



Advanced Computer Science Course

Lecture 3

Advanced Language Models

Tishreen University
Computer and automatic control
engineering dept.

Master Program- 2024

1st year

Dr. Ali Mahmoud Mayya

Language Models

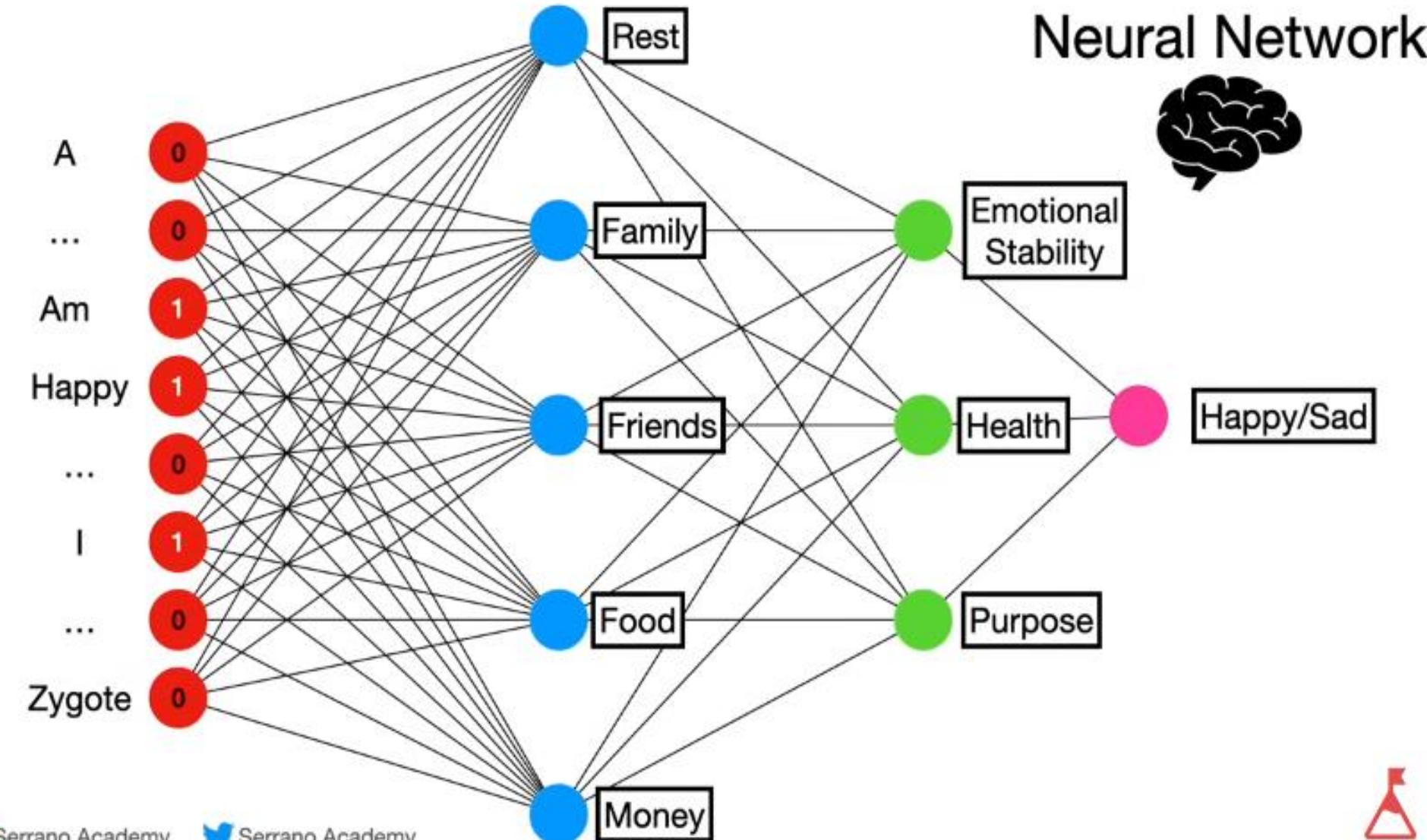
Sentiment analysis

Feed Forward Fully-connected Neural Networks

Model	
Happy	3 pts.
Sad	-4 pts.
Magnificent	4 pts.
Hello	0 pts.
Very	0 pts.
So	0 pts.
...	...

I am very happy!! 3 😊
0 0 0 3

I had a very bad day -4 😢
0 0 0 0 -4 0



Language Models

Sentiment analysis
Predicting next word in the sentence (simple start)

Feed Forward Fully-connected Neural Networks

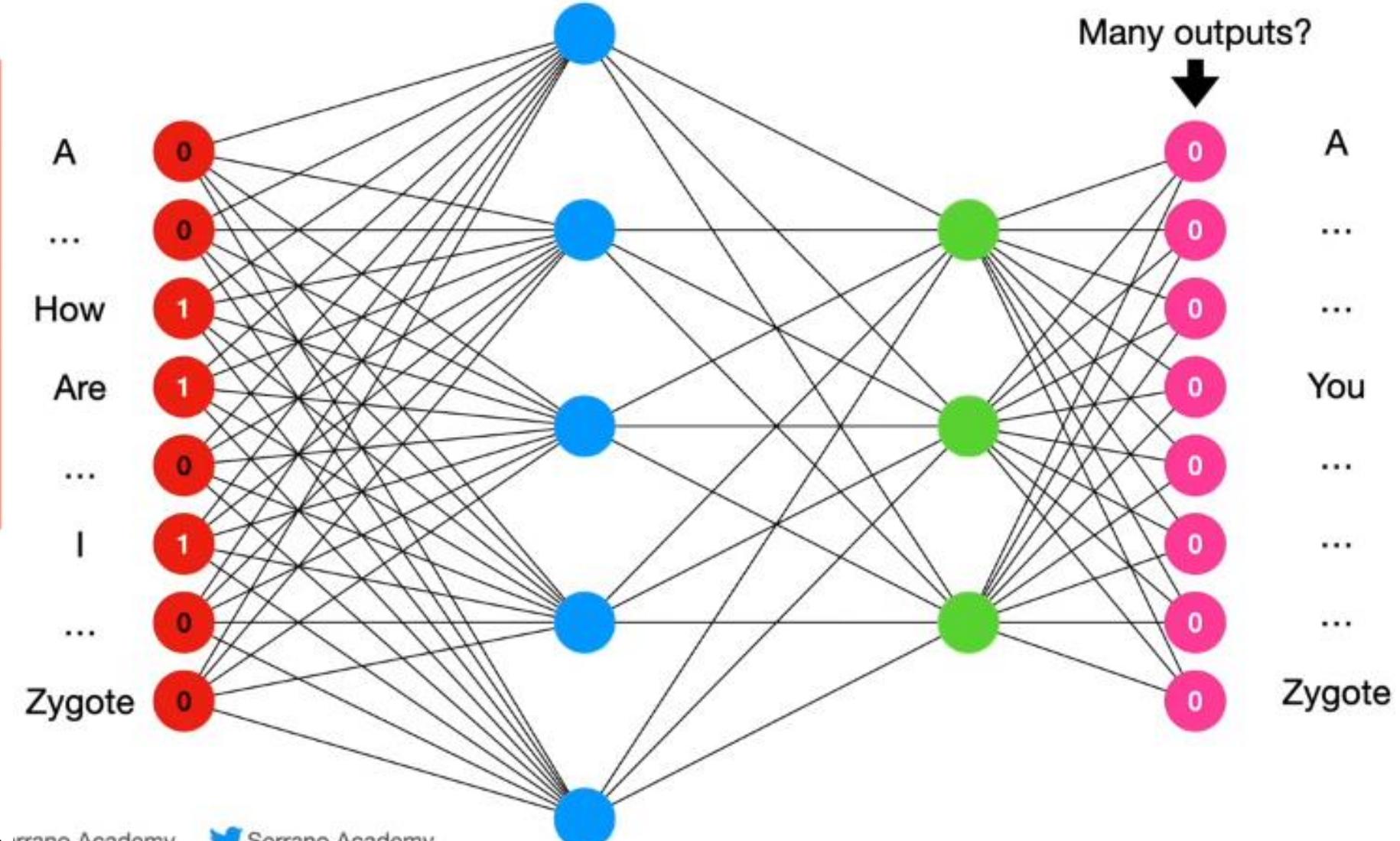
Model	
Happy	3 pts.
Sad	-4 pts.
Magnificent	4 pts.
Hello	0 pts.
Very	0 pts.
So	0 pts.
...	...

I am very happy!! 3 😊

0 0 0 3

I had a very bad day -4 😢

0 0 0 0 -4 0

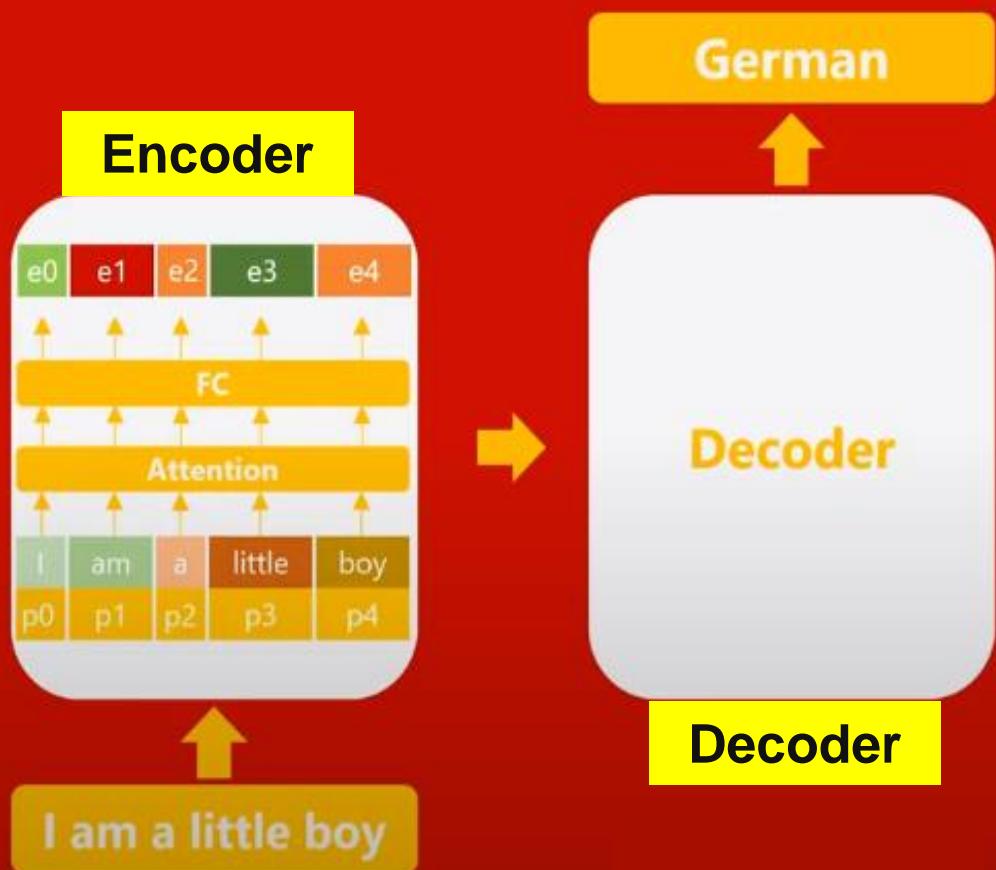


Language Models

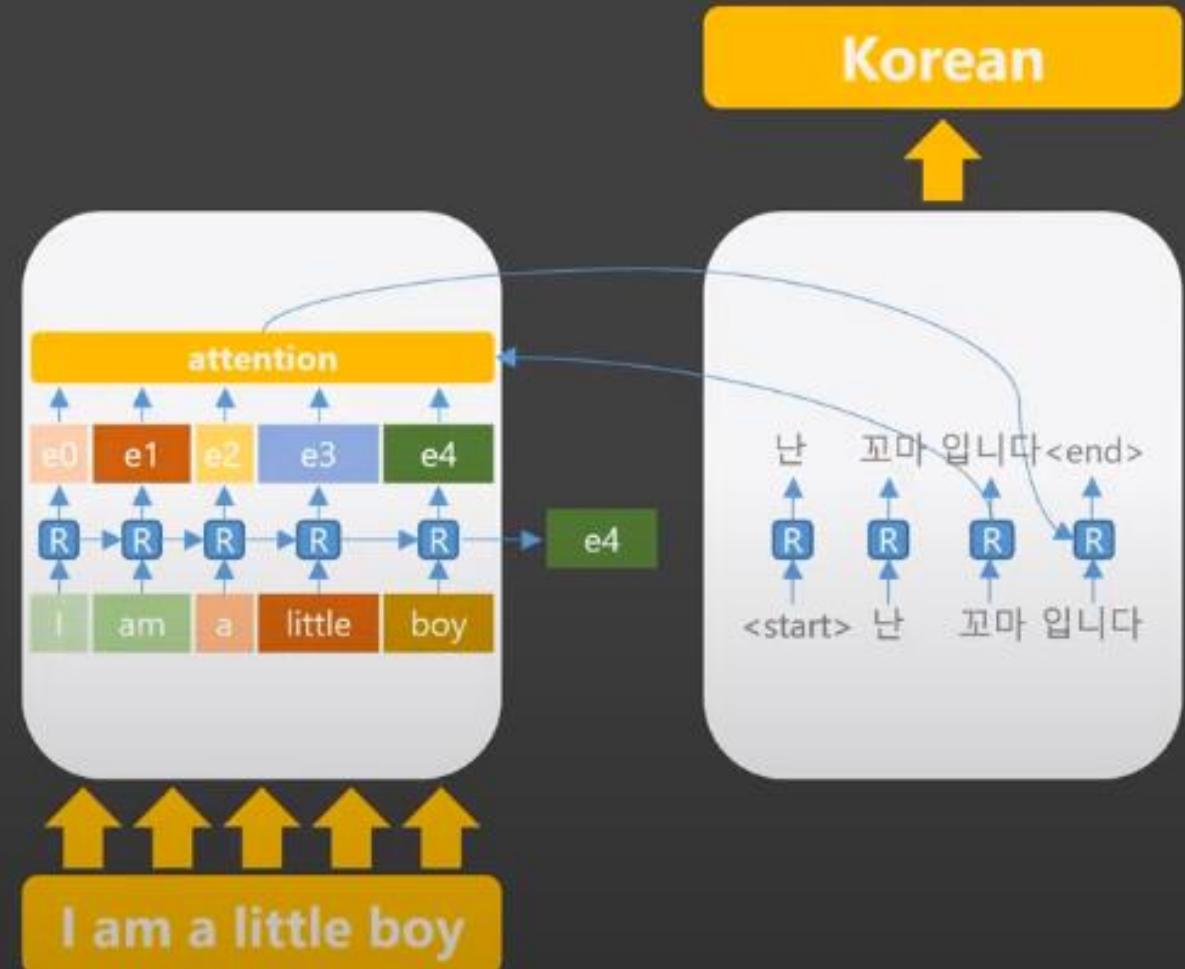
Transformer models

English to Korean Translation example

Transformer



Trend before Transformer



Language Models

English to Korean Translation ex)

the **attention mechanism** enables the model to assign different attention weights/scores to different tokens in a sequence. This way, the model can give more importance to relevant information, ignore irrelevant information, and effectively capture long-range dependencies in the data



Transformer models

Language Models

Transformer models

Encoder part:

Multi-Head attention layers:

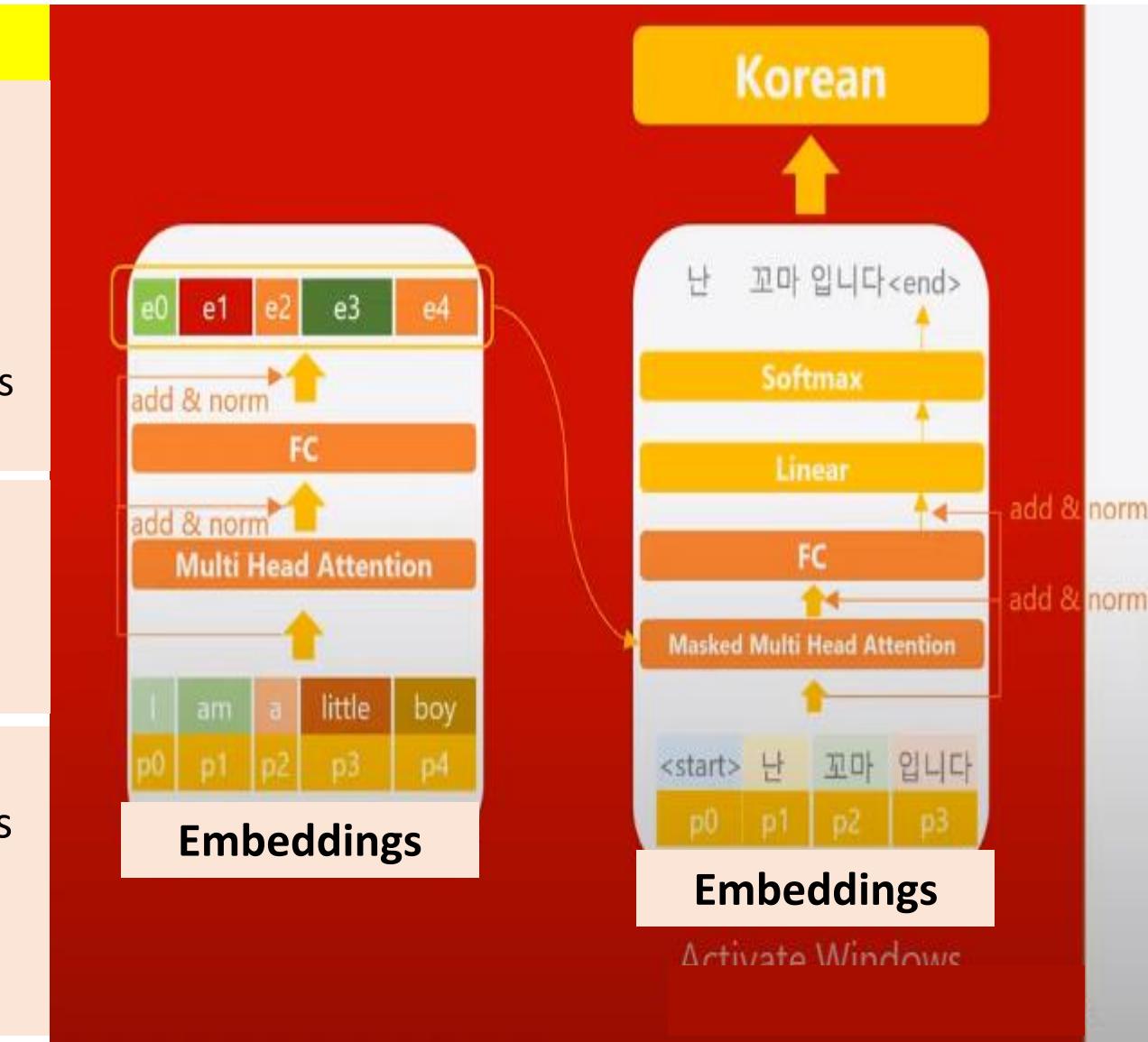
1. Computes self-attention scores for each position in the input sequence.
2. Employs multiple attention heads to capture different aspects of context.
3. The output is a weighted sum of values, where the weights are determined by the attention scores.

Feed-Forward Network:

1. Fully connected feed-forward network.
2. Two linear transformations followed by a Relu function.
3. Non-linearity and helps model complex relationships.

Add&Norm skip connections:

- 1- Deliver lower layer's information to the higher layers very efficiently.
- 2- Enhance training process especially with too deep model (Eliminate vanishing gradient problem)



Language Models

Transformer models

Decoder part:

Masked Multi-Head attention layers:

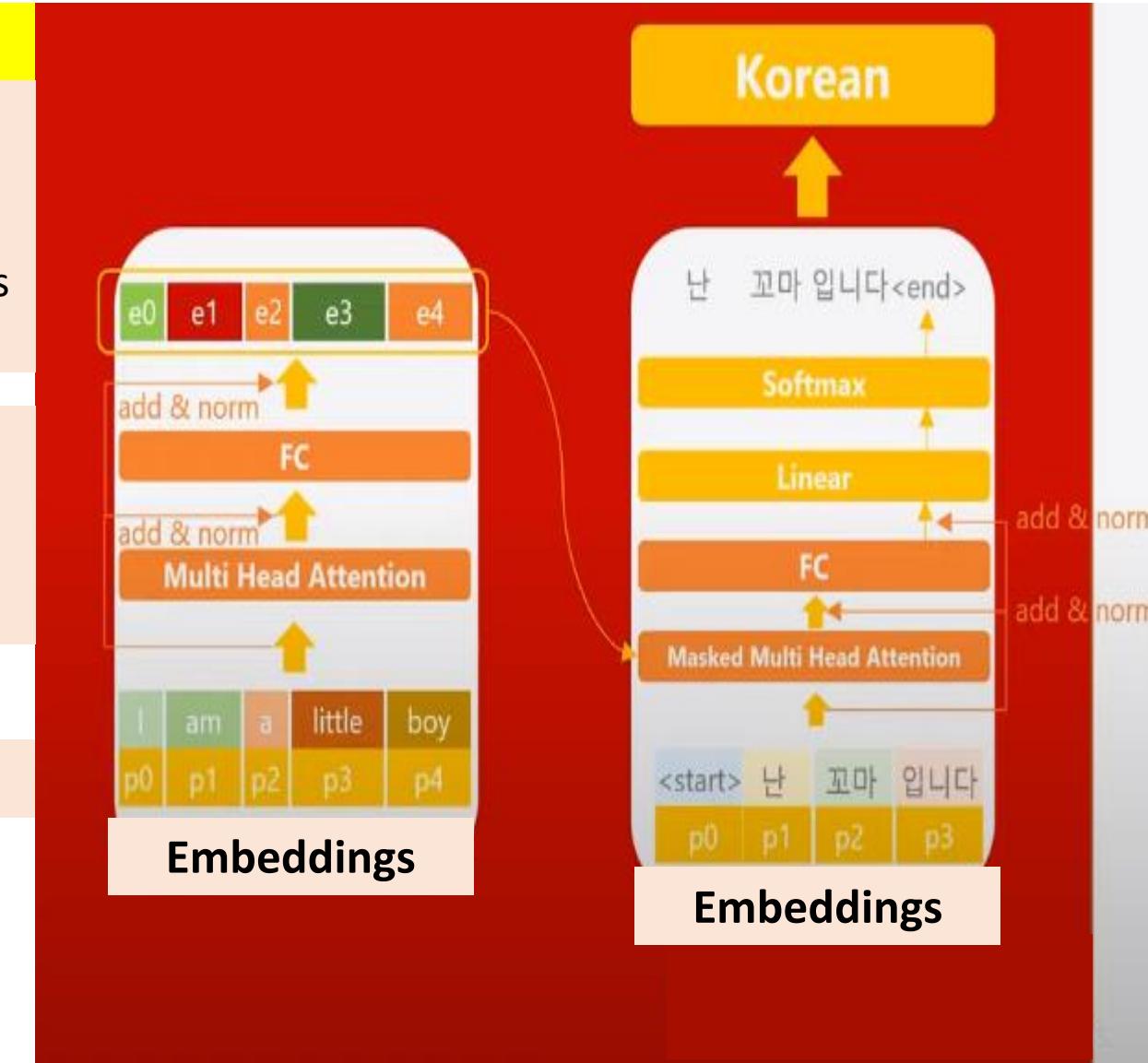
Like the encoder's self-attention, but with a mask to prevent attending to future positions.

- Ensures that the decoder only uses information from previous positions during prediction.

Feed-Forward Network:

1. Fully connected feed-forward network.
2. Two linear transformations followed by a Relu function.
3. Helps model complex relationships.

Linear layer with softmax activation for classification

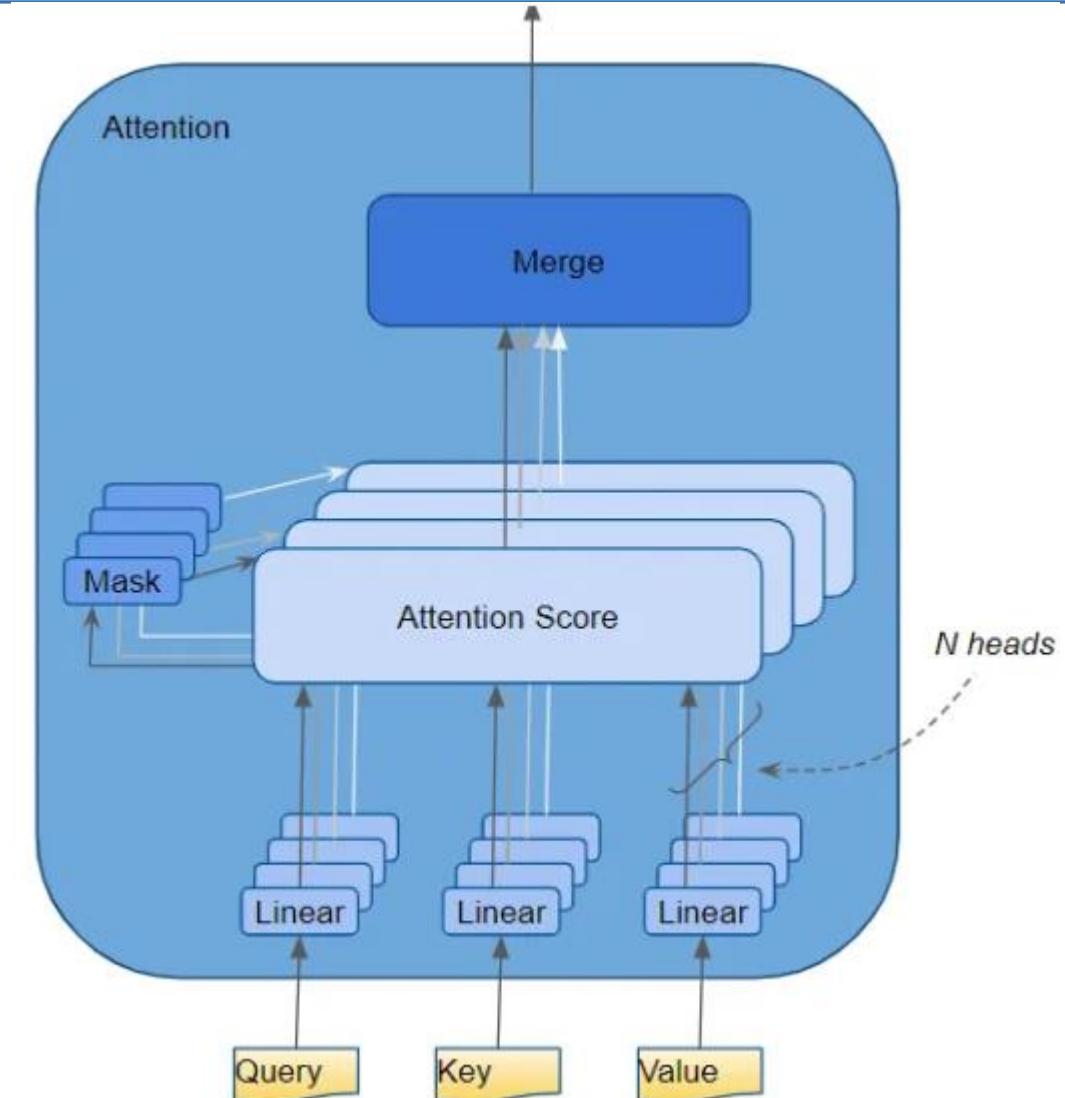


Language Models

Transformer models

Multi-head attention layer architecture

- The Attention module splits its Query, Key, and Value parameters N-ways
- Passes each split independently through a separate Head.
- Attention calculations are then combined together to produce a final Attention score.

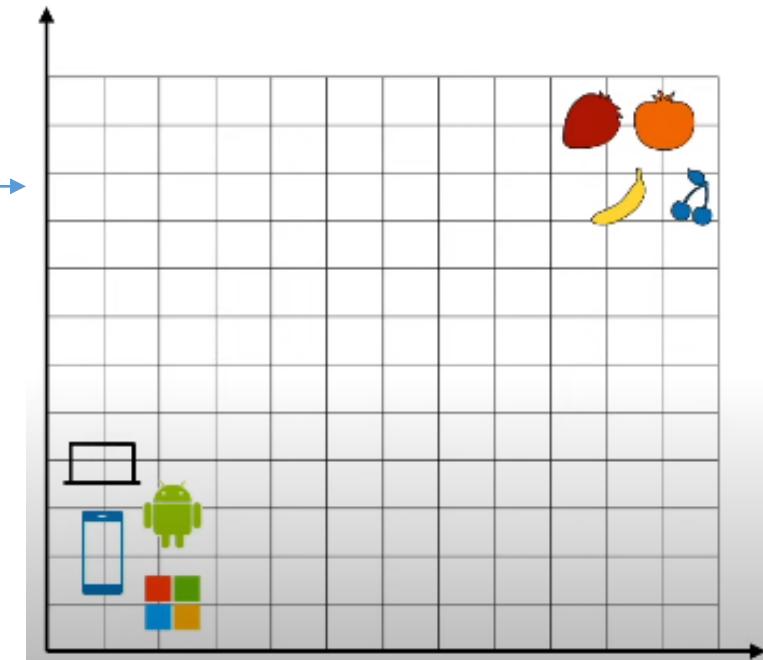
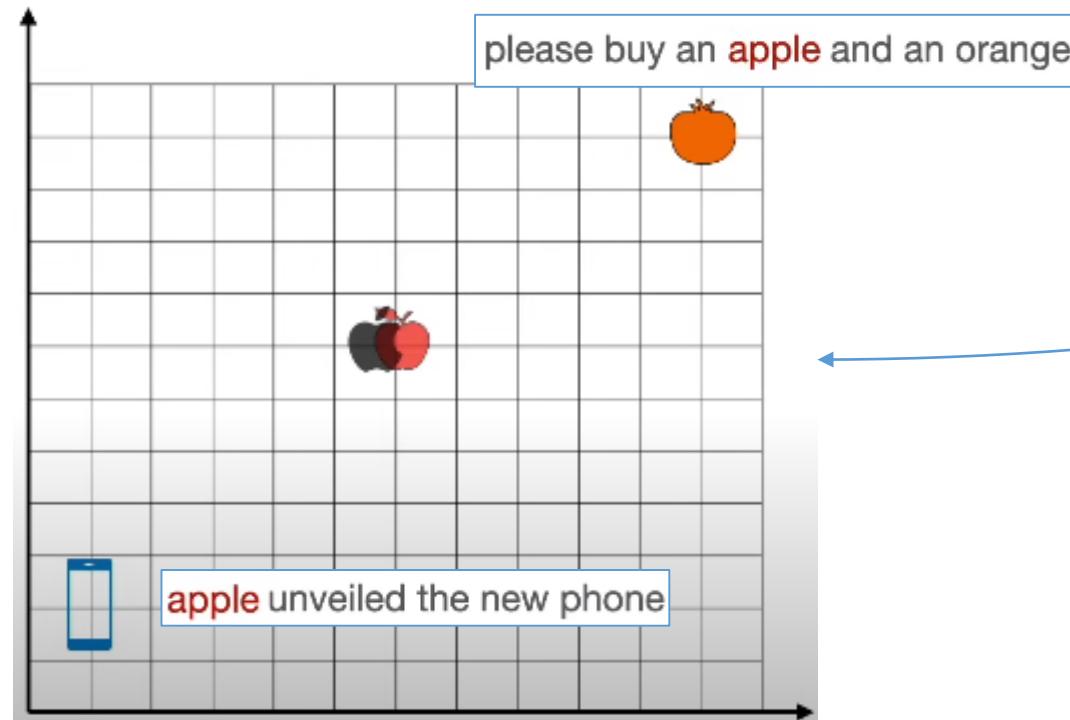


Language Models

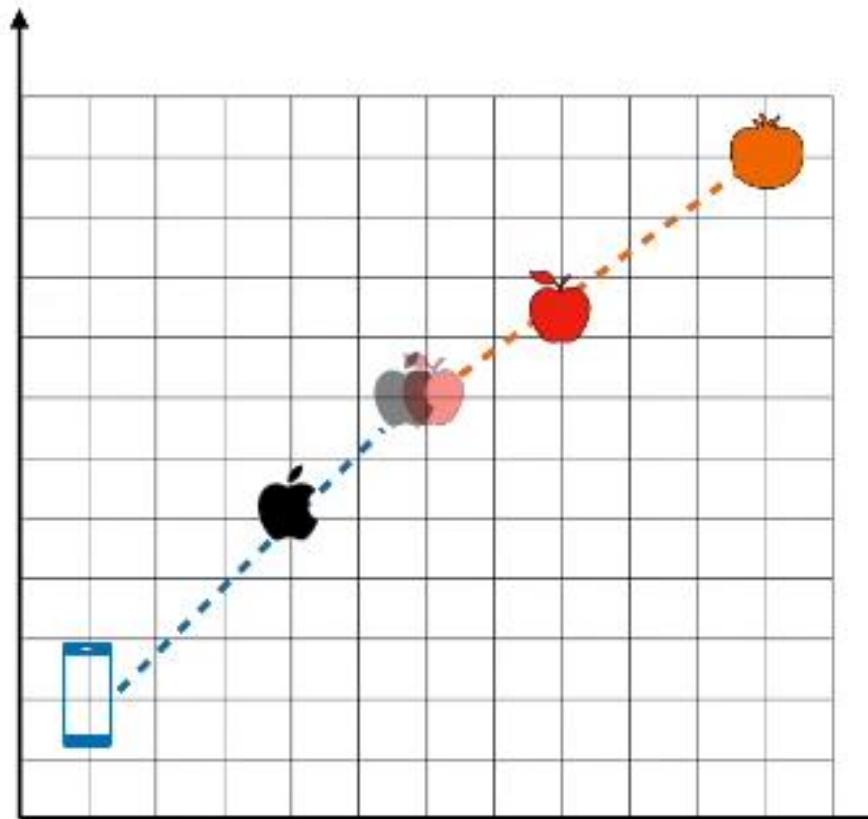
Transformer models

What is embedding?

Mechanism to put words or (long pieces of text) in such a way that similar words are placed closed to each others.



Words pulling words



please buy an **apple** and an **orange**

apple unveiled the new **phone**

please buy an **apple** and an **orange**

Language Models

Transformer models



Language Models

Transformer models

Similarity

Measure 1: Dot product

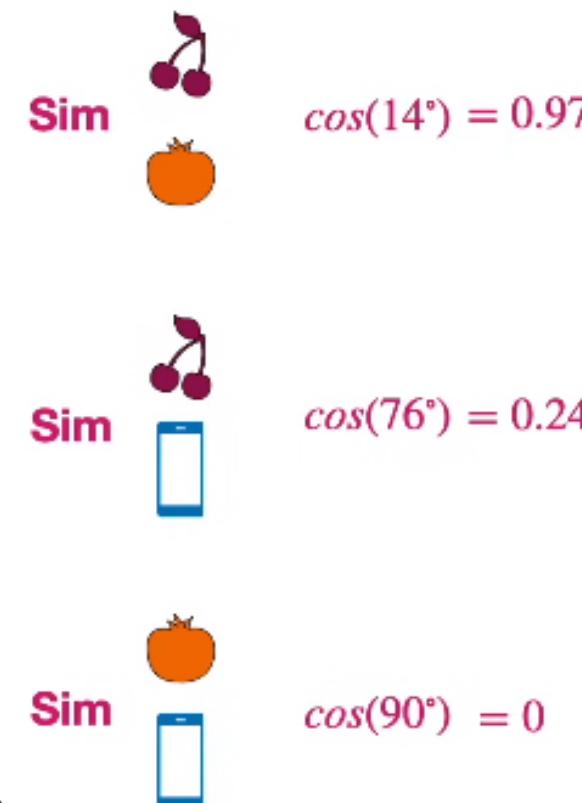
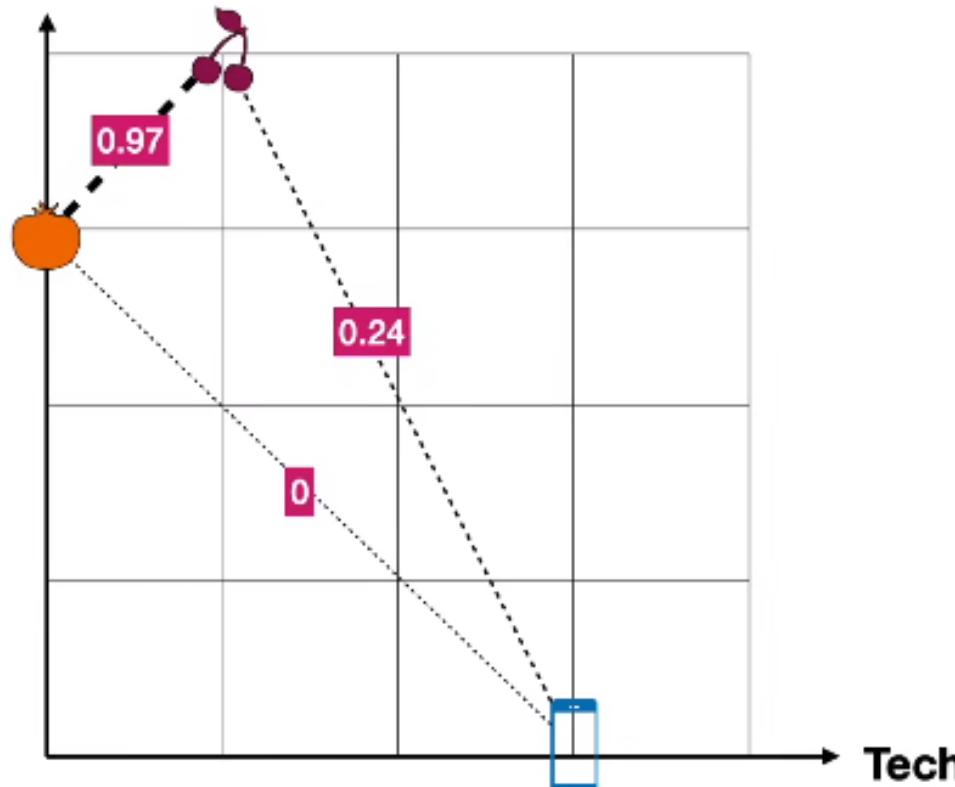


	Tech	Fruitness	
Sim	1	4	$1 \cdot 0 + 4 \cdot 3 = 12$
Sim	0	3	
Sim	1	4	$1 \cdot 3 + 4 \cdot 0 = 3$
Sim	3	0	
Sim	0	3	$0 \cdot 3 + 3 \cdot 0 = 0$
Sim	3	0	

Similarity

Measure 2: Cosine similarity

Fruitness



Similarity

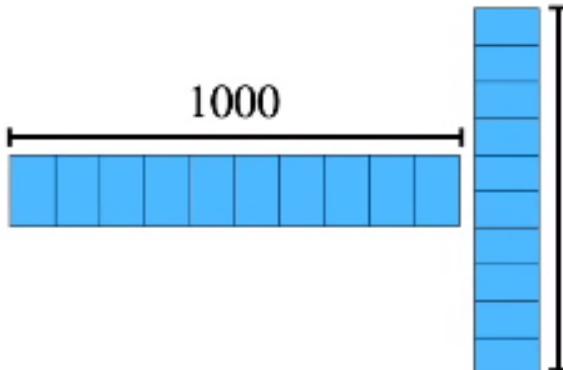
Measure 3: Scaled dot product

Dot product divided by the square root of the length of the vector



1	4
---	---

$$1 \cdot 0 + 4 \cdot 3 = 12 \longrightarrow \frac{12}{\sqrt{2}} = 8.49$$



0	3
---	---

Sim



1	4
---	---

$$1 \cdot 3 + 4 \cdot 0 = 3 \longrightarrow \frac{3}{\sqrt{2}} = 2.12$$

Sim



3	0
---	---

Sim



0	3
---	---

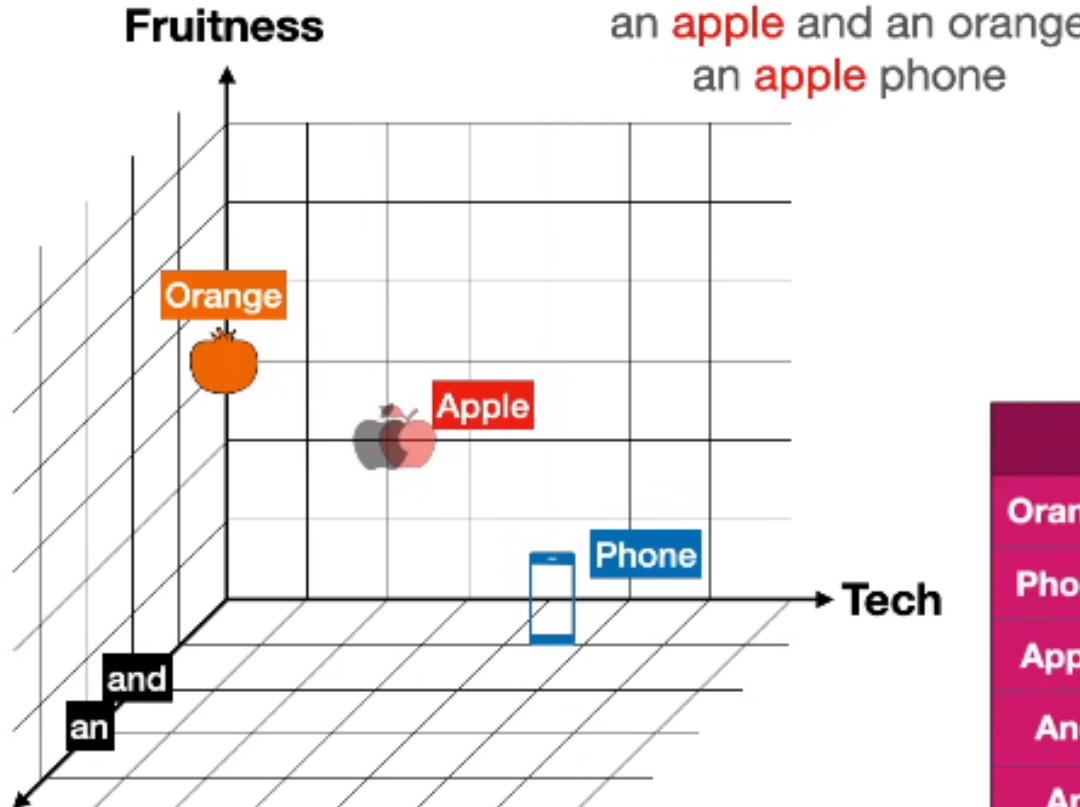
$$0 \cdot 3 + 3 \cdot 0 = 0 \longrightarrow \frac{0}{\sqrt{2}} = 0$$

3	0
---	---



Similarity

Cosine similarity

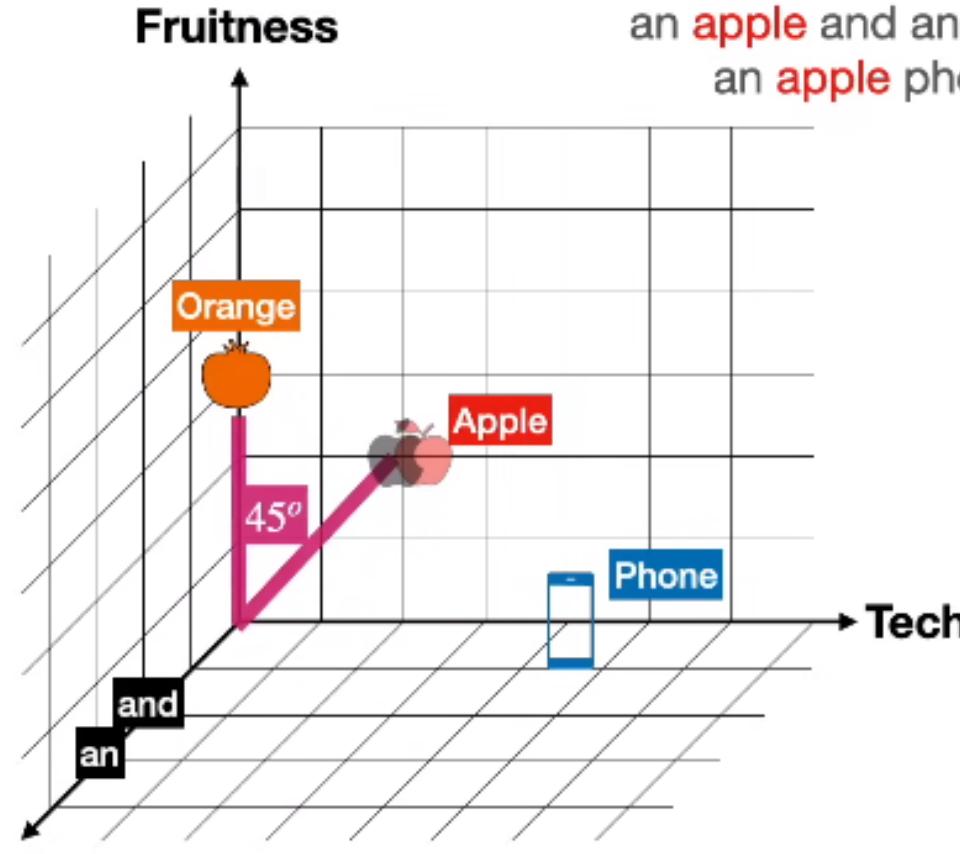


	Tech	Fruitness	Other
Orange	0	3	0
Phone	4	0	0
Apple	2	2	0
And	0	0	2
An	0	0	3

	Orange	Phone	Apple	And	An
Orange	1				
Phone		1			
Apple			1		
And				1	
An					1

Similarity

Cosine similarity



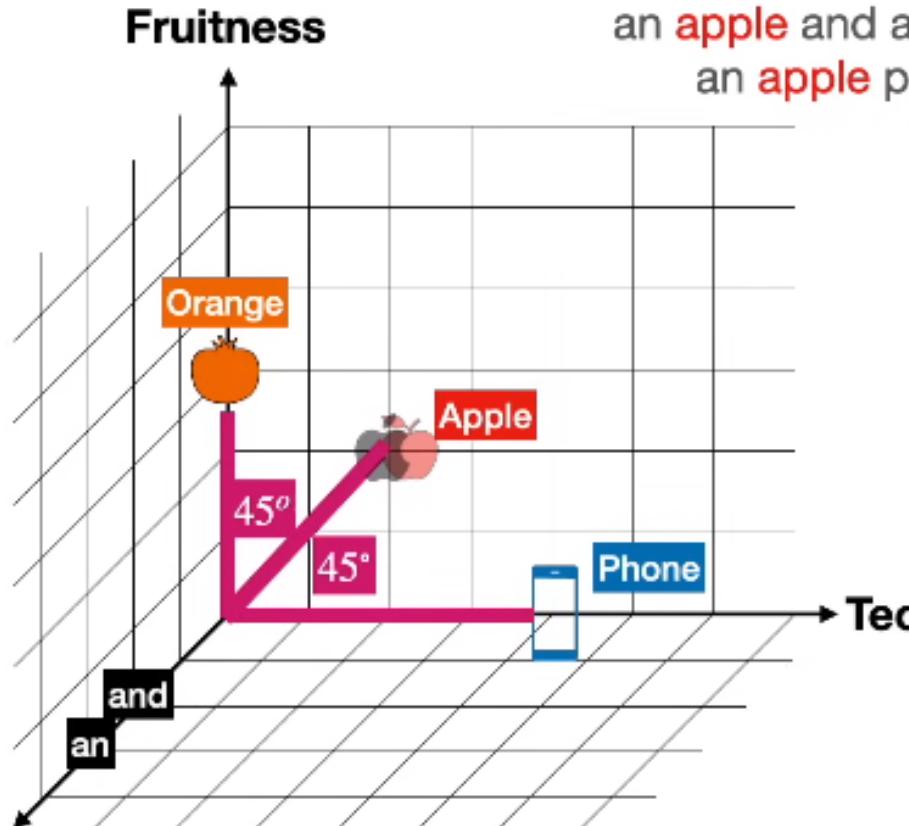
an **apple** and an orange
an **apple** phone

	Tech	Fruitness	Other
Orange	0	3	0
Phone	4	0	0
Apple	2	2	0
And	0	0	2
An	0	0	3

	Orange	Phone	Apple	And	An
Orange	1	0			
Phone	0	1			
Apple			1		
And				1	
An					1

Similarity

Cosine similarity



an **apple** and an orange
an **apple** phone

	Tech	Fruitness	Other
Orange	0	3	0
Phone	4	0	0
Apple	2	2	0
And	0	0	2
An	0	0	3

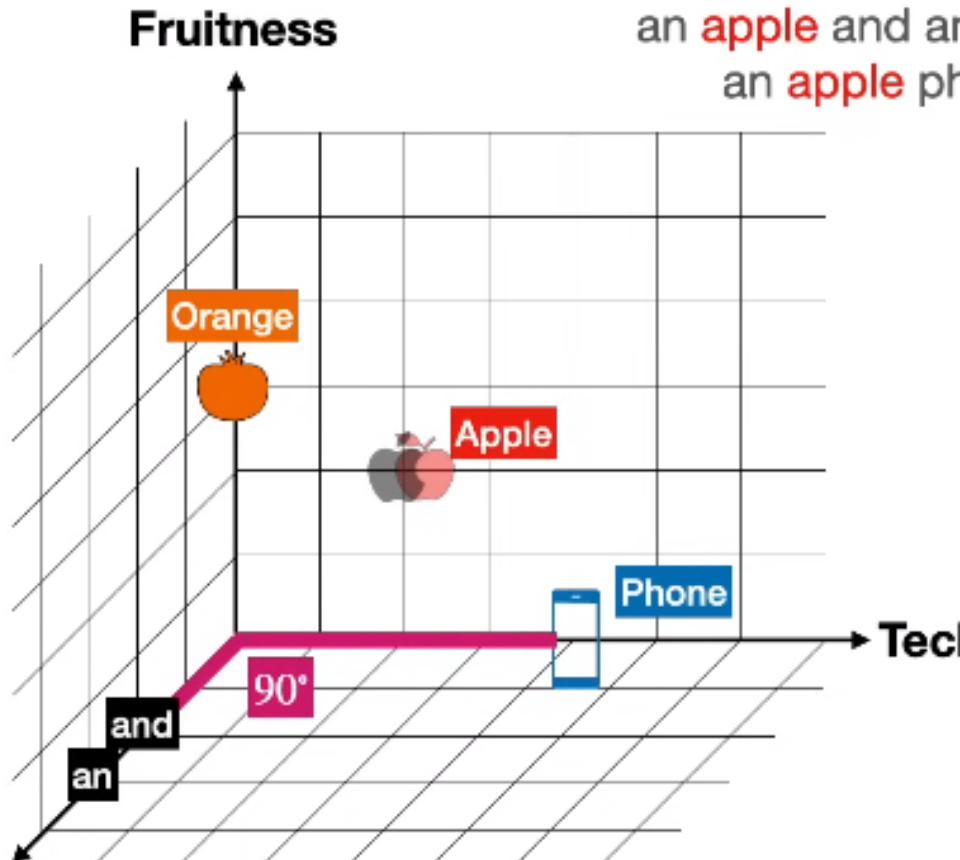
	Orange	Phone	Apple	And	An
Orange	1	0	0.71		
Phone	0	1	0.71		
Apple	0.71	0.71	1		
And				1	
An					1

Language Models

Transformer models

Similarity

Cosine similarity



	Tech	Fruitness	Other
Orange	0	3	0
Phone	4	0	0
Apple	2	2	0
And	0	0	2
An	0	0	3

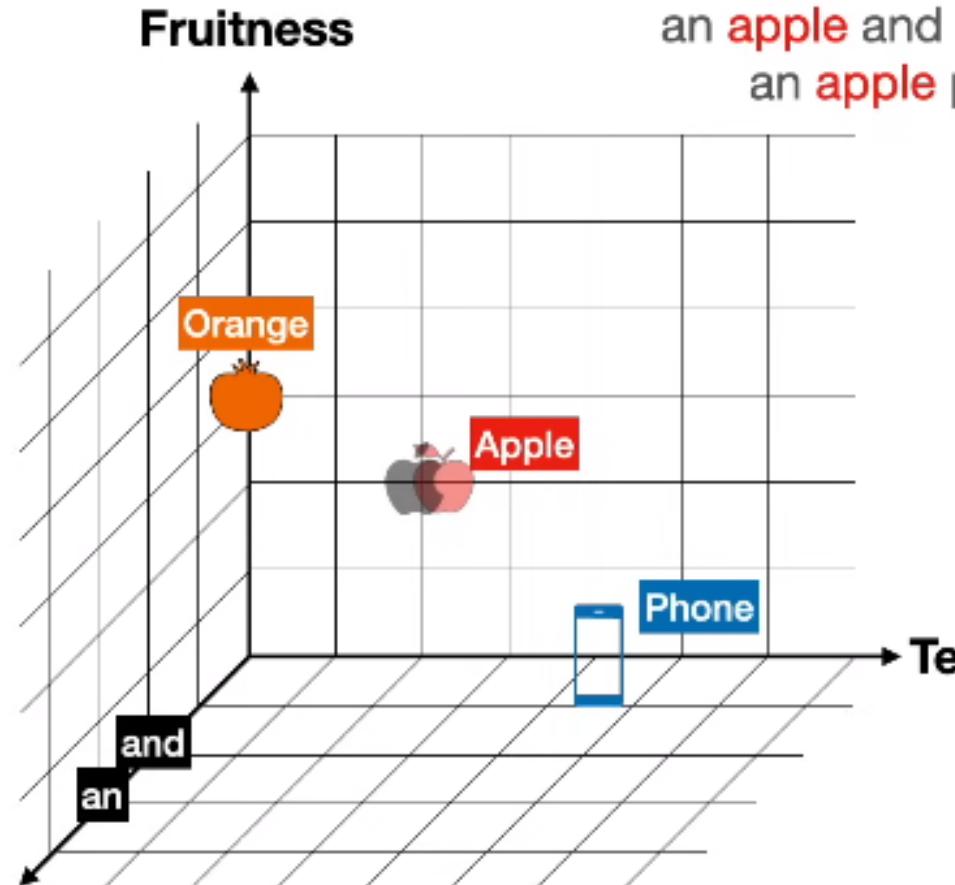
	Orange	Phone	Apple	And	An
Orange	1	0	0.71		
Phone	0	1	0.71		
Apple	0.71	0.71	1		
And				1	
An					1

Language Models

Transformer models

Similarity

Cosine similarity



an **apple** and an orange
an **apple** phone

	Tech	Fruitness	Other
Orange	0	3	0
Phone	4	0	0
Apple	2	2	0
And	0	0	2
An	0	0	3

	Orange	Phone	Apple	And	An
Orange	1	0	0.71	0	0
Phone	0	1	0.71	0	0
Apple	0.71	0.71	1	0	0
And	0	0	0	1	1
An	0	0	0	1	1

Similarity

Word math

an **apple** and an orange

	Orange	Apple	And	An
Orange	1	0.71	0	0
Apple	0.71	1	0	0
And	0	0	1	1
An	0	0	1	1

$$\text{Orange} \rightarrow 1 \text{ Orange} + 0.71 \text{ Apple}$$

$$\text{Apple} \rightarrow 0.71 \text{ Orange} + 1 \text{ Apple}$$

$$\text{And} \rightarrow 1 \text{ And} + 1 \text{ An}$$

$$\text{An} \rightarrow 1 \text{ An} + 1 \text{ And}$$

Language Models

Transformer models

Similarity

an apple phone

	Phone	Apple	An
Phone	1	0.71	0
Apple	0.71	1	0
An	0	0	1

Phone → 1 Phone + 0.71 Apple

Apple → 0.71 Phone + 1 Apple

An → 1 An

Similarity

Normalization

Want coefficients to add to 1

$$\text{Orange} \rightarrow \frac{1 \text{ Orange} + 0.71 \text{ Apple}}{1 + 0.71} = 0.58 \text{ Orange} + 0.42 \text{ Apple}$$

Need coefficients to be positive

!

$$\text{Orange} \rightarrow \frac{1 \text{ Orange} - 1 \text{ Motorcycle}}{1 - 1} = \times$$

Similarity

Softmax

$$x \rightarrow e^x$$

$$\text{Orange} \rightarrow \frac{e^1 \text{Orange} + e^{0.71} \text{Apple}}{e^1 + e^{0.71}} = 0.57 \text{Orange} + 0.43 \text{Apple}$$



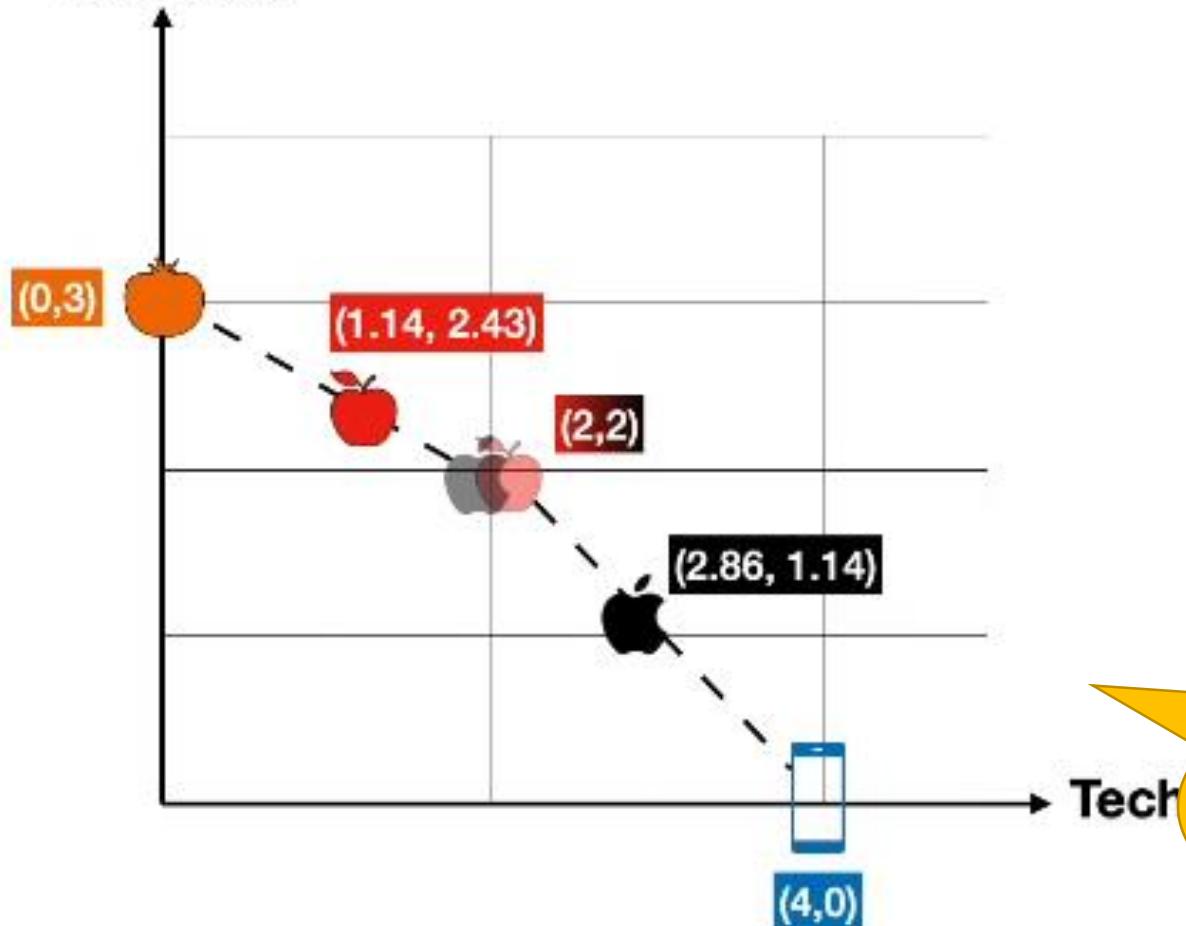
$$\text{Orange} \rightarrow \frac{e^1 \text{Orange} + e^{-1} \text{Motorcycle}}{e^1 + e^{-1}} = 0.88 \text{Orange} + 0.12 \text{Motorcycle}$$

Language Models

Transformer models

Similarity

Fruitness



an **apple** and an orange

Apple → 0.43 **Orange** + 0.57 **Apple**

an **apple** phone

Apple → 0.43 **Phone** + 0.57 **Apple**

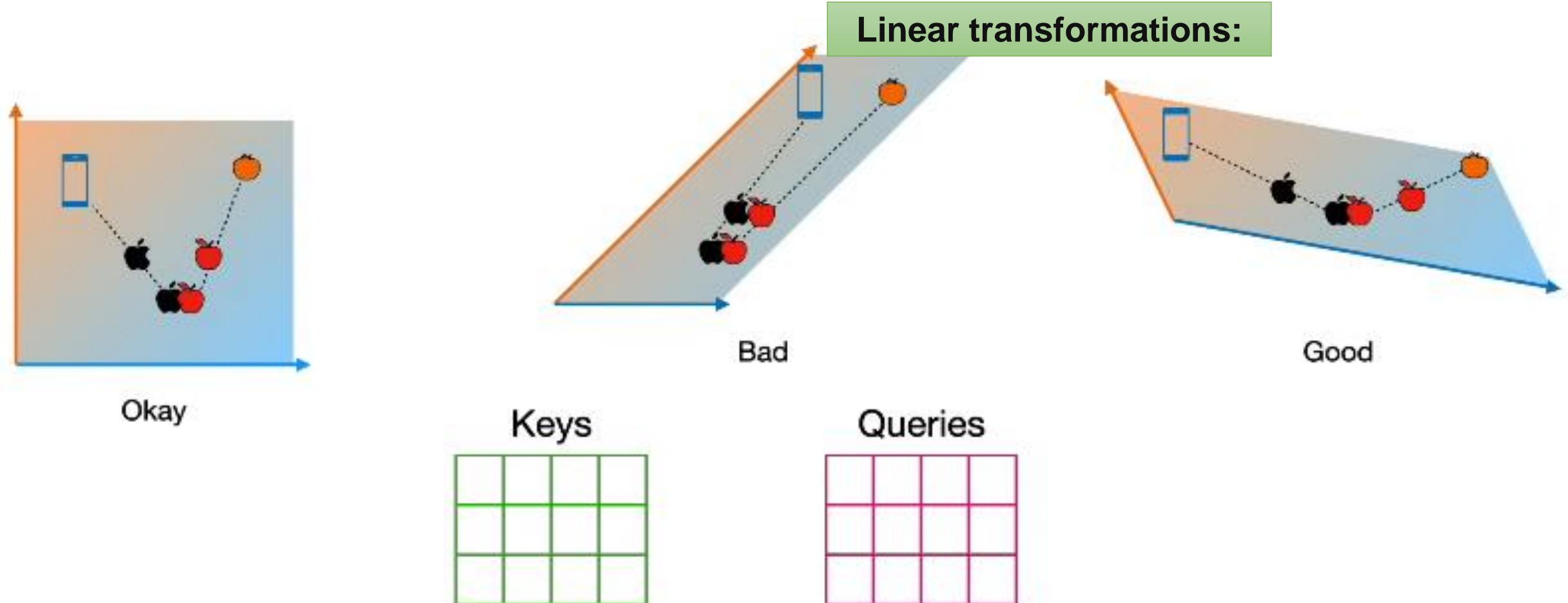
Two different
coordinates of the same
word with different
context

In transformer model,
exactly with attention
mechanism, this
operation is performed
many times

Language Models

Transformer models

Keys, Query and Values

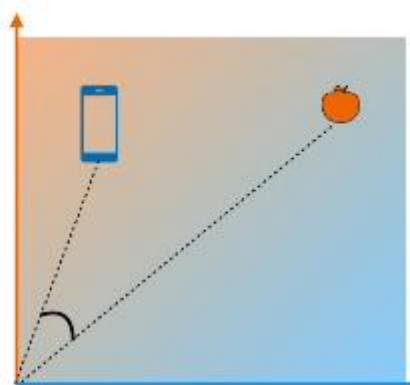


Language Models

Keys, Query and Values

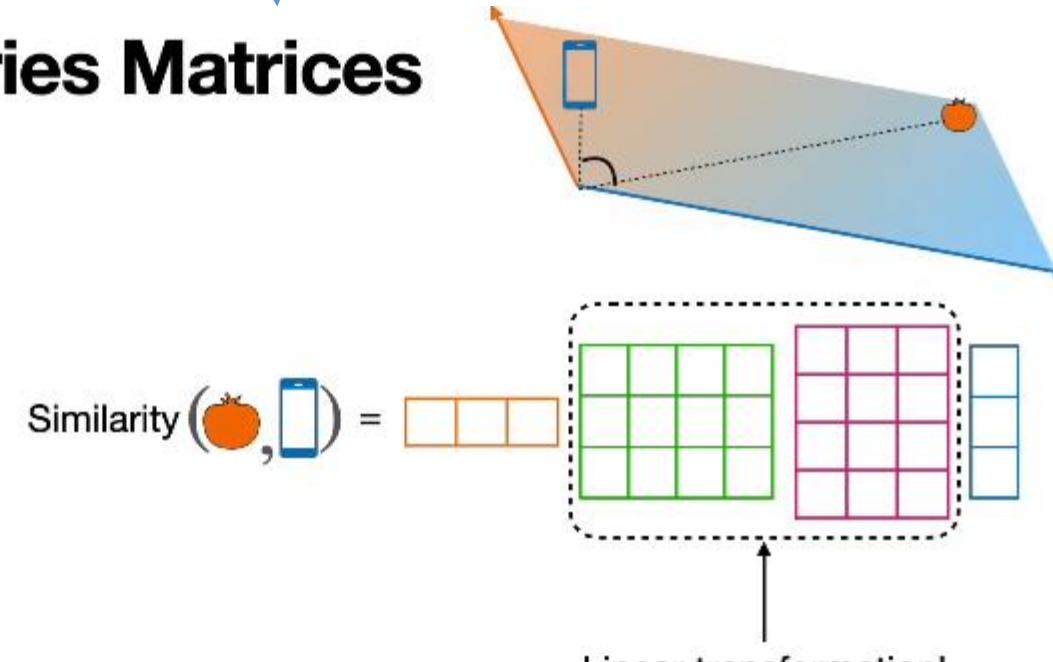
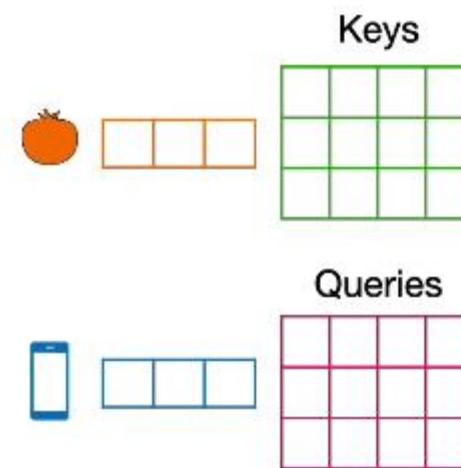
Transformer models

Similarity

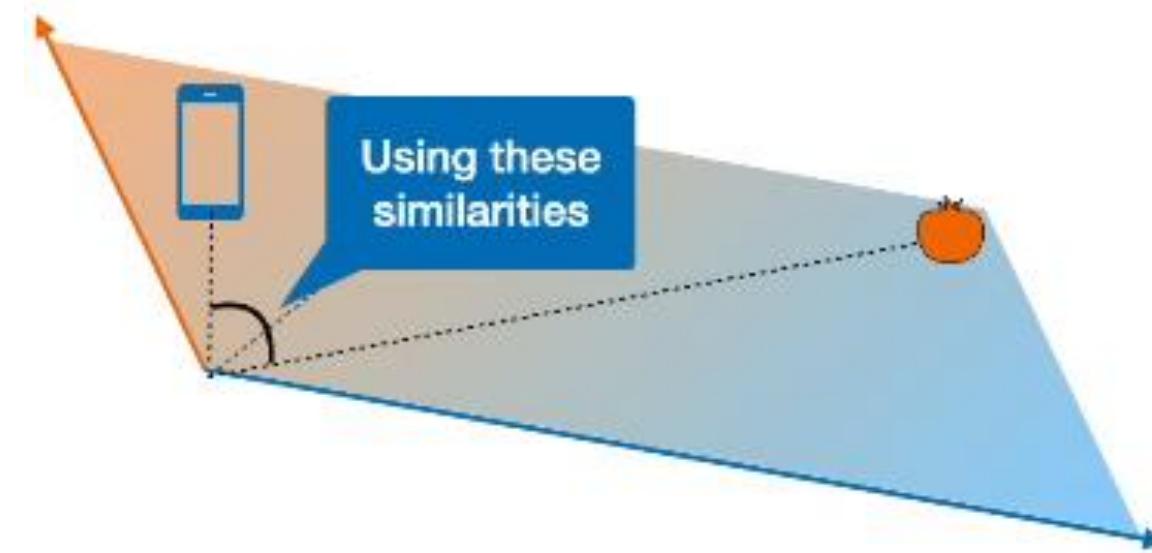


$$\text{Similarity}(\text{apple}, \text{phone}) = \begin{matrix} \text{apple} \\ \text{phone} \end{matrix} \cdot \begin{matrix} \text{apple} \\ \text{phone} \end{matrix}$$

Keys and Queries Matrices



Values matrix

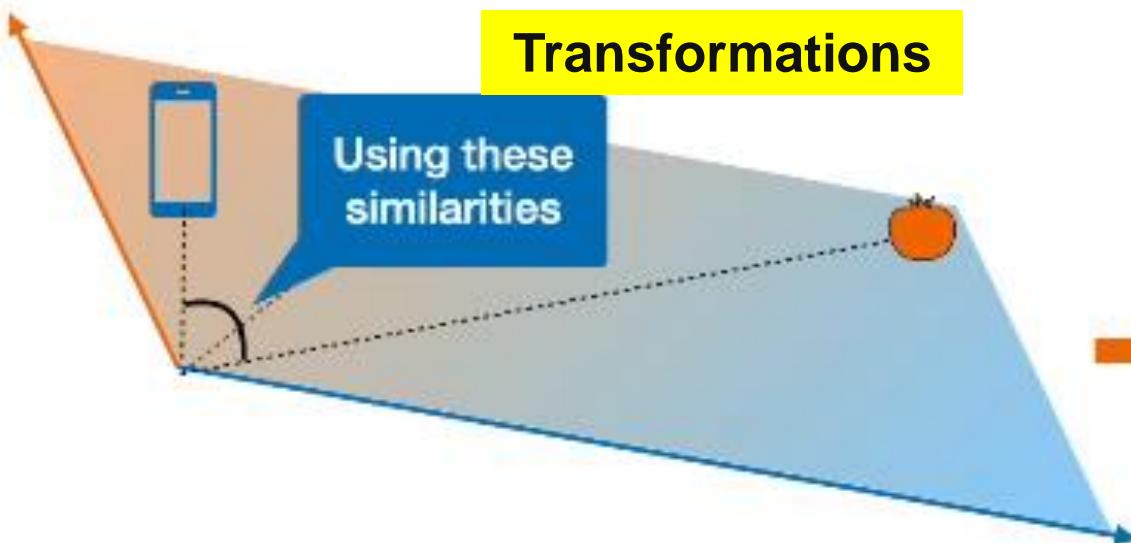


Best embedding for finding similarities

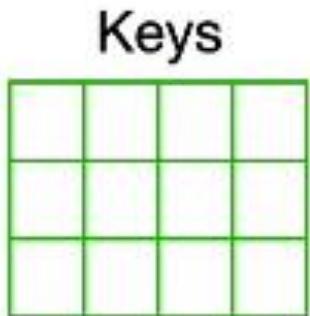


Best embedding for finding the next word

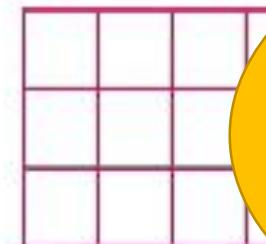
Values matrix



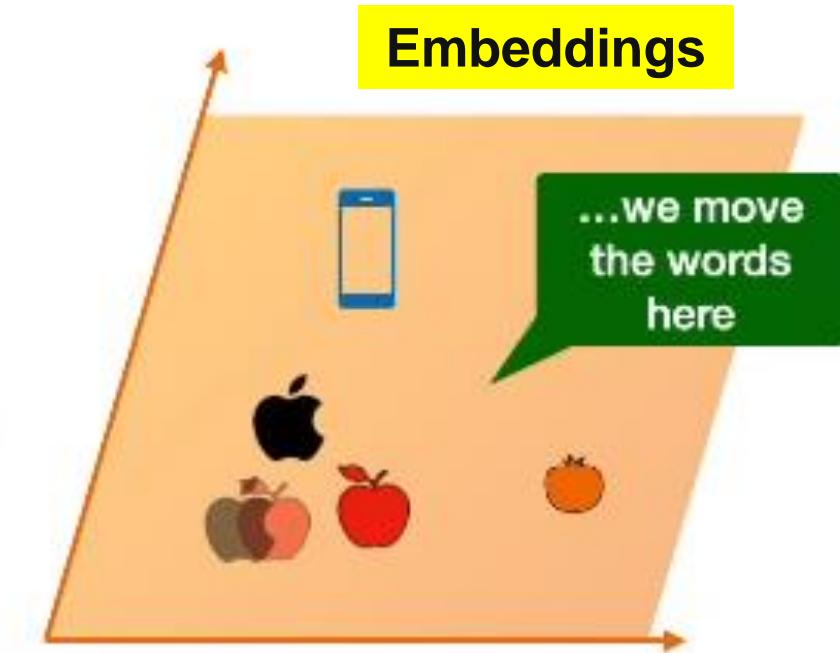
Best embedding for finding similarities



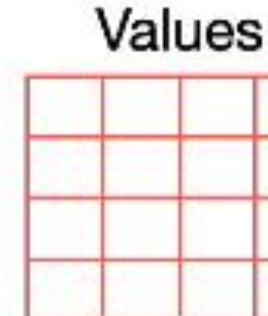
Queries



Embeddings here is responsible of capturing the color, size, technology, flavour, etc. (Features)



Best embedding for finding the next word



Values

Embeddings here is responsible of capturing the next word in the sentence

Value matrix

an **apple** and an orange

	Orange	Apple	And	An
Orange	0.4	0.3	0.15	0.15
Apple	0.3	0.4	0.15	0.15
And	0.15	0.15	0.5	0.5
An	0.15	0.15	0.5	0.5

Value matrix

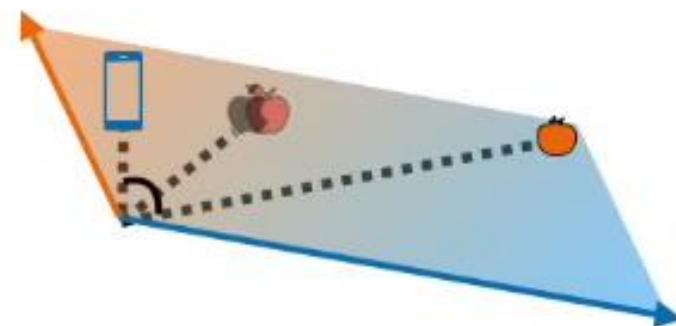
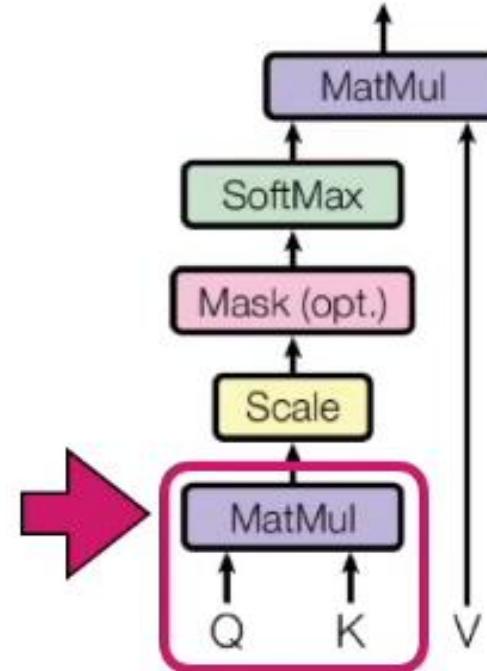
=

	Orange	Apple	And	An
Orange	v_{11}	v_{12}	v_{13}	v_{14}
Apple	v_{21}	v_{22}	v_{23}	v_{24}
And	v_{31}	v_{32}	v_{33}	v_{34}
An	v_{41}	v_{42}	v_{43}	v_{44}

apple \longrightarrow $0.3 \cdot \text{orange}$
 $+0.4 \cdot \text{apple}$
 $+0.15 \cdot \text{and}$
 $+0.15 \cdot \text{an}$

apple \longrightarrow $v_{21} \cdot \text{orange}$
 $+v_{22} \cdot \text{apple}$
 $+v_{23} \cdot \text{and}$
 $+v_{24} \cdot \text{an}$

Scaled Dot-Product Attention



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

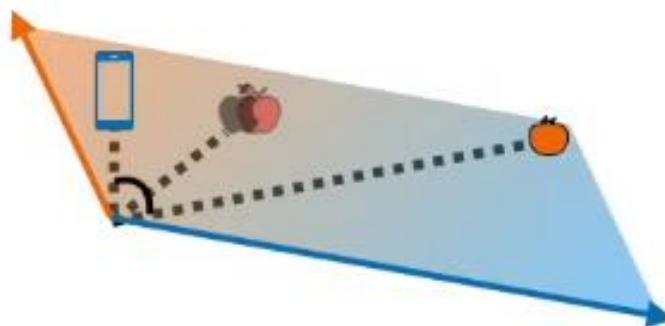
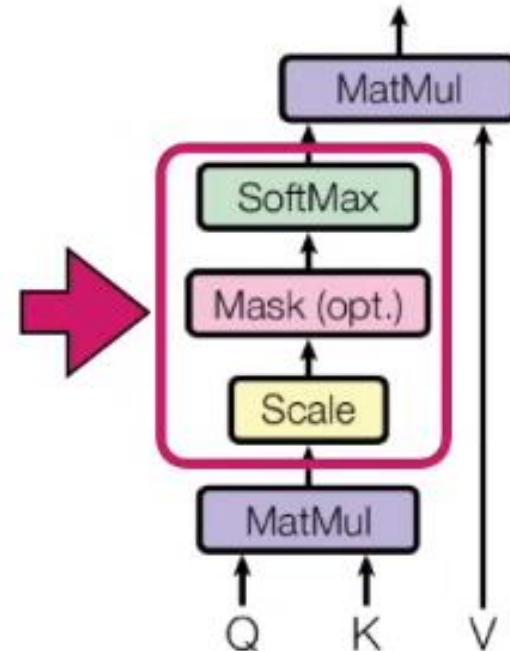
Two square matrices are shown side-by-side. The left matrix, labeled K , has a green grid pattern. The right matrix, labeled Q , has a pink grid pattern.

Language Models

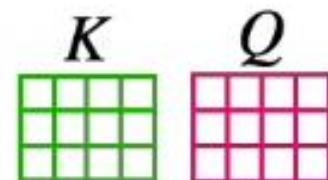
Attention Layer

Transformer models

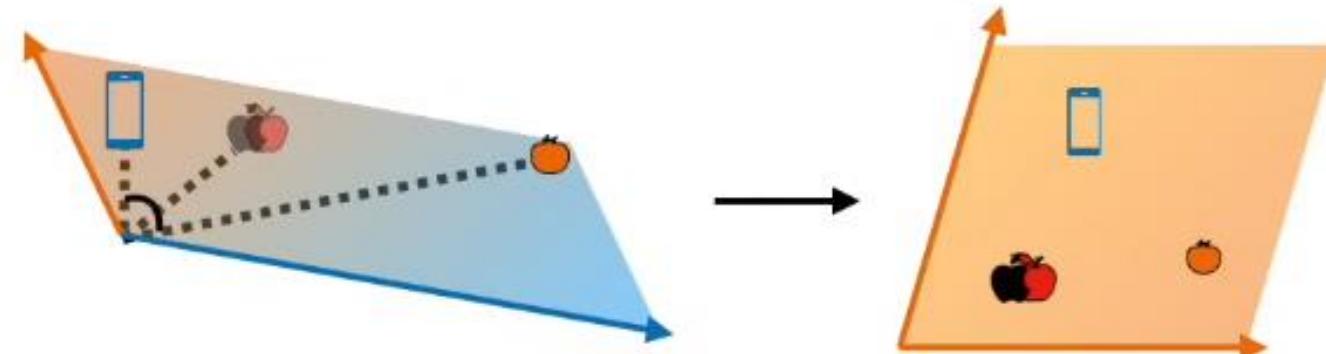
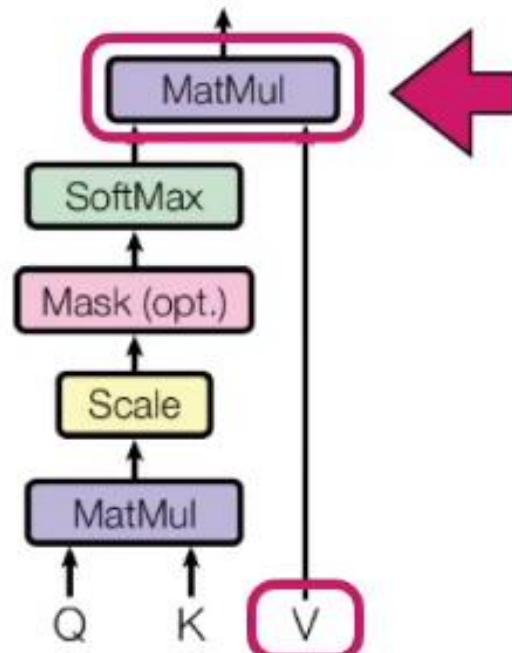
Scaled Dot-Product Attention



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



Scaled Dot-Product Attention



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

$$K \quad Q$$

Two square matrices, K and Q, are shown side-by-side. Matrix K is green and matrix Q is pink, representing key and query vectors respectively.

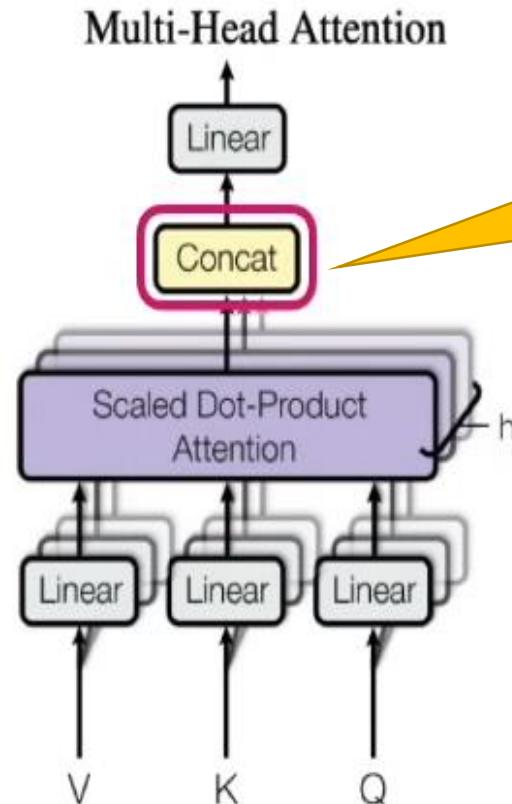
$$V$$

A single square matrix, V, is shown in red, representing the value vector.

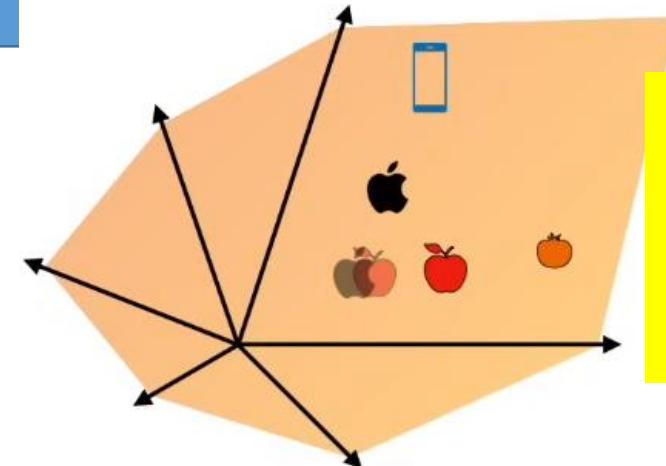
Language Models

Multi-Head Attention Layer

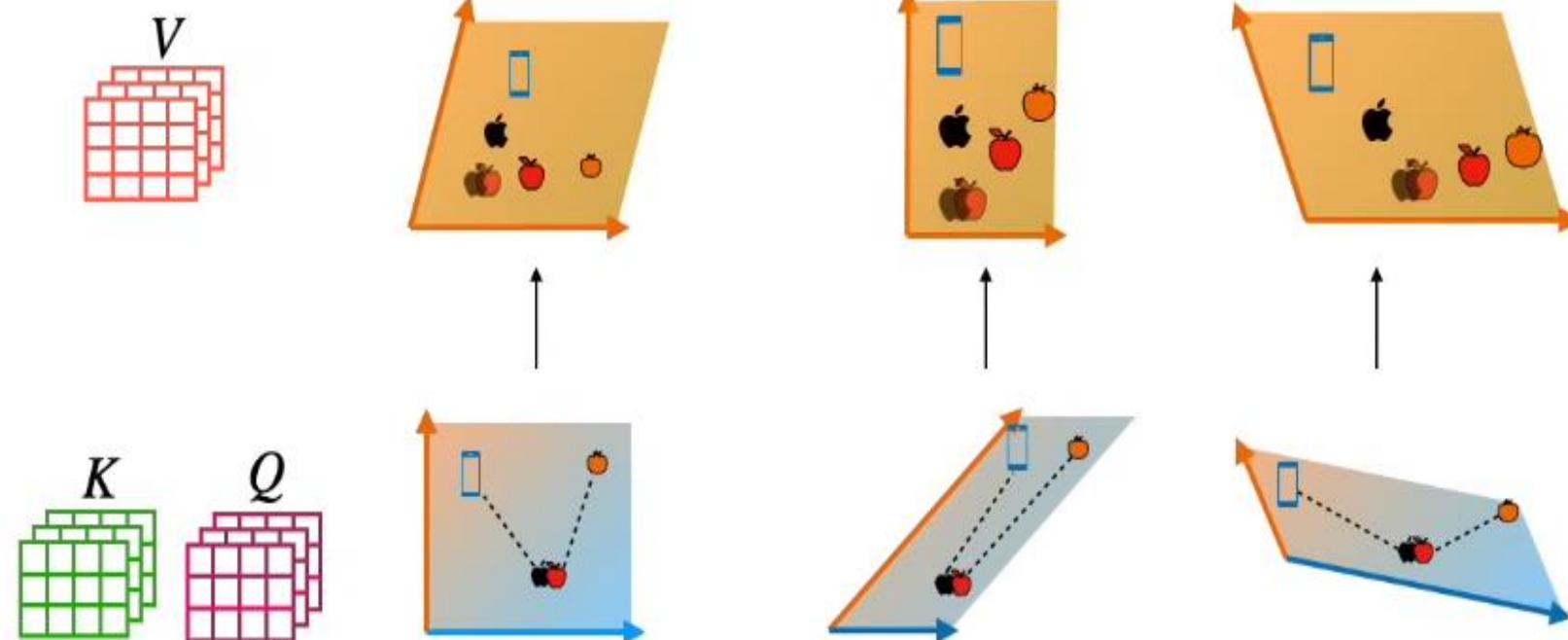
Transformer models



Concatenation of all embeddings since we don't know which one is better!



Getting higher dimensionality because of concatenation layer



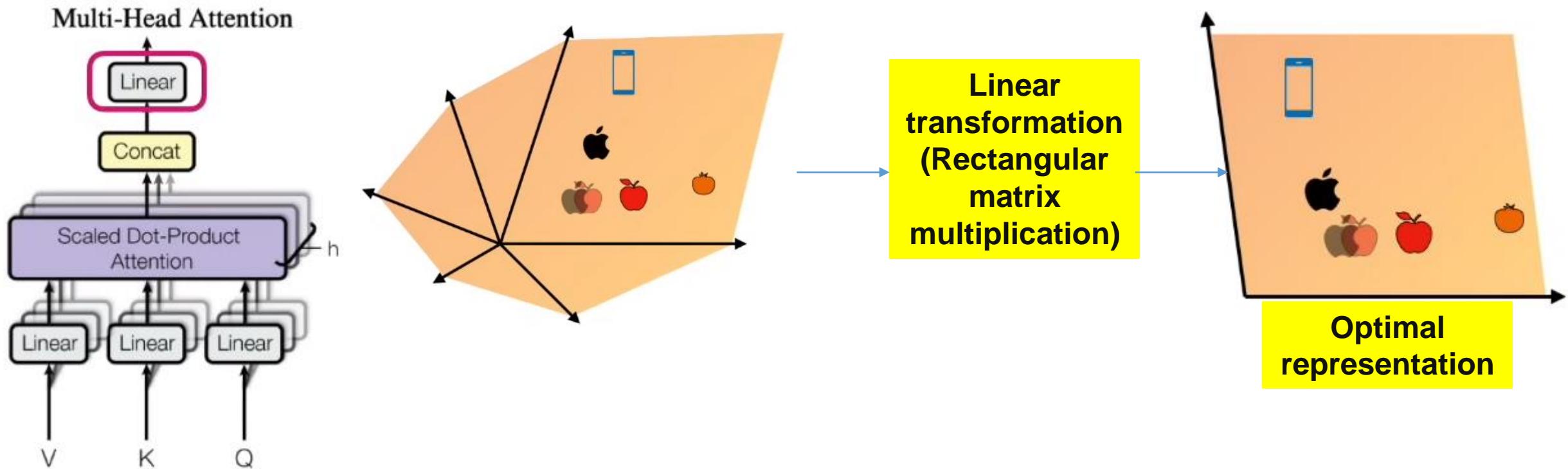
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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

Language Models

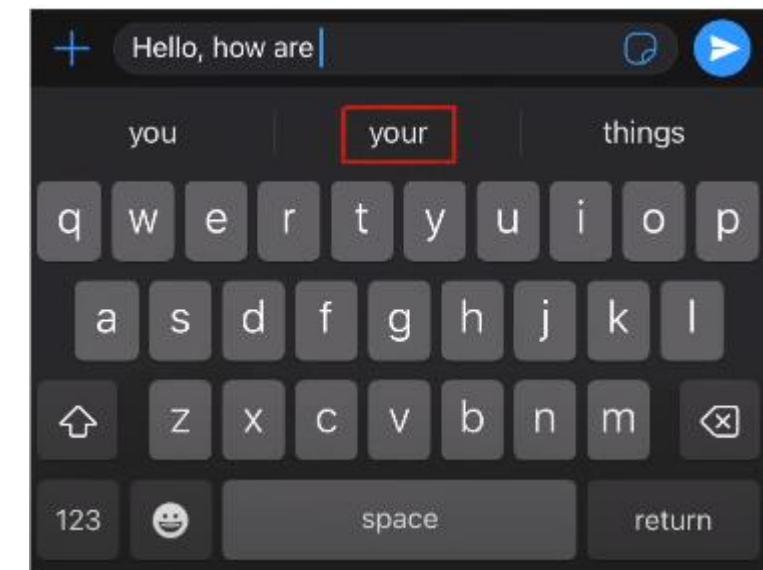
Multi-Head Attention Layer

Transformer models



Language Models

Transformer models



Language Models

ngrams

Ngram: method of Looking at N words
that appears before!

3grams example

Hello, How are



Data

Hello, how are you?
Hello, how are **things** going?
Hello, how are **things** today?
Hello, how are **the** kids?
Hello, how are **the** others?
Hello, how are **they** doing?
Hello, how are **things** happening?

1gram example

Hello, How are



Transformer models

Data

... are you ...

... are sad ...

... are?...

... are happy ...

... are ready ...

... are happy ...

... are free ...

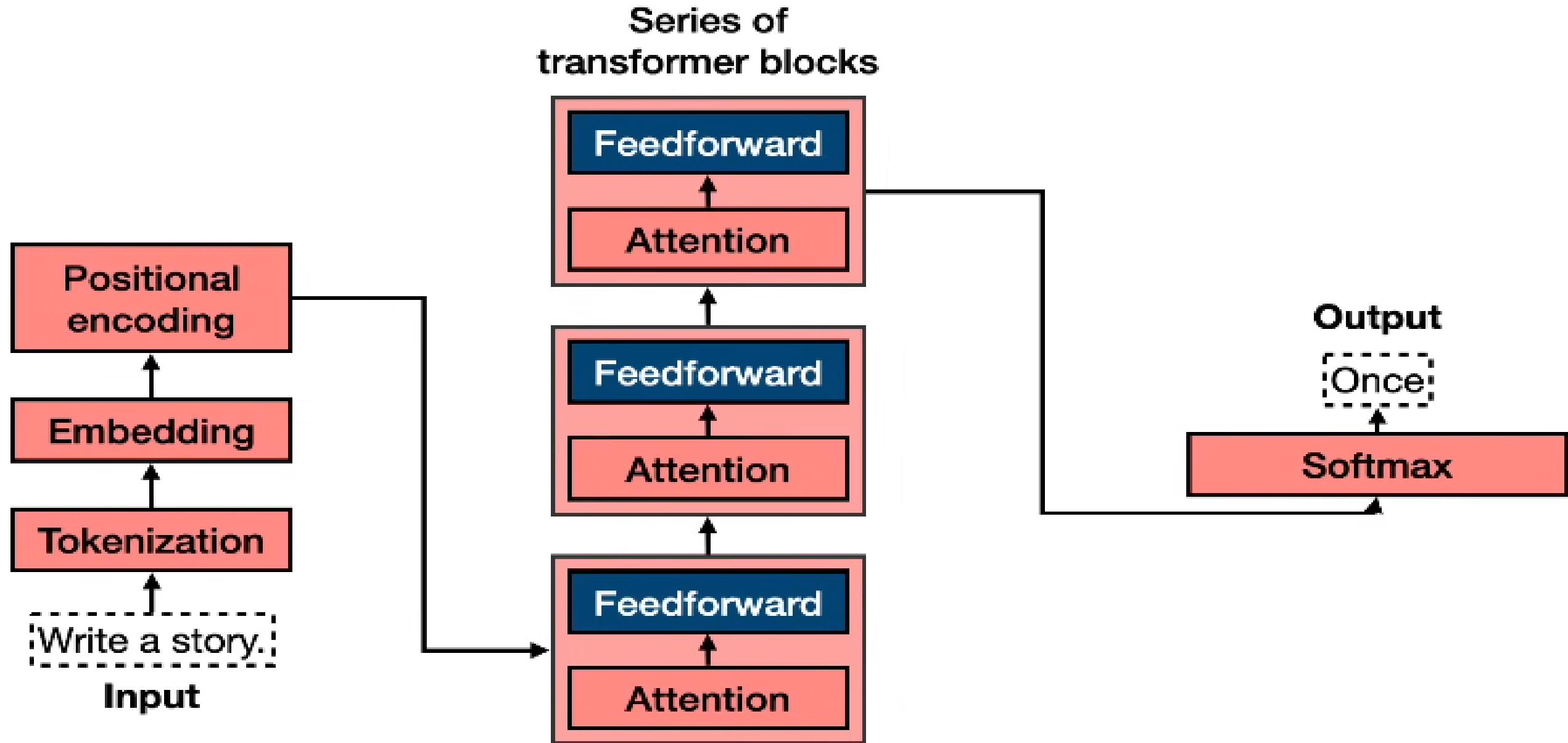
Data

That day where
I was singing
and two clouds
appeared

???

Language Models

Transformer models



Language Models

Transformer models

Step 1: Tokenization

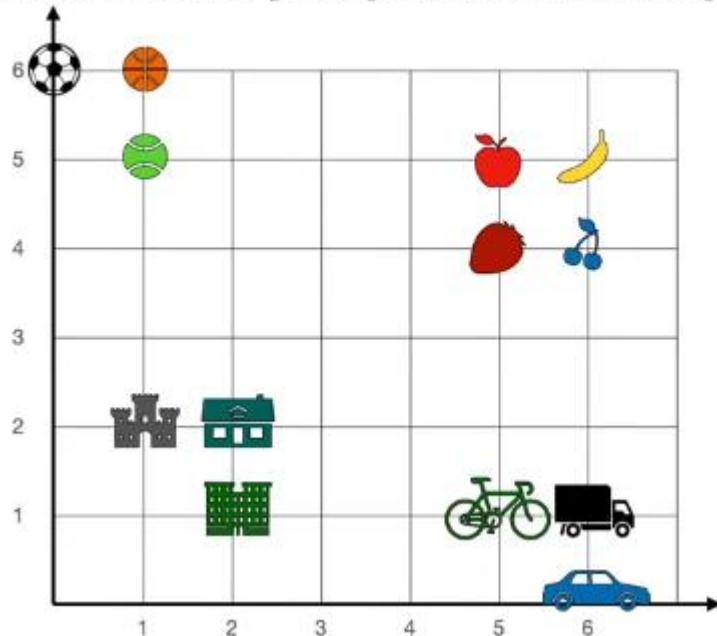
Split sentence into segments (words or tokens)

Write a story. → Write A story .

doesn't → does n't

Step 2: Embeddings

Where would you put the word apple?

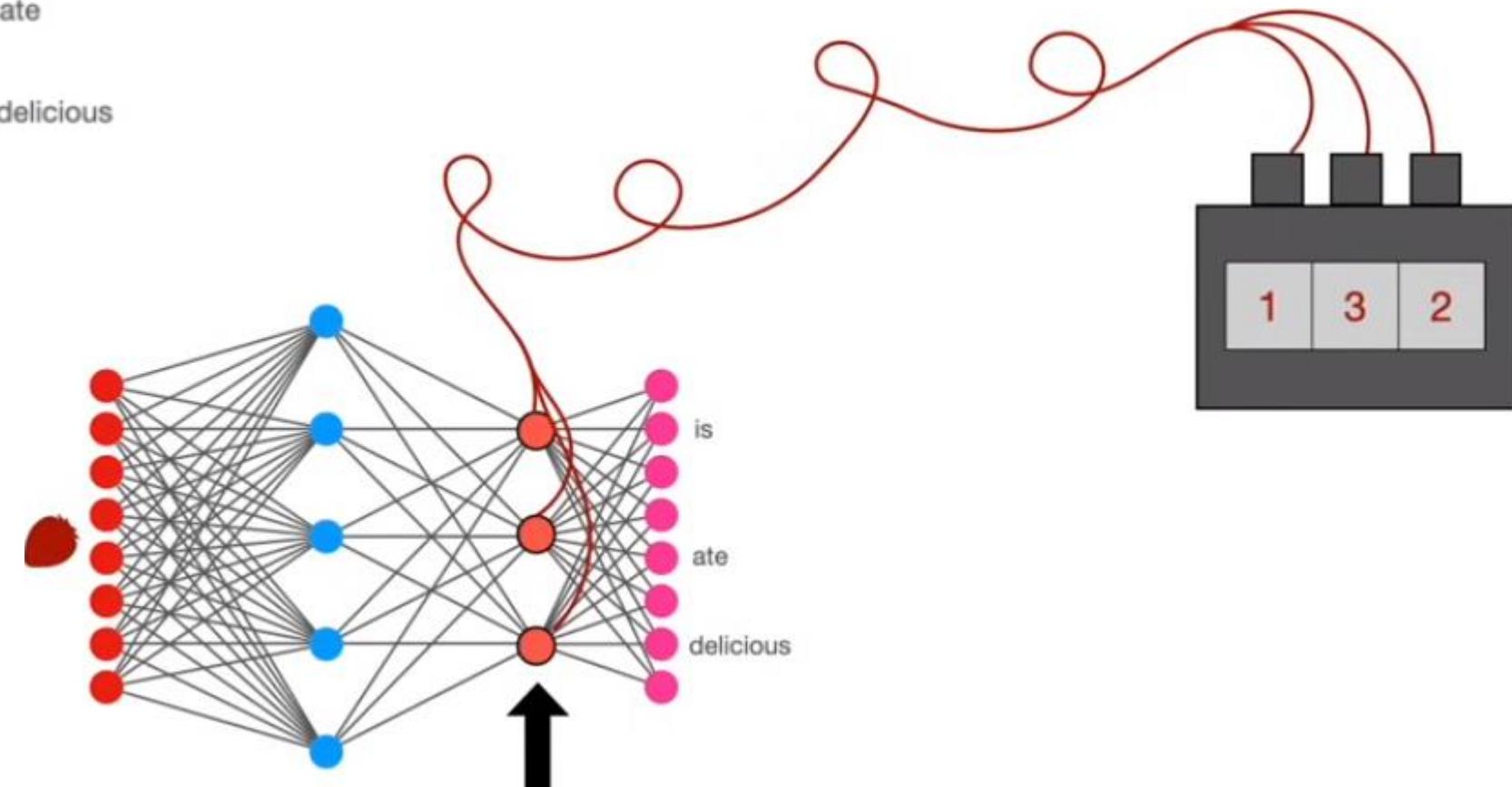
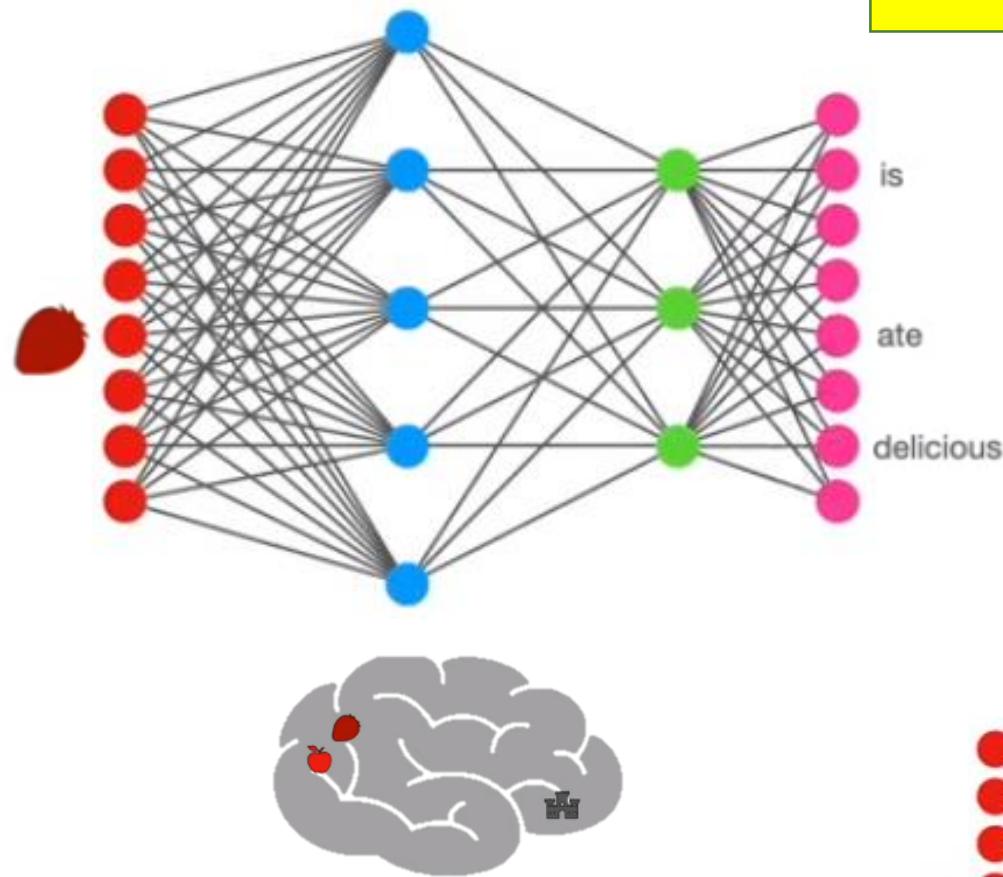


Word	Numbers	
Apple	5	5
Banana	6	5
Strawberry	5	4
Cherry	6	4
Soccer	0	6
Basketball	1	6
Tennis	1	5
Castle	1	2
House	2	2
Building	2	1
Bicycle	5	1
Truck	6	1
Car	6	0

Language Models

Step 2: Embeddings

Word2Vec embedding method



Transformer models

Word	Numbers		
Strawberry	1	3	2
Apple	1.1	2.9	2.2
Castle	7	-5.4	0.4

Language Models

Step 3: Positional Encoding

Transformer models

What is positional encoding (first part of encoder)

Index value, can't be used to represent an item's position in transformer models (For long sequences, the indices can grow large in magnitude).

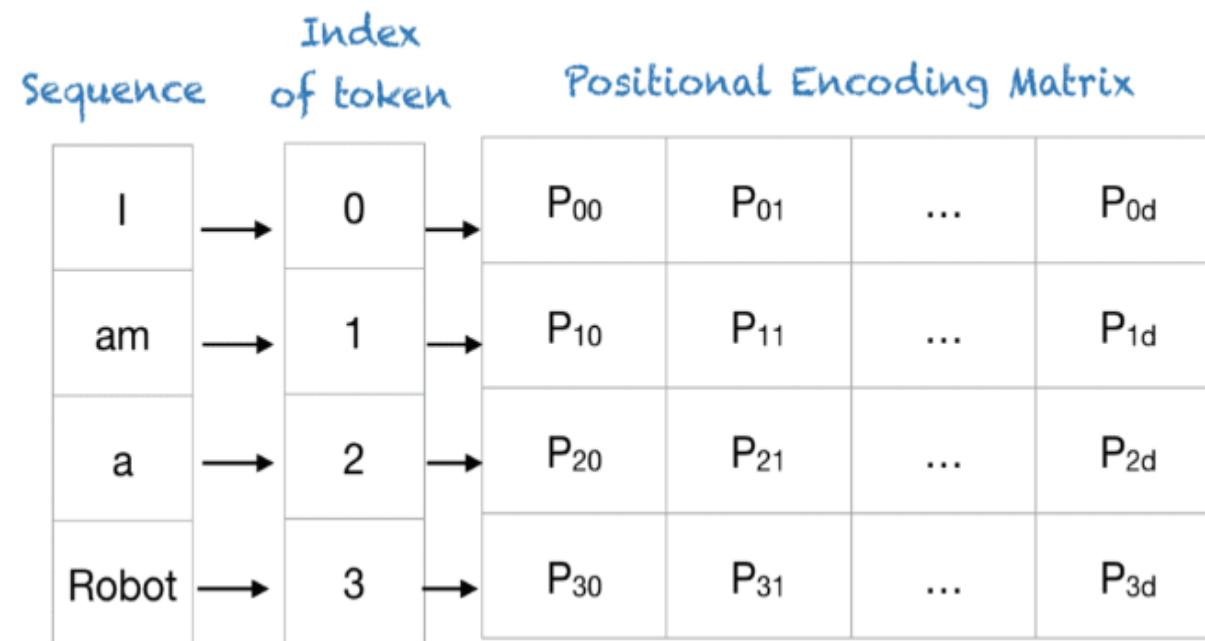
Normalize the index value to lie between 0 and 1 can create problems for variable length sequences as they would be normalized differently.

Transformer used smart positional encoding.

Describes the location or **position** of an entity in a sequence so that each position is assigned a **unique representation**.

Each position/index is mapped to a vector.

output of the positional encoding layer is a matrix, where each row of the matrix represents an encoded object of the sequence summed with its positional information.



Positional Encoding Concept

Taking into account the **order** or words in the context.

Let's assume that we have an input sequence of length L

We need to know the position of the Kth object (word).

The positional encoding is given by sine and cosine functions of varying frequencies:

P(k,j): **Position function**
for mapping a
position **K** in the input
sequence to index **(k,j)** of
the positional matrix

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

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

$$0 \leq i < d/2$$

d: Dimension of the output
embedding space

n: User-defined scalar, set to
10,000 by the authors
of [Attention Is All You Need](#).

Language Models

Step 3: Positional Encoding

Transformer models

Positional Encoding Example “I am a robot” with n=100 and d=4

Sequence	Index of token,	Positional Encoding			
		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

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

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

Positional Encoding *Python implementation*

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def getPositionEncoding(seq_len, d, n=10000):
5     P = np.zeros((seq_len, d))
6     for k in range(seq_len):
7         for i in np.arange(int(d/2)):
8             denominator = np.power(n, 2*i/d)
9             P[k, 2*i] = np.sin(k/denominator)
10            P[k, 2*i+1] = np.cos(k/denominator)
11    return P
12
13 P = getPositionEncoding(seq_len=4, d=4, n=100)
14 print(P)
```

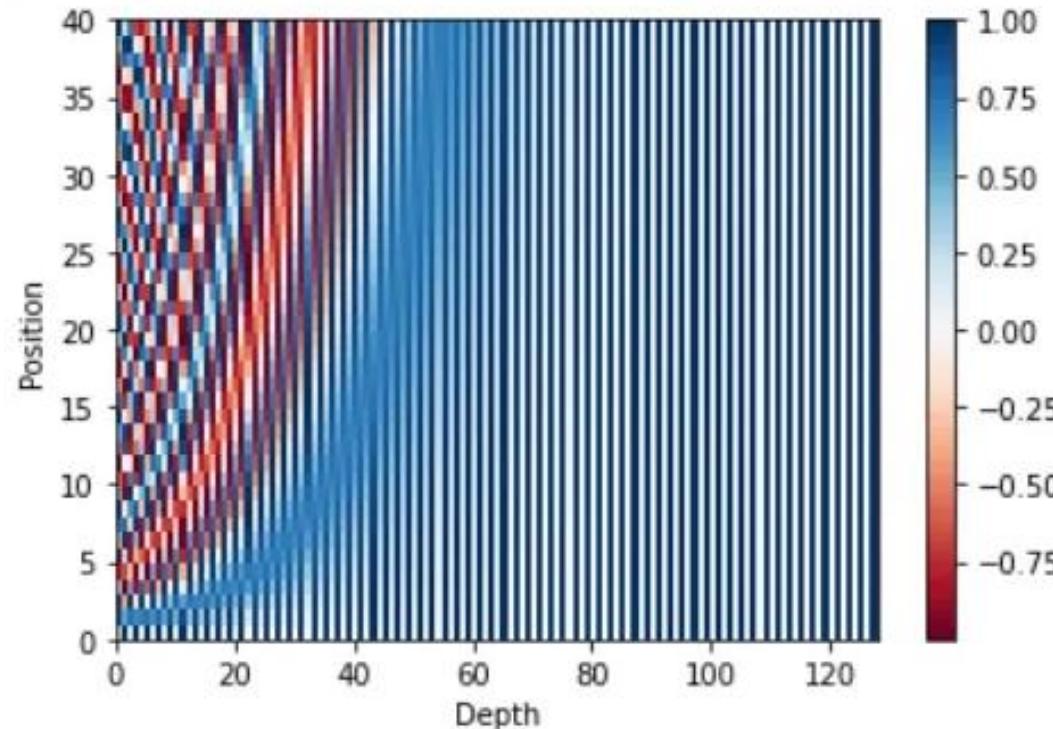
$$P(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$

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

```
1 [[ 0.          1.          0.          1.          ]
2  [ 0.84147098  0.54030231  0.09983342  0.99500417]
3  [ 0.90929743 -0.41614684  0.19866933  0.98006658]
4  [ 0.14112001 -0.9899925   0.29552021  0.95533649]]
```

Positional Encoding (*Why sin and cosine*)

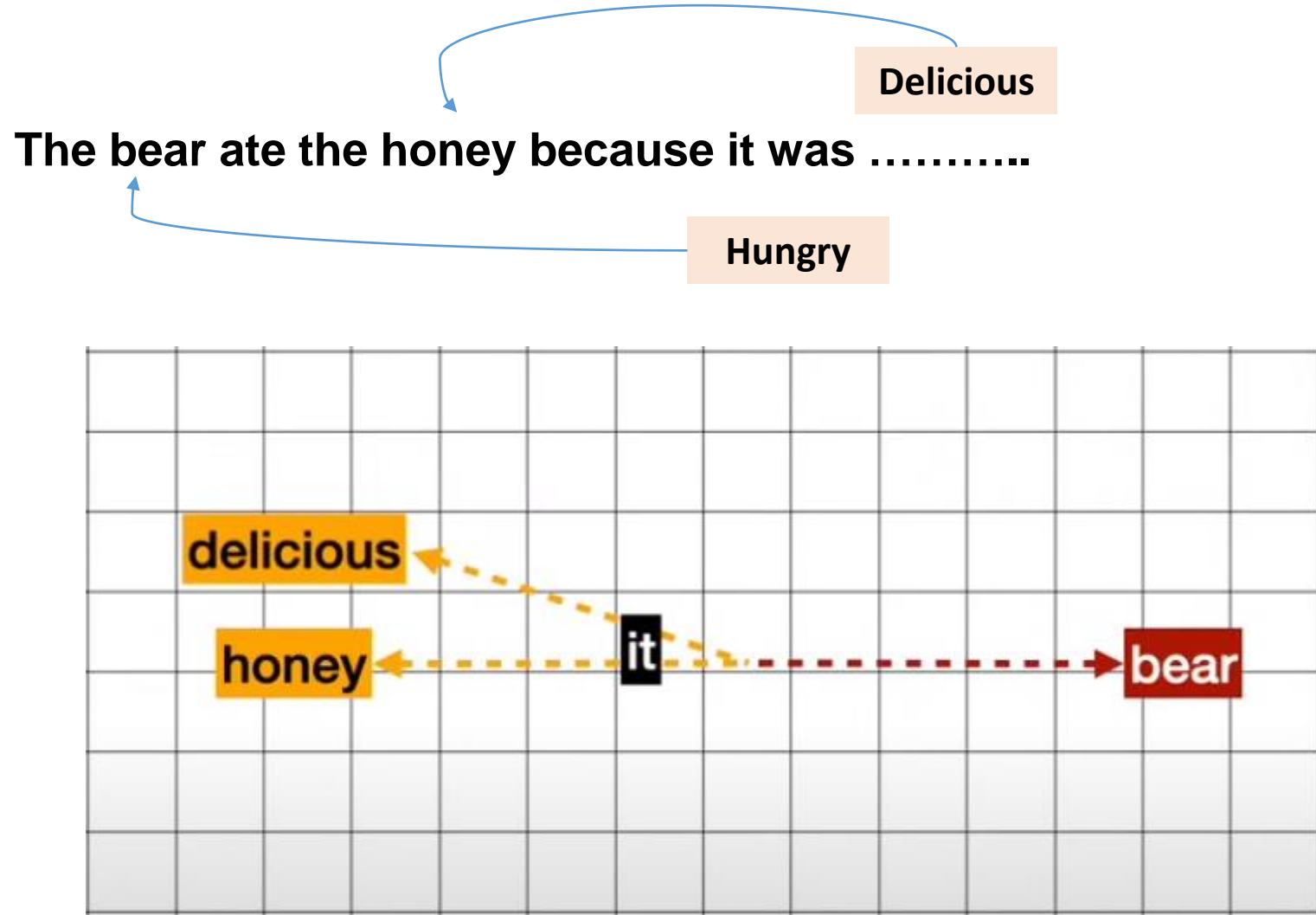
1. The sine and cosine functions have values in $[-1, 1]$, which keeps the values of the positional encoding matrix in a normalized range.
2. As the sinusoid for each position is different, you have a unique way of encoding each position.
3. You have a way of measuring or quantifying the similarity between different positions, hence enabling you to encode the relative positions of words.



Language Models

Transformer models

Capturing the context by make attention to specific parts

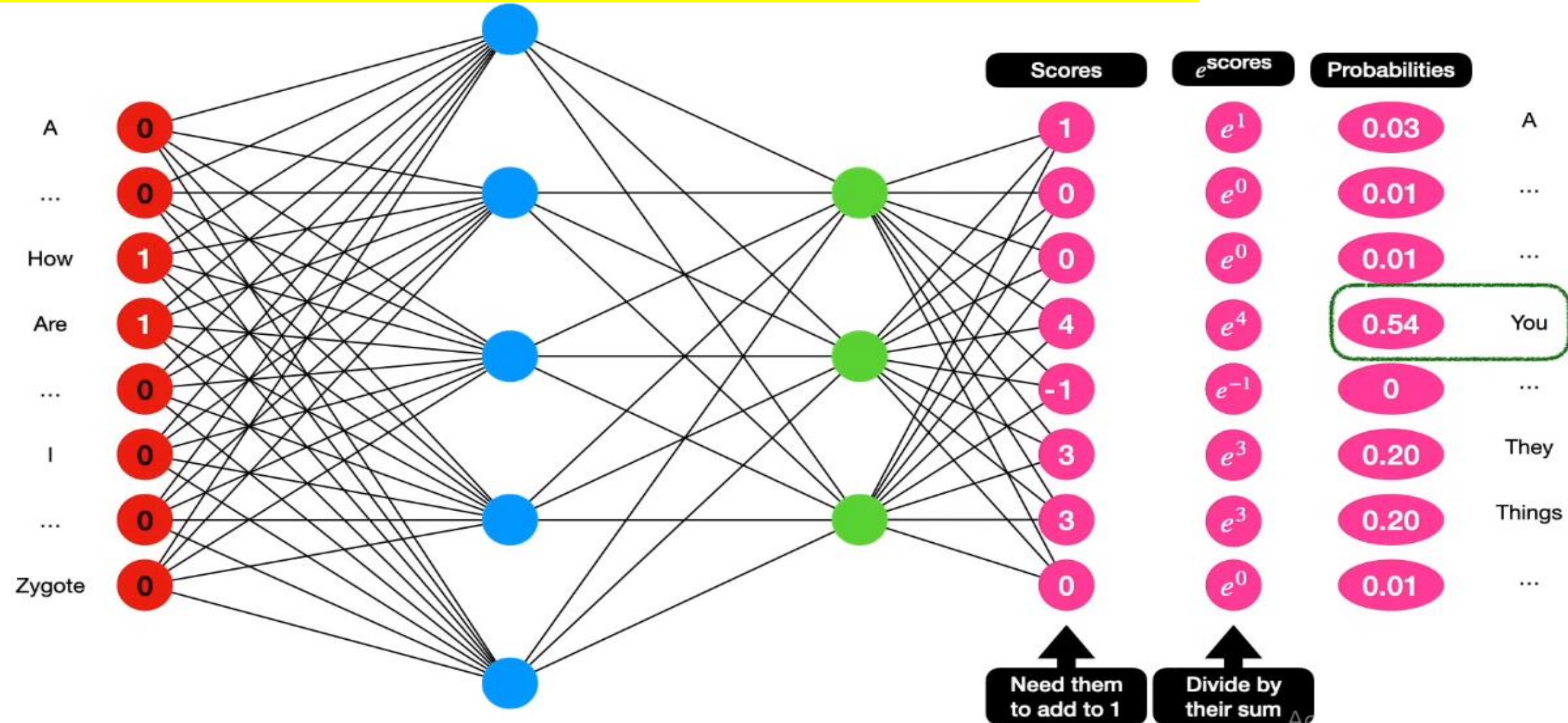


Language Models

Step 5: Softmax

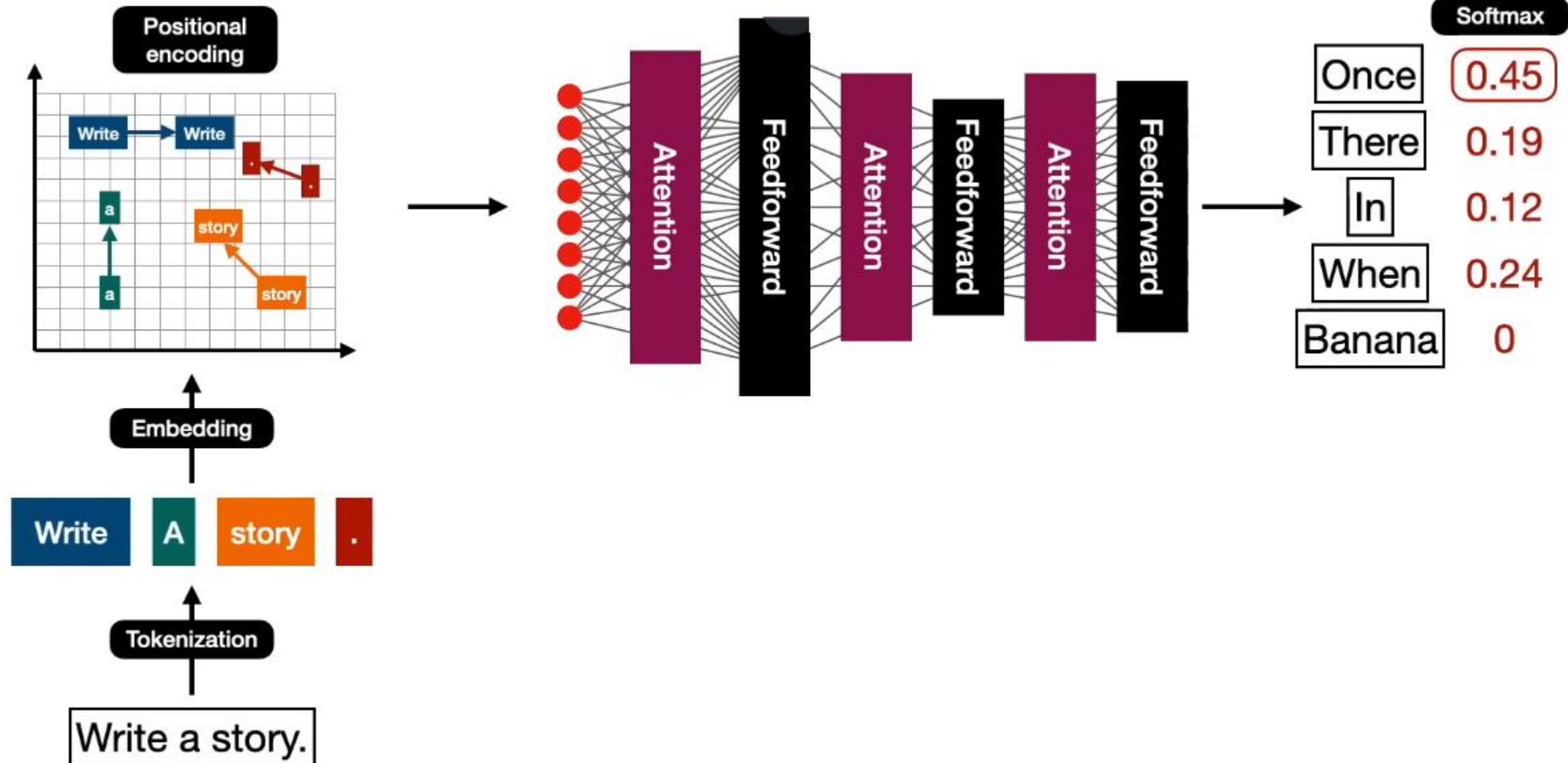
Transformer models

Classification based on Probability



Language Models

Transformer models



Language Models

Transformer models

The transformer models needs to be trained on large datasets to learn how to predict the next word in the sentence.

Now what if we want the transformer model to answer a question?

What is the capital of Nigeria? Abuja

Quiz

What is the capital of Nigeria?

What is the capital of Chad?

What is the capital of Lebanon?

Story

What is the capital of Nigeria? She asked.

Chat

What is the capital of Nigeria?

That is a good question

History

What is the capital of Nigeria?

Since 1991, it's Abuja, but before, it used to be Lagos

Language Models

Transformer models

How to enforce transformer model to answer questions?

Post training (Fine tuning)

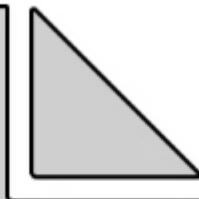
Training the transformer model on a dataset of questions and answers

What is the capital of Nigeria? | Abuja

Q/A

What is the capital of
Nigeria?

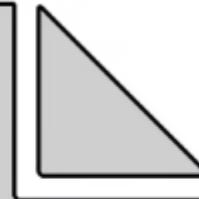
Abuja



Q/A

What is the capital of
Colombia??

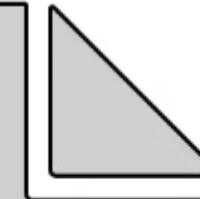
Bogotá



Q/A

Who discovered
algebra?

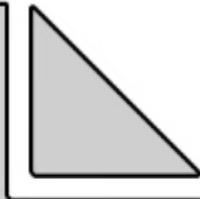
Al-Khwarizmi



Q/A

Who discovered
abstract algebra?

Emmy Noether



Language Models

Transformer models

How to enforce transformer model to chat with us?

Post training (Fine tuning)

Training the transformer model on a dataset of **commands and actions**

Hello, how are you?

Good, and you?

...

Chat

Hello, how are you?

I'm good, and you?

Great, thank you!

Chat

Good morning, how
can I help you?

Thank you, can you
connect me with...

Chat

Hi mom!

Hello dear!

Chat

Hello, please
connect me with
customer support.

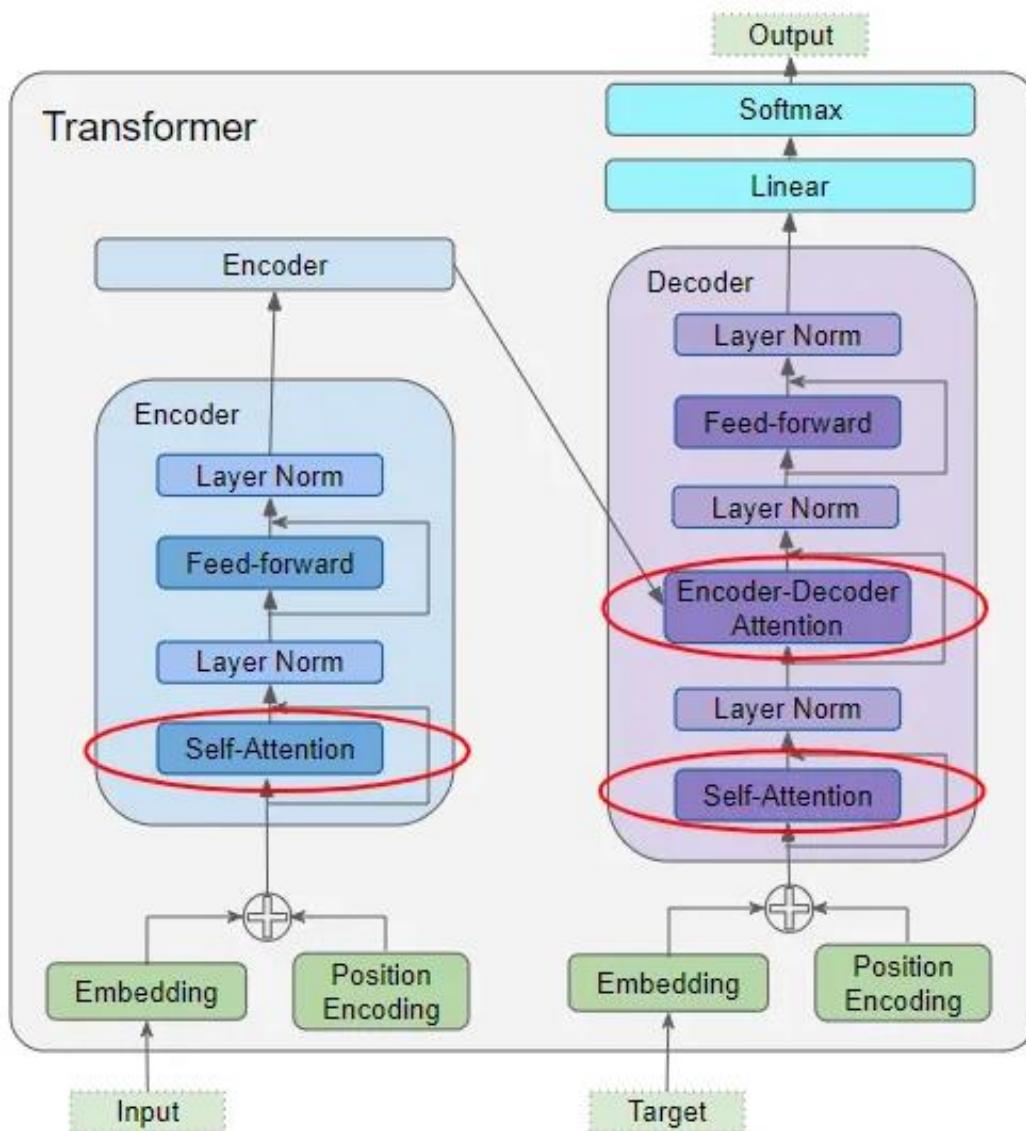
Of course!

Thank you!

Language Models

Transformer models

When to use both parts of transformer model (Encoder & Decoder)



In transformation applications:
Translation from French to English

In this case, we need an extra Encoder-Decoder attention layer
We also need to connect the output of the last encoder part with
this Encoder-Decoder attention layer of the first decoder part.