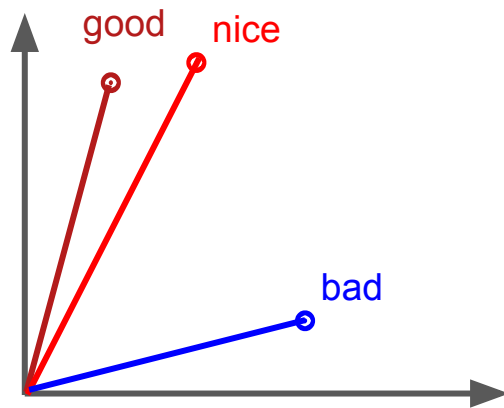


CS 4782 - Midterm Review

Snehal, Adhitya, Sean, Žiga, Lucas

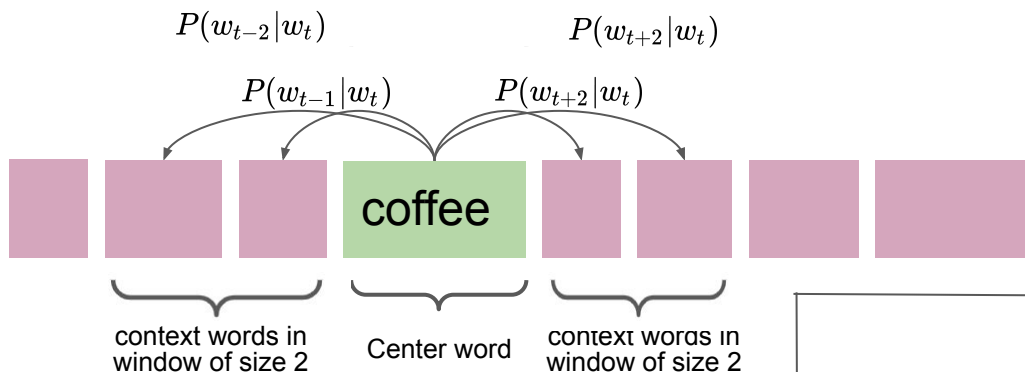
Word Embeddings

Motivated by semantic similarity

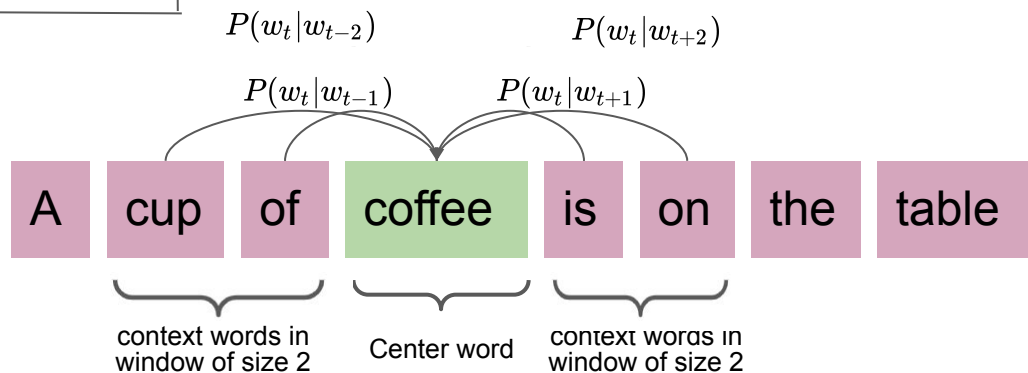


Word2Vec

Skipgram - Predict context from target



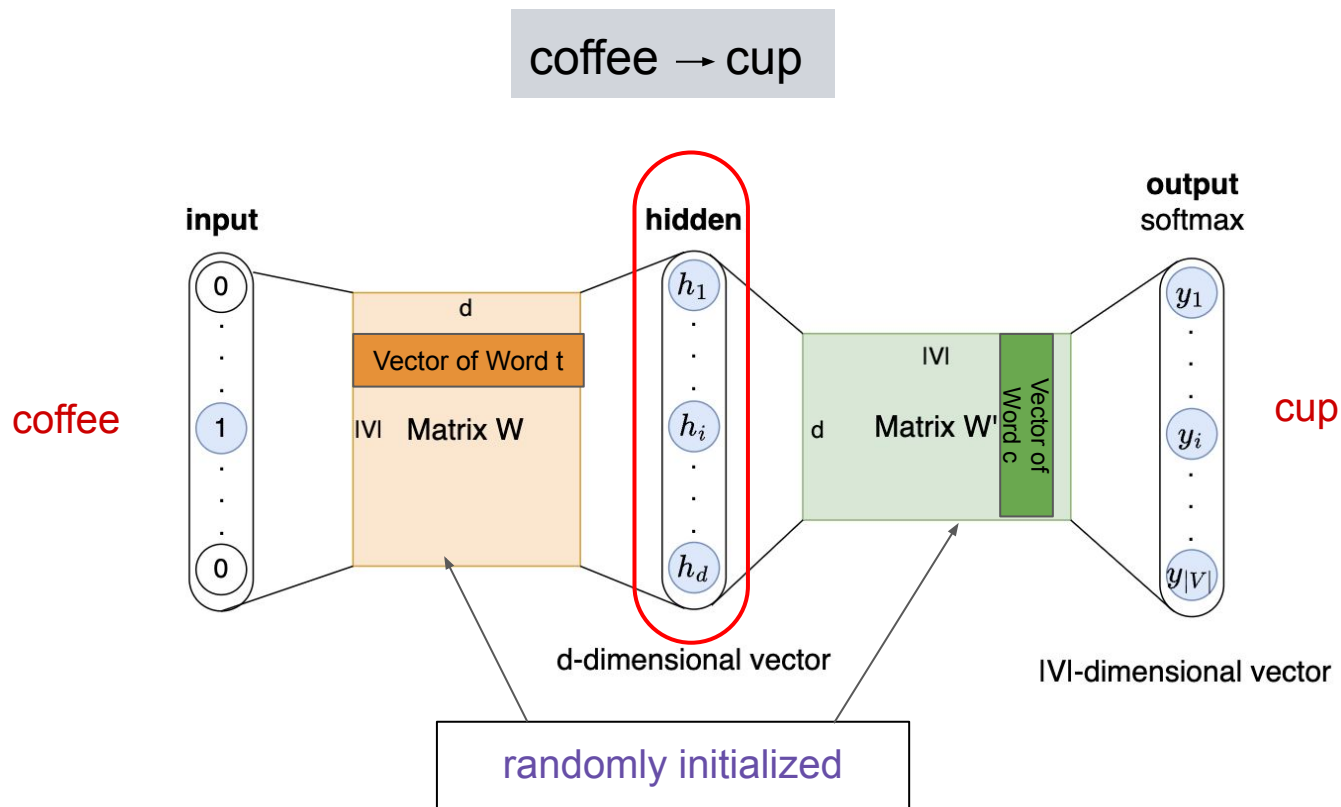
Continuous Bag of Words (CBOW) - predict target from context



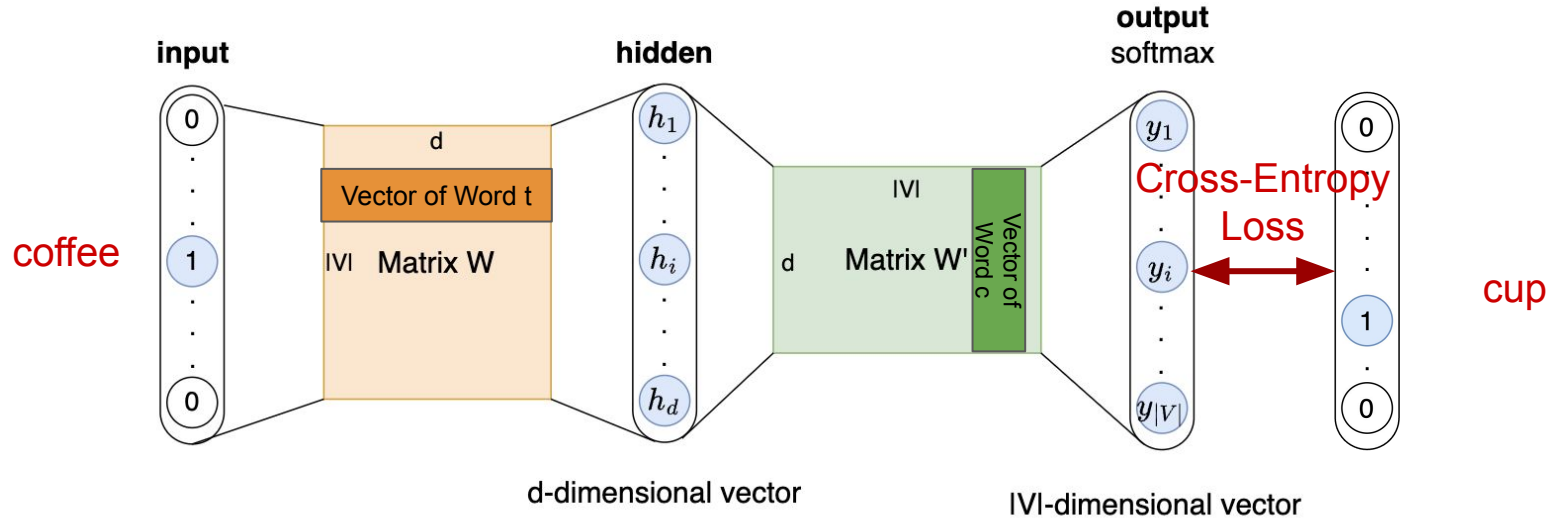
Training Data

One-hot encoded vectors of all the words in your vocabulary

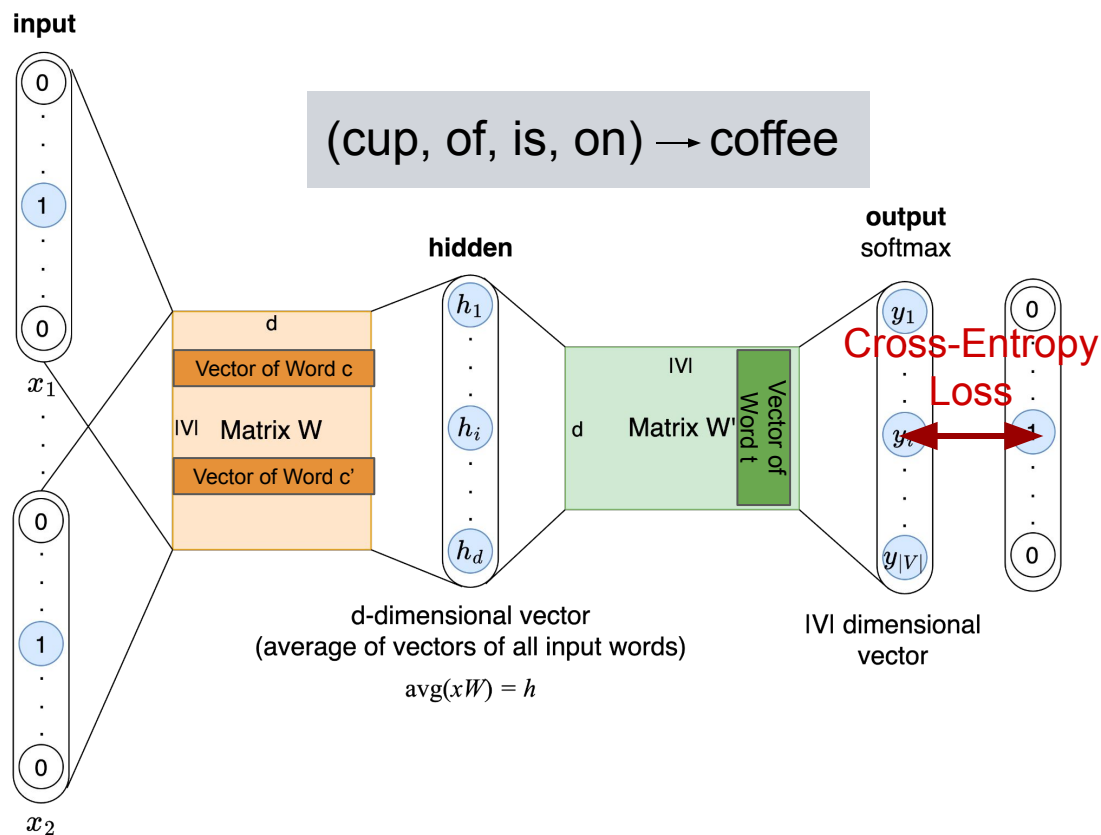
SkipGram



Calculate Loss and BackPropagate till Minima



CBOW

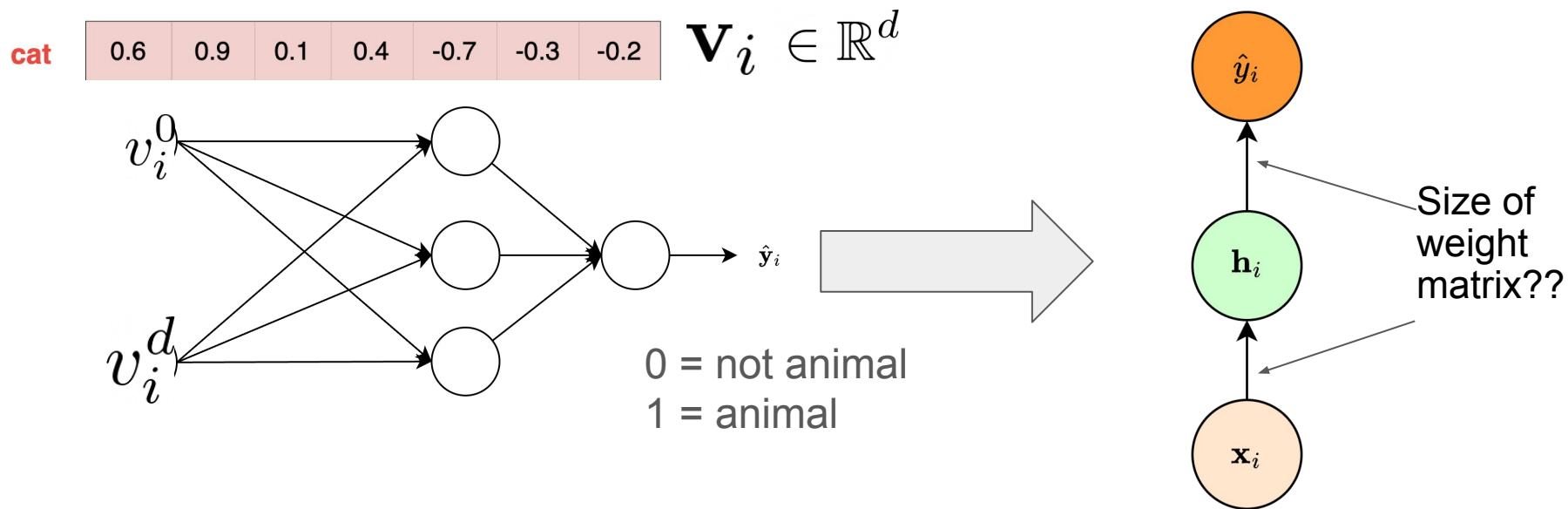


Calculate Loss
and
BackPropagate
till Minima

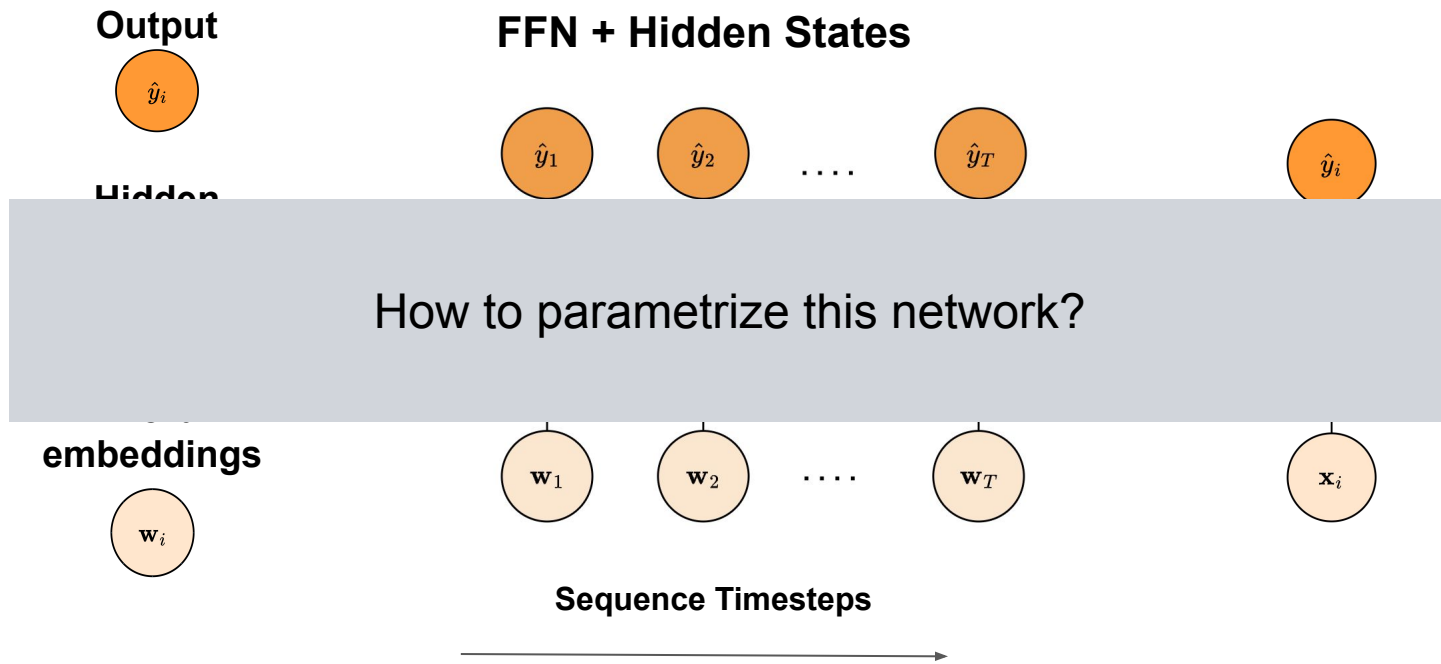
Transitioning to RNNs

Let's simplify!

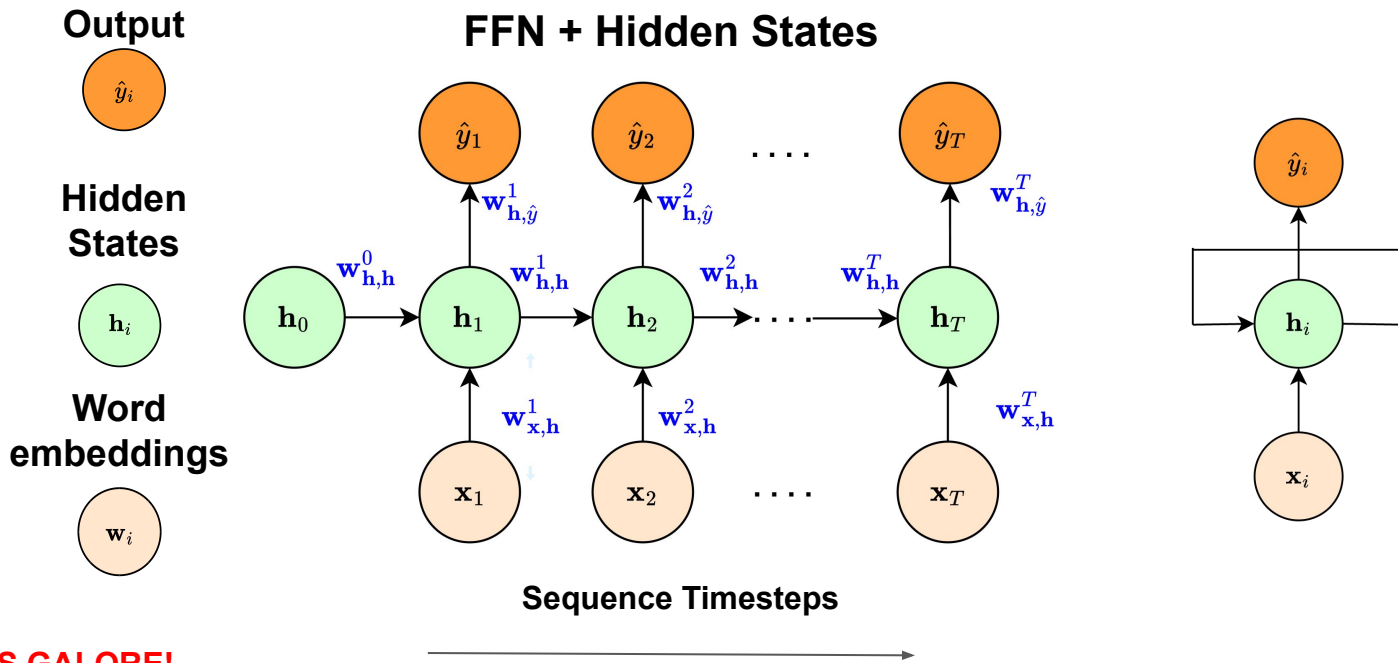
What if we have a single word and a single output?



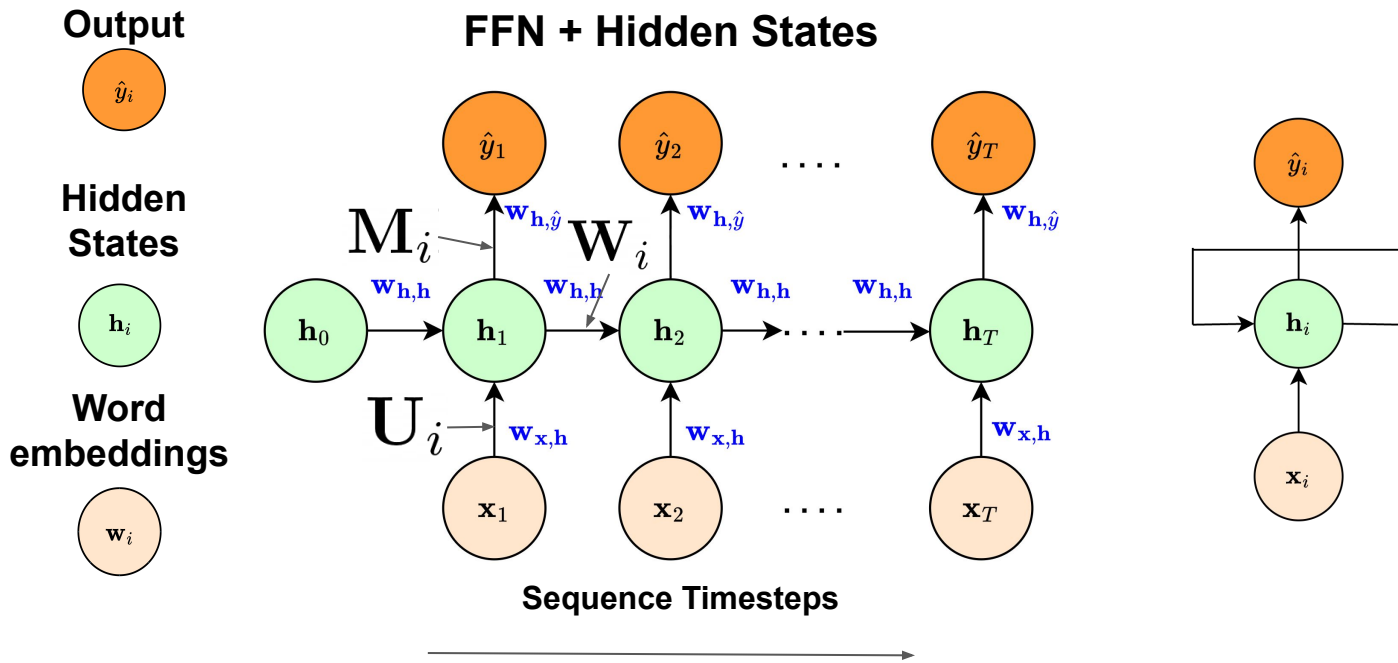
Recurrent neural network (RNN)



Maybe we add a weight matrix between every state??

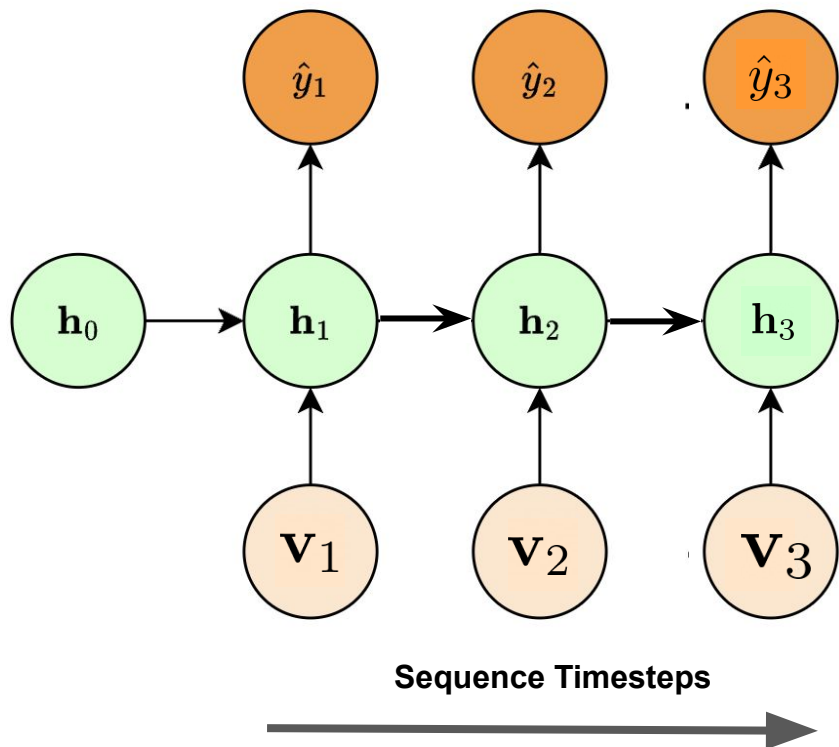


Solution: Reuse same weight matrices - *wherever possible*



Recurrent neural network (RNN)

Use the same parameters across different timesteps.



Output

$$\hat{y}_i = \mathbf{M} \mathbf{h}_i$$

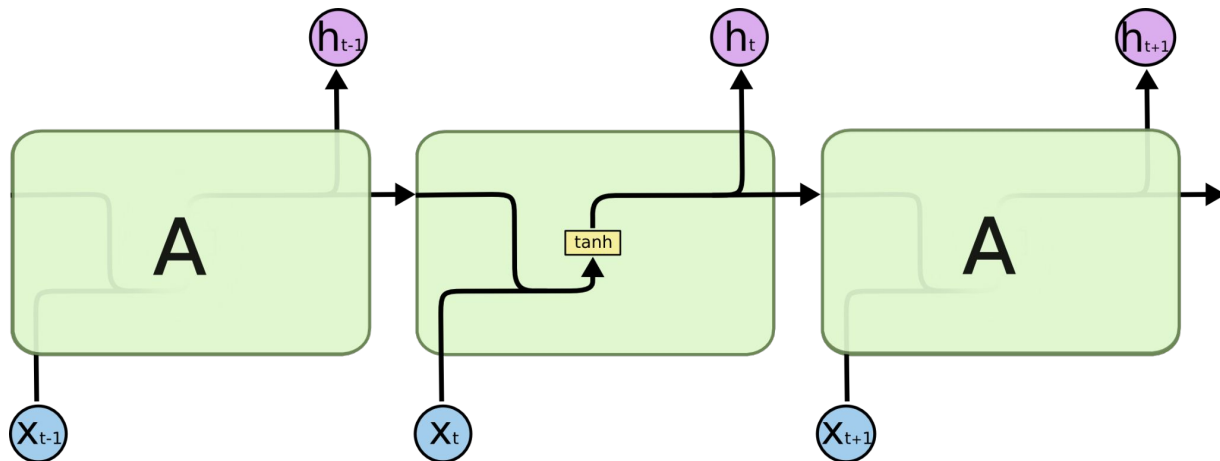
Hidden State

$$\mathbf{h}_i = \sigma(\mathbf{U} \mathbf{v}_i + \mathbf{W} \mathbf{h}_{i-1})$$

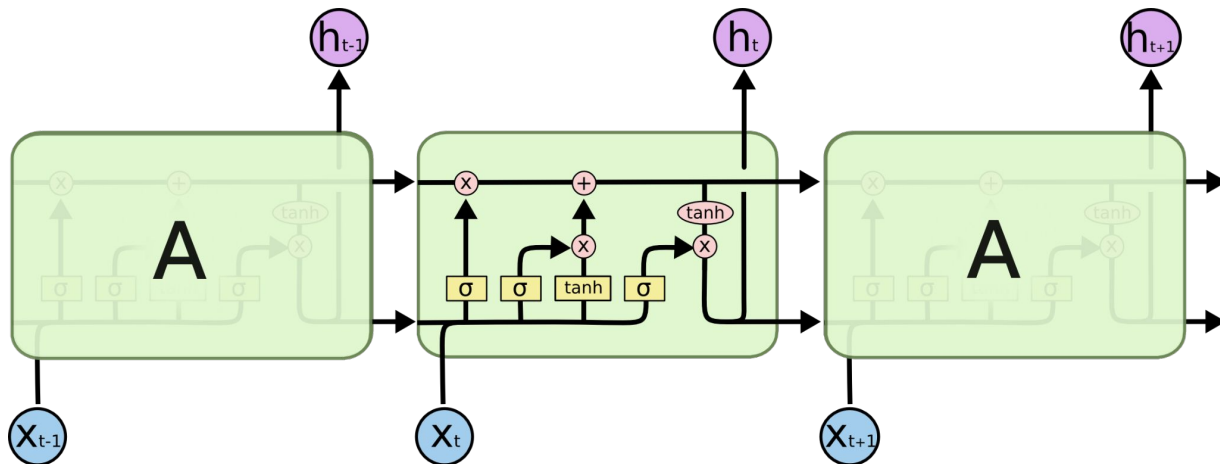
Problems

- **Looooooooong** Context Issues
- Sequential - slow

RNN



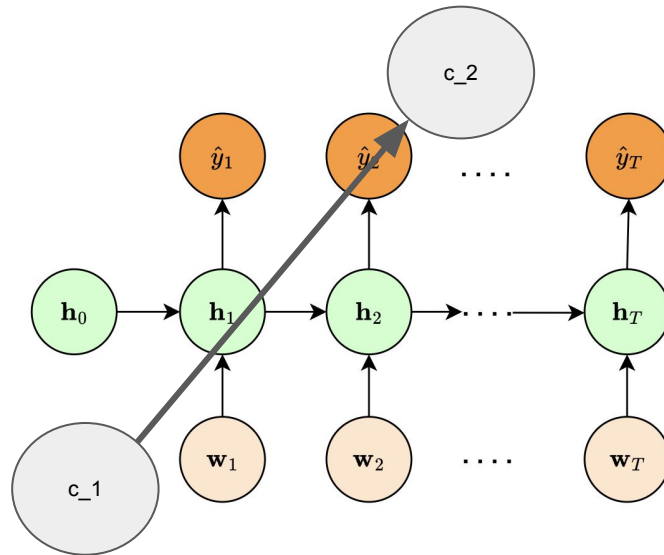
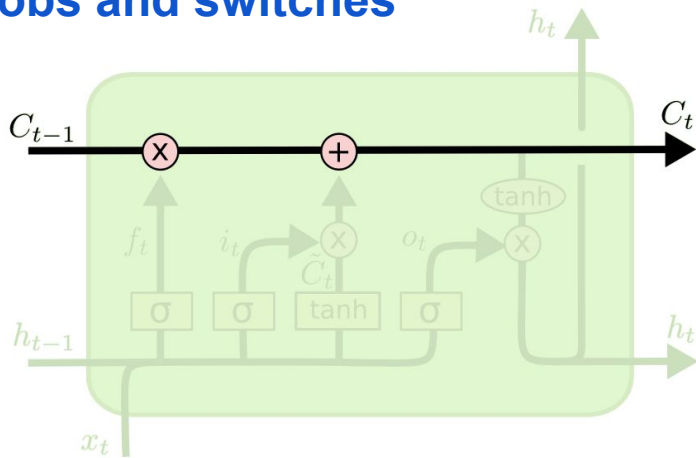
Long-short Term Memory (LSTM)



Long-short Term Memory (LSTM)

- Main idea: add a “cell” state that allows information to flow easily
 - Similar to residual connections
 - No repeated matrix multiplications!

Control flow by adding knobs and switches

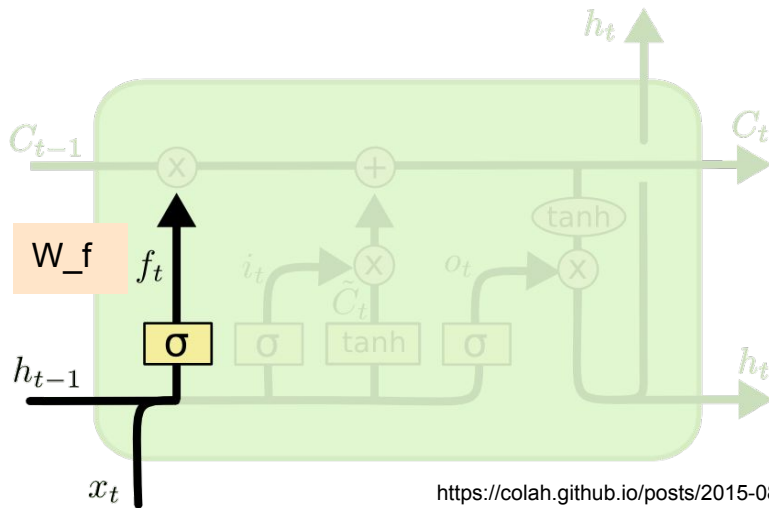


LSTMs- Forget Gate

- Forget gate- function of current input and previous hidden state
- Controls what should be remembered in the cell state

multiply

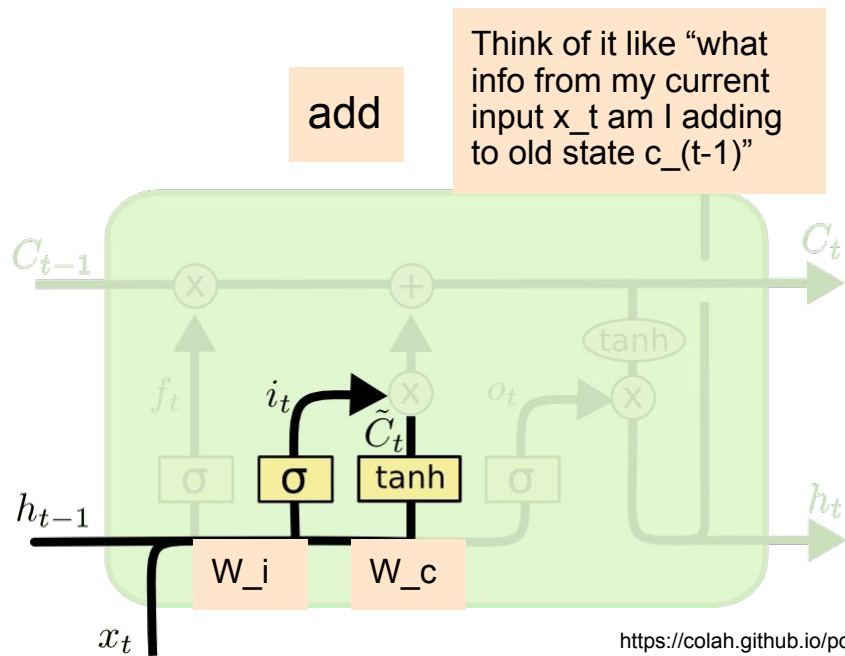
Think of it like blocking certain dimensions and letting others flow



$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

LSTMs- Input Gate

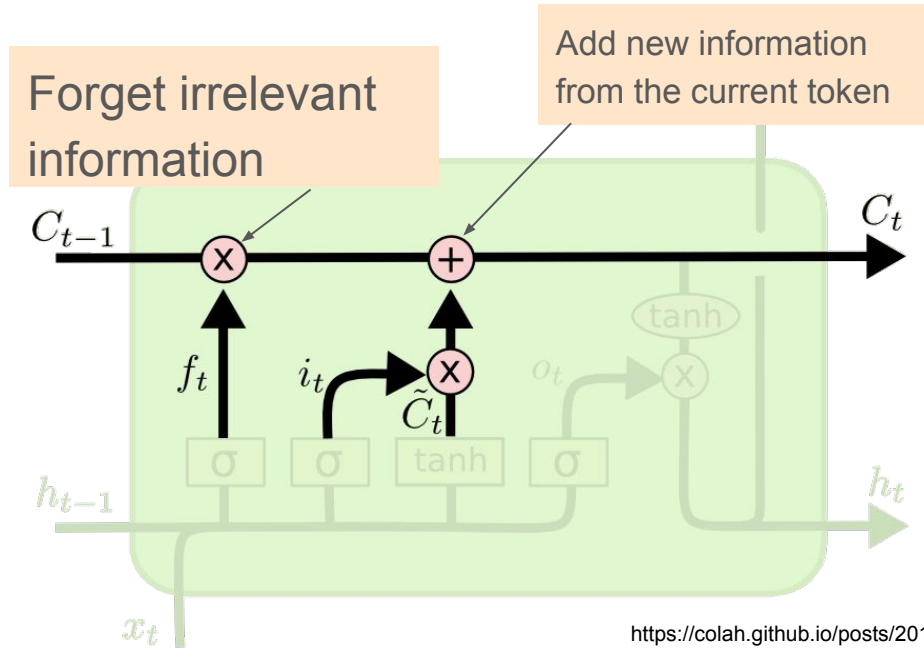
- Input gate- function of current input and previous hidden state
- Decides what information to write to the cell state



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

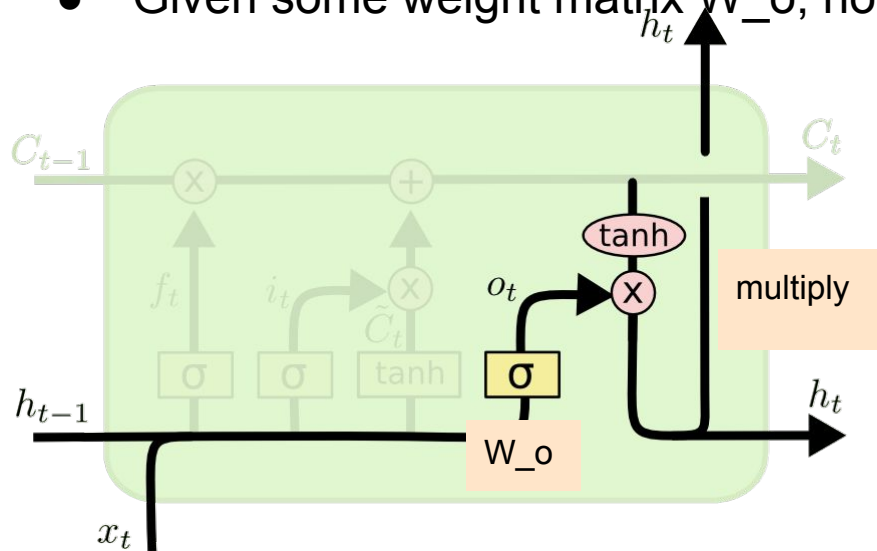
LSTM- Cell Update



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

LSTM- Output Gate

- Output gate- function of current input and previous hidden state
- Controls flow of information from the cell state to the hidden state
- Given some weight matrix W_o , how do we write to o_t and h_t ?



$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

Think of it like
“wanting to maintain
a latent hidden
space”

LSTMs

- Performs better with long sequences
- But still sequential!!

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

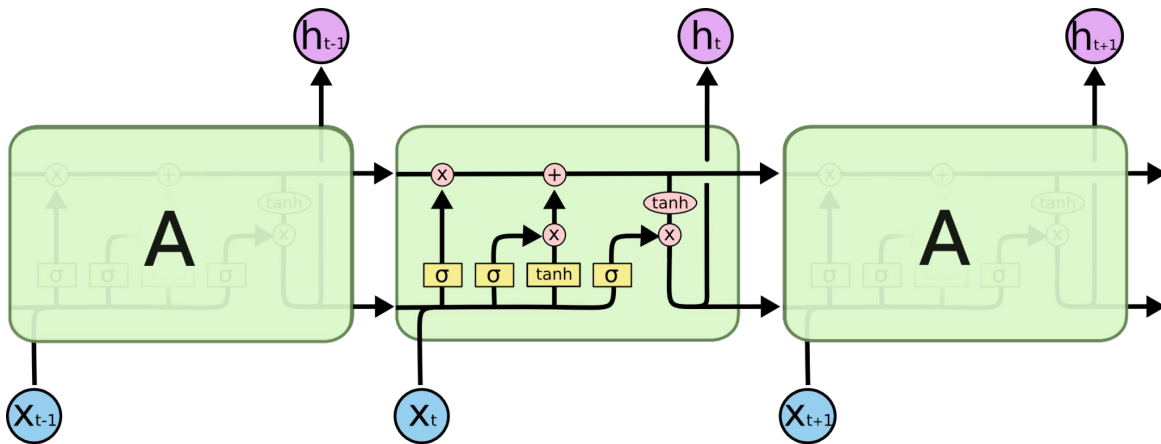
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

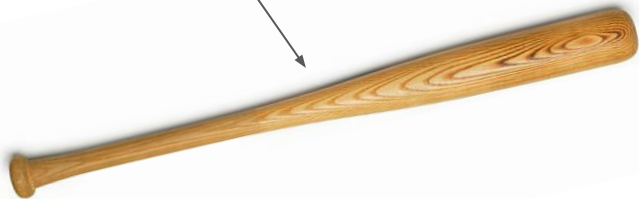
$$h_t = o_t * \tanh(C_t)$$



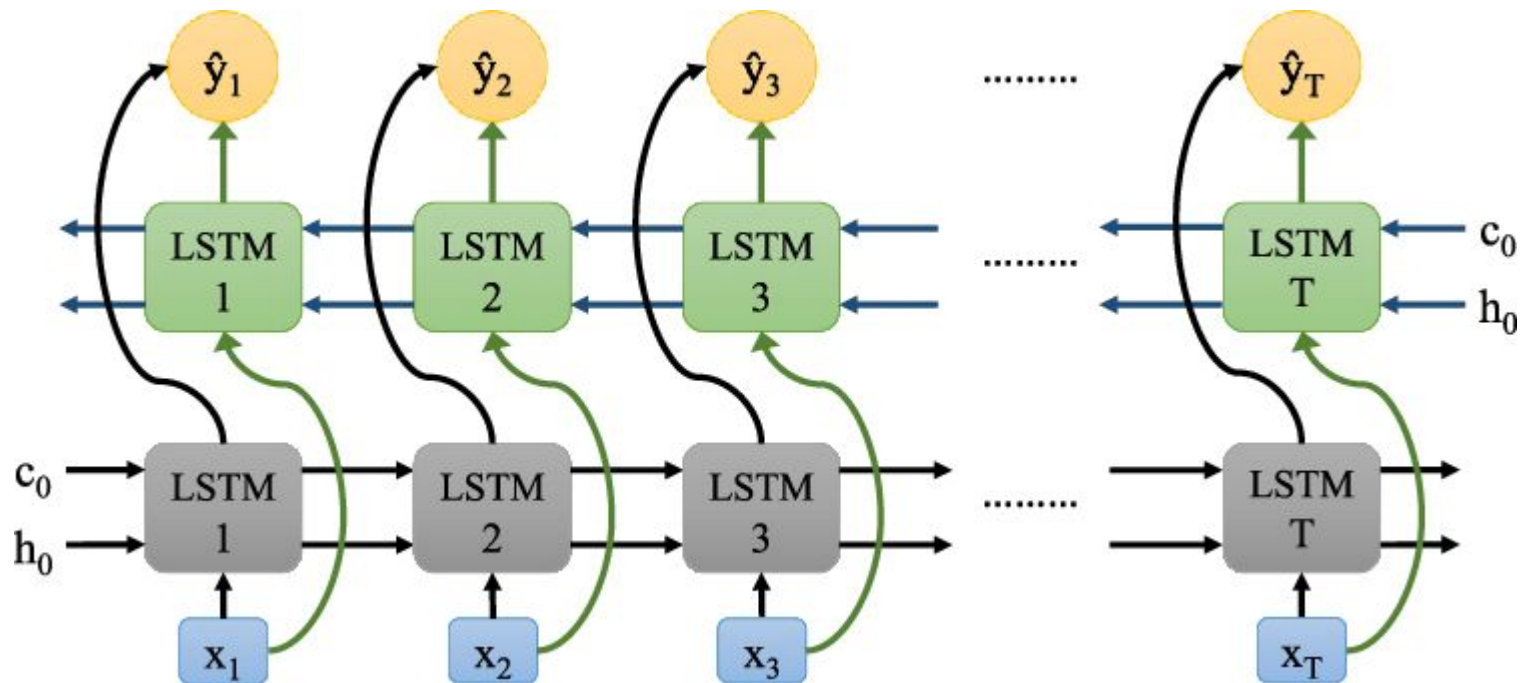
Context influences word meaning.



A **bat** flew out of the dugout, startling the baseball player and making him drop his **bat**.

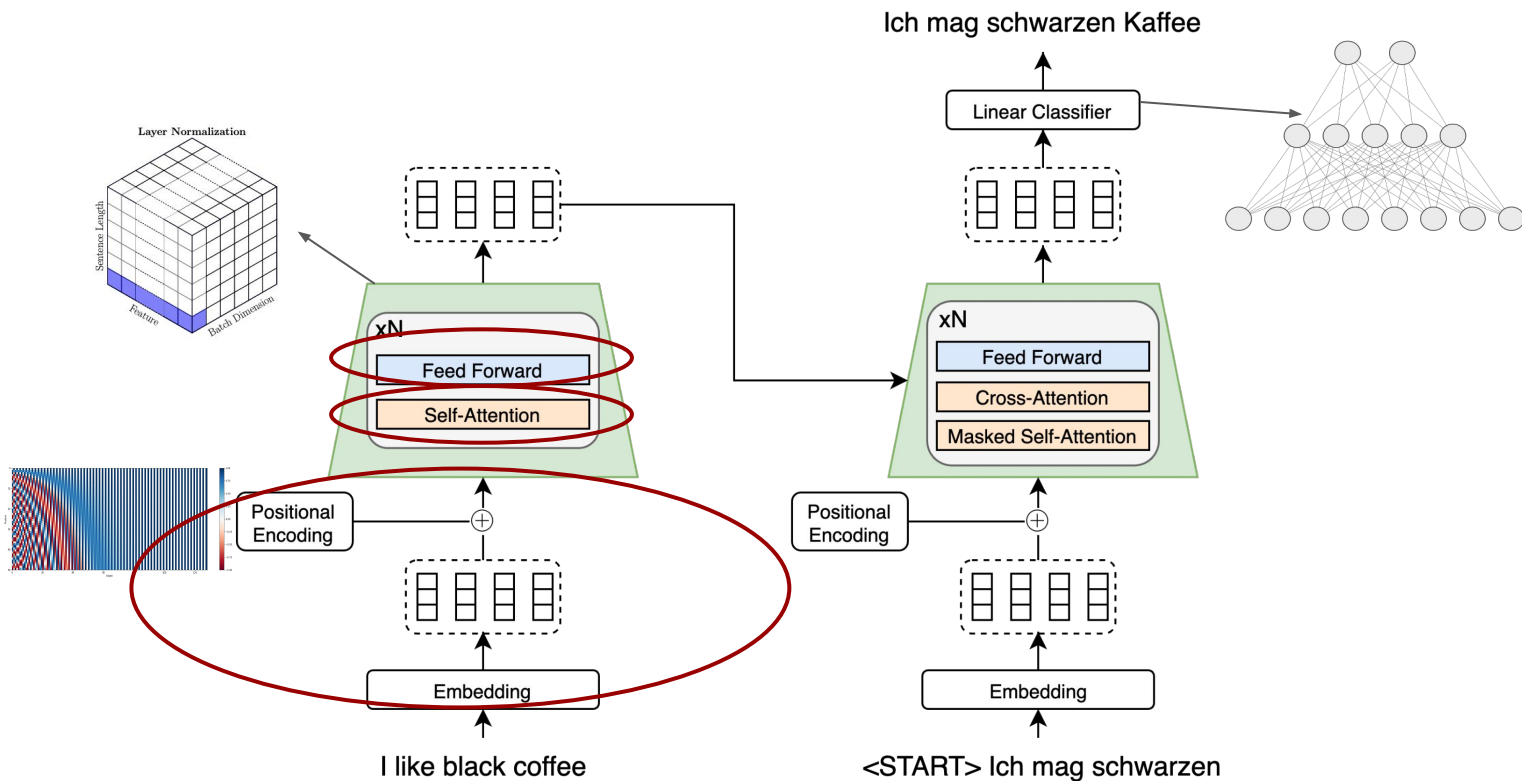


Bidirectional LSTM



Transformers

What we cover



RNNs and LSTMs were sequential - Transformers are parallelized!

Token Position

Integer	0	1	2	3	4	5	6	7
Binary	0	0	0	0	1	1	1	1
	0	0	1	1	0	0	1	1
	0	1	0	1	0	1	0	1

Positional Embeddings

For a position pos and embedding dimension i :

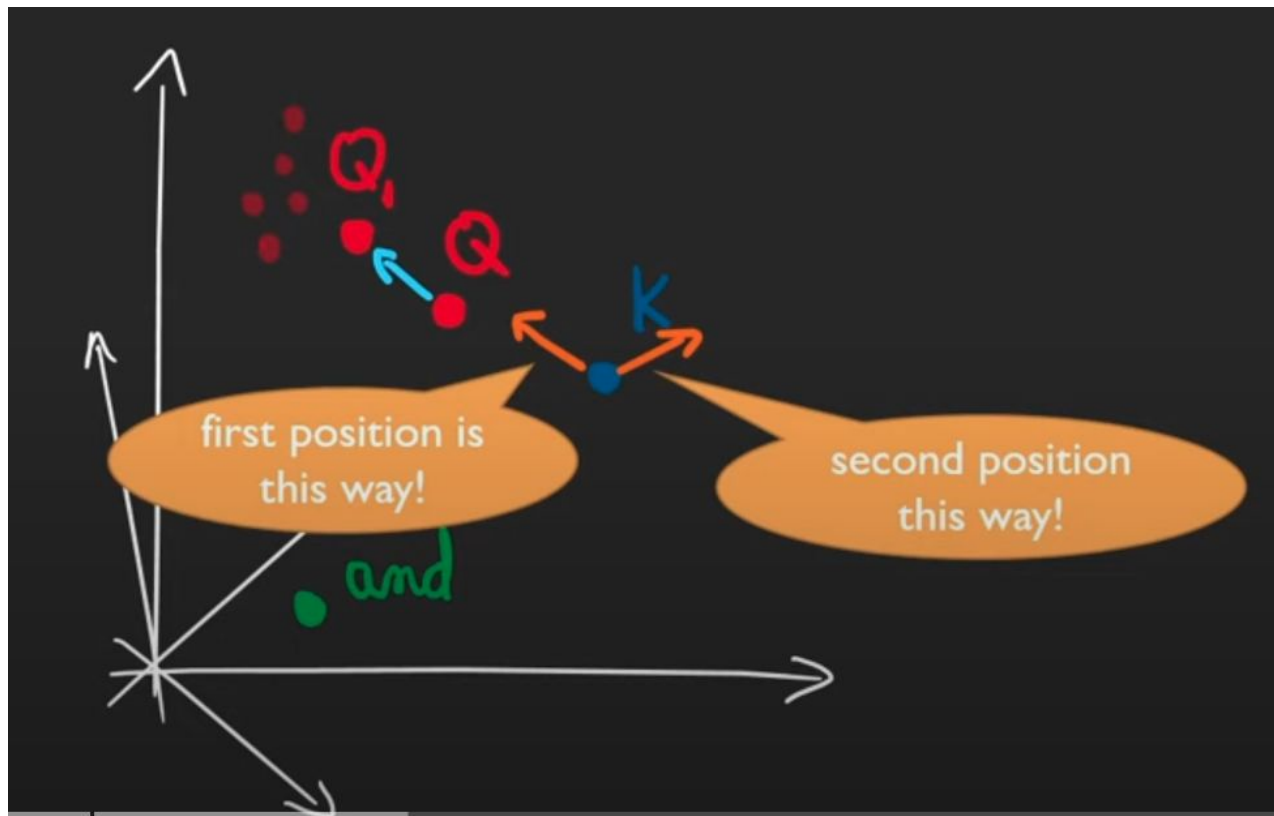
- $PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$
- $PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$

Token Position

Position	0	1	2	3	4	5	6	7
Sinusoid	0.5	0.3	0.3	0.5	0.7	0.9	0.9	0.7
	0.5	0.5	0.7	0.7	0.5	0.5	0.7	0.7
	0	1	0	1	0	1	0	1

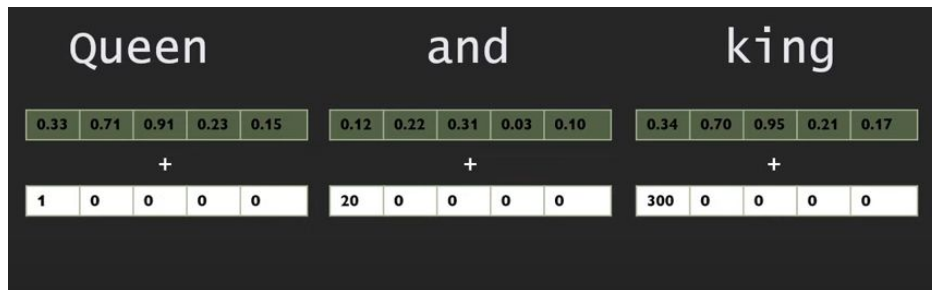
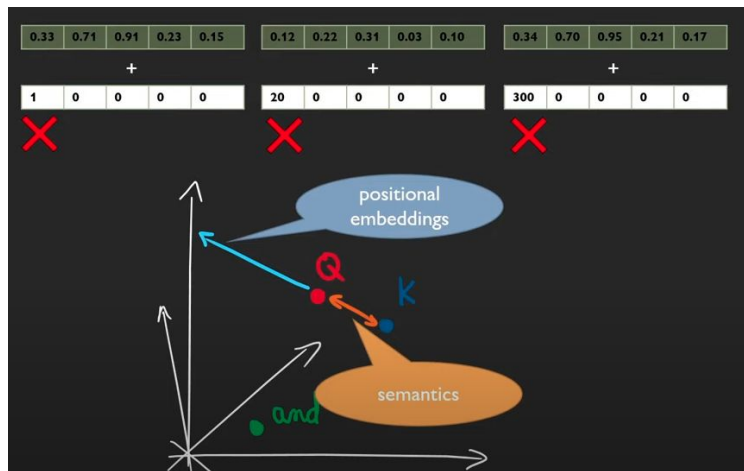
Positional Embeddings

Sentence: “Queen and King”



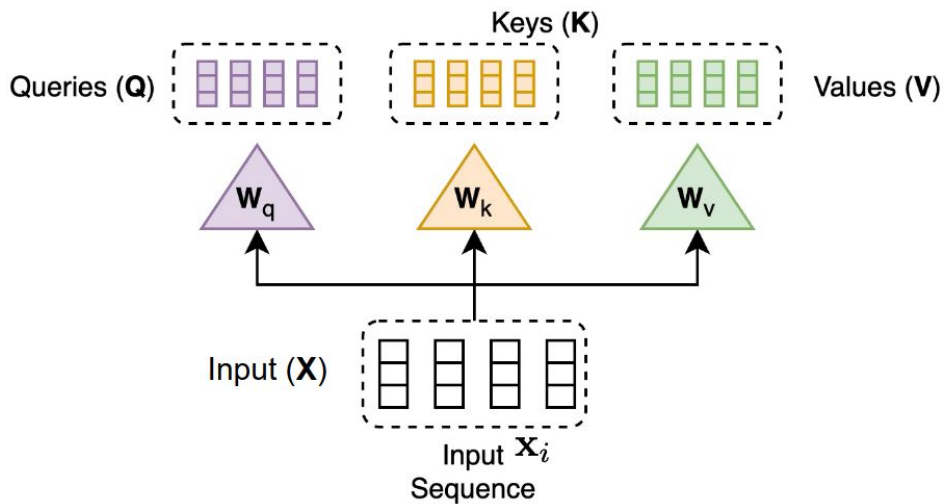
Some requirements

- Every position should have a unique identifier
- Independent of input
- Independent of sequence length
- Numbers shouldn't be too large



Self-Attention

$$\text{Attention}(\overbrace{Q, K, V}^{(N \times D) \text{ input Matrices}}) = \overbrace{\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V}^{(N \times D) \text{ output Matrix}}$$



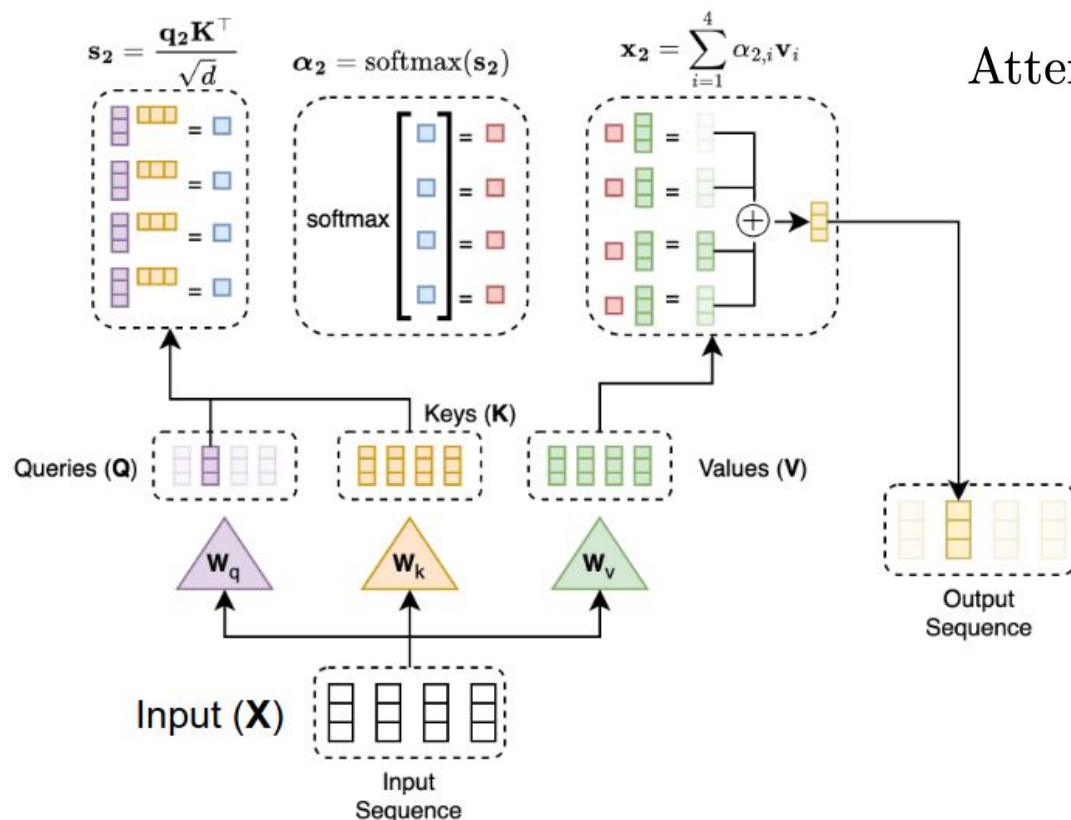
1. Get Q, K, V from the N input token using the weights.

$$w_q, w_k, w_v \quad (N,) \rightarrow (N, D) \begin{cases} Q \\ K \\ V \end{cases}$$

2. Use attention to transform this token representation with the other tokens in the sequence.

$$(N, D) \rightarrow (N, D) \{A\}$$

Self-Attention



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

1. Get Q, K, V from the N input token using the weights.

$$W_q, W_k, W_v \quad (N,) \rightarrow (N, D) \begin{cases} Q, \\ K, \\ V \end{cases}$$

2. Use attention to transform this token representation with the other tokens in the sequence.

$$(N, D) \rightarrow (N, D) \{A\}$$

Self-Attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Implementation:

$$\text{MultiHead}(\mathbf{X}) = \mathbf{W}_O \text{Concat}(\text{head}_1, \dots, \text{head}_h)$$

$$\text{where head}_i = \text{Attention}(\mathbf{W}_Q^i \mathbf{X}, \mathbf{W}_K^i \mathbf{X}, \mathbf{W}_V^i \mathbf{X})$$

Multi-Head-Attention(input x):

Split input into query, key, and value vectors

q = split_heads($\mathbf{W}_q(x)$) # (b, h, n, d)

k = split_heads($\mathbf{W}_k(x)$) # (b, h, n, d)

v = split_heads($\mathbf{W}_v(x)$) # (b, h, n, d)

$$\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i$$

$$\mathbf{k}_i = \mathbf{W}_k \mathbf{x}_i$$

$$\mathbf{v}_i = \mathbf{W}_v \mathbf{x}_i$$

Compute attention scores and apply them to values

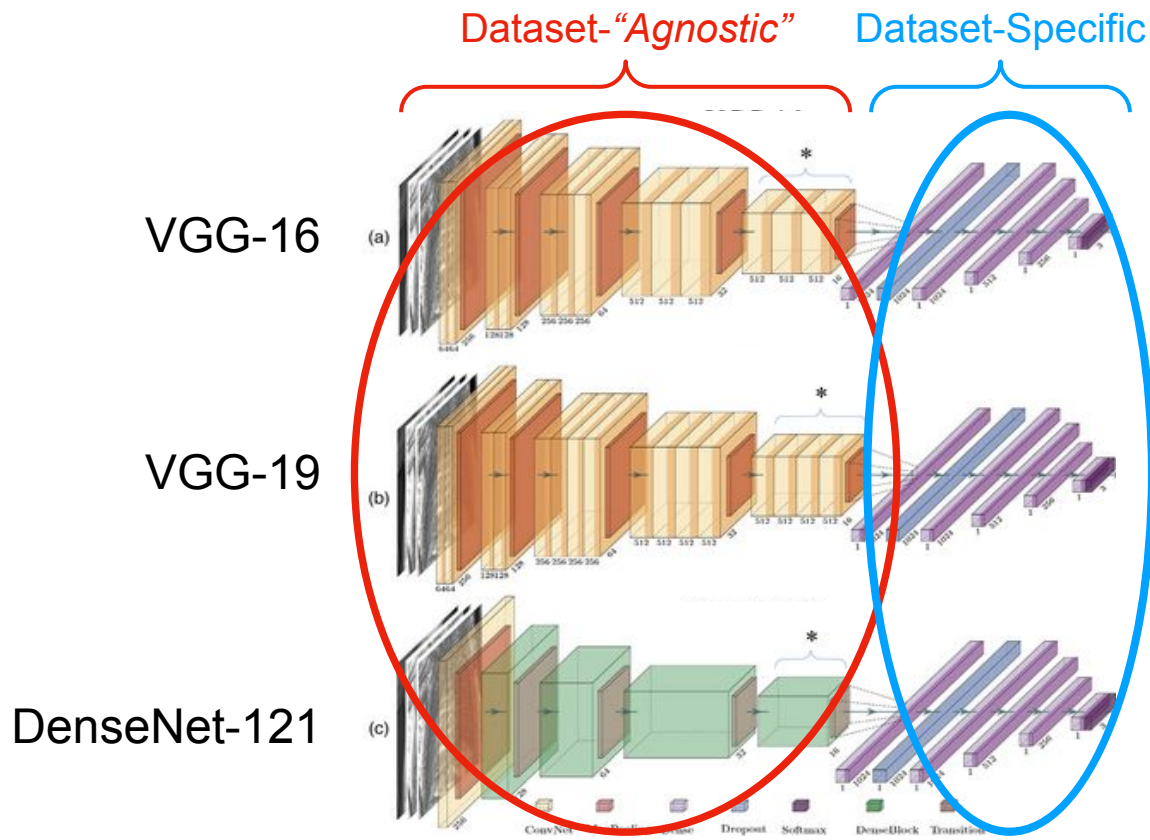
attn_output = compute_attention(q, k, v)

Combine attention heads and apply output transformation

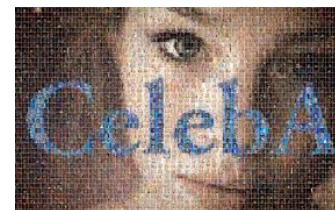
output = $\mathbf{W}_o(\text{combine_heads}(\text{attn_output}))$

return output

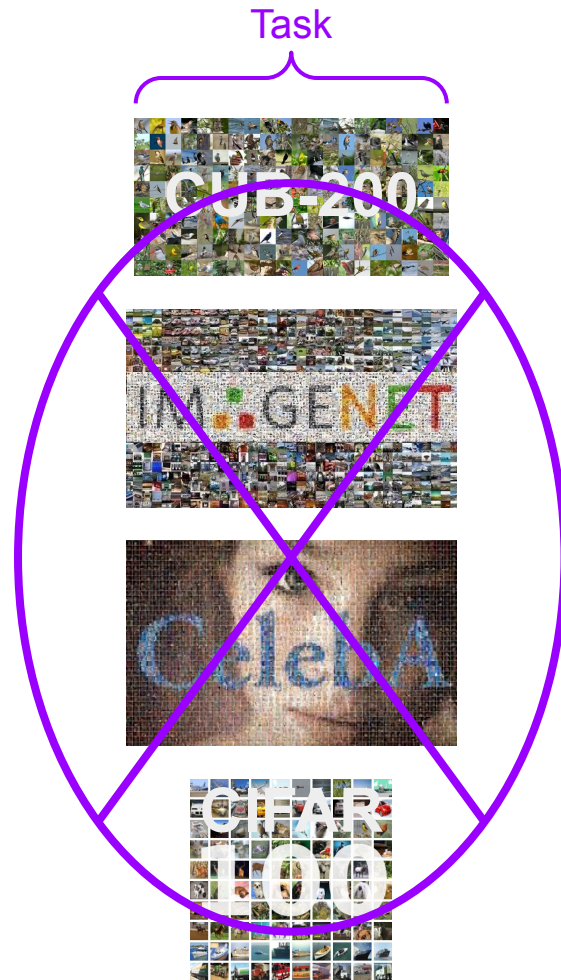
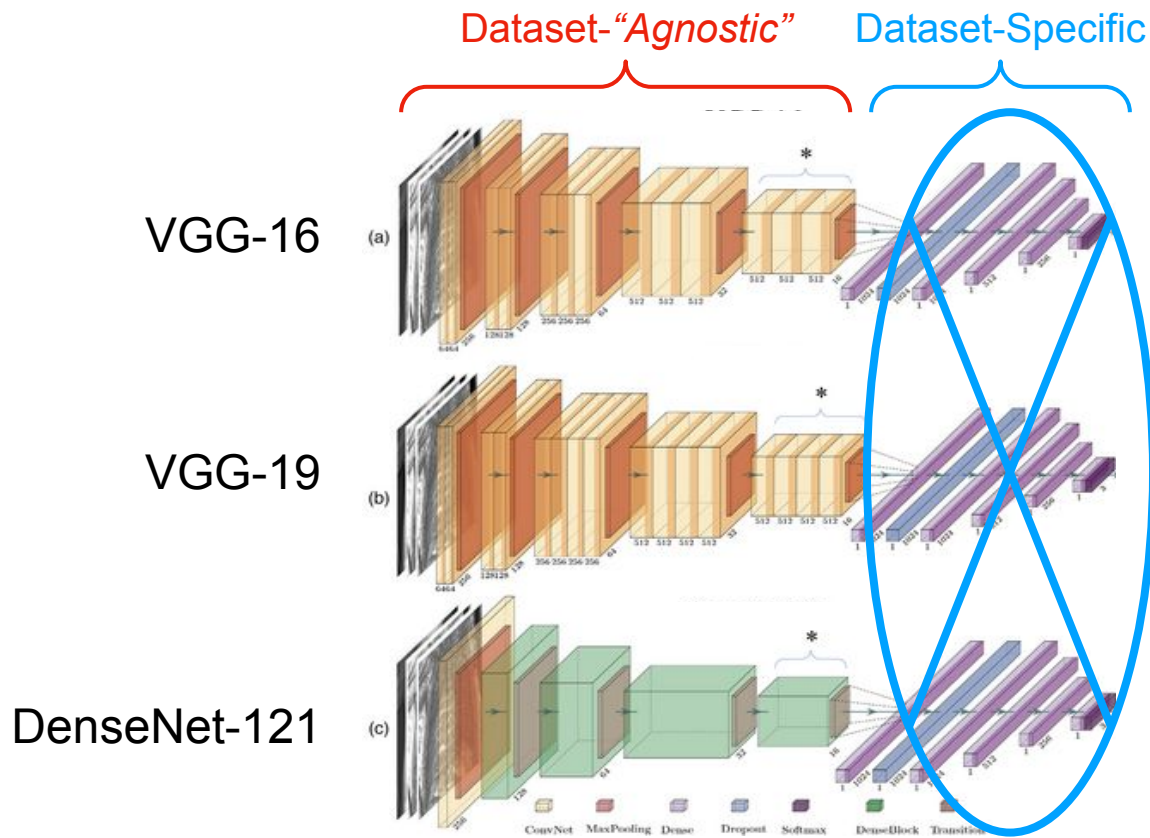
Representation Learning



Task



Representation Learning



Representation Learning

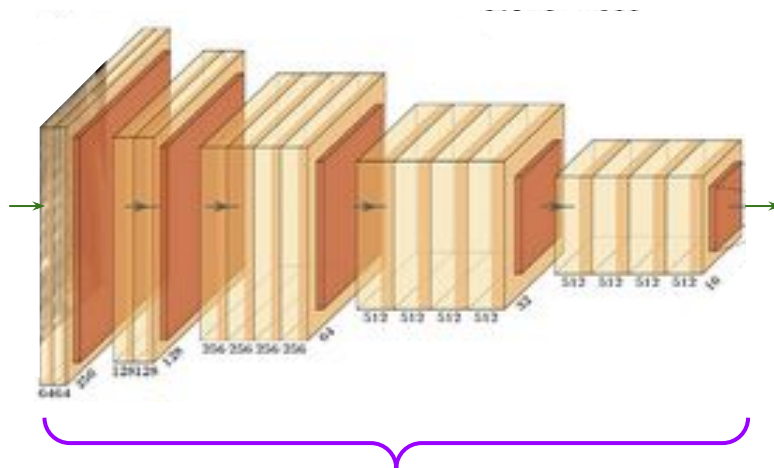
Unlabeled Image Data

$$\mathbb{R}^{224 \times 224}$$

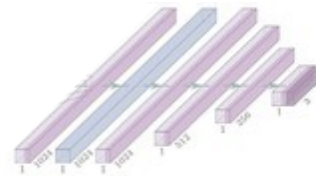


Image Representation

$$f(x) : \mathbb{R}^{224 \times 224} \rightarrow \mathbb{R}^{768}$$



How to get this useful
image embedder?



B.ring
Y.our
O.wn
C.lassifier

Representation Learning

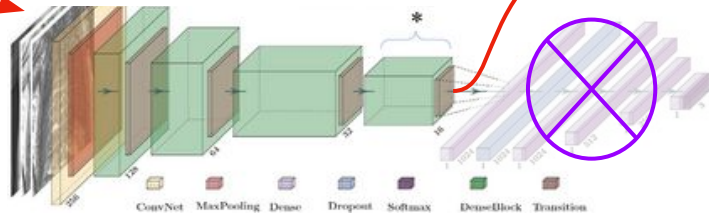


Q: How to get an image embedder?

Train on ImageNet and throw away the classifier.

→ Supervised Model Features

Neural Net Features: Nearest Neighbors



Representation Learning

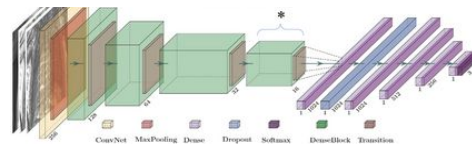
Q: How to get an image embedder?

→ Supervised Model Features

→ Self-Supervised Learning

◆ Pre-training Task

Pick a Task Correlated with Image Understanding
(Make it easy to label!)



- ☐ 0°
- ☒ 90°
- ☐ 180°
- ☐ 270°

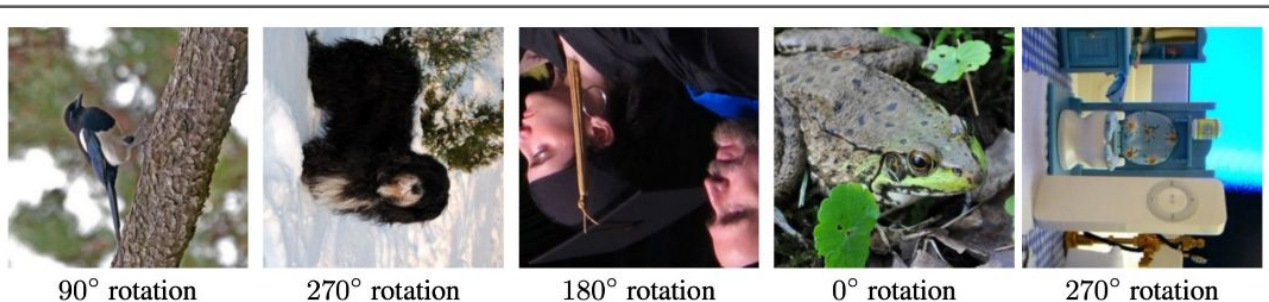


Figure 1: Images rotated by random multiples of 90 degrees (e.g., 0, 90, 180, or 270 degrees). The core intuition of our self-supervised feature learning approach is that if someone is not aware of the concepts of the objects depicted in the images, he cannot recognize the rotation that was applied to them.

Representation Learning

Q: How to get an image embedder?

→ Supervised Model Features

→ Self-Supervised Learning

◆ Pre-training Task

◆ Contrastive Learning

\mathbf{X} : Sample

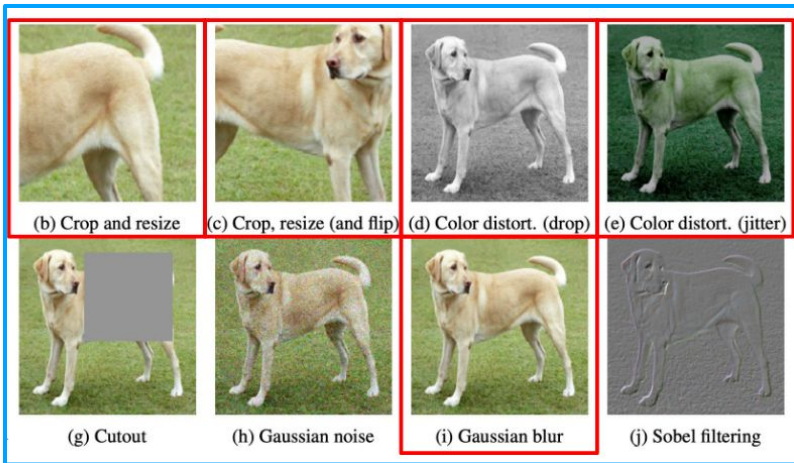


(a) Original

(+negative augmentations)
 \mathbf{X}^- : **Negative** examples (*different images*)



\mathbf{X}^+ : **Positive** examples (*augmentations*)



Representation Learning

Q: How to get an image embedder?

→ Supervised Model Features

→ Self-Supervised Learning

◆ Pre-training Task

◆ Contrastive Learning

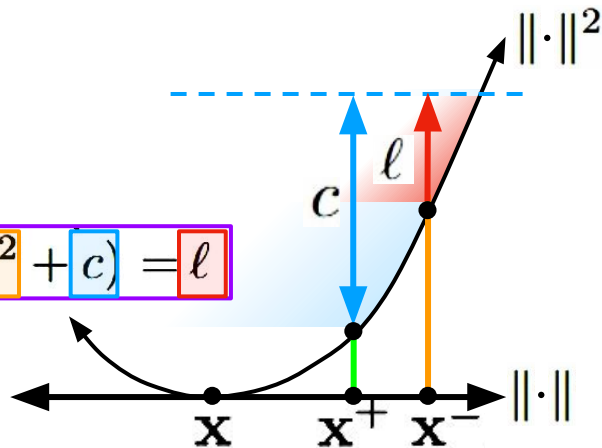
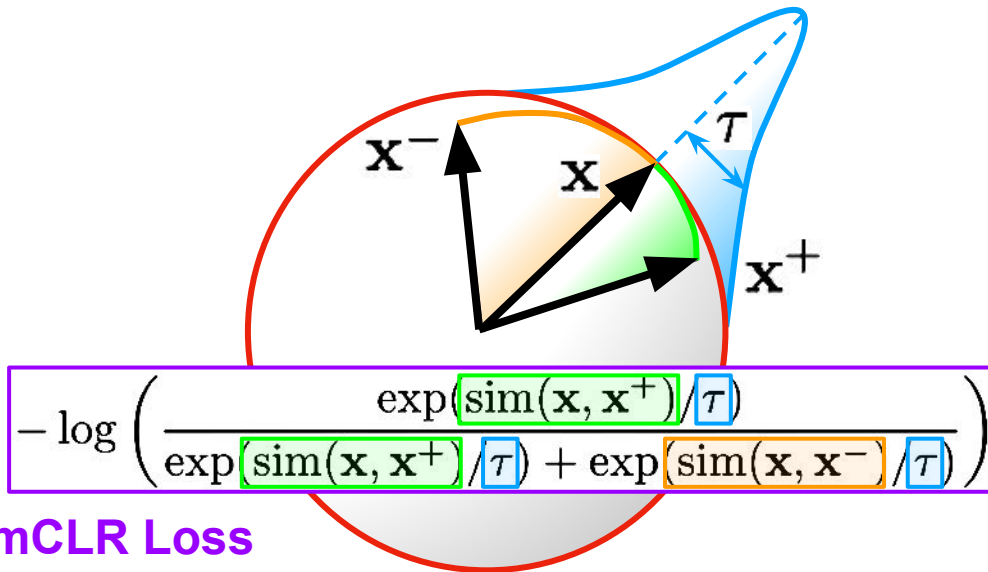
\mathbf{x} : Sample

\mathbf{x}^- : Negative examples

\mathbf{x}^+ : Positive examples

$$\max(0, \|f(\mathbf{x}_i) - f(\mathbf{x}^+)\|^2 - \|f(\mathbf{x}_i) - f(\mathbf{x}^-)\|^2 + c) = \ell$$

Triplet Loss



Representation Learning (Multimodal)

a **multimodal**

Q: How to get ~~an image~~ **embedder**?

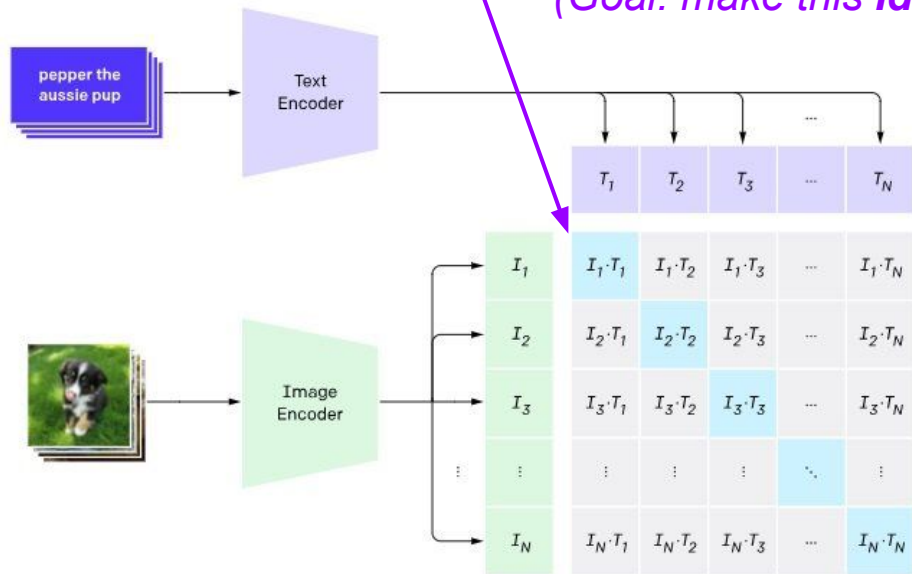
→ Supervised Model Features

→ Self-Supervised Learning

◆ Pre-training Task

◆ **Contrastive Learning**

1. Contrastive pre-training



Representation Learning (Multimodal)

a multimodal

Q: How to get ~~an image~~ embedder?

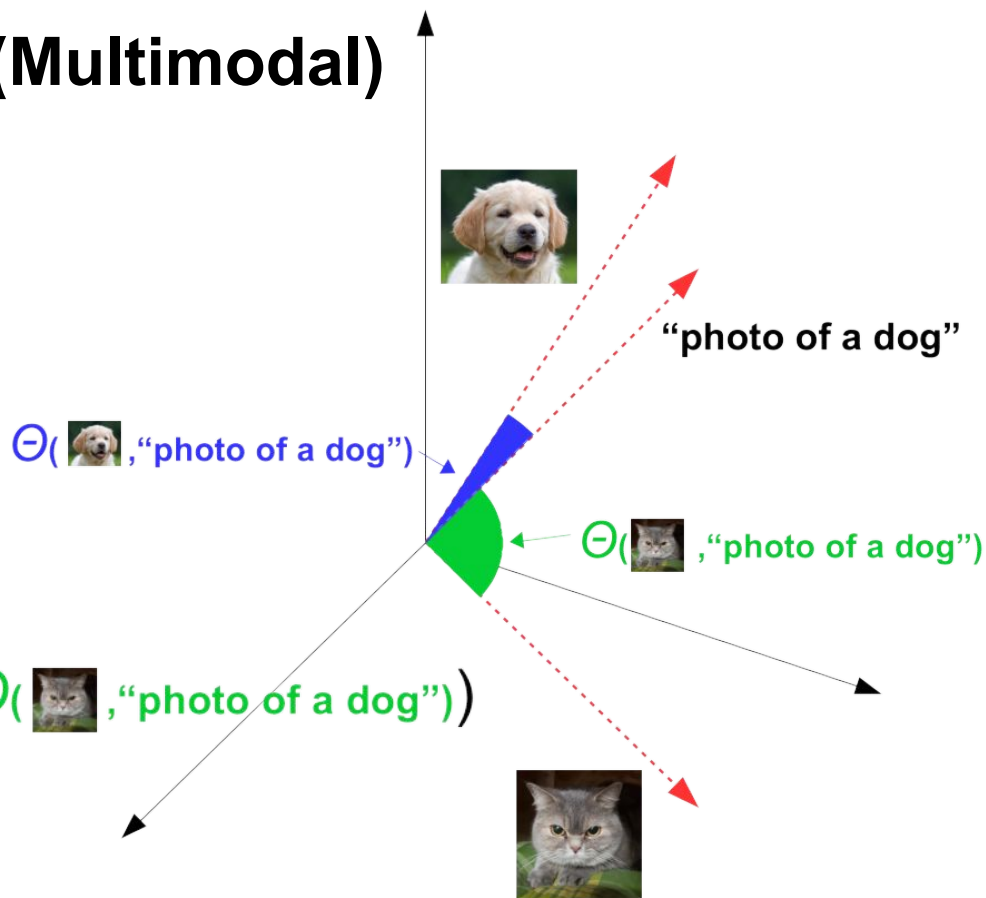
→ Supervised Model Features

→ Self-Supervised Learning

◆ Pre-training Task

◆ Contrastive Learning

$$\cos(\Theta(\text{dog_img}, \text{"photo of a dog"})) > \cos(\Theta(\text{cat_img}, \text{"photo of a dog"}))$$



Representation Learning

Q: How to get an image embedder?

- Supervised Model Features
- Self-Supervised Learning
 - ◆ Pre-training Task
 - ◆ Contrastive Learning
 - ◆ **Teacher-Student**

What is this an image of?



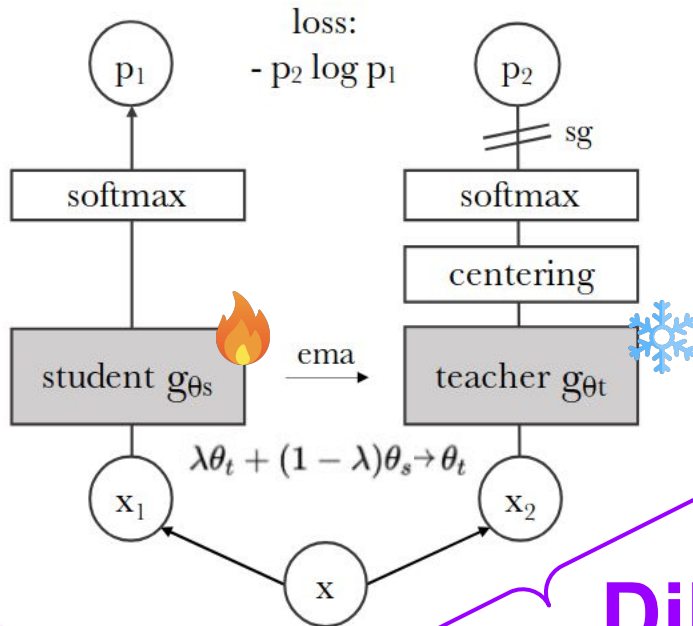
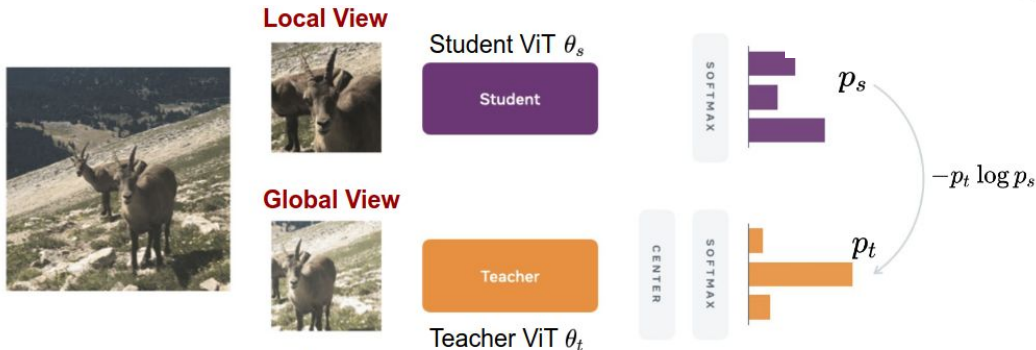
Representation Learning

Q: How to get an image embedder?

→ Supervised Model Features

→ Self-Supervised Learning

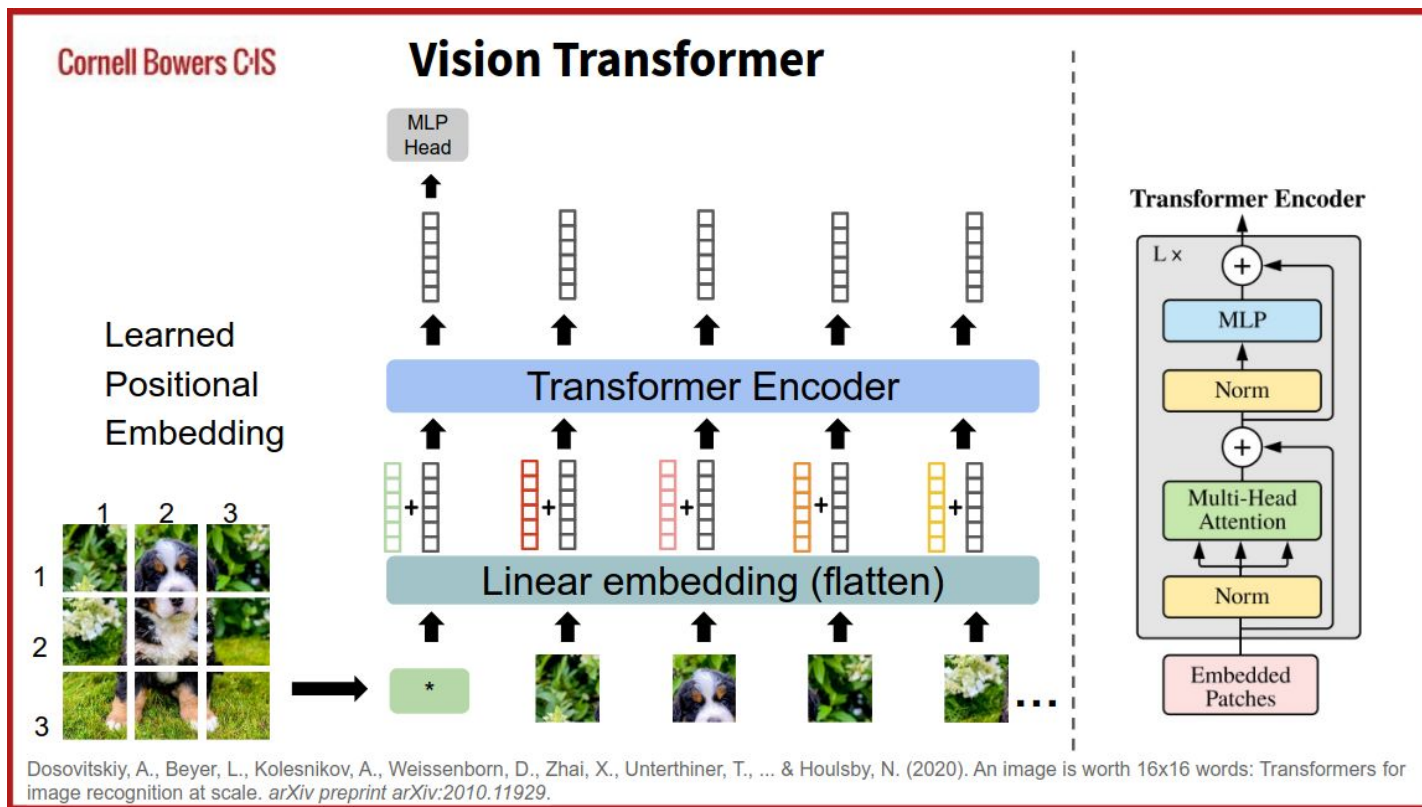
- ◆ Pre-training Task
- ◆ Contrastive Learning
- ◆ **Teacher-Student**



DiNO
“Distillation No Labels”

Vision Transformers (ViT)

Frequently pretrained with DiNO, you'll probably see "*DiNO-ViTs*" a lot.

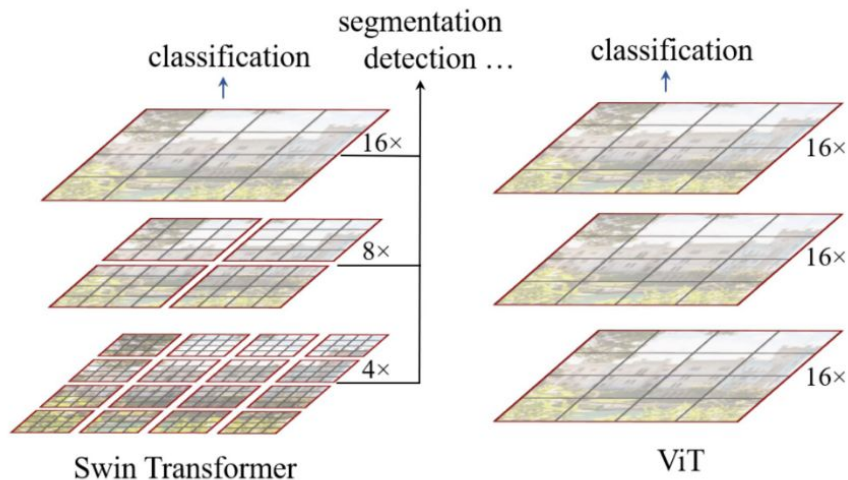


Vision Transformers (ViT)

SWIN “Scaled Window”

- Popular for large images and segmentation because of hierarchical features.
- Attention is per-window.
 - Shifted windows allow long-range connections at higher layers.

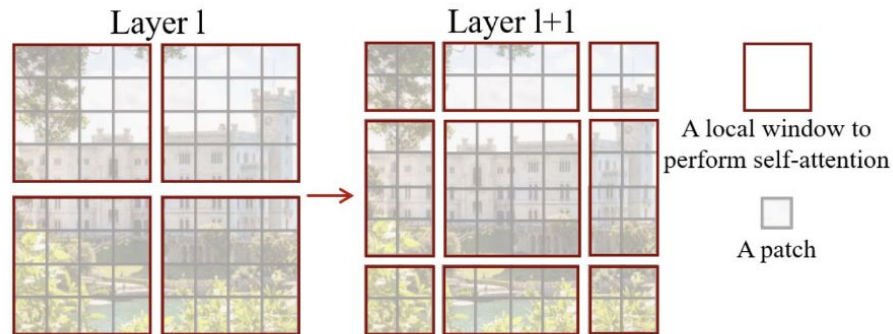
Cornell Bowers CIS



Cornell Bowers CIS

Shifted Window attention

Linear computational complexity with respect to image size



Vision Transformers (ViT)

DiNO v2

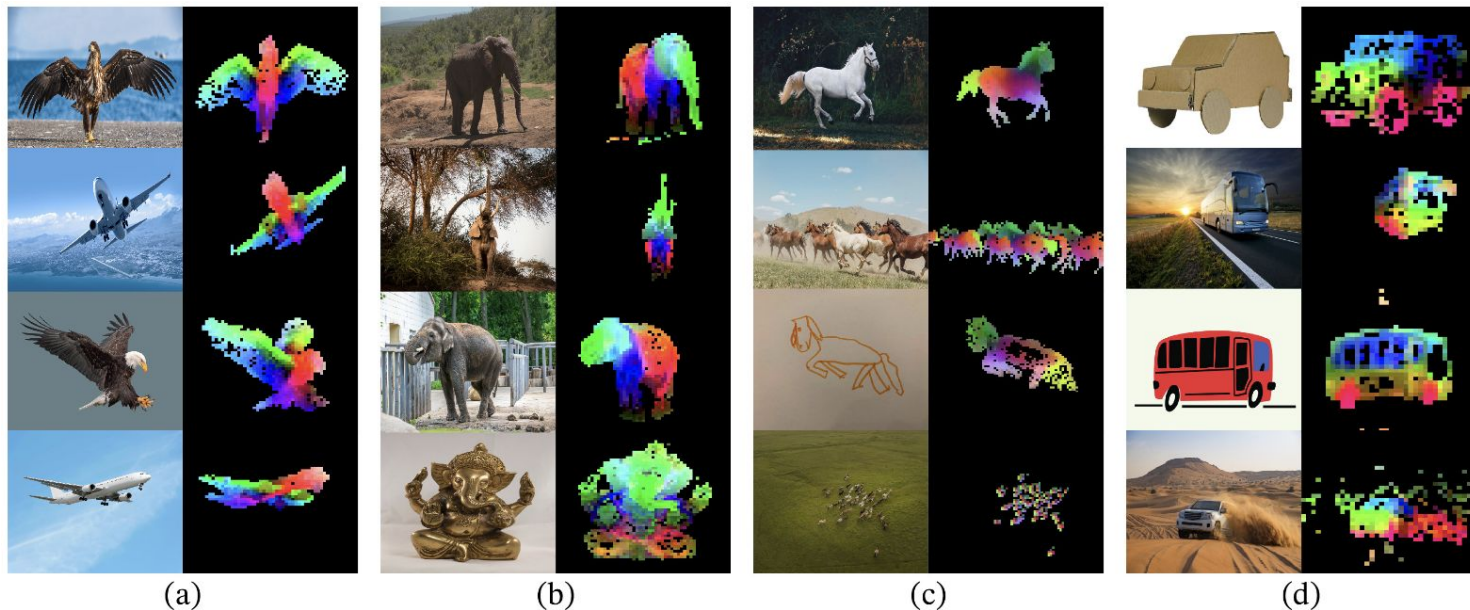
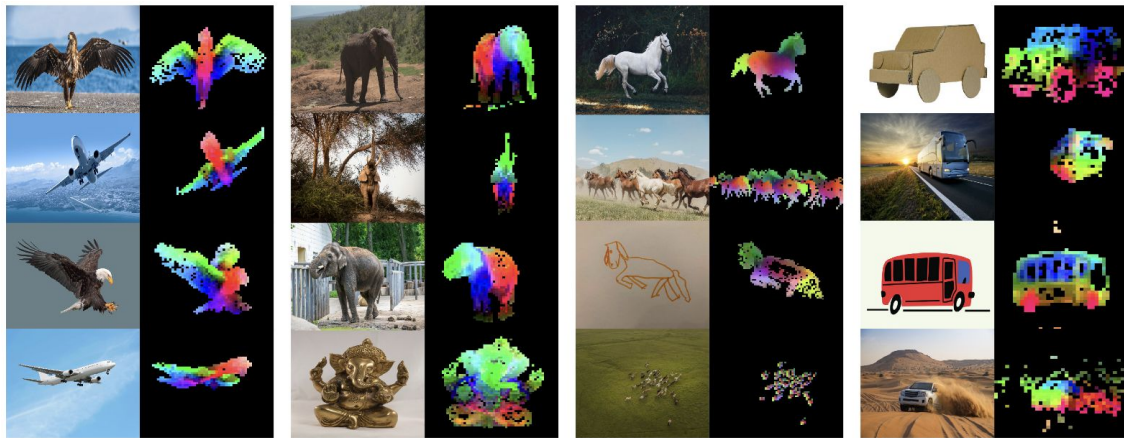


Figure 1: **Visualization of the first PCA components.** We compute a PCA between the patches of the images from the same column (a, b, c and d) and show their first 3 components. Each component is matched to a different color channel. Same parts are matched between related images despite changes of pose, style or even objects. Background is removed by thresholding the first PCA component.

Vision Transformers (ViT)

Figure 1: Visualization of the first PCA components of the image patches.

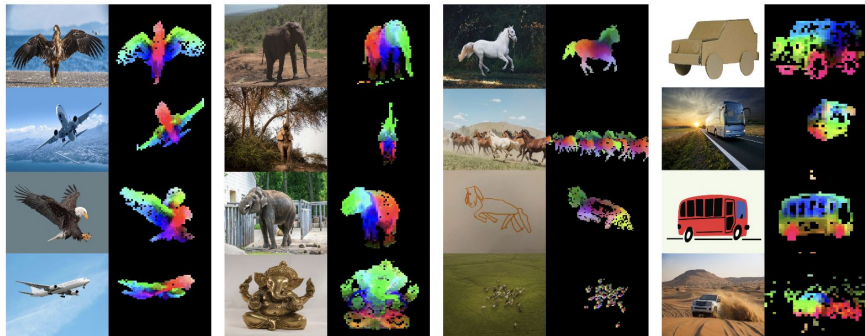


DiNO v2

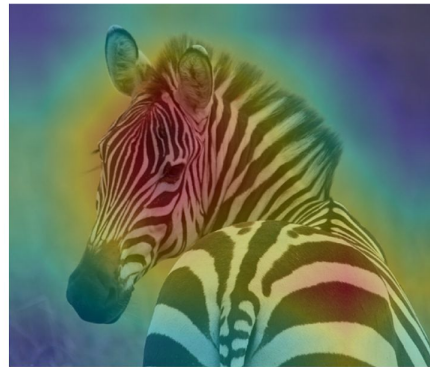
Q: Why don't we get these visualizations from deep CNNs?

(NOTE: You can train ResNet-50 with a DiNO, and it works great!)

Vision Transformers (ViT) *(vs deep CNNs)*



class: zebra



class: convertible



ViTs start with a global receptive fields. CNNs take a while to combine global features.

- Interpretable long-range dependencies from attention are visible early on.
- DiNO's objective encourages memorizing “parts-of-whole”.

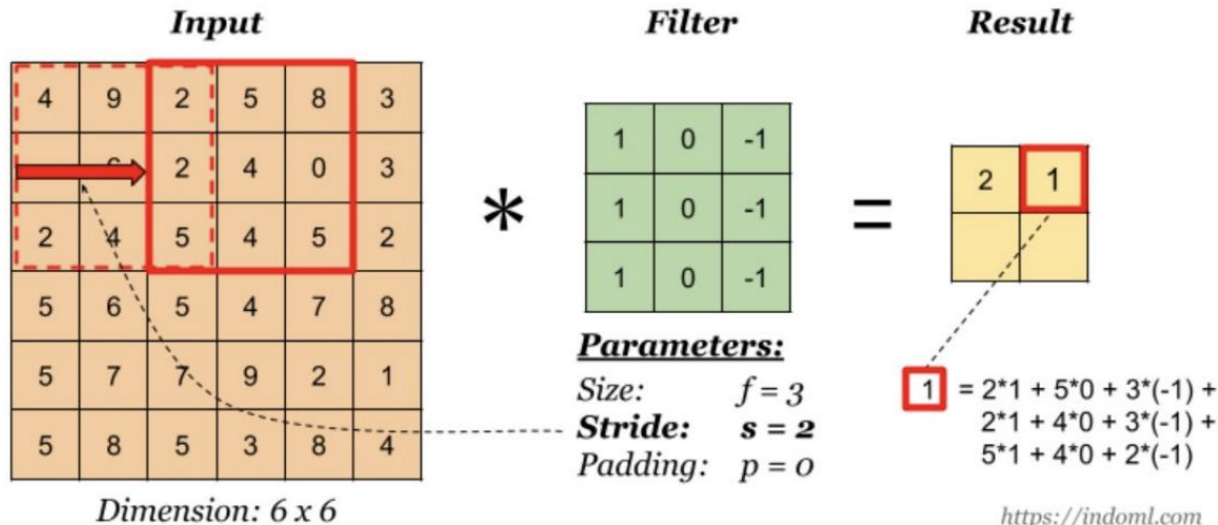
- Creates high level features, but it is hard to determine which pixels they correspond to.
- Must use methods like **Class Activation Mapping** up the network to get to pixel-level effects.

CNN review

HW1 Q4

Output dimension:

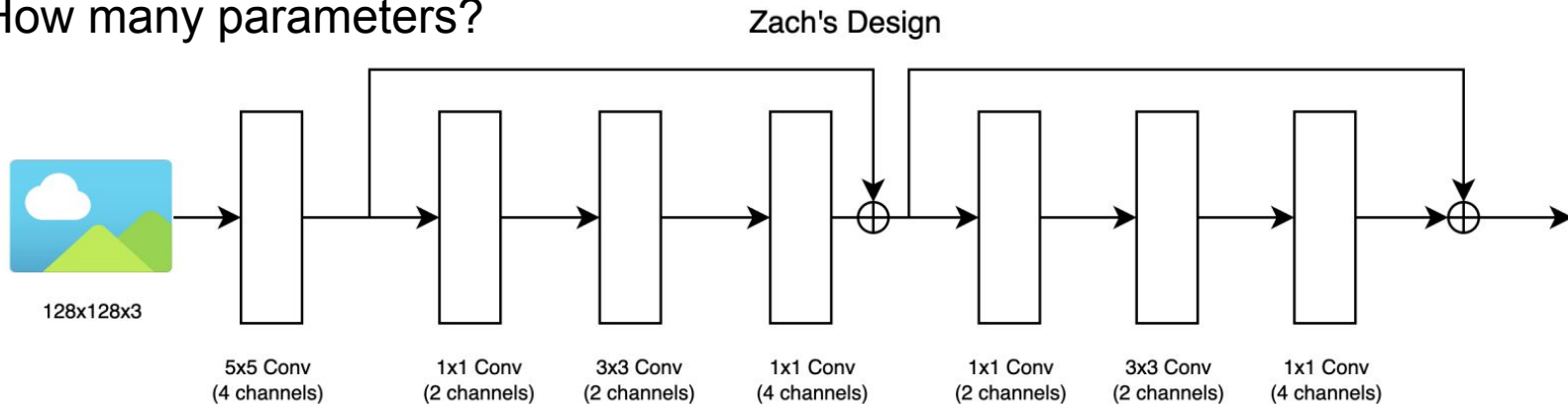
$$(\lfloor \frac{n - k + 2p}{s} \rfloor + 1) \times (\lfloor \frac{n - k + 2p}{s} \rfloor + 1) \times l$$



Filter with stride (s) = 2

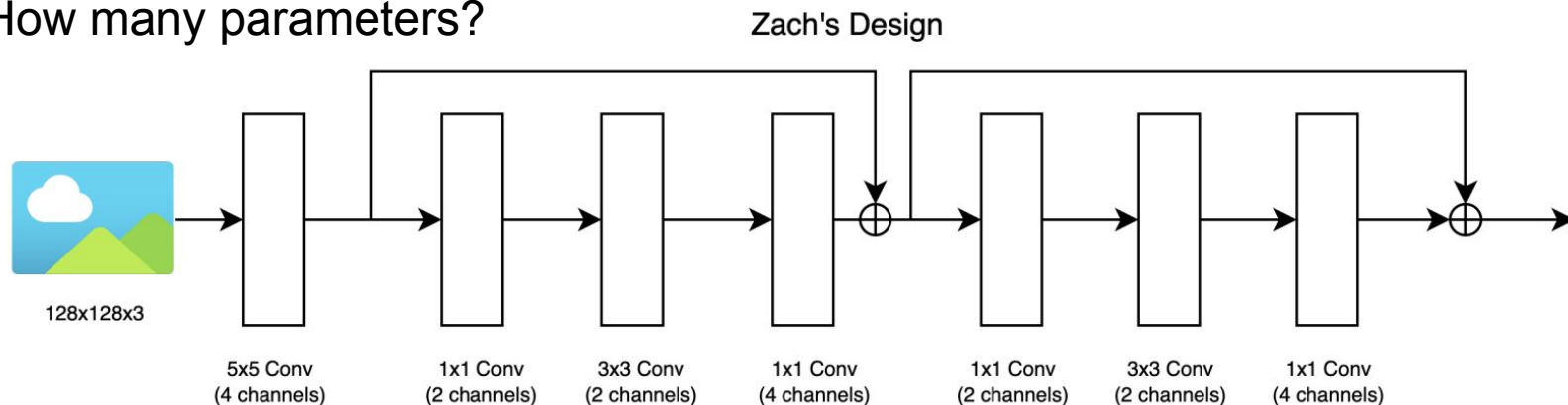
SP24 Prelim Q7

How many parameters?



SP24 Prelim Q7

How many parameters?



Layer 1: $(5 \times 5) \times 3 \times 4 = 300$

Layer 2: $(1 \times 1) \times 4 \times 2 = 8$

Layer 3: $(3 \times 3) \times 2 \times 2 = 36$

Layer 4: $(1 \times 1) \times 2 \times 4 = 8$

Layer 5: $(1 \times 1) \times 4 \times 2 = 8$

Layer 6: $(3 \times 3) \times 2 \times 2 = 36$

Layer 7: $(1 \times 1) \times 2 \times 4 = 8$

Total: 404

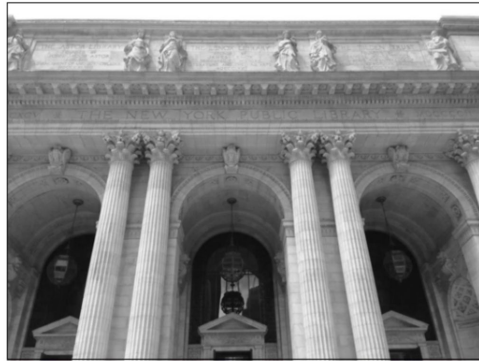
The set of filters is:

$$\left\{ \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}, \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}, \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} \right\}$$

CNN review

Which output comes
from which filter?

Input image:



SP24 Prelim Q5

<https://setosa.io/ev/image-kernels/>

CNN review

Which output comes from which filter?

Input image:



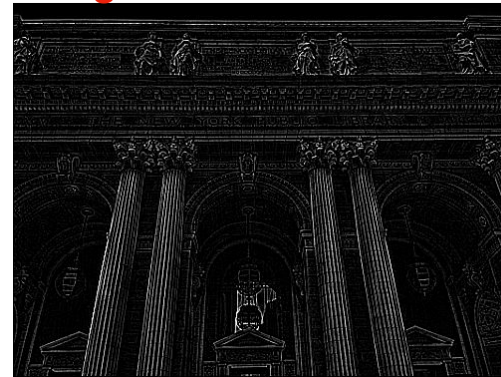
The set of filters is:

$$\left\{ \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}, \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}, \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} \right\}$$

Sharpen



Edge



Left sobel



Bottom sobel

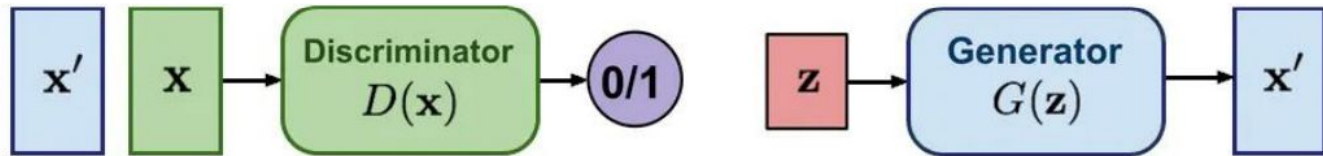


SP24 Prelim Q5

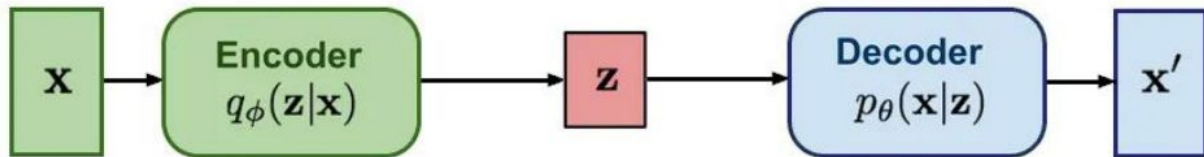
<https://setosa.io/ev/image-kernels/>

GANs

GAN: Adversarial training



VAE: maximize variational lower bound



Diffusion models:
Gradually add Gaussian noise and then reverse



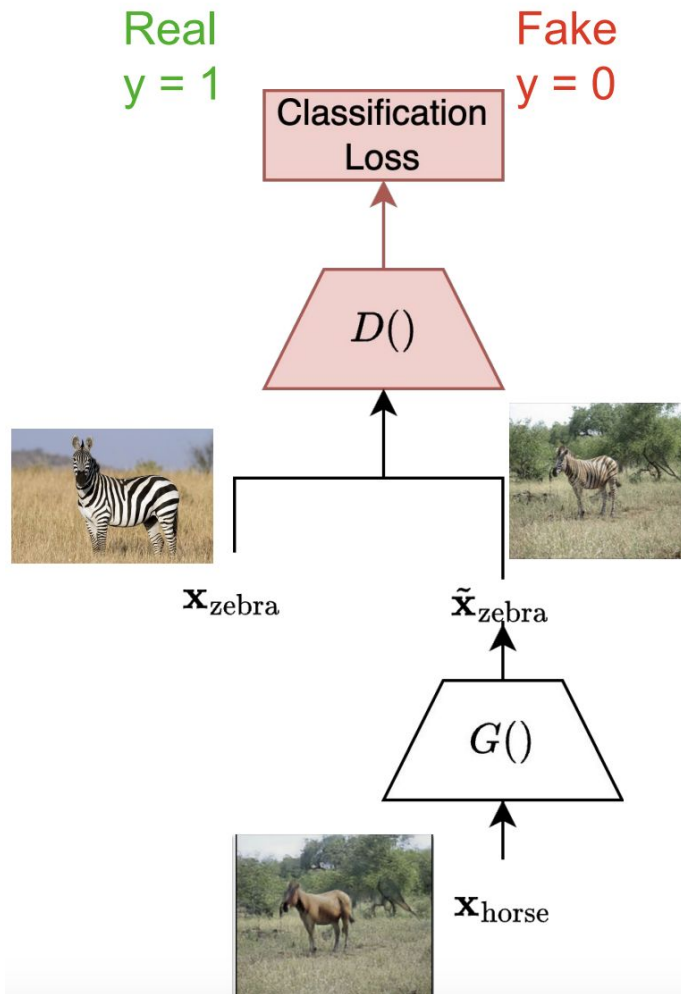
GANs

- Motivation: Conditional image generation without labelled pairs (Why? Data scarcity)
 - Eg. translate horses into zebras



[12 Astonishing Facts About Horses](#)

[Zebra Facts | Live Science](#)



GANs

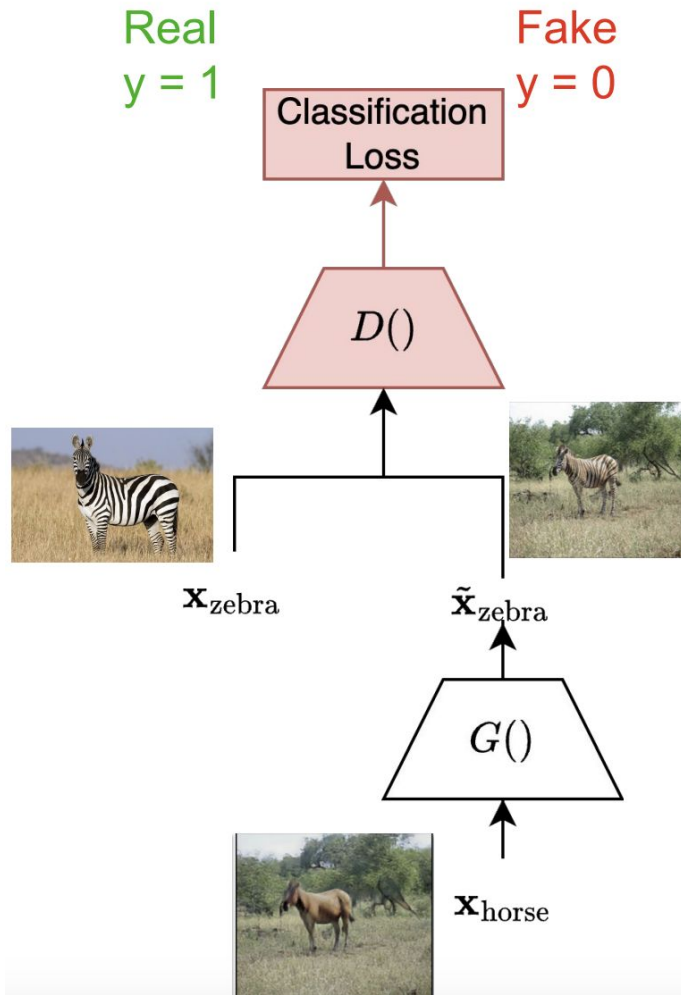
Adversarial Networks

- Generator (G): performs the image translation
 - tries to *fool* the discriminator
- Discriminator (D): predicts whether the image was generated by the generator
 - measures how good the generator did

To train, iteratively:

1. Update discriminator on real/fake images
2. Update generator using feedback from *new* discriminator

Why can't we train them simultaneously?



GANs

Positive example: real
image from labelled dataset

Negative example: fake
image from generator



Discriminator Loss:

$$L_D = \min_D [-\log(D(\mathbf{x})) - \log(1 - D(G(\mathbf{z})))]$$

Generator Loss:

$$\max_G \min_D [-\log(D(\mathbf{x})) - \log(1 - D(G(\mathbf{z})))]$$

$$L_G = -L_D$$

GANs

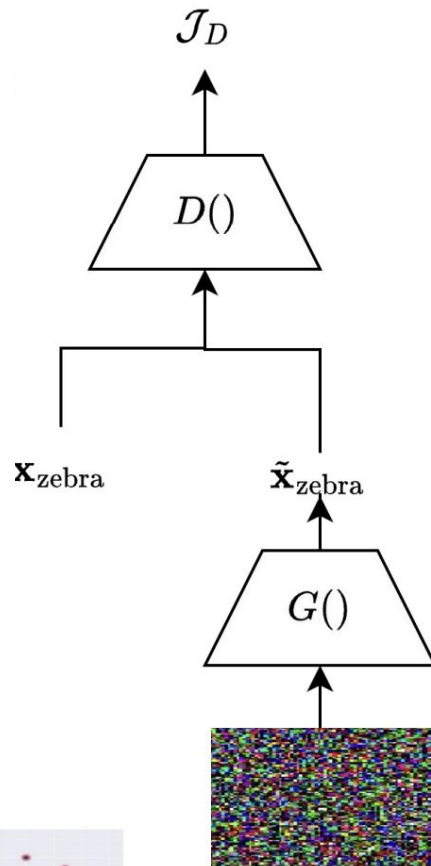
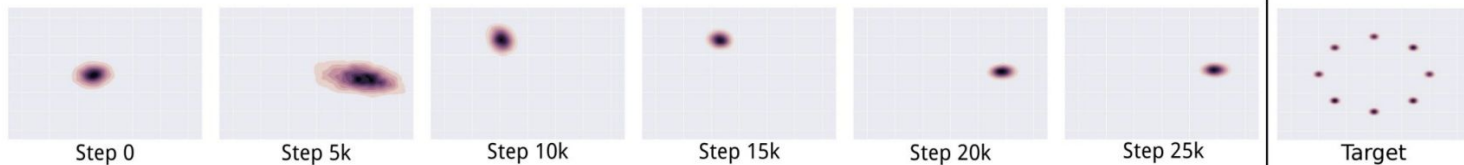
GANs can also be used for *unconditional* image generation

- Sample generator using Gaussian noise
- Results in very high quality images, quickly!

However, GANs experience *Mode Collapse*

- Learn only what best fools the discriminator
- Generated images do not represent full distribution of image class

Mode collapse



Prelim Q2 (Only focus on GANs)

For each of the following statements, indicate whether the statement describes VAEs, GANs, or Diffusion models. You should list all that apply.

1. Trained with a likelihood-based objective.
2. Has a learnable encoder that maps the data to a Gaussian distribution.
3. Has latent variables that must be the same size as the input data.
4. Learns to generate images by fooling a discriminator.
5. Generates high-quality, realistic images.
6. Generates images with a single forward pass of the network.
7. Often suffers from poor diversity among generations.
8. Transforms samples from a unit normal distribution to samples from the data distribution.

Prelim Q2 (Only focus on GANs)

For each of the following statements, indicate whether the statement describes VAEs, GANs, or Diffusion models. You should list all that apply.

1. Trained with a likelihood-based objective.

N

2. Has a learnable encoder that maps the data to a Gaussian distribution.

N

3. Has latent variables that must be the same size as the input data.

N

4. Learns to generate images by fooling a discriminator.

Y

5. Generates high-quality, realistic images.

Y

6. Generates images with a single forward pass of the network.

Y

7. Often suffers from poor diversity among generations.

Y

8. Transforms samples from a unit normal distribution to samples from the data distribution.

Y

Brief U-Net Review

✓ Convolutions

Allow parallelization to extract latent vector for each pixel

✓ Hourglass

Improve efficiency by reducing computations with downsampling

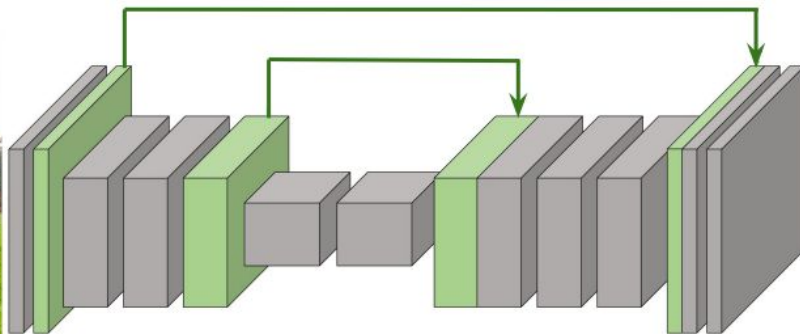
Increase receptive field size by convolving on downsampled feature maps

✓ Skip Connections

Improve prediction quality by combining low-level image features



Input Image



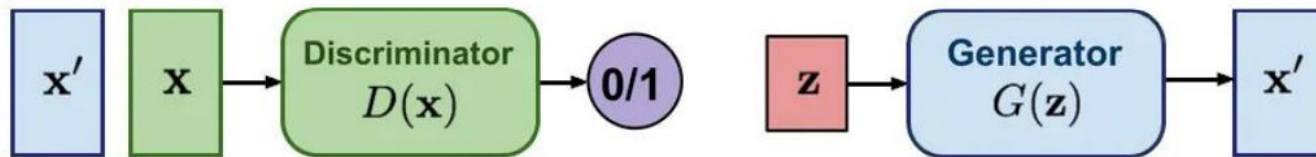
Hourglass CNN with Skip Connections



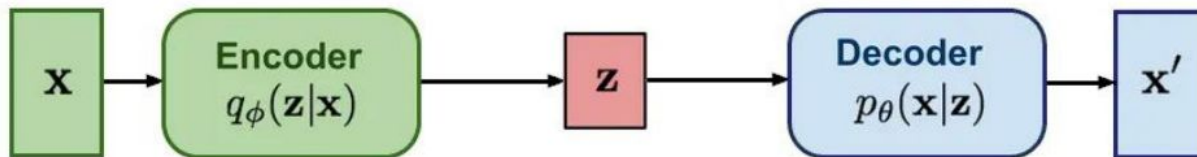
Prediction

VAEs

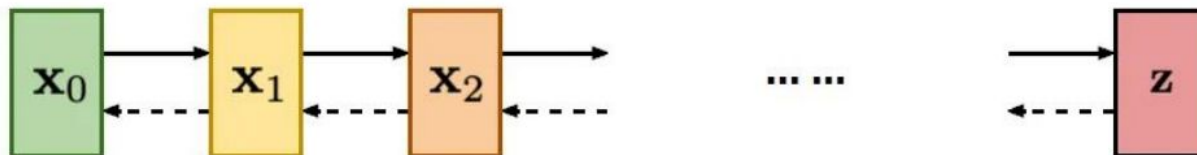
GAN: Adversarial training



VAE: maximize variational lower bound

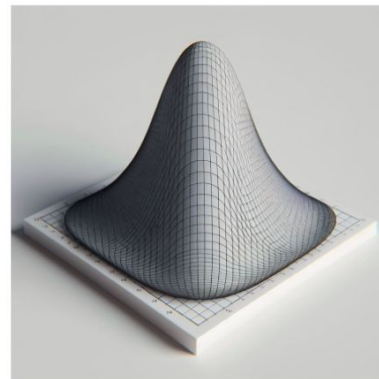
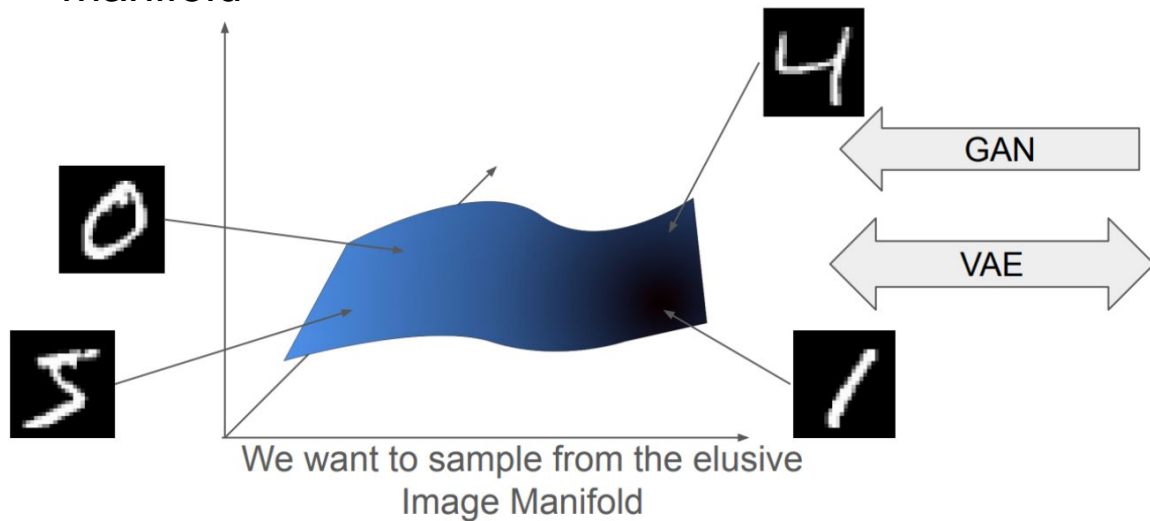


Diffusion models:
Gradually add Gaussian noise and then reverse



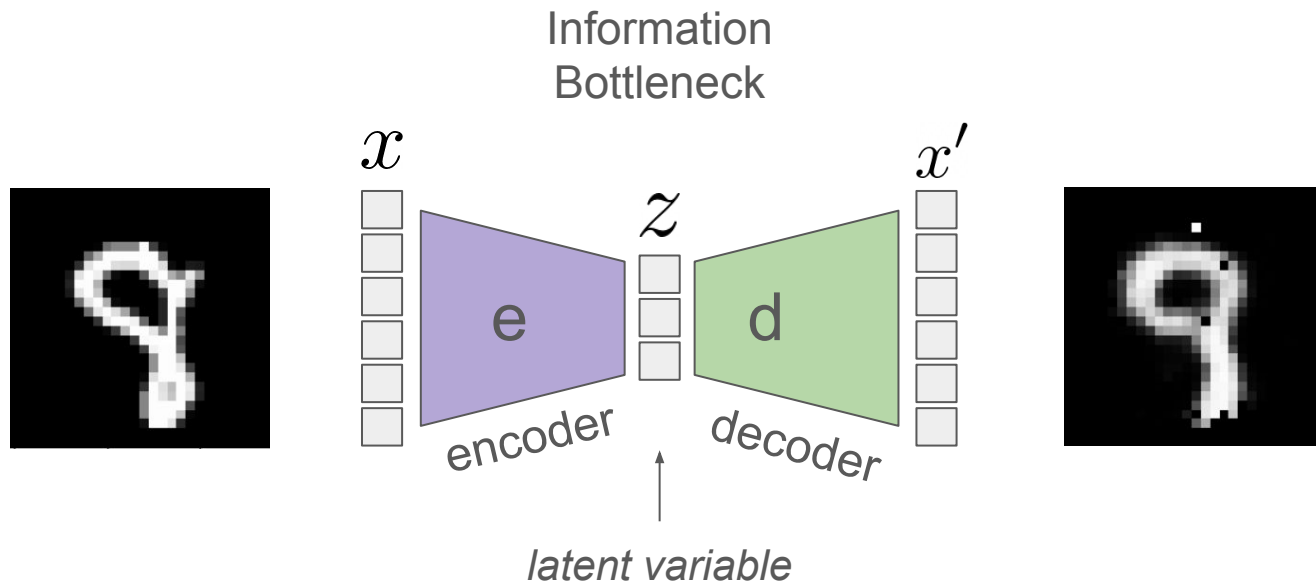
Data Manifold

- Data distribution $\mathbf{P}(\mathbf{x})$ defines a low-dimensional manifold
- Naive random sampling in this space will almost certainly not be on the manifold



We **can** sample from a Gaussian Distribution

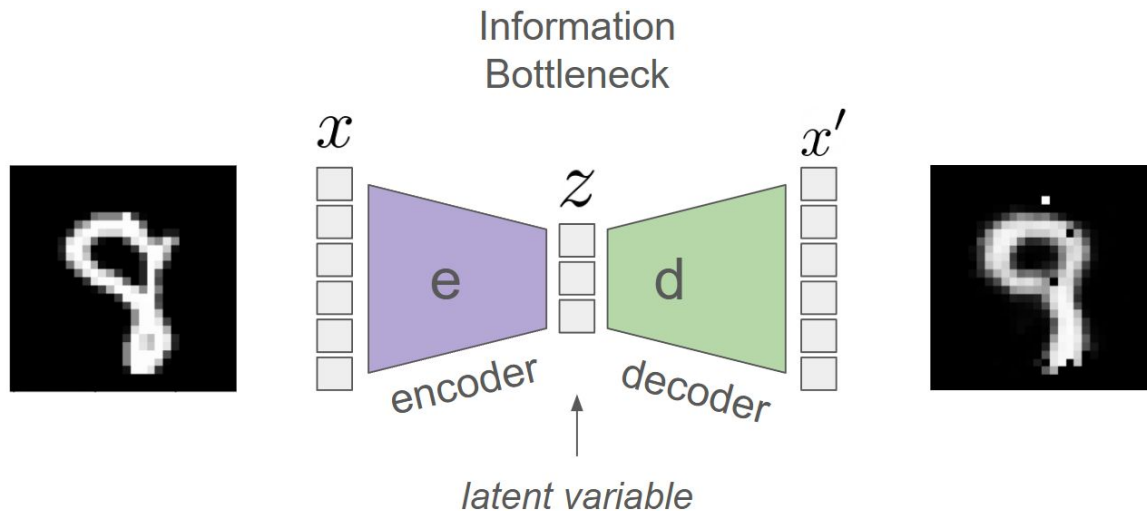
Autoencoders



We typically use MSE or MAE to compute reconstruction loss

Autoencoders -> VAEs

Q: We want to generate images - how can we sample from the latent space?



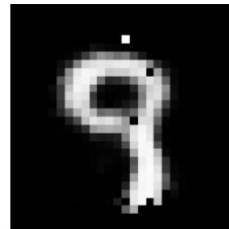
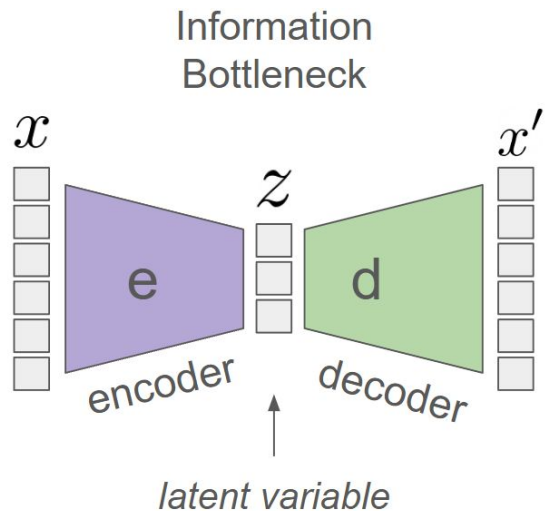
Autoencoders -> VAEs

Q: We want to generate images - how can we sample from the latent space?

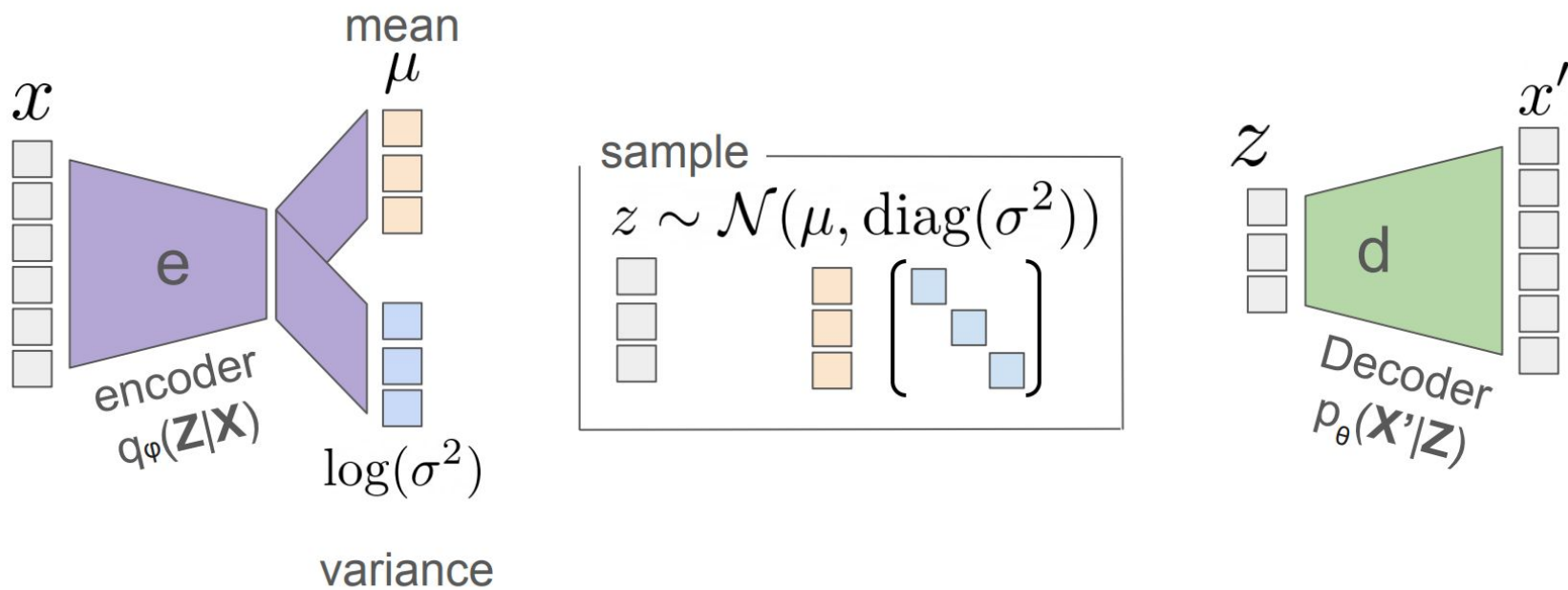
A: Add regularization to encourage the latent space to approximate a Gaussian.

Foreshadowing ELBO...

$$\underbrace{\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))}_{\text{prior matching term}}$$

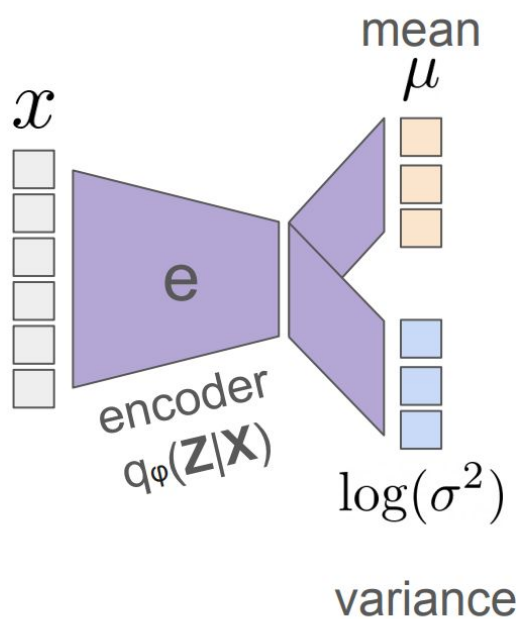


Probabilistic **Encoder** (Gaussian)



$$\max_{\phi, \theta} \mathbb{E}_{z \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log(p_{\theta}(\mathbf{x}|\mathbf{z}))]$$

Probabilistic **Encoder** (Gaussian)



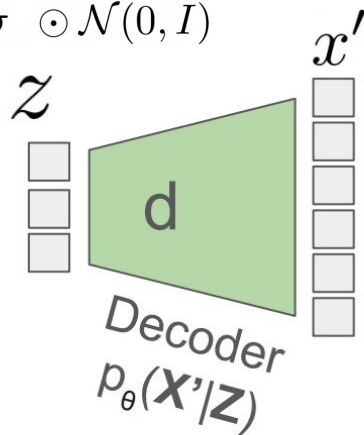
Use the reparameterization trick

to sample $\mathcal{N}(\mu, \text{diag}(\sigma^2)) = \mu + \sigma \odot \mathcal{N}(0, I)$

sample

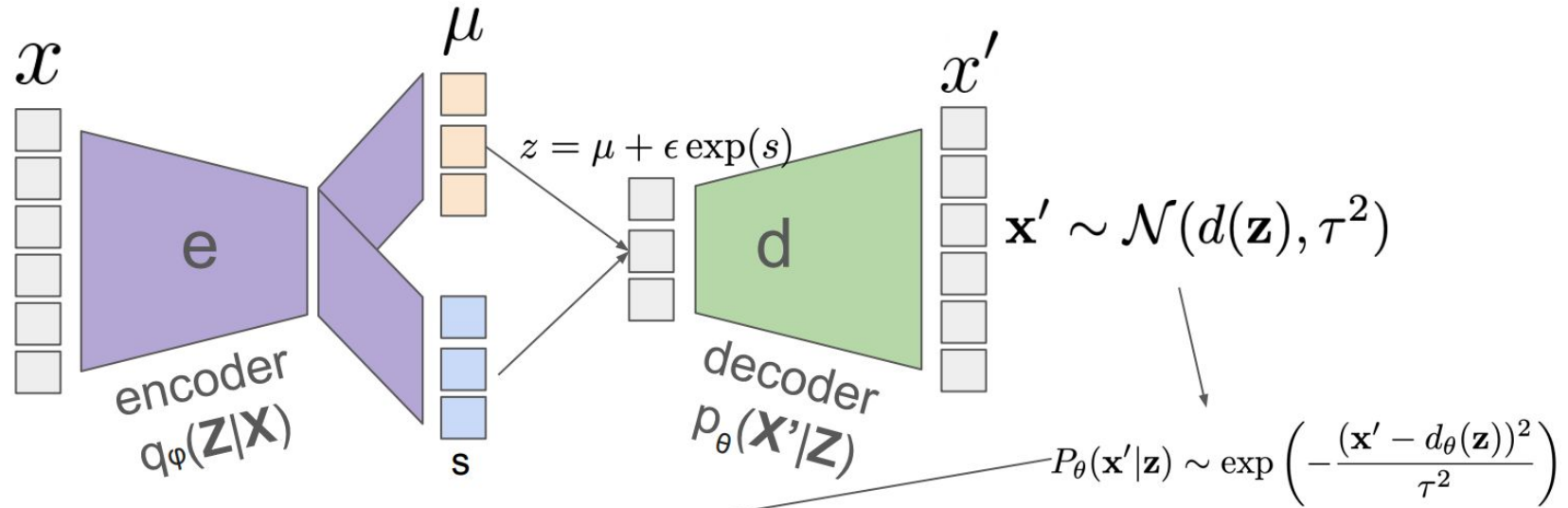
$$z \sim \mathcal{N}(\mu, \text{diag}(\sigma^2))$$

Diagram illustrating the sampling process. It shows the mean μ (3 orange squares) and the log variance $\log(\sigma^2)$ (3 blue squares) being used to sample z . The sampling is represented by a box containing the equation $z \sim \mathcal{N}(\mu, \text{diag}(\sigma^2))$ and a matrix representation of the Gaussian distribution: $\begin{bmatrix} \square & & \\ & \square & \\ & & \square \end{bmatrix}$.



$$\max_{\phi, \theta} \mathbb{E}_{z \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log(p_\theta(\mathbf{x}|\mathbf{z}))]$$

Probabilistic **decoder** (Gaussian)

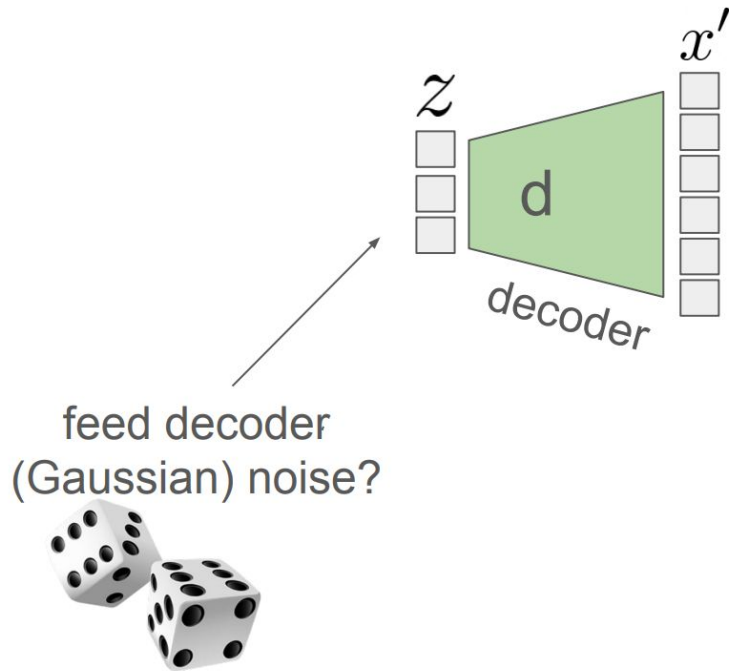
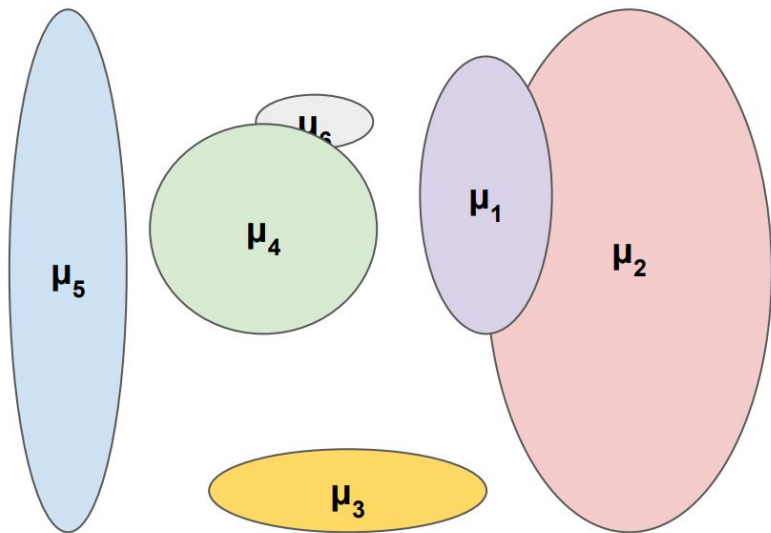


$$\max_{\phi, \theta} \mathbb{E}_{z \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log(p_\theta(\mathbf{x}|\mathbf{z}))] = \min_{\phi, \theta} \mathbb{E}_{z \sim q_\phi(\mathbf{z}|\mathbf{x})} [(\mathbf{x} - d_\theta(\mathbf{z}))^2]$$

Upshot: Reconstruction formulation results in squared loss

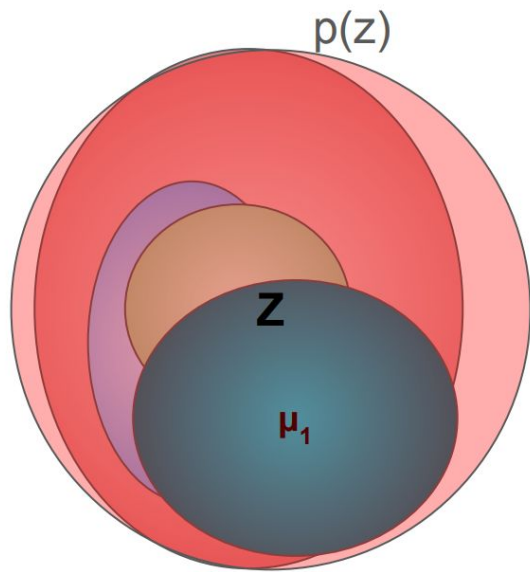
Sampling

How can we sample, if each sample has its own latent distribution?

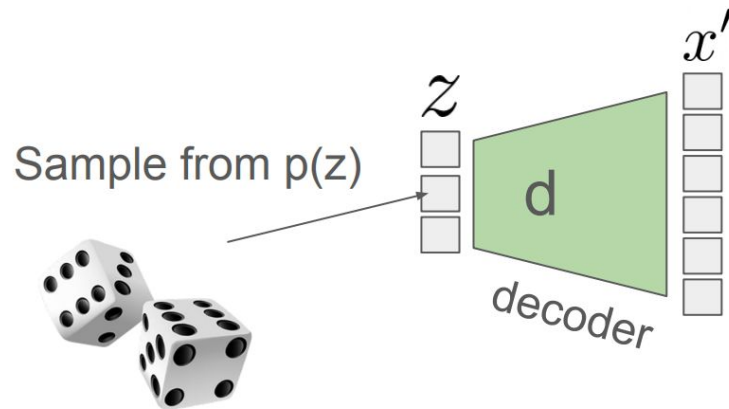


Sampling

Solution: Regularize all distributions to be close to the standard normal $\mathbf{N}(\mathbf{0}; \mathbf{I})$.



$$\begin{array}{cc} \text{maximize} & \\ \underbrace{\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{\text{reconstruction term}} & - \underbrace{D_{\text{KL}}(q_{\phi}(z|x) \parallel p(z))}_{\text{prior matching term}} \end{array}$$



ELBO

$$\log p(\mathbf{x}) = \log p(\mathbf{x}) \int q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z}$$

(Multiply by $1 = \int q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z}$)

$$= \int q_{\phi}(\mathbf{z}|\mathbf{x}) (\log p(\mathbf{x})) d\mathbf{z}$$

(Bring evidence into integral)

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x})]$$

(Definition of Expectation)

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \right]$$

(Apply Equation 2)

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}) q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x}) q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

(Multiply by $1 = \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})}$)

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x})} \right]$$

(Split the Expectation)

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] + D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}))$$

(Definition of [KL Divergence](#))

$$\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

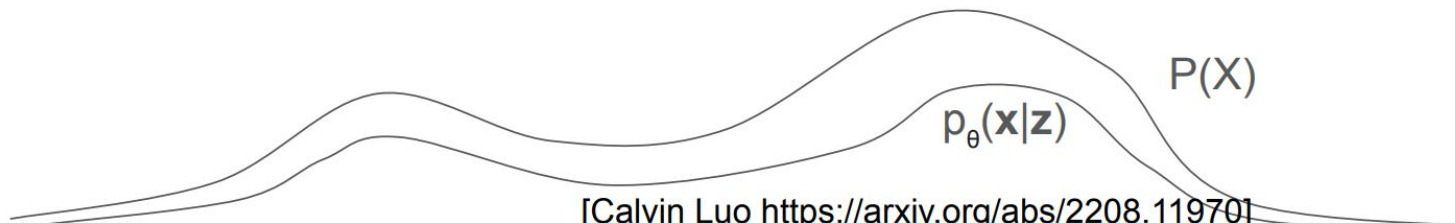
(KL Divergence always ≥ 0)

ELBO

$$\begin{aligned}\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] && \text{(Chain Rule of Probability)} \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] && \text{(Split the Expectation)} \\ &= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{prior matching term}} && \text{(Definition of KL Divergence)}\end{aligned}$$

(We are **maximizing** this lower bound.)

If we maximize $p_\theta(\mathbf{x}|\mathbf{z})$ and minimize the D_{KL} we get close to $P(\mathbf{x})$.



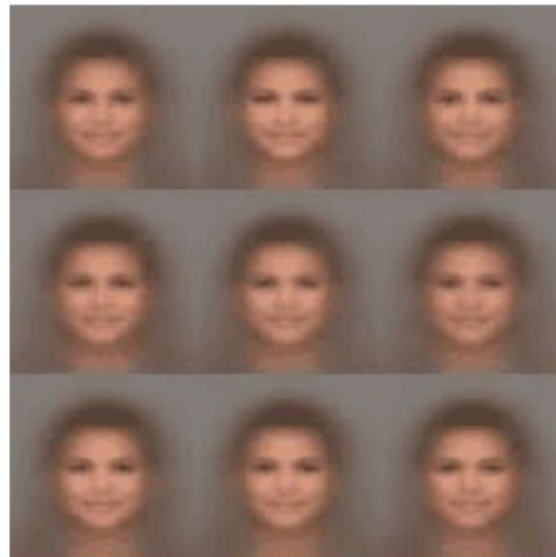
Drawbacks

- Out of the box, generated images can be blurry.

Question: Why? How do GANs fix this problem?



[VAE v. GAN](#)



<https://borisburkov.net/2022-12-31-1/>

Summary

- Generative Image models learn a mapping from the **Standard Normal Gaussian** to the **Image Manifold**
 - GANs learn this through a **discriminator**.
 - VAEs learn it through **variational autoencoders**
- AutoEncoders learn to **compress** and **reconstruct** data
- VAEs make these AutoEncoders **probabilistic**
 - Minimize the **reconstruction loss**
 - Latent space is sampled from Gaussian distributions
 - Sampling is made differentiable with the **Reparameterization Trick**
 - Deviations from the Prior (Standard Normal Gaussian) is penalized by **KL divergence**
- The ELBO is a **lower bound** of $P(X)$
 - Maximizing the ELBO, and minimizing the KL divergence makes $P(x|z)$ close to $P(x)$

Prelim Problem (Only focus on VAEs)

For each of the following statements, indicate whether the statement describes VAEs, GANs, or Diffusion models. You should list all that apply.

1. Trained with a likelihood-based objective.
2. Has a learnable encoder that maps the data to a Gaussian distribution.
3. Has latent variables that must be the same size as the input data.
4. Learns to generate images by fooling a discriminator.
5. Generates high-quality, realistic images.
6. Generates images with a single forward pass of the network.
7. Often suffers from poor diversity among generations.
8. Transforms samples from a unit normal distribution to samples from the data distribution.

Prelim Problem (Only focus on VAEs)

For each of the following statements, indicate whether the statement describes VAEs, GANs, or Diffusion models. You should list all that apply.

1. Trained with a likelihood-based objective.

Y

2. Has a learnable encoder that maps the data to a Gaussian distribution.

Y

3. Has latent variables that must be the same size as the input data.

N

4. Learns to generate images by fooling a discriminator.

N

5. Generates high-quality, realistic images.

N

6. Generates images with a single forward pass of the network.

Y

7. Often suffers from poor diversity among generations.

N

8. Transforms samples from a unit normal distribution to samples from the data distribution.

Y

Prelim Problem: ELBO

In class, we saw that the ELBO is given as:

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] = \log p(\mathbf{x}) - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})),$$

where the left-hand-side is the ELBO, \mathbf{x} represents the observed data, and \mathbf{z} represents the latent variables.

a) Show that the ELBO is a lower bound of the log-likelihood of the data.

$$\begin{aligned} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] &= \log p(\mathbf{x}) - D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})) \\ &\leq \log p(\mathbf{x}) \end{aligned} \quad (\text{Non-negativity of KL})$$

Prelim Problem: ELBO

b) Rearrange the ELBO to show that that

$$\mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(x, z)}{q_\phi(z|x)} \right] = E_{q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x) || p(z)).$$

Recall that the KL-Divergence is defined as:

$$\begin{aligned} D_{KL}(P||Q) &= \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] \\ \log p(x) &\geq \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(x, z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p_\theta(x|z)p(z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] + \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p(z)}{q_\phi(z|x)} \right] \\ &= E_{q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x) || p(z)) \end{aligned}$$

Prelim Problem: ELBO

c) Interpret the two terms in the final ELBO expression in the context of VAEs. What effect does each term have? Write 2-3 sentences.

d) Consider a situation where your variational distribution perfectly matches the true posterior. In other words, $q_\phi(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$. Given this, can you simplify the the ELBO further? Recall that the ELBO is given as

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] = \log p(\mathbf{x}) - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})).$$

If your variational distribution matches the true posterior, what does maximizing the ELBO accomplish?

Prelim Problem: ELBO

c) Interpret the two terms in the final ELBO expression in the context of VAEs. What effect does each term have? Write 2-3 sentences.

The KL divergence term helps regularize the distribution to be close to the prior distribution. The first term is the expected reconstruction error.

d) Consider a situation where your variational distribution perfectly matches the true posterior. In other words, $q_\phi(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$. Given this, can you simplify the the ELBO further? Recall that the ELBO is given as

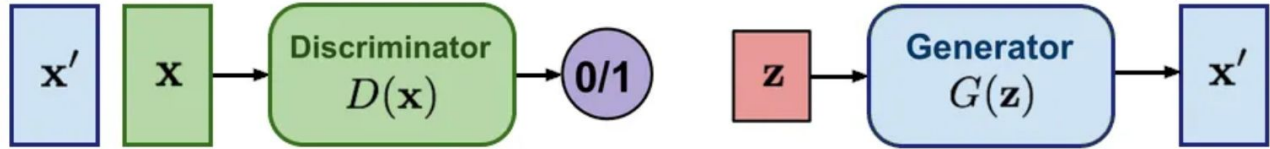
$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] = \log p(\mathbf{x}) - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})).$$

If your variational distribution matches the true posterior, what does maximizing the ELBO accomplish?

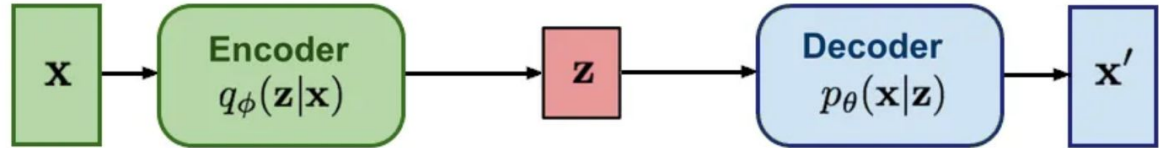
When the variational distribution perfectly matches the true posterior, the KL divergence term in the ELBO becomes zero since $D_{\text{KL}}(p(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x})) = 0$. Therefore, the ELBO simplifies to the log-likelihood of the data, $\log p(\mathbf{x})$. Maximizing the ELBO is then equivalent to maximizing the log-likelihood of the data.

Diffusion Models

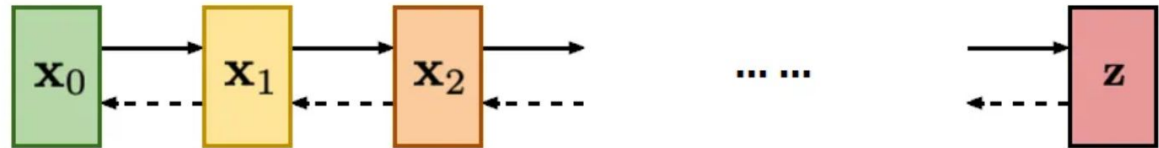
GAN: Adversarial training



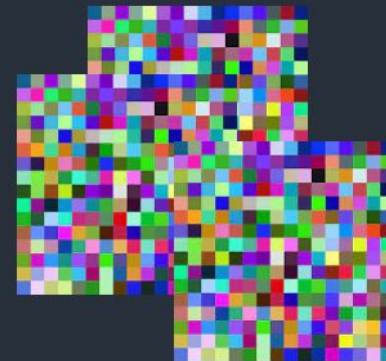
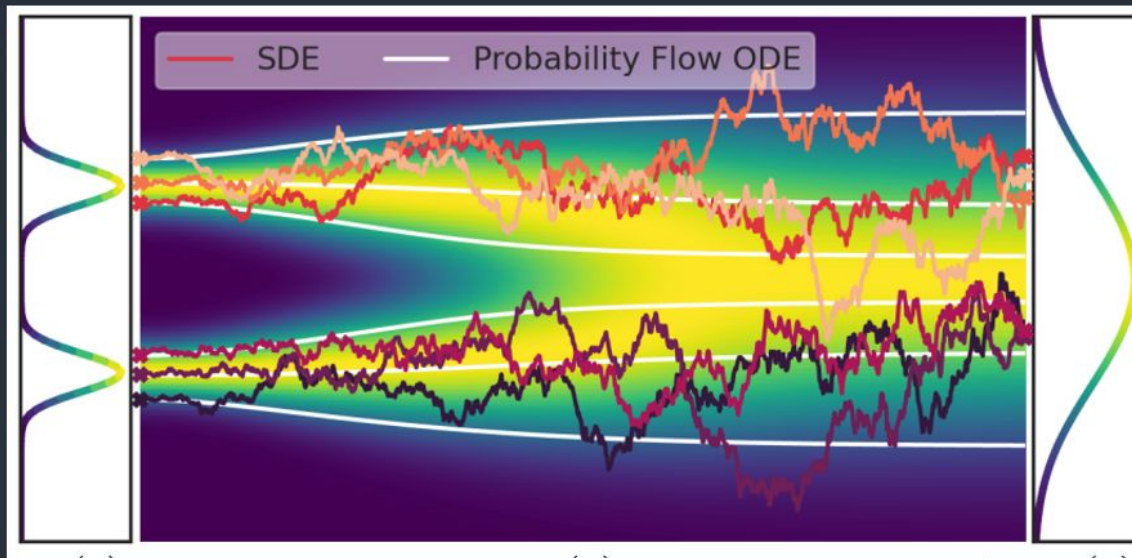
VAE: maximize variational lower bound



Diffusion models:
Gradually add Gaussian noise and then reverse



Forward diffusion



Reverse Process (generative)

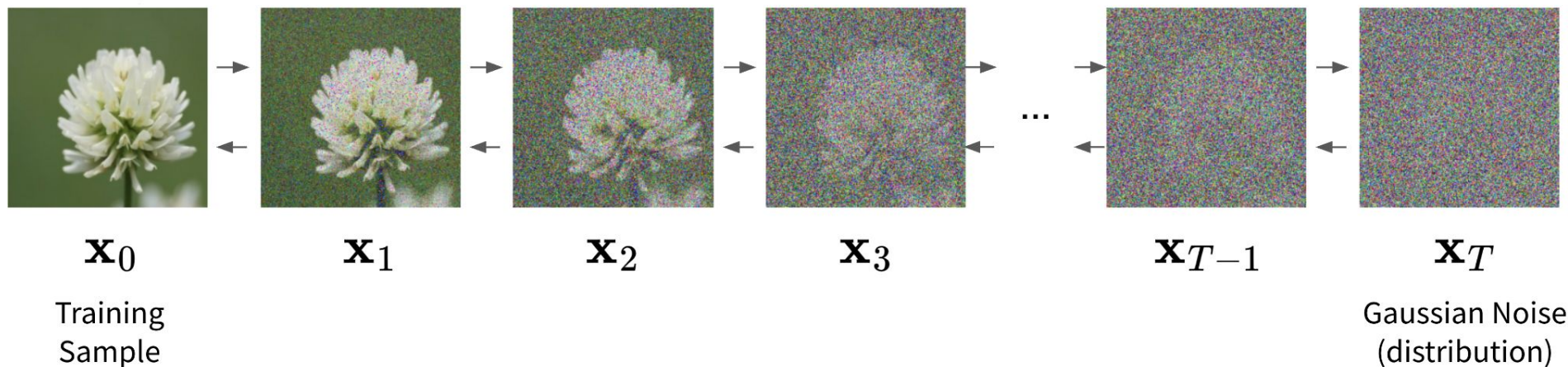
Diffusion Models - Main point

We DEFINE a mapping between a source distribution and Gaussian Distribution (Gaussian Noise)

We LEARN a *reverse* mapping from the Gaussian Distribution to the source distribution.

Forward Process

- “Destroying Data”
- Adding noise forms a **Markov Chain** → **Markov Property**
- **Recall:** $q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0) = q(\mathbf{x}_t | \mathbf{x}_{t-1})$



Forward Process

- Given a *sampling schedule*, predict the noisy image at timestep \mathbf{t} from timestep $\mathbf{0}$:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

- Usually, we can afford a larger update step when the sample gets noisier so:

$$\bar{\alpha}_1 > \dots > \bar{\alpha}_T.$$

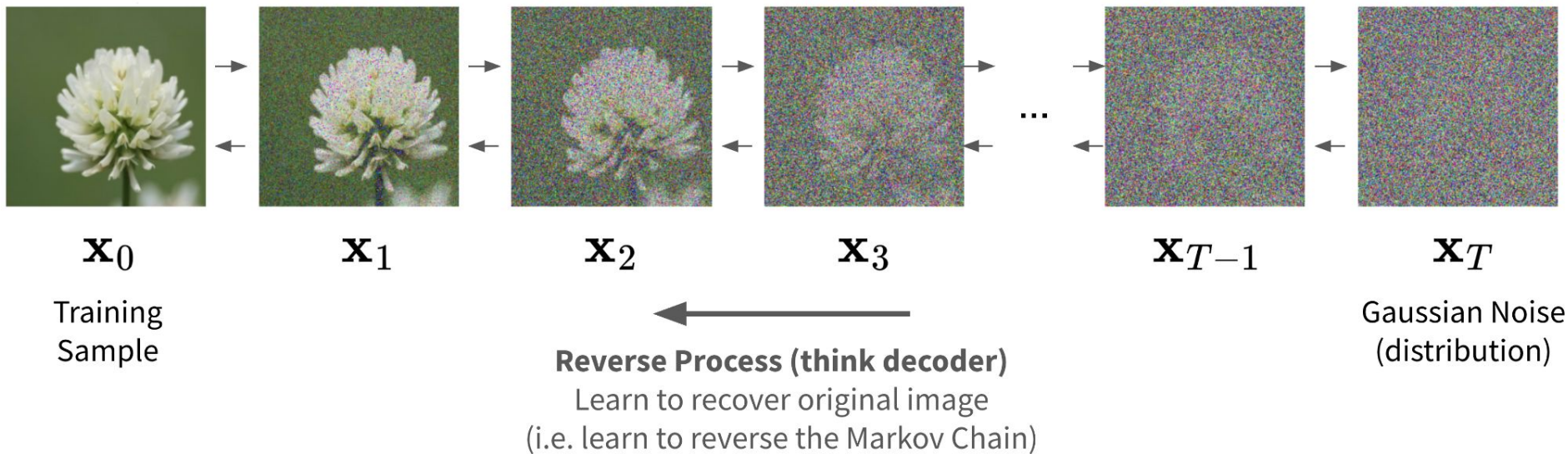
- Want:

$$q(\mathbf{x}_T) \approx \mathcal{N}(0, \mathbf{I})$$

Backward Process

Goal: given an image that was “*noised*” for t steps, predict the original image

- Markov chains are not always invertible.
- We *learn* the *inverse (Markov) process*



How do we sample from reverse?

- Sure, try Bayes?
- There's a **problem...**
- **Intractable:** $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})q(\mathbf{x}_{t-1})}{q(\mathbf{x}_t)}$$



$$q(\mathbf{x}_t) = \int q(\mathbf{x}_0)q(\mathbf{x}_1|\mathbf{x}_0)...q(\mathbf{x}_t|\mathbf{x}_{t-1})d\mathbf{x}_0d\mathbf{x}_1...d\mathbf{x}_{t-1}$$

Learn the reverse

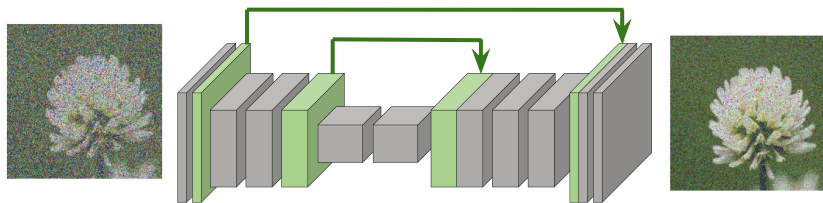
We have \mathbf{x}_0 during **training**; train a **generative model**

$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ is **tractable**

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

$$\mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))]$$



How do we actually train for the reverse process?

Find the model that **maximizes the likelihood**
of the training data

$$\mathbf{max} \log p(\mathbf{x})$$

$$\log p(\mathbf{x}) \geq \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

Find the model that **maximizes the likelihood** of the training data

$\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t))]$ is a *denoising matching term*. We learn desired denoising transition step $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$ as an approximation to tractable, ground-truth denoising transition step $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$. The $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ transition step can act as a ground-truth signal, since it defines how to denoise a noisy image \mathbf{x}_t with access to what the final, completely denoised image \mathbf{x}_0 should be. This term is therefore minimized when the two denoising steps match as closely as possible, as measured by their KL Divergence.

$$\log p(\mathbf{x}) \geq \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{reconstruction term}} - \underbrace{\mathcal{D}_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) || p(\mathbf{x}_T))}_{\text{prior matching term}} - \sum_{t=2} \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}}$$

With large T , the prior matching goes to 0

Reparametrization of the noise prediction

Recall that we need to learn a neural network to approximate

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

Want to train $\boldsymbol{\mu}_{\theta}$ to get:

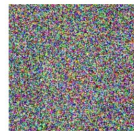
$$\tilde{\boldsymbol{\mu}}_t = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right)$$

Since we have \mathbf{x}_t during training, we can reparametrize to predict

$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t)$



$\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)$



$\boldsymbol{\epsilon}_t$

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right)$$

$$\text{Thus } \mathbf{x}_{t-1} = \mathcal{N}(\mathbf{x}_{t-1}; \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

Reparametrization of the Loss

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right)$$

$$\text{Thus } \mathbf{x}_{t-1} = \mathcal{N}(\mathbf{x}_{t-1}; \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$

$$\sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{denoising matching term}} \longrightarrow L(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)\|^2]$$

Loss is **MSE** of actual predicted loss!

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_{\theta}(\mathbf{x}_t, t)) \approx q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$$

Training Algo

Repeat until convergence

1. $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ \leftarrow Sample original image from image distribution
2. $t \sim U\{1, 2, \dots, T\}$ \leftarrow Sample random time step uniformly
3. $\epsilon \sim \mathcal{N}(0, 1)$ \leftarrow Sample Gaussian noise
4. Optimizer step on $L(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2]$
 - \leftarrow Model predicts noise applied at time step t and calculate loss

Sampling Algo


$\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ \leftarrow Sample pure Gaussian noise

For $t = T, T - 1 \dots, 1$

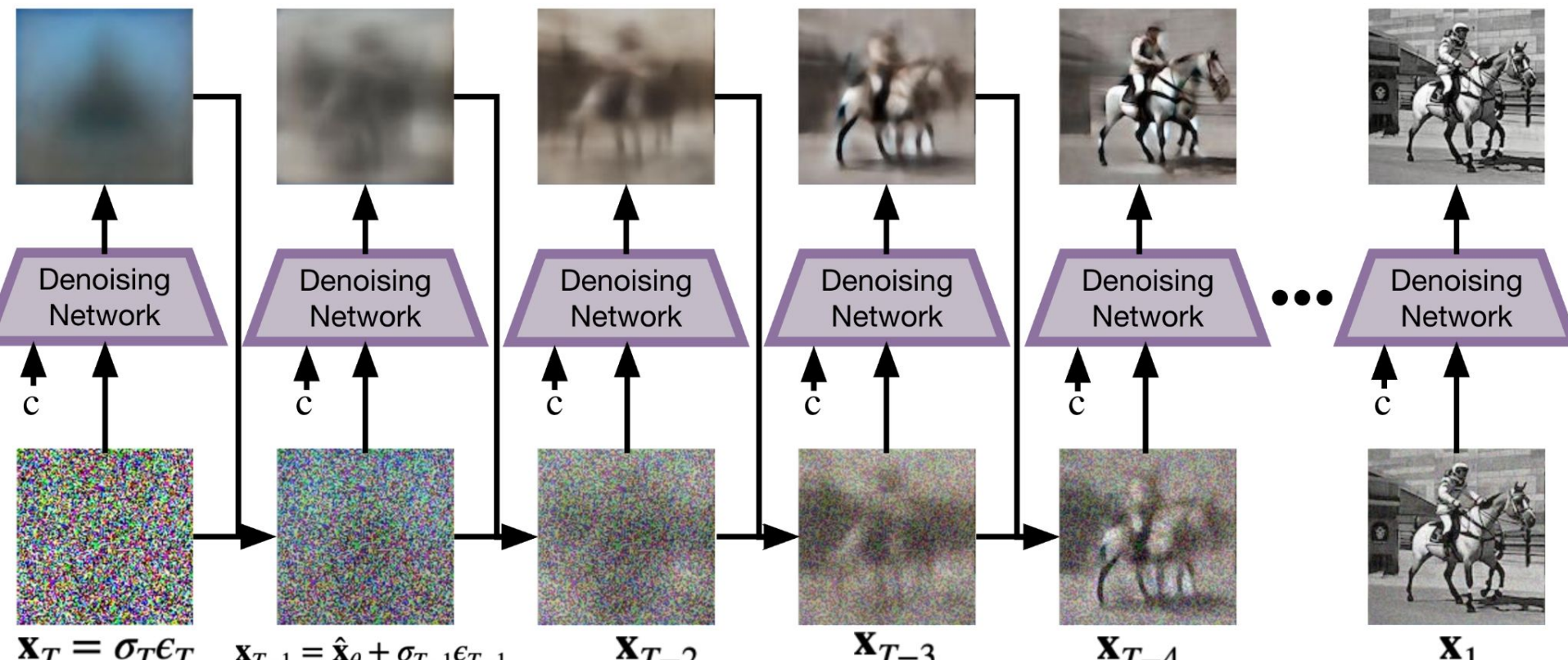
$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = \mathbf{0}$ \leftarrow Sample Gaussian noise to apply to image

$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ \leftarrow Predict noise applied to image and remove that noise

Return \mathbf{x}_0


$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_\theta(\mathbf{x}_t, t))$$

Diffusion Sampling



Alternative Perspective

So far:

Denoising Diffusion Probabilistic Models

Alternative:

**Generative Modeling by Estimating Gradients of the
Data Distribution**

Score-based Models

Langevin dynamics allow you to **sample** from distribution - **even if it is not normalized**.

Would like to model the probability density function:

$$p_{\theta}(\mathbf{x}) = \frac{e^{-f_{\theta}(\mathbf{x})}}{Z_{\theta}}, \text{ where } Z_{\theta} > 0 \text{ is a normalizing constant s.t. } \int p_{\theta}(\mathbf{x}) d\mathbf{x} = 1$$

Still want to **maximize the log-likelihood**

$$\max_{\theta} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i).$$

Problem: Normalization constant intractable \rightarrow Approximate the score function:

$$\mathbf{s}_{\theta}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) = -\nabla_{\mathbf{x}} f_{\theta}(\mathbf{x}) - \underbrace{\nabla_{\mathbf{x}} \log Z_{\theta}}_{=0} = -\nabla_{\mathbf{x}} f_{\theta}(\mathbf{x})$$

What is the score anyway?

- What direction in data space to move in order to further increase its likelihood

Loss?

$$\begin{aligned} & \arg \min_{\theta} \mathcal{D}_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)) \\ &= \arg \min_{\theta} \mathcal{D}_{\text{KL}}(\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \parallel \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t))) \end{aligned} \quad (83)$$

$$= \arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{\alpha_t} \left[\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla \log p(\mathbf{x}_t)\|_2^2 \right] \quad (84)$$

The gradient of x in dataspace, for arbitrary noise level t

Add noise for more accurate scores

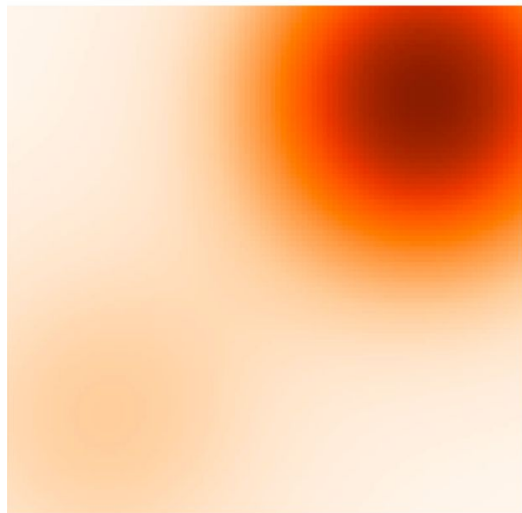
Fisher Divergence

$$\mathbb{E}_{\mathbf{x}} \left[\left\| \mathbf{s}_{\theta}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p(x) \right\|_2^2 \right]$$

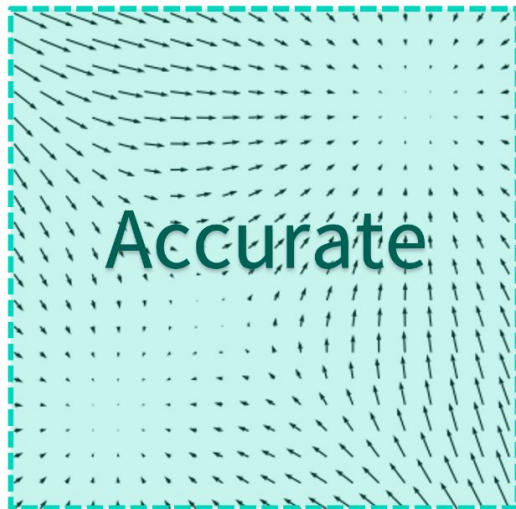
Predicted score!

unknown!

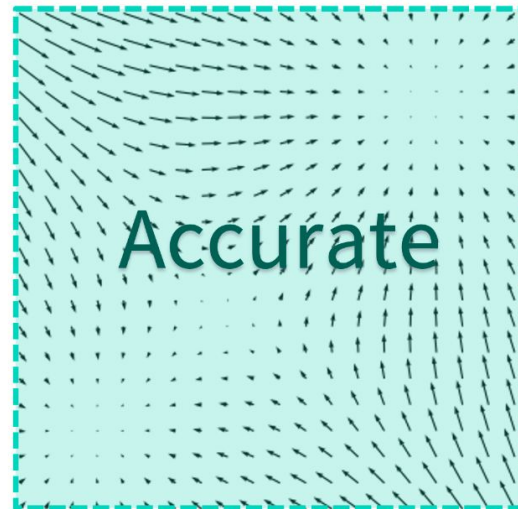
Perturbed density



Perturbed scores



Estimated scores



Training

Training Objective for noise level t :

$$\sum_{t=1}^T \lambda(t) \mathbb{E}_{\mathbf{x}, t} [\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|_2^2]$$

Using results from denoising score matching [1]:

$$\sum_{t=1}^T \lambda(t) \mathbb{E}_{\mathbf{x}, t} [\|\mathbf{s}_{\theta}(\tilde{\mathbf{x}}, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | \mathbf{x})\|_2^2]$$

Using the definition of the pdf of a Gaussian,

$$\sum_{t=1}^T \lambda(t) \mathbb{E}_{\mathbf{x}, t} [\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \frac{\mathbf{x}_t - \mathbf{x}}{\sigma_t^2}\|_2^2]$$

Looks familiar?

Okay, saw we learned this cool and fancy score function, what now?

- Just sample through Langevin dynamics!
- starting at any *arbitrary point* in the same space, iteratively follow the score until a mode is reached :)

Conditioning? Guidance?

- Just condition at each step!

$$p(\mathbf{x}_{0:T} \mid y) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, y)$$

$$\mathbf{s}_{\theta}(\mathbf{x}_t, t, y) \approx \nabla \log p(\mathbf{x}_t \mid y)$$

Conditional Diffusion — Classifier Guidance

- Use Bayes' rule to decompose the conditional score into the unconditional score and a likelihood term

$$\begin{aligned}
 \nabla \log p(\mathbf{x}_t | y) &= \nabla \log \left(\frac{p(\mathbf{x}_t)p(y | \mathbf{x}_t)}{p(y)} \right) \\
 &= \nabla \log p(\mathbf{x}_t) + \nabla \log p(y | \mathbf{x}_t) - \nabla \log p(y) \\
 &= \underbrace{\nabla \log p(\mathbf{x}_t)}_{\text{unconditional score}} + \underbrace{\nabla \log p(y | \mathbf{x}_t)}_{\text{adversarial gradient}}
 \end{aligned}$$

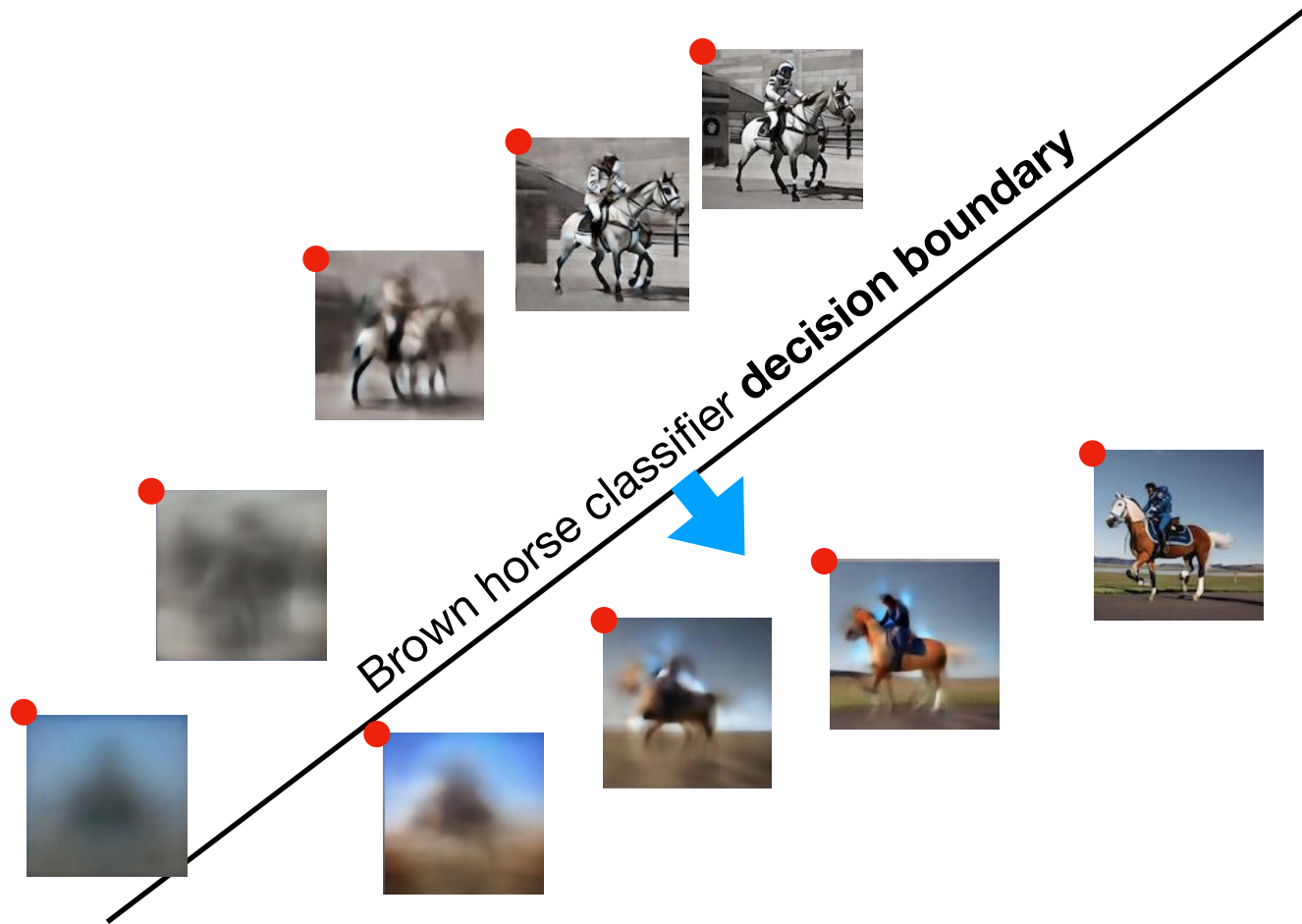
- Only need to train a classifier on noised data
- Use classifier to guide noise!

Can add \gamma

$$\mathbf{x}_{t-1} = \hat{\mathbf{x}}_{\theta}(\mathbf{x}_t) + \sigma_{t-1}\epsilon_{t-1} + \alpha \nabla_{\mathbf{x}_t} \log (P(y | \mathbf{x}_t))$$

de-noise
guide noise

add noise



Classifier-Free Guidance

- Train a joint conditional and unconditional diffusion model
- Conditioning information is added by concatenating to input or cross attending
- Modified conditional distribution

Randomly drops the condition during training and linearly combines the condition and unconditional output during sampling

$$\log p_t(\mathbf{x}_t|\mathbf{y}) \propto p_t(\mathbf{x}_t|\mathbf{y}) p_t(\mathbf{y}|\mathbf{x}_t)^w$$

- Conditional sampling

$$\nabla_{\mathbf{x}_t} \log \tilde{p}_t(\mathbf{x}_t|\mathbf{y}) = \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) + w(\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t))$$

Significantly improves quality of conditional models, but decreases diversity

Equivalence of the two views

We have therefore derived three equivalent objectives to optimize a VDM: learning a neural network to predict the original image \mathbf{x}_0 , the source noise ϵ_0 , or the score of the image at an arbitrary noise level $\nabla \log p(\mathbf{x}_t)$ [2, 10]. The VDM can be scalably trained by stochastically sampling timesteps t and minimizing the norm of the prediction with the ground truth target.

Problem 2

For each of the following statements, indicate whether the statement describes VAEs, GANs, or Diffusion models. You should list all that apply.

1. Trained with a likelihood-based objective.
2. Has a learnable encoder that maps the data to a Gaussian distribution.
3. Has latent variables that must be the same size as the input data.
4. Learns to generate images by fooling a discriminator.
5. Generates high-quality, realistic images.
6. Generates images with a single forward pass of the network.
7. Often suffers from poor diversity among generations.
8. Transforms samples from a unit normal distribution to samples from the data distribution.

Problem 2

For each of the following statements, indicate whether the statement describes VAEs, GANs, or Diffusion models. You should list all that apply.

- | | |
|---|--|
| Y | 1. Trained with a likelihood-based objective. |
| N | 2. Has a learnable encoder that maps the data to a Gaussian distribution. |
| Y | 3. Has latent variables that must be the same size as the input data. |
| N | 4. Learns to generate images by fooling a discriminator. |
| Y | 5. Generates high-quality, realistic images. |
| N | 6. Generates images with a single forward pass of the network. |
| N | 7. Often suffers from poor diversity among generations. |
| Y | 8. Transforms samples from a unit normal distribution to samples from the data distribution. |

9. CLIP is a constrastive learning algorithm that learns to map paired text and images close together in a shared embedding space.

10. Vision transformers typically require less training data than convolutional neural networks (CNNs) to achieve comparable performance on image classification tasks.

9. CLIP is a constrastive learning algorithm that learns to map paired text and images close together in a shared embedding space.

Solution: True.

10. Vision transformers typically require less training data than convolutional neural networks (CNNs) to achieve comparable performance on image classification tasks.

Solution: False. Vision transformers often require larger amounts of training data compared to CNNs to achieve comparable performance on image classification tasks. CNNs have inductive biases such as translation invariance and local receptive fields, which allow them to learn effective features with less training data. However, when trained on sufficiently large datasets, vision transformers have shown impressive performance and have even outperformed CNNs on various benchmarks.

Problem 9: Image Captioning [9 pts]

Adhitya wants to train a model for image captioning. He plans to train a decoder-only transformer model that generates a natural language description by attending to visual features extracted from an input image. He has access to a pretrained ViT model to extract image features. Adhitya has written most of the network architecture and needs a little help with the cross-attention layer in the decoder (the cross-attention layer is used to attend to the image features).

a) Adhitya has already written code to divide the input images of shape $(h \times w \times 3)$ into k patches and feed the patches to a ViT to produce 128-dimensional features for every patch. The hidden dimension of the decoder is 128. What are the dimensions of the weight matrices W_Q, W_K, W_V in the cross-attention layer?

Problem 9: Image Captioning [9 pts]

Adhitya wants to train a model for image captioning. He plans to train a decoder-only transformer model that generates a natural language description by attending to visual features extracted from an input image. He has access to a pretrained ViT model to extract image features. Adhitya has written most of the network architecture and needs a little help with the cross-attention layer in the decoder (the cross-attention layer is used to attend to the image features).

a) Adhitya has already written code to divide the input images of shape $(h \times w \times 3)$ into k patches and feed the patches to a ViT to produce 128-dimensional features for every patch. The hidden dimension of the decoder is 128. What are the dimensions of the weight matrices W_Q, W_K, W_V in the cross-attention layer?

The dimensions of the weight matrices W_Q, W_K, W_V in the cross-attention layer are:

$W_Q : (128, 128)$ $W_K : (128, 128)$ $W_V : (128, 128)$

Problem 9: Image Captioning [9 pts]

Adhitya wants to train a model for image captioning. He plans to train a decoder-only transformer model that generates a natural language description by attending to visual features extracted from an input image. He has access to a pretrained ViT model to extract image features. Adhitya has written most of the network architecture and needs a little help with the cross-attention layer in the decoder (the cross-attention layer is used to attend to the image features).

b) When you are predicting the first word, what are the shapes of the output matrices obtained after multiplying W_Q, W_K, W_V with their appropriate inputs.

Problem 9: Image Captioning [9 pts]

Adhitya wants to train a model for image captioning. He plans to train a decoder-only transformer model that generates a natural language description by attending to visual features extracted from an input image. He has access to a pretrained ViT model to extract image features. Adhitya has written most of the network architecture and needs a little help with the cross-attention layer in the decoder (the cross-attention layer is used to attend to the image features).

b) When you are predicting the first word, what are the shapes of the output matrices obtained after multiplying W_Q, W_K, W_V with their appropriate inputs.

The shapes of the output matrices after multiplying W_Q, W_K, W_V with their appropriate inputs are:

$Q : (1, 128)$ - The query matrix is obtained by multiplying the decoder's input embedding (or the previous layer's output) with W_Q . $K : (k, 128)$ - The key matrix is obtained by multiplying the ViT's output features with W_K . $V : (k, 128)$ - The value matrix is obtained by multiplying the ViT's output features with W_V .

Problem 9: Image Captioning [9 pts]

Adhitya wants to train a model for image captioning. He plans to train a decoder-only transformer model that generates a natural language description by attending to visual features extracted from an input image. He has access to a pretrained ViT model to extract image features. Adhitya has written most of the network architecture and needs a little help with the cross-attention layer in the decoder (the cross-attention layer is used to attend to the image features).

c) Suppose the ViT model produced 64-dimensional features for every patch (instead of 128-dimensional features), what would you have to modify to attend to these image features?

Problem 9: Image Captioning [9 pts]

Adhitya wants to train a model for image captioning. He plans to train a decoder-only transformer model that generates a natural language description by attending to visual features extracted from an input image. He has access to a pretrained ViT model to extract image features. Adhitya has written most of the network architecture and needs a little help with the cross-attention layer in the decoder (the cross-attention layer is used to attend to the image features).

c) Suppose the ViT model produced 64-dimensional features for every patch (instead of 128-dimensional features), what would you have to modify to attend to these image features?

Add a linear layer to project the embeddings to 128 dimensions.