

GPT (Generative Pre-trained Transformer)

2023. 8

Yongjin Jeong, KwangWoon University

[참고] 본 자료에는 인터넷에서 다운받아 사용한 그림이나 수식들이 있으니 다른 용도로 사용하거나 외부로 유출을 금해 주시기 바랍니다.

Seq2seq based on RNN

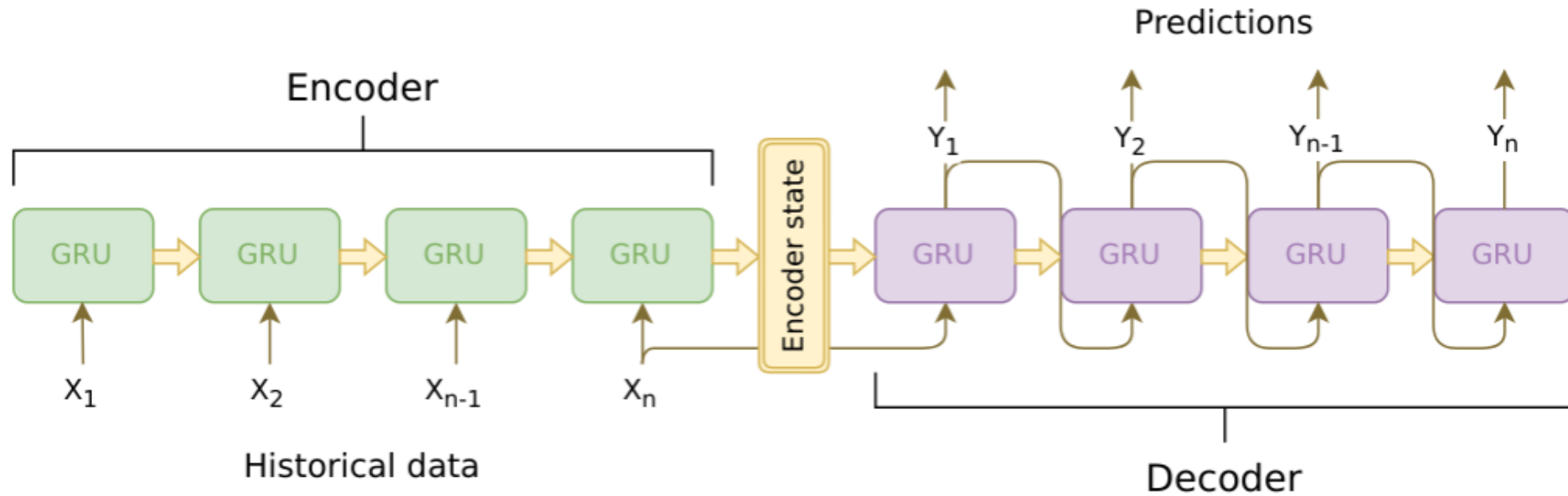


Image source: https://jeddy92.github.io/JEddy92.github.io/ts_seq2seq_intro/

- One fixed context vector (the last hidden state) is transferred to Decoder.

Attention with RNN

- **General RNN Encoder-decoder**

- Encoder

$$h_t = f(x_t, h_{t-1})$$

$$c = q(\{h_1, \dots, h_{T_x}\})$$

$\mathbf{x}=(x_1, \dots, x_{T_x})$: input sequence

$h_t \in R^n$: a hidden state at time t ,

c : context vector generated from sequence of hidden states

← simply, weighted average of some values

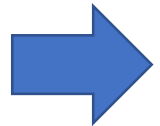
- Decoder:

- trained to predict the next word $y_{t'}$ given the context vector (c) and all the previously predicted words $\{y_{t'} \dots, y_{t'-1}\}$

$$p(\mathbf{y}) = \prod_{t=1}^T p(y_t \mid \{y_1, \dots, y_{t-1}\}, c)$$

$\mathbf{y}=(y_1, \dots, y_{T_y})$: output sequence

with RNN, modeled as



$$p(y_t \mid \{y_1, \dots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c)$$

S_t : hidden state of RNN

Attention with RNN – Align and Translate

- **Decoder with Attention**

In a new model architecture, we define each conditional probability in Eq. (2) as:

$$p(y_i | y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i), \quad (4)$$

where s_i is an RNN hidden state for time i , computed by

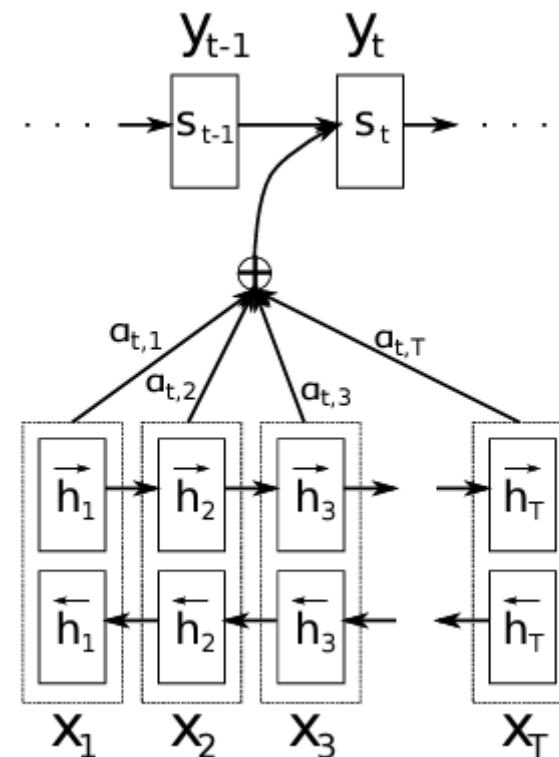
$$s_i = f(s_{i-1}, y_{i-1}, c_i).$$

Unlike the encoder–decoder approach (see Eq.(2)), here the probability is conditioned on a distinct context vector c_i for each target word y_i .

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j \quad (\text{weighted sum of annotations } h_j) \quad h_j = \begin{bmatrix} \vec{h}_j^\top; \overleftarrow{h}_j^\top \end{bmatrix}^\top$$

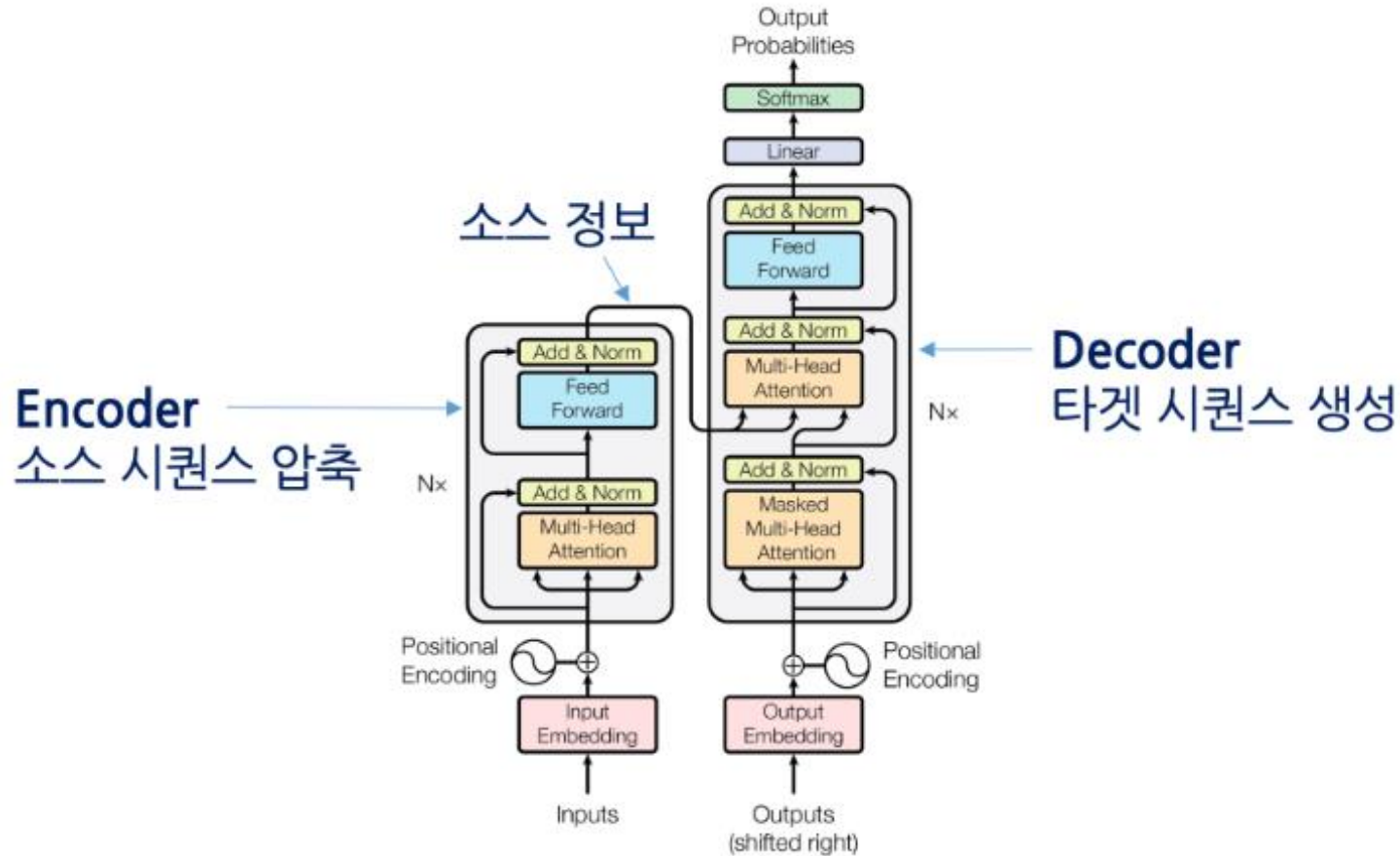
weights $\longrightarrow \alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}, \text{ where } e_{ij} = a(s_{i-1}, h_j)$

alignment model ; it scores how well the inputs around position j and the output at position i match. The score is based on the RNN hidden states $i-1$ (just before emitting y_i) and the j -th annotation h_j of the input sentence. (can be parameterized as a Forward Neural Network)



The Transformer

- Without RNN



Encoder:

- Multi-head self-attention mechanism
- A simple position-wise fully connected feed-forward network

Decoder:

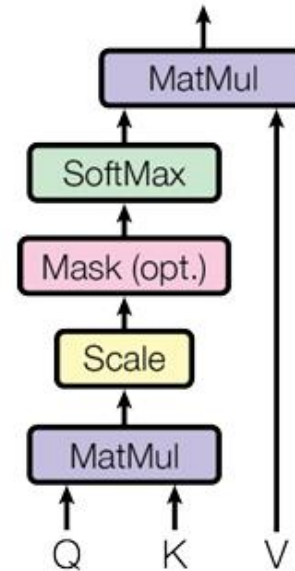
- **Masked** Multi-head self-attention mechanism (ensures that the predictions for position i can depend only on the known outputs at positions less than i)
- A simple position-wise fully connected feed-forward network
- Multi-head attention over the output of the encoder stack

Multi-Head Attention

- **Attention**

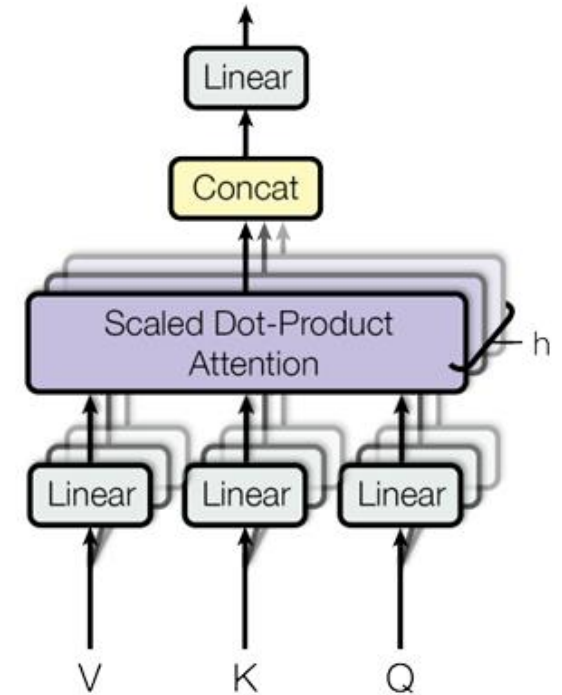
- Mapping a **query** and a set of key-value pairs to an output
- The output is computed as a weighted sum of the values, where the weight is computed by a compatibility function of the query with the corresponding key
- Query, key, value are packed into matrices Q , K , V respectively.

Scaled Dot-Product Attention



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

Multi-Head Attention



$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

$$\text{attention}(q, k, v) = \sum_i \text{similarity}(q, k_i) * v_i$$

$$W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$$

$$W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$$

The Transformer - Attention

- **Multi-head attention**

- In "encoder-decoder attention" layers, the queries come from the previous decoder layer, and the memory keys and values come from the output of the encoder. (This allows every position in the decoder to attend overall positions in the input sequences. This mimics the typical encoder-decoder attention mechanisms in sequence-to-sequence models.)
- **Self-attention layers in Encoder:** all of the keys, values, and queries come from the same place (the output of the previous layer in the encoder). Each position in the encoder can attend to all positions in the previous layer of the encoder.
- **Self-attention layers in Decoder:** we allow each position in the decoder to attend to all positions in the decoder up to and including that position. We need to prevent leftward information flow in the decoder to preserve the auto-regressive property. We implement this inside of scaled dot-product attention by masking out(setting to $-\infty$) all values in the input of the softmax which correspond to illegal connections.

Attention by Query, Key, and Value

- Query, Key, Value
 - Simply generated with three different matrices (linear **transformation**) <- why?
 - **Query**: the input or prompt provided by the user.
 - **Key**: a piece of information used to retrieve relevant knowledge. (In the case of ChatGPT, the key is typically generated from the query.)
 - **Value**: the information retrieved by using the key
- What is the **Weight** for Value (Query * Key):
 - to calculate the (cosine) similarity or dot product between two vectors (a scalar value that represents the **similarity or relevance between the key and the query**)
 - Multiplying the key and query is a common mathematical operation used to calculate the similarity or **dot product** between two vectors.

- ❖ Intuition: for the sentence "Jane visits Africa."

 - When your eyes see ***jane***, your brain looks for **the most related word** in the rest of the sentence to understand what **jane** is about (query). Your brain focuses or attends to the word **visit** (key).

Attention by Query, Key, and Value

- Why **transformation** for Query, Key, Value?
 - we could still use the original encoder state vectors as the queries, keys, and values. So, **why we need the transformation?**
 - Simply a matrix multiplication: (I is the input (encoder) state vector)
 $\text{Query} = I \times W^Q$, $\text{Key} = I \times W^K$, $\text{Value} = I \times W^V$
 - (1) If we do not transform input vectors, the dot product for computing the weight for each input's value will always yield a maximum weight score for the individual input token itself.
(input i -번째 위치에서 n 개의 attention weight 를 계산하면 그 자신에서 가장 높은 값이 나옴)
 - For example, for the **pronoun** token, we need it to attend to its referent, not the pronoun token itself.
 - (2) The transformation **may yield better representations for Query, Key, and Value**. (see the next slide)
 - (3) conversion of the input vector into a space with a desired dimension (practically useful)

Attention by Query, Key, and Value

- transformation may yield better representation for computing similarity between two vectors. (<https://youtu.be/K38wVcdNuFc?t=10>)

Latent Semantic Indexing

- Assume there are $n=5$ movies viewed and rated by $m=7$ persons, as below, the Singular Value Decomposition (SVD) of the input matrix yields 2 majors groups (topics): Data are from: <https://www.youtube.com/watch?v=K38wVcdNuFc>

Reduced
SVD

$$A_{m \times n} = U_{m \times r} \Sigma_{r \times r} (V_{n \times r})^T \approx U_{m \times d} \Sigma_{d \times d} (V_{n \times d})^T$$

$$\begin{array}{l}
 \text{P1} \\ \text{P2} \\ \text{P3} \\ \text{P4} \\ \text{P5} \\ \text{P6} \\ \text{P7}
 \end{array}
 \begin{bmatrix}
 1 & 1 & 1 & 0 & 0 \\
 3 & 3 & 3 & 0 & 0 \\
 4 & 4 & 4 & 0 & 0 \\
 5 & 5 & 5 & 0 & 0 \\
 0 & 2 & 0 & 4 & 4 \\
 0 & 0 & 0 & 5 & 5 \\
 0 & 1 & 0 & 2 & 2
 \end{bmatrix}
 =
 \begin{bmatrix}
 0.13 & 0.02 & -0.01 \\
 0.41 & 0.07 & -0.03 \\
 0.55 & 0.09 & -0.04 \\
 0.68 & 0.11 & -0.05 \\
 0.15 & -0.59 & 0.65 \\
 0.07 & -0.73 & -0.67 \\
 0.07 & -0.29 & 0.32
 \end{bmatrix}
 \times
 \begin{bmatrix}
 12.4 & 0 & 0 \\
 0 & 9.5 & 0 \\
 0 & 0 & 1.3
 \end{bmatrix}$$

$$\begin{array}{l}
 \text{Star Wars} \\ \text{The Matrix} \\ \text{Iron man} \\ \text{U got mail} \\ \text{Titanic}
 \end{array}
 \begin{bmatrix}
 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\
 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \\
 0.40 & -0.80 & 0.40 & 0.09 & 0.09
 \end{bmatrix}$$

$$\text{P8} = [5 \ 0 \ 0 \ 0 \ 0] \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \end{bmatrix}^T = [2.81 \ 0.63]$$

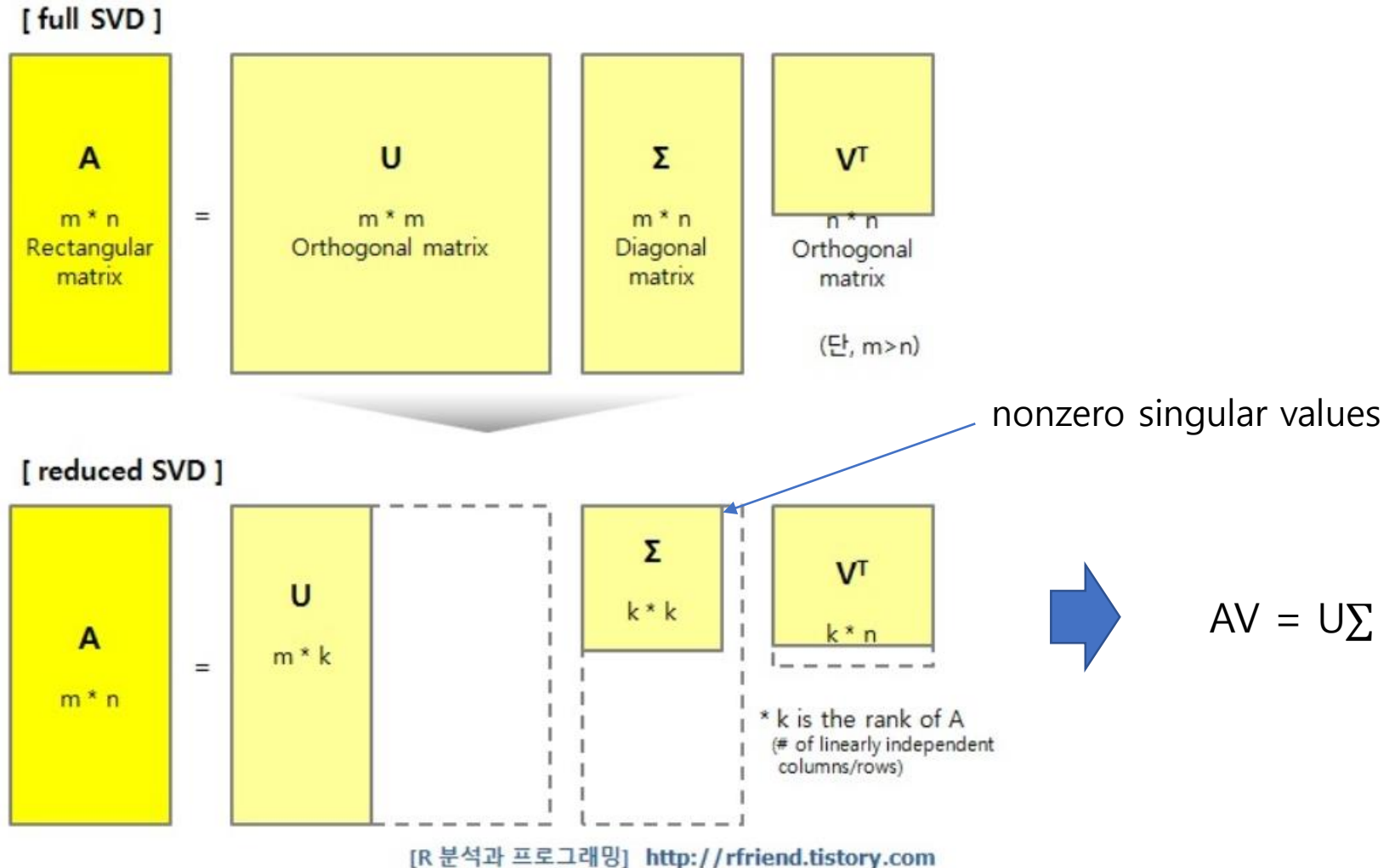
$$\text{P9} = [0 \ 4 \ 5 \ 0 \ 0] \times \begin{bmatrix} 0.56 & 0.59 & 0.56 & 0.09 & 0.09 \\ 0.12 & -0.02 & 0.12 & -0.69 & -0.69 \end{bmatrix}^T = [5.18 \ 0.52]$$

Cosine similarity is 0.0 in 5-d vector Cosine similarity is now 0.99

- By multiplying an input vector with a matrix V (from the SVD), we get a better representation for computing the compatibility between two vectors, if these two vectors are similar in the topic space as shown in the example in the figure.
- And, these matrices for transformation can be learned in a neural network!
- (before transformation)
similarity(P8,P9)=0
- (after transformation)
similarity(P8',P9') = 0.99

Attention by Query, Key, and Value

- Reduced SVD



GPT and BERT

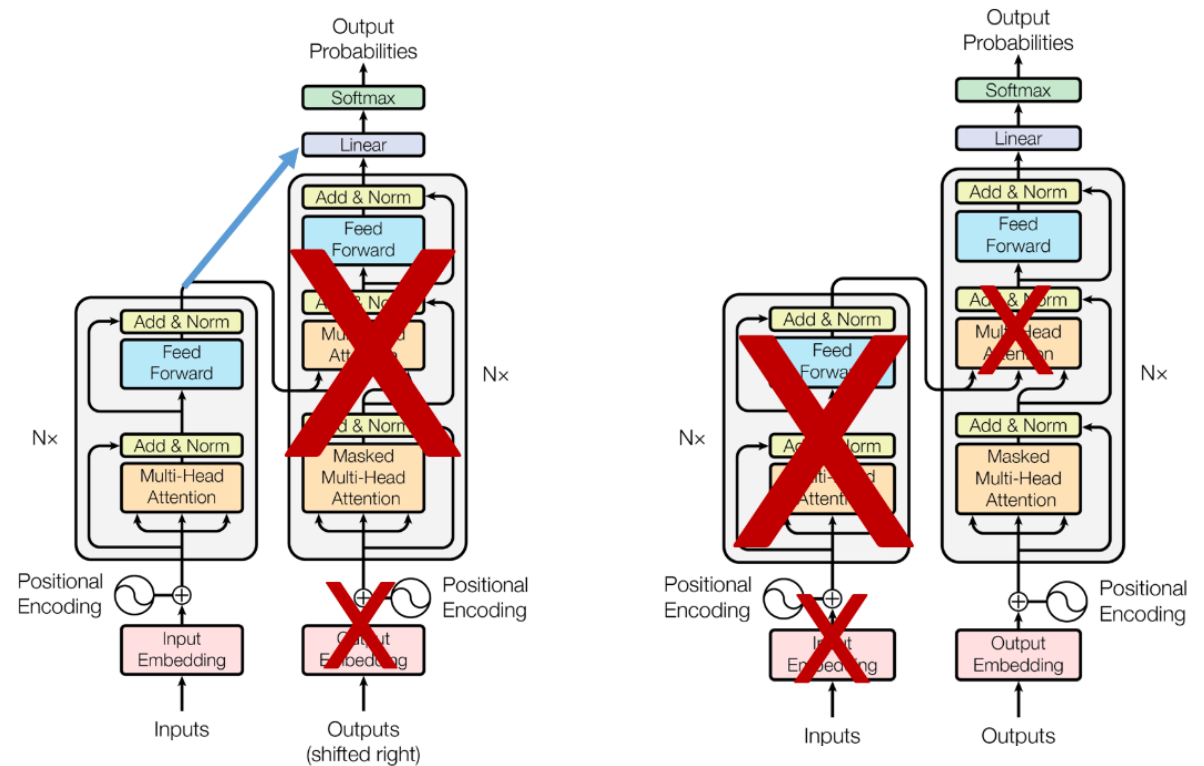
- GPT:
 - Language model (predicts the next word)
 - Unidirectional
- BERT(Bidirectional Encoder Representations from Transformers)
 - Masked Language Model
 - Predicts the blanks in the middle of sentences (bidirectional)

GPT

어제 카페 갔었어 거기 사람 많더라

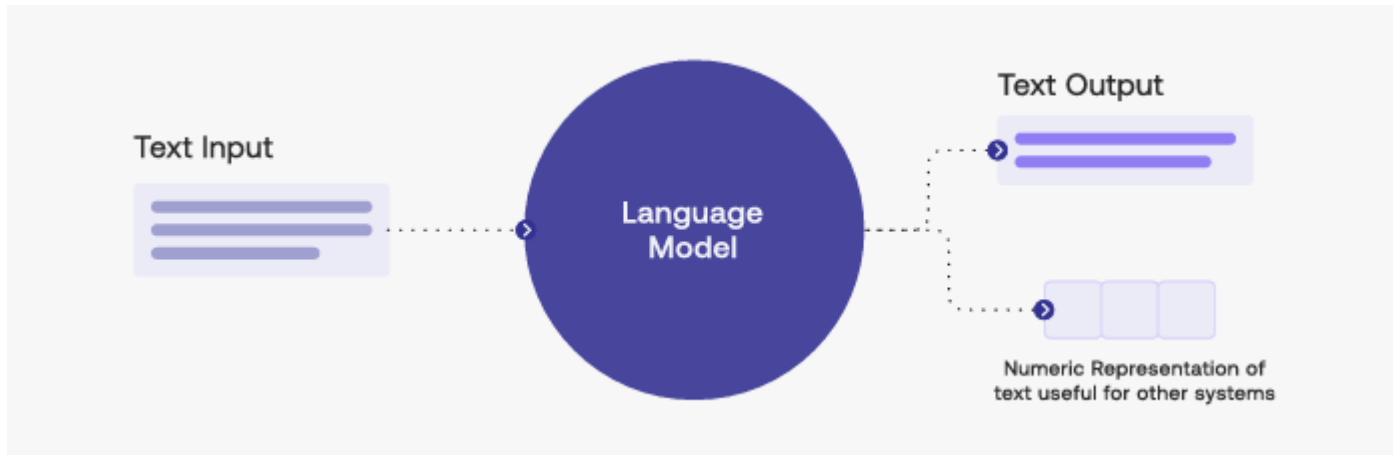
BERT

어제 카페 갔었어 사람 많더라



What is GPT?

- Large Language Model (LLM)
 - Predicts text which should follow from a given prompt using probabilistic methods
 - Machine learning model trained on massive amounts of text



- GPT-3 Language Model
 - 499 Billion tokens (words) drawn from the Web, books, and Wikipedia
 - Trained to predict the next word in a sequence

GPT-3

- 175 Billion parameters
- Trained with 45 TB of text data which includes sources from Wikipedia and books
- It has 96 decoder layers and is built on a system with 285k CPU cores, 10K GPUs, and 400 Gbps network connectivity for each GPU server
- Input sequence
 - Up to 2048 tokens (words encoded as numeric values)
 - Encoding allows word meanings to be compared (e.g., similar words have similar numeric values)
- Output sequence
 - Up to 2048 “guesses” at the most likely next word in the sequence
- Key advantages
 - **Attention** (a technique for determining how important each token is to the output)
- Risk of LLMs (by Dr. Josh Locklair)
 - LLMs absorb the worldview and biases of their human creators
 - LLMs “blindly” propagate information without consideration of human impact
 - LLMs cannot consider the truth of the text they generate

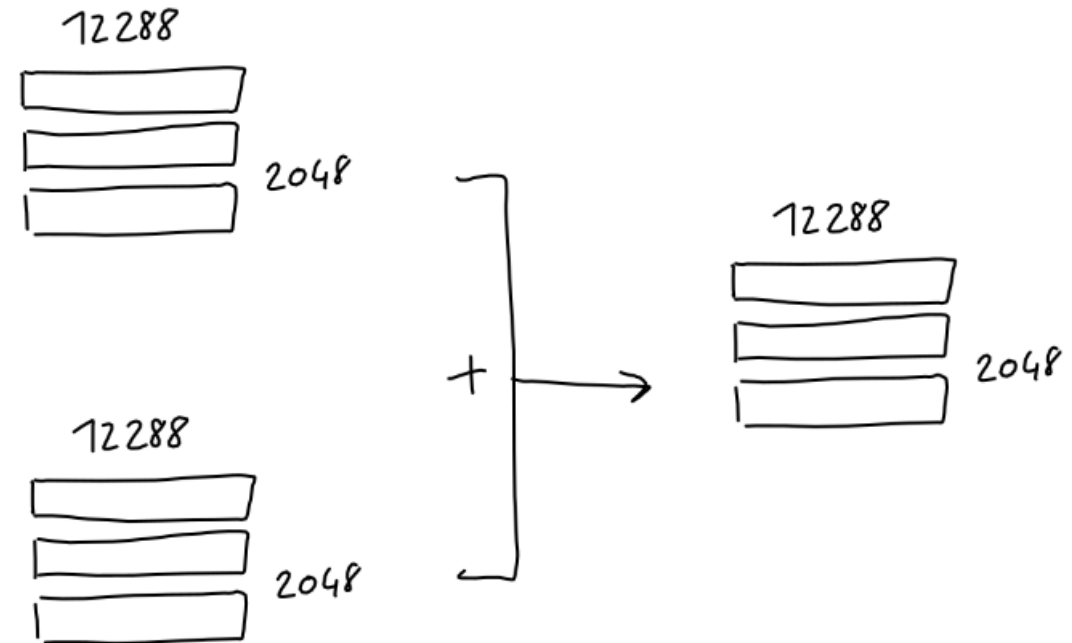
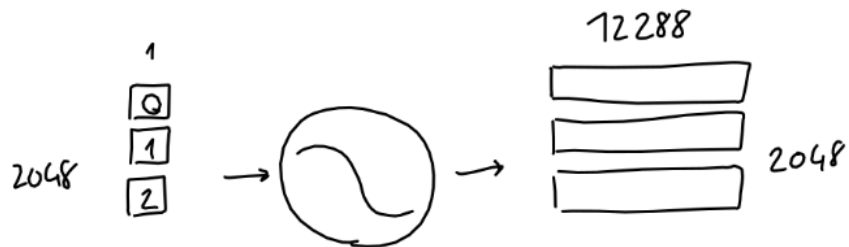
How Transformers are used in ChatGPT?

- **Tokenization**
 - The input text is divided into individual tokens and each token is assigned a unique numerical representation.
- **Embedding**
 - Each token is then transformed into a high-dimensional numerical vector, called an embedding
- **Positional Encoding**
- **Transformer Layers** (Attention Mechanism + Feedforward Network)
- **Decoder**
 - It takes the transformed input sequence and generates the output sequence token by token.
 - At each step, the model uses a combination of the previously generated tokens and the contextual information from the input sequence to predict the next token.
- **Generation**
 - The output tokens are decoded back into readable text, and the process can be repeated to generate a response or continuation.

GPT-3 (Encoding)

- One-hot encoding, Embedding, Positional encoding

$$\begin{matrix} 50 & 257 \\ 0 & \begin{bmatrix} 0 & 0 & 0 & 1 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & 1 & \dots & 0 \\ \vdots & & & & & & \\ 0 & 0 & 1 & 0 & 0 & \dots & 0 \end{bmatrix} \\ 2048 & \end{matrix} \times \begin{matrix} 50 & 257 \\ W_E \\ 12 & 288 \end{matrix} = \begin{matrix} 12 & 288 \\ 2048 \\ \begin{bmatrix} 0.1 & \dots & -0.2 \\ \vdots & \ddots & \vdots \\ 0.3 & \dots & -2.5 \end{bmatrix} \end{matrix}$$



GPT-3

- Multi-head Attention

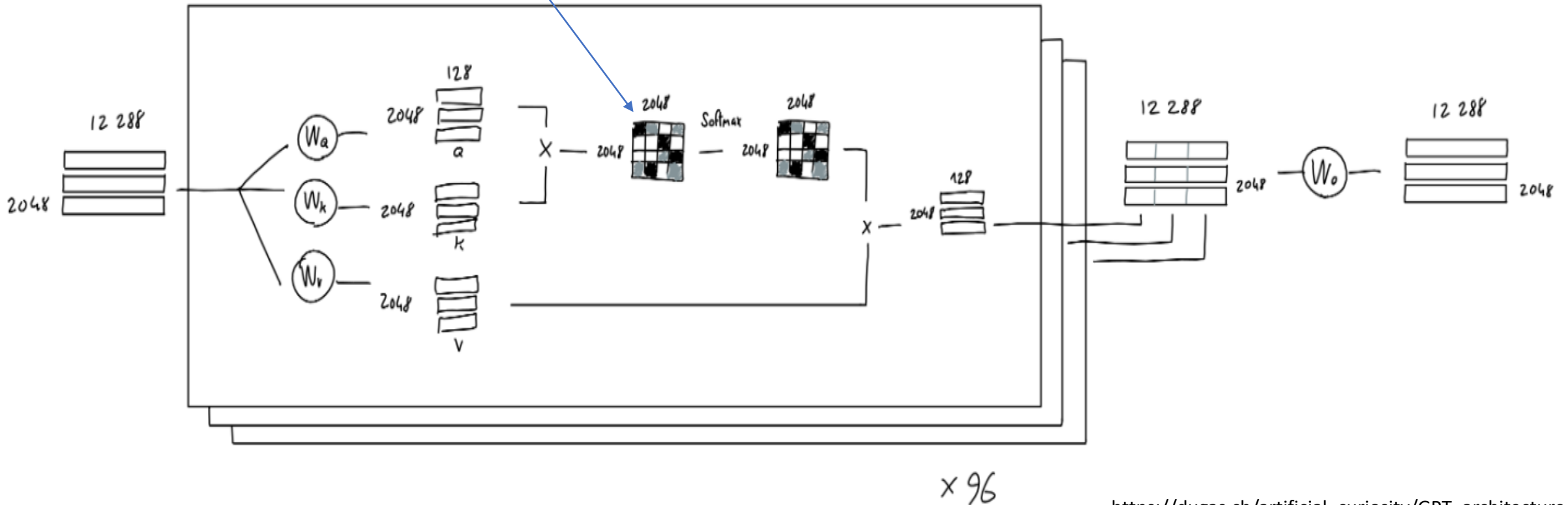
(q, K, V)

q = the vector representing a word

K and V = your memory, thus all the words that have been generated before.

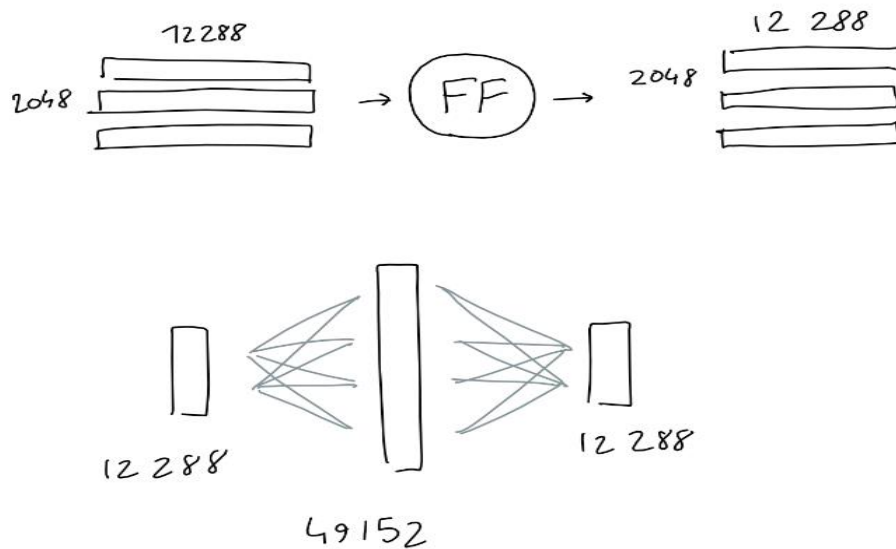
So what you do with attention is that you take your current query (word in most cases) and look in your memory for similar keys.

weights

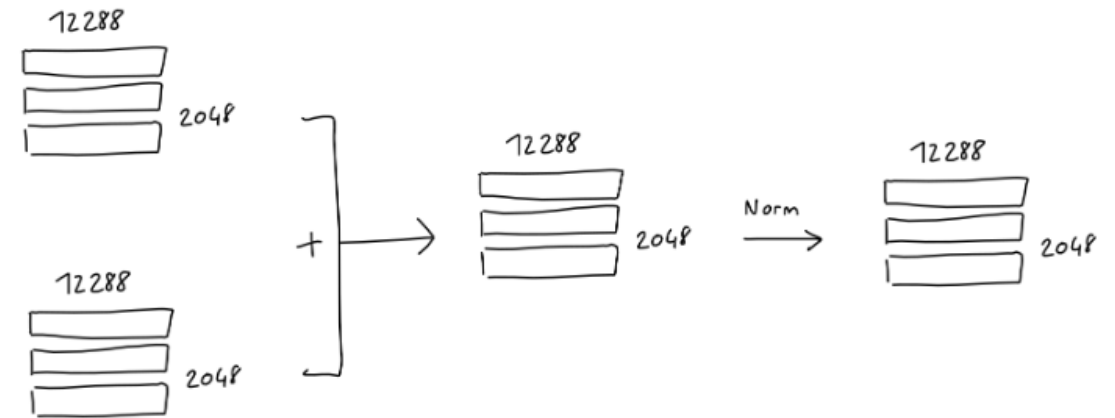


GPT-3

- Feed Forward
 - Hidden layer size = $4 * 12,288 = 49,152$

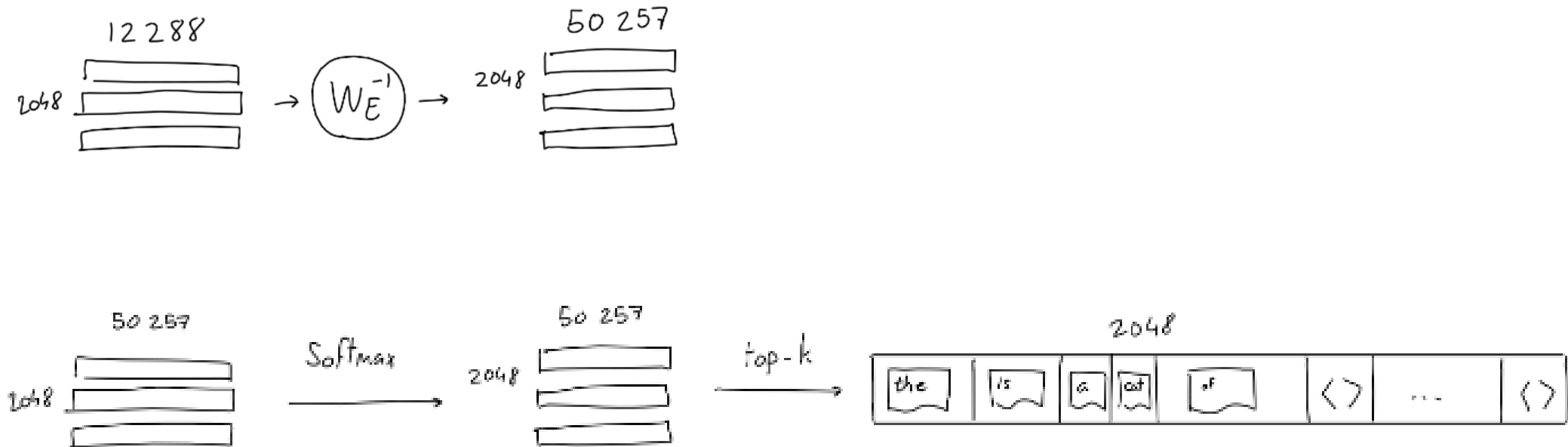


- Add and Norm



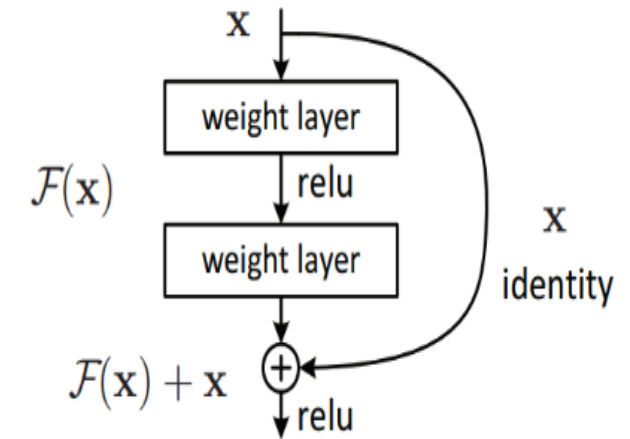
GPT-3 (Decoding)

- The input has been processed into a 2048 x 12288 matrix. This matrix is supposed to contain, for each of the 2048 output positions in the sequence, a 12288-vector of information about which word should appear. But how do we extract this information?



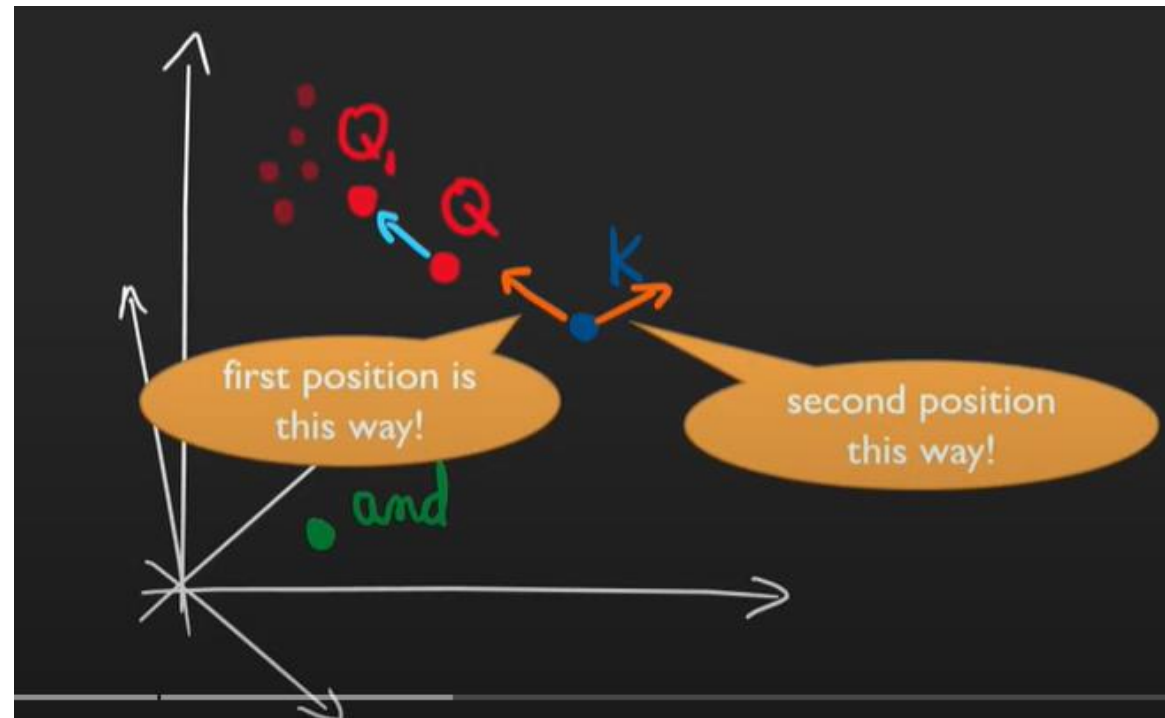
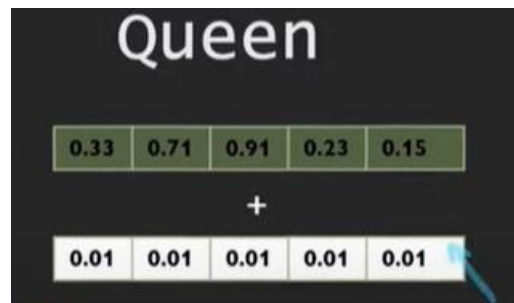
Residual Neural Networks (ResNets)

- Why ResNets?
 - **Addressing Vanishing Gradients**: Skip connections allow gradients to flow directly to earlier layers, and the gradients have a shorter path to propagate, making it easier to update the weights of earlier layers and improve the model's overall performance.
 - **Facilitating Training of Deeper Networks**: By allowing information to flow through shortcut connections, ResNets enable the training of much deeper neural networks.
 - **Promoting Feature Reuse**: The skip connections in ResNets facilitate the direct flow of information from earlier layers to later layers. This allows the later layers to access the activations and features from earlier layers, promoting feature reuse.
 - **Easier Optimization**: provide a form of identity mapping that makes it easier for the network to optimize its weights. This can speed up the convergence during training and enable more efficient optimization.



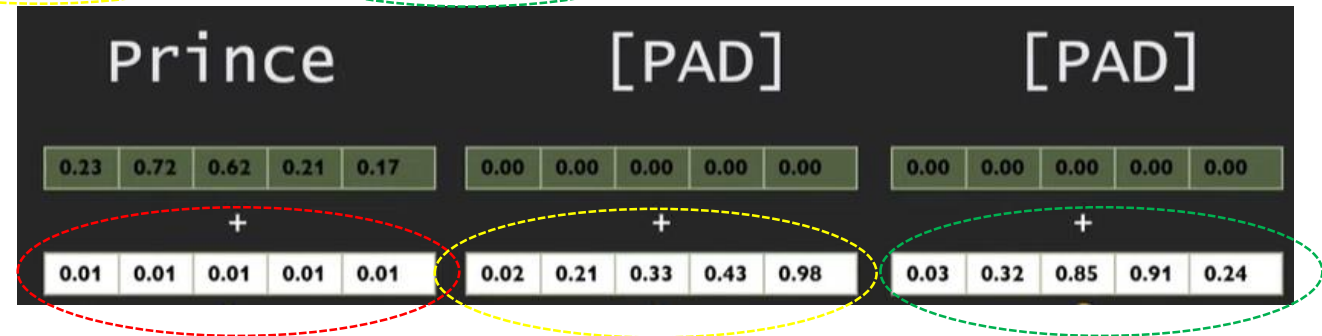
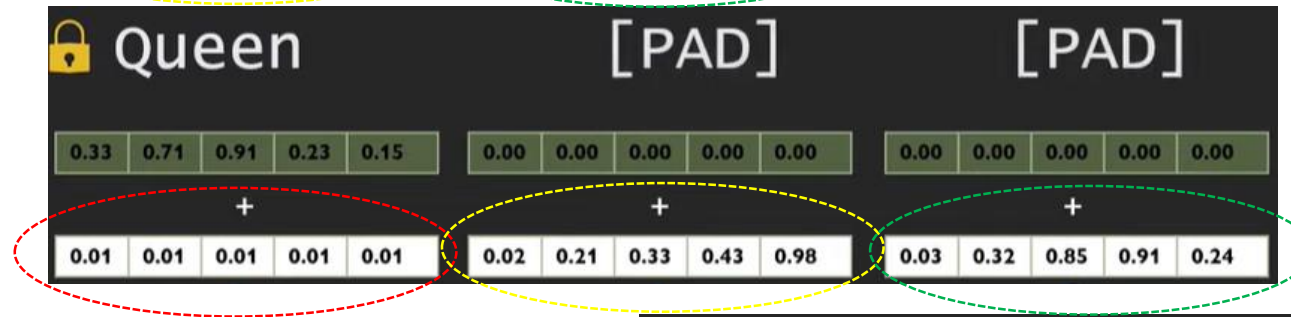
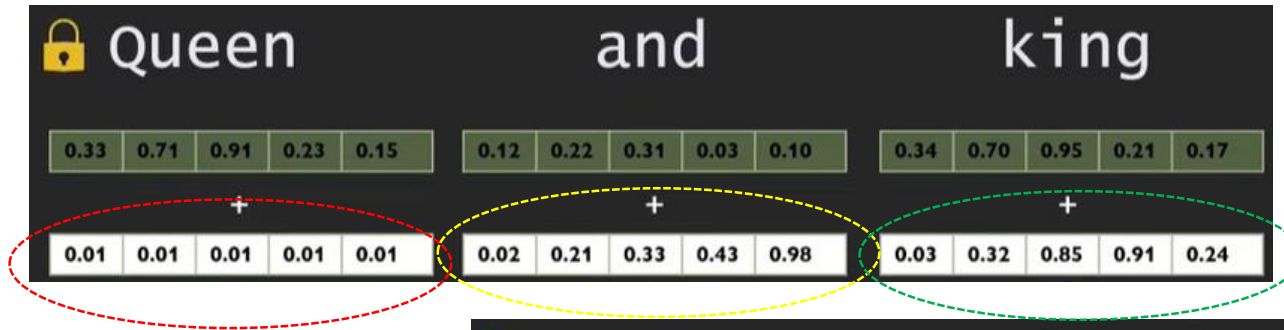
Positional Encoding

- What is positional encoding?
 - It incorporates the sequential order or position information of the tokens in a sequence.



Positional Encoding

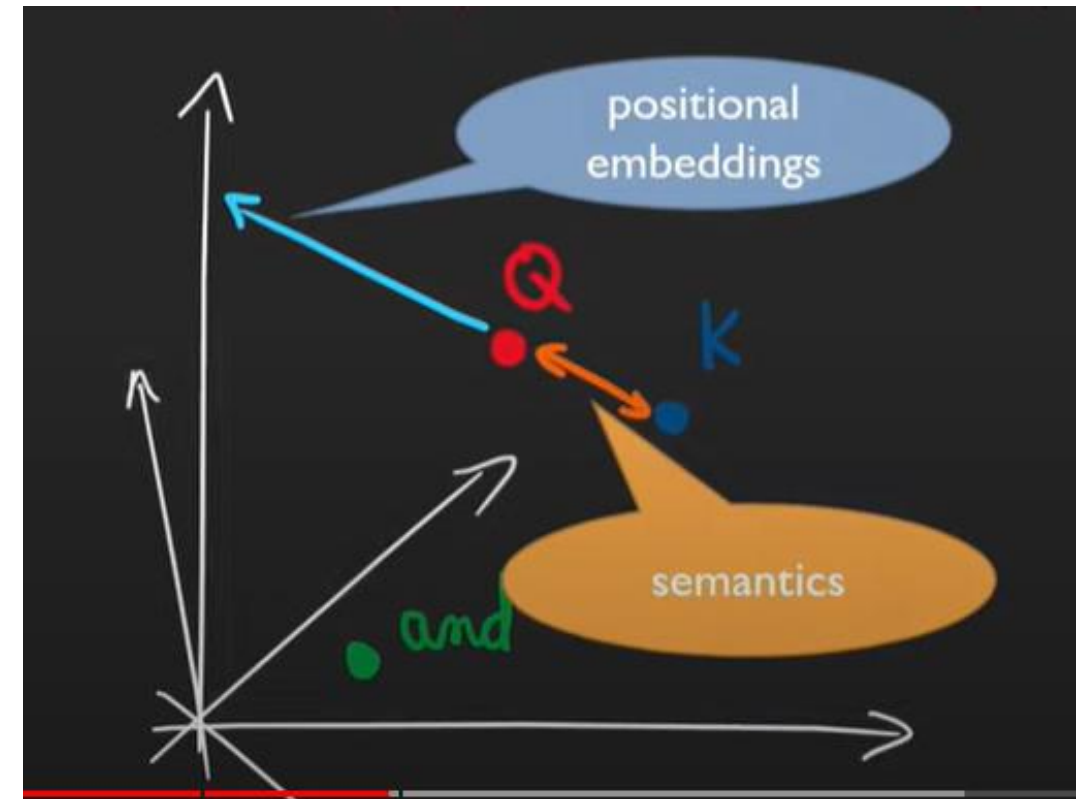
- Every position should have the same identifier, irrespective of the sequence



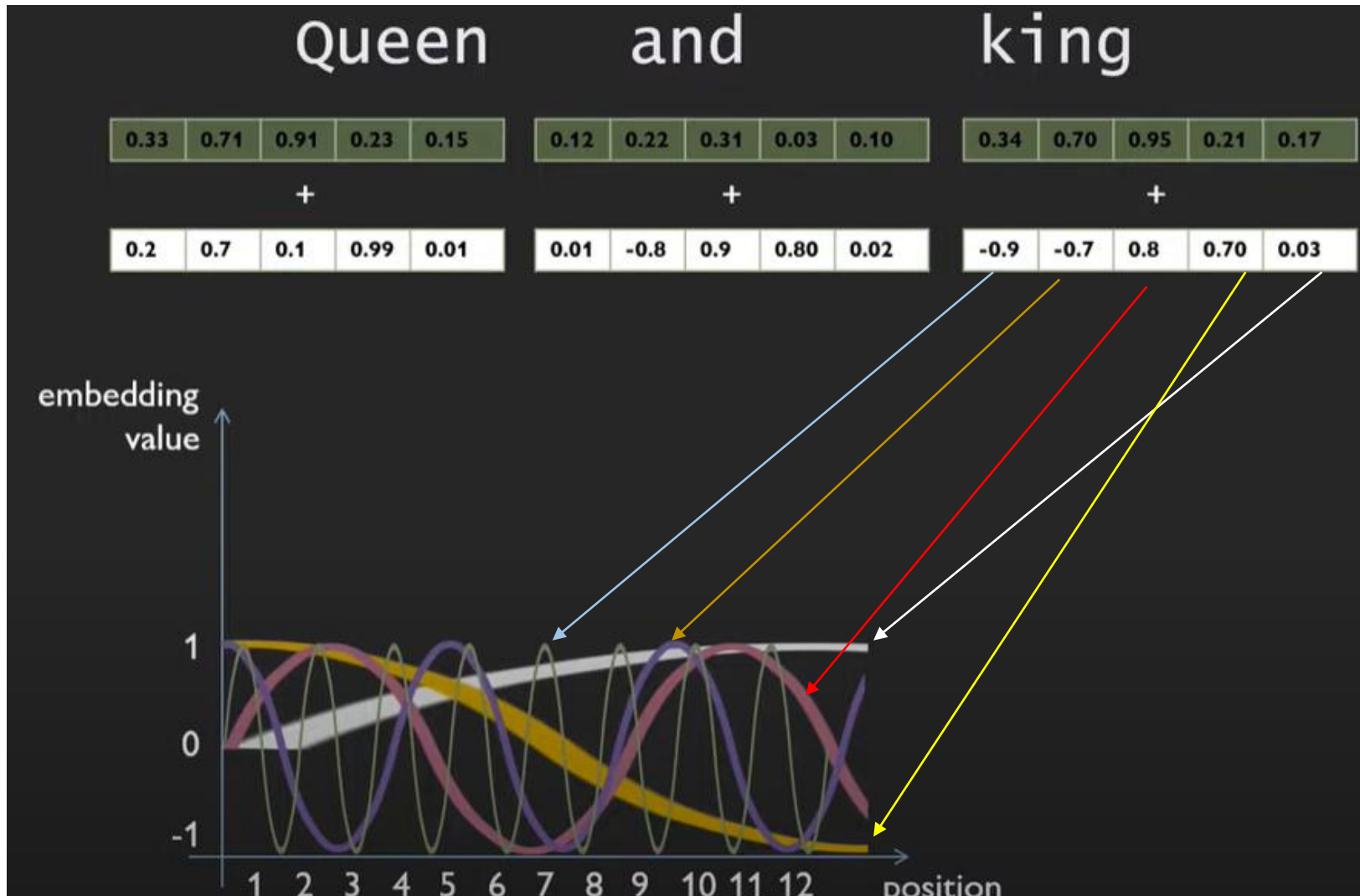
Positional Encoding

- Because the embeddings push the original vector by a bit into the “first token club”, they should not be too large. Otherwise, they push the vectors into very distinct subspaces where the positional similarity or dissimilarity overtones the semantic similarity.

| Queen | | | | | and | | | | | king | | | | |
|-------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| 0.33 | 0.71 | 0.91 | 0.23 | 0.15 | 0.12 | 0.22 | 0.31 | 0.03 | 0.10 | 0.34 | 0.70 | 0.95 | 0.21 | 0.17 |
| + | | | | | + | | | | | + | | | | |
| 1 | 0 | 0 | 0 | 0 | 20 | 0 | 0 | 0 | 0 | 300 | 0 | 0 | 0 | 0 |
| ✗ | | | | | ✗ | | | | | ✗ | | | | |



Positional Encoding



벡터의 각 element 마다 (다른 컬러로 표기) 주기를 달리 함으로써 단어 위치로 인한 positional vector 값들의 차이가 나타나게 함

Positional Encoding

- Positional Encoding
 - Usually designed to have a small magnitude compared to the token embedding values, which allows the model to retain the meaningful information from the original token.
- Sinusoidal Positional Encoding

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

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

- k: position of an object in the input sequence ($0 \leq k < L$)
- d: dimension of the output embedding
- $P(k, j)$: position matrix
- n: user-defined scalar (e.g. set to 10,000)
- i: used for mapping to column indices $0 \leq i < d/2$, with a single value of i maps to both sine and cosine functions

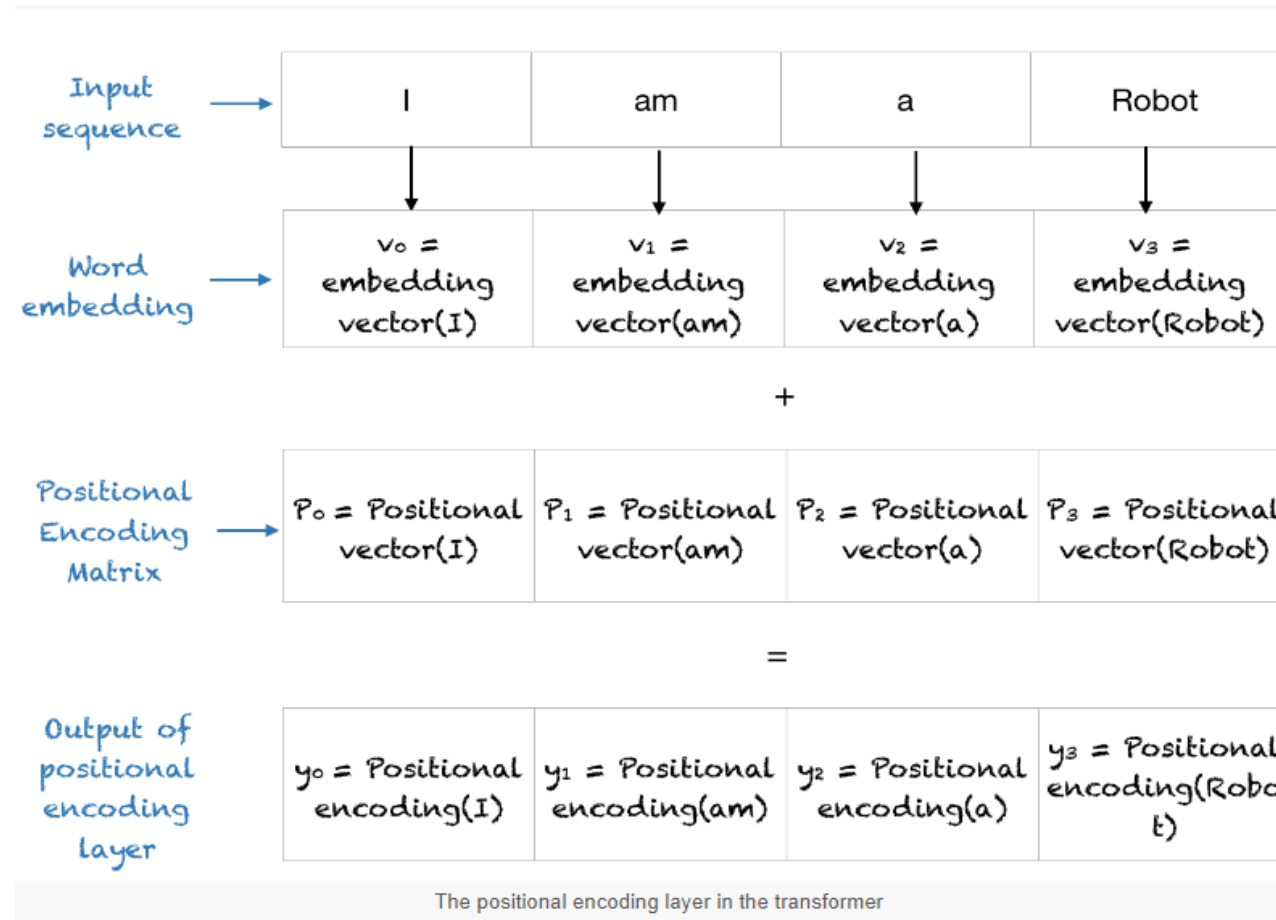
Positional Encoding Matrix with $d=4$, $n=100$

| Sequence | Index of token, k | $i=0$ | $i=0$ | $i=1$ | $i=1$ |
|----------|---------------------|------------------------------|-------------------------------|-------------------------------|-------------------------------|
| 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 |

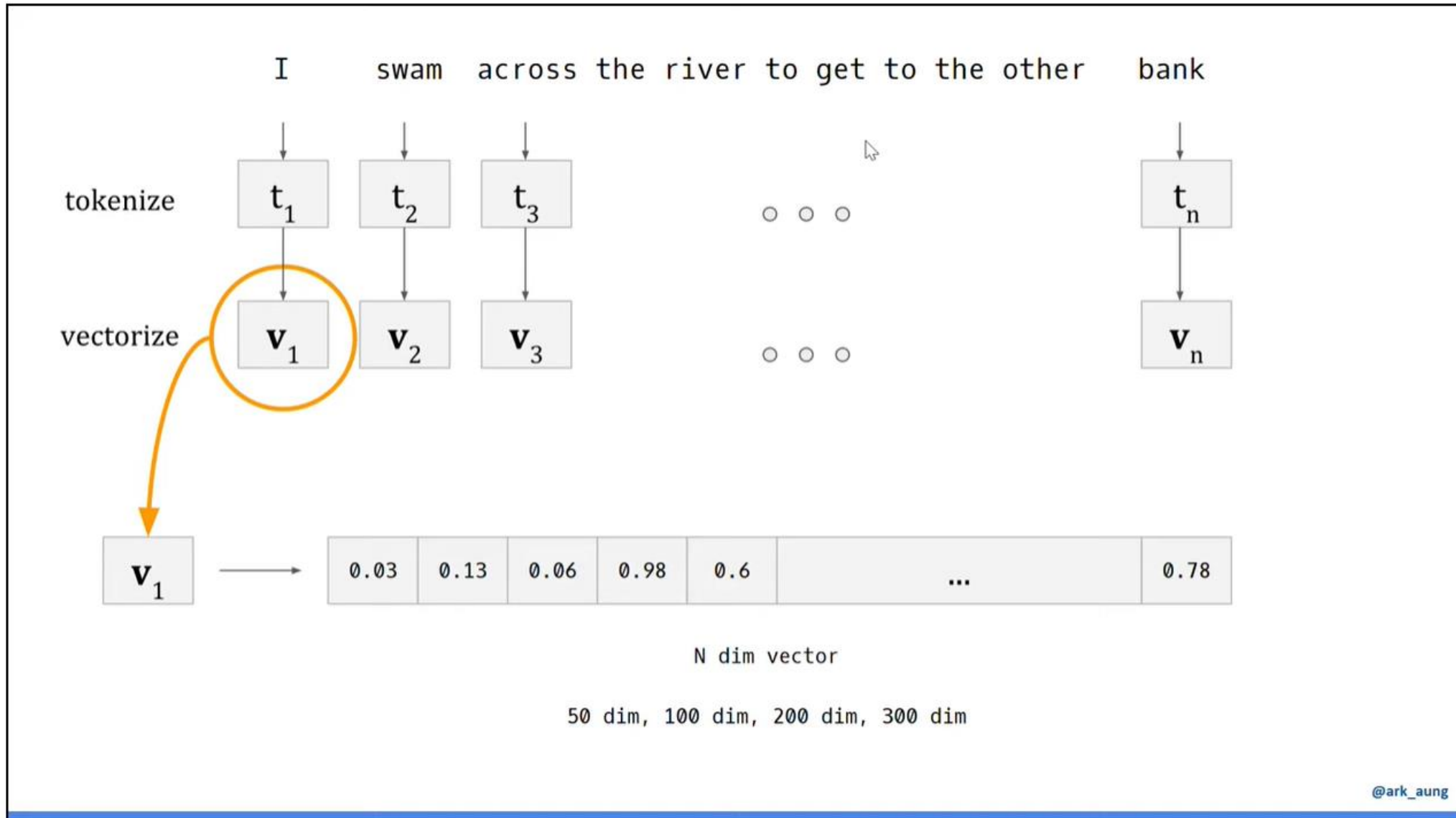
Positional Encoding Matrix for the sequence 'I am a robot'

Positional Encoding

- Final output of the Positional Encoding Layer

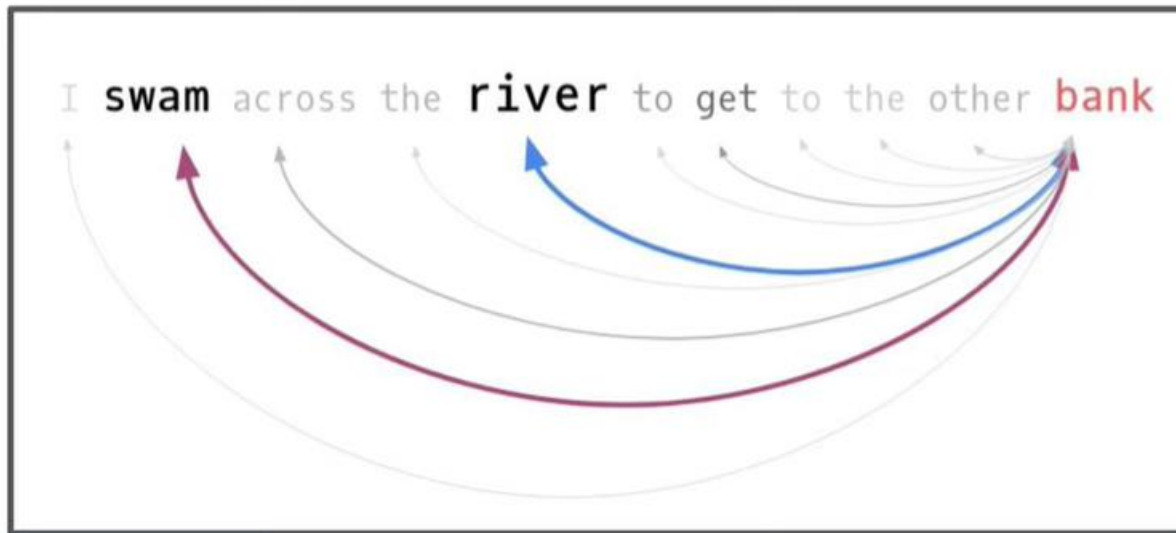
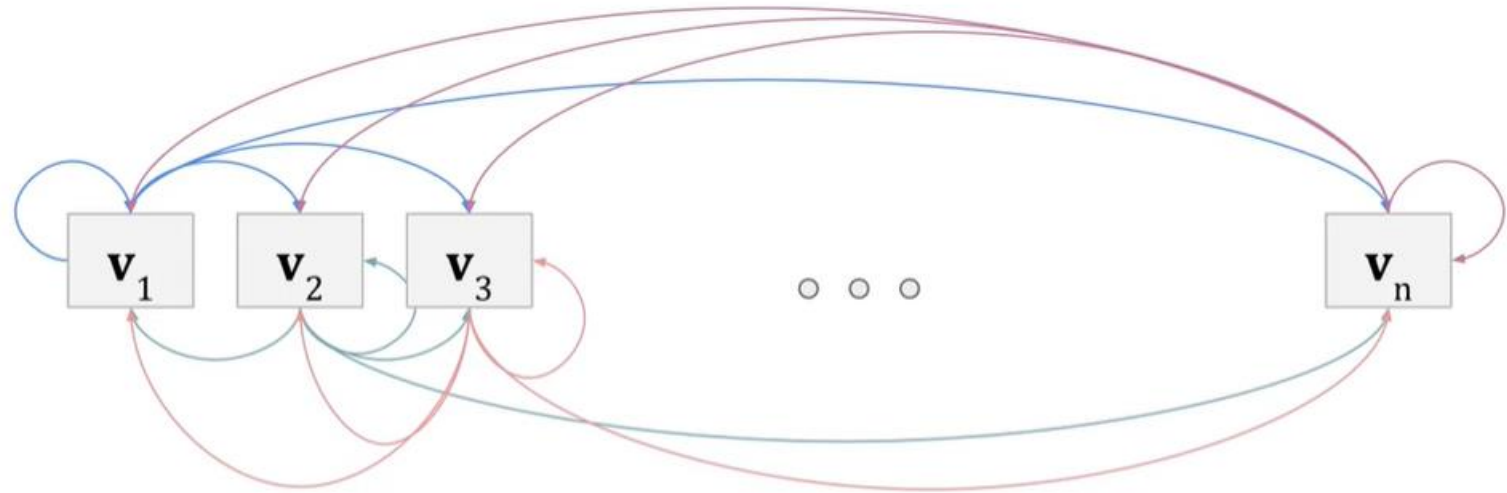


Self-Attention



Self-Attention

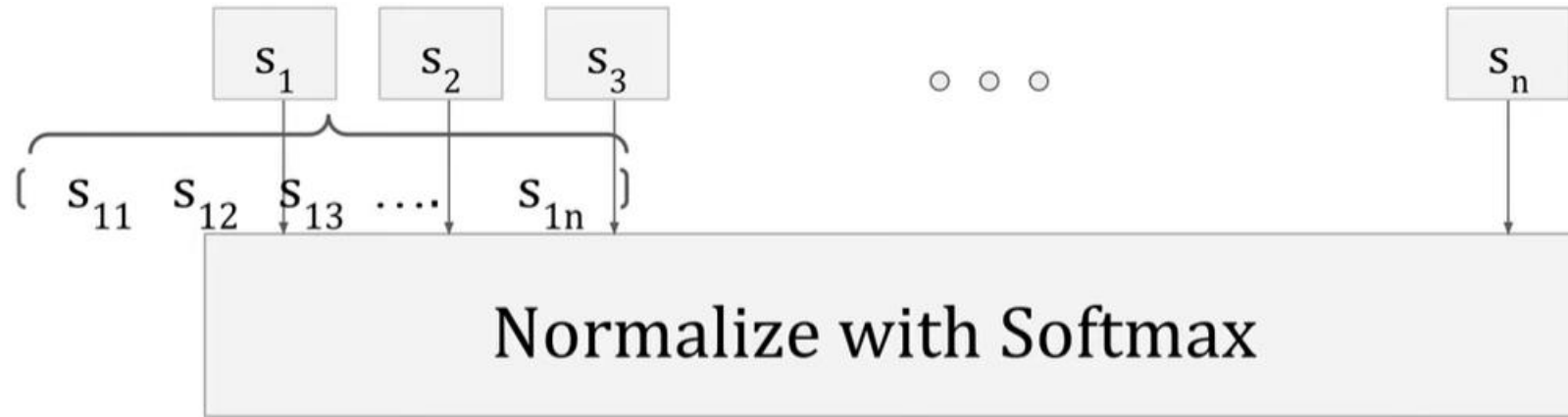
- Similarity
 - Dot product
 - Cosine similarity



$$\begin{aligned} S_{n1} &= \mathbf{v}_n \mathbf{v}_1^T \\ S_{n2} &= \mathbf{v}_n \mathbf{v}_2^T \\ S_{n3} &= \mathbf{v}_n \mathbf{v}_3^T \end{aligned}$$

$$S_{nn} = \mathbf{v}_n \mathbf{v}_n^T$$

Self-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

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

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

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

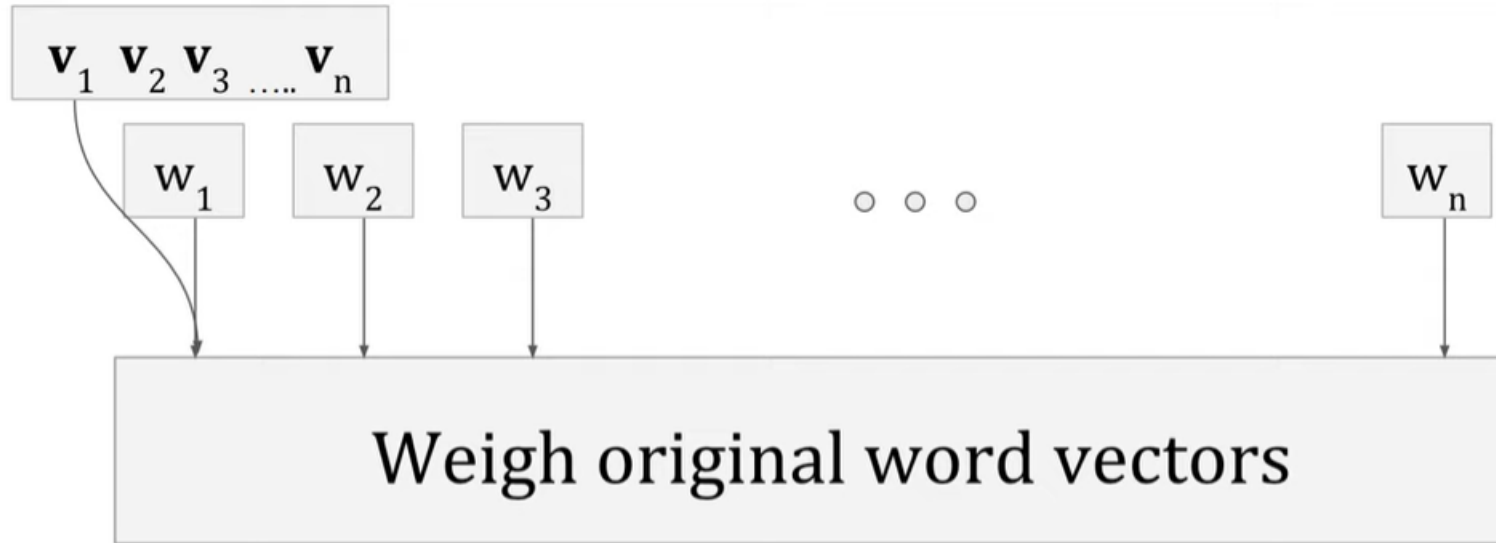
$$\begin{aligned} w_{n1} &= \text{softmax}(s_{n1}) \\ w_{n2} &= \text{softmax}(s_{n2}) \\ w_{n3} &= \text{softmax}(s_{n3}) \end{aligned}$$

0
0
0

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

Self-Attention

3



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

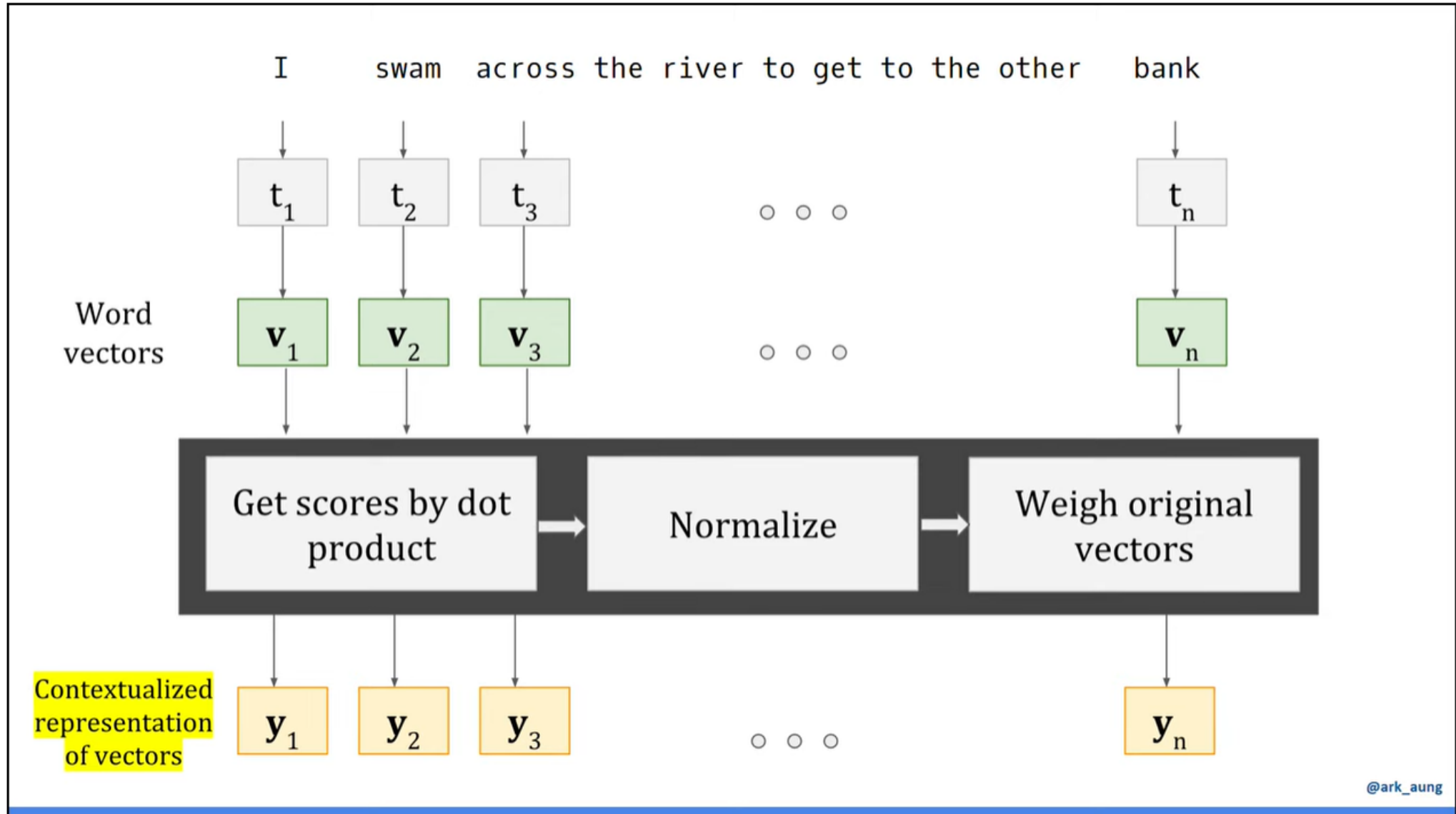
$$\mathbf{y}_2 = w_{21} \mathbf{v}_1 + w_{22} \mathbf{v}_2 + w_{23} \mathbf{v}_3 + \dots + w_{2n} \mathbf{v}_n$$

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

0
0
0

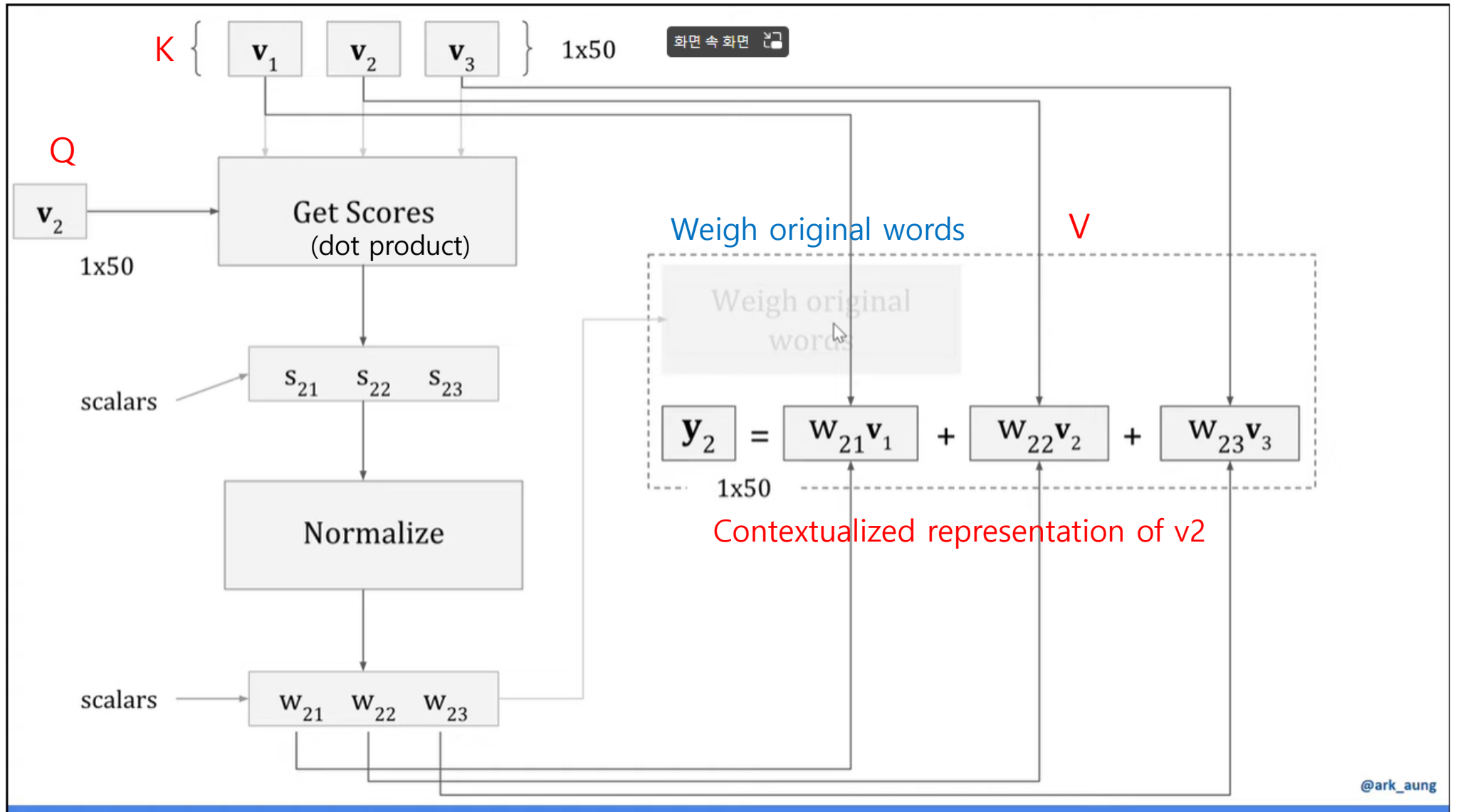
$$\mathbf{y}_n = w_{n1} \mathbf{v}_1 + w_{n2} \mathbf{v}_2 + w_{n3} \mathbf{v}_3 + \dots + w_{nn} \mathbf{v}_n$$

Self-Attention



Self-Attention mechanism

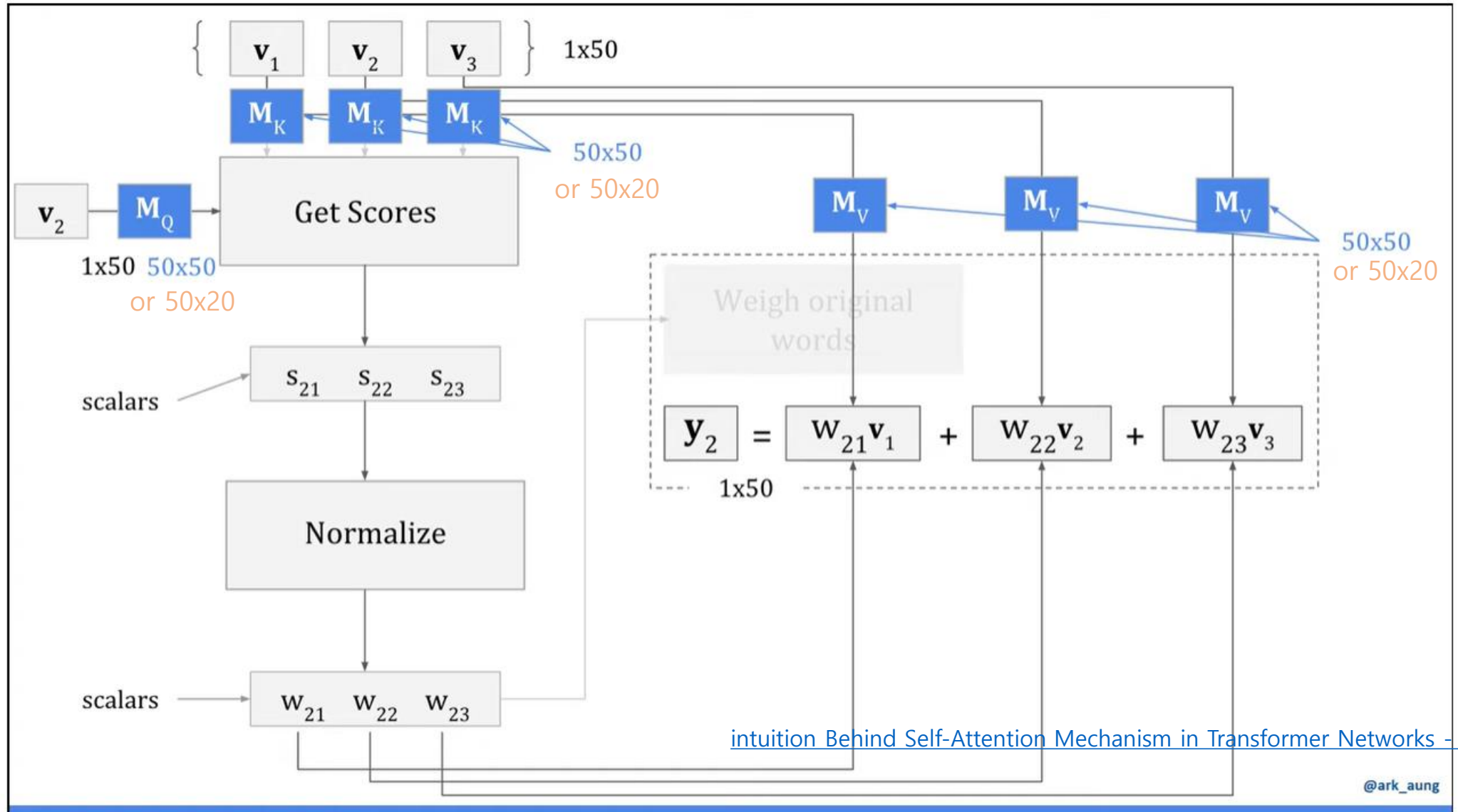
- Example
 - 3 words
 - 50-d vectors



Self-Attention

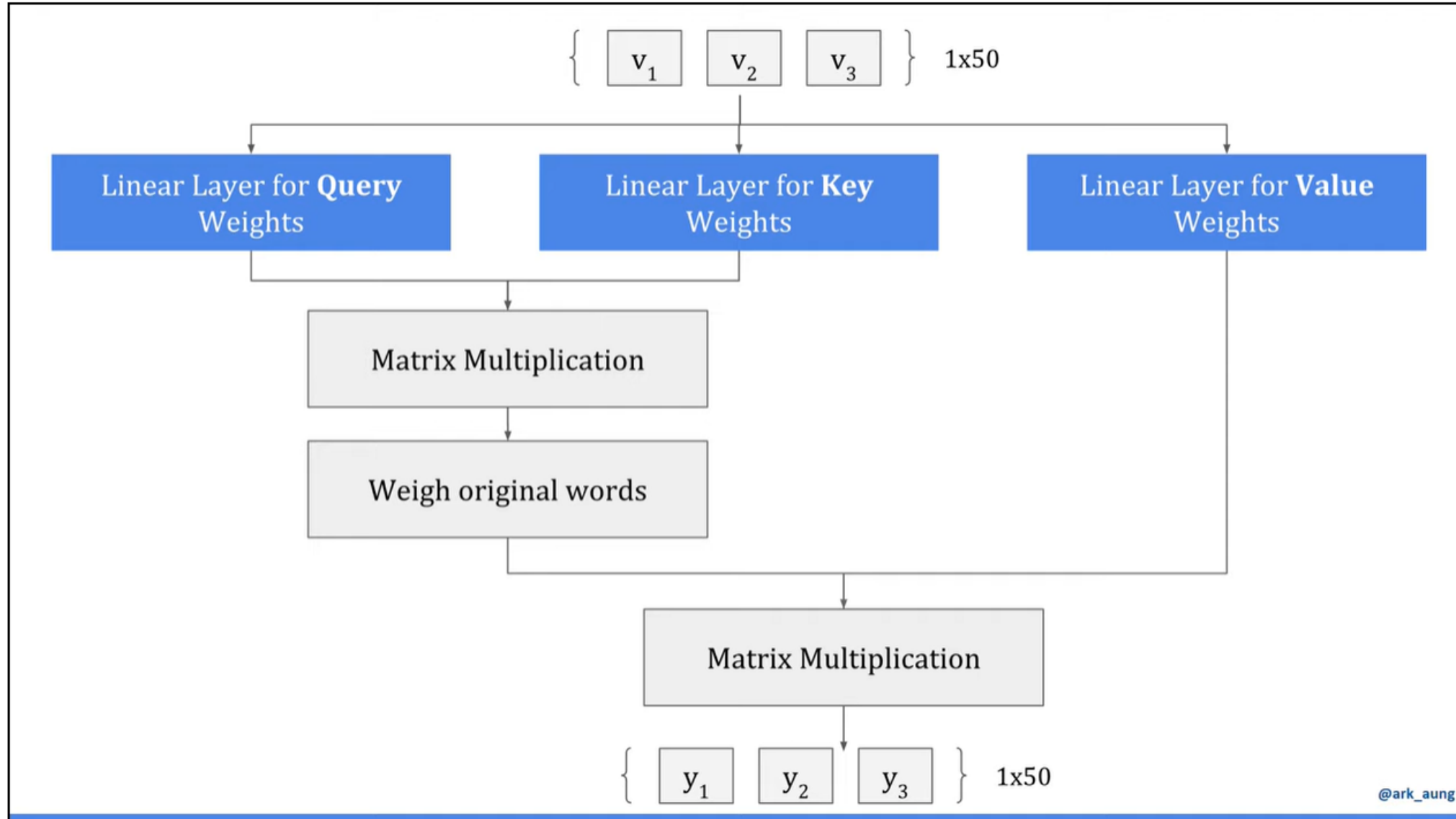
- Injecting learnable parameters

(*) may have different size

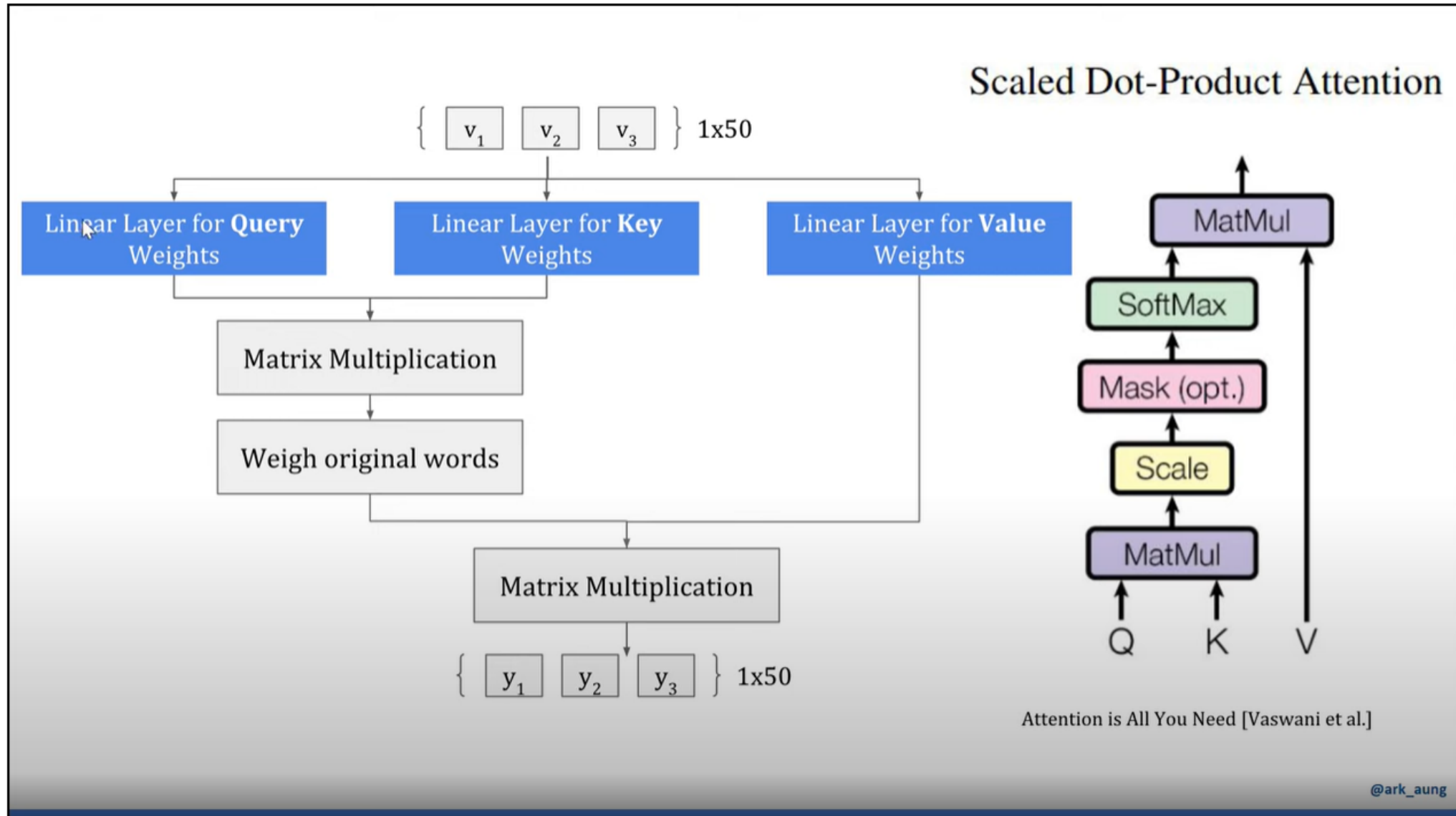


[intuition Behind Self-Attention Mechanism in Transformer Networks - YouTube](#)

Self-Attention

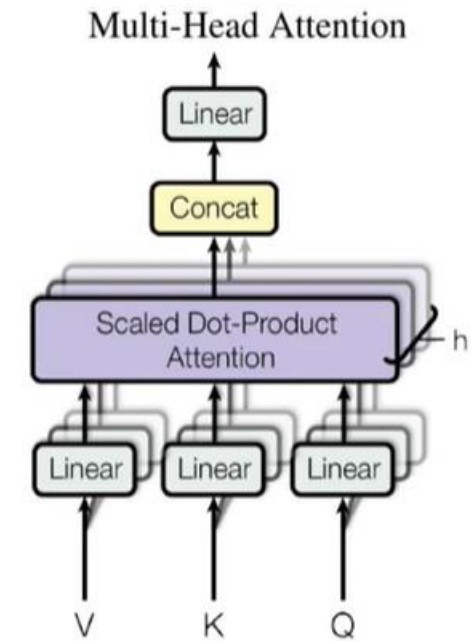
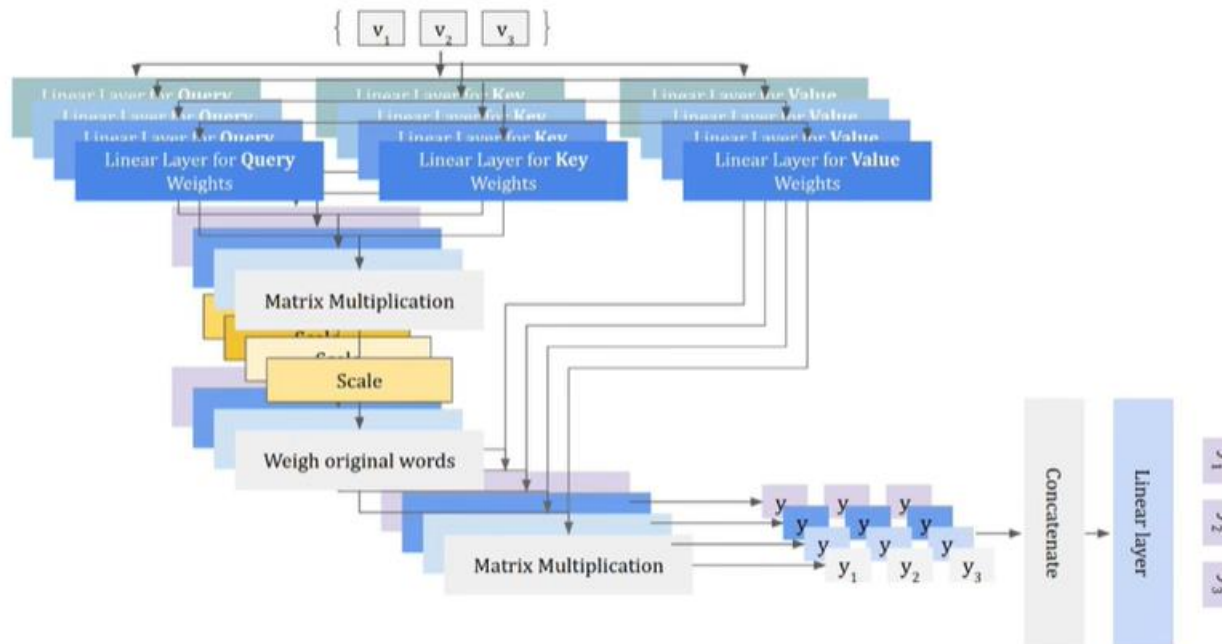


Self-Attention



Self-Attention

- Multi-head attention
 - Diversity



Attention is All You Need [Vaswani et al.]

@ark_aung

Transformer

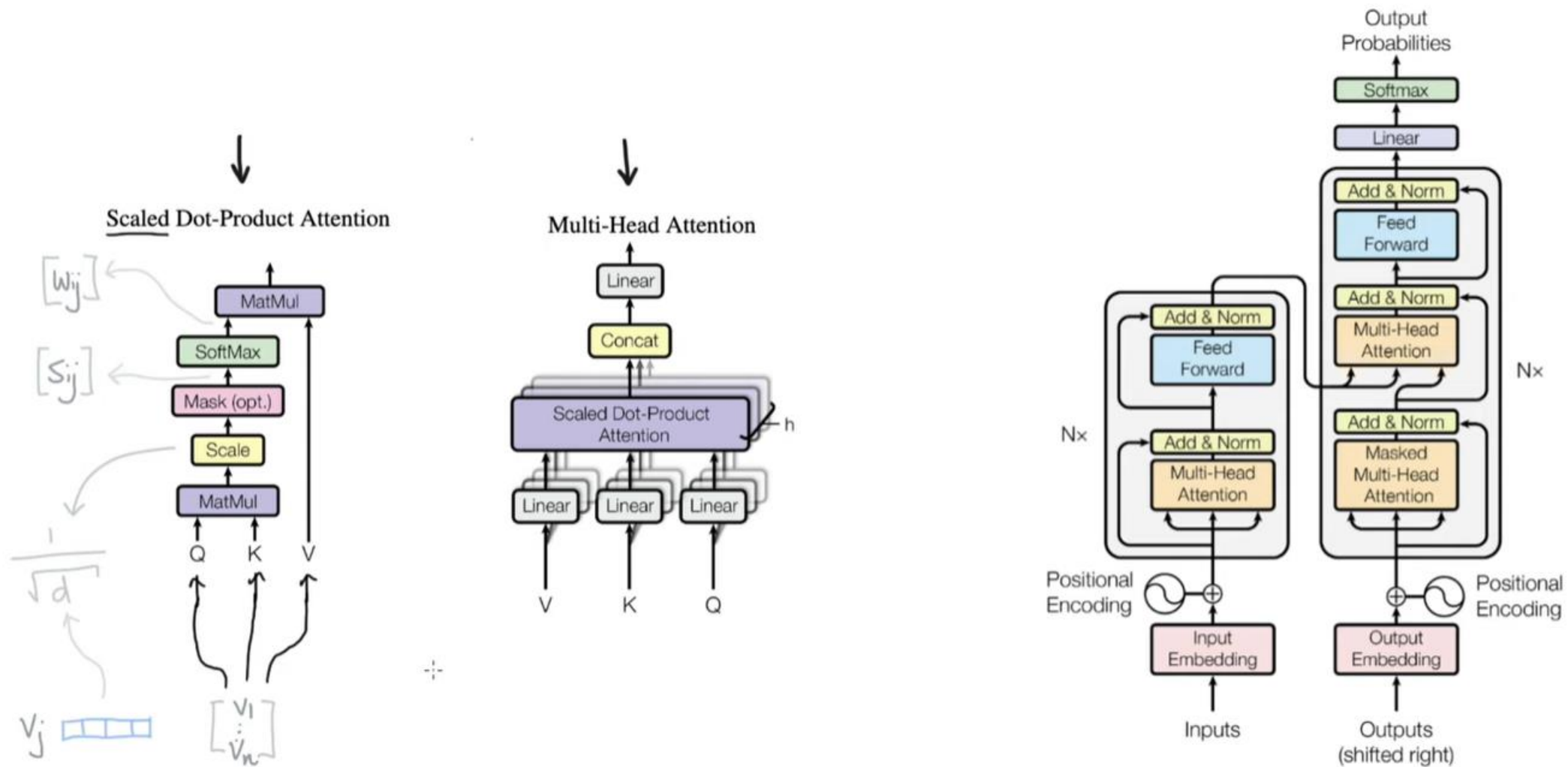


Figure 1: The Transformer - model architecture.