



Transformers

CISC 7026 - Introduction to Deep Learning

Steven Morad

University of Macau

Review	2
Going Deeper	17
Transformers	34
Positional Encoding	38
Applications	59
Text Transformers	61
Image Transformers	66
Unsupervised Training	69
Closing Remarks	87
Course Evaluation	94

Review

Review

Last time, we derived various forms of **attention**

We started with composite memory

Review

Last time, we derived various forms of **attention**

We started with composite memory

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i$$

Review

Last time, we derived various forms of **attention**

We started with composite memory

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i$$

Given large enough T , we will eventually run out of storage space

Review

Last time, we derived various forms of **attention**

We started with composite memory

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i$$

Given large enough T , we will eventually run out of storage space

The sum is a **lossy** operation that can store a limited amount of information

Review

Last time, we derived various forms of **attention**

We started with composite memory

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i$$

Given large enough T , we will eventually run out of storage space

The sum is a **lossy** operation that can store a limited amount of information

Review

So we introduced a forgetting term γ

Review

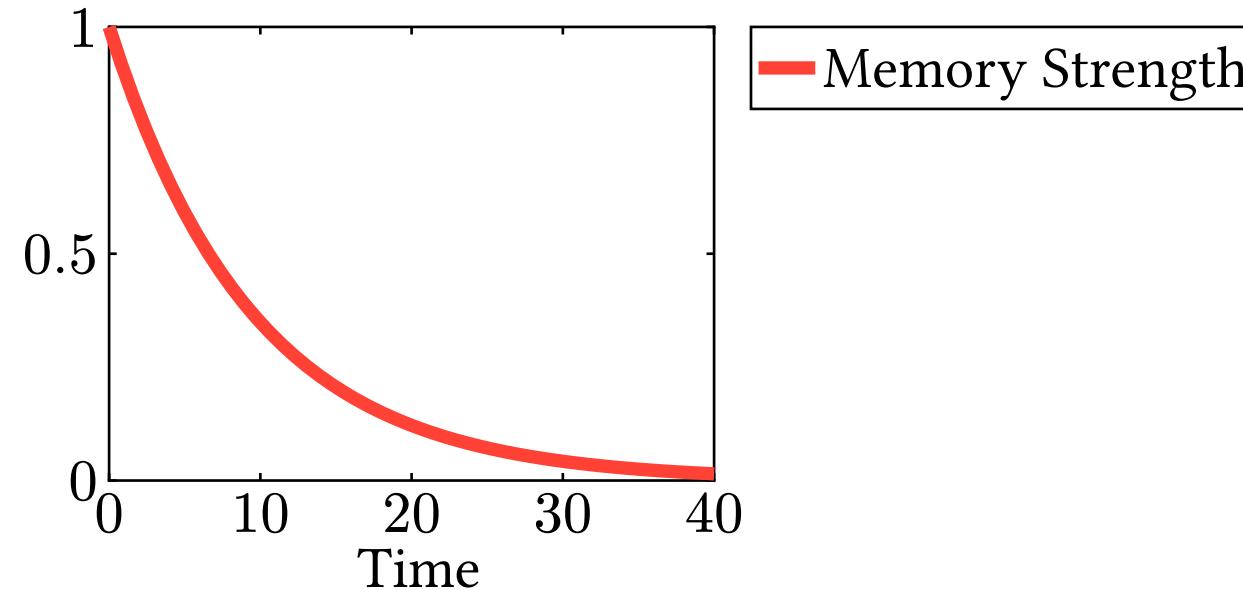
So we introduced a forgetting term γ

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \gamma^{T-i} \cdot \boldsymbol{\theta}^\top \overline{\mathbf{x}}_i$$

Review

So we introduced a forgetting term γ

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^T \gamma^{T-i} \cdot \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i$$



Review

We went to a party and the forgetting seemed ok

Review

We went to a party and the forgetting seemed ok



Review

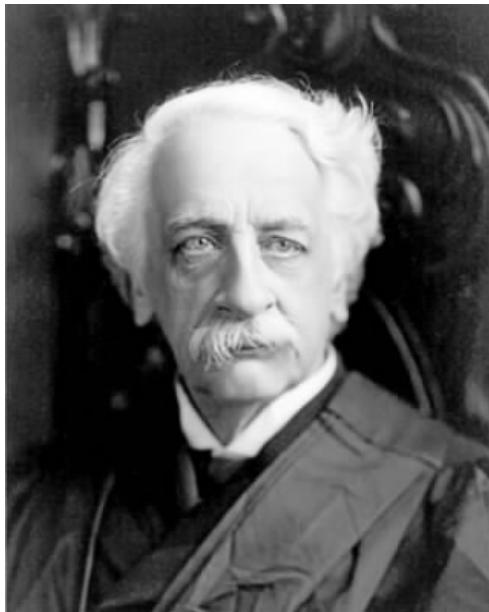
We went to a party and the forgetting seemed ok



10 PM

Review

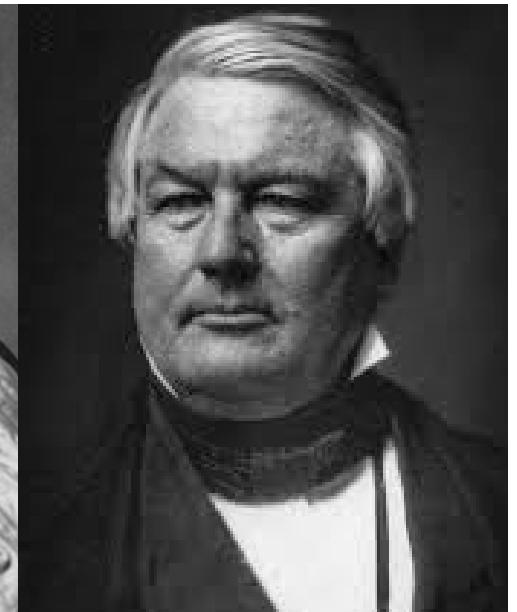
We went to a party and the forgetting seemed ok



10 PM

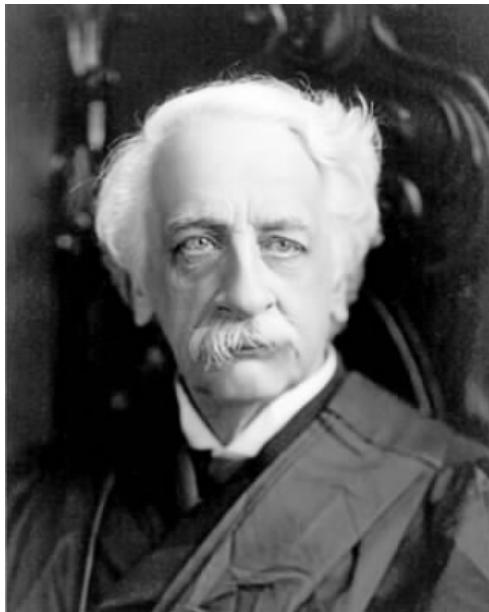


11 PM



Review

We went to a party and the forgetting seemed ok



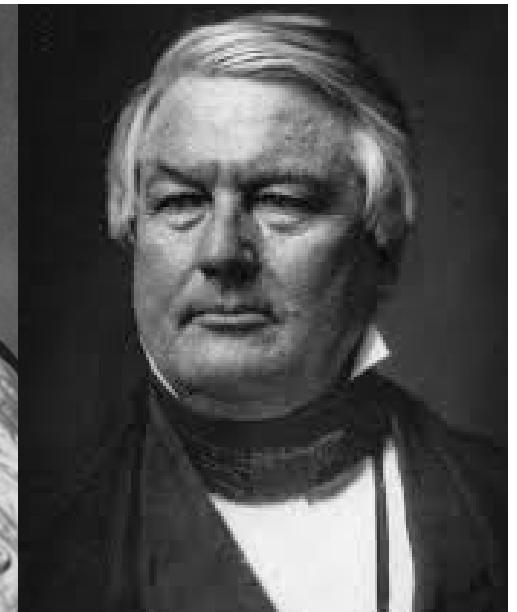
10 PM



11 PM

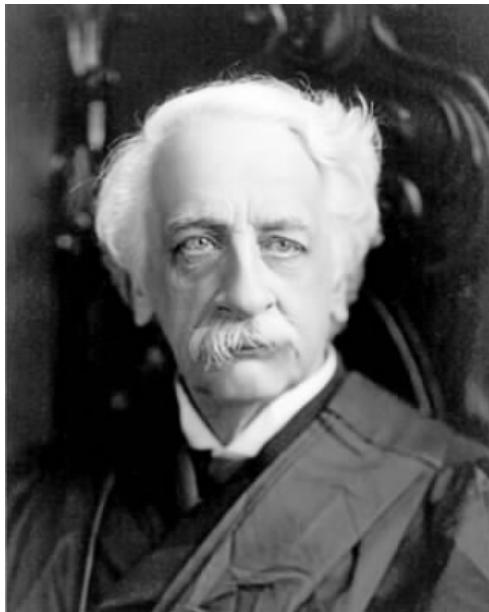


12 AM



Review

We went to a party and the forgetting seemed ok



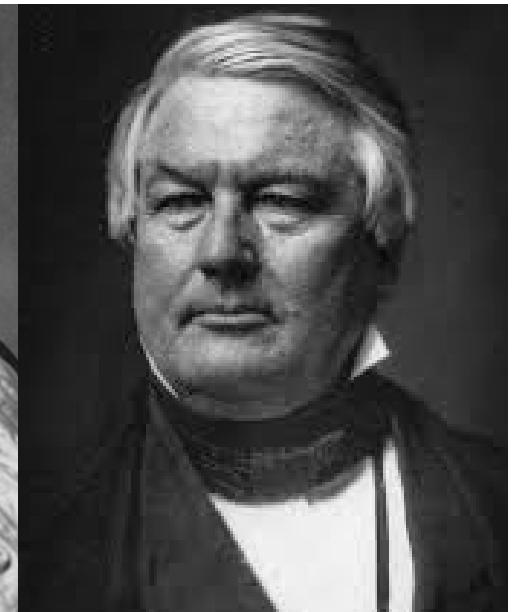
10 PM



11 PM



12 AM



1 AM

Review



Review



$$\gamma^3 \theta^\top \bar{x}_1$$

Review



$$\gamma^3 \theta^\top \bar{x}_1$$

$$\gamma^2 \theta^\top \bar{x}_2$$

Review



$$\gamma^3 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_1$$

$$\gamma^2 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_2$$

$$\gamma^1 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_3$$

Review



$$\gamma^3 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_1$$

$$\gamma^2 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_2$$

$$\gamma^1 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_3$$

$$\gamma^0 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_4$$

Review



$$\gamma^3 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_1$$

$$\gamma^2 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_2$$

$$\gamma^1 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_3$$

$$\gamma^0 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_4$$

Review

But we encountered problems when Taylor Swift arrived at the party

Review

But we encountered problems when Taylor Swift arrived at the party



Review

But we encountered problems when Taylor Swift arrived at the party



$$\gamma^4 \theta^\top \bar{x}_1$$

$$\gamma^3 \theta^\top \bar{x}_2$$

$$\gamma^2 \theta^\top \bar{x}_3$$

$$\gamma^1 \theta^\top \bar{x}_4$$

$$\gamma^0 \theta^\top \bar{x}_5$$

Review

But we encountered problems when Taylor Swift arrived at the party



$$\gamma^4 \theta^\top \bar{x}_1$$

$$\gamma^3 \theta^\top \bar{x}_2$$

$$\gamma^2 \theta^\top \bar{x}_3$$

$$\gamma^1 \theta^\top \bar{x}_4$$

$$\gamma^0 \theta^\top \bar{x}_5$$

Review



$$\gamma^4 \theta^\top \bar{x}_1$$

$$\gamma^3 \theta^\top \bar{x}_2$$

$$\gamma^2 \theta^\top \bar{x}_3$$

$$\gamma^1 \theta^\top \bar{x}_4$$

$$\gamma^0 \theta^\top \bar{x}_5$$

Review



$$\gamma^4 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_1$$

$$\gamma^3 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_2$$

$$\gamma^2 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_3$$

$$\gamma^1 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_4$$

$$\gamma^0 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_5$$

With composite memory, we forget Taylor Swift!

Review



$$\gamma^4 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_1$$

$$\gamma^3 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_2$$

$$\gamma^2 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_3$$

$$\gamma^1 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_4$$

$$\gamma^0 \boldsymbol{\theta}^\top \bar{\boldsymbol{x}}_5$$

With composite memory, we forget Taylor Swift!

Our model of human memory was incomplete

Review

So we introduced **attention**

Review

So we introduced **attention**

The attention we pay to person i is

$$\lambda \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda \right)_i = \text{softmax} \left(\begin{bmatrix} \boldsymbol{\theta}_\lambda^\top \bar{\mathbf{x}}_1 \\ \vdots \\ \boldsymbol{\theta}_\lambda^\top \bar{\mathbf{x}}_T \end{bmatrix} \right)_i = \frac{\exp(\boldsymbol{\theta}_\lambda^\top \bar{\mathbf{x}}_i)}{\sum_{j=1}^T \exp(\boldsymbol{\theta}_\lambda^\top \bar{\mathbf{x}}_j)}$$

Review



Review



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \cdot \theta^\top \bar{x}_1$$

Review



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_2 \\ \cdot \theta^\top \bar{x}_1 \quad \cdot \theta^\top \bar{x}_2$$

Review



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_2 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_3 \\ \cdot \theta^\top \bar{x}_1 \quad \cdot \theta^\top \bar{x}_2 \quad \cdot \theta^\top \bar{x}_3$$

Review



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_2 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_3 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_4 \\ \cdot \theta^\top \bar{x}_1 \quad \cdot \theta^\top \bar{x}_2 \quad \cdot \theta^\top \bar{x}_3 \quad \cdot \theta^\top \bar{x}_4$$

Review



$$\lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_1 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_2 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_3 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_4 \lambda \left(\begin{bmatrix} x_1 \\ \vdots \\ x_5 \end{bmatrix}, \theta_\lambda \right)_5 \\ \cdot \theta^\top \bar{x}_1 \quad \cdot \theta^\top \bar{x}_2 \quad \cdot \theta^\top \bar{x}_3 \quad \cdot \theta^\top \bar{x}_4 \quad \cdot \theta^\top \bar{x}_5$$

Review



Review



$$0.70 \cdot \theta^\top \bar{x}_1$$

Review



$$0.70 \cdot \theta^\top \bar{x}_1 - 0.04 \cdot \theta^\top \bar{x}_2$$

Review



$$0.70 \cdot \theta^\top \bar{x}_1 \quad 0.04 \cdot \theta^\top \bar{x}_2 \quad 0.03 \cdot \theta^\top \bar{x}_3$$

Review



$$0.70 \cdot \theta^\top \bar{x}_1 \quad 0.04 \cdot \theta^\top \bar{x}_2 \quad 0.03 \cdot \theta^\top \bar{x}_3 \quad 0.20 \cdot \theta^\top \bar{x}_4$$

Review



$$0.70 \cdot \theta^\top \bar{x}_1 \quad 0.04 \cdot \theta^\top \bar{x}_2 \quad 0.03 \cdot \theta^\top \bar{x}_3 \quad 0.20 \cdot \theta^\top \bar{x}_4 \quad 0.03 \cdot \theta^\top \bar{x}_5$$

Review



$$0.70 \cdot \theta^\top \bar{x}_1 \quad 0.04 \cdot \theta^\top \bar{x}_2 \quad 0.03 \cdot \theta^\top \bar{x}_3 \quad 0.20 \cdot \theta^\top \bar{x}_4 \quad 0.03 \cdot \theta^\top \bar{x}_5$$

$$0.70 + 0.04 + 0.03 + 0.20 + 0.03 = 1.0$$

Review

Then, we introduced **keys** and **queries**

Review

Then, we introduced **keys** and **queries**

Query: Which person will help me on my exam?

Review

Then, we introduced **keys** and **queries**

Query: Which person will help me on my exam?



Review

Then, we introduced **keys** and **queries**

Query: Which person will help me on my exam?



Musician

Lawyer

Shopkeeper

Chef

Scientist

Review

Then, we introduced **keys** and **queries**

Query: Which person will help me on my exam?

Musician

Lawyer

Shopkeeper

Chef

Scientist



Review

$$\mathbf{K} = [k_1 \ k_2 \ \dots \ k_T] = [\theta_K^\top x_1 \ \theta_K^\top x_2 \ \dots \ \theta_K^\top x_T], \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

Review

$$\mathbf{K} = [k_1 \ k_2 \ \dots \ k_T] = [\theta_K^\top x_1 \ \theta_K^\top x_2 \ \dots \ \theta_K^\top x_T], \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

$$\mathbf{q} = \theta_Q^\top x_q, \quad \theta_Q \in \mathbb{R}^{d_x \times d_h}, \quad \mathbf{q} \in \mathbb{R}^{d_h}$$

Review

$$\mathbf{K} = [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\boldsymbol{\theta}_K^\top \mathbf{x}_1 \ \boldsymbol{\theta}_K^\top \mathbf{x}_2 \ \dots \ \boldsymbol{\theta}_K^\top \mathbf{x}_T], \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

$$\mathbf{q} = \boldsymbol{\theta}_Q^\top \mathbf{x}_q, \quad \boldsymbol{\theta}_Q \in \mathbb{R}^{d_x \times d_h}, \quad \mathbf{q} \in \mathbb{R}^{d_h}$$

$$\begin{aligned} \lambda\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta} \right) &= \text{softmax}(\mathbf{q}^\top \mathbf{K}) = \text{softmax}(\mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T]) \\ &= \text{softmax}([\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T]) \end{aligned}$$

Review

$$\mathbf{K} = [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\theta_K^\top \mathbf{x}_1 \ \theta_K^\top \mathbf{x}_2 \ \dots \ \theta_K^\top \mathbf{x}_T], \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

$$\mathbf{q} = \theta_Q^\top \mathbf{x}_q, \quad \theta_Q \in \mathbb{R}^{d_x \times d_h}, \quad \mathbf{q} \in \mathbb{R}^{d_h}$$

$$\begin{aligned} \lambda \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta \right) &= \text{softmax}(\mathbf{q}^\top \mathbf{K}) = \text{softmax}(\mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T]) \\ &= \text{softmax}([\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T]) \end{aligned}$$

We call this **dot-product attention**

Review

$$\mathbf{K} = [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\theta_K^\top \mathbf{x}_1 \ \theta_K^\top \mathbf{x}_2 \ \dots \ \theta_K^\top \mathbf{x}_T], \quad \mathbf{K} \in \mathbb{R}^{d_h \times T}$$

$$\mathbf{q} = \theta_Q^\top \mathbf{x}_q, \quad \theta_Q \in \mathbb{R}^{d_x \times d_h}, \quad \mathbf{q} \in \mathbb{R}^{d_h}$$

$$\begin{aligned} \lambda \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta \right) &= \text{softmax}(\mathbf{q}^\top \mathbf{K}) = \text{softmax}(\mathbf{q}^\top [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T]) \\ &= \text{softmax}([\mathbf{q}^\top \mathbf{k}_1 \ \mathbf{q}^\top \mathbf{k}_2 \ \dots \ \mathbf{q}^\top \mathbf{k}_T]) \end{aligned}$$

We call this **dot-product attention**

Then we add attention back to the composite model

Review

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \sum_{i=1}^T \theta^\top x_i \cdot \lambda\left(\begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta_\lambda\right)_i$$

Review

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}^\top \mathbf{x}_i \cdot \lambda\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda\right)_i$$

We relabel $\boldsymbol{\theta}$ to $\boldsymbol{\theta}_V$

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}_{\textcolor{red}{V}}^\top \mathbf{x}_i \cdot \lambda\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda\right)_i$$

Review

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}^\top \mathbf{x}_i \cdot \lambda\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda\right)_i$$

We relabel $\boldsymbol{\theta}$ to $\boldsymbol{\theta}_V$

$$f\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}\right) = \sum_{i=1}^T \boldsymbol{\theta}_{\textcolor{red}{V}}^\top \mathbf{x}_i \cdot \lambda\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta}_\lambda\right)_i$$

In dot-product attention, we call $\boldsymbol{\theta}_V^\top \mathbf{x}_i$ the **value**

Review

In **dot product self attention** we create queries for all inputs

Review

In **dot product self attention** we create queries for all inputs

$$\mathbf{Q} = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_T] = [\theta_Q^\top \mathbf{x}_1 \ \theta_Q^\top \mathbf{x}_2 \ \dots \ \theta_Q^\top \mathbf{x}_T], \quad \mathbf{Q} \in \mathbb{R}^{T \times d_h}$$

Review

In **dot product self attention** we create queries for all inputs

$$\mathbf{Q} = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_T] = [\theta_Q^\top \mathbf{x}_1 \ \theta_Q^\top \mathbf{x}_2 \ \dots \ \theta_Q^\top \mathbf{x}_T], \quad \mathbf{Q} \in \mathbb{R}^{T \times d_h}$$

Review

Attention looks messy, but we can rewrite it in matrix form

Review

Attention looks messy, but we can rewrite it in matrix form

$$Q = [q_1 \ q_2 \ \dots \ q_T] = [\theta_Q^\top x_1 \ \theta_Q^\top x_2 \ \dots \ \theta_Q^\top x_T]$$

$$K = [k_1 \ k_2 \ \dots \ k_T] = [\theta_K^\top x_1 \ \theta_K^\top x_2 \ \dots \ \theta_K^\top x_T]$$

$$V = [v_1 \ v_2 \ \dots \ v_T] = [\theta_V^\top x_1 \ \theta_V^\top x_2 \ \dots \ \theta_V^\top x_T]$$

Review

Attention looks messy, but we can rewrite it in matrix form

$$\mathbf{Q} = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_T] = [\theta_Q^\top \mathbf{x}_1 \ \theta_Q^\top \mathbf{x}_2 \ \dots \ \theta_Q^\top \mathbf{x}_T]$$

$$\mathbf{K} = [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\theta_K^\top \mathbf{x}_1 \ \theta_K^\top \mathbf{x}_2 \ \dots \ \theta_K^\top \mathbf{x}_T]$$

$$\mathbf{V} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_T] = [\theta_V^\top \mathbf{x}_1 \ \theta_V^\top \mathbf{x}_2 \ \dots \ \theta_V^\top \mathbf{x}_T]$$

$$\text{attn}\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}}\right)\mathbf{V}$$

Review

Attention looks messy, but we can rewrite it in matrix form

$$\mathbf{Q} = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_T] = [\theta_Q^\top \mathbf{x}_1 \ \theta_Q^\top \mathbf{x}_2 \ \dots \ \theta_Q^\top \mathbf{x}_T]$$

$$\mathbf{K} = [\mathbf{k}_1 \ \mathbf{k}_2 \ \dots \ \mathbf{k}_T] = [\theta_K^\top \mathbf{x}_1 \ \theta_K^\top \mathbf{x}_2 \ \dots \ \theta_K^\top \mathbf{x}_T]$$

$$\mathbf{V} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_T] = [\theta_V^\top \mathbf{x}_1 \ \theta_V^\top \mathbf{x}_2 \ \dots \ \theta_V^\top \mathbf{x}_T]$$

$$\text{attn}\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}}\right)\mathbf{V}$$

This operation powers today's biggest models

Going Deeper

Going Deeper

Modern transformers can be very deep

Going Deeper

Modern transformers can be very deep

Very deep networks require two new training tricks

Going Deeper

Modern transformers can be very deep

Very deep networks require two new training tricks

We must understand these tricks before implementing the transformer

Going Deeper

Modern transformers can be very deep

Very deep networks require two new training tricks

We must understand these tricks before implementing the transformer

Trick 1: Residual connections

Going Deeper

Modern transformers can be very deep

Very deep networks require two new training tricks

We must understand these tricks before implementing the transformer

Trick 1: Residual connections

Trick 2: Layer normalization

Going Deeper

Modern transformers can be very deep

Very deep networks require two new training tricks

We must understand these tricks before implementing the transformer

Trick 1: Residual connections

Trick 2: Layer normalization

We will start with the **residual connection**

Going Deeper

Remember that a two-layer MLP is a universal function approximator

Going Deeper

Remember that a two-layer MLP is a universal function approximator

$$| f(\mathbf{x}, \boldsymbol{\theta}) - g(\mathbf{x}) | < \varepsilon$$

Going Deeper

Remember that a two-layer MLP is a universal function approximator

$$| f(\mathbf{x}, \boldsymbol{\theta}) - g(\mathbf{x}) | < \varepsilon$$

This is only as the width of the network goes to infinity

Going Deeper

Remember that a two-layer MLP is a universal function approximator

$$| f(\mathbf{x}, \boldsymbol{\theta}) - g(\mathbf{x}) | < \varepsilon$$

This is only as the width of the network goes to infinity

For certain problems, we need deeper networks

Going Deeper

Remember that a two-layer MLP is a universal function approximator

$$| f(\mathbf{x}, \boldsymbol{\theta}) - g(\mathbf{x}) | < \varepsilon$$

This is only as the width of the network goes to infinity

For certain problems, we need deeper networks

But there is a limit!

Going Deeper

Remember that a two-layer MLP is a universal function approximator

$$| f(\mathbf{x}, \boldsymbol{\theta}) - g(\mathbf{x}) | < \varepsilon$$

This is only as the width of the network goes to infinity

For certain problems, we need deeper networks

But there is a limit!

Going Deeper

Making the network too deep can hurt performance

Going Deeper

Making the network too deep can hurt performance

The theory is that the input information is **lost** somewhere in the network

Going Deeper

Making the network too deep can hurt performance

The theory is that the input information is **lost** somewhere in the network

$$\mathbf{y} = f_k(\dots f_2(f_1(\mathbf{x}, \boldsymbol{\theta}_1), \boldsymbol{\theta}_2), \dots, \boldsymbol{\theta}_k)$$

Going Deeper

$$\mathbf{y} = f_k(\dots f_2(f_1(\mathbf{x}, \boldsymbol{\theta}_1), \boldsymbol{\theta}_2), \dots, \boldsymbol{\theta}_k)$$

Claim: If the input information is present in all layers of the network, then we should be able to learn the identity function $f(x) = x$

Going Deeper

$$y = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Claim: If the input information is present in all layers of the network, then we should be able to learn the identity function $f(x) = x$

$$x = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Question: We have seen a similar model, what was it?

Going Deeper

$$\mathbf{y} = f_k(\dots f_2(f_1(\mathbf{x}, \boldsymbol{\theta}_1), \boldsymbol{\theta}_2), \dots, \boldsymbol{\theta}_k)$$

Claim: If the input information is present in all layers of the network, then we should be able to learn the identity function $f(\mathbf{x}) = \mathbf{x}$

$$\mathbf{x} = f_k(\dots f_2(f_1(\mathbf{x}, \boldsymbol{\theta}_1), \boldsymbol{\theta}_2), \dots, \boldsymbol{\theta}_k)$$

Question: We have seen a similar model, what was it?

Answer: Autoencoder! But it was just two functions, not k functions

Going Deeper

$$y = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Claim: If the input information is present in all layers of the network, then we should be able to learn the identity function $f(x) = x$

$$x = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Question: We have seen a similar model, what was it?

Answer: Autoencoder! But it was just two functions, not k functions

Question: Do you agree or disagree with the claim?

Going Deeper

$$y = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Claim: If the input information is present in all layers of the network, then we should be able to learn the identity function $f(x) = x$

$$x = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Question: We have seen a similar model, what was it?

Answer: Autoencoder! But it was just two functions, not k functions

Question: Do you agree or disagree with the claim?

<https://colab.research.google.com/drive/1qVlbQKpTuBYIa7FvC4IH-kJq-E0jmc0d#scrollTo=bg74S-AvbmJz>

Going Deeper

$$x = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Going Deeper

$$x = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Very deep networks struggle to learn the identity function

Going Deeper

$$x = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Very deep networks struggle to learn the identity function

If the input information is available, then learning the identity function should be very easy!

Going Deeper

$$x = f_k(\dots f_2(f_1(x, \theta_1), \theta_2), \dots, \theta_k)$$

Very deep networks struggle to learn the identity function

If the input information is available, then learning the identity function should be very easy!

Question: How can we prevent the input from getting lost?

Going Deeper

We can feed the input to each layer

Going Deeper

We can feed the input to each layer

The first approach is called the **DenseNet** approach

Going Deeper

We can feed the input to each layer

The first approach is called the **DenseNet** approach

$$z_1 = f_1(x, \theta_1)$$

$$z_2 = f_2\left(\begin{bmatrix} x \\ z_1 \end{bmatrix}, \theta_2\right)$$

⋮

$$z_k = f_k\left(\begin{bmatrix} x \\ z_1 \\ \vdots \\ z_{k-1} \end{bmatrix}, \theta_k\right)$$

Going Deeper

$$z_1 = f_1(x, \theta_1)$$

$$z_2 = f_2\left(\begin{bmatrix} x \\ z_1 \end{bmatrix}, \theta_2\right)$$

⋮

$$z_k = f_k\left(\begin{bmatrix} x \\ z_1 \\ \vdots \\ z_{k-1} \end{bmatrix}, \theta_k\right)$$

Going Deeper

$$z_1 = f_1(x, \theta_1)$$

$$z_2 = f_2\left(\begin{bmatrix} x \\ z_1 \end{bmatrix}, \theta_2\right)$$

⋮

$$z_k = f_k\left(\begin{bmatrix} x \\ z_1 \\ \vdots \\ z_{k-1} \end{bmatrix}, \theta_k\right)$$

Question: Any issues with the DenseNet approach?

Going Deeper

$$z_1 = f_1(x, \theta_1)$$

$$z_2 = f_2\left(\begin{bmatrix} x \\ z_1 \end{bmatrix}, \theta_2\right)$$

⋮

$$z_k = f_k\left(\begin{bmatrix} x \\ z_1 \\ \vdots \\ z_{k-1} \end{bmatrix}, \theta_k\right)$$

Question: Any issues with the DenseNet approach?

Answer: Very deep networks require too many parameters!

Going Deeper

The next method is called the **ResNet** approach

$$z_1 = f_1(x, \theta_1)$$

$$z_2 = f_2(x, \theta_2) + z_1$$

$$z_3 = f_2(x, \theta_3) + z_2$$

⋮

$$z_k = f_k(x, \theta_k) + z_{k-1}$$

Going Deeper

The next method is called the **ResNet** approach

$$z_1 = f_1(x, \theta_1)$$

$$z_2 = f_2(x, \theta_2) + z_1$$

$$z_3 = f_2(x, \theta_3) + z_2$$

⋮

$$z_k = f_k(x, \theta_k) + z_{k-1}$$

This allows information to flow around the layers

Going Deeper

The next method is called the **ResNet** approach

$$z_1 = f_1(x, \theta_1)$$

$$z_2 = f_2(x, \theta_2) + x$$

$$z_3 = f_2(x, \theta_3) + z_2$$

⋮

$$z_k = f_k(x, \theta_k) + z_{k-1}$$

This allows information to flow around the layers

Fewer parameters than a DenseNet, but performs slightly worse

Going Deeper

The next method is called the **ResNet** approach

$$z_1 = f_1(x, \theta_1)$$

$$z_2 = f_2(x, \theta_2) + z_1$$

$$z_3 = f_2(x, \theta_3) + z_2$$

⋮

$$z_k = f_k(x, \theta_k) + z_{k-1}$$

This allows information to flow around the layers

Fewer parameters than a DenseNet, but performs slightly worse

For very deep networks, we use ResNets over DenseNets

Going Deeper

We call $f(x) + x$ a **residual connection**

Going Deeper

We call $f(x) + x$ a **residual connection**

$$\underbrace{f(x)}_{\text{Residual}} + x$$

Going Deeper

We call $f(x) + x$ a **residual connection**

$$\underbrace{f(x)}_{\text{Residual}} + x$$

Instead of learning how to change x , f learns the residual of x

Going Deeper

We call $f(x) + x$ a **residual connection**

$$\underbrace{f(x)}_{\text{Residual}} + x$$

Instead of learning how to change x , f learns the residual of x

Think of the residual as a small change to f

$$f(x, \theta) + x = \varepsilon + x$$

Going Deeper

We call $f(x) + x$ a **residual connection**

$$\underbrace{f(x)}_{\text{Residual}} + x$$

Instead of learning how to change x , f learns the residual of x

Think of the residual as a small change to f

$$f(x, \theta) + x = \varepsilon + x$$

This prevents the input x from getting lost in very deep networks

Going Deeper

The second trick is called **layer normalization**

Going Deeper

The second trick is called **layer normalization**

With parameter initialization and weight decay, we ensure the parameters are fairly small

Going Deeper

The second trick is called **layer normalization**

With parameter initialization and weight decay, we ensure the parameters are fairly small

But we can still have very small or very large outputs from each layer

Going Deeper

The second trick is called **layer normalization**

With parameter initialization and weight decay, we ensure the parameters are fairly small

But we can still have very small or very large outputs from each layer

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{d_x} \theta_{1,i} x_i$$

Going Deeper

The second trick is called **layer normalization**

With parameter initialization and weight decay, we ensure the parameters are fairly small

But we can still have very small or very large outputs from each layer

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{d_x} \theta_{1,i} x_i$$

Question: If all $x_i = 1$, $\theta_{1,i} = 0.01$ and $d_x = 1000$, what is the output?

Going Deeper

$$f_1(\boldsymbol{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{d_x} \theta_{1,i} x_i$$

Going Deeper

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{d_x} \theta_{1,i} x_i$$

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{1000} 0.01 \cdot 1 = 10$$

Going Deeper

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{d_x} \theta_{1,i} x_i$$

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{1000} 0.01 \cdot 1 = 10$$

What if we add another layer with the same d_x and θ ?

Going Deeper

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{d_x} \theta_{1,i} x_i$$

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{1000} 0.01 \cdot 1 = 10$$

What if we add another layer with the same d_x and θ ?

$$f_2(\mathbf{z}, \boldsymbol{\theta}_2) = \sum_{i=1}^{1000} 0.01 \cdot 10 = 100$$

Going Deeper

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{d_x} \theta_{1,i} x_i$$

$$f_1(\mathbf{x}, \boldsymbol{\theta}_1) = \sum_{i=1}^{1000} 0.01 \cdot 1 = 10$$

What if we add another layer with the same d_x and θ ?

$$f_2(\mathbf{z}, \boldsymbol{\theta}_2) = \sum_{i=1}^{1000} 0.01 \cdot 10 = 100$$

Question: What is the problem?

Going Deeper

Let us look at the gradient

Going Deeper

Let us look at the gradient

$$\nabla_{\theta_1} f_2(f_1(x, \theta_1), \theta_2) =$$

Going Deeper

Let us look at the gradient

$$\nabla_{\theta_1} f_2(f_1(x, \theta_1), \theta_2) = \nabla_{\theta_1}[f_2](f_1(x, \theta_1)) \cdot \nabla_{\theta_1}[f_1](x, \theta_1)$$

Going Deeper

Let us look at the gradient

$$\begin{aligned}\nabla_{\theta_1} f_2(f_1(x, \theta_1), \theta_2) &= \nabla_{\theta_1}[f_2](f_1(x, \theta_1)) \cdot \nabla_{\theta_1}[f_1](x, \theta_1) \\ &\approx 100 \cdot 10\end{aligned}$$

Can cause exploding or vanishing gradient

Going Deeper

Let us look at the gradient

$$\begin{aligned}\nabla_{\theta_1} f_2(f_1(x, \theta_1), \theta_2) &= \nabla_{\theta_1}[f_2](f_1(x, \theta_1)) \cdot \nabla_{\theta_1}[f_1](x, \theta_1) \\ &\approx 100 \cdot 10\end{aligned}$$

Can cause exploding or vanishing gradient

Deeper network \Rightarrow worse exploding/vanishing issues

Going Deeper

Let us look at the gradient

$$\begin{aligned}\nabla_{\theta_1} f_2(f_1(x, \theta_1), \theta_2) &= \nabla_{\theta_1}[f_2](f_1(x, \theta_1)) \cdot \nabla_{\theta_1}[f_1](x, \theta_1) \\ &\approx 100 \cdot 10\end{aligned}$$

Can cause exploding or vanishing gradient

Deeper network \Rightarrow worse exploding/vanishing issues

Question: What can we do?

Going Deeper

We can use **layer normalization**

Going Deeper

We can use **layer normalization**

First, layer normalization **centers** the output of the layer

Going Deeper

We can use **layer normalization**

First, layer normalization **centers** the output of the layer

$$\mu = d_y \sum_{i=1}^{d_y} f(x, \theta)_i$$

$$f(x, \theta) - \mu$$

Question: What does this do?

Going Deeper

We can use **layer normalization**

First, layer normalization **centers** the output of the layer

$$\mu = d_y \sum_{i=1}^{d_y} f(x, \theta)_i$$

$$f(x, \theta) - \mu$$

Question: What does this do?

Answer: Creates zero-mean output (both positive and negative values)

Going Deeper

$$\mu = d_y \sum_{i=1}^{d_y} f(\mathbf{x}, \theta)_i \quad f(\mathbf{x}, \theta) - \mu$$

Going Deeper

$$\mu = d_y \sum_{i=1}^{d_y} f(\mathbf{x}, \theta)_i \quad f(\mathbf{x}, \theta) - \mu$$

Then, layer normalization **rescales** the outputs

Going Deeper

$$\mu = d_y \sum_{i=1}^{d_y} f(\mathbf{x}, \boldsymbol{\theta})_i \quad f(\mathbf{x}, \boldsymbol{\theta}) - \mu$$

Then, layer normalization **rescales** the outputs

$$\sigma = \frac{\sqrt{\sum_{i=1}^{d_y} f(\mathbf{x}, \boldsymbol{\theta}_i - \mu)^2}}{d_y}$$

Going Deeper

$$\mu = d_y \sum_{i=1}^{d_y} f(\mathbf{x}, \boldsymbol{\theta})_i \quad f(\mathbf{x}, \boldsymbol{\theta}) - \mu$$

Then, layer normalization **rescales** the outputs

$$\sigma = \frac{\sqrt{\sum_{i=1}^{d_y} f(\mathbf{x}, \boldsymbol{\theta}_i - \mu)^2}}{d_y}$$

$$\text{LN}(f(\mathbf{x}, \boldsymbol{\theta})) = \frac{f(\mathbf{x}, \boldsymbol{\theta}) - \mu}{\sigma}$$

Going Deeper

$$\mu = d_y \sum_{i=1}^{d_y} f(\mathbf{x}, \boldsymbol{\theta})_i \quad f(\mathbf{x}, \boldsymbol{\theta}) - \mu$$

Then, layer normalization **rescales** the outputs

$$\sigma = \frac{\sqrt{\sum_{i=1}^{d_y} f(\mathbf{x}, \boldsymbol{\theta}_i - \mu)^2}}{d_y}$$

$$\text{LN}(f(\mathbf{x}, \boldsymbol{\theta})) = \frac{f(\mathbf{x}, \boldsymbol{\theta}) - \mu}{\sigma}$$

Now, the output of the layer is normally distributed

Going Deeper

If the output is normally distributed:

Going Deeper

If the output is normally distributed:

- 99.7% of outputs $\in [-3, 3]$

Going Deeper

If the output is normally distributed:

- 99.7% of outputs $\in [-3, 3]$
- 99.99% of outputs $\in [-4, 4]$

Going Deeper

If the output is normally distributed:

- 99.7% of outputs $\in [-3, 3]$
- 99.99% of outputs $\in [-4, 4]$
- 99.9999% of outputs $\in [-5, 5]$

Going Deeper

If the output is normally distributed:

- 99.7% of outputs $\in [-3, 3]$
- 99.99% of outputs $\in [-4, 4]$
- 99.9999% of outputs $\in [-5, 5]$

This helps prevent vanishing and exploding gradients

Going Deeper

Now, let's combine residual connections and layer norm and try our very deep network again

<https://colab.research.google.com/drive/1qVlbQKpTuBYIa7FvC4IH-kJq-E0jmc0d#scrollTo=iQtXjGYiz5CD>

Transformers

Transformers

Now we have everything we need to implement a transformer

Transformers

Now we have everything we need to implement a transformer

- Attention

Transformers

Now we have everything we need to implement a transformer

- Attention
- MLP

Transformers

Now we have everything we need to implement a transformer

- Attention
- MLP
- Residual connections

Transformers

Now we have everything we need to implement a transformer

- Attention
- MLP
- Residual connections
- Layer normalization

Transformers

Now we have everything we need to implement a transformer

- Attention
- MLP
- Residual connections
- Layer normalization

A transformer consists of many **transformer layers**

Transformers

```
class TransformerLayer(nn.Module):
    def __init__(self):
        self.attn = Attention()
        self.mlp = Sequential(
            Linear(d_h, d_h), LeakyReLU(), Linear(d_h, d_h))
        self.norm = nn.LayerNorm(
            d_h, elementwise_affine=False)

    def forward(self, x):
        # Residual connection and layer norm
        x = self.norm(self.attn(x) + x)
        x = self.norm(self.mlp(x) + x)
        return x
```

Transformers

```
class Transformer(nn.Module):
    def __init__(self):
        self.block1 = TransformerLayer()
        self.block2 = TransformerLayer()
        self.block3 = TransformerLayer()

    def forward(self, x):
        x = self.block1(x)
        x = self.block2(x)
        x = self.block3(x)
        return x
```

Transformers

```
class Transformer(nn.Module):
    def __init__(self):
        self.block1 = TransformerLayer()
        self.block2 = TransformerLayer()
        self.block3 = TransformerLayer()

    def forward(self, x):
        x = self.block1(x)
        x = self.block2(x)
        x = self.block3(x)
        return x
```

Question: What are the input/output shapes of the transformer?

Transformers

```
class Transformer(nn.Module):  
    def __init__(self):  
        self.block1 = TransformerLayer()  
        self.block2 = TransformerLayer()  
        self.block3 = TransformerLayer()  
  
    def forward(self, x):  
        x = self.block1(x)  
        x = self.block2(x)  
        x = self.block3(x)  
        return x
```

Question: What are the input/output shapes of the transformer?

Answer: $f : \mathbb{R}^{T \times d_x} \times \Theta \mapsto \mathbb{R}^{T \times d_y}$

Positional Encoding

Positional Encoding

Our transformer maps T inputs to T outputs

Positional Encoding

Our transformer maps T inputs to T outputs

$$f : \mathbb{R}^{T \times d_x} \times \Theta \mapsto \mathbb{R}^{T \times d_y}$$

Positional Encoding

Our transformer maps T inputs to T outputs

$$f : \mathbb{R}^{T \times d_x} \times \Theta \mapsto \mathbb{R}^{T \times d_y}$$

- T words in a sentence

Positional Encoding

Our transformer maps T inputs to T outputs

$$f : \mathbb{R}^{T \times d_x} \times \Theta \mapsto \mathbb{R}^{T \times d_y}$$

- T words in a sentence
- T pixels in an image

Positional Encoding

Our transformer maps T inputs to T outputs

$$f : \mathbb{R}^{T \times d_x} \times \Theta \mapsto \mathbb{R}^{T \times d_y}$$

- T words in a sentence
- T pixels in an image
- T amino acids in a protein

Positional Encoding

Our transformer maps T inputs to T outputs

$$f : \mathbb{R}^{T \times d_x} \times \Theta \mapsto \mathbb{R}^{T \times d_y}$$

- T words in a sentence
- T pixels in an image
- T amino acids in a protein

Question: Do we care about the order of the T inputs?

Positional Encoding

Our transformer maps T inputs to T outputs

$$f : \mathbb{R}^{T \times d_x} \times \Theta \mapsto \mathbb{R}^{T \times d_y}$$

- T words in a sentence
- T pixels in an image
- T amino acids in a protein

Question: Do we care about the order of the T inputs?

Answer: In some tasks, yes! In others, no

Positional Encoding

Question: Detecting birds in an image of T pixels, does order matter?

Positional Encoding

Question: Detecting birds in an image of T pixels, does order matter?



Positional Encoding

Question: Detecting birds in an image of T pixels, does order matter?



Answer: Yes!

Positional Encoding

Question: T robots searching for an object. Does order matter?

Positional Encoding

Question: T robots searching for an object. Does order matter?



Positional Encoding

Question: T robots searching for an object. Does order matter?



Answer: No!

Positional Encoding

Question: We translate a sentence containing T words. Does order matter?

Positional Encoding

Question: We translate a sentence containing T words. Does order matter?

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{dog} \\ \text{licks} \\ \text{the} \\ \text{owner} \end{bmatrix}$$

Positional Encoding

Question: We translate a sentence containing T words. Does order matter?

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{dog} \\ \text{licks} \\ \text{the} \\ \text{owner} \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ \textcolor{red}{x_5} \\ x_3 \\ x_4 \\ \textcolor{red}{x_2} \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{owner} \\ \text{licks} \\ \text{the} \\ \text{dog} \end{bmatrix}$$

Positional Encoding

Question: We translate a sentence containing T words. Does order matter?

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{dog} \\ \text{licks} \\ \text{the} \\ \text{owner} \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ \textcolor{red}{x_5} \\ x_3 \\ x_4 \\ \textcolor{red}{x_2} \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{owner} \\ \text{licks} \\ \text{the} \\ \text{dog} \end{bmatrix}$$

Answer: Yes!

Positional Encoding

For some tasks, we must know the order of the inputs

Positional Encoding

For some tasks, we must know the order of the inputs

Question: Does the transformer know the order of the T inputs?

Positional Encoding

For some tasks, we must know the order of the inputs

Question: Does the transformer know the order of the T inputs?

Define a **permutation matrix** $P \in \{0, 1\}^{T \times T}$ that reorders the inputs

Positional Encoding

For some tasks, we must know the order of the inputs

Question: Does the transformer know the order of the T inputs?

Define a **permutation matrix** $P \in \{0, 1\}^{T \times T}$ that reorders the inputs

Positional Encoding

Example 1:

$$P = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}; \quad a = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix};$$

Positional Encoding

Example 1:

$$P = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}; \quad a = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}; \quad Pa = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}$$

Positional Encoding

Example 1:

$$P = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}; \quad a = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}; \quad Pa = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}$$

Example 2:

$$P = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}; \quad a = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix};$$

Positional Encoding

Example 1:

$$P = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}; \quad a = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}; \quad Pa = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}$$

Example 2:

$$P = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}; \quad a = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}; \quad Pa = \begin{bmatrix} 4 \\ 3 \\ 5 \end{bmatrix}$$

Positional Encoding

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) \neq Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

Positional Encoding

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) \neq Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

f is equivariant, order **does** matter

Positional Encoding

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) \neq Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

f is equivariant, order **does** matter

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) = Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

Positional Encoding

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) \neq Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

f is equivariant, order **does** matter

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) = Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

f is equivariant, order **does not** matter

Which is a transformer?

Positional Encoding

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) \neq Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

f is equivariant, order **does** matter

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) = Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

f is equivariant, order **does not** matter

Which is a transformer?

MLP is equivariant, but what about attention?

Positional Encoding

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) \neq Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

f is equivariant, order **does** matter

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) = Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

f is equivariant, order **does not** matter

Which is a transformer?

MLP is equivariant, but what about attention?

Positional Encoding

Recall dot product self attention

Positional Encoding

Recall dot product self attention

$$Q = [q_1 \ q_2 \ \dots \ q_T] = [\theta_Q^\top x_1 \ \theta_Q^\top x_2 \ \dots \ \theta_Q^\top x_T]$$

$$K = [k_1 \ k_2 \ \dots \ k_T] = [\theta_K^\top x_1 \ \theta_K^\top x_2 \ \dots \ \theta_K^\top x_T]$$

$$V = [v_1 \ v_2 \ \dots \ v_T] = [\theta_V^\top x_1 \ \theta_V^\top x_2 \ \dots \ \theta_V^\top x_T]$$

Positional Encoding

Recall dot product self attention

$$Q = [q_1 \ q_2 \ \dots \ q_T] = [\theta_Q^\top x_1 \ \theta_Q^\top x_2 \ \dots \ \theta_Q^\top x_T]$$

$$K = [k_1 \ k_2 \ \dots \ k_T] = [\theta_K^\top x_1 \ \theta_K^\top x_2 \ \dots \ \theta_K^\top x_T]$$

$$V = [v_1 \ v_2 \ \dots \ v_T] = [\theta_V^\top x_1 \ \theta_V^\top x_2 \ \dots \ \theta_V^\top x_T]$$

$$\text{attn}\left(\begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_h}}\right)V$$

Positional Encoding

Permuting the inputs reorders Q, K, V

Positional Encoding

Permuting the inputs reorders Q, K, V

$$PQ = [q_T \ q_1 \ \dots \ q_2] = [\theta_Q^\top x_T \ \theta_Q^\top x_1 \ \dots \ \theta_Q^\top x_2]$$

$$PK = [k_T \ k_1 \ \dots \ k_2] = [\theta_K^\top x_T \ \theta_K^\top x_1 \ \dots \ \theta_K^\top x_2]$$

$$PV = [v_T \ v_1 \ \dots \ v_2] = [\theta_V^\top x_T \ \theta_V^\top x_1 \ \dots \ \theta_V^\top x_2]$$

Positional Encoding

Permuting the inputs reorders Q, K, V

$$PQ = [q_T \ q_1 \ \dots \ q_2] = [\theta_Q^\top x_T \ \theta_Q^\top x_1 \ \dots \ \theta_Q^\top x_2]$$

$$PK = [k_T \ k_1 \ \dots \ k_2] = [\theta_K^\top x_T \ \theta_K^\top x_1 \ \dots \ \theta_K^\top x_2]$$

$$PV = [v_T \ v_1 \ \dots \ v_2] = [\theta_V^\top x_T \ \theta_V^\top x_1 \ \dots \ \theta_V^\top x_2]$$

$$\text{attn}\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{(PQ)(PK)^\top}{\sqrt{d_h}}\right)(PV)$$

Positional Encoding

Permuting the inputs reorders Q, K, V

$$PQ = [q_T \ q_1 \ \dots \ q_2] = [\theta_Q^\top x_T \ \theta_Q^\top x_1 \ \dots \ \theta_Q^\top x_2]$$

$$PK = [k_T \ k_1 \ \dots \ k_2] = [\theta_K^\top x_T \ \theta_K^\top x_1 \ \dots \ \theta_K^\top x_2]$$

$$PV = [v_T \ v_1 \ \dots \ v_2] = [\theta_V^\top x_T \ \theta_V^\top x_1 \ \dots \ \theta_V^\top x_2]$$

$$\text{attn}\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{(PQ)(PK)^\top}{\sqrt{d_h}}\right)(PV)$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{(\mathbf{PQ})(\mathbf{PK})^\top}{\sqrt{d_h}}\right)(\mathbf{PV})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{(\mathbf{PQ})(\mathbf{PK})^\top}{\sqrt{d_h}}\right)(\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}(\mathbf{PK})^\top}{\sqrt{d_h}}\right)(\mathbf{PV})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{(\mathbf{PQ})(\mathbf{PK})^\top}{\sqrt{d_h}}\right)(\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}(\mathbf{PK})^\top}{\sqrt{d_h}}\right)(\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}(\mathbf{K}^\top \mathbf{P}^\top)}{\sqrt{d_h}}\right)(\mathbf{PV})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\frac{(\mathbf{PQ})(\mathbf{PK})^\top}{\sqrt{d_h}}\right)(\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}(\mathbf{PK})^\top}{\sqrt{d_h}}\right)(\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}(\mathbf{K}^\top \mathbf{P}^\top)}{\sqrt{d_h}}\right)(\mathbf{PV})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}(\mathbf{K}^\top \mathbf{P}^\top)}{\sqrt{d_h}}\right)(\mathbf{P}\mathbf{V})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}(\mathbf{K}^\top \mathbf{P}^\top)}{\sqrt{d_h}}\right)(\mathbf{P}\mathbf{V})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right)(\mathbf{P}\mathbf{V})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}(\mathbf{K}^\top \mathbf{P}^\top)}{\sqrt{d_h}}\right)(\mathbf{P}\mathbf{V})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right)(\mathbf{P}\mathbf{V})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right)(\mathbf{P}\mathbf{V})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}(\mathbf{K}^\top \mathbf{P}^\top)}{\sqrt{d_h}}\right)(\mathbf{P}\mathbf{V})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \text{softmax}\left(\mathbf{P} \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right)(\mathbf{P}\mathbf{V})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right)(\mathbf{P}\mathbf{V})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right) (\mathbf{PV})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right) (\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}}\right) \mathbf{P}^\top (\mathbf{PV})$$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right) (\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}}\right) \mathbf{P}^\top (\mathbf{PV})$$

\mathbf{P} swaps row i and j

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right) (\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}}\right) \mathbf{P}^\top (\mathbf{PV})$$

\mathbf{P} swaps row i and j \mathbf{P}^T swaps row j and i

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right) (\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}}\right) \mathbf{P}^\top (\mathbf{PV})$$

\mathbf{P} swaps row i and j \mathbf{P}^T swaps row j and i $\mathbf{P}^\top \mathbf{P} = \mathbf{I}$

Positional Encoding

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}} \mathbf{P}^\top\right) (\mathbf{PV})$$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}}\right) \mathbf{P}^\top (\mathbf{PV})$$

\mathbf{P} swaps row i and j \mathbf{P}^T swaps row j and i $\mathbf{P}^\top \mathbf{P} = \mathbf{I}$

$$\text{attn}\left(\mathbf{P} \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = \mathbf{P} \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_h}}\right) \mathbf{V}$$

Positional Encoding

$$\text{attn}\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = P \text{softmax}\left(\frac{QK^\top}{\sqrt{d_h}}\right)V$$

Positional Encoding

$$\text{attn}\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = P \text{softmax}\left(\frac{QK^\top}{\sqrt{d_h}}\right)V$$

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) = Pf\begin{pmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \end{pmatrix}$$

Positional Encoding

$$\text{attn}\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = P \text{softmax}\left(\frac{QK^\top}{\sqrt{d_h}}\right)V$$

$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) = Pf\begin{pmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \end{pmatrix}$$

Question: What does this mean?

Positional Encoding

$$\text{attn}\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix}, \theta\right) = P \text{softmax}\left(\frac{QK^\top}{\sqrt{d_h}}\right)V$$
$$f\left(P \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) = Pf\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right)$$

Question: What does this mean?

Answer: Attention is equivariant, order **does not** matter

Positional Encoding

This makes sense, in our party example we never consider the order

Positional Encoding

This makes sense, in our party example we never consider the order



Positional Encoding

This makes sense, in our party example we never consider the order



Transformer cannot determine order of inputs! **Equivariant** to ordering

Positional Encoding

The following sentences are the same to a transformer

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \mathbf{x}_5 \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{dog} \\ \text{licks} \\ \text{the} \\ \text{owner} \end{bmatrix}$$

Positional Encoding

The following sentences are the same to a transformer

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \mathbf{x}_5 \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{dog} \\ \text{licks} \\ \text{the} \\ \text{owner} \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_5 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{owner} \\ \text{licks} \\ \text{the} \\ \text{dog} \end{bmatrix}$$

Positional Encoding

The following sentences are the same to a transformer

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \mathbf{x}_5 \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{dog} \\ \text{licks} \\ \text{the} \\ \text{owner} \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_5 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \text{The} \\ \text{owner} \\ \text{licks} \\ \text{the} \\ \text{dog} \end{bmatrix}$$

This is a problem! For some tasks, we care about input order

Positional Encoding

Question: What are some ways we can introduce ordering?

Positional Encoding

Question: What are some ways we can introduce ordering?

Answer 1: We can introduce forgetting

Positional Encoding

Question: What are some ways we can introduce ordering?

Answer 1: We can introduce forgetting

ALiBi: Press, Ofir, Noah Smith, and Mike Lewis. “Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation.” *International Conference on Learning Representations*.

Positional Encoding

Question: What are some ways we can introduce ordering?

Answer 1: We can introduce forgetting

ALiBi: Press, Ofir, Noah Smith, and Mike Lewis. “*Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation.*” *International Conference on Learning Representations*.

RoPE: Su, Jianlin, et al. “*Roformer: Enhanced transformer with rotary position embedding.*” *Neurocomputing*.

Positional Encoding

Question: What are some ways we can introduce ordering?

Answer 1: We can introduce forgetting

ALiBi: Press, Ofir, Noah Smith, and Mike Lewis. “*Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation.*” *International Conference on Learning Representations*.

RoPE: Su, Jianlin, et al. “*Roformer: Enhanced transformer with rotary position embedding.*” *Neurocomputing*.

Answer 2: We can modify the inputs based on their ordering

Positional Encoding

Question: What are some ways we can introduce ordering?

Answer 1: We can introduce forgetting

ALiBi: Press, Ofir, Noah Smith, and Mike Lewis. “*Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation.*” *International Conference on Learning Representations*.

RoPE: Su, Jianlin, et al. “*Roformer: Enhanced transformer with rotary position embedding.*” *Neurocomputing*.

Answer 2: We can modify the inputs based on their ordering

We will focus on answer 2 because it is more common

Positional Encoding

$$\text{attn} \left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_T \end{bmatrix}, \theta \right)$$

Positional Encoding

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta \right)$$

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(1) \\ f_{\text{pos}}(2) \\ \vdots \\ f_{\text{pos}}(3) \end{bmatrix}, \theta \right)$$

Positional Encoding

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta \right)$$

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(1) \\ f_{\text{pos}}(2) \\ \vdots \\ f_{\text{pos}}(3) \end{bmatrix}, \theta \right)$$

Even if we permute the inputs, we still know the order!

Positional Encoding

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta \right)$$

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(1) \\ f_{\text{pos}}(2) \\ \vdots \\ f_{\text{pos}}(3) \end{bmatrix}, \theta \right)$$

Even if we permute the inputs, we still know the order!

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_T \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_2 \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(T) \\ f_{\text{pos}}(1) \\ \vdots \\ f_{\text{pos}}(2) \end{bmatrix}, \theta \right)$$

Positional Encoding

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \theta \right)$$

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_T \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(1) \\ f_{\text{pos}}(2) \\ \vdots \\ f_{\text{pos}}(3) \end{bmatrix}, \theta \right)$$

Even if we permute the inputs, we still know the order!

$$\text{attn} \left(\begin{bmatrix} \mathbf{x}_T \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_2 \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(T) \\ f_{\text{pos}}(1) \\ \vdots \\ f_{\text{pos}}(2) \end{bmatrix}, \theta \right)$$

Positional Encoding

What is f_{pos} ?

Positional Encoding

What is f_{pos} ?

Easiest solution is just some random parameters

Positional Encoding

What is f_{pos} ?

Easiest solution is just some random parameters

$$f(i) = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_T \end{bmatrix}_i$$

Positional Encoding

What is f_{pos} ?

Easiest solution is just some random parameters

$$f(i) = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_T \end{bmatrix}_i$$

If ordering matters, learn θ to be different

Positional Encoding

What is f_{pos} ?

Easiest solution is just some random parameters

$$f(i) = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_T \end{bmatrix}_i$$

If ordering matters, learn θ to be different

If order does not matter, learn θ to be the same/zero

Positional Encoding

What is f_{pos} ?

Easiest solution is just some random parameters

$$f(i) = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_T \end{bmatrix}_i$$

If ordering matters, learn θ to be different

If order does not matter, learn θ to be the same/zero

In torch, this is called `nn.Embedding`

Positional Encoding

So, our final transformer is

```
class Transformer(nn.Module):
    def __init__(self):
        self.position_embedding = nn.Embedding(d_x)
        self.block1 = TransformerLayer()
        self.block2 = TransformerLayer()

    def forward(self, x):
        x = x + embedding(torch.arange(x.shape[0]))
        x = self.block1(x)
        x = self.block2(x)
        return x
```

Positional Encoding

Let us code up the transformer in colab

<https://colab.research.google.com/drive/1qVlIbQKpTuBYIa7FvC4IH-kJq-E0jmc0d#scrollTo=iQtXjGYiz5CD>

Applications

Applications

Now that we understand the transformer, how do we use it?

Applications

Now that we understand the transformer, how do we use it?

We can apply it to many domains, but today we will focus on text and images

Applications

Now that we understand the transformer, how do we use it?

We can apply it to many domains, but today we will focus on text and images

Text Transformers

Text Transformers

In text transformers, we create a parameter for each word

Text Transformers

In text transformers, we create a parameter for each word

We call these parameters **tokens**

Text Transformers

In text transformers, we create a parameter for each word

We call these parameters **tokens**

The simple way is to create one parameter for each unique word in the dataset

Text Transformers

In text transformers, we create a parameter for each word

We call these parameters **tokens**

The simple way is to create one parameter for each unique word in the dataset

Find unique words, and assign each one a parameter

Text Transformers

In text transformers, we create a parameter for each word

We call these parameters **tokens**

The simple way is to create one parameter for each unique word in the dataset

Find unique words, and assign each one a parameter

$$\text{unique} \left(\begin{bmatrix} \text{John likes movies} \\ \text{Mary likes movies} \\ \text{I like dogs} \end{bmatrix} \right) = [\text{John likes movies Mary I dogs}]$$

Text Transformers

Then assign a parameter to each word

$$\begin{bmatrix} \text{John} \\ \text{likes} \\ \text{movies} \\ \text{Mary} \\ I \\ \text{dogs} \end{bmatrix} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \theta_5 \\ \theta_6 \end{bmatrix}$$

$$[\theta_1 \ \theta_2 \ \theta_3] =$$

Text Transformers

Then assign a parameter to each word

$$\begin{bmatrix} \text{John} \\ \text{likes} \\ \text{movies} \\ \text{Mary} \\ I \\ \text{dogs} \end{bmatrix} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \theta_5 \\ \theta_6 \end{bmatrix}$$

$$[\theta_1 \ \theta_2 \ \theta_3] = \text{John likes movies}$$

$$[\theta_4 \ \theta_2 \ \theta_1] =$$

Text Transformers

Then assign a parameter to each word

$$\begin{bmatrix} \text{John} \\ \text{likes} \\ \text{movies} \\ \text{Mary} \\ I \\ \text{dogs} \end{bmatrix} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \theta_5 \\ \theta_6 \end{bmatrix}$$

$$[\theta_1 \ \theta_2 \ \theta_3] = \text{John likes movies}$$

$$[\theta_4 \ \theta_2 \ \theta_1] = \text{Mary likes John}$$

Text Transformers

```
unique_words = set(sentence.split(" ")) for sentence in xs)
tokens = {word: i for i, word in enumerate(unique_words)}
embeddings = nn.Embedding(len(tokens), d_x)
# Convert from words to parameters
xs = []
for sentence in sentences:
    xs.append([])
    for word in sentence:
        index = embeddings[word]
        token = tokens[index]
        xs.append(token)

print(xs)
>>> [[Tensor(...), Tensor(...), ...]]
```

Text Transformers

```
model = Transformer()
for tokenized_sentence in xs:
    # Convert list to tensor
    x = torch.stack(tokenized_sentence)
    y = model(x)
```

Image Transformers

Image Transformers

In image transformers, we treat a **patch** of pixels as an x

Image Transformers

In image transformers, we treat a **patch** of pixels as an x

$$X \in [0, 1]^{16 \times 16}$$

Image Transformers

In image transformers, we treat a **patch** of pixels as an x

$$X \in [0, 1]^{16 \times 16}$$

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
x_9	...						

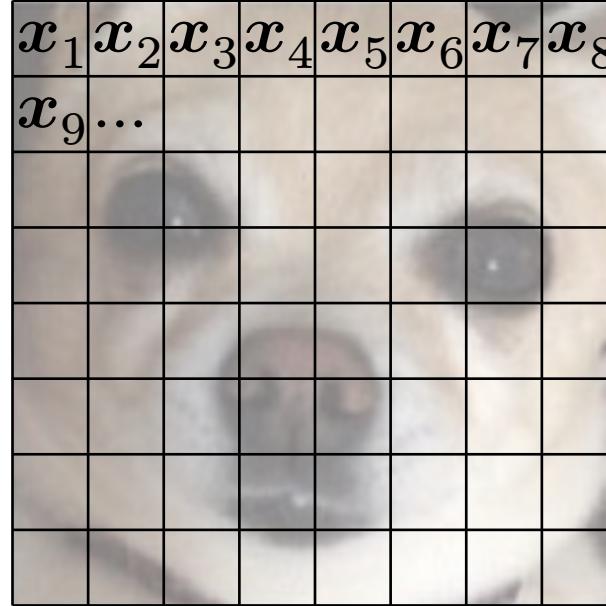


Image Transformers

```
# Convert image into patches
patches = []
for i in range(0, x.shape[0] - k + 1, k):
    for j in range(0, x.shape[1] - k + 1, k):
        patches.append(
            x[i: i + k , j: j + k]
)
patches = stack(patches, axis=0)
print(patches.shape)
>>> (T, k, k)

model = Transformer()
y = model(patches)
```

Unsupervised Training

Unsupervised Training

Question: How do we train transformers?

Unsupervised Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

Unsupervised Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

- Classification loss

Unsupervised Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

- Classification loss
- Regression loss

Unsupervised Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

- Classification loss
- Regression loss

In practice, transformers require lots of training data

Unsupervised Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

- Classification loss
- Regression loss

In practice, transformers require lots of training data

There is almost infinite empirical scaling – add more data, model becomes stronger

Unsupervised Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

- Classification loss
- Regression loss

In practice, transformers require lots of training data

There is almost infinite empirical scaling – add more data, model becomes stronger

Transformers are limited by the size of datasets today

Unsupervised Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

- Classification loss
- Regression loss

In practice, transformers require lots of training data

There is almost infinite empirical scaling – add more data, model becomes stronger

Transformers are limited by the size of datasets today

There are not enough students to label training data!

Unsupervised Training

Question: How can we train transformers if we cannot create big enough datasets?

Unsupervised Training

Question: How can we train transformers if we cannot create big enough datasets?

Answer: We can use **unsupervised learning**

Unsupervised Training

Question: How can we train transformers if we cannot create big enough datasets?

Answer: We can use **unsupervised learning**

The internet contains billions of unlabeled sentences and images

Unsupervised Training

Question: How can we train transformers if we cannot create big enough datasets?

Answer: We can use **unsupervised learning**

The internet contains billions of unlabeled sentences and images

We can mask/hide part of the input, and train the model to predict the missing parts

Unsupervised Training

Question: How can we train transformers if we cannot create big enough datasets?

Answer: We can use **unsupervised learning**

The internet contains billions of unlabeled sentences and images

We can mask/hide part of the input, and train the model to predict the missing parts

Another name for this is **generative pre-training** (GPT)

Unsupervised Training

Question: How can we train transformers if we cannot create big enough datasets?

Answer: We can use **unsupervised learning**

The internet contains billions of unlabeled sentences and images

We can mask/hide part of the input, and train the model to predict the missing parts

Another name for this is **generative pre-training** (GPT)

This method is **extremely** powerful

Unsupervised Training

We pose the following objective

Unsupervised Training

We pose the following objective

$$\arg \min_{\theta} \mathcal{L} \left(\begin{bmatrix} \boldsymbol{x}_1 \\ \vdots \\ \boldsymbol{x}_{T-1} \end{bmatrix}, \theta \right) = \arg \min_{\theta} [-\log P(\boldsymbol{x}_T \mid \boldsymbol{x}_1, \dots, \boldsymbol{x}_{T-1}; \theta)]$$

Unsupervised Training

We pose the following objective

$$\arg \min_{\theta} \mathcal{L} \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{T-1} \end{bmatrix}, \theta \right) = \arg \min_{\theta} [-\log P(\mathbf{x}_T \mid \mathbf{x}_1, \dots, \mathbf{x}_{T-1}; \theta)]$$

$$P \left(\begin{bmatrix} \text{movies} \\ \text{mary} \\ \text{dogs} \end{bmatrix} \mid \text{John, likes; } \theta \right) = \begin{bmatrix} 0.5 \\ 0.3 \\ 0.2 \end{bmatrix}$$

Unsupervised Training

We pose the following objective

$$\arg \min_{\theta} \mathcal{L}\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{T-1} \end{bmatrix}, \theta\right) = \arg \min_{\theta} [-\log P(\mathbf{x}_T \mid \mathbf{x}_1, \dots, \mathbf{x}_{T-1}; \theta)]$$

$$P\left(\begin{bmatrix} \text{movies} \\ \text{mary} \\ \text{dogs} \end{bmatrix} \mid \text{John, likes; } \theta\right) = \begin{bmatrix} 0.5 \\ 0.3 \\ 0.2 \end{bmatrix}$$

$$P(\text{movies} \mid \text{John, likes; } \theta) = 0.5$$

Unsupervised Training

We pose the following objective

$$\arg \min_{\theta} \mathcal{L}\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{T-1} \end{bmatrix}, \theta\right) = \arg \min_{\theta} [-\log P(\mathbf{x}_T \mid \mathbf{x}_1, \dots, \mathbf{x}_{T-1}; \theta)]$$

$$P\left(\begin{bmatrix} \text{movies} \\ \text{mary} \\ \text{dogs} \end{bmatrix} \mid \text{John, likes; } \theta\right) = \begin{bmatrix} 0.5 \\ 0.3 \\ 0.2 \end{bmatrix}$$

$$P(\text{movies} \mid \text{John, likes}; \theta) = 0.5$$

Update θ so that $P(\text{movies} \mid \text{John, likes}; \theta) = 1.0$

Unsupervised Training

We pose the following loss function

Unsupervised Training

We pose the following loss function

$$\arg \min_{\theta} \mathcal{L} \left(\begin{bmatrix} \boldsymbol{x}_1 \\ \vdots \\ \boldsymbol{x}_{T-1} \end{bmatrix}, \theta \right) = \arg \min_{\theta} [-\log P(\boldsymbol{x}_T \mid \boldsymbol{x}_1, \dots, \boldsymbol{x}_{T-1}; \theta)]$$

Unsupervised Training

We pose the following loss function

$$\arg \min_{\theta} \mathcal{L} \left(\begin{bmatrix} \boldsymbol{x}_1 \\ \vdots \\ \boldsymbol{x}_{T-1} \end{bmatrix}, \theta \right) = \arg \min_{\theta} [-\log P(\boldsymbol{x}_T \mid \boldsymbol{x}_1, \dots, \boldsymbol{x}_{T-1}; \theta)]$$

For pixels, the square error represents a normal distribution over pixel values

Unsupervised Training

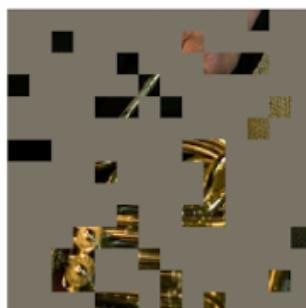
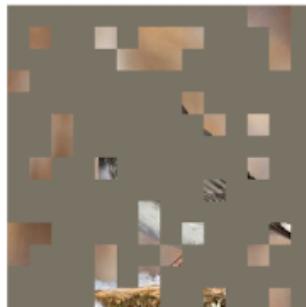
We pose the following loss function

$$\arg \min_{\theta} \mathcal{L}\left(\begin{bmatrix} x_1 \\ \vdots \\ x_{T-1} \end{bmatrix}, \theta\right) = \arg \min_{\theta} [-\log P(x_T \mid x_1, \dots, x_{T-1}; \theta)]$$

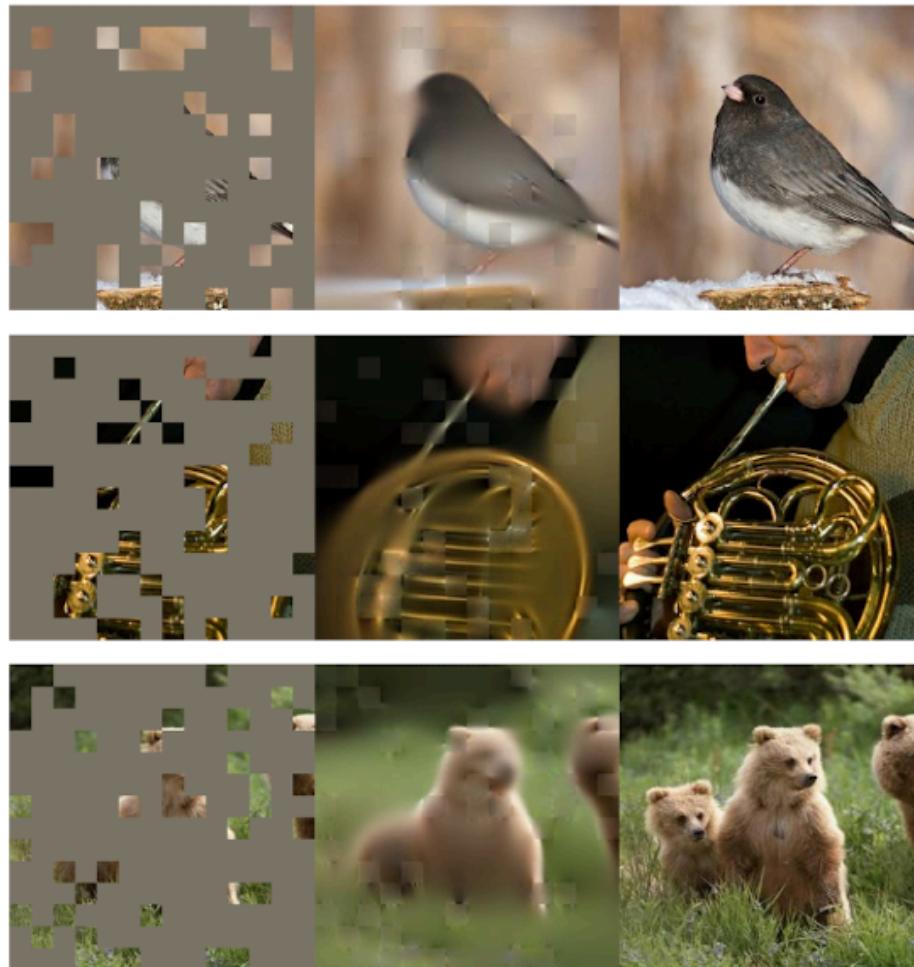
For pixels, the square error represents a normal distribution over pixel values

$$\begin{aligned} & \arg \min_{\theta} [-\log P(\text{pixel_3} \mid \text{pixel_1}, \text{pixel_2}; \theta)] \\ &= \arg \min_{\theta} \left(f\left(\begin{bmatrix} \text{pixel_1} \\ \text{pixel_2} \end{bmatrix}, \theta\right) - \text{pixel_3} \right)^2 \end{aligned}$$

Unsupervised Training

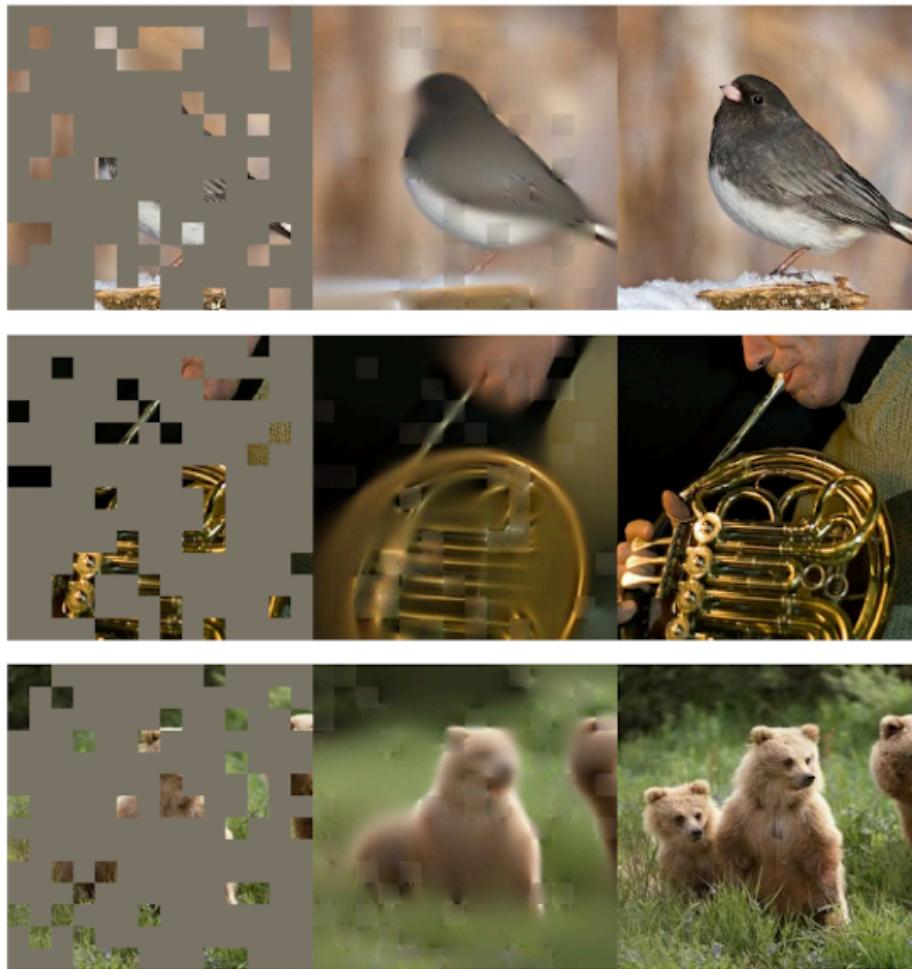


Unsupervised Training



He, Kaiming, et al. “Masked autoencoders are scalable vision learners.” Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2022.

Unsupervised Training



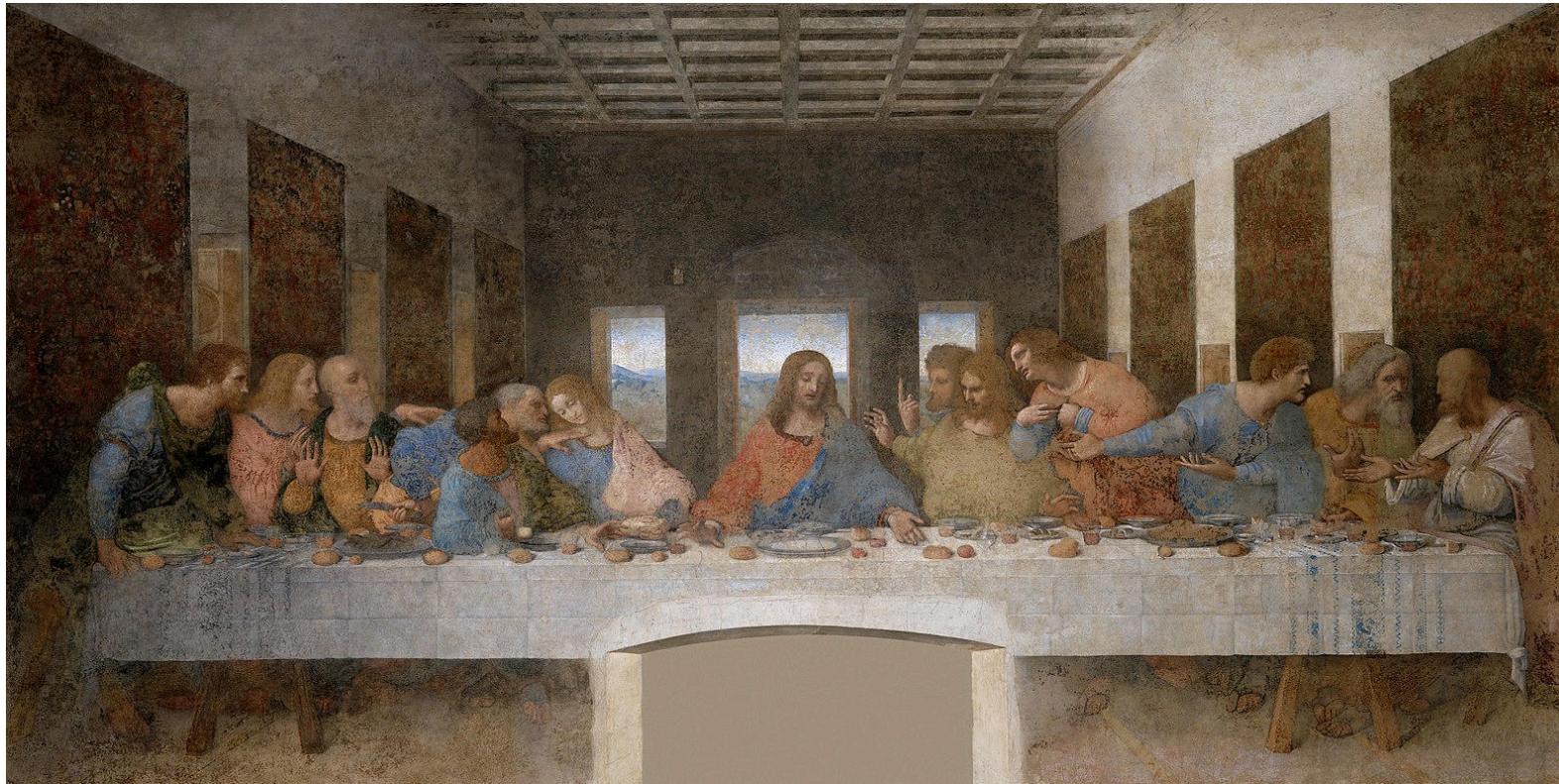
He, Kaiming, et al. “Masked autoencoders are scalable vision learners.” Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2022.

Unsupervised Training

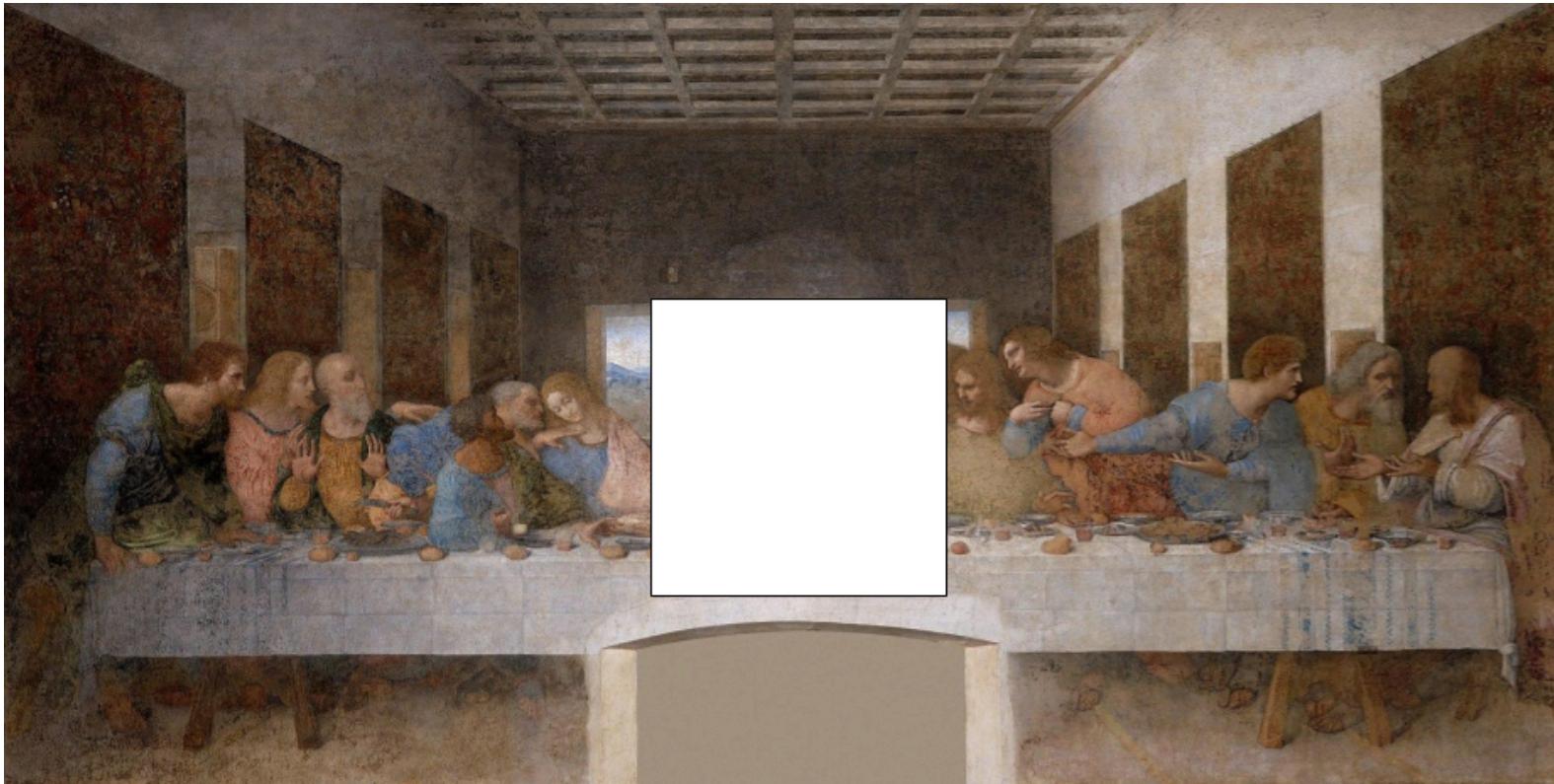
Anyone familiar with Da Vinci's painting *The Last Supper*?

Unsupervised Training

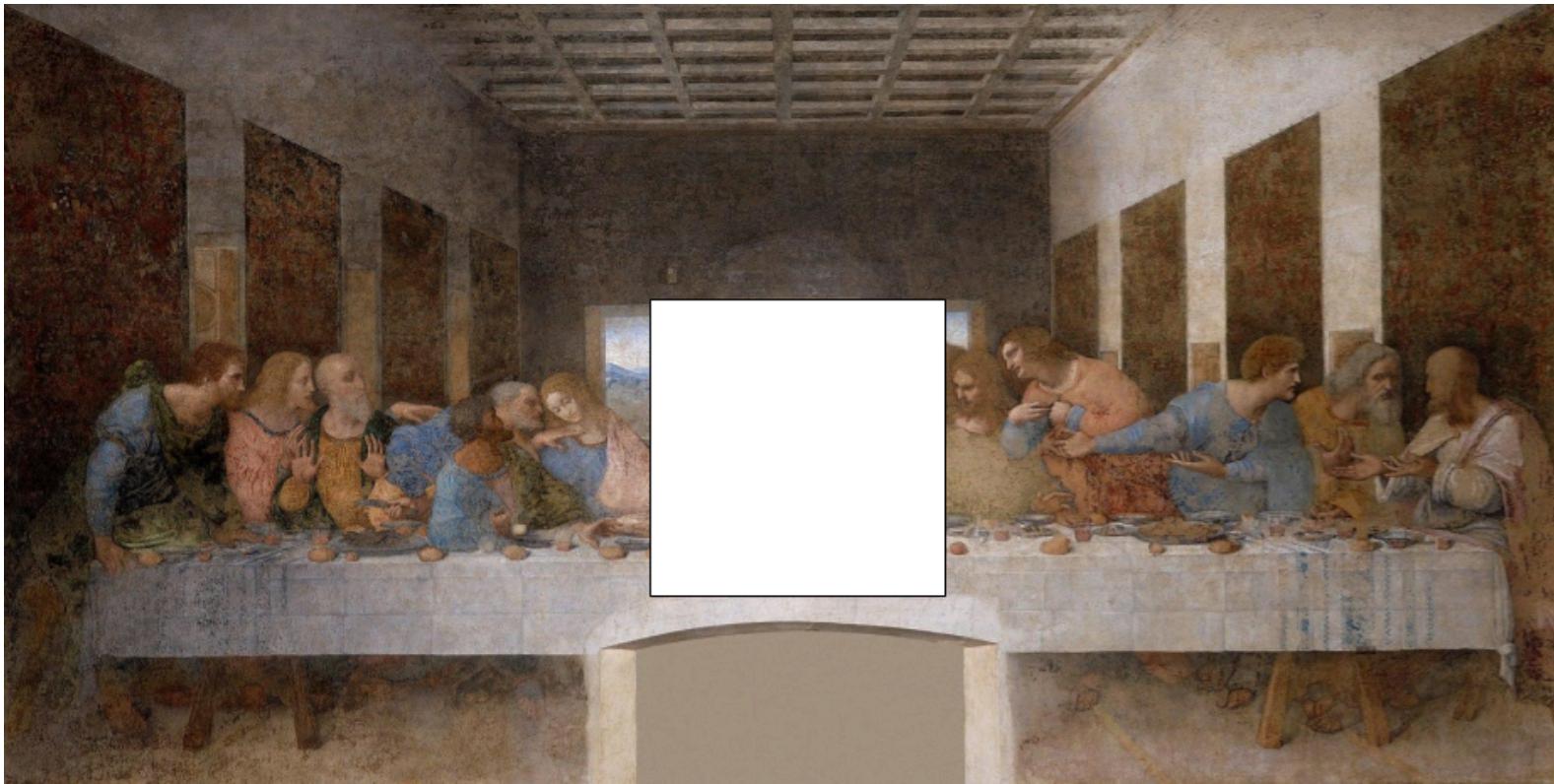
Anyone familiar with Da Vinci's painting *The Last Supper*?



Unsupervised Training



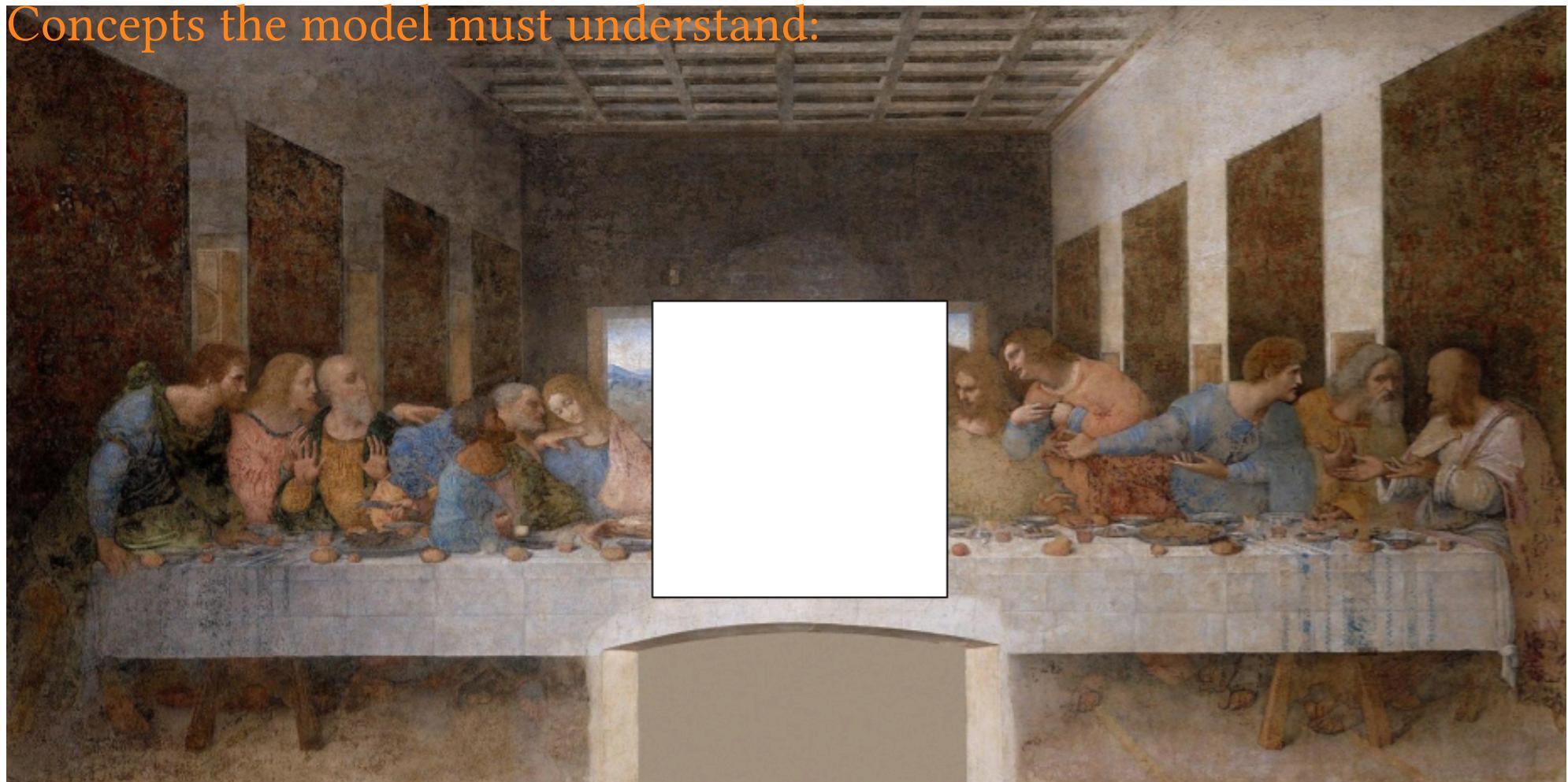
Unsupervised Training



Question: What concepts does the vision transformer need to understand to predict the missing pixels?

Unsupervised Training

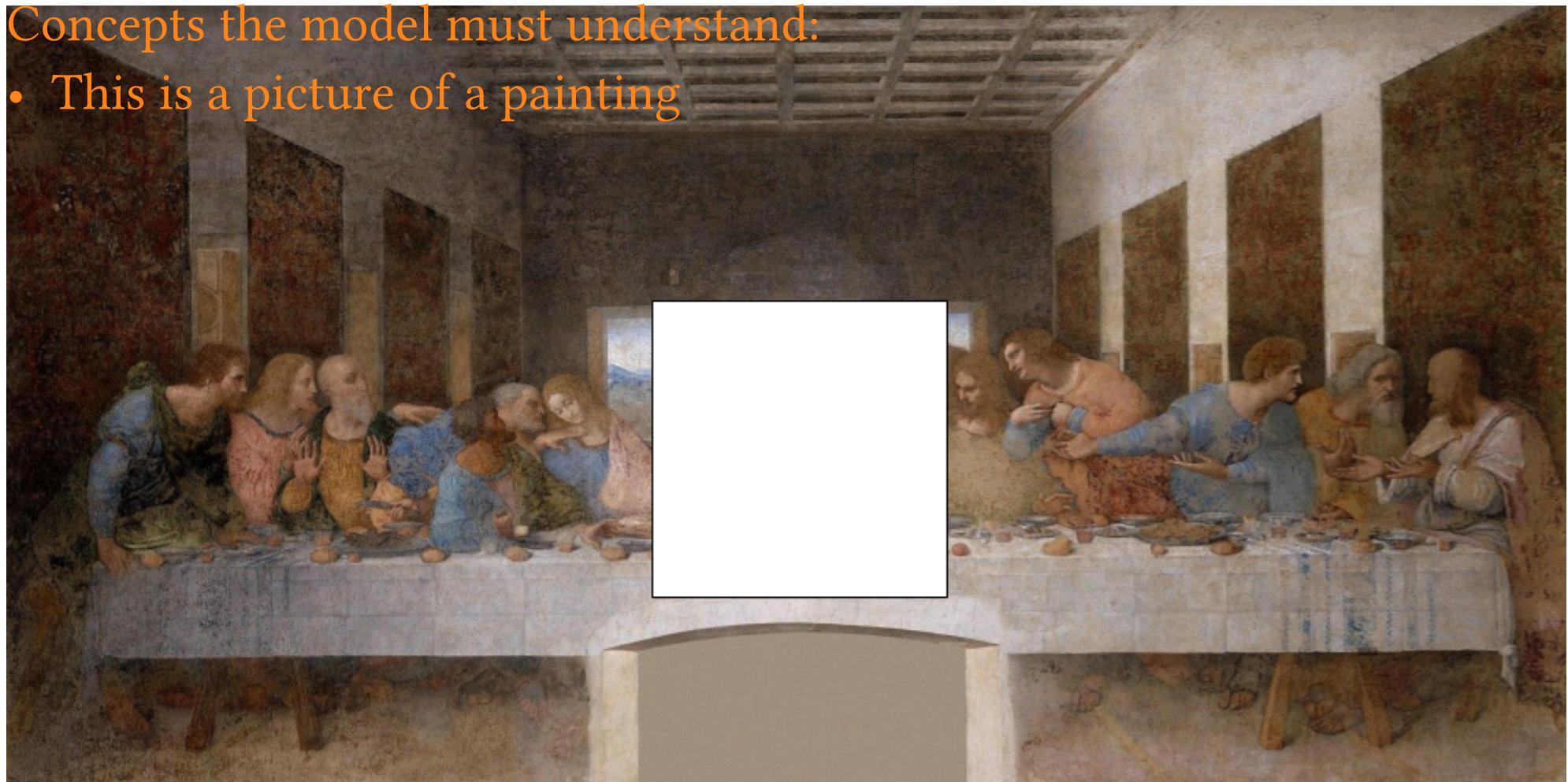
Concepts the model must understand:



Unsupervised Training

Concepts the model must understand:

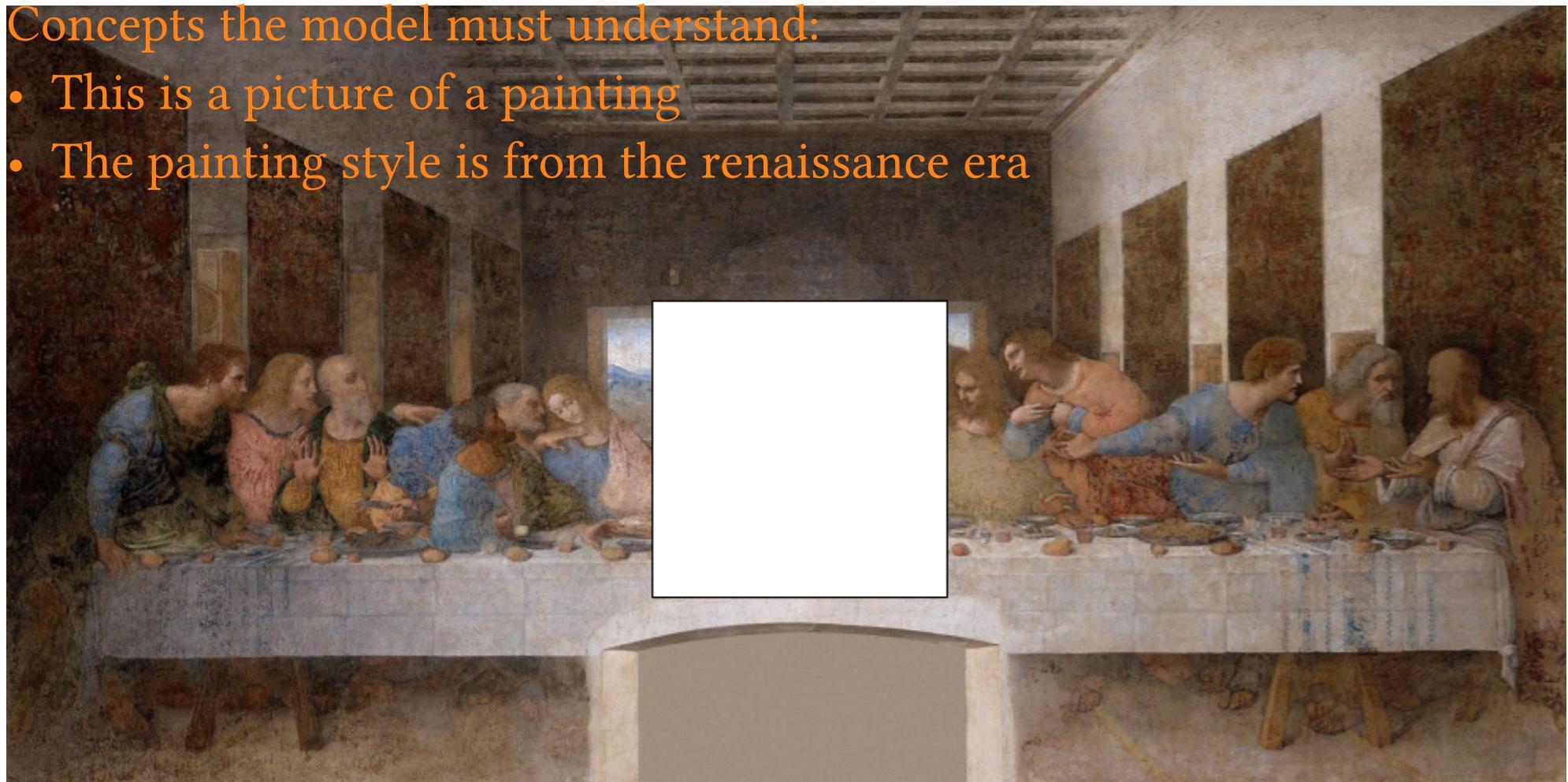
- This is a picture of a painting



Unsupervised Training

Concepts the model must understand:

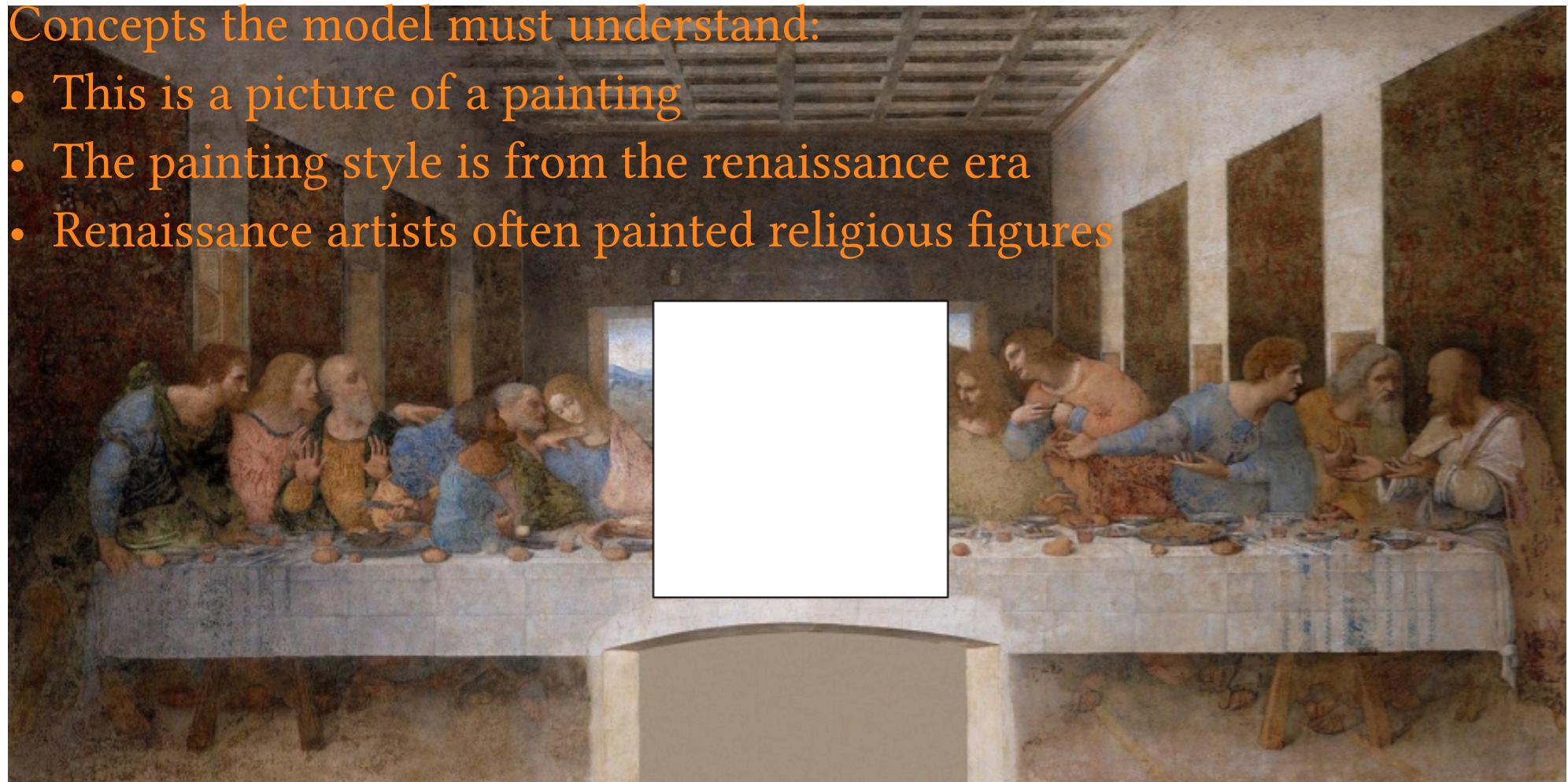
- This is a picture of a painting
- The painting style is from the renaissance era



Unsupervised Training

Concepts the model must understand:

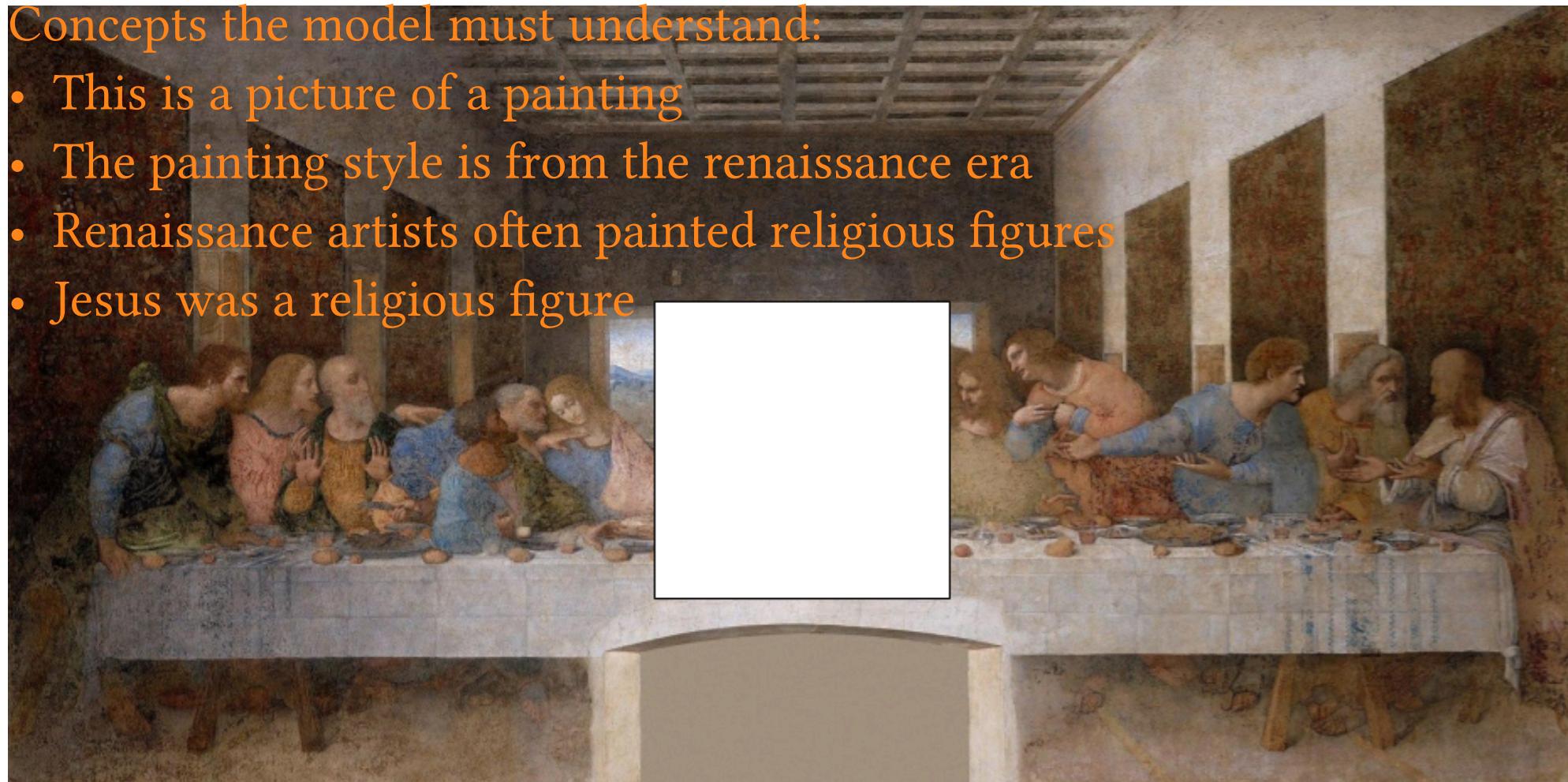
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures



Unsupervised Training

Concepts the model must understand:

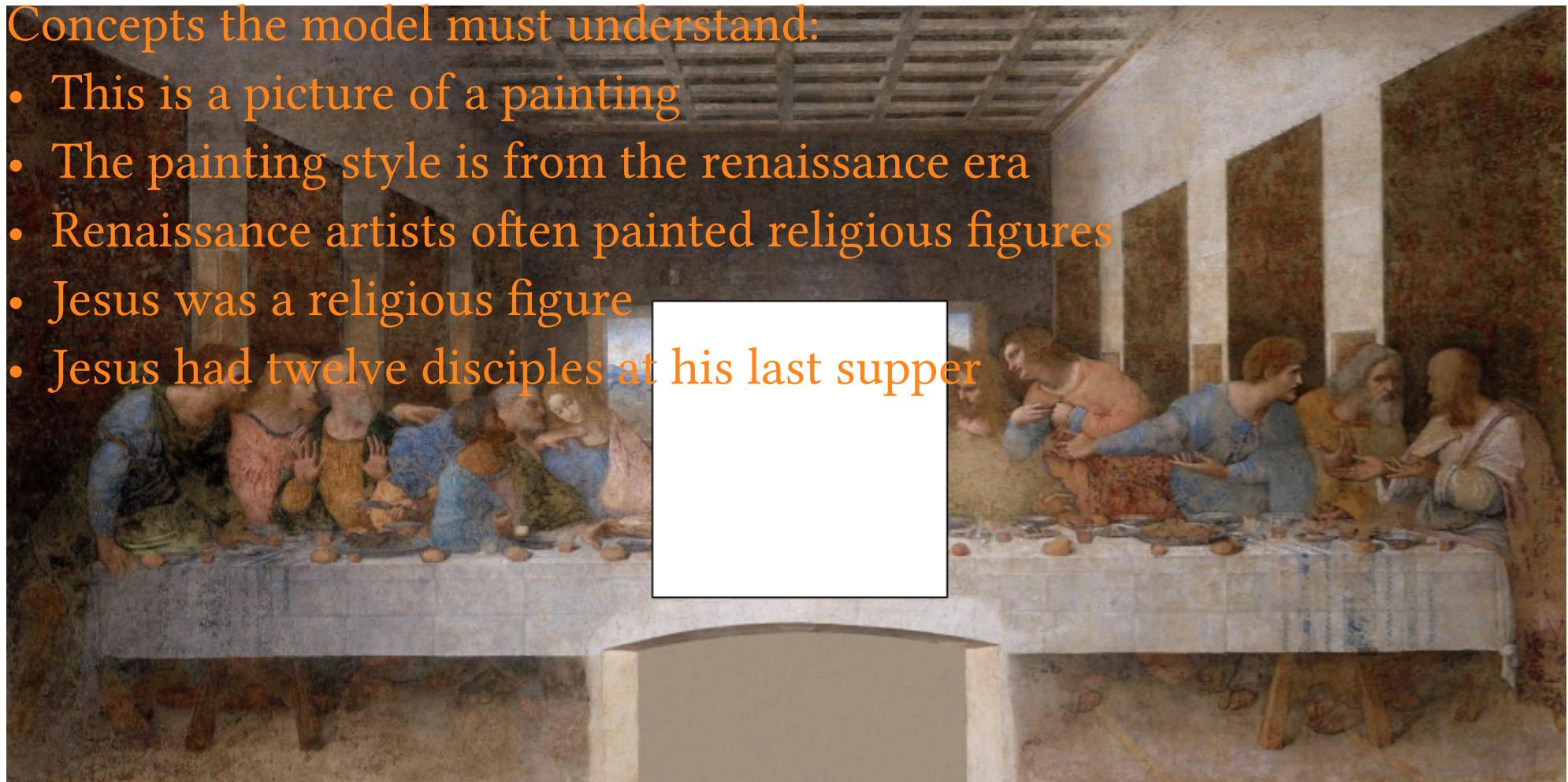
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures
- Jesus was a religious figure



Unsupervised Training

Concepts the model must understand:

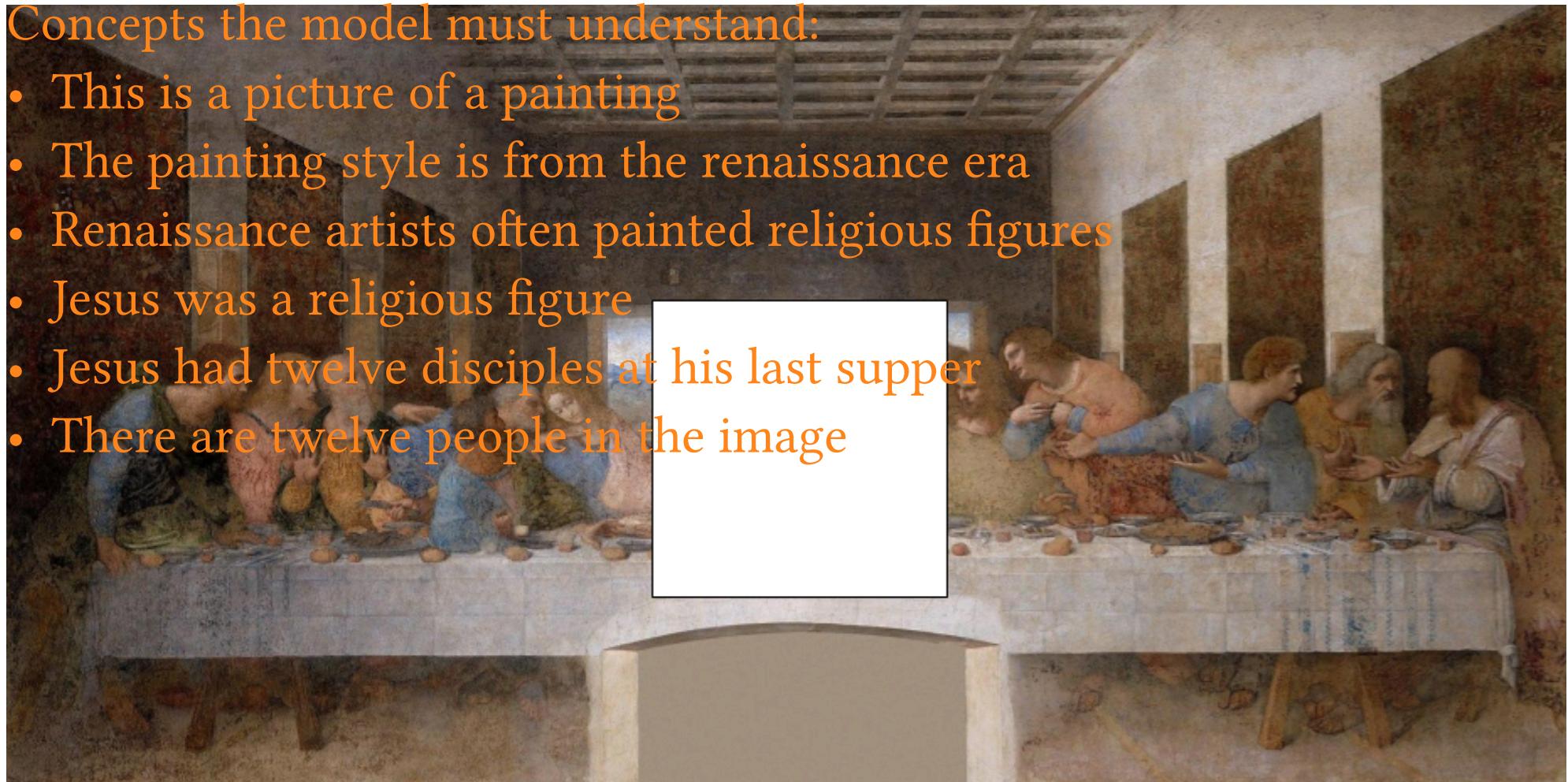
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures
- Jesus was a religious figure
- Jesus had twelve disciples at his last supper



Unsupervised Training

Concepts the model must understand:

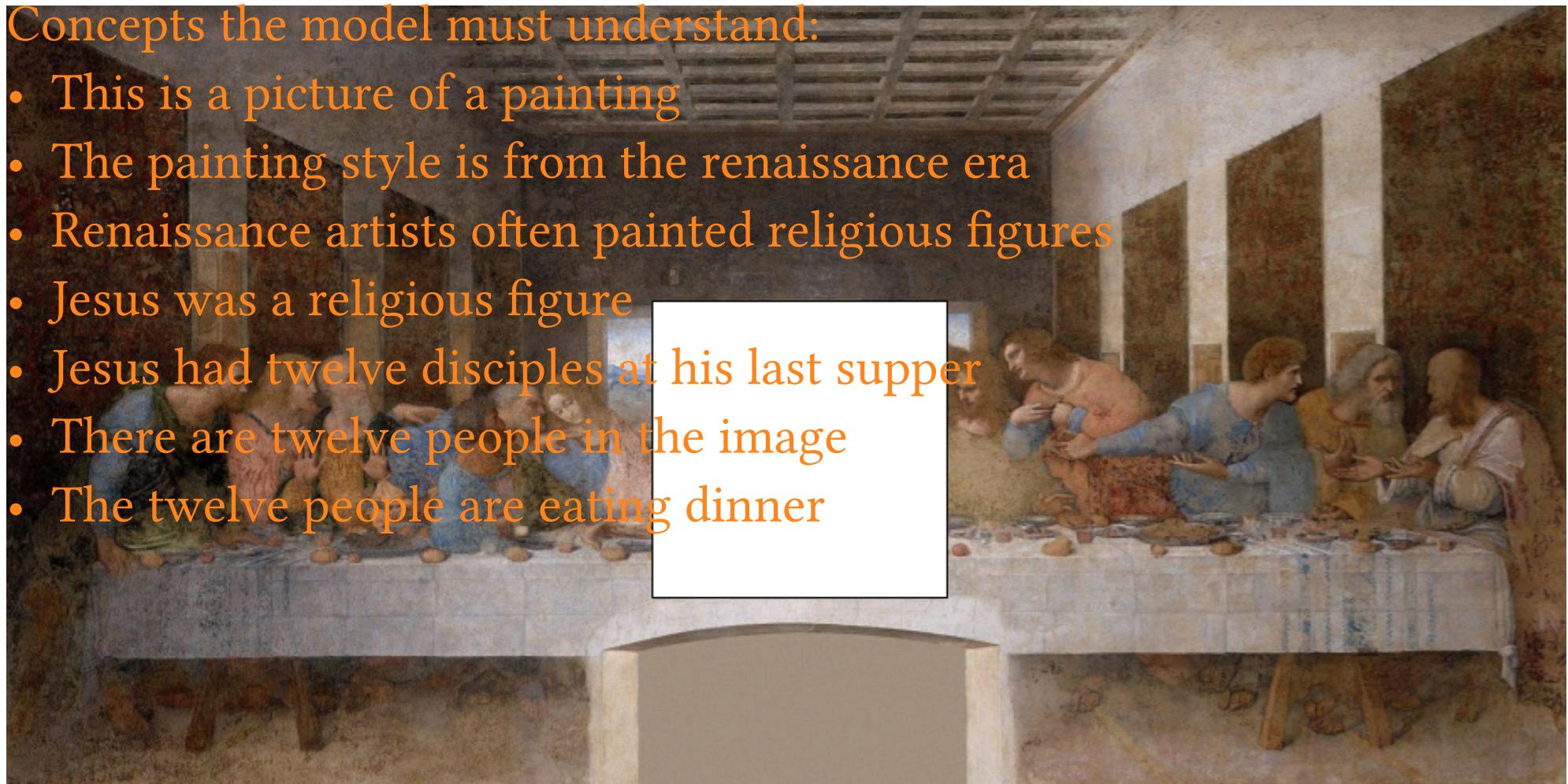
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures
- Jesus was a religious figure
- Jesus had twelve disciples at his last supper
- There are twelve people in the image



Unsupervised Training

Concepts the model must understand:

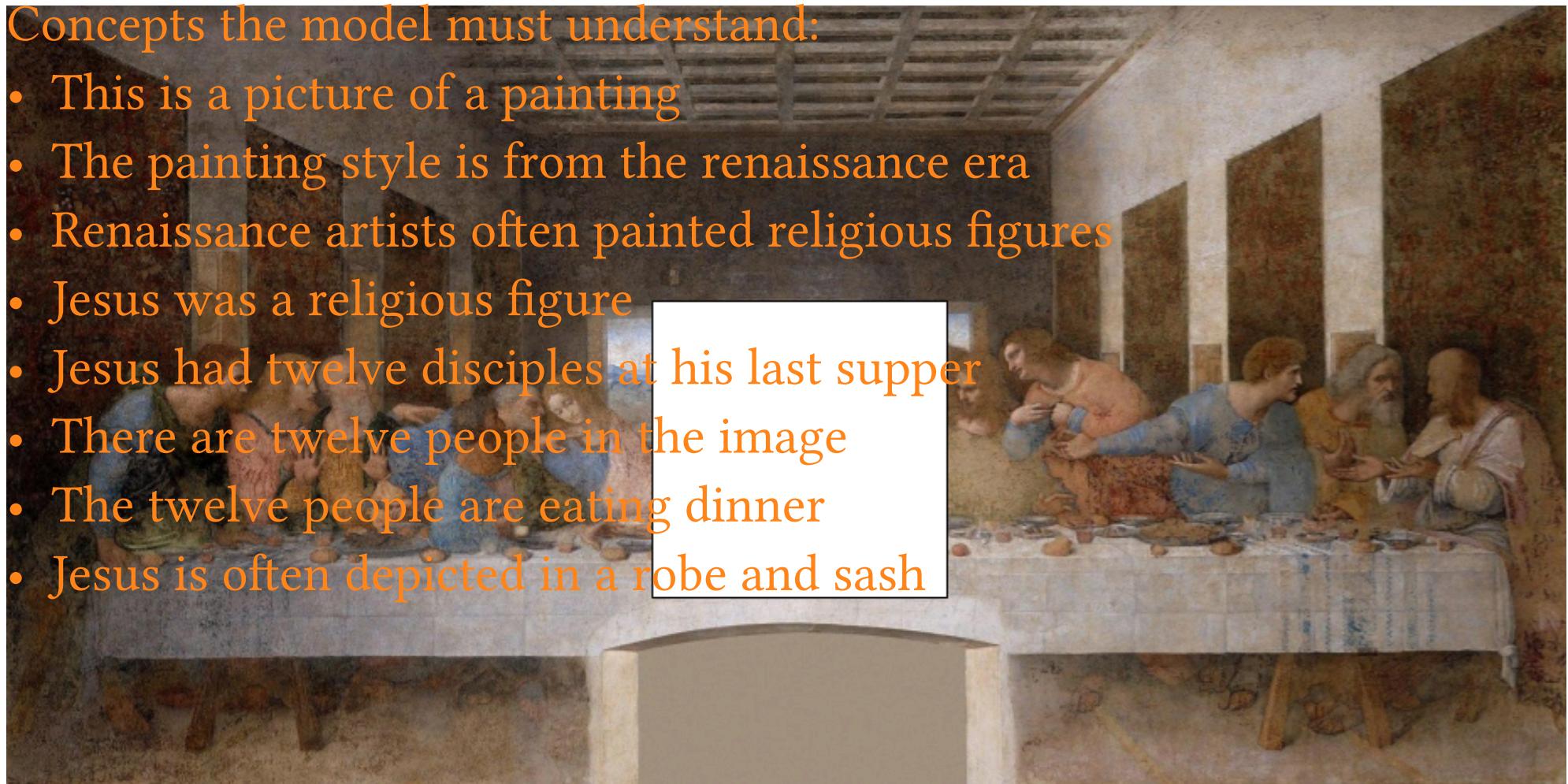
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures
- Jesus was a religious figure
- Jesus had twelve disciples at his last supper
- There are twelve people in the image
- The twelve people are eating dinner



Unsupervised Training

Concepts the model must understand:

- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures
- Jesus was a religious figure
- Jesus had twelve disciples at his last supper
- There are twelve people in the image
- The twelve people are eating dinner
- Jesus is often depicted in a robe and sash



Unsupervised Training

We train the model to fix the image

Unsupervised Training

We train the model to fix the image

To fix the image, the model must understand so much of our world

Unsupervised Training

We train the model to fix the image

To fix the image, the model must understand so much of our world

This is the power of generative pre-training

Unsupervised Training

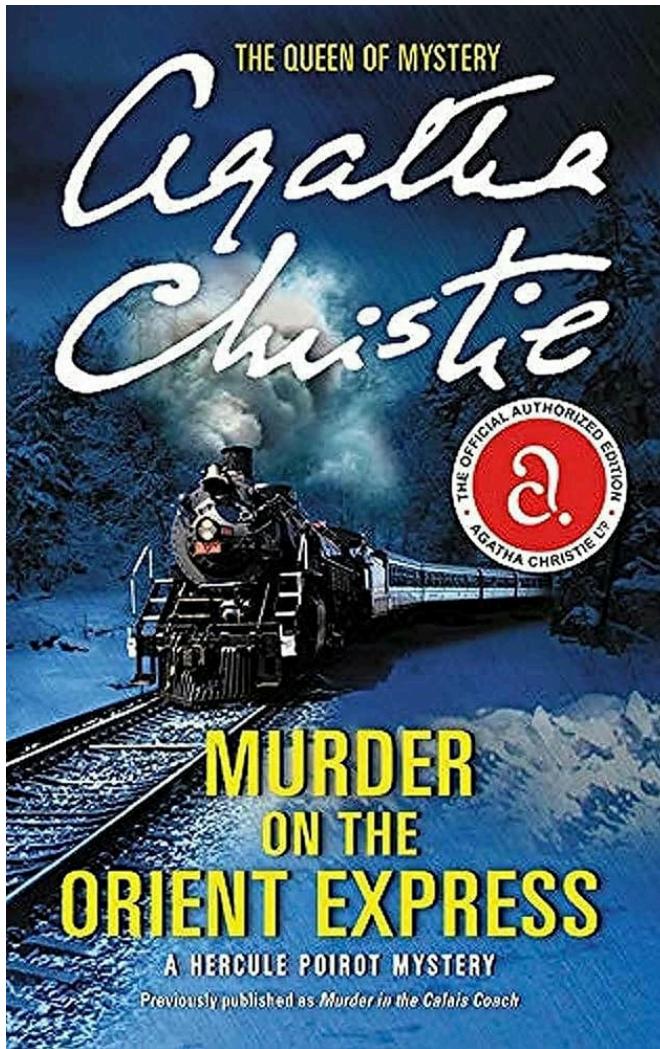
We train the model to fix the image

To fix the image, the model must understand so much of our world

This is the power of generative pre-training

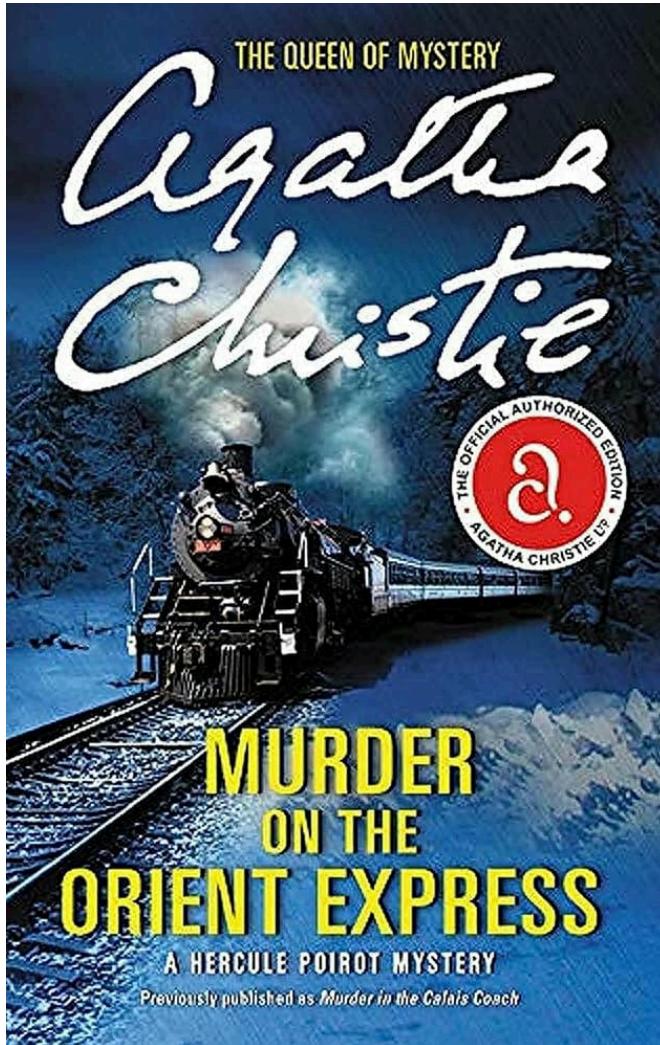
What about text transformers?

Unsupervised Training



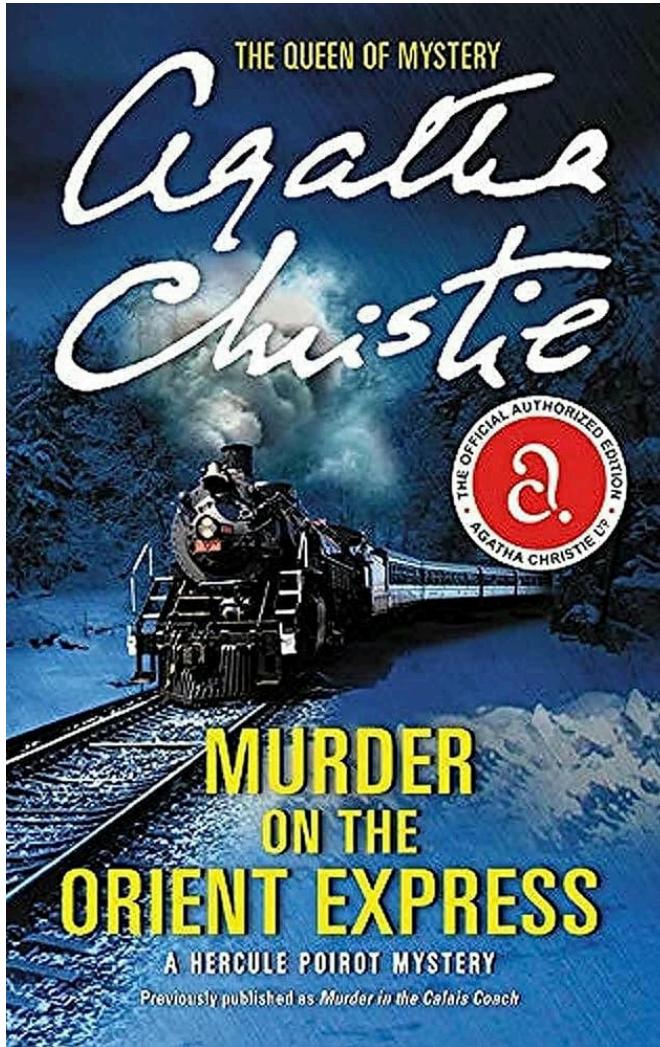
This is a mystery novel

Unsupervised Training



This is a mystery novel
Clues, intrigue, murder, etc

Unsupervised Training

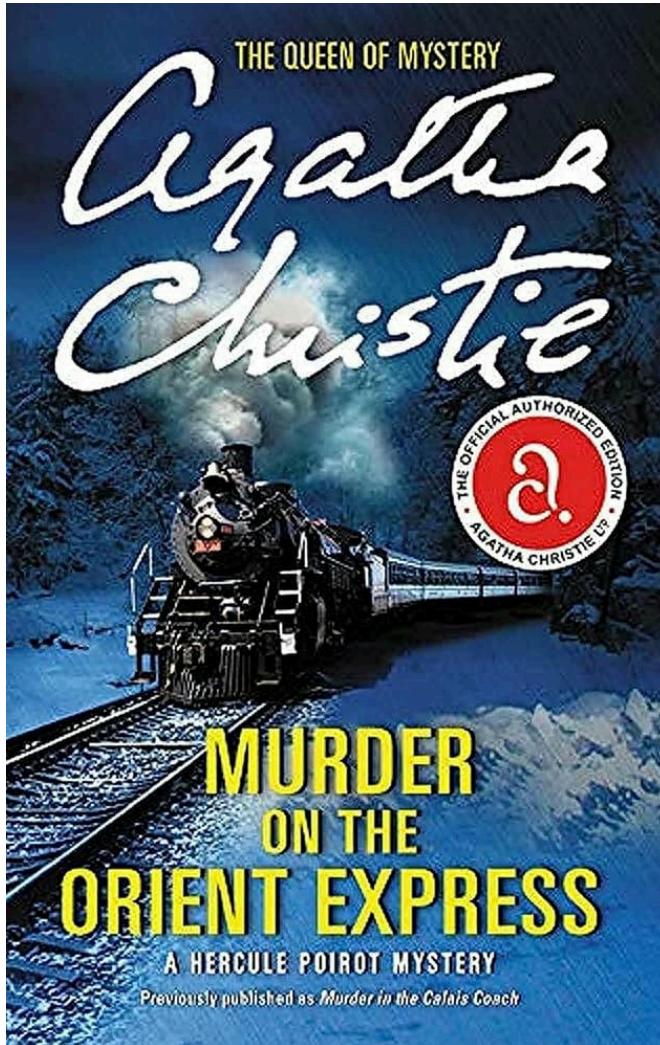


This is a mystery novel

Clues, intrigue, murder, etc

“Ah, said inspector Poirot, the murderer must be _____.”

Unsupervised Training



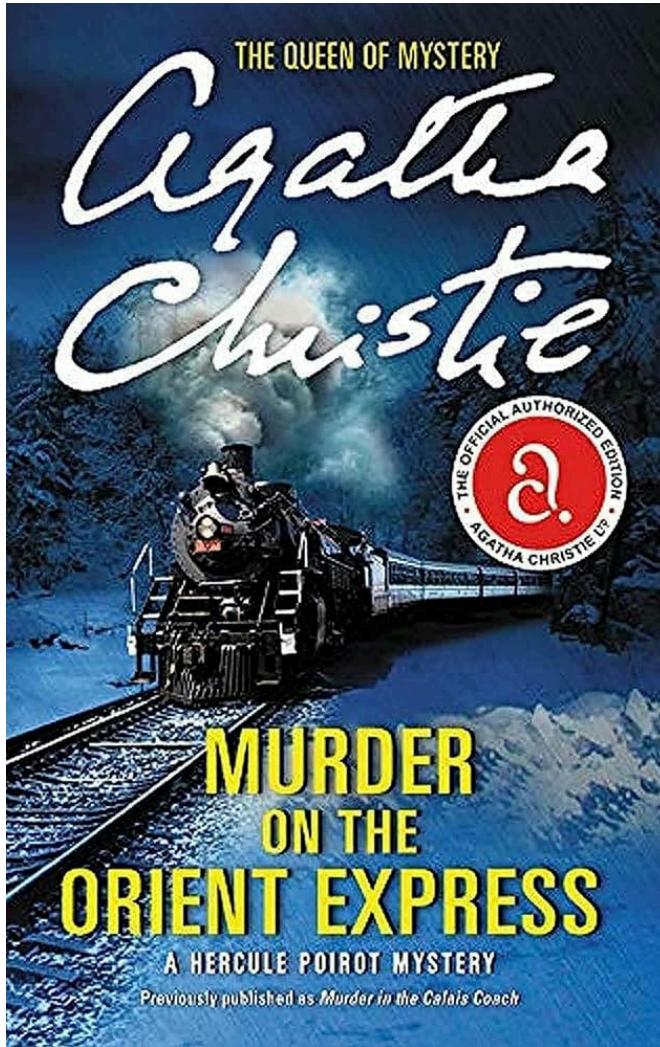
This is a mystery novel

Clues, intrigue, murder, etc

“Ah, said inspector Poirot, the murderer must be _____.”

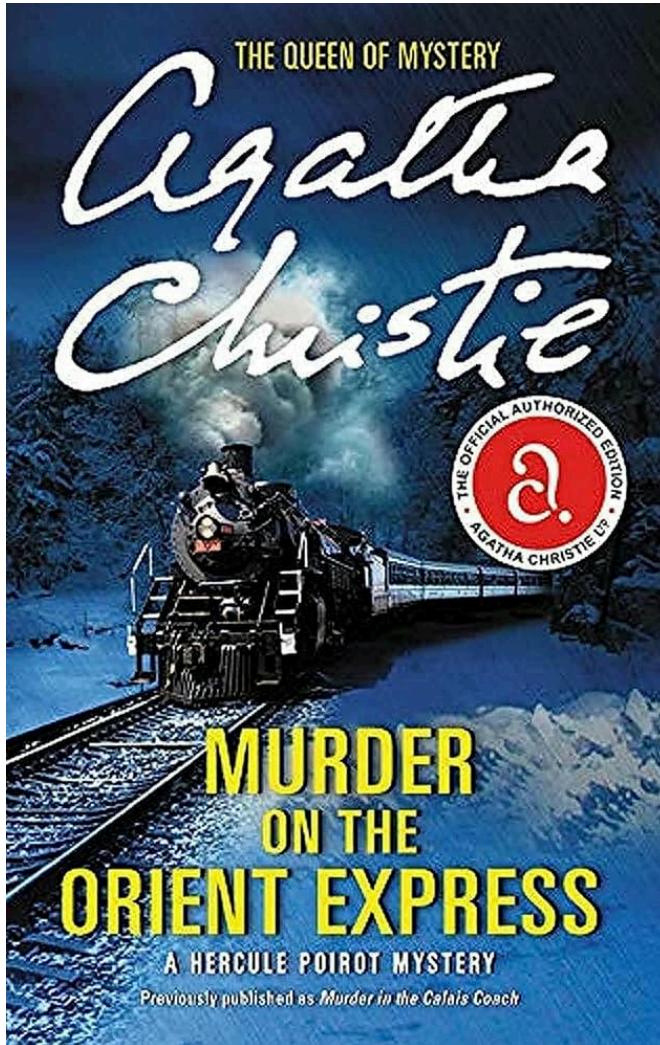
To complete the sentence, the model must understand:

Unsupervised Training



To complete the sentence, the model must understand:

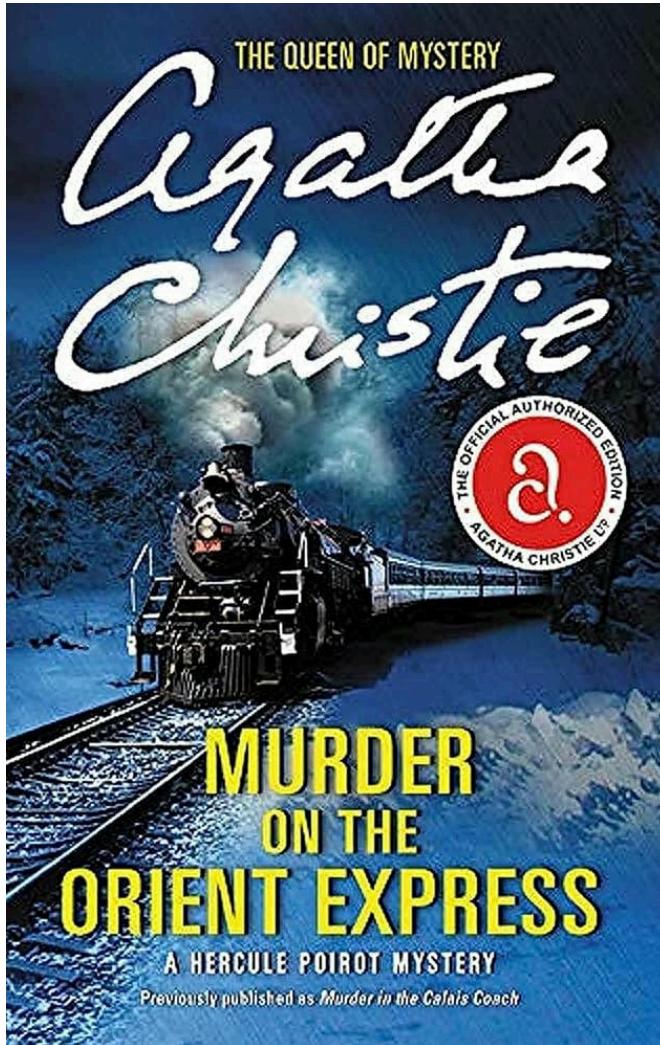
Unsupervised Training



To complete the sentence, the model must understand:

- What a murder is

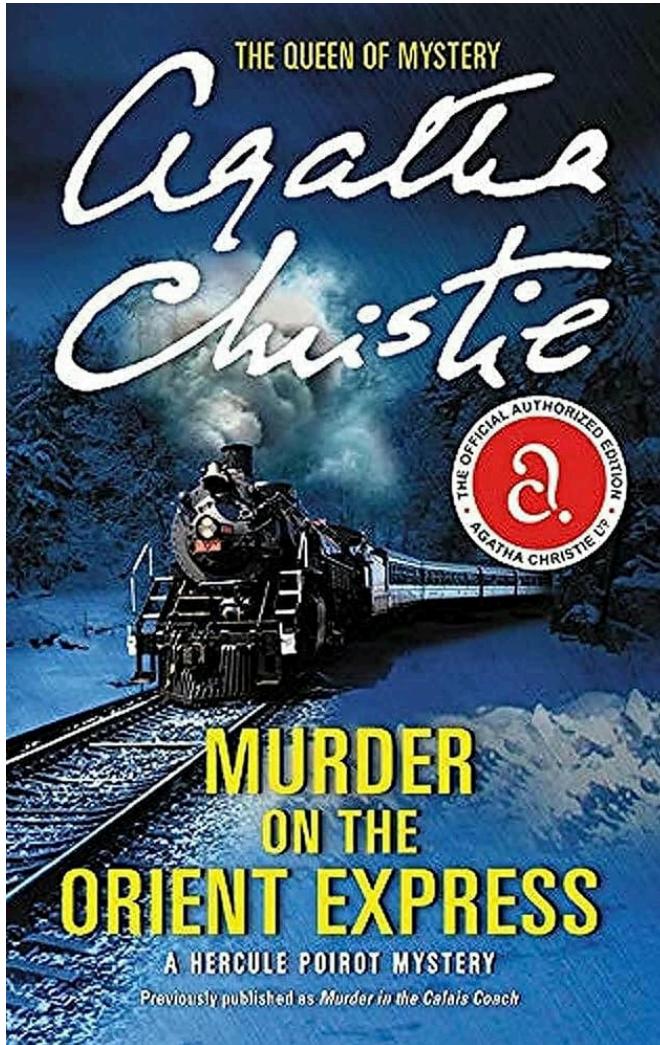
Unsupervised Training



To complete the sentence, the model must understand:

- What a murder is
- What it means to be alive

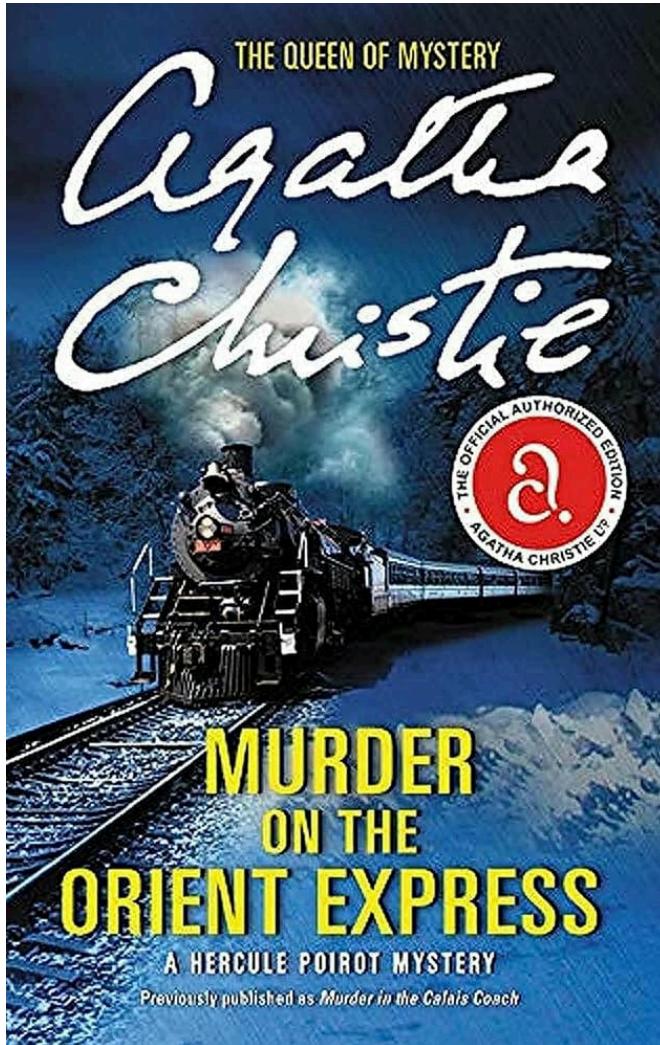
Unsupervised Training



To complete the sentence, the model must understand:

- What a murder is
- What it means to be alive
- Emotions like anger, jealousy, betrayal, love

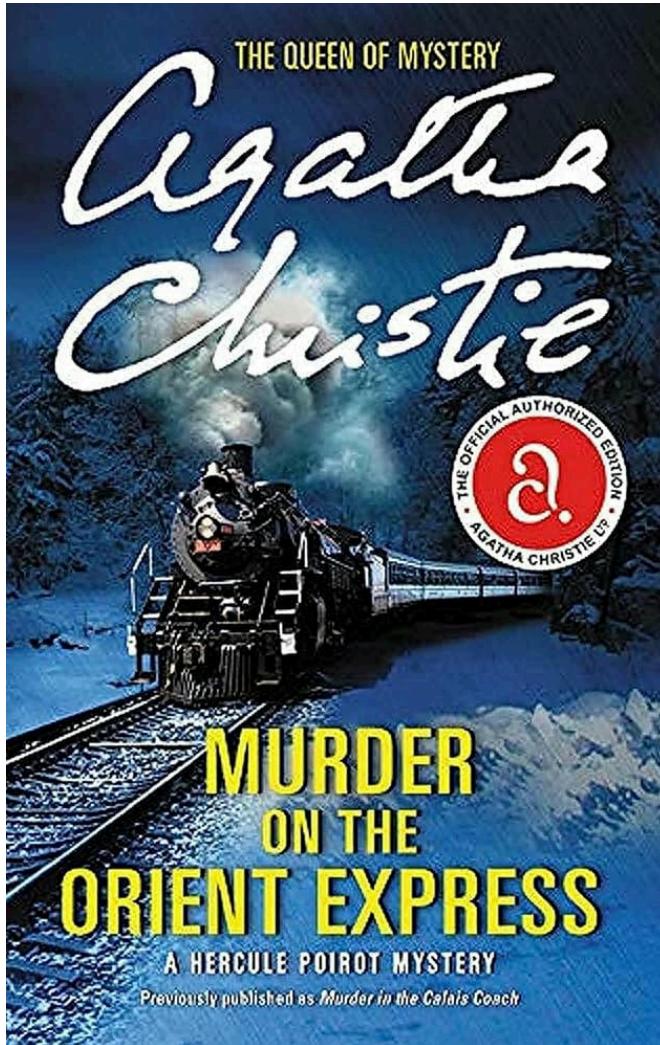
Unsupervised Training



To complete the sentence, the model must understand:

- What a murder is
- What it means to be alive
- Emotions like anger, jealousy, betrayal, love
- Personalities of each character

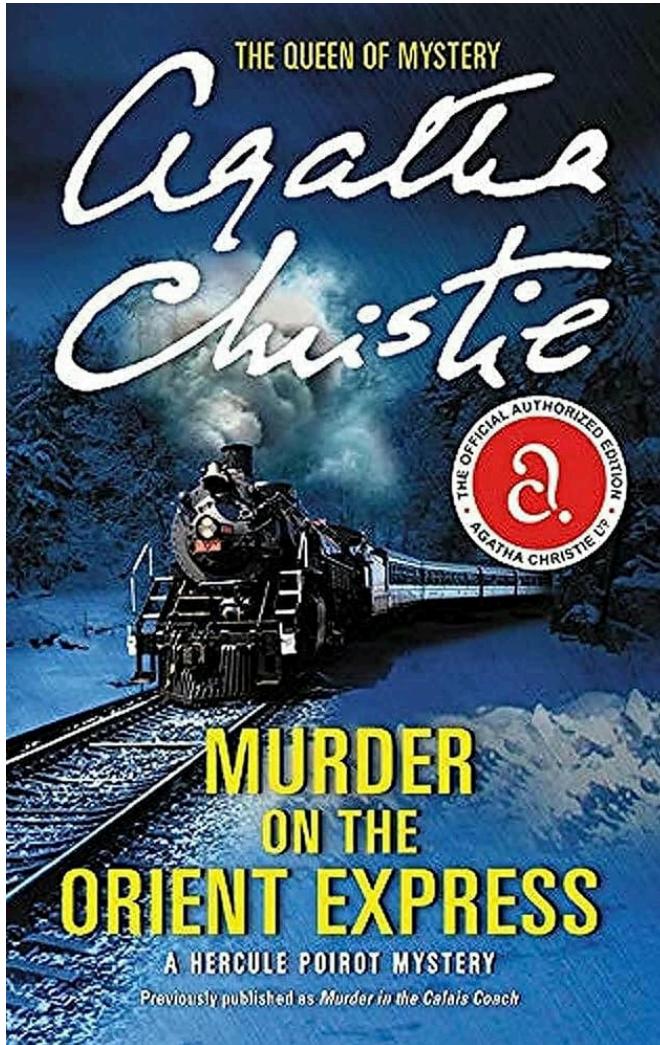
Unsupervised Training



To complete the sentence, the model must understand:

- What a murder is
- What it means to be alive
- Emotions like anger, jealousy, betrayal, love
- Personalities of each character
- Why a human would murder another human

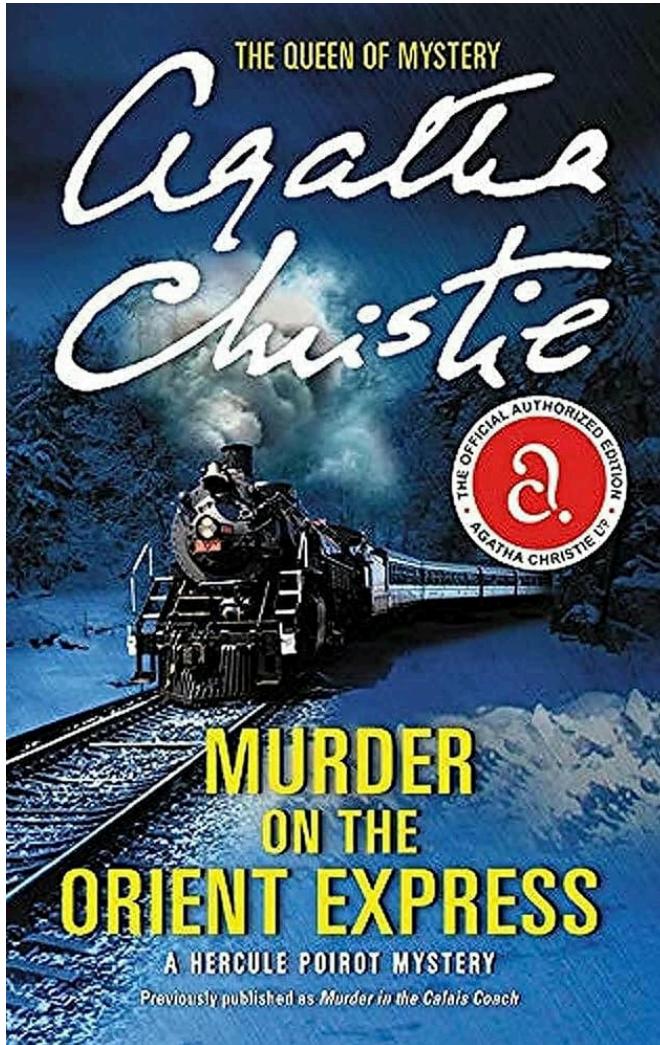
Unsupervised Training



To complete the sentence, the model must understand:

- What a murder is
- What it means to be alive
- Emotions like anger, jealousy, betrayal, love
- Personalities of each character
- Why a human would murder another human
- How humans react to emotions

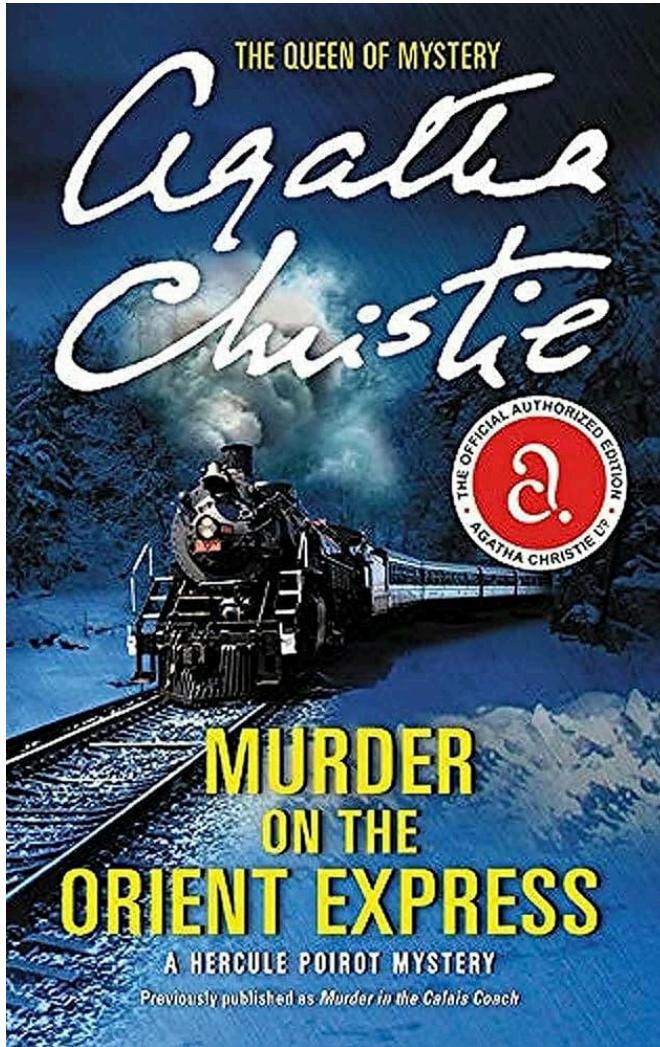
Unsupervised Training



To complete the sentence, the model must understand:

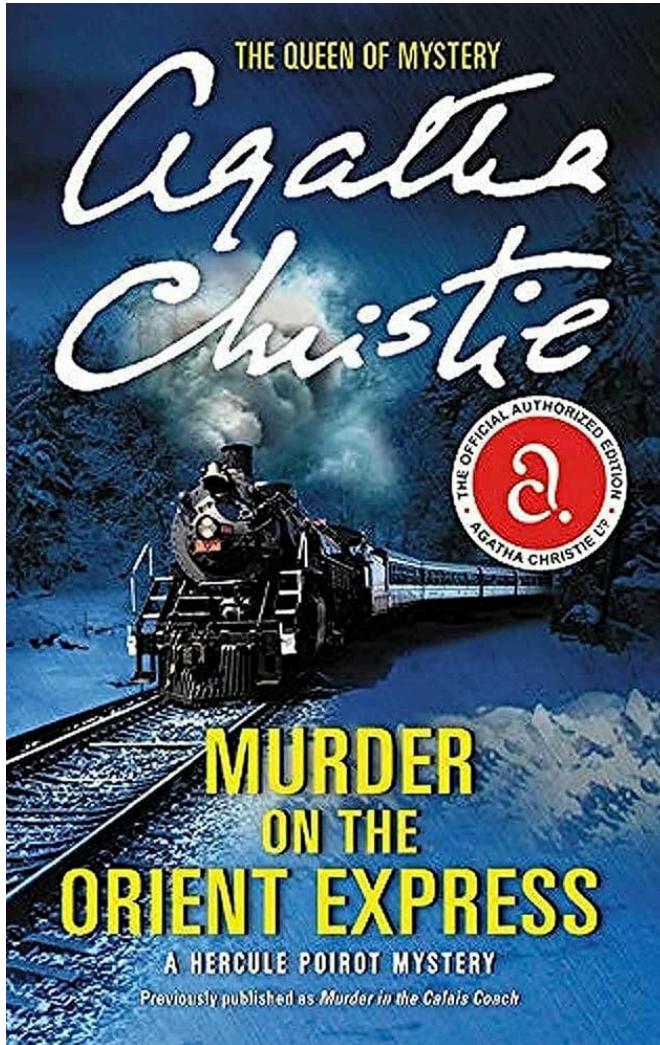
- What a murder is
- What it means to be alive
- Emotions like anger, jealousy, betrayal, love
- Personalities of each character
- Why a human would murder another human
- How humans react to emotions
- How to tell if someone lies

Unsupervised Training



To predict the murderer, the model must understand so much about humans and our society

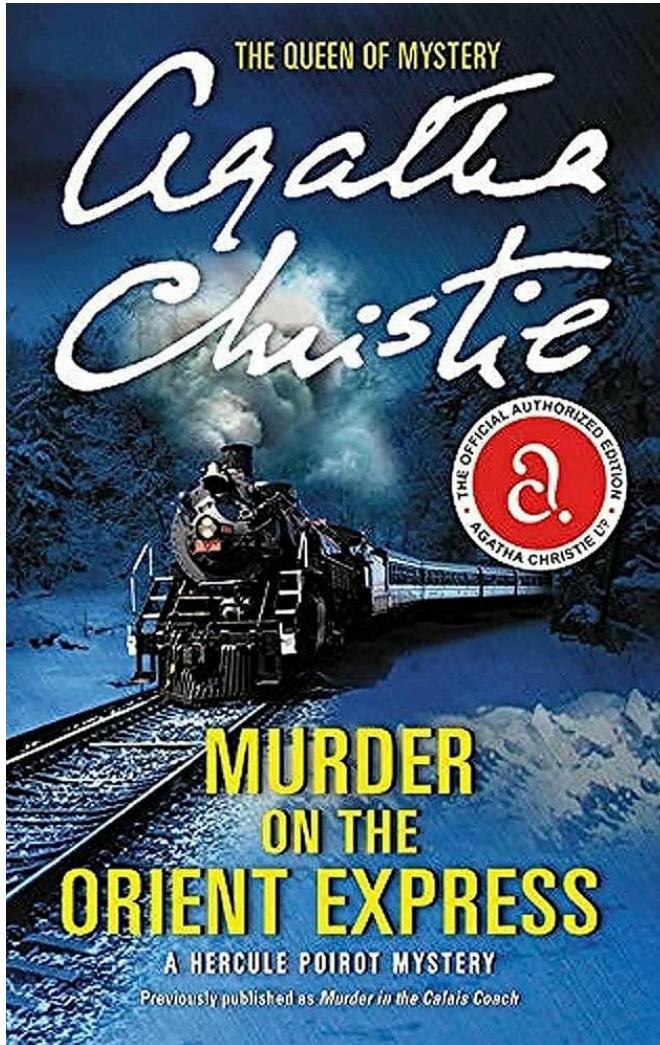
Unsupervised Training



To predict the murderer, the model must understand so much about humans and our society

The Books3 dataset contains 200,000 books

Unsupervised Training



To predict the murderer, the model must understand so much about humans and our society

The Books3 dataset contains 200,000 books

We train the model to predict the ending of all these books

Unsupervised Training

We can apply this same concept to:

Unsupervised Training

We can apply this same concept to:

- Predict missing base pairs in a strand of DNA

Unsupervised Training

We can apply this same concept to:

- Predict missing base pairs in a strand of DNA
- Predict missing audio from a song

Unsupervised Training

We can apply this same concept to:

- Predict missing base pairs in a strand of DNA
- Predict missing audio from a song
- Predict the outcome of particle collisions at the Large Hadron Supercollider

Unsupervised Training

We can apply this same concept to:

- Predict missing base pairs in a strand of DNA
- Predict missing audio from a song
- Predict the outcome of particle collisions at the Large Hadron Supercollider

All we need is a large enough dataset!

Unsupervised Training

What if we put the model in a robot?

Unsupervised Training

What if we put the model in a robot?

Give the model what the robot sees and what the robot does

Unsupervised Training

What if we put the model in a robot?

Give the model what the robot sees and what the robot does

Predict what the robot will see next

Unsupervised Training

What if we put the model in a robot?

Give the model what the robot sees and what the robot does

Predict what the robot will see next

We call this a **world model**

Unsupervised Training

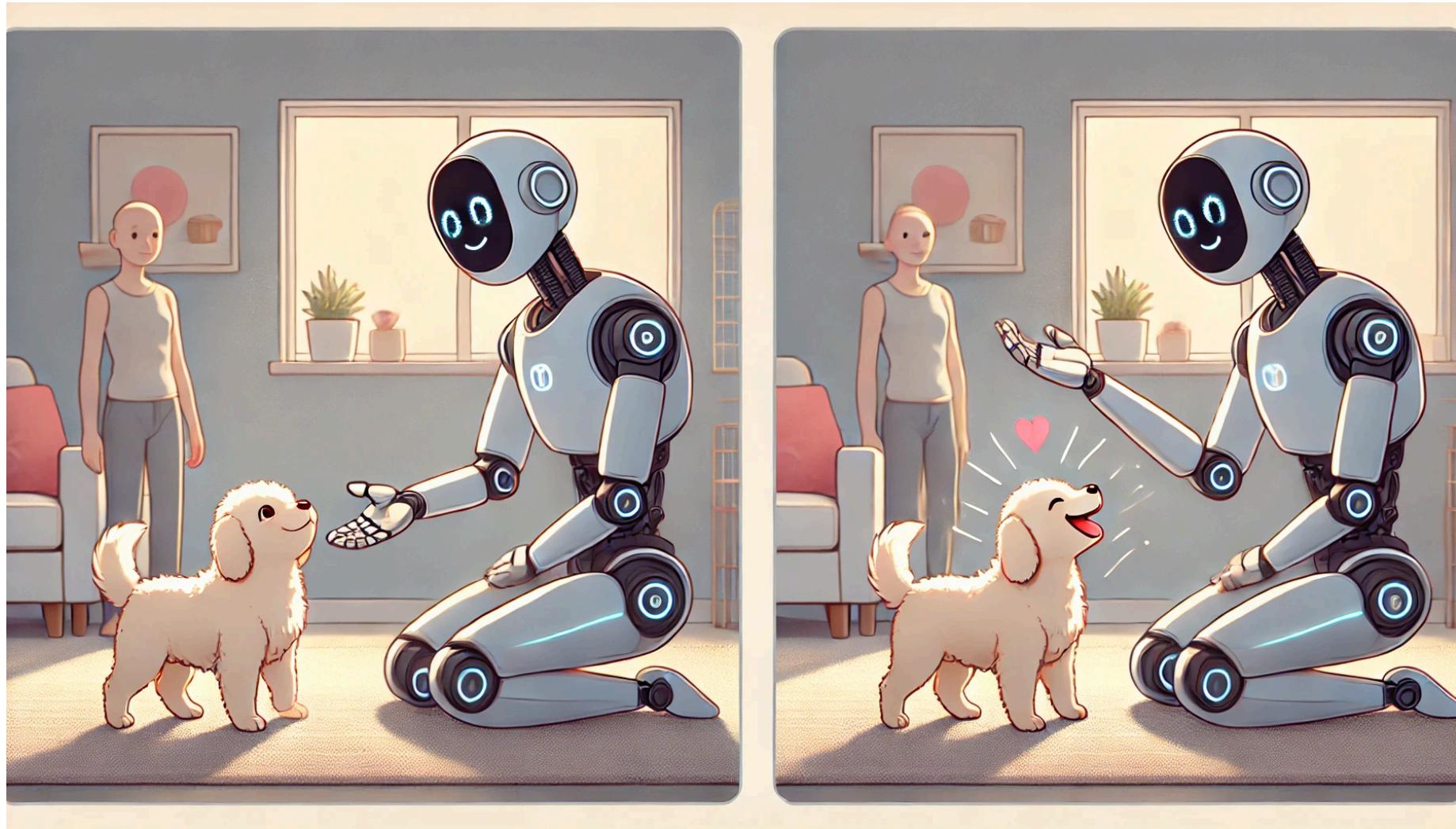
What if we put the model in a robot?

Give the model what the robot sees and what the robot does

Predict what the robot will see next

We call this a **world model**

Unsupervised Training



Unsupervised Training

Soon, I will apply for a grant to train a world model

Unsupervised Training

Soon, I will apply for a grant to train a world model

If I win, I will need help creating a robot dataset

Unsupervised Training

Soon, I will apply for a grant to train a world model

If I win, I will need help creating a robot dataset

I will need some humans to control our robots in the world

Unsupervised Training

Soon, I will apply for a grant to train a world model

If I win, I will need help creating a robot dataset

I will need some humans to control our robots in the world

If you are interested, give me your email after class

Unsupervised Training

These transformers learn and understand the structure of our world

Unsupervised Training

These transformers learn and understand the structure of our world

But their understanding is trapped

Unsupervised Training

These transformers learn and understand the structure of our world

But their understanding is trapped

They can only finish sentences or complete pictures

Unsupervised Training

These transformers learn and understand the structure of our world

But their understanding is trapped

They can only finish sentences or complete pictures

How can we use this strong understanding to help humans?

Unsupervised Training

These transformers learn and understand the structure of our world

But their understanding is trapped

They can only finish sentences or complete pictures

How can we use this strong understanding to help humans?

- Identify pictures of cancer

Unsupervised Training

These transformers learn and understand the structure of our world

But their understanding is trapped

They can only finish sentences or complete pictures

How can we use this strong understanding to help humans?

- Identify pictures of cancer
- Make scientific discoveries

Unsupervised Training

These transformers learn and understand the structure of our world

But their understanding is trapped

They can only finish sentences or complete pictures

How can we use this strong understanding to help humans?

- Identify pictures of cancer
- Make scientific discoveries
- Minimize human suffering

Unsupervised Training

These transformers learn and understand the structure of our world

But their understanding is trapped

They can only finish sentences or complete pictures

How can we use this strong understanding to help humans?

- Identify pictures of cancer
- Make scientific discoveries
- Minimize human suffering

Today, we use **reinforcement learning**

We will formally introduce reinforcement learning next lecture

Closing Remarks

Closing Remarks

This is the last in-person lecture

Closing Remarks

This is the last in-person lecture

I will record a video on reinforcement learning next week

Closing Remarks

This is the last in-person lecture

I will record a video on reinforcement learning next week

I will be here from 7:00PM on December 2 for questions/discussin on reinforcement learning

Closing Remarks

In this course, we started from Gauss in 1795

Closing Remarks

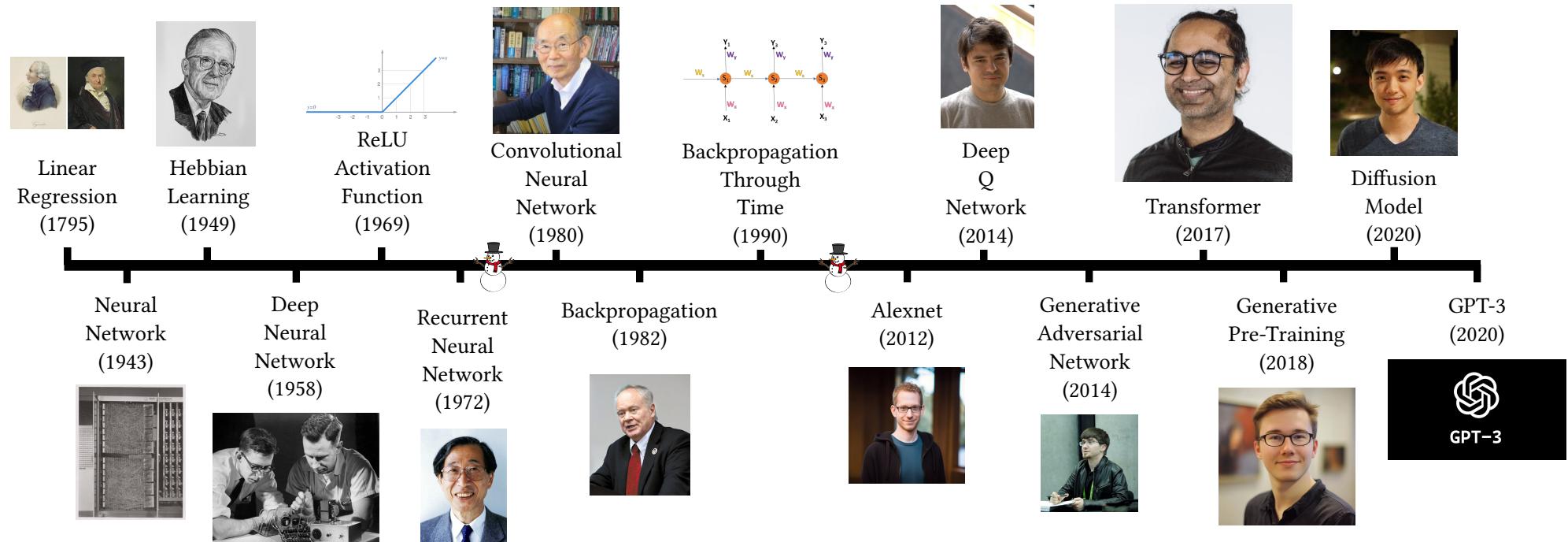
In this course, we started from Gauss in 1795

We built up concepts until we reached the modern age

Closing Remarks

In this course, we started from Gauss in 1795

We built up concepts until we reached the modern age



Closing Remarks

We learned about:

Closing Remarks

We learned about:

- Linear regression

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent
- RMSProp and Adam

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent
- RMSProp and Adam
- Convolutional neural networks

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent
- RMSProp and Adam
- Convolutional neural networks
- Composite memory

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent
- RMSProp and Adam
- Convolutional neural networks
- Composite memory
- Recurrent neural networks

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent
- RMSProp and Adam
- Convolutional neural networks
- Composite memory
- Recurrent neural networks
- Autoencoders

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent
- RMSProp and Adam
- Convolutional neural networks
- Composite memory
- Recurrent neural networks
- Autoencoders
- Variational autoencoders

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent
- RMSProp and Adam
- Convolutional neural networks
- Composite memory
- Recurrent neural networks
- Autoencoders
- Variational autoencoders
- Graph neural networks

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent
- RMSProp and Adam
- Convolutional neural networks
- Composite memory
- Recurrent neural networks
- Autoencoders
- Variational autoencoders
- Graph neural networks
- Attention and transformers

Closing Remarks

We learned about:

- Linear regression
- Polynomial regression
- Biological neurons
- Artificial neurons
- Perceptrons
- Backpropagation
- Gradient descent
- Classification
- Parameter initialization
- Many activation functions
- Stochastic gradient descent
- RMSProp and Adam
- Convolutional neural networks
- Composite memory
- Recurrent neural networks
- Autoencoders
- Variational autoencoders
- Graph neural networks
- Attention and transformers
- Generative pre-training

Closing Remarks

I hope you enjoyed the course!

Closing Remarks

I hope you enjoyed the course!

But there are many more topics to learn!

Closing Remarks

I hope you enjoyed the course!

But there are many more topics to learn!

Now, you have the tools to study deep learning on your own

Closing Remarks

I hope you enjoyed the course!

But there are many more topics to learn!

Now, you have the tools to study deep learning on your own

You have the tools to train neural networks for real problems

Closing Remarks

I hope you enjoyed the course!

But there are many more topics to learn!

Now, you have the tools to study deep learning on your own

You have the tools to train neural networks for real problems

Closing Remarks

In the first lecture, I asked everyone in this class for something

Closing Remarks

In the first lecture, I asked everyone in this class for something

Question: Do you remember what it was?

Closing Remarks

Deep learning is a powerful tool

Closing Remarks

Deep learning is a powerful tool

Like all powerful tools, deep learning can be used for good or evil

Closing Remarks

Deep learning is a powerful tool

Like all powerful tools, deep learning can be used for good or evil

- COVID-19 vaccine

Closing Remarks

Deep learning is a powerful tool

Like all powerful tools, deep learning can be used for good or evil

- COVID-19 vaccine
- Creating art

Closing Remarks

Deep learning is a powerful tool

Like all powerful tools, deep learning can be used for good or evil

- COVID-19 vaccine
- Creating art
- Autonomous driving

Closing Remarks

Deep learning is a powerful tool

Like all powerful tools, deep learning can be used for good or evil

- COVID-19 vaccine
- Creating art
- Autonomous driving
- Making DeepFakes

Closing Remarks

Deep learning is a powerful tool

Like all powerful tools, deep learning can be used for good or evil

- COVID-19 vaccine
- Creating art
- Autonomous driving
- Making DeepFakes
- Weapon guidance systems

Closing Remarks

Deep learning is a powerful tool

Like all powerful tools, deep learning can be used for good or evil

- COVID-19 vaccine
- Creating art
- Autonomous driving
- Making DeepFakes
- Weapon guidance systems
- Discrimination

Closing Remarks

Deep learning is a powerful tool

Like all powerful tools, deep learning can be used for good or evil

- COVID-19 vaccine
- Creating art
- Autonomous driving
- Making DeepFakes
- Weapon guidance systems
- Discrimination

Before training a model, think about whether it is good or bad for the world

Course Evaluation

Course Evaluation

Department instructed me to ask you for course feedback

Course Evaluation

Department instructed me to ask you for course feedback

We take this feedback seriously

Course Evaluation

Department instructed me to ask you for course feedback

We take this feedback seriously

Your feedback will impact future courses (and my job)

Course Evaluation

Department instructed me to ask you for course feedback

We take this feedback seriously

Your feedback will impact future courses (and my job)

If you like the course, please say it!

Course Evaluation

Department instructed me to ask you for course feedback

We take this feedback seriously

Your feedback will impact future courses (and my job)

If you like the course, please say it!

Please be specific on what you like and do not like

Course Evaluation

Department instructed me to ask you for course feedback

We take this feedback seriously

Your feedback will impact future courses (and my job)

If you like the course, please say it!

Please be specific on what you like and do not like

Your likes/dislikes will change your future courses

Course Evaluation

I must leave the room to let you fill out this form

Course Evaluation

I must leave the room to let you fill out this form

Please scan the QR code and complete the survey

Course Evaluation

I must leave the room to let you fill out this form

Please scan the QR code and complete the survey

Department has suggested 10 minutes

Course Evaluation

I must leave the room to let you fill out this form

Please scan the QR code and complete the survey

Department has suggested 10 minutes

I will return in 10 minutes

Course Evaluation

I must leave the room to let you fill out this form

Please scan the QR code and complete the survey

Department has suggested 10 minutes

I will return in 10 minutes

<https://isw.um.edu.mo/siaweb>

Course Evaluation

If you participated in class come see me after class

Course Evaluation

If you participated in class come see me after class

- Answered a question

Course Evaluation

If you participated in class come see me after class

- Answered a question
- Asked a question