# Tutorial 5

Fei Gao & Yuliia Frund

# Word Embeddings

# Word Embedding Types

- Static
  - (e.g.: Word2Vec, GloVe)

- Contextual
  - (e.g.: ELMo, BERT)

# Skip-gram embeddings

Train a model that predicts context words:
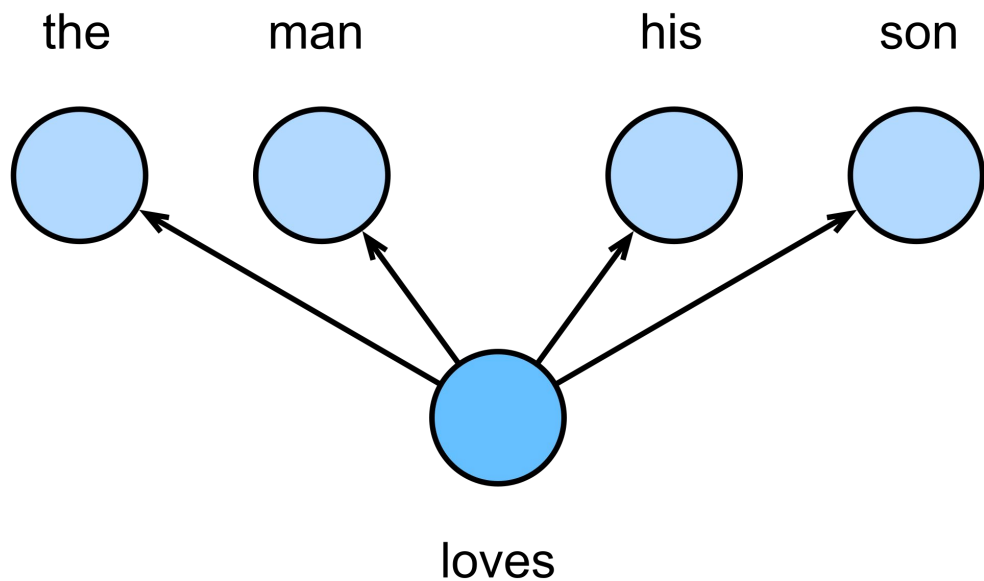
[CONTEXT TARGET CONTEXT]



Image: https://en.wikipedia.org/wiki/Word2vec

# Skip-gram embeddings

Train a model that predicts context words:

[CONTEXT TARGET CONTEXT]

These are **static** embeddings!

Image: https://en.wikipedia.org/wiki/Word2vec
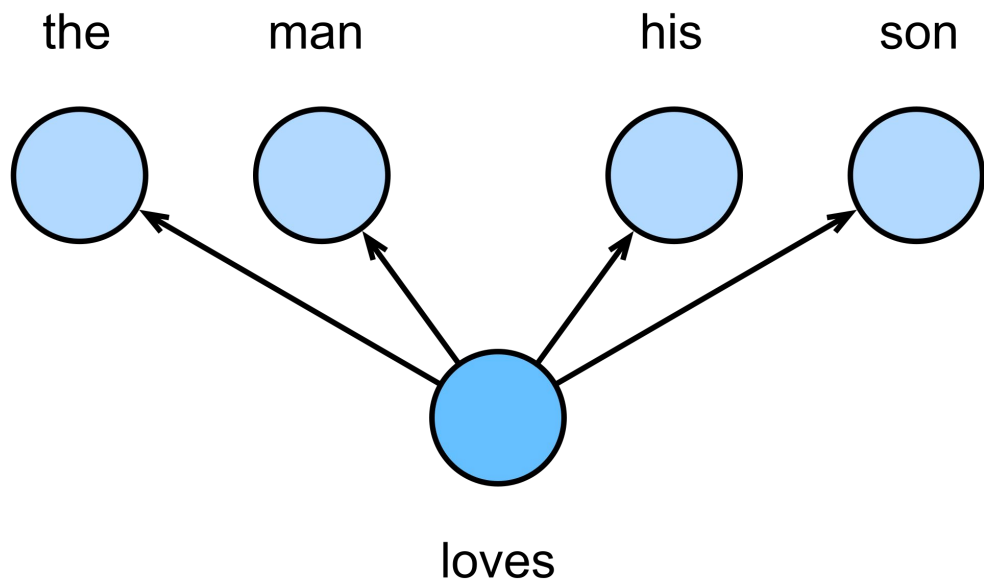
the        man        his        son

loves

# Skip-gram embeddings

In other words...

Suppose this is our training data:

The quick brown fox jumps over the lazy dog.

Target word: FOX

A model with a context size 1 will predict:

FOX, brown

FOX, jumps

# Byte Pair Encoding

# Lexicon

- Let's imagine that you're training a machine translation system.
- Will it see **ALL** possible words in target and source language?
  - No! There will always be words your system hasn't seen!
- How do you deal with **unseen words**?
  - Subword units!

# Byte Pair Encoding



**function** BYTE-PAIR ENCODING(strings $C$, number of merges $k$) **returns** vocab $V$

$V \leftarrow$ all unique characters in $C$      # initial set of tokens is characters
**for** $i = 1$ **to** $k$ **do**      # merge tokens $k$ times
    $t_L, t_R \leftarrow$ Most frequent pair of adjacent tokens in $C$
    $t_{NEW} \leftarrow t_L + t_R$      # make new token by concatenating
    $V \leftarrow V + t_{NEW}$      # update the vocabulary
    Replace each occurrence of $t_L, t_R$ in $C$ with $t_{NEW}$      # and update the corpus
**return** $V$

**Figure 2.13**     The token learner part of the BPE algorithm for taking a corpus broken up into individual characters or bytes, and learning a vocabulary by iteratively merging tokens. Figure adapted from Bostrom and Durrett (2020).

Sennrich, Rico; Birch, Alexandra; Haddow, Barry (2015-08-31). "Neural Machine Translation of Rare Words with Subword Units"
https://web.stanford.edu/~jurafsky/slp3/2.pdf

# Byte Pair Encoding

In practice:

- Very frequent words are likely to be stored whole
- Rare and unseen words can still be handled
- Manageable vocabulary size

# Attention

# Class Slide 15

## Attention Head (contd.)

- Now we have projected the inputs with three transformations

- We use the query and the key to compute attention

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}} \quad \text{(9.11)}$$

$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \ \forall j \le i \quad \text{(9.12)}$$

$$\mathbf{a}_i = \sum_{j \le i} \alpha_{ij} \mathbf{v}_j \quad \text{(9.13)}$$

**Causal attention**

Because we are only looking at the past

**Full attention**

Looks at both past and the future

# Attention Head (contd.)

## Calculate α3

- Now we have projected the inputs with three transformations

- We use the query and the key to compute attention

How to calculate the weight row $\alpha_3$?

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}} \quad \text{(9.11)}$$

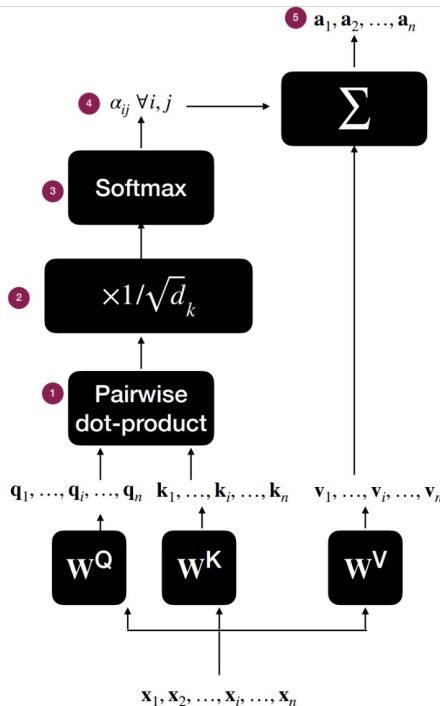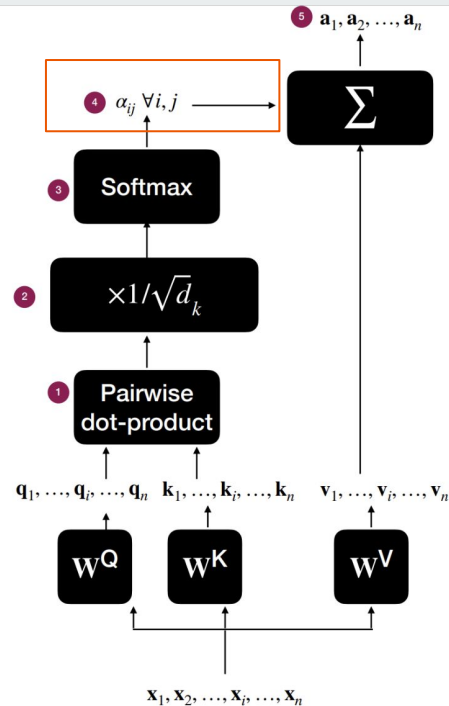$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \ \forall j \leq i \quad \text{(9.12)}$$

$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j \quad \text{(9.13)}$$

$\alpha_3$ is a row. Think: What is its individual components?

Hint: what are j's that are smaller than i (i = 3, the query number), starting from 1?

**Causal attention**

Because we are only looking at the past

**Full attention**

Looks at both past and the future



$\mathbf{a}_1, \mathbf{a}_2, \ldots, \mathbf{a}_n$

⑤

④ $\alpha_{ij} \ \forall i, j$ → $\Sigma$

③ Softmax

② $\times 1/\sqrt{d}_k$

① Pairwise dot-product

$\mathbf{q}_1, \ldots, \mathbf{q}_i, \ldots, \mathbf{q}_n$ $\mathbf{k}_1, \ldots, \mathbf{k}_i, \ldots, \mathbf{k}_n$ $\mathbf{v}_1, \ldots, \mathbf{v}_i, \ldots, \mathbf{v}_n$

$\mathbf{W}^Q$ $\mathbf{W}^K$ $\mathbf{W}^V$

$\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_i, \ldots, \mathbf{x}_n$

# Calculate α3

How to calculate the weight row

$\alpha_3$?

$\alpha_3$ is a combination of

$\alpha_{31}$, $\alpha_{32}$, and $\alpha_{33}$

# Attention Head (contd.)

- Now we have projected the inputs with three transformations

- We use the query and the key to compute attention

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}} \quad (9.11)$$

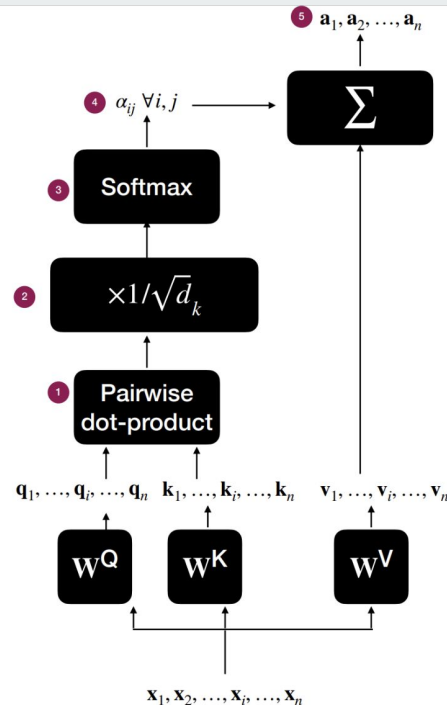$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \ \forall j \leq i \quad (9.12)$$

$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j \quad (9.13)$$

**Causal attention**
Because we are only looking at the past

**Full attention**
Looks at both past and the future

# Attention Head (contd.)

## Calculate α3

Now, calculate the score for each pair of $q_i$ and $k_j$

Then, do softmax for these three together.

Then we get $\alpha_{31}$, $\alpha_{32}$, and $\alpha_{33}$

- Now we have projected the inputs with three transformations

- We use the query and the key to compute attention

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}} \quad (9.11)$$

$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \ \forall j \leq i \quad (9.12)$$

$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j \quad (9.13)$$

**Causal attention**

Because we are only looking at the past

**Full attention**

Looks at both past and the future

## Calculate α3

Q: query
K: key

$$\frac{e^{\text{Score}(q_3,k_1)}}{e^{\text{Score}(q_3,k_1)} + e^{\text{Score}(q_3,k_2)} + e^{\text{Score}(q_3,k_3)}}$$

$$\frac{e^{\text{Score}(q_3,k_2)}}{e^{\text{Score}(q_3,k_1)} + e^{\text{Score}(q_3,k_2)} + e^{\text{Score}(q_3,k_3)}}$$

$$\frac{e^{\text{Score}(q_3,k_3)}}{e^{\text{Score}(q_3,k_1)} + e^{\text{Score}(q_3,k_2)} + e^{\text{Score}(q_3,k_3)}}$$

- Now we have projected the inputs with three transformations

- We use the query and the key to compute attention

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}} \quad \text{(9.11)}$$

$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \ \forall j \le i \quad \text{(9.12)}$$

$$\mathbf{a}_i = \sum_{j \le i} \alpha_{ij} \mathbf{v}_j \quad \text{(9.13)}$$

**Causal attention**
Because we are only looking at the past

**Full attention**
Looks at both past and the future

# Calculate α3

$$\text{Softmax}\left(\frac{QK^{\top}}{\sqrt{d_K}}\right) =$$

|      | K I  | am    | good  |
|------|------|-------|-------|
| I    | 0.90 | 0.07  | 0.03  |
| am   | 0.025| 0.95  | 0.025 |
| good | 0.21 | 0.03  | 0.76  |

$$\frac{e^{\text{Score}(q_3,k_1)}}{e^{\text{Score}(q_3,k_1)} + e^{\text{Score}(q_3,k_2)} + e^{\text{Score}(q_3,k_3)}}$$

$$\frac{e^{\text{Score}(q_3,k_2)}}{e^{\text{Score}(q_3,k_1)} + e^{\text{Score}(q_3,k_2)} + e^{\text{Score}(q_3,k_3)}}$$

$$\frac{e^{\text{Score}(q_3,k_3)}}{e^{\text{Score}(q_3,k_1)} + e^{\text{Score}(q_3,k_2)} + e^{\text{Score}(q_3,k_3)}}$$

# Attention Head (contd.)

## Query and Key

- Now we have projected the inputs with three transformations

- We use the query and the key to compute attention

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}} \tag{9.11}$$

$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \; \forall j \leq i \tag{9.12}$$
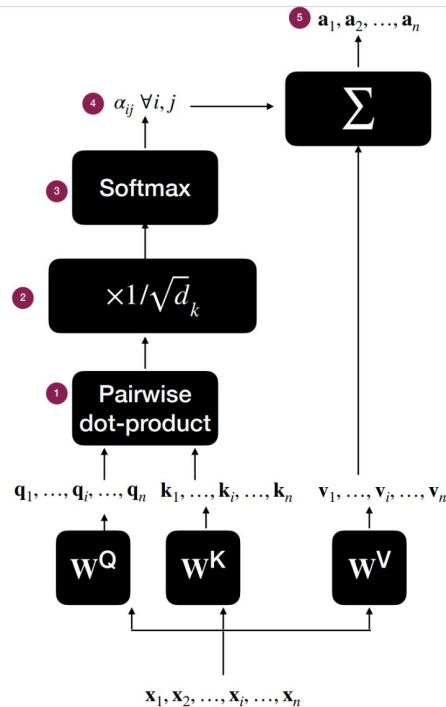
$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j \tag{9.13}$$

**Causal attention**

Because we are only looking at the past

**Full attention**

Looks at both past and the future

Note how it is written here:
For any j that is smaller than i…
This means
1. i is bigger than j
2. i is your current position, and j is the information before

# Query and Key

Query (Q): a search request in a database.
You send a request - Find me all books about dragons - to the library database.

Key (K): labels or tags attached to items in the database.
In the library, each book might have keys like "genre", "author", "subject"
*How to Kill a Dragon: 'dragon', 'Indo-European linguistics', 'Calvert Watkins'*

How does Q work with K？
The system compares the Q to the K.

# Calculate α3

What this graph means:

You calculate for x3



Output of self-attention  $a_3$

6. Sum the weighted value vectors

5. Weigh each **value** vector  $\alpha_{3,1}$  $\alpha_{3,2}$  $\alpha_{3,3}$

4. Turn into  $\alpha_{i,j}$  weights via softmax

3. Divide score by  $\sqrt{d_k}$   $\sqrt{d_k}$  $\sqrt{d_k}$  $\sqrt{d_k}$

2. Compare x3's **query** with the **keys** for x1, x2, and x3

1. Generate **key**, **query**, **value** vectors

$X_1$  $X_2$  $X_3$

# Calculate α3



Output of self-attention $a_3$

6. Sum the weighted value vectors

5. Weigh each **value** vector

$\alpha_{3,1}$ $\alpha_{3,2}$ $\alpha_{3,3}$

4. Turn into $\alpha_{i,j}$ weights via softmax

3. Divide score by $\sqrt{d_k}$ $\sqrt{d_k}$ $\sqrt{d_k}$ $\sqrt{d_k}$

2. Compare x3's **query** with the **keys** for x1, x2, and x3

1. Generate **key**, **query**, **value** vectors

$w^k$ k $w^k$ k $w^k$ k
$w^q$ q $w^q$ q $w^q$ q
$w^v$ v $w^v$ v $w^v$ v

$x_1$ $x_2$ $x_3$

# Calculate α3



Output of self-attention $a_3$

6. Sum the weighted value vectors

5. Weigh each **value** vector $\alpha_{3,1}$ $\alpha_{3,2}$ $\alpha_{3,3}$

4. Turn into $\alpha_{i,j}$ weights via softmax

3. Divide score by $\sqrt{d_k}$ $\sqrt{d_k}$ $\sqrt{d_k}$ $\sqrt{d_k}$

2. Compare x3's **query** with the **keys** for x1, x2, and x3

1. Generate **key**, **query**, **value** vectors

k q v k q v k q v

$w^k$ $w^q$ $w^v$

$x_1$ $x_2$ $x_3$

$$\text{Softmax}\left(\frac{QK^{\top}}{\sqrt{d_K}}\right) = \begin{array}{c|ccc} & \text{I} & \text{am} & \text{good} \\ \hline \text{I} & 0.90 & 0.07 & 0.03 \\ \text{am} & 0.025 & 0.95 & 0.025 \\ \text{good} & \boxed{0.21} & 0.03 & 0.76 \end{array}$$

$$\frac{e^{\text{Score}(q_3,k_1)}}{e^{\text{Score}(q_3,k_1)} + e^{\text{Score}(q_3,k_2)} + e^{\text{Score}(q_3,k_3)}}$$

$$\frac{e^{\text{Score}(q_3,k_2)}}{e^{\text{Score}(q_3,k_1)} + e^{\text{Score}(q_3,k_2)} + e^{\text{Score}(q_3,k_3)}}$$

$$\frac{e^{\text{Score}(q_3,k_3)}}{e^{\text{Score}(q_3,k_1)} + e^{\text{Score}(q_3,k_2)} + e^{\text{Score}(q_3,k_3)}}$$

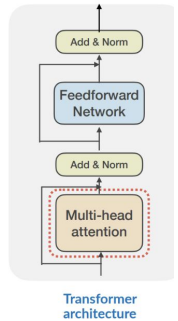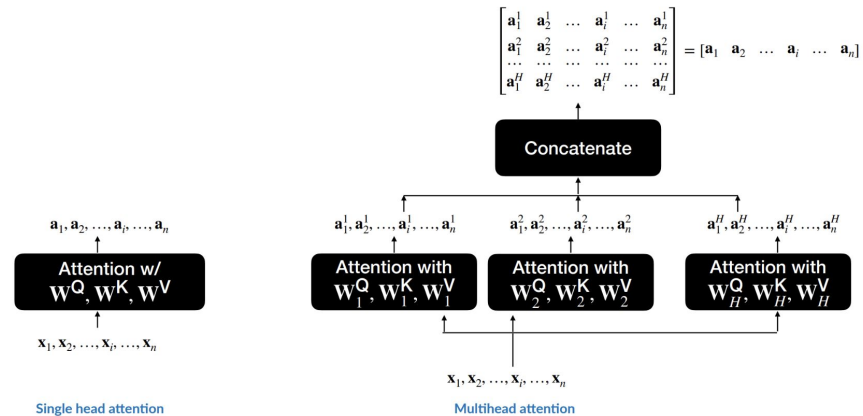# Attention

Masked-attention
Bidirectional Self Attention

## Multihead attention

$$\begin{bmatrix} \mathbf{a}_1^1 & \mathbf{a}_2^1 & \dots & \mathbf{a}_i^1 & \dots & \mathbf{a}_n^1 \\ \mathbf{a}_1^2 & \mathbf{a}_2^2 & \dots & \mathbf{a}_i^2 & \dots & \mathbf{a}_n^2 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \mathbf{a}_1^H & \mathbf{a}_2^H & \dots & \mathbf{a}_i^H & \dots & \mathbf{a}_n^H \end{bmatrix} = [\mathbf{a}_1 \quad \mathbf{a}_2 \quad \dots \quad \mathbf{a}_i \quad \dots \quad \mathbf{a}_n]$$

Concatenate

$\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_i, \dots, \mathbf{a}_n$

Attention w/
$\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V$

$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n$

**Single head attention**

$\mathbf{a}_1^1, \mathbf{a}_2^1, \dots, \mathbf{a}_i^1, \dots, \mathbf{a}_n^1$

$\mathbf{a}_1^2, \mathbf{a}_2^2, \dots, \mathbf{a}_i^2, \dots, \mathbf{a}_n^2$

$\mathbf{a}_1^H, \mathbf{a}_2^H, \dots, \mathbf{a}_i^H, \dots, \mathbf{a}_n^H$

Attention with
$\mathbf{W}_1^Q, \mathbf{W}_1^K, \mathbf{W}_1^V$

Attention with
$\mathbf{W}_2^Q, \mathbf{W}_2^K, \mathbf{W}_2^V$

Attention with
$\mathbf{W}_H^Q, \mathbf{W}_H^K, \mathbf{W}_H^V$

$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n$

**Multihead attention**

Add & Norm

Feedforward
Network

Add & Norm

Multi-head
attention

**Transformer
architecture**

# Attention

This formula in class is: masked self attention

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}} \; ①②$$

$$④ \; \alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \;\; \forall j \leq i \;\; ③$$

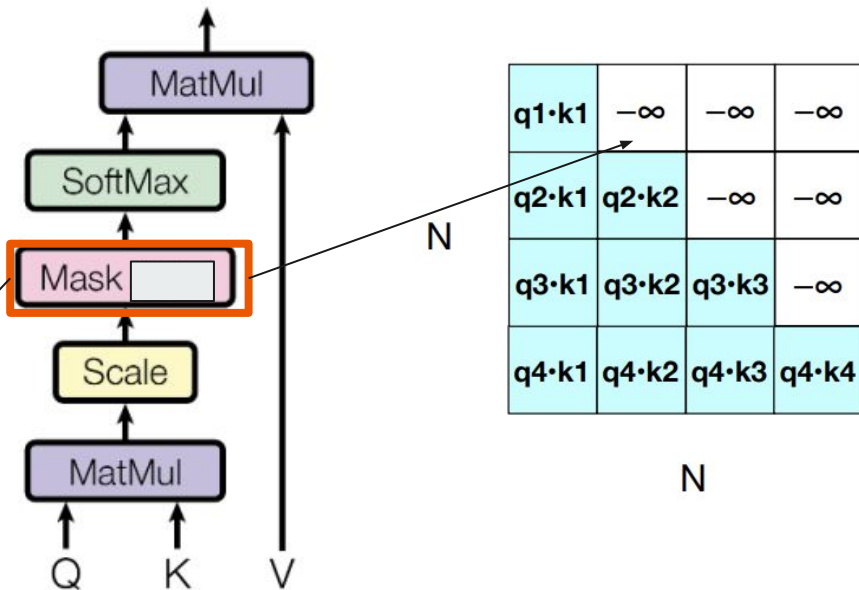$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$
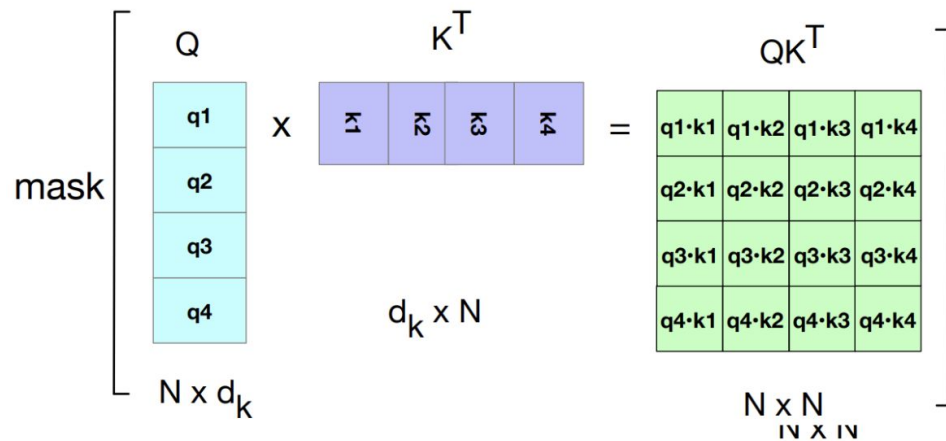
# Attention

Multi-head attention (see class slides)
Masked self-attention
Bidirectional self-attention

# Attention

Multi-head attention (see class slides)
Masked self-attention
Bidirectional self-attention

$$\text{Softmax}\left(\frac{QK^{\top}}{\sqrt{d_K}}\right) =$$

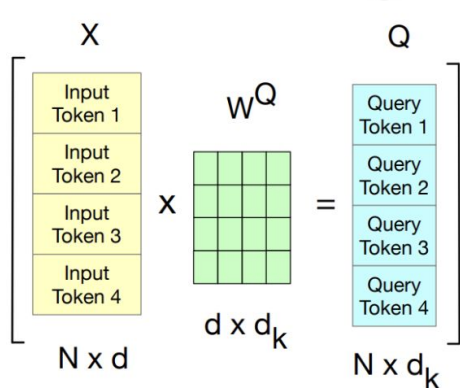|      | I     | am   | good  |
|------|-------|------|-------|
| I    | 0.90  | 0.07 | 0.03  |
| am   | 0.025 | 0.95 | 0.025 |
| good | 0.21  | 0.03 | 0.76  |

# Attention

Multi-head attention (see class slides)
Masked self-attention
Bidirectional self-attention
Encoder-decoder attention (not self-attention)

# Attention again

# More resources?

https://www.youtube.com/watch?v=iDulhoQ2pro (Yannic Kilcher)

Or search for "attention calculation" "self-attention math" etc. on Youtube or medium

# End