

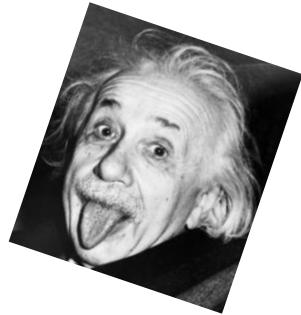
(some)
Generative Models
(which are not necessarily GANs)

December 23rd, 2021

Niv Haim

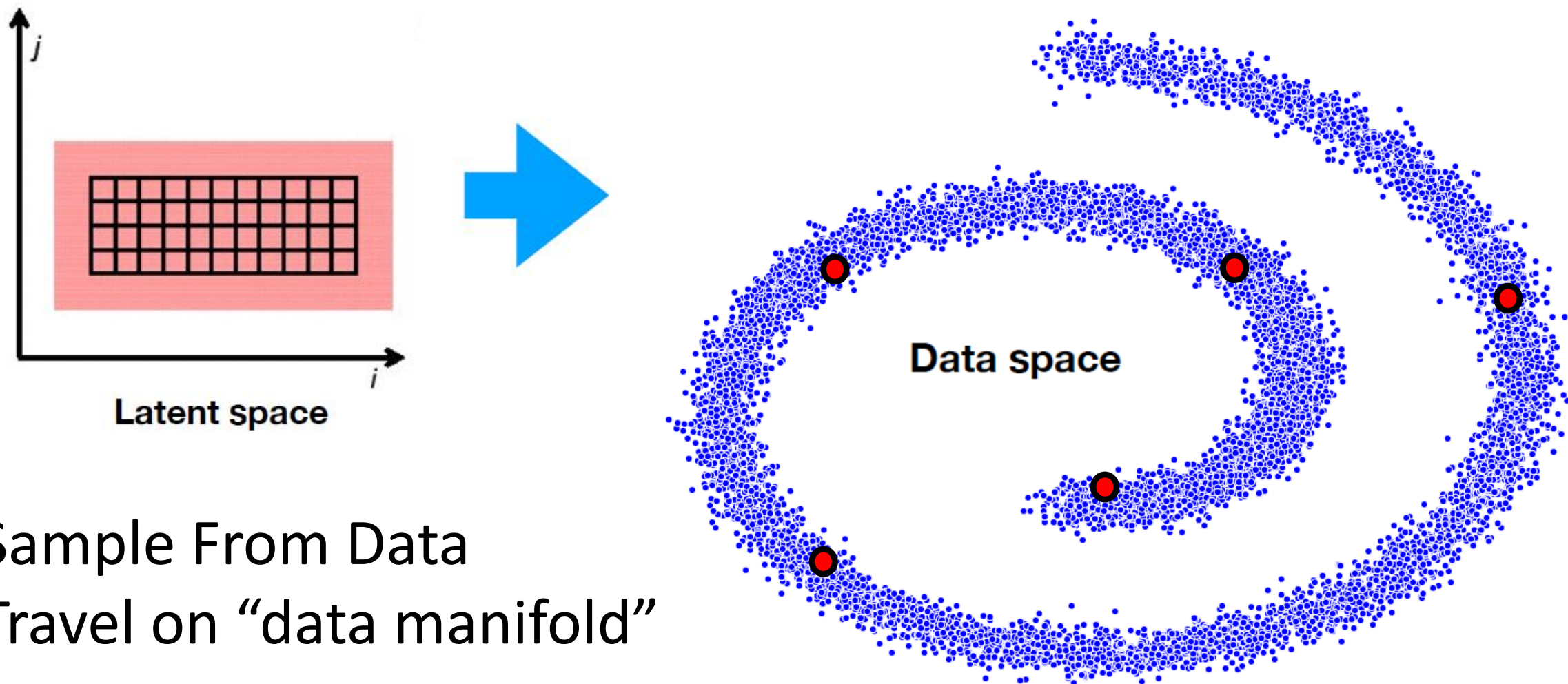
Outline

- Autoencoders (AE)
- Variational Autoencoders (VAE)
- Vector Quantized VAE (VQ-VAE)
- Diffusion Models



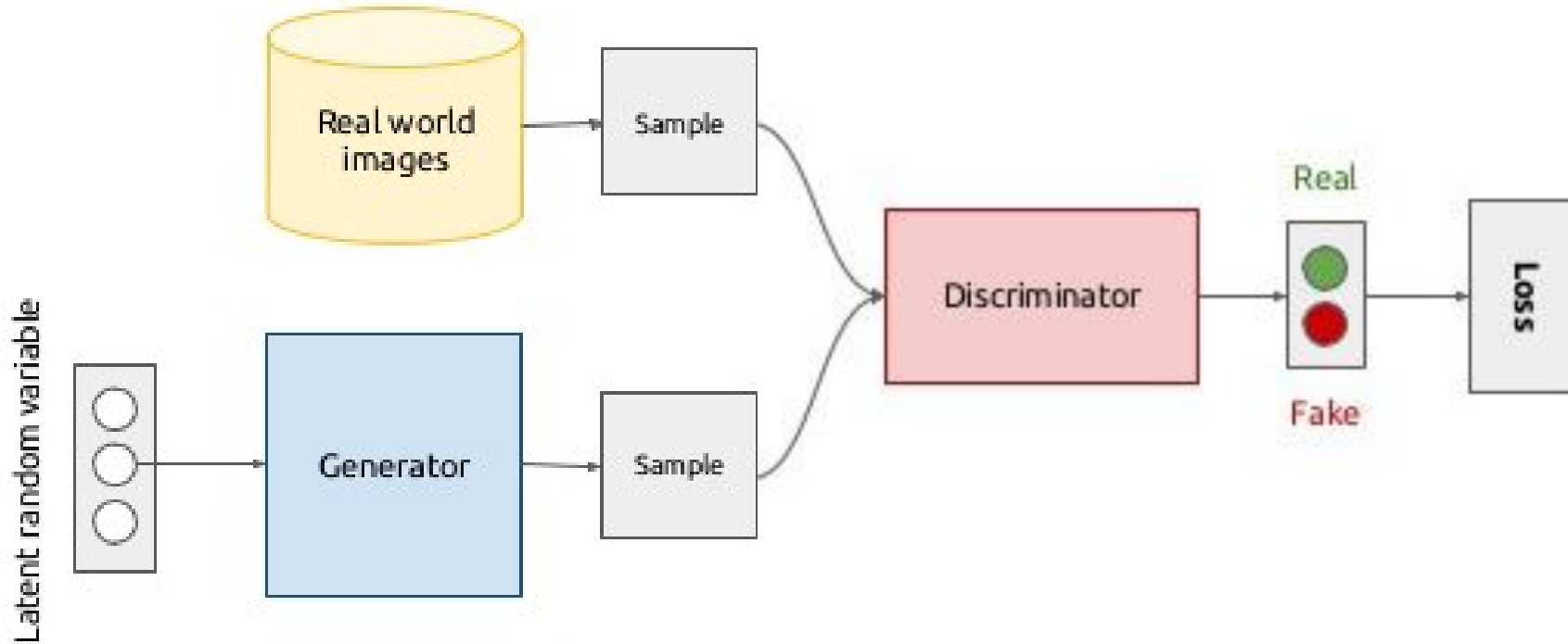
Objective - Reminder

Goal: map simple (known) distribution to the data distribution



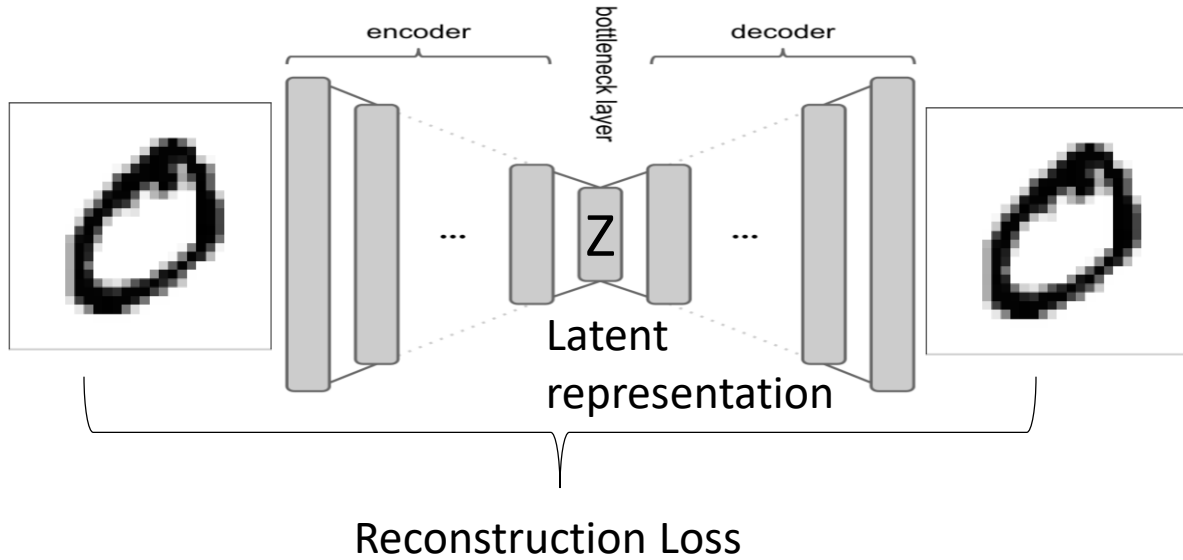
Sample From Data
Travel on “data manifold”

GANs - Reminder

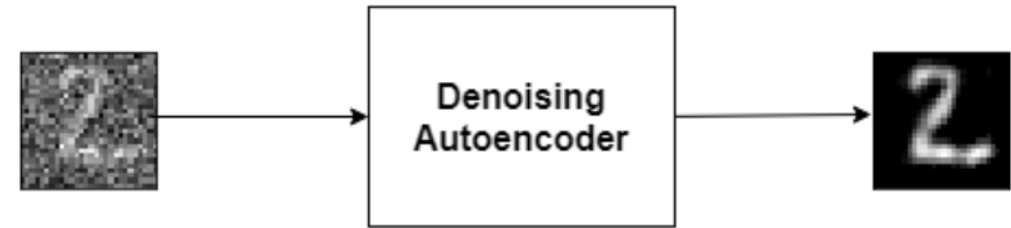


$$\mathcal{L}_{GAN} = \min_G \max_D \mathbb{E}_{x \sim p_{data}} [\log D(x)] + \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))]$$

Autoencoders



Dimensionality Reduction Denoising Autoencoder



source: <https://theailearner.com/2018/11/10/denoising-autoencoders>

“Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion” [Vincent et al, 2010]

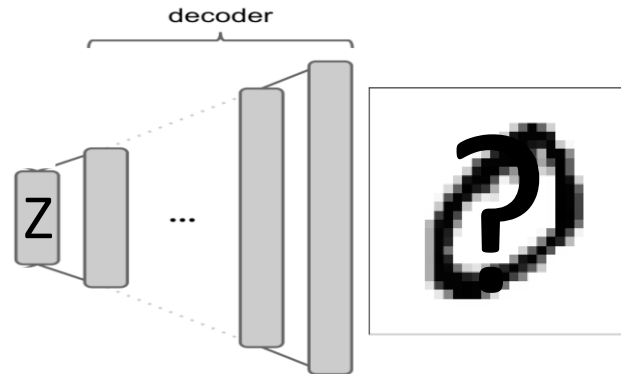
What happens in Latent Space?

- Generative model?

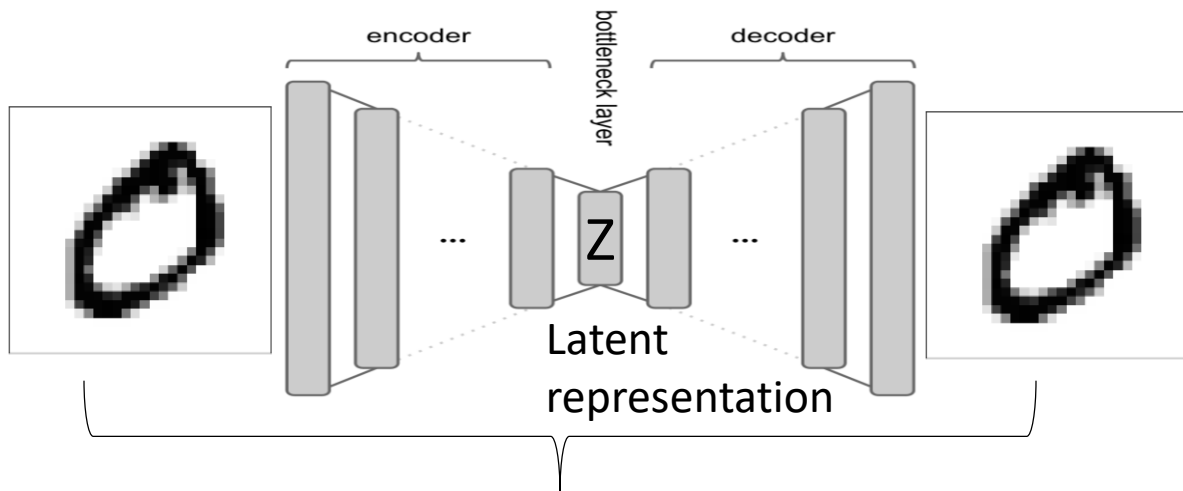
$$Z_1 = E(\text{img1})$$

$$Z_3 = Z_1 + \frac{1}{2} Z_2$$

$$Z_2 = E(\text{img2})$$



Autoencoders



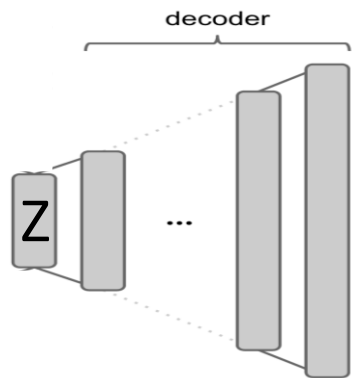
Reconstruction Loss

- Generative model?

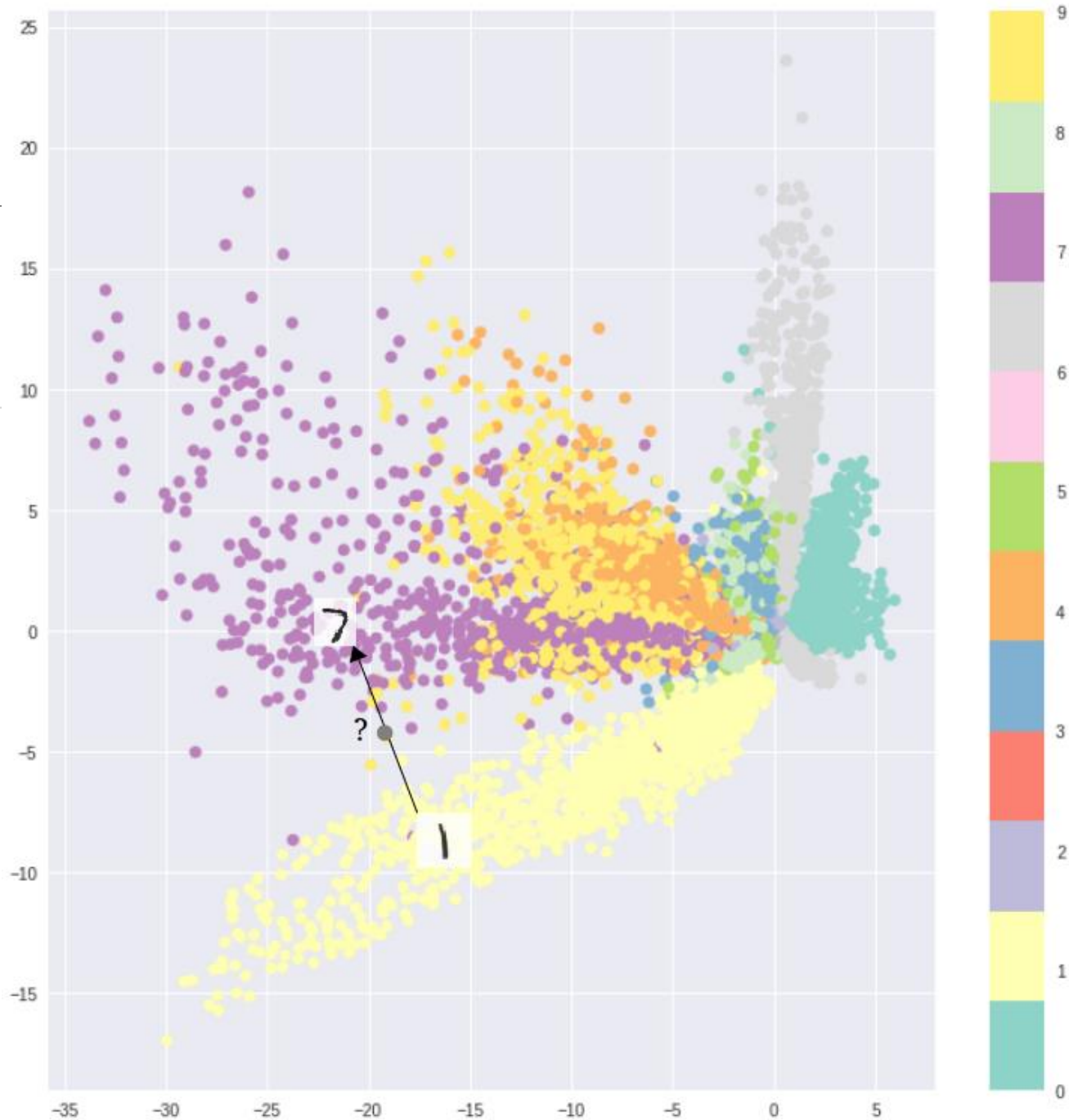
$$Z_1 = E(\text{img1})$$

$$Z_3 = Z_1 + \frac{1}{2} Z_2$$

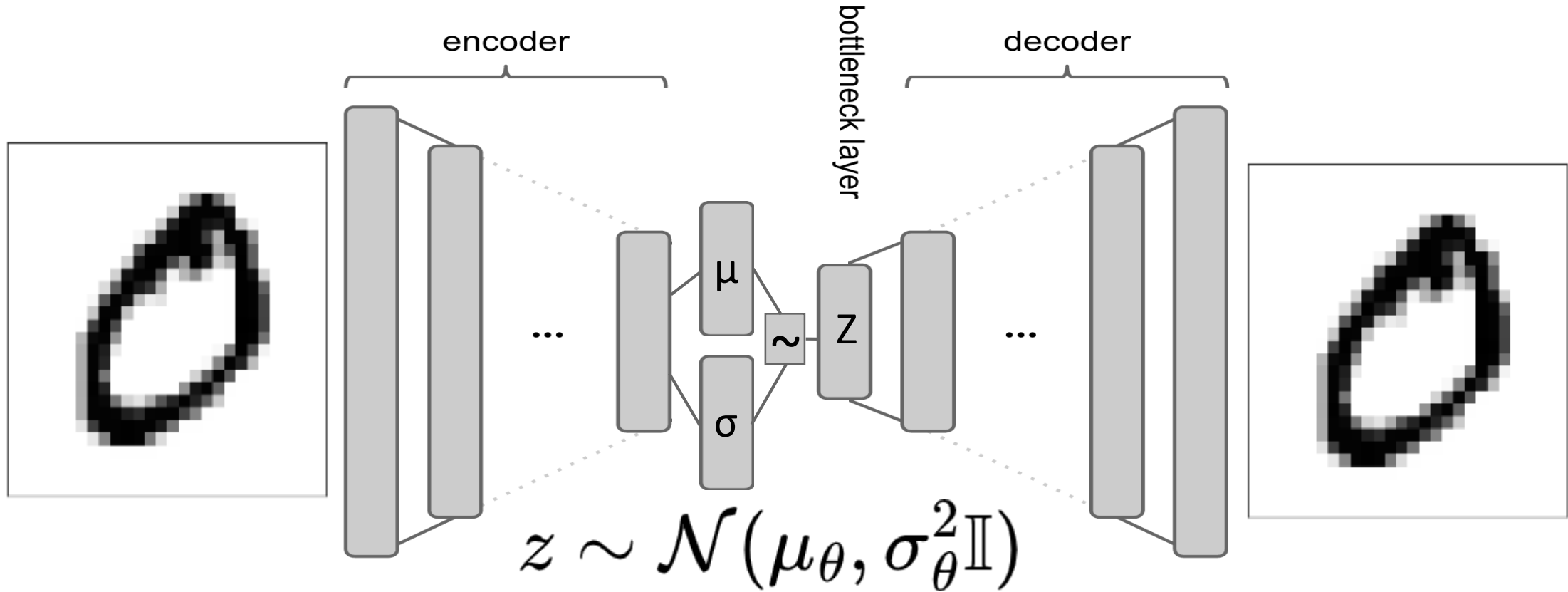
$$Z_2 = E(\text{img2})$$



?



Variational Autoencoders (VAE)

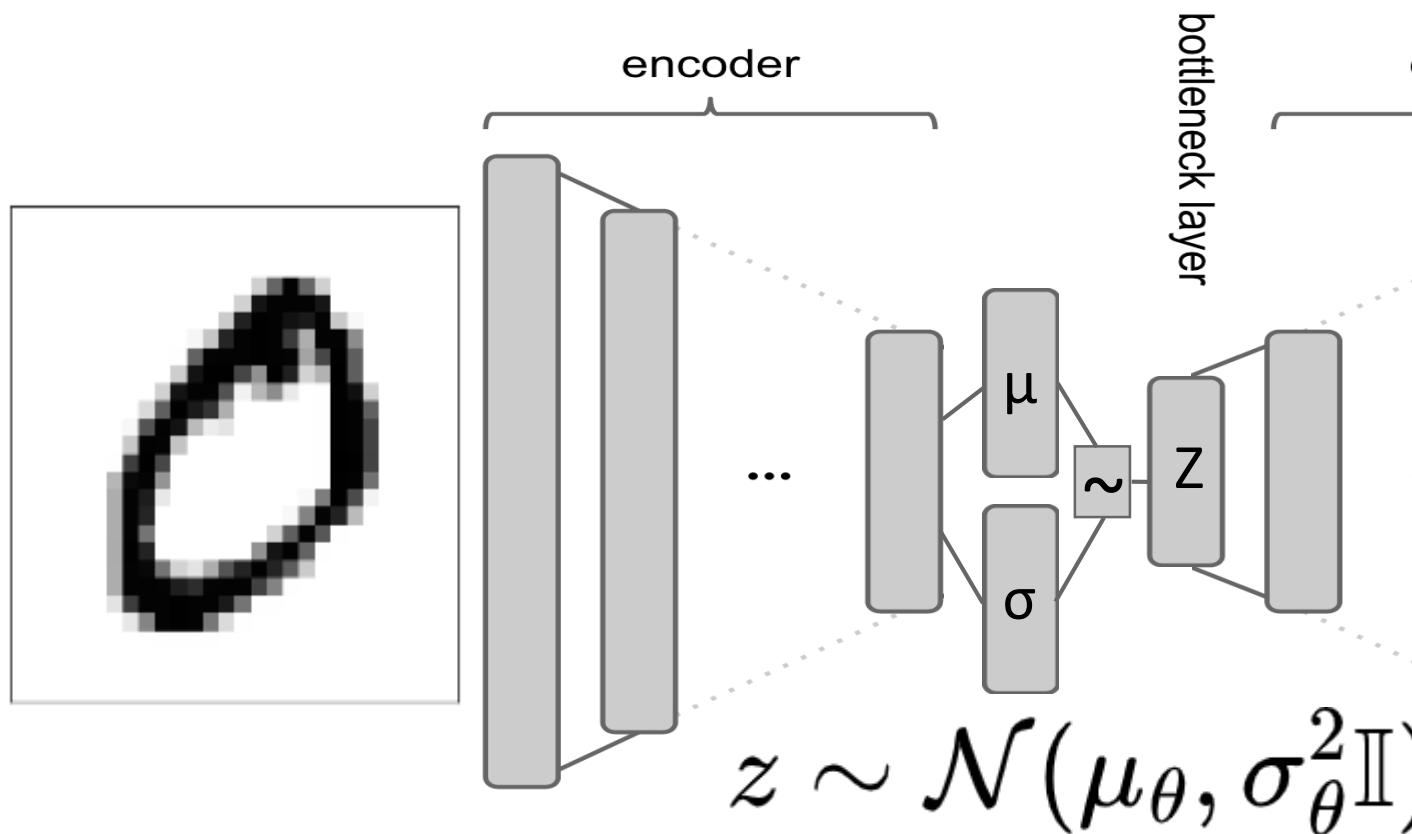


“Reparametrization Trick”:

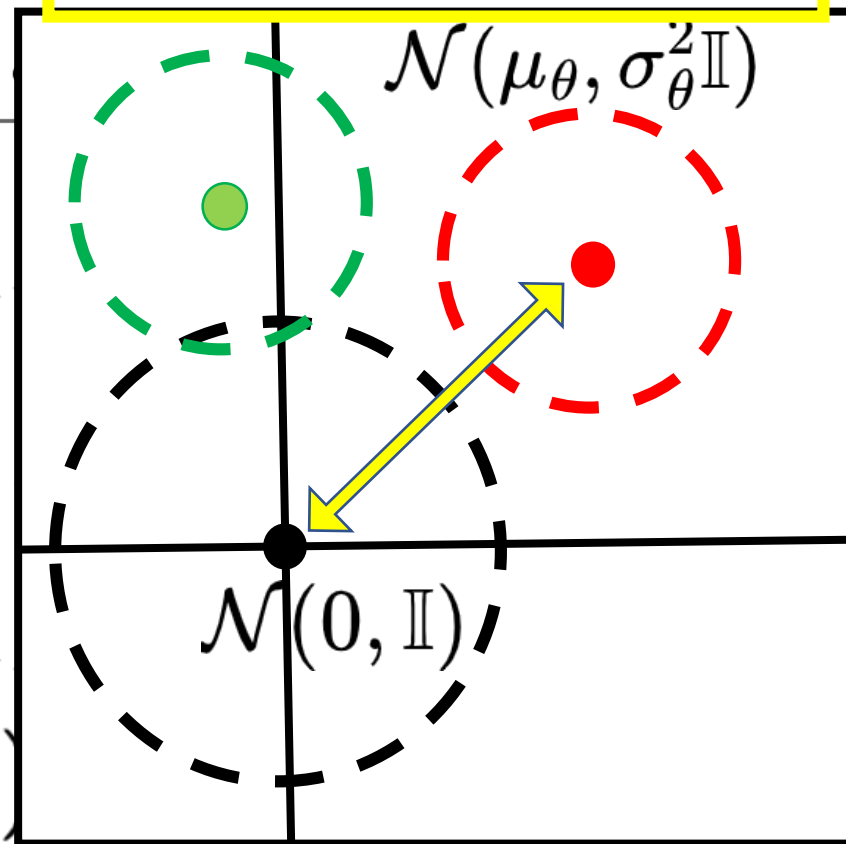
$$z_{\theta} = \mu_{\theta} + \sigma_{\theta} \cdot \mathcal{N}(0, \mathbb{I})$$

`torch.randn_like(std)`

Variational Autoencoders (VAE)



KL Between 2 Gaussians



Encourage $p(z) \sim \mathcal{N}(0,1)$:

$$\sum_{i=1}^n \sigma_i^2 + \mu_i^2 - \log(\sigma_i) - 1$$

“Reparametrization Trick”:

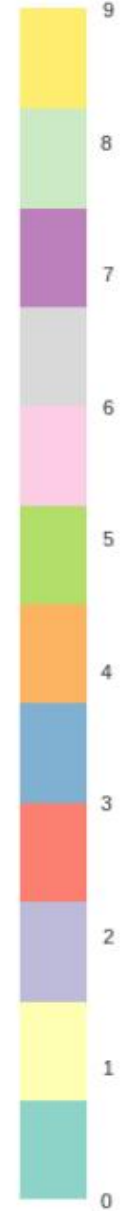
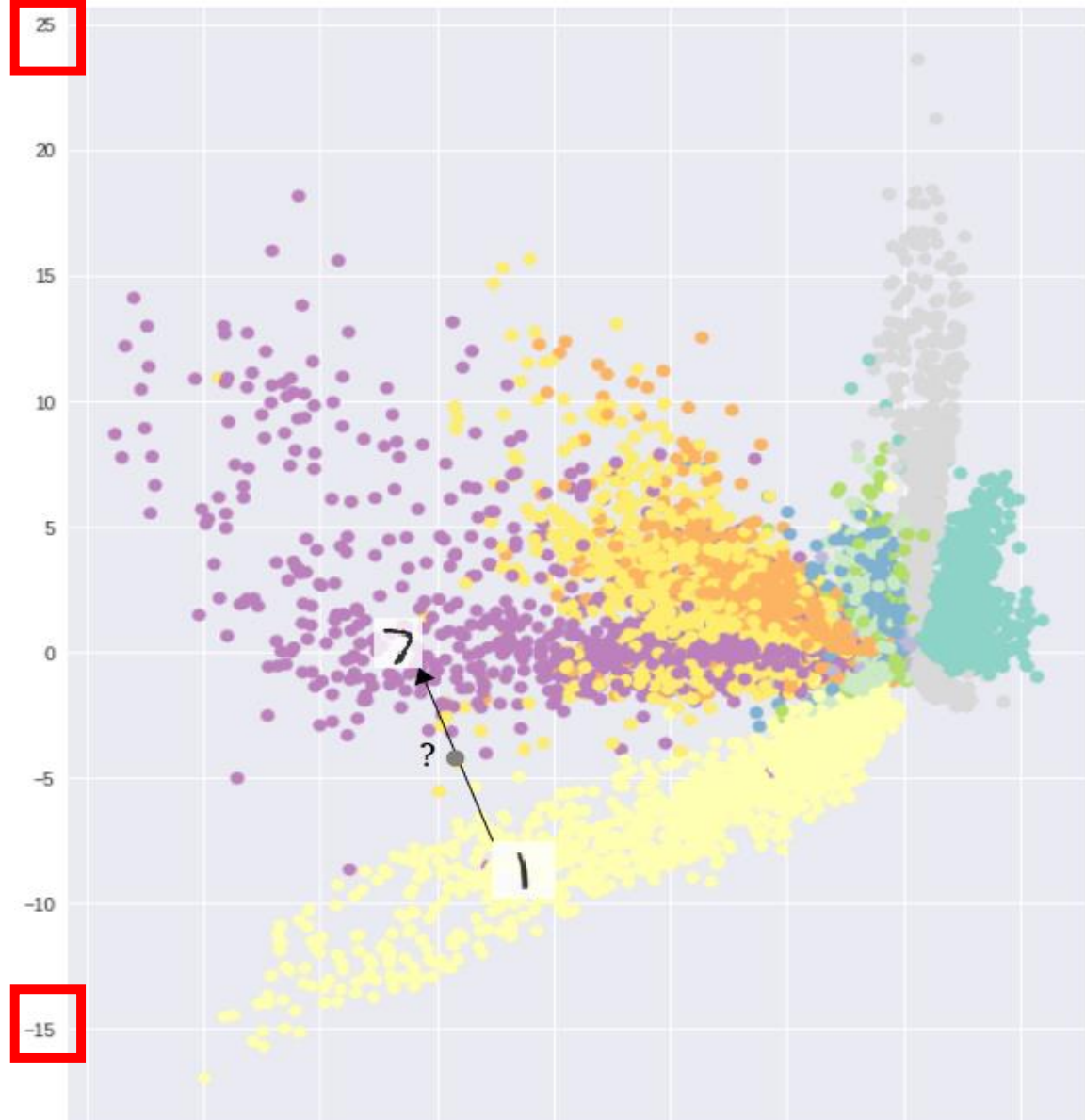
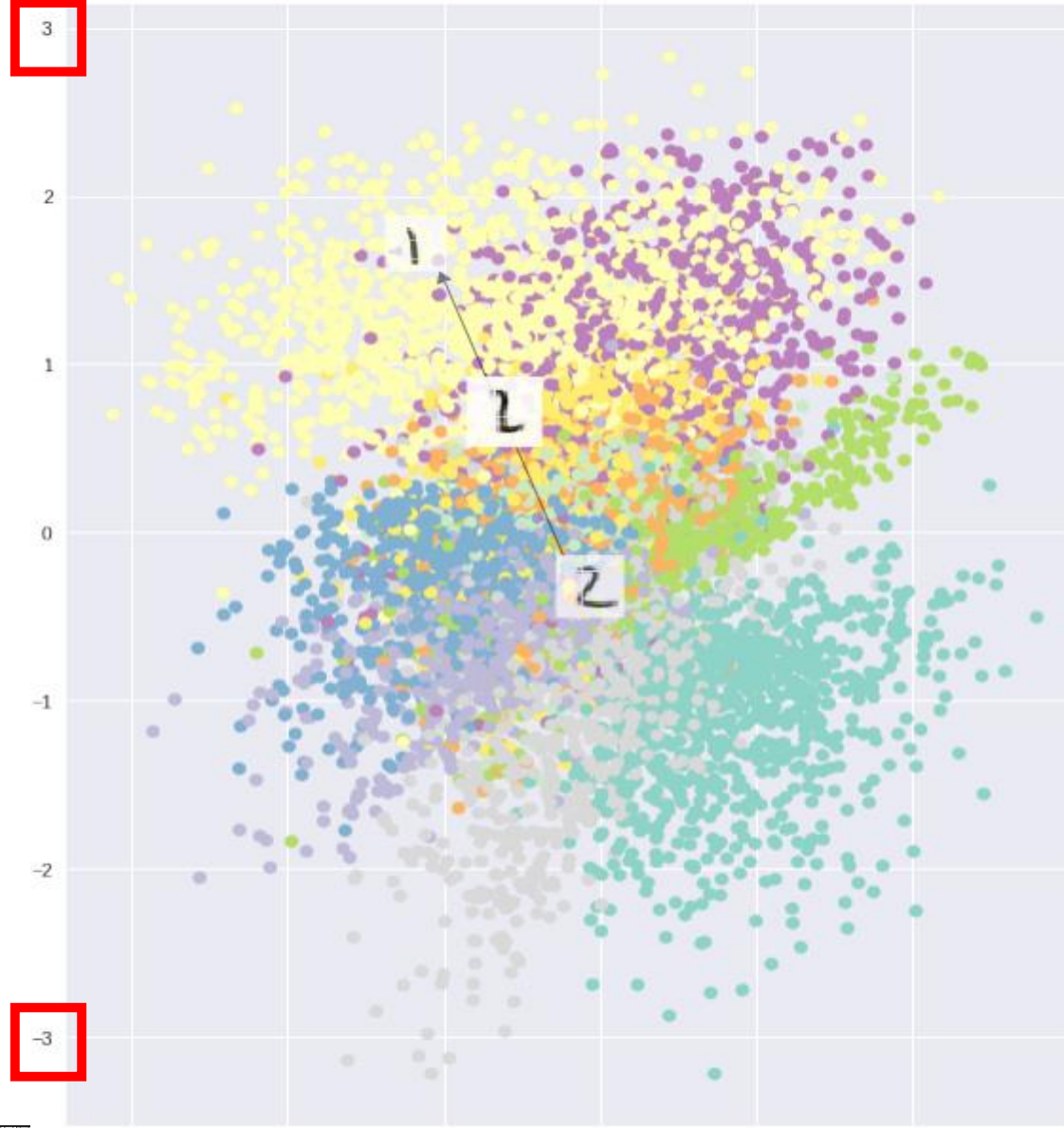
$$z_\theta = \mu_\theta + \sigma_\theta \cdot \mathcal{N}(0, \mathbb{I})$$

`torch.randn_like(std)`

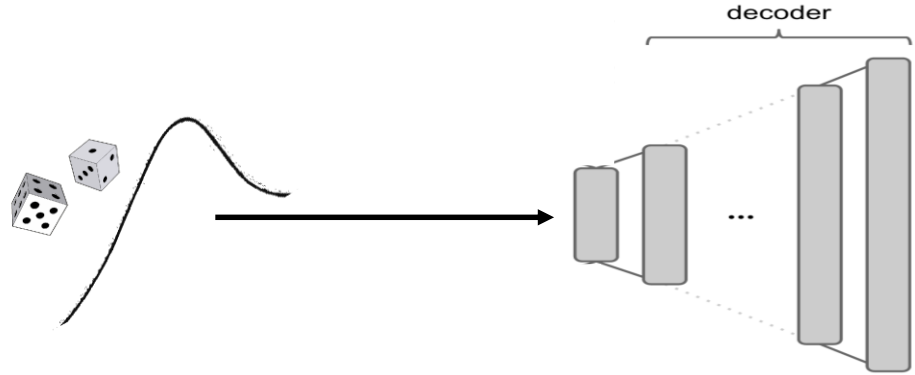
VAE

Also check out the scale!

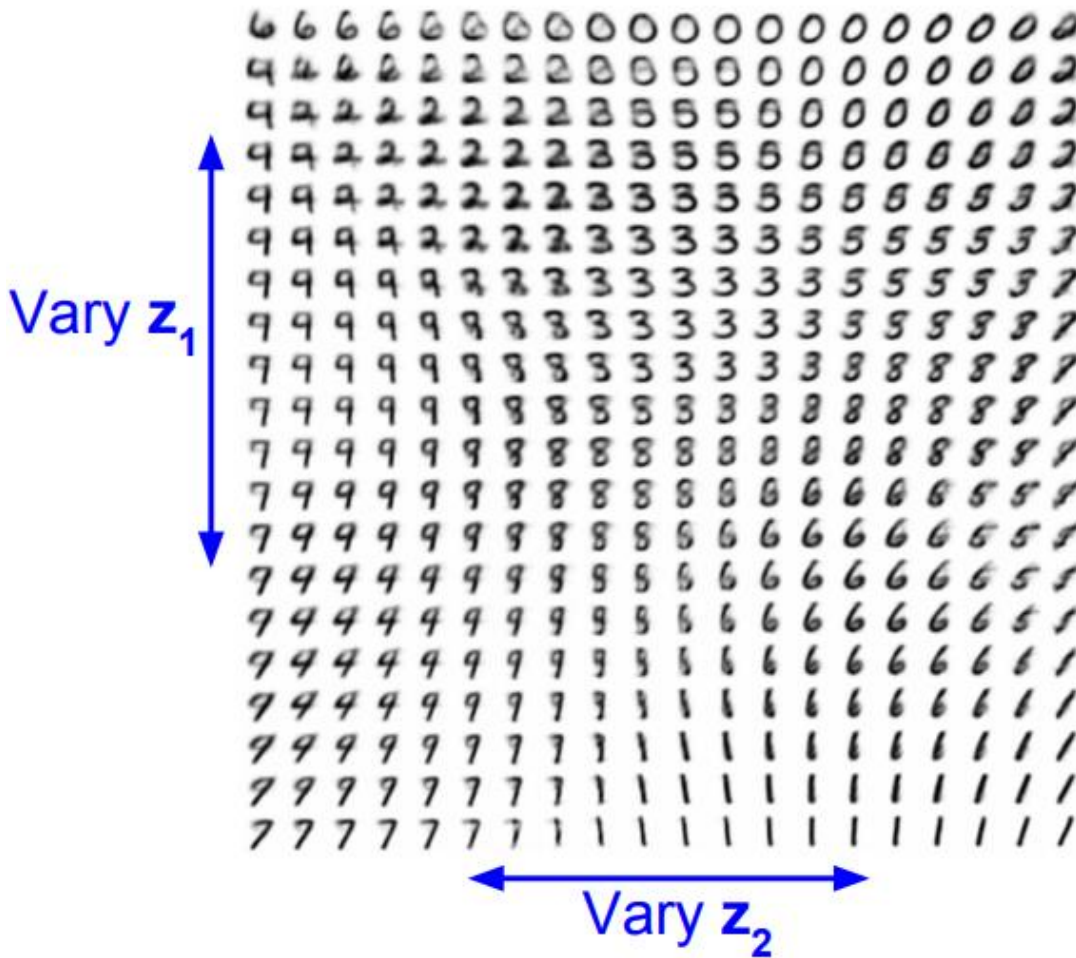
AE



Generate Data



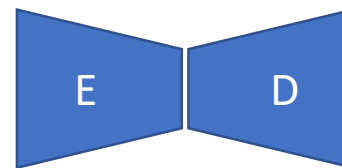
Data manifold for 2-d z



From VAE paper



Probabilistic Interpretation

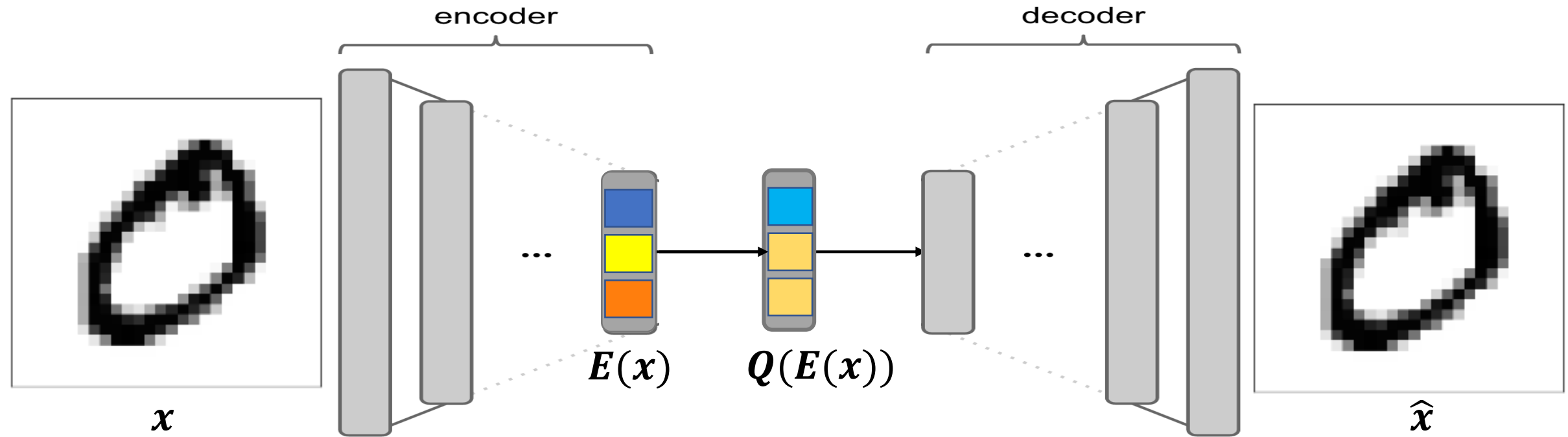
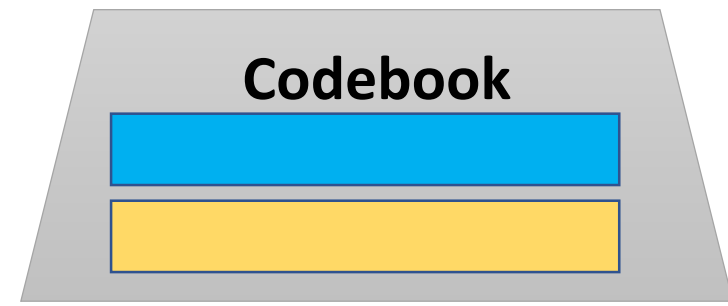


Reconstruction
+ Latent Prior

$$\log p_{\theta}(x^{(i)}) = \dots$$

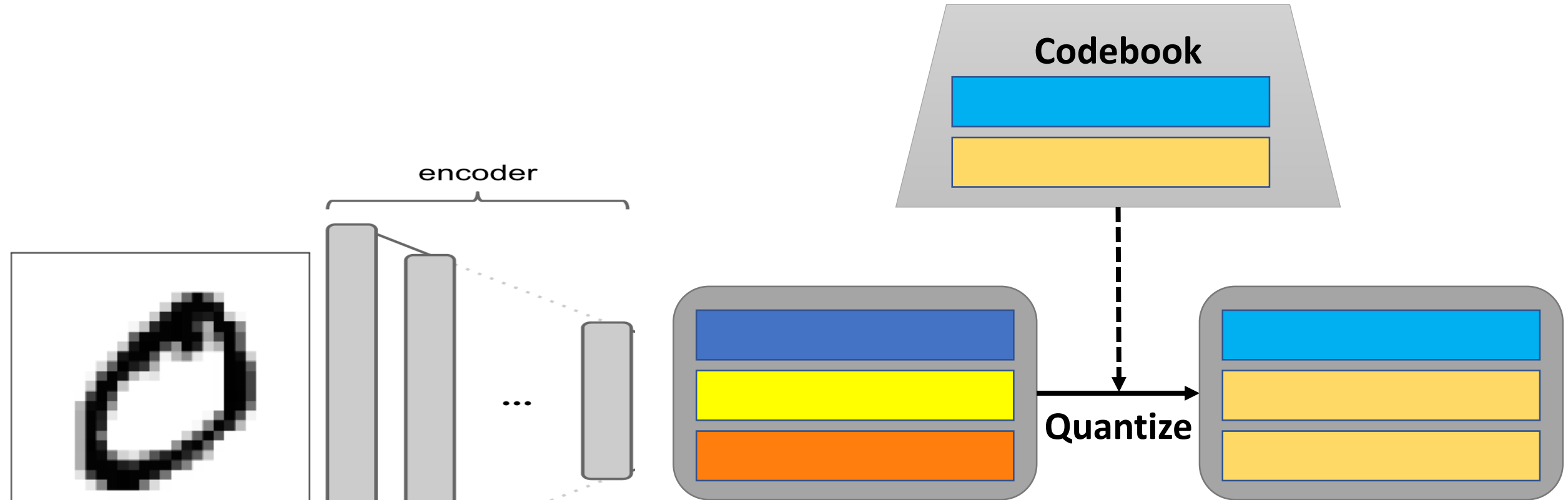
VAEs minimize a lower
bound of the (minus)
log likelihood

Vector-Quantized (VQ) VAE



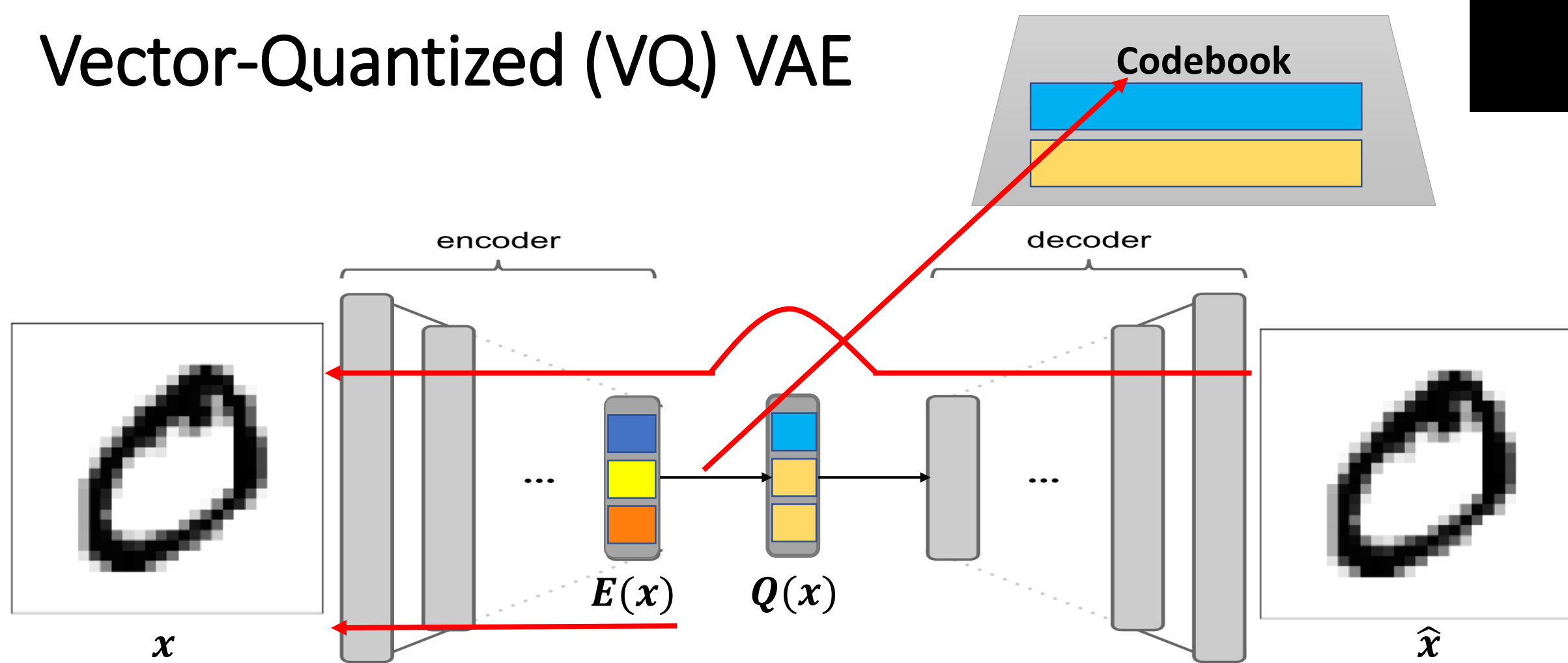
Quantization / Discretization

Vector-Quantized (VQ) VAE



$$Q(x)_i = \min_{z_j \in \text{Codebook}} \|E(x)_i - z_j\|$$

Vector-Quantized (VQ) VAE



$$\mathcal{L}_{rec} = \|\hat{x} - x\|_2^2 \quad \mathcal{L}_{commit} = \|E(x) - sg(Q(x))\|_2^2 \quad \mathcal{L}_{codebook} = \|sg(E(x)) - Q(x)\|_2^2$$

Quantization is non-differentiable!

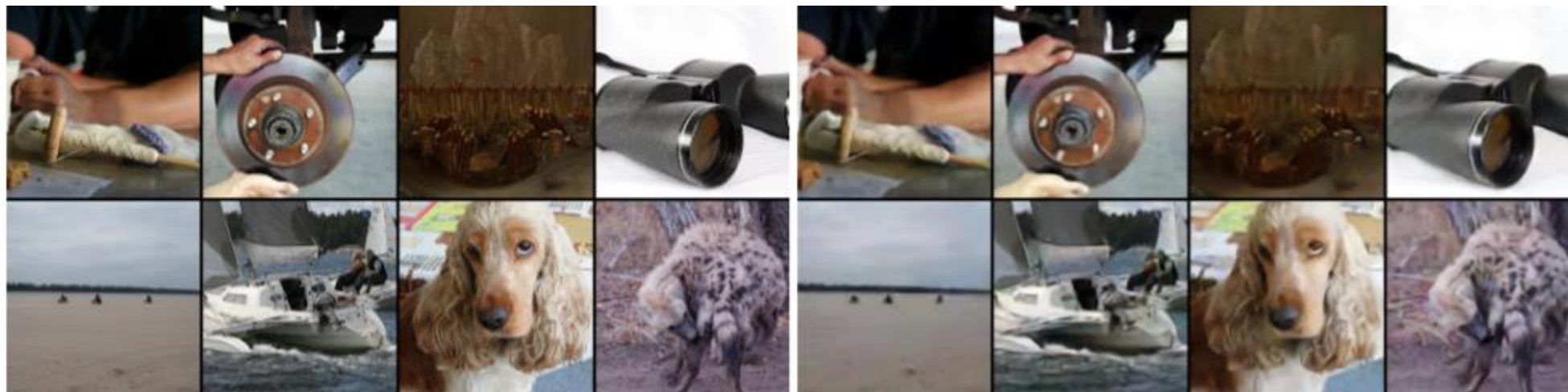
sg = stop-gradient



VQ-VAE Reconstructions

Real

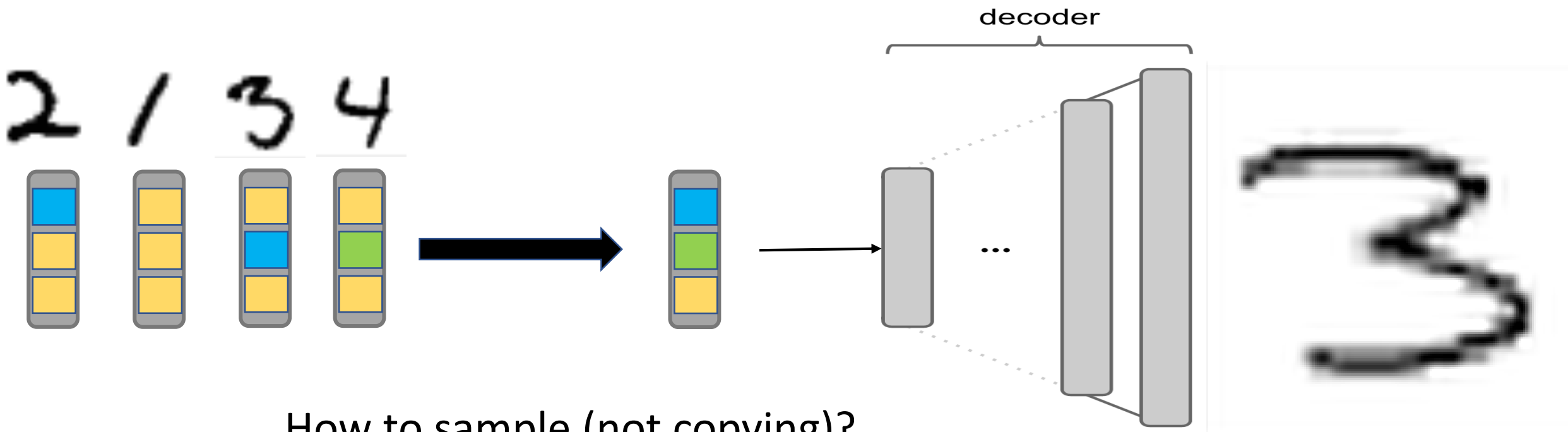
Reconstructed



“Neural Discrete Representation Learning “ [van den Oord et al., 2017]

Sampling New Instances

How to sample new “sentence”?

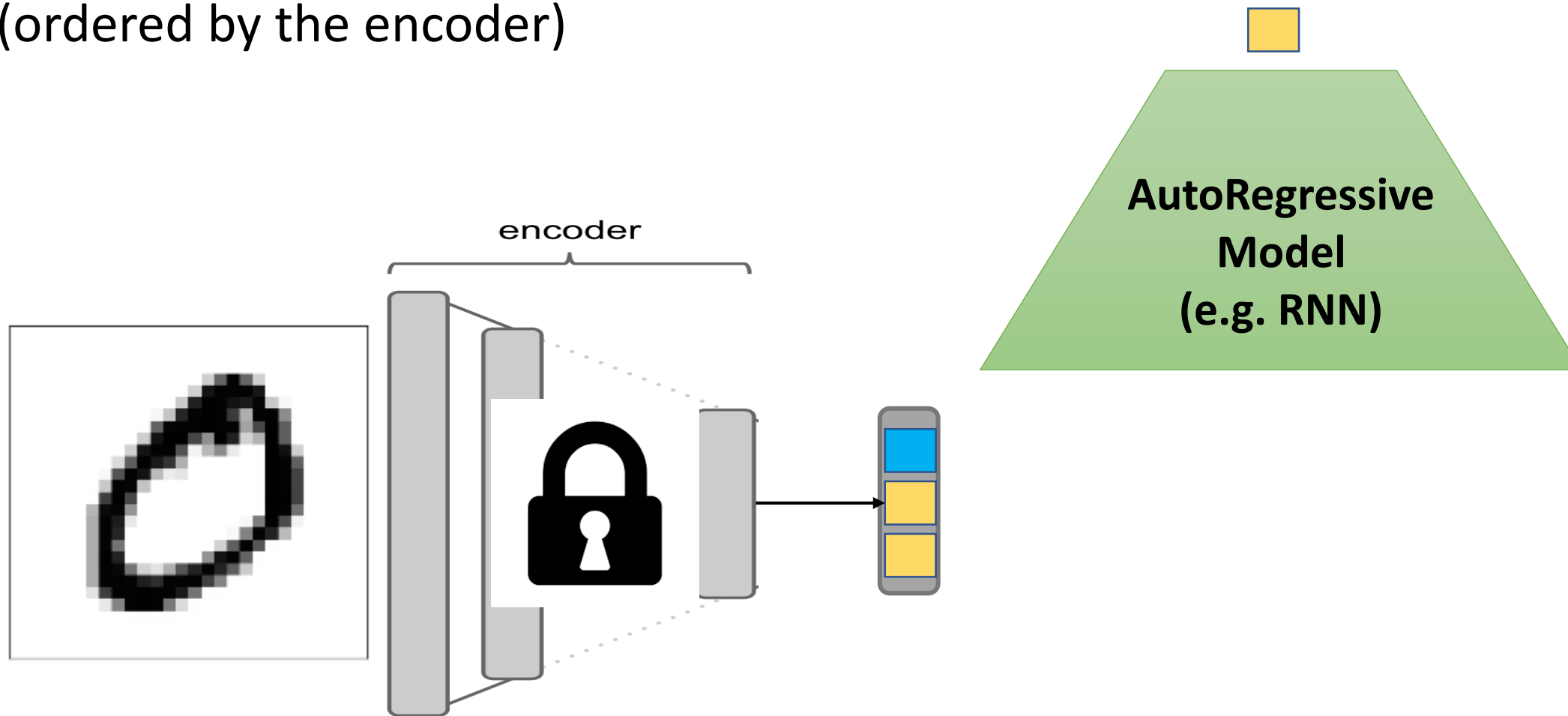


How to sample (not copying)?

Train an autoregressive model to predict “words”

Sampling New Instances

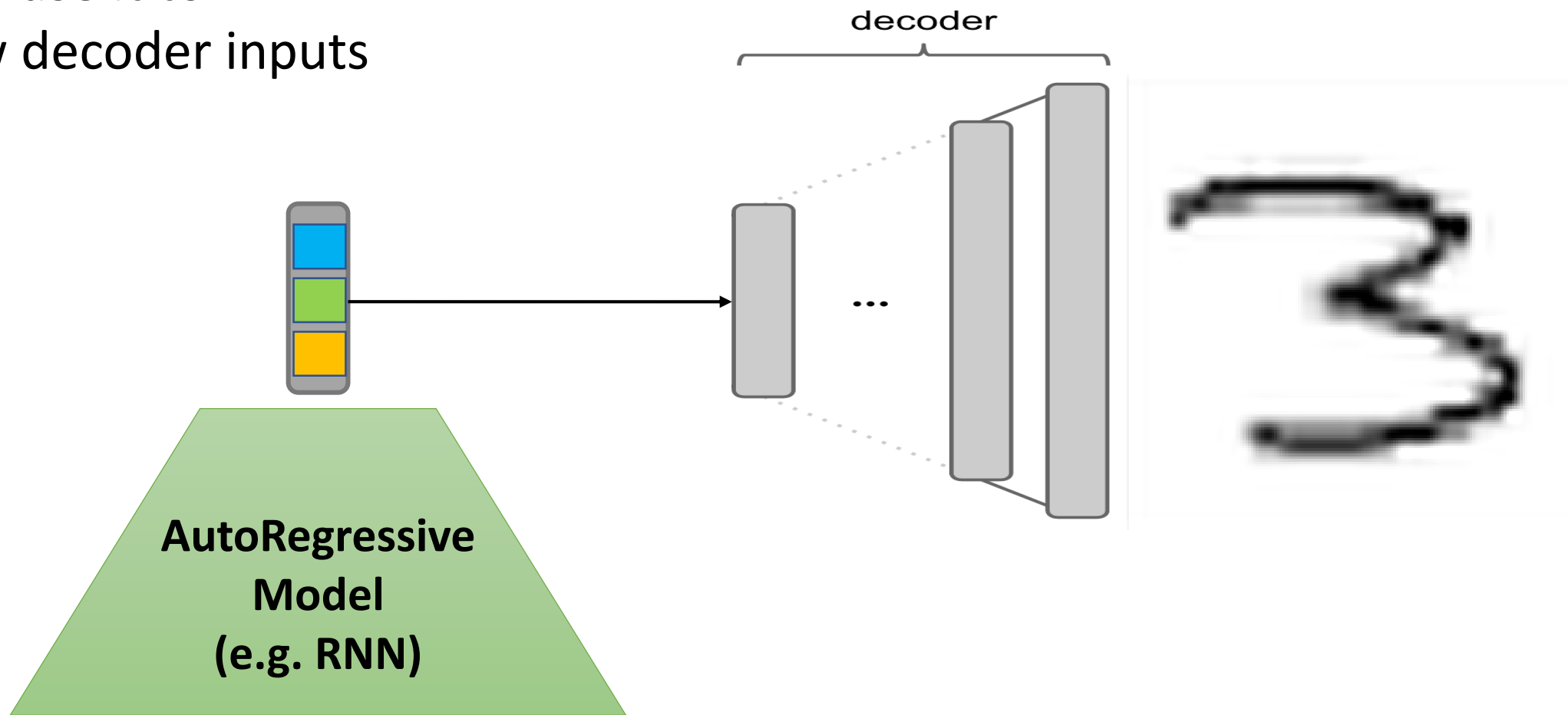
Training an autoregressive
on codebook elements
(ordered by the encoder)



“Neural Discrete Representation Learning “ [van den Oord et al., 2017]

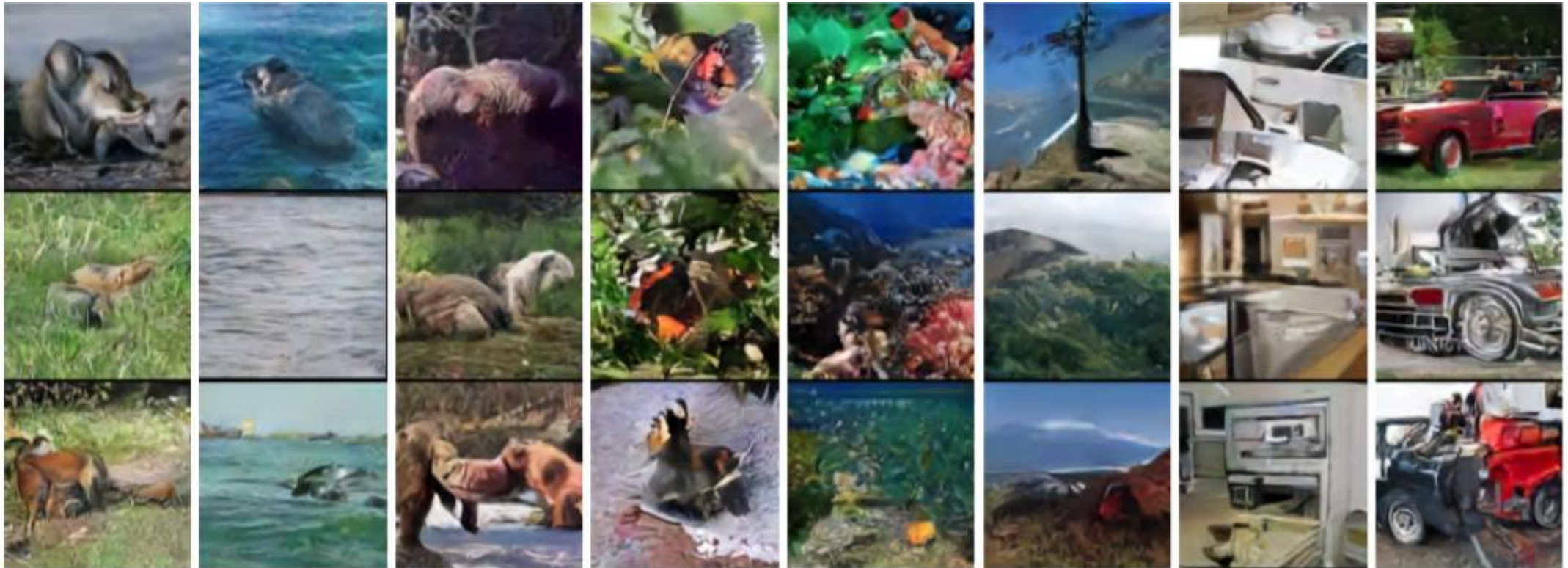
Sampling New Instances

Once trained, use it to generate new decoder inputs



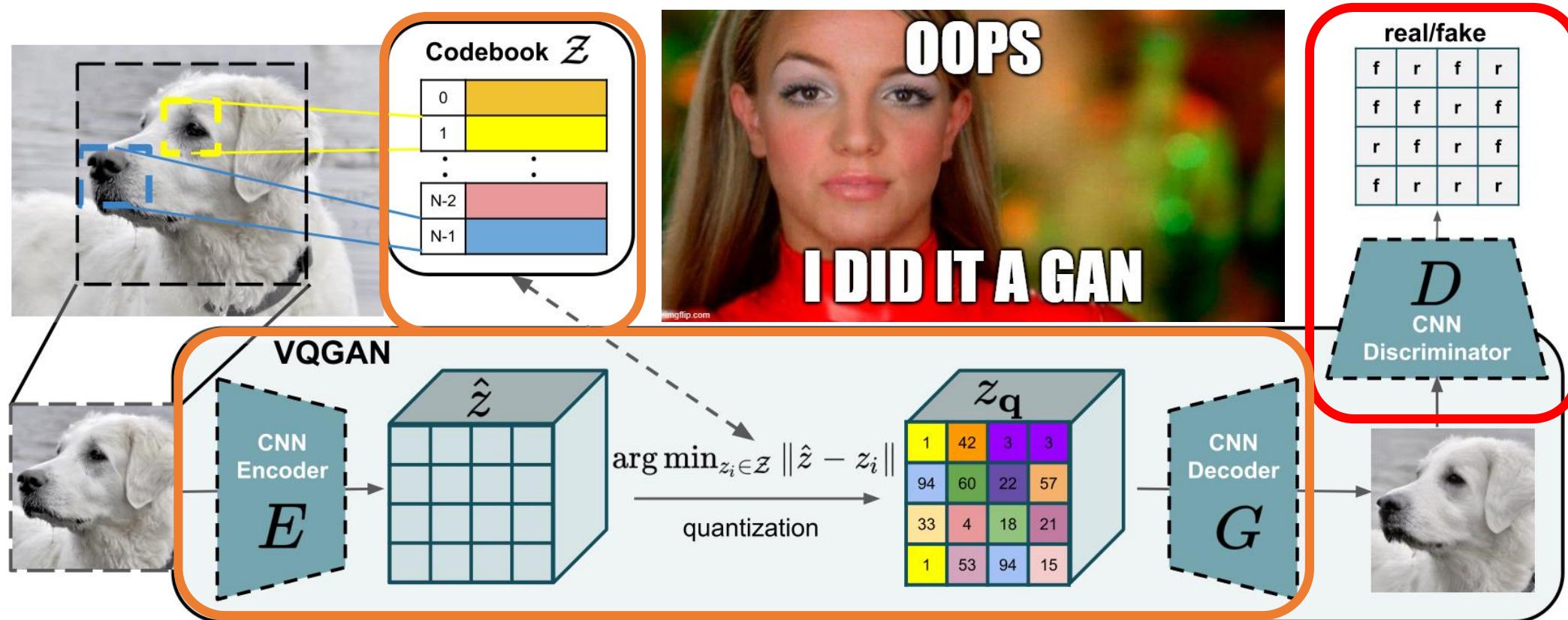
“Neural Discrete Representation Learning “ [van den Oord et al., 2017]

Sampling New Instances - Results



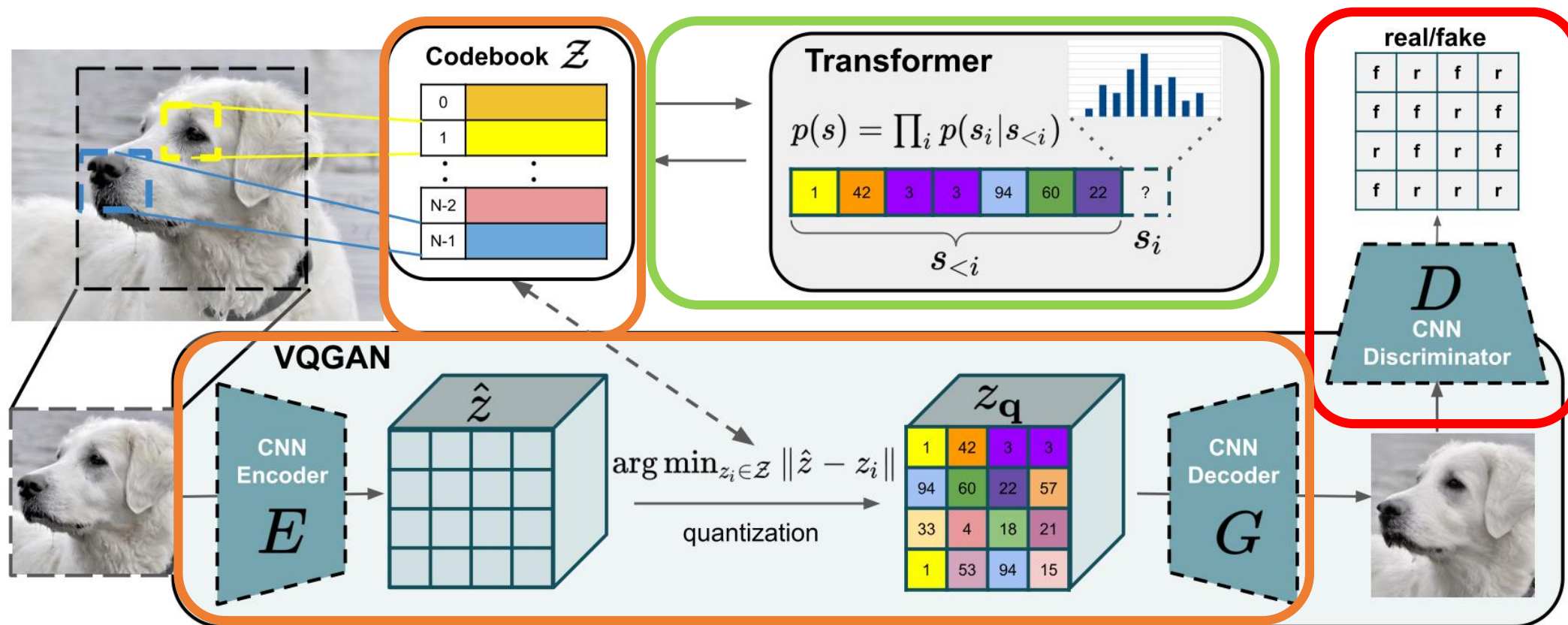
“Neural Discrete Representation Learning “ [van den Oord et al., 2017]

VQGAN (Taming Transformers)



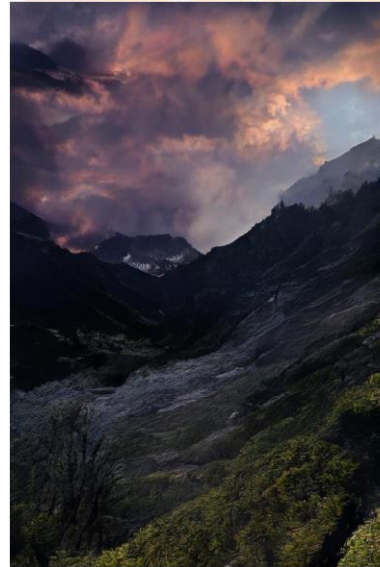
VQ-VAE + GAN

VQGAN (Taming Transformers)



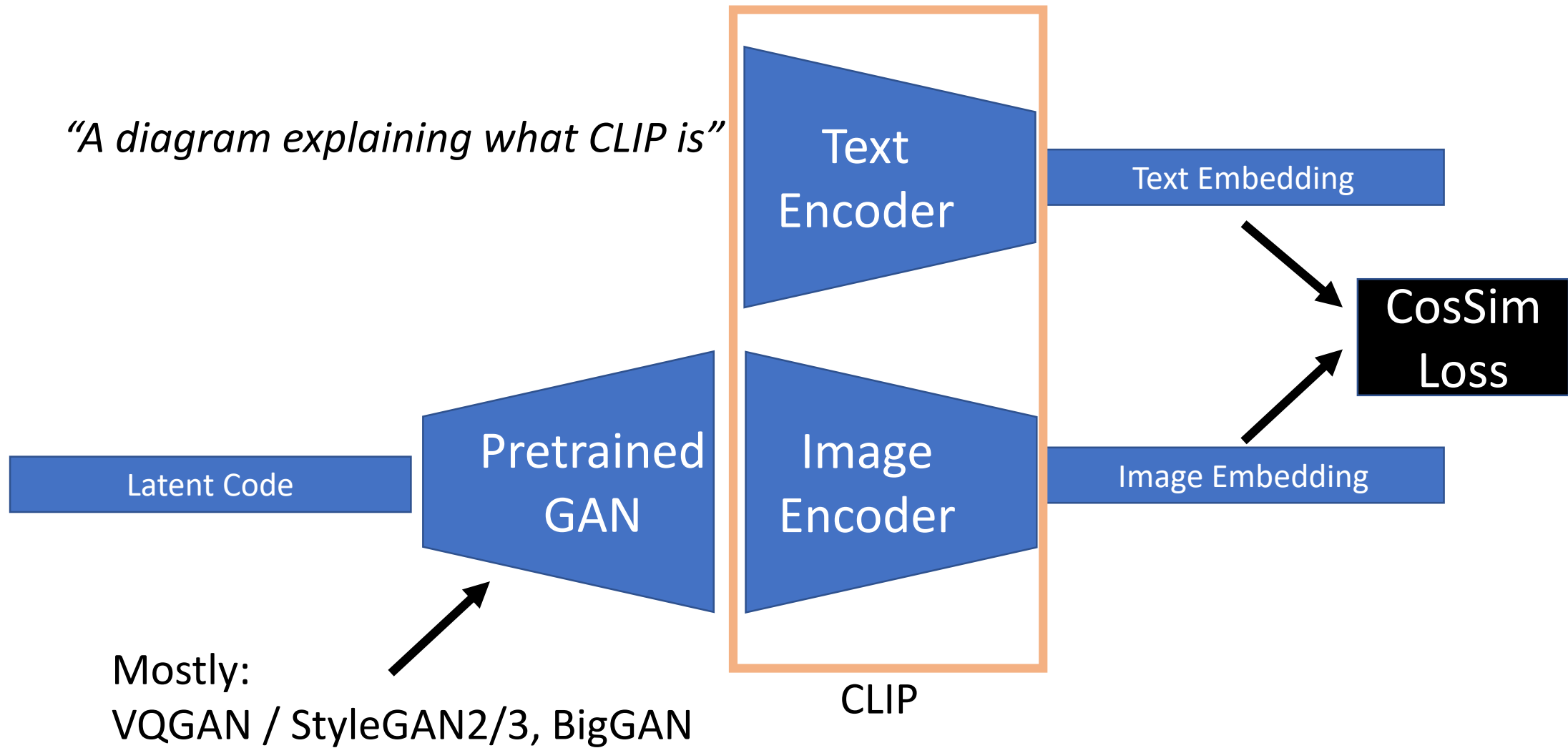
VQ-VAE + GAN Transformer

VQGAN (Taming Transformers)



Taming Transformers for High-Resolution Image Synthesis [Esser et al. 2020]

CLIP: Contrastive Language-Image Pre-Training



CLIP: Contrastive Language-Image Pre-Training



Nerdy
@Nerdy



Rikkar.tez
@socalpathy

working on some more #pixelart animations, this one is a fantasy valley created entirely with #AI.

#generativeart #pixels #vqganclip #tezos 
#hicetnunc



<https://twitter.com/RiversHaveWings>

...
outerart
aper



3985932



(~)VQ-VAE + Transformers: DALL-E

Ramesh et al., 2021

Teapot in the shape of a rubik's cube



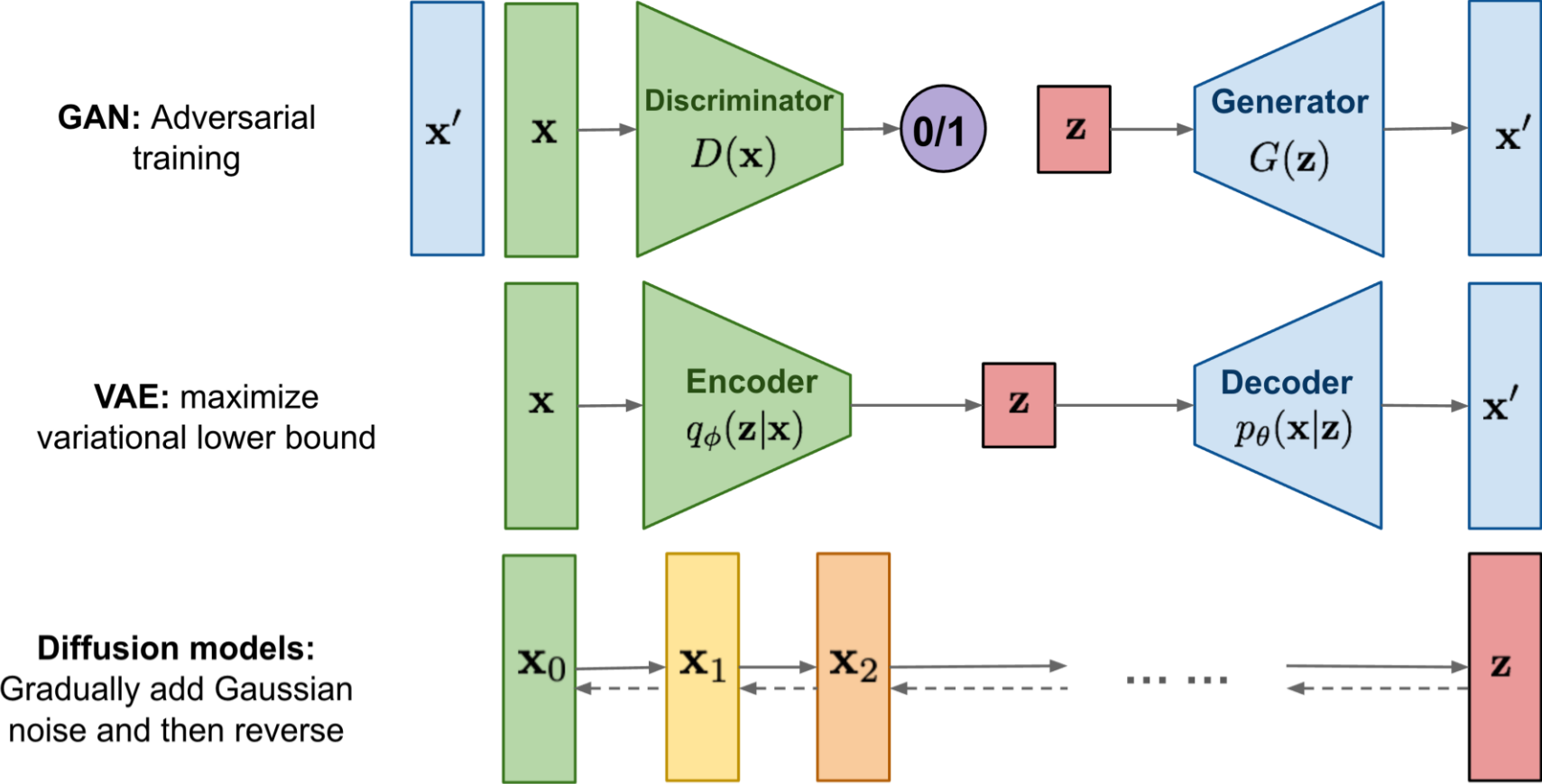
Soap-dispenser in the shape of a doughnut



Store front with 'pytorch'

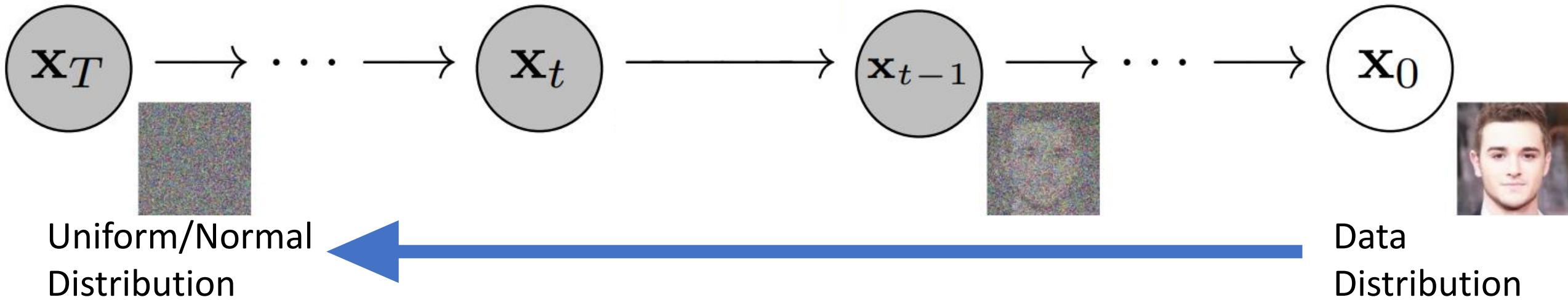


Diffusion Models



Diffusion Models

1. How are we going to do that?!
2. Why should it even work?



Dye represents data density

Observation:

Diffusion destroys structure

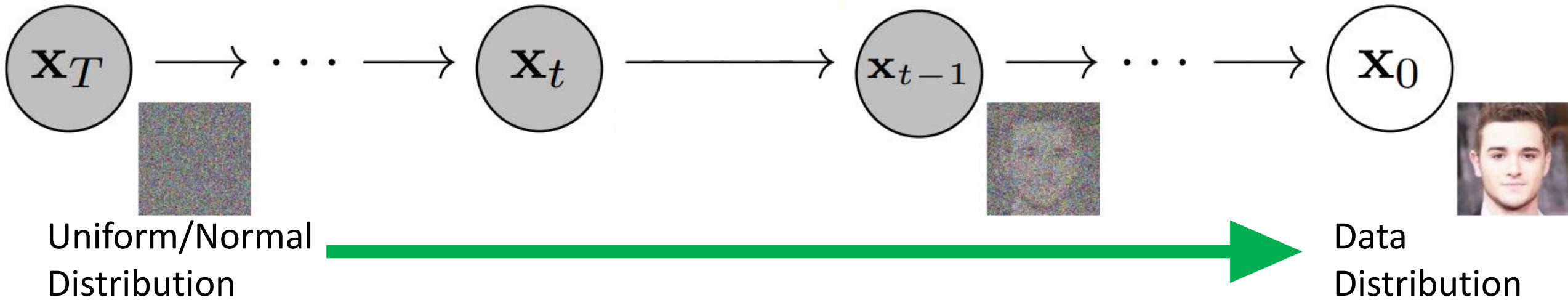
Core Idea:

Recover structure by reversing time



Diffusion Models

1. How are we going to do that?!
2. Why should it even work?



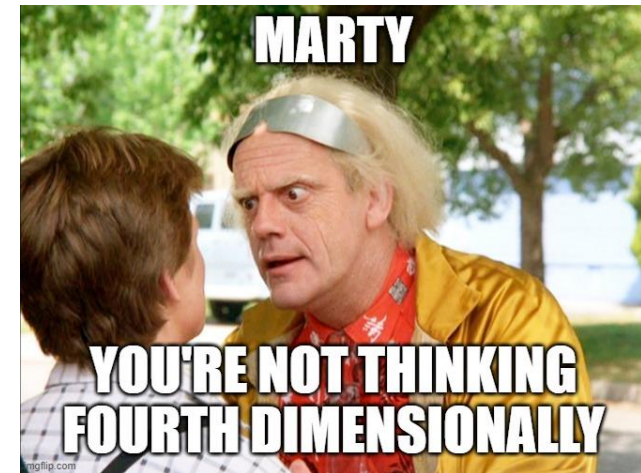
Dye represents data density

Observation:

Diffusion destroys structure

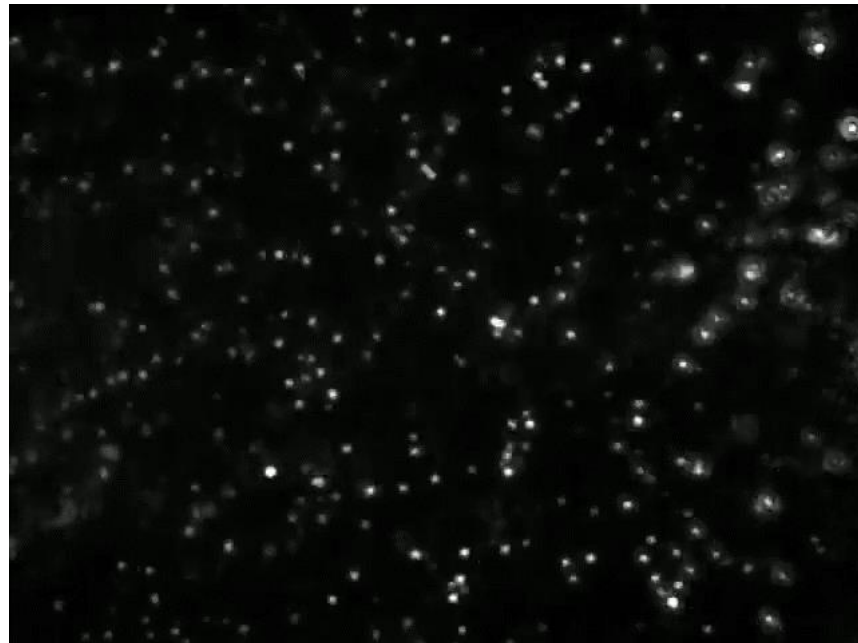
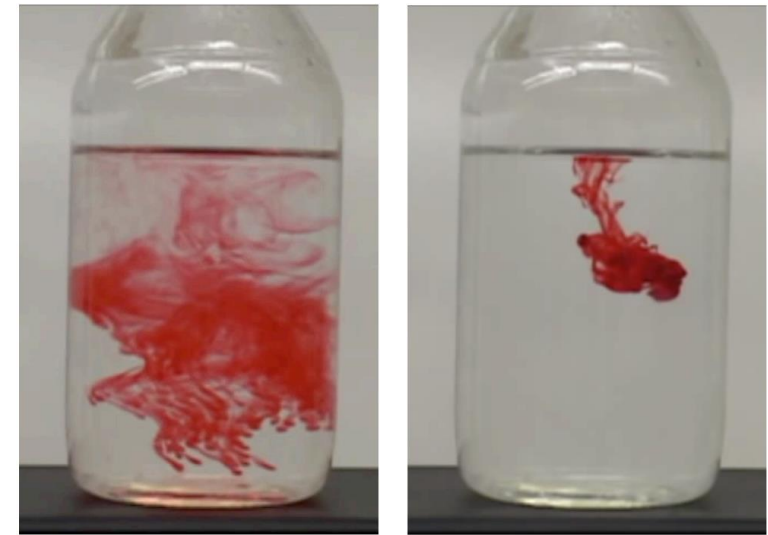
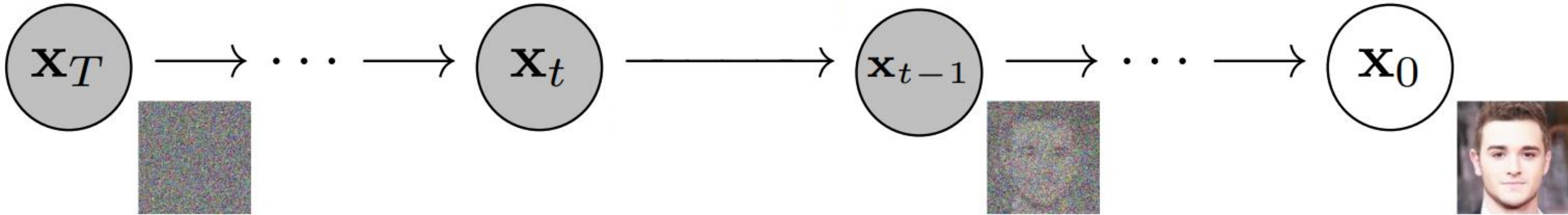
Core Idea:

Recover structure by reversing time



Diffusion Models

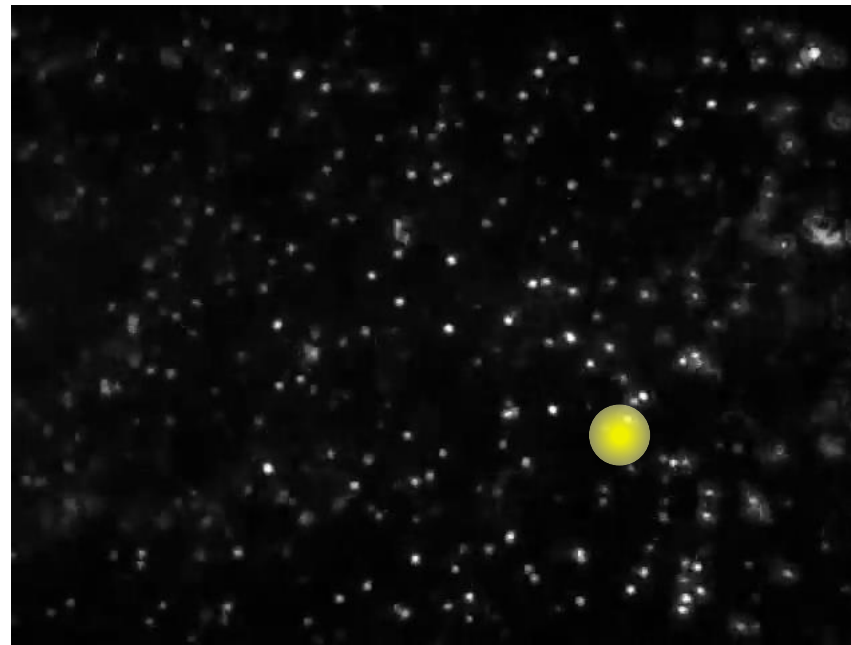
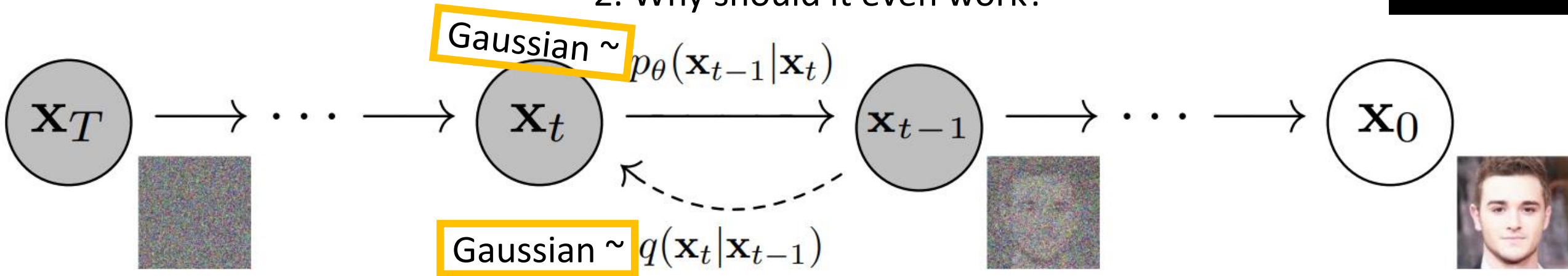
1. How are we going to do that?!
2. Why should it even work?



© [Rutger Saly](#)

Diffusion Models

1. How are we going to do that?!
2. Why should it even work?



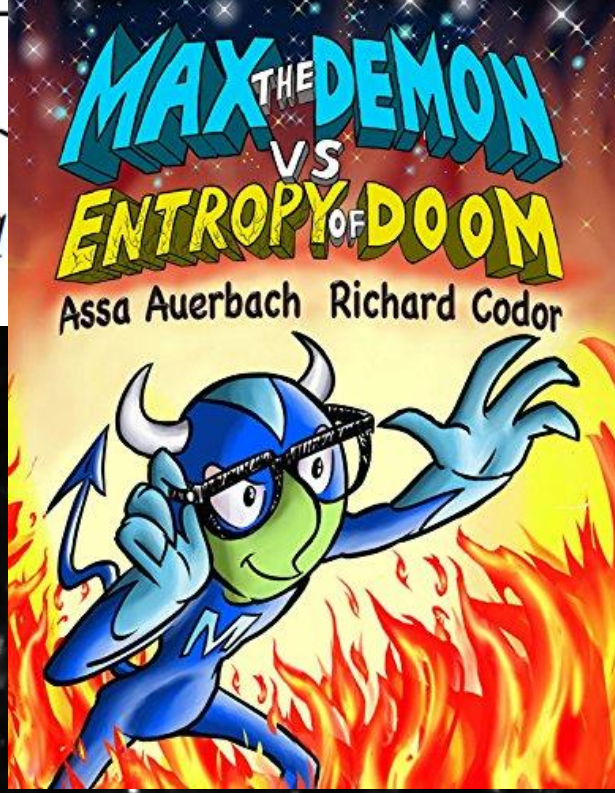
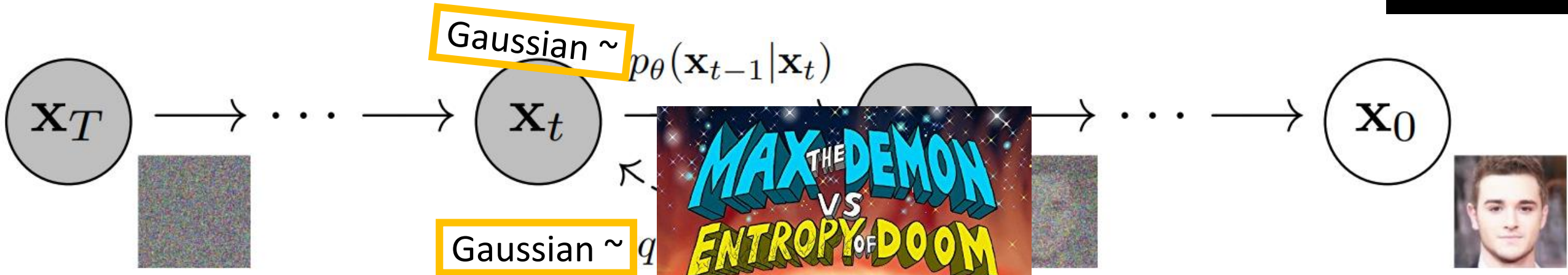
Position updates
are small Gaussian

For Small step size:
Forward and Reverse
Have similar functional
shape [Feller, 1949]

© [Rutger Saly](#)

Diffusion Models

1. How on earth are we going to do that?!
2. Why on earth should it even work?



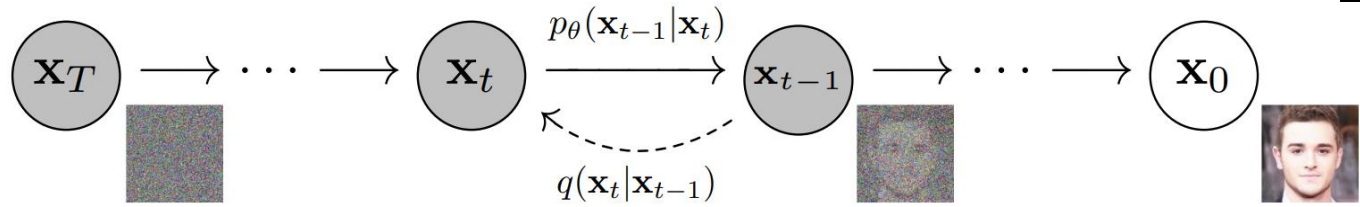
Position updates are small Gaussian

For Small step size:
Forward and Reverse
Have similar functional shape [Feller, 1949]

© [Rutger Saly](#)



Diffusion Models



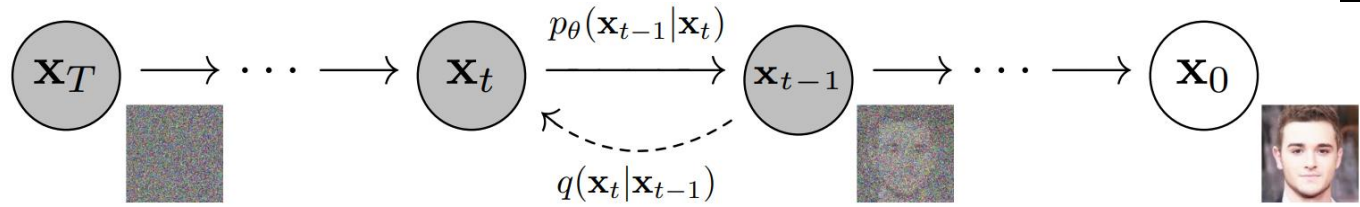
$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

add small noise

decay to origin

Forward process:
Gradually adding noise
According to beta schedule

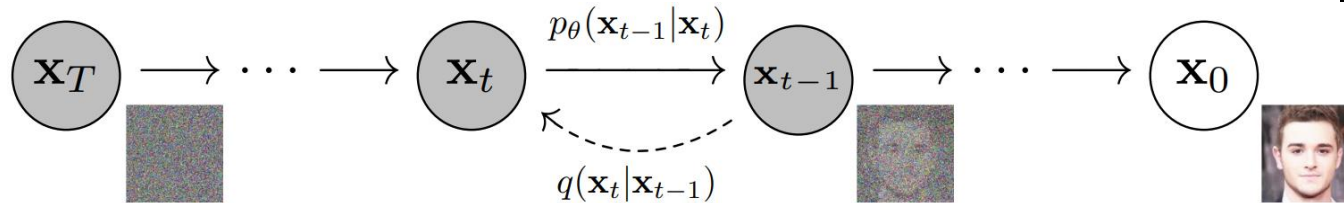
Diffusion Models



$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ Intractable ☹️ (depends on entire dataset)

Diffusion Models



$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ Intractable ☹️ (depends on entire dataset)

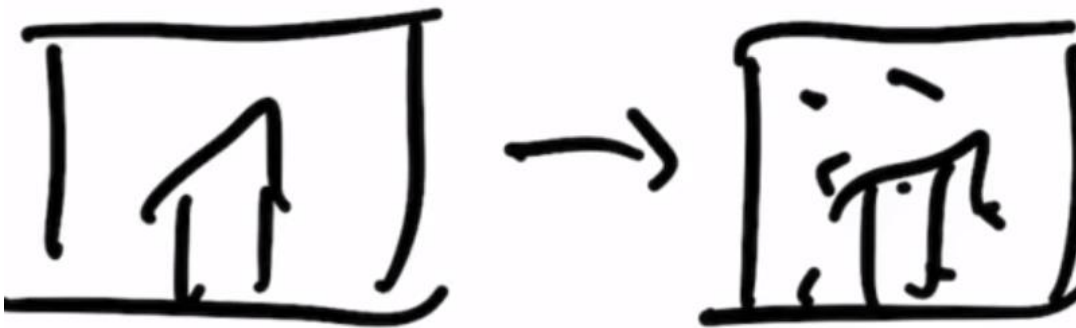
Reverse:

Needs a lot (!) of priors about the data/world

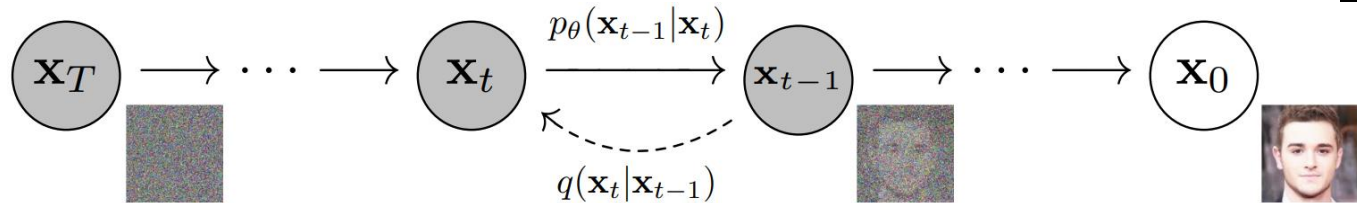


Forward:

Easy



Diffusion Models



$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ Intractable ☹️ (depends on entire dataset)

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

$$\boldsymbol{\mu}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right)$$

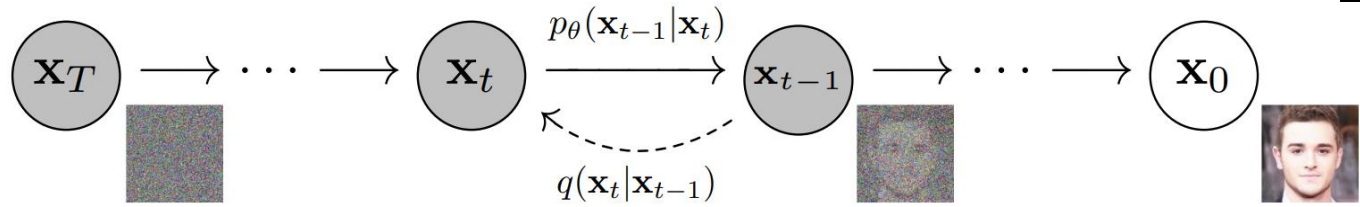
$$\boldsymbol{\Sigma}(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$$

$$\alpha_t = 1 - \beta_t$$

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

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

Diffusion Models



$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ Intractable ☹️ (depends on entire dataset)

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2\mathbf{I})$$

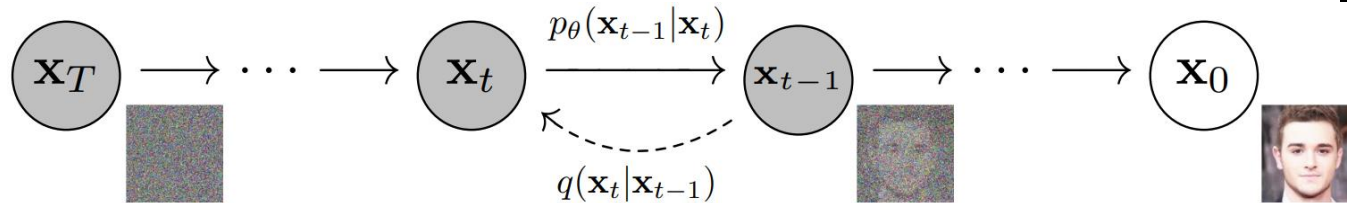
$$\boldsymbol{\mu}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right)$$

$$\alpha_t = 1 - \beta_t$$

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

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

Diffusion Models



$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ Intractable ☹️ (depends on entire dataset)

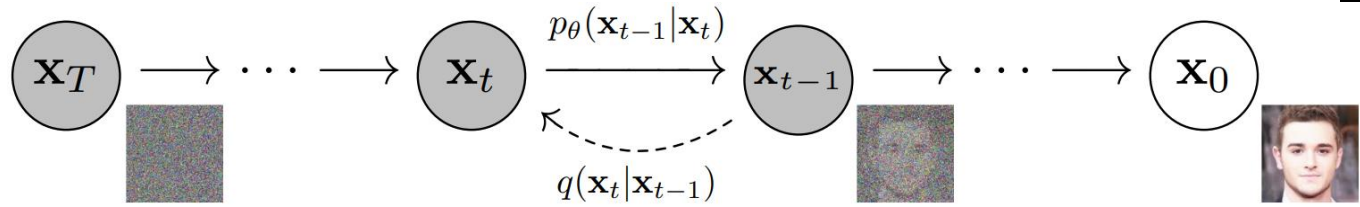
$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

$$\boldsymbol{\mu}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t \right)$$

Algorithm 1 Training

- 1: **repeat**
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on
$$\nabla_\theta \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$
- 6: **until** converged

Diffusion Models



$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ Intractable ☹️ (depends on entire dataset)

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

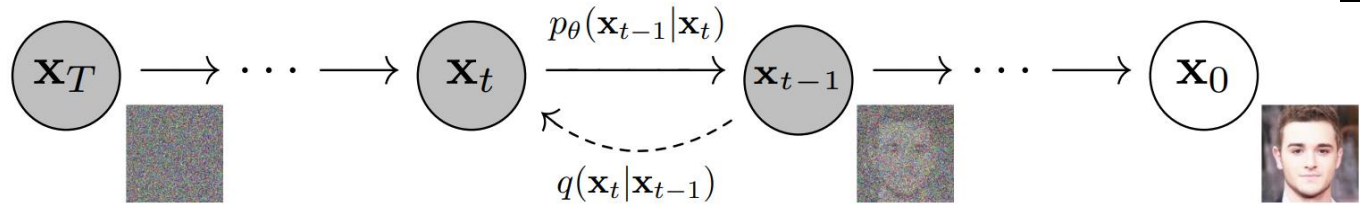
Train by maximizing log-likelihood:

$$L_{\text{CE}} = \mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)]$$

= ... tons of algebra ...

$$\leq [\text{terms of KL/entropy of Gaussians}](\theta)$$

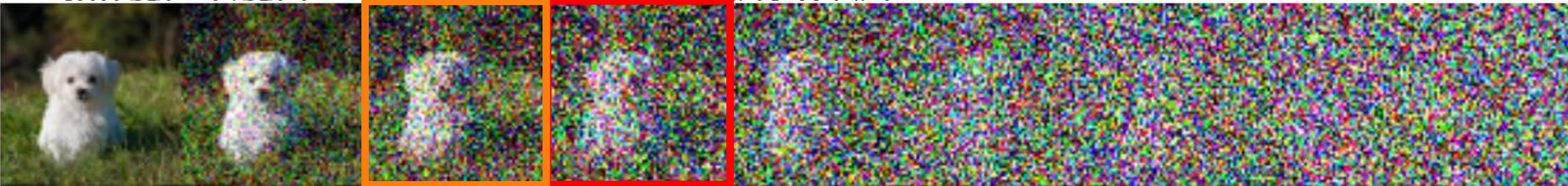
Diffusion Models



$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ Intractable ☹️ (depends on entire dataset)

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma^2 \mathbf{I})$$



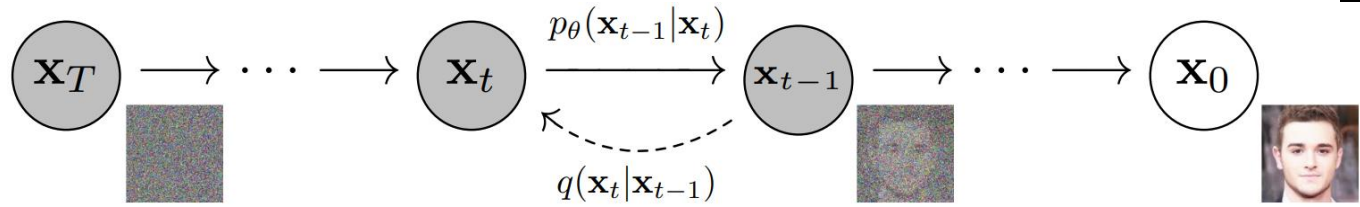
= ... tons of algebra ...

$$\leq \dots + D_{KL}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))$$

$$= q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0) \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)}$$

Tractable!

Diffusion Models



$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

$q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ Intractable ☹️ (depends on entire dataset)

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2\mathbf{I})$$

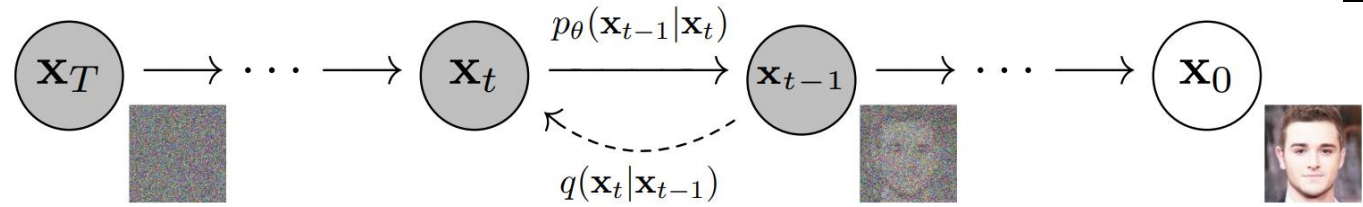
Train by maximizing log-likelihood:

$$L_{\text{CE}} = \mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)]$$

= ... tons of algebra ...

$$\leq \dots + D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))$$

Diffusion Models



$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ Intractable ☹️ (depends on entire dataset)

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

$$\boldsymbol{\mu}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right)$$

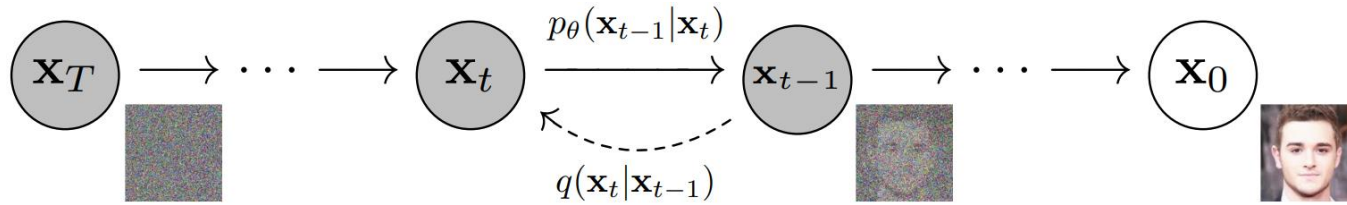
$$\alpha_t = 1 - \beta_t$$

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

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

$$L_{\text{CE}} = \mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{t, x_0, \boldsymbol{\epsilon}} [\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2]$$

Diffusion Models



Algorithm 1 Training

1: **repeat**

2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$

3: $t \sim \text{Uniform}(\{1, \dots, T\})$

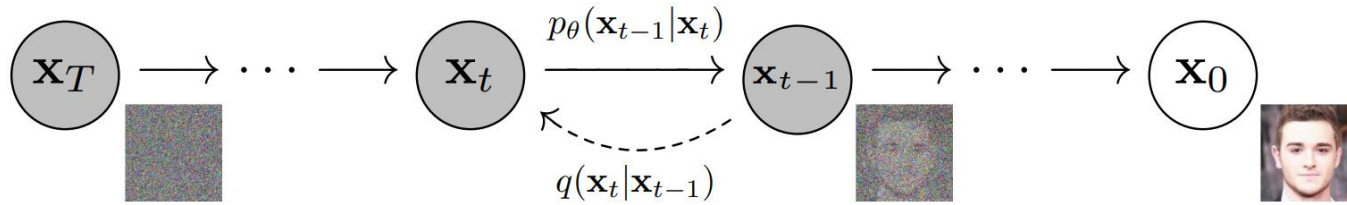
4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

5: Take gradient descent step on

$$\nabla_\theta \left\| \epsilon - \epsilon_\theta \left(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t \right) \right\|^2$$

6: **until** converged

Diffusion Models



Algorithm 1 Training

- 1: **repeat**
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on
$$\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, t)\|^2$$
- 6: **until** converged

Algorithm 2 Sampling

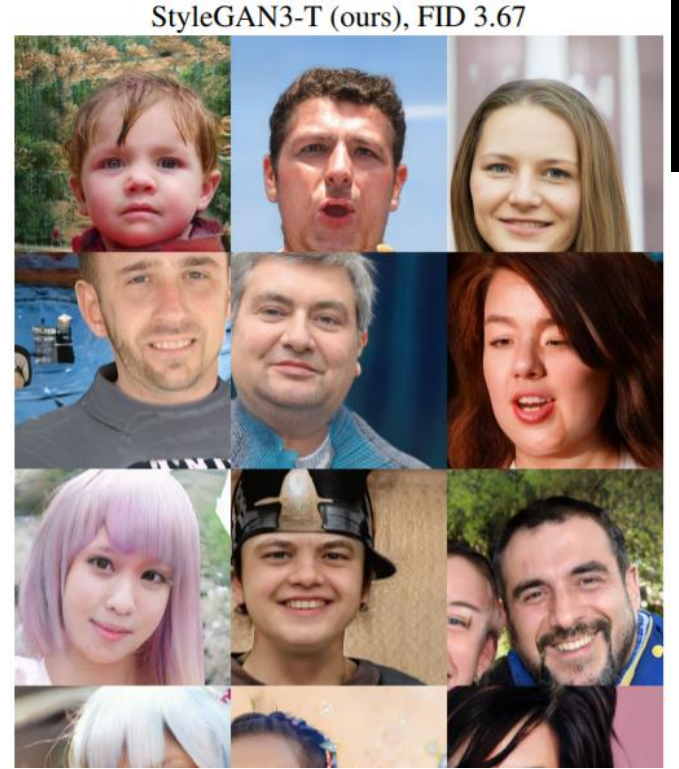
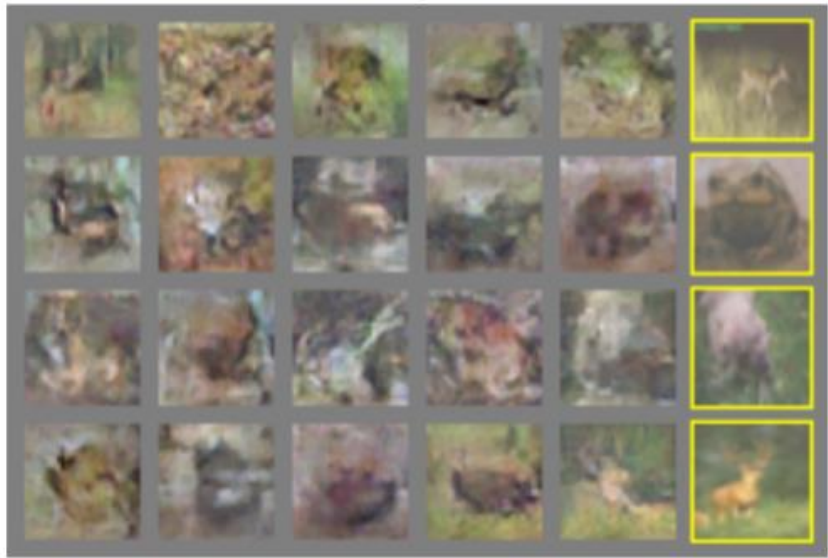
- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for** $t = T, \dots, 1$ **do**
- 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$
- 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return** \mathbf{x}_0

Computed from \mathbf{x}_0
(by adding noise - analytically)

GANs – 2014 → 2021



b)



Diffusion Models

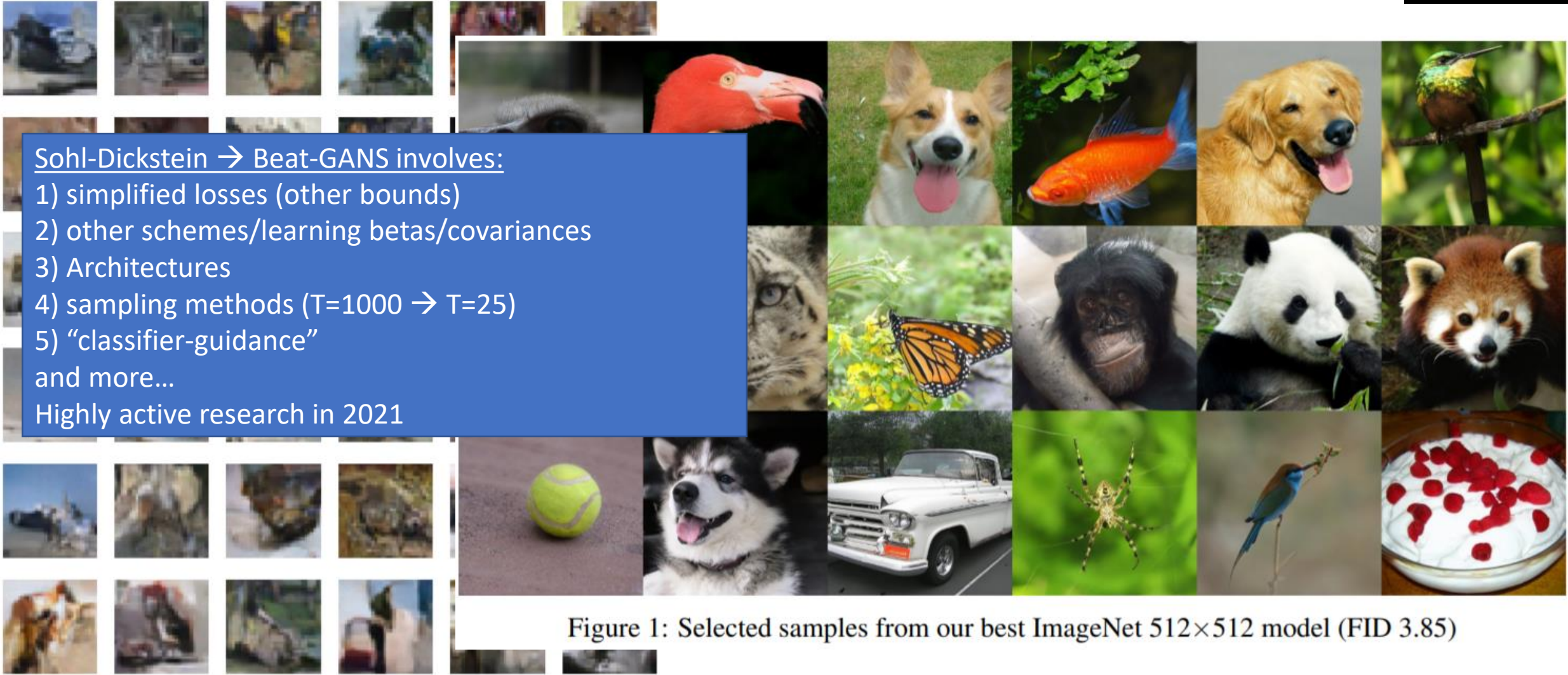


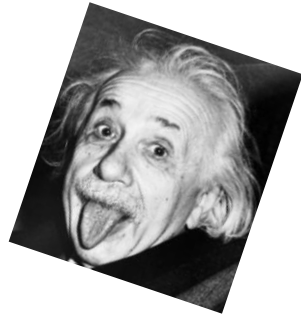
Figure 1: Selected samples from our best ImageNet 512×512 model (FID 3.85)

“Deep Unsupervised Learning using Nonequilibrium Thermodynamics” [Sohl-Dickstein et al. 2015]

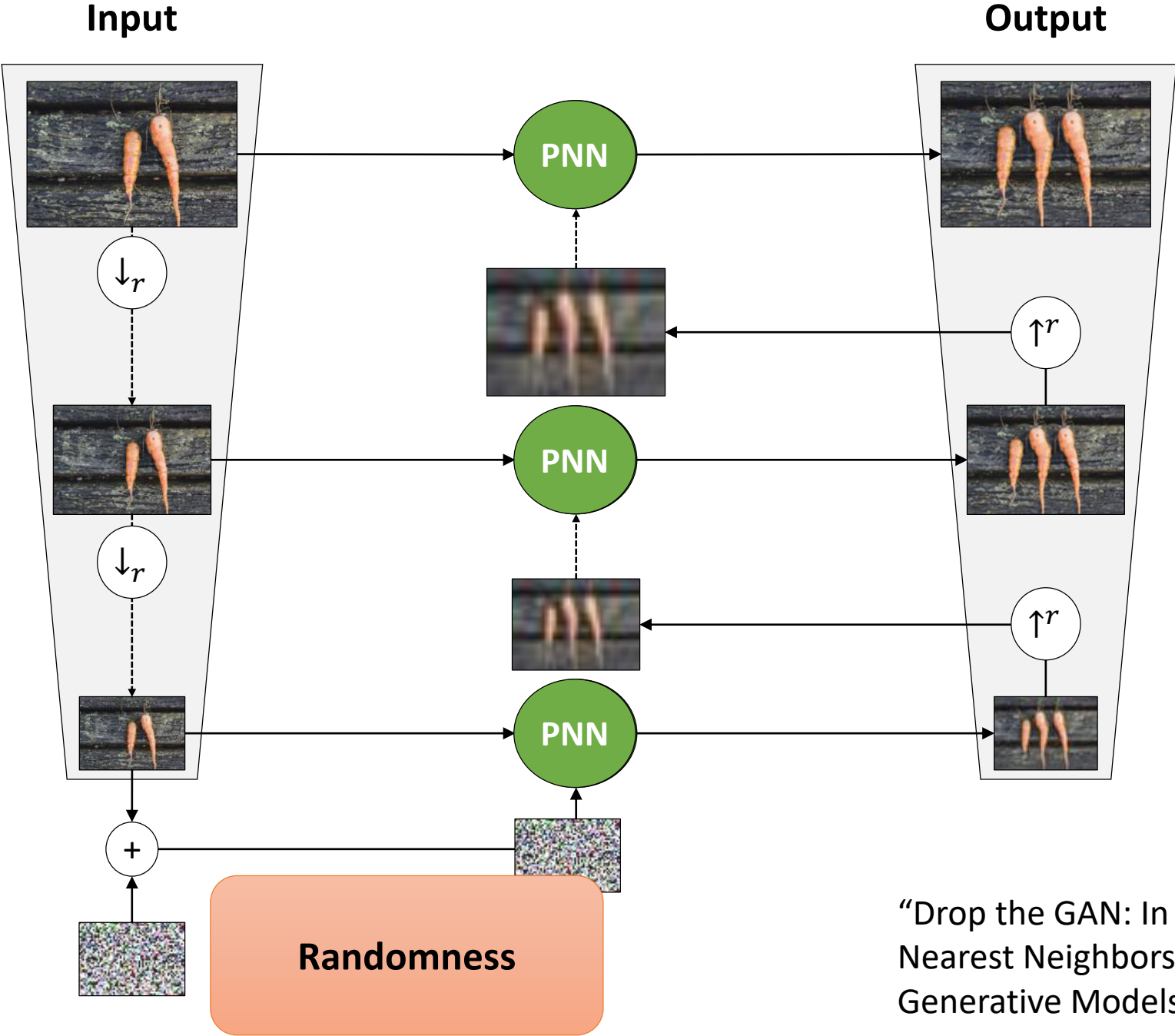
“Diffusion Models Beat GANs on Image Synthesis “ [Dhariwal & Nichol, 2021]

Summary

- Autoencoders (AE)
- Variational Autoencoders (VAE)
- Vector Quantized VAE (VQ-VAE)
- VQGAN
- Diffusion Models
- Epilogue: PNNs
- Bonus: ..



GPNN



“Drop the GAN: In Defense of Patches Nearest Neighbors as Single Image Generative Models” [Granot et al., 2021]

Diverse images generated from a single image

Input



Generated



Source Image

GPNN

SinGAN



“Drop the GAN: In Defense of Patches Nearest Neighbors as Single Image Generative Models” [Granot et al., 2021]

Diverse videos generated from a single video

↓ Original Video (the rest are generated)



Diffusion Models + CLIP



“a hedgehog using a calculator”



“a corgi wearing a red bowtie and a purple party hat”



“robots meditating in a vipassana retreat”



“a fall landscape with a small cottage next to a lake”



“a surrealist dream-like oil painting by salvador dali of a cat playing checkers”



“a professional photo of a sunset behind the grand canyon”



“a high-quality oil painting of a psychedelic hamster dragon”



“an illustration of albert einstein wearing a superhero costume”



Diffusion Models + CLIP

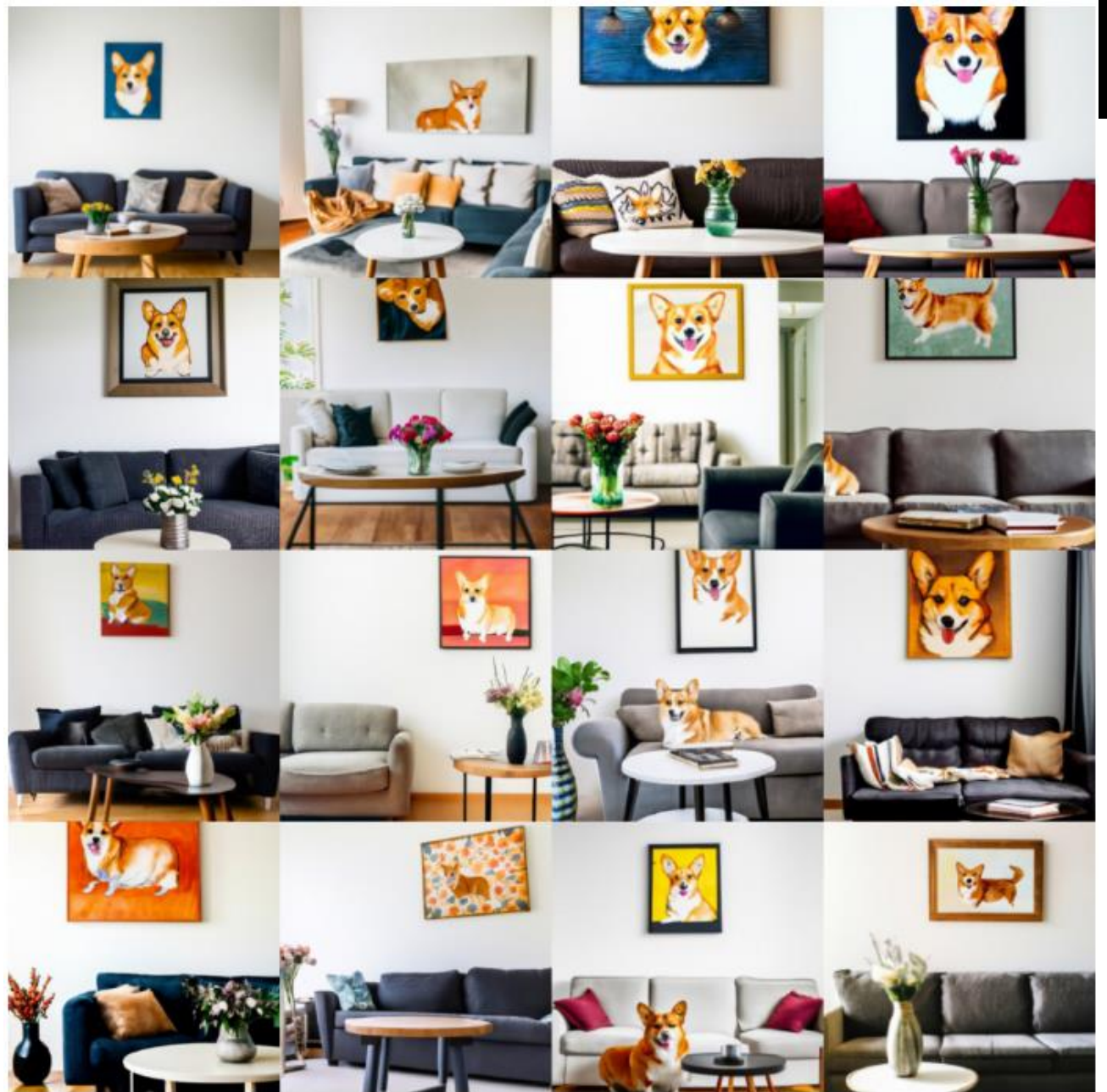


on the prompt “a stained glass window of a panda eating bamboo”.

Diffusion Models + CLIP

Prompt:

“A cozy living room with a painting of a corgi on the wall above a couch and a round coffee table in front of a couch and a vase of flowers on a coffee table”



Next Week:



Learning From Videos