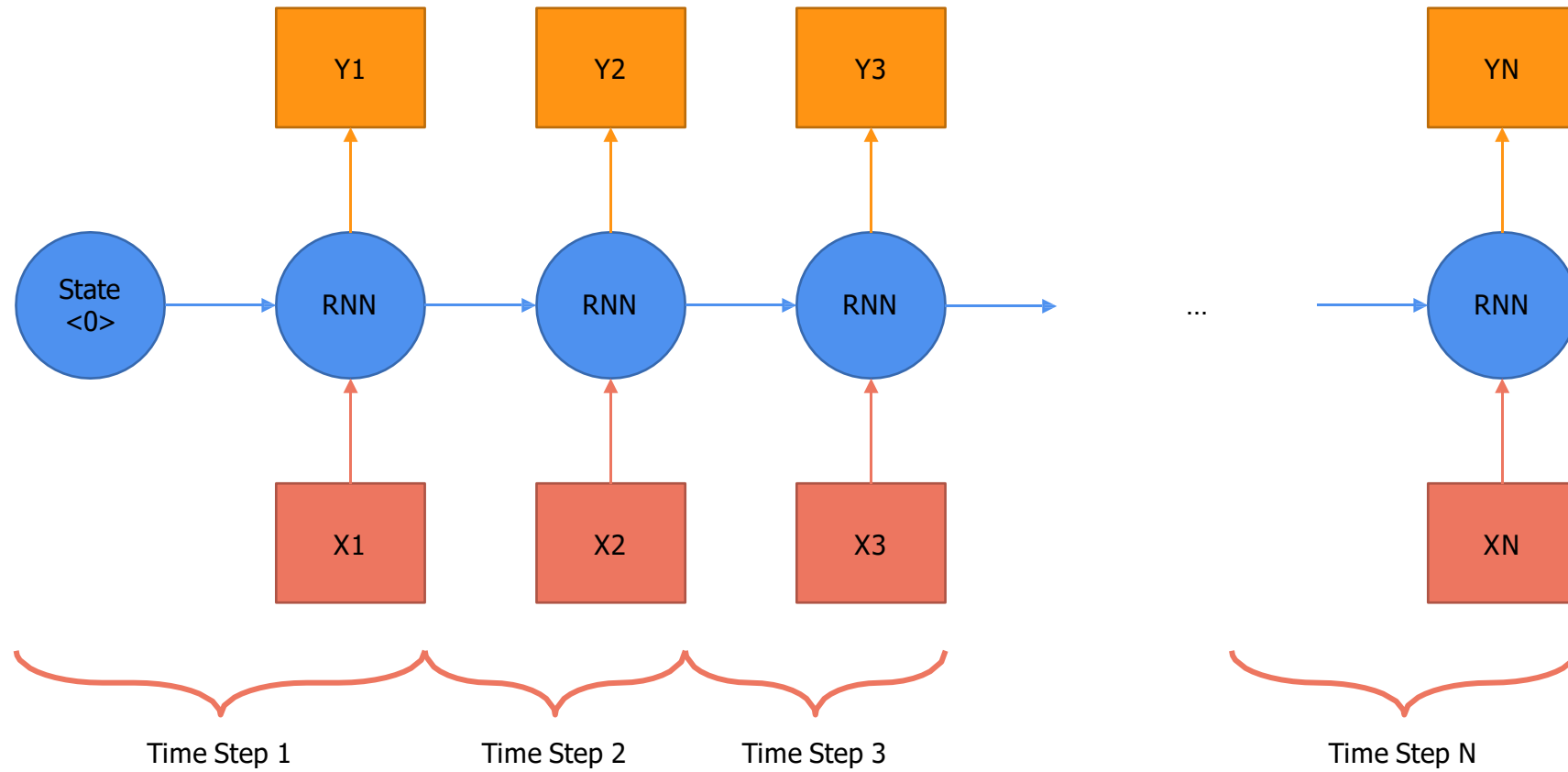


# Transformers

Paper: Attention is all you need

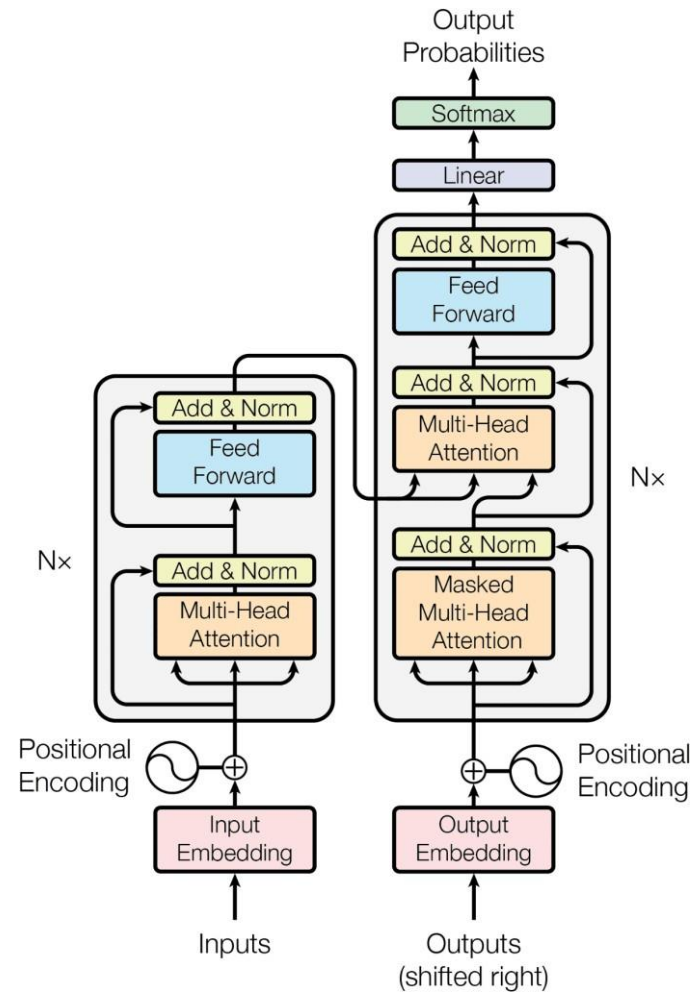
# Recurrent Neural Networks (RNN)



# Problems with RNN (among others)

1. Slow computation for long sequences
2. Vanishing or exploding gradients
3. Difficulty in accessing information from long time ago

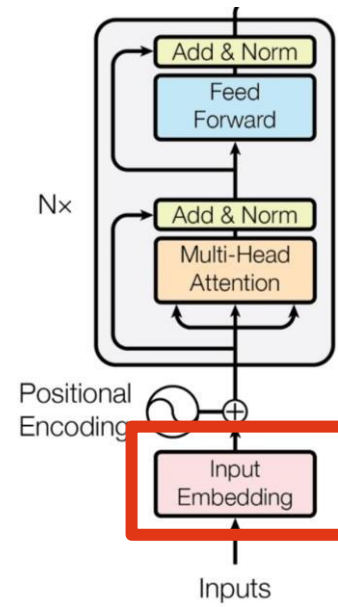
# Transformers



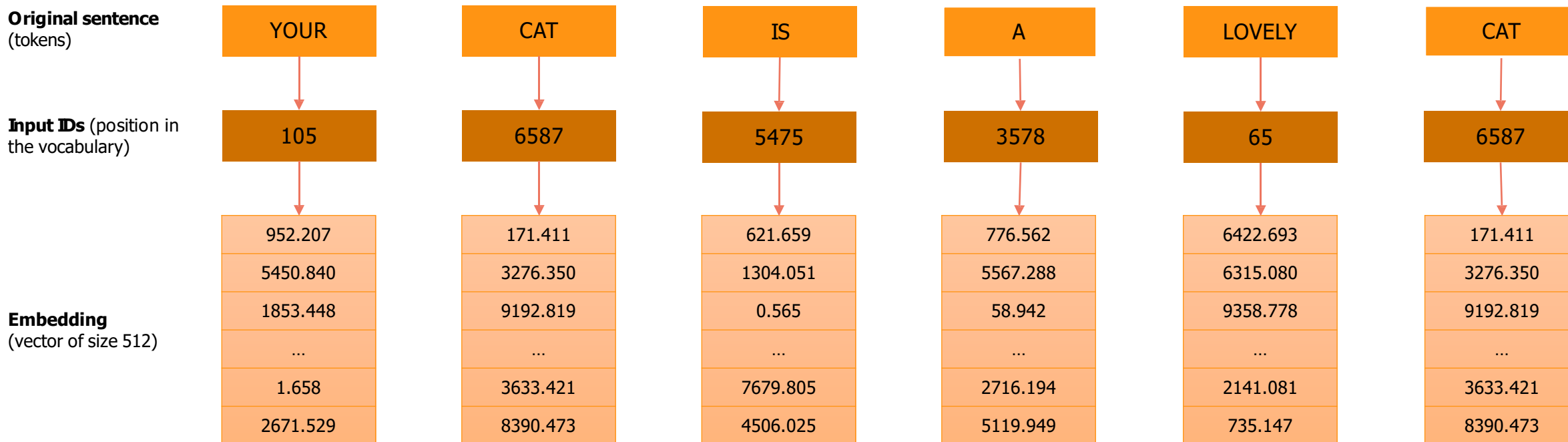
# Notations

Input matrix (sequence,  $d_{\text{model}}$ )

# Encoder

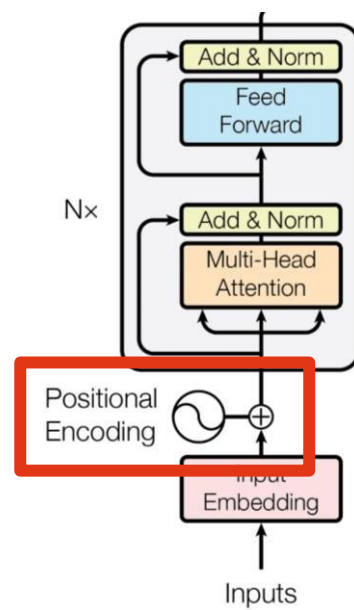


# What is an input embedding?



We define  $d_{\text{model}} = 512$ , which represents the size of the embedding vector of each word

# Encoder





# What is positional encoding?

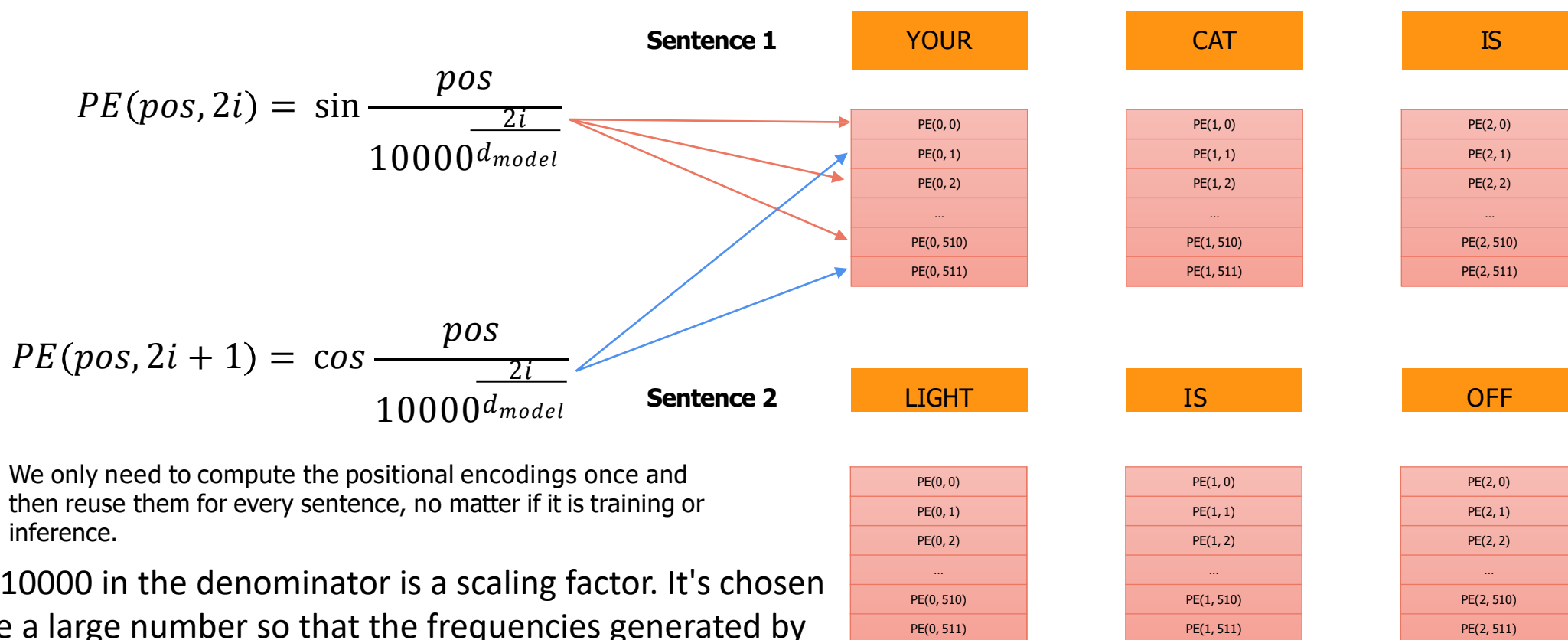
- We want each word to carry some information about its position in the sentence.
  - To inject information about the position or order of elements in a sequence
  - We want the model to treat words that appear close to each other as “close” and words that are distant as “distant”.
  - We want the positional encoding to represent a pattern that can be learned by the model.
- 
- Why sine cosine?
    1. Sine and cosine are periodic functions, so they’ll repeat after some time
    2. Constrained values: lets the value b/w a certain range
    3. Easier to extrapolate for longer sequences, even if we haven't seen certain lengths, we can interpret them during inference

# What is positional encoding?

Original sentence	YOUR	CAT	IS	A	LOVELY	CAT
Embedding (vector of size 512)	952.207	171.411	621.659	776.562	6422.693	171.411
	5450.840	3276.350	1304.051	5567.288	6315.080	3276.350
	1853.448	9192.819	0.565	58.942	9358.778	9192.819
	...	...	...	...	...	...
	1.658	3633.421	7679.805	2716.194	2141.081	3633.421
	2671.529	8390.473	4506.025	5119.949	735.147	8390.473
Position Embedding (vector of size 512). Only computed once and reused for every sentence during training and inference.	+	+	+	+	+	+
	...	1664.068	...	...	...	1281.458
	...	8080.133	...	...	...	7902.890
	...	2620.399	...	...	...	912.970
	...	...	...	...	...	3821.102
	...	9386.405	...	...	...	1659.217
Encoder Input (vector of size 512)	...	3120.159	...	...	...	7018.620
	=	=	=	=	=	=
	...	1835.479	...	...	...	1452.869
	...	11356.483	...	...	...	11179.24
	...	11813.218	...	...	...	10105.789
	...	...	...	...	...	...
	...	13019.826	...	...	...	5292.638
	...	11510.632	...	...	...	15409.093

# What is positional encoding?

These sinusoidal functions are designed to create unique and distinct positional embeddings for different positions in a sequence, allowing the model to distinguish between tokens based on their position. The formula is fixed and does not adapt to the specific content of the tokens.



We only need to compute the positional encodings once and then reuse them for every sentence, no matter if it is training or inference.

The 10000 in the denominator is a scaling factor. It's chosen to be a large number so that the frequencies generated by the exponent are relatively small. This ensures that the positional encodings don't dominate the word embeddings

- **i = 0:**

- $PE(pos, 2i) = PE(pos, 0) = \sin(pos / 10000^{(2*0/d\_model)}) = \sin(pos / 10000^0)$
- $PE(pos, 2i+1) = PE(pos, 1) = \cos(pos / 10000^{(2*0/d\_model)}) = \cos(pos / 10000^0)$

- **i = 1:**

- $PE(pos, 2i) = PE(pos, 2) = \sin(pos / 10000^{(2*1/d\_model)}) = \sin(pos / 10000^{(2/d\_model)})$
- $PE(pos, 2i+1) = PE(pos, 3) = \cos(pos / 10000^{(2*1/d\_model)}) = \cos(pos / 10000^{(2/d\_model)})$

- **i = 2:**

- $PE(pos, 2i) = PE(pos, 4) = \sin(pos / 10000^{(2*2/d\_model)}) = \sin(pos / 10000^{(4/d\_model)})$
- $PE(pos, 2i+1) = PE(pos, 5) = \cos(pos / 10000^{(2*2/d\_model)}) = \cos(pos / 10000^{(4/d\_model)})$

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

$$\begin{aligned} PE(0) &= \left( \sin\left(\frac{0}{10000^{\frac{1}{d_{model}}}}\right), \cos\left(\frac{0}{10000^{\frac{1}{d_{model}}}}\right), \sin\left(\frac{0}{10000^{\frac{1}{d_{model}}}}\right), \cos\left(\frac{0}{10000^{\frac{1}{d_{model}}}}\right), \dots, \sin\left(\frac{0}{10000^{\frac{1}{d_{model}}}}\right), \cos\left(\frac{0}{10000^{\frac{1}{d_{model}}}}\right) \right) \\ &= (\sin(0), \cos(0), \sin(0), \cos(0), \dots, \sin(0), \cos(0)) \\ &= (0, 1, 0, 1, \dots, 0, 1) \end{aligned}$$

Therefore, position 0 is a vector with alternatively repeating 0s and 1s.

Position 1 is as follows:

$$\begin{aligned} PE(1) &= \left( \sin\left(\frac{1}{10000^{\frac{1}{d_{model}}}}\right), \cos\left(\frac{1}{10000^{\frac{1}{d_{model}}}}\right), \sin\left(\frac{1}{10000^{\frac{1}{d_{model}}}}\right), \cos\left(\frac{1}{10000^{\frac{1}{d_{model}}}}\right), \dots, \sin\left(\frac{1}{10000^{\frac{1}{d_{model}}}}\right), \cos\left(\frac{1}{10000^{\frac{1}{d_{model}}}}\right) \right) \\ &= (0.8414, 0.5403, 0.8218, 0.5696, \dots, 0.0001, 0.9999) \end{aligned}$$

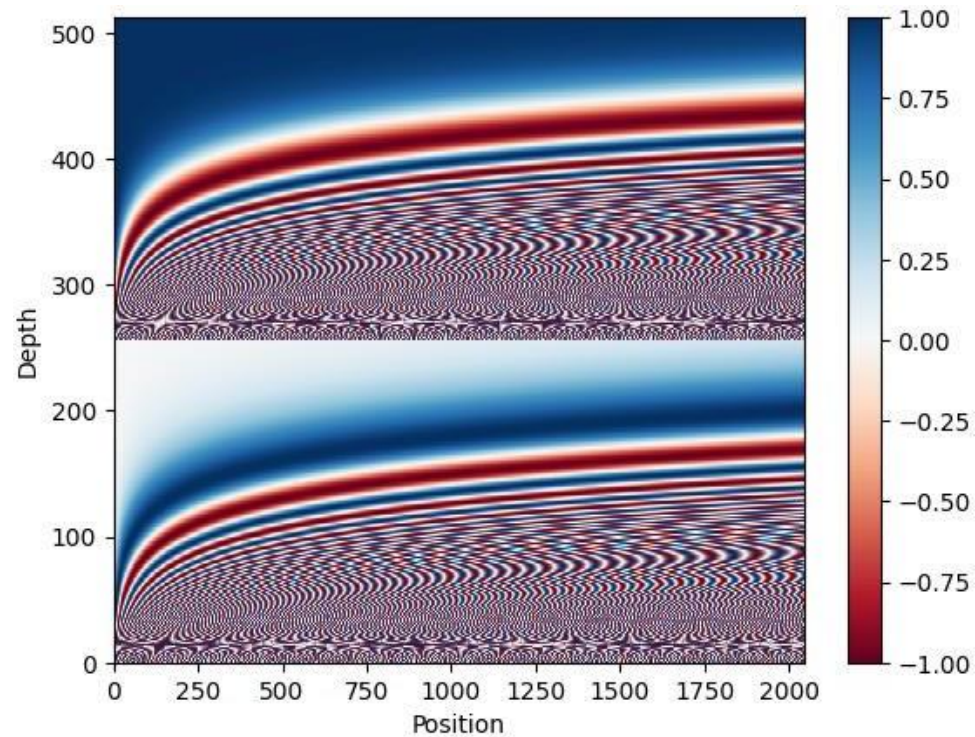
# Pre Computation of PE

positional encodings are typically pre-computed. The sine and cosine functions are deterministic, meaning that for a given position and dimension, they will always produce the same value. Therefore, you only need to calculate them once for each possible position up to a maximum sequence length. These pre-computed values are then stored and added to the word embeddings as needed. You don't need to recalculate them for every new word, unless the word is at a position beyond your pre-computed range.

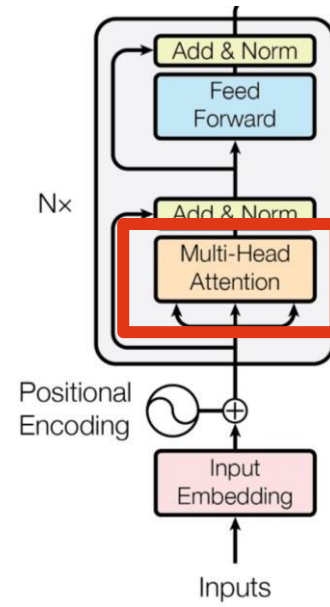
This pre-computation makes positional encoding very efficient. It's a one-time cost that significantly improves the self-attention mechanism's ability to handle word order without adding significant computational overhead during training or inference.

# Why trigonometric functions?

Trigonometric functions like **cos** and **sin** naturally represent a pattern that the model can recognize as continuous, so relative positions are easier to see for the model. By watching the plot of these functions, we can also see a regular pattern, so we can hypothesize that the model will see it too.



# Encoder



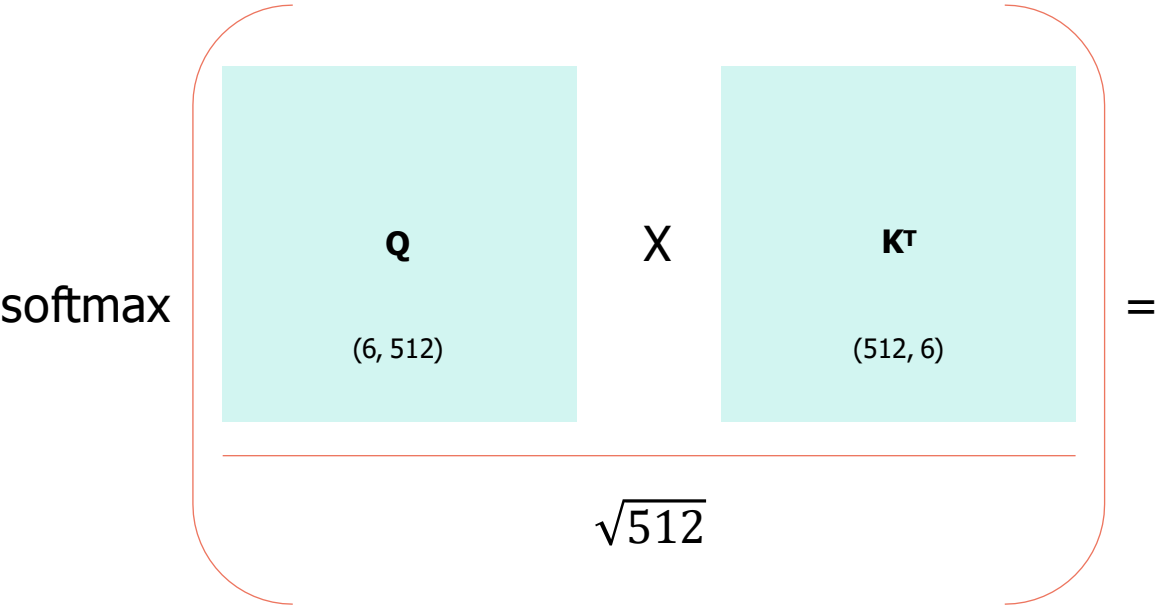
# What is Self-Attention?

Self-Attention allows the model to relate words to each other.

In this simple case we consider the sequence length **seq** = 6 and **d<sub>model</sub>** = **d<sub>k</sub>** = 512.

The matrices **Q**, **K** and **V** are just the input sentence.

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



	YOUR	CAT	IS	A	LOVELY	CAT	$\Sigma$
YOUR	0.268	0.119	0.134	0.148	0.179	0.152	1
CAT	0.124	0.278	0.201	0.128	0.154	0.115	1
IS	0.147	0.132	0.262	0.097	0.218	0.145	1
A	0.210	0.128	0.206	0.212	0.119	0.125	1
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174	1
CAT	0.195	0.114	0.203	0.103	0.157	0.229	1

\* all values are random.

\* for simplicity I considered only one head, which makes  $d_{\text{model}} = d_k$ .

(6, 6)



# How to compute Self-Attention?

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

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

(6, 6)

X

V

(6, 512)

=

Attention

(6, 512)

Each row in this matrix captures not only the meaning (given by the embedding) or the position in the sentence (represented by the positional encodings) but also each word's interaction with other words.

# Self-Attention in detail

- Self-Attention is permutation invariant since it does not cater the order of words, and treats each word independently  
(if you do not consider the contribution of the positional encoding)

Order does not affect the output of the self attention

If u rearrange the inputs, attention values will remain the same

- Self-Attention requires no parameters. Up to now the interaction between words has been driven by their embedding and the positional encodings. This will change later.
- We expect values along the diagonal to be the highest.
- If we don't want some positions to interact, we can always set their values to  $-\infty$  before applying the *softmax* in this matrix and the model will not learn those interactions. We will use this in the decoder.

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

# Self-Attention in detail

Permutation invariance is a problem because word order is crucial for understanding language. The sentence "The cat sat" is very different from "Sat cat the." Standard self-attention doesn't inherently understand this.

## **How is this solved?**

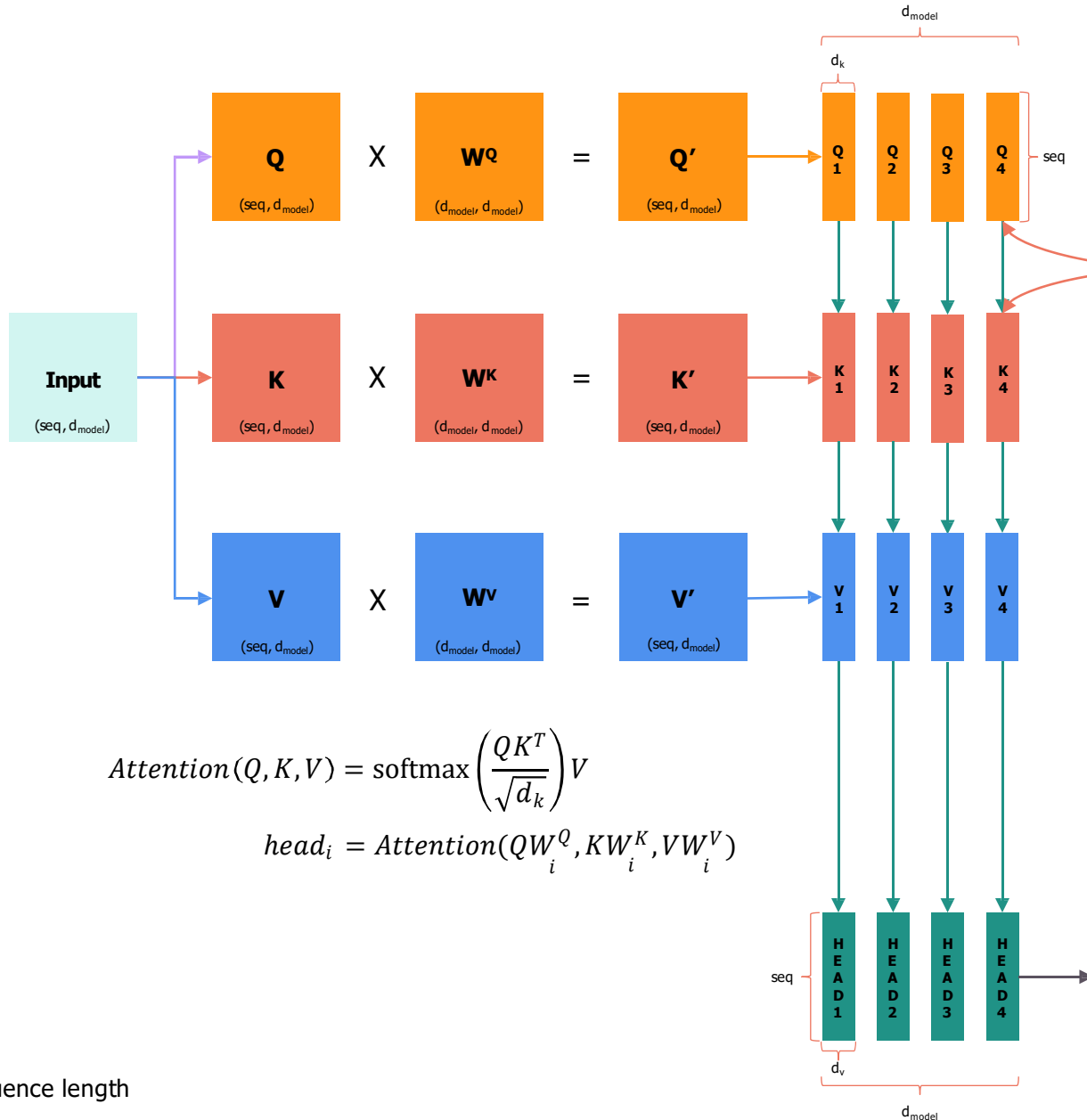
Positional encodings are added to the input embeddings to address this limitation. These encodings provide information about the absolute position of each word, allowing the self-attention mechanism to consider word order when calculating attention weights. With positional encodings, the attention weights for the permuted sentence would be *different* from the original sentence, reflecting the change in meaning due to the change in word order.

# Multi-head Attention

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

$$MultiHead(Q, K, V) = \text{Concat}(head_1 \dots head_h)W^O$$

$$head_i = Attention(QW^Q, KW_i^K, VW^V)$$



**query:** The "query" is like a word looking for other words to pay attention to.

**key:** The "key" is like a word being looked at by other words.

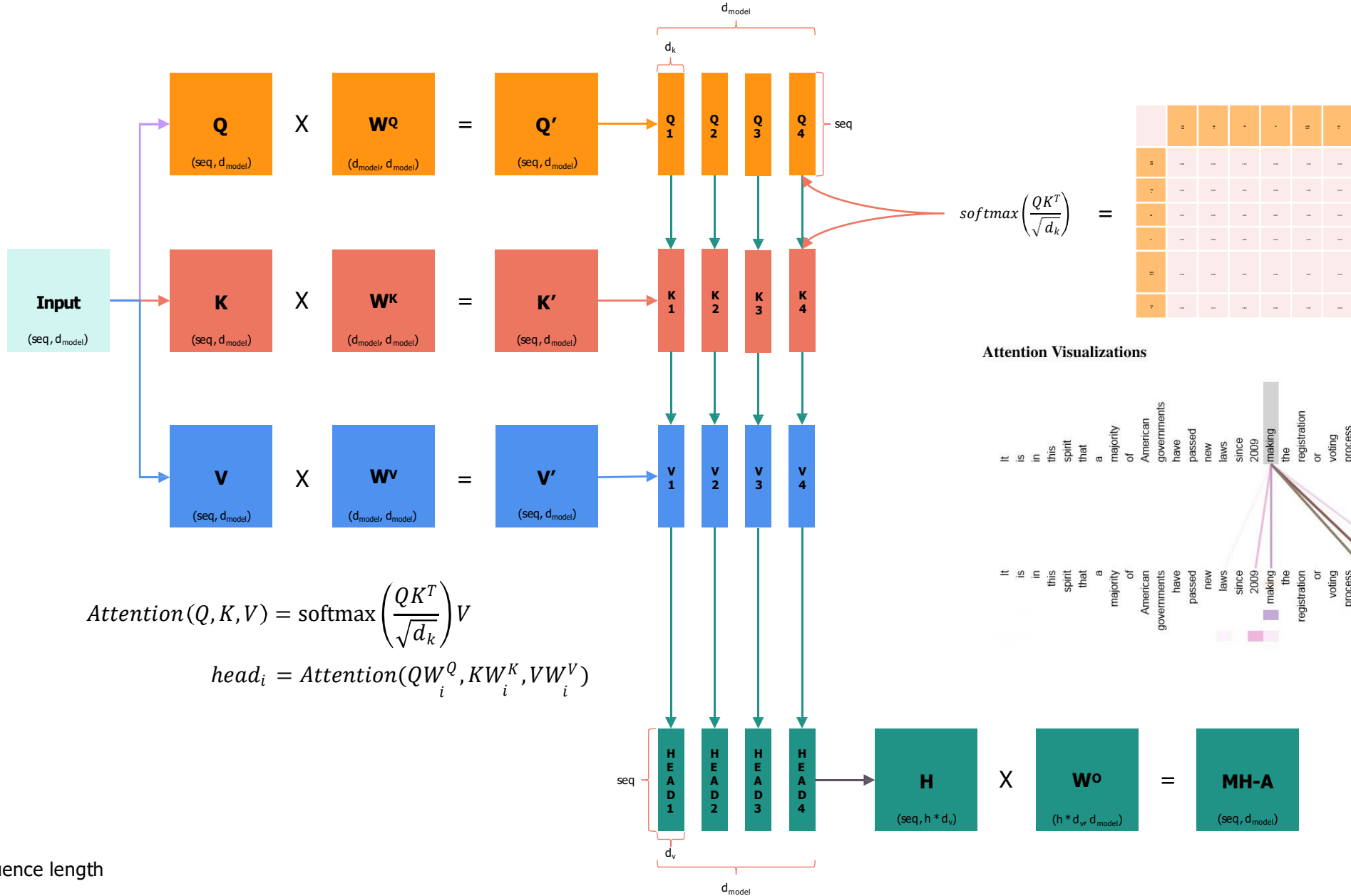
**value:** the "value" is like the information or meaning of a word

$seq$  = sequence length

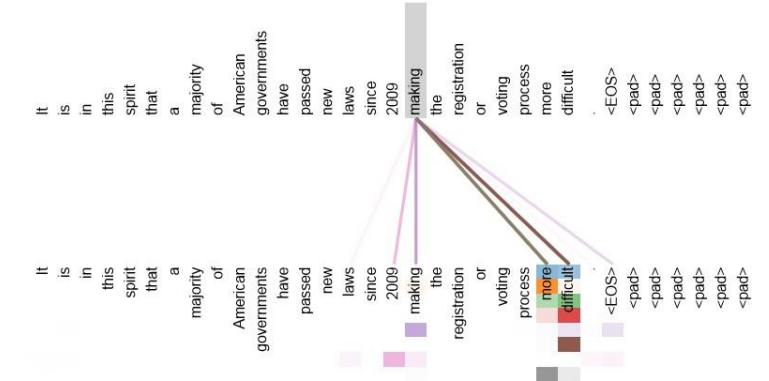
$d_{model}$  = size of the embedding vector

$h$  = number of heads

$d_k = d_v$  =  $d_{model} / h$



Attention Visualizations



- $seq$  = sequence length
- $d_{model}$  = size of the embedding vector
- $h$  = number of heads
- $d_k = d_v$  =  $d_{model} / h$

The purpose of using multiple heads is to allow the model to focus on different aspects or patterns in the input sequence in parallel.

$$MultiHead(Q, K, V) = Concat(head_1 \dots head_h)W^O$$

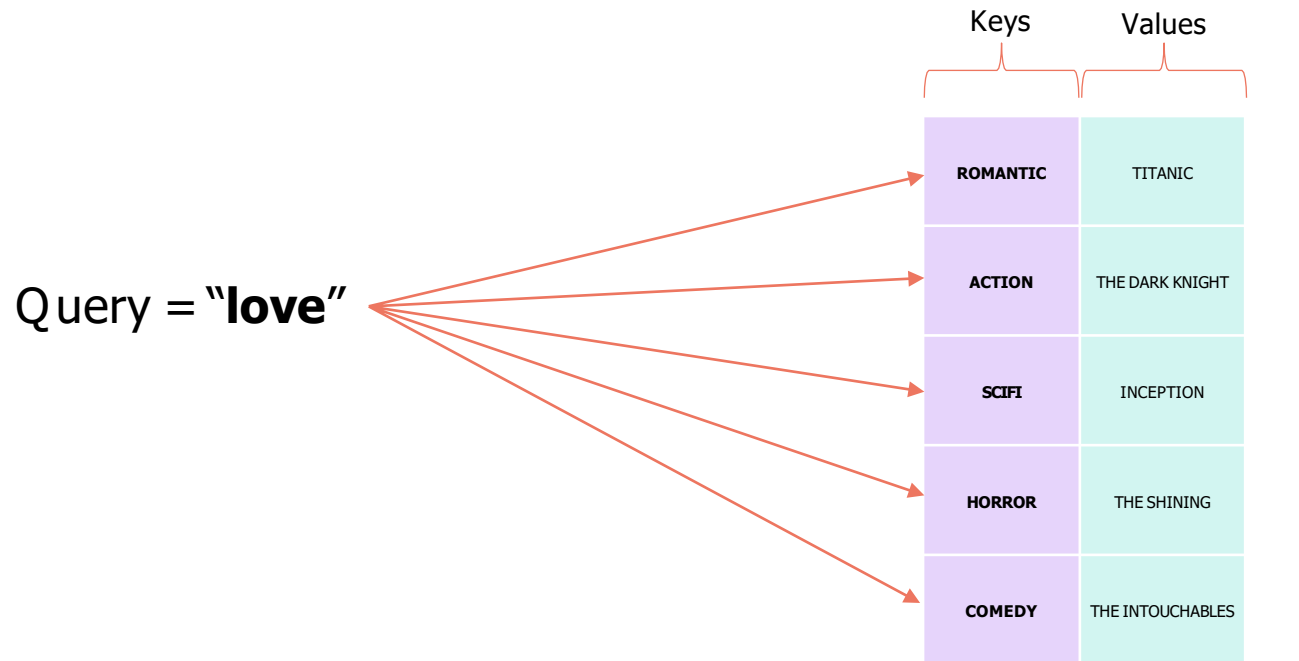
**1.Attention:** Attention mechanism is used in deep learning models to selectively focus on important parts of input data while ignoring irrelevant parts. In the context of natural language processing (NLP), attention mechanism has been used in sequence-to-sequence models to align relevant parts of the input sentence with each word in the output sentence.

**2.Self-attention:** Self-attention is a special case of attention mechanism where the input sequence is the same as the output sequence. In other words, each word in the input sequence attends to all other words in the same sequence to generate a representation of itself. Self-attention is used in transformers to process input sequences by generating a set of weighted vectors for each word in the sequence.

**3.Multi-head attention:** Multi-head attention is an extension of self-attention where the attention mechanism is applied multiple times in parallel with different weight matrices to capture different aspects of the input sequence. This allows the model to attend to multiple parts of the input sequence at the same time and learn more complex representations. Multi-head attention is used in transformers to process input sequences and generate the context vectors that are used to generate the output sequence

# Why query, keys and values?

The Internet says that these terms come from the database terminology or the Python-like dictionaries.



\* this could be a Python dictionary or a database table.

**query:** The "query" is like a word looking for other words to pay attention to.

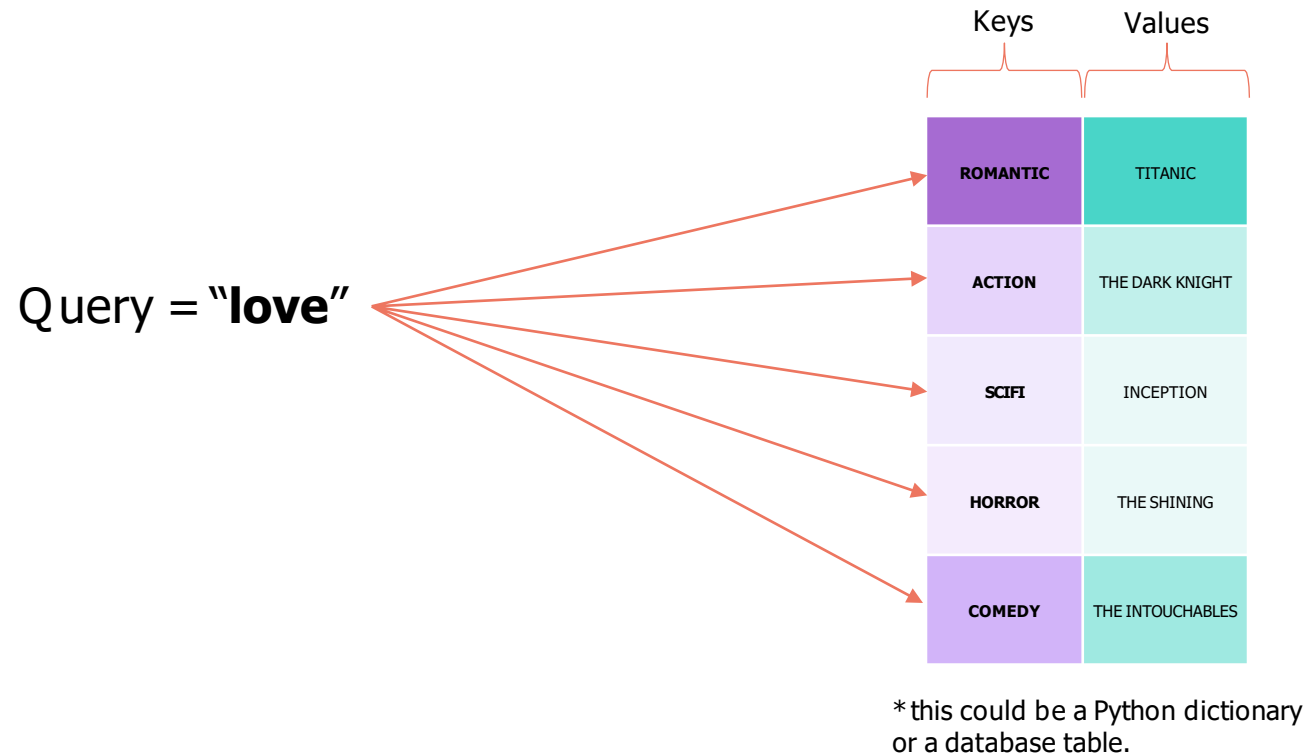
**key:** The "key" is like a word being looked at by other words.

**value:** the "value" is like the information or meaning of a word

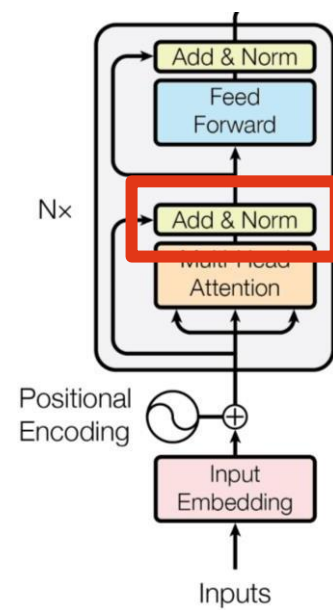


# Why query, keys and values?

The Internet says that these terms come from the database terminology or the Python-like dictionaries.



# Encoder



# What is layer normalization?

Layer normalization ensures that each part of the sequence (each word or token) is normalized independently, providing stability and ensuring the model generalizes well across the entire sequence

Batch of 3 items

ITEM 1

ITEM 2

ITEM 3

50.147
3314.825
...
...
8463.361
8.021

1242.223
688.123
...
...
434.944
149.442

9.370
4606.674
...
...
944.705
21189.444

- **Batch normalization** looks at all the data samples in the batch and normalizes for each feature.

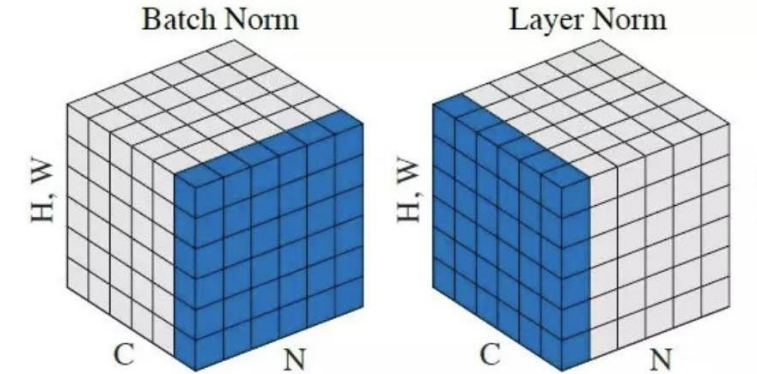
- **Layer normalization** normalizes the values for each instance, taking into account each instance separately

$$\mu_1$$
$$\sigma_1^2$$

$$\mu_2$$
$$\sigma_2^2$$

$$\mu_3$$
$$\sigma_3^2$$

$$\hat{x}_j = \frac{x_j - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$



We also introduce two parameters, usually called **gamma** (multiplicative) and **beta** (additive) that introduce some fluctuations in the data, because maybe having all values between 0 and 1 may be too restrictive for the network. The network will learn to tune these two parameters to introduce fluctuations when necessary.