Topics:

- Variational Autoencoders

# CS 4803-DL / 7643-A
# ZSOLT KIRA

- **A4 due March 30th (grace until April 1st)**

- **Projects!**
  - Project Check-in extended to **March 24th (grace 26th)**
  - Make sure to contribute equally with your teammates!!!
  - We will have optional team peer review, and reduce scores if necessary

- **Meta OH today 3pm ET**

Back to Generative Models

**Supervised Learning**

- Train Input: $\{X, Y\}$
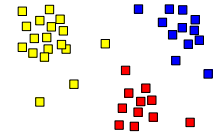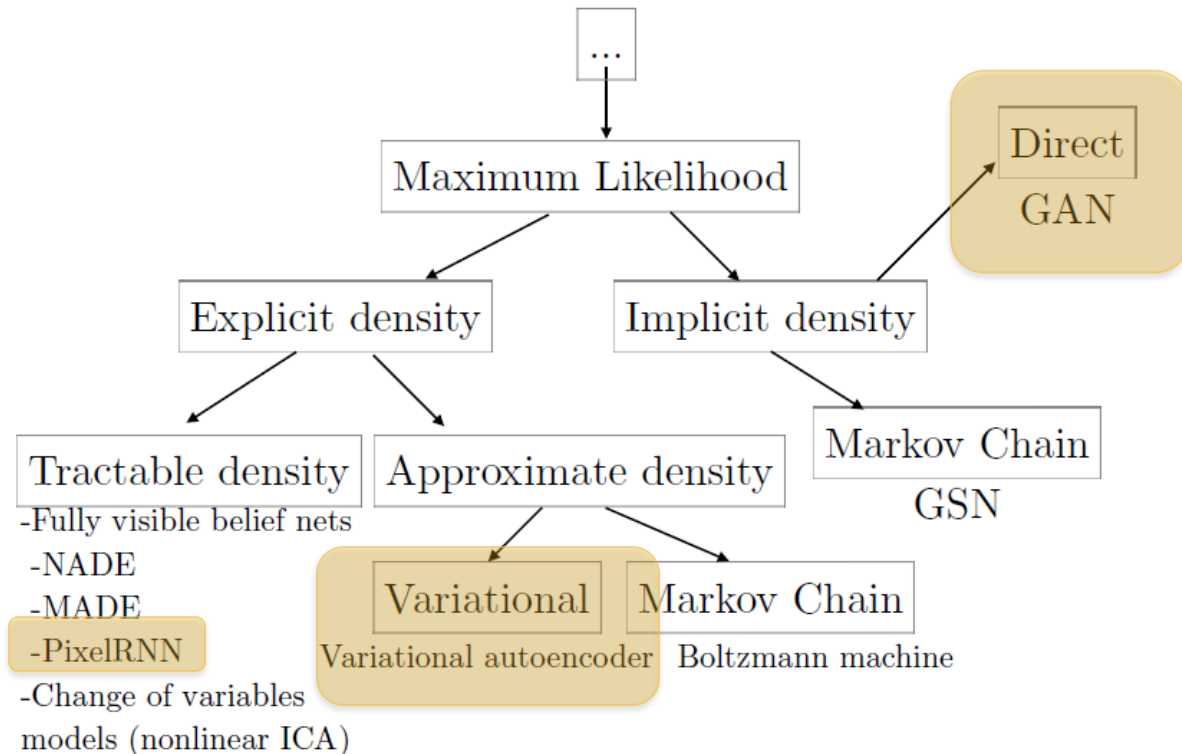- Learning output: $f : X \rightarrow Y, P(y|x)$
- e.g. classification

Sheep
Dog
→ Cat
Lion
Giraffe

**Less Labels**

**Unsupervised Learning**

- Input: $\{X\}$
- Learning output: $P(x)$
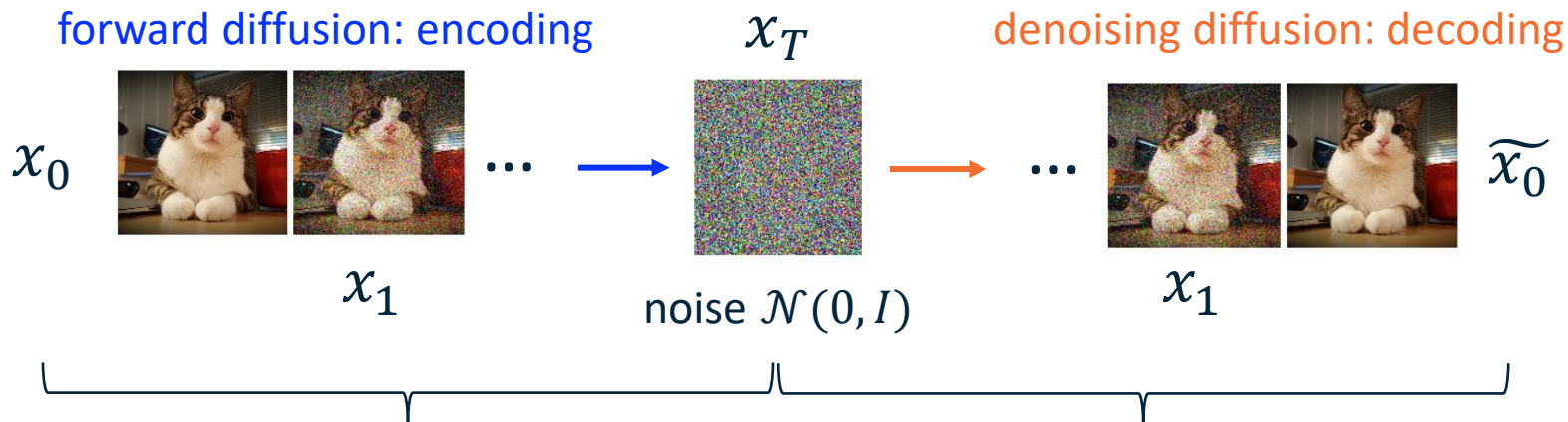- Example: Clustering, density estimation, etc.

Supervised Learning

Georgia Tech

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Generative Models**

# Forward/Reverse Processes

forward diffusion: encoding

$x_T$

denoising diffusion: decoding

$x_0$

$x_1$

noise $\mathcal{N}(0, I)$

$x_1$

$\widetilde{x_0}$

Known / predefined:

$q(x_{1:T}|x_0)$

Output: Noise mean to remove, sample & use w/ $x_t$ to get $x_{t-1}$

Input:

$x_t$

U-Net

Unknown / learned:

$$p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^{T} p_\theta(x_{t-1}|x_t)$$

# "StableDiffusion"



Rombach and Blattmann *et al.*, 2022

10 months ago



6-7 second videos

**Video Generation**

Georgia Tech

# Transformers!



$$q_s \in \mathbb{R}^{(B.(1+t)) \times 1 \times h_p \times w_p \times C}$$
$$q_t \in \mathbb{R}^{(B.W_n) \times (1+t) \times h'_p \times w'_p \times c}$$

$$k = v = \text{concat}[q_s, \text{t5\_xl (text)}]$$

**Gupta et al., Photorealistic Video Generation with Diffusion Models**

Now



**Video Generation**                https://openai.com/sora    Georgia Tech

**Video Generation**

https://openai.com/sora

Georgia Tech

**Video Generation**

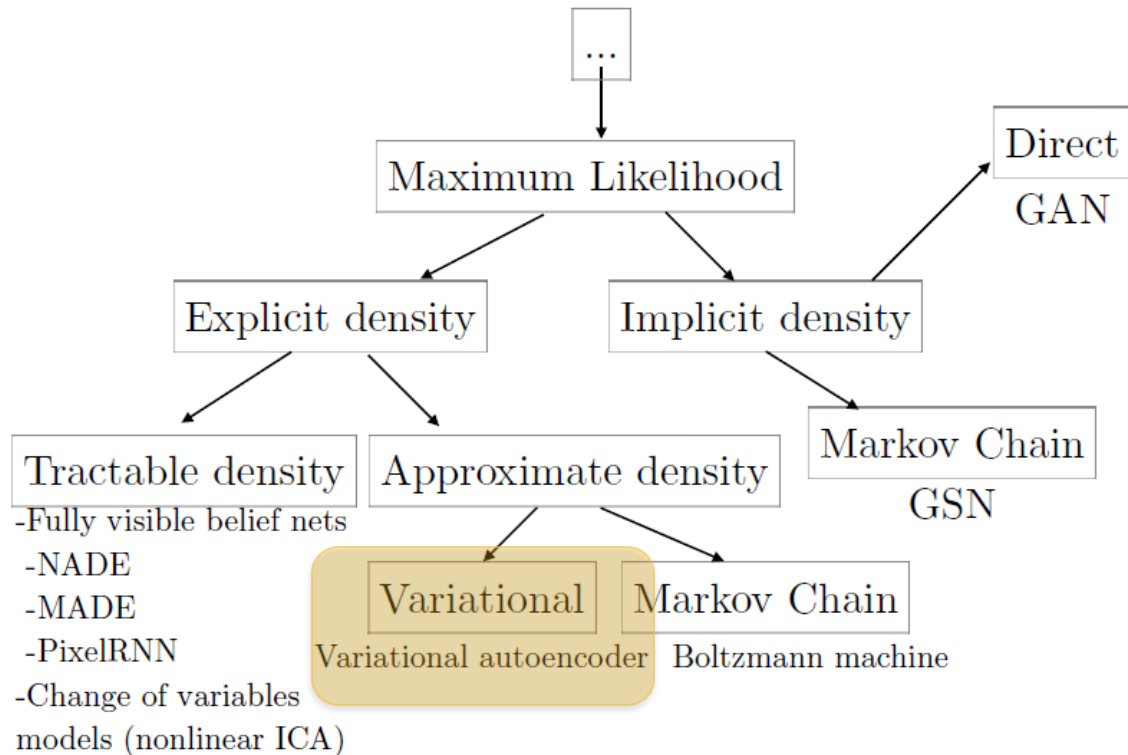https://openai.com/sora

Georgia Tech

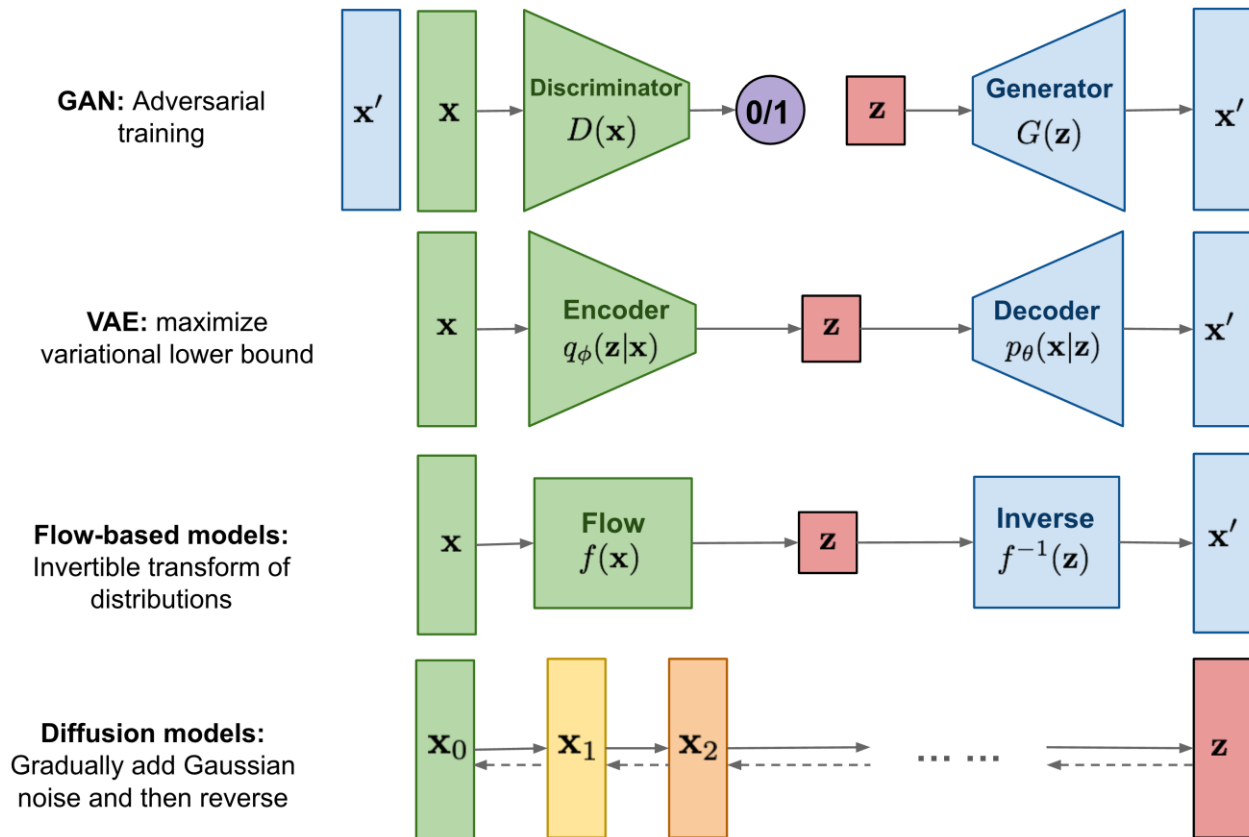**Video Generation – Failure Cases**

Georgia Tech

# Variational Autoencoders (VAEs)

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Generative Models**

# Comparison



**GAN:** Adversarial training

**VAE:** maximize variational lower bound

**Flow-based models:** Invertible transform of distributions

**Diffusion models:** Gradually add Gaussian noise and then reverse

Minimize the difference (with MSE)

Encoder

Decoder
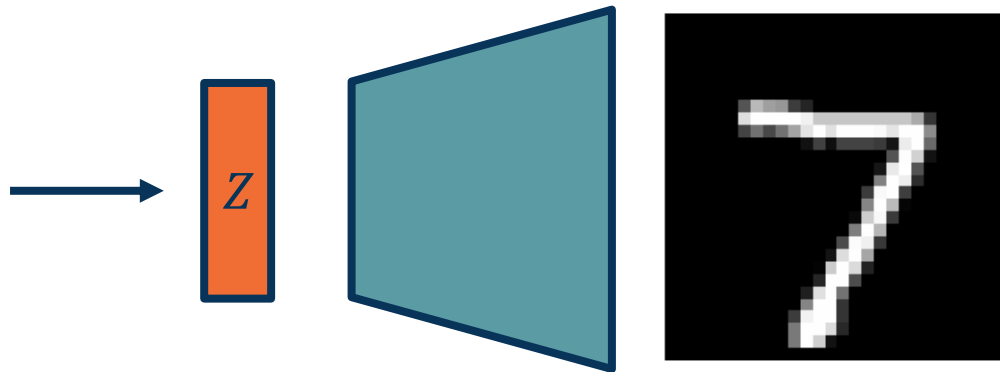
**Low dimensional embedding**

Linear layers with reduced dimension or Conv-2d layers with stride

Linear layers with increasing dimension or Conv-2d layers with bilinear upsampling

**Autoencoders**

Georgia Tech

**What is this?**
**Hidden/Latent variables**
**Factors of variation that**
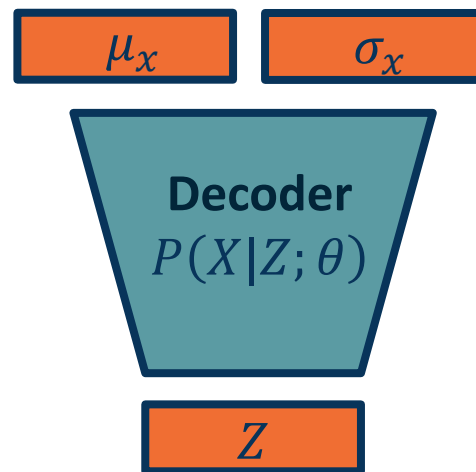**produce an image:**
**(digit, orientation, scale, etc.)**

$$P(X) = \int P(X|Z;\theta)P(Z)dZ$$

- ⬡ We cannot maximize this likelihood due to the integral

- ⬡ Instead we maximize a variational *lower bound* (VLB) that we *can* compute

*Kingma & Welling, Auto-Encoding Variational Bayes*

**Formalizing the Generative Model**

- We can combine the probabilistic view, sampling, autoencoders, and approximate optimization

- Just as before, sample $Z$ from simpler distribution

- We can also output parameters of a probability distribution!
  - **Example**: $\mu, \sigma$ of Gaussian distribution
  - For multi-dimensional version output diagonal covariance

- How can we maximize $P(X) = \int P(X|Z; \theta) P(Z) dZ$



$\mu_x$    $\sigma_x$

**Decoder**
$P(X|Z; \theta)$

$Z$

- We can combine the probabilistic view, sampling, autoencoders, and approximate optimization



$$\mu_z \quad \sigma_z$$

**Encoder**
$Q(Z|X;\phi)$

$X$

- Given an image, estimate $Z$

- Again, output *parameters of a distribution*

- We can tie the encoder and decoder together into a probabilistic autoencoder
  - Given data (X), estimate $\mu_z, \sigma_z$ and sample from $N(\mu_z, \sigma_z)$
  - Given $Z$, estimate $\mu_x, \sigma_x$ and sample from $N(\mu_x, \sigma_x)$



| $\mu_z$ | $\sigma_z$ | | $\mu_x$ | $\sigma_x$ |

**Encoder**
$Q(Z|X; \phi)$

**Decoder**
$P(X|Z; \theta)$

X

Z

Georgia Tech

◆ How can we optimize the parameters of the two networks?

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

**Maximizing Likelihood**

Georgia Tech

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \qquad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \qquad (\text{Bayes' Rule})$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \qquad (\text{Multiply by constant})$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \qquad (\text{Logarithms})$$

**Maximizing Likelihood**

Georgia Tech

Aside: KL Divergence (distance measure for distributions), always >= 0

$$KL(p||q) = H_c(p, q) - H(p) = \sum p(x)\log p(x) - \sum p(x)\log q(x)$$

Definition of Expectation

$$\mathbb{E}[f] = \mathbb{E}_{x \sim q}[f(x)] = \sum_{x \in \Omega} q(x)f(x)$$

$$KL(a||b) = E[\log a(x)] - E[\log b(x)] = E\left[\log \frac{a(x)}{b(x)}\right]$$

**KL-Divergence**

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})}\left[\log p_\theta(x^{(i)})\right] \qquad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z)p_\theta(z)}{p_\theta(z \mid x^{(i)})}\right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z)p_\theta(z)}{p_\theta(z \mid x^{(i)})}\frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})}\right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - \mathbf{E}_z\left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)}\right] + \mathbf{E}_z\left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})}\right] \quad \text{(Logarithms)}$$

$$= \mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z)) + D_{KL}(q_\phi(z \mid x^{(i)}) \,||\, p_\theta(z \mid x^{(i)}))$$

The expectation wrt. z (using encoder network) let us write nice KL terms

*From CS231n, Fei-Fei Li, Justin Johnson, Serena Yeung*

**Maximizing Likelihood**

Georgia Tech

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z)) + D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z \mid x^{(i)}))$$

Decoder network gives $p_\theta(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)

This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!

$p_\theta(z|x)$ intractable (saw earlier), can't compute this KL term :( But we know KL divergence always >= 0.

**Maximizing Likelihood**

Georgia Tech

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \qquad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \qquad (\text{Bayes' Rule})$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \qquad (\text{Multiply by constant})$$

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \qquad (\text{Logarithms})$$

$$= \underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z \mid x^{(i)}))}_{> 0}$$

$$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$
Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg\max_{\theta, \phi} \sum_{i=1}^{N} \mathcal{L}(x^{(i)}, \theta, \phi)$$
Training: Maximize lower bound

**Maximizing Likelihood**

Georgia Tech

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z\left[\log p_\theta(x^{(i)}\mid z)\right] - D_{KL}(q_\phi(z\mid x^{(i)})\,\|\,p_\theta(z))}_{\mathcal{L}(x^{(i)},\theta,\phi)}$$

Make approximate posterior distribution close to prior



$\mu_z$  $\sigma_z$

**Encoder**
$Q(Z|X;\phi)$

X

*From CS231n, Fei-Fei Li, Justin Johnson, Serena Yeung*
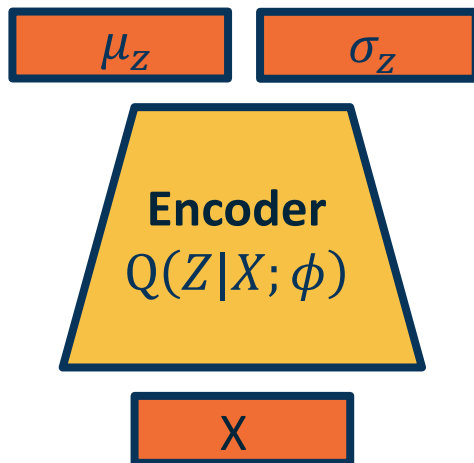
**Forward and Backward Passes**

Georgia Tech

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$
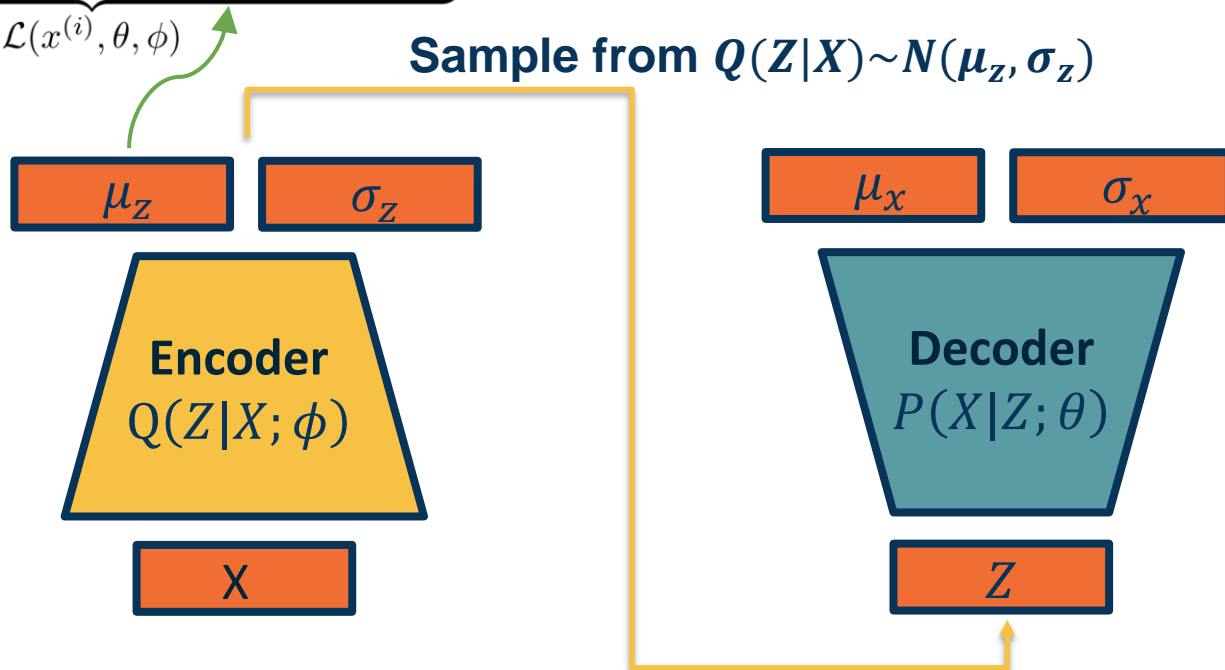
**Sample from $Q(Z|X) \sim N(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z)$**

| $\mu_z$ | $\sigma_z$ |
|---------|------------|

**Encoder**
$Q(Z|X; \phi)$

X

| $\mu_x$ | $\sigma_x$ |
|---------|------------|

**Decoder**
$P(X|Z; \theta)$

Z

*From CS231n, Fei-Fei Li, Justin Johnson, Serena Yeung*

**Forward and Backward Passes**

Georgia Tech

Putting it all together: maximizing the likelihood lower bound

Maximize likelihood of original input being reconstructed

$$\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z))$$
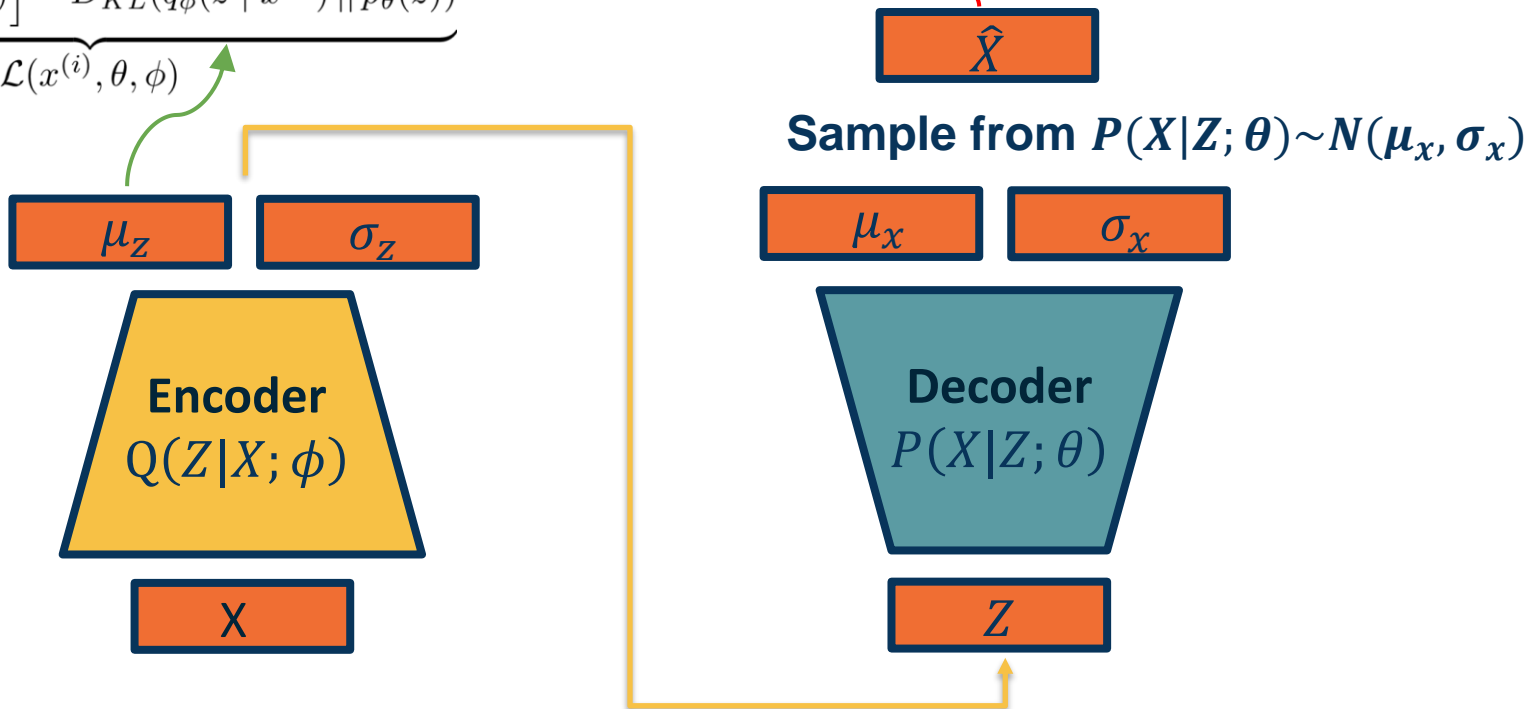
$$\mathcal{L}(x^{(i)}, \theta, \phi)$$

$\hat{X}$

**Sample from** $P(X|Z; \boldsymbol{\theta}) \sim N(\boldsymbol{\mu_x}, \boldsymbol{\sigma_x})$

$\mu_z$    $\sigma_z$       $\mu_x$    $\sigma_x$

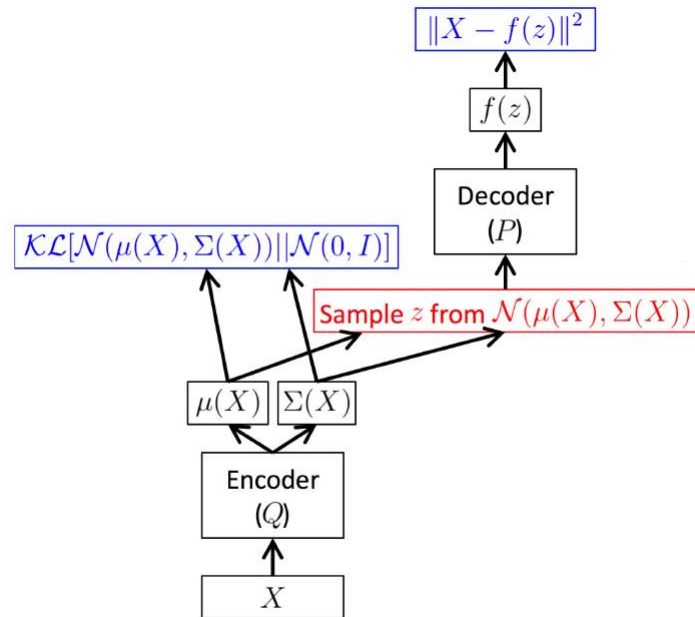**Encoder** $Q(Z|X; \phi)$

**Decoder** $P(X|Z; \theta)$

X

Z

*From CS231n, Fei-Fei Li, Justin Johnson, Serena Yeung*

**Forward and Backward Passes**

Georgia Tech

- Problem with respect to the VLB: updating $\phi$

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})} \left[ \log \frac{p_\theta(\boldsymbol{z}, \boldsymbol{x})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})} \right]$$

$$= -D_{\text{KL}}(q_\phi(\boldsymbol{z}|\boldsymbol{x}) \| p_\theta(\boldsymbol{z})) + \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})$$
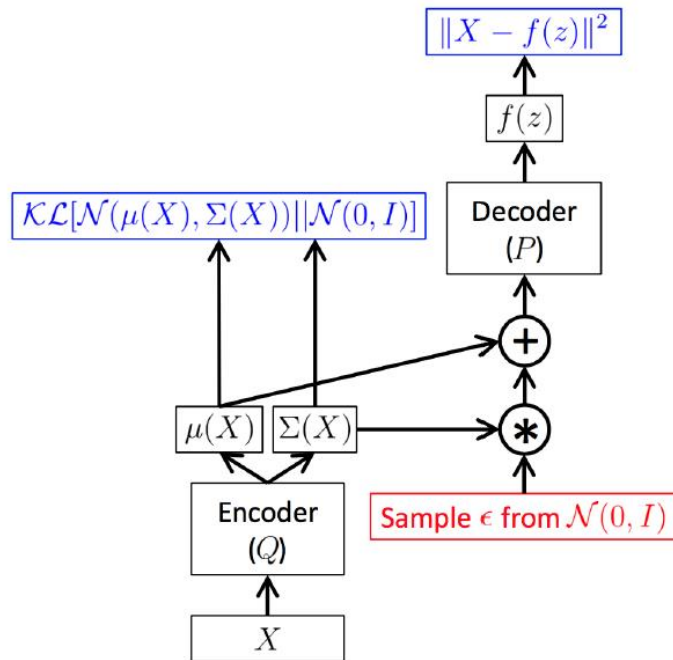
- $Z \sim Q(Z|X; \phi)$ : need to differentiate through the sampling process w.r.t $\phi$ (encoder is probabilistic)



*From: Tutorial on Variational Autoencoders*
*https://arxiv.org/abs/1606.05908*

*From: http://gokererdogan.github.io/2016/07/01/reparameterization-trick/*

**Problem**

- Solution: make the randomness independent of encoder output, making the encoder deterministic

- Gaussian distribution example:
  - Previously: encoder output = random variable $z \sim N(\mu, \sigma)$
  - Now encoder output = distribution parameter $[\mu, \sigma]$
  - $z = \mu + \epsilon * \sigma, \epsilon \sim N(0,1)$



*From: Tutorial on Variational Autoencoders*
*https://arxiv.org/abs/1606.05908*
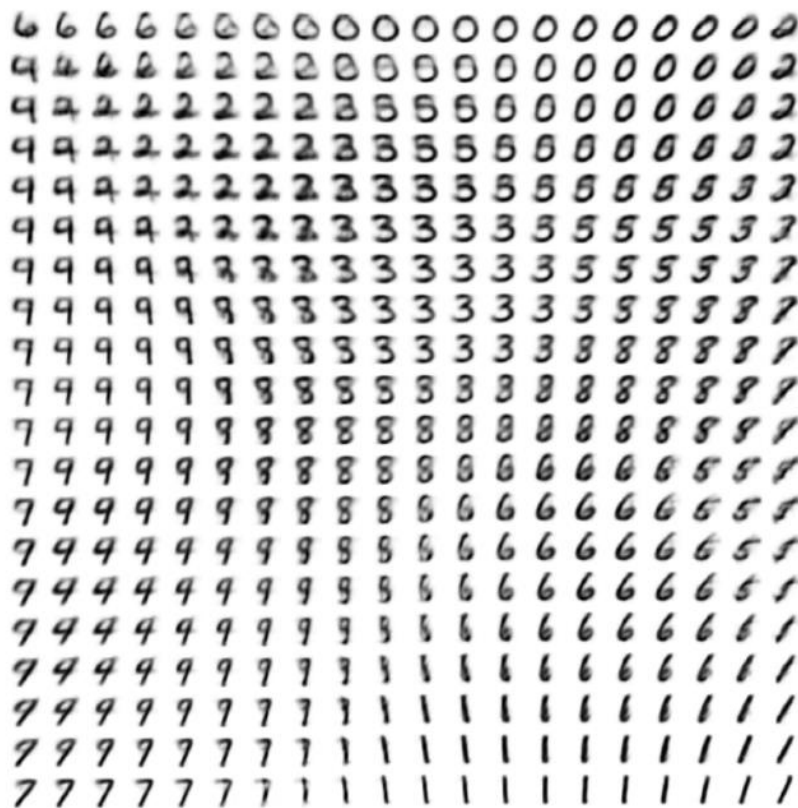
*From: http://gokererdogan.github.io/2016/07/01/reparameterization-trick/*

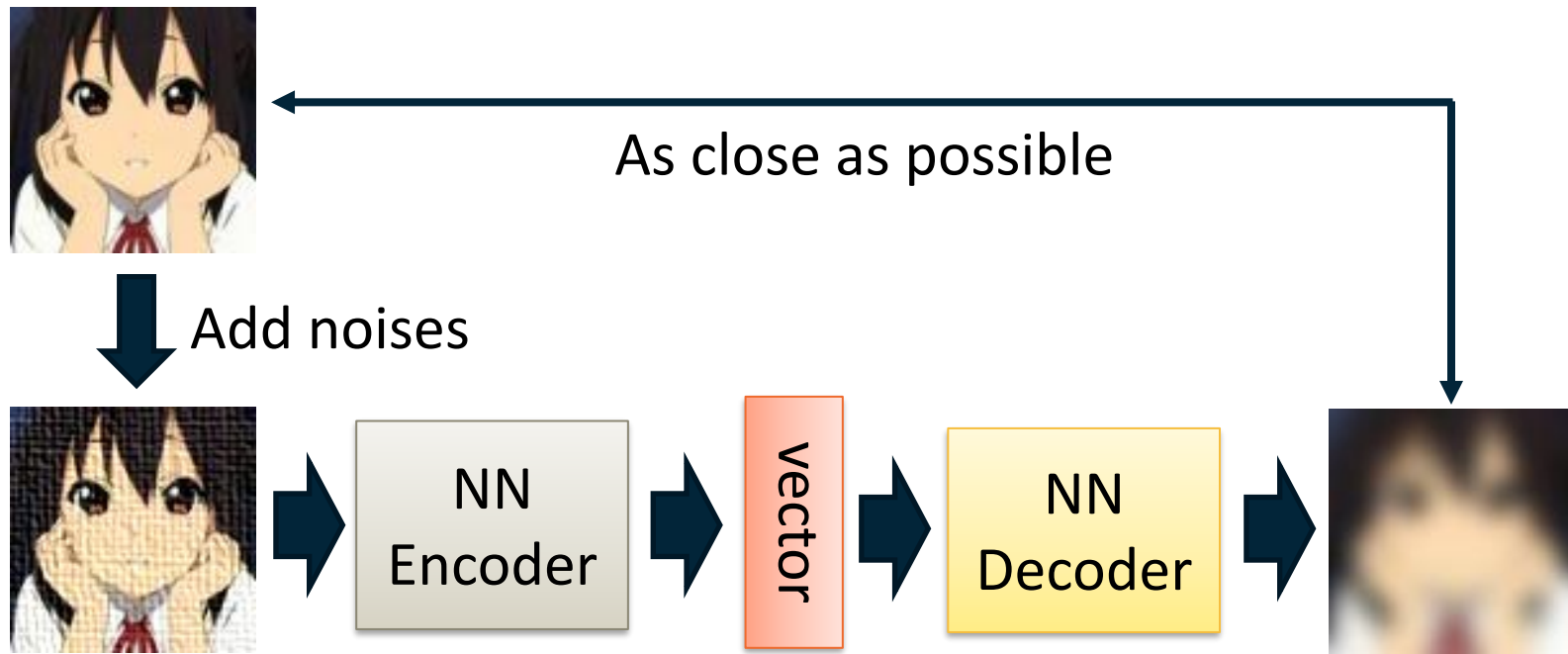**Reparameterization Trick: Solution**

$z_1$

$z_2$

*Kingma & Welling, Auto-Encoding Variational Bayes*

**Interpretability of Latent Vector**

Georgia Tech

- Variational Autoencoders (VAEs) provide a principled way to perform approximate maximum likelihood optimization
  - Requires some assumptions (e.g. Gaussian distributions)

- Samples are often not as competitive as diffusion models or GANs

- Latent features (learned in an unsupervised way!) often good for downstream tasks:
  - Example: World models for reinforcement learning (Ha et al., 2018)

*Ha & Schmidhuber, World Models, 2018*

**Summary**

# De-noising Auto-encoder



As close as possible

Add noises

NN Encoder → Vector → NN Decoder

Vincent, Pascal, et al. "Extracting and composing robust features with denoising autoencoders." *ICML,* 2008.

Georgia Tech

# Discrete Representation

- Vector Quantized Variational Auto-encoder (VQVAE)



Codebook
(a set of vectors)
Learn from data

(c.f. attention)
Compute similarity

The most similar one
is the input of decoder.

Slide by Hung-yi Lee

Georgia Tech

- Variational Autoencoders (VAEs) provide a principled way to perform approximate maximum likelihood optimization
  - Requires some assumptions (e.g. Gaussian distributions)

- Samples are often not as competitive as GANs

- Latent features (learned in an unsupervised way!) often good for downstream tasks:
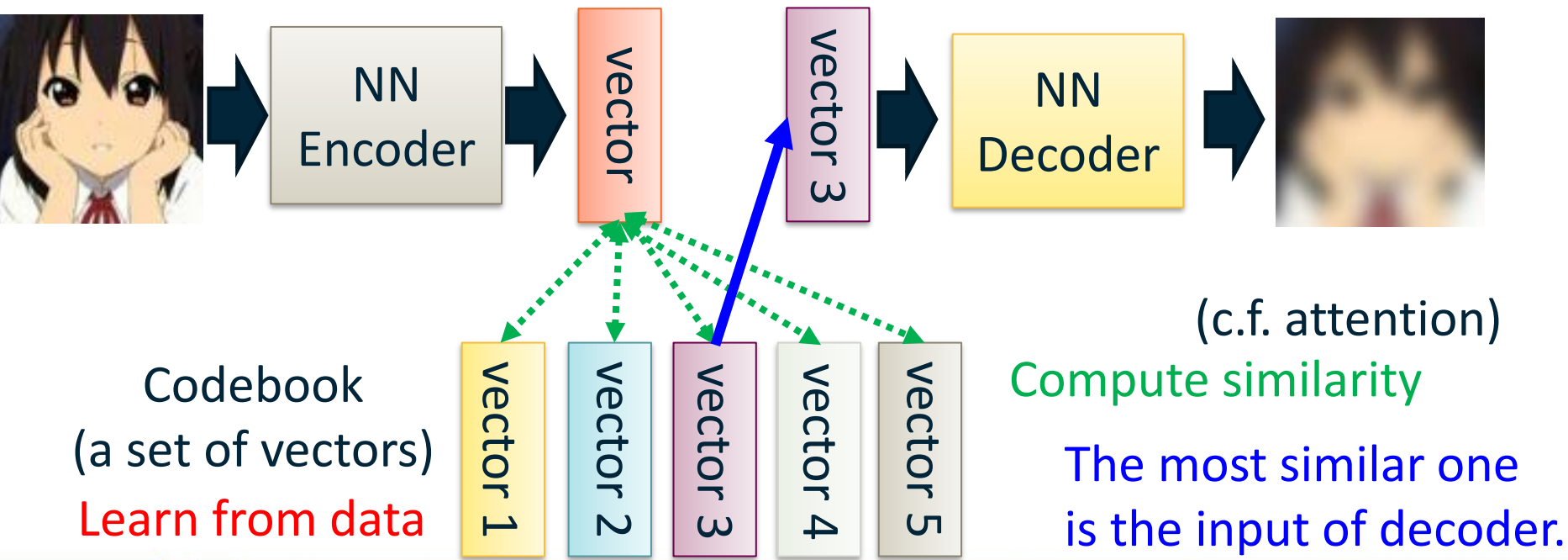  - Example: World models for reinforcement learning (Ha et al., 2018)

*Ha & Schmidhuber, World Models, 2018*

**Summary**

Georgia Tech

# Comparing the different generative models

Q. Which ones are VAEs good at?

| | Autoregressive (VAEs) | GANs | Diffusion |
|---|---|---|---|
| Mode coverage / diversity of generations | | | |
| Fast sampling | | | |
| High quality samples | | | |

# Comparing the different generative

VAEs are bad at generating high quality samples

|  | Autoregressive (VAEs) | GANs | Diffusion |
|---|---|---|---|
| Mode coverage / diversity of generations | ✅ | | |
| Fast sampling | ✅ | | |
| High quality samples | ❌ | | |

# Comparing the different generative models

Q. Which ones are GANs good at?

| | Autoregressive (VAEs) | GANs | Diffusion |
|---|---|---|---|
| Mode coverage / diversity of generations | ✅ | | |
| Fast sampling | ✅ | | |
| High quality samples | ❌ | | |

# Comparing the different generative models

GANs suffer from mode collapse

| | Autoregressive (VAEs) | GANs | Diffusion |
|---|---|---|---|
| Mode coverage / diversity of generations | ✅ | ❌ | |
| Fast sampling | ✅ | ✅ | |
| High quality samples | ❌ | ✅ | |

# Comparing the different generative models

Q. Which ones are Diffusion models good at?

| | Autoregressive (VAEs) | GANs | Diffusion |
|---|---|---|---|
| Mode coverage / diversity of generations | ✅ | ❌ | |
| Fast sampling | ✅ | ✅ | |
| High quality samples | ❌ | ✅ | |

# Comparing the different generative models

Diffusion models are bad at sampling fast.

| | Autoregressive (VAEs) | GANs | Diffusion |
|---|:---:|:---:|:---:|
| Mode coverage / diversity of generations | ✅ | ❌ | ✅ |
| Fast sampling | ✅ | ✅ | ❌ |
| High quality samples | ❌ | ✅ | ✅ |

- Several ways to learn *generative* models via deep learning

- **Generative Adversarial Networks (GANs):**
  - Pro: Amazing results across many image modalities
  - Con: Unstable/difficult training process, computationally heavy for good results
  - Con: Limited success for discrete distributions (language)
  - Con: Hard to evaluate (implicit model)
- **Variational Autoencoders:**
  - Pro: Principled mathematical formulation
  - Pro: Results in disentangled latent representations
  - Con: Approximation inference, results in somewhat lower quality reconstructions
- **Diffusion Models**
  - Pro: Great results and diversity!
  - Con: Slow generation (though lots of tricks to address)

*Ha & Schmidhuber, World Models, 2018*

## Overall Summary

Georgia Tech

# Comparison

# Plan Moving Forward

- Spring break!

- Guest lecture by Will Held on large language models!

- Reinforcement learning

- Open to other topics after:
  - Visualization and interpretability
  - Vision-language models
  - 3D / NeRFs
  - Robotics



The role of RLHF in ChatGPT