

Attention & Self-Attention

Zhongsheng Wang
University of Auckland

Mar 15, 2024



THE UNIVERSITY OF
AUCKLAND
Te Whare Wānanga o Tāmaki Makaurau
NEW ZEALAND

Paper link: <https://arxiv.org/pdf/1706.03762.pdf>

Table of Contents

- Problems with RNN/LSTM
- Attention
 - Attention is all you need
- Self-attention
 - Better understanding of $\text{Softmax}(XX^T)X$
 - Scaled dot-product attention
 - Size of each variable
 - Calculation demo
 - Position encoding
- Multi-head Attention
- Interview questions

Problems with RNN/LSTM

- Memory length is limited (RNN) Long-term memory
- Forgot key information (LSTM)
- Cannot be parallelized (Calculate t_0 , then t_1 , $t_2 \dots$)
 - Training efficiency is relatively low

Attention



- What do you see **at first glance** in the picture on the left?
- **Human:** Limited attention span, focusing on one thing in a short period (finding the highest priority thing to do among the mess of things)
- **Computer:** Limited computing resources, and we hope it can efficiently capture critical information from a large amount of information.

Attention is all you need

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
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

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

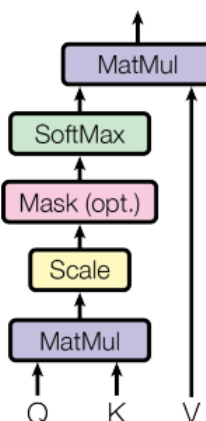
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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

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

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

Scaled Dot-Product Attention



Multi-Head Attention

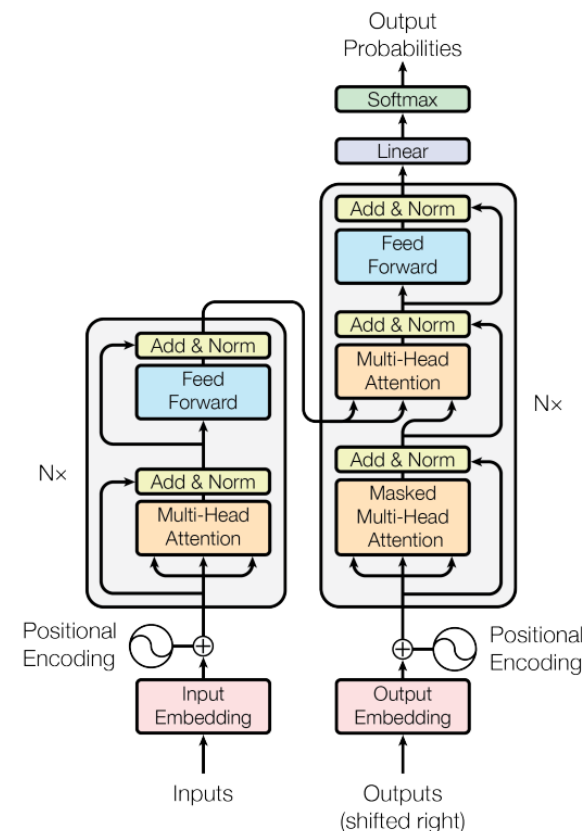
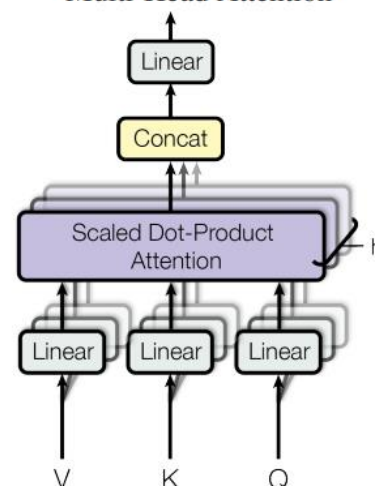


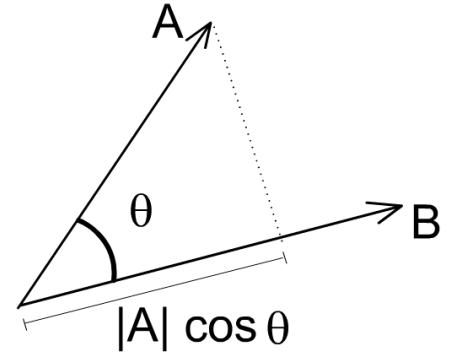
Figure 1: The Transformer - model architecture.

Self-attention

- **Input:** A sentence with several words. **Output:** Word vector for each word.
- Create Matrix QKV from original input x
- Q: Query
 - Will match each k
- K: Key
 - Will be matched by every q
- V: Value
 - Extract value Information from x
- Judge the similarity using QK and rewrite the representation of a word vector using V (match the most similar result)

Better understanding of $\text{Softmax}(XX^T)X$

- Dot-product: Similarity between 2 word vectors
- Word vec1: $[1, 3, 2]^T$ Word vec2: $[4, 1, 1]^T$ (Column vector)
- Dot-product = $1*4 + 3*1 + 2*1 = 9$
- Self dot-product for vec1 = $1*1 + 3*3 + 2*2 = 14$
- In matrix, dot-product $\Rightarrow \text{vec} * \text{vec}^T$ (Matrix transpose) = $a_1*a_1 + a_2*a_2 + \dots$



Cont.

- If we combine all word vectors:

- $X = [x_1, x_2, x_3] = \begin{bmatrix} x_{11} & x_{21} & x_{31} \\ x_{12} & x_{22} & x_{32} \\ x_{13} & x_{23} & x_{33} \end{bmatrix}$ (word vector dimension is 3)

- $XX^T = [x_1, x_2, x_3] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} x_{11} & x_{21} & x_{31} \\ x_{12} & x_{22} & x_{32} \\ x_{13} & x_{23} & x_{33} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix}$

- Calculate the similarity between two pairs (including themselves)

Cont.

$$\text{Softmax}(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

- $\text{Softmax}(XX^T)X$: convert the dot-product weight value back to the presentation of an average word vector
- $\text{Softmax}(XX^T)$: weight value X : word vec representation
- So the result of $\text{Softmax}(XX^T)X$ is a new word vector representation that is **weighted and summed by the attention mechanism**

Scaled Dot-Product Attention

$$\text{Softmax}(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

Key

$$K = W_k X$$

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

$$Q = W_q X$$

Query

Value

$$K = W_k X$$

$$Q = W_q X$$

$$V = W_v X$$

$$V = W_v X$$

$$A = K^T Q$$

Dimension of K

$$A' = \text{softmax}(A)$$

$$Y = VA'$$

Size of each variable

- Size of X : (a, b) **column vector**

$$Q = W_q X$$

- a : the dimension of a word vector

$$K = W_k X$$

- b : number of words in a corpus(a sentence)

$$V = W_v X$$

$$A = K^T Q$$

- Size of W_{qkv} : $(c, a) \Rightarrow$ Size of QKV: (c, b)

$$A' = \text{softmax}(A)$$

- $W = QX^T$ (Random Start)

$$Y = VA'$$

- Size of A : $(b, b) = (b, c) * (c, b)$

Cont.

$$\begin{aligned}Q &= W_q X \\K &= W_k X \\V &= W_v X \\A &= K^T Q \\A' &= \text{softmax}(A) \\Y &= V A'\end{aligned}$$

- $A: [a_1, a_2, a_3, \dots, a_b]$ $a_1: [a_{11}, a_{12}, \dots, a_{1b}]$
- Softmax for each column
 - $A: [a'_1, a'_2, a'_3, \dots, a'_b]$
 - $a_1: [a'_{11}, a'_{12}, \dots, a'_{1b}]$ (b: number of words in a sentence)
- $Y: (c, b) \Rightarrow (c, b) * (b, b)$
- So the output is still **the representation of each word**
- c is a hyperparameter (set whatever u want)

Calculation demo

- Suppose:
 - 4 words in a sentence (x_1, x_2, x_3), dimension of the word vector is 4
 - Size of X : $(a, b) \Rightarrow a=4(\text{dimension})$ $b=3(\text{num of words})$
 - For each x in X : size of x : $(a, 1) \Rightarrow (4, 1)$
 - $x_1 = [1, 2, 3, 4]^T$
 - $x_2 = [3, 2, 1, 3]^T$
 - $x_3 = [0, 1, 1, 4]^T$
 - All of them are column vectors

Cont.

- Random start W_q W_k W_v : (Size of W : (c, a) suppose c is also 4) $\Rightarrow (4, 4)$ so d_k is also 4

$$\bullet W_q = \begin{bmatrix} 3 & 1 & 1 & 3 \\ 3 & 5 & 0 & 0 \\ 4 & 7 & 1 & 1 \\ 2 & 1 & 2 & 3 \end{bmatrix} \quad W_k = \begin{bmatrix} 2 & 2 & 0 & 1 \\ 4 & 2 & 3 & 1 \\ 0 & 2 & 4 & 1 \\ 1 & 4 & 2 & 1 \end{bmatrix} \quad W_v = \begin{bmatrix} 1 & 3 & 2 & 4 \\ 2 & 0 & 1 & 1 \\ 3 & 1 & 3 & 4 \\ 2 & 1 & 1 & 0 \end{bmatrix}$$

Cont.

- For x_1 , calculate q_1, k_1, v_1 :

$$\bullet q_1 = W_X * x_1 = \begin{bmatrix} 2 & 2 & 0 & 1 \\ 4 & 2 & 3 & 1 \\ 0 & 2 & 4 & 1 \\ 1 & 4 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 2 * 1 + 2 * 2 + 1 * 4 \\ 4 * 1 + 2 * 2 + 3 * 3 + 1 * 4 \\ 2 * 2 + 3 * 4 + 1 * 4 \\ 1 * 1 + 4 * 2 + 2 * 3 + 1 * 4 \end{bmatrix} = \begin{bmatrix} 10 \\ 18 \\ 15 \\ 19 \end{bmatrix}$$

- $Q = [q_1, q_2, q_3, q_4] = W_X * [x_1, x_2, x_3, x_4] = W_X * X$ (size:(4, 4))
- Each col represents a query vector of a word
- Same for x_2, x_3 and x_4

Cont.

- To calculate the similarity of x_1 and x_2 :
• Use $q_1 * k_2$ and then Softmax

$$q_1 = \begin{bmatrix} 10 \\ 18 \\ 15 \\ 19 \end{bmatrix}$$

$$\bullet \quad k_2 = Wk * x_2 = \begin{bmatrix} 2 & 2 & 0 & 1 \\ 4 & 2 & 3 & 1 \\ 0 & 2 & 4 & 1 \\ 1 & 4 & 2 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 1 \\ 3 \end{bmatrix} = \begin{bmatrix} 2 * 3 + 2 * 2 + 1 * 3 \\ 4 * 3 + 2 * 2 + 1 * 3 + 1 * 3 \\ 2 * 2 + 1 * 4 + 1 * 3 \\ 1 * 3 + 4 * 2 + 2 * 1 + 1 * 3 \end{bmatrix} = \begin{bmatrix} 13 \\ 22 \\ 11 \\ 16 \end{bmatrix}$$

$$\bullet \quad \text{So } \alpha_{12} = \begin{bmatrix} 10 \\ 18 \\ 15 \\ 19 \end{bmatrix} * \begin{bmatrix} 13 \\ 22 \\ 11 \\ 16 \end{bmatrix} = 10*13 + 18*22 + 15*11 + 19*16 = 995$$

Cont.

- Softmax after calculating all the similarity of x1 with others (contain itself) ($\alpha_{11} \alpha_{12} \alpha_{13} \alpha_{14}$)
- $\alpha_{11} = \alpha_{11} / \sqrt{d_k}$ (according to the formula)
 - Value after dot-product becomes very large, and the gradient disappears after softmax
- $\alpha'_{11} = \frac{e^{\alpha_{11}}}{e^{\alpha_{11}} + e^{\alpha_{12}} + e^{\alpha_{13}} + e^{\alpha_{14}}}$ (Softmax)
- Same way to calculate $\alpha'_{11} \alpha'_{12} \alpha'_{13} \alpha'_{14}$
- $y1 = [v1, v2, v3, v4] * [\alpha'_{11}, \alpha'_{12}, \alpha'_{13}, \alpha'_{14}]^T$

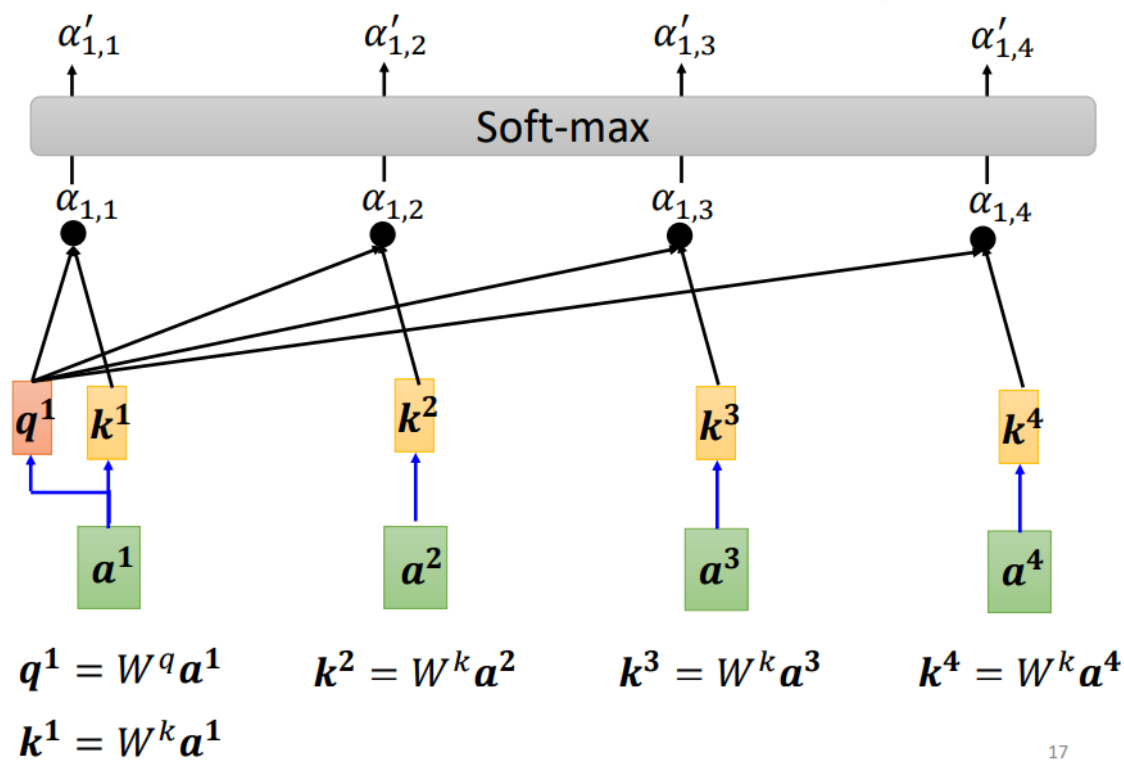
Cont.

- $y1 = [v1, v2, v3, v4] * [\alpha'_{11}, \alpha'_{12}, \alpha'_{13}, \alpha'_{14}]^T$
- Size of v1: (4, 1) V: (4, 4) y1: (4, 1)*(4, 1)[dot-product] = (4, 1)
- So y1 has the same size as x1, we think this is a new word vector that contains all structural information

Cont.

Self-attention

$$\alpha'_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$

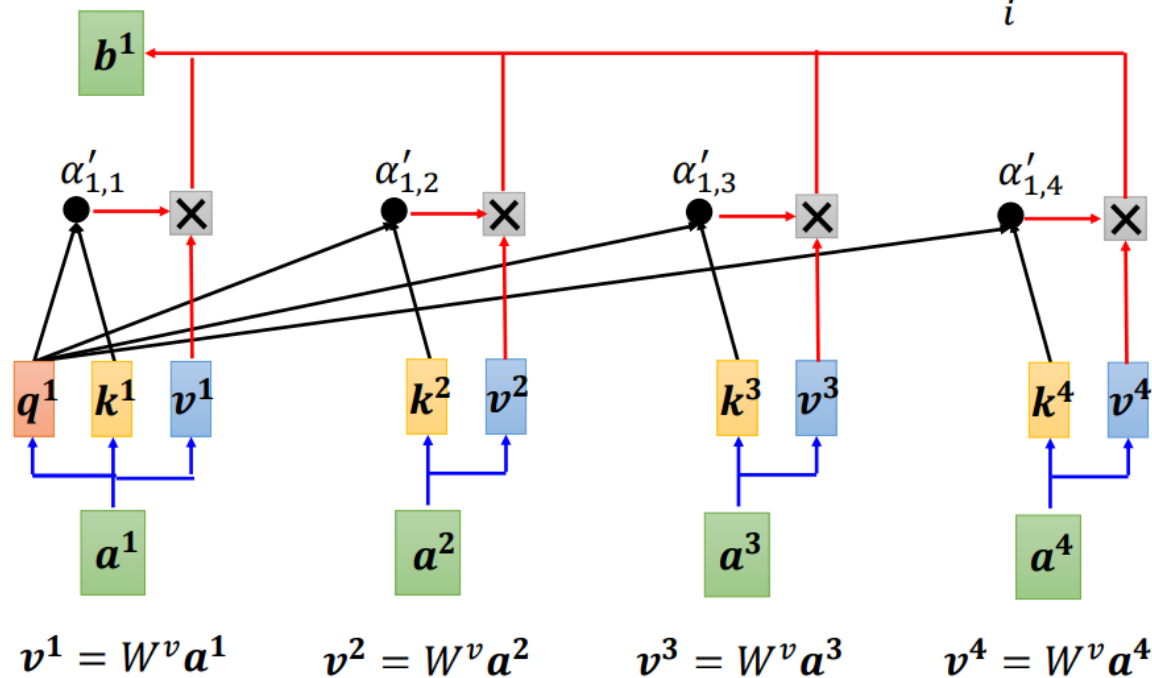


17

Self-attention

Extract information based on attention scores

$$b^1 = \sum_i \alpha'_{1,i} v^i$$



18

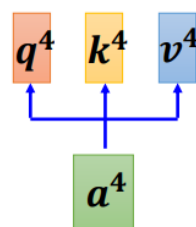
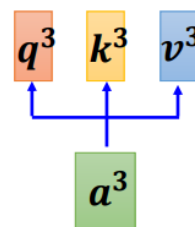
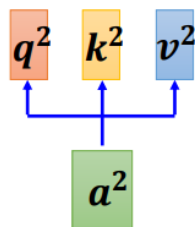
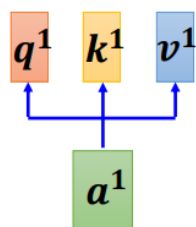
Cont.

Self-attention

$$q^i = W^q a^i \quad \underbrace{q^1 q^2 q^3 q^4}_Q = \underbrace{W^q}_I \underbrace{a^1 a^2 a^3 a^4}_I$$

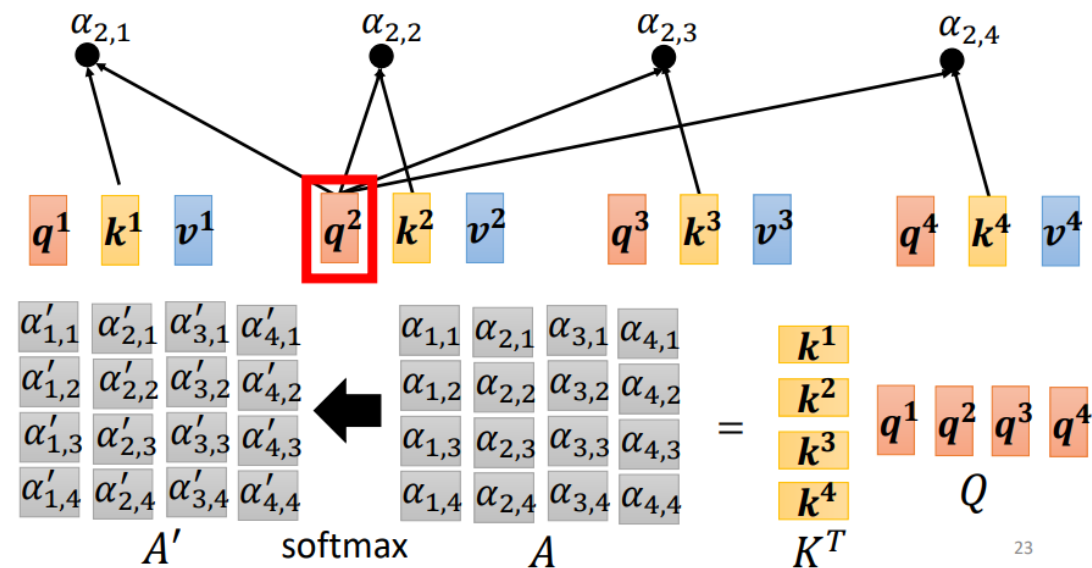
$$k^i = W^k a^i \quad \underbrace{k^1 k^2 k^3 k^4}_K = \underbrace{W^k}_I \underbrace{a^1 a^2 a^3 a^4}_I$$

$$v^i = W^v a^i \quad \underbrace{v^1 v^2 v^3 v^4}_V = \underbrace{W^v}_I \underbrace{a^1 a^2 a^3 a^4}_I$$



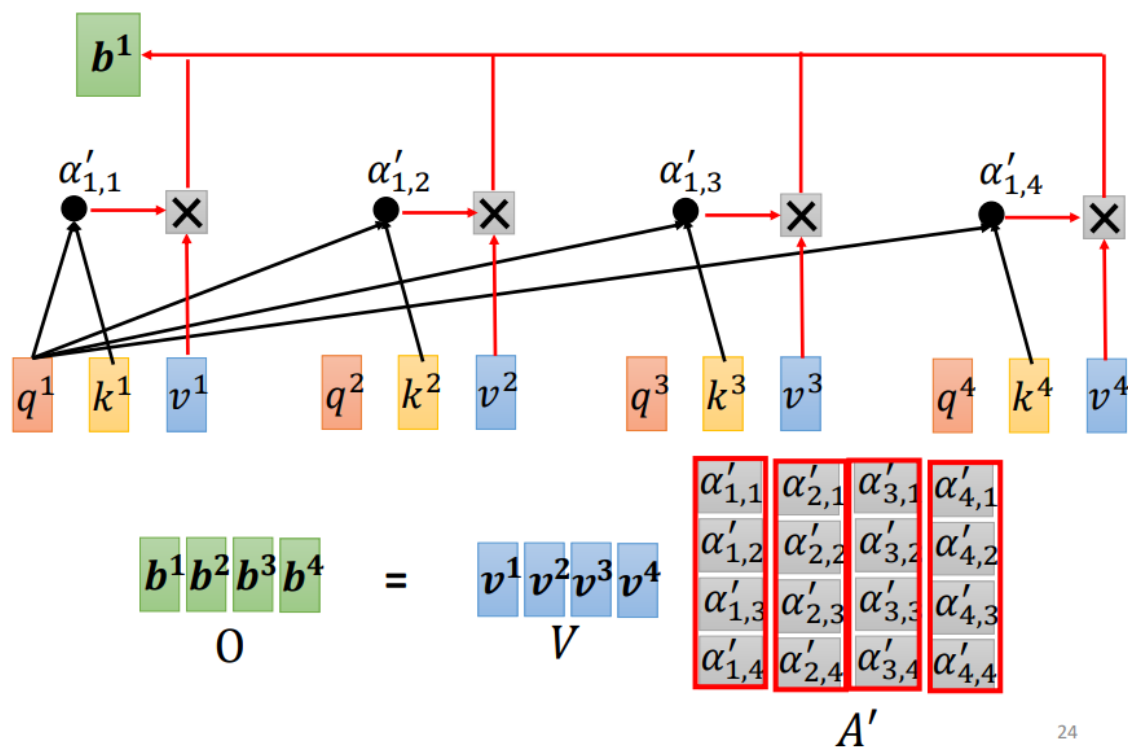
Self-attention

$$\begin{array}{l} \alpha_{1,1} = k^1 q^1 \\ \alpha_{1,3} = k^3 q^1 \end{array} \quad \begin{array}{l} \alpha_{1,2} = k^2 q^1 \\ \alpha_{1,4} = k^4 q^1 \end{array} \quad \begin{array}{c} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{array} = \begin{array}{c} k^1 \\ k^2 \\ k^3 \\ k^4 \end{array} q$$



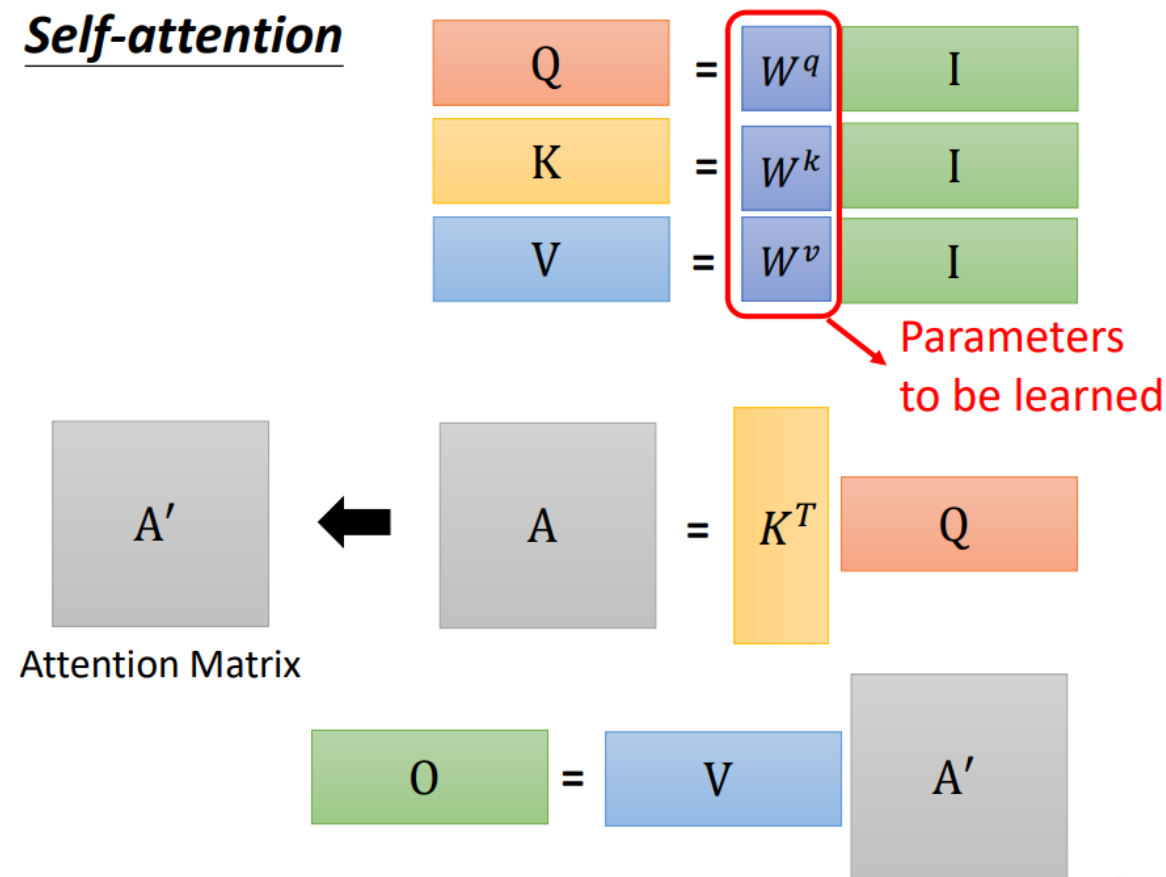
Cont.

Self-attention



24

Self-attention



25

Position encoding

- Self-attention only considers the content of the context and does not consider the relationship between text positions
- “I love you” and “you love I” have the same output using Self-attention
- Solution:
 - 1. add a position embedding to the original input x
 - 2. Using a trainable position embedding parameter

Cont.

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

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

- Add position embedding to the sentence “I love you”
- Suppose: variable “pos” starts from 0, dmodel = 4 (that is, dimension of the word vector)

word	Dimension 1	Dimension 2	Dimension 3	Dimension 4
I	0	0	0	0
love	$\sin(1/10000^{2/4})$	$\cos(1/10000^{2/4})$		
you				

- For word “I” :
 - pos = 0
 - i = 0: $PE(0, 0) = \sin(0/\text{whatever}) = 0$; $PE(0, 0) = 0$
 - i = 1: $PE(0, 2) = 0$; $PE(0, 3) = 0$

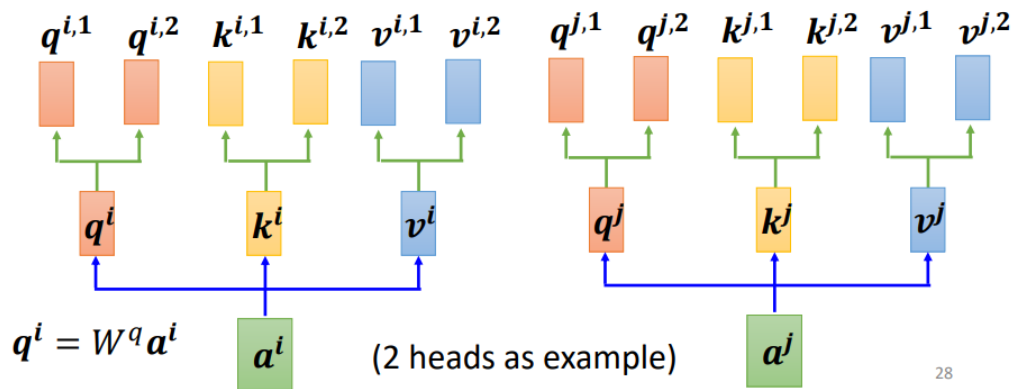
Multi-head Attention

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

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

Multi-head Self-attention Different types of relevance

$$b^i = W^O \begin{bmatrix} b^{i,1} \\ b^{i,2} \end{bmatrix}$$



- Suppose:
- Size $X = (4, 2)$ $[x_1, x_2]^T$
- Size $QKV = (2, 2)$ $W = (2, 4)$ Dimension reduction
- If num of head $\Rightarrow 2$
- For q_1 split into 2 sub-matrix q_{11} and q_{12}
- $q_{11}: (2, 1)$ $q_{12}: (2, 1)$ $k_{11}: (2, 1)$ $k_{12}: (2, 1)$
- $a_{11} = k_{11}^T * q_{11} = (1, 2) * (2, 1) = (1, 1)$ then Softmax
- $b_{11} = v_{11} \cdot a'_{11} = (2, 1) \cdot (1, 1) = (2, 1)$
- Combine b_{11} b_{12} we get $\Rightarrow (4, 1)$
 - Also for b_{21} and b_{22} we get $(4, 1)$ add, then $(4, 2)$
- Optimize through a fully connected (FC) layer

Interview questions

- What is the core of Self-Attention?
- Why does Transformer require additional positional encoding?
- What exactly is QKV matrix in Self-attention? Why do they exist?
- Why does Transformer need Multi-head Attention? What are its benefits?
- Why must the dot product model be scaled before Softmax normalization?
- Why does each head in multi-head attention in Transformer need to be dimensionally reduced?