Topics:

- Image Segmentation and Decoders
- Generative Adversarial Networks

# CS 4644-DL / 7643-A
# ZSOLT KIRA

- **Assignment 3**
  - Due **March 9th 11:59pm EST**
  - **Oops.** Diffusion models accidentally included. No need to do it by Mar 9th! @258

- **Projects**
  - Project proposal due **March ~~15th~~ 17th**
  - Proposal description out on canvas @256

- Meta office hours today 3pm ET on language models

# Self-Attention Layer

One **query** per **input vector**

**Inputs**:
**Input vectors**: $X$ (Shape: $N_X \times D_X$)
**Key matrix**: $W_K$ (Shape: $D_X \times D_Q$)
**Value matrix:** $W_V$ (Shape: $D_X \times D_V$)
**Query matrix**: $W_Q$ (Shape: $D_X \times D_Q$)

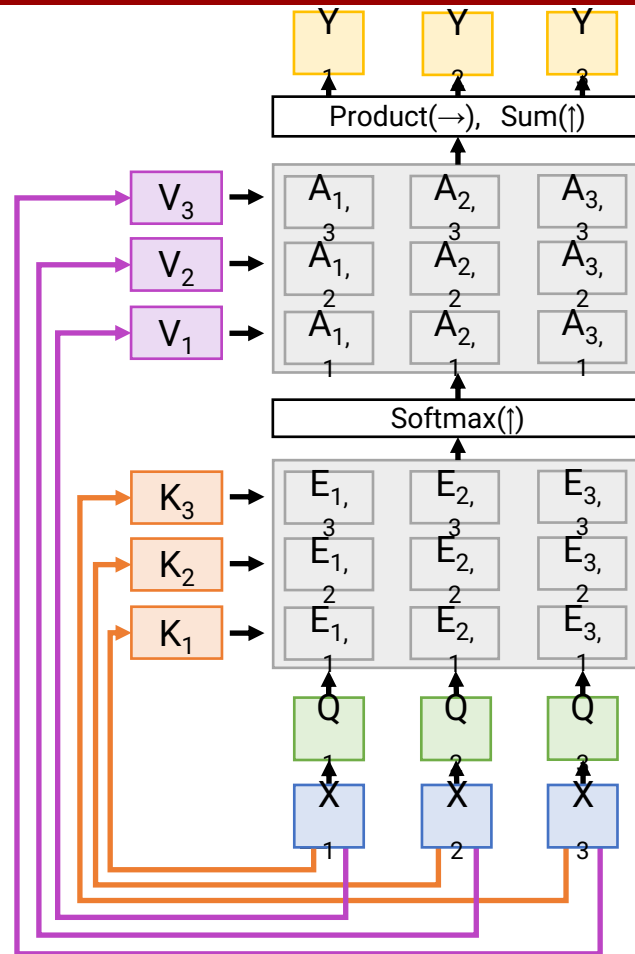**Computation**:
**Query vectors**: $Q = XW_Q$
**Key vectors**: $K = XW_K$ (Shape: $N_X \times D_Q$)
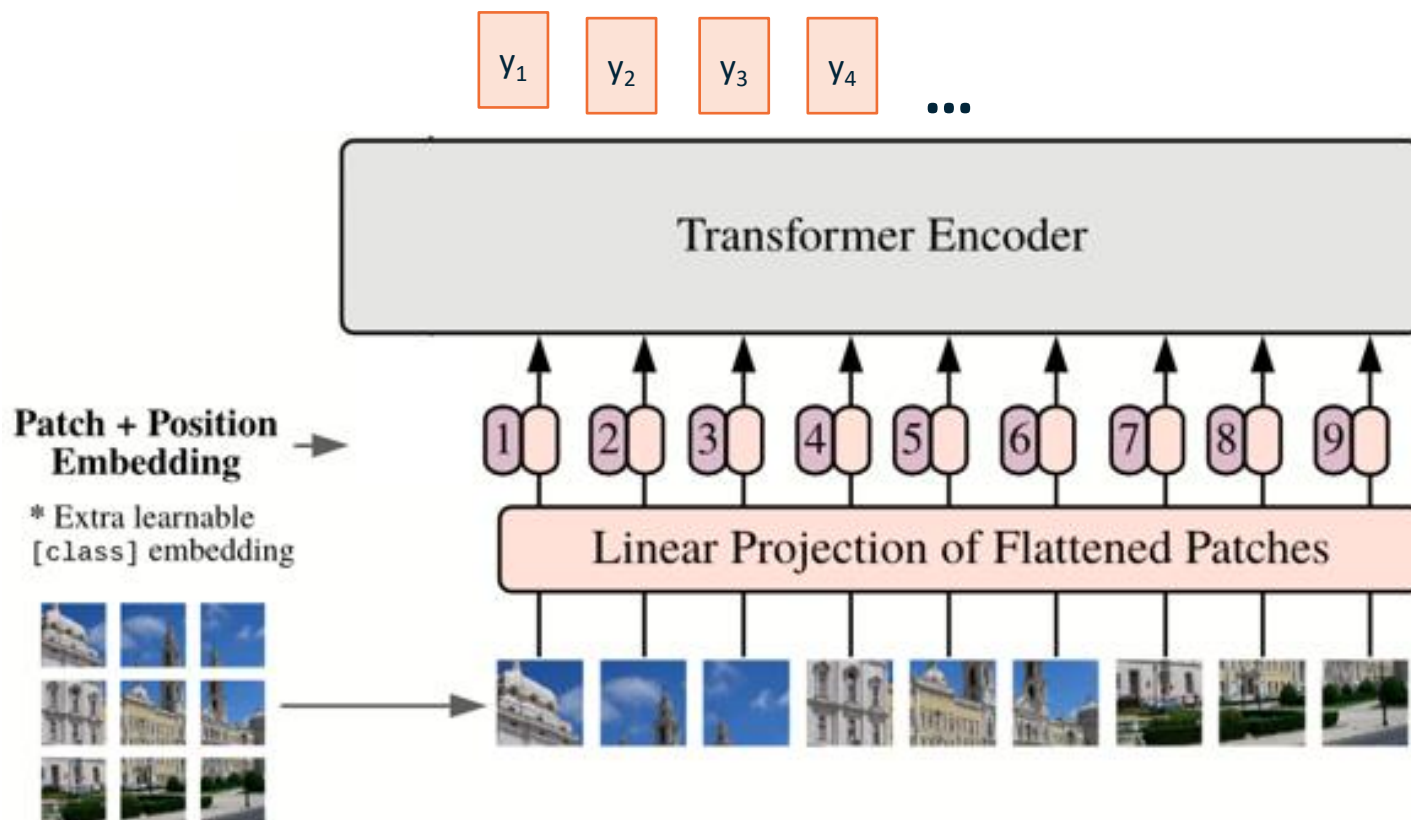**Value vectors**: $V = XW_V$ (Shape: $N_X \times D_V$)
**Similarities**: $E = QK^T$ (Shape: $N_X \times N_X$) $E_{i,j} = Q_i \cdot K_j / \text{sqrt}(D_Q)$
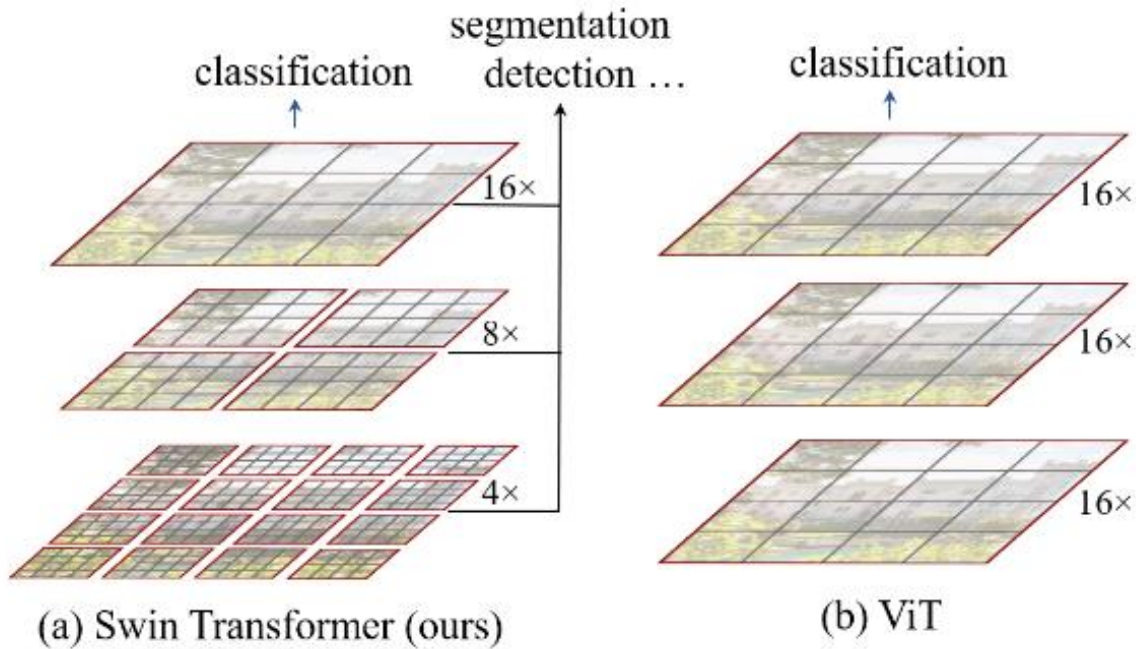**Attention weights**: $A = \text{softmax}(E, \text{dim}=1)$ (Shape: $N_X \times N_X$)
**Output vectors**: $Y = AV$ (Shape: $N_X \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

- **How do we do classification?**

**Patches as input to Self-Attention**

Georgia Tech

(a) Swin Transformer (ours)  (b) ViT

**Ideas:**
- Use smaller patches (4x4x3)
- Project them to lower dimension (4)
- Merge tokens at deeper levels
- Full attention => Window attention
  - => Shifted window attention

Swin Transformer: Hierarchical Vision Transformer using Shifted Windows
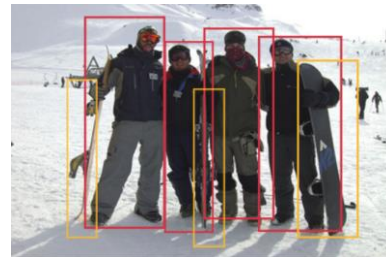Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, Baining Guo

**Swin Transformers**

Georgia Tech

**Classification**
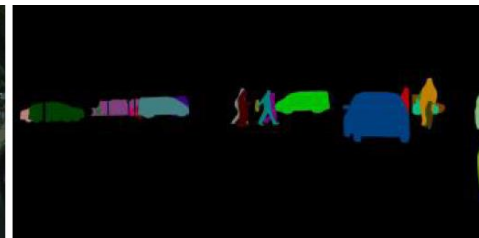(Class distribution per image)

**Object Detection**
(List of bounding boxes with class distribution per box)

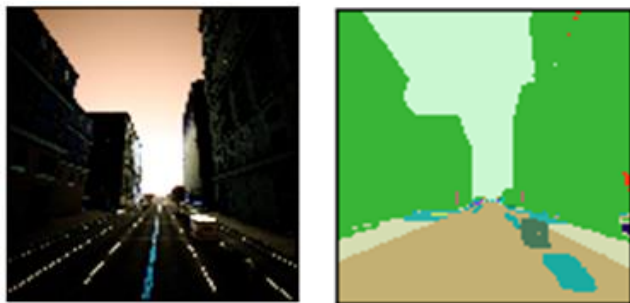**Semantic Segmentation**
(Class distribution per pixel)

**Instance Segmentation**
(Class distribution per pixel with unique ID)

Car    Coffee Cup    Bird

**Computer Vision Tasks**

Georgia Tech

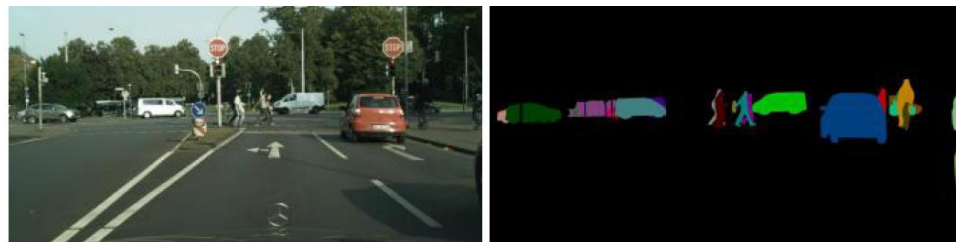# Given an image, output another image

- Each output contains class distribution per pixel

- More generally an image-to-image problem



**Semantic Segmentation**
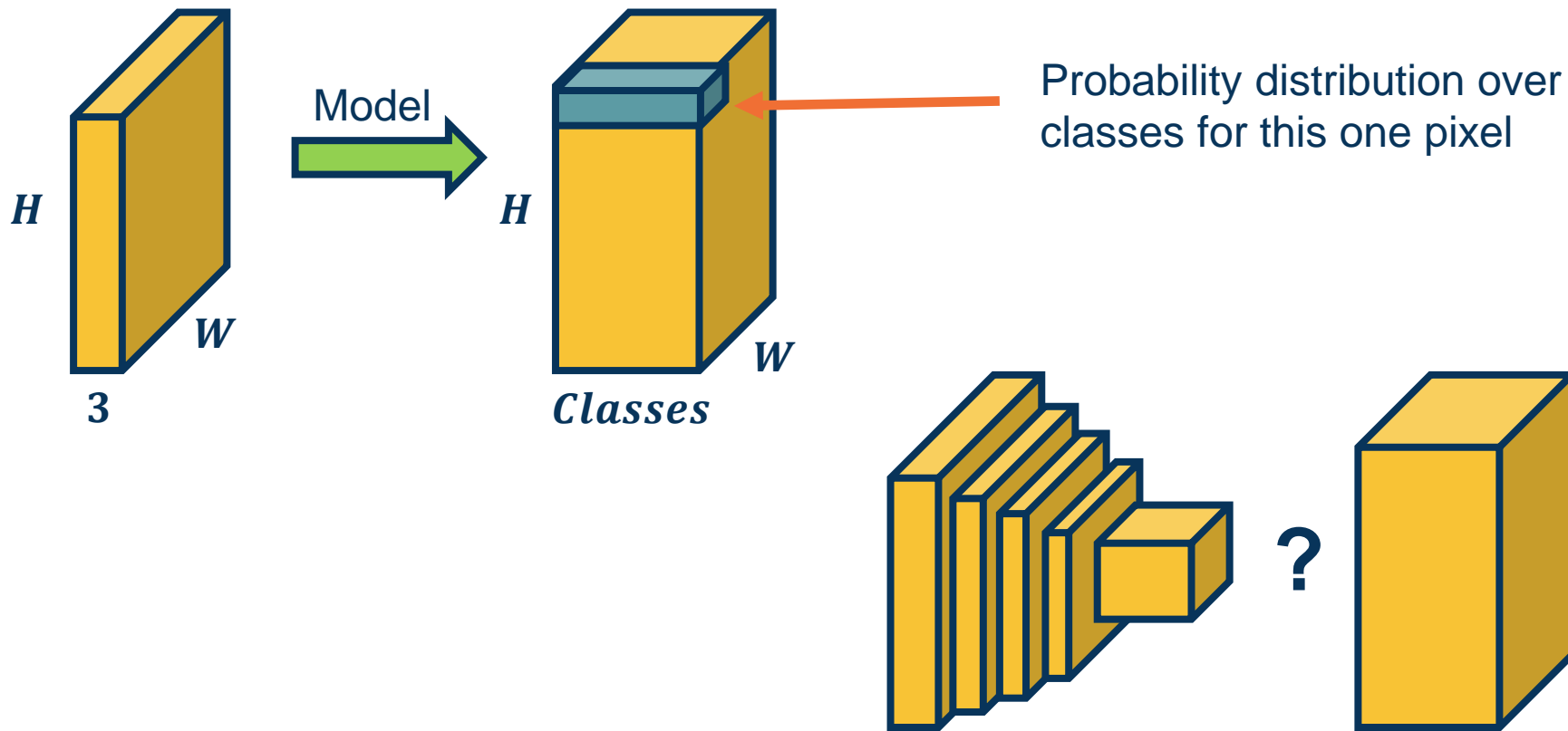(Class distribution per pixel)



**Instance Segmentation**
(Class distribution per pixel with unique ID)

**Segmentation Tasks**

Georgia Tech

Model

H · W · 3

H · W · Classes

Probability distribution over classes for this one pixel

?

Convolutional Neural Network (CNN)

Image

Convolution + Non-Linear Layer

Pooling Layer

Convolution + Non-Linear Layer

Useful, lower-dimensional features

Encoder

We can develop learnable or non-learnable upsampling layers!

Decoder

(De)Convolution + Non-Linear Layer

(Un)Pooling Layer

(De)Convolution + Non-Linear Layer

"Image"

Useful, lower-dimensional features

Idea: "De"Convolution and UnPooling

Georgia Tech

**Example :** Max pooling

⬡ Stride window across image but perform per-patch **max operation**

$$X(0:1, 0:1) = \begin{bmatrix} 100 & 150 \\ 100 & 200 \end{bmatrix} \implies \max(0:1,0:1) = 200$$

Copy value to position chosen as max in encoder, fill reset of this window with zeros



**Pooling**

$H = 5$

$W = 5$

**UnPooling**

$W = 5$

$H = 5$

**Idea:** Remember max elements in encoder! Copy value from equivalent position, rest are zeros

**Max Unpooling**

$$X = \begin{bmatrix} 120 & 150 & 120 \\ 100 & 50 & 110 \\ 25 & 25 & 10 \end{bmatrix} \quad \Rightarrow \quad Y = \begin{bmatrix} 150 & 150 \\ 100 & 110 \end{bmatrix}$$

2x2 max pool

**Encoder**

**Decoder**

2x2 max unpool

$$X = \begin{bmatrix} 300 & 450 \\ 100 & 250 \end{bmatrix} \quad \Rightarrow \quad Y = \begin{bmatrix} 0 & 300 & - \\ 0 & 0 & - \\ - & - & - \end{bmatrix}$$

**Max Unpooling Example (one window)**

Georgia Tech

$$X_{enc} = \begin{bmatrix} 120 & 150 & 120 \\ 100 & 50 & 110 \\ 25 & 25 & 10 \end{bmatrix}$$ ➡ **2x2 max pool** $$Y_{enc} = \begin{bmatrix} 150 & 150 \\ 100 & 110 \end{bmatrix}$$

**Contributions from multiple windows are summed**

**Encoder**

**Decoder**

**2x2 max unpool**

$$X_{dec} = \begin{bmatrix} 300 & 450 \\ 100 & 250 \end{bmatrix}$$ ➡ $$Y_{dec} = \begin{bmatrix} 0 & 300+450 & 0 \\ 100 & 0 & 250 \\ 0 & 0 & 0 \end{bmatrix}$$

**Max Unpooling Example**

Georgia Tech

**Convolutional Neural Network (CNN)**

Image

Convolution + Non-Linear Layer

Pooling Layer

Convolution + Non-Linear Layer

Useful, lower-dimensional features

**Encoder**

We pull max indices from corresponding layers (requires symmetry in encoder/decoder)

**Decoder**

(De)Convolution + Non-Linear Layer

Useful, lower-dimensional features

(Un)Pooling Layer

(De)Convolution + Non-Linear Layer

"Image"

# How can we *upsample* using convolutions and learnable kernel?

**Normal Convolution**



$H = 5$

$W = 5$

$k_1 = 3$

$k_2 = 3$

$H - k_1 + 1$

$W - k_2 + 1$

**Transposed Convolution (also known as "deconvolution", fractionally strided conv)**

Idea: Take each input pixel, multiply by learnable kernel, "stamp" it on output



$k_1 = 3$

$k_2 = 3$

$H = 5$

Georgia Tech

$$X = \begin{bmatrix} 120 & 150 & 120 \\ 100 & 50 & 110 \\ 25 & 25 & 10 \end{bmatrix} \qquad K = \begin{bmatrix} 1 & -1 \\ 2 & -2 \end{bmatrix}$$

**Contributions from multiple windows are summed**

$$\begin{bmatrix} 120 & -120 & 0 & 0 \\ 240 & -240 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

**Incorporate**
$X(0,0)$

$$\begin{bmatrix} 120 & -120 + 150 & -150 & 0 \\ 240 & -240 + 300 & -300 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

**Incorporate**
$X(1,0)$

**Transposed Convolution Example**

Georgia Tech

Convolutional Neural Network (CNN)

We can either learn the kernels, or take corresponding encoder kernel and rotate 180 degrees (no decoder learning)

Useful, lower-dimensional features

Image

Convolution + Non-Linear Layer

Pooling Layer

Convolution + Non-Linear Layer

Decoder

(De)Convolution + Non-Linear Layer

(Un)Pooling Layer

(De)Convolution + Non-Linear Layer

"Image"

Encoder

Useful, lower-dimensional features

**Symmetry in Encoder/Decoder**

Georgia Tech

We can start with a pre-trained trunk/backbone (e.g. network pretrained on ImageNet)!

# U-Net

**You can have skip connections to bypass bottleneck!**



*Ronneberger, et al., "U-Net: Convolutional Networks for Biomedical Image Segmentation", 2015*

Georgia Tech

## Summary

- Various ways to get **image-like outputs,** for example to predict segmentations of input images

- We can have various upsampling layers that actually increase the size

- Encoder/decoder architectures are popular ways to leverage these to perform general image-to-image tasks

- Other methods exist:

    - Fully convolutional neural networks

Georgia Tech

**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.

Georgia Tech

# Supervised Learning

x → **Classification** → y    **Discrete**

x → **Regression** → y    **Continuous**

# Unsupervised Learning

x → **Clustering** → c    **Discrete**

x → **Dimensionality Reduction** → z    **Continuous**

x → **Density Estimation** → p(x)    **On simplex**

Unsupervised Learning

Georgia Tech

# Traditional unsupervised learning methods:



Density estimation

**Modeling $P(x)$**

Deep Generative Models



Clustering

**Comparing/ Grouping**

Metric learning & clustering



Principal Component Analysis

**Representation Learning**

Almost all deep learning!

Similar in deep learning, but **from neural network/learning perspective**

**What to Learn?**

Georgia Tech

# Discriminative vs. Generative Models

- Discriminative models model $P(y|x)$

  - Example: Model this via neural network, SVM, etc.

- Generative models model $P(x)$

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Generative Models**

# Discriminative vs. Generative Models

- Discriminative models model $P(y|x)$

  - Example: Model this via neural network, SVM, etc.

- Generative models model $P(x)$

- We can parameterize our model as $P(x, \theta)$ and use maximum likelihood to optimize the parameters given an unlabeled dataset:

$$\boldsymbol{\theta}^* = \arg\max_{\boldsymbol{\theta}} \prod_{i=1}^{m} p_{\text{model}}\left(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}\right)$$

$$= \arg\max_{\boldsymbol{\theta}} \log \prod_{i=1}^{m} p_{\text{model}}\left(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}\right)$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{m} \log p_{\text{model}}\left(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}\right).$$

- They are called generative because they can often generate *samples*

  - Example: Multivariate Gaussian with estimated parameters $\boldsymbol{\mu}, \boldsymbol{\sigma}$

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Generative Models**

Georgia Tech

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Generative Models**

Generative Adversarial Networks (GANs)

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Generative Models**

- *Implicit* generative models do not actually learn an explicit model for $p(x)$

- Instead, learn to *generate samples* from $p(x)$
  - Learn good feature representations
  - Perform data augmentation
  - Learn world models (a simulator!) for reinforcement learning

- How?
  - **Learn to sample** from a neural network output
  - **Adversarial training** that uses one network's predictions to train the other (dynamic loss function!)
  - **Lots of tricks** to make the optimization more stable

**Implicit Models**

- We would like to *sample* from $p(x)$ using a neural network
- **Idea:**
  - Sample from a simple distribution (Gaussian)
  - Transform the sample to $p(x)$



$N(\mu, \sigma)$        **Neural Network**        $p(x)$

Samples        Samples

- Input can be a vector with (independent) Gaussian random numbers
- We can use a CNN to generate images!

**Generator**

**Vector of Random Numbers**

$N(\mu, \sigma)$

**Neural Network**

$p(x)$

- **Goal:** We would like to generate *realistic* images. How can we drive the network to learn how to do this?

- **Idea:** Have *another* network try to distinguish a real image from a generated (fake) image

  - **Why?** Signal can be used to determine how well it's doing at generation



**Generator**

**Discriminator**

**Vector of Random Numbers**

**Real or Fake?**

$N(\mu, \sigma)$

**Neural Network**

$p(x)$

Georgia Tech

**Generator:** Update weights to improve realism of generated images

**Discriminator:** Update weights to better discriminate

Generator

Discriminator

Vector of Random Numbers

Mini-batch of real & fake data

Cross-entropy (Real or Fake?) We know the answer (self-supervised)

**Question: What loss functions can we use (for each network)?**

**Generative Adversarial Networks (GANs)**

- Since we have two networks competing, this is a mini-max two player game
  - Ties to game theory
  - Not clear what (even local) Nash equilibria are for this game

**Mini-max Two Player Game**

- Since we have two networks competing, this is a mini-max two player game
  - Ties to game theory
  - Not clear what (even local) Nash equilibria are for this game

- The full mini-max objective is:

$$\min_G \max_D V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$$

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Mini-max Two Player Game**

$$\min_{G} \max_{D} V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$$

- where $D(x)$ is the discriminator outputs probability ([0,1]) of **real** image
- $x$ is a **real image** and $G(z)$ is a **generated** image

- The discriminator wants to **maximize** this:
  - $D(x)$ is pushed up (to 1) because $x$ is a real image
  - $1 - D(G(z))$ is also pushed up to 1 (so that $\text{D}(\text{G}(\text{z}))$ is pushed down to 0)
  - In other words, discriminator wants to classify real images as real (1) and fake images as fake (0)

**Discriminator Perspective**

Georgia Tech

$$\min_{G} \max_{D} V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$$

- where $D(x)$ is the discriminator outputs probability ([0,1]) of **real** image
- $x$ is a **real image** and $G(z)$ is a **generated** image

- The generator wants to **minimize** this:
  - $1 - D(G(z))$ is pushed down to 0 (so that $\text{D}(\text{G}(z))$ is pushed up to 1)
  - This means that the generator is **fooling** the discriminator, i.e. succeeding at generating images that the discriminator can't discriminate from real

- Since we have two networks competing, this is a mini-max two player game
  - Ties to game theory
  - Not clear what (even local) Nash equilibria are for this game

- The full mini-max objective is:

**Sample from fake**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$$

**Generator *minimizes***

**How well discriminator does (0 for fake)**

  - where $D(x)$ is the discriminator outputs probability ([0,1]) of **real** image
  - $x$ is a **real image** and $G(z)$ is a **generated** image

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Mini-max Two Player Game**

Georgia Tech

- Since we have two networks competing, this is a mini-max two player game
  - Ties to game theory
  - Not clear what (even local) Nash equilibria are for this game

- The full mini-max objective is:

**Sample from real**　　　　　　**Sample from fake**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$$
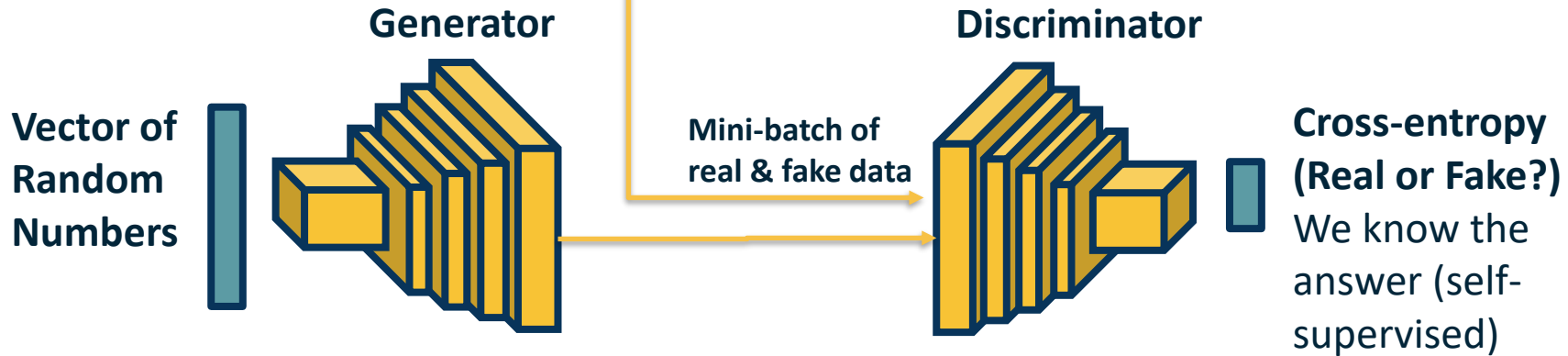
**Discriminator *maximizes***　　**How well discriminator does (1 for real)**　　**How well discriminator does (0 for fake)**

- where $D(x)$ is the discriminator outputs probability ([0,1]) of **real** image
- $x$ is a **real image** and $G(z)$ is a **generated** image

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Mini-max Two Player Game**

Georgia Tech

**Generator**

**Vector of Random Numbers**

**Mini-batch of real & fake data**

**Discriminator**

**Cross-entropy (Real or Fake?)** We know the answer (self-supervised)

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

**Generator Loss**

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

**Discriminator Loss**

**Generative Adversarial Networks (GANs)**

- The generator part of the objective does not have good gradient properties

$$\min_G \max_D V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$$

- High gradient when $D\big(G(z)\big)$ is high (that is, discriminator is wrong)
- We want it to improve when samples are *bad* (discriminator is right)



- Alternative objective, ***maximize***:

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(D_{\theta_d}(G_{\theta_g}(z)))$$

*Plot from CS231n, Fei-Fei Li, Justin Johnson, Serena Yeung*

**Converting to Max-Max Game**

Georgia Tech

**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 $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.

        • Update the discriminator by ascending its stochastic gradient:

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

    **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:

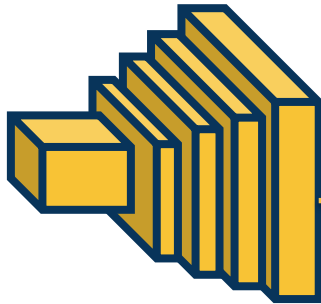$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right).$$

**end for**

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

---

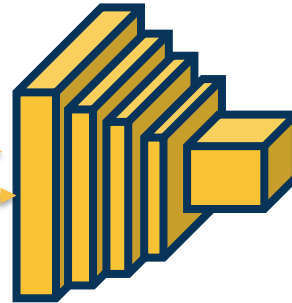*Goodfellow, NeurIPS 2016 Generative Adversarial Nets*

# Final Algorithm

Georgia Tech

At the end, we have:
- An *implicit* generative model!
- Features from discriminator

**Generator**

**Discriminator**

**Vector of Random Numbers**

**Mini-batch of real & fake data**

**Cross-entropy (Real or Fake?)** We know the answer (self-supervised)
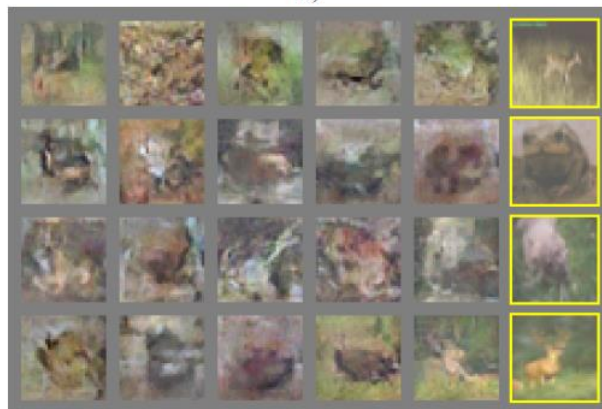
Georgia Tech

a)

b)

c)

d)

- Low-resolution images but look decent!

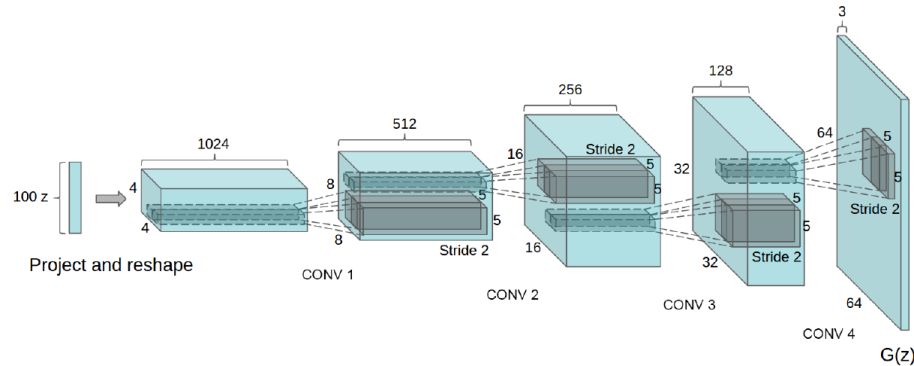- Last column are nearest neighbor matches in dataset

**Early Results**

- GANs are very difficult to train due to the mini-max objective

- Advancements include:
  - More stable architectures
  - Regularization methods to improve optimization
  - Progressive growing/training and scaling

*Goodfellow, NeurIPS 2016 Generative Adversarial Nets*

**Difficulty in Training**

Architecture guidelines for stable Deep Convolutional GANs

- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

*Radford et al., Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks*

**DCGAN**

- Training GANs is difficult due to:
  - Minimax objective – For example, what if generator learns to memorize training data (no variety) or only generates part of the distribution?
  - Mode collapse – Capturing only some modes of distribution

- Several theoretically-motivated regularization methods
  - Simple example: Add noise to real samples!

$$\lambda \cdot \mathbb{E}_{x \sim P_{real}, \delta \sim N_d(0, cI)} \left[ \| \nabla_{\mathbf{x}} D_\theta(x + \delta) \| - k \right]^2$$

*Kodali et al., On Convergence and Stability of GANs (also known as How to Train your DRAGAN)*

**Regularization**

Georgia Tech

# Generative Adversarial Nets: Convolutional Architectures

Samples
from the
model look
much
better!

Radford et al,
ICLR 2016

Georgia
Tech

# Generative Adversarial Nets: Convolutional Architectures

Interpolating
between
random
points in
latent space



Radford et al,
 ICLR 2016

Georgia
Tech

*Brock et al., Large Scale GAN Training for High Fidelity Natural Image Synthesis*

**Example Generated Images - BigGAN**

(a) 128×128          (b) 256×256          (c) 512×512          (d)

Figure 4: Samples from our model with truncation threshold 0.5 (a-c) and an example of class leakage in a partially trained model (d).

*Brock et al., Large Scale GAN Training for High Fidelity Natural Image Synthesis*

**Failure Examples - BigGAN**

Georgia Tech

https://www.youtube.com/watch?v=PCBTZh41Ris

**Video Generation**

- Generative Adversarial Networks (GANs) can produce amazing images!

- Several drawbacks
  - High-fidelity generation heavy to train
  - Training can be unstable
  - No explicit model for distribution

- Larger number of extensions:
  - GANs conditioned on labels or other information
  - Adversarial losses for other applications

**Summary**

# Comparison of Methods