

# Tutorial on Variational Autoencoders

**David Nagy** @dvgnagy

**David Szepesvari**

**EEML 2021, Budapest**



# Why are VAEs interesting?

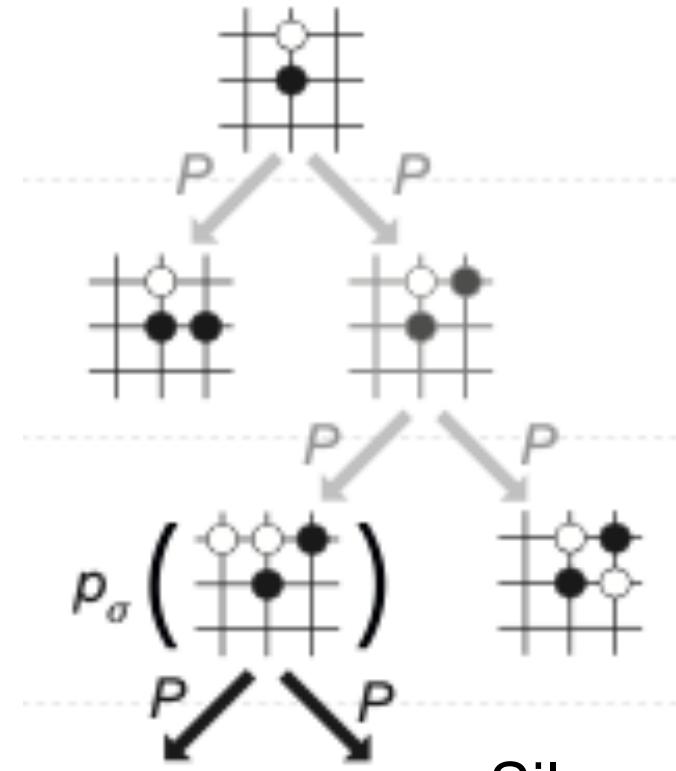
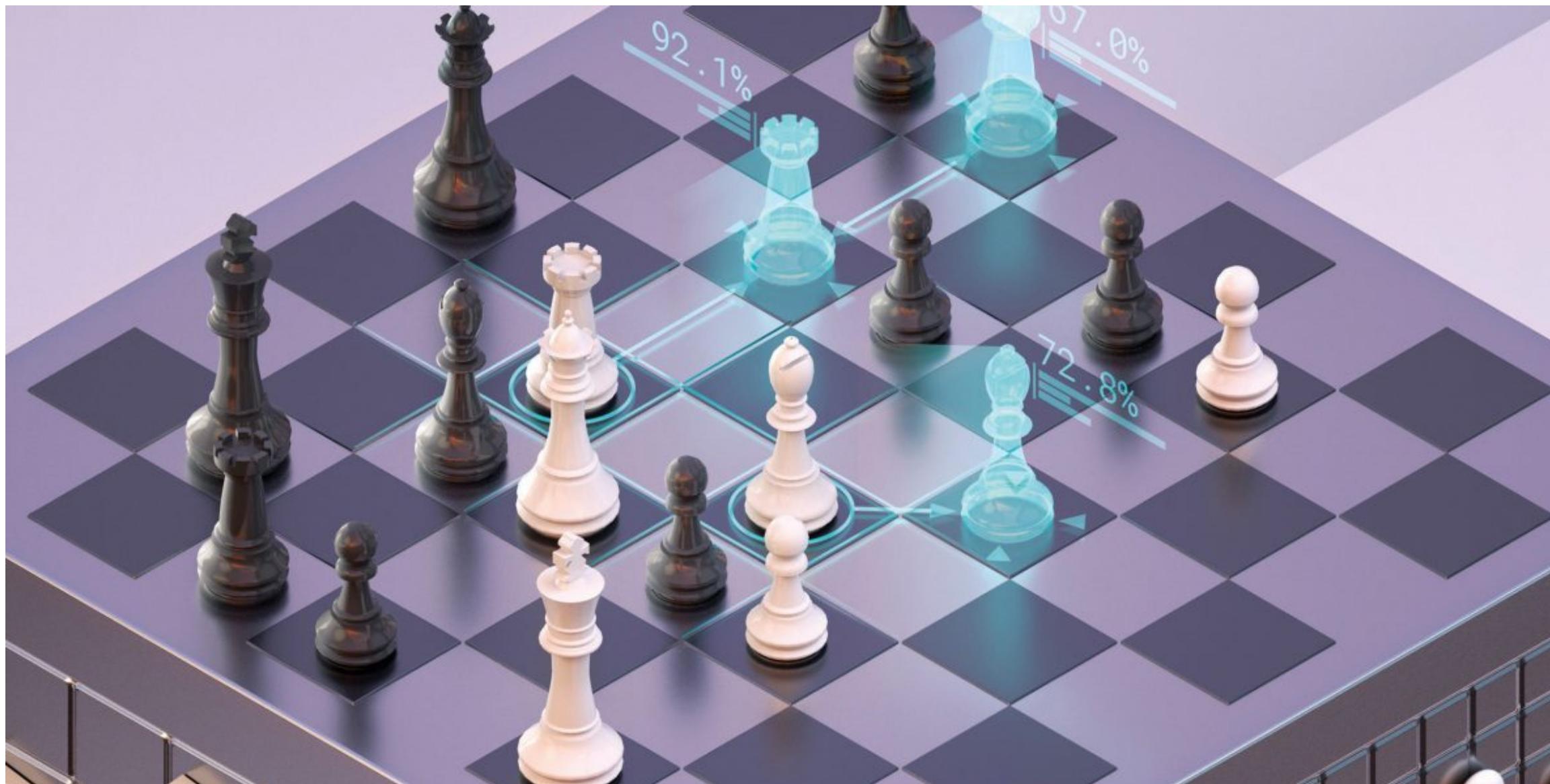
- a single model that can both
  - work as a generative model
  - infer latent variables
- both
  - theoretical interpretation in probabilistic modeling framework
  - but also works on naturalistic data
- a building block in many larger models

# generative models

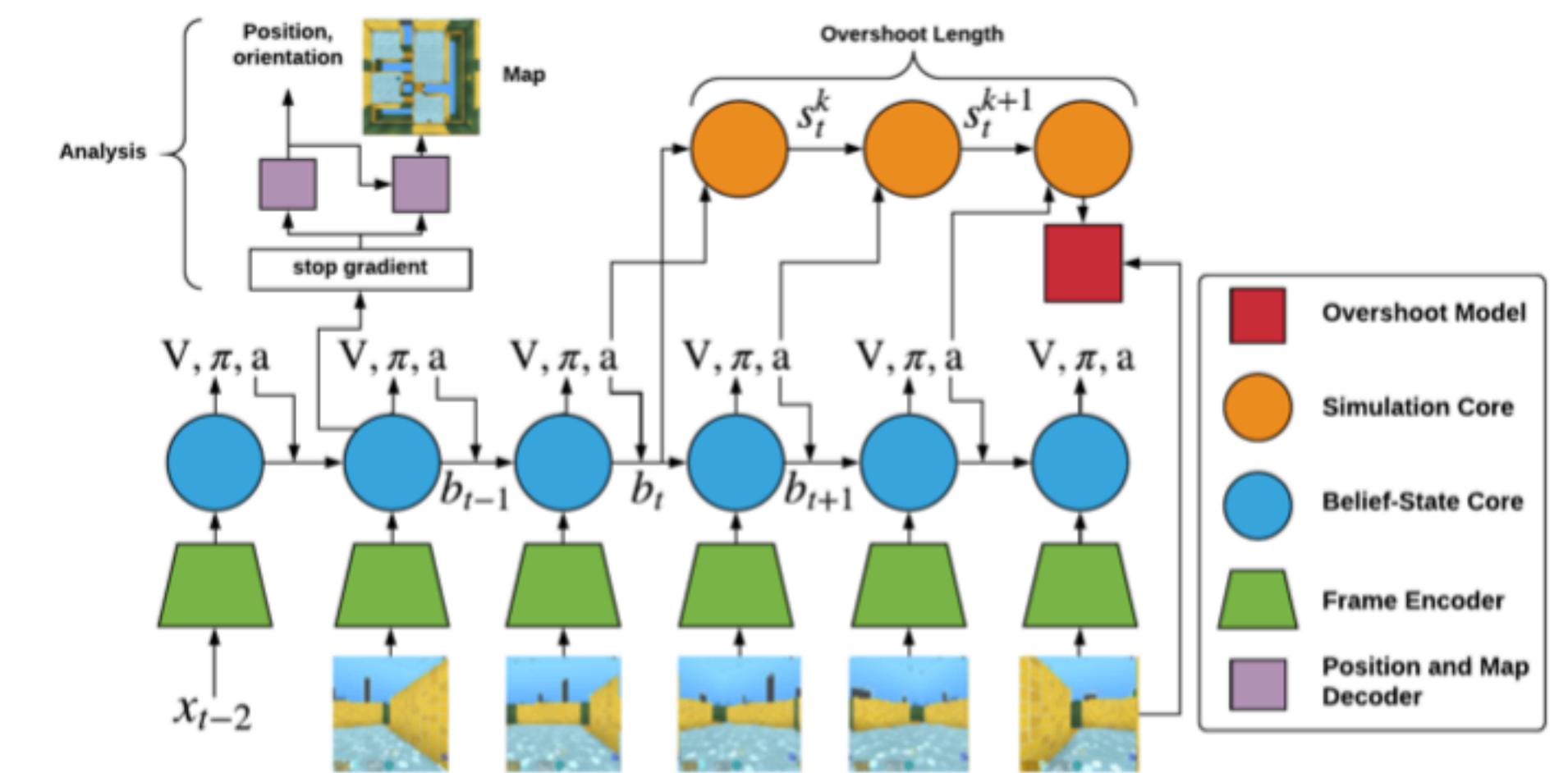
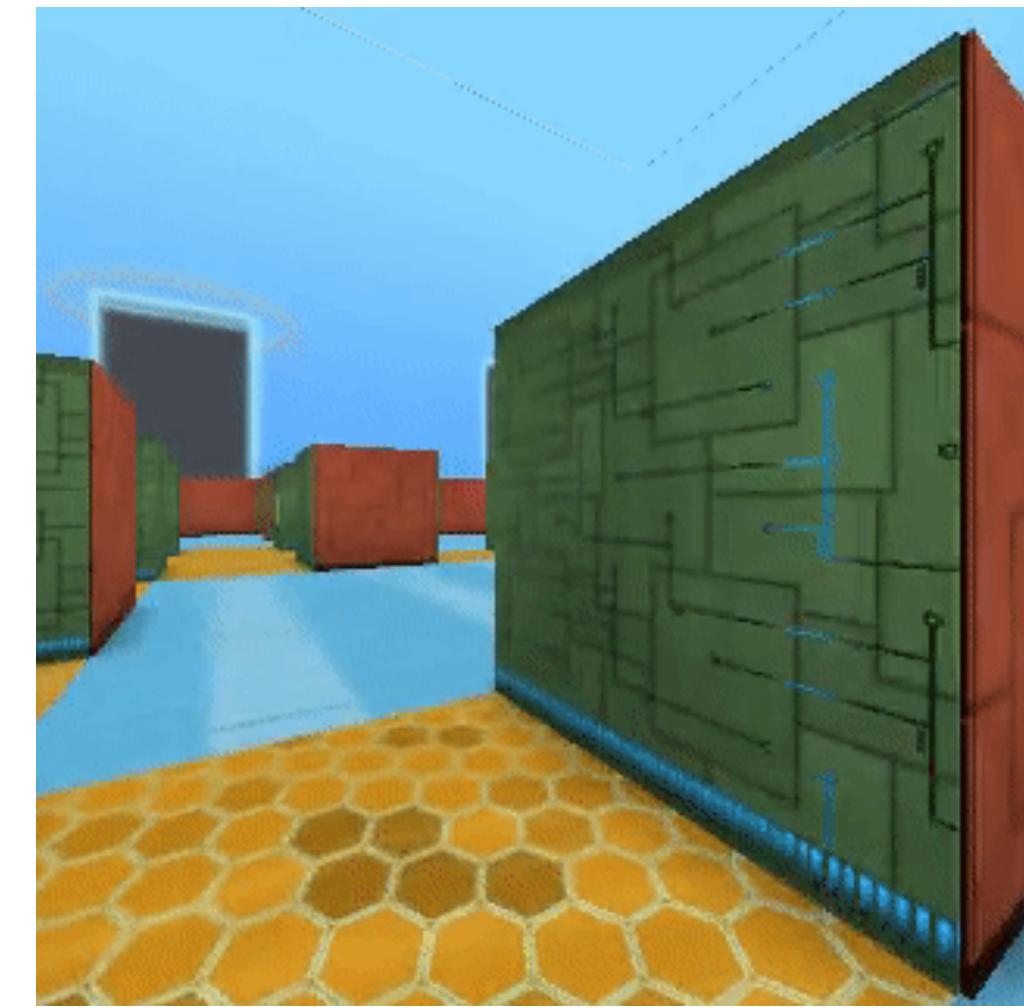


Karras et al. 2017

# generative models

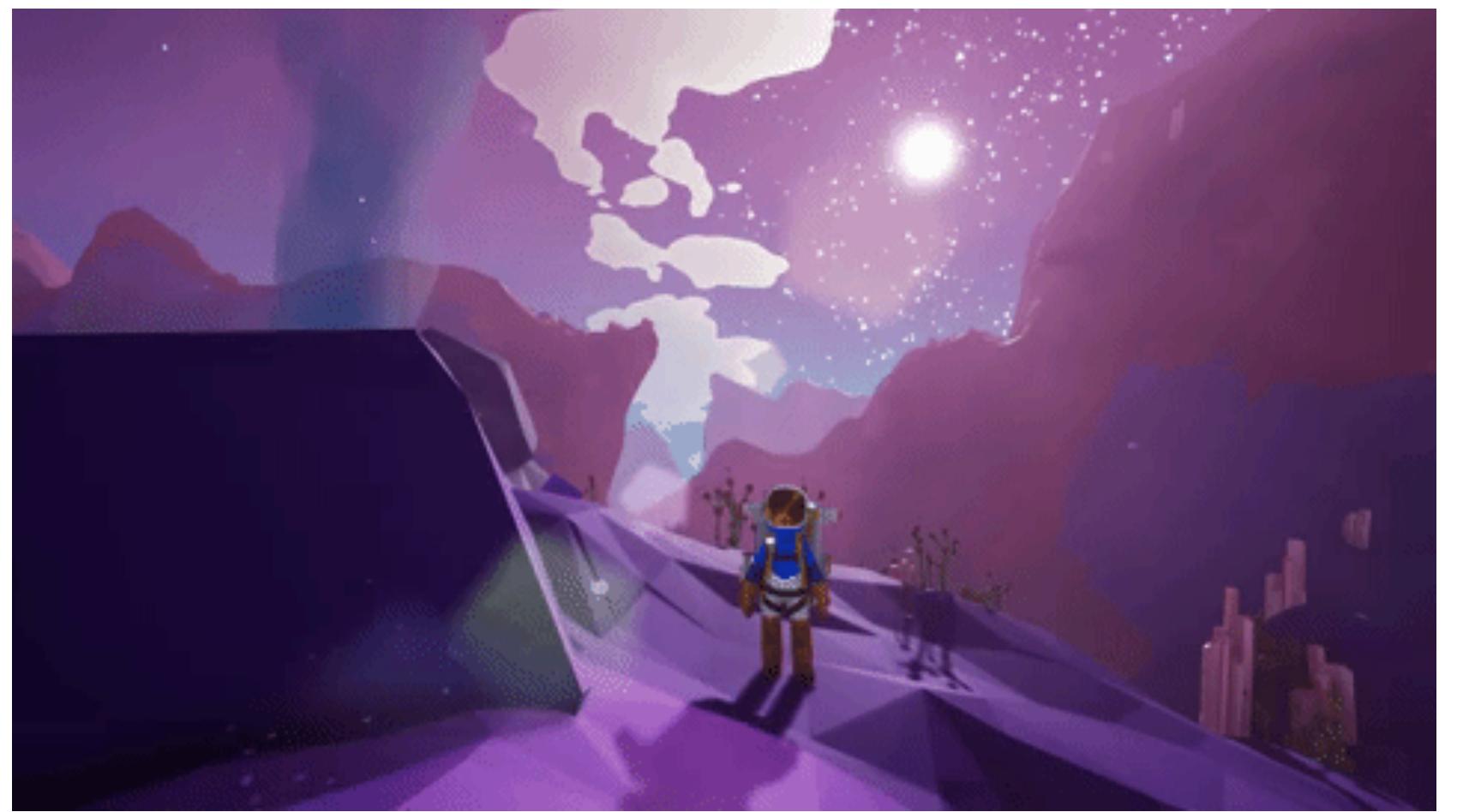
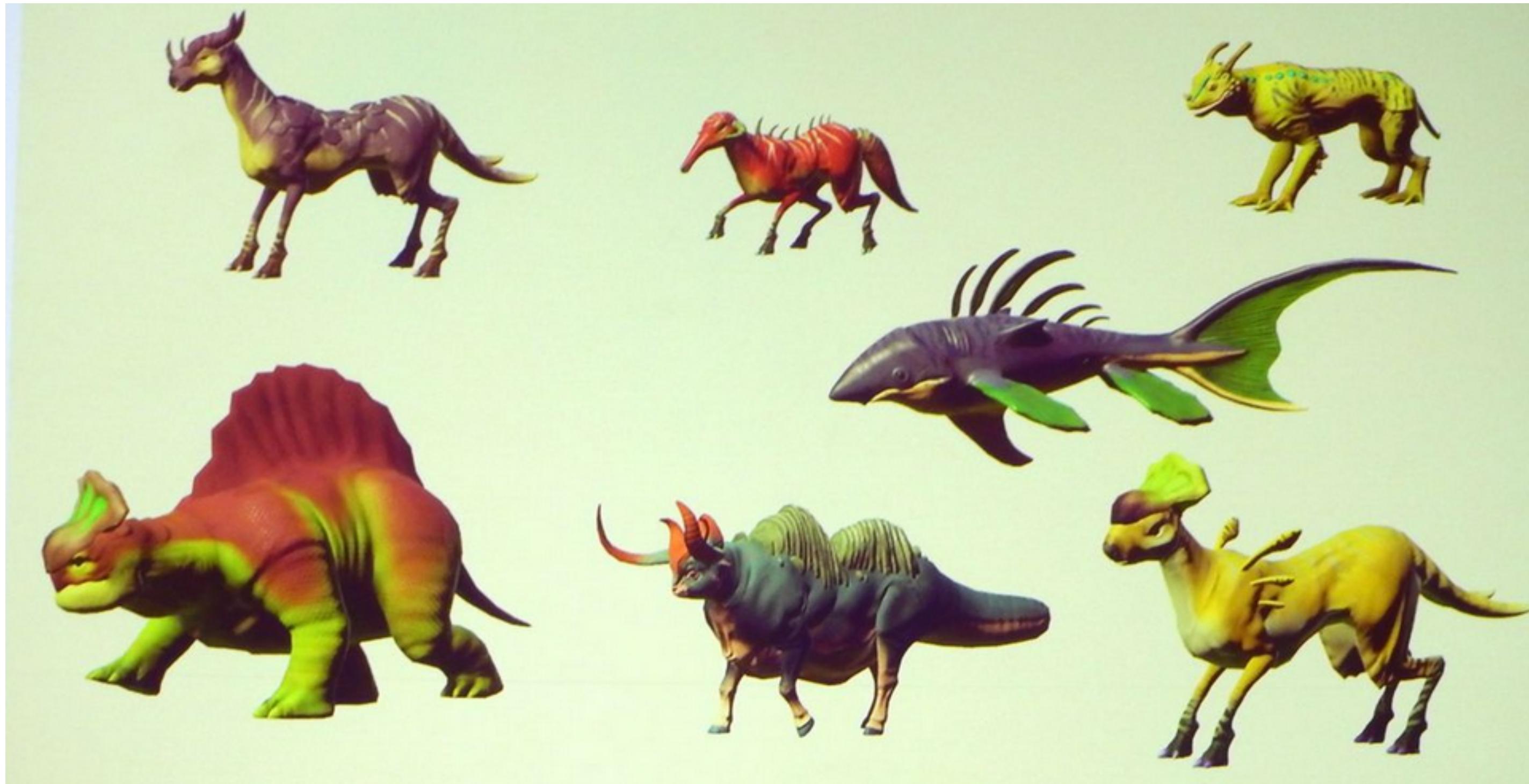


Silver et al 2016

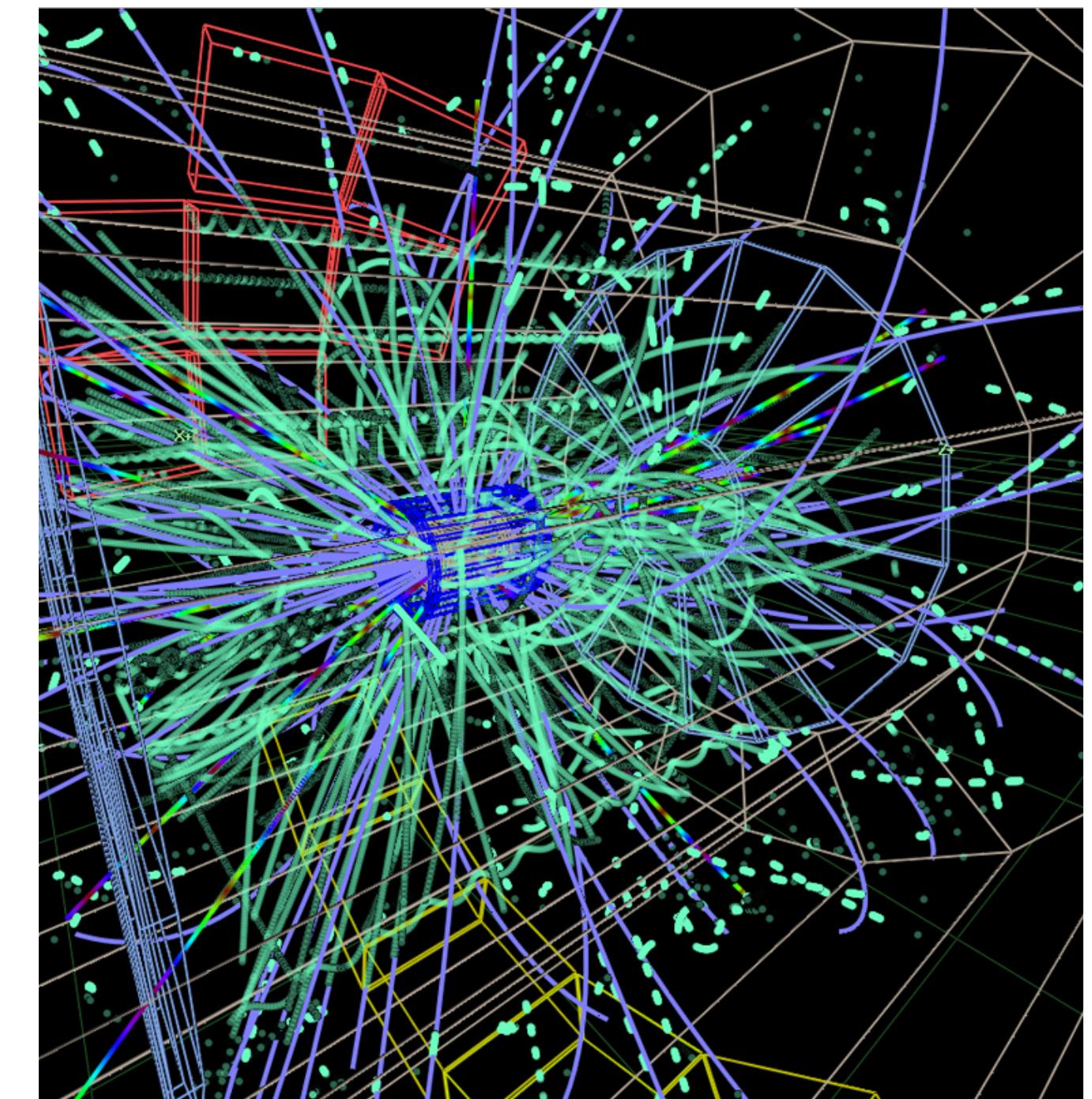
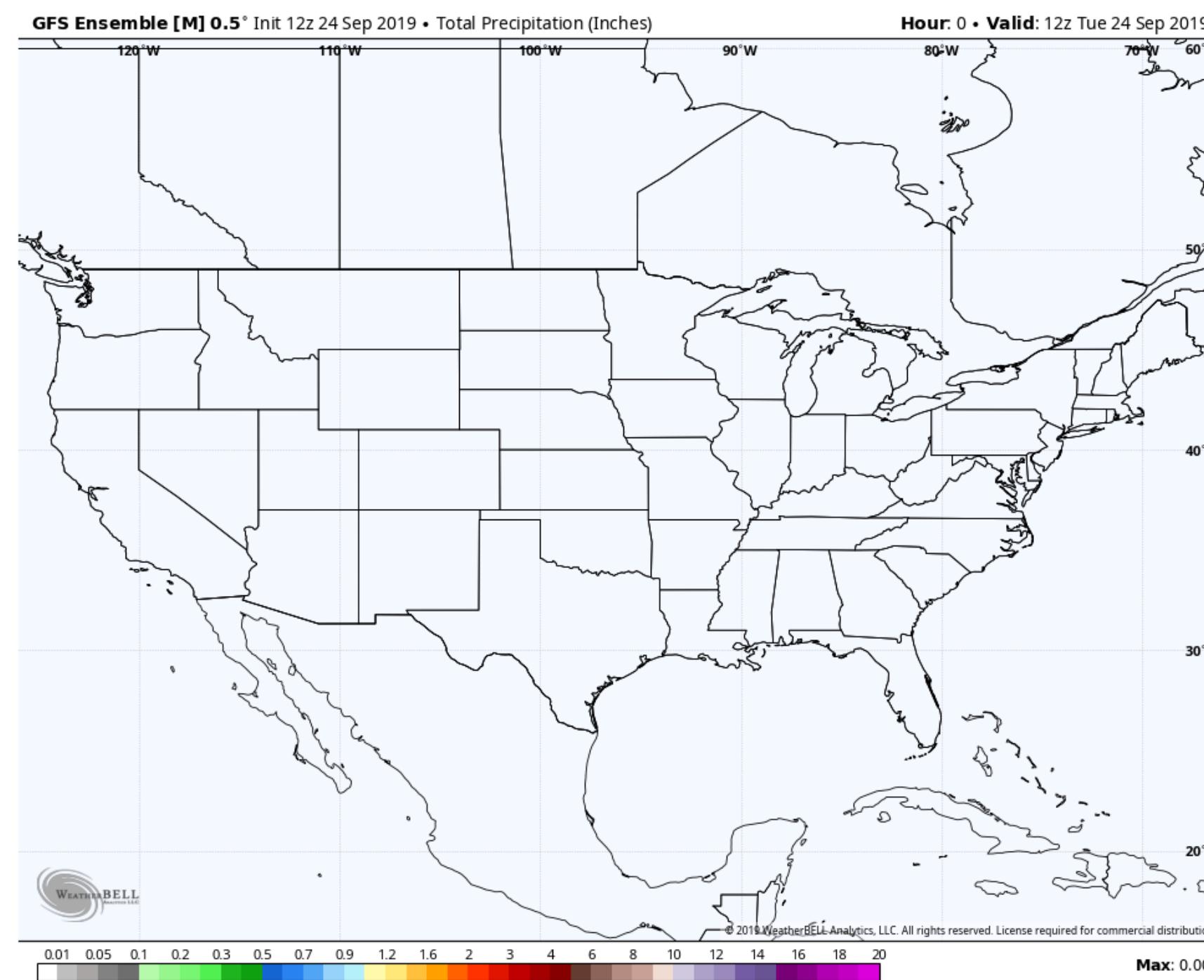
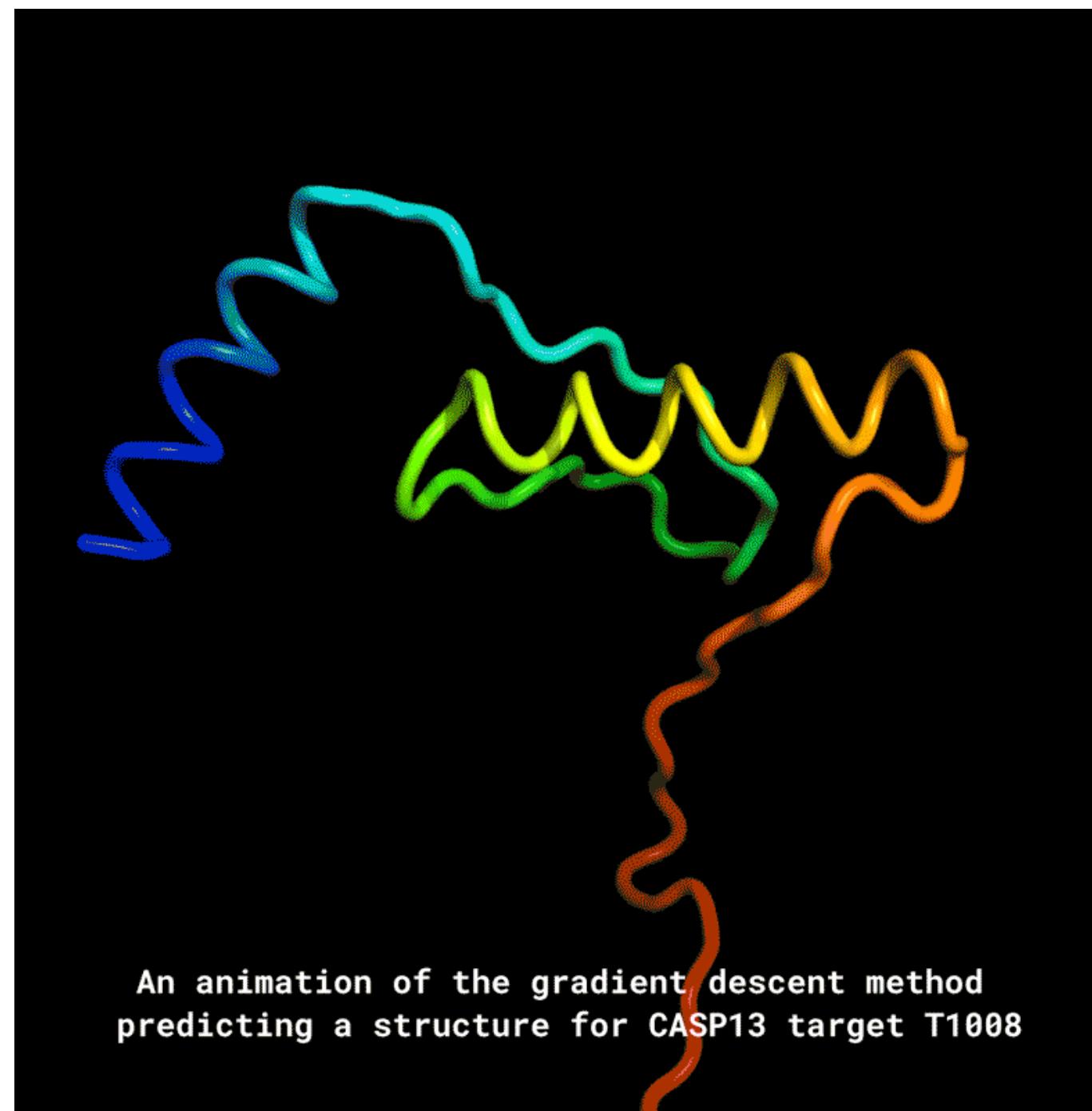


Besse et al 2019

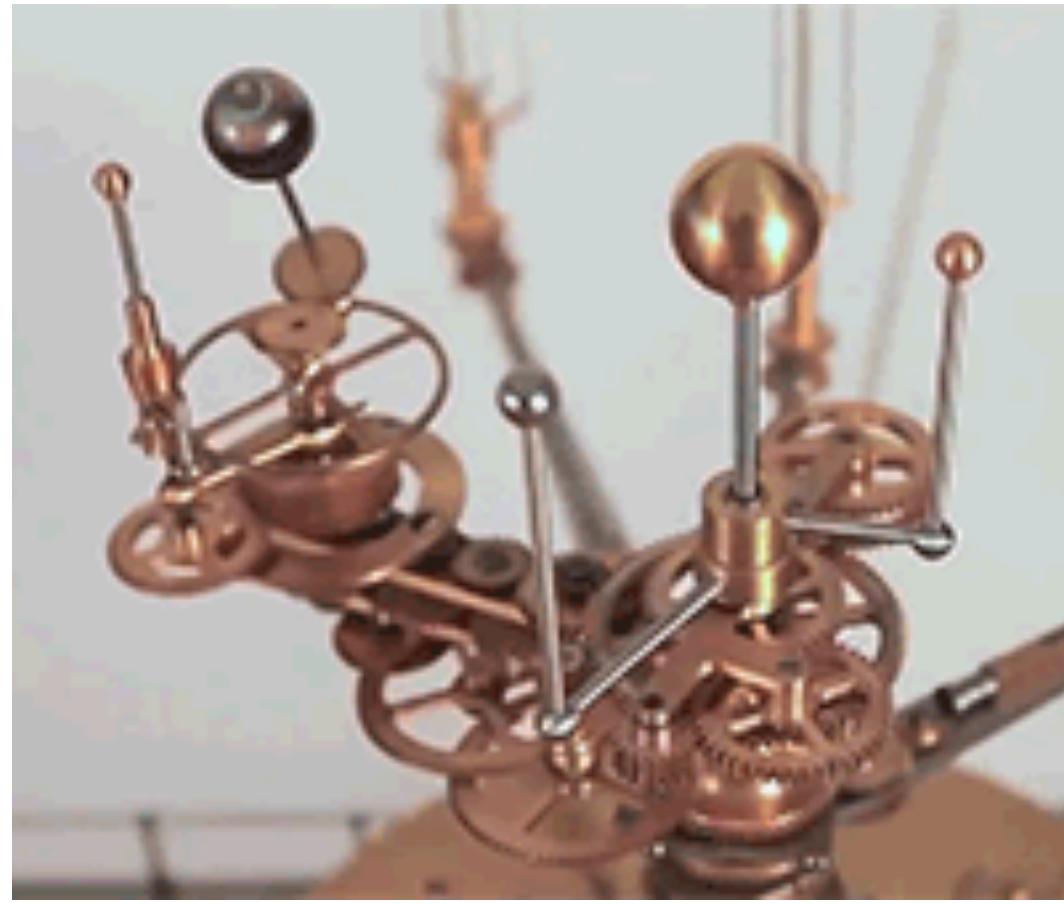
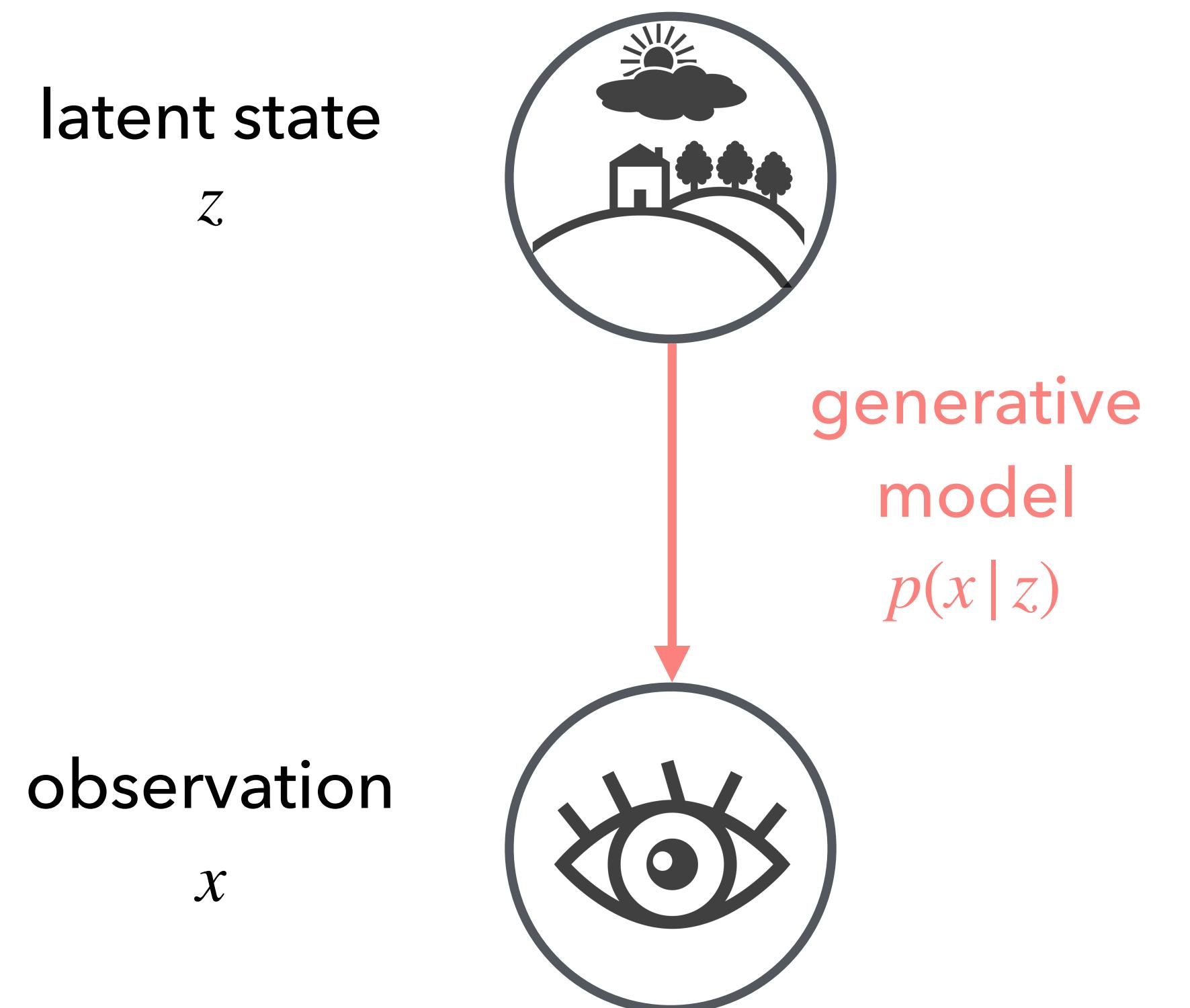
# generative models



# generative models



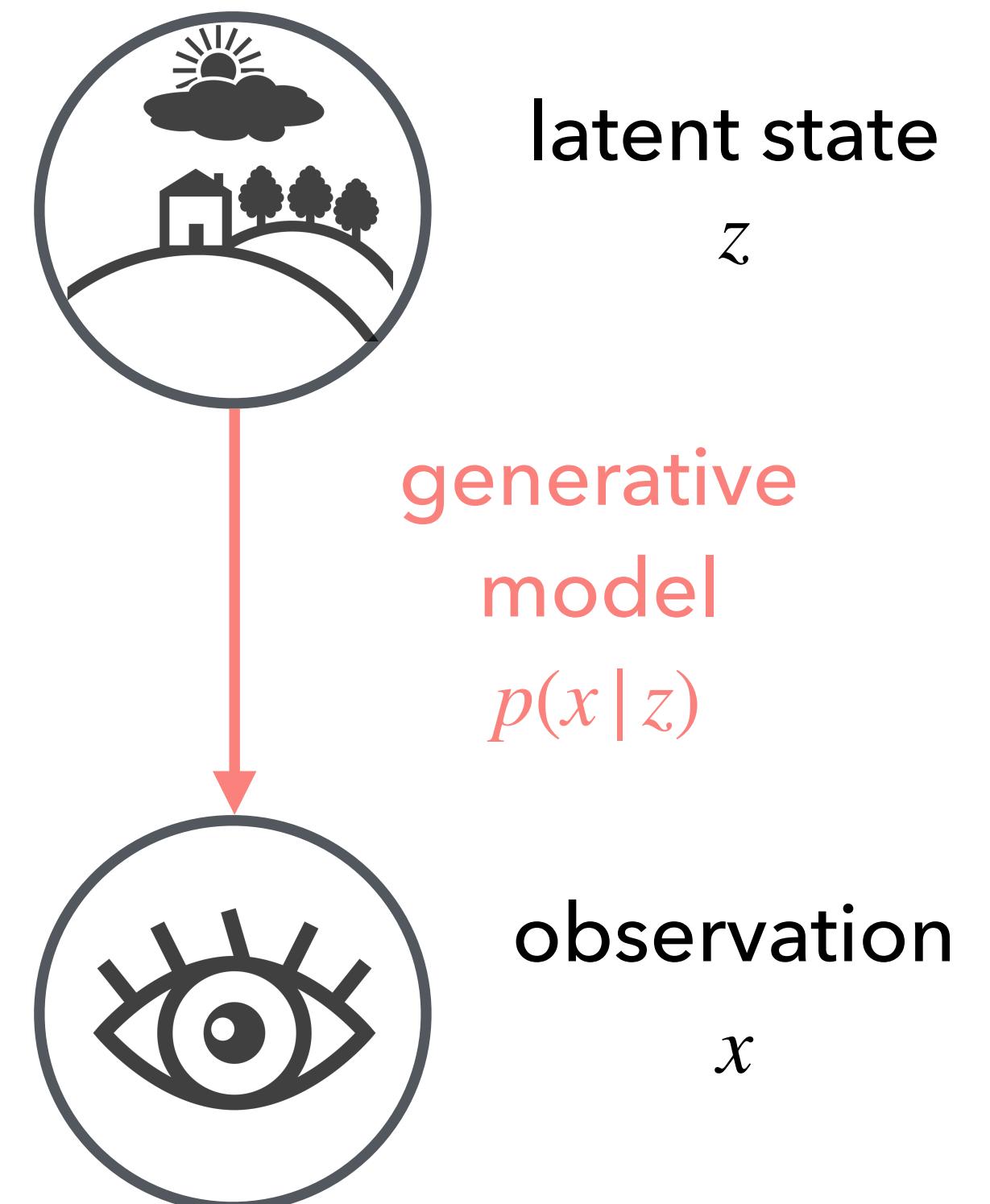
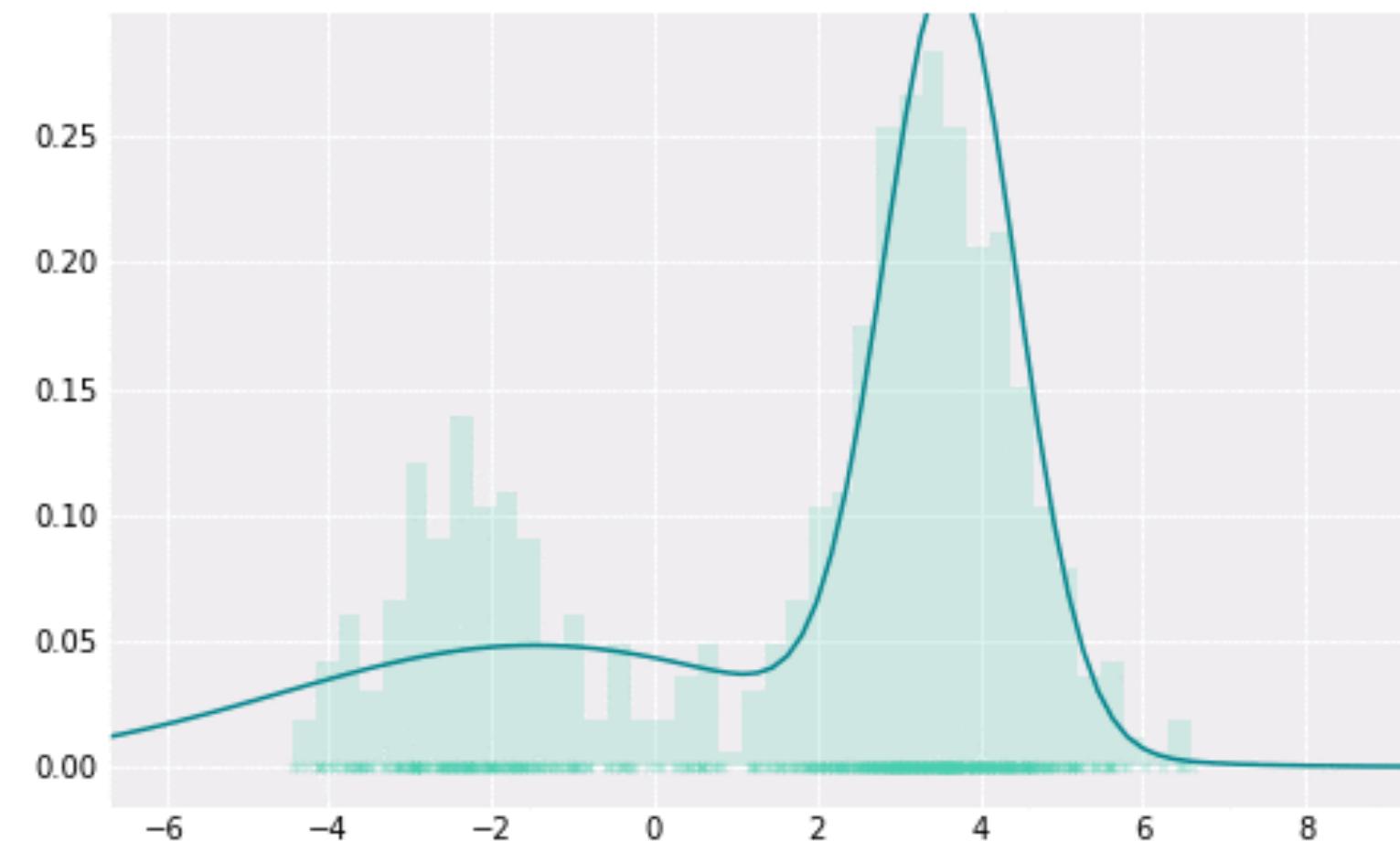
# generative models



# how do we learn a generative model?

- match empirical distribution to predictive distribution

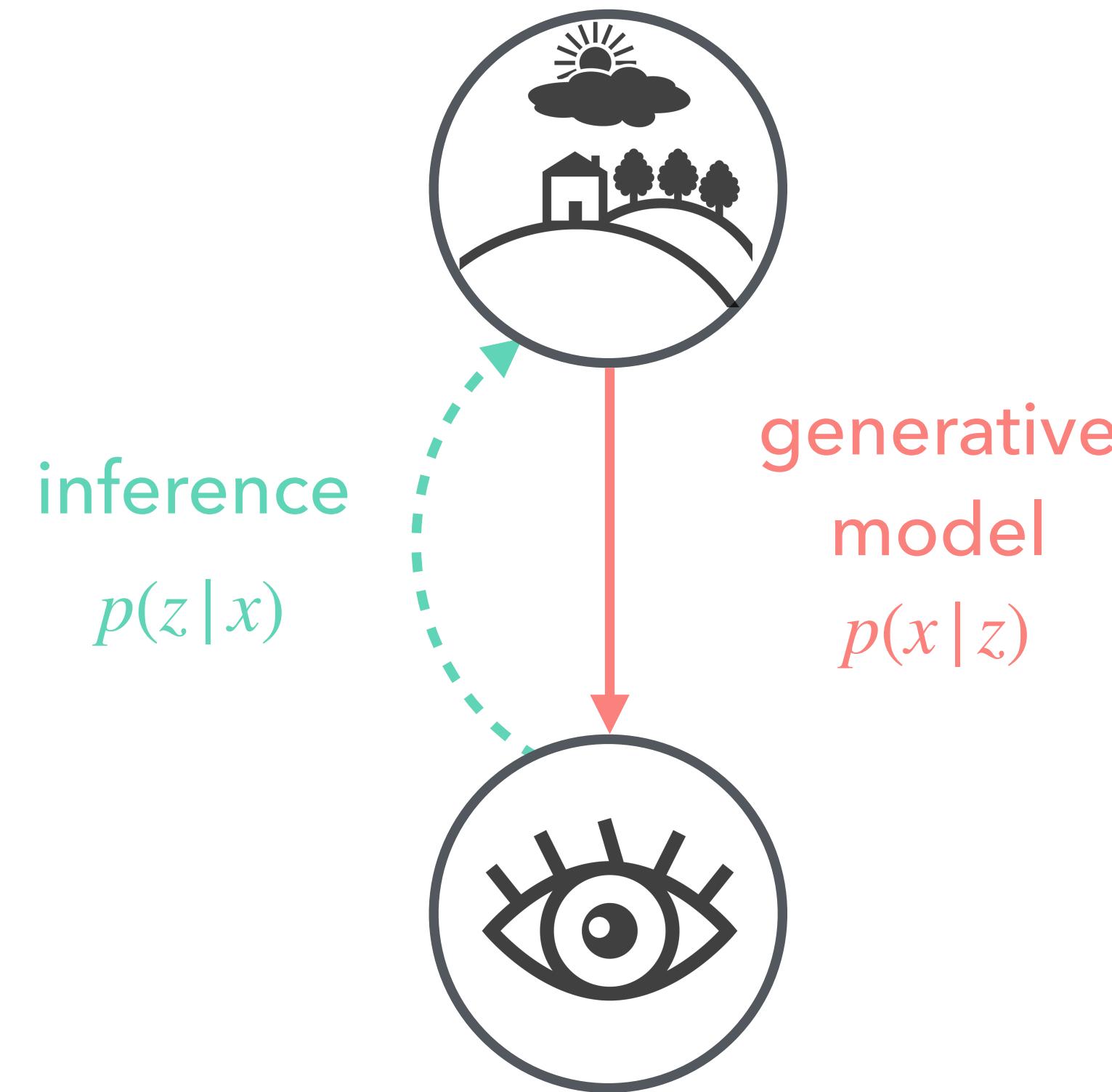
$$\arg \min_{model} KL [ p_{data}(x) || p_{model}(x) ]$$



- equivalent to maximum likelihood

$$\arg \max_{\theta} \log p_{\theta}(x)$$

# inference



ruise

13

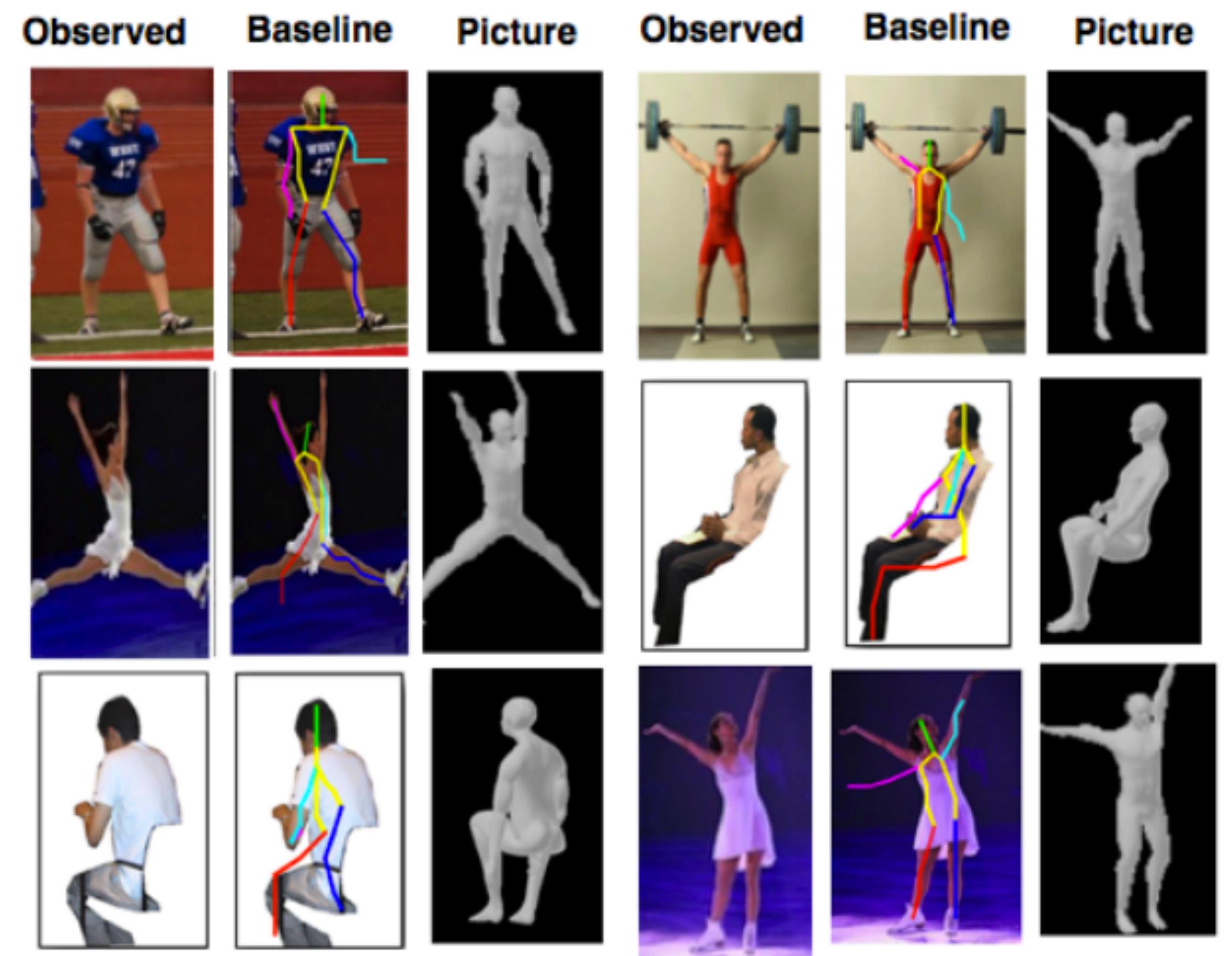
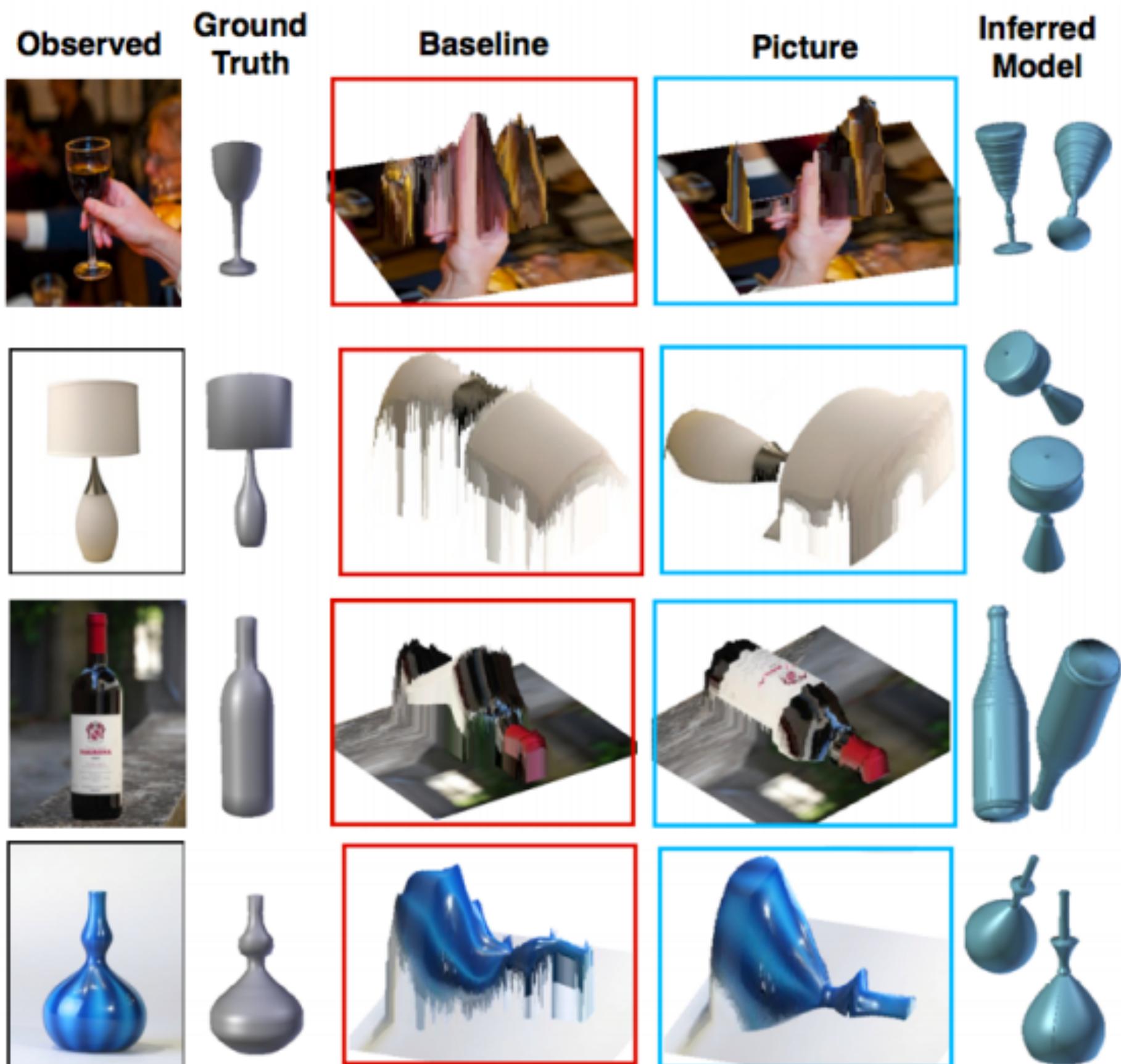
# inference

SPEED  
LIMIT  
**30**

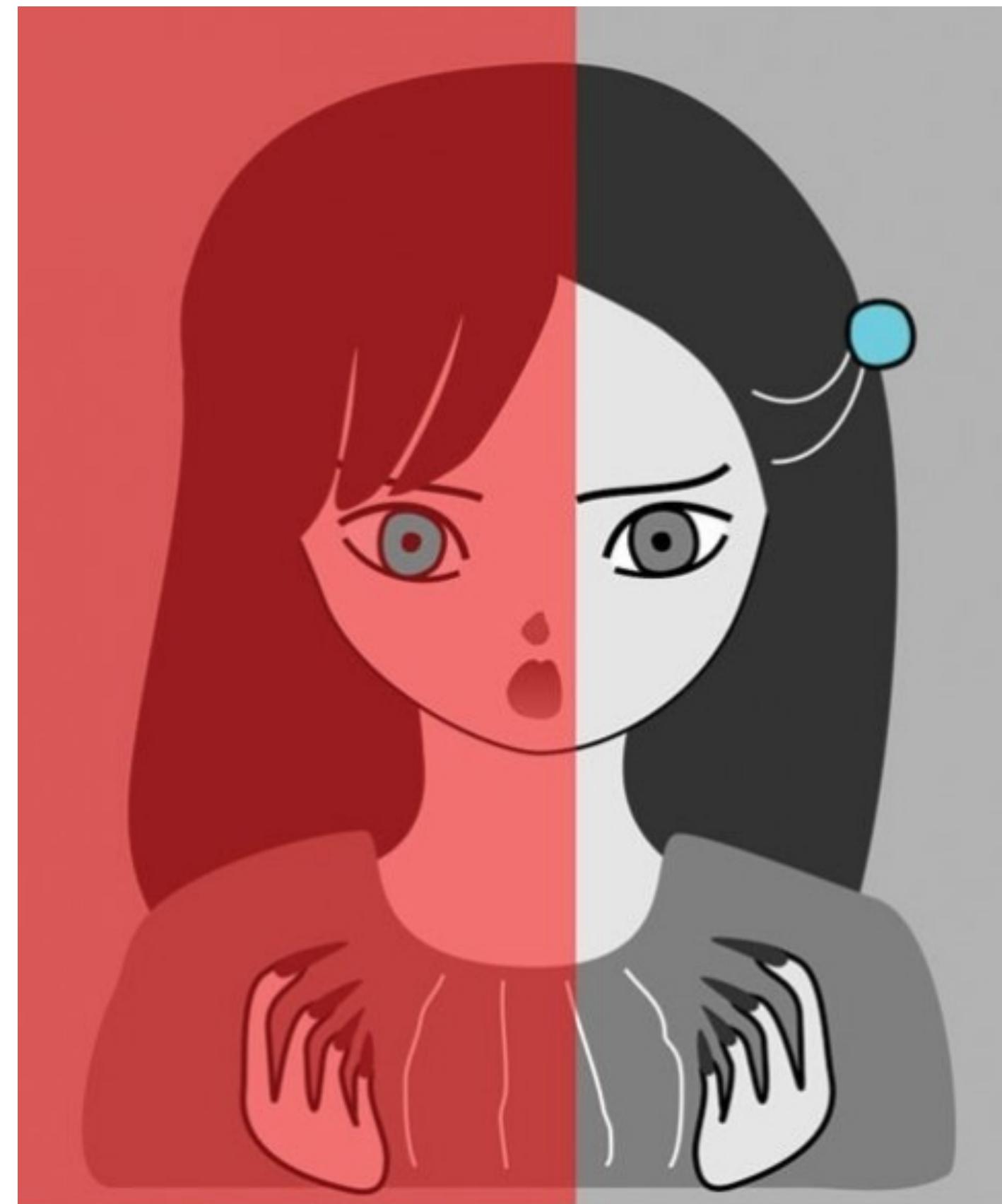


Google

# inference



# inference

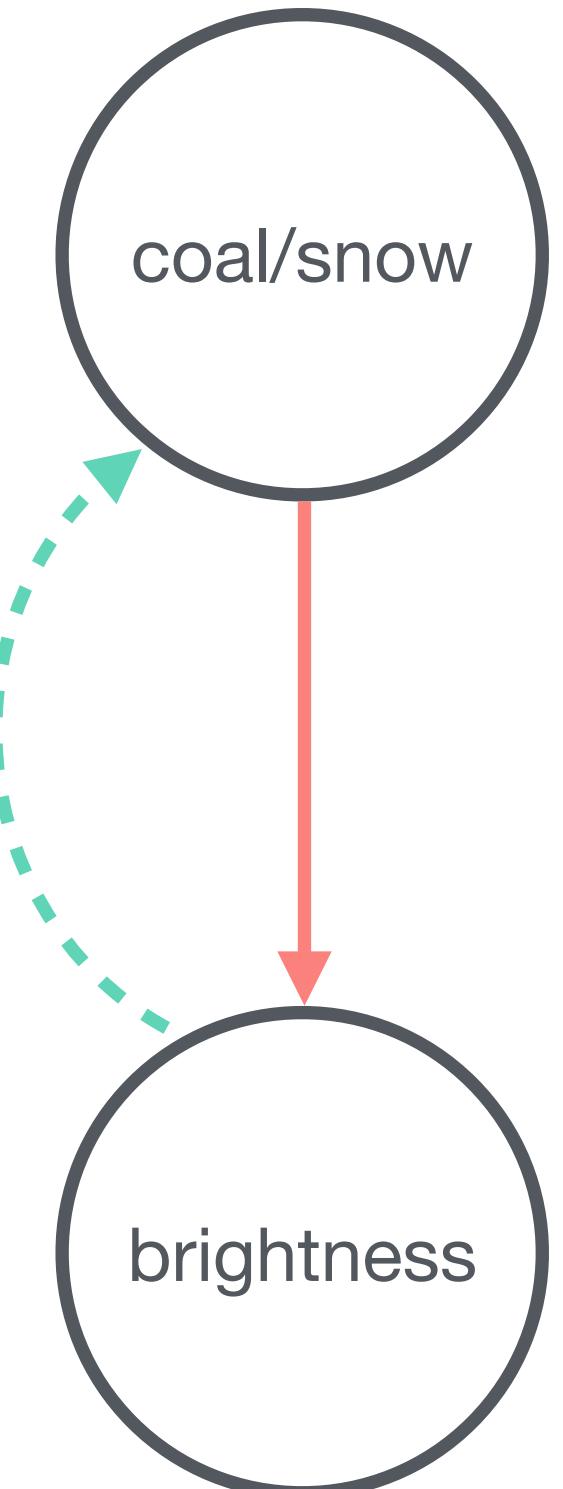


# inference



- what we see is not the data but an interpretation of the data
- ‘unconscious inference’ over latent variables

# how do we do inference?

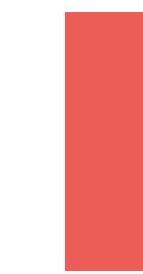




# how do we do inference?

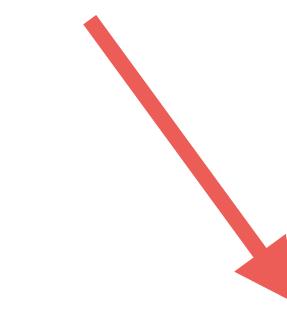


# how do we do inference?

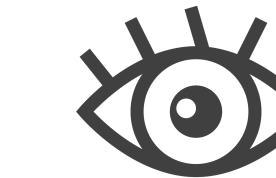


$\text{lighting} \times \text{reflectance} = \text{brightness}$

lighting



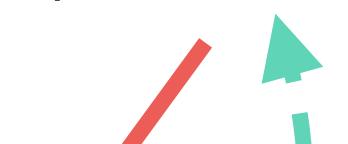
brightness



reflectance  
(material)



*inference*

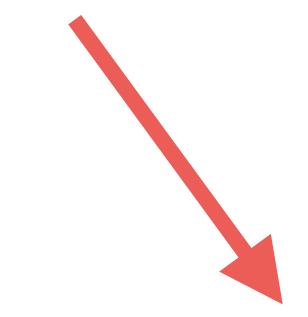


# how do we do inference?

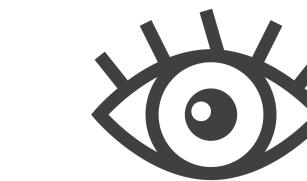


$\text{lighting} \times \text{reflectance} = \text{brightness}$

lighting



brightness



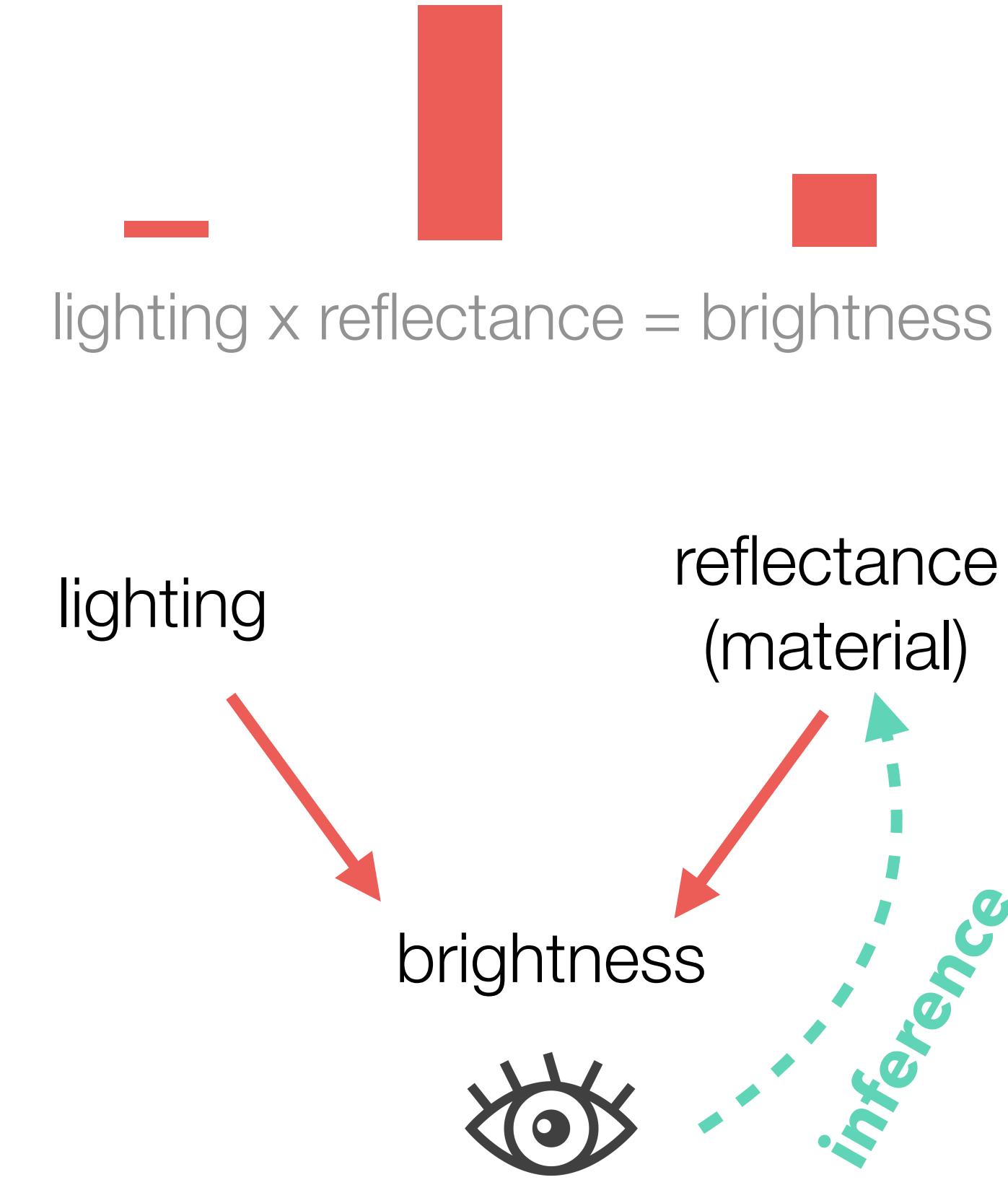
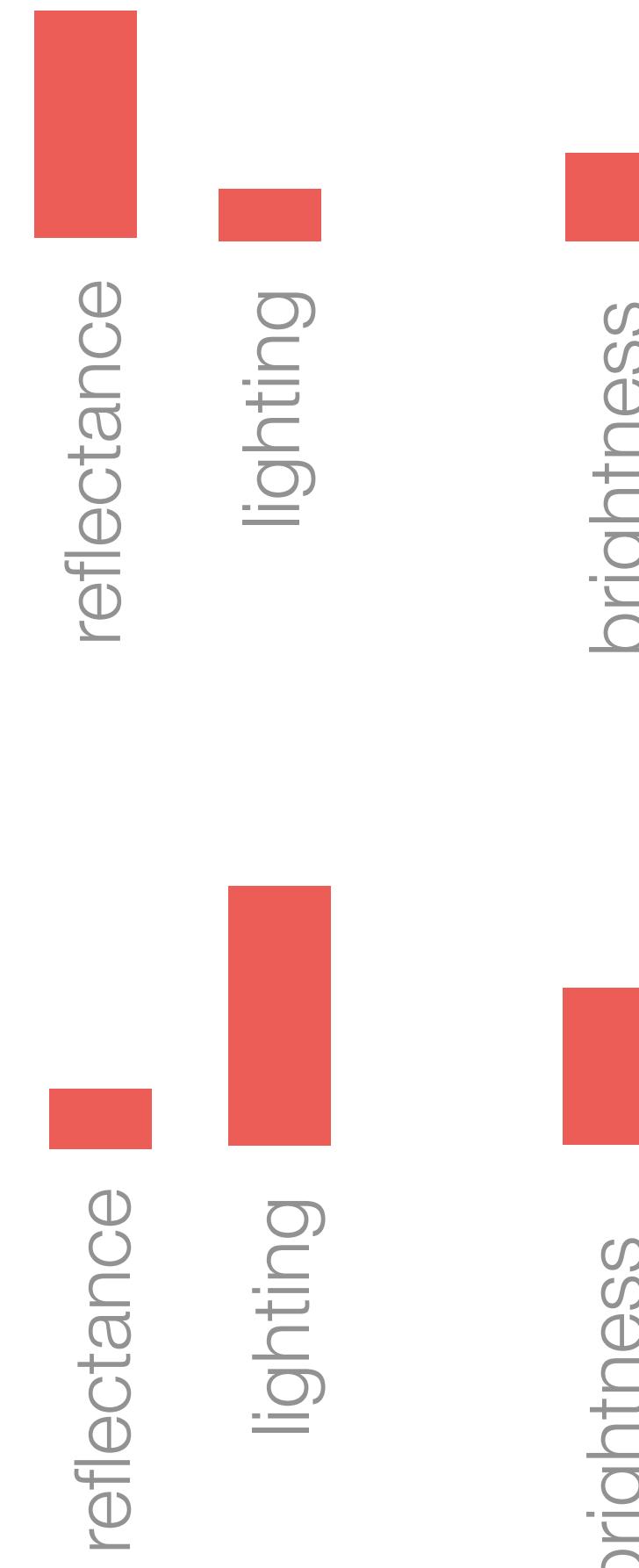
reflectance  
(material)



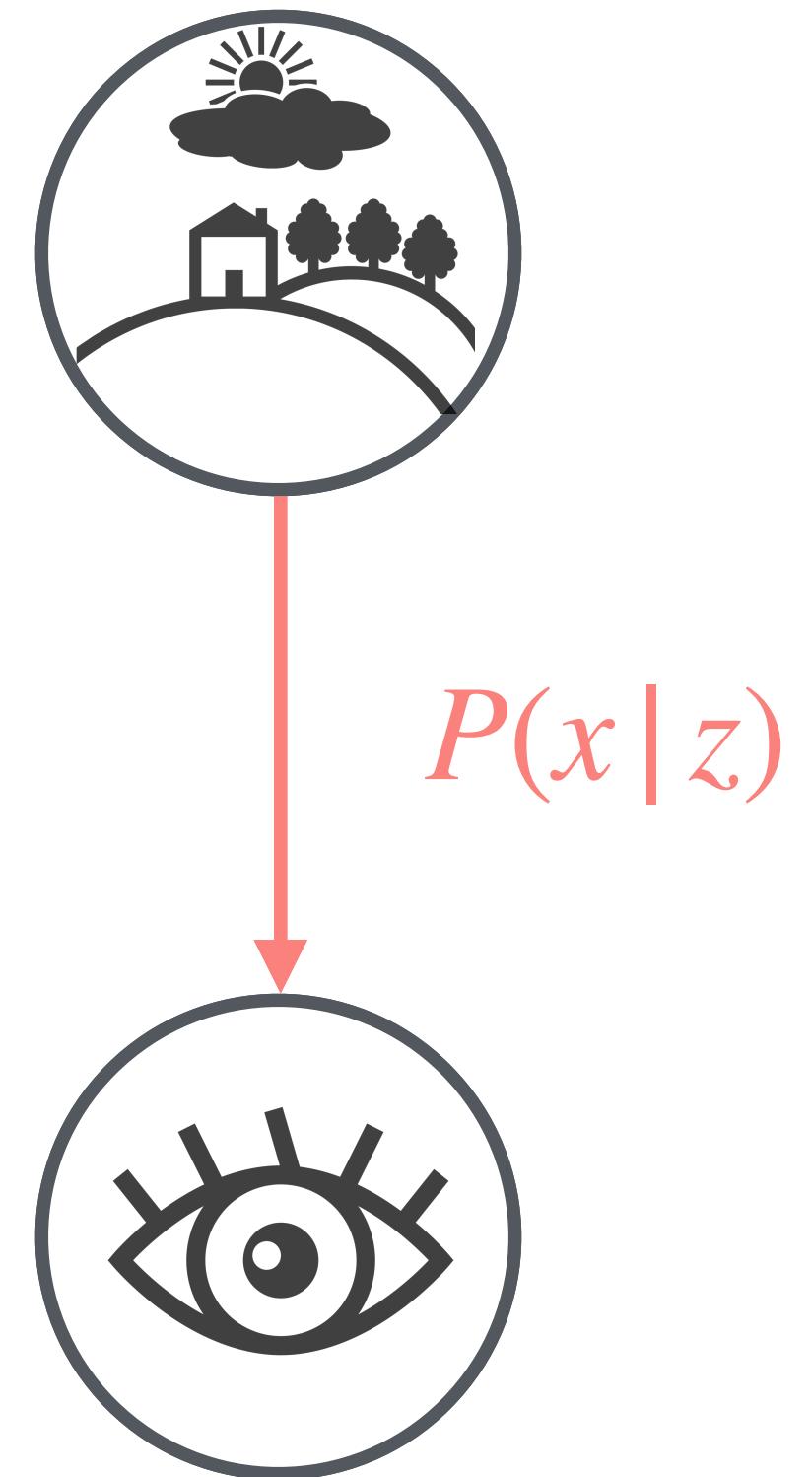
*inference*



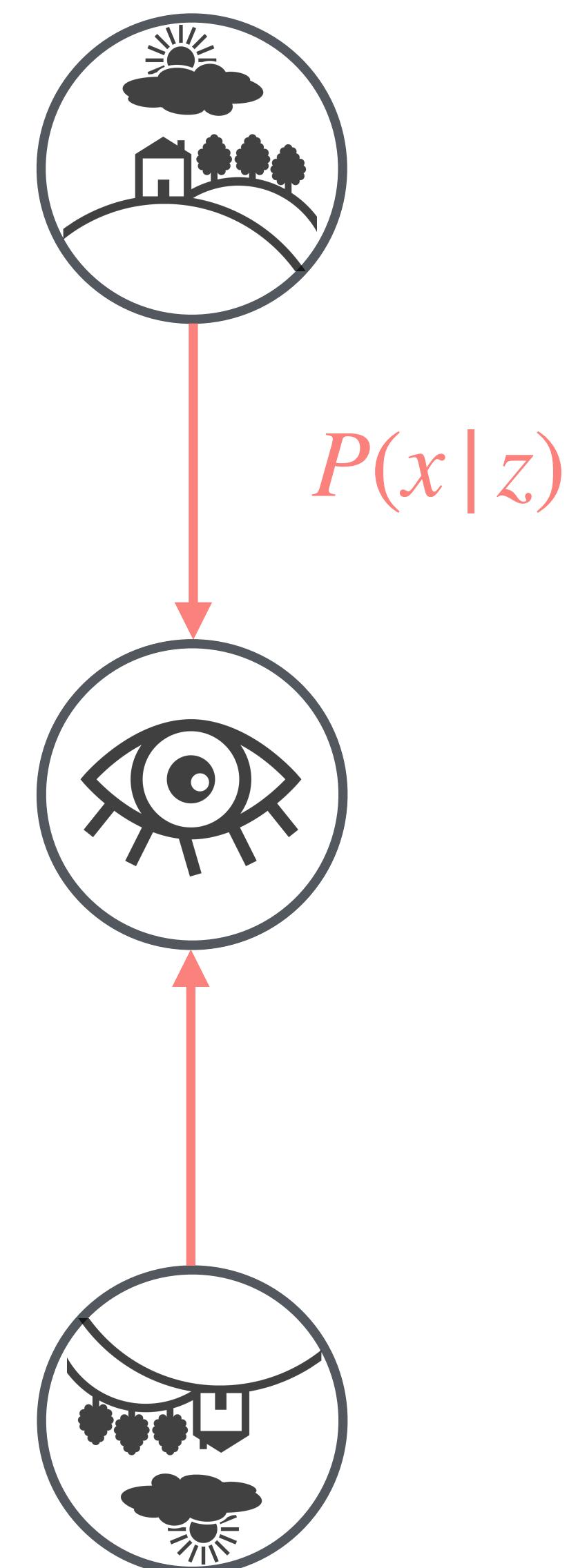
# how do we do inference?



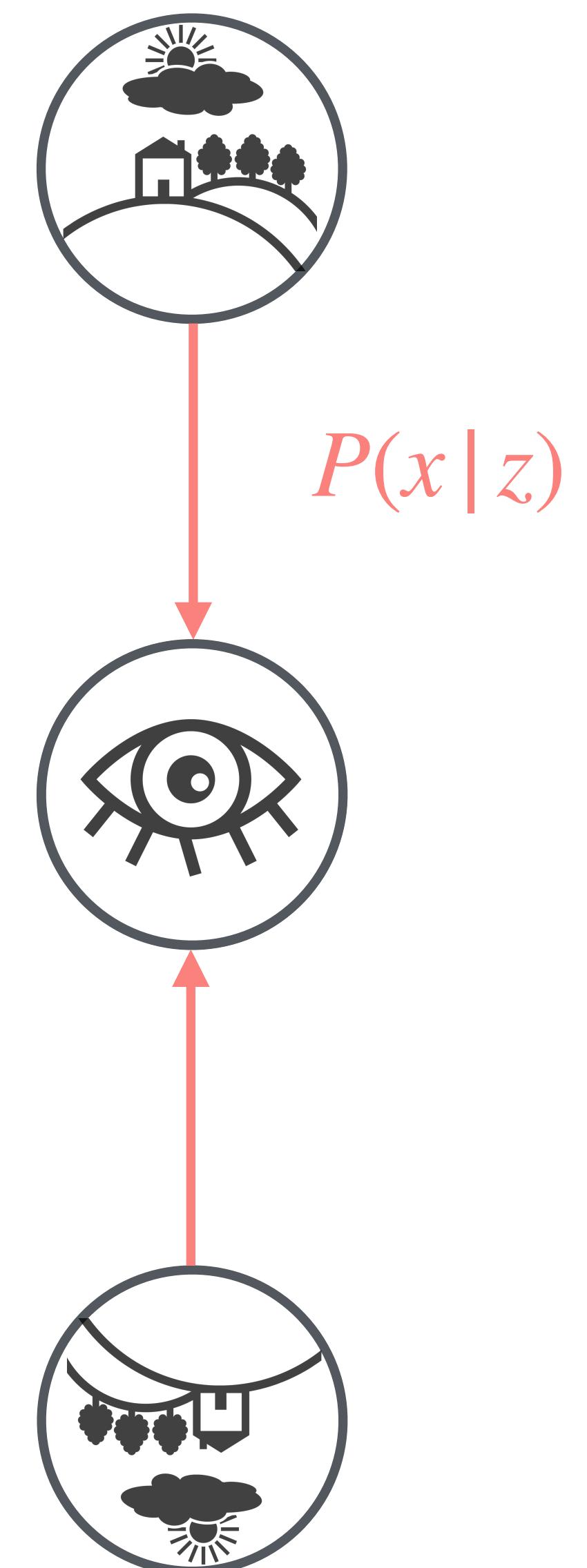
# how do we do inference?



# how do we do inference?

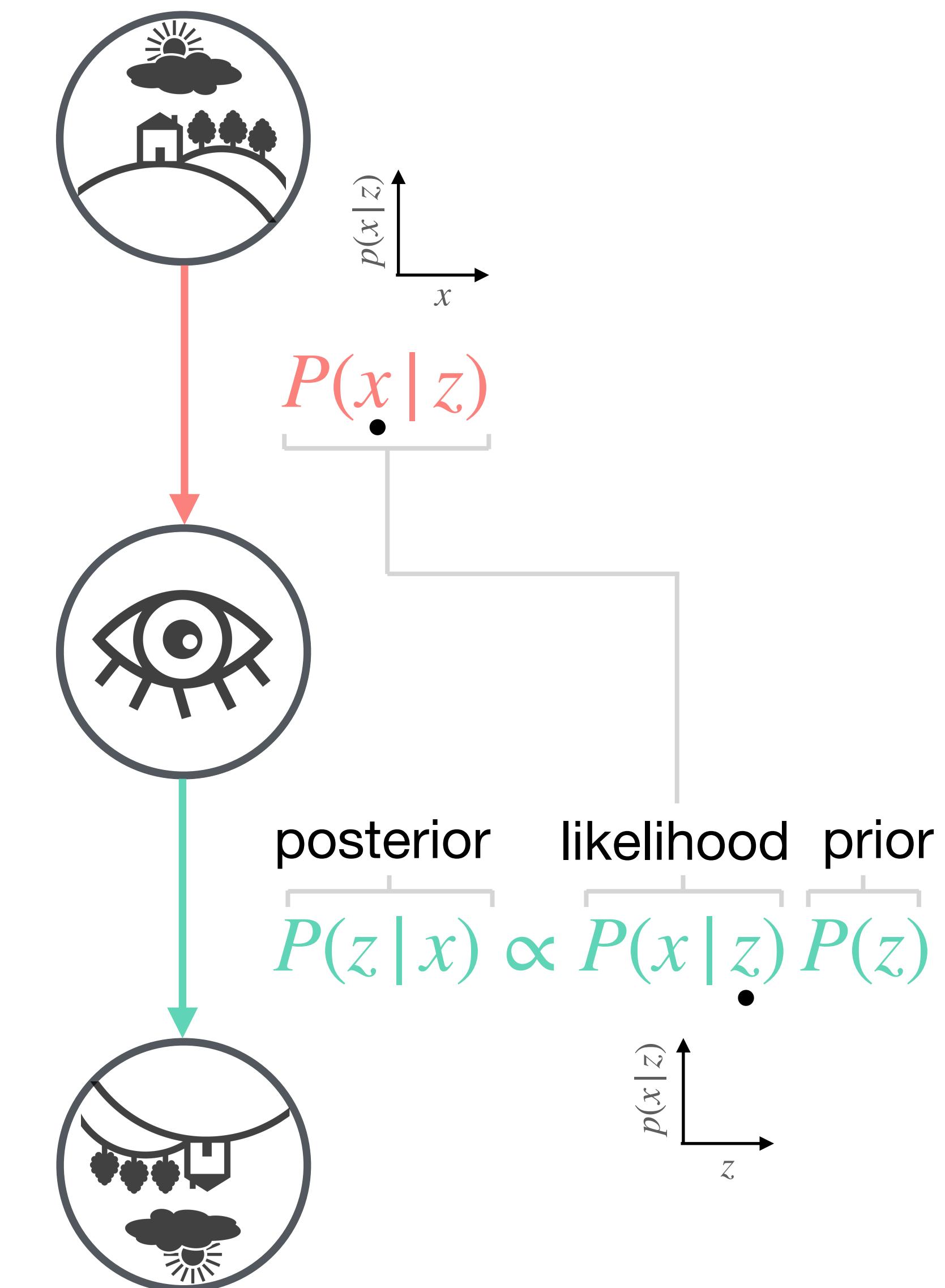


# how do we do inference?



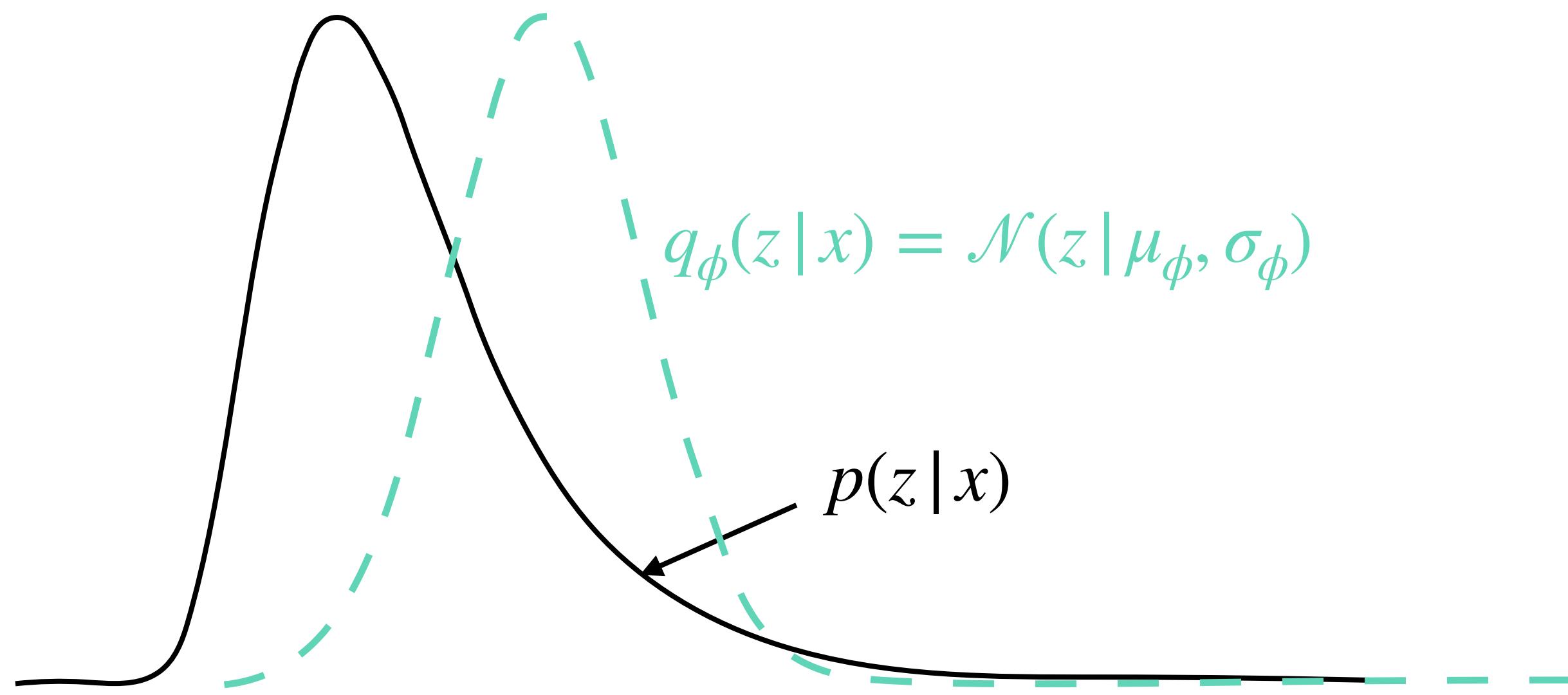
# how do we do inference?

- inference can be performed by inverting the generative model
- “vision is inverse graphics”
- Bayesian inference



# how do we do inference?

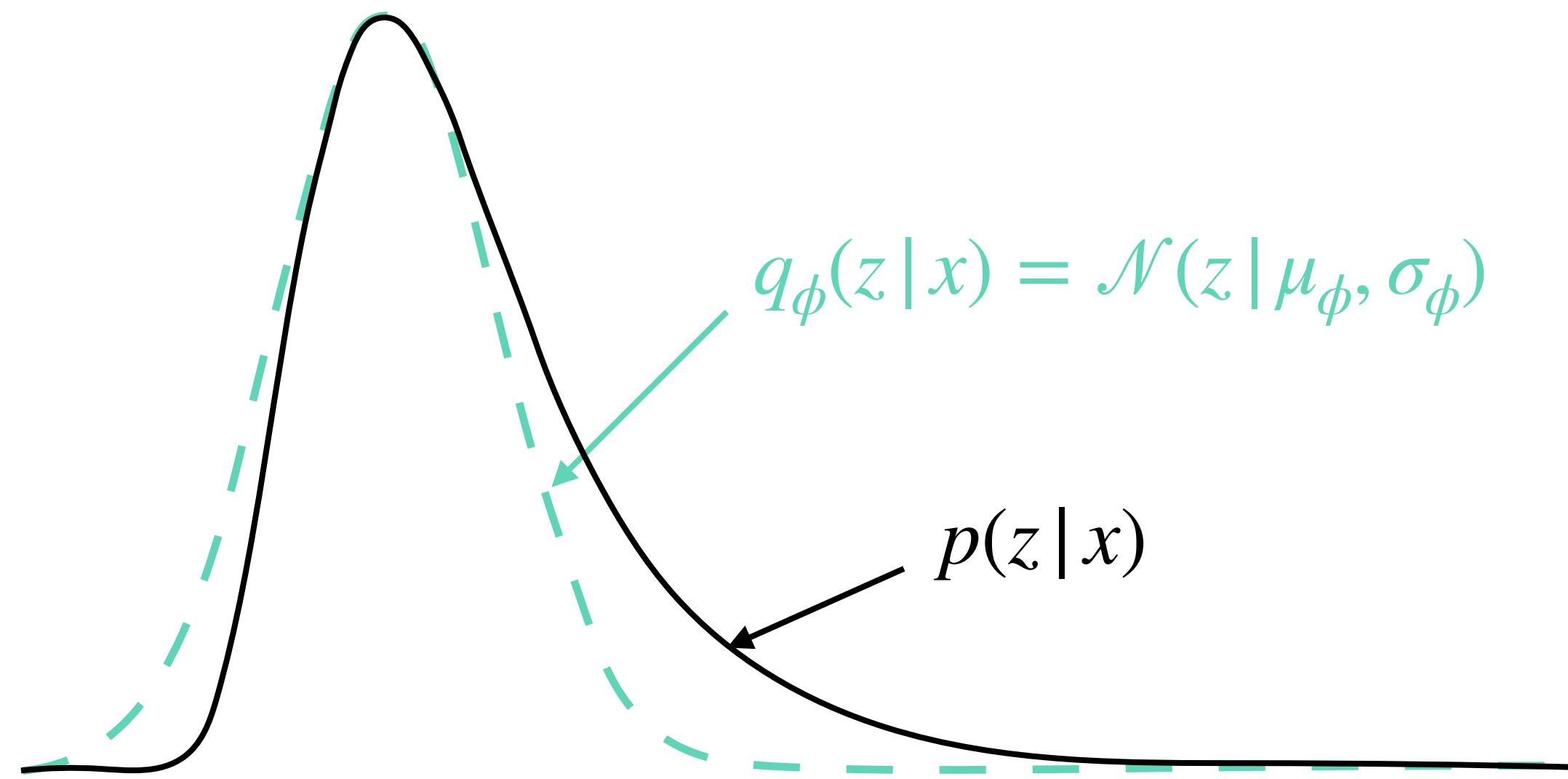
## approximation 1: variational bayes



$$\arg \min_{\phi} KL[q_\phi(z|x) || p(z|x)]$$

# how do we do inference?

## approximation 1: variational bayes



$$\arg \min_{\phi} KL[q_\phi(z|x) || p(z|x)]$$

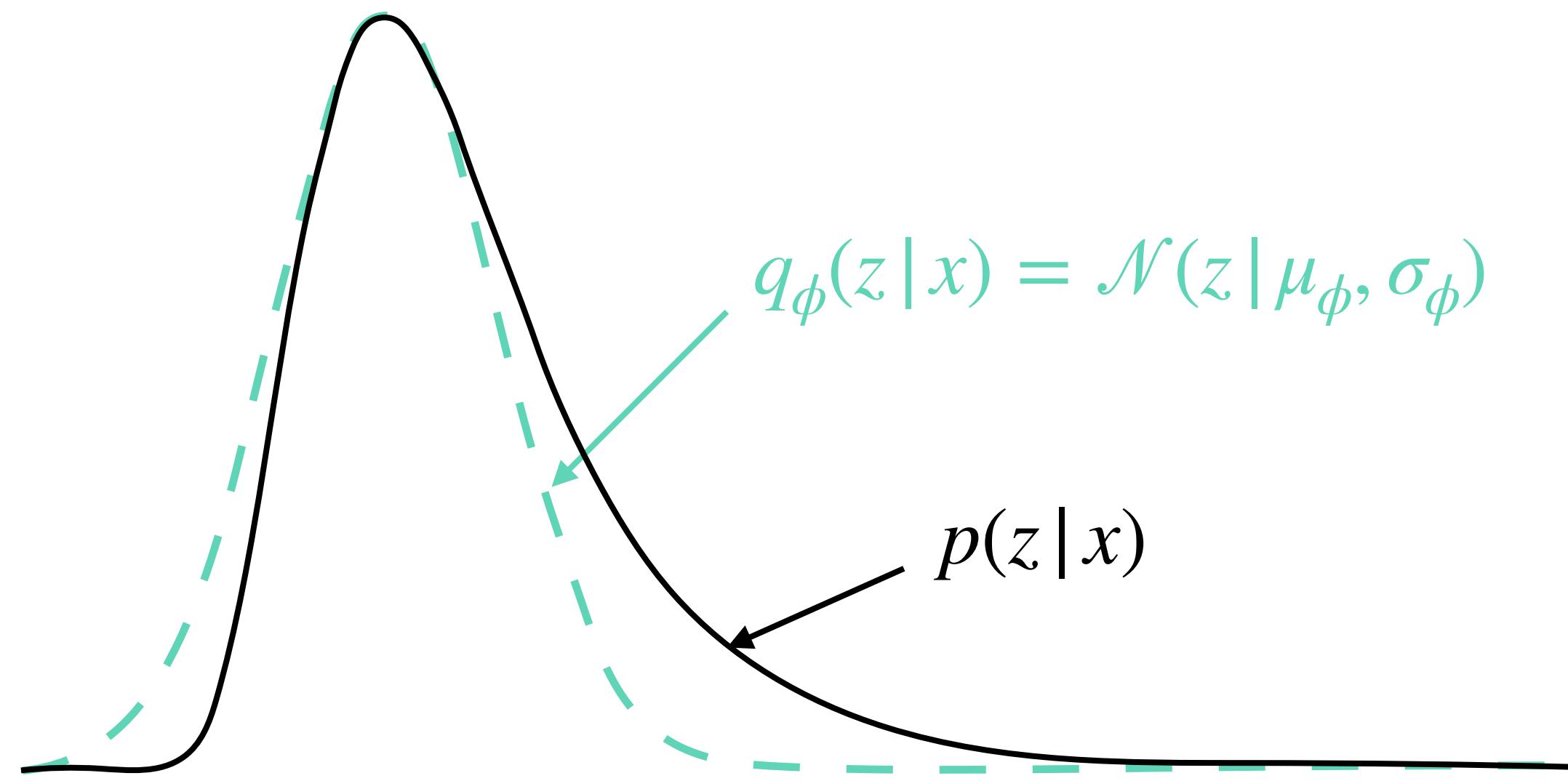
KL[ | ]

Ferenc notation ([inference.vc](#))

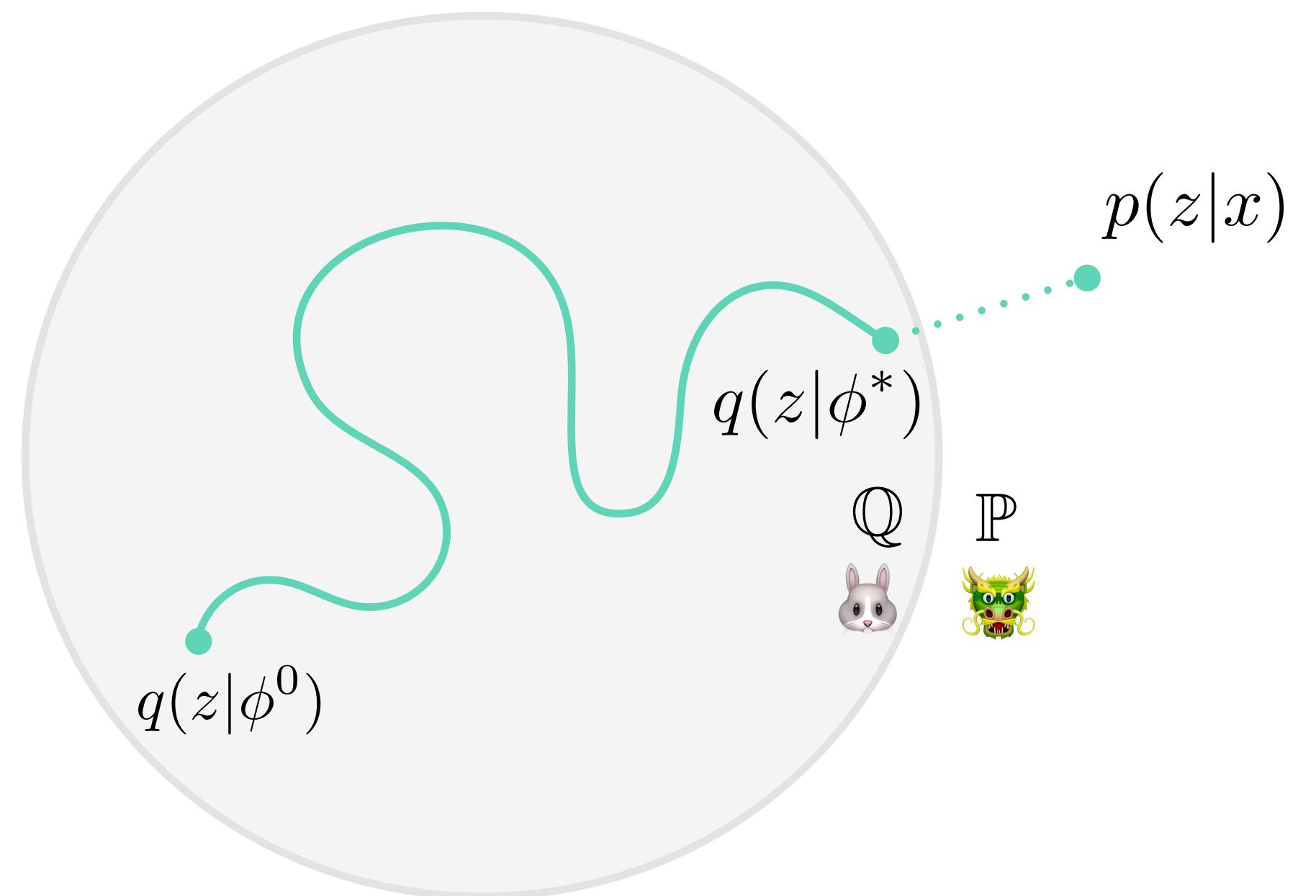
- $p_\theta(z)$  is very easy ,
- $p_\theta(x|z)$  is easy ,
- $p_\theta(x, z)$  is easy ,
- $p_\theta(x)$  is super-hard ,
- $p_\theta(z|x)$  is mega-hard

# how do we do inference?

## approximation 1: variational bayes



$$\arg \min_{\phi} KL[q_\phi(z|x) || p(z|x)]$$



# how do we do inference?

## approximation 2: amortised inference

$f$

$$x_1 \rightarrow q(z | x_1)$$

$$x_2 \rightarrow q(z | x_2)$$

:

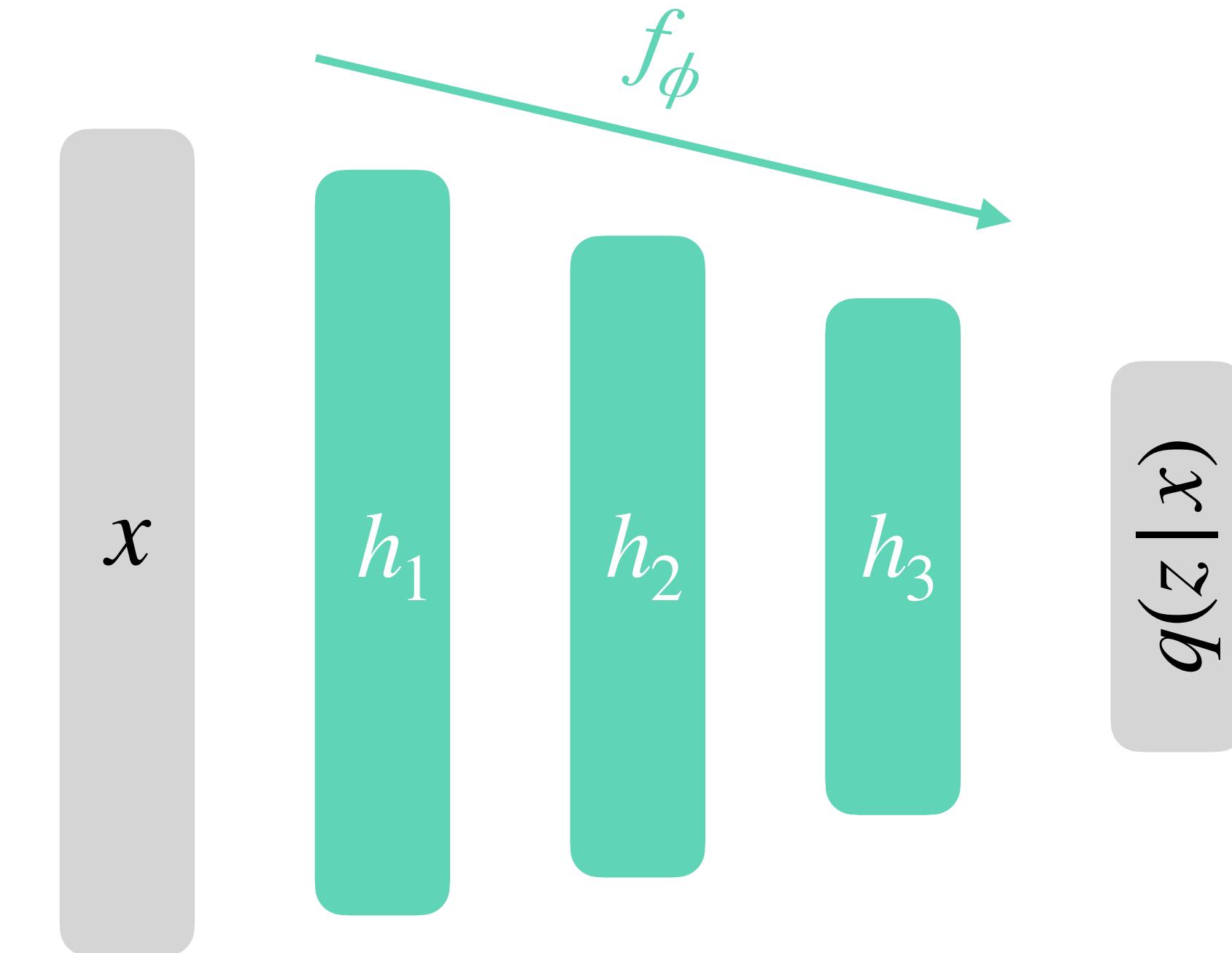
- $f$  mapping between datapoint and posterior
- if we have  $x, q$  pairs for a set of datapoints,
- we can treat learning  $f$  as a supervised regression problem

$$x_N \rightarrow q(z | x_N)$$

# how do we do inference?

## approximation 2: amortised inference

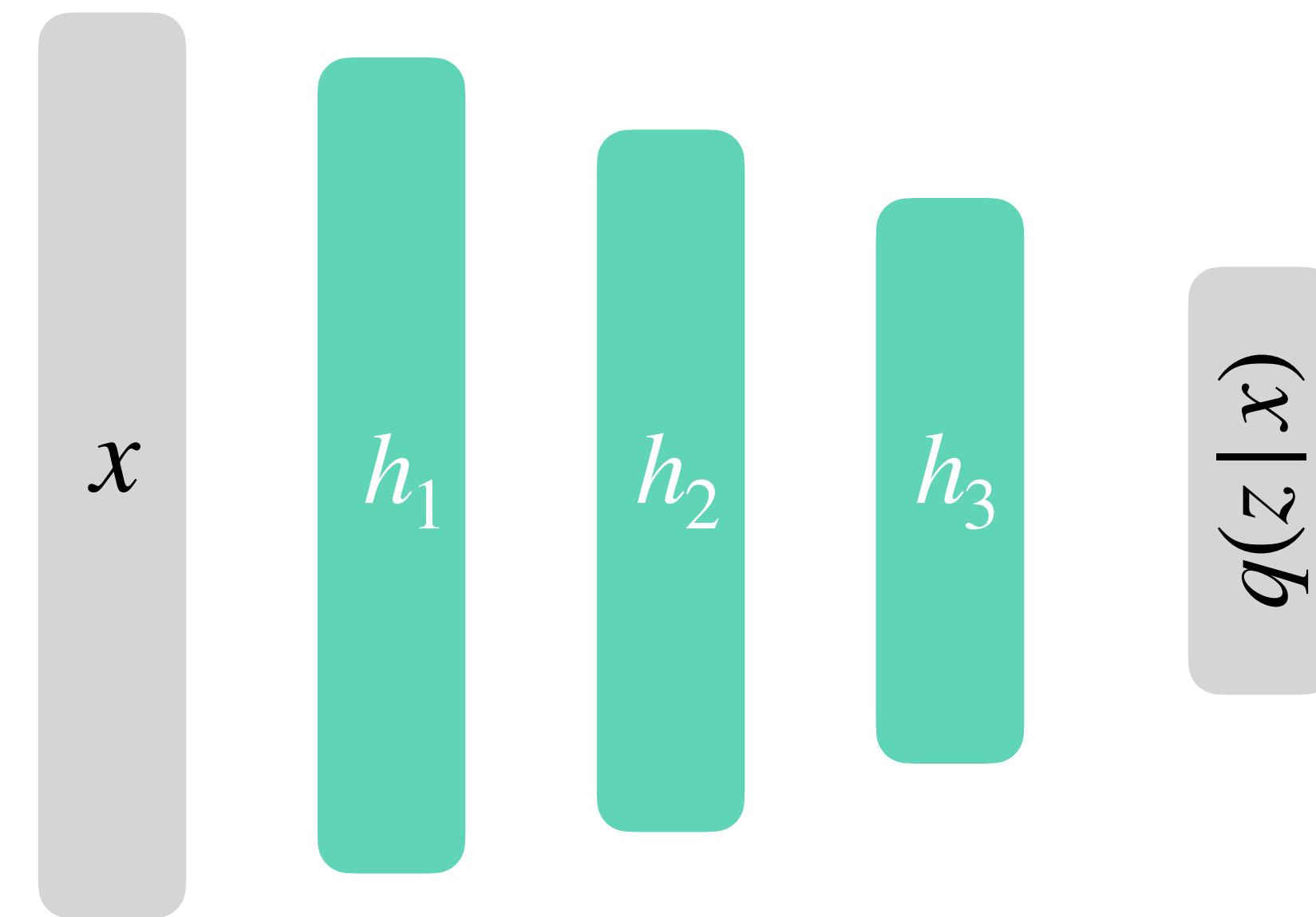
$$\begin{array}{ll} x_1 & \xrightarrow{f \approx f_\phi} q(z | x_1) \\ x_2 & \xrightarrow{f \approx f_\phi} q(z | x_2) \\ \vdots & \\ x_N & \xrightarrow{f \approx f_\phi} q(z | x_N) \end{array}$$



- treat learning  $f$  as a supervised regression problem
- approximate the mapping with a simpler function (NN)
- inference network

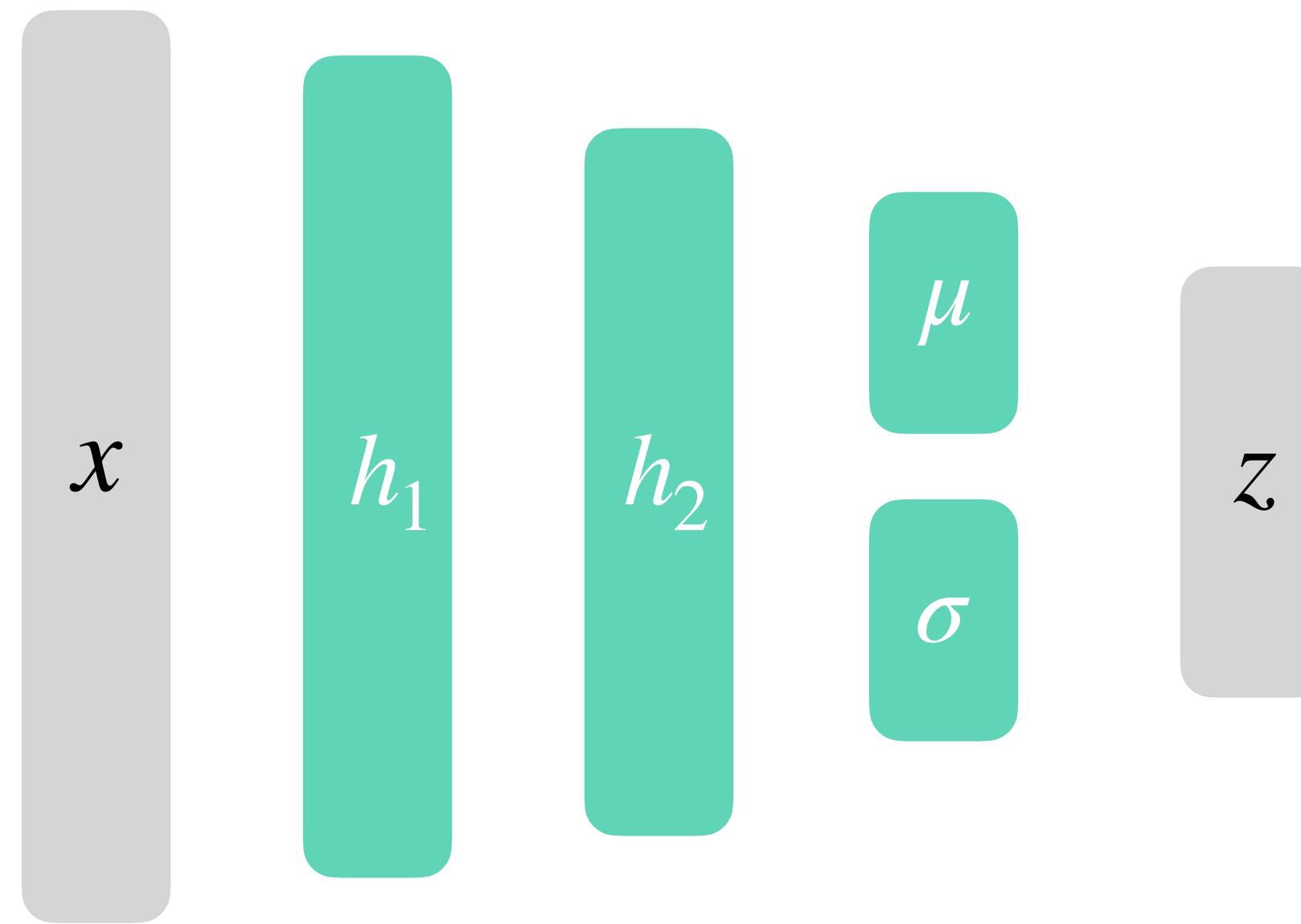
# how do we do inference?

reparameterisation trick



# how do we do inference?

## reparameterisation trick

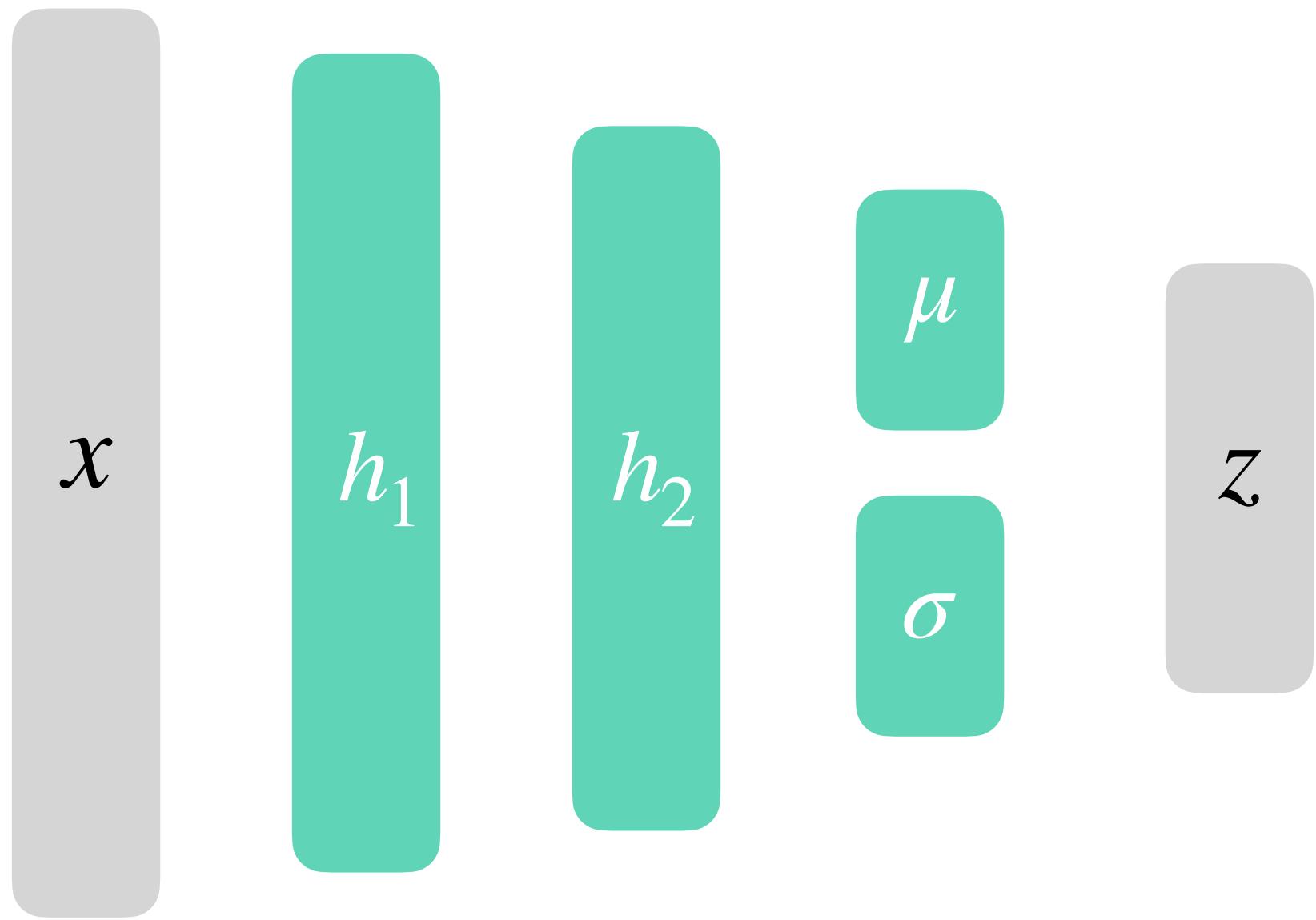


$$q_\phi(z|x) = \mathcal{N}(z|f_\phi^\mu(x), f_\phi^\sigma(x))$$

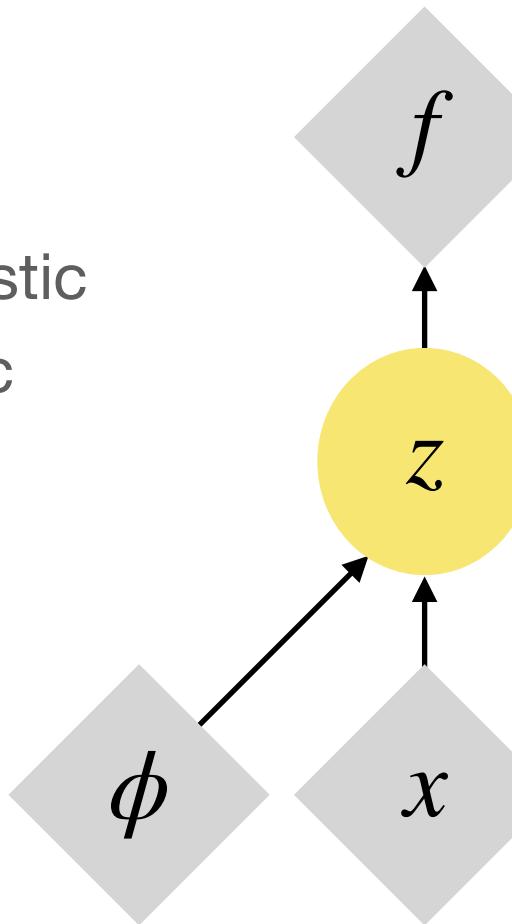
$$z \sim q_\phi(z|x)$$

# how do we do inference?

## reparameterisation trick

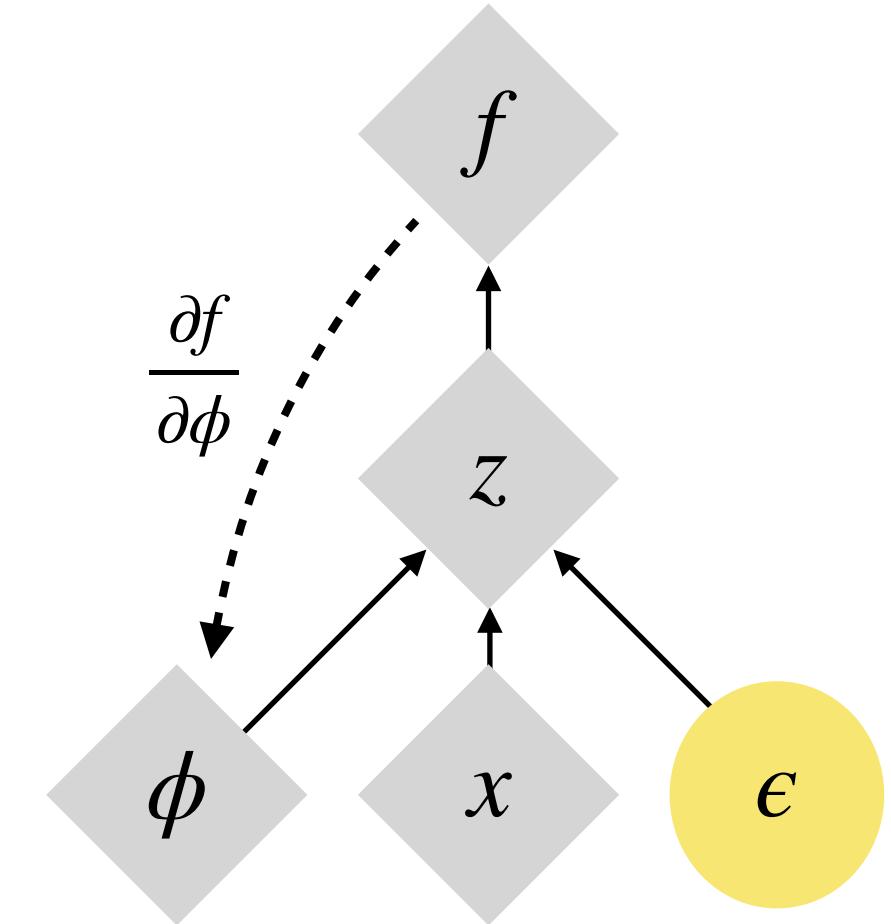


◆ deterministic  
● stochastic



$$q_\phi(z|x) = \mathcal{N}(z | f_\phi^\mu(x), f_\phi^\sigma(x))$$

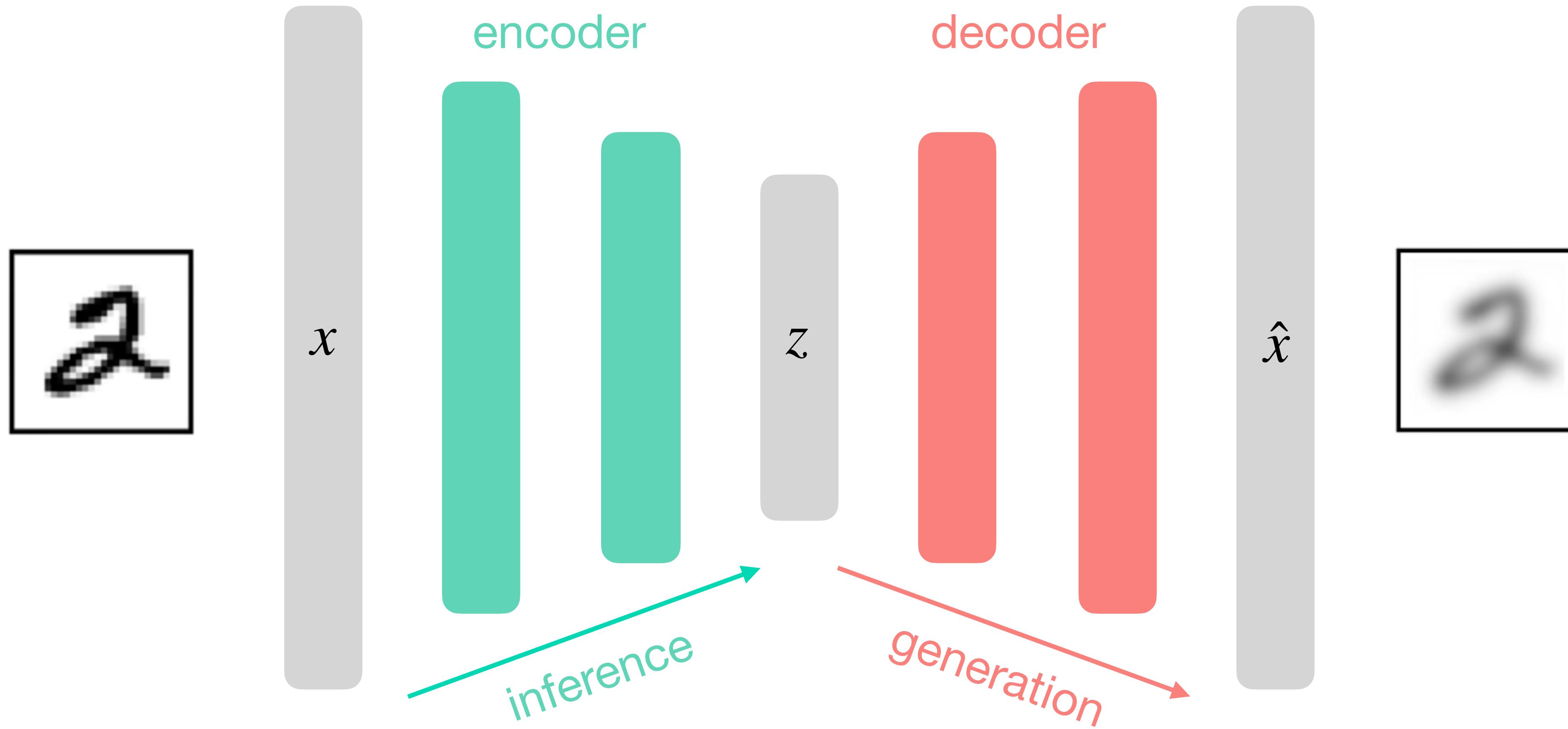
$$z \sim q_\phi(z|x)$$



$$z = \mu + \sigma \odot \epsilon$$

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

# Variational Autoencoder



# Combined objective evidence lower bound (ELBO)

$\log$  🐍

$KL$ [🐰 | 🐉]

$$\log p_{\theta}(x) - KL[q_{\phi}(z|x) || p(z|x)] = \mathcal{L}_{ELBO}(\theta, \phi, x)$$

# Combined objective evidence lower bound (ELBO)

$\log$  🐍

$\text{KL}$ [🐰 | 🐉]

$$\log p_\theta(x) - KL[q_\phi(z|x) || p(z|x)] = \mathcal{L}_{ELBO}(\theta, \phi, x)$$

$$\mathcal{L}_{ELBO}(\theta, \phi, x) = \mathbb{E}_{🐰} \log 🐹 - \mathbb{E}_{🐰} \text{KL}[🐰 | 🐥]$$

# Combined objective evidence lower bound (ELBO)

$\log$  🐍

$\text{KL}$ [🐰 | 🎭]

$$\log p_\theta(x) - \text{KL}[q_\phi(z|x) || p(z|x)] = \mathcal{L}_{ELBO}(\theta, \phi, x)$$

$\mathbb{E}_{\text{🐰}}$   $\log$  🐹

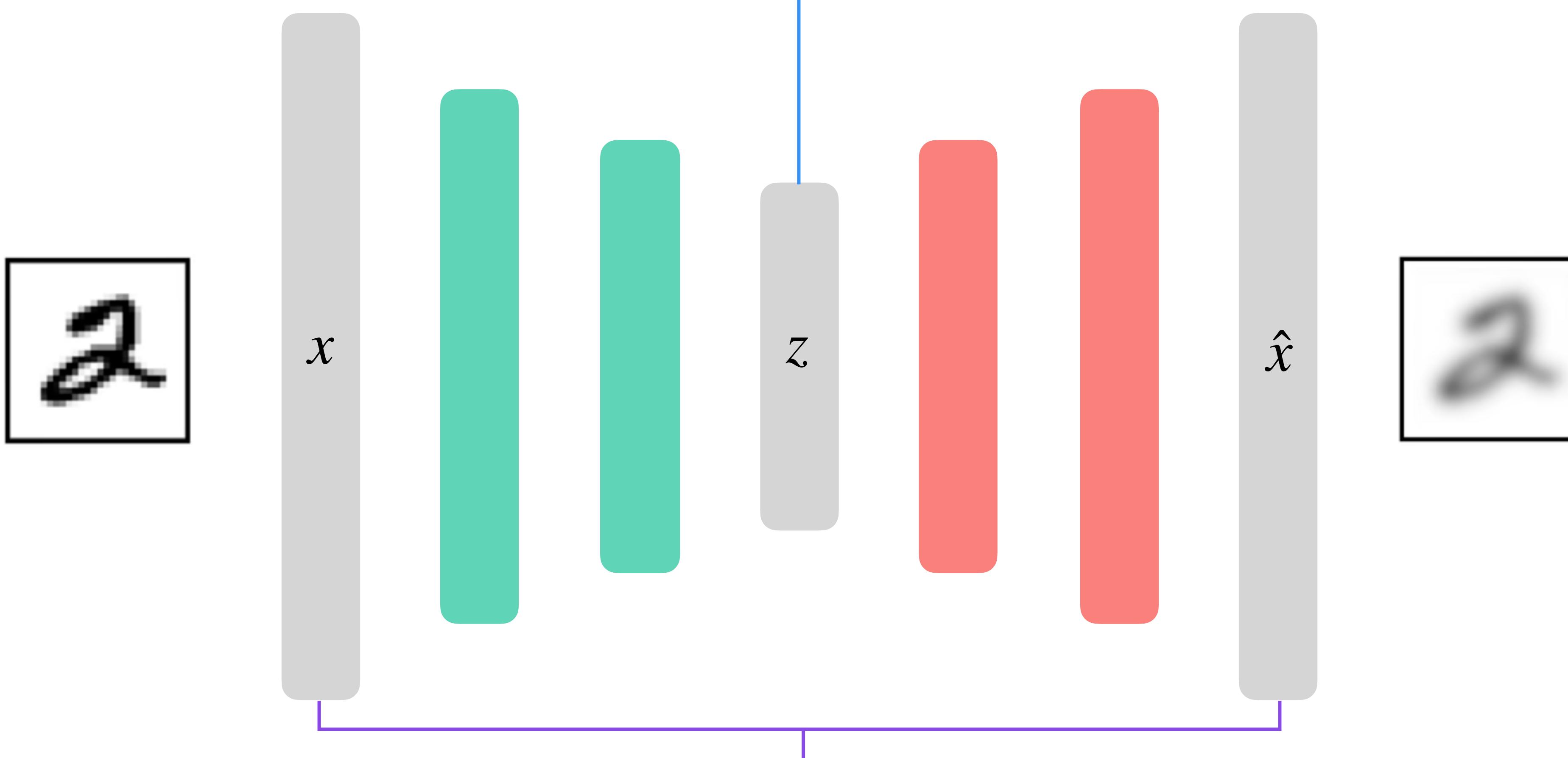
$\mathbb{E}_{\text{🐰}}$   $\text{KL}$ [🐰 | 🐥]

$$\mathcal{L}_{ELBO}(\theta, \phi, x) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \text{KL}[q_\phi(z|x) || p(z)]$$

reconstruction      regularisation

$$-KL[q_\phi(z|x) || p(z)]$$

regularisation



reconstruction

$$\mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] = \|x - \hat{x}\|^2$$

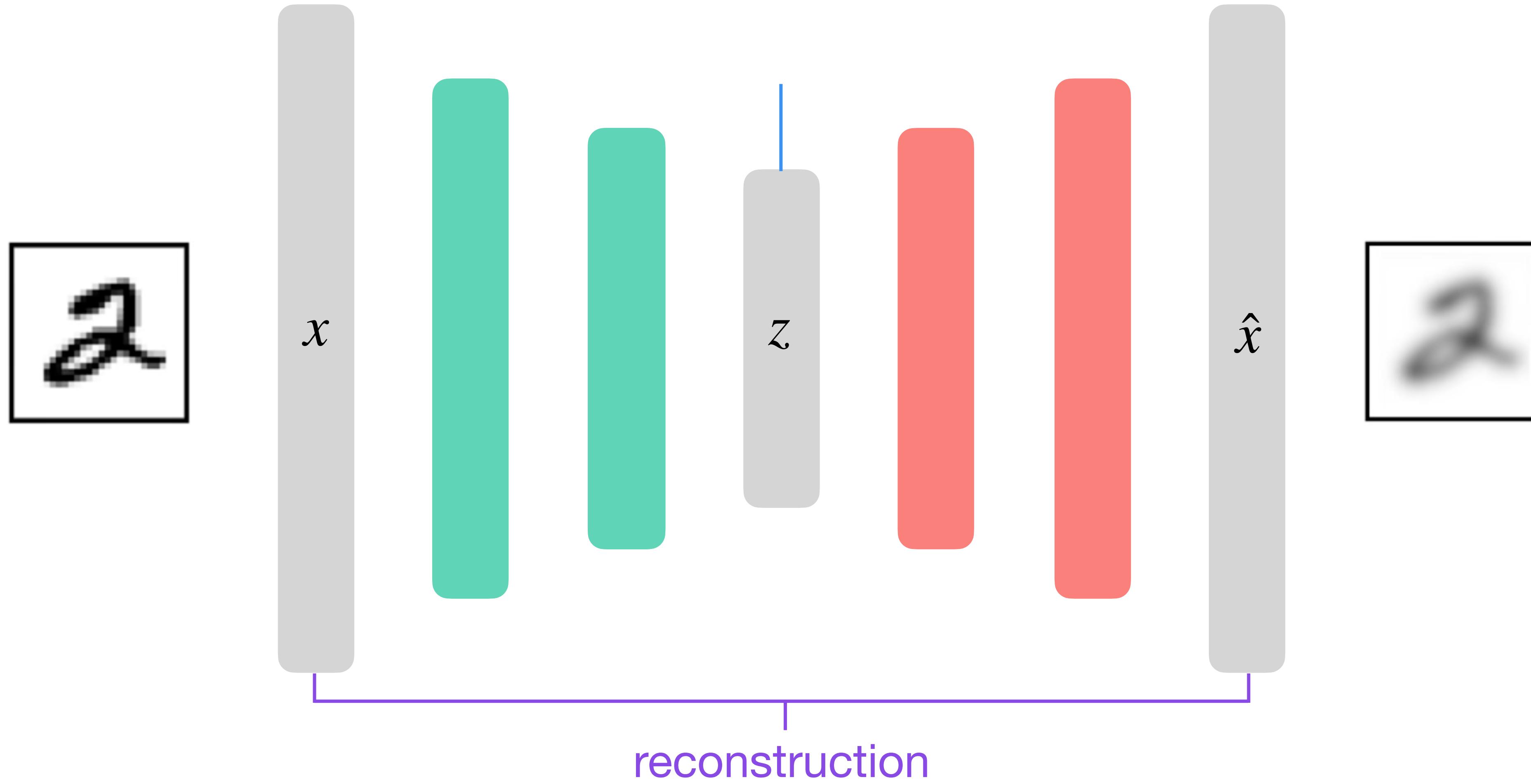
for gaussian noise model

$-KL[q_\phi(z|x) || p(z)]$   
regularisation

$$q_\phi(z|x) = \mathcal{N}(z|f_\phi^\mu(x), f_\phi^\sigma(x))$$
$$p(z) = \mathcal{N}(z|0, I)$$

analytical result for KL if Gaussian prior and posterior:

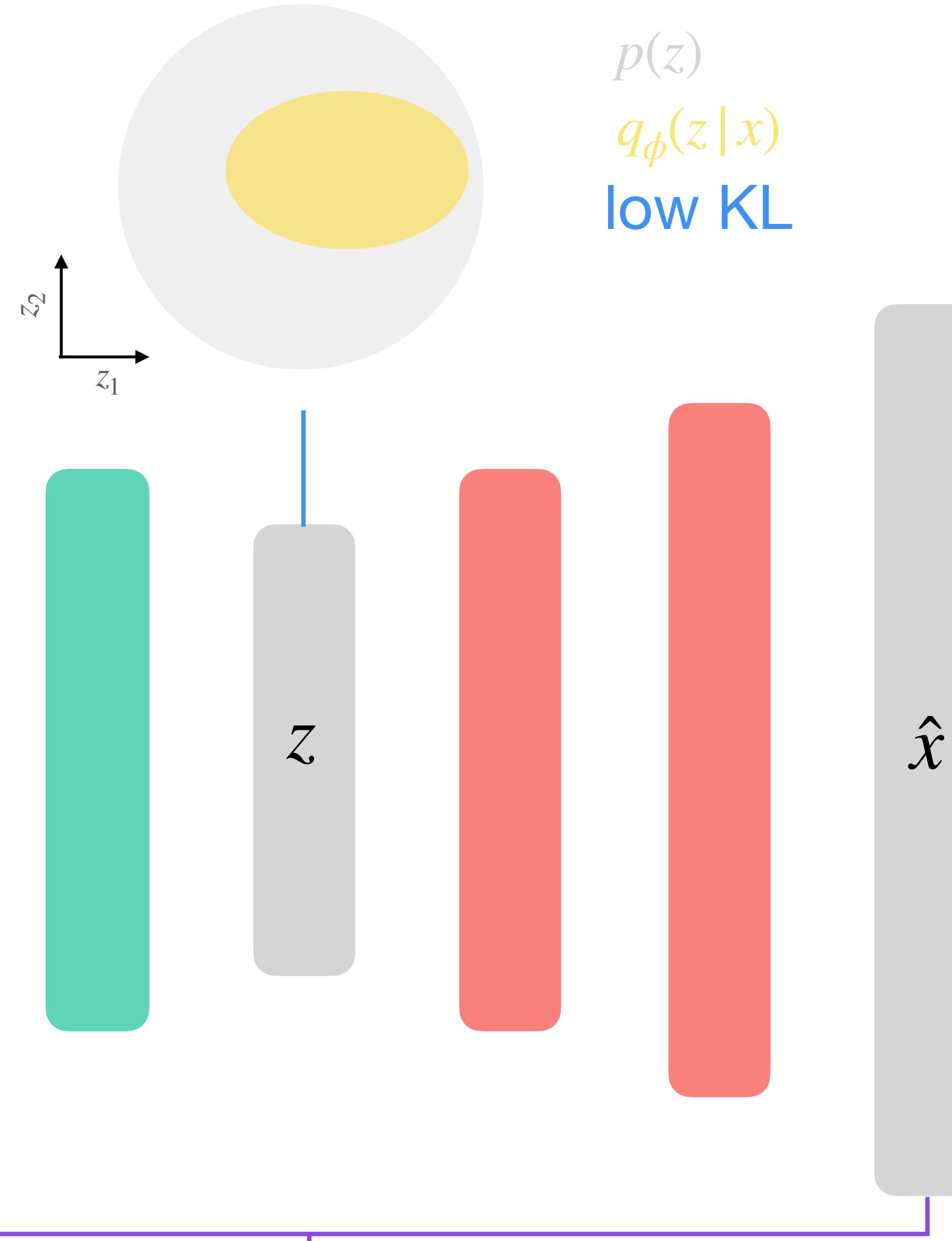
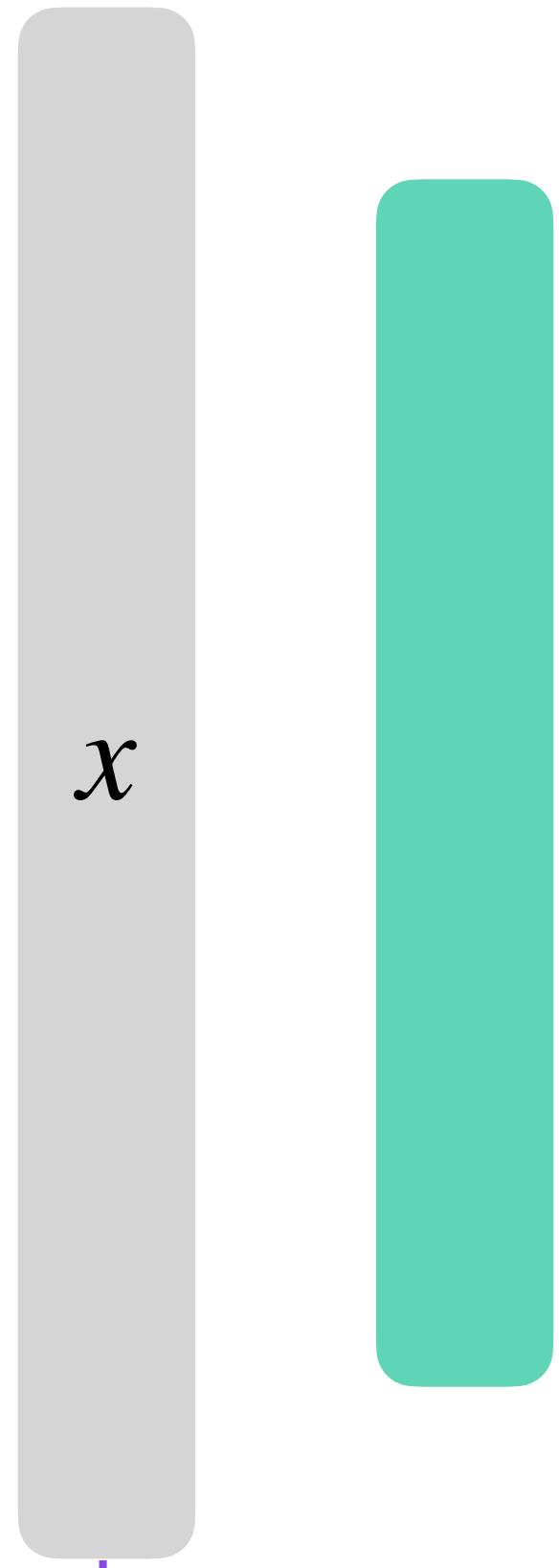
$$\frac{1}{2}(1 + \log \sigma^2 - \mu^2 - \sigma^2)$$



$$\mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)]$$

$$-KL[q_\phi(z|x) || p(z)]$$

regularisation



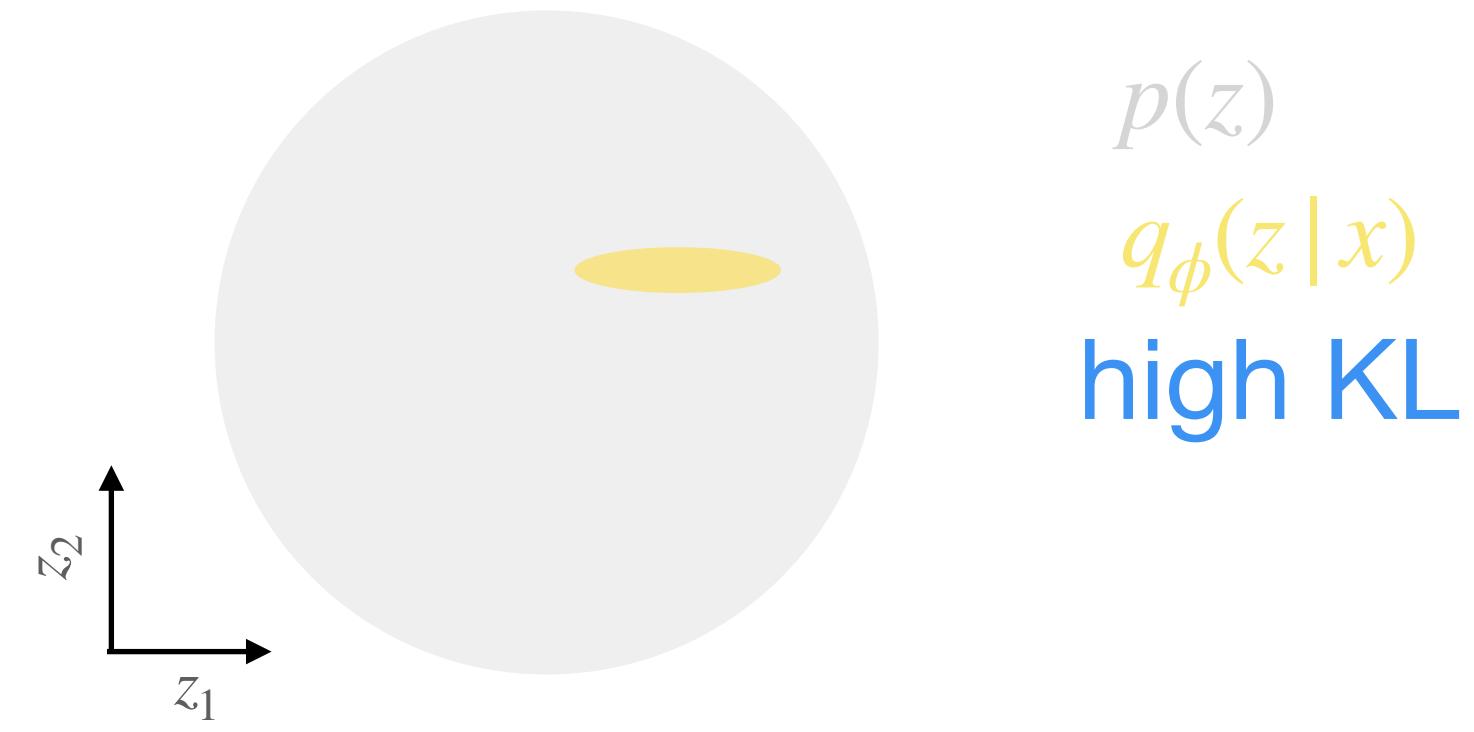
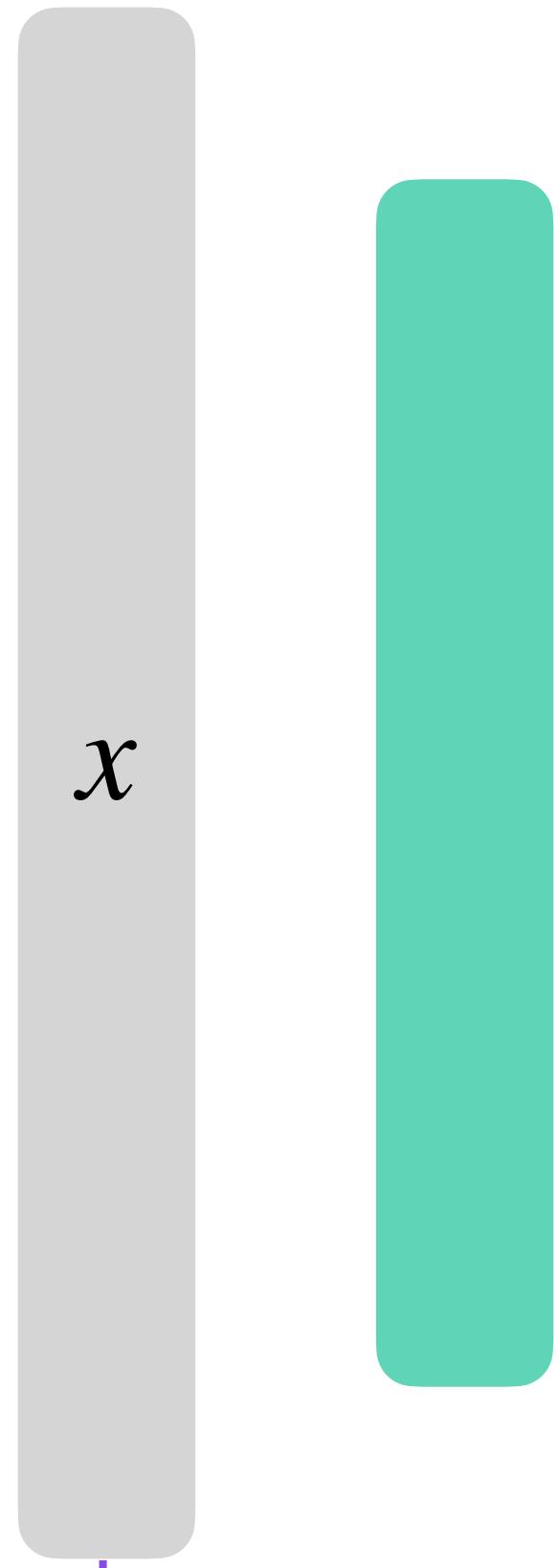
$p(z)$   
 $q_\phi(z|x)$   
low KL



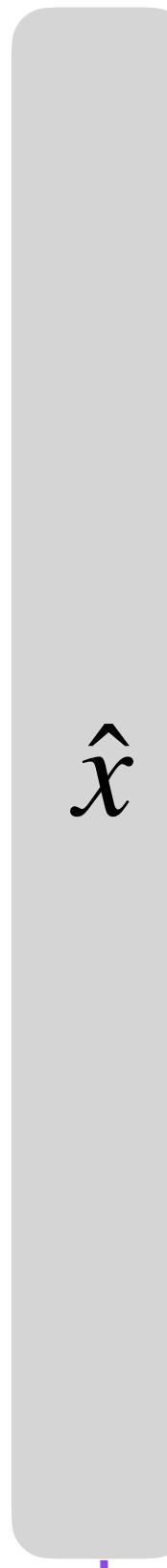
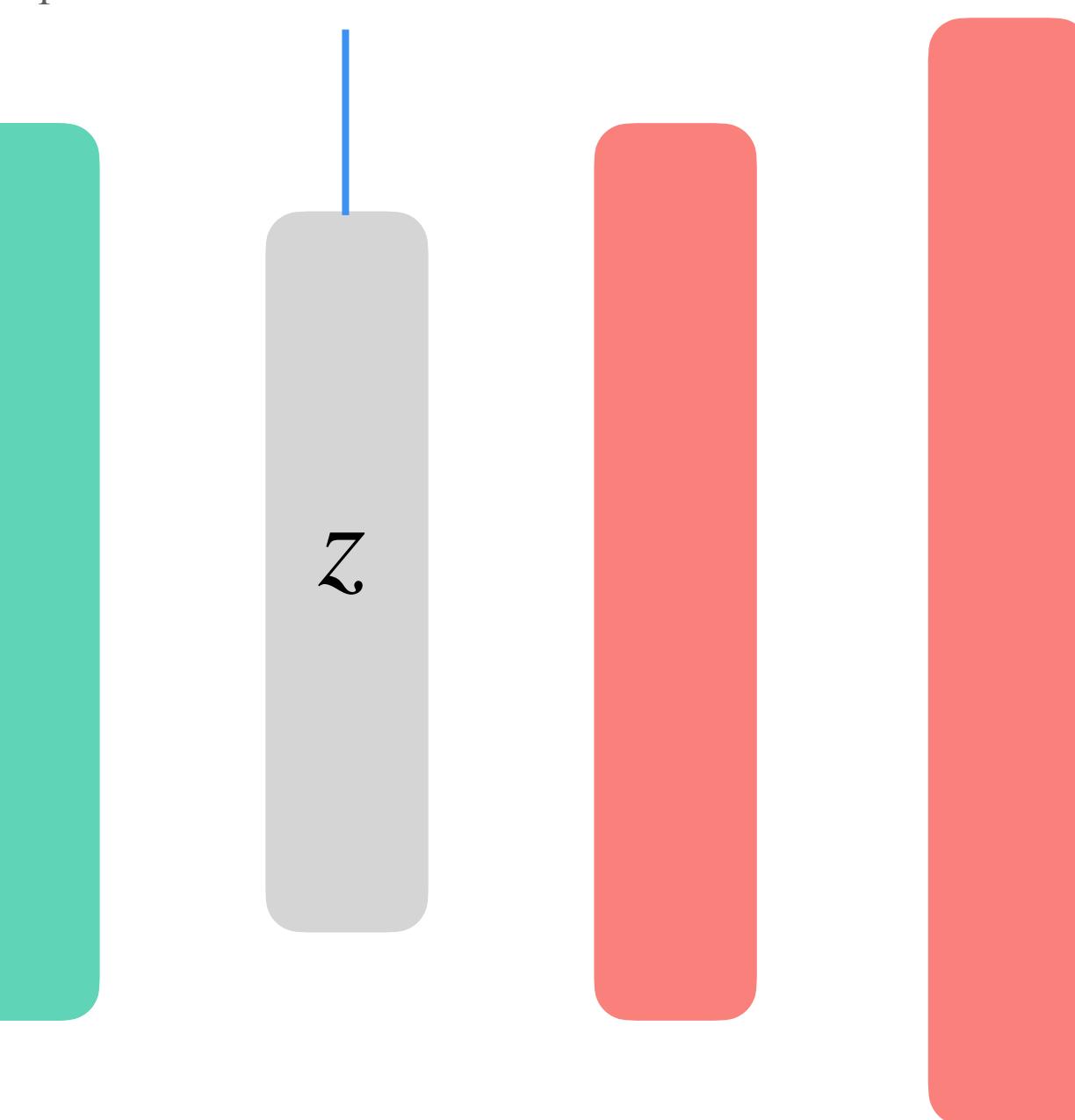
$$\mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)]$$

$$-KL[q_\phi(z|x) || p(z)]$$

regularisation



$p(z)$   
 $q_\phi(z|x)$   
high KL

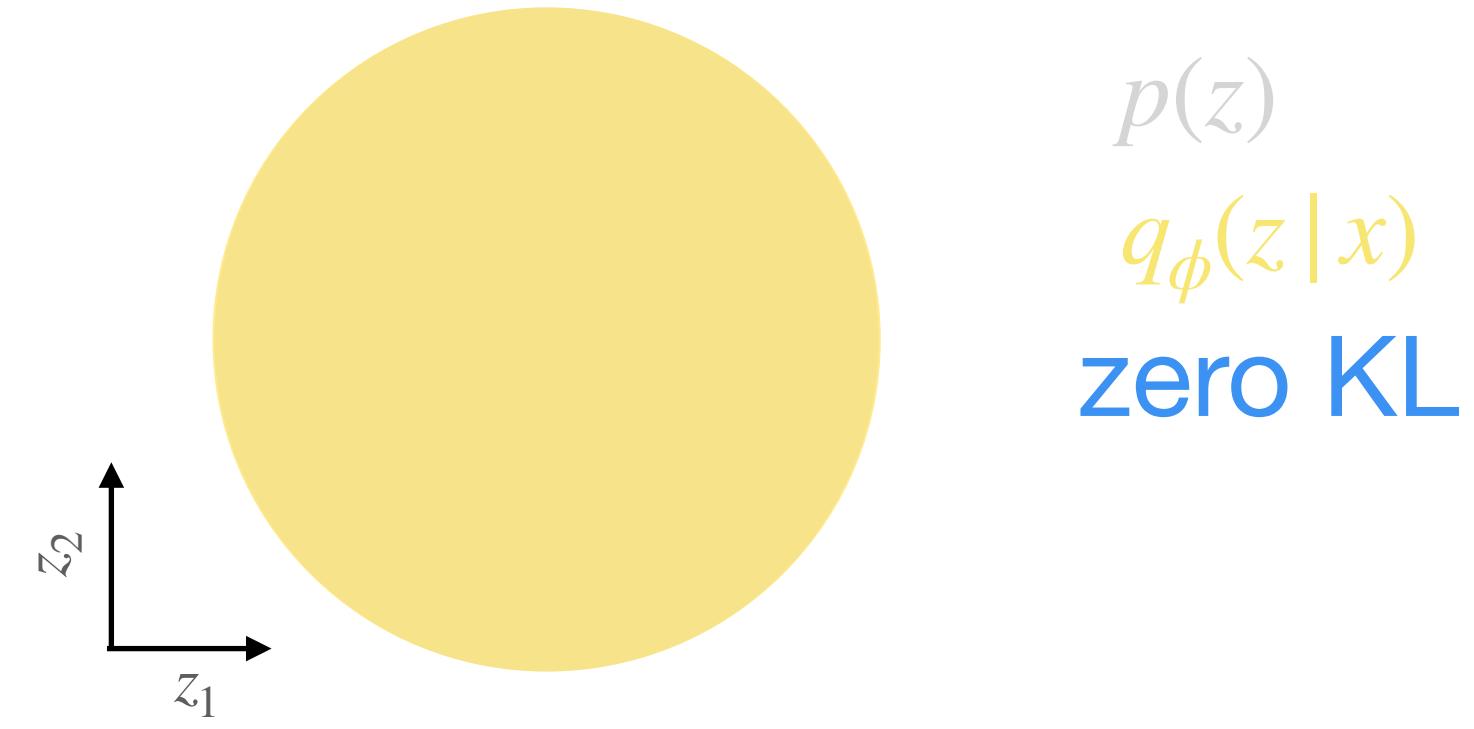
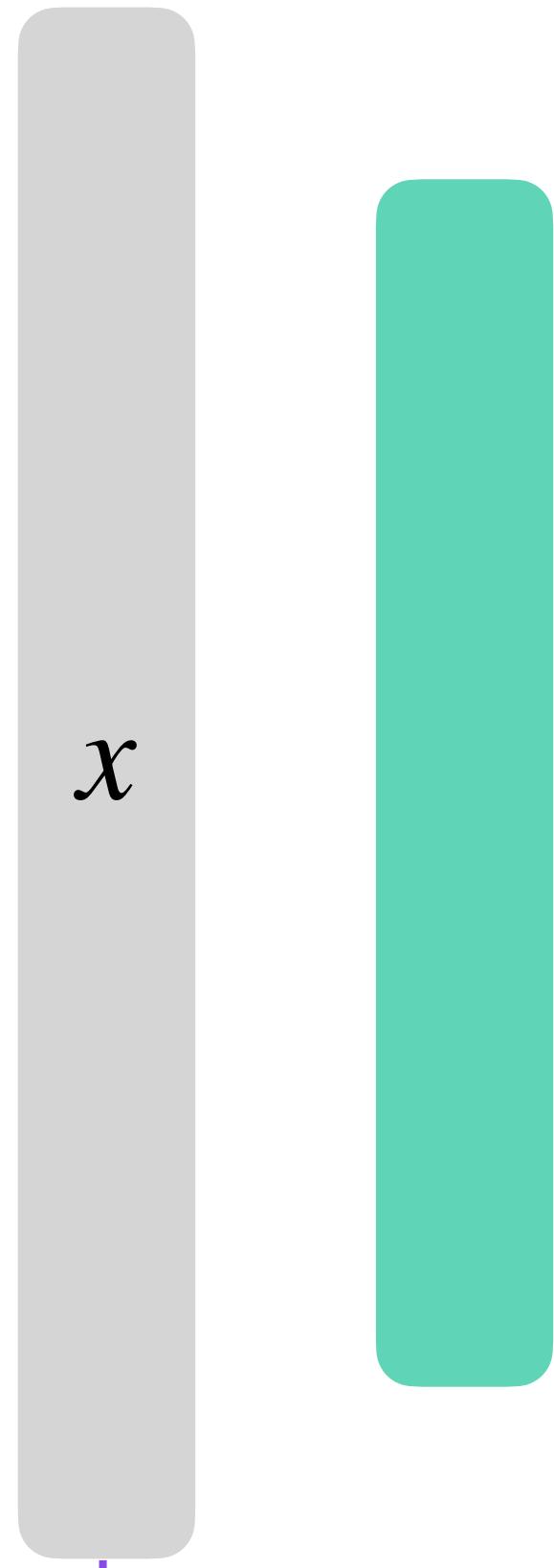


reconstruction

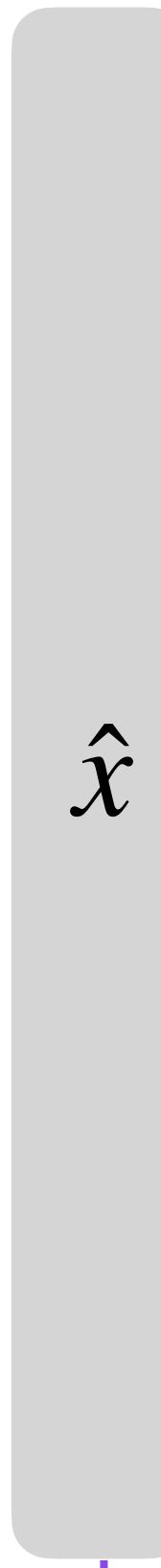
$$\mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)]$$

$$-KL[q_\phi(z|x) || p(z)]$$

regularisation



$p(z)$   
 $q_\phi(z|x)$   
zero KL



reconstruction

$$\mathbb{E}_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)]$$

# Demo

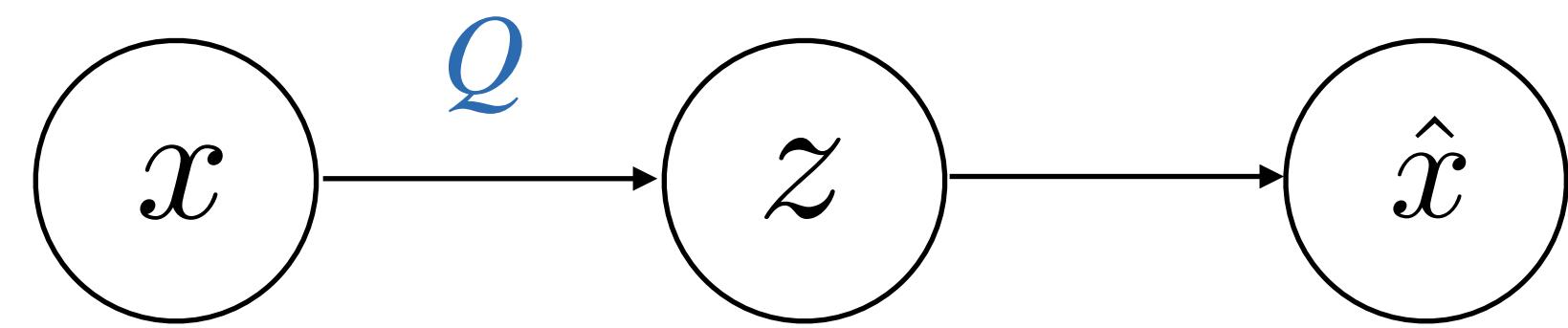
<https://github.com/eemlcommunity/PracticalSessions2021/tree/main/generative>

# **Part II**

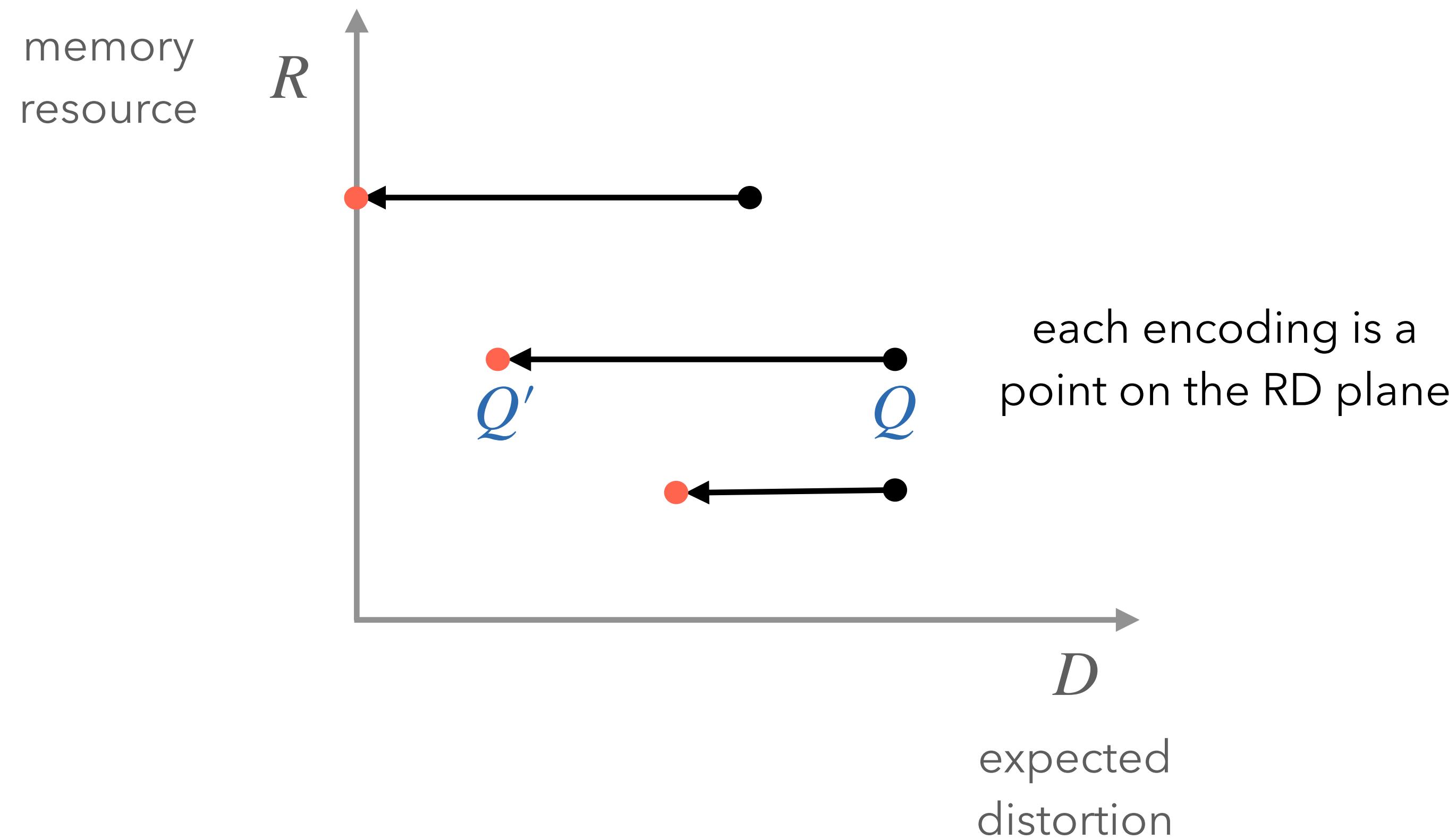
# **β-VAE**

$$\mathcal{L}(\theta, \phi, x) = \mathbb{E}_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - \textcolor{red}{\beta} \text{ } KL[q_\phi(z|x) || p(z)]$$

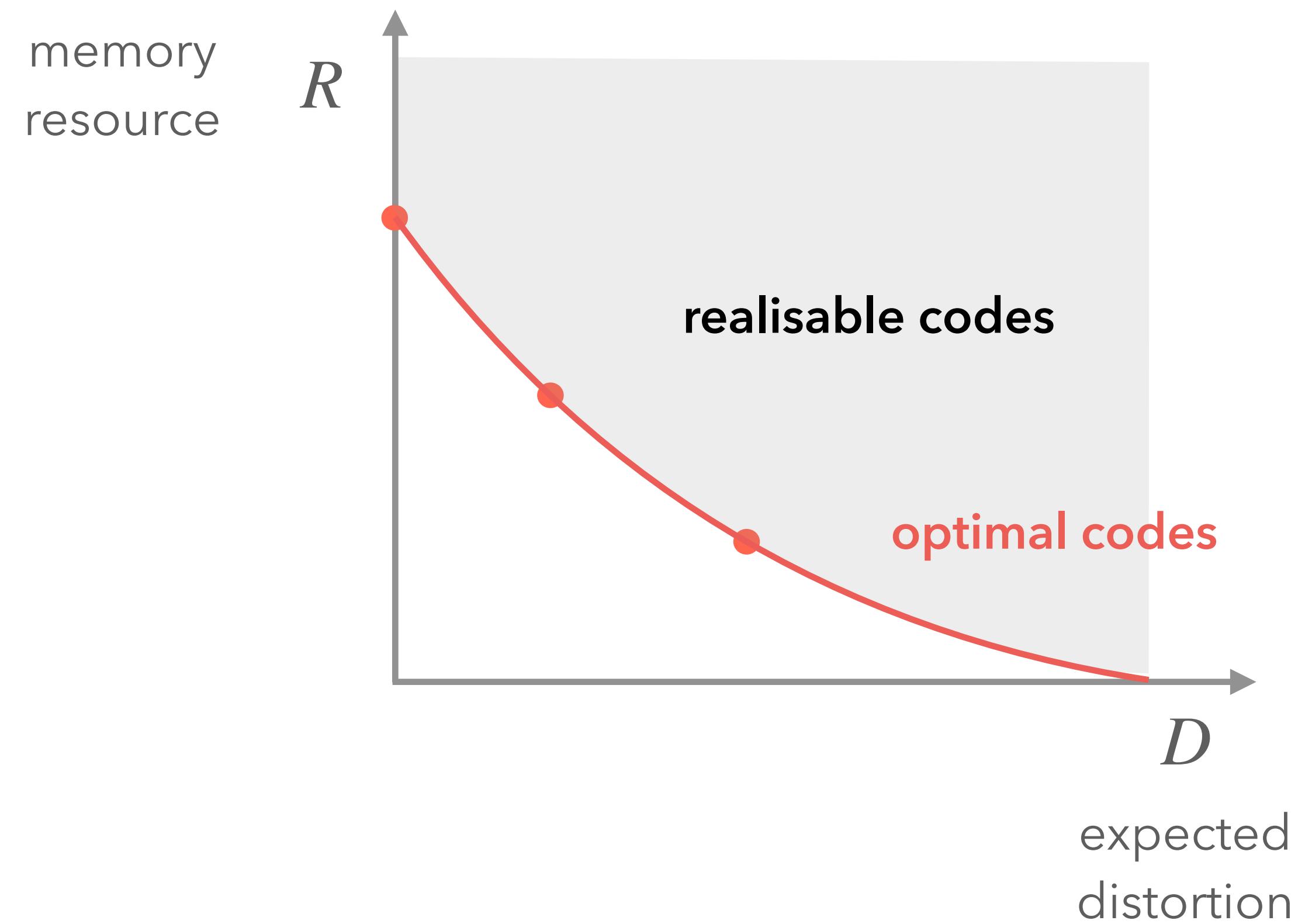
# Compression view



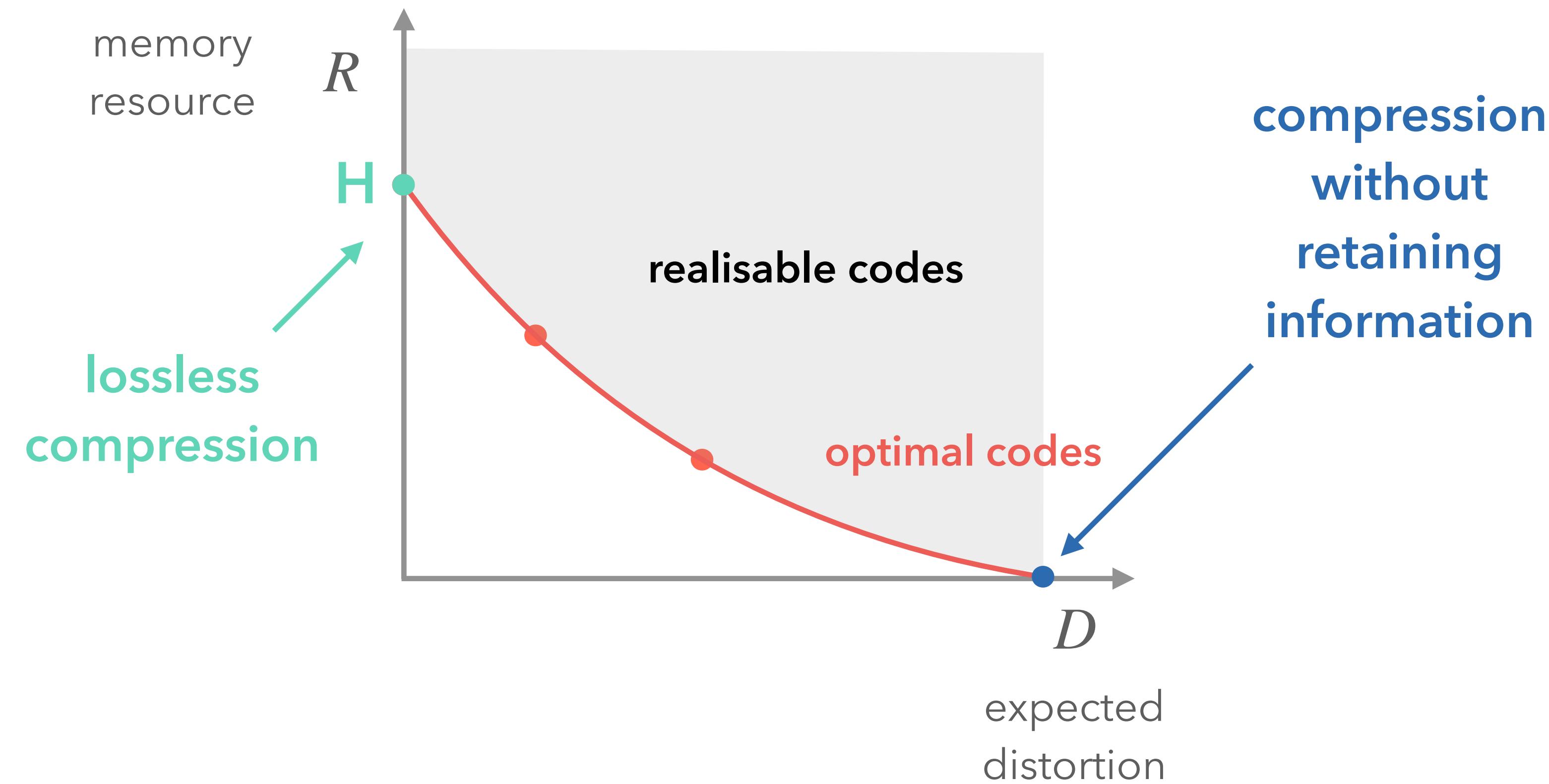
# Optimal compression



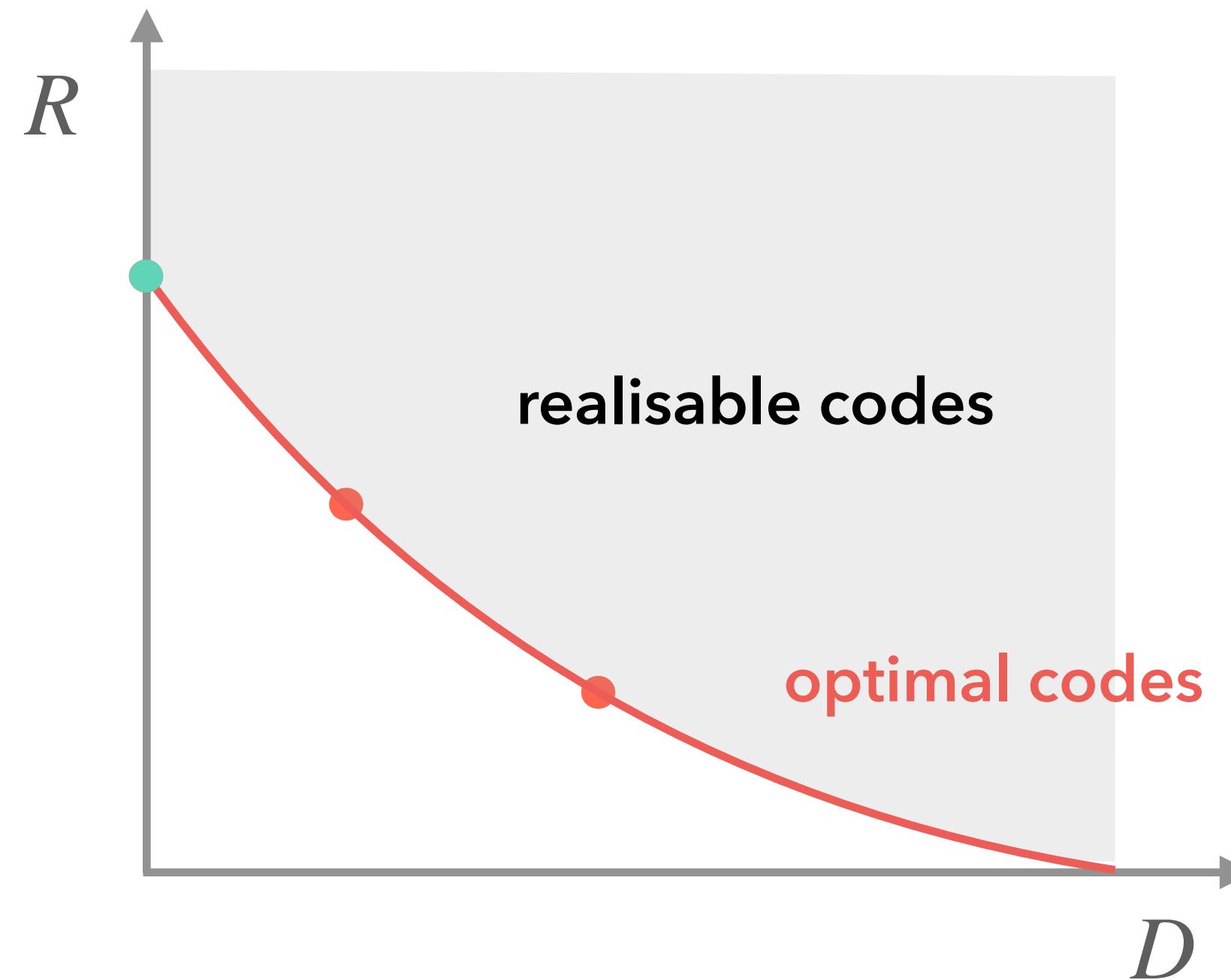
# Optimal compression



# Optimal compression



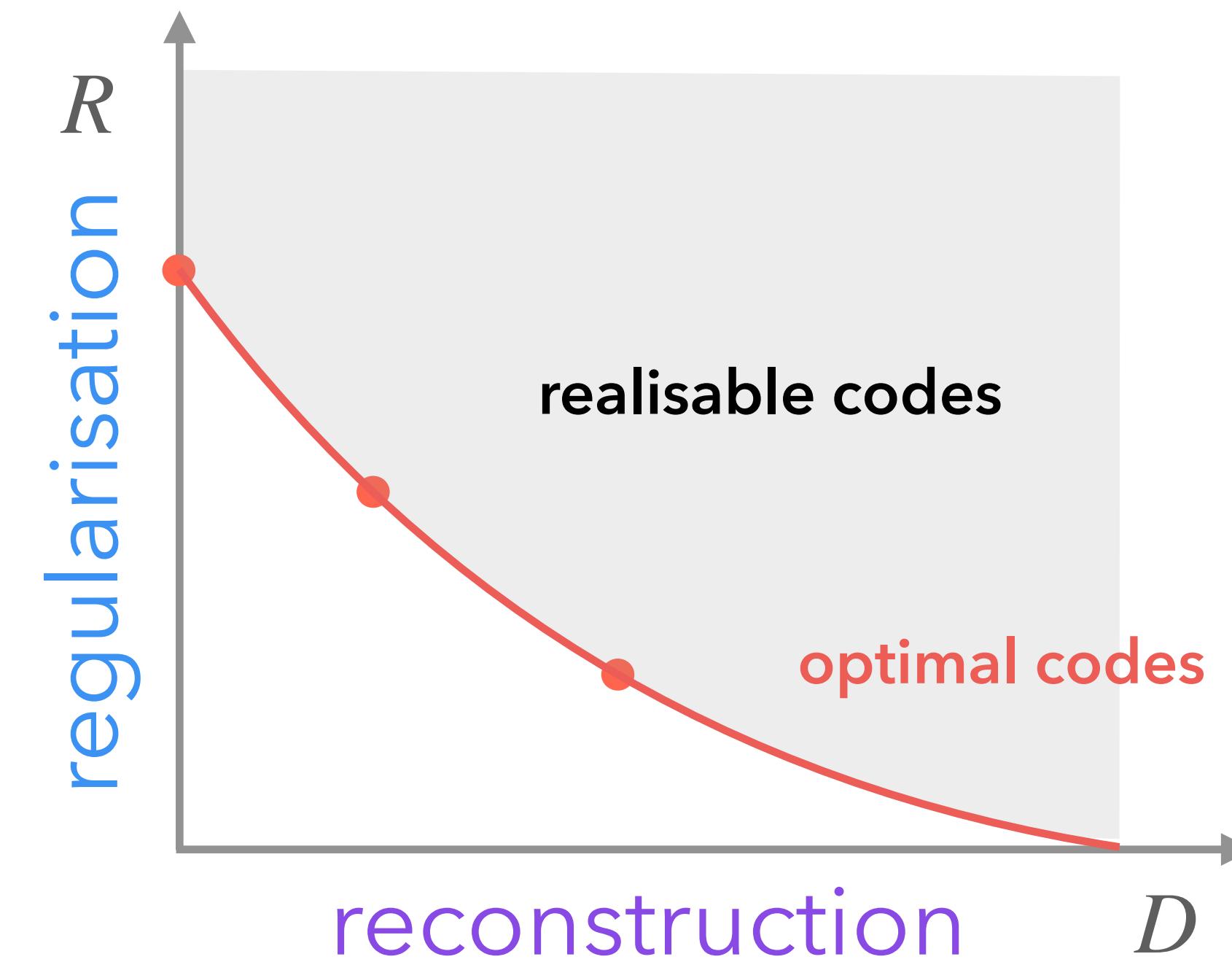
# Compression view of ELBO



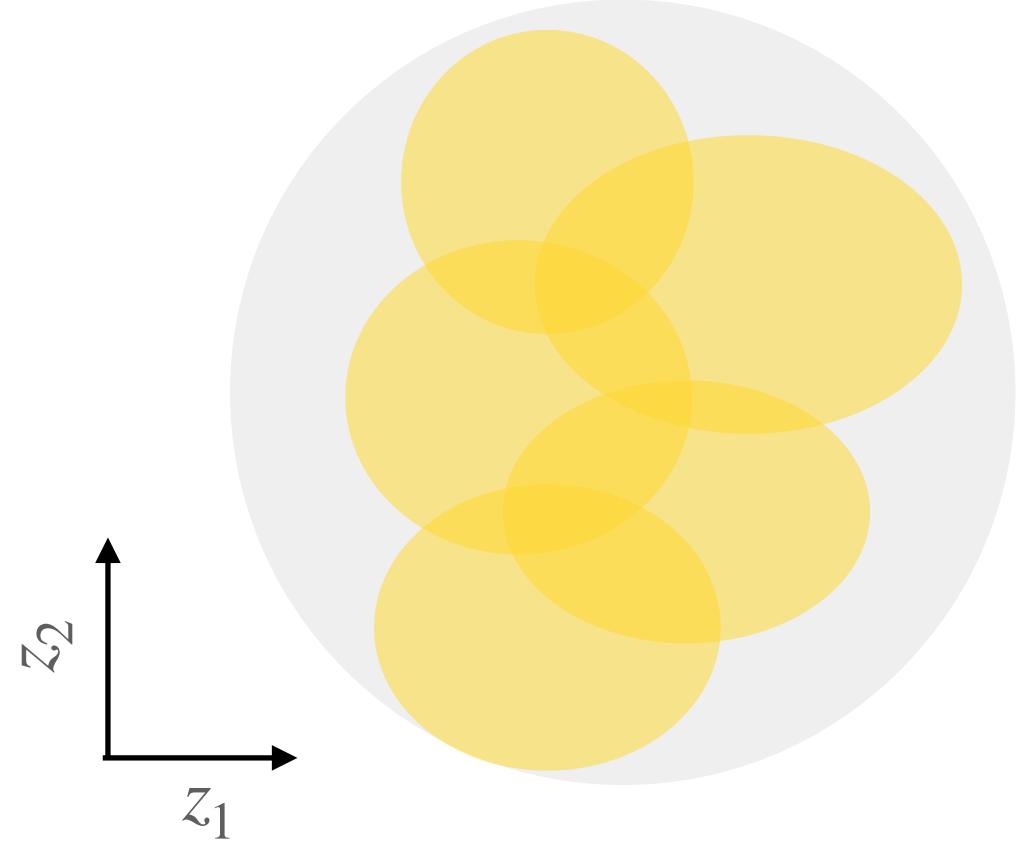
$$\mathcal{L}(\theta, \phi, x) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - KL[q_\phi(z|x) || p(z)]$$

reconstruction                      regularisation

# Compression view of ELBO

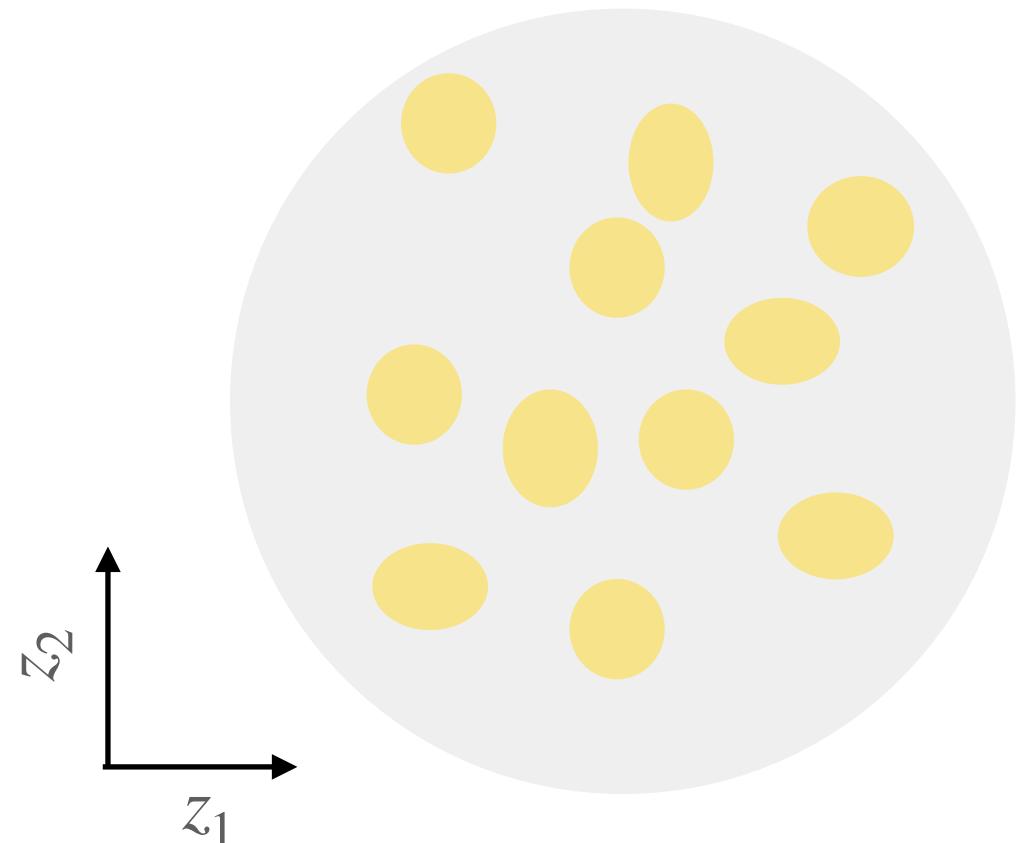


# Compression view of ELBO

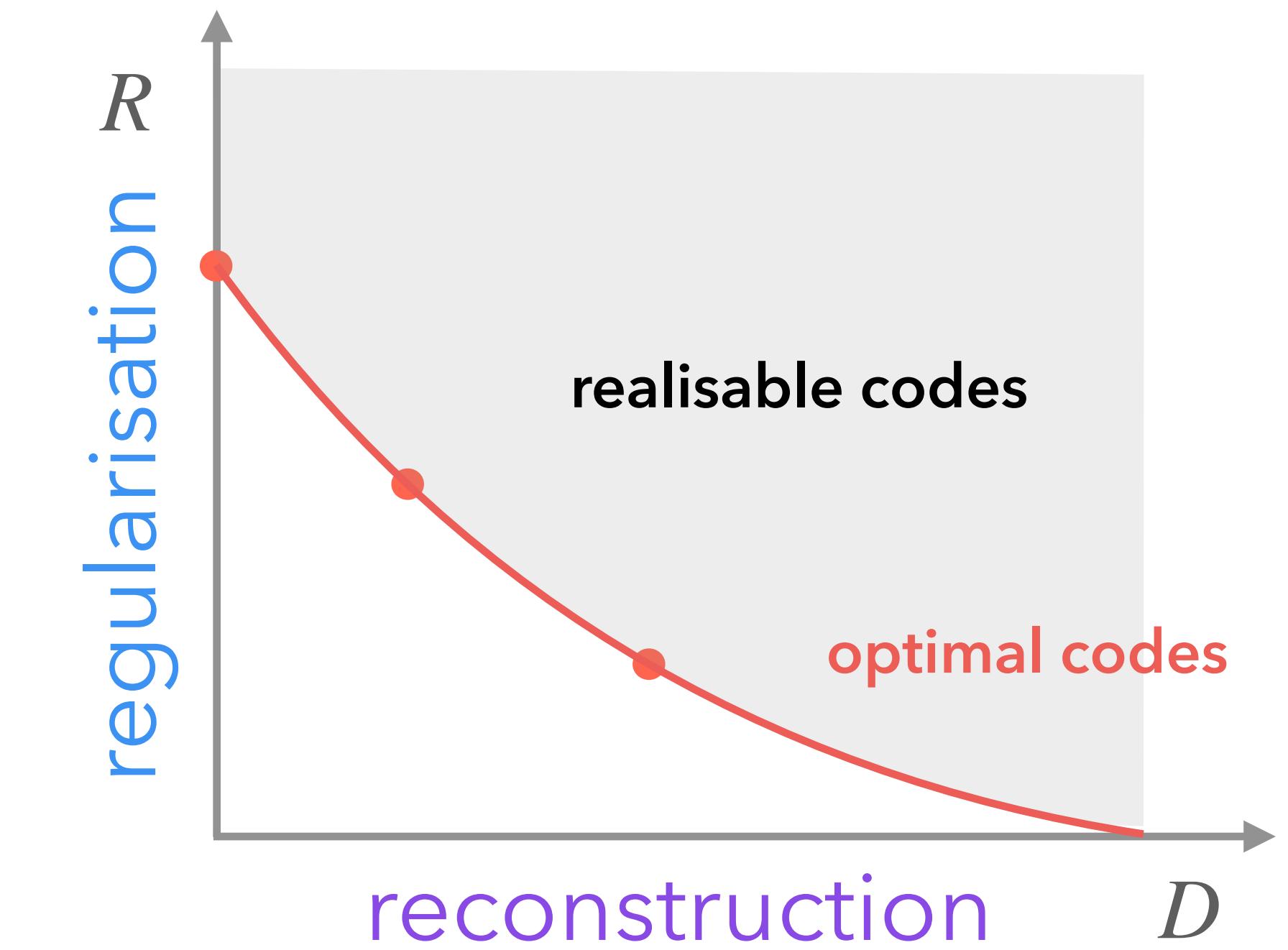


$$p(z)$$
$$q_{\phi}(z|x)$$

- low expected KL
- low mutual information



- high expected KL
- high mutual information



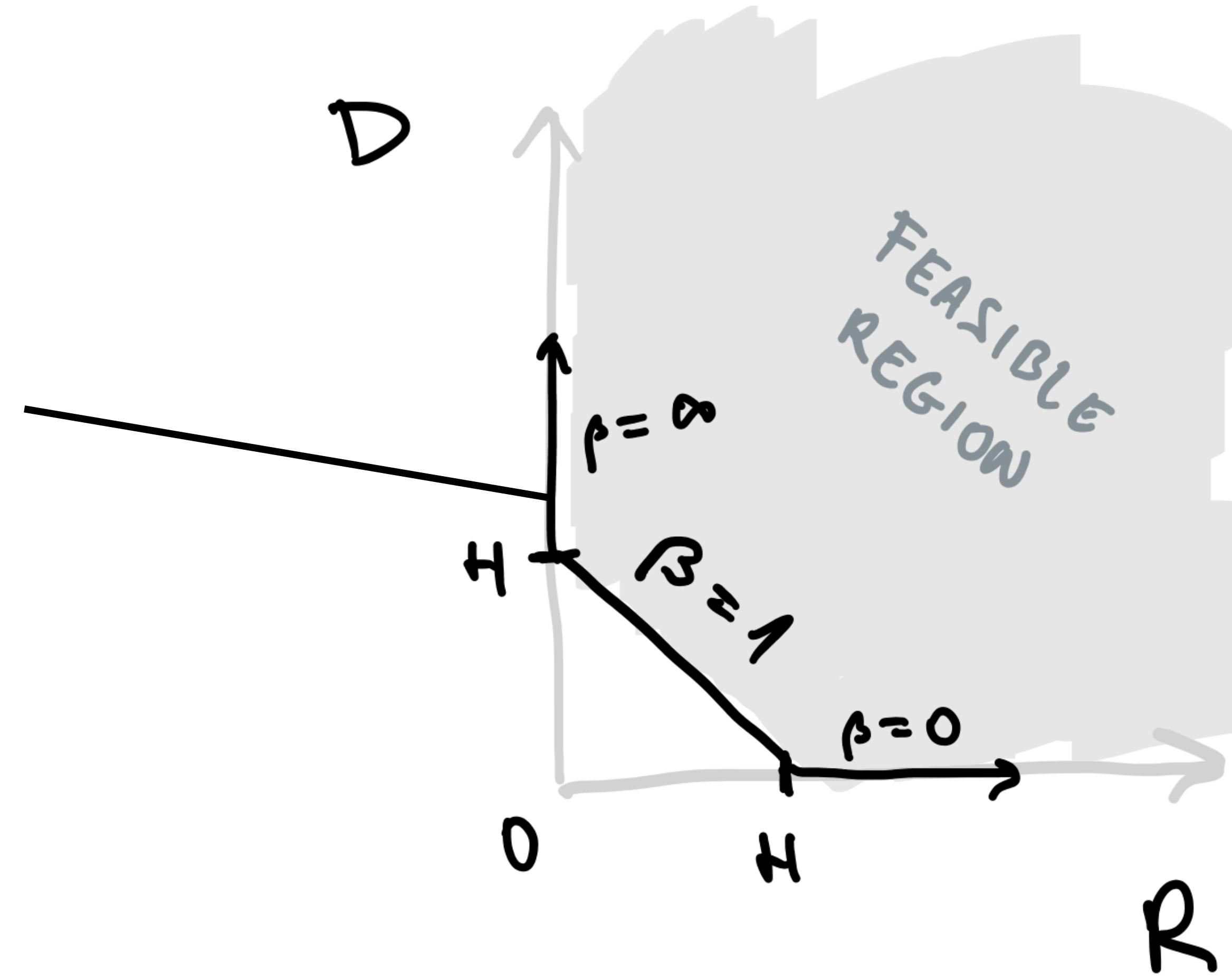
$$|x - \hat{x}|^2$$



# RD plane

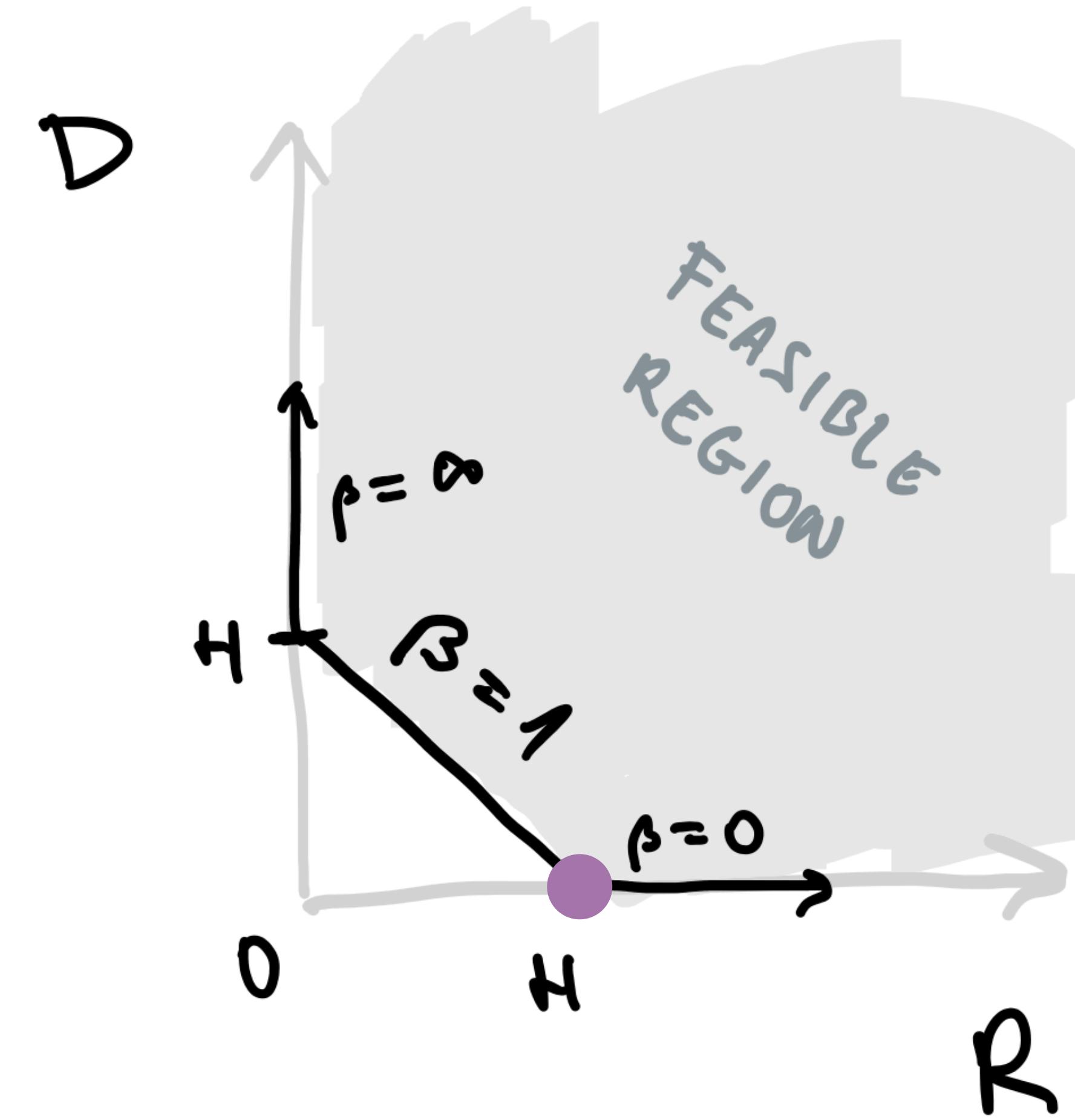
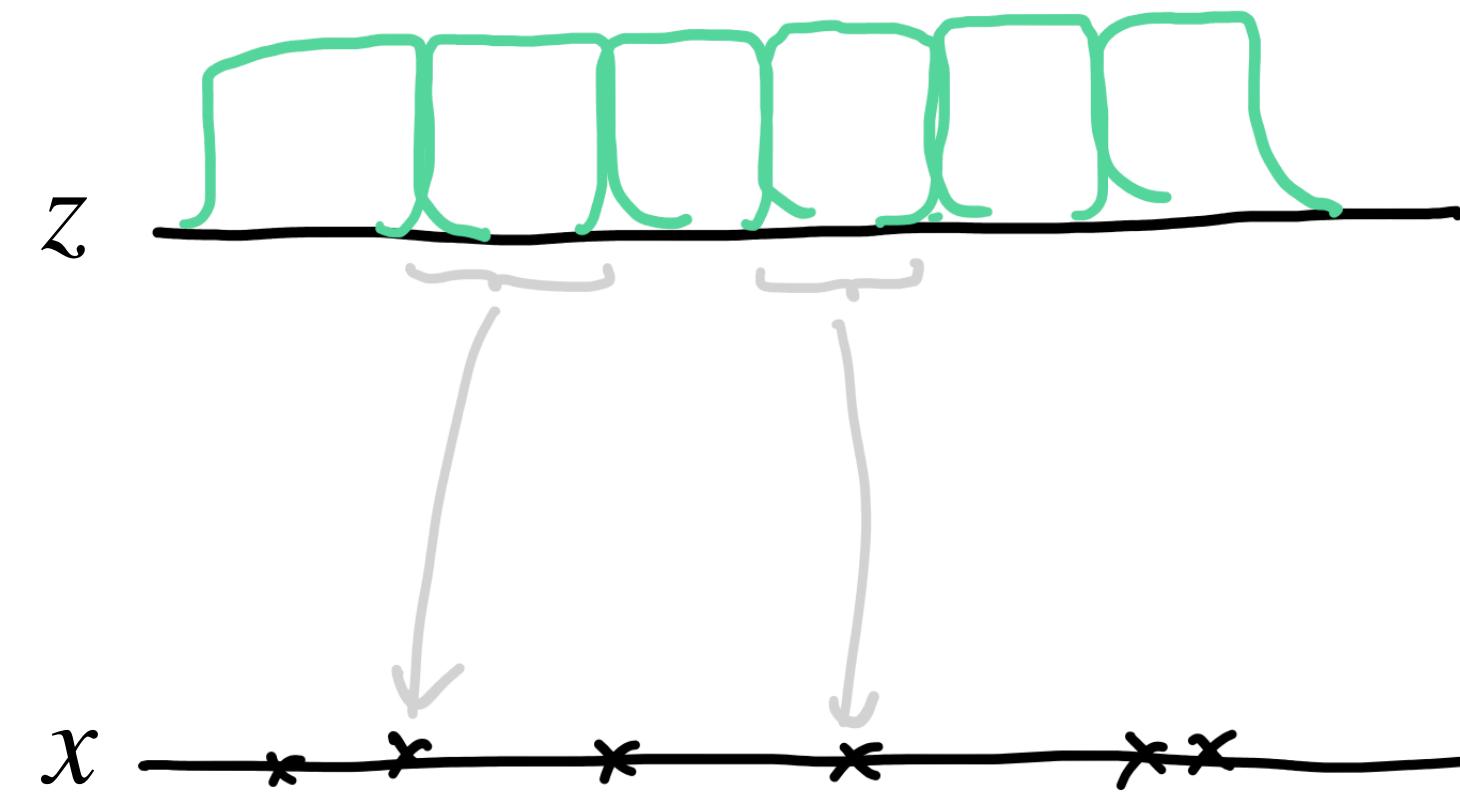
**no constraints on  
model family**

- encoder
- decoder
- prior



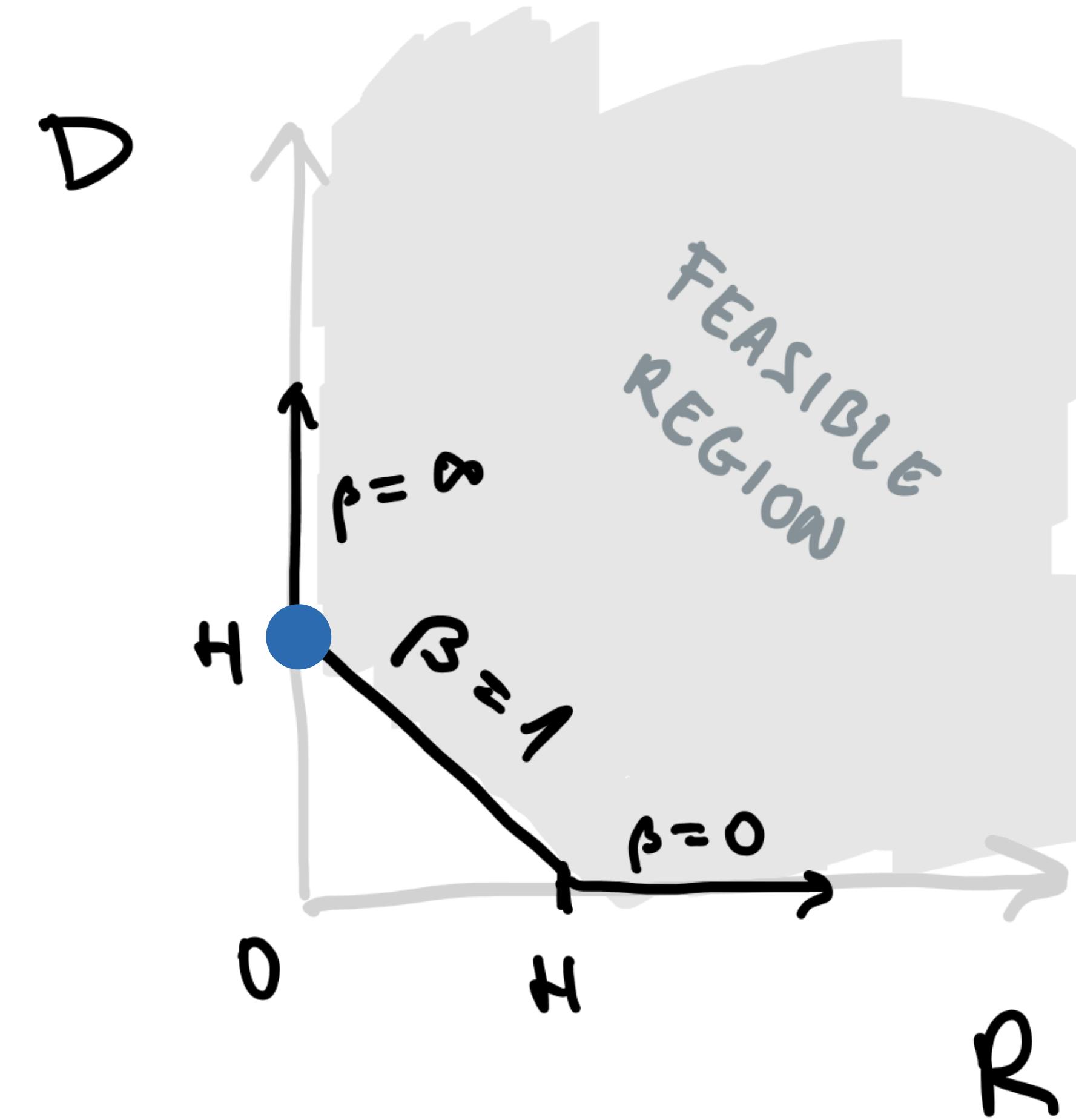
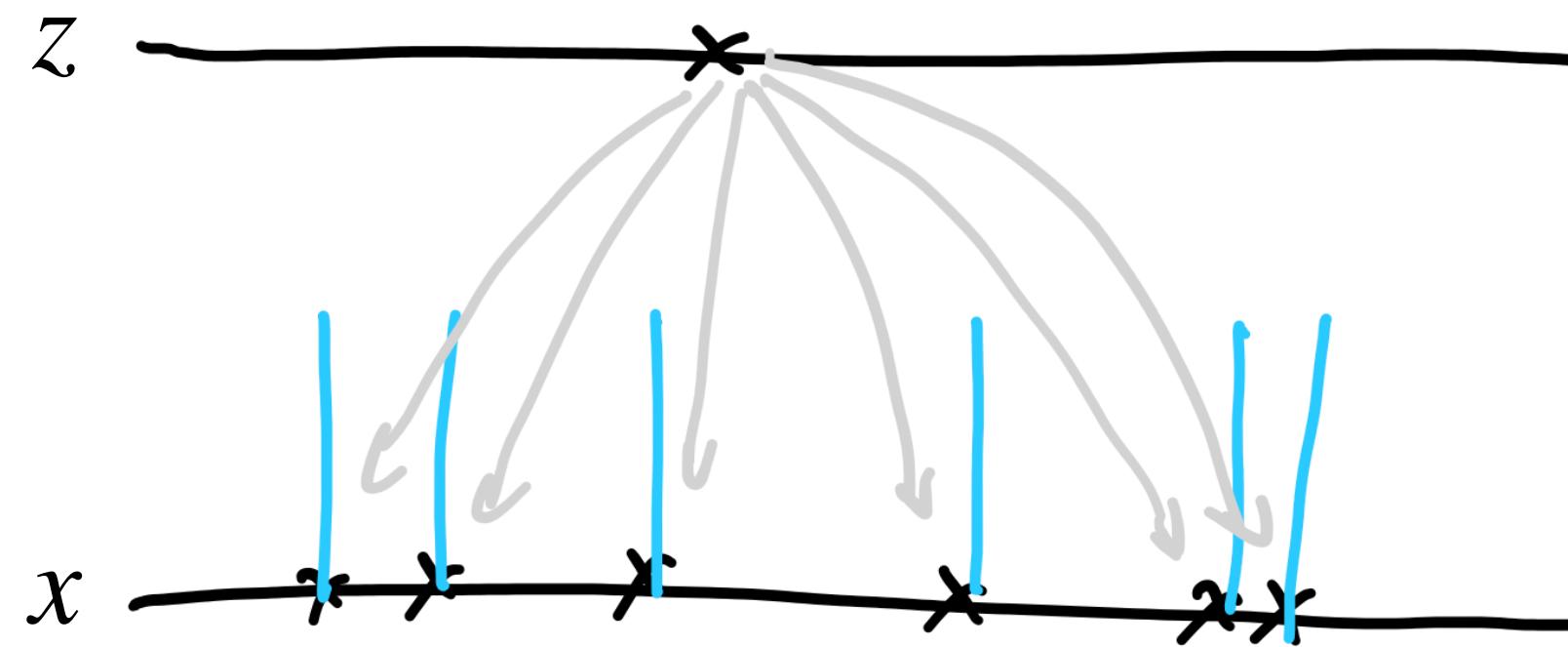
# Auto-encoding limit

- lossless compression
- store and return data points
- each region of  $Z$  corresponds to a datapoint



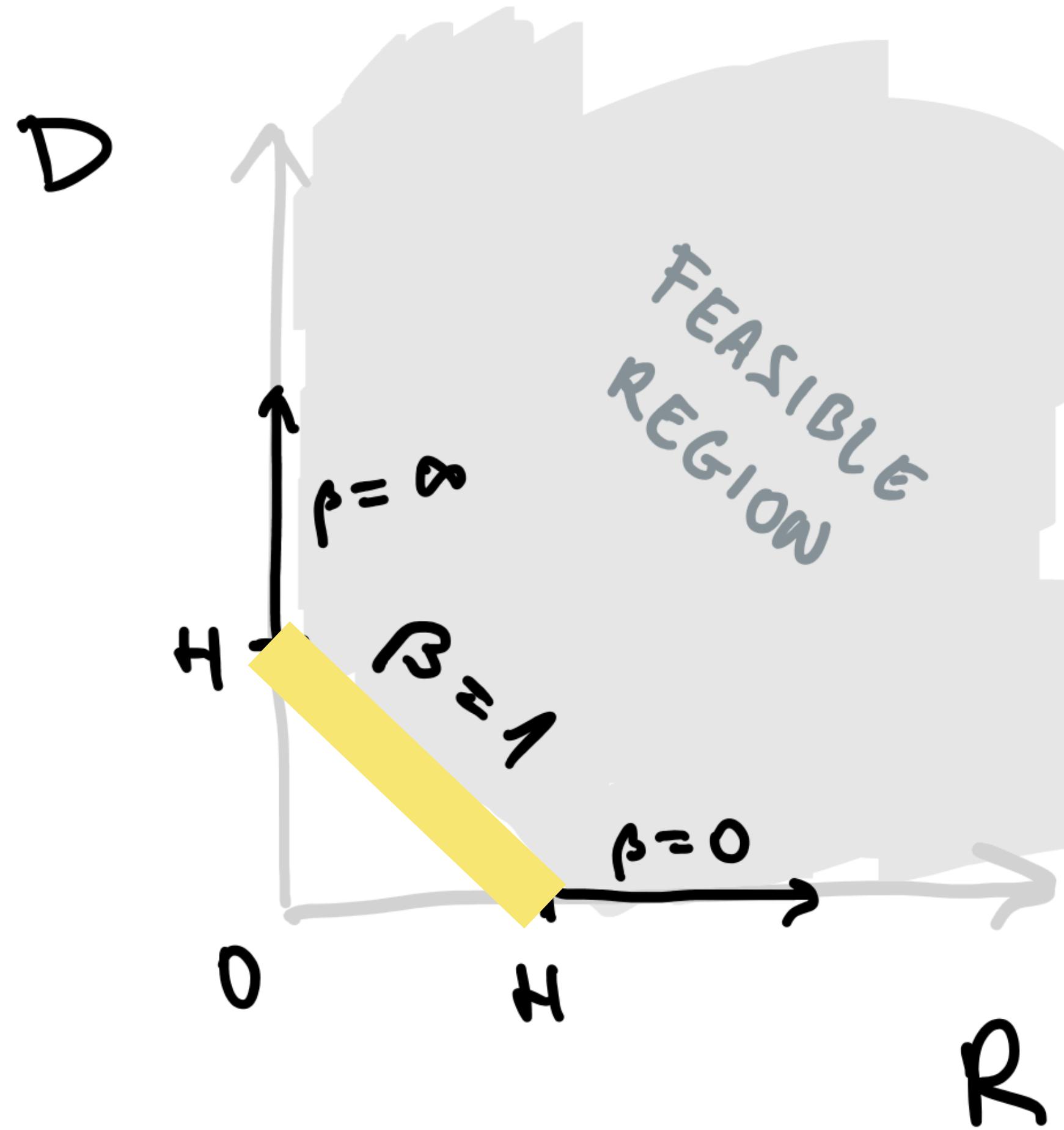
# Auto-decoding limit

- compression without retaining information
- $z$  independent of  $x$
- density estimation without representation learning

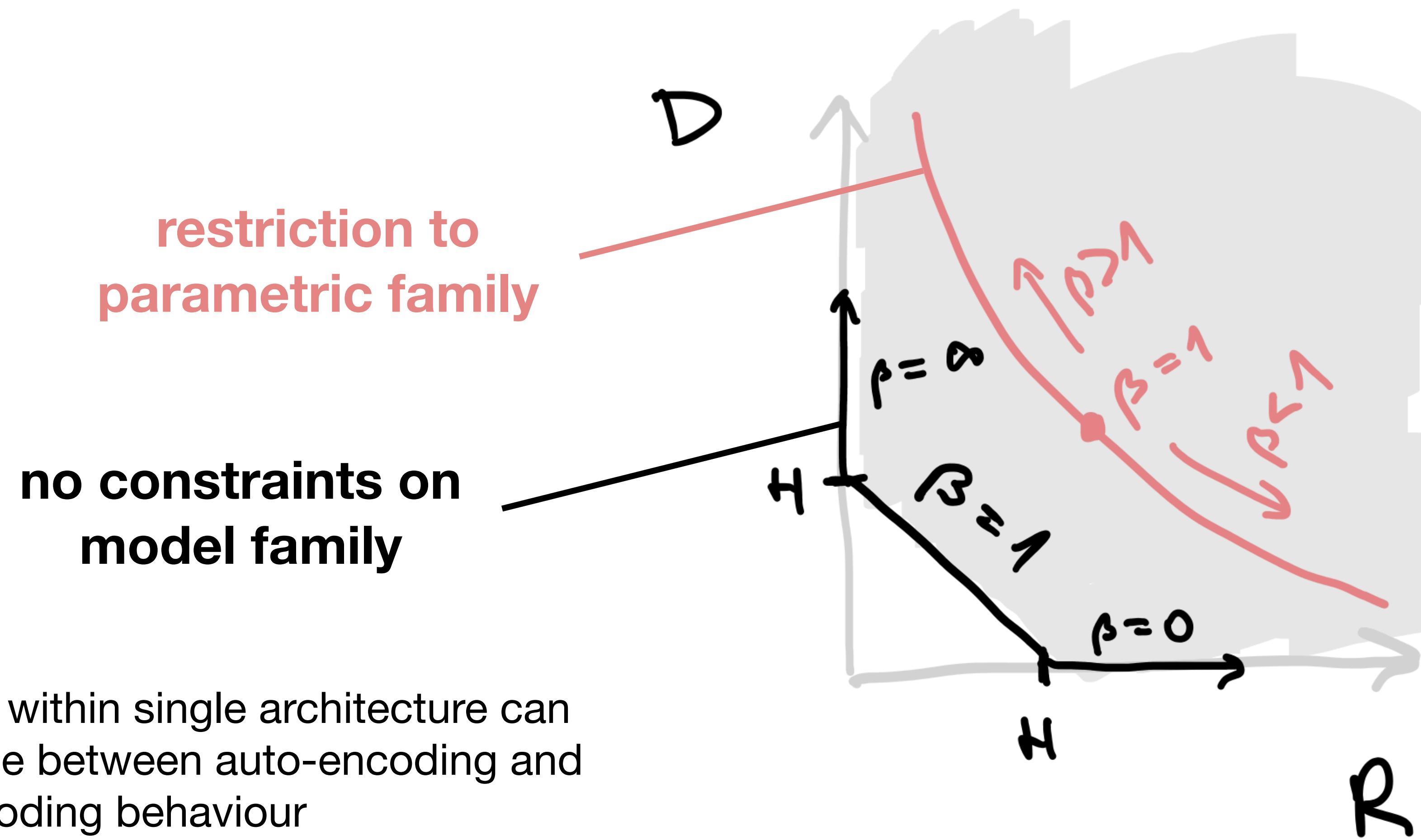


# Intermediate points

- intermediate trade-offs between rate and distortion
- these points can't be targeted through  $\beta$ -VAE objective since they all belong to  $\beta=1$
- point is selected implicitly through model architecture and initial conditions

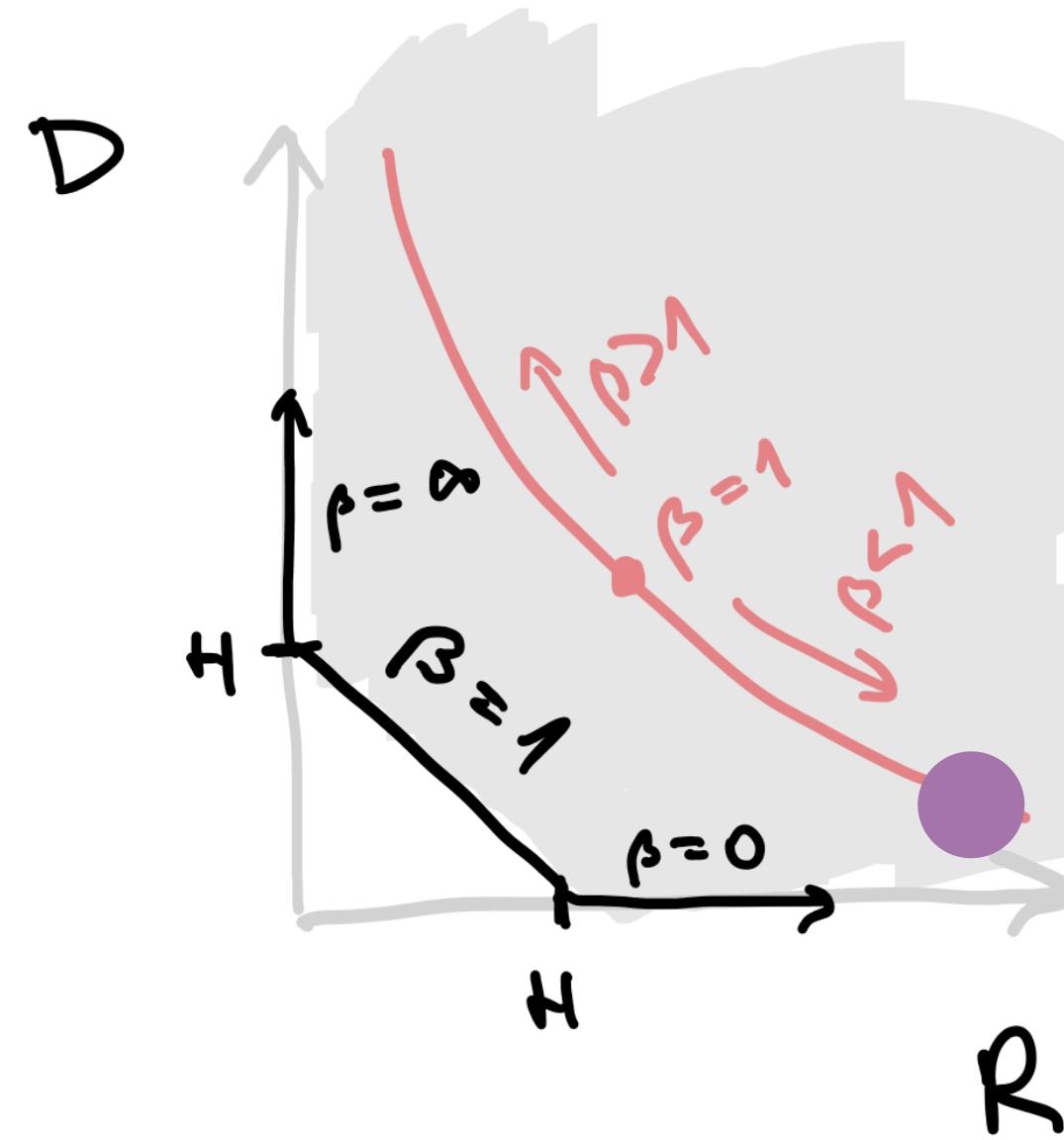
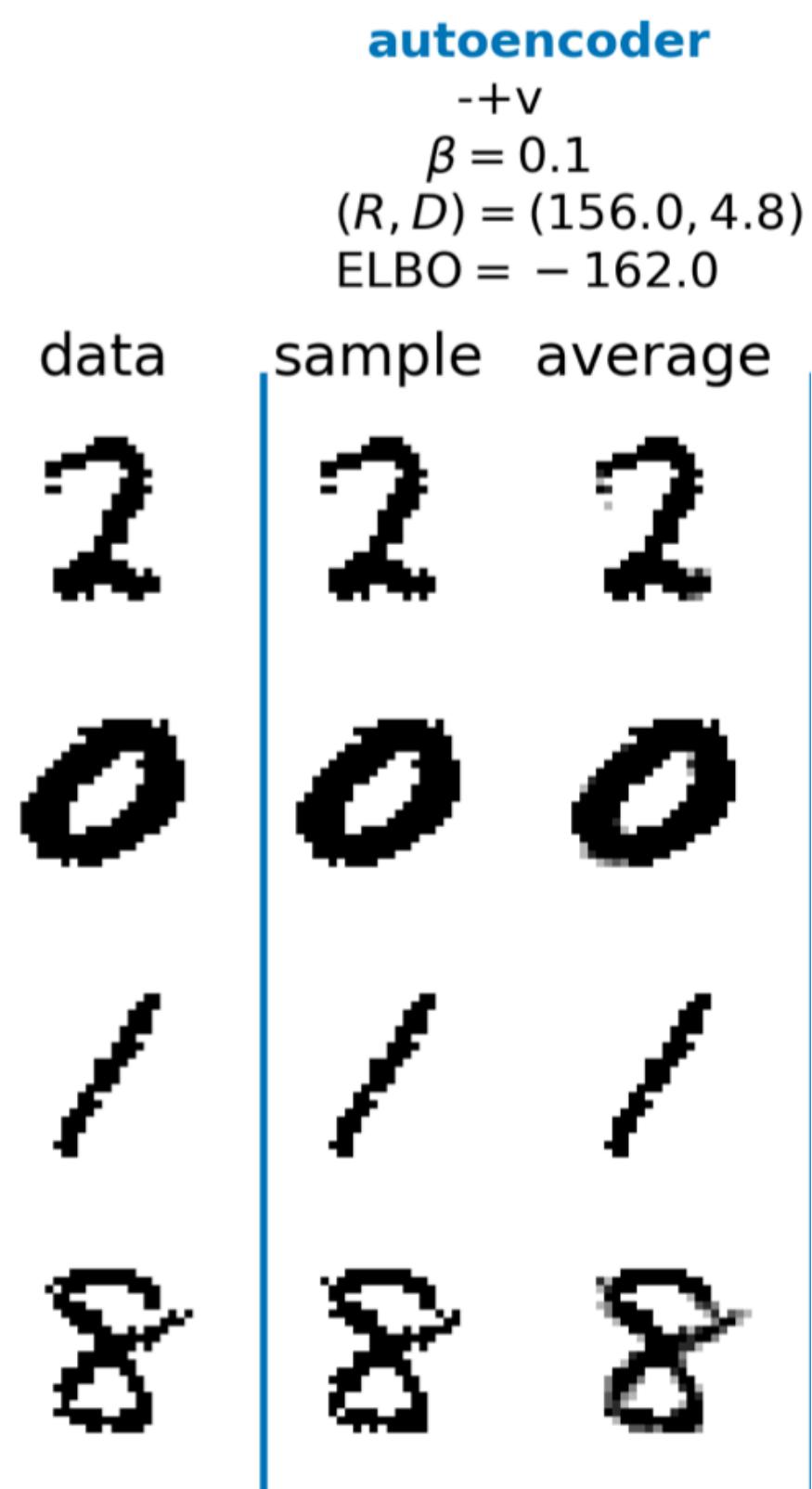


# Constrained model family



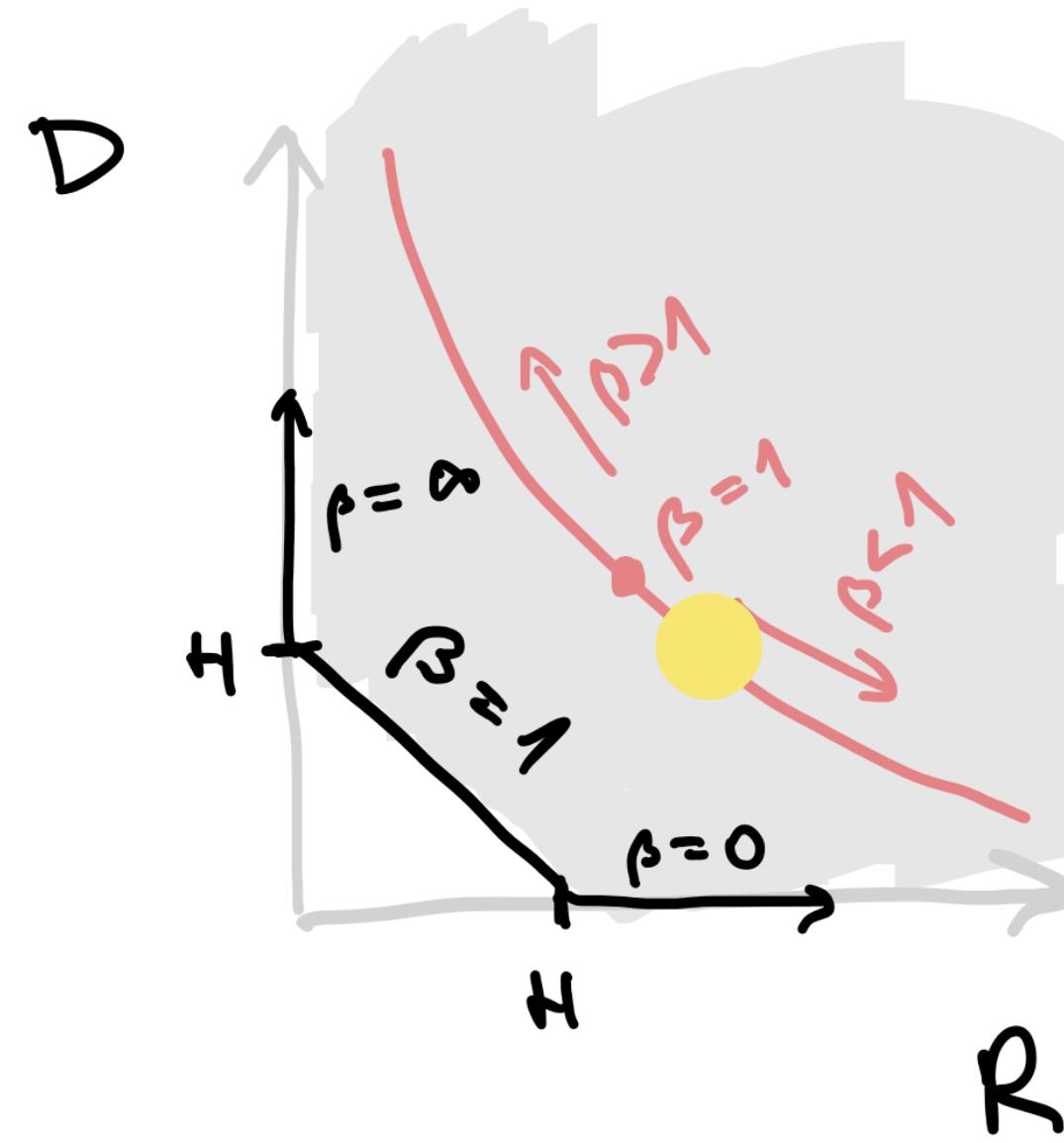
- varying  $\beta$  within single architecture can interpolate between auto-encoding and auto-decoding behaviour

# Reconstruction

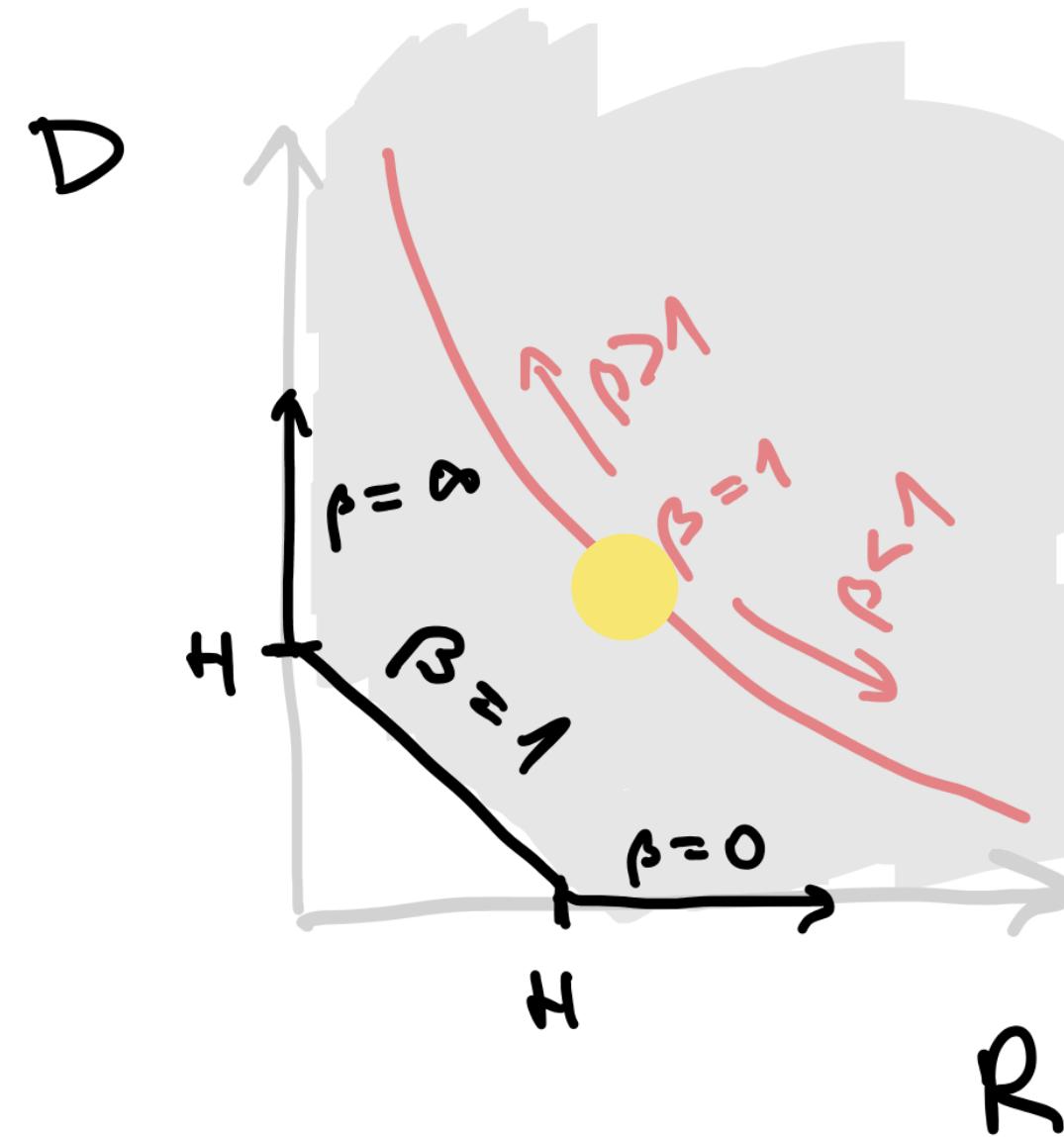
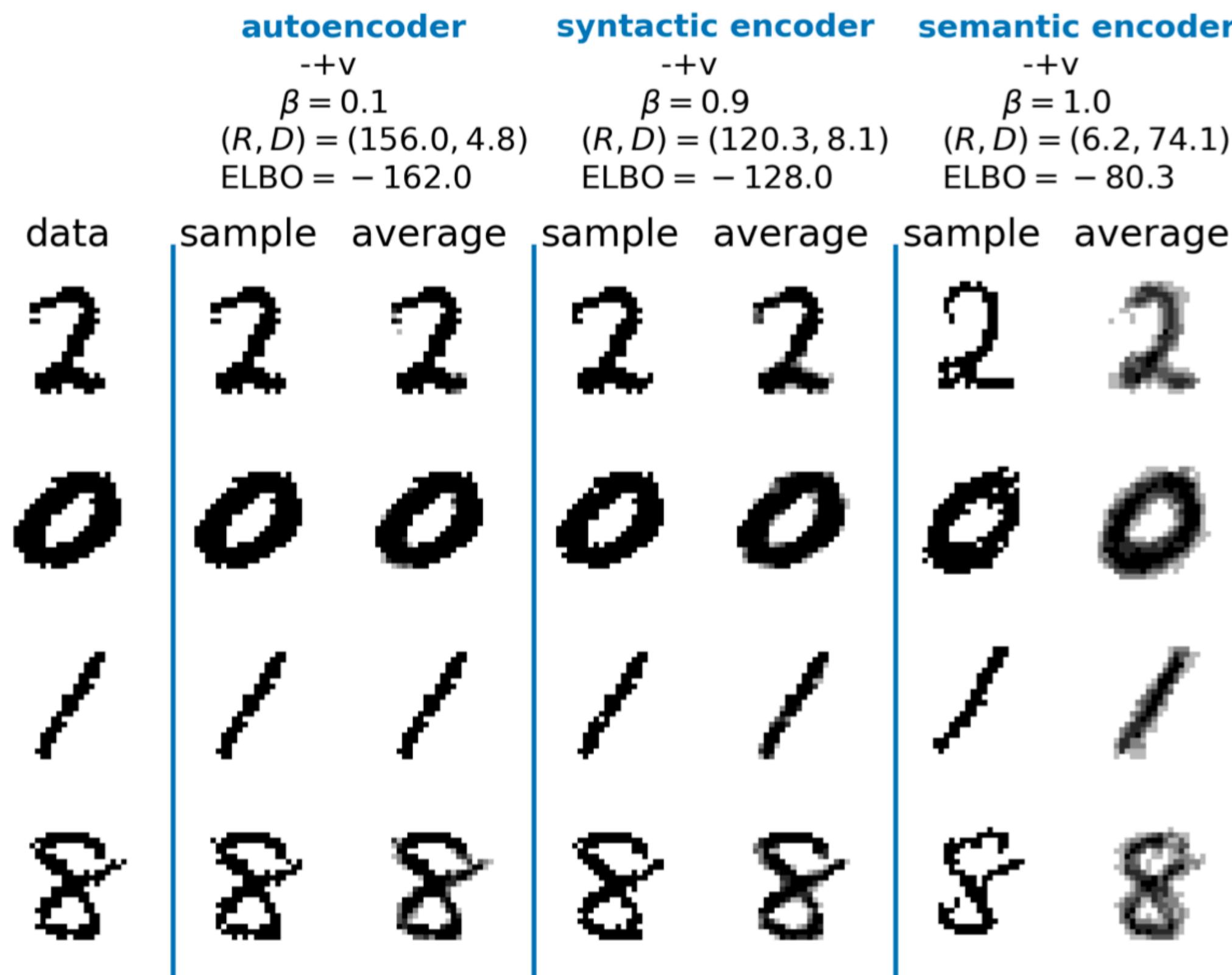


# Reconstruction

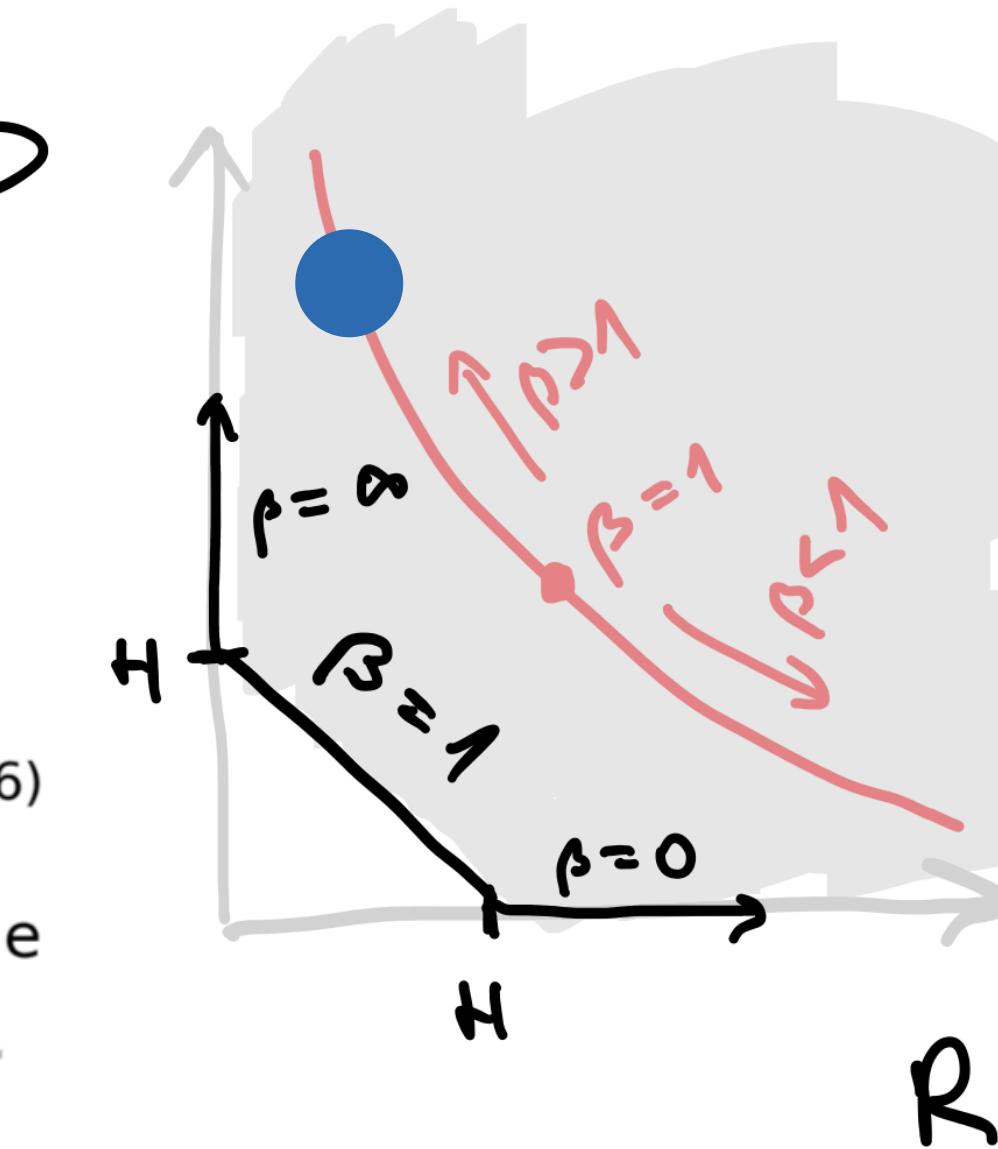
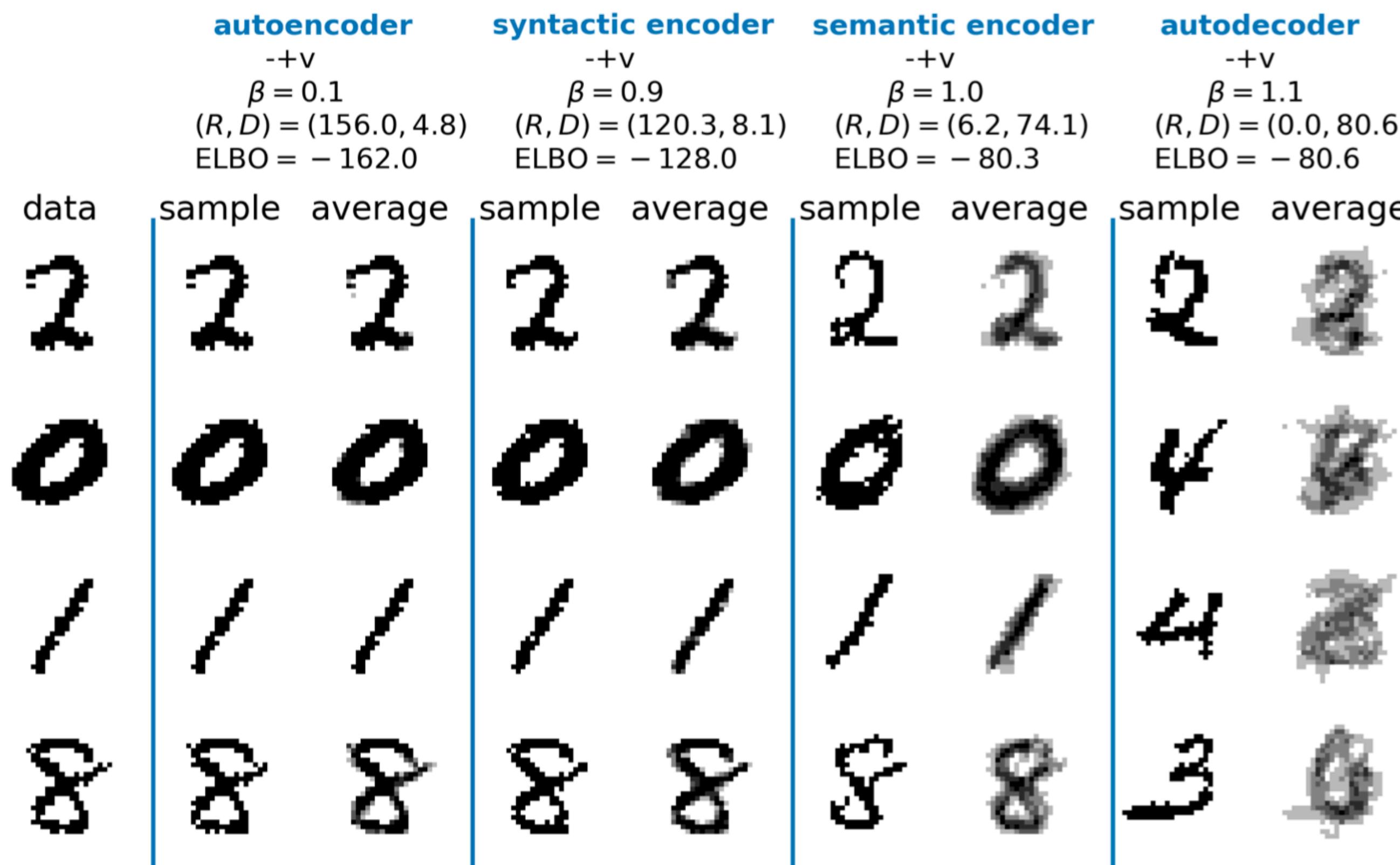
data	autoencoder		syntactic encoder	
	-+v $\beta = 0.1$ $(R, D) = (156.0, 4.8)$ ELBO = -162.0	-+v $\beta = 0.9$ $(R, D) = (120.3, 8.1)$ ELBO = -128.0	sample	average
2	2 2	2 2		
0	0 0	0 0		
/	/ /	/ /		
8	8 8	8 8		



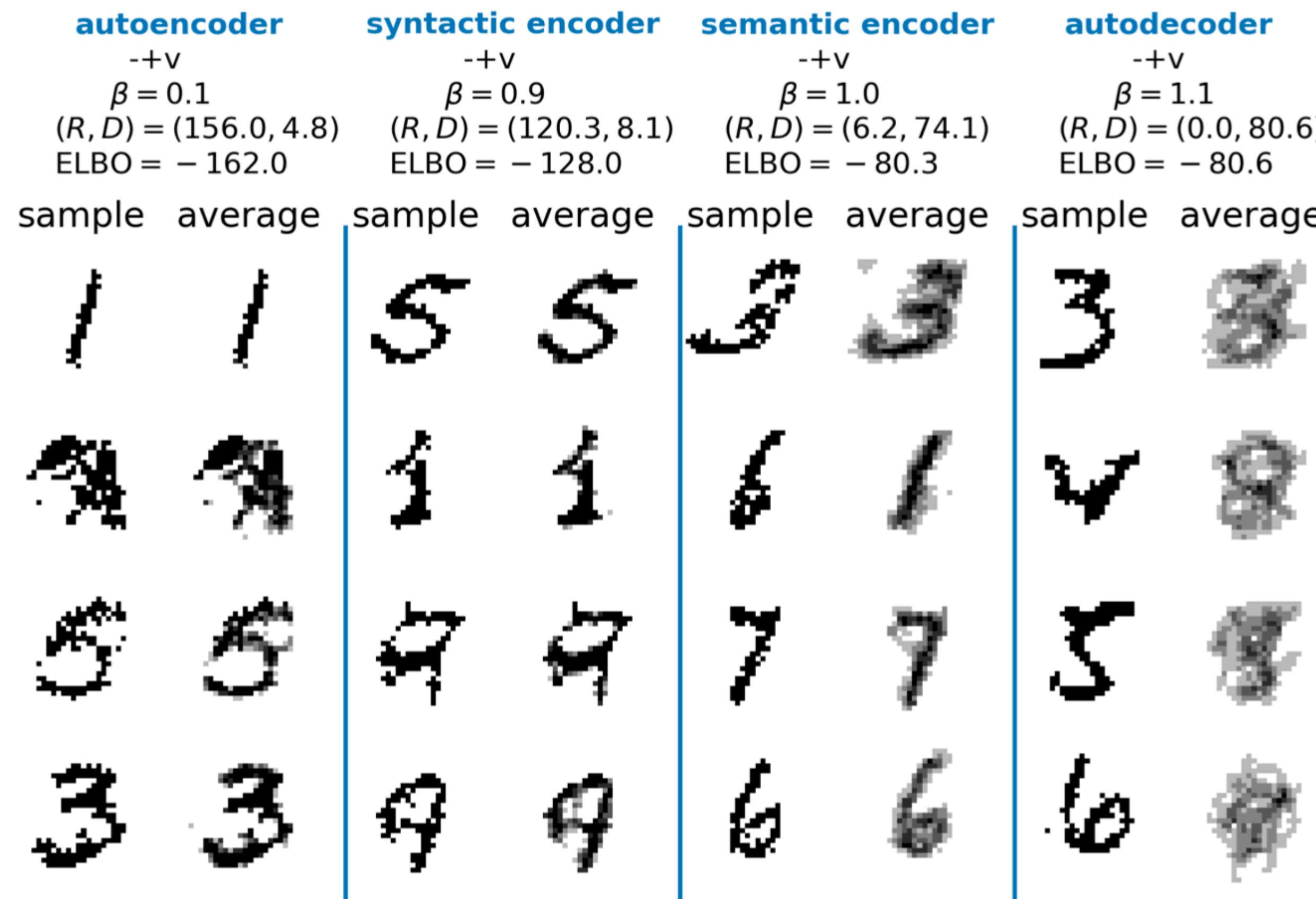
# Reconstruction



# Reconstruction

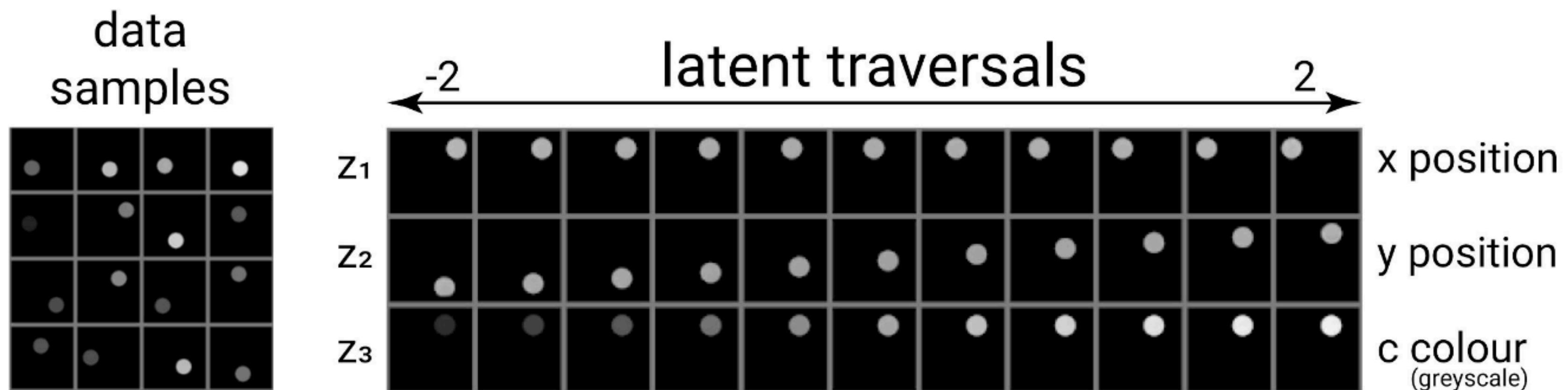


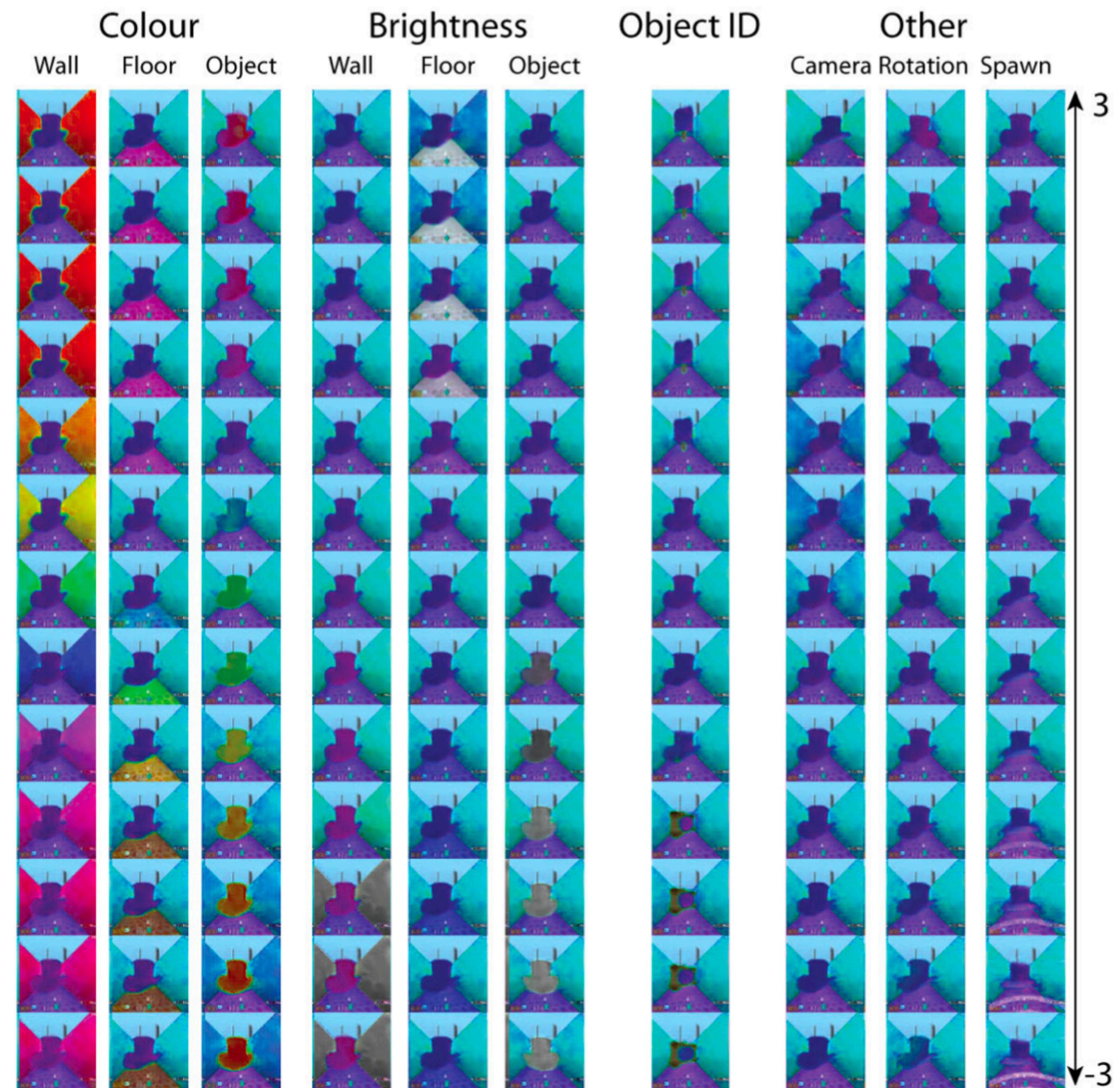
# Generation



# Disentangled representations

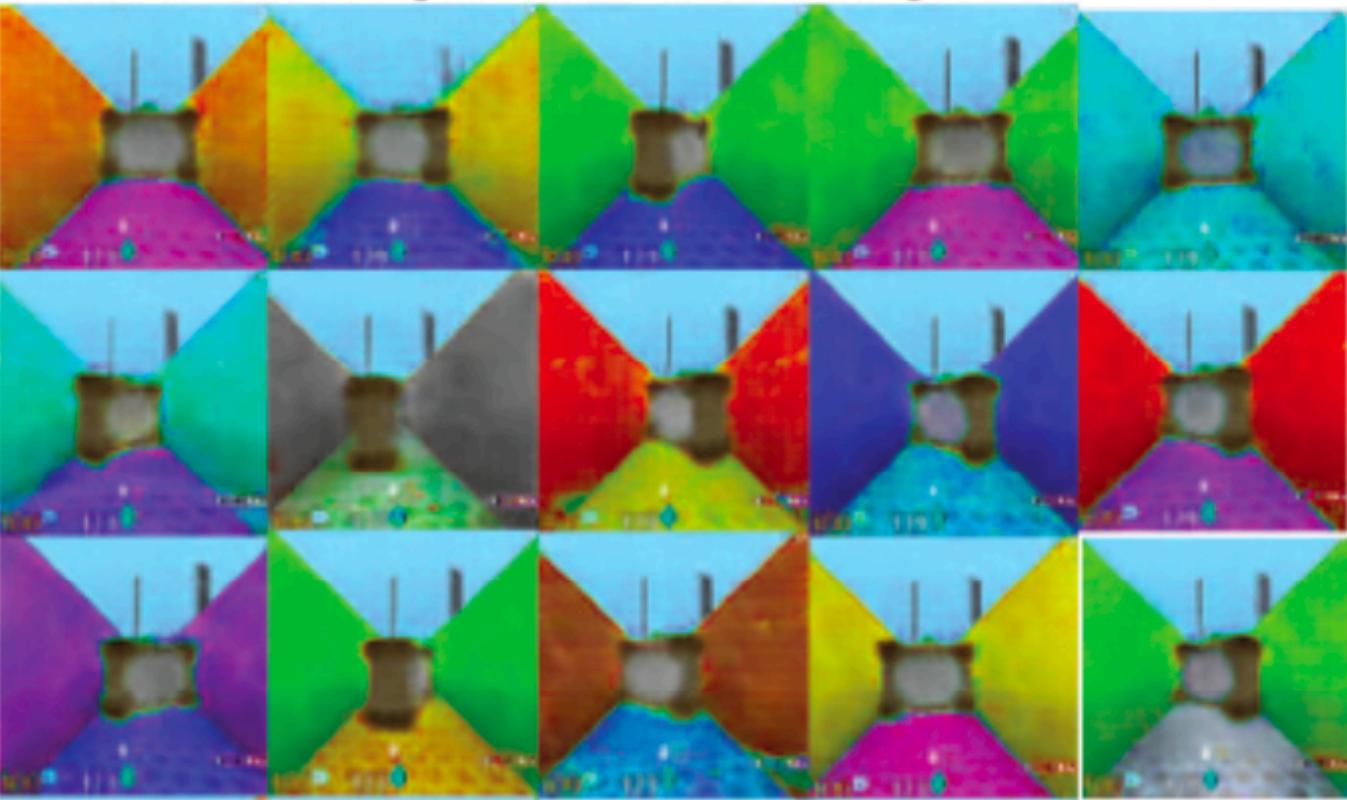
- we were able to control the information content of the latent representation with  $\beta$
  - we would also prefer representations that are **disentangled**
    - disentangled = factorised + interpretable (Bengio, 2013)
    - one generative factor per latent dimension
    - capture symmetry transformations



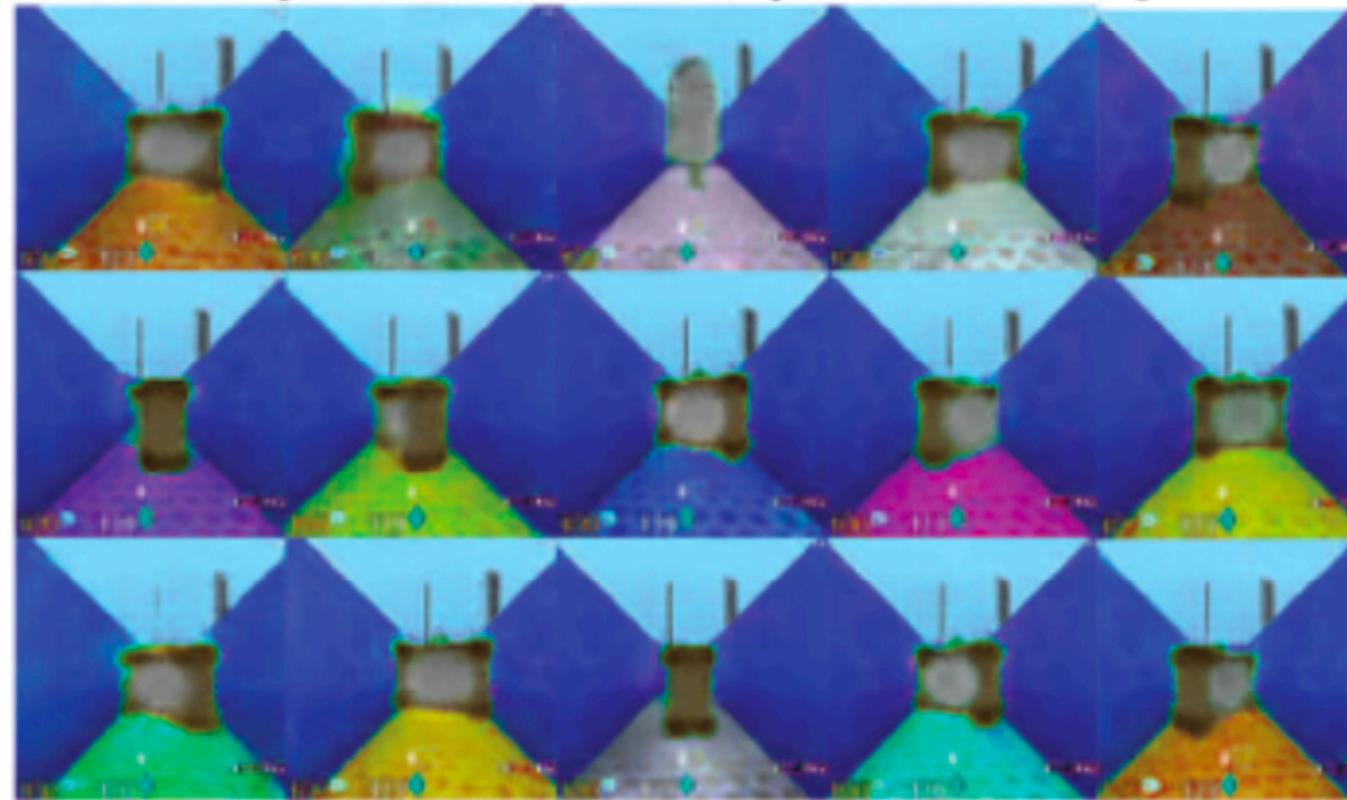


**A**

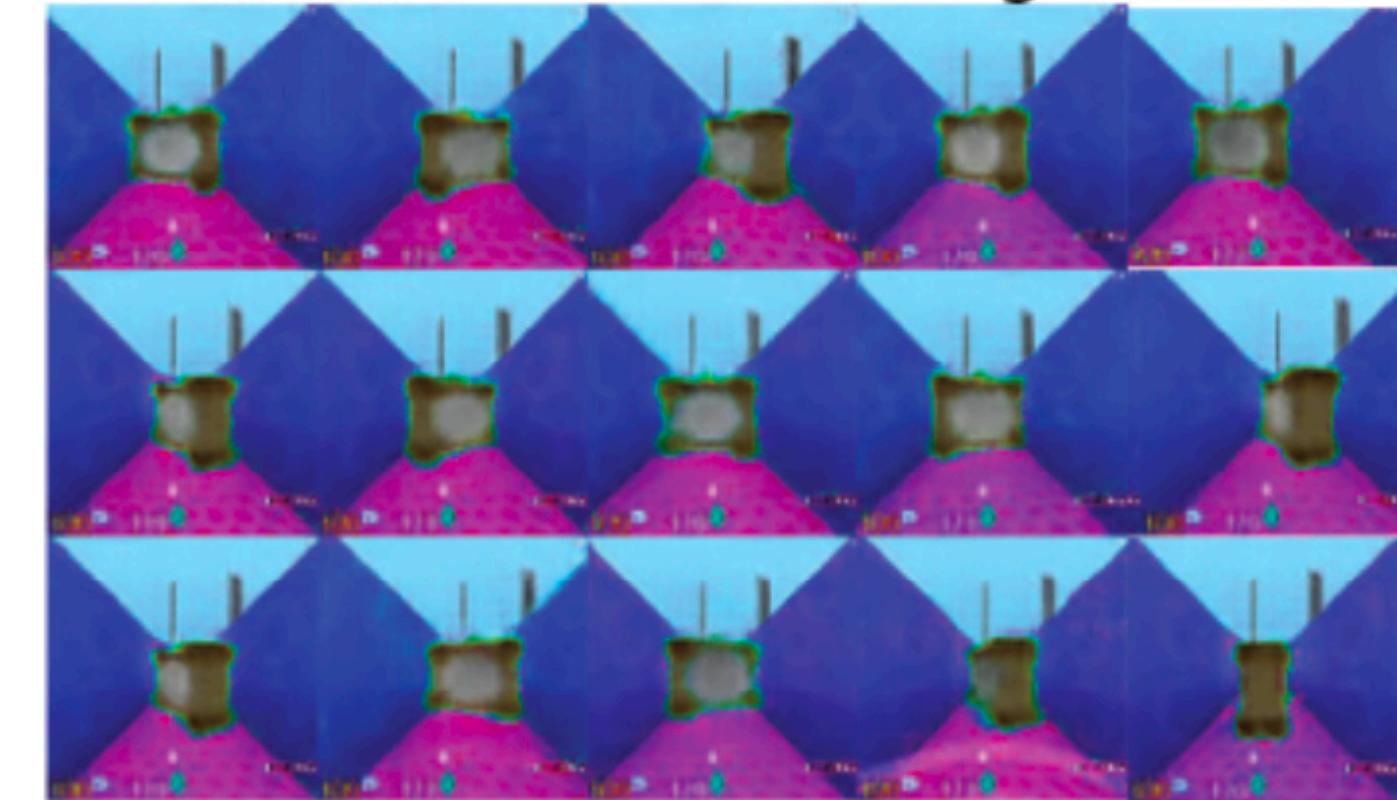
{white suitcase}



{white suitcase, blue wall}



{white suitcase, blue wall, magenta floor}

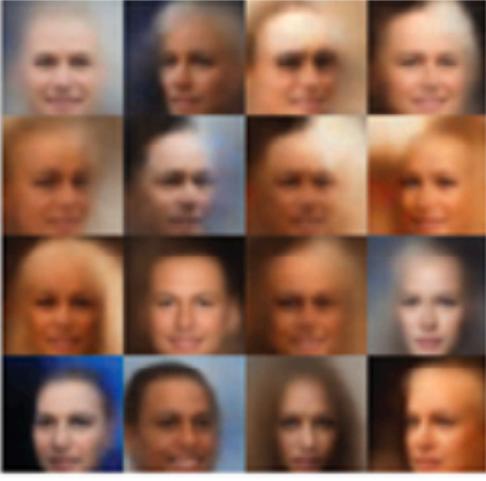


SCAN

{smiling}



{bald}



{bangs}



{male}



{pale skin}



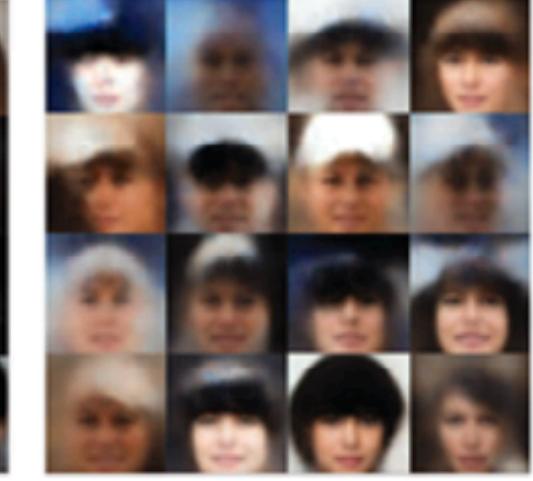
{eyeglasses}



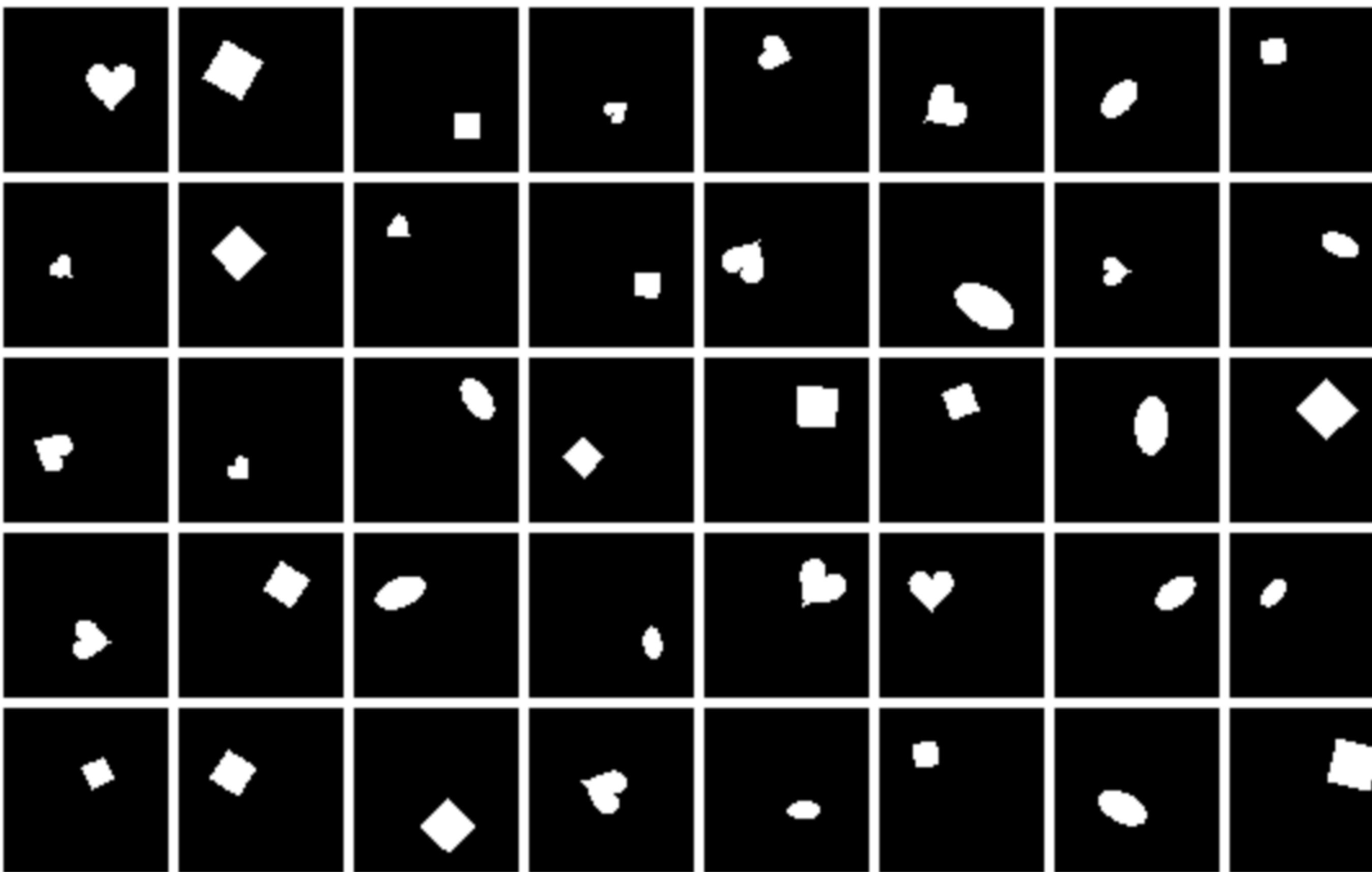
{black hair}



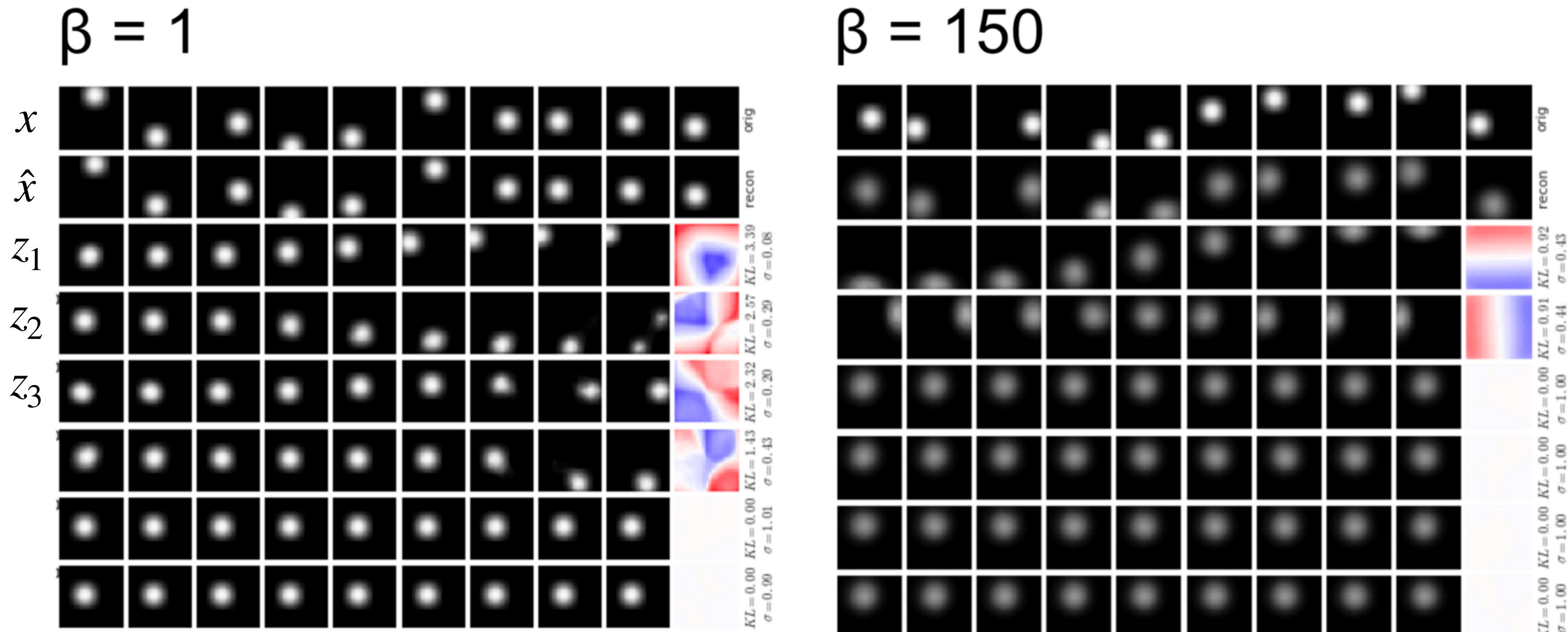
{hat}



# dsprites dataset

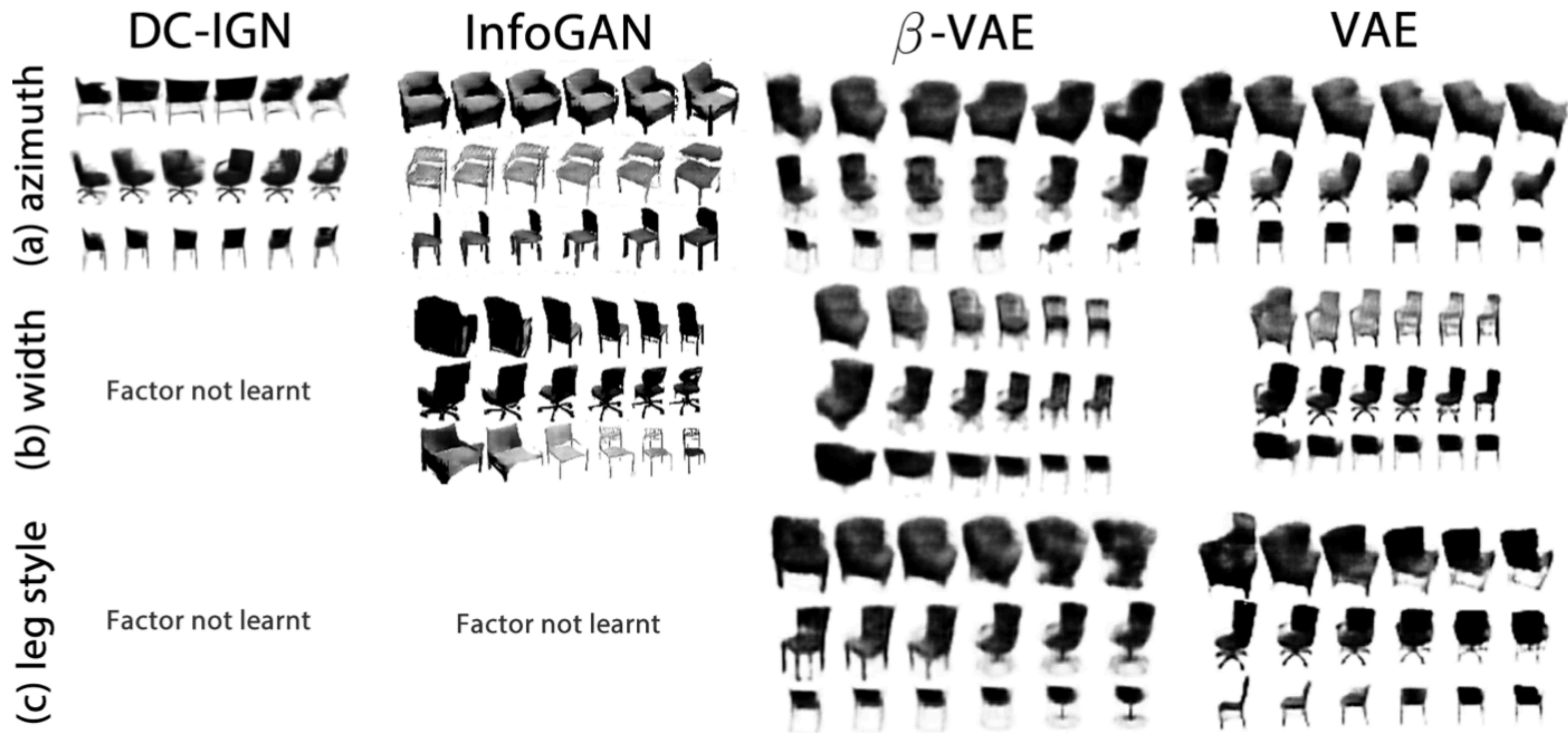


# $\beta$ -VAE



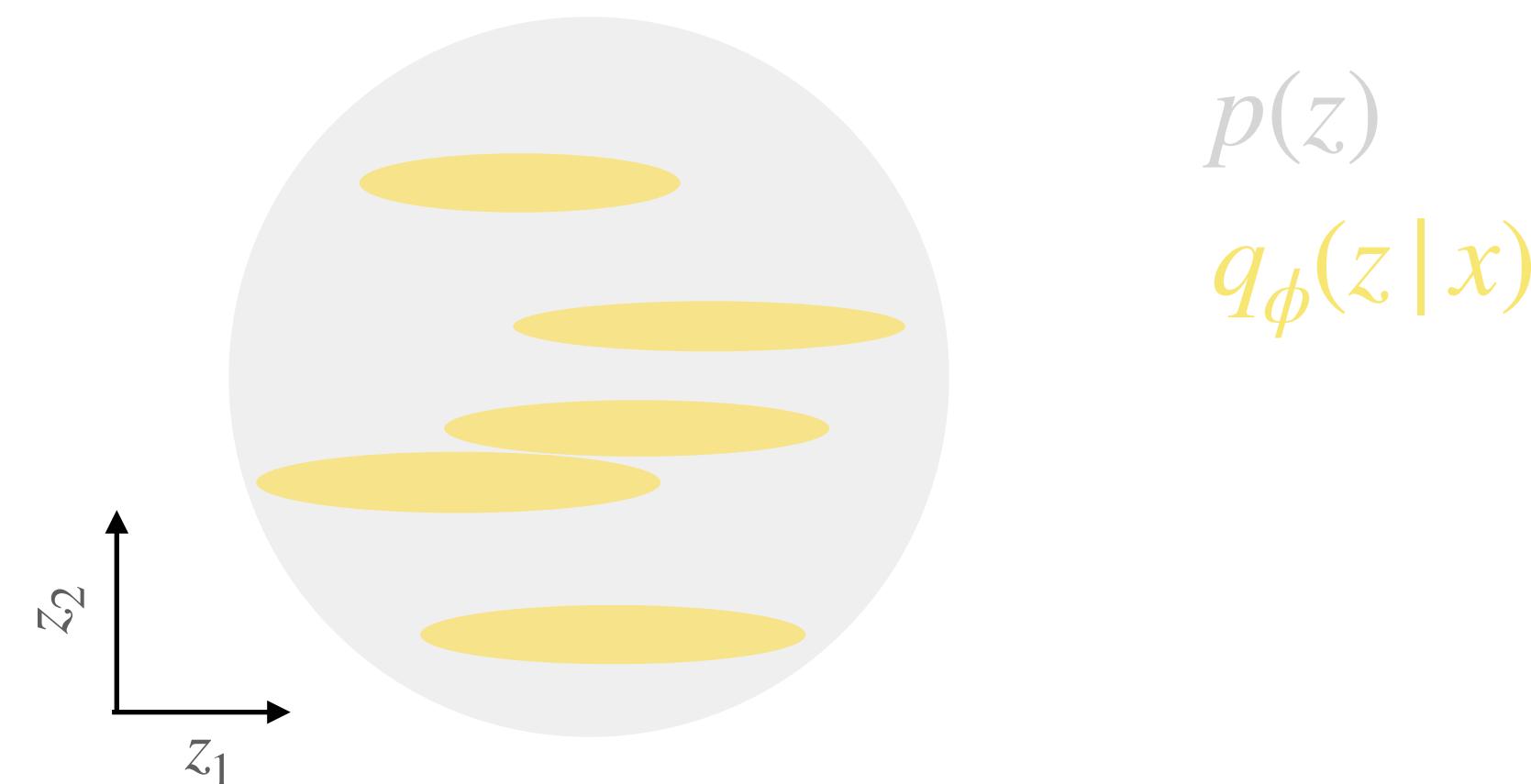
(Burgess et al. 2017)

# $\beta$ -VAE



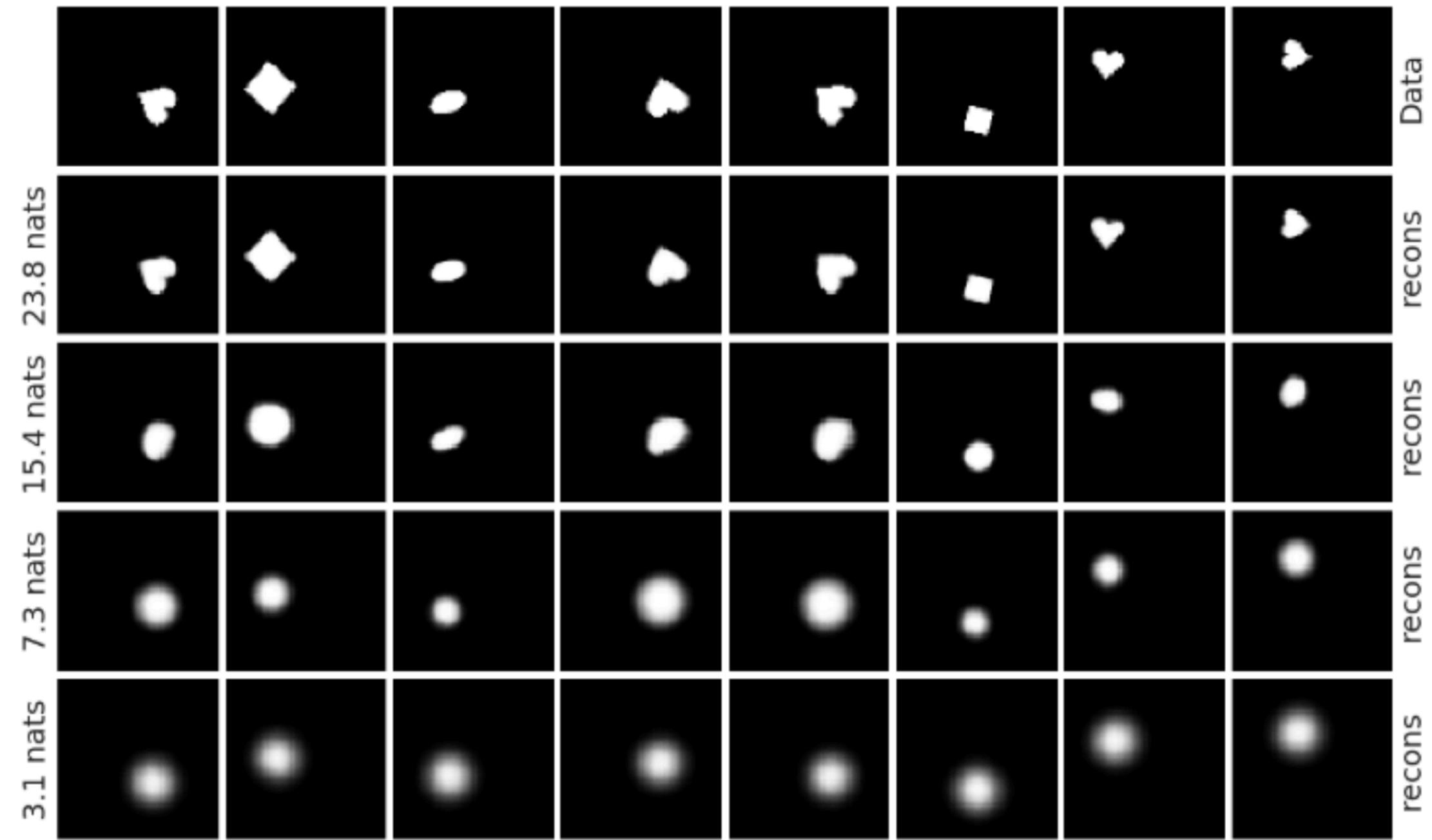
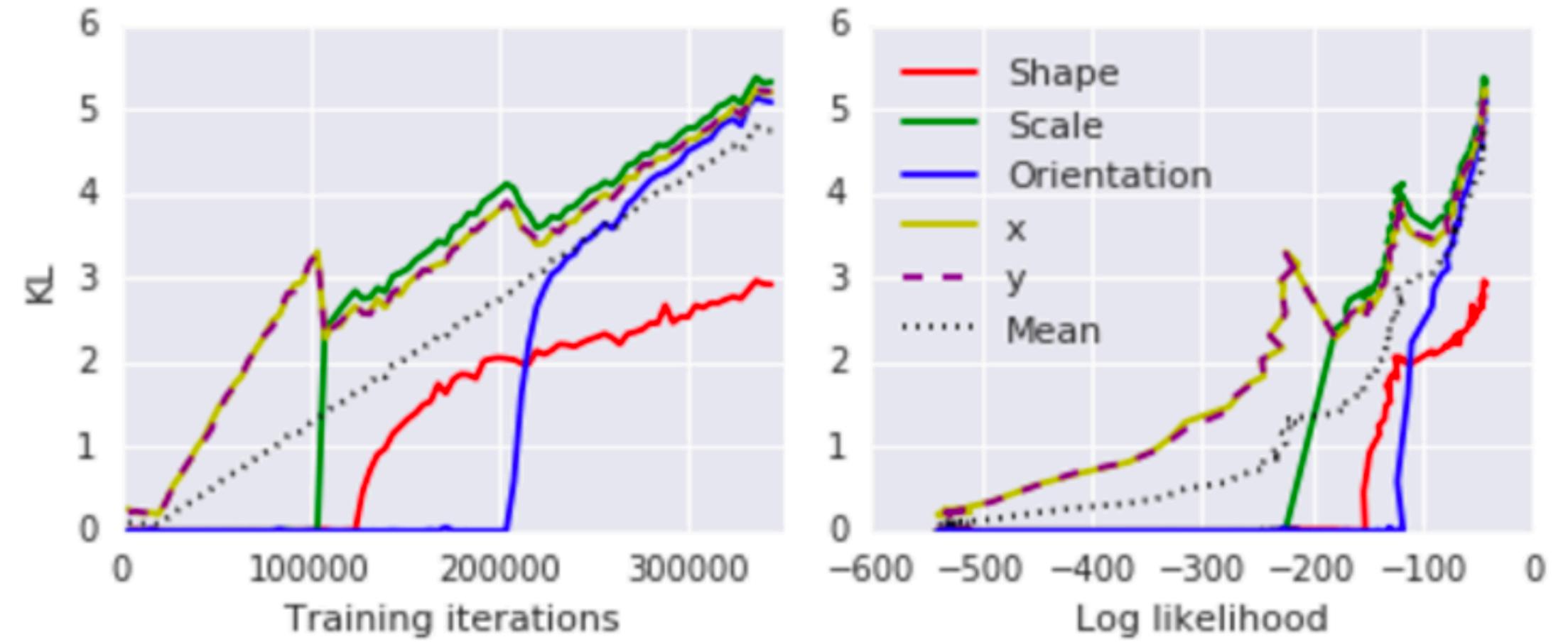
# Why does $\beta$ -VAE disentangle?

- mainly because of the diagonal normal approximate posterior
  - different generative factors give different average contributions to reconstruction
  - these factors should have different capacities allocated to them
  - with a diagonal posterior, this necessitates using different dimensions for each factor



# Why does $\beta$ -VAE disentangle?

- mainly because of the diagonal normal approximate posterior
  - different generative factors give different average contributions to reconstruction
  - these factors should have different capacities allocated to them
  - with a diagonal posterior, this necessitates using different dimensions for each factor
- Burgess et al., 2017. “**Understanding Disentangling in Beta-VAE.**” <http://arxiv.org/abs/1804.03599>
- Rolinek et al., 2019. “**Variational Autoencoders Pursue Pca Directions (by Accident).**” <https://doi.org/10.1109/CVPR.2019.01269>
- Zietlow, Dominik, Michal Rolínek, and Georg Martius. 2020. “**Demystifying Inductive Biases for  $\beta$ -VAE Based Architectures.**”



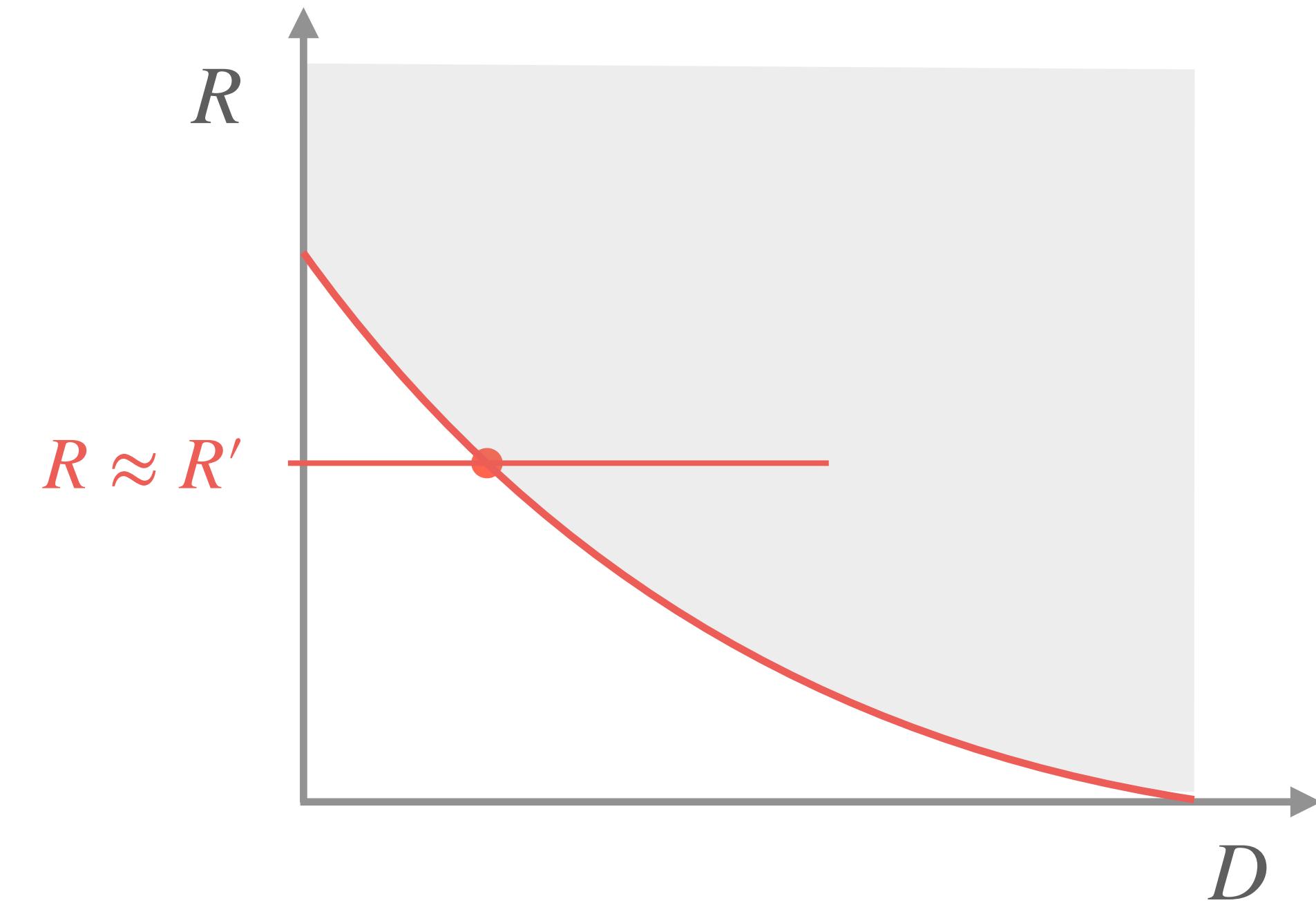
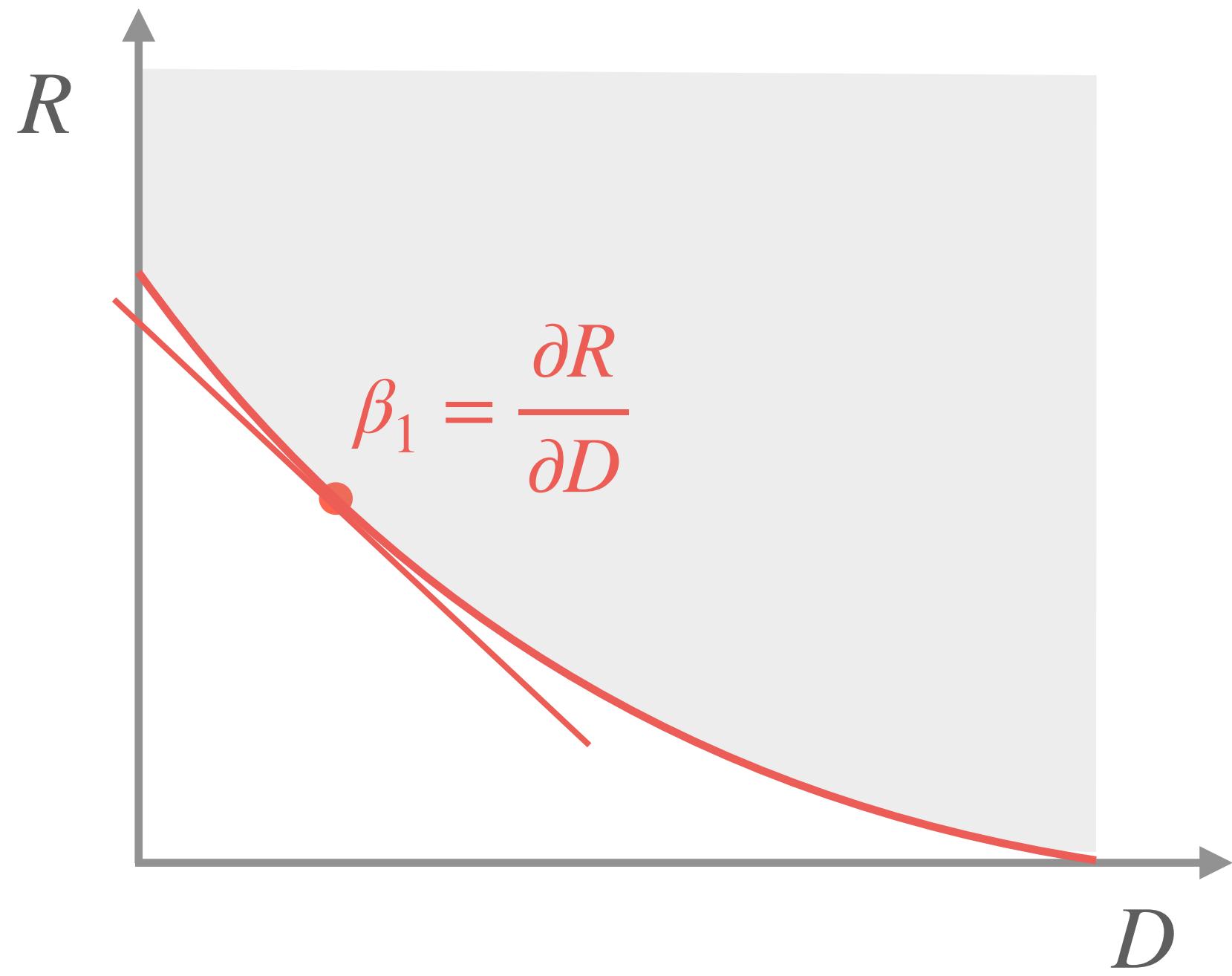
$$\mathcal{L}(\theta, \phi; \mathbf{x}(\mathbf{f}), \mathbf{z}, C) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{f})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \gamma |D_{KL}(q_\phi(\mathbf{z}|\mathbf{f}) \parallel p(\mathbf{z})) - C|$$

# Modified objective

$$\mathcal{L}(\theta, \phi, x) = \mathbb{E}_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - \beta \textcolor{red}{KL}[q_\phi(z|x) || p(z)]$$

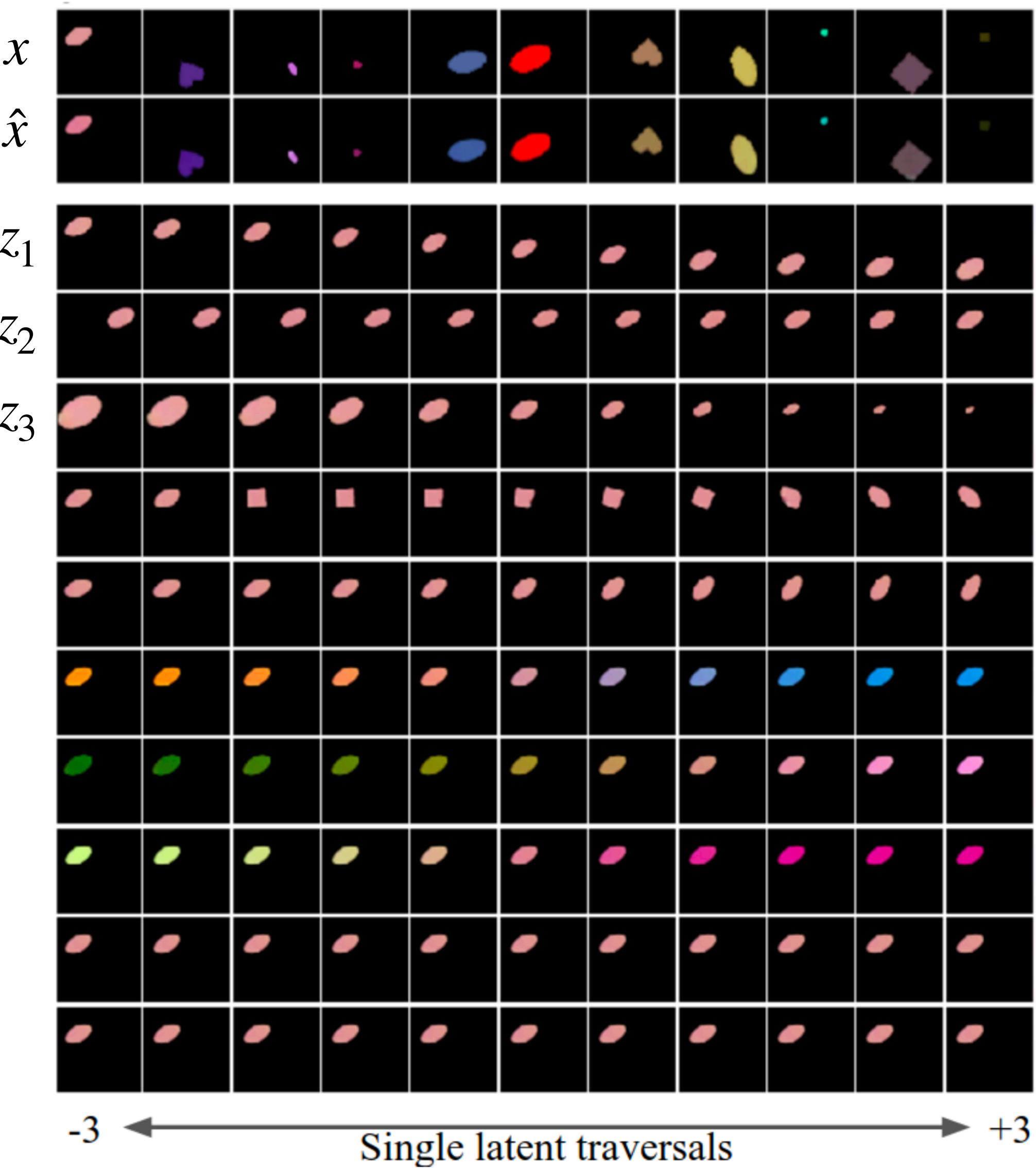
$$\mathcal{L}'(\theta, \phi, x, C) = \mathbb{E}_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - \gamma |KL[q_\phi(z|x) || p(z)] - C|$$

# Modified objective



$$\mathcal{L}(\theta, \phi, x) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \beta \text{KL}[q_\phi(z|x) || p(z)]$$

$$\mathcal{L}'(\theta, \phi, x, C) = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \gamma |\text{KL}[q_\phi(z|x) || p(z)] - C|$$



# Demo

<https://github.com/eemlcommunity/PracticalSessions2021/tree/main/generative>