

# 10-301/601: Introduction to Machine Learning

## Lecture 32 – Generative Models for Vision

Henry Chai

6/16/25

# Front Matter

- Announcements
  - HW8 released 6/13, due 6/17 (tomorrow) at 11:59 PM
  - Final on 6/20 (next Friday) at **8:30 AM** in BH A36 (here!)
    - **We will not use the full 3-hour window**
    - All topics from Lectures 17 to 30 are in-scope
    - **The final is *not* cumulative:** pre-midterm content may be referenced but will not be the *primary* focus of any question
    - You are allowed to bring one letter-size sheet of notes; you may put *whatever* you want on *both sides*

# Image Generation

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

# Image Generation

sea anemone

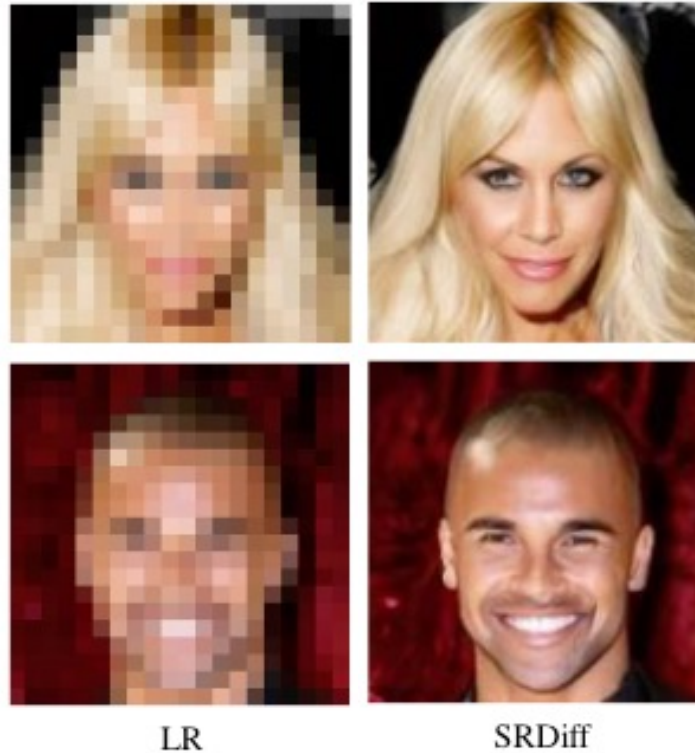
brain coral

slug



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

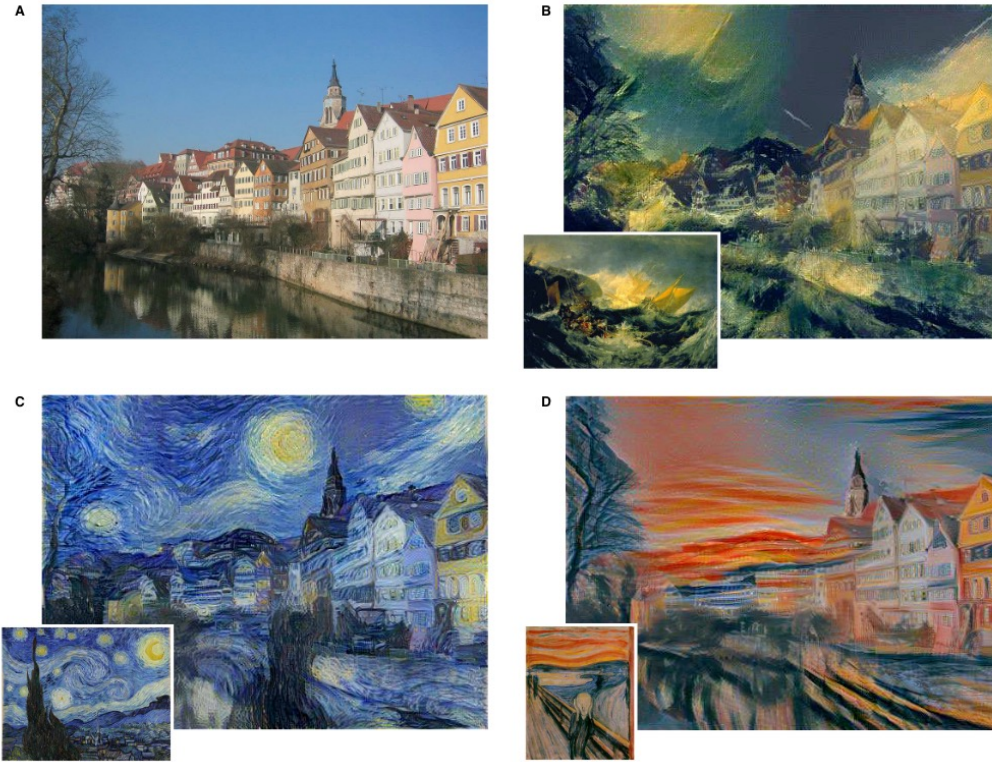
- Class-conditional generation
- **Super resolution**
- Image Editing
- Style transfer
- Text-to-image (TTI) generation

# Image Generation



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

# 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.

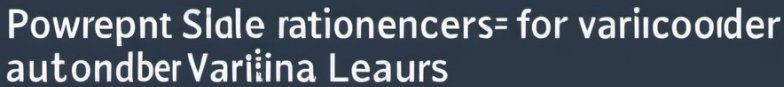


- Class-conditional generation
  - Super resolution
  - Image Editing
  - Style transfer
  - **Text-to-image (TTI) generation**
- Given a text description, sample an image that depicts the prompt



# Slide Generation?

*Prompt:* powerpoint slide explaining



- Class-conditional generation

- Super resolution

- Image Editing

- Style transfer

- **Text-to-image (TTI) generation**

# Slide Generation!

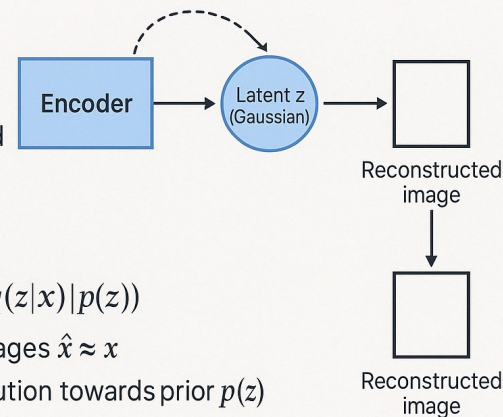
*Prompt:* powerpoint slide explaining variational autoencoders for an intro to ML course, easy to follow, with an explanation of the evidence lower bound

## Variational Autoencoders (VAE)

- Unsupervised generative model that learns a latent space
- Encoder: maps input  $x$  to distribution  $q(z|x)$
- Decoder: maps latent  $z$  to reconstructed  $\hat{x}$  from  $p(x|z)$
- Training objective: maximize Evidence Lower Bound (ELBO)

$$\log p(x) \geq E_{q(z|x)} [\log p(x|z) - \text{KL}q(z|x)|p(z))]$$

- Reconstruction term (likelihood) encourages  $\hat{x} \approx x$
- KL divergence regularizes latent distribution towards prior  $p(z)$



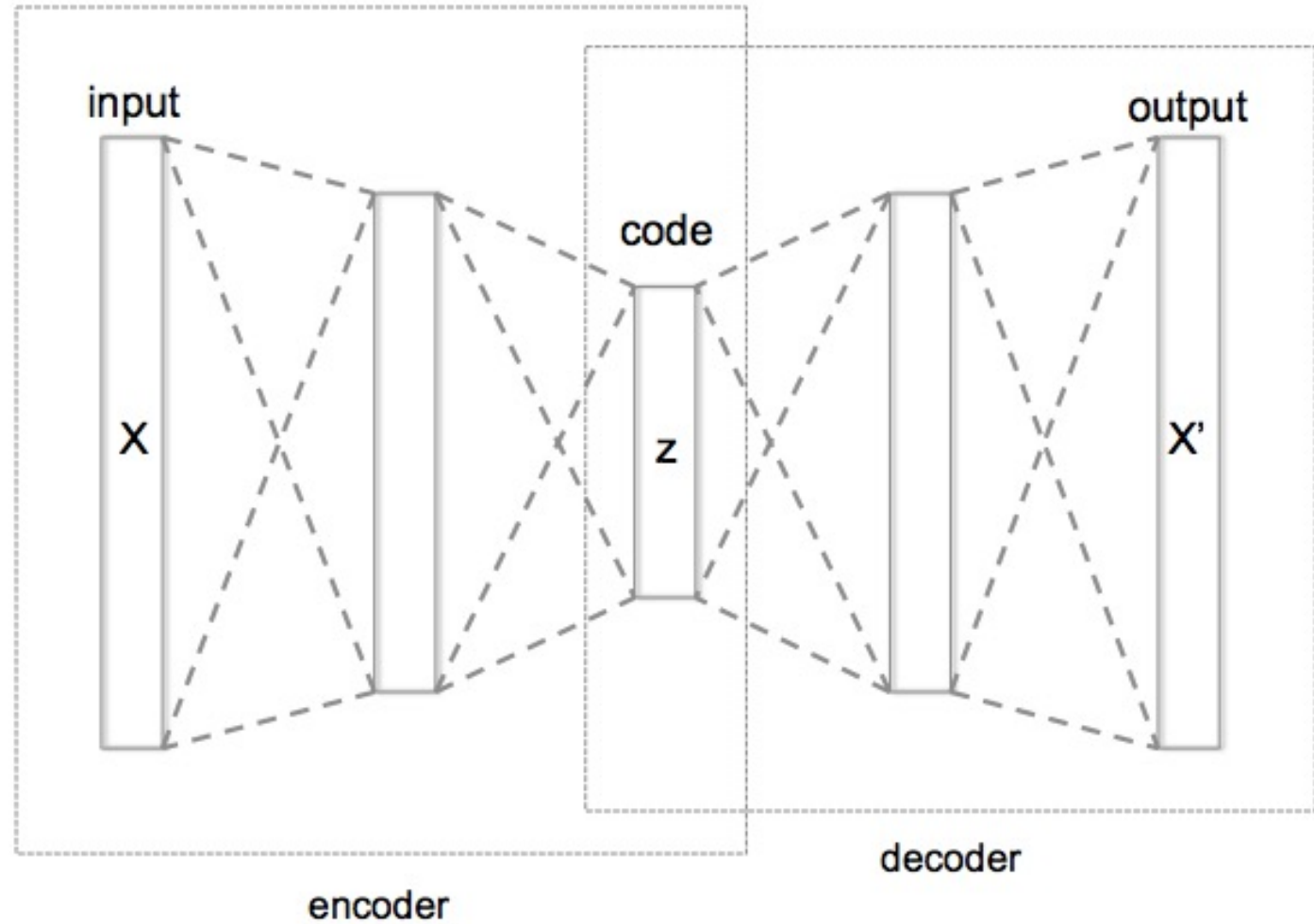
Intro to Machine Learning | Your Name | Date

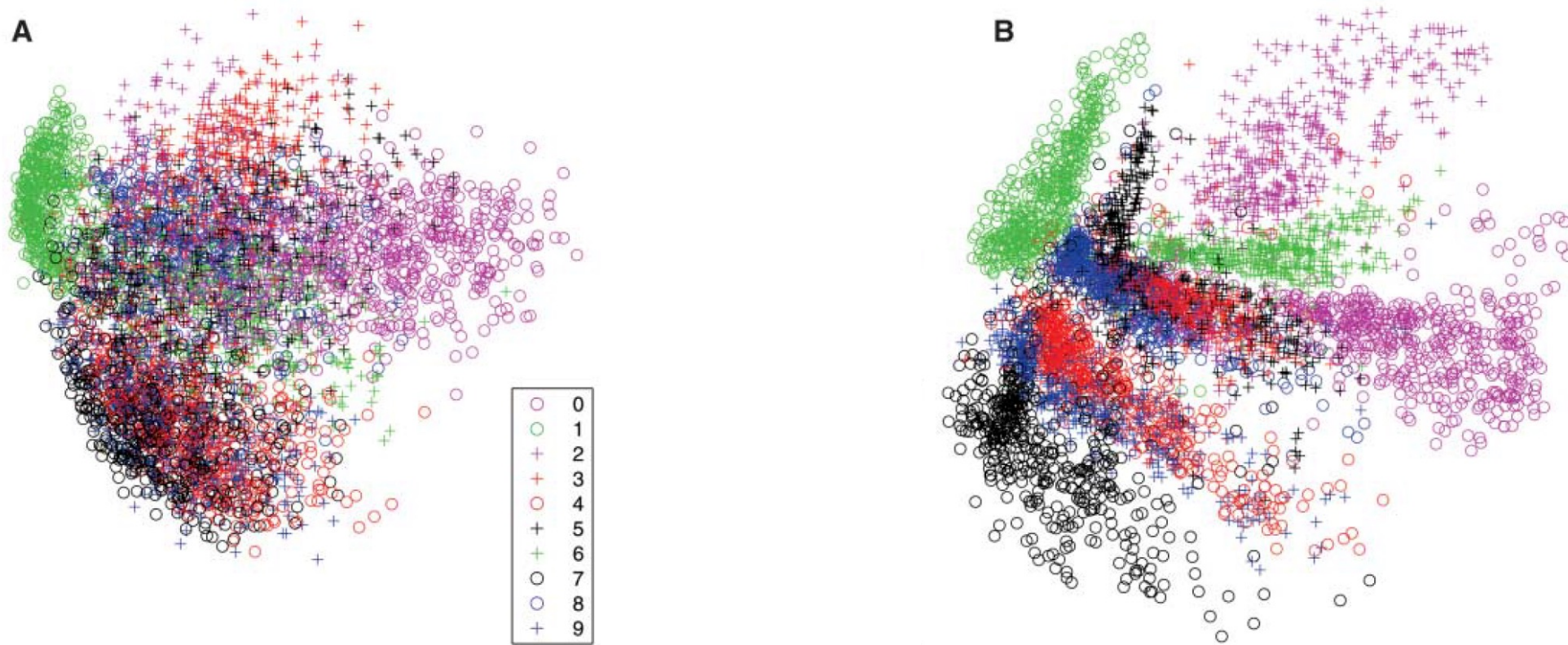
- Class-conditional generation
- Super resolution
- Image Editing
- Style transfer
- **Text-to-image (TTI) generation**

# 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: Deep Autoencoders

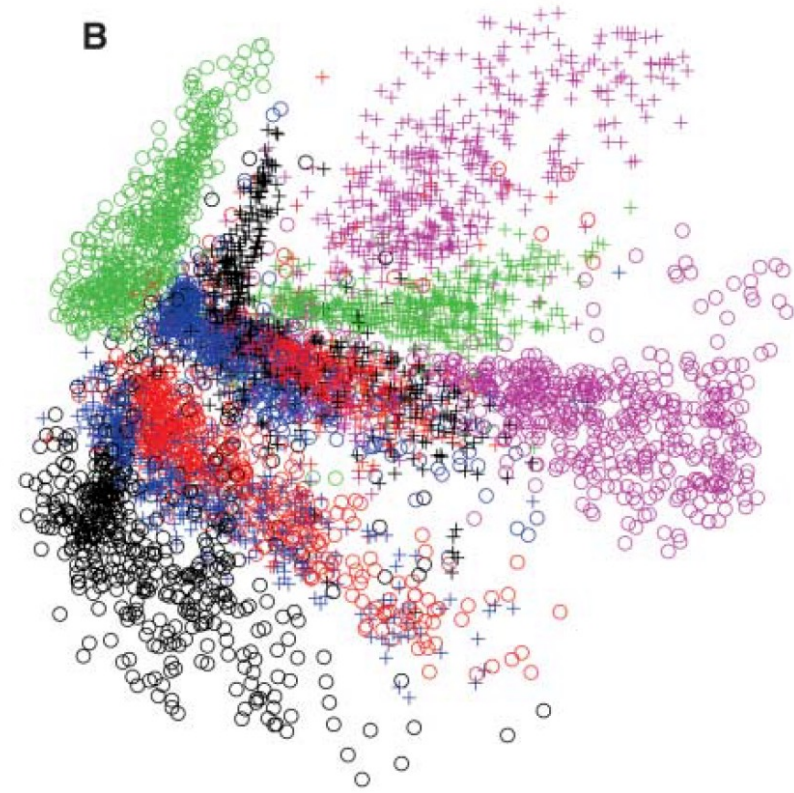




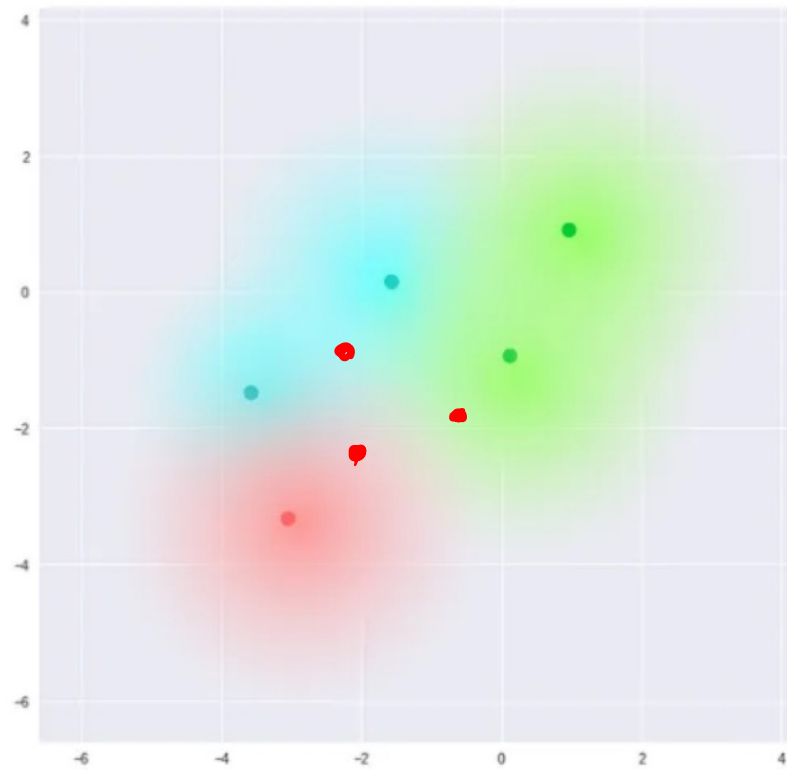
# PCA (A) vs. Autoencoders (B) (Hinton and Salakhutdinov, 2006)



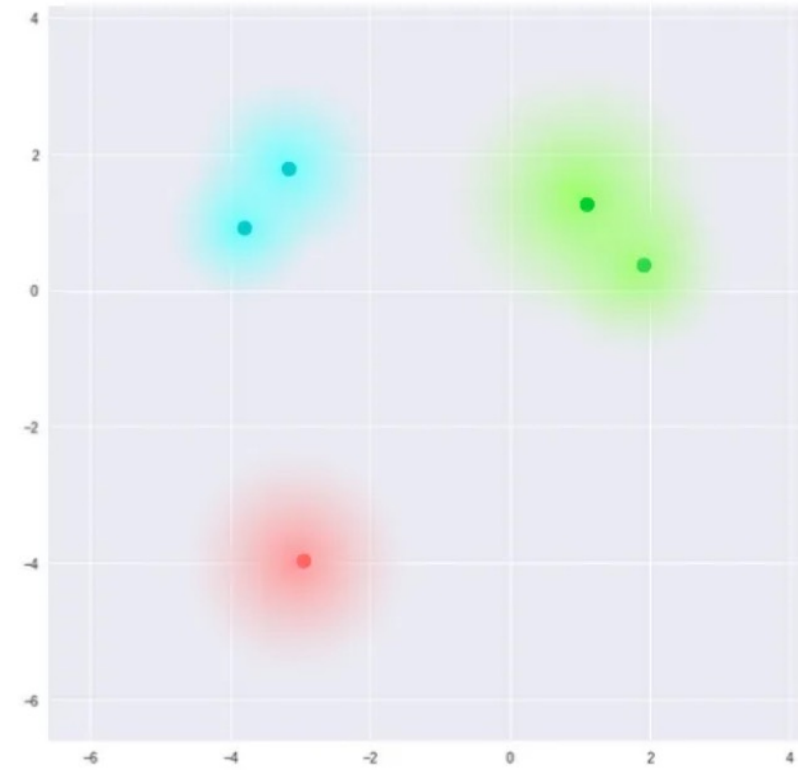
- Issue: latent space is sparse...
  - Sampling from latent space of an autoencoder creates outputs that are effectively identical to images in the training dataset



# Autoencoder Latent Space



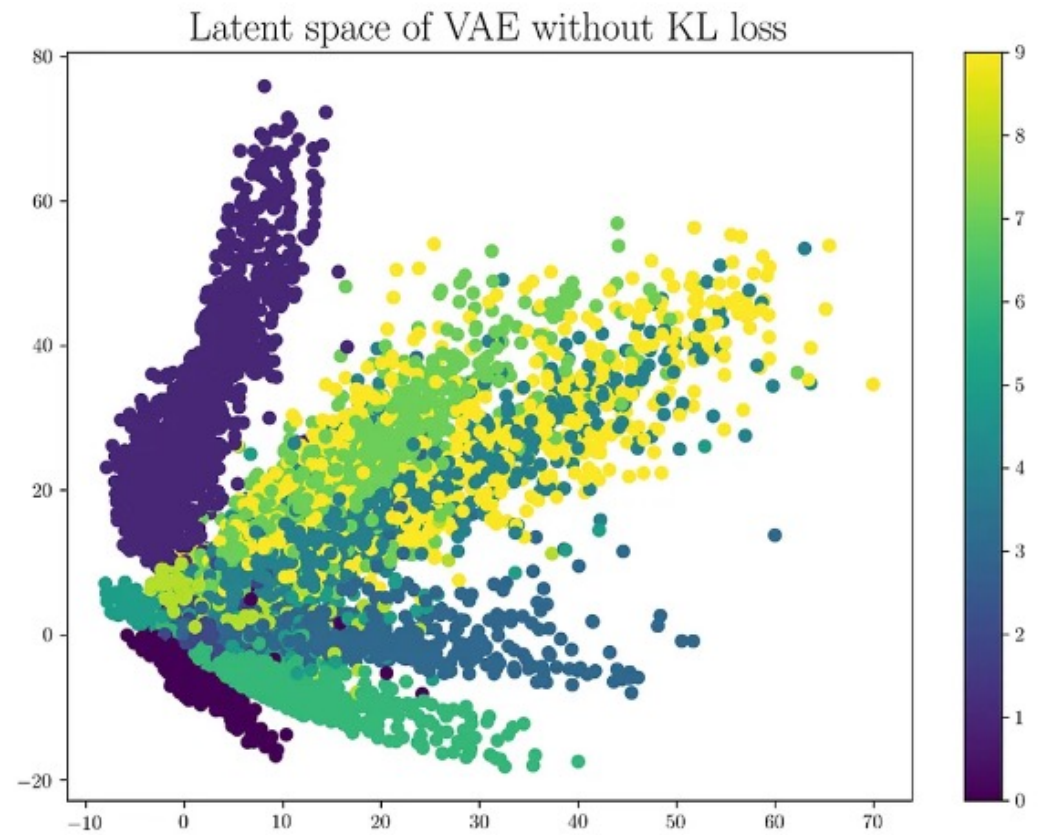
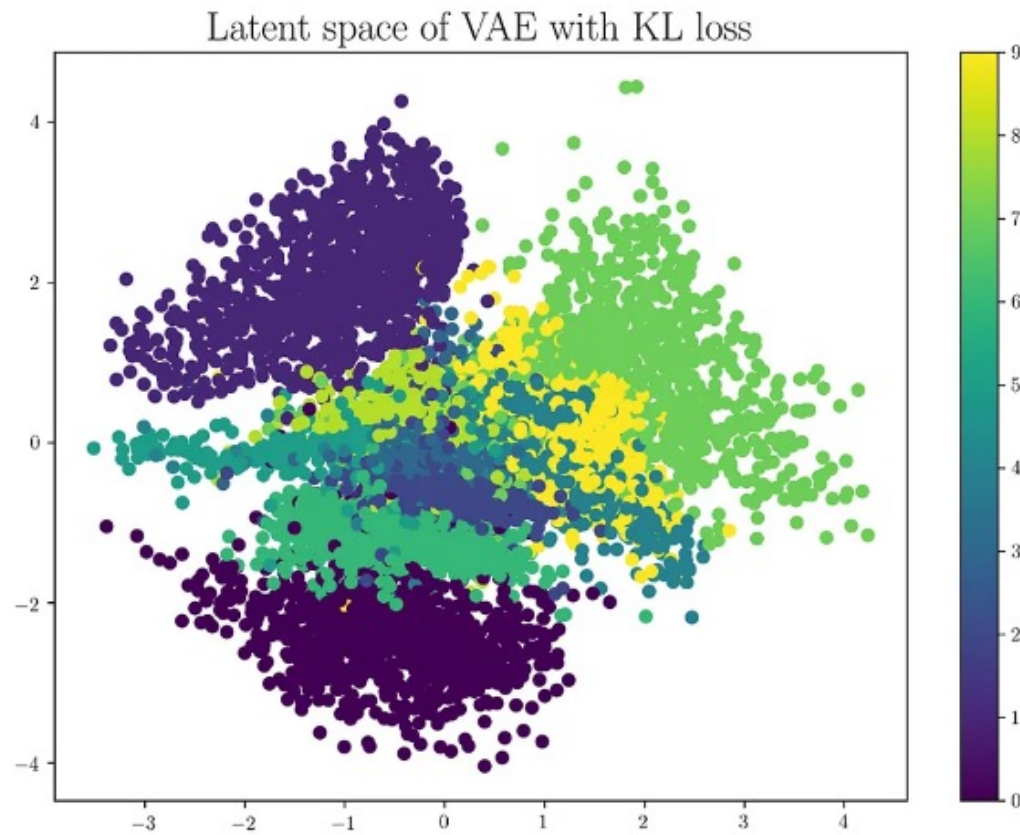
What we require



What we may inadvertently end up with

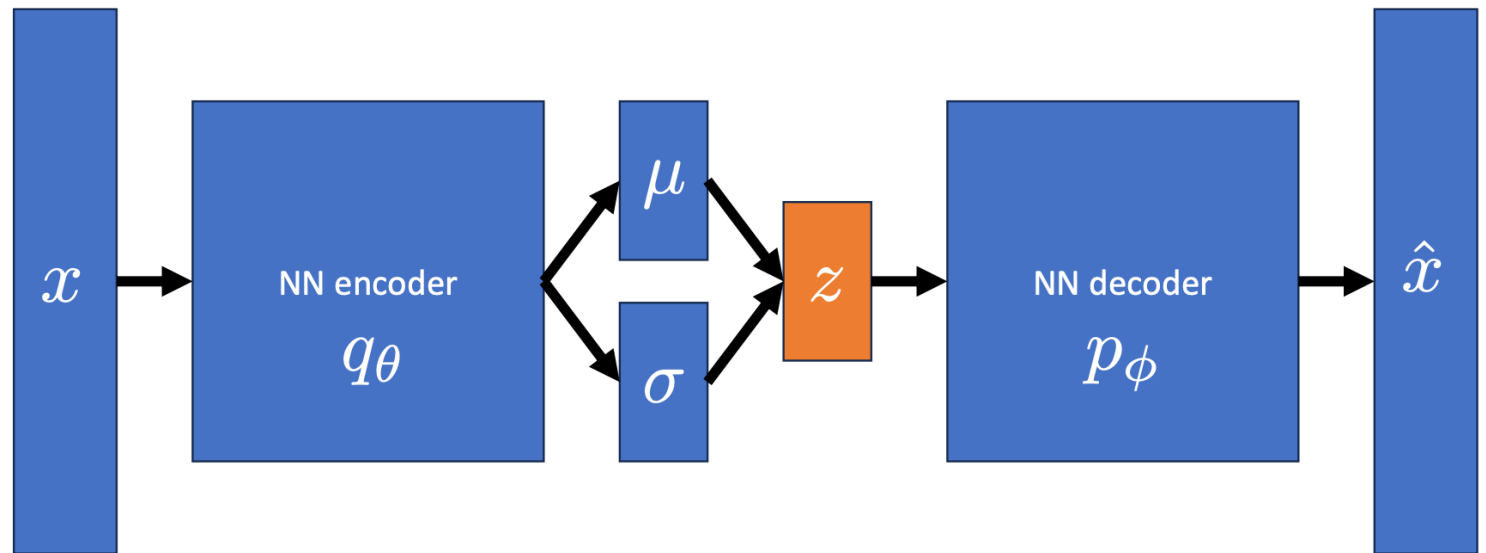
# Autoencoder Latent Space



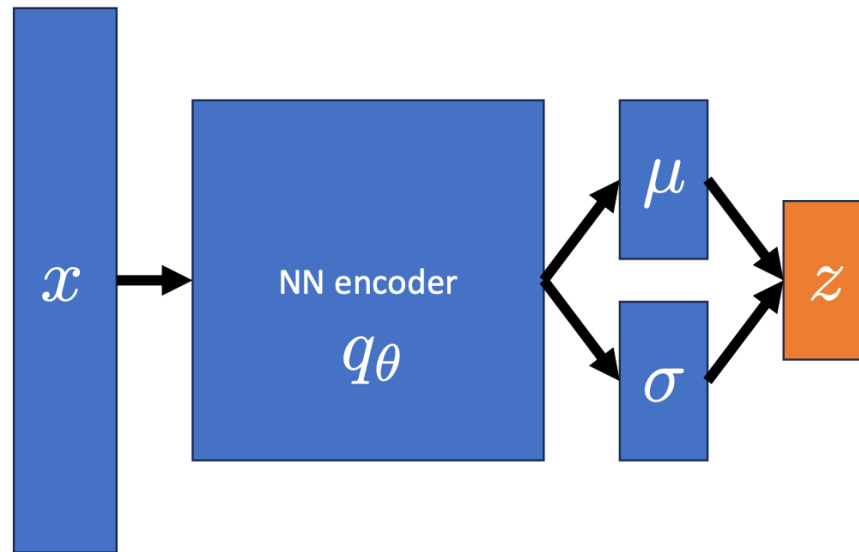


# Variational Autoencoder Latent Space

# Variational Autoencoder: Network Perspective



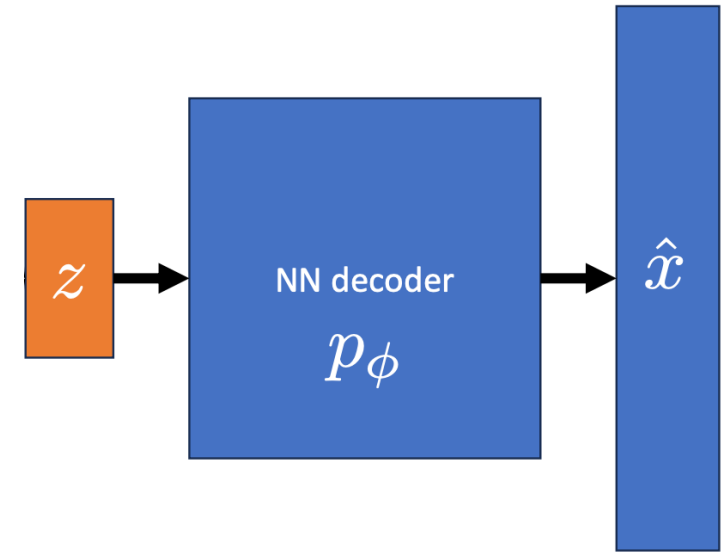
# Variational Autoencoder: Network Perspective



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

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

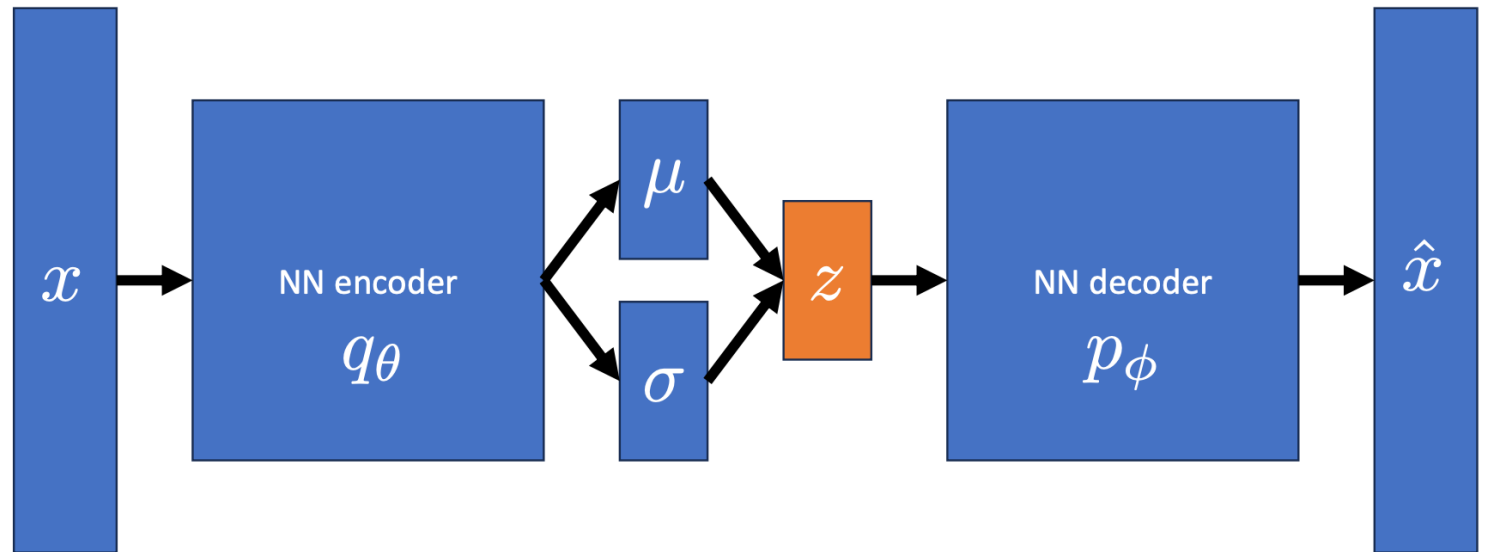
# Variational Autoencoder: Network Perspective



- Decoder tries to minimize the reconstruction error *in expectation* between  $\mathbf{x}^{(i)}$  and a sample from another (conditional) distribution e.g.,

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

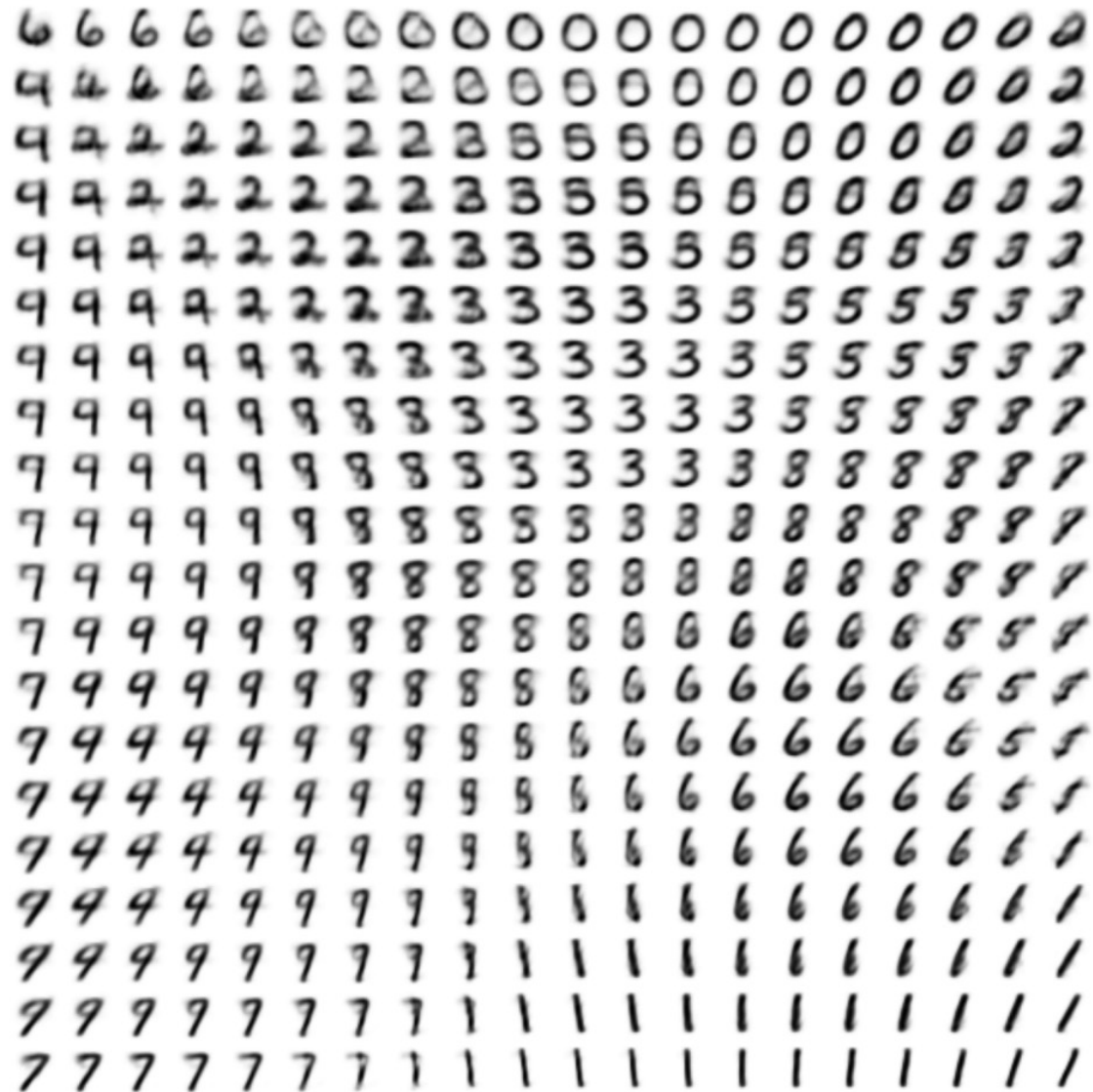
# Variational Autoencoder: Network Perspective



- Objective: minimize the expected reconstruction error plus a *regularizer* that encourages a dense latent space

$$\mathcal{L}(\theta, \phi) = \sum_{i=1}^N \left( -\mathbb{E}_{q_\theta(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\phi(\mathbf{x}^{(i)}|\mathbf{z})] \right) + KL \left( q_\theta(\mathbf{z}|\mathbf{x}^{(i)}) \parallel p(\mathbf{z}) \right)$$

# Variational Autoencoder: Latent Space Visualization



# Variational Autoencoder: Generated Samples





# Variational Autoencoder: Generated Samples?



Can we encode  
the idea that  
samples should  
be  
indistinguishable  
from real  
observations into  
the objective  
function?



Source: <https://arxiv.org/pdf/1312.6114.pdf>

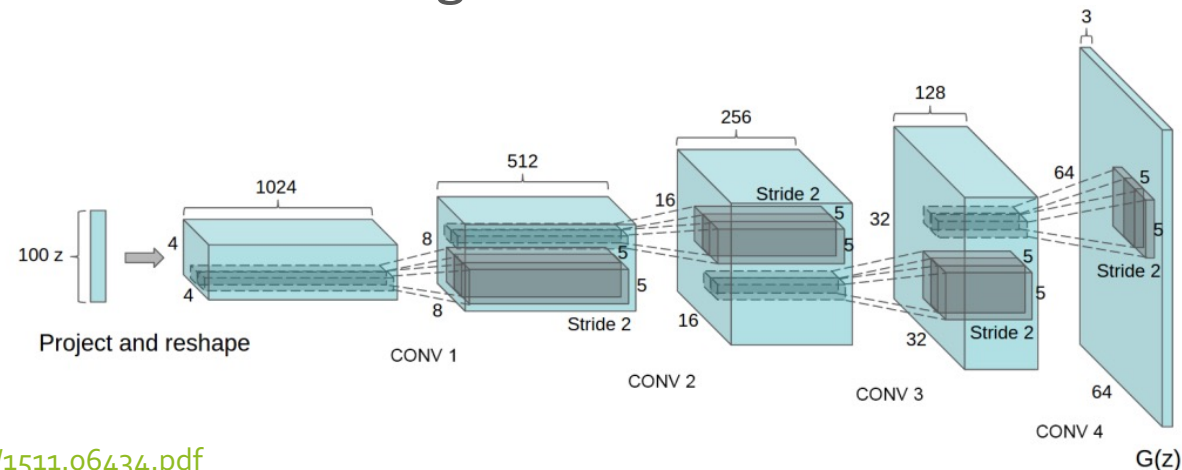
Source: MNIST

# 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
- 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

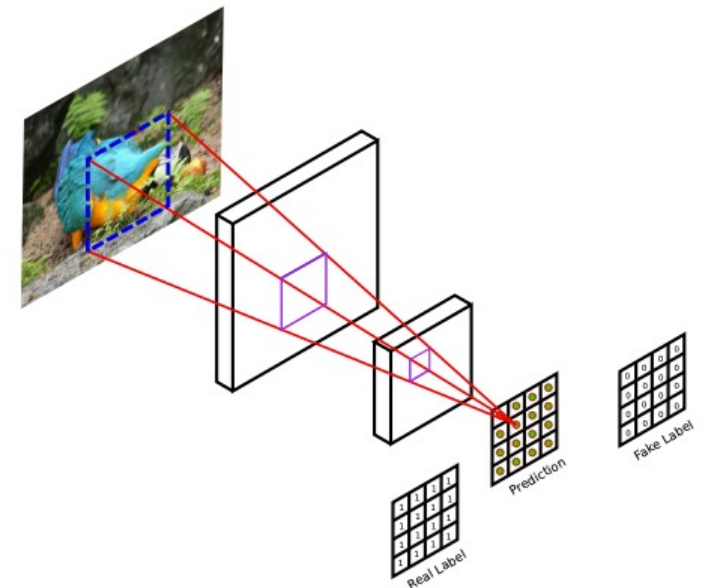
# Generative Adversarial Networks (GANs)

- 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



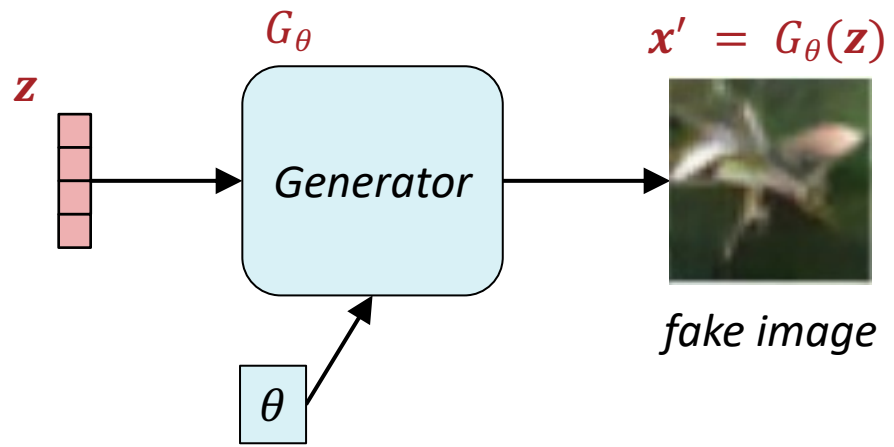
# Generative Adversarial Networks (GANs)

- 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 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 to generator to avoid creating blurry images



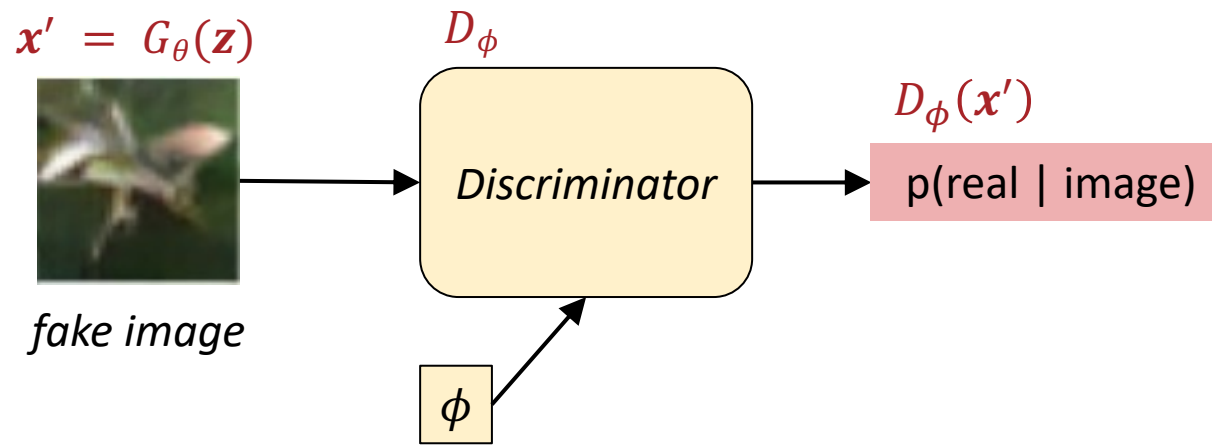
# Generative Adversarial Networks (GANs): Training

- 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

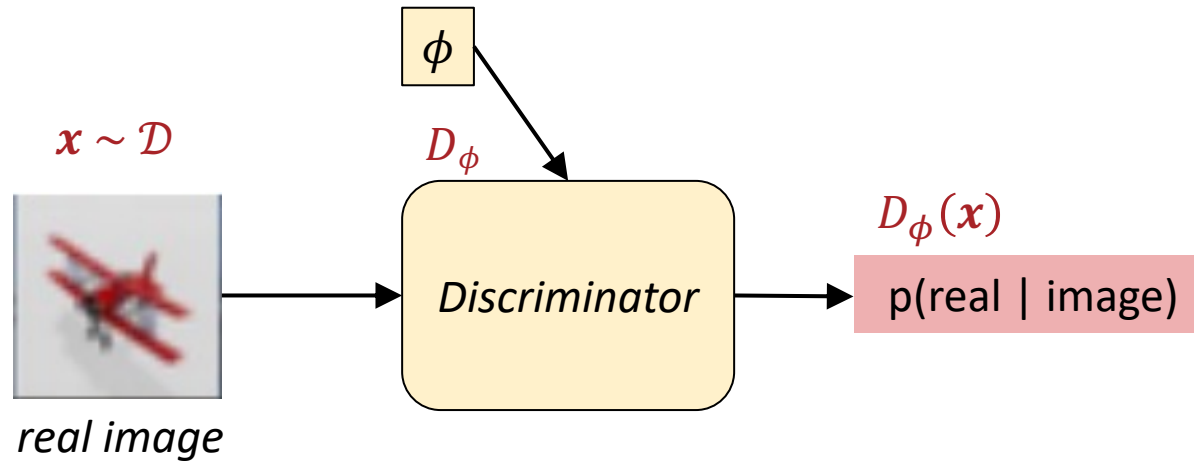


# GANs: Architecture

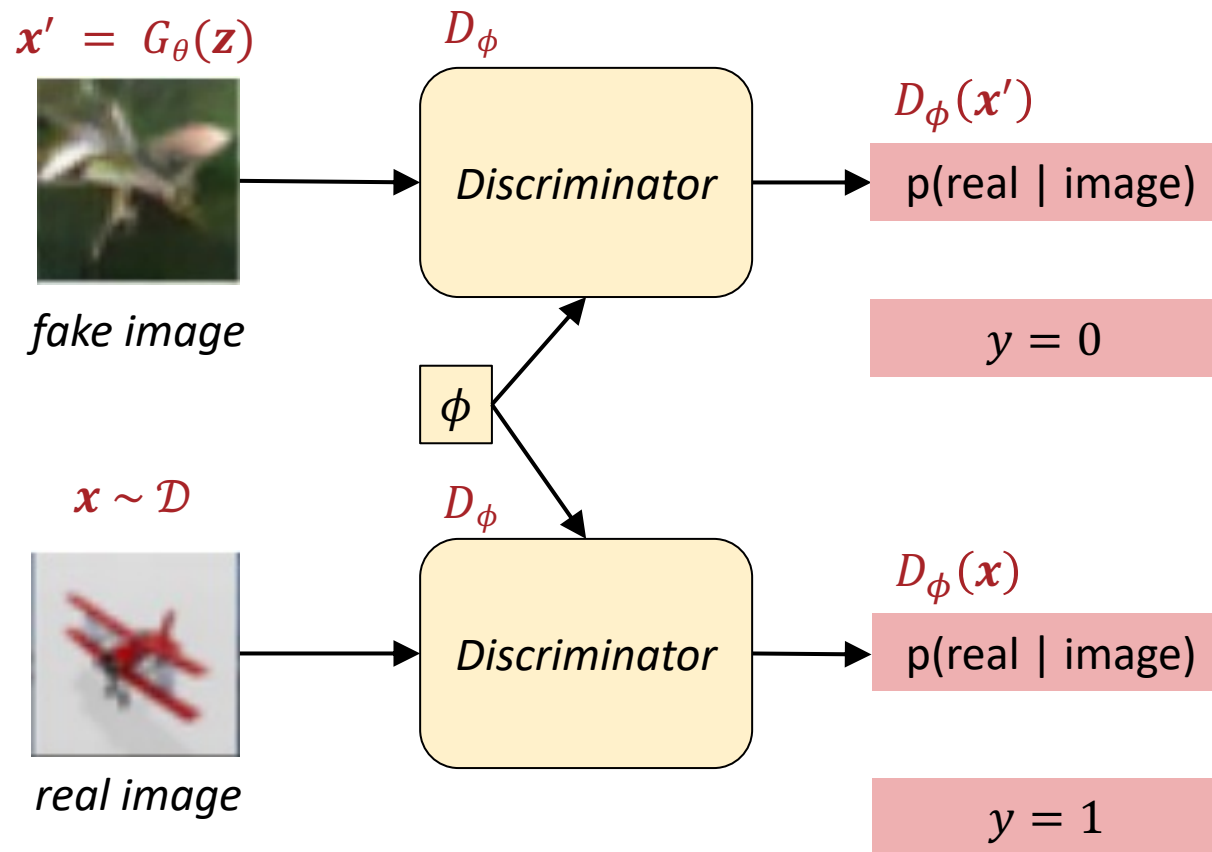




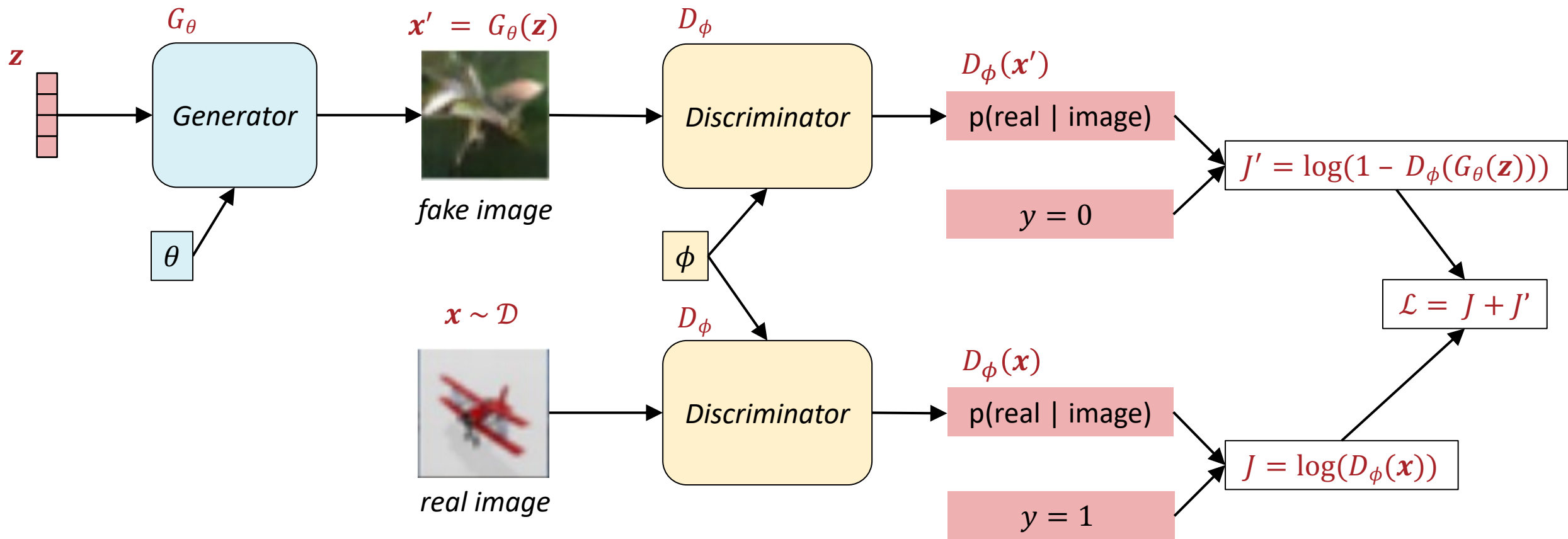
# GANs: Architecture



# GANs: Architecture



# GANs: Architecture

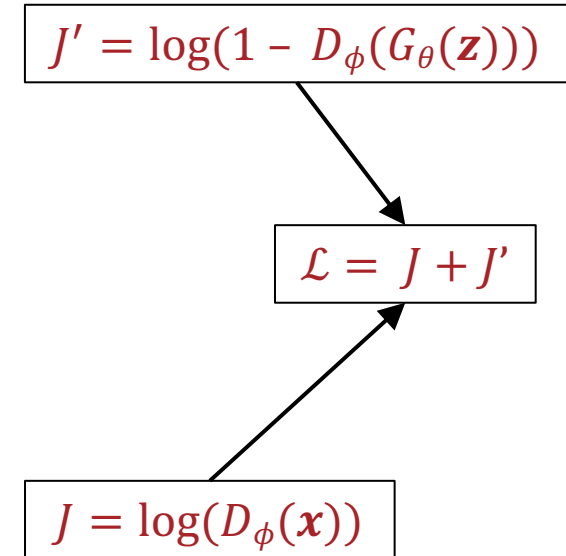


# GANs: Architecture

The discriminator is trying to maximize 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: Architecture

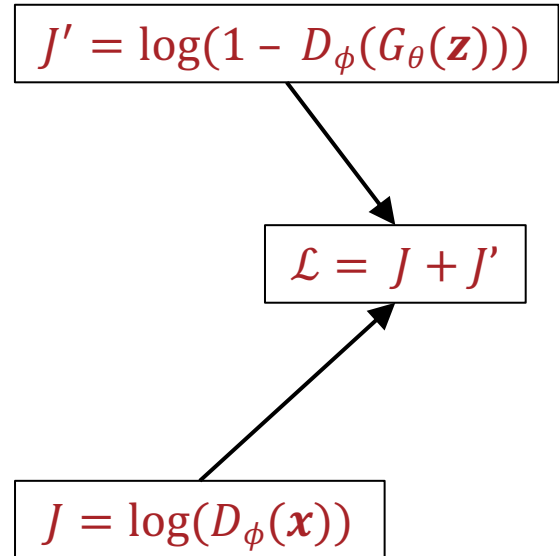
Both objectives (and hence, their sum) are differentiable

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

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

Training alternates between:

1. Keeping  $\theta$  fixed and backpropagating through  $D_{\phi}$
2. Keeping  $\phi$  fixed and backpropagating through  $G_{\theta}$



# GANs: Architecture

# 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:

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

**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:

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

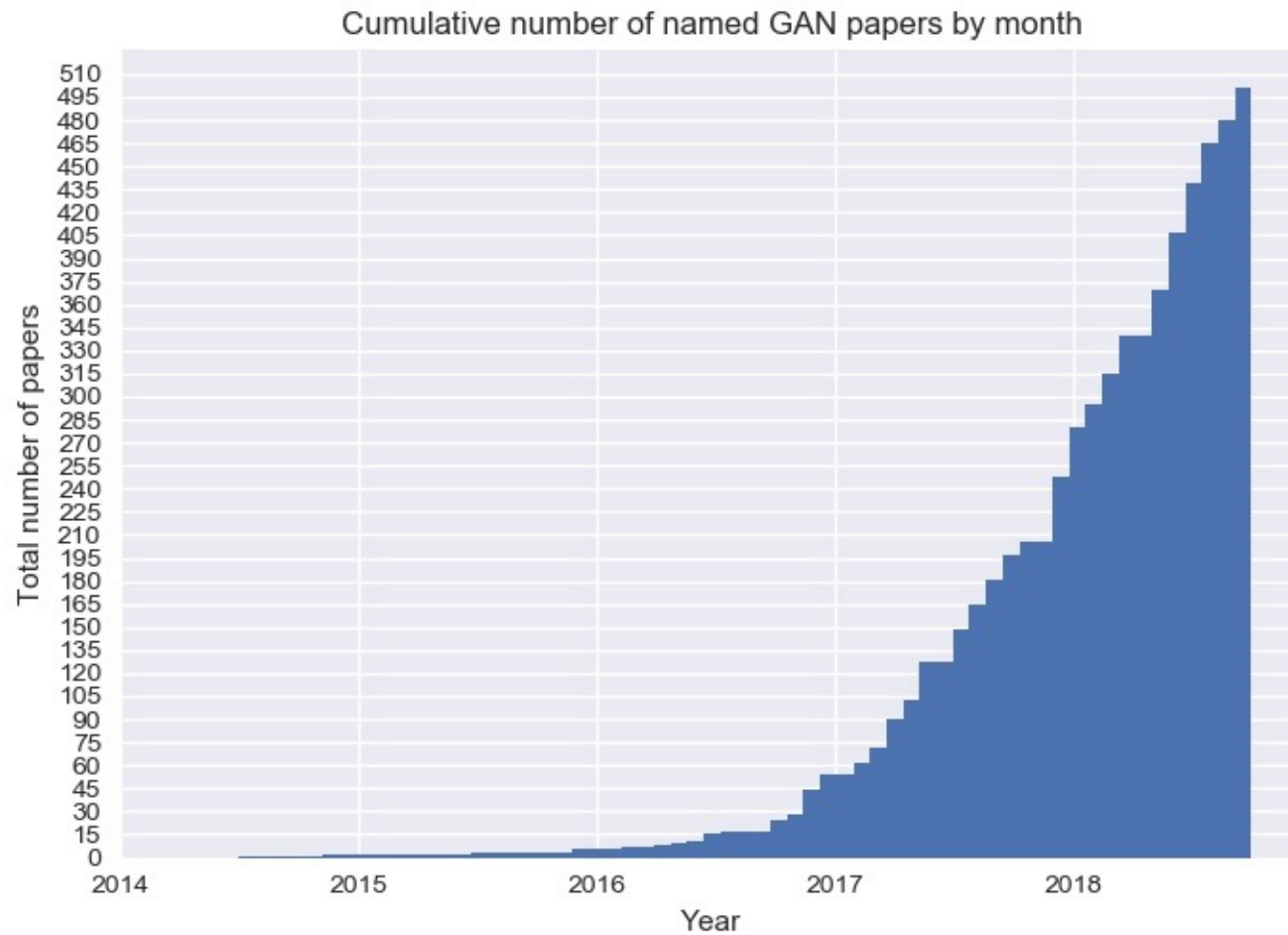
**end for**

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



# GANs Everywhere!



# The rise of vision transformers and diffusion models

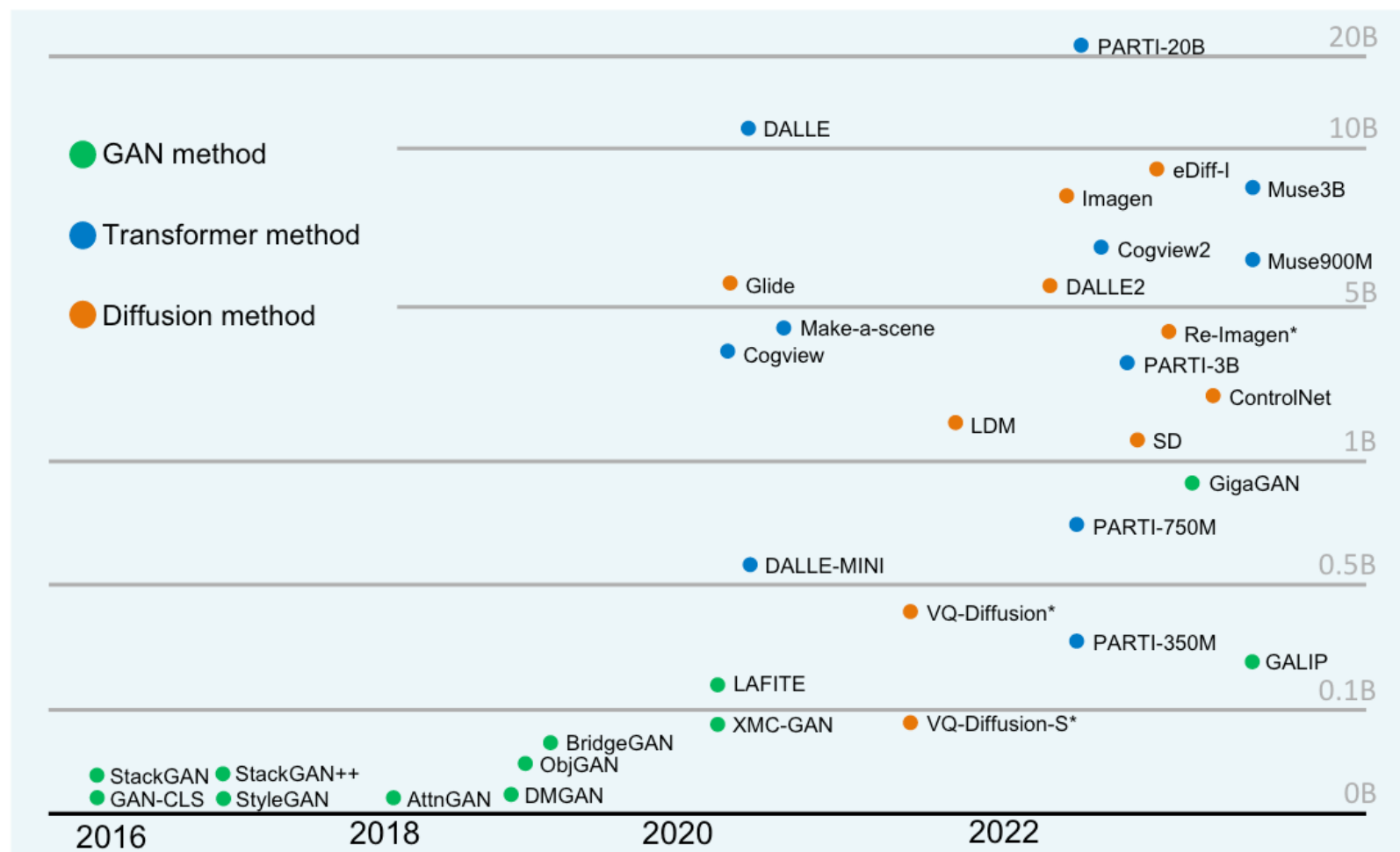


Fig. 5. Timeline of TTI model development, where green dots are GAN TTI models, blue dots are autoregressive Transformers and orange dots are Diffusion TTI models. Models are separated by their parameter, which are in general counted for all their components. Models with asterisk are calculated without the involvement of their text encoders.

But wait, what the heck are “vision transformers” and “diffusion”?

Take 10-423/623 next semester to find out!

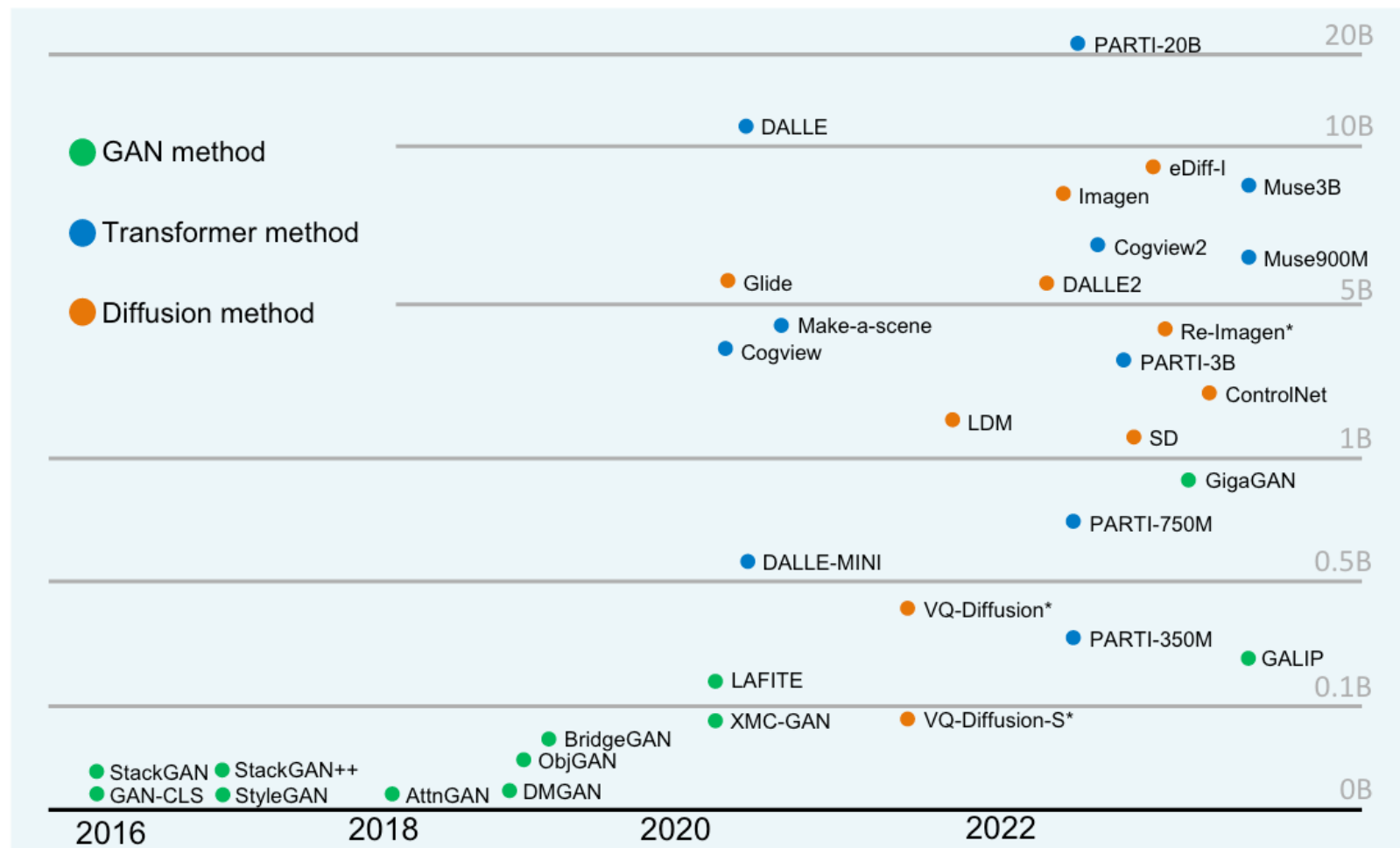


Fig. 5. Timeline of TTI model development, where green dots are GAN TTI models, blue dots are autoregressive Transformers and orange dots are Diffusion TTI models. Models are separated by their parameter, which are in general counted for all their components. Models with asterisk are calculated without the involvement of their text encoders.