

# Attention Is All You Need

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar\*  
Google Research  
nikip@google.com

Jakob Uszkoreit\*  
Google Research  
usz@google.com

Llion Jones\*  
Google Research  
llion@google.com

Aidan N. Gomez\* †  
University of Toronto  
aidan@cs.toronto.edu

Lukasz Kaiser\*  
Google Brain  
lukasz.kaiser@google.com

Illia Polosukhin\* ‡  
illia.polosukhin@gmail.com

## Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

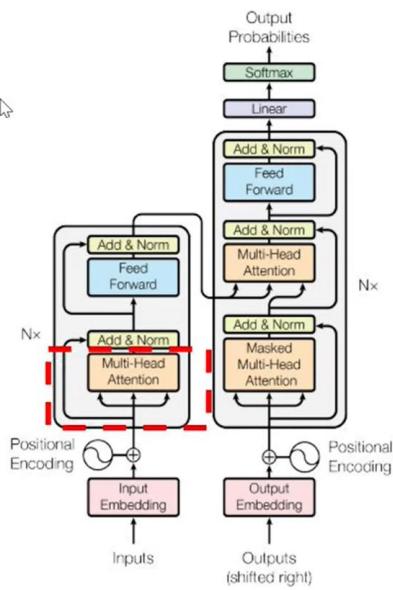


Figure 1: The Transformer - model architecture.

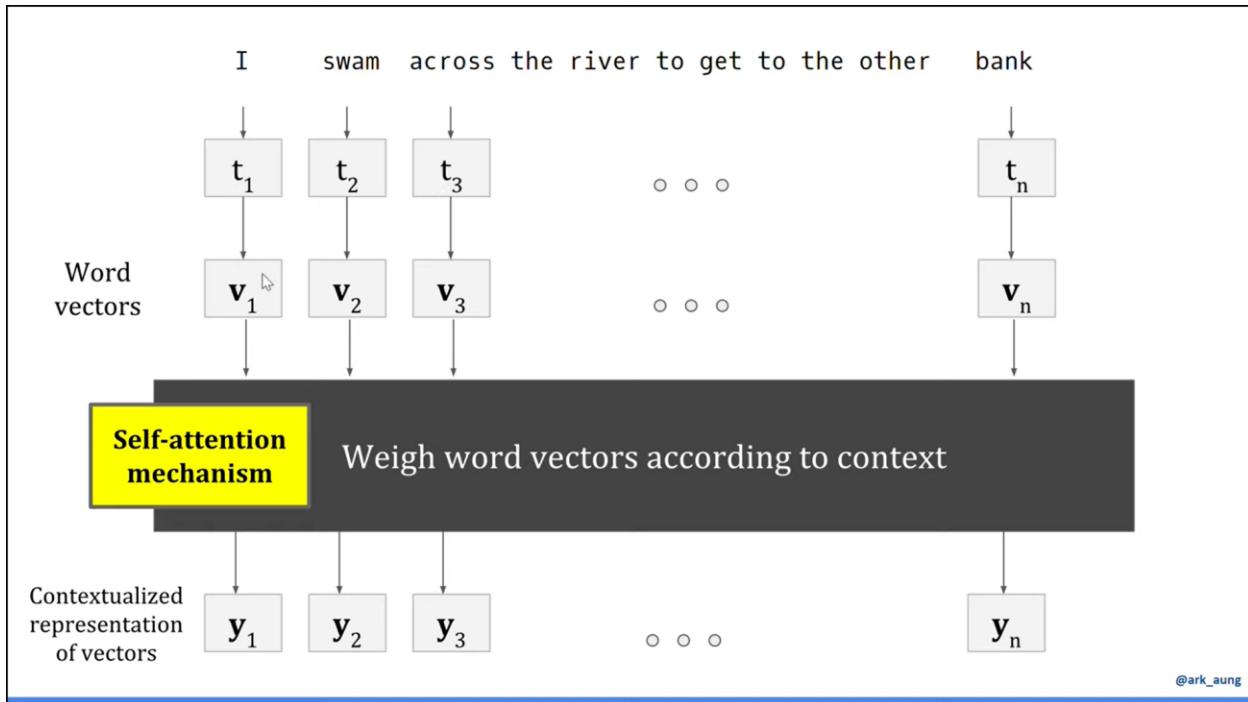
Attention is All You Need [Vaswani et al.]

@ark\_aung

Self attention mechanism add the contextual information in the sentence.

This contextual information added the weight to each words accroding to the context of sentence. Due to this similar word/token/vector have similar direction. NN weigh the words according to their context

It show the word vector that have same semantic meaning are very close to each other.

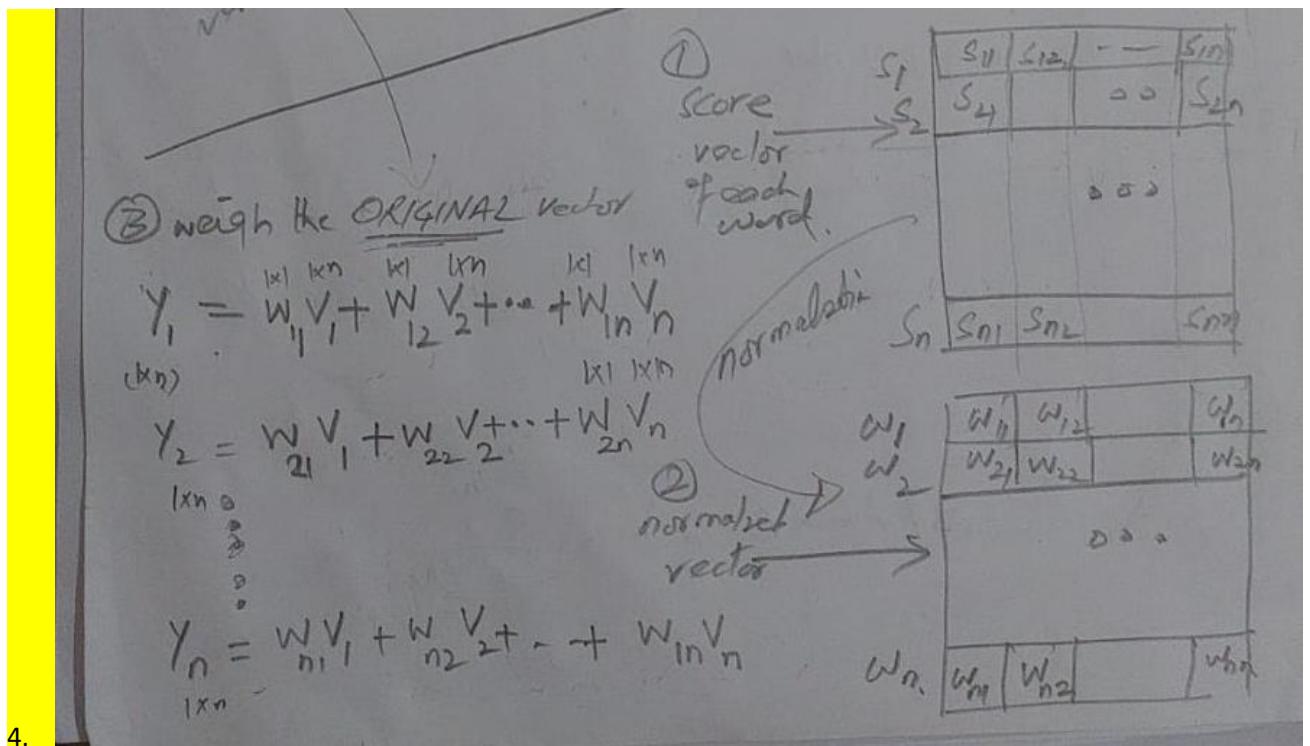


After weighing the word vectors according to their context, we go to the contextual representation of the words.

Here after weighing the word (applying the self attention mechanism), bank, swam, and river words have similarity/ closer values.

Three steps:

1. Dot product between vectors to get the scalar values
2. Normalized the values to get the weights
3. Weigh the vectors according the weight and get the contextualized representation of every vectors.



I swam across the river to get to the other **bank**

Bank == financial institution?  
Bank == sloping raised land?

I swam across the river to get to the other **bank**

Bank == financial institution?  
Bank == sloping raised land?

You shall know a word by the company it keeps!  
J.R.Firth(1957)

@ark\_aung

Q. What is the meaning of work bank?

A. You shall know a word by the company it keep. .JR Firth(1957)

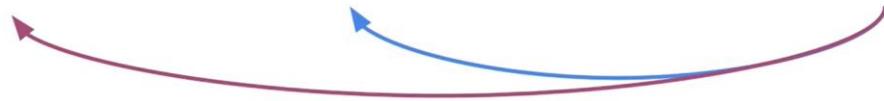
Therefore, bank is not about the financial intitute as if we see their neighborhood words.

If the sentence is

"I deposited my paycheck at the **bank** yesterday."

In this sentence, "bank" refers to a financial institution where money is stored, managed, and deposited.

I swam across the river to get to the other **bank**



Bank == financial institution?

Bank == sloping raised land?

@ark\_aung

I swam across the river to get to the other **bank**



~~Bank == financial institution~~

Bank == sloping raised land

@ark\_aung

I deposited my paycheck at the **bank** yesterday.



Bank == financial institution  
~~Bank == sloping raised land~~

Things changes correpsonding to the contax of the sentences



I swam across the river to get to the other **bank**

## Context matters!

@ark\_aung

Context is important to define the meaning of word



in Neural network

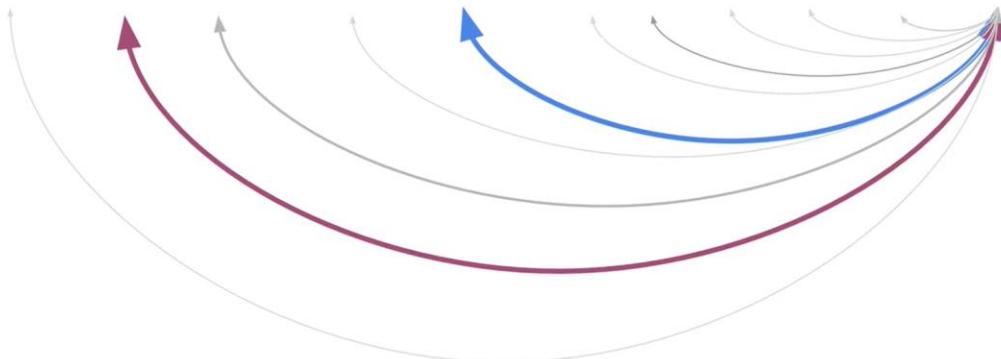
Can we build a **mechanism** that **weights** neighboring words to enhance the meaning of the word of interest?



@ark\_aung

Can we build the mechanism in the NN that weigh the words according to their context

I **swam** across the **river** to get to the other **bank**



@ark\_aung

\*\* If we weigh the words according to their context, then the words/vectors like bank, river and swam are very close to each other

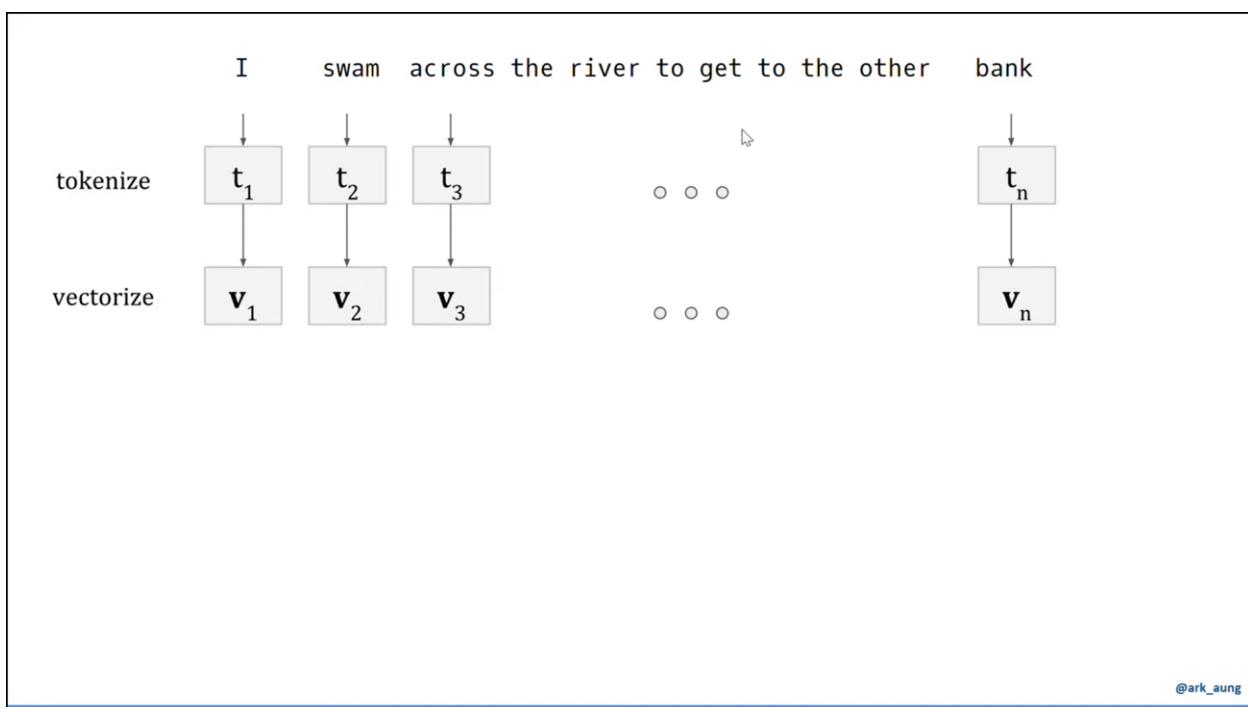
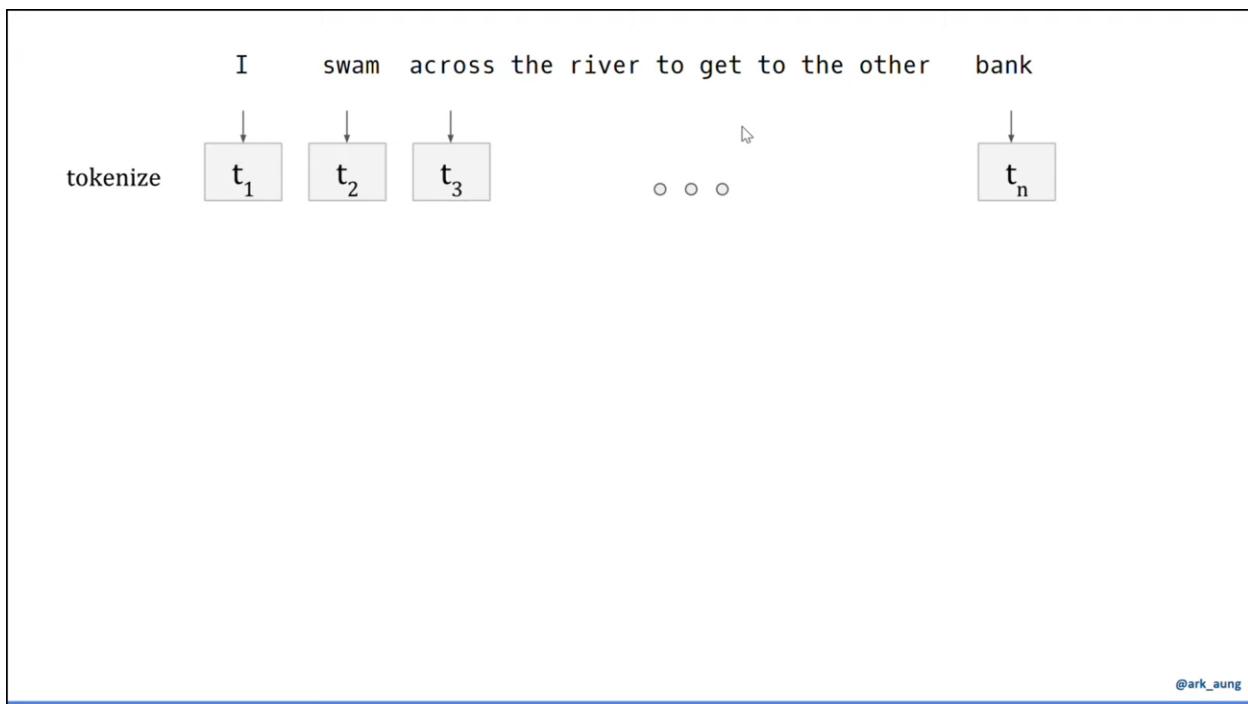
Swam and river are seems to be more closely related words to the bank, compare to the other words.

I **swam** across the **river** to get to the other **bank**

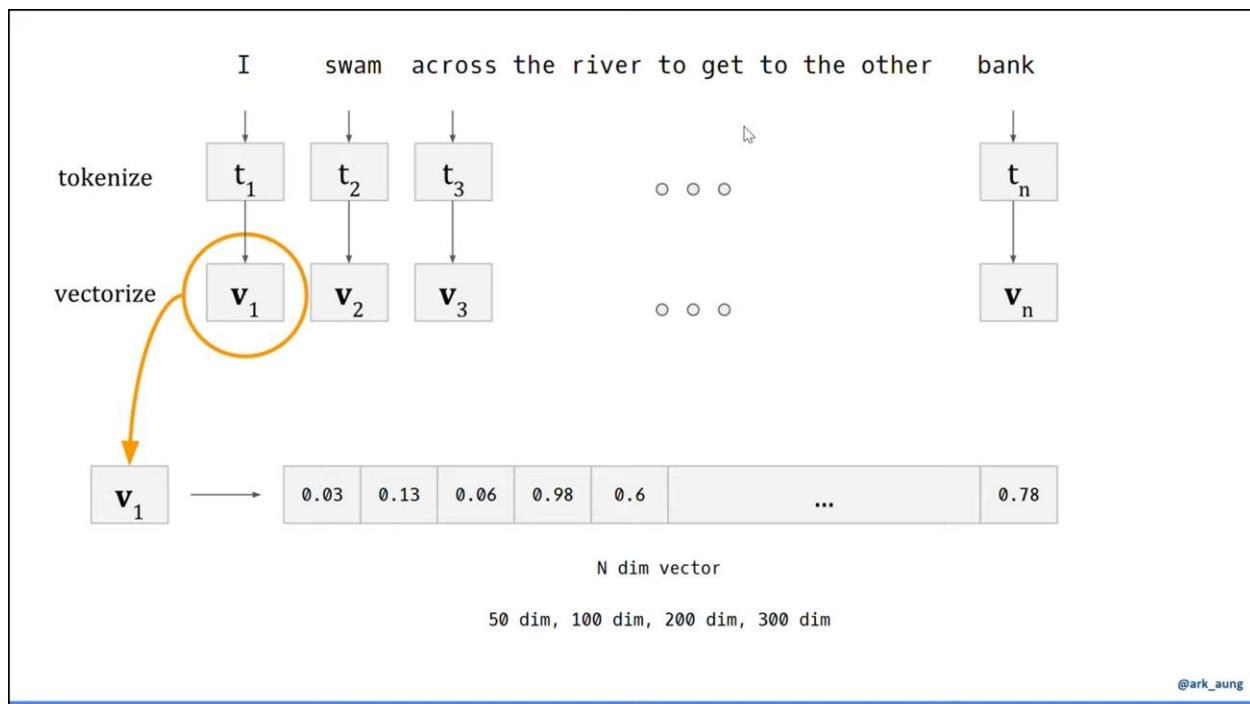
The main purpose of self-attention mechanism is to add **contextual information** to words in the sentence.

@ark\_aung

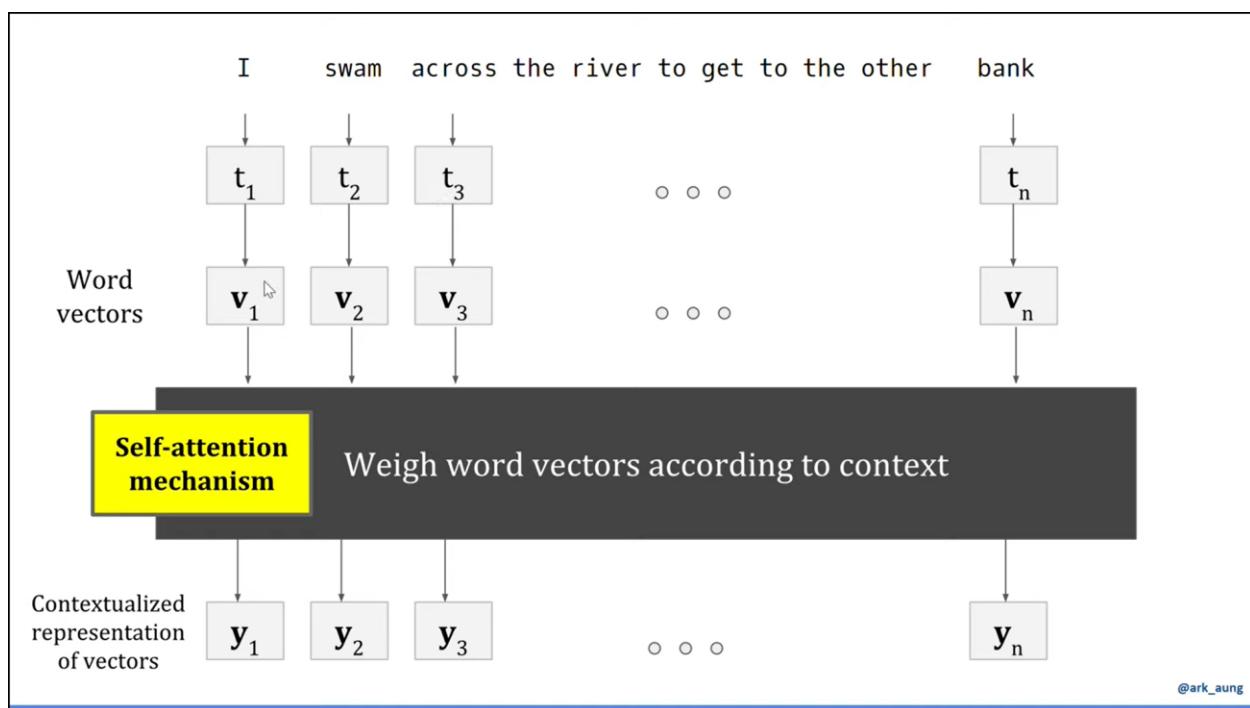
let see more tangible example



**Word to vector conversion : as the NN doesn't understand the words, so need to convert them into vectors. NN can understand the vectors.**



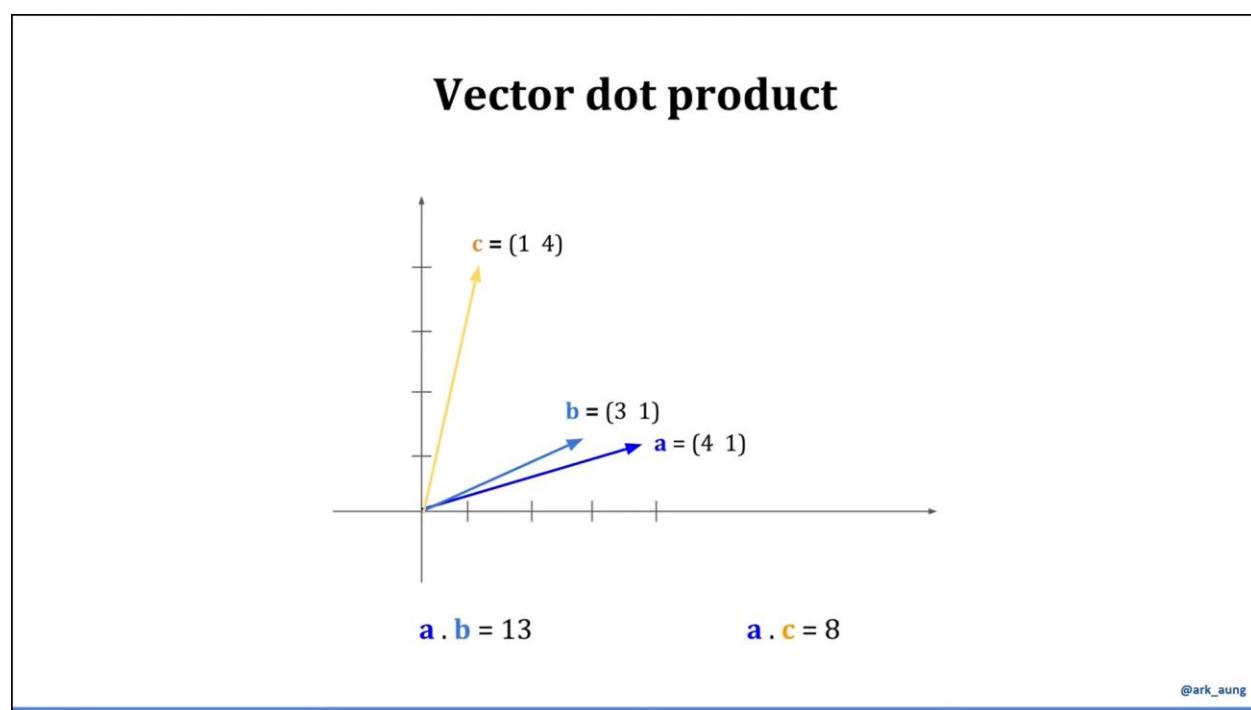
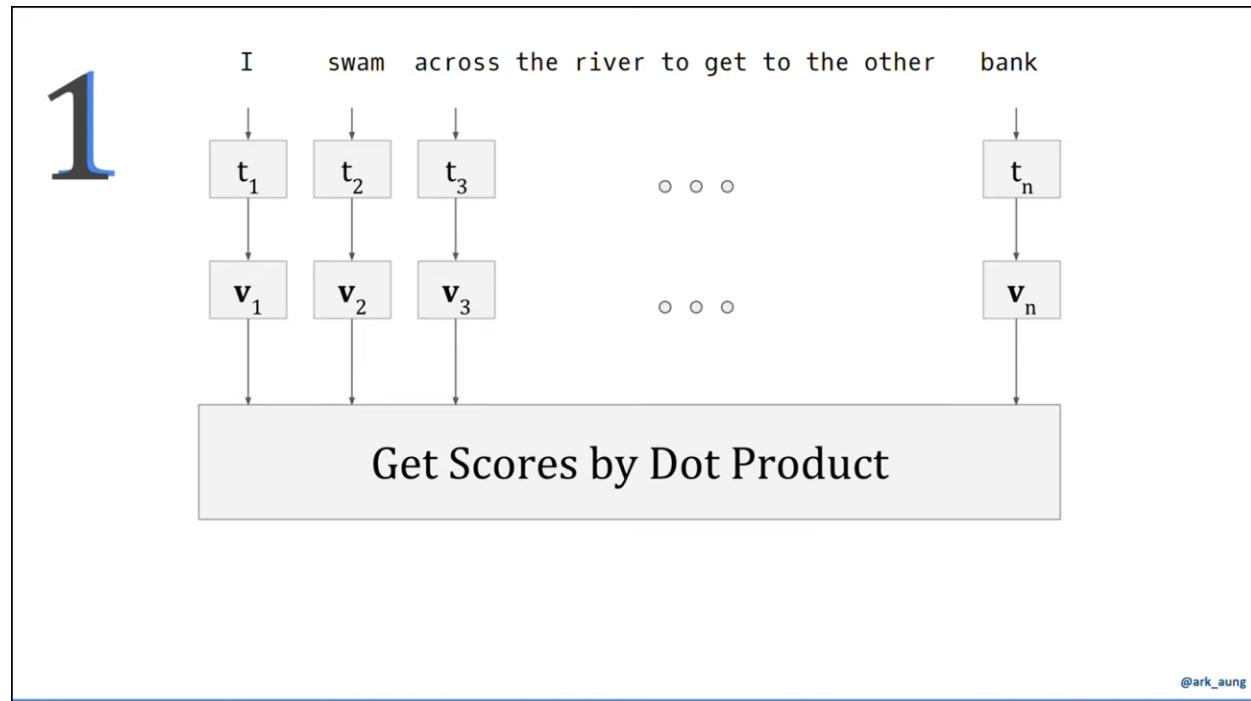
Word with same semantic meaning have the same vector values as shown in the following slides



After weighing the word vectors according to their context, we go to the contextual representation of the words.

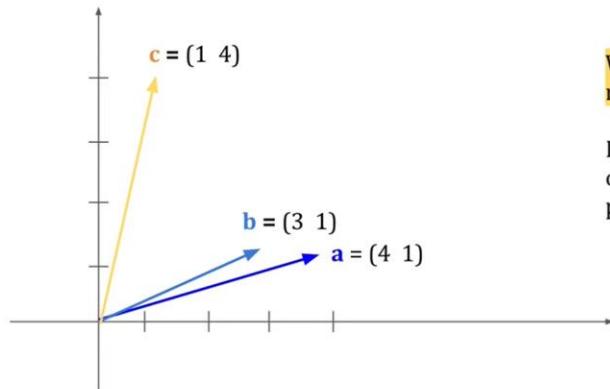
Here after weighing the word (applying the self attention mechanism), bank, swam, and river words have similiy/ closer values.

Folloing is the strategy we adapt for getting the contextual representation of vectors



Here the words are in the form of vectors. Dot produce between perpendicular vector is zero if there are not close to each other.

## Vector dot product



What I want you to remember:

If two vectors are close to one another, their dot product is going to be big

$$\mathbf{a} \cdot \mathbf{b} = 13$$

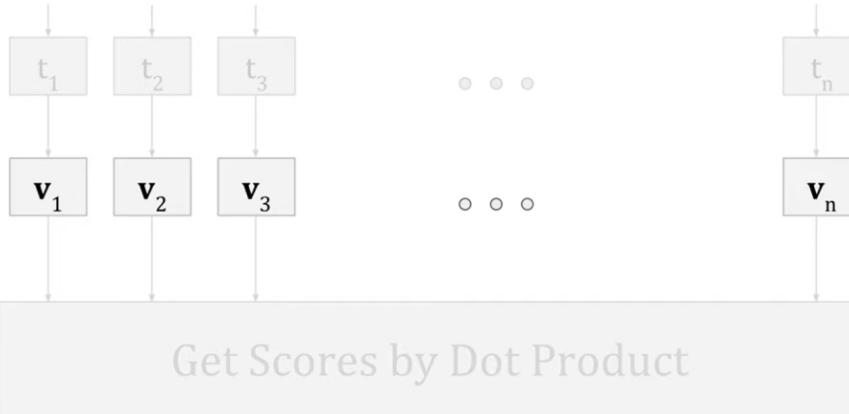
$$\mathbf{a} \cdot \mathbf{c} = 8$$

@ark\_aung

1

I swam across the river to get to the other bank

1x50



Get Scores by Dot Product

@ark\_aung



$$S_{11} = \mathbf{v}_1 \mathbf{v}_1^T$$

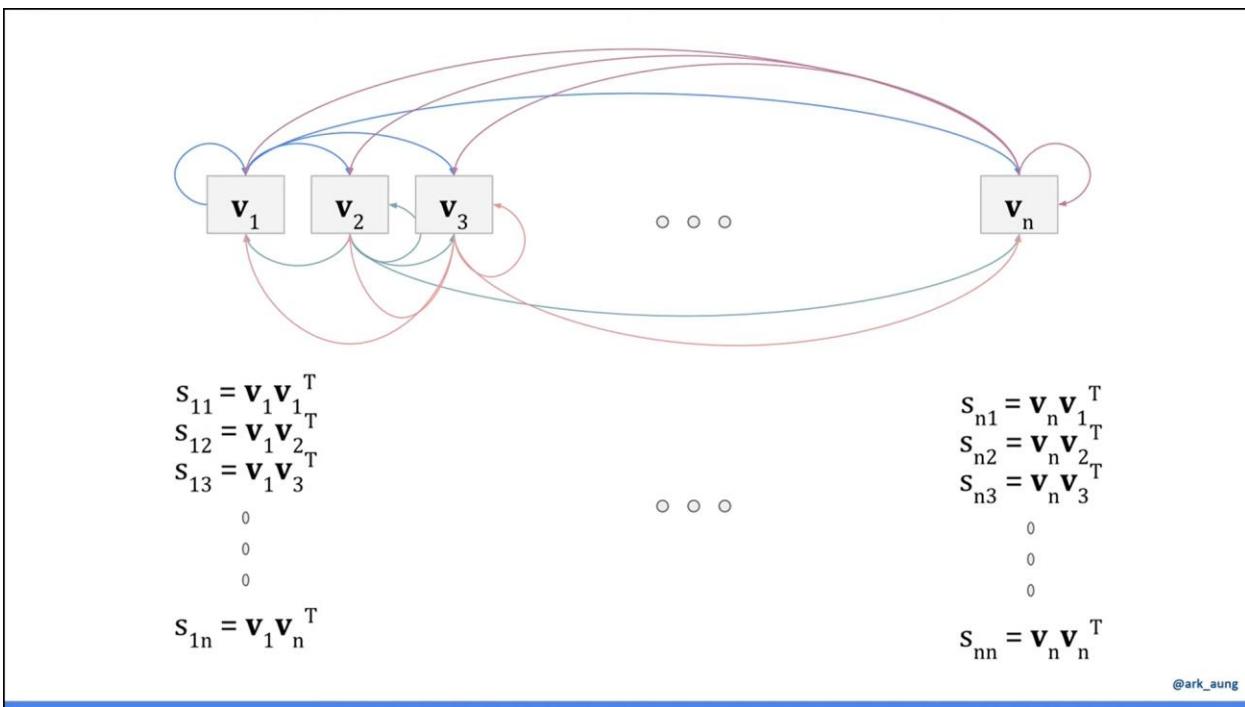
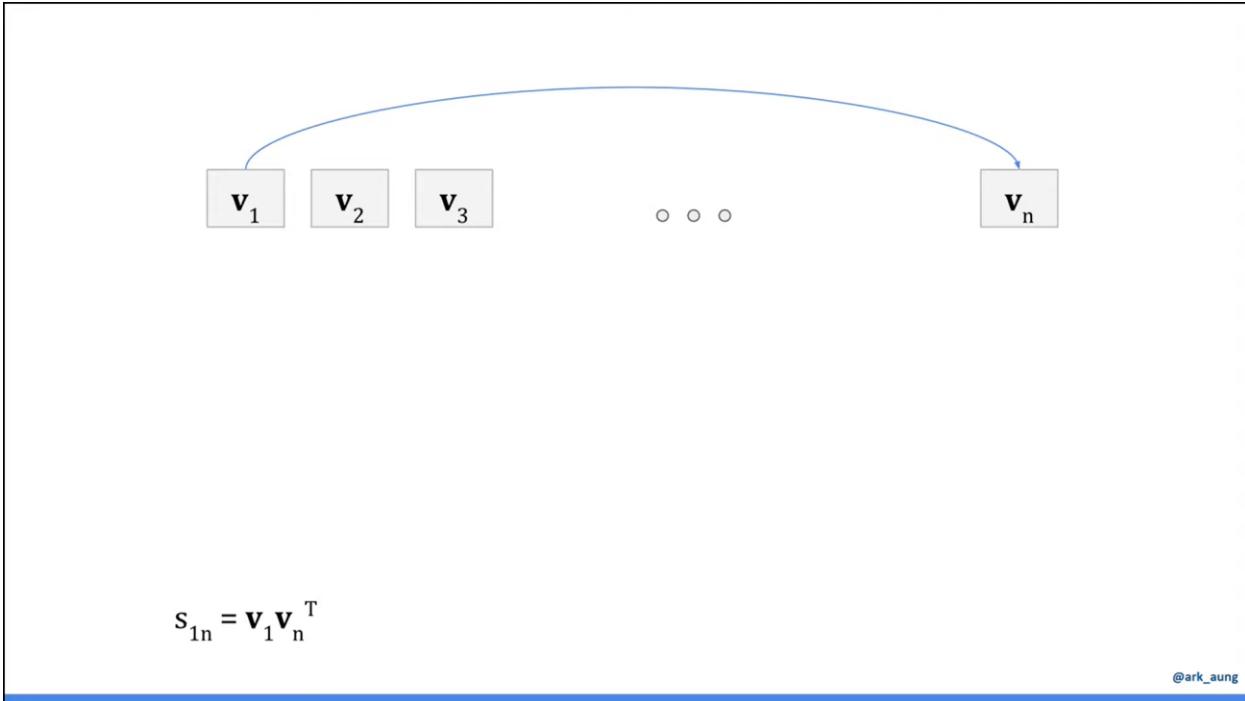
1x50   50x1

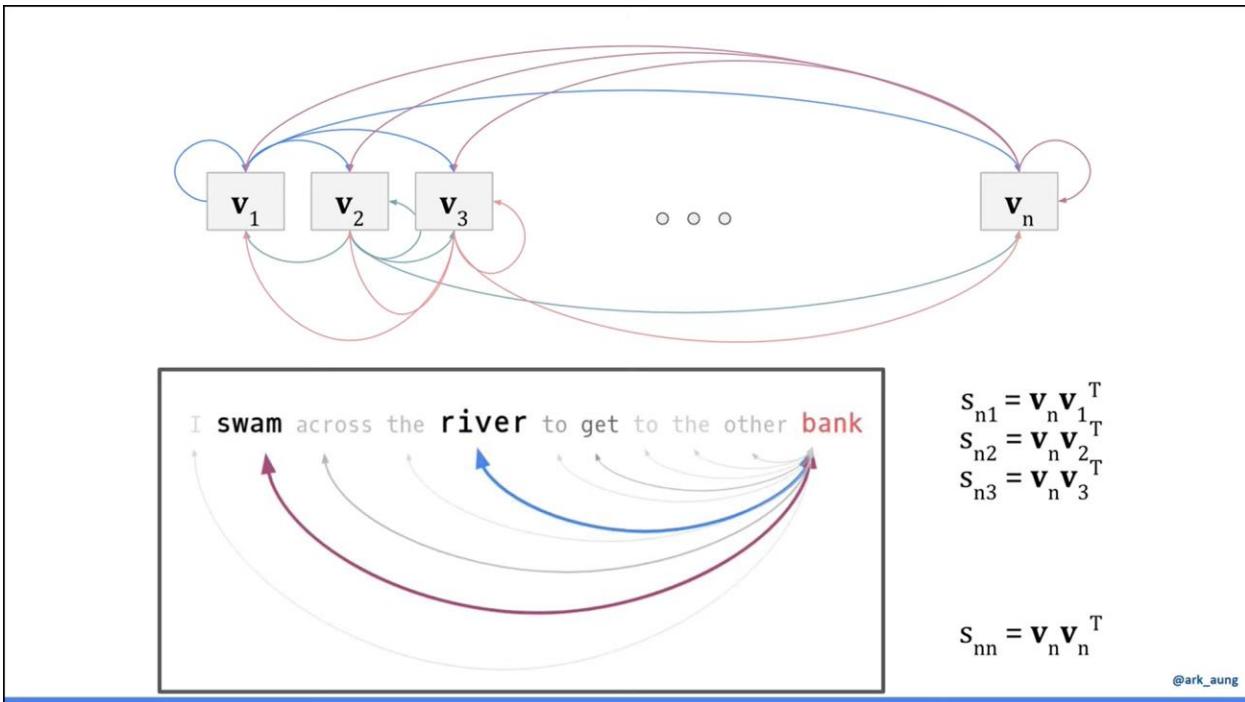
@ark\_aung



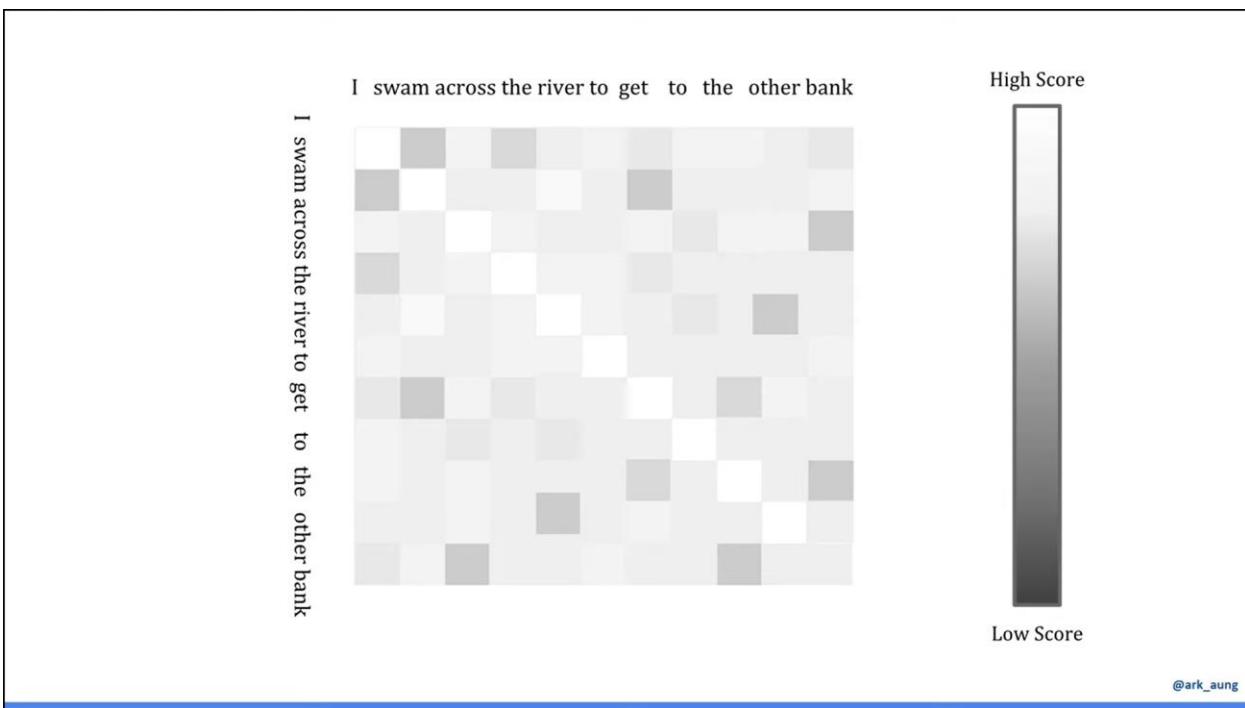
$$S_{12} = \mathbf{v}_1 \mathbf{v}_2^T$$

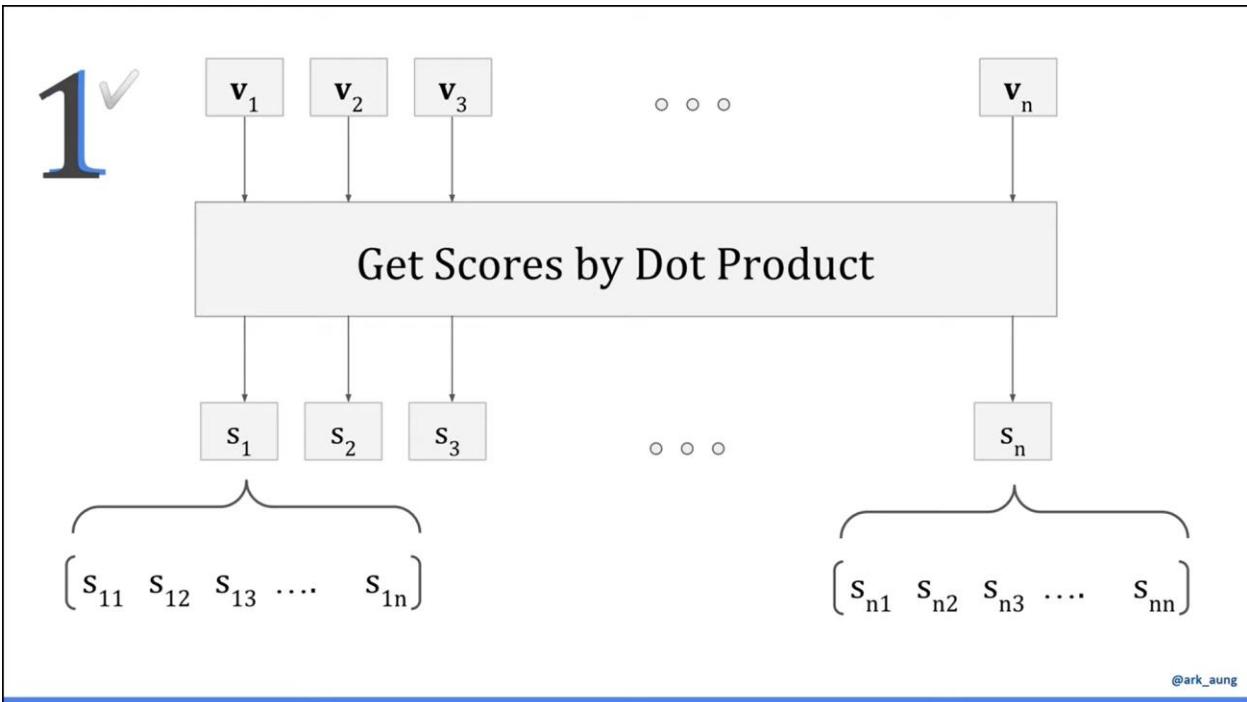
@ark\_aung





**Graphically we can see the higher values are at diagonal**





# 2



@ark\_aung

Normalized because range can be very large between the dot product values, so we need to normalized it using softmax

The softmax normalization ensures that:

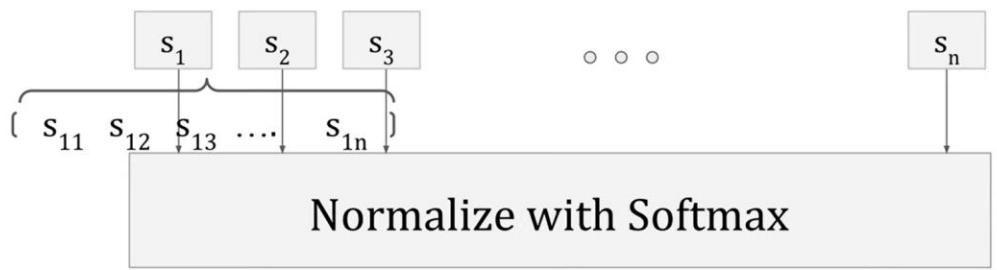
- Scores are interpretable as probabilities.

### 1. Convert Scores into Probabilities:

The raw similarity scores (dot products between query and key vectors) can be arbitrary real numbers. These raw values don't directly represent how much focus should be placed on each token. The softmax function converts these scores into a probability distribution, where all attention scores sum to 1.

- This makes it easy to interpret the results: higher probability means more attention is given to that token.

For example, if we have raw similarity scores like [3,1,-2], applying softmax converts them into probabilities like [0.84,0.13,0.03], where the first token gets the most attention.



$$\begin{aligned} w_{11} &= \text{softmax}(s_{11}) \\ w_{12} &= \text{softmax}(s_{12}) \\ w_{13} &= \text{softmax}(s_{13}) \end{aligned}$$

0  
0  
0

$$w_{1n} = \text{softmax}(s_{1n})$$

$$\sigma(s_i) \triangleq \frac{e^{s_i}}{\sum_{j=1}^n e^{s_j}}$$

0 0 0

$$\sum_{i=1}^n w_{1i} = 1$$

$$\begin{aligned} w_{n1} &= \text{softmax}(s_{11}) \\ w_{n2} &= \text{softmax}(s_{12}) \\ w_{n3} &= \text{softmax}(s_{13}) \end{aligned}$$

0  
0  
0

$$w_{nn} = \text{softmax}(s_{nn})$$

@ark\_aung

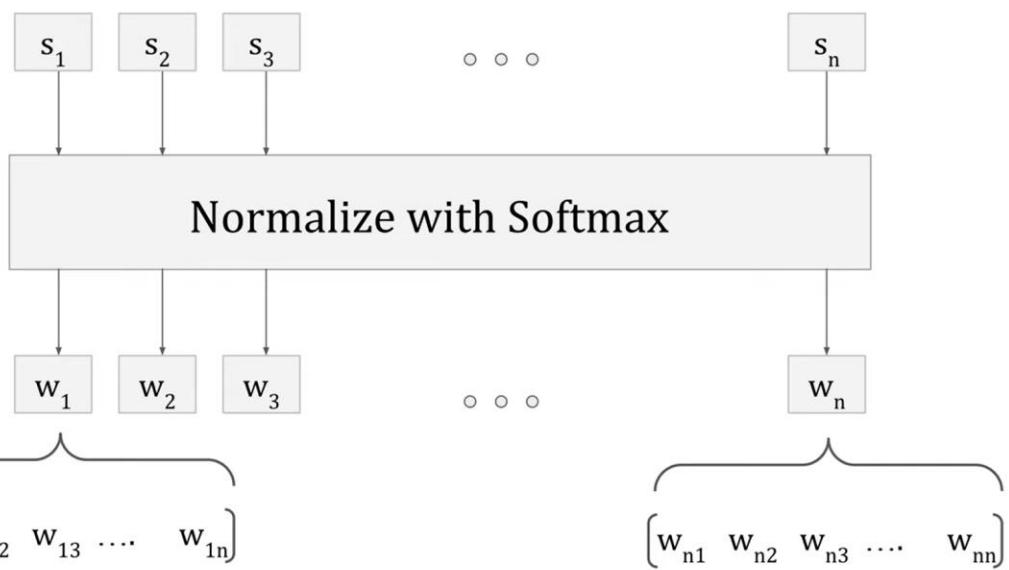
$$\mathbf{S} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$\mathbf{W} = \frac{\mathbf{e}^{\cdot} \mathbf{s}}{\sum_{i=1}^n e^i} \mathbf{s}$$

$$\mathbf{W} = \frac{\begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}}{(e^5 + e^2 + e^{-1} + e^3)}$$

$$\mathbf{W} = \frac{\begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}}{176.3} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

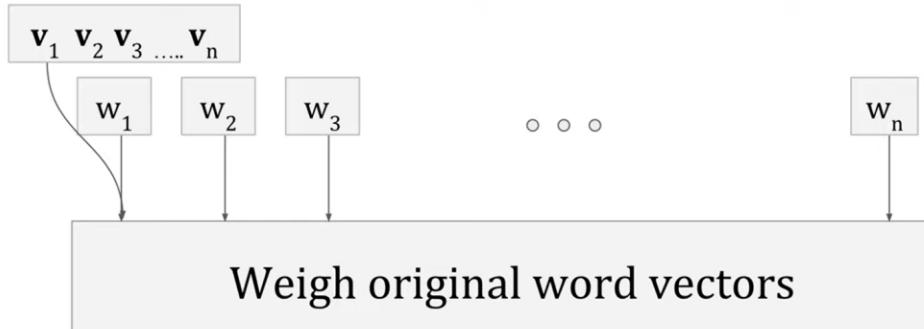
2 ✓



@ark\_aung

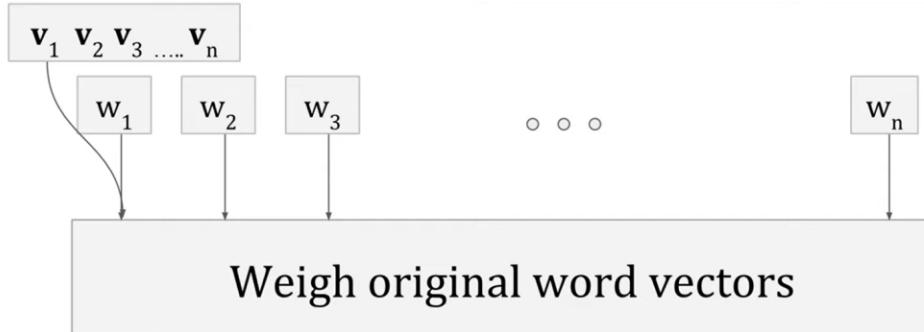
We get weights from the dot produce scores

3



@ark\_aung

# 3



$$y_1 = w_{11}v_1 + w_{12}v_2 + w_{13}v_3 + \dots + w_{1n}v_n$$

$$y_2 = w_{21}v_1 + w_{22}v_2 + w_{23}v_3 + \dots + w_{2n}v_n$$

$$y_3 = w_{31}v_1 + w_{32}v_2 + w_{33}v_3 + \dots + w_{3n}v_n$$

0

0

0

$$y_n = w_{n1}v_1 + w_{n2}v_2 + w_{n3}v_3 + \dots + w_{nn}v_n$$

@ark\_aung

W's are weights and v are the vectors. Y1 is the contextual representation of vecto v1.

Y1 is got by gving w11 weight to itself (v1)

w12 v2 is the influence of v2 on the y1

w13 v3 is the influence of v3 on the y1

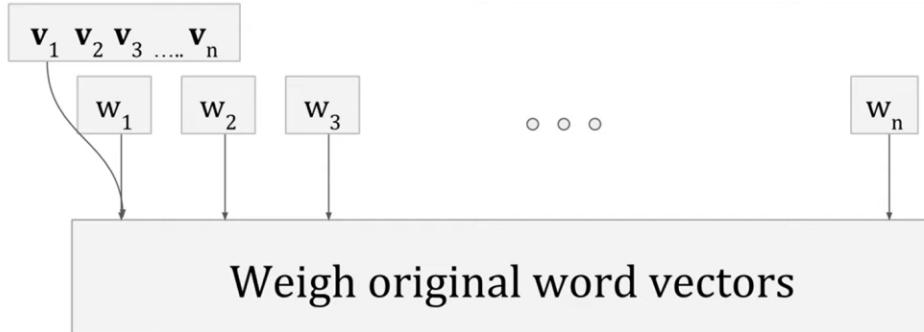
....

**Now y1 is the contextual respresentation of v1. That is it is no longer pointing towards the same direction as before. It is pointing towards the different direction, which is influenced by its neighbours**

**V1 is the orignal vector, while y1 is the enhanced version v1/ augmented version of v1 by weighing its sorrounding to itself by following process:**

1. Get score by produce
2. Normalize
3. Weigh the original vectors.

# 3



$$\mathbf{y}_1 = w_{11} \mathbf{v}_1 + w_{12} \mathbf{v}_2 + w_{13} \mathbf{v}_3 + \dots + w_{1n} \mathbf{v}_n$$

$$\mathbf{y}_2 = w_{21} \mathbf{v}_1 + w_{22} \mathbf{v}_2 + w_{23} \mathbf{v}_3 + \dots + w_{2n} \mathbf{v}_n$$

$$\mathbf{y}_3 = w_{31} \mathbf{v}_1 + w_{32} \mathbf{v}_2 + w_{33} \mathbf{v}_3 + \dots + w_{3n} \mathbf{v}_n$$

0

0

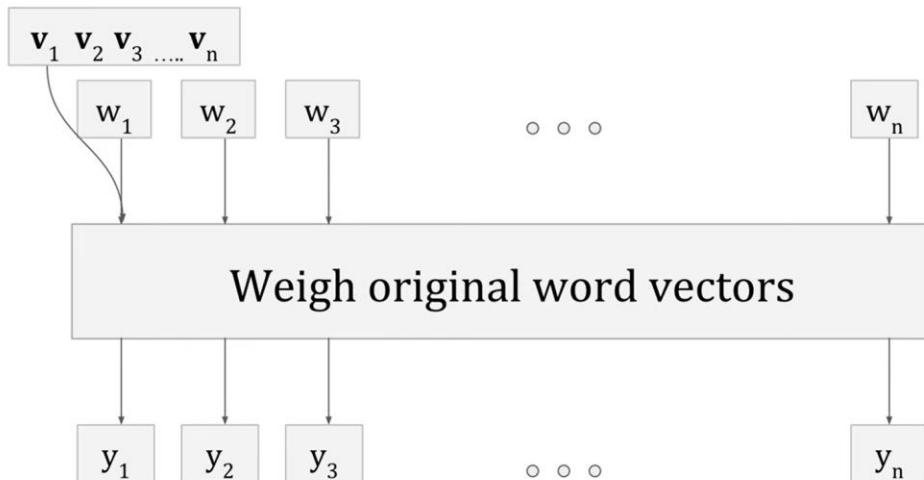
0

$$\mathbf{y}_n = w_{n1} \mathbf{v}_1 + w_{n2} \mathbf{v}_2 + w_{n3} \mathbf{v}_3 + \dots + w_{nn} \mathbf{v}_n$$

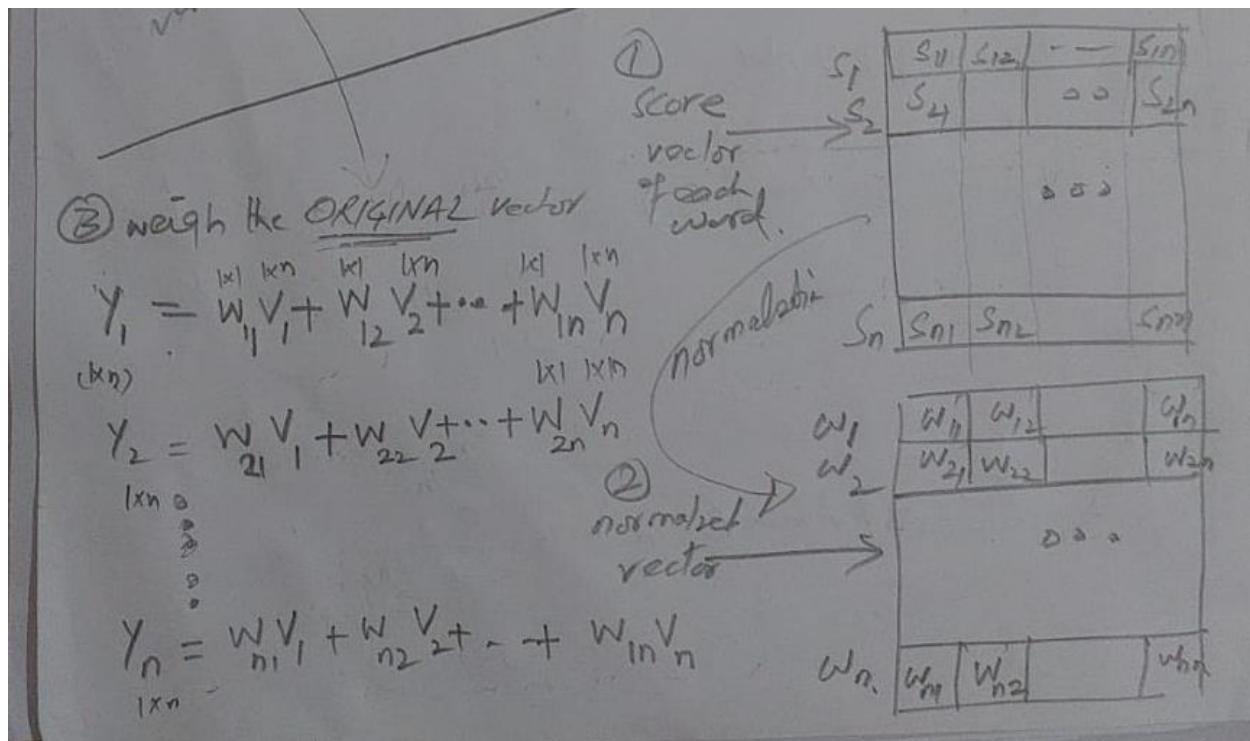
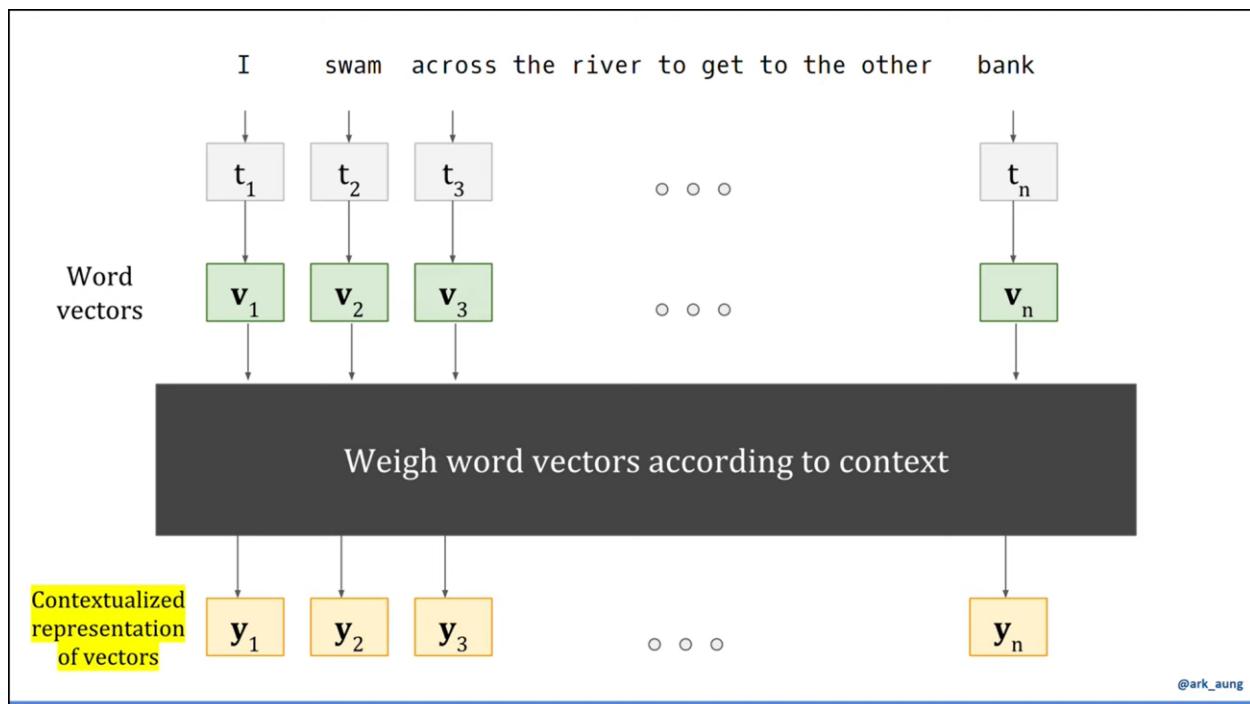
@ark\_aung

**Contextual Representation:** The result of the weighted sum of the value vectors produces a new representation for the query word that is informed by the context of other words in the sequence. This representation captures not just the identity of the query word but also relevant information from other words, effectively enriching its context.

# 3 ✓

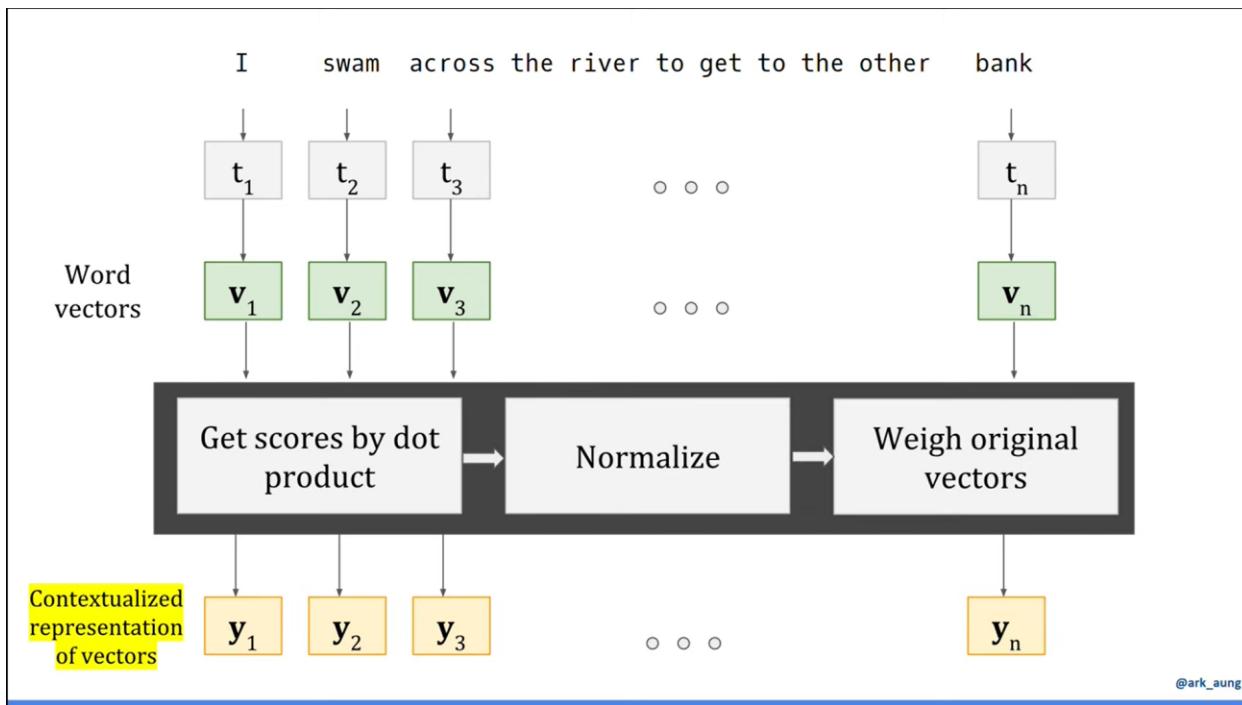


@ark\_aung



## Lecture#2

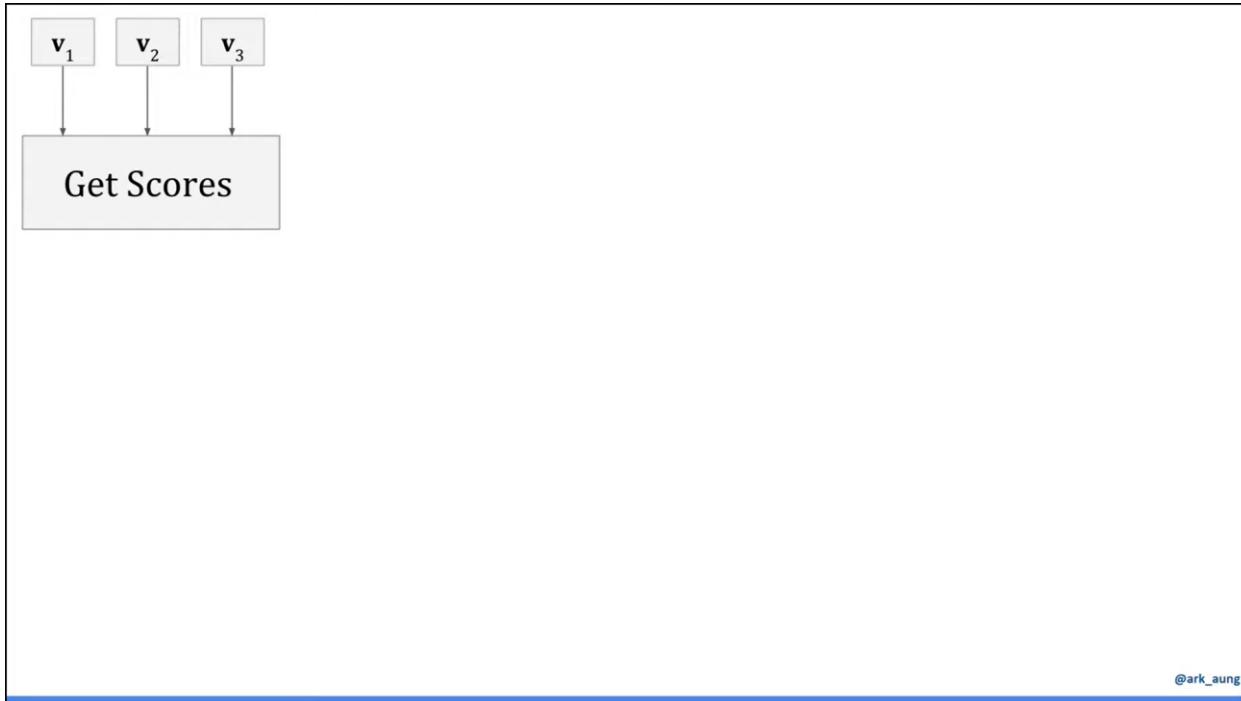
### 1. recap



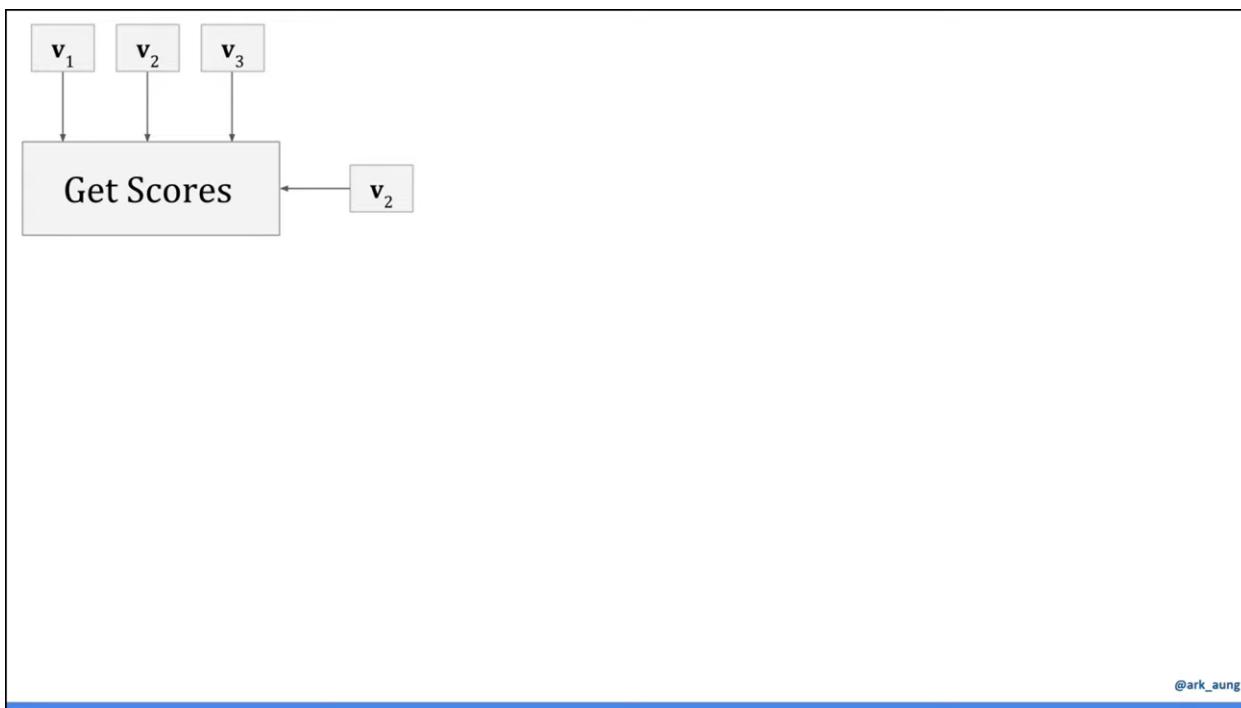
### 2. Database terms:

- **Query/question**= which word is close to  $v_2$  or am I close to  $v_1$  or am I close to  $v_2$  etc
- **Key**:  $v_1, v_2, v_3$  are database value that needs to be compare
- **Values**:  $v_1, v_2, v_3$  are the values that used to get the contextulazed reprsentation to  $v_2$

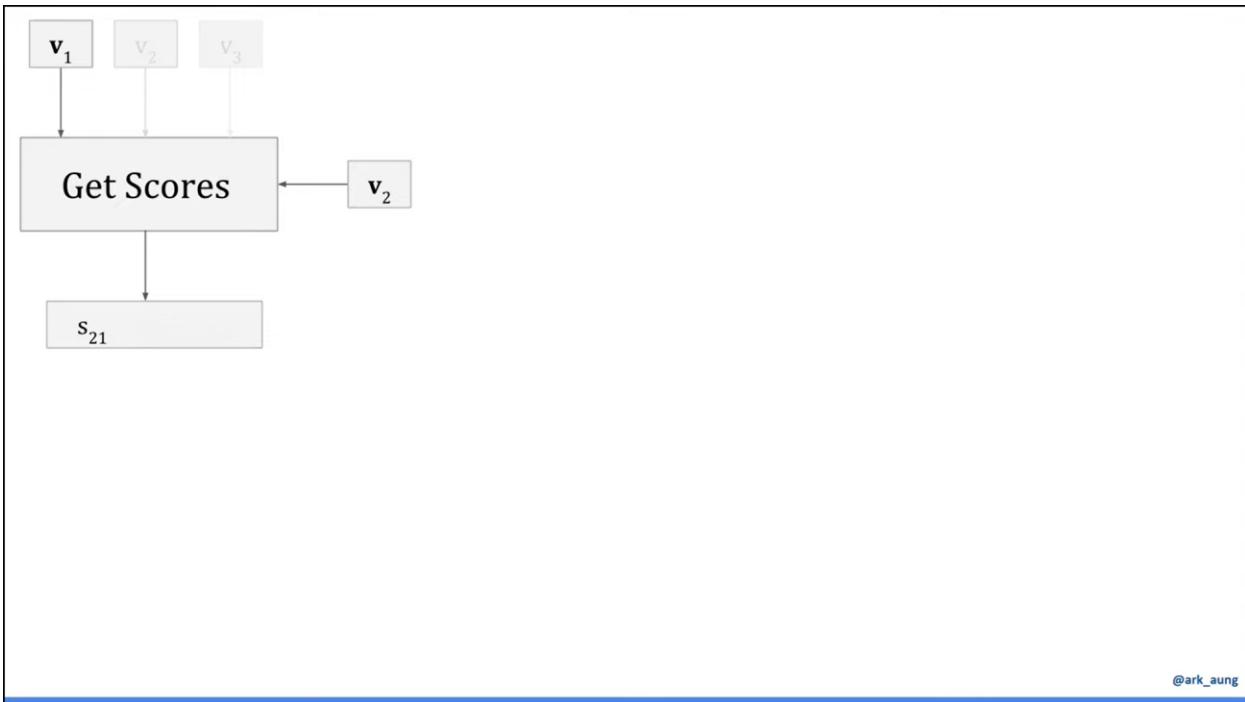
3. Example: Lets see more concrete example of three words and want to find the contextualized representation of  $v_2$ .



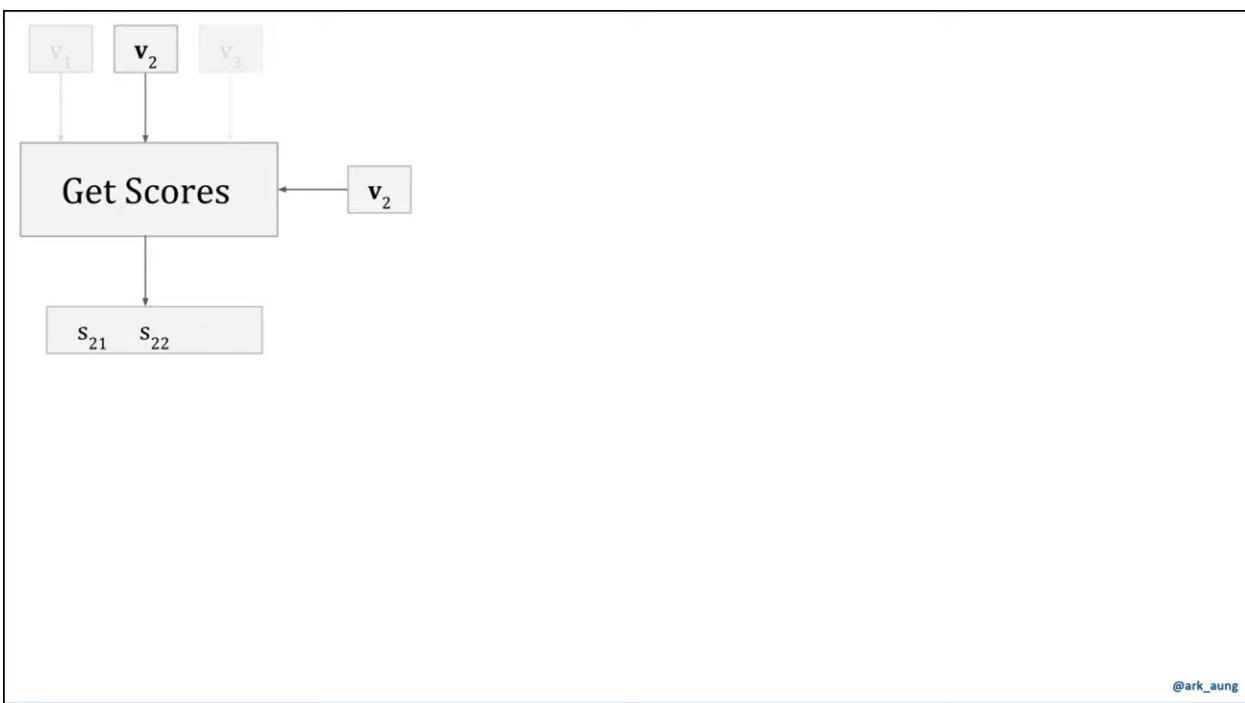
@ark\_aung



@ark\_aung



@ark\_aung



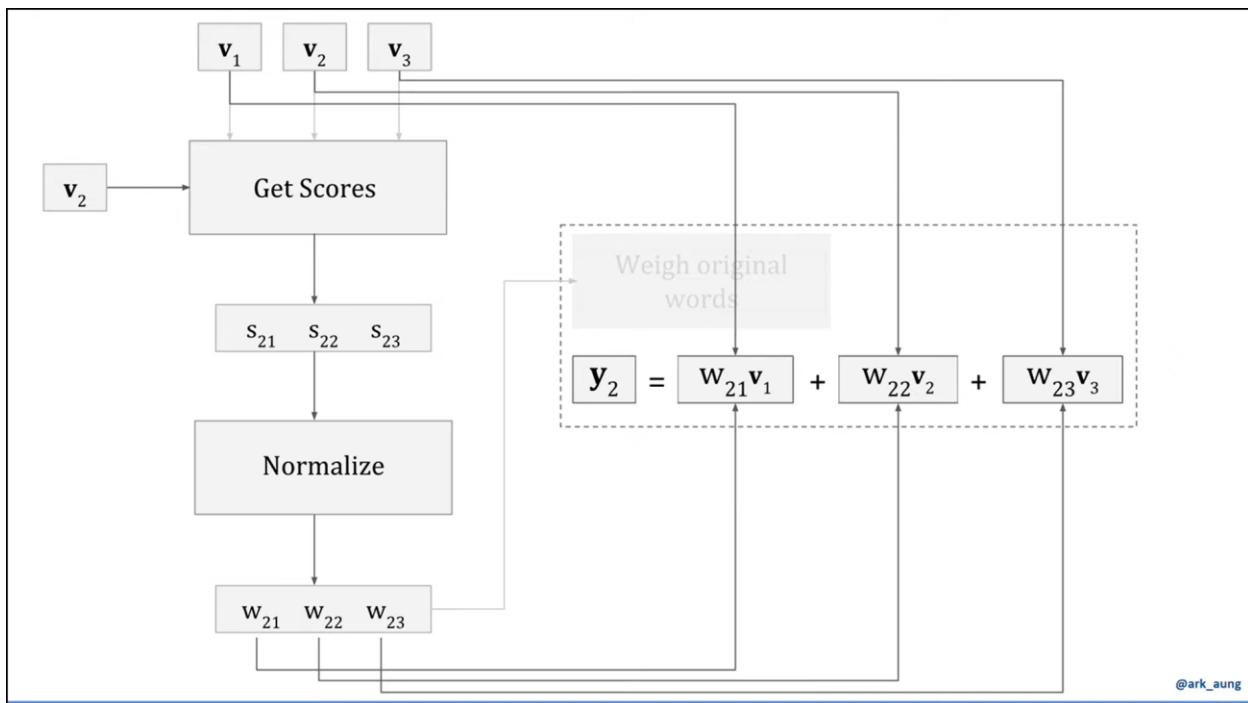
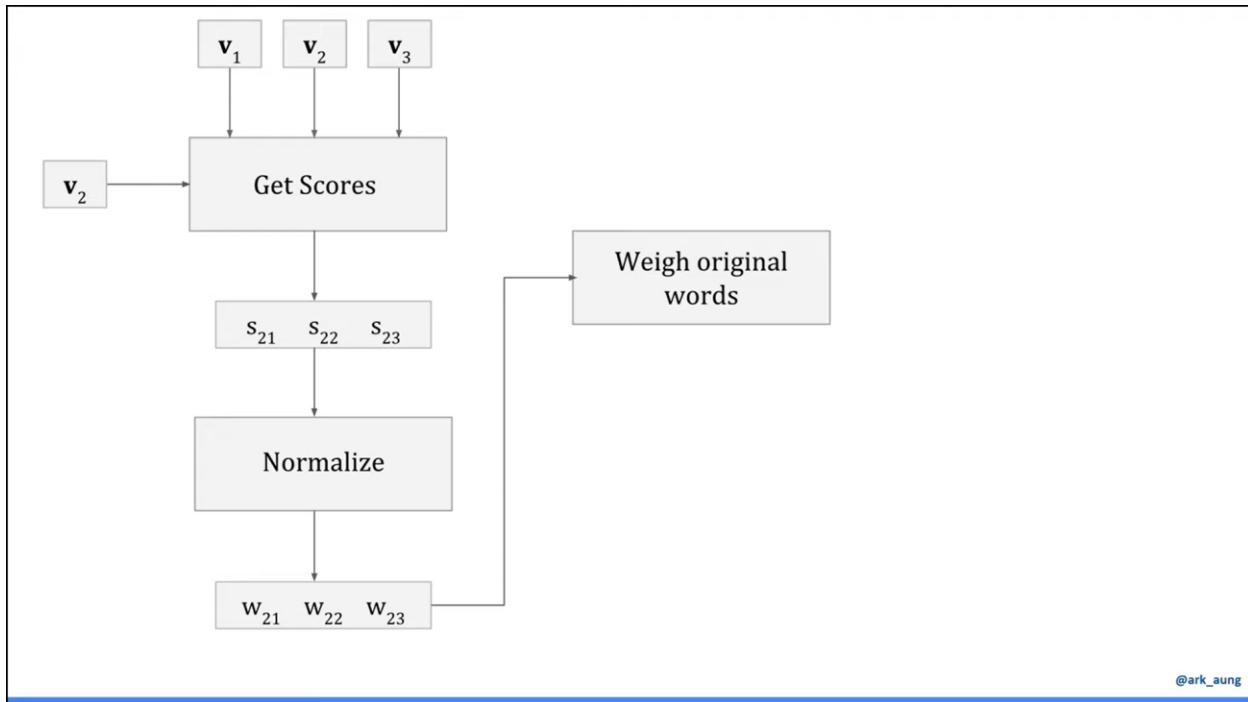
@ark\_aung



@ark\_aung

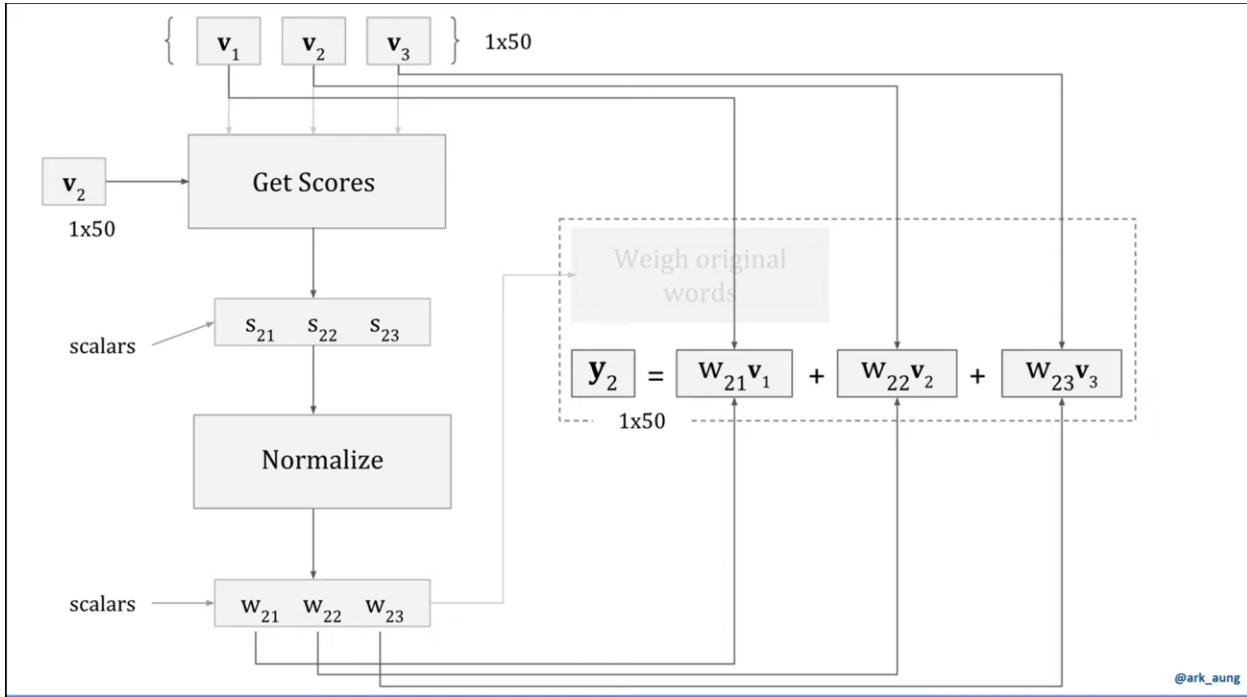


@ark\_aung

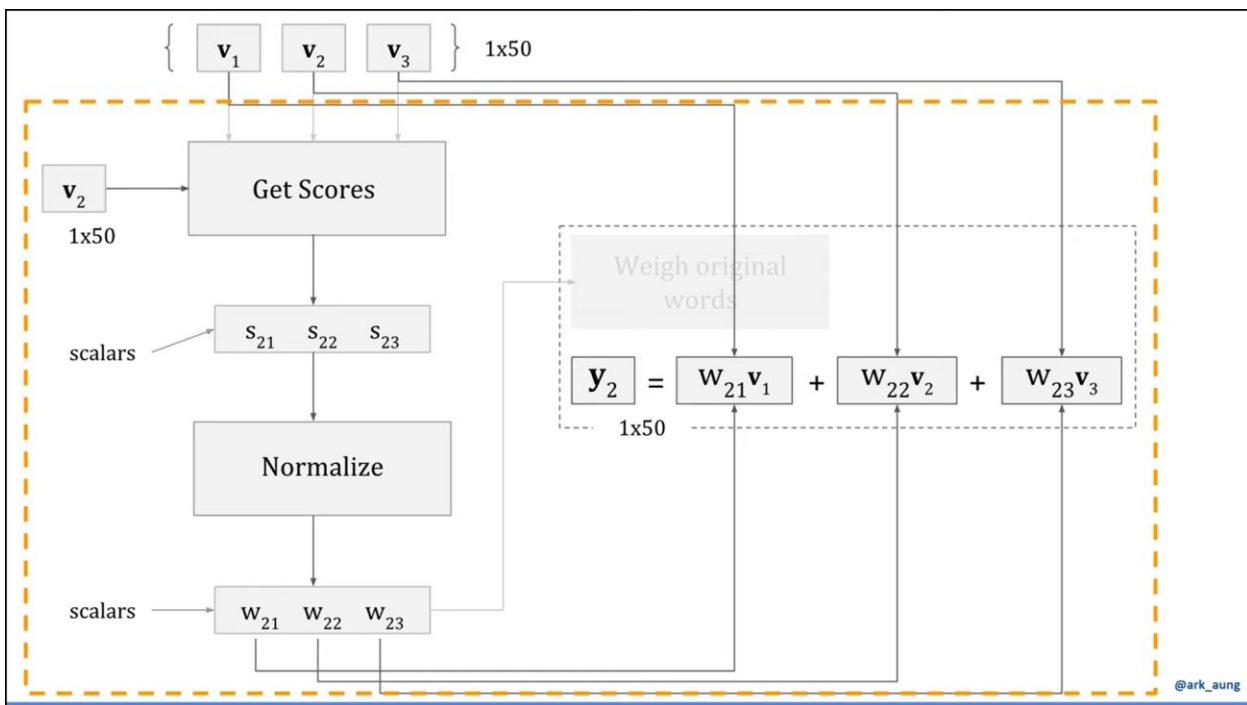


y2 is contextulized represntation of orignal vector v2

#### 4. dimension



@ark\_aung



@ark\_aung

#### 5. Orange box is he self attentionor attention mechanism

Which word is closer to  $v_2$ , With respect to database, it is query/question , and keys are other words that need to check with  $v_2$  to check either they are similar to  $v_2$  or not, values..... we need to inject the weights to preserve the diemension

## 6. Where are the learning????

for example of translation where the model output the contextualized vectors and now the decoder generate the french sentence using the contextualized information but if the prediction is not good then loss will be higher.... So there is need to update that could also be corresponding to query, key, and value vectors..... so that we able to produce the good results.

Self attention mechanism learn to focus on important parts of a sentence though training.

### Why Learning is Necessary

- Without learning, the attention mechanism would not adapt to different contexts, resulting in poor performance.
- With learning, the model dynamically adjusts the attention weights based on the meaning and context of the sentence, ensuring the model attends to the most relevant parts for the task at hand (e.g., translating "bank" correctly as a financial institution).

Learning allows the model to capture context-specific relationships, ensuring that it performs well across various tasks and contexts.

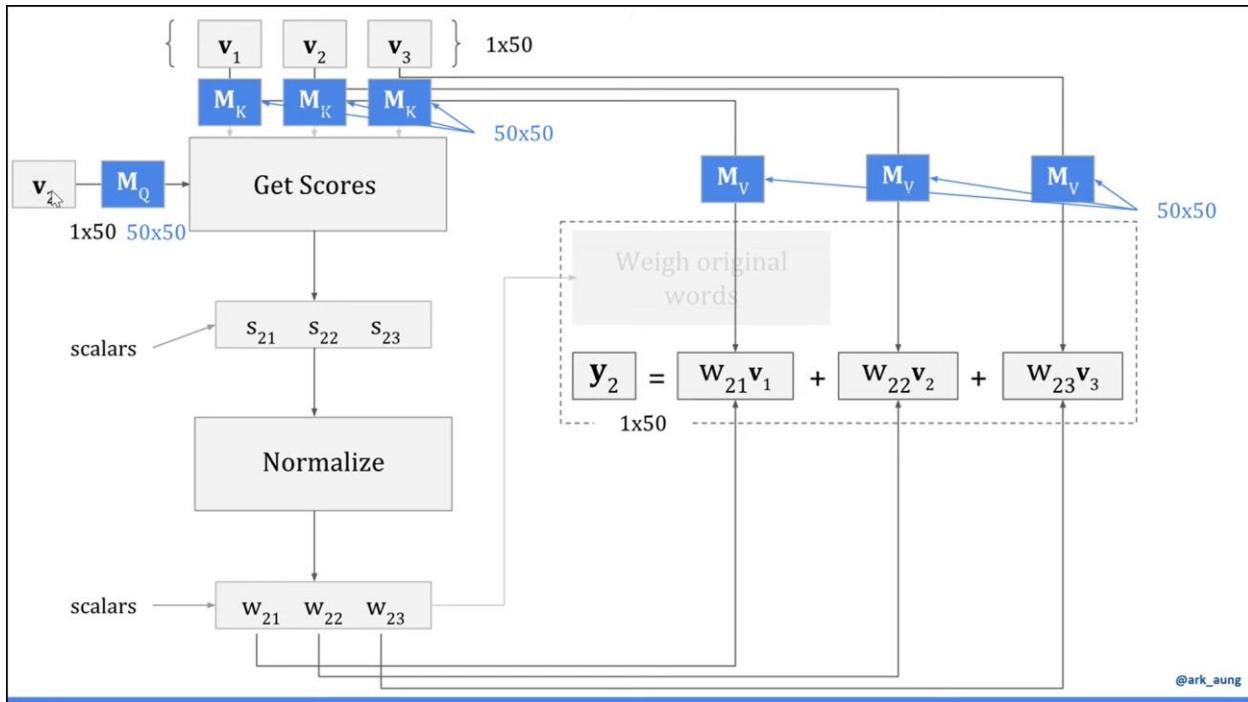
For example, in the sentence "I deposited money at the bank", after training, the model learns to:

- Pay attention to "money" when processing "bank" to understand that this is a financial institution.
- Ignore irrelevant words like "at" and "the" for understanding the meaning of "bank."

### Step 5: Training and Backpropagation

The learning process occurs through training. When the model translates sentences, it compares its output to the correct translation (or output). A loss function is used to measure the difference between the model's prediction and the actual correct answer.

- The model adjusts the query, key, and value weights through backpropagation, updating the weights to improve the attention scores and overall performance on the task.
  - Over time, the model learns how much attention to pay to various words based on context. It learns to focus on the right words to improve accuracy.
-



These weights are learned during BP, while the dimension of outputs are same .... These weights are the dense layers. We have 3 linear layers one for query, one is for value and one for key

V1 is 1x50 dimensional vector if there is single linear layer having 50 neurons than no of parameters are 50x50. And matrix size is 50x50

Can you draw dense having 10 neurons and show forward and backward pass.

## 7. Toy example

### Example Contextualized representation of a vector:

Let's go through a simple, **toy example** to demonstrate how the **self-attention mechanism** works in the context of a Transformer model.

#### Example Sentence:

We'll use the short sentence: "The cat eats fish."

#### Objective:

Let's understand how each word in this sentence attends to other words using **self-attention**.

#### Step 1: Assign Embeddings

For simplicity, let's assume we have the following **embedding vectors** for each word (these are made-up numbers to keep the math simple):

- "The" → [1, 0]
- "cat" → [0, 1]
- "eats" → [1, 1]
- "fish" → [0, 0]

### Step 2: Create Query, Key, and Value Matrices

Let's also assume that we have the following learnable weight matrices for query ( $W_Q$ ), key ( $W_K$ ), and value ( $W_V$ ):

$$W_Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad W_K = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad W_V = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

### Step 3: Calculate Query, Key, and Value Vectors for Each Word

We calculate the query, key, and value vectors for each word using the weight matrices:

For the word "The":

- Query:  $Q_{\text{The}} = W_Q \cdot \text{Embedding}_{\text{The}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$
- Key:  $K_{\text{The}} = W_K \cdot \text{Embedding}_{\text{The}} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$
- Value:  $V_{\text{The}} = W_V \cdot \text{Embedding}_{\text{The}} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$

For the word "cat":

- Query:  $Q_{\text{cat}} = W_Q \cdot \text{Embedding}_{\text{cat}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$
- Key:  $K_{\text{cat}} = W_K \cdot \text{Embedding}_{\text{cat}} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$
- Value:  $V_{\text{cat}} = W_V \cdot \text{Embedding}_{\text{cat}} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$

For the word "eats":

- Query:  $Q_{\text{eats}} = W_Q \cdot \text{Embedding}_{\text{eats}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$
- Key:  $K_{\text{eats}} = W_K \cdot \text{Embedding}_{\text{eats}} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$
- Value:  $V_{\text{eats}} = W_V \cdot \text{Embedding}_{\text{eats}} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 7 \end{bmatrix}$

For the word "fish":

- Query:  $Q_{\text{fish}} = W_Q \cdot \text{Embedding}_{\text{fish}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$
- Key:  $K_{\text{fish}} = W_K \cdot \text{Embedding}_{\text{fish}} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$
- Value:  $V_{\text{fish}} = W_V \cdot \text{Embedding}_{\text{fish}} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

#### Step 4: Calculate Attention Scores (Dot Products between Query and Key)

We calculate the attention score for each pair of words by taking the dot product between the query of one word and the key of another.

Let's focus on the word "eats" and compute how much it attends to other words:

1. Attention score between "eats" and "The":

$$\text{Attention score} = Q_{\text{eats}} \cdot K_{\text{The}} = [1, 1] \cdot [0, 1] = 0 + 1 = 1$$

2. Attention score between "eats" and "cat":

$$\text{Attention score} = Q_{\text{eats}} \cdot K_{\text{cat}} = [1, 1] \cdot [1, 0] = 1 + 0 = 1$$

3. Attention score between "eats" and "eats":

$$\text{Attention score} = Q_{\text{eats}} \cdot K_{\text{eats}} = [1, 1] \cdot [1, 1] = 1 + 1 = 2$$

4. Attention score between "eats" and "fish":

$$\text{Attention score} = Q_{\text{eats}} \cdot K_{\text{fish}} = [1, 1] \cdot [0, 0] = 0 + 0 = 0$$

## Step 5: Apply Softmax to Get Attention Weights

The attention weights are calculated by applying the softmax function to the attention scores:

$$\text{Softmax}(1, 1, 2, 0)$$

1. Exponentiate each score:

- $e^1 = 2.718, e^1 = 2.718, e^2 = 7.389, e^0 = 1$

2. Sum the exponentiated scores:

$$2.718 + 2.718 + 7.389 + 1 = 13.825$$

3. Calculate the attention weights:

- For "The":  $\frac{2.718}{13.825} \approx 0.196$
- For "cat":  $\frac{2.718}{13.825} \approx 0.196$

.....

- We computed the attention weights with respect to "eats" as follows:
  - Attention weight for "The": 0.196
  - Attention weight for "cat": 0.196
  - Attention weight for "eats": 0.534
  - Attention weight for "fish": 0.072

These weights represent how much the word "eats" pays attention to each of the other words, including itself.

## Step 6: Calculate the Weighted Sum of Value Vectors

Recap of Value Vectors:

- Value vector for "The": [1, 3]
- Value vector for "cat": [2, 4]
- Value vector for "eats": [3, 7]
- Value vector for "fish": [0, 0]

Contextual Vector Calculation for "eats":

The contextual vector is the weighted sum of all the value vectors using the attention weights. This is calculated as follows:

$$\text{Contextual Vector} = (0.196 \cdot [1, 3]) + (0.196 \cdot [2, 4]) + (0.534 \cdot [3, 7]) + (0.072 \cdot [0, 0])$$

## Step 6: Calculate the Weighted Sum of Value Vectors

Recap of Value Vectors:

- Value vector for "The": [1, 3]
- Value vector for "cat": [2, 4]
- Value vector for "eats": [3, 7]
- Value vector for "fish": [0, 0]

Contextual Vector Calculation for "eats":

The contextual vector is the weighted sum of all the value vectors using the attention weights. This is calculated as follows:

$$\text{Contextual Vector} = (0.196 \cdot [1, 3]) + (0.196 \cdot [2, 4]) + (0.534 \cdot [3, 7]) + (0.072 \cdot [0, 0])$$

Step-by-Step Calculation:

1. Multiply each value vector by its corresponding attention weight:

- For "The":  $0.196 \times [1, 3] = [0.196, 0.588]$
- For "cat":  $0.196 \times [2, 4] = [0.392, 0.784]$
- For "eats":  $0.534 \times [3, 7] = [1.602, 3.738]$
- For "fish":  $0.072 \times [0, 0] = [0, 0]$

2. Add these weighted value vectors together:

$$\text{Contextual Vector} = [0.196, 0.588] + [0.392, 0.784] + [1.602, 3.738] + [0, 0]$$

3. Perform the addition:

$$\begin{aligned}\text{Contextual Vector} &= [0.196 + 0.392 + 1.602 + 0, 0.588 + 0.784 + 3.738 + 0] \\ &= [2.19, 5.11]\end{aligned}$$

## What Does This Mean?

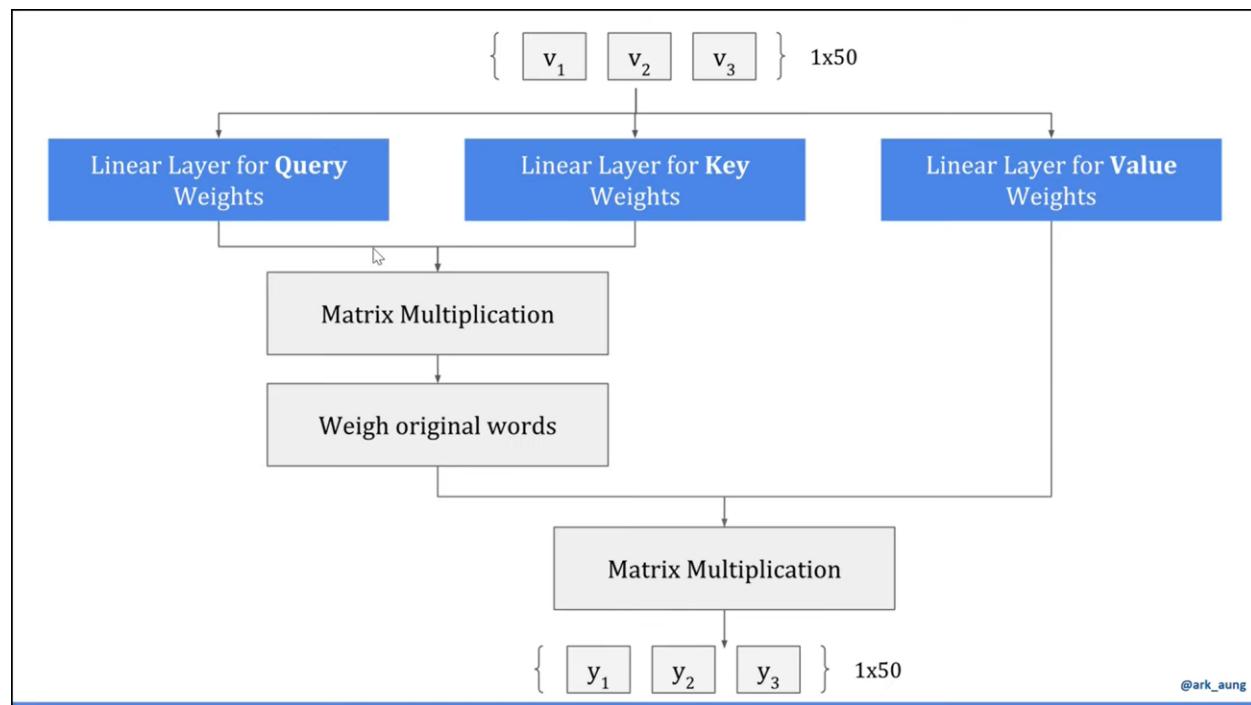
The **contextual vector** for the word "eats" now incorporates information from all the words in the sentence based on how much attention "eats" paid to each of them.

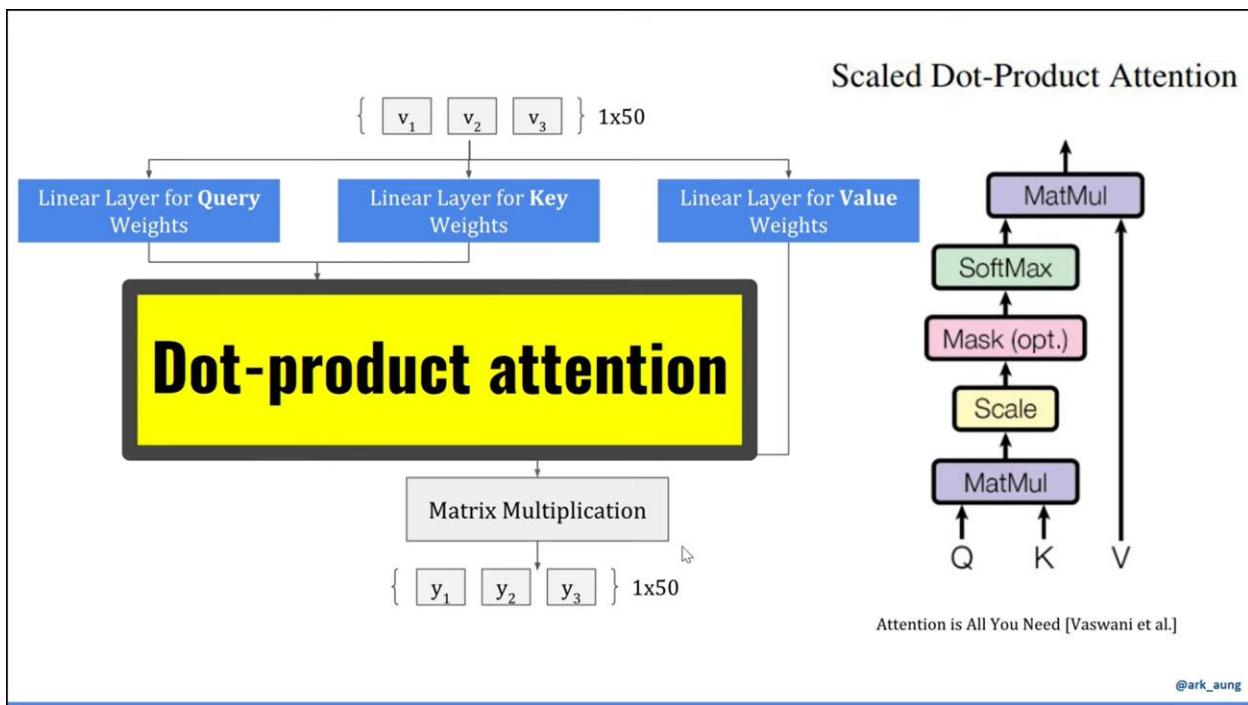
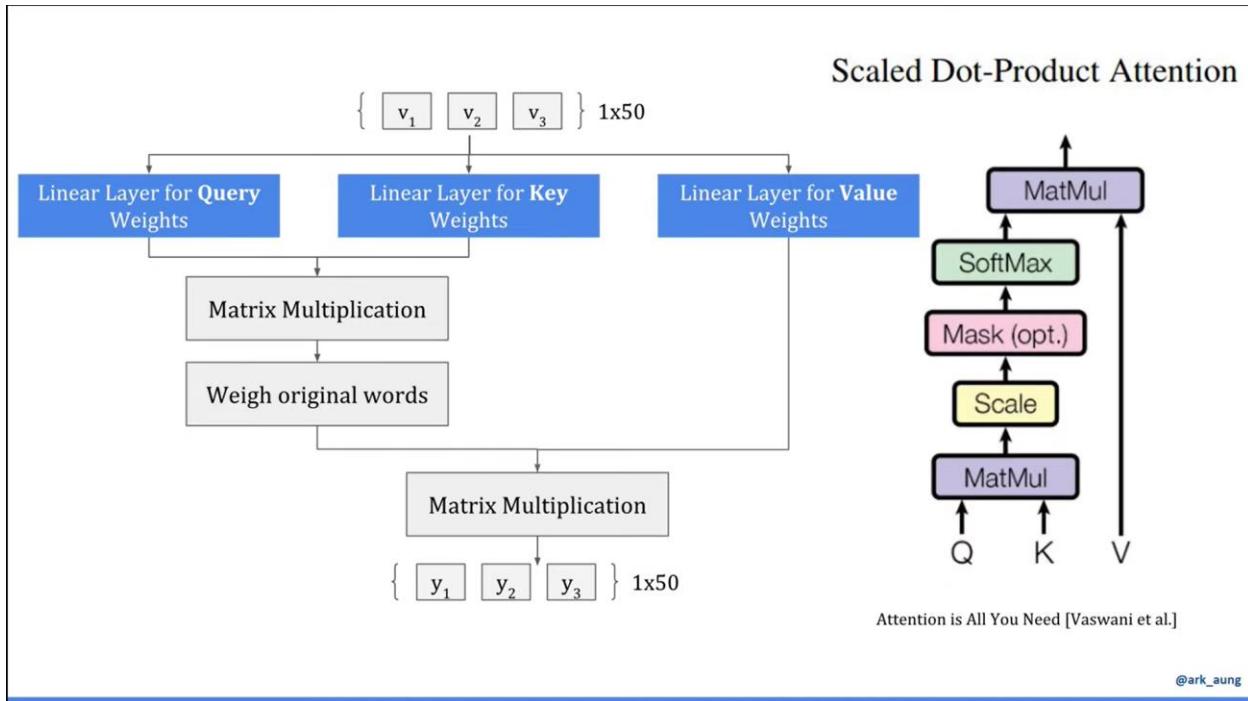
This new vector is a **context-aware representation** of the word "eats" that reflects its relationship with the words "The," "cat," and "fish."

### Recap of the Process:

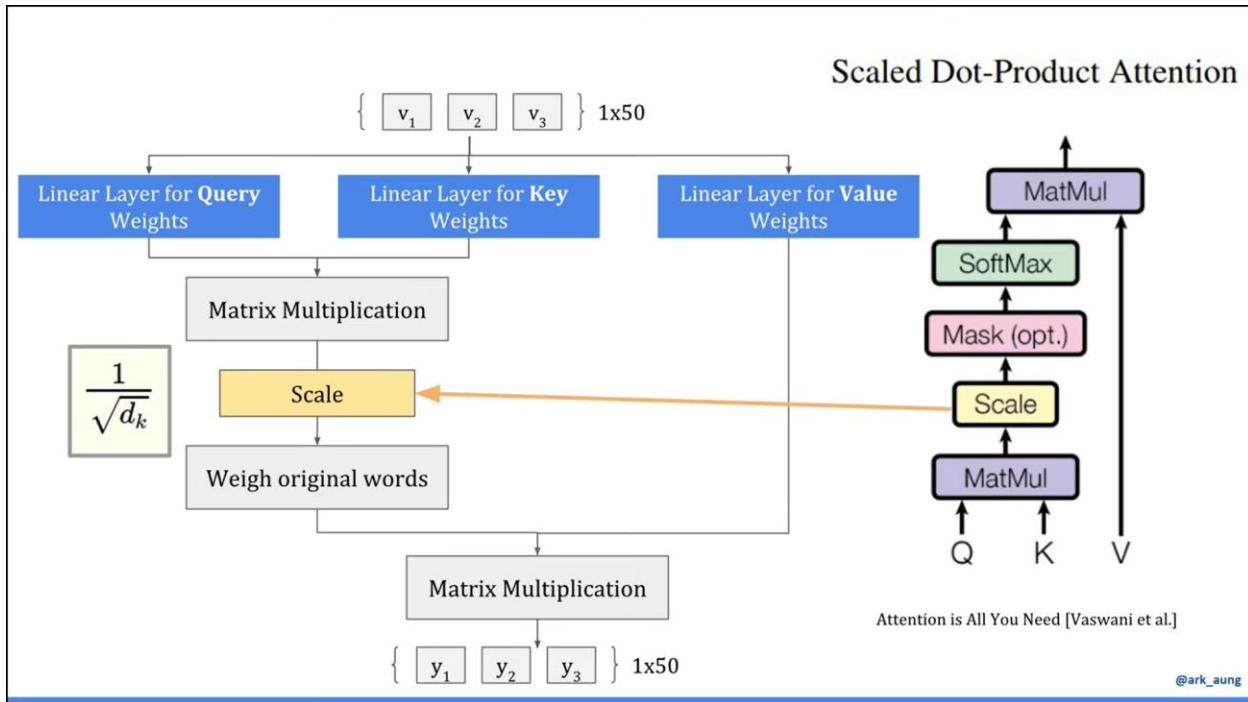
1. Calculate Query, Key, and Value vectors for each word.
2. Compute attention scores (dot product of Query and Key).
3. Apply softmax to obtain attention weights.
4. Use the attention weights to create a weighted sum of the Value vectors to produce the final contextual vector.

This toy example demonstrates how **self-attention** works in a simple context. By following this process for each word in the sentence, the model can create **context-aware representations** for all words, which improves its understanding of the relationships between them.





Here is the scaled version:



### Intuition behind the scaling:

Assume we have query vector  $v_1 = [1 \ 1 \ 1]$

Key vectors:  $v_1 = [1 \ 1 \ 1]$ ,  $v_2 = [1 \ 0 \ 0]$ ,  $v_3 = [0 \ 1 \ 0]$ ,  $v_4 = [0 \ 0 \ 1]$

$$S_{11} = v_1 v_1.T, S_{12} = v_1 v_2.T, S_{13} = v_1 v_3.T, S_{14} = v_1 v_4.T$$

$$S_1 = [S_{11} \ S_{12} \ S_{13} \ S_{14}] = [3 \ 1 \ 1 \ 1]$$

Without scaling:

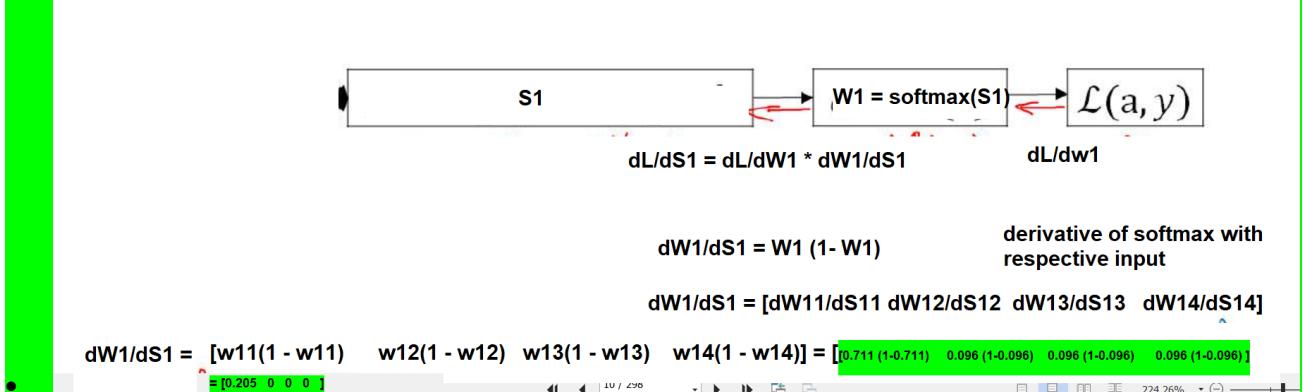
$$W_1 = \text{softmax}(S_1) = [0.711 \ 0.096 \ 0.096 \ 0.096]$$

- Without scaling, the softmax outputs are skewed more towards the largest value, giving 71.1% attention to the first token/vector.

During gradient computation, we have:

$$\begin{aligned} dW_1/dS_1 &= W_1 (1-w_1) = [W_{11} (1-w_{11}) \quad W_{12} (1-w_{12}) \quad W_{13} (1-w_{13}) \quad W_{14} (1-w_{14})] = \\ &= [0.711 (1-0.711) \quad 0.096 (1-0.096) \quad 0.096 (1-0.096) \quad 0.096 (1-0.096)] \\ &= [0.711 (1-0.711) \quad 0.096 (1-0.096) \quad 0.096 (1-0.096) \quad 0.096 (1-0.096)] \\ &= [0.205 \ 0.087 \ 0.087 \ 0.087] \\ &= [0.205 \ 0 \ 0 \ 0] \end{aligned}$$

- In this case, the gradient for the last three token/vectors is negligible, and only first tokens contribute to the model's learning, leading to lesser meaningful updates, shown below:



With Scaling (scaling factor as dimension of vector is 4, so  $\sqrt{4} = 2$ ):

$$S1 = [S11 \ S12 \ S13 \ S14] = [3 \ 1 \ 1 \ 1]$$

$$S1' = S1 / \sqrt{4} = [3 \ 1 \ 1 \ 1] / 2 = [1.5 \ 0.5 \ 0.5 \ 0.5]$$

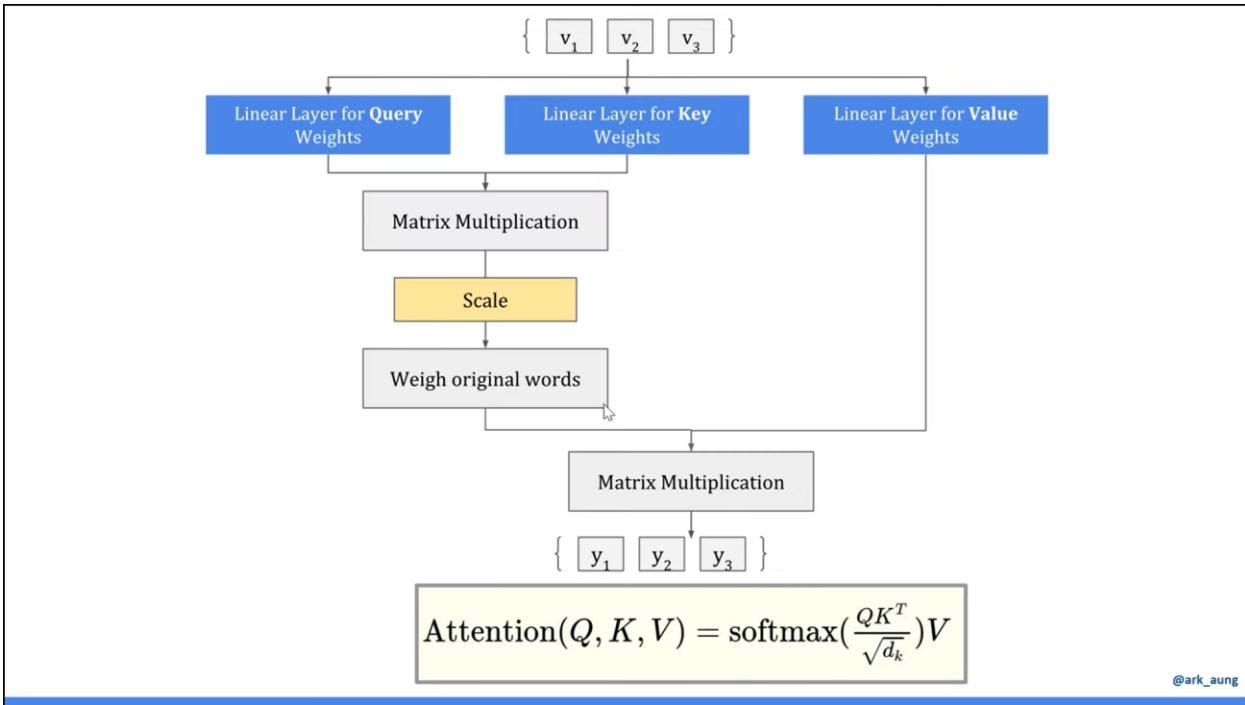
$$W1 = \text{softmax}(S1) = [0.47 \ 0.17 \ 0.17 \ 0.17]$$

With scaling, the dot product might result in a more balanced softmax output like

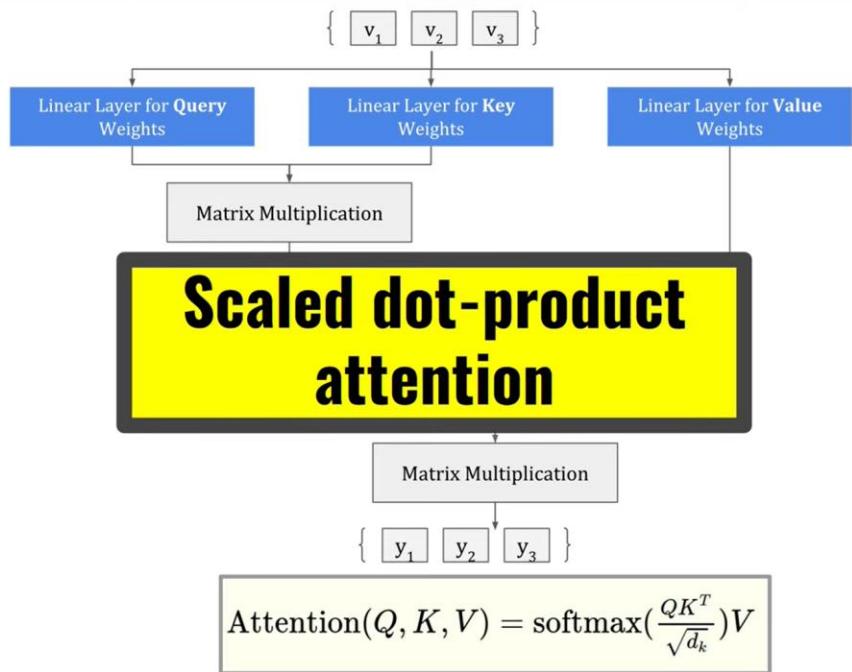
[0.47 0.17 0.17 0.17]. In this case, the gradient for the last three token is no longer negligible, and all tokens contribute to the model's learning, leading to more meaningful updates.

#### Effect of Scaling:

- Without scaling, the softmax outputs are skewed more towards the largest value, giving 71.1% attention to the first key.
- With scaling, the attention is more balanced, distributing as 47%, 17%, and 17%, and 17%, which prevents overly sharp peaks in attention distribution



This is similar to the diagram in paper



@ark\_aung

## Lecture # 03

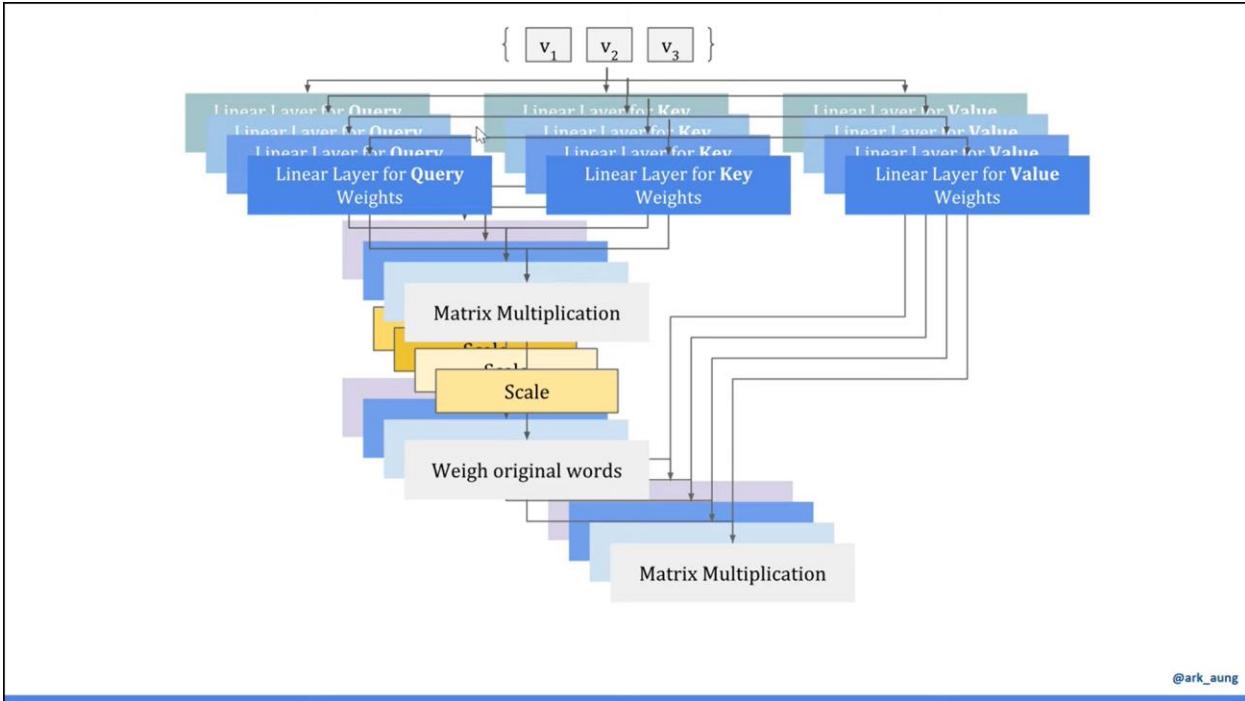
Two new idea in transformer compare to other NN: positional encoding + multi head attention

Question, can we use the same architecture of transformer

1. Translation
  2. Next word generation.
  3. How the chat gpt work.... Is it work like next word generation
  4. Is the translation just like the generation of a next work???
- 
1. Multihead attention
  2. WordtoVec encoding
  3. Positional encoding,
  4. Encoder layers
  5. Deconder layers
  6. Cross attention:
  7. Mask multihead layers etc.

1. Multihead attention

When you are looking at a word in a sentence, you are **giving attention to more than one thing**



## 1. Analogous to Convolutional Neural Networks (CNNs)

**Multiple Feature Maps:** The concept is similar to using multiple kernels in CNNs, where each kernel learns different features from the input. In transformers, multiple heads allow for learning various "feature maps" of word relationships, enhancing the model's ability to understand context

## 2. Reasons for Using Multiple Attention Heads:

It expands the model's ability to focus on different positions

**A. Diverse Relationships:** Each attention head can focus on different aspects of the input, enabling the model to learn various relationships among words. This is particularly important in natural language processing, where words can relate in complex ways (e.g., negation, moderation)

**B. Parallel Processing:** By running multiple attention heads in parallel, the transformer can process information more efficiently and capture richer interpretations of the input data

## E. Capture Different Aspects of Relationships

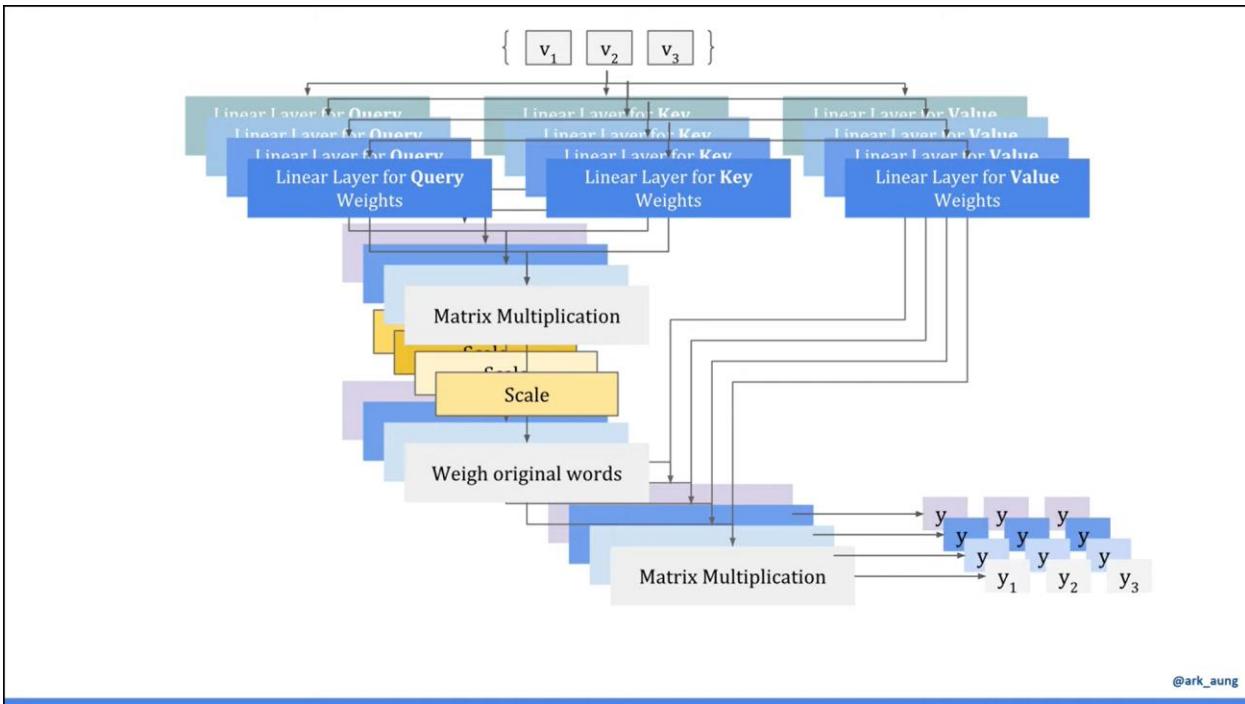
- Each attention head computes **self-attention** independently.

- By doing this, different heads can **focus on different parts of the input sequence** or capture different types of relationships between words (or tokens). This means that:
  - One head might focus on **short-range dependencies** (words close to each other).
  - Another head might capture **long-range dependencies** (words far apart).
  - A different head might focus on **syntactic information**, while another might focus on **semantic information**.

## Summary:

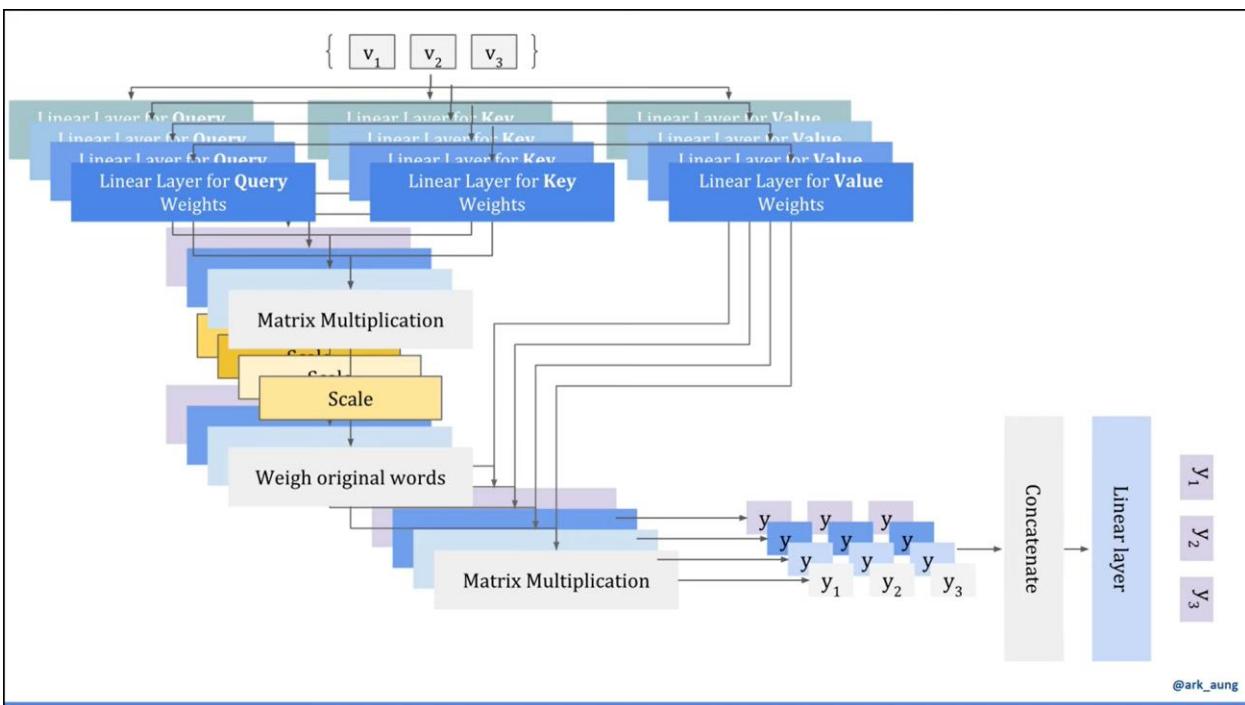
- The **training data and task-specific requirements** guide each head to learn different **attention patterns** based on the feedback from the loss function.
- Different heads in **multi-head attention** learn to focus on **short-range vs. long-range dependencies**, as well as **syntactic vs. semantic relationships**, because each head has **independent weight matrices** that are learned through training.
- Some heads specialize in **grammatical structure**, while others focus on **meaning**.
- This diversity of attention patterns makes the model more powerful and capable of capturing complex relationships in the input sequence.

By allowing different **attention heads to specialize in capturing different types of relationships**, the Transformer model can better understand the complexities of natural language, improving performance on tasks like translation, summarization, and classification.

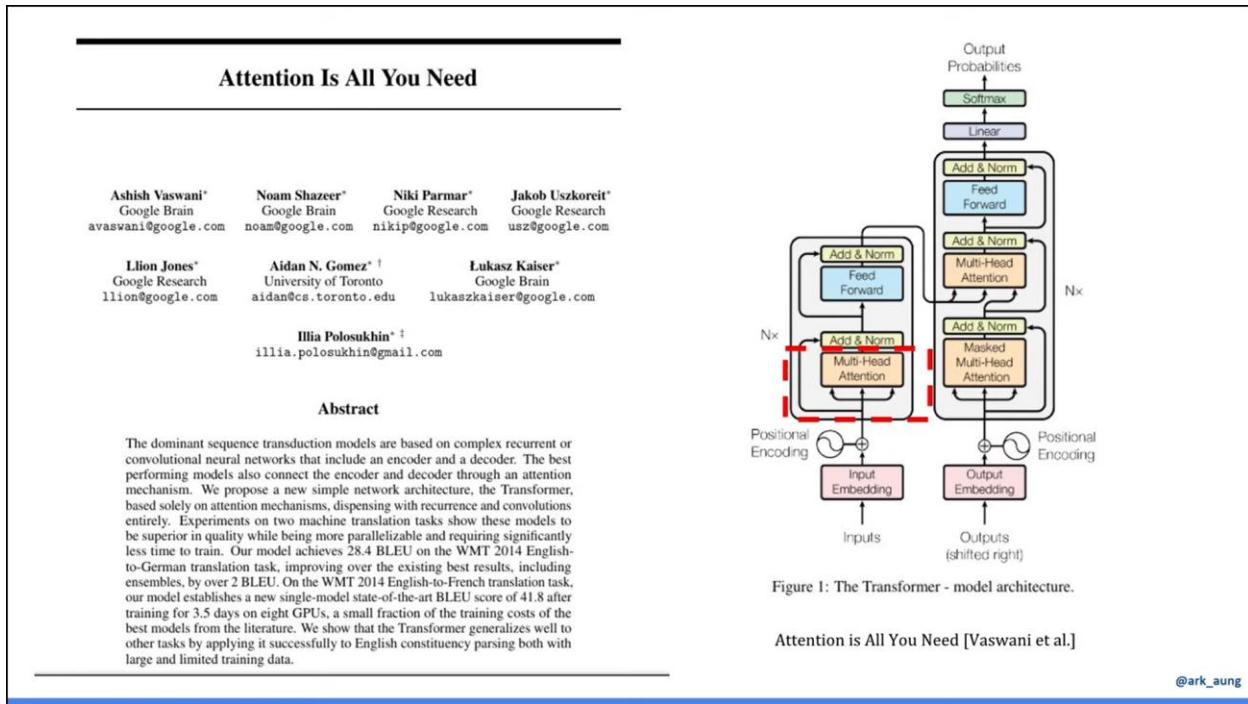
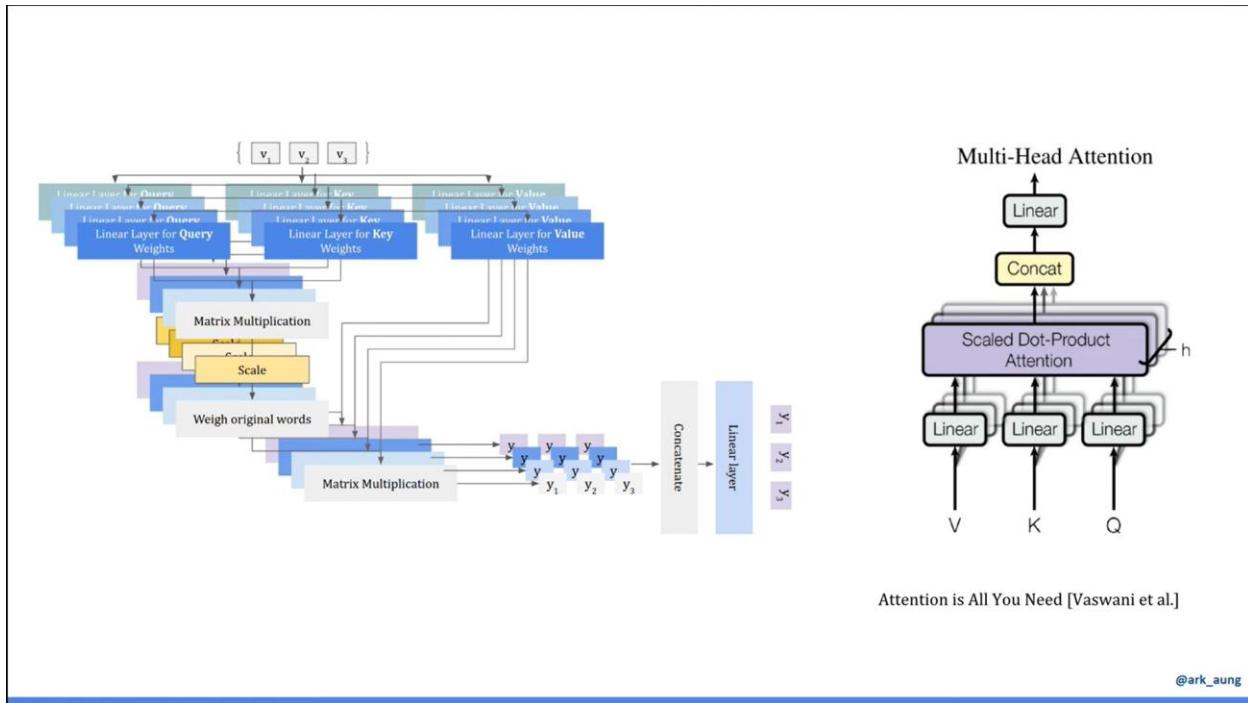


@ark\_aung

We concatenate ally contextual vectors to generate the single vectors



@ark\_aung

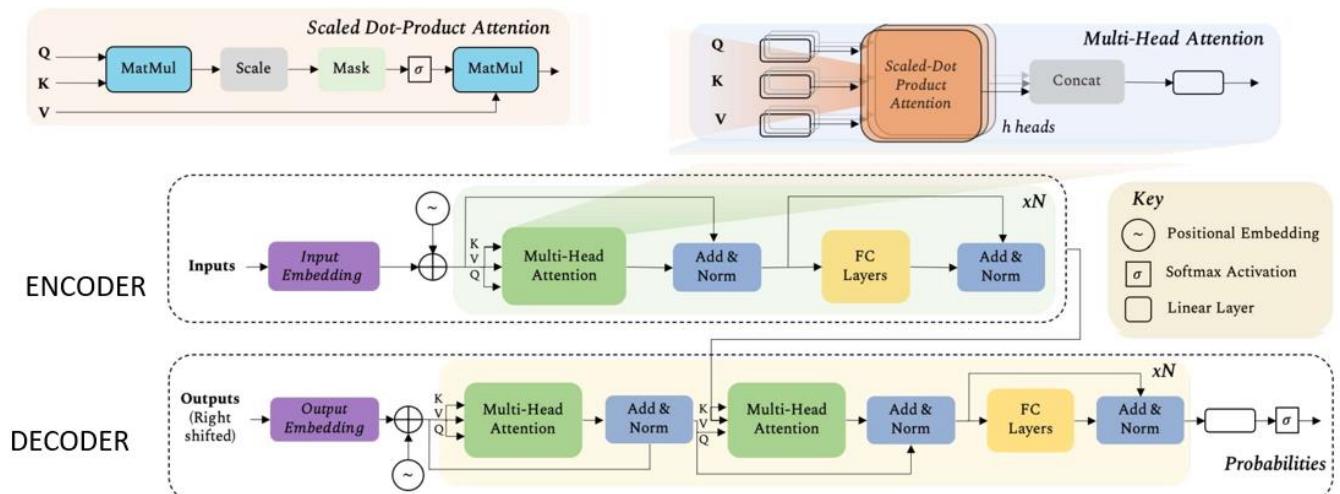


# Transformer

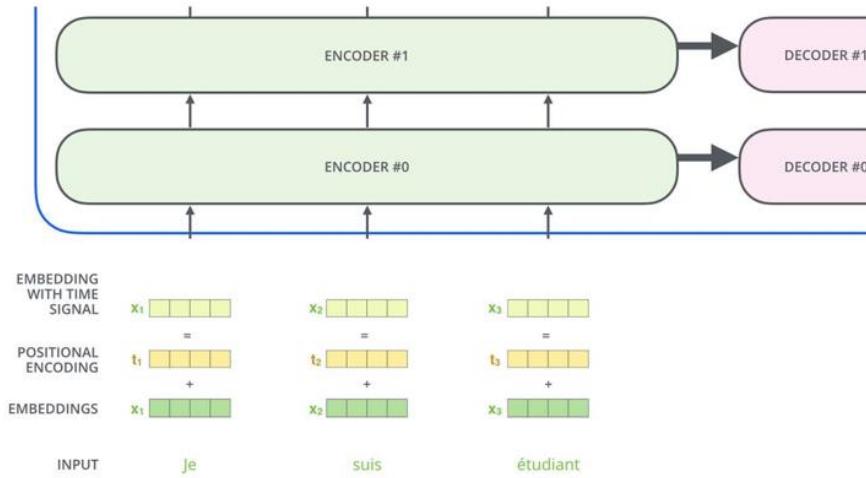
- It consists of Encoder and Decoder Blocks
- Main components of each block:
  - Self-Attention
  - Layer Normalization
  - Feed Forward Network

Transformer uses self attention mechanism to capture the weighted sum of all the tokens in a sequence at each layer.

# Transformer



## 2. Word and positional embeddings



### a. Word Embeddings

Word encoding in transformers refers to the process of converting textual data into numerical representations that can be processed by the model. This involves several key components:

- Definition: Word embeddings are dense vector representations of words that capture semantic meanings. In transformers, each word or token is mapped to a vector in a high-dimensional space, where similar meanings are represented by vectors that are close together.

### b. Positional Encoding:

- Since the Transformer model has no inherent sense of the order of words due to its lack of recurrence or convolution (RNN/LSTM), it requires a way to incorporate positional information into the input embeddings. Positional encodings are added to the input embeddings to allow the model to understand the order of words in a sequence.
- SINCE transformer process all token simultaneously, so it need some additional mechanism to encode the postional information
- Purpose: Since transformers do not utilize recurrence or convolution, they need a way to incorporate the order of words in a sequence. Positional encoding adds information about the position of each word in the input sequence, ensuring that the model understands the sequence's structure.
- Mechanism: Positional encodings are generated using sine and cosine functions, creating unique vectors for each position in the input sequence. These positional vectors are added to the word embeddings, allowing the model to retain both meaning and order

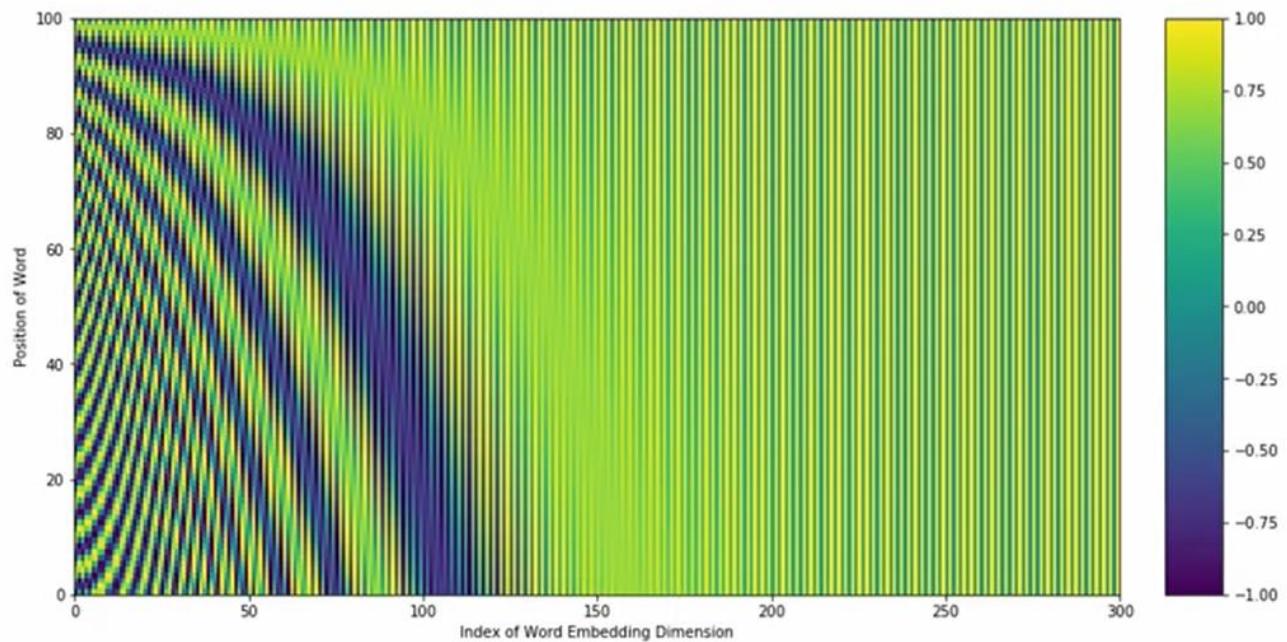
$$P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$

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

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'

- **n is a constant set to 100, which determines the wavelength scaling**



The sine and cosine functions with varying frequencies or wavelengths create distinctive patterns. The colors indicate values ranging from -1 to 1.

- The X-axis represents the length of the word embedding (512 dimensions, as used by the Transformer model).
  - The Y-axis represents the position of a token within a sequence.

If a word is at the 20th position, the corresponding horizontal line showcases a unique combination of sine and cosine values for that position. This positional encoding is a distinctive set of values, ensuring that no other word shares the same encoding.

Positional encoding is just like the binary representation of decimal number, which is unique for all numbers.

0:	0	0	0	0	7:	1	0	0	0
1:	0	0	0	1	8:	1	0	0	1
2:	0	0	1	0	9:	1	0	1	0
3:	0	0	1	1	10:	1	0	1	1
4:	0	1	0	0	11:	1	1	0	0
5:	0	1	0	1	12:	1	1	0	1
6:	0	1	1	0	13:	1	1	1	0
7:	0	1	1	1	14:	1	1	1	1

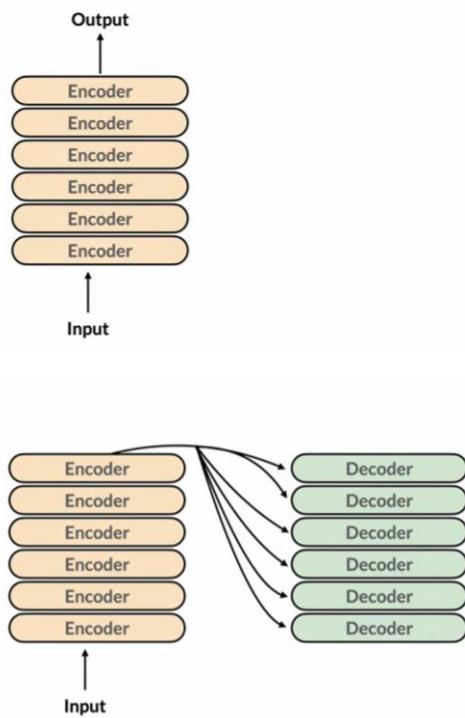
Another way of use binary reparsent. ...but sin and cosine are better

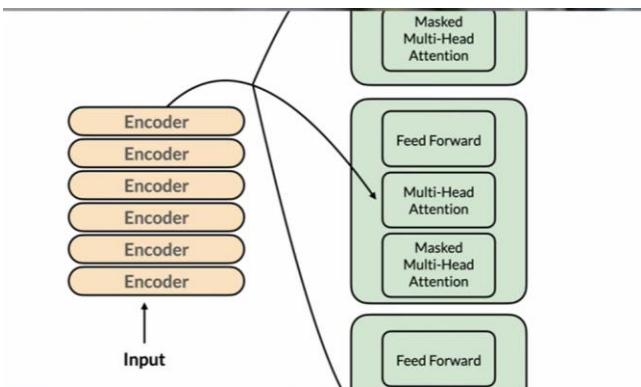
## Position Encoding

- Can also be learned
- Learn like other parameters

### Summey: Encoding Process

- The encoding process typically involves:
  - **Input Preparation:** Text is tokenized into words or subwords.
  - **Embedding Lookup:** Each token is converted into its corresponding embedding vector.
  - **Adding Positional Encoding:** The positional encoding vector is added to each word embedding to form a final input representation for the transformer



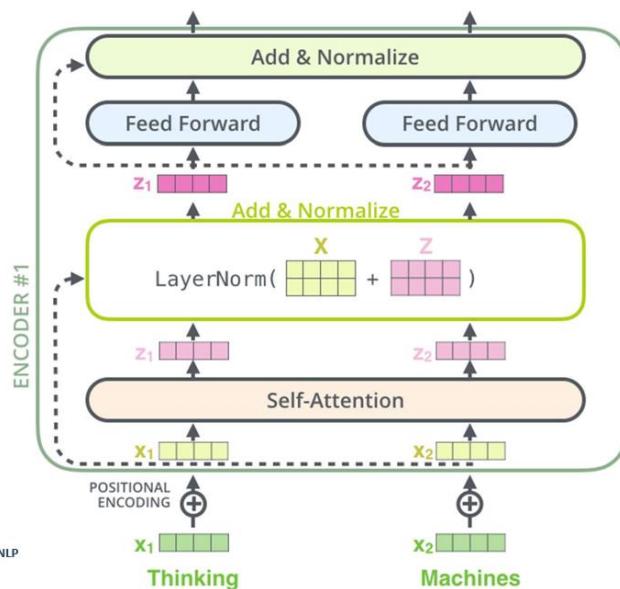


### 3. Add and normalize



Attention and Transformers

#### Add and Normalize



**Add: Skip connection just like resnet**

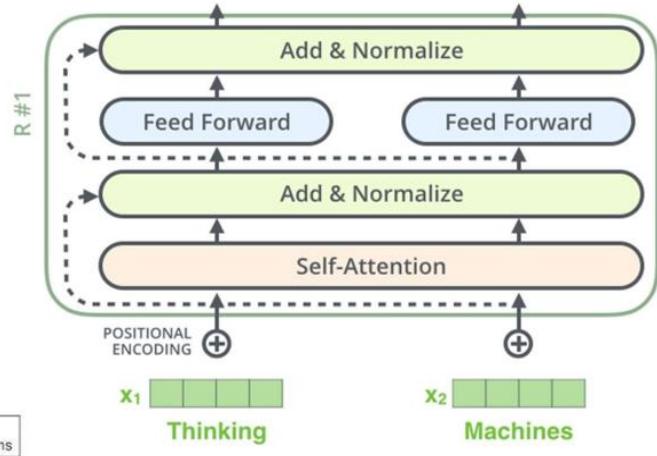
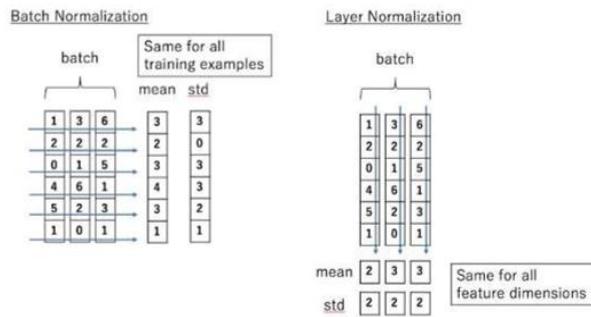
## Normalize: layer normalization just like batch normalization



## Attention and Transformers

### Layer Normalization (Hinton)

Layer normalization normalizes the inputs across the features.



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

$$\text{Normalized Activations} = \frac{\text{Activations} - \text{Mean}}{\text{Standard Deviation}}$$





## Covariant shift problem

### Batch Normalization

- It occurs in mini batches  $\mathcal{D}$
- It is due to the difference in distribution of data among the batches,
- It makes difficult for optimizers that learn a decision boundary during training of batch #1 work well to make prediction for batch #2 due to covariant shift.
- So ultimately it slows down the training process.

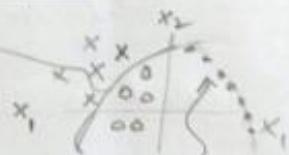


Batch I  
= 100 images

Rose	-	-	-
Not Rose	-	-	-

(close roses)

optimizer for  
Batch I learns  
this decision  
boundary



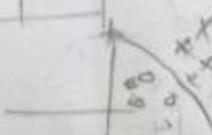
although Good Train function is same.  
(Rose v/s Not Rose)

Batch II  
= 100

Rose	-	-	-
Not Rose	-	-	-

(open roses)

optimizer for batch I is unable  
to make correct prediction for  
batch II → to learn Batch II DR & incorporate with Batch I



Normalization: helps to overcome the covariance shift problem.

- covariate shift in neural networks refers to a change in the distribution of input data, which can lead to a decrease in model performance.
- Normalization refers to a change in the input data distribution that can occur when the input features have different statistics (e.g., mean and variance/close rose vs open rose) in different parts of the training data. Because the input distributions are more stable, the model learns more quickly. This often allows for the use of higher learning rates, which can significantly reduce training time.

The slide explains *Layer Normalization*, a technique introduced by Geoffrey Hinton, and provides a visual comparison with *Batch Normalization*. Here's a breakdown of the concepts presented:

### 1. Layer Normalization vs. Batch Normalization:

- **Batch Normalization:** This technique normalizes each input feature by computing the mean and standard deviation (std) across the batch. In essence, normalization occurs over the entire batch for each feature separately. This means that the mean and std are computed across all training examples for each feature dimension, **leading to a consistent scale within each feature across the batch**.
- **Layer Normalization:** Layer normalization normalizes the inputs across the feature dimension for each individual sample, making it independent of batch size
- In contrast, layer normalization normalizes across all feature dimensions for a single example, rather than across a batch. This means that the mean and std are calculated independently for each input example, not across the batch. **This technique is beneficial in scenarios where batch statistics may not be stable or applicable, such as when working with sequential data like in Natural Language Processing (NLP) or with Transformers.**

### Key Differences Highlighted in the Slide:

- **Batch Normalization (Left Side Diagram):** Mean and standard deviation are computed column-wise (across all samples in a batch). The normalization is done independently for each feature.
- **Layer Normalization (Middle Diagram):** Mean and standard deviation are computed across all features of a given sample. This normalization is done per sample, which is why layer normalization is better suited for tasks like text or sequential data processing.

### 2. Transformer Architecture:

The slide provides a simplified view of a Transformer encoder, with the following key components:

- **Self-Attention:** The self-attention mechanism processes each input word or token (e.g., x1x\_1x1 and x2x\_2x2 in this example) by considering the relationships between all words/tokens in the input sequence. It captures dependencies regardless of their positions within the sequence.
- **Add & Normalize Layers:** This is where the Layer Normalization occurs. In the diagram, you can see two "Add & Normalize" boxes. In the first, normalization is applied after self-attention. In the

second, it's applied after the feed-forward network. This normalization ensures that the data is scaled appropriately at different stages of the network.

- **Feed-Forward Networks:** These networks apply non-linear transformations to the output of the self-attention mechanism. Each feed-forward block is followed by an "Add & Normalize" operation.
- **Positional Encoding:** Since Transformer models have no inherent understanding of word order, positional encodings are added to input embeddings to provide information about the positions of tokens in the sequence (labeled as 'Thinking' and 'Machines' in the diagram).

### Conclusion:

The slide emphasizes how Layer Normalization is applied across features and how it fits into the broader context of Transformer architectures. It visually illustrates the use of normalization in handling dependencies between elements in sequential data (like in NLP) and the importance of normalization in stabilizing and improving deep learning models.

If you have specific questions about any part of the slide or want more details on any section, let me know!

Example:



## 1. Batch Normalization (BN):

Batch Normalization works by normalizing the activations of a layer across the entire batch. Let's consider a simple neural network example:

### Example:

Imagine a fully connected neural network layer with 3 input features and a batch size of 4. Let's assume that the activations for a particular layer look like this:

$$\text{Activations} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix} \quad \begin{array}{l} \text{batch size} = 4 \\ \text{input features} = 3 \end{array}$$

Here, each row represents an example, and each column represents a feature. We want to apply batch normalization to these activations.

### Steps for Batch Normalization:

1. Compute the Mean and Standard Deviation for each feature across the batch:

$$\text{Mean} = \left[ \frac{1+4+7+10}{4}, \frac{2+5+8+11}{4}, \frac{3+6+9+12}{4} \right] = [5.5, 6.5, 7.5]$$

$$\text{Standard Deviation} = [\text{std}(1, 4, 7, 10), \text{std}(2, 5, 8, 11), \text{std}(3, 6, 9, 12)] \approx [3.5, 3.5, 3.5]$$

2. Normalize each feature in the batch using the mean and standard deviation:

$$\text{Normalized Activations} = \frac{\text{Activations} - \text{Mean}}{\text{Standard Deviation}}$$

Applying this to our activations:

$$\begin{bmatrix} \frac{1-5.5}{3.5}, & \frac{2-6.5}{3.5}, & \frac{3-7.5}{3.5} \\ \frac{4-5.5}{3.5}, & \frac{5-6.5}{3.5}, & \frac{6-7.5}{3.5} \\ \frac{7-5.5}{3.5}, & \frac{8-6.5}{3.5}, & \frac{9-7.5}{3.5} \\ \frac{10-5.5}{3.5}, & \frac{11-6.5}{3.5}, & \frac{12-7.5}{3.5} \end{bmatrix} = \begin{bmatrix} -1.29, & -1.29, & -1.29 \\ -0.43, & -0.43, & -0.43 \\ 0.43, & 0.43, & 0.43 \\ 1.29, & 1.29, & 1.29 \end{bmatrix}$$

3. Apply learnable parameters (scale and shift): After normalizing, we typically apply scaling and shifting parameters ( $\gamma$  and  $\beta$ ), which are learnable during training.



## 2. Layer Normalization (LN):

Layer Normalization normalizes the activations across all features for each individual example in a batch.

### Example:

Let's consider the same activation matrix as before:

$$\text{Activations} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix}$$

### Steps for Layer Normalization:

1. Compute the Mean and Standard Deviation for each example (row-wise calculation):

- For the first example: [1, 2, 3]:

$$\text{Mean} = \frac{1+2+3}{3} = 2, \quad \text{Standard Deviation} = \sqrt{\frac{(1-2)^2 + (2-2)^2 + (3-2)^2}{3}} = 1$$

- For the second example: [4, 5, 6]:

$$\text{Mean} = \frac{4+5+6}{3} = 5, \quad \text{Standard Deviation} = 1$$

Similarly, we find the mean and standard deviation for the remaining examples.

2. Normalize each feature for each example using the example's mean and standard deviation:

- For the first example:

$$\text{Normalized Activations} = \left[ \frac{1-2}{1}, \frac{2-2}{1}, \frac{3-2}{1} \right] = [-1, 0, 1]$$

- For the second example:

$$\text{Normalized Activations} = \left[ \frac{4-5}{1}, \frac{5-5}{1}, \frac{6-5}{1} \right] = [-1, 0, 1]$$

Applying these to all examples:

$$\text{Normalized Activations} = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$$

3. Apply learnable parameters (scale and shift): As with batch normalization, we typically apply scaling ( $\gamma$ ) and shifting ( $\beta$ ) parameters here too.

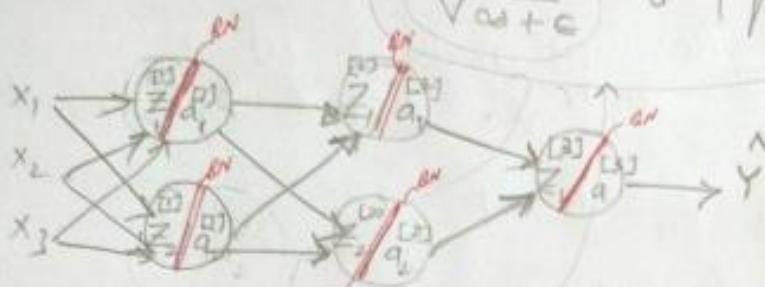
## Summary of Key Differences:

- Batch Normalization: Normalizes across a batch for each feature. Primarily used in CNNs and dense layers.
- Layer Normalization: Normalizes across features for each example individually. Typically used in RNNs, LSTMs, and Transformers.



## Adding Batch Norm to a NN

$$\tilde{Z} = \frac{Z - \mu}{\sqrt{\sigma^2 + \epsilon}} \times \gamma + \beta$$



$\cdot X \xrightarrow{w, b} Z \xrightarrow{\mu, \sigma^2} \text{Batch (BN) Norm} \xrightarrow{\beta, \gamma} Z \xrightarrow{g(Z)} a \xrightarrow{w, b} Z \xrightarrow{\mu, \sigma^2} \text{BN} \xrightarrow{\beta, \gamma} Z \xrightarrow{g(Z)} a \xrightarrow{w, b} \dots$

- Compute  $\mu$  &  $\sigma^2$
- Each example is subtracted from mean & scale by variance
- rescale by  $\gamma$  &  $\beta$

• Parameters:  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$

$\beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]}$

(additional parameters)

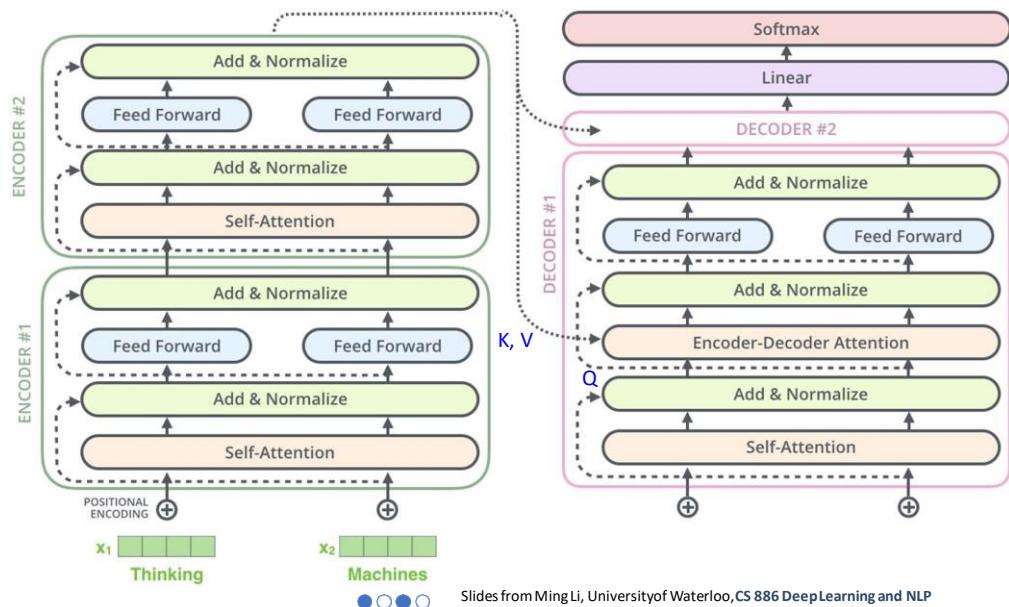
( $\beta$  initializations are different from BN)

- compute  $d\mu^{[l]}$  &  $d\gamma^{[l]}$  during back propagation

$$\beta^{[l+1]} = \beta^{[l]} - \alpha d\beta^{[l]}$$

- will will use PyTorch framework that do it for us (`torch.nn.BatchNorm`)

# The complete transformer



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

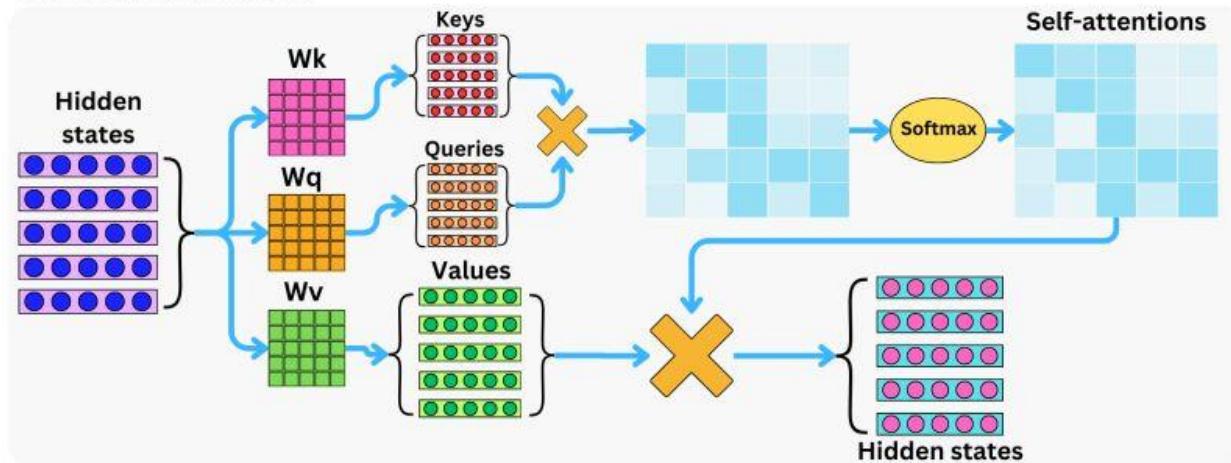


Lecture 04:

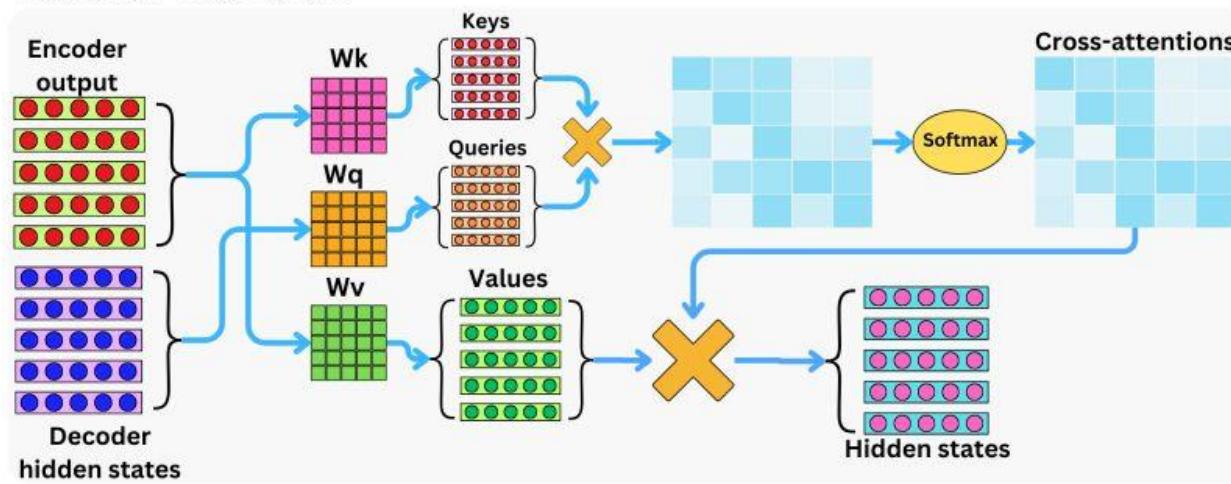
# Self-Attention VS Cross-Attention

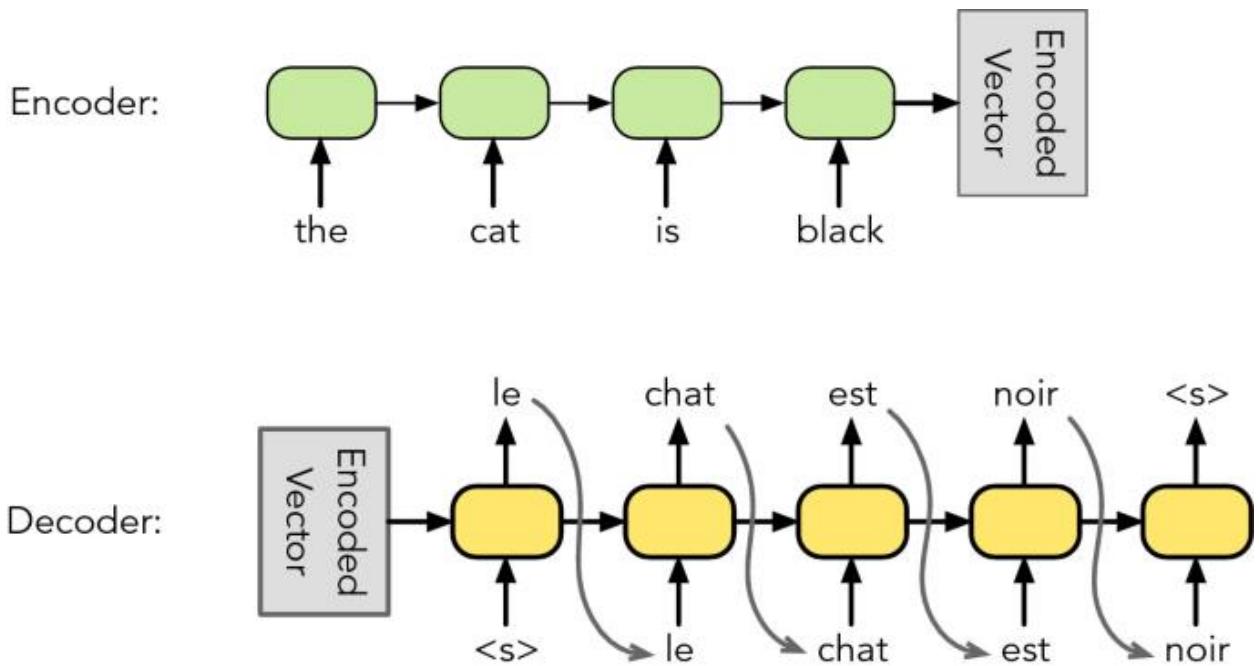
[TheAiEdge.io](https://TheAiEdge.io)

## The Self-Attentions



## The Cross-Attentions





### Toy Example

Let's consider a simple example where we translate "The cat sits on the mat" from English to French.

#### Step 1: Encoding

- The encoder processes the input sentence and generates hidden states for each word:
  - "The" → \$h\_1\$
  - "cat" → \$h\_2\$
  - "sits" → \$h\_3\$
  - "on" → \$h\_4\$
  - "the" → \$h\_5\$
  - "mat" → \$h\_6\$

**Step 2: Decoding with Cross Attention** [During cross-attention, this query vector interacts with all key vectors (K\_1 to K\_7). The key vector for "deposited" (K\_2) receives a high attention weight because it closely matches the decoder's query vector, allowing the decoder to generate "déposé"]

- Suppose we are generating the first French word "Le" (the equivalent of "The"). The decoder creates a query vector based on its current state.
- As the decoder processes each token, it uses cross attention to focus on relevant parts of the encoder's output. For generating "Le," the decoder creates a query vector based on its current state (which starts with <SOS>). [This is often represented by a special token, such as <SOS> (start of sentence). This token is fed into the decoder to initiate the generation process]
- It computes **attention scores** using this query with all encoder hidden states:
  - Score for h\_1 (The): high
  - Score for h\_2 (cat): medium
  - Score for h\_3 (sits): low
  - Score for h\_4 (on): low
  - Score for h\_5 (the): low
  - Score for h\_6 (mat): low

### Step 3: Weighted Sum

- After applying softmax to these scores, we get attention weights:
  - High weight on "The", lower weights on others.
- The context vector is computed as:

$$\text{Context} = \text{weight1} \cdot h_1 + \text{weight2} \cdot h_2 + \dots + \text{weight6} \cdot h_6$$

### Step 4: Generating Output

- The decoder uses this context vector to generate "Le". This process continues for subsequent words like "chat" (cat), where it will again use cross attention to focus on relevant parts of the encoder's output based on its current query.

- Iterative Process: The generated word "Le" is then fed back into the decoder as part of its input for generating subsequent words in the translation (e.g., "chat" for cat). This process continues until an end-of-sentence token is produced.

#### 4. Cross-Attention Mechanism (Encoder-Decoder Attention) in Detail

This cross-attention layer is what enables the model to “look back” at the input sequence while generating the output sequence, ensuring the output words are contextually aligned with the input.

##### Summary of Key Points:

- **Query from Decoder:** Represents what the decoder is currently focusing on.
- **Keys and Values from Encoder:** Capture the processed information from the input sequence.
- **Cross-Attention Mechanism:** Matches the query from the decoder with the keys from the encoder to determine which input elements to focus on.
- **Context Vector:** A weighted combination of the values based on attention scores, guiding the decoder in generating output.

Let's go through a simple translation example using the Transformer model and see how the **cross-attention mechanism** (encoder-decoder attention) works step-by-step.

##### Example:

- Imagine translating the
- English sentence "**She ate an apple**" into French as "**Elle a mangé une pomme**".

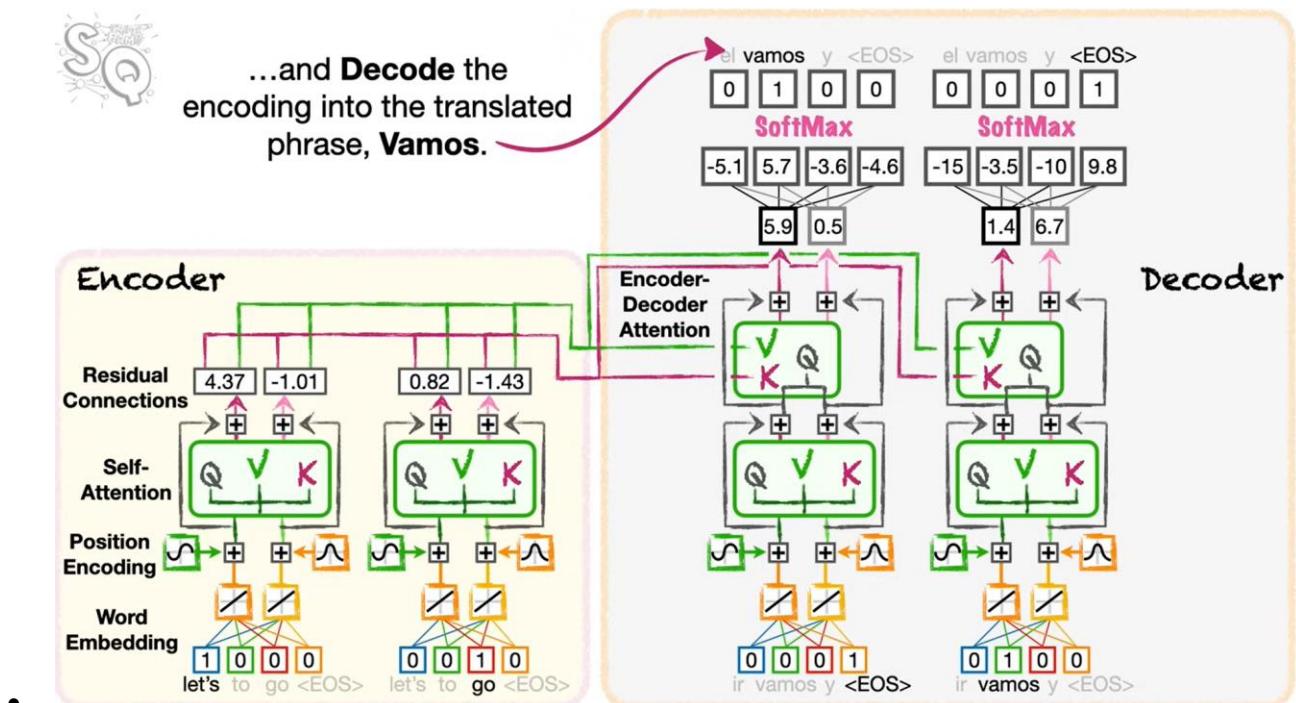
- The key goal of the cross-attention mechanism is to help the decoder pay attention to the most relevant parts of the encoded input sequence while generating each word in the translated output.

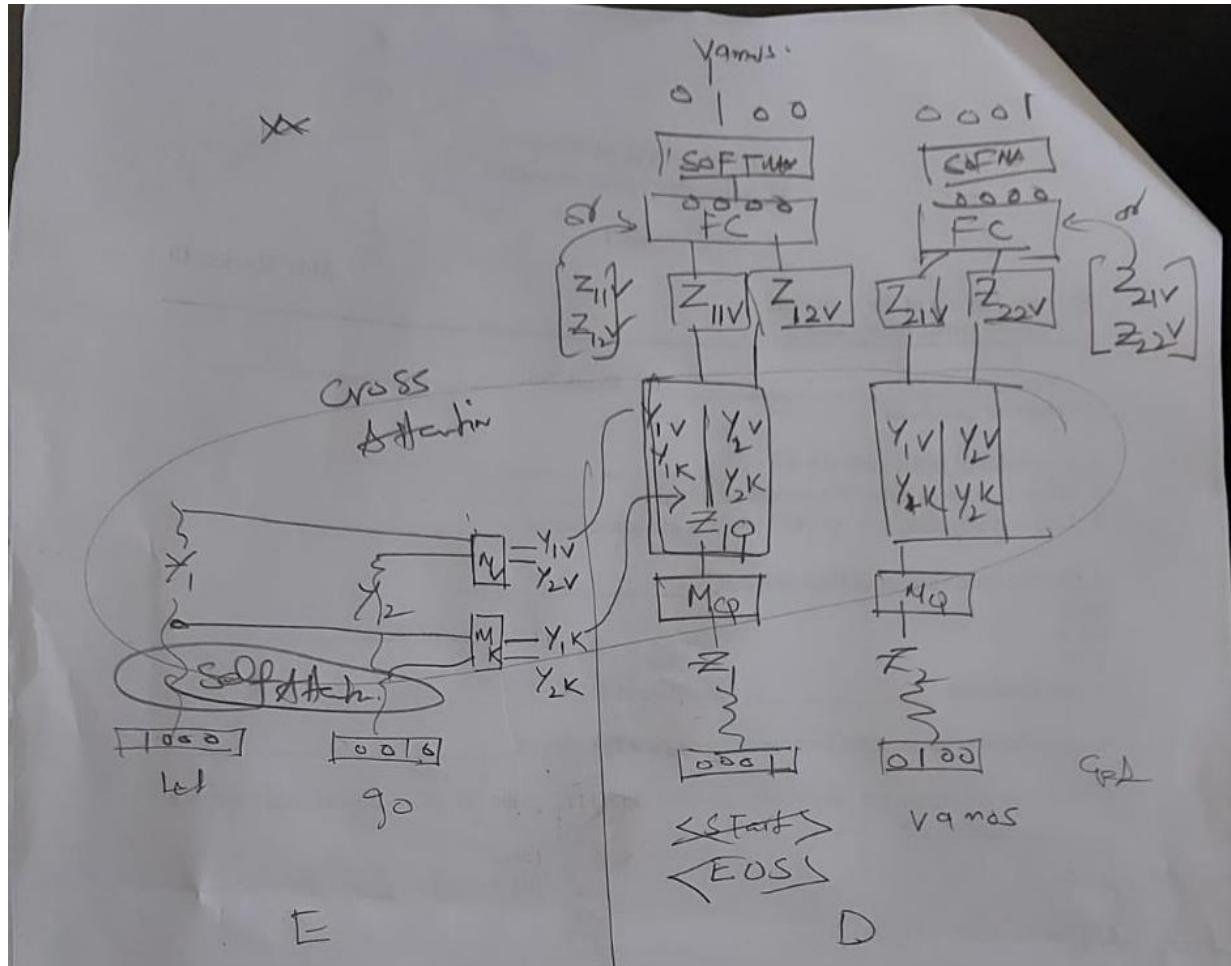
Let's go through a simple translation example using the Transformer model and see how the **cross-attention mechanism** (encoder-decoder attention) works step-by-step.

### Example:

Imagine translating the English sentence "She ate an apple" into French as "Elle a mangé une pomme".

The key goal of the cross-attention mechanism is to help the decoder pay attention to the most relevant parts of the encoded input sequence while generating each word in the translated output.





$\Rightarrow$  Decoder decode the  
encoding int translated  
phrase, Vamos

Let's  $g_2 \Rightarrow \text{vamos}$

Next time:

- WordtoVec encoding
- Positional encoding,
- Encoder layers
- Deconder layers



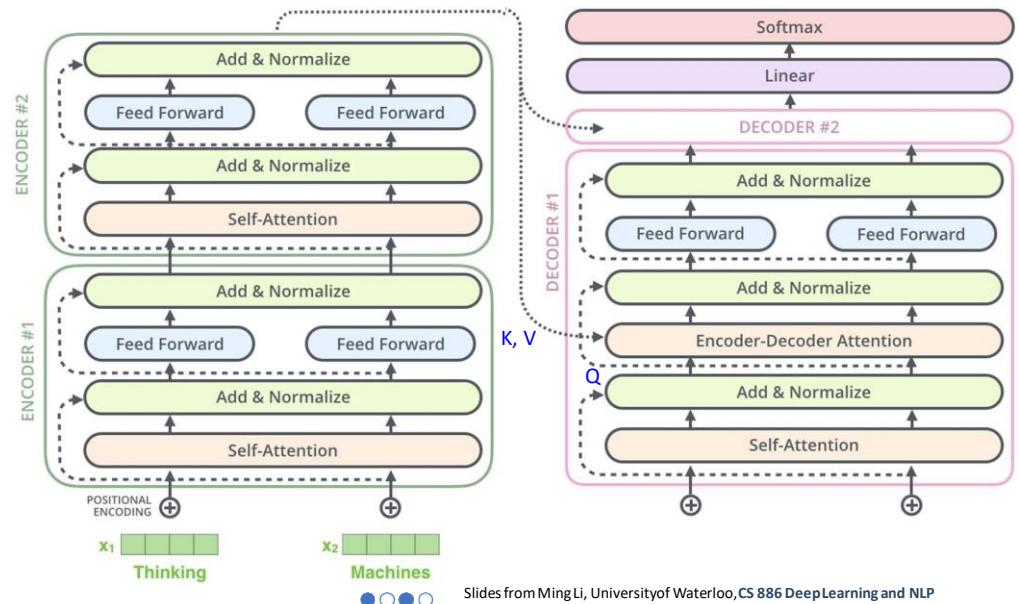
## Lecture 4

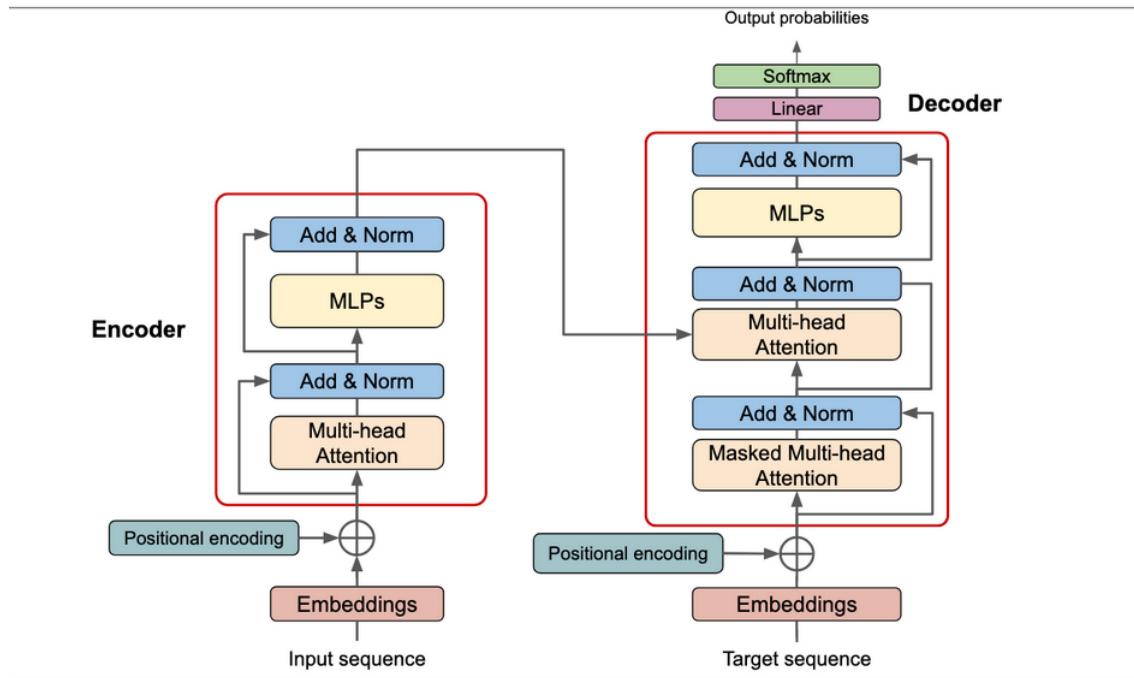
- Cross attention
- Mask multihead layers



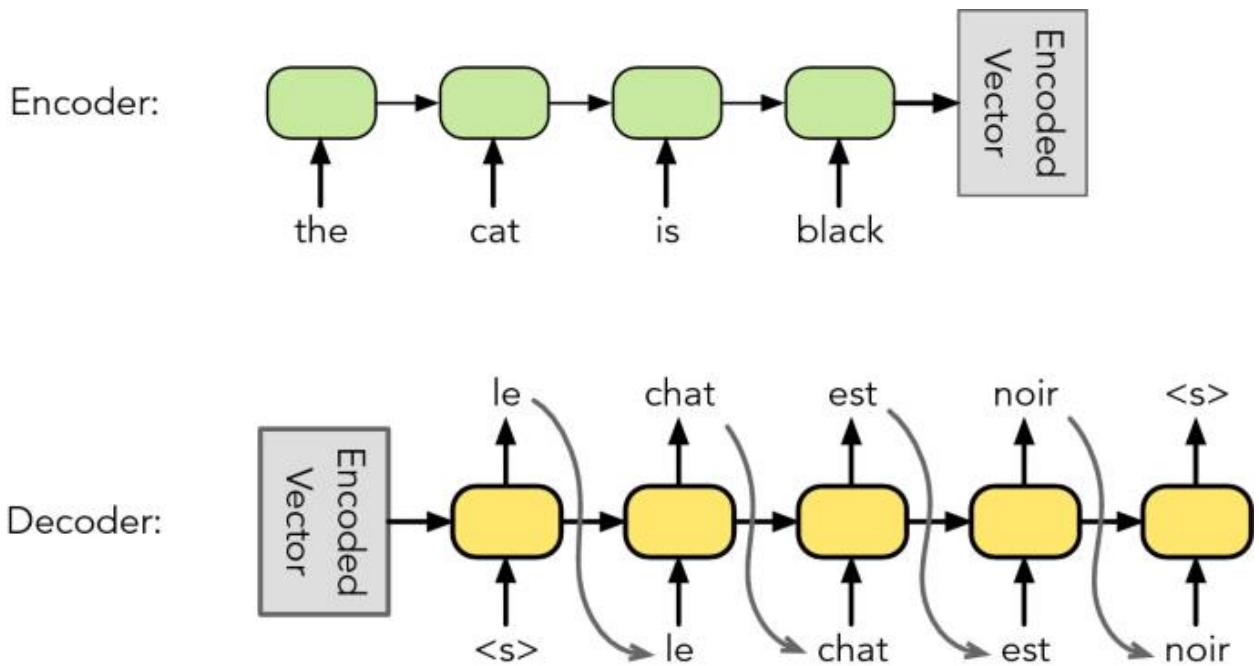
### Attention and Transformers

## The complete transformer





<https://deepprevious.github.io/posts/001-transformer/>



#### Deconder: Output in auto regressive way

#### What is Masked Self-Attention?

- Masked self-attention is a mechanism used in the decoder of transformer models that allows the model to focus on previously generated words while generating the next word.
- It prevents the model from attending to future tokens, ensuring that the generation is autoregressive (i.e., it only considers words that have already been generated).

#### Cross Attention

- Cross Attention is a mechanism used in transformer models, particularly in the decoder, to allow the model to focus on relevant parts of the input sequence (provided by the encoder) while generating output.

#### Context Vector:

- Represents the relationship between the current query (from the decoder) and the relevant information from the encoder's output.
- Provides crucial context from the input sentence to aid in generating a coherent output.

#### Hidden State: Tells the model what has been generated so far.

- Helps guide the model in predicting the next word based on previously generated words.
- `Hidden_state_current = Hidden_state_previous + Context_vector_current`

#### Example:

Translating "The cat eats fish" into French ("Le chat mange du poisson")

Suppose each word's embedding (from the encoder) is represented as a vector with three dimensions (for simplicity). We'll also assign the decoder's hidden state a similar 3-dimensional vector representation.

a. Encoder Output

Let's say the encoder generates embeddings for each word in the English sentence "The cat eats fish":

- `E\_the = [0.5, 0.1, 0.3]`

- `E\_cat = [0.6, 0.4, 0.2]`

- `E\_eats = [0.7, 0.3, 0.5]`

- `E\_fish = [0.2, 0.8, 0.6]`

[they are contextual vectors that contain, semantic, positional, and contextual information]

b. Decoder Hidden State Evolution (Example Values)

At each decoding step, the decoder will update its hidden state based on self-attention, cross-attention context provided by the encoder's embeddings

#### Step 1: Generating the First Word ("Le")

1. Initialization with `<start>` Token:

- The decoder starts with a hidden state initialized from the `<start>` token, say:

Hidden\_state\_0 = [0.1, 0.1, 0.1]

2. Cross-Attention Calculation:

- Attention weights are calculated for each encoder embedding to focus on the most relevant words for generating the first word. Suppose the weights for "The" and "cat" are higher:

Attention weights: [Weight\_the = 0.6, Weight\_cat = 0.4, Weight\_eats = 0.0, Weight\_fish = 0.0]

Attention weights: [0.6, 0.4, 0.0, 0.0] // <start> cross attention with encoder embedding

- Compute the context vector by applying these weights to the encoder embeddings:

Context\_vector\_1 = (0.6 \* [0.5, 0.1, 0.3]) + (0.4 \* [0.6, 0.4, 0.2]) = [0.54, 0.22, 0.26]

Context\_vector\_1 = [0.54, 0.22, 0.26]

3. Update Hidden State:

- The initial hidden state and context vector combine to form the new hidden state for predicting the first word. Suppose a simple addition operation (for illustration):

$$\text{Hidden\_state\_1} = \text{Hidden\_state\_0} + \text{Context\_vector\_1} = [0.1, 0.1, 0.1] + [0.54, 0.22, 0.26] =$$

$$\text{Hidden\_state\_1} = [0.64, 0.32, 0.36]$$

- This hidden state represents the context needed to generate the word "Le"

**Hidden\_state\_1** gives the model the context needed to understand what word to generate next. The softmax function then helps convert this context into a probability for each possible word, allowing the model to choose the best word, which in this case is "Le"

#### #### Step 2: Generating the Second Word ("chat")

##### 1. Masked Self-Attention on "Le":

- The hidden state now includes the word "Le," so the decoder applies masked self-attention to reference this.

##### 2. Cross-Attention Calculation: (here we only take "Le" embeddings as query vector)

- New cross-attention weights focus more on `E\_cat` for the next word:

Attention weights: [Weight\_the = 0.1, Weight\_cat = 0.9, Weight\_eats = 0.0, Weight\_fish = 0.0]

Attention weights: [ 0.1, 0.9, 0.0, 0.0 ]

- Compute the Context Vector:

$$\text{Context\_vector\_2} = (0.1 * [0.5, 0.1, 0.3]) + (0.9 * [0.6, 0.4, 0.2]) = [0.59, 0.37, 0.21]$$

$$\text{Context\_vector\_2} = [0.59, 0.37, 0.21]$$

##### 3. Update Hidden State:

- Combine this with the previous hidden state:

$$\text{Hidden\_state\_2} = \text{Hidden\_state\_1} + \text{Context\_vector\_2} = [0.64, 0.32, 0.36] + [0.59, 0.37, 0.21] = [1.23, 0.69, 0.57]$$

$$\text{Hidden\_state\_2} = [1.23, 0.69, 0.57]$$

- This hidden state is now used to generate "chat."

---

#### #### Step 3: Generating the Third Word ("mange")

##### 1. Masked Self-Attention on "Le chat":

- The hidden state includes both "Le" and "chat," so masked self-attention allows it to reference these without "looking ahead."

## 2. Cross-Attention Calculation:

- Attention now shifts to `E\_eats`:

Attention weights[Weight\_the = 0.0, Weight\_cat = 0.0, Weight\_eats = 1.0, Weight\_fish = 0.0]

Attention weights[0.0, 0.0, 1.0, 0.0]

- Context vector:

Context\_vector\_3 = [0.7, 0.3, 0.5]

## 3. \*\*Update Hidden State\*\*:

- Combine with the previous hidden state:

Hidden\_state\_3 = Hidden\_state\_2 + Context\_vector\_3 = [1.23, 0.69, 0.57] + [0.7, 0.3, 0.5] = [1.93, 0.99, 1.07]

Hidden\_state\_3 = [1.93, 0.99, 1.07]

- This hidden state is used to predict "mange."

---

This process continues for each word in the sequence. Each hidden state is a combination of prior context and relevant encoder information, allowing the decoder to maintain a fluent and contextually accurate translation.

These hidden states, which evolve with each generation step, ensure that the translation flows naturally and aligns with the meaning of the original sentence.