

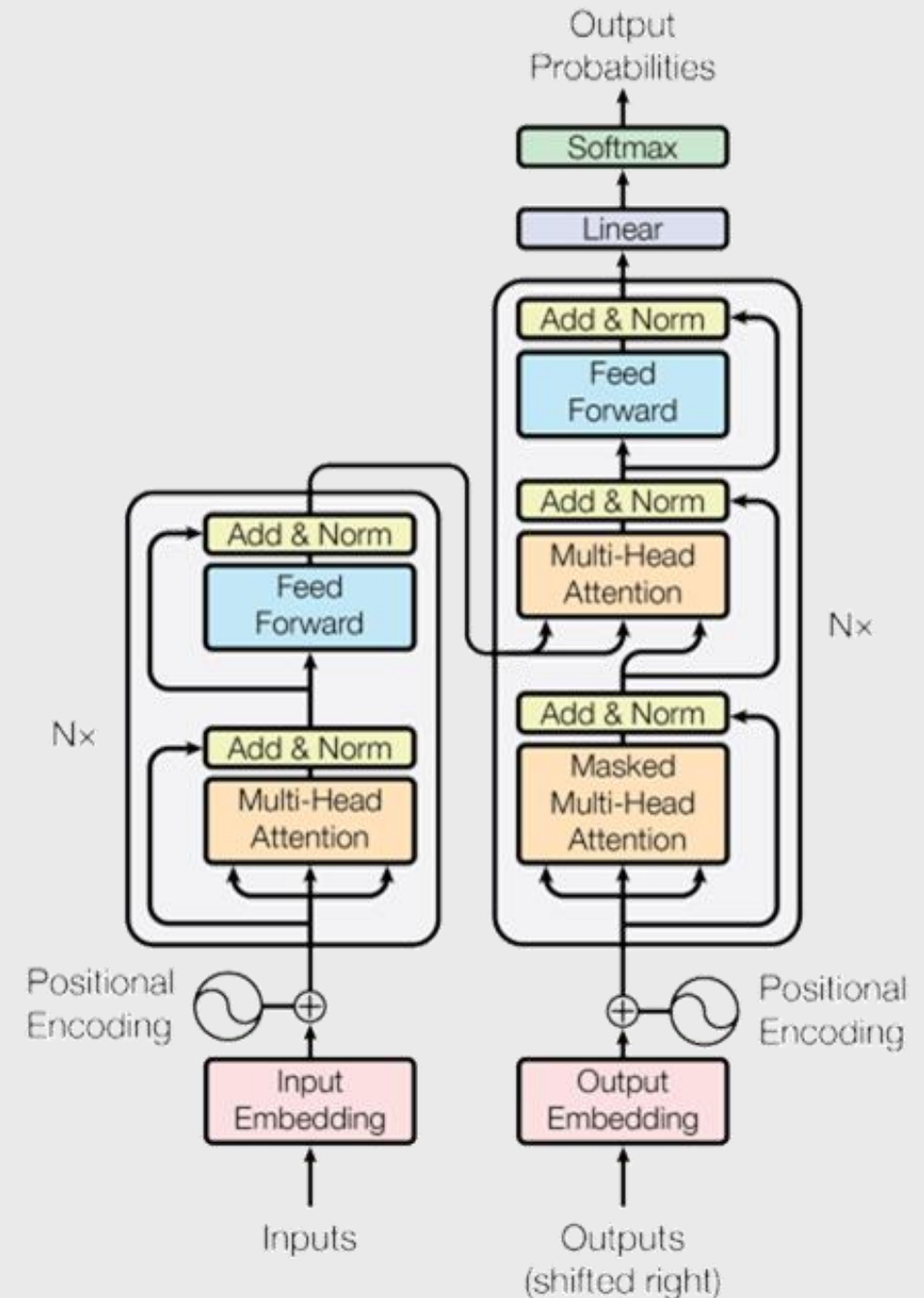
(Vision) Transformers

Computer Vision and Artificial
Intelligence

Dr Varun Ojha

varun.ojha@ncl.ac.uk

School of Computing
Newcastle University



Learning objectives (Vision Transformer)

By the end of this week, you will be able to:

- Learn the concepts of Transformer Models
- Understand the Self-Attention Mechanisms (the basic building block of Transformers)
- Understand an image classification using Vision Transformers

Content of this week (ViT)

- **Part 1: Self-Attention Block**
 - Basic concept of Transformers
 - Word Embedding
 - Self-Attention Mechanism
- **Part 2: Vision Transformer**
 - Basic concept of Vision Transformers
 - Architecture of a Vision Transformer
 - Performance of Vision Transformer

Word Embedding

Given a word W (e.g. “intelligence”) we want W to be a real vector of dimension d . Dimension d is also called word embedding dimension.

$$\mathbf{W}: \text{words} \rightarrow \mathbb{R}^d$$

$$\text{“intelligence”} \rightarrow (w_1, w_2, \dots, w_d) \rightarrow (0.1, -0.8, \dots, 0.9)$$

Word Embedding

W_1 = “I love artificial intelligence”

W_2 = “I like computational intelligence”

We create a vocabulary V collecting all unique words.

$V = \{\text{“I”, “love”, “like”, “artificial”, “computational”, “intelligence”}\}$

For this example, vocabulary size $|V| = 6$

Word Embedding: Word \rightarrow Integer

$V = \{\text{"I"}, \text{"love"}, \text{"artificial"}, \text{"computational"}, \text{"intelligence"}, \text{"like"}\}$

I \rightarrow 0

love \rightarrow 1

like \rightarrow 2

artificial \rightarrow 3

computational \rightarrow 4

intelligence \rightarrow 5

Word Embedding: Integer \rightarrow Word

$V = \{\text{"I"}, \text{"love"}, \text{"artificial"}, \text{"computational"}, \text{"intelligence"}, \text{"like"}\}$

$0 \rightarrow \text{I}$

$1 \rightarrow \text{love}$

$2 \rightarrow \text{like}$

$3 \rightarrow \text{artificial}$

$4 \rightarrow \text{computational}$

$5 \rightarrow \text{intelligence}$

One-Hot Encoding

$V = \{\text{"I"}, \text{"love"}, \text{"like"}, \text{"artificial"}, \text{"computational"}, \text{"intelligence"}\}$

$$\text{"I"} = \begin{bmatrix} \mathbf{1} \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{"love"} = \begin{bmatrix} 0 \\ \mathbf{1} \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

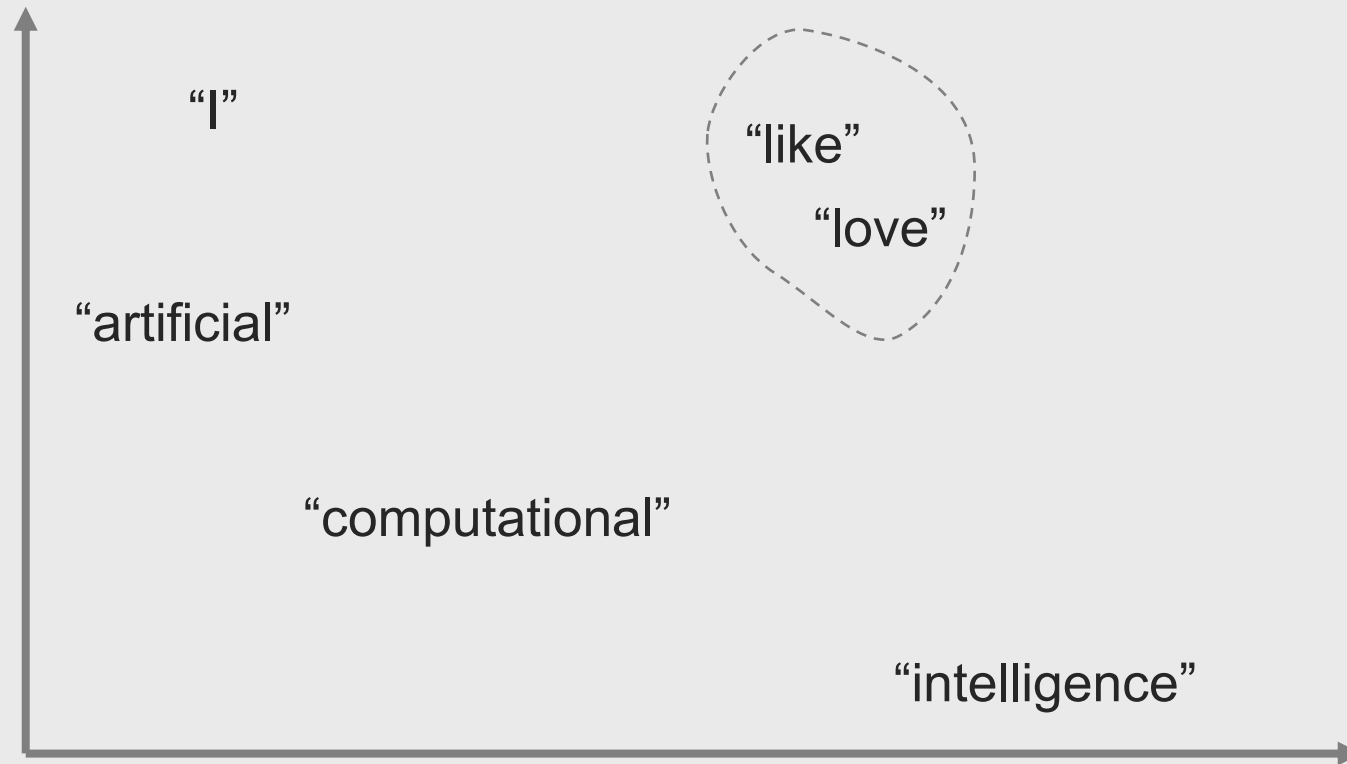
$$\text{"like"} = \begin{bmatrix} 0 \\ 0 \\ \mathbf{1} \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{"artificial"} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \mathbf{1} \\ 0 \\ 0 \end{bmatrix}$$

$$\text{"computational"} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \mathbf{1} \\ 0 \end{bmatrix}$$

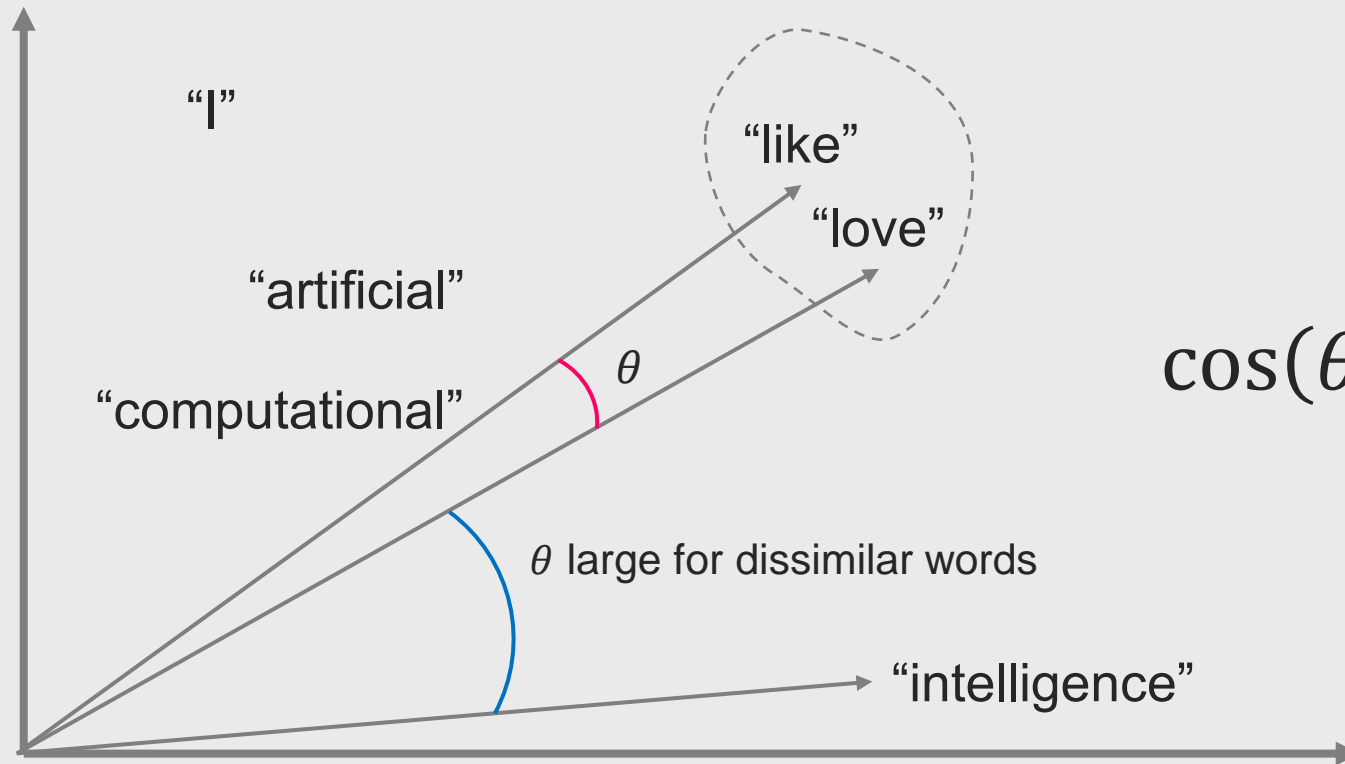
$$\text{"intelligence"} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \mathbf{1} \end{bmatrix}$$

Similarity between words?



Objective is to place similar words close to each other

Similarity between words?



$$\cos(\theta) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

Objective is to place similar words close to each other

t-SNE visualisation of words

Turian *et al.* (2010)

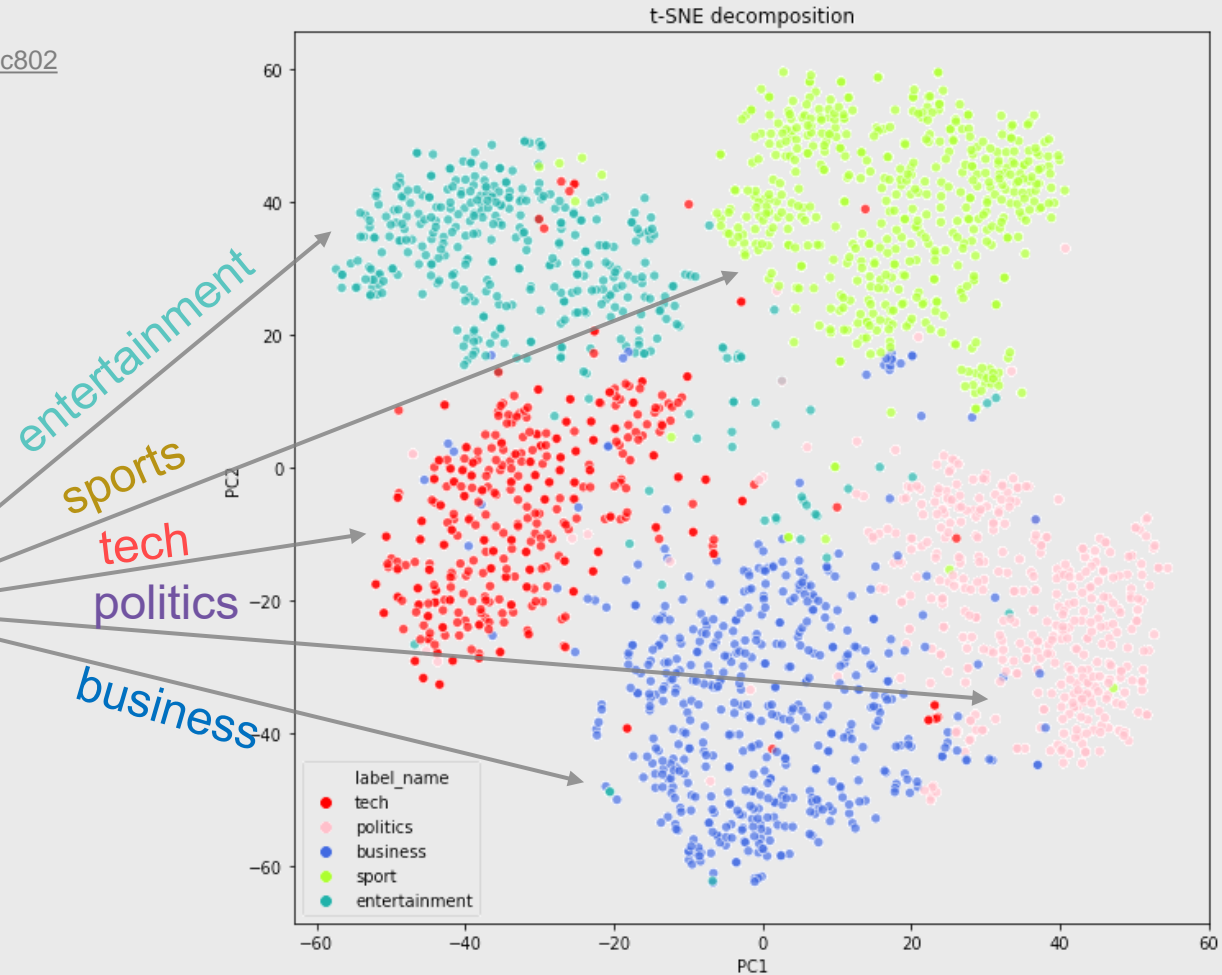


t-SNE visualisation

Example Source:

<https://towardsdatascience.com/text-classification-in-python-dd95d264c802>

All similar topics
are closer to
each other



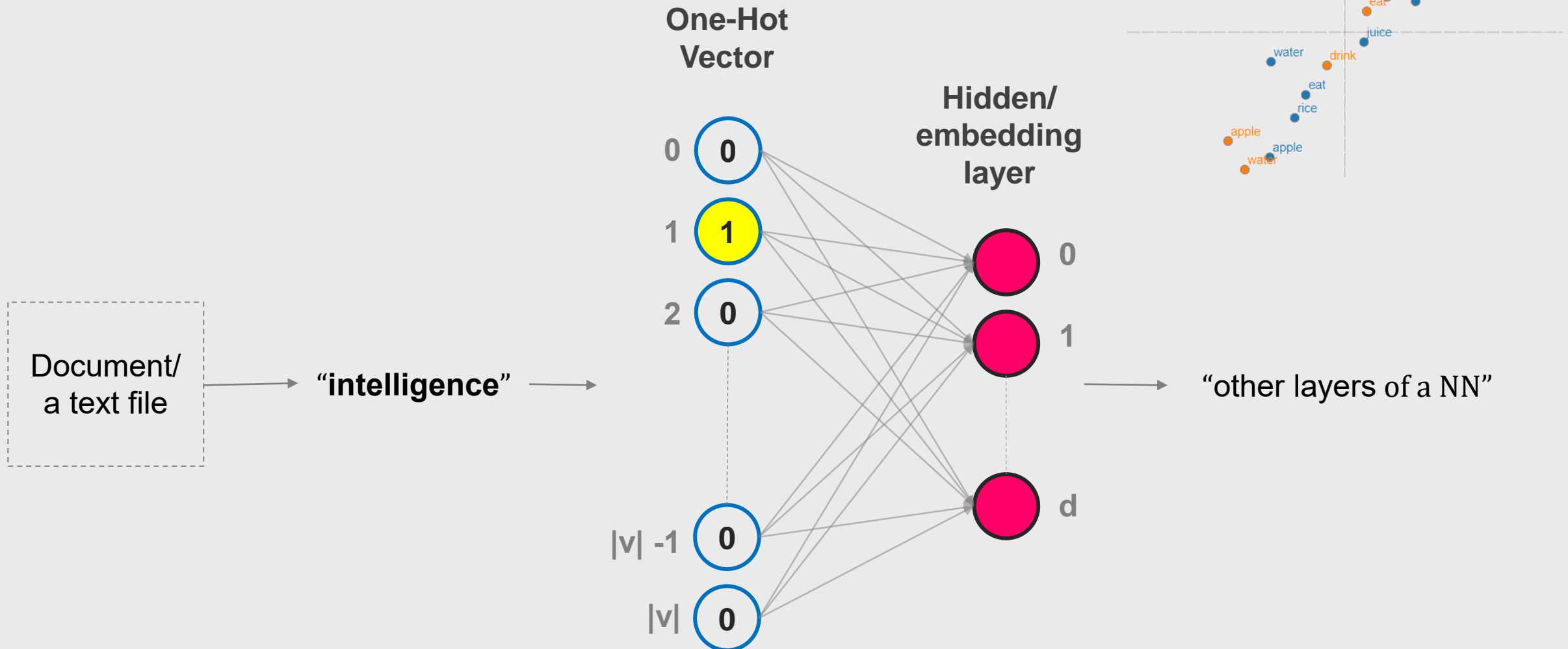
Word Embedding: Objective

Given a word W (e.g. “intelligence”) we want to W a real vector of dimension d

$$\mathbf{W}: \text{words} \rightarrow \mathbb{R}^d$$

$$\text{“intelligence”} \rightarrow (w_1, w_2, \dots, w_d) \rightarrow (0.1, -0.8, \dots, 0.9)$$

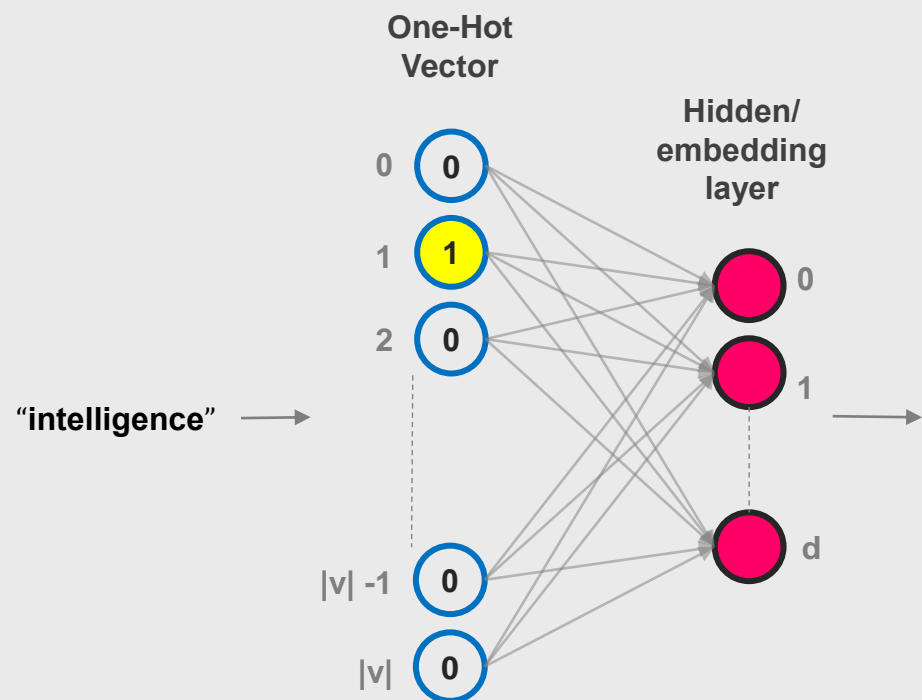
Word Embedding



Check online here: <https://ronxin.github.io/wevi/>

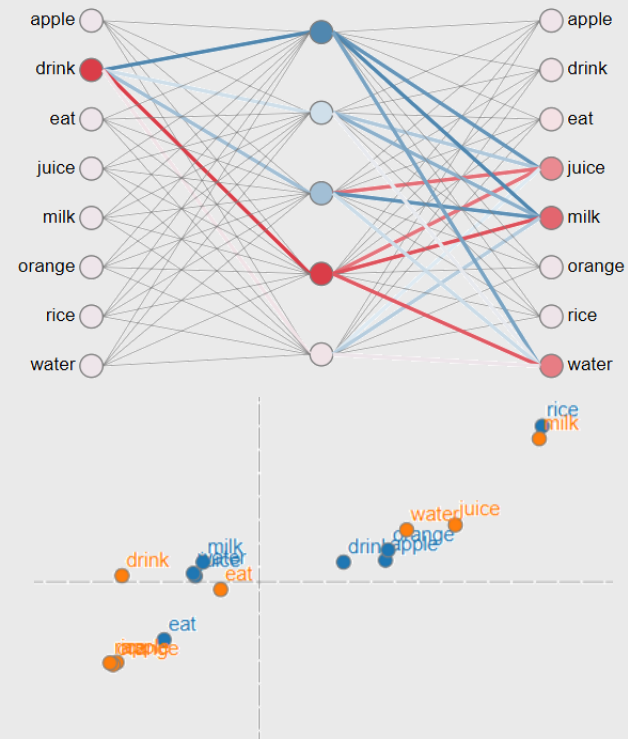
Word Embedding

Check online here: <https://ronxin.github.io/wevi/>



$$\text{“intelligence”} = [0 \quad 1 \quad 0 \quad \dots \quad 0 \quad 0] \times \begin{bmatrix} 0.5 & 4.6 & \dots & 0.9 \\ 0.1 & -0.8 & \dots & 0.7 \\ \vdots & \vdots & \ddots & \vdots \\ 0.6 & 0.8 & \dots & 0.3 \\ 0.3 & -0.6 & \dots & -0.8 \\ -0.5 & 0.5 & \dots & 0.1 \end{bmatrix}$$

$$(1 \times |V|) \cdot (|V| \times d) \Rightarrow (1 \times d)$$



Word Embedding

$$\begin{aligned} \text{"intelligence"} &= [0 \quad 1 \quad 0 \quad \dots \quad 0 \quad 0] \times \begin{bmatrix} 0.5 & 4.6 & \dots & 0.7 \\ 0.1 & -0.8 & \dots & 0.9 \\ 0.6 & 0.8 & \dots & 0.3 \\ 0.3 & -0.6 & \dots & -0.8 \\ -0.5 & 0.5 & \dots & 0.1 \end{bmatrix} \\ &= [0.1, -0.8, \dots, 0.9] \end{aligned}$$

Word Embedding

← Embedding dimension d →

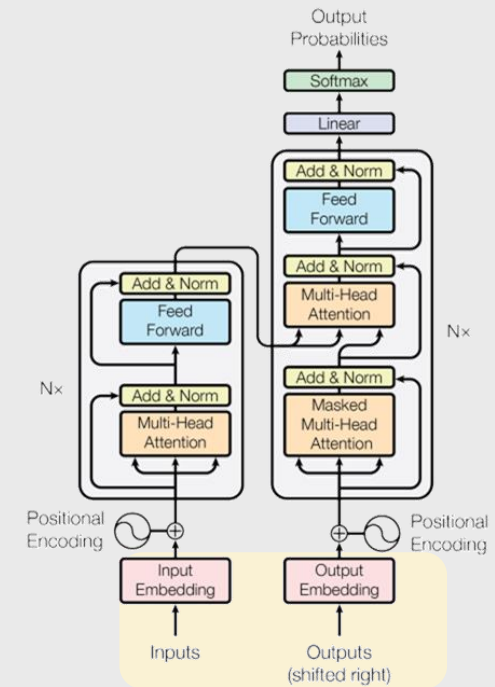
“intelligence”

Input token

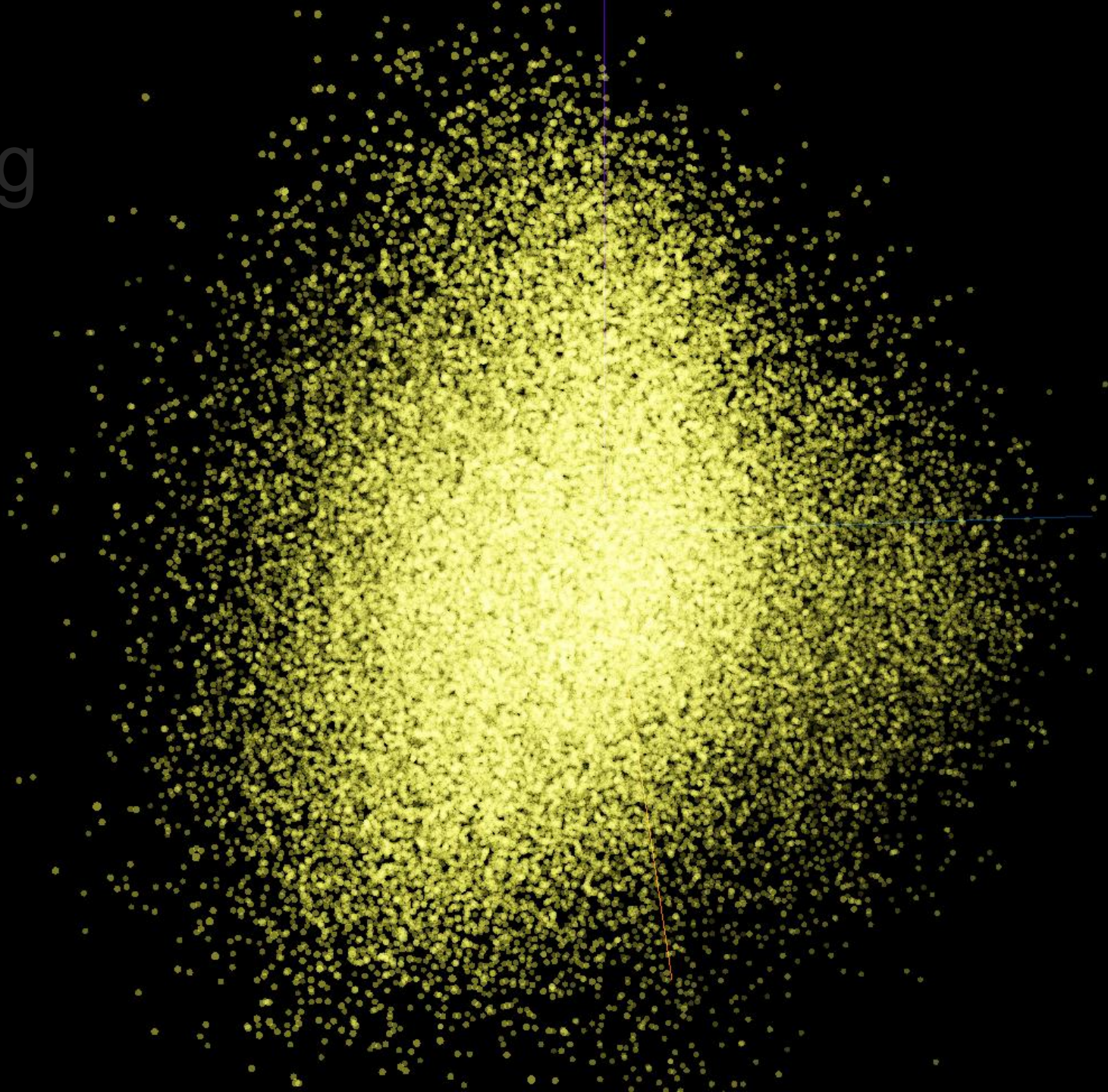
0.5	4.6	...	0.9
0.1	-0.8	...	0.7
0.6	0.8	...	0.3
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.
0.3	-0.6	...	-0.8
-0.5	0.5	...	0.1

Lookup Table
Embedding Weight Matrix

↑ Vocabulary size ↓



Word Embedding



Positional Encoding

Let's have a sentence "This is Computer Vision Class" of $n = 5$ sequence length

And each word x_t (e.g., "Computer") is represented by an embedding vector of size for example $d = 10$ (this could be very large number)

That is mathematically t-th word is represented as

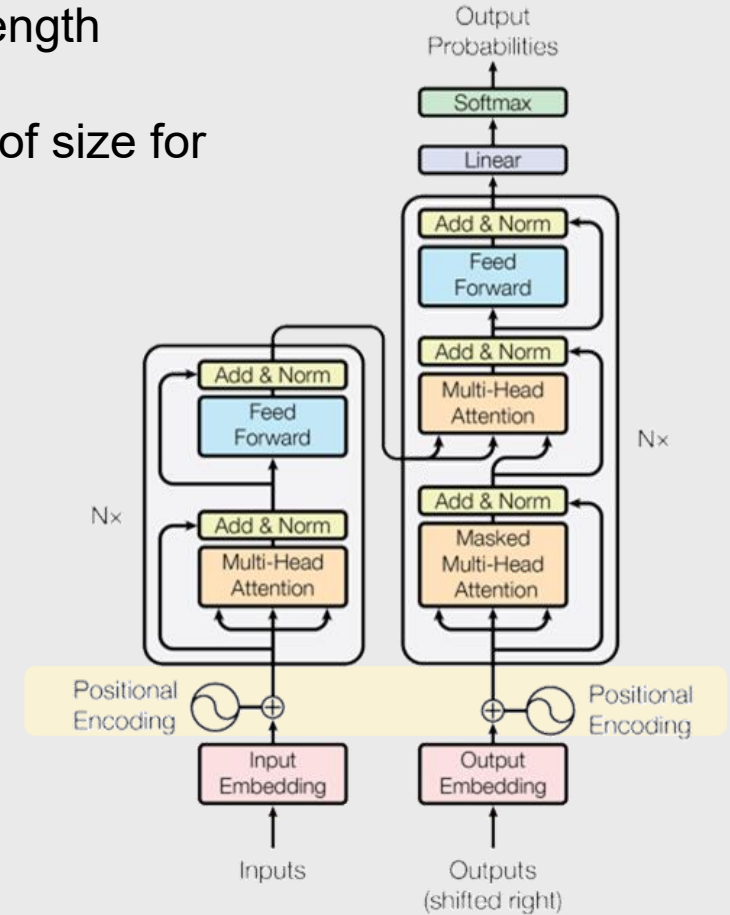
$$x_t \in \mathbb{R}^d$$

Then the positional encoding will be presented as:

$$p(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$

$$p(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$

For $pos = 0, 1, \dots, n$ and $i = 0, 1, \dots, \frac{d}{2}$



Positional Encoding of word: $x_t \in \mathbb{R}^d$

It assign a value relevant to the position of the word in the sentence

$$p_{pos} = \begin{bmatrix} \sin(\omega_1) \\ \cos(\omega_1) \\ \sin(\omega_2) \\ \cos(\omega_2) \\ \vdots \\ \vdots \\ \sin(\omega_{d/2}) \\ \cos(\omega_{d/2}) \end{bmatrix}$$

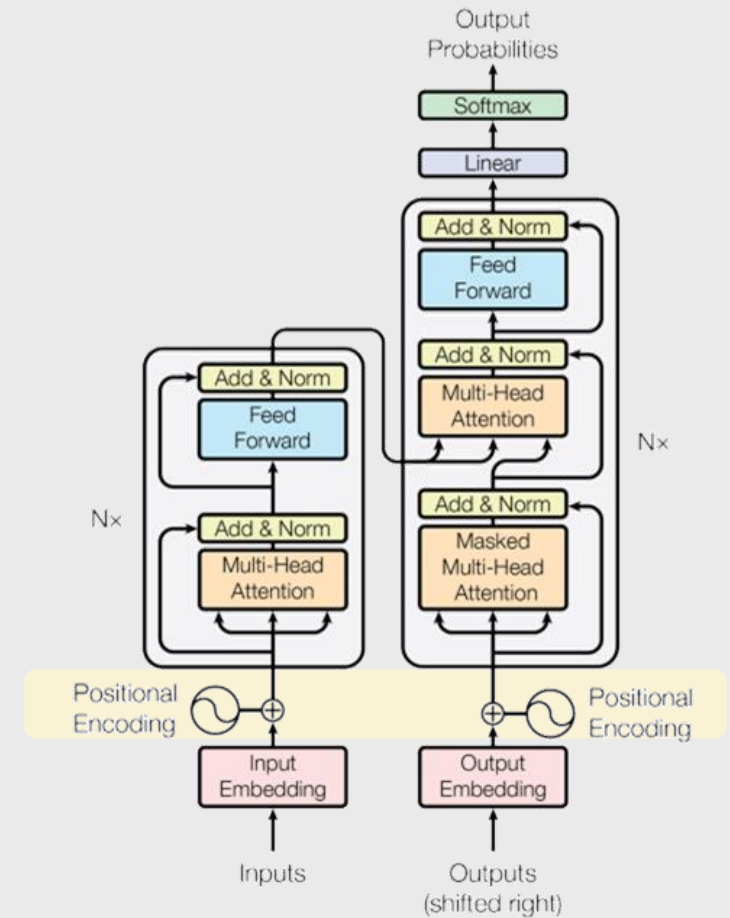
Alternating function and values

where

$$\omega_t = \frac{pos}{10000^{2i/d}}$$

$$P = \begin{bmatrix} p_0 \\ p_1 \\ \vdots \\ p_n \end{bmatrix} \quad n \times d$$

Positional encoding matrix

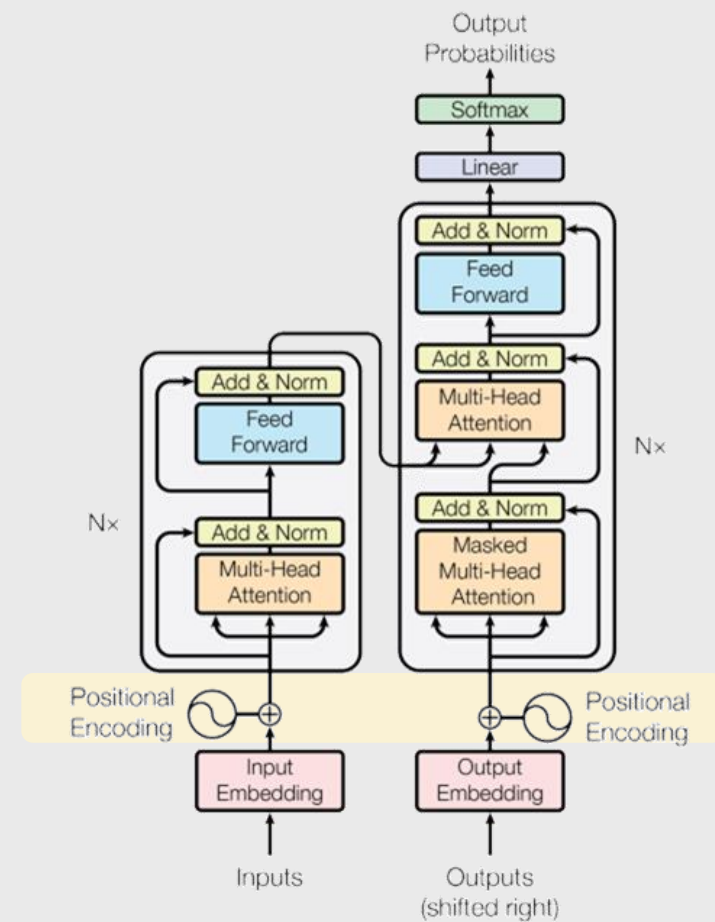
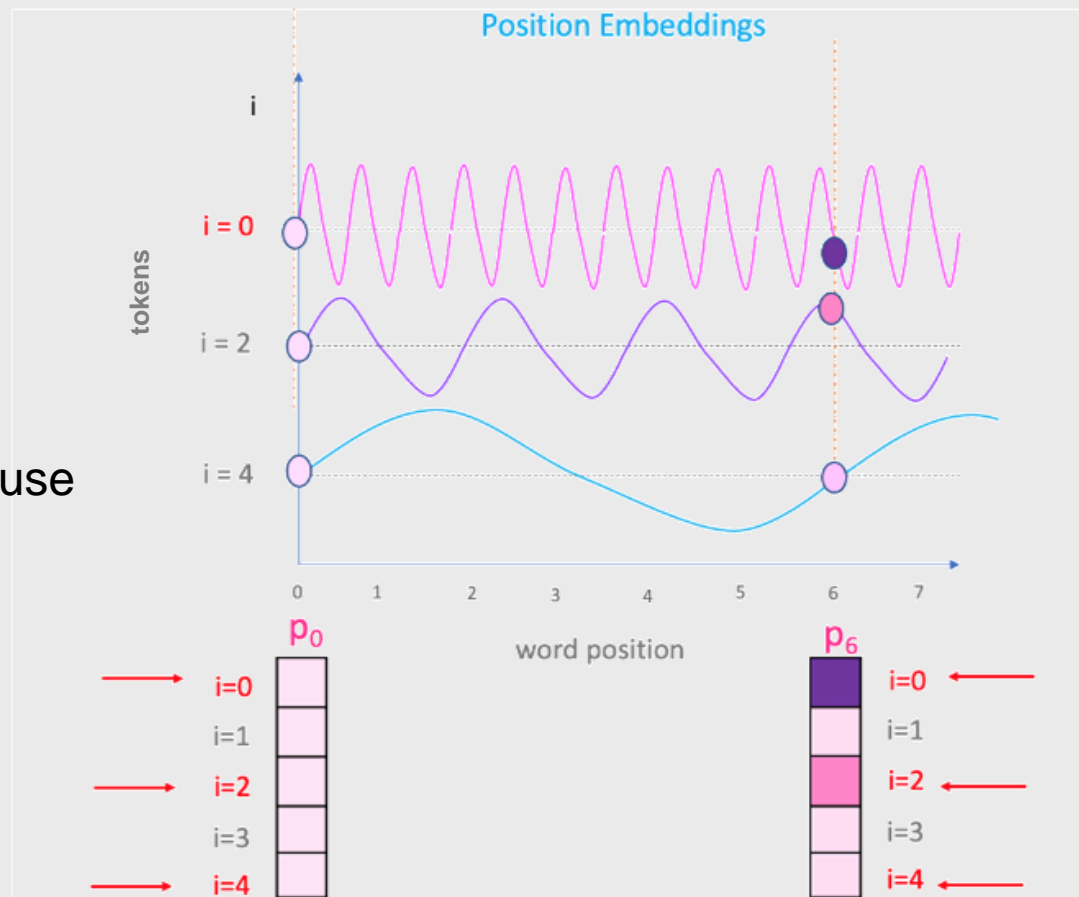


Positional Encoding of word: $x_t \in \mathbb{R}^d$

It assign a value relevant to the position of the word in the sentence

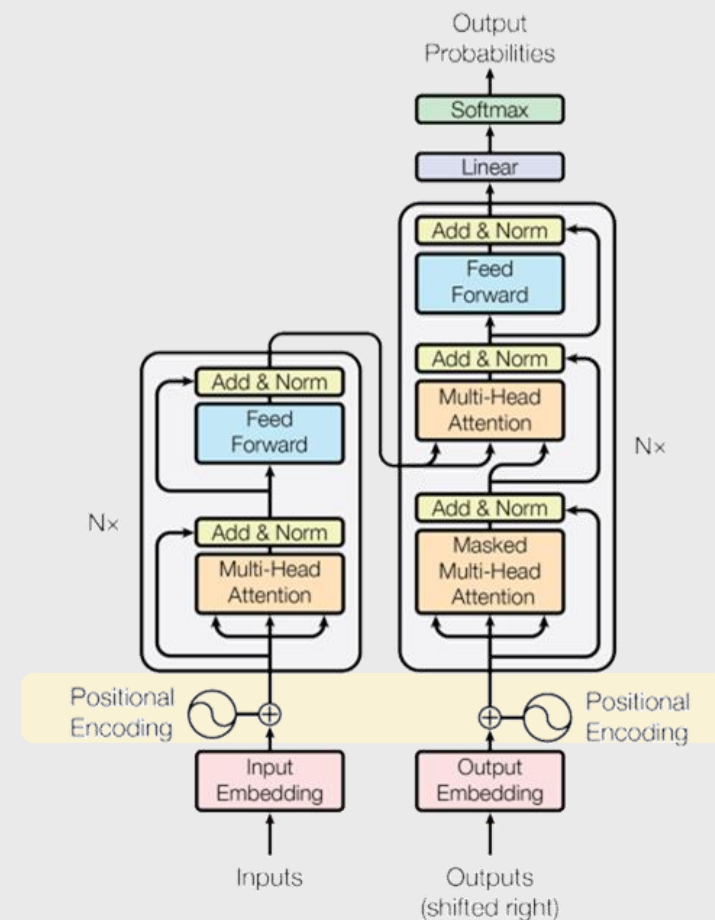
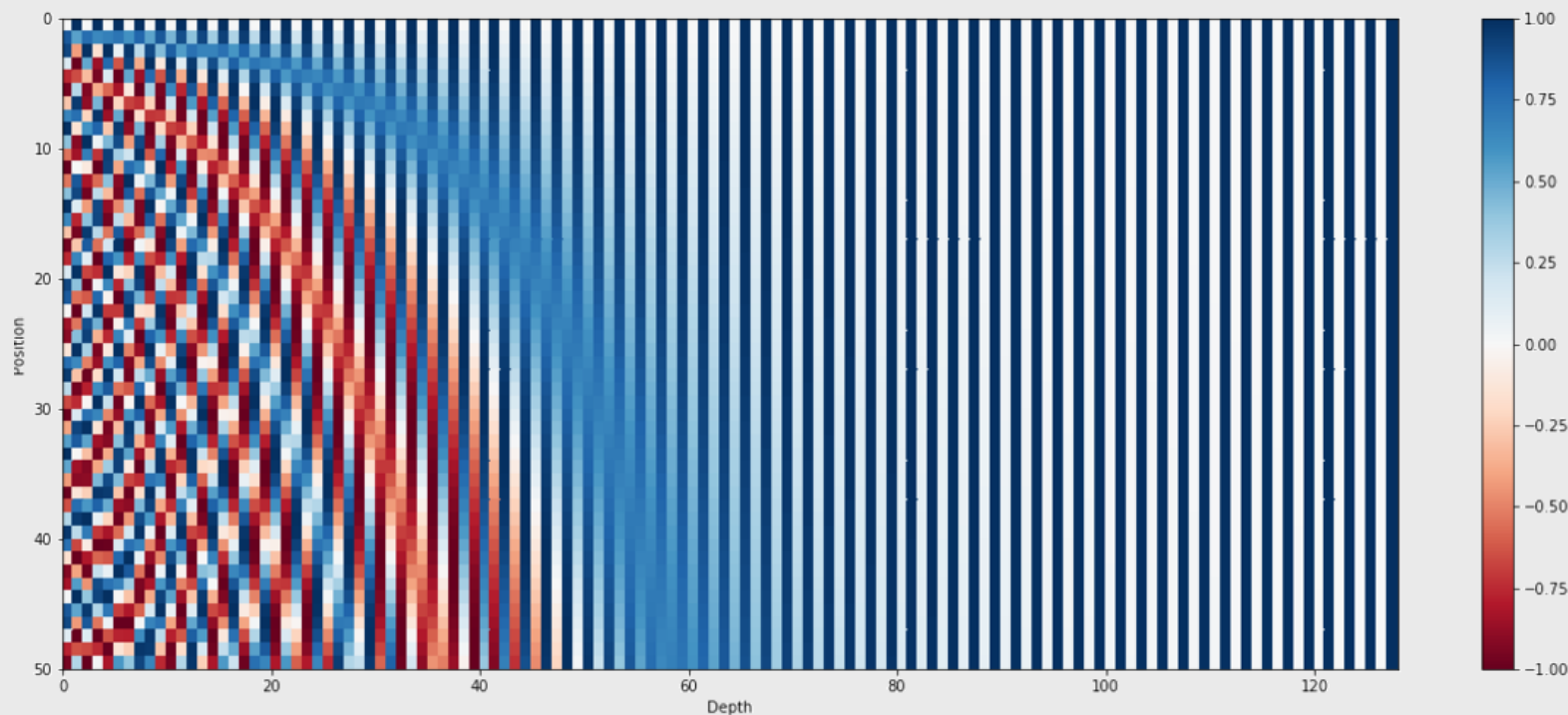
$$\sin(\omega_1 \cdot t)$$

only Sinusoidal
function here because
 $i = 2k$

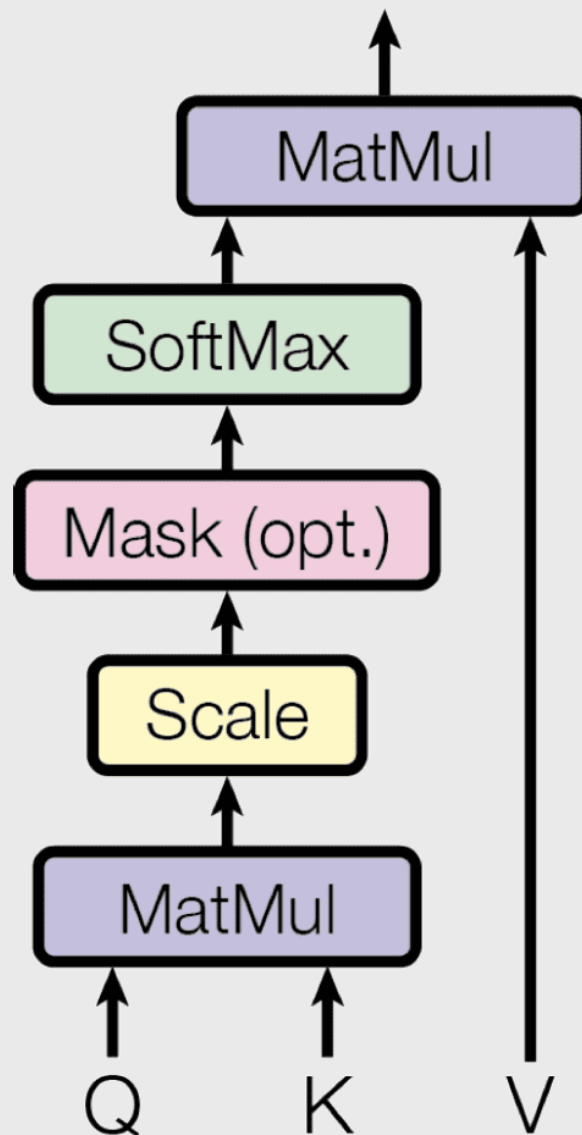


Positional Encoding of word: $x_t \in \mathbb{R}^d$

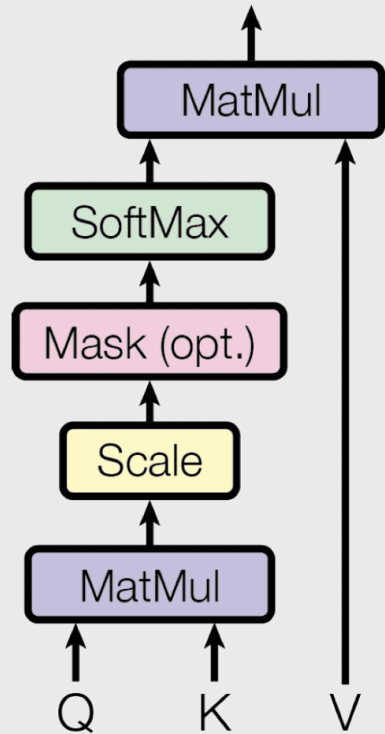
It assign a value relevant to the position of the word in the sentence



Self-Attention



Self-Attention



$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

Self-Attention

Let's have a sentence "This is Computer Vision Class" of $n = 5$ sequence length

And each word x (e.g., "Computer") is represented by an embedding vector of size for example $d = 10$ (this could be very large number)

That is mathematically a word is presented as

$$x^j \in \mathbb{R}^d$$

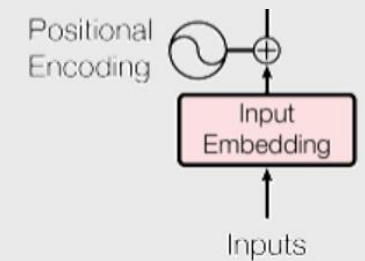
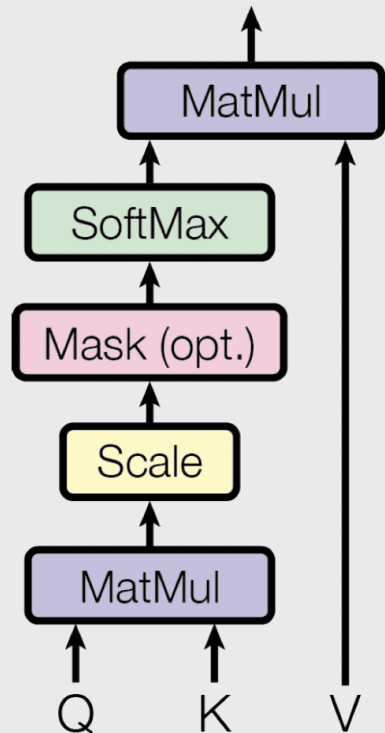
And we have weight matrices

$$\text{Query } \mathbf{W}_Q = d \times d_k, \quad \text{Key } \mathbf{W}_K = d \times d_k, \quad \text{Value } \mathbf{W}_V = d \times d_v$$

Then we perform a linear transformation of the input of x^j via Query, Key and Value matrices to obtain Query, Key and Value vectors as:

$$q_i^{1 \times d_k} = x_i^{1 \times d} \times \mathbf{W}_Q^{d \times d_k}, \quad k_i^{1 \times d_k} = x_i^{1 \times d} \times \mathbf{W}_K^{d \times d_k}, \quad \text{and} \quad v_i^{1 \times d_v} = x_i^{1 \times d} \times \mathbf{W}_V^{d \times d_v}$$

For all words $i = 1, \dots, n$ in the sentence.

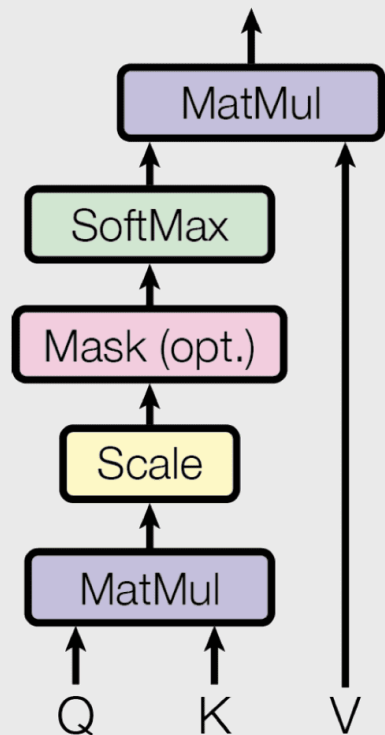


Self-Attention

We can pack the following Query, Key and Value vectors into a matrix forms:

$$\overset{1 \times 2}{q_i^{1 \times d_k}} = x_i^{1 \times d} \times \overset{d \times d_k}{W_Q}, \quad \overset{1 \times 2}{k_i^{1 \times d_k}} = x_i^{1 \times d} \times \overset{d \times d_k}{W_K}, \quad \text{and} \quad \overset{1 \times 2}{v_i^{1 \times d_k}} = x_i^{1 \times d} \times \overset{d \times d_k}{W_V}$$

For all words $i = 1, \dots, n$ in the sentence.

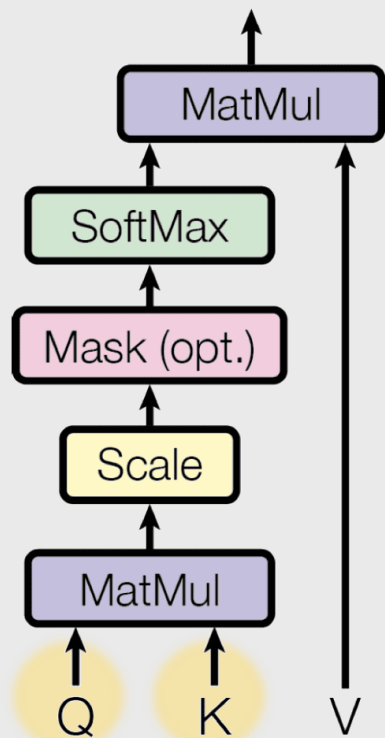


$$Q = \begin{bmatrix} q_1 & q_2 & \dots & q_n \\ 1 & 2 & & 5 \end{bmatrix} \quad K = \begin{bmatrix} k_1 & k_2 & \dots & k_n \\ 1 & 2 & & 5 \end{bmatrix} \quad V = \begin{bmatrix} v_1 & v_2 & \dots & v_n \\ 1 & 2 & & 5 \end{bmatrix}$$

Self-Attention

Vaswani et al. Attention Is All You Need (NIPS 2017)

n is the number of tokens in a sentence



Q is a matrix of size $n \times d_k$

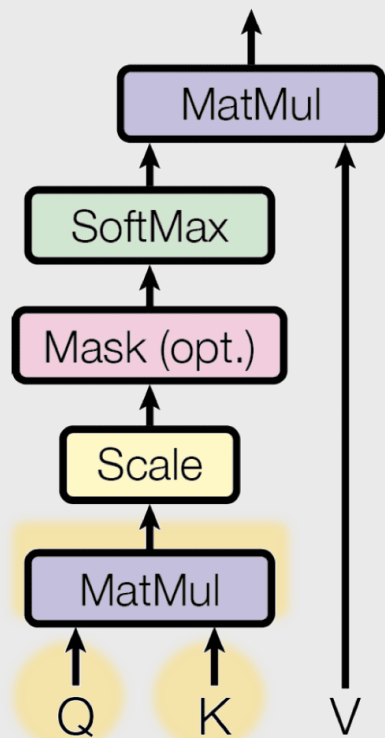
$$Q = \begin{bmatrix} q_{1,1} & \cdots & q_{n,1} \\ q_{1,2} & \ddots & q_{n,2} \\ \vdots & & \vdots \\ q_{1,d_k} & \cdots & q_{n,d_k} \end{bmatrix}_{5 \times 2}$$

$$\text{Attention} (Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

$$K = \begin{bmatrix} k_{1,1} & \cdots & k_{n,1} \\ k_{1,2} & \ddots & k_{n,2} \\ \vdots & & \vdots \\ k_{1,d_k} & \cdots & k_{n,d_k} \end{bmatrix}_{5 \times 2}$$

K is a matrix of size $n \times d_k$

Self-Attention



\mathbf{Q} is a matrix of size $n \times d_k$

$$\mathbf{Q} = \begin{bmatrix} q_{1,1} & \cdots & q_{1,d_k} \\ q_{2,1} & \cdots & q_{2,d_k} \\ \vdots & \ddots & \vdots \\ q_{n,1} & \cdots & q_{n,d_k} \end{bmatrix}_{5 \times 2}$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

$$\mathbf{K} = \begin{bmatrix} k_{1,1} & \cdots & k_{1,d_k} \\ k_{2,1} & \cdots & k_{2,d_k} \\ \vdots & \ddots & \vdots \\ k_{n,1} & \cdots & k_{n,d_k} \end{bmatrix}_{5 \times 2}$$

\mathbf{K} is a matrix of size $n \times d_k$

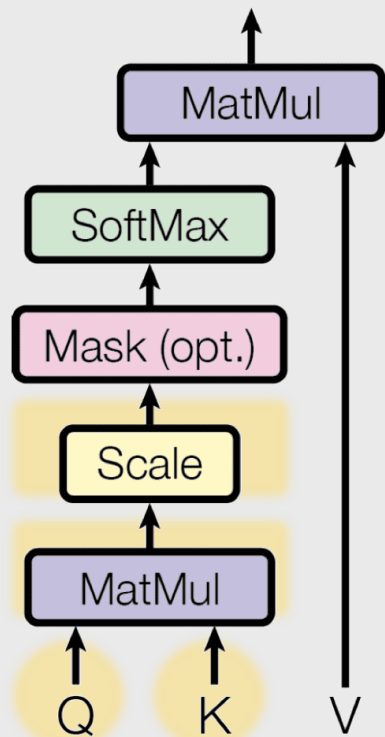
$$\begin{bmatrix} e_{1,1} & e_{1,2} & \cdots & e_{1,n} \\ e_{2,1} & e_{2,2} & \cdots & e_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ e_{n,1} & e_{n,2} & \cdots & e_{n,n} \end{bmatrix}_{5 \times 5}$$

alignment scores

$$\begin{bmatrix} q_{1,1} & \cdots & q_{1,d_k} \\ q_{2,1} & \cdots & q_{2,d_k} \\ \vdots & \ddots & \vdots \\ q_{n,1} & \cdots & q_{n,d_k} \end{bmatrix}_{5 \times 2} \cdot \begin{bmatrix} k_{1,1} & \cdots & k_{1,d_k} \\ k_{2,1} & \cdots & k_{2,d_k} \\ \vdots & \ddots & \vdots \\ k_{n,1} & \cdots & k_{n,d_k} \end{bmatrix}_{2 \times 5}$$

dot product

Self-Attention



\mathbf{Q} is a matrix of size $n \times d_k$

$$\mathbf{Q} = \begin{bmatrix} q_{1,1} & \cdots & q_{n,1} \\ q_{1,2} & \ddots & q_{n,2} \\ \vdots & & \vdots \\ q_{1,d_k} & \cdots & q_{n,d_k} \end{bmatrix}_{5 \times 2}$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

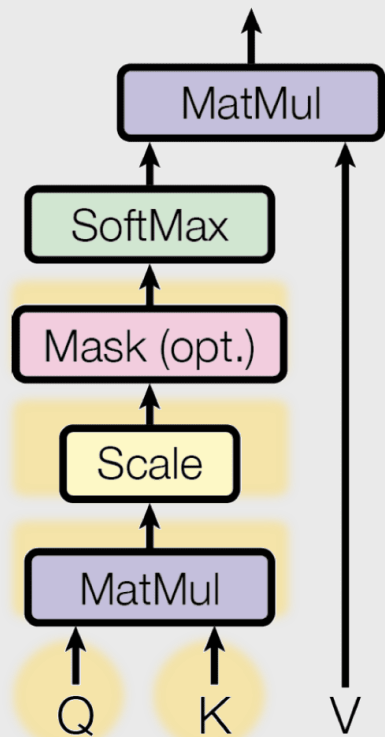
$$\mathbf{K} = \begin{bmatrix} k_{1,1} & \cdots & k_{n,1} \\ k_{1,2} & \ddots & k_{n,2} \\ \vdots & & \vdots \\ k_{1,d_k} & \cdots & k_{n,d_k} \end{bmatrix}_{5 \times 2}$$

\mathbf{K} is a matrix of size $n \times d_k$

$$\begin{bmatrix} \frac{e_{1,1}}{\sqrt{d_k}} & \frac{e_{1,2}}{\sqrt{d_k}} & \cdots & \frac{e_{1,n}}{\sqrt{d_k}} \\ \frac{e_{2,1}}{\sqrt{d_k}} & \frac{e_{2,2}}{\sqrt{d_k}} & \cdots & \frac{e_{2,n}}{\sqrt{d_k}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{e_{n,1}}{\sqrt{d_k}} & \frac{e_{n,2}}{\sqrt{d_k}} & \cdots & \frac{e_{n,n}}{\sqrt{d_k}} \end{bmatrix}_{5 \times 5}$$

Scaling of alignment scores

Self-Attention



\mathbf{Q} is a matrix of size $n \times d_k$

$$\mathbf{Q} = \begin{bmatrix} q_{1,1} & \cdots & q_{n,1} \\ q_{1,2} & \ddots & q_{n,2} \\ \vdots & & \vdots \\ q_{1,d_k} & \cdots & q_{n,d_k} \end{bmatrix}_{5 \times 2}$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

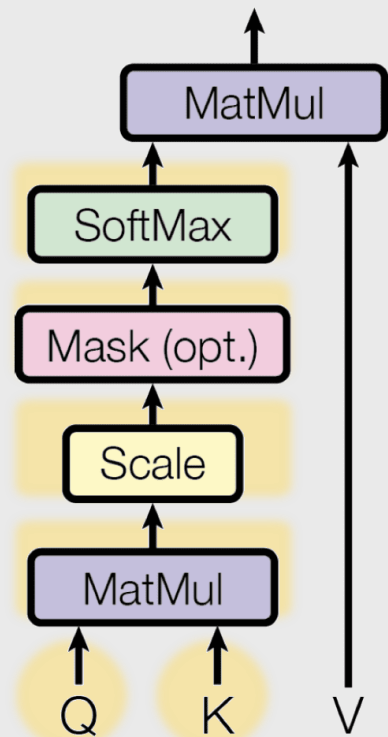
$$\mathbf{K} = \begin{bmatrix} k_{1,1} & \cdots & k_{n,1} \\ k_{1,2} & \ddots & k_{n,2} \\ \vdots & & \vdots \\ k_{1,d_k} & \cdots & k_{n,d_k} \end{bmatrix}_{5 \times 2}$$

\mathbf{K} is a matrix of size $n \times d_k$

$$\begin{bmatrix} \frac{e_{1,1}}{\sqrt{d_k}} & -\infty & -\infty \\ \frac{e_{2,1}}{\sqrt{d_k}} & \frac{e_{2,2}}{\sqrt{d_k}} & \cdots & -\infty \\ \vdots & \vdots & \ddots & \vdots \\ \frac{e_{n,1}}{\sqrt{d_k}} & \frac{e_{n,2}}{\sqrt{d_k}} & \cdots & \frac{e_{n,n}}{\sqrt{d_k}} \end{bmatrix}_{5 \times 5}$$

Apply Mask (optional)

Self-Attention



\mathbf{Q} is a matrix of size $n \times d_k$

$$\mathbf{Q} = \begin{bmatrix} q_{1,1} & \cdots & q_{n,1} \\ q_{1,2} & \ddots & q_{n,2} \\ \vdots & & \vdots \\ q_{1,d_k} & \cdots & q_{n,d_k} \end{bmatrix}_{5 \times 2}$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

$$\mathbf{K} = \begin{bmatrix} k_{1,1} & \cdots & k_{n,1} \\ k_{1,2} & \ddots & k_{n,2} \\ \vdots & & \vdots \\ k_{1,d_k} & \cdots & k_{n,d_k} \end{bmatrix}_{5 \times 2}$$

\mathbf{K} is a matrix of size $n \times d_k$

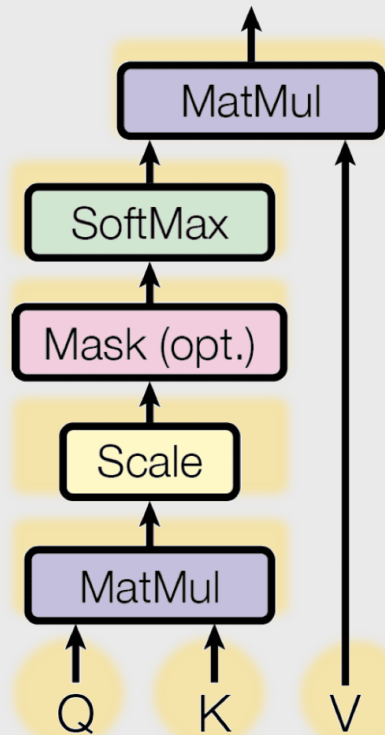
$$\begin{bmatrix} \text{softmax} \left(\frac{e_{1,1}}{\sqrt{d_k}} \right) & 0 & \cdots & 0 \\ \text{softmax} \left(\frac{e_{2,1}}{\sqrt{d_k}} \right) & \frac{e_{2,2}}{\sqrt{d_k}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \text{softmax} \left(\frac{e_{n,1}}{\sqrt{d_k}} \right) & \frac{e_{n,2}}{\sqrt{d_k}} & \cdots & \frac{e_{n,n}}{\sqrt{d_k}} \end{bmatrix}_{5 \times 5}$$

SoftMax of Scaled dot product

$\frac{1}{\sqrt{d_k}}$ is the scaling factor

Self-Attention

Vaswani et al. Attention Is All You Need (NIPS 2017)



\mathbf{Q} is a matrix of size $n \times d_k$

$$\mathbf{Q} = \begin{bmatrix} q_{1,1} & \cdots & q_{n,1} \\ q_{1,2} & \ddots & q_{n,2} \\ \vdots & & \vdots \\ q_{1,d_k} & \cdots & q_{n,d_k} \end{bmatrix} 5 \times 2$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

$$\mathbf{K} = \begin{bmatrix} k_{1,1} & \cdots & k_{n,1} \\ k_{1,2} & \ddots & k_{n,2} \\ \vdots & & \vdots \\ k_{1,d_k} & \cdots & k_{n,d_k} \end{bmatrix} 5 \times 2$$

\mathbf{K} is a matrix of size $n \times d_k$

$$\mathbf{V} = \begin{bmatrix} v_{1,1} & \cdots & v_{1,d_v} \\ v_{2,1} & \ddots & v_{2,d_v} \\ \vdots & & \vdots \\ v_{n,1} & \cdots & v_{n,d_v} \end{bmatrix} 5 \times 2$$

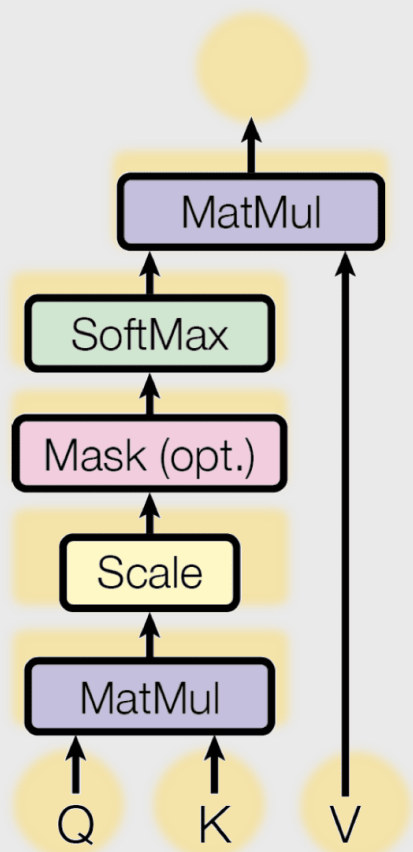
\mathbf{V} is a matrix of size $n \times d_v$

$$\begin{bmatrix} \text{softmax} \left(\frac{e_{1,1}}{\sqrt{d_k}} \right) & 0 & \cdots & 0 \\ \text{softmax} \left(\frac{e_{2,1}}{\sqrt{d_k}} \right) & \frac{e_{2,2}}{\sqrt{d_k}} & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ \text{softmax} \left(\frac{e_{d_k,1}}{\sqrt{d_k}} \right) & \frac{e_{d_k,2}}{\sqrt{d_k}} & \cdots & \frac{e_{d_k,d_k}}{\sqrt{d_k}} \end{bmatrix} 5 \times 5 \cdot \begin{bmatrix} v_{1,1} & \cdots & v_{1,d_v} \\ v_{2,1} & \ddots & v_{2,d_v} \\ \vdots & & \vdots \\ v_{n,1} & \cdots & v_{n,d_v} \end{bmatrix} 5 \times 2$$

resulting self-attention weights

Self-Attention

Vaswani et al. Attention Is All You Need (NIPS 2017)



Q is a matrix of size $n \times d_k$

$$Q = \begin{bmatrix} q_{1,1} & \cdots & q_{n,1} \\ q_{1,2} & \ddots & q_{n,2} \\ \vdots & & \vdots \\ q_{1,d_k} & \cdots & q_{n,d_k} \end{bmatrix}$$

$$\text{Head}_1 = \begin{bmatrix} a_{1,1} & \cdots & a_{1,d_v} \\ a_{2,1} & \ddots & a_{2,d_v} \\ \vdots & & \vdots \\ a_{n,1} & \cdots & a_{n,d_v} \end{bmatrix} 5 \times 2$$

resulting self-attention weights

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

$$\mathbf{K} = \begin{bmatrix} k_{1,1} & \cdots & k_{n,1} \\ k_{1,2} & \ddots & k_{n,2} \\ \vdots & & \vdots \\ k_{1,d_k} & \cdots & k_{n,d_k} \end{bmatrix}$$

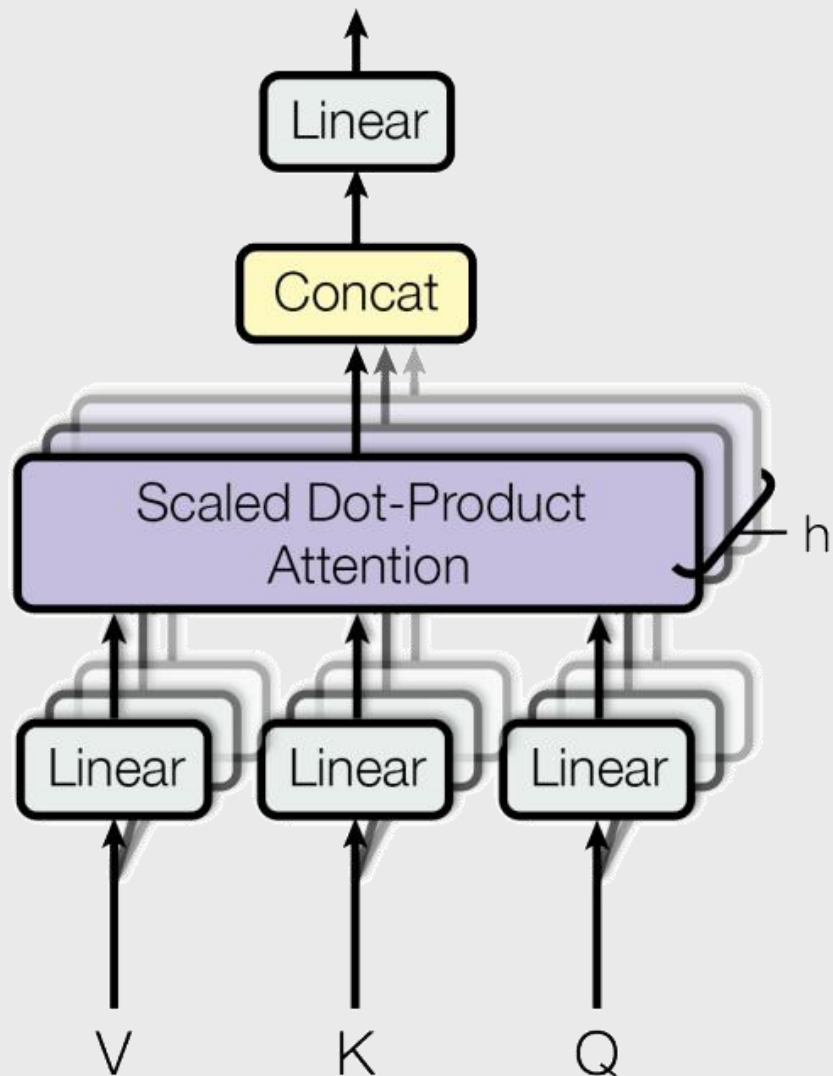
$$\mathbf{V} = \begin{bmatrix} v_{1,1} & \cdots & v_{1,d_v} \\ v_{2,1} & \ddots & v_{2,d_v} \\ \vdots & & \vdots \\ v_{d,1} & \cdots & v_{d,d_v} \end{bmatrix}$$

\mathbf{K} is a matrix of size $n \times d_k$

\mathbf{V} is a matrix of size $d \times d_v$

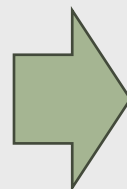
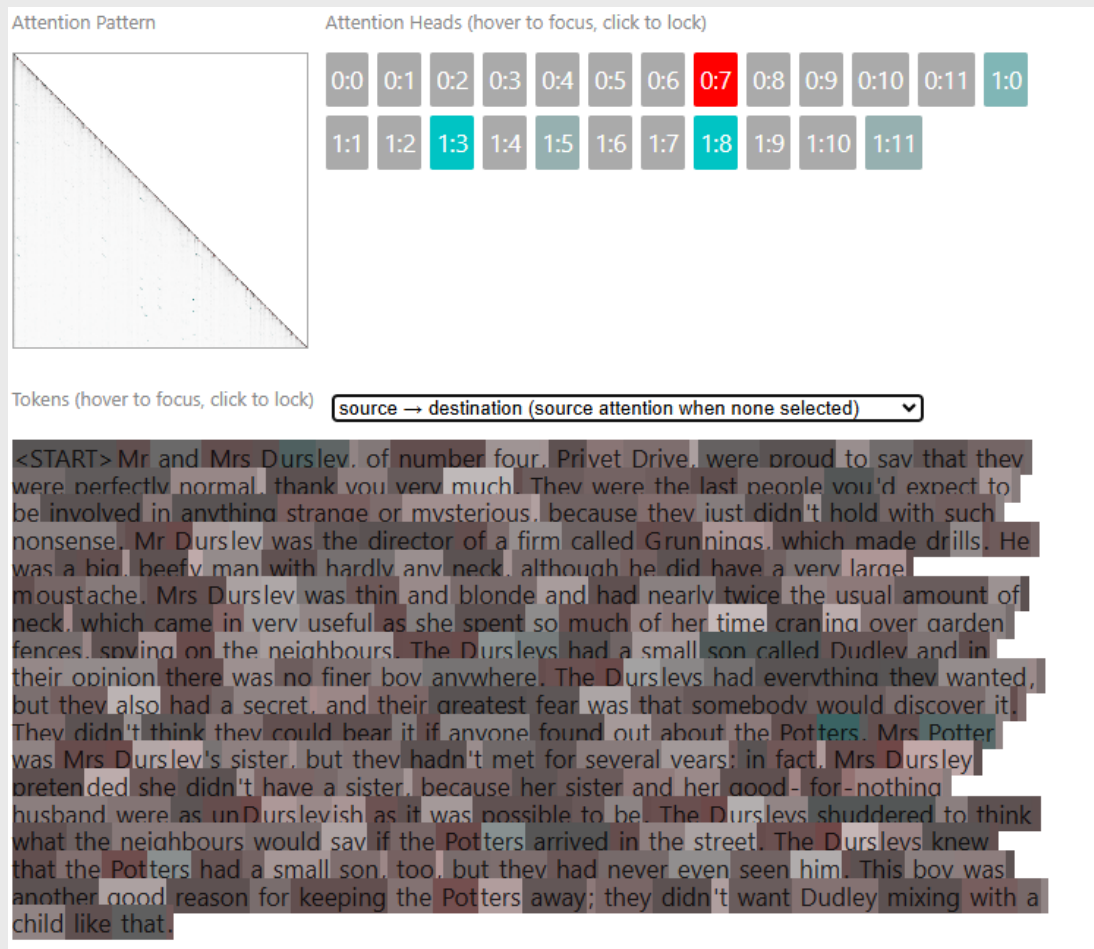
Multi Head Self-Attention

Vaswani et al. Attention Is All You Need (NIPS 2017)



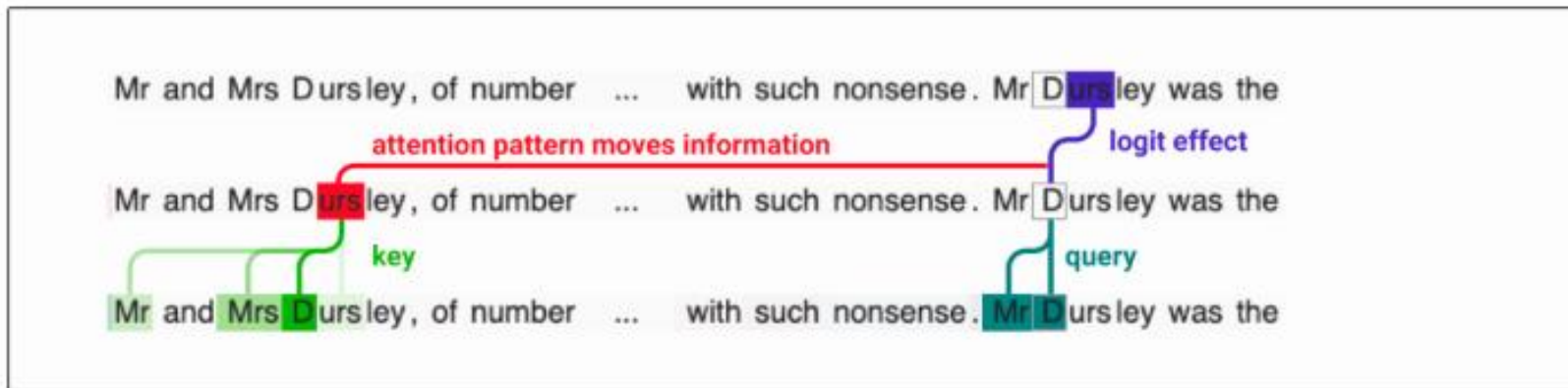
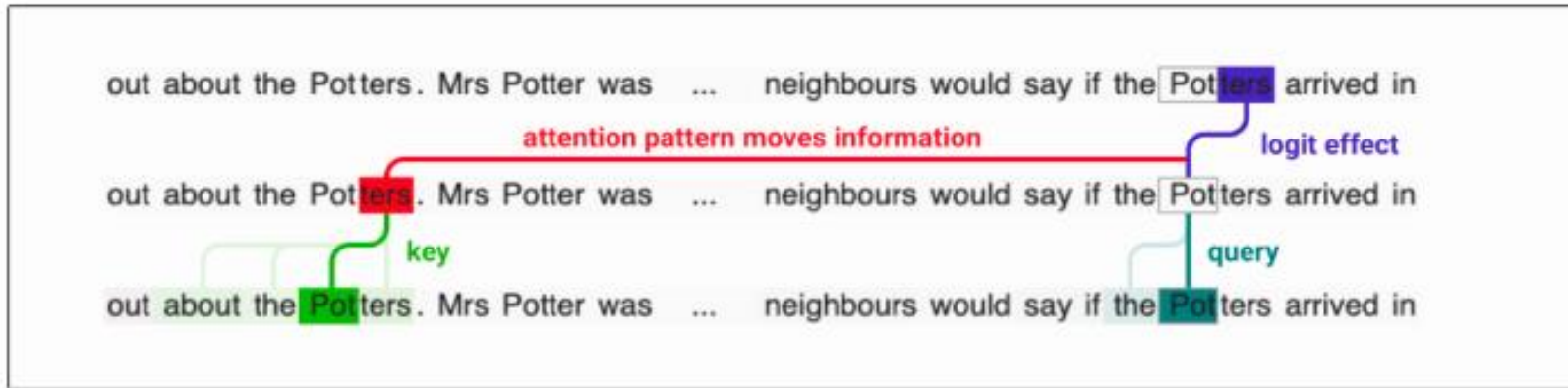
$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \\ = \text{Concat}(\text{Head}_1, \text{Head}_2, \dots, \text{Head}_h) \mathbf{W}_{\text{Out}}$$

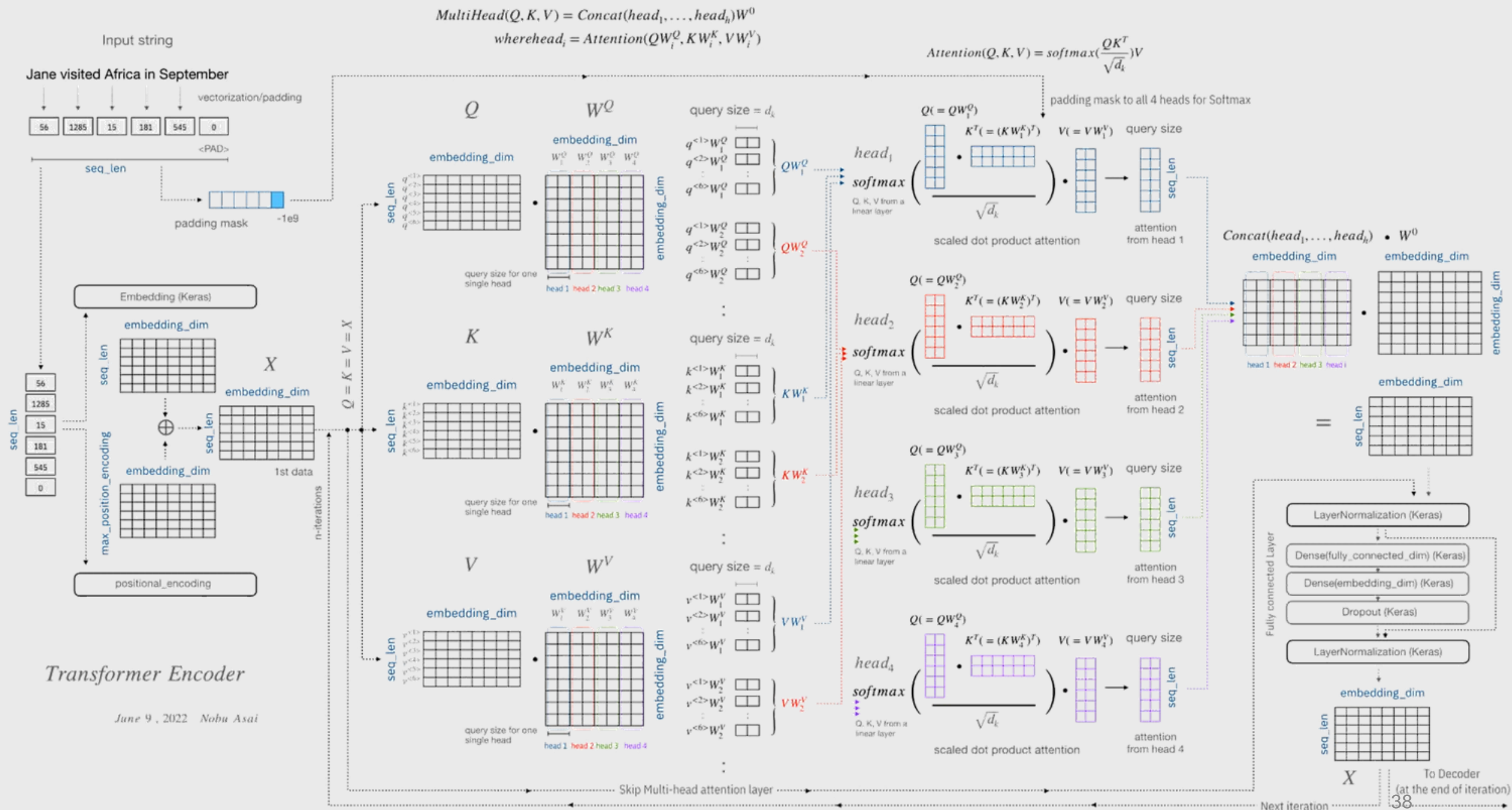
Attention Map



How Query and Key might work

The query searches for "similar" key vectors, but because keys are shifted, it finds the next token.

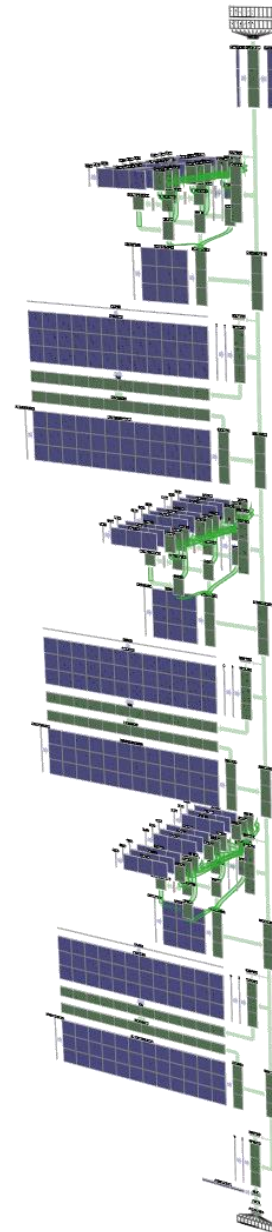
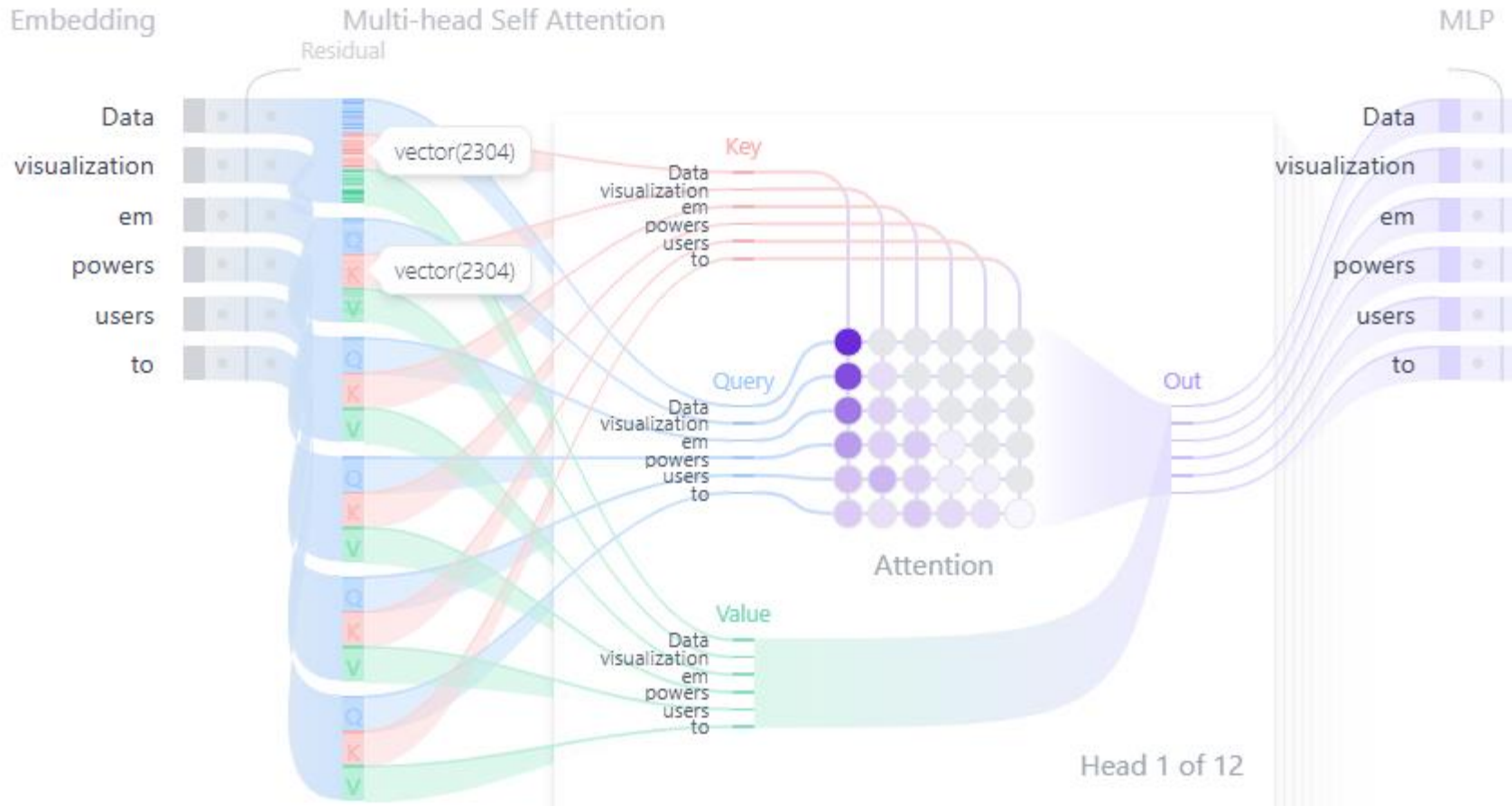




Transformer Visualisations & Explainers (Online Resources)

nano-gpt

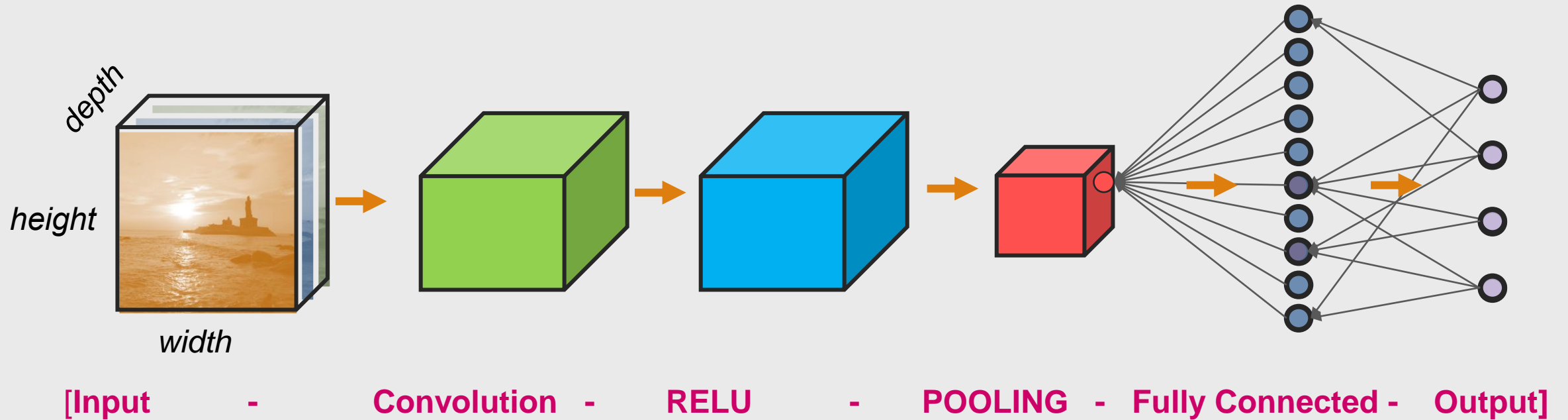
n_params = 85,584



- <https://bbycroft.net/llm>
- <https://poloclub.github.io/transformer-explainer/>
- <https://jalammar.github.io/illustrated-transformer/>

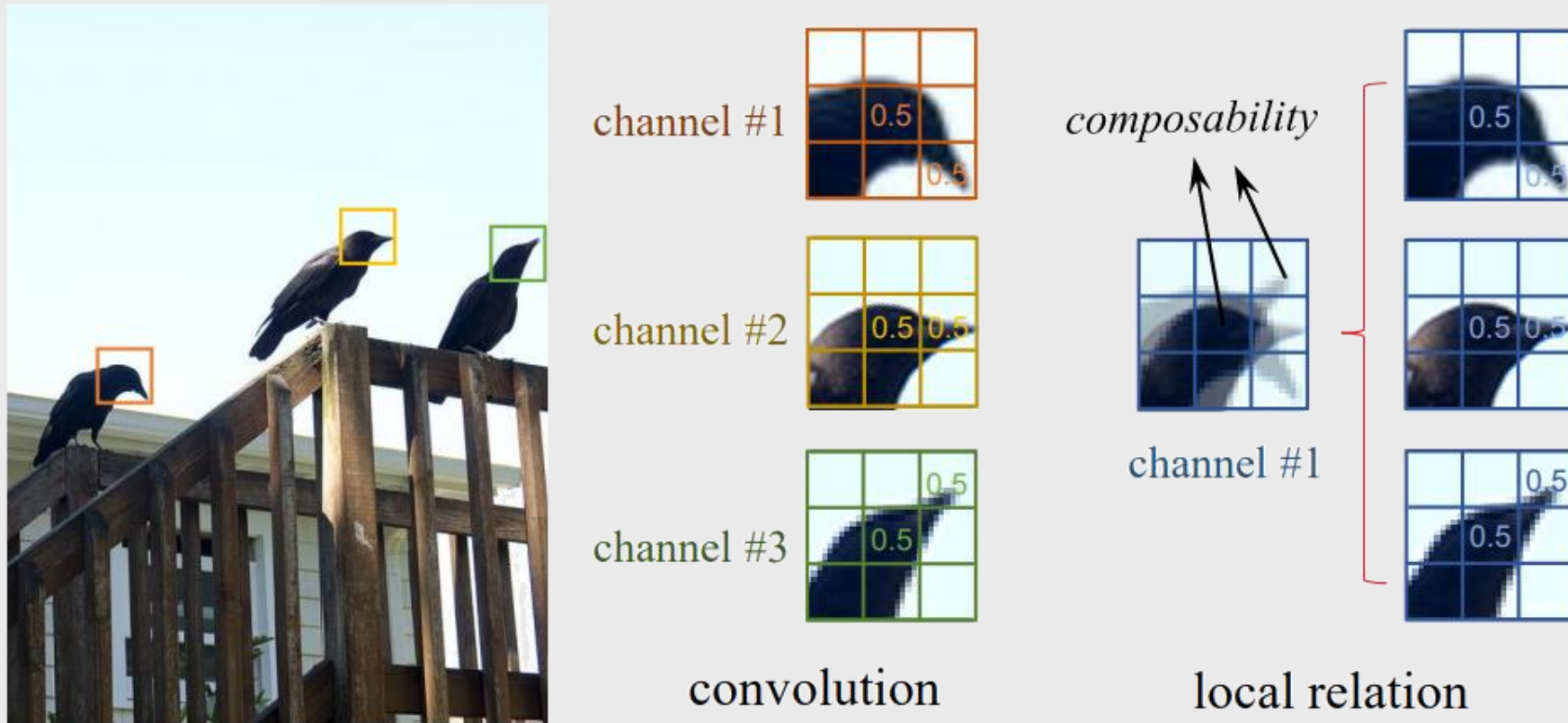
40

Convolutional Neural Nets



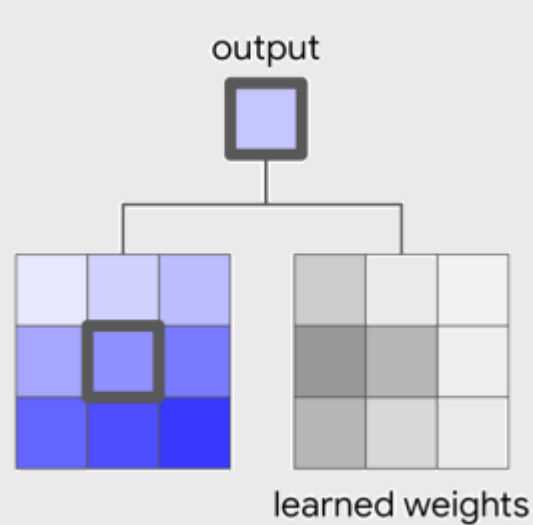
3×3 convolution layer and the 3×3 local relation layer

Hu et al. (2019). Local relation networks for image recognition. ICCV



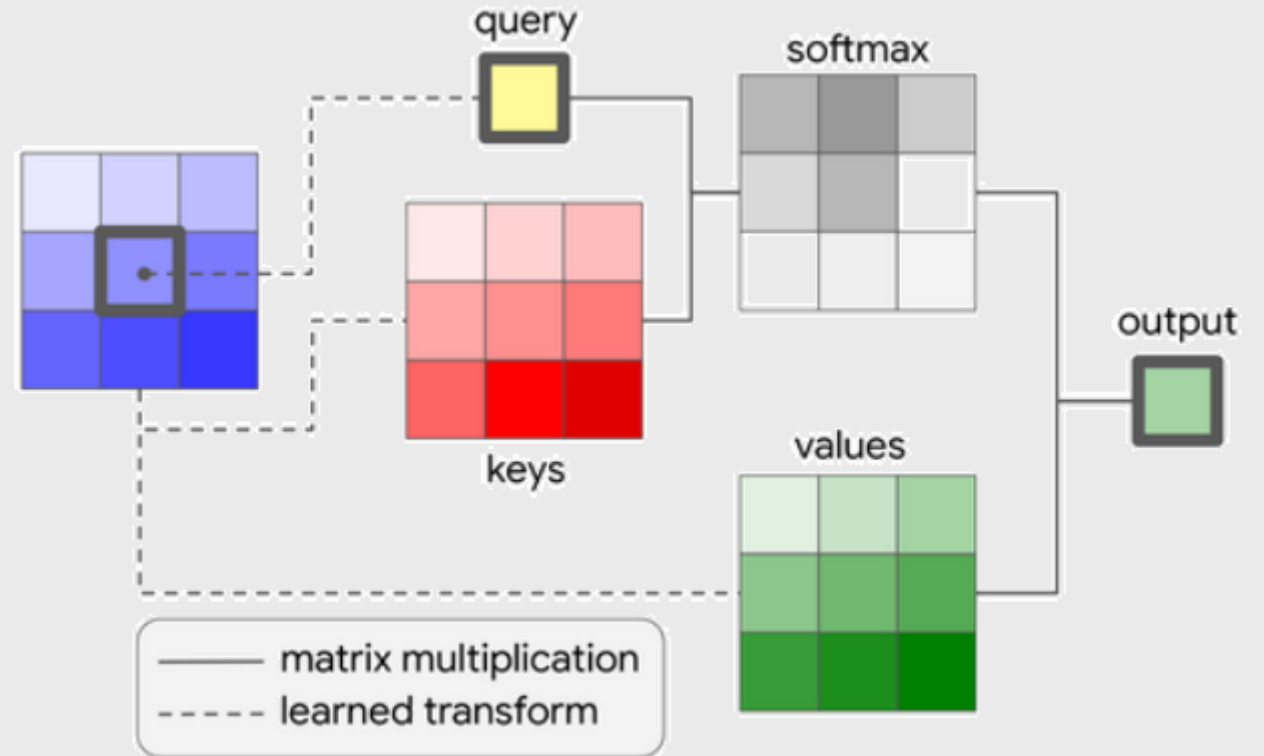
Convolutional Nets Vs Transformer

Ramachandran et al. Stand-alone self-attention in vision models. NIPS 2019



3×3 convolution.

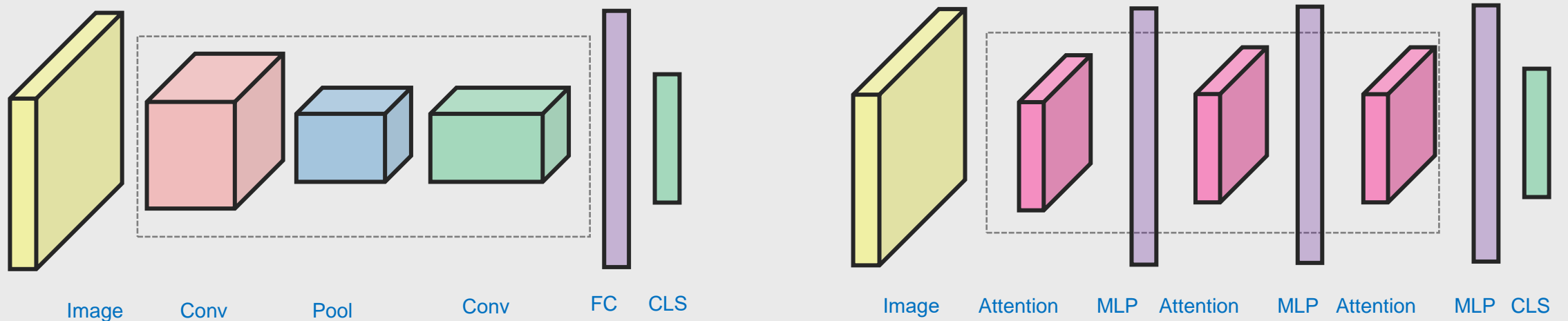
The output is the inner product between the local window and the learned weights



Self-attention around image local region
The output is local self attention

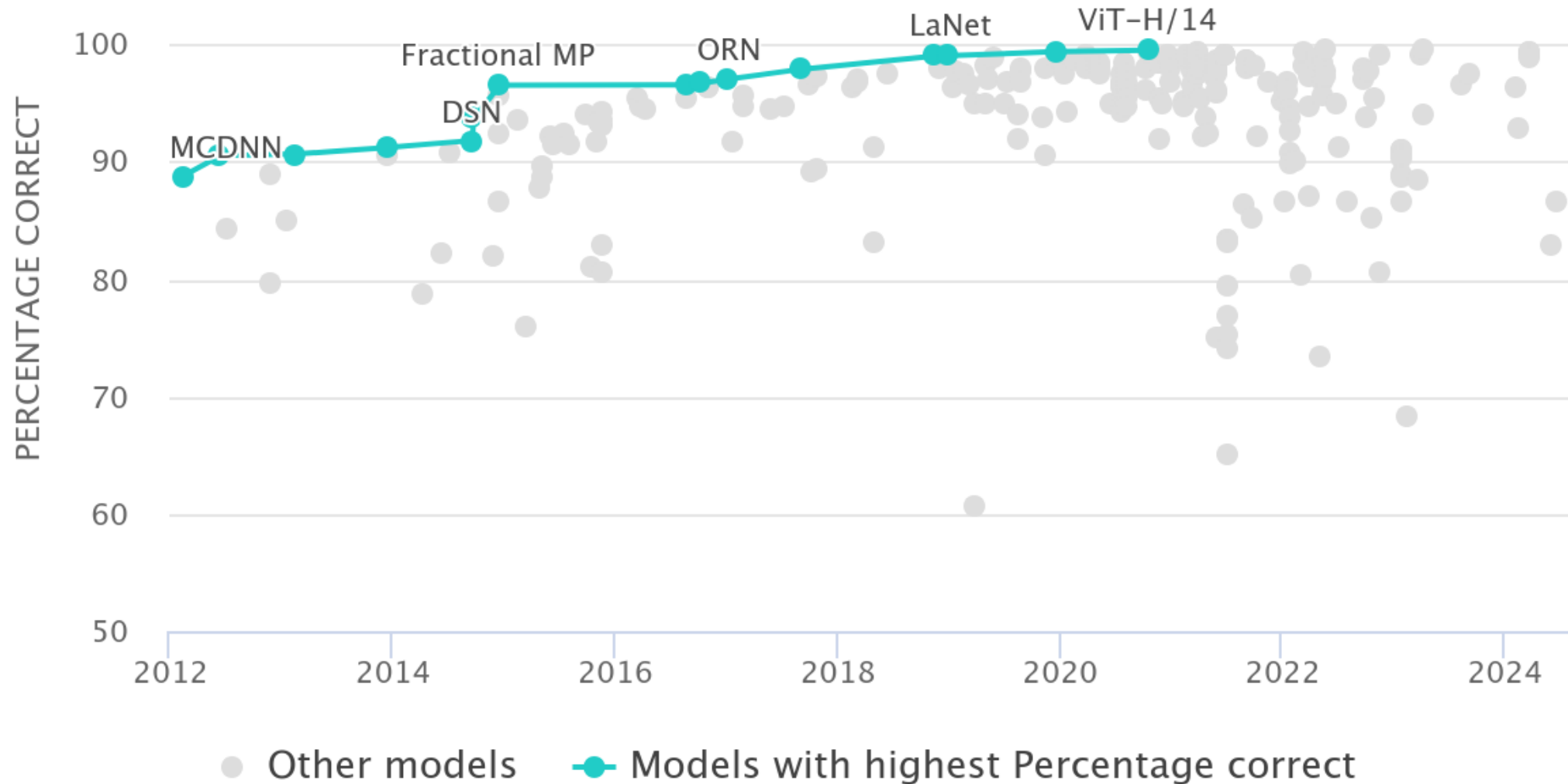
Standalone Self-attention in Vision Models

Ramachandran et al. Stand-alone self-attention in vision models. NIPS 2019

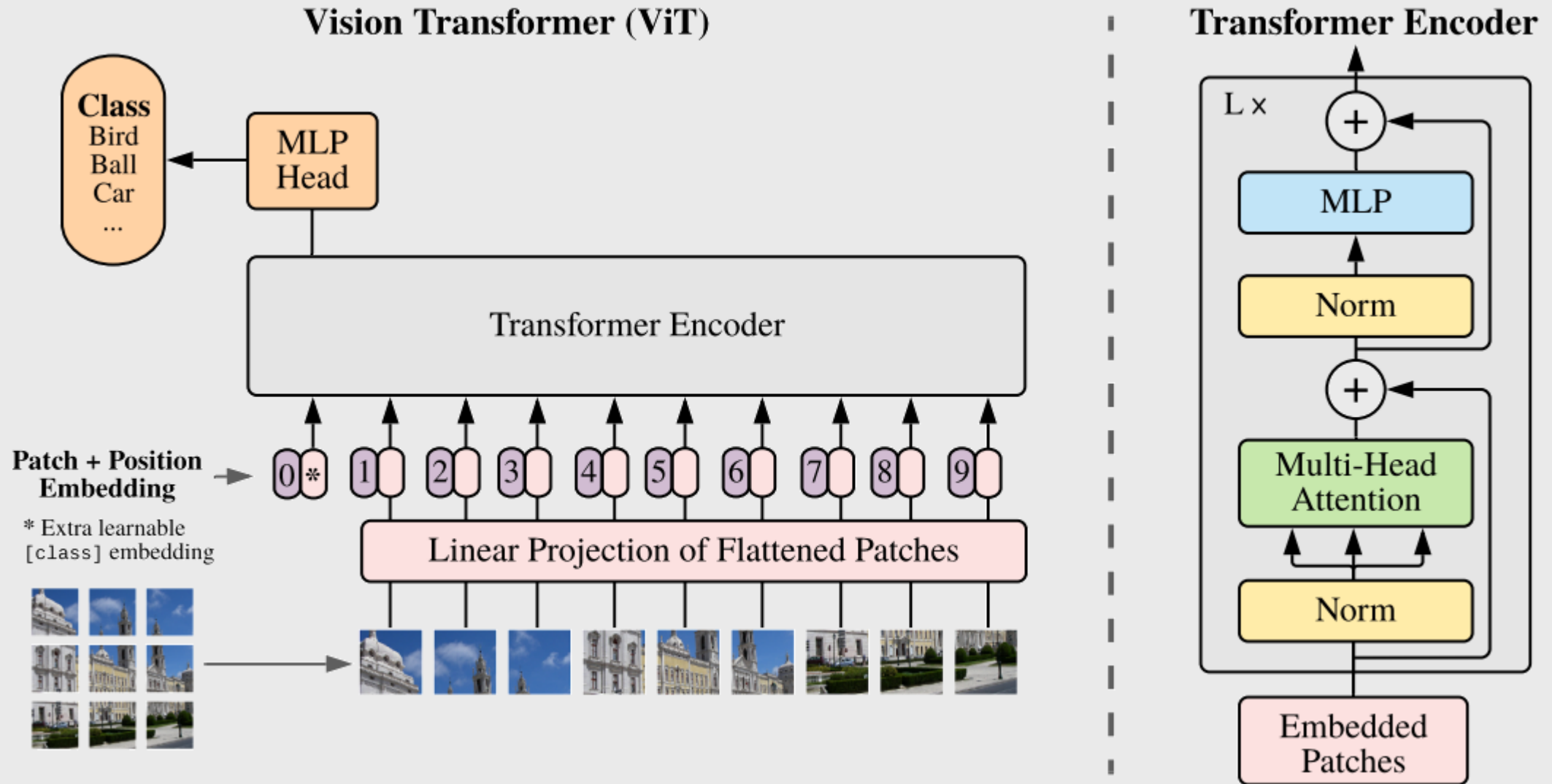


	ResNet-26			ResNet-38			ResNet-50		
	FLOPS (B)	Params (M)	Acc. (%)	FLOPS (B)	Params (M)	Acc. (%)	FLOPS (B)	Params (M)	Acc. (%)
Baseline	4.7	13.7	74.5	6.5	19.6	76.2	8.2	25.6	76.9
Conv-stem + Attention	4.5	10.3	75.8	5.7	14.1	77.1	7.0	18.0	77.4
Full Attention	4.7	10.3	74.8	6.0	14.1	76.9	7.2	18.0	77.6

Convolutional Nets Vs Transformers

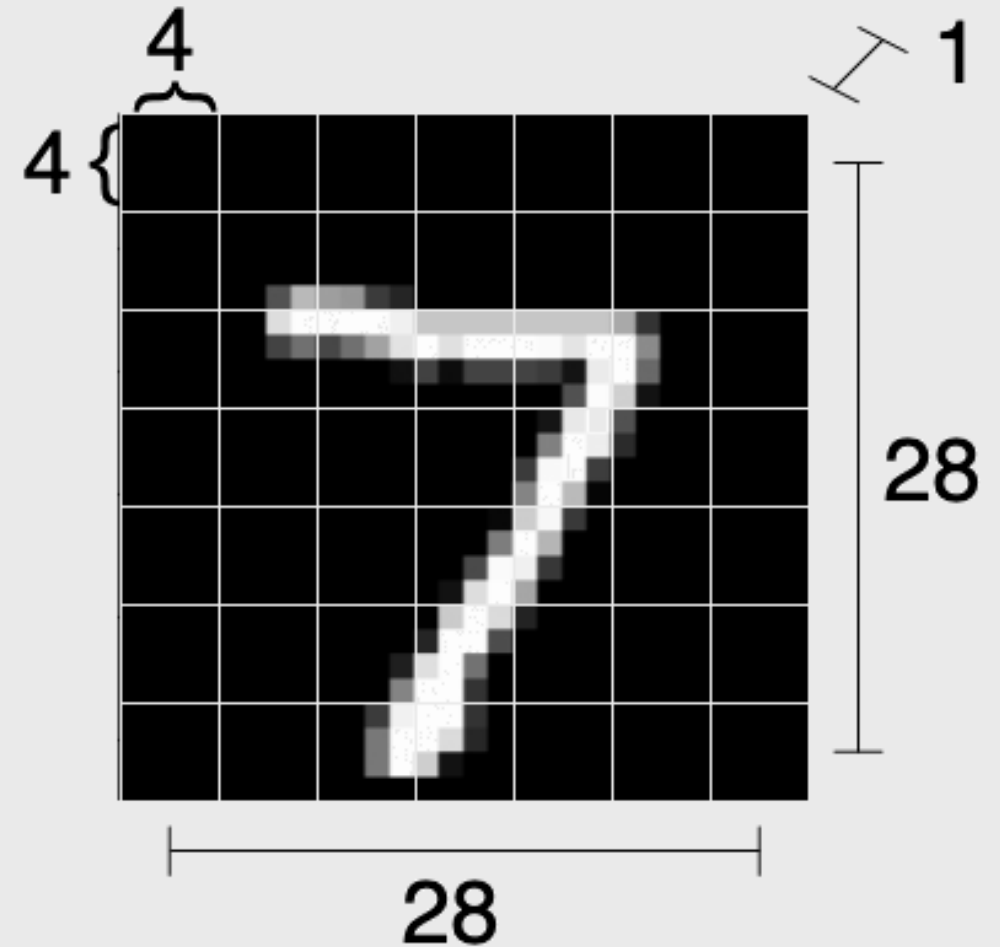
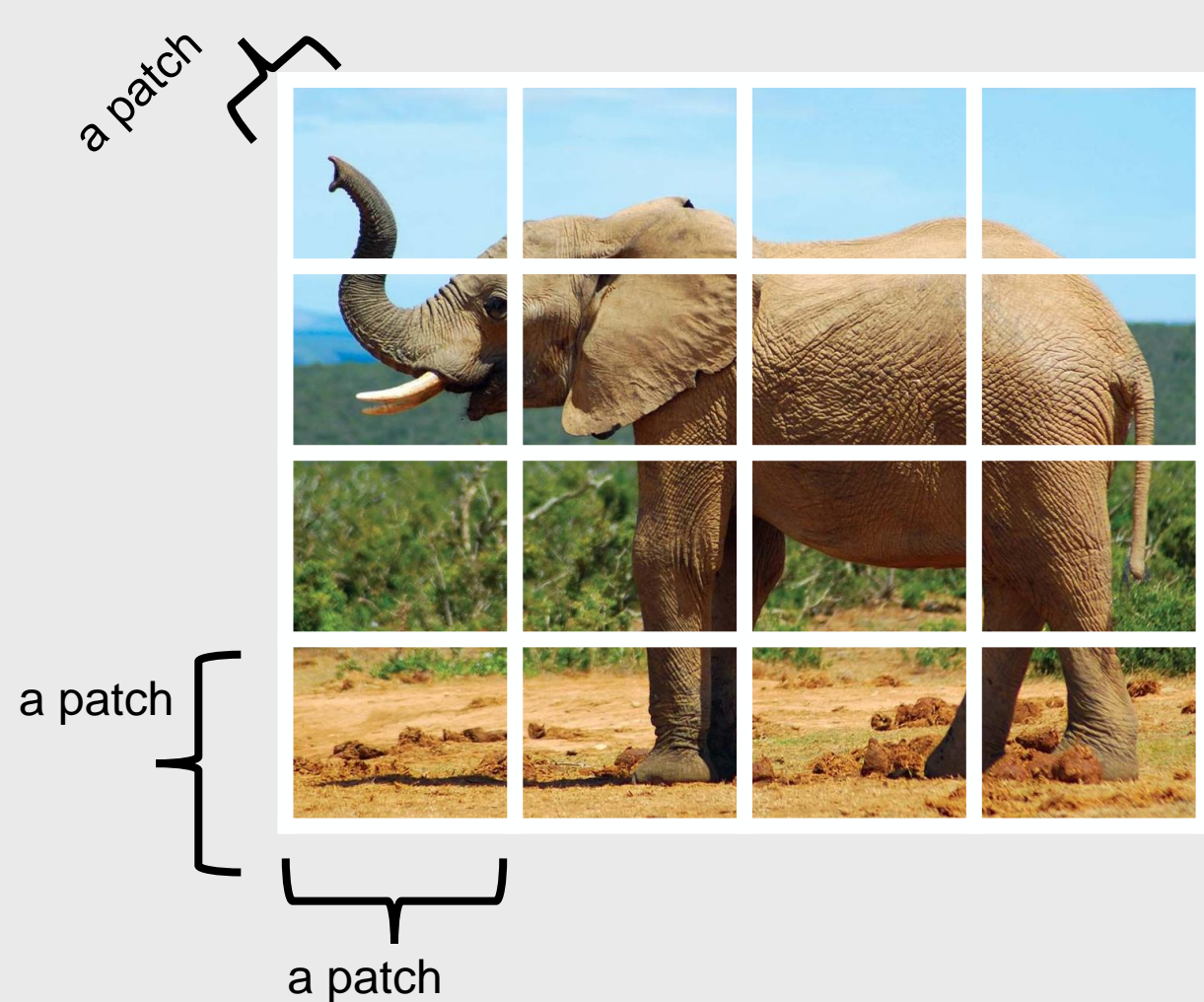


How Vision Transformer Models Works



Splitting an Image into Patches

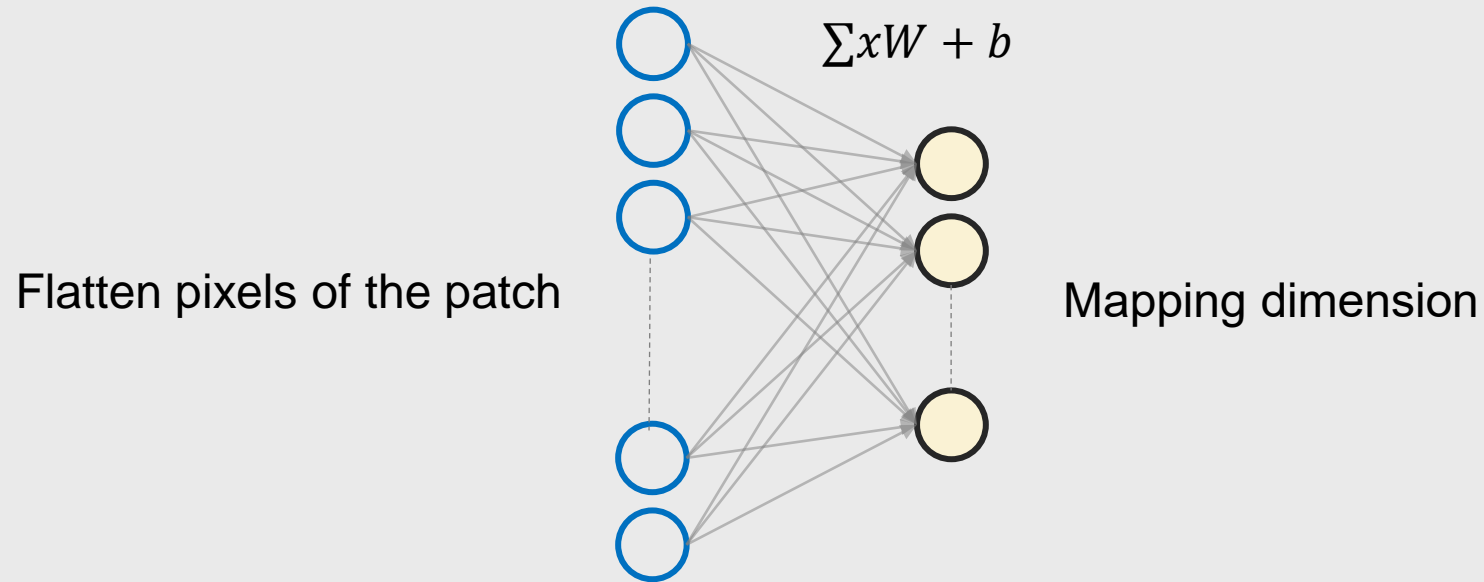
Split the image into patches, each of size (H'xW'xD)



MIST dataset

Linear mapping

Linear projection to D-dimensional vector



Linear Mapping
h-dim



Flatten
array (n-
dim)
Patches

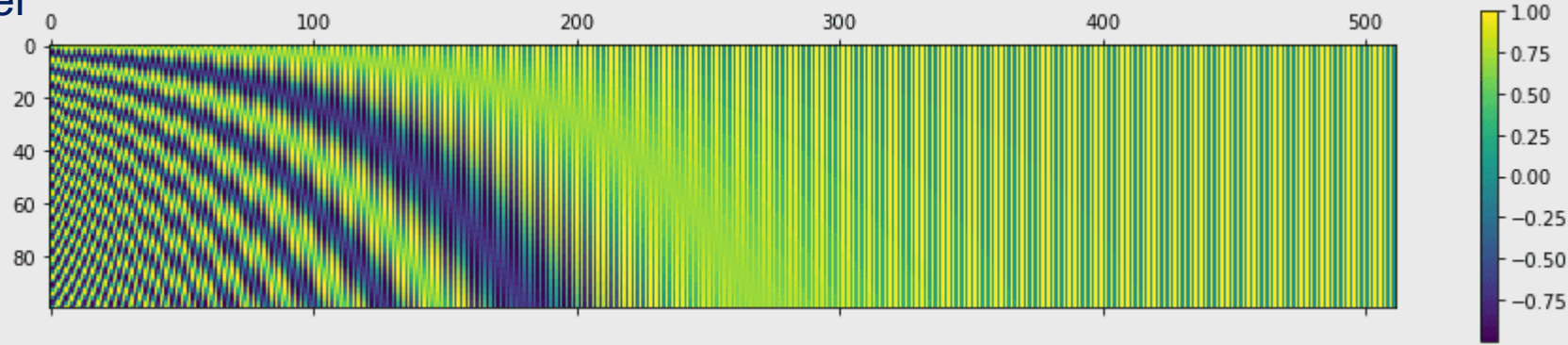


— —



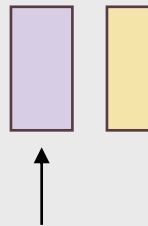
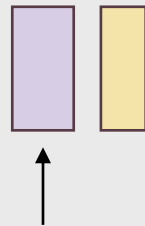
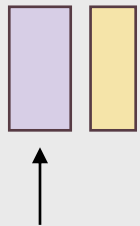
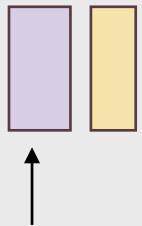
Positional Encoding

Inform the model where the patch's position in the image is. In other word use sine and cosine values for respective patch number

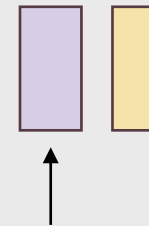
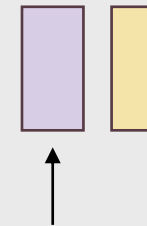
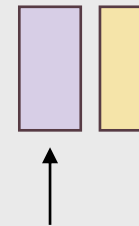


$$p_{i,j} = \begin{cases} \sin\left(\frac{i}{10000^{\frac{j}{d_{emb-dim}}}}\right) & \text{if } j \text{ is even} \\ \cos\left(\frac{i}{10000^{\frac{j-1}{d_{emb-dim}}}}\right) & \text{if } j \text{ is odd} \end{cases}$$

Linear
Mapping
h-dim



— —



Flatten
array (n-
dim)
Patches



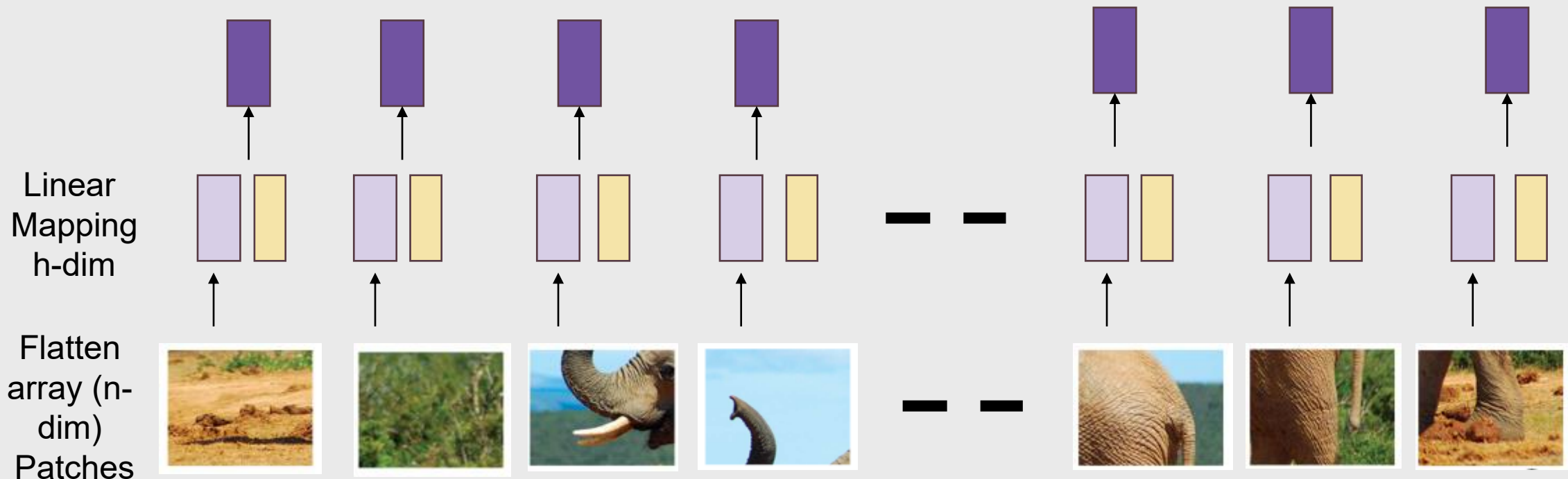
— —



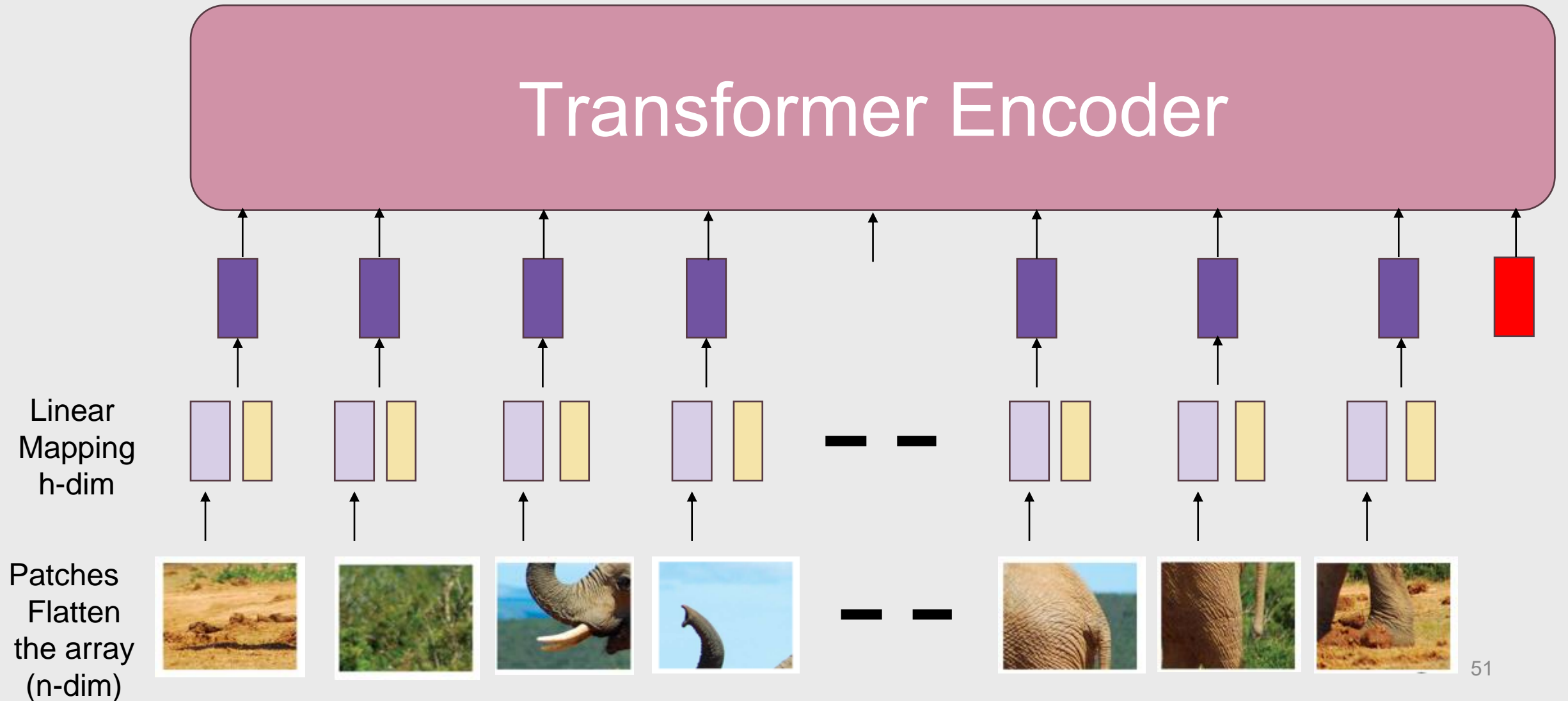
Positional Encoding and Vectors

Inform the model where the patch's position in the image is. In other word use sine and cosine values for respective patch number

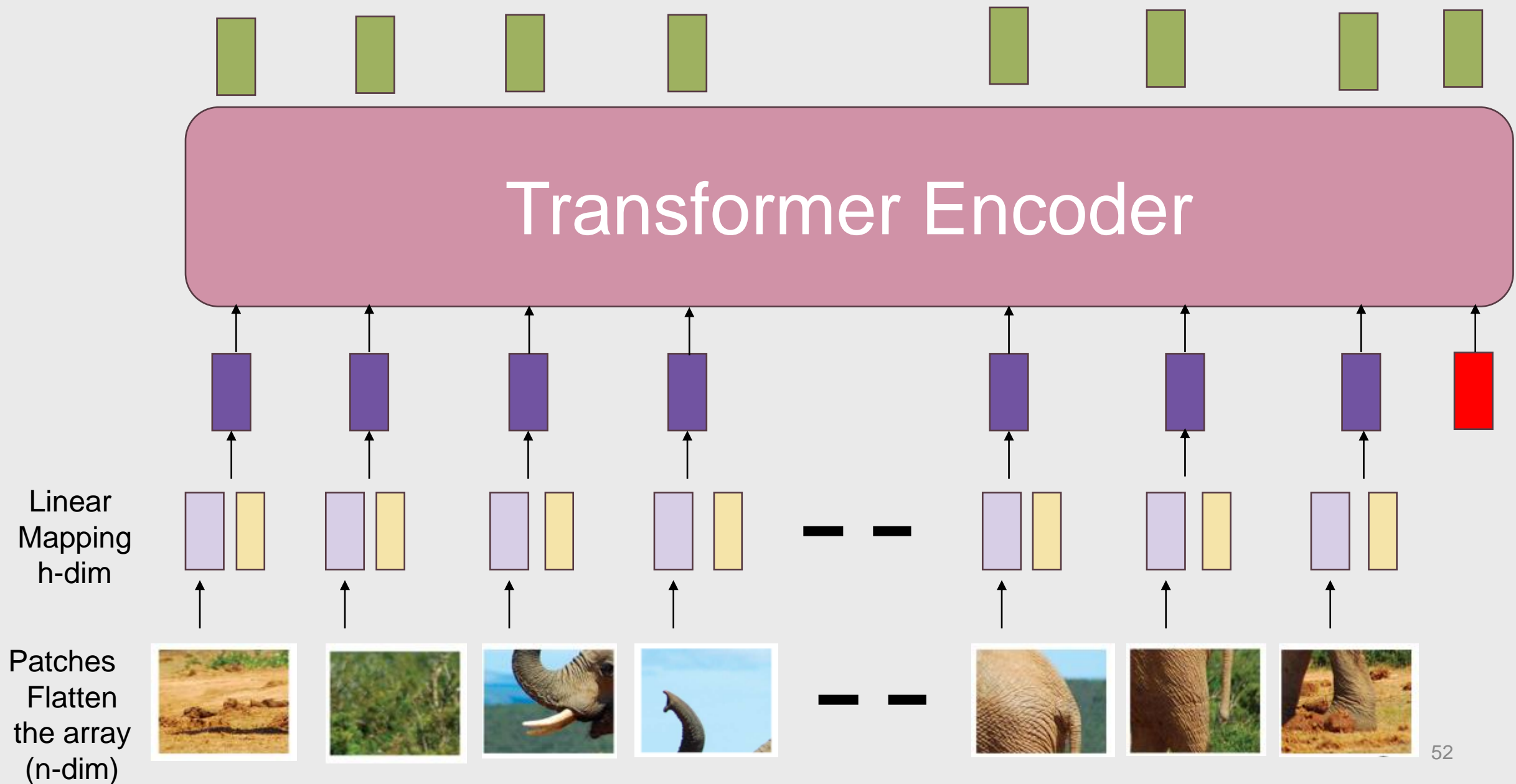
$$p_{i,j} = \begin{cases} \sin\left(\frac{i}{10000^{\frac{j}{d_{emb-dim}}}}\right) & \text{if } j \text{ is even} \\ \cos\left(\frac{i}{10000^{\frac{j-1}{d_{emb-dim}}}}\right) & \text{if } j \text{ is odd} \end{cases}$$



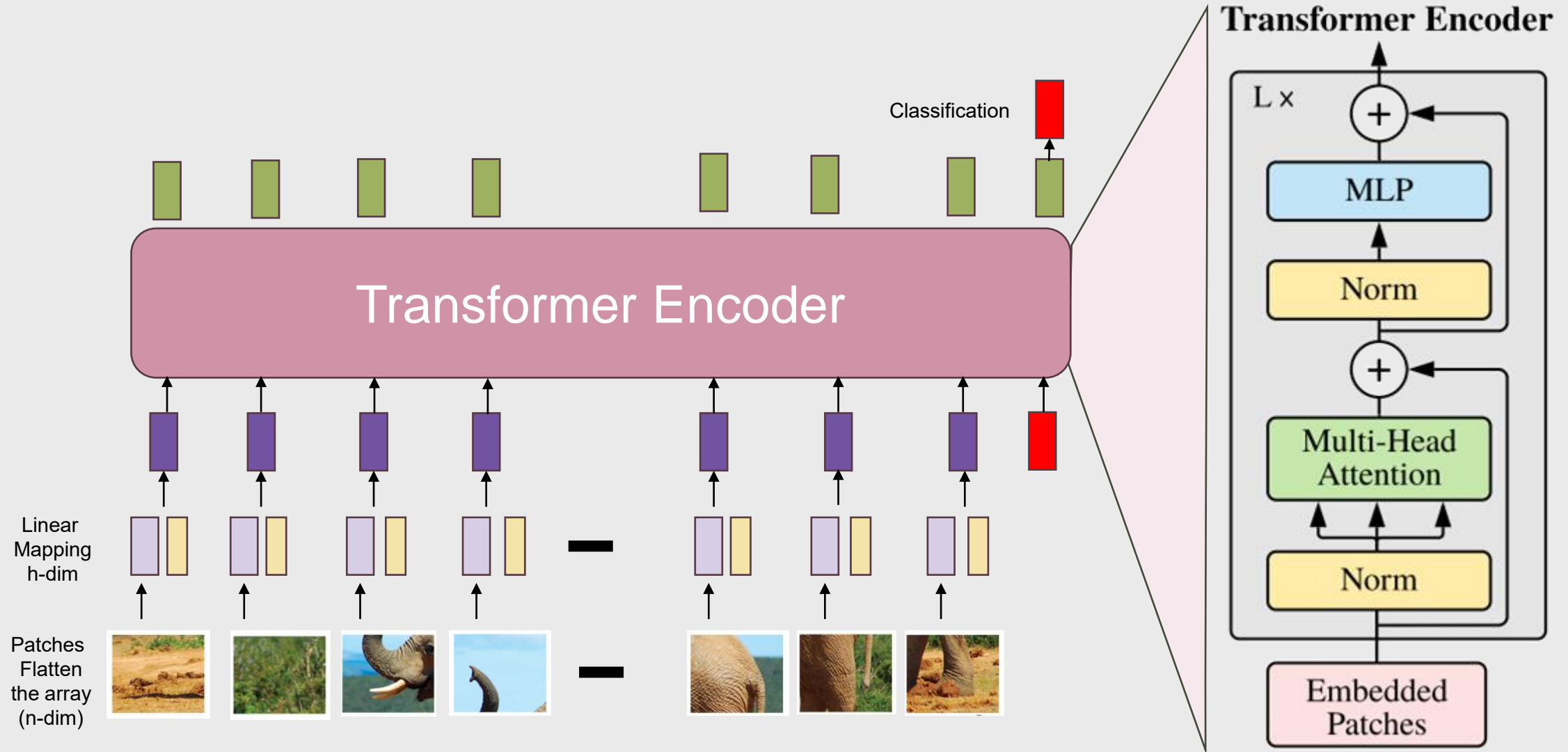
Add a Learnable Classification Token



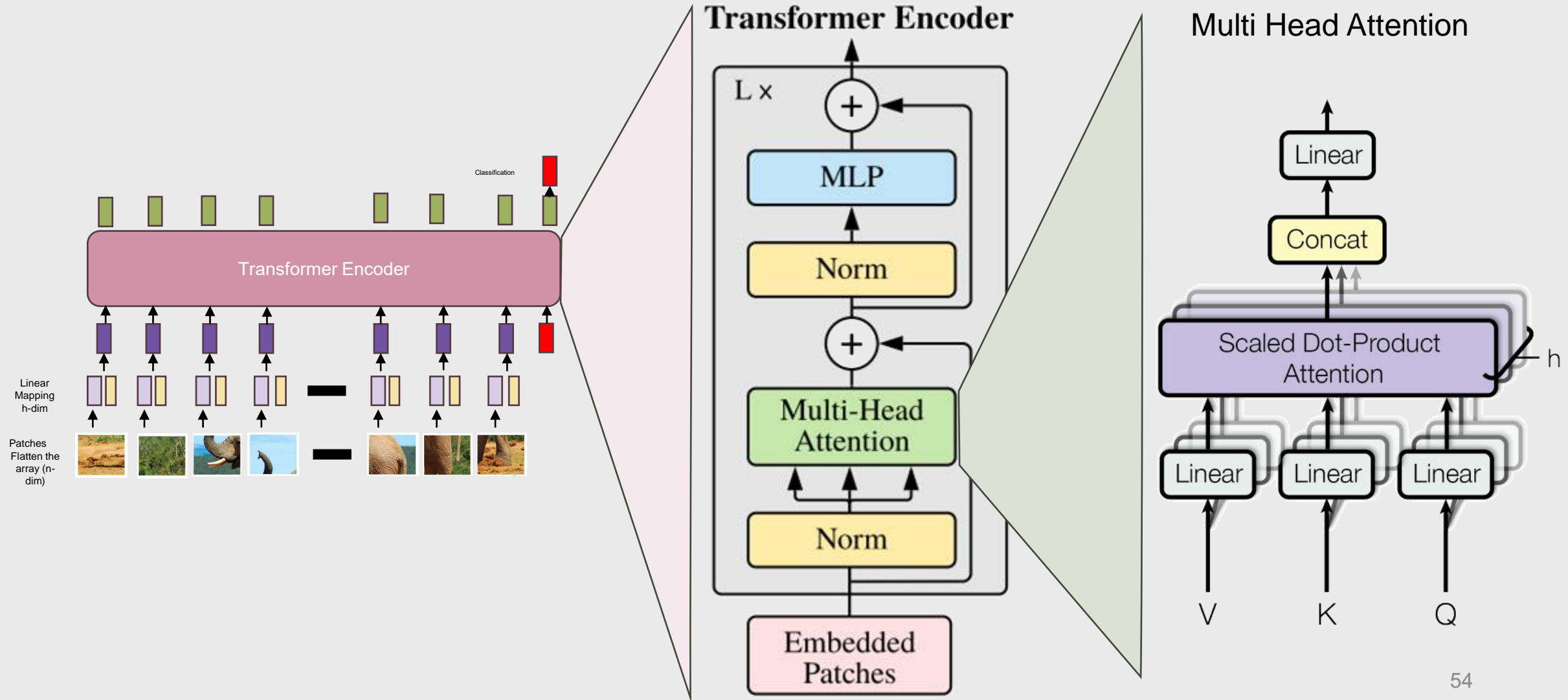
Output vector of Transformer Encoder



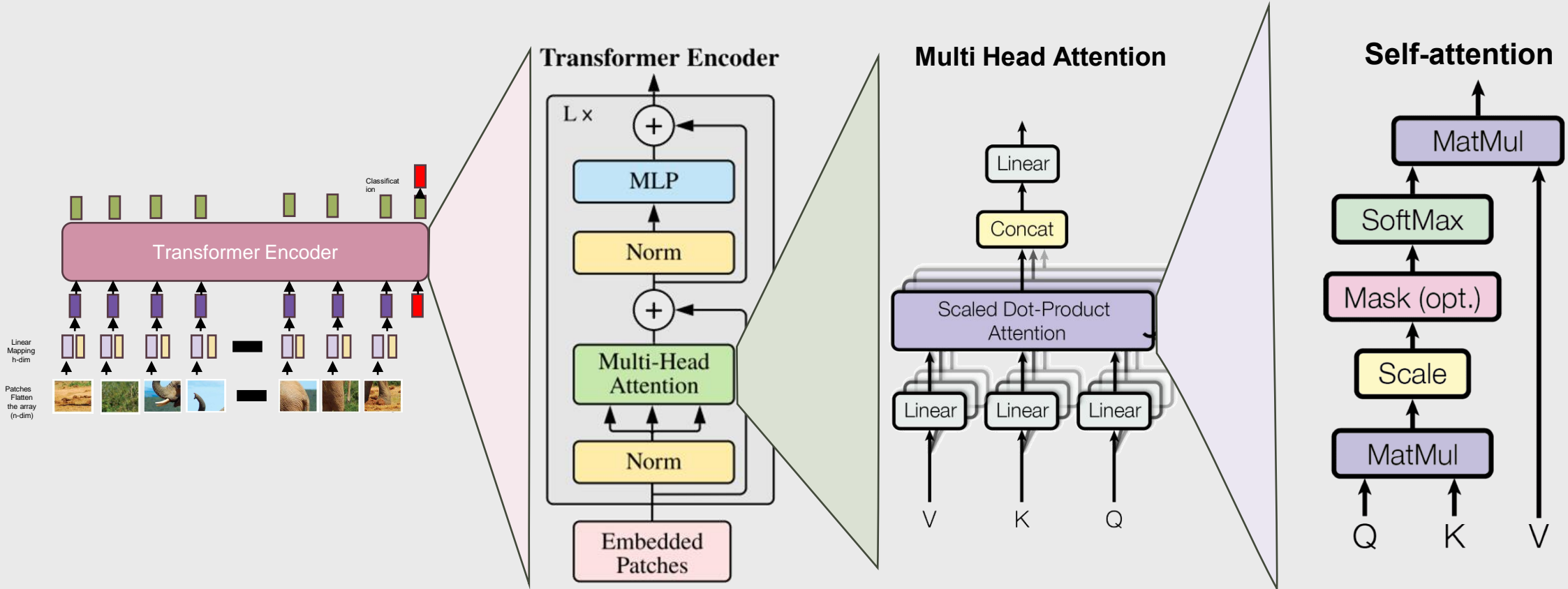
Same as ChatGPT Transformer



Same as ChatGPT Transformer



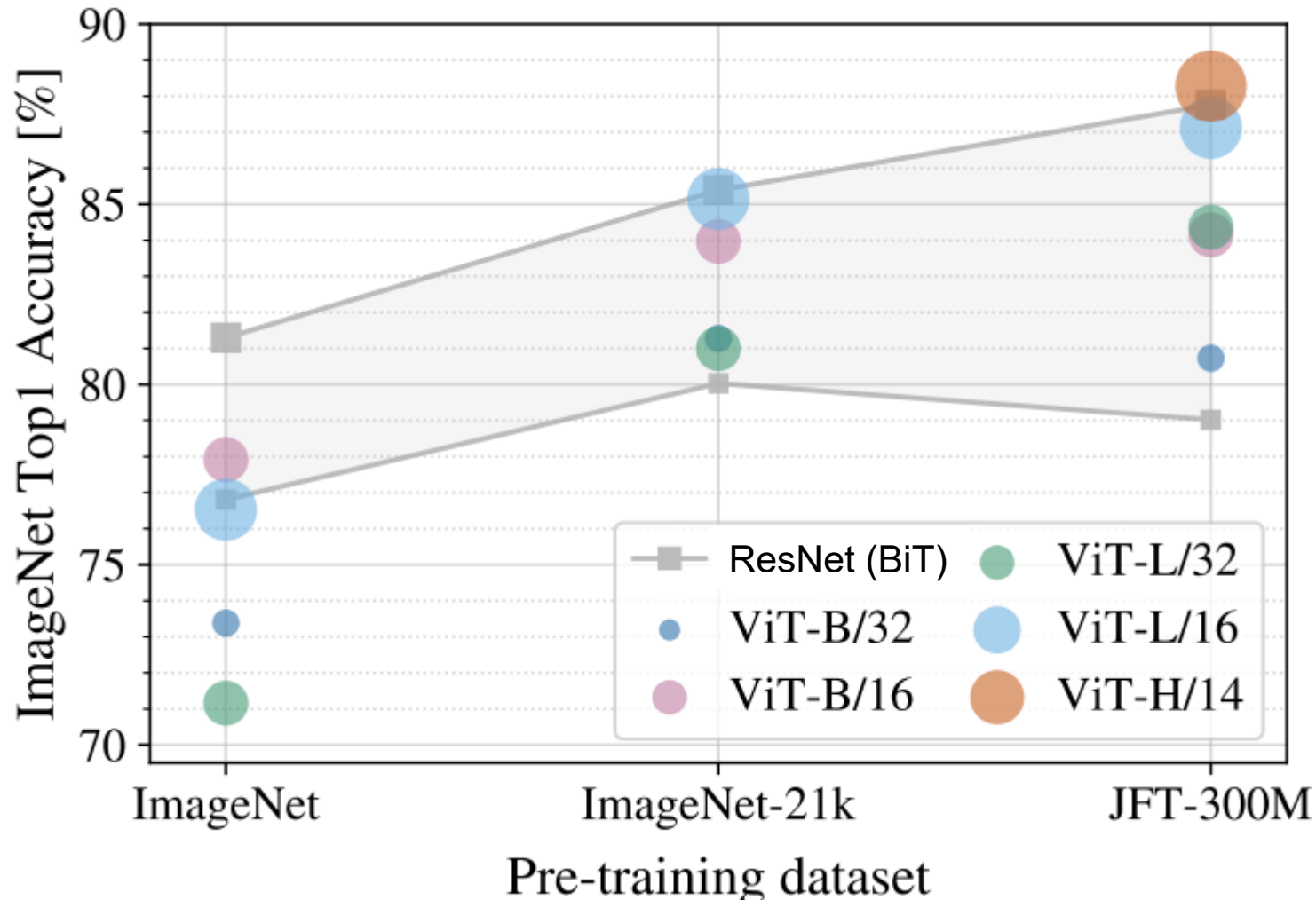
Same as ChatGPT Transformer



$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

ViT Performance

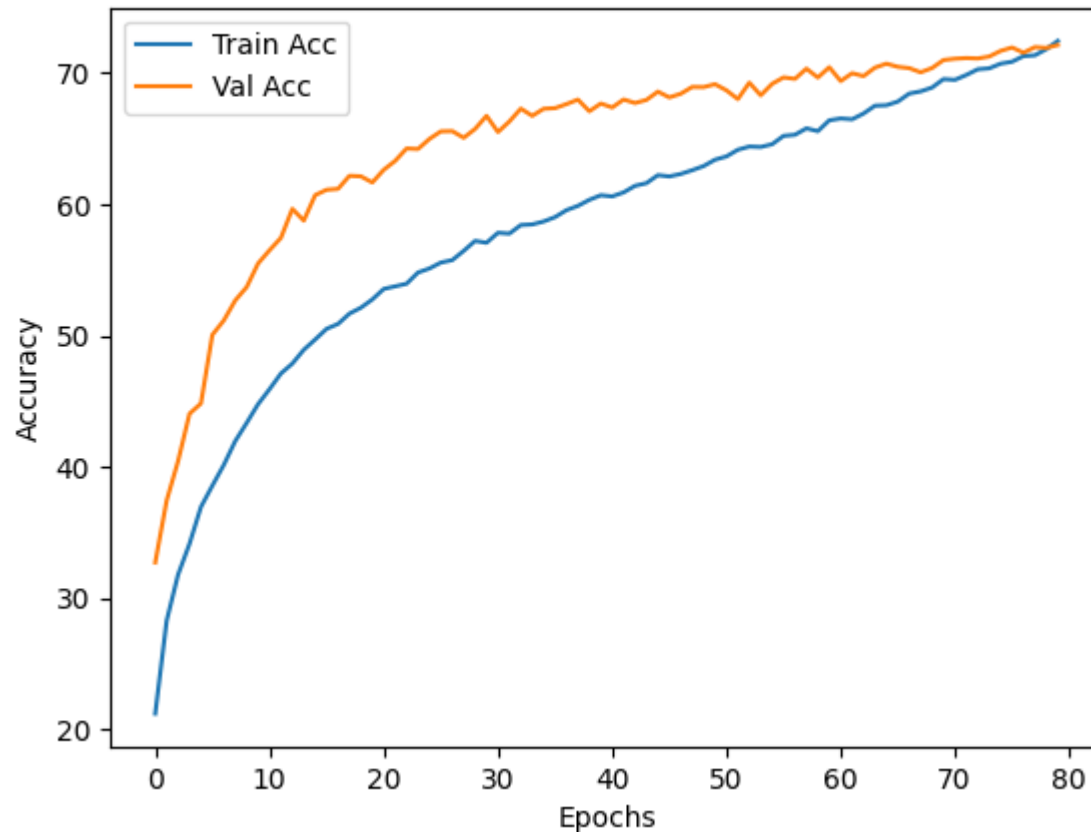
Dosovitskiy et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, ICLR 2017



Note that this performance is only achieved when ViT is pre-trained on large dataset (in this case JFT-300M is a 300 million image dataset of Google)

ViT on CIFAR-10 (without Pre-Training)

Dosovitskiy et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, ICLR 2017



Note that the performance depends on hyper parameter tuning models' size etc.

ViT Performance

- Worse than ResNet when trained just on ImageNet
- Performance improved when pre-trained on very large dataset
- Pretrained outperforms much bigger CNNs
- You need large GPUs (Computational Cost is very high)

Coursework Brief (Part III)

More details to be released this week (before Practical Session)

Implement Convolutional Neural Networks (specifically using VGG16) on CIFAR-10 dataset and solve following three problems:

- For the training **use early stopping** and **save the model** that produce best validation results. (you will need to use some of training data as **validation set**) [Marks 10: 5+3+2]
- What would be the performance of VGG16 **with or without batch normalization** to it. Show using a **convergence graph** [Marks 10: 5+5]
- **Visualise the Convolutional Features / Filters**. This could be done by using **imshow** or similar methods. Show how filters features **changes over different layers over a test image**. [Marks 20: 10+10]