

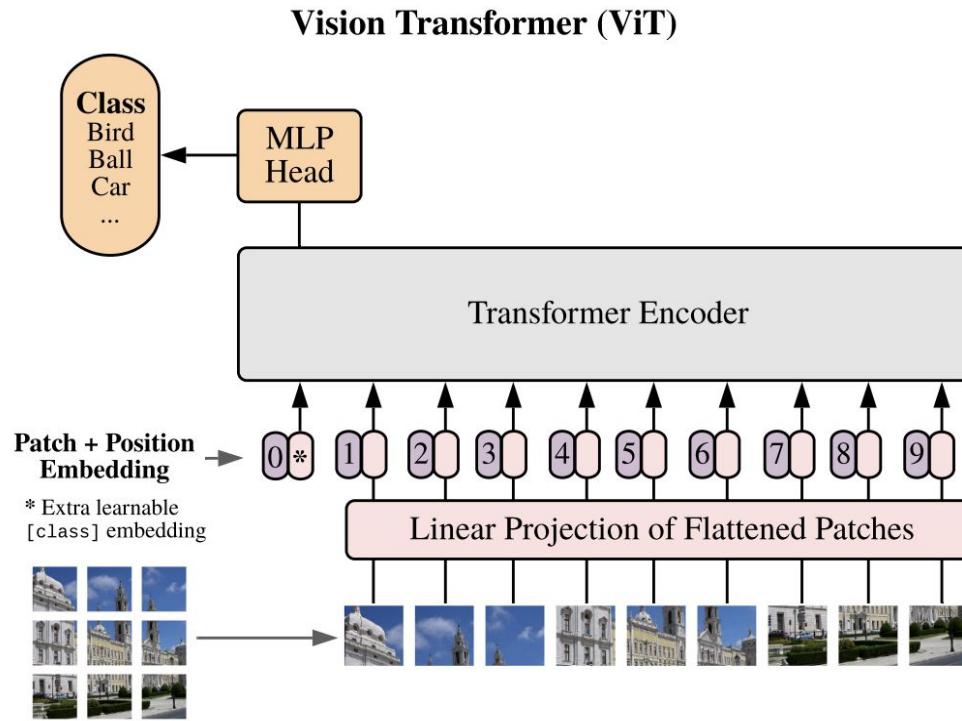
Announcements

- Challenge released.
- No laptops during class.
- Special access code to TorchStack: QYHBEGH
 - <https://torchstack.dev/>

Last week

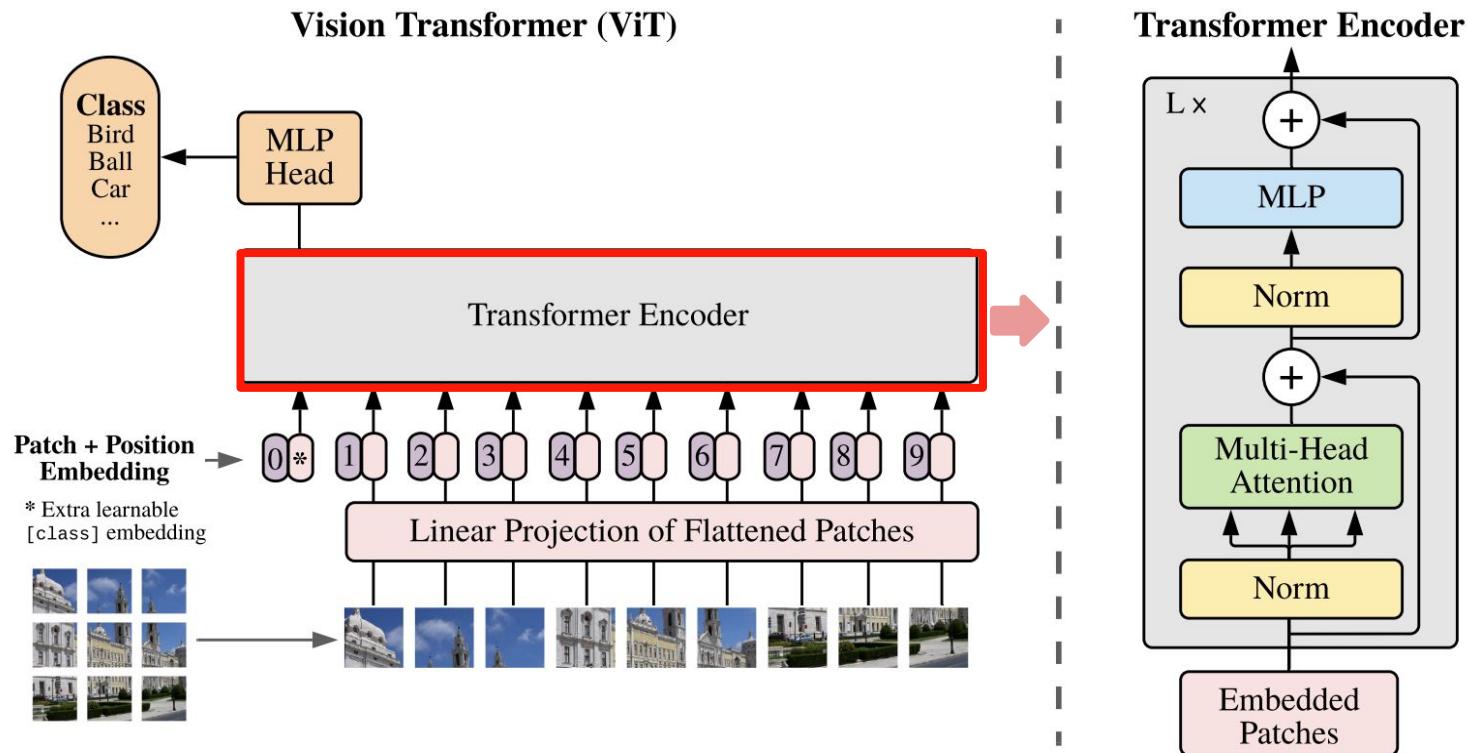
- Scaling up of transformers
- Extending transformers to computer vision

Recall: Transformers for vision



AN IMAGE IS WORTH 16X16 WORDS:
TRANSFORMERS FOR IMAGE RECOGNITION
AT SCALE, ICLR 21

Recall: Transformers for vision



AN IMAGE IS WORTH 16X16 WORDS:
TRANSFORMERS FOR IMAGE RECOGNITION
AT SCALE, ICLR 21

Today

- Positional embeddings in transformers
- Model interpretability

Today

- **Positional embeddings in transformers**
- Model interpretability

Self-Attention Layer

One query per input vector

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $Q = XW_Q$

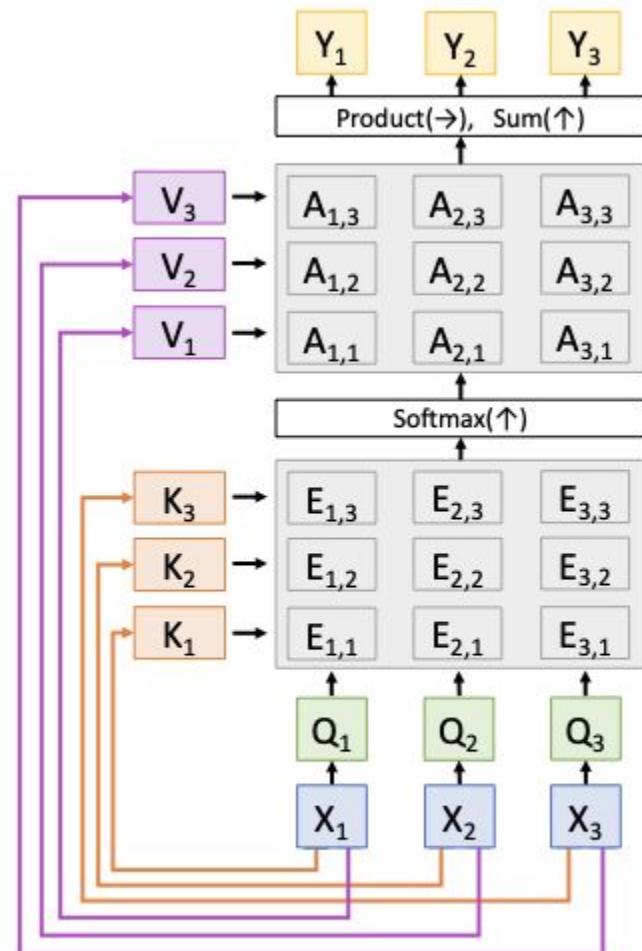
Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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



Self-Attention Layer

Consider permuting
the input vectors:

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $Q = XW_Q$

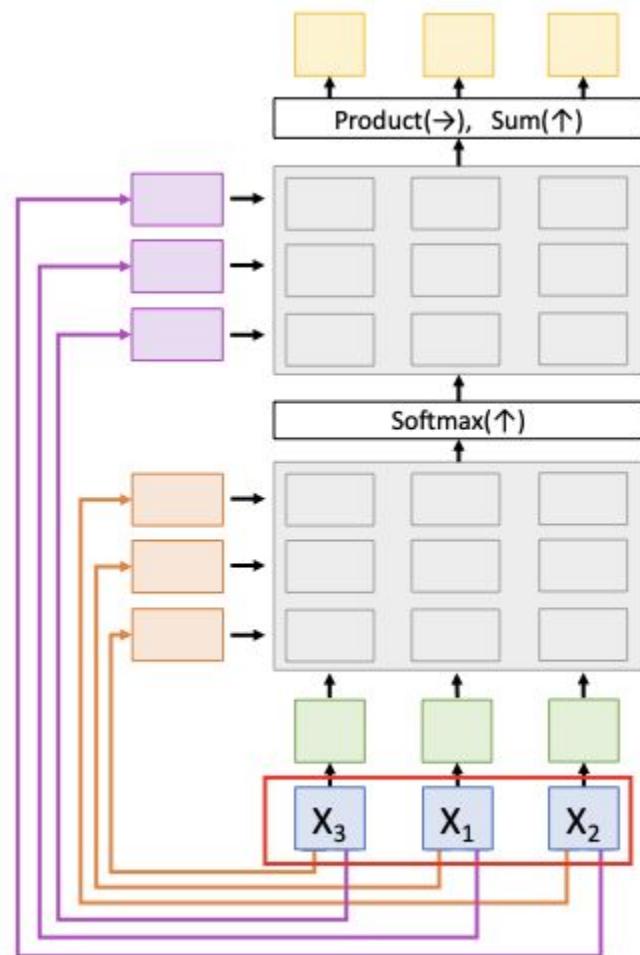
Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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



Self-Attention Layer

Inputs:

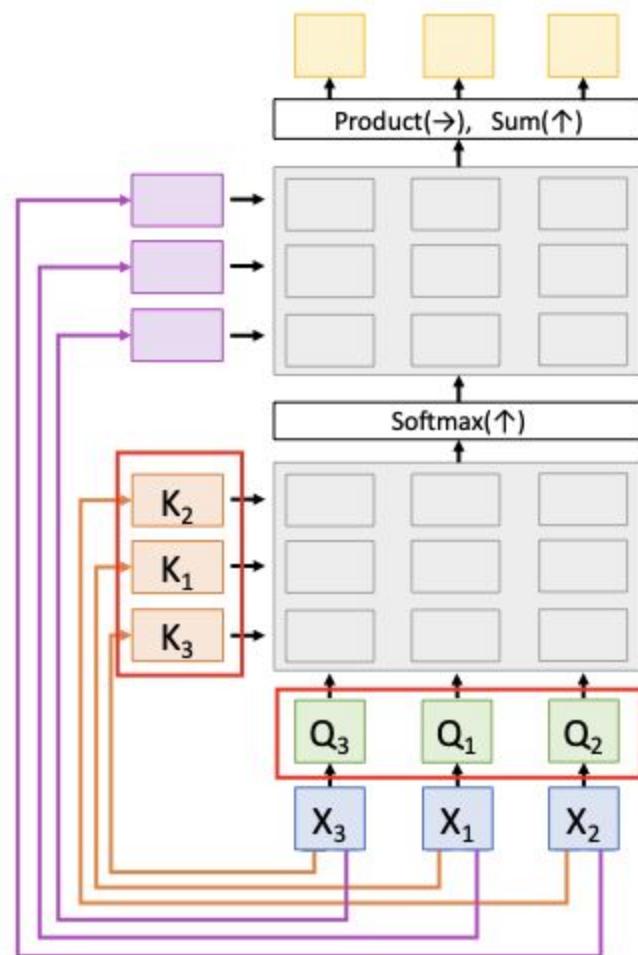
Input vectors: X (Shape: $N_x \times D_x$)
Key matrix: W_K (Shape: $D_x \times D_Q$)
Value matrix: W_V (Shape: $D_x \times D_V$)
Query matrix: W_Q (Shape: $D_x \times D_Q$)

Consider **permuting** the input vectors:

Queries and Keys will be the same, but permuted

Computation:

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



Self-Attention Layer

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Consider permuting
the input vectors:

Similarities will be the
same, but permuted

Computation:

Query vectors: $Q = XW_Q$

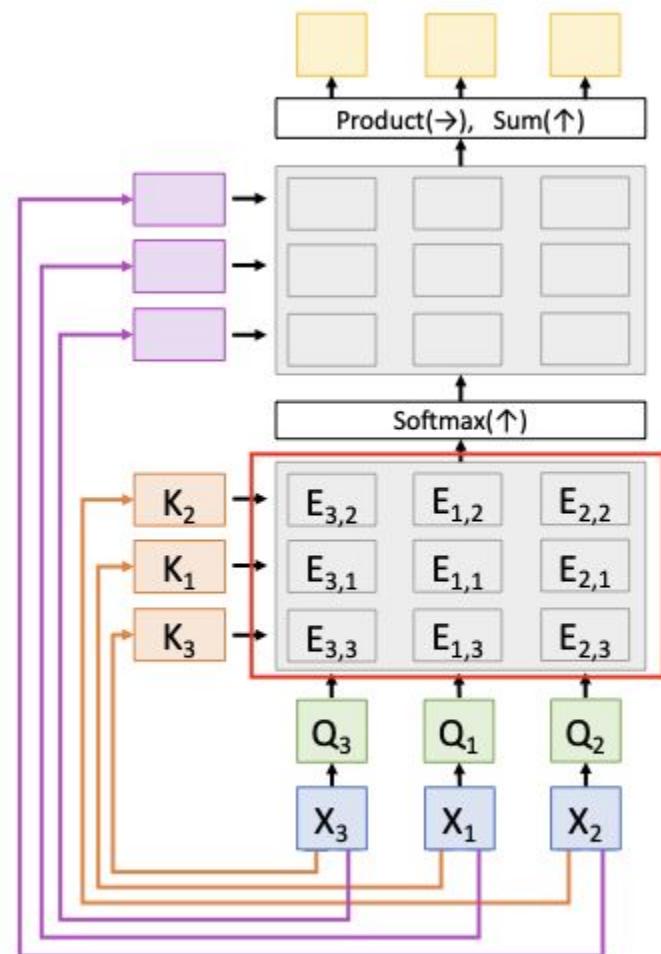
Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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



Self-Attention Layer

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Consider **permuting**
the input vectors:

Attention weights will be
the same, but permuted

Computation:

Query vectors: $Q = XW_Q$

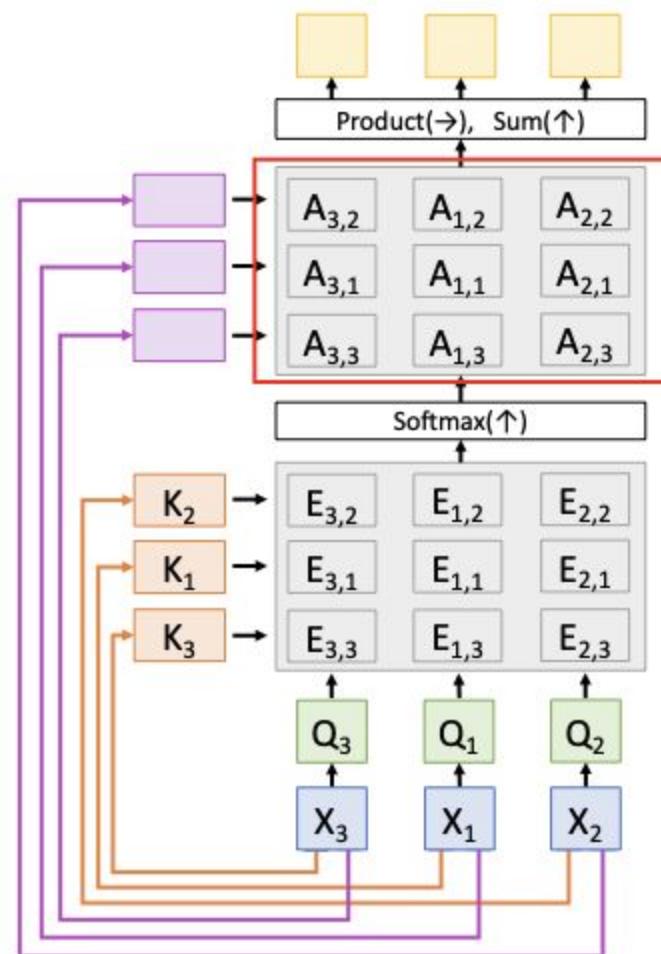
Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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



Self-Attention Layer

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Consider permuting
the input vectors:

Values will be the
same, but permuted

Computation:

Query vectors: $Q = XW_Q$

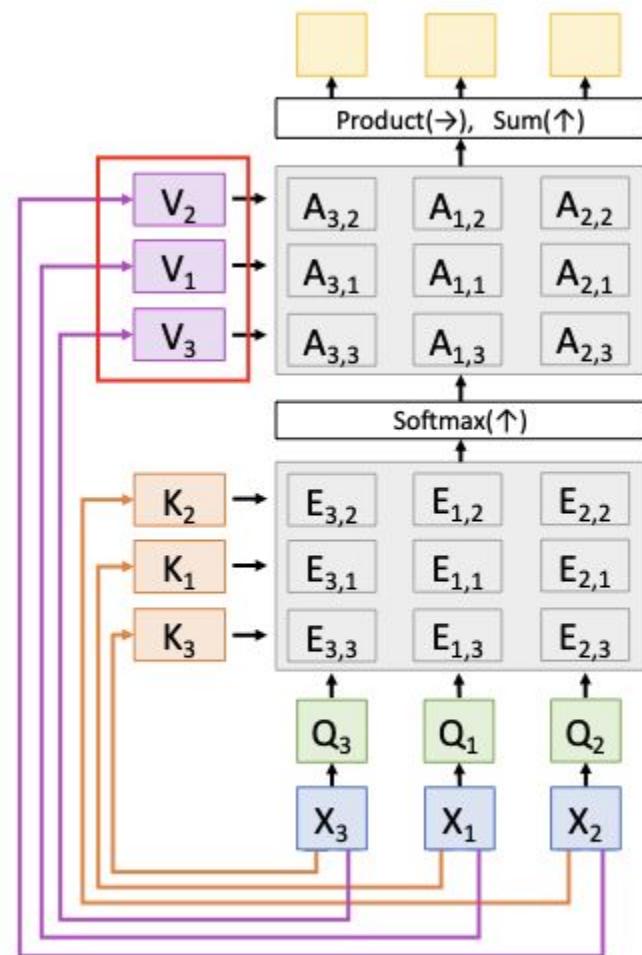
Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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



Self-Attention Layer

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Consider **permuting** the input vectors:

Outputs will be the same, but permuted

Computation:

Query vectors: $Q = XW_Q$

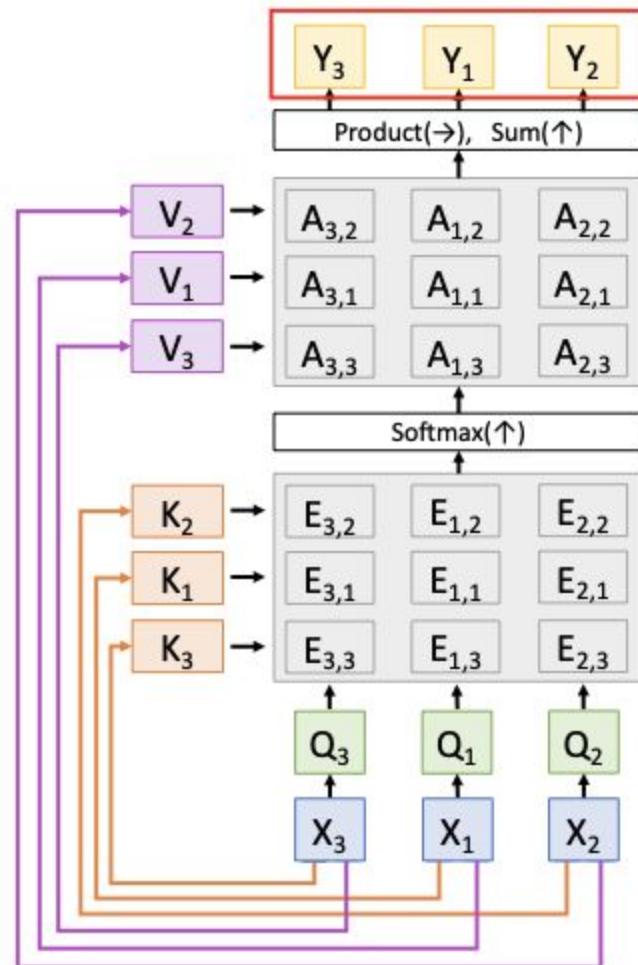
Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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



Self-Attention Layer

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $Q = XW_Q$

Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

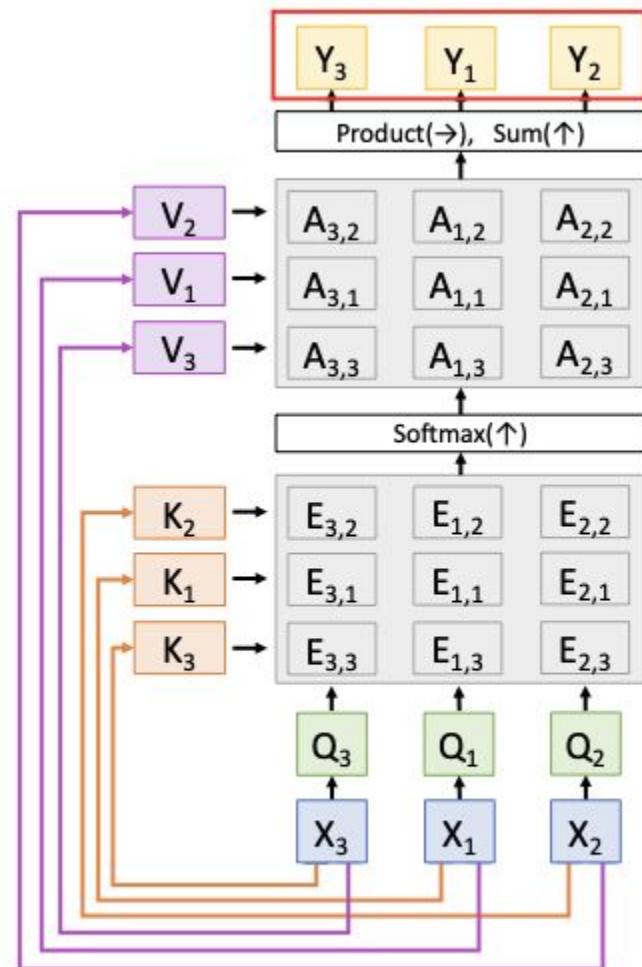
Output vectors: $Y = AV$ (Shape: $N_x \times D_V$) $Y_i = \sum_j A_{i,j} V_j$

Consider **permuting** the input vectors:

Outputs will be the same, but permuted

Self-attention layer is **Permutation Equivariant**
 $f(s(x)) = s(f(x))$

Self-Attention layer works on **sets** of vectors



Self-Attention Layer

Self attention doesn't
“know” the order of the
vectors it is processing!

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $Q = XW_Q$

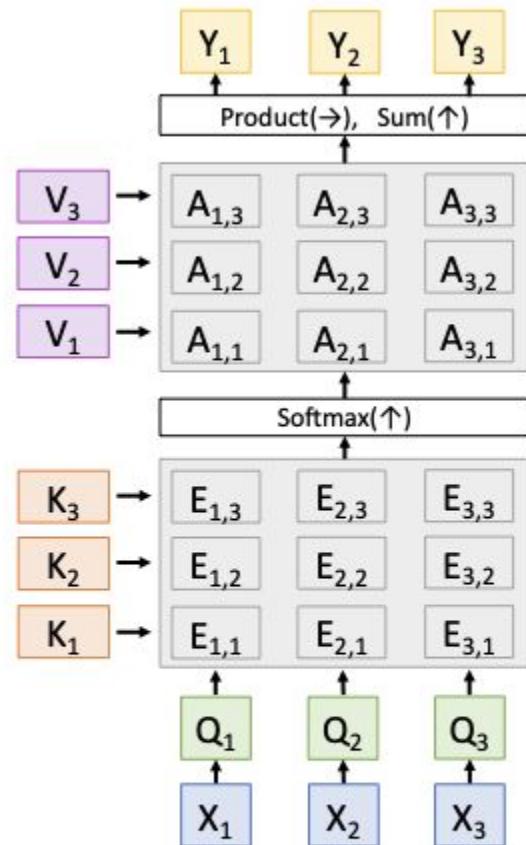
Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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



Images...



One ordering



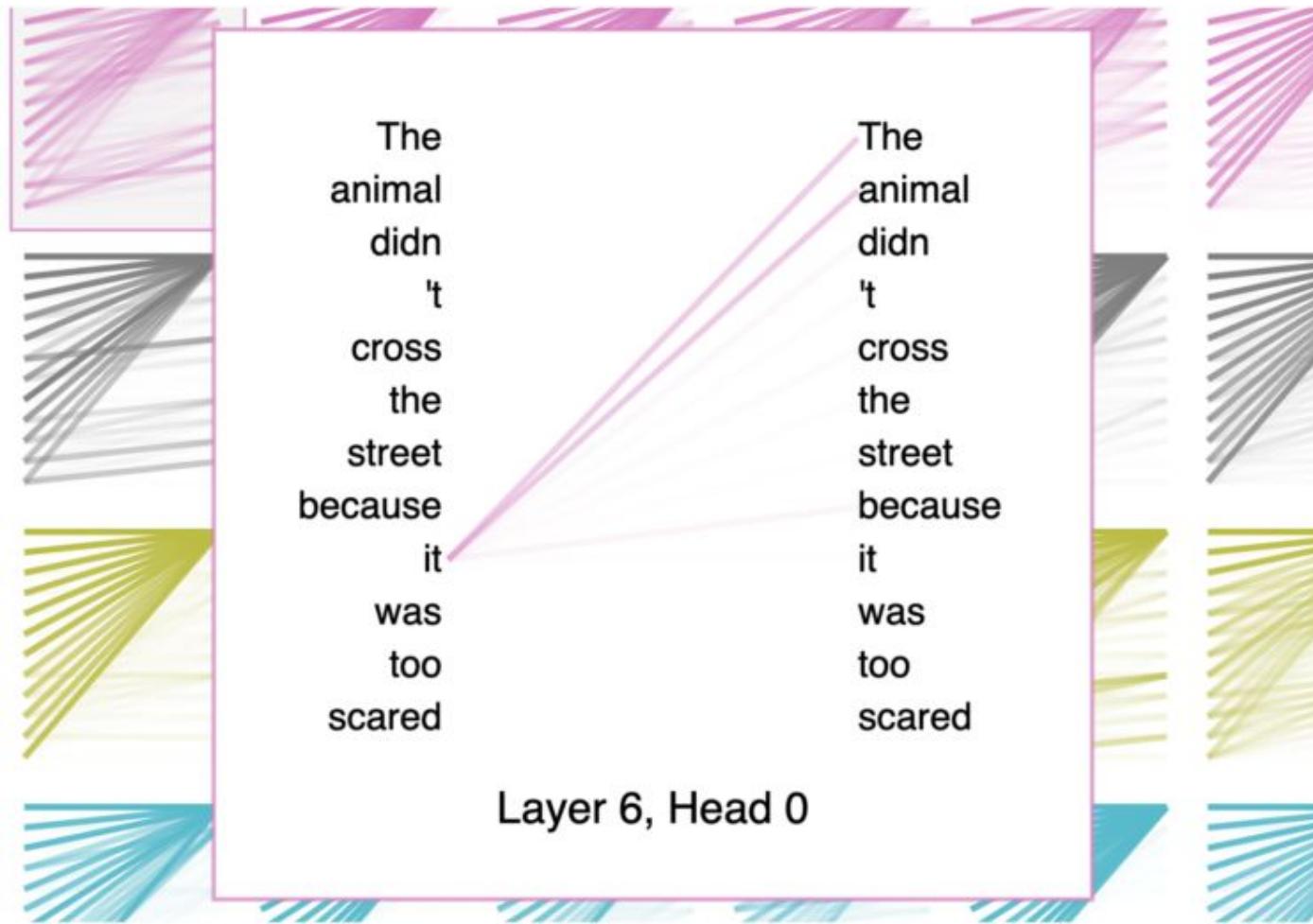
Shuffled



- Would the model know that the input is shuffled? No



Should the model know the order of the input tokens? Why or why not?



Input sentence: The animal didn't cross the street because it was too scared.

Shuffled input: didn't street it it too The cross animal the because scared too

Summary so far

- *Ordering of the sequence captures inherent semantic information in text, audio, image, and video.*
- *It is crucial to preserve the ordering.*

Self-Attention Layer

Self attention doesn't
“know” the order of the
vectors it is processing!

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Computation:

Query vectors: $Q = XW_Q$

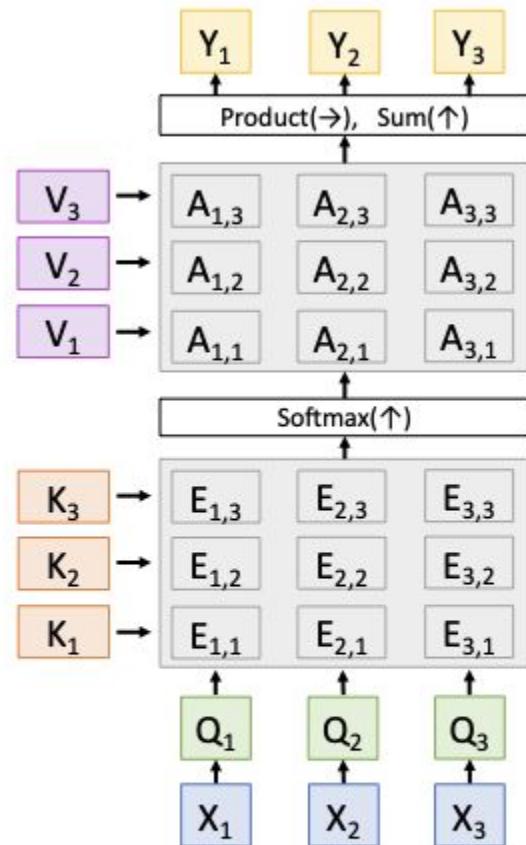
Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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



Self-Attention Layer

Inputs:

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

Key matrix: W_K (Shape: $D_x \times D_Q$)

Value matrix: W_V (Shape: $D_x \times D_V$)

Query matrix: W_Q (Shape: $D_x \times D_Q$)

Self attention doesn't
“know” the order of the
vectors it is processing!

In order to make
processing position-
aware, concatenate or
add **positional encoding**
to the input

Computation:

Query vectors: $Q = XW_Q$

E can be learned lookup
table, or fixed function

Key vectors: $K = XW_K$ (Shape: $N_x \times D_Q$)

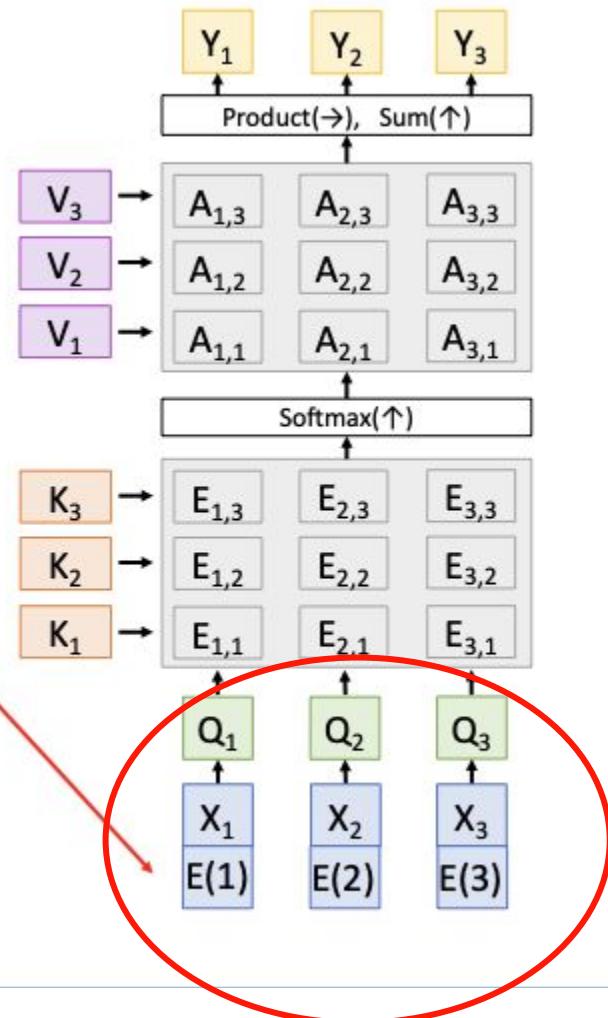
(Shape: $N_x \times N_x$)

Value Vectors: $V = XW_V$ (Shape: $N_x \times D_V$)

Similarities: $E = QK^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (Q_i \cdot K_j) / \sqrt{D_Q}$

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

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

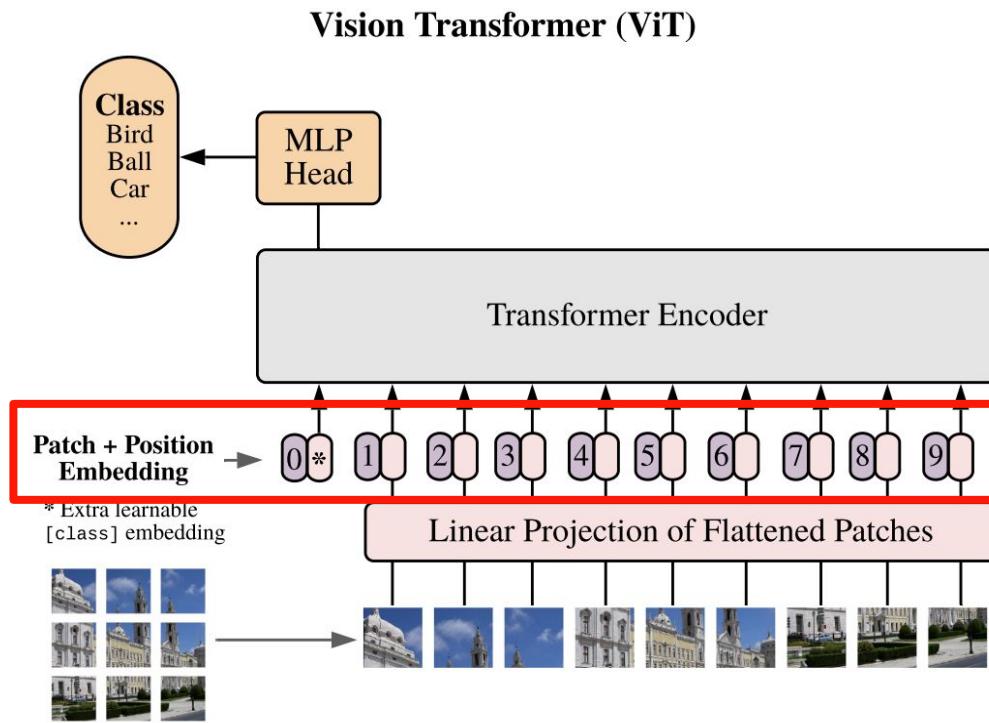


Positional encodings == identifiers that preserve the ordering.

Summary so far

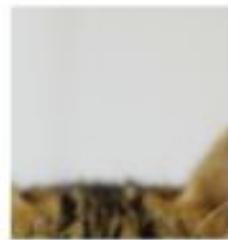
- *Ordering of the sequence captures inherent semantic information in text, audio, image, and video.*
- *It is crucial to preserve the ordering.*
- ***Positional embeddings: Help remember the order of the input.***

How to incorporate positional embedding?



AN IMAGE IS WORTH 16X16 WORDS:
TRANSFORMERS FOR IMAGE RECOGNITION
AT SCALE, ICLR 21

Add to the input token embeddings



0.33	0.71	0.91	0.23	0.15
------	------	------	------	------

+

0.01	0.01	0.01	0.01	0.01
------	------	------	------	------

0.12	0.22	0.31	0.03	0.10
------	------	------	------	------

+

0.02	0.21	0.33	0.43	0.98
------	------	------	------	------

+

0.34	0.70	0.95	0.21	0.17
------	------	------	------	------

0.03	0.32	0.85	0.91	0.24
------	------	------	------	------

Add to the input token embeddings

- Retain the feature dimensionality.



0.33	0.71	0.91	0.23	0.15
------	------	------	------	------

+

0.01	0.01	0.01	0.01	0.01
------	------	------	------	------

0.12	0.22	0.31	0.03	0.10
------	------	------	------	------

+

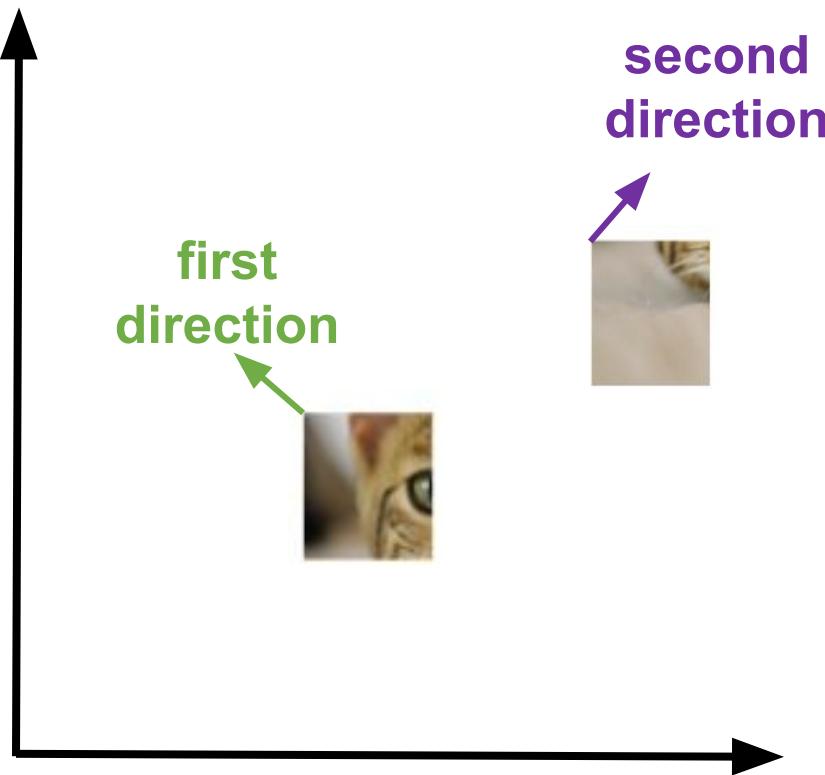
0.02	0.21	0.33	0.43	0.98
------	------	------	------	------

+

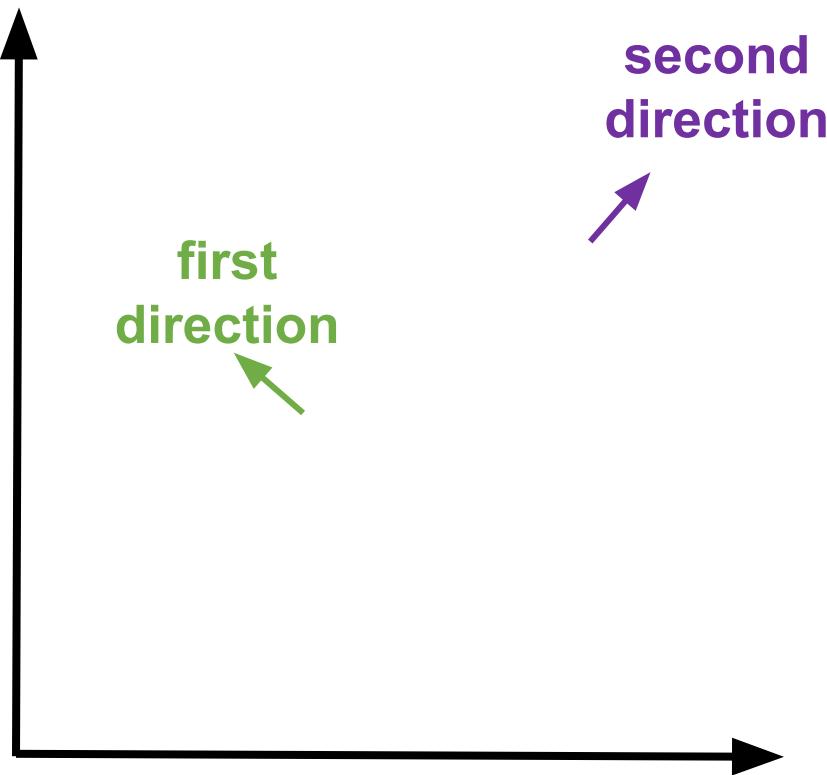
0.34	0.70	0.95	0.21	0.17
------	------	------	------	------

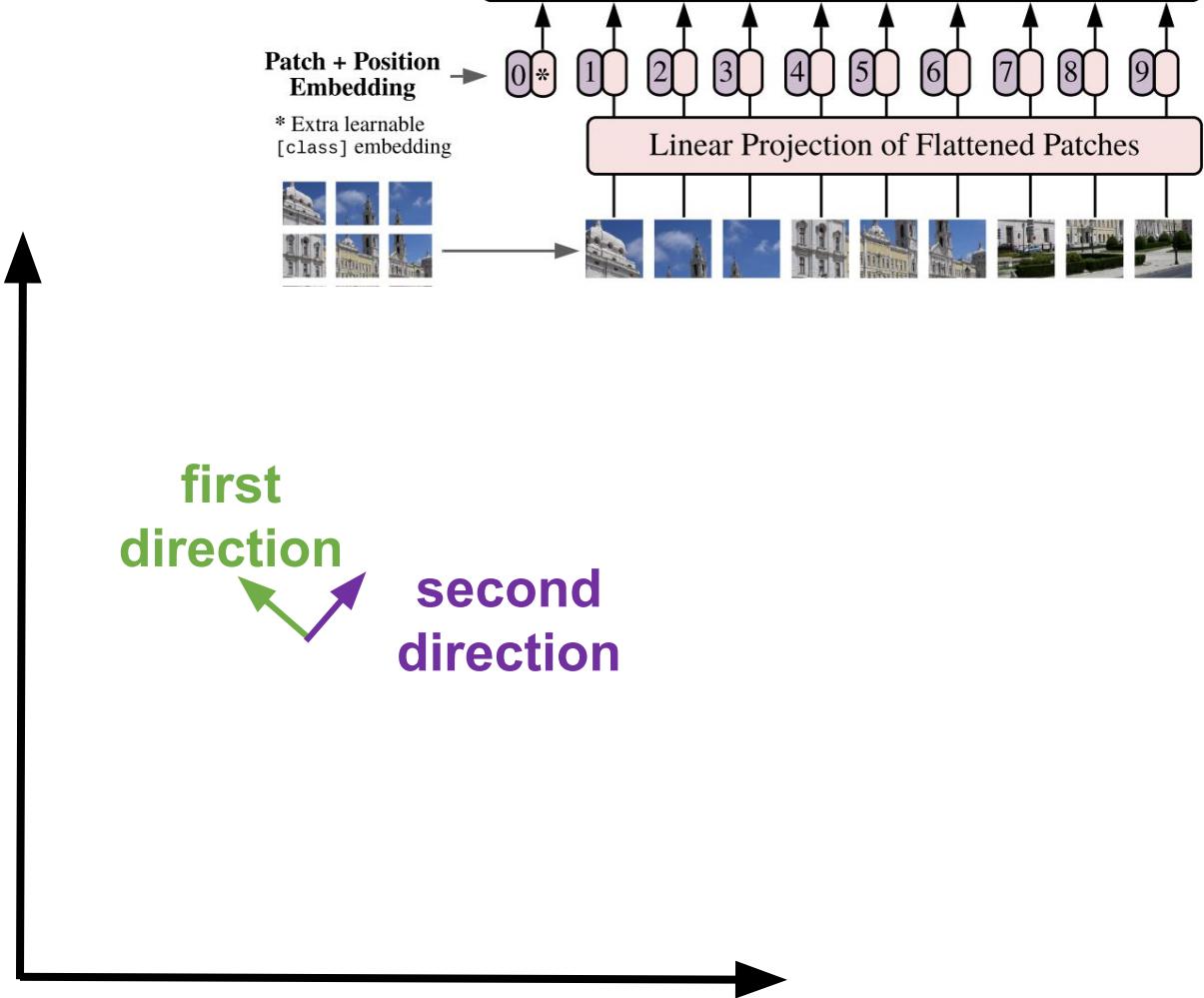
0.03	0.32	0.85	0.91	0.24
------	------	------	------	------

Intuition behind addition of positional encodings



Positional embeddings

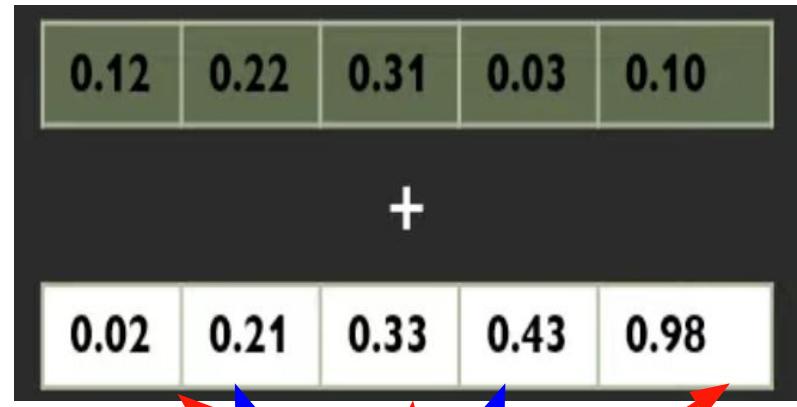




- Retain the dimensionality of input token.
- Not too high or too low of a magnitude.

Commonly used positional encodings

- Sinusoidal functions
- Cosine functions

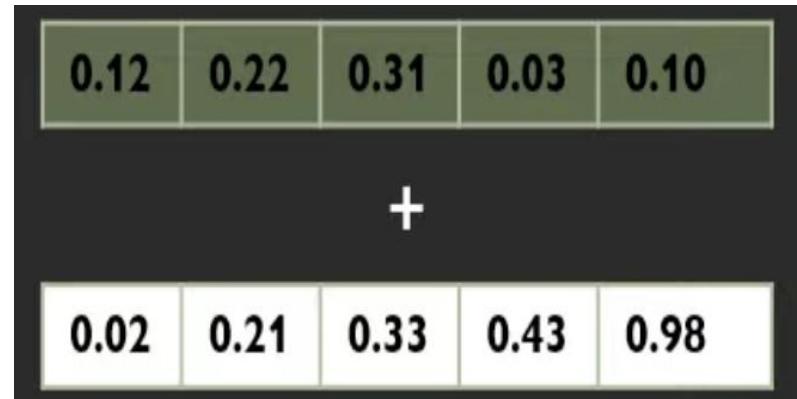


$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

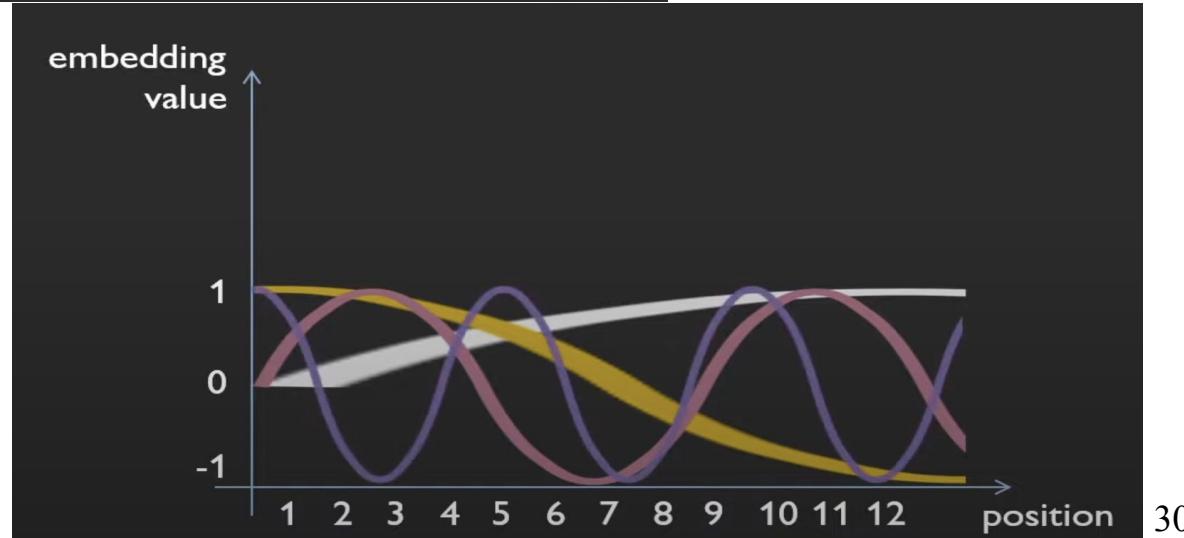
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

Commonly used positional encodings

- Sinusoidal functions
- Cosine functions



$$PE_{(pos, 2i)}$$
$$PE_{(pos, 2i+1)}$$





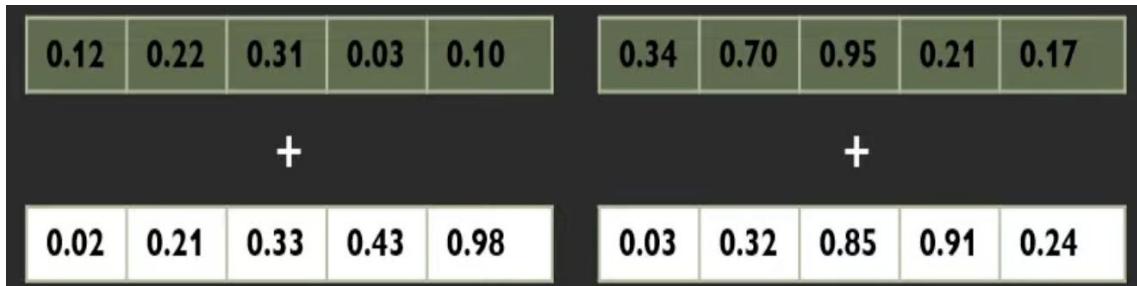
What are the upsides of using sinusoidal / cosine functions as positional embeddings?

- ⓘ Presenting with animations, GIFs or speaker notes? Enable our [Chrome extension](#)



Are there any downsides of using sine and cosine functions as positional embeddings?

Downsides of adding positional embeddings

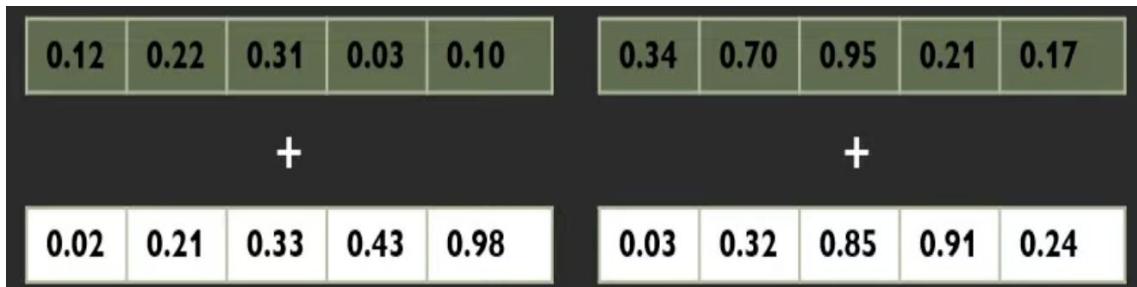


- ✓ Retain the dimensionality of input token.
- ✓ Not too high or too low of a magnitude.
- ✗ Mixing up semantic information with positional information

Summary so far

- ***Positional embeddings:*** *Help remember the order of the input.*
 - **Addition:**
 - Retains the feature dim
 - Mixes up semantic information with positional information

Downsides of adding positional embeddings



- ✓ Retain the dimensionality of input token.
- ✓ Not too high or too low of a magnitude.
- ✗ Mixing up semantic information with positional information



How about concatenation?

Concatenation of positional embeddings

0.12	0.22	0.31	0.03	0.10
------	------	------	------	------

+

v/s

0.12	0.22	0.31	0.03	0.10	0.02	0.21	0.33	0.43	0.98
------	------	------	------	------	------	------	------	------	------

0.02	0.21	0.33	0.43	0.98
------	------	------	------	------



**Which of the following is true about concatenating positional embeddings?
Select all that apply**

- ⓘ Presenting with animations, GIFs or speaker notes? Enable our [Chrome extension](#)



Which of the following is true about concatenating positional embeddings?
Select all that apply

Concatenation prevents mixing up semantic and positional information.

48%

Concatenation of positional embeddings leads to an increase in the memory and computational time.

54%

In order to concatenate positional embeddings, it is sufficient to reduce the token embeddings to half its size.

15%

All of the above

51%

Join at

slido.com
#2812 878

Concatenation of positional embeddings



Studied heavily in computer vision

Pros:

- Add an additional independent direction.
- Does not mix up semantic and positional information.

Cons:

- Increases compute / memory footprint.

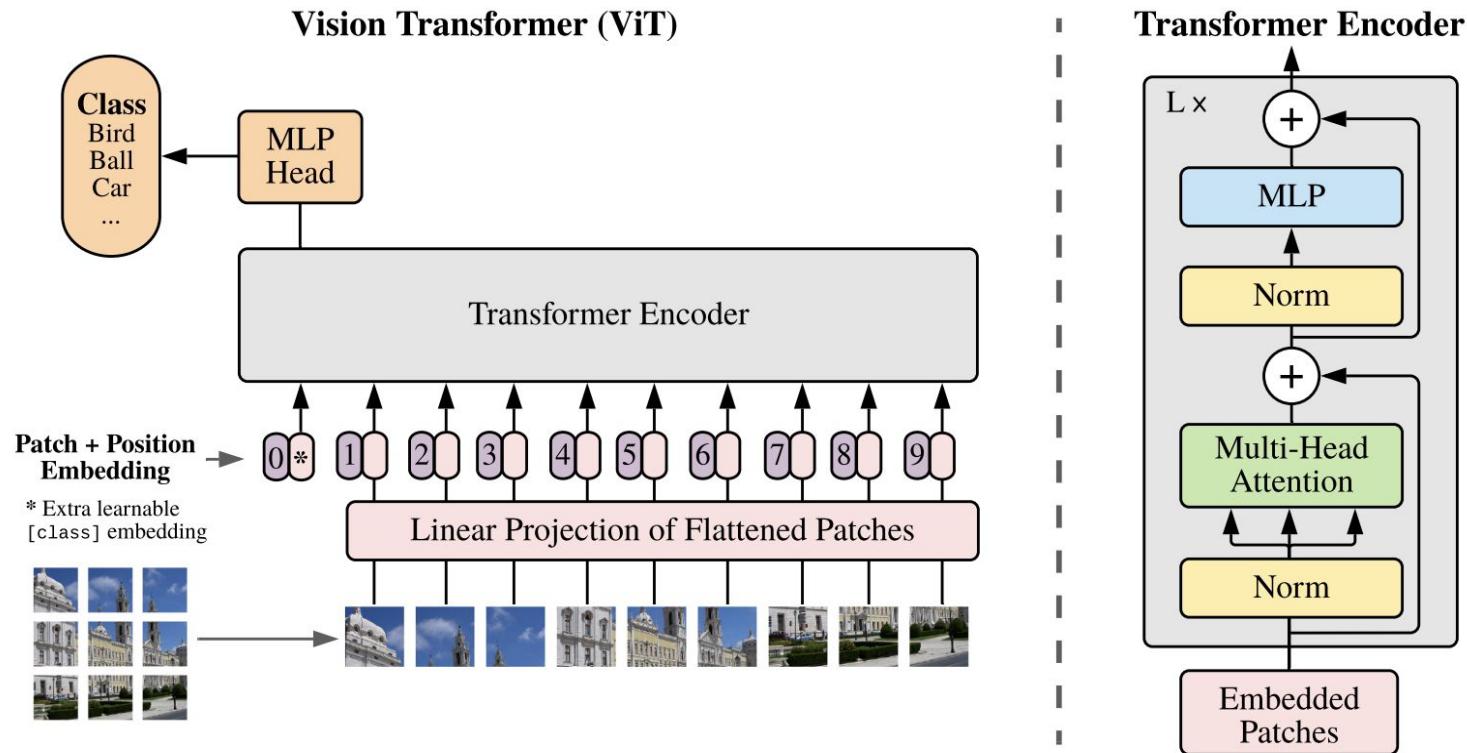
Summary so far

- ***Positional embeddings:*** *Help remember the order of the input.*
 - **Addition:**
 - Mixes up semantic information with positional information
 - **Concatenation:**
 - Does not mix up the signals
 - Increases compute / memory footprint.

Other options

- Learn your positional embeddings along with the task.
- Relative positional embeddings
 - based on semantic distance between the tokens
- RoPE
 - embedding + rotated variant of the embedding.
 -
 -
 -
 -
- A very active area of research!

Summary



- ✓ Very powerful architecture.
- ⚠️ Very data hungry - will not recommend unless there is sufficient data even for fine tuning tasks.

Today

- Positional embeddings in transformers
- **Model interpretability**

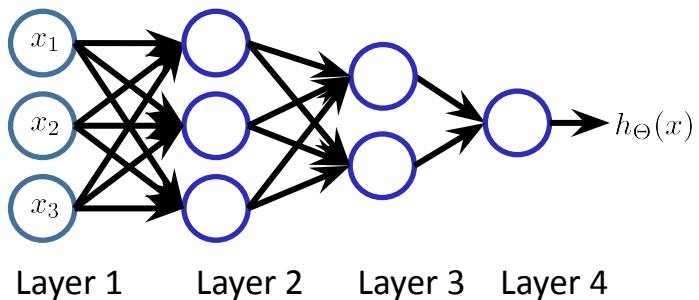
So far: Network architectures

Recurrent

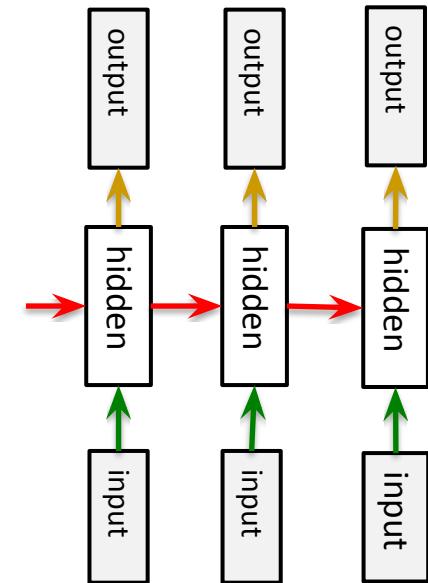
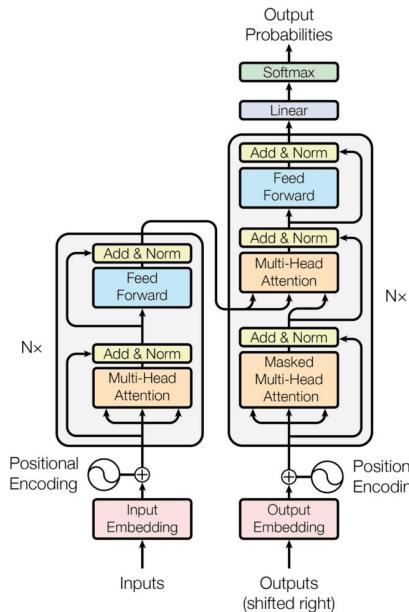
time □

Feed-forward

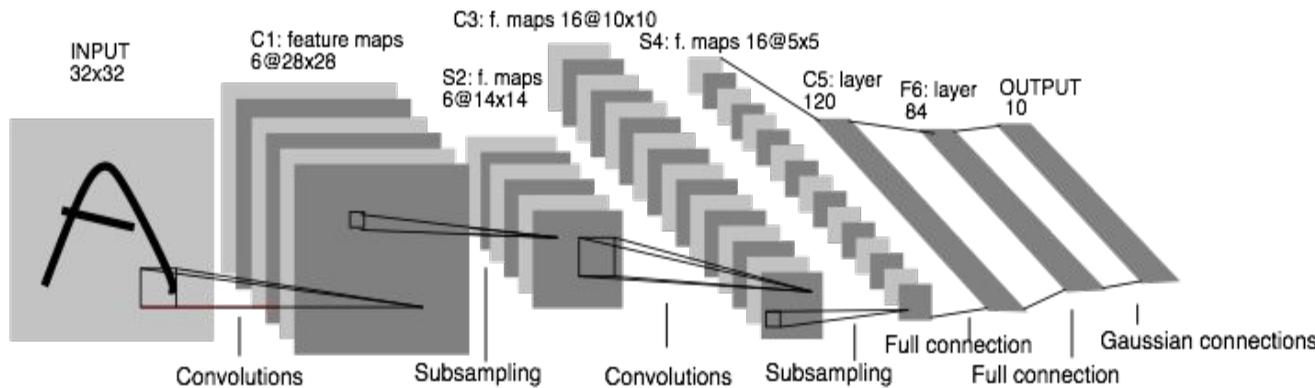
Fully connected



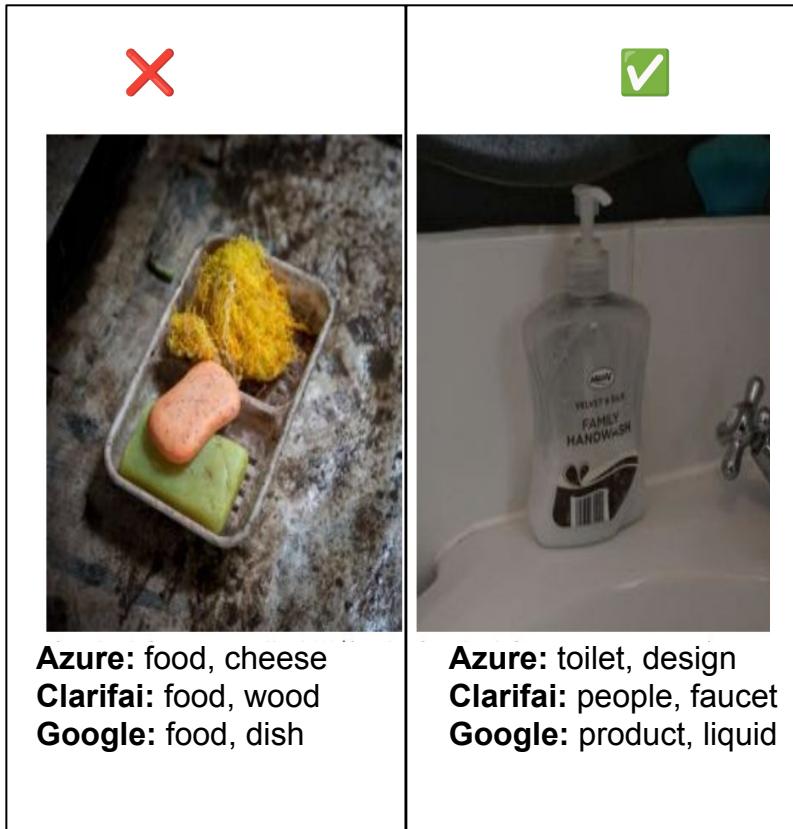
Transformer



Convolutional



Why do we need model understanding?



- Model debugging
- Establish user trust

Why do we need model understanding?

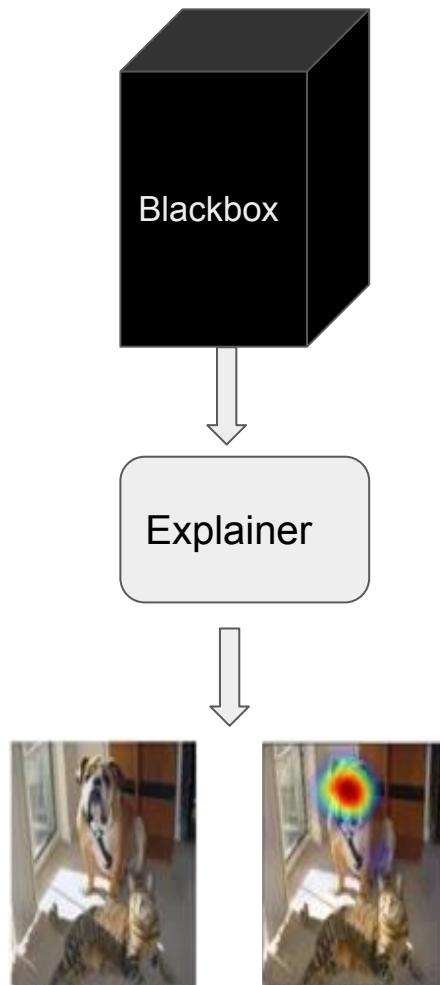


- *What makes this image look like a “lion” to the model?*
- *Does the lion cub contribute as much?*
- *What about the reflections?*

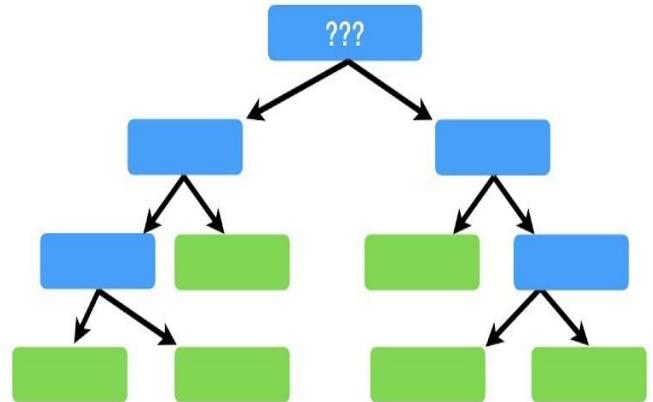
Build a deeper understanding of the model

How to achieve model understanding?

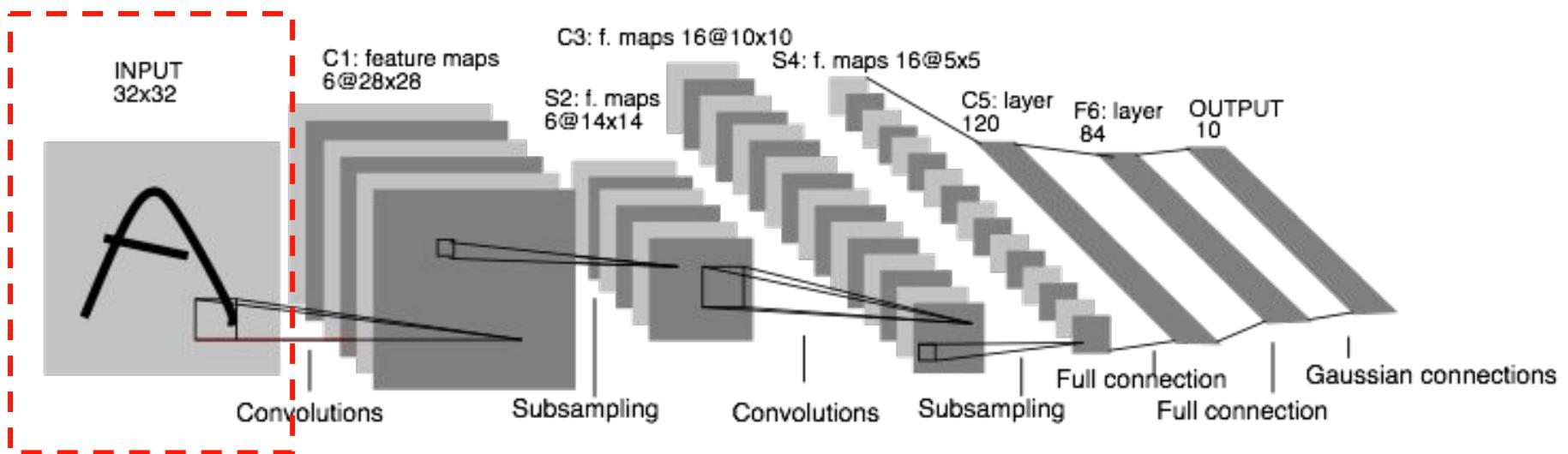
Posthoc explanations



Inherently interpretable models



End to end pipeline

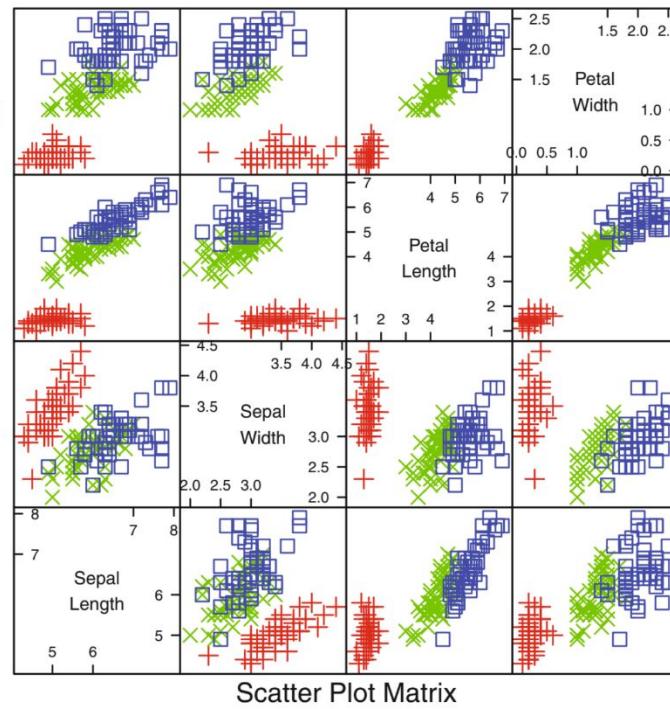
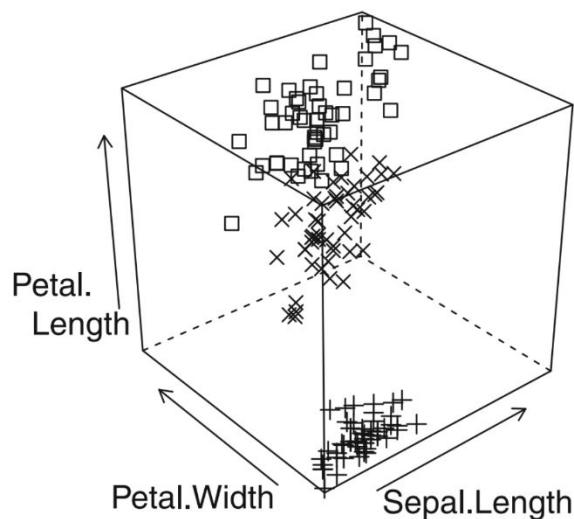


Recall: Visualizing features



Iris dataset-

- 4 features
 - Sepal length/width
 - Petal length/width
- 3 classes



Visualizing high dimensional data is hard!

Recall: Visualizing features

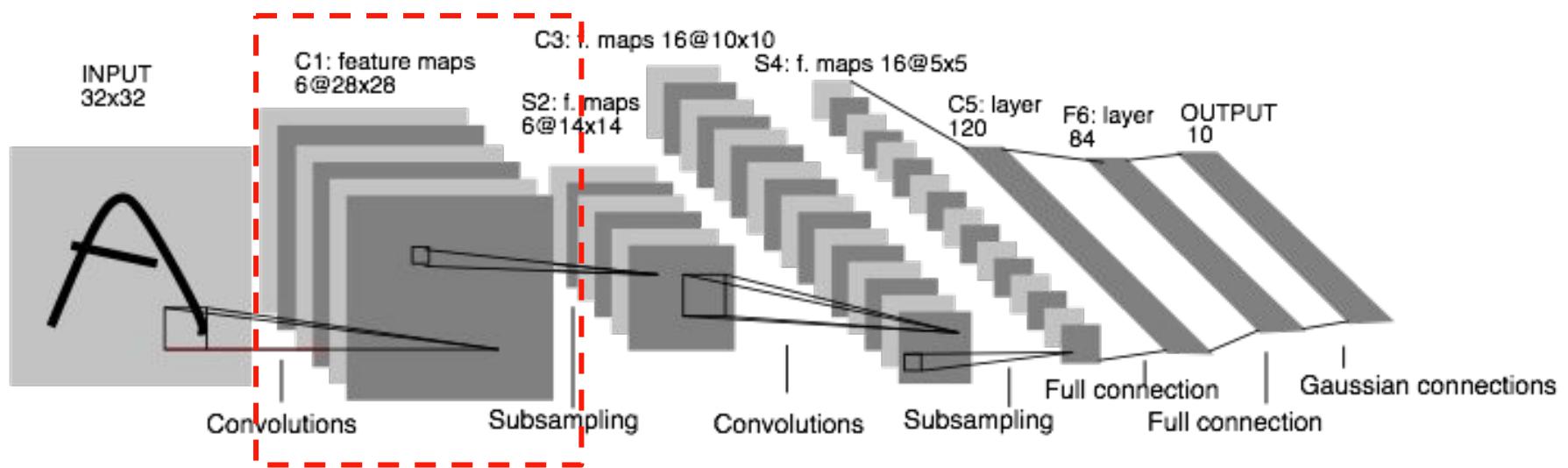
Goal: Projecting data to fewer dimensions

- Retain relative information about the feature higher dimensions
- Many project down to 2D

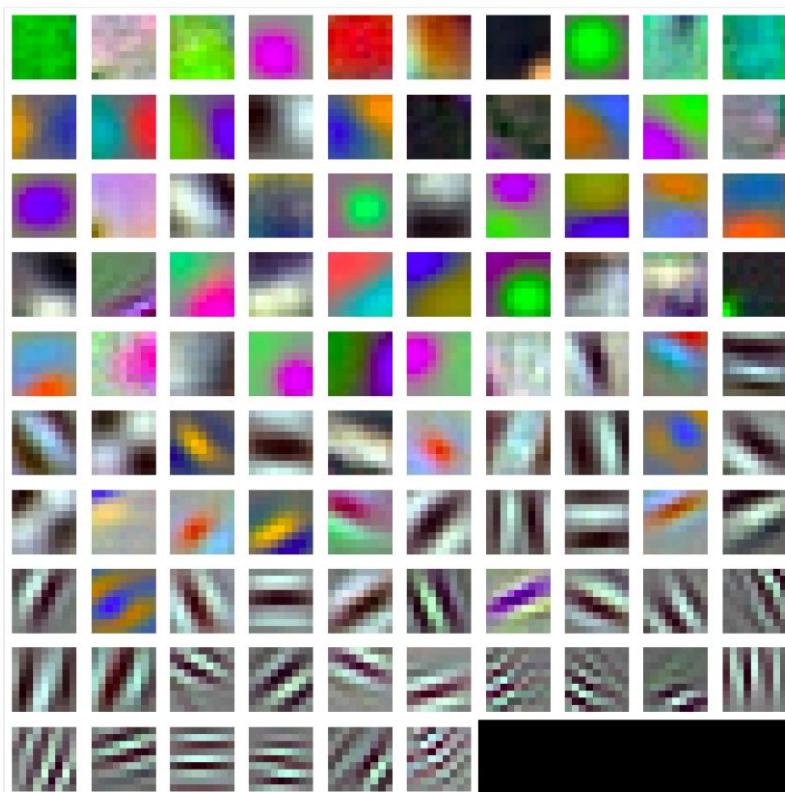
Potential Solutions:

- PCA
- Convert into a (sort-of) classification problem

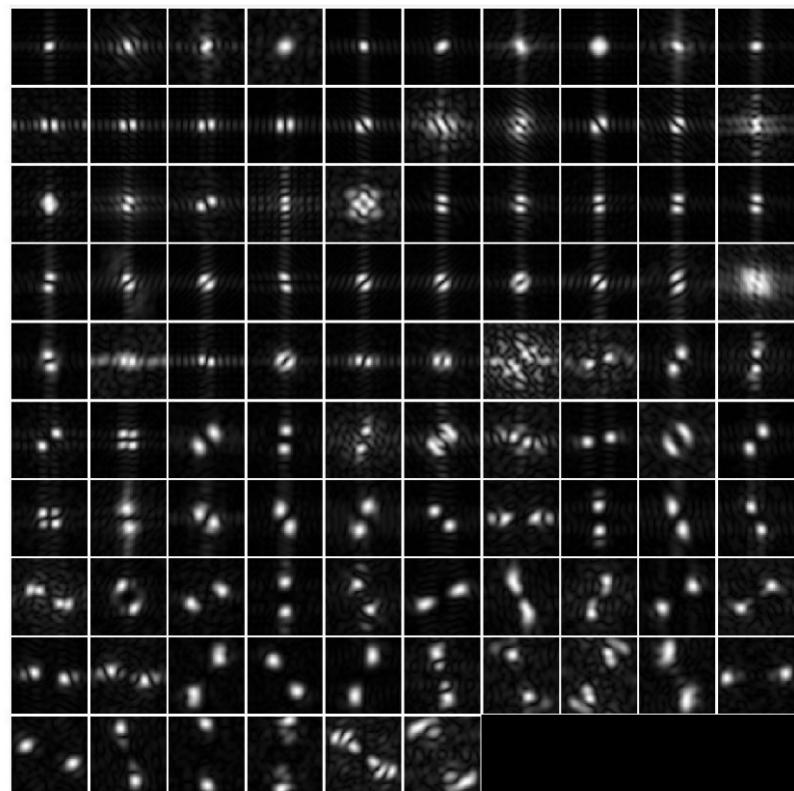
End to end pipeline

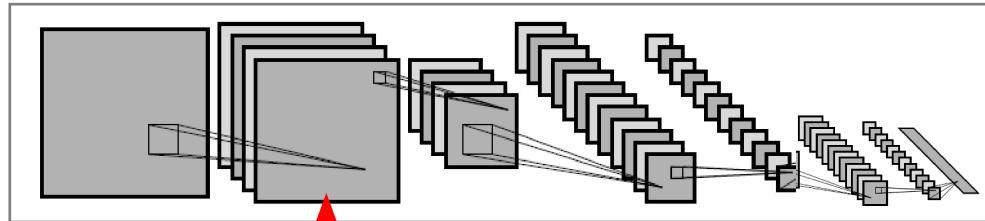


Convolutional filters



96 Units in conv1





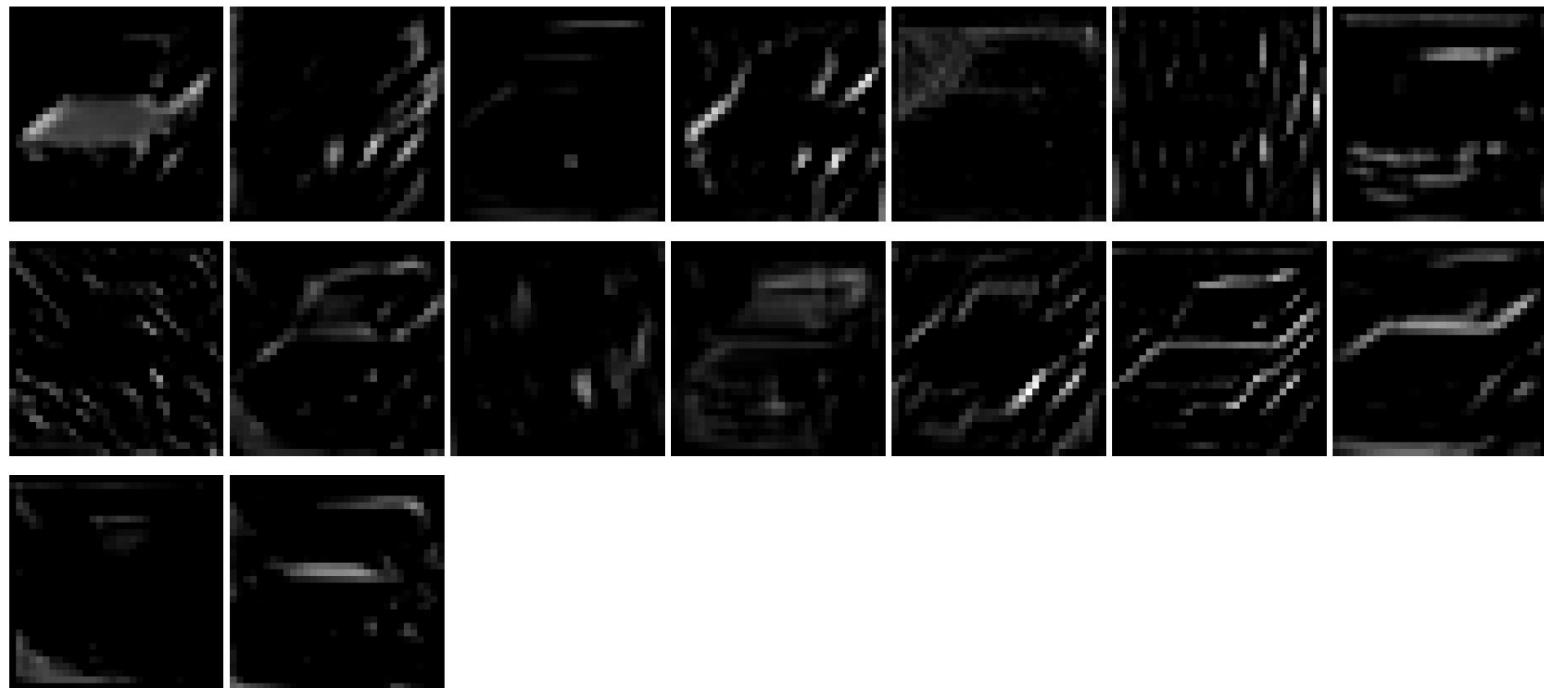
input (32x32x3)



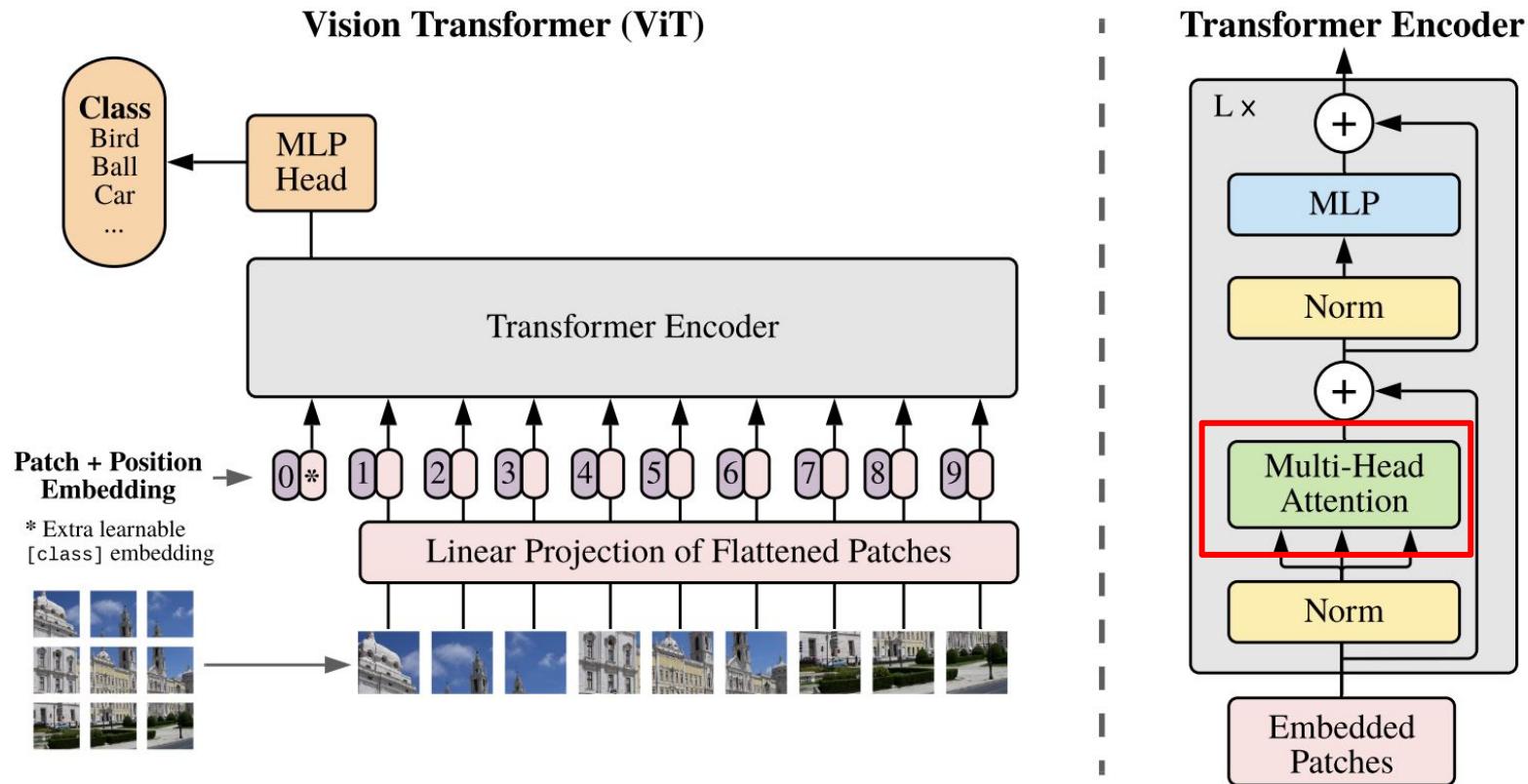
Convolution

conv (32x32x16) params: $16 \times 5 \times 5 \times 3 + 16 = 1216$

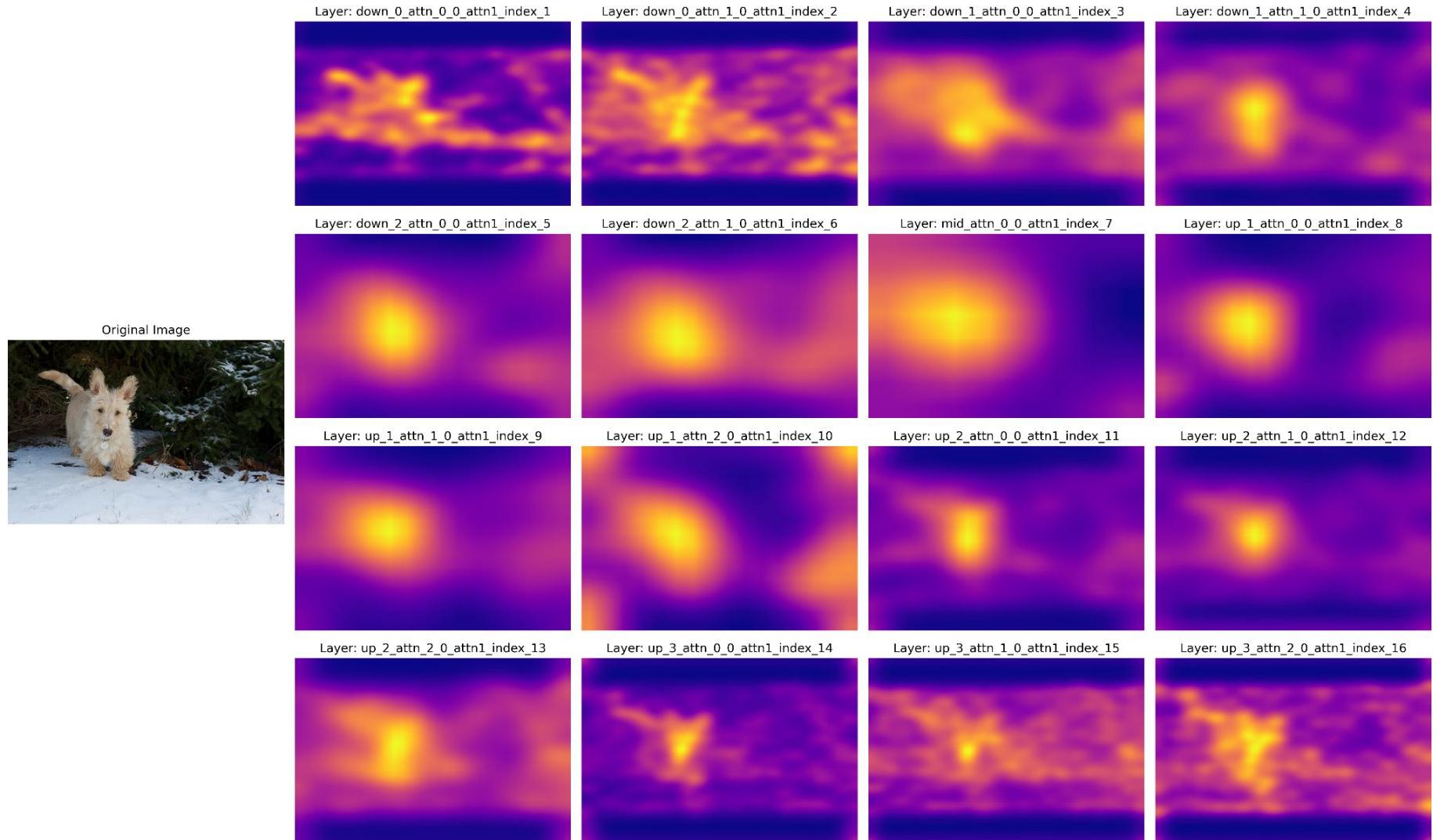
filter size 5x5x3, stride 1



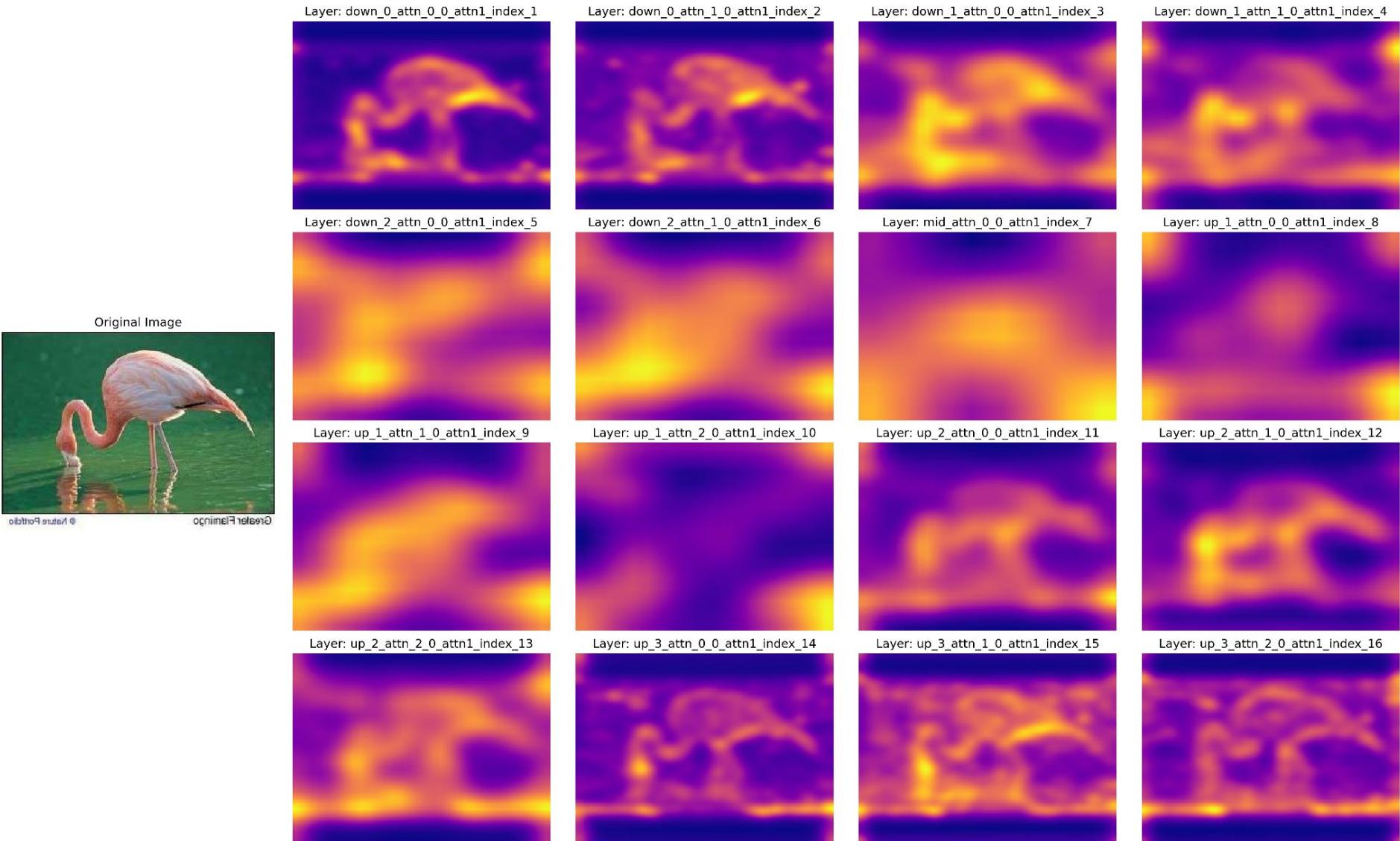
Recall: Attention in transformers



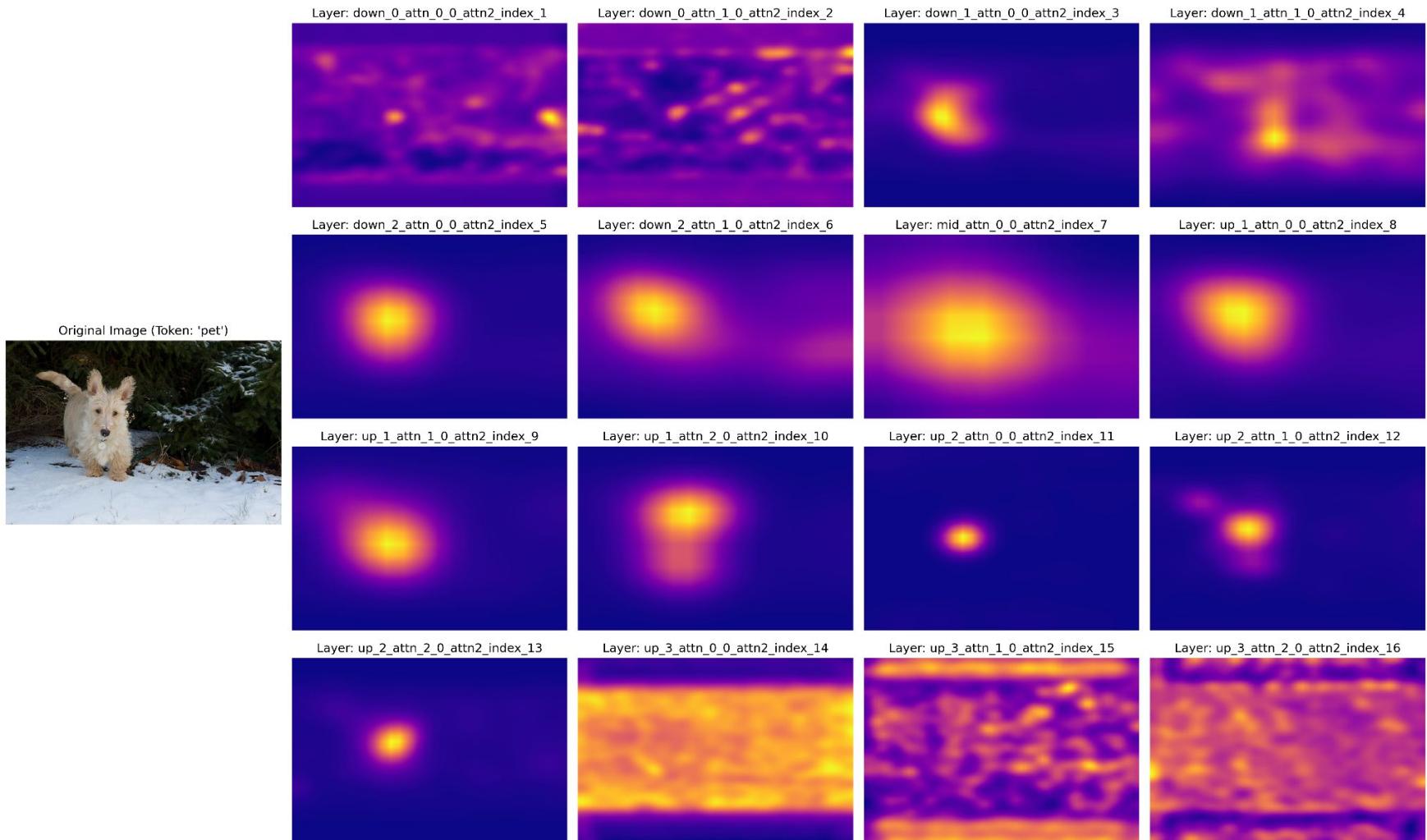
Self-Attention Map visualizations



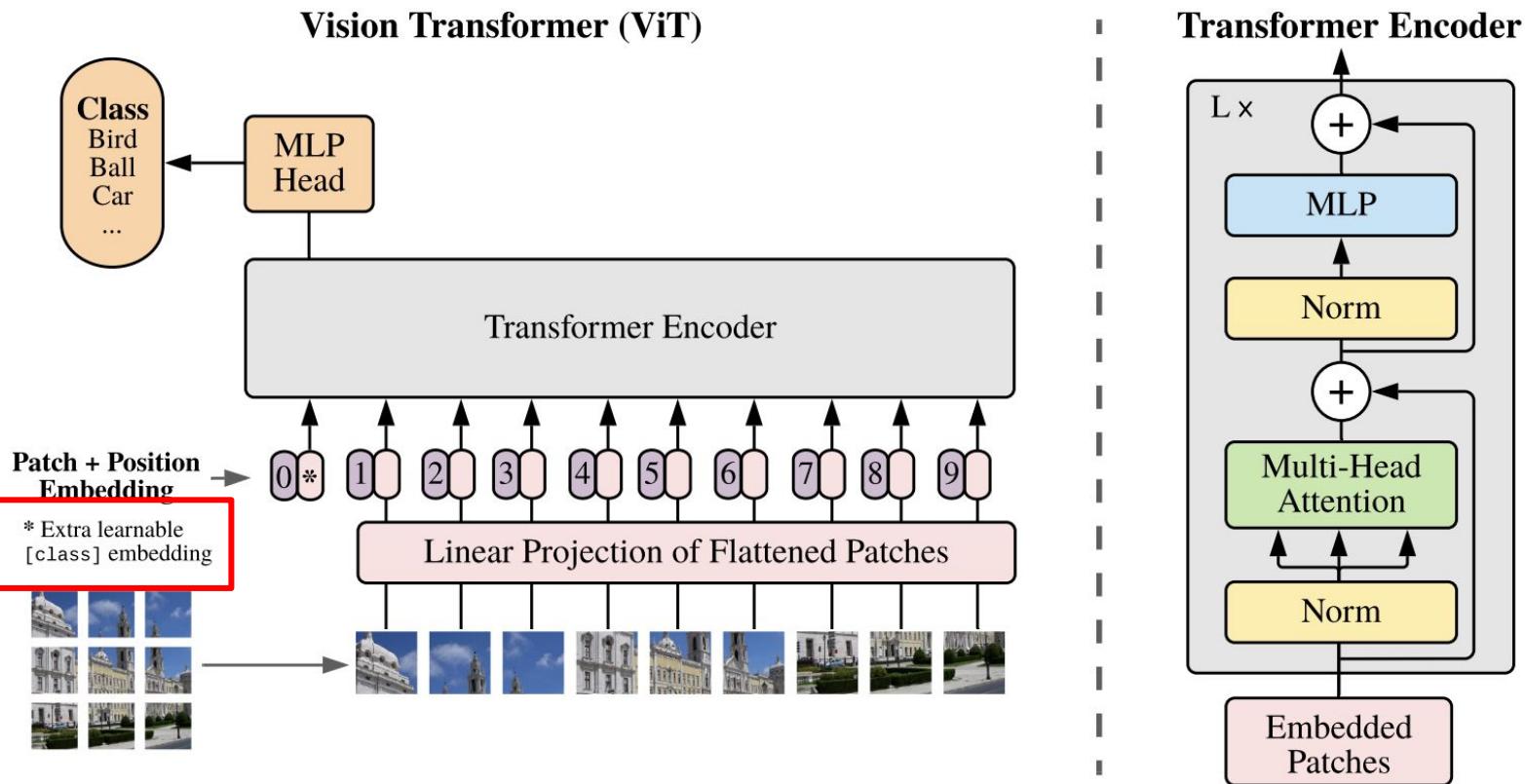
Self-Attention Map visualizations



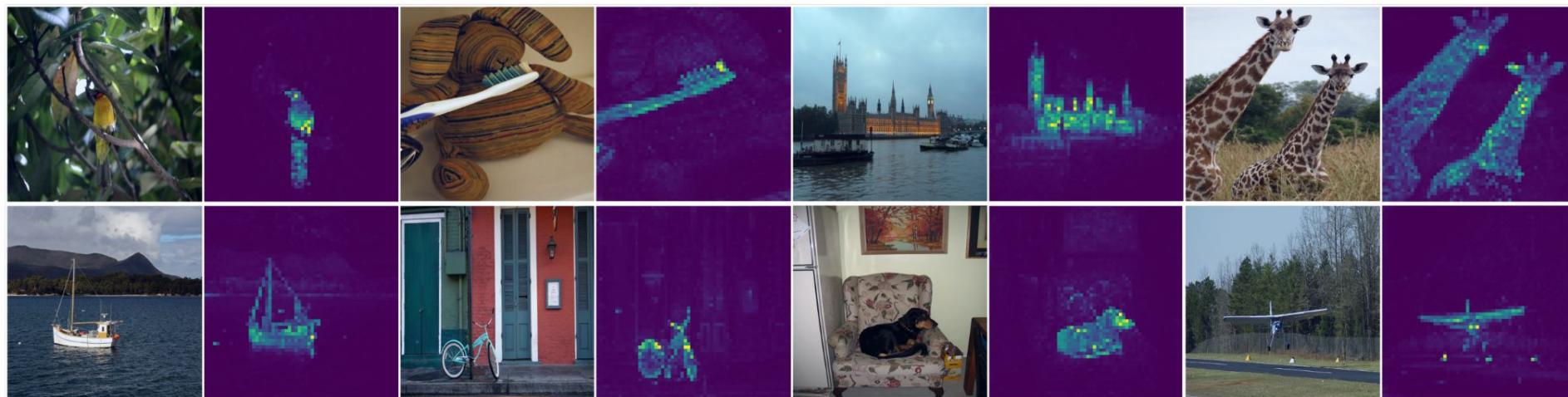
Cross attention visualization, when value = “pet”



Recall: “CLS” token in transformers



What is captured in CLS token?



Caron et.al., Emerging Properties in Self-Supervised Vision Transformers

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: $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q$

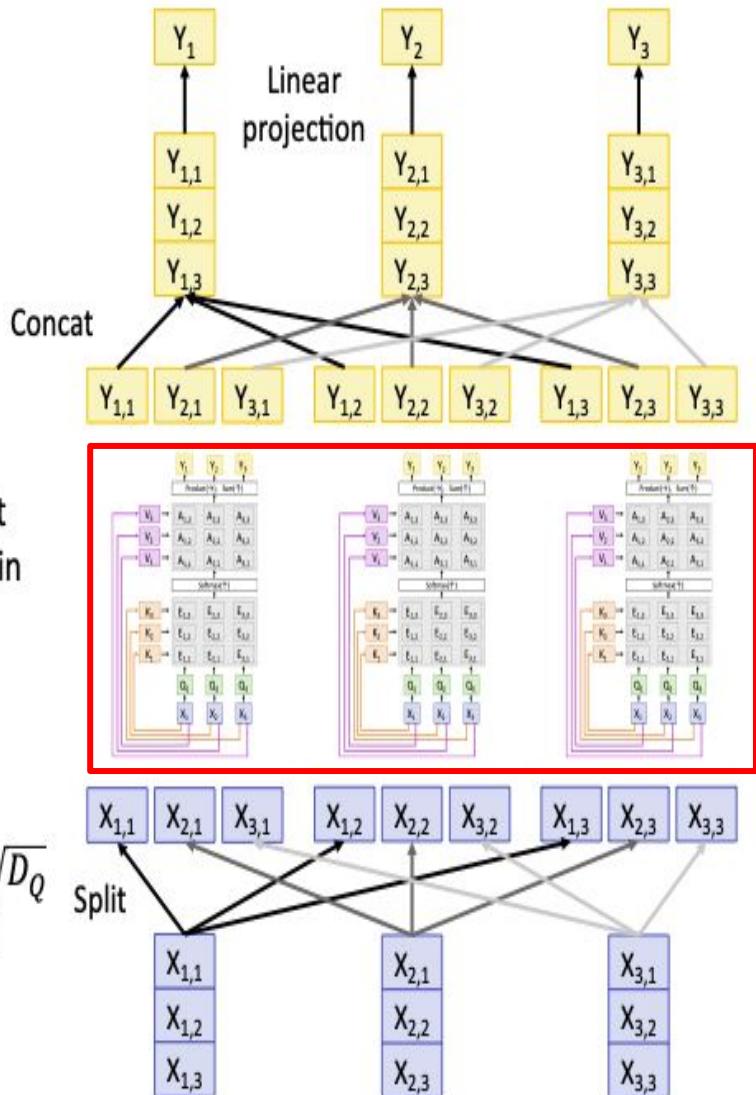
Key vectors: $\mathbf{K} = \mathbf{X}\mathbf{W}_K$ (Shape: $N_x \times D_Q$)

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

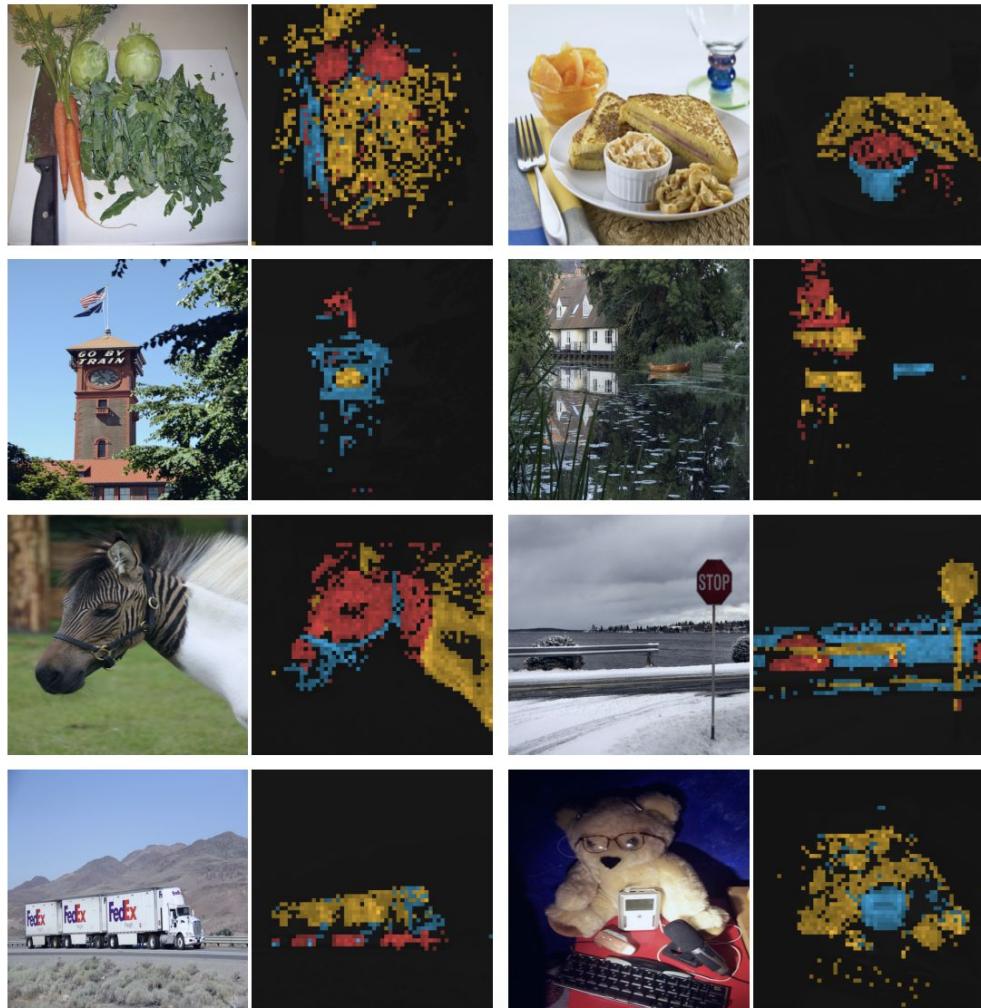
Similarities: $E = \mathbf{Q}\mathbf{K}^T / \sqrt{D_Q}$ (Shape: $N_x \times N_x$) $E_{i,j} = (\mathbf{Q}_i \cdot \mathbf{K}_j) / \sqrt{D_Q}$

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

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



What is captured in CLS token in the intermediate attention layers?



*Each color corresponds to different attention head.