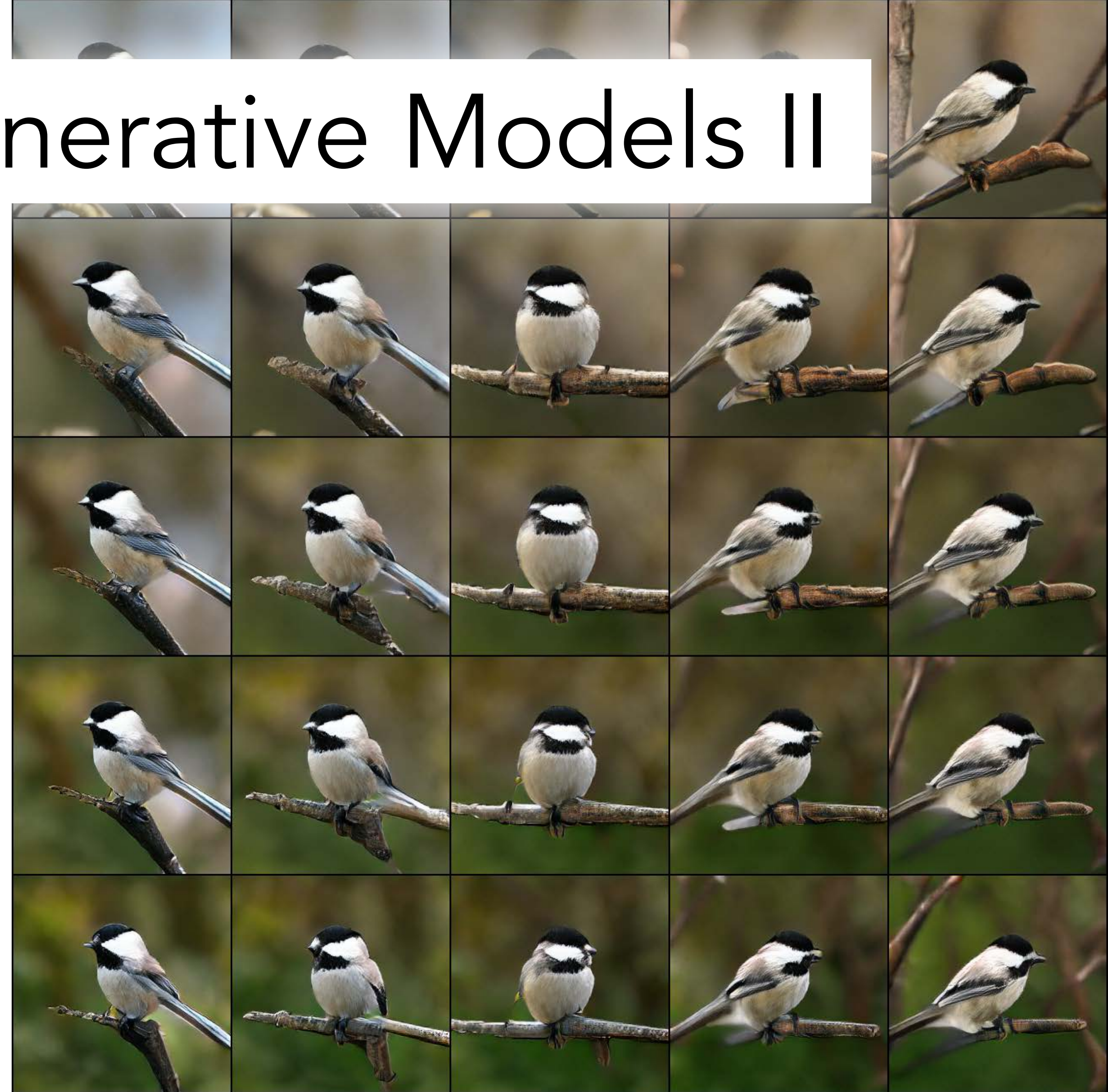


# Lecture 15: Deep Generative Models II

Speaker: Phillip Isola

© Torralba, Isola, and Freeman. All rights reserved.  
This content is excluded from our Creative Commons  
license. For more information, see  
<https://ocw.mit.edu/help/faq-fair-use/>





# Deep generative models II

- Generative modeling as inverse representation learning
- Variational autoencoders
- Disentanglement

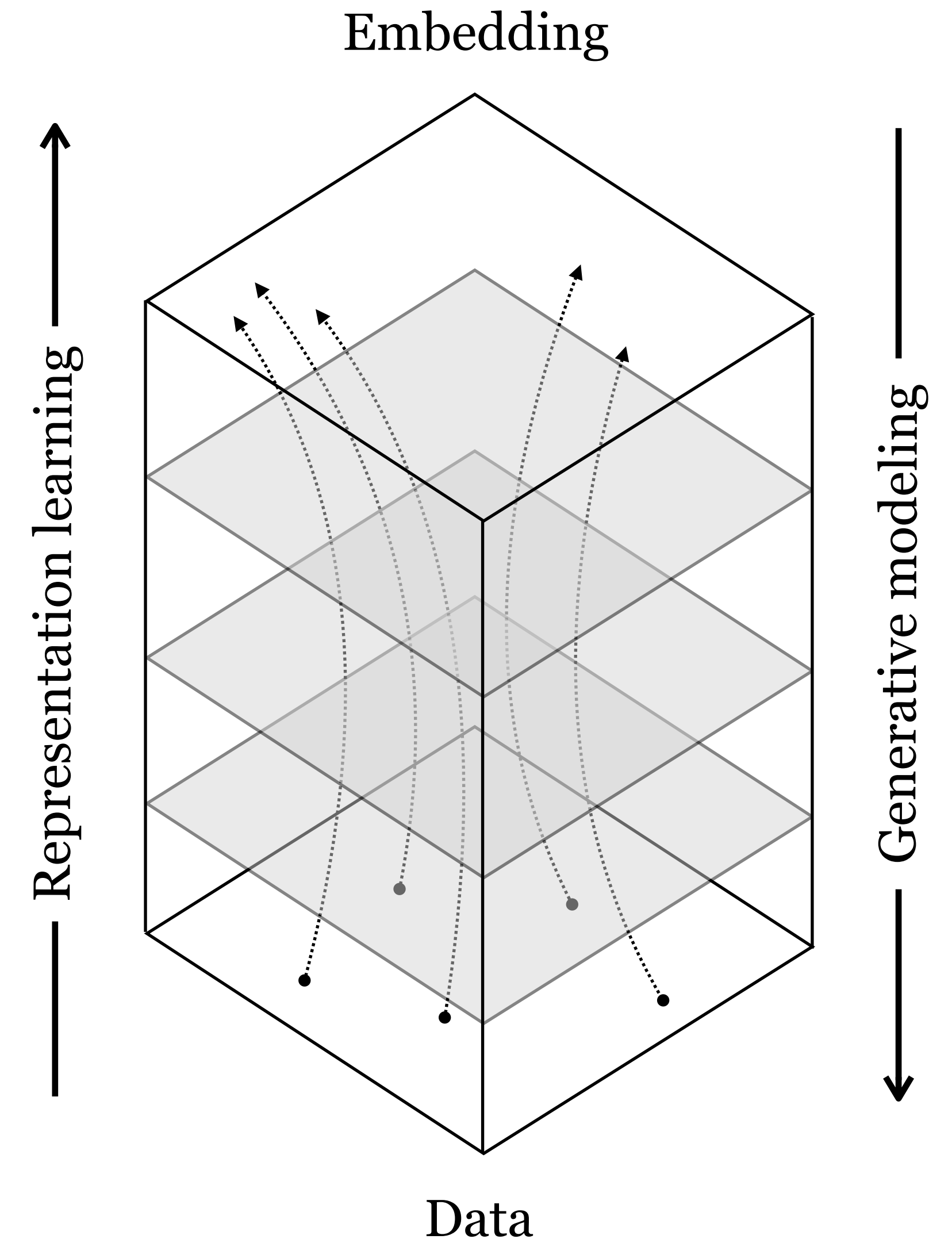
# Generative modeling vs Representation learning

## **Representation learning:**

mapping data to abstract representations  
(analysis)

## **Generative modeling:**

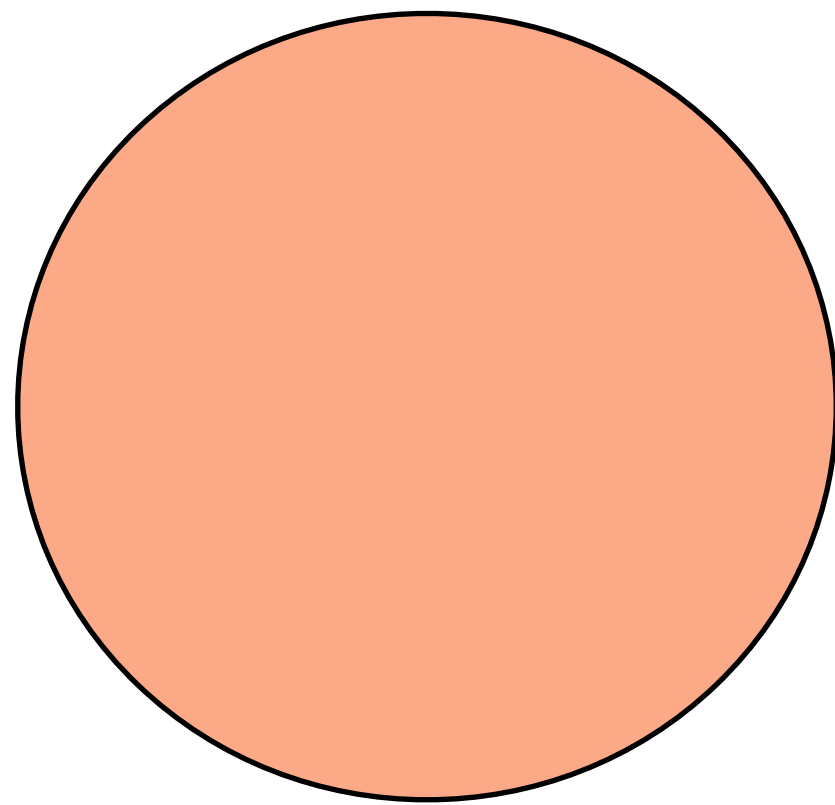
mapping abstract representations to data  
(*synthesis*)



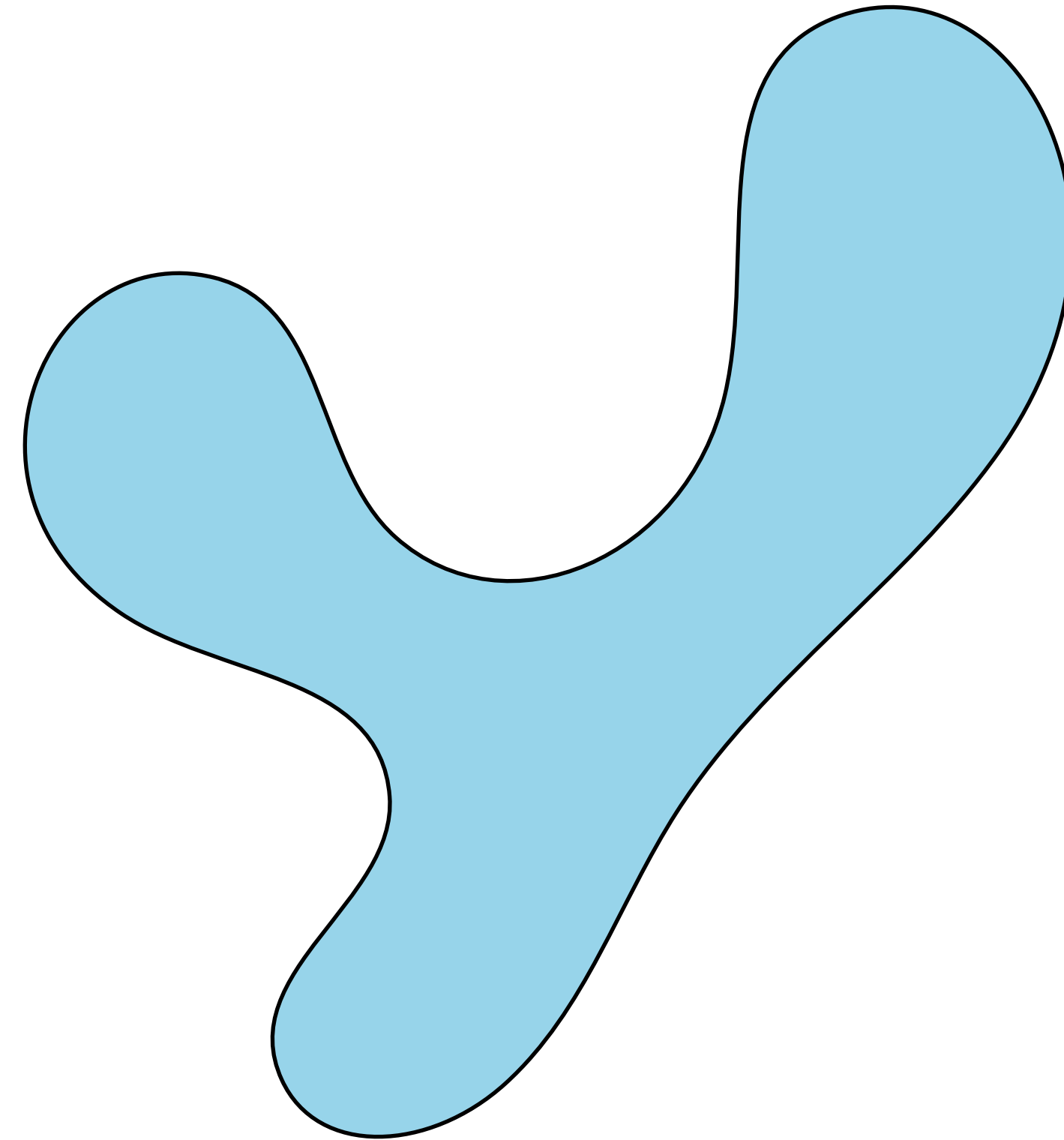
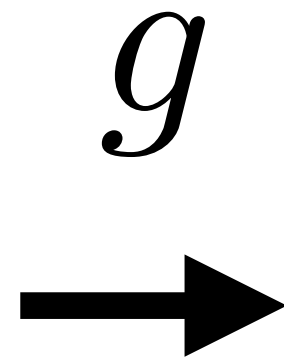
# Deep generative models are distribution transformers

Prior distribution

Target distribution

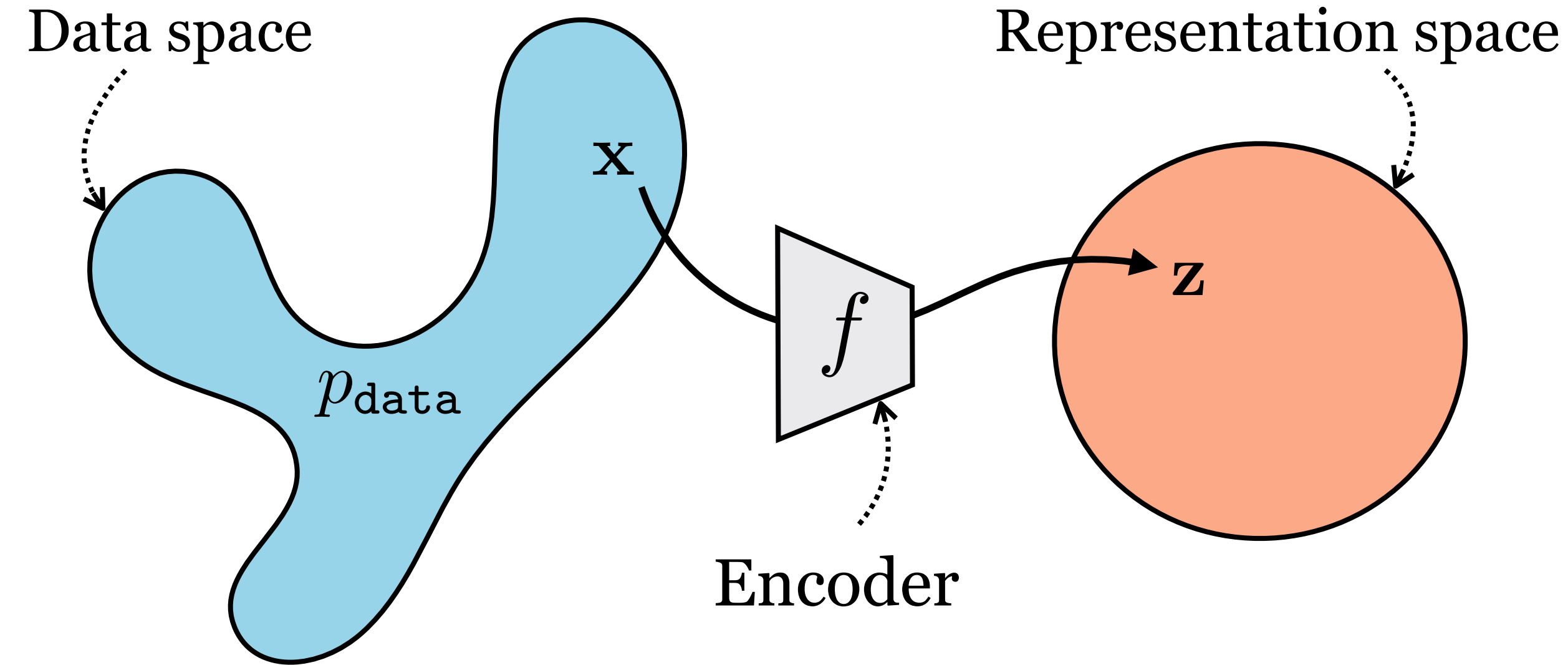


$p(z)$

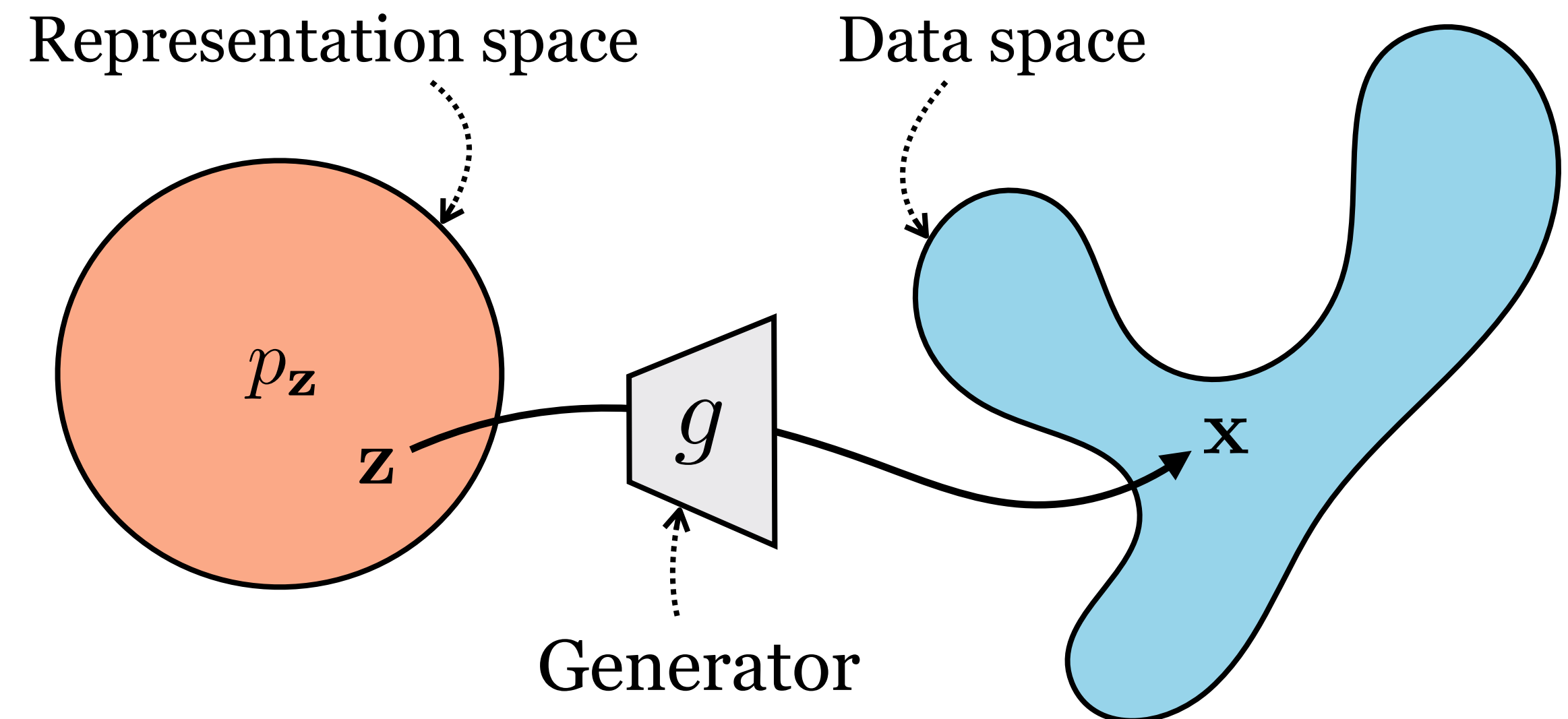


$p(x)$

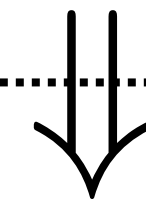
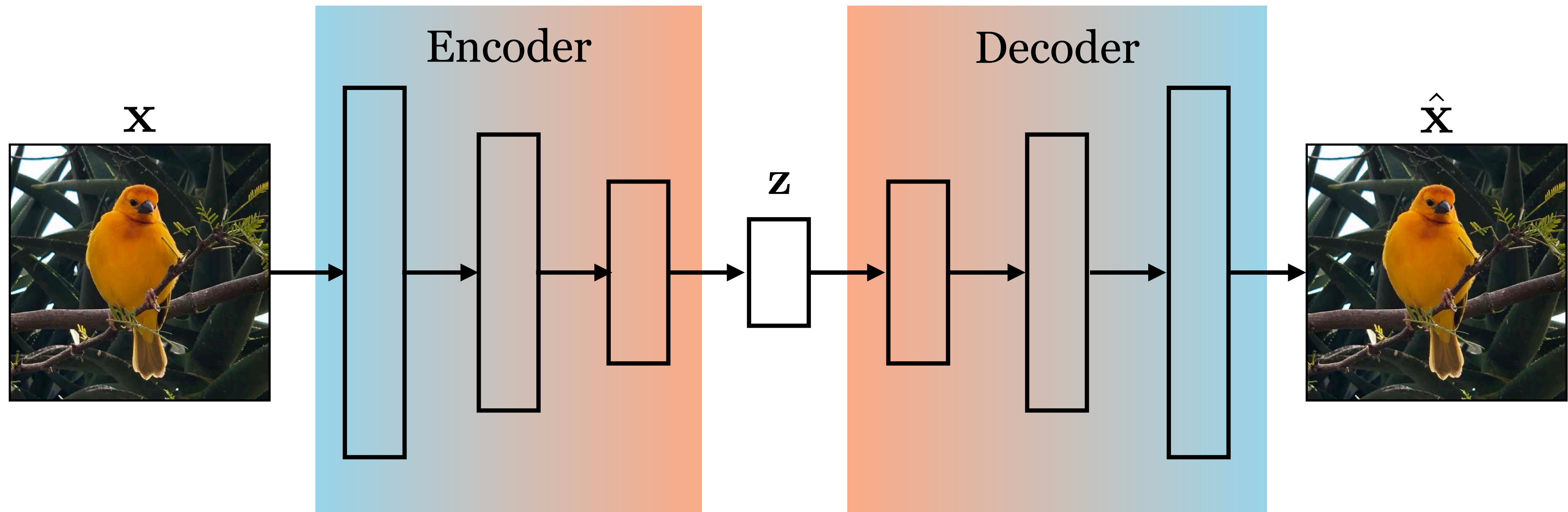
## Representation learning



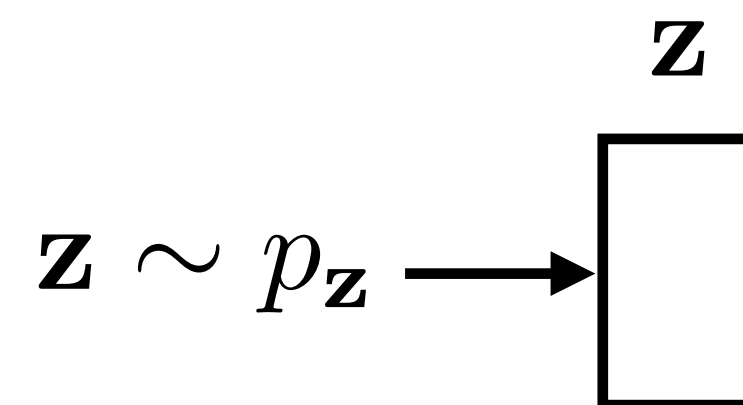
## Generative modeling



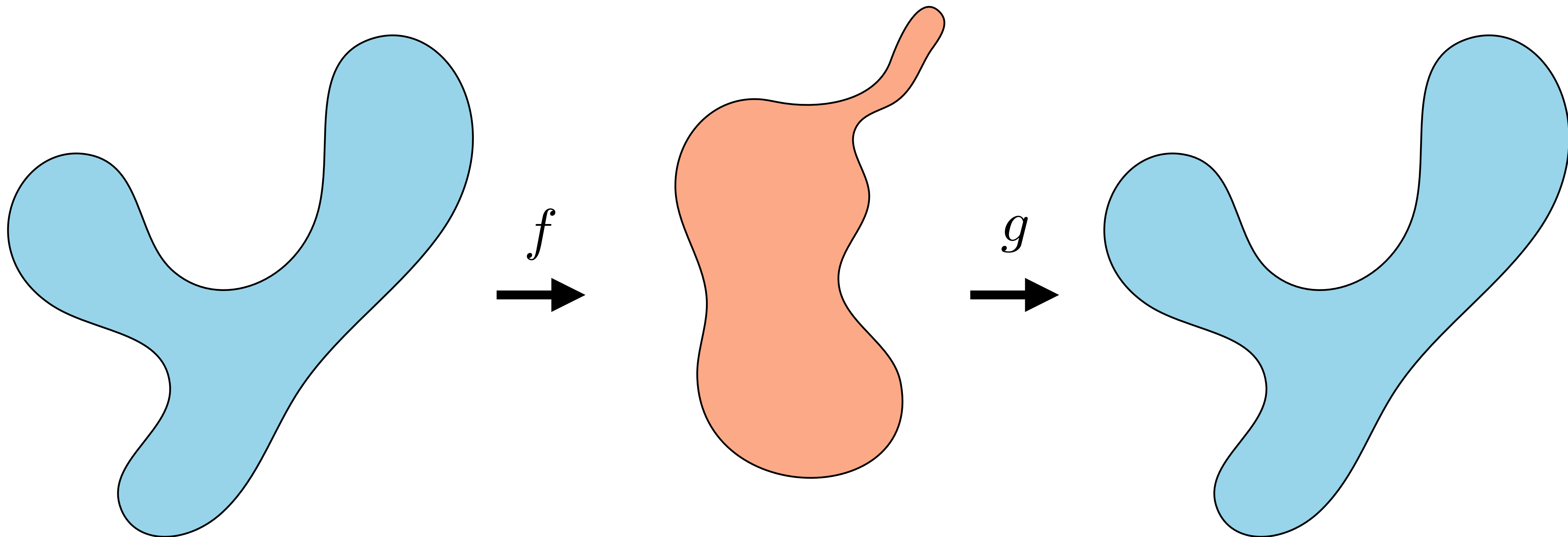
# Autoencoder



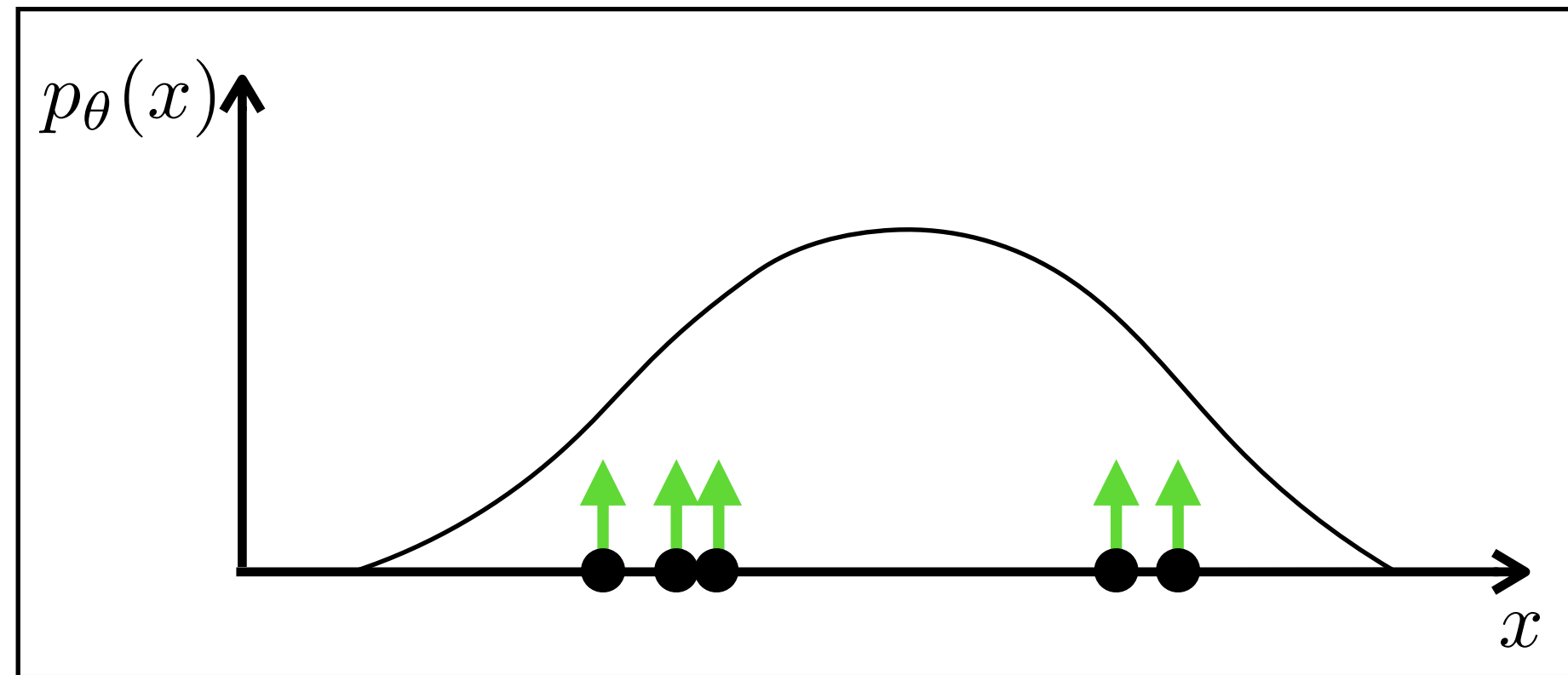
# Generative model



⋮



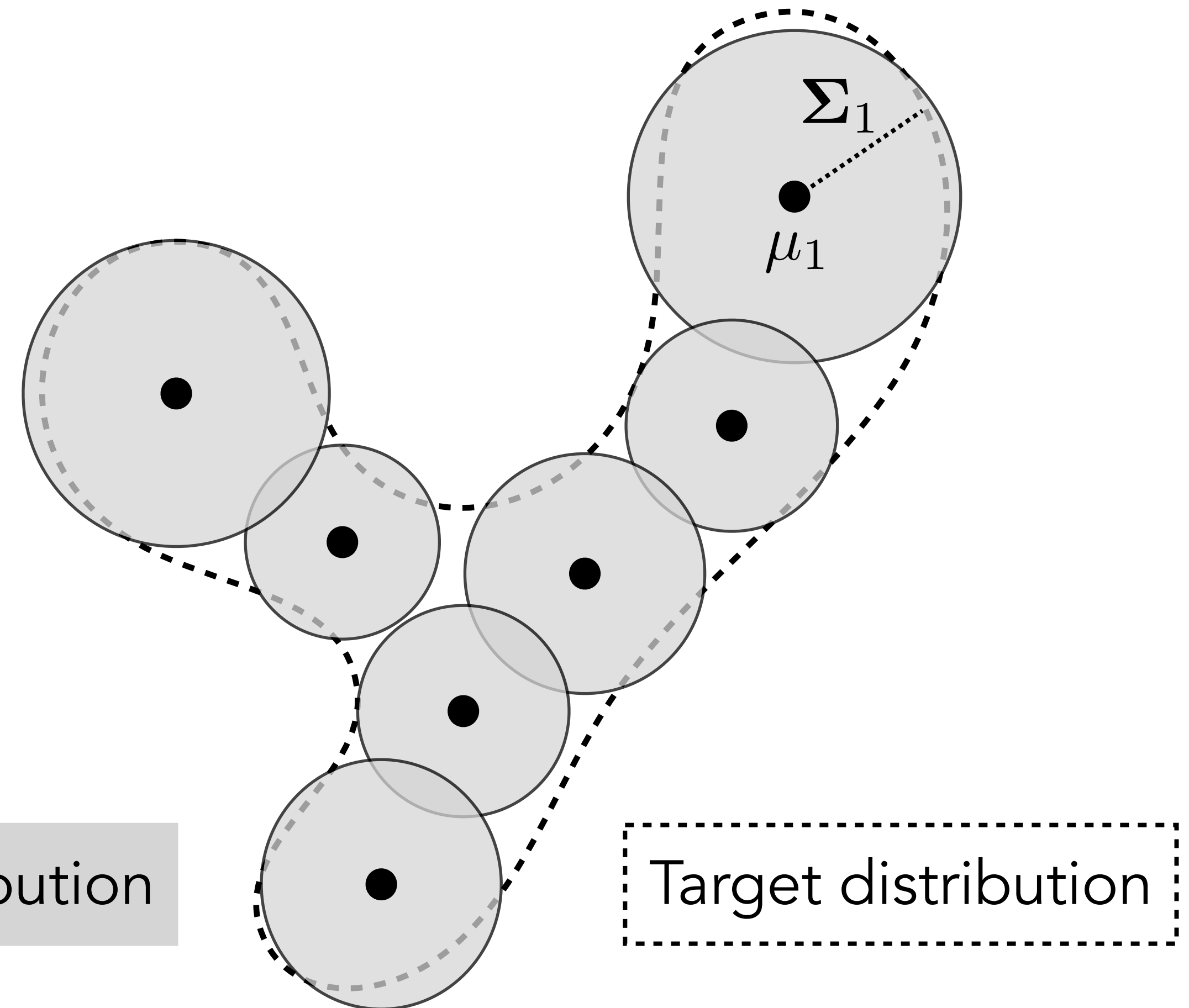
# Mixture of Gaussians



Parameters:  $\{w_i, \mu_i, \Sigma_i\}_{i=1}^k$

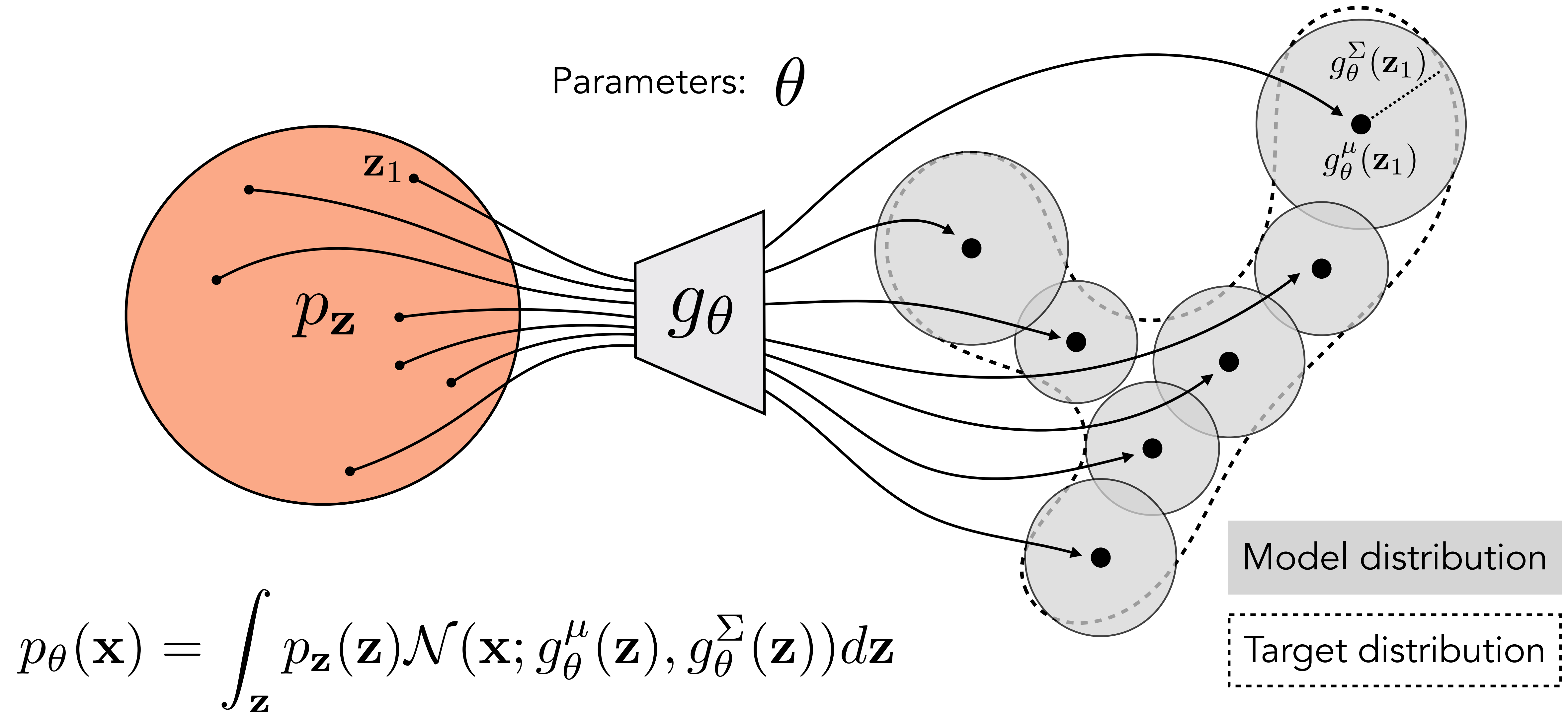
$$p_{\theta}(\mathbf{x}) = \sum_{i=1}^k w_i \mathcal{N}(\mathbf{x}; \mu_i, \Sigma_i)$$

Model distribution

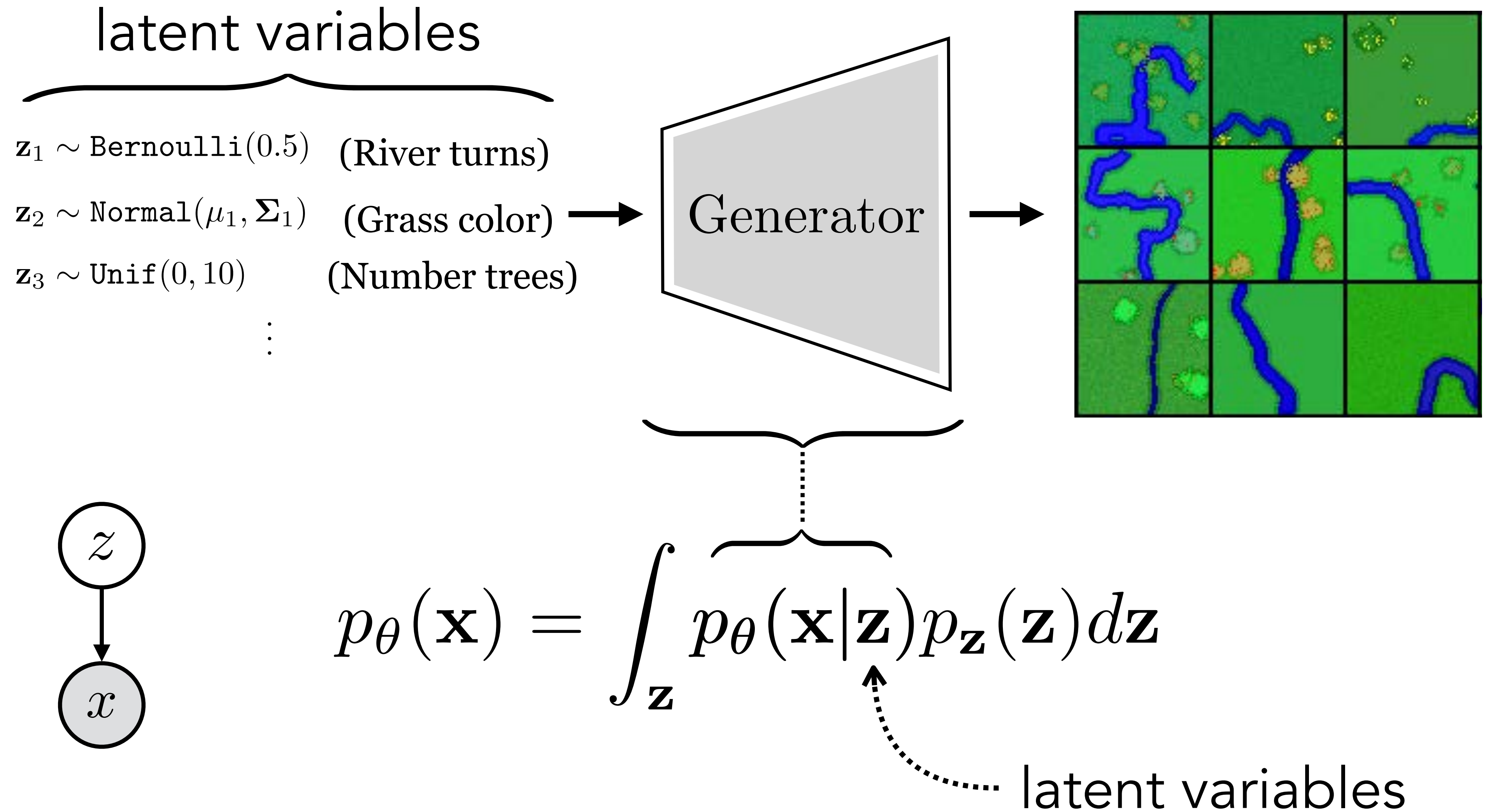




# Infinite Mixture of Gaussians



# Latent variable models



# Infinite GMM / VAE

Learner

Objective

$$\max_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\theta}(\mathbf{x})]$$

Hypothesis space

$$p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} \underbrace{\mathcal{N}(\mathbf{x}; g_{\theta}^{\mu}(\mathbf{z}), g_{\theta}^{\Sigma}(\mathbf{z}))}_{p_{\theta}(\mathbf{x} | \mathbf{z})} \underbrace{\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})}_{p_{\mathbf{z}}(\mathbf{z})} d\mathbf{z}$$

$$x = g_{\theta}^{\mu}(\mathbf{z}) + g_{\theta}^{\sigma}(\mathbf{z})\epsilon \quad \epsilon \sim \mathcal{N}(0, 1)$$

“Reparameterization trick”

Density

$$p_{\theta} : \mathcal{X} \rightarrow [0, \infty)$$

Sampler

$$g_{\theta} : \mathcal{Z} \rightarrow \mathcal{X}$$

Data  $\rightarrow$

$\rightarrow$

# VAEs — learning the model parameters

$$\theta^* = \arg \max_{\theta} L(\{\mathbf{x}^{(i)}\}_{i=1}^N, \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^N \log \underbrace{\int_{\mathbf{z}} \overbrace{\mathcal{N}(\mathbf{x}^{(i)}; g_{\theta}^{\mu}(\mathbf{z}), g_{\theta}^{\Sigma}(\mathbf{z}))}^{p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z})} \overbrace{\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})}^{p_{\mathbf{z}}(\mathbf{z})} d\mathbf{z}}_{p_{\theta}(\mathbf{x}^{(i)})}$$



Trick #1: Estimate the integral via sampling (Monte Carlo)

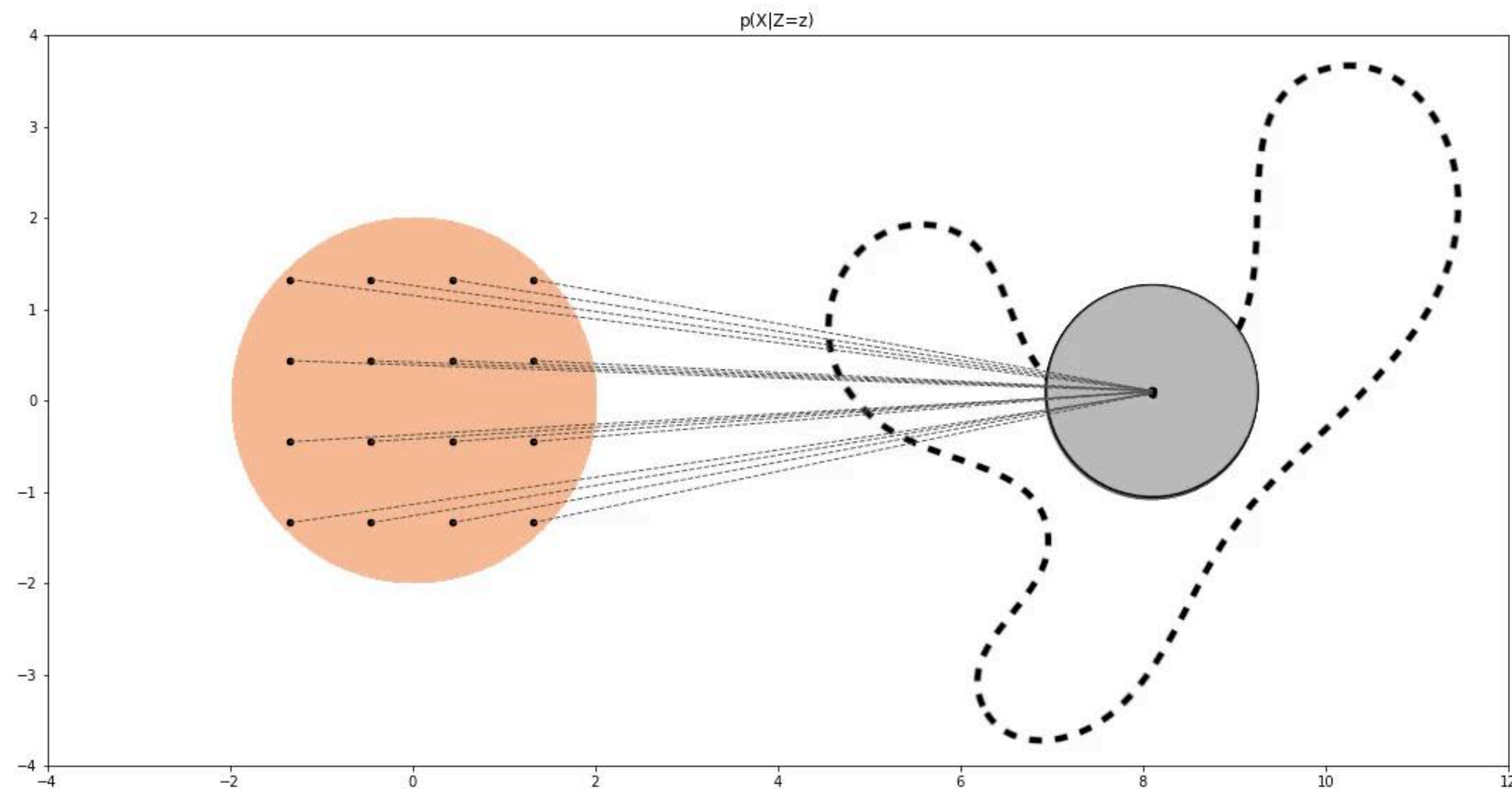
$$p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}|\mathbf{z})p_{\mathbf{z}}(\mathbf{z})d\mathbf{z}$$

$$= \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [p_{\theta}(\mathbf{x}|\mathbf{z})]$$

$$\approx \frac{1}{M} \sum_{i=1}^M p_{\theta}(\mathbf{x}|\mathbf{z}^{(i)}), \quad \mathbf{z}^{(j)} \sim p_{\mathbf{z}}$$

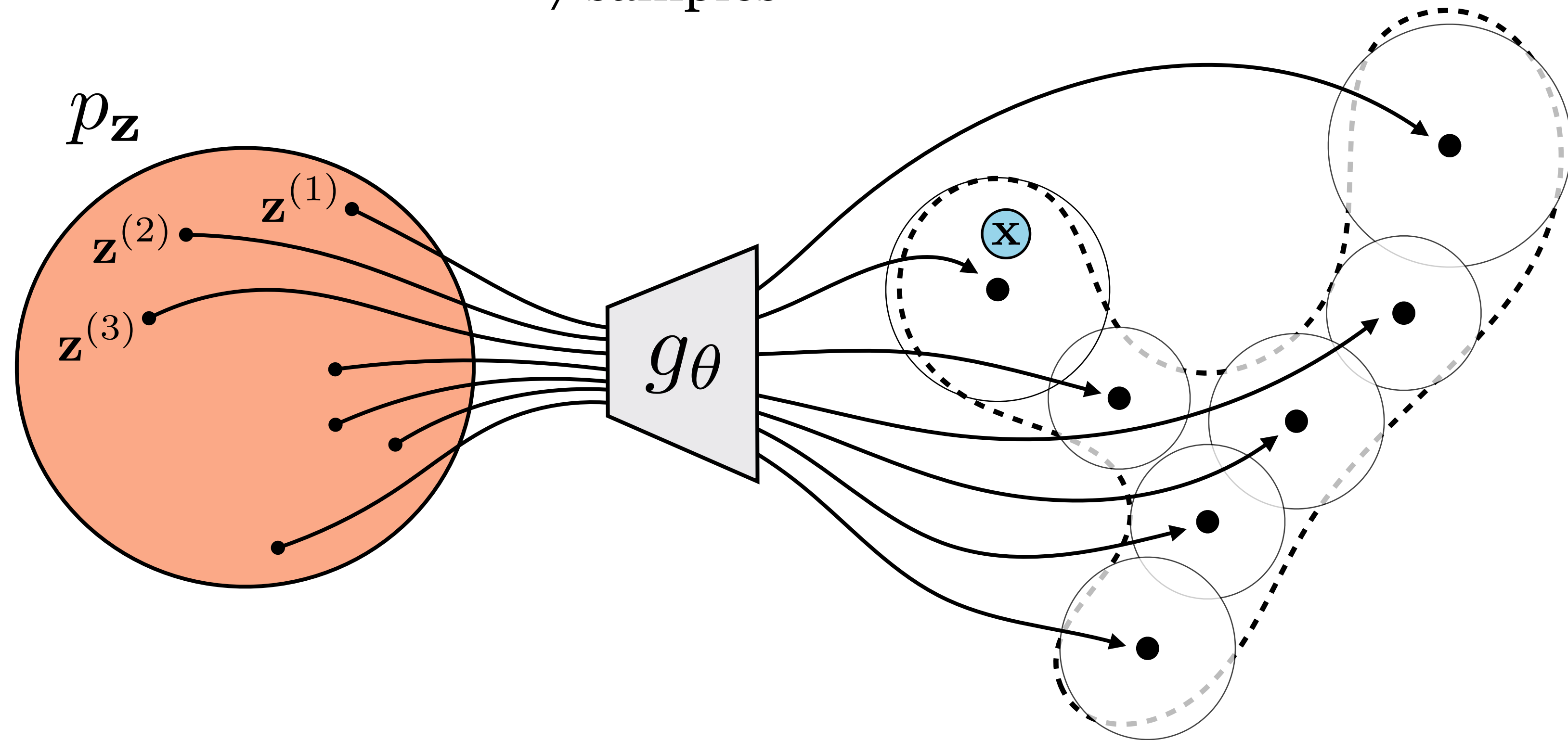
# Fitting an infinite mixture of Gaussians via Monte Carlo

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^N \log \frac{1}{M} \sum_{j=1}^M \overbrace{\mathcal{N}(\mathbf{x}^{(i)}; g_{\theta}^{\mu}(\mathbf{z}^{(j)}), g_{\theta}^{\Sigma}(\mathbf{z}^{(j)}))}^{p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(j)})}$$



## Trick #2: Importance sampling

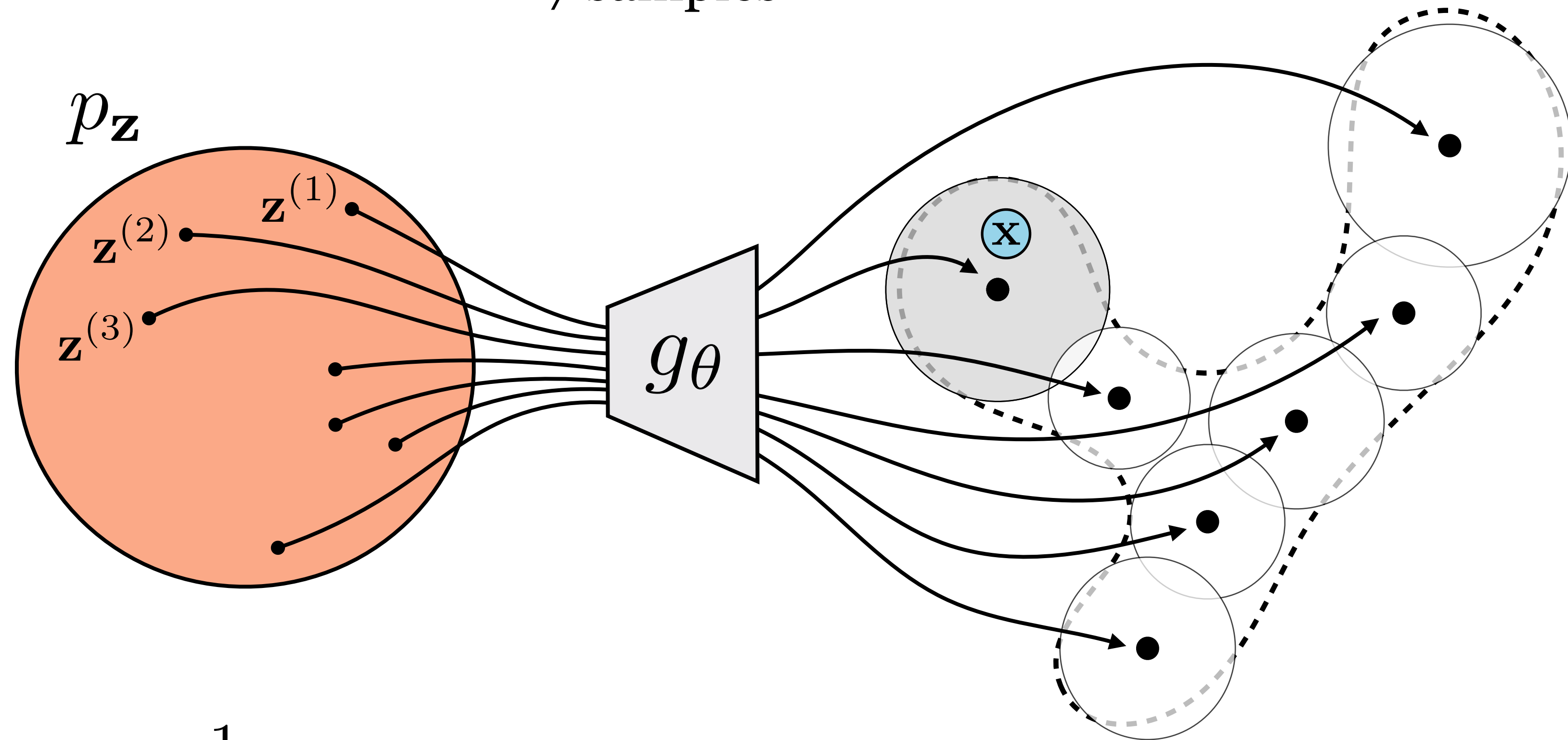
Random sampling     $\mathbf{z}^{(j)} \sim p_{\mathbf{z}}$   
7 samples



$$p_{\theta}(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [p_{\theta}(\mathbf{x}|\mathbf{z})]$$

## Trick #2: Importance sampling

Random sampling     $\mathbf{z}^{(j)} \sim p_{\mathbf{z}}$   
7 samples



$$p_{\theta}(\mathbf{x}) \approx \frac{1}{M} (p_{\theta}(\mathbf{x}|\mathbf{z}^{(1)}) + p_{\theta}(\mathbf{x}|\mathbf{z}^{(2)}) + p_{\theta}(\mathbf{x}|\mathbf{z}^{(3)}) + \dots)$$

$$p_{\theta}(\mathbf{x}) \approx \frac{1}{M} (0 + p_{\theta}(\mathbf{x}|\mathbf{z}^{(2)}) + 0 + \dots)$$



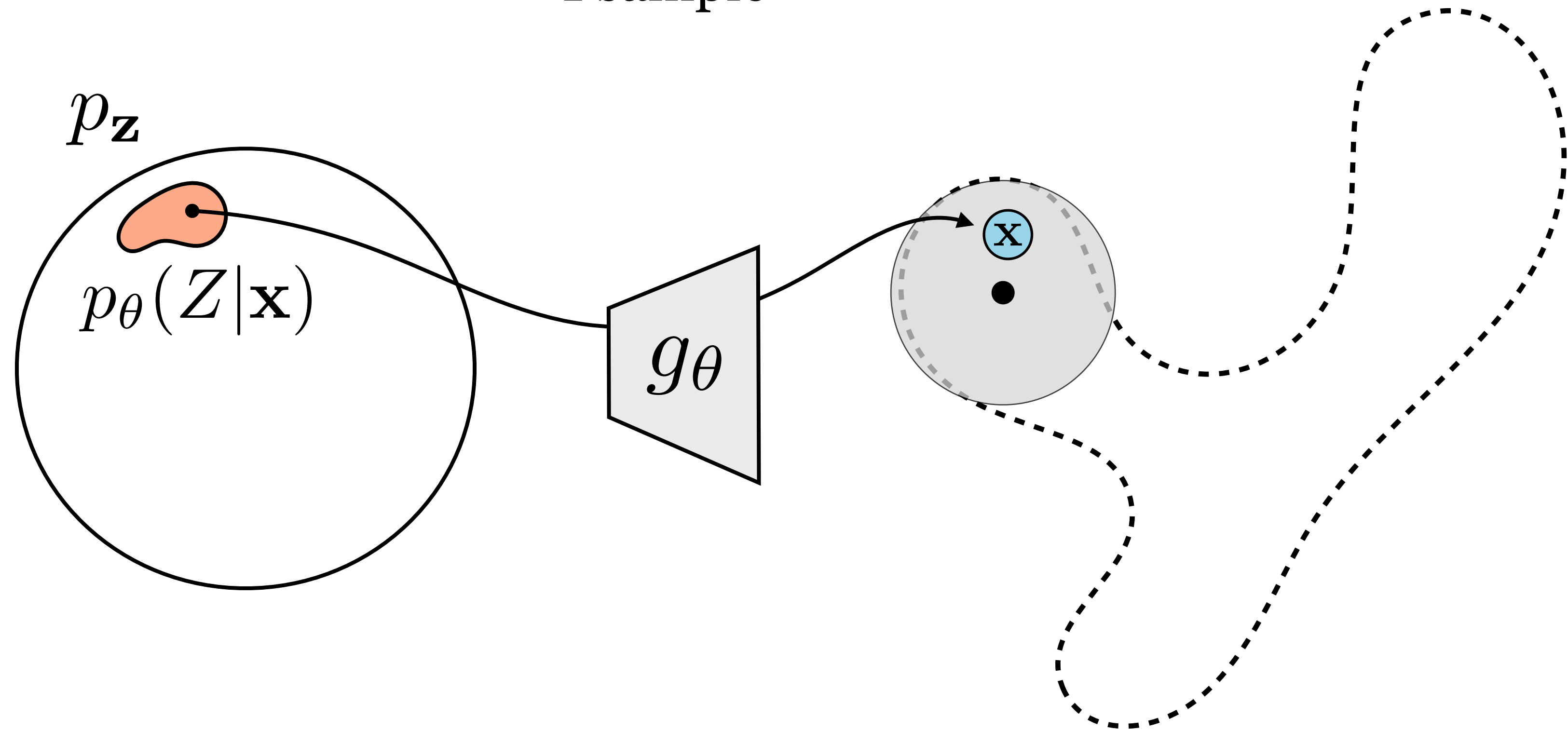
## Trick #2: Importance sampling

$$\begin{aligned} p_{\theta}(\mathbf{x}) &= \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} \left[ p_{\theta}(\mathbf{x} | \mathbf{z}) \right] = \int_{\mathbf{z}} p_{\mathbf{z}}(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z}) d\mathbf{z} = \int_{\mathbf{z}} q_{\mathbf{z}}(\mathbf{z}) \frac{p_{\mathbf{z}}(\mathbf{z})}{q_{\mathbf{z}}(\mathbf{z})} p_{\theta}(\mathbf{x} | \mathbf{z}) d\mathbf{z} \\ &= \mathbb{E}_{\mathbf{z} \sim q_{\mathbf{z}}} \left[ \frac{p_{\mathbf{z}}(\mathbf{z})}{q_{\mathbf{z}}(\mathbf{z})} p_{\theta}(\mathbf{x} | \mathbf{z}) \right] \end{aligned}$$

Set  $q_{\mathbf{z}} = p_{\theta}(Z | \mathbf{x})$

## Trick #2: Importance sampling

Importance sampling    $\mathbf{z}^{(j)} \sim p_\theta(Z|\mathbf{x})$   
1 sample

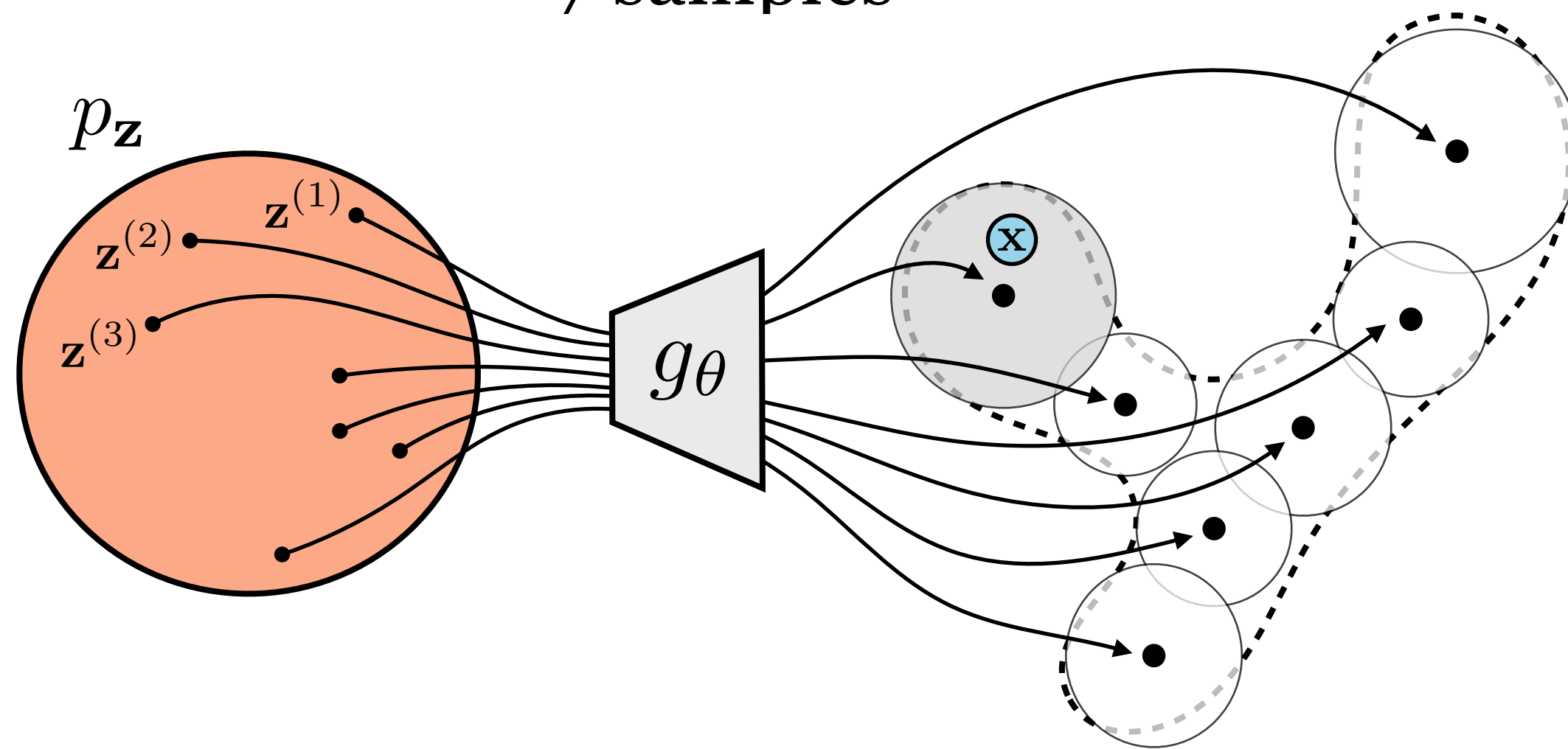


$$p_\theta(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [p_\theta(\mathbf{x}|\mathbf{z})] \quad \longrightarrow \quad p_\theta(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim p_\theta(Z|\mathbf{x})} \left[ \frac{p_{\mathbf{z}}(\mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} p_\theta(\mathbf{x}|\mathbf{z}) \right]$$

# Trick #2: Importance sampling

## Random sampling

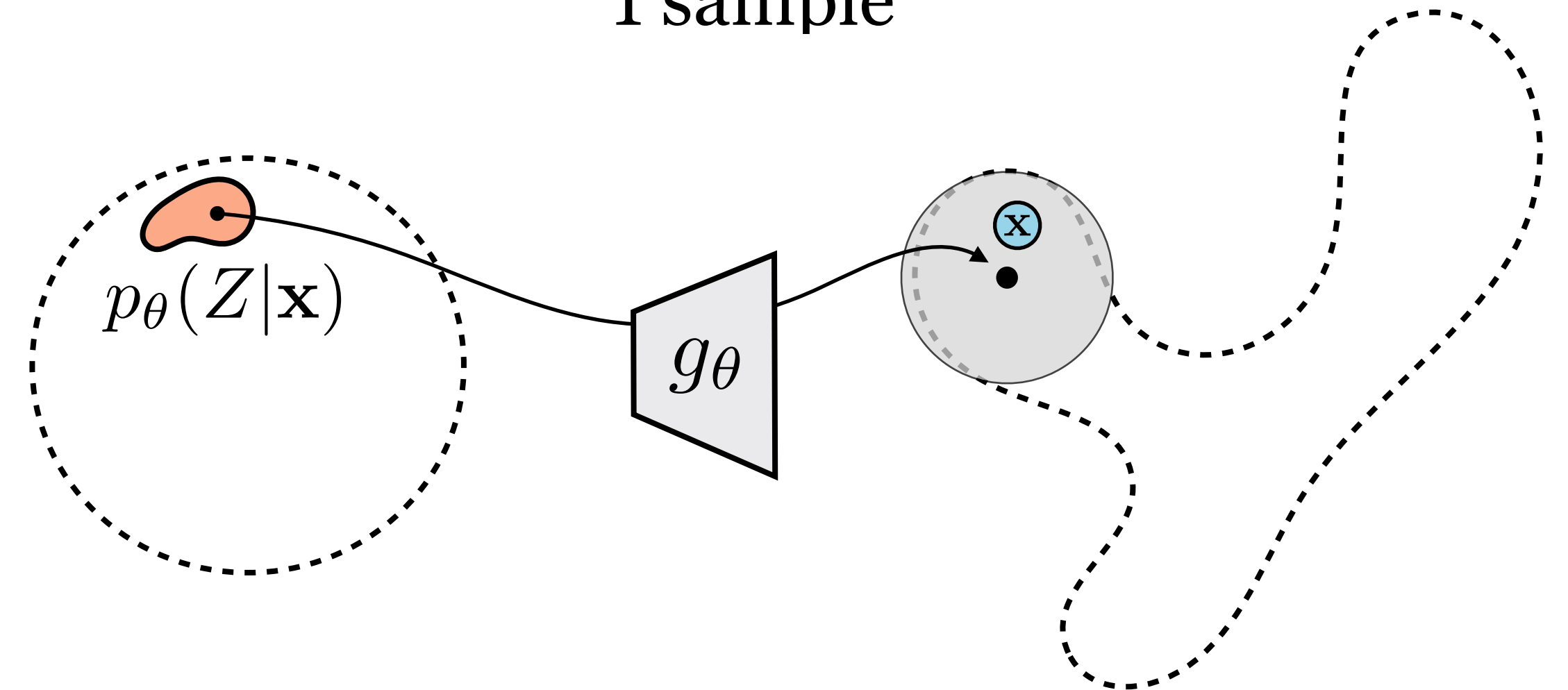
7 samples



$$p_{\theta}(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [p_{\theta}(\mathbf{x}|\mathbf{z})]$$

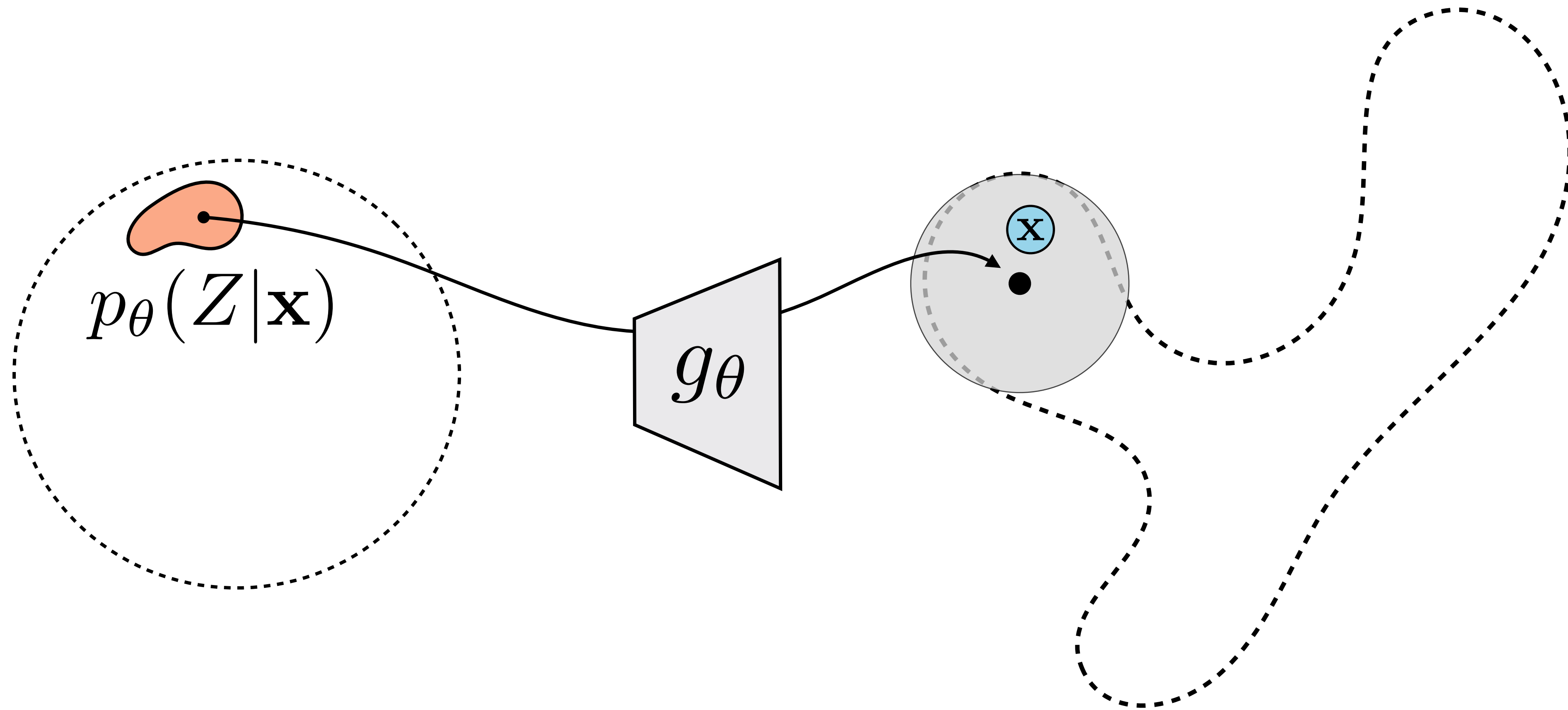
## Importance sampling

1 sample



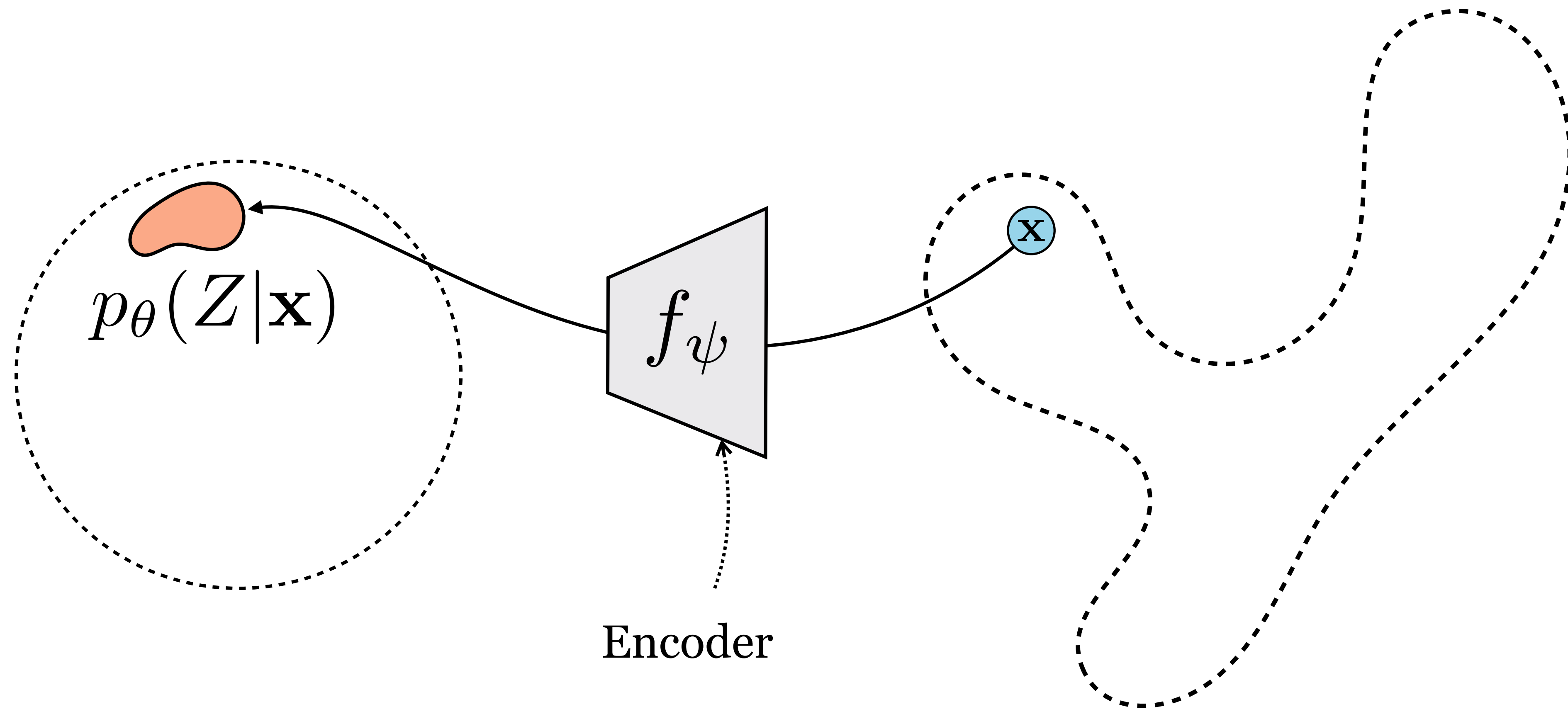
$$p_{\theta}(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim p_{\theta}(Z|\mathbf{x})} \left[ \frac{p_{\mathbf{z}}(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} p_{\theta}(\mathbf{x}|\mathbf{z}) \right]$$

Trick #3: Predict the optimal sampling distribution for each given datapoint



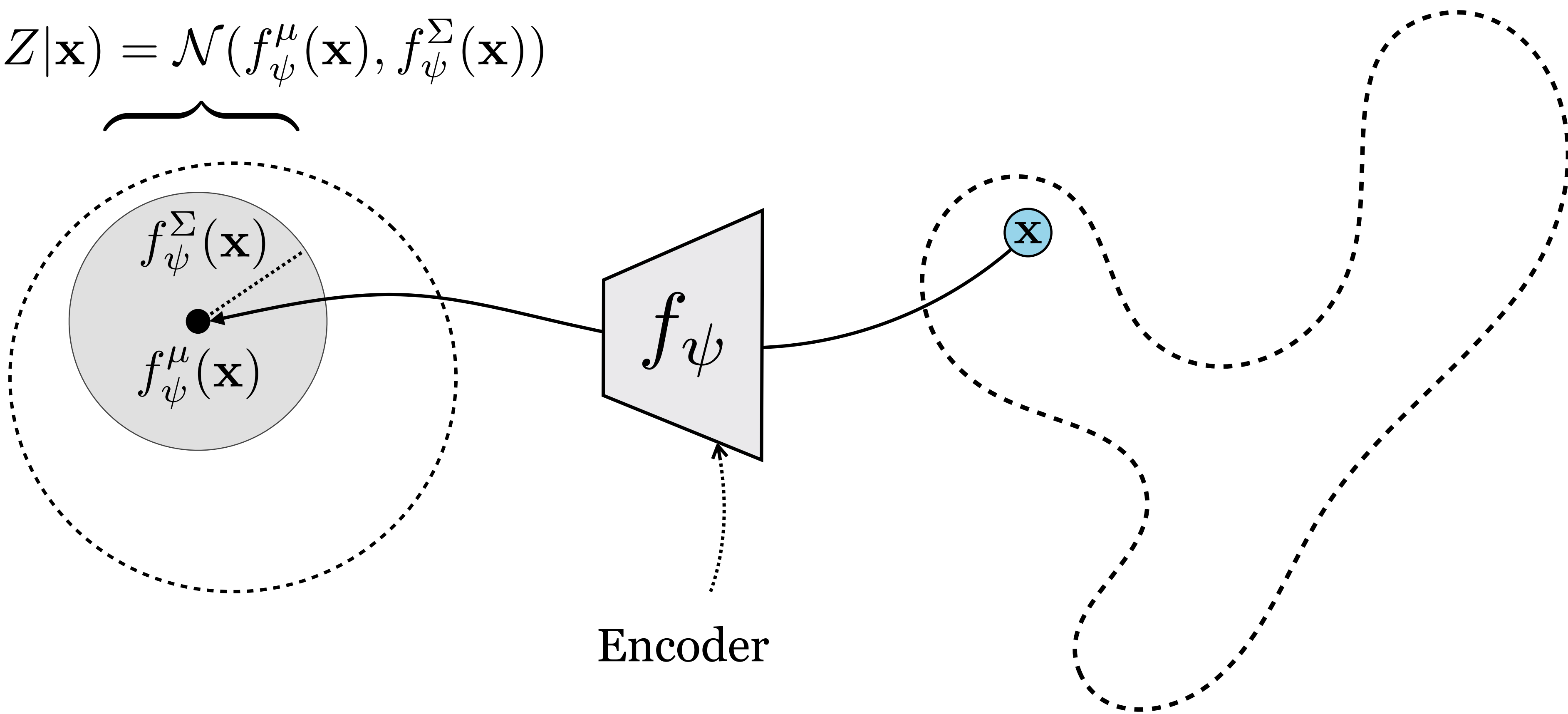


Trick #3: Predict the optimal sampling distribution for each given datapoint



# Trick #3: Predict the optimal sampling distribution for each given datapoint

$$q_{\psi}(Z|\mathbf{x}) = \mathcal{N}(\underbrace{f_{\psi}^{\mu}(\mathbf{x}), f_{\psi}^{\Sigma}(\mathbf{x})})$$



# Trick #3: Predict the optimal sampling distribution for each given datapoint

We are fitting this ..... to model this

$$J_q(\mathbf{x}, \psi) = -\text{KL}(\boxed{q_\psi(Z | \mathbf{x})} \| p_\theta(Z | \mathbf{x}))$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\psi(Z | \mathbf{x})} [-\log q_\psi(\mathbf{z} | \mathbf{x}) + \log p_\theta(\mathbf{x} | \mathbf{z}) + \log p_{\mathbf{z}}(\mathbf{z})] - \log p_\theta(\mathbf{x})$$

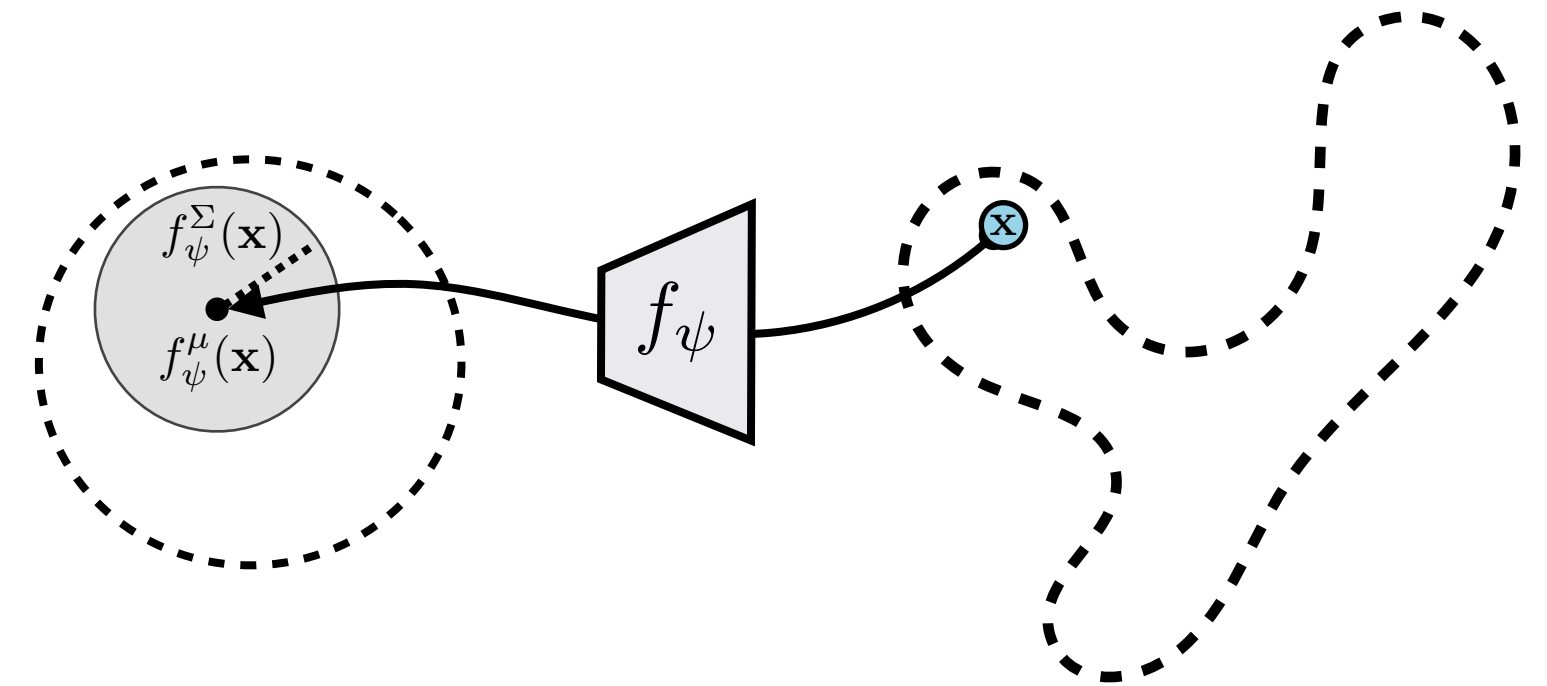
Our learning problem for q:

$$\begin{aligned} \psi^* &= \arg \max_{\psi} \frac{1}{N} \sum_{i=1}^N J_q(\mathbf{x}^{(i)}, \psi) \\ &= \arg \max_{\psi} \frac{1}{N} \sum_{i=1}^N \underbrace{\mathbb{E}_{\mathbf{z} \sim q_\psi(Z | \mathbf{x}^{(i)})} [-\log q_\psi(\mathbf{z} | \mathbf{x}^{(i)}) + \log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) + \log p_{\mathbf{z}}(\mathbf{z})]}_J \end{aligned}$$

# Putting everything together

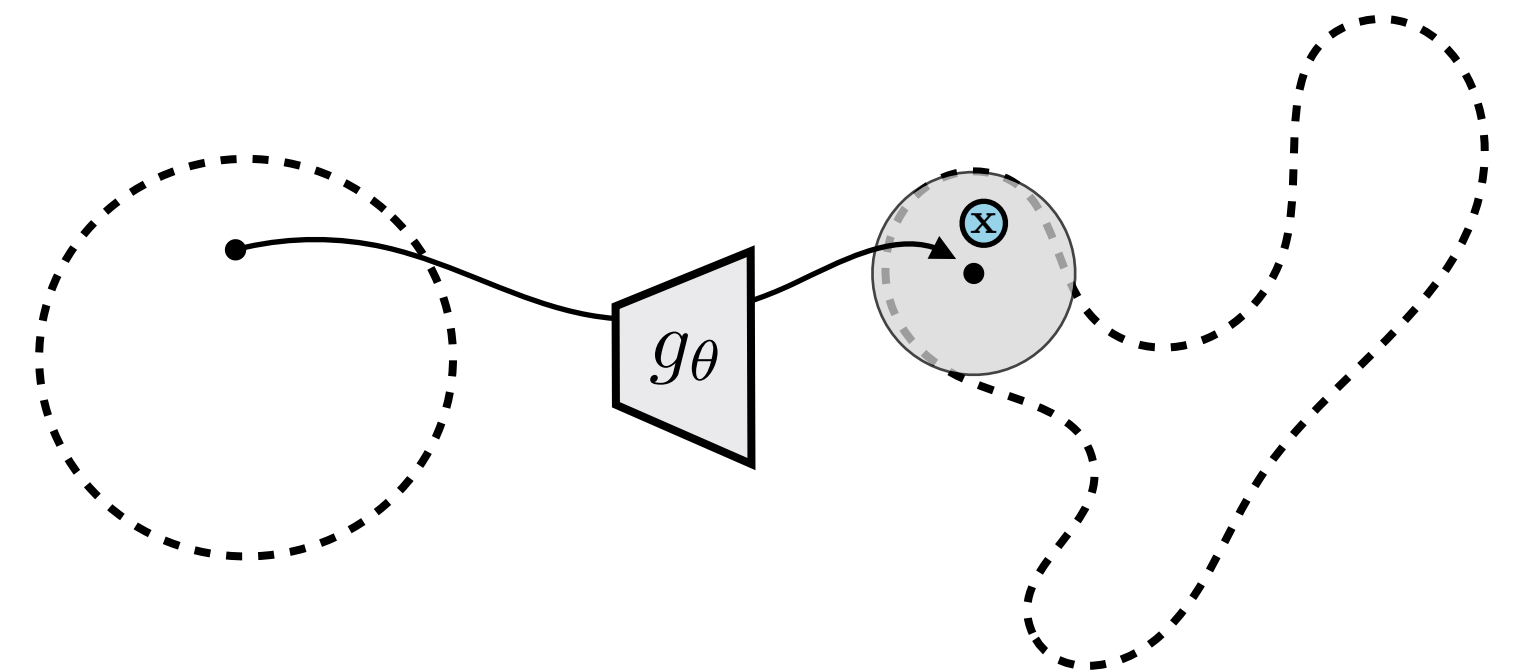
Our learning problem for q:

$$\begin{aligned}\psi^* &= \arg \max_{\psi} \frac{1}{N} \sum_{i=1}^N J_q(\mathbf{x}^{(i)}, \psi) \\ &= \arg \max_{\psi} \frac{1}{N} \sum_{i=1}^N \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x}^{(i)})} [-\log q_{\psi}(\mathbf{z} | \mathbf{x}^{(i)}) + \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) + \log p_{\mathbf{z}}(\mathbf{z})]}_J\end{aligned}$$



Our learning problem for p:

$$\begin{aligned}J_p(\mathbf{x}, \theta) &= \log \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z} | \mathbf{x})} \left[ \frac{p_{\mathbf{z}}(\mathbf{z})}{q_{\psi}(\mathbf{z} | \mathbf{x})} p_{\theta}(\mathbf{x} | \mathbf{z}) \right] \\ \theta^* &= \arg \max_{\theta} \frac{1}{N} \sum_{i=1}^N J_p(\mathbf{x}^{(i)}, \theta)\end{aligned}$$





# Putting everything together

Improving our learning problem for p:

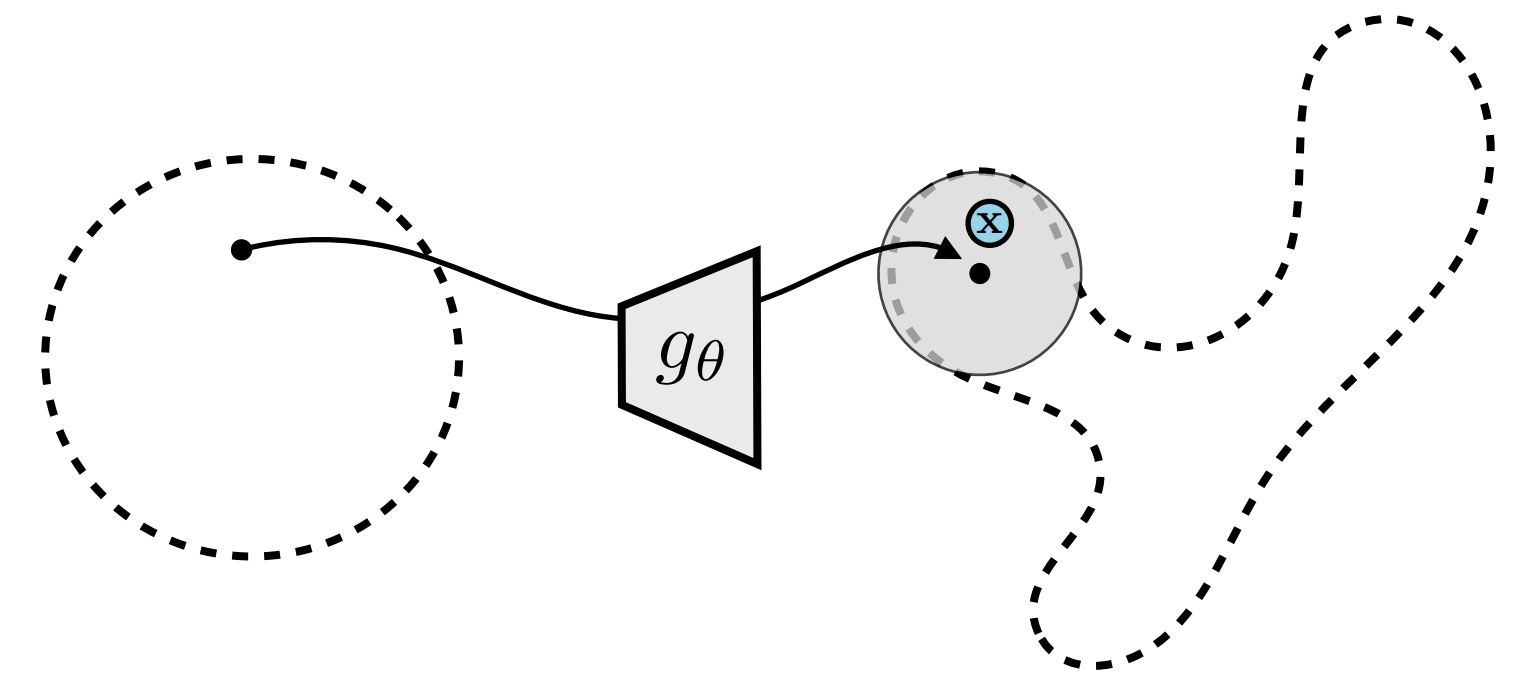
$$J_p(\mathbf{x}, \theta) = \log \mathbb{E}_{\mathbf{z} \sim q_\psi(\mathbf{z} | \mathbf{x})} \left[ \frac{p_{\mathbf{z}}(\mathbf{z})}{q_\psi(\mathbf{z} | \mathbf{x})} p_\theta(\mathbf{x} | \mathbf{z}) \right]$$

$$\geq \mathbb{E}_{\mathbf{z} \sim q_\psi(\mathbf{z} | \mathbf{x})} \left[ \log \left( \frac{p_{\mathbf{z}}(\mathbf{z})}{q_\psi(\mathbf{z} | \mathbf{x})} p_\theta(\mathbf{x} | \mathbf{z}) \right) \right] \quad \triangleleft \text{ Jensen's inequality}$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\psi(\mathbf{z} | \mathbf{x})} \left[ -\log q_\psi(\mathbf{z} | \mathbf{x}) + \log p_\theta(\mathbf{x} | \mathbf{z}) + \log p_{\mathbf{z}}(\mathbf{z}) \right]$$

$$= J \quad \triangleleft \text{ VAE objective}$$

$$= \mathbb{E}_{\mathbf{z} \sim q_\psi(\mathbf{z} | \mathbf{x})} \left[ \log p_\theta(\mathbf{x} | \mathbf{z}) \right] - \text{KL}(q_\psi(\mathbf{z} | \mathbf{x}) \parallel p_{\mathbf{z}})$$

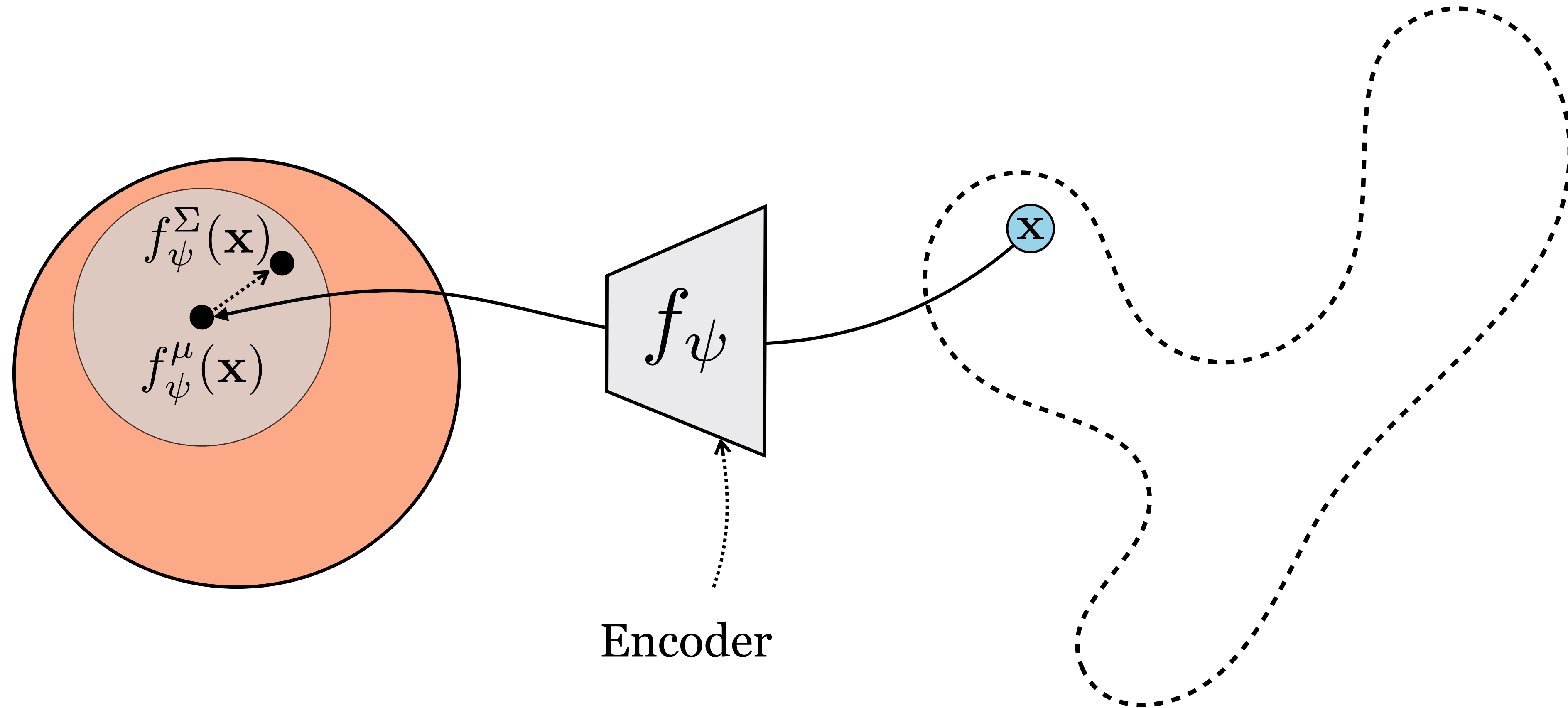


← **J = Evidence Lower-Bound (ELBO)**

Now, p and q share the same exact objective!

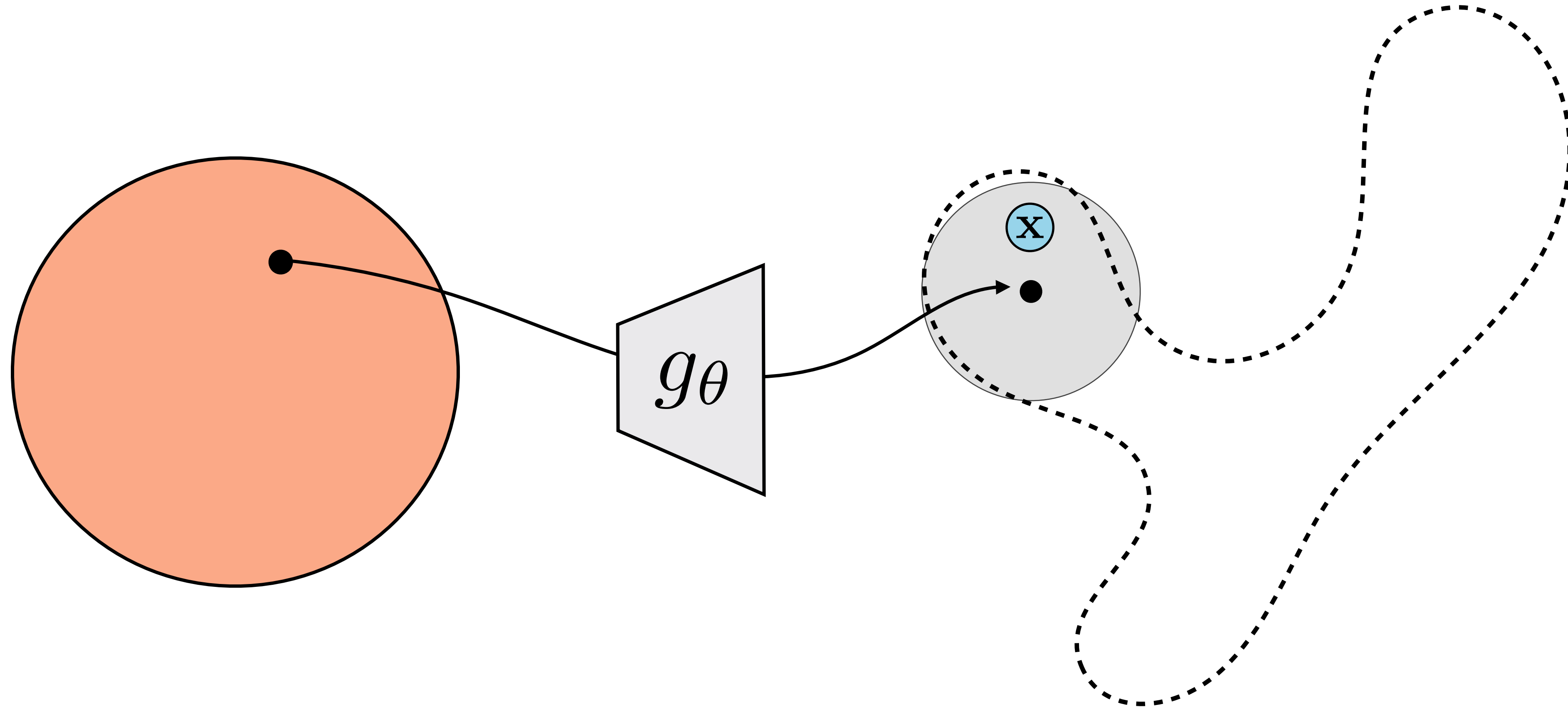
$$\theta^*, \psi^* = \arg \max_{\theta, \phi} \frac{1}{N} \sum_{i=1}^N J(\mathbf{x}^{(i)}, \theta, \phi)$$

Trick #3: Predict the optimal sampling distribution for each given datapoint



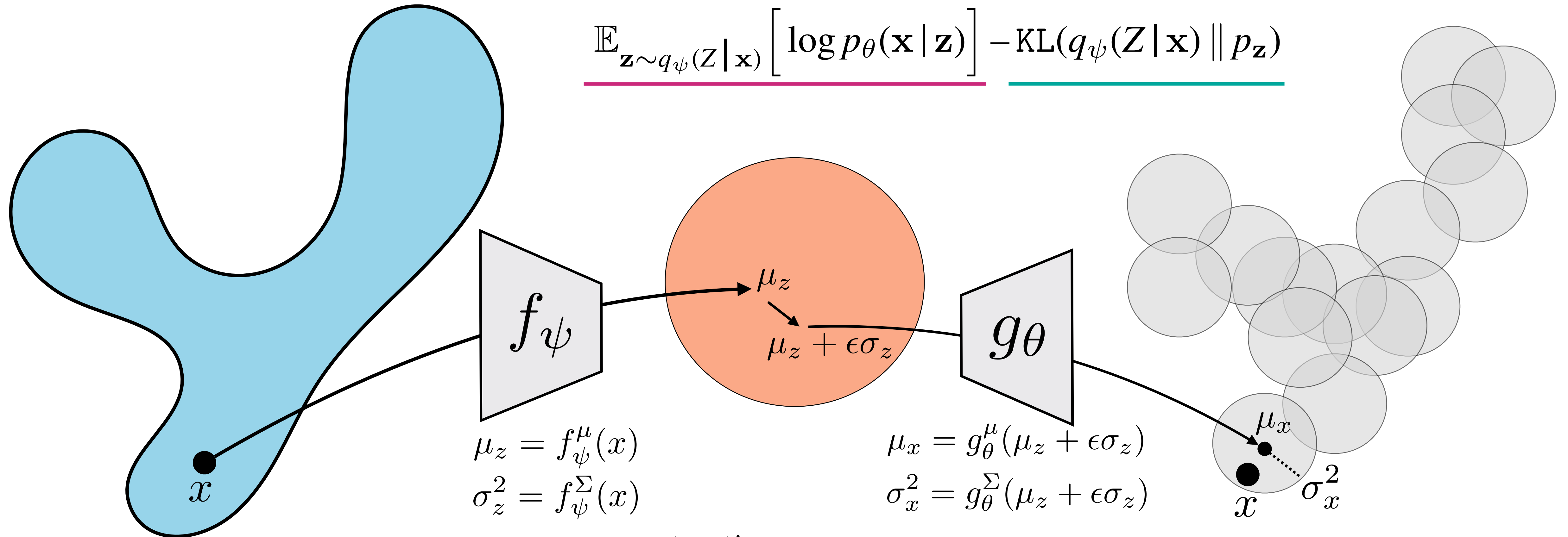
$$p_{\theta}(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim q_{\psi}(\mathbf{z}|\mathbf{x})} \left[ \frac{p_{\mathbf{z}}(\mathbf{z})}{q_{\psi}(\mathbf{z}|\mathbf{x})} p_{\theta}(\mathbf{x}|\mathbf{z}) \right]$$

Trick #3: Predict the optimal sampling distribution for each given datapoint



$$p_\theta(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim q_\psi(Z|\mathbf{x})} \left[ \frac{p_{\mathbf{z}}(\mathbf{z})}{q_\psi(\mathbf{z}|\mathbf{x})} p_\theta(\mathbf{x}|\mathbf{z}) \right]$$

# Looks like an autoencoder!



$$\mathbb{E}_{\mathbf{z} \sim q_\psi(Z | \mathbf{x})} \left[ \log p_\theta(\mathbf{x} | \mathbf{z}) \right] - \text{KL}(q_\psi(Z | \mathbf{x}) \parallel p_{\mathbf{z}})$$

$$\mu_z = f_\psi^\mu(x)$$

$$\sigma_z^2 = f_\psi^\Sigma(x)$$

$$\mu_x = g_\theta^\mu(\mu_z + \epsilon \sigma_z)$$

$$\sigma_x^2 = g_\theta^\Sigma(\mu_z + \epsilon \sigma_z)$$

reconstruction error

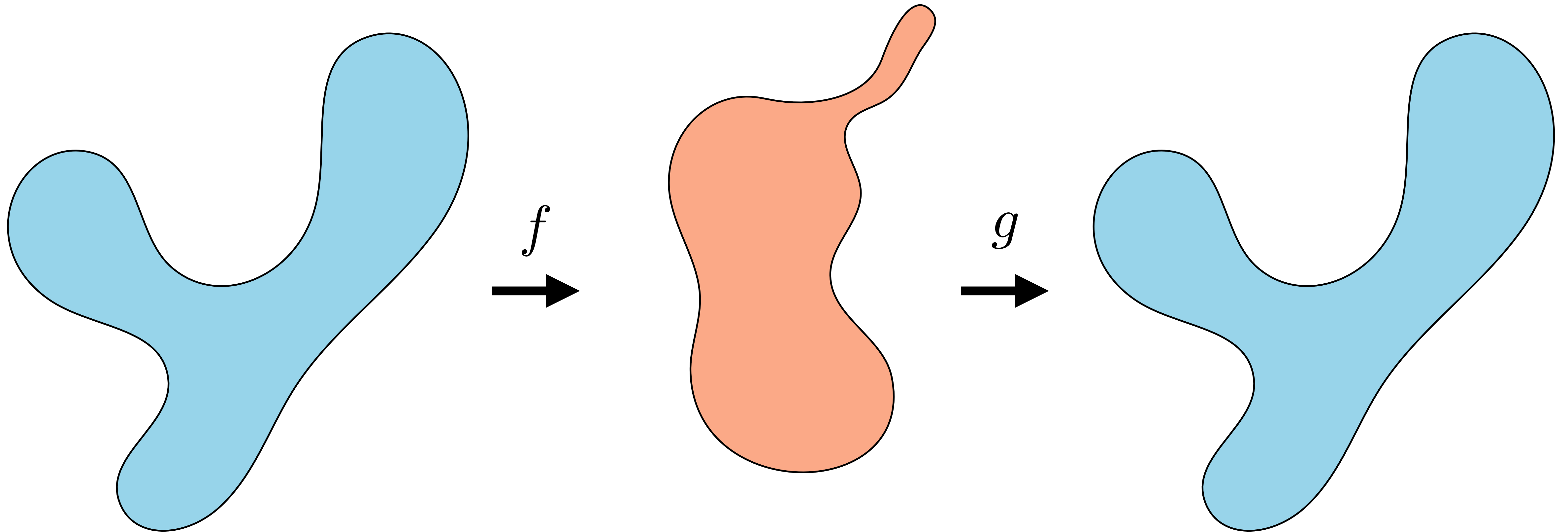
$$\log \frac{1}{\sigma_x \sqrt{2\pi}} - \frac{(x - \mu_x)^2}{2\sigma_x^2} - \frac{1}{2}(\mu_z^2 + \sigma_z^2 - \log(\sigma_z^2) - 1)$$

(eqns for 1D  $x$  and  $z$ )

data likelihood under 1 importance sample

squash encodings toward origin

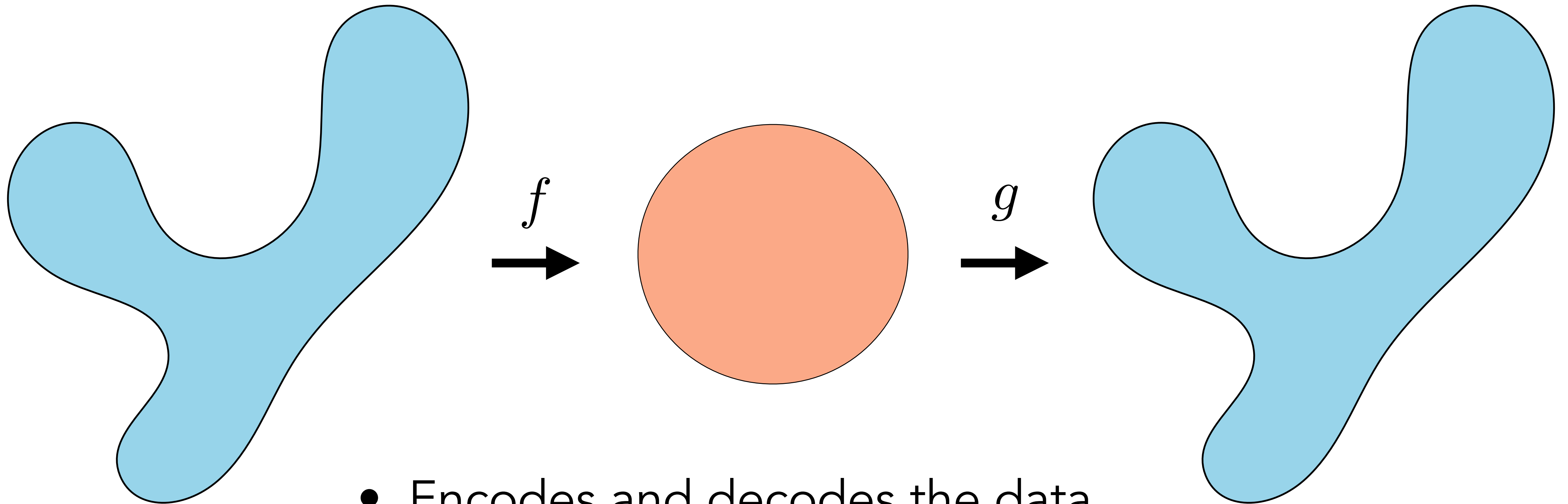
# Autoencoder



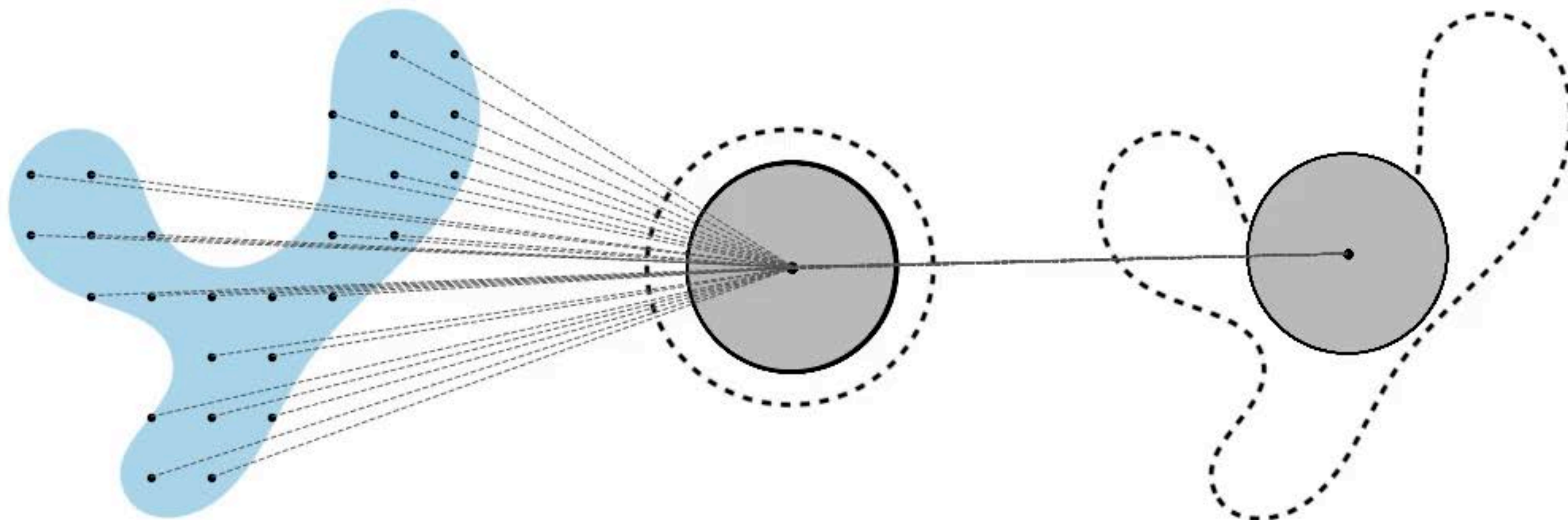
- Encodes and decodes the data
- Low-dimensional bottleneck

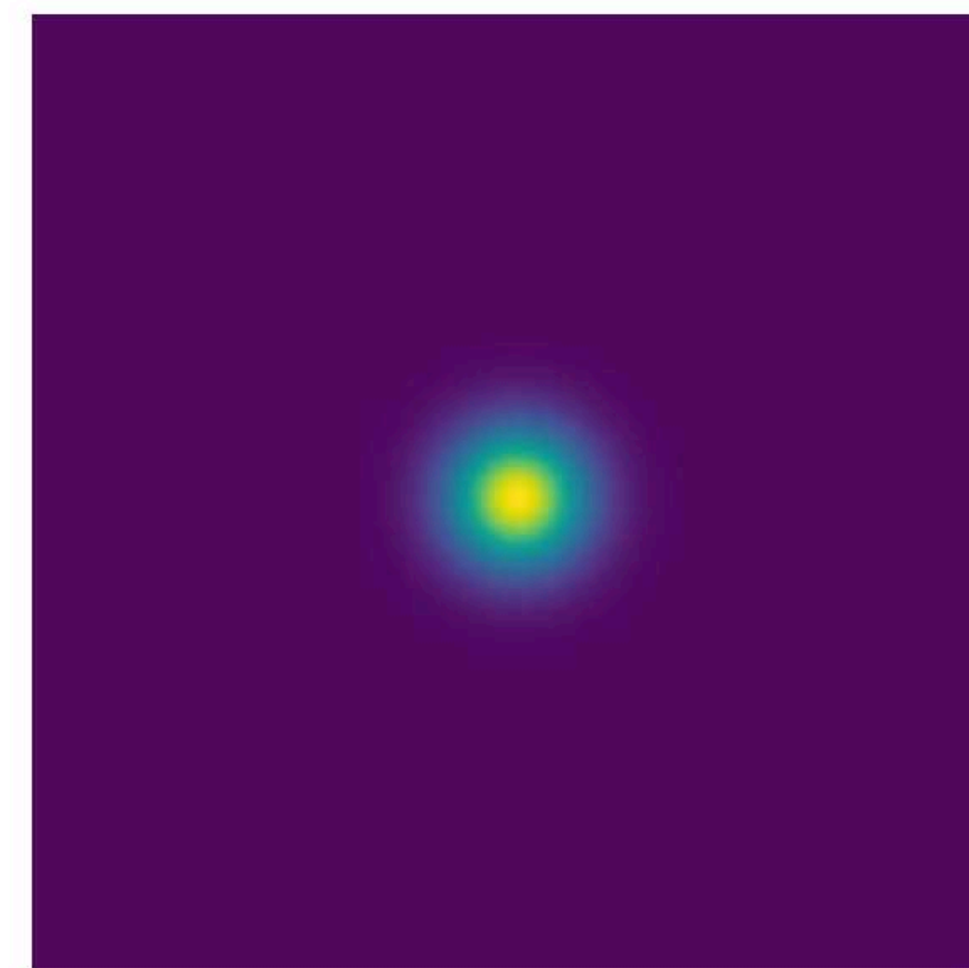
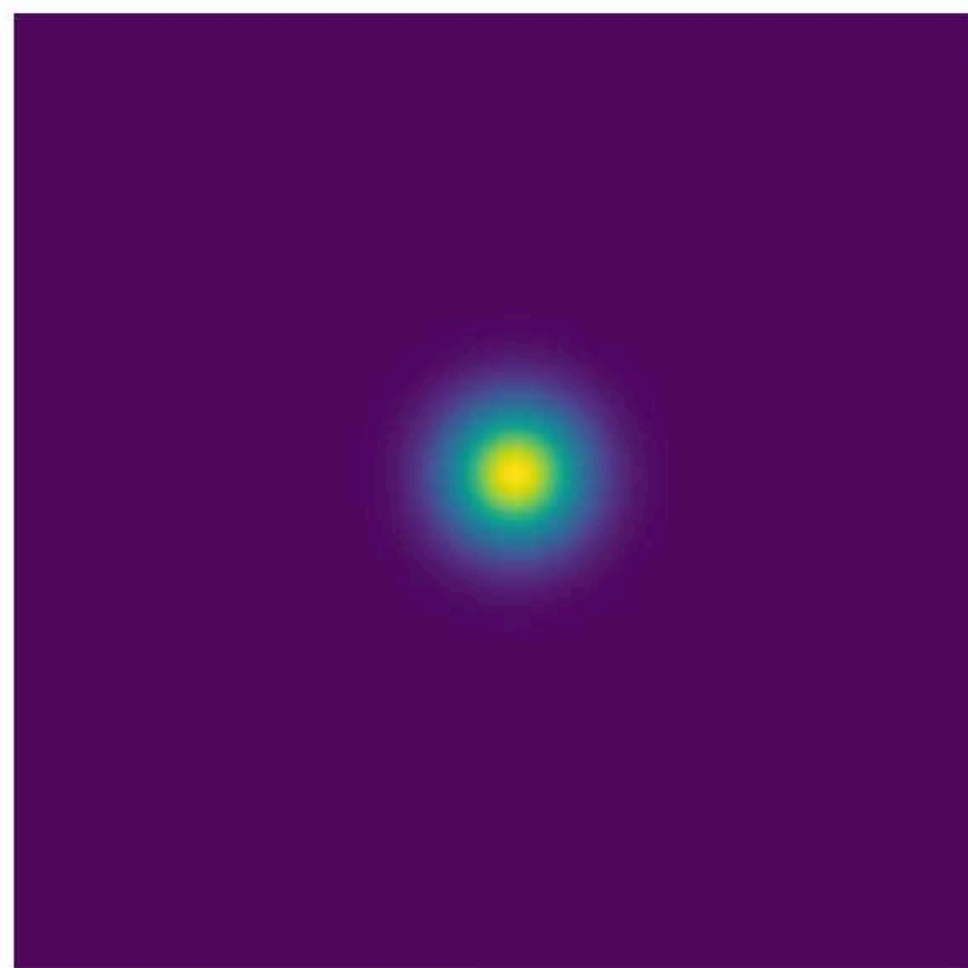
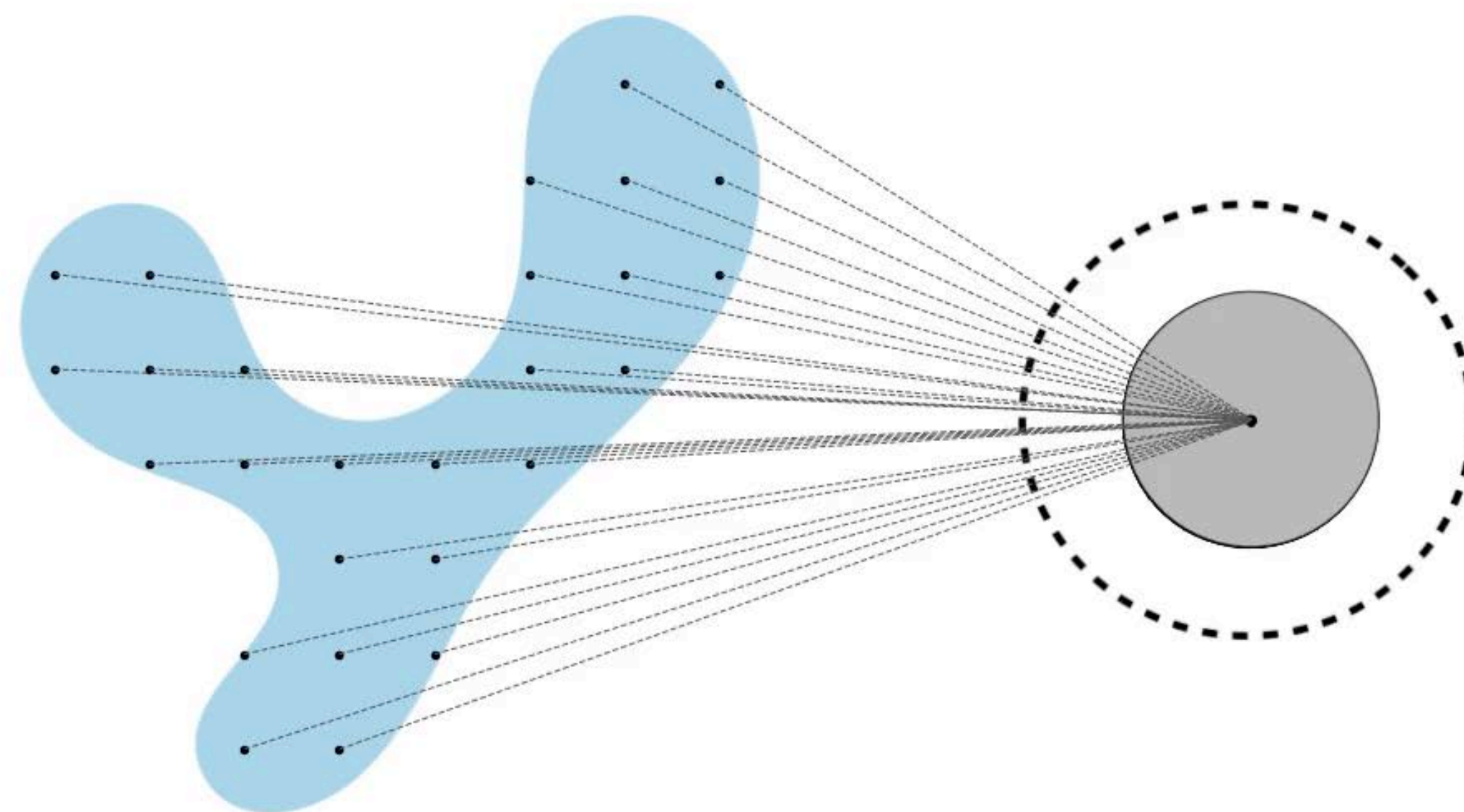
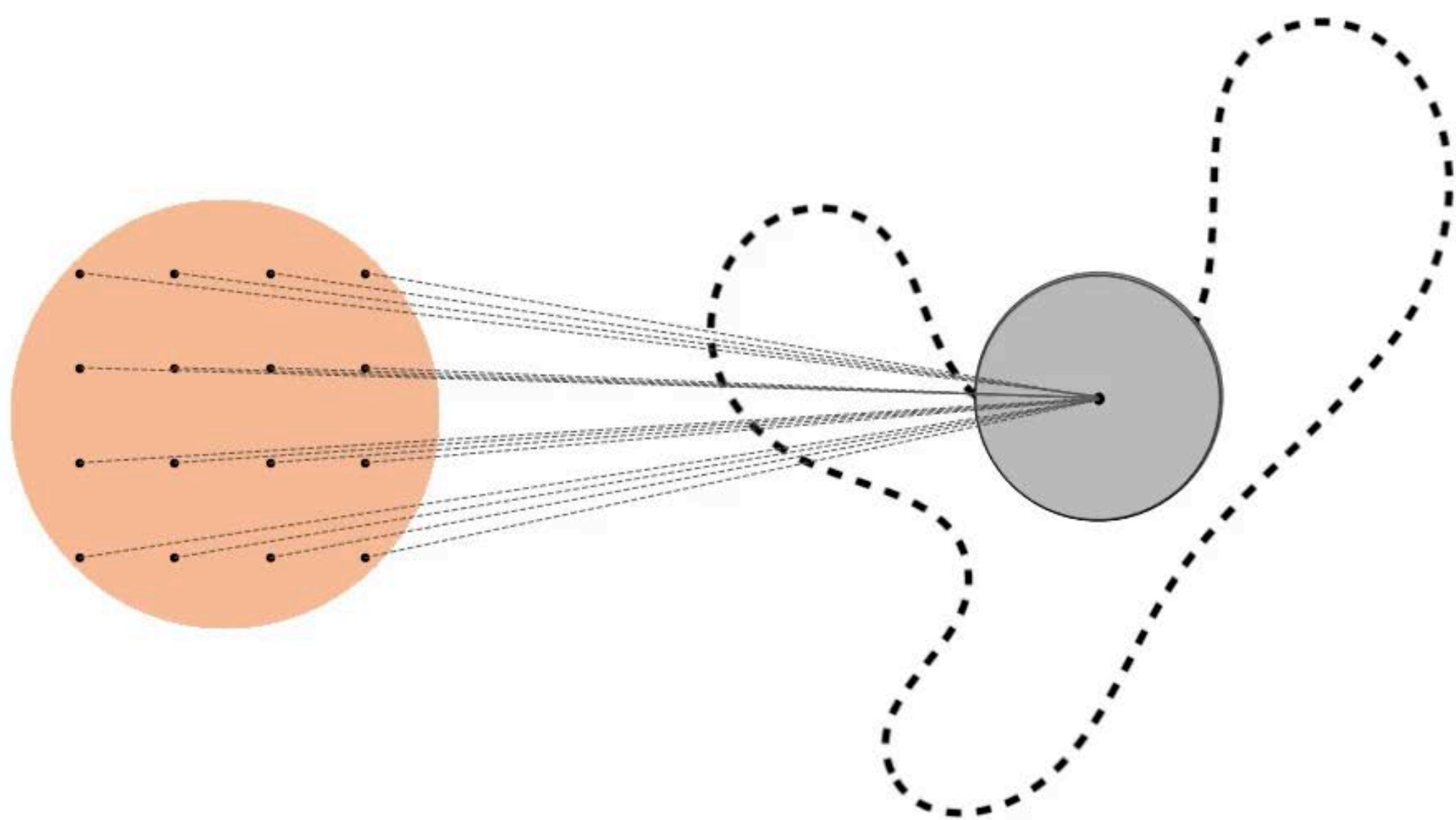


# Variational Autoencoder



- Encodes and decodes the data
- Low-dimensional bottleneck
- Gaussian bottleneck (**can sample; disentangled**)





# VAE — three tricks

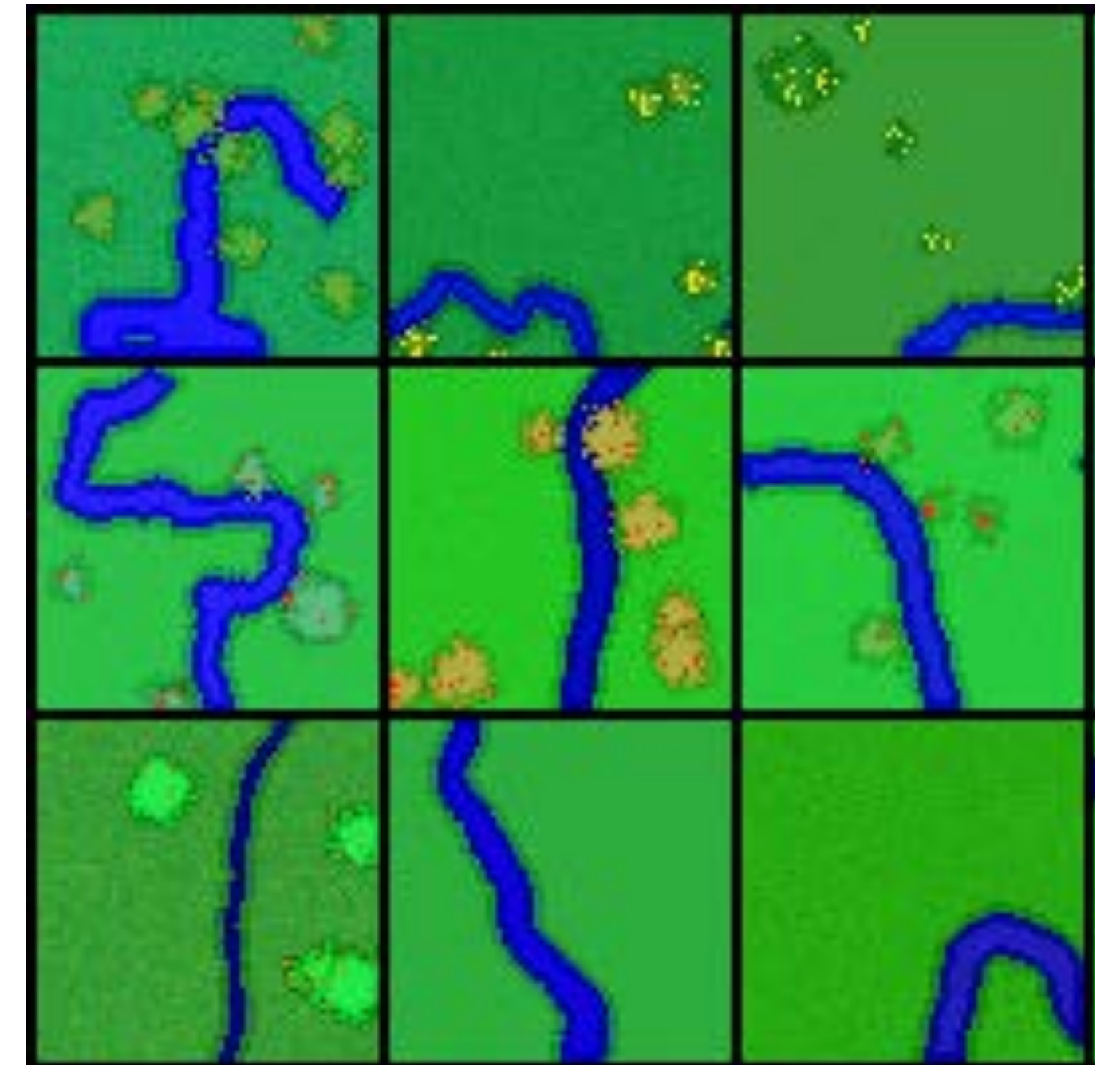
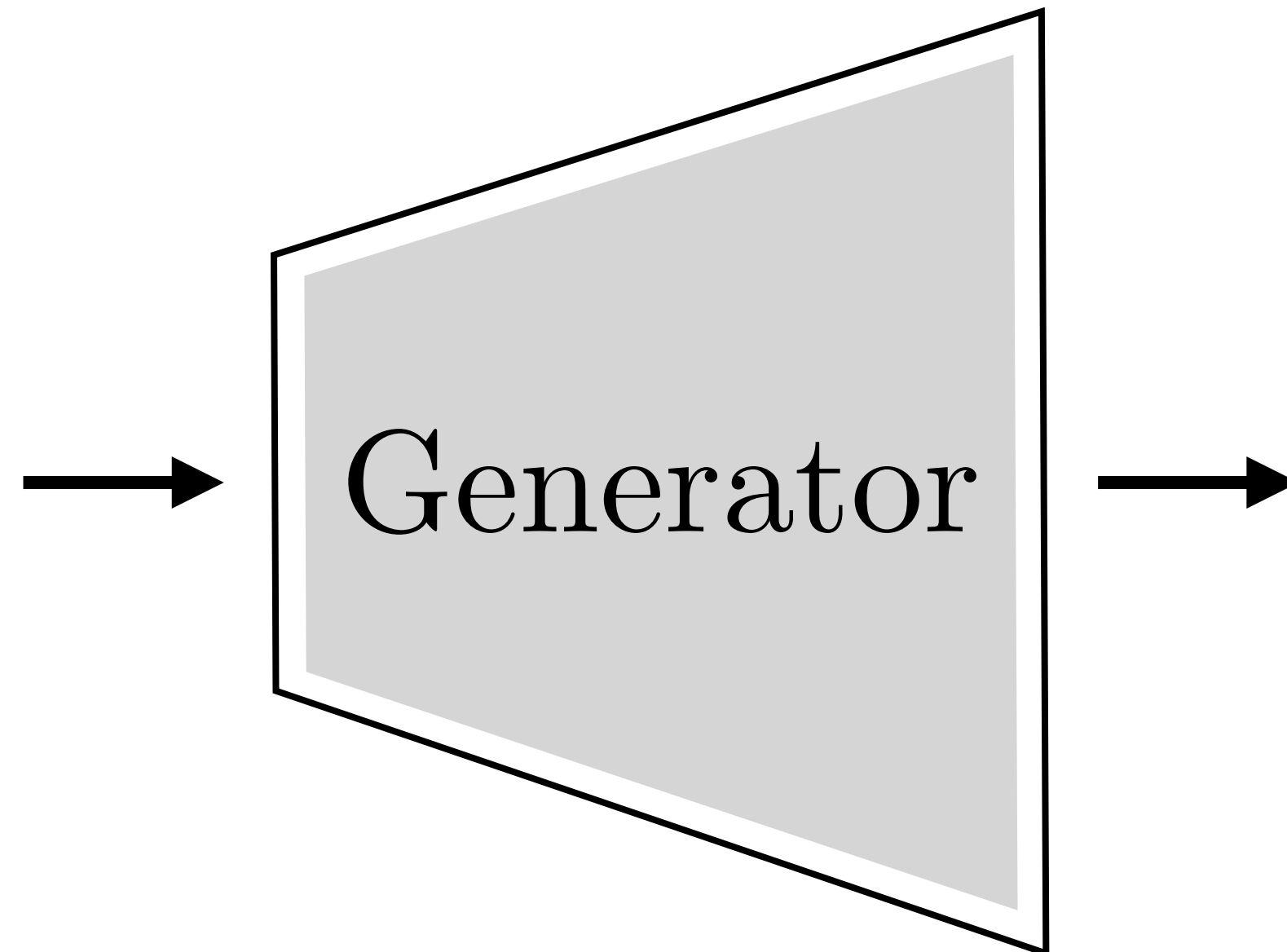
- Trick #1: approximate an infinite mixture with samples
- Trick #2: sample efficiently via importance sampling
- Trick #3: predict the optimal sampling distribution for each datapoint

# Training a VAE — Step by step

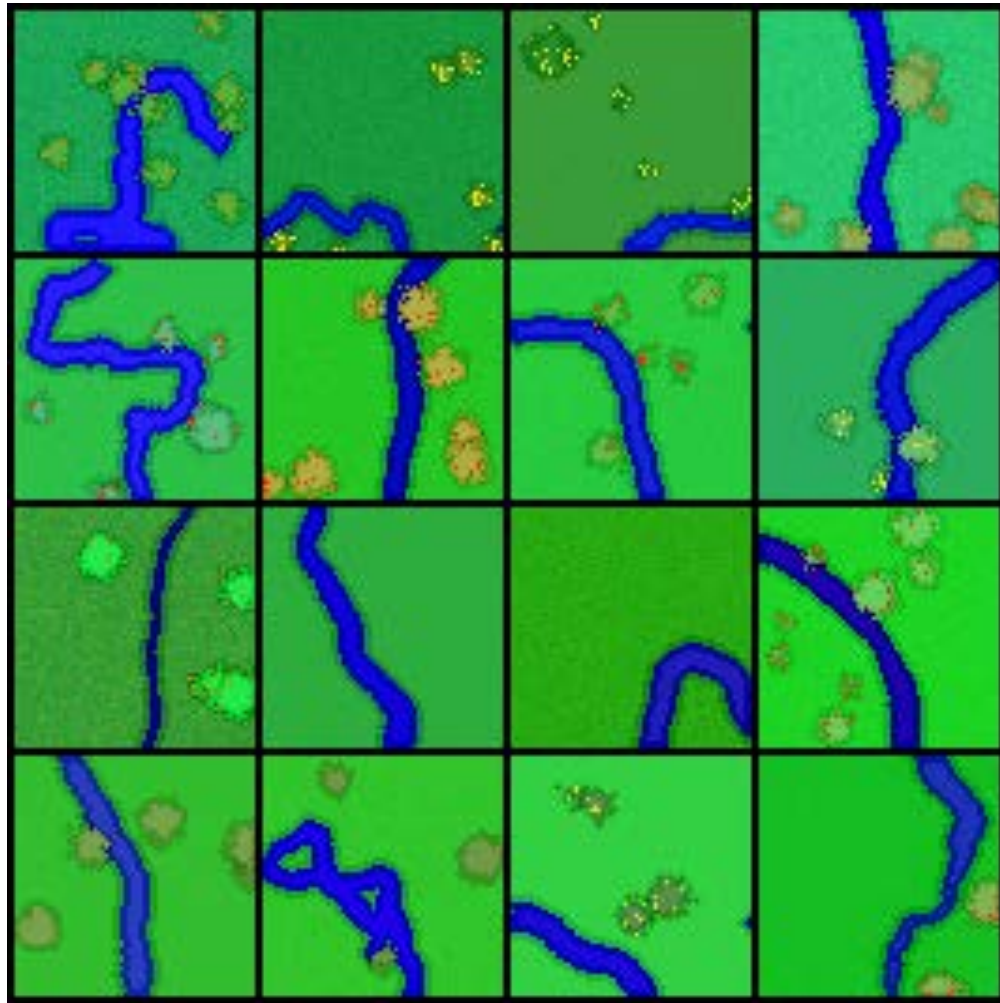
- Sample one or more  $\mathbf{x} \sim \{\mathbf{x}^{(i)}\}_{i=1}^N$
- *Encode* the data with a forward pass through  $f_\psi$
- For each datapoint, create one or more noisy latent codes using the distribution parameterized by the encoder
- *Decode* the data by passing the noisy latent codes through the  $g_\theta$
- Compute the losses and backprop to update  $\theta$  and  $\psi$



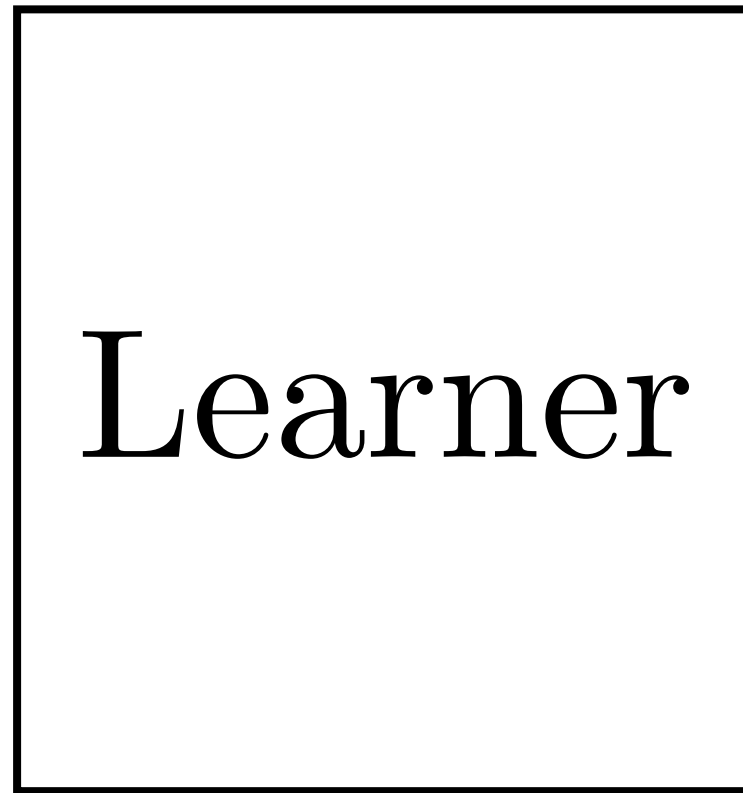
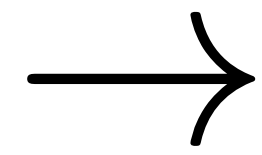
$\mathbf{z}_1 \sim \text{Bernoulli}(0.5)$     (River turns)  
 $\mathbf{z}_2 \sim \text{Normal}(\mu_1, \Sigma_1)$     (Grass color)  
 $\mathbf{z}_3 \sim \text{Unif}(0, 10)$     (Number trees)  
 $\vdots$



Training



Data



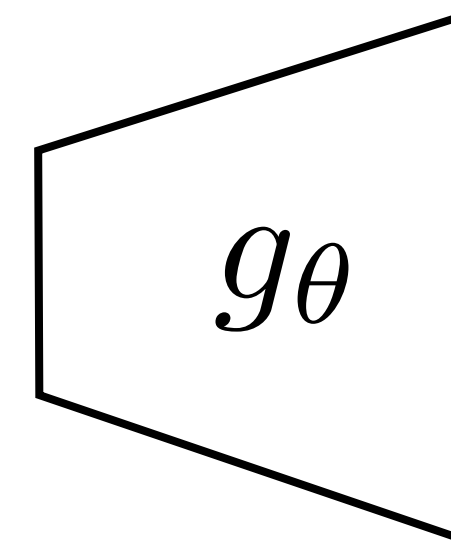
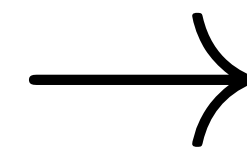
Learner

$\rightarrow f_\psi, g_\theta$

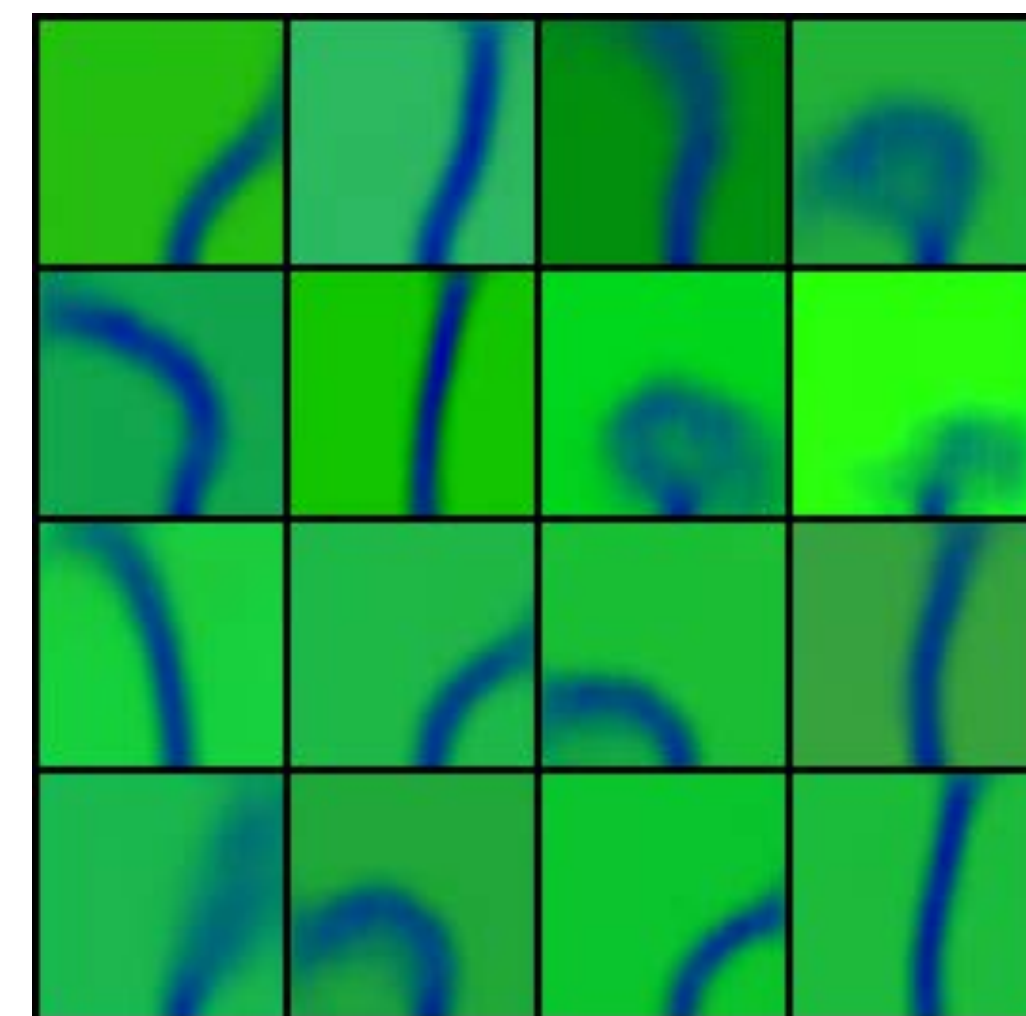
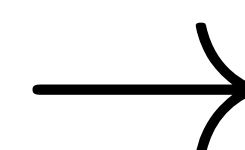


Sampling

$\mathbf{z}$



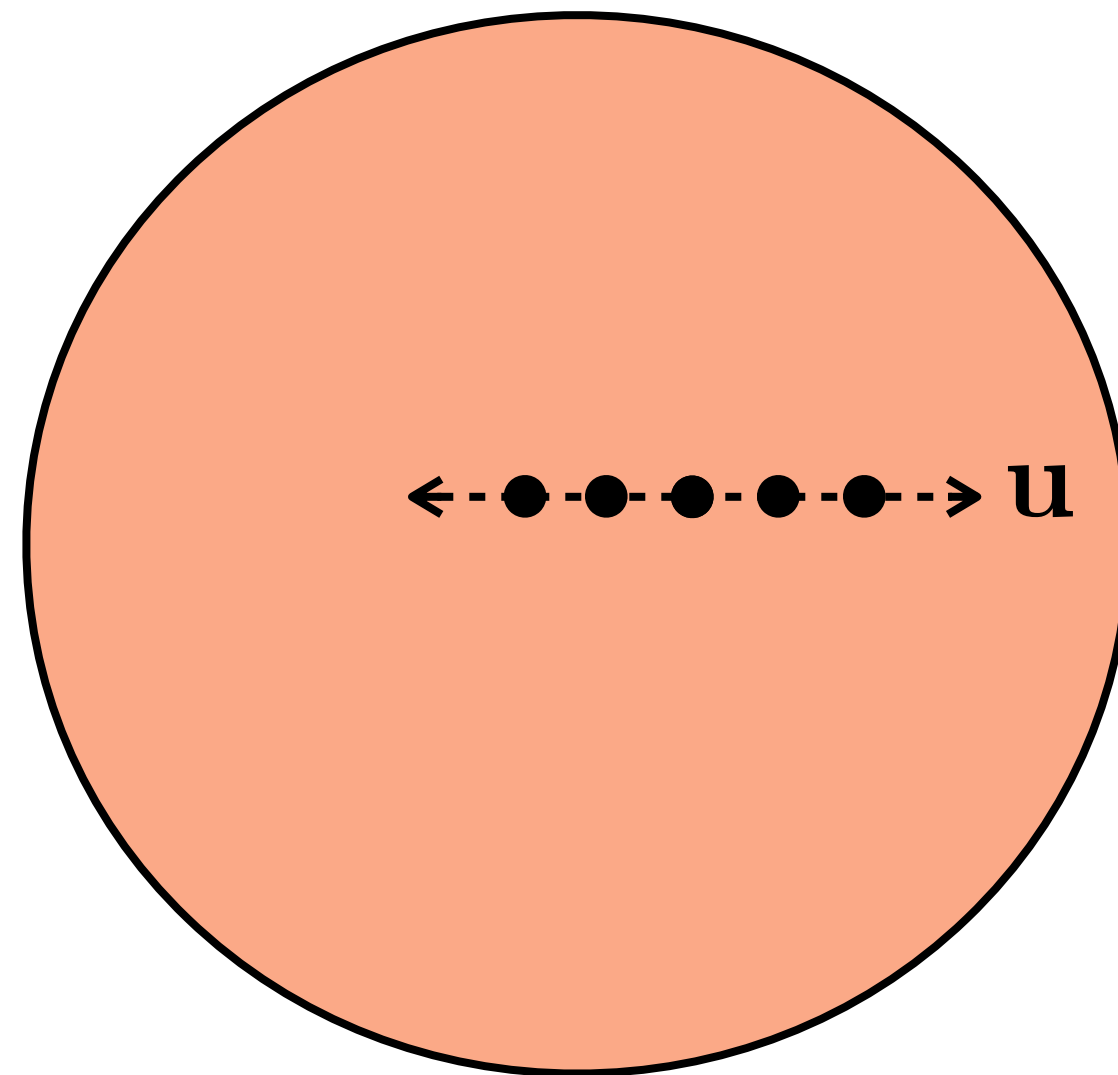
$g_\theta$



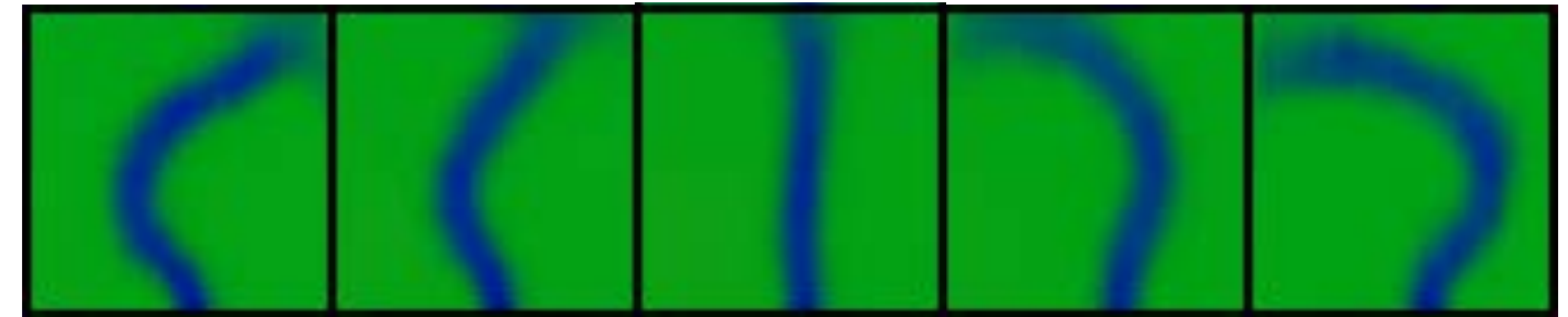
Samples

# Synthesized Images

Latent variables  
(controls)



$\mathbf{v}$



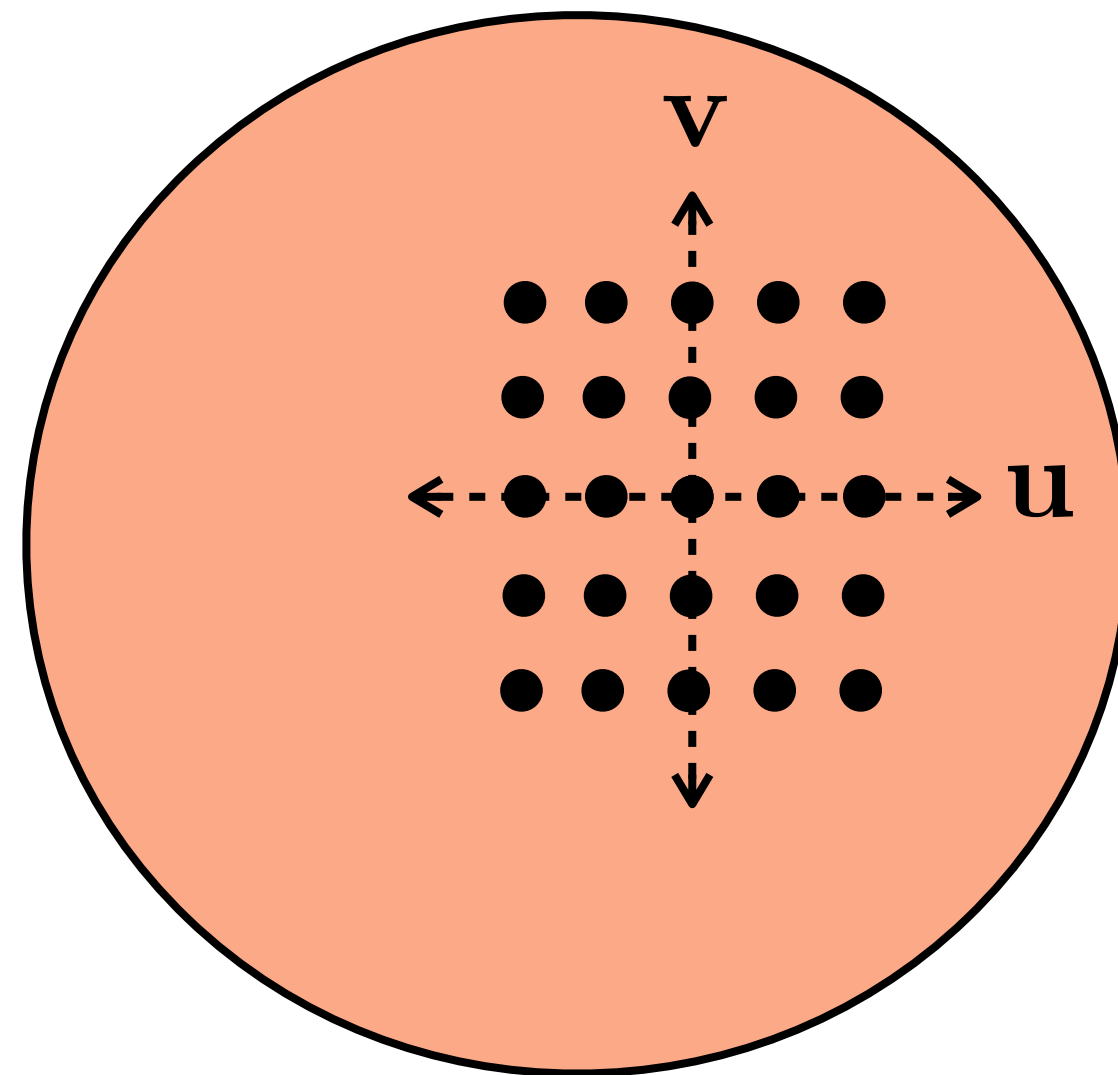
$\mathbf{u}$

“River curvature”



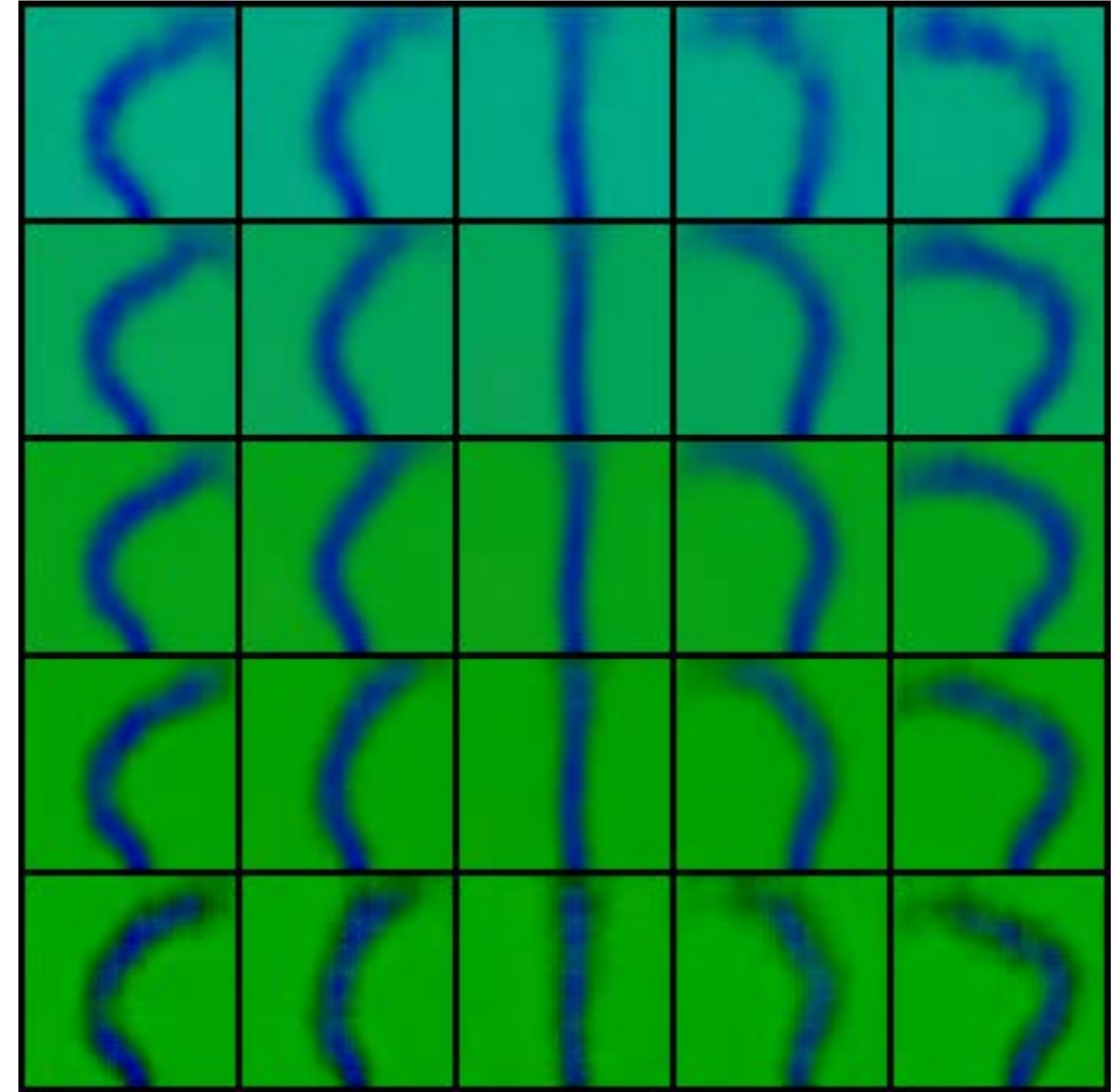
# Synthesized Images

Latent variables  
(controls)



“Background color”

$v$

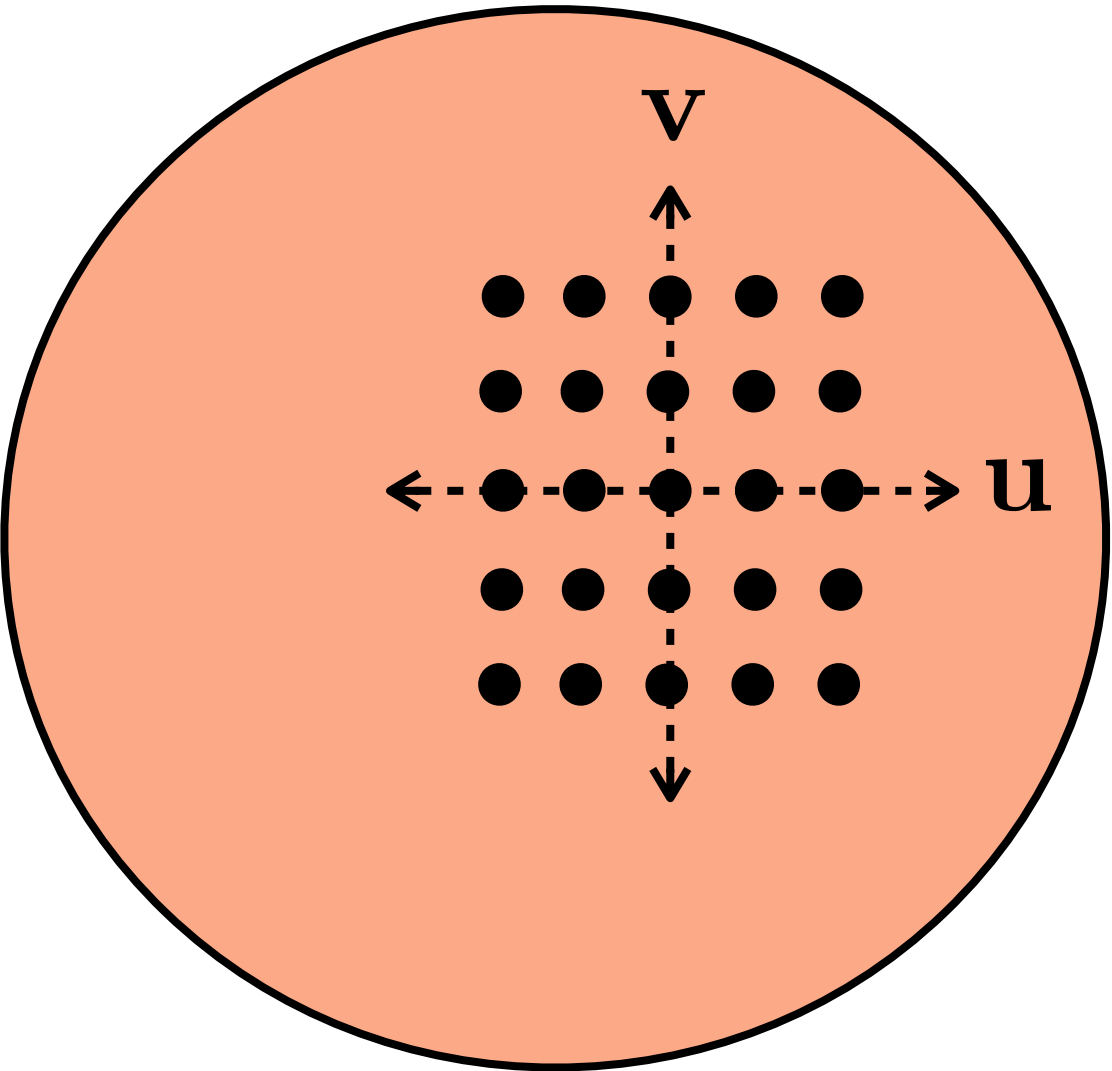


“River curvature”

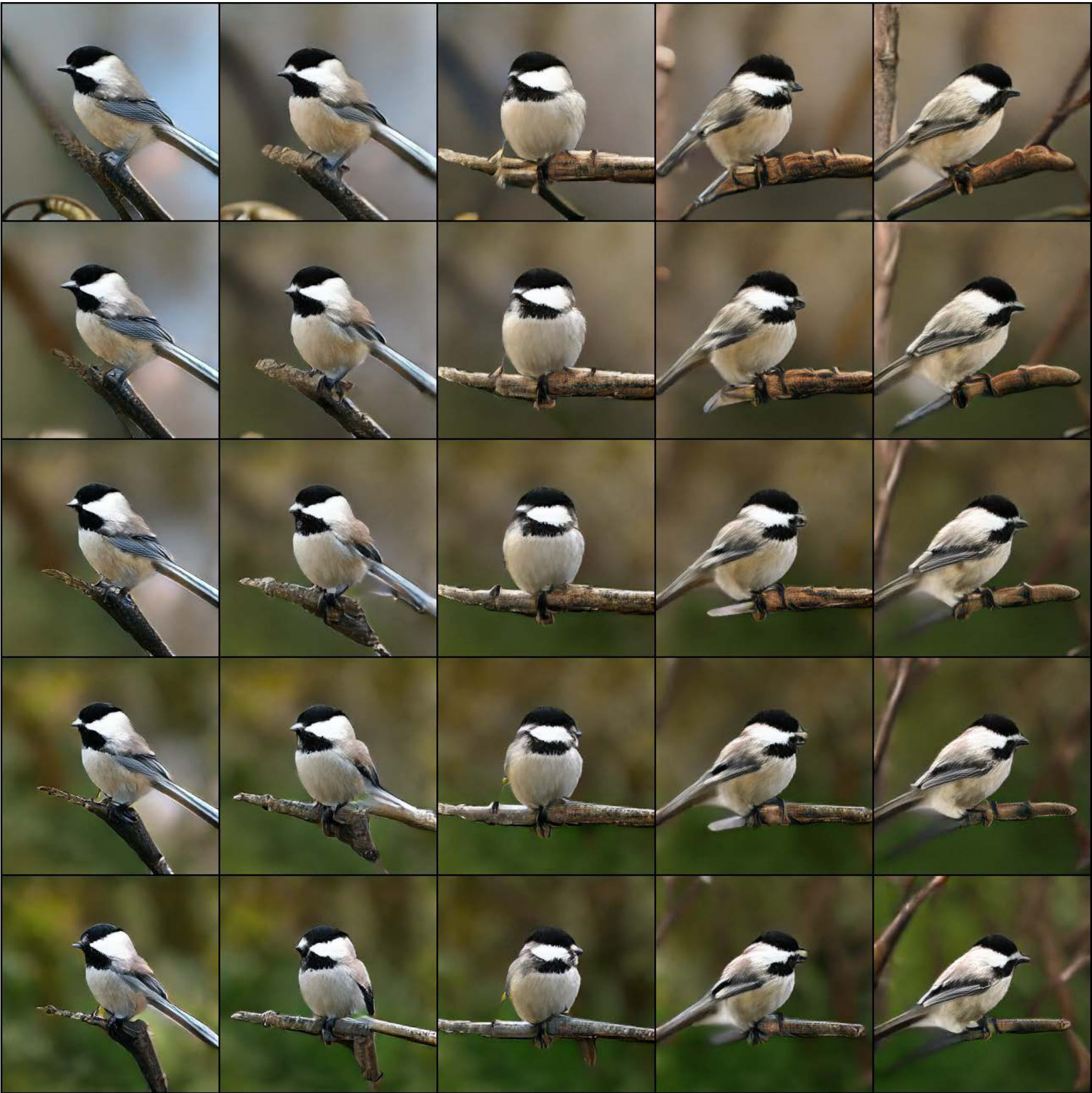
$u$



Latent variables  
(controls)

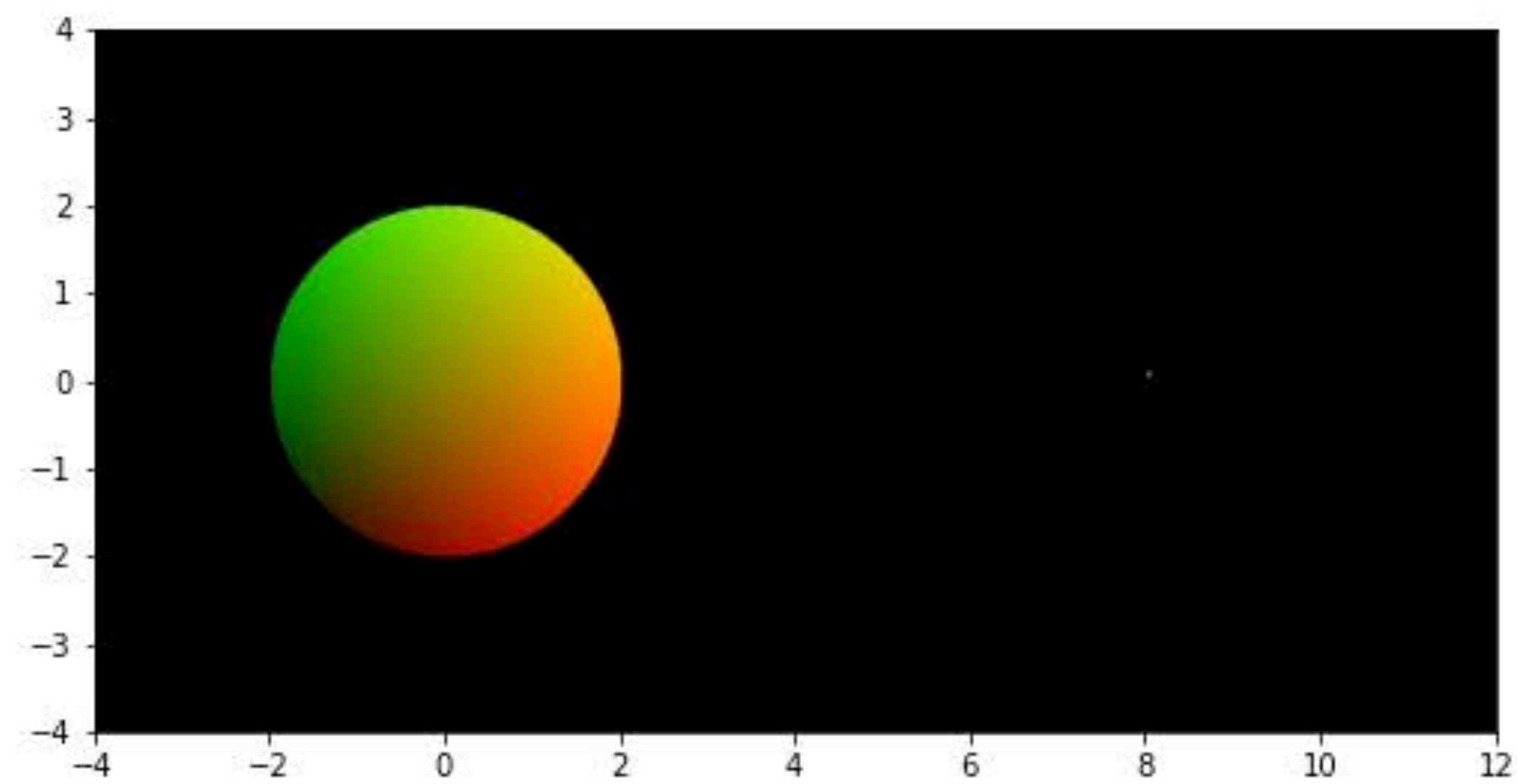
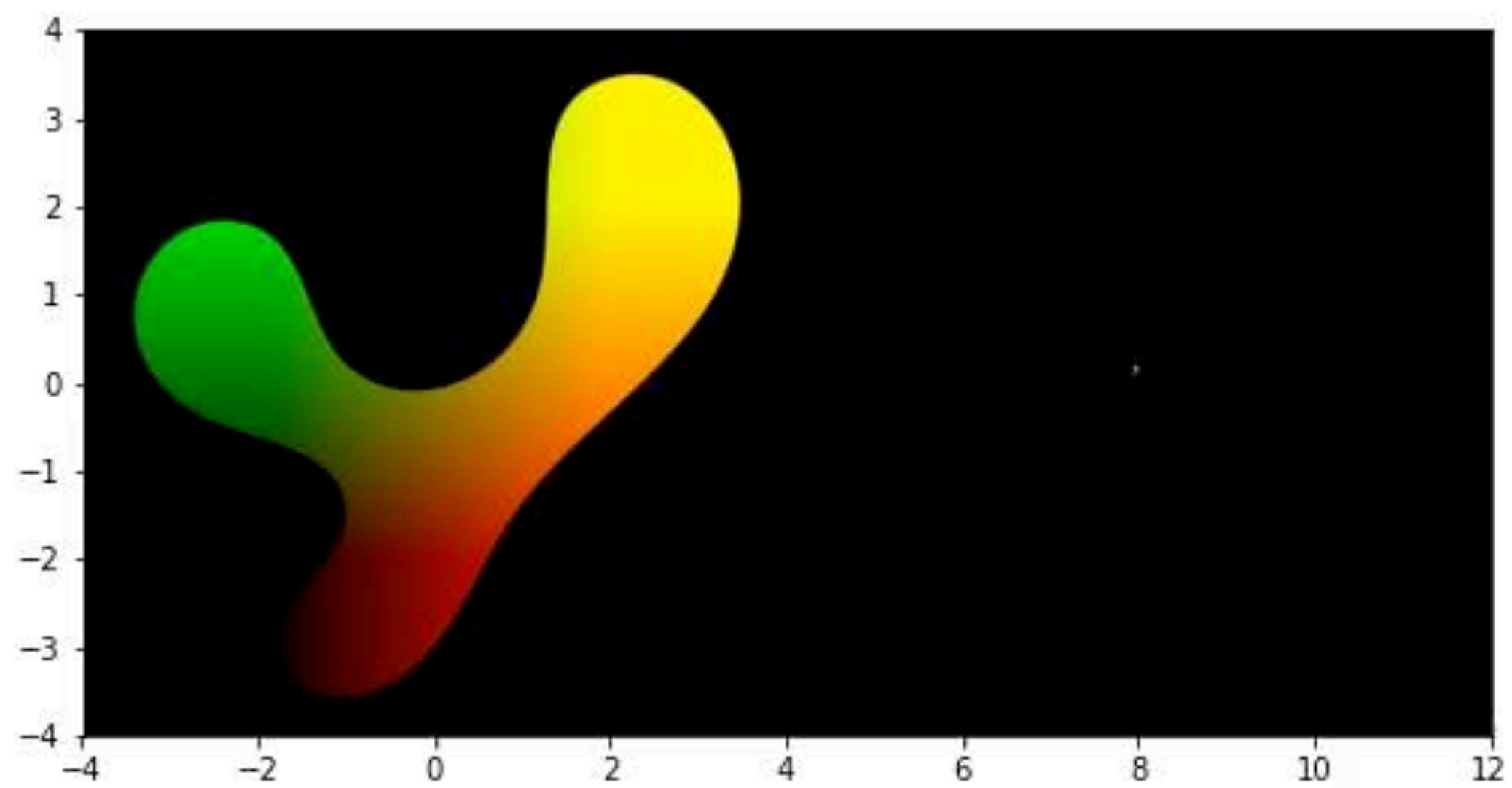
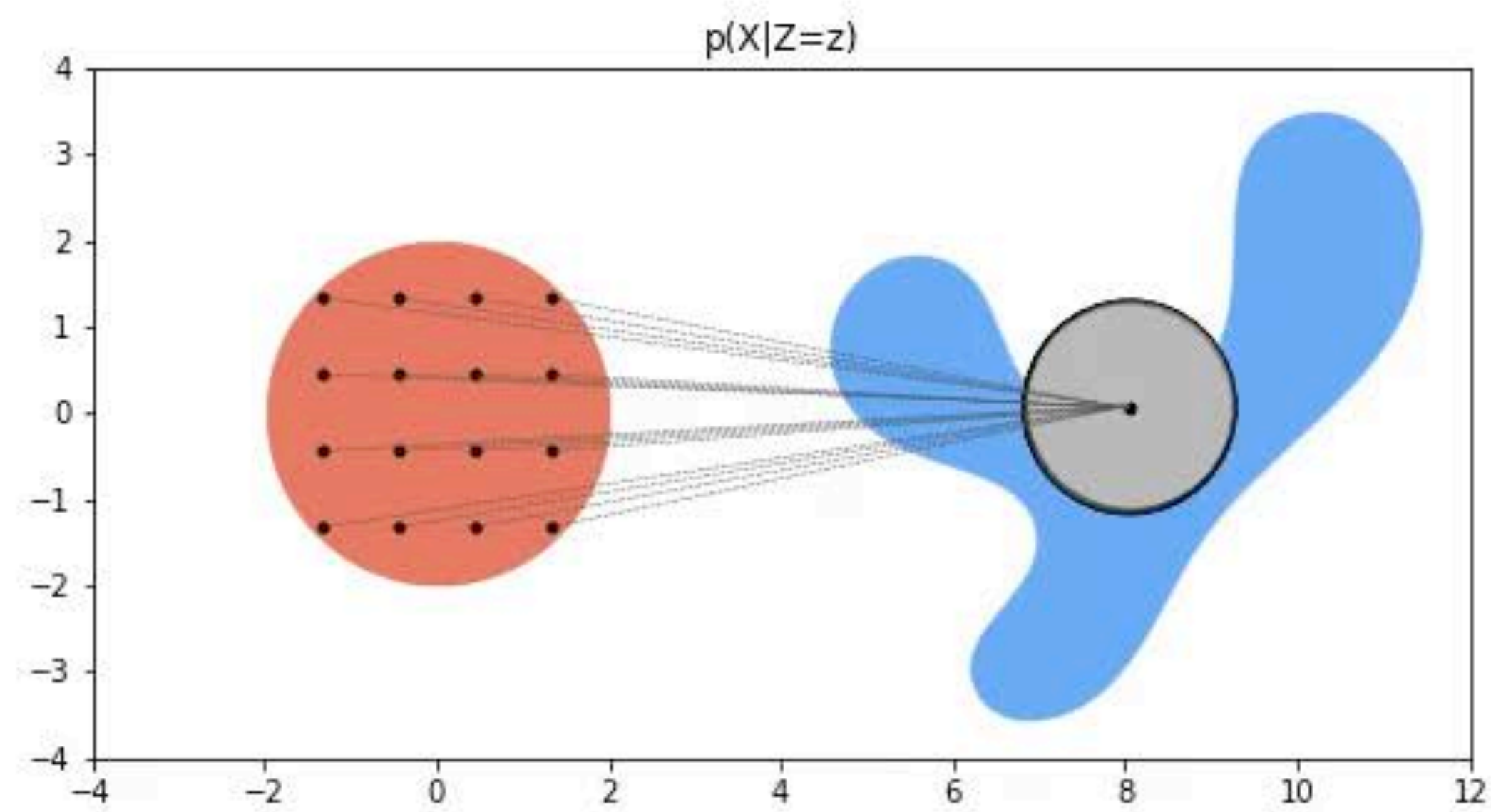
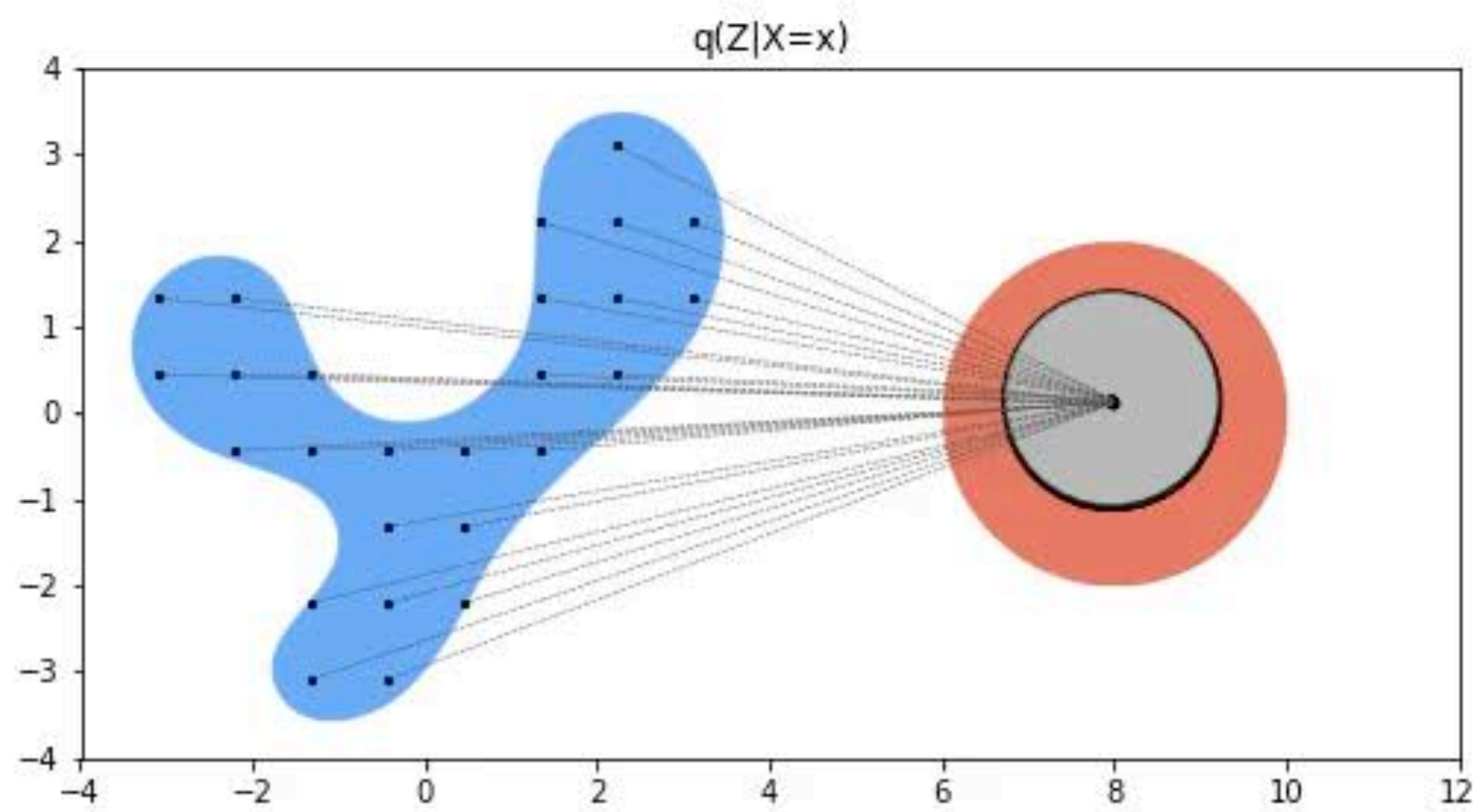


“Background color”



“Bird orientation”







# Comparing popular model types

Method	Latents?	Density/Energy?	Generator?
Energy-based models	<b>X</b>	✓ (energy only)	<b>X</b>
Gaussian	<b>X</b>	✓	<b>X</b>
Autoregressive models	<b>X</b>	✓	✓ (slow)
Diffusion models	✓ (high-dimensional)	<b>X</b>	✓ (slow)
GANs	✓	<b>X</b>	✓
VAEs	✓	<b>X</b>	✓

\* rough categorization of vanilla versions of each model

MIT OpenCourseWare

<https://ocw.mit.edu>

6.7960 Deep Learning

Fall 2024

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>