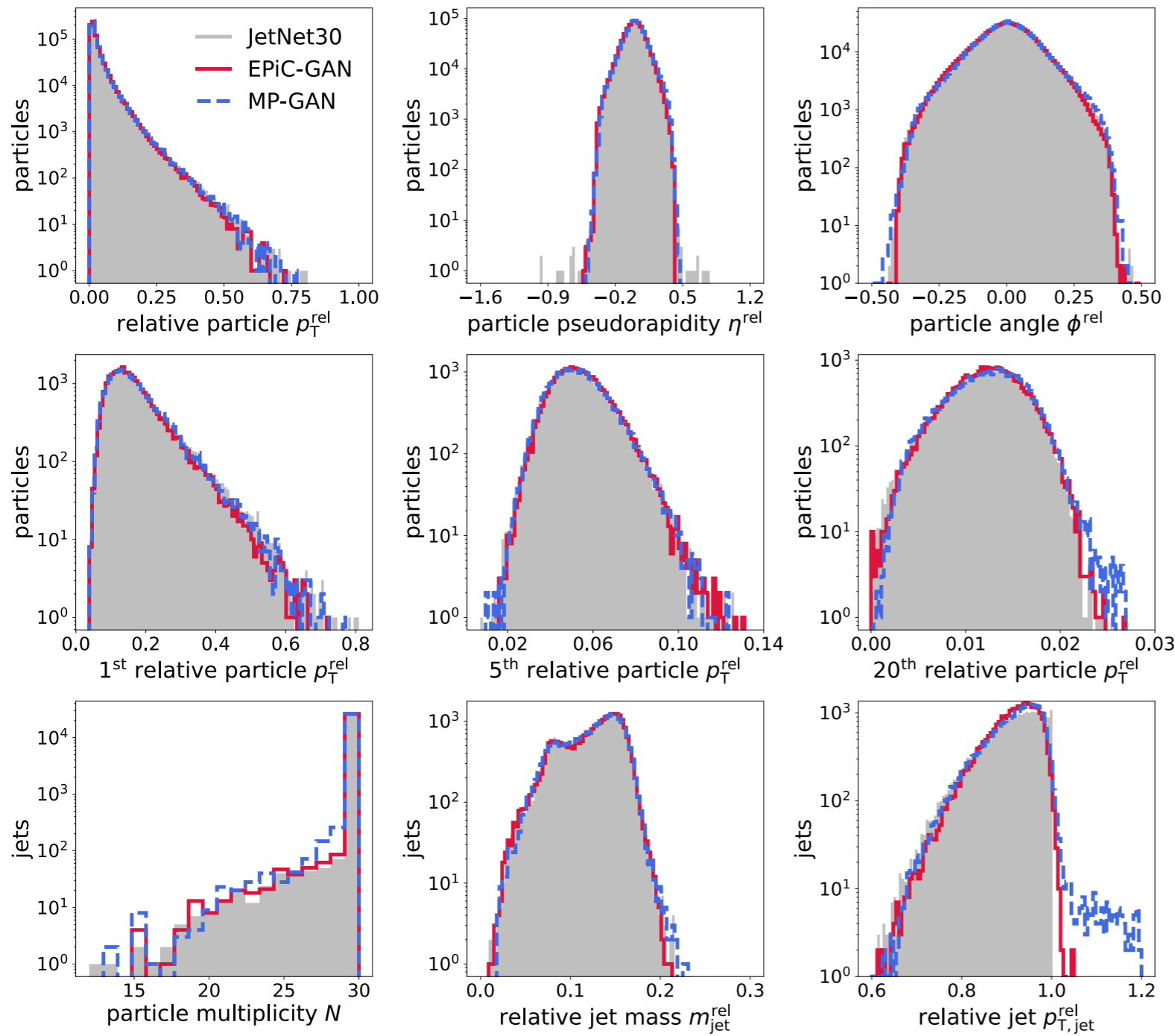
(a) Generator

Add **permutation symmetry** to GAN architecture



(b) Discriminator

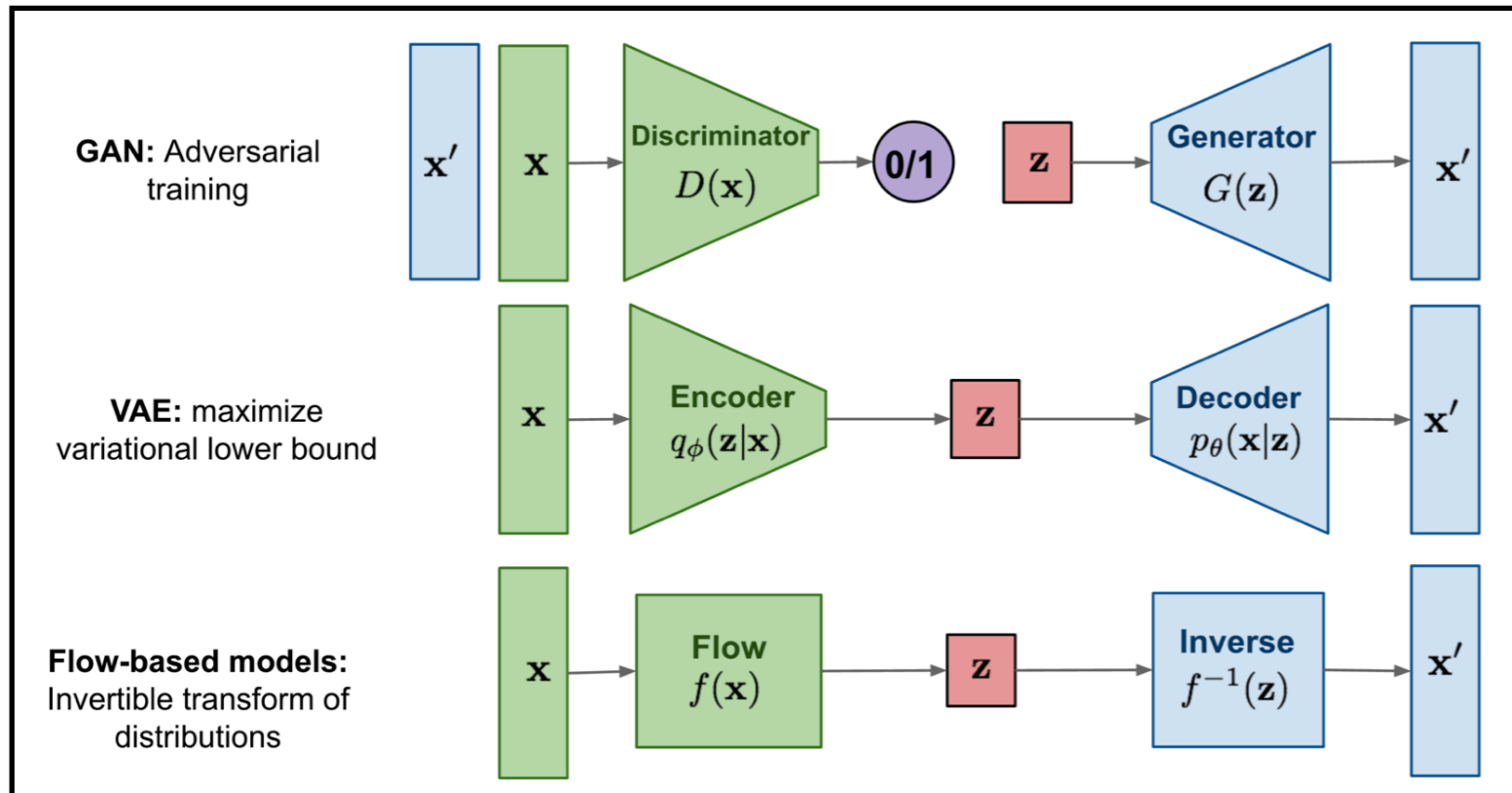
Generative results III



Point clouds/graph GANs
successfully generate jet
constituents

Closing I

Closing



Understand three basic generative architectures

First look at simulating fixed grid and point-cloud data

