10-301/601: Introduction to Machine Learning Lecture 26 – Generative Models for Vision

Matt Gormley & Henry Chai 4/21/25

Front Matter

- Announcements
 - HW9 released 4/17, due 4/24 (Thursday) at 11:59 PM
 - You may only use at most 2 late days on HW9
 - Exam 3 on 5/1 from 1 PM to 3 PM
 - We will not use the full 3-hour window
 - All topics from Lectures 17 to 25 (inclusive) are inscope, excluding the MLE/MAP portion of Lecture 17
 - Exam 1 and 2 content may be referenced but will not be the primary focus of any question
 - Please watch Piazza carefully for your room and seat assignments
 - You are allowed to bring one letter-size sheet of notes; you may put whatever you want on both sides

Q: So how do we actually solve for principal components?

• A: By maximizing the variance of our projections (which is equivalent is to minimizing the reconstruction error)

$$\widehat{\boldsymbol{v}} = \underset{\boldsymbol{v}: \|\boldsymbol{v}\|_{2}^{2}=1}{\operatorname{argmax}} \boldsymbol{v}^{T}(X^{T}X)\boldsymbol{v}$$

$$\boldsymbol{v}: \|\boldsymbol{v}\|_{2}^{2}=1$$

$$\mathcal{L}(\boldsymbol{v}, \lambda) = \boldsymbol{v}^{T}(X^{T}X)\boldsymbol{v} - \lambda(\|\boldsymbol{v}\|_{2}^{2}-1)$$

$$= \boldsymbol{v}^{T}(X^{T}X)\boldsymbol{v} - \lambda(\boldsymbol{v}^{T}\boldsymbol{v}-1)$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{v}} = (X^{T}X)\boldsymbol{v} - \lambda\boldsymbol{v}$$

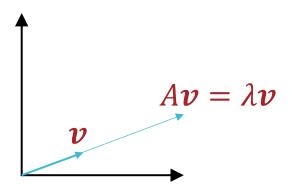
$$\rightarrow (X^{T}X)\widehat{\boldsymbol{v}} - \lambda\widehat{\boldsymbol{v}} = 0 \rightarrow (X^{T}X)\widehat{\boldsymbol{v}} = \lambda\widehat{\boldsymbol{v}}$$

• $\widehat{\boldsymbol{v}}$ is an eigenvector of X^TX and λ is the corresponding eigenvalue!

Background: Eigenvectors & Eigenvalues

• Given a square matrix $A \in \mathbb{R}^{N \times N}$, a vector $v \in \mathbb{R}^{N \times 1}$ is an eigenvector of A iff there exists some scalar λ such that

$$A\mathbf{v} = \lambda \mathbf{v}$$



Intuition: A scales or stretches

v but does not rotate it

Key property: the eigenvectors of <u>symmetric</u> matrices
 (e.g., the covariance matrix of a data set) are orthogonal!

11/25/24

Maximizing the Variance

$$\widehat{\boldsymbol{v}} = \underset{\boldsymbol{v}: \|\boldsymbol{v}\|_{2}^{2}=1}{\operatorname{argmax}} \, \boldsymbol{v}^{T}(X^{T}X) \boldsymbol{v}$$

$$(X^{T}X) \widehat{\boldsymbol{v}} = \lambda \widehat{\boldsymbol{v}} \, \rightarrow \, \widehat{\boldsymbol{v}}^{T}(X^{T}X) \widehat{\boldsymbol{v}} = \lambda \widehat{\boldsymbol{v}}^{T} \widehat{\boldsymbol{v}} = \lambda$$

- The first principal component is the eigenvector $\widehat{m v}_1$ that corresponds to the largest eigenvalue λ_1
- The second principal component is the eigenvector $\widehat{m v}_2$ that corresponds to the second largest eigenvalue λ_1
 - $\widehat{oldsymbol{v}}_1$ and $\widehat{oldsymbol{v}}_2$ are orthogonal
- Etc ...
- λ_i is a measure of how much variance falls along $\widehat{m{v}}_i$

11/25/24

Singular Value Decomposition (SVD) for PCA

• Every real-valued matrix $X \in \mathbb{R}^{N \times D}$ can be expressed as

$$X = USV^T$$

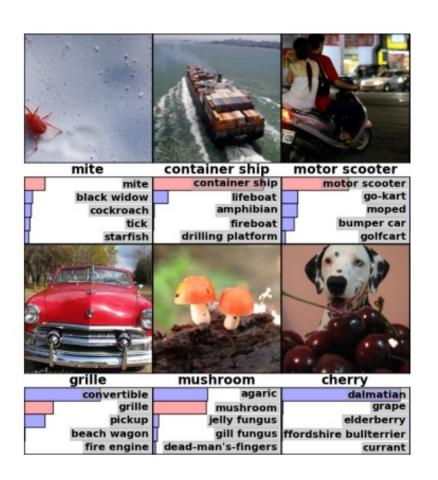
where:

- 1. $U \in \mathbb{R}^{N \times N}$ columns of U are eigenvectors of XX^T
- 2. $V \in \mathbb{R}^{D \times D}$ columns of V are eigenvectors of $X^T X$
- 3. $S \in \mathbb{R}^{N \times D}$ diagonal matrix whose entries are the eigenvalues of $X \to \text{squared entries}$ are the eigenvalues of XX^T and X^TX

- Image Classification
- Object Localization
- Object Detection
- Semantic Segmentation
- Instance Segmentation
- Image Captioning
- Image Generation

4/21/25

- Image Classification
- Object Localization
- Object Detection
- Semantic Segmentation
- Instance Segmentation
- Image Captioning
- Image Generation



- Image Classification
- Object Localization
- Object Detection
- Semantic Segmentation
- Instance Segmentation
- Image Captioning
- Image Generation



Given an image, predict a single label and a bounding box,
 represented as position (x, y)
 and height/width (h, w).

- Image Classification
- Object Localization
- Object Detection

Given an image, for each
 object predict a bounding box
 and a label, l: (x, y, w, h, l)

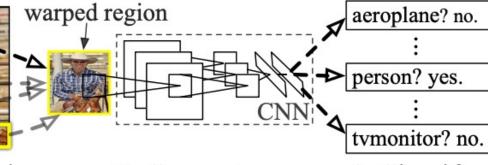
R-CNN: Regions with CNN features



1. Input image



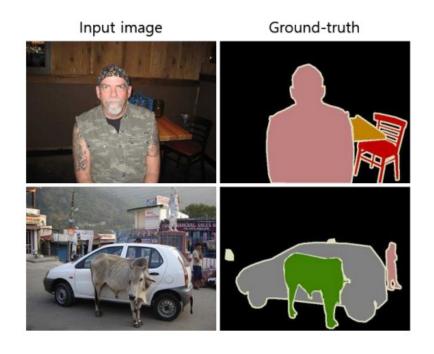
2. Extract region proposals (~2k)



3. Compute CNN features

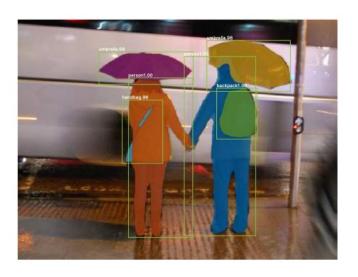
4. Classify regions

- Image Classification
- Object Localization
- Object Detection
- Semantic Segmentation
- Instance Segmentation
- Image Captioning
- Image Generation



 Given an image, predict a label for every pixel in the image

- Image Classification
- Object Localization
- Object Detection
- Semantic Segmentation
- Instance Segmentation
- Image Captioning
- Image Generation



Predict per-pixel labels as in semantic segmentation, but differentiate between different instances of the same label e.g., given two people, one should be labeled person-1 and one should be labeled person-2

- Image Classification
- Object Localization
- Object Detection
- Semantic Segmentation
- Instance Segmentation
- Image Captioning
- Image Generation



Ground Truth Caption: A little boy runs away from the approaching waves of the ocean.

Generated Caption: A young boy is running on the beach.



Ground Truth Caption: A brunette girl wearing sunglasses and a yellow shirt.

Generated Caption: A woman in a black shirt and sunglasses smiles.

- Take an image as input, and generate a sentence describing it as output
 - Dense captioning
 generates one description
 per bounding box
 - Typical methods use both a CNN and some sort of

language model

Source: https://dl.acm.org/doi/pdf/10.1145/3295748

- 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

4/21/25

sea anemone

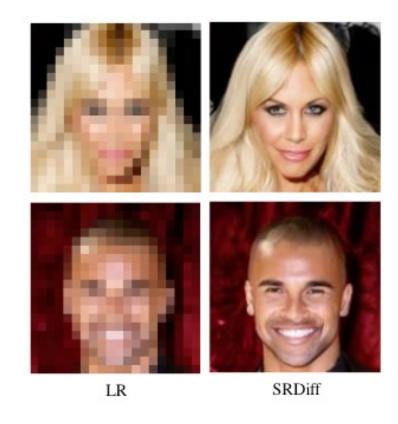
brain coral

slug



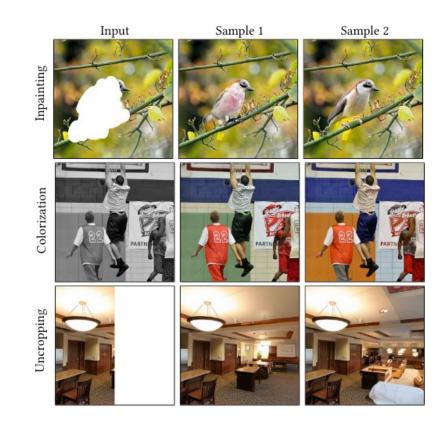
- Given a class label, sample a new image from that class
 - Image classification takes an image and predicts its label p(y|x)
 - Class-conditional generation is doing this in reverse p(x|y)

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



 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



- 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









 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

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

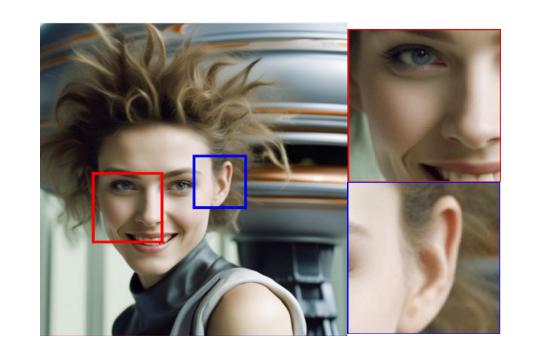
Prompt: Epic long distance cityscape photo of New York City flooded by the ocean and overgrown buildings and jungle ruins in rainforest, at sunset, cinematic shot, highly detailed, 8k, golden light



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

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

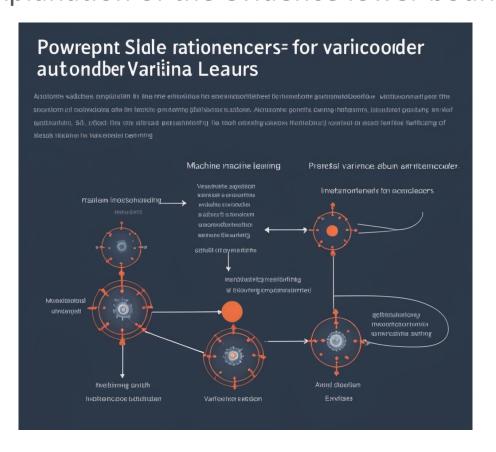
Image Generation



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

Slide Generation?

Prompt: powerpoint slide explainingvariational autoencoders for an intro toML course, easy to follow, with anexplanation of the evidence lower bound

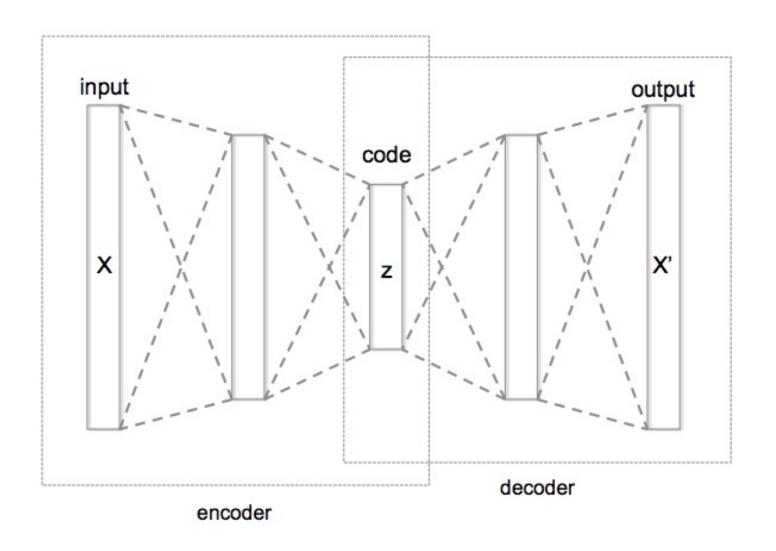


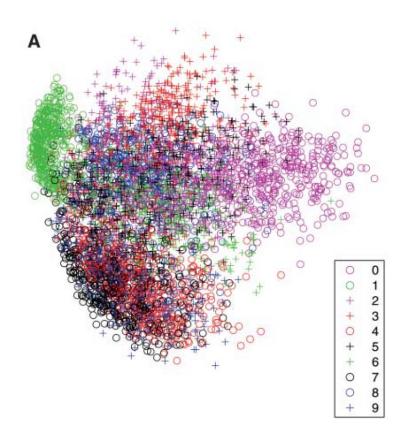
- Class-conditional generation
- Super resolution
- Image Editing
- Style transfer
- Text-to-image (TTI)generation

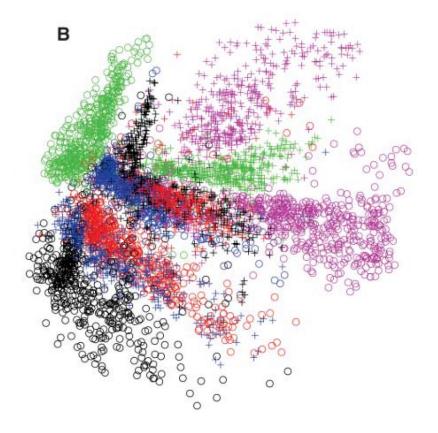
- Fundamental challenge: images are incredibly highdimensional objects with complex relationships between elements
- Idea: learn a low-dimensional representation of images, sample points in the low-dimensional space and project them up to the original image space

4/21/25 **23**

Deep Autoencoders

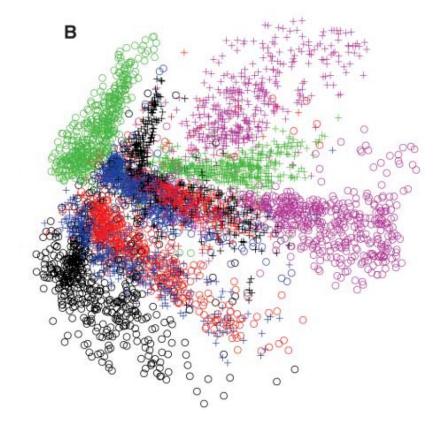




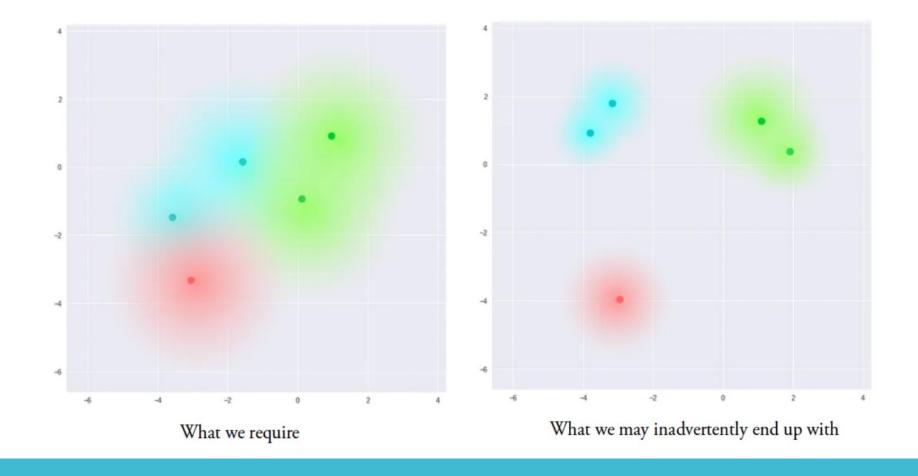


PCA (A) vs. Autoencoders (B) (Hinton and Salakhutdinov, 2006)

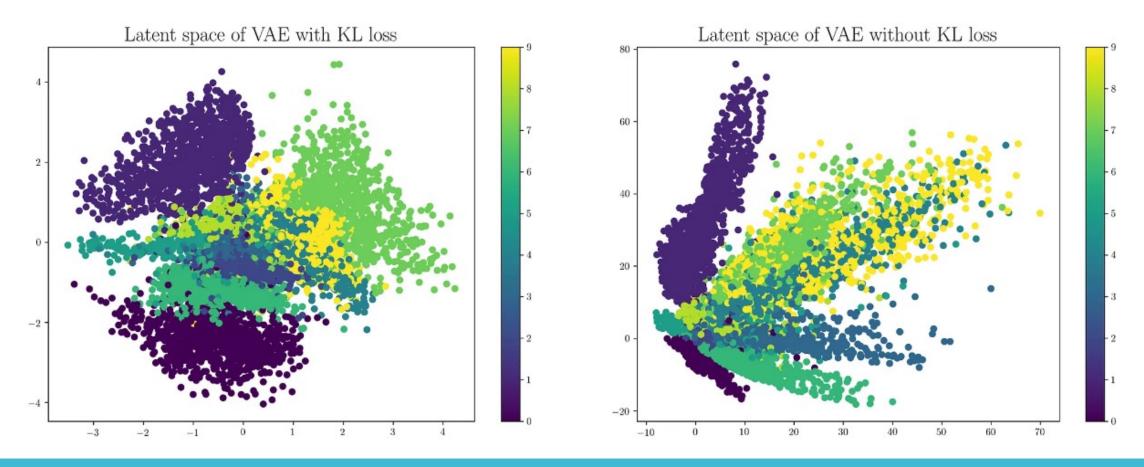
- Issue: latent space is sparse...
 - Sampling from latent space of an autoencoder creates outputs that are effectively identical to images in the training dataset



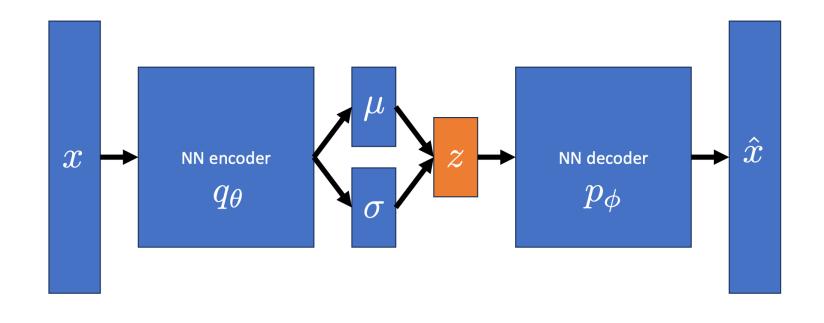
Autoencoder Latent Space

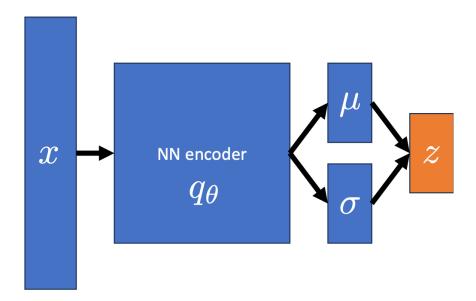


Autoencoder Latent Space



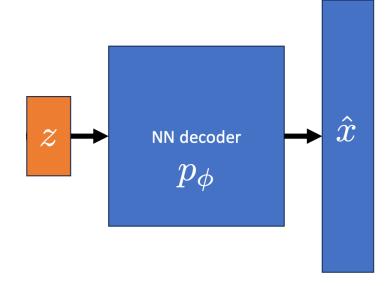
Variational Autoencoder Latent Space





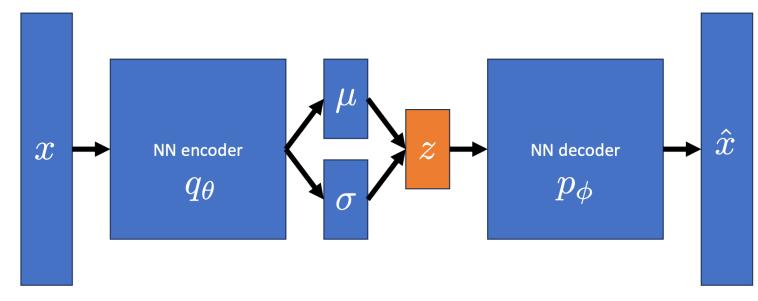
- Encoder learns a mean vector and a (diagonal) covariance matrix for each input
- These are used to sample a latent representation e.g.,

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



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

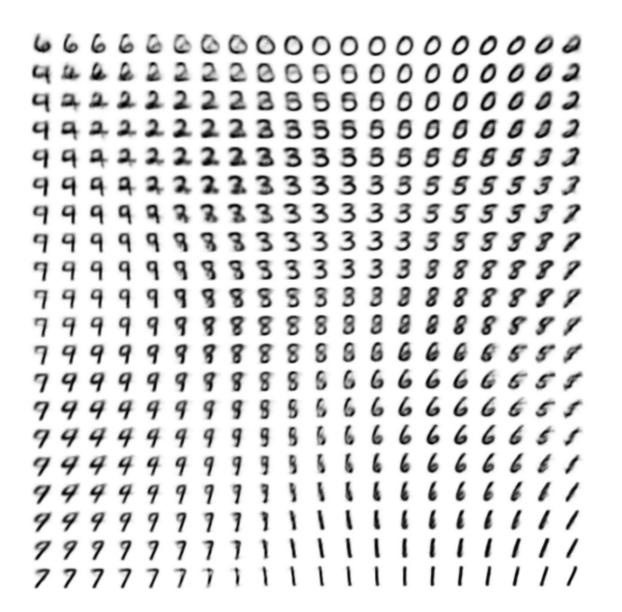
$$\widehat{\boldsymbol{x}}^{(i)} \mid \boldsymbol{z}^{(i)} \sim \mathcal{N}\left(\mu_{\phi}(\boldsymbol{z}^{(i)}), \sigma_{\phi}^{2}(\boldsymbol{z}^{(i)})\right)$$



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

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

Variational
Autoencoder:
Latent Space
Visualization



Variational Autoencoder: Generated Samples



Variational Autoencoder: Generated Samples?

3681796691 6757863485 2179712845 4819018894 7618641560 7592658197 222234480 0 2 3 8 0 7 3 8 5 7 0146460243 7128169861

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

```
28383857383681796691
83827935386757863485
35991395132179712845
19189334924819018894
27364302637618641560
59700838457592658197
69436285922222234480
849050000660238073857
74163036010146460243
2 + 20 4 3 1 9 5 0 7 / 28 1 6 9 8 6 /
```

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

Source: MNIST

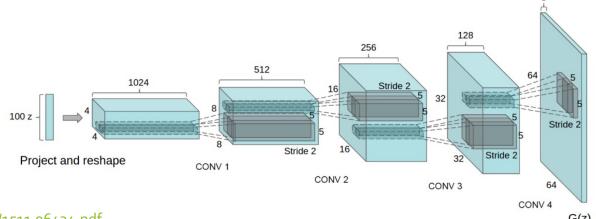
4/21/25

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

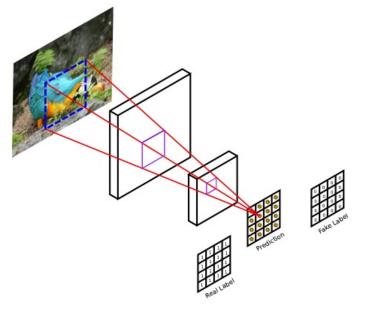
Generative Adversarial Networks (GANs)

- 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



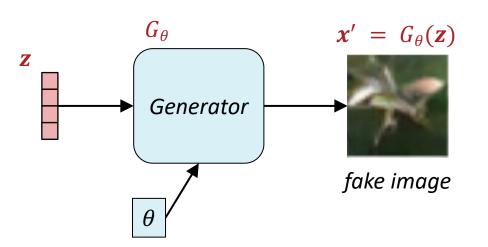
Generative Adversarial Networks (GANs)

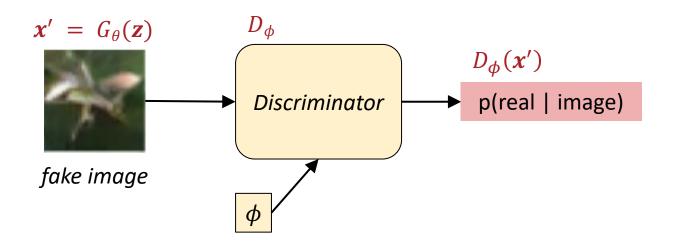
- A GAN consists of two (deterministic) models:
 - a generator that takes a vector of random noise as input, and generates an image
 - a **discriminator** that takes in an image classifies whether it is real (label = 1) or fake (label = 0)
- Example discriminator: PatchGAN
 - Traditional CNN that looks
 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

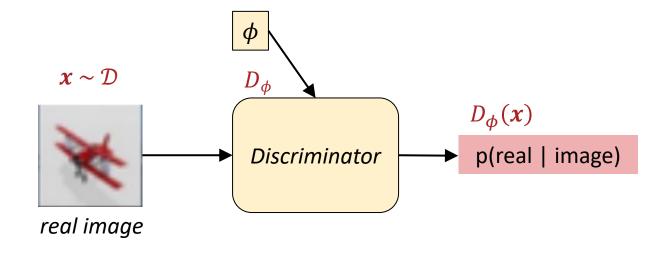


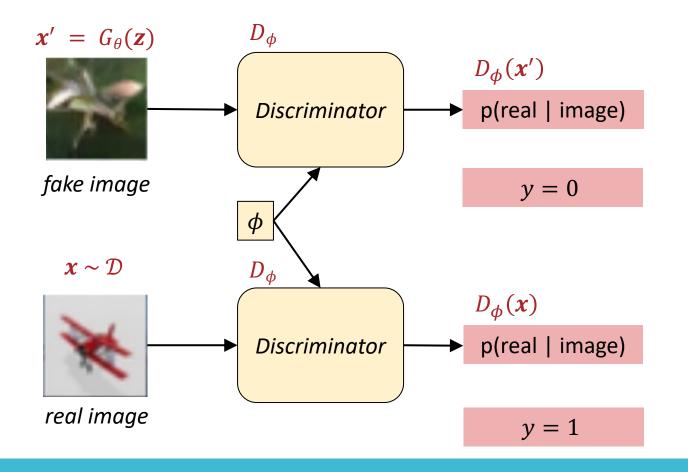
Generative Adversarial Networks (GANs): Training

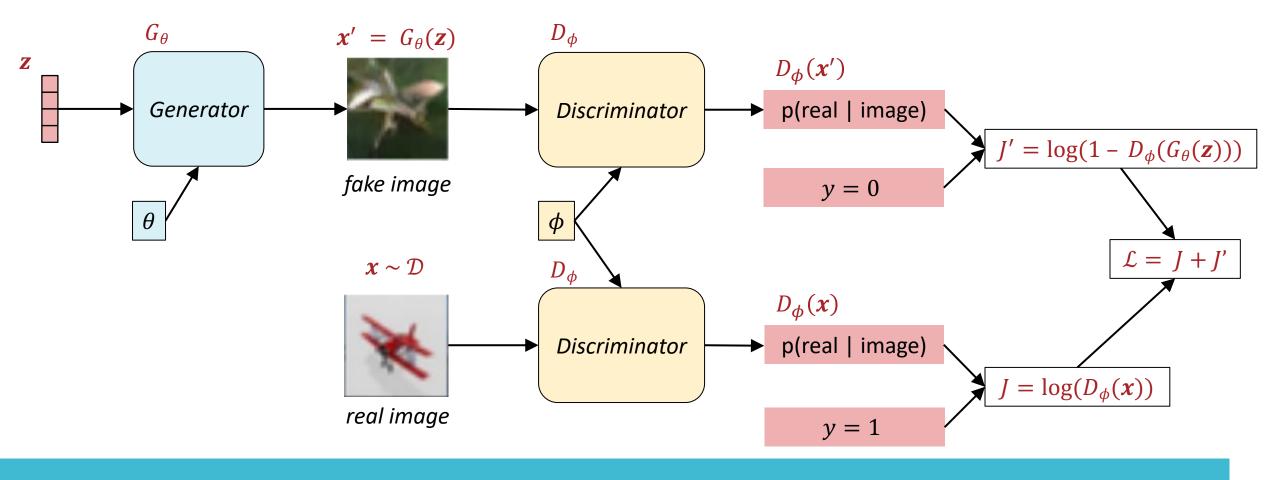
- A GAN consists of two (deterministic) models:
 - a generator that takes a vector of random noise as input, and generates an image
 - a discriminator that takes in an image classifies
 whether it is real (label = 1) or fake (label = 0)
 - Both models are typically (but not necessarily) neural networks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











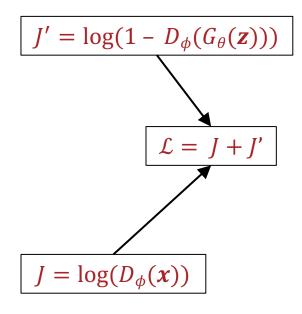
The discriminator is trying to maximize the usual cross-entropy

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

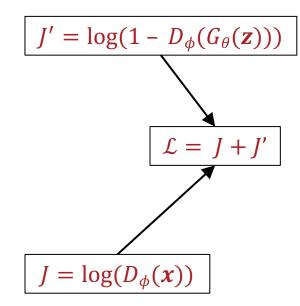
Both objectives (and hence, their sum) are differentiable

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

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

Training alternates between:

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



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 $\{z^{(1)}, \ldots, z^{(m)}\}$ 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:

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

end for

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

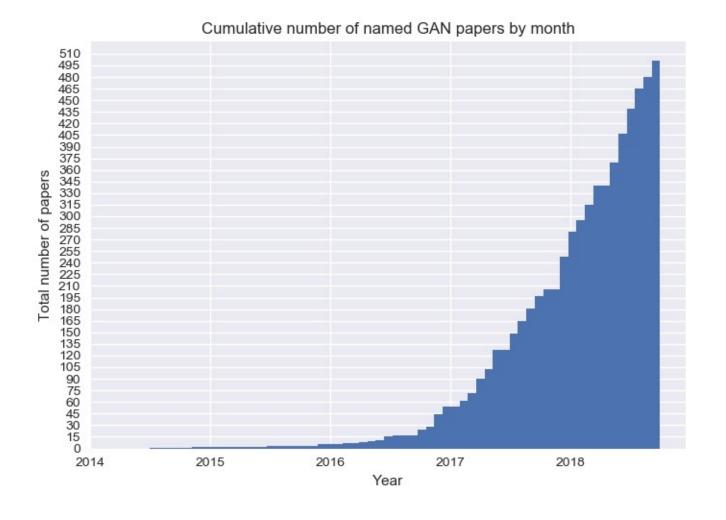
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log \left(1 - D \left(G \left(\boldsymbol{z}^{(i)} \right) \right) \right).$$

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

GANs Everywhere!



The rise of vision transformers and diffusion models

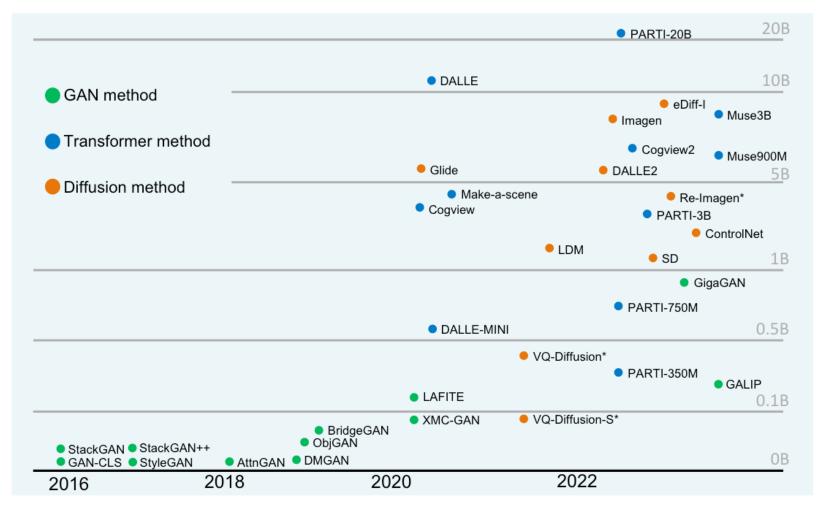


Fig. 5. Timeline of TTI model development, where green dots are GAN TTI models, blue dots are autoregressive Transformers and orange dots are Diffusion TTI models. Models are separated by their parameter, which are in general counted for all their components. Models with asterisk are calculated without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

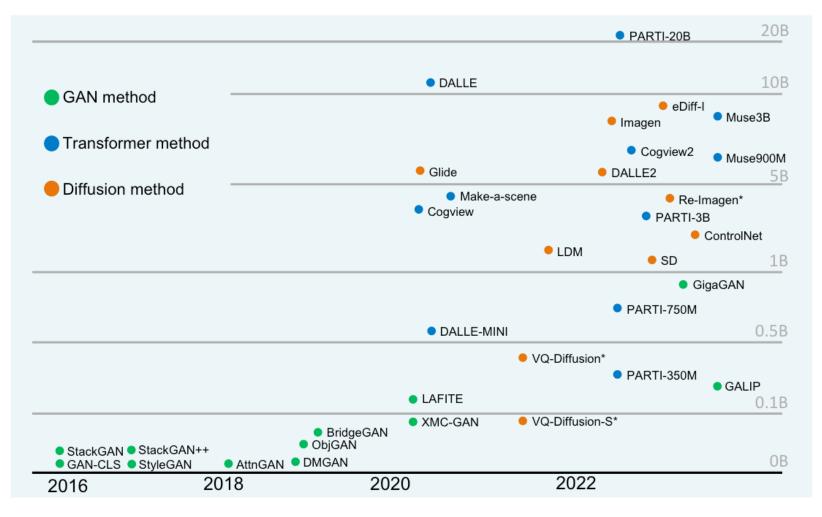


Fig. 5. Timeline of TTI model development, where green dots are GAN TTI models, blue dots are autoregressive Transformers and orange dots are Diffusion TTI models. Models are separated by their parameter, which are in general counted for all their components. Models with asterisk are calculated without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