



EECS 230 Deep Learning

Lecture 16: Generative Models

Some slides from Frank Noe, Simon Prince and Pascal Poupart

Outline

- ❑ Generative model
- ❑ Variational autoencoder
 - ❑ Autoencoder
 - ❑ Variational autoencoder
- ❑ Generative Adversarial Network
- ❑ Diffusion Model



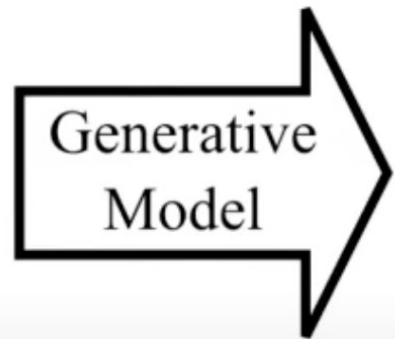
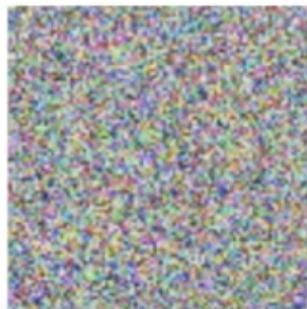
Generative model

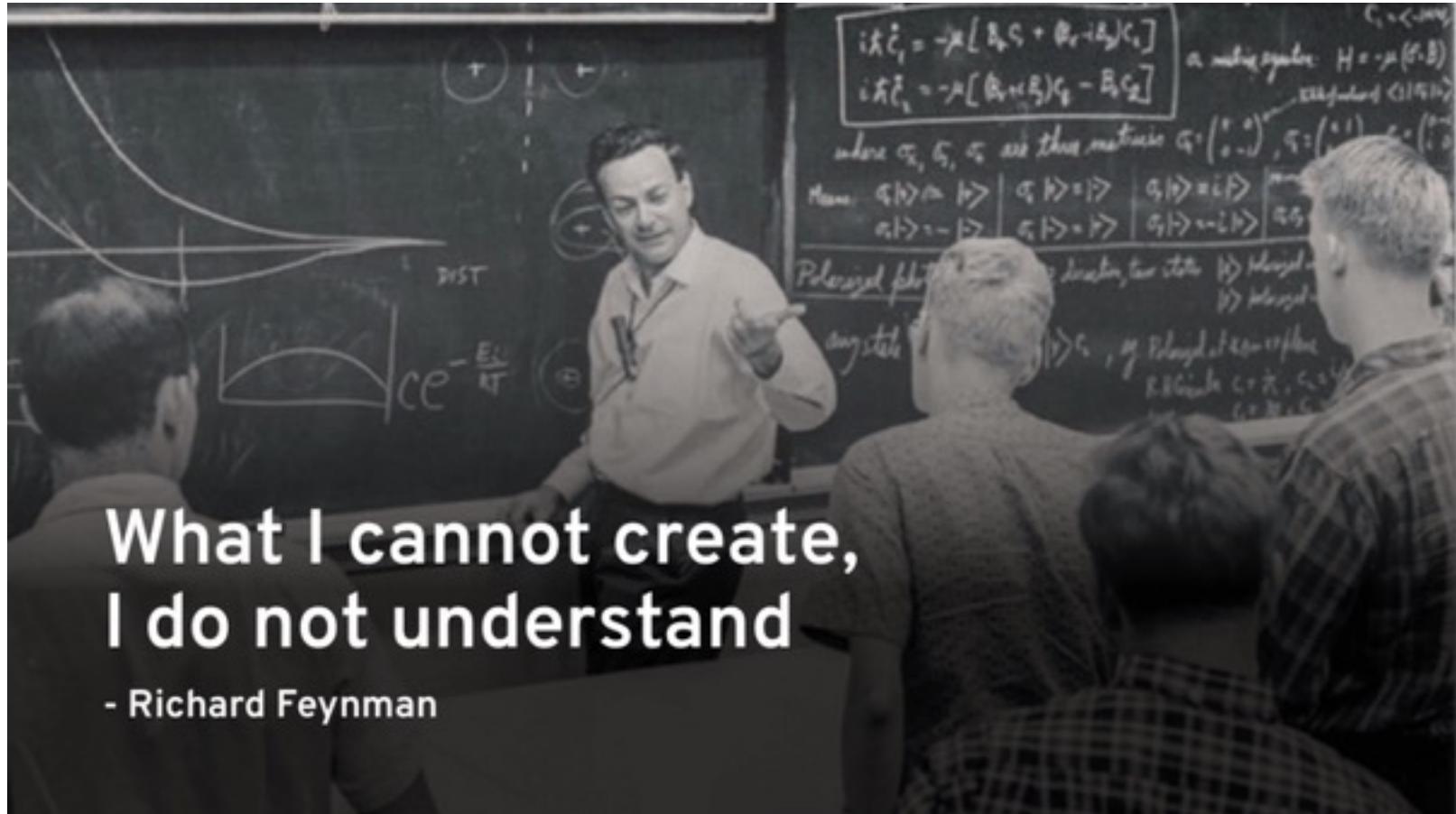
So far...

- ❑ Discriminative model $P(y|x)$
 - ❑ Given input data x , predict y
 - ❑ E.g., classification, regression
- ❑ Generative model
 - ❑ Model data distribution $P(x)$
 - ❑ Sample from $P(x)$ to generate new data

Generative model

Noise $\sim N(0,1)$





What I cannot create,
I do not understand

- Richard Feynman

Types of generative neural networks

- ❑ Boltzmann machines
- ❑ Sigmoid belief networks
- ❑ **Variational autoencoders (inference net + generator net)**
- ❑ **Generative adversarial networks (generator net + discriminator net)**
- ❑ Normalizing flows
- ❑ **Diffusion Models**
- ❑ ...



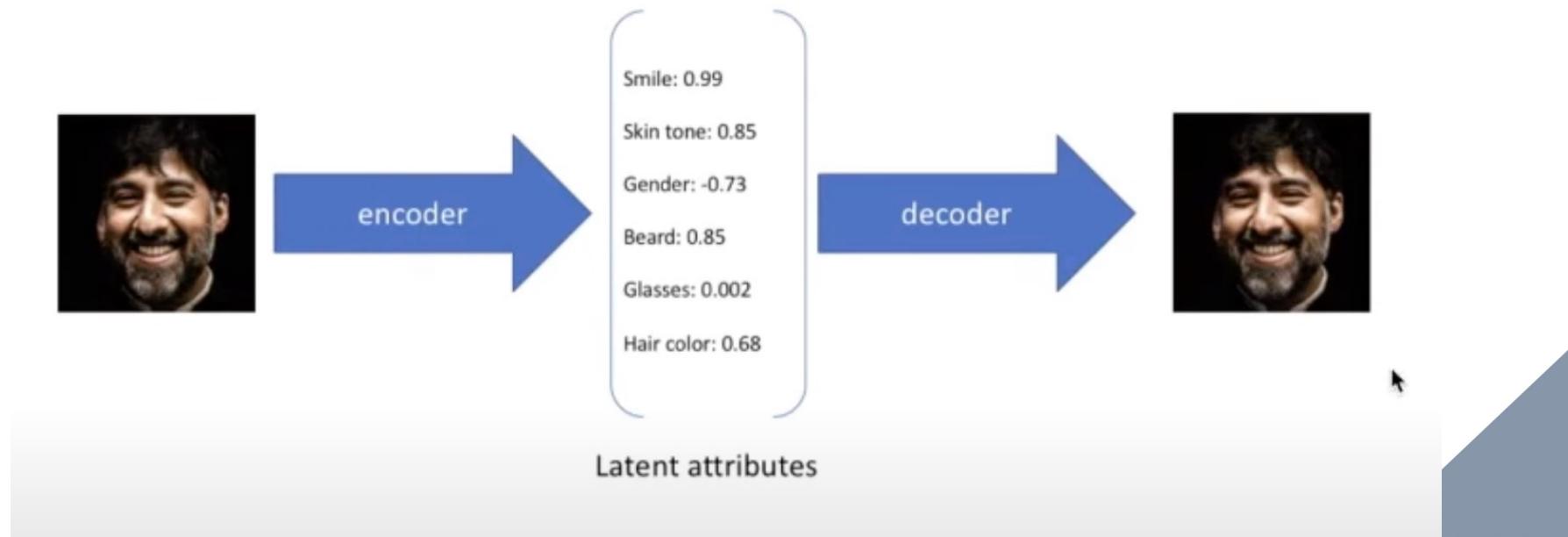
Variational Autoencoder

Autoencoder

- ❑ Special type of feed forward network for
 - ❑ Compression
 - ❑ Denoising
 - ❑ Sparse representation
 - ❑ Data generation

Autoencoder

- ❑ Encoder: $f(\cdot)$
- ❑ Decoder: $g(\cdot)$
- ❑ Autoencoder: $g(f(x)) = x$



Latent variable encodes “essential” information about input data

Linear autoencoder

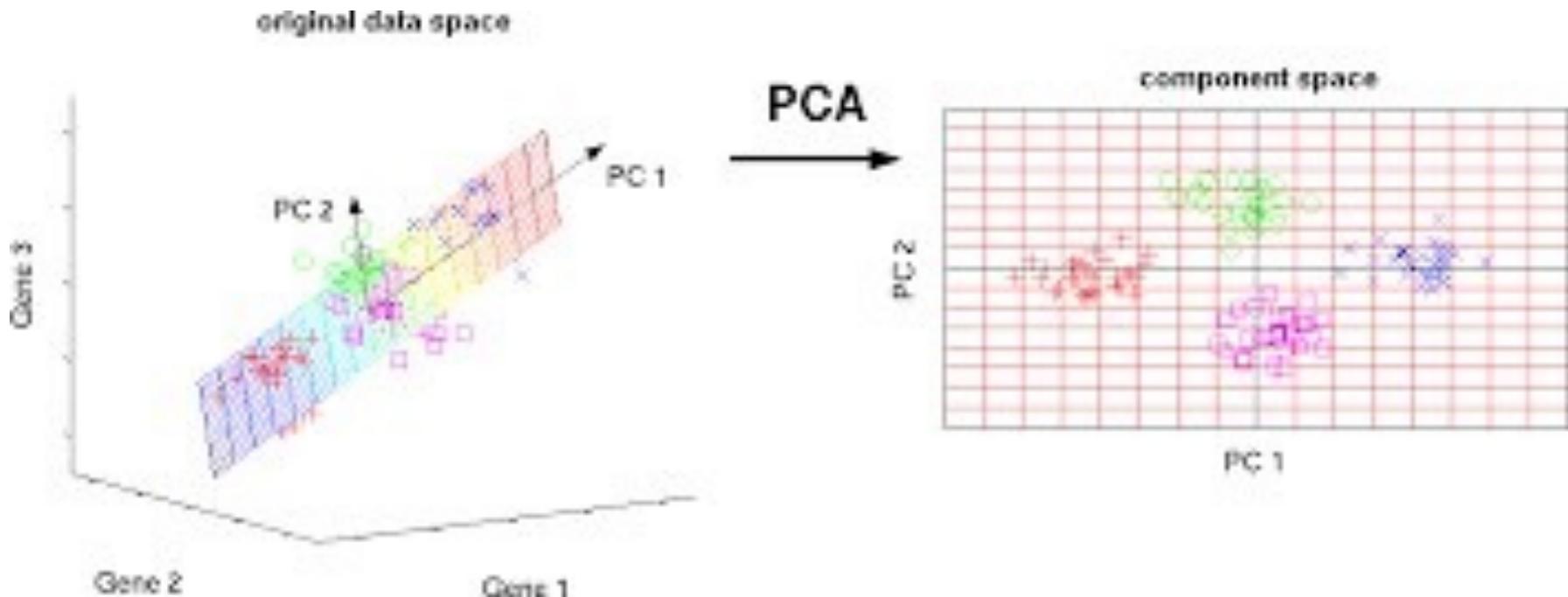
- ❑ Objective: find weights W_f and W_g that minimize reconstruction error

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| \mathbf{W}_g \mathbf{W}_f \mathbf{x}_n - \mathbf{x}_n \right\|_2^2$$

- ❑ When using Euclidean norm (i.e., squared loss), solution is the same as principal component analysis (PCA)

Recap: principle component analysis

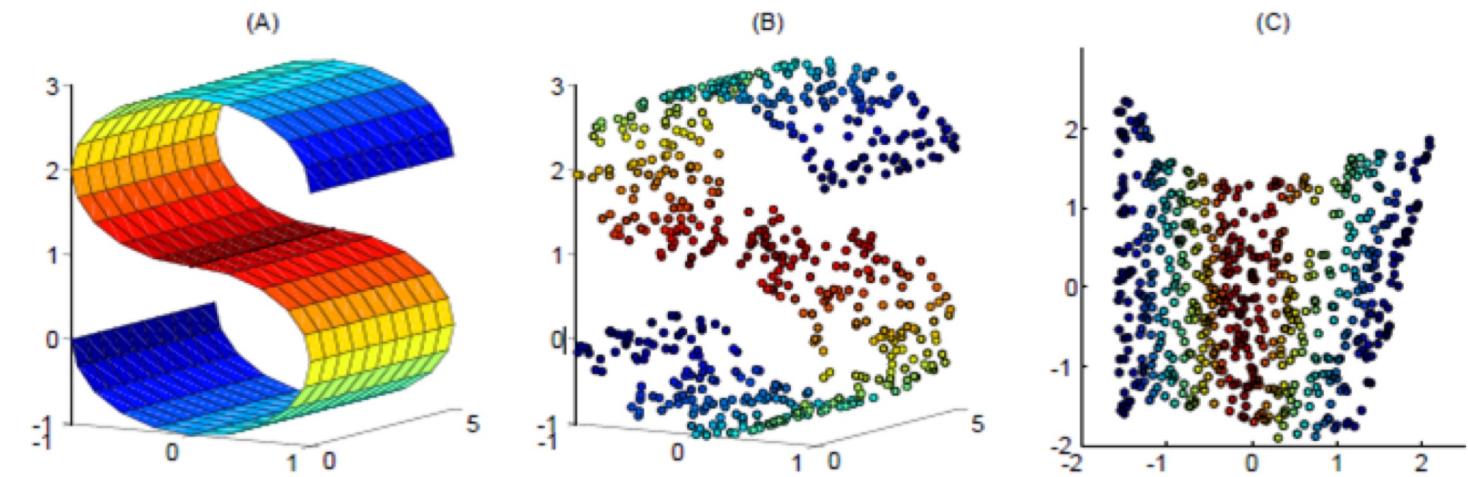
- Components with maximum variance



Non-linear autoencoder

◻ $f(\cdot)$ and $g(\cdot)$ are both non-linear functions

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| g(f(x_n; \mathbf{W}_f); \mathbf{W}_g) - x_n \right\|_2^2$$



Sparse representation

- When more hidden nodes than inputs, use regularization to constrain autoencoder
- Example: force hidden nodes to be sparse

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| g(f(\mathbf{x}_n; \mathbf{W}_f); \mathbf{W}_g) - \mathbf{x}_n \right\|_2^2 + c \underbrace{\text{nnz}\left(f(\mathbf{x}_n; \mathbf{W}_f)\right)}_{\text{Sparse hidden nodes}}$$

where $\text{nnz}\left(f(\mathbf{x}_n; \mathbf{W}_f)\right)$ is the number of non-zero entries in the vector produced by f .

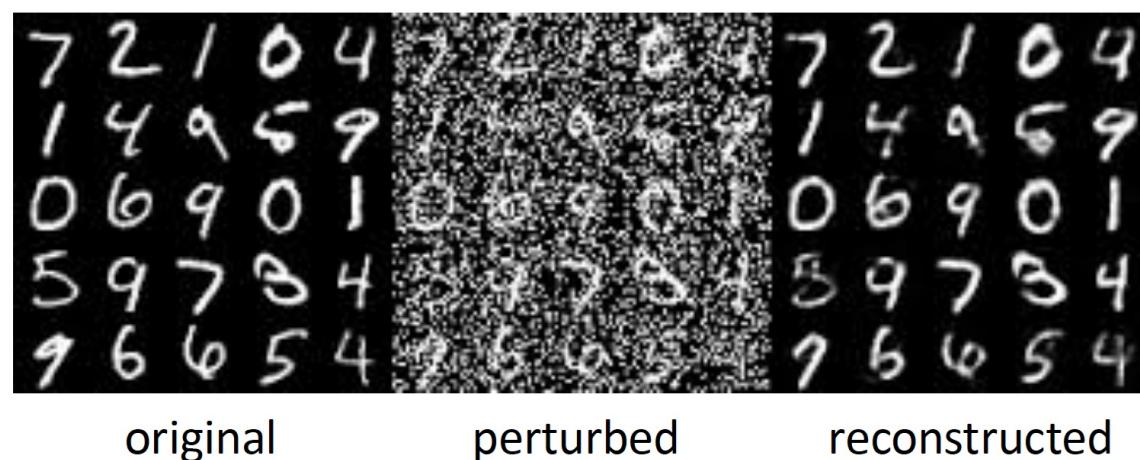
- Approximate objective: L_1 regularization

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| g(f(\mathbf{x}_n; \mathbf{W}_f); \mathbf{W}_g) - \mathbf{x}_n \right\|_2^2 + c \left\| f(\mathbf{x}_n; \mathbf{W}_f) \right\|_1$$

Denoising autoencoder

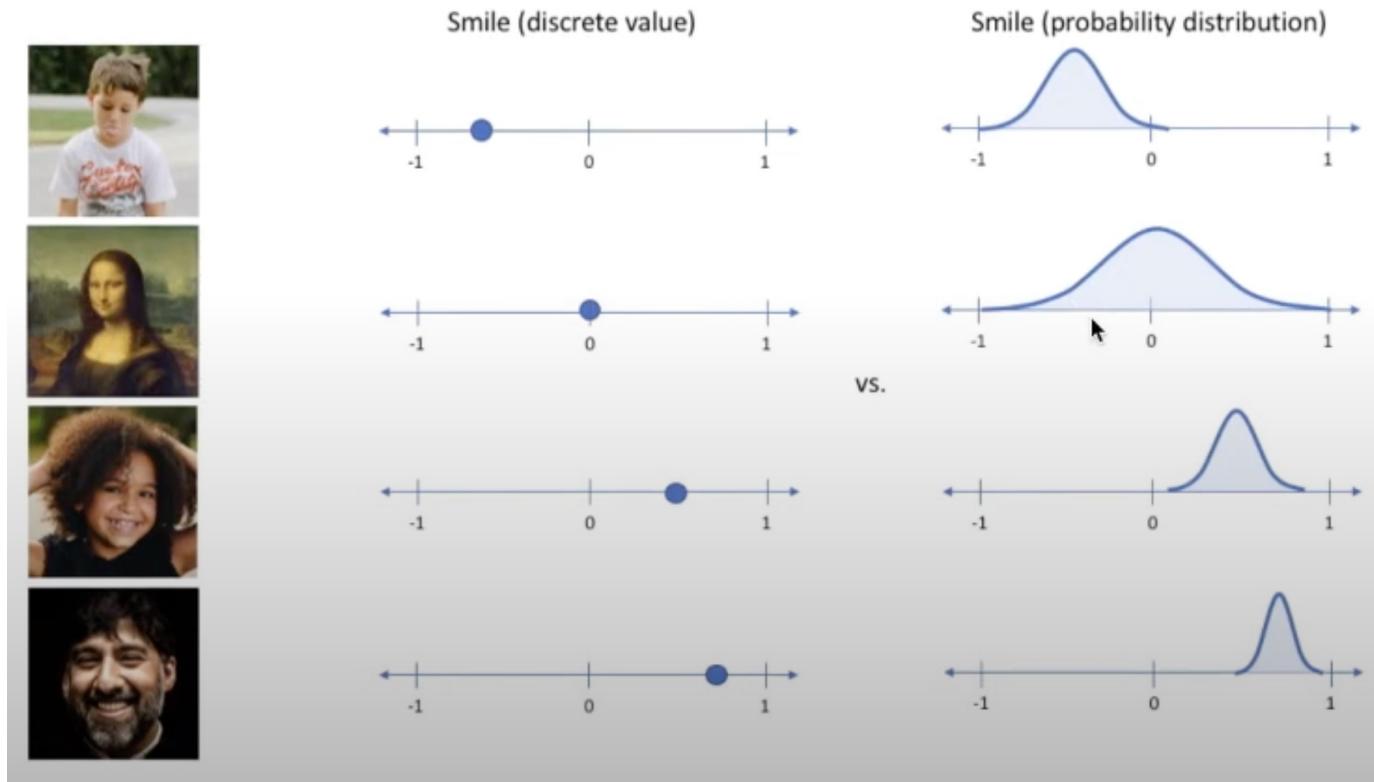
- Consider noisy version \tilde{x} of the input x

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| g(f(\tilde{\mathbf{x}}_n; \mathbf{W}_f); \mathbf{W}_g) - \mathbf{x}_n \right\|_2^2 + c \left\| f(\tilde{\mathbf{x}}_n; \mathbf{W}_f) \right\|_1$$

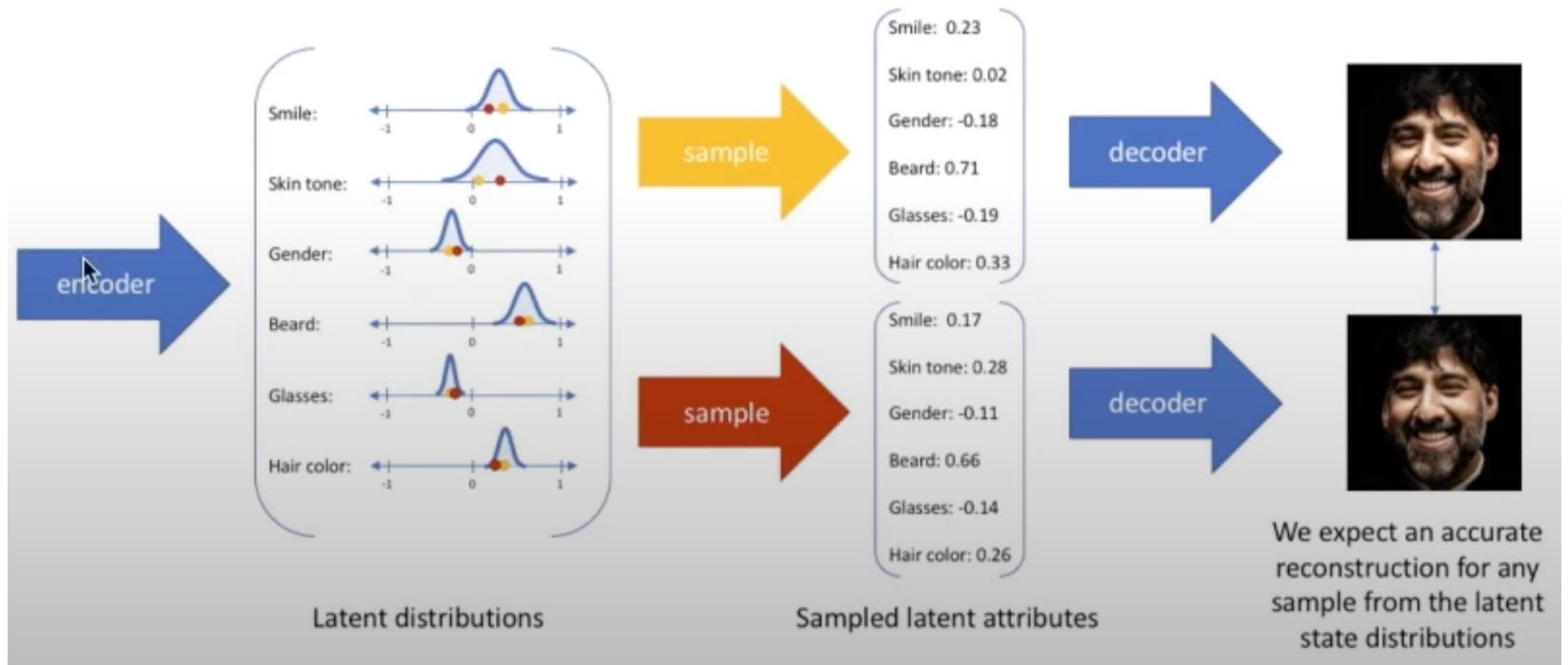


Variational/Probabilistic autoencoder

- ❑ Instead of a single value for each attribute, represent each attribute as a range of values
- ❑ VAE: Describe latent attribute in probabilistic terms



Variational/Probabilistic autoencoder



Variational/Probabilistic autoencoder

- ❑ Let $f()$ and $g()$ represent conditional distributions
 - ❑ $f: \Pr(h|x; w_f)$
 - ❑ $g: \Pr(x|h; w_g)$
- ❑ The decoder $g()$ can be treated as a generative model
 - ❑ First sample h from $\Pr(h)$
 - ❑ Then sample x from $\Pr(x|h; w_g)$

Variational autoencoder

- ❑ Idea: train encoder $\Pr(\mathbf{h}|\mathbf{x}; \mathbf{w}_f)$ to approach a simple and fixed distribution, e.g., $N(\mathbf{h}; \mathbf{0}, \mathbf{I})$

$$\max_{\mathbf{W}} \sum_n \log \Pr(\mathbf{x}_n; \mathbf{W}_f, \mathbf{W}_g) - c \underbrace{\text{KL}(\Pr(\mathbf{h}|\mathbf{x}_n; \mathbf{W}_f) || N(\mathbf{h}; \mathbf{0}, \mathbf{I}))}_{\text{Kullback-Leibler divergence}}$$

Kullback-Leibler divergence
Distance measure for distributions

Variational Autoencoder Likelihood

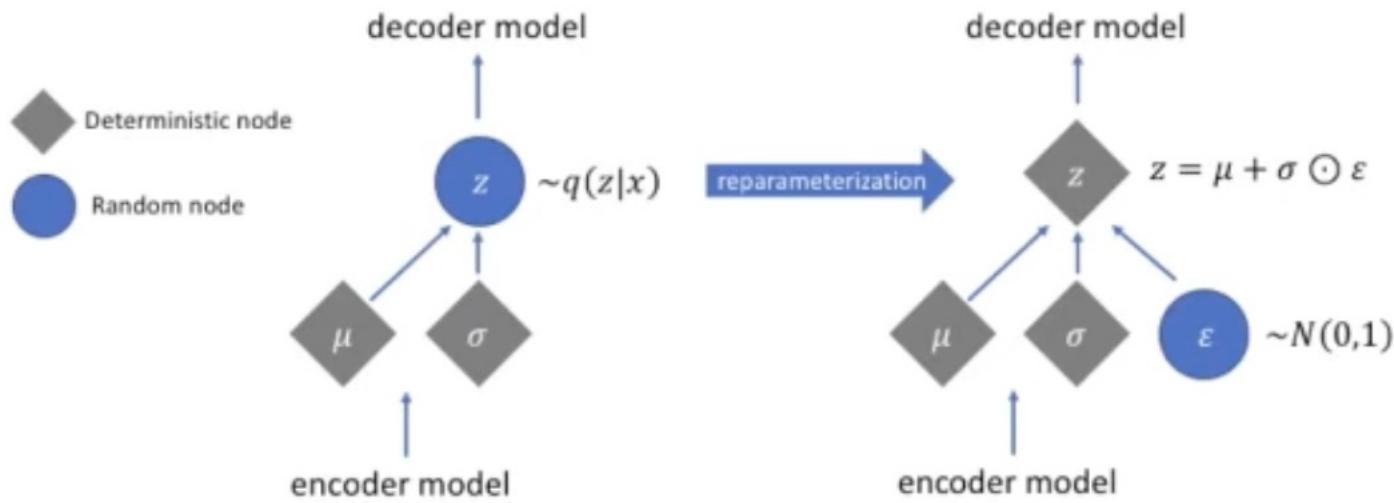
- ❑ How to compute $\Pr(x_n; W_f, W_g)$?

$$\Pr(x_n; W_f, W_g) = \int_h \Pr(x_n | h; W_g) \Pr(h | x_n; W_f) dh$$

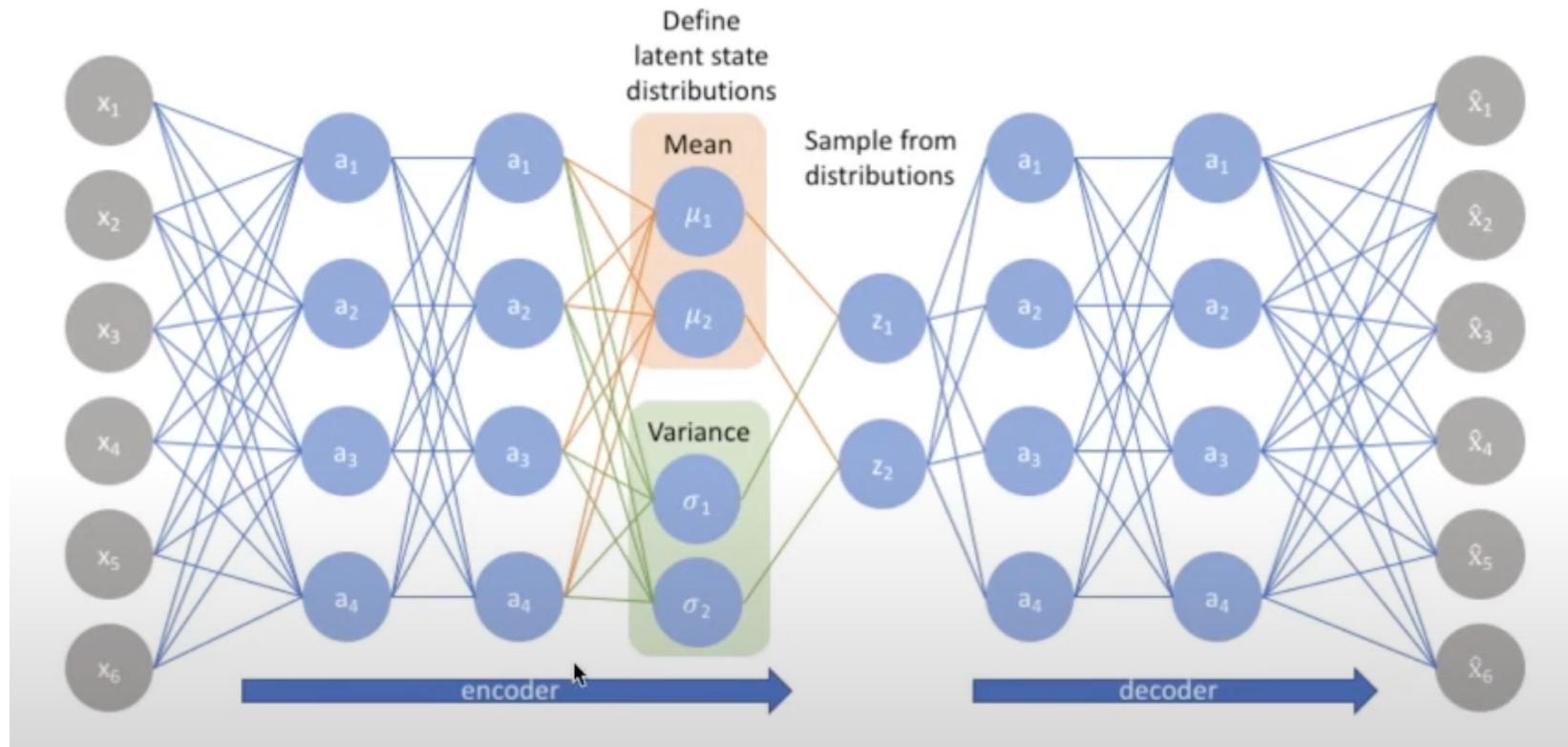
- ❑ Obtain mean and variance of $\Pr(h)$ by neural network

$$\Pr(h | x_n; W_f) = N(h; \mu_n(x_n; W_f), \sigma_n(x_n; W_f)I)$$

Reparameterization trick



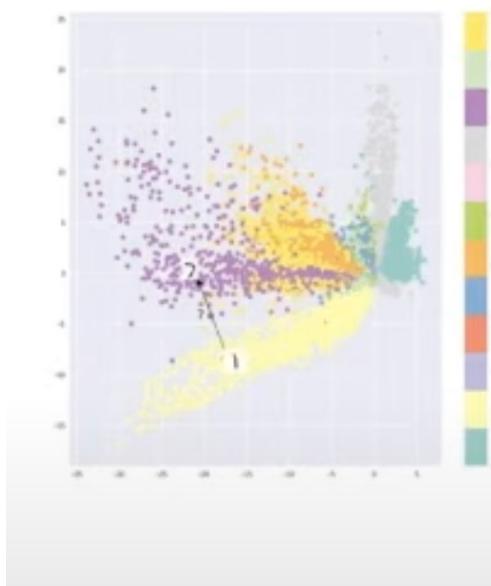
VAE implementation



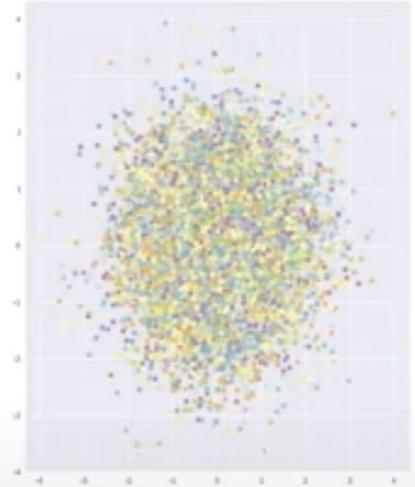
MNIST VAE

$$\max_{\mathbf{W}} \sum_n \log \Pr(x_n; \mathbf{W}_f, \mathbf{W}_g) - c \underbrace{\text{KL}(\Pr(\mathbf{h}|x_n; \mathbf{W}_f) || N(\mathbf{h}; \mathbf{0}, \mathbf{I}))}_{\text{Kullback-Leibler divergence}}$$

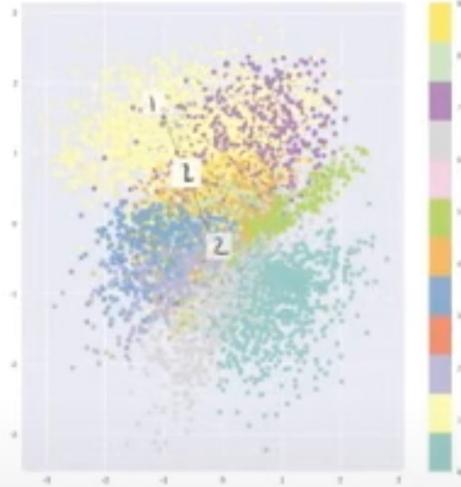
Only reconstruction loss



Only KL divergence



Combination



Examples from VAE

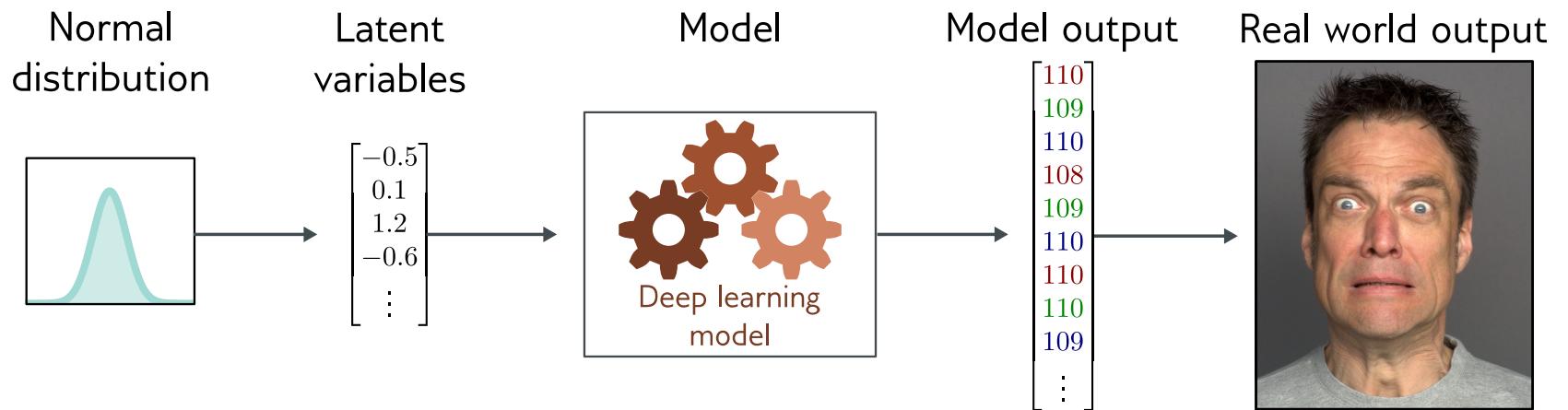


6 6 6 6 6 6 6 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
4 4 4 4 2 2 2 2 2 2 2 0 0 0 0 0 0 0 0 0 0 0 0 0
4 2 2 2 2 2 2 2 2 2 2 2 5 5 5 5 0 0 0 0 0 0 0 0
4 9 2 2 2 2 2 2 2 2 2 2 3 3 3 3 5 5 5 5 5 5 5 5 3
4 9 4 2 2 2 2 2 2 3 3 3 3 5 5 5 5 5 5 5 5 5 3
4 9 9 2 2 2 3 3 3 3 3 3 5 5 5 5 5 5 5 5 5 3
4 9 9 9 2 2 3 3 3 3 3 3 5 5 5 5 5 5 5 5 3
7 9 9 9 9 9 8 3 3 3 3 3 3 3 5 5 5 5 8 7
7 9 9 9 9 9 8 3 3 3 3 3 3 3 3 3 3 3 8 8 8 8 7
7 9 9 9 9 9 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 7
7 9 9 9 9 9 8 8 8 8 8 8 8 8 8 6 6 6 6 6 5 7
7 9 9 9 9 9 9 8 8 8 8 8 8 8 8 6 6 6 6 6 5 7
7 9 9 9 9 9 9 9 8 5 5 5 5 6 6 6 6 6 6 5 7
7 9 9 9 9 9 9 9 9 5 5 5 5 6 6 6 6 6 6 6 5 7
7 9 9 9 9 9 9 9 9 9 5 5 5 6 6 6 6 6 6 6 5 7
7 9 9 9 9 9 9 9 9 9 9 1 1 6 6 6 6 6 6 6 6 7
7 9 9 9 9 9 9 9 9 9 9 1 1 6 6 6 6 6 6 6 6 7
7 9 9 9 9 9 9 9 9 9 1 1 1 1 1 1 1 1 1 1 1 1 1
7 9 9 9 9 9 9 9 9 9 1 1 1 1 1 1 1 1 1 1 1 1 1
7 9 9 9 9 9 9 9 9 9 1 1 1 1 1 1 1 1 1 1 1 1 1
7 7



Generative Adversarial Network

Latent variable models



Latent variable models map a random “latent” variable to create a new data sample

In GANs, this is called the **generator**:

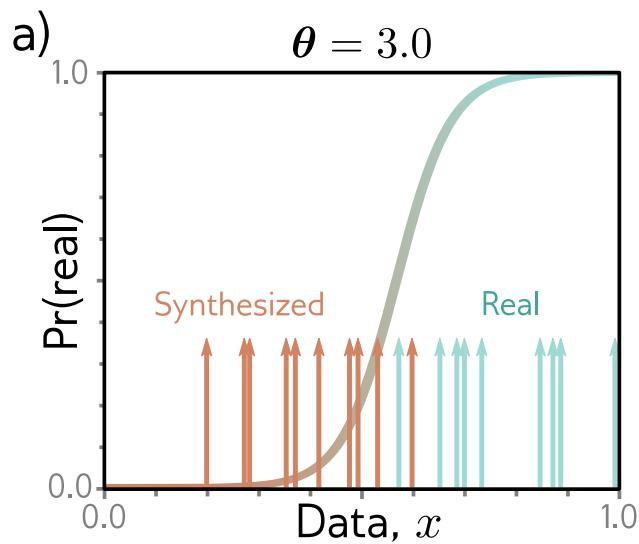
$$\mathbf{x}_* = \mathbf{g}[\mathbf{z}_i, \theta]$$

Generative adversarial models

- GAN loss function
- DCGAN results and problems
- Tricks for improving performance
- Conditional GANs
- Image translation models
- StyleGAN

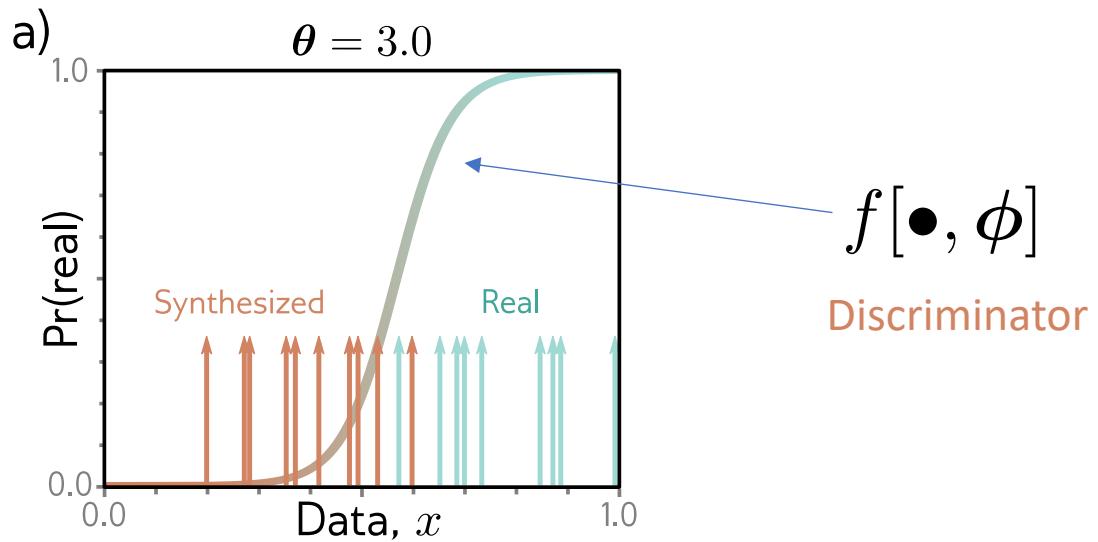
GAN example

$$x_j^* = g[z_j, \theta] = z_j + \theta$$



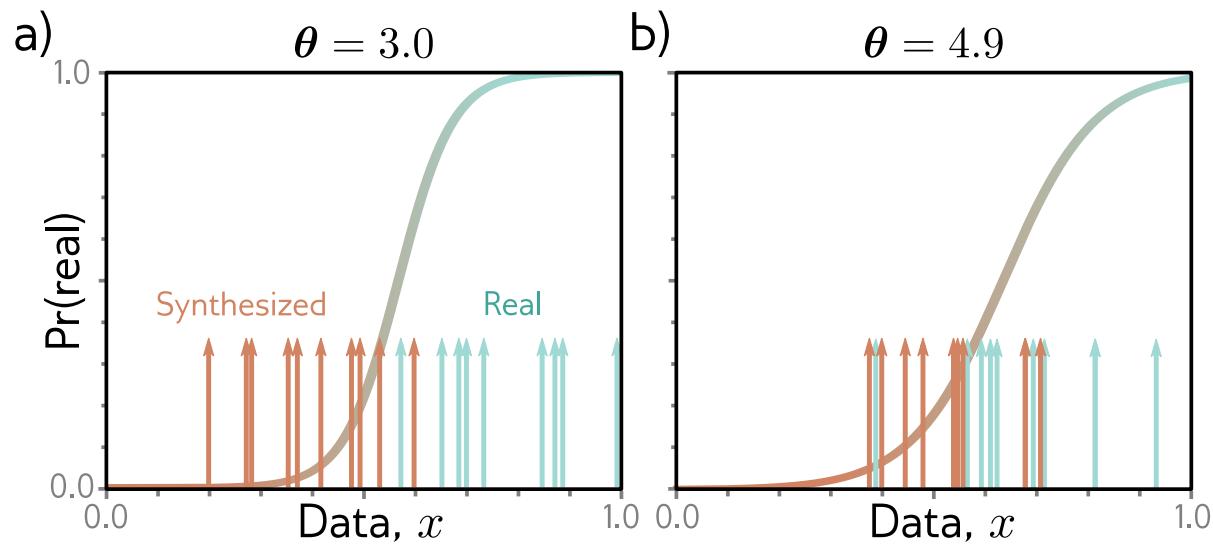
GAN example

$$x_j^* = g[z_j, \theta] = z_j + \theta$$



GAN example

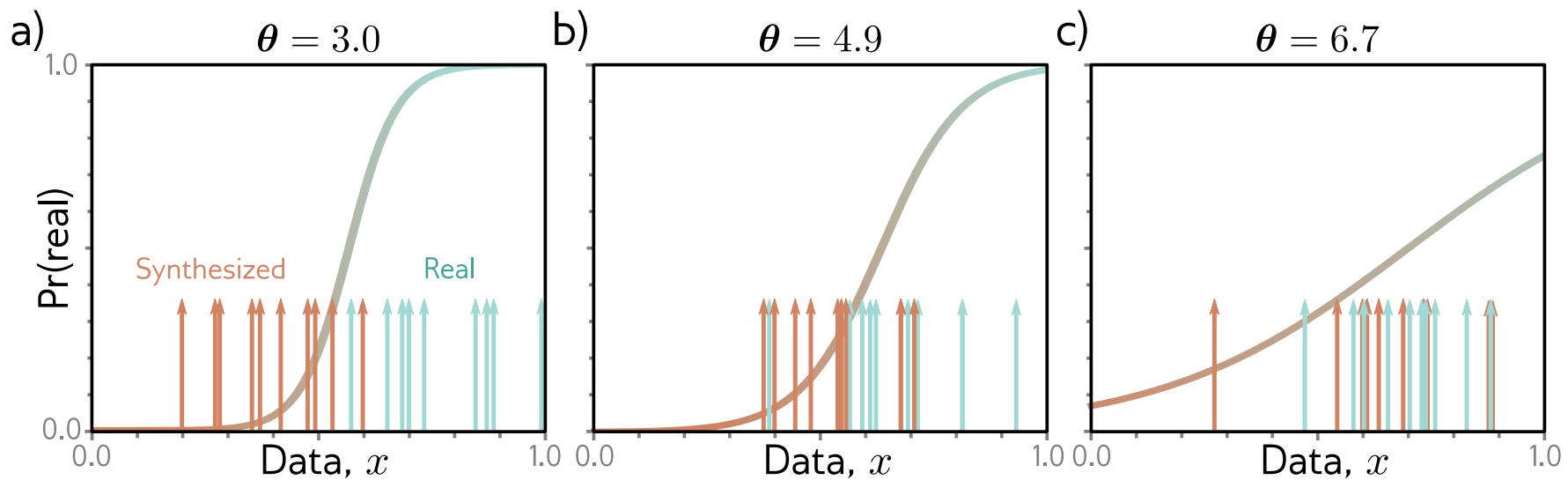
$$x_j^* = g[z_j, \theta] = z_j + \theta$$



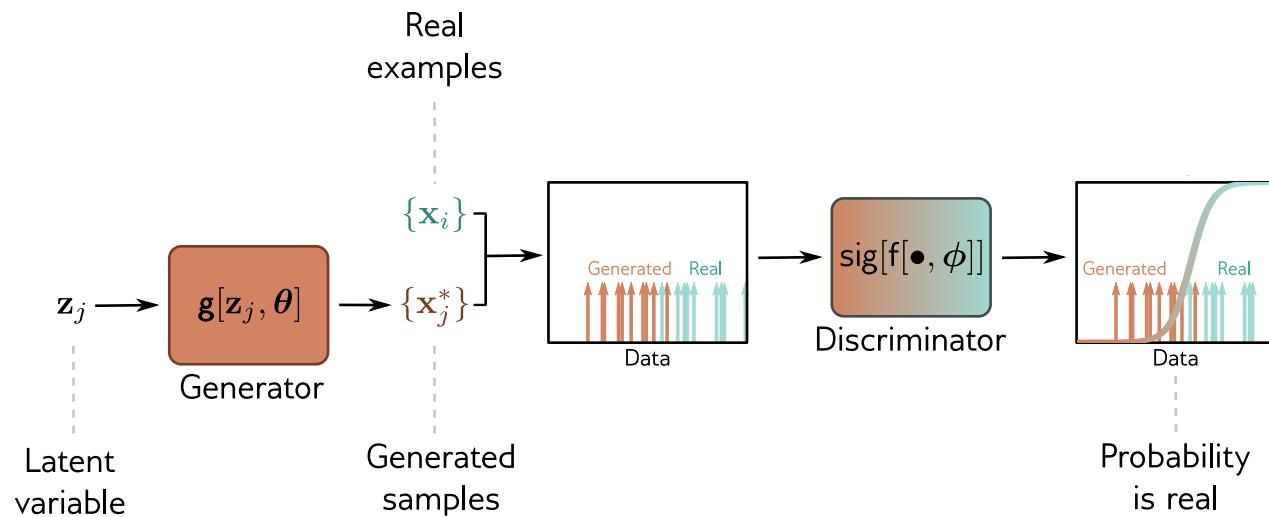
Slide credit: Simon Prince

GAN example

$$x_j^* = g[z_j, \theta] = z_j + \theta$$



GAN



Slide credit: Simon Prince

GAN cost function

Discriminator uses standard cross entropy loss:

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[\sum_i -(1 - y_i) \log [1 - \operatorname{sig}[f[\mathbf{x}_i, \phi]]] - y_i \log [\operatorname{sig}[f[\mathbf{x}_i, \phi]]] \right]$$

Discriminator: generated samples, $y = 0$, real examples, $y = 1$:

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[\sum_j -\log [1 - \operatorname{sig}[f[\mathbf{x}_j^*, \phi]]] - \sum_i \log [\operatorname{sig}[f[\mathbf{x}_i, \phi]]] \right]$$

Generator loss: make generated samples more likely under discriminator (i.e. make discriminator loss larger)

$$\hat{\phi}, \hat{\theta} = \operatorname{argmin}_{\phi} \left[\operatorname{argmax}_{\theta} \left[\sum_j -\log [1 - \operatorname{sig}[f[\mathbf{g}[\mathbf{z}_j, \theta], \phi]]] - \sum_i \log [\operatorname{sig}[f[\mathbf{x}_i, \phi]]] \right] \right]$$

GAN Cost function

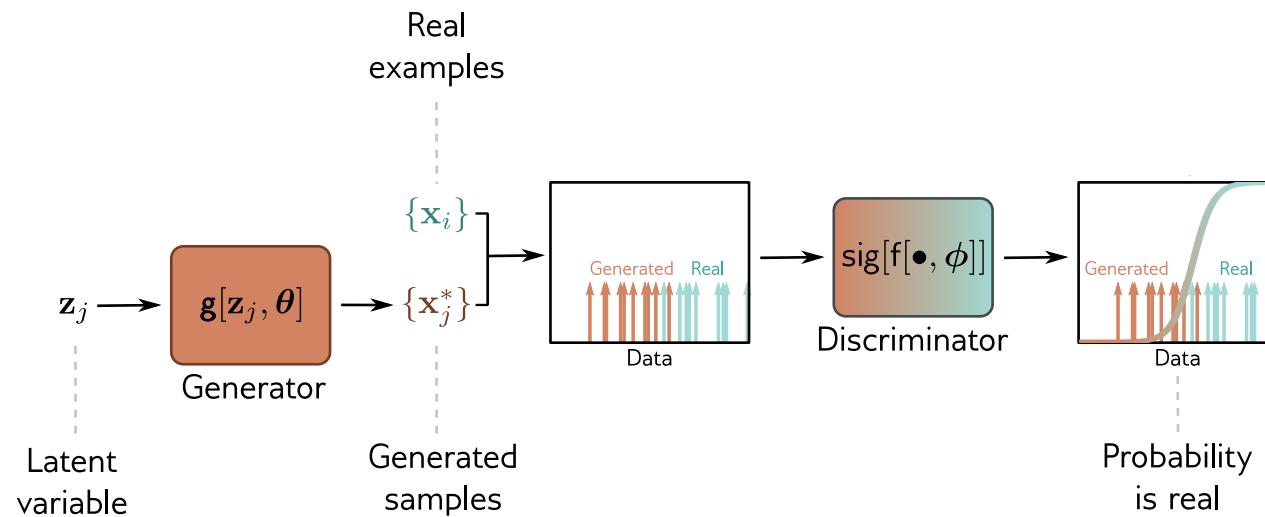
$$\hat{\phi}, \hat{\theta} = \operatorname{argmin}_{\phi} \left[\operatorname{argmax}_{\theta} \left[\sum_j -\log \left[1 - \operatorname{sig}[f[\mathbf{g}[\mathbf{z}_j, \theta], \phi]] \right] - \sum_i \log \left[\operatorname{sig}[f[\mathbf{x}_i, \phi]] \right] \right] \right]$$

Can divide into two parts:

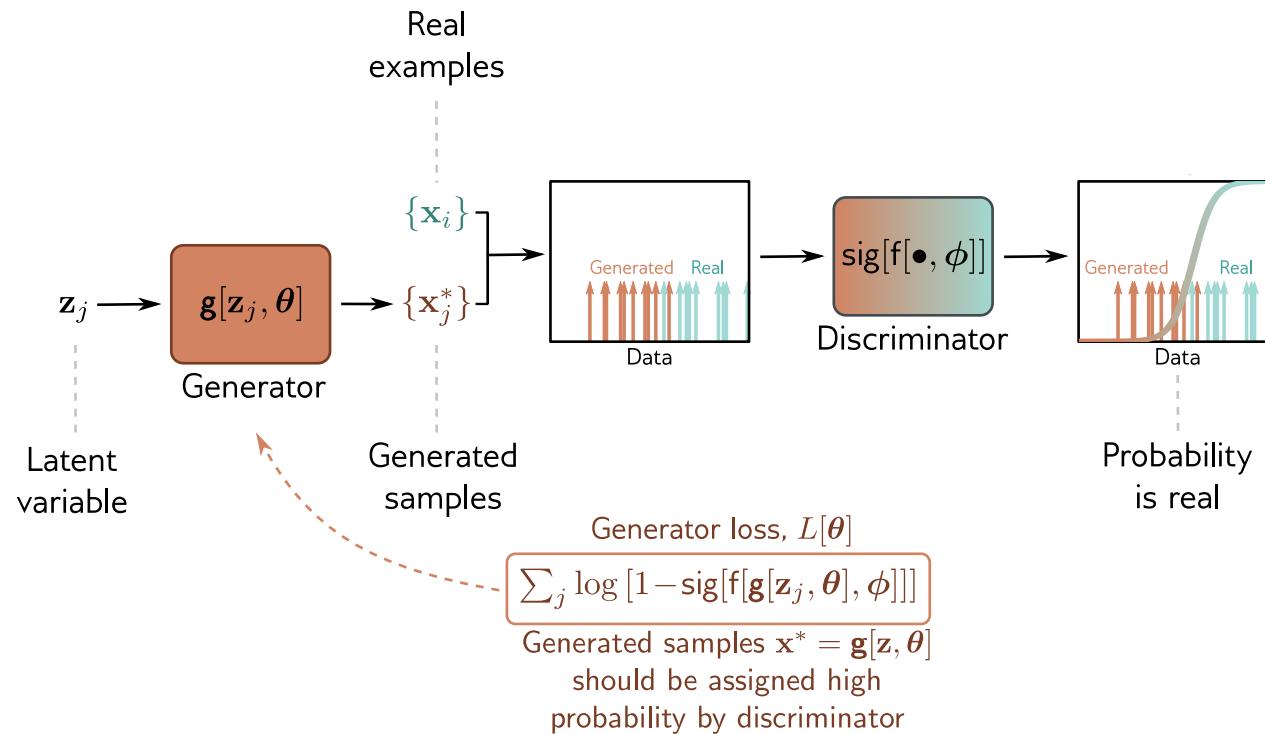
$$L[\phi] = \sum_j -\log \left[1 - \operatorname{sig}[f[\mathbf{g}[\mathbf{z}_j, \theta], \phi]] \right] - \sum_i \log \left[\operatorname{sig}[f[\mathbf{x}_i, \phi]] \right]$$

$$L[\theta] = \sum_j \log \left[1 - \operatorname{sig}[f[\mathbf{g}[\mathbf{z}_j, \theta], \phi]] \right]$$

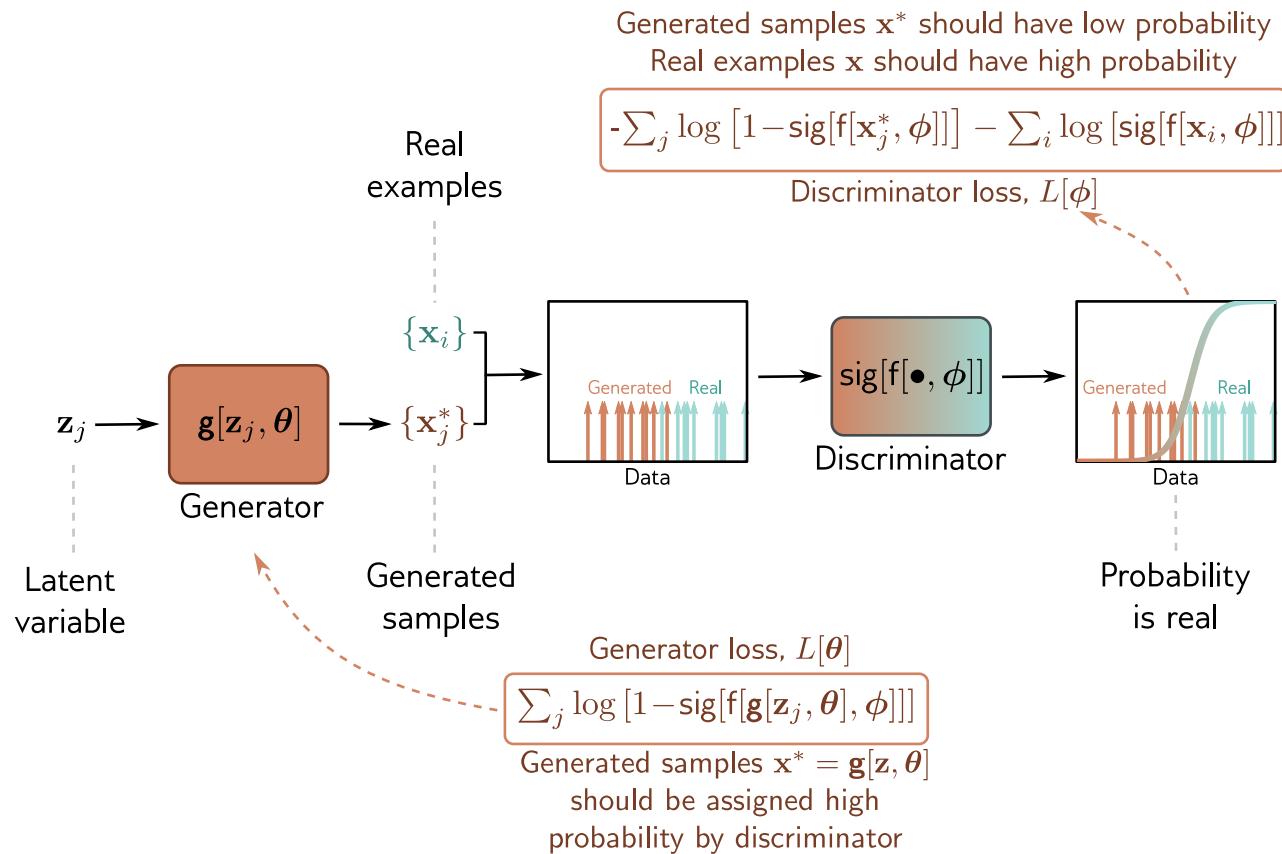
GAN loss function



GAN loss function



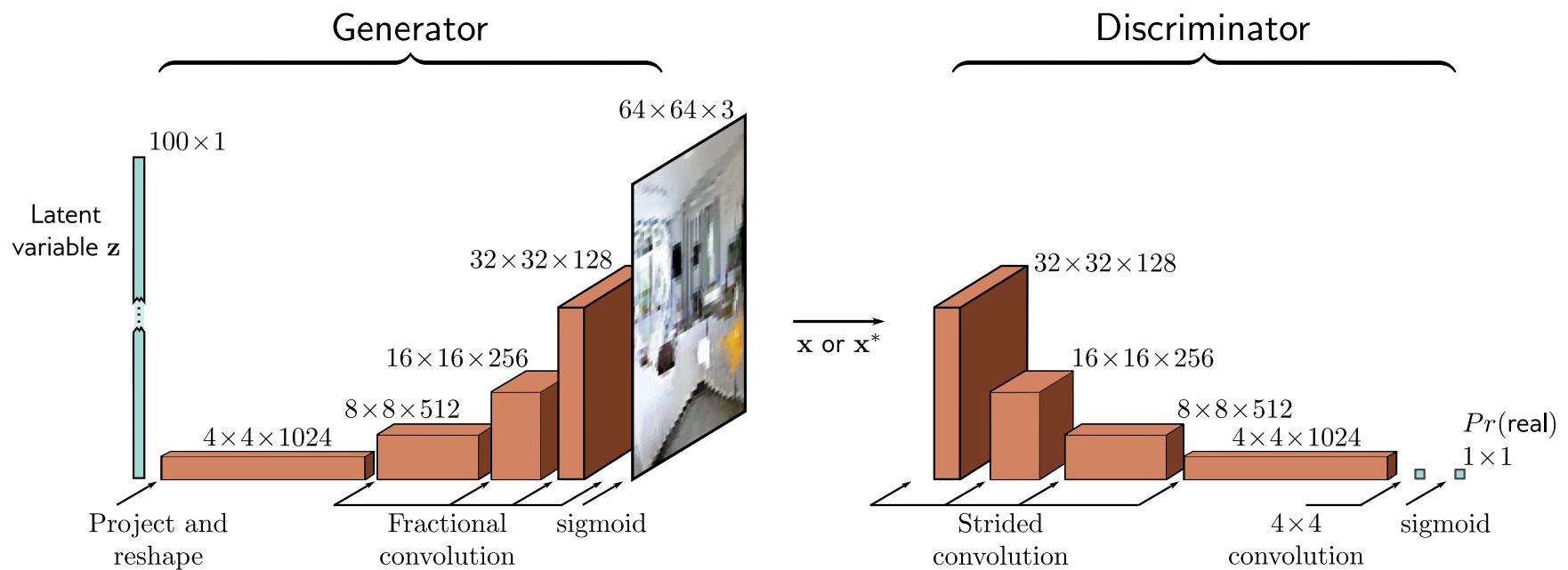
GAN loss function



Generative adversarial models

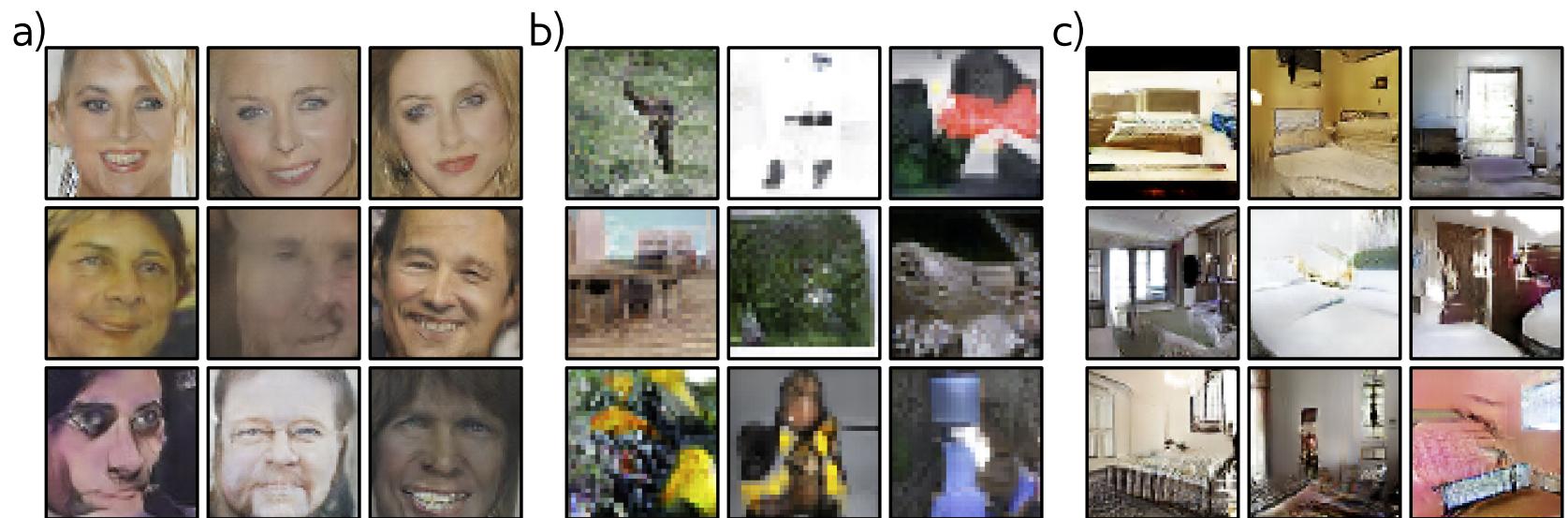
- GAN loss function
- DCGAN results and problems
- Tricks for improving performance
- Conditional GANs
- Image translation models

Deep Convolutional (DC) GAN



Slide credit: Simon Prince

DC GAN Results



Slide credit: Simon Prince

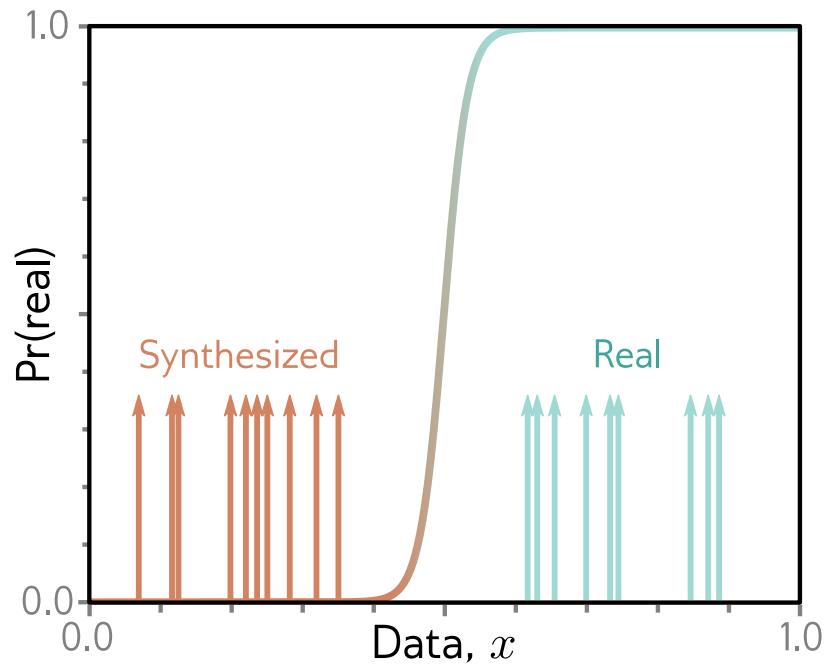
Mode collapse

Training GANs is difficult. Often fails.



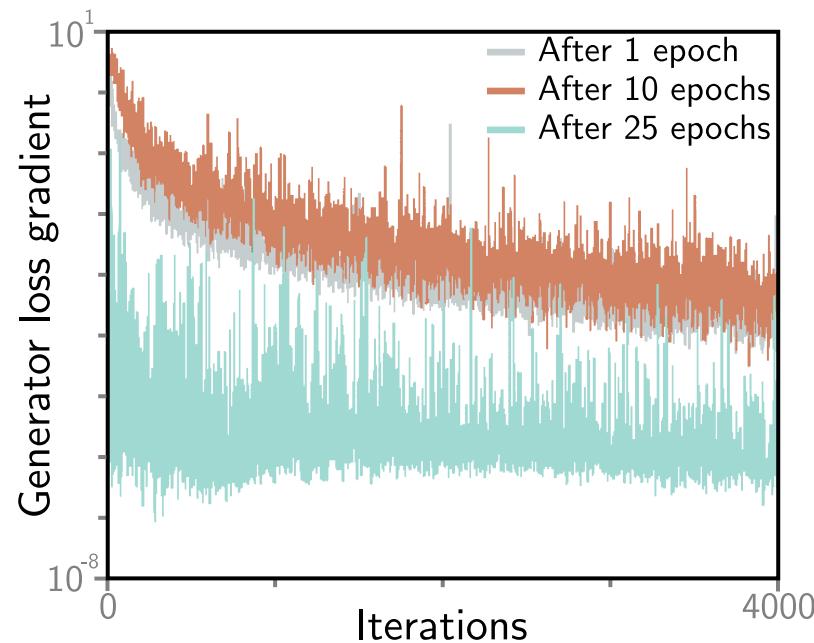
Slide credit: Simon Prince

GAN training instability



Slide credit: Simon Prince

GAN training instability



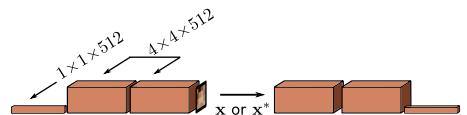
Fix generator and continue training discriminator, then generator gradients disappear.

Generative adversarial models

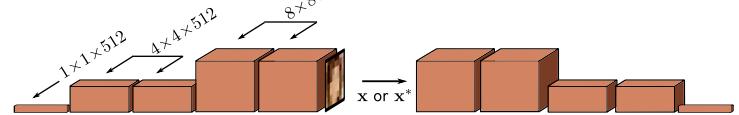
- GAN loss function
- DCGAN results and problems
- Tricks for improving performance
- Conditional GANs
- Image translation models

Trick 1: Progressive growing

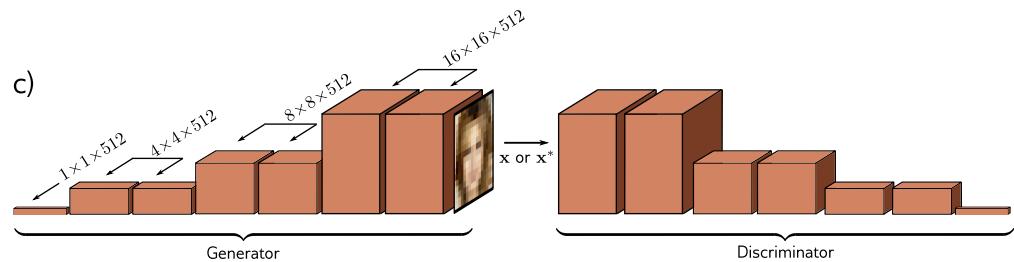
a)



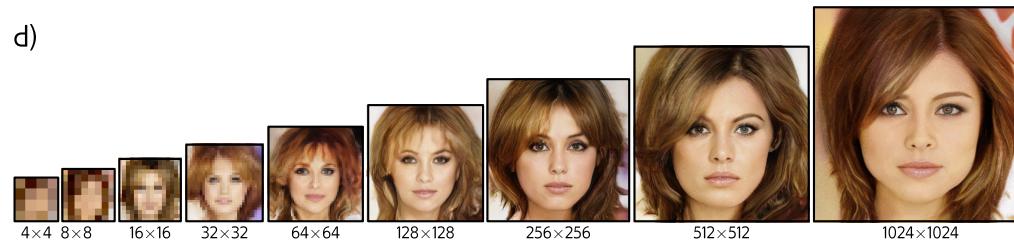
b)



c)



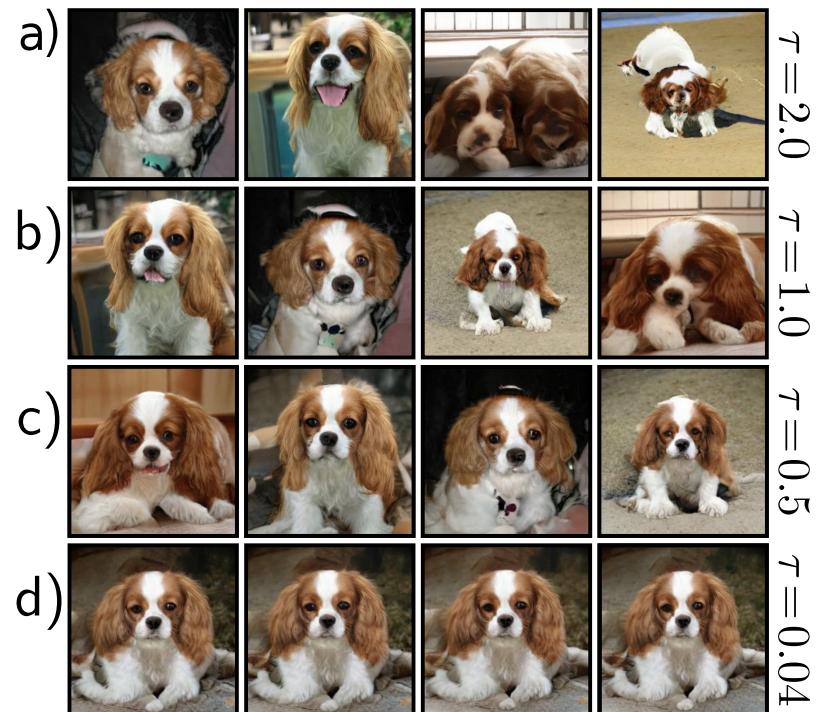
d)



Slide credit: Simon Prince

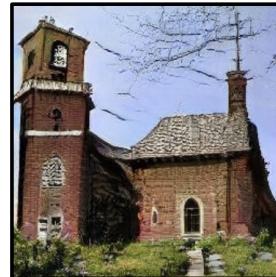
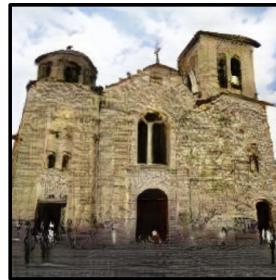
Trick 2: Truncation

Only choose random values of latent variables that are less than a threshold τ .



Slide credit: Simon Prince

Example results



Slide credit: Simon Prince

Interpolation



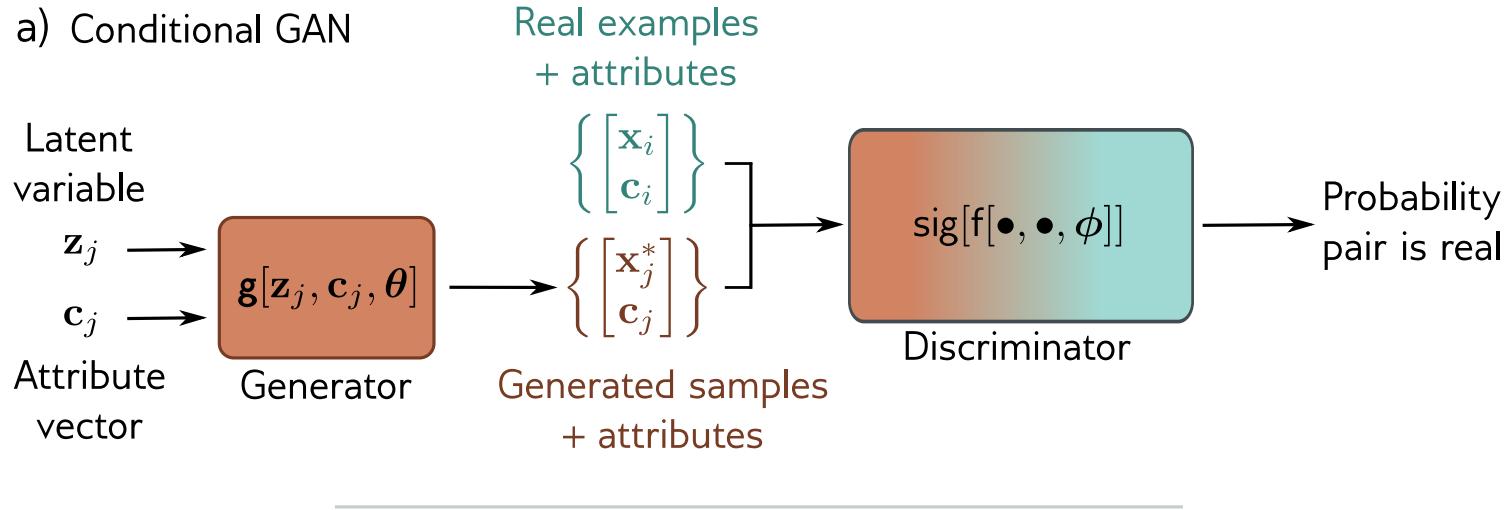
Slide credit: Simon Prince

Generative adversarial models

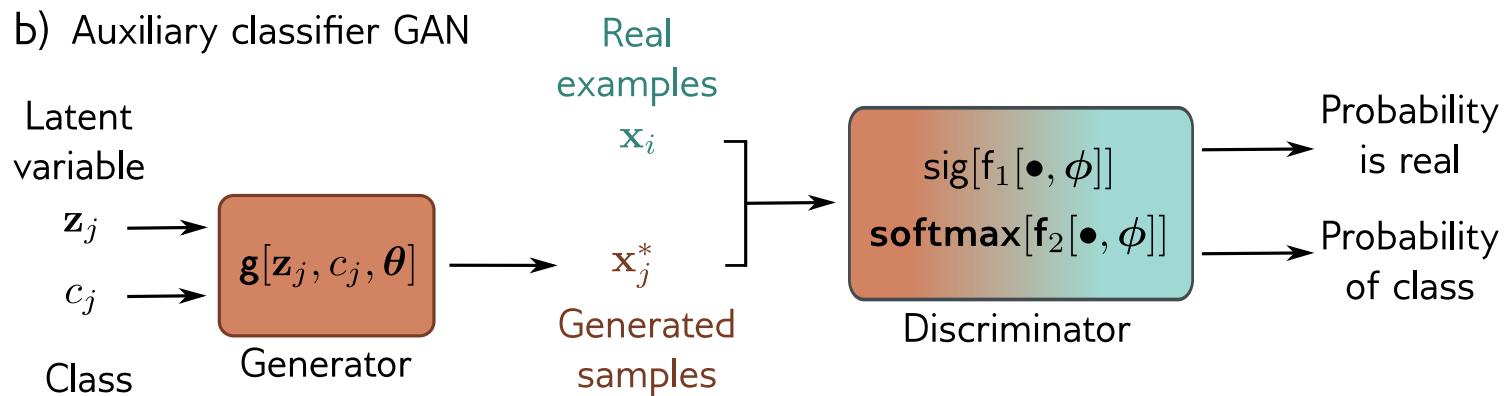
- GAN loss function
- DCGAN results and problems
- Tricks for improving performance
- Conditional GANs
- Image translation models

Conditional GAN models

a) Conditional GAN



b) Auxiliary classifier GAN



Slide credit: Simon Prince

ACGAN results

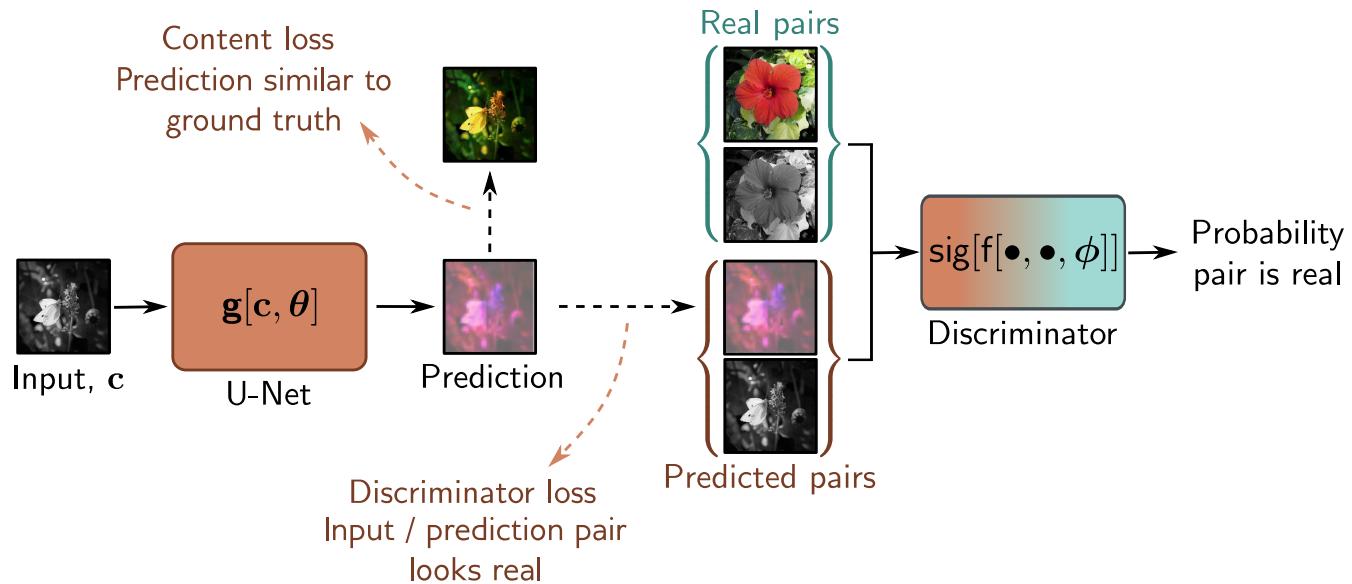


Slide credit: Simon Prince

Generative adversarial models

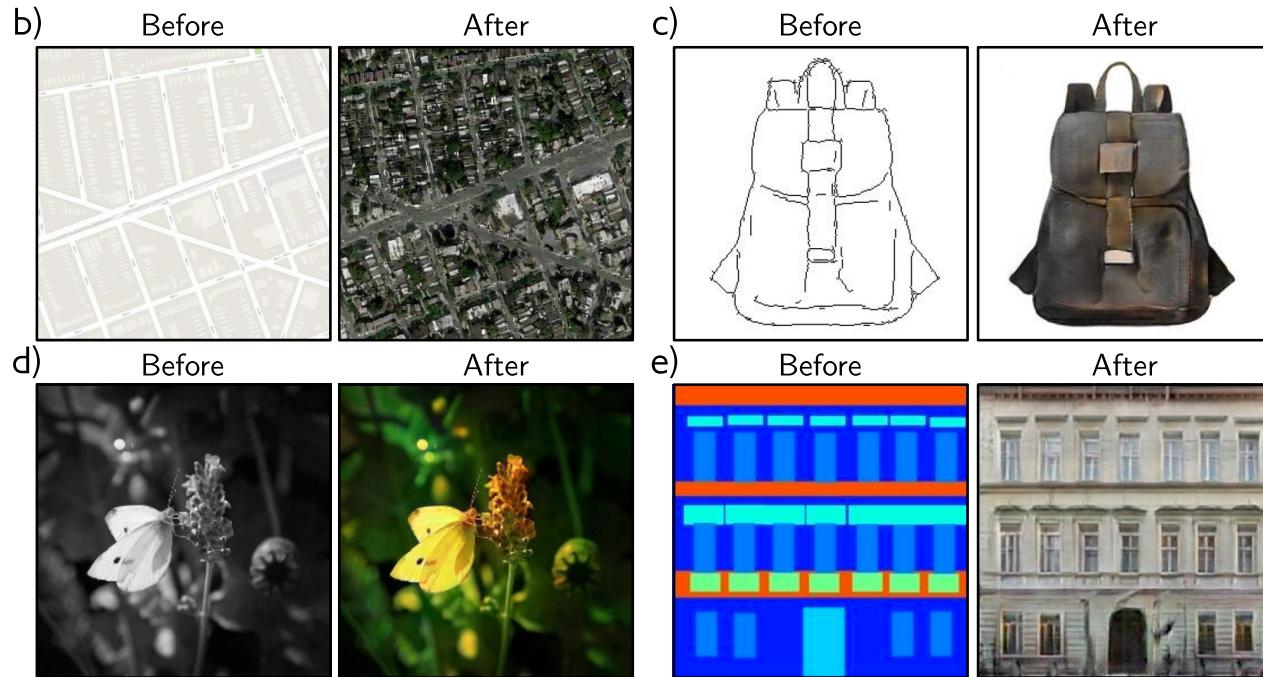
- GAN loss function
- DCGAN results and problems
- Tricks for improving performance
- Conditional GANs
- Image translation models

Image translation: Pix2Pix



Slide credit: Simon Prince

Image translation: Pix2Pix



Slide credit: Simon Prince

Image translation: SRGAN

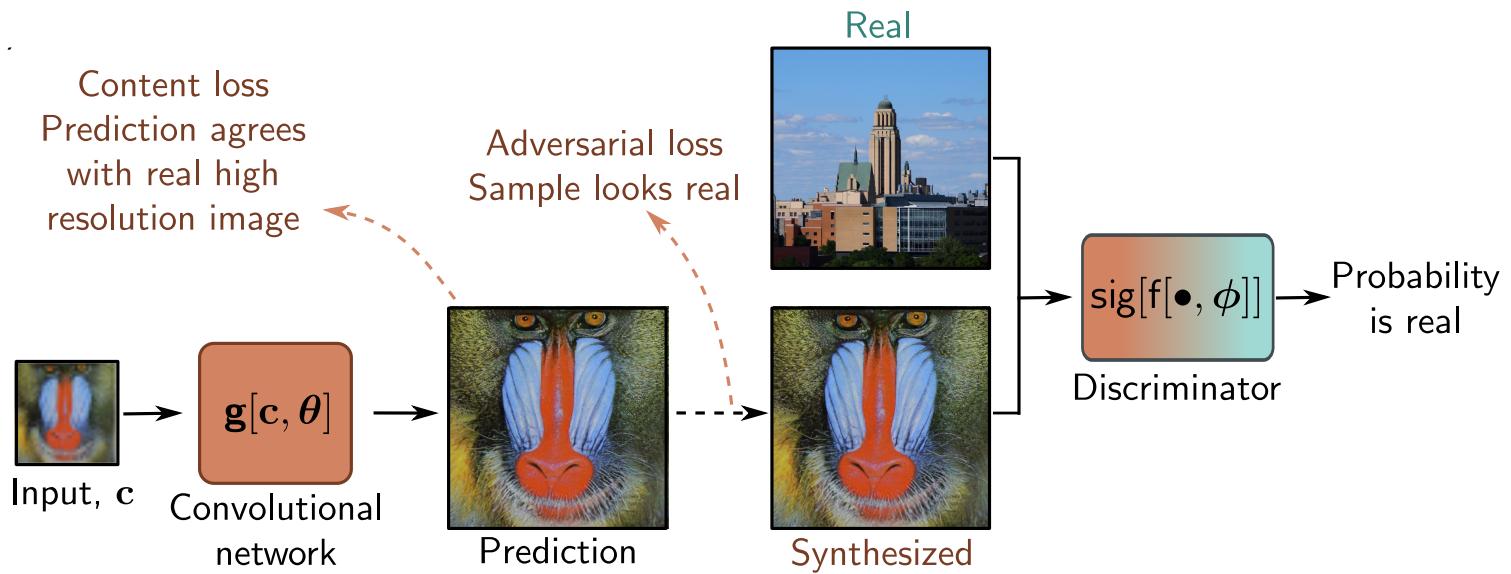
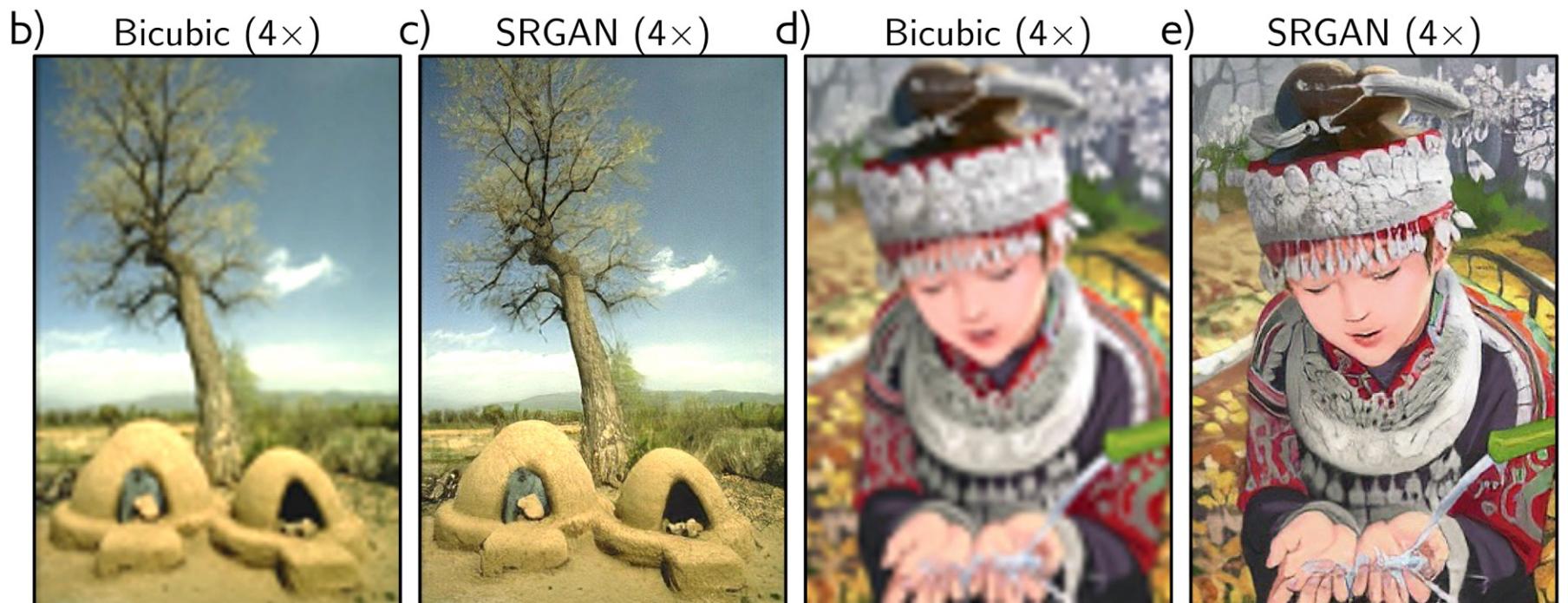
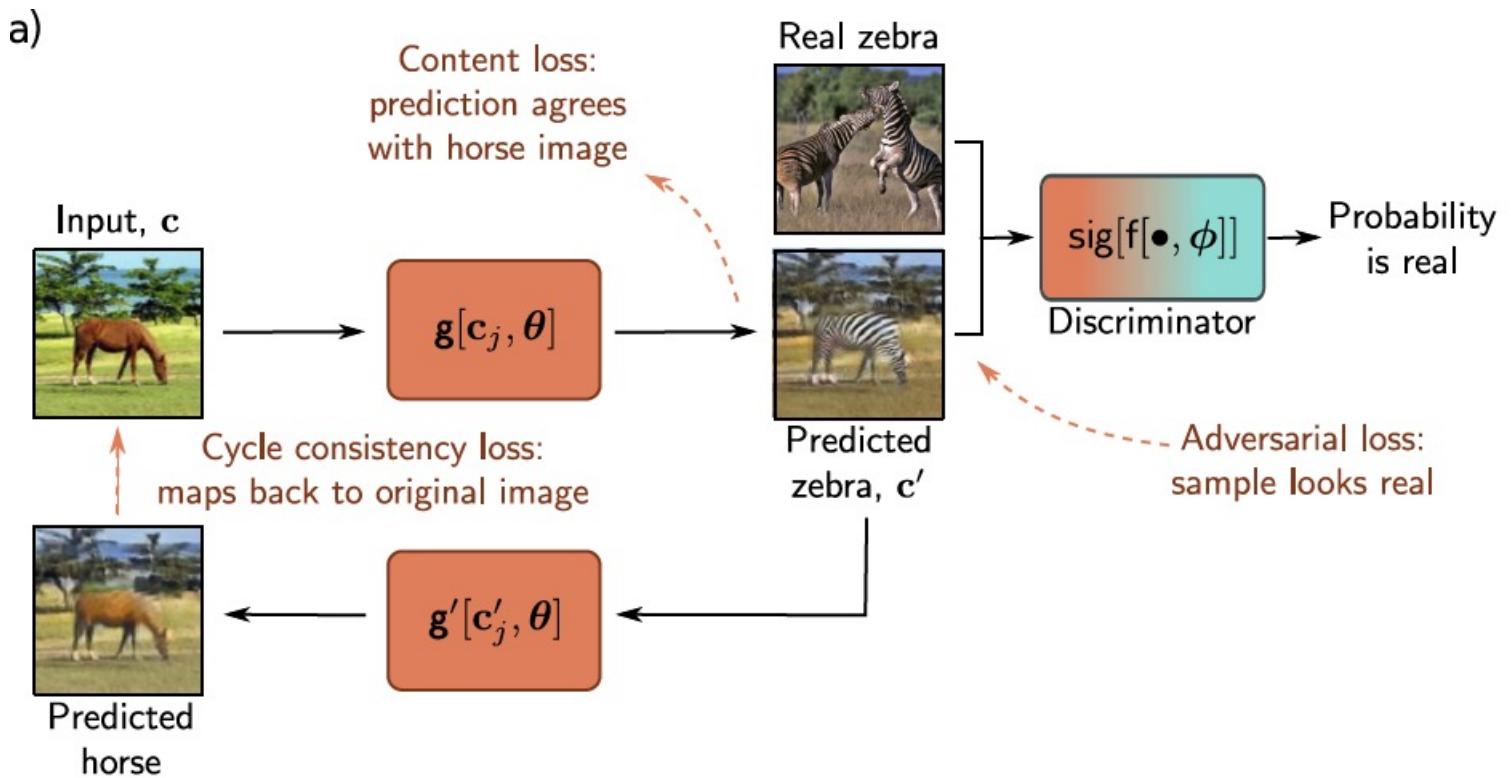


Image translation: SRGAN



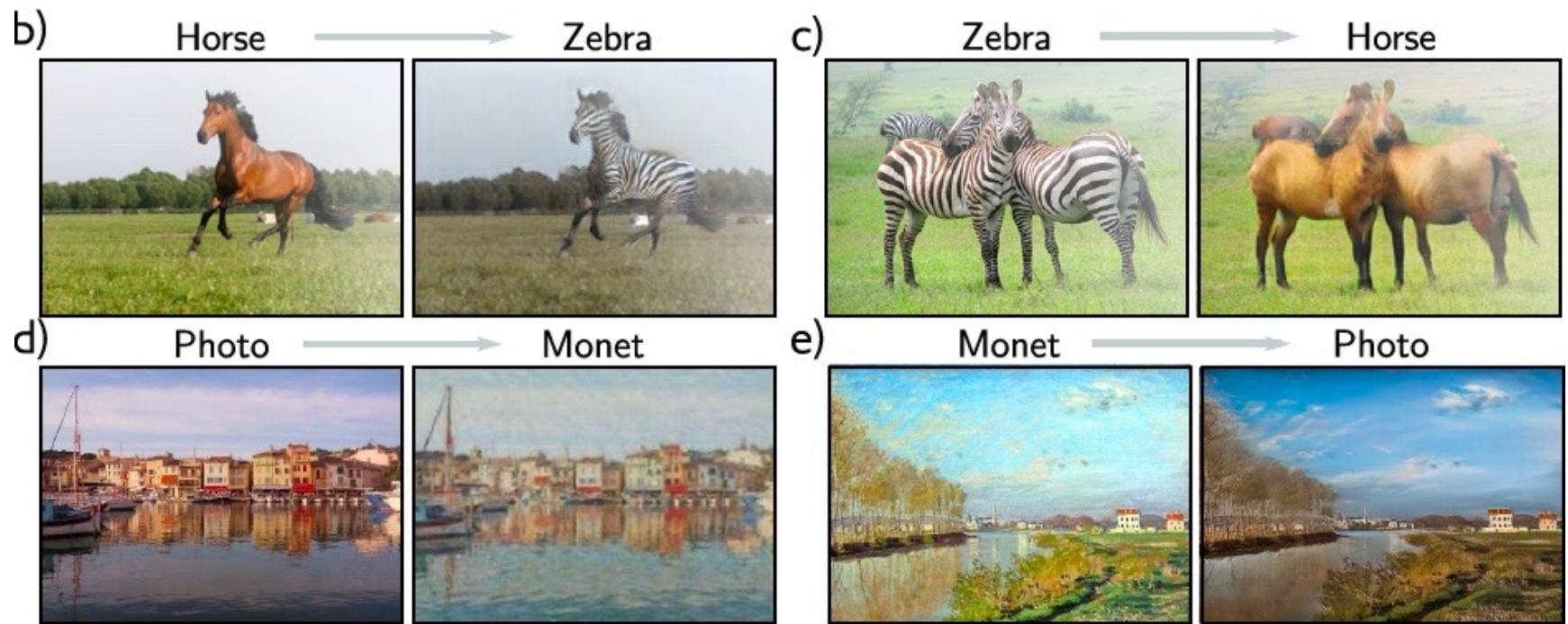
Slide credit: Simon Prince

Image translation: CycleGAN



Slide credit: Simon Prince

Image translation: CycleGAN



Slide credit: Simon Prince



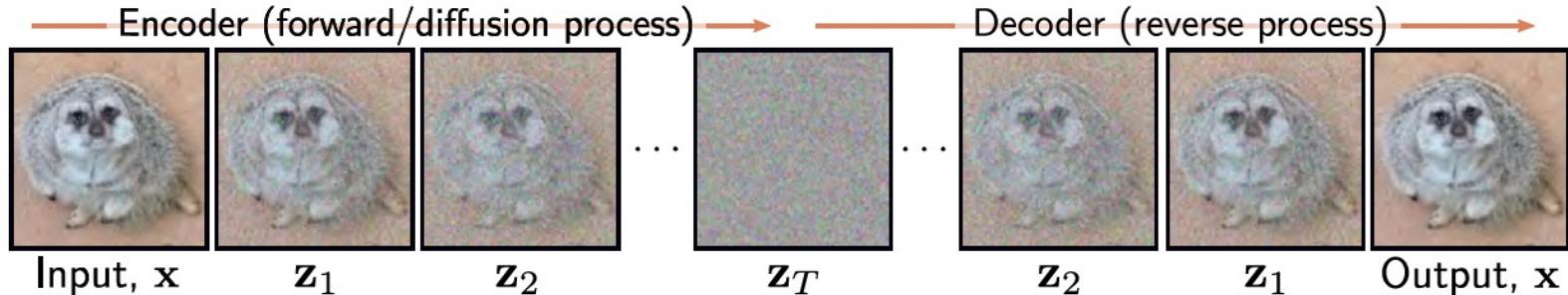
Diffusion Model

Video Generation – OpenAI Sora



Denoising Diffusion models

- ❑ Forward process/diffusion gradually adds noise to data
- ❑ Backward process generates data by denoising



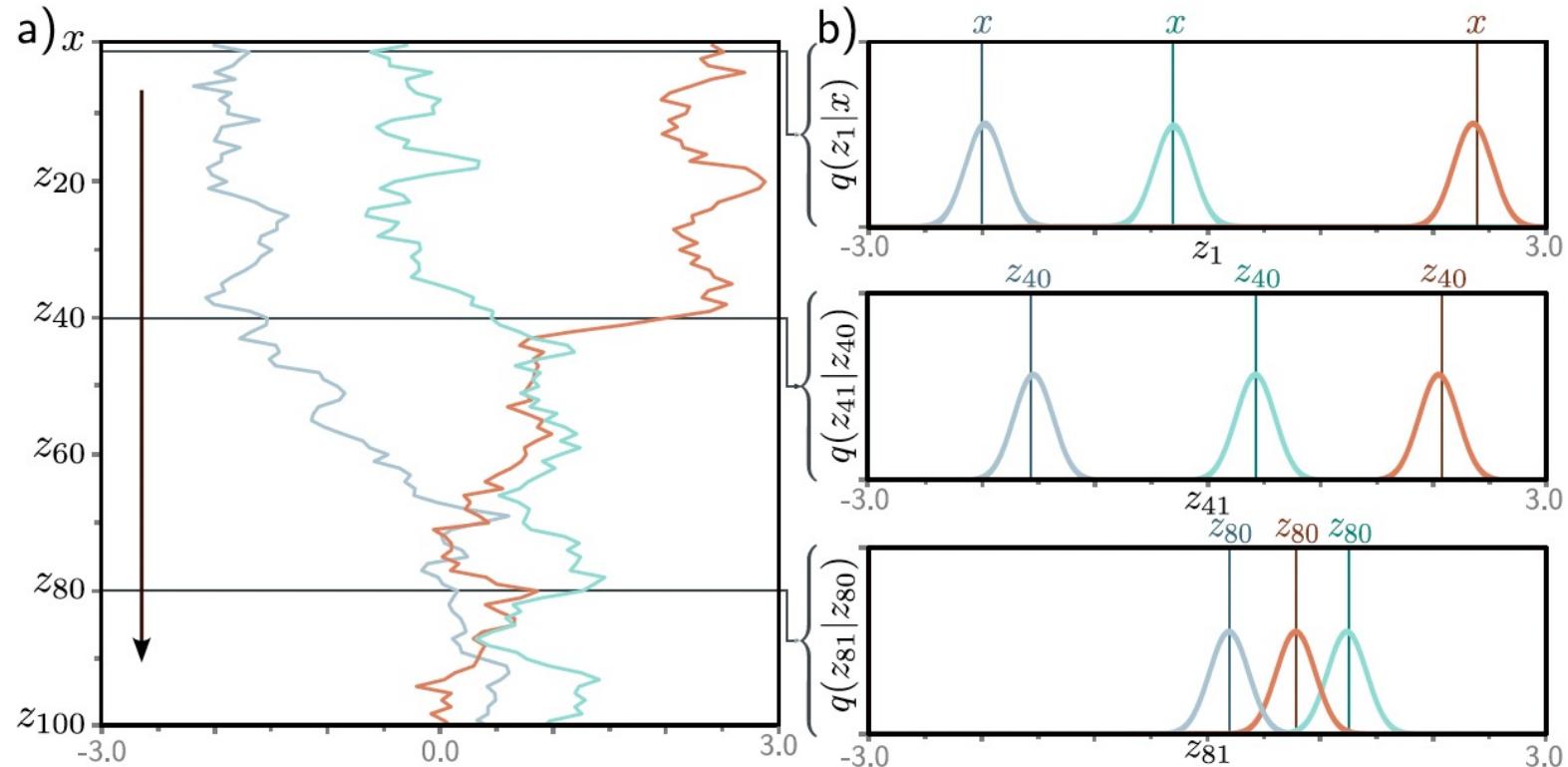
Forward process/diffusion

$$\begin{aligned}\mathbf{z}_1 &= \sqrt{1 - \beta_1} \cdot \mathbf{x} + \sqrt{\beta_1} \cdot \boldsymbol{\epsilon}_1 \\ \mathbf{z}_t &= \sqrt{1 - \beta_t} \cdot \mathbf{z}_{t-1} + \sqrt{\beta_t} \cdot \boldsymbol{\epsilon}_t\end{aligned}\quad \forall t \in 2, \dots, T,$$

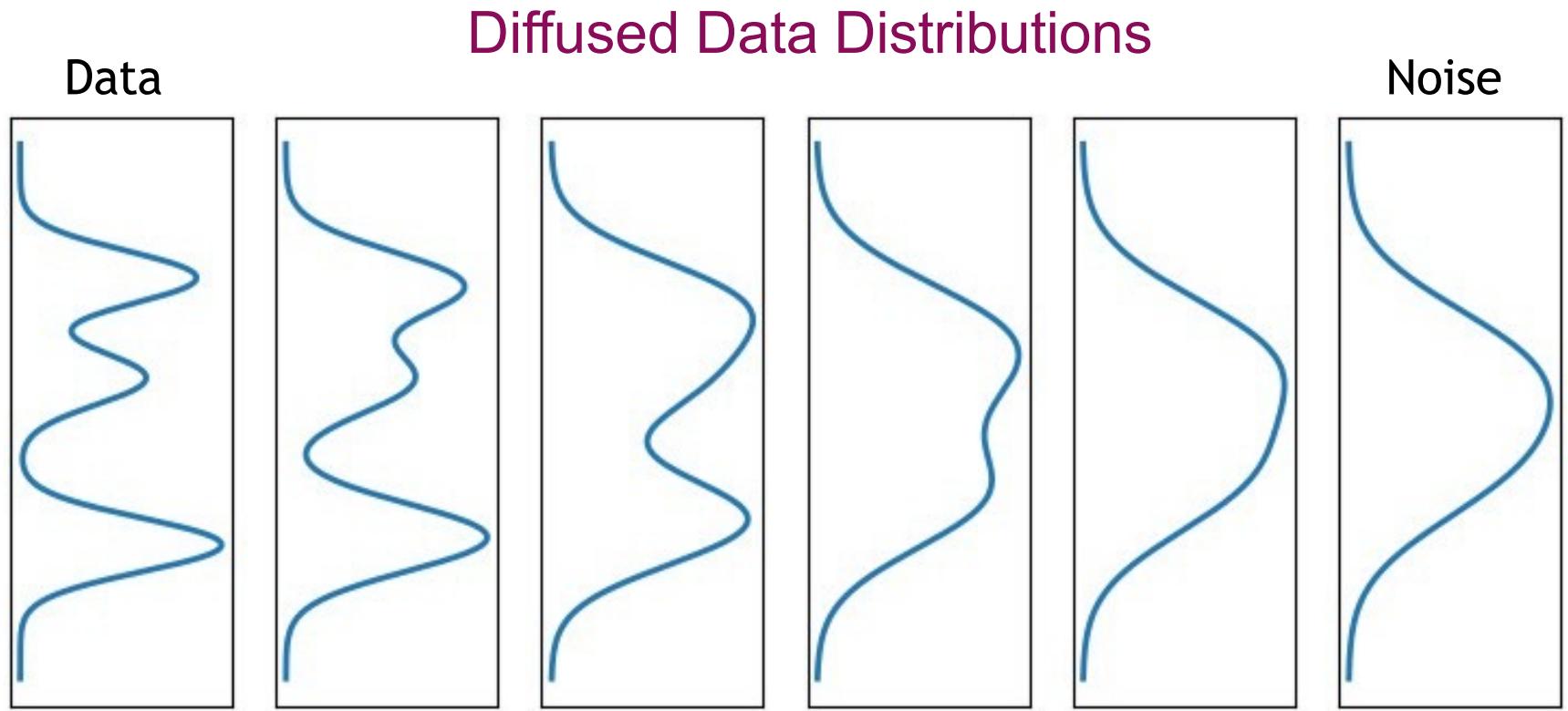
$$\begin{aligned}q(\mathbf{z}_1 | \mathbf{x}) &= \text{Norm}_{\mathbf{z}_1} \left[\sqrt{1 - \beta_1} \mathbf{x}, \beta_1 \mathbf{I} \right] \\ q(\mathbf{z}_t | \mathbf{z}_{t-1}) &= \text{Norm}_{\mathbf{z}_t} \left[\sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I} \right]\end{aligned}\quad \forall t \in \{2, \dots, T\}.$$

$$q(\mathbf{z}_{1\dots T} | \mathbf{x}) = q(\mathbf{z}_1 | \mathbf{x}) \prod_{t=2}^T q(\mathbf{z}_t | \mathbf{z}_{t-1}).$$

Forward process



What happens to data distribution in forward process?



Diffusion kernel

- ❑ Can we directly compute \mathbf{z}_t without computing previous \mathbf{z} ?

$$\mathbf{z}_t = \sqrt{\alpha_t} \cdot \mathbf{x} + \sqrt{1 - \alpha_t} \cdot \boldsymbol{\epsilon},$$

where $\alpha_t = \prod_{s=1}^t 1 - \beta_s$. We can equivalently write this in probabilistic form:

$$q(\mathbf{z}_t | \mathbf{x}) = \text{Norm}_{\mathbf{z}_t} \left[\sqrt{\alpha_t} \cdot \mathbf{x}, (1 - \alpha_t) \mathbf{I} \right].$$

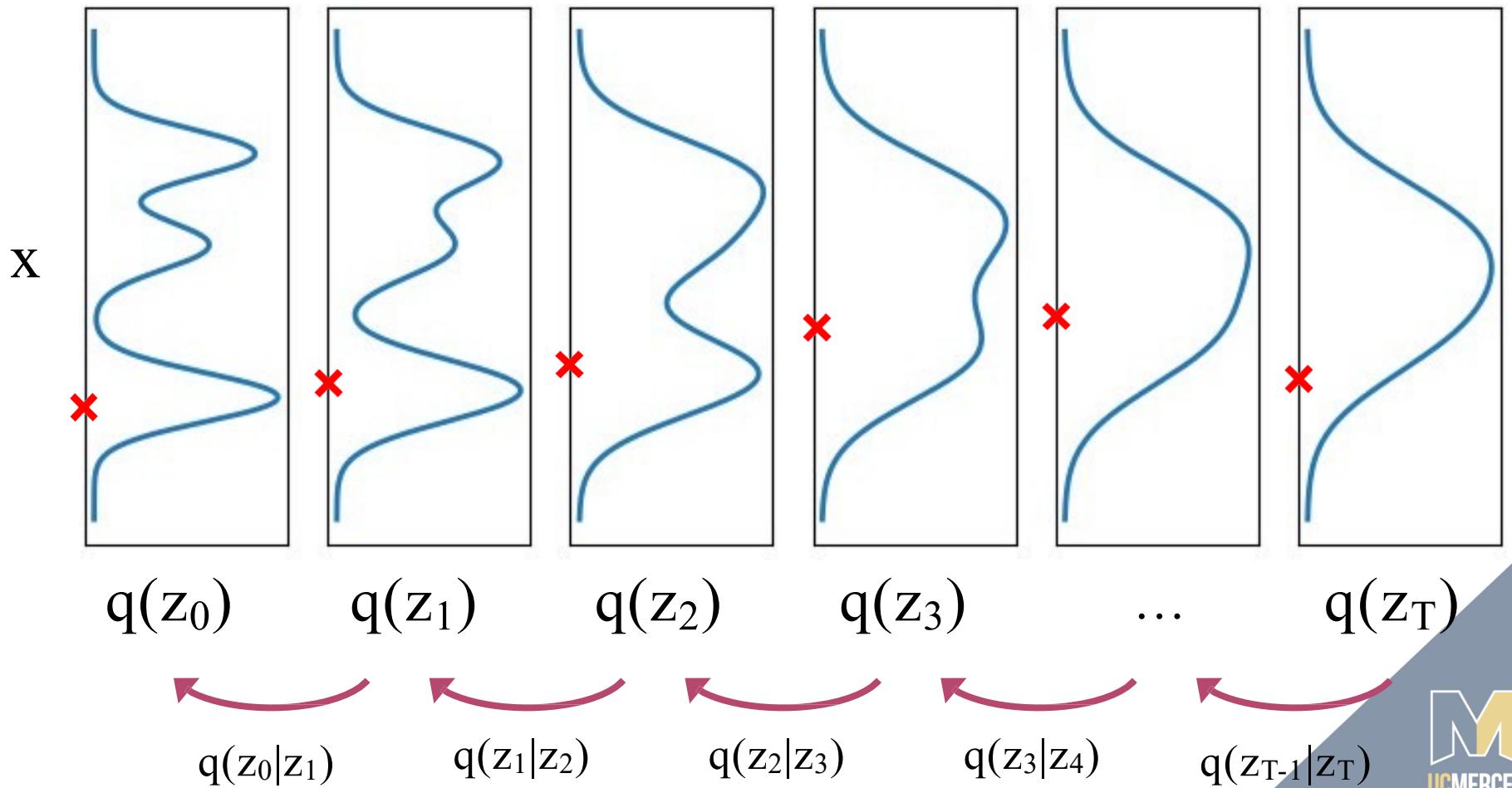
- ❑ Marginal distribution

$$q(\mathbf{z}_t) = \int q(\mathbf{z}_t | \mathbf{x}) Pr(\mathbf{x}) d\mathbf{x}.$$

- ❑ Quiz: What is the distribution of \mathbf{z}_T ?

Generative learning by denoising

Diffused Data Distributions



Conditional distribution in reverse process

- To reverse the process, we apply Bayes' rule

$$q(\mathbf{z}_{t-1} | \mathbf{z}_t) = \frac{q(\mathbf{z}_t | \mathbf{z}_{t-1}) q(\mathbf{z}_{t-1})}{q(\mathbf{z}_t)}.$$

- This is intractable since we cannot compute marginal $q(\mathbf{z}_t)$
- KEY assumption: Approximate conditional distribution using a Gaussian distribution
- Note: closed form Gaussian if conditioned on \mathbf{X}

$$q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) = \text{Norm}_{\mathbf{z}_{t-1}} \left[\frac{(1 - \alpha_{t-1})}{1 - \alpha_t} \sqrt{1 - \beta_t} \mathbf{z}_t + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1 - \alpha_t} \mathbf{x}, \frac{\beta_t (1 - \alpha_{t-1})}{1 - \alpha_t} \mathbf{I} \right]$$

Decoder model (reverse process)

- Let a neural network predict mean of Gaussian

$$\begin{aligned} Pr(\mathbf{z}_T) &= \text{Norm}_{\mathbf{z}_T} [\mathbf{0}, \mathbf{I}] \\ Pr(\mathbf{z}_{t-1} | \mathbf{z}_t, \phi_t) &= \text{Norm}_{\mathbf{z}_{t-1}} \left[\mathbf{f}_t[\mathbf{z}_t, \phi_t], \sigma_t^2 \mathbf{I} \right] \\ Pr(\mathbf{x} | \mathbf{z}_1, \phi_1) &= \text{Norm}_{\mathbf{x}} \left[\mathbf{f}_1[\mathbf{z}_1, \phi_1], \sigma_1^2 \mathbf{I} \right], \end{aligned}$$

- What about variance of Gaussian? – Fixed schedule

Training

- ❑ To train the model, we maximize the log-likelihood of the training data $\{\mathbf{x}_i\}$ with respect to the parameters

$$\hat{\boldsymbol{\phi}}_{1\dots T} = \operatorname{argmax}_{\boldsymbol{\phi}_{1\dots T}} \left[\sum_{i=1}^I \log \left[Pr(\mathbf{x}_i | \boldsymbol{\phi}_{1\dots T}) \right] \right]$$

Evidence lower bound (ELBO)

$$\begin{aligned}\log [Pr(\mathbf{x}|\phi_{1...T})] &= \log \left[\int Pr(\mathbf{x}, \mathbf{z}_{1...T}|\phi_{1...T}) d\mathbf{z}_{1...T} \right] \\ &= \log \left[\int q(\mathbf{z}_{1...T}|\mathbf{x}) \frac{Pr(\mathbf{x}, \mathbf{z}_{1...T}|\phi_{1...T})}{q(\mathbf{z}_{1...T}|\mathbf{x})} d\mathbf{z}_{1...T} \right] \\ &\geq \int q(\mathbf{z}_{1...T}|\mathbf{x}) \log \left[\frac{Pr(\mathbf{x}, \mathbf{z}_{1...T}|\phi_{1...T})}{q(\mathbf{z}_{1...T}|\mathbf{x})} \right] d\mathbf{z}_{1...T}.\end{aligned}$$

$$= \mathbb{E}_{q(\mathbf{z}_1|\mathbf{x})} \left[\log [Pr(\mathbf{x}|\mathbf{z}_1, \phi_1)] \right] - \sum_{t=2}^T \mathbb{E}_{q(\mathbf{z}_t|\mathbf{x})} \left[D_{KL} \left[q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) \middle\| Pr(\mathbf{z}_{t-1}|\mathbf{z}_t, \phi_t) \right] \right]$$

Analyzing ELBO

$$\mathbb{E}_{q(\mathbf{z}_1|\mathbf{x})} \left[\log [Pr(\mathbf{x}|\mathbf{z}_1, \phi_1)] \right] - \sum_{t=2}^T \mathbb{E}_{q(\mathbf{z}_t|\mathbf{x})} \left[D_{KL} \left[q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) \middle\| Pr(\mathbf{z}_{t-1}|\mathbf{z}_t, \phi_t) \right] \right]$$

$$Pr(\mathbf{x}|\mathbf{z}_1, \phi_1) = \text{Norm}_{\mathbf{x}} \left[\mathbf{f}_1[\mathbf{z}_1, \phi_1], \sigma_1^2 \mathbf{I} \right]$$

$$Pr(\mathbf{z}_{t-1}|\mathbf{z}_t, \phi_t) = \text{Norm}_{\mathbf{z}_{t-1}} \left[\mathbf{f}_t[\mathbf{z}_t, \phi_t], \sigma_t^2 \mathbf{I} \right]$$

$$q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) = \text{Norm}_{\mathbf{z}_{t-1}} \left[\frac{(1-\alpha_{t-1})}{1-\alpha_t} \sqrt{1-\beta_t} \mathbf{z}_t + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1-\alpha_t} \mathbf{x}, \frac{\beta_t(1-\alpha_{t-1})}{1-\alpha_t} \mathbf{I} \right]$$

Diffusion Loss Function

$$L[\phi_{1\dots T}] = \sum_{i=1}^I \underbrace{\left(-\log \left[\text{Norm}_{\mathbf{x}_i} [\mathbf{f}_1[\mathbf{z}_{i1}, \phi_1], \sigma_1^2 \mathbf{I}] \right] \right)}_{\text{reconstruction term}} + \sum_{t=2}^T \frac{1}{2\sigma_t^2} \left\| \underbrace{\frac{1-\alpha_{t-1}}{1-\alpha_t} \sqrt{1-\beta_t} \mathbf{z}_{it} + \frac{\sqrt{\alpha_{t-1}} \beta_t}{1-\alpha_t} \mathbf{x}_i - \mathbf{f}_t[\mathbf{z}_{it}, \phi_t]}_{\text{target, mean of } q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x})} - \underbrace{\mathbf{f}_t[\mathbf{z}_{it}, \phi_t]}_{\text{predicted } \mathbf{z}_{t-1}} \right\|^2 \right), \quad (18.29)$$

where \mathbf{x}_i is the i^{th} data point, and \mathbf{z}_{it} is the associated latent variable at diffusion step t .

Reparameterization of target

□ Target (mean) can be reparameterized using noise

$$\mathbf{z}_t = \sqrt{\alpha_t} \cdot \mathbf{x} + \sqrt{1 - \alpha_t} \cdot \boldsymbol{\epsilon}.$$

$$\mathbf{x} = \frac{1}{\sqrt{\alpha_t}} \cdot \mathbf{z}_t - \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}} \cdot \boldsymbol{\epsilon}$$

□ Reparameterization of network

$$\mathbf{f}_t[\mathbf{z}_t, \phi_t] = \frac{1}{\sqrt{1 - \beta_t}} \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t} \sqrt{1 - \beta_t}} \mathbf{g}_t[\mathbf{z}_t, \phi_t].$$

□ Simple loss function for diffusion model:

$$\begin{aligned} L[\phi_{1\dots T}] &= \sum_{i=1}^I \sum_{t=1}^T \left\| \mathbf{g}_t[\mathbf{z}_{it}, \phi_t] - \boldsymbol{\epsilon}_{it} \right\|^2 \\ &= \sum_{i=1}^I \sum_{t=1}^T \left\| \mathbf{g}_t \left[\sqrt{\alpha_t} \cdot \mathbf{x}_i + \sqrt{1 - \alpha_t} \cdot \boldsymbol{\epsilon}_{it}, \phi_t \right] - \boldsymbol{\epsilon}_{it} \right\|^2. \end{aligned}$$

Diffusion Model Training

Algorithm 18.1: Diffusion model training

Input: Training data \mathbf{x}

Output: Model parameters ϕ_t

repeat

for $i \in \mathcal{B}$ **do** // For every training example index in batch
 $t \sim \text{Uniform}[1, \dots, T]$ // Sample random timestep
 $\epsilon \sim \text{Norm}[\mathbf{0}, \mathbf{I}]$ // Sample noise
 $\ell_i = \left\| \mathbf{g}_t \left[\sqrt{\alpha_t} \mathbf{x}_i + \sqrt{1 - \alpha_t} \epsilon, \phi_t \right] - \epsilon \right\|^2$ // Compute individual loss

Accumulate losses for batch and take gradient step

until converged

Diffusion Model Inference

Algorithm 18.2: Sampling

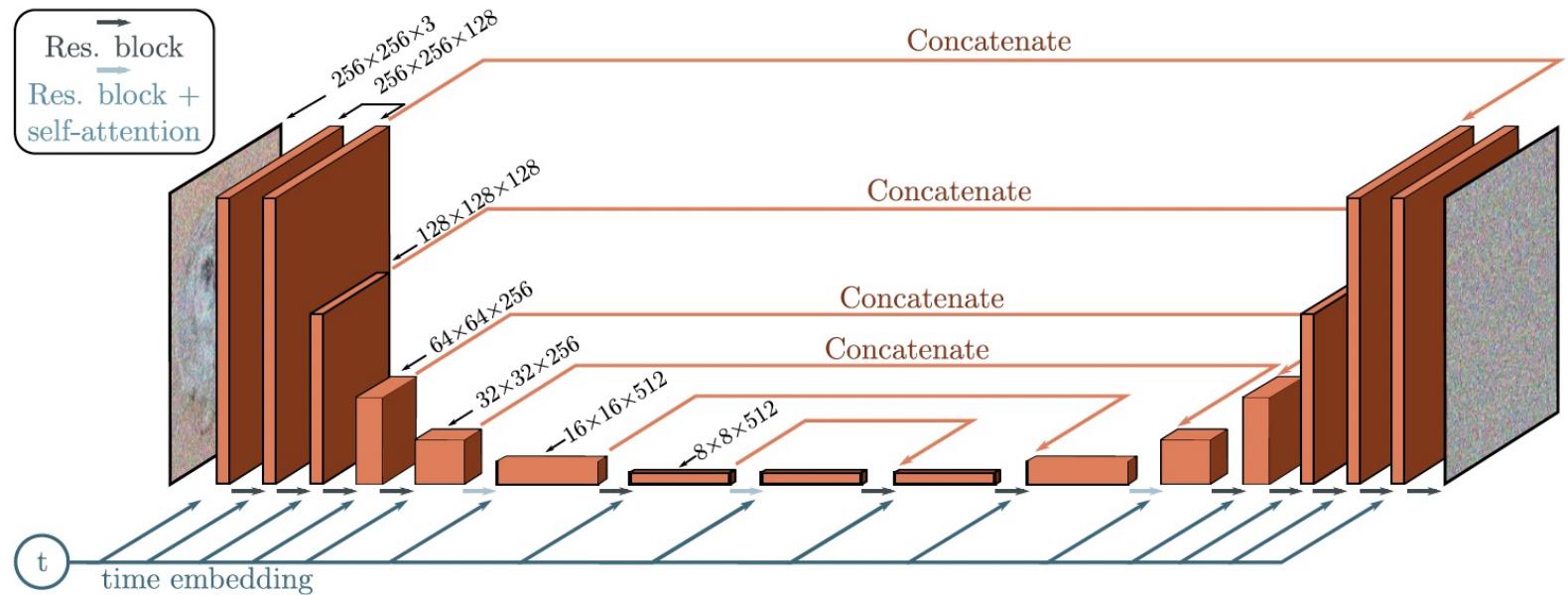
Input: Model, $\mathbf{g}_t[\bullet, \phi_t]$

Output: Sample, \mathbf{x}

```
 $\mathbf{z}_T \sim \text{Norm}_{\mathbf{z}}[\mathbf{0}, \mathbf{I}]$  // Sample last latent variable
for  $t = T \dots 2$  do
     $\hat{\mathbf{z}}_{t-1} = \frac{1}{\sqrt{1-\beta_t}} \mathbf{z}_t - \frac{\beta_t}{\sqrt{1-\alpha_t} \sqrt{1-\beta_t}} \mathbf{g}_t[\mathbf{z}_t, \phi_t]$  // Predict previous latent variable
     $\epsilon \sim \text{Norm}_{\epsilon}[\mathbf{0}, \mathbf{I}]$  // Draw new noise vector
     $\mathbf{z}_{t-1} = \hat{\mathbf{z}}_{t-1} + \sigma_t \epsilon$  // Add noise to previous latent variable
 $\mathbf{x} = \frac{1}{\sqrt{1-\beta_1}} \mathbf{z}_1 - \frac{\beta_1}{\sqrt{1-\alpha_1} \sqrt{1-\beta_1}} \mathbf{g}_1[\mathbf{z}_1, \phi_1]$  // Generate sample from  $\mathbf{z}_1$  without noise
```

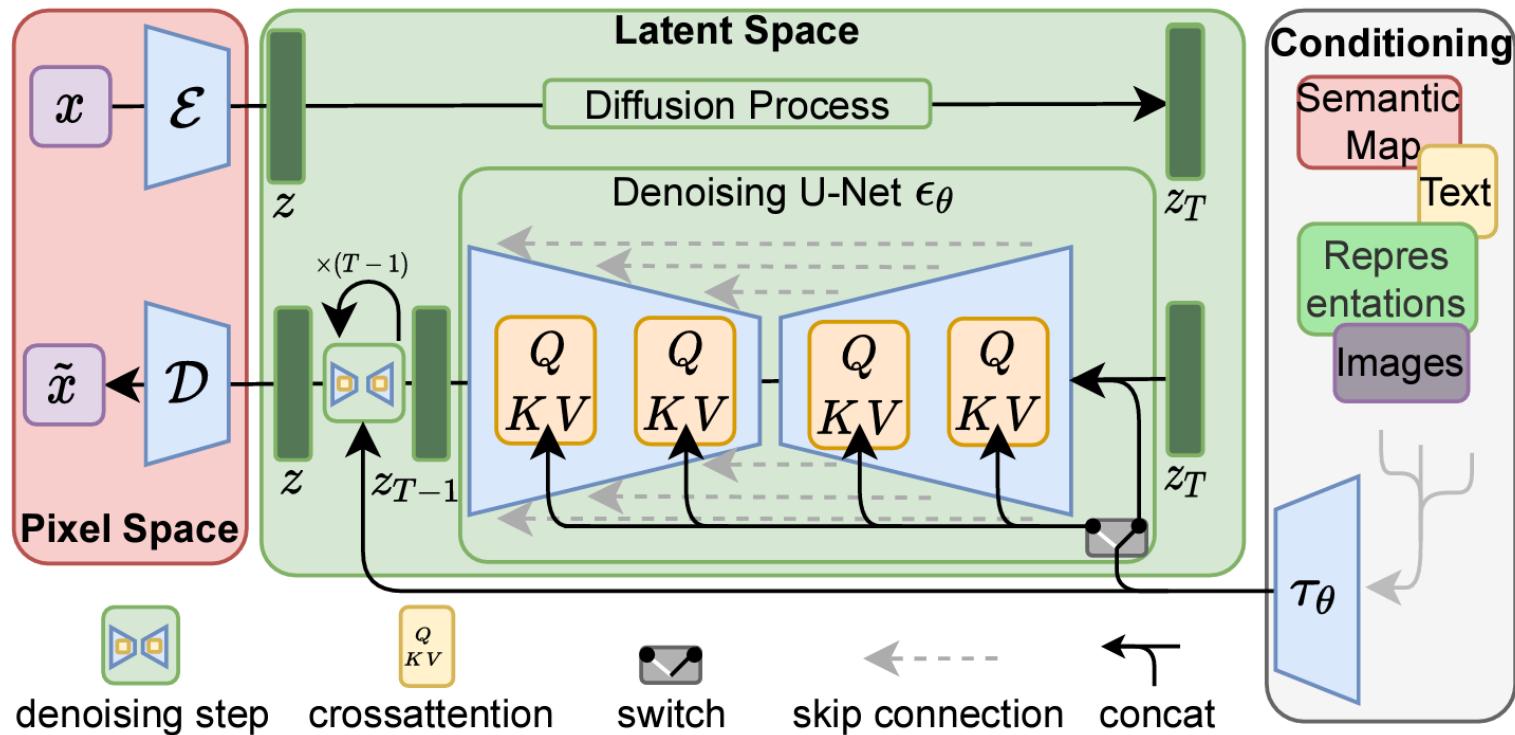
Implementation

- ❑ U-Net based denoising network
- ❑ Noisy image as input
- ❑ Noise as output



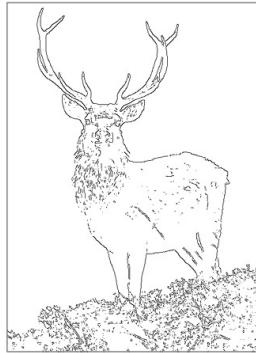
Latent diffusion model

- Reduce dimensionality of image by VAE
- Diffusion process in latent space of VAE
- Conditional generation using text, image, semantic map etc.



Parameter-efficient Conditonal Generation

- ControlNet: Controlling diffusion model with various control



Input Canny edge



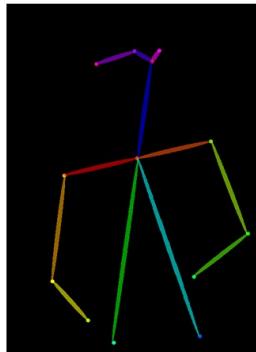
Default



"masterpiece of fairy tale, giant deer, golden antlers"



"..., quaint city Galic"



Input human pose



Default



"chef in kitchen"

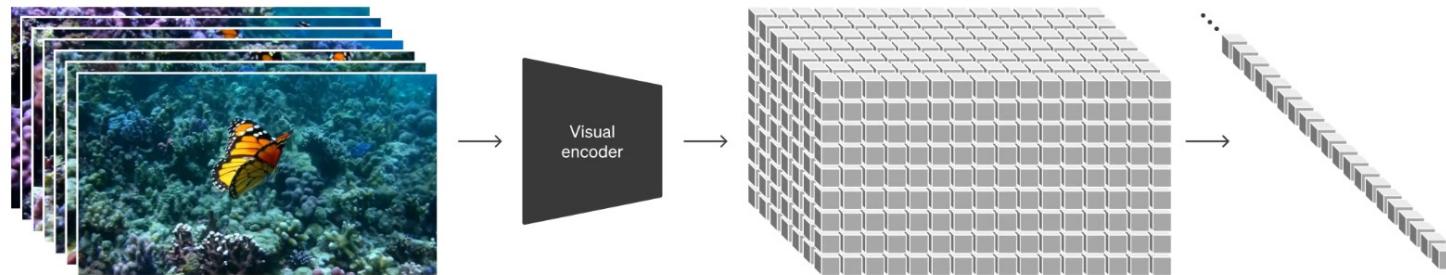


"Lincoln statue"

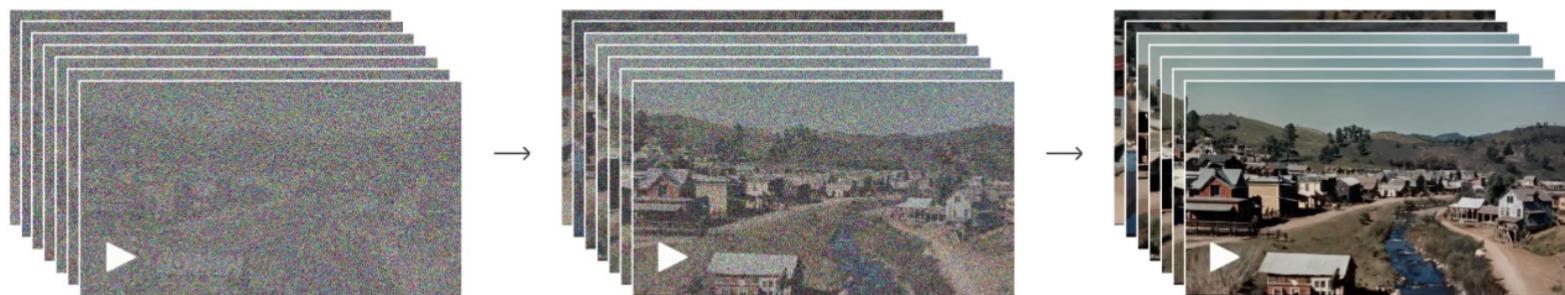
- Key idea: Zero convolution

How Sora works?

- ❑ Latent diffusion model + U-ViT



Turning video into patches and compressed



Diffusion model for video