

Generative modeling

Siddharth Mishra-Sharma

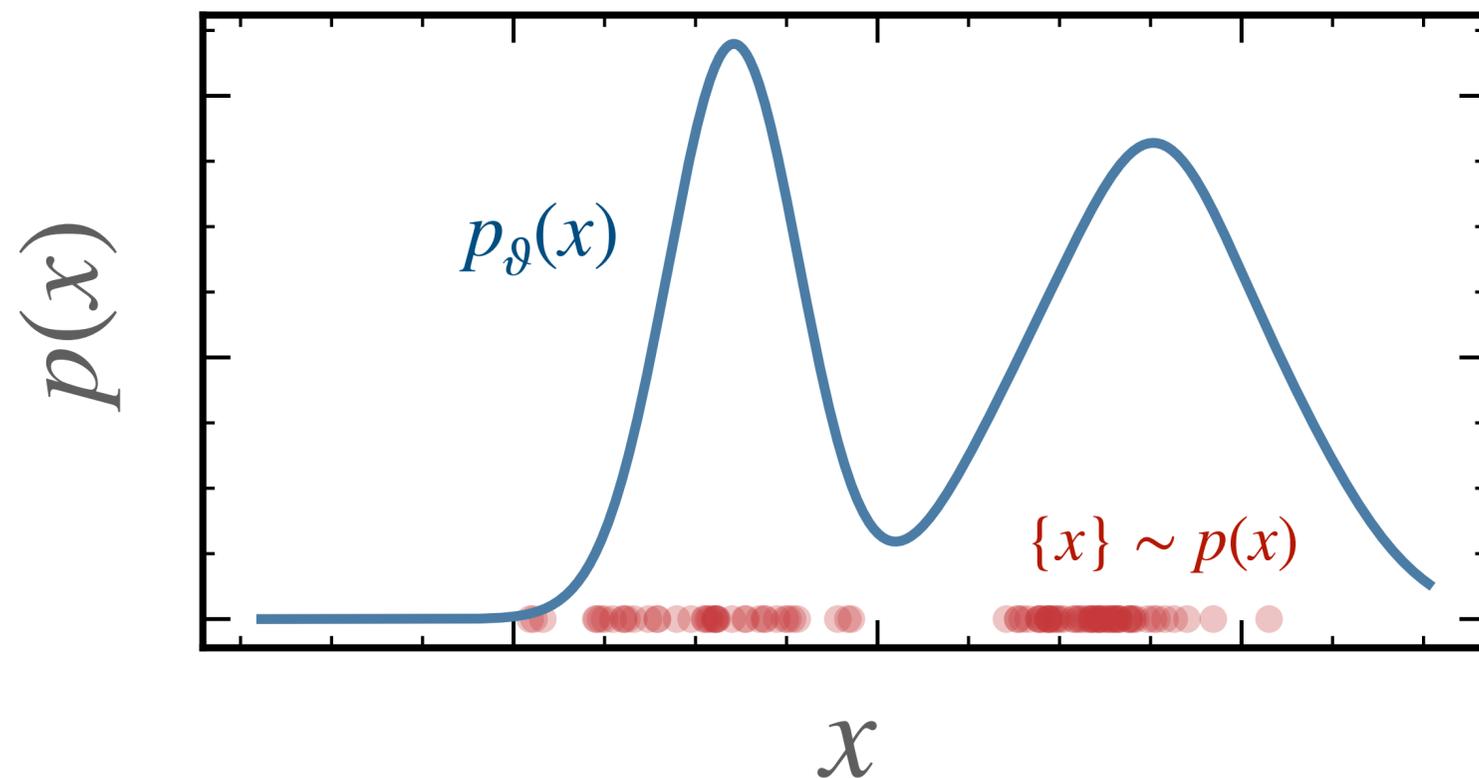
CDS DS 595

<https://smsharma.io/teaching/ds595-ai4science.html>

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 $p(x)$



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*)



Seedance 2.0 (Feb 2026)

https://x.com/lexx_aura/status/2022994297881280557?s=20



The landscape of deep generative models

[Karsten Kreis; [CVPR 2022 Tutorial](#)]



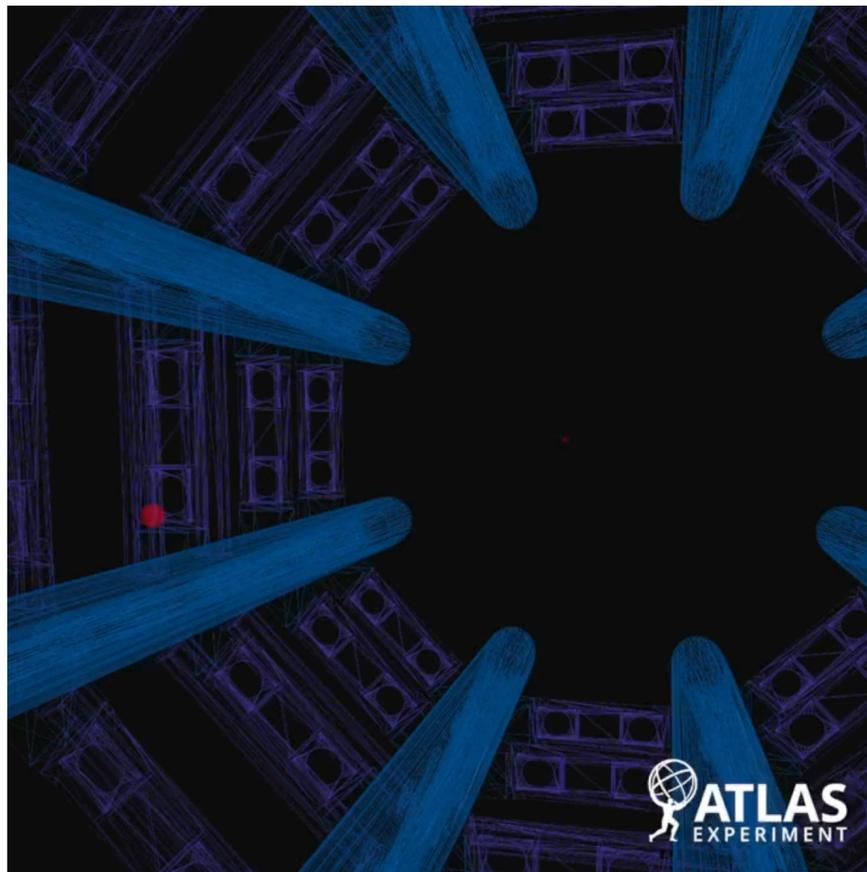
Simulators

$$x \sim p(x)$$

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

Collider data

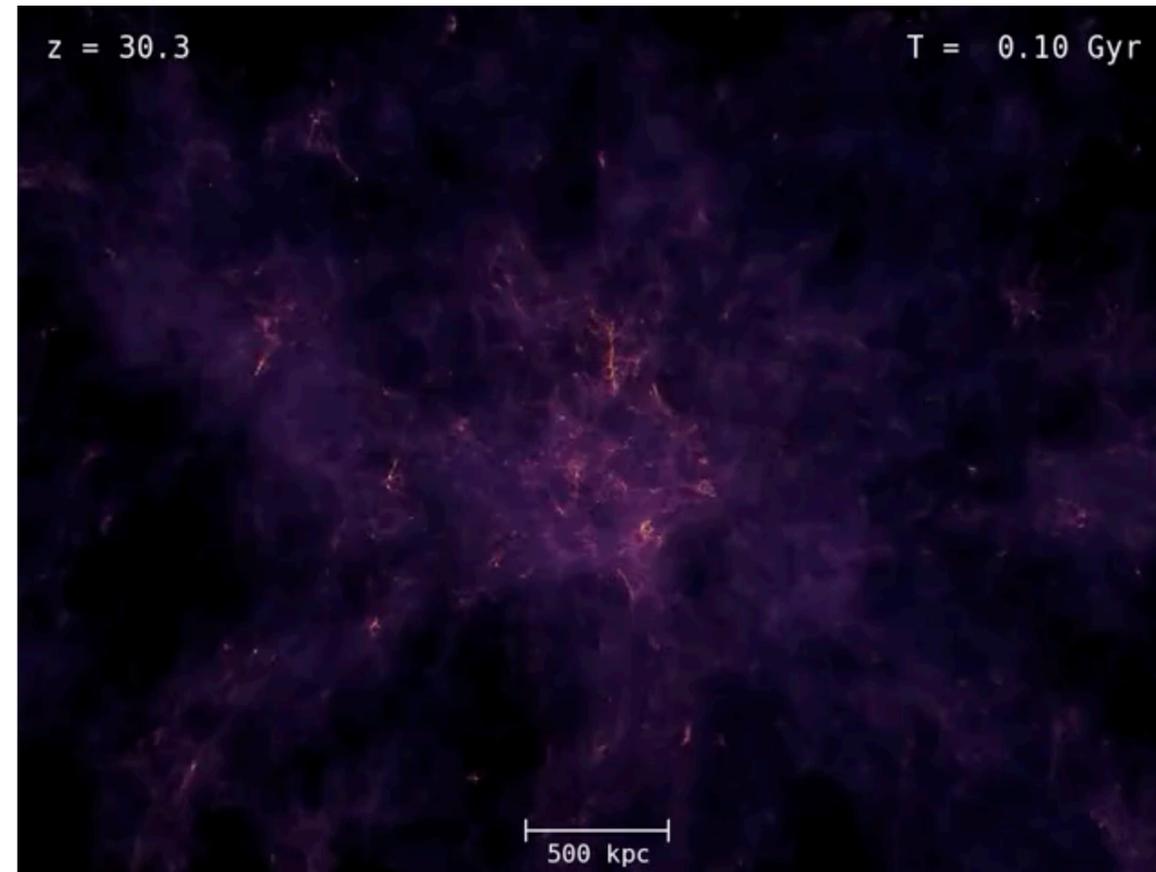
particles $\sim p(\text{particles})$



[C. Cesarotti with ATLAS]

Cosmology data

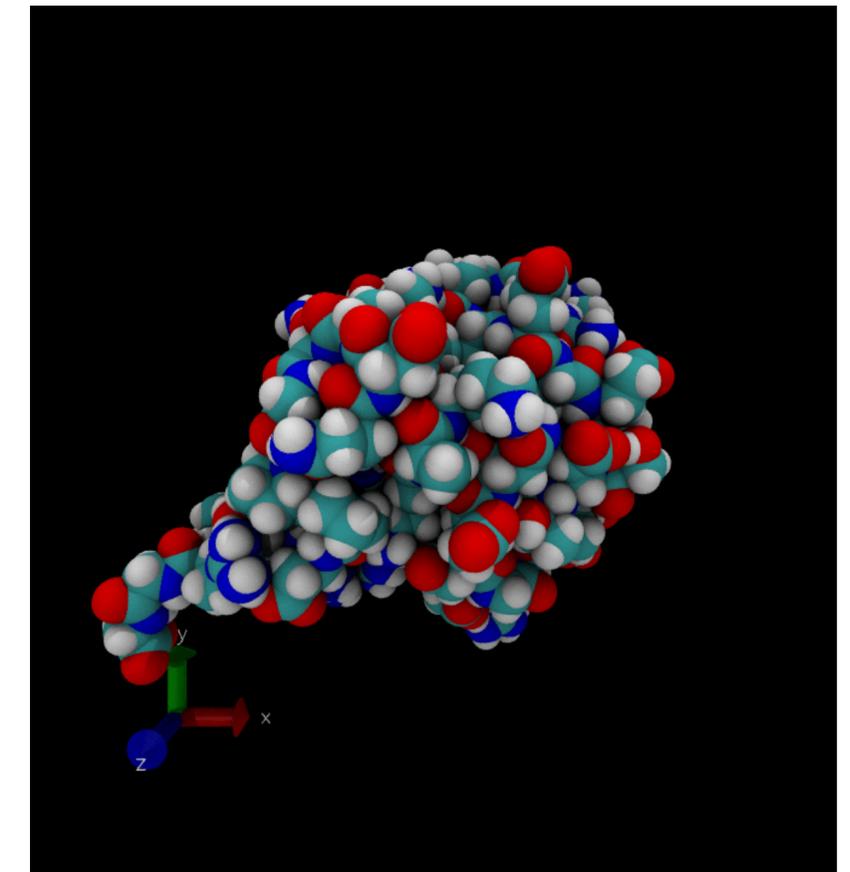
particles $\sim p(\text{particles})$



[Aquarius simulation]

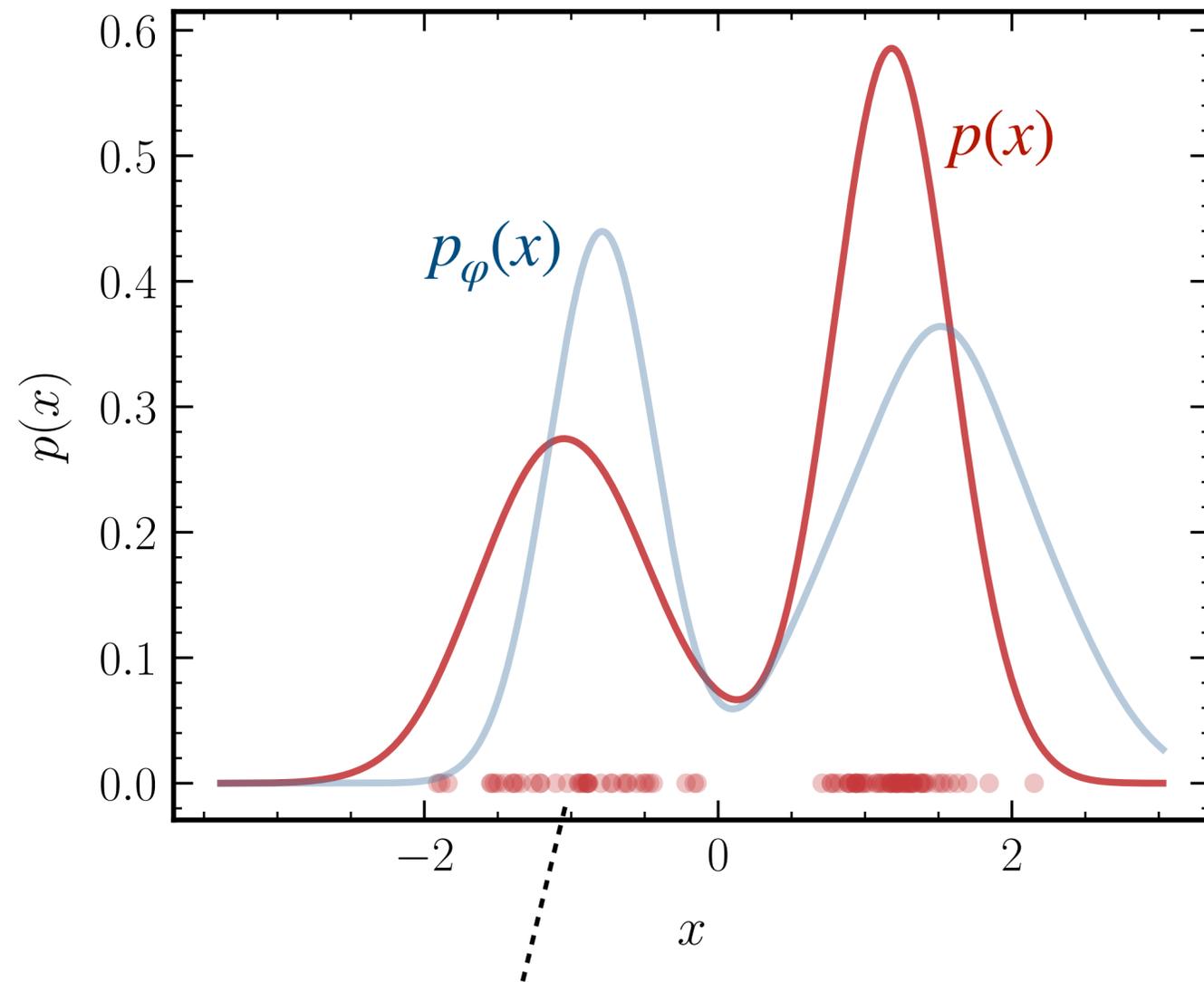
Molecular dynamics

configurations $\sim p(\text{configurations})$



[E. Cancès et al]

Learning the data distribution



$$\{x\}_{\text{train}} \sim p(x)$$

1. Ingredients:

- A parameterized distribution $p_\phi(x)$
- 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

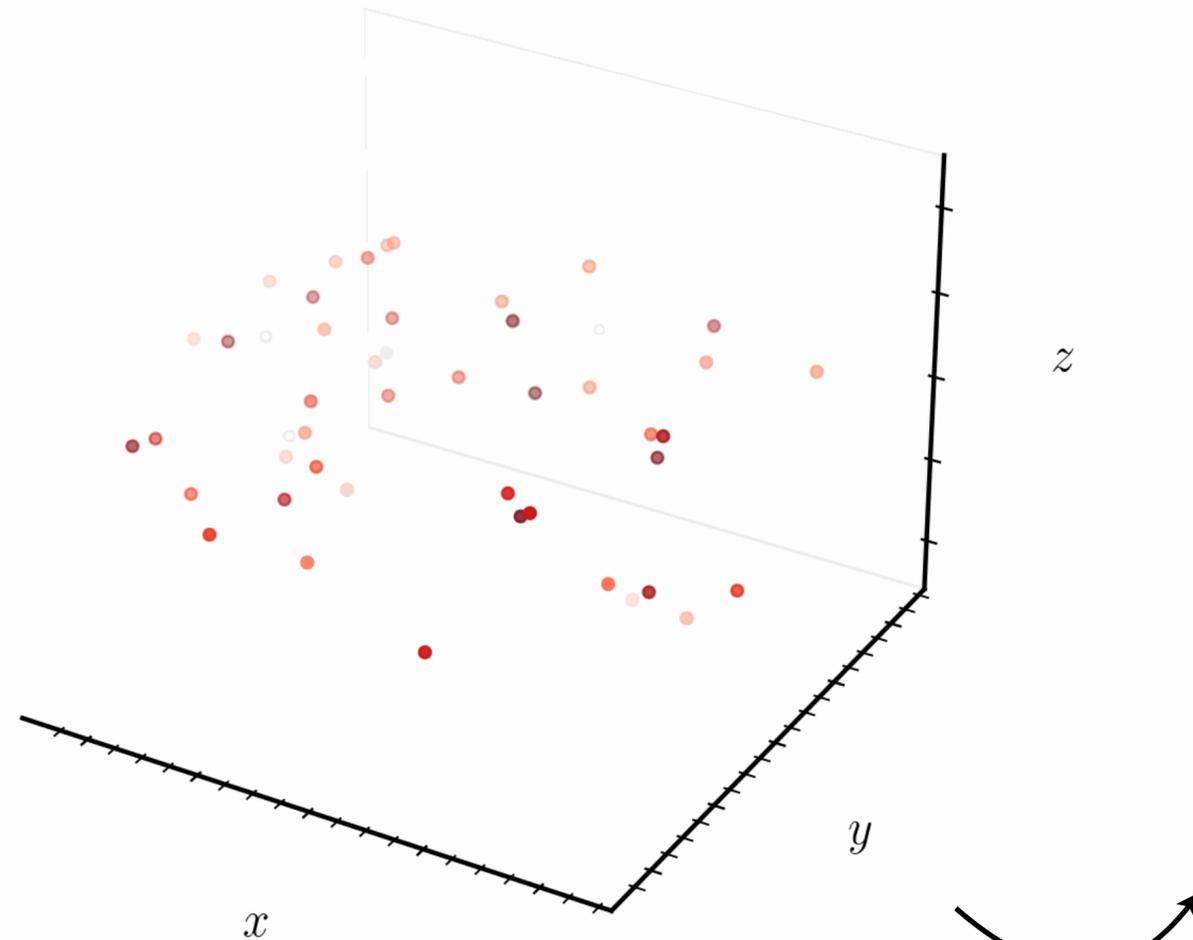
$$\hat{\phi} = \arg \max_{\phi} \left[\log p_\phi(\{x\}_{\text{train}}) \right]$$

Not so fast...

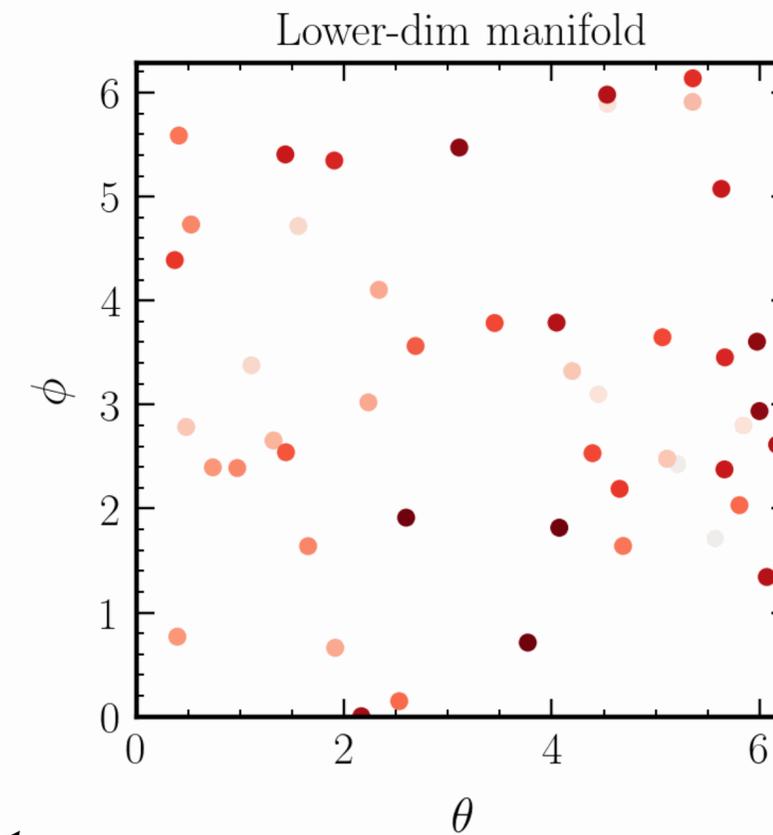
The pursuit of low-dimensional structure

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

“Difficult to model” x

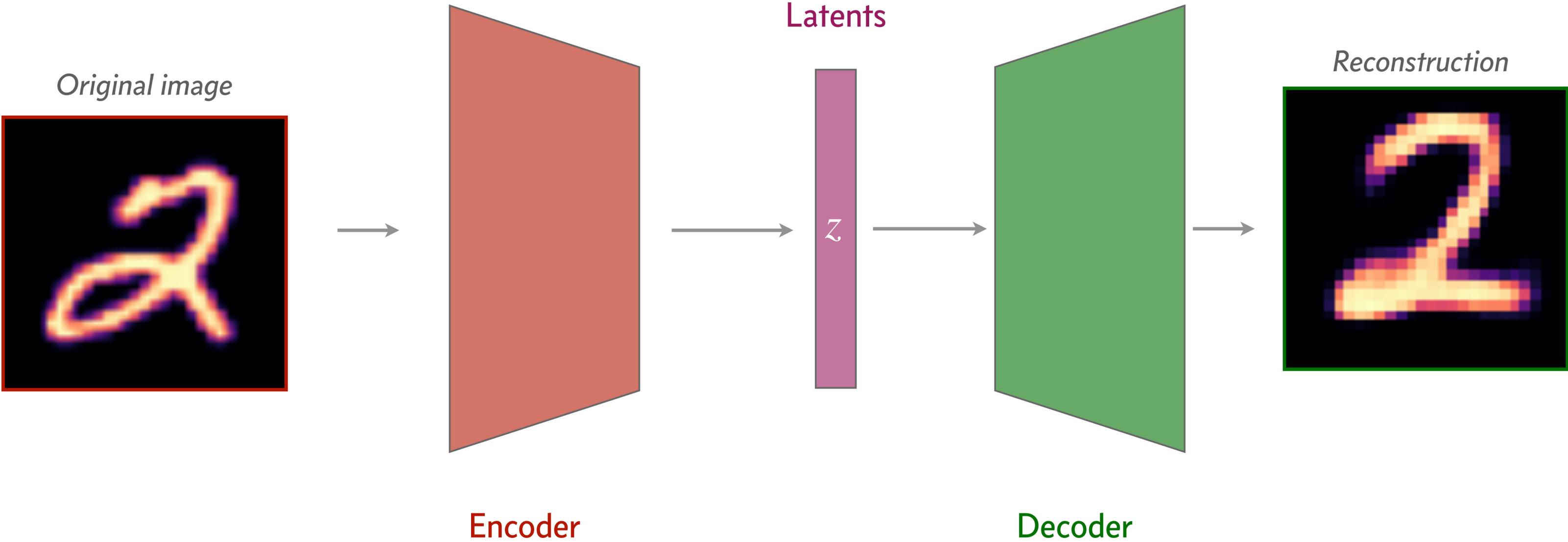


“Easy to model” z



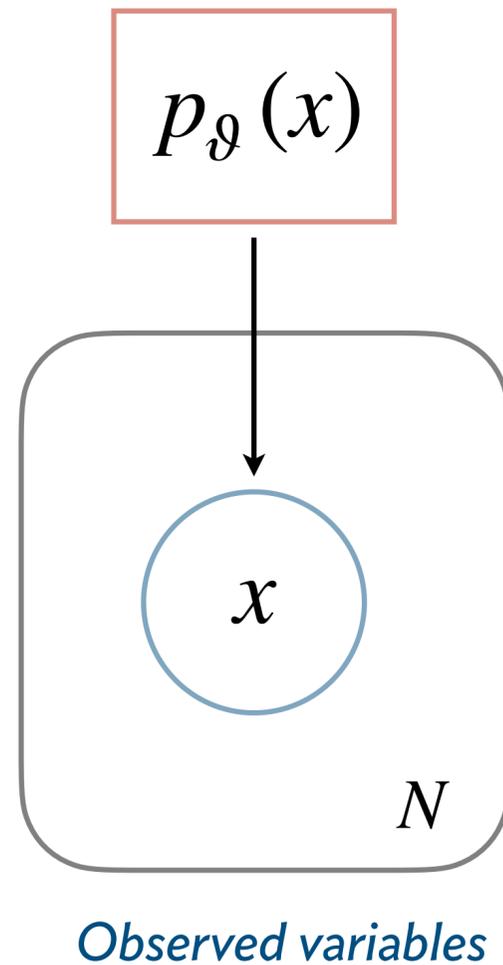
$p(z | x)$

Autoencoders

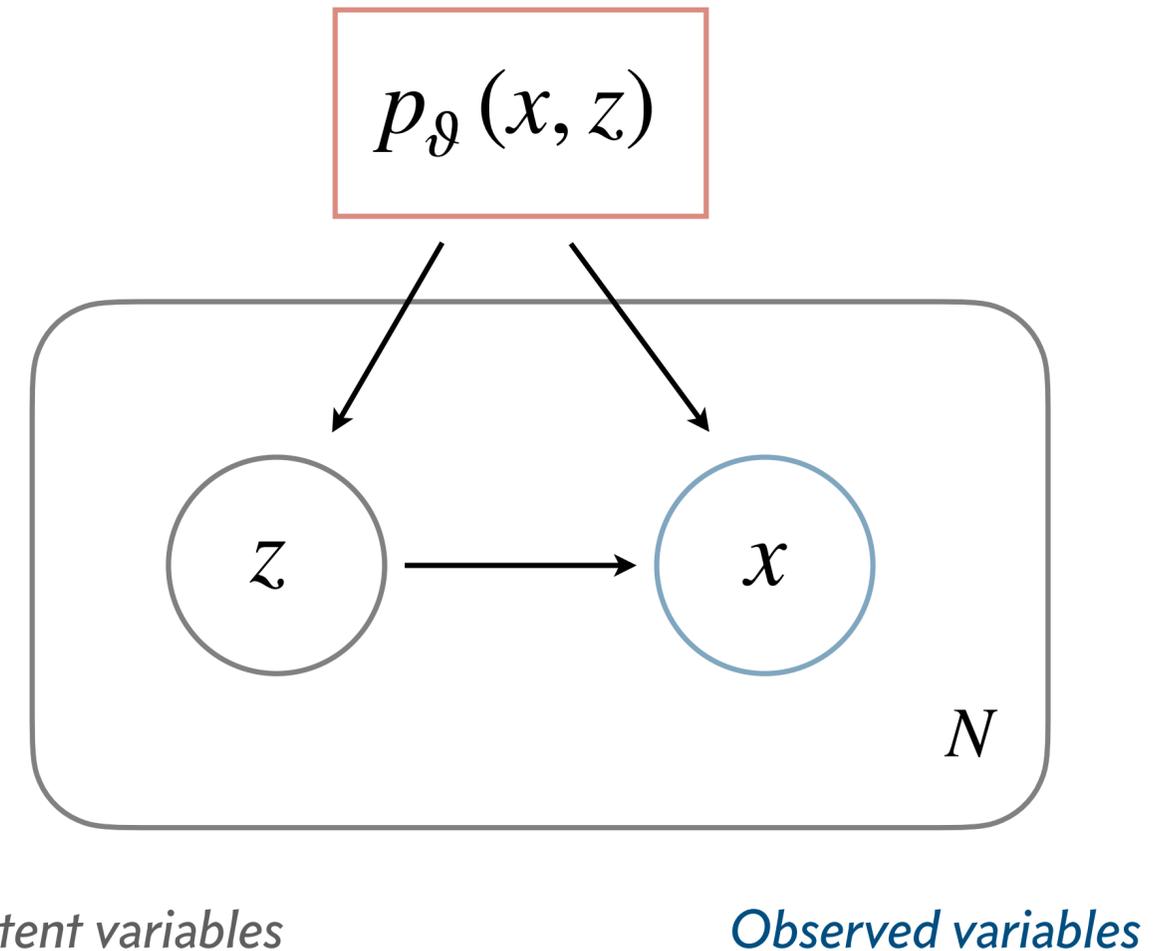


Latent-variable modeling

Learn lower-dimensional structure in the data distribution



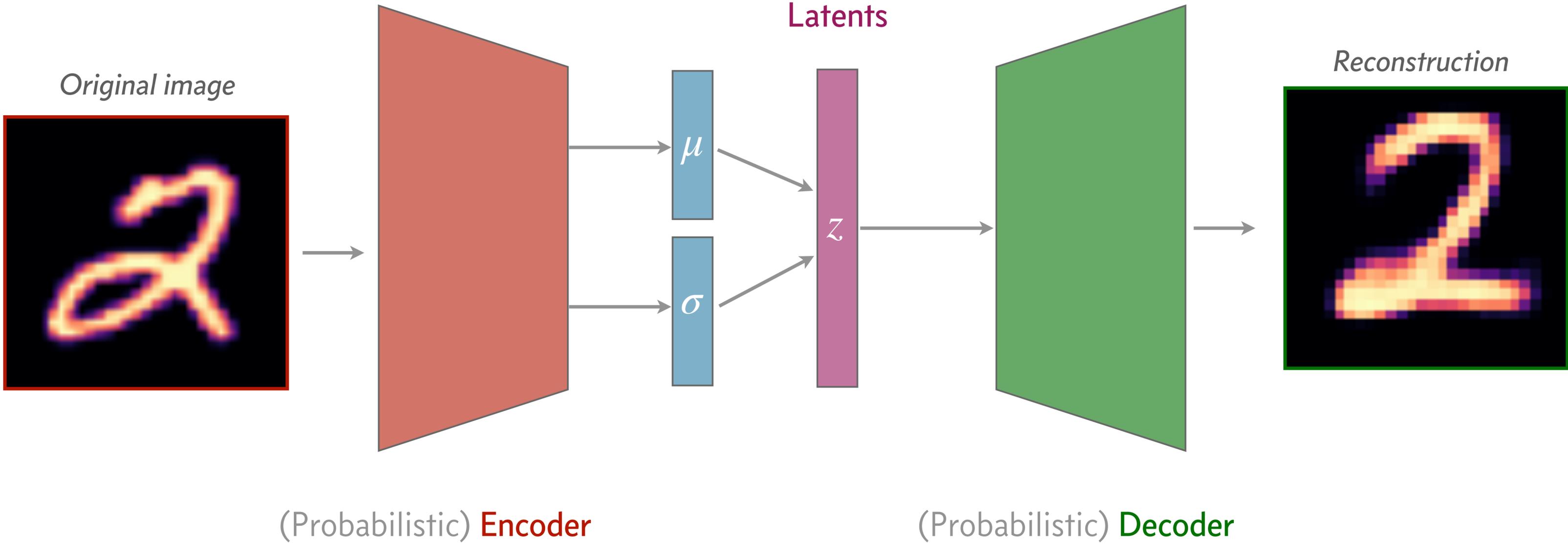
Make the problem easier by making it "harder":
introduce *joint distribution* $p_{\theta}(x, z)$



Common factorization:

$$p_{\theta}(x, z) = p(z) \cdot p_{\theta}(x | z)$$

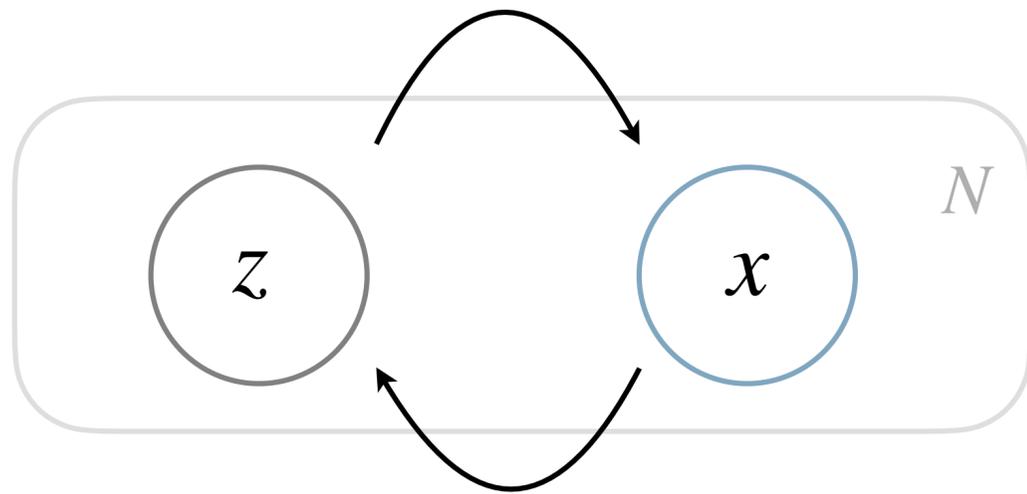
A bird-eye view



A Bayesian latent-variable model optimized with variational inference

Reverse process

$$p_{\theta}(x | z) \cdot p(z)$$



$$q_{\phi}(z | x) \cdot p(x)$$

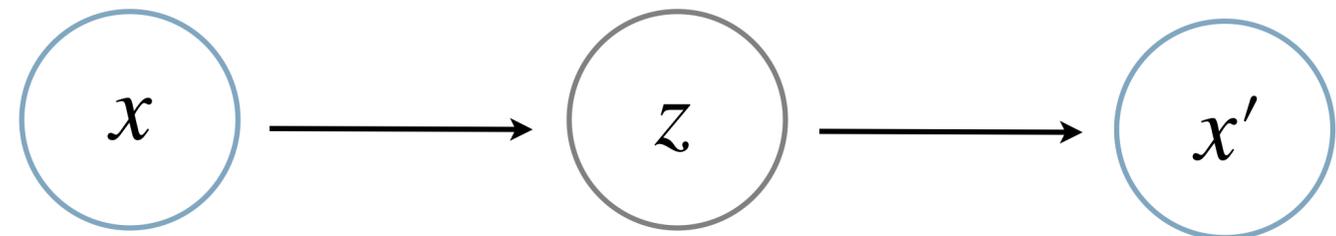
Forward process

Maximizing ELBO
 \equiv Minimizing reverse KL
 \equiv "Aligning the forward and reverse processes"

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

Forward process

Reverse process



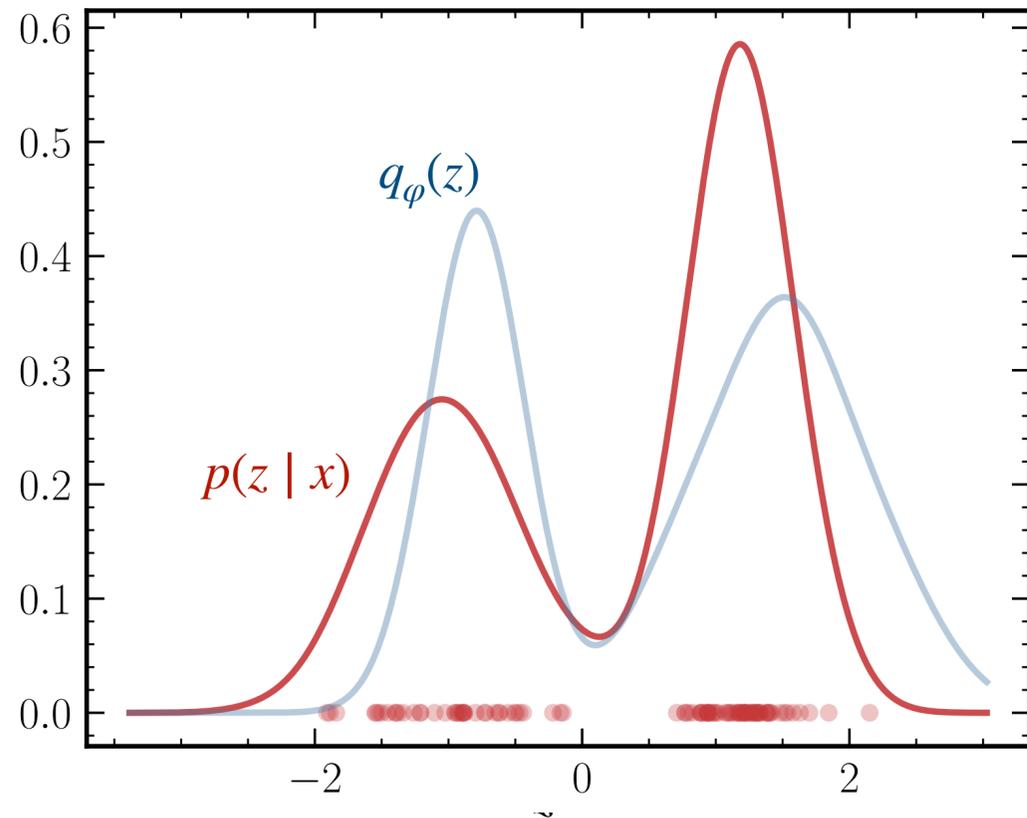
Latents

Variational inference

A general-purpose technique for posterior estimation

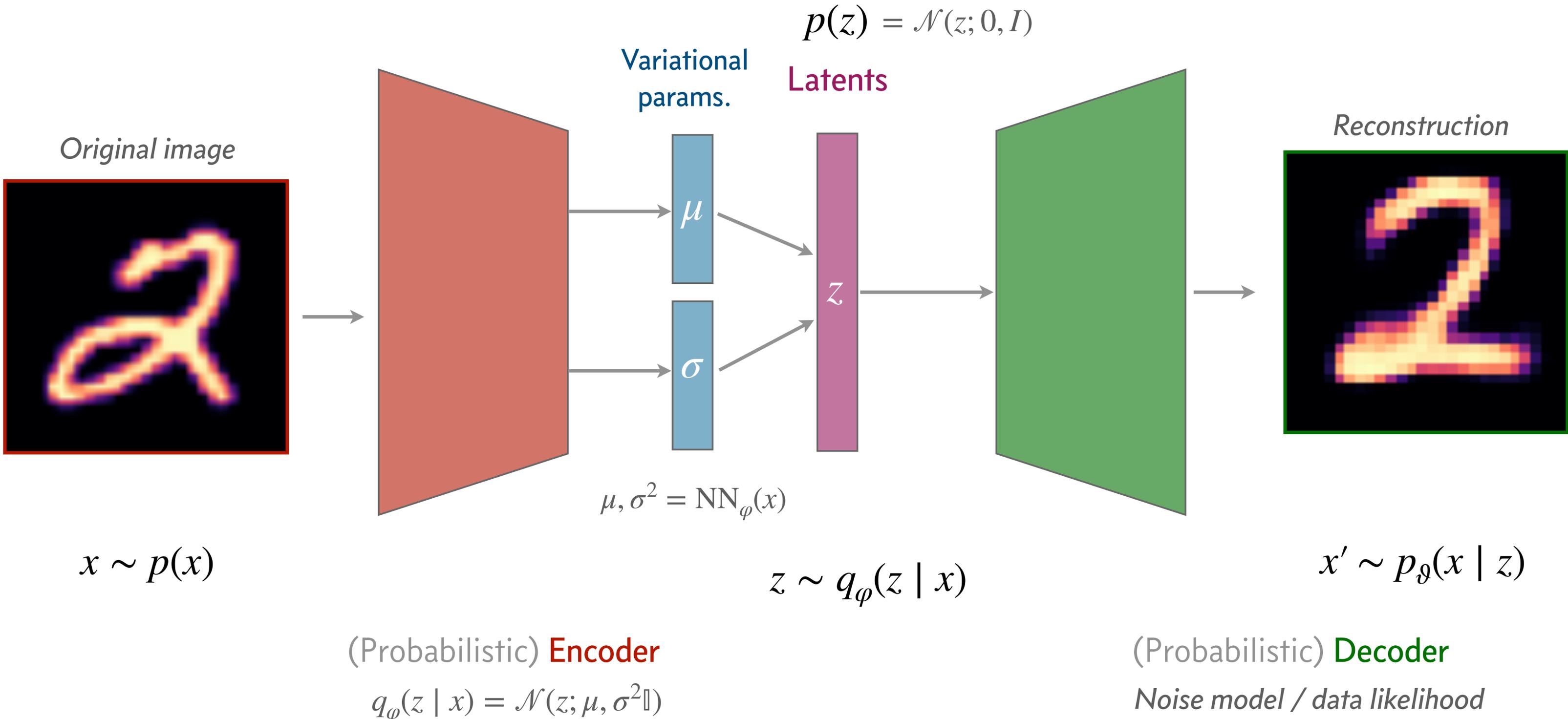
$$D_{\text{KL}}(q_\varphi(z) \| p(z | x)) \geq 0 \quad \text{Evidence} - \text{Evidence Lower BOund (ELBO)}$$

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



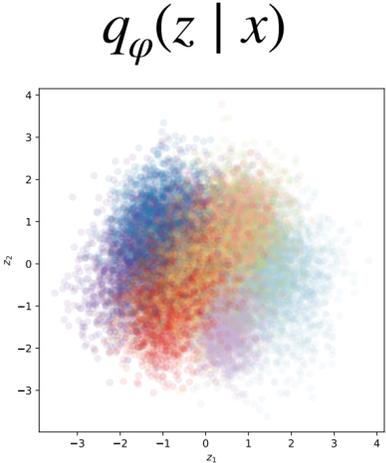
$$\begin{aligned} \text{ELBO} &= \left\langle \log p_\vartheta(x, z) - \log q_\varphi(z | x) \right\rangle_{q_\varphi} \\ &= \left\langle \log p_\vartheta(x | z) + \log p(z) - \log q_\varphi(z | x) \right\rangle_{q_\varphi} \\ &= \underbrace{\left\langle \log p_\vartheta(x | z) \right\rangle_{q_\varphi}}_{\text{"Reconstruction"}} - \underbrace{D_{\text{KL}}(q_\varphi(z | x) \| p(z))}_{\text{"Regularization"}} \end{aligned}$$

VAEs in practice



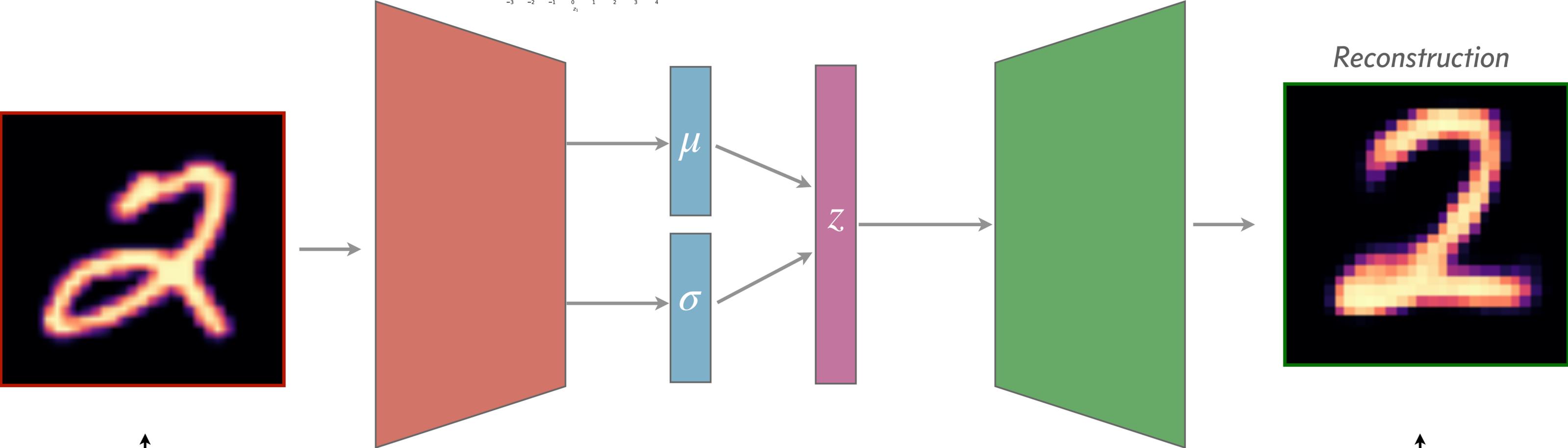
VAEs in practice

$$\text{ELBO} = \underbrace{\langle \log p_{\theta}(x | z) \rangle_{q_{\phi}}}_{\text{Reconstruction}} - \underbrace{D_{\text{KL}}(q_{\phi}(z | x) \| p(z))}_{\text{Regularization}}$$



$$D_{\text{KL}}(q_{\phi}(z | x) \| p(z))$$

Regularization



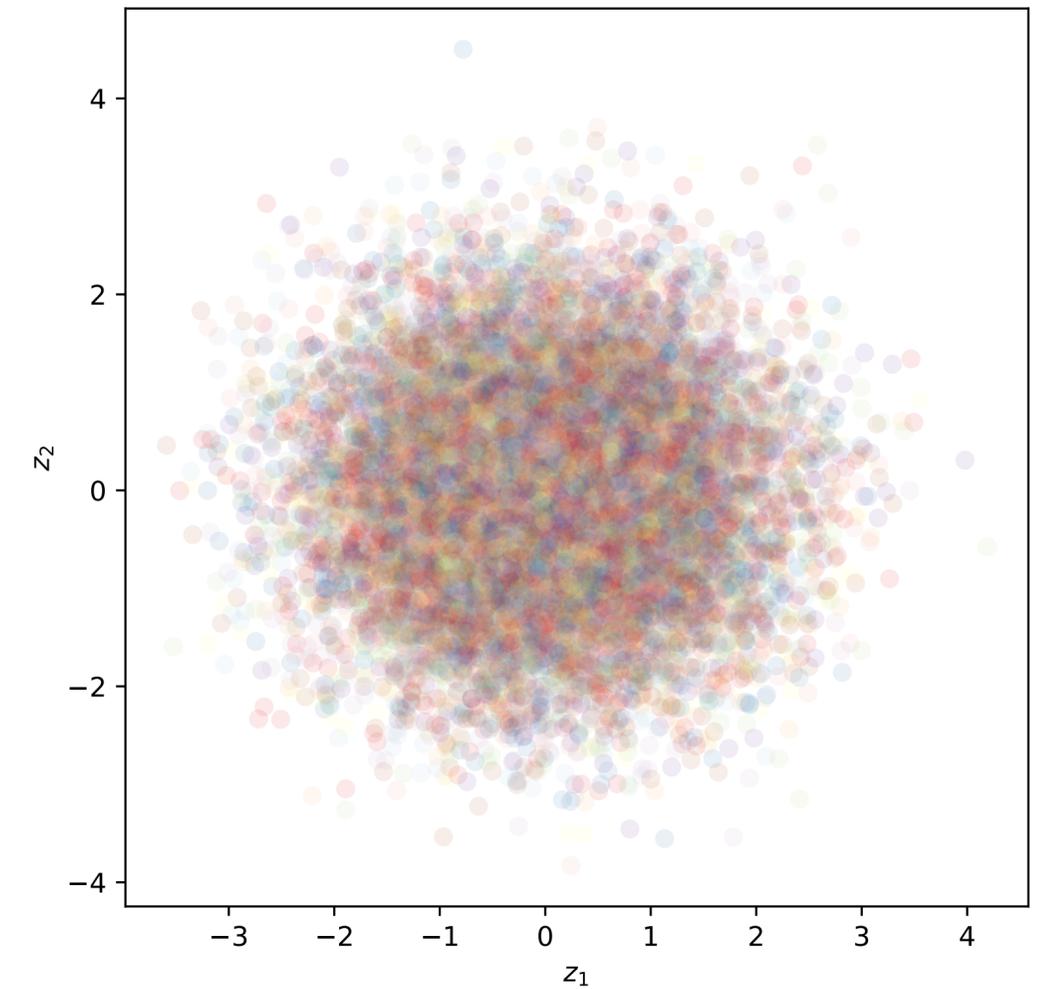
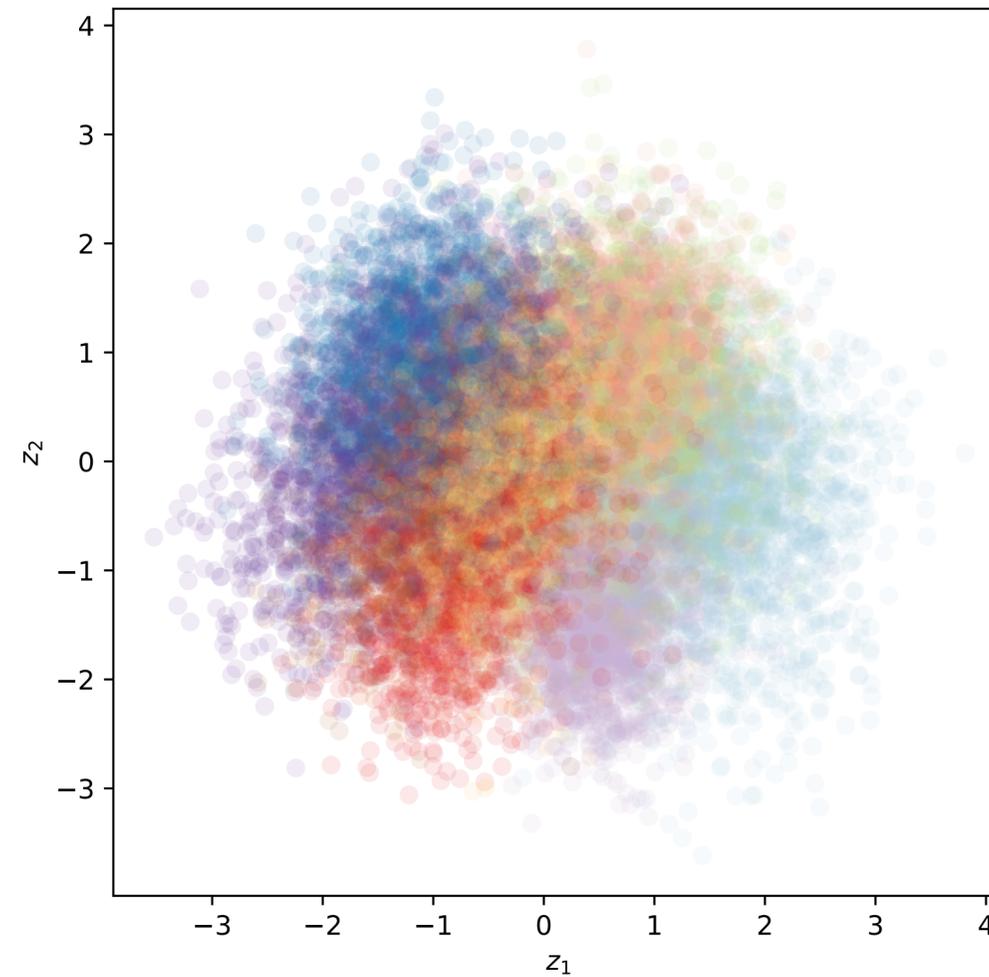
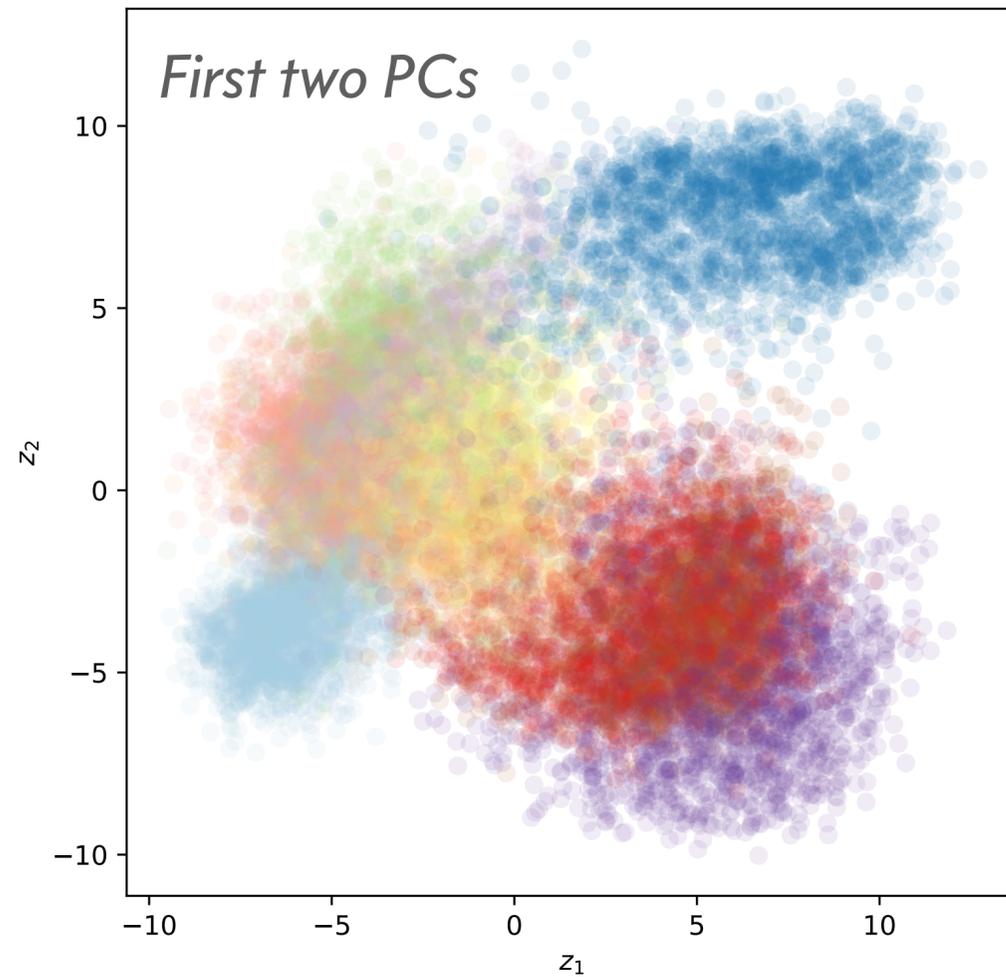
$$\langle \log p_{\theta}(x | z) \rangle_{q_{\phi}} \quad \text{Reconstruction (e.g., MSE, ...)} \quad \|x - x'\|_2^2$$

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



Neural compression: *Rate-distortion theory*

Autoencoding is a form of (neural) compression!

$$-\text{ELBO} = \underbrace{-\langle \log p_{\theta}(x | z) \rangle_{q_{\phi}}}_{\text{Distortion}} + \underbrace{D_{\text{KL}}(q_{\phi}(z | x) || p(z))}_{\text{Rate}}$$

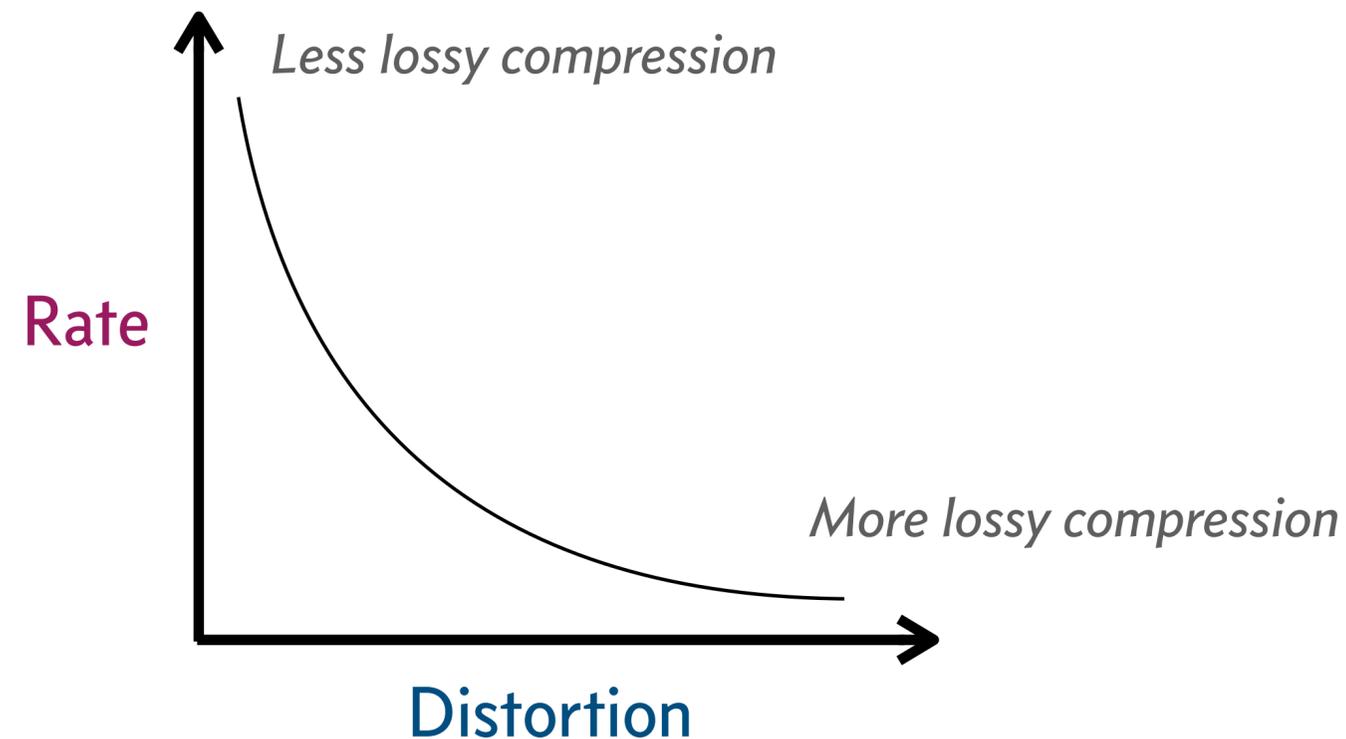
Distortion

Rate

“Reconstruction loss”

“Amount of compression”

*Rate-distortion curve
quantifies this tradeoff*



Controlling compression and disentanglement: β -VAEs

$$-\text{ELBO} = \underbrace{-\langle \log p_{\vartheta}(x | z) \rangle_{q_{\varphi}}}_{\text{Distortion}} + \underbrace{\beta \cdot D_{\text{KL}}(q_{\varphi}(z | x) \| p(z))}_{\text{Rate}}$$

Distortion

Rate

- Larger σ : More of the data variation is attributed to the likelihood \rightarrow larger " β ", more compression
- Smaller σ : Latents z try to capture more of the variation in the data (e.g. small perceptual features)

