

10-423/623: Generative AI Lecture 6 – Generative Adversarial Networks and Variational Autoencoders

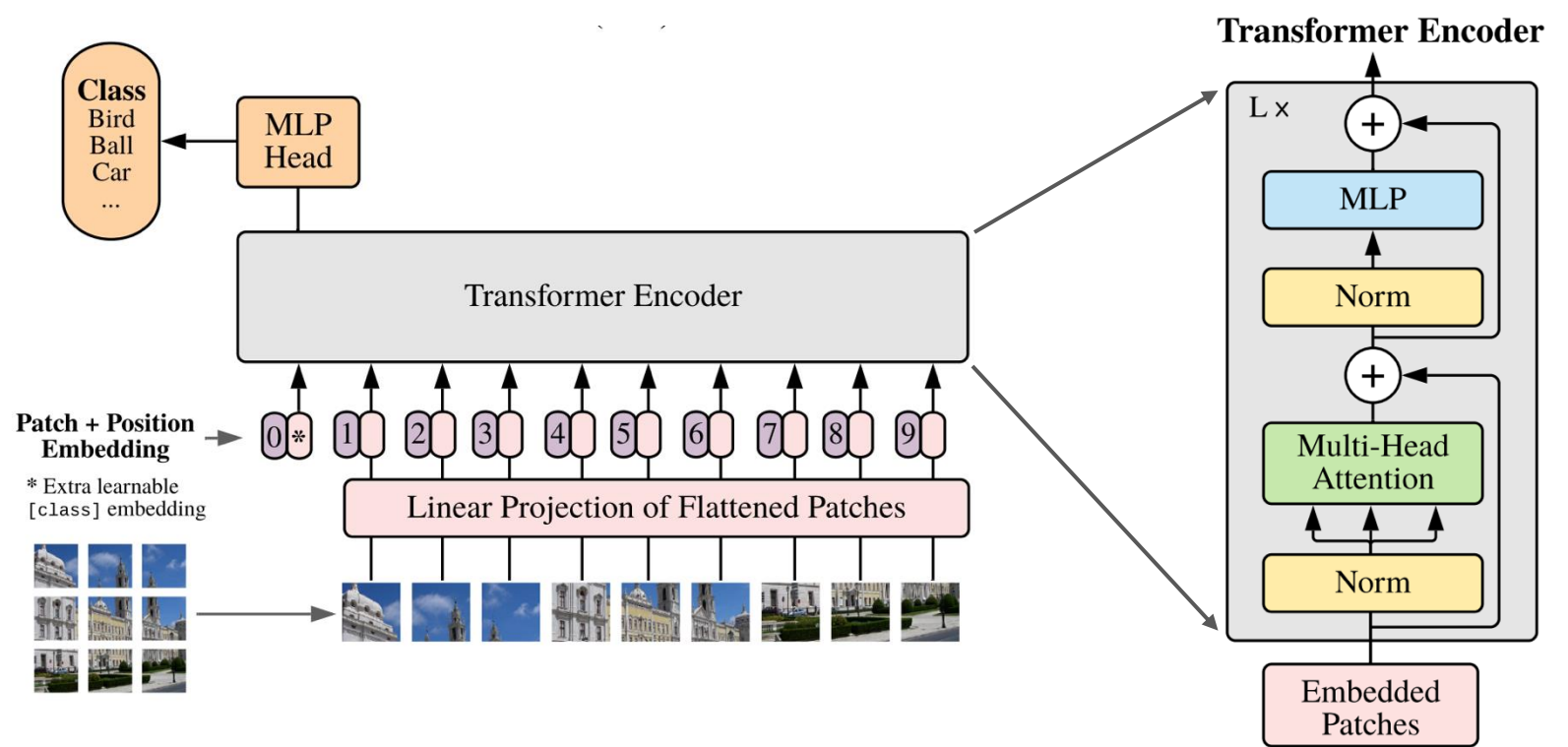
Henry Chai & Matt Gormley

9/16/24

Front Matter

- Announcements:
 - HW1 released 9/9, due 9/23 at 11:59 PM

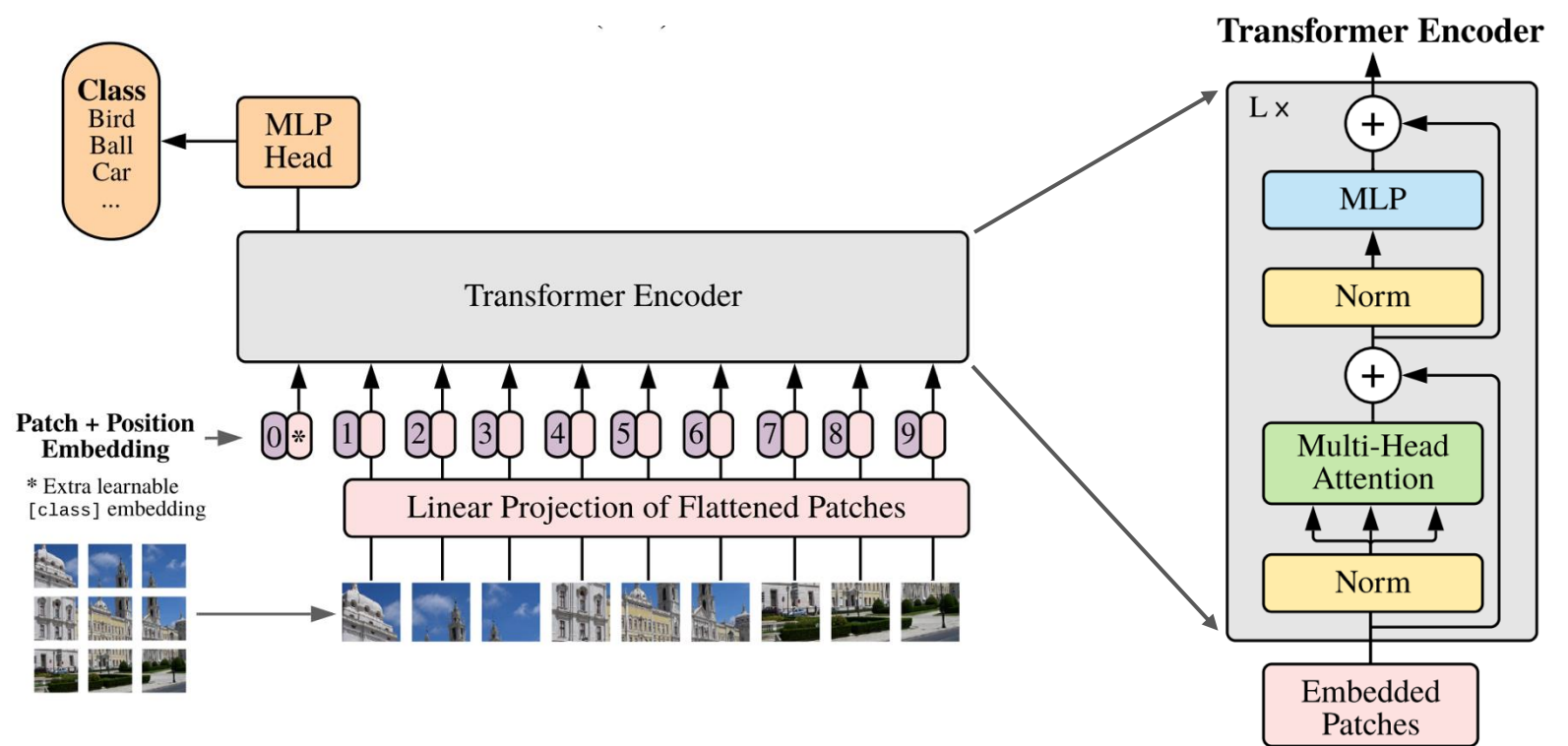
Recall: Vision Transformer (ViT)



- Instead of words as input, the inputs are $P \times P$ pixel *patches*
- Each patch is embedded linearly into a vector of size 1024
- Uses 1D positional embeddings
- Pre-trained on a large, supervised dataset (e.g., ImageNet 21K, JFT-300M)

Is this even a
generative
model?

Not
inherently...



- Instead of words as input, the inputs are $P \times P$ pixel *patches*
- Each patch is embedded linearly into a vector of size 1024
- Uses 1D positional embeddings
- Pre-trained on a large, supervised dataset (e.g., ImageNet 21K, JFT-300M)

Common Tasks in Computer Vision

- Image Classification
- Object Localization
- Object Detection
- Semantic Segmentation
- Instance Segmentation
- Image Captioning
- **Image Generation**
- Class-conditional generation
- Super resolution
- Image Editing
- Style transfer
- Text-to-image (TTI) generation

Image Generation

sea anemone

brain coral

slug



- Given a class label, sample a new image from that class
 - Image classification takes an image and predicts its label using $p(y|x)$
 - Class-conditional generation does this in reverse with $p(x|y)$

- **Class-conditional generation**
- Super resolution
- Image Editing
- Style transfer
- Text-to-image (TTI) generation

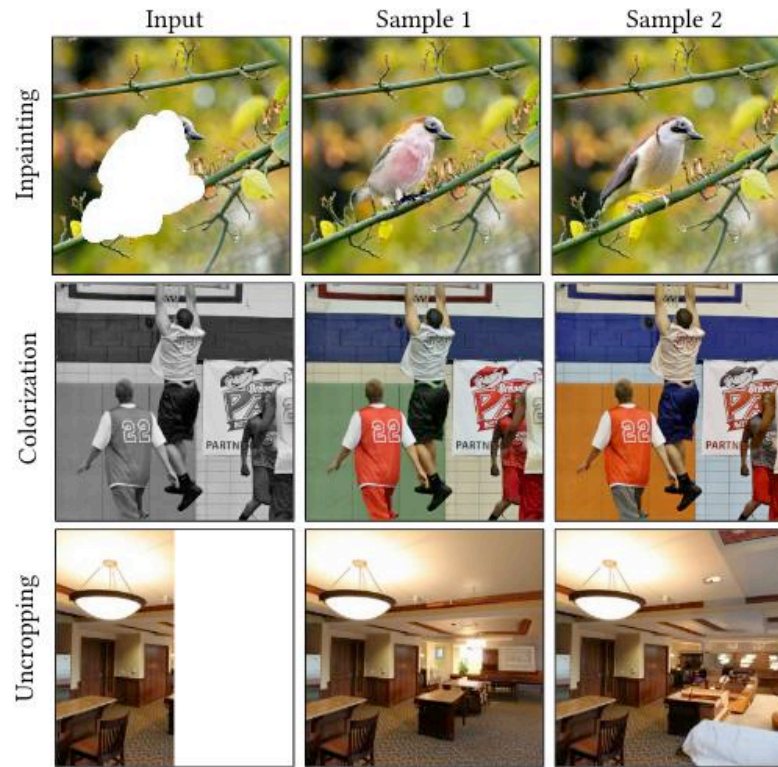
Image Generation



- Given a low-resolution image, generate a high-resolution reconstruction of the image

- Class-conditional generation
- **Super resolution**
- Image Editing
- Style transfer
- Text-to-image (TTI) generation

Image Generation



- **Inpainting** fills in the (pre-specified) missing pixels
- **Colorization** restores color to a greyscale image
- **Uncropping** creates a photo-realistic reconstruction of a missing side of an image

- Class-conditional generation
- Super resolution
- **Image Editing**

Image Generation



- Given two images, present the semantic content of the *source* image in the style of the *reference* image

- Class-conditional generation
- Super resolution
- Image Editing
- **Style transfer**
- Text-to-image (TTI) generation

Image Generation

Prompt: A propaganda poster depicting a cat dressed as french emperor napoleon holding a piece of cheese.



- Given a text description, sample an image that depicts the prompt

- Class-conditional generation
- Super resolution
- Image Editing
- Style transfer
- **Text-to-image (TTI) generation**

Image Generation

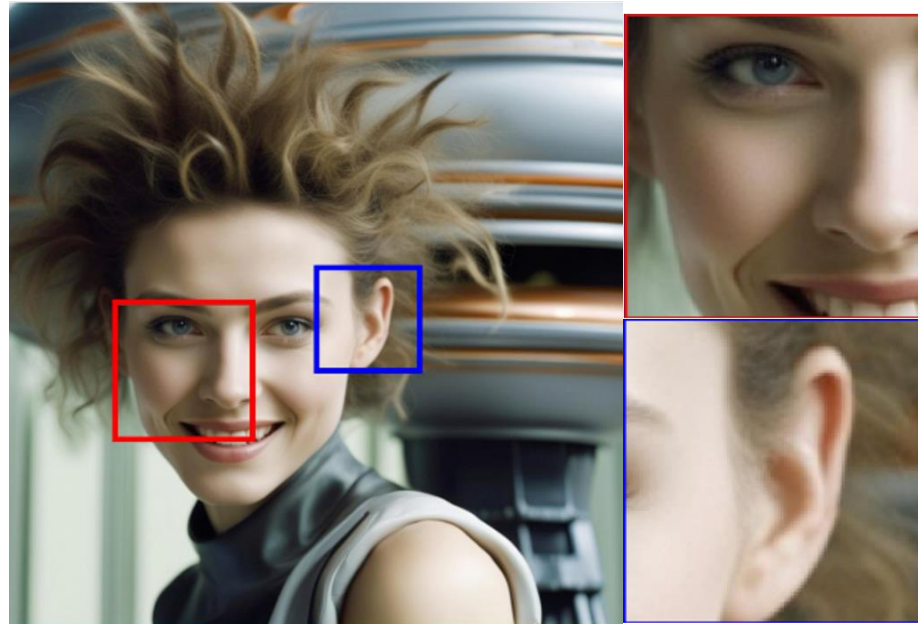
Prompt: Epic long distance cityscape photo of New York City flooded by the ocean and overgrown buildings and jungle ruins in rainforest, at sunset, cinematic shot, highly detailed, 8k, golden light



- Class-conditional generation
- Super resolution
- Image Editing
- Style transfer
- **Text-to-image (TTI) generation**

Image Generation

Prompt: close up headshot, futuristic young woman, wild hair sly smile in front of gigantic UFO, dslr, sharp focus, dynamic composition



- Class-conditional generation
- Super resolution
- Image Editing
- Style transfer
- **Text-to-image (TTI) generation**

Slide Generation?

Prompt: powerpoint slide explaining generative adversarial networks for a generative AI course, easy to follow, with a clear explanation of the objective function



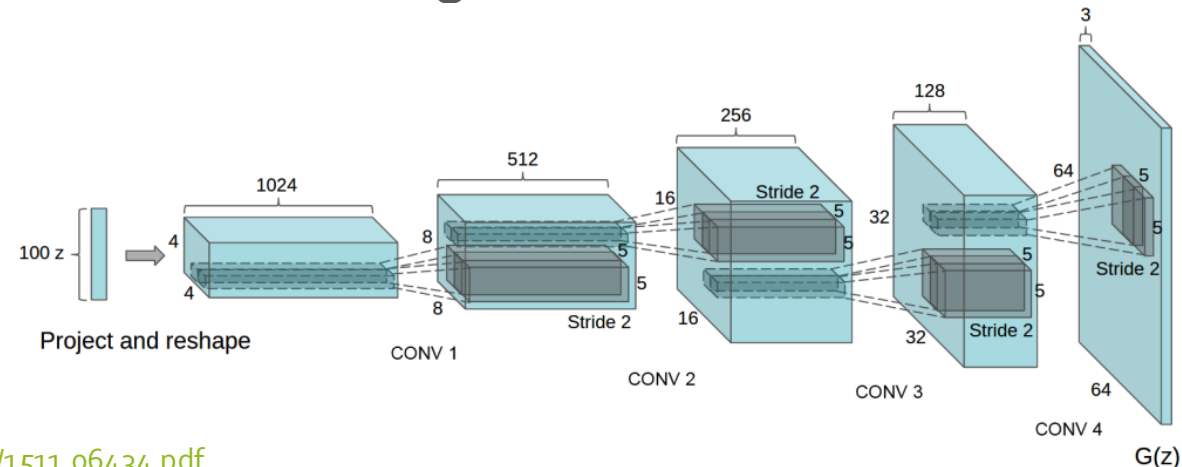
- Class-conditional generation
- Super resolution
- Image Editing
- Style transfer
- **Text-to-image (TTI) generation**

Generative Adversarial Networks (GANs)

- A GAN consists of two (deterministic) models:
 - a **generator** that takes a vector of random noise as input, and generates an image
 - a **discriminator** that takes in an image classifies whether it is real (label = 1) or fake (label = 0)
 - Both models are typically (but not necessarily) neural networks

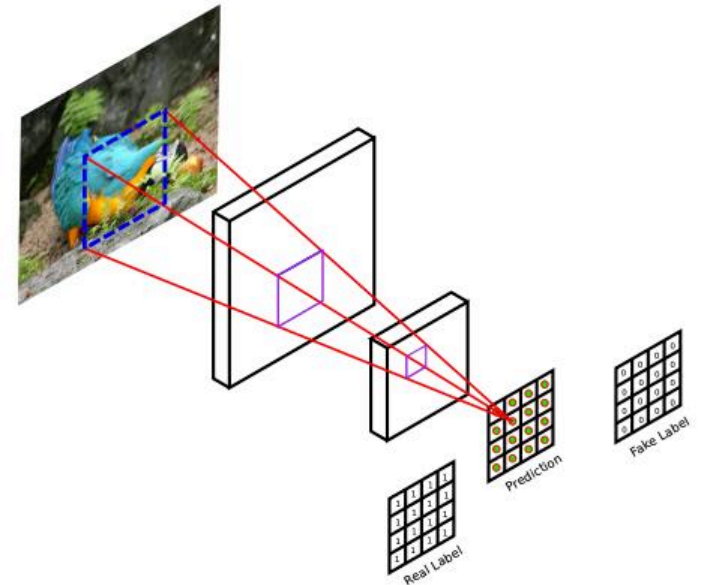
Generative Adversarial Networks (GANs)

- A GAN consists of two (deterministic) models:
 - a **generator** that takes a vector of random noise as input, and generates an image
- Example generator: DCGAN
 - An inverted CNN with four *fractionally-strided* convolution layers that grow the size of the image from layer to layer; final layer has three channels to generate color images



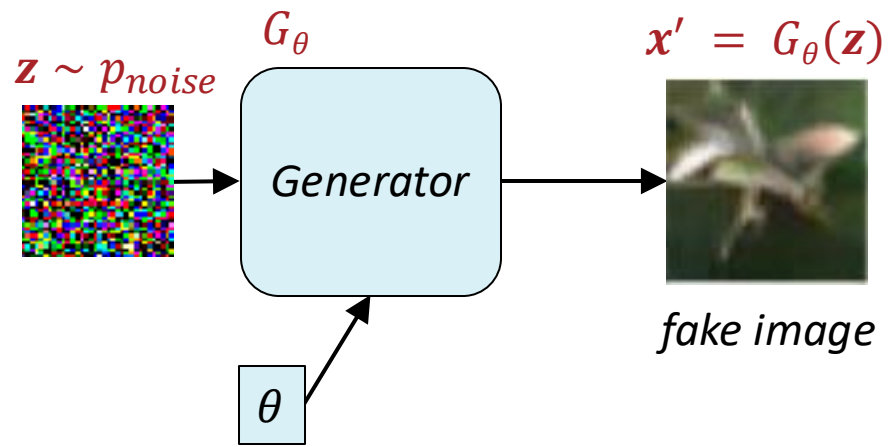
Generative Adversarial Networks (GANs)

- A GAN consists of two (deterministic) models:
 - a **generator** that takes a vector of random noise as input, and generates an image
 - a **discriminator** that takes in an image and classifies whether it is real (label = 1) or fake (label = 0)
- Example discriminator: PatchGAN
 - Traditional CNN that looks at each patch of the image and tries to predict whether it is real or fake; can help encourage the generator to avoid creating blurry images



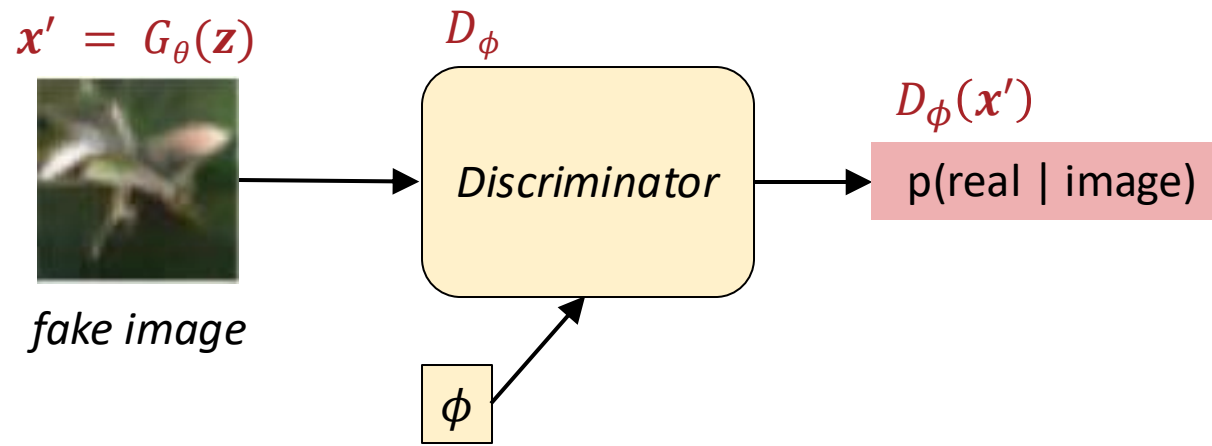
Generative Adversarial Networks (GANs): Training

- A GAN consists of two (deterministic) models:
 - a **generator** that takes a vector of random noise as input, and generates an image
 - a **discriminator** that takes in an image classifies whether it is real (label = 1) or fake (label = 0)
 - Both models are typically (but not necessarily) neural networks
- During training, the GAN plays a two-player minimax game: the generator tries to create realistic images to fool the discriminator and the discriminator tries to identify the real images from the fake ones

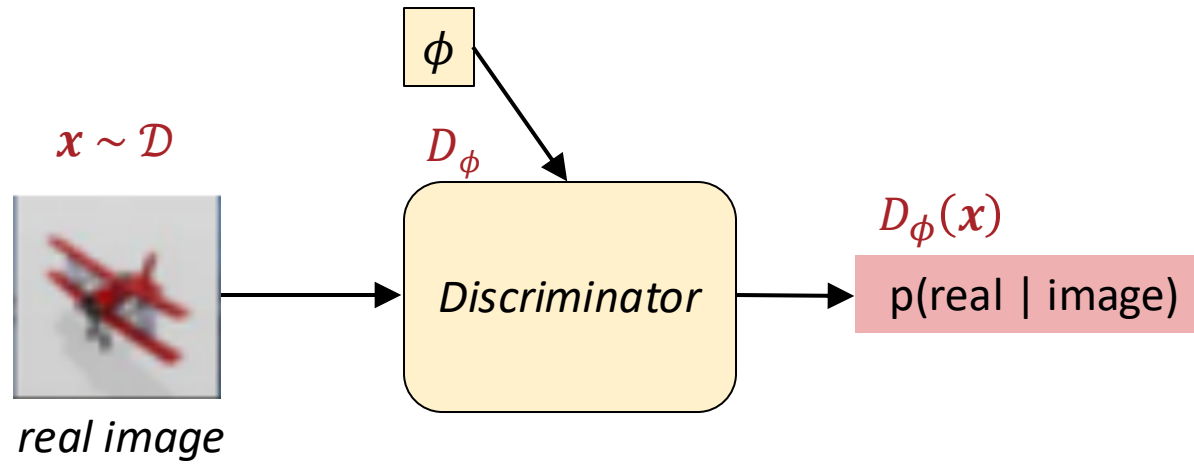


Typically, p_{noise} is a standard Gaussian i.e., $\mathcal{N}(\mathbf{0}, \sigma^2 I)$

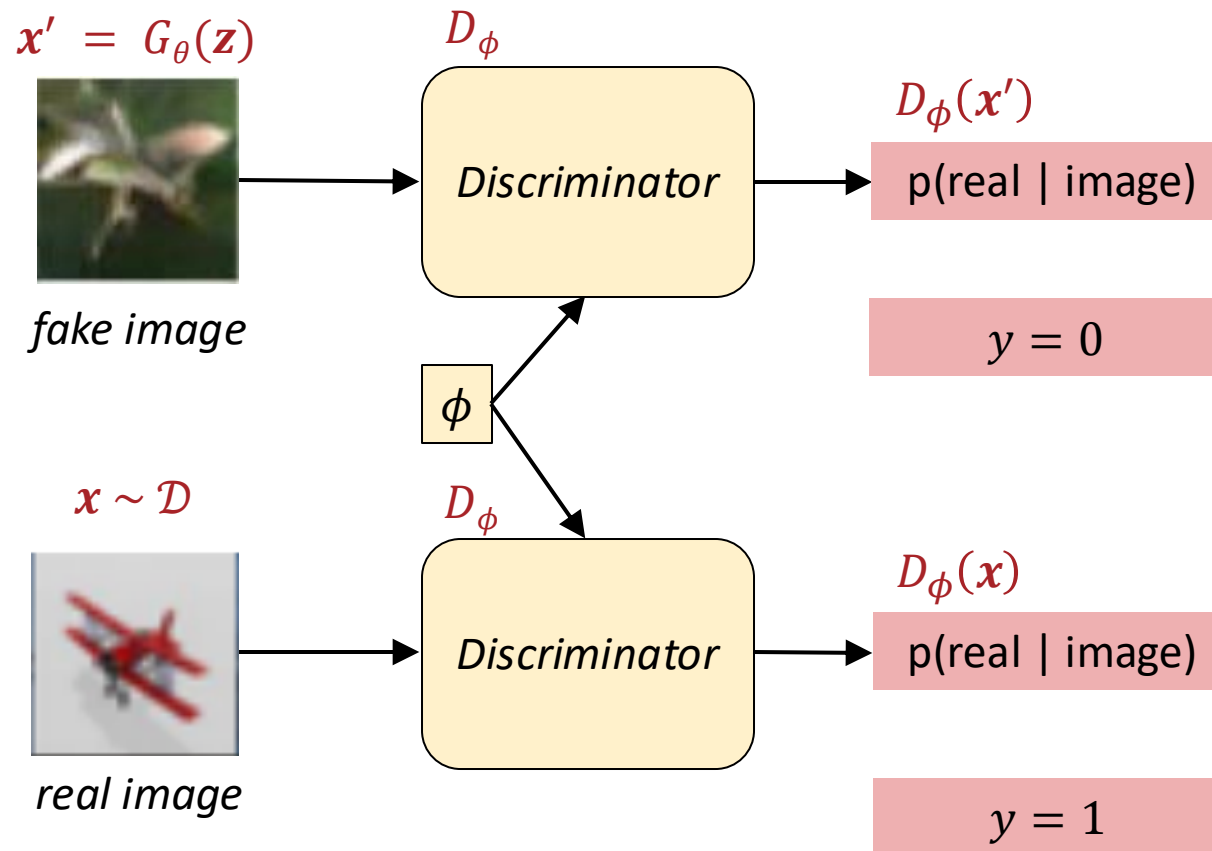
GANs: Architecture



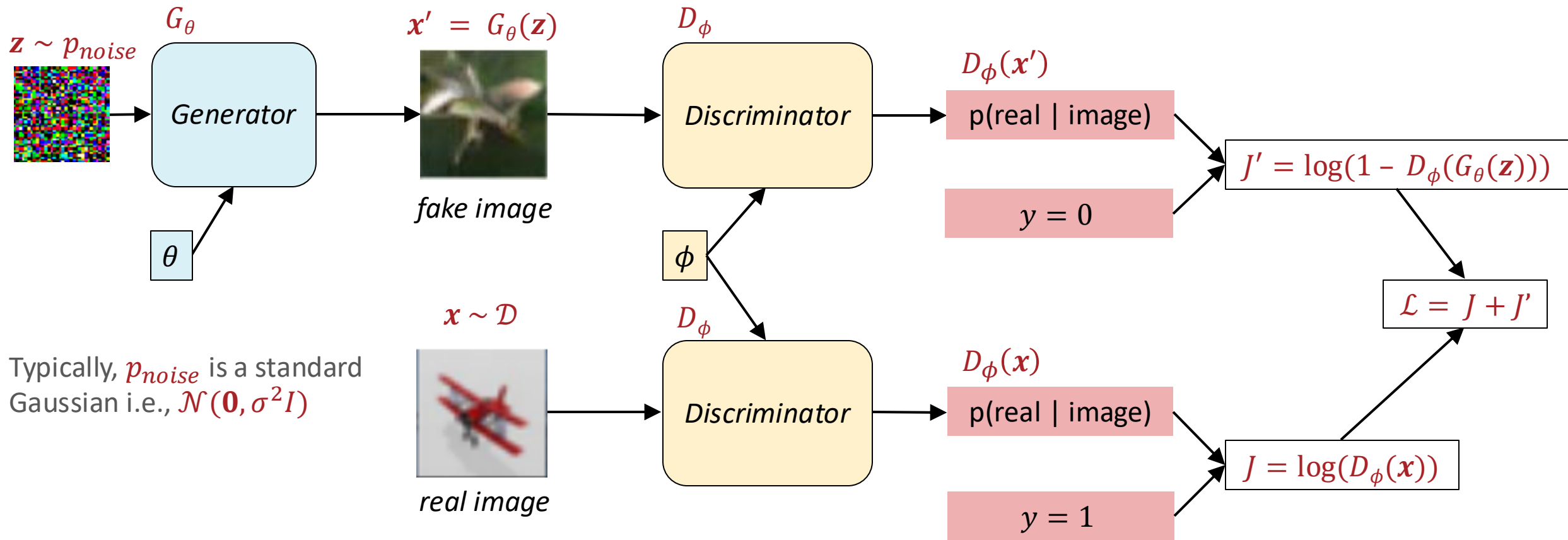
GANs: Architecture



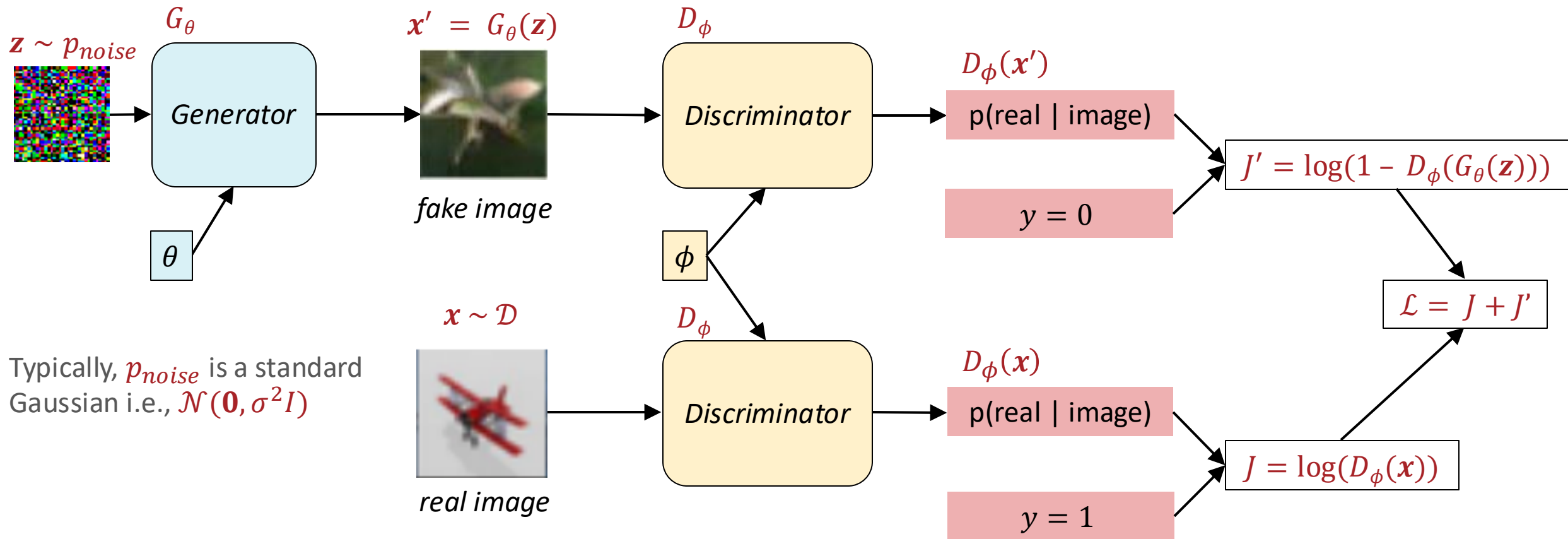
GANs: Architecture



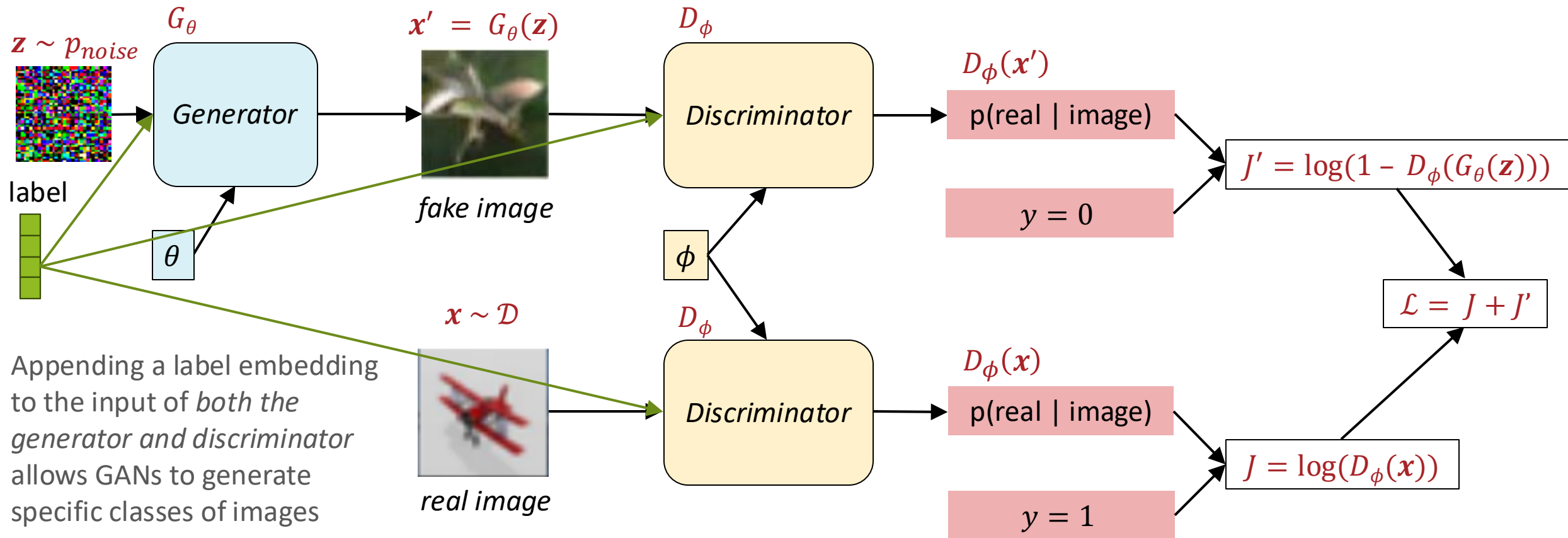
GANs: Architecture



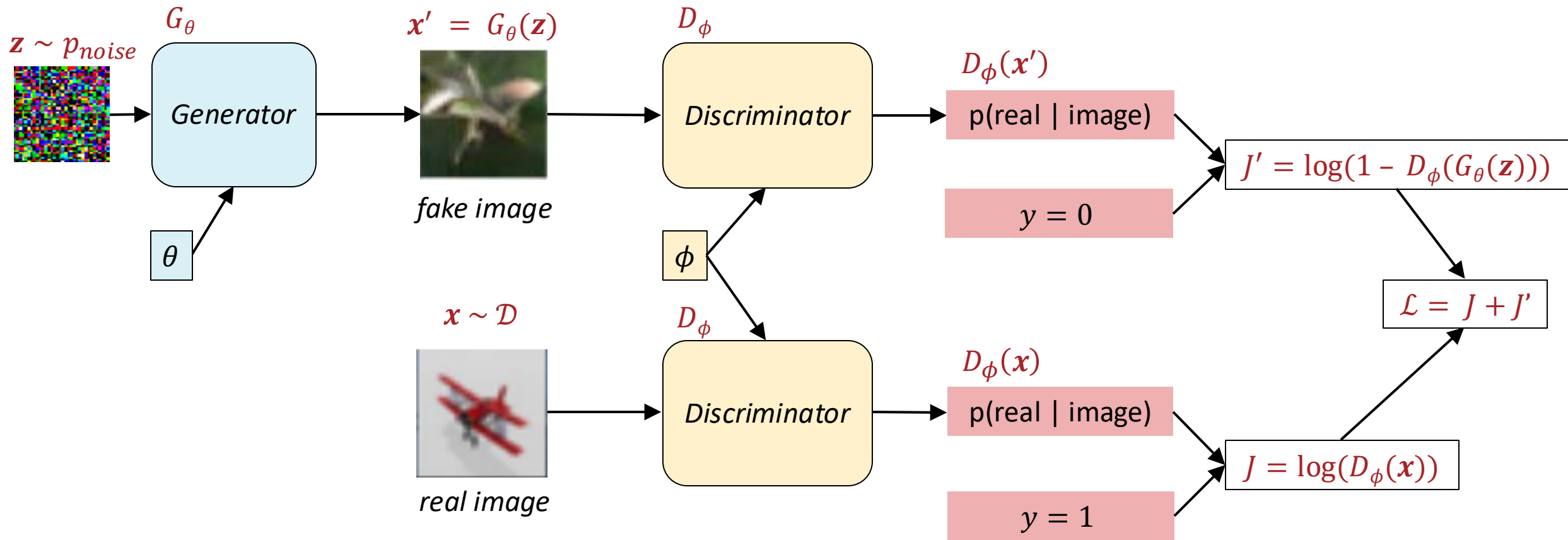
GANs: Architecture



Can we backpropagate through G_θ given that \mathbf{z} is stochastic?



Class-conditional GANs

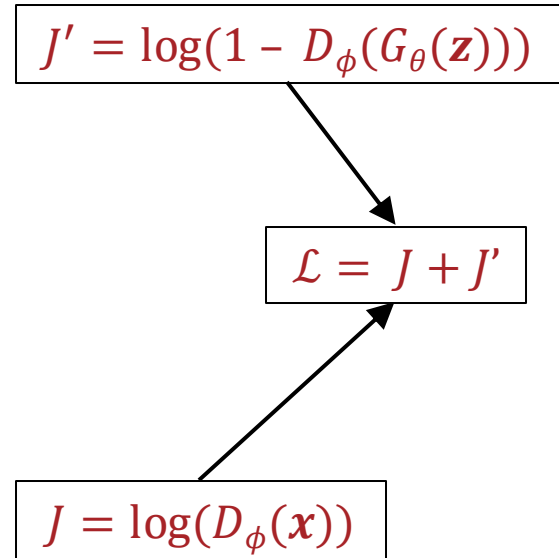


So how do we go about training one of these things?

The discriminator is trying to minimize the usual cross-entropy loss for binary classification with labels {real = 1, fake = 0}

$$\min_{\phi} \log \left(D_{\phi}(\mathbf{x}^{(i)}) \right) + \log \left(1 - D_{\phi}(G_{\theta}(\mathbf{z}^{(i)})) \right)$$
$$\max_{\theta} \log \left(1 - D_{\phi}(G_{\theta}(\mathbf{z}^{(i)})) \right)$$

The generator is trying to maximize the likelihood of its generated (fake) image being classified as real, according to a fixed discriminator



GANs: Training

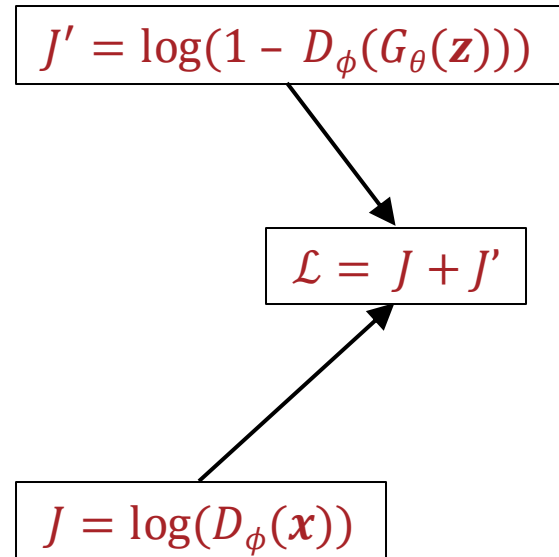
Both objectives (and hence, their sum) are differentiable!

$$\min_{\phi} \log \left(D_{\phi}(\mathbf{x}^{(i)}) \right) + \log \left(1 - D_{\phi}(G_{\theta}(\mathbf{z}^{(i)})) \right)$$

$$\max_{\theta} \log \left(1 - D_{\phi}(G_{\theta}(\mathbf{z}^{(i)})) \right)$$

Training alternates between:

1. Keeping θ fixed and backpropagating through D_{ϕ}
2. Keeping ϕ fixed and backpropagating through G_{θ}



GANs: Training

GANs: Training

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log \left(1 - D(G(z^{(i)})) \right) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D(G(z^{(i)})) \right).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

- Optimization is like block coordinate descent but instead of exact optimization, we take a step of mini-batch SGD

But what about those Vision Transformer things we talked about last week?

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log \left(1 - D(G(z^{(i)})) \right) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

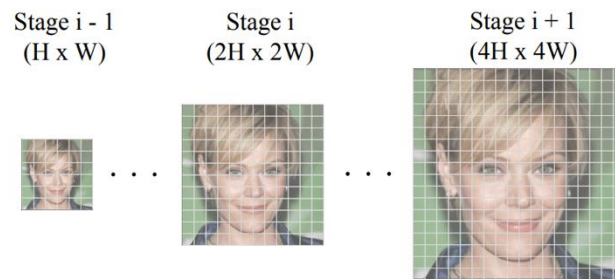
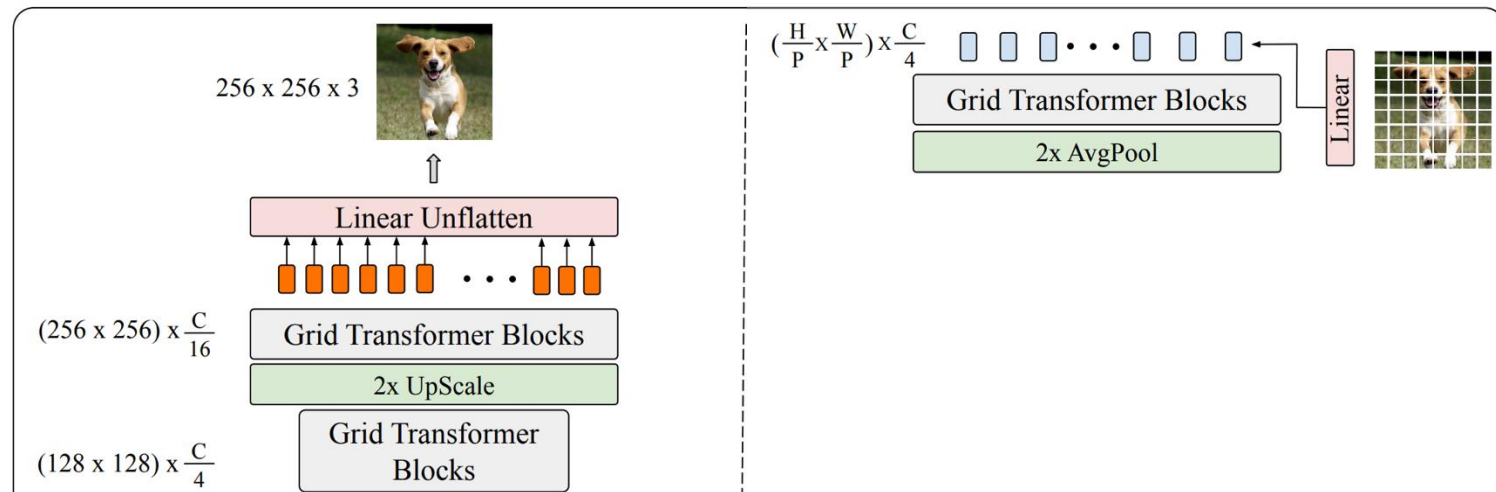
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D(G(z^{(i)})) \right).$$

end for

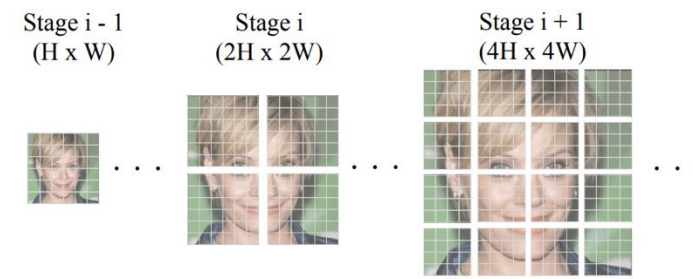
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

- Optimization is like block coordinate descent but instead of exact optimization, we take a step of mini-batch SGD

TransGANs



(a) Standard Self-Attention



(b) Grid Self-Attention

Figure 3: Grid Self-Attention across different transformer stages. We replace Standard Self-Attention with Grid Self-Attention when the resolution is higher than 32×32 and the grid size is set to be 16×16 by default.

ViTGANs

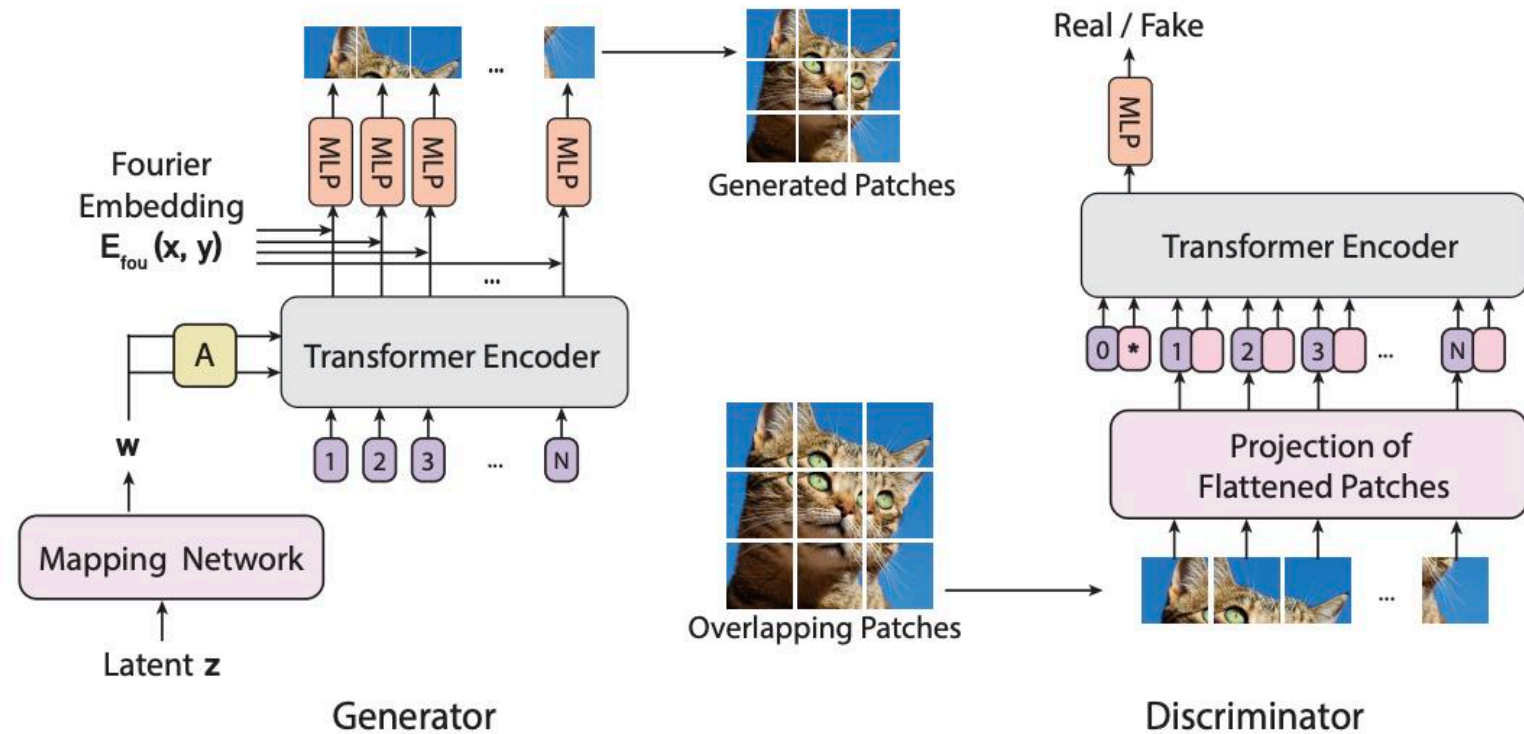
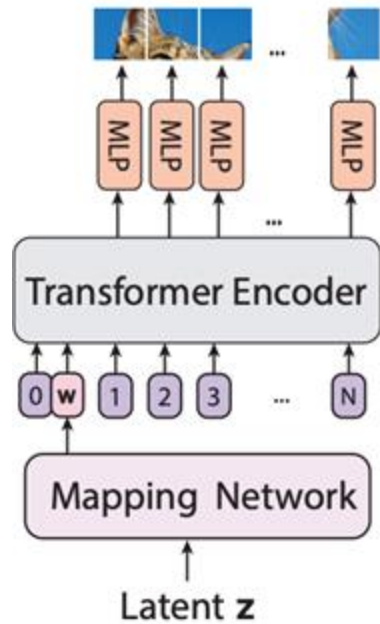
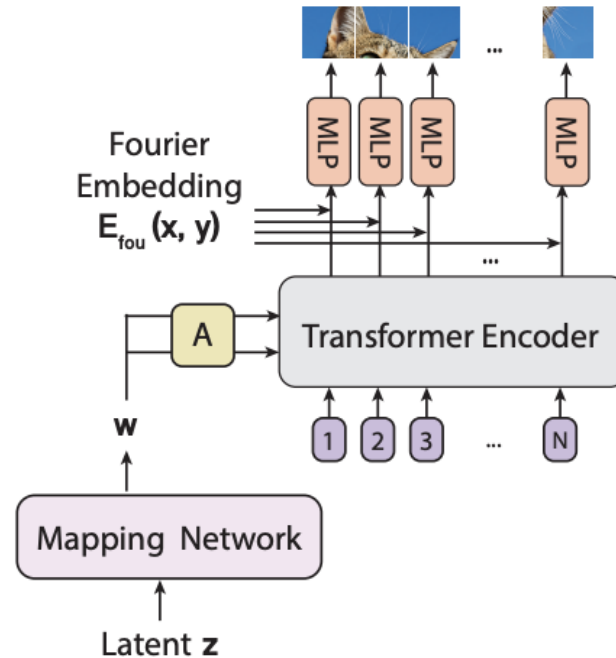


Figure 1: **Overview of the proposed ViTGAN framework.** Both the generator and the discriminator are designed based on the Vision Transformer (ViT). Discriminator score is derived from the classification embedding (denoted as [*] in the Figure). The generator generates pixels patch-by-patch based on patch embeddings.

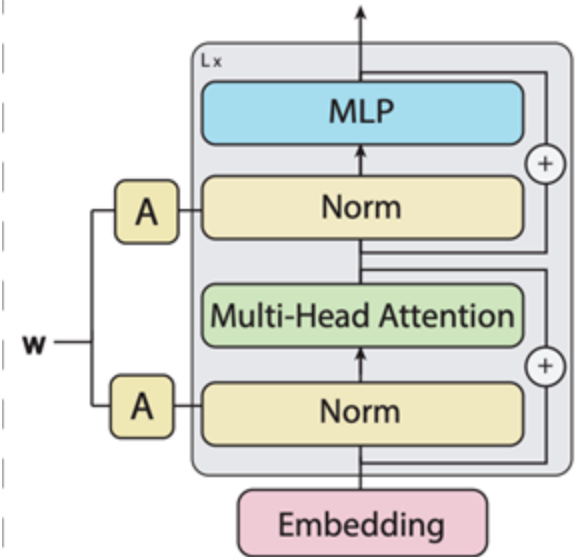
ViTGANs



"Plausible Baseline"

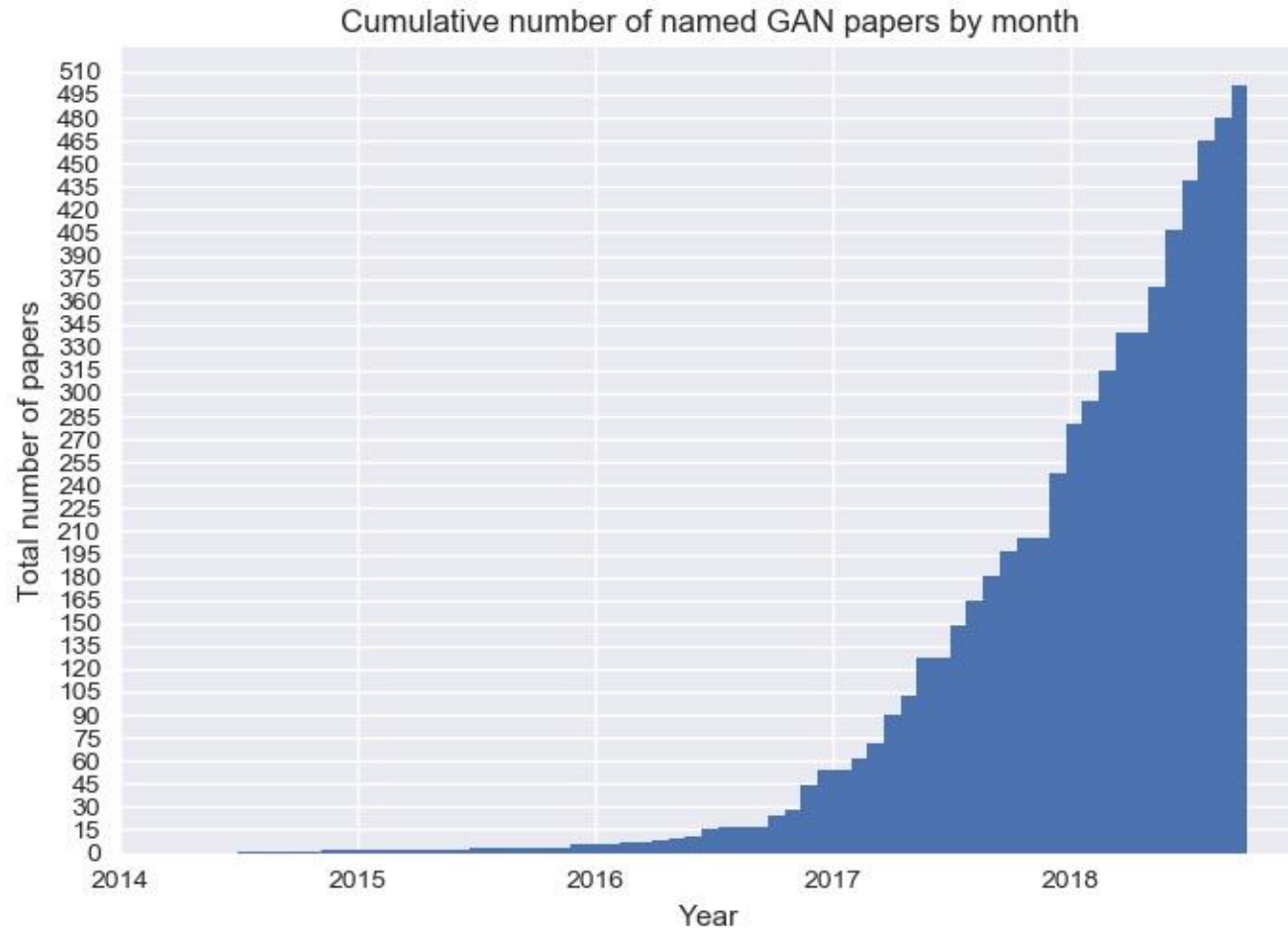


ViTGAN generator

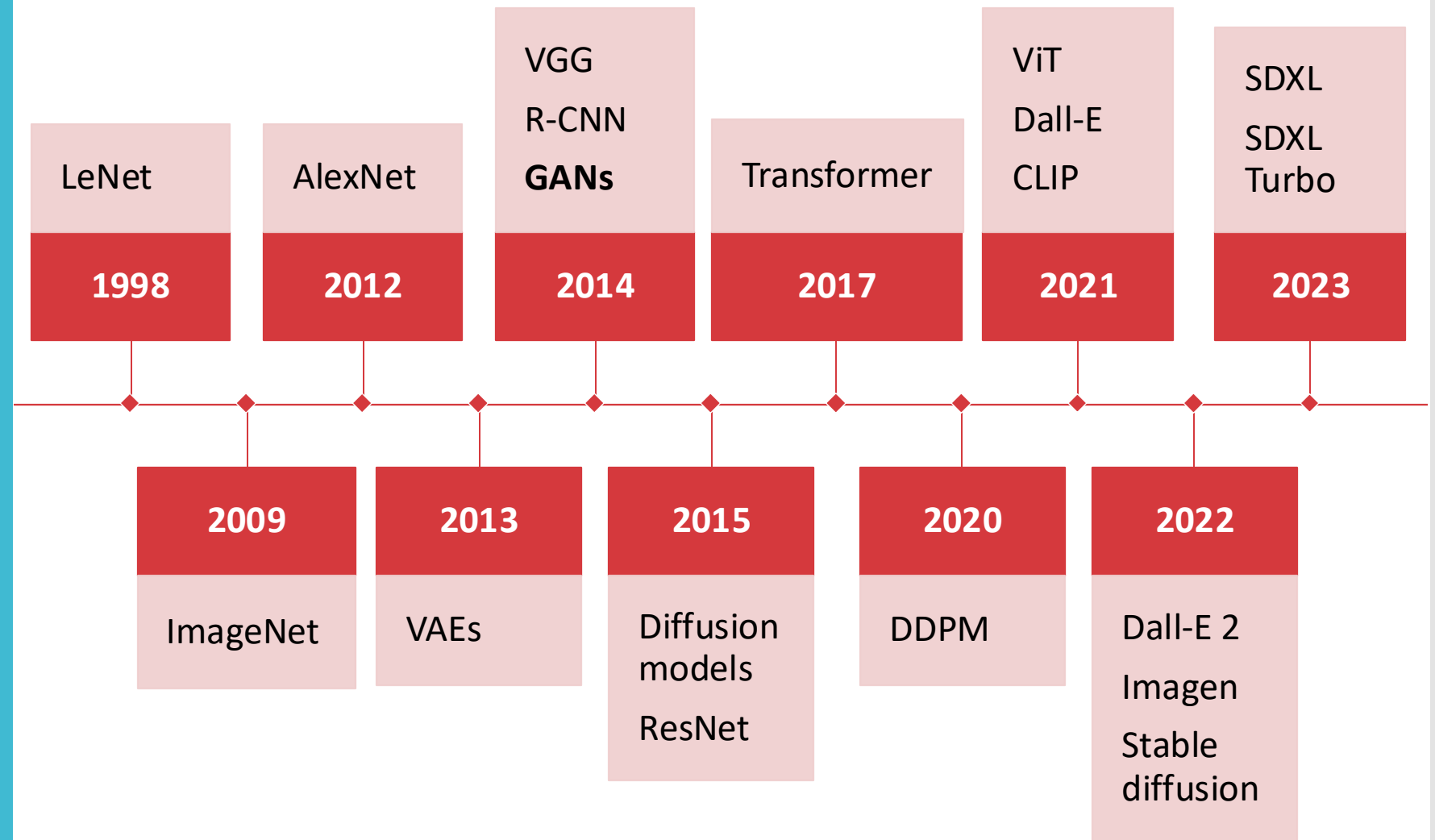


Self-Modulated LayerNorm
in Transformer Encoder

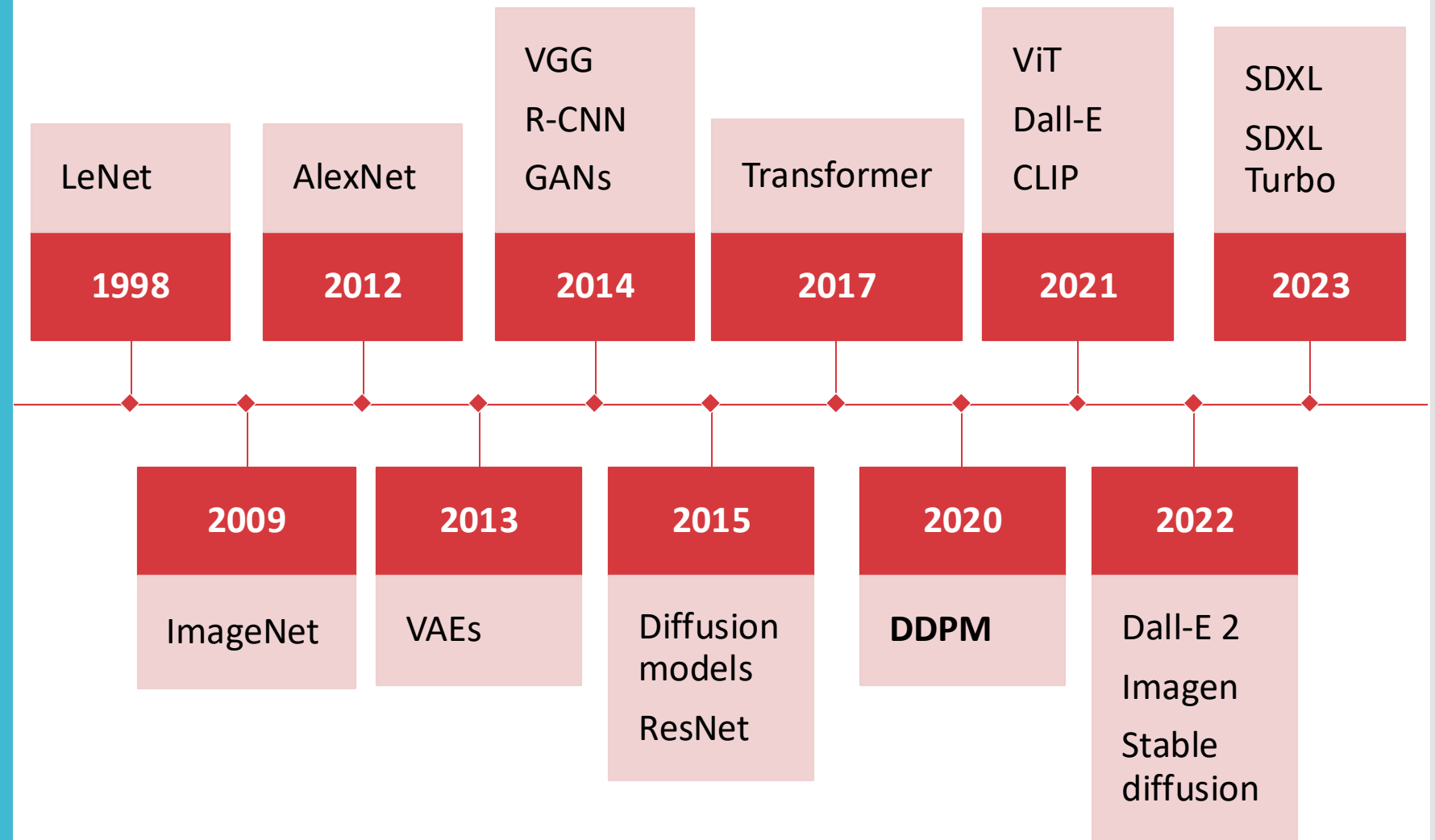
GANs Everywhere!



Recall: Computer Vision Timeline



Recall: Computer Vision Timeline



GANs vs. Diffusion



GAN generated images



Training images



Diffusion generated images

Recall: Computer Vision Timeline

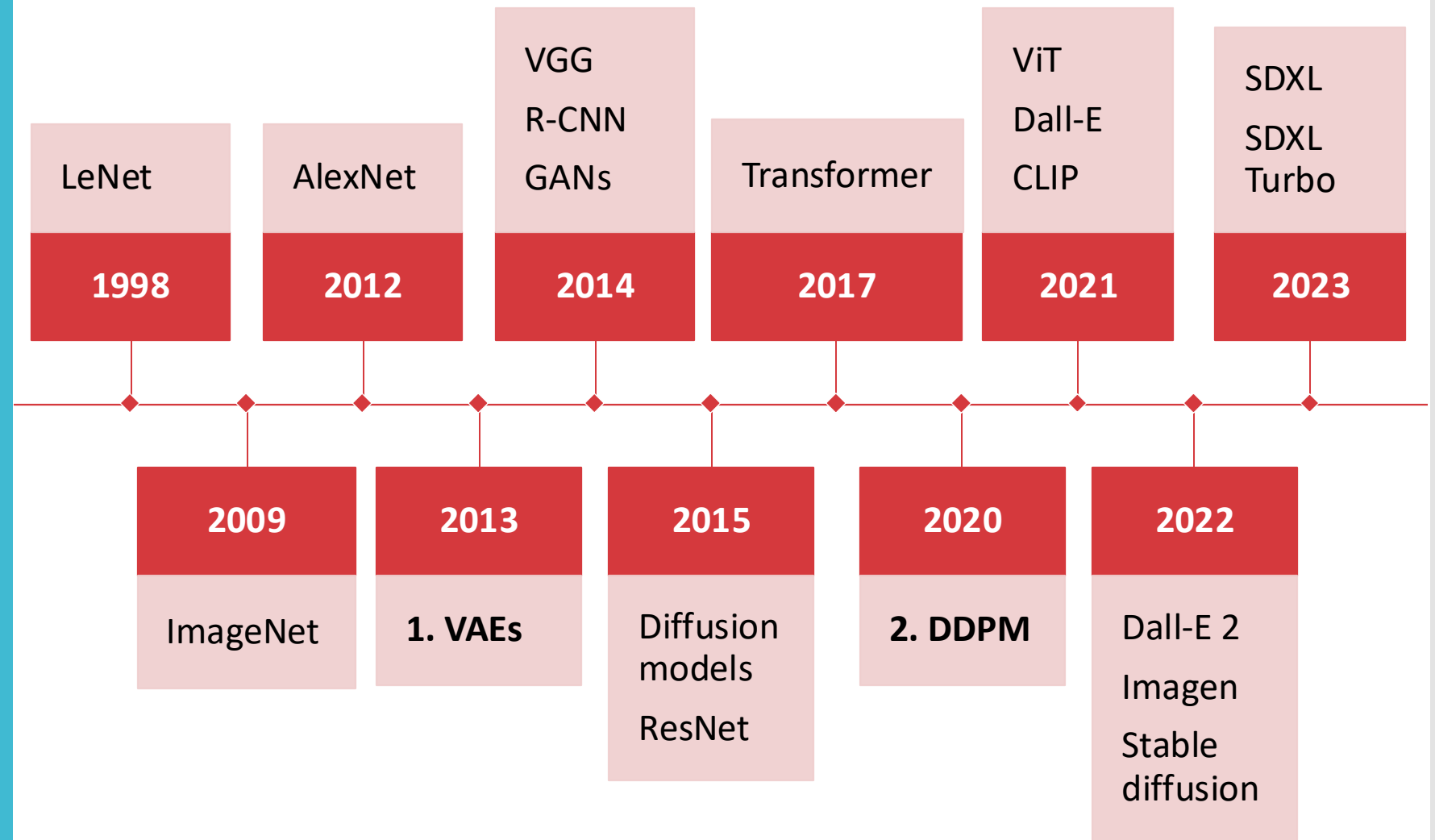
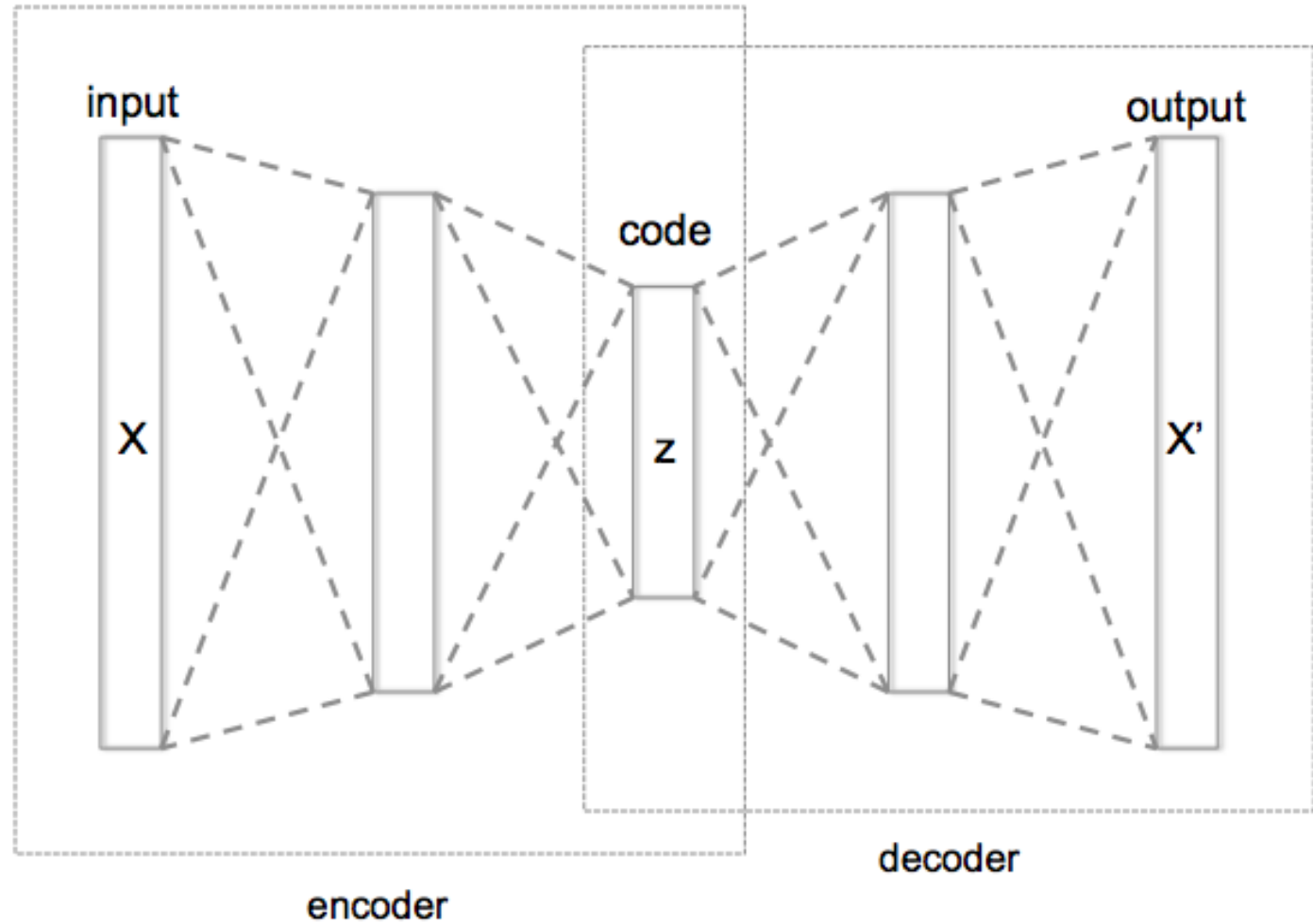


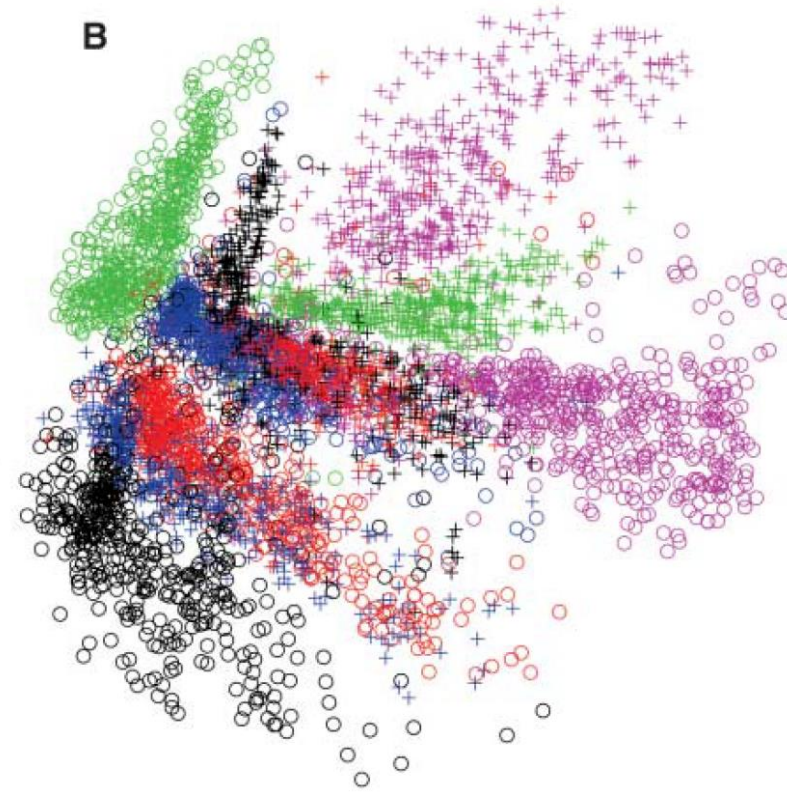
Image Generation

- Fundamental challenge: images are incredibly high-dimensional objects with complex relationships between elements
- Idea: learn a low-dimensional representation of images, sample points in the low-dimensional space and project them up to the original image space

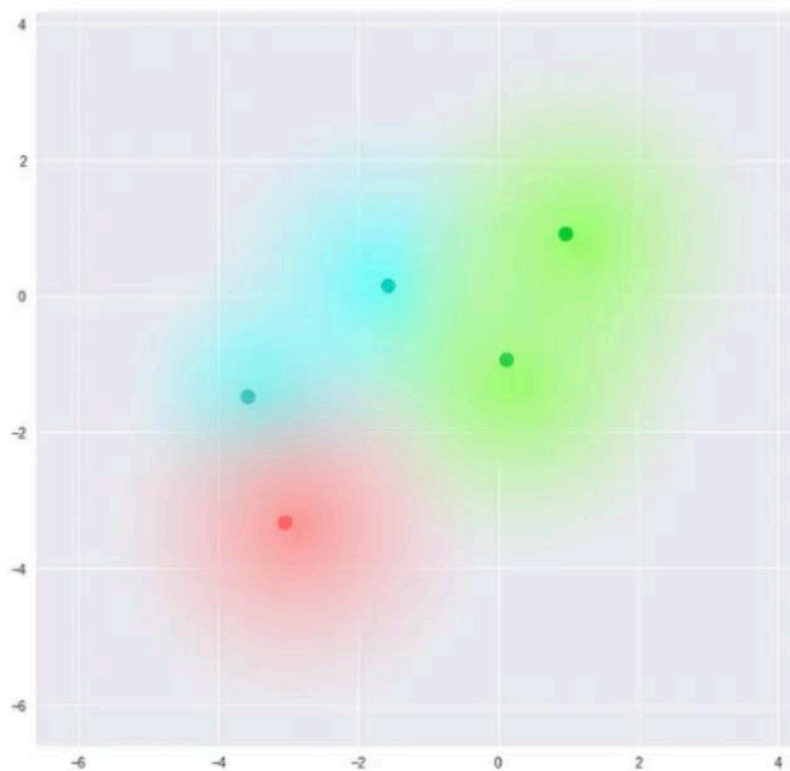
Recall: Autoencoders



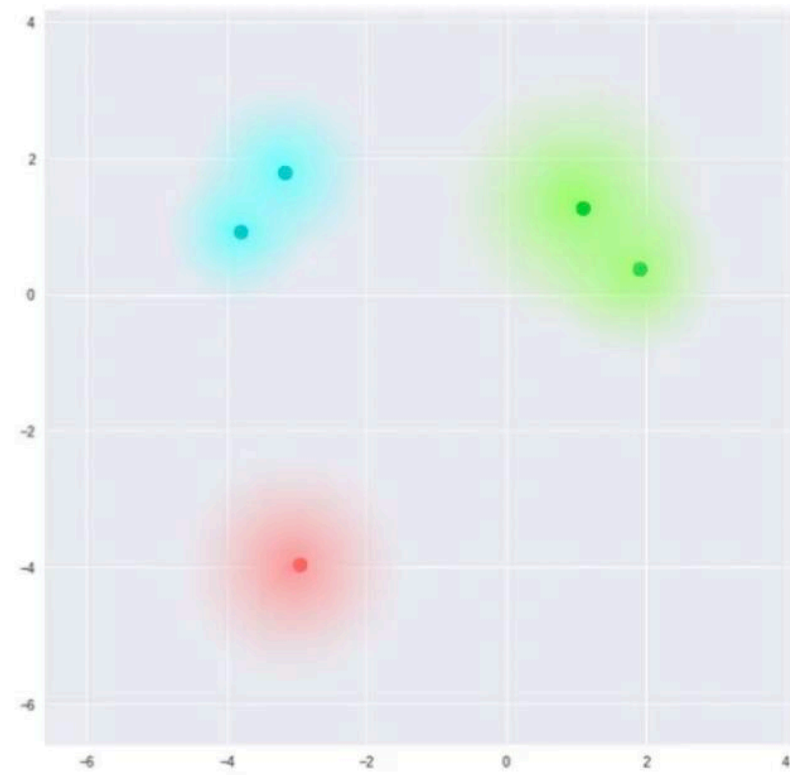
- Issue: latent space is sparse...
 - Sampling from latent space of an autoencoder creates outputs that are effectively identical to images in the training dataset



Autoencoder Latent Space



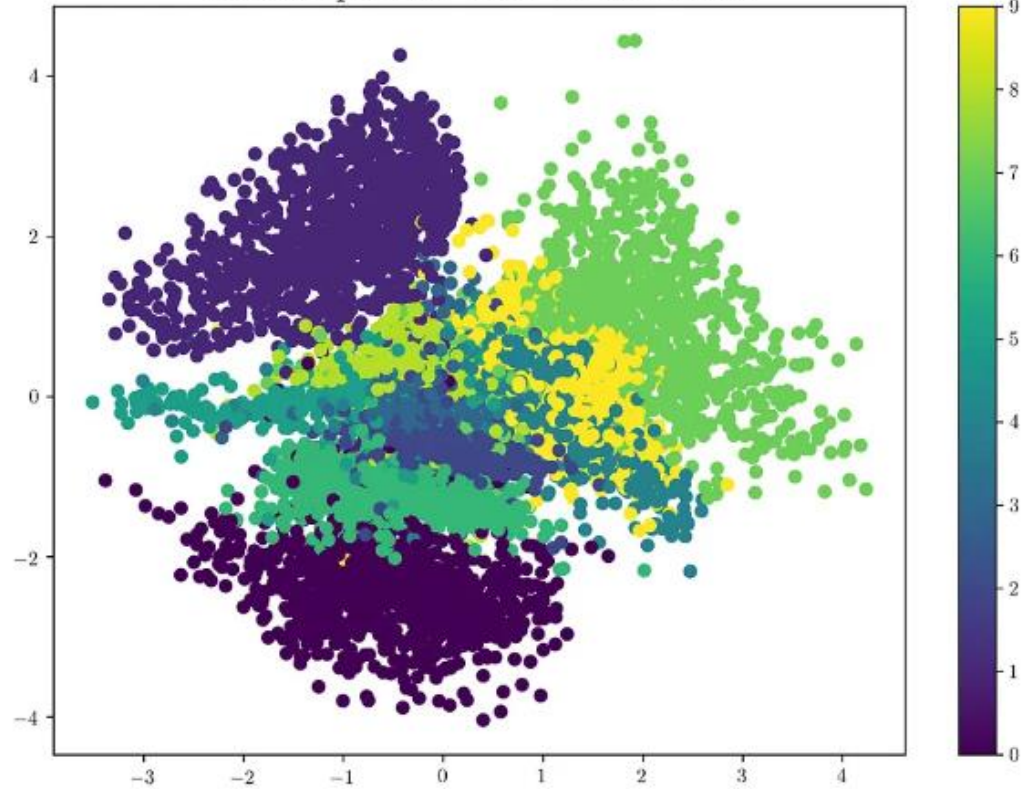
What we require



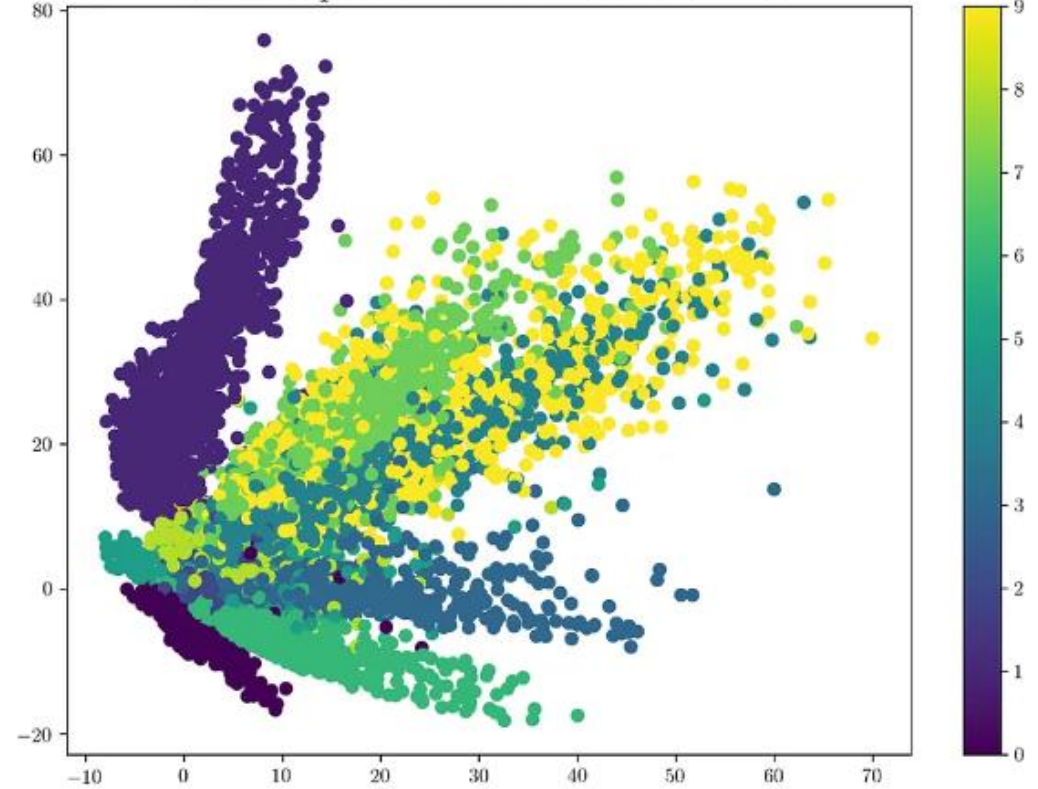
What we may inadvertently end up with

Autoencoder Latent Space

Latent space of VAE with KL loss

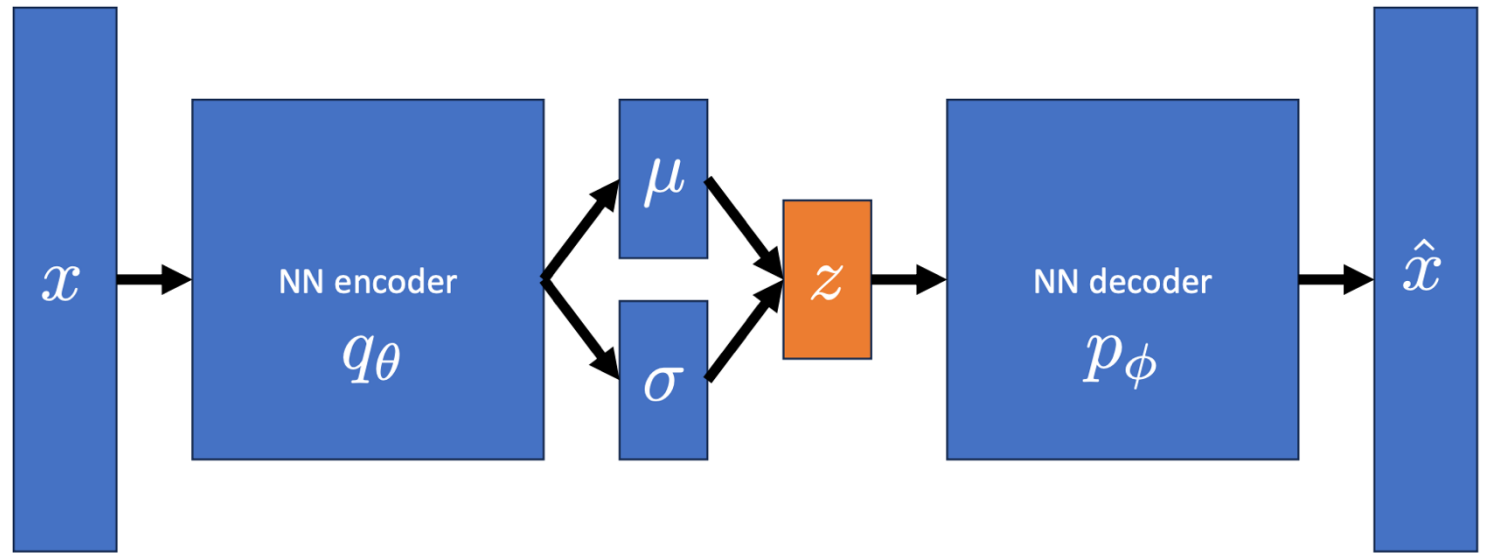


Latent space of VAE without KL loss

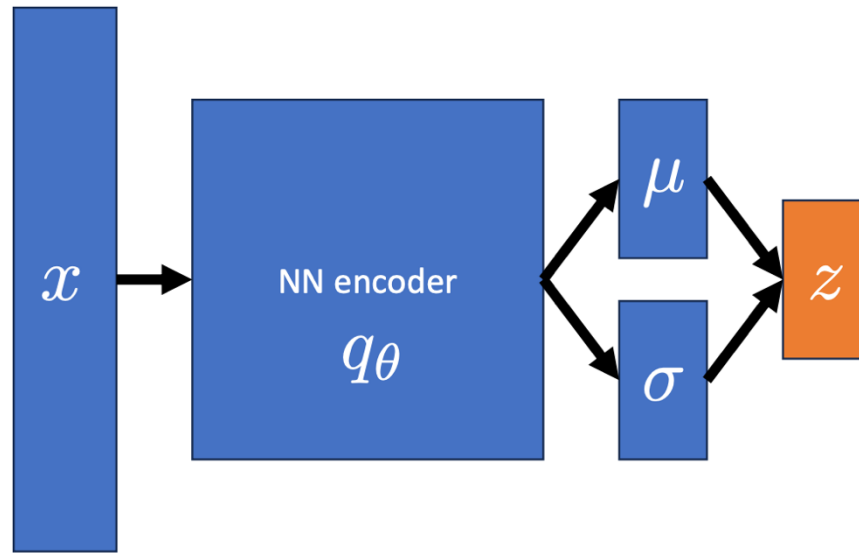


Variational Autoencoder Latent Space

Variational Autoencoder: Network Perspective



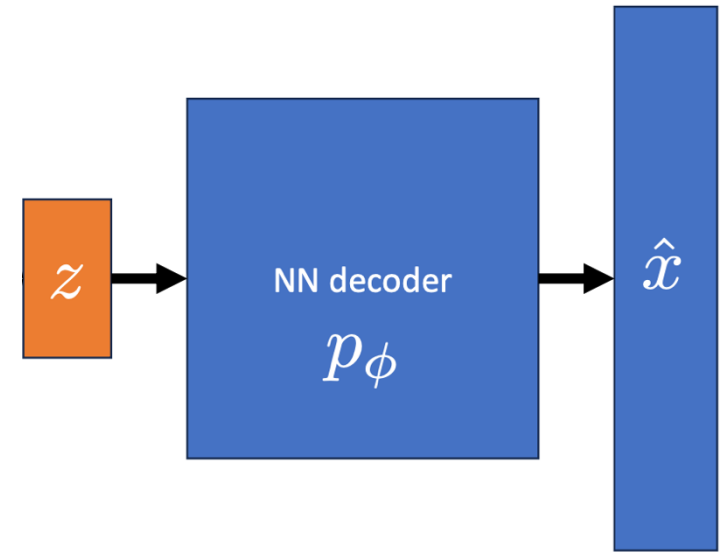
Variational Autoencoder: Network Perspective



- Encoder learns a mean vector and a (diagonal) covariance matrix for each input
- These are used to *sample* a latent representation e.g.,

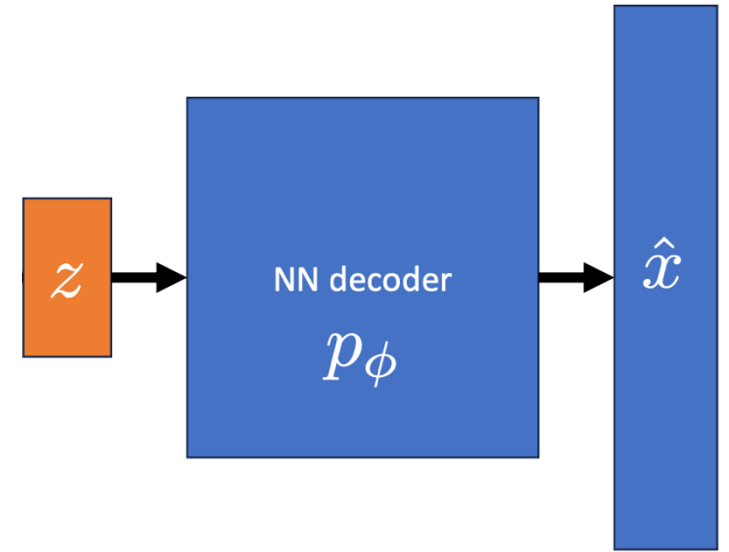
$$\mathbf{z}^{(i)} \mid \mathbf{x}^{(i)} \sim \mathcal{N} \left(\mu_{\theta}(\mathbf{x}^{(i)}), \sigma_{\theta}^2(\mathbf{x}^{(i)}) \right)$$

Variational Autoencoder: Network Perspective



- Decoder tries to minimize the reconstruction error *in expectation* between $\mathbf{x}^{(i)}$ and a sample from another learned (conditional) distribution e.g.,
$$\hat{\mathbf{x}}^{(i)} \mid \mathbf{z}^{(i)} \sim \mathcal{N} \left(\mu_\phi(\mathbf{z}^{(i)}), \sigma_\phi^2(\mathbf{z}^{(i)}) \right)$$

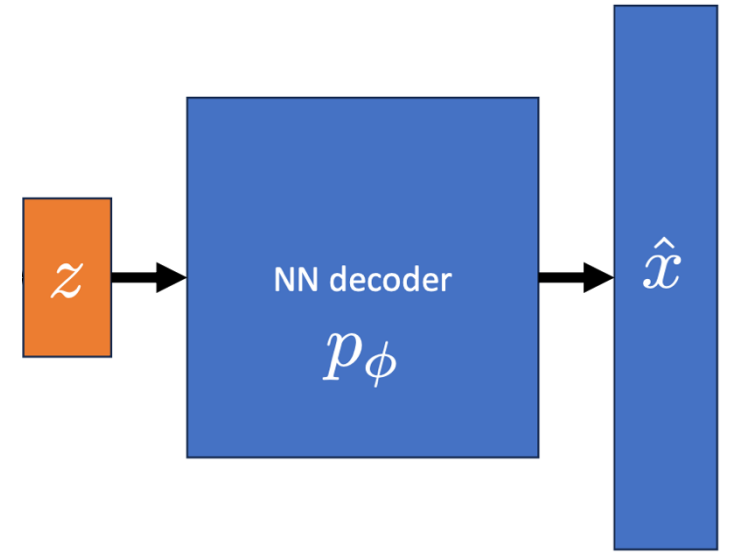
Variational Autoencoder: Network Perspective



- Decoder tries to maximize the likelihood of the true $\mathbf{x}^{(i)}$ under another learned (conditional) distribution e.g.,

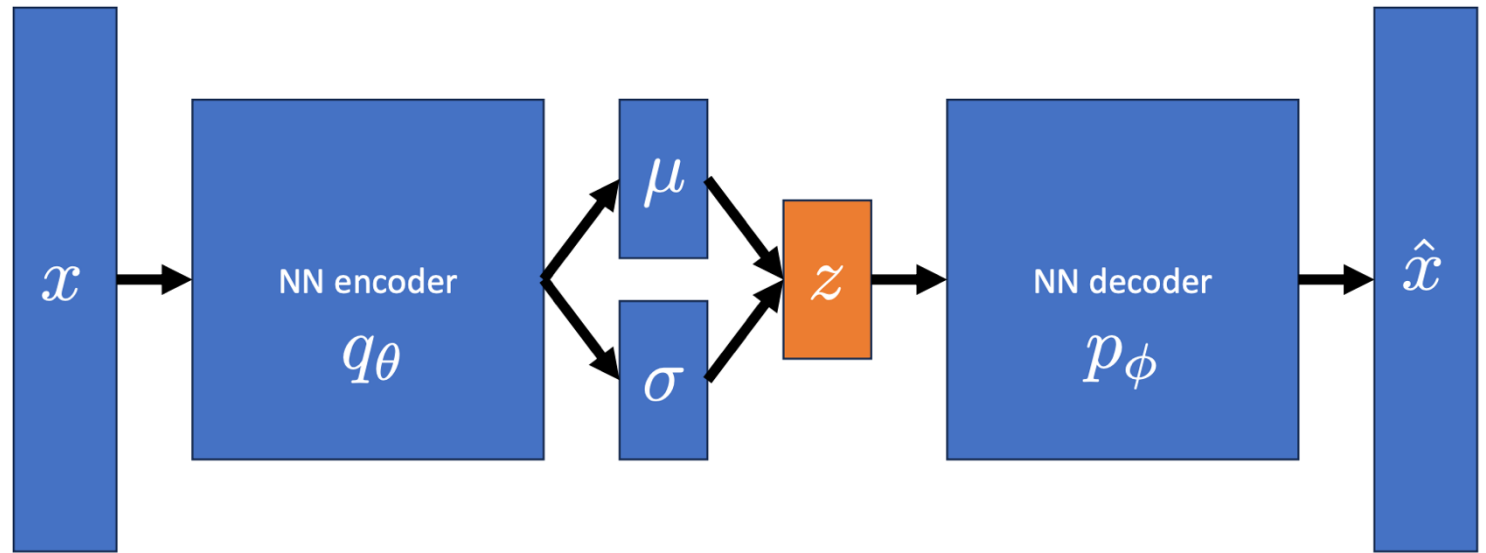
$$\hat{\mathbf{x}}^{(i)} \mid \mathbf{z}^{(i)} \sim \mathcal{N} \left(\mu_{\phi}(\mathbf{z}^{(i)}), \sigma_{\phi}^2(\mathbf{z}^{(i)}) \right)$$

Variational Autoencoder: Network Perspective



- Decoder tries to minimize the negative log-likelihood of the true $\mathbf{x}^{(i)}$ under another learned (conditional) distribution e.g.,
$$\hat{\mathbf{x}}^{(i)} \mid \mathbf{z}^{(i)} \sim \mathcal{N}(\mu_{\phi}(\mathbf{z}^{(i)}), \sigma_{\phi}^2(\mathbf{z}^{(i)}))$$

Variational Autoencoder: Network Perspective



- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that encourages a dense latent space

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

$$\ell_i(\theta, \phi) = -\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\phi(\mathbf{x}^{(i)}|\mathbf{z})] + KL(q_\theta(\mathbf{z}|\mathbf{x}^{(i)}) \parallel p(\mathbf{z}))$$

KL Divergence

- For two distributions $q(x)$ and $p(x)$ over $x \in \mathcal{X}$, the **Kullback-Leibler (KL) divergence** is

$$KL(q||p) = \mathbb{E}_q \left[\log \frac{q(x)}{p(x)} \right] = \sum_{x \in \mathcal{X}} q(x) \log \frac{q(x)}{p(x)}$$

KL Divergence

- For two distributions $q(x)$ and $p(x)$ over $x \in \mathcal{X}$, the **Kullback-Leibler (KL) divergence** is

$$KL(q||p) = \mathbb{E}_q \left[\log \frac{q(x)}{p(x)} \right] = \int_{x \in \mathcal{X}} q(x) \log \frac{q(x)}{p(x)} dx$$

- The KL divergence
 1. measures the **proximity** of two distributions q and p
 2. is minimized when $q(x) = p(x)$ for all $x \in \mathcal{X}$
 3. is **not** symmetric: $KL(q || p) \neq KL(p || q)$

KL Divergence: Example

- Keeping all else constant, consider the effect of differences between p and q for certain x' on $KL(q || p)$

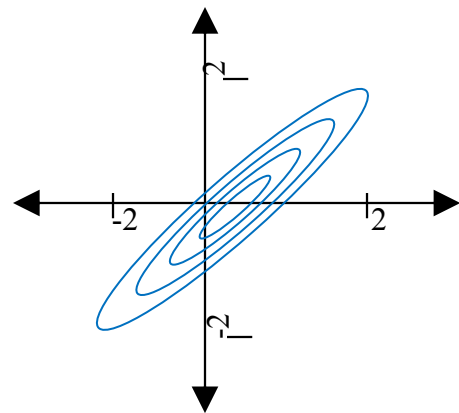
x'	$q(x')$	$p(x')$	$q(x') \log \left(\frac{q(x')}{p(x')} \right)$	effect on $KL(q p)$
1	0.9	0.9	0	no increase
2	0.9	0.1	1.97	big increase
3	0.1	0.9	-0.21	little decrease
4	0.1	0.1	0	little decrease

- KL divergence wants good approximations for values with high probability under q
- KL divergence does not really care about values with low probability under q

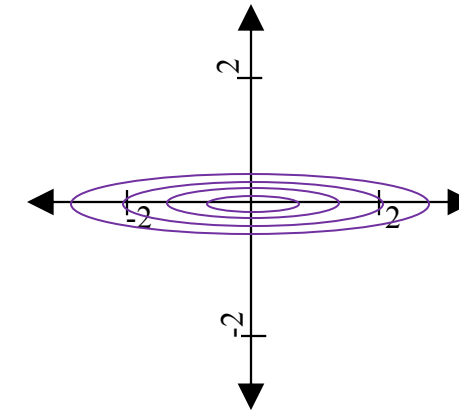
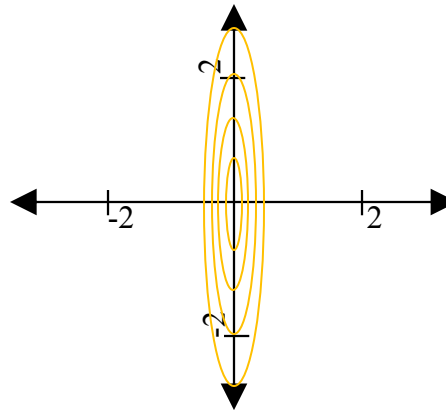
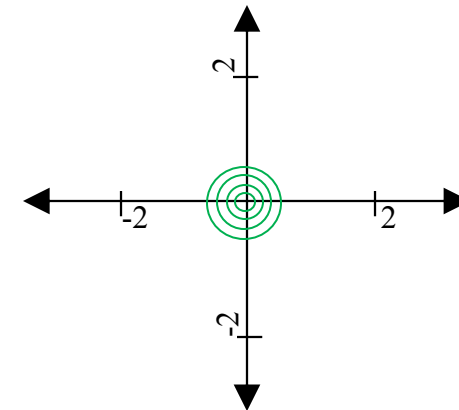
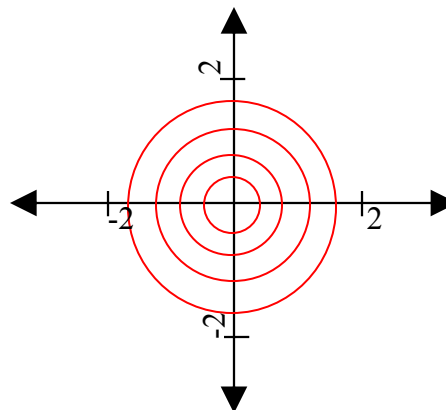
KL Divergence: In-class Exercise

- Which q minimizes $KL(q || p)$ for the given p ?

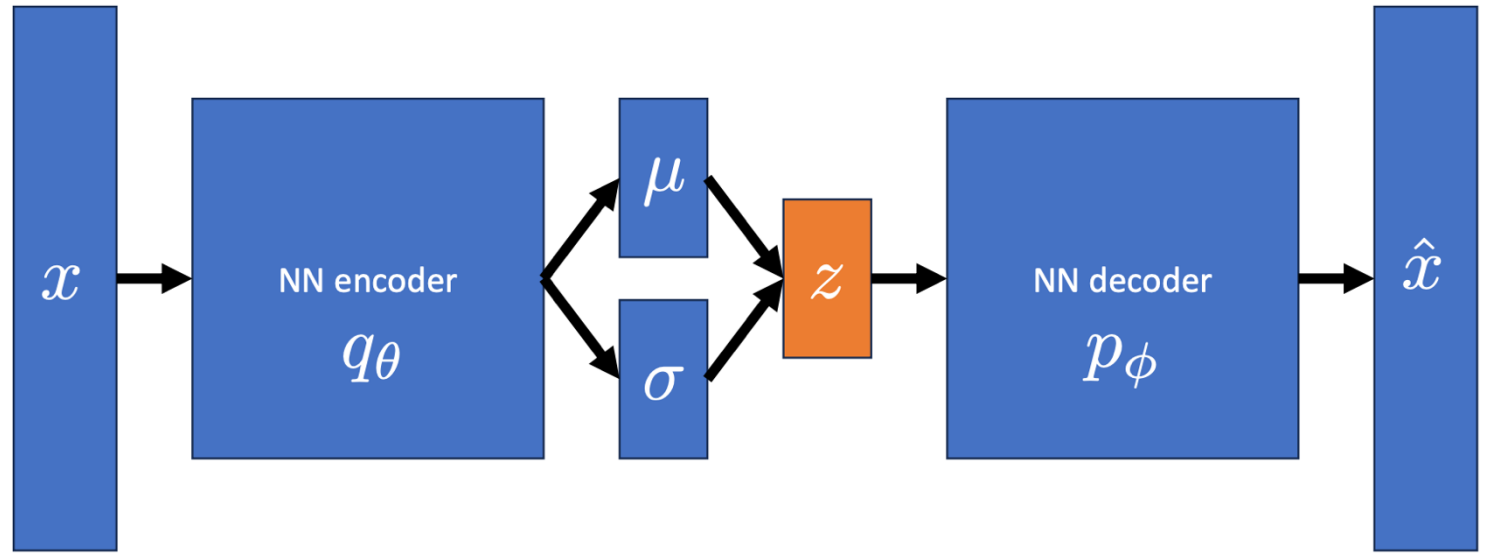
$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu} = [0, 0]^T, \boldsymbol{\Sigma})$$



$$q(x_1, x_2) = \mathcal{N}_1(x_1 | \mu_1, \sigma_1^2) \mathcal{N}_2(x_2 | \mu_2, \sigma_2^2)$$



Variational Autoencoder: Network Perspective

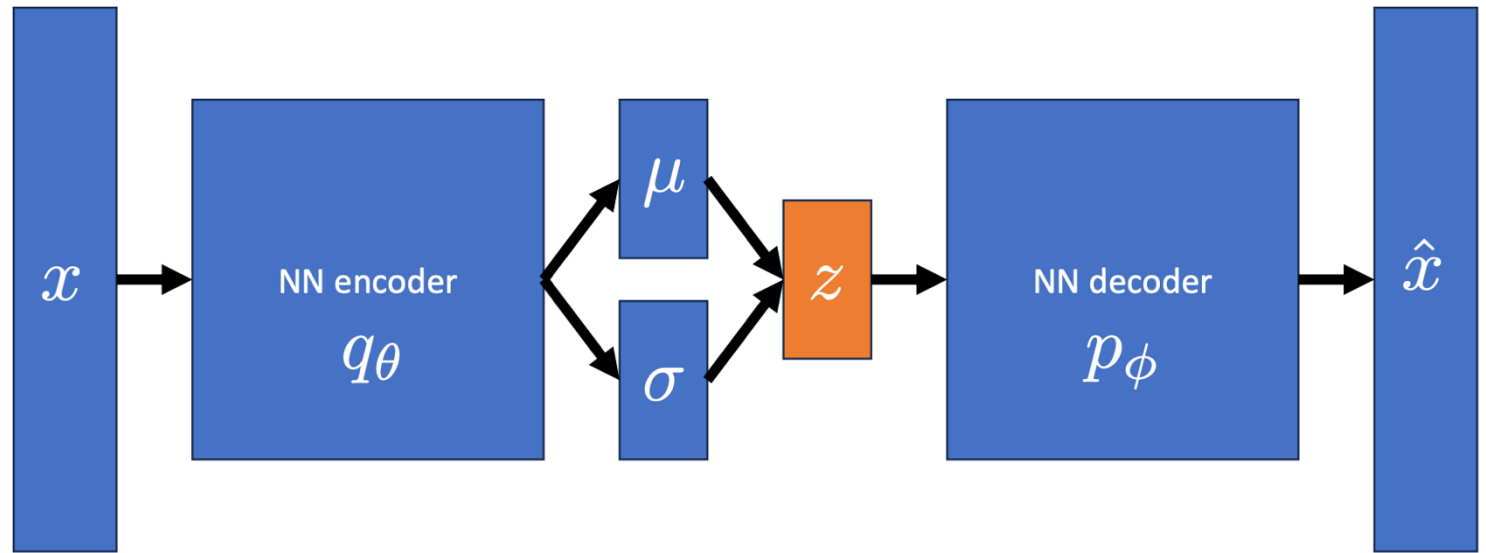


- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that encourages a dense latent space

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

$$\ell_i(\theta, \phi) = -\mathbb{E}_{q_{\theta}(z|x^{(i)})} [\log p_{\phi}(x^{(i)}|z)] + KL(q_{\theta}(z|x^{(i)}) \parallel p(z))$$

So what should we set p to?

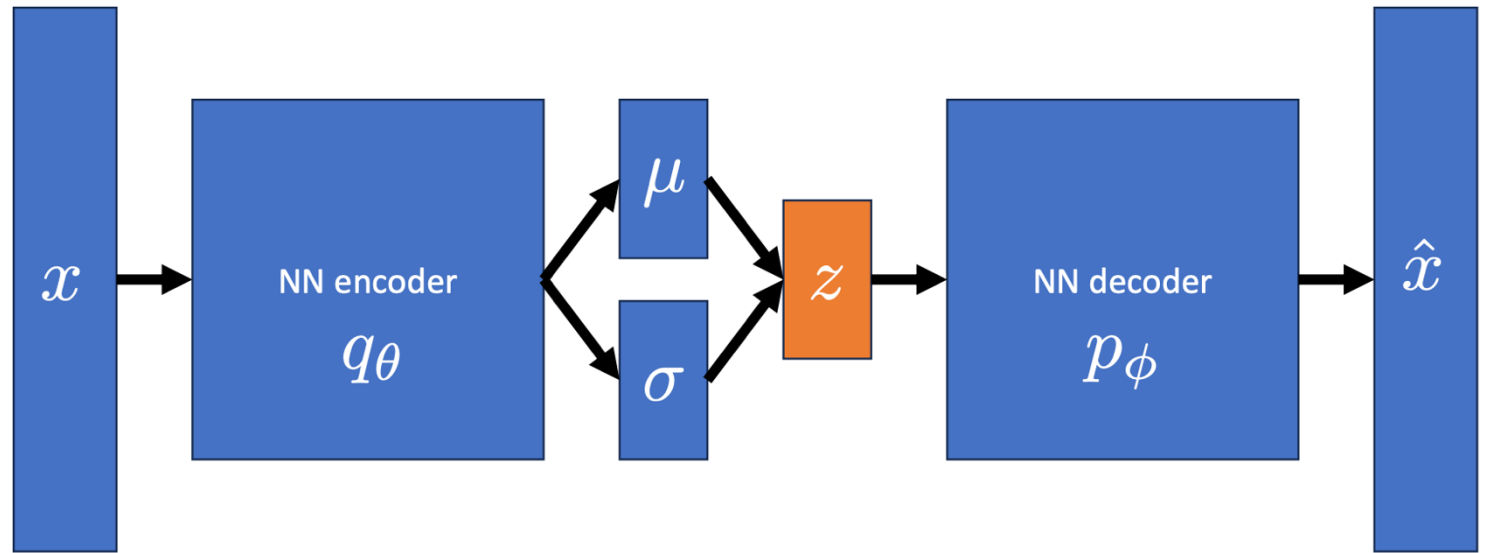


- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that encourages a **dense latent space**

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

$$\ell_i(\theta, \phi) = -\mathbb{E}_{q_\theta(z|x^{(i)})} [\log p_\phi(x^{(i)}|z)] + KL(q_\theta(z|x^{(i)}) \parallel p(z))$$

Variational Autoencoder: Network Perspective



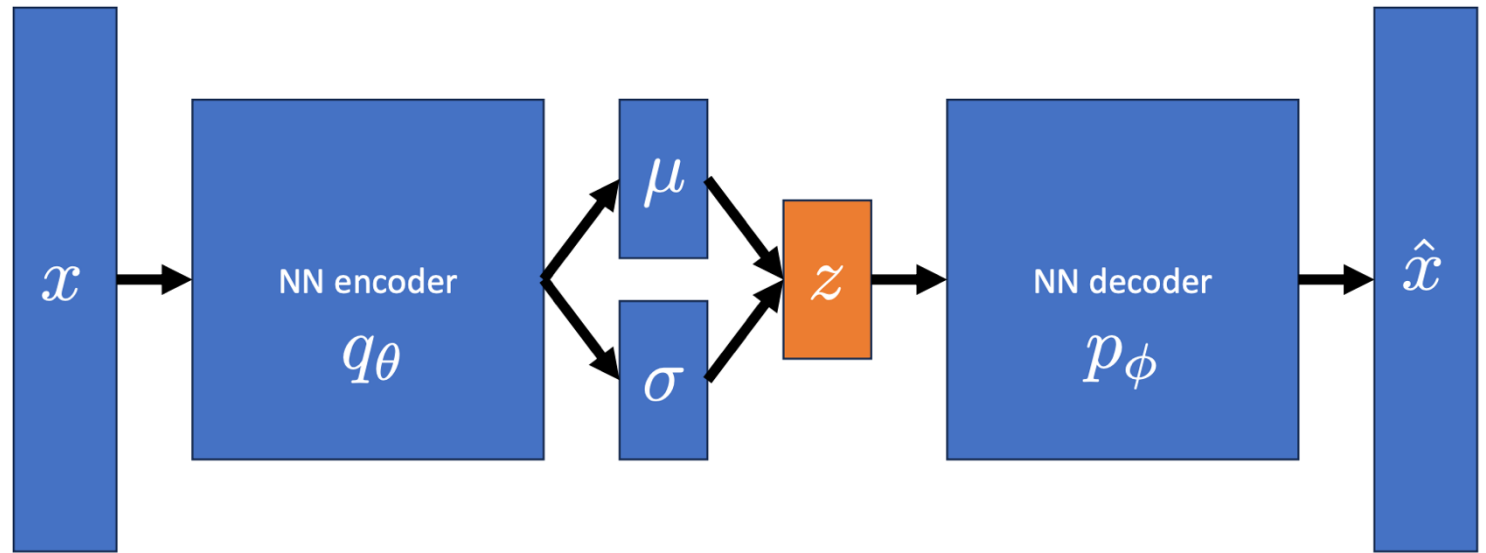
- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that encourages a dense latent space

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

$$\ell_i(\theta, \phi) \approx - \left(\frac{1}{S} \sum_{s=1}^S \log p_{\phi}(x^{(i)} | z_s) \right) + KL(q_{\theta}(z | x^{(i)}) \| p(z))$$

for samples $\mathbf{z}_1, \dots, \mathbf{z}_S \sim q_{\theta}(\mathbf{z} | x^{(i)})$

Can we
backpropagate
through q_θ
given that
samples of \mathbf{z}
are stochastic?



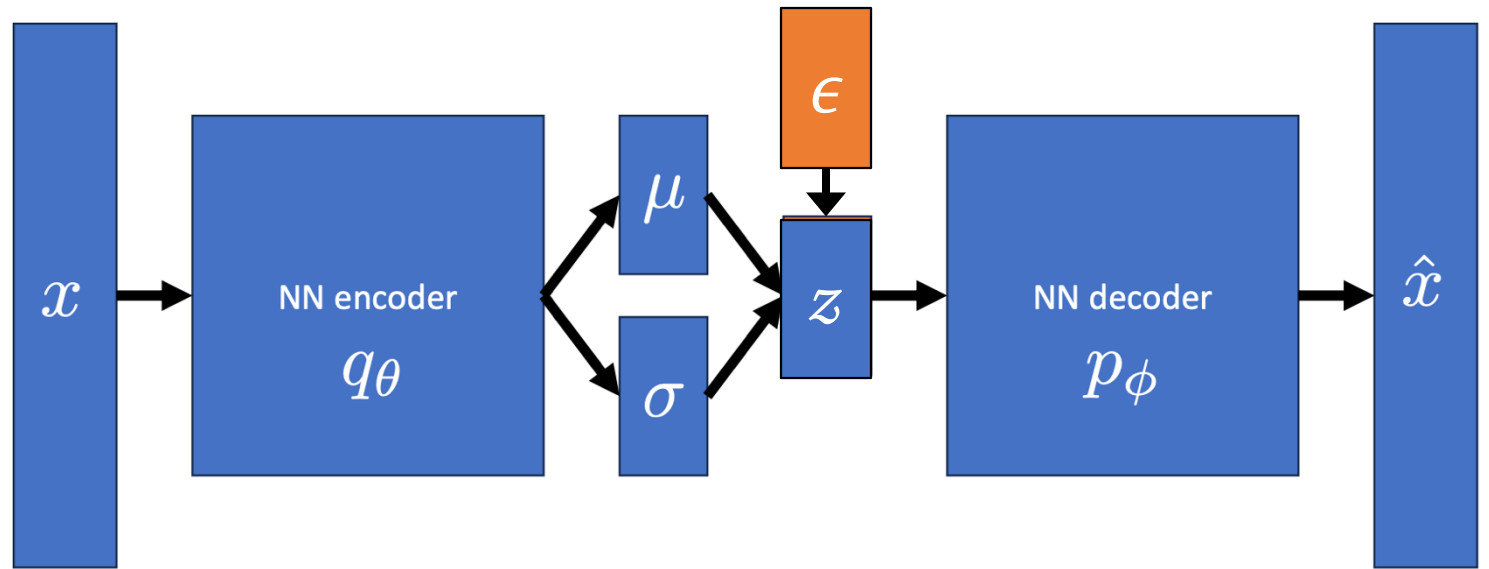
- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that encourages a dense latent space

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

$$\ell_i(\theta, \phi) \approx - \left(\frac{1}{S} \sum_{s=1}^S \log p_\phi(\mathbf{x}^{(i)} | \mathbf{z}_s) \right) + KL(q_\theta(\mathbf{z} | \mathbf{x}^{(i)}) \parallel p(\mathbf{z}))$$

for samples $\mathbf{z}_1, \dots, \mathbf{z}_S \sim q_\theta(\mathbf{z} | \mathbf{x}^{(i)})$

Reparameterization Trick



- Objective: minimize the negative log-likelihood of the dataset plus a *regularization term* that encourages a dense latent space

$$J(\theta, \phi) = \sum_{i=1}^N \ell_i(\theta, \phi)$$

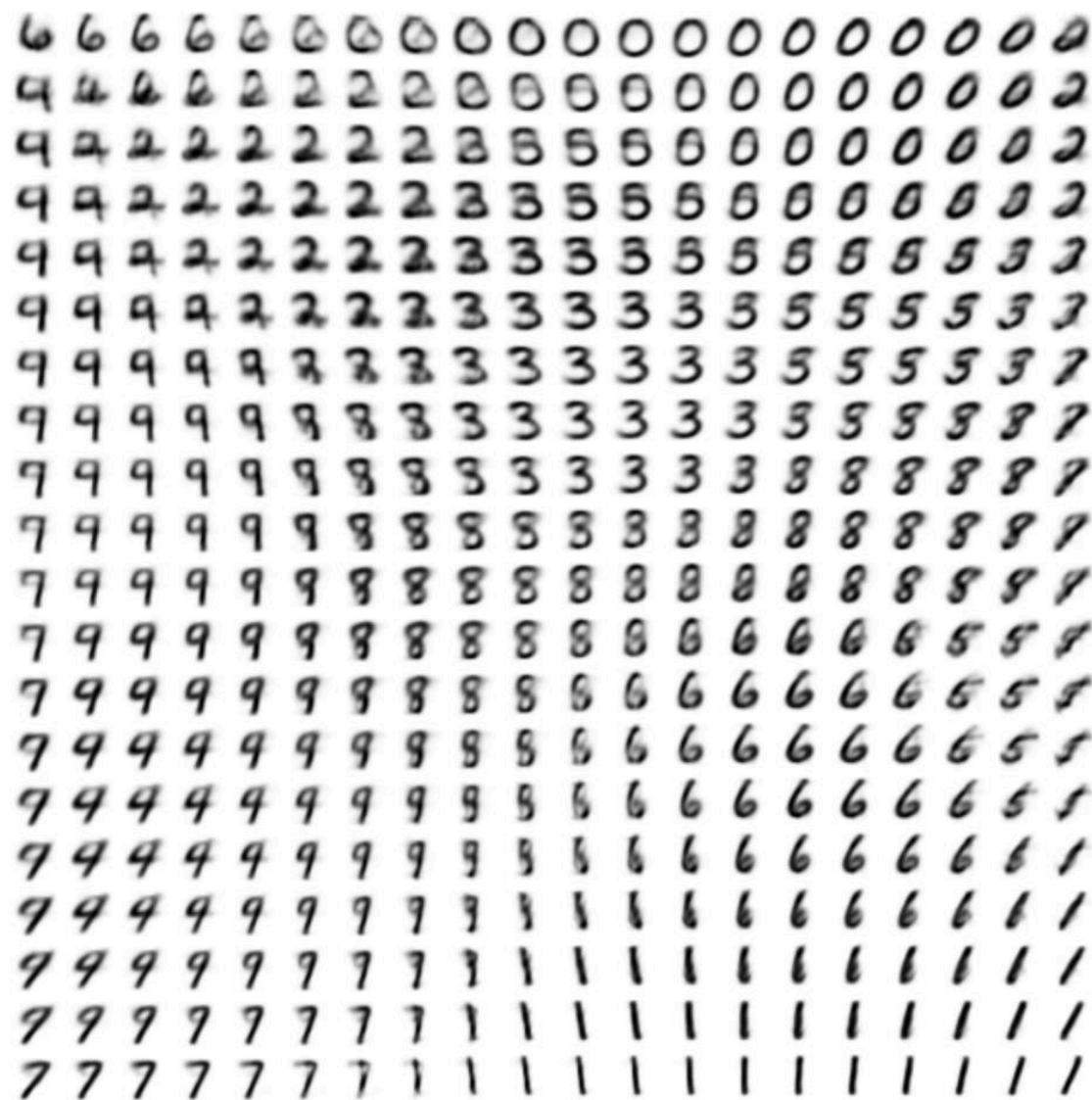
$$\ell_i(\theta, \phi) \approx - \left(\frac{1}{S} \sum_{s=1}^S \log p_\phi(x^{(i)} | z_s) \right) + KL(q_\theta(z | x^{(i)}) \| p(z))$$

for $z_s = \mu_\theta(x^{(i)}) + \sigma_\theta(x^{(i)}) \odot \epsilon_s$ where $\epsilon_s \sim N(\mathbf{0}, I)$

$$\begin{aligned}
\ell_i(\boldsymbol{\theta}, \boldsymbol{\phi}) &= -\mathbb{E}_{q_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)})}[\log p_{\boldsymbol{\phi}}(\mathbf{x}^{(i)}|\mathbf{z})] + KL(q_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)}) \parallel p(\mathbf{z})) \\
&\approx -\left(\frac{1}{S} \sum_{s=1}^S \log p_{\boldsymbol{\phi}}(\mathbf{x}^{(i)}|\mathbf{z}_s(\boldsymbol{\theta}))\right) + KL(q_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)}) \parallel p(\mathbf{z})) \\
&= -\left(\frac{1}{S} \sum_{s=1}^S \log \mathcal{N}(\mathbf{x}^{(i)}; \boldsymbol{\mu}_{\boldsymbol{\phi}}(\mathbf{z}_s(\boldsymbol{\theta})), \boldsymbol{\sigma}_{\boldsymbol{\phi}}^2(\mathbf{z}_s(\boldsymbol{\theta})))\right) \\
&\quad + KL(\mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), \boldsymbol{\sigma}_{\boldsymbol{\theta}}^2(\mathbf{x}^{(i)})) \parallel \mathcal{N}(\mathbf{0}, I)) \\
&= -\left(\frac{1}{S} \sum_{s=1}^S \log \mathcal{N}(\mathbf{x}^{(i)}; \boldsymbol{\mu}_{\boldsymbol{\phi}}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + \boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \odot \boldsymbol{\epsilon}_s), \boldsymbol{\sigma}_{\boldsymbol{\phi}}^2(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + \boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \odot \boldsymbol{\epsilon}_s))\right) \\
&\quad + KL(\mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), \boldsymbol{\sigma}_{\boldsymbol{\theta}}^2(\mathbf{x}^{(i)})) \parallel \mathcal{N}(\mathbf{0}, I))
\end{aligned}$$

Variational Autoencoder: Objective Function

Variational Autoencoder: Latent Space Visualization

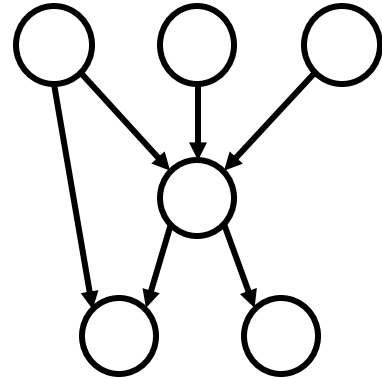


Variational Autoencoder: Generated Samples...

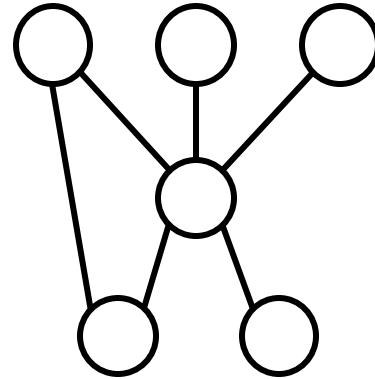


Three Types of Graphical Models

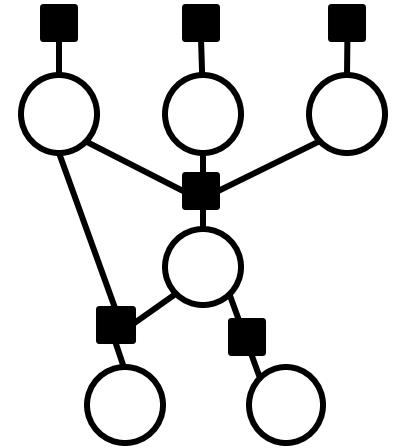
Directed Graphical Model



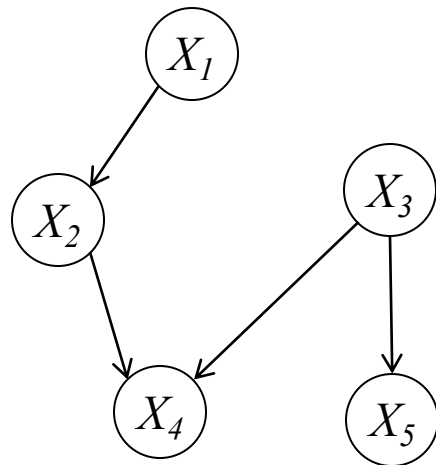
Undirected Graphical Model



Factor Graph

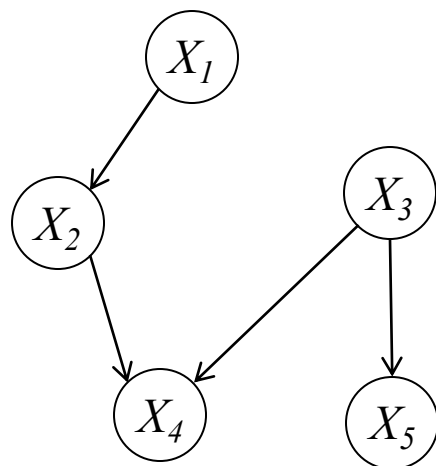


Directed Graphical Models a.k.a. Bayesian Networks



$$P(X_1, X_2, X_3, X_4, X_5) = P(X_1) \\ * P(X_2|X_1) \\ * P(X_3) \\ * P(X_4|X_2, X_3) \\ * P(X_5|X_3)$$

Directed Graphical Models a.k.a. Bayesian Networks



$$P(X_1, \dots, X_D) = \prod_{d=1}^D P(X_d | \text{parents}(X_d))$$

A Bayesian Network consists of:

- a graph G (the *qualitative specification*), which can be
 - specified using prior knowledge / domain expertise
 - learned from the training data (model selection)
- conditional probabilities (the *quantitative specification*)
 - these will depend on the relative types of the variables