

# Announcements

- Challenge released.
- Quiz-4 releasing today. 24 hours to work on it.
- No laptops during class.

# Last time

Practical scenarios where RNN is used

RNN Gradients

Attention

Types of attention

Transformers for language

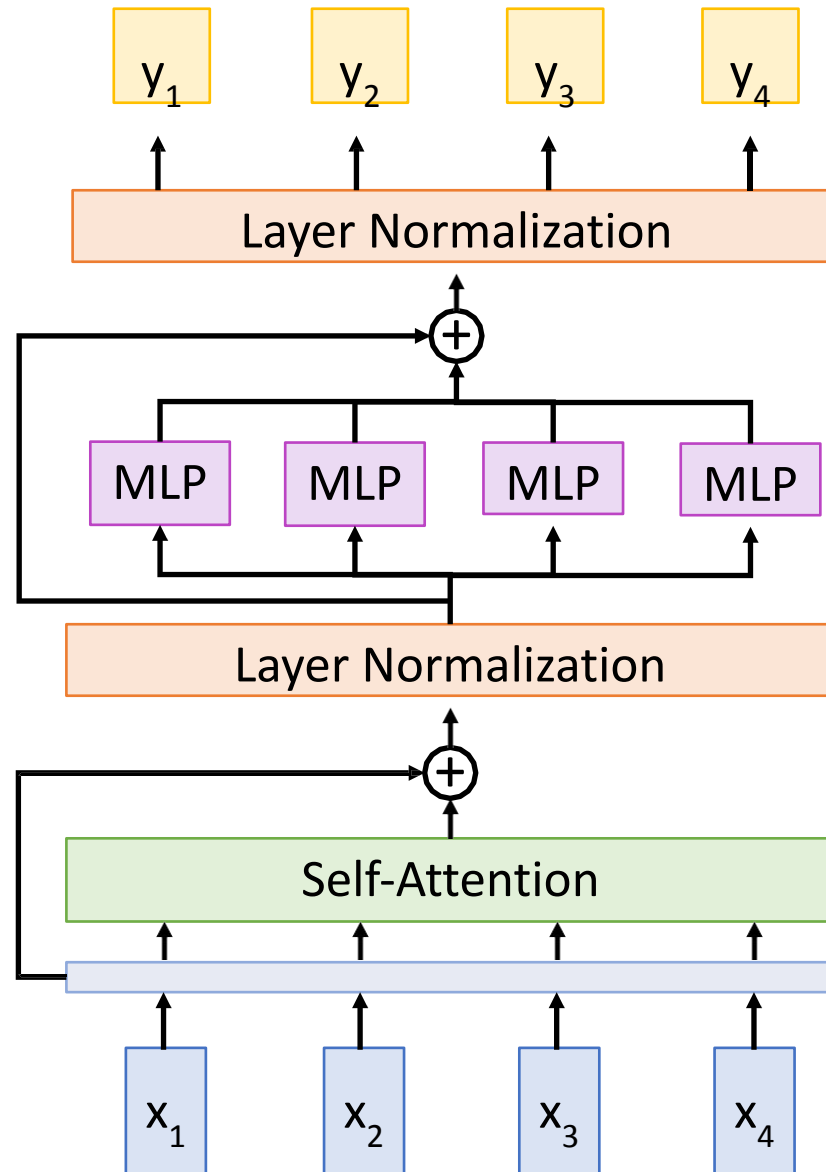
# The Transformer

Residual connection

MLP independently  
on each vector

Residual connection

All vectors interact  
with each other



# The Transformer: Transfer Learning

“ImageNet Moment for Natural Language Processing”

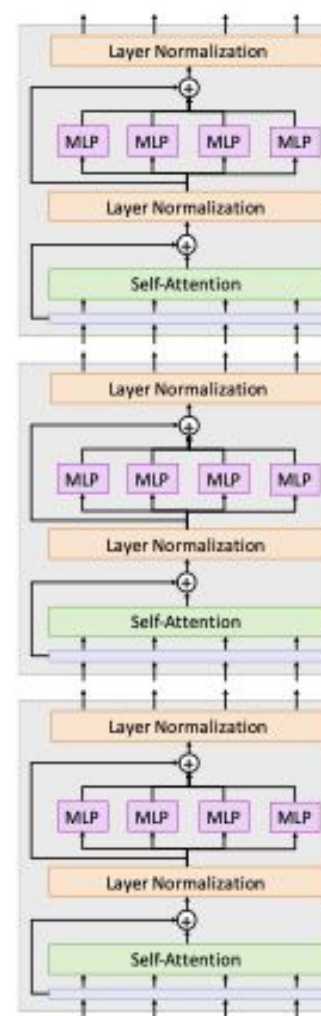
## Pretraining:

Download a lot of text from the internet

Train a giant Transformer model for language modeling

## Finetuning:

Fine-tune the Transformer on your own NLP task



# Today

- A few questions from last class
- Scaling up of transformers
- Extending transformers to computer vision

# Today

- **A few questions from last class**
- Scaling up of transformers
- Extending transformers to computer vision

# What is being learnt as Q, K, V?

Q = query

K = key

V = value

# Attention Layer

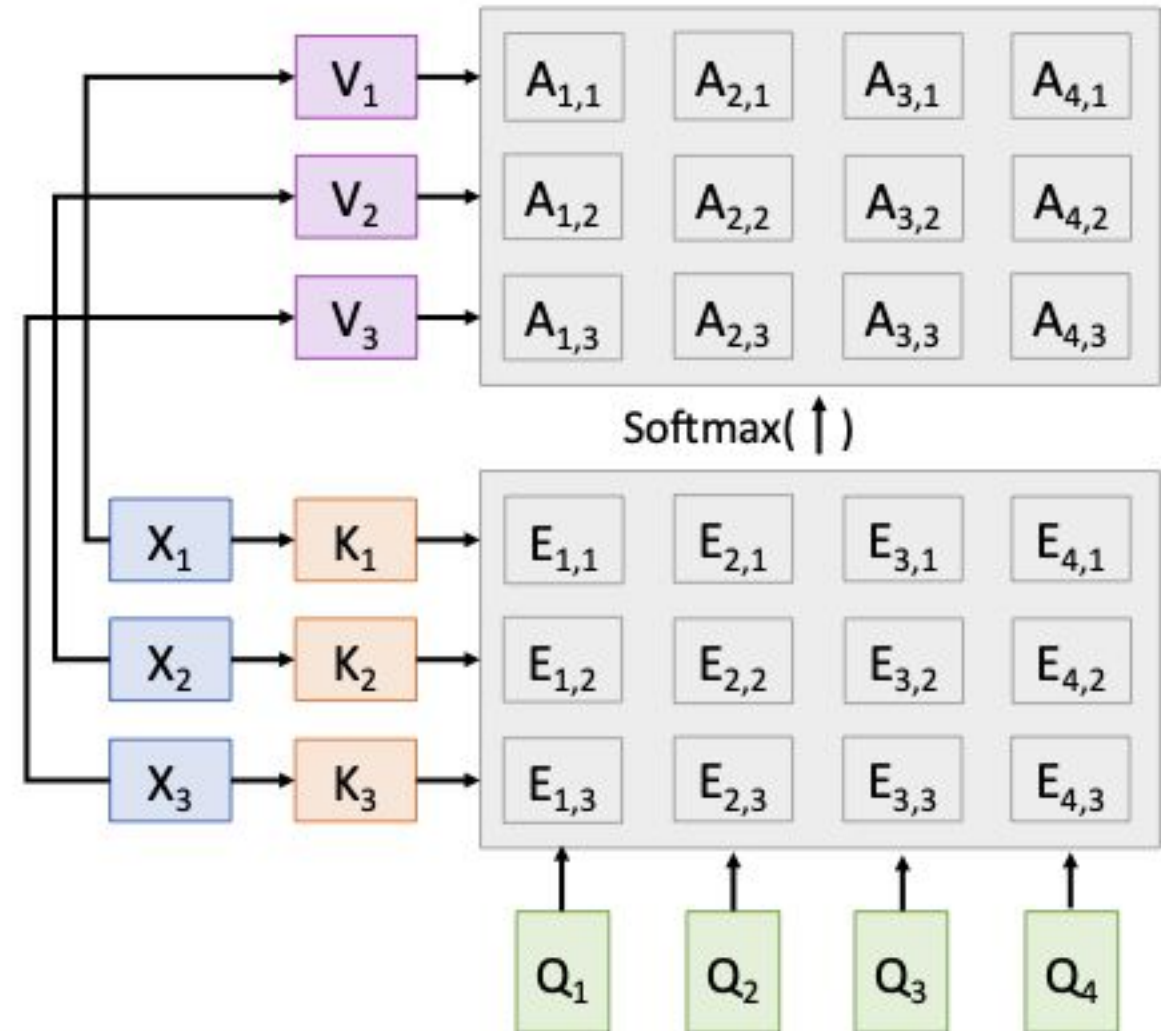
## Inputs:

Query vectors:  $\mathbf{Q}$  (Shape:  $N_Q \times D_Q$ )

Input vectors:  $\mathbf{X}$  (Shape:  $N_X \times D_X$ )

Key matrix:  $\mathbf{W}_K$  (Shape:  $D_X \times D_Q$ )

Value matrix:  $\mathbf{W}_V$  (Shape:  $D_X \times D_V$ )





# Attention Layer

## Inputs:

Query vectors:  $\mathbf{Q}$  (Shape:  $N_Q \times D_Q$ )

Input vectors:  $\mathbf{X}$  (Shape:  $N_X \times D_X$ )

Key matrix:  $\mathbf{W}_K$  (Shape:  $D_X \times D_Q$ )

Value matrix:  $\mathbf{W}_V$  (Shape:  $D_X \times D_V$ )

## Computation:

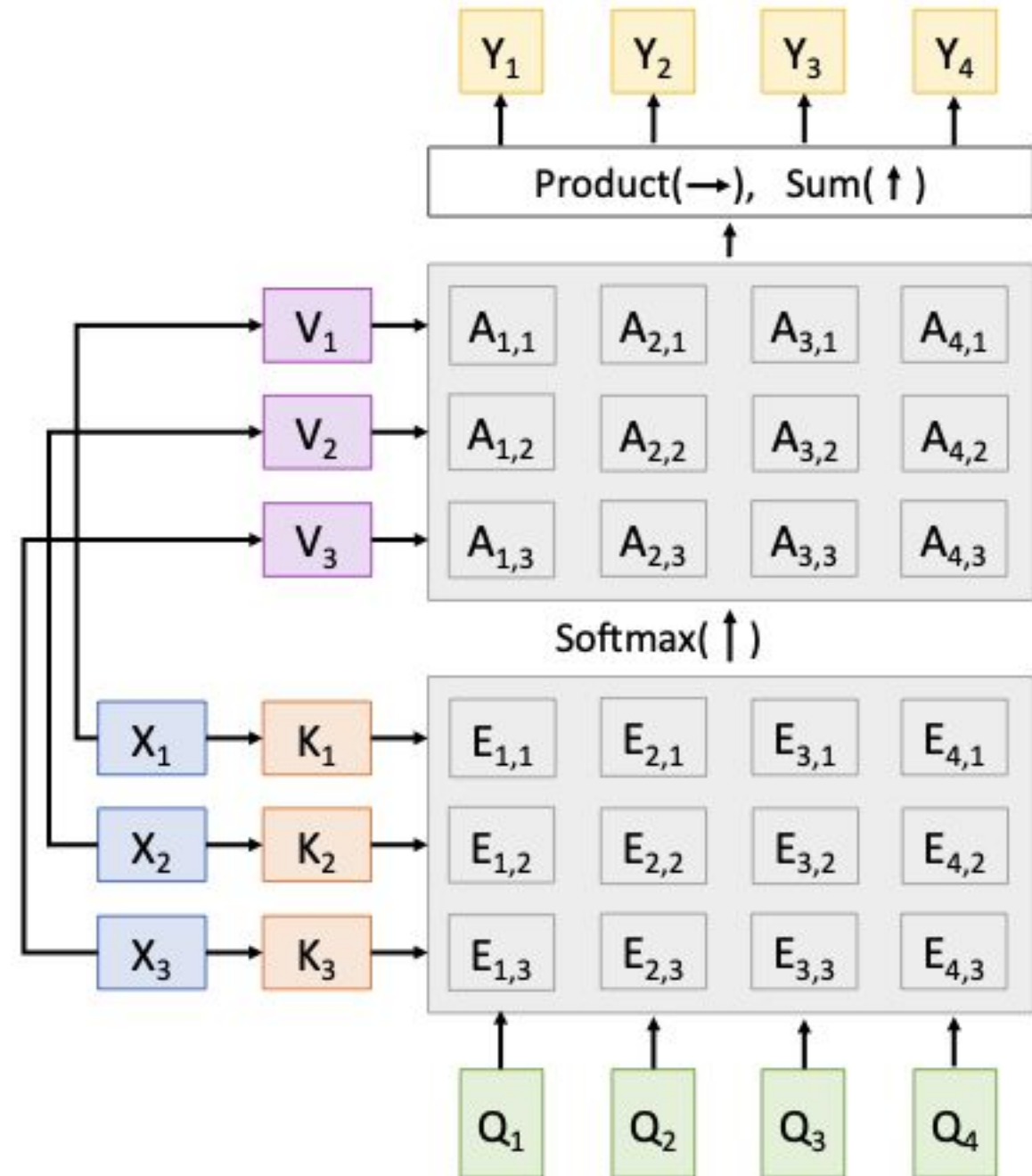
Key vectors:  $\mathbf{K} = \mathbf{XW}_K$  (Shape:  $N_X \times D_Q$ )

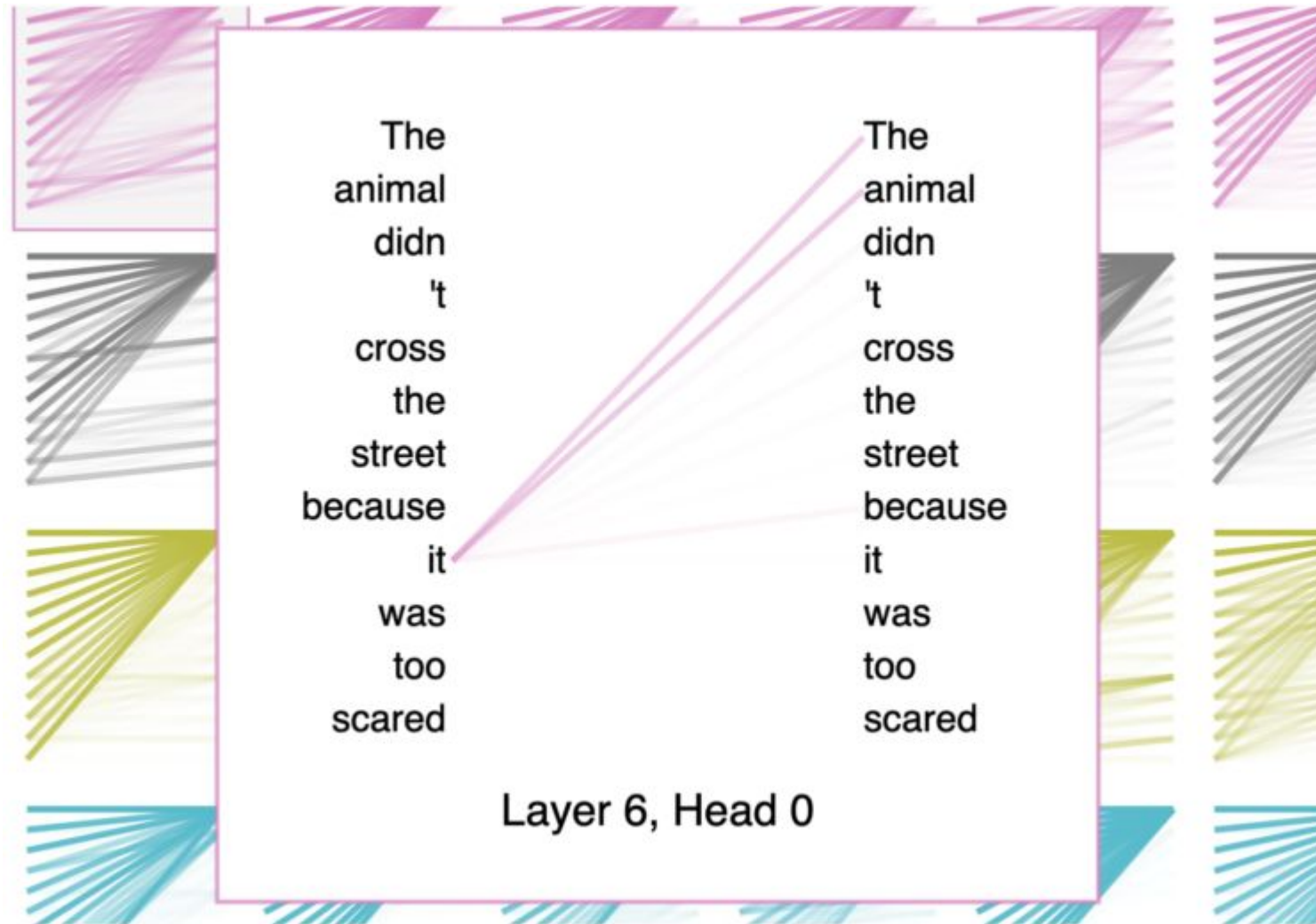
Value Vectors:  $\mathbf{V} = \mathbf{XW}_V$  (Shape:  $N_X \times D_V$ )

Similarities:  $\mathbf{E} = \mathbf{QK}^T / D_Q$  (Shape:  $N_Q \times N_X$ )  $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / D_Q$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$  (Shape:  $N_Q \times N_X$ )

Output vectors:  $\mathbf{Y} = \mathbf{AV}$  (Shape:  $N_Q \times D_V$ )  $Y_i = \sum_j A_{i,j} \mathbf{V}_j$





# Parallelization in attention

**Q = query**

**K = key**

**V = value**

- Matrix multiplication is hardware friendlier than convolution
- For a fixed set of GFLOPs, transformers can be trained faster than Convolutional Networks

# Parallelization in attention

**Q = query**

**K = key**

**V = value**

# Multihead Self-Attention

## Inputs:

Input vectors:  $\mathbf{X}$  (Shape:  $N_x \times D_x$ )

Key matrix:  $\mathbf{W}_k$  (Shape:  $D_x \times D_k$ )

Value matrix:  $\mathbf{W}_v$  (Shape:  $D_x \times D_v$ )

Query matrix:  $\mathbf{W}_q$  (Shape:  $D_x \times D_q$ )

Use H independent  
“Attention Heads” in  
parallel

## Computation:

Query vectors:  $\mathbf{Q} = \mathbf{XW}_q$

Key vectors:  $\mathbf{K} = \mathbf{XW}_k$  (Shape:  $N_x \times D_k$ )

Value Vectors:  $\mathbf{V} = \mathbf{XW}_v$  (Shape:  $N_x \times D_v$ )

Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_q}$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$  (Shape:  $N_x \times N_x$ )

Output vectors:  $\mathbf{Y} = \mathbf{AV}$  (Shape:  $N_x \times D_v$ )  $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

$X_1$

$X_2$

$X_3$

# Multihead Self-Attention

## Inputs:

Input vectors:  $\mathbf{X}$  (Shape:  $N_x \times D_x$ )

Key matrix:  $\mathbf{W}_k$  (Shape:  $D_x \times D_Q$ )

Value matrix:  $\mathbf{W}_v$  (Shape:  $D_x \times D_V$ )

Query matrix:  $\mathbf{W}_Q$  (Shape:  $D_x \times D_Q$ )

Use H independent  
“Attention Heads” in  
parallel

## Computation:

Query vectors:  $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors:  $\mathbf{K} = \mathbf{XW}_k$  (Shape:  $N_x \times D_Q$ )

Value Vectors:  $\mathbf{V} = \mathbf{XW}_v$  (Shape:  $N_x \times D_V$ )

Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_Q}$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$  (Shape:  $N_x \times N_x$ )

Output vectors:  $\mathbf{Y} = \mathbf{AV}$  (Shape:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Split

$X_{1,1}$
$X_{1,2}$
$X_{1,3}$

$X_{2,1}$
$X_{2,2}$
$X_{2,3}$

$X_{3,1}$
$X_{3,2}$
$X_{3,3}$

# Multihead Self-Attention

## Inputs:

Input vectors:  $\mathbf{X}$  (Shape:  $N_x \times D_x$ )

Key matrix:  $\mathbf{W}_K$  (Shape:  $D_x \times D_Q$ )

Value matrix:  $\mathbf{W}_V$  (Shape:  $D_x \times D_V$ )

Query matrix:  $\mathbf{W}_Q$  (Shape:  $D_x \times D_Q$ )

Use H independent  
“Attention Heads” in  
parallel

## Computation:

Query vectors:  $\mathbf{Q} = \mathbf{XW}_Q$

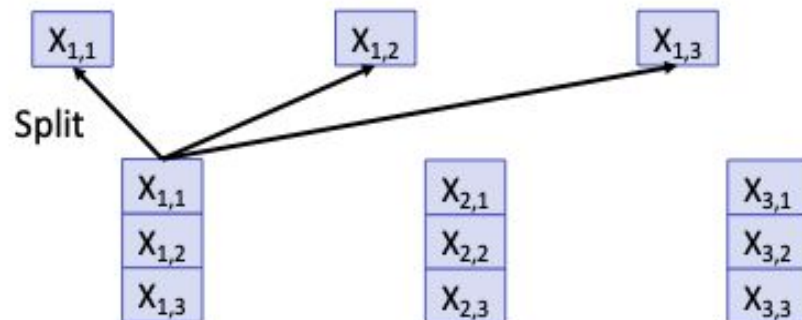
Key vectors:  $\mathbf{K} = \mathbf{XW}_K$  (Shape:  $N_x \times D_Q$ )

Value Vectors:  $\mathbf{V} = \mathbf{XW}_V$  (Shape:  $N_x \times D_V$ )

Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_Q}$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$  (Shape:  $N_x \times N_x$ )

Output vectors:  $\mathbf{Y} = \mathbf{AV}$  (Shape:  $N_x \times D_V$ )  $Y_i = \sum_j A_{i,j} \mathbf{V}_j$





# Multihead Self-Attention

## Inputs:

Input vectors:  $\mathbf{X}$  (Shape:  $N_x \times D_x$ )

Key matrix:  $\mathbf{W}_k$  (Shape:  $D_x \times D_k$ )

Value matrix:  $\mathbf{W}_v$  (Shape:  $D_x \times D_v$ )

Query matrix:  $\mathbf{W}_q$  (Shape:  $D_x \times D_q$ )

Use H independent  
“Attention Heads” in  
parallel

## Computation:

Query vectors:  $\mathbf{Q} = \mathbf{XW}_q$

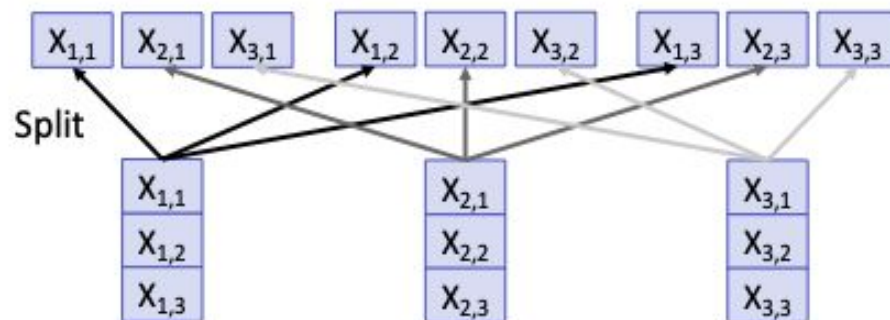
Key vectors:  $\mathbf{K} = \mathbf{XW}_k$  (Shape:  $N_x \times D_k$ )

Value Vectors:  $\mathbf{V} = \mathbf{XW}_v$  (Shape:  $N_x \times D_v$ )

Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_q}$  (Shape:  $N_x \times N_x$ )  $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_q}$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$  (Shape:  $N_x \times N_x$ )

Output vectors:  $\mathbf{Y} = \mathbf{AV}$  (Shape:  $N_x \times D_v$ )  $Y_i = \sum_j A_{i,j} \mathbf{V}_j$



# Multihead Self-Attention

## Inputs:

Input vectors:  $\mathbf{X}$  (Shape:  $N_X \times D_X$ )

Key matrix:  $\mathbf{W}_K$  (Shape:  $D_X \times D_Q$ )

Value matrix:  $\mathbf{W}_V$  (Shape:  $D_X \times D_V$ )

Query matrix:  $\mathbf{W}_Q$  (Shape:  $D_X \times D_Q$ )

Use  $H$  independent  
“Attention Heads” in  
parallel

## Computation:

Query vectors:  $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors:  $\mathbf{K} = \mathbf{XW}_K$  (Shape:  $N_X \times D_Q$ )

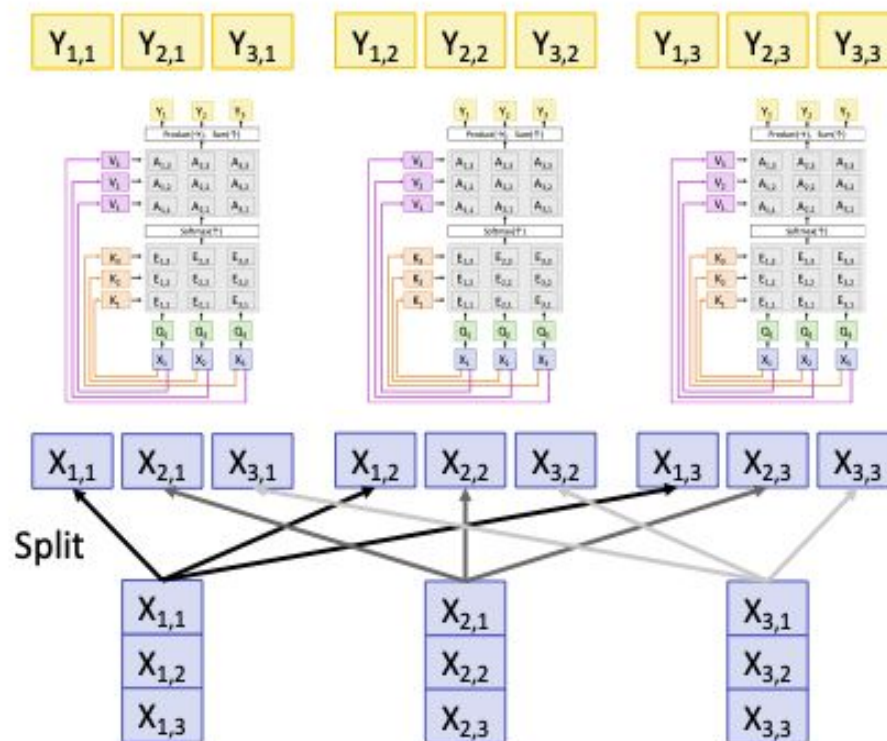
Value Vectors:  $\mathbf{V} = \mathbf{XW}_V$  (Shape:  $N_X \times D_V$ )

Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q}$  (Shape:  $N_X \times N_X$ )  $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_Q}$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$  (Shape:  $N_X \times N_X$ )

Output vectors:  $\mathbf{Y} = \mathbf{AV}$  (Shape:  $N_X \times D_V$ )  $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Run self-attention in parallel on each set of  
input vectors (different weights per head)



# Multihead Self-Attention

## Inputs:

Input vectors:  $\mathbf{X}$  (Shape:  $N_X \times D_X$ )

Key matrix:  $\mathbf{W}_K$  (Shape:  $D_X \times D_Q$ )

Value matrix:  $\mathbf{W}_V$  (Shape:  $D_X \times D_V$ )

Query matrix:  $\mathbf{W}_Q$  (Shape:  $D_X \times D_Q$ )

## Computation:

Query vectors:  $\mathbf{Q} = \mathbf{XW}_Q$

Key vectors:  $\mathbf{K} = \mathbf{XW}_K$  (Shape:  $N_X \times D_Q$ )

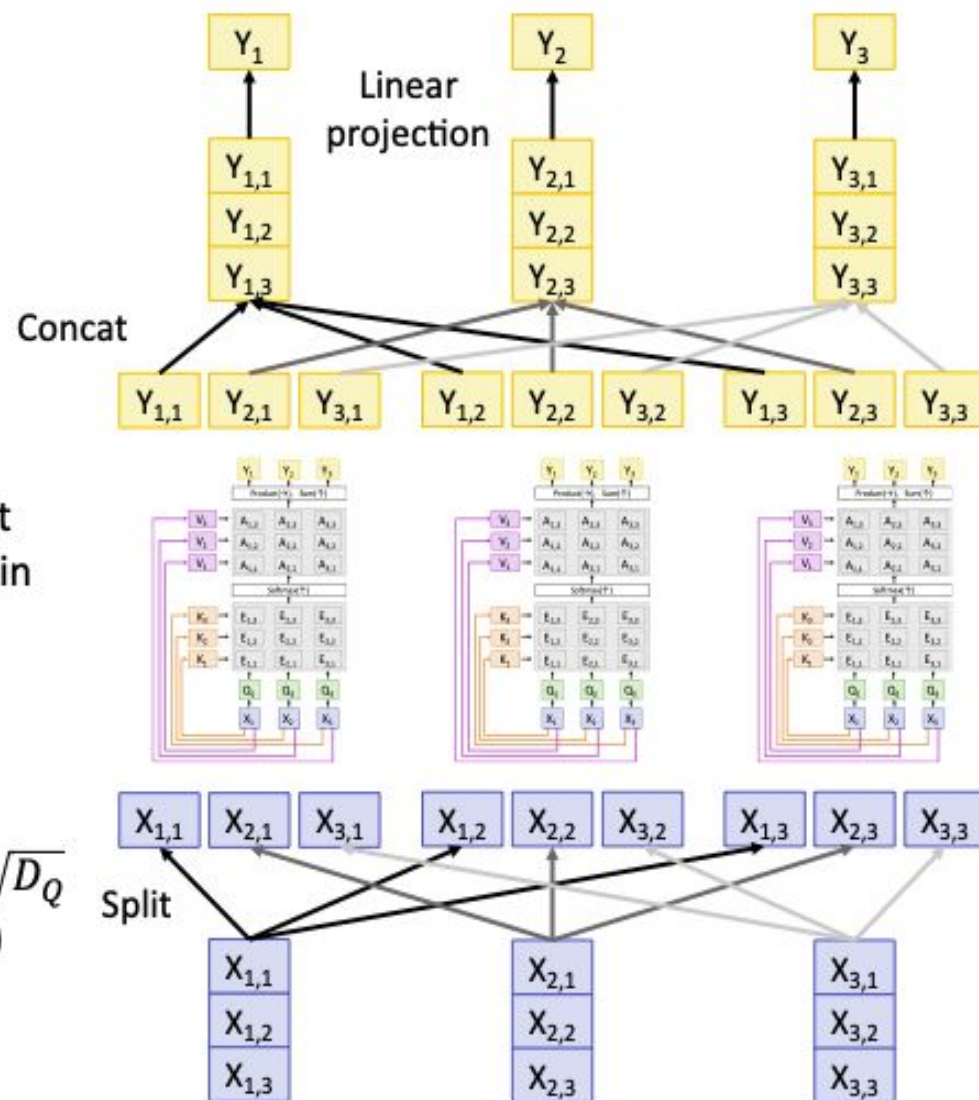
Value Vectors:  $\mathbf{V} = \mathbf{XW}_V$  (Shape:  $N_X \times D_V$ )

Similarities:  $\mathbf{E} = \mathbf{QK}^T / \sqrt{D_Q}$  (Shape:  $N_X \times N_X$ )  $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_Q}$

Attention weights:  $\mathbf{A} = \text{softmax}(\mathbf{E}, \text{dim}=1)$  (Shape:  $N_X \times N_X$ )

Output vectors:  $\mathbf{Y} = \mathbf{AV}$  (Shape:  $N_X \times D_V$ )  $Y_i = \sum_j A_{i,j} \mathbf{V}_j$

Use H independent  
“Attention Heads” in  
parallel



# Summary so far

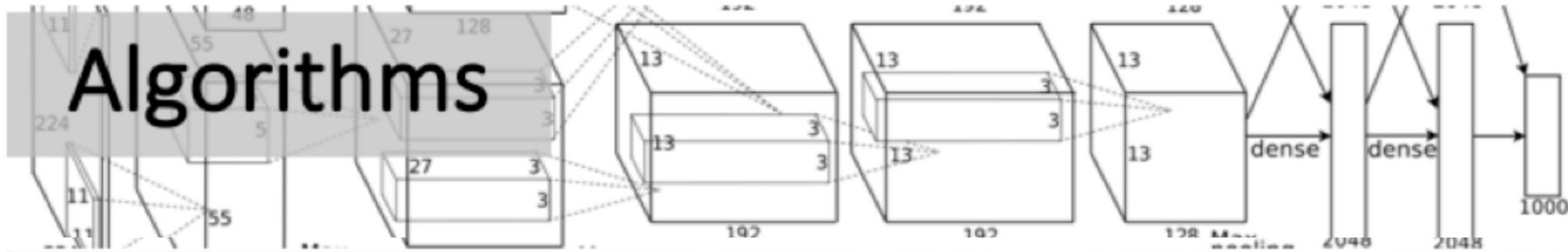
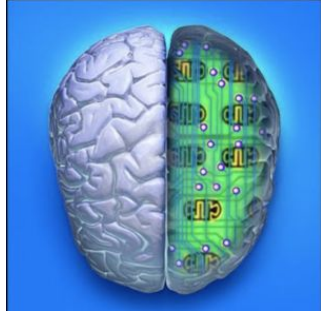
- GPUs favor matrix multiplications over convolutions
- Multi-head attention allows parallel processing of tokens

# Today

- A few questions from last class
- **Scaling up of transformers**
- Extending transformers to computer vision



# What helped us get here?



- CNNs
- Transformers



- Corpus of webdata
- Weak and self supervision

# Scaling up Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160 GB	1024x V100 GPU (1 day)

# Scaling up Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160 GB	1024x V100 GPU (1 day)
GPT-2	48	1600	?	1.5B	40 GB	
Megatron-LM	72	3072	32	8.3B	174 GB	512x V100 GPU (9 days)
Turing-NLG	78	4256	28	17B	?	256x V100 GPU
GPT-3	96	12,288	96	175B	694GB	?
Gopher	80	16,384	128	280B	10.55 TB	4096x TPUv3 (38 days)

Rae et al, "Scaling Language Models: Methods, Analysis, & Insights from Training Gopher", arXiv 2021

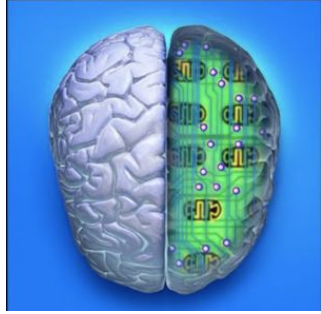


# Scaling up Transformers

\$3,768,320 on Google Cloud (eval price)

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160 GB	1024x V100 GPU (1 day)
GPT-2	48	1600	?	1.5B	40 GB	
Megatron-LM	72	3072	32	8.3B	174 GB	512x V100 GPU (9 days)
Turing-NLG	78	4256	28	17B	?	256x V100 GPU
GPT-3	96	12,288	96	175B	694GB	?
Gopher	80	16,384	128	280B	10.55 TB	4096x TPUv3 (38 days)

# Summary: what helped us get here?



- CNNs
- Transformers



- Corpus of webdata
- Weak and self supervision



- A big reason for AI advancements

Examples of how generic these models are

# GPT-3: Programming by prompt

**Input / Output examples**

**Test example**

*Italics: Completion by GPT-3*

**Poor English input:** I eated the purple berries.

**Good English output:** I ate the purple berries.

**Poor English input:** Thank you for picking me as your designer. I'd appreciate it.

**Good English output:** Thank you for choosing me as your designer. I appreciate it.

**Poor English input:** The mentioned changes have done. or I did the alteration that you requested. or I changed things you wanted and did the modifications.

**Good English output:** The requested changes have been made. or I made the alteration that you requested. or I changed things you wanted and made the modifications.

**Poor English input:** I'd be more than happy to work with you in another project.

**Good English output:**



# GPT-3: Programming by prompt

**Input / Output examples**

**Test example**

*Italics: Completion by GPT-3*

**Poor English input:** I eated the purple berries.

**Good English output:** I ate the purple berries.

**Poor English input:** Thank you for picking me as your designer. I'd appreciate it.

**Good English output:** Thank you for choosing me as your designer. I appreciate it.

**Poor English input:** The mentioned changes have done. or I did the alteration that you requested. or I changed things you wanted and did the modifications.

**Good English output:** The requested changes have been made. or I made the alteration that you requested. or I changed things you wanted and made the modifications.

**Poor English input:** I'd be more than happy to work with you in another project.

**Good English output:** *I'd be more than happy to work with you on another project.*

<https://www.gwern.net/GPT-3#successes>

## GPT-3: Literary Style Transfer

**Topic:** humor: literary parodies and homages, pastiches, style parodies.

**Parodies of the fantasy novel series *Harry Potter* in the style of various famous authors:**

**By Ernest Hemingway:**

## GPT-3: Literary Style Transfer

**Topic:** humor: literary parodies and homages, pastiches, style parodies.

**Parodies of the fantasy novel series *Harry Potter* in the style of various famous authors:**

**By Ernest Hemingway:**

**"It was a cold day on Privet Drive. A child cried. Harry felt nothing.** He was dryer than dust. He had been silent too long. He had not felt love. He had scarcely felt hate. Yet the Dementor's Kiss killed nothing. Death didn't leave him less dead than he had been a second before. It wasn't about living or dying really. It wasn't about death. It was about the Dark Lord, his creator, his engenderer. He was not resurrected. He was created again. He was whole again in the courtyard of Malfoy Manor."

## GPT-3: Literary Style Transfer

**Topic:** humor: literary parodies and homages, pastiches, style parodies.

**Parodies of the fantasy novel series *Harry Potter* in the style of various famous authors:**

**By Ernest Hemingway:**

"It was a cold day on Privet Drive. A child cried. Harry felt nothing. He was dryer than dust. He had been silent too long. He had not felt love. He had scarcely felt hate. Yet the Dementor's Kiss killed nothing. Death didn't leave him less dead than he had been a second before. It wasn't about living or dying really. It wasn't about death. It was about the Dark Lord, his creator, his engenderer. He was not resurrected. He was created again. He was whole again in the courtyard of Malfoy Manor."

**By Jane Austen:**

"It is a truth universally acknowledged, that a broken Harry is in want of a book. This he knows to his cost; pressing that wretched nose against the window of a bookstore in Diagon Alley the last thing before he goes into hiding for several years, where upon, he goes straight to Azkaban."



## GPT-3: Literary Style Transfer

**Topic:** humor: literary parodies and homages, pastiches, style parodies.

**Parodies of the fantasy novel series *Harry Potter* in the style of various famous authors:**

**By Ernest Hemingway:**

**"It was a cold day on Privet Drive. A child cried. Harry felt nothing.** He was dryer than dust. He had been silent too long. He had not felt love. He had scarcely felt hate. Yet the Dementor's Kiss killed nothing. Death didn't leave him less dead than he had been a second before. It wasn't about living or dying really. It wasn't about death. It was about the Dark Lord, his creator, his engenderer. He was not resurrected. He was created again. He was whole again in the courtyard of Malfoy Manor."

**By Jane Austen:**

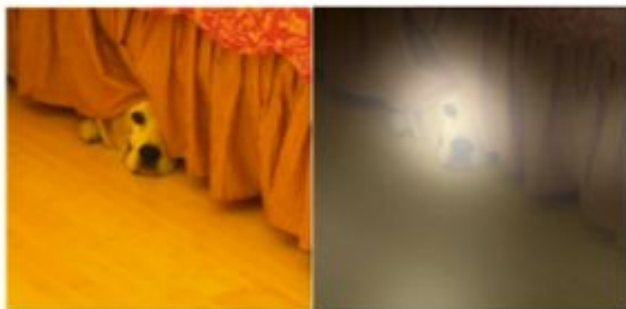
"It is a truth universally acknowledged, that a broken Harry is in want of a book. This he knows to his cost; pressing that wretched nose against the window of a bookstore in Diagon Alley the last thing before he goes into hiding for several years, where upon, he goes straight to Azkaban."

**By Arthur Conan Doyle:**

"Harry pushed at the swinging doors of the bookshop hard, and nearly knocked himself unconscious. He staggered in with his ungainly package, his cheeks scarlet with cold and the shame of having chosen the wrong month to go Christmas shopping. The proprietor of the store, however, didn't cast even a cursory glance at him, being far more interested in an enormous hunk of wood lying in the middle of the floor, which certainly looked like a gallows. Yes, the proprietor said to a reedy old man wearing a bowler hat and a forlorn expression that made Harry want to kick him, I can rent you such a gallows for a small fee."

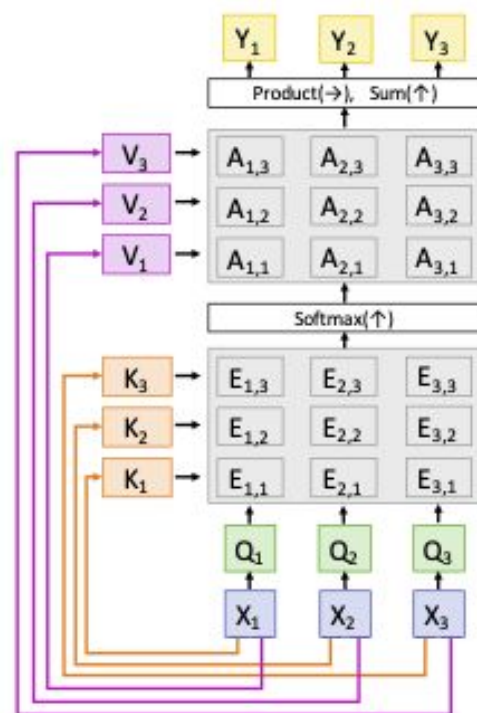
# Summary

Adding **Attention** to RNN models lets them look at different parts of the input at each timestep

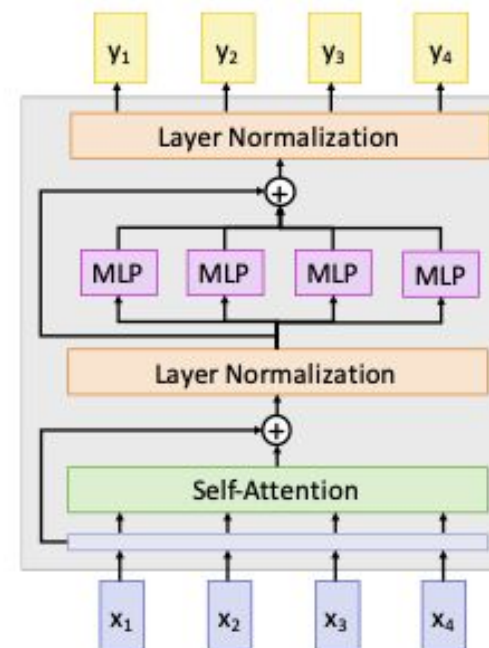


A dog is standing on a hardwood floor.

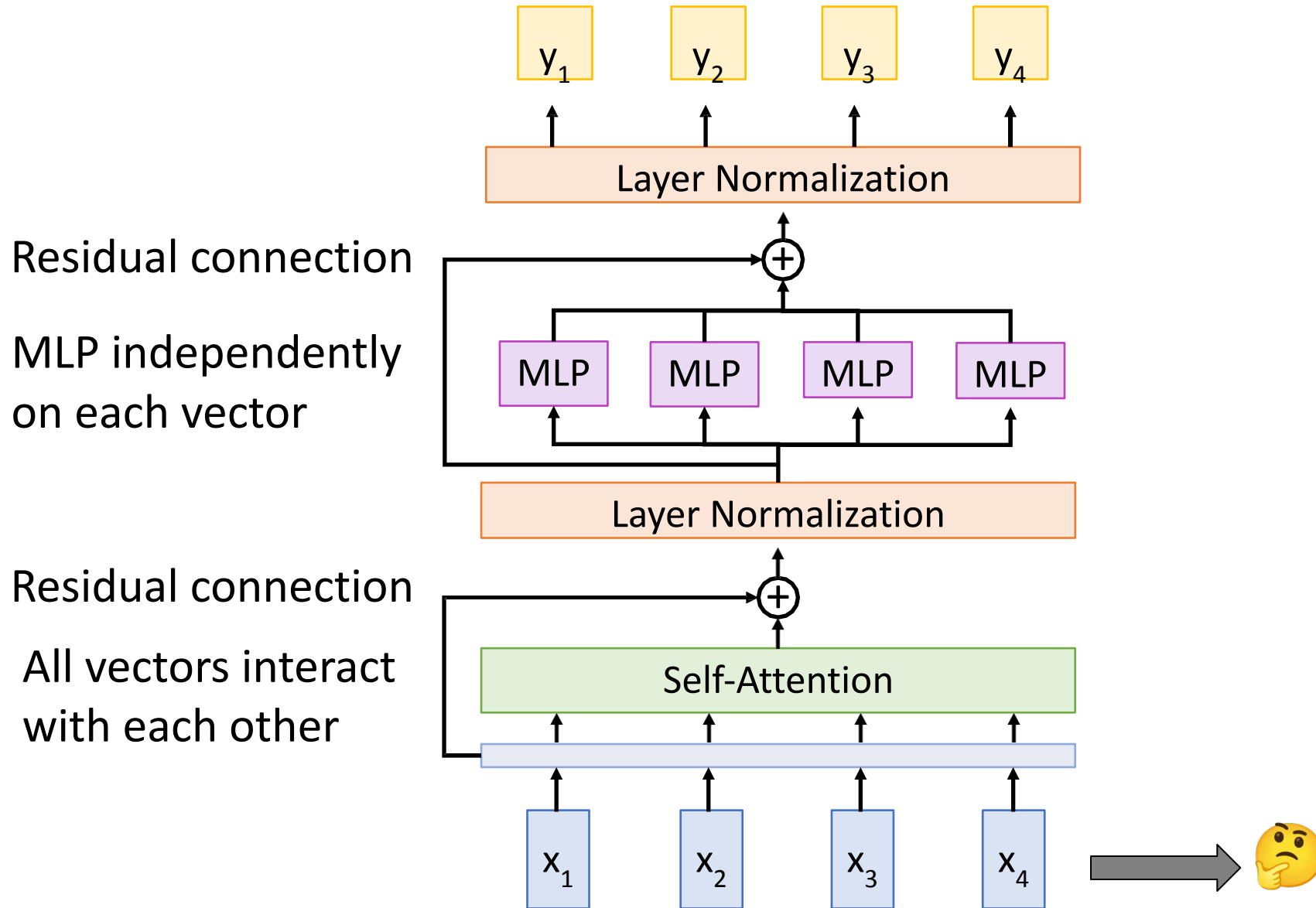
Generalized **Self-Attention** is new, powerful neural network primitive



**Transformers** are a new neural network model that only uses attention



# Words or characters as input



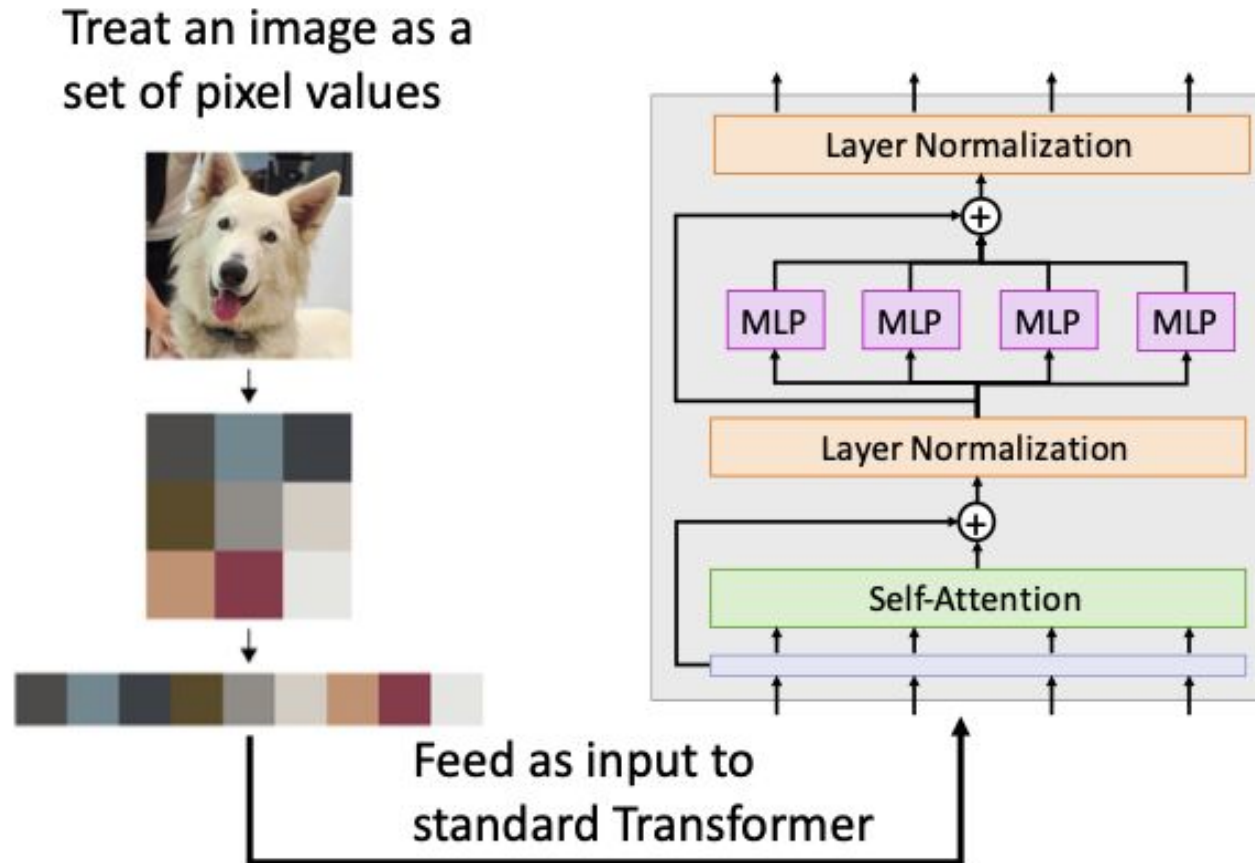
# Today

- A few questions from last class
- Scaling up of transformers
- **Extending transformers to computer vision**



# How do we extend transformers to images?

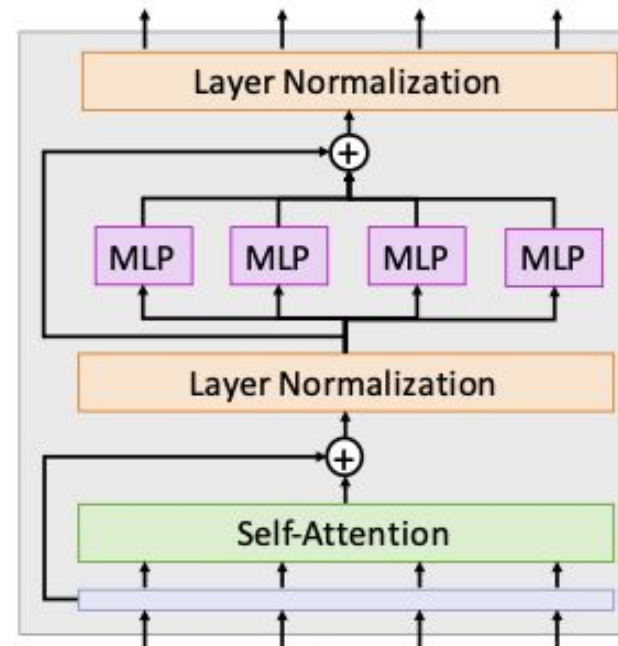
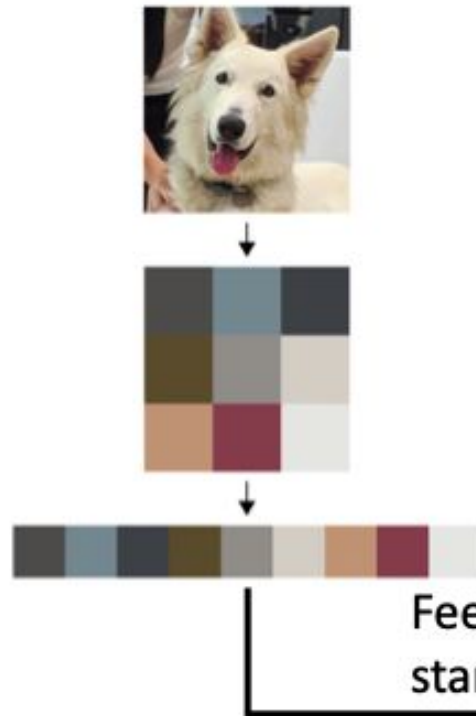
# One option: Pixels as tokens



Chen et al, "Generative Pretraining from Pixels", ICML 2020

# One option: Pixels as tokens

Treat an image as a set of pixel values



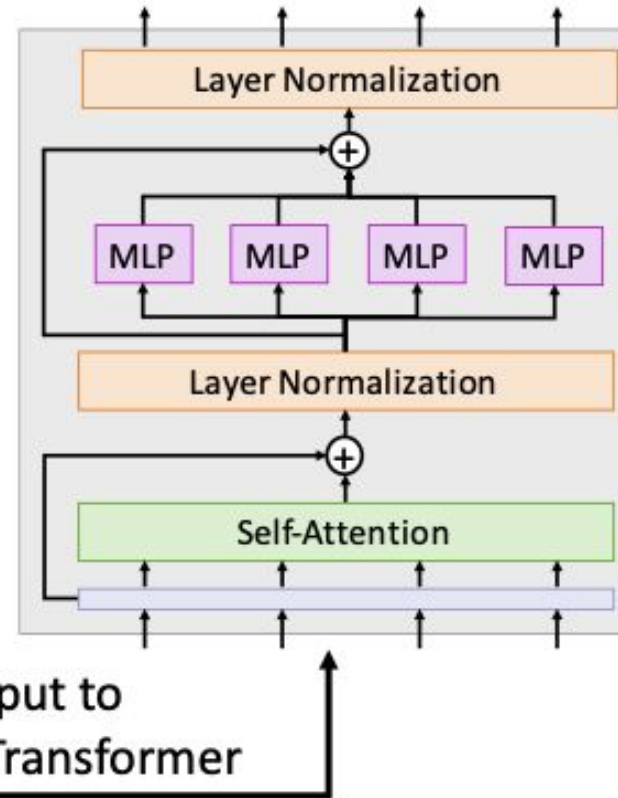
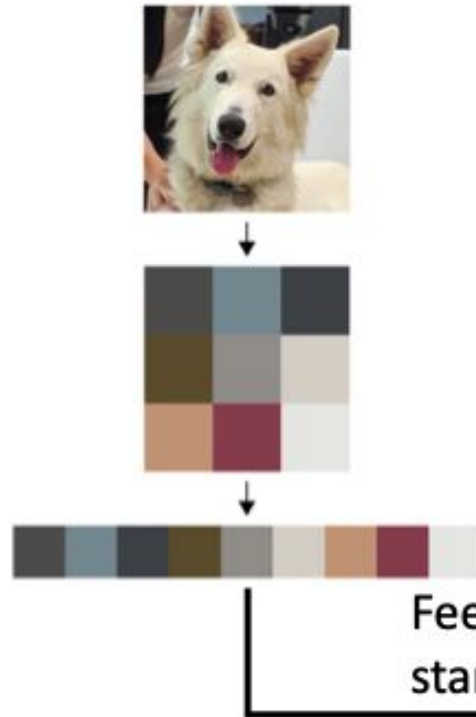
Problem: Memory use!

$R \times R$  image needs  $R^4$  elements per attention matrix



# One option: Pixels as tokens

Treat an image as a set of pixel values

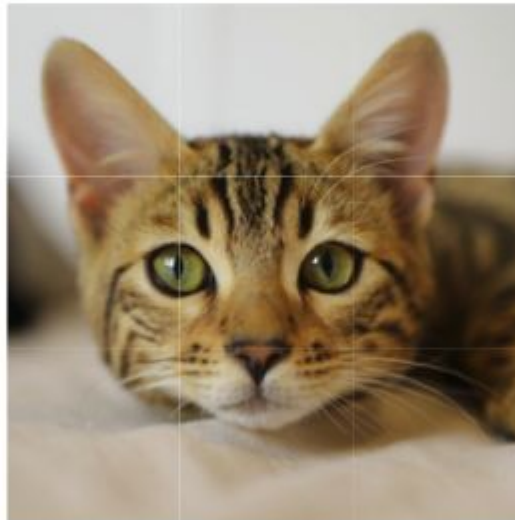


Problem: Memory use!

$R \times R$  image needs  $R^4$  elements per attention matrix

$R=128$ , 48 layers, 16 heads per layer takes 768GB of memory for attention matrices for a single example...

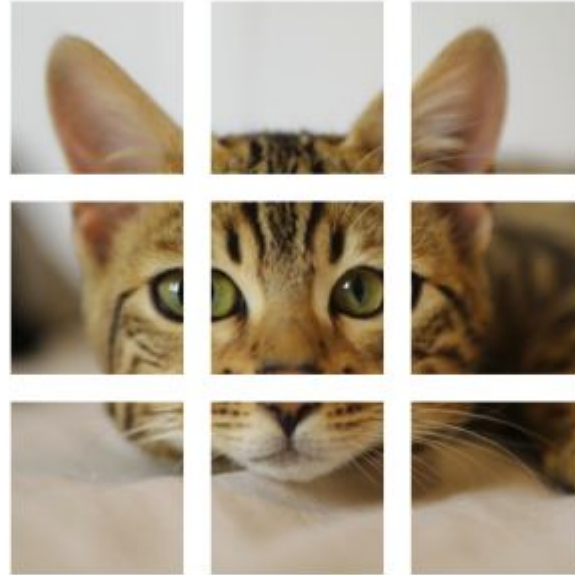




Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

# Second option: Patches as tokens



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

N input patches, each  
of shape 3x16x16



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

Linear projection to  
D-dimensional vector

N input patches, each  
of shape 3x16x16



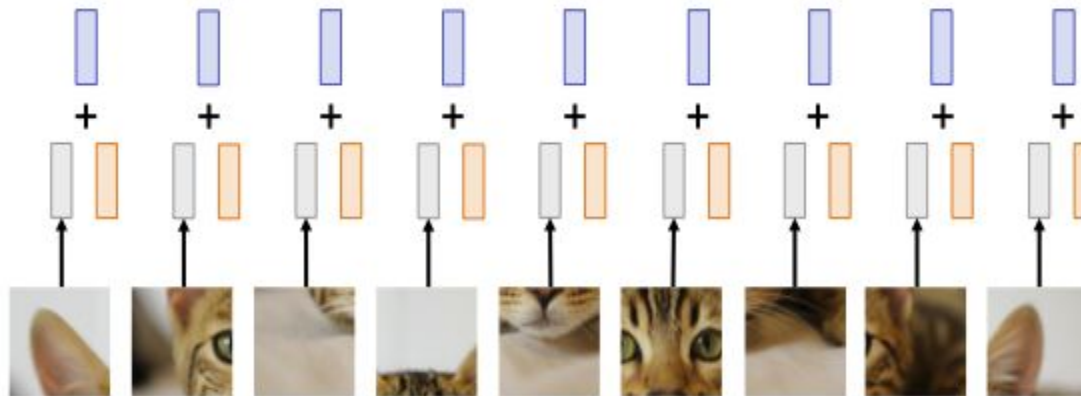
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

Add positional  
embedding: learned D-  
dim vector per position

Linear projection to  
D-dimensional vector

N input patches, each  
of shape 3x16x16



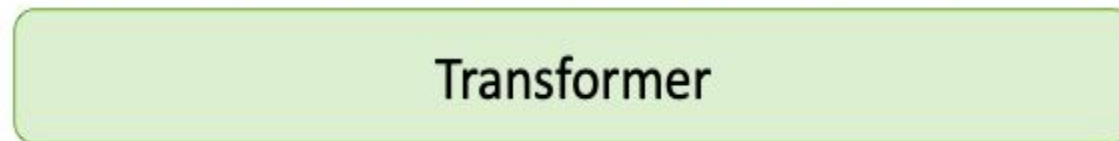
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)

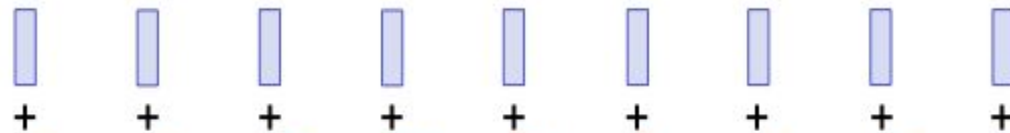
Output vectors



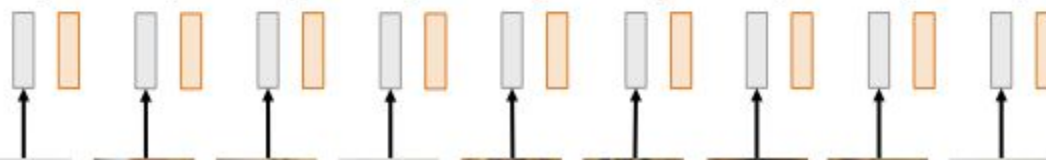
Exact same as  
NLP Transformer!



Add positional  
embedding: learned D-  
dim vector per position



Linear projection to  
D-dimensional vector



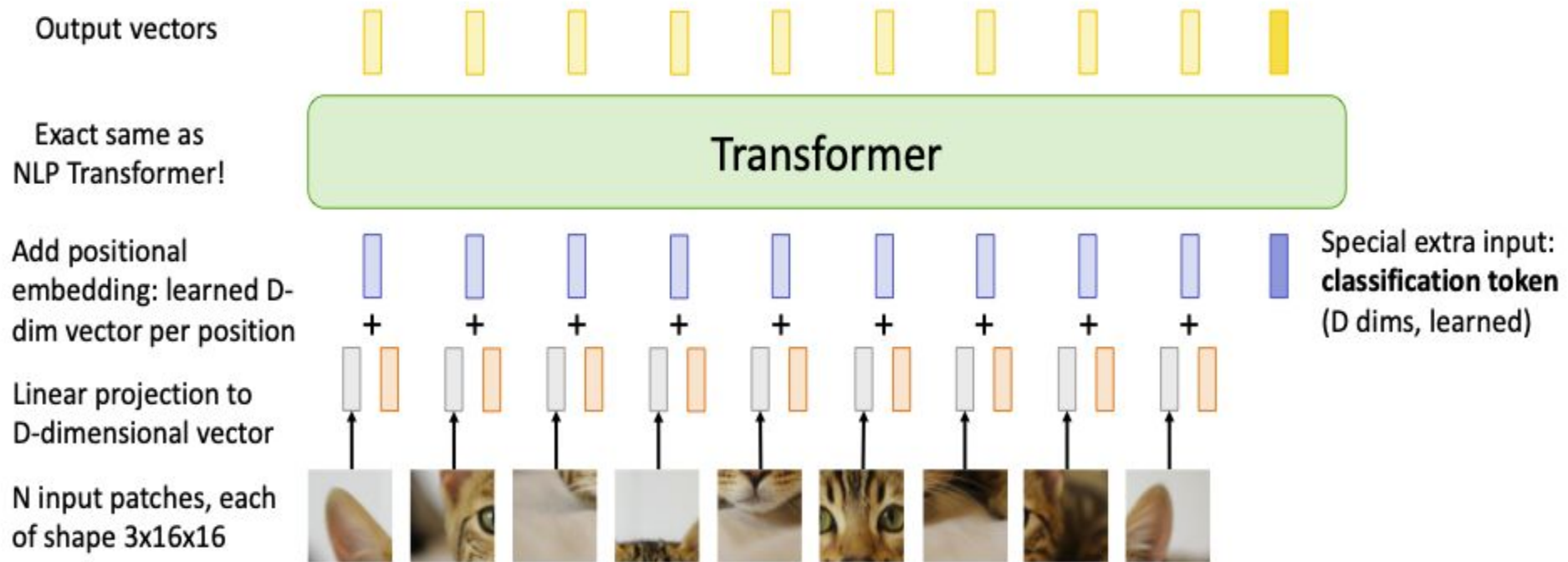
N input patches, each  
of shape 3x16x16



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

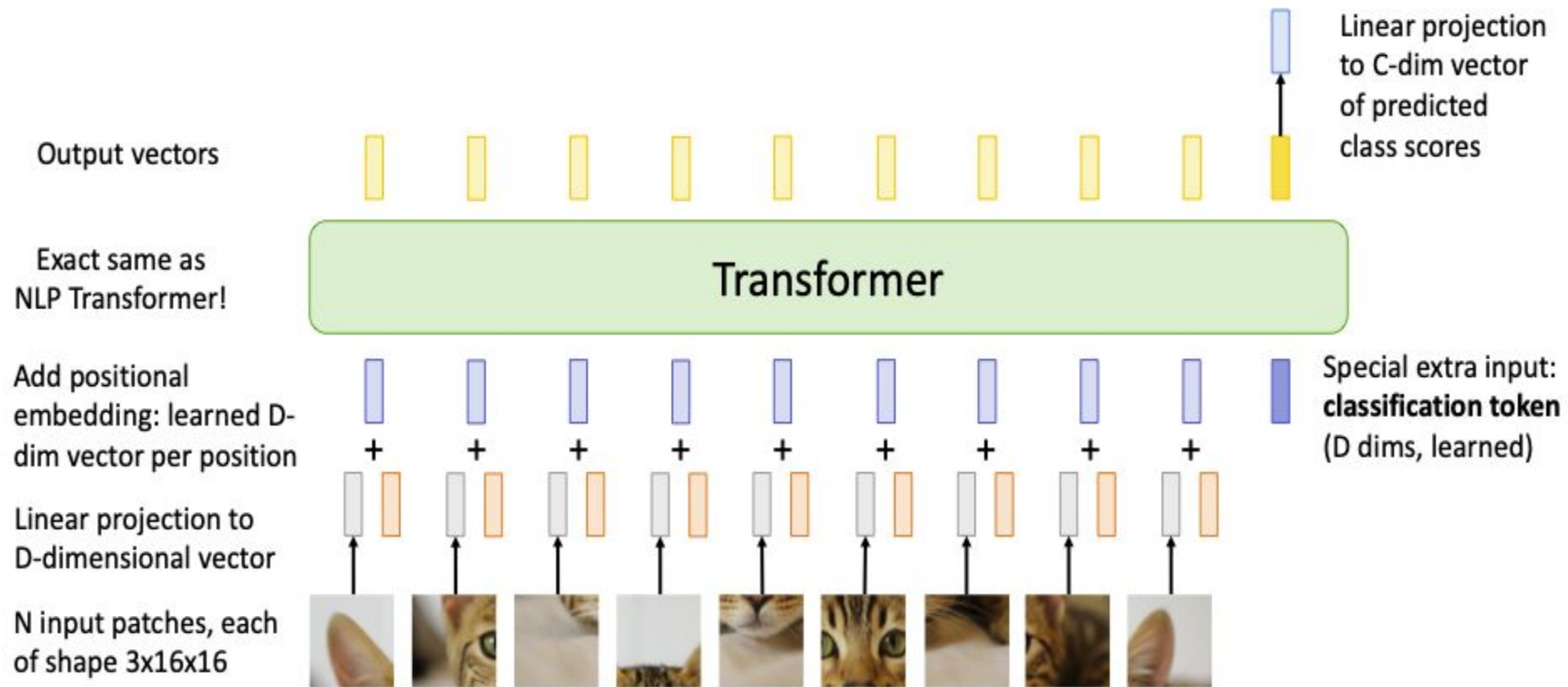
[Cat image](#) is free for commercial  
use under a [Pixabay license](#)





Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

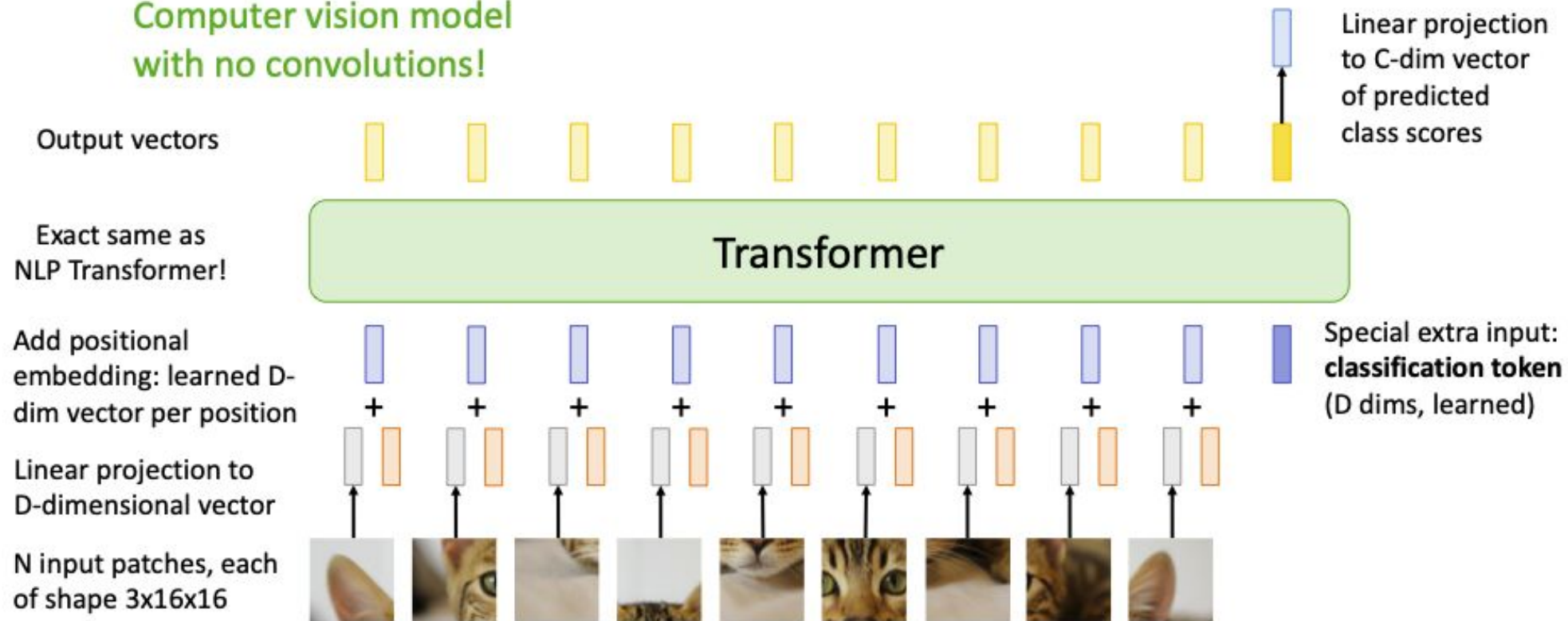


Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformer (ViT)

Computer vision model  
with no convolutions!



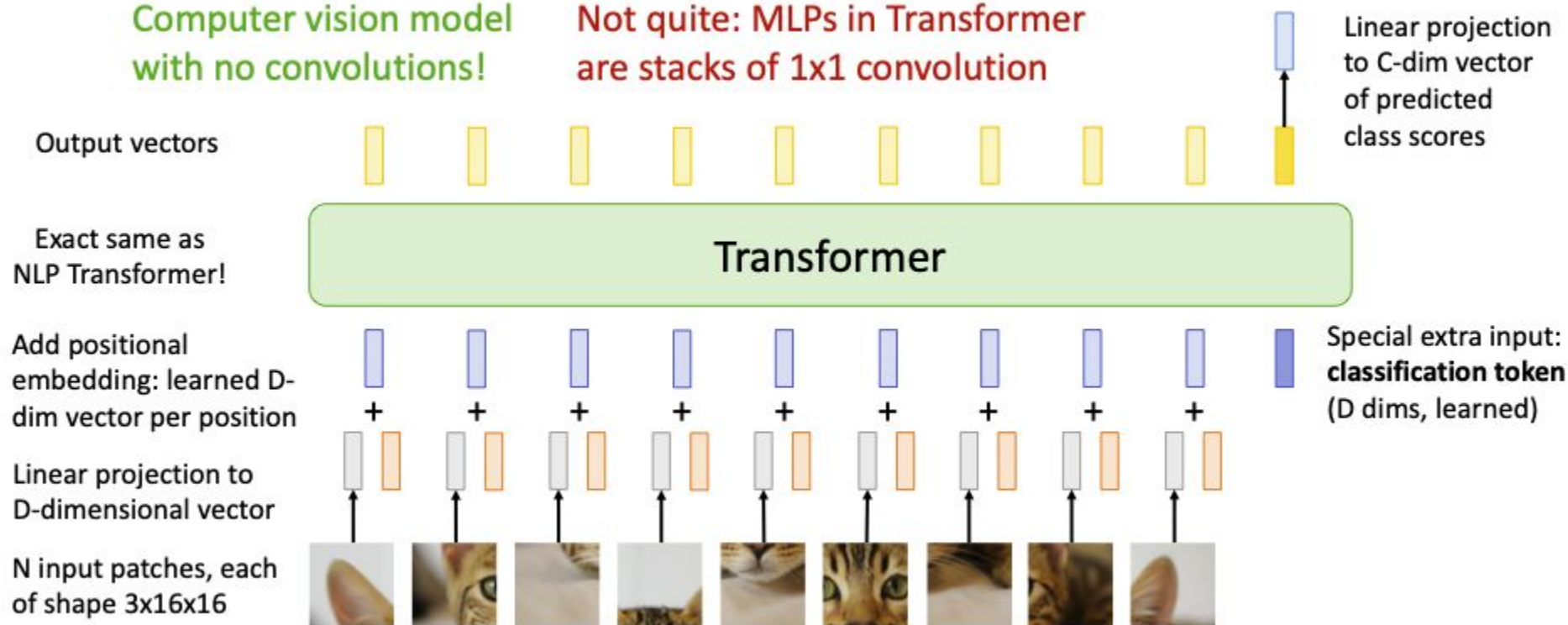
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformer (ViT)

Computer vision model  
with no convolutions!

Not quite: MLPs in Transformer  
are stacks of 1x1 convolution



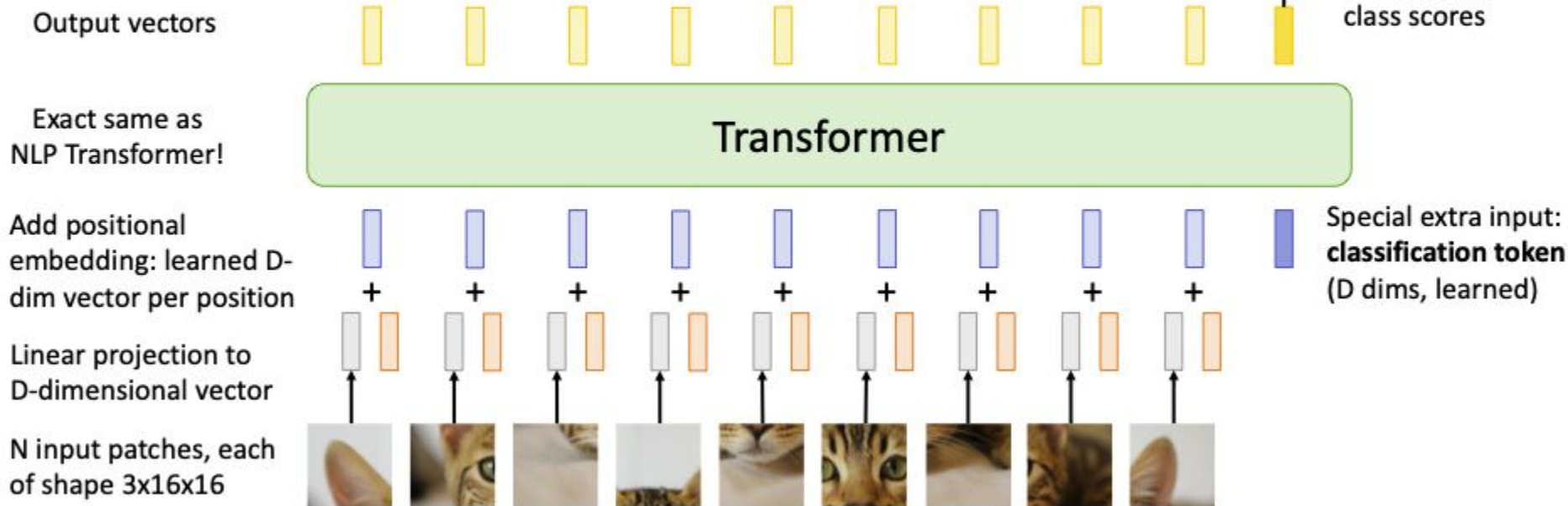
Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial use under a [Pixabay license](#)

# Vision Transformer (ViT)

In practice: take 224x224 input image,  
divide into 14x14 grid of 16x16 pixel  
patches (or 16x16 grid of 14x14 patches)

With 48 layers, 16 heads per  
layer, all attention matrices  
take 112 MB (or 192MB)

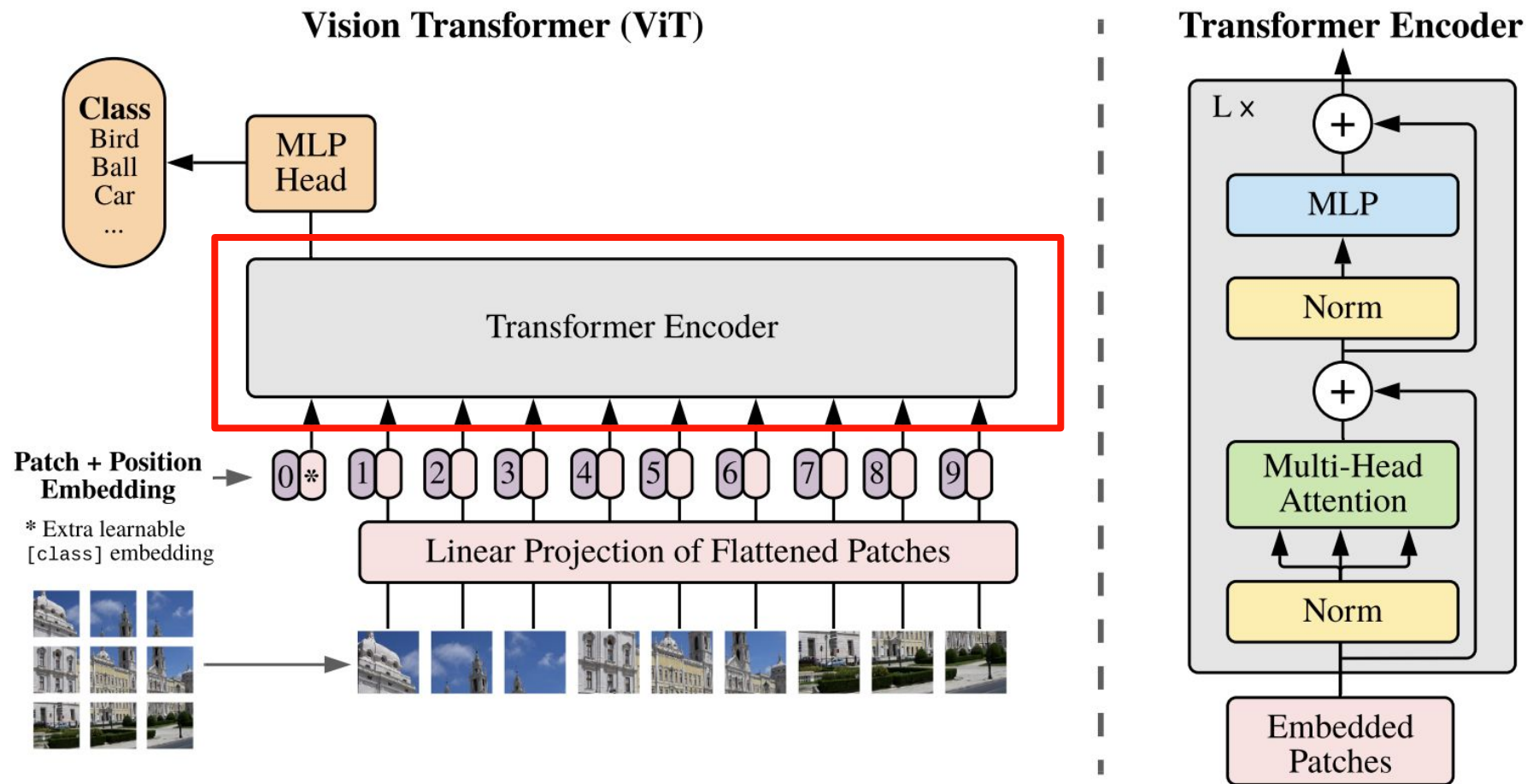


Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

[Cat image](#) is free for commercial  
use under a [Pixabay license](#)



# Putting this all together...

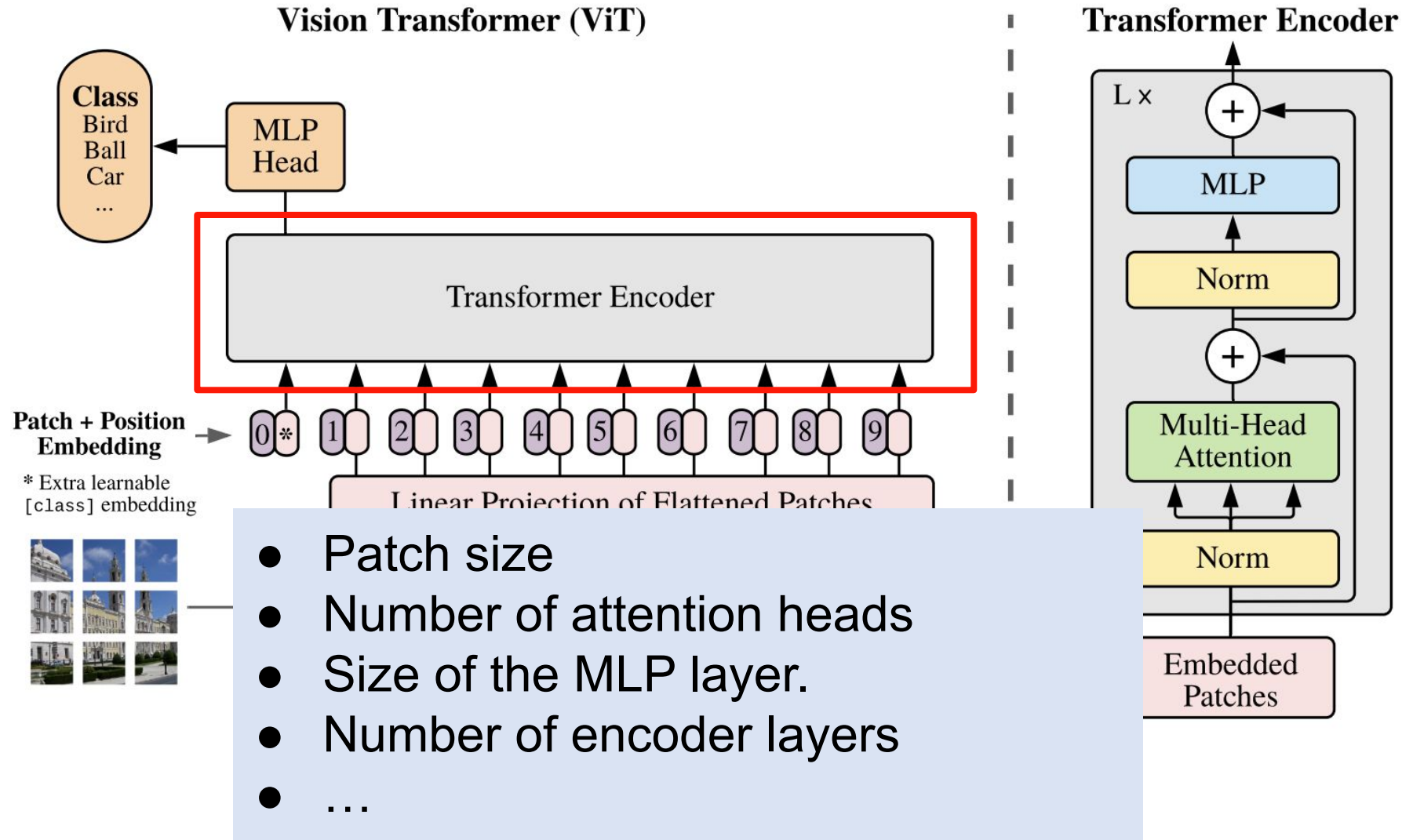




# What are some of the parameter choices involved in designing a ViT?



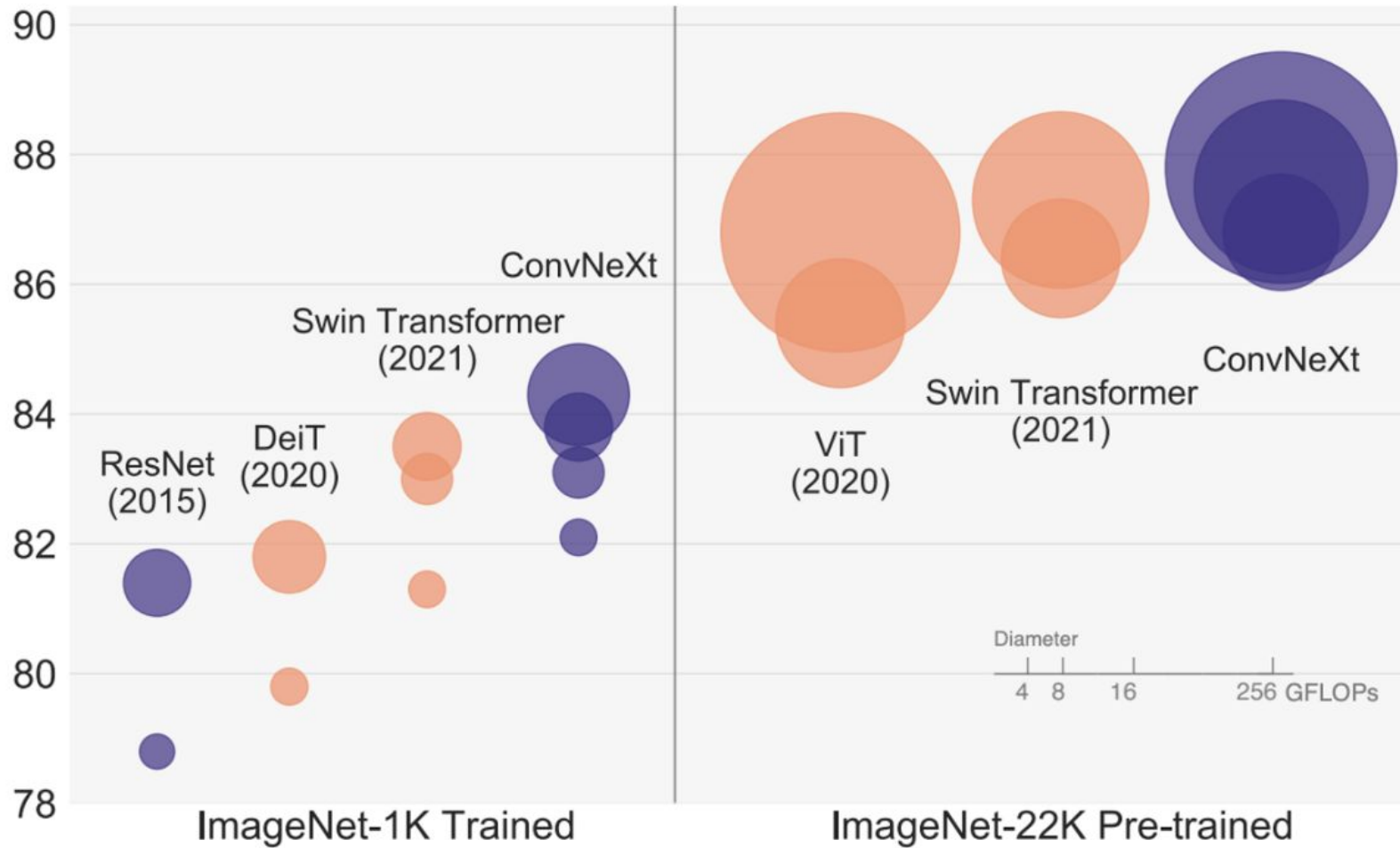
# Putting this all together...



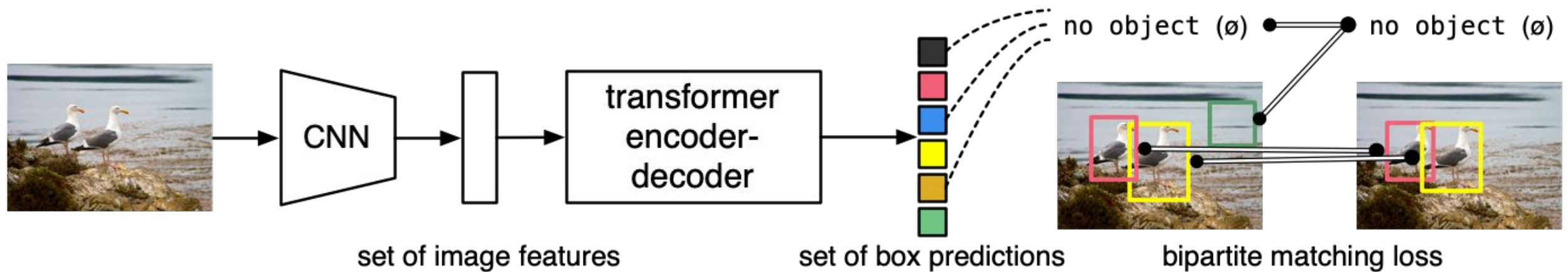
# Convolution

## Transformers

ImageNet-1K Acc.

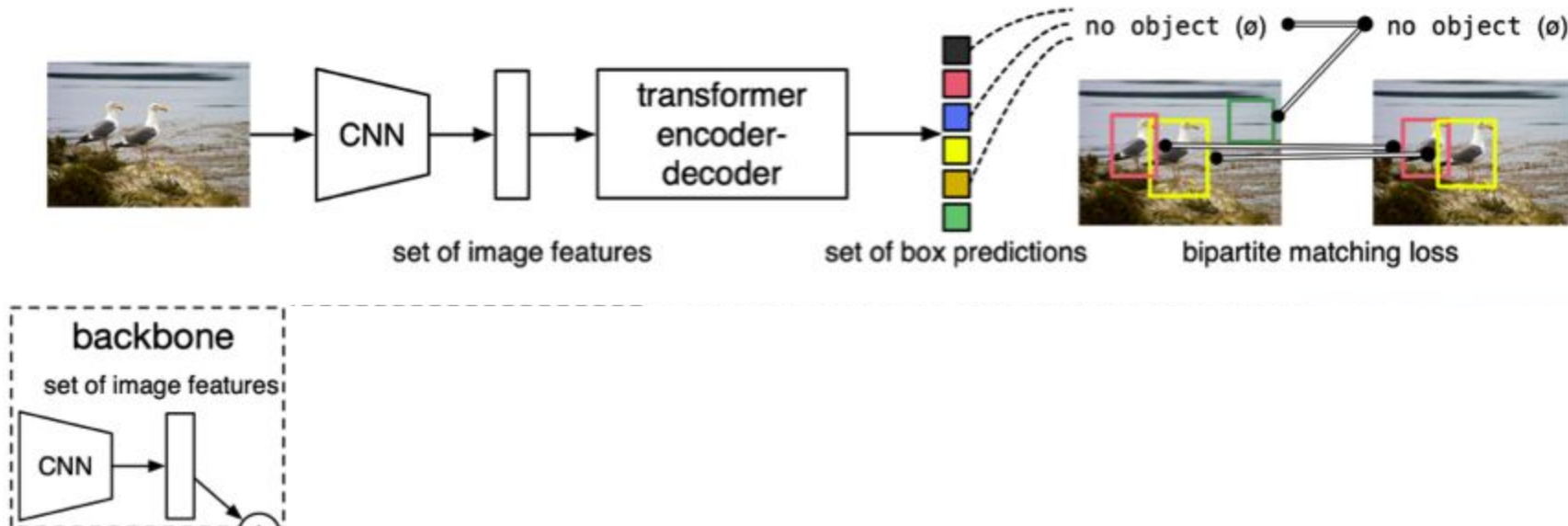


# Convolution 🤝 Transformers



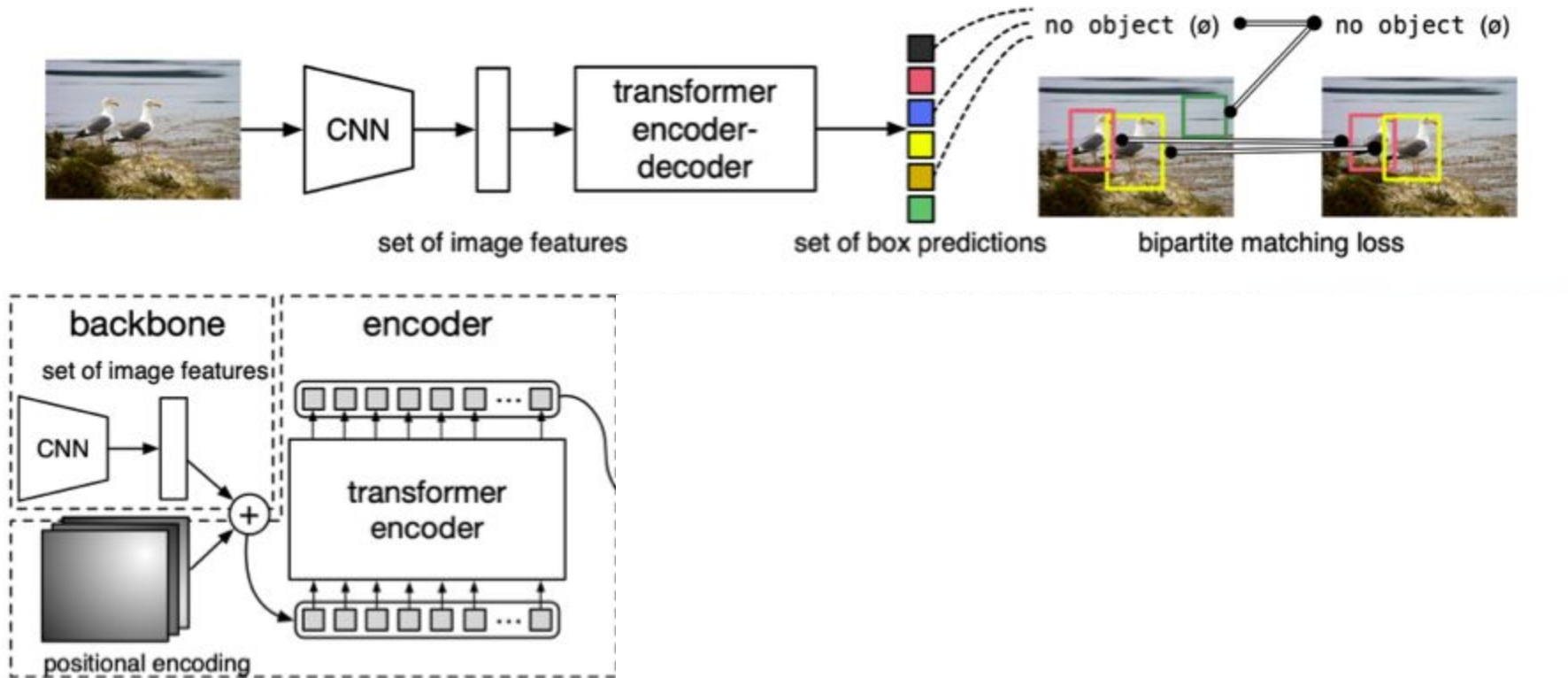
# Convolution 🤝 Transformers

## Object Detection with Transformers: DETR



# Convolution 🤝 Transformers

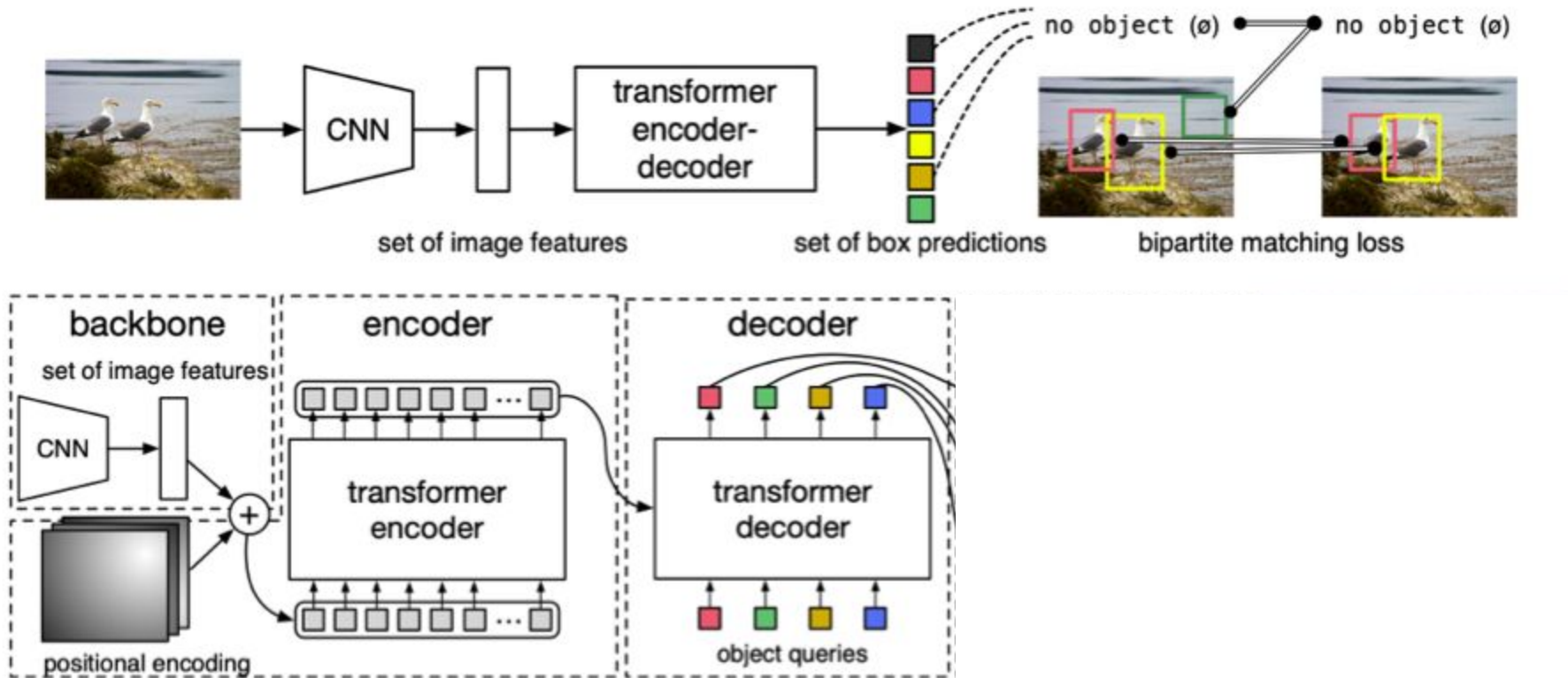
## Object Detection with Transformers: DETR



Carion et al, "End-to-End Object Detection with Transformers", ECCV 2020

# Convolution 🤝 Transformers

## Object Detection with Transformers: DETR

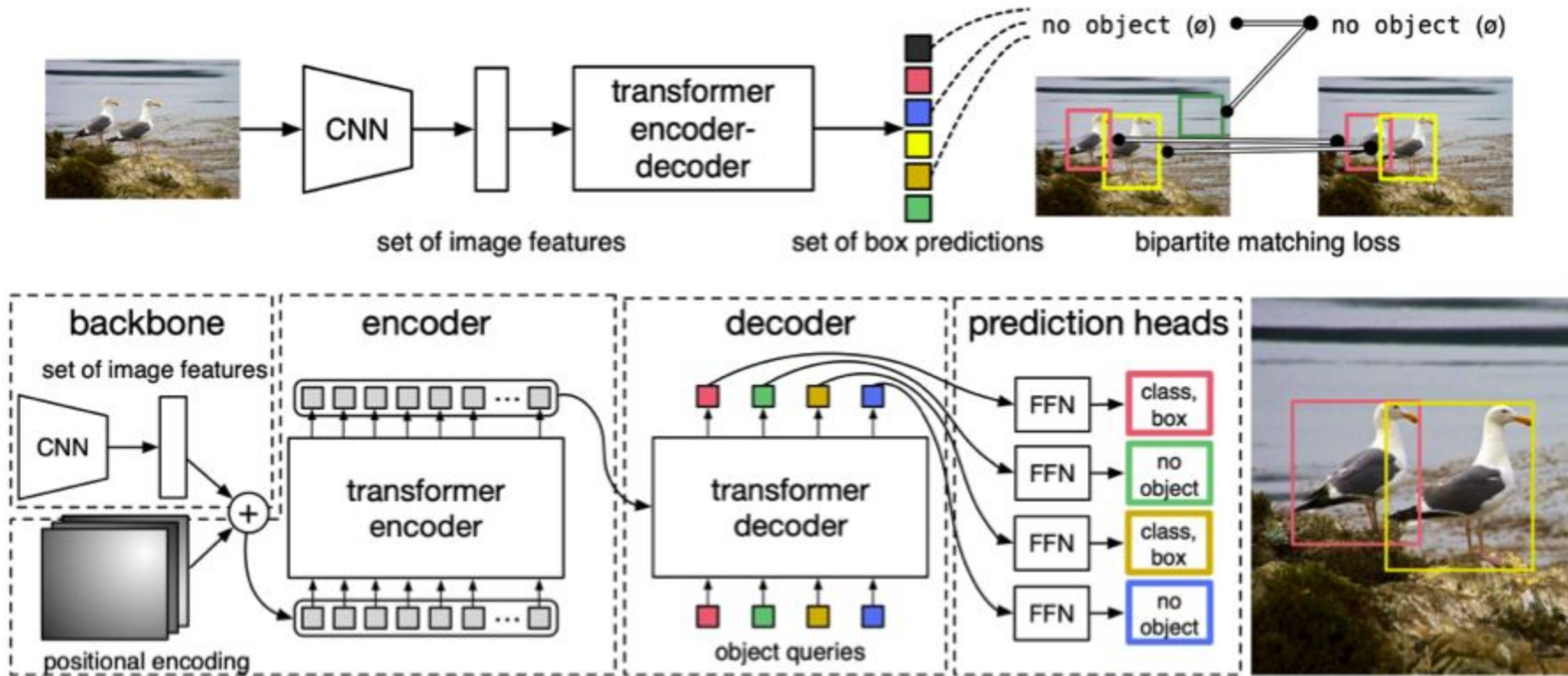


Carion et al, "End-to-End Object Detection with Transformers", ECCV 2020



# Convolution 🤝 Transformers

## Object Detection with Transformers: DETR



Carion et al, "End-to-End Object Detection with Transformers", ECCV 2020



Transformer architectures are now the basic backbone architectures for most vision, audio, and language **generative** and **discriminative** models.

- Transformer architectures are now the basic backbone architectures for most vision, audio, and language **generative** and **discriminative** models.
- Allows interactions of multiple modalities through cross-attention.
  - Eg: In text to image generative models,
    - Cross-attention
      - Key = Value = image tokens
      - Query = text prompt (text tokens)