



Transformers

CISC 7026 - Introduction to Deep Learning

Steven Morad

University of Macau

Review	2
Going Deeper	18
Transformers	35
Positional Encoding	39
Applications	60
Text Transformers	62
Image Transformers	69
Generative Pre-Training	73
Closing Remarks	93
Course Evaluation	100

Review

Review

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

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

So we introduced a forgetting term γ

Review

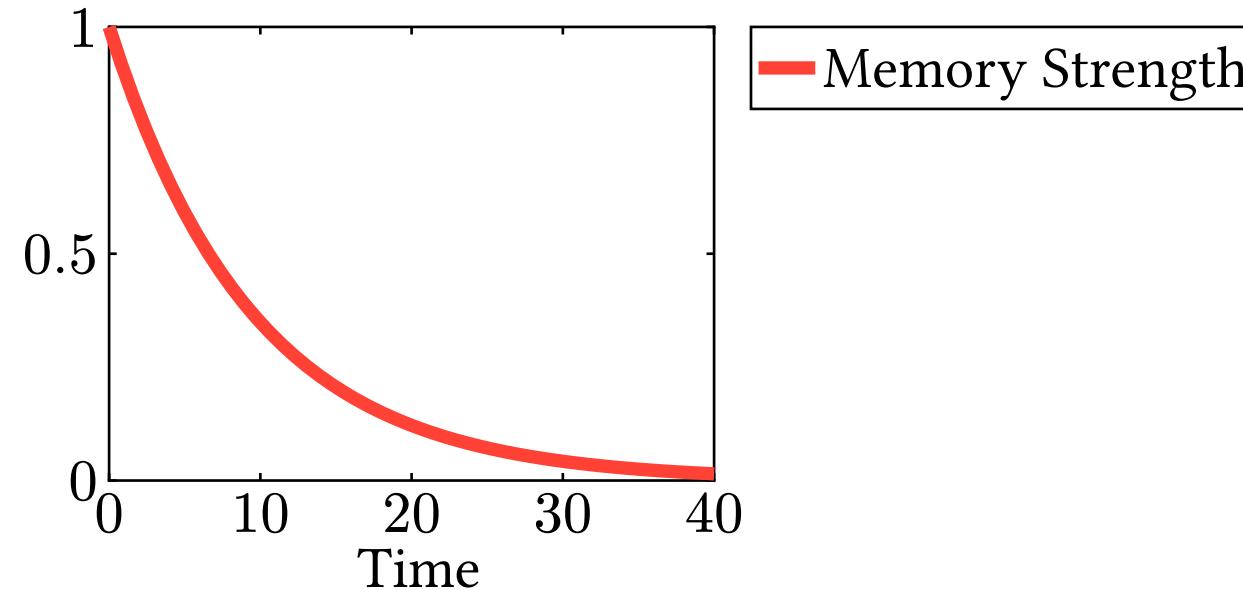
So we introduced a forgetting term γ

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

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

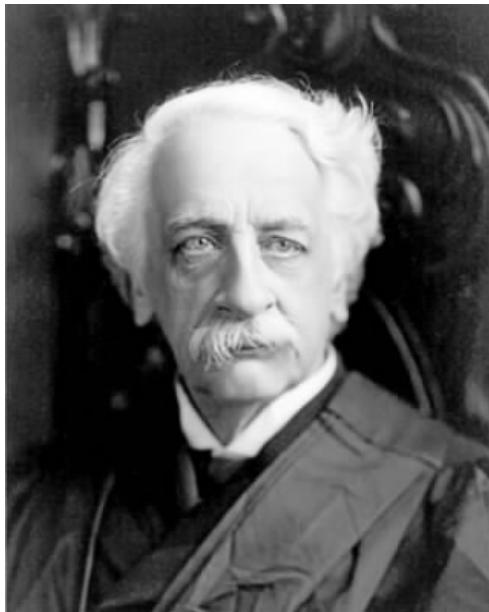
We went to a party and the forgetting seemed ok



10 PM

Review

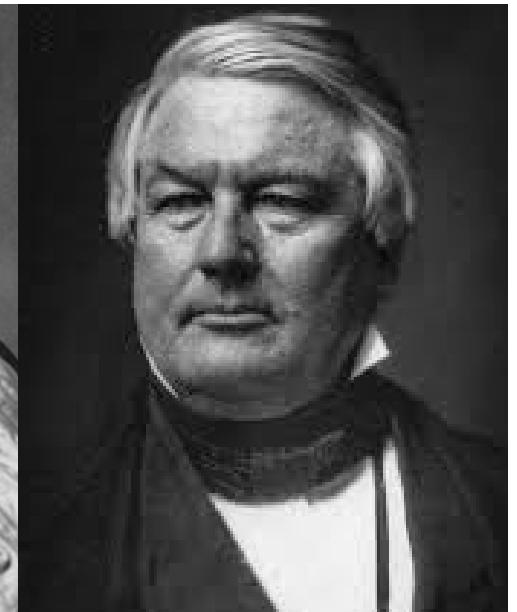
We went to a party and the forgetting seemed ok



10 PM

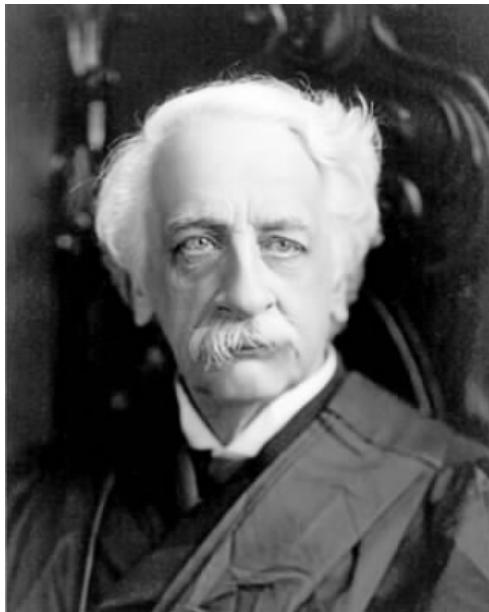


11 PM



Review

We went to a party and the forgetting seemed ok



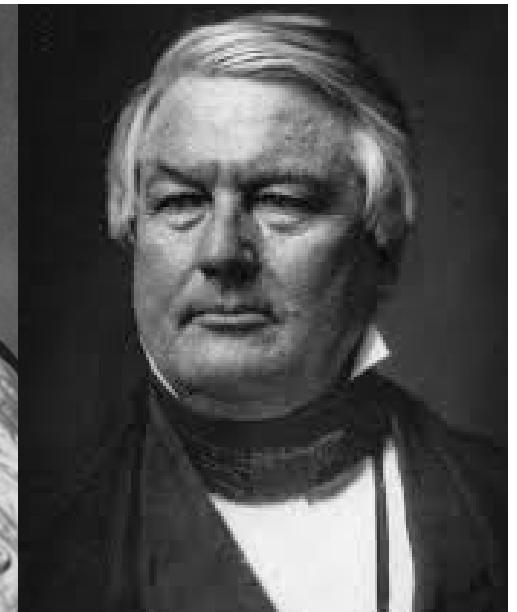
10 PM



11 PM

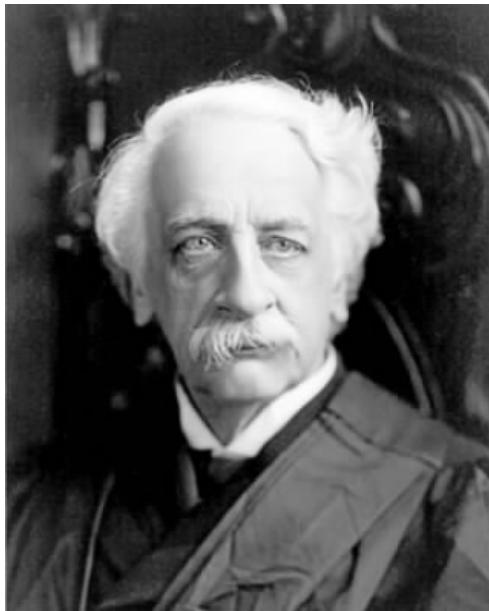


12 AM



Review

We went to a party and the forgetting seemed ok



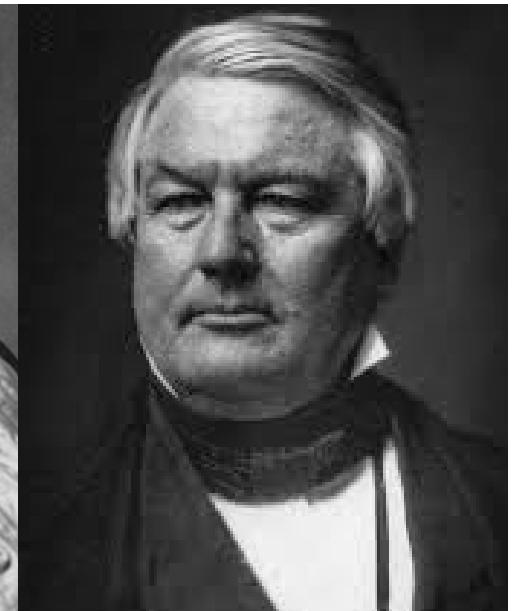
10 PM



11 PM



12 AM



1 AM

Review



Review



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

Review



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

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

Review



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

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

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

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

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



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

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

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

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

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

Review



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

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

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

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

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

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

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

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

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

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

With attention, we can create the **transformer**

Review

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

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

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

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

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

With attention, we can create the **transformer**

Review

Why should we care about the transformer?

Review

Why should we care about the transformer?

It is arguably the most powerful neural network architecture today

Review

Why should we care about the transformer?

It is arguably the most powerful neural network architecture today

- AlphaFold (Nobel prize)

Review

Why should we care about the transformer?

It is arguably the most powerful neural network architecture today

- AlphaFold (Nobel prize)
- ChatGPT, Qwen, LLaMA, etc

Review

Why should we care about the transformer?

It is arguably the most powerful neural network architecture today

- AlphaFold (Nobel prize)
- ChatGPT, Qwen, LLaMA, etc
- DinoV2

Review

Why should we care about the transformer?

It is arguably the most powerful neural network architecture today

- AlphaFold (Nobel prize)
- ChatGPT, Qwen, LLaMA, etc
- DinoV2

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

It is often more efficient to create deeper networks instead

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

It is often more efficient to create deeper networks instead

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

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

Going Deeper

Making the network too deep can hurt performance

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

At each layer, we lose a little bit of information

Going Deeper

Making the network too deep can hurt performance

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

At each layer, we lose a little bit of information

With enough layers, all the information in \mathbf{x} is lost!

Going Deeper

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

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 $f(x) = x$

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

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 $f(x) = x$

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

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

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 information from x is present in every layer, 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 information from x is present in every layer, then learning the identity function should be very easy!

Question: How can we prevent x from getting lost?

Going Deeper

We can feed x to each layer

Going Deeper

We can feed x to each layer

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

Going Deeper

We can feed x 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

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

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

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

⋮

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

This allows information to flow around the layers

Fewer parameters than a DenseNet, but performs slightly worse

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

ResNets use a **residual connection**

Going Deeper

ResNets use a **residual connection**

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

Going Deeper

ResNets use a **residual connection**

$$z = \underbrace{f(x, \theta)}_{\text{Residual}} + x$$

Going Deeper

ResNets use a **residual connection**

$$z = \underbrace{f(x, \theta)}_{\text{Residual}} + x$$

Think of the residual as a small change to x

Going Deeper

ResNets use a **residual connection**

$$z = \underbrace{f(x, \theta)}_{\text{Residual}} + x$$

Think of the residual as a small change to x

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

Going Deeper

ResNets use a **residual connection**

$$z = \underbrace{f(x, \theta)}_{\text{Residual}} + x$$

Think of the residual as a small change to x

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

Each layer is a small perturbation of x

Going Deeper

ResNets use a **residual connection**

$$z = \underbrace{f(x, \theta)}_{\text{Residual}} + x$$

Think of the residual as a small change to x

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

Each layer is a small perturbation of x

This prevents 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

And the magnitude of the outputs impacts the gradient

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

And the magnitude of the outputs impacts the gradient

$$f_1(x, \theta_1) = \sum_{i=1}^{d_x} \theta_{1,i} \cdot x_i$$

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

And the magnitude of the outputs impacts the gradient

$$f_1(x, \theta_1) = \sum_{i=1}^{d_x} \theta_{1,i} \cdot x_i$$

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

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

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

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?

Answer: We can normalize the output of each layer

Going Deeper

Layer normalization normalizes the output of a layer

Going Deeper

Layer normalization normalizes the output of a layer

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

Going Deeper

Layer normalization normalizes the output of a layer

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

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

Going Deeper

Layer normalization normalizes the output of a layer

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

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

$$f(\mathbf{x}, \boldsymbol{\theta}) - \mu$$

Going Deeper

Layer normalization normalizes the output of a layer

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

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

$$f(\mathbf{x}, \boldsymbol{\theta}) - \mu$$

Question: What does this do?

Going Deeper

Layer normalization normalizes the output of a layer

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

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

$$f(\mathbf{x}, \boldsymbol{\theta}) - \mu$$

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 output by standard deviation

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 output by standard deviation

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

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 output by standard deviation

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

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

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 output by standard deviation

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

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

Layer norm makes the output normally distributed, $y_i \sim \mathcal{N}(0, 1)$

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.layer1 = TransformerLayer()
        self.layer2 = TransformerLayer()
        self.layer3 = TransformerLayer()

    def forward(self, x):
        x = self.layer1(x)
        x = self.layer2(x)
        x = self.layer3(x)
        return x
```

Transformers

```
class Transformer(nn.Module):
    def __init__(self):
        self.layer1 = TransformerLayer()
        self.layer2 = TransformerLayer()
        self.layer3 = TransformerLayer()

    def forward(self, x):
        x = self.layer1(x)
        x = self.layer2(x)
        x = self.layer3(x)
        return x
```

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

Transformers

```
class Transformer(nn.Module):
    def __init__(self):
        self.layer1 = TransformerLayer()
        self.layer2 = TransformerLayer()
        self.layer3 = TransformerLayer()

    def forward(self, x):
        x = self.layer1(x)
        x = self.layer2(x)
        x = self.layer3(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?

$$f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \theta\right) \stackrel{?}{=} f\left(\begin{bmatrix} x_2 \\ x_1 \end{bmatrix}, \theta\right)$$

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?

$$f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \theta\right) \stackrel{?}{=} f\left(\begin{bmatrix} x_2 \\ x_1 \end{bmatrix}, \theta\right)$$

Answer: It depends. 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 certain tasks, the order of inputs matters

Positional Encoding

For certain tasks, the order of inputs matters

Question: Does the order matter to the transformer?

$$f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \theta\right) \stackrel{?}{=} f\left(\begin{bmatrix} x_2 \\ x_1 \end{bmatrix}, \theta\right)$$

Positional Encoding

For certain tasks, the order of inputs matters

Question: Does the order matter to the transformer?

$$f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \theta\right) \stackrel{?}{=} f\left(\begin{bmatrix} x_2 \\ x_1 \end{bmatrix}, \theta\right)$$

Let us find out!

Positional Encoding

For certain tasks, the order of inputs matters

Question: Does the order matter to the transformer?

$$f\left(\begin{bmatrix}x_1 \\ x_2\end{bmatrix}, \theta\right) \stackrel{?}{=} f\left(\begin{bmatrix}x_2 \\ x_1\end{bmatrix}, \theta\right)$$

Let us find out!

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

Positional Encoding

For certain tasks, the order of inputs matters

Question: Does the order matter to the transformer?

$$f\left(\begin{bmatrix}x_1 \\ x_2\end{bmatrix}, \theta\right) \stackrel{?}{=} f\left(\begin{bmatrix}x_2 \\ x_1\end{bmatrix}, \theta\right)$$

Let us find out!

First, 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

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

Order **does** matter (not equivariant)

Positional Encoding

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

Order **does** matter (not equivariant)

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

Positional Encoding

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

Order **does** matter (not equivariant)

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

Order **does not** matter (equivariant)

Positional Encoding

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

Order **does** matter (not equivariant)

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

Order **does not** matter (equivariant)

Which is a transformer?

Positional Encoding

Recall dot product self attention

Positional Encoding

Recall dot product self attention

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

Positional Encoding

Recall dot product self attention

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

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

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

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

Positional Encoding

Permuting the inputs reorders Q, K, V

Positional Encoding

Permuting the inputs reorders Q, K, V

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

Positional Encoding

Permuting the inputs reorders Q, K, V

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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/transformer **does not** understand order.
Equivariant, order **does not** matter to the transformer.

Positional Encoding

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

Positional Encoding

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



Positional Encoding

This makes sense, in our party attention 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

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

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

Can we make the transformer care about 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 (function of time)

Positional Encoding

Question: What are some ways we can introduce ordering?

Answer 1: We can introduce forgetting (function of time)

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 (function of time)

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 (function of time)

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 (function of time)

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

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

Positional Encoding

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

$$\text{attn} \left(\begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(1) \\ \vdots \\ f_{\text{pos}}(T) \end{bmatrix}, \theta \right)$$

Positional Encoding

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

$$\text{attn} \left(\begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(1) \\ \vdots \\ f_{\text{pos}}(T) \end{bmatrix}, \theta \right)$$

Add ordering information to the inputs

Positional Encoding

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

$$\text{attn} \left(\begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(1) \\ \vdots \\ f_{\text{pos}}(T) \end{bmatrix}, \theta \right)$$

Add ordering information to the inputs

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

Positional Encoding

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

$$\text{attn} \left(\begin{bmatrix} x_1 \\ \vdots \\ x_T \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(1) \\ \vdots \\ f_{\text{pos}}(T) \end{bmatrix}, \theta \right)$$

Add ordering information to the inputs

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

$$\text{attn} \left(\begin{bmatrix} x_T \\ x_1 \\ \vdots \\ x_2 \end{bmatrix} + \begin{bmatrix} f_{\text{pos}}(T) \\ f_{\text{pos}}(1) \\ \vdots \\ f_{\text{pos}}(2) \end{bmatrix}, \theta \right)$$

Positional Encoding

What is f_{pos} ?

Positional Encoding

What is f_{pos} ?

Easy solution is just some random parameters

Positional Encoding

What is f_{pos} ?

Easy solution is just some random parameters

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

Positional Encoding

What is f_{pos} ?

Easy solution is just some random parameters

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

We call this a **positional encoding**

Positional Encoding

What is f_{pos} ?

Easy solution is just some random parameters

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

We call this a **positional encoding**

In torch, this is called `nn.Embedding`

Positional Encoding

What is f_{pos} ?

Easy solution is just some random parameters

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

We call this a **positional encoding**

In torch, this is called `nn.Embedding`

Now, let us rewrite the transformer with the positional encoding

Positional Encoding

So, our final transformer is

```
class Transformer(nn.Module):
    def __init__(self):
        self.f_pos = nn.Embedding(d_x, T)
        self.layer1 = TransformerLayer()
        self.layer2 = TransformerLayer()

    def forward(self, x):
        x = x + f_pos(torch.arange(x.shape[0]))
        x = self.layer1(x)
        x = self.layer2(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 a transformer to many types of problems

Applications

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

We can apply a transformer to many types of problems

But today we will focus on text and images

Applications

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

We can apply a transformer to many types of problems

But today we will focus on text and images

Let us see how to create inputs for our transformers

Text Transformers

Text Transformers

Consider a dataset of sentences

$$\begin{bmatrix} \text{John likes movies} \\ \text{Mary likes movies} \\ \text{I like dogs} \end{bmatrix}$$

Text Transformers

Consider a dataset of sentences

$$\begin{bmatrix} \text{John likes movies} \\ \text{Mary likes movies} \\ \text{I like dogs} \end{bmatrix}$$

This is a vector of sentences, but a transformer input is $\mathbb{R}^{T \times d_x}$

Text Transformers

Consider a dataset of sentences

$$\begin{bmatrix} \text{John likes movies} \\ \text{Mary likes movies} \\ \text{I like dogs} \end{bmatrix}$$

This is a vector of sentences, but a transformer input is $\mathbb{R}^{T \times d_x}$

What if we represent a sentence like this

$$\underbrace{\begin{bmatrix} \text{John} \\ \text{likes} \\ \text{movies} \end{bmatrix}}_{d_x} \Big\} T$$

Text Transformers

$$\underbrace{\begin{bmatrix} \text{John} \\ \text{likes} \\ \text{movies} \end{bmatrix}}_{d_x} \Big\} T$$

Text Transformers

$$\underbrace{\begin{bmatrix} \text{John} \\ \text{likes} \\ \text{movies} \end{bmatrix}}_{d_x} \Big\} T$$

But these are words, the transformer needs the input to be numbers

Text Transformers

$$\underbrace{\begin{bmatrix} \text{John} \\ \text{likes} \\ \text{movies} \end{bmatrix}}_{d_x} \Big\} T$$

But these are words, the transformer needs the input to be numbers

What if we create a vector to represent each word?

Text Transformers

$$\underbrace{\begin{bmatrix} \text{John} \\ \text{likes} \\ \text{movies} \end{bmatrix}}_{d_x} \Big\} T$$

But these are words, the transformer needs the input to be numbers

What if we create a vector to represent each word?

$$\begin{bmatrix} \text{John} \\ \text{likes} \\ \text{movies} \end{bmatrix} = \underbrace{\begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}}_{d_x} \Big\} T$$

Text Transformers

Step 1: Find all unique words in
the dataset

Text Transformers

Step 1: Find all unique words in the dataset

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

Text Transformers

Step 1: Find all unique words in the dataset

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

= [John likes movies Mary I dogs]

Step 2: Create a vector representation for each unique word

Text Transformers

Step 1: Find all unique words in the dataset

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

Step 2: Create a vector representation for each unique word

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

Text Transformers

Step 1: Find all unique words in the dataset

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

Step 2: Create a vector representation for each unique word

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

Step 3: Replace words with vector representations

Text Transformers

Example: Convert the sentence to vector representations

Text Transformers

Example: Convert the sentence to vector representations

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

Text Transformers

Example: Convert the sentence to vector representations

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

John likes movies =

Text Transformers

Example: Convert the sentence to vector representations

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

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

Text Transformers

Example: Convert the sentence to vector representations

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

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

$$\text{Mary likes John} =$$

Text Transformers

Example: Convert the sentence to vector representations

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

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

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

Text Transformers

Example: Convert the sentence to vector representations

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

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

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

Now, let us write some pseudocode

Text Transformers

```
unique_words = set(sentence.split(" ") for sentence in xs)
word_index = {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 = word_index[word]
        representation = embeddings[word_index]
        xs.append(representation)

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

Text Transformers

Now, feed our dataset to the transformer

Text Transformers

Now, feed our dataset to the transformer

```
model = Transformer()
for sentence_representation 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

$$X \in [0, 1]^{3 \times 16 \times 16}$$

Image Transformers

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

$$X \in [0, 1]^{3 \times 16 \times 16}$$

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

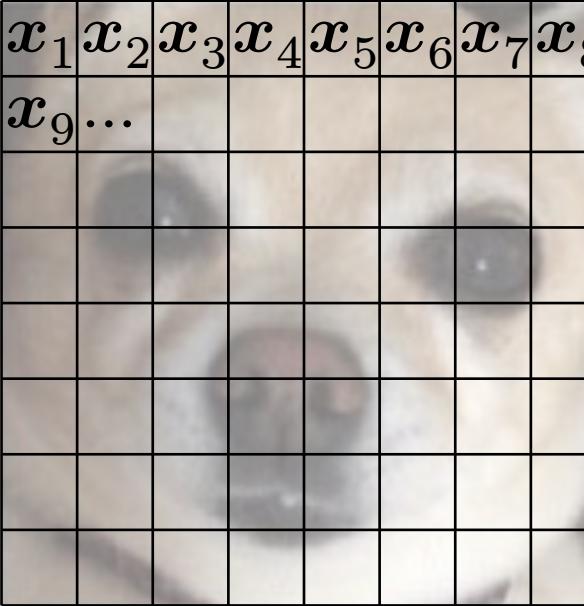
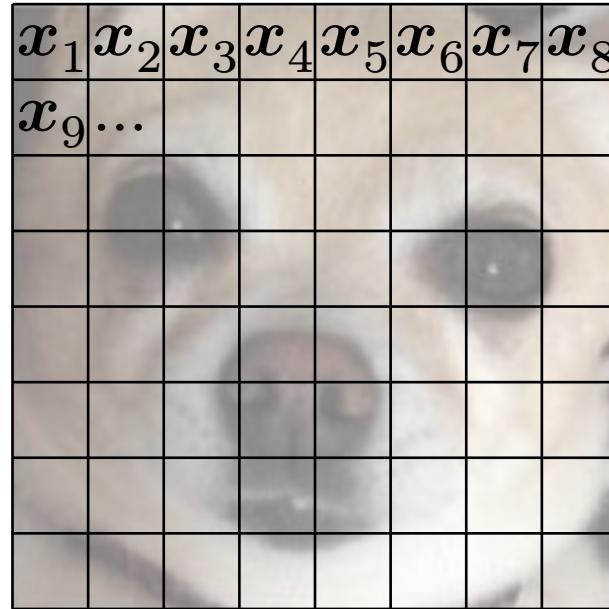


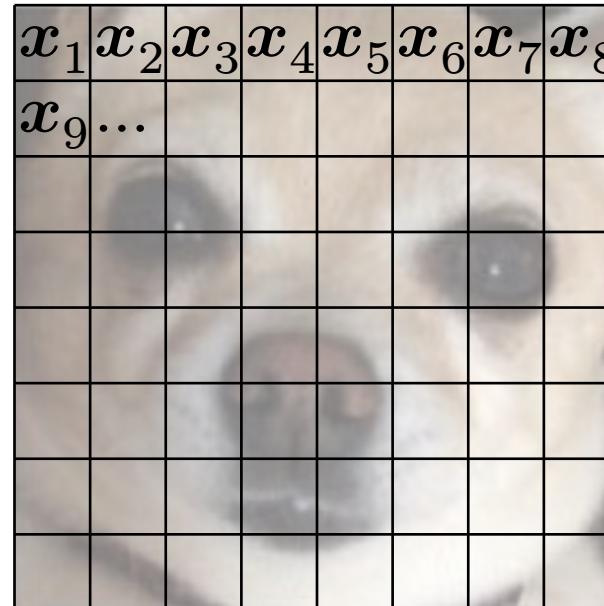
Image Transformers

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



Then, feed a sequence of patches
to the transformer

Image Transformers



Then, feed a sequence of patches
to the transformer

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

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)
```

Generative Pre-Training

Generative Pre-Training

Question: How do we train transformers?

Generative Pre-Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

Generative Pre-Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

- Classification loss

Generative Pre-Training

Question: How do we train transformers?

Answer: Can train just like other neural networks

- Classification loss
- Regression loss

Generative Pre-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

Generative Pre-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

Transformer scale very well – add more data, model becomes stronger

Generative Pre-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

Transformer scale very well – add more data, model becomes stronger

Today, dataset size limits transformers

Generative Pre-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

Transformer scale very well – add more data, model becomes stronger

Today, dataset size limits transformers

There are not enough graduate students to label training data!

Generative Pre-Training

Question: How can we train transformers with finite students/datasets?

Generative Pre-Training

Question: How can we train transformers with finite students/datasets?

Answer: We can use **unsupervised learning**

Generative Pre-Training

Question: How can we train transformers with finite students/datasets?

Answer: We can use **unsupervised learning**

The internet contains billions of unlabeled sentences and images

Generative Pre-Training

Question: How can we train transformers with finite students/datasets?

Answer: We can use **unsupervised learning**

The internet contains billions of unlabeled sentences and images

Mask/hide part of the input, train the model to predict the missing part

Generative Pre-Training

Question: How can we train transformers with finite students/datasets?

Answer: We can use **unsupervised learning**

The internet contains billions of unlabeled sentences and images

Mask/hide part of the input, train the model to predict the missing part

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

Generative Pre-Training

Question: How can we train transformers with finite students/datasets?

Answer: We can use **unsupervised learning**

The internet contains billions of unlabeled sentences and images

Mask/hide part of the input, train the model to predict the missing part

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

We **generate** the missing data

Generative Pre-Training

Question: How can we train transformers with finite students/datasets?

Answer: We can use **unsupervised learning**

The internet contains billions of unlabeled sentences and images

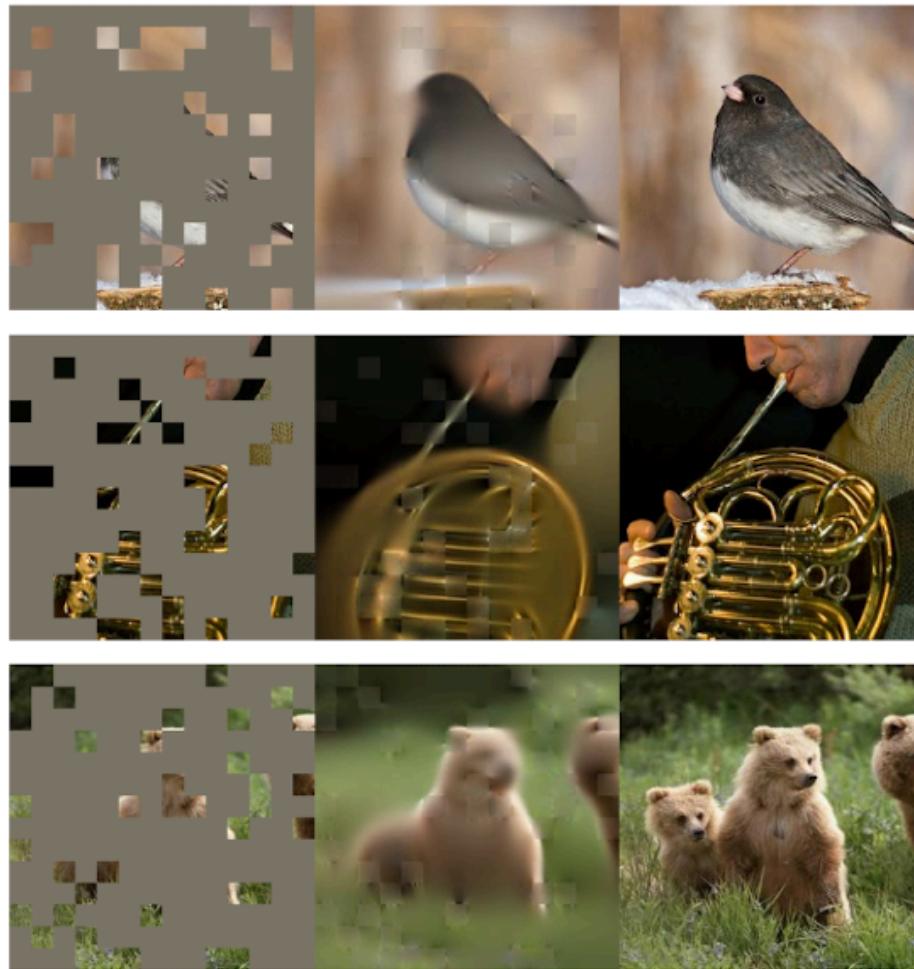
Mask/hide part of the input, train the model to predict the missing part

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

We **generate** the missing data

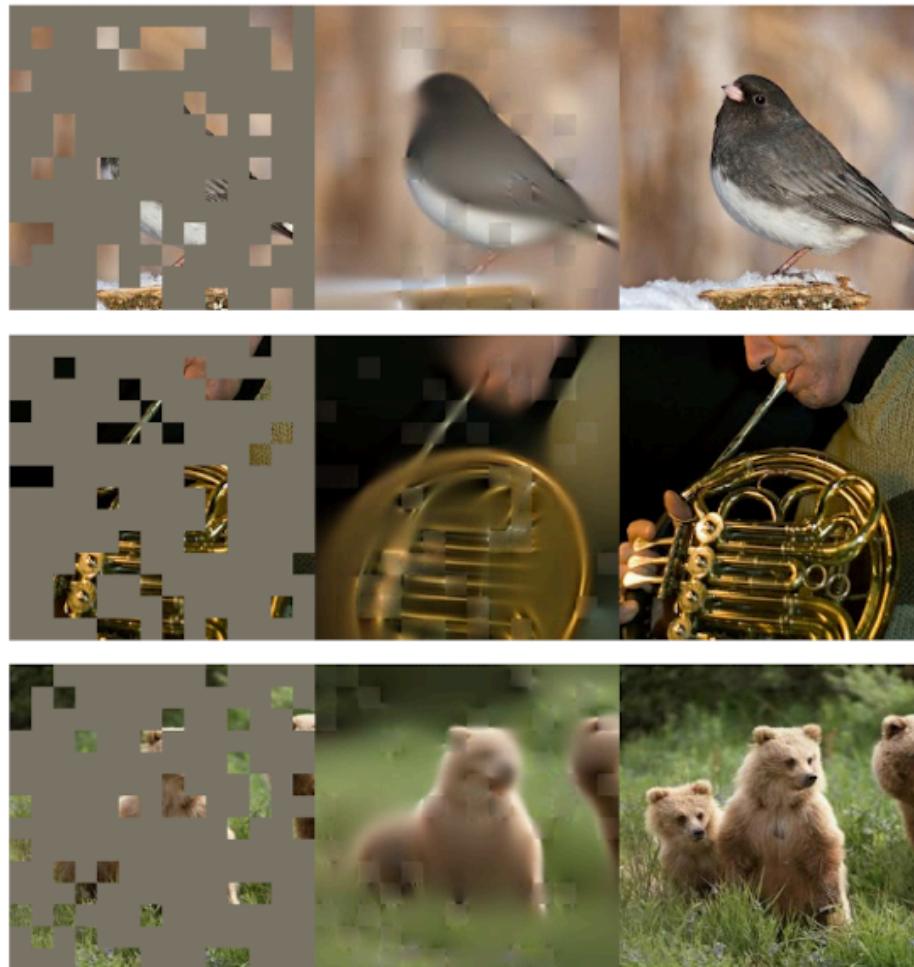
This method is **extremely** powerful

Generative Pre-Training



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

Generative Pre-Training



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

Generative Pre-Training

In GPT, we optimize the following objective

Generative Pre-Training

In GPT, we optimize the following objective

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

Generative Pre-Training

In GPT, we optimize the following objective

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

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

Generative Pre-Training

In GPT, we optimize the following objective

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

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

$$\boldsymbol{x}_T = \text{"movies"}$$

Generative Pre-Training

In GPT, we optimize the following objective

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

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

$$\mathbf{x}_T = \text{"movies"}$$

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

Generative Pre-Training

In GPT, we optimize the following objective

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

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

$$\mathbf{x}_T = \text{"movies"}$$

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

Update θ so that $P(\text{movies} \mid \text{John, likes; } \theta) = 0.6$

Generative Pre-Training

What about for images?

Generative Pre-Training

What about for images?

Use the same objective

Generative Pre-Training

What about for images?

Use the same objective

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

Generative Pre-Training

What about for images?

Use the same objective

$$\begin{aligned} \arg \min_{\theta} \mathcal{L} \left(\begin{bmatrix} \boldsymbol{x}_1 \\ \vdots \\ \boldsymbol{x}_T \end{bmatrix}, \boldsymbol{\theta} \right) &= \arg \min_{\theta} [-\log P(\boldsymbol{x}_T \mid \boldsymbol{x}_1, \dots, \boldsymbol{x}_{T-1}; \boldsymbol{\theta})] \\ &= \arg \min_{\theta} [-\log P(\text{pixel_3} \mid \text{pixel_1}, \text{pixel_2}; \boldsymbol{\theta})] \end{aligned}$$

Generative Pre-Training

What about for images?

Use the same objective

$$\begin{aligned} \arg \min_{\theta} \mathcal{L}\left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_T \end{bmatrix}, \boldsymbol{\theta} \right) &= \arg \min_{\theta} [-\log P(\mathbf{x}_T \mid \mathbf{x}_1, \dots, \mathbf{x}_{T-1}; \boldsymbol{\theta})] \\ &= \arg \min_{\theta} [-\log P(\text{pixel_3} \mid \text{pixel_1}, \text{pixel_2}; \boldsymbol{\theta})] \end{aligned}$$

Minimizing the square error is equivalent to minimizing the (Gaussian) negative log probability

$$= \arg \min_{\theta} \left(f\left(\begin{bmatrix} \text{pixel_1} \\ \text{pixel_2} \end{bmatrix}, \boldsymbol{\theta} \right) - \text{pixel_3} \right)^2$$

Generative Pre-Training

Why does this work so well?

Generative Pre-Training

Why does this work so well?

Let us see what the models learn with a GPT loss

Generative Pre-Training

Why does this work so well?

Let us see what the models learn with a GPT loss

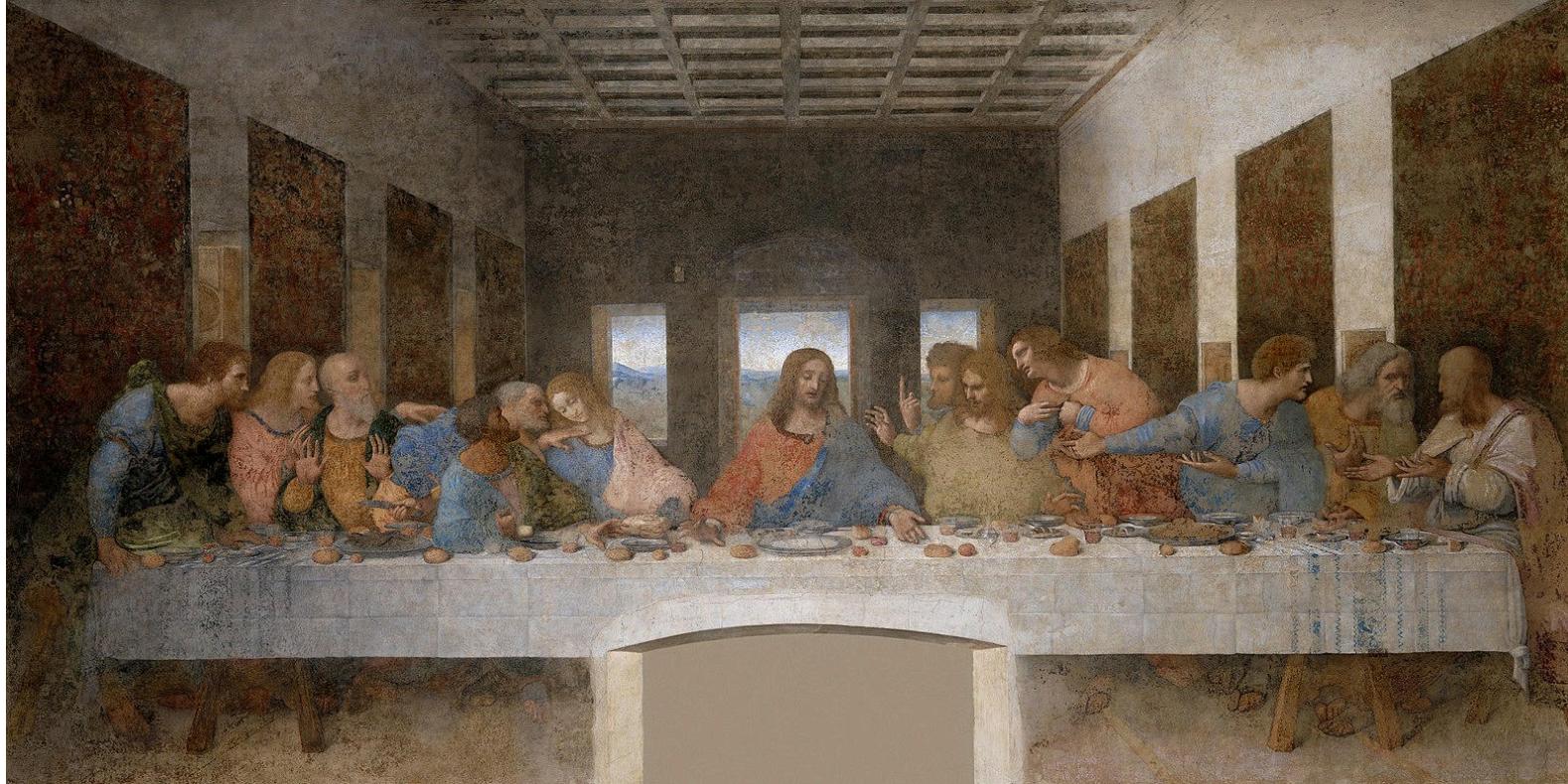
Let us start with the image transformer

Generative Pre-Training

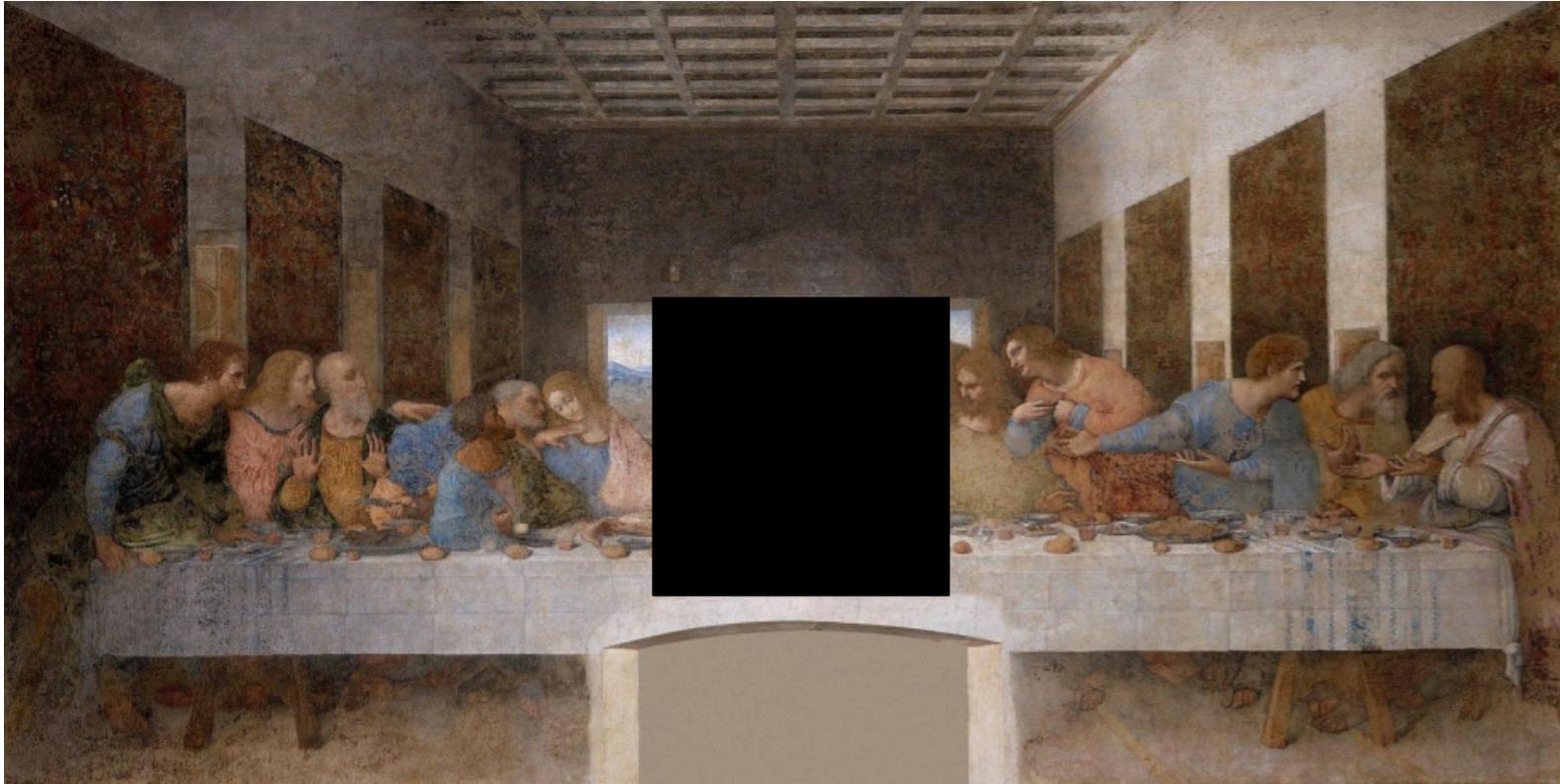
Anyone familiar with Da Vinci's painting *The Last Supper*?

Generative Pre-Training

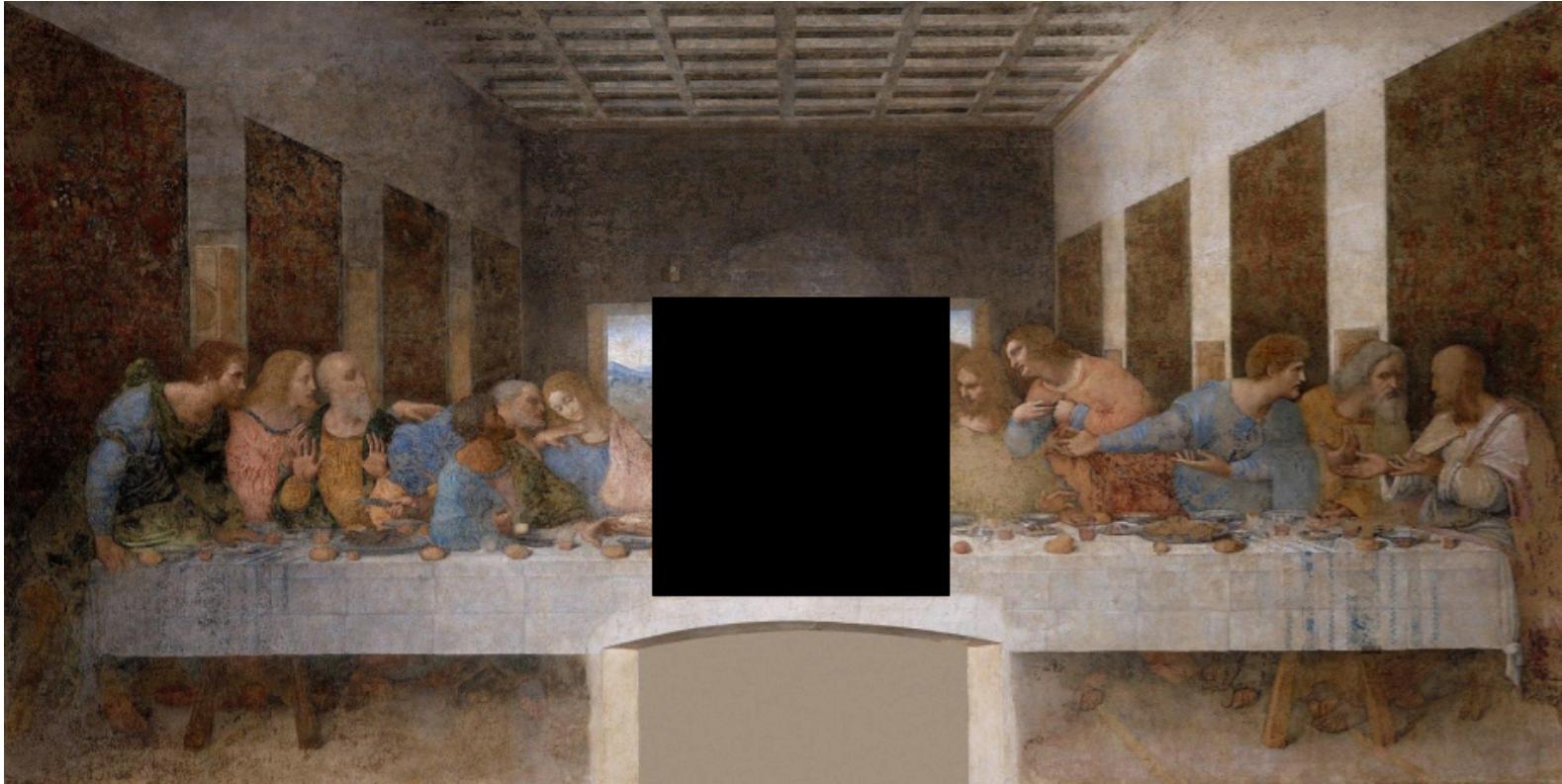
Anyone familiar with Da Vinci's painting *The Last Supper*?



Generative Pre-Training



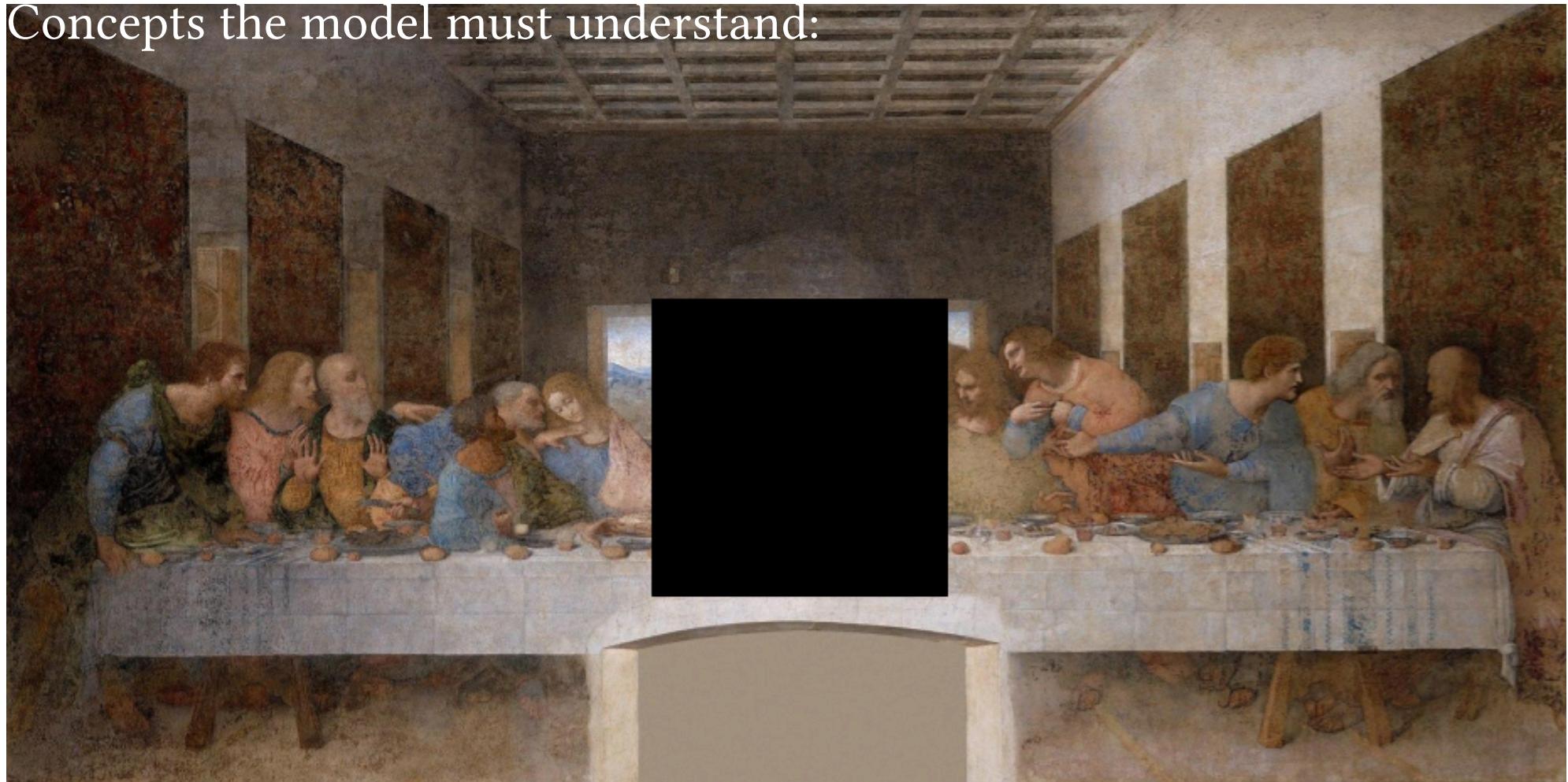
Generative Pre-Training



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

Generative Pre-Training

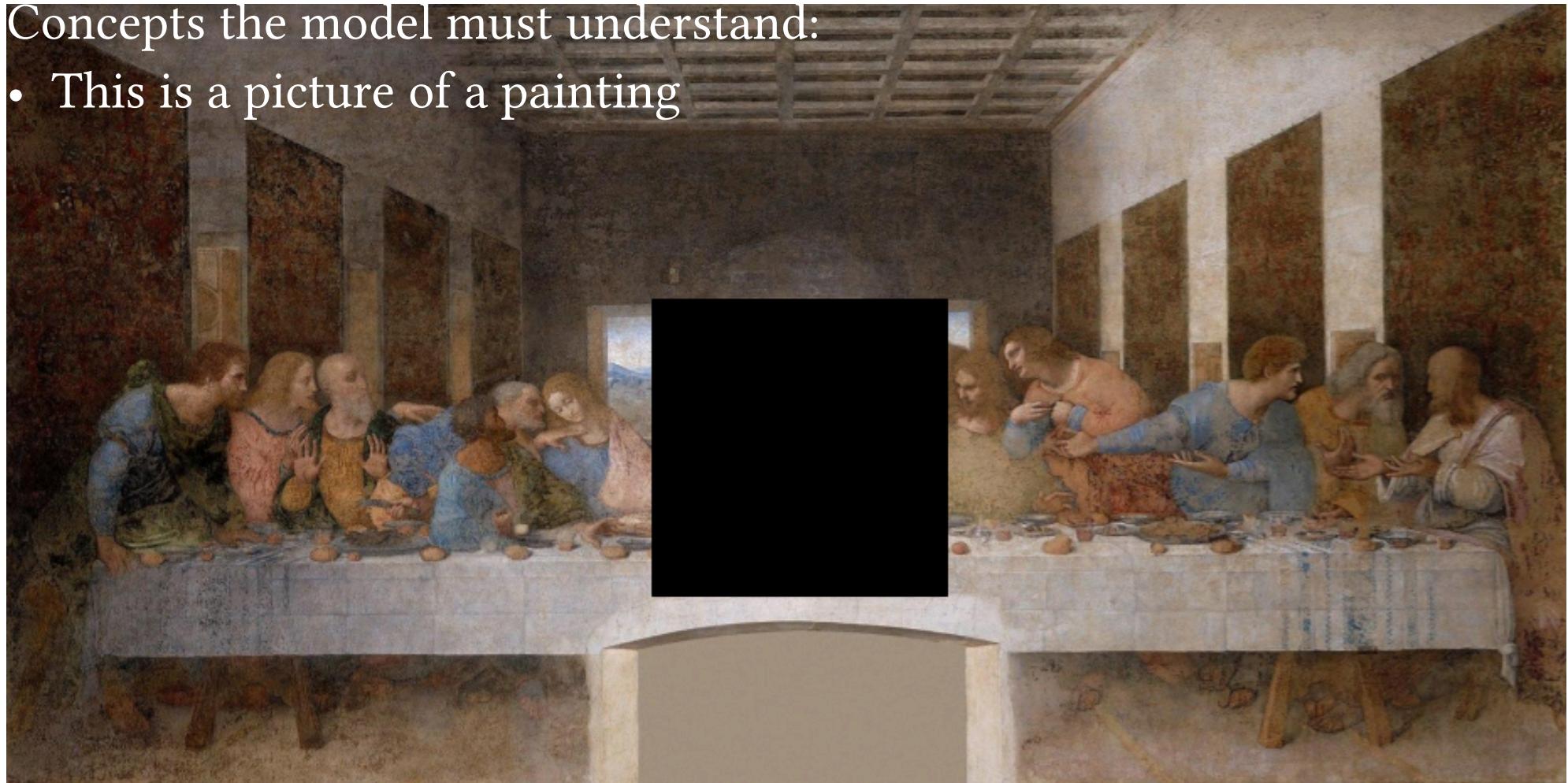
Concepts the model must understand:



Generative Pre-Training

Concepts the model must understand:

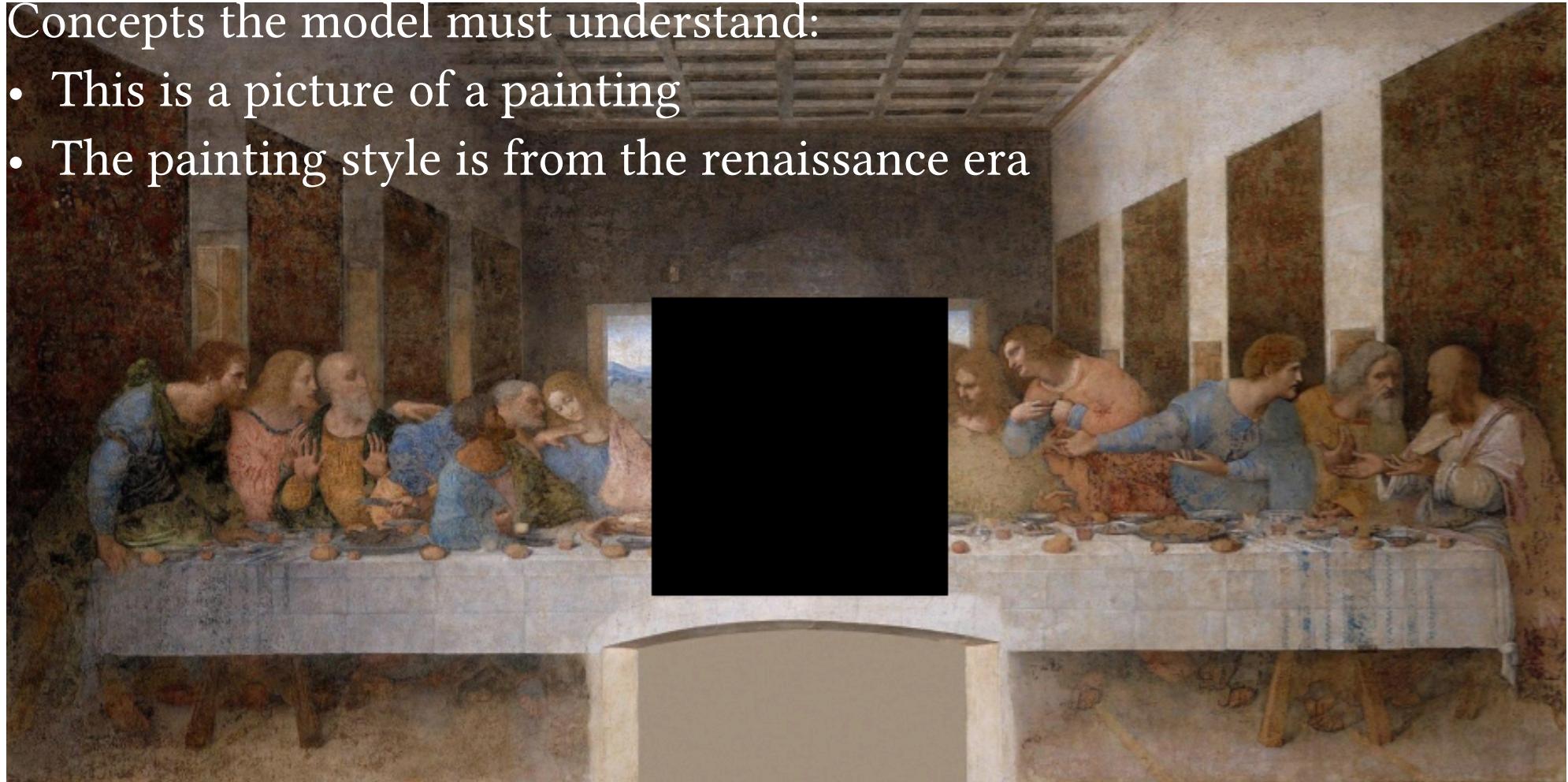
- This is a picture of a painting



Generative Pre-Training

Concepts the model must understand:

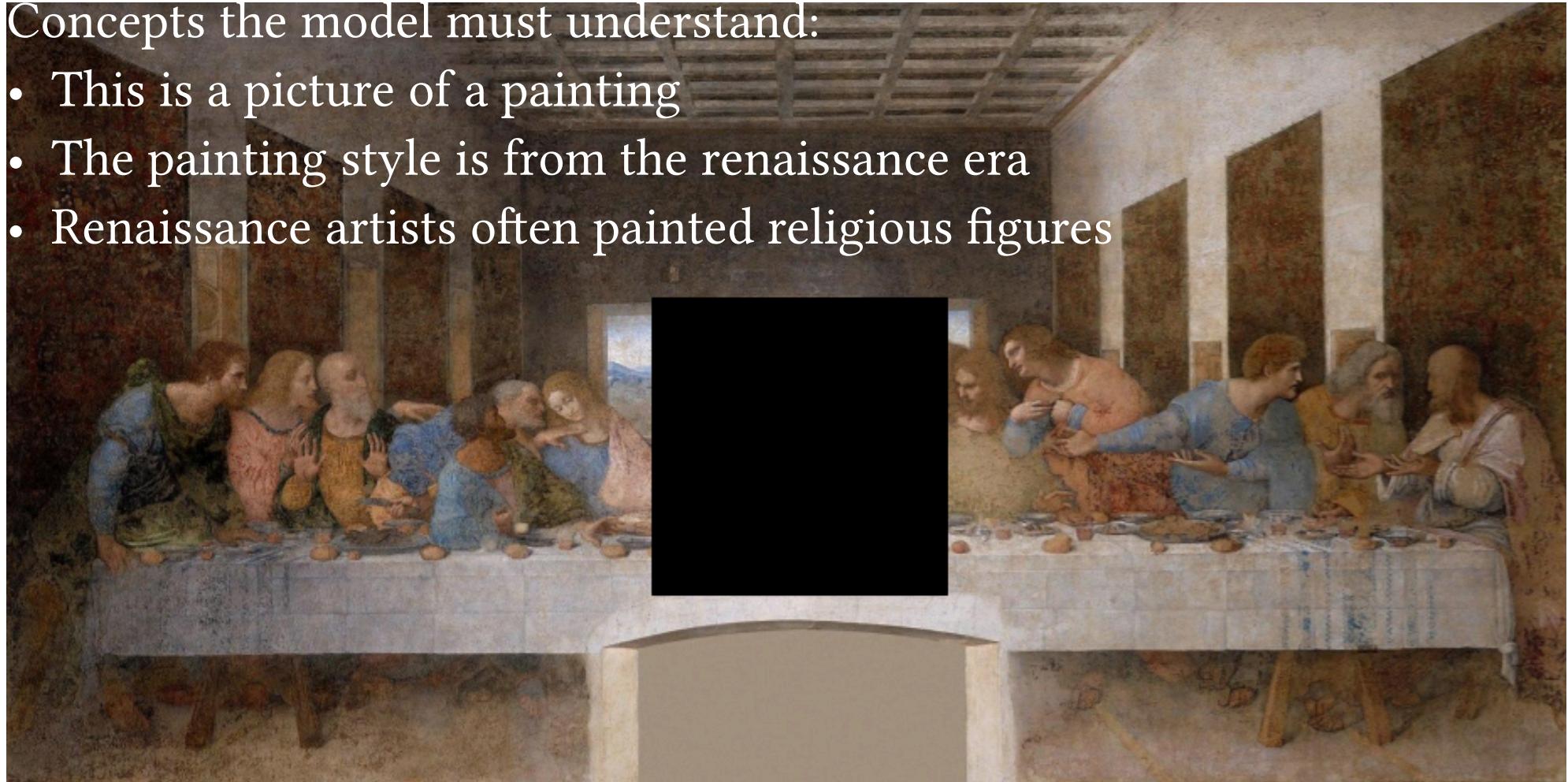
- This is a picture of a painting
- The painting style is from the renaissance era



Generative Pre-Training

Concepts the model must understand:

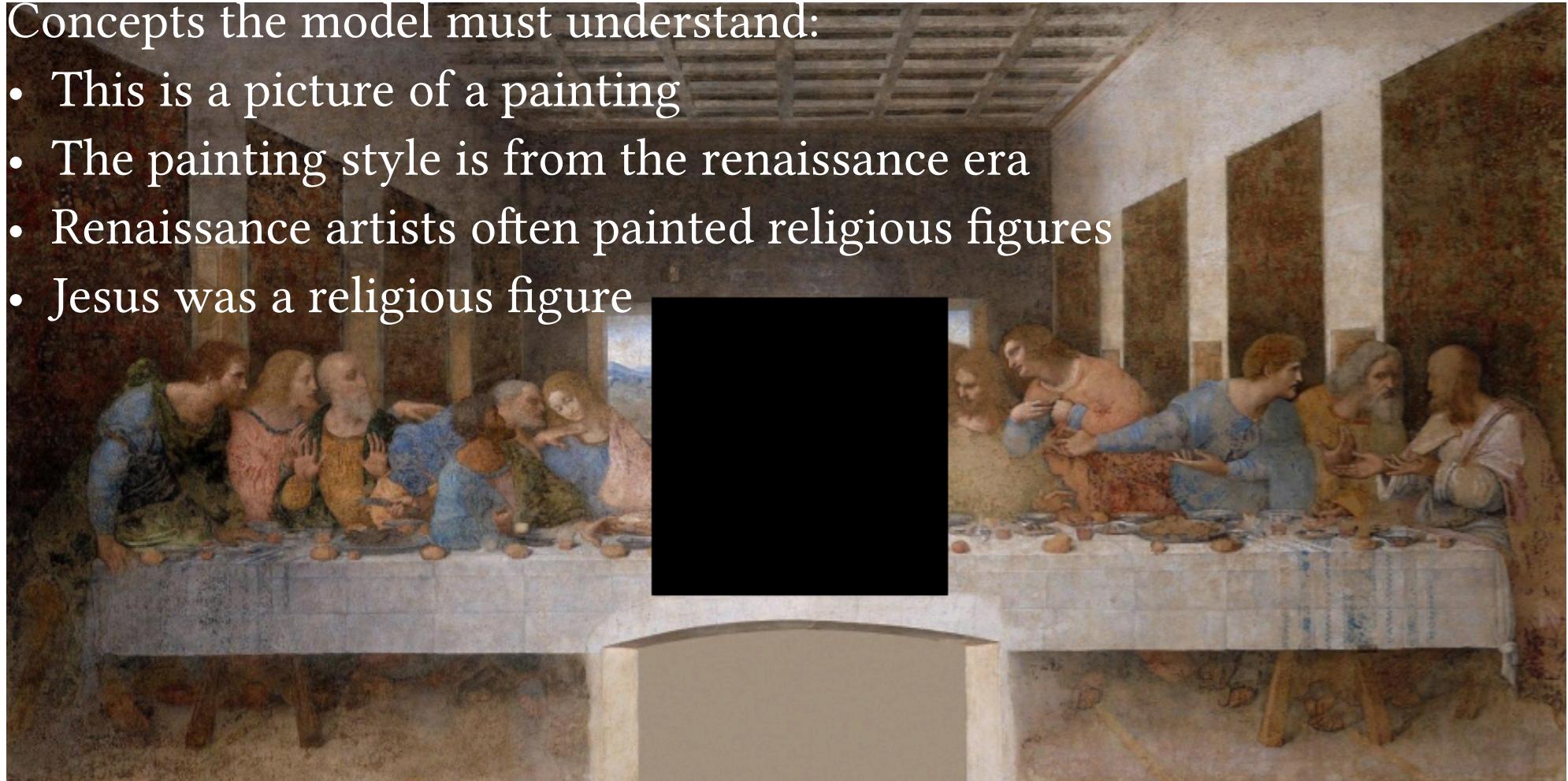
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures



Generative Pre-Training

Concepts the model must understand:

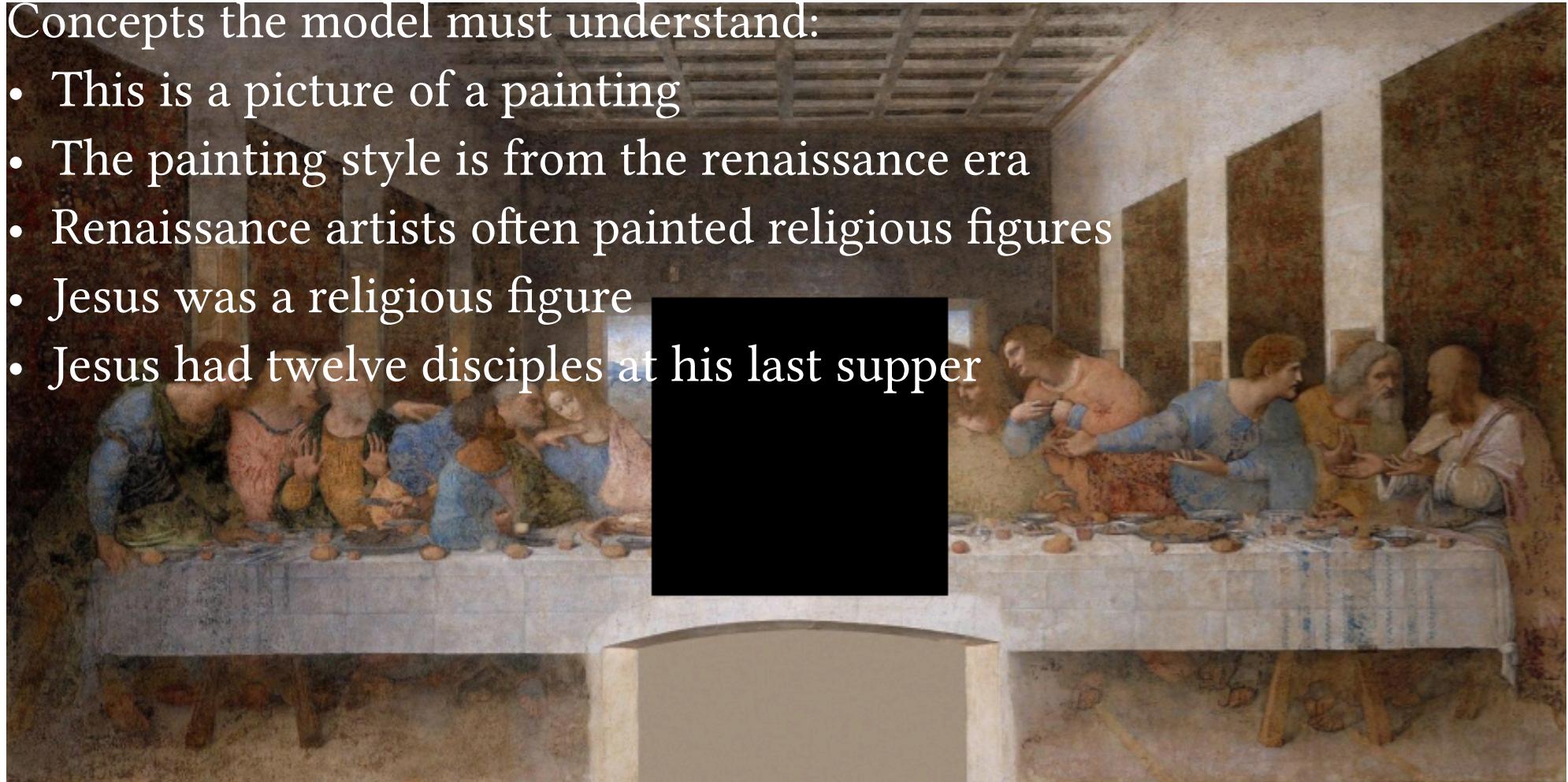
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures
- Jesus was a religious figure



Generative Pre-Training

Concepts the model must understand:

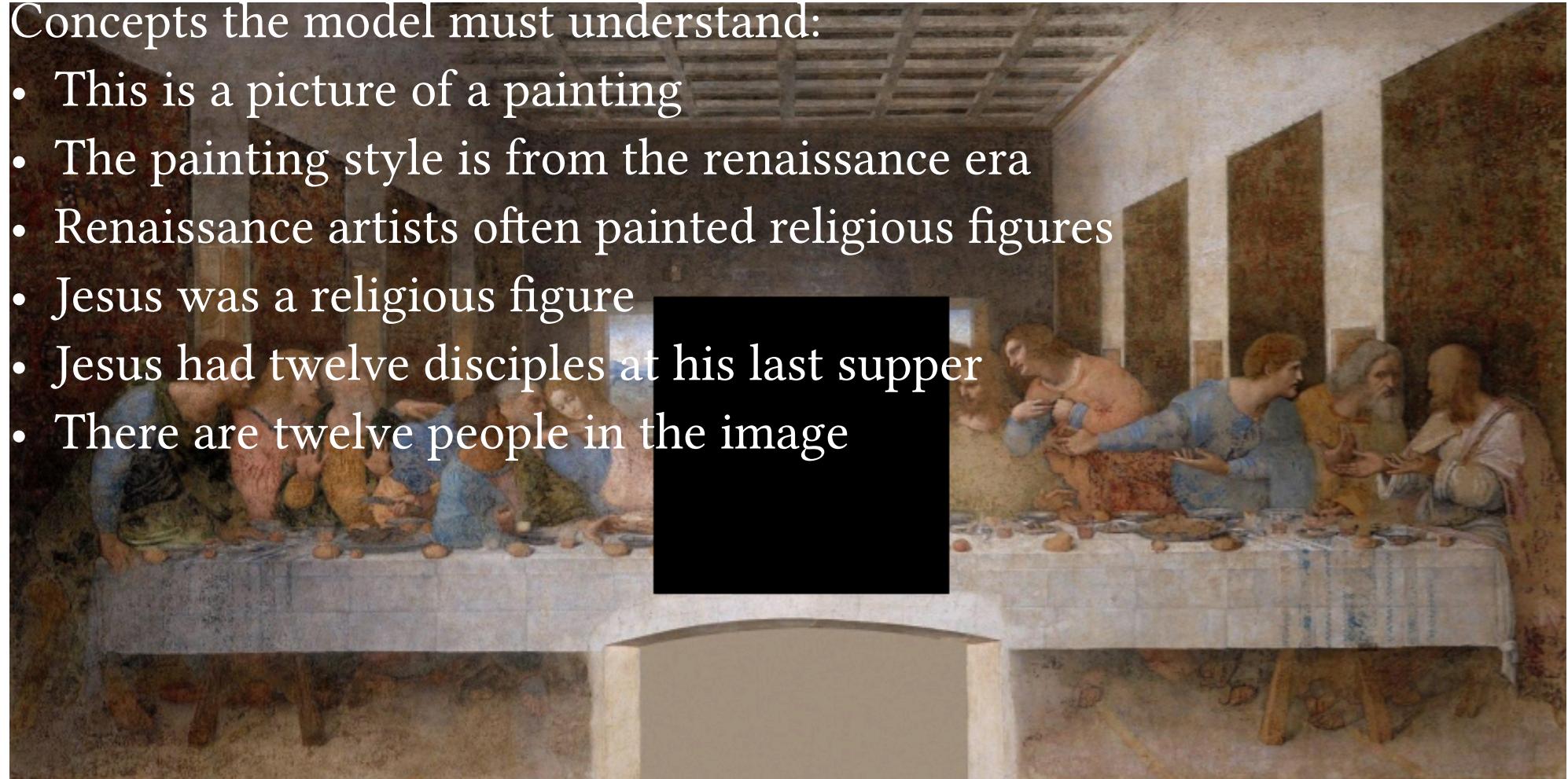
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures
- Jesus was a religious figure
- Jesus had twelve disciples at his last supper



Generative Pre-Training

Concepts the model must understand:

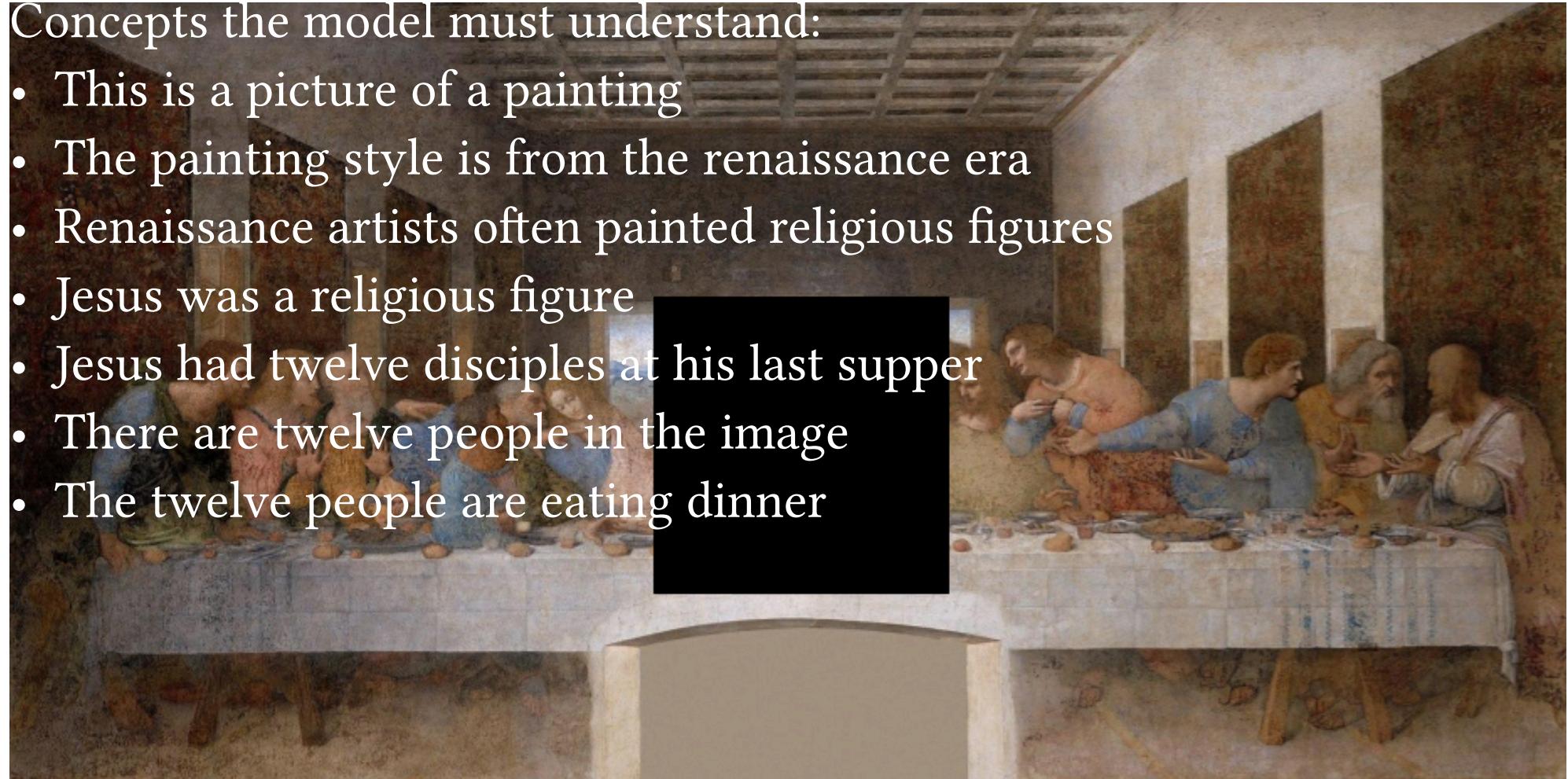
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures
- Jesus was a religious figure
- Jesus had twelve disciples at his last supper
- There are twelve people in the image



Generative Pre-Training

Concepts the model must understand:

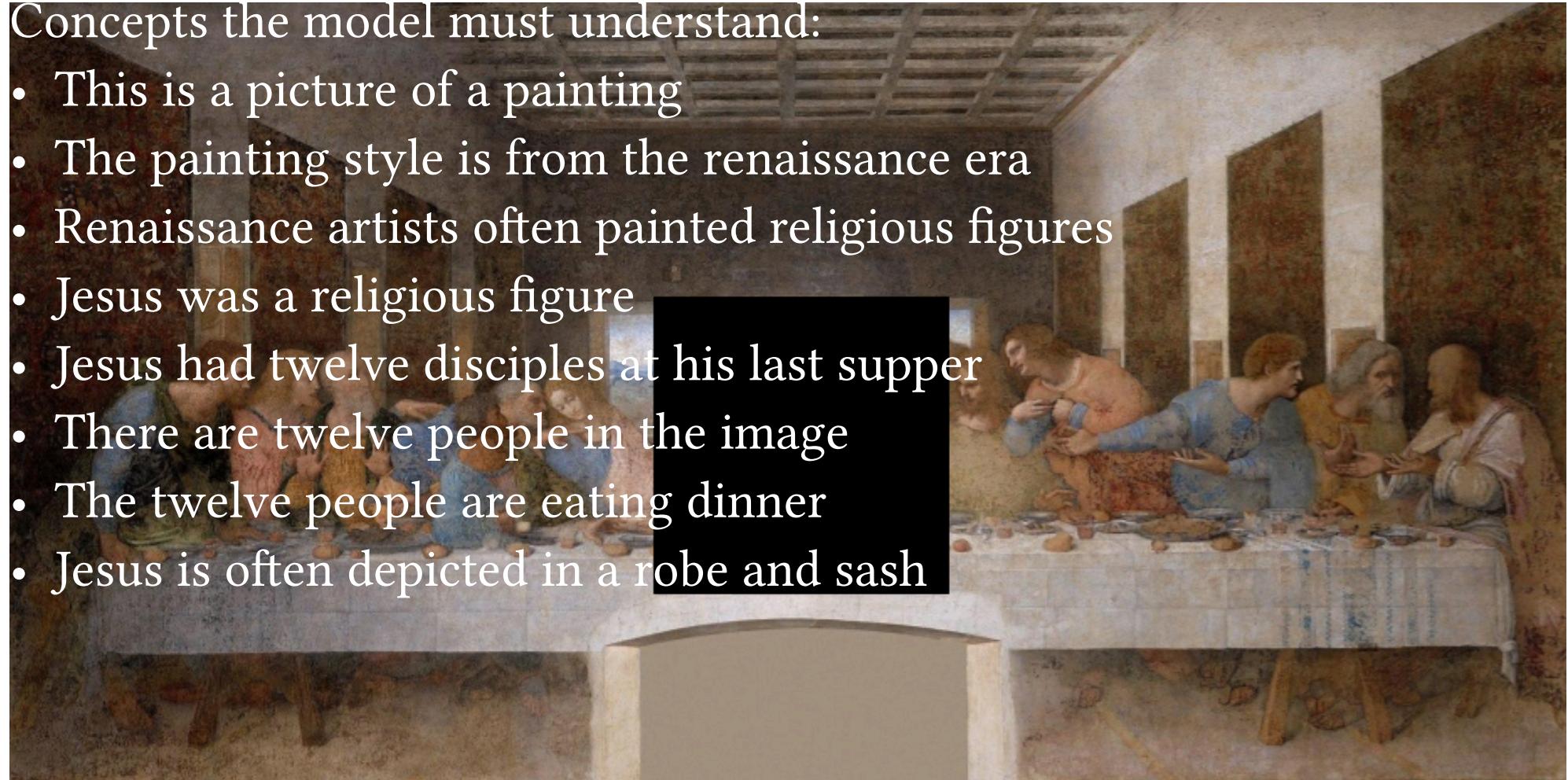
- This is a picture of a painting
- The painting style is from the renaissance era
- Renaissance artists often painted religious figures
- Jesus was a religious figure
- Jesus had twelve disciples at his last supper
- There are twelve people in the image
- The twelve people are eating dinner



Generative Pre-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



Generative Pre-Training

We train the model to fix the image

Generative Pre-Training

We train the model to fix the image

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

Generative Pre-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

Generative Pre-Training

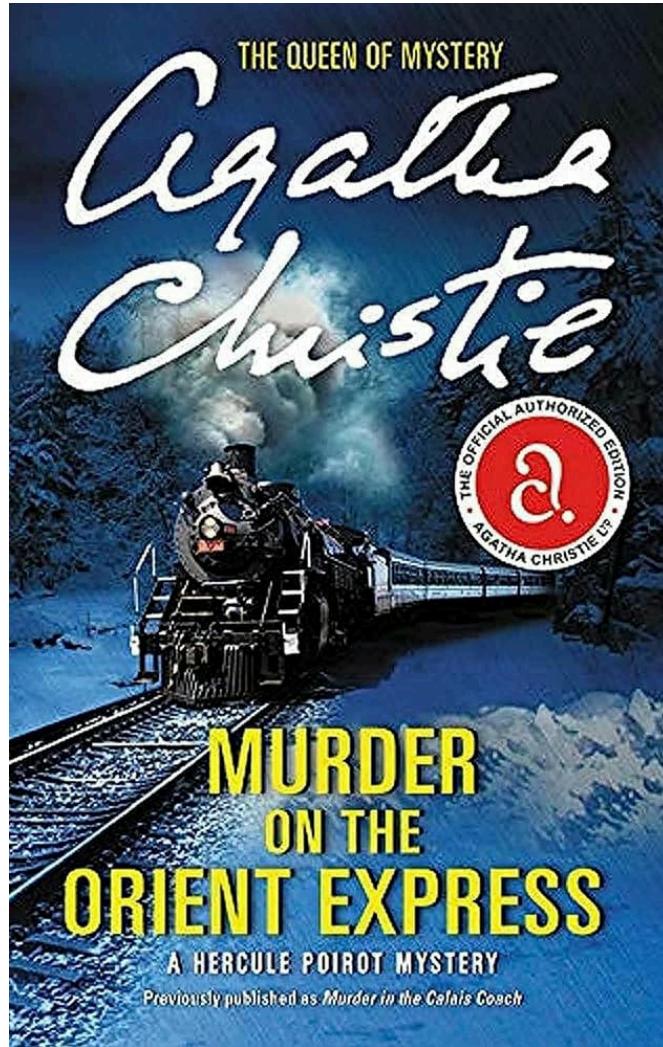
We train the model to fix the image

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

This is the power of generative pre-training

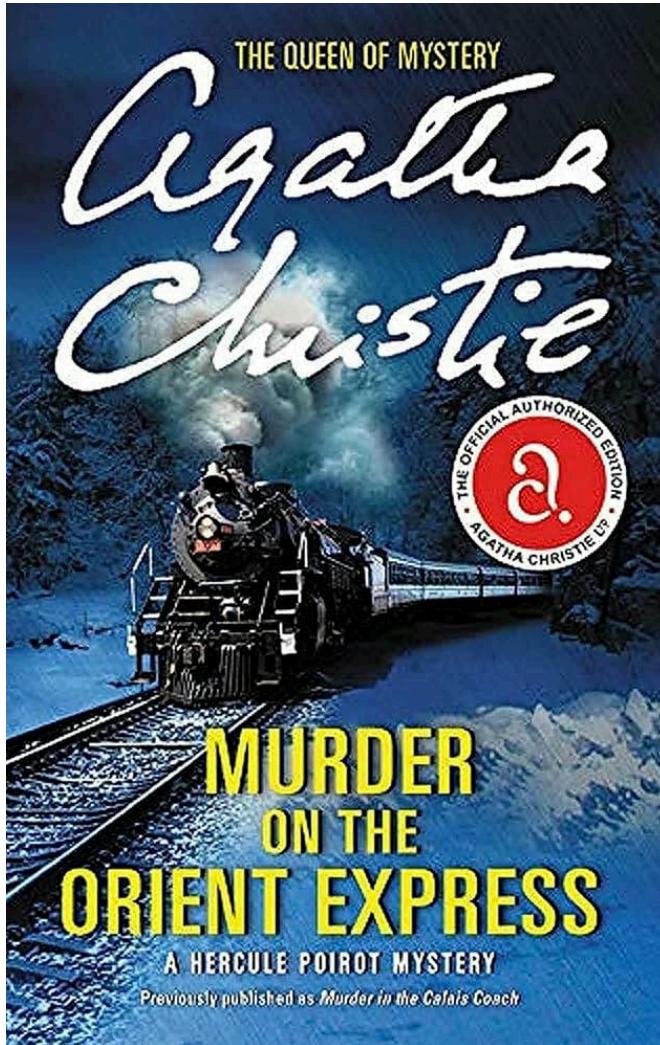
What about text transformers?

Generative Pre-Training



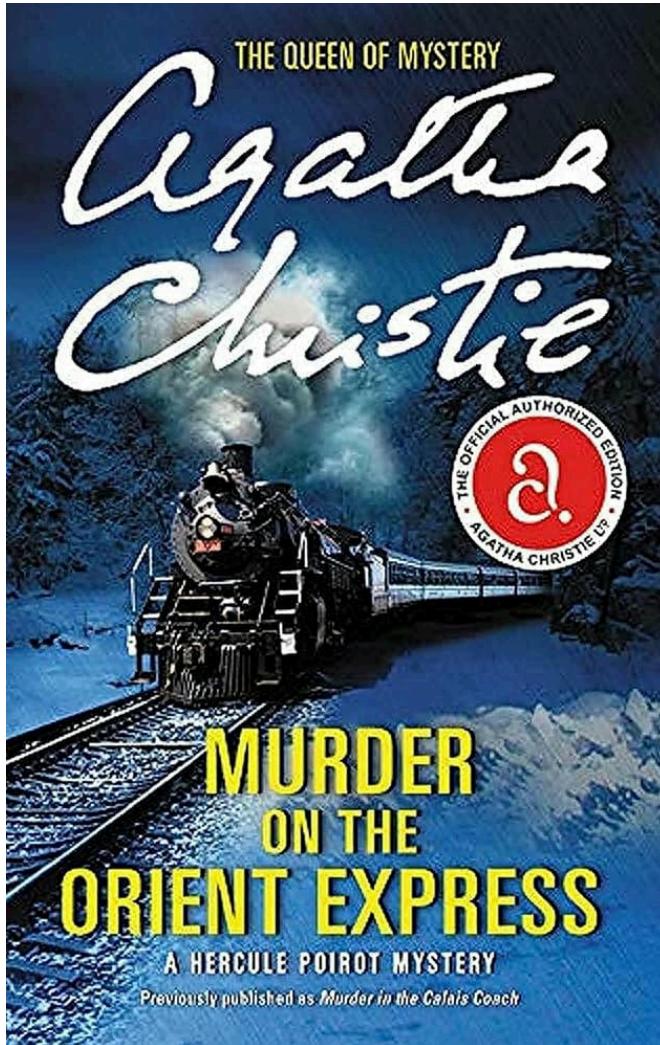
This is a mystery novel

Generative Pre-Training



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

Generative Pre-Training

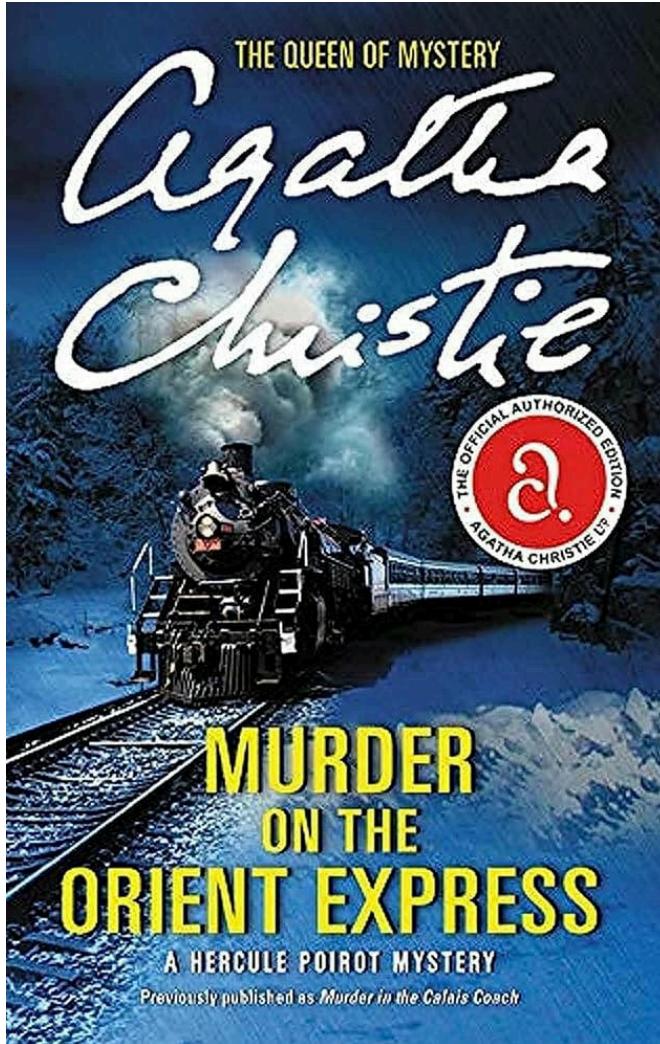


This is a mystery novel

Clues, intrigue, murder, etc

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

Generative Pre-Training



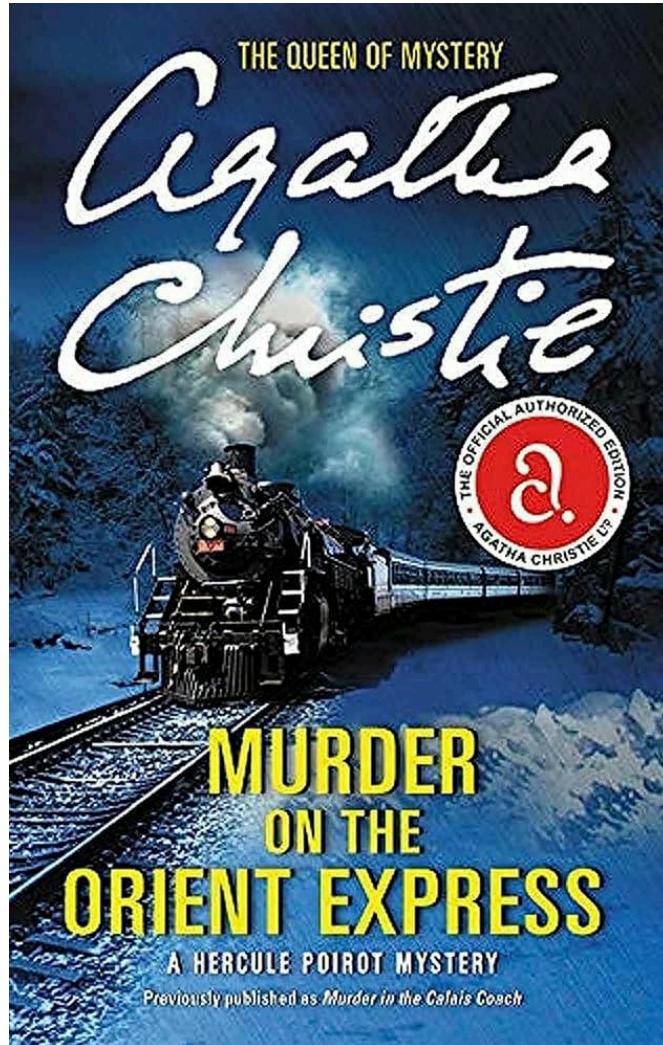
This is a mystery novel

Clues, intrigue, murder, etc

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

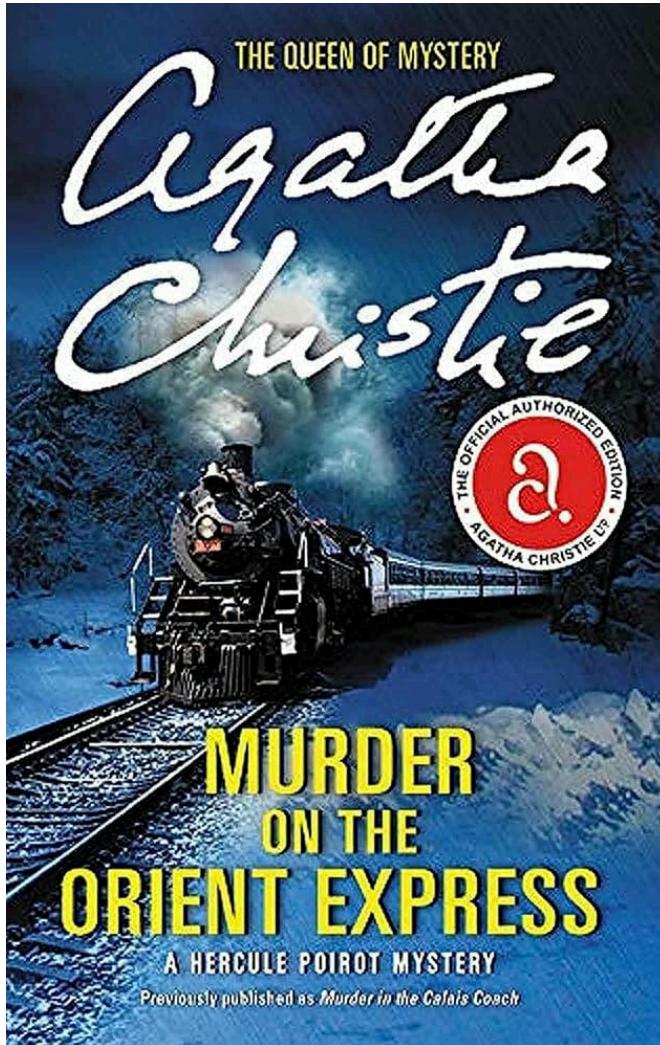
To complete the sentence, the model must understand:

Generative Pre-Training



To complete the sentence, the model must understand:

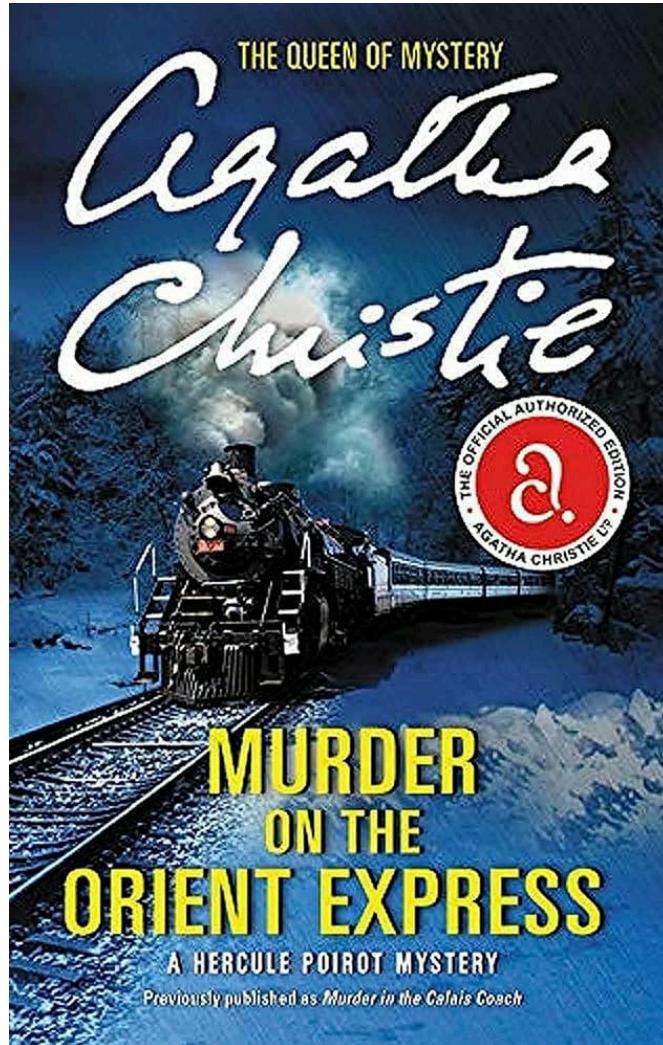
Generative Pre-Training



To complete the sentence, the model must understand:

- What a murder is

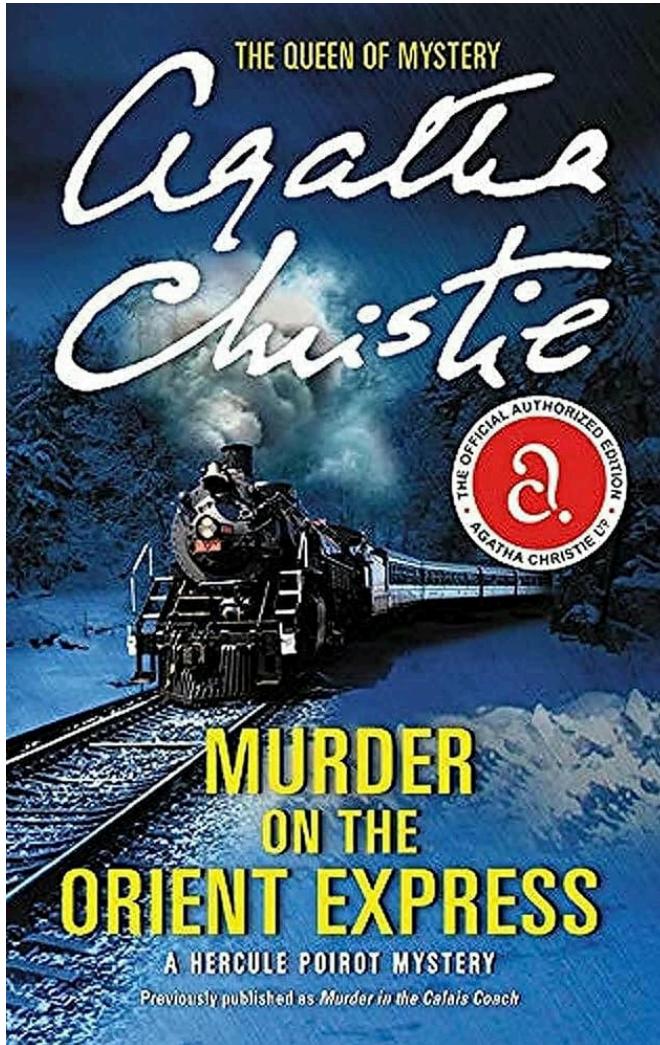
Generative Pre-Training



To complete the sentence, the model must understand:

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

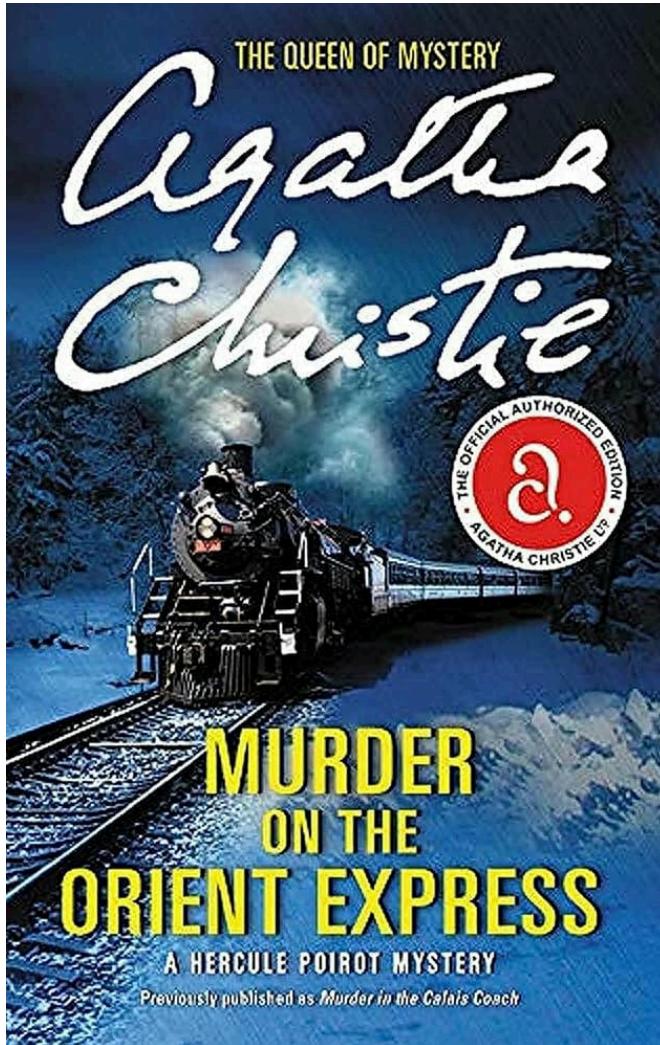
Generative Pre-Training



To complete the sentence, the model must understand:

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

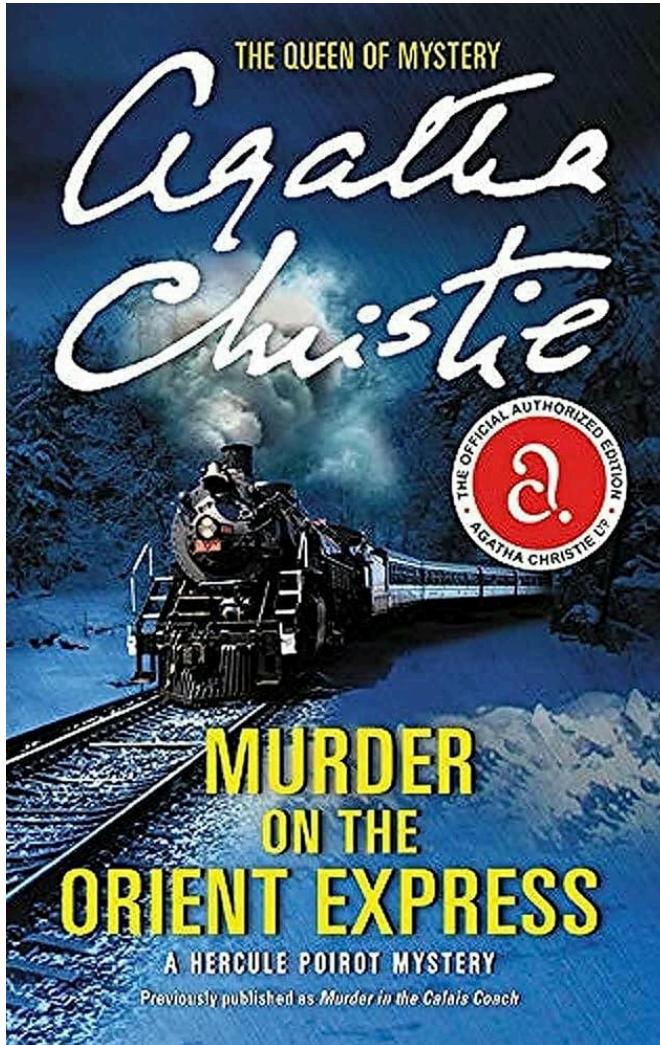
Generative Pre-Training



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

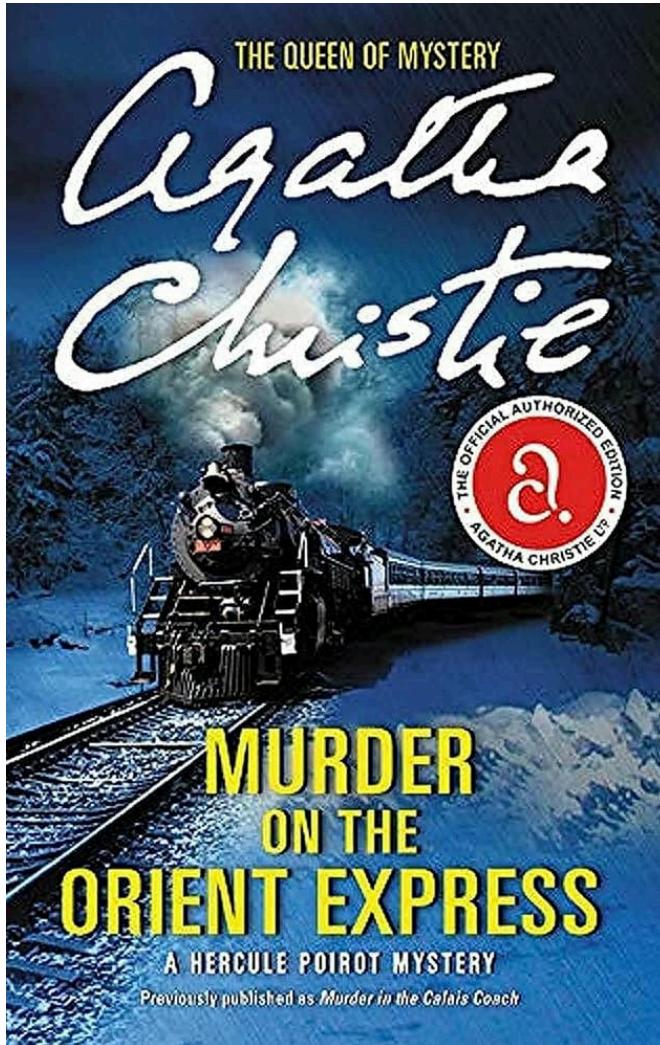
Generative Pre-Training



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

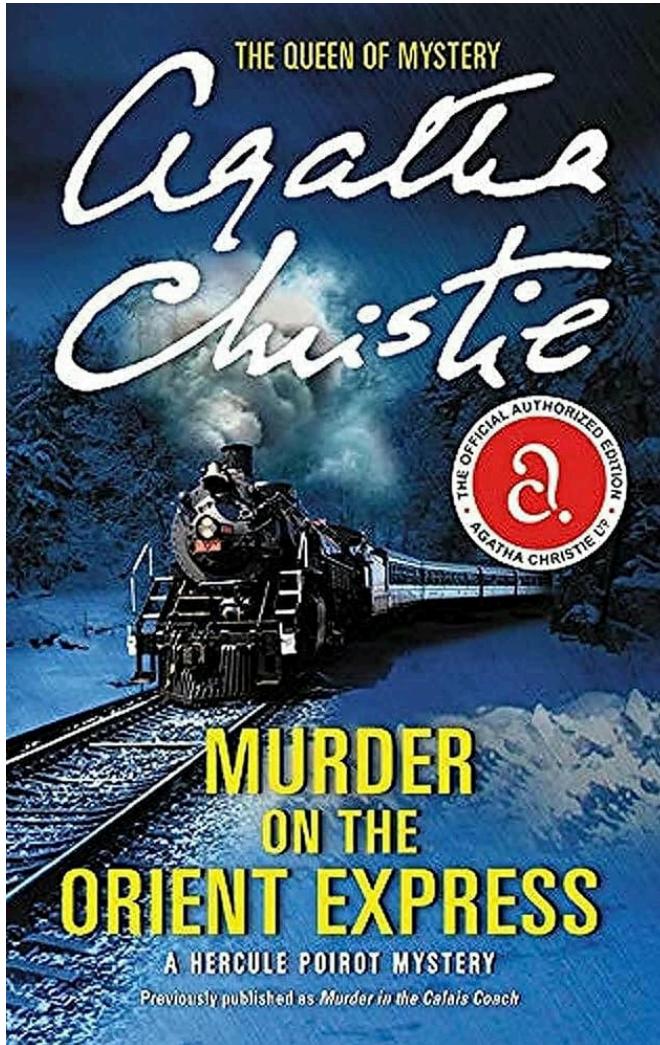
Generative Pre-Training



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

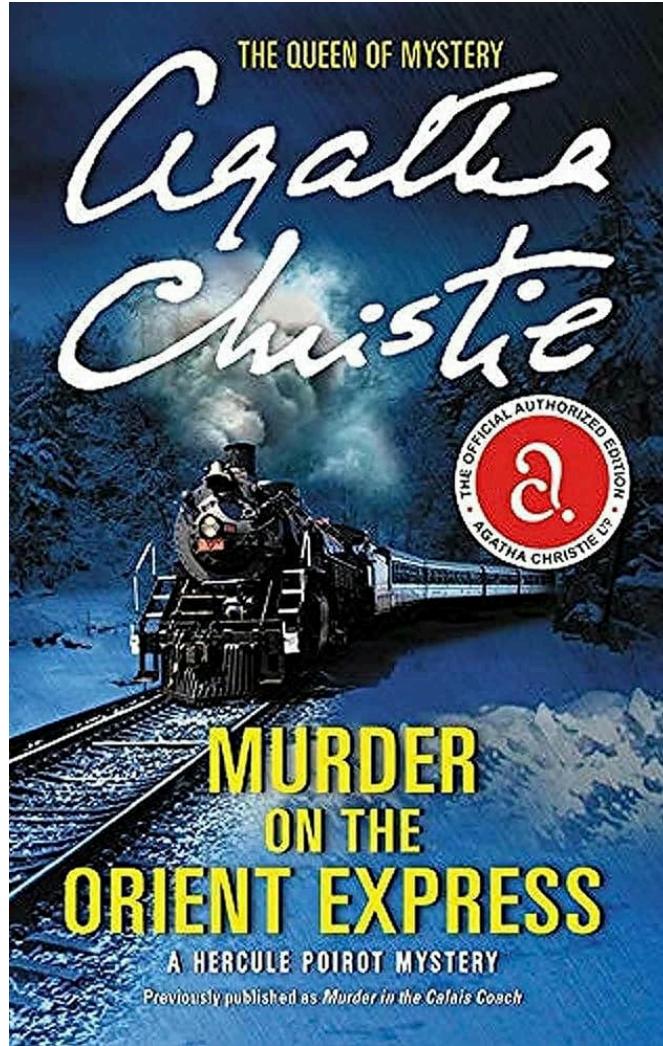
Generative Pre-Training



To complete the sentence, the model must understand:

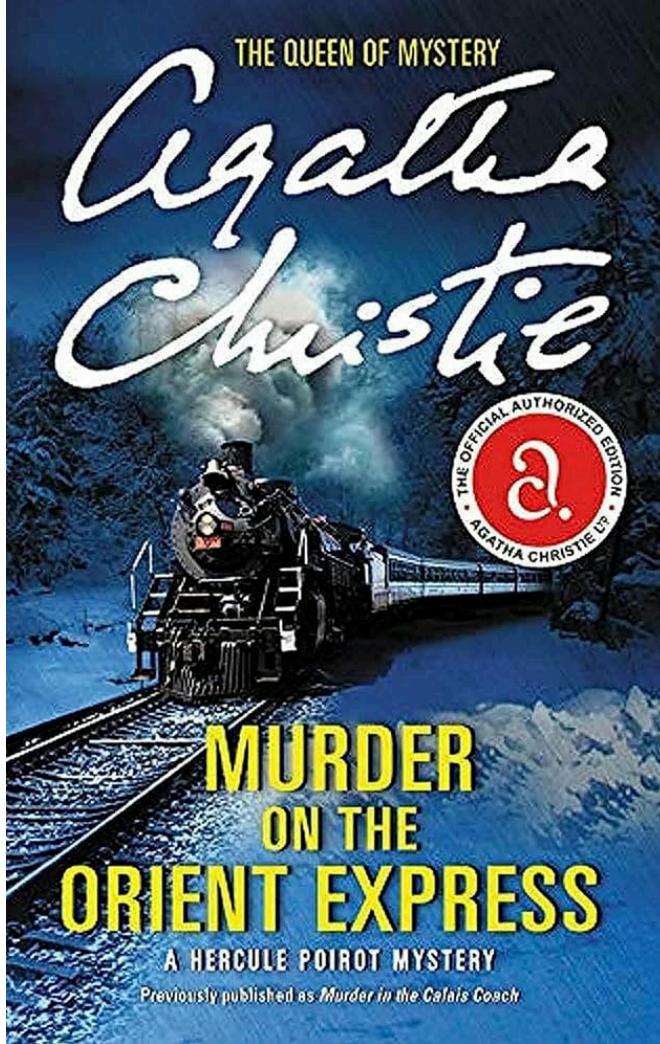
- What a murder is
- What it means to be alive
- Emotions like anger, jealousy, betrayal, love
- Personalities of each character
- Why a human would murder another human
- How humans react to emotions
- How to tell if someone lies

Generative Pre-Training



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

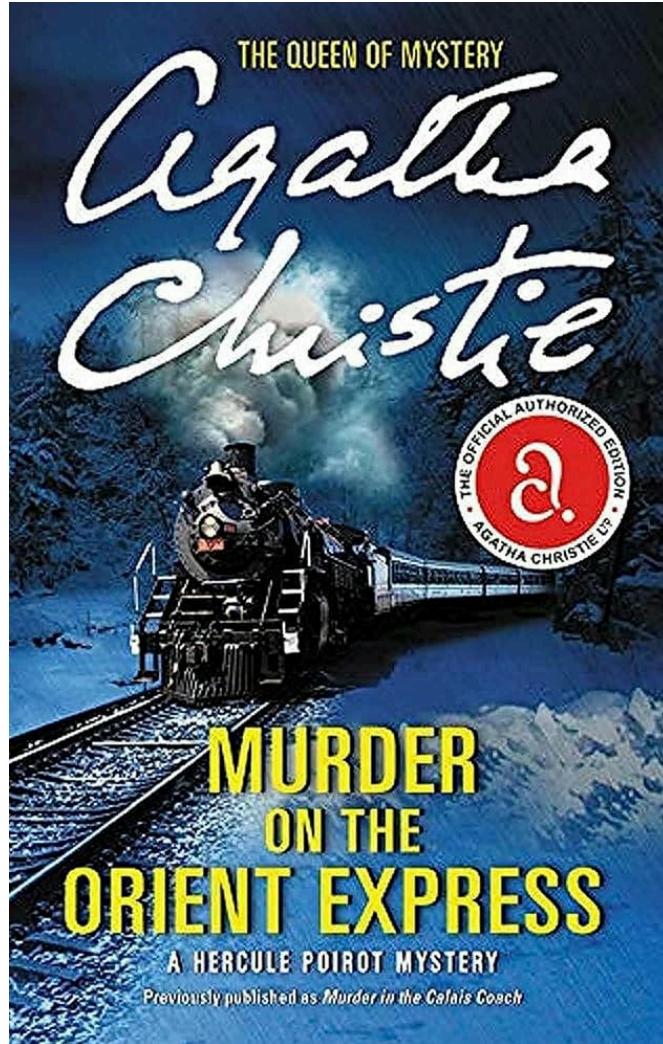
Generative Pre-Training



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

The Books3 dataset contains 200,000 books

Generative Pre-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

Generative Pre-Training

We can apply this same concept to:

Generative Pre-Training

We can apply this same concept to:

- Predict missing base pairs in a strand of DNA

Generative Pre-Training

We can apply this same concept to:

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

Generative Pre-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

Generative Pre-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!

Generative Pre-Training

What if we put the model in a robot?

Generative Pre-Training

What if we put the model in a robot?

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

Generative Pre-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

Generative Pre-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

Call this a **world model** because it models the structure of our world

Generative Pre-Training

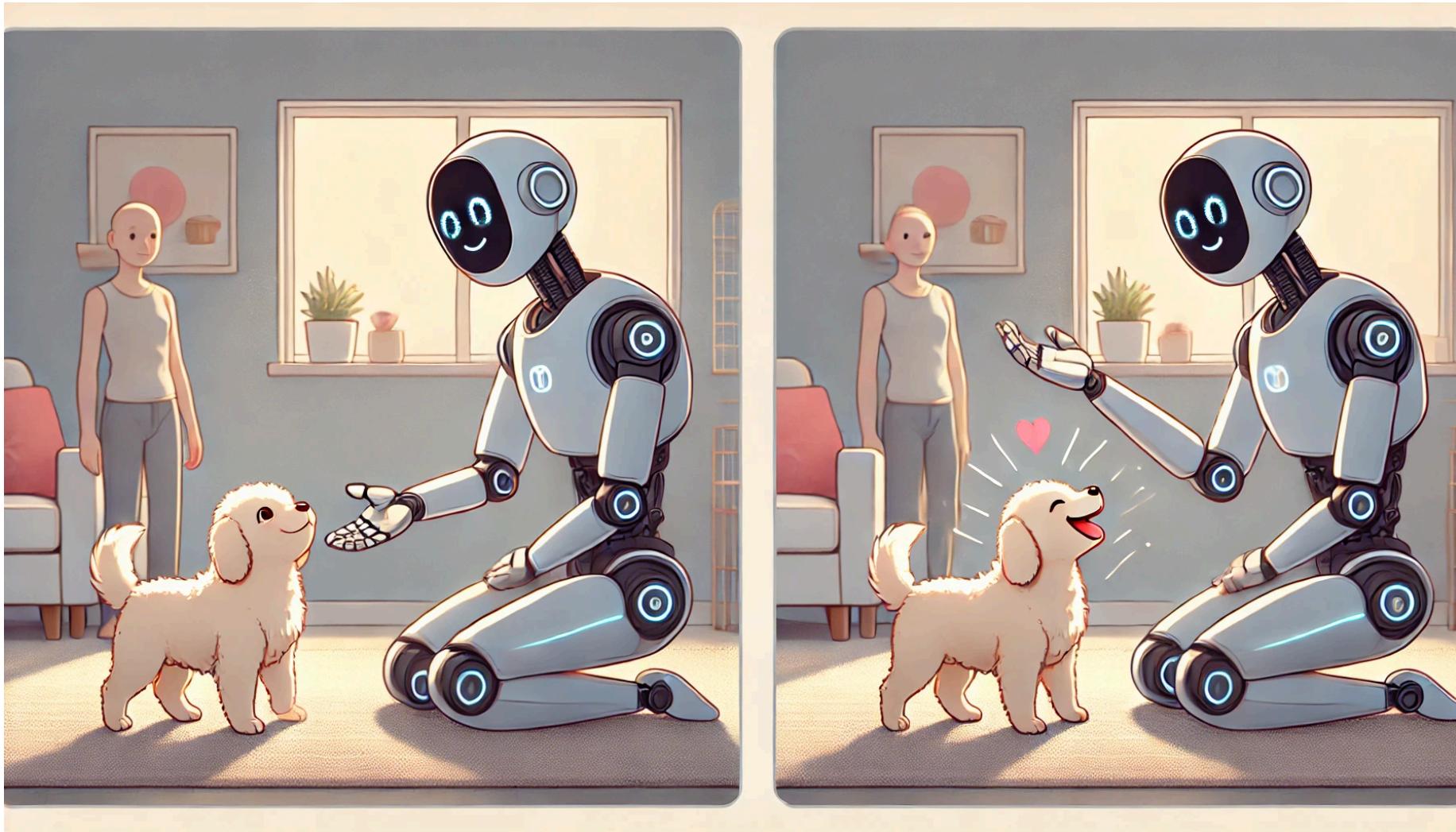
What if we put the model in a robot?

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

Predict what the robot will see next

Call this a **world model** because it models the structure of our world

Generative Pre-Training



Generative Pre-Training

The world model must understand:

Generative Pre-Training

The world model must understand:

- How to control the robot

Generative Pre-Training

The world model must understand:

- How to control the robot
- How to pet a dog

Generative Pre-Training

The world model must understand:

- How to control the robot
- How to pet a dog
- Dogs have feelings

Generative Pre-Training

The world model must understand:

- How to control the robot
- How to pet a dog
- Dogs have feelings
- Petting makes dogs happy

Generative Pre-Training

The world model must understand:

- How to control the robot
- How to pet a dog
- Dogs have feelings
- Petting makes dogs happy
- Dogs smile when happy

Generative Pre-Training

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

Generative Pre-Training

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

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

Generative Pre-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

Generative Pre-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

Generative Pre-Training

These transformers learn and understand the structure of our world

Generative Pre-Training

These transformers learn and understand the structure of our world

But their understanding is trapped

Generative Pre-Training

These transformers learn and understand the structure of our world

But their understanding is trapped

They can only finish sentences or complete pictures

Generative Pre-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?

Generative Pre-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

Generative Pre-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

Generative Pre-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

Generative Pre-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/discussion on reinforcement learning

Closing Remarks

In this course, we started from Gauss in 1795

Closing Remarks

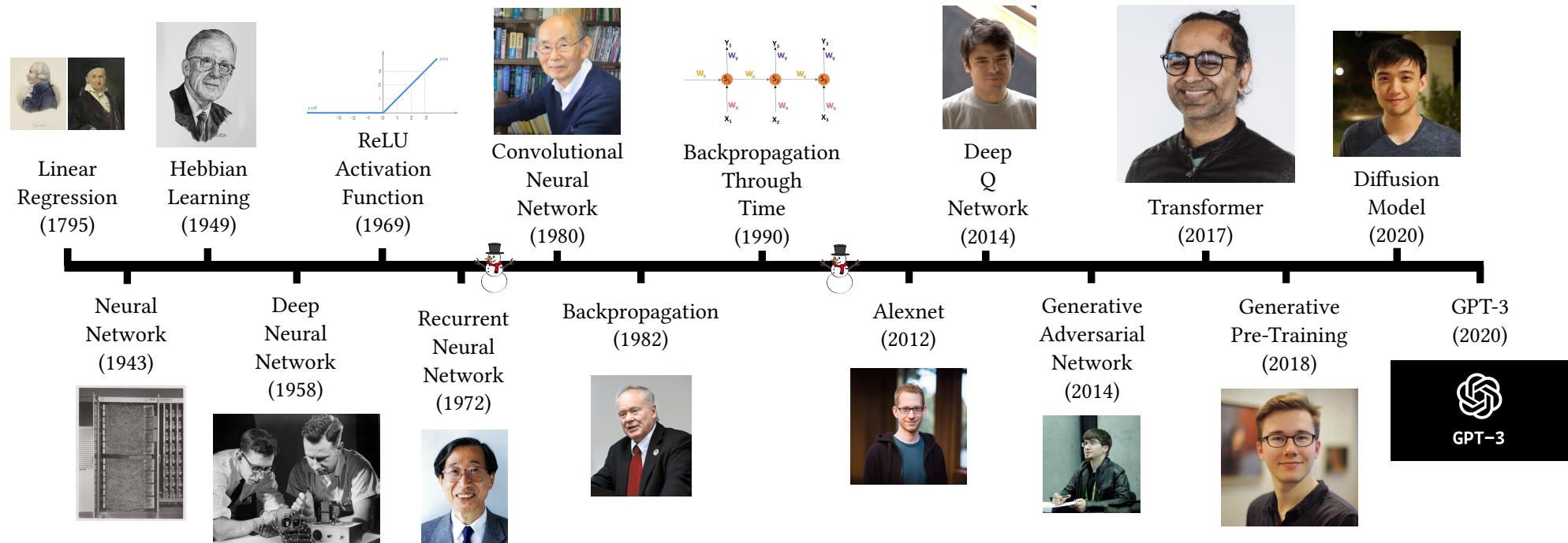
In this course, we started from Gauss in 1795

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

I will return in 10 minutes to see if everyone has finished

Course Evaluation

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

Please scan the QR code and complete the survey

I will return in 10 minutes to see if everyone has finished

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

Course Evaluation

If you participated in class, come see me after class

- Answered a question
- Asked a question

If you would like to collect data for a world model, come see me after class