



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386



# TOBIAS BUCK

27.9.2024 [TOBIAS.BUCK@IWR.UNI-HEIDELBERG.DE](mailto:TOBIAS.BUCK@IWR.UNI-HEIDELBERG.DE)

# GENERATIVE AI

# OUTLINE

1. Day 2: Generative AI
  - A. Variational Autoencoder
  - B. Normalizing Flows
  - C. Score matching / diffusion models
  - D. Simulation-Based Inference
  - E. Model Comparison & Model Misspecification

# Motivation

# MODELING COMPLEX PROBABILITY DISTRIBUTIONS

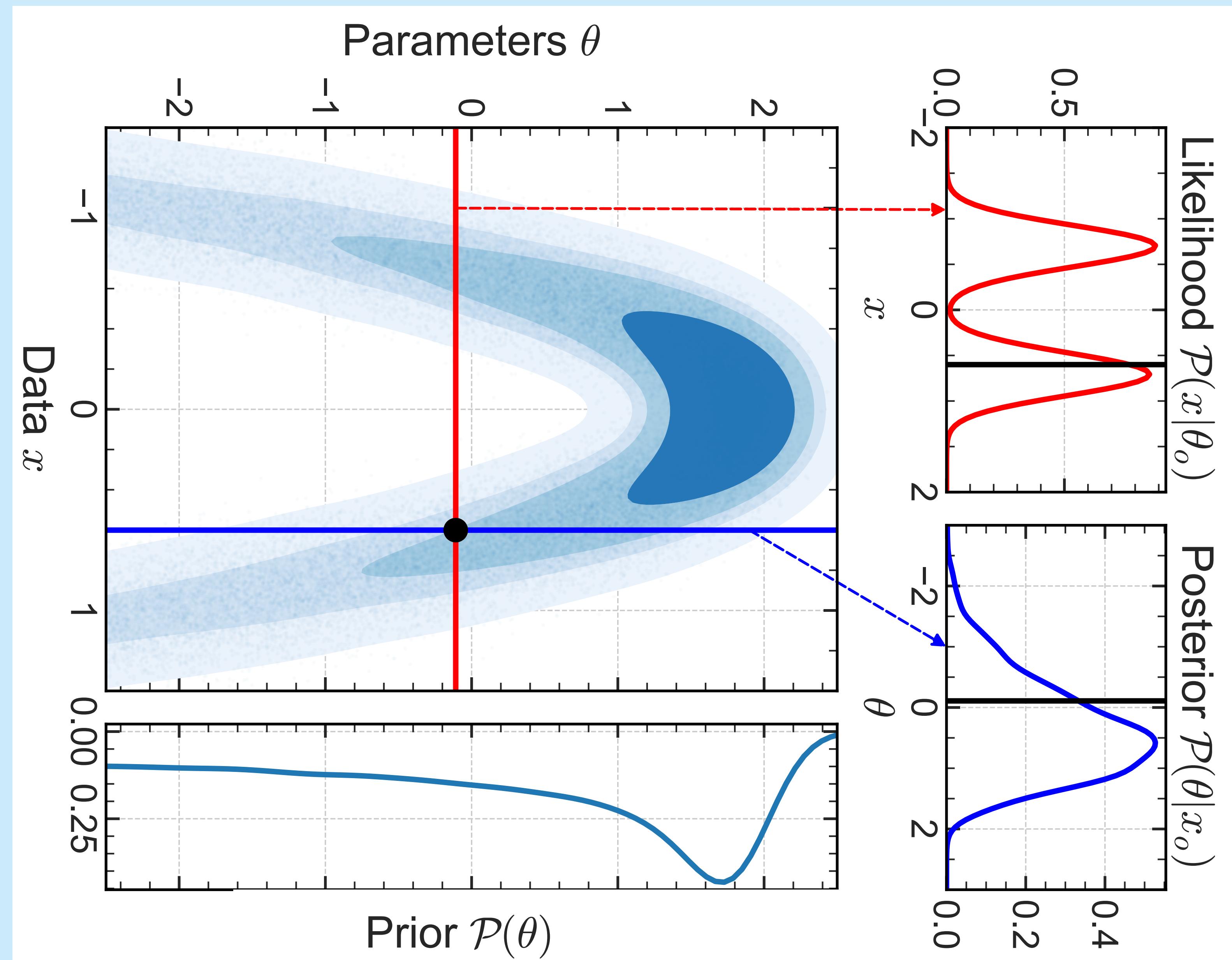
Bayesian inference aims at determining  $p(\theta | \mathbf{x}_0)$

$$p(\theta | \mathbf{x}_0) = \frac{p(\mathbf{x}_0 | \theta)p(\theta)}{p(\mathbf{x}_0)} \propto p(\mathbf{x}_0 | \theta)p(\theta)$$

In astrophysics,  $\mathbf{x}_0$  typically results from a large number of mechanisms/effects that transform the data and involve a large number of latent variables  $z$ , hence the marginal likelihood  $p(\mathbf{x}_0 | \theta)$  is intractable.

$$p(\mathbf{x}_0 | \theta) = \int p(\mathbf{x}_0 | \theta, z)p(z)dz$$

# LIKELIHOODS, PRIORS AND POSTERIORS



# MODELING COMPLEX PROBABILITY DISTRIBUTIONS

How can we approximate high-dimensional, complex probability distributions  $p(\theta | \mathbf{x}_0)$ ?

Goal:

- effectively: learn a model from the data!
- model  $p(\theta | \mathbf{x}_0)$  explicitly or implicitly
- sample and evaluate  $p(\theta | \mathbf{x}_0)$

Options:

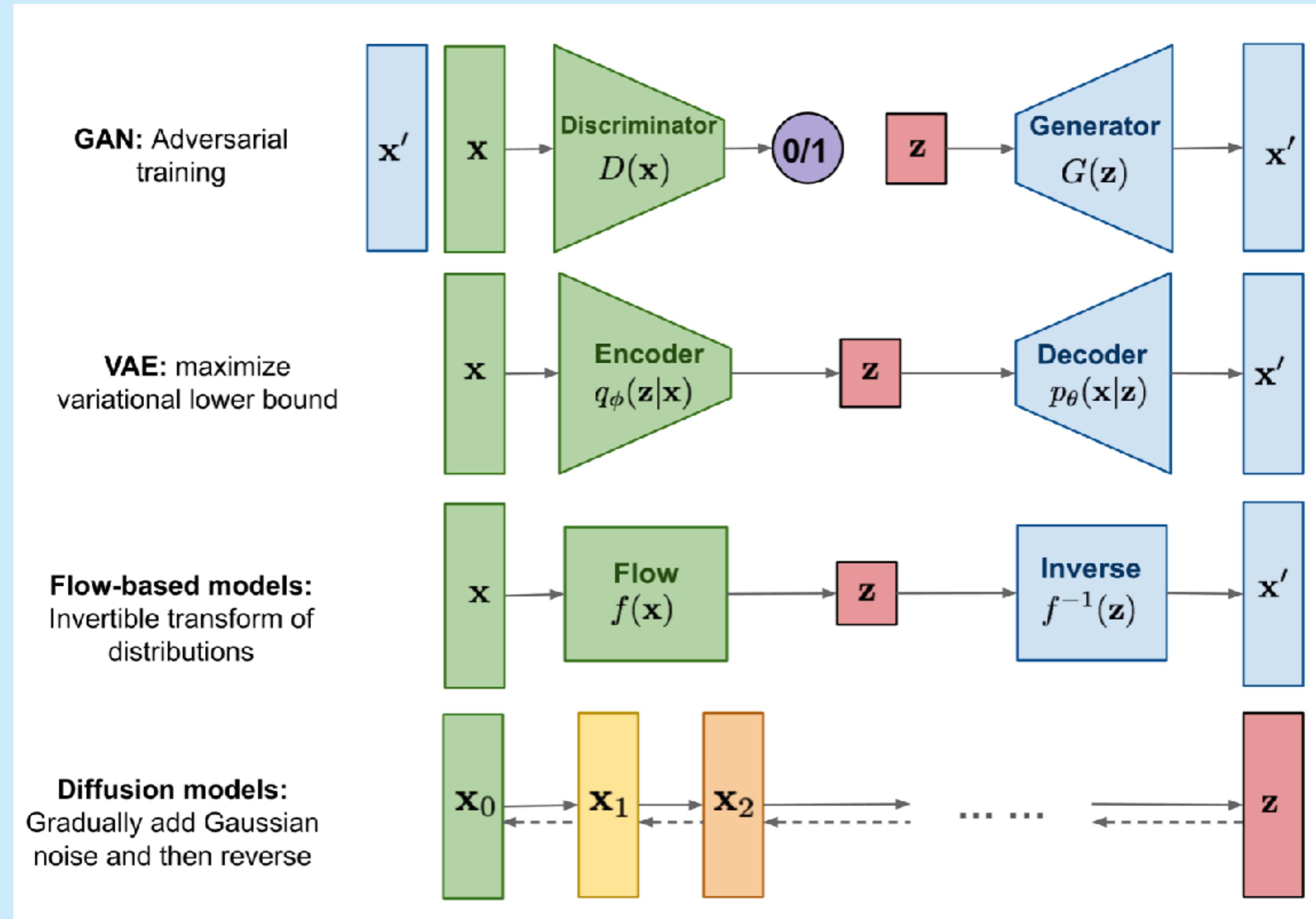
- normalizing flows
- VAEs
- GANs
- score matching / flow matching
- and possibly more

implicit model: architectural constraints

explicit model: prone to mode collapse

# GENERATIVE AI FLAVOURS

"Creating noise from data is easy; creating data from noise is generative modeling." (Song+2020)



# Likelihood-based models

# Normalizing Flows

# NORMALIZING FLOWS

How can we approximate high-dimensional, complex probability distributions  $p(\theta | \mathbf{x}_0)$ ?

Goal:

- model  $p(\theta | \mathbf{x}_0)$  explicitly
- sample and evaluate  $p(\theta | \mathbf{x}_0)$

Advantages of NF:

- exact log-likelihood calculation
- density evaluation possible

Idea:

Transform a simple base distribution through a series of invertible transformations.

# MODELING COMPLEX PROBABILITY DISTRIBUTIONS - NORMALIZING FLOWS

$$T = T_1 \circ T_2$$



# NORMALISING FLOW THEORY

Normalizing Flows: Change of variables

$$1 = \int p_Z(z) dz$$

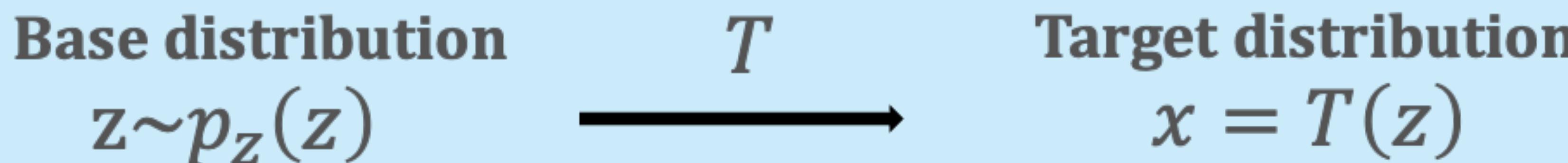
**Transform:**  $x = T(z)$

$$\int p_X(x) dx = 1 = \int p_Z(z) dz = \int p_Z(T^{-1}(x)) |\det J_{T^{-1}}(x)| dx$$

$$p_X(x) = p_Z(T^{-1}(x)) |\det J_{T^{-1}}(x)|$$

# NORMALISING FLOW THEORY

## Normalizing Flows: Definition



### Defining Properties:

- $T$  must be invertible
- $T, T^{-1}$  differentiable

### Diffeomorphism

$$T = T_N \circ T_{N-1} \circ \dots \circ T_1$$

$$T^{-1} = T_1^{-1} \circ T_2^{-1} \circ \dots \circ T_N^{-1}$$

$$\det J_{T_2 \circ T_1}(z) = \det J_{T_2}(T_1(z)) \cdot \det J_{T_1}(z)$$

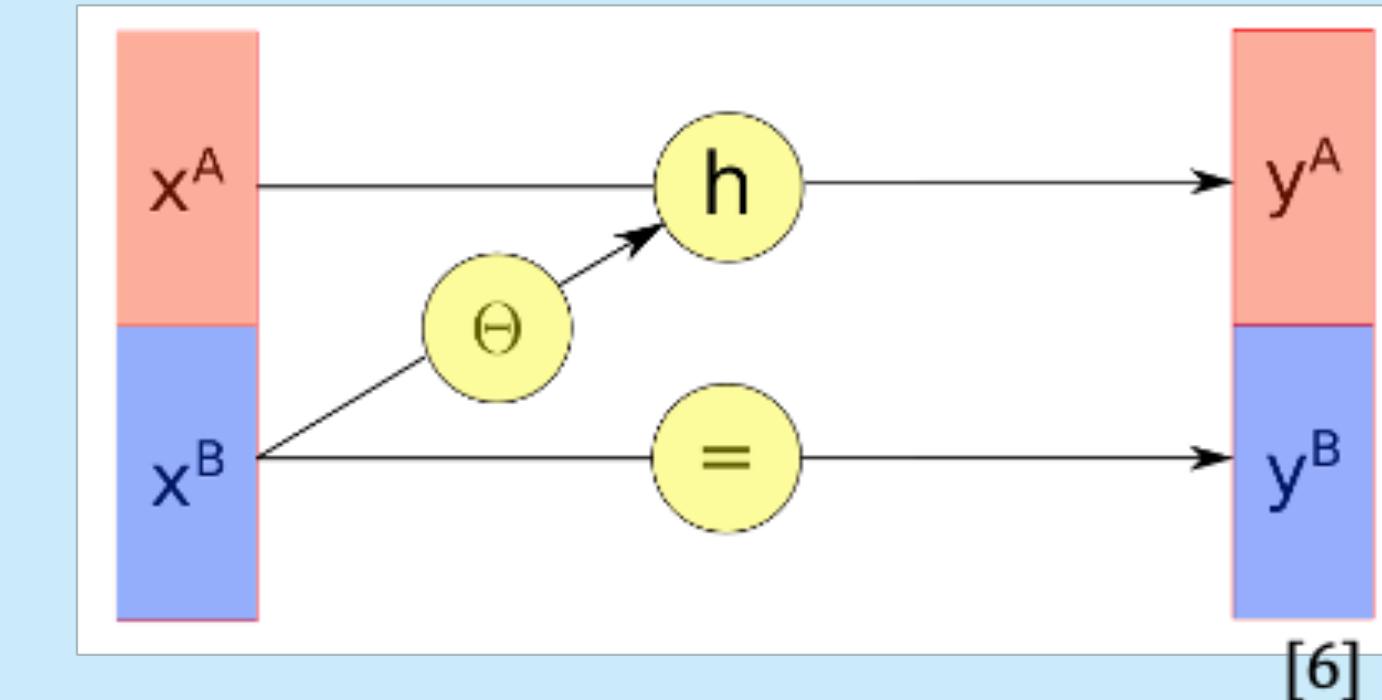
# NORMALISING FLOW: IMPLEMENTING TRANSFORMATIONS

## Coupling Flows

- Two partitions:

$$y^A = h(x^A, \theta(x^B))$$

$$y^B = x^B$$



- Conditioner  $\theta(x^B)$  does not have to be invertible  
→ use Neural Networks for good expressivity
- Plethora of different coupling functions

# NORMALISING FLOW: IMPLEMENTING TRANSFORMATIONS

## Coupling Functions (How to choose $h$ ?)

- Additive  $h(x; \theta) = x + \theta$
- Affine  $h(x; \theta) = \theta_1 x + \theta_2$  (NICE, RealNVP, Glow, MAF, IAF)

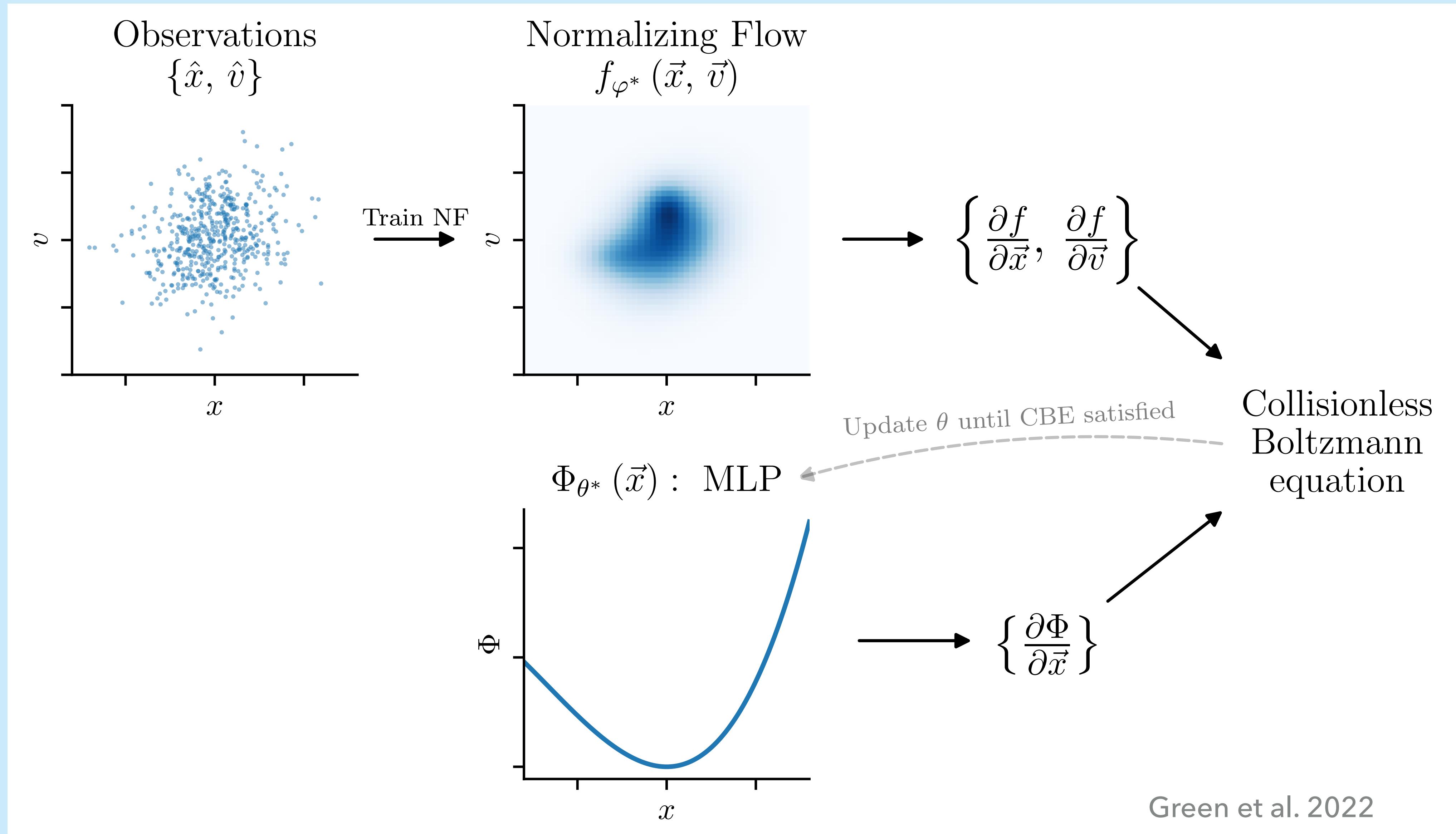
$$y_{1:d-1} = x_{1:d-1}$$

$$y_{d:n} = \exp(s_{d:n}(x_{1:d-1})) \cdot y_{d:n} + t_{d:n}(x_{1:d-1})$$

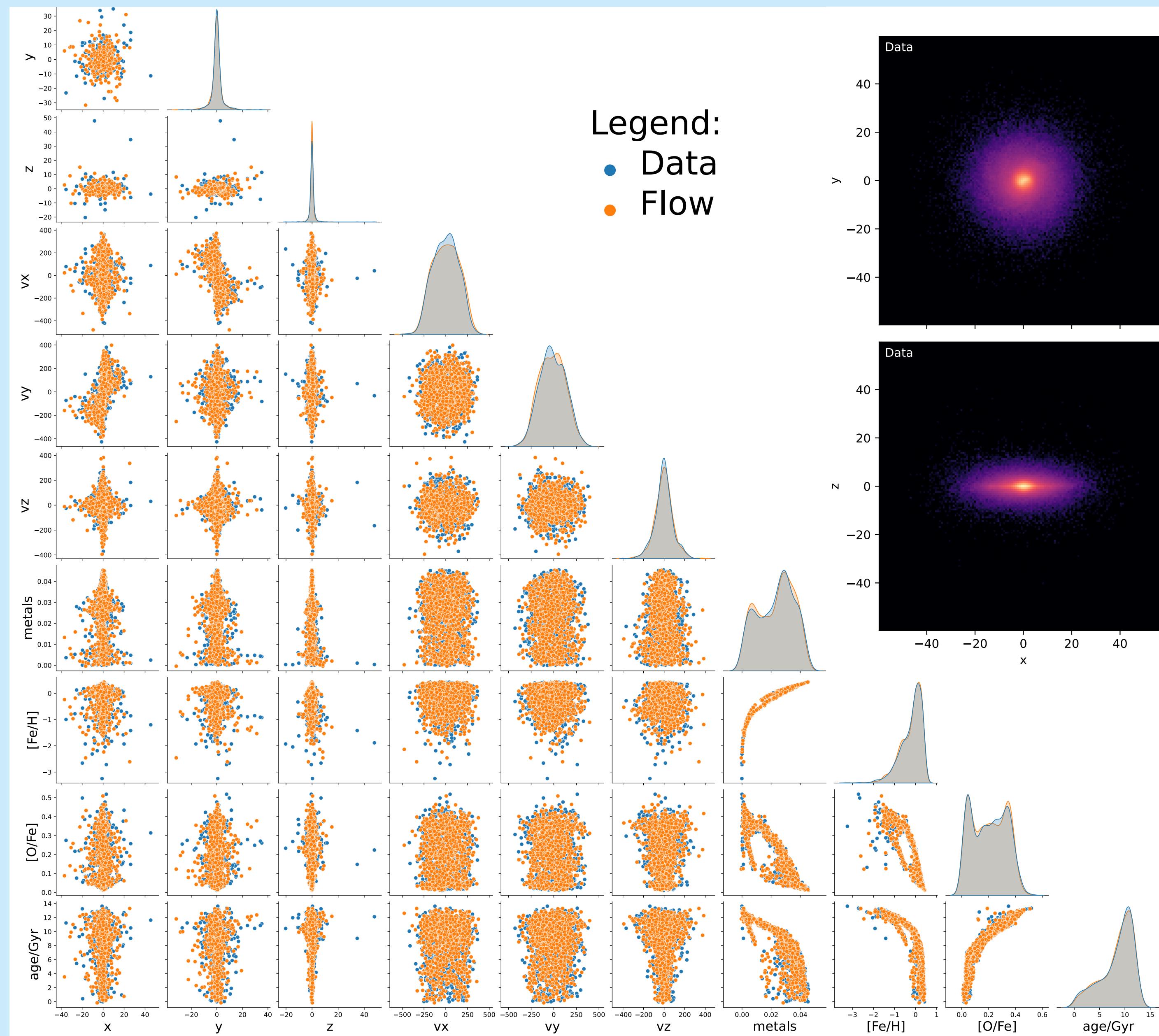
- Continuous mixture CDFs  $h(x; \theta) = \theta_1 F(x, \theta_3) + \theta_2$
- Splines (RQ-NSF)
- Neural Networks

Important: after each coupling layer permute dimensions (or use 1x1 convolution like GLOW)

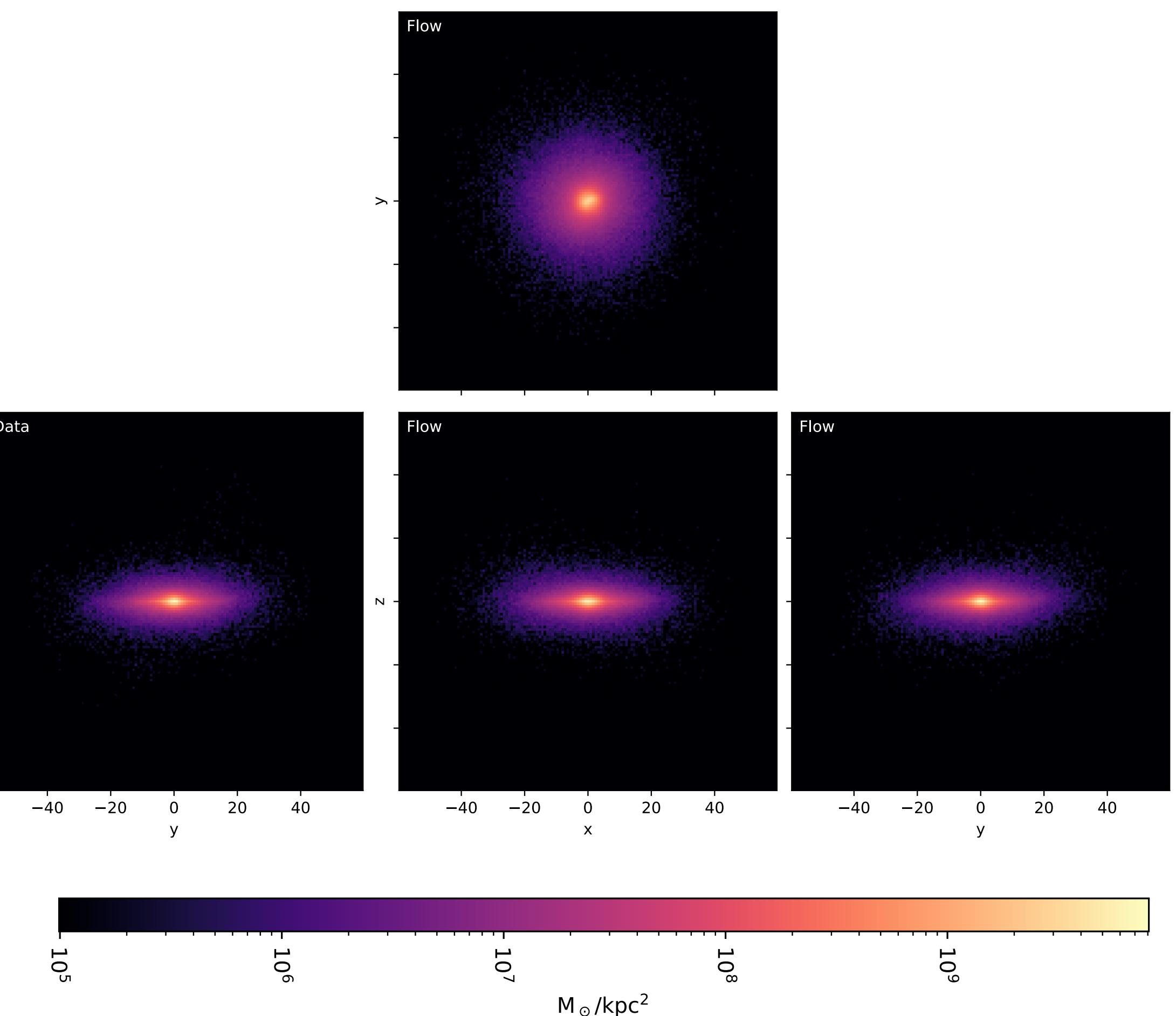
# NORMALISING FLOW: APPLICATION I



# NORMALISING FLOW: APPLICATION II



2D mass density cornerplot. Left: data, right: sample



Wolf+Buck 2023

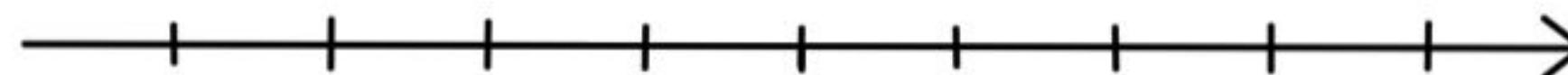
# NORMALISING FLOW: APPLICATION III

- Normalizing flows for random fields in cosmology (Rouhiainen+2021)
- Bayesian Stokes inversion with normalizing flows (Baso+2022)
- A Hierarchy of Normalizing Flows for Modelling the Galaxy-Halo Relationship (Lovell+2023)
- HiFlow: Generating Diverse Hi Maps and Inferring Cosmology while Marginalizing over Astrophysics Using Normalizing Flows (Hassan+2022)
- Normalizing Flows as an Avenue to Studying Overlapping Gravitational Wave Signals (Langendorff+2023)
- Charting Galactic Accelerations: When and How to Extract a Unique Potential from the Distribution Function (An+2021)
- Charting galactic accelerations II: how to ‘learn’ accelerations in the solar neighbourhood (Naik+2021)
- many many more...

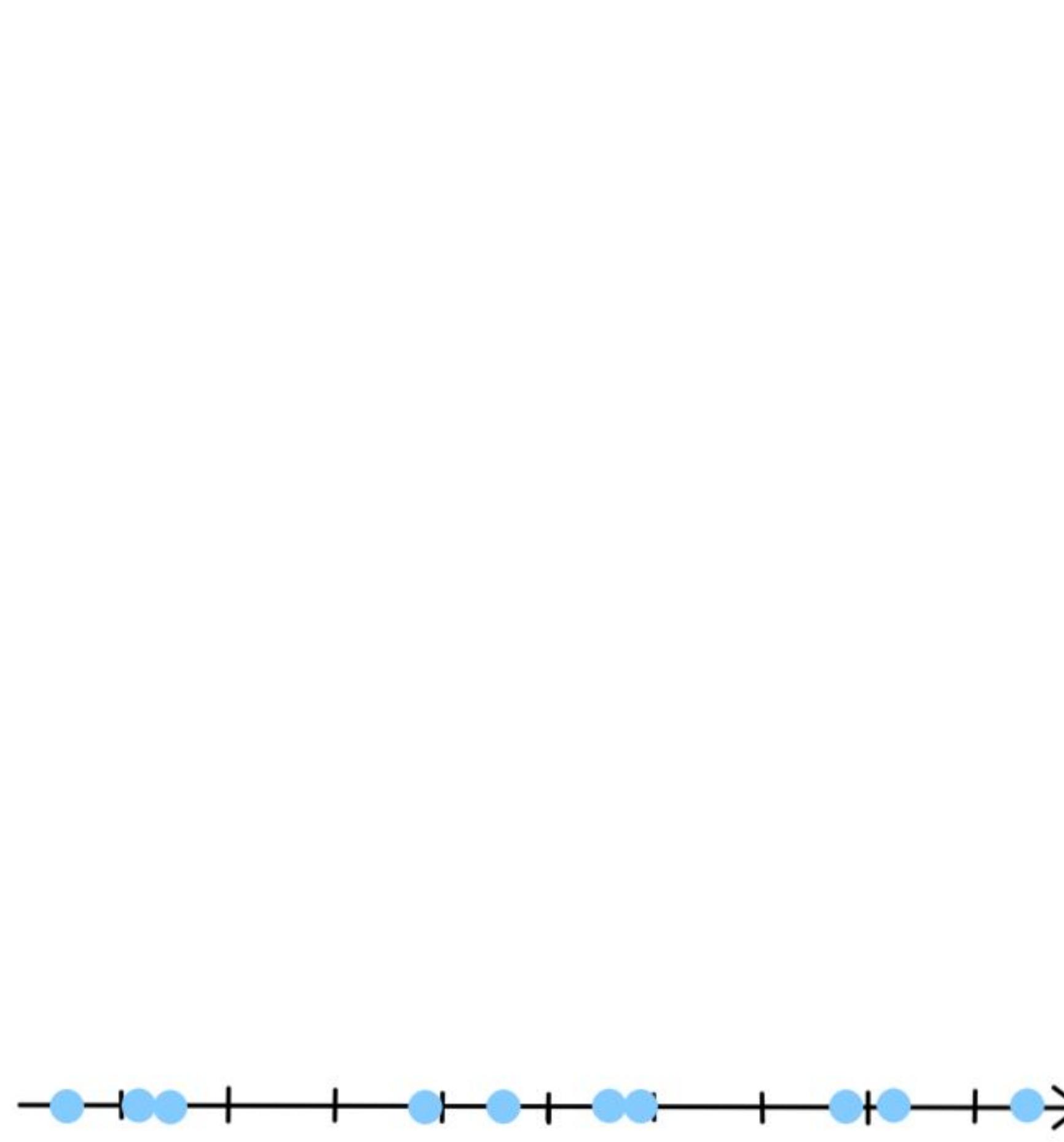
# Variational Autoencoder

# THE CURSE OF DIMENSIONALITY

Divide range into  $K$  intervals  $\Delta x_1$



# THE CURSE OF DIMENSIONALITY



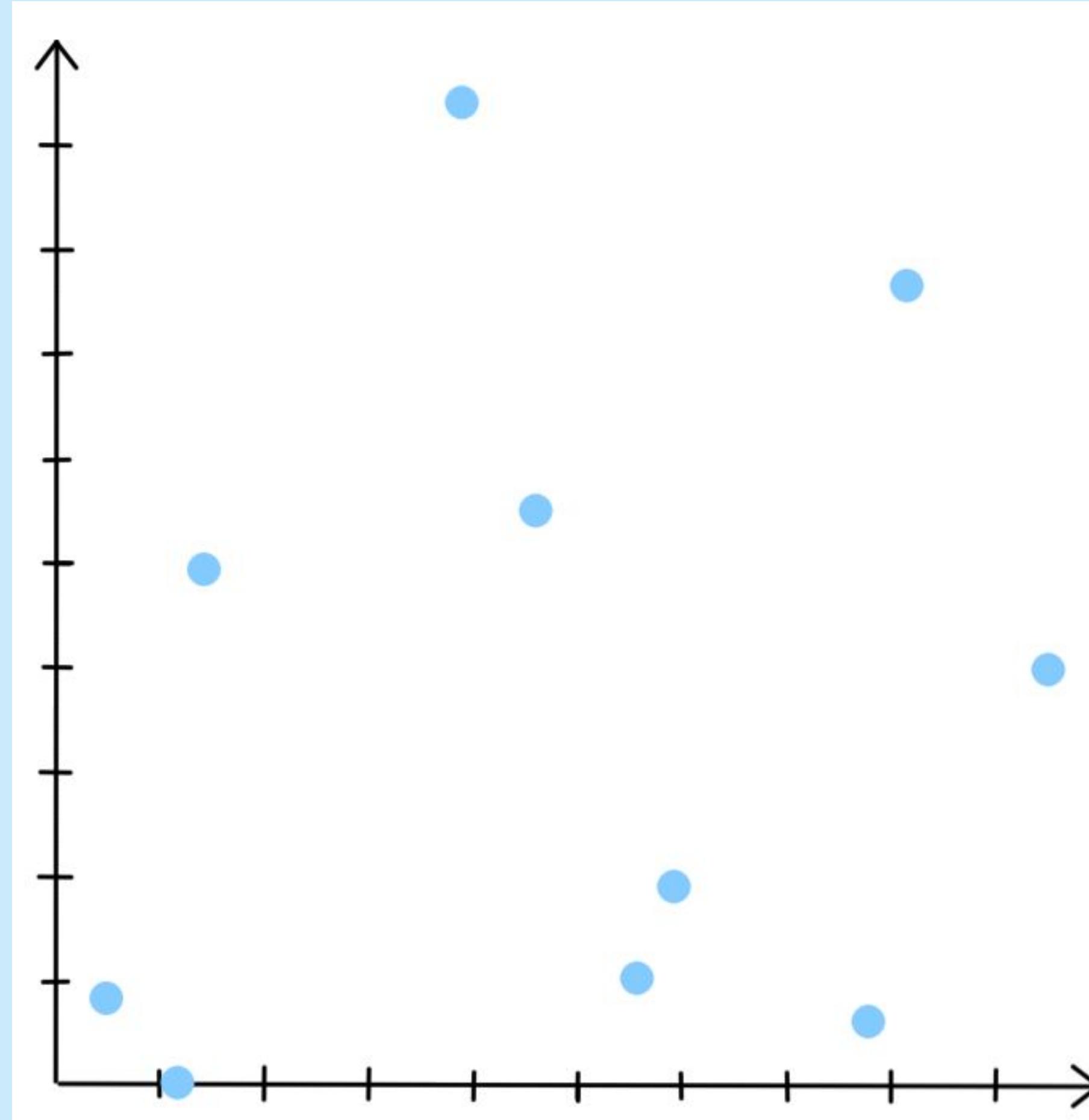
Divide range into  $K$  intervals  $\Delta x_1$

Take data set of size  $N$

On average:  $n = N/K$  points per interval

→ Here:  $n = 10/10 = 1$

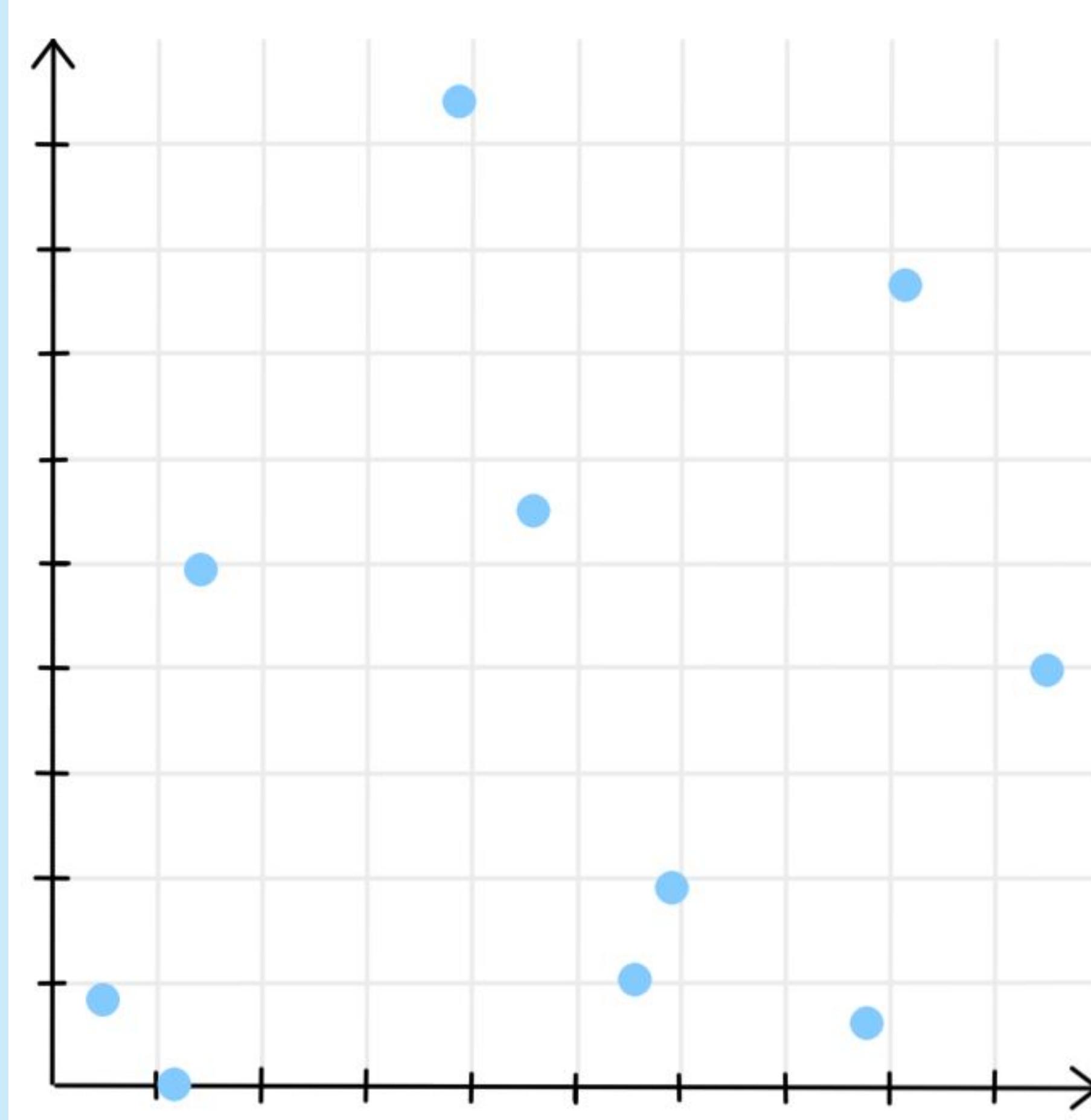
# THE CURSE OF DIMENSIONALITY



Divide range into  $K$  intervals  $\Delta x_1$   
Take data set of size  $N$   
On average:  $n = N/K$  points per interval  
→ Here:  $n = 10/10 = 1$

For 2D:  $K^2$  intervals of size  $\Delta x_1 \Delta x_2$

# THE CURSE OF DIMENSIONALITY



Divide range into  $K$  intervals  $\Delta x_1$

Take data set of size  $N$

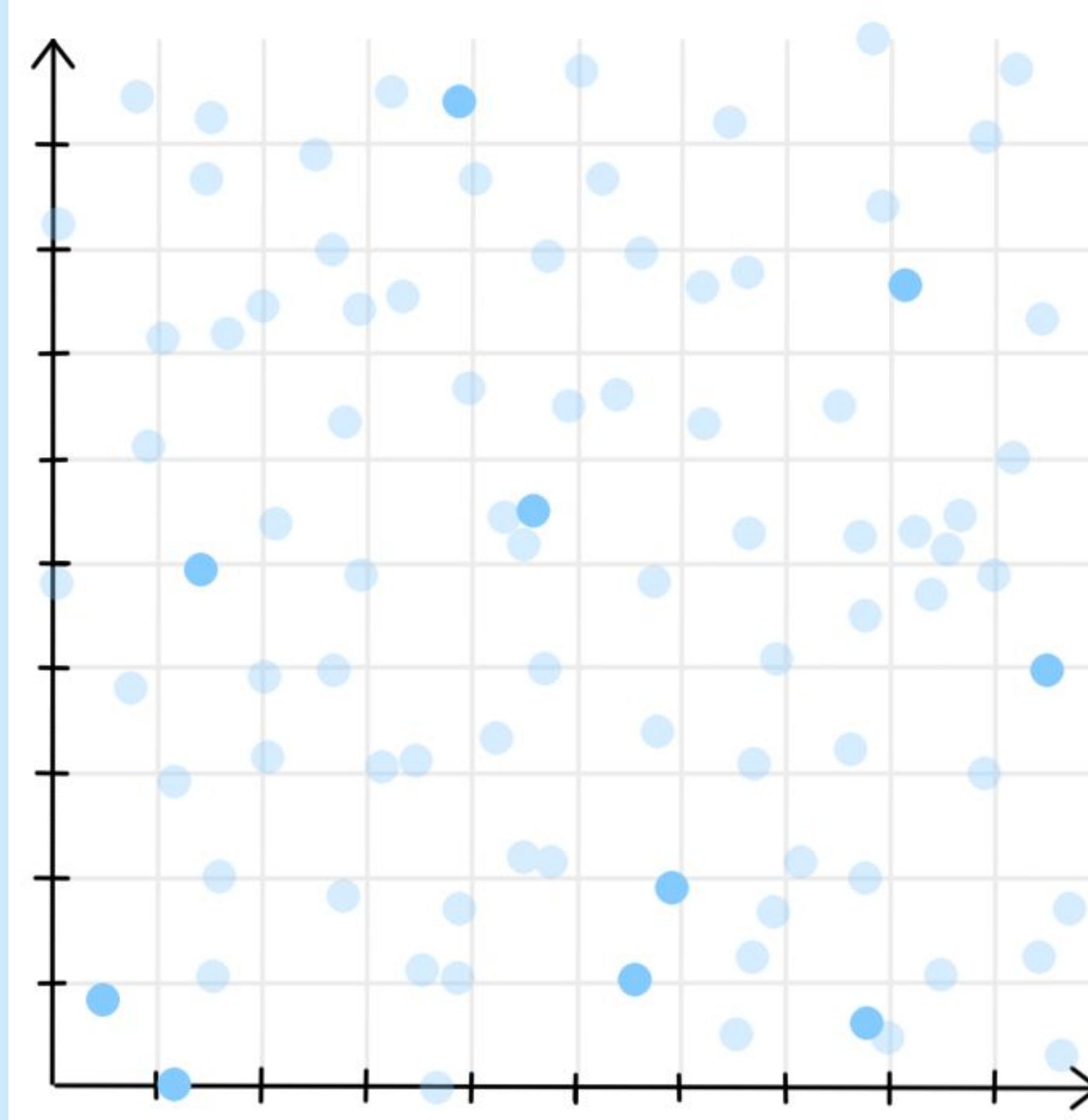
On average:  $n = N/K$  points per interval

→ Here:  $n = 10/10 = 1$

For 2D:  $K^2$  intervals of size  $\Delta x_1 \Delta x_2$

On average:  $n = 10/100 = 0.1$

# THE CURSE OF DIMENSIONALITY



Divide range into  $K$  intervals  $\Delta x_1$

Take data set of size  $N$

On average:  $n = N/K$  points per interval

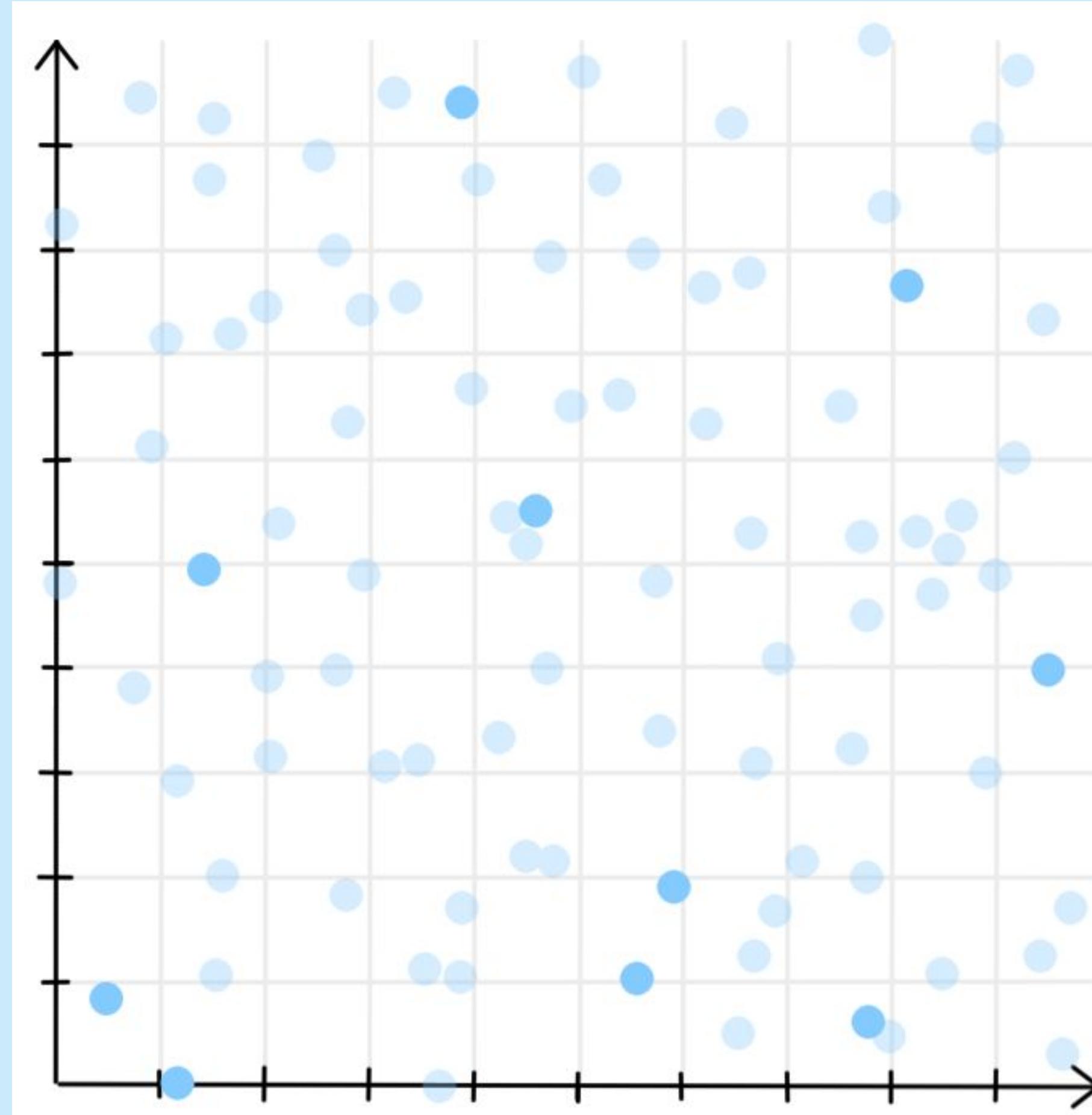
→ Here:  $n = 10/10 = 1$

For 2D:  $K^2$  intervals of size  $\Delta x_1 \Delta x_2$

On average:  $n = 10/100 = 0.1$

For same  $n$ ,  $nK^2 = N^2 = 100$  data points needed

# THE CURSE OF DIMENSIONALITY



Divide range into  $K$  intervals  $\Delta x_1$

Take data set of size  $N$

On average:  $n = N/K$  points per interval

→ Here:  $n = 10/10 = 1$

For 2D:  $K^2$  intervals of size  $\Delta x_1 \Delta x_2$

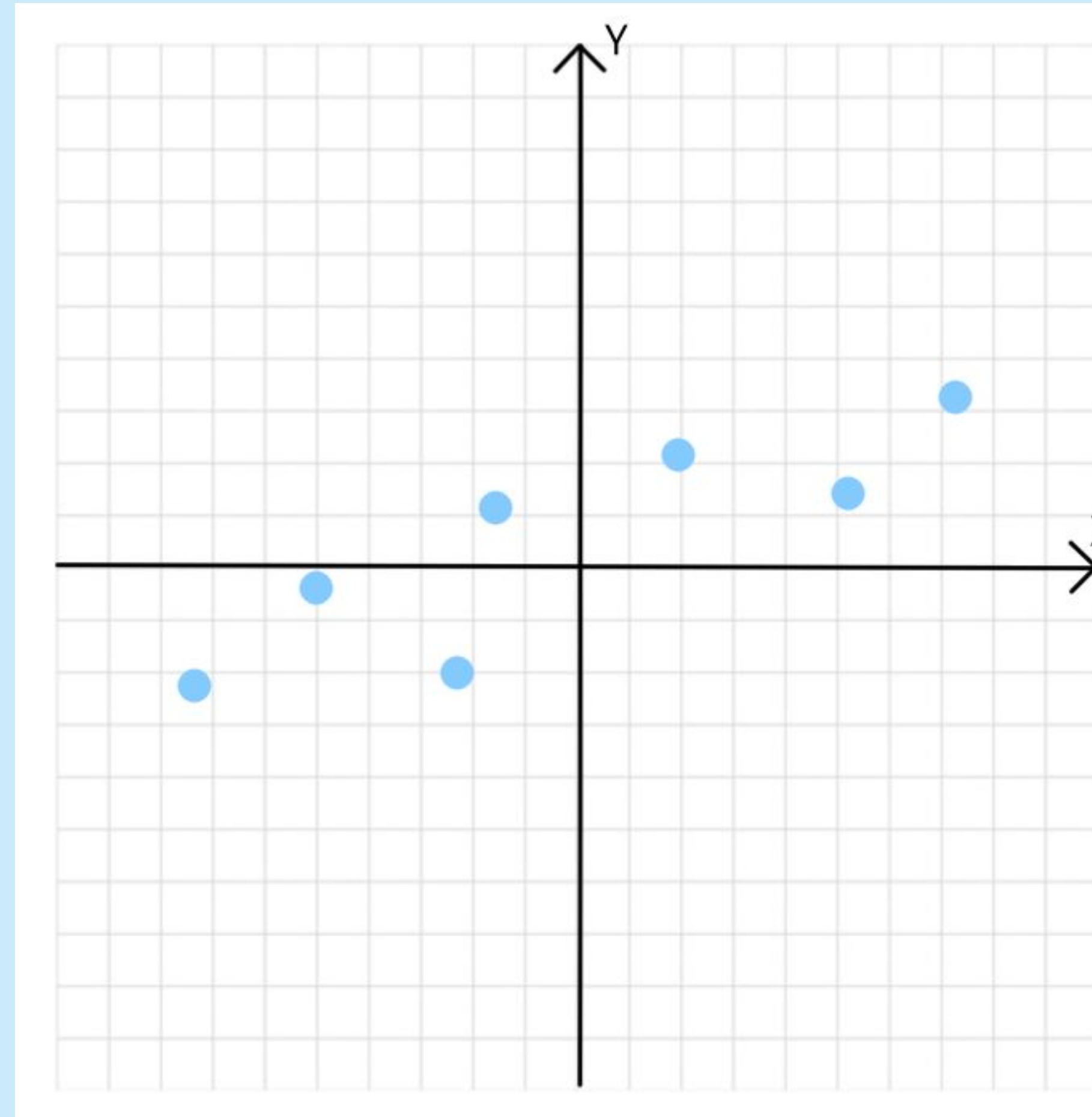
On average:  $n = 10/100 = 0.1$

For same  $n$ ,  $nK^2 = N^2 = 100$  data points needed

For  $p$  dimensions:  $N^p$  data points required

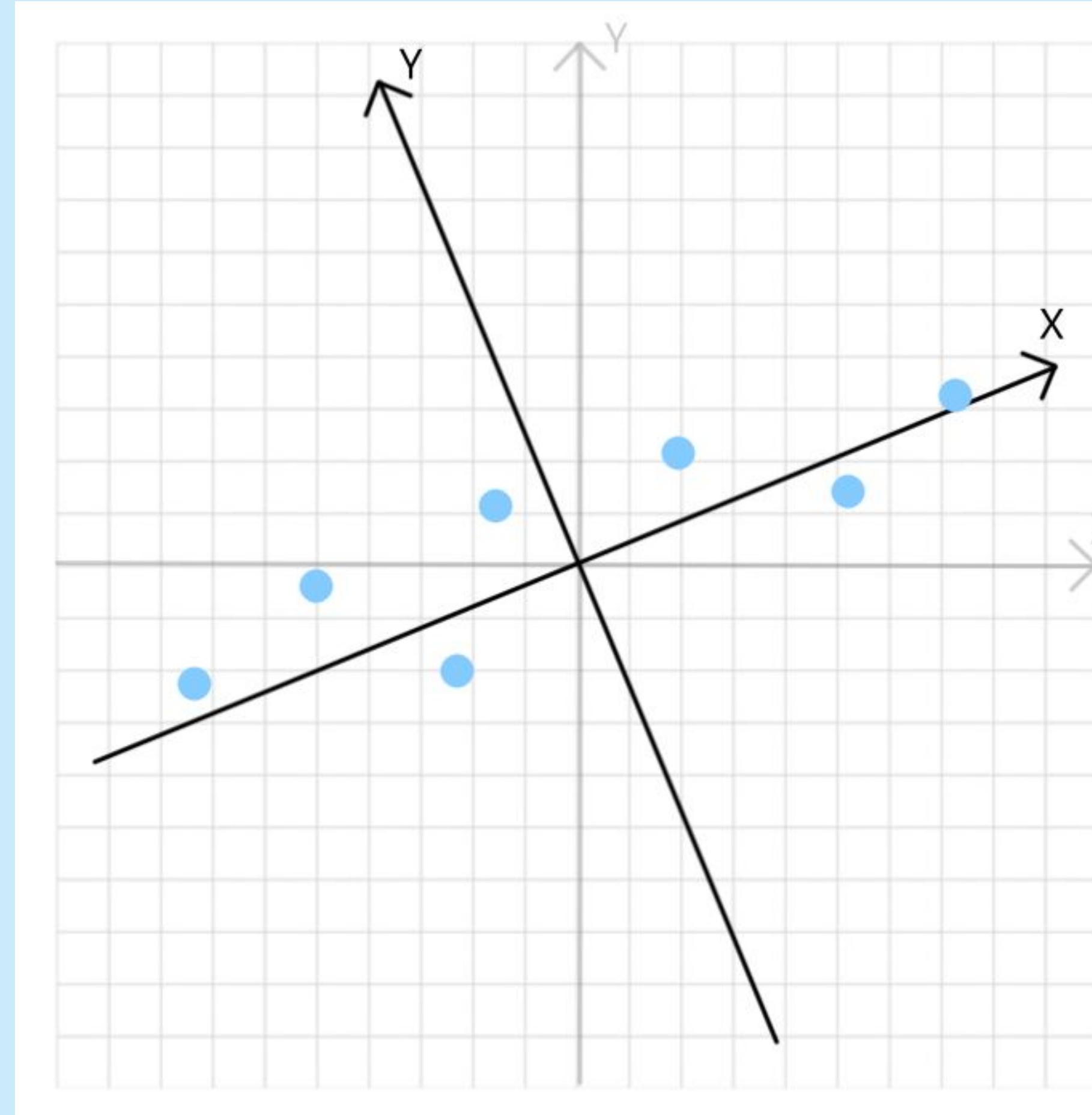
→  $p = 4$ : 10.000 points

# PRINCIPAL COMPONENT ANALYSIS



Maximise the variance along a new axis

# PRINCIPAL COMPONENT ANALYSIS

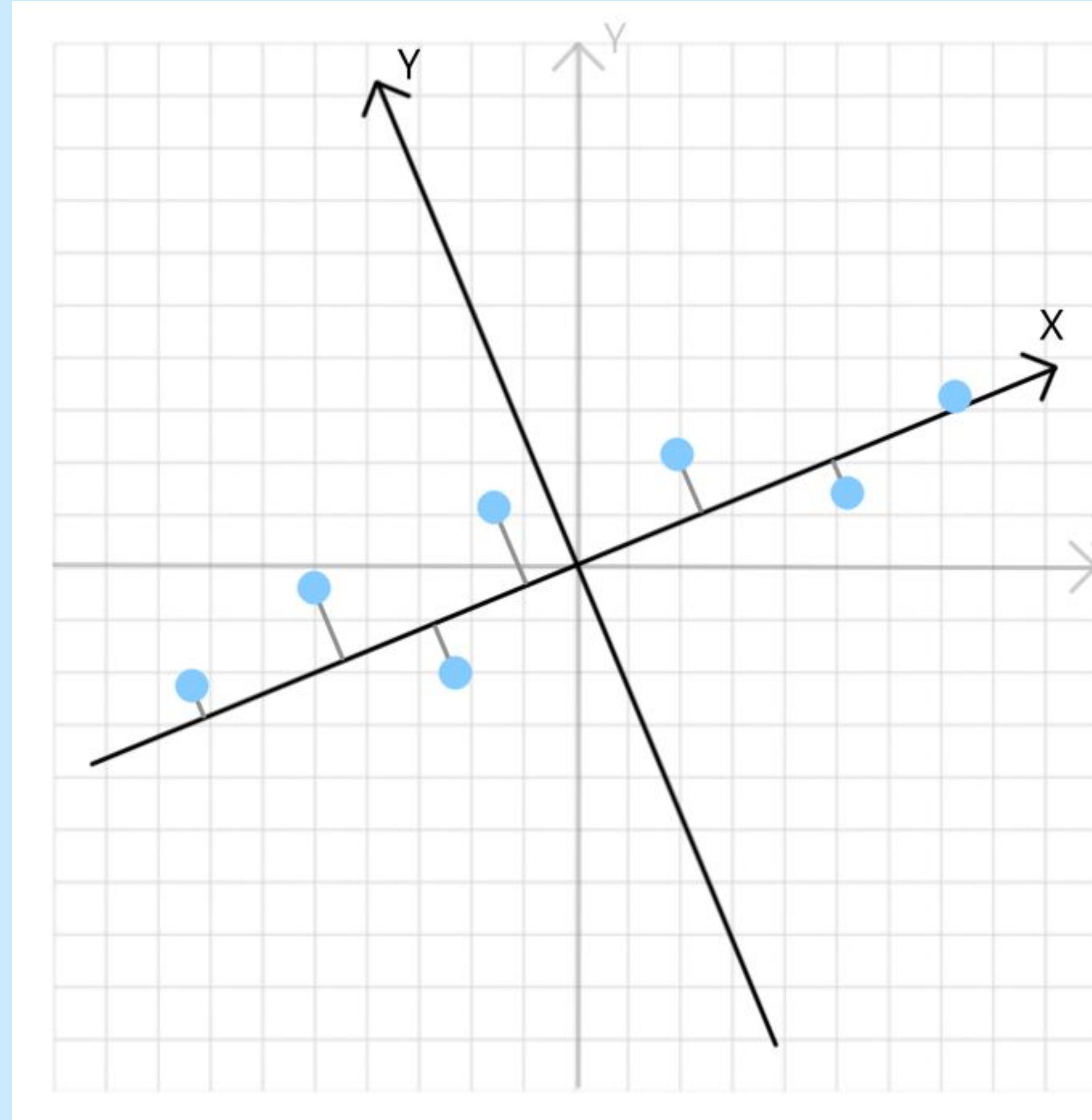


Maximise the variance along a new axis

Variance is maximised along  $x$ , and minimised along  $y$

Neglect  $y$  to reduce dimensionality

# PRINCIPAL COMPONENT ANALYSIS

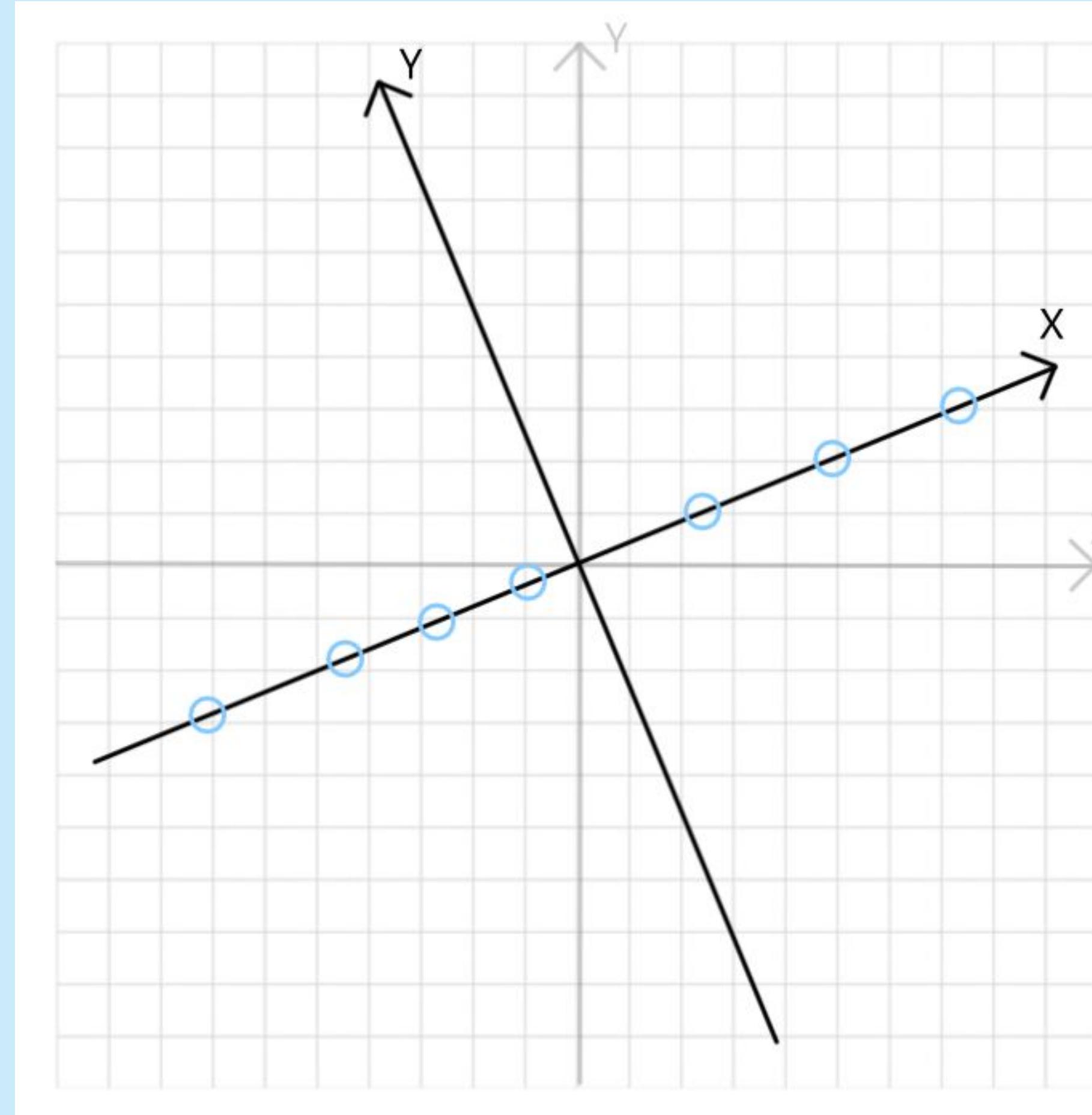


Maximise the variance along a new axis

Variance is maximised along x, and minimised  
along y

Neglect y to reduce dimensionality  
→ Project data on x

# PRINCIPAL COMPONENT ANALYSIS



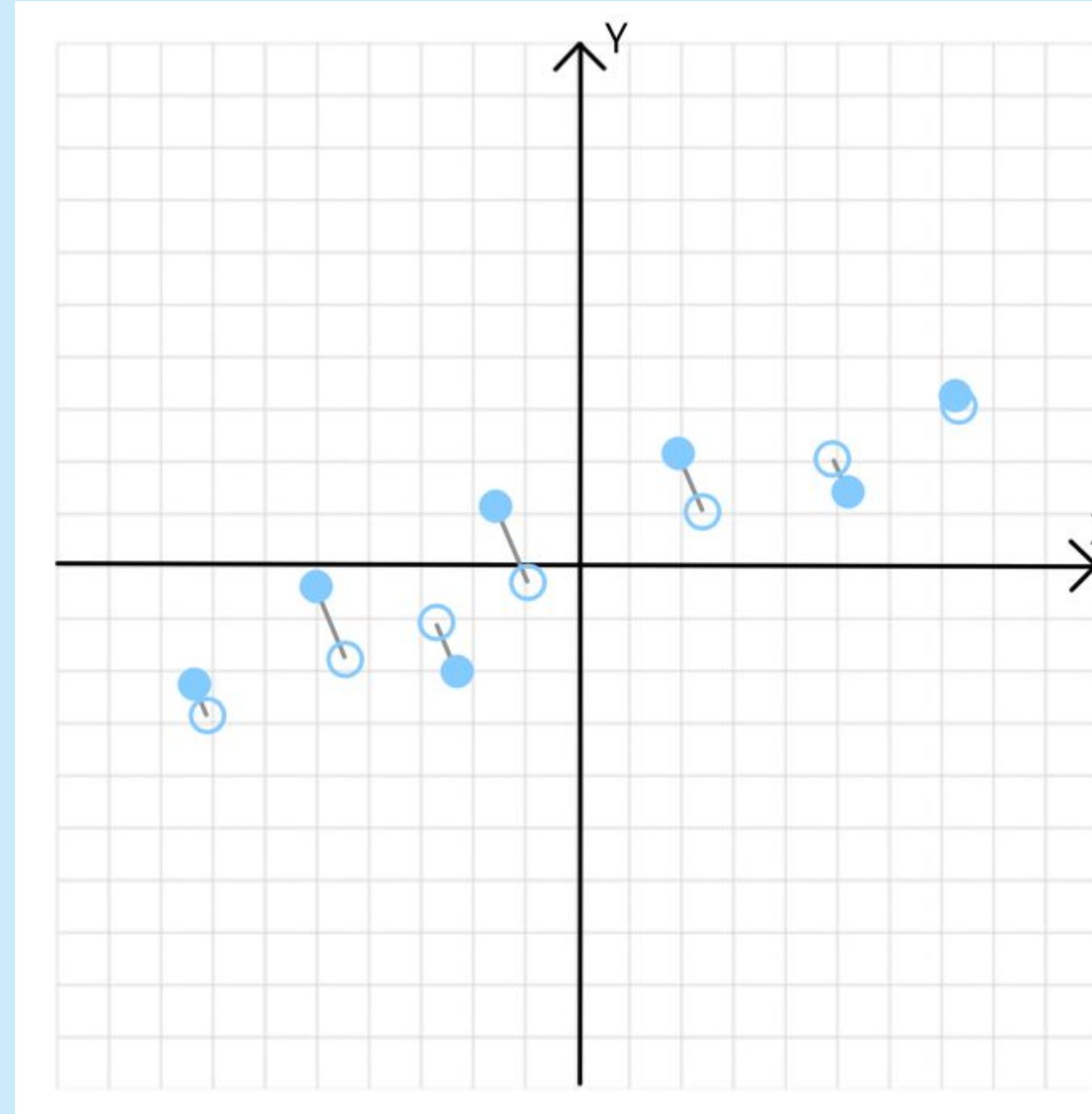
Maximise the variance along a new axis

Variance is maximised along x, and minimised along y

Neglect y to reduce dimensionality  
→ Project data on x

That's the new dimensionally reduced data set

# PRINCIPAL COMPONENT ANALYSIS



Maximise the variance along a new axis

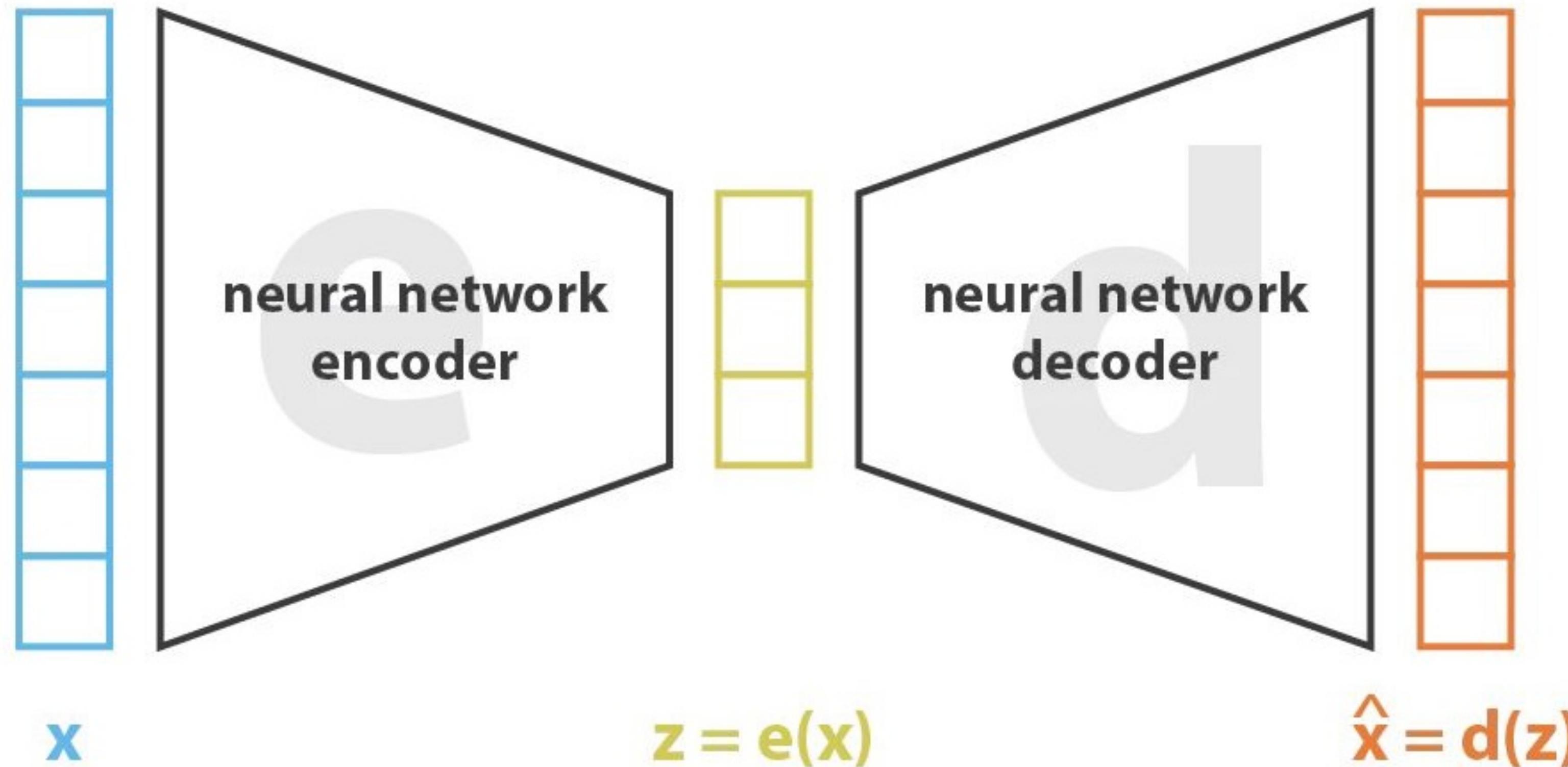
Variance is maximised along x, and minimised along y

Neglect y to reduce dimensionality  
→ Project data on x

That's the new dimensionally reduced data set

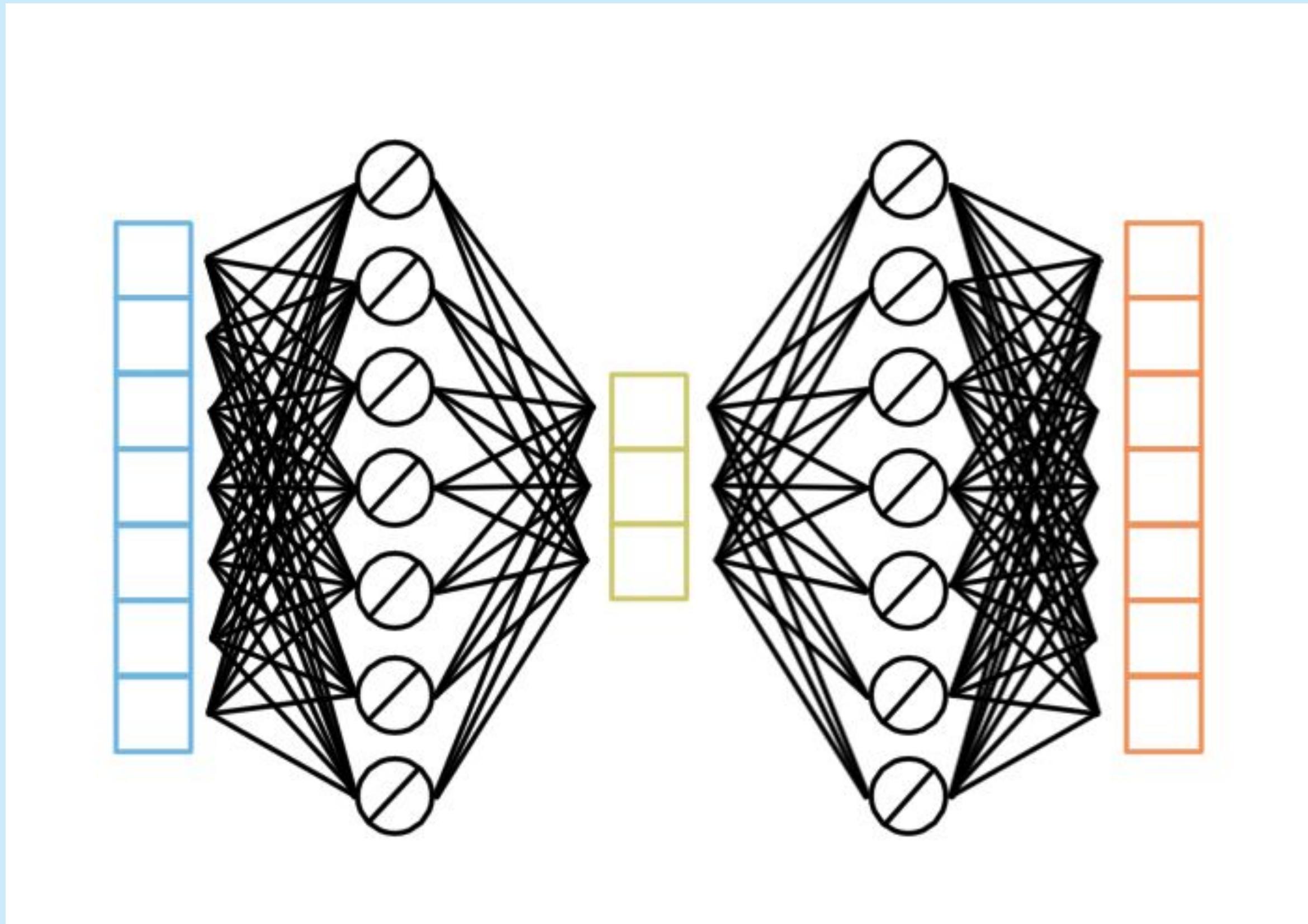
Information (Reconstruction) loss is minimised

# AUTOENCODER

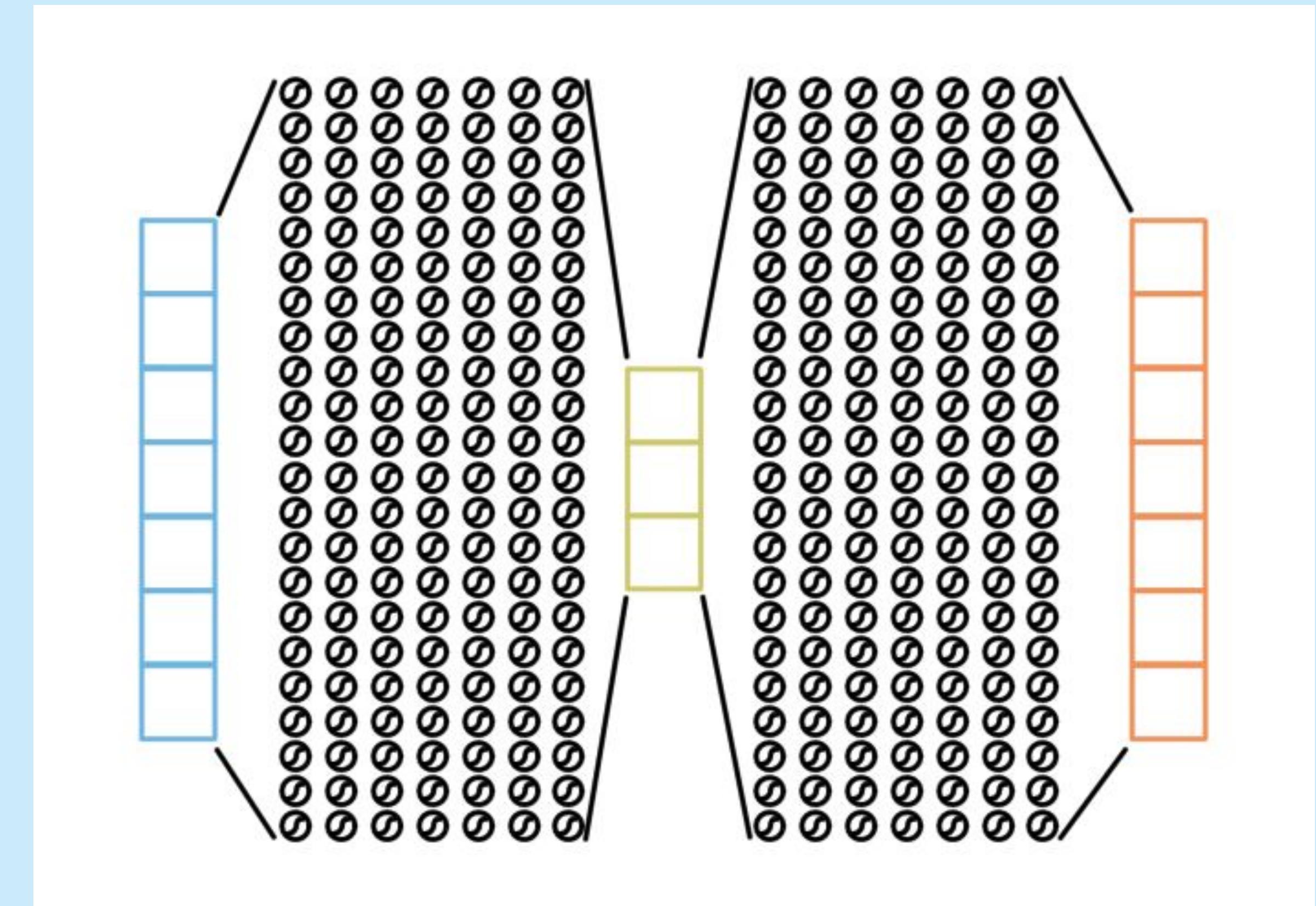


$$\text{loss} = \| \mathbf{x} - \hat{\mathbf{x}} \|^2 = \| \mathbf{x} - \mathbf{d}(z) \|^2 = \| \mathbf{x} - \mathbf{d}(e(\mathbf{x})) \|^2$$

# AUTOENCODER

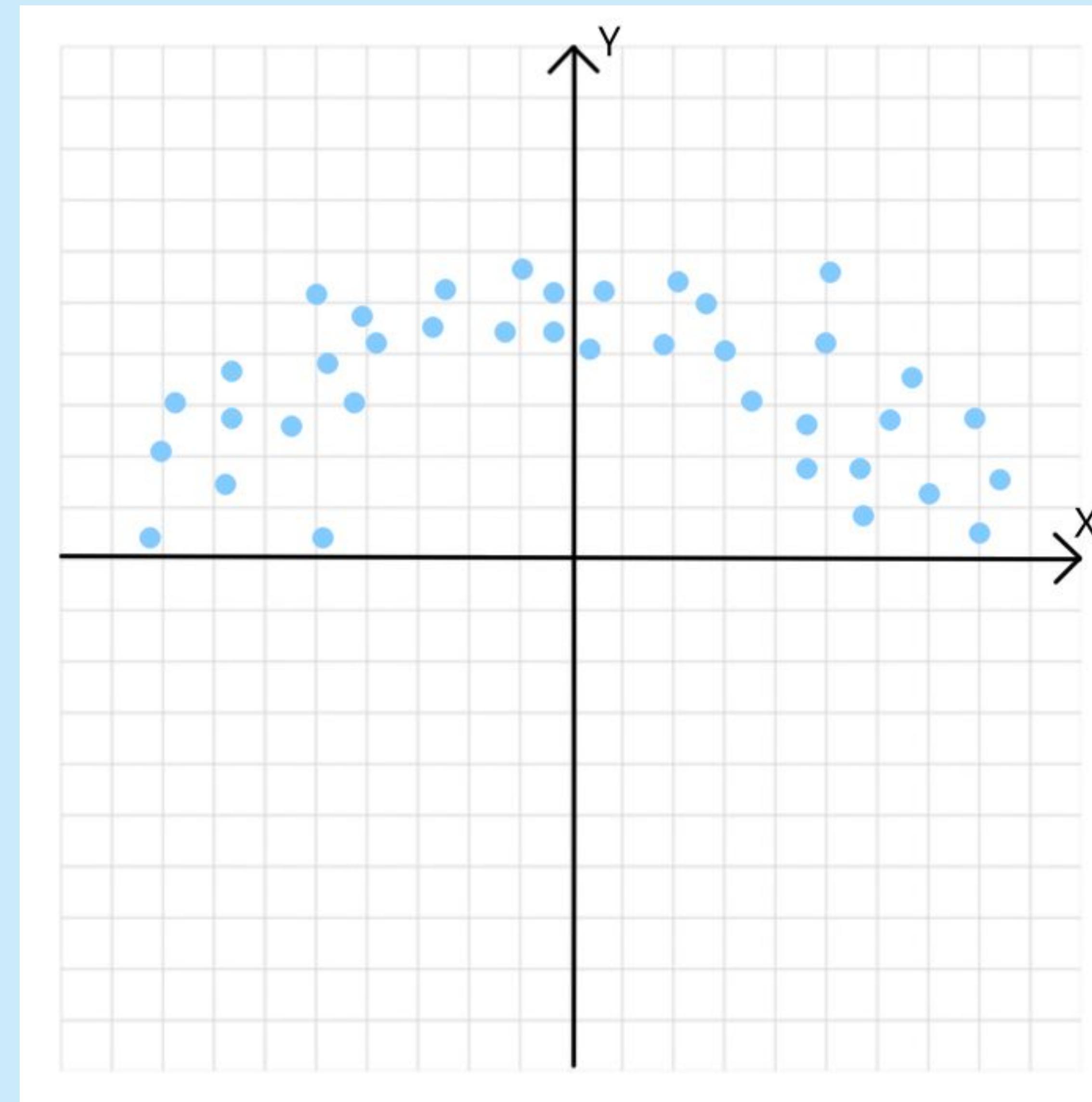


Linear transformation of inputs; very similar to PCA

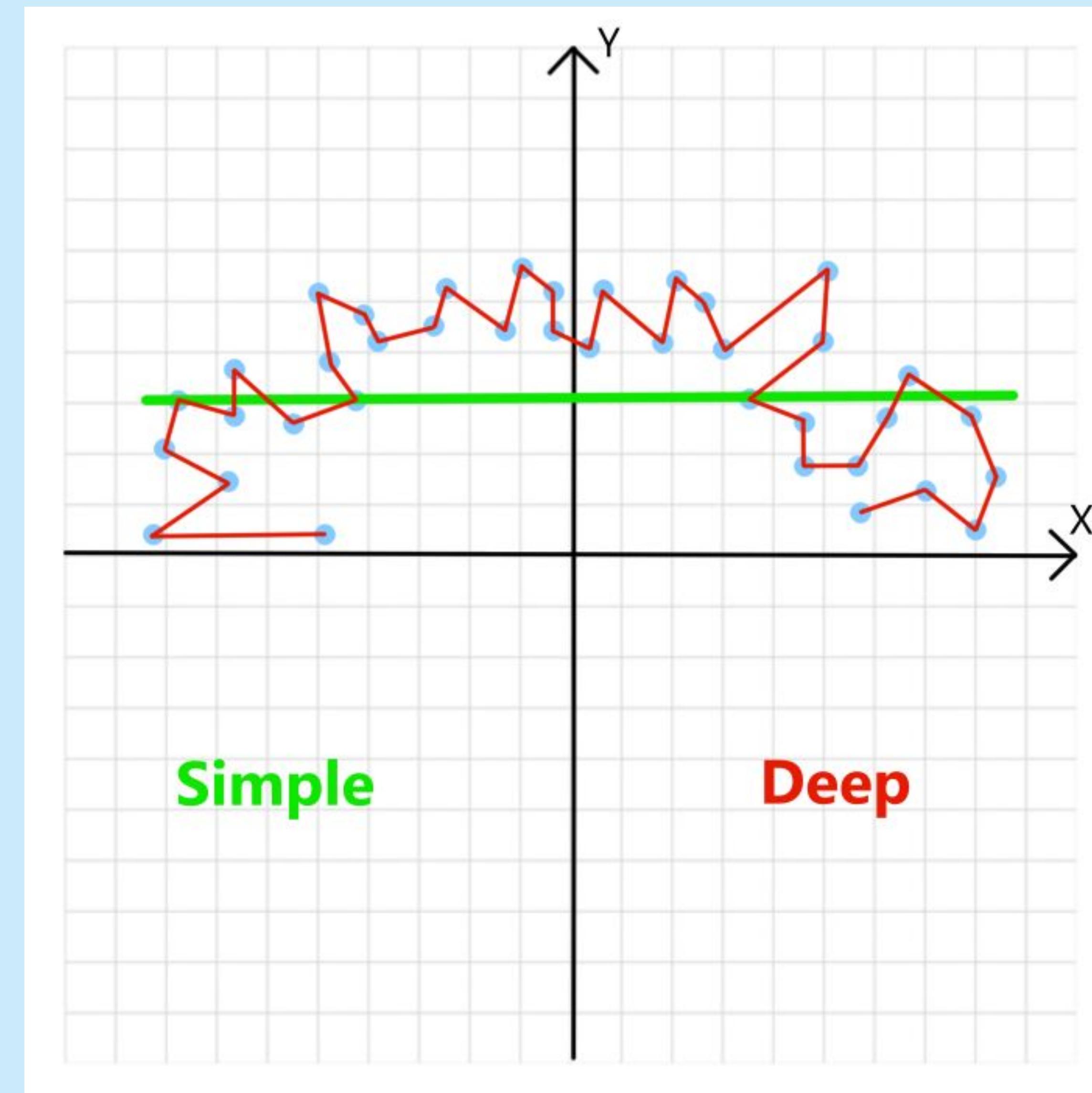


Infinitely deep network has the capacity for lossless compression into a single dimension

# AUTOENCODER: LOSSLESS COMPRESSION



# AUTOENCODER: LOSSLESS COMPRESSION

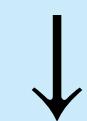


# AUTOENCODER: LOSSLESS COMPRESSION

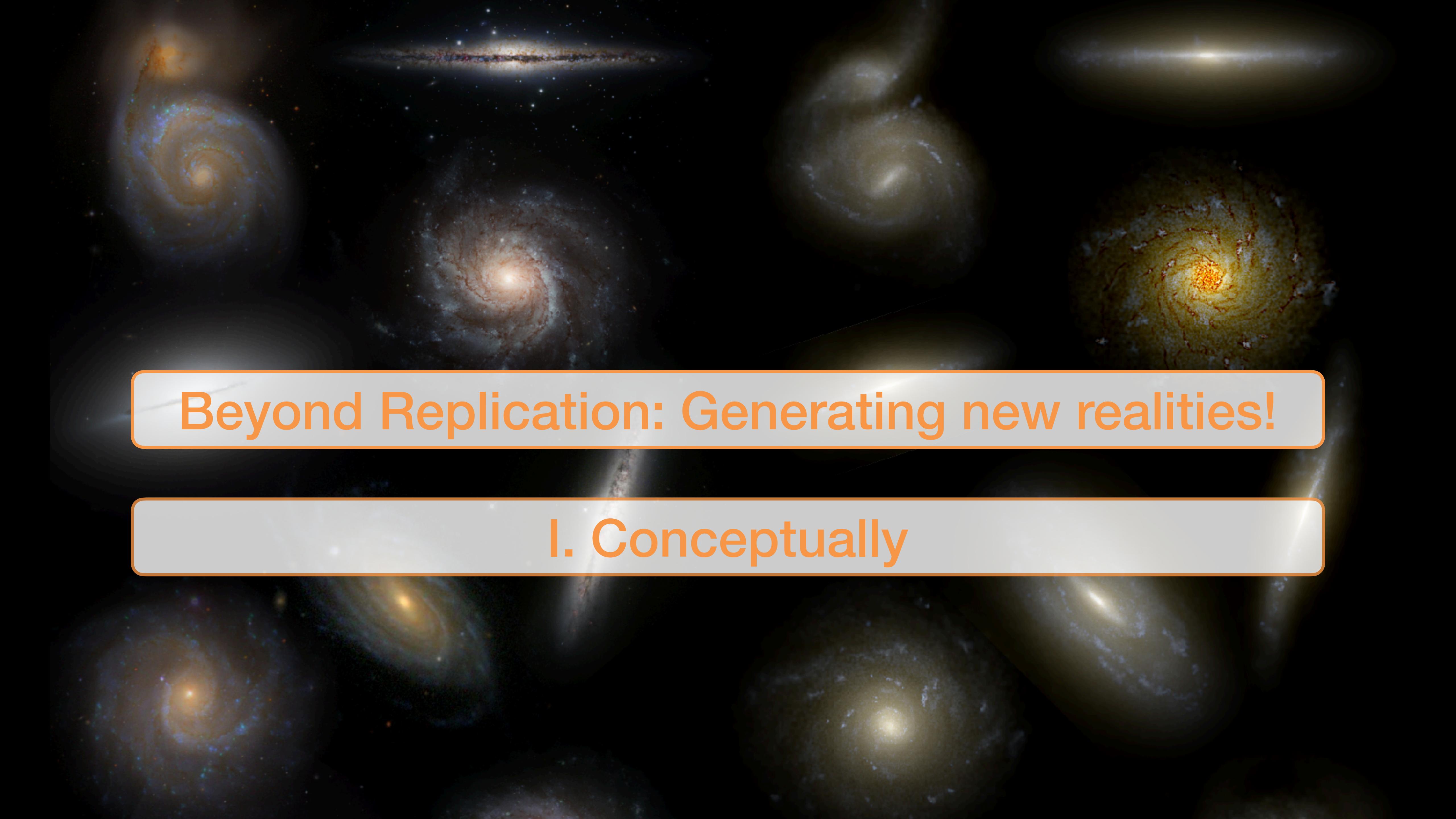
Is lossless compression to a single dimension always what we want?



Strong compression suffers from a lack of interpretability and generalisation



Reduce dimensionality, while also keeping the the data structure interpretable

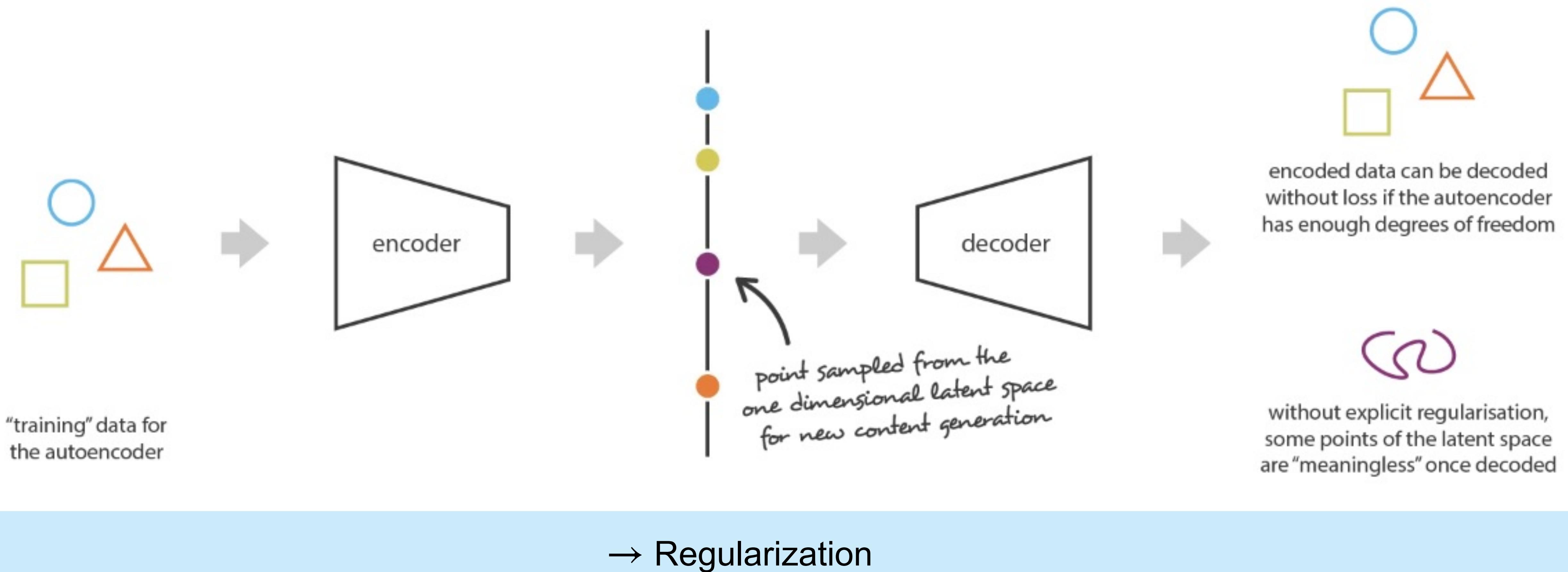


Beyond Replication: Generating new realities!

I. Conceptually

# GENERATION WITH AUTOENCODERS

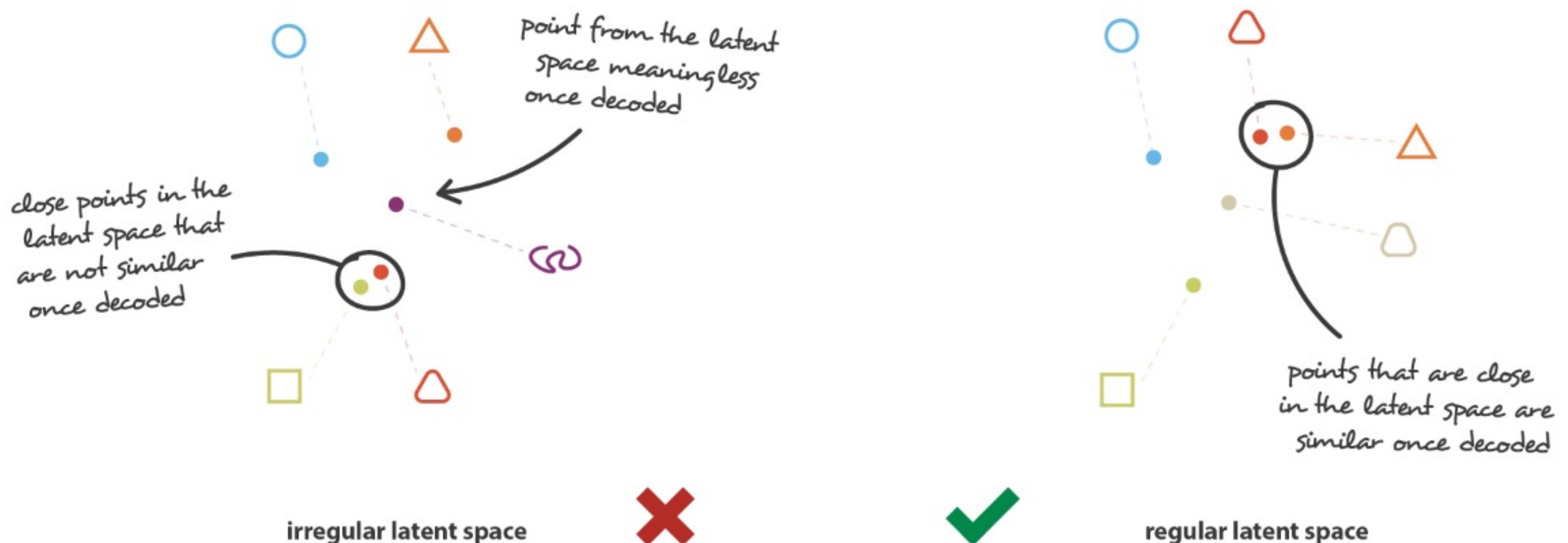
The Autoencoder is trained to maximise reconstruction, not to result in a meaningful latent space.



# GENERATION WITH AUTOENCODERS

## WHAT DO WE WANT TO ACHIEVE?

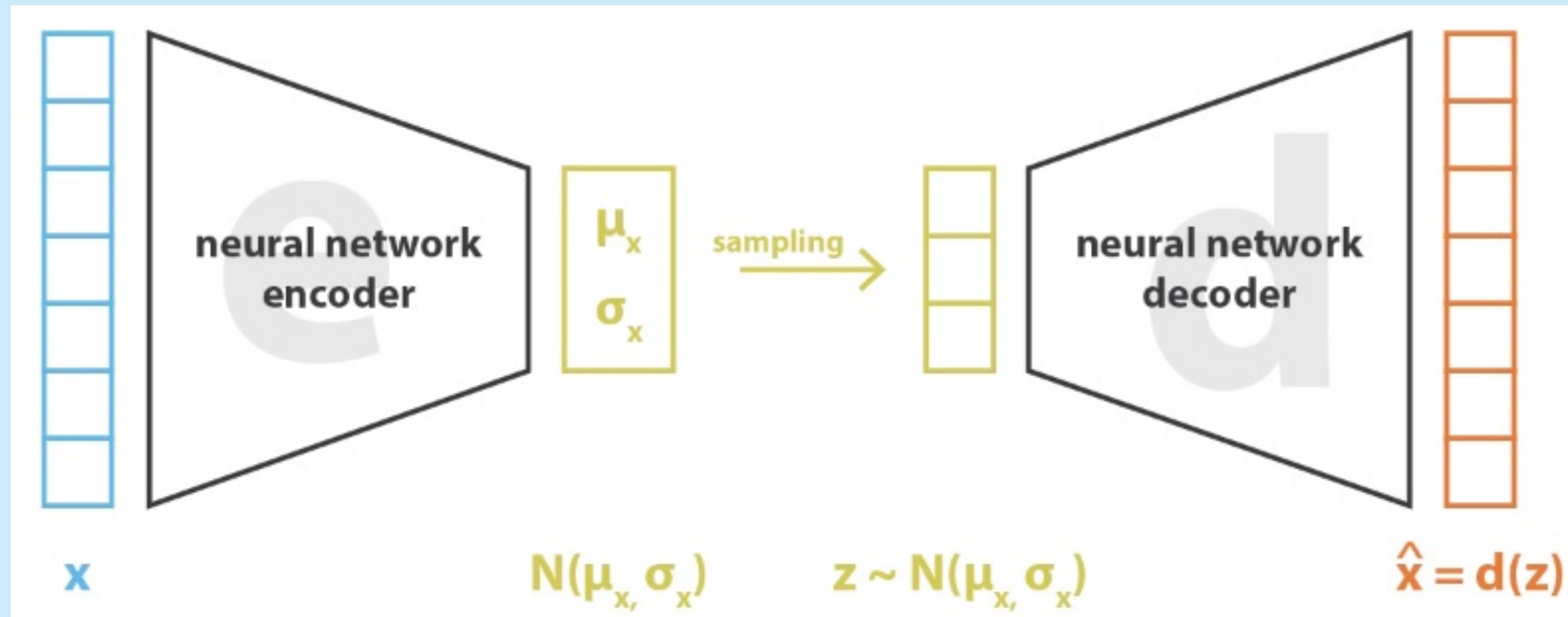
Continuity and Completeness



# GENERATION WITH AUTOENCODERS

## HOW DO WE ACHIEVE OUR GOALS?

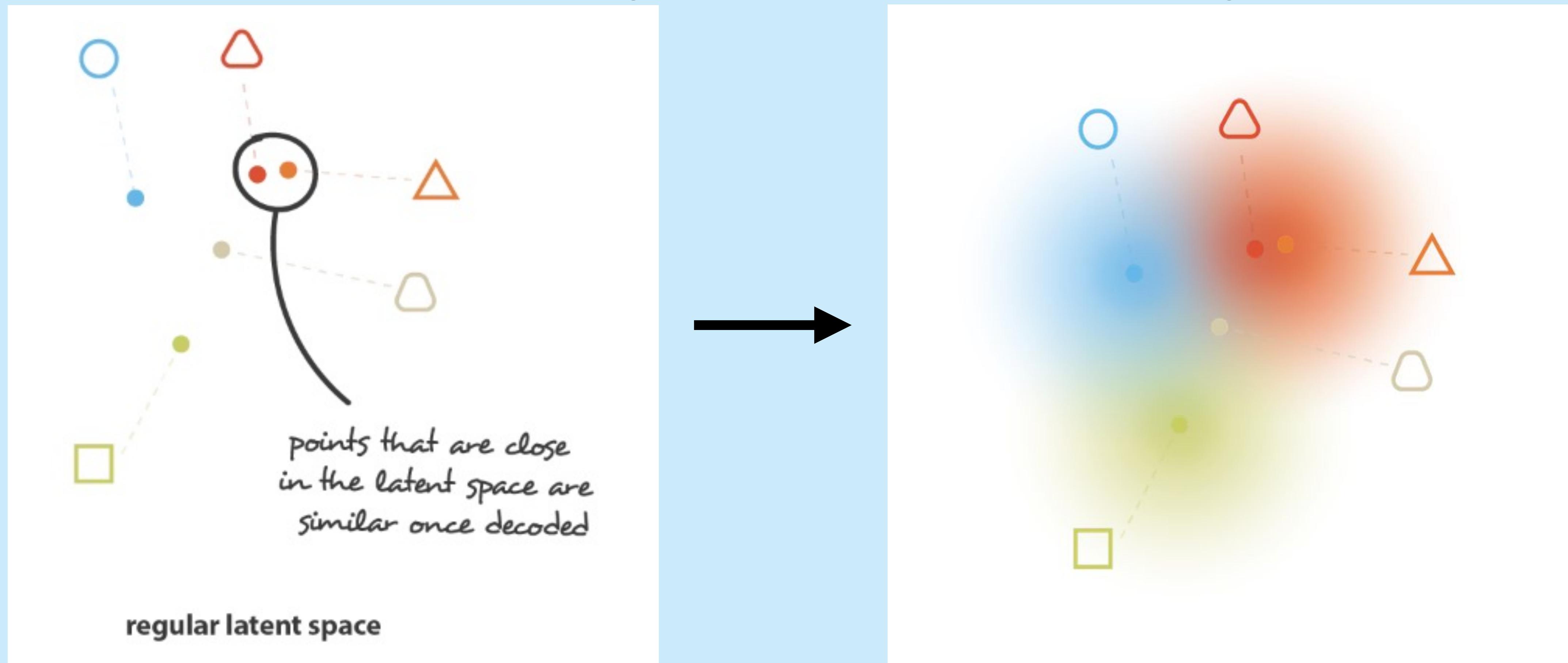
Solution: Encoder gives parameters used for sampling



# GENERATION WITH AUTOENCODERS

## HOW DO WE ACHIEVE OUR GOALS?

Solution: Encoder gives parameters used for sampling



# GENERATION WITH AUTOENCODERS

## HOW DO WE ACHIEVE OUR GOALS?

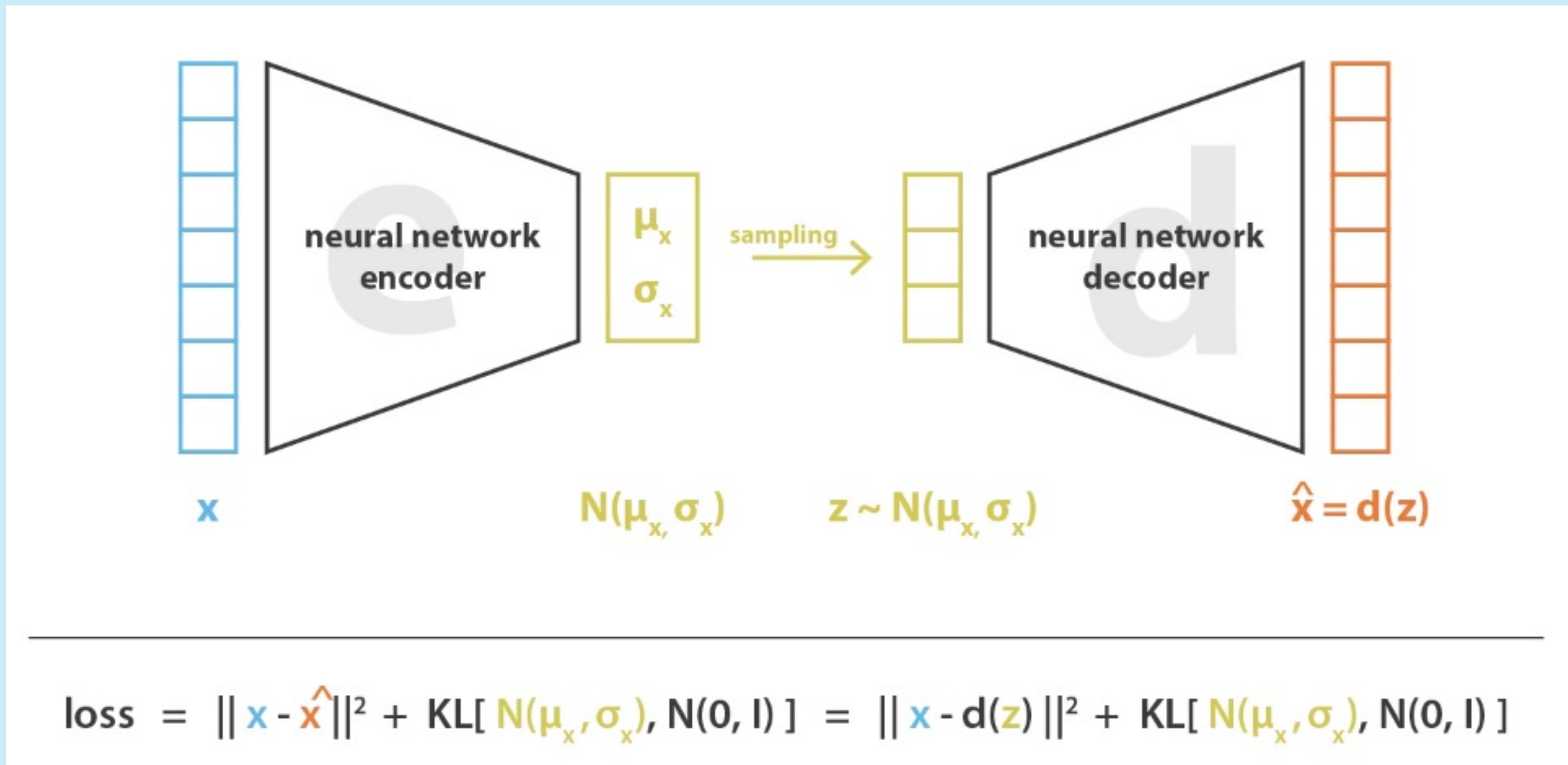
Problem: Still only trained for reconstruction...



# GENERATION WITH AUTOENCODERS

## HOW DO WE ACHIEVE OUR GOALS?

Solution: Adjust the loss term, to enforce a certain distribution



# GENERATION WITH AUTOENCODERS

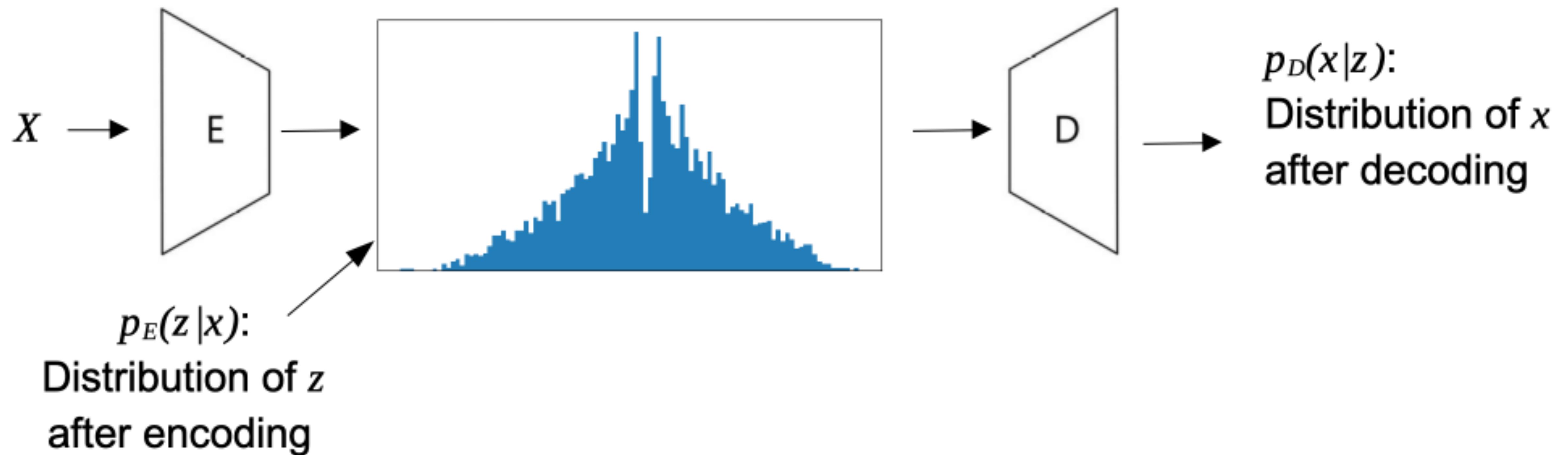
## HOW DO WE ACHIEVE OUR GOALS?

Solution: Adjust the loss term, to enforce a certain distribution

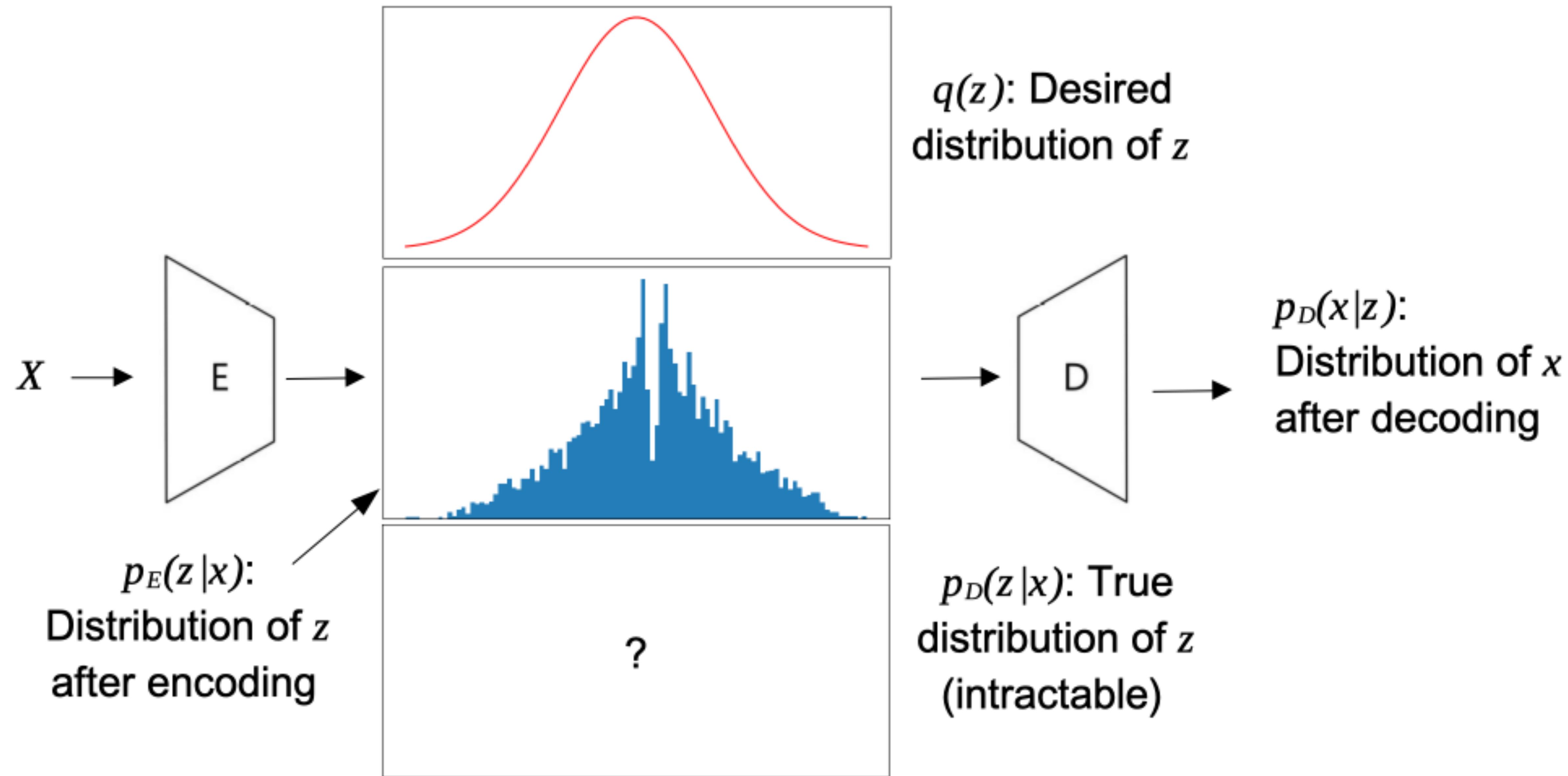


## II. Mathematically

# VARIATIONAL AUTOENCODERS: BASICS



# VARIATIONAL AUTOENCODERS: BASICS



# VARIATIONAL AUTOENCODERS: MATHEMATICAL BACKGROUND

Kullback-Leibler (KL) Divergence: a measure for the similarity of two distributions  $p$  and  $q$

$$\text{KL}(p(x) \parallel q(x)) = \int_{-\infty}^{+\infty} p(x) \log \left( \frac{p(x)}{q(x)} \right) dx$$

## Properties:

Always non-negative, due to Jensen's inequality

Exactly 0 only if the argument of the logarithm is 1, i.e. the two distributions are equal

# VARIATIONAL AUTOENCODERS: MATHEMATICAL BACKGROUND

Goal: Find optimal encoder and decoder to minimise the difference between original and approximation:

$$\hat{E}, \hat{D} = \arg \min_{E,D} \left( \text{KL} (p^*(x) || p(x)) \right)$$

# VARIATIONAL AUTOENCODERS: MATHEMATICAL BACKGROUND

Goal: Find optimal encoder and decoder to minimise the difference between original and approximation:

$$\hat{E}, \hat{D} = \arg \min_{E,D} (\text{KL}(p^*(x) || p(x)))$$

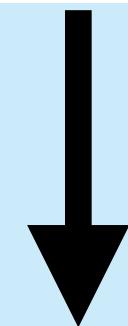
$$\begin{aligned}\text{KL}(p^*(x) || p(x)) &= \int_{-\infty}^{+\infty} p^*(x) \log \left( \frac{p^*(x)}{p(x)} \right) dx \\ &= \int_{-\infty}^{\infty} p^*(x) \log(p^*(x)) dx - \int_{-\infty}^{\infty} p^*(x) \log(p(x)) dx\end{aligned}$$

# VARIATIONAL AUTOENCODERS: MATHEMATICAL BACKGROUND

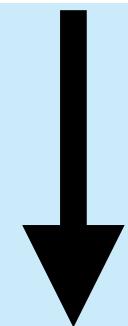
Goal: Find optimal encoder and decoder to minimise the difference between original and approximation:

$$\hat{E}, \hat{D} = \arg \min_{E,D} (\text{KL} (p^*(x) || p(x)))$$

$$\begin{aligned}\text{KL} (p^*(x) || p(x)) &= \int_{-\infty}^{+\infty} p^*(x) \log \left( \frac{p^*(x)}{p(x)} \right) dx \\ &= \int_{-\infty}^{\infty} p^*(x) \log (p^*(x)) dx - \int_{-\infty}^{\infty} p^*(x) \log (p(x)) dx\end{aligned}$$



Independent of model!



$\mathbb{E}_{x \sim p^*} [-\log (p(x))]$

# NEED TO MAXIMISE THE LOG-LIKELIHOOD

$$\begin{aligned} & \log(p(x)) \\ &= \log(p(x)) \int_{-\infty}^{+\infty} p_E(z|x) dz \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log(p(x)) dz \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log \left( p(x) \cdot \frac{p_E(z|x)}{p_E(z|x)} \cdot \frac{q(z)}{q(z)} \cdot \frac{p_D(x|z)}{p_D(x|z)} \right) dz \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \left[ \log \left( \frac{p_E(z|x)}{q(z)p_D(x|z)/p(x)} \right) - \log \left( \frac{p_E(z|x)}{q(z)} \right) + \log(p_D(x|z)) \right] dz \\ &= \text{KL}(p_E(z|x) || p_D(z|x)) - \text{KL}(p_E(z|x) || q(z)) + \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] \end{aligned}$$

# NEED TO MAXIMISE THE LOG-LIKELIHOOD

$$\begin{aligned} & \log(p(x)) \\ &= \log(p(x)) \int_{-\infty}^{+\infty} p_E(z|x) dz \quad \xrightarrow{\text{Insert 1}} \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log(p(x)) dz \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log \left( p(x) \cdot \frac{p_E(z|x)}{p_E(z|x)} \cdot \frac{q(z)}{q(z)} \cdot \frac{p_D(x|z)}{p_D(x|z)} \right) dz \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \left[ \log \left( \frac{p_E(z|x)}{q(z)p_D(x|z)/p(x)} \right) - \log \left( \frac{p_E(z|x)}{q(z)} \right) + \log(p_D(x|z)) \right] dz \\ &= \text{KL}(p_E(z|x) || p_D(z|x)) - \text{KL}(p_E(z|x) || q(z)) + \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] \end{aligned}$$

# NEED TO MAXIMISE THE LOG-LIKELIHOOD

$$\begin{aligned} & \log(p(x)) \\ &= \log(p(x)) \int_{-\infty}^{+\infty} p_E(z|x) dz \quad \text{Insert 1} \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log(p(x)) dz \quad \text{Log independent of } z \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log \left( p(x) \cdot \frac{p_E(z|x)}{p_E(z|x)} \cdot \frac{q(z)}{q(z)} \cdot \frac{p_D(x|z)}{p_D(x|z)} \right) dz \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \left[ \log \left( \frac{p_E(z|x)}{q(z)p_D(x|z)/p(x)} \right) - \log \left( \frac{p_E(z|x)}{q(z)} \right) + \log(p_D(x|z)) \right] dz \\ &= \text{KL}(p_E(z|x) || p_D(z|x)) - \text{KL}(p_E(z|x) || q(z)) + \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] \end{aligned}$$

# NEED TO MAXIMISE THE LOG-LIKELIHOOD

$$\begin{aligned} & \log(p(x)) \\ &= \log(p(x)) \int_{-\infty}^{+\infty} p_E(z|x) dz \quad \text{Insert 1} \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log(p(x)) dz \quad \text{Log independent of } z \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log \left( p(x) \cdot \frac{p_E(z|x)}{p_E(z|x)} \cdot \frac{q(z)}{q(z)} \cdot \frac{p_D(x|z)}{p_D(x|z)} \right) dz \quad \text{Insert 1 again} \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \left[ \log \left( \frac{p_E(z|x)}{q(z)p_D(x|z)/p(x)} \right) - \log \left( \frac{p_E(z|x)}{q(z)} \right) + \log(p_D(x|z)) \right] dz \\ &= \text{KL}(p_E(z|x) || p_D(z|x)) - \text{KL}(p_E(z|x) || q(z)) + \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] \end{aligned}$$

# NEED TO MAXIMISE THE LOG-LIKELIHOOD

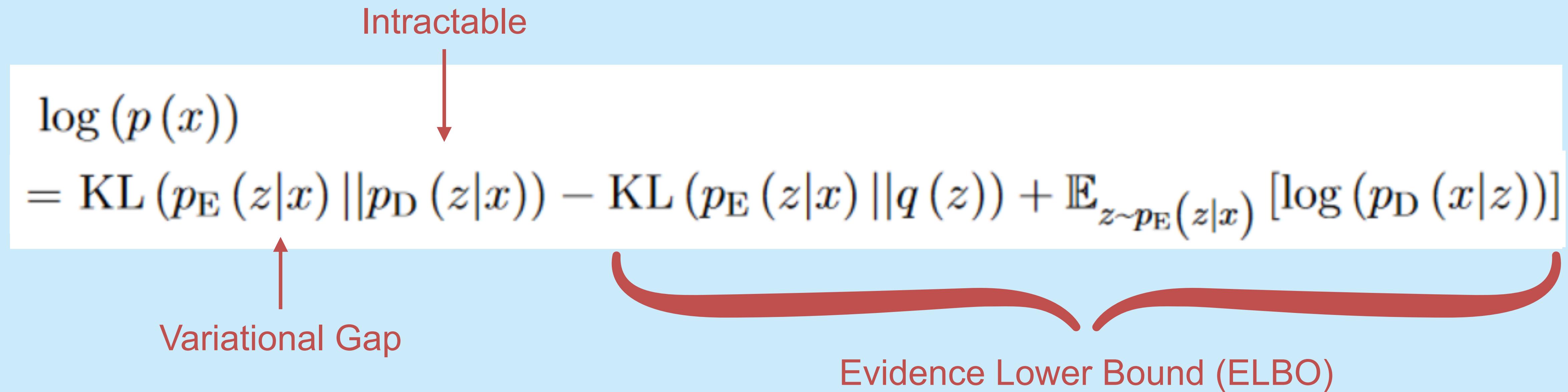
$$\begin{aligned} & \log(p(x)) \\ &= \log(p(x)) \int_{-\infty}^{+\infty} p_E(z|x) dz \quad \text{Insert 1} \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log(p(x)) dz \quad \text{Log independent of } z \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \log \left( p(x) \cdot \frac{p_E(z|x)}{p_E(z|x)} \cdot \frac{q(z)}{q(z)} \cdot \frac{p_D(x|z)}{p_D(x|z)} \right) dz \quad \text{Insert 1 again} \\ &= \int_{-\infty}^{+\infty} p_E(z|x) \left[ \log \left( \frac{p_E(z|x)}{q(z)p_D(x|z)/p(x)} \right) - \log \left( \frac{p_E(z|x)}{q(z)} \right) + \log(p_D(x|z)) \right] dz \\ &= \text{KL}(p_E(z|x) || p_D(z|x)) - \text{KL}(p_E(z|x) || q(z)) + \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] \end{aligned}$$

Bayes Theorem

# DERIVING THE ELBO

$$\begin{aligned} & \log(p(x)) \\ &= \text{KL}(p_E(z|x) || p_D(z|x)) - \text{KL}(p_E(z|x) || q(z)) + \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] \end{aligned}$$

# DERIVING THE ELBO



# MAXIMISING THE ELBO

$$\text{ELBO} = \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] - \text{KL}(p_E(z|x) || q(z))$$



Reconstruction Error:  
Original data  $x$  should be a highly  
probable output



Optimum reached for perfect  
reconstruction without variability, i.e.  
deterministic model

→ No generative abilities



How close is the latent space  
distribution to the desired distribution



Optimum reached for  $z$  independent of  $x$

→ No reconstruction abilities

# EXPLICIT LOSS EXPRESSION FOR VAE

$$\hat{E}, \hat{D} = \arg \min_{E,D} (-\text{ELBO}) \quad \text{with}$$

$$\text{ELBO} = \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] - \text{KL}(p_E(z|x) || q(z))$$

# EXPLICIT LOSS EXPRESSION FOR VAE

$$\hat{E}, \hat{D} = \arg \min_{E, D} (-\text{ELBO}) \quad \text{with}$$

$$\text{ELBO} = \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] - \text{KL}(p_E(z|x) || q(z))$$

Assumptions:

Standard Normal Prior:

$$q(z) = \mathcal{N}(0, \mathbf{1})$$

Gaussian Encoder Distribution:

$$p_E(z|x) = \mathcal{N}\left(z|\mu_E(x), \text{Diag}(\sigma(x)^2)\right)$$

Gaussian Decoder with Uncertainty:

$$p_D(x|z) = \mathcal{N}(x|\mu_D(z), \beta^2 \mathbf{1})$$

# EXPLICIT LOSS EXPRESSION FOR VAE

$$\hat{E}, \hat{D} = \arg \min_{E,D} (-\text{ELBO}) \quad \text{with} \quad \text{ELBO} = \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] - \text{KL}(p_E(z|x) || q(z))$$

$$\begin{aligned}\mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] &= \mathbb{E}_{z \sim p_E(z|x)} [\log(\mathcal{N}(x|\mu_D(z), \beta^2 \mathbf{I}))] \\ &= \mathbb{E}_{z \sim p_E(z|x)} \left[ -\frac{(x - \mu_D(z))^2}{2\beta^2} + \log(2\pi\beta^2)^{d/2} \right]\end{aligned}$$

# EXPLICIT LOSS EXPRESSION FOR VAE

$$\hat{E}, \hat{D} = \arg \min_{E,D} (-\text{ELBO}) \quad \text{with} \quad \text{ELBO} = \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] - \text{KL}(p_E(z|x) || q(z))$$

$$\begin{aligned}\mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] &= \mathbb{E}_{z \sim p_E(z|x)} [\log(\mathcal{N}(x|\mu_D(z), \beta^2 \mathbf{1}))] \\ &= \mathbb{E}_{z \sim p_E(z|x)} \left[ -\frac{(x - \mu_D(z))^2}{2\beta^2} + \log(\cancel{\beta^2})^{d/2} \right]\end{aligned}$$

$$\mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] = -\frac{1}{M} \sum_{k=1}^M \left( \frac{(x_k - \mu_{D,k}(z))^2}{2\beta^2} \right)$$

# EXPLICIT LOSS EXPRESSION FOR VAE

$$\hat{E}, \hat{D} = \arg \min_{E,D} (-\text{ELBO}) \quad \text{with}$$

$$\text{ELBO} = \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] - \text{KL}(p_E(z|x) || q(z))$$

$$\begin{aligned} -\text{KL}(p_E(z|x) || q(z)) &= -\text{KL}\left(\mathcal{N}\left(z|\mu_E(x), \text{Diag}(\sigma(x)^2)\right) || \mathcal{N}(0, I)\right) \\ &= -\frac{1}{2} \sum_{j=1}^d \left( \mu_j(x)^2 + \sigma_j(x)^2 - \log(\sigma_j(x)^2) - 1 \right) \end{aligned}$$

# EXPLICIT LOSS EXPRESSION FOR VAE

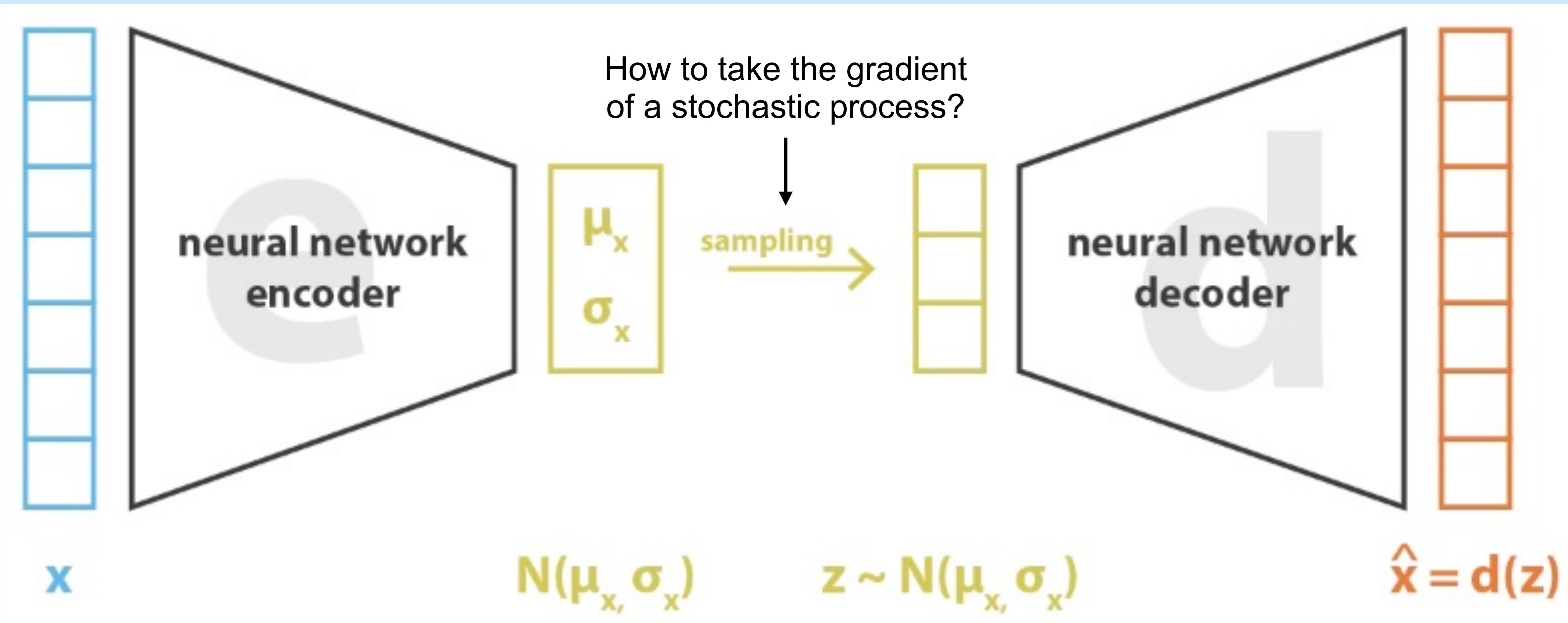
$$\hat{E}, \hat{D} = \arg \min_{E, D} (-\text{ELBO}) \quad \text{with}$$

$$\text{ELBO} = \mathbb{E}_{z \sim p_E(z|x)} [\log(p_D(x|z))] - \text{KL}(p_E(z|x) || q(z))$$

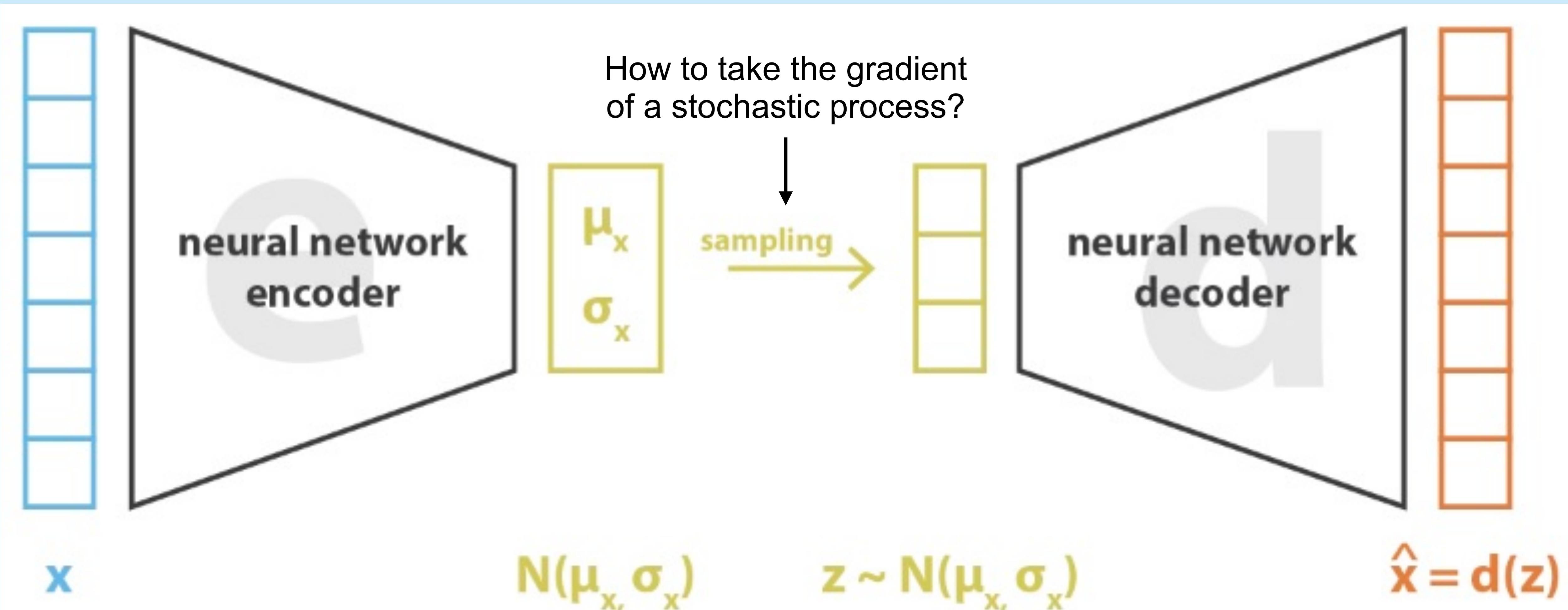
Loss Function for training a VAE

$$-\text{ELBO} = \frac{1}{2\beta^2} \cdot \frac{1}{M} \sum_{k=1}^M \left( (x_k - \mu_{D,k}(z))^2 \right) + \frac{1}{2} \sum_{j=1}^d \left( \mu_j(x)^2 + \sigma_j(x)^2 - \log(\sigma_j(x)^2) - 1 \right)$$

# BUT... ONE MORE PROBLEM



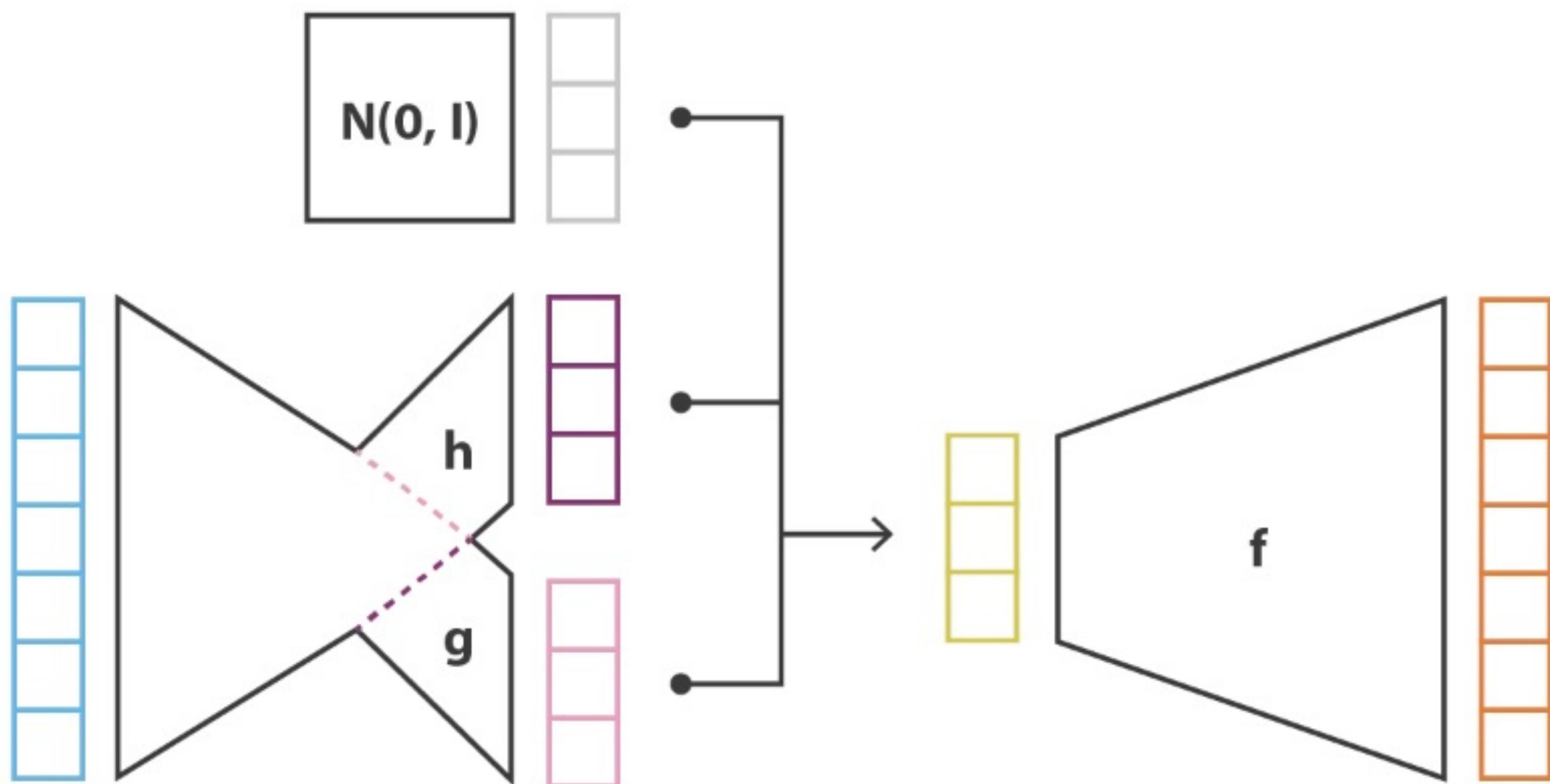
# BUT... ONE MORE PROBLEM



Reparameterisation  
Trick:

$$z \sim \mathcal{N}(\mu, \sigma) \rightarrow z = \mu + \zeta \cdot \sigma \text{ with } \zeta \sim \mathcal{N}(0, 1)$$

# BUT... ONE MORE PROBLEM



$$\mu_x = g(x)$$

$$\sigma_x = h(x)$$

$$\zeta \sim N(0, I)$$

$$z = \sigma_x \zeta + \mu_x$$

$$\hat{x} = f(z)$$

$$\text{loss} = C \| x - \hat{x} \|^2 + \text{KL}[ N(\mu_x, \sigma_x), N(0, I) ] = C \| x - f(z) \|^2 + \text{KL}[ N(g(x), h(x)), N(0, I) ]$$

### III. Practically: Code examples

# DIFFERENT VAE FLAVOURS

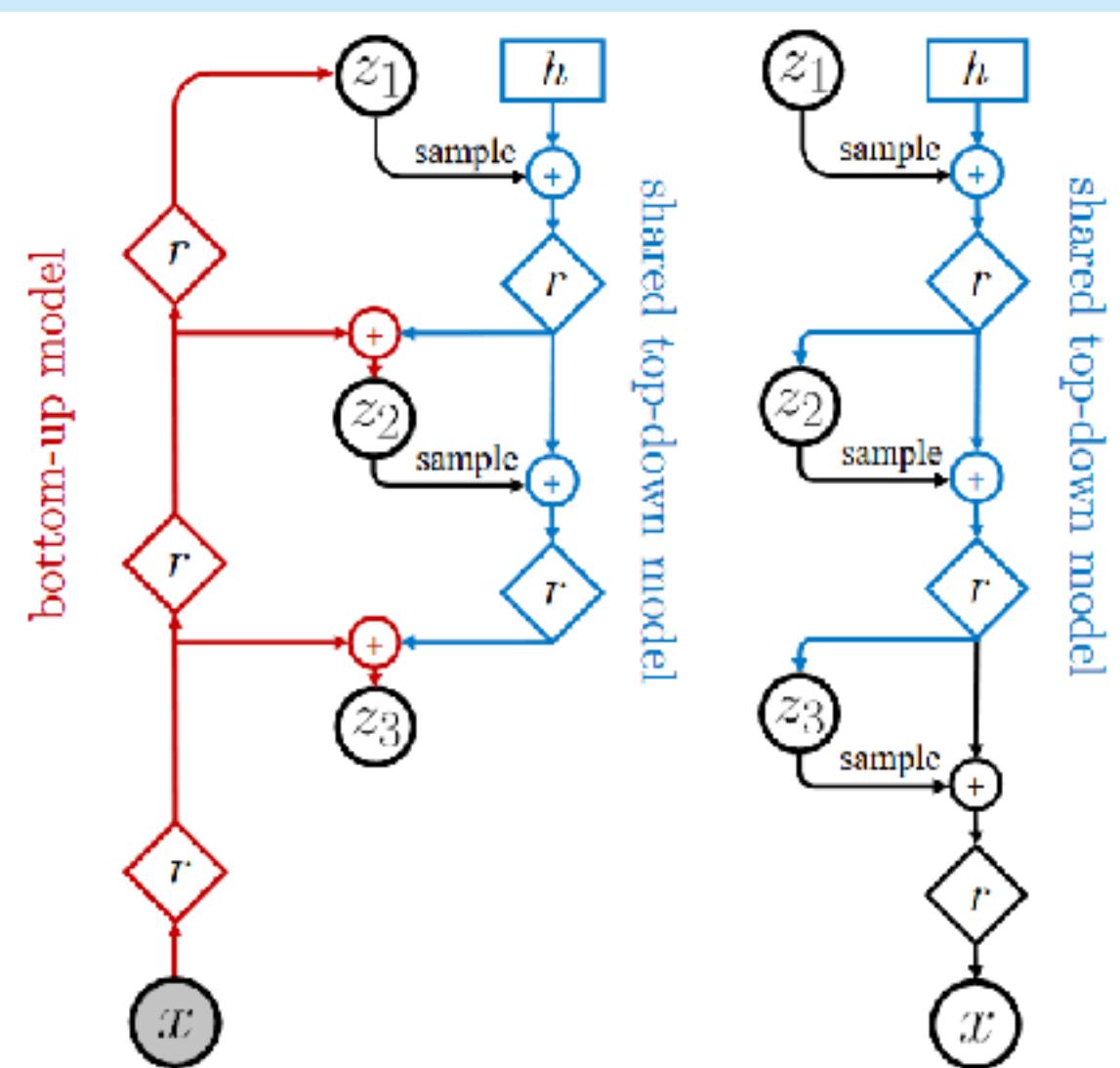
## 1. AVAE (Adversarial Variational Autoencoder)

- Combines VAEs and GANs to minimise problems with each model
- Avoid GAN model collapse with latent space
- Create more realistic samples with additional “realism” loss

## 2. $\beta$ -VAE

## 3. k-sparses VAE

## 4. NVEA (Nouveau Variational Autoencoder)



## Astro Examples of VAE

# DIFFERENT VAE APPLICATIONS IN ASTRO

1. Dimensionality Reduction of SDSS Spectra with Variational Autoencoders (Portillo+2020)
2. Deblending galaxies with variational autoencoders: A joint multiband, multi-instrument approach (Arcelin+2020)
3. Bayesian decomposition of the Galactic multi-frequency sky using probabilistic autoencoders (Milosevic+2021)
4. Bayesian parameter estimation using conditional variational autoencoders for gravitational-wave astronomy (Gabbard+2022)
5. PELS-VAE: Physics-Enhanced Latent Space Variational Autoencoder for physically informed light curves of periodic variable stars (Martínez-Palomera+2020)
5. Model comparison and model misspecification for galaxy simulations (Zhou+2024 subm.)

# Bayesian Model Comparison

# MODEL COMPARISON & MODEL MISSPECIFICATION

**Bayesian model selection** assigns posterior probabilities  $p(\mathcal{M}_k \mid \mathbf{d})$  to models  $\mathcal{M}_k \in \{\mathcal{M}_1, \dots, \mathcal{M}_N\}$  (instead of to values of their parameters  $\boldsymbol{\theta}_k$ ), conditional on observed data  $\mathbf{d}$ . The conventional approach is to compute the marginal likelihood (or *evidence*)  $p(\mathbf{d} \mid \mathcal{M}_k)$ , which is the average likelihood  $p(\mathbf{d} \mid \boldsymbol{\theta}_k)$  of parameters distributed according to the prior  $p(\boldsymbol{\theta}_k)$ :

$$p(\mathbf{d} \mid \mathcal{M}_k) = \int p(\mathbf{d} \mid \boldsymbol{\theta}_k) p(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k \quad (1)$$

(where the presence of  $\mathcal{M}_k$ 's parameters  $\boldsymbol{\theta}_k$  implies conditioning on  $\mathcal{M}_k$  in the right-hand side). The prior belief in the model,  $p(\mathcal{M}_k)$ , is then updated to its posterior probability in accordance with Bayes' theorem:  $p(\mathcal{M}_k \mid \mathbf{d}) \propto p(\mathbf{d} \mid \mathcal{M}_k) p(\mathcal{M}_k)$ , normalised over all models considered.

# MODEL COMPARISON & MODEL MISSPECIFICATION

## Methods

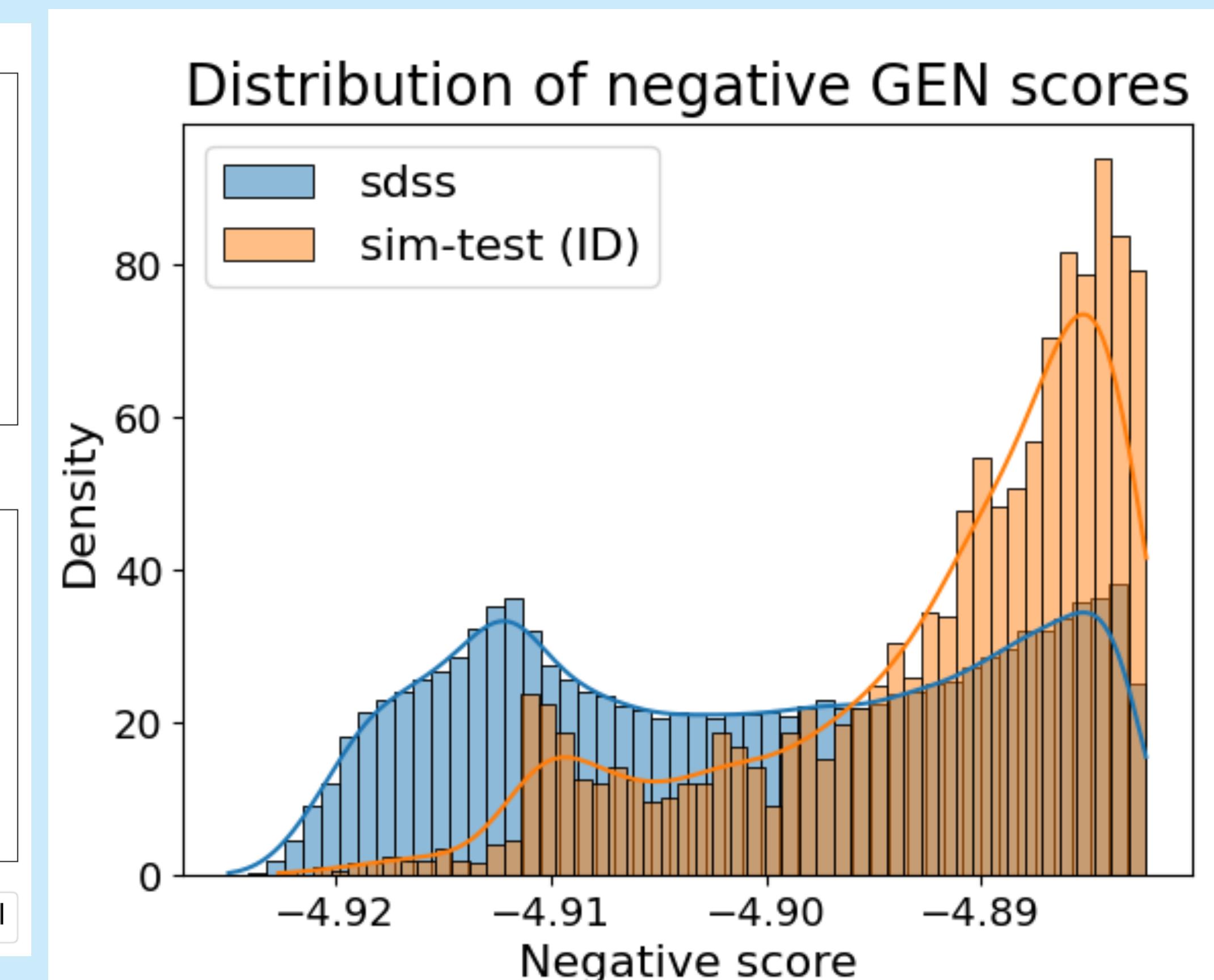
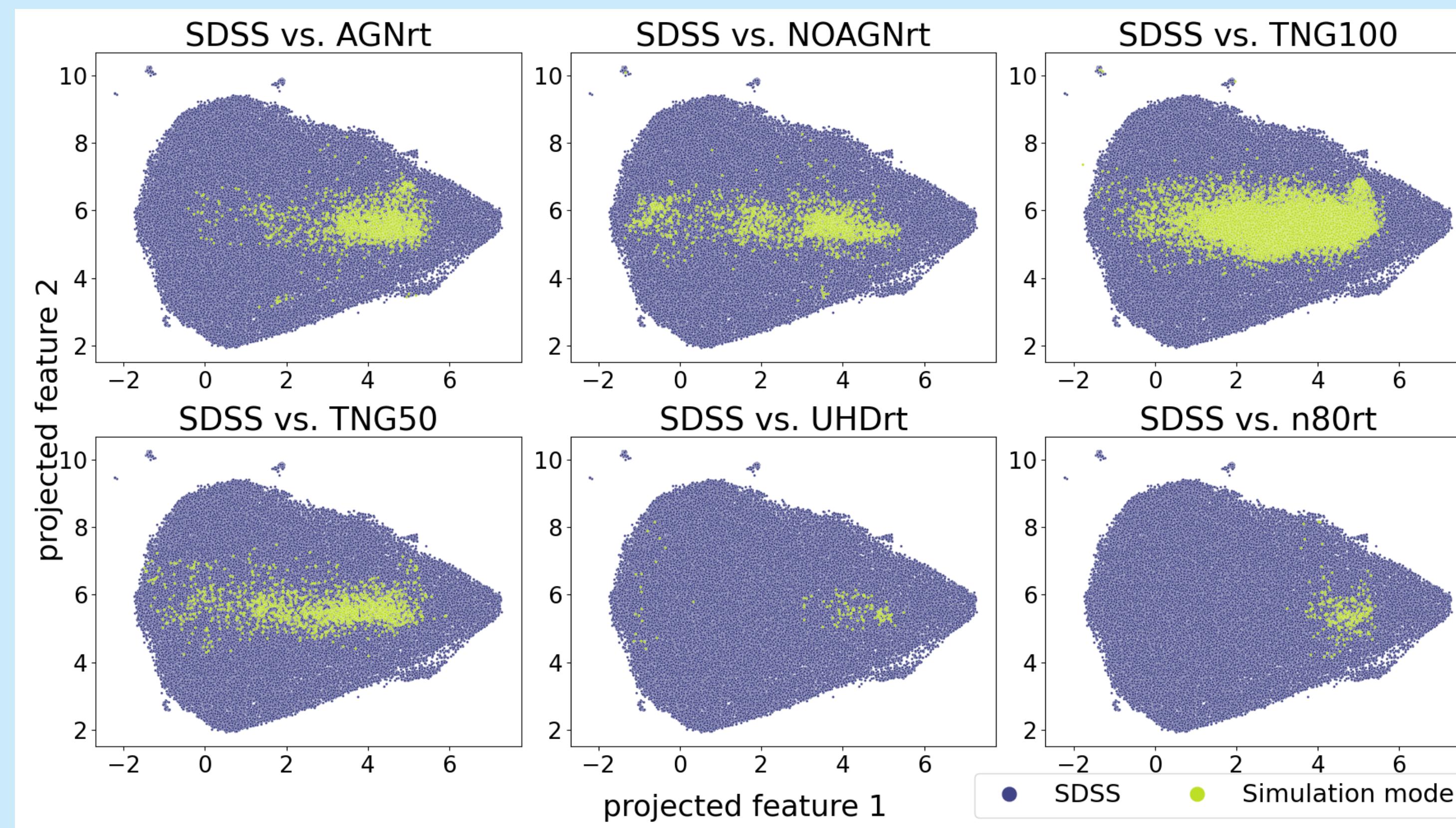
- Jeffrey+Wandelt 2023: loss functions for two-way model comparison with an emphasis on recovering accurate extreme Bayes factors
- Radev+2021: estimate a Dirichlet distribution over an arbitrary number of models using a NN and variational optimisation
- Elsemüller+2023 and Karchev+2024: use a cross-entropy loss for multi-class posterior probabilities
- Macciò+2022: Model selection for star formation prescriptions in cosmo sims
- Zhou+2024: Model misspecification plus model comparison for low simulation budget applications
- Jin+2024: Model comparison of cosmo sims via GANomaly scores

# MODEL COMPARISON & MODEL MISSPECIFICATION

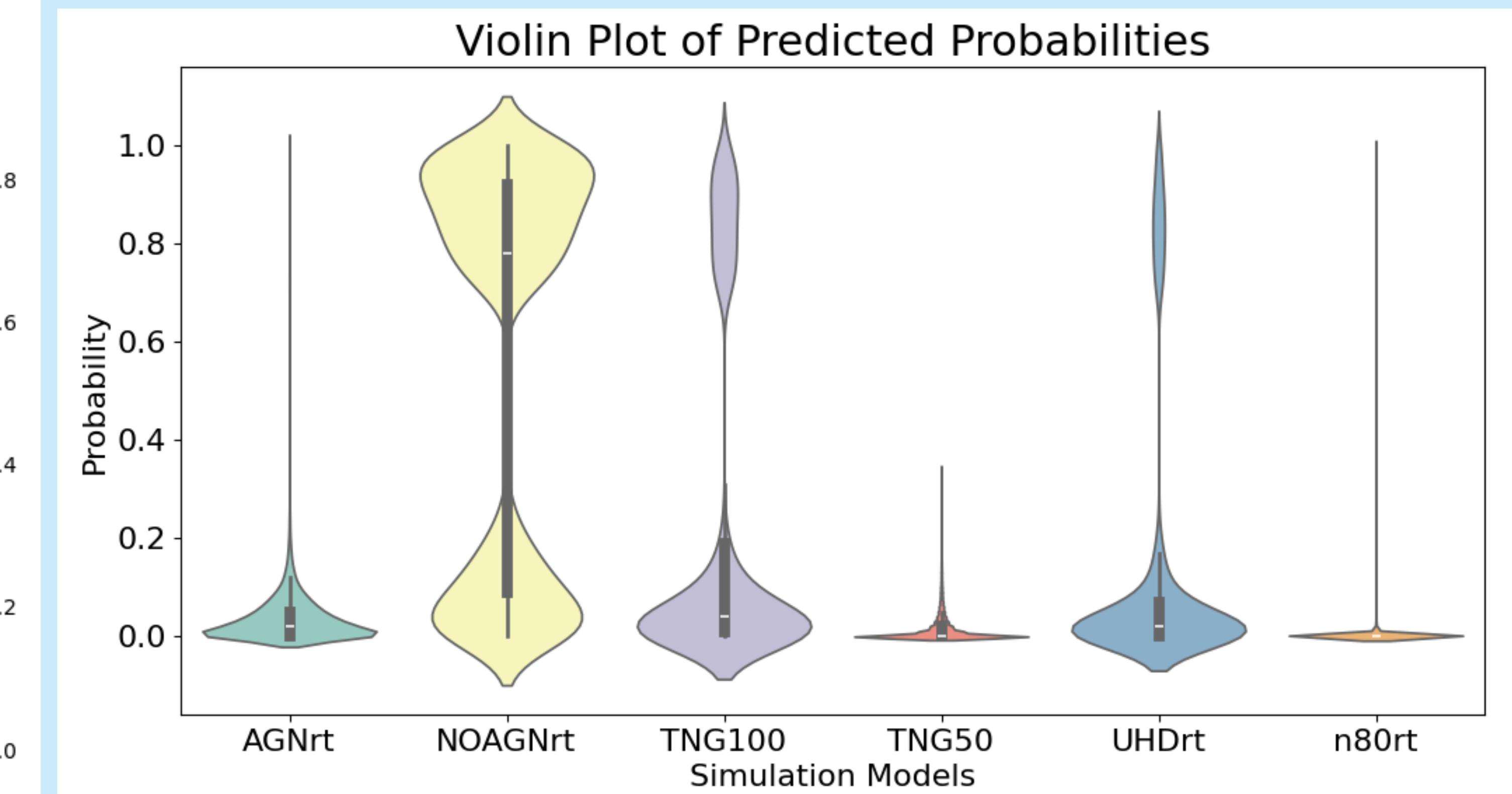
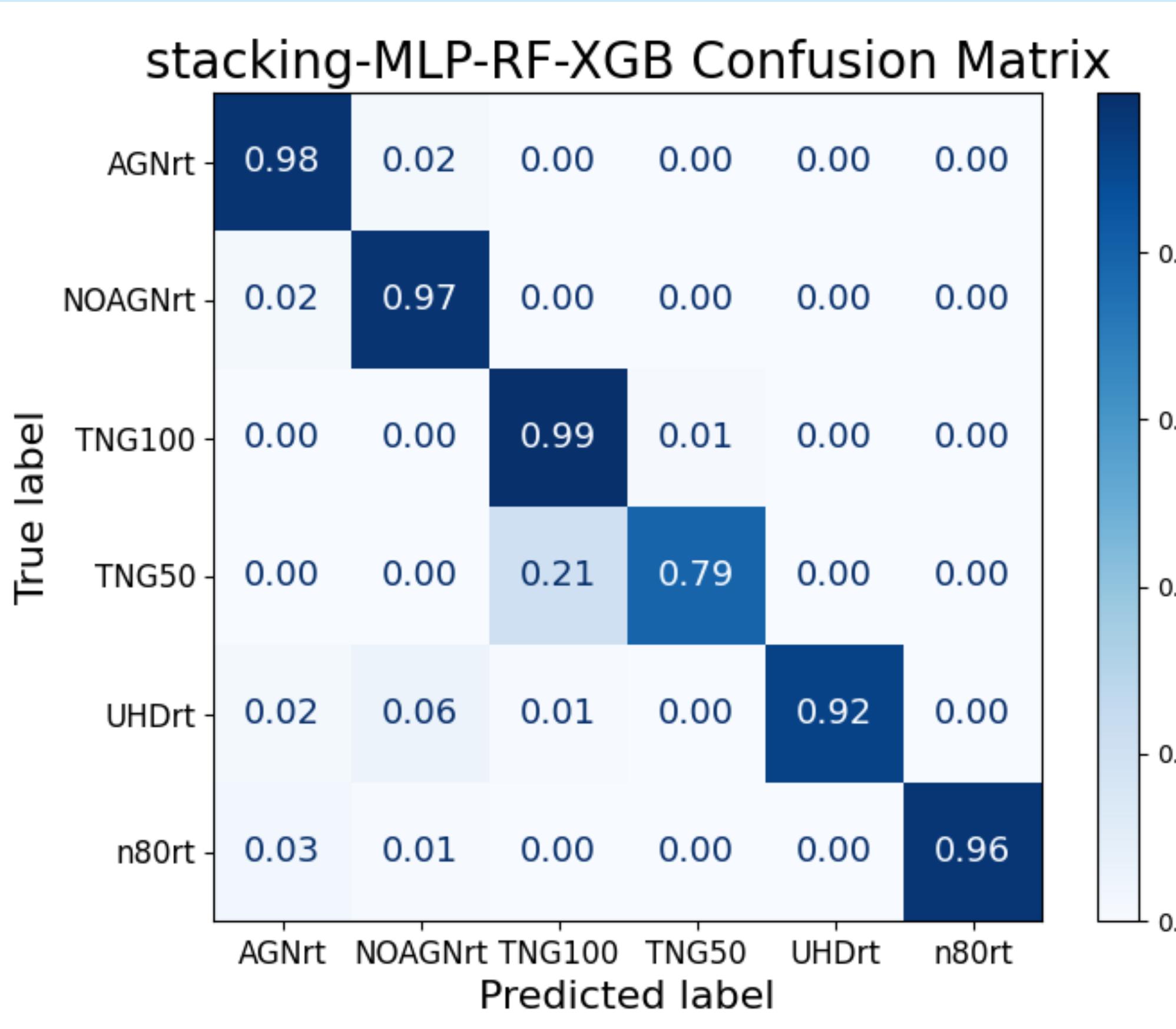
Idea:

- Train embedding network on ~600.000 SDSS images, then encode simulated SDSS images
- Train simulation classifiers on embeddings, apply to real SDSS images

UMAP projection of k-sparse encoding of sims and Jobs

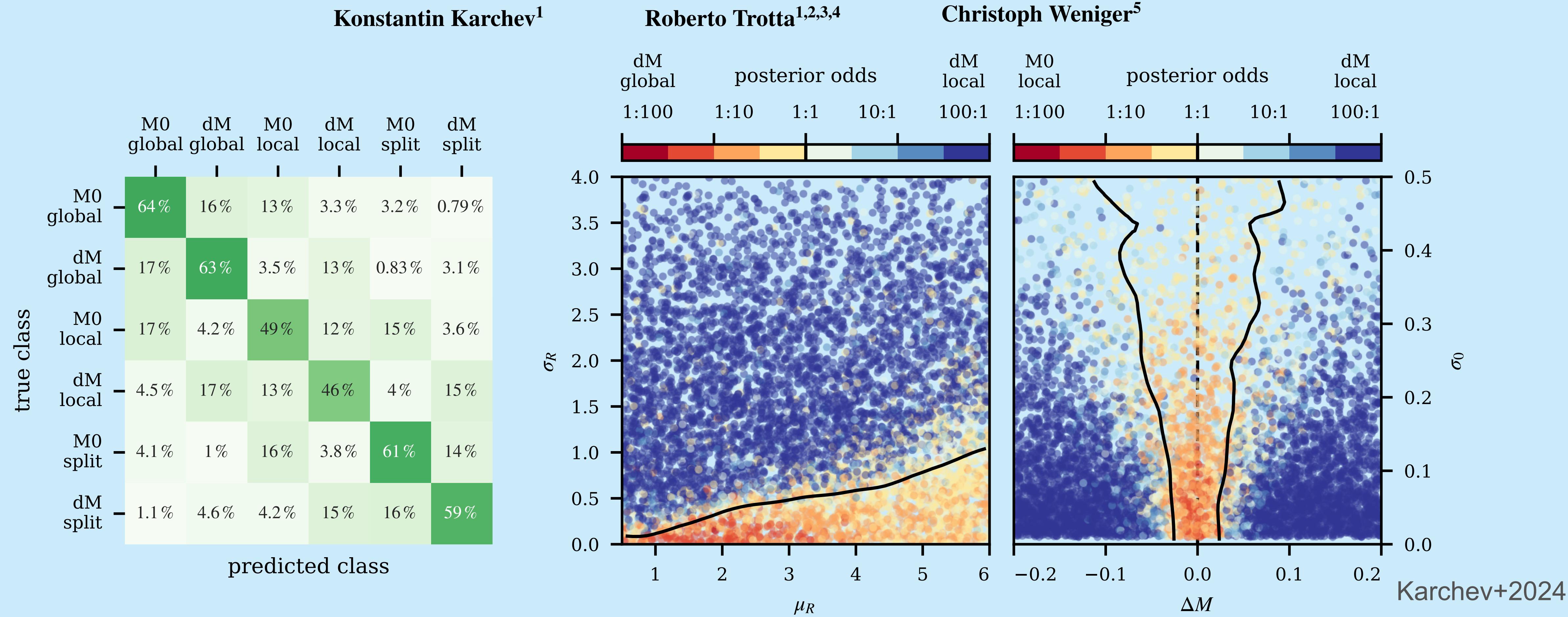


# MODEL COMPARISON & MODEL MISSPECIFICATION



# MODEL COMPARISON/SELECTION

## SimSIMS: Simulation-based Supernova Ia Model Selection with thousands of latent variables



# WHERE TO FIND THE LECTURE MATERIAL



<https://github.com/TobiBu/graddays>

# Score Matching and diffusion models

# SCORE MATCHING AND DIFFUSION MODELS

What is the score  $s(\mathbf{x})$  of a pdf  $p(\mathbf{x})$ ?

$$s(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x})$$

Gibbs measure: 
$$p(x) = \frac{1}{\mathcal{Z}} \exp(-\beta E(x))$$

$$\log p(x) = -\beta E(x) - \log \mathcal{Z}$$

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}) = -\beta \nabla_x E(x)$$

# SCORE MATCHING AND DIFFUSION MODELS

What is the score  $s(\mathbf{x})$  of a pdf  $p(\mathbf{x})$ ?

$$s(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x})$$

Gibbs measure: 
$$p(x) = \frac{1}{\mathcal{Z}} \exp(-\beta E(x))$$

$$\log p(x) = -\beta E(x) - \log \mathcal{Z}$$

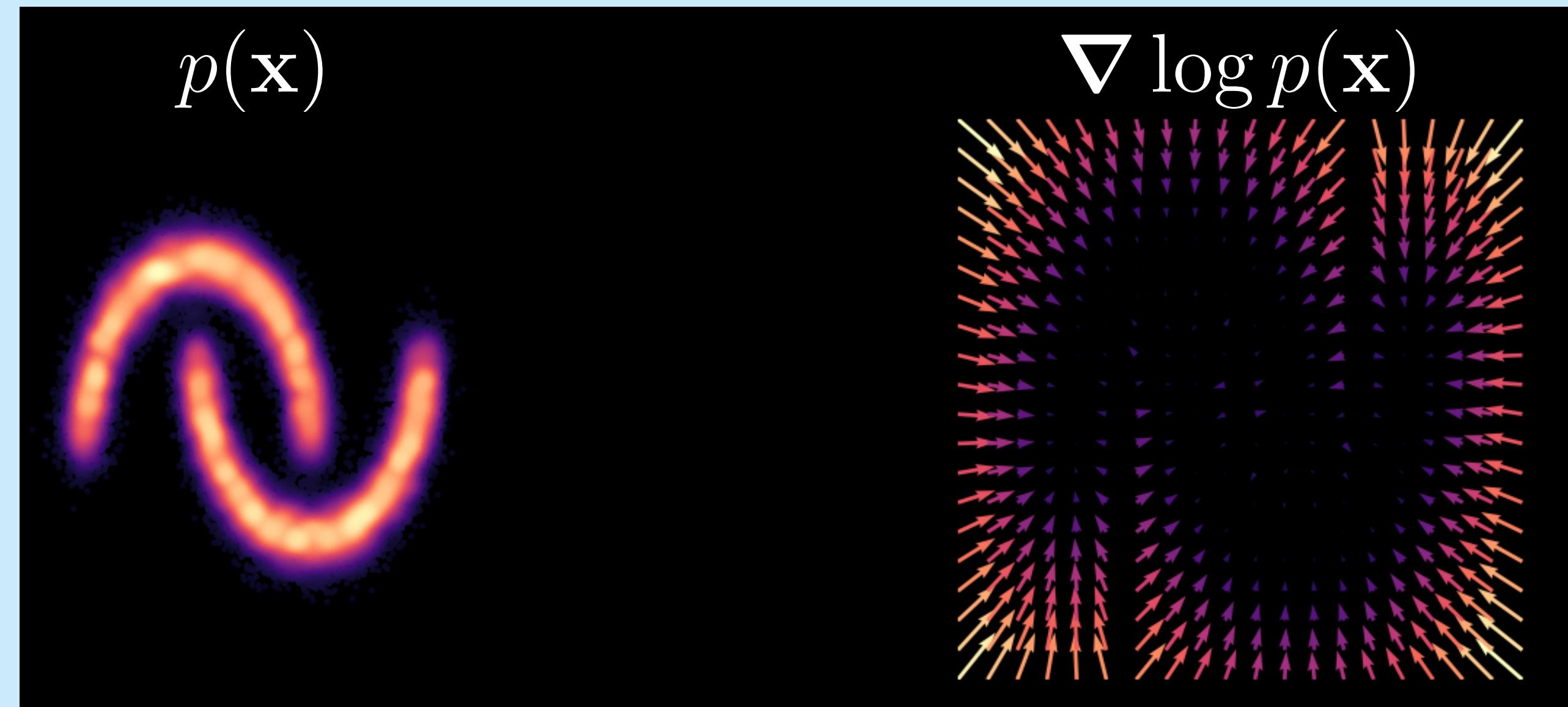
$$\nabla_{\mathbf{x}} \log p(\mathbf{x}) = -\beta \nabla_x E(x)$$

local and easier to learn

# SCORE MATCHING AND DIFFUSION MODELS

What is the score  $s(\mathbf{x})$  of a pdf  $p(\mathbf{x})$ ?

$$s(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x})$$



Idea: Learn  $p(\mathbf{x})$  solely from data samples, then sample new instances.

# TRAINING OBJECTIVE FOR SCORE MATCHING

train via Fisher divergence

$$\mathbb{E}_{p(x)}[\|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - s_\theta(x)\|_2^2]$$

however, this needs access to the data score  $\nabla_x \log p(x)$

score matching circumvents this!

additionally, score matching does not require  $s_\theta(x)$  to be an actual score function of any normalised distribution.

only requirement: vector valued function of same dimensional input/output

But how does this work?

# TRAINING OBJECTIVE FOR SCORE MATCHING

train via Fisher divergence

$$\mathbb{E}_{p(x)}[\|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - s_{\theta}(x)\|_2^2]$$

using integration by parts and some math magic we end up with (Hyvärinen+2005)

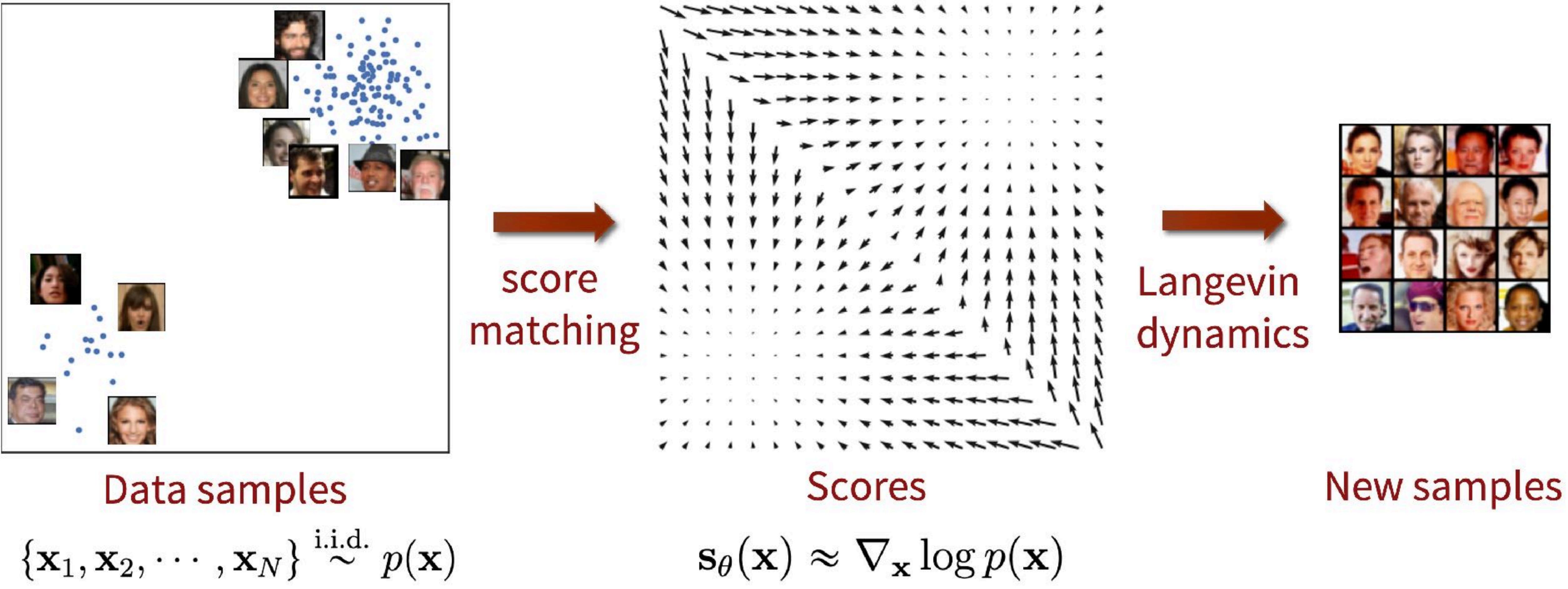
$$\begin{aligned} & \mathbb{E}_{p(x)}\left[\left\|\frac{1}{2}s_{\theta}(x)\right\|_2^2 + \text{trace}(\nabla_{\mathbf{x}} \log p(\mathbf{x}))\right] \\ & \approx \sum_{i=1}^N \left[\frac{1}{2}\|s_{\theta}(x)\|_2^2 + \text{trace}(\nabla_{\mathbf{x}} \log p(\mathbf{x}))\right] \end{aligned}$$

Trace of the Jacobian of  $s$   
Not scalable,  
but can be efficiently  
approximated in high d

But how does sampling then work? → Langevin dynamics (Welling&Teh 2011)

# Langevin dynamics

# LANGEVIN DYNAMICS CONCEPT



$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

$$\mathbf{s}_\theta(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$$

# MATH OF LANGEVIN DYNAMICS

Langevin dynamics provides an MCMC procedure to sample from a distribution  $p(x)$  using only its score function  $\nabla_x \log p(x)$

specifically, it initialises the chain from an arbitrary prior distribution  $x_0 \sim \pi(x)$  and iterates

$$x_{i+1} \leftarrow x_i + \epsilon \nabla_x \log p(x) + \sqrt{2\epsilon} z_i, \quad i = 0, 1, \dots, K$$

where  $z_i \sim \mathcal{N}(0, I)$

When  $\epsilon \rightarrow 0$  and  $K \rightarrow \infty$  the  $x_K$  from above converge to samples from  $p(x)$

# MATH OF LANGEVIN DYNAMICS

Langevin dynamics provides an MCMC procedure to sample from a distribution  $p(x)$  using only its score function  $\nabla_x \log p(x)$

specifically, it initialises the chain from an arbitrary prior distribution  $x_0 \sim \pi(x)$  and iterates

$$x_{i+1} \leftarrow x_i + \epsilon \boxed{\nabla_x \log p(x)} + \sqrt{2\epsilon} z_i, \quad i = 0, 1, \dots, K$$

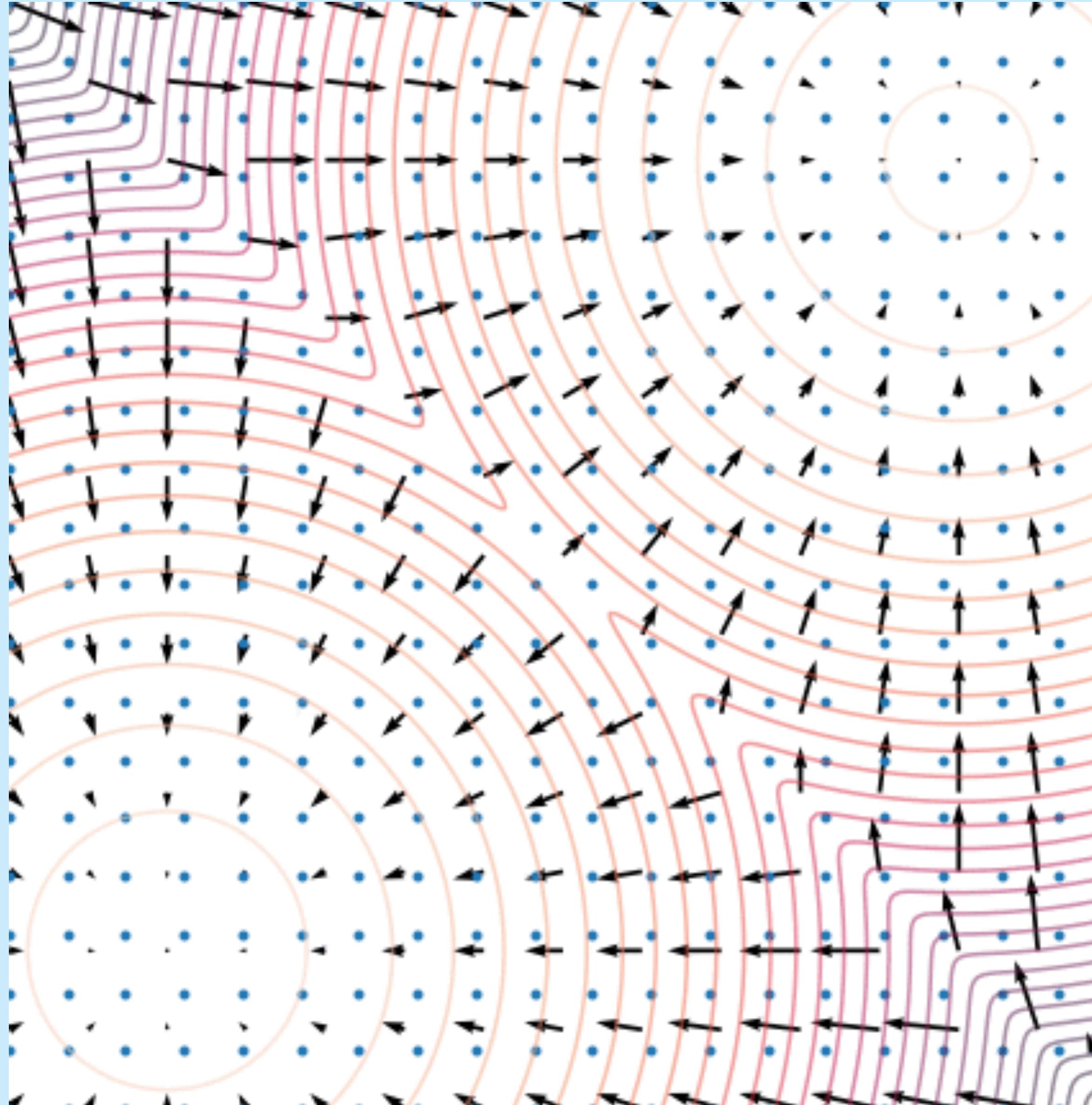
where  $z_i \sim \mathcal{N}(0, I)$

Learned from data

When  $\epsilon \rightarrow 0$  and  $K \rightarrow \infty$  the  $x_K$  from above converge to samples from  $p(x)$

Note: Langevin dynamics accesses  $p(x)$  only through  $\nabla_x \log p(x) \approx s_\theta(x)$

# LANGEVIN DYNAMICS IN ACTION



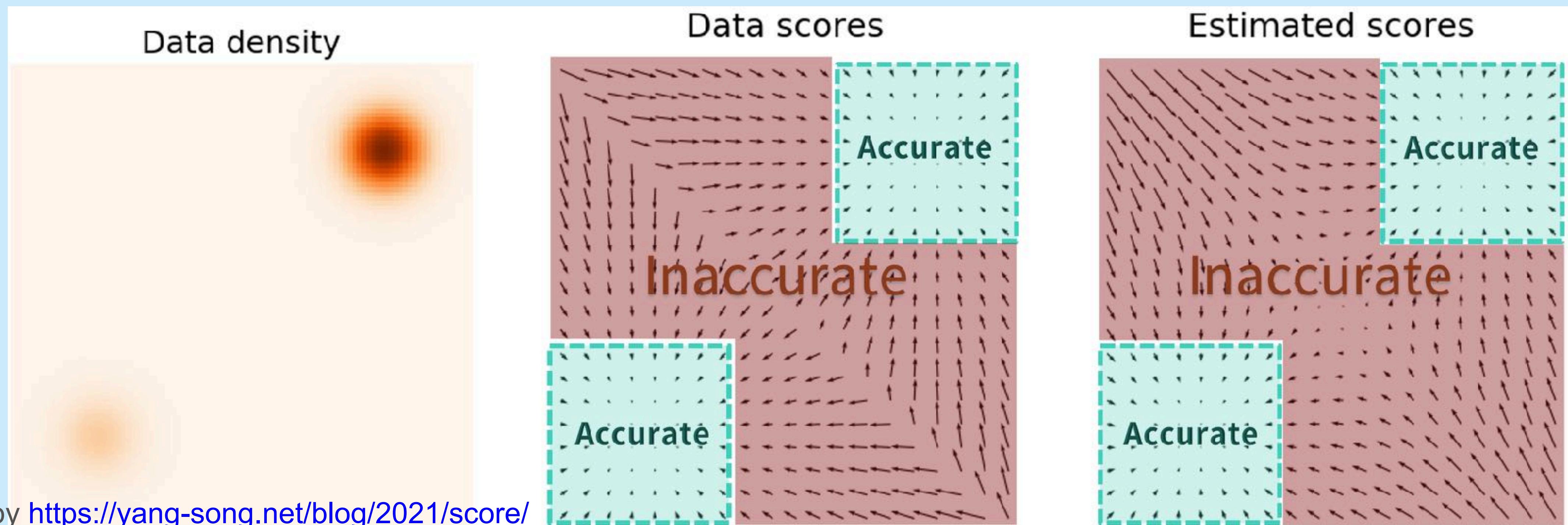
blog post by <https://yang-song.net/blog/2021/score/>

# LANGEVIN DYNAMICS PROBLEM SOLVED?

NO! Estimated score function only accurate for high density regions  
this is expected since score matching minimises the Fisher divergence

$$\mathbb{E}_{p(x)}[\|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - s_\theta(x)\|_2^2] = \int p(x) \|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - s_\theta(x)\|_2^2 dx$$

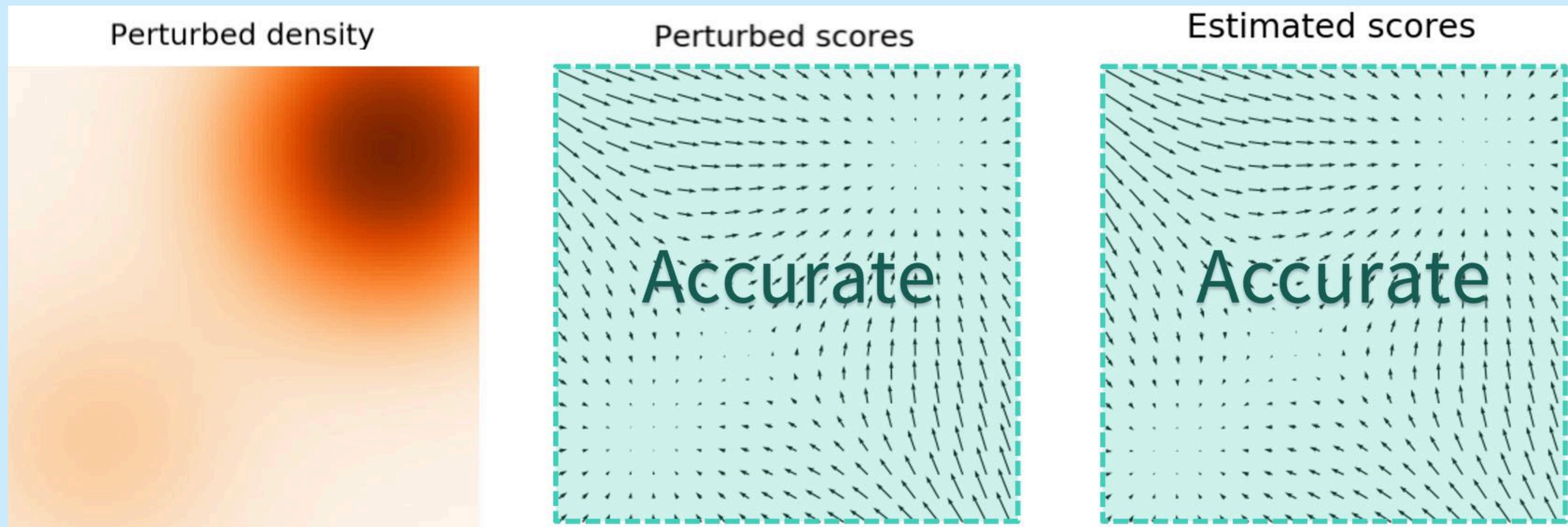
since here  $l_2$  norm is weighted by  $p(x)$  its value is largely ignored in low density regions



# SOLUTION: SCORE MATCHING WITH MULTIPLE NOISE PERTURBATIONS

Perturb the data with noise and train score model on noisy data

But: How to choose suitable noise perturbations?

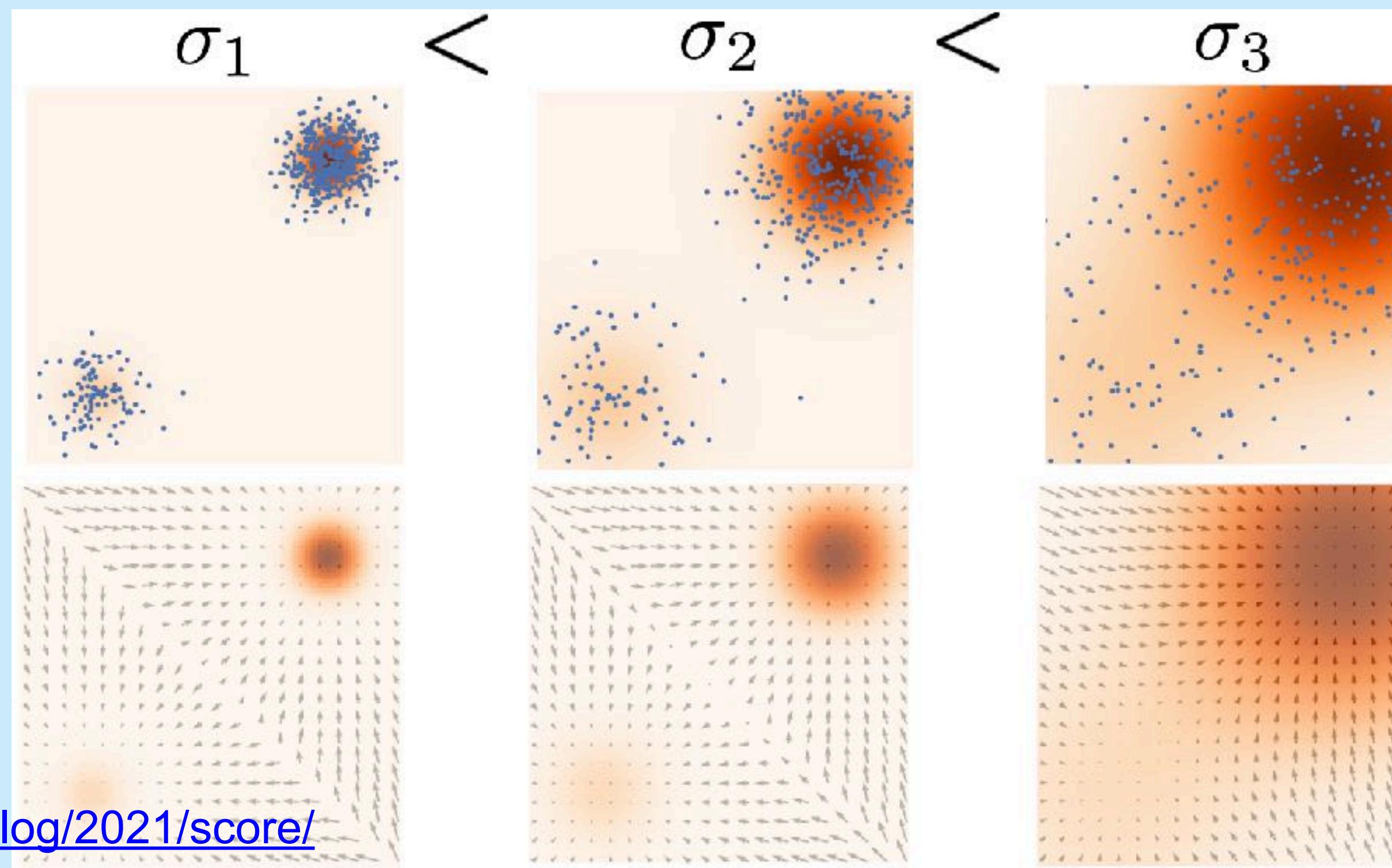


# SOLUTION: SCORE MATCHING WITH MULTIPLE NOISE PERTURBATIONS

Perturb the data with noise of increasing  $\sigma_1 < \sigma_2 < \sigma_3 \dots < \sigma_L$

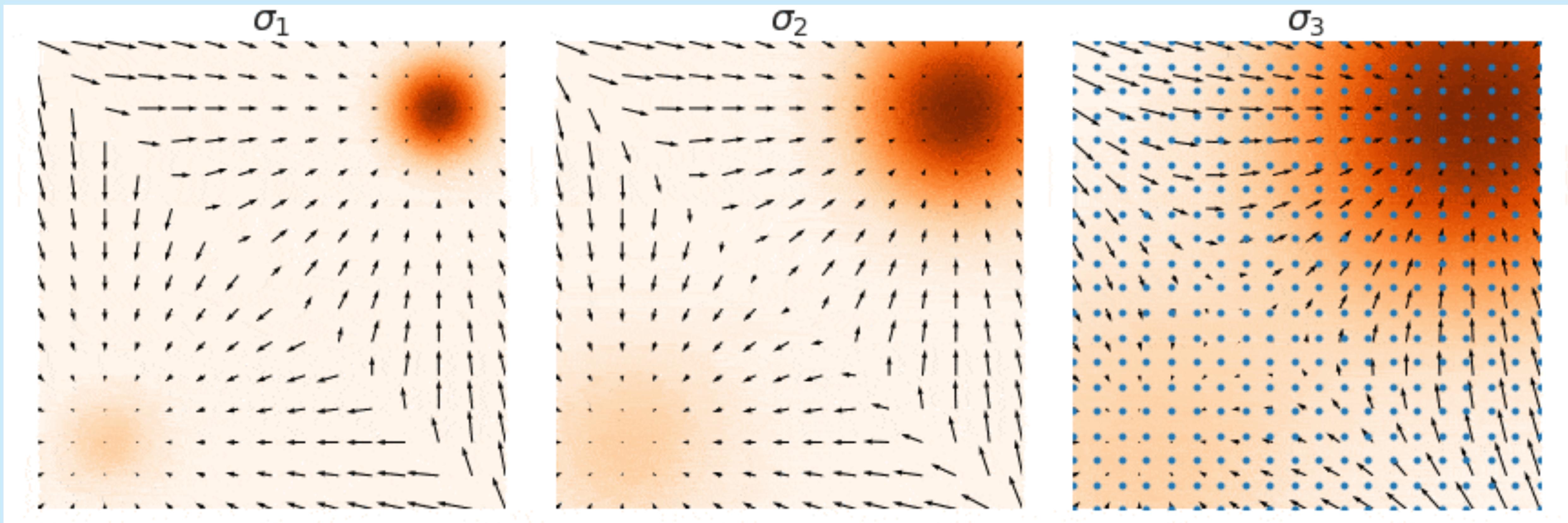
sample  $x \sim p_{\sigma_i}(x)$  by sampling  $x \sim p(x)$  and computing  $x + \sigma_i z$  with  $z \sim \mathcal{N}(0, I)$

then train **noise conditional score-based model**  $s_{\theta}(x, i)$



# SAMPLING VIA ANNEALED LANGEVIN DYNAMICS

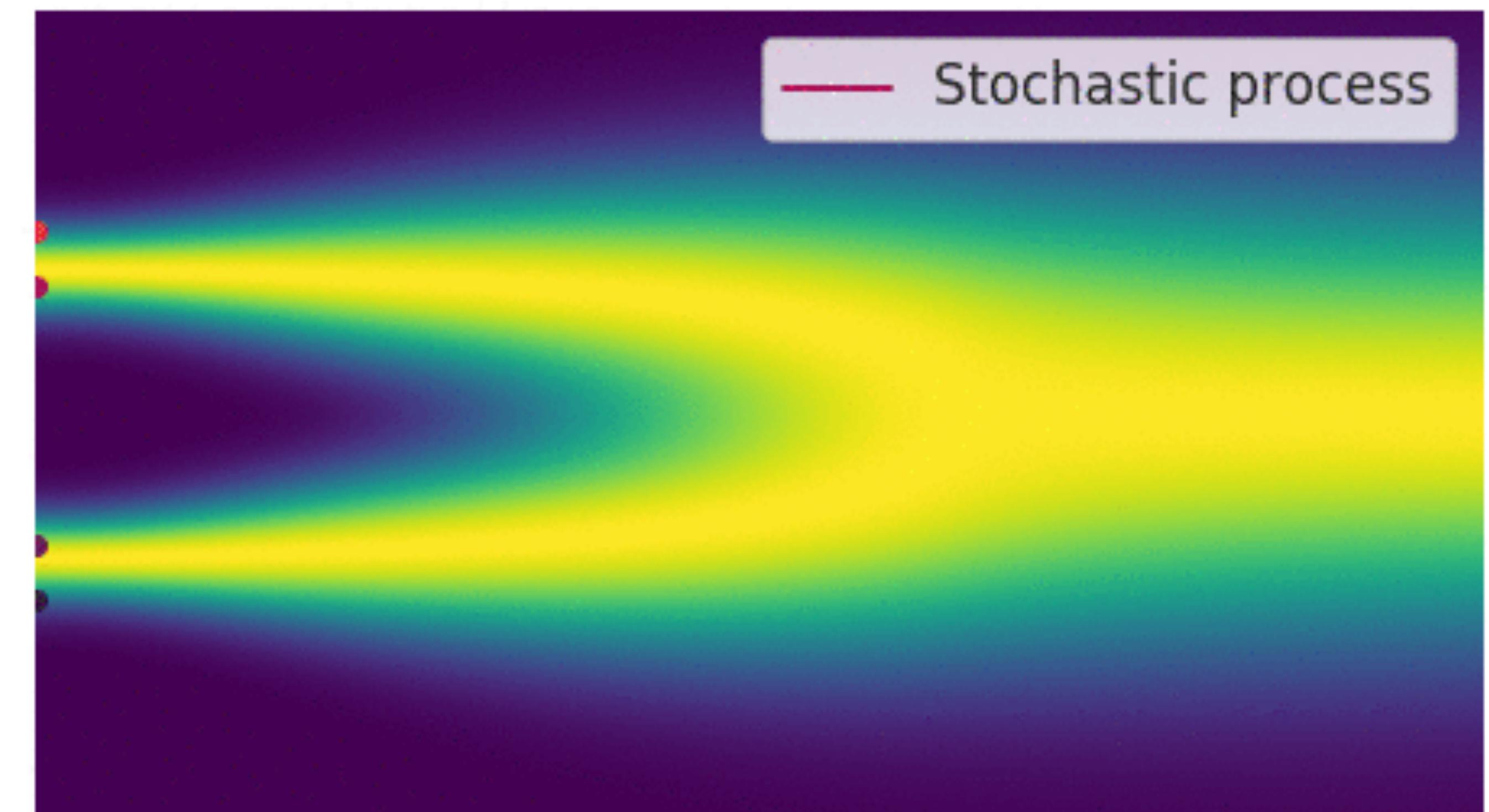
In the limit of infinite noise schedules  $\rightarrow$  higher quality samples, exact log-likelihood computation and more important: **controllable generation for inverse problem solving**



# SCORE-BASED GENERATIVE MODELING WITH STOCHASTIC DIFFERENTIAL EQUATIONS

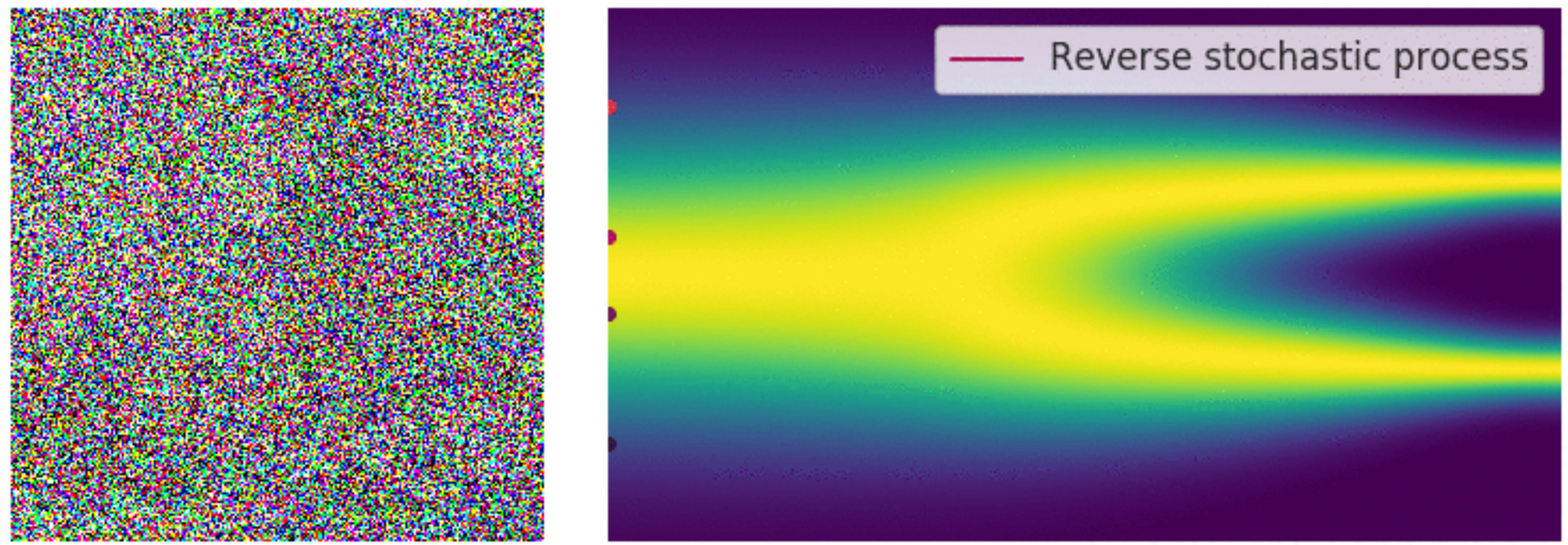
In the limit of infinite noise schedules → noise perturbation procedure is a continuous-time stochastic process

In general, the exact form of the SDE is part of the model, e.g. Variance Exploding SDE (VE SDE), the Variance Preserving SDE (VP SDE), and the sub-VP SDE (Song+2021)



# HOW TO SAMPLE WITH SCORE-BASED MODELS BASED ON SDES?

Note! any SDE has a corresponding reverse SDE!



# THE CONCEPT BEHIND SDE BASED SCORE MATCHING

Sohl-Dickstein et al. (2015), Ho et al. (2020), Song et al. (2021)

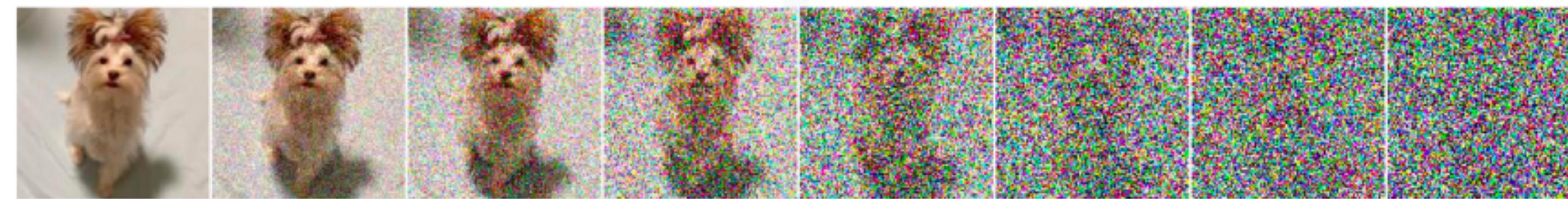
**Forward SDE:**

$$d\mathbf{z}_t = -\frac{1}{2}\beta(t)\mathbf{z}_t dt + \sqrt{\beta(t)} d\mathbf{w}_t$$

Easy

$t = 0$

$\mathbf{z}_0 \sim p_{\text{target}}$



$t = 1$

$\mathbf{z}_1 \sim \mathcal{N}(0, I_d)$

**Backward SDE:**

$$d\mathbf{z}_t = \left[ -\frac{1}{2}\beta(t)\mathbf{z}_t - \underbrace{\beta(t) \nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t)}_{\text{Score function}} \right] dt + \sqrt{\beta(t)} d\bar{\mathbf{w}}_t$$

Hard

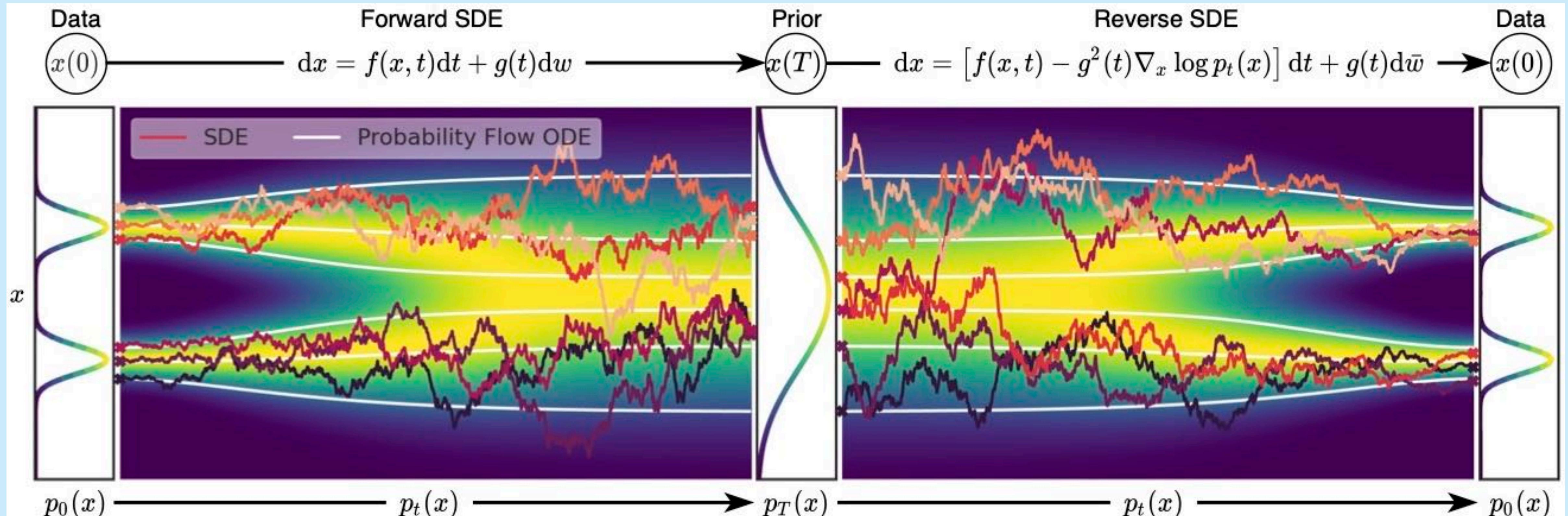
Score function

Solve with your favourite SDE solver, e.g. Euler-Maruyama, stochastic Runge-Kutta, etc.  
Jolicœur-Martineau+2021 for adaptive step size SDE solver

# SIDE NOTE: PROBABILITY FLOW ODES

Any SDE can be converted into an ordinary differential equation (ODE) without changing its marginal distribution

$$dx = \left[ f(x, t) - \frac{1}{2} g^2(t) \nabla_x \log p_t(x) \right] dt$$



# SIDE NOTE: PROBABILITY FLOW ODES

Any SDE can be converted into an ordinary differential equation (ODE) without changing its marginal distribution

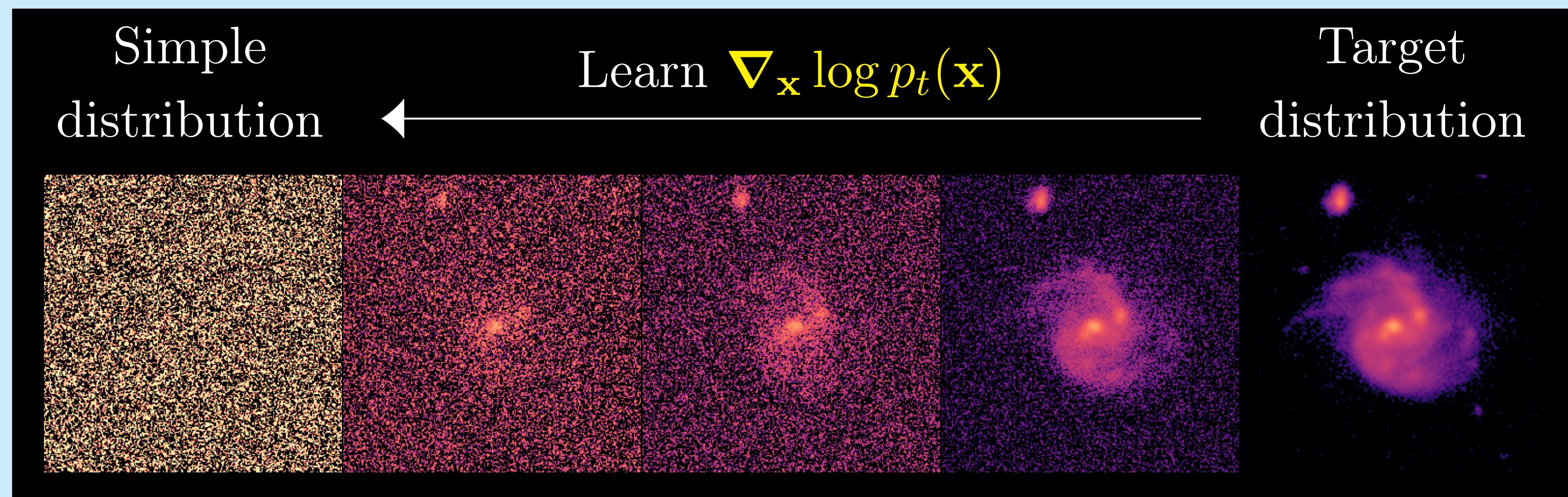
$$dx = \left[ f(x, t) - \frac{1}{2} g^2(t) \nabla_x \log p_t(x) \right] dt$$

Why is this interesting?

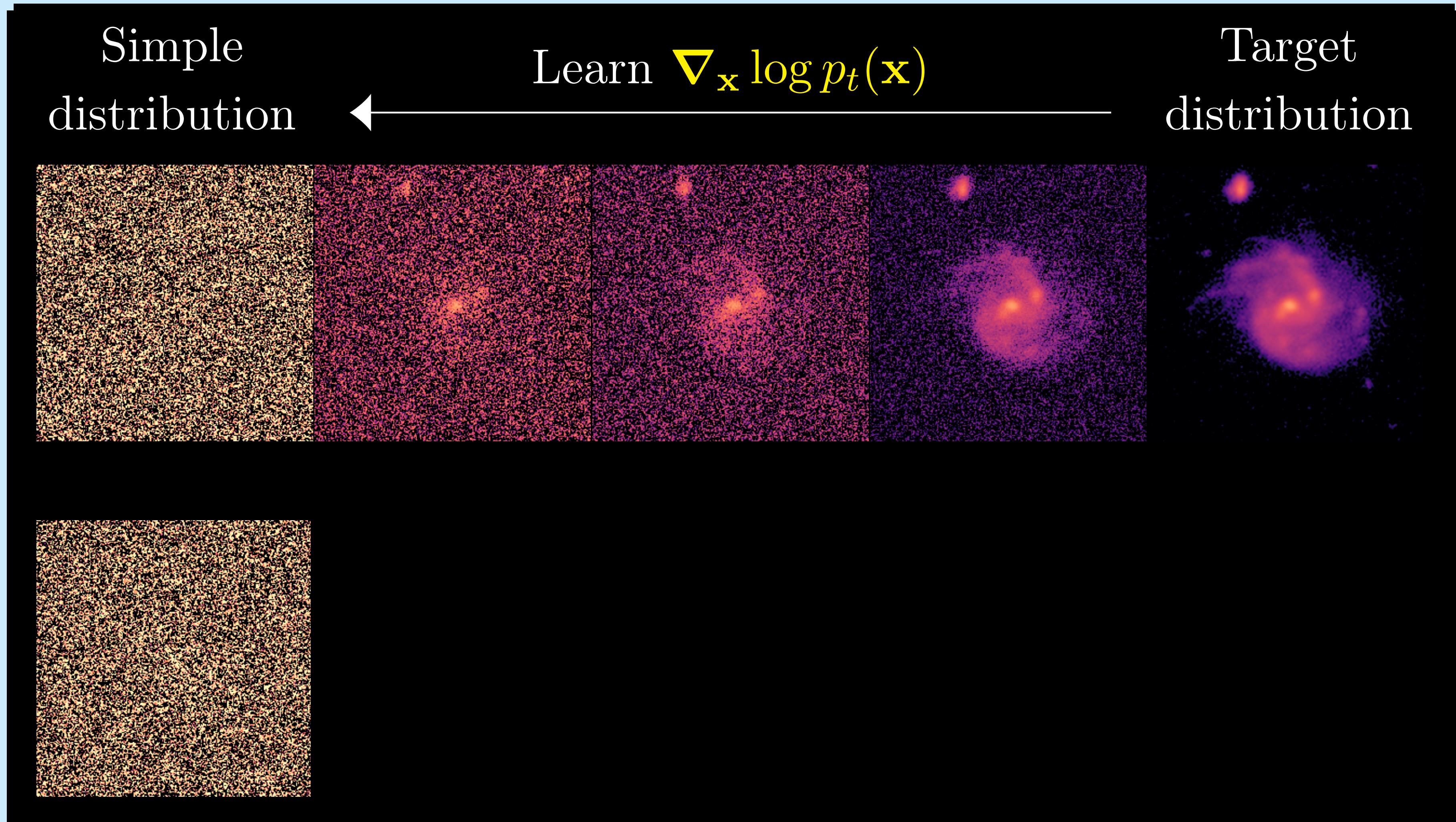
when replacing  $\nabla_x \log p_t(x)$  with approximation  $s_\theta(x, t)$  the probability flow ODE becomes a special case of a neural ODE (Chen+2018)

in particular its an example of a continuous normalizing flow (Gerathwohl+2019)

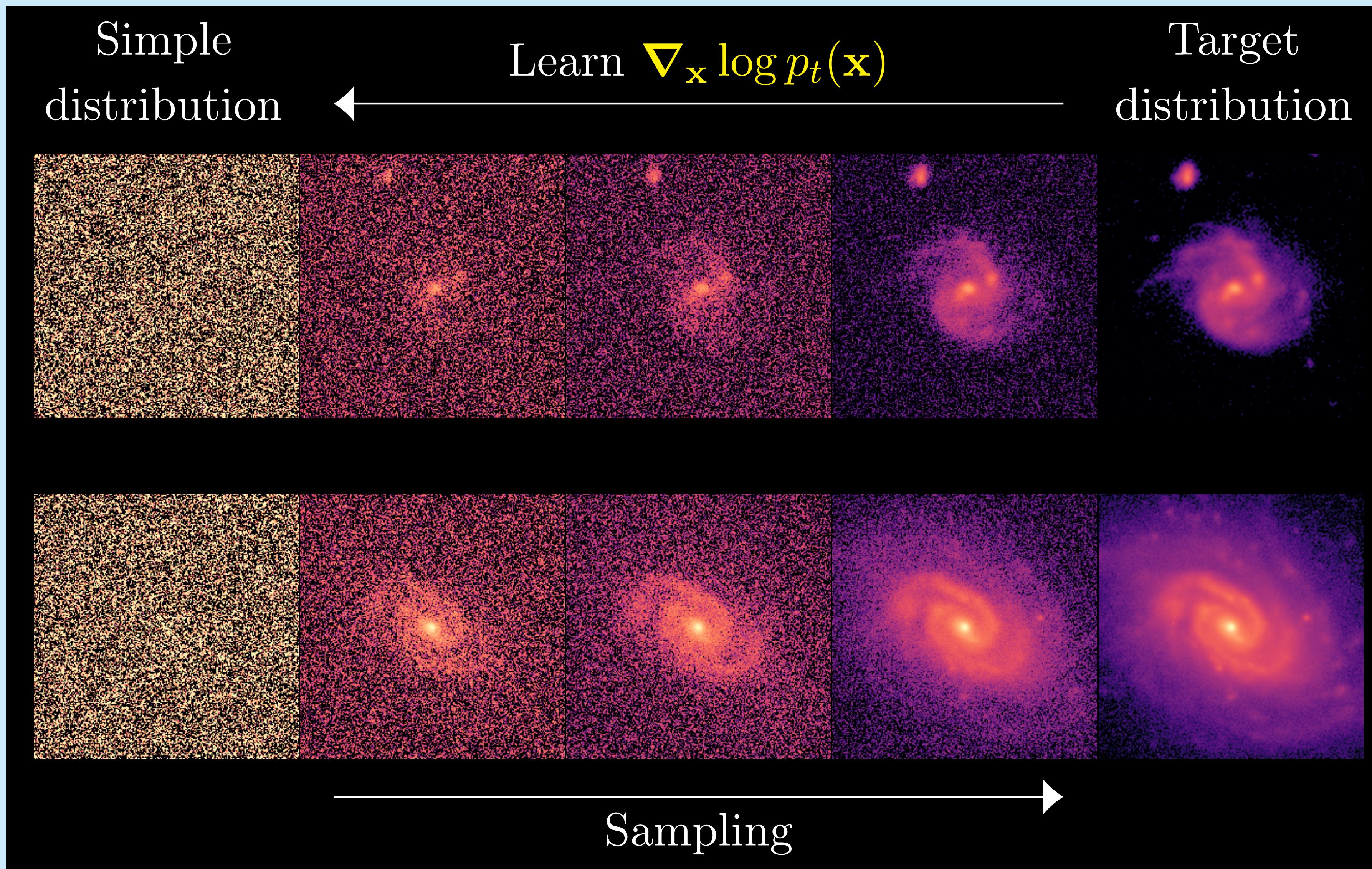
# TRAINING SCORE MATCHING MODELS



# TRAINING SCORE MATCHING MODELS



# TRAINING SCORE MATCHING MODELS

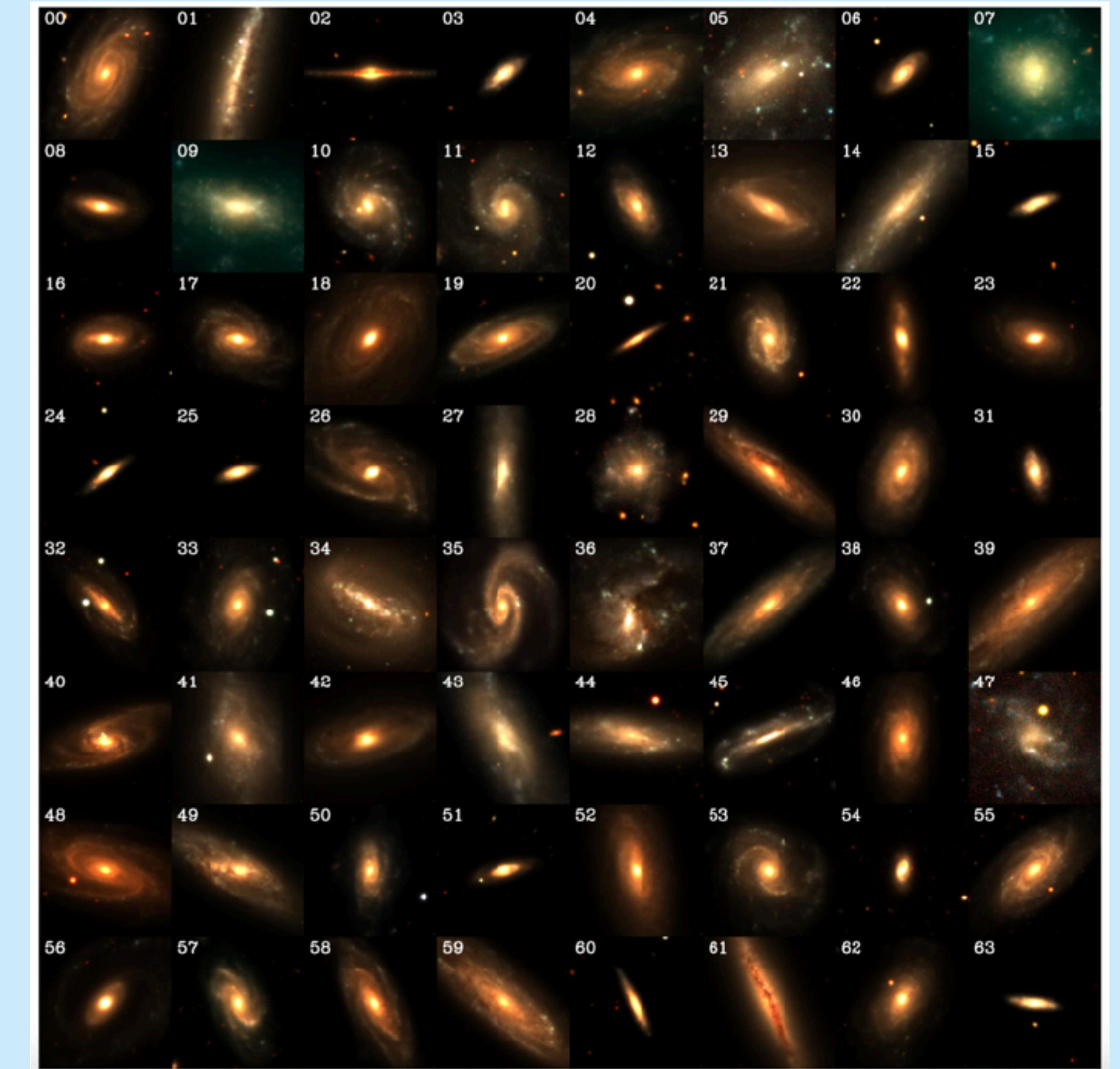


# Applications to Astrophysics

# APPLICATION OF SCORE MATCHING AND DIFFUSION MODELS

Learn generative model purely from data!

Smith+2021



# DIFFUSION MODELS AS A MAPPING $p_{\text{dust}} \leftrightarrow p_{\text{CMB}}$

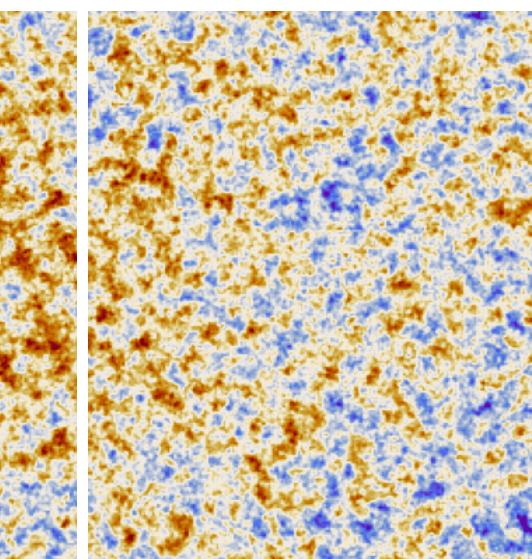
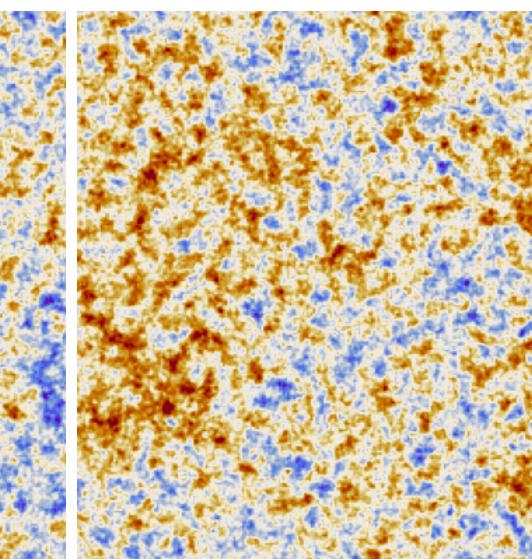
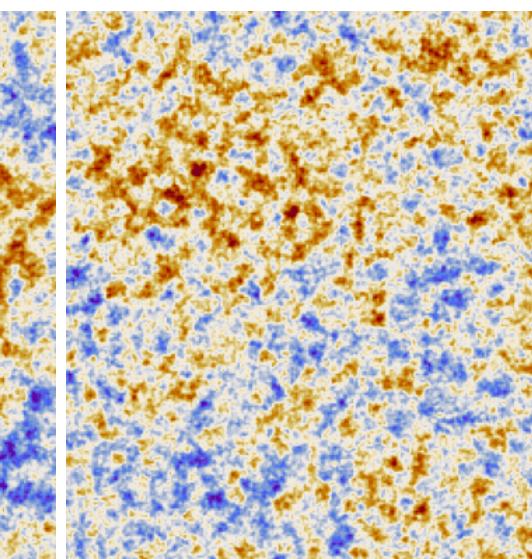
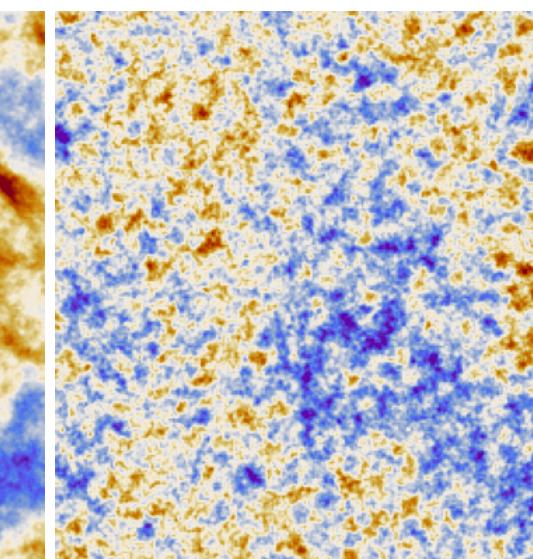
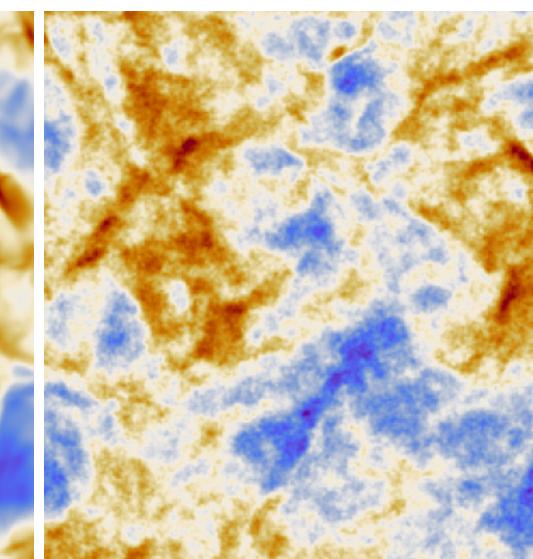
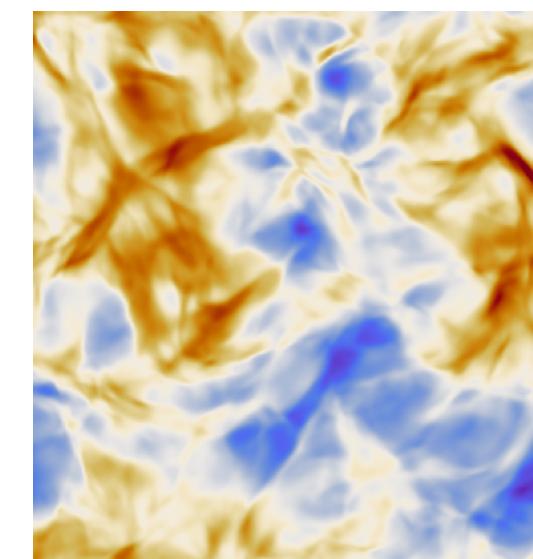
**Forward SDE:**

$$dz_t = -\frac{1}{2}\beta(t)z_t dt + \sqrt{\beta(t)}\Sigma_{\phi}^{\frac{1}{2}} d\mathbf{w}_t$$

Easy

$t = 0$

$\mathbf{z}_0 \sim p_{\text{dust}}$



$t = 1$

$\mathbf{z}_1 \sim p_{\text{CMB}}$

$\mathcal{T} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$

**Backward SDE:**

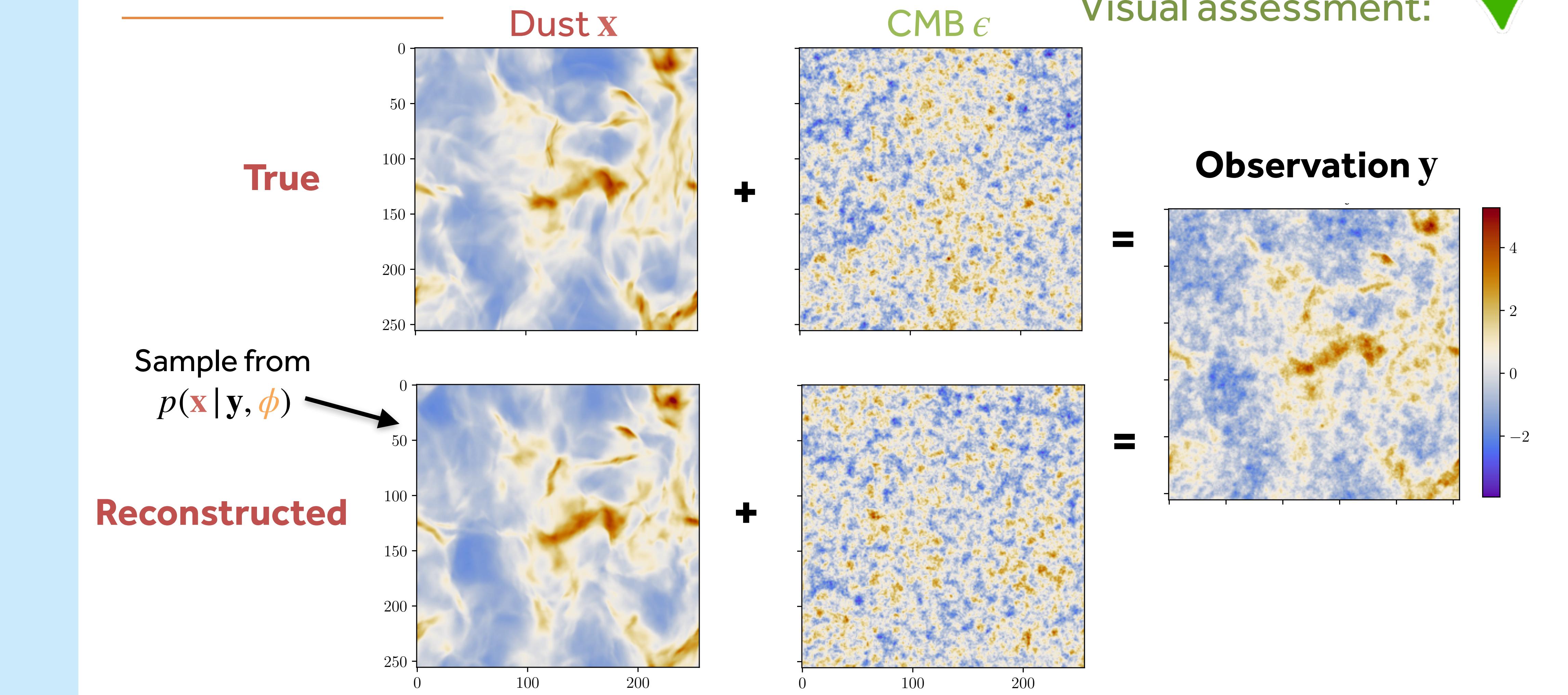
$$dz_t = \left[ -\frac{1}{2}\beta(t)z_t - \beta(t)\Sigma_{\phi} \underbrace{\nabla_{\mathbf{z}_t} \log p_t(\mathbf{z}_t)}_{\text{Score function}} \right] dt + \sqrt{\beta(t)}\Sigma_{\phi}^{\frac{1}{2}} d\bar{\mathbf{w}}_t$$

Hard

Score function

# DIFFUSION MODELS AS A MAPPING $p_{\text{dust}} \leftrightarrow p_{\text{CMB}}$

Component Separation ?



# POSTERIOR SAMPLES WITH DIFFUSION MODELS

Posterior  $p(x | y)$  with observation  $y$  is given by Bayes' theorem:

$$\log p(x | y) = \log p(y | x) + \log p(x) - \log p(y)$$

with  $p(y | x)$  being the likelihood  
and  $p(x)$  the prior.

hence the score is given by:

$$\nabla_x \log p(x | y) = \nabla_x \log p(y | x) + \nabla_x \log p(x) - \nabla_x \log p(y)$$

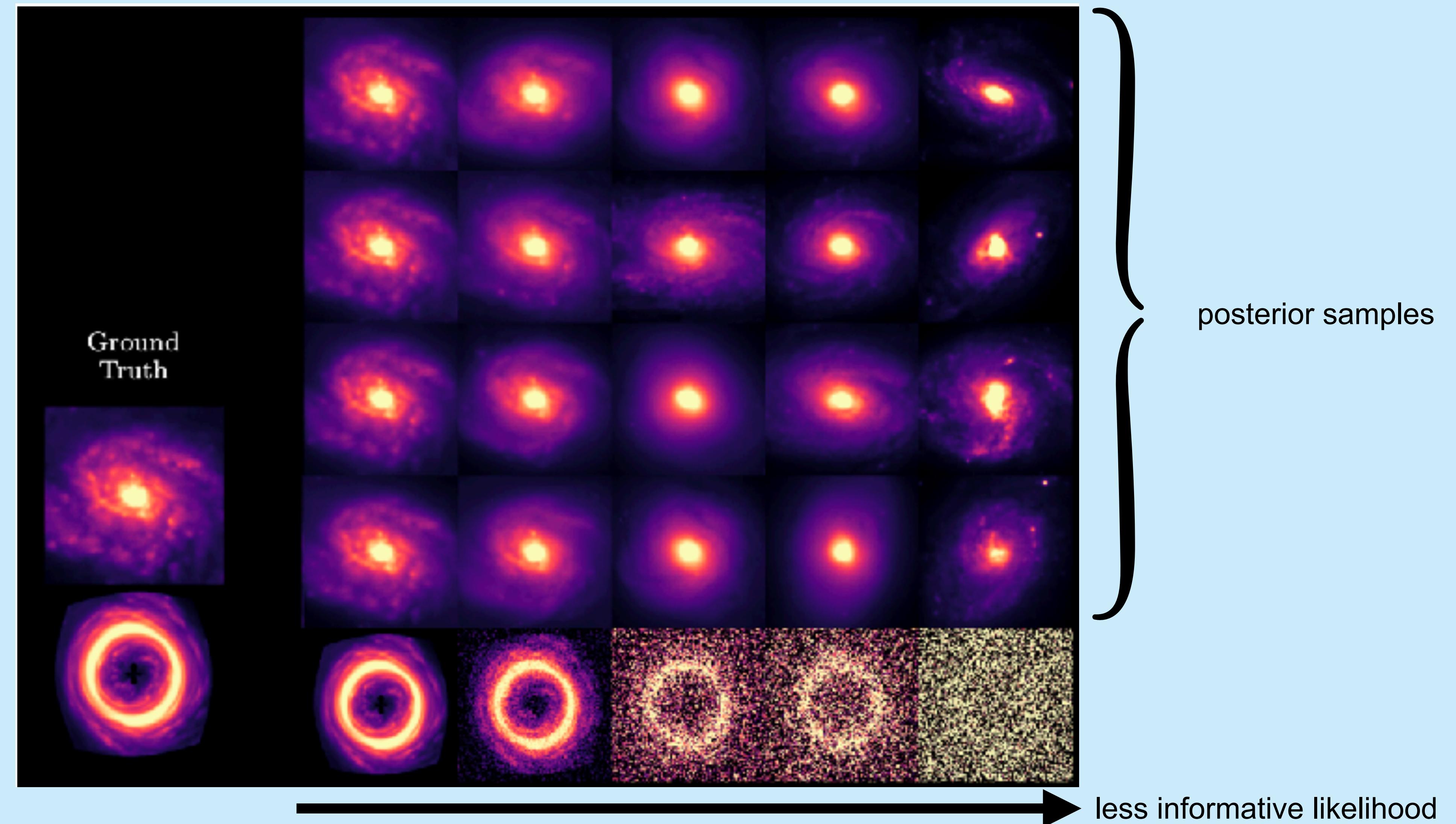
To a good approximation, we can calculate the likelihood score analytically or numerically from the known forward process.

This is the score we learnt with Score Matching from data!

=0 because it does not depend on x.

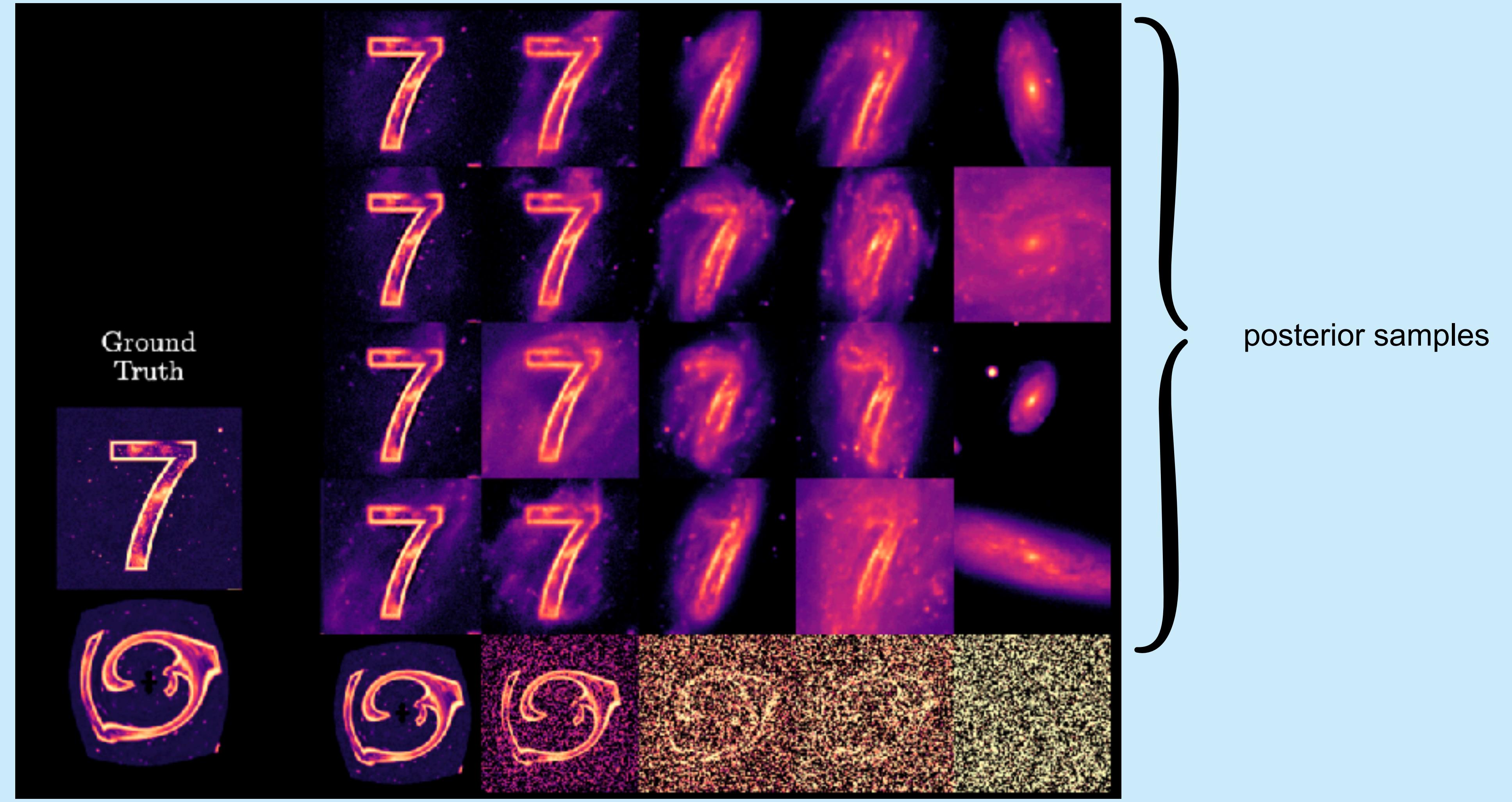
# POSTERIOR SAMPLES WITH DIFFUSION MODELS

Posterior samples of source galaxies in strong gravitational lenses with score-based priors



# POSTERIOR SAMPLES WITH DIFFUSION MODELS

Posterior samples for out-of-distribution galaxies

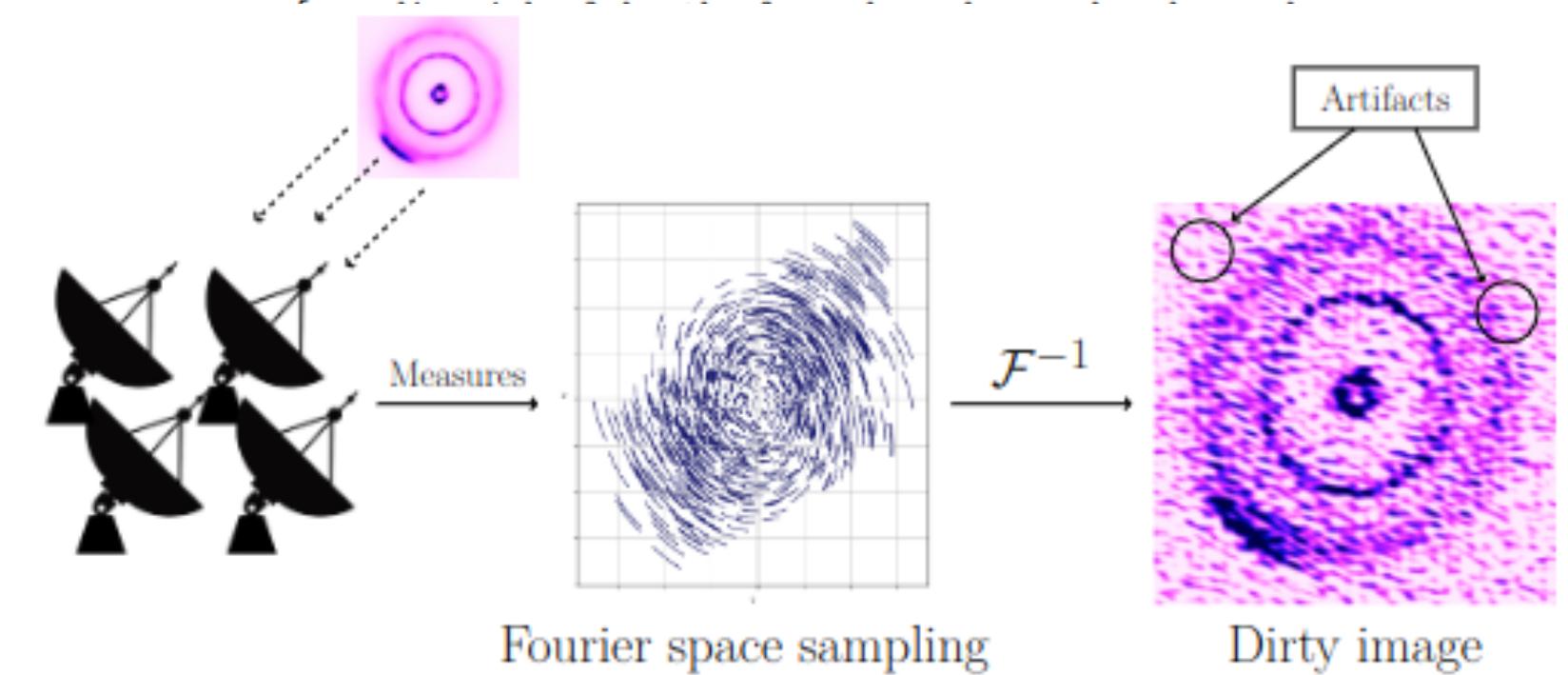


# MORE APPLICATIONS OF DIFFUSION MODELS

## Bayesian Imaging for Radio Interferometry with Score-Based Priors

Noé Dia<sup>1,2,4</sup>  M. J. Yantovski-Barth<sup>1,2,4</sup>  Alexandre Adam<sup>1,2,4</sup>  Micah Bowles<sup>5</sup>   
Pablo Lemos<sup>1,2,3,4</sup> Anna M. M. Scaife<sup>5,6</sup> Yashar Hezaveh<sup>1,2,3,4,7,8</sup>  
Laurence Perreault-Levasseur<sup>1,2,3,4,7,8</sup>

<sup>1</sup>Université de Montréal <sup>2</sup>Cielo Institute <sup>3</sup>Flatiron Institute <sup>4</sup>Mila <sup>5</sup>University of Manchester  
<sup>6</sup>The Alan Turing Institute <sup>7</sup>Trottier Space Institute <sup>8</sup>Perimeter Institute



## THE ASTROPHYSICAL JOURNAL LETTERS

### OPEN ACCESS

#### Beyond Gaussian Noise: A Generalized Approach to Likelihood Analysis with Non-Gaussian Noise

Ronan Legin<sup>5,1,2,3</sup>  Alexandre Adam<sup>5,1,2,3</sup>  Yashar Hezaveh<sup>1,2,3,4</sup> 

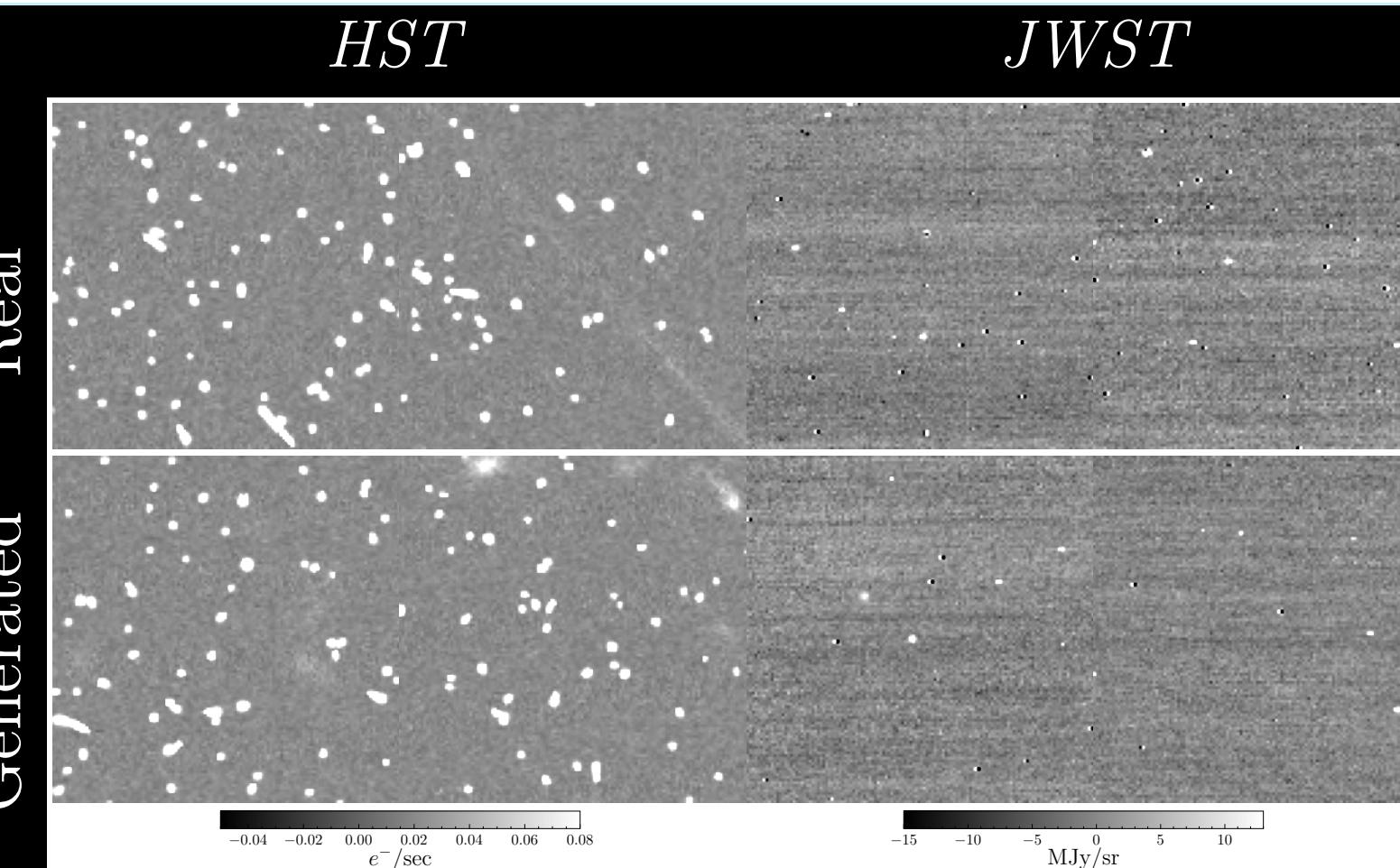
Laurence Perreault-Levasseur<sup>1,2,3,4</sup> 

Published 2023 June 6 • © 2023. The Author(s). Published by the American Astronomical Society.

[The Astrophysical Journal Letters, Volume 949, Number 2](#)

Citation Ronan Legin et al 2023 ApJL 949 L41

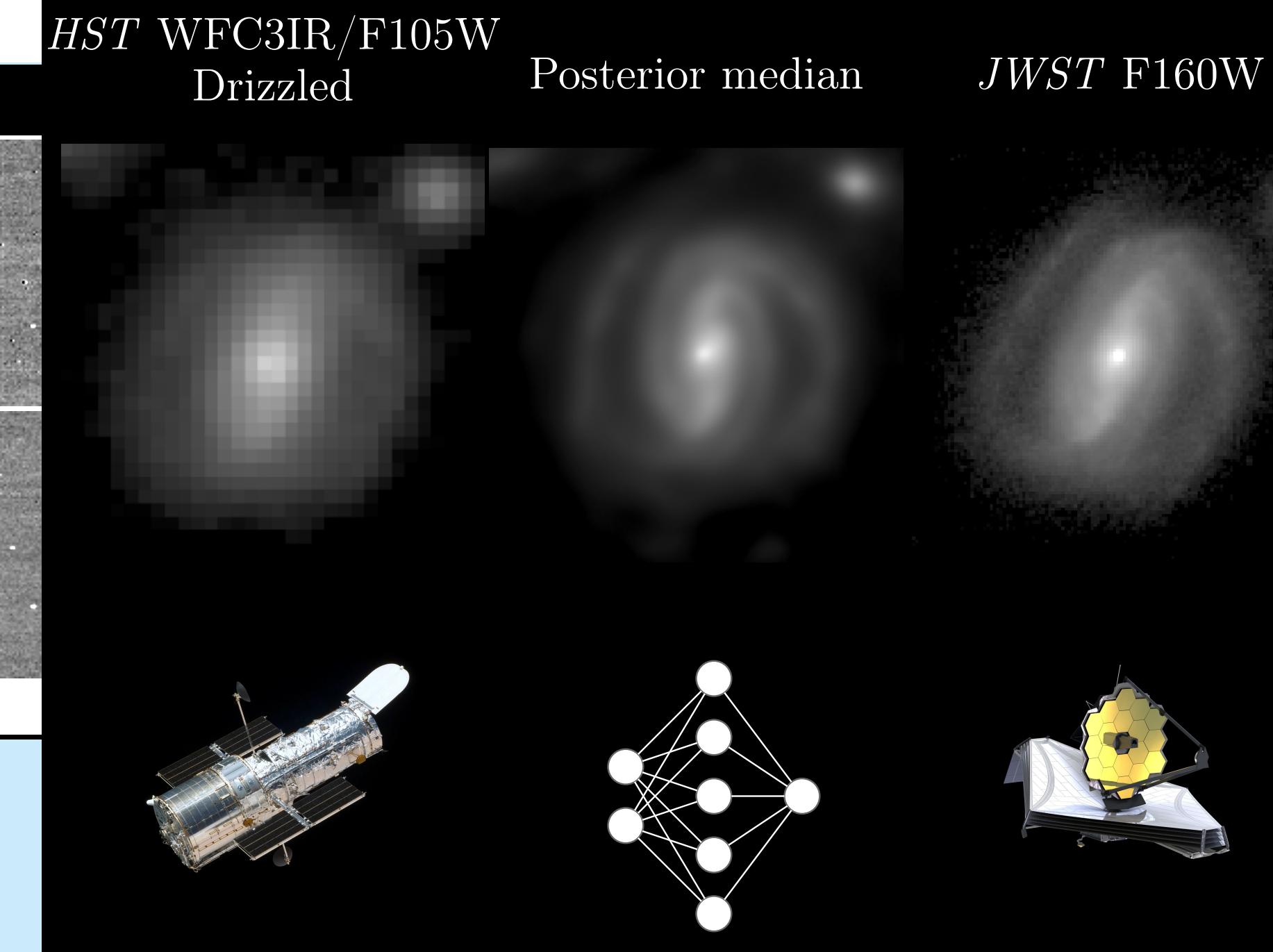
DOI 10.3847/2041-8213/acd645



#### Echoes in the Noise: Posterior Samples of Faint Galaxy Surface Brightness Profiles with Score-Based Likelihoods and Priors

Alexandre Adam<sup>1,2,4</sup> Connor Stone<sup>1,2,4</sup> Connor Bottrell<sup>5,6</sup> Ronan Legin<sup>1,2,4</sup>  
Yashar Hezaveh<sup>1,2,3,4,7,8</sup> Laurence Perreault-Levasseur<sup>1,2,3,4,7,8</sup>

<sup>1</sup>Université de Montréal <sup>2</sup>Cielo Institute <sup>3</sup>CCA, Flatiron Institute <sup>4</sup>Mila  
<sup>5</sup>ICRAR <sup>6</sup>Kavli IPMU <sup>7</sup>Trottier Space Institute <sup>8</sup>Perimeter Institute



# WHERE TO FIND THE LECTURE MATERIAL



<https://github.com/TobiBu/graddays>

# SUMMARY & CONCLUSION

Main take away:

scientific motivated inductive bias helps to be more robust, more data efficient  
and better interpretable

My personal message:

Write better code!  
Share more data!  
Build more open-source software!

This will accelerate research cycles and lets you engage with peers early on!

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG  
PILE OF LINEAR ALGEBRA, THEN COLLECT  
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL  
THEY START LOOKING RIGHT.

