# An introduction to generative modeling

## *with applications in physics*

smsharma.io/iaifi-summer-school-2023

## Siddharth Mishra-Sharma

(smsharma.io/@kdqg1)

NSF AI Institute for Artificial Intelligence
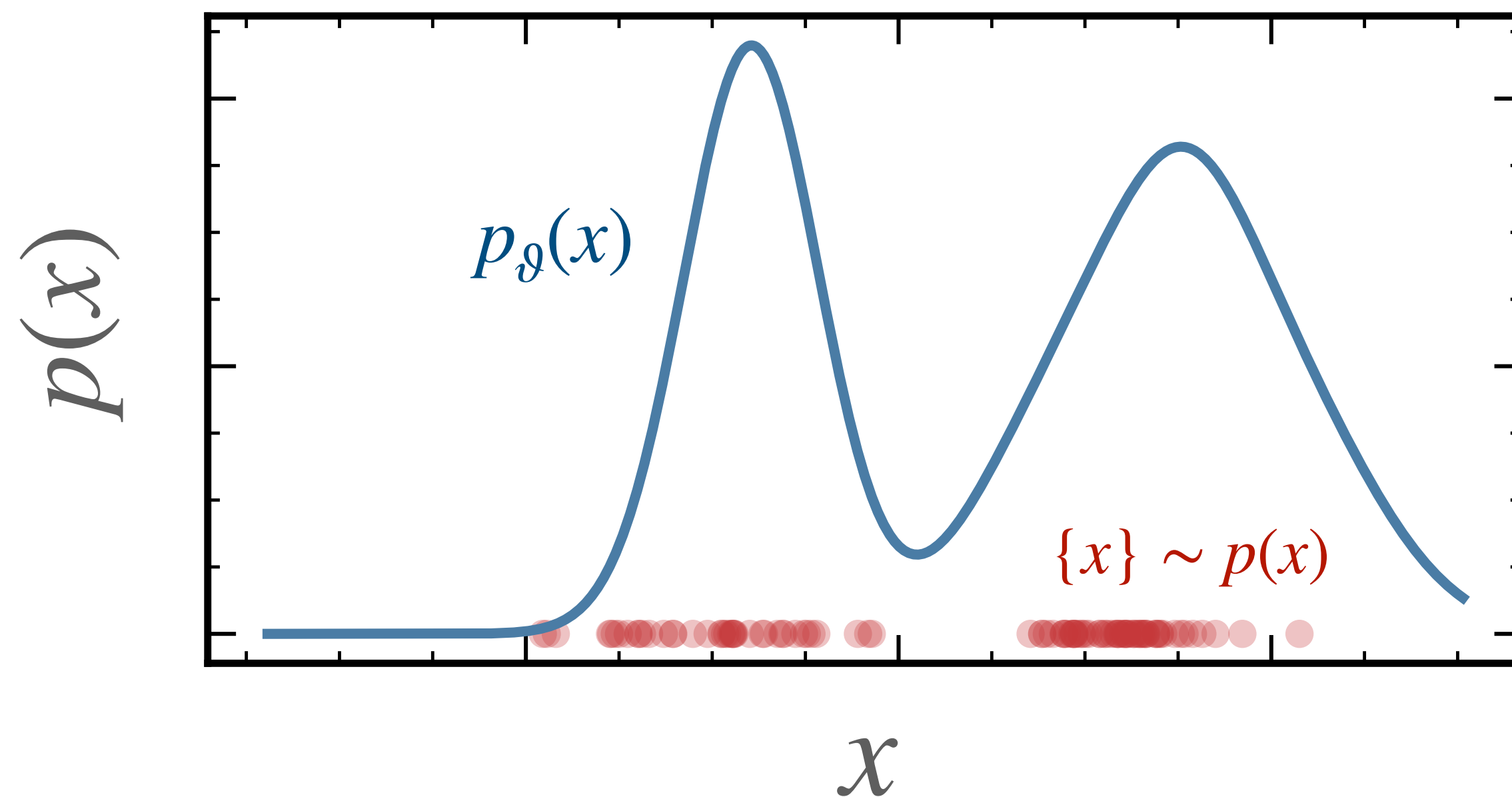and Fundamental Interactions (IAIFI)

IAIFI Summer School

August 8, 2023

# Generative models

*Generative models are simulators of the data*

Goal: learn a probability distribution $p_\vartheta(x)$ that is as close as possible to the true underlying data distribution $p(x)$



$p_\vartheta(x)$

$\{x\} \sim p(x)$

*1. Sampling*

$$x \sim p_\vartheta(x)$$

*2. Density estimation*

$$\log p_\vartheta(x)$$

# Evolution of deep generative models
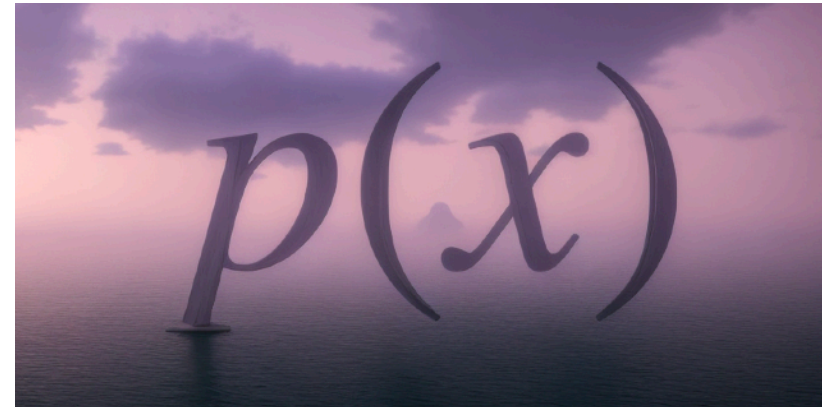


Variational autoencoders
(from Kingma et al 2013)



Diffusion models
(*Midjourney* 2023)
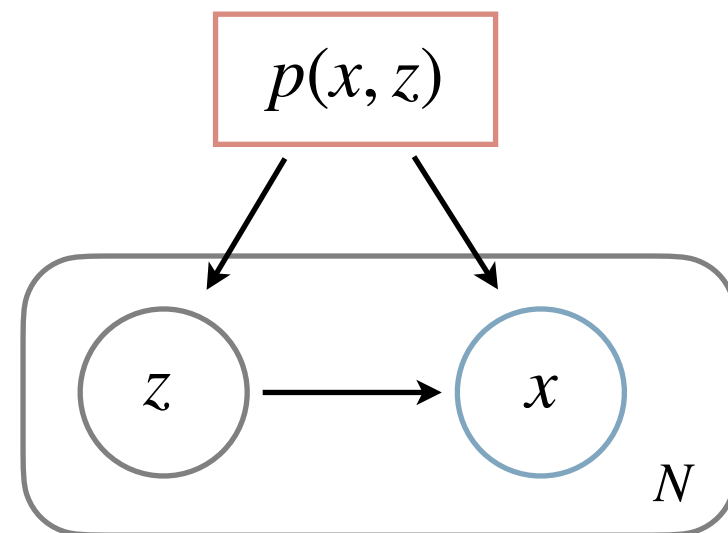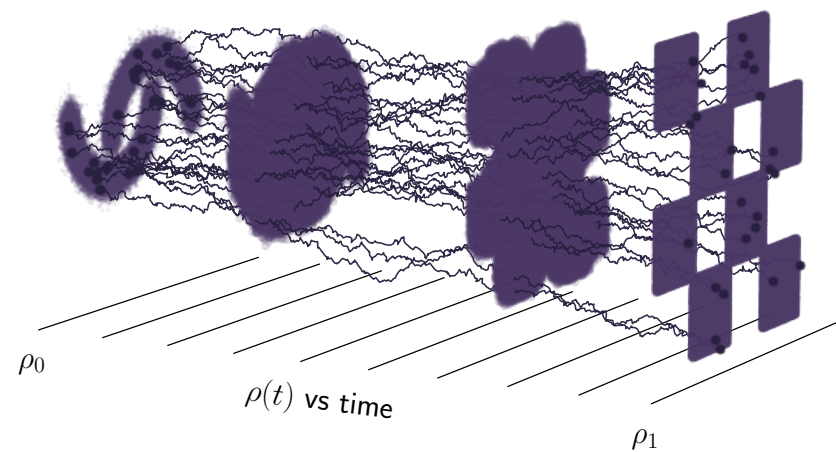
The landscape of deep generative models

[Karsten Kreis; CVPR 2022 Tutorial]

Energy-based models

Diffusion models

Variational autoencoders

Normalizing flows

Autoregressive models

Generative adversarial networks

# Outline



## Why (deep) generative modeling?

*What is it, and what can it do for you?*



## Variational auto encoders

*Latent-variable modeling, and compression is all you need*



## Diffusion models

*Models based on iterative refinement*

## Normalizing flows

*Invertible transformations*

# *Simulators*

$$x \sim p(x)$$

Simulators are ubiquitous: *they prescribe a way to sample from the data distribution*
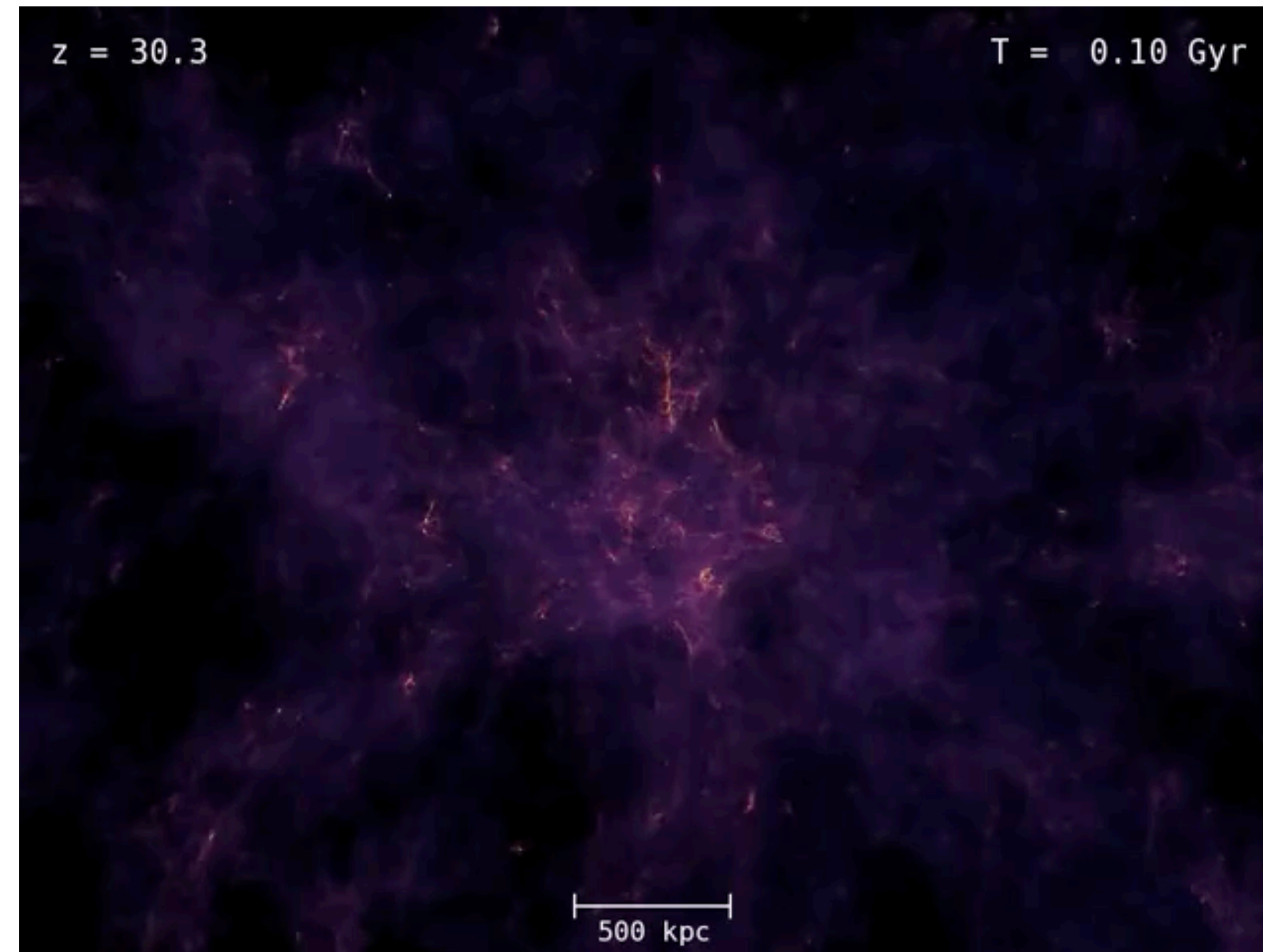
## Collider data

particles $\sim p$(particles)
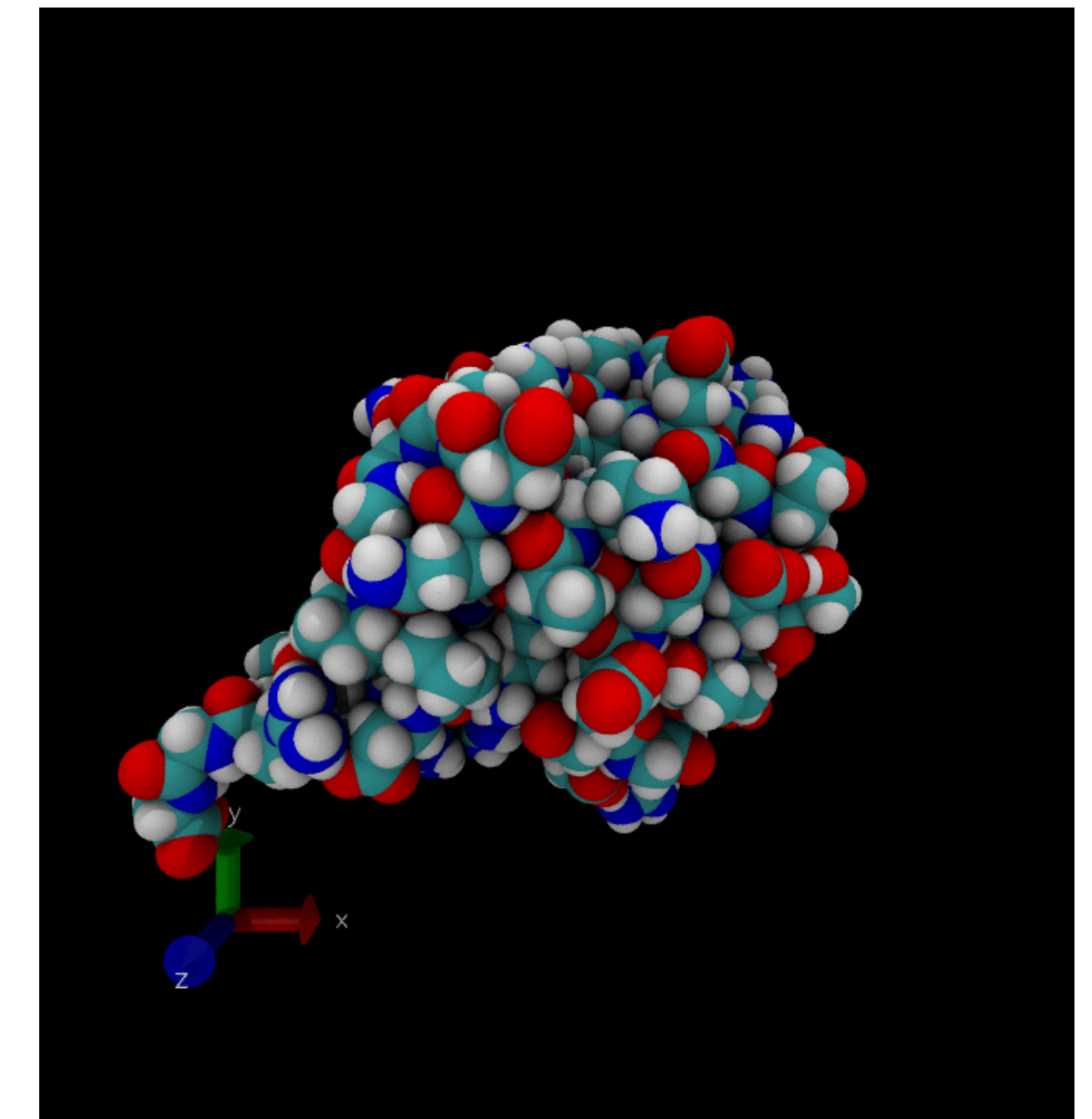


[C. Cesarotti with ATLAS]

## Cosmology data

particles $\sim p$(particles)



[Aquarius simulation]

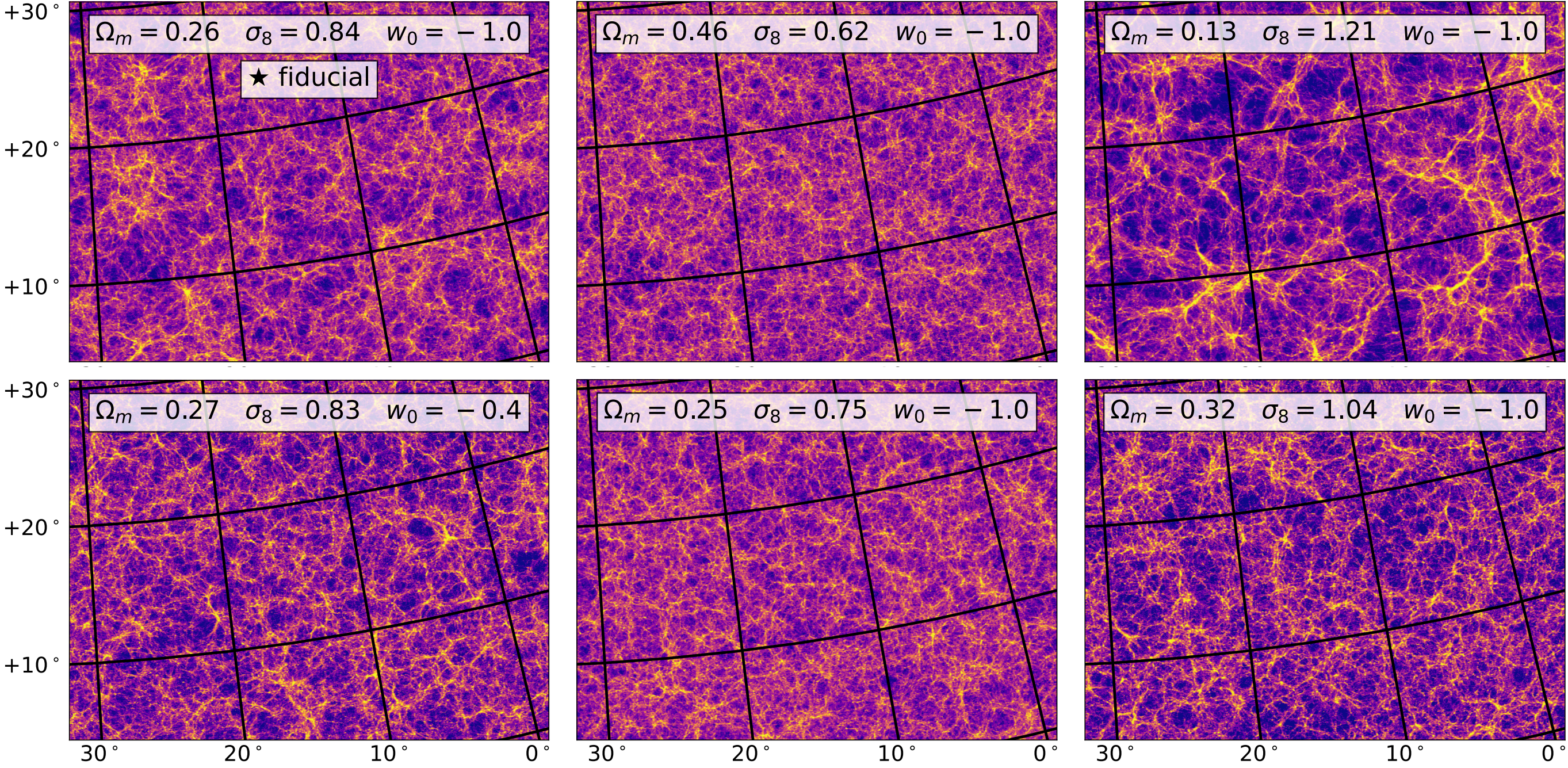## Molecular dynamics

configurations $\sim p$(configurations)



[E. Cances et al]

# *Conditional* simulators

Conditional simulations *sample from the likelihood $p(x \mid \theta)$*

## Cosmology data

$$\text{map} \sim p(\text{map} \mid \{\Omega_m, \sigma_8, w_0\})$$



[Kacprzak et al 2022]

$$x \sim p(x; \mathcal{M})$$
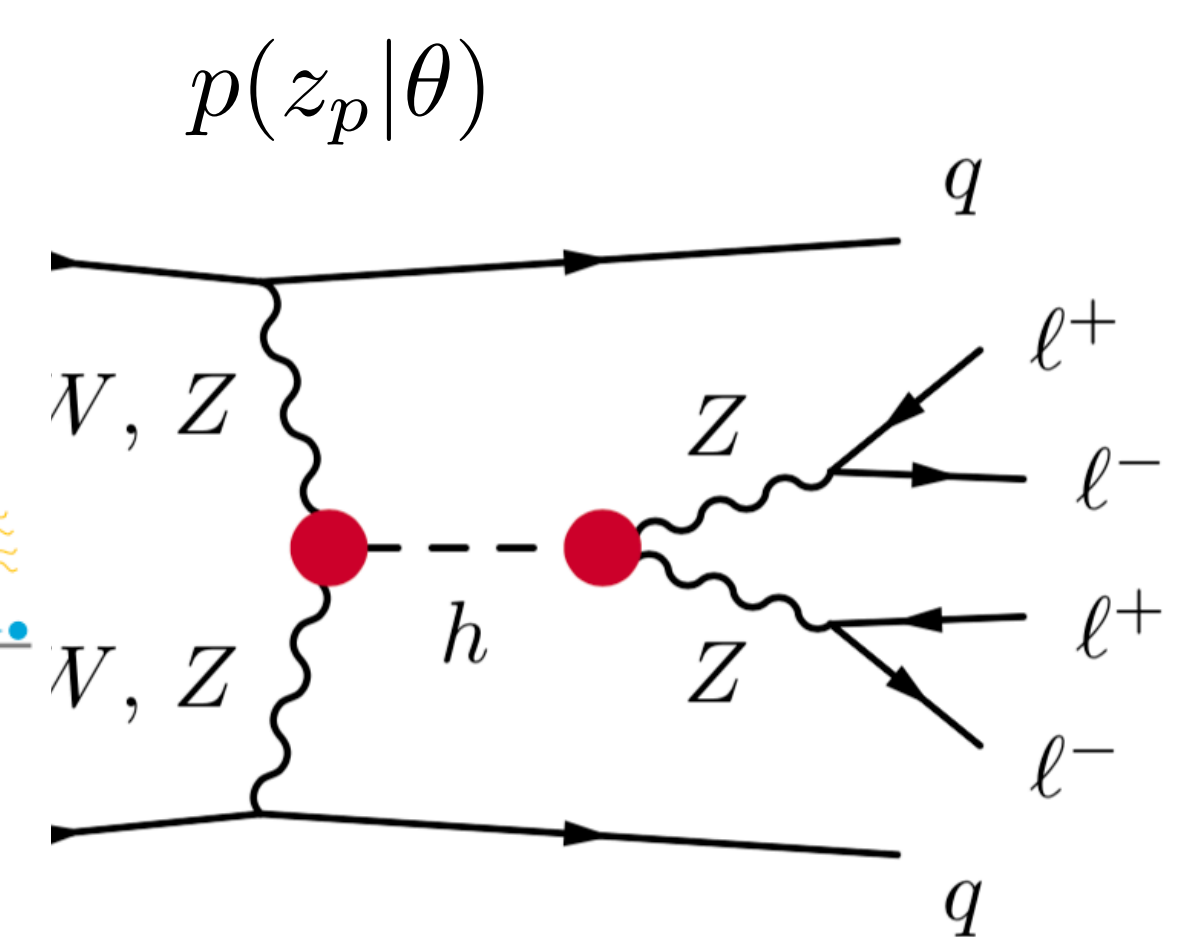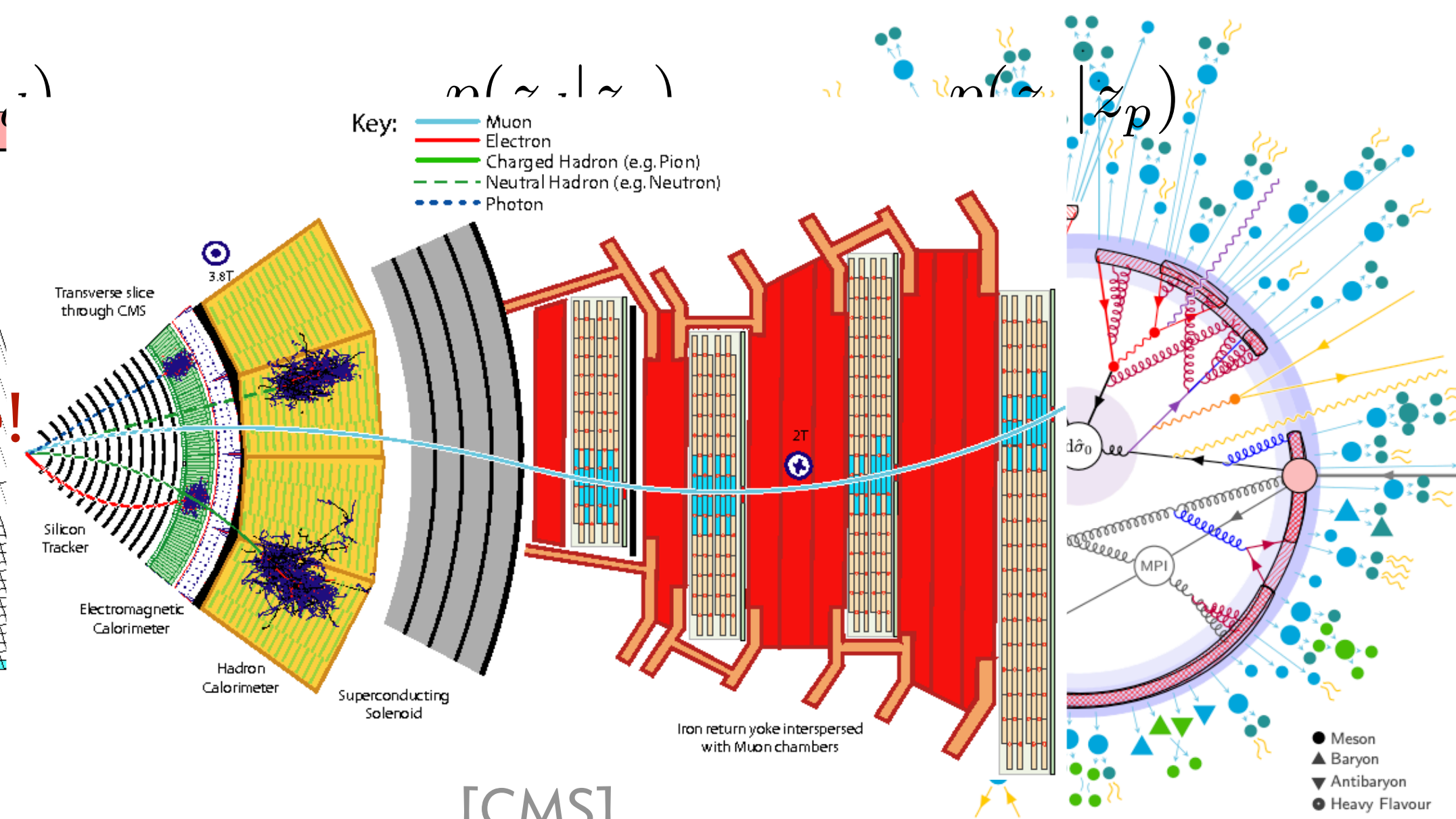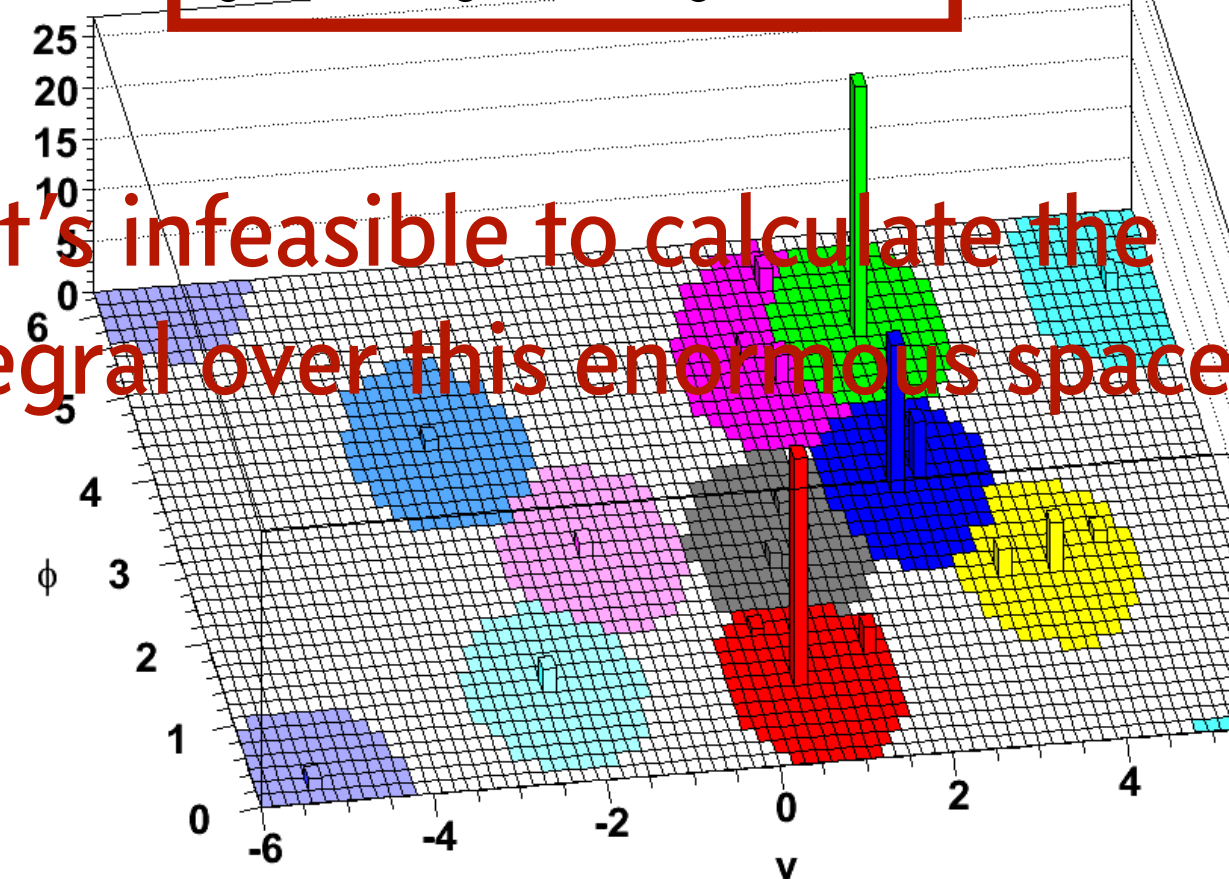
Model

*or*

$$x \sim p(x \mid \theta)$$

Model parameters

# Are simulators all you need?

*Latent variables*

| Observables | Detector interactions | Shower splittings | Parton-level momenta | Theory parameters |
|---|---|---|---|---|
| $x$ | $z_d$ | $z_s$ | $z_p$ | $\theta$ |

$$p(x|\theta) = \int \mathrm{d}z_d \int \mathrm{d}z_s \int \mathrm{d}z_p \; p(x|z_d) \cdots \quad p(z_d|z_s) \quad p(z_s|z_p) \quad p(z_p|\theta)$$

It's infeasible to calculate the integral over this enormous space!

[M. Cacciari, G. Salam, G. Soyez 0802.1189]

[CMS]

[Pythia]
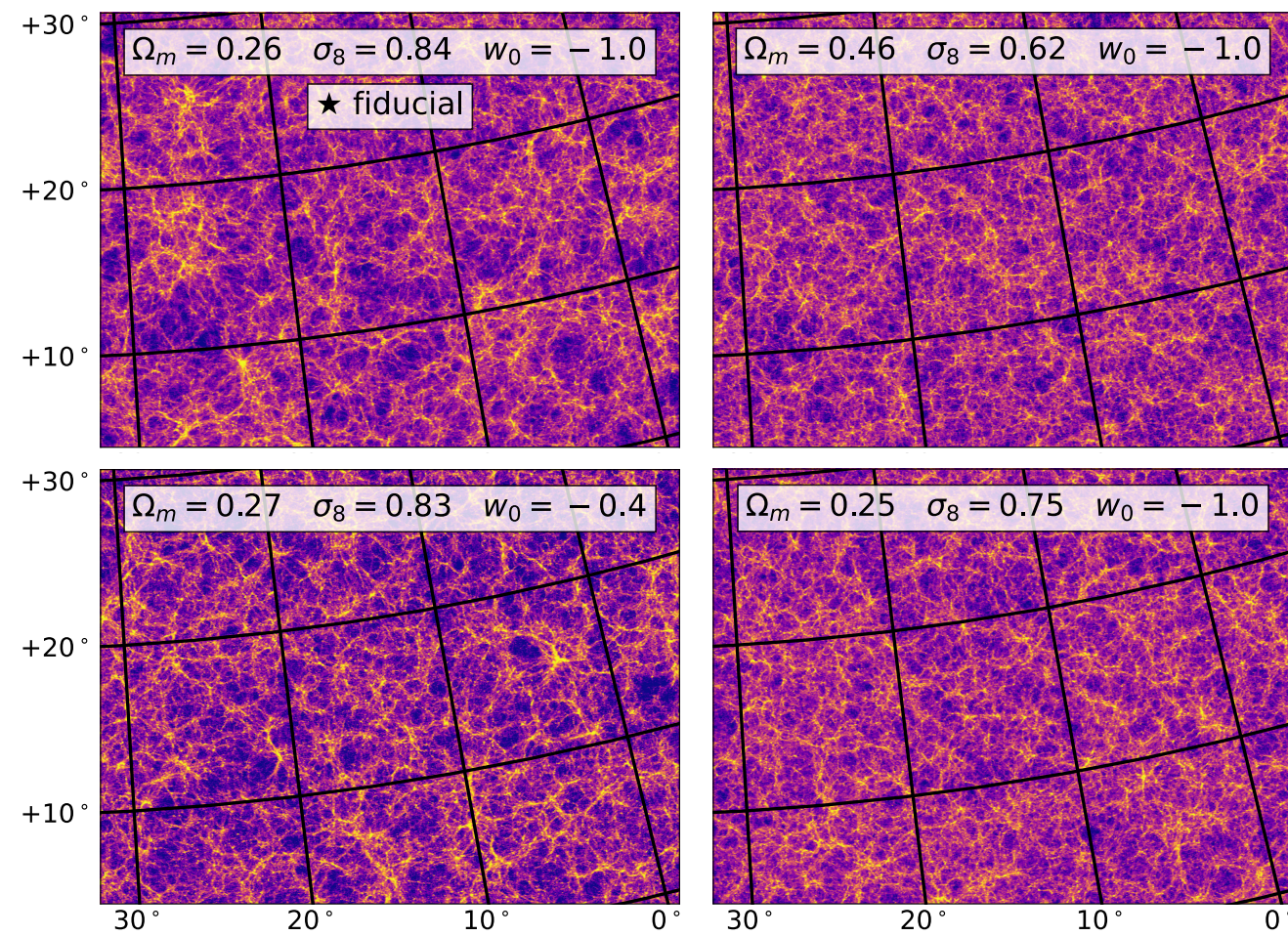
Prediction (simulation)

Inference

# What could we do with $p(x)$?

## Produce samples

*for downstream applications: a fast simulator/emulator*

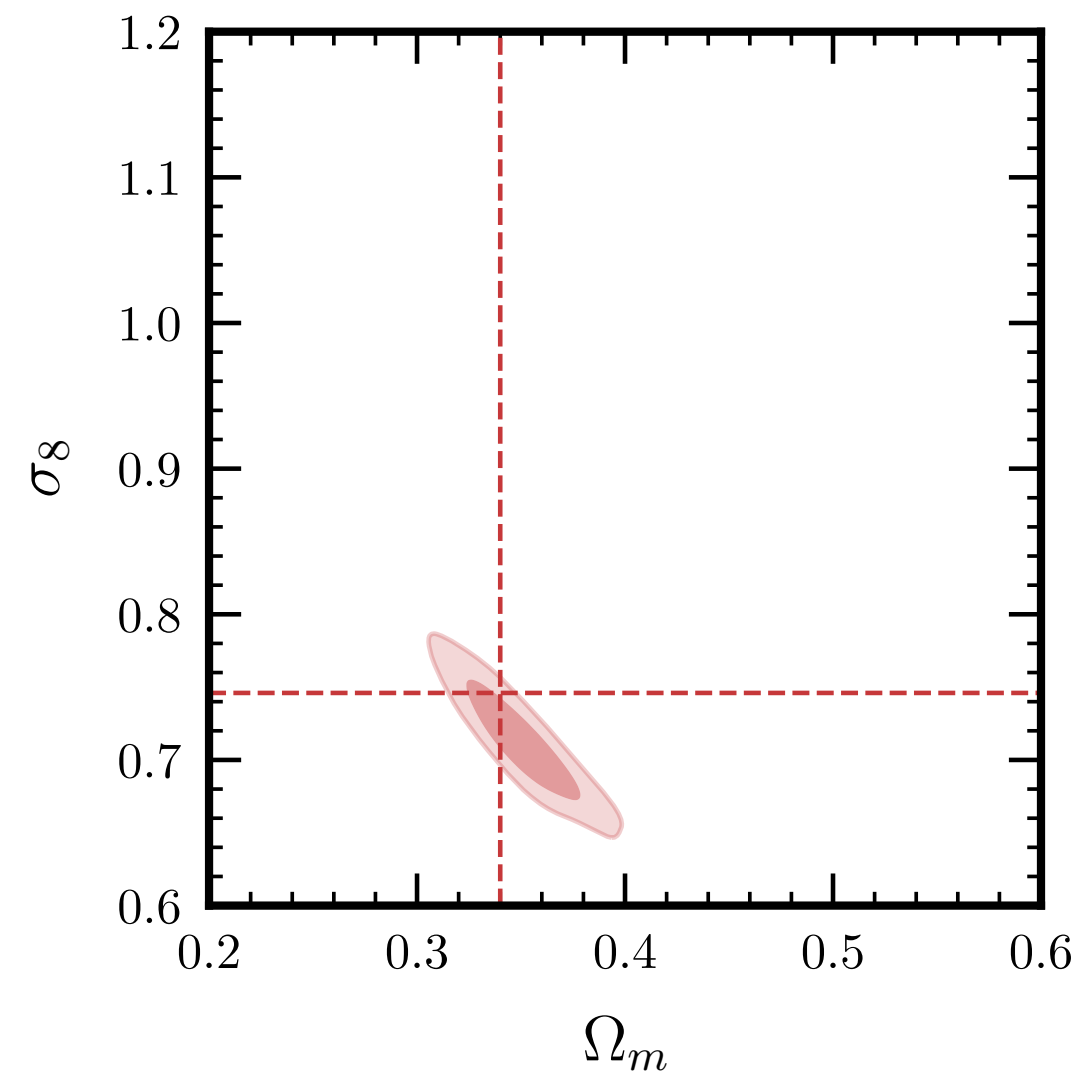$$x \sim p(x \mid \theta)$$
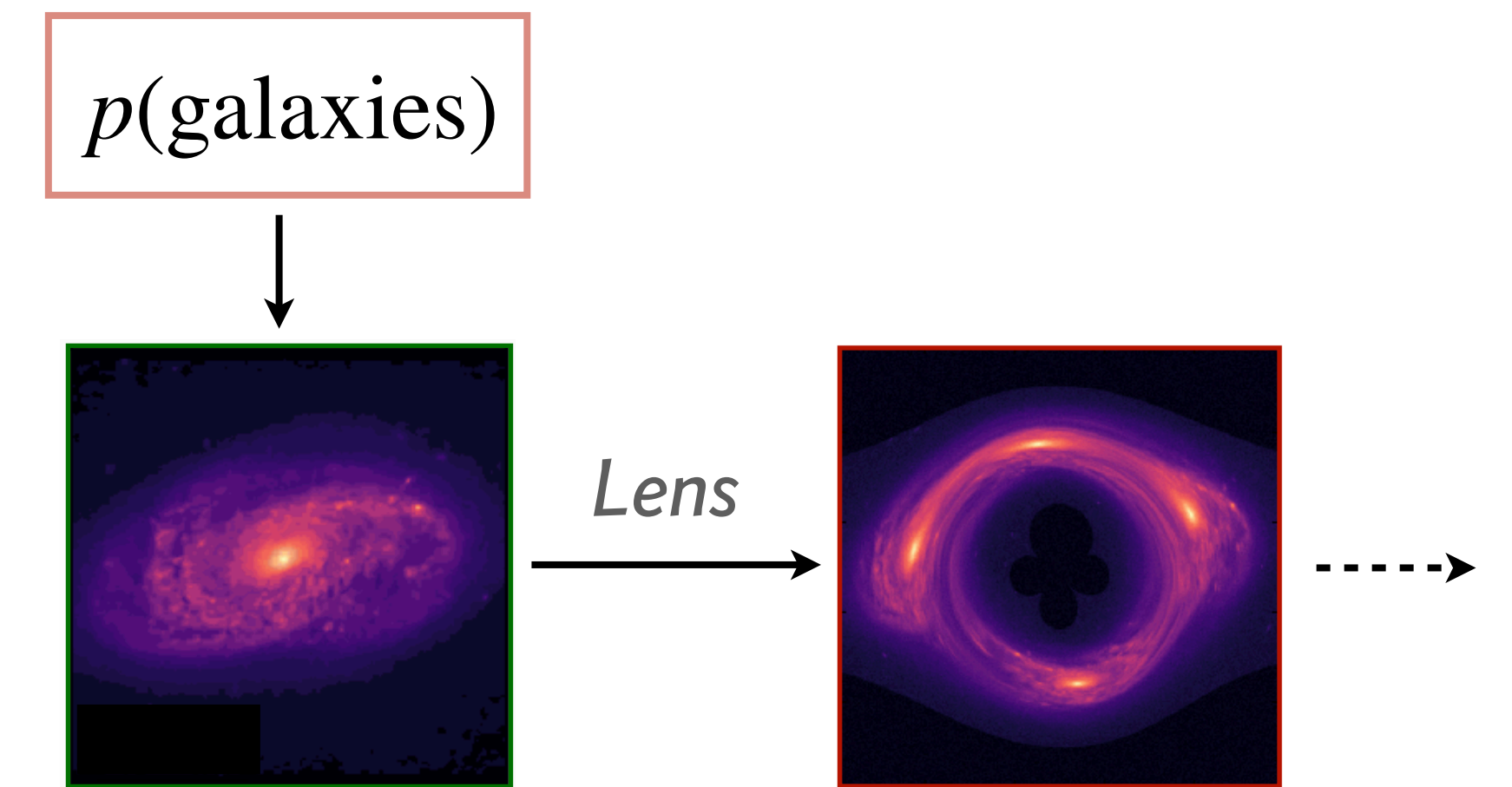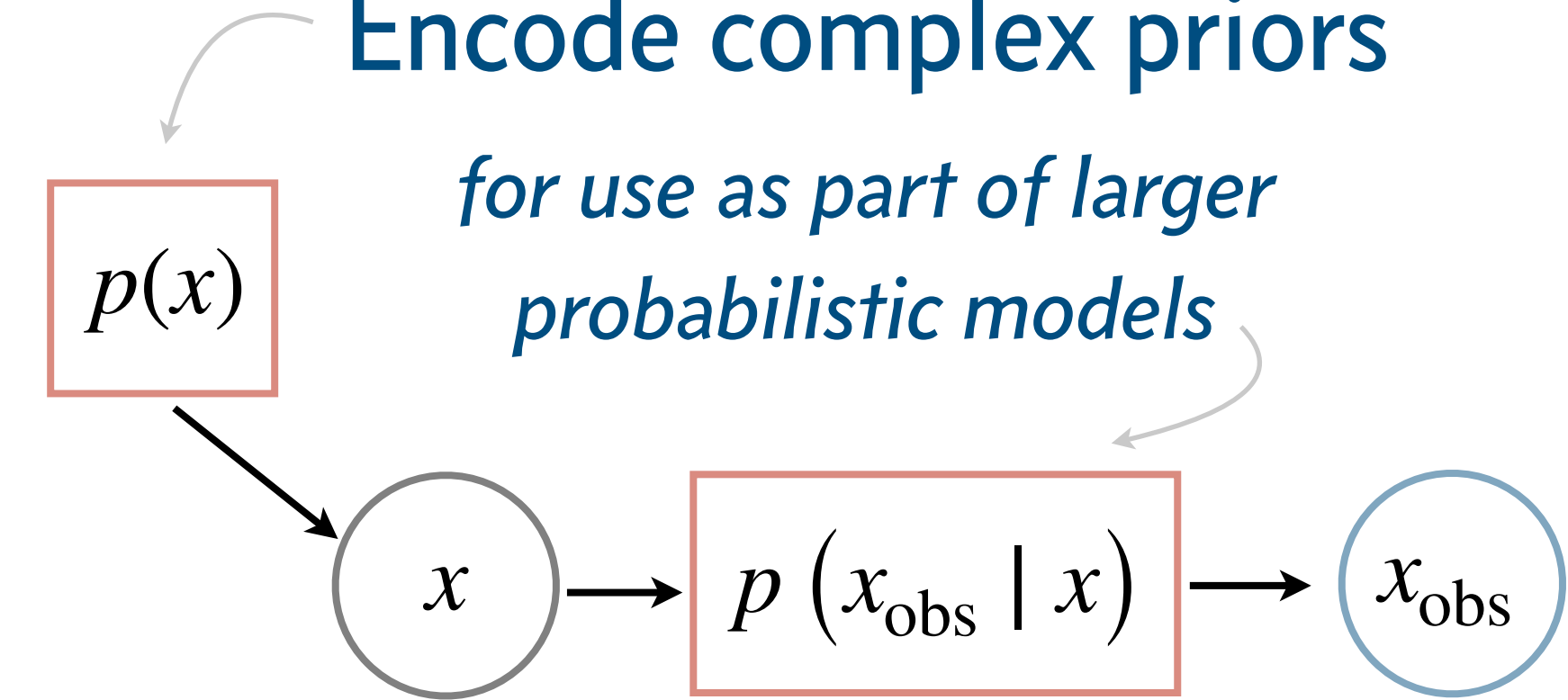


[Kacprzak et al 2022]

## Evaluate likelihood

*for model selection, parameter inference, outlier detection, …*

$$p(\theta \mid x) = \frac{p(x \mid \theta) \cdot p(\theta)}{p(x)}$$



## Encode complex priors

*for use as part of larger probabilistic models*

$$p(x)$$

$$x \longrightarrow p\left(x_{\text{obs}} \mid x\right) \longrightarrow x_{\text{obs}}$$

$$p(\text{galaxies})$$



*Predicted reconstruction*    *Observed lensed i*    *True source image*    *Predicted reconstruction*    *Predicted sou*
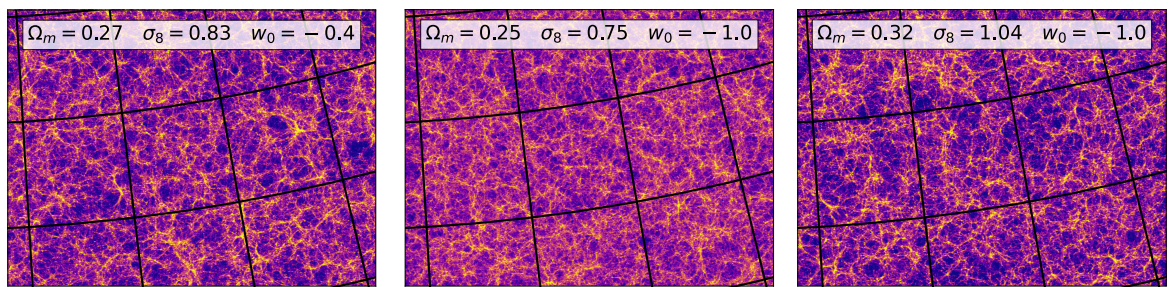
*Lens*

0.5″

*Generative modeling can efficiently enable these for a wide variety of scientific data/models!*

# Learning the data distribution

*I'm sold! How do I learn a generative model for my data?*



1. Ingredients:
   - A parameterized distribution $p_\varphi(x)$
   - Samples from the data distribution $\{x\}_{\text{train}} \sim p(x)$

   (empirical or simulated)

2. Maximize the likelihood of the model under the training data samples
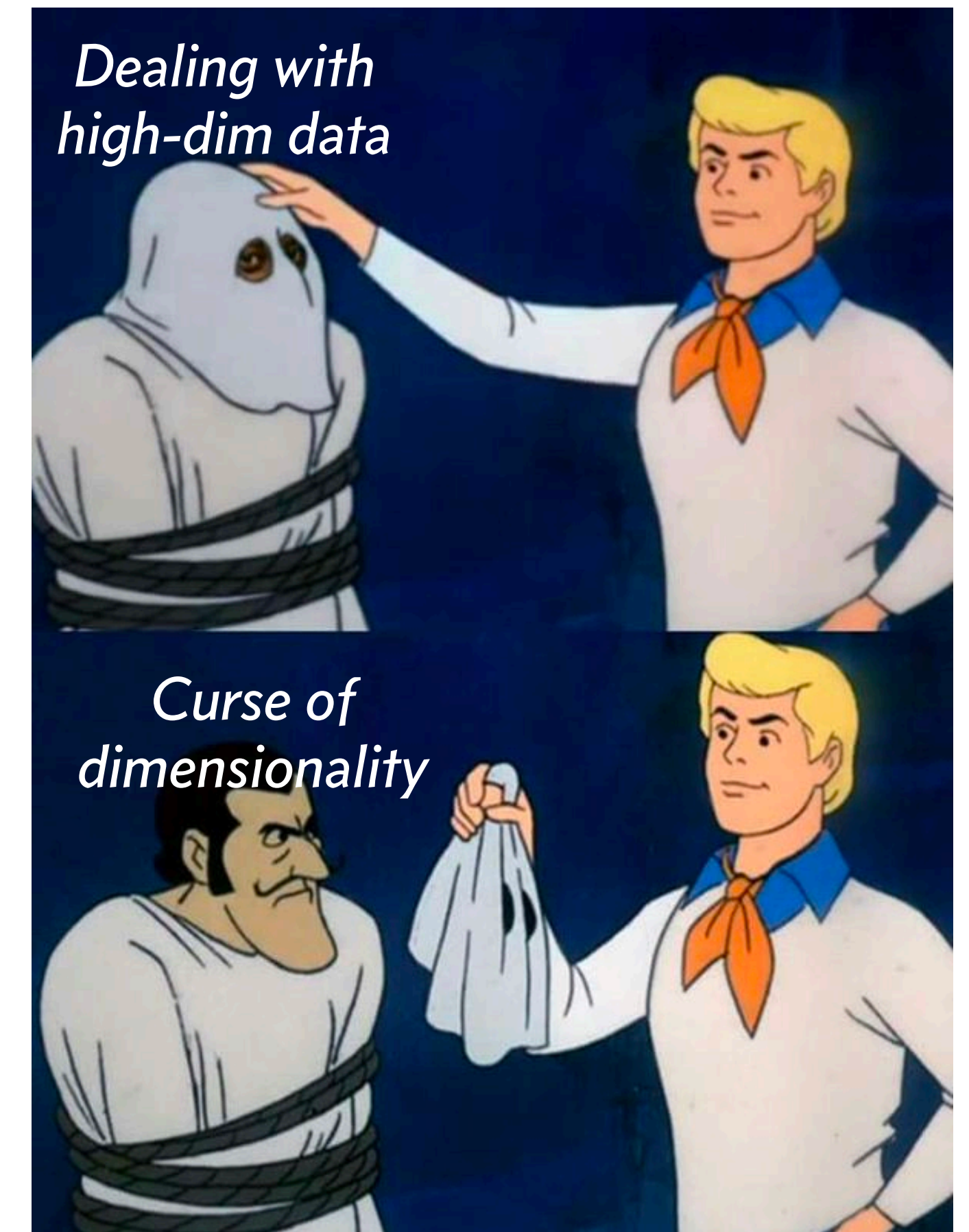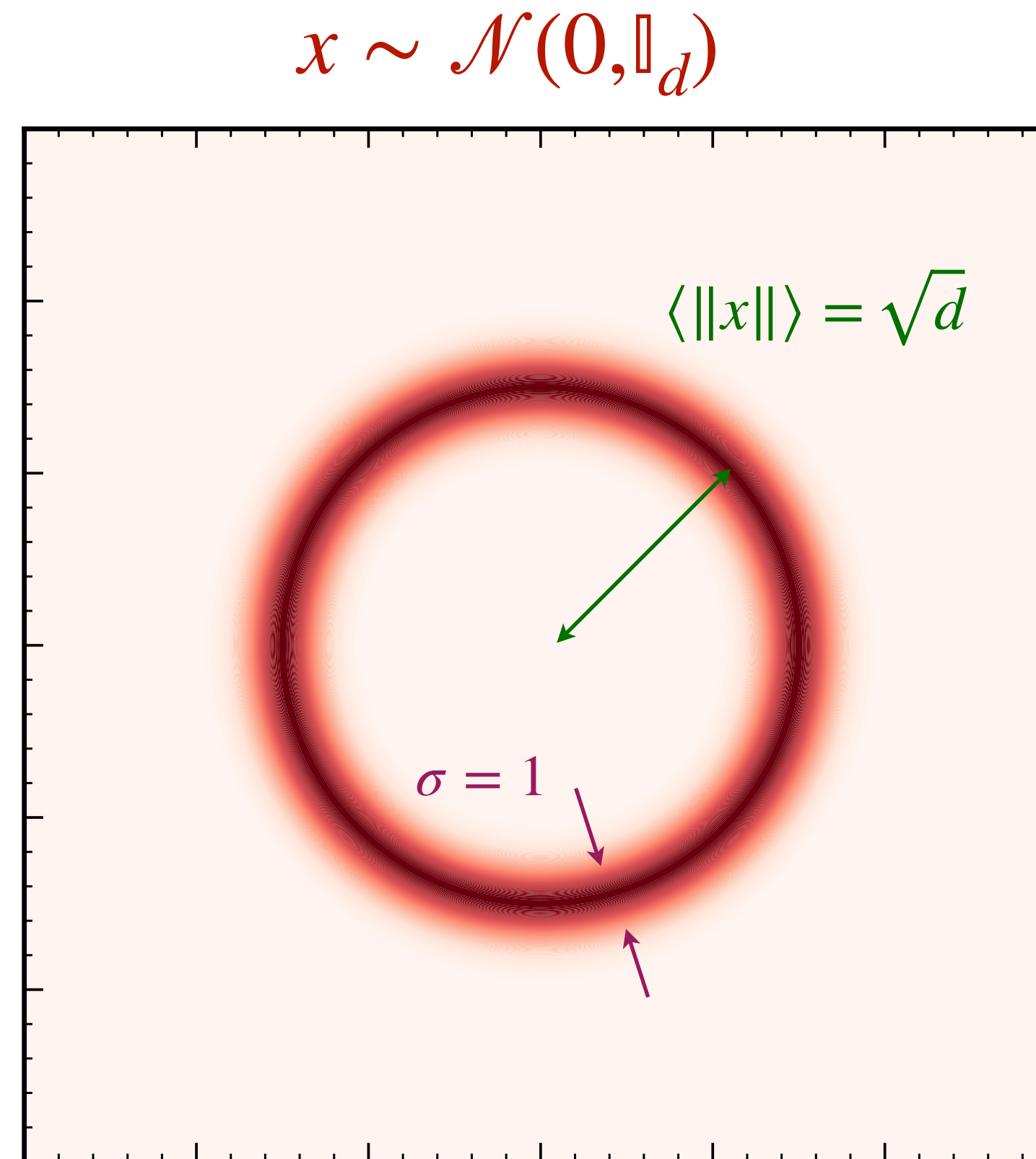
$$\hat{\varphi} = \arg\max_{\varphi} \left[ \log p_\varphi \left( \{x\}_{\text{train}} \right) \right]$$

*Not so fast...*

# The curse of dimensionality

Where is most of the probability mass concentrated in high dimensions?

$$x \sim \mathcal{N}(0, \mathbb{I}_d)$$



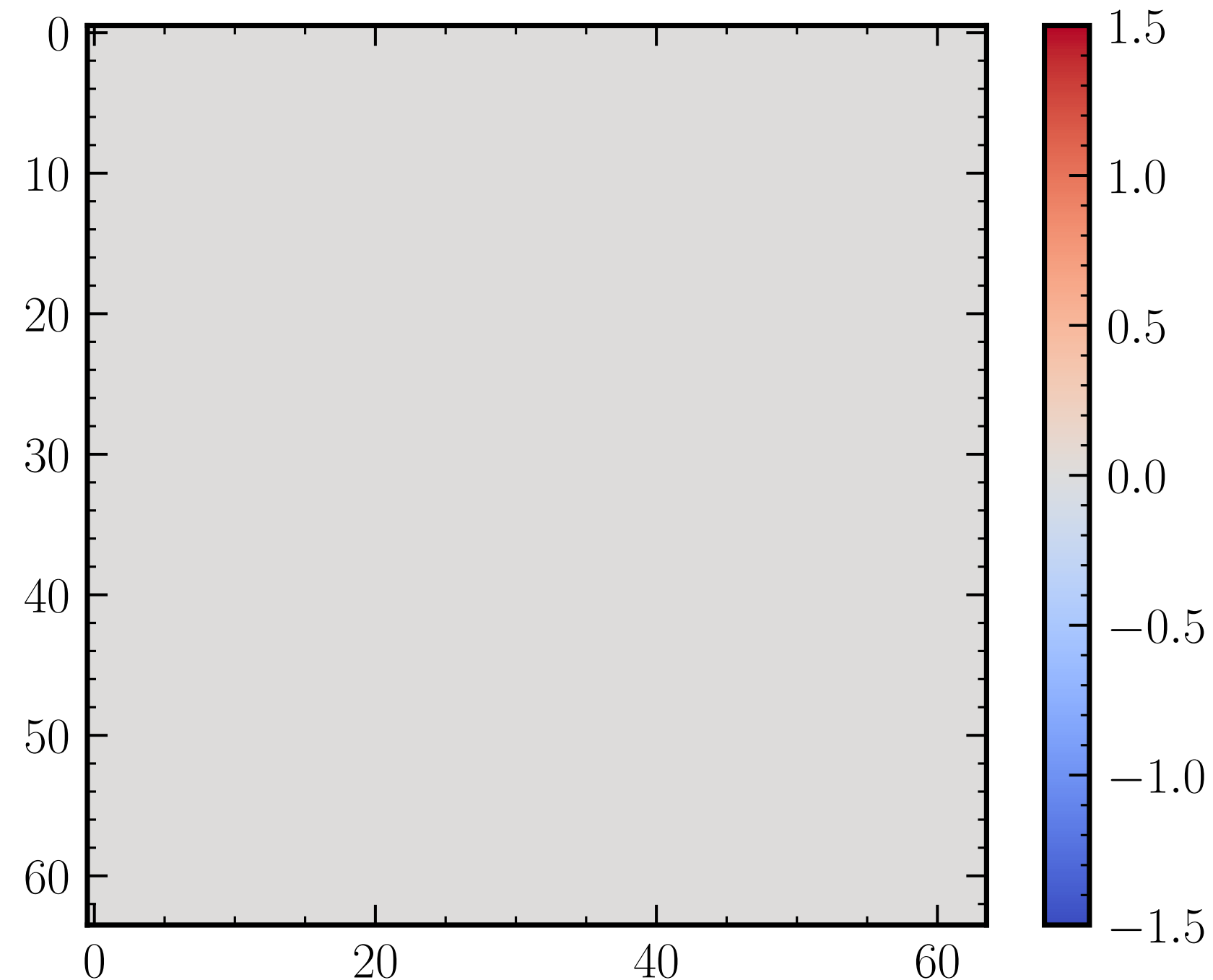$\langle \|x\| \rangle = \sqrt{d}$

$\sigma = 1$

- In high dimensions, most of the probability density of a Gaussian distribution lies in a thin shell at distance $\sqrt{d}$ from the center

- Vanishingly small fraction of distribution support is actually occupied.



Dealing with high-dim data

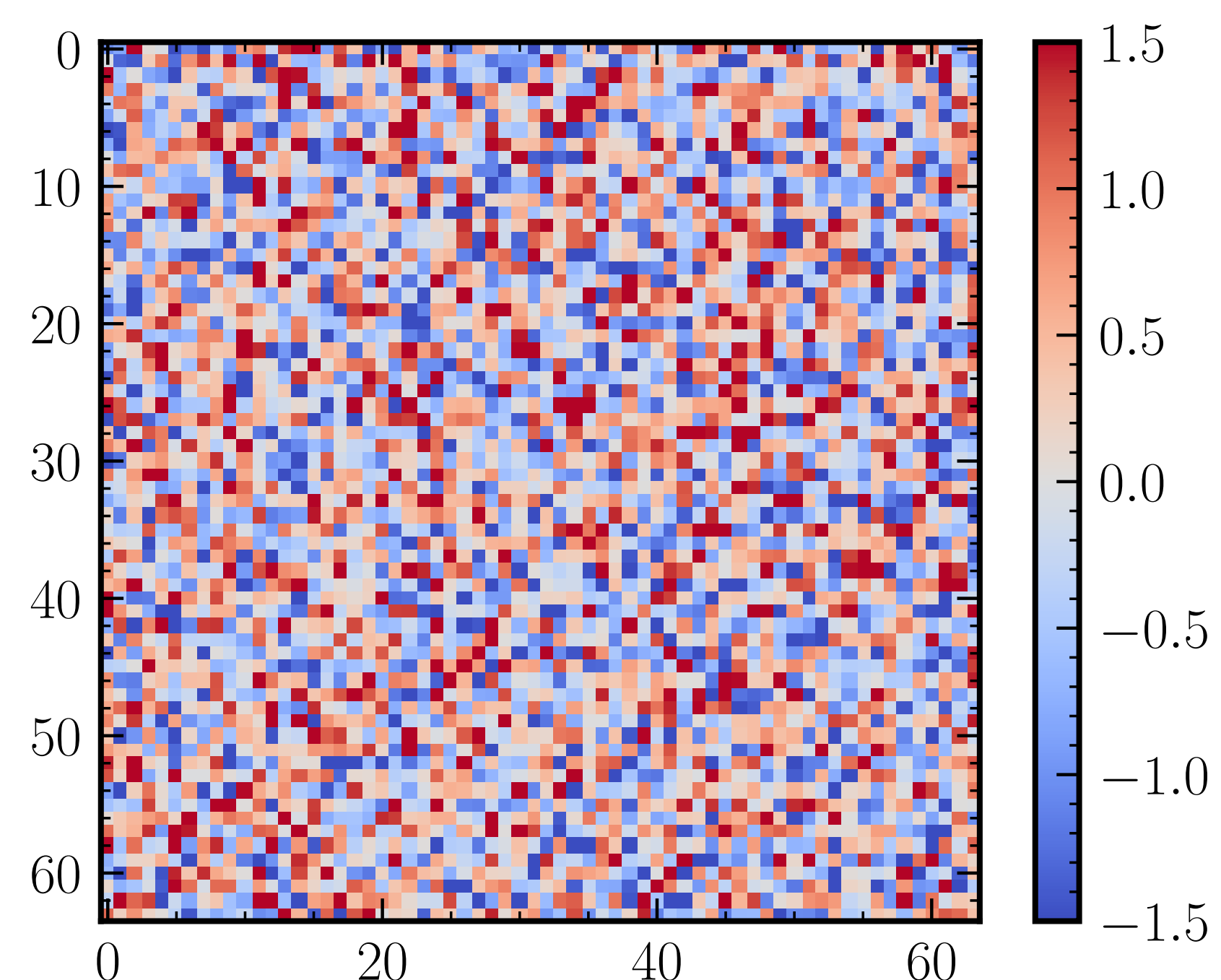Curse of dimensionality

*Learning high-dimensional distributions is challenging!*

# *Typicality* and likelihood of samples

Which of these samples have a higher likelihood under $\mathscr{L} = \mathscr{N}(0, \mathbb{I}_d)$?



$$\log \mathscr{L} \approx -0.92 \text{ nats/dim}$$
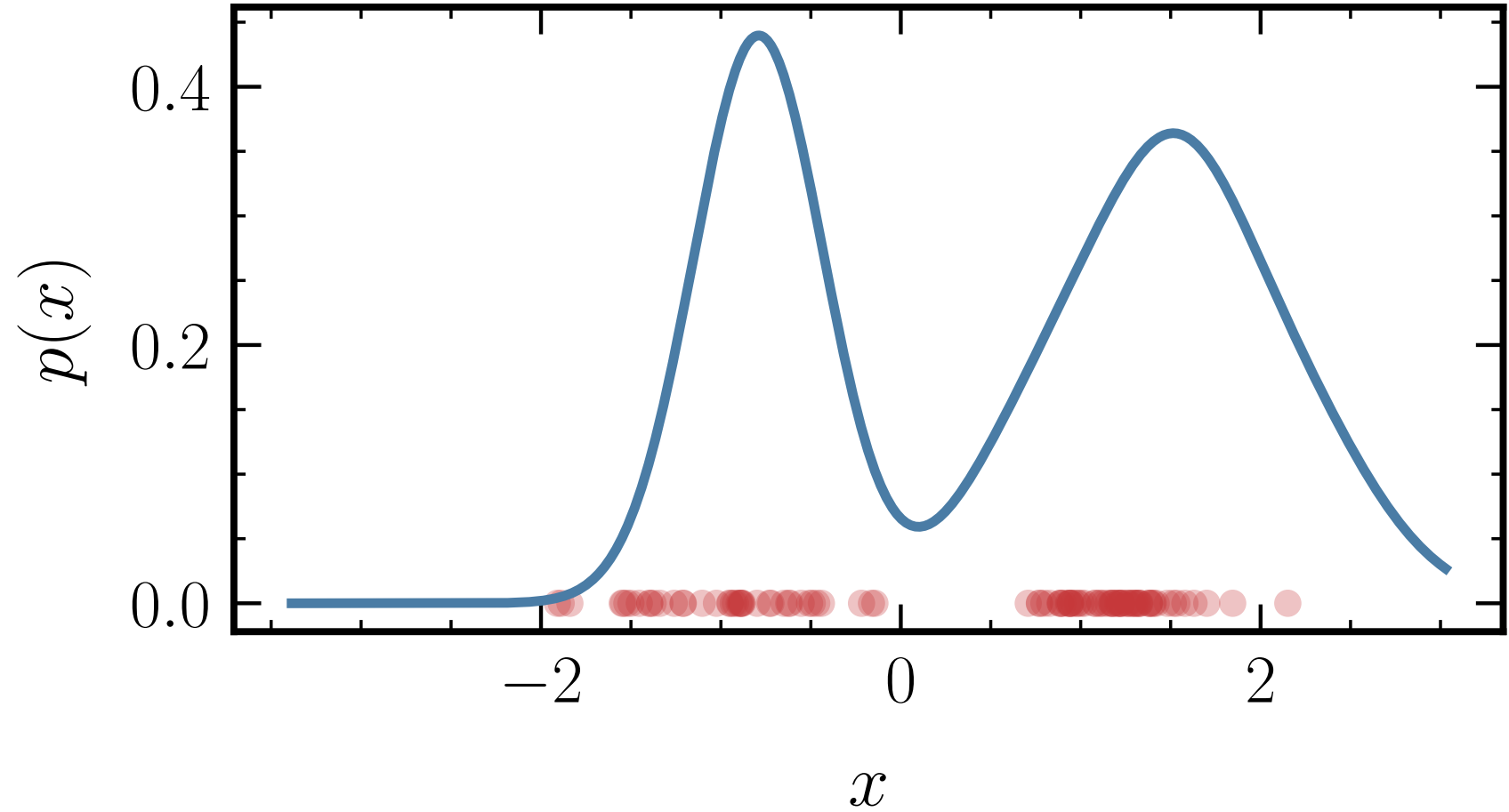
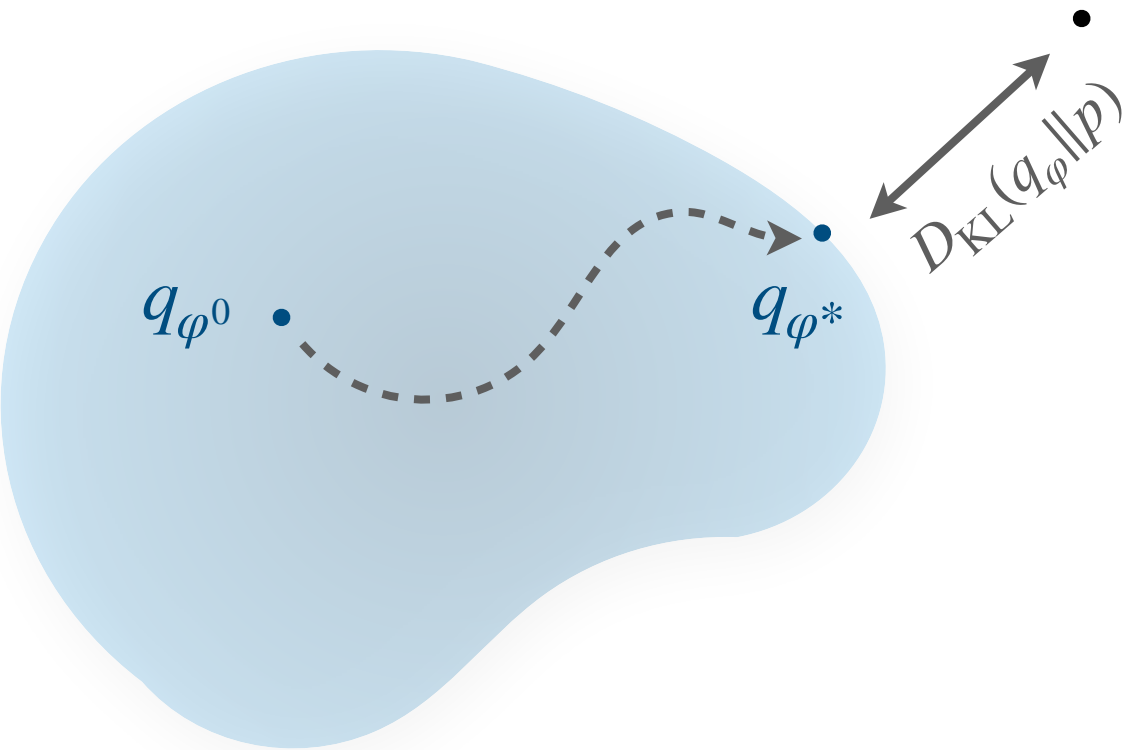$$\log \mathscr{L} \approx -1.43 \text{ nats/dim}$$

*Evaluation of high-dimensional distributions is challenging!*

# (Some) Ways of training deep generative models

**Maximum-likelihood**



**Optimizing a bound on the likelihood**



$D_{\mathrm{KL}}(q_\varphi \| p)$

$q_{\varphi^0}$

$q_{\varphi^*}$

**Score-matching**



**Optimal transport**



**Adversarial training**

# Outline



## Why (deep) generative modeling?

*What is it, and what can it do for you?*



## Variational auto encoders

*Latent-variable modeling, and compression is all you need*



## Diffusion models

*Models based on iterative refinement*

Nonlinear flows

*Invertible transformations*

# A bird-eye view



Original image

Latents

Reconstruction

$\mu$

$\sigma$

$z$

(Probabilistic) **Encoder**

(Probabilistic) **Decoder**

# The pursuit of low-dimensional structure

Real-world datasets often live in structured low-dimensional manifolds

"Difficult to model" $x$                            "Easy to model" $z$



Lower-dim manifold

$$p(z \mid x)$$

# Latent-variable modeling

Make the problem easier by making it "harder":

introduce *joint distribution* $p_\theta(x, z)$

Learn lower-dimensional structure in the data distribution



$$p_\vartheta(x)$$

$$x$$

$$N$$

*Observed variables*



$$p_\vartheta(x, z)$$

$$z \longrightarrow x$$

$$N$$

*Latent variables*          *Observed variables*

Common factorization:

$$p_\vartheta(x, z) = p(z) \cdot p_\vartheta(x \mid z)$$

# Latent-variable modeling

**Maximum-likelihood training?**

$$\vartheta^* = \arg\max_{\vartheta} p_{\vartheta}(x)$$

$$= \arg\max_{\vartheta} \int p_{\vartheta}(x \mid z) p(z) \, \mathrm{d}z$$

$$= \arg\max_{\vartheta} \left\langle p_{\vartheta}(x \mid z) \right\rangle_{p(z)}$$

*Difficult to build a good estimator!*



*Dealing with high-dim data*

The intractability of $p(x)$ is closely related to the intractability of the *posterior* $p(z \mid x)$

*Curse of dimensionality*

$$p(z \mid x) = \frac{p(x, z)}{p(x)} \quad \text{(Bayes' theorem)}$$

# Kullback–Leibler (KL) divergence

A measure of similarity between two probability distributions

$$D_{\mathrm{KL}}(P\|Q) = \int_{-\infty}^{\infty} \mathrm{d}x\, p(x)\, \log\left(\frac{p(x)}{q(x)}\right)$$

$$= \left\langle \log \frac{p(x)}{q(x)} \right\rangle_{x \sim p(x)}$$

$$= -\left\langle \log q(x) \right\rangle_{p(x)} + \left\langle \log p(x) \right\rangle_{p(x)}$$

Cross-entropy $\mathbb{H}(P, Q)$ $-$ (Self-)entropy $\mathbb{H}(P)$

Formally: *expected excess "surprise" from using Q as a model when the actual distribution is P*

# KL-divergence

$$D_{\mathrm{KL}}(Q\|P) = \int_{-\infty}^{\infty} \mathrm{d}x \, q(x) \, \log\left(\frac{q(x)}{p(x)}\right)$$

A measure of similarity between two probability distributions

Not symmetric! $\quad D_{\mathrm{KL}}(Q\|P) \neq D_{\mathrm{KL}}(P\|Q)$

"True" distribution

"Forward" KL $D_{\mathrm{KL}}(P\|Q)$     "Reverse" KL $D_{\mathrm{KL}}(Q\|P)$



$q_\varphi(z)$

$p(z)$

Forward KL

$$D_{\mathrm{KL}}(P_{\mathscr{D}}\|Q_\varphi) = -\Big\langle \log q_\varphi(z) \Big\rangle_{z \sim p_{\mathscr{D}}(z)} + \mathrm{const}.$$

Maximum-likelihood inference is equivalent to minimizing the *forward* KL

Non-negative! $\quad D_{\mathrm{KL}}(Q\|P) \geq 0$

# Variational inference

Infer the posterior over the latent parameters



$p(z \mid x)$

$D_{\mathrm{KL}}(q_\varphi \| p)$

Gradient ascent

$q_{\varphi^0}$

$q_{\varphi^*}$

$q_\varphi(z \mid x)$

**A two-for-one!**

☑ Estimate approximate posterior $q_\varphi(z \mid x) \approx p(z \mid x)$

☑ Estimate likelihood/evidence $\mathrm{ELBO} \approx p(x)$

$\geq 0$      Evidence    $-$    Evidence Lower BOund (ELBO)

$$D_{\mathrm{KL}}\left(q_\varphi(z \mid x)\|p(z \mid x)\right) = \log p(x) - \left\langle \log p_\vartheta(x, z) - \log q_\varphi(z) \right\rangle_{q_\varphi(z)}$$

# Variational inference

A general-purpose technique for posterior estimation: *optimization* instead of *sampling*



[EuCAPT White Paper 2021]

# A Bayesian latent-variable model optimized with variational inference

*We're so back*

**Reverse process**

$$p_\vartheta\left(x \mid z\right) \cdot p(z)$$

$z$    $x$    $N$

$$q_\varphi(z \mid x) \cdot p(x)$$

**Forward process**

*It's so over*

Maximizing ELBO

$\equiv$ Minimizing *reverse* KL

$\equiv$ "Aligning the forward and reverse processes"

$$\text{Minimize} \quad \left\langle \log \frac{q(x, z)}{p(x, z)} \right\rangle$$

**Forward process**    **Reverse process**

$x \longrightarrow z \longrightarrow x'$

**Latents**

# Variational inference

A general-purpose technique for posterior estimation

$\geq 0$      Evidence    $-$    Evidence Lower BOund (ELBO)

$$D_{\mathrm{KL}}\Big(q_\varphi(z)\|p(z\mid x)\Big) = \log p(x) - \Big\langle \log p_\vartheta(x,z) - \log q_\varphi(z)\Big\rangle_{q_\varphi(z)}$$



$p(z\mid x)$

$q_\varphi(z)$

$$\mathrm{ELBO} = \Big\langle \log p_\vartheta(x,z) - \log q_\varphi(z\mid x)\Big\rangle_{q_\varphi}$$

$$= \Big\langle \log p_\vartheta(x\mid z) + \log p(z) - \log q_\varphi(z\mid x)\Big\rangle_{q_\varphi}$$

$$= \Big\langle \log p_\vartheta(x\mid z)\Big\rangle_{q_\varphi} - D_{\mathrm{KL}}\Big(q_\varphi(z\mid x)\,\|\,p(z)\Big)$$

"Reconstruction"             "Regularization"

# VAEs in practice

$$p(z) = \mathcal{N}(z; 0, I)$$

Variational
params.

Latents

*Original image*

*Reconstruction*

$$\mu$$

$$\sigma$$

$$z$$

$$\mu, \sigma^2 = \mathrm{NN}_\varphi(x)$$

$$x \sim p(x)$$

$$z \sim q_\varphi(z \mid x)$$

$$x' \sim p_\vartheta(x \mid z)$$

(Probabilistic) **Encoder**

$$q_\varphi(z \mid x) = \mathcal{N}(z; \mu, \sigma^2 \mathbb{I})$$

(Probabilistic) **Decoder**

*Noise model / data likelihood*

# VAEs in practice

$$q_\varphi(z \mid x)$$



$$\text{ELBO} = \left\langle \log p_\vartheta(x \mid z) \right\rangle_{q_\varphi} - D_{\text{KL}} \left( q_\varphi(z \mid x) \parallel p(z) \right)$$

$$D_{\text{KL}} \left( q_\varphi(z \mid x) \parallel p(z) \right)$$

Regularization

$\mu$

$\sigma$

$z$

*Reconstruction*

$$\left\langle \log p_\vartheta(x \mid z) \right\rangle_{q_\varphi}$$

Reconstruction (e.g., MSE, …)

$$\|x - x'\|_2^2$$

# A semantically meaningful latent space

The KL-term enforces simplicity in the latent space, encouraging learned semantic structure and *disentanglement*

*Pure reconstruction*

*More latent regularization*



*First two PCs*

# Neural compression: *Rate-distortion theory*

Autoencoding is a form of (neural) compression!

$$-\mathrm{ELBO} = -\left\langle \log p_\vartheta(x \mid z) \right\rangle_{q_\varphi} + D_{\mathrm{KL}}\left( q_\varphi(z \mid x) \,\|\, p(z) \right)$$

Distortion

Rate

"Reconstruction loss"

"Amount of compression"



*Rate-distortion curve quantified this tradeoff*

# Controlling compression and disentanglement: $\beta$-VAEs

$$-\text{ELBO} = -\Big\langle \log p_\vartheta(x \mid z) \Big\rangle_{q_\varphi} + \textcolor{red}{\beta} \cdot D_{\text{KL}}\Big( q_\varphi(z \mid x) \,\|\, p(z) \Big)$$

<span style="color:#2a6db0">Distortion</span>　　　　　　　　　　<span style="color:#a83268">Rate</span>

If the data-generating process is associated with a principled noise model, by using it (the *likelihood*) as the reconstruction loss we are aiming to reconstruct the mean data.

$$\log p(x \mid z; x') = -\frac{1}{2}\left(\frac{x - x'}{\sigma}\right)^2 + \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right)$$

- Larger $\sigma$: More of the data variation is attributed to the likelihood $\rightarrow$ larger "$\beta$", more compression

- Smaller $\sigma$: Latents $z$ try to capture more of the variation in the data (e.g. small perceptual features)

# Tutorials 1 and 2: variational inference and VAEs

Lead: Carol Cuesta-Lazaro



- Implement the ELBO objective for variational inference

- Construct a VAE and use it to build a generative model of galaxy images using samples form the HST COSMOS dataset

- Boilerplate code for training/reconstruction/ sampling for quick iteration

- Experiment with trade-offs between pure reconstruction and a latent space regularization

[Mandelbaum et al; https://zenodo.org/record/3242143]

# Outline



## Why (deep) generative modeling?

*What is it, and what can it do for you?*



## Variational auto encoders

*Latent-variable modeling, and compression is all you need*



## Diffusion models

*Models based on iterative refinement*

N...li...fl...

*Invertible transformations*

# Diffusion models: overview

**Forward process** *(adding noise)*

$\longrightarrow$

$x(t=0) \sim p(x)$

$x(t=1) \sim \mathcal{N}(0,1)$



$x_{t-1} \sim p(x_{t-1} \mid x_t)$

$\longleftarrow$

**Reverse process** *(denoising)*

Prompt: *"A cat perched on an Ikea Poang chair"*

image $\sim p(\text{image} \mid \text{text prompt})$

# Towards diffusion: hierarchical VAEs

# Towards diffusion: (Markovian) hierarchical VAEs

Reverse process
$$p\left(x, z_1, z_2, \cdots, z_T\right) = p\left(z_T\right) p\left(z_{T-1} \mid z_T\right) \cdots p\left(z_1 \mid z_2\right) p\left(x \mid z_1\right)$$



$$p\left(z_{T-1} \mid z_T\right) \quad p\left(z_3 \mid z_4\right) \quad p\left(z_2 \mid z_3\right) \quad p\left(z_1 \mid z_2\right) \quad p\left(x \mid z_1\right)$$

$$p(z) \longrightarrow \boxed{z_T} \quad \bullet\bullet\bullet \quad \boxed{z_3} \quad \boxed{z_2} \quad \boxed{z_1} \quad \boxed{x} \longleftarrow q(x)$$

$$q(z_T \mid z_{T-1}) \quad q(z_4 \mid z_3) \quad q(z_3 \mid z_2) \quad q(z_2 \mid z_1) \quad q(z_1 \mid x)$$

Forward process
$$q\left(x, z_1, z_2, \cdots, z_T\right) = q(x) q\left(z_1 \mid x\right) q\left(z_2 \mid z_1\right) \cdots q\left(z_T \mid z_{T-1}\right)$$

# Towards diffusion: (Markovian) hierarchical VAEs

Diffusion models can be seen as hierarchical VAEs with a few restrictions:

- The forward *(encoding)* distribution prescribed as a Markov chain of Gaussians; it is not learned

$$q\left(z_t \mid z_{t-1}\right) = \mathcal{N}\left(z_t; \alpha_t z_{t-1}, \beta_t\right)$$

- Distributions of latents at the final timestep $T$ is a standard (unit) Gaussian

$$q(z_T \mid z_{T-1}, \dots x) = \mathcal{N}(z_T; 0, \mathbb{I})$$

- The dimensionality of latents is the same as the data dimensionality

$$\dim(z_t) = \dim(x)$$

# Variational diffusion models

Align the forward and reverse distributions;
variational lower bound (ELBO) as before

$$L = \left\langle \log \frac{q\left(x, z_1, z_2, \cdots, z_T\right)}{p\left(x, z_1, z_2, \cdots, z_T\right)} \right\rangle_{q(x)}$$



Some calculations later

$$L = \left\langle \log p_\vartheta\left(x \mid z_1\right) \right\rangle_{q(z_1|x)} - D_{\mathrm{KL}}\left(q\left(z_T \mid x\right) \| p\left(z_T\right)\right) - \sum_{t=2}^{T} \left\langle D_{\mathrm{KL}}\left(q\left(z_{t-1} \mid z_t, x\right) \| p_\vartheta\left(z_{t-1} \mid z_t\right)\right) \right\rangle_{q(z_t|x)}$$

**Reconstruction**

*(Noise model; no trainable parameters)*

**Prior regularization**

*(No trainable parameters)*

**Denoising matching**

ELBO! Bound on $p(x)$

# The forward process and diffusion kernel

<span style="color:red">Predict arbitrary timestep without Markovian sampling</span>

$$\sum_{t=2}^{T} \left\langle D_{\mathrm{KL}} \left( q \left( z_{t-1} \mid z_t, x \right) \| p_{\vartheta} \left( z_{t-1} \mid z_t \right) \right) \right\rangle_{q(z_t|x)}$$

$q \left( z_t \mid x \right)$



$z_T$ ••• $z_3$ $z_2$ $z_1$ $x$ ⟵ $q(x)$

**Variance-preserving noise schedule**

$$q \left( z_t \mid z_{t-1} \right) = \mathcal{N} \left( \sqrt{1 - \beta_t} \cdot z_{t-1}, \beta_t \right)$$

$$z_t = \sqrt{1 - \beta_t} \cdot z_{t-1} + \sqrt{\beta_t} \cdot \varepsilon$$

**Diffusion kernel** $\qquad \alpha_t = 1 - \beta_t$

$$q \left( z_t \mid x \right) = \mathcal{N} \left( \sqrt{\bar{\alpha}_t} \cdot x, \sqrt{1 - \bar{\alpha}_t} \right)$$

$$\bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i$$

# The denoising objective

Given the nature of the forward (noising) process, $q\left(z_{t-1} \mid z_t, x\right)$ can be computed analytically

$$q\left(z_{t-1} \mid z_t, x\right) = \mathcal{N}\left(z_{t-1}; \mu_q(x_t, x_0), \sigma_q(t)\mathbb{I}\right)$$

*Match*      *Set equal*

$$p_\vartheta\left(z_{t-1} \mid z_t, x\right) = \mathcal{N}\left(z_{t-1}; \mu_\vartheta(x_t, x_0), \sigma_\vartheta(t)\mathbb{I}\right)$$

Learnable denoising distribution; assume Gaussian

$$\sum_{t=2}^{T} \left\langle D_{\mathrm{KL}}\left(q\left(z_{t-1} \mid z_t, x\right) \| p_\vartheta\left(z_{t-1} \mid z_t\right)\right)\right\rangle_{q(z_t|x)}$$


Some more calculations later

Denoising loss

$$\frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}\left(1 - \alpha_t\right)^2}{\left(1 - \bar{\alpha}_t\right)^2} \left[\left\| \hat{x}_\theta\left(z_t, t\right) - x \right\|^2\right]$$

# The denoising objectives

$x$-prediction; MLE

$$\frac{1}{2\sigma_q^2(t)} \frac{\bar{\alpha}_{t-1}\left(1-\alpha_t\right)^2}{\left(1-\bar{\alpha}_t\right)^2} \left[ \left\| \hat{x}_\theta\left(z_t, t\right) - x \right\|^2 \right]$$

$\epsilon$-prediction; MLE

$$\frac{1}{2\sigma_q^2(t)} \frac{\left(1-\alpha_t\right)^2}{\left(1-\bar{\alpha}_t\right)\alpha_t} \left[ \left\| \epsilon - \hat{\epsilon}_\theta\left(x_t, t\right) \right\|^2 \right]$$



$\epsilon$-prediction; "simple"

$$\left\| \epsilon - \hat{\epsilon}_\theta\left(x_t, t\right) \right\|^2$$

Typical objective for training image diffusion models: *SOTA on many tasks!*

# Simple objectives as a weighted sum of ELBOs

Kingma and Gao (2023) showed that common objectives can be written as a weighted sum (across different noise levels) of ELBOs

$$L_w(x) = \left\langle \textcolor{green}{w(t)} \cdot \textcolor{red}{w_{\mathrm{ML}}(t)} \left\| \epsilon - \hat{\epsilon}_\theta(z_t, t) \right\|^2 \right\rangle_{t, \epsilon}$$

Additional weighting ($w_{\mathrm{ML}}^{-1}$ for $\epsilon$-prediction)

Weighting for ELBO/ML objective



Monotonic weighting functions

$$L(t; x) \equiv D_{\mathrm{KL}}\left( q\left(z_{t:1} \mid x\right) \| p\left(z_{t:1}\right) \right)$$

$$L_w(x) \propto \int_0^1 \textcolor{green}{\frac{d}{dt} w(t)} L(t; x)\, dt + \textcolor{green}{w(t_0)} L(0; x)$$

Interpretation: *data augmentation with additive Gaussian noise / data-distribution smoothing*

ELBOs at different noise levels

# Continuous-time/SDE formulation

$$x_t = \sqrt{1 - \beta_t} \cdot x_{t-1} + \sqrt{\beta_t} \cdot \varepsilon$$

In the limit of infinite time steps, $\Delta_t \to 0$ and the forward diffusion process can be written as

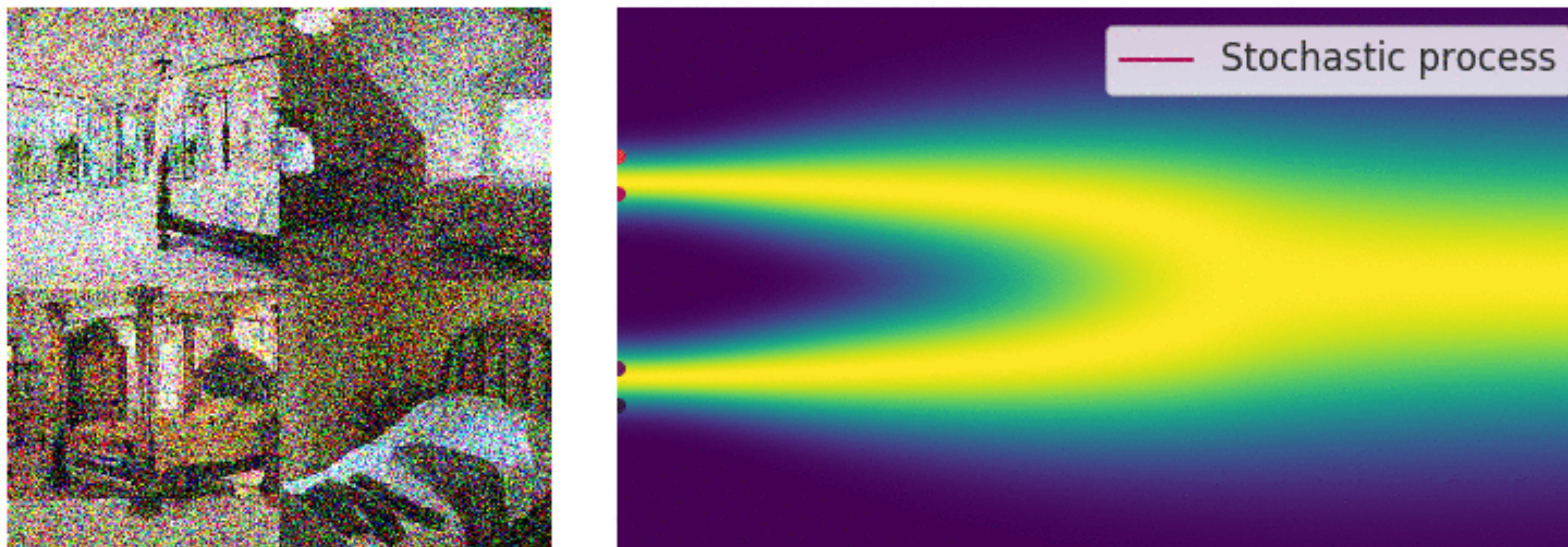$$x_t = \sqrt{1 - \beta(t)\Delta_t} \, x_{t-1} + \sqrt{\beta(t)\Delta_t} \, \mathcal{N}(0, \mathbb{I})$$

$$\approx x_{t-1} - \frac{\beta(t)\Delta_t}{2} x_{t-1} + \sqrt{\beta(t)\Delta_t} \, \mathcal{N}(0, \mathbb{I})$$

Which is an update rule corresponding to the Euler-Murayama discretization of the stochastic differential equation (SDE)

$$\mathrm{d}x_t = -\frac{1}{2}\beta(t)x_t \, \mathrm{d}t + \sqrt{\beta(t)} \mathrm{d}w_t$$

# Continuous-time/SDE formulation

The forward diffusion process defined by an SDE



$$\mathrm{d}x_t = -\frac{1}{2}\beta(t)x_t\ \mathrm{d}t + \sqrt{\beta(t)}\mathrm{d}w_t$$

# The reverse SDE

The reverse process satisfies a reverse-time SDE that can be derived from the forward SDE and the score of the marginal distribution, $\nabla_{x_t} \log q(x_t)$



$$\mathrm{d}x_t = \left[-\frac{1}{2}\beta(t)x_t - \beta(t)\nabla_{x_t}\log q\left(x_t\right)\right]\mathrm{d}t + \sqrt{\beta(t)}\mathrm{d}w_t$$

# Denoising score matching

Need to compute the score $\nabla_{x_t} \log q(x_t)$

The *conditional* score $\nabla_{x_t} \log q(x_t \mid x)$ can be computed using the diffusion kernel

$$\nabla_{x_t} \log q(x_t \mid x) = -\frac{(x_t - x)}{\sigma_t^2} = -\frac{\epsilon}{\sigma_t}$$

Noise-prediction

$$\frac{1}{2\sigma_q^2(t)} \frac{\left(1 - \alpha_t\right)^2}{\left(1 - \bar{\alpha}_t\right) \alpha_t} \left[ \left\| \epsilon - \hat{\epsilon}_\theta \left(x_t, t\right) \right\|^2 \right]$$

$\Longleftrightarrow$

Score-matching *

$$\frac{1}{2\sigma_q^2(t)} \frac{\left(1 - \alpha_t\right)^2}{\alpha_t} \left[ \left\| s_\theta \left(x_t, t\right) - \nabla \log q(x_t) \right\|^2 \right]$$

*The noise- and score-prediction networks are equivalent up to a std-scaling*

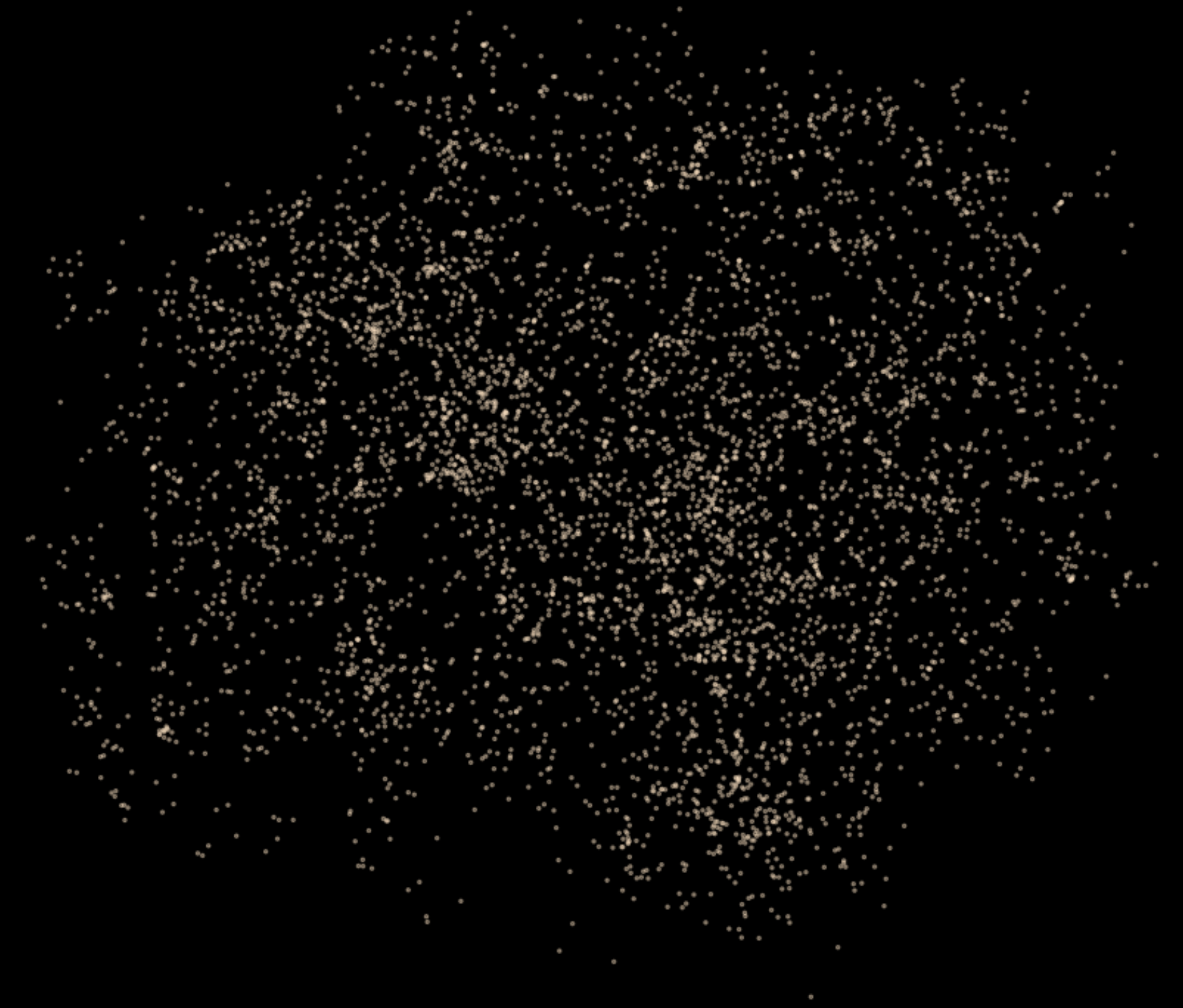# The noise/score-prediction model and *latent diffusion*
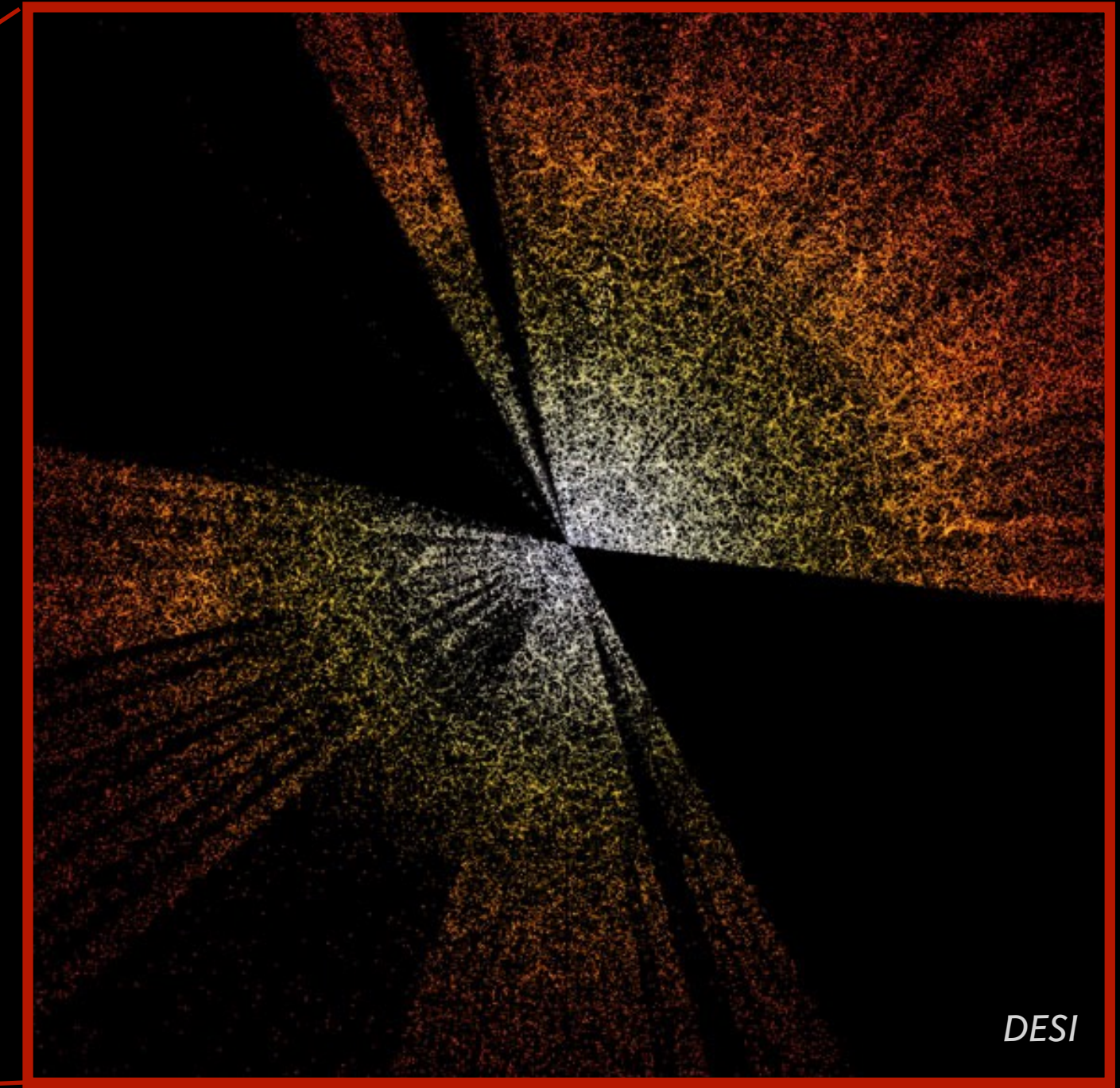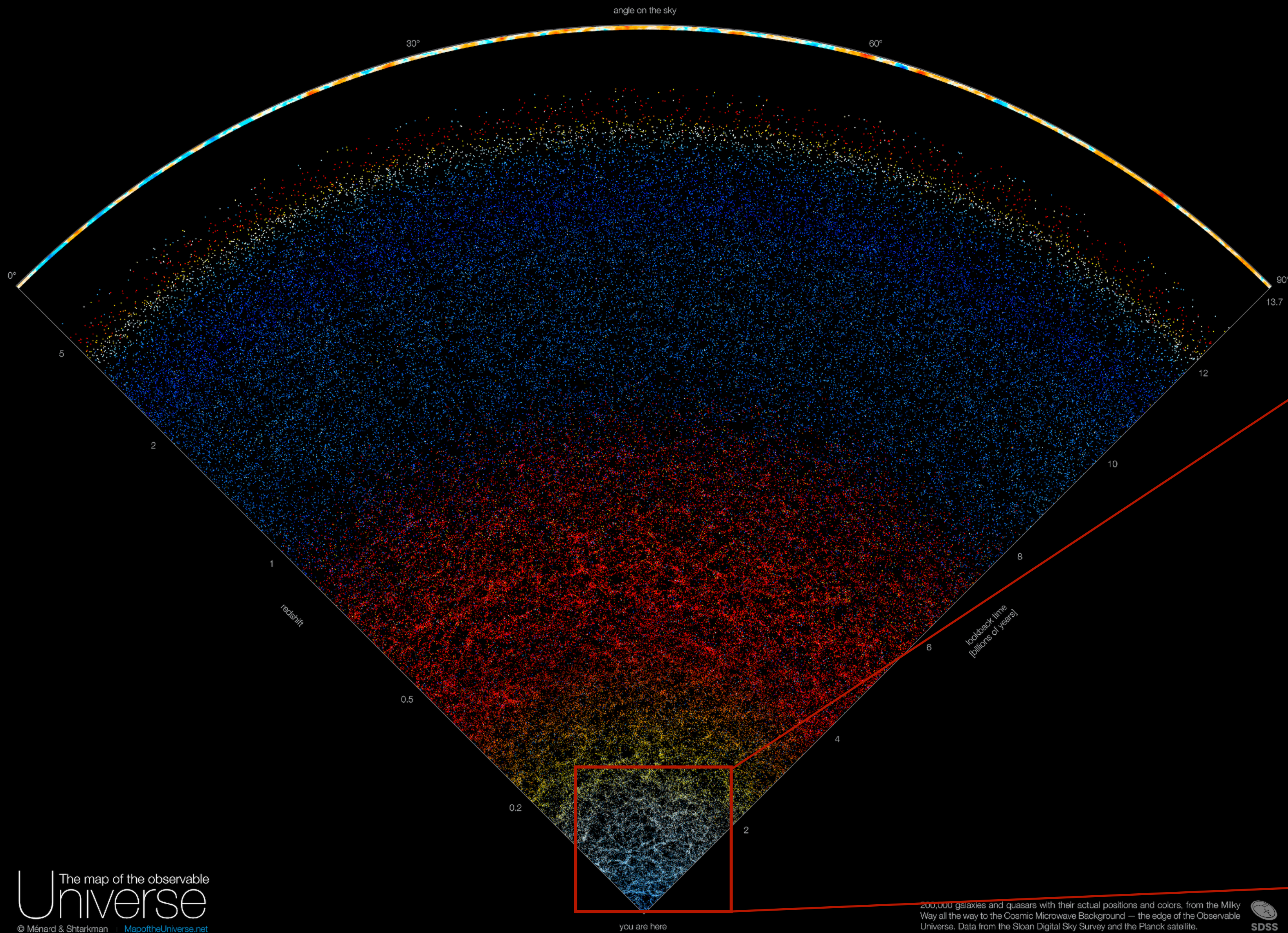
[Rombach et al 2021]



*Perceptual compression while retaining semantically meaningful information*

$\Omega_m = 0.10$

An application to *galaxy clustering*

Work with Carol Cuesta-Lazaro [MIT/IAIFI]

angle on the sky

30°          60°

0°                                                    90°

13.7

5

12

2

10

redshift

1

8

0.5

lookback time
[billions of years]

6

0.2

4

2

DESI

you are here

The map of the observable
# Universe
© Ménard & Shtarkman    |    MapoftheUniverse.net

200,000 galaxies and quasars with their actual positions and colors, from the Milky
Way all the way to the Cosmic Microwave Background — the edge of the Observable
Universe. Data from the Sloan Digital Sky Survey and the Planck satellite.

SDSS

$$p(x \mid \theta)$$

$$p\left( \boxed{\text{DESI}} \;\middle|\; \boxed{\text{Cosmology parameters}} \right)$$

*Simplify this*      *and/or this*

Intractable!

# Emulation and inference



$$\sim p(\ \blacksquare\ |\ \text{Cosmology}\ )$$

*Emulation/ sampling*

$$p(\ \begin{array}{c}\text{Cosmology}\\ \Omega_m, \sigma_8\end{array}\ |\ \blacksquare\ )$$

$\nabla_{\{\Omega_m,\sigma_8\}}p$

*Differentiable likelihood → even better!*
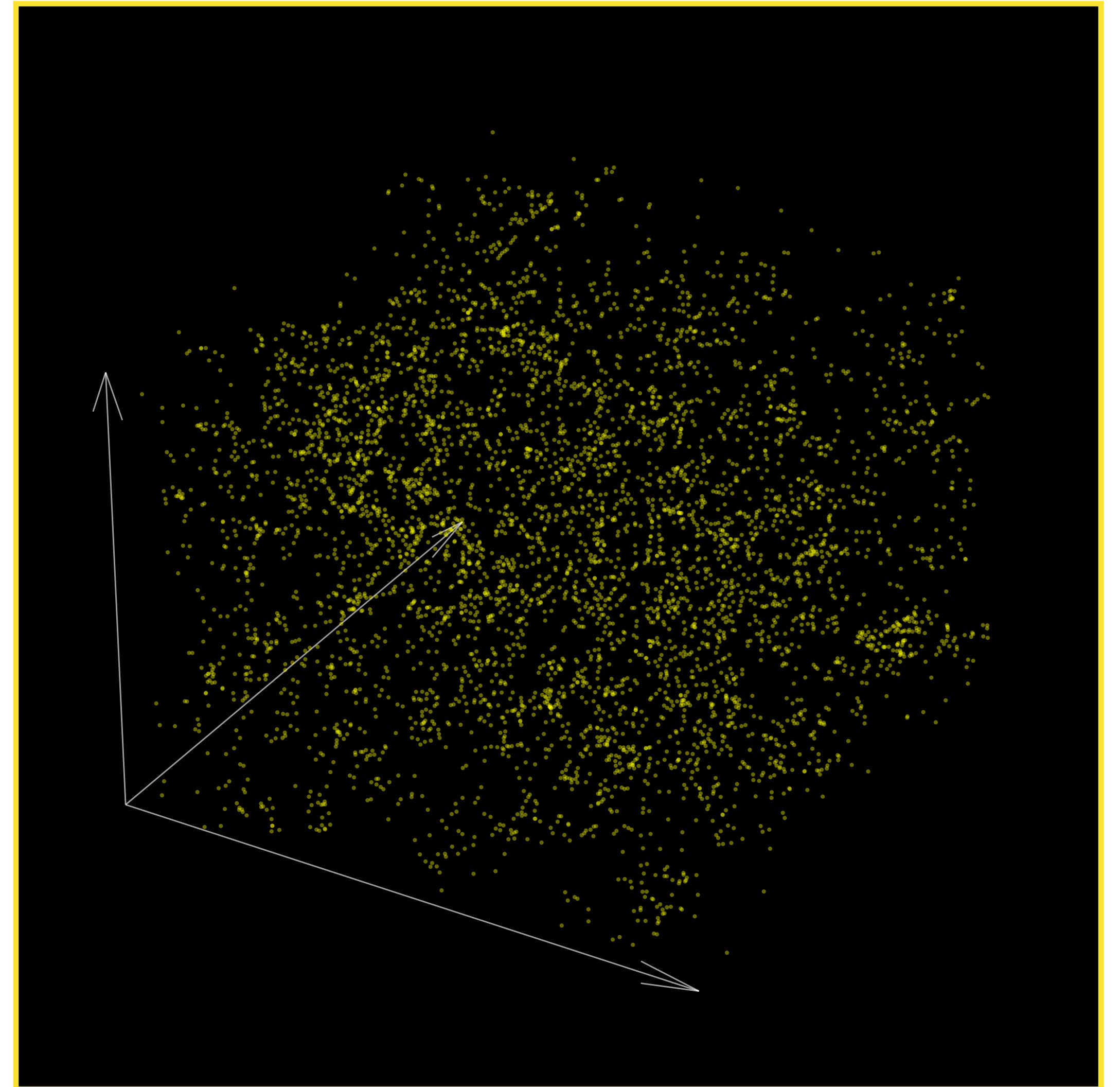


*Parameter estimation*

# The diffusion score model

## Want a score model that

- Operates on *sets* of varying cardinality
- Is permutation equivariant
- Efficiently captures correlation structure of point cloud
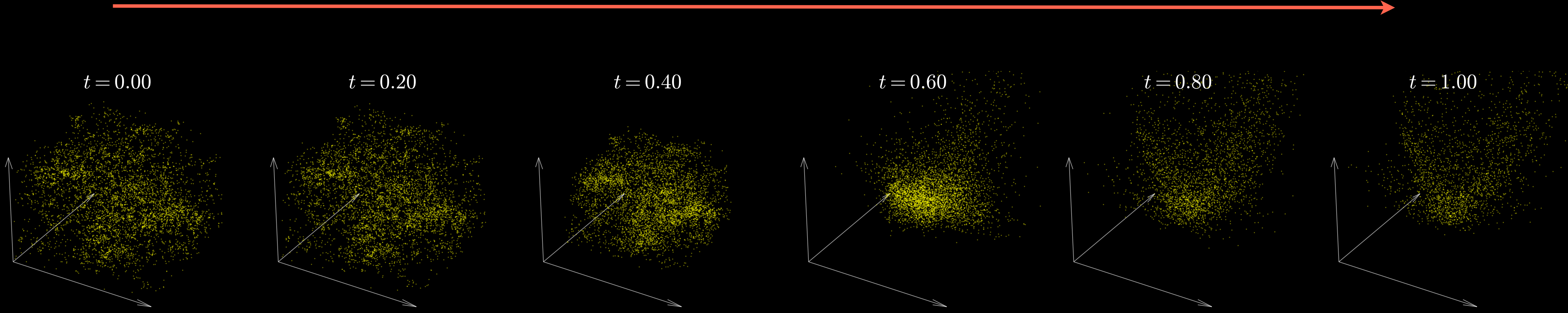
✓ Tranformers
✓ Graph neural networks

# *Graph neural network-guided diffusion* on galaxies

Samples from *Quijote* [Villaescusa-Navarro et al, 2021].

Forward process *(adding noise)*

$t = 0.00$      $t = 0.20$      $t = 0.40$      $t = 0.60$      $t = 0.80$      $t = 1.00$



Reverse process *(denoising)*
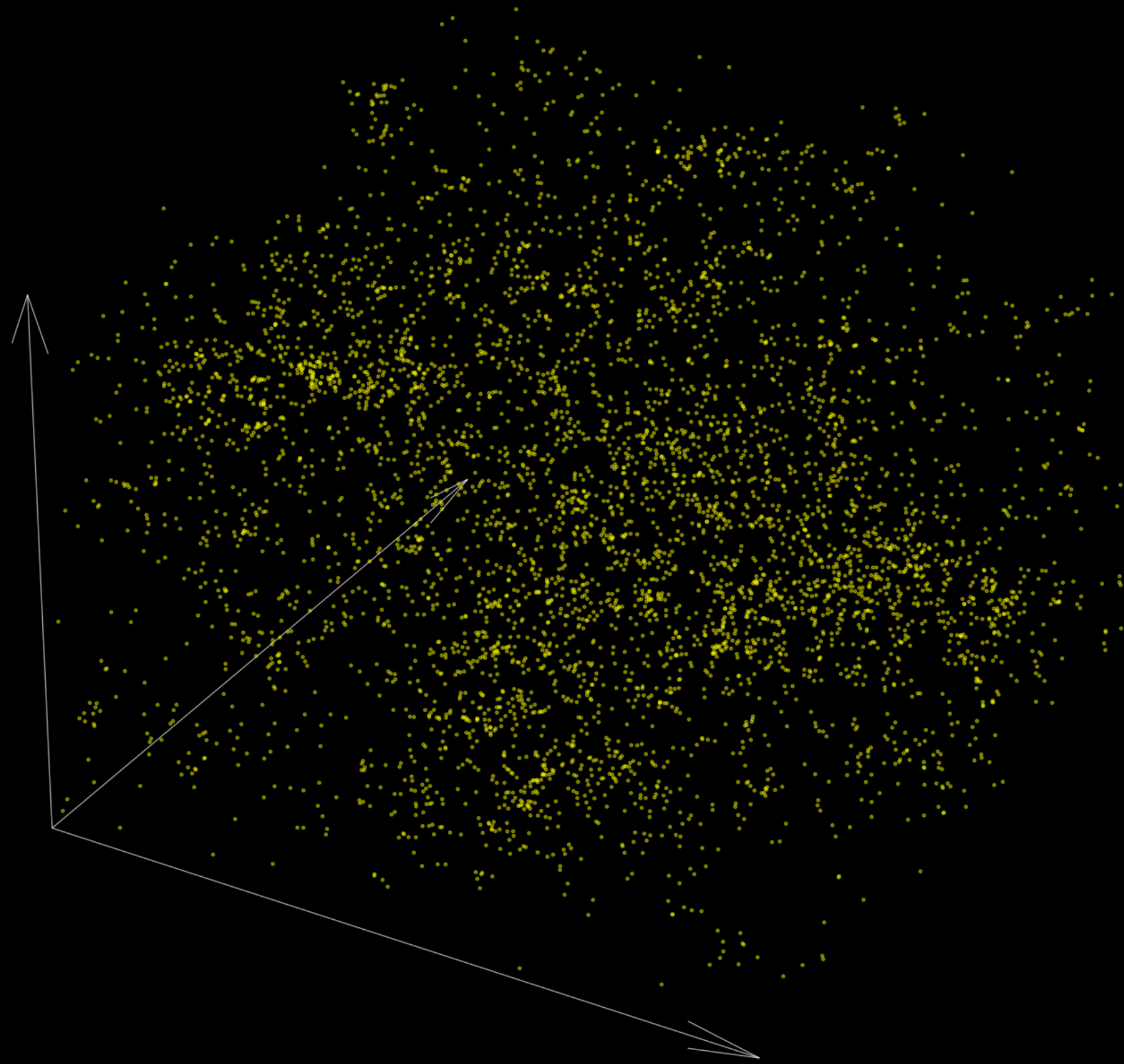
Cosmology parameters $\{\Omega_m, \sigma_8, \ldots\}$
+ Galaxy-halo parameters
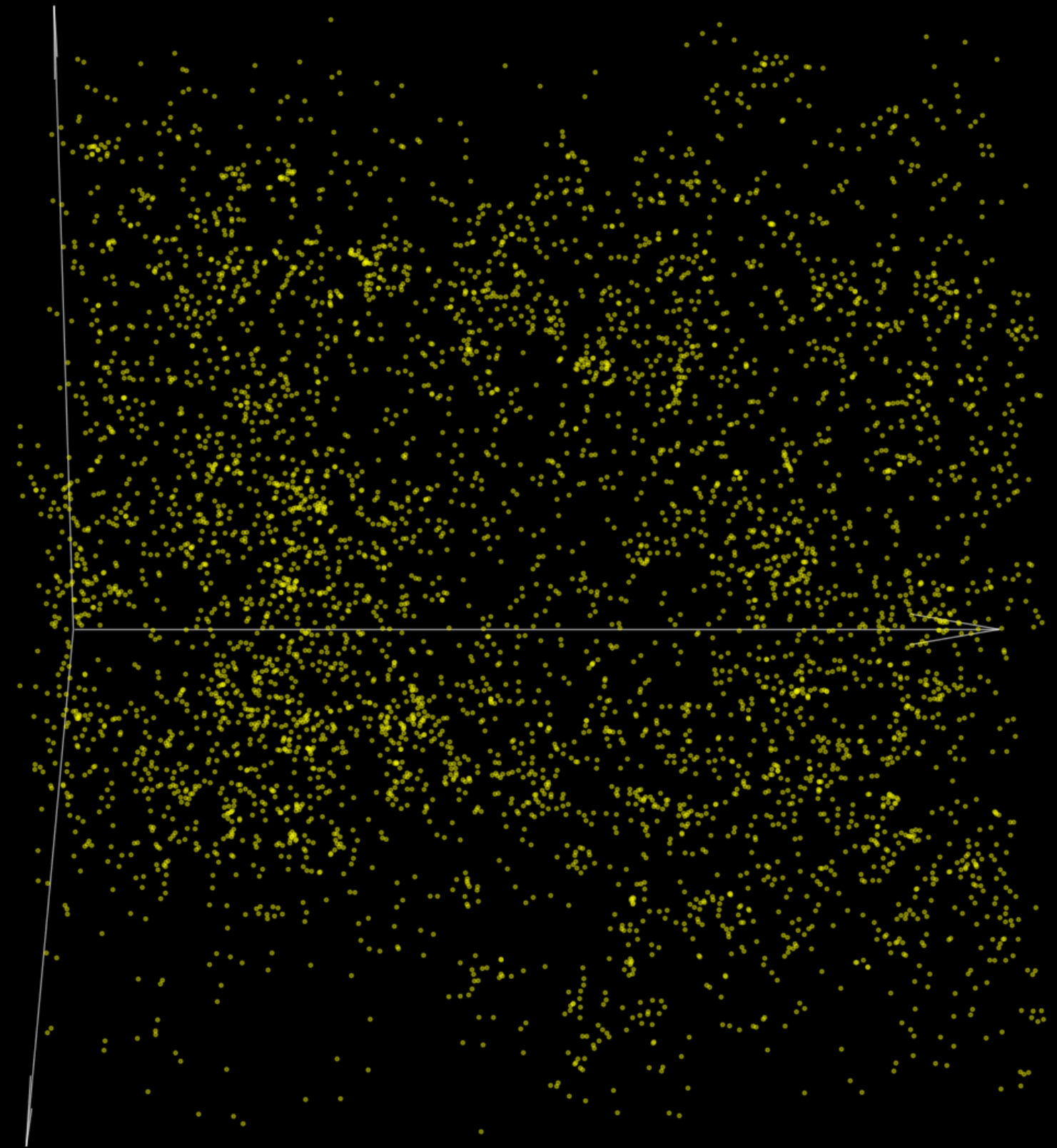
# Diffusion on galaxies

Diffusion process

Conditional generation $x \sim p(x \mid \Omega_m, \sigma_8)$
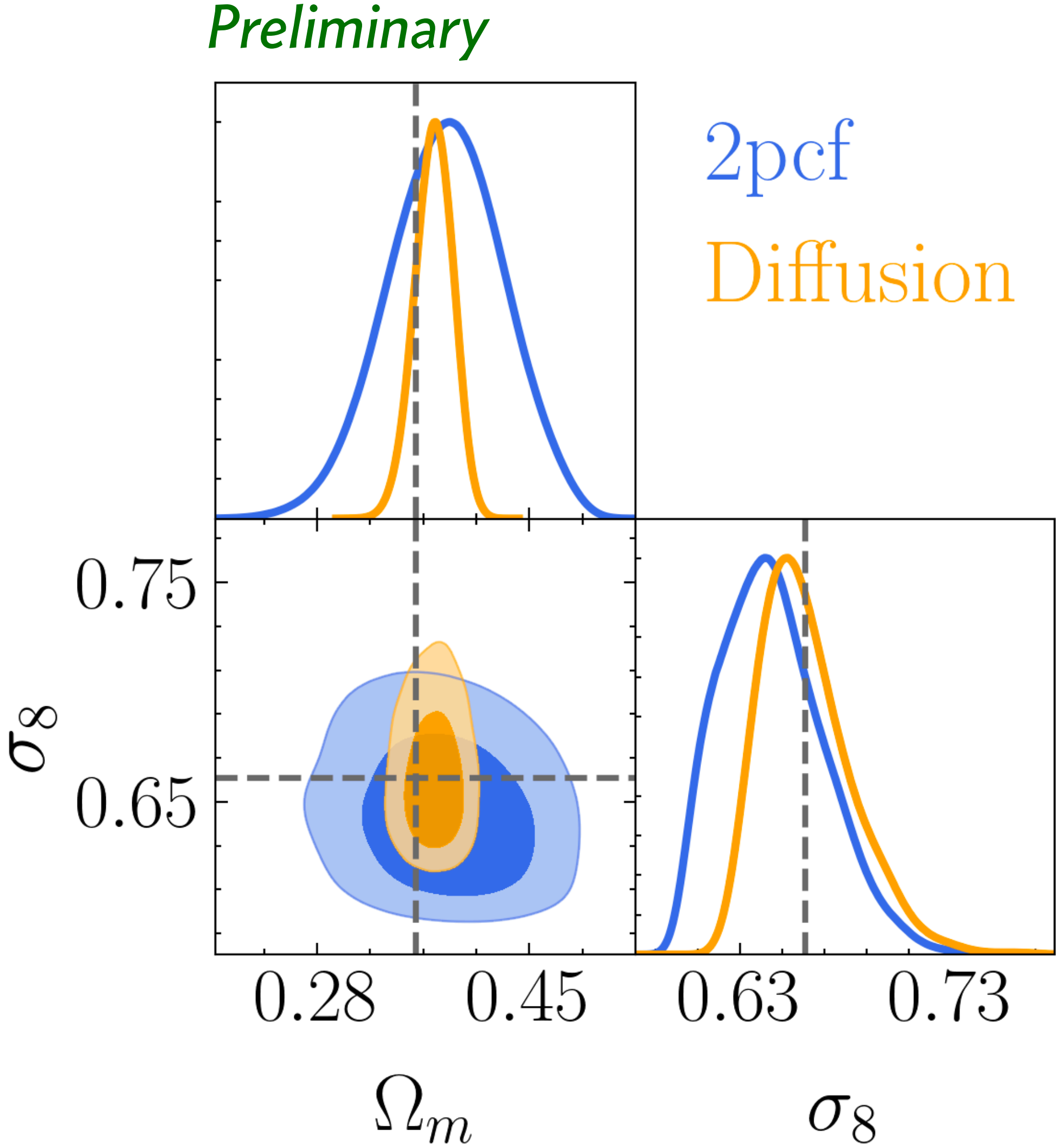
$t = 0.00$

$\Omega_m = 0.10, \sigma_8 = 0.60$
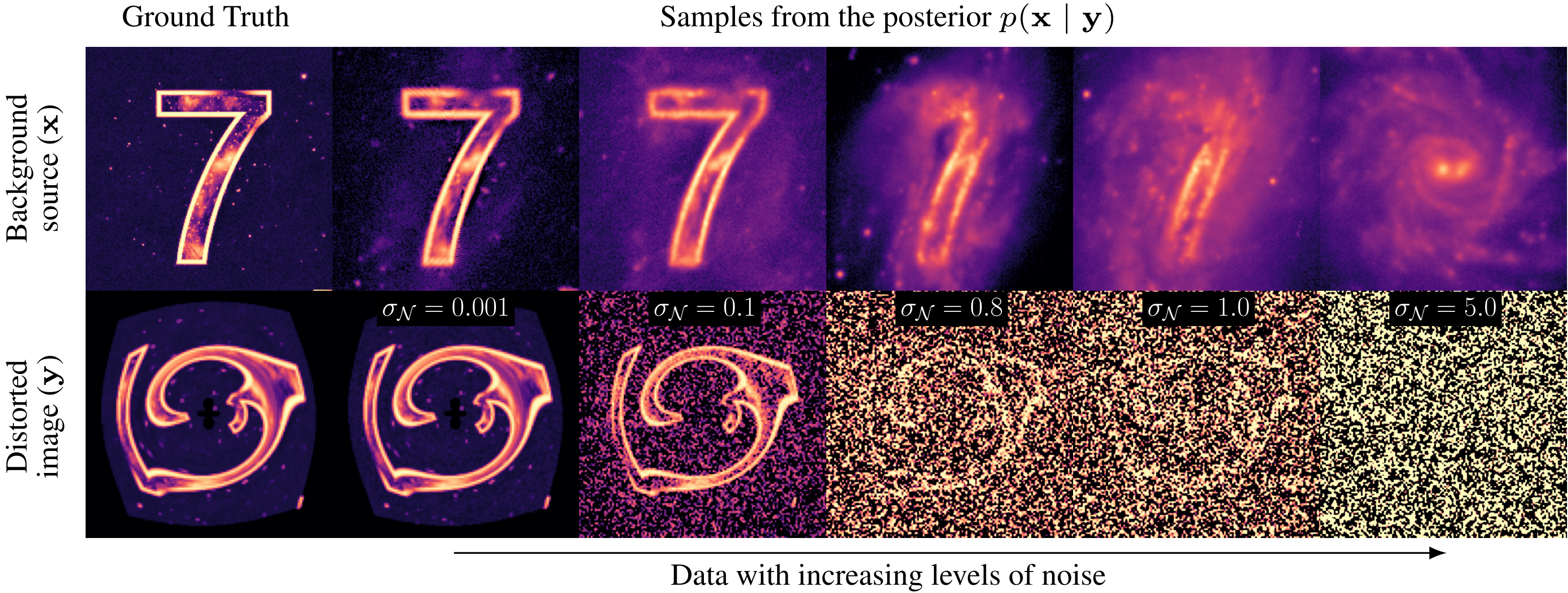
# Likelihoods and parameter inference

For a given dataset, can use the likelihood $p(x \mid \theta)$ for posterior parameter inference

- Monte Carlo sampling (MCMC, nested sampling, HMC…)
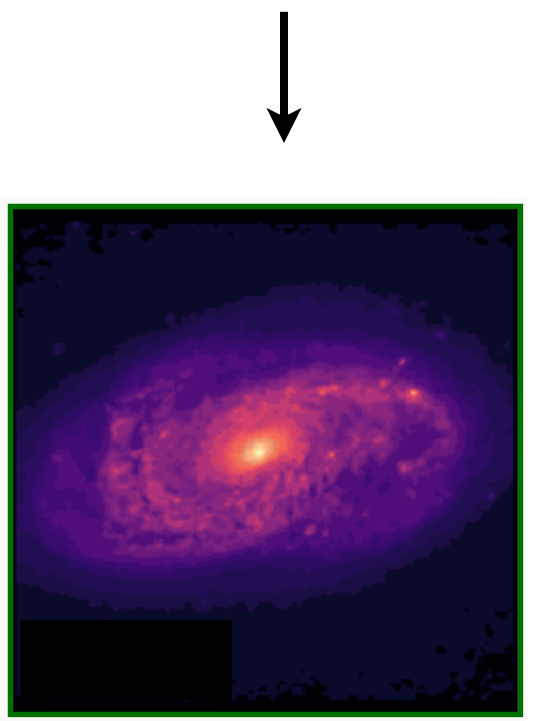- Variational inference



*Preliminary*

2pcf
Diffusion

# Another application: *as a galaxy prior for gravitational lensing*
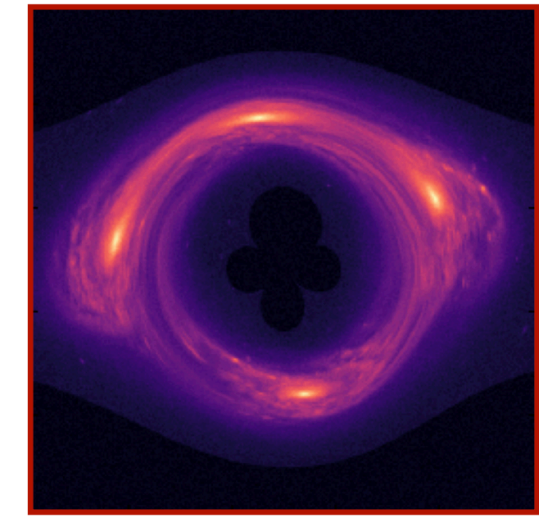
[Adam, Coogan, Malkin et al 2022]
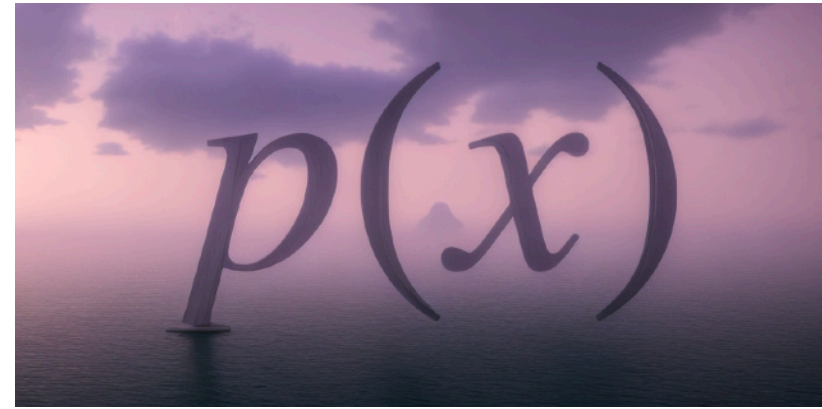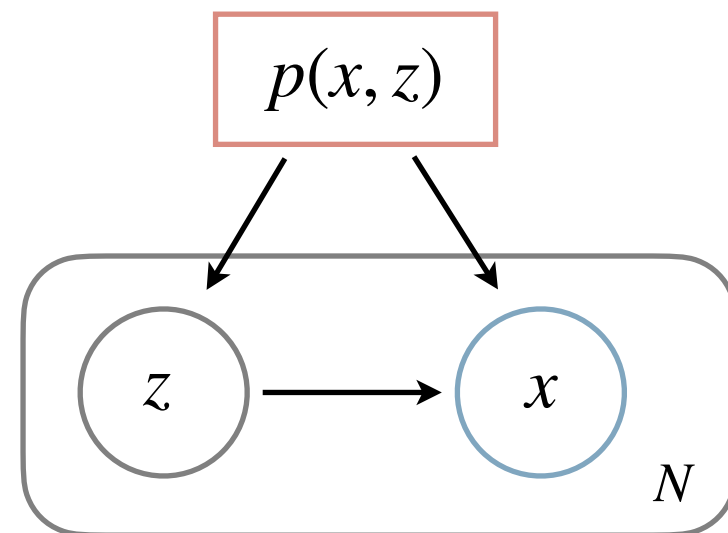


Figure 2: Application of the method to a lensing system with a highly out-of-distribution source. The ground truth is given in the leftmost panel. Other panels show increasingly noisy data (lower row) and a sample from their corresponding source posterior (upper row). As the likelihood becomes less informative, the prior dominates, making the sources increasingly look like galaxies.
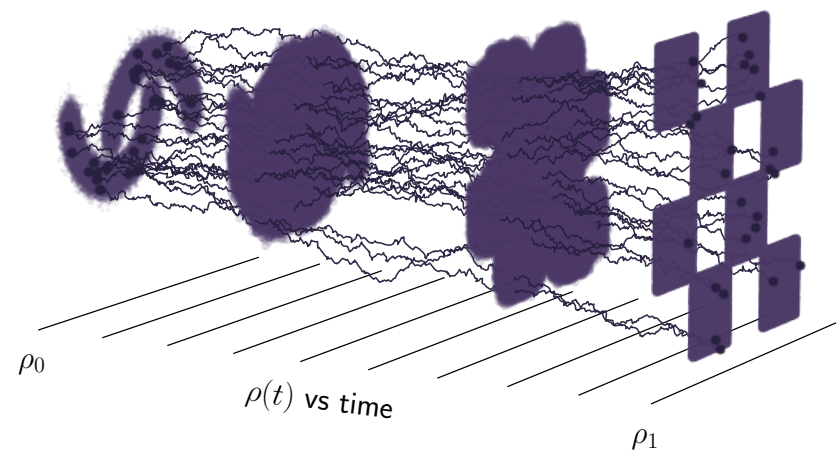
# Outline



## Why (deep) generative modeling?

*What is it, and what can it do for you?*



## Variational auto encoders

*Latent-variable modeling, and compression is all you need*



## Diffusion models

*Models based on iterative refinement*

## Normalizing flows

*Invertible transformations*

# Normalizing flows



Base density

*(IAIFI logo)*

Target density

$$\theta = f(z)$$

One-to-one transformation
Tractable $f^{-1}$ and $\det \nabla f$

$$p_z(z)$$

$$p(\theta) = p_z \left( f^{-1}(\theta) \right) |\det \nabla f|^{-1}$$

*Efficient density estimation:* $\log \hat{p}(\theta)$



*Efficient sampling:* $\theta \sim \hat{p}(\theta)$

# Normalizing flows



The distribution $p(z)$ should

- *Have an easy-to-evaluate density*
- *Be easy to sample from $z \sim p(z)$*

Typically
$$p(z) = \mathcal{N}(0, \mathbb{I})$$

The function $f$ should be

- *One-to-one*
- *Differentiable* } Diffeomorphism
- *Invertible*
- *Tractable $f^{-1}$ and $\det \nabla f$*

- Constrained form of the transformation can limit the expressivity of flows compared to e.g. diffusion models.

- However, for certain physics applications the transformation can be restricted in a specific, desired way; *see Miranda Cheng's lectures on Wednesday!*

# Normalizing flows



$f(z_2)$     $f(z_1)$

*Multiple flow transformation can be easily composed for e.g. expressivity*

$$p(z_T) \longrightarrow \quad z_T \qquad z_{T-1} \quad \cdots \quad z_1 \qquad x \quad \longleftarrow \quad p(x)$$

$f^{-1}(z_1)$     $f^{-1}(x)$

Computing $p(x)$: *change-of-variables formula*

$$\int p(x)dx = \int p(z)dz = 1$$

$$p(x) = p(z)\left|\frac{dz}{dx}\right| = p\left(f^{-1}(x)\right)\left|\frac{df^{-1}}{dx}\right| = p\left(f^{-1}(x)\right)|\det \nabla f|^{-1}$$

Train using maximum-likelihood objective

$$\varphi^* = \left\langle \arg\max_{\varphi} p\left(f_{\varphi}^{-1}(x)\right)|\det \nabla f_{\varphi}|^{-1}\right\rangle_{x\sim p(x)}$$

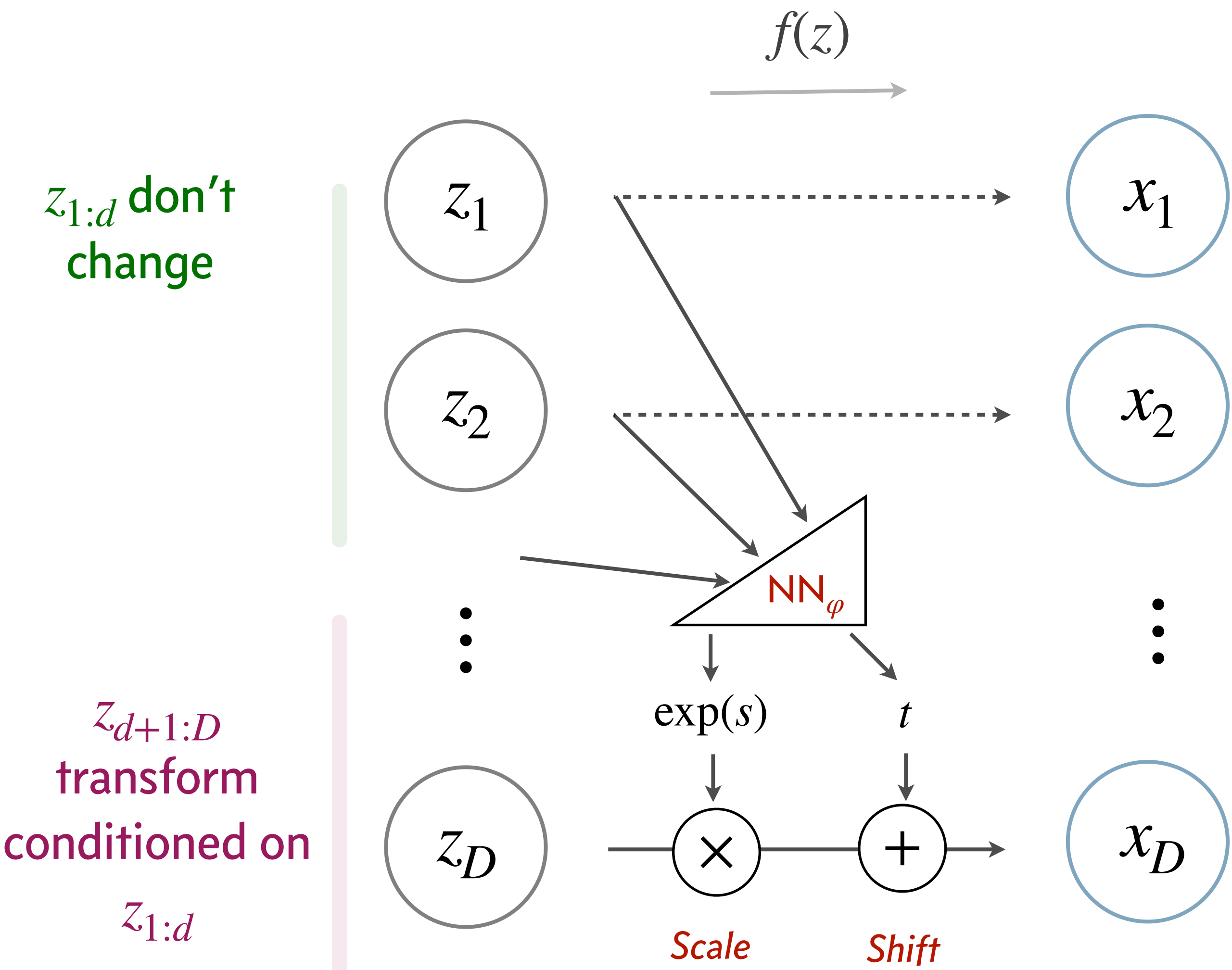# Simple flow transformations

Example: *Affine coupling flow* [RealNVP; Dinh et al 2016]

$f(z)$

$z_{1:d}$ don't change

$z_2$

$z_{d+1:D}$ transform conditioned on $z_{1:d}$

$z_1$

NN$_\varphi$

$\exp(s)$     $t$

$z_D$

$\times$   $+$

*Scale*    *Shift*

$x_1$

$x_2$

$x_D$

*Transformation* ✅

$$x_{d+1:D} = z_{d+1:D} \odot \exp\left( s\left( x_{1:d} \right) \right) + t\left( x_{1:d} \right)$$

*Inverse* ✅

$$z_{d+1:D} = \left( x_{d+1:D} - t\left( x_{1:d} \right) \right) \odot \exp\left( -s\left( x_{1:d} \right) \right)$$

*Jacobian determinant* ✅

$$\det(\nabla f) = \prod_{j=1}^{D-d} \exp\left( s\left( z_{1:d} \right) \right)_j = \exp\left( \sum_{j=1}^{D-d} s\left( z_{1:d} \right)_j \right)$$

+ Switch up order of transformed variables at every transformation

# Continuous-time normalizing flows

**Parameterize the transformation by a neural ODE**



ODE with reversible dynamics

$$\frac{\mathrm{d}x}{\mathrm{d}t} = f\left(x(t)\right)$$

Instantaneous change-of-variable formula

$$\frac{\mathrm{d}\log p(x(t))}{\mathrm{d}t} = -\operatorname{Tr}\left(\frac{\mathrm{d}f}{\mathrm{d}x(t)}\right)$$
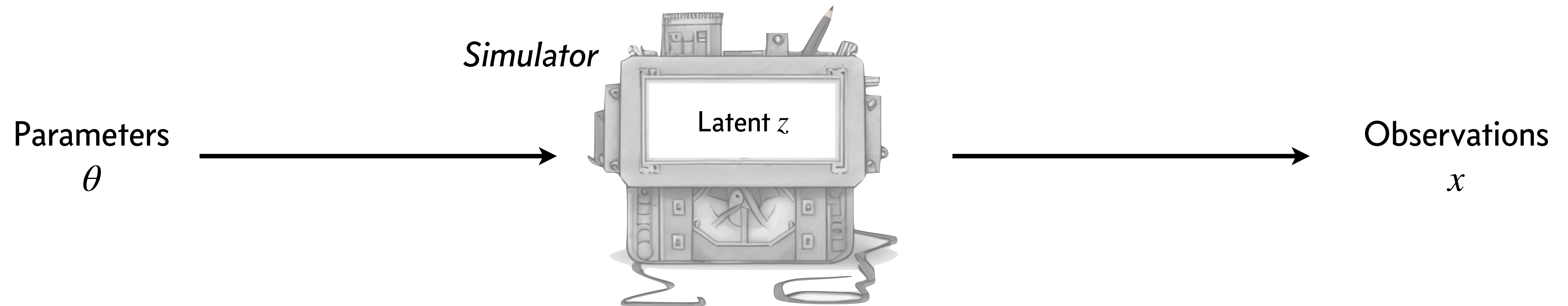
Pro ✅

Unrestricted form of transformation $f(x)$!

Cons 😵‍💫

- Need for efficient trace calculation
- Solving an ODE and backpropping through the solution can make for cumbersome training

# Simulation-based inference (SBI)



*Simulator*

Latent $z$

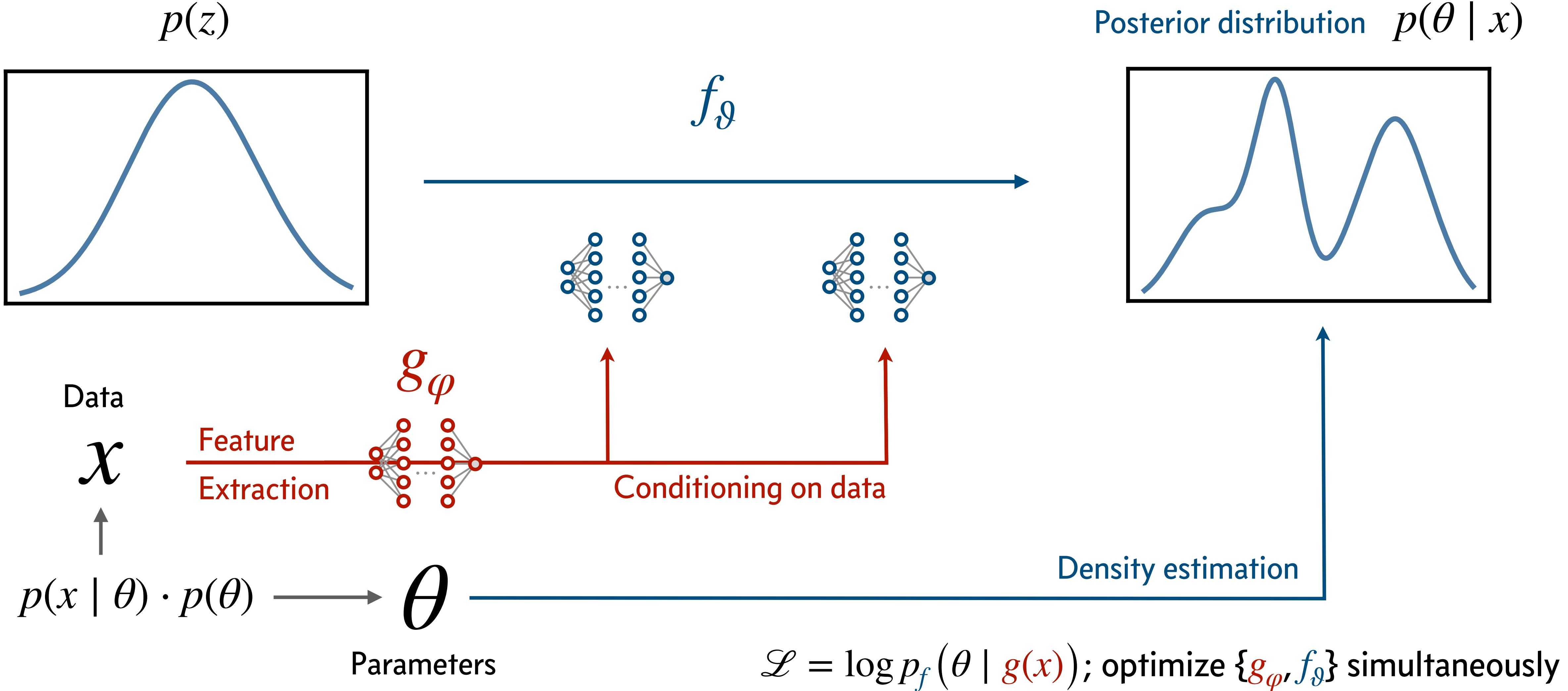Parameters
$\theta$

Observations
$x$

Prediction:
- Well-motivated mechanistic, causal model
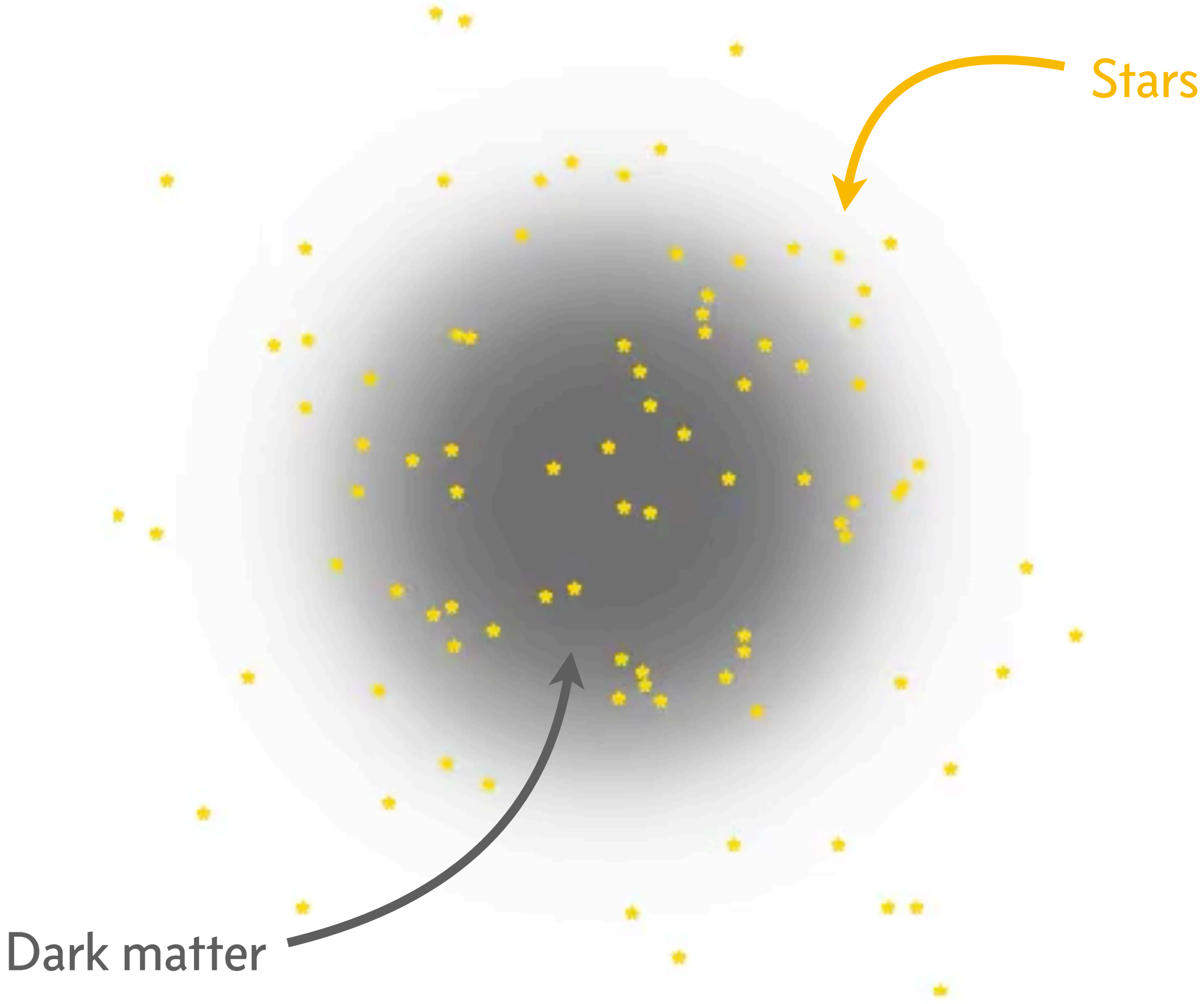- Simulator can generate samples $x \sim p(x \mid \theta)$

Inference:
- Likelihood $p(x \mid \theta) = \int \mathrm{d}z \, p(x, z \mid \theta)$ is intractable
- *Inference is challenging*

# Flows in *simulation-based inference*

Flows are commonly employed as *conditional posterior density estimators* in simulation-based inference

$p(z)$

$f_\vartheta$

Posterior distribution $p(\theta \mid x)$

$g_\varphi$

Data

$x$

Feature Extraction

Conditioning on data

$p(x \mid \theta) \cdot p(\theta)$

$\theta$

Parameters

Density estimation

$\mathscr{L} = \log p_f\left(\theta \mid g(x)\right)$; optimize $\{g_\varphi, f_\vartheta\}$ simultaneously
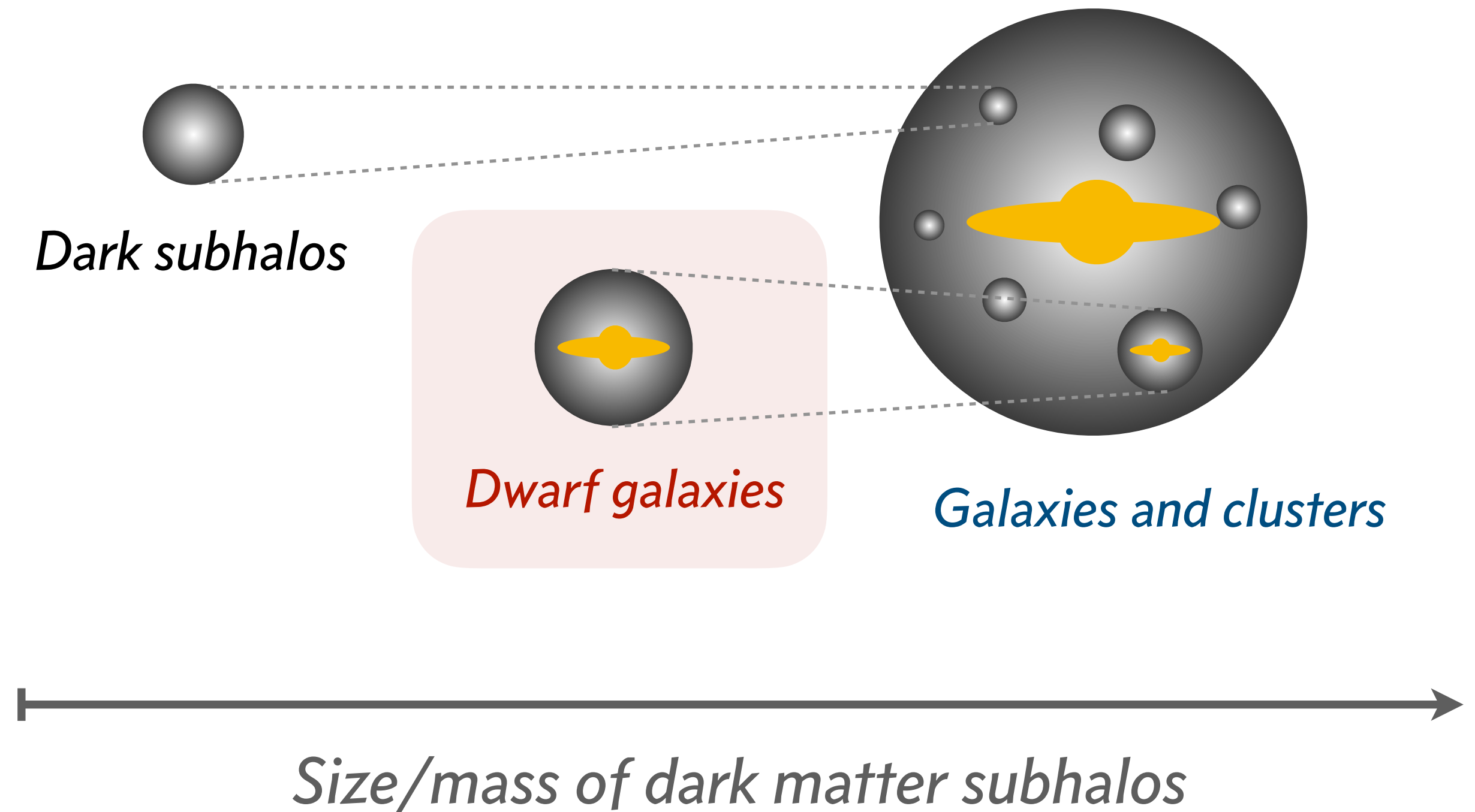
Stars

# An application to *dwarf galaxies*

Dark matter

# Another application: extracting the dark matter distribution from dwarf galaxies

*Dwarf galaxies* are intermediate-sized galaxies well-suited for studying dark matter
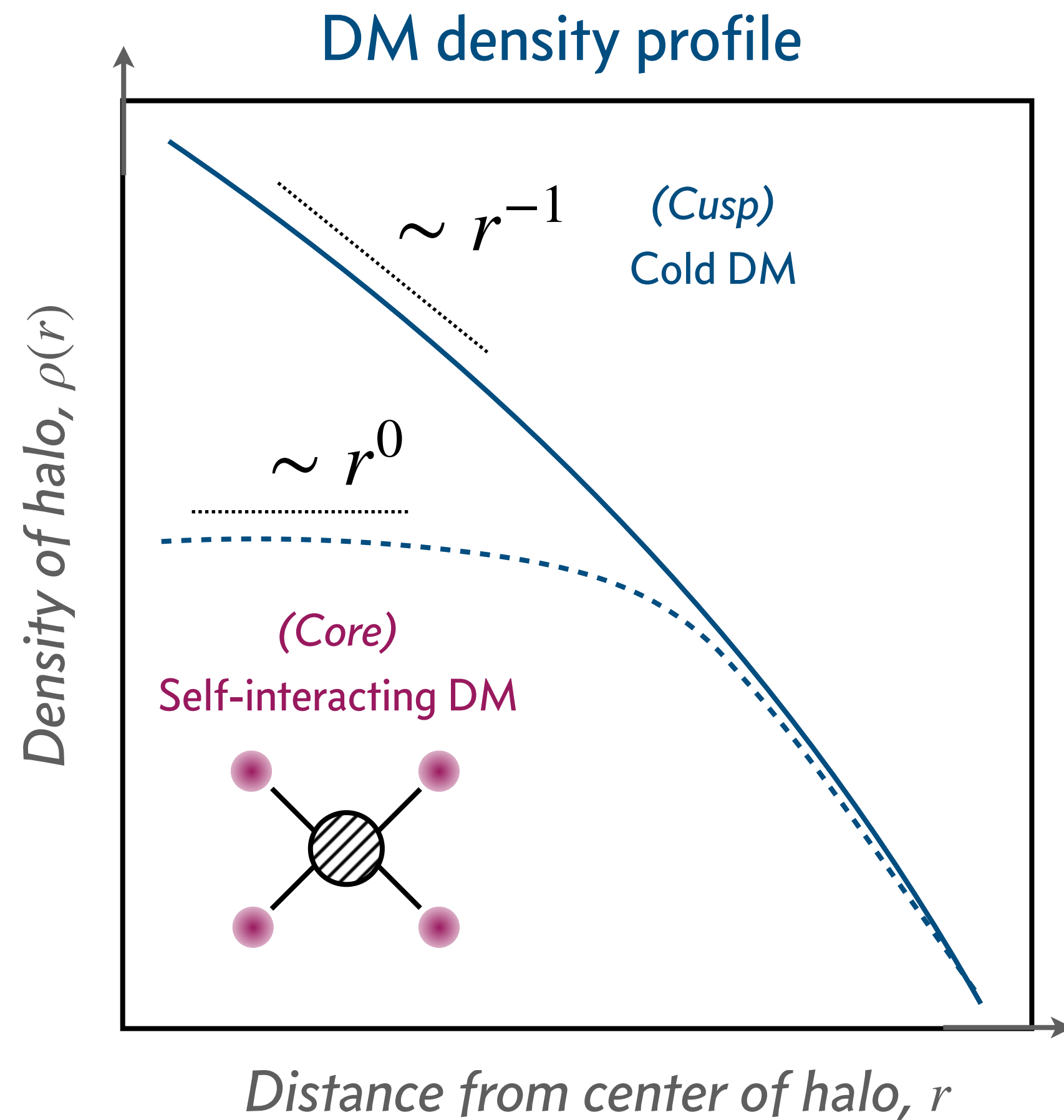
[Nguyen, SM et al 2022]



Dark subhalos

Dwarf galaxies

Galaxies and clusters

Size/mass of dark matter subhalos

*Fornax dwarf galaxy*



*ESO*

# Extracting the dark matter distribution from dwarf galaxies

## DM density profile

$\sim r^{-1}$ *(Cusp)*
Cold DM

$\sim r^0$

*(Core)*
Self-interacting DM

Density of halo, $\rho(r)$

Distance from center of halo, $r$

Prediction
(🤗 with a simulator!)

Inference
(💀)

# Extracting the dark matter distribution   [Nguyen, SM et al 2022]



$\theta \sim p(\theta \mid x)$

Conditional flow transformation $\theta = f(u)$

Base distribution $\mathcal{N}(u)$

Posterior distribution $p(\theta \mid x)$

$\gamma$

$r_{\mathrm{s}}$

$\rho_0$

DM + stellar parameters $\theta$

Sim

$\{\vec{r}, v_{\mathrm{los}}\}$

$x$

$g(x)$

Conditioning

Cored profile $\gamma = 0$   Cuspy profile $\gamma = 1$

$\log_{10} \rho \, [\mathrm{M_\odot/kpc^3}]$

Posterior   Truth

$\log_{10}(M/\mathrm{M_\odot})$

$\log_{10}(r/r_\star)$   $\log_{10}(r/r_\star)$

# Back to diffusion: the *probability flow ODE*

For any diffusion process, there exists a corresponding deterministic process whose trajectories share the same marginal probability densities $p(x_t)$ as the SDE [Song et al 2021]



Probability flow ODE

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

*Diffusion mo...*
*"simulation-f...*
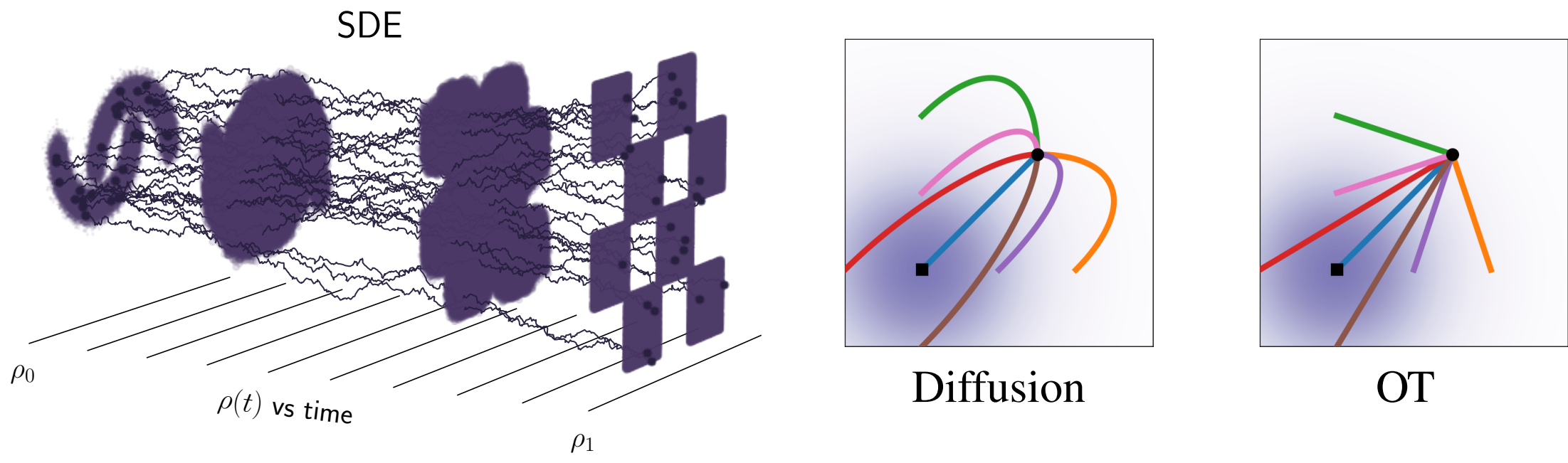
*Wait, it's just a normalizing flow?*

*Always has been*

# Iterative refinement, interpolants, consistency, …

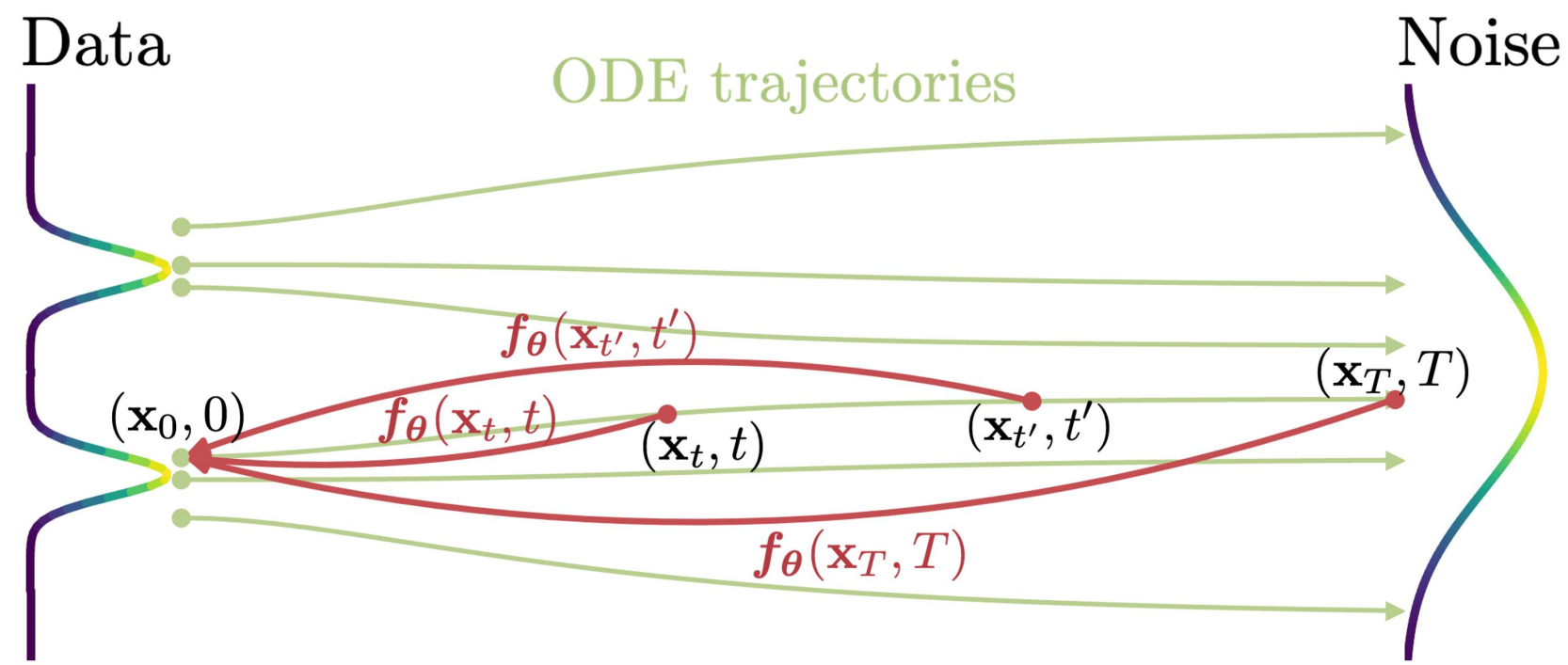**Empirical success of diffusion models has led to many new formulations and methods of training**

### Stochastic interpolants and flow matching

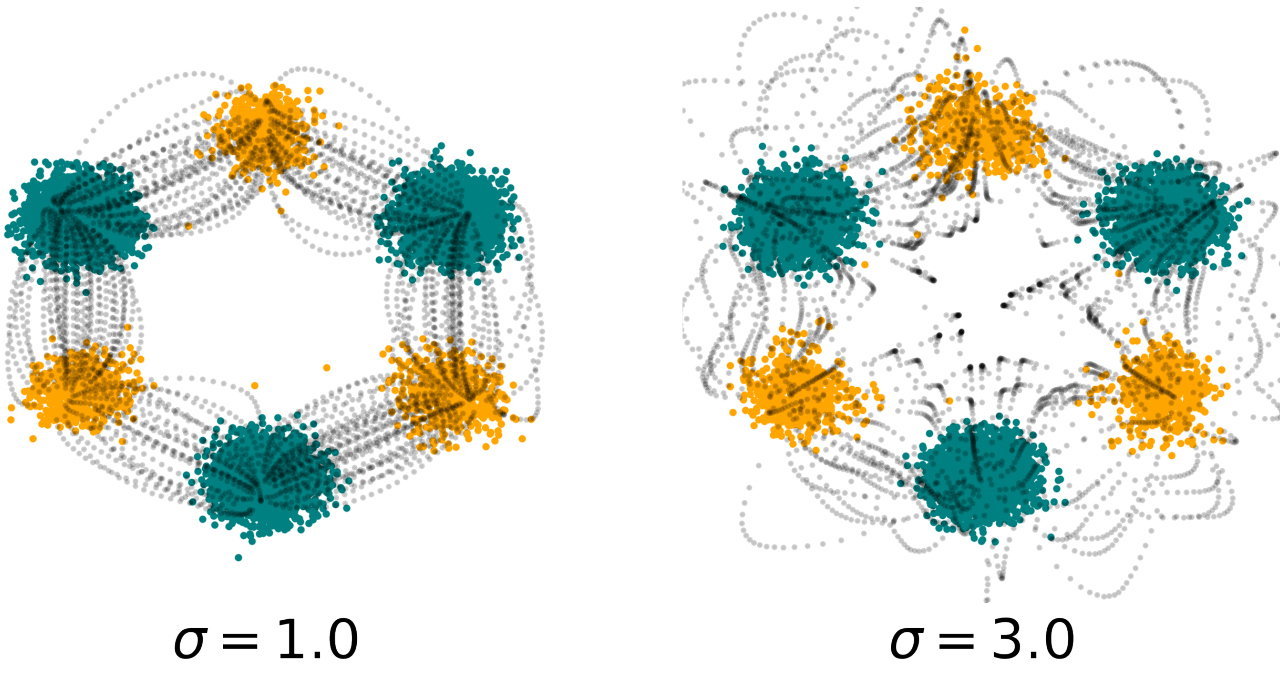[Albergo et al 2023, Lipman et al 2022, …]



### ODE trajectory consistency, Fokker-Planck regularization, …

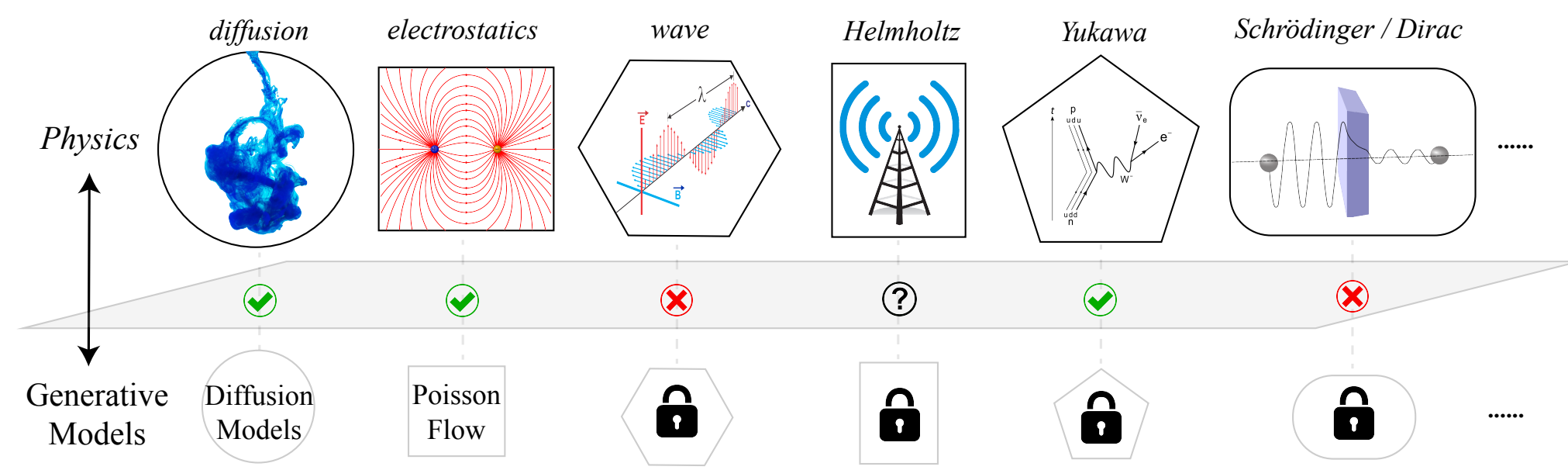[Song et al 2023, Lai et al 2023, …]



### Diffusion Schrödinger's bridges

[Shi et al 2023, …]



### Generative models inspired by other physical processes

[Xu et al 2022, Liu et al 2023, …]

End.