

# Introduction to Transformers

## Lecture-2

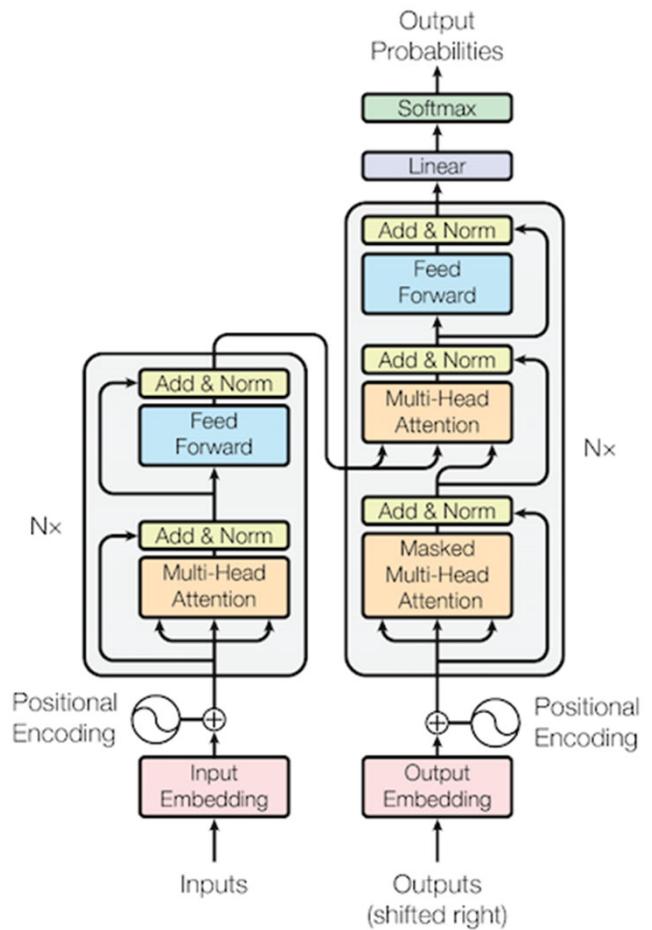
### CAP6412 Spring 2024

Mubarak Shah

[shah@crcv.ucf.edu](mailto:shah@crcv.ucf.edu)

# Contents

- Basics
  - What is Transformer?
  - Self-Attention
    - Query, Key, Value
  - Position encoding
  - Encoder-Decoder
- BERT
- Vision Transformers
  - ViT
  - SWIN



# Transformer

- Used for modeling long dependencies between input sequence elements
- Supports parallel processing of sequence as compared to RNN (e.g. LSTM)
- Allows processing multiple modalities
  - (e.g., images, videos, text and speech) using similar processing blocks
- Typically, pre-trained using pretext tasks on largescale (unlabeled) datasets
- Demonstrates excellent scalability to very large networks and huge datasets.
- **GPT-4** (Generative Pretrained Transformer)

# Vision Applications

- Recognition tasks (e.g., image classification, object detection, action recognition, and segmentation),
- Generative modeling, multi-modal tasks (e.g., visual-question answering, visual reasoning, and visual grounding),
- Video processing (e.g., activity recognition, video forecasting),
- Low-level vision (e.g., image super-resolution, image enhancement, and colorization)
- 3D analysis (e.g., point cloud classification and segmentation)

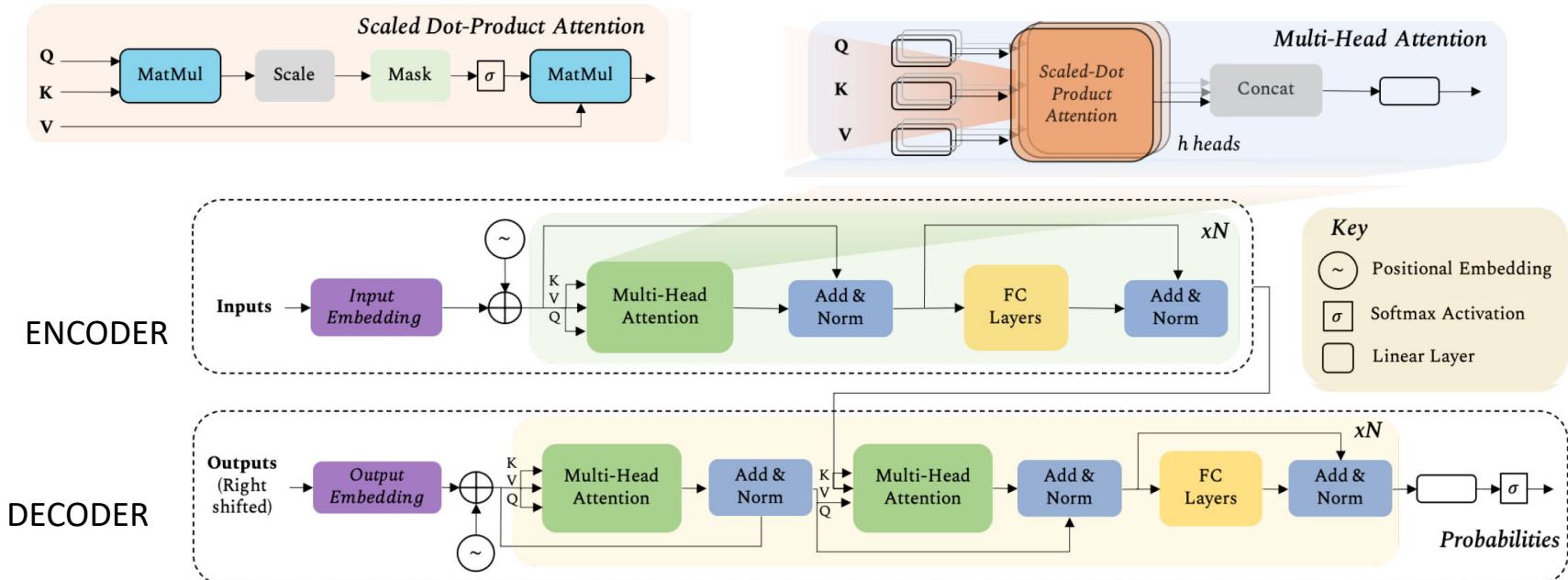
# Natural Language Processing

- BERT (Bidirectional Encoder Representations from Transformers),
- GPT (Generative Pre-trained **Transformer**) v1-4,
- RoBERTa (Robustly Optimized BERT Pre-training)
- T5 (Text-to-Text Transfer Transformer)

# Transformer Basics

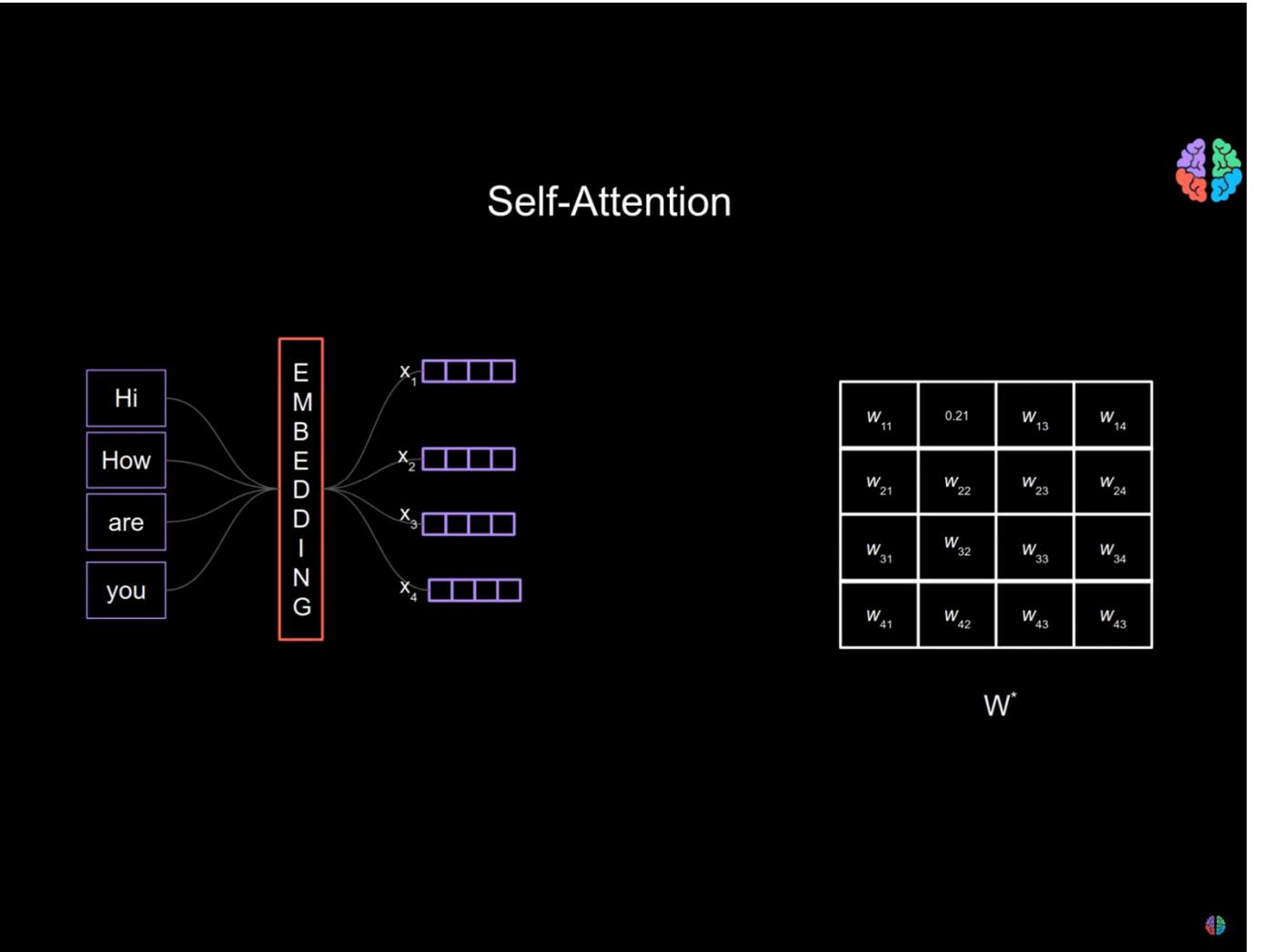
- It consists of Encoder and Decoder Blocks
- Main components of each block:
  - Self-Attention
  - Layer Normalization
  - Feed Forward Network

# Transformer



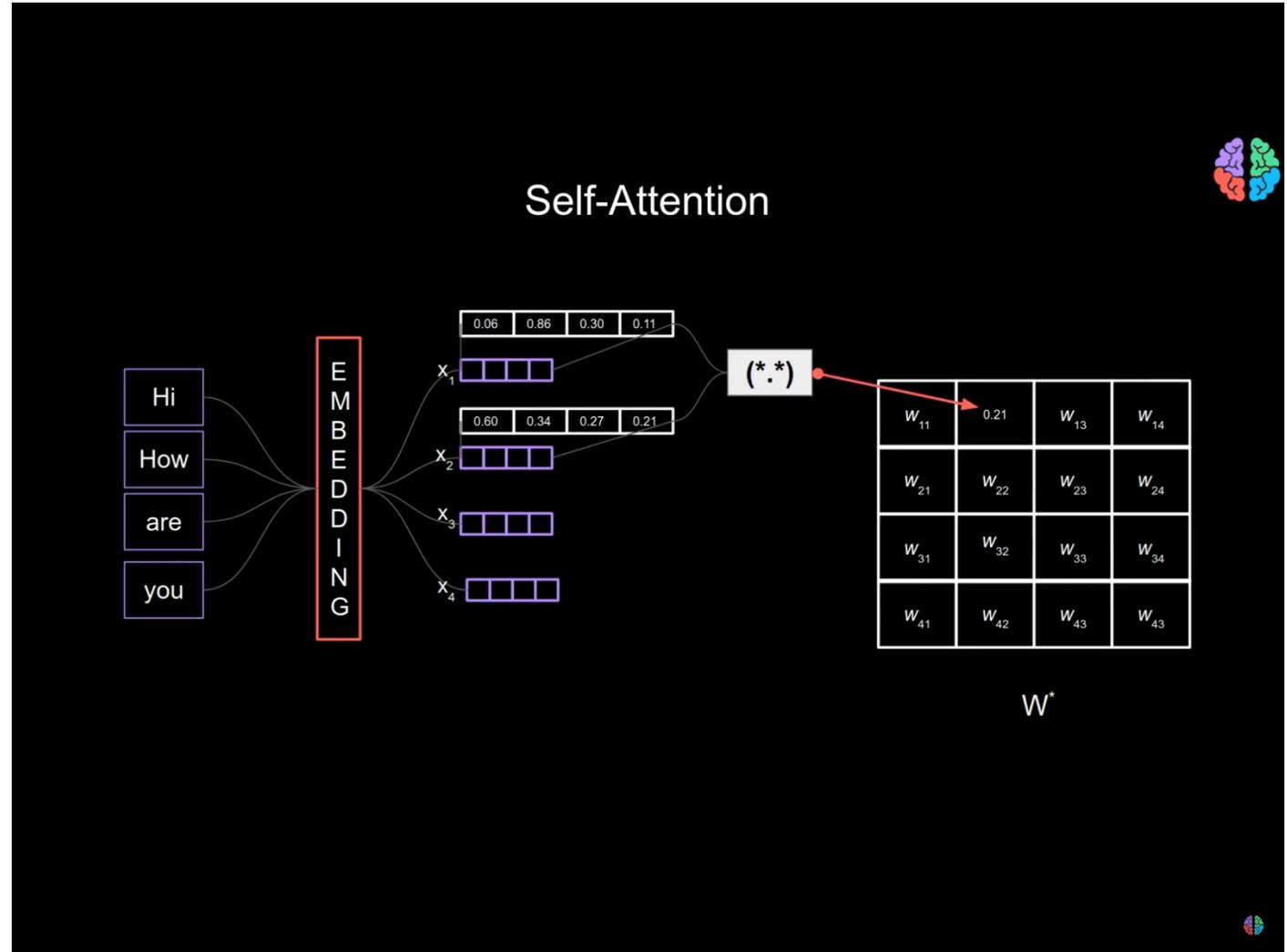


Slide courtesy of AI Bites, Youtube Channel:  
<https://www.youtube.com/c/AIBites>



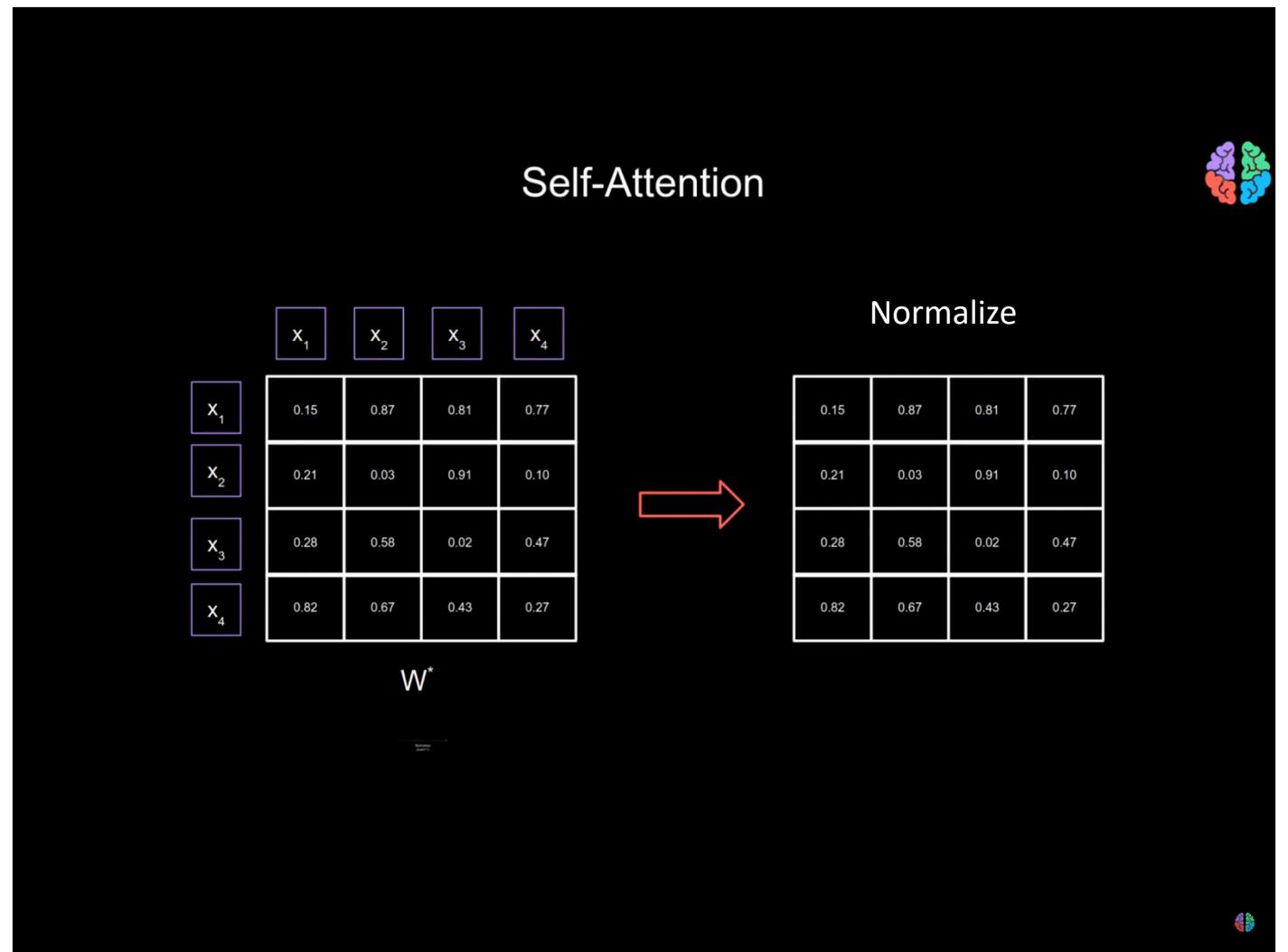


Slide courtesy of AI Bites, Youtube Channel:  
<https://www.youtube.com/c/AIBites>





Slide courtesy of AI Bites, Youtube Channel:  
<https://www.youtube.com/c/AIBites>



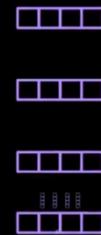


## Self-Attention

0.15	0.87	0.81	0.77
0.21	0.03	0.91	0.10
0.28	0.58	0.02	0.47
0.82	0.67	0.43	0.27

W

\*



Slide courtesy of AI Bites, Youtube Channel:  
<https://www.youtube.com/c/AIBites>





## Self-Attention

0.15	0.87	0.81	0.77
0.21	0.03	0.91	0.10
0.28	0.58	0.02	0.47
0.82	0.67	0.43	0.27

W



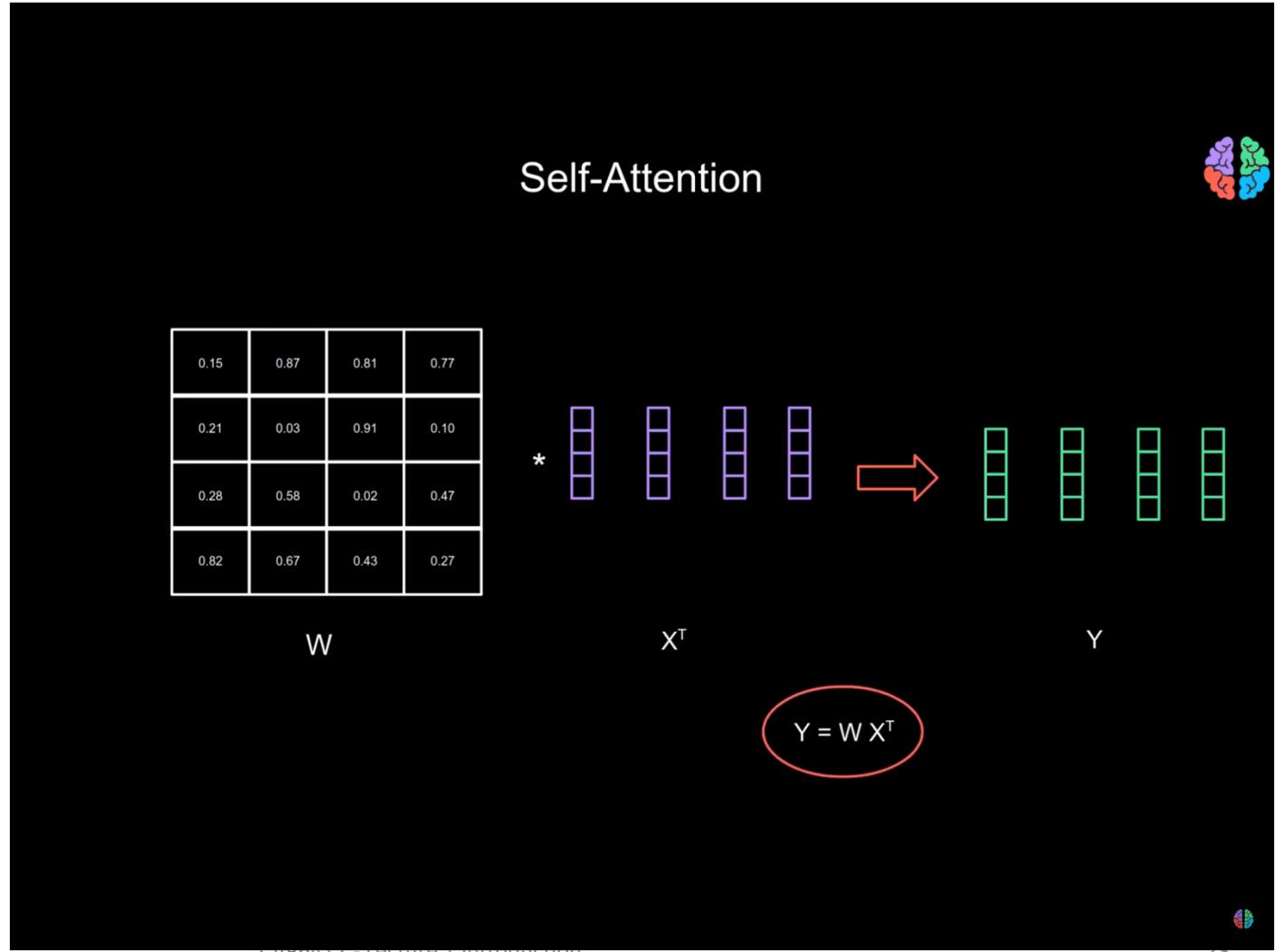
$X^T$

Y



Slide courtesy of AI Bites, Youtube Channel:  
<https://www.youtube.com/c/AIBites>

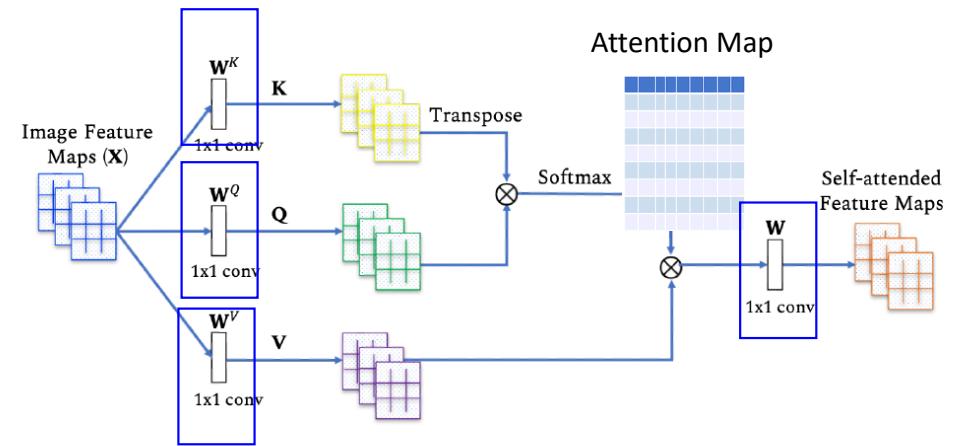
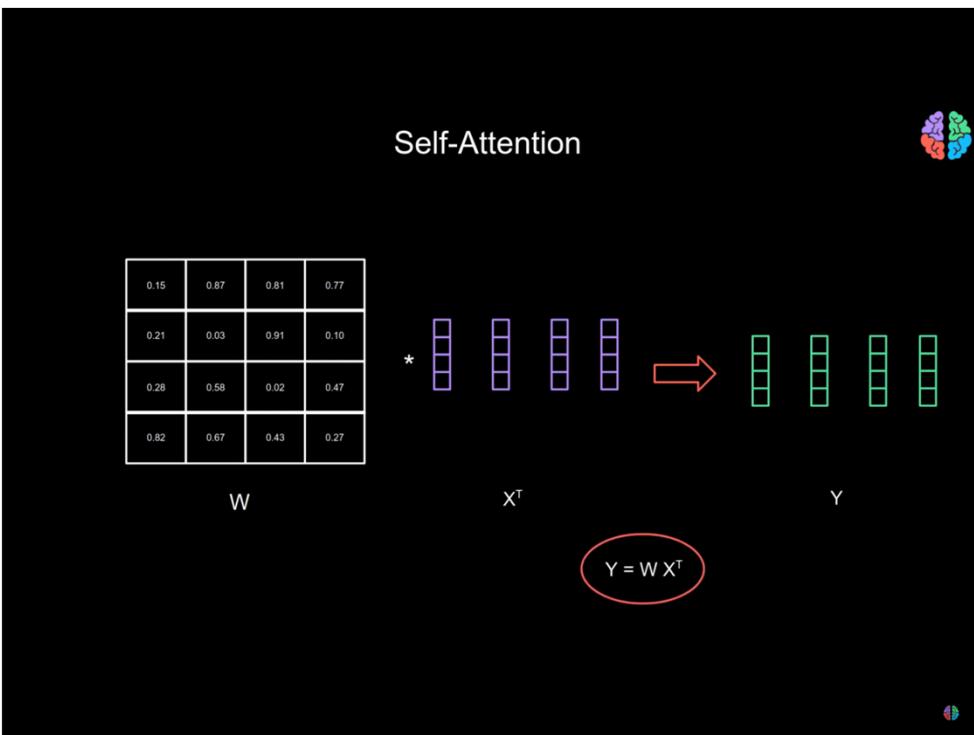
Slide courtesy of AI Bites, Youtube  
Channel:  
<https://www.youtube.com/c/AIBites>



# Self-Attention

- So far no learning!

# Self-Attention of Image Features



# Self-Attention (Matrix Explanation)

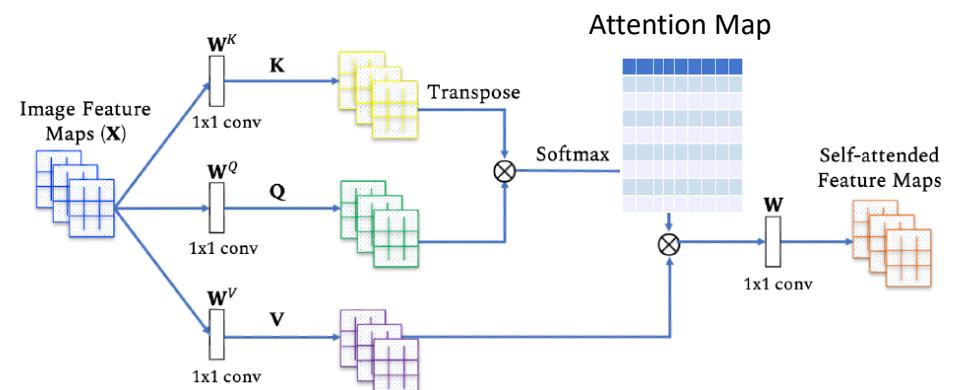
Lets denote a sequence of  $n$  entities  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$

by  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,

Queries ( $\mathbf{W}^Q \in \mathbb{R}^{d \times d_q}$ ), Keys ( $\mathbf{W}^K \in \mathbb{R}^{d \times d_k}$ ) and Values ( $\mathbf{W}^V \in \mathbb{R}^{d \times d_v}$ ), where  $d_q = d_k$ . The input sequence  $\mathbf{X}$  is

first projected onto these weight matrices to get  $\mathbf{Q} = \mathbf{X}\mathbf{W}^Q$ ,  $\mathbf{K} = \mathbf{X}\mathbf{W}^K$  and  $\mathbf{V} = \mathbf{X}\mathbf{W}^V$ . The output  $\mathbf{Z} \in \mathbb{R}^{n \times d_v}$  of the

$$\mathbf{Z} = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_q}} \right) \mathbf{V}.$$



# Self-Attention

Lets denote a sequence of  $n$  entities  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$

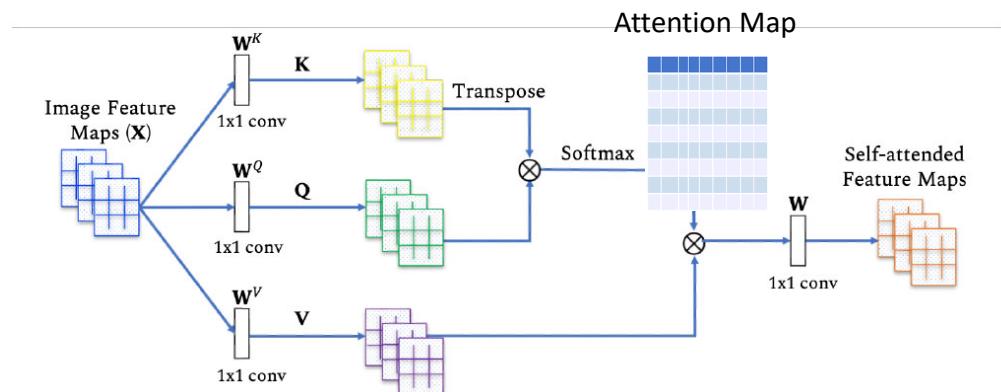
by  $\mathbf{X} \in \mathbb{R}^{n \times d}$   
 $9 \times 3$

Queries ( $\mathbf{W}^Q \in \mathbb{R}^{d \times d_q}$ ), Keys ( $\mathbf{W}^K \in \mathbb{R}^{d \times d_k}$ ) and Values ( $\mathbf{W}^V \in \mathbb{R}^{d \times d_v}$ ), where  $d_q = d_k$ . The input sequence  $\mathbf{X}$  is  
 $3 \times 3$

first projected onto these weight matrices to get  $\mathbf{Q} = \mathbf{X}\mathbf{W}^Q$ ,  
 $\mathbf{K} = \mathbf{X}\mathbf{W}^K$  and  $\mathbf{V} = \mathbf{X}\mathbf{W}^V$ . The output  $\mathbf{Z} \in \mathbb{R}^{n \times d_v}$  of the  
 $9 \times 3$

$$\mathbf{Z} = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_q}} \right) \mathbf{V}.$$

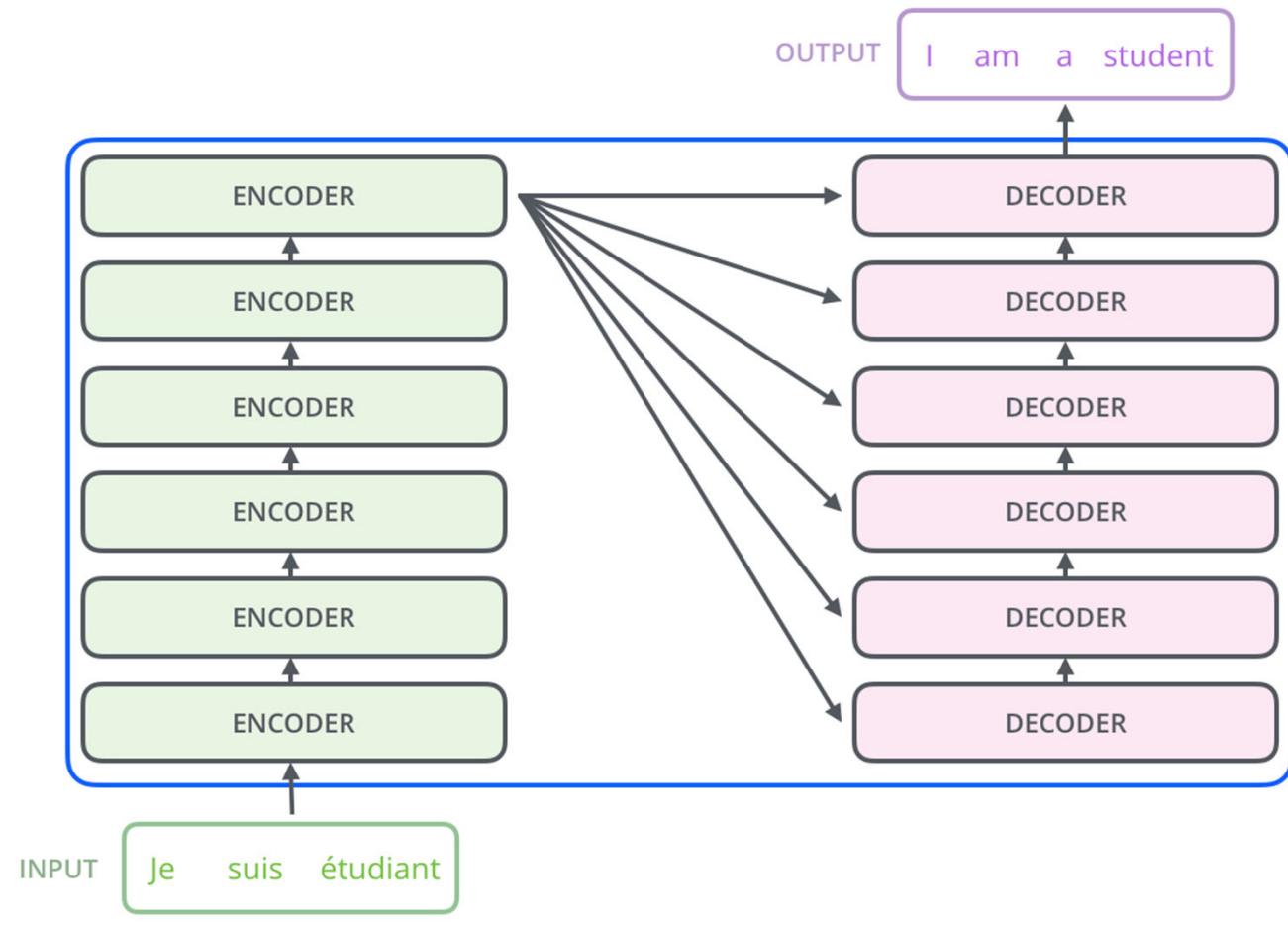
$$9 \times 3 \quad 9 \times 3 \quad 9 \times 9 \quad 9 \times 3$$



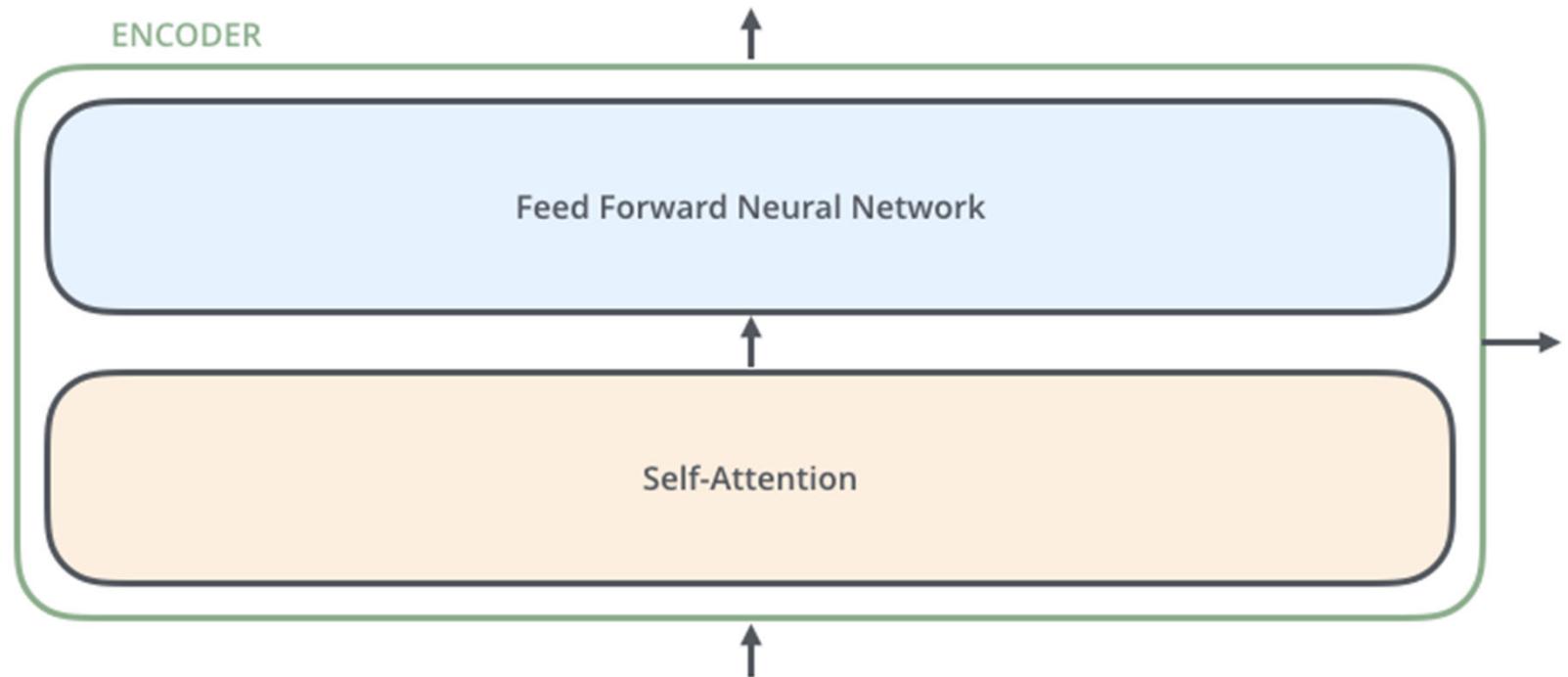
# Transformers (Attention is all you need 2017)

- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in NeurIPS, 2017. (102,618 citations)
- Two valuable sources
  - <http://nlp.seas.harvard.edu/2018/04/03/attention.html>
  - <https://jalammar.github.io/illustrated-transformer/> (slides come from this source)
- Slides from Ming Li, University of Waterloo, **CS 886 Deep Learning and NLP**

## Transformer

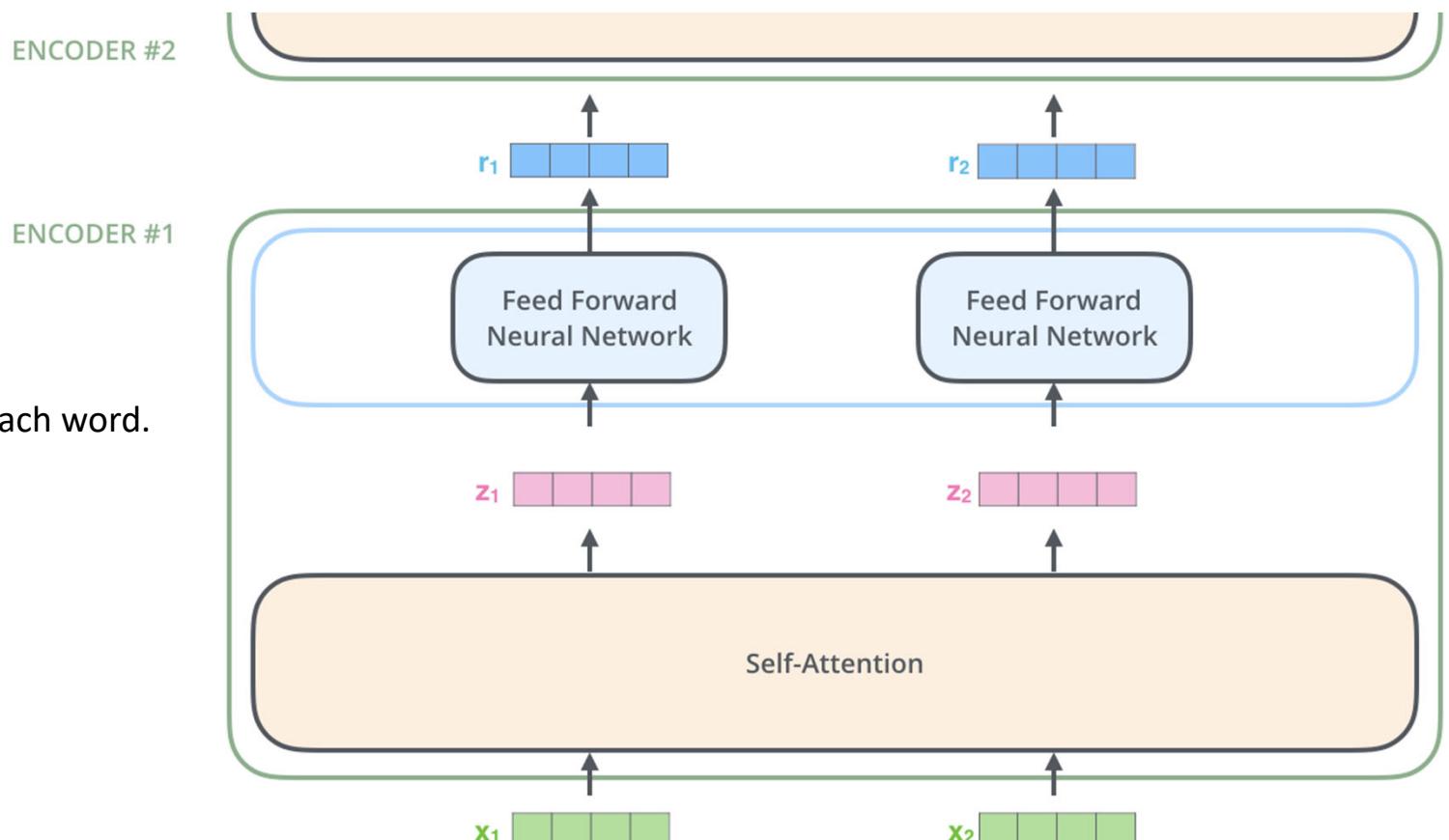


An Encoder Block: same structure, different parameters



## Encoder

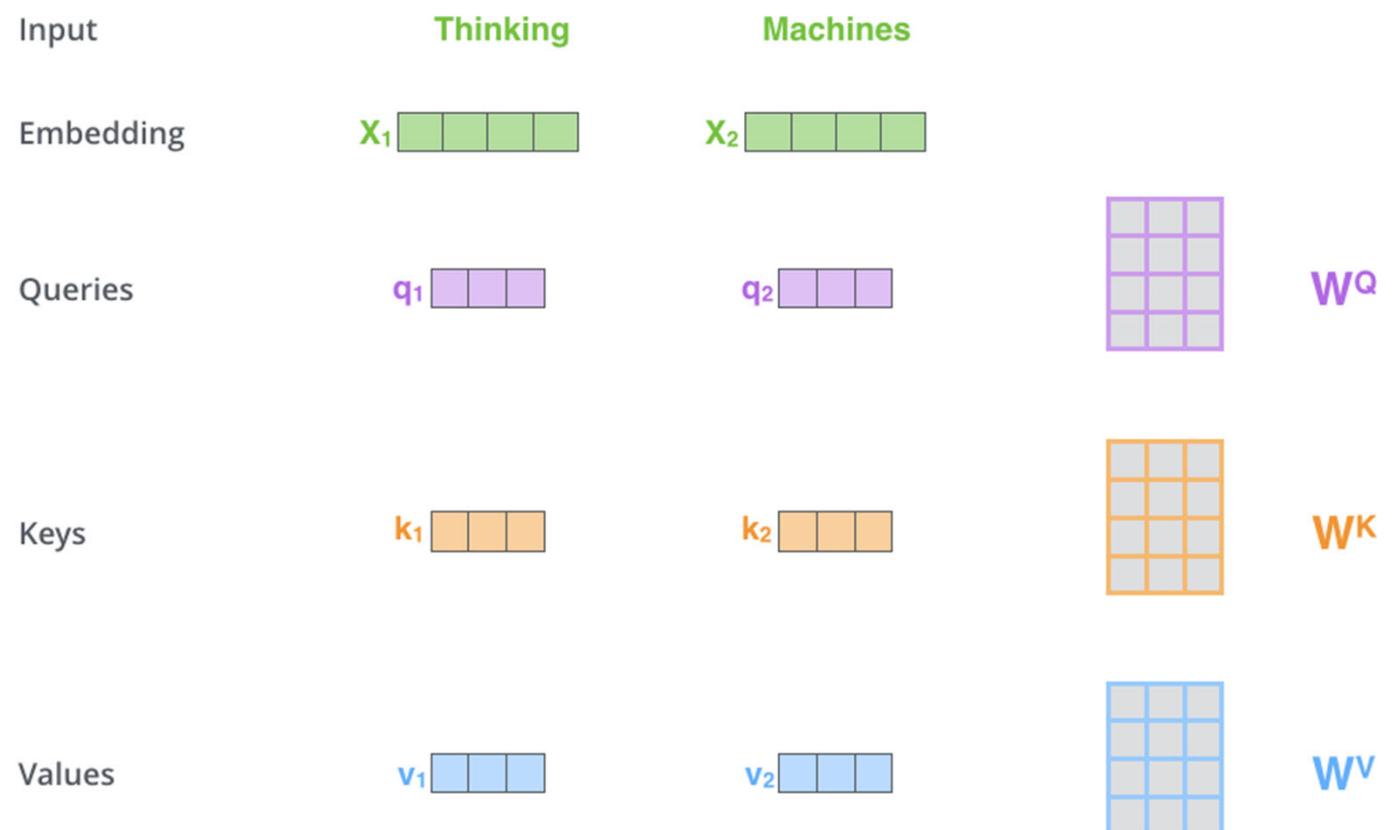
The ffnn is independent for each word.  
Hence can be parallelized.



## Self Attention

- First, we create three vectors by multiplying input embedding  $x_i$  with three matrices

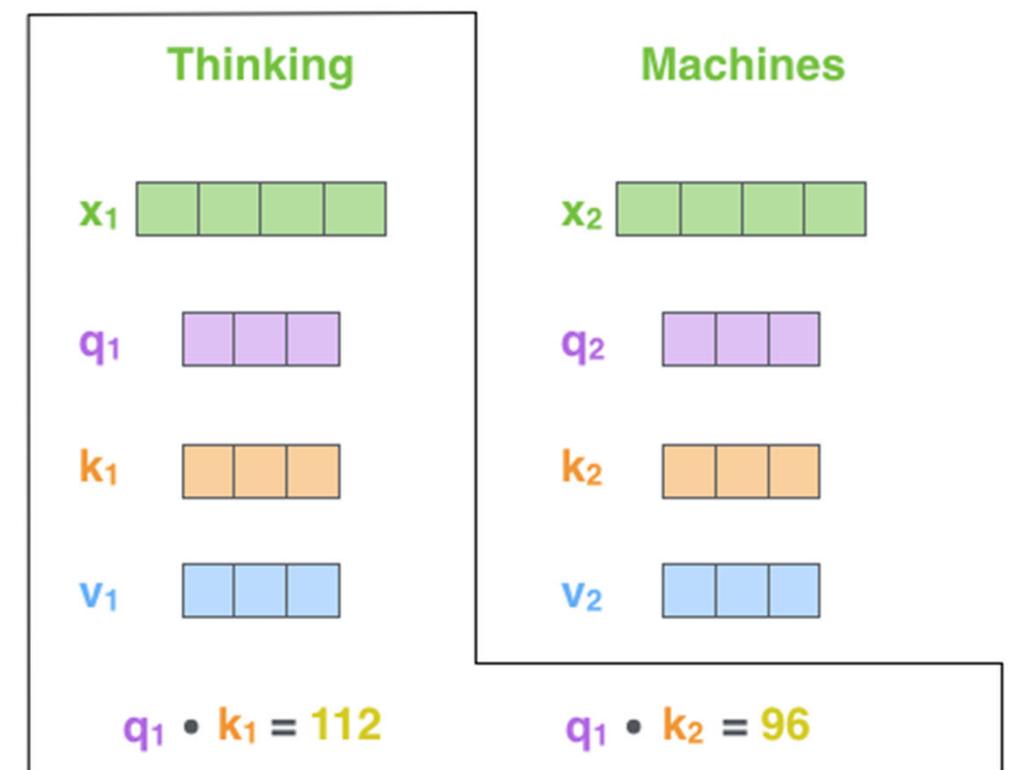
- $q_i = x_i W^Q$
- $K_i = x_i W^K$
- $V_i = x_i W^V$



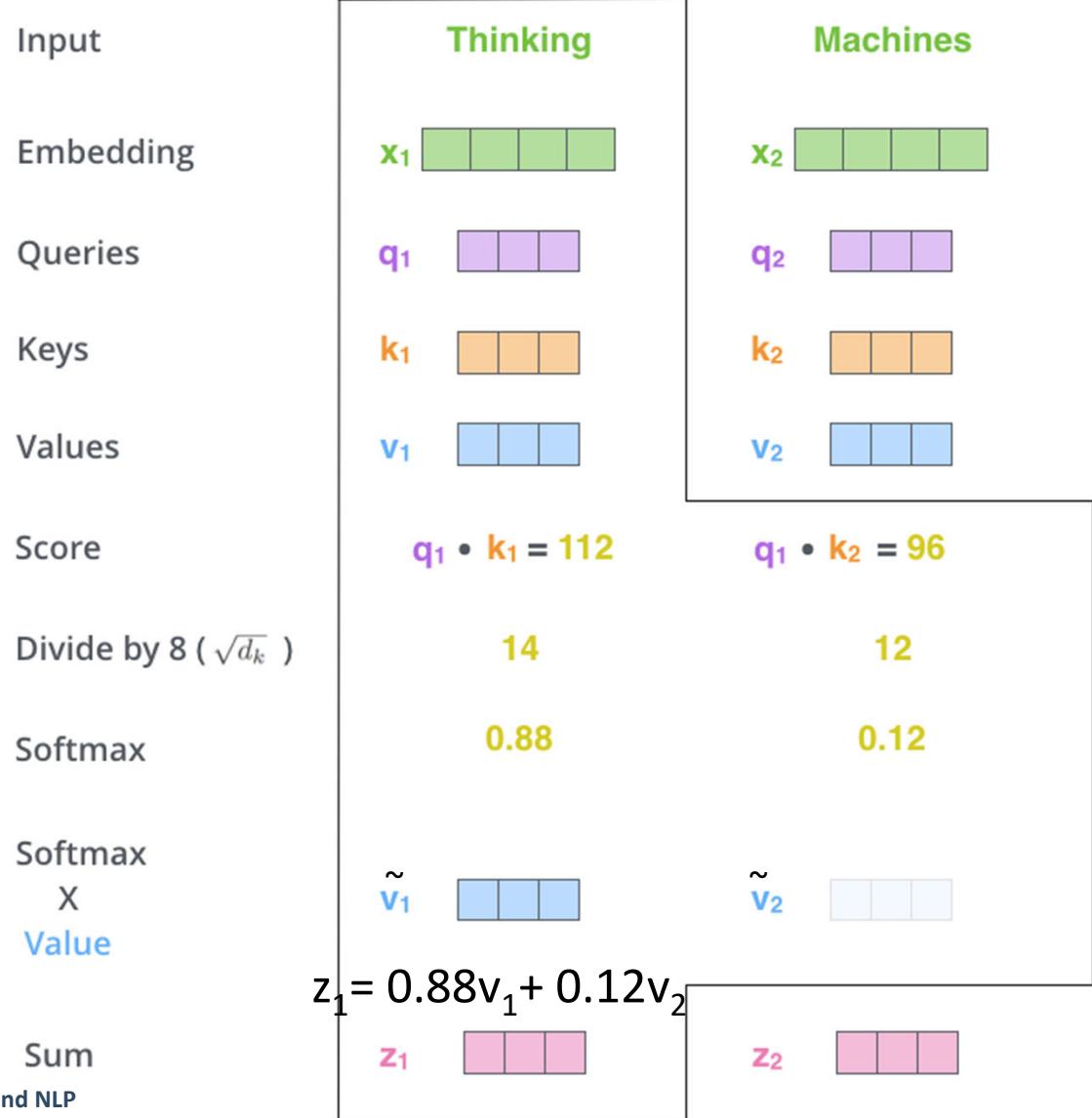
## Self Attention

Now we calculate a score to determine how much focus to place on other parts of the input.

Input  
Embedding  
Queries  
Keys  
Values  
Score

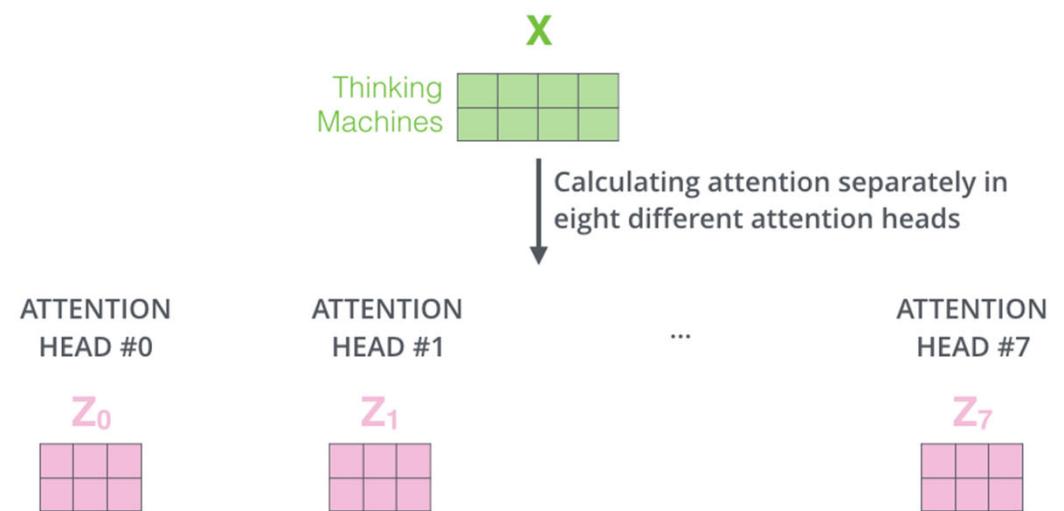


## Self Attention



## Multiple heads

1. It expands the model's ability to focus on different positions.
2. It gives the attention layer multiple "representation subspaces"

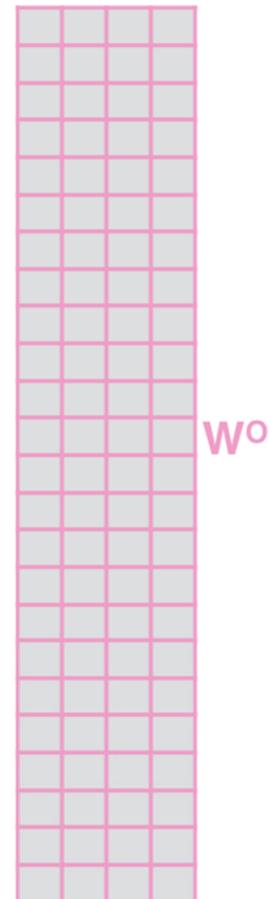


## Attention and Transformers

1) Concatenate all the attention heads



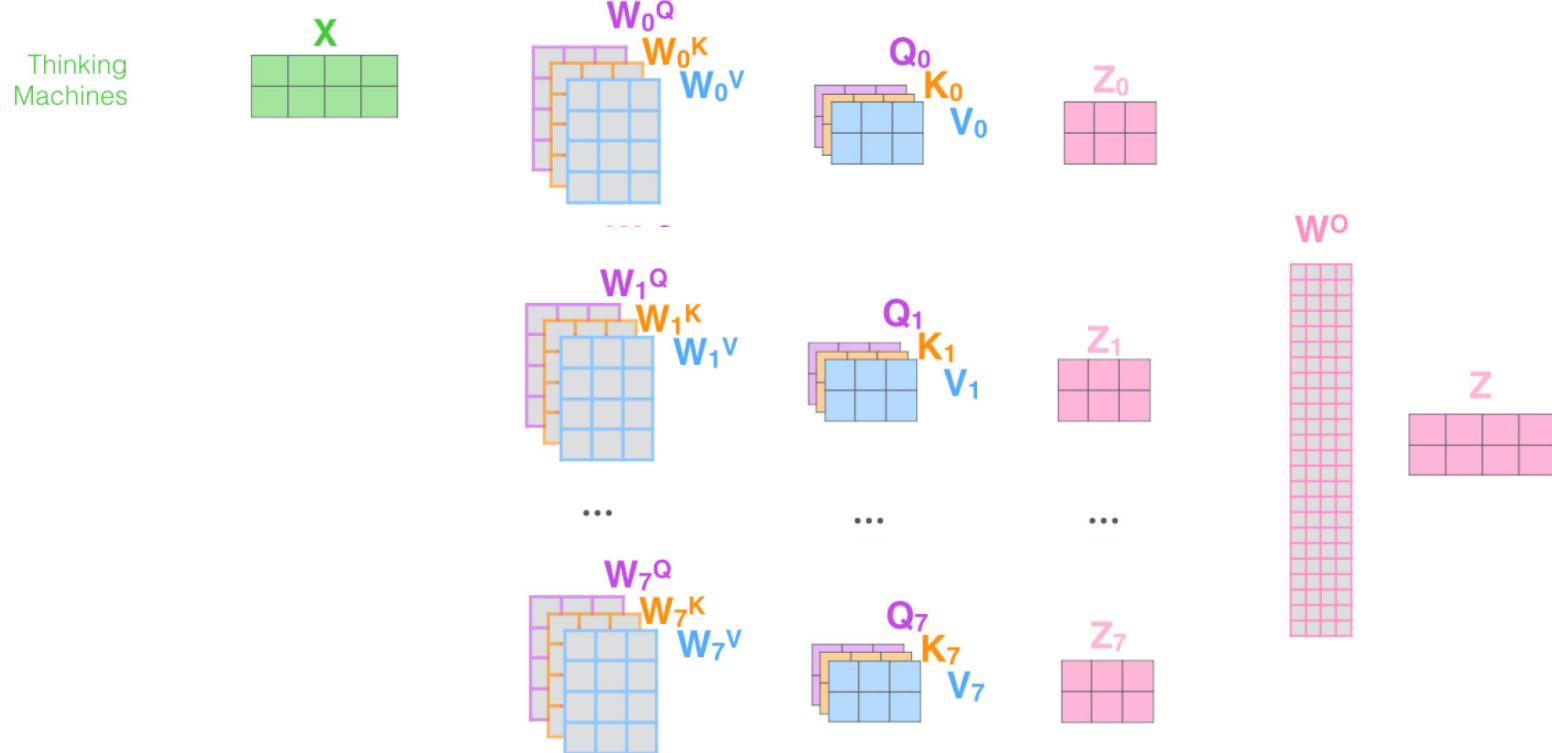
2) Multiply with a weight matrix  $W^o$  that was trained jointly with the model



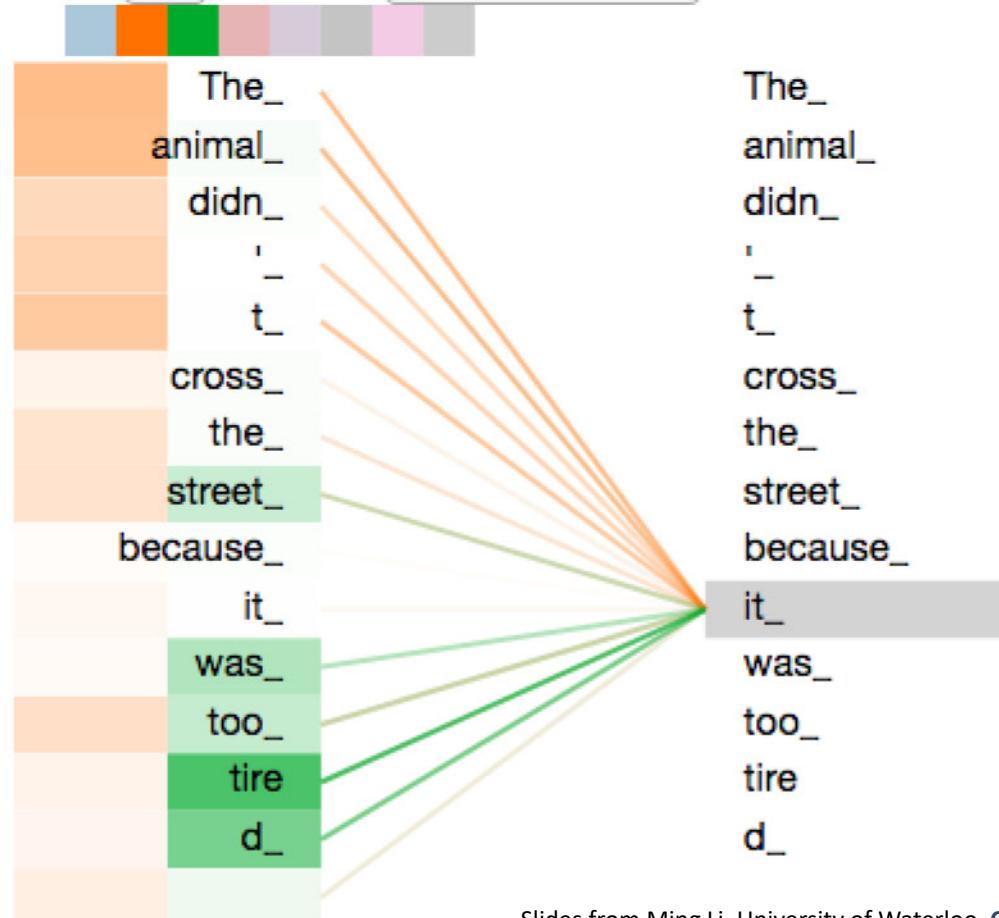
3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= \begin{matrix} Z \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array} \end{matrix}$$

- 1) This is our input sentence\*
- 2) We embed each word\*
- 3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices
- 4) Calculate attention using the resulting  $Q/K/V$  matrices
- 5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^o$  to produce the output of the layer

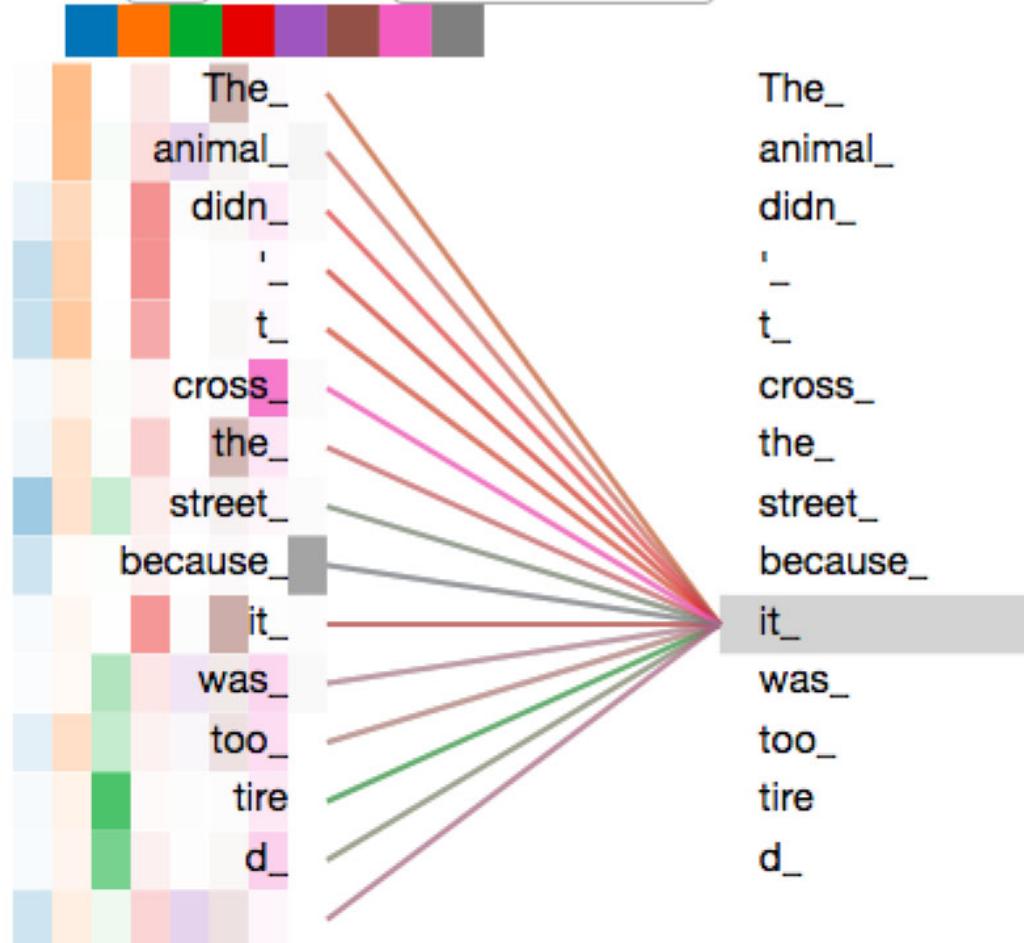


Layer: 5    Attention: Input - Input



Slides from Ming Li, University of Waterloo, CS 886 Deep Learning and NLP

Layer: 5 Attention: Input - Input





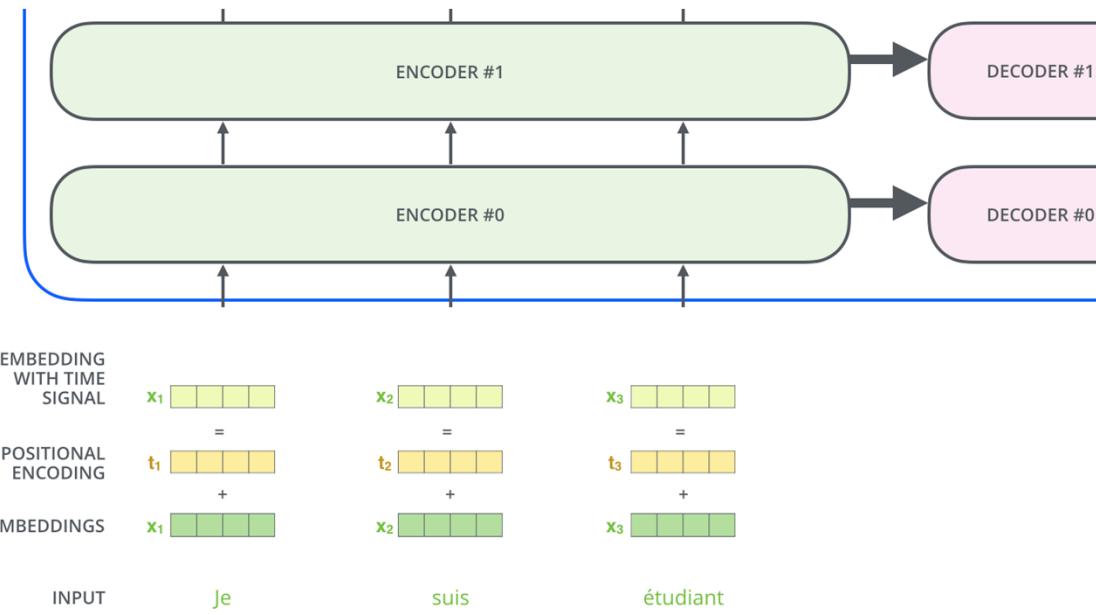
## Attention and Transformers

### Representing the input order (positional encoding)

- Transformer is permutation invariant
- The transformer adds a vector to each input embedding.
- These vectors follow a specific pattern that the model learns
- Learned pattern helps model
  - to determine the position of each word, or
  - the distance between different words.

## Attention and Transformers

### Representing the input order (positional encoding)



# Position Encoding

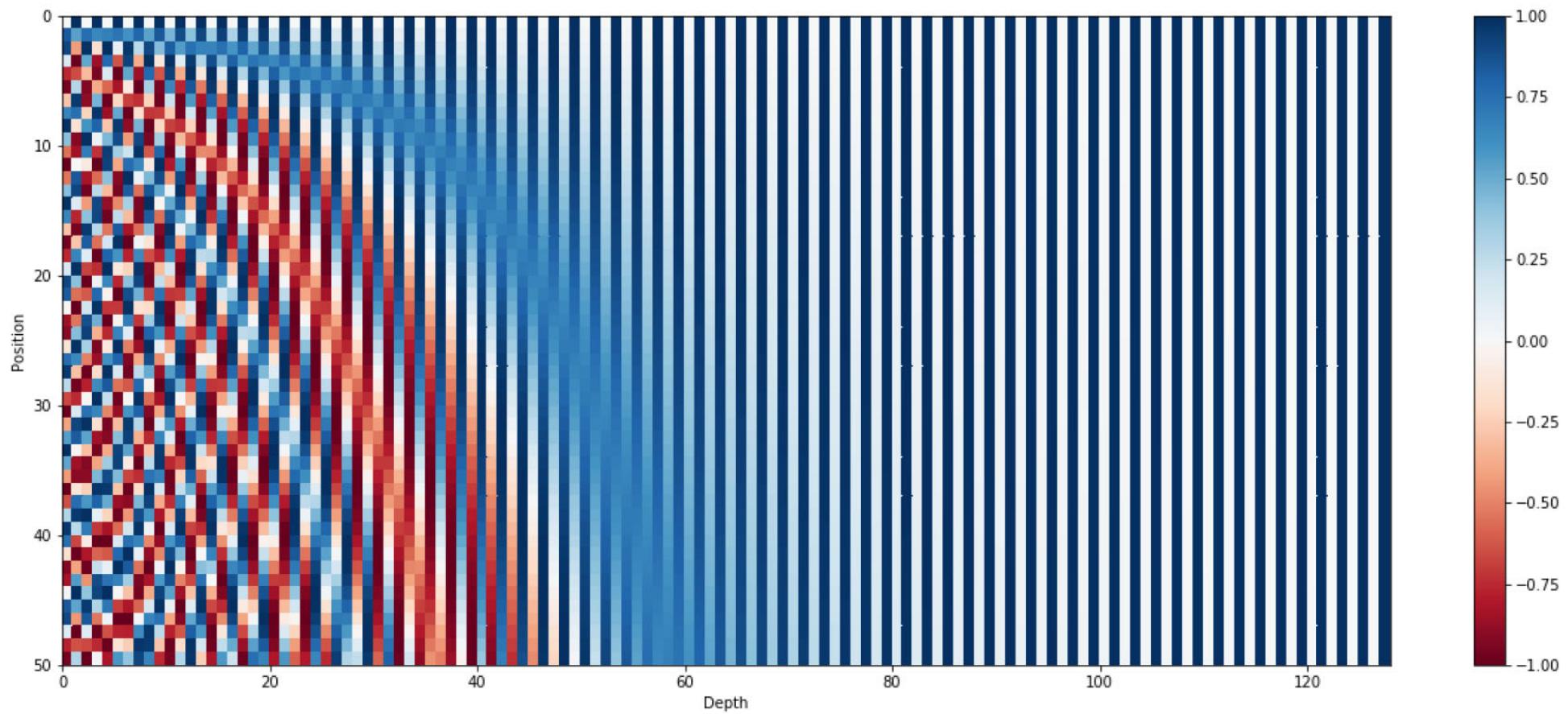
$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

where

$$\omega_k = \frac{1}{10000^{2k/d}}$$

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

# Position Encoding



$$P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$

$$P(k, 2i+1) = \cos\left(\frac{k}{n^{2i/d}}\right)$$

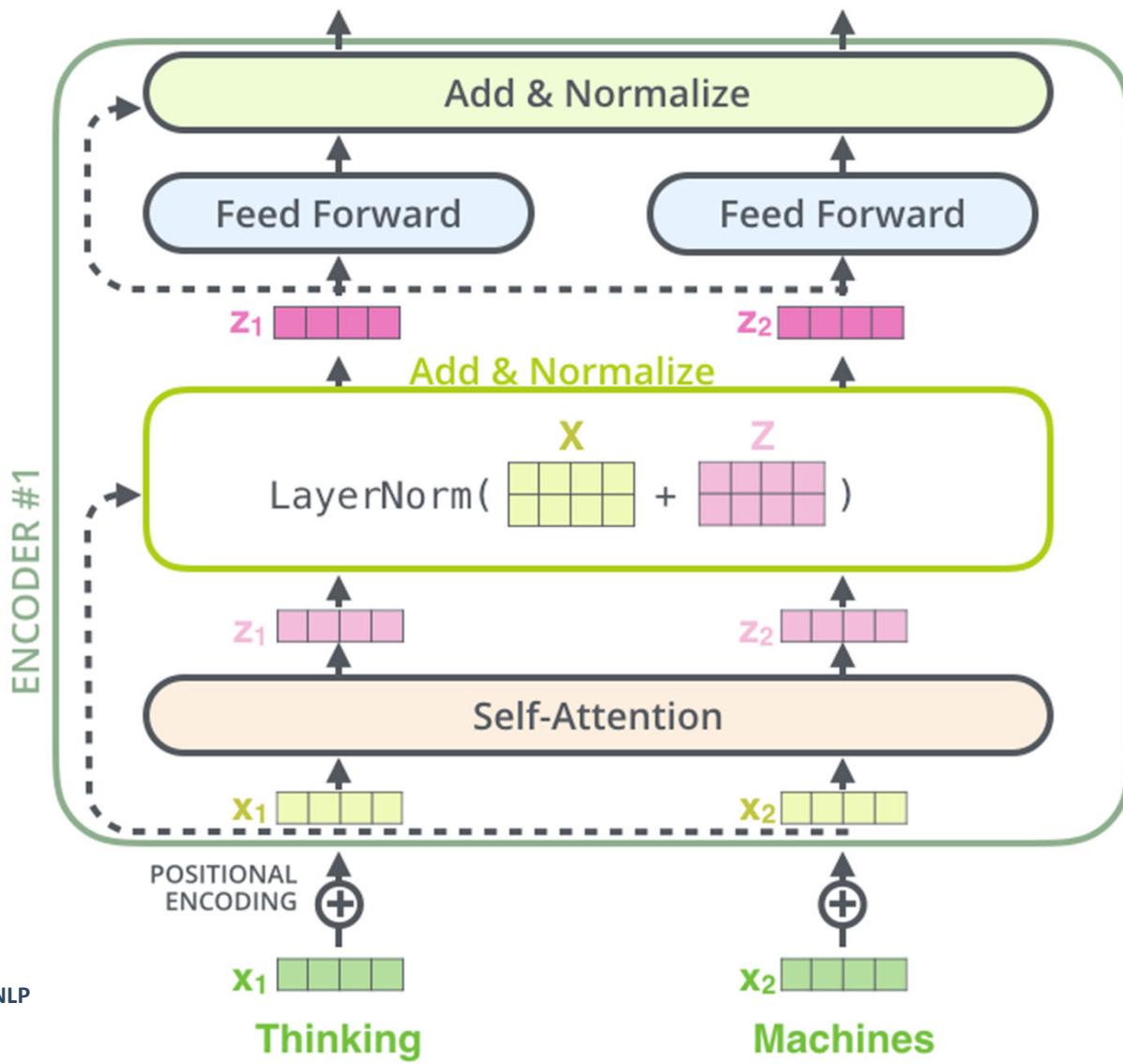
Sequence	Index of token, $k$	Positional Encoding Matrix with $d=4, n=100$			
		$i=0$	$i=0$	$i=1$	$i=1$
I	0	$P_{00}=\sin(0) = 0$	$P_{01}=\cos(0) = 1$	$P_{02}=\sin(0) = 0$	$P_{03}=\cos(0) = 1$
am	1	$P_{10}=\sin(1/1) = 0.84$	$P_{11}=\cos(1/1) = 0.54$	$P_{12}=\sin(1/10) = 0.10$	$P_{13}=\cos(1/10) = 1.0$
a	2	$P_{20}=\sin(2/1) = 0.91$	$P_{21}=\cos(2/1) = -0.42$	$P_{22}=\sin(2/10) = 0.20$	$P_{23}=\cos(2/10) = 0.98$
Robot	3	$P_{30}=\sin(3/1) = 0.14$	$P_{31}=\cos(3/1) = -0.99$	$P_{32}=\sin(3/10) = 0.30$	$P_{33}=\cos(3/10) = 0.96$

Positional Encoding Matrix for the sequence 'I am a robot'

# Position Encoding

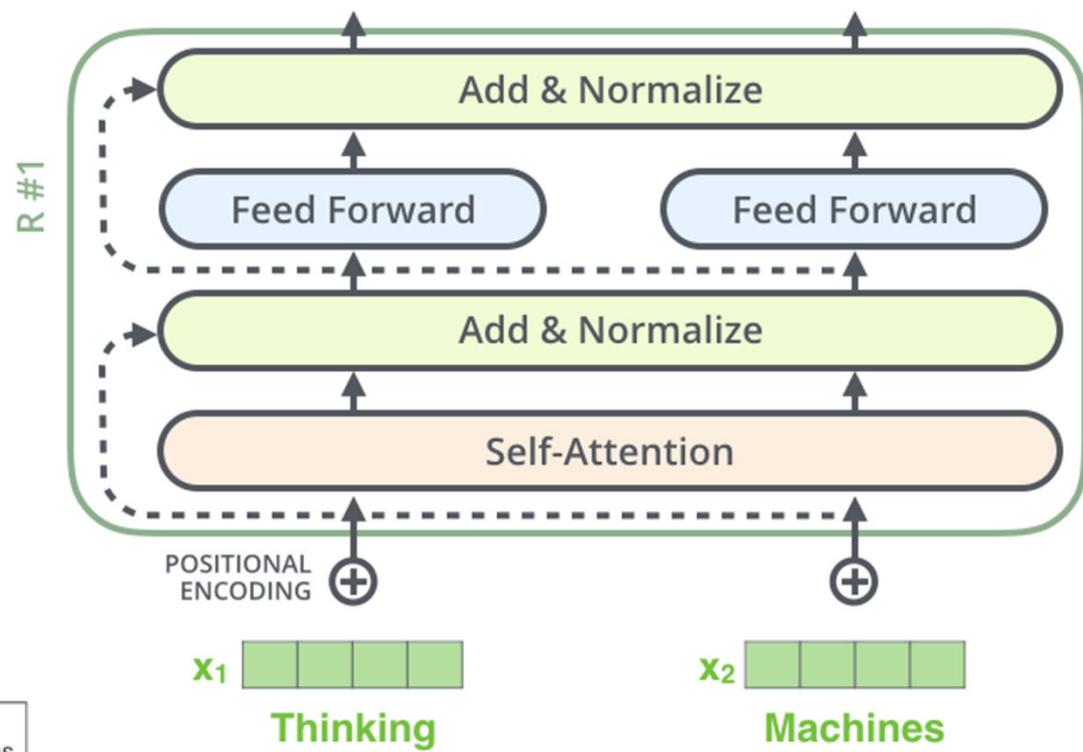
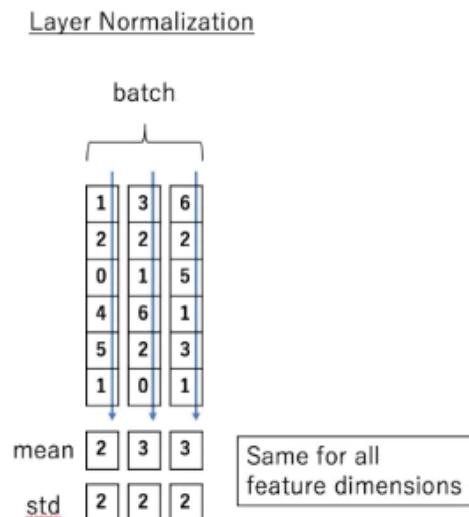
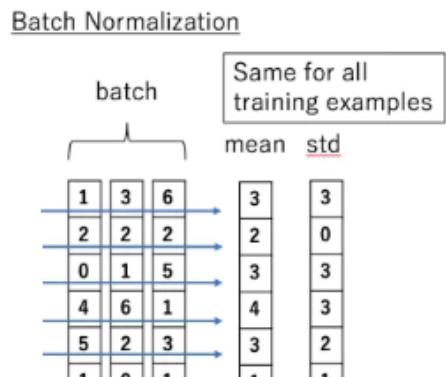
- Can also be learned
- Learn like other parameters

## Add and Normalize



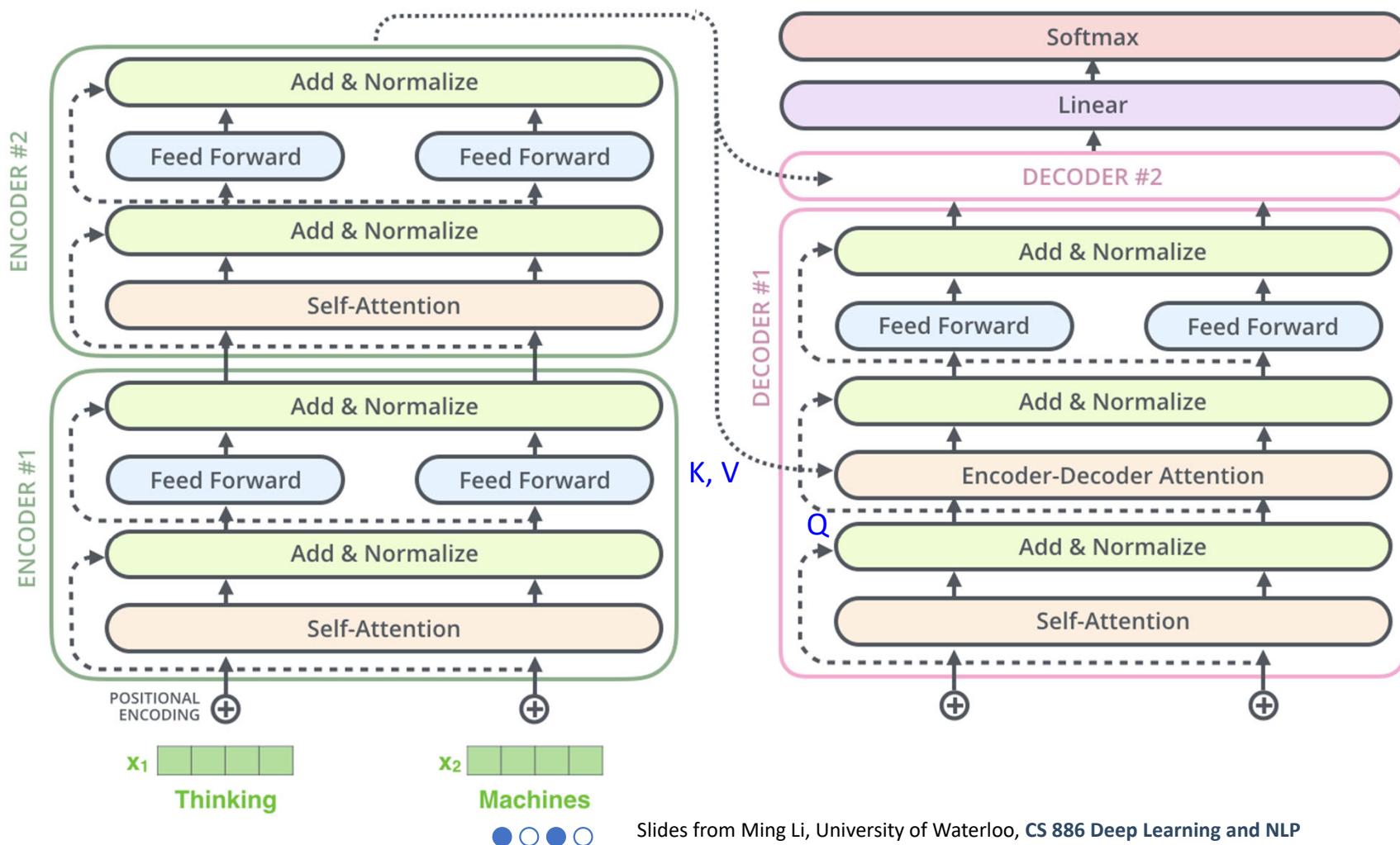
## Layer Normalization (Hinton)

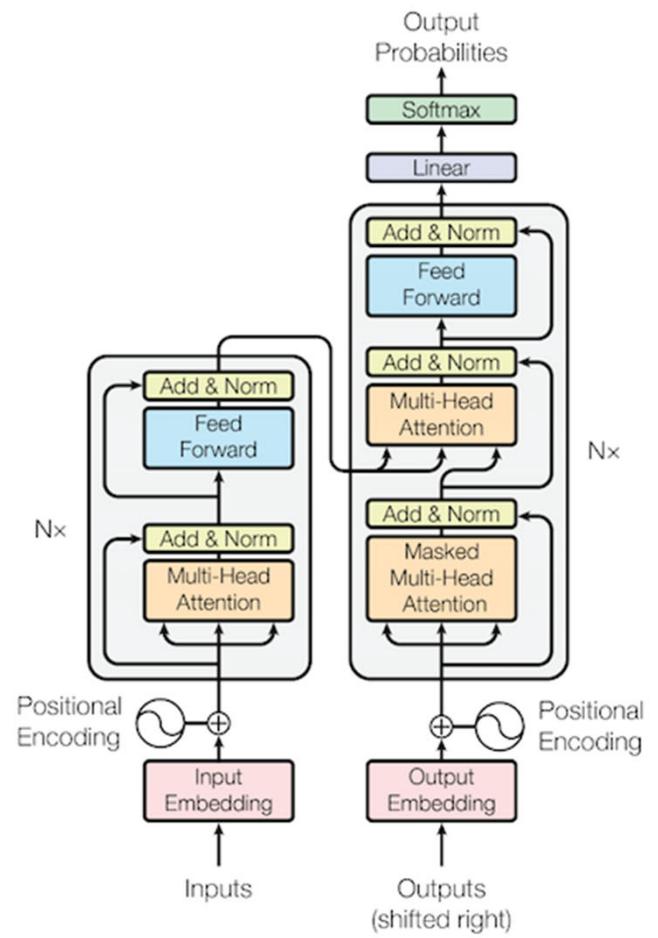
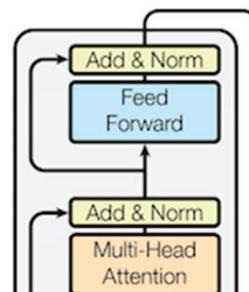
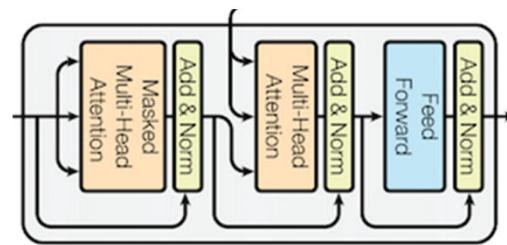
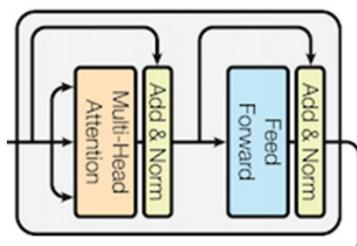
Layer normalization normalizes the inputs across the features.



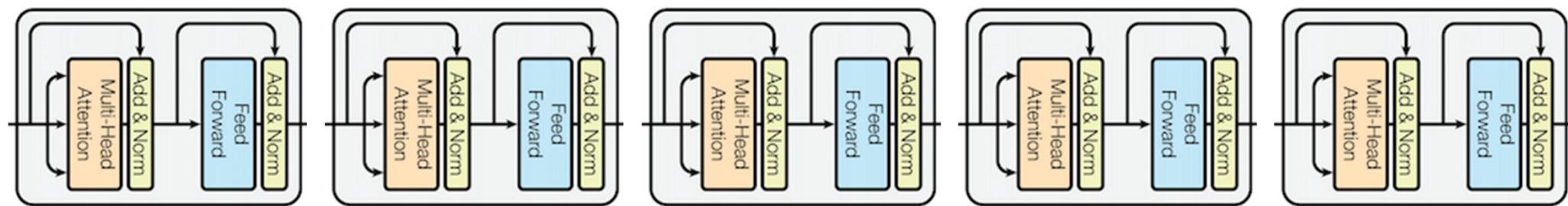
Slides from Ming Li, University of Waterloo, CS 886 Deep Learning and NLP

# The complete transformer





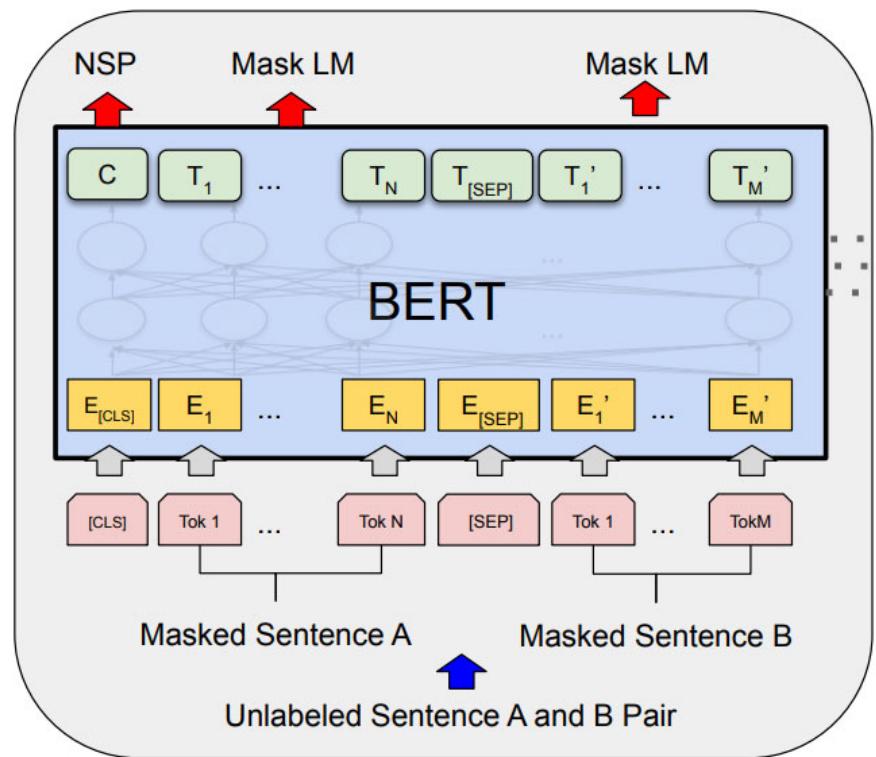
# BERT (Stack Encoder Blocks)

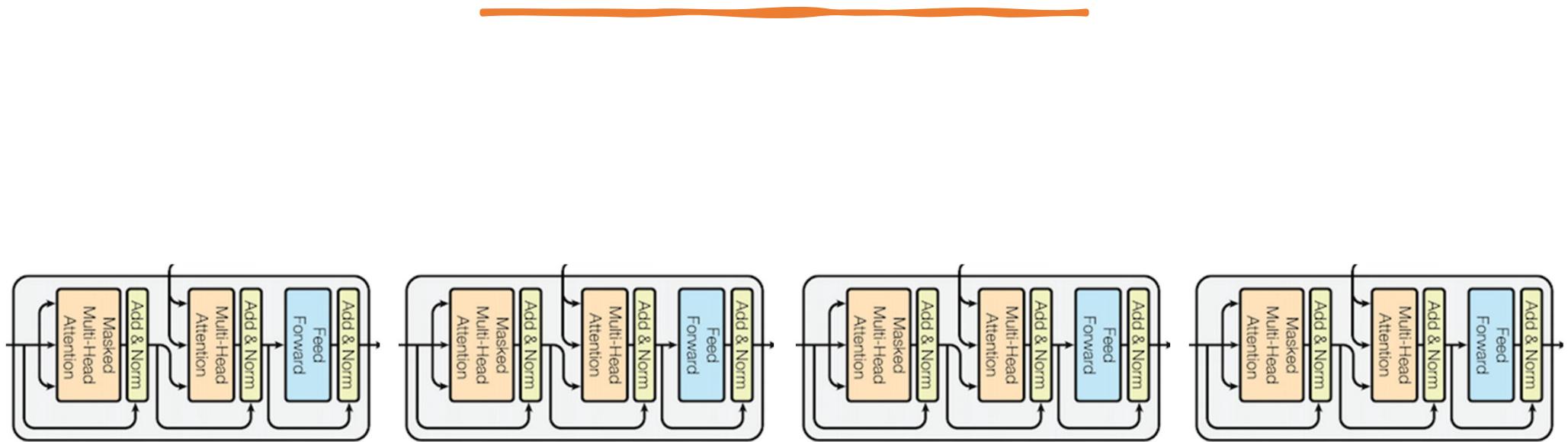


## BERT (Bidirectional Encoder Representations from Transformers)

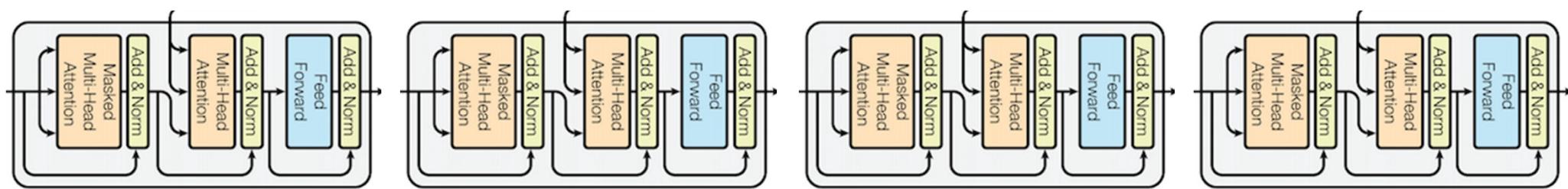
- BERT jointly encodes the right and left context of a word in a sentence to improve the learned feature representations
- BERT is trained on two pre-text tasks in self-supervised manner
  - Masked Language Model (MLM)
    - Mask fixed percentage (15%) of words in a sentence predict these masked words
    - In predicting the masked words, the model learns the bidirectional context.
  - Next Sentence Prediction (NSP)
    - Given a pair of sentences A and B the model predicts a binary label i.e., whether the pair is valid or not from the original document
    - Pair is formed such that B is the actual sentence (next to A) 50% of the time, and B is a random sentence for other 50% of the time.

# BERT





# GPT (Stack Decoder Blocks)



# BERT and GPT

[https://www.youtube.com/shorts/BEt\\_BACGw6g](https://www.youtube.com/shorts/BEt_BACGw6g)

# Vision Transformers

Mubarak Shah

[shah@crcv.ucf.edu](mailto:shah@crcv.ucf.edu)

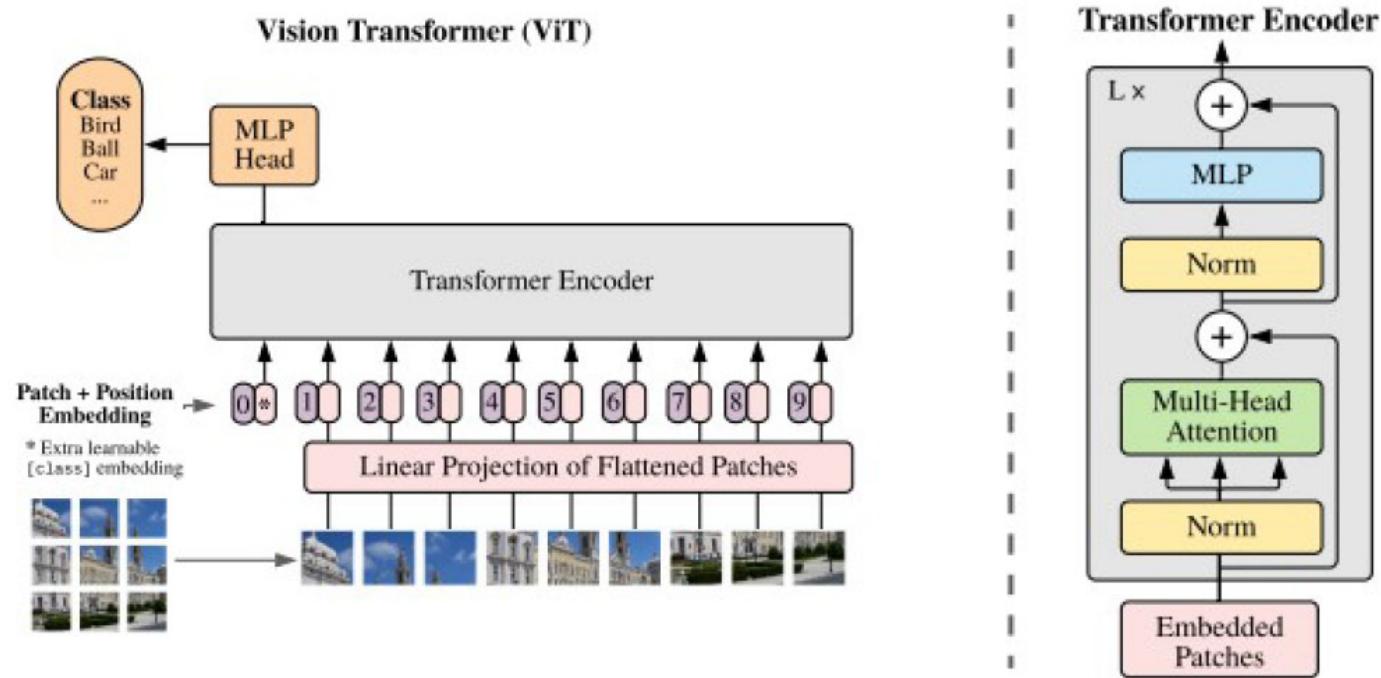
# Vision Transformer (ViT)

A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., “An image is worth 16x16 words: Transformers for image recognition at scale,” arXiv preprint arXiv:2010.11929, 2020. ([ICLR 2021; 26,938 citations!](#))

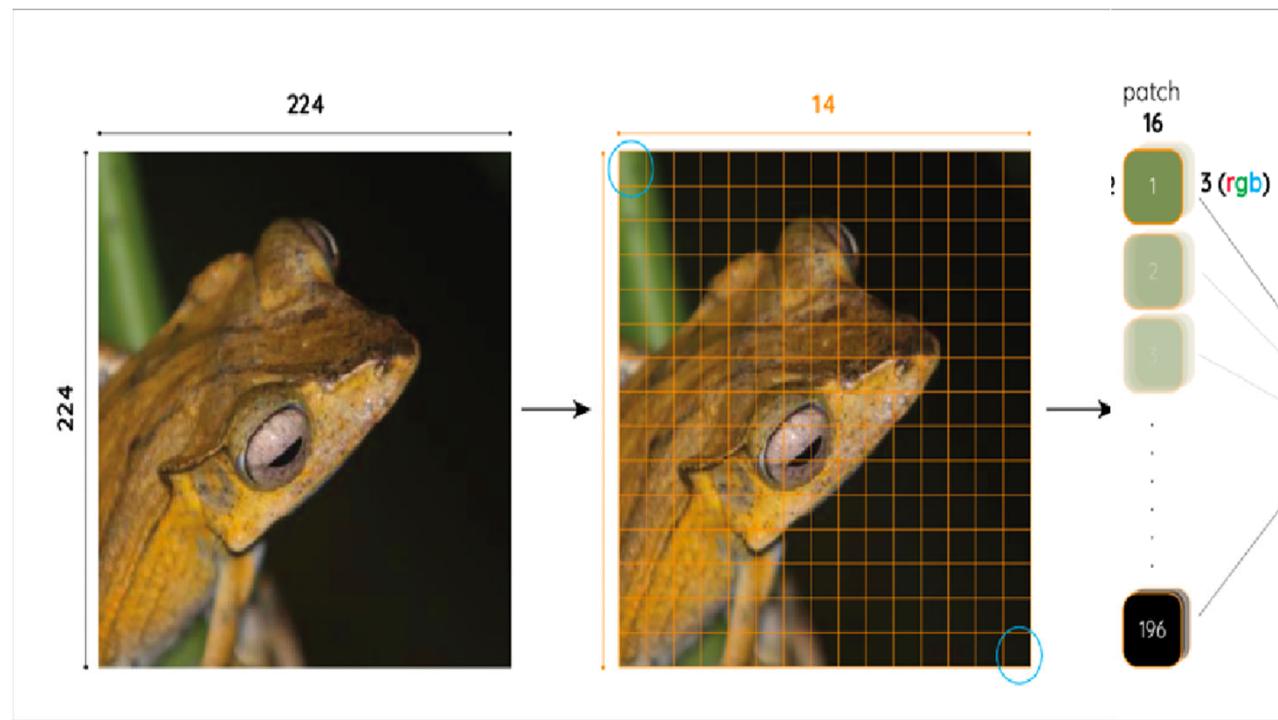
# Vision Transformer (ViT)

- Naive application of self-attention to images requires high computation
- Divide an image into 16x16 patches (tokens)
- Transformers need to be trained on large datasets
- ViT attains excellent results when pre-trained on JFT-300M
  - 88:55% on ImageNet,
  - 90:72% on ImageNet-Real,
  - 94:55% on CIFAR-100,
  - 77:63% on the VTAB suite of 19 tasks

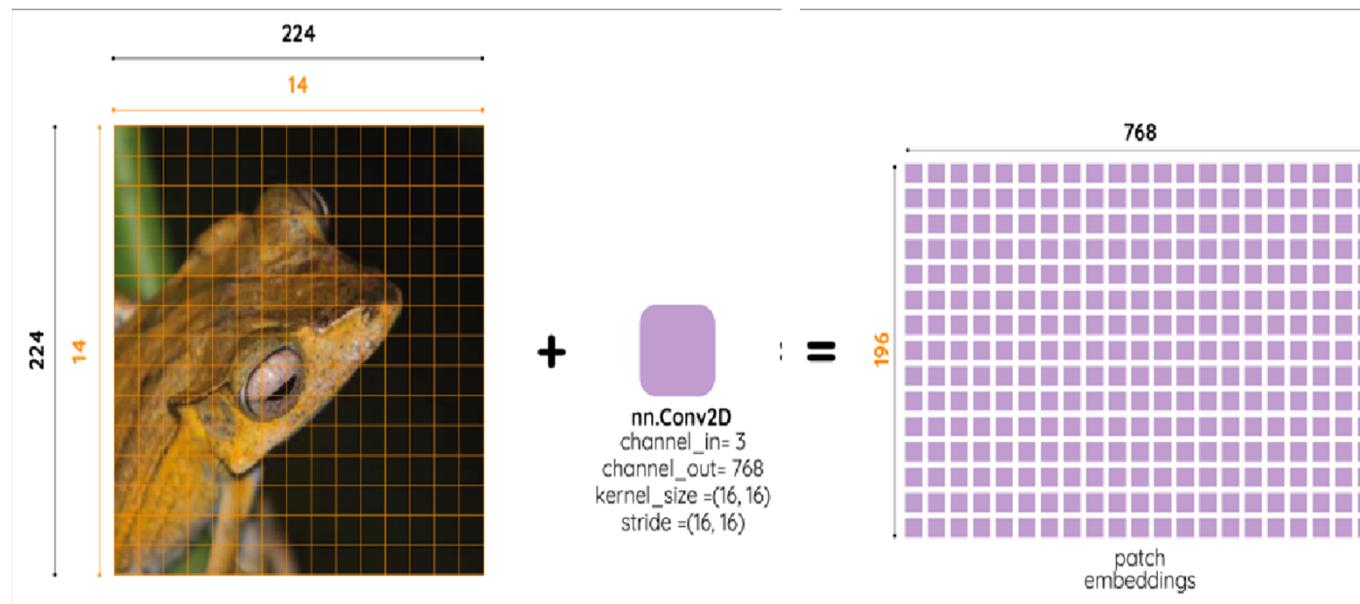
# Vision Transformer (ViT)



# Divide image into 16x16 patches

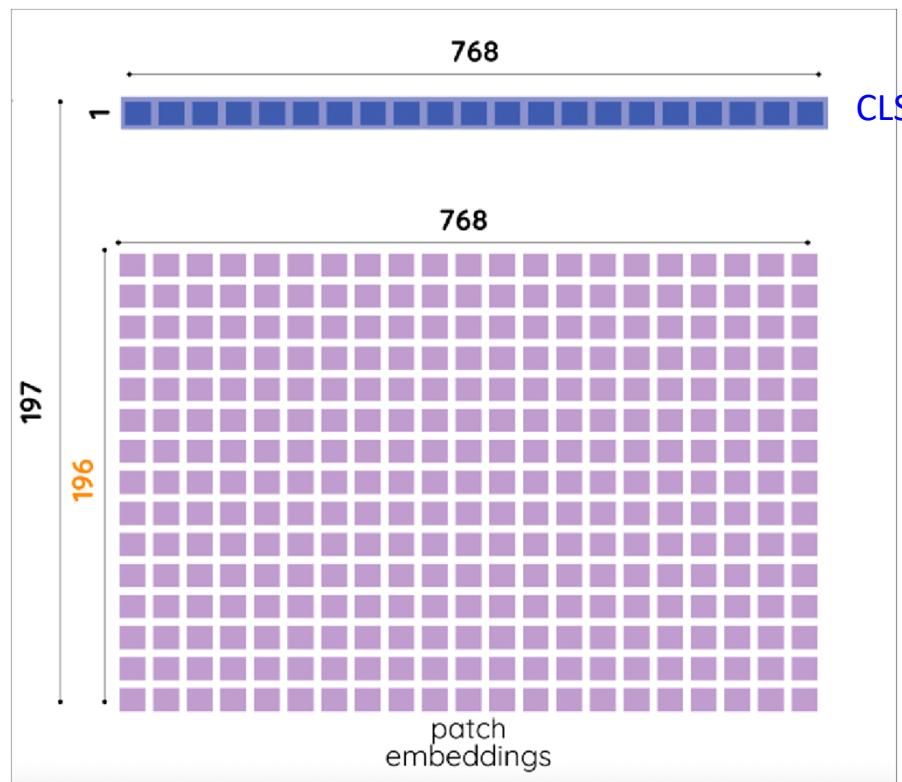


# Generate embedding for each patch

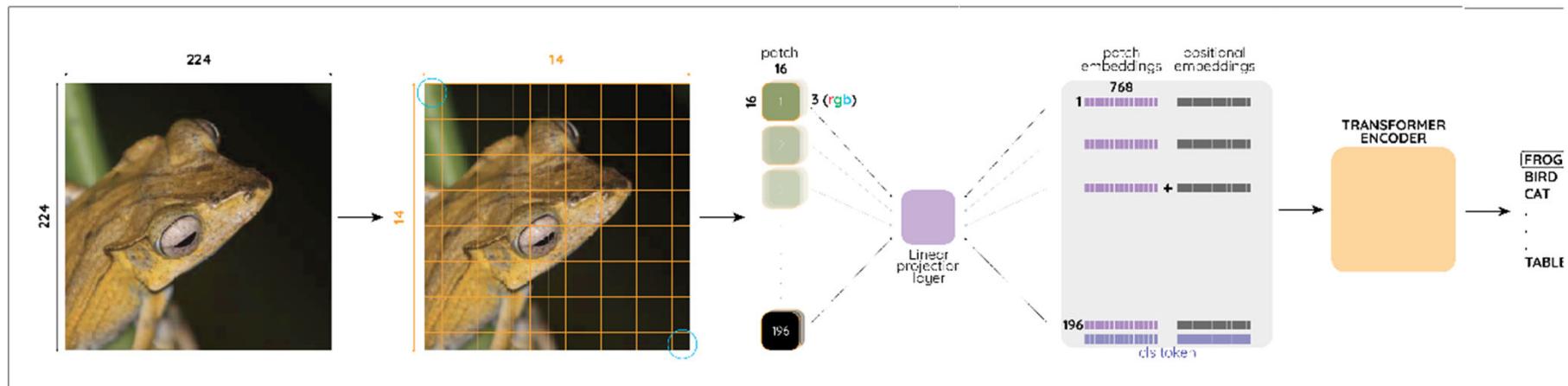


Slide credit: Piotr Mazurek

# CLS (Classification) Token



# Complete ViT



Slide credit: Piotr Mazurek

# VIT Model Variants

Model	Layers	Hidden size $D$	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

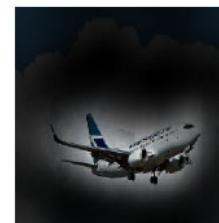
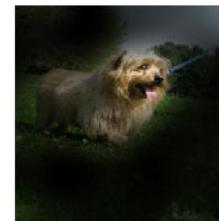
Table 1: Details of Vision Transformer model variants.

# Results

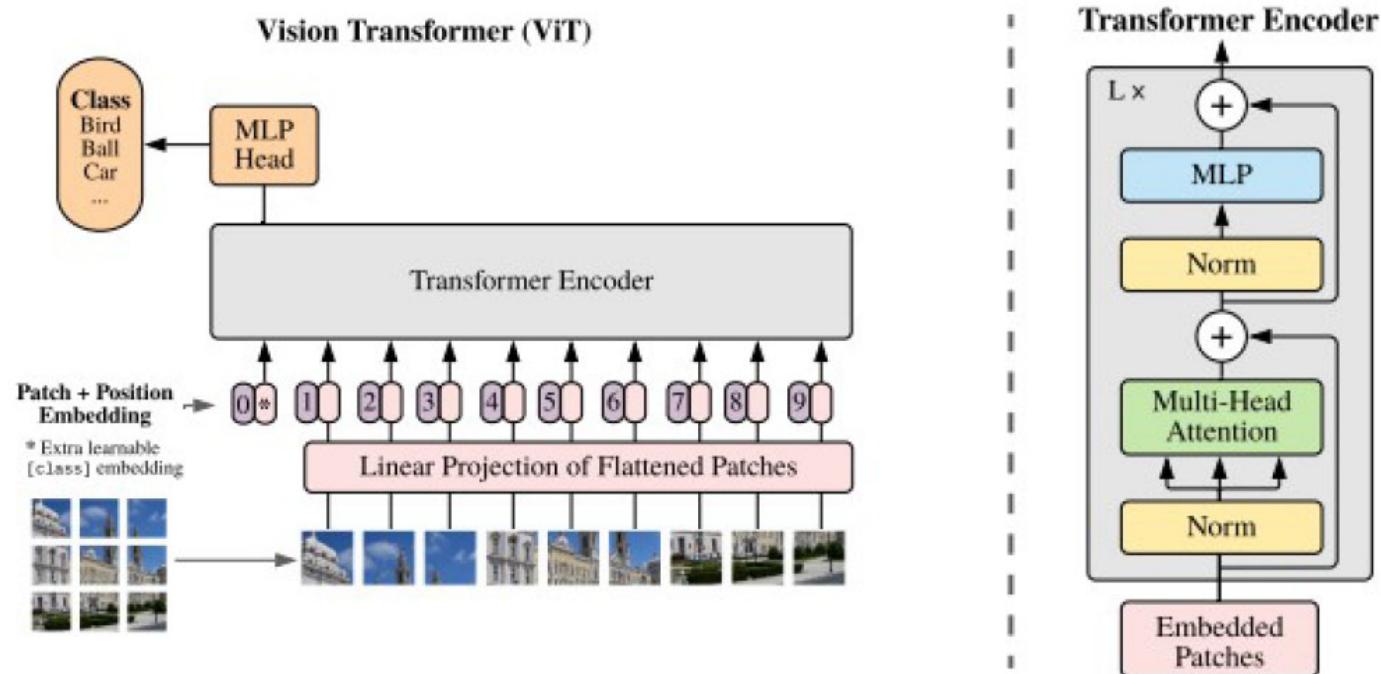
	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	<b>88.55</b> ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4/88.5*
ImageNet ReaL	<b>90.72</b> ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	<b>99.50</b> ± 0.06	99.42 ± 0.03	99.15 ± 0.03	99.37 ± 0.06	—
CIFAR-100	<b>94.55</b> ± 0.04	93.90 ± 0.05	93.25 ± 0.05	93.51 ± 0.08	—
Oxford-IIIT Pets	<b>97.56</b> ± 0.03	97.32 ± 0.11	94.67 ± 0.15	96.62 ± 0.23	—
Oxford Flowers-102	99.68 ± 0.02	<b>99.74</b> ± 0.00	99.61 ± 0.02	99.63 ± 0.03	—
VTAB (19 tasks)	<b>77.63</b> ± 0.23	76.28 ± 0.46	72.72 ± 0.21	76.29 ± 1.70	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

# Attention

Input      Attention



# Vision Transformer (ViT)



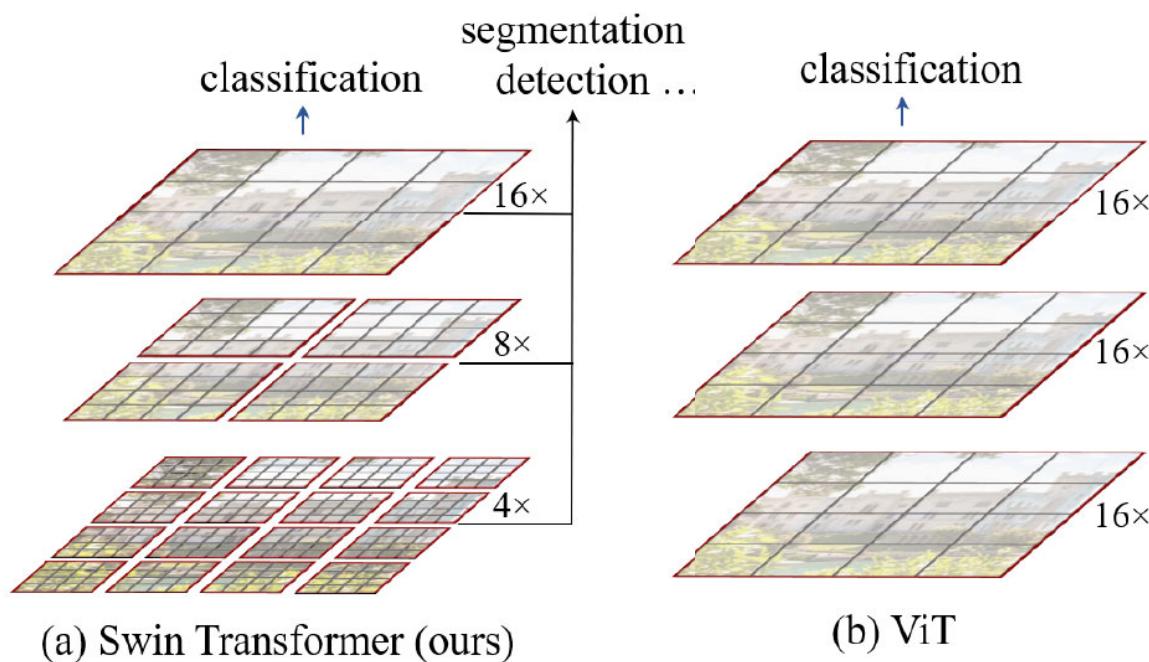
# SWIN

- Z Liu, Y Lin, Y Cao, H Hu, Y Wei, Z Zhang, S Lin, “Swin transformer: Hierarchical vision transformer using shifted windows”, ICCV-2021.  
**(Marr Prize) 13,035 Citations**

# SWIN

- Adapting Transformer from language to vision is challenging
- Differences between language and vision Domains
  - Large variations in the scale of visual entities
  - High resolution of pixels in images compared to words in text
- To address these differences, SWIN proposes
  - A hierarchical Transformer whose representation is computed with Shifted windows
  - Shifted Windowing limit attention
    - To local windows and
    - Allowing cross window connections

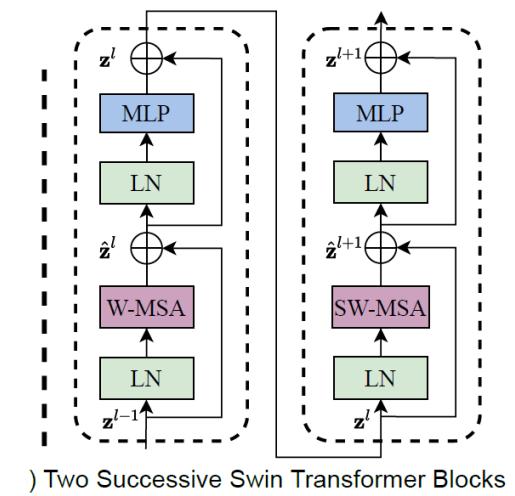
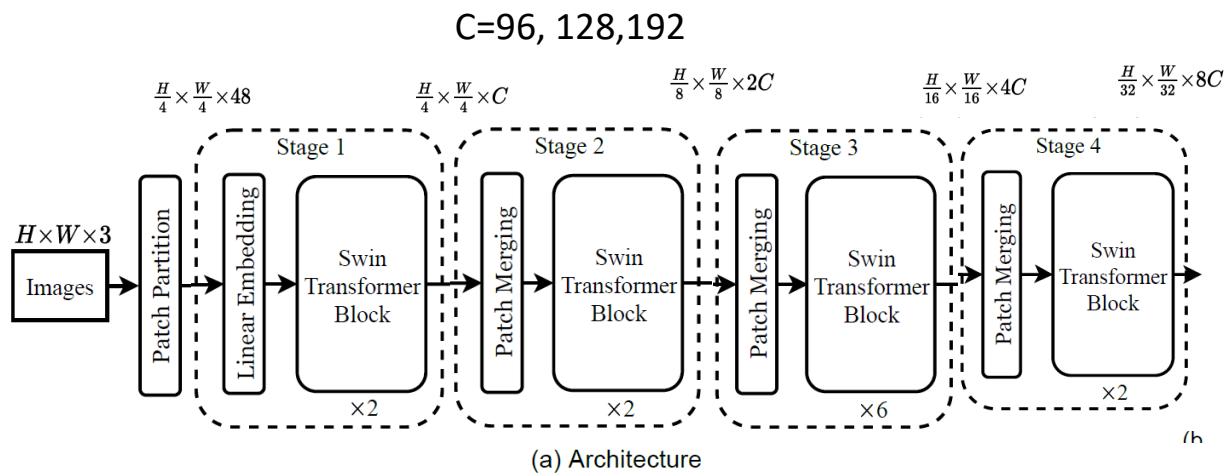
# Hierarchical Feature Maps and Local Attention



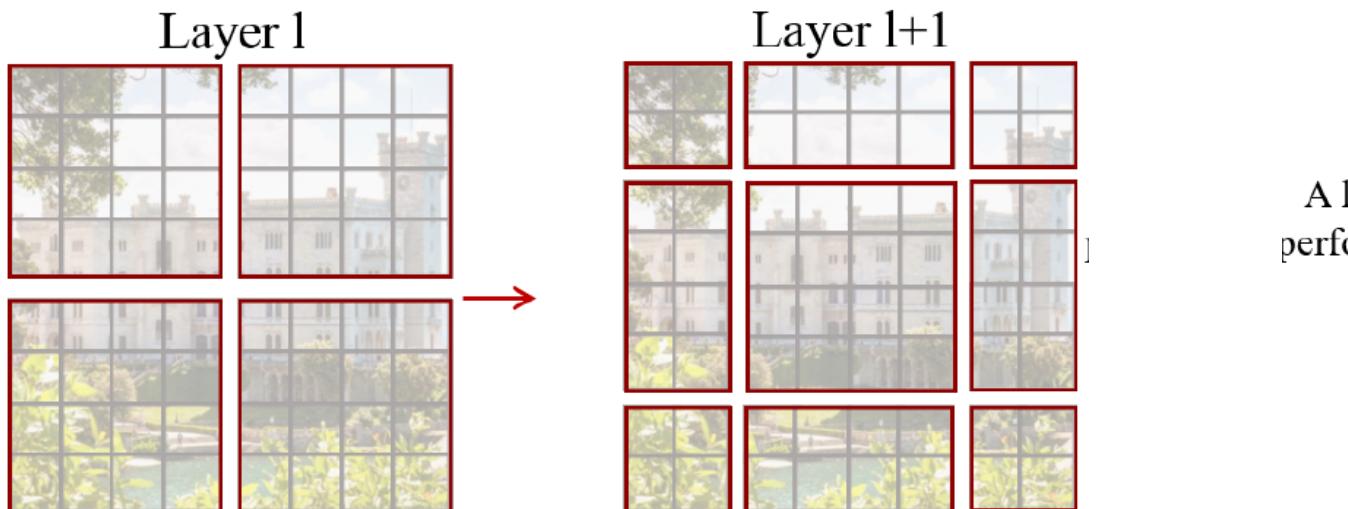
# SWIN

- 87.3 top-1 accuracy on ImageNet-1K
- Dense prediction tasks such as
  - Object detection (58.7 box AP and 51.1 mask AP on COCO)
  - Semantic segmentation (53.5 mIoU on ADE20K )
- Performance surpasses the previous state-of-the art by
  - +2.7 box AP and +2.6 mask AP on COCO, and
  - +3.2 mIoU on ADE20K,

# SWIN Architecture



# Self-Attention within each window and shifted windows



A local window to perform self-attention  
A patch

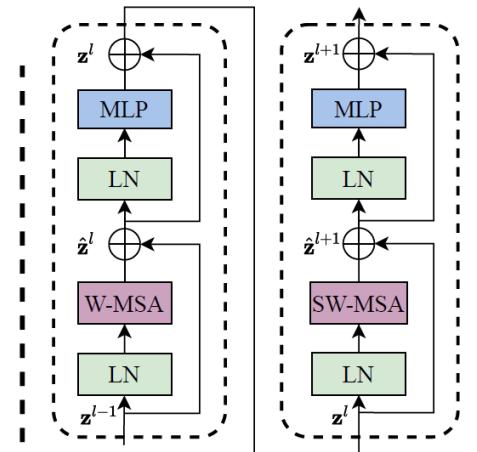
# Transformer Blocks

$$\hat{\mathbf{z}}^l = \text{W-MSA}(\text{LN}(\mathbf{z}^{l-1})) + \mathbf{z}^{l-1},$$

$$\mathbf{z}^l = \text{MLP}(\text{LN}(\hat{\mathbf{z}}^l)) + \hat{\mathbf{z}}^l,$$

$$\hat{\mathbf{z}}^{l+1} = \text{SW-MSA}(\text{LN}(\mathbf{z}^l)) + \mathbf{z}^l,$$

$$\mathbf{z}^{l+1} = \text{MLP}(\text{LN}(\hat{\mathbf{z}}^{l+1})) + \hat{\mathbf{z}}^{l+1},$$



) Two Successive Swin Transformer Blocks

# Different Configurations

- Swin-T:  $C = 96$ , layer numbers =  $\{2, 2, 6, 2\}$
- Swin-S:  $C = 96$ , layer numbers =  $\{2, 2, 18, 2\}$
- Swin-B:  $C = 128$ , layer numbers =  $\{2, 2, 18, 2\}$
- Swin-L:  $C = 192$ , layer numbers =  $\{2, 2, 18, 2\}$

# Results

(a) Regular ImageNet-1K trained models					
method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
RegNetY-4G [44]	224 <sup>2</sup>	21M	4.0G	1156.7	80.0
RegNetY-8G [44]	224 <sup>2</sup>	39M	8.0G	591.6	81.7
RegNetY-16G [44]	224 <sup>2</sup>	84M	16.0G	334.7	82.9
ViT-B/16 [19]	384 <sup>2</sup>	86M	55.4G	85.9	77.9
ViT-L/16 [19]	384 <sup>2</sup>	307M	190.7G	27.3	76.5
DeiT-S [57]	224 <sup>2</sup>	22M	4.6G	940.4	79.8
DeiT-B [57]	224 <sup>2</sup>	86M	17.5G	292.3	81.8
DeiT-B [57]	384 <sup>2</sup>	86M	55.4G	85.9	83.1
Swin-T	224 <sup>2</sup>	29M	4.5G	755.2	81.3
Swin-S	224 <sup>2</sup>	50M	8.7G	436.9	83.0
Swin-B	224 <sup>2</sup>	88M	15.4G	278.1	83.5
Swin-B	384 <sup>2</sup>	88M	47.0G	84.7	84.5

# Results

## (b) ImageNet-22K pre-trained models

method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
R-101x3 [34]	$384^2$	388M	204.6G	-	84.4
R-152x4 [34]	$480^2$	937M	840.5G	-	85.4
ViT-B/16 [19]	$384^2$	86M	55.4G	85.9	84.0
ViT-L/16 [19]	$384^2$	307M	190.7G	27.3	85.2
Swin-B	$224^2$	88M	15.4G	278.1	85.2
Swin-B	$384^2$	88M	47.0G	84.7	86.4
Swin-L	$384^2$	197M	103.9G	42.1	87.3

# Results

(b) Various backbones w. Cascade Mask R-CNN										
	AP <sup>box</sup>	AP <sup>box</sup> <sub>50</sub>	AP <sup>box</sup> <sub>75</sub>	AP <sup>mask</sup>	AP <sup>mask</sup> <sub>50</sub>	AP <sup>mask</sup> <sub>75</sub>	param	FLOPs	FPS	
DeiT-S <sup>†</sup>	48.0	67.2	51.7	41.4	64.2	44.3	80M	889G	10.4	
R50	46.3	64.3	50.5	40.1	61.7	43.4	82M	739G	18.0	
Swin-T	<b>50.5</b>	<b>69.3</b>	<b>54.9</b>	<b>43.7</b>	<b>66.6</b>	<b>47.1</b>	86M	745G	15.3	
X101-32	48.1	66.5	52.4	41.6	63.9	45.2	101M	819G	12.8	
Swin-S	<b>51.8</b>	<b>70.4</b>	<b>56.3</b>	<b>44.7</b>	<b>67.9</b>	<b>48.5</b>	107M	838G	12.0	
X101-64	48.3	66.4	52.3	41.7	64.0	45.1	140M	972G	10.4	
Swin-B	<b>51.9</b>	<b>70.9</b>	<b>56.5</b>	<b>45.0</b>	<b>68.4</b>	<b>48.7</b>	145M	982G	11.6	

Method	Backbone	val mIoU	test score	#param.	FLOPs	FPS
DLab.v3+ [11]	ResNet-101	44.1	-	63M	1021G	16.0
DNL [65]	ResNet-101	46.0	56.2	69M	1249G	14.8
OCRNet [67]	ResNet-101	45.3	56.0	56M	923G	19.3
UperNet [63]	ResNet-101	44.9	-	86M	1029G	20.1
OCRNet [67]	HRNet-w48	45.7	-	71M	664G	12.5
DLab.v3+ [11]	ResNeSt-101	46.9	55.1	66M	1051G	11.9
DLab.v3+ [11]	ResNeSt-200	48.4	-	88M	1381G	8.1
SETR [73]	T-Large <sup>‡</sup>	50.3	61.7	308M	-	-
UperNet	DeiT-S <sup>†</sup>	44.0	-	52M	1099G	16.2
UperNet	Swin-T	46.1	-	60M	945G	18.5
UperNet	Swin-S	49.3	-	81M	1038G	15.2
UperNet	Swin-B <sup>‡</sup>	51.6	-	121M	1841G	8.7
UperNet	Swin-L <sup>‡</sup>	<b>53.5</b>	<b>62.8</b>	234M	3230G	6.2

Table 3. Results of semantic segmentation on the ADE20K val

# Summary

- ViT is first Vision Transformer, but trained on huge dataset of 300M
- SWIN employs window attention
- Performs well on other tasks: object detection, semantic segmentation