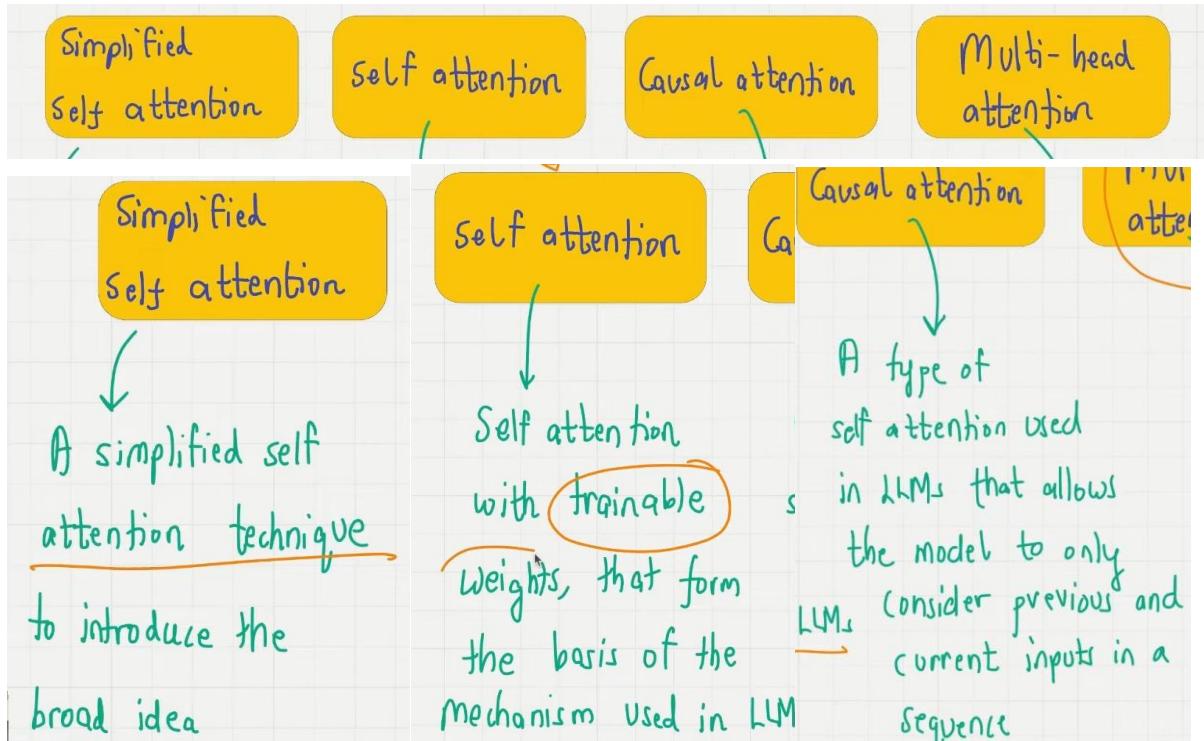


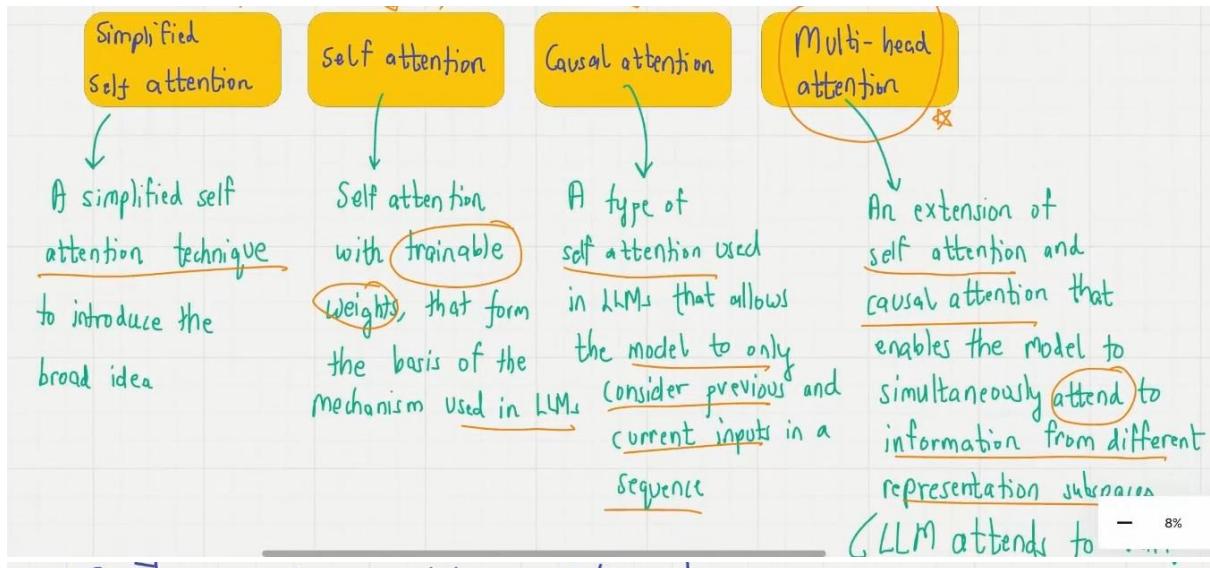
"The cat that was sitting on the mat,
which was next to the dog, jumped"

① The 4 types of attention mechanisms:



An extension of self attention and causal attention that enables the model to simultaneously attend to information from different representation subspaces

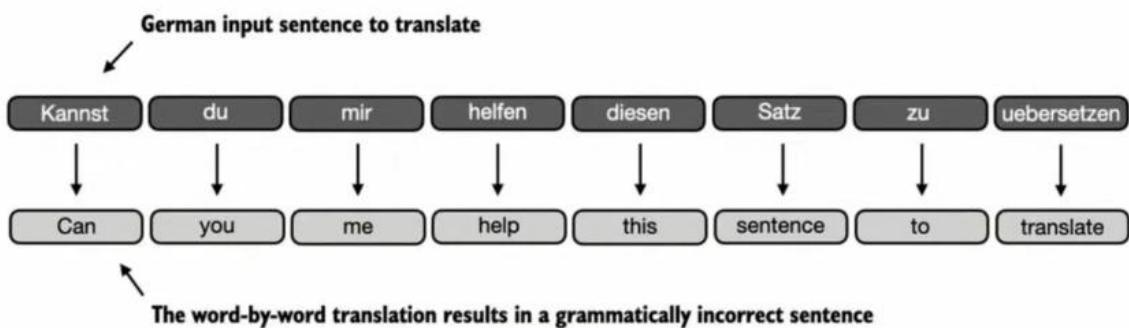
Multi-head attention
LLM attends to various input data in parallel)



② The problem with modeling long sequences:

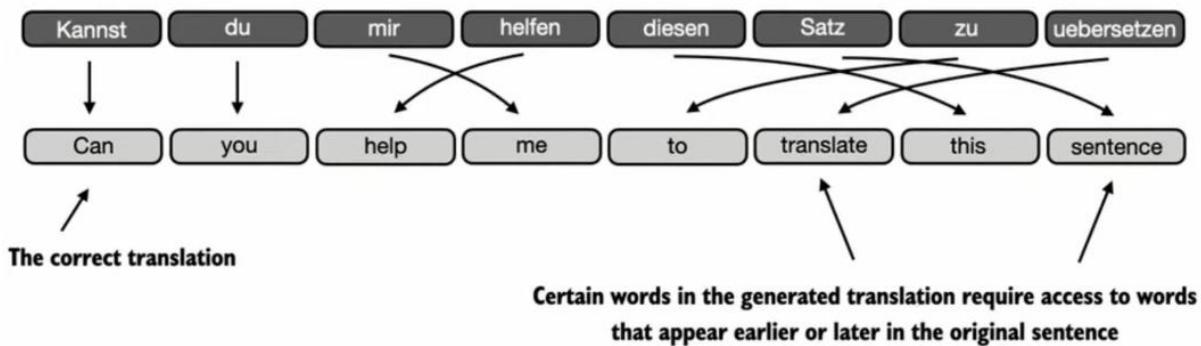
(a) What is the problem with architectures without the attention mechanism which came before LLMs?

* let's consider a language translation model *



Word by word translation does not work!

Can You Help me?
 કણ હુમ મિરી હેલ્પ મારી કોરોગો?



The translation process requires **(contextual)**
understanding and **grammar alignment**.

(b) To address this issue that we cannot translate text word by word, it is common to use a neural network with two submodules:

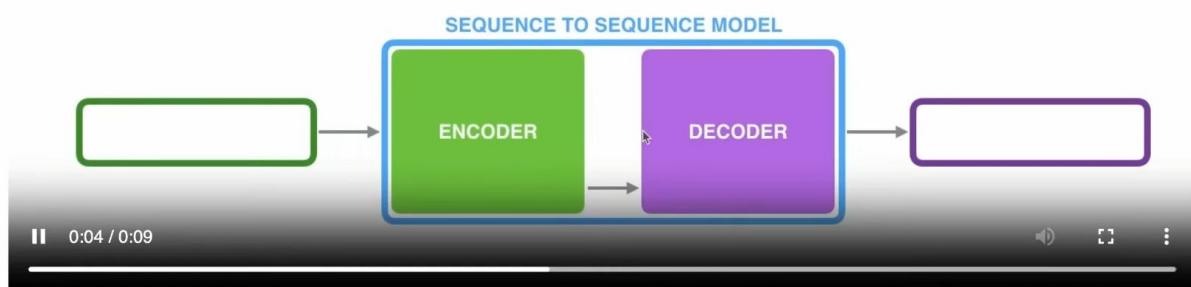
Encoder Decoder

[jalamar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/](http://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/)

Looking under the hood

Under the hood, the model is composed of an **encoder** and a **decoder**.

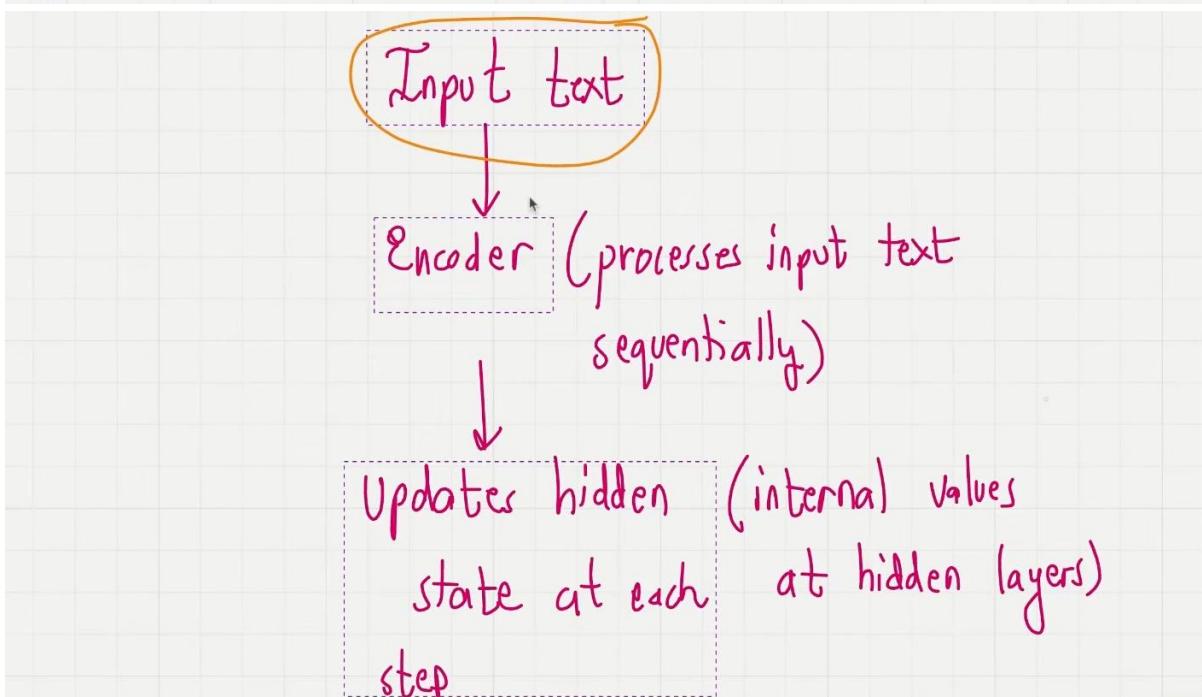
The **encoder** processes each item in the input sequence, it compiles the information it captures into a vector (called the **context**). After processing the entire input sequence, the **encoder** sends the **context** over to the **decoder**, which begins producing the output sequence item by item.

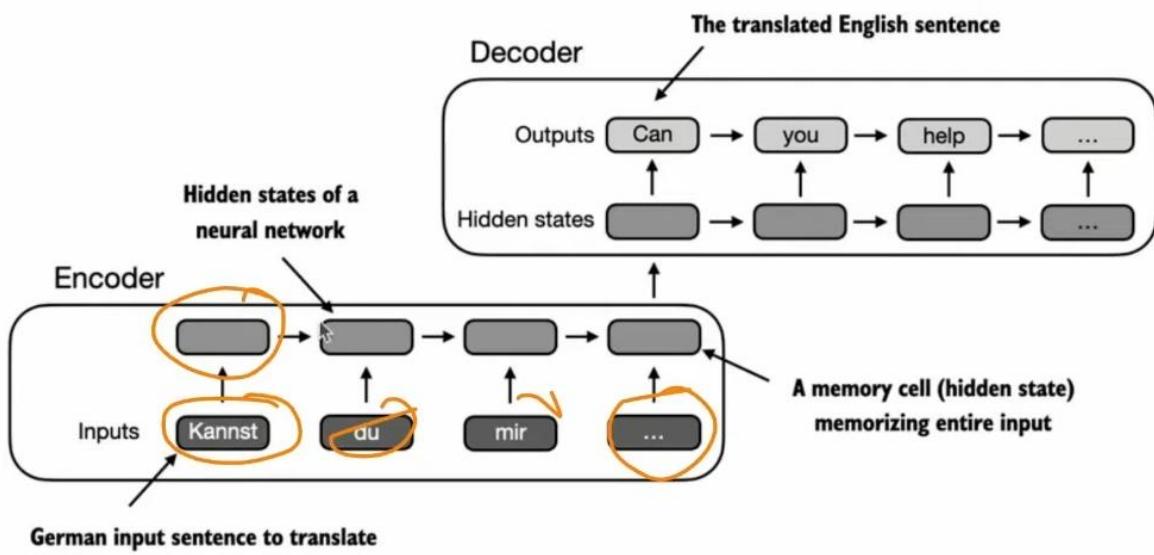
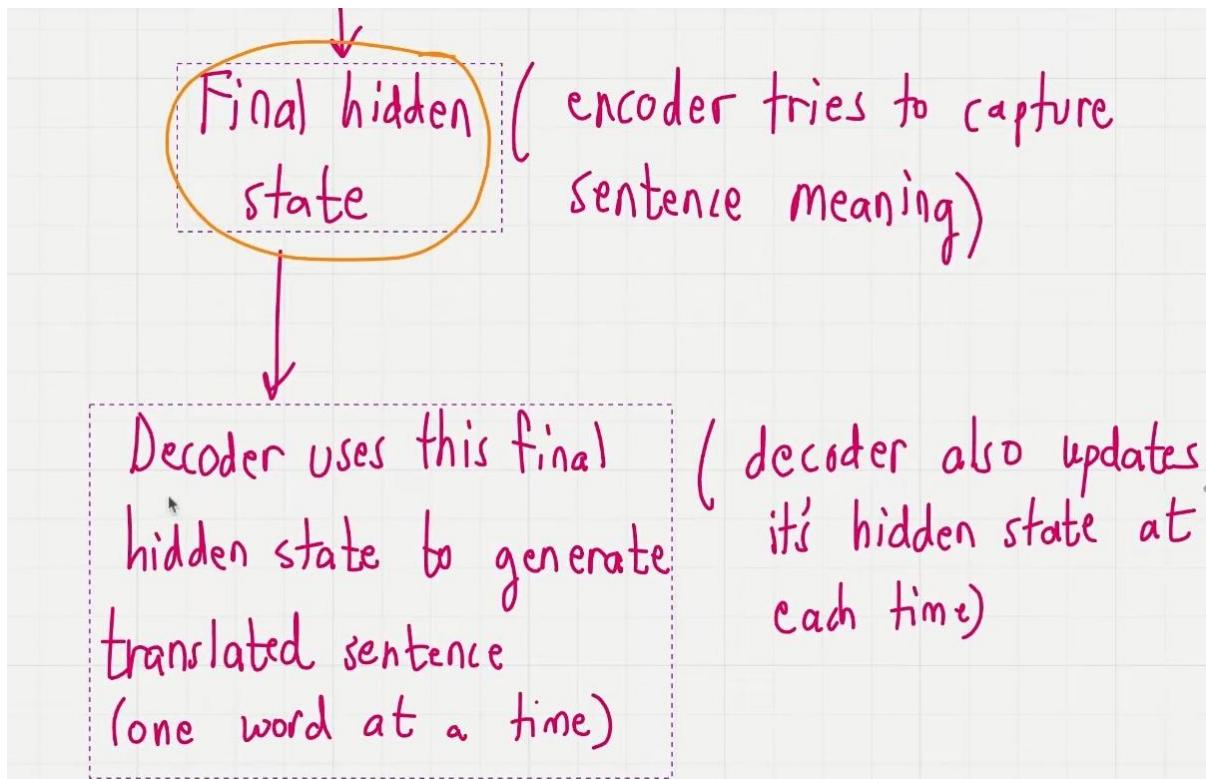


(c) Before transformers, Recurrent Neural Networks (RNNs) were the most popular encoder-decoder architecture for language translation.

(d) **RNN:** Output from previous step is fed as input to current text

(e) Here's how the encoder-decoder RNN works:





(f) The encoder processes the entire input text into a hidden state (memory cell). Decoder takes this hidden state to produce an output.
 think of this as an embedding vector.

Decoder takes this hidden state to produce an output.

den
- (g) * Big issue: RNN can't directly access earlier hidden states from the encoder, during the decoding phase.

It relies solely on current hidden state.

This leads to a loss of context, especially in complex sentences where dependencies might span long distances.

Encoder
compresses entire
input sequence into a single hidden
state vector

If the sentence is very long, it becomes very difficult for the RNN to capture all information in a single vector

"The cat that was sitting on the mat,
which was next to the dog, jumped"



"Le chat qui était assis sur le tapis,
qui était à côté du chien, a sauté"

Here, the key action "jumped" depends

Here, the key action "jumped" depends

on the subject "cat" but also on understanding

the longer dependencies ("that was sitting on the
mat, next to the dog")

The RNN decoder might struggle with this

③ Capturing data dependencies with attention mechanisms

(a) RNNs work fine for translating short sentences, but
don't work for long texts as they don't have
direct access to previous words in the input.

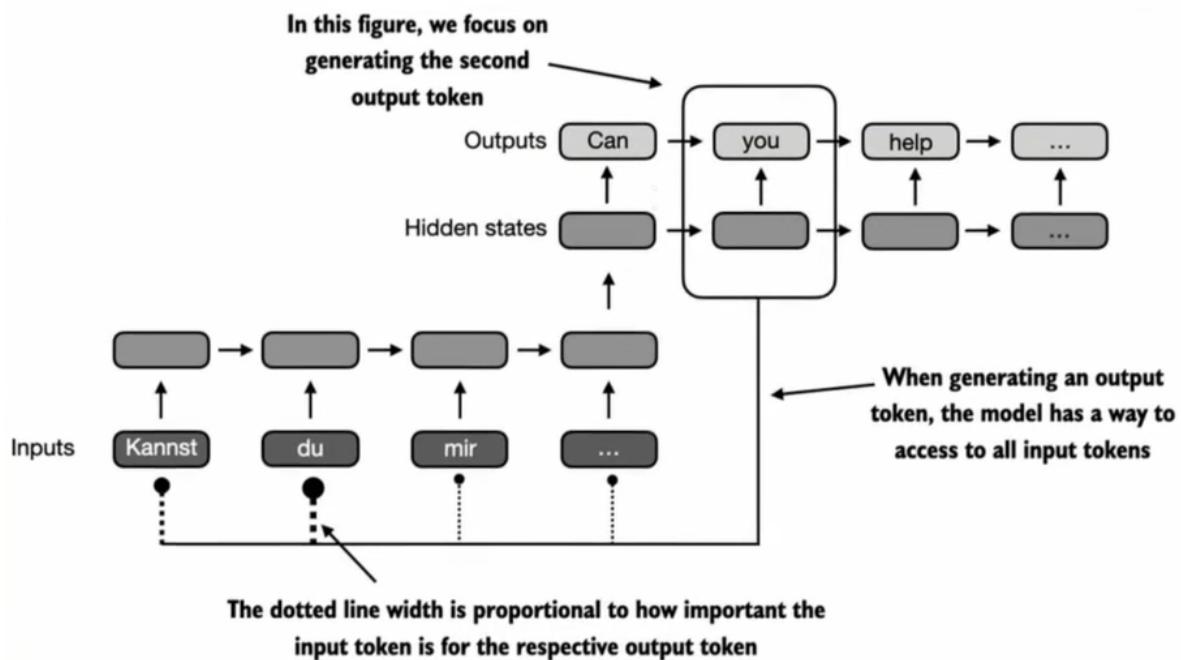
(b) One major shortcoming in this approach is that:

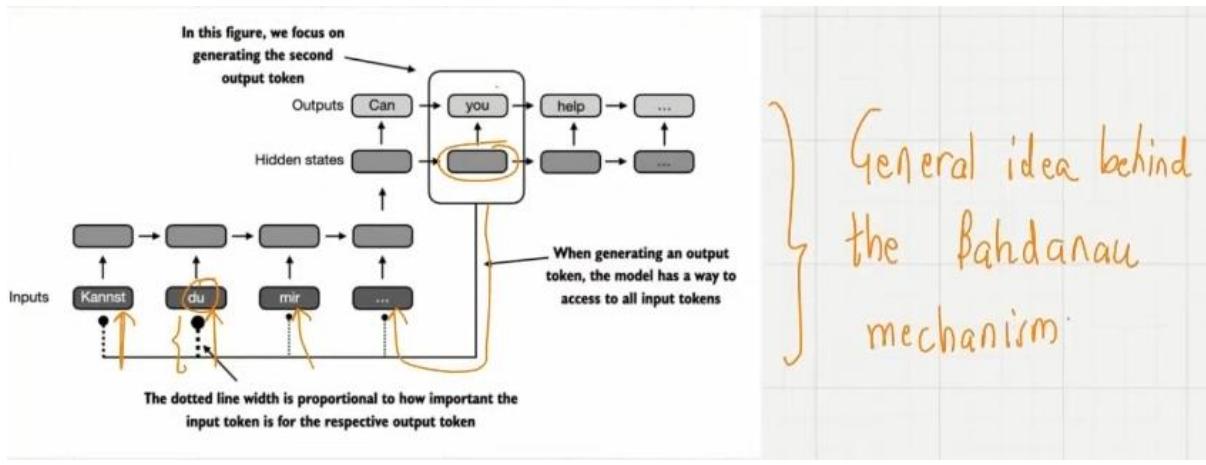
RNN must remember the entire encoded input in a single hidden state before passing it to the decoder

(c) In 2014, researchers developed the so called

"Bahdanau, attention mechanism for RNNs"

Modifies the encoder-decoder RNN such that the decoder can selectively access different parts of the input sequence at each decoding step.





- Using an attention mechanism, the text generating decoder part of the network can access all input tokens selectively.
- This means that some input tokens are more important than others for generating a given output token.
- This importance is determined by the so called attention weights.

(d) Only 3 years later, researchers found that RNN architectures are not required for building deep neural networks for natural language processing and proposed the original transformer architecture; with a self attention mechanism inspired by the Bahdanau attention mechanism.

"The cat that was sitting on the mat, which was next to the dog, jumped"

"Le chat qui était assis sur le tapis, qui était à côté du chien, a sauté"

At each decoding step, the model can look back at the entire input sequence and decide which parts are most relevant to generate current word

When decoder is predicting "sauté",
 the attention mechanism allows
 it to focus on part of input
 that corresponds to "jumped".
 Dynamic focus on different parts
 of input sequence allows models
 to learn long range dependencies
 more effectively.

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



RNNs

LSTMs

Attention

Transformers

1980

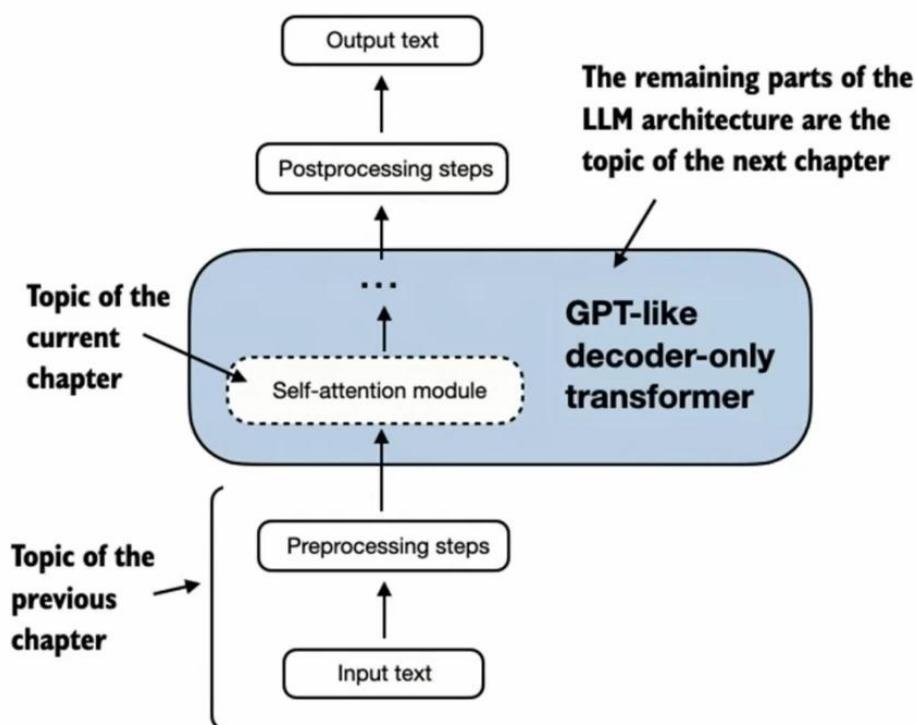
1997

2014

2017

(e) Self attention is a mechanism that allows each position of the input sequence to attend to all positions in the same sequence when computing the representation of a sequence.

(f) Self attention is a key component of contemporary LLMs based on the transformer architecture, such as the GPT series.



④ Attending to different parts of the input with self attention

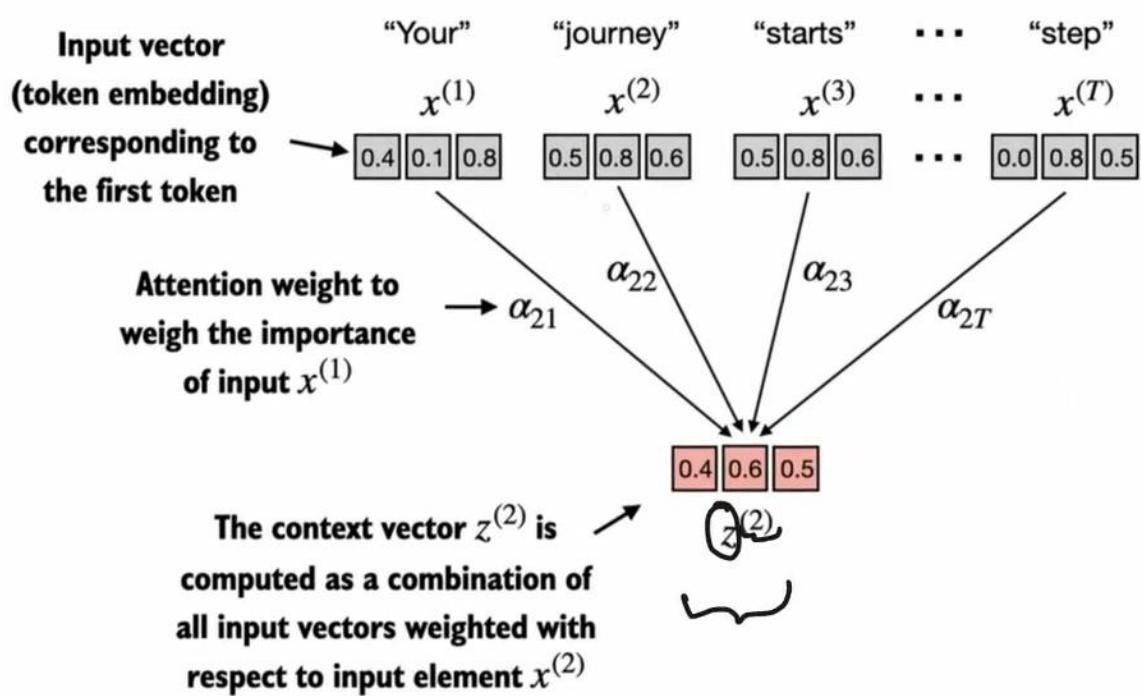
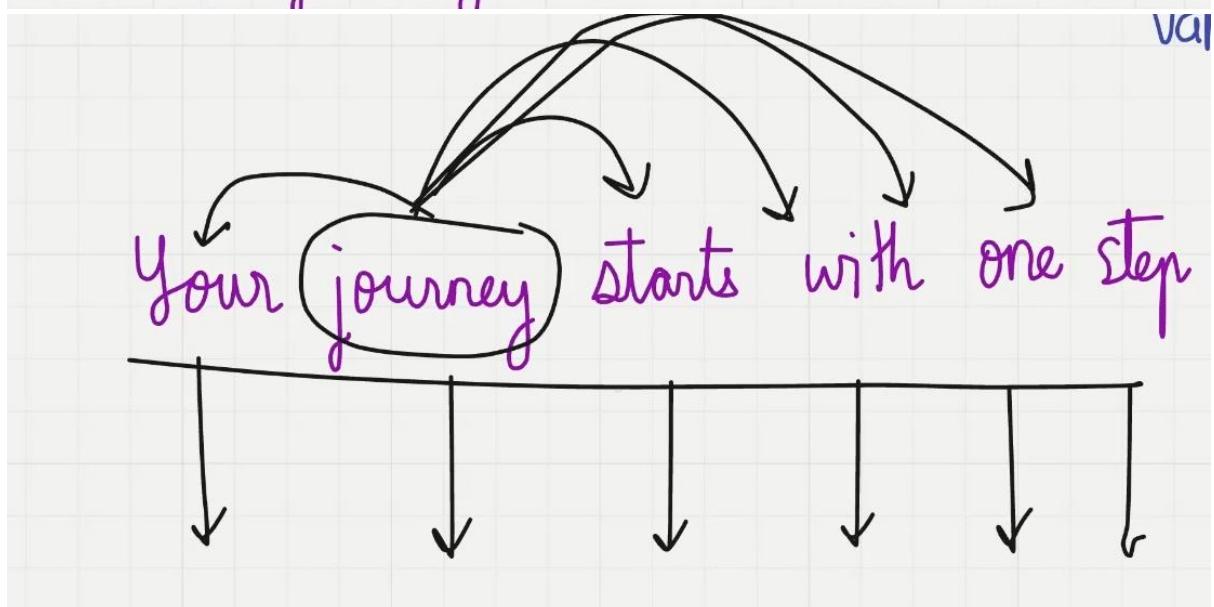
(a) In "self attention", the "self" refers to the mechanisms' ability to compute attention weights by relating different positions in a single input sequence.

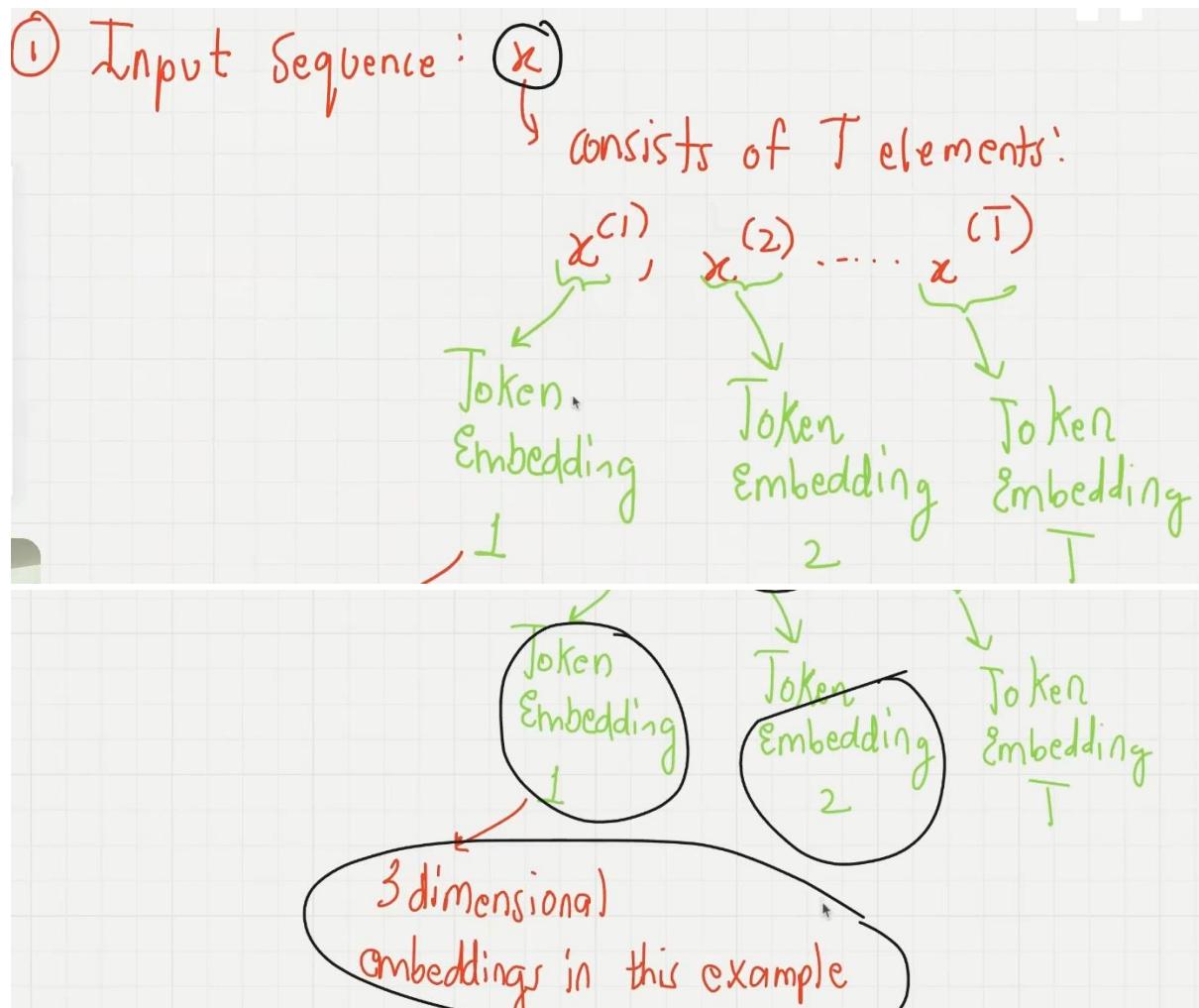
- b) It learns the relationship between various parts of the input itself, such as words in a sentence.
- c) This is in contrast to traditional attention mechanisms, where the focus is on relationships between elements of 2 different sequences.

A simple self attention mechanism without trainable weights

In this section, we will implement a simple variant of self-attention, free from any trainable weights.

Your journey starts with one step





② Goal of self attention:

"Calculate a context vector $z^{(i)}$ for each element $x^{(i)}$ "

Context vector: Enriched embedding vector

Let us focus on 2nd element ($x^{(2)}$)

Corresponding context vector is $\underline{z}^{(2)}$.

$\underline{z}^{(2)}$ is an embedding which contains information about $\underline{x}^{(2)}$ and all other input elements $x^{(1)}$ to $x^{(T)}$

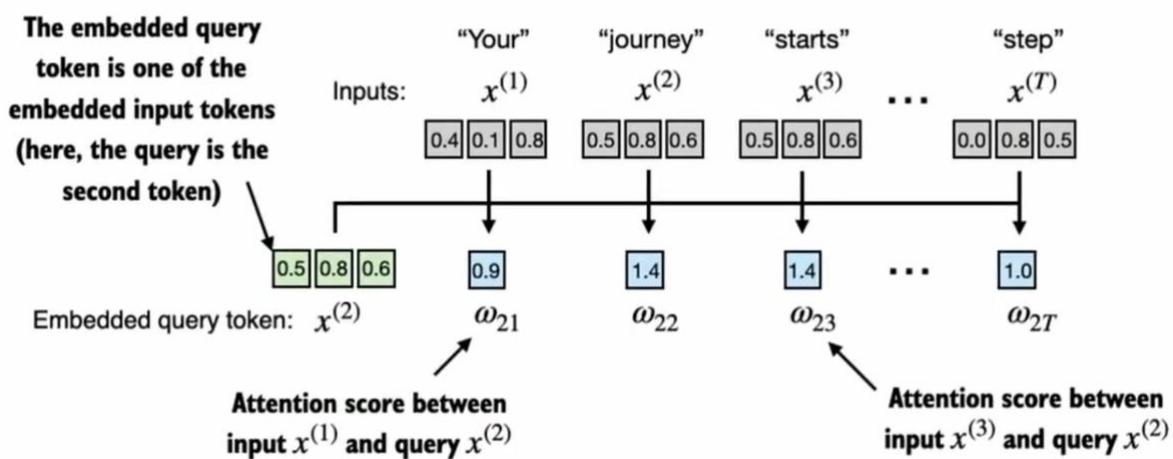
2nd element ($x^{(2)}$)

context vector is $\underline{z}^{(2)}$.

query

Later, we will add trainable weights that help the LLM learn to construct these context vectors; so that they are relevant for the LLM to generate the next token.

First step of implementing self-attention is to compute the intermediate values w , also referred to as attention scores.



6. Intermediate attention scores calculated between the query token and each input taken

Query: Dot products quantifies how much two vectors aligned

→ In the context of self attention mechanisms, dot product determines the extent to which elements of a sequence attend to one another → higher the dot product, higher the similarity and attention scores between two elements.

In the next step, we normalize each of the attention scores that we computed previously.

In practice, it's more common and advisable to use the softmax function for normalization.

This approach is better at managing extreme values and offers more favorable gradient properties.

Below is a basic implementation of the softmax function for normalizing the attention scores:

$$\begin{array}{c}
 \xrightarrow{\hspace{1cm}} [x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ x_6] \\
 \downarrow \text{Softmax} \\
 \left[\frac{e^{x_1}}{\text{sum}} \ \frac{e^{x_2}}{\text{sum}} \ \frac{e^{x_3}}{\text{sum}} \ \frac{e^{x_4}}{\text{sum}} \ \frac{e^{x_5}}{\text{sum}} \ \frac{e^{x_6}}{\text{sum}} \right]
 \end{array}$$

$$\left[\begin{array}{cccccc} e^{x_1} & e^{x_2} & e^{x_3} & e^{x_4} & e^{x_5} & e^{x_6} \\ \text{sum} & \text{sum} & \text{sum} & \text{sum} & \text{sum} & \text{sum} \end{array} \right]$$

↓

$$e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4} + e^{x_5} + e^{x_6}$$

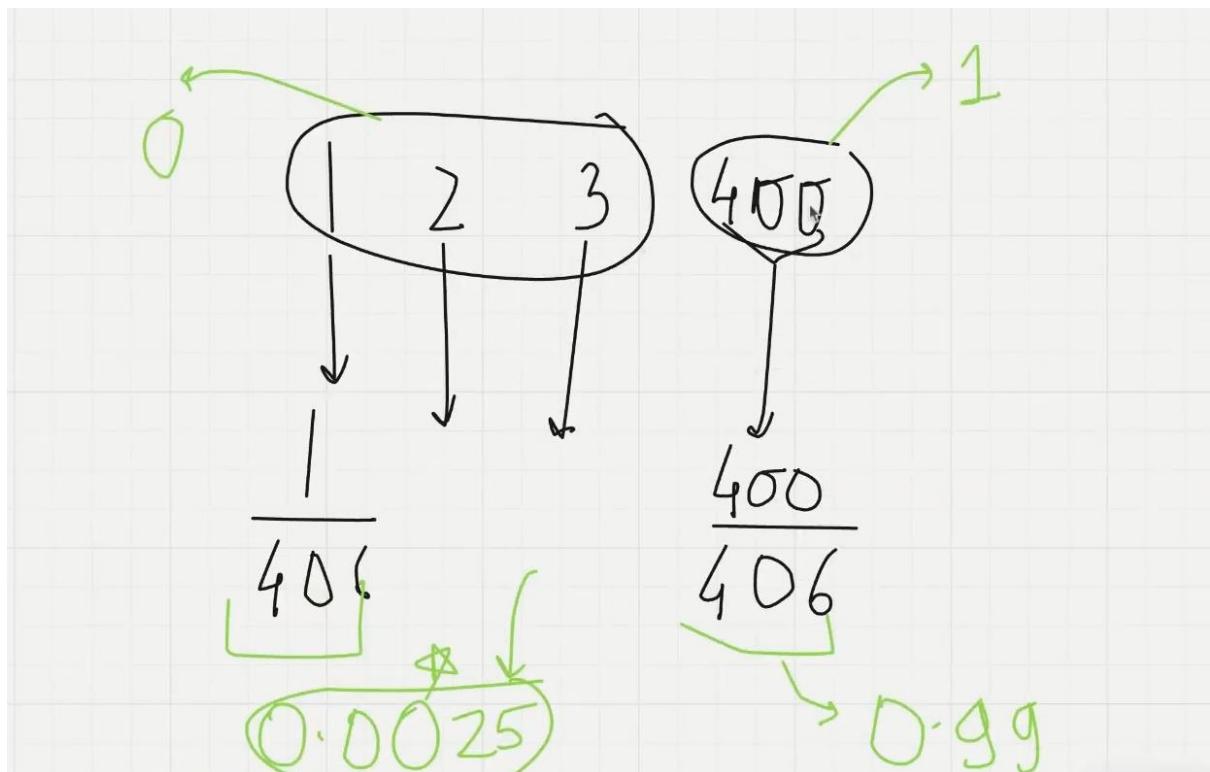
PyTorch

Softmax

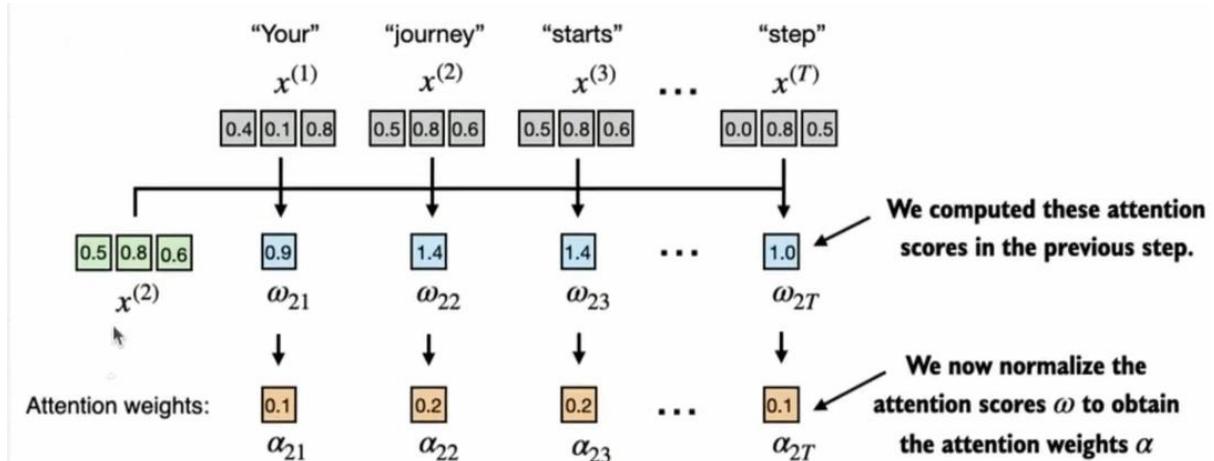
$$\left[\begin{array}{cccccc} x_1 & x_2 & x_3 & x_4 & x_5 & x_6 \end{array} \right]$$

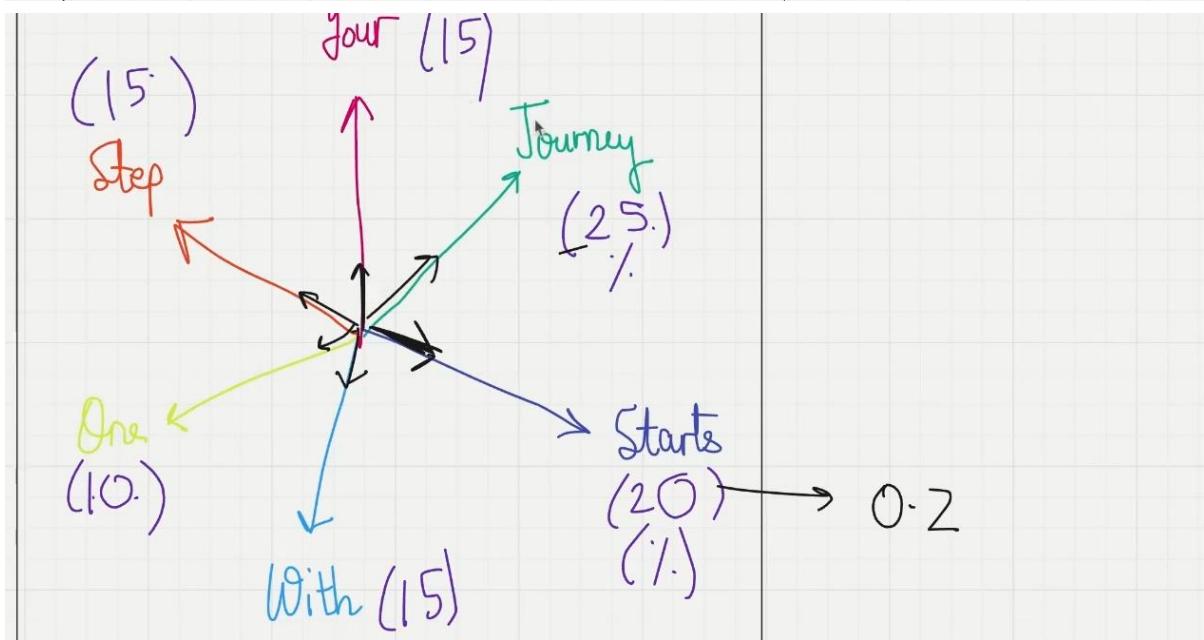
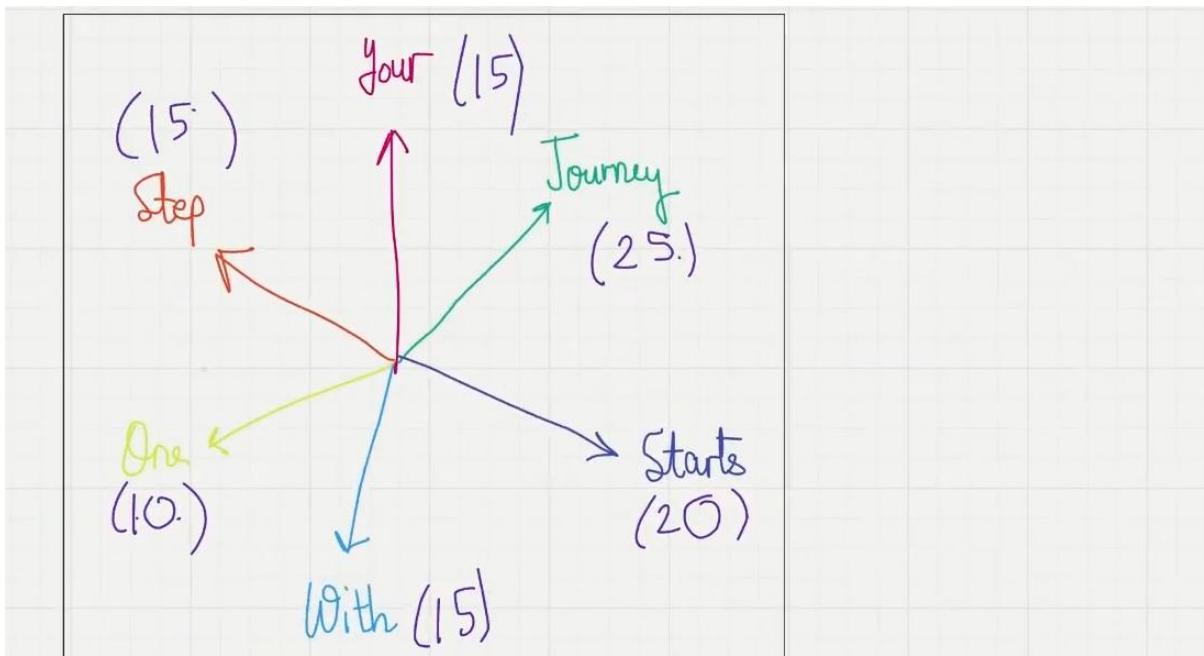
↓

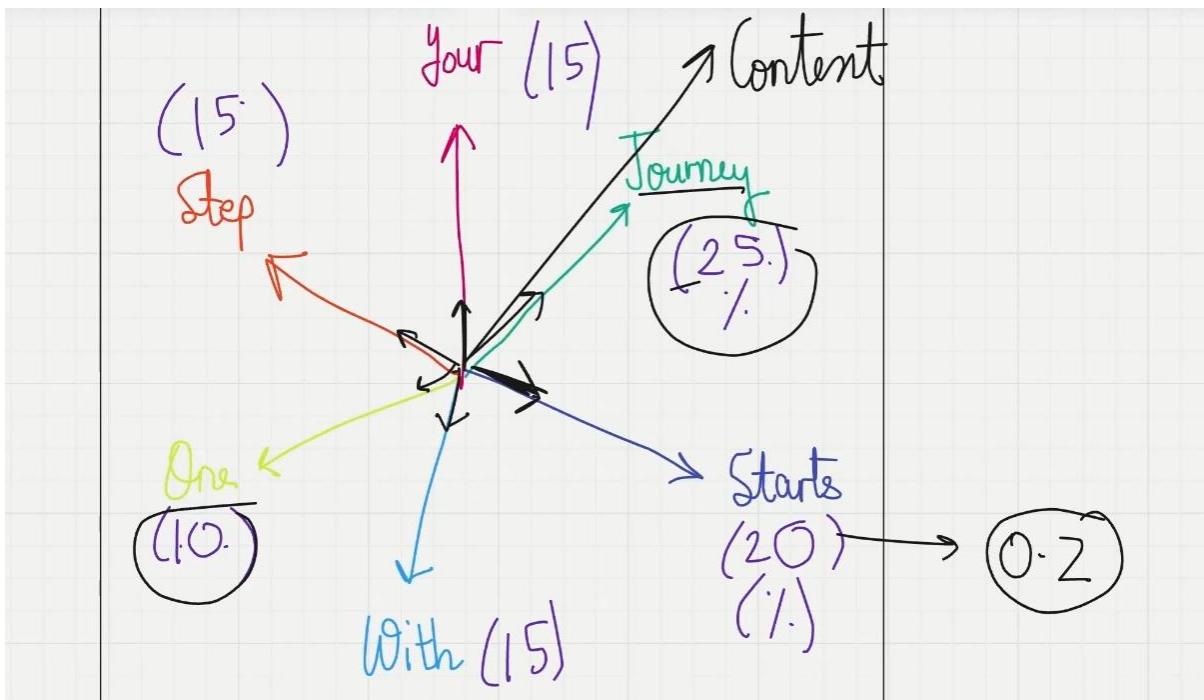
$$\left[\begin{array}{cccccc} e^{x_1 - \max} & e^{x_2 - \max} & e^{x_3 - \max} & e^{x_4 - \max} \\ \text{sum} & \text{sum} & \text{sum} & \text{sum} \end{array} \right]$$



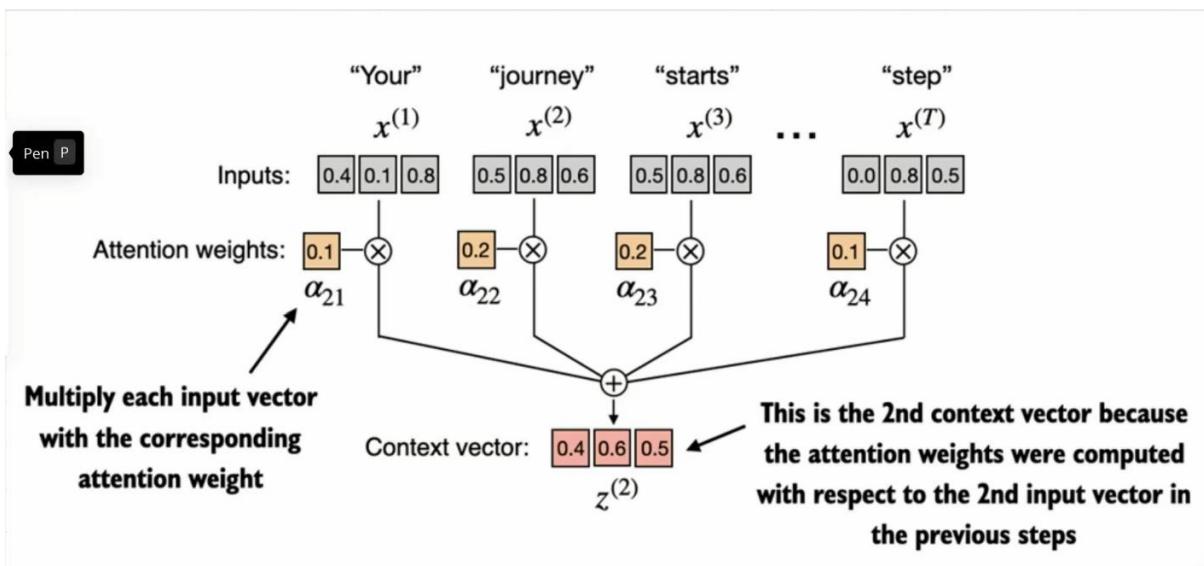
$$\begin{aligned}
 & \left[\frac{e^{x_1 - \max}}{\text{sum}}, \frac{e^{x_2 - \max}}{\text{sum}}, \frac{e^{x_3 - \max}}{\text{sum}}, \frac{e^{x_4 - \max}}{\text{sum}} \right] \\
 & \text{sum} = e^{x_1 - \max} + e^{x_2 - \max} + \dots + e^{x_T - \max}
 \end{aligned}$$







(8) After computing the normalized attention weights, we calculate the context vector $z^{(2)}$ by multiplying the embedded input tokens $x^{(i)}$, with the corresponding attention weights and then summing the resultant vectors.



⑨ Computing attention weights for all input tokens.



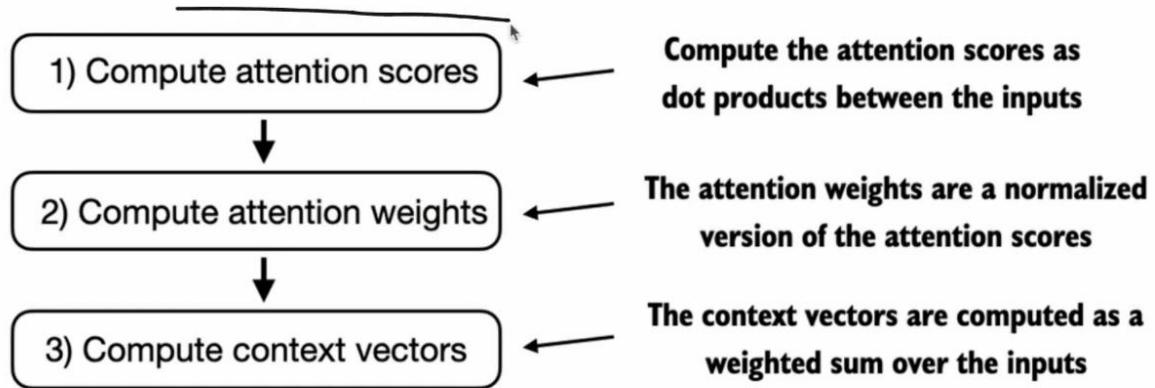
We can extend similar computation as before to calculate attention weights and context vectors for all inputs

	Your	journey	starts	with	one	step
Your	0.20	0.20	0.19	0.12	0.12	0.14
journey	0.13	0.23	0.23	0.12	0.10	0.15
starts	0.13	0.23	0.23	0.12	0.11	0.15
with	0.14	0.20	0.20	0.14	0.12	0.17
one	0.15	0.19	0.19	0.13	0.18	0.12
step	0.13	0.21	0.21	0.14	0.09	0.18

← This row contains the attention weights (normalized attention scores) computed previously

⑩

We will follow the same 3 steps as before:



Attention weights:

$$r_2 \left\{ \begin{bmatrix} 0.2098 & 0.20 & 0.19 & 0.12 & 0.12 & 0.14 \\ 0.138 & 0.237 & 0.233 & 0.124 & 0.108 & 0.158 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \times \begin{bmatrix} 0.43 & 0.15 & 0.89 \\ 0.55 & 0.87 & 0.66 \\ \vdots & \vdots & \vdots \\ c_1 & c_2 & c_3 \end{bmatrix} \right\}$$

Inputs: $\begin{bmatrix} 0.43 & 0.15 & 0.89 \\ 0.55 & 0.87 & 0.66 \\ \vdots & \vdots & \vdots \\ c_1 & c_2 & c_3 \end{bmatrix}$

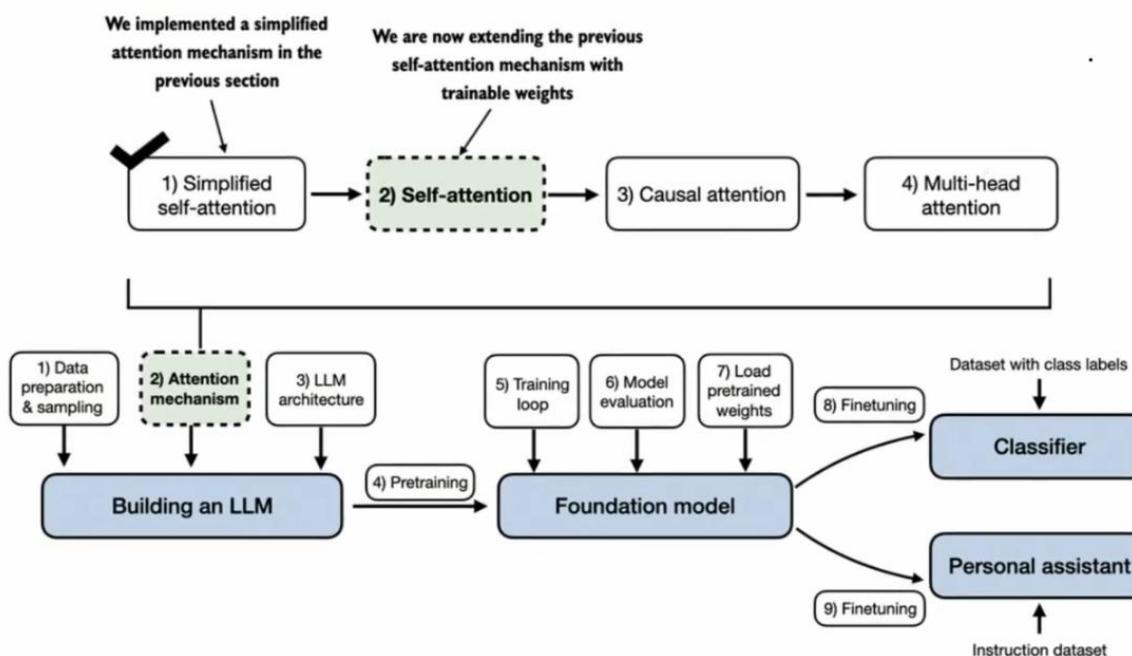
$$\begin{bmatrix} r_2 \cdot c_1 & r_2 \cdot c_2 & r_2 \cdot c_3 \end{bmatrix} = \begin{bmatrix} \square & \square & \square \\ \square & \square & \square \\ \vdots & \vdots & \vdots \\ \square & \square & \square \end{bmatrix} \quad 6 \times 3$$

$$= 0.138 * [0.43 \ 0.15 \ 0.89] \\ + 0.237 * [0.55 \ 0.87 \ 0.66] \\ + \dots$$

Lecture → Self attention with trainable weights

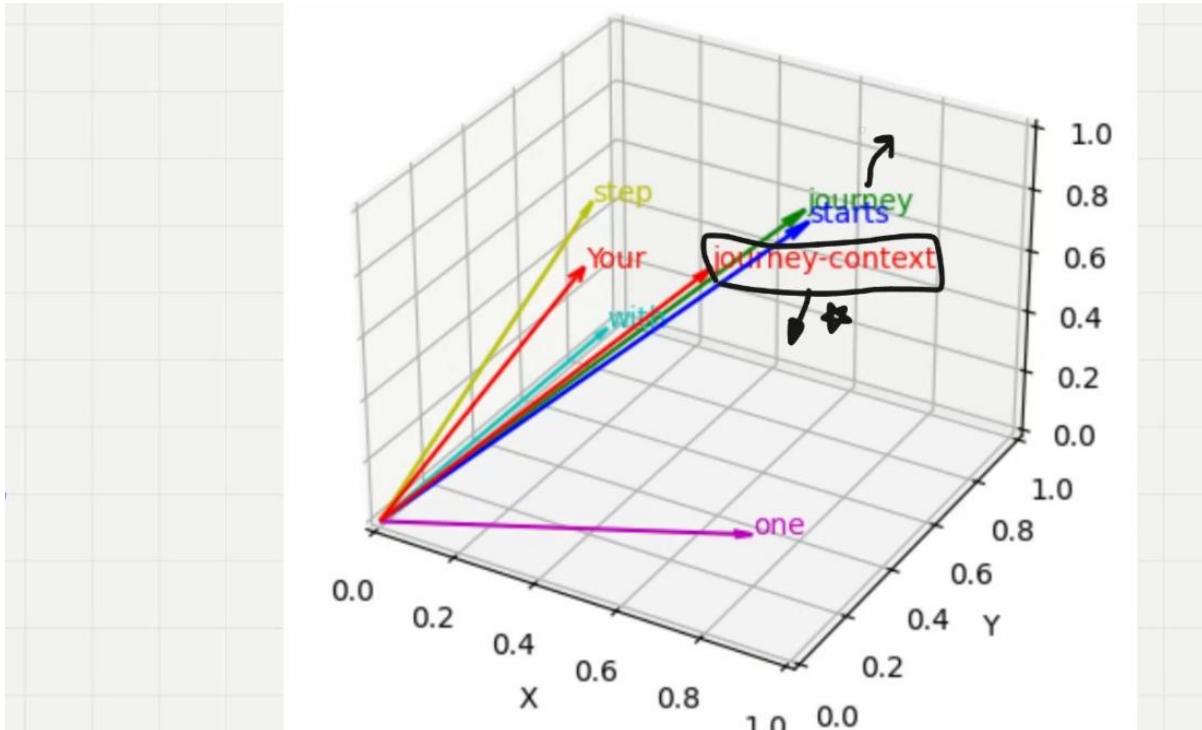
↳ In this section, we will learn about the self attention mechanism used in the original transformer architecture, the GPT models and most other popular LMs.

↳ This self attention mechanism is also called as "scaled dot-product attention"



- ① We want to compute context vectors as weighted sums over the input vectors specific to a certain input element.
- ② We will introduce weight matrices that are updated during model training.

- ② We will introduce weight matrices that are updated during model training.
- ③ These trainable weight matrices are crucial so that the model can learn to produce "good" context vectors.

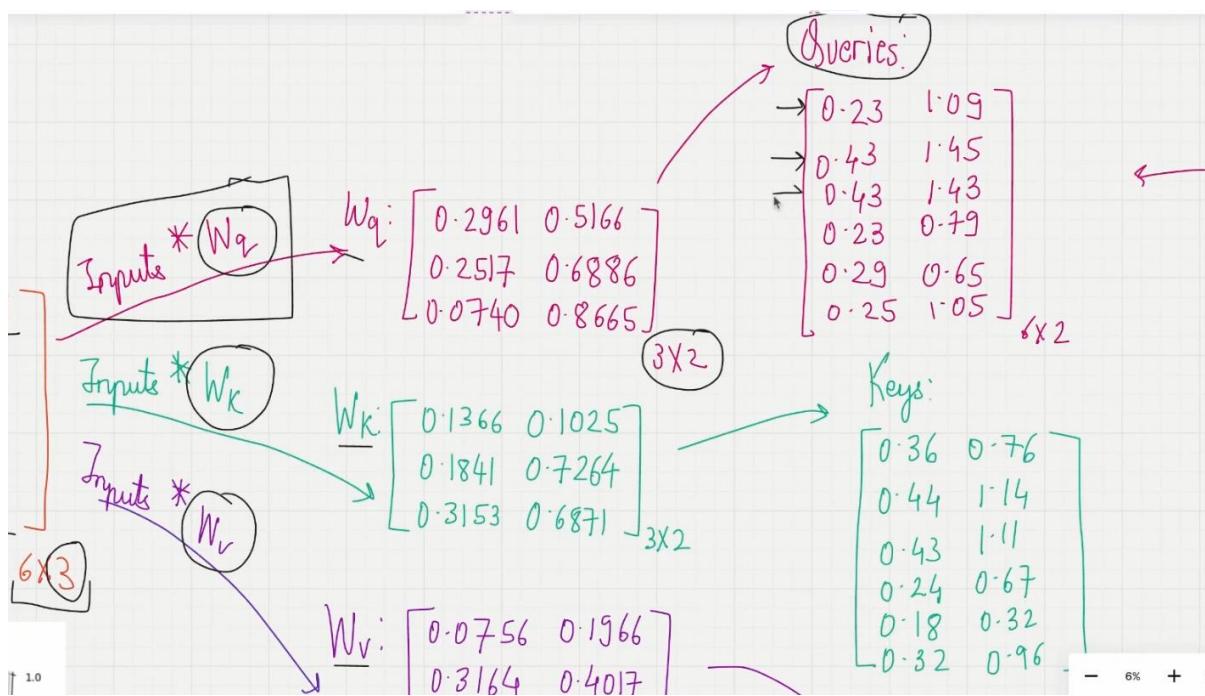
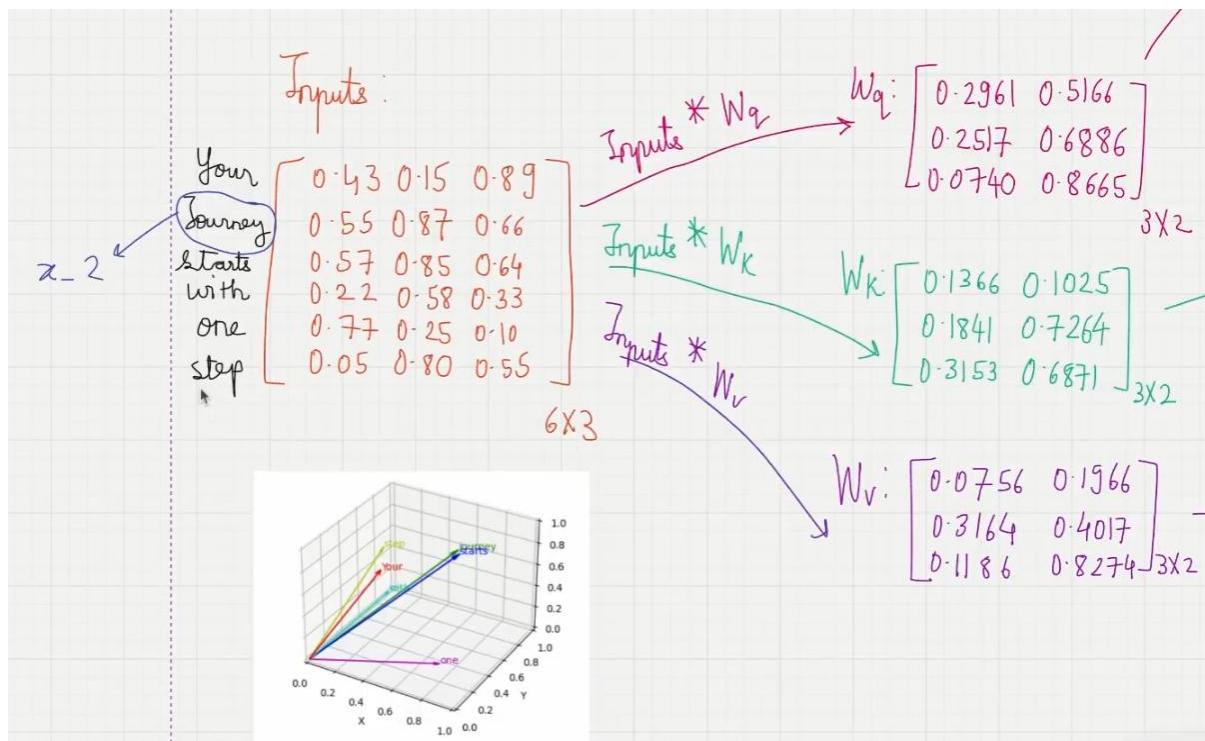


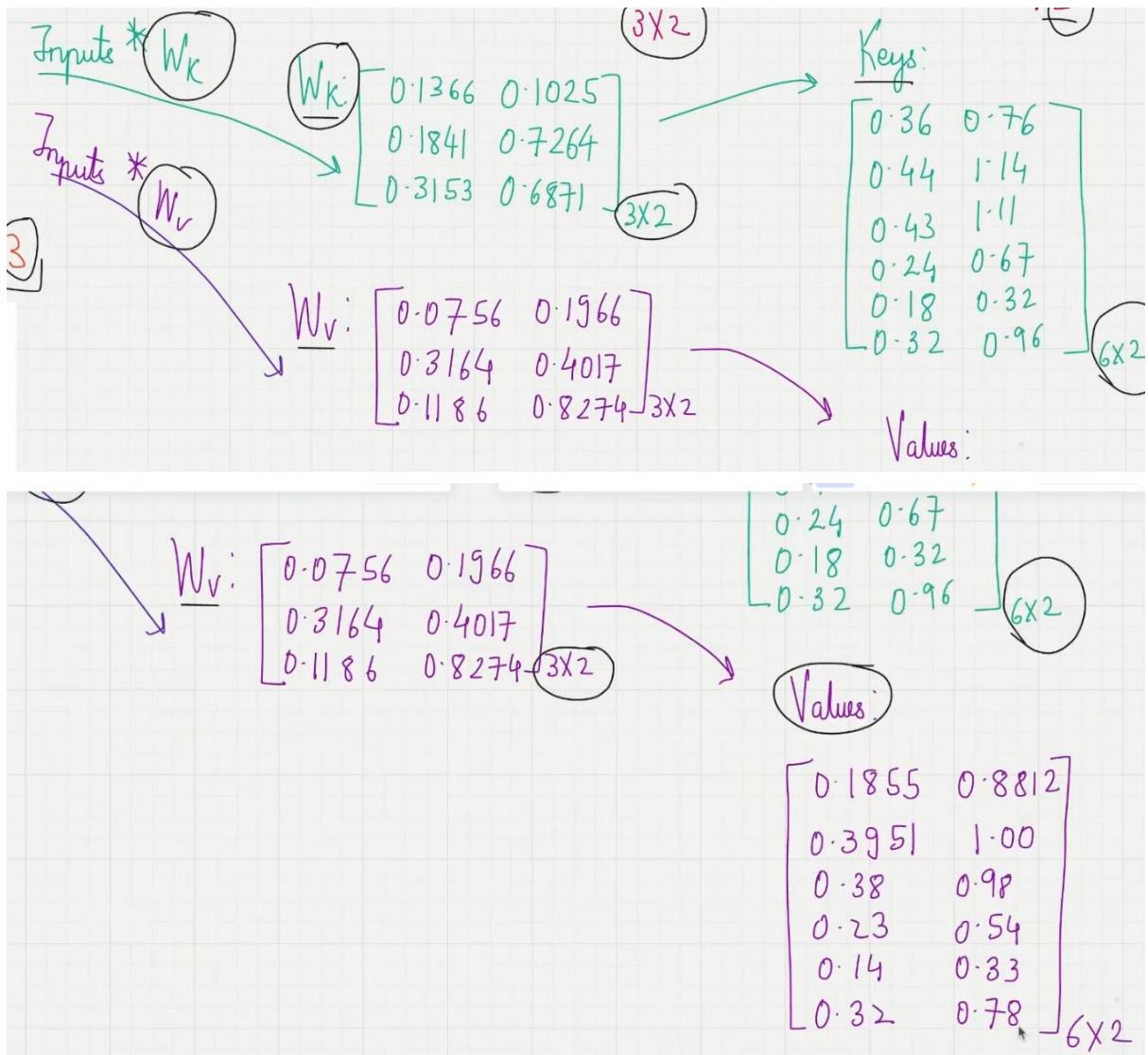
- ④ We will implement the self attention mechanism step by step by introducing 3 trainable weight matrices: W_q , W_k and W_v

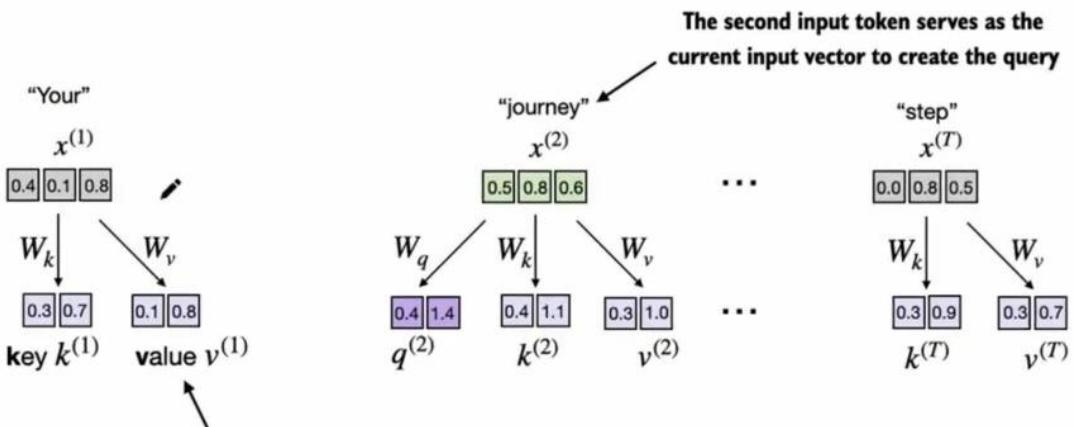
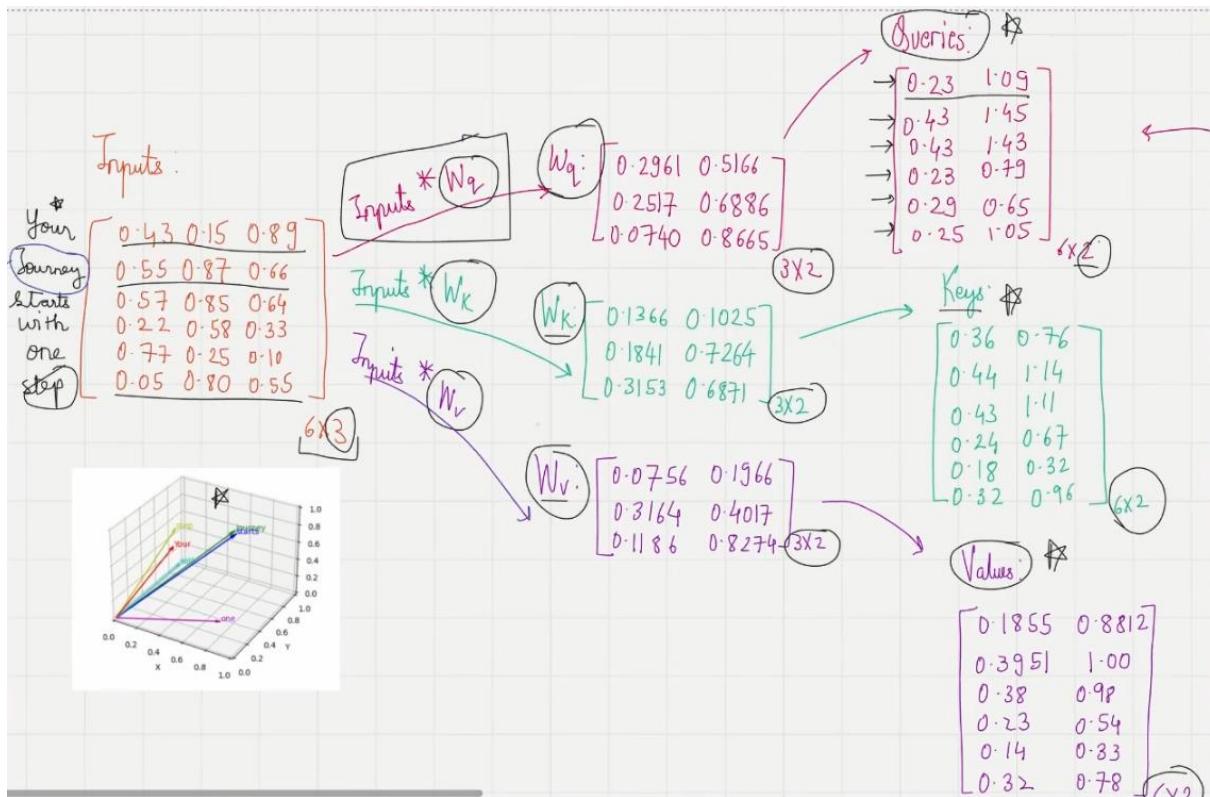
These 3 matrices are used to project the embedded input tokens $x^{(i)}$ into query, key and value vectors



① Converting input embeddings into key, query, value vectors

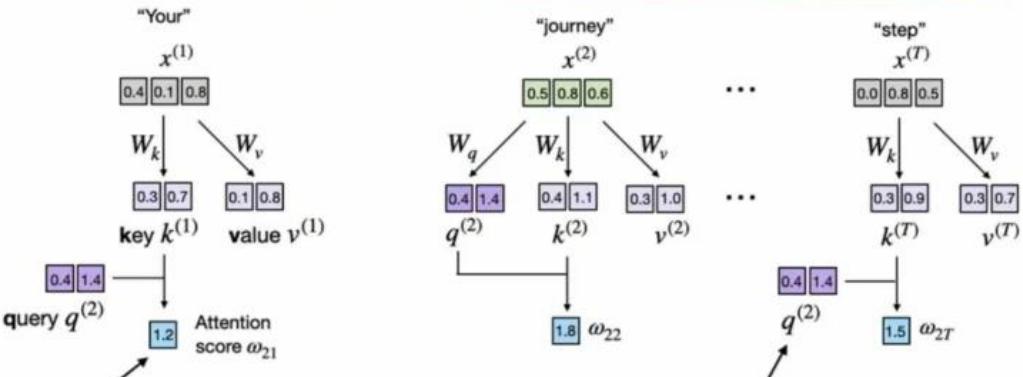






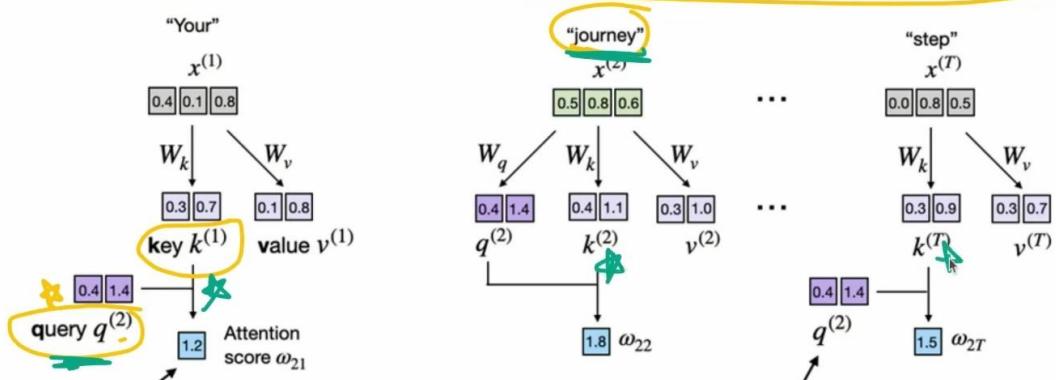
This is the value vector corresponding to the first input token obtained via matrix multiplication between the weight matrix W_v and input token $x^{(1)}$

- ⑤ We will start by computing only one context vector $z^{(2)}$ for illustration purposes
- Computing the attention scores:



The unscaled attention score is computed as a dot product between the query and the key vectors

Since we want to compute the context vector for the second input token, the query is derived from that second input token



The unscaled attention score is computed as a dot product between the query and the key vectors

Since we want to compute the context vector for the second input token, the query is derived from that second input token

Queries:

0.23	1.09
0.43	1.45
0.43	1.43
0.23	0.79
0.29	0.65
0.25	1.05

6x2

Keys:

0.36	0.76
0.44	1.14
0.43	1.11
0.24	0.67
0.18	0.32
0.32	0.96

6x2

Queries * Keys^T

$$\begin{bmatrix}
 0.23 & 1.09 \\
 0.43 & 1.45 \\
 0.43 & 1.43 \\
 0.23 & 0.79 \\
 0.29 & 0.65 \\
 0.25 & 1.05
 \end{bmatrix}_{6 \times 2} *
 \begin{bmatrix}
 0.36 & 0.44 & 0.43 & 0.24 & 0.18 & 0.32 \\
 0.76 & 1.14 & 0.24 & 0.67 & 0.32 & 0.96
 \end{bmatrix}_{2 \times 6}$$

`tensor([[0.9231, 1.3545, 1.3241, 0.7910, 0.4032, 1.1330],
 [1.2705, 1.8524, 1.8111, 1.0795, 0.5577, 1.5440],
 [1.2544, 1.8284, 1.7877, 1.0654, 0.5508, 1.5238],
 [0.6973, 1.0167, 0.9941, 0.5925, 0.3061, 0.8475],
 [0.6114, 0.8819, 0.8626, 0.5121, 0.2707, 0.7307],
 [0.8995, 1.3165, 1.2871, 0.7682, 0.3937, 1.0996]])`

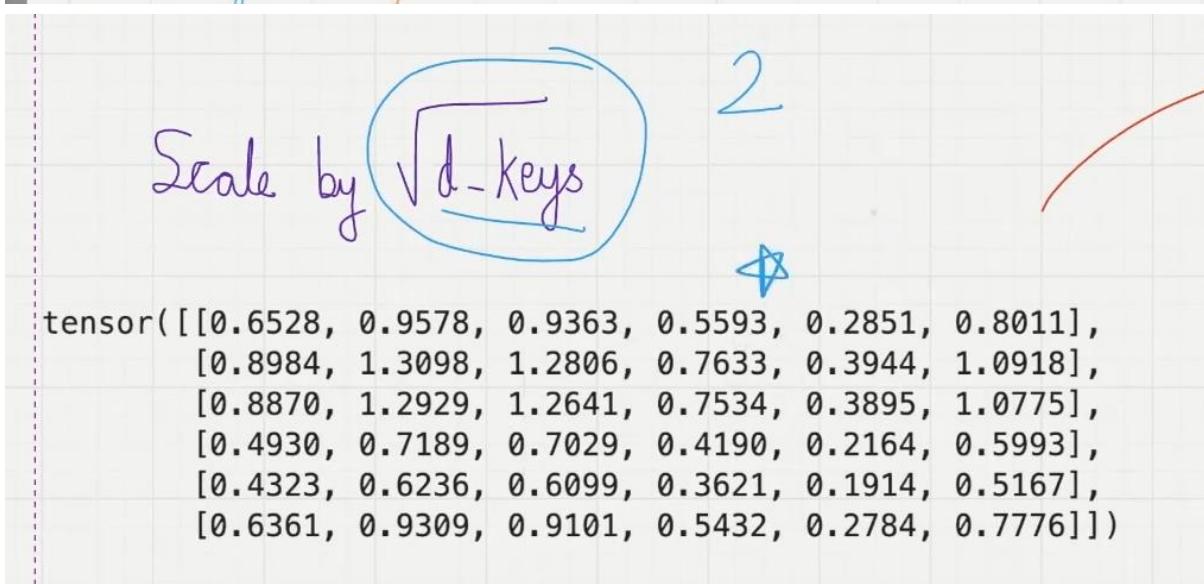
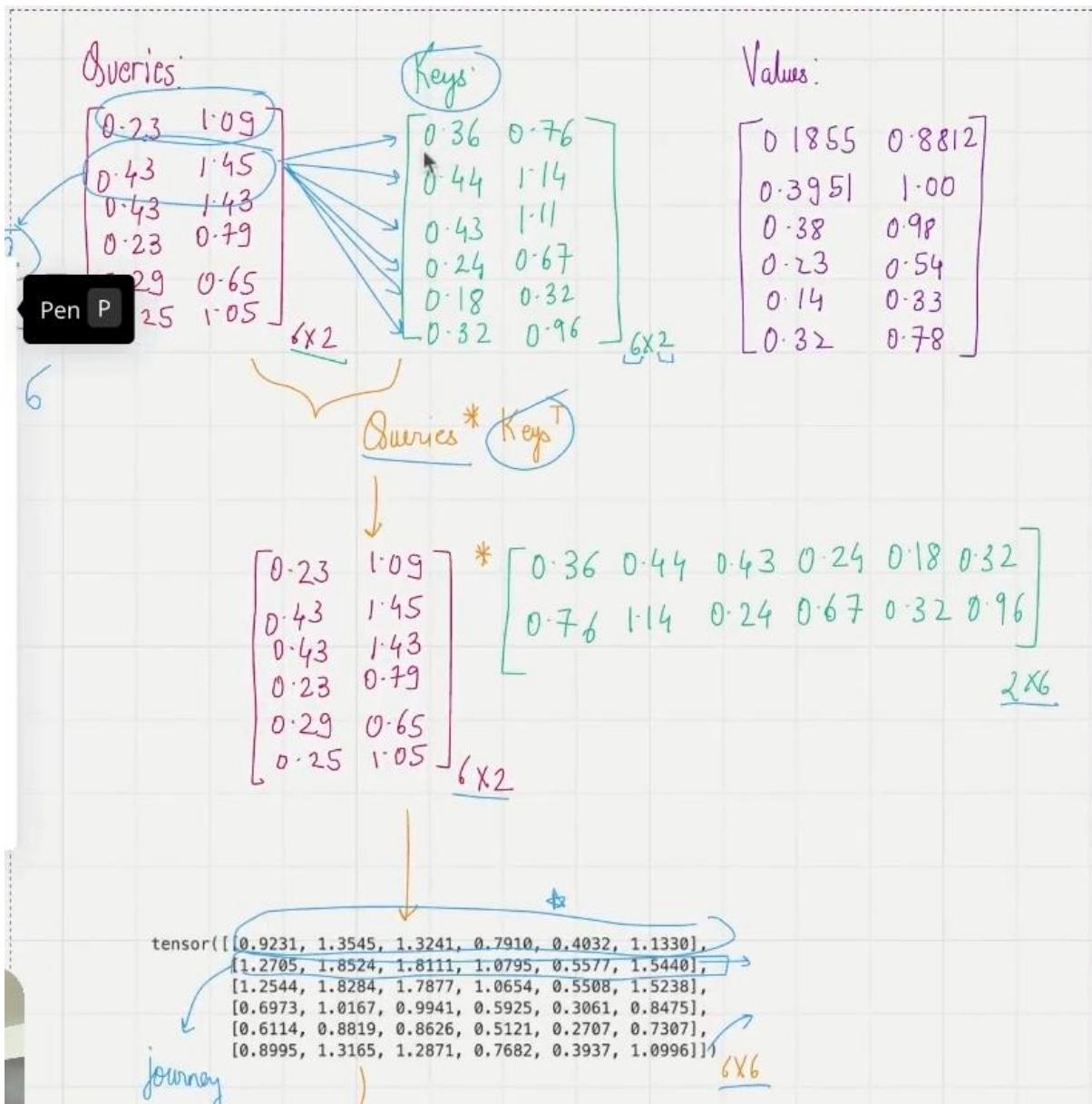
6×6

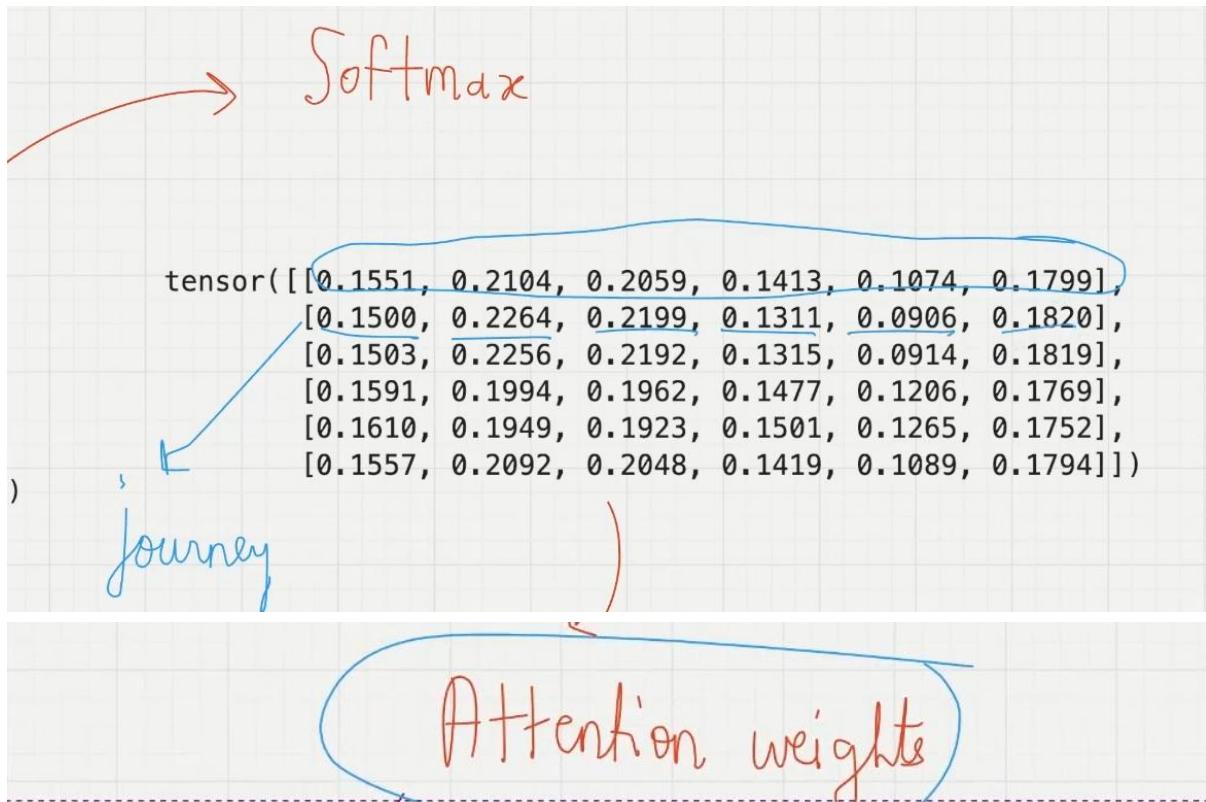
Attention scores

③ Attention weights

Scale by $\sqrt{d_{\text{keys}}}$

`tensor([[0.6528, 0.9578, 0.9363, 0.5593, 0.2851, 0.8011],
 [0.8984, 1.3098, 1.2806, 0.7633, 0.3944, 1.0918],
 [0.8870, 1.2929, 1.2641, 0.7534, 0.3895, 1.0775],
 [0.4930, 0.7189, 0.7029, 0.4190, 0.2164, 0.5993],
 [0.4323, 0.6236, 0.6099, 0.3621, 0.1914, 0.5167],
 [0.6361, 0.9309, 0.9101, 0.5432, 0.2784, 0.7776]])`





WHY DIVIDE BY SQRT (DIMENSION)

Reason 1: For stability in learning

The softmax function is sensitive to the magnitudes of its inputs. When the inputs are large, the differences between them become much more pronounced. This causes the softmax output to become "peaky," where the highest value reaches 1.0 and others receive very little.

In attention mechanisms, particularly in transformers, if the dot products between query and key vectors are large, the resulting attention scores can become very large. This results in a very sharp softmax distribution, making the model overly confident. Such distributions can make learning unstable,

- The **softmax function** is sensitive to the magnitudes of its inputs.
- When inputs are large, the exponential values become more pronounced.
- This causes the softmax output to become "**peaky**," where one value dominates almost all the probability mass, and others receive very little.
- In **attention mechanisms** (especially in transformers), large dot products between query and key vectors can produce large attention scores.
- Large scores lead to **sharp softmax distributions**, making the model **overly confident** in a single key.
- This can cause **training instability**.
- Dividing by the **square root of the dimension** normalizes these scores.
- This normalization helps prevent scores from becoming **excessively large** and results in **more stable learning**.

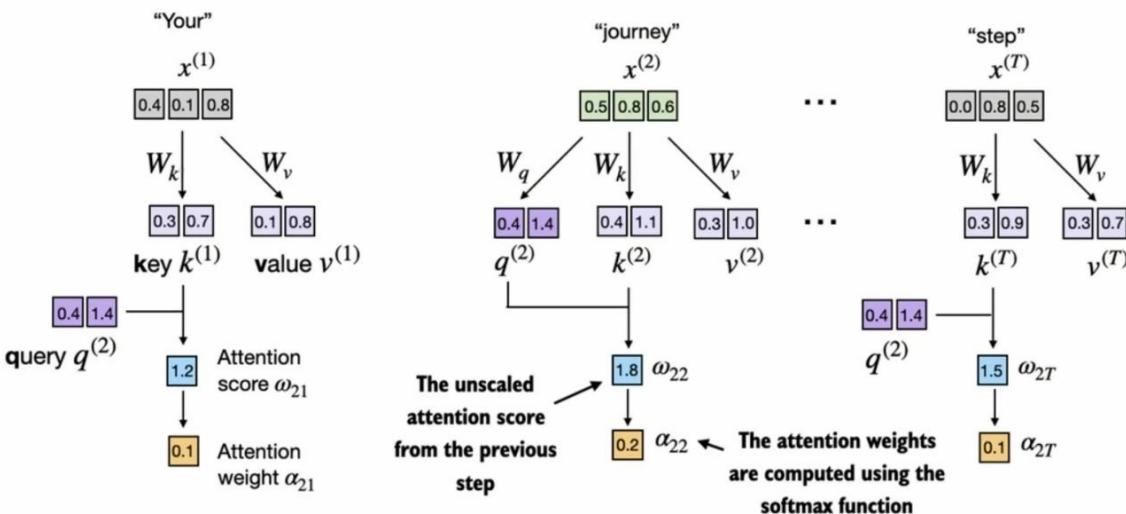
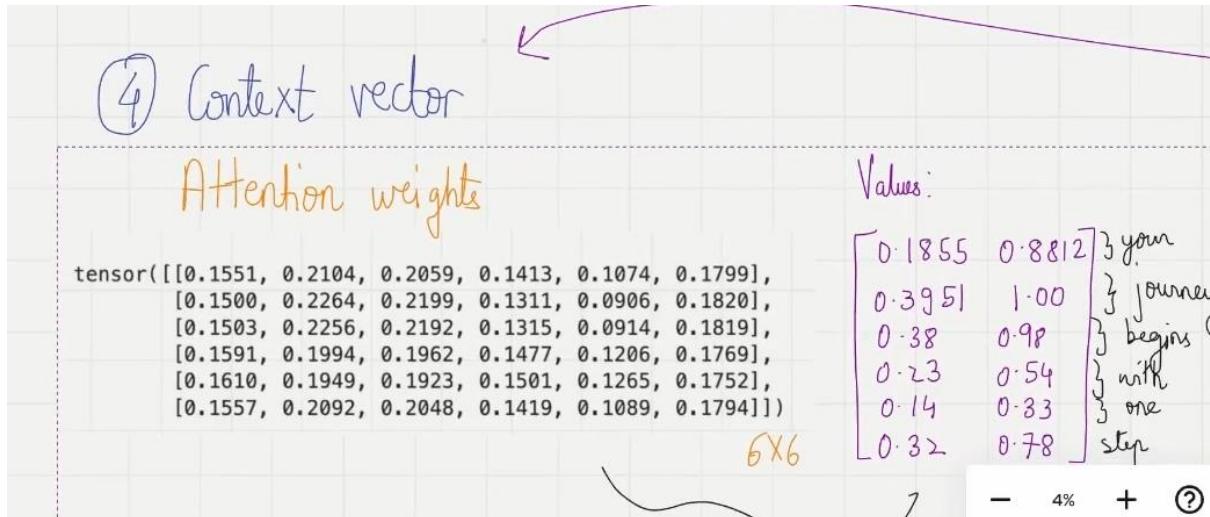
BUT WHY SQRT?

Reason 2: To make the variance of the dot product stable

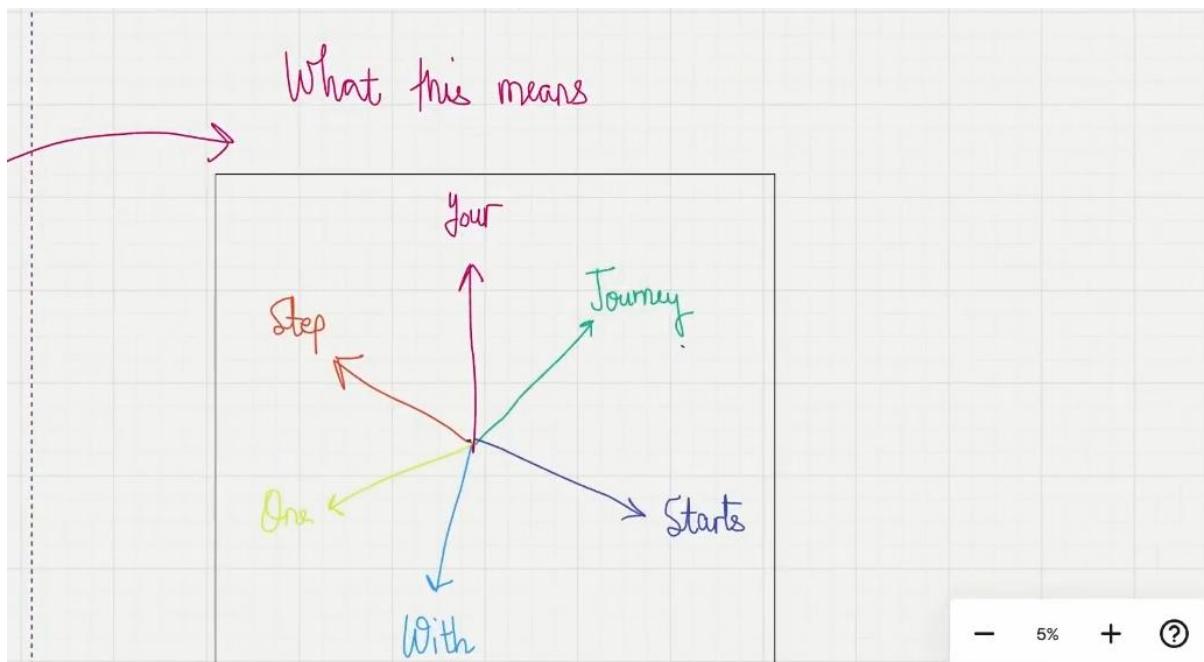
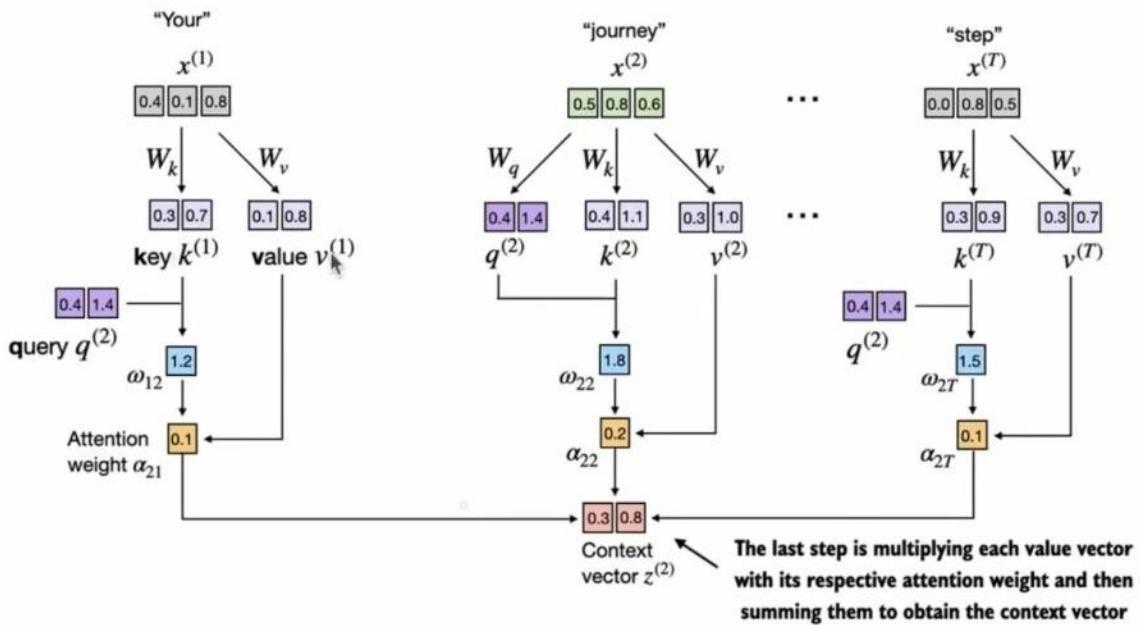
The dot product of Q and K increases the variance because multiplying two random numbers

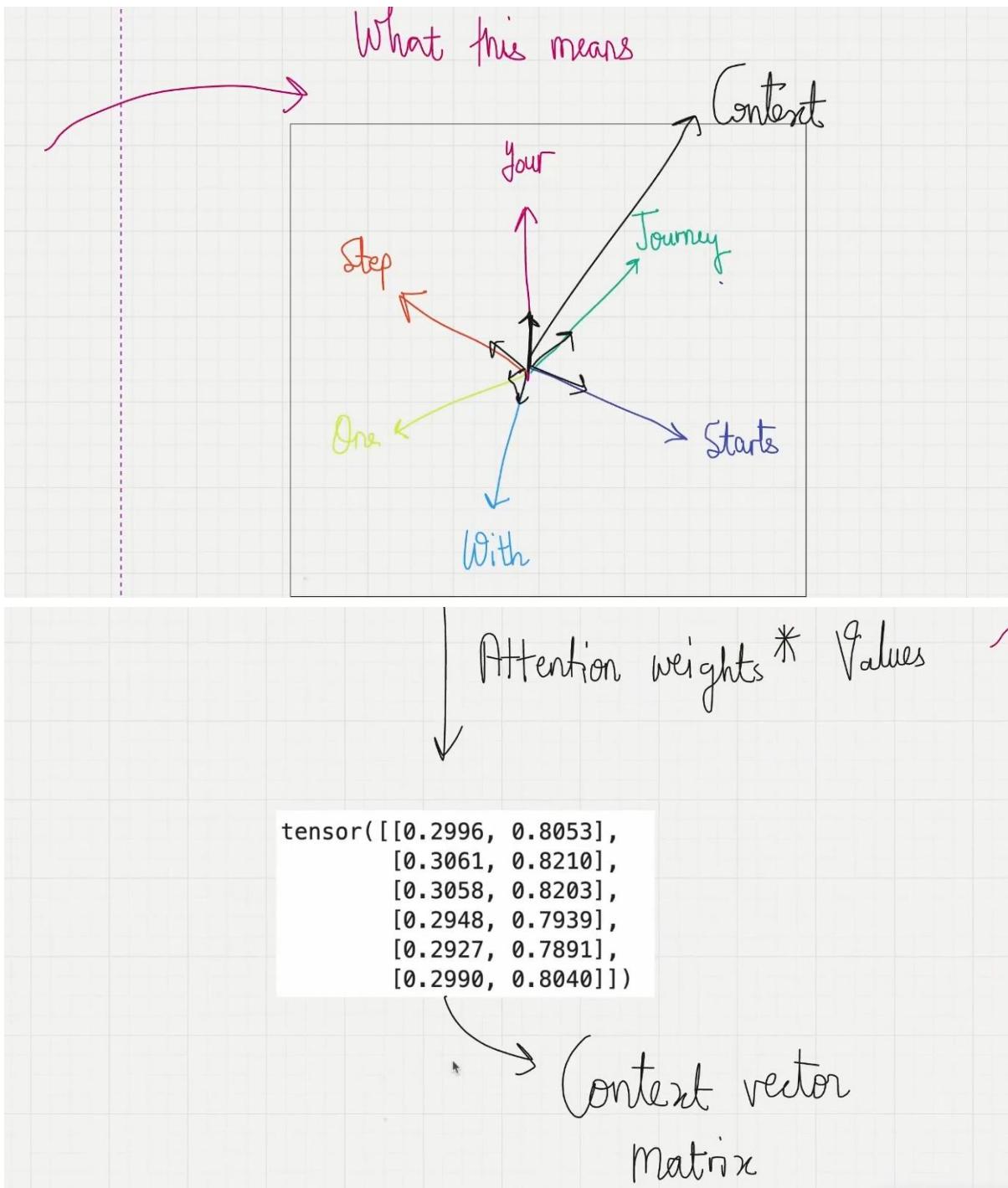
The increase in variance grows with the dimension.

Dividing by sqrt (dimension) keeps the variance close to 1



Final step: Compute context vectors





* Why do we use the terms: key, query and value?

Query: Analogous to search query in a database.

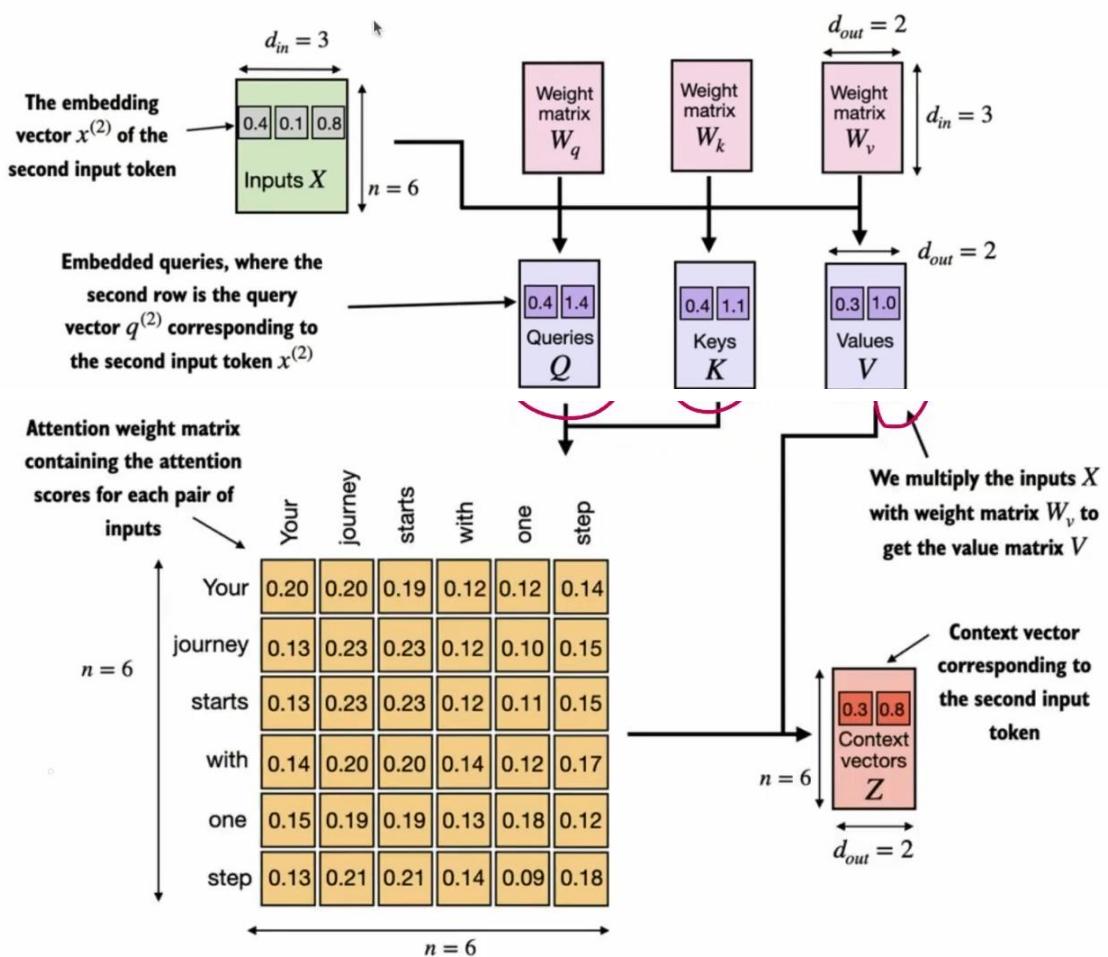
It represents the current token the model

focuses on

Key: In attention mechanism, each item in input sequence has a key. Keys are used to match with the query.

Value: It represents the actual context or representation of the input items. Once the model determines which keys (which parts of the input) are most relevant to the query (current focus item), it retrieves the corresponding values.

⑥ Implementing a compact self attention Python class



* Self attention involves the trainable

Weight matrices: W_q , W_k and W_v

* These matrices transform input data into queries, keys and values; which are crucial components of the attention mechanism.

* Why do we use the terms: key, query and value?

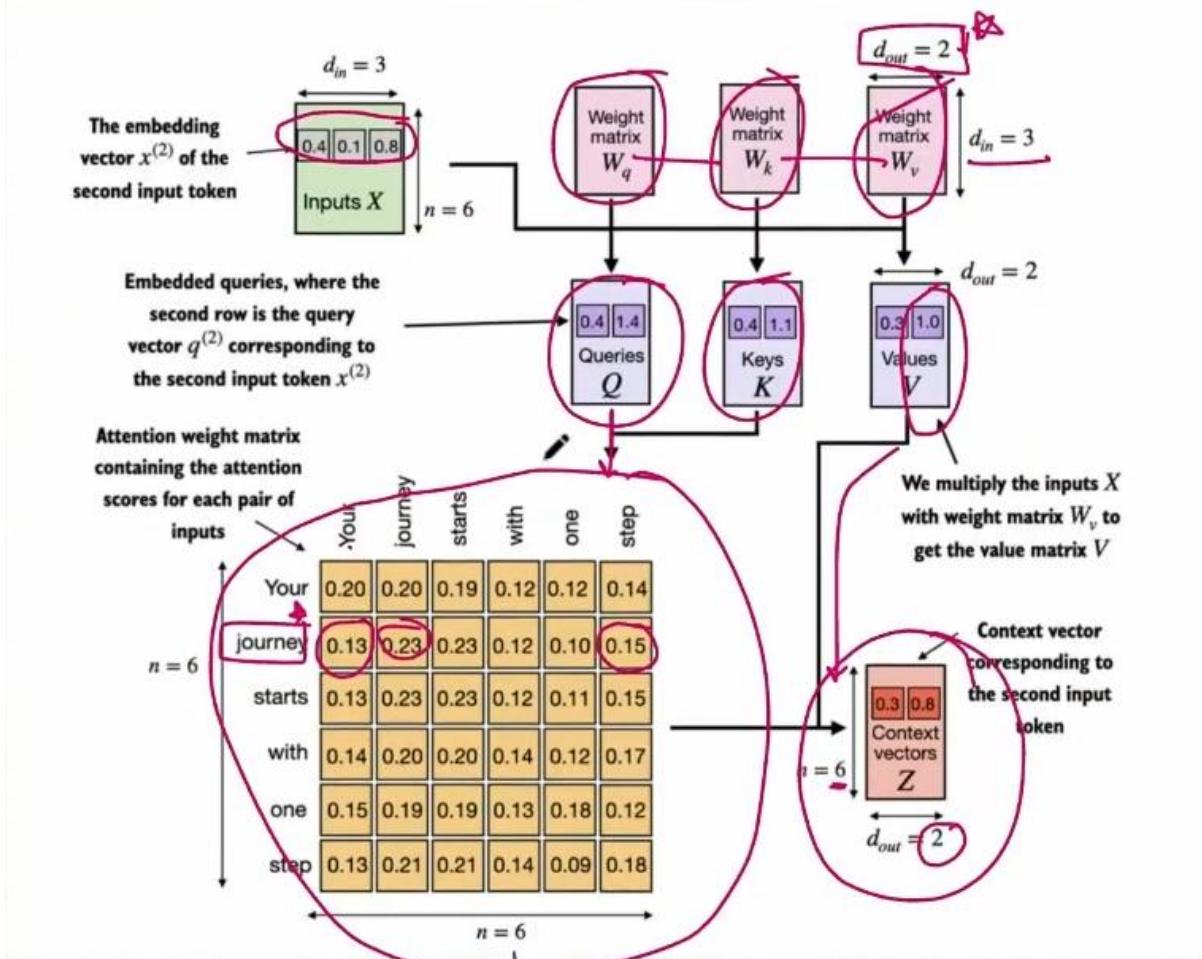
Query: Analogous to search query in a database.

It represents the current token the model focuses on

Key: In attention mechanism, each item in input sequence has a key. Keys are used to match with the query

Value: It represents the actual content or representation of the input items.

Once the model determines which keys (which parts of the input) are most relevant to the query (current focus item), it retrieves the corresponding values.



* Self attention involves the trainable

Weight matrices: W_q , W_k and W_v

* These matrices transform input data into queries, keys and values; which are crucial components of the attention mechanism.

In the next section, we will look at:

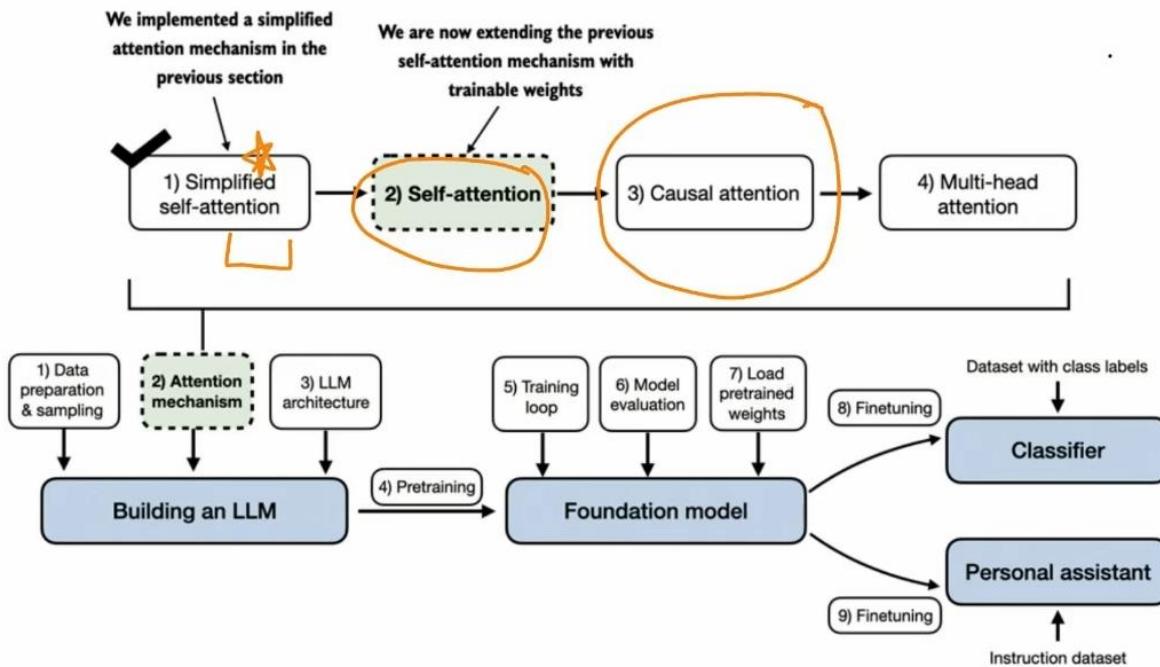
Causal attention
/

Multi head
attention

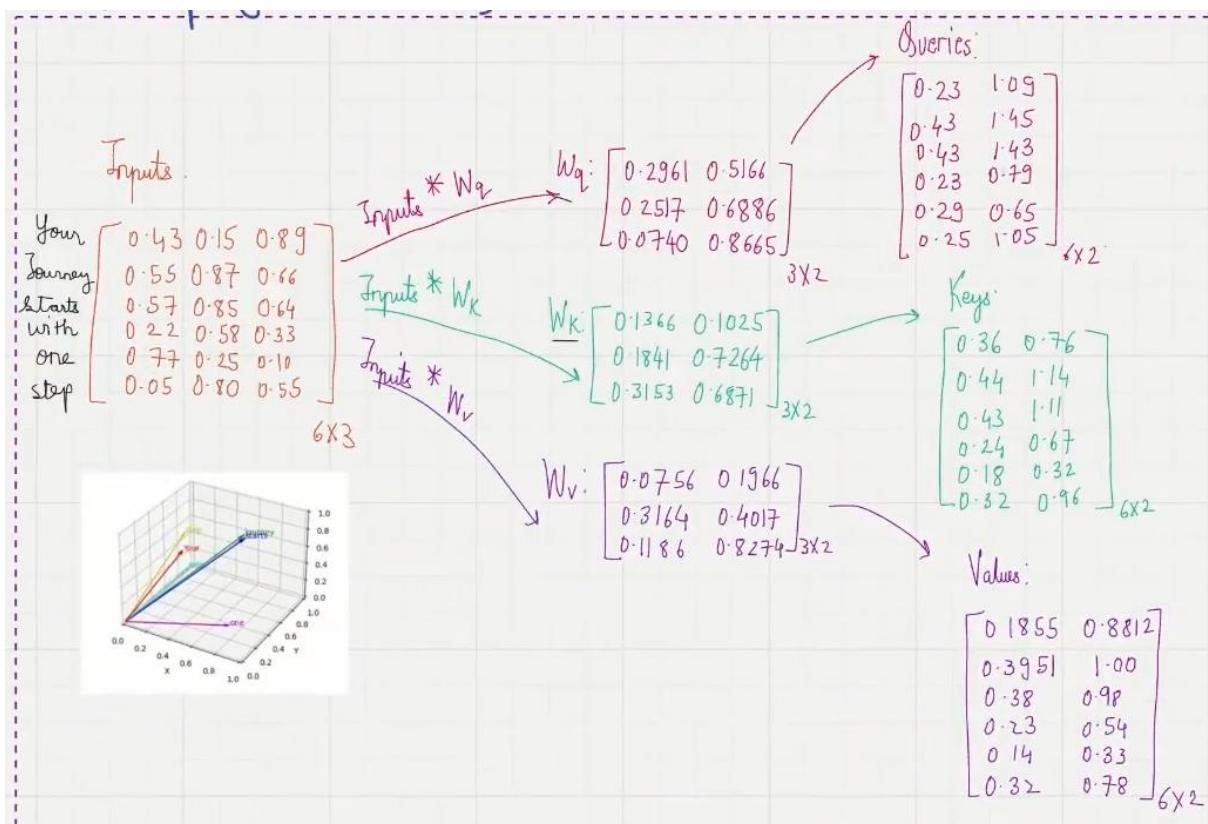
Modify attention
mechanism → prevents
model from accessing
future information in a
sequence, which is crucial
for tasks like language
modeling.

↓
Split the attention
mechanism into
multiple "heads".
Each head learns
different aspects of
the data → improves model
performance in complex tasks.

Lecture → Causal attention



Recap (until now)



$\rightarrow \text{Queries}^* \text{Keys}^T$

```
tensor([[0.9231, 1.3545, 1.3241, 0.7910, 0.4032, 1.1330],  
       [1.2705, 1.8524, 1.8111, 1.0795, 0.5577, 1.5440],  
       [1.2544, 1.8284, 1.7877, 1.0654, 0.5508, 1.5238],  
       [0.6973, 1.0167, 0.9941, 0.5925, 0.3061, 0.8475],  
       [0.6114, 0.8819, 0.8626, 0.5121, 0.2707, 0.7307],  
       [0.8995, 1.3165, 1.2871, 0.7682, 0.3937, 1.0996]])
```

6x6

Attention scores

\rightarrow Scaling by $\sqrt{d_k}$ -keys
 $+ \text{Softmax}$ \rightarrow Attention weights^{*}
Values

```
tensor([[0.1551, 0.2104, 0.2059, 0.1413, 0.1074, 0.1799],  
       [0.1500, 0.2264, 0.2199, 0.1311, 0.0906, 0.1820],  
       [0.1503, 0.2256, 0.2192, 0.1315, 0.0914, 0.1819],  
       [0.1591, 0.1994, 0.1962, 0.1477, 0.1206, 0.1769],  
       [0.1610, 0.1949, 0.1923, 0.1501, 0.1265, 0.1752],  
       [0.1557, 0.2092, 0.2048, 0.1419, 0.1089, 0.1794]])
```

```
tensor([[0.2996, 0.8053],  
       [0.3061, 0.8210],  
       [0.3058, 0.8203],  
       [0.2948, 0.7939],  
       [0.2927, 0.7891],  
       [0.2990, 0.8040]])
```

Attention weights

Content Vector
Matrix

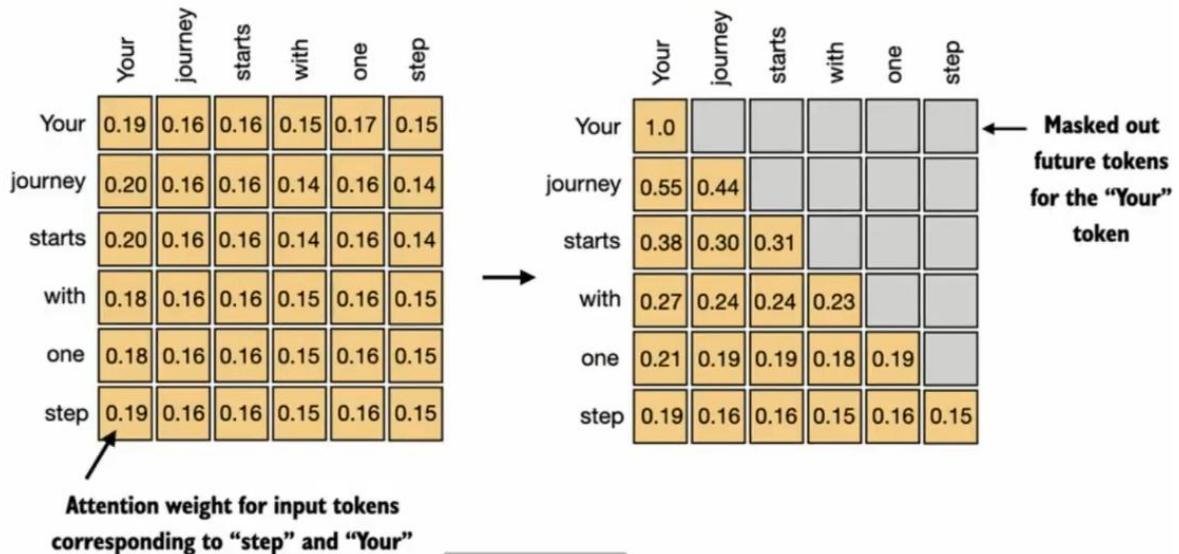
① Causal attention, also known as masked attention
is a special form of self attention.

② It restricts the model to only consider previous
and current inputs in a sequence, when
processing any given token.

③ This is in contrast to the self attention mechanism,
which allows access to the entire input sequence

④ When computing attention scores, the causal attention
mechanism ensures that the model only factors in
tokens that occur at or before the current token
in the sequence.

⑤ To achieve this in GPT like LLMs, for each
token processed, we mask out the future tokens,
which come after the current token in the
input text.

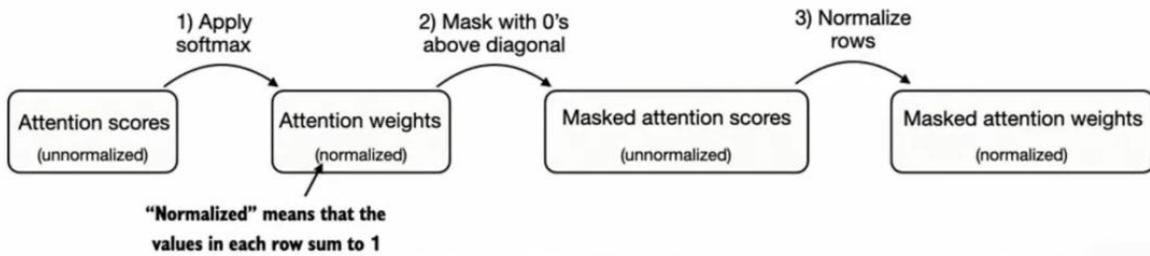


↓ We mask out the attention weights above the diagonal, and we normalize the non masked attention weights, such that the attention weights sum upto 1 in each row.

* Applying a Causal Attention Mask *

Strategy: Get attention weights

↓ Zero out elements above the diagonal and normalize the resulting matrix.



More efficient
way

Attention → Upper Triangular → Softmax
scores infinity mask

Masking additional attention weights with dropout *

↳ Dropout is a deep learning technique where randomly selected hidden layer units are ignored during training

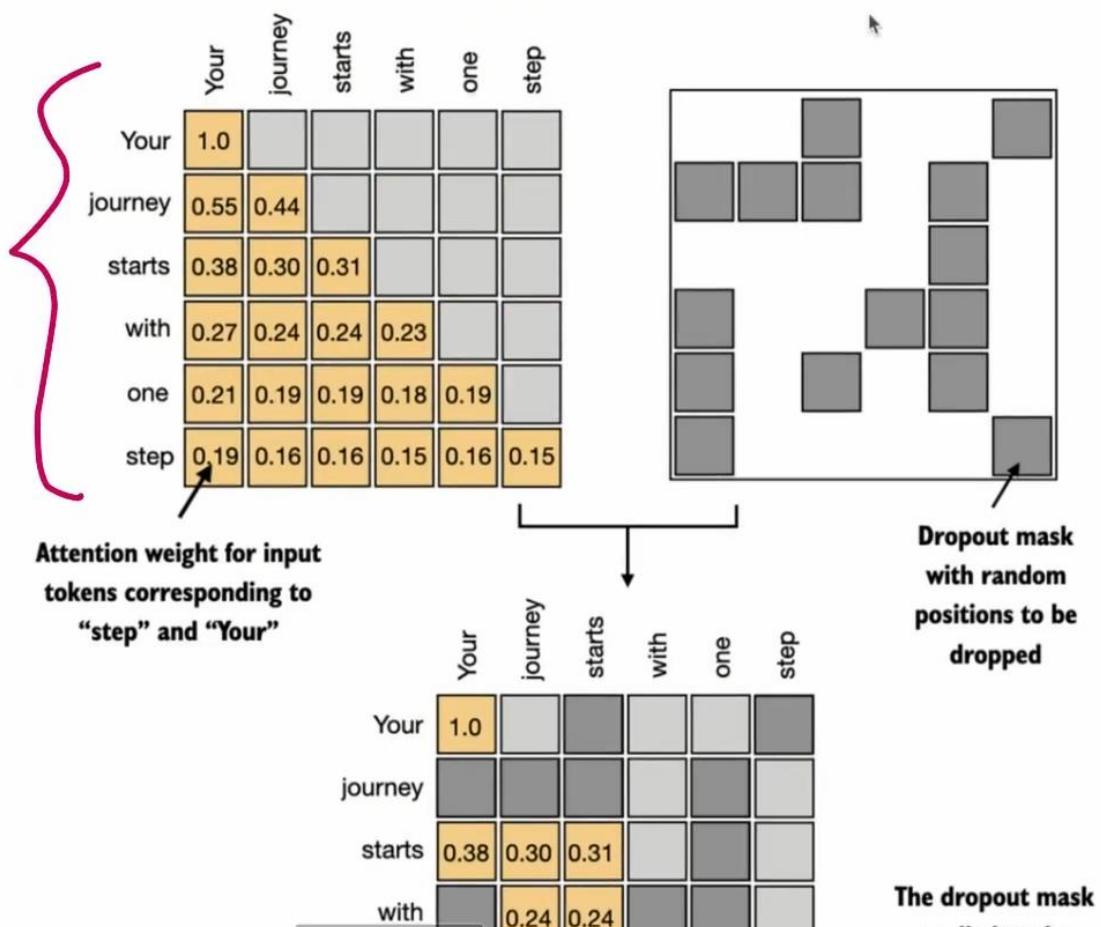
This prevents overfitting and improves generalization performance.

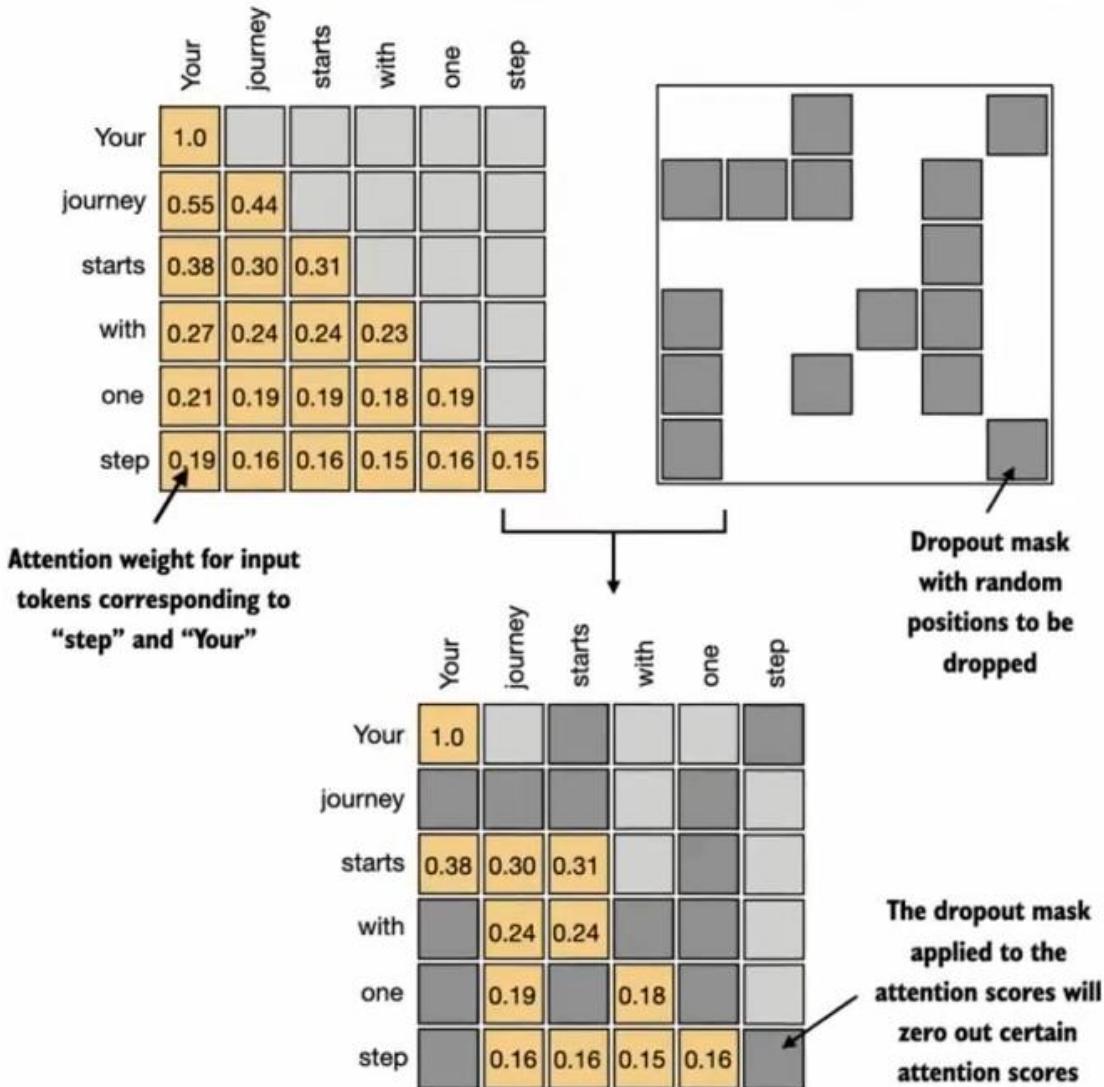
↳ In transformer architecture, including models like GPT, dropout in the attention mechanism is applied in 2 specific areas

after calculating
attention scores

after applying attention
weights to value vectors

→ Applying dropout after calculating attention
weights is more common, and we consider that.





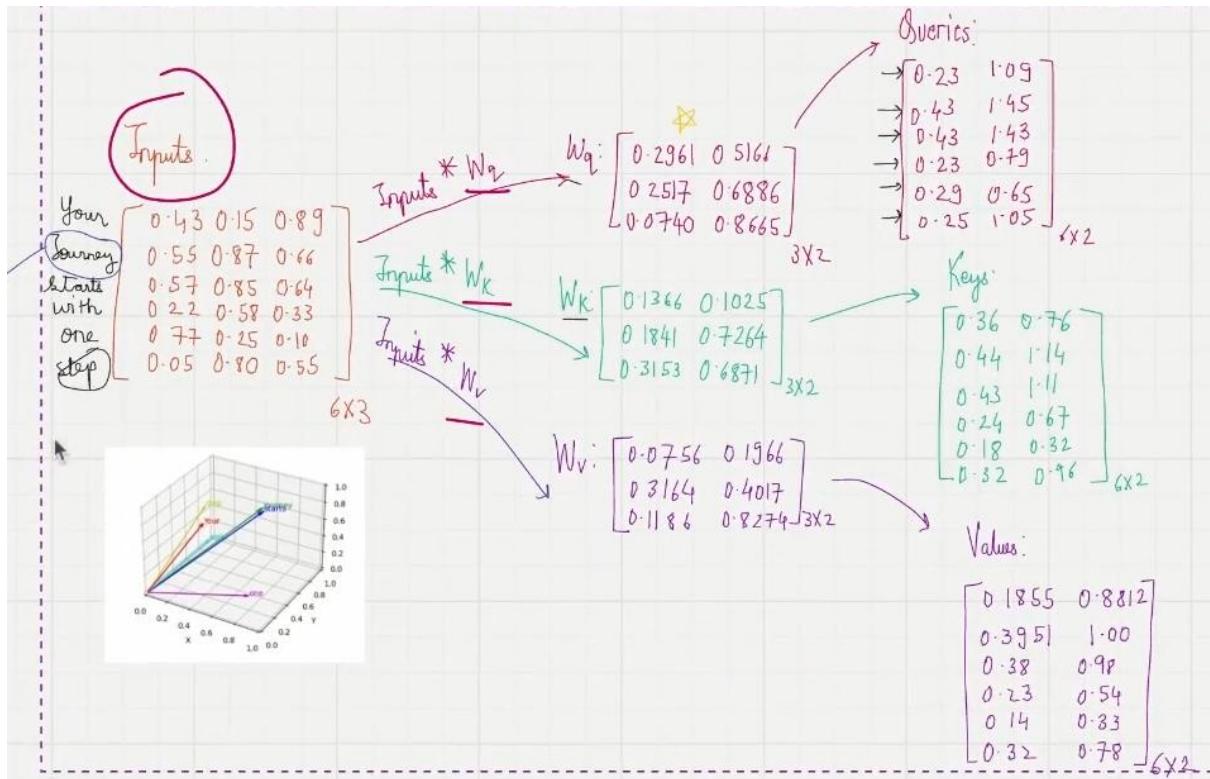
* Now, we can implement a Causal Attention class, which incorporates Causal Attention and Dropout modifications into the Self Attention class we implemented earlier.

MASKING ADDITIONAL ATTENTION WEIGHTS WITH DROPOUT

In the following code example, we use a dropout rate of 50%, which means masking out half of the attention weights.

When we train the GPT model in later chapters, we will use a lower dropout rate, such as 0.1 or 0.2.

In the following code, we apply PyTorch's dropout implementation first to a 6x6 tensor consisting of ones for illustration.



Batch 1

→ Queries * Keys^T → Scaling by $\sqrt{d_k}$ -keys + Dropout + softmax → Attention weight Values

```
tensor([[0.1555, 1.3545, 1.3241, 0.7910, 0.4032, 1.1330],
       [1.2785, 1.8929, 0.8111, 1.0795, 0.5577, 1.5440],
       [1.2344, 1.8284, 1.1054, 0.5568, 1.5320],
       [0.8462, 0.9841, 0.5924, 0.4874, 0.4872],
       [0.6114, 0.8819, 0.8826, 0.5121, 0.2707, 0.3071],
       [0.8995, 1.3165, 1.2871, 0.7682, 0.3937, 1.0399]]).t
```

Attention scores

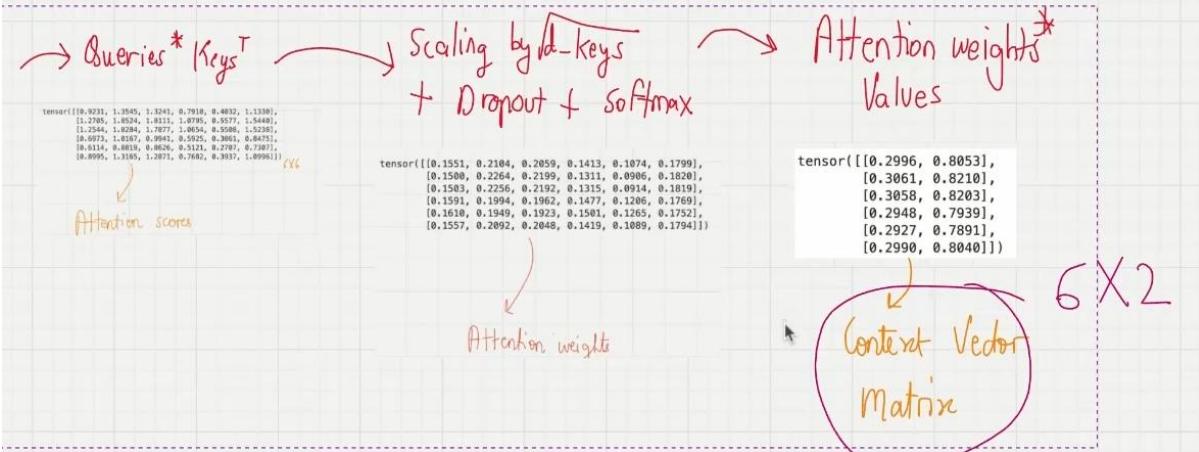
```
tensor([[0.1551, 0.2104, 0.2059, 0.1413, 0.1074, 0.1799],
       [0.1500, 0.2264, 0.2199, 0.1311, 0.0906, 0.1820],
       [0.1503, 0.2256, 0.2192, 0.1315, 0.0914, 0.1819],
       [0.1591, 0.1994, 0.1962, 0.1477, 0.1206, 0.1769],
       [0.1610, 0.1949, 0.1923, 0.1501, 0.1265, 0.1752],
       [0.1557, 0.2092, 0.2048, 0.1419, 0.1089, 0.1794]])
```

```
tensor([[0.2996, 0.8053],
       [0.3061, 0.8210],
       [0.3058, 0.8203],
       [0.2948, 0.7939],
       [0.2927, 0.7891],
       [0.2990, 0.8040]])
```

Attention weights

Context Vector Matrix

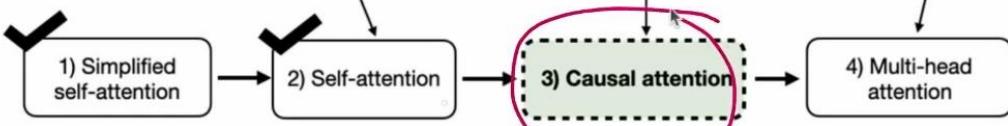
Batch 2



In the previous section, we implemented a self-attention mechanism with trainable weights

In this section, we extended the self-attention mechanism with a causal mask and dropout mask

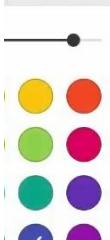
In the next section, we extend causal attention to multi-head attention



What all we have achieved so far.

Lecture → Multi-head attention mechanism

The term "multi-head" refers to dividing the attention mechanism into multiple heads.

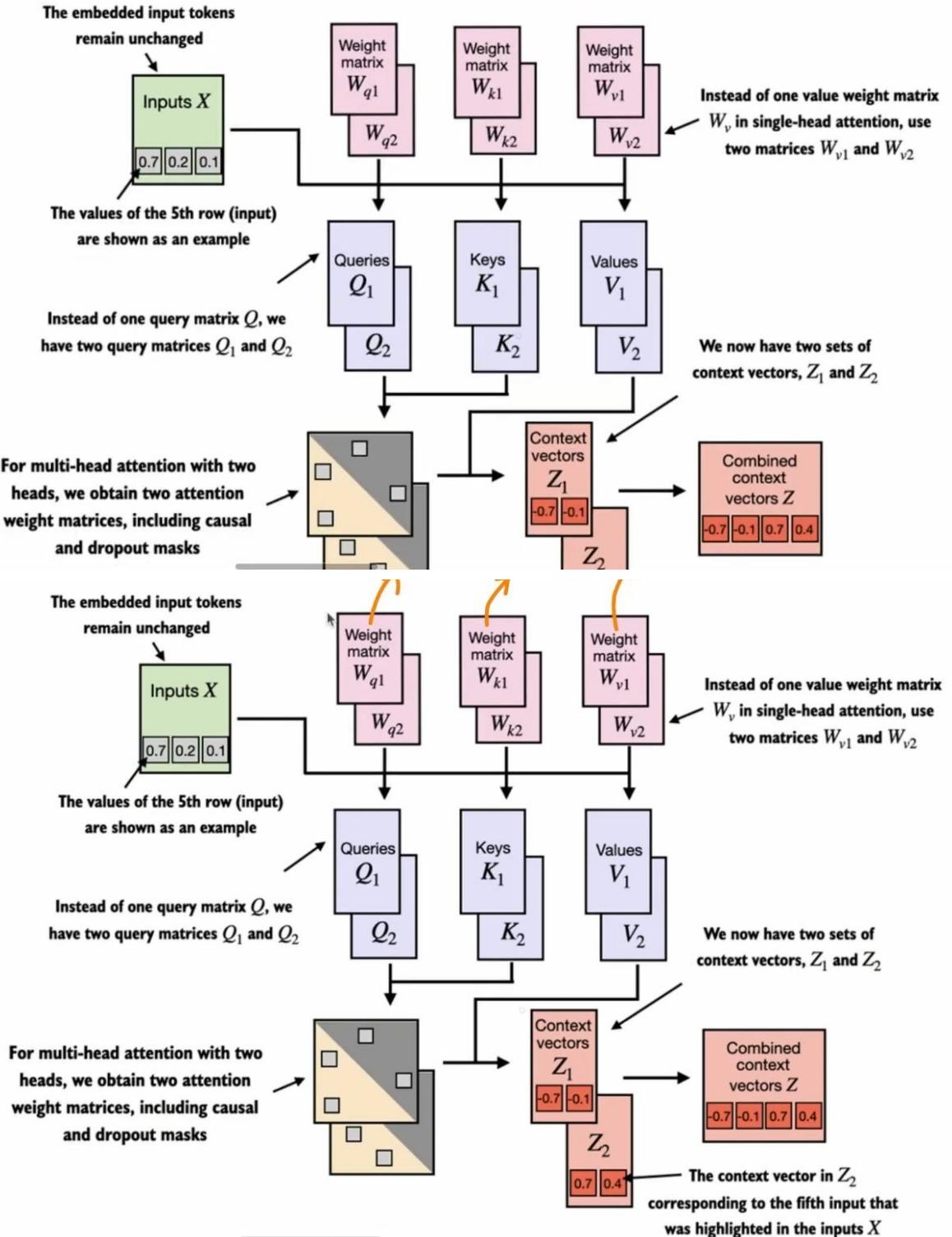


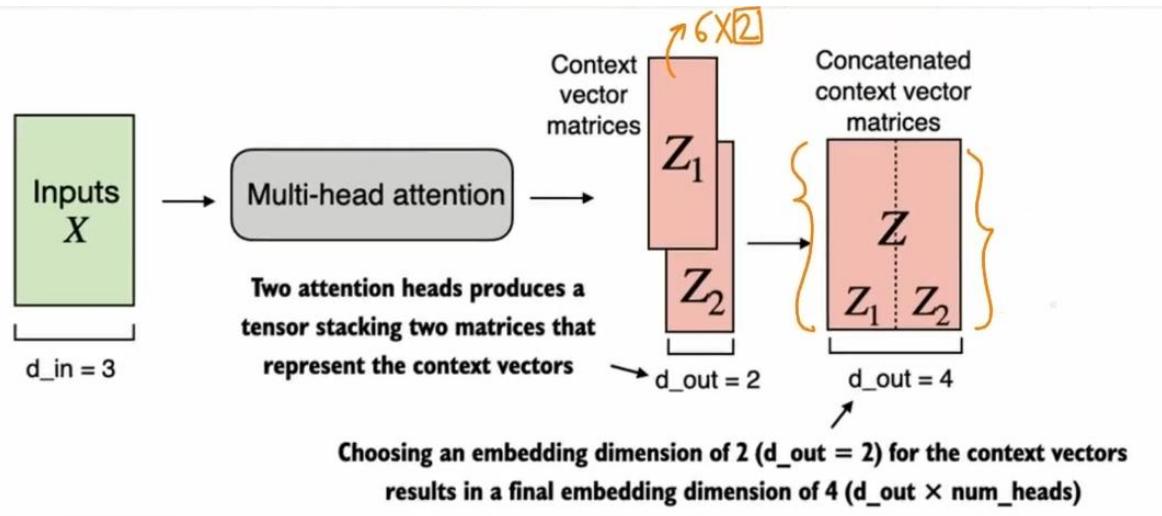
each operating independently

(a) Stacking multiple single head attention layers

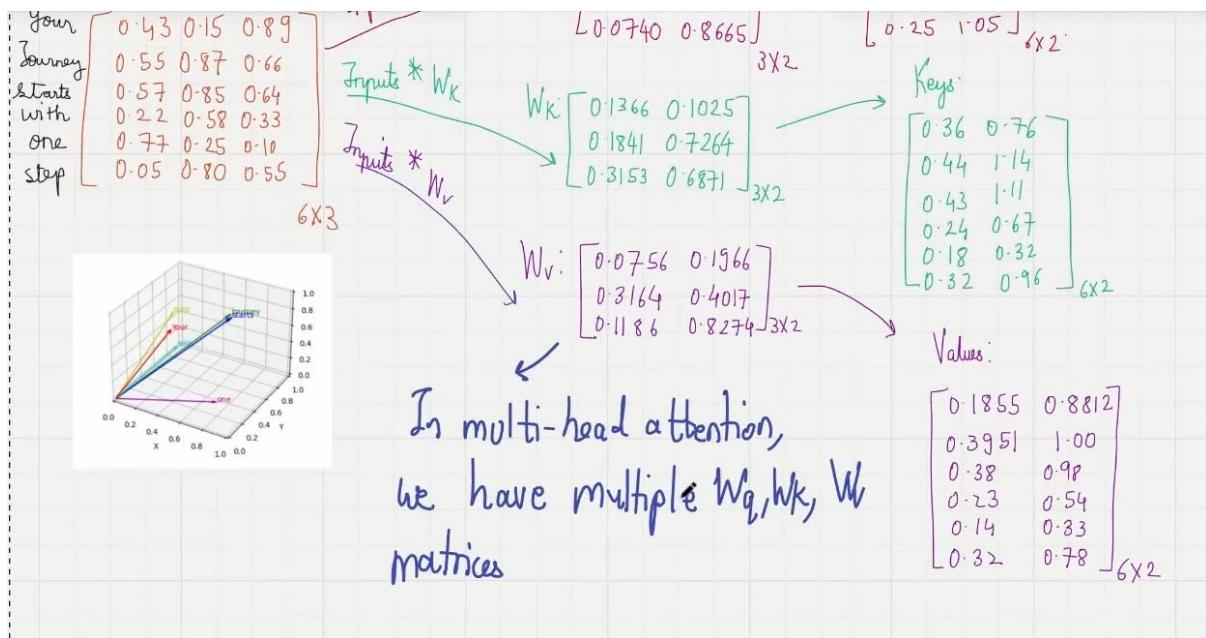
① Implementing multi-head attention involves creating multiple instances of the self-attention mechanism; each with its own weights; and then combining their outputs.

② This can be computationally intensive, but it makes LLMs powerful at complex pattern recognition tasks.





↳ Main idea: Run the attention mechanism multiple times (in parallel) with different, learned linear projections: the results of multiplying input data (like query, key value vectors) by a weight matrix.



(b) Implementing multi-head attention with weight splits

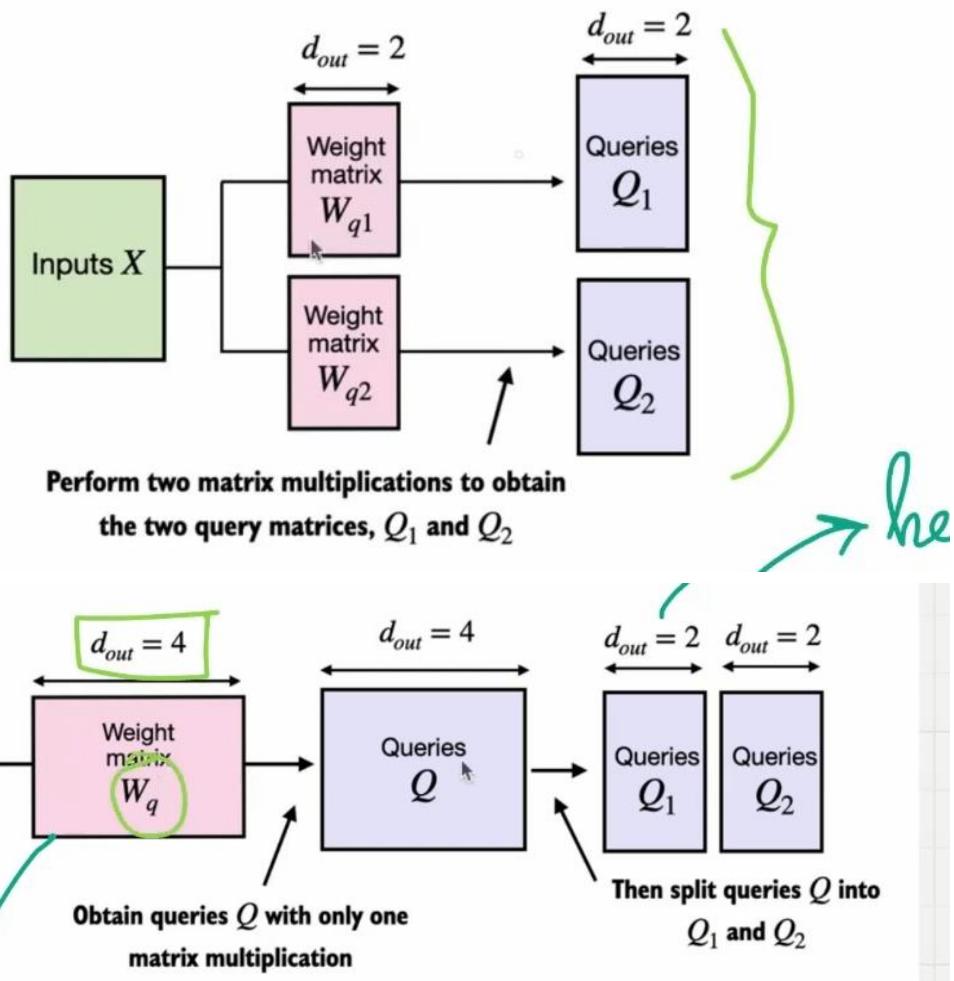
(i) Instead of maintaining 2 separate classes:
 MultiHeadAttentionWrapper and CausalAttention,
 we combined both of these into a single
 MultiHeadAttention class.

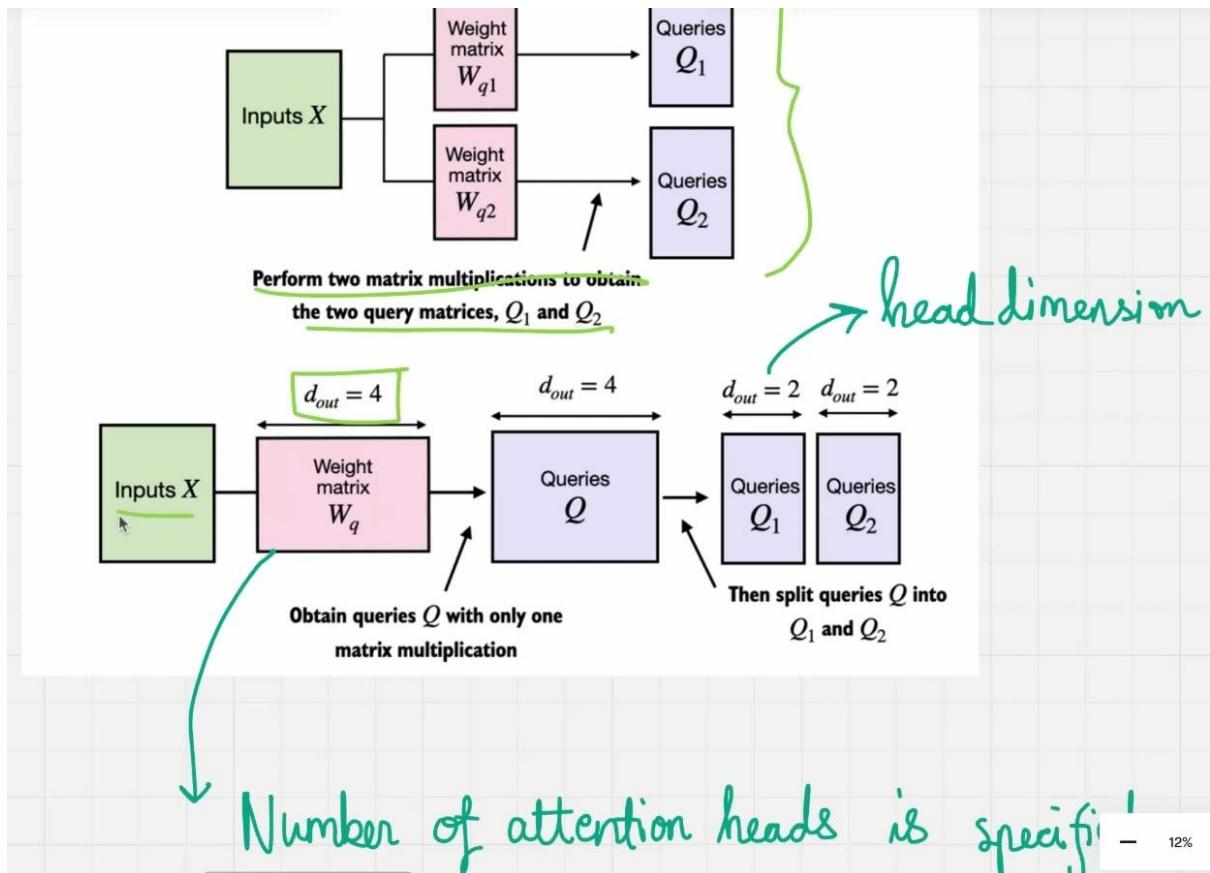
Implementing multi-head attention with weight splits

(i) Instead of maintaining 2 separate classes:

MultiHeadAttentionWrapper and Causal Attention
we combined both of these into a single

MultiHeadAttention class.





IMPLEMENTING MULTI-HEAD ATTENTION WITH WEIGHT SPLITS

Instead of maintaining two separate classes, `MultiHeadAttentionWrapper` and `CausalAttention`, we can combine both of these concepts into a single `MultiHeadAttention` class.

Also, in addition to just merging the `MultiHeadAttentionWrapper` with the `CausalAttention` code, we will make some other modifications to implement multi-head attention more efficiently.

In the `MultiHeadAttention` class, multiple heads are implemented by creating a list of `CausalAttention` objects (`self.heads`), each representing a separate attention head.

The `CausalAttention` class independently performs the attention mechanism, and the results from each head are concatenated.

In contrast, the following `MultiHeadAttention` class integrates the multi-head functionality within a single class.

It splits the input into multiple heads by reshaping the projected query, key, and value tensors and then combines the results from these heads after computing attention.

Number of attention heads is specified.
 Head dimension = d_{out}
 n_{heads}

* Let us take a simple example:

Step ①: Start with the input

b, num_tokens, d-in = (1, 3, 6)

the
cat
sleeps

```
...  
  
x = torch.tensor([[[1.0, 2.0, 3.0, 4.0, 5.0, 6.0],  
                  [6.0, 5.0, 4.0, 3.0, 2.0, 1.0],  
                  [1.0, 1.0, 1.0, 1.0, 1.0, 1.0]]])
```

batch = 1,
num_tokens = 3, d-in = 6

Step ②: Decide d-out, num-heads

↓ ↓
6 2

Step ②: Decide d-out, num-heads
↓ ↓
6 2 head_dim = $\frac{6}{2} = 3$

Step ③: Initialize trainable weight matrices

for key, query, value (W_k, W_q, W_v)
 (6×6)

$(6 \times 6) \rightarrow d_{in} \times d_{out}$

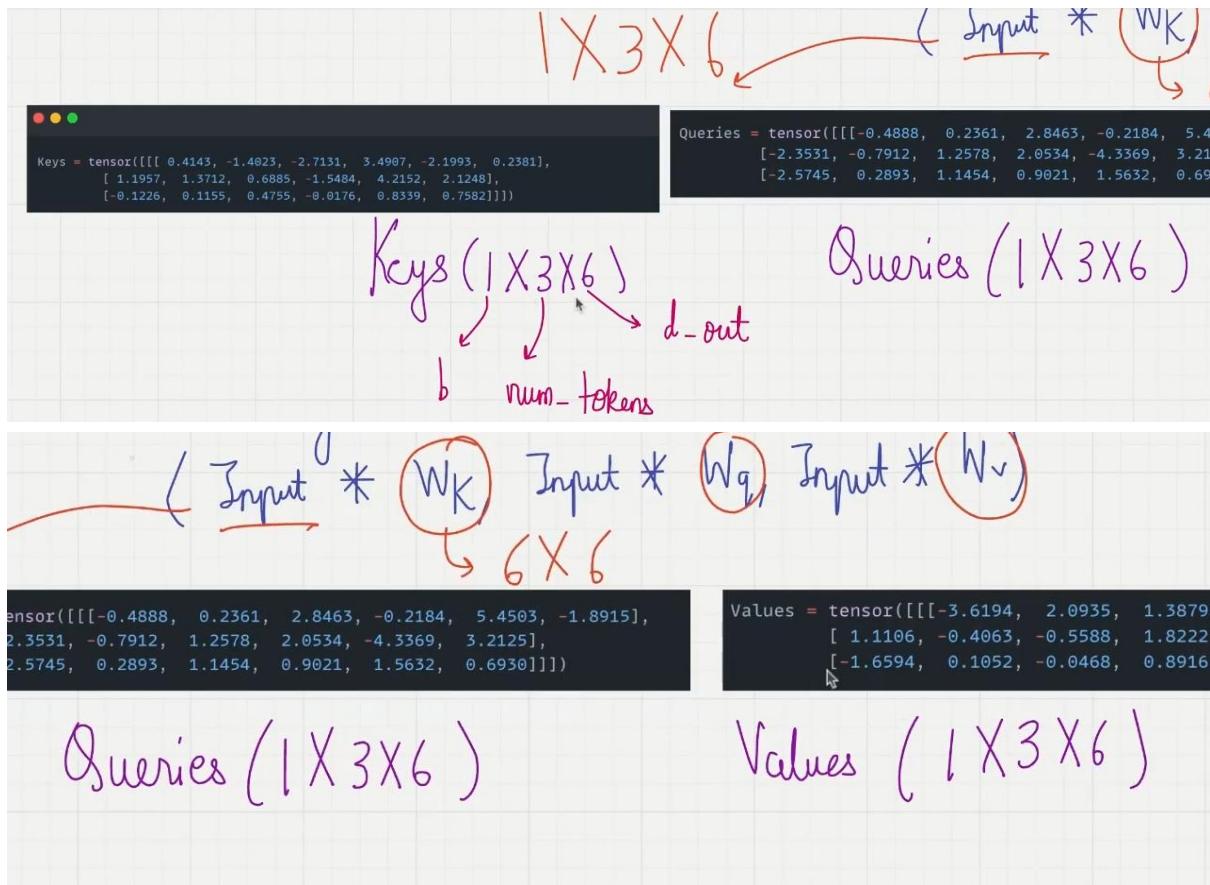
```
Wq = tensor([[ 0.6323, -0.2366,  1.2455,  0.3465,  1.2458,  0.3229],
             [ 0.6571, -0.2378, -0.5311, -0.2610, -1.4819, -1.6418],
             [-0.2990,  0.4216,  0.2114, -0.0271, -0.5682,  0.6937],
             [-1.1291, -1.0102,  0.6946,  0.1094,  0.5130, -0.8669],
             [ 0.3480,  0.2593,  0.4412,  1.0017, -0.3913, -0.2878],
             [ 0.2484,  0.2846, -0.3386, -0.6164,  1.2722,  0.5754]])
```

```
Wk = tensor([[-0.3703,  0.5431, -0.0372, -0.4406,  0.4103, -0.1773],
            [ 1.5993, -0.2777, -1.1909, -0.4301,  0.6927, -1.3304],
            [ 1.2470, -0.1872, -0.1670,  1.4302,  1.2927,  0.4822],
            [-0.0984, -0.8983,  0.3334, -0.6312,  0.1022, -1.0715],
            [-0.7647,  0.1734,  0.6305,  1.0155,  0.8474,  0.1454],
            [-1.5085, -0.4529,  0.0997, -0.1084,  0.8046,  0.3459]])
```

(6^{10})

```
Wv = tensor([[ 1.6395,  1.1234, -0.1001,  0.5021, -1.0590,  0.1412,
              [-0.4271,  0.5681,  0.4164, -1.2534,  1.3061,  0.3610],
              [-0.2824, -0.4314,  1.2358,  0.1181, -1.2467,  0.1893],
              [ 1.3440,  0.1487, -0.6174,  0.8890, -0.3282,  1.4662],
              [ 0.1814, -0.4761, -0.0402,  0.7326,  0.7654, -0.1080],
              [-0.8974,  0.6786,  0.5602, -0.2443, -0.4883,  1.3996]]))
```

Step ④: Calculate Keys, Queries, Value Matrix
 $(Input * W_k, Input * W_q, Input * W_v)$



Step ⑤: Unroll last dimension of Keys, Queries and Values
to include num_heads and head_dim

$$\text{head_dim} = d\text{-out} / \text{num_heads} = 6 / 2 = 3$$

$$(b, \text{num_tokens}, d\text{-out}) \rightarrow (b, \text{num_tokens}, \text{num_heads}, \text{head_dim})$$

$$(1, 3, 6) \rightarrow (1, 3, 2, 3)$$

$$(b, \text{num_tokens}, d\text{-out}) \rightarrow (b, \overset{2}{\text{num_tokens}}, \text{num_heads}, \text{head_dim})$$

$$(1, 3, 6) \rightarrow (1, 3, 2, 3)$$

```

Reshaped Queries matrix:
tensor([[[[-0.4888,  0.2361,  2.8463],
        [-0.2184,  5.4503, -1.8915], ],
       [[-2.3531, -0.7912,  1.2578],
        [ 2.0534, -4.3369,  3.2125]], ],
      [[-2.5745,  0.2893,  1.1454],
       [ 0.9021,  1.5632,  0.6930]]]], dtype=torch.float64)

```

$(\overbrace{1}, \overbrace{3}, \overbrace{2}, \overbrace{3})$

```

Reshaped Keys matrix:
tensor([[[[ 0.4143, -1.4023, -2.7131],
        [ 3.4907, -2.1993,  0.2381]], ,
       [[ 1.1957,  1.3712,  0.6885],
        [-1.5484,  4.2152,  2.1248]], ,
       [[-0.1226,  0.1155,  0.4755],
        [-0.0176,  0.8339,  0.7582]]]], dtype=torch.float64)

```

$(\overbrace{1}, \overbrace{3}, \overbrace{2}, \overbrace{3})$

```

Reshaped Values matrix:
tensor([[[[-3.6194,  2.0935,  1.3879],
        [ 2.1231, -1.2262, -0.2556]], ,
       [[ 1.1957,  1.3712,  0.6885],
        [-1.5484,  4.2152,  2.1248]], ,
       [[-0.1226,  0.1155,  0.4755],
        [-0.0176,  0.8339,  0.7582]]]], dtype=torch.float64)

```

$(\overbrace{1}, \overbrace{3}, \overbrace{2}, \overbrace{3})$

Step⑥: Group matrices by "number of heads"

$(b, \text{num_tokens}, \text{num_heads}, \text{head_dim}) \rightarrow (b, \text{num_heads}, \text{num_tokens}, \text{head_dim})$

$(\overbrace{1}, \overbrace{3}, \overbrace{2}, \overbrace{3}) \rightarrow (\overbrace{1}, \overbrace{2}, \overbrace{3}, \overbrace{3})$

(1, 2, 3, 3)

Transposed Queries matrix:

```
tensor([[[[-0.4888,  0.2361,  2.8463],
          [-2.3531, -0.7912,  1.2578],
          [-2.5745,  0.2893,  1.1454]]])
```

→ Head 1


```
[[[-0.2184,  5.4503, -1.8915],
  [ 2.0534, -4.3369,  3.2125],
  [ 0.9021,  1.5632,  0.6930]]])
```

→ Head 2

Transposed Keys matrix:

```
tensor([[[[ 0.4143, -1.4023, -2.7131],
          [ 1.1957,  1.3712,  0.6885],
          [-0.1226,  0.1155,  0.4755]],

         [[ 3.4907, -2.1993,  0.2381],
          [-1.5484,  4.2152,  2.1248],
          [-0.0176,  0.8339,  0.7582]]]])
```

→ Head 1

→ Head 2

Transposed Keys matrix:

```
tensor([[[[ 0.4143, -1.4023, -2.7131],
          [ 1.1957,  1.3712,  0.6885],
          [-0.1226,  0.1155,  0.4755]],

         [[ 3.4907, -2.1993,  0.2381],
          [-1.5484,  4.2152,  2.1248],
          [-0.0176,  0.8339,  0.7582]]]])
```

→ Head 1

→ Head 2

Transposed Values matrix:

```
tensor([[[[-3.6194,  2.0935,  1.3879],
          [ 1.1106, -0.4063, -0.5588],
          [-1.6594,  0.1052, -0.0468]],

         [[ 2.1231, -1.2262, -0.2556],
          [ 1.8222,  1.8721,  0.4929],
          [ 0.8916, -1.4384, -0.5651]]]])
```

→ Head 1

→ Head 2

Step 7: Find attention scores:

Queries * Keys Transpose (2, 3)

Queries * Keys Transpose (2, 3)

```
Queries matrix:  
tensor([[[[-0.4888,  0.2361,  2.8463],  
        [-2.3531, -0.7912,  1.2578],  
        [-2.5745,  0.2893,  1.1454]],  
  
       [[-0.2184,  5.4503, -1.8915],  
        [ 2.0534, -4.3369,  3.2125],  
        [ 0.9021,  1.5632,  0.6930]]])
```

```
Keys.transpose(2, 3) =  
tensor([[[[ 0.4143,  3.4907, -0.1226],  
        [-1.4023, -2.1993,  0.1155],  
        [-2.7131,  0.2381,  0.4755]],  
  
       [[ 1.1957, -1.5484, -0.0176],  
        [ 1.3712,  4.2152,  0.8339],  
        [ 0.6885,  2.1248,  0.7582]]])
```

Head 1
Head 2
(1, 2, 3, 3)

X
Head 2
(1, 2, 3, 3)

(b, num_heads, num_tokens, head_dim)

(b, num_heads, head_dim, num_tokens)

Attention Scores

```
queries * keys.transpose(2, 3):  
tensor([[[[-8.2252,  2.0863,  1.2960],  
        [-1.3722, -7.0720, -1.4290],  
        [-5.8961, -2.7236, -1.0160]],  
  
       [[ 4.6567, -7.1312,  4.7274],  
        [-1.3167,  1.3964, -0.6018],  
        [ 2.3820,  2.7213,  0.8448]]])
```

(1, 2, 3, 3)

(b, num_heads, num_tokens, num_tokens)

Step ⑧: Find attention weights



Mask attention scores to implement
causal attention

A
A
G
H

```
tensor([[[[-8.2252,      -inf,      -inf],  
         [-1.3722,      -7.0720,      -inf],  
         [-5.8961,      -2.7236,      -1.0160]],  
  
        [[ 4.6567,      -inf,      -inf],  
         [-1.3167,      1.3964,      -inf],  
         [ 2.3820,      2.7213,      0.8448]]]])
```

Divide by $\sqrt{\text{head_dim}} = \sqrt{d_{\text{out}}/\text{num_heads}} = \sqrt{6/2} = \sqrt{3}$

```
tensor([[[[-4.7496,      -inf,      -inf],  
         [-0.7920,      -4.0838,      -inf],  
         [-3.4033,      -1.5727,      -0.5862]],  
  
        [[ 2.6880,      -inf,      -inf],  
         [-0.7603,      0.8060,      -inf],  
         [ 1.3757,      1.5708,      0.4877]]]])
```

↓ Apply softmax

```
tensor([[[[1.0000, 0.0000, 0.0000],  
         [0.9641, 0.0359, 0.0000],  
         [0.0417, 0.2603, 0.6980]],  
  
        [[1.0000, 0.0000, 0.0000],  
         [0.1727, 0.8273, 0.0000],  
         [0.3807, 0.4627, 0.1566]]]])
```

Attention weights
(1, 2, 3, 3)
(b, num-heads, num-tokens)

Attention weights
(1, 2, 3, 3)
(b, num-heads, num-tokens, num-tokens)

→ We can also implement Dropout
after this.

Content Vector = $\text{Attention Weights} \times \text{Values}$

Attention Weights
(b, num-heads, num-tokens, num-tokens)

Values
(b, num-heads, num-tokens, head-dim)

```
# Final attention scores
attn_scores = torch.tensor([[[[1.0000, 0.0000, 0.0000],
                            [0.9641, 0.0359, 0.0000],
                            [0.0417, 0.2603, 0.6980]],

                           [[1.0000, 0.0000, 0.0000],
                            [0.1727, 0.8273, 0.0000],
                            [0.3807, 0.4627, 0.1566]]]]])
```

(1, 2, 3, 3)

```
values = torch.tensor([[[[-3.6194, 2.0935, 1.3879],
                        [1.1106, -0.4063, -0.5588],
                        [-1.6594, 0.1052, -0.0468]],

                       [[2.1231, -1.2262, -0.2556],
                        [1.8222, 1.8721, 0.4929],
                        [0.8916, -1.4384, -0.5651]]]])
```

(1, 2, 3, 3)

```
Context vectors:
tensor([[[[-3.6194, 2.0935, 1.3879],
          [-1.8728, 1.6494, 0.7163],
          [-0.3553, -0.5607, -0.1181]],

         [[1.5961, -0.9367, -0.3400],
          [1.3712, 0.7805, 0.2233],
          [0.6997, -0.4487, -0.0530]]]])
```

(1, 2, 3, 3)

(b , num-heads, num-tokens, head-dim)

Step ⑩ Reformatted context vectors:

(b , num-heads, num-tokens, head-dim) \rightarrow (b , num-tokens, num-heads, head-dim)

```
tensor([[[[-3.6194, 2.0935, 1.3879],
          [1.5961, -0.9367, -0.3400]],

         [[-1.8728, 1.6494, 0.7163],
          [1.3712, 0.7805, 0.2233]],

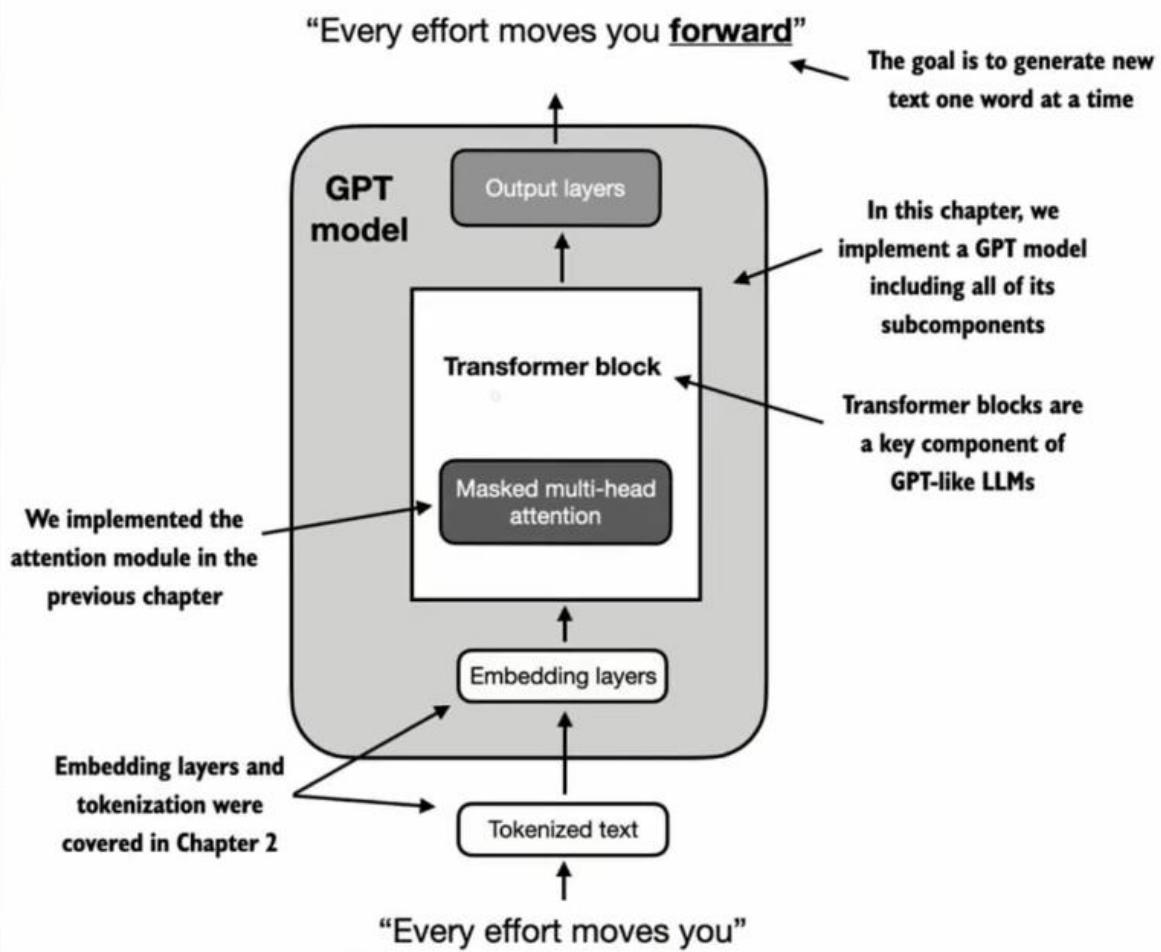
         [[-0.3553, -0.5607, -0.1181],
          [0.6997, -0.4487, -0.0530]]]])
```

(1, 3, 2, 3)

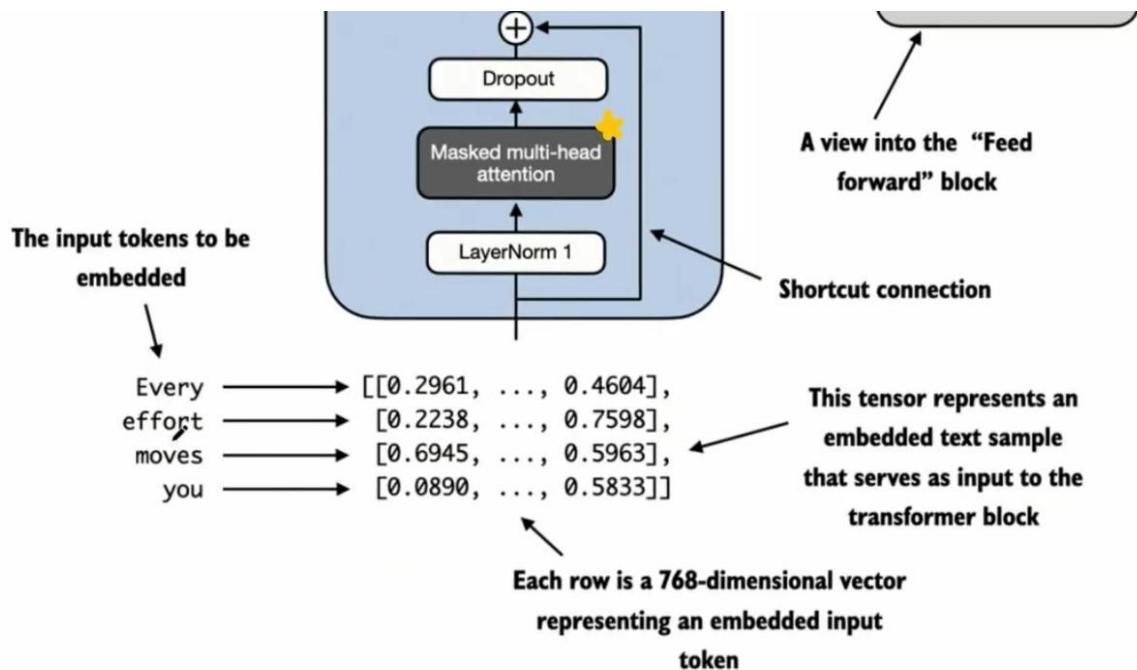
(b , num-tokens, num-heads, head-dim)

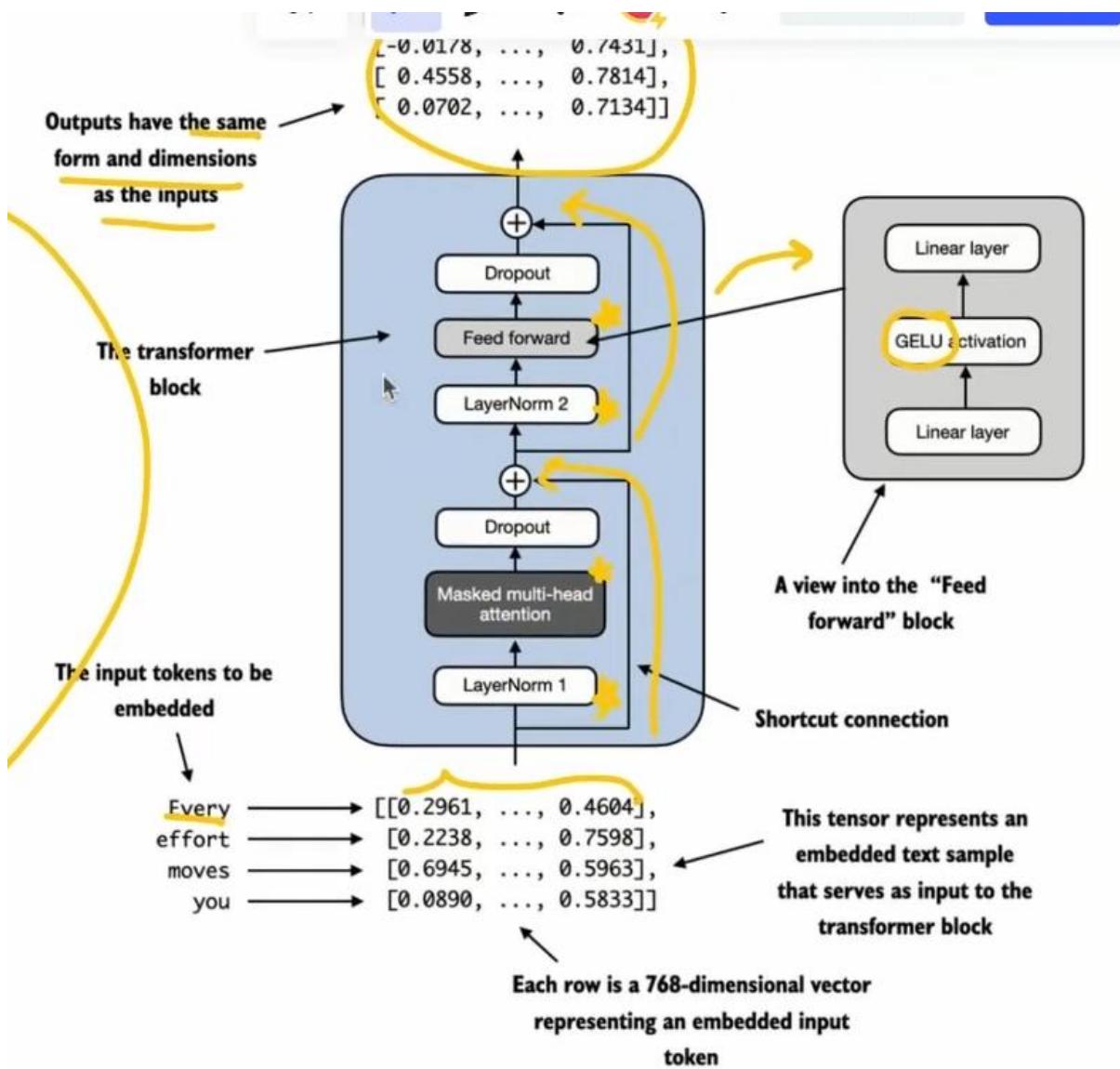
After learning about the attention mechanism, let us learn about the LLM architecture now.

) Coding an LLM architecture:



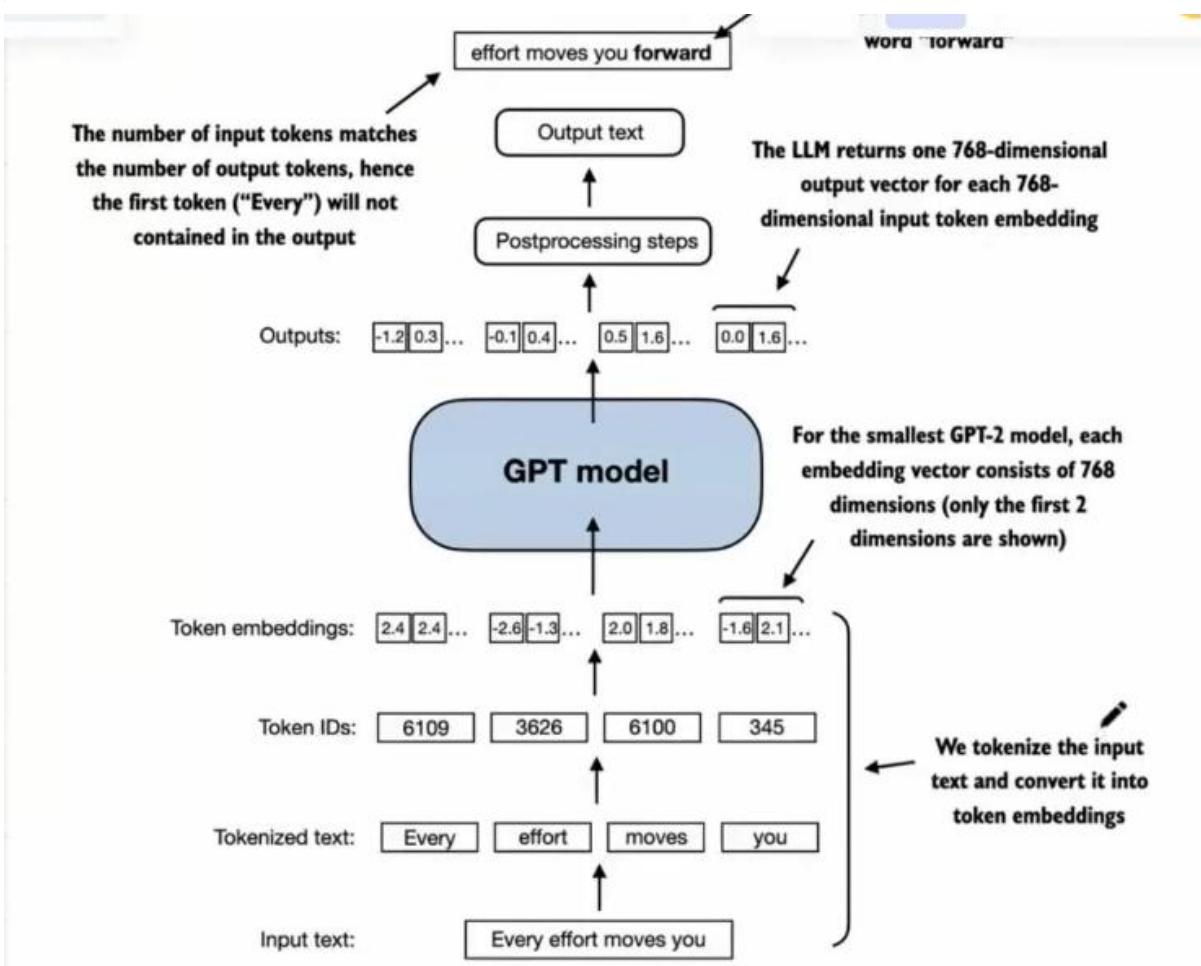
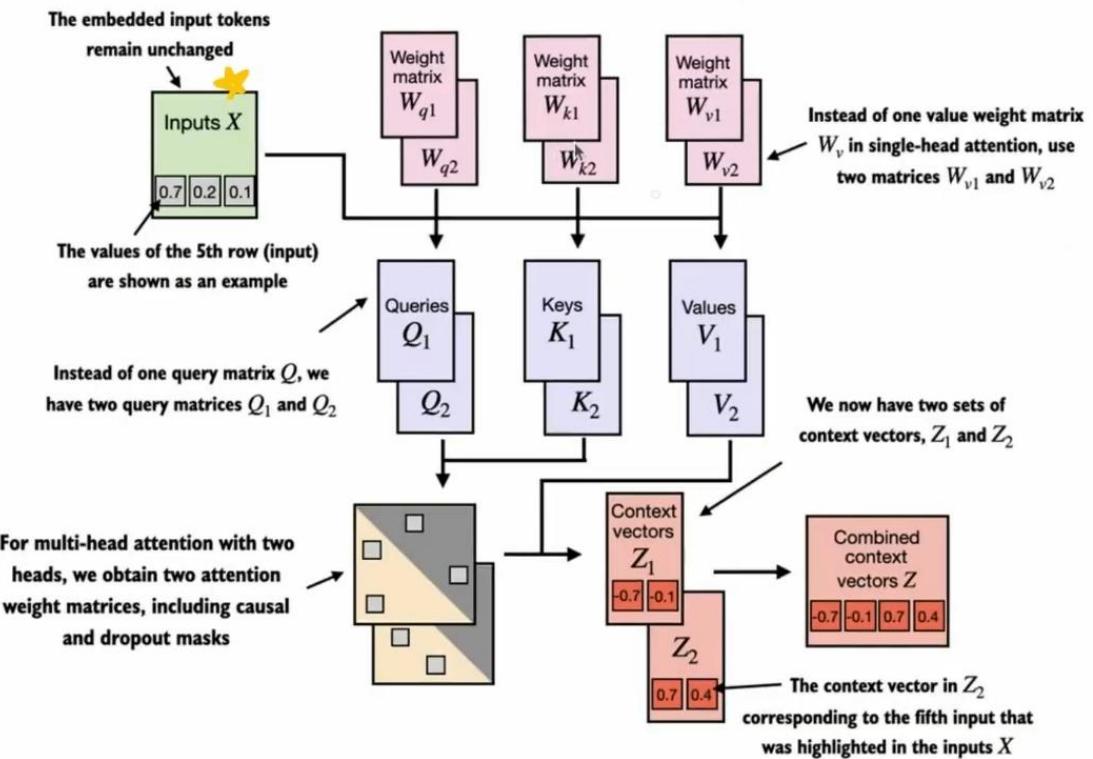
Zoom In Transformer Block





② What we have learnt so far:

- Ⓐ Input tokenization
- Ⓑ Embedding (Token + Positional)
- Ⓒ Masked multi-head attention



③ What we are yet to learn:

ⓐ Transformer blocks

④ We will scale upto the size of a small GPT-2 model → 124 million parameters

Language Models are Unsupervised Multitask Learners

Parameters	Layers	d_{model}
117M	12	768
345M	24	1024
762M	36	1280
1542M	48	1600

→ This was corrected to 124 later

⑤ OpenAI has made GPT-2 weights public.

GPT-3,4 weights have not yet been made public.

We will use these paramet