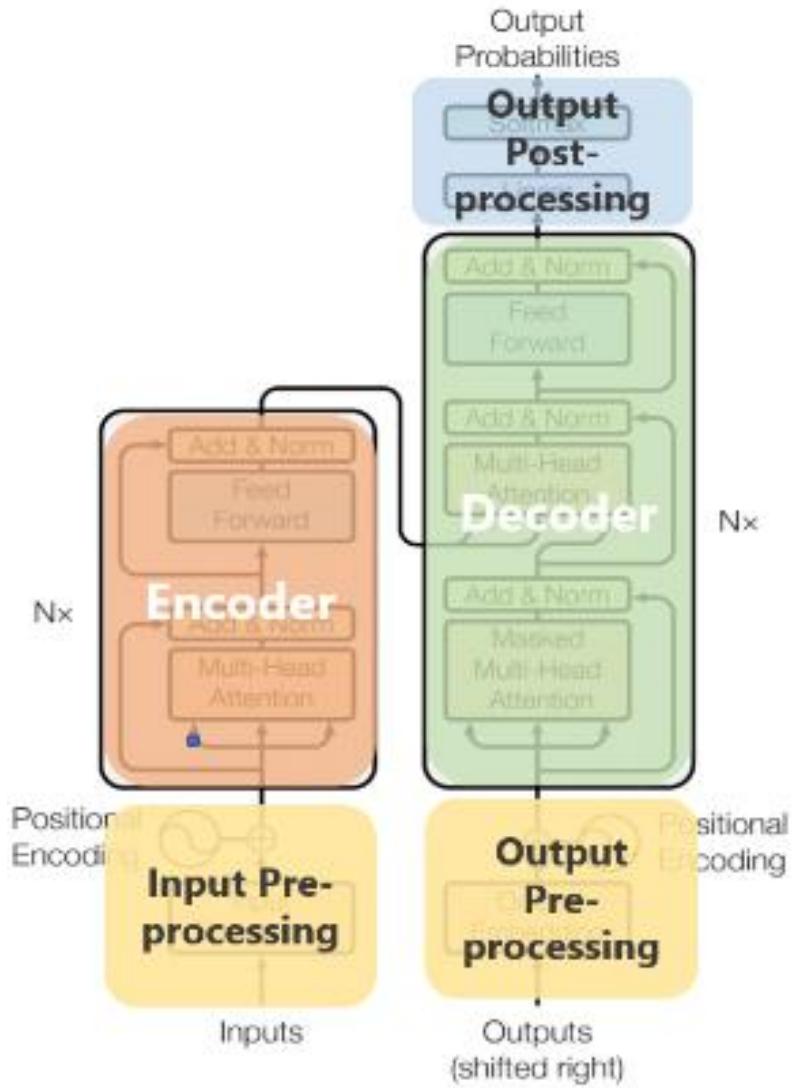
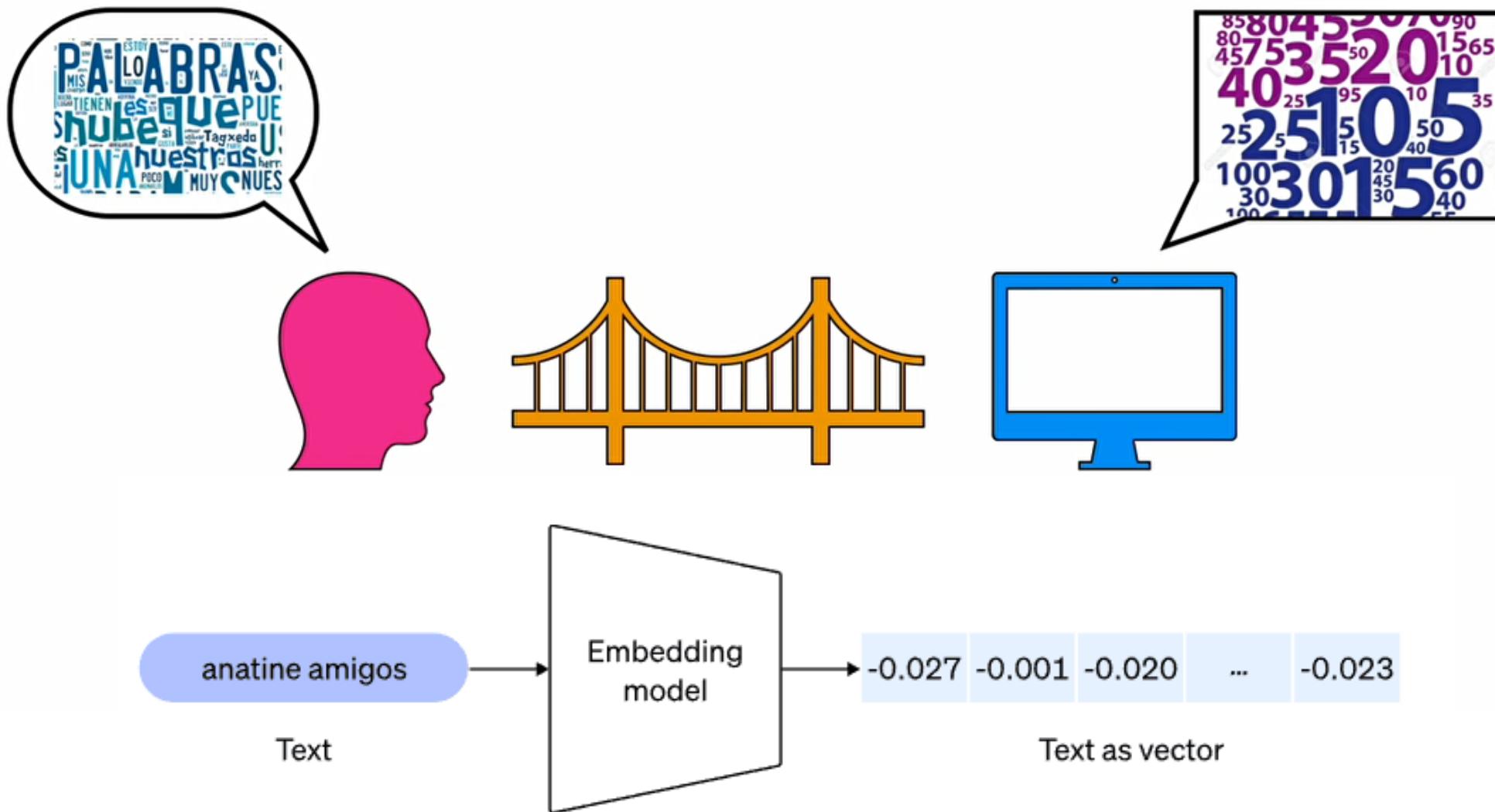


INTRO TO TRANSFORMERS



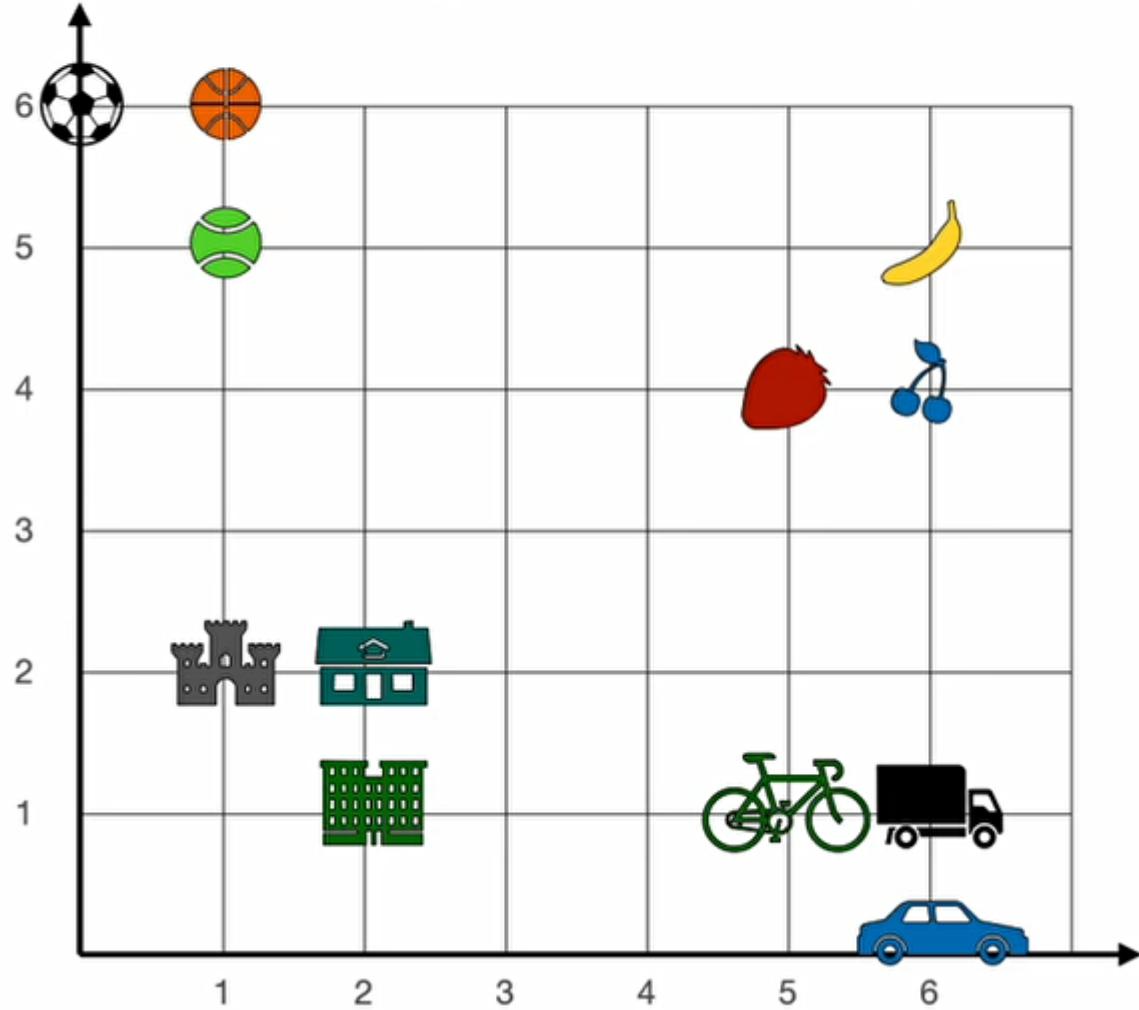
Embeddings

Embeddings



Embeddings Quiz 1

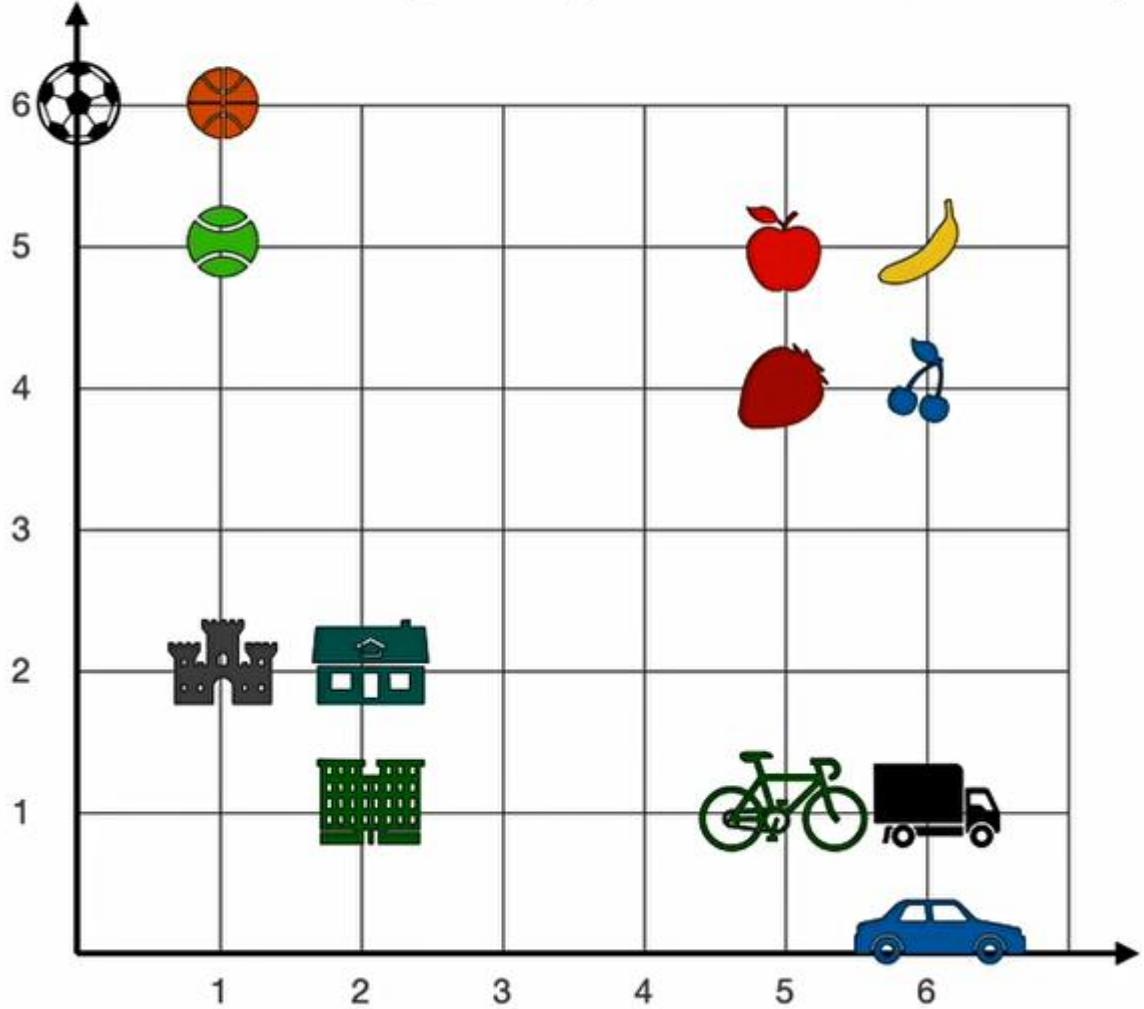
Where would you put the word apple?



Word	Numbers	
Apple	?	?
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

Embeddings Quiz 1

Where would you put the word apple?

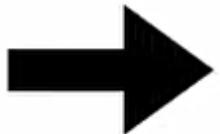


Word	Numbers	
Apple	5	5
Banana	6	5
Strawberry	5	4
Cherry	6	4
Soccer	0	6
Basketball	1	6
Tennis	1	5
Castle	1	2
House	2	2
Building	2	1
Bicycle	5	1
Truck	6	1
Car	6	0

Word embeddings

Many more columns

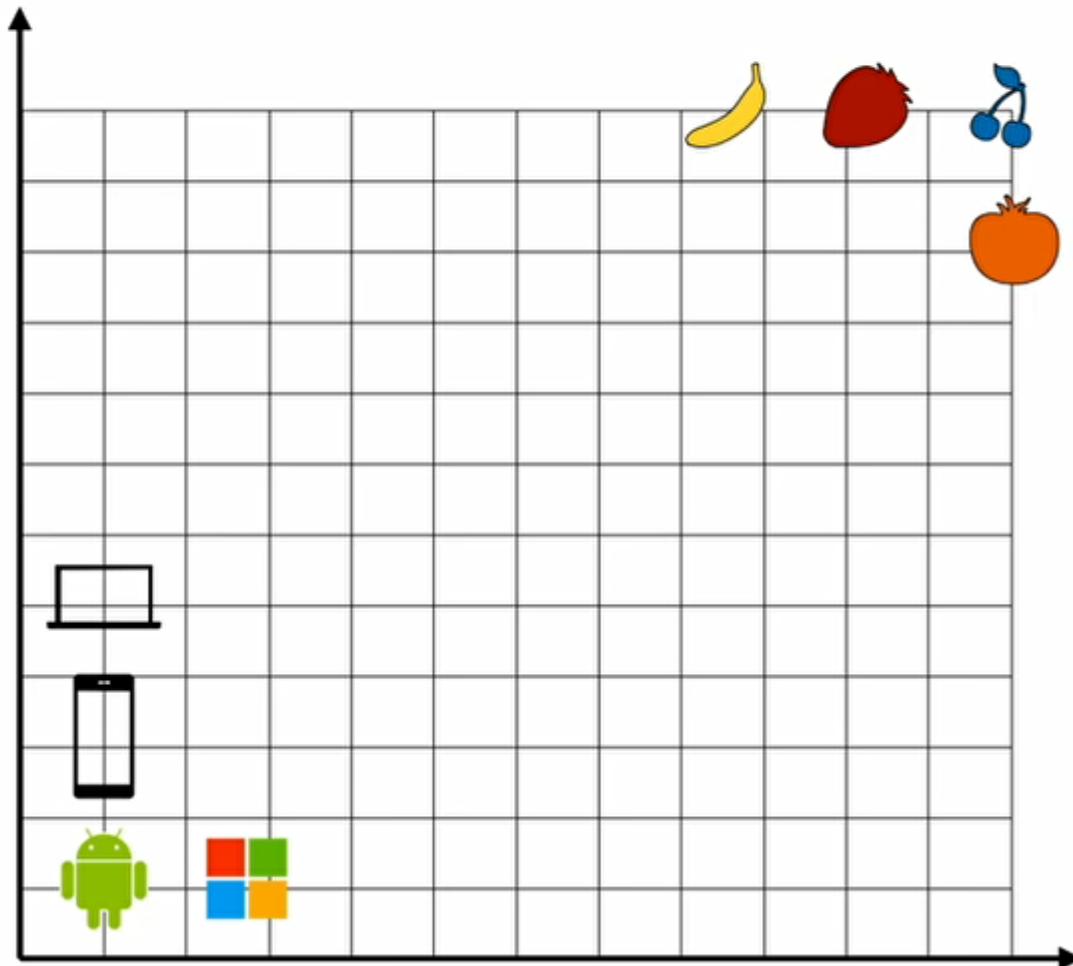
Word	Numbers	
Apple	5	5
Soccer	0	6
House	2	2
Car	6	0



Word	Numbers			
A	-0.82	-0.32	...	0.23
Aardvark	0.419	1.28	...	-0.06
...			...	
Zygote	-0.74	-1.02	...	1.35

4096

Embeddings Quiz 2



Top right or bottom left?

Cherry 

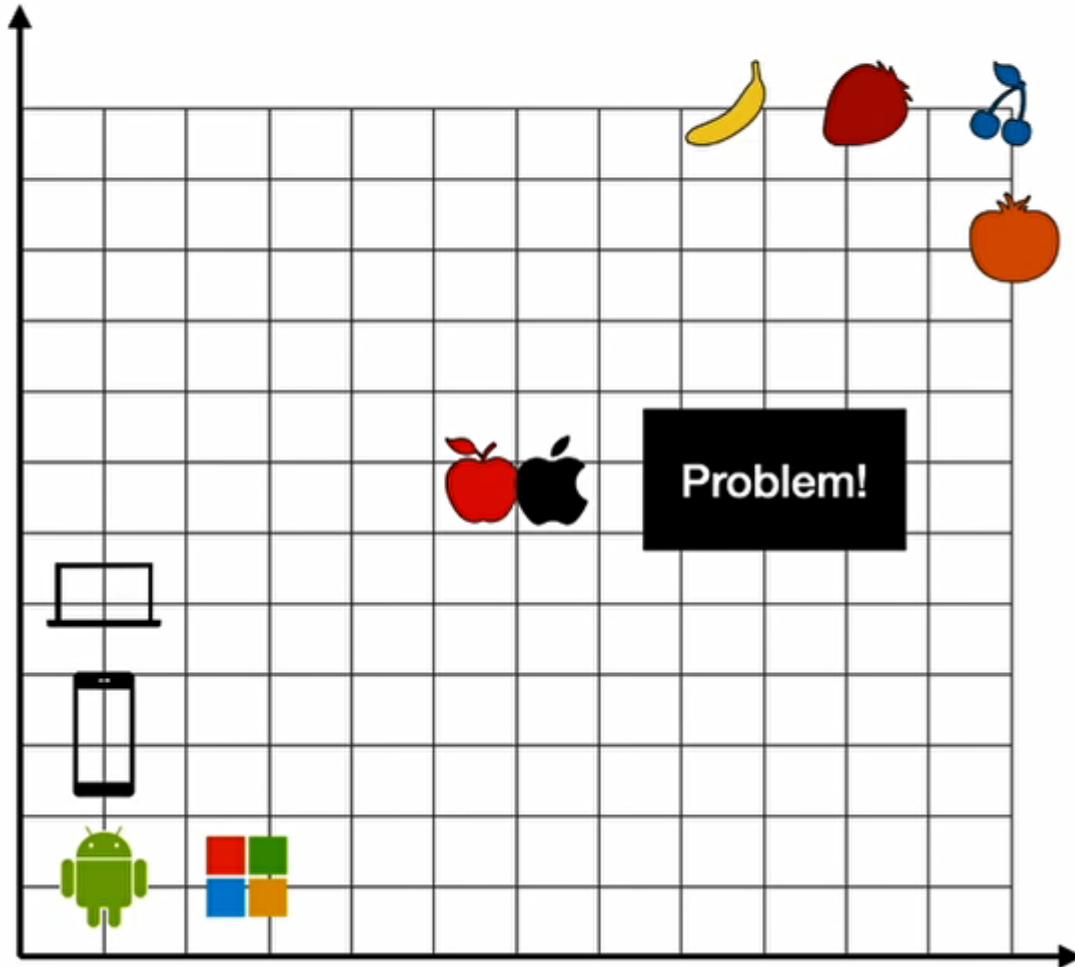
Android 

Laptop 

Banana 

Apple?

Embeddings Quiz 2



Top right or bottom left?

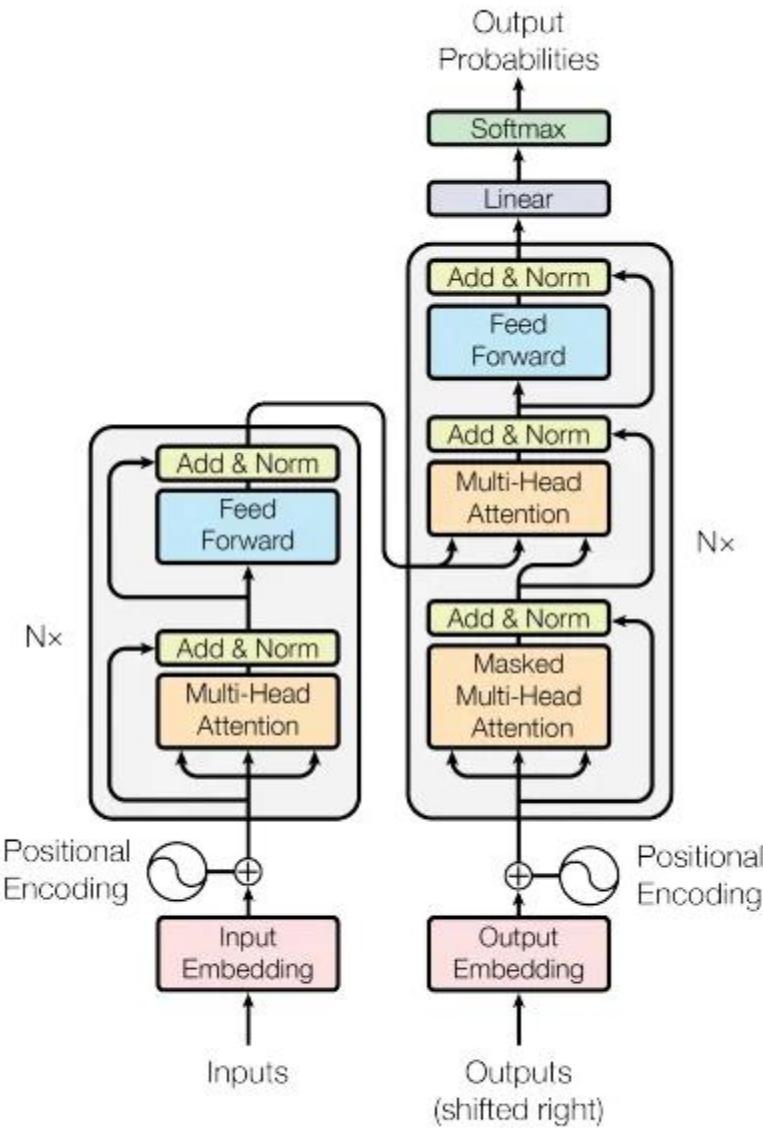
Cherry

Android

Laptop

Banana

Apple?

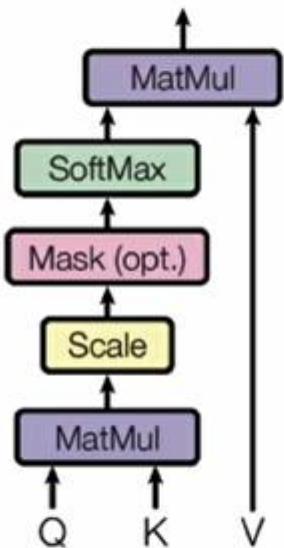


WE CAN SOLVE THIS TYPE OF PROBLEM WITH ATTENTION, SELF ATTENTION, IT USE THE CONTENT OF THE SENTENCE IN ORDER TO HELP THE EMBEDDINGS TO RESOLVE THIS TYPE OF AMBIGUOUS

Self- Attention

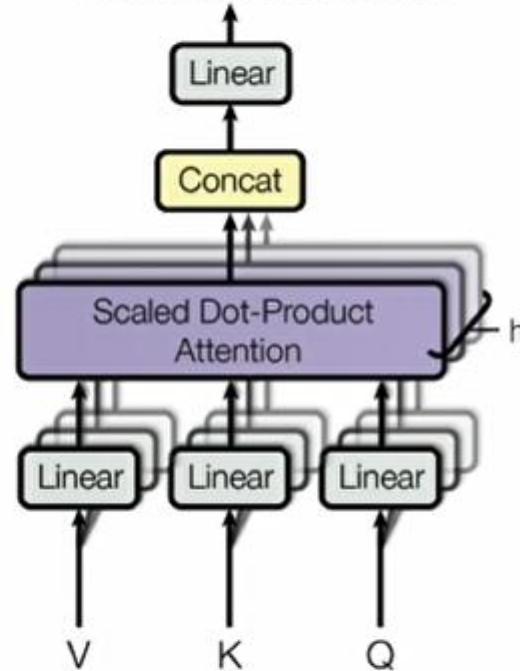
Attention

Scaled Dot-Product Attention



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

Multi-Head Attention



$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \\ \text{where } \text{head}_i &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned}$$

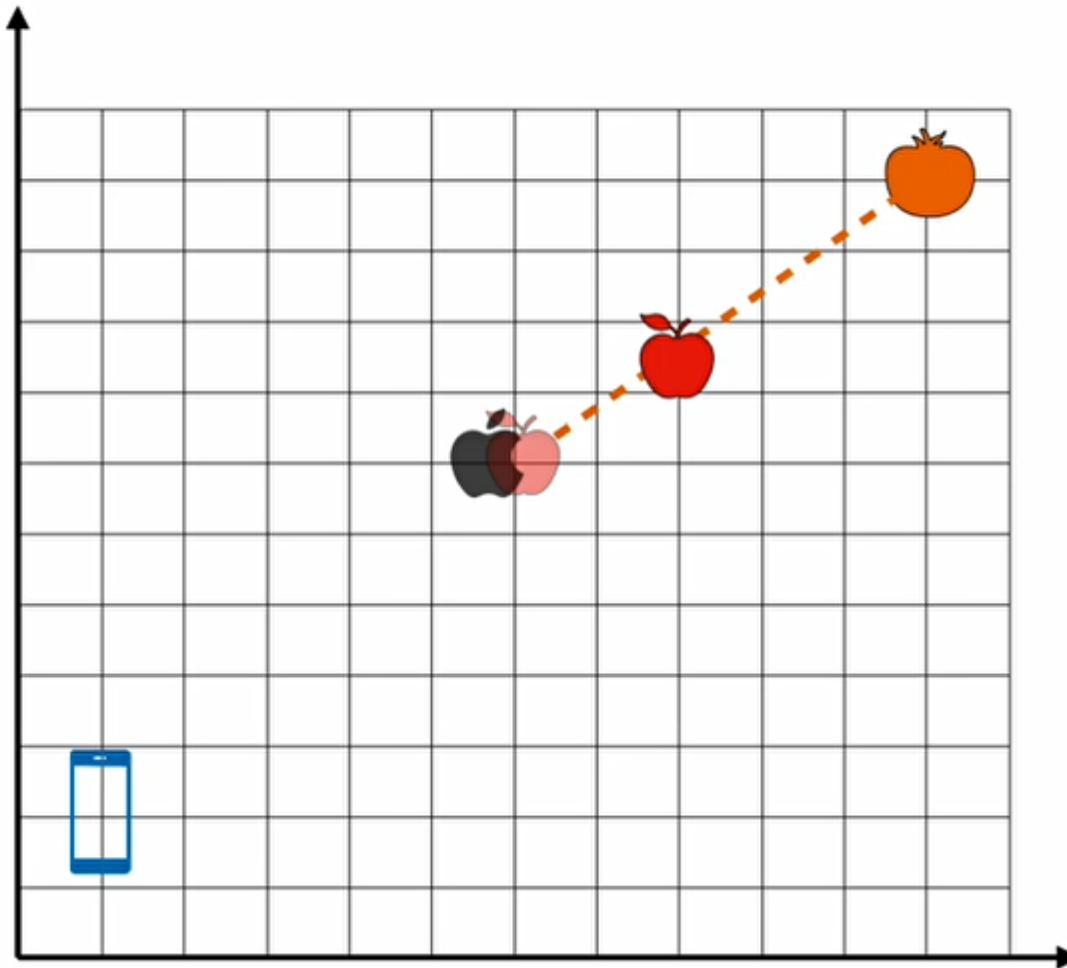


WE NEED A CONTEXT, WE NEED TO USE OTHER WORDS IN A SENTENCE TO
HELP THE MAIN UNDERSTANDING

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

apple unveiled the new **phone**

Attention

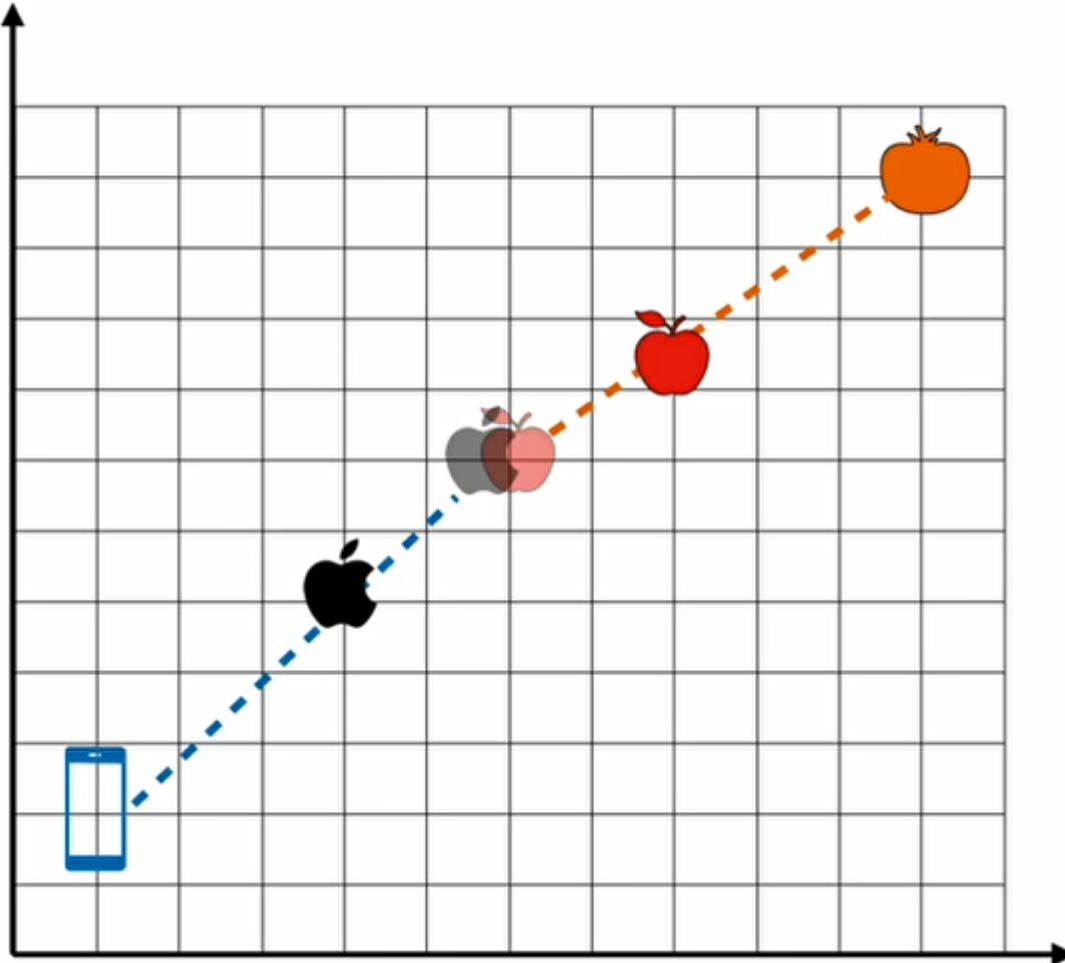


please buy an **apple** and an **orange**

apple unveiled the new phone

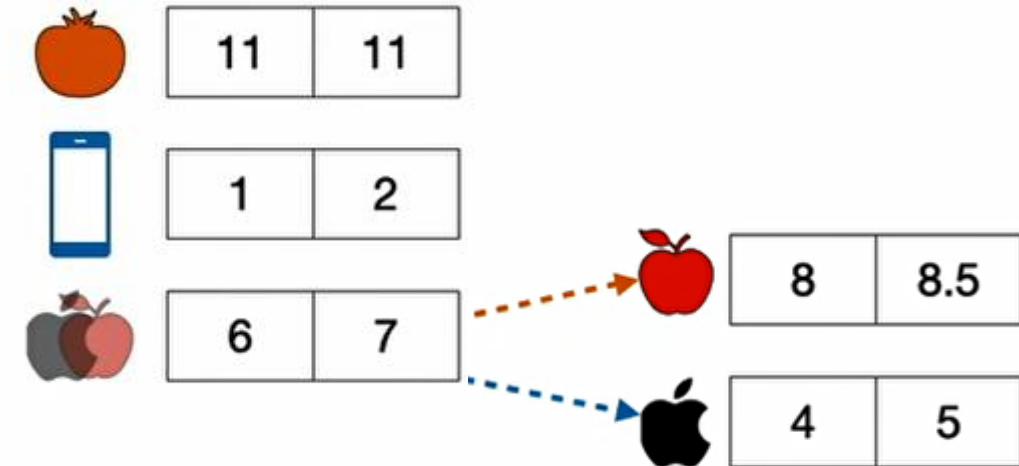
Attention

OF COURSE THIS IS REPEATED SO MANY MANY TIMES !!!

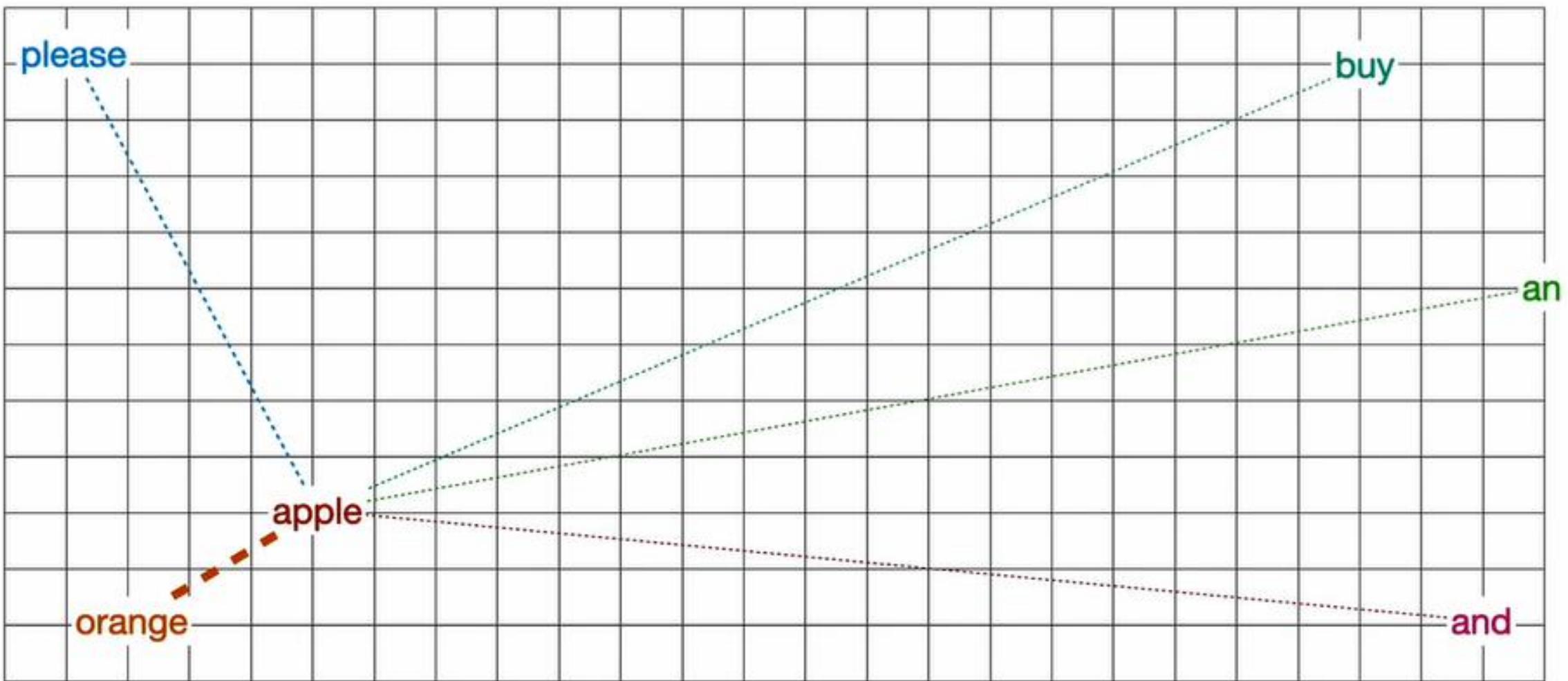


please buy an **apple** and an **orange**

apple unveiled the new **phone**



What about the other words?



please buy an apple and an orange

You apply attention to all the words

please

buy

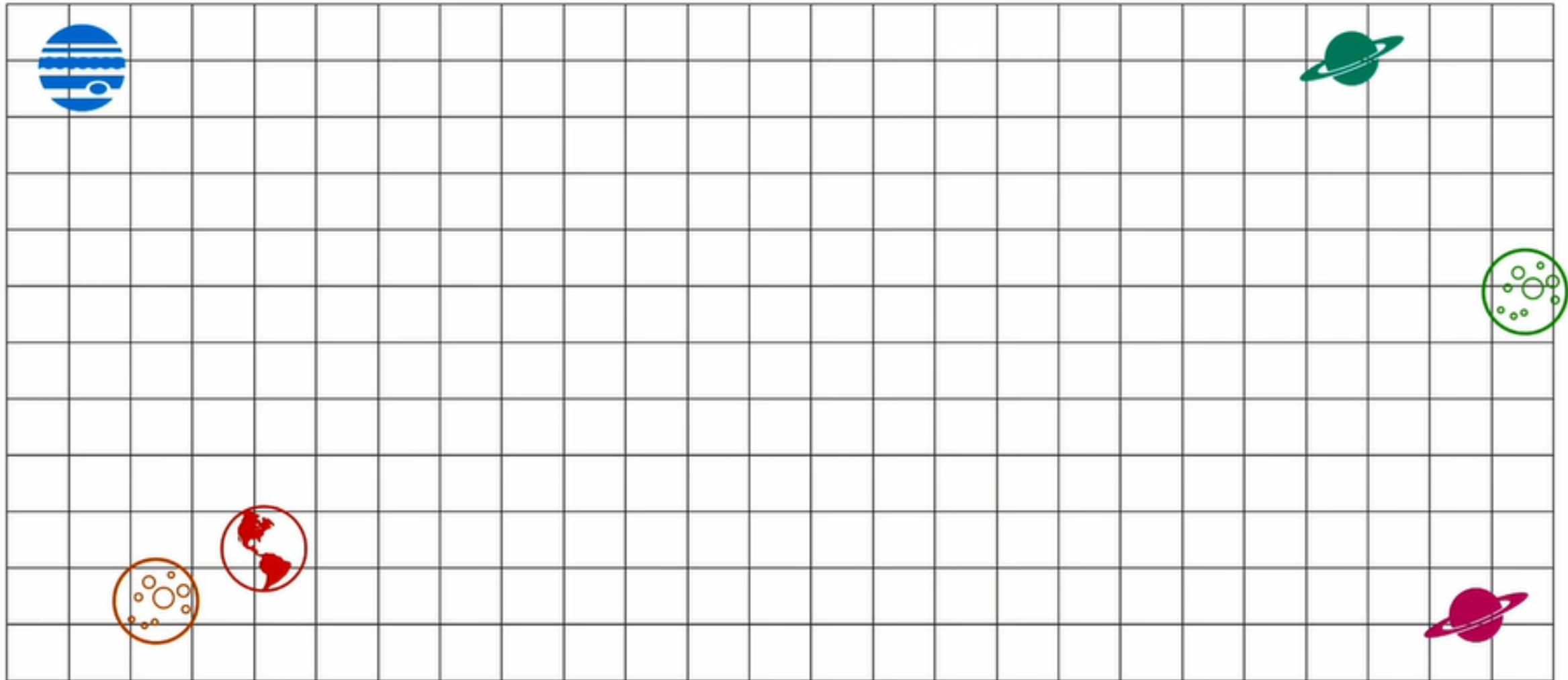
an

apple

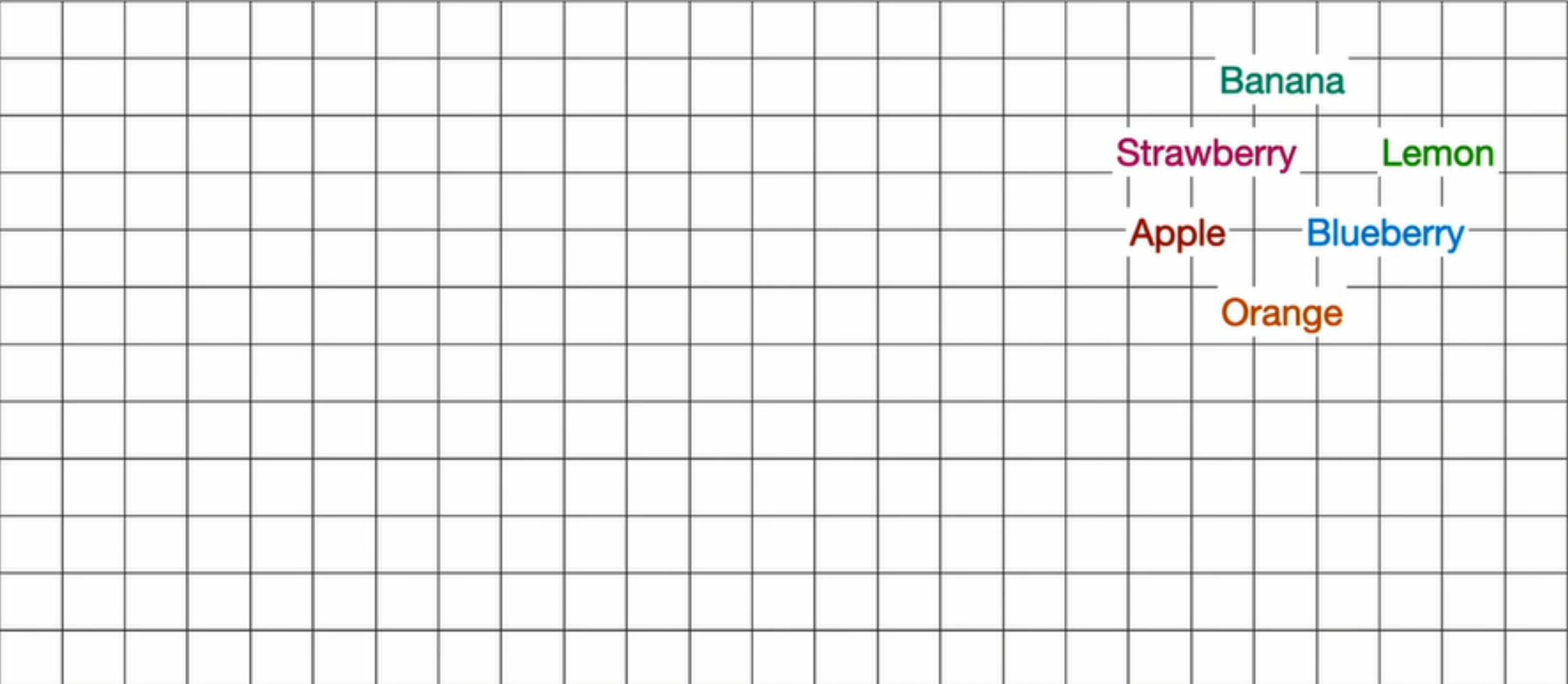
orange

and

It's kind of like gravity...



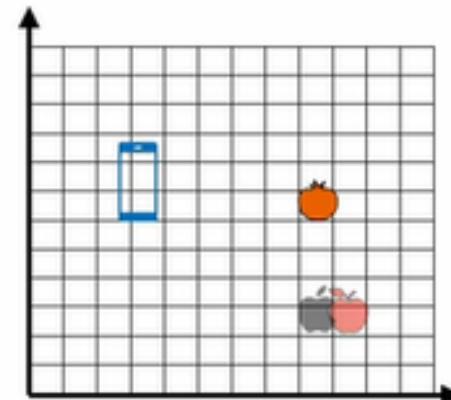
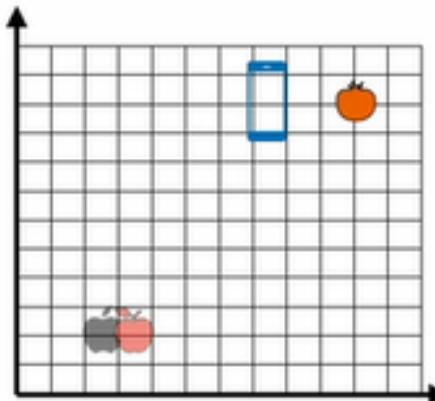
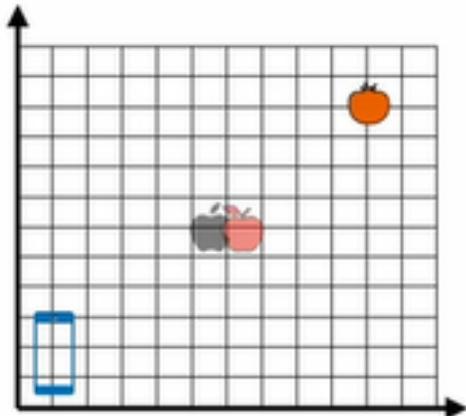
Context pulls



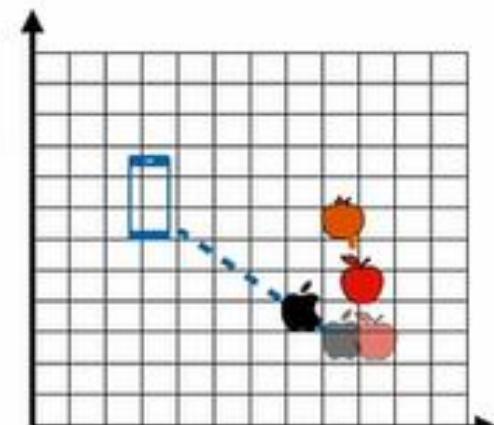
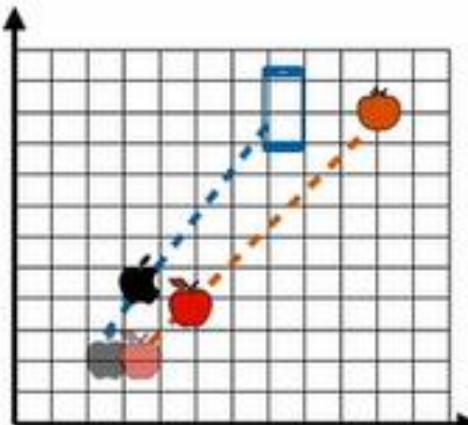
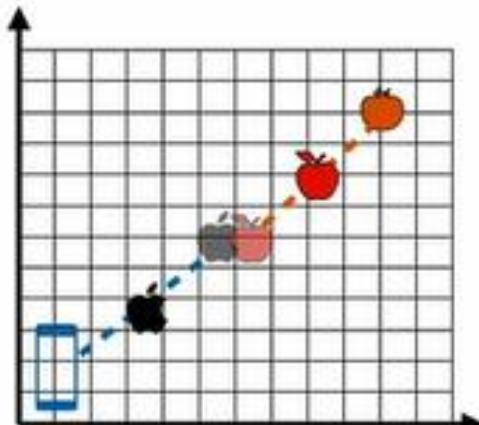
Banana
Strawberry Lemon
Apple Blueberry
Orange

Multi-head attention

Ideally, we'd like to have lots of embeddings



Ideally, we'd like to have lots of embeddings

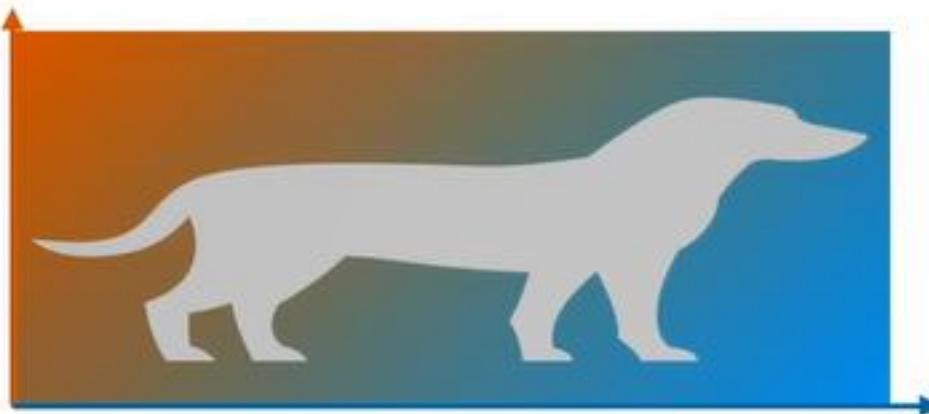


Solution: We'll build embeddings by modifying existing embeddings

Linear transformations



Rotate



Stretch



Stretch

Linear transformation preserves the vector structure of its arguments, as it conserves vector addition and scalar multiplication.

For example, a rotation, a stretching, a shear, and a projection are all examples of linear transformations if they satisfy the properties described above

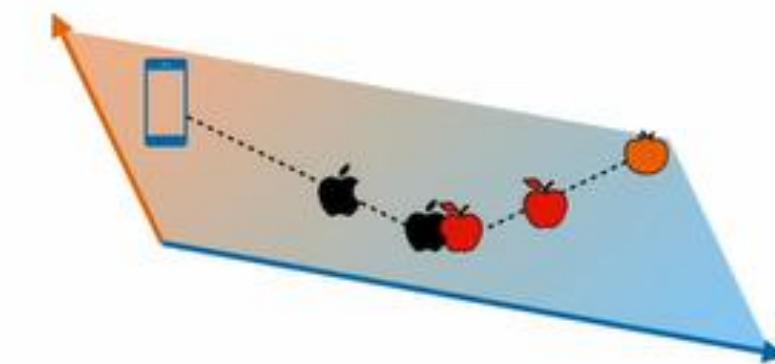
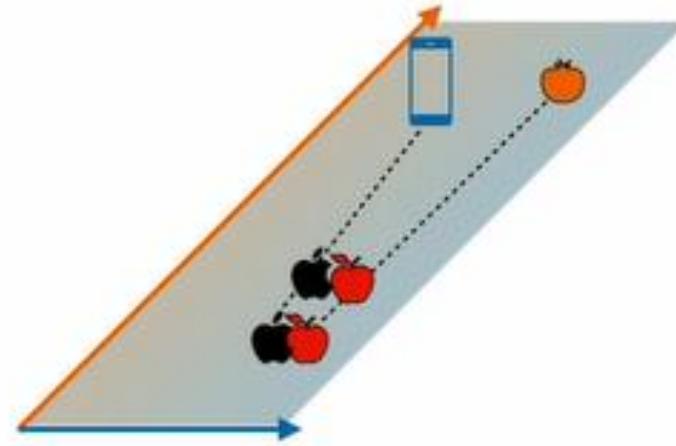
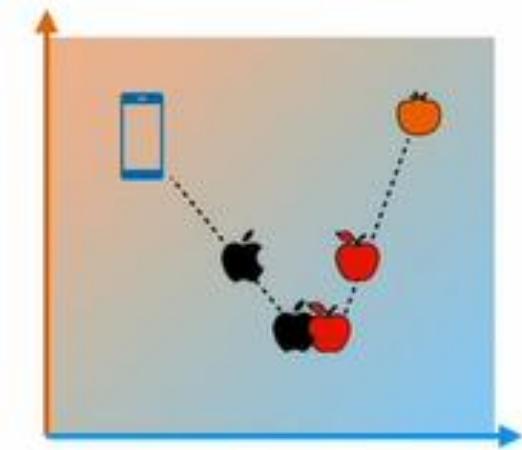
Linear transformations



Combination
----->

A dashed arrow pointing from the original dog image to the transformed one, labeled "Combination".

Get new embeddings from existing ones

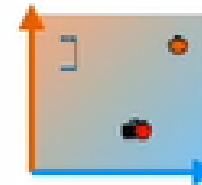
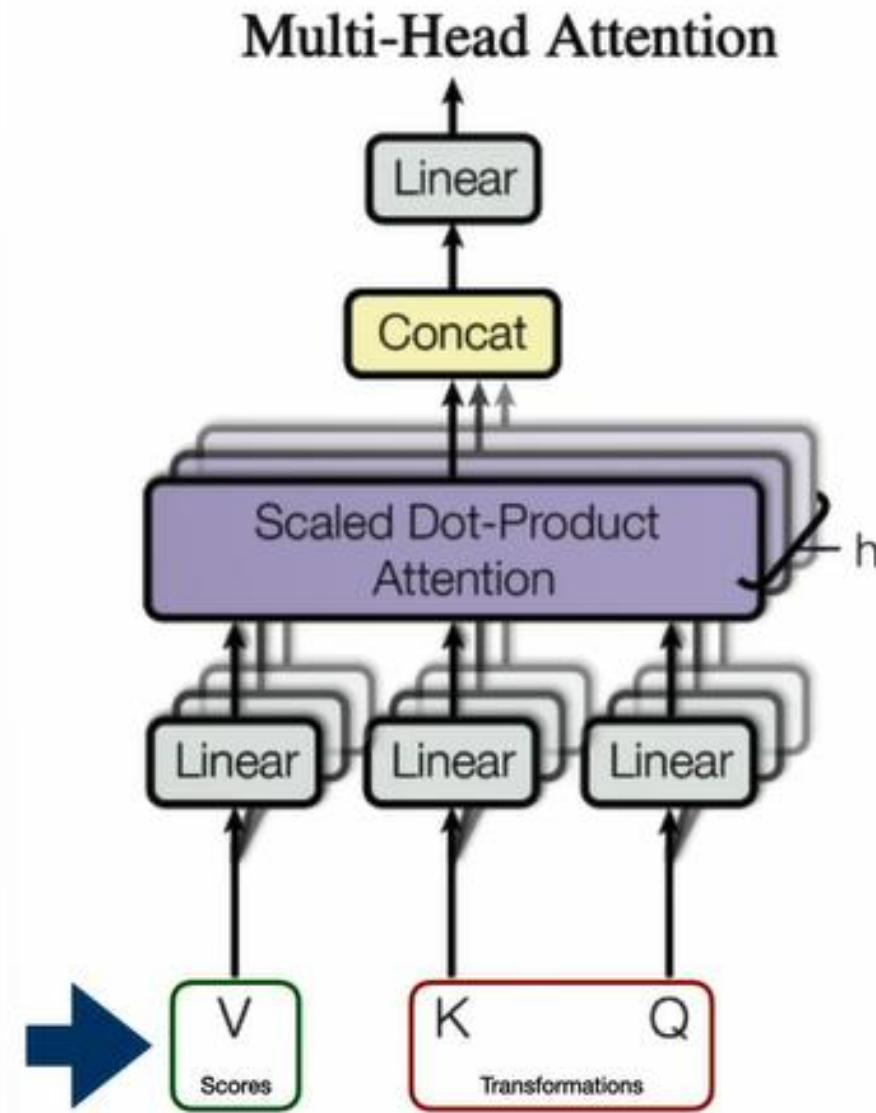


Score: 1

Score: 0.1

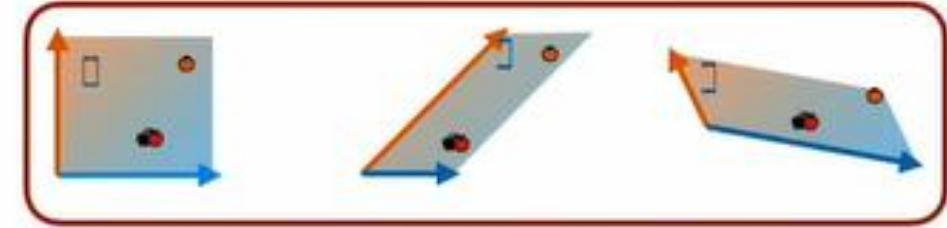
Score: 4

WITH THE TRAIN OF NEAURAL NETWORK, WE TRAIN THE EMBEDDINGS SCORE, SO HERE THE “Q” AND “K” MATRIX COME IN:

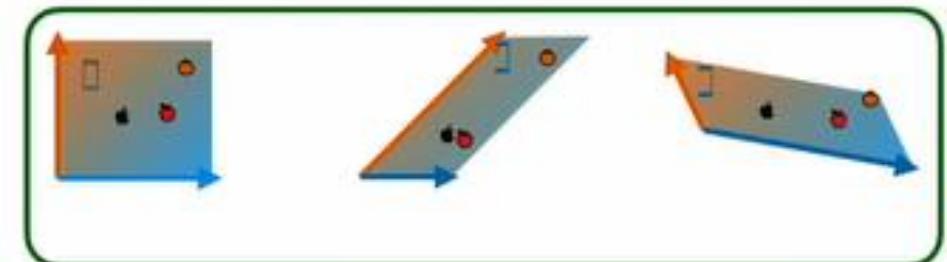


please buy an **apple** and an orange
apple unveiled the new phone

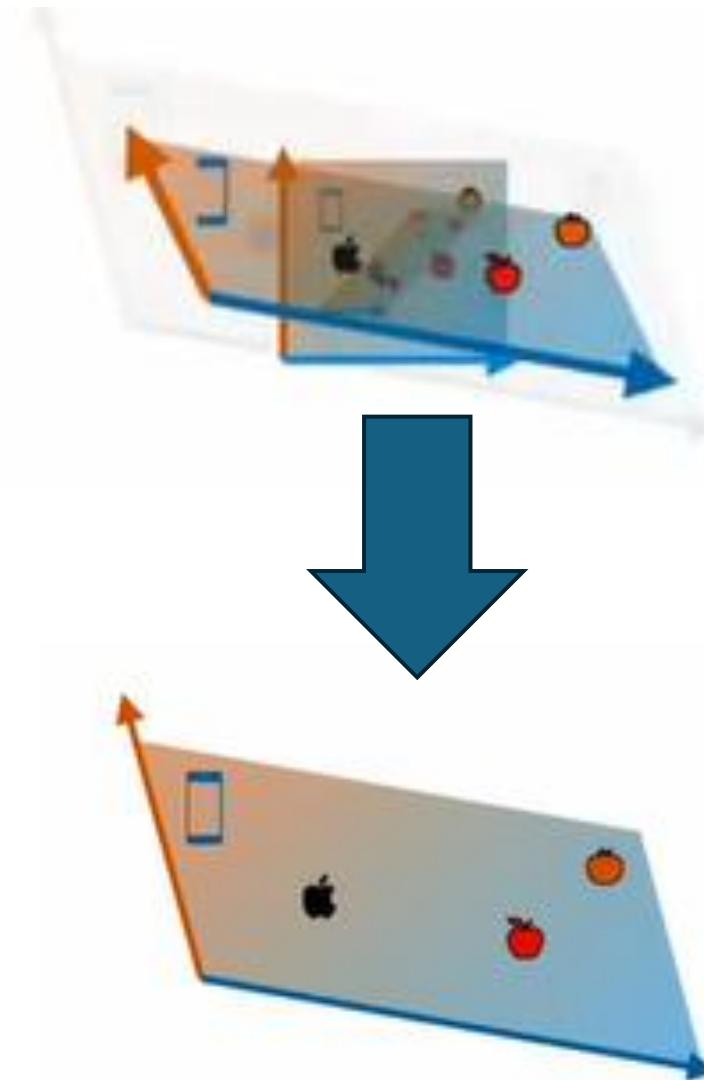
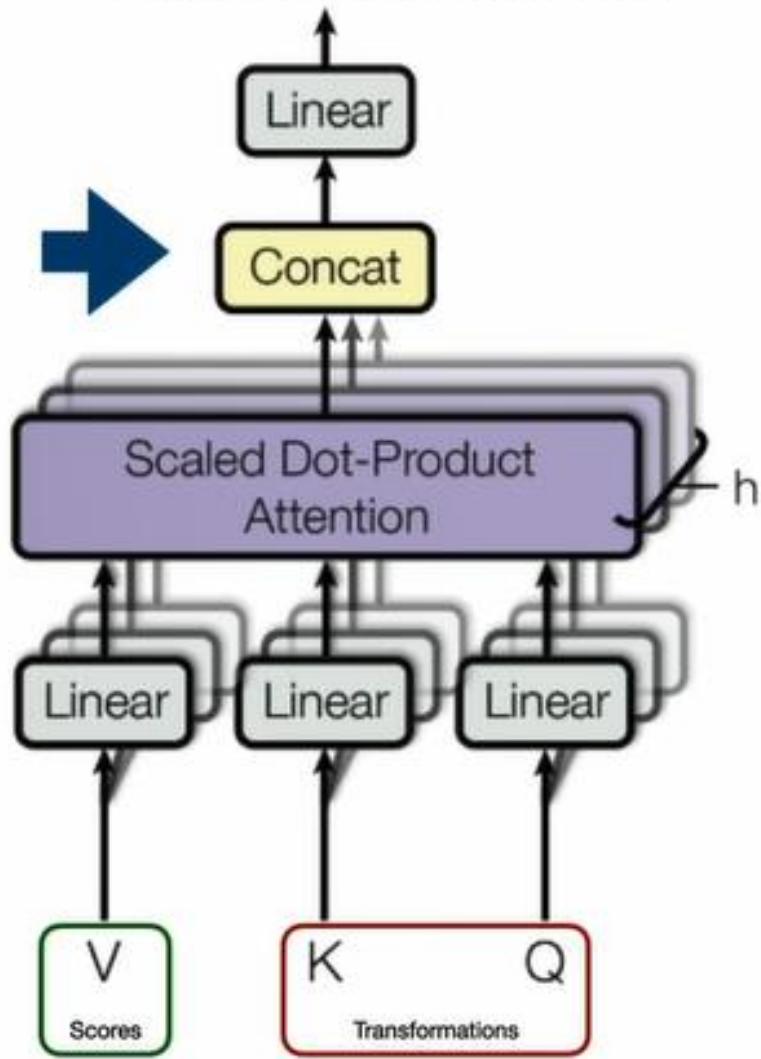
K and Q matrix help the transformation to enhance the similarity



V Matrix help to visualize get better Embeddings

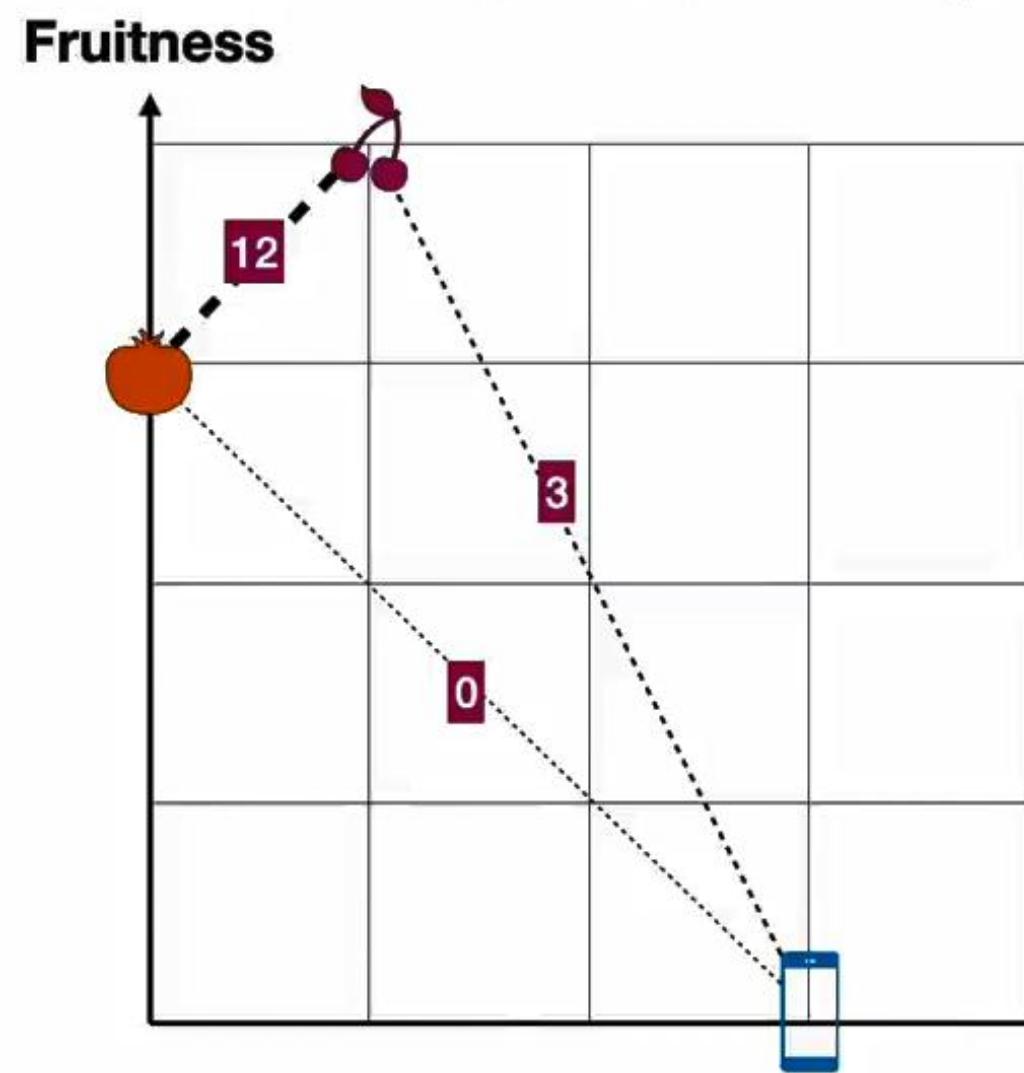


Multi-Head Attention



Similarity

Measure 1: Dot product



Sim	Tech	Fruitness	Sim	Tech	Fruitness
Cherries	1	4	Tomato	0	3
Tomato	1	4	Smartphone	3	0
Smartphone	0	3	Cherries	3	0

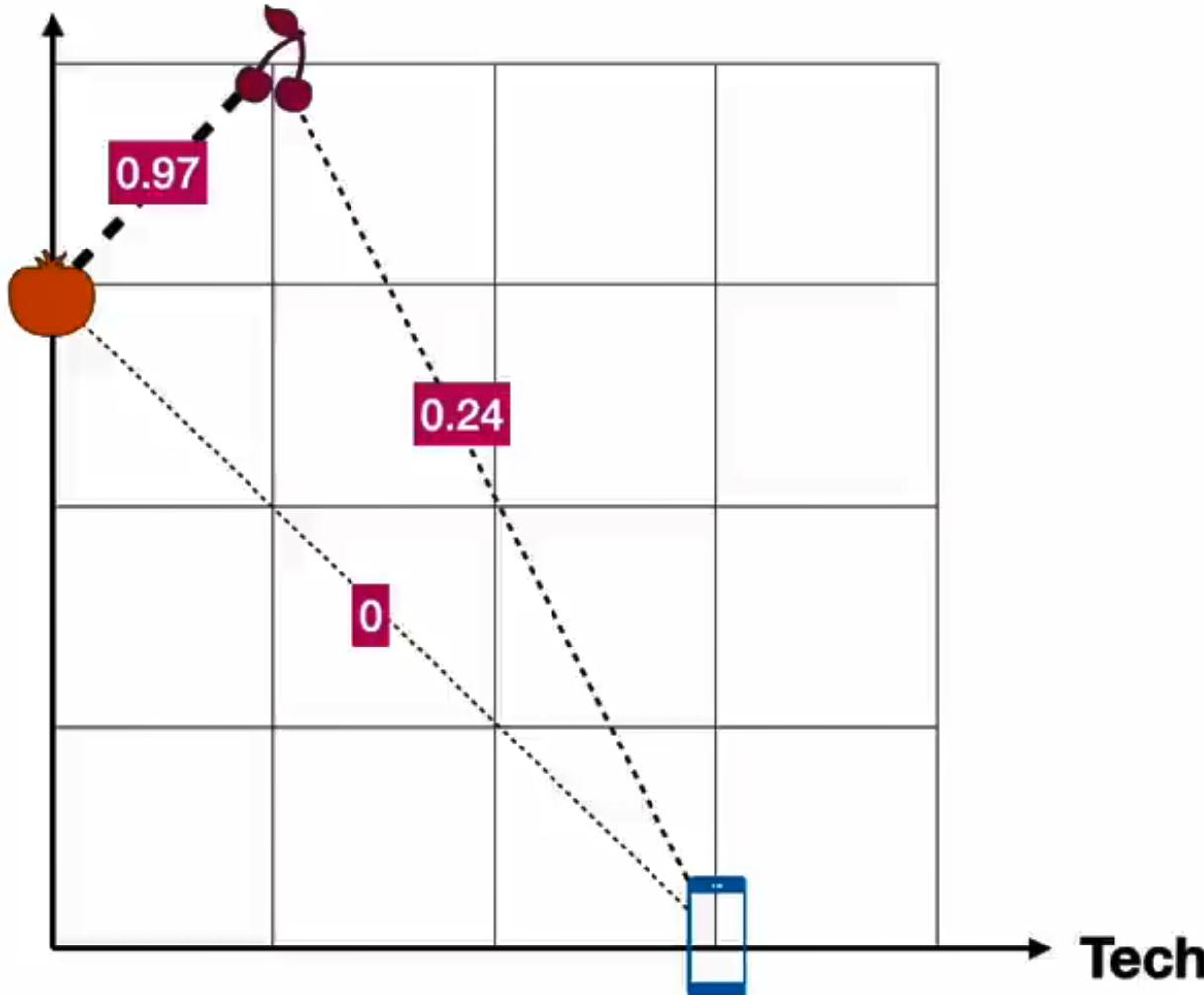
$1 \cdot 0 + 4 \cdot 3 = 12$

$1 \cdot 3 + 4 \cdot 0 = 3$

$0 \cdot 3 + 3 \cdot 0 = 0$

Measure 2: Cosine similarity

Fruitness



Sim

$$\cos(14^\circ) = 0.97$$



Sim

$$\cos(76^\circ) = 0.24$$

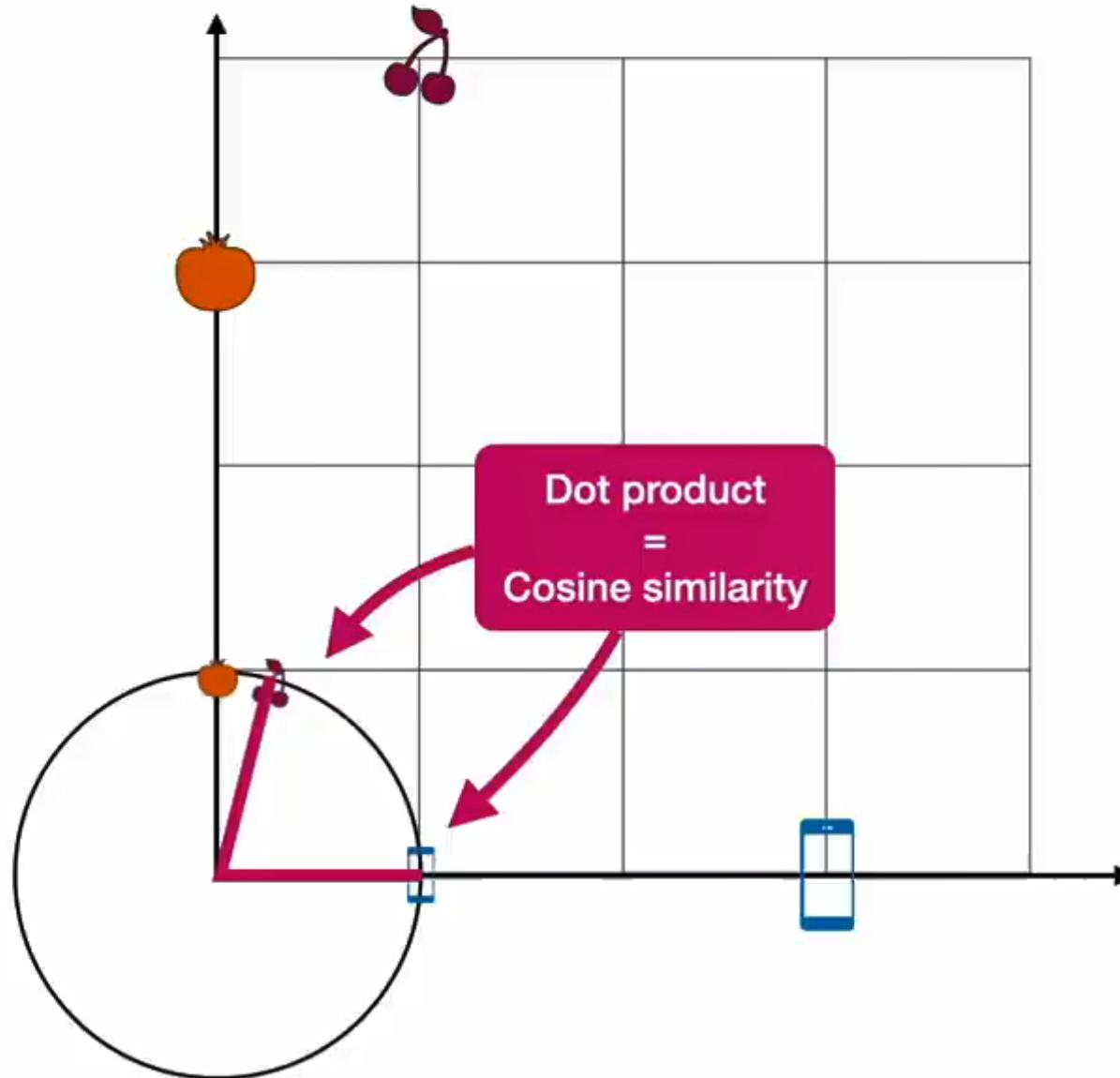


Sim

$$\cos(90^\circ) = 0$$

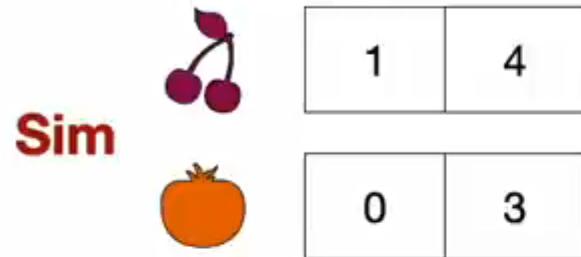


Dot product and cosine similarity

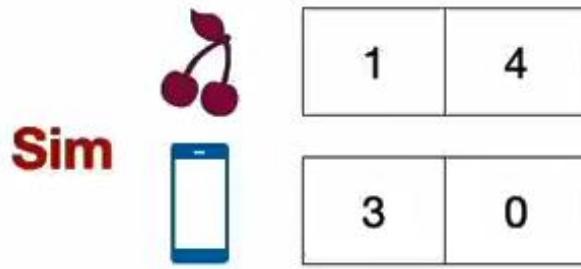


Measure 3: Scaled dot product

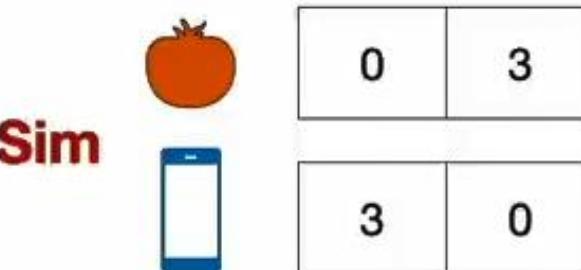
Dot product divided by the square root of the length of the vector



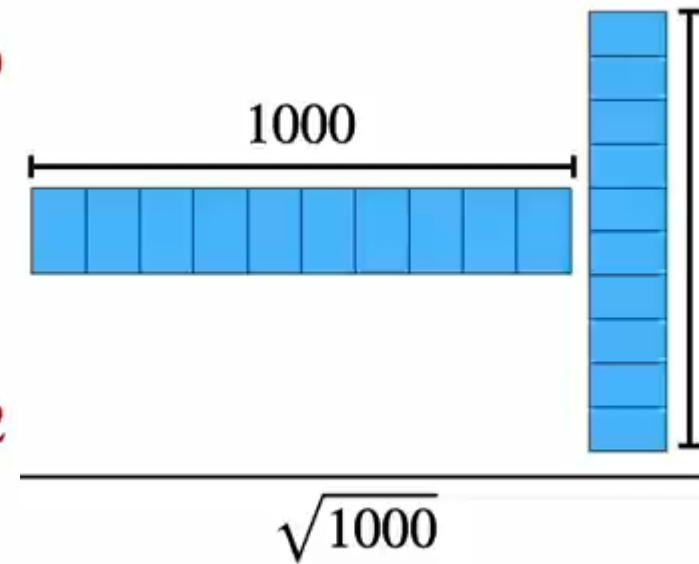
$$1 \cdot 0 + 4 \cdot 3 = 12 \longrightarrow \frac{12}{\sqrt{2}} = 8.49$$



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



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

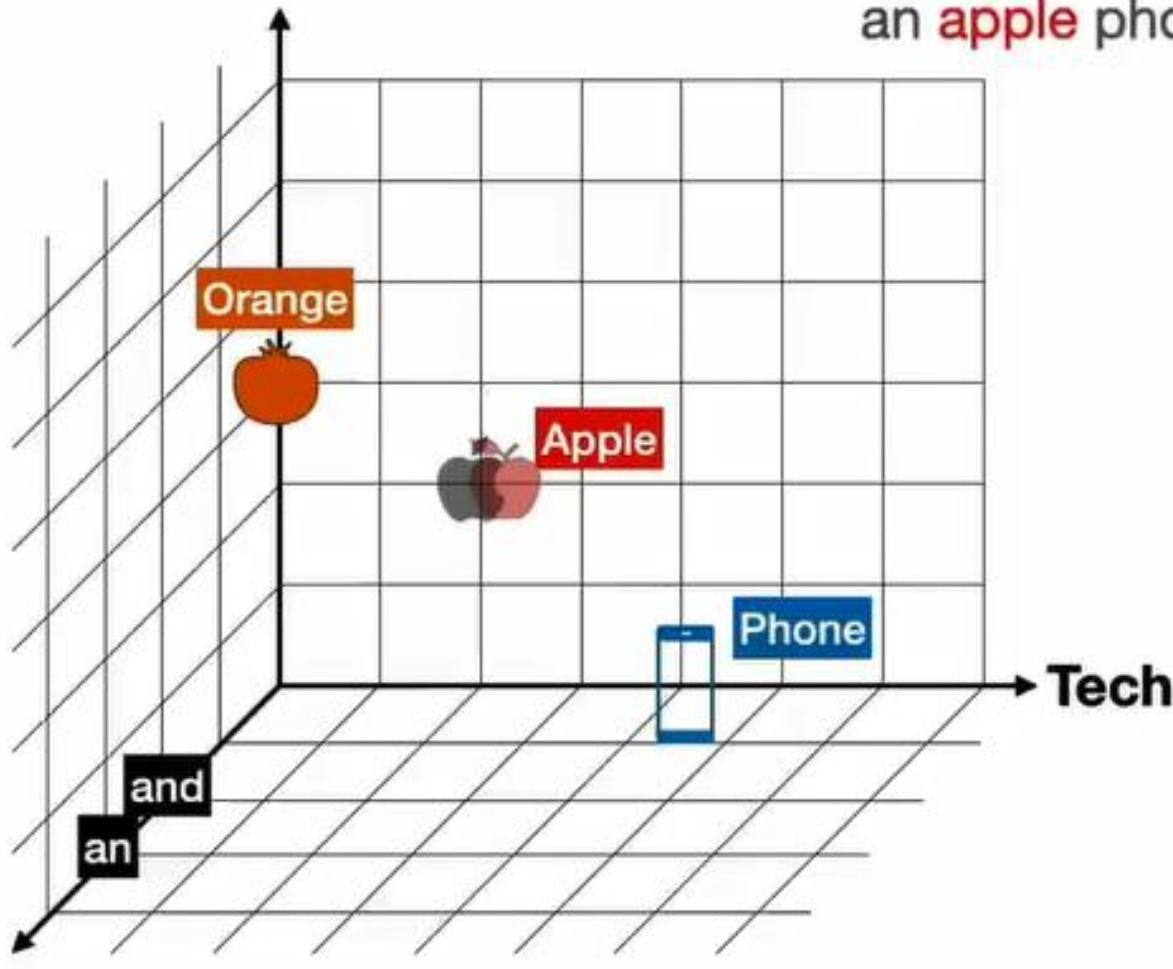


NOW, LET'S GET TO ATTENTION!

Cosine similarity

Fruitness

an apple and an orange
an apple phone



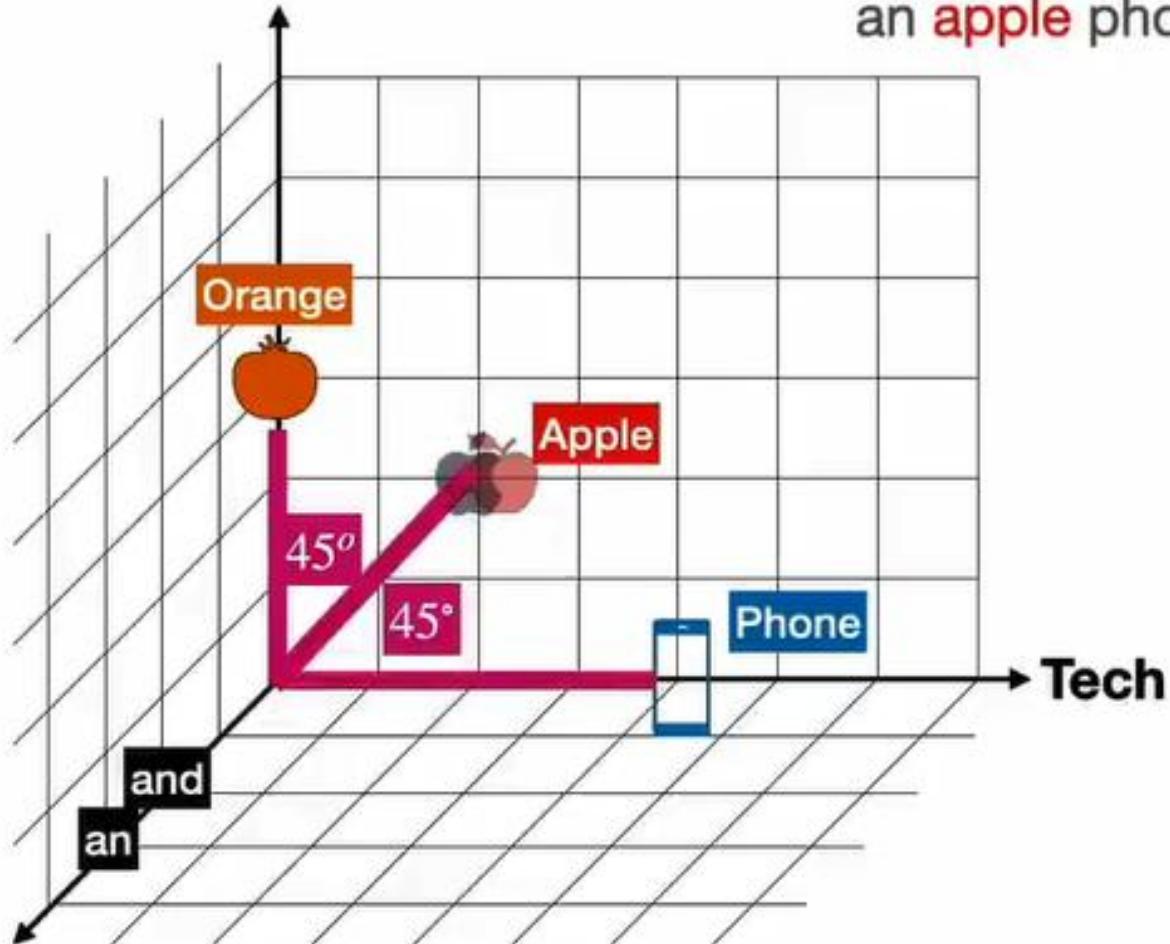
	Tech	Fruitness	Other
Orange	0	3	0
Phone	4	0	0
Apple	2	2	0
And	0	0	2
An	0	0	3

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

Other

Cosine similarity

Fruitness



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

Other

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

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

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

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

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

Word math

an apple phone

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

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

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

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

Normalization

Want coefficients to add to 1

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

Solution?

$$x \longrightarrow e^x$$

Softmax

$$x \longrightarrow e^x$$

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



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

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

$$\begin{aligned}\text{Orange} &\rightarrow 0.57 \text{ Orange} + 0.43 \text{ Apple} \\ \text{Apple} &\rightarrow 0.43 \text{ Orange} + 0.57 \text{ Apple} \\ \text{And} &\rightarrow 0.5 \text{ And} + 0.5 \text{ An} \\ \text{An} &\rightarrow 0.5 \text{ An} + 0.5 \text{ And}\end{aligned}$$

an **apple** phone

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

$$\begin{aligned}\text{Phone} &\rightarrow 0.57 \text{ Phone} + 0.43 \text{ Apple} \\ \text{Apple} &\rightarrow 0.43 \text{ Phone} + 0.57 \text{ Apple} \\ \text{An} &\rightarrow 1 \text{ An}\end{aligned}$$

an **apple** and an orange

	Orange	Apple	And	An
Orange	1	0.71	0	0

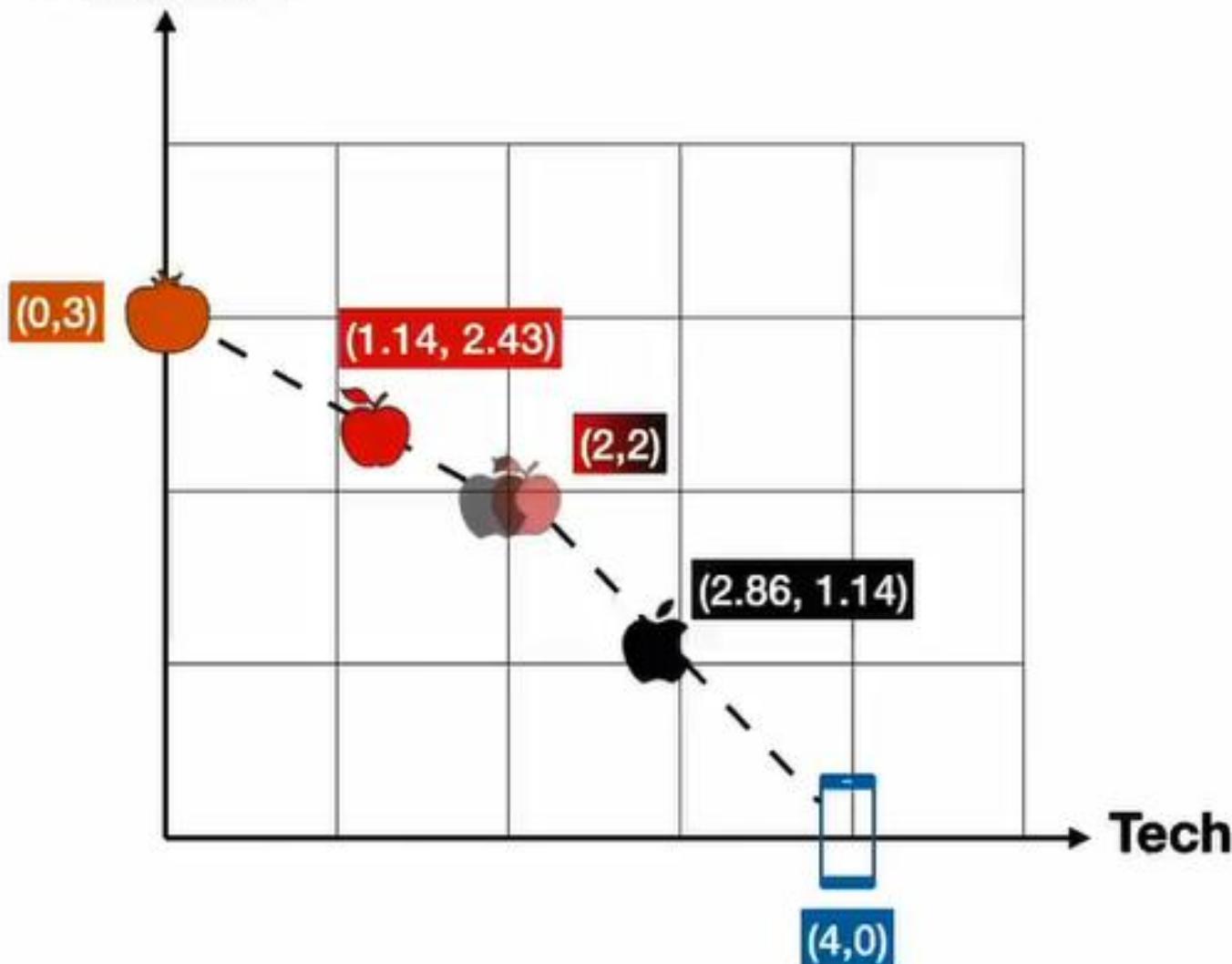
Ok, I'm kinda
lying to you...

$$\text{Orange} \rightarrow 0.57 \text{ Orange} + 0.43 \text{ Apple}$$

$$\frac{e^1 \text{Orange} + e^{0.71} \text{Apple} + e^0 \text{And} + e^0 \text{An}}{e^1 + e^{0.71} + e^0 + e^0}$$

$$\text{Orange} \rightarrow 0.4 \text{ Orange} + 0.3 \text{ Apple} + 0.15 \text{ And} + 0.15 \text{ An}$$

Fruitness



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

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

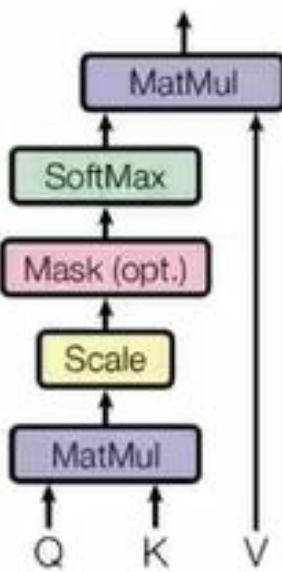
an **apple** phone

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

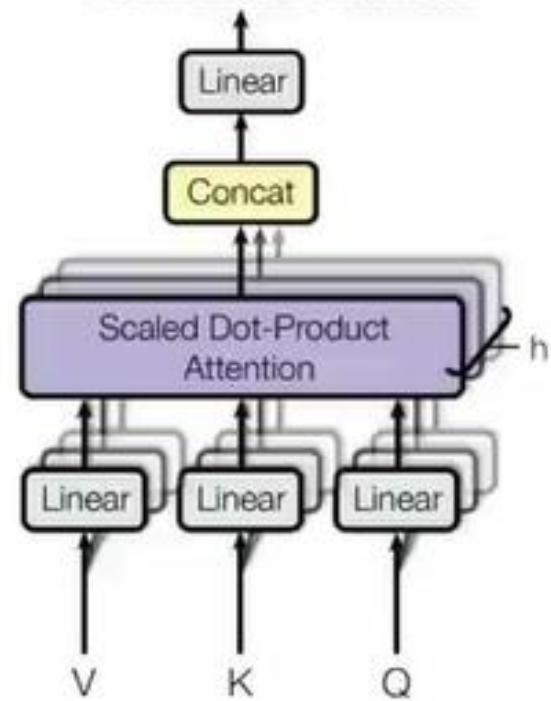
The keys, queries, and values matrices

Attention

Scaled Dot-Product Attention



Multi-Head Attention



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

$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \\ \text{where } \text{head}_i &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned}$$

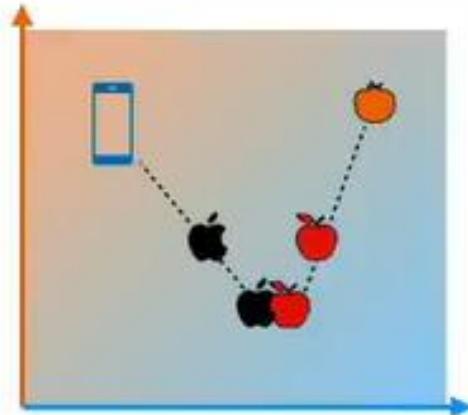
The Key and Query matrices

Keys

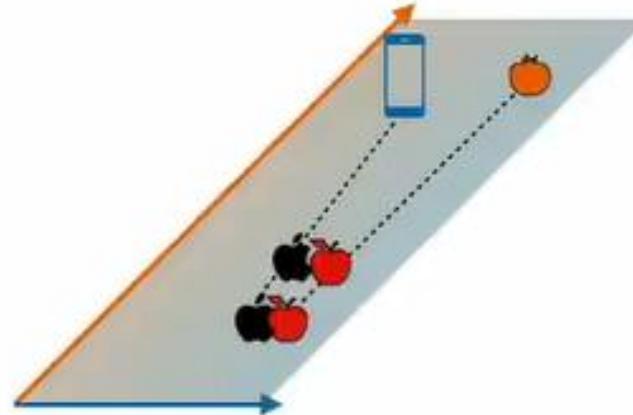
Queries

Values

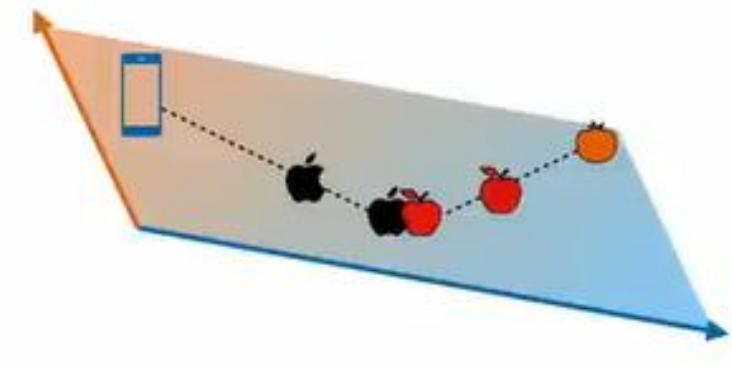
Get new embeddings from existing ones



Okay

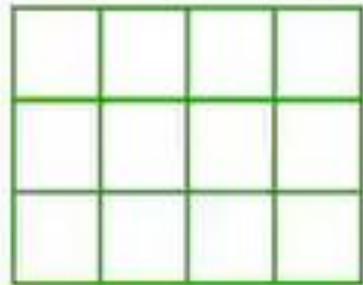


Bad

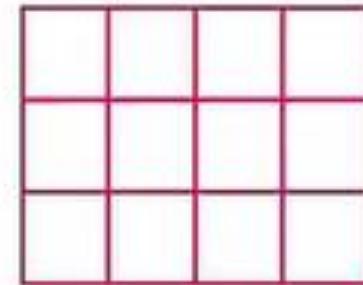


Good

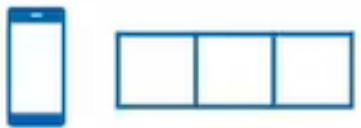
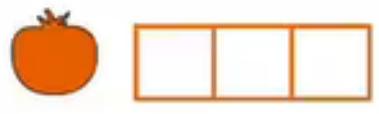
Keys



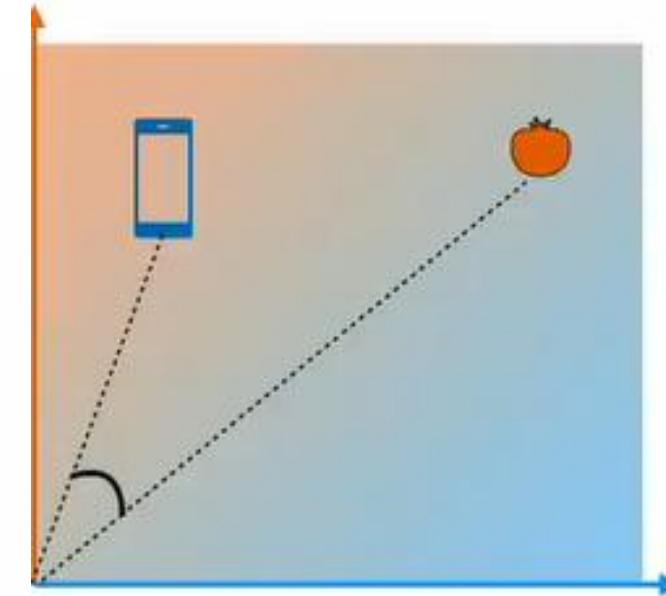
Queries



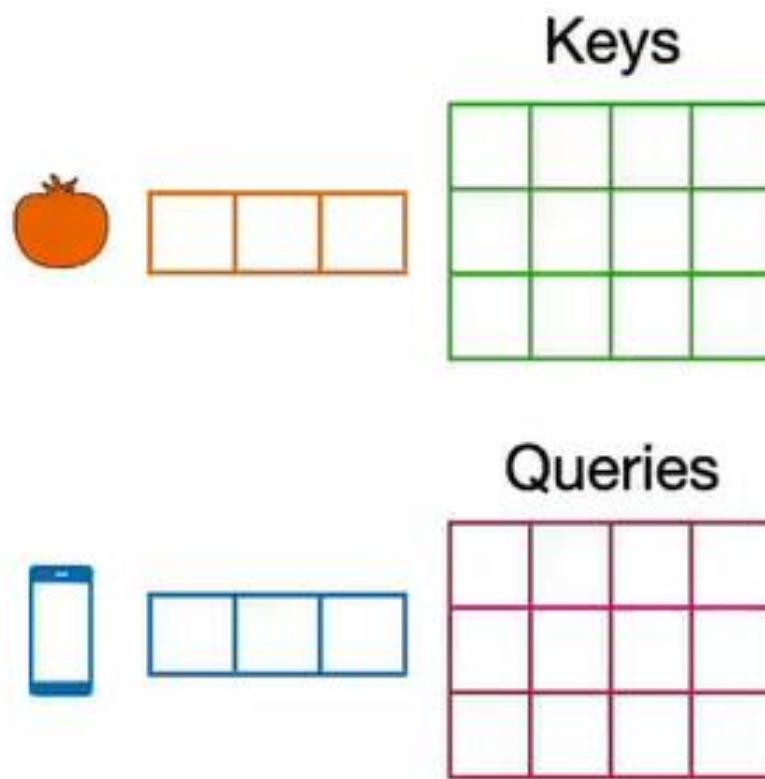
Similarity



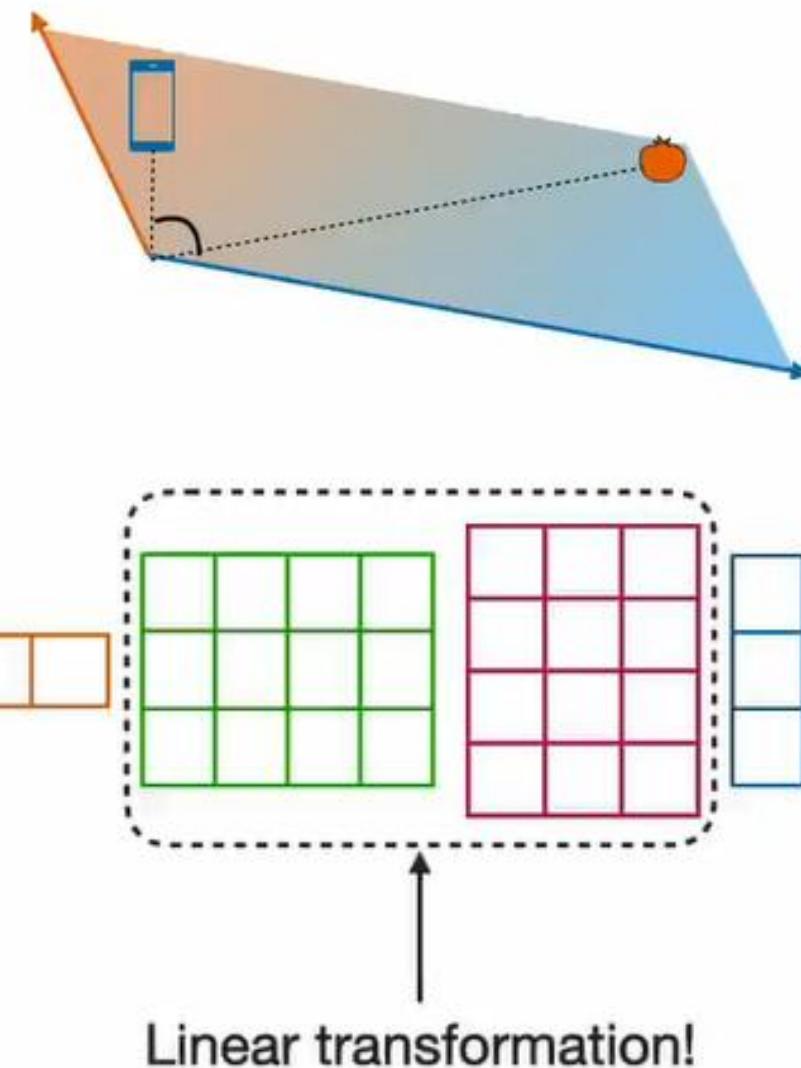
Similarity (,) =  



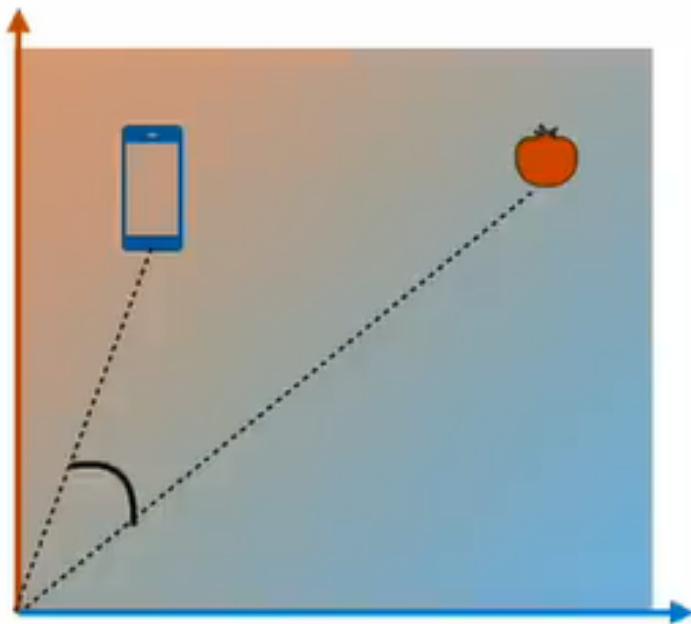
Keys and Queries Matrices



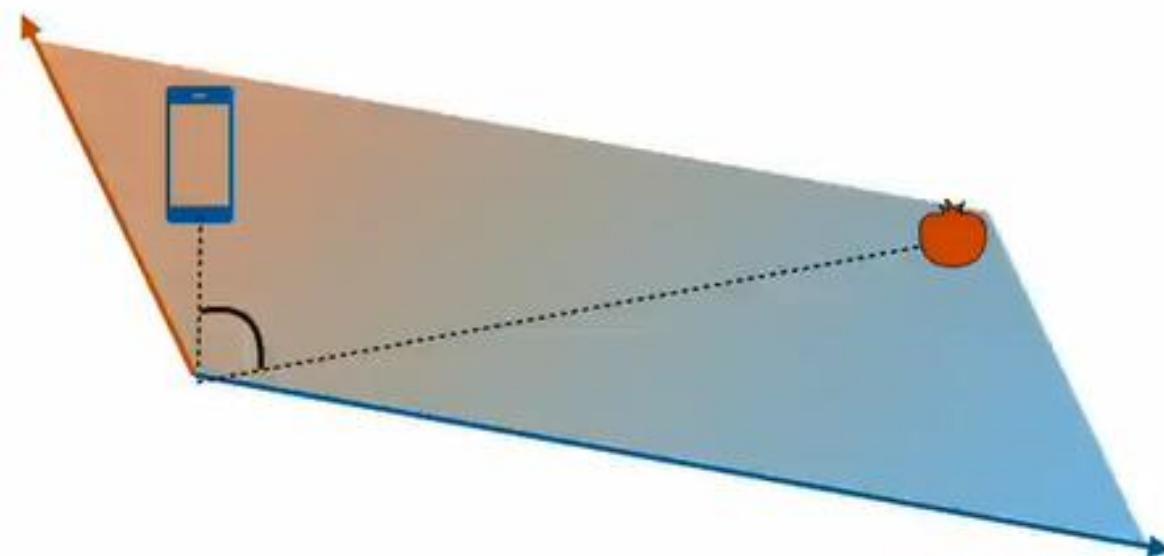
Similarity $(\text{Tomato}, \text{Smartphone}) =$



Similarity on a transformed embedding



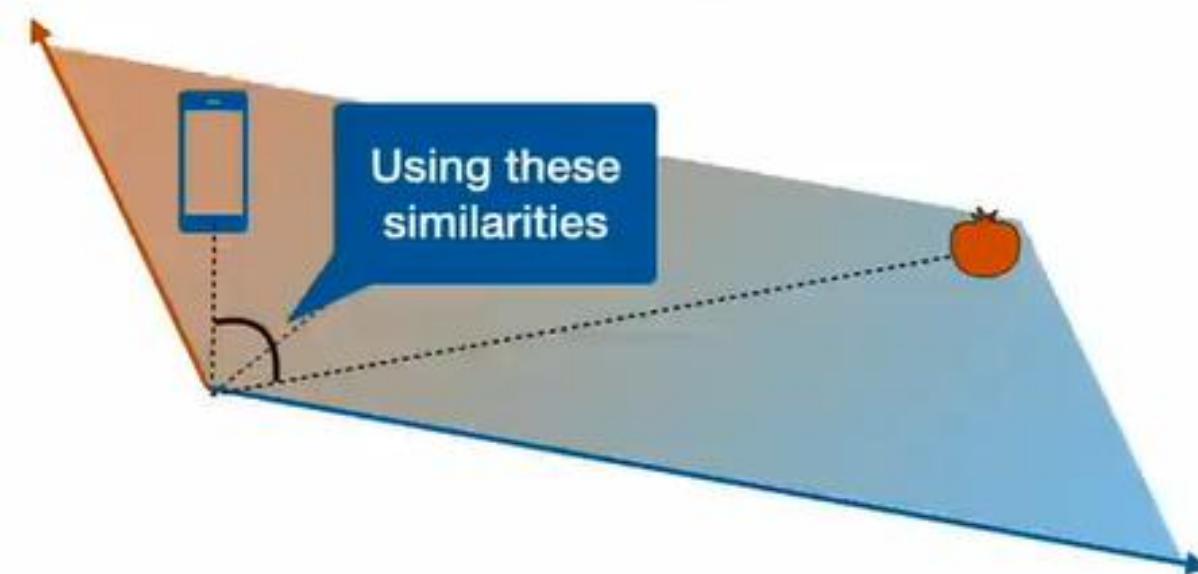
$$\text{Similarity}(\text{apple}, \text{phone}) = \begin{matrix} \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \end{matrix}$$



$$\text{Similarity}(\text{apple}, \text{phone}) = \begin{matrix} \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \end{matrix} \quad \begin{matrix} \boxed{} & \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \boxed{} & \boxed{} & \boxed{} \end{matrix} \quad \begin{matrix} \boxed{} & \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \boxed{} & \boxed{} & \boxed{} \\ \boxed{} & \boxed{} & \boxed{} & \boxed{} \end{matrix} \quad \begin{matrix} \boxed{} & \boxed{} \\ \boxed{} & \end{matrix}$$

The values matrix

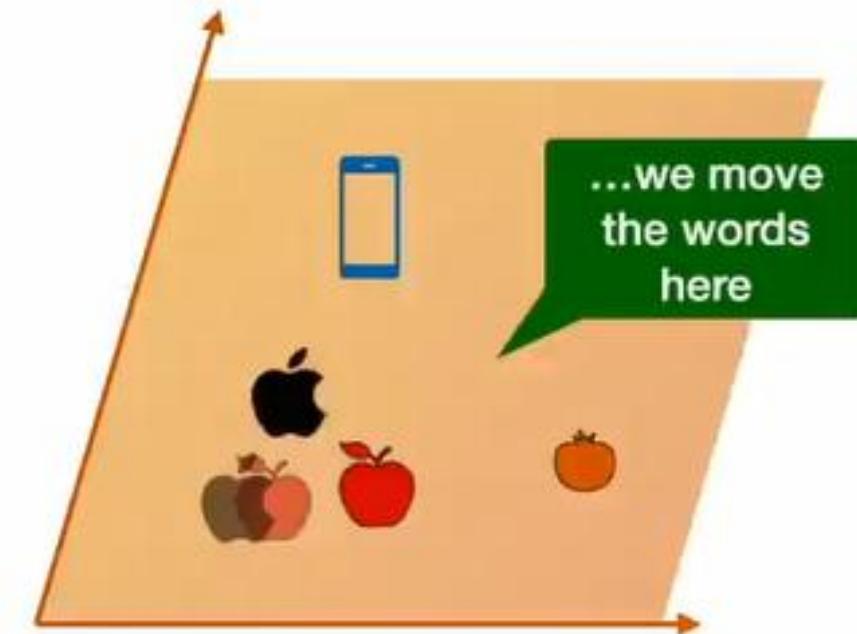
Values matrix



Best embedding for finding similarities

Keys

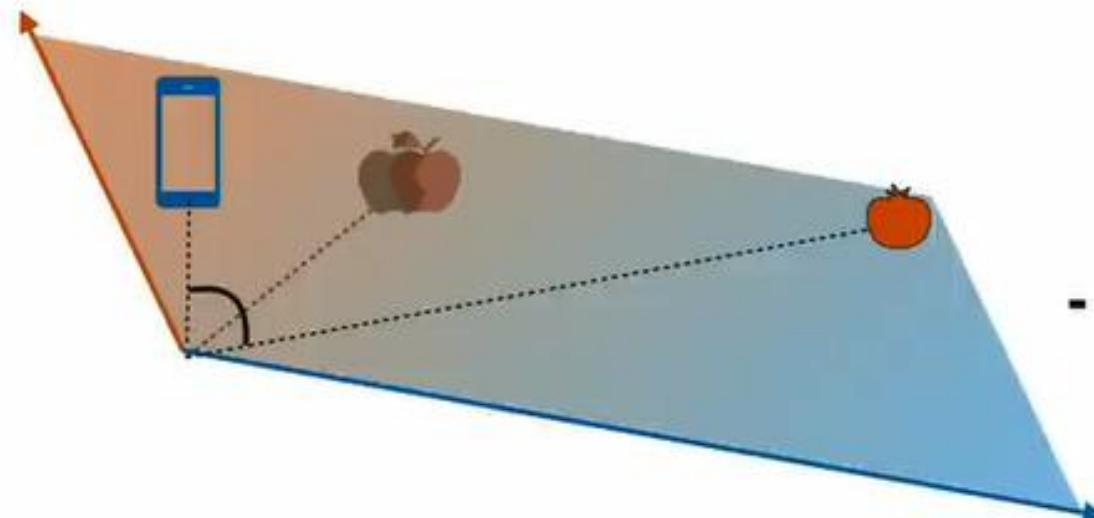
Queries



Best embedding for finding the next word

Values

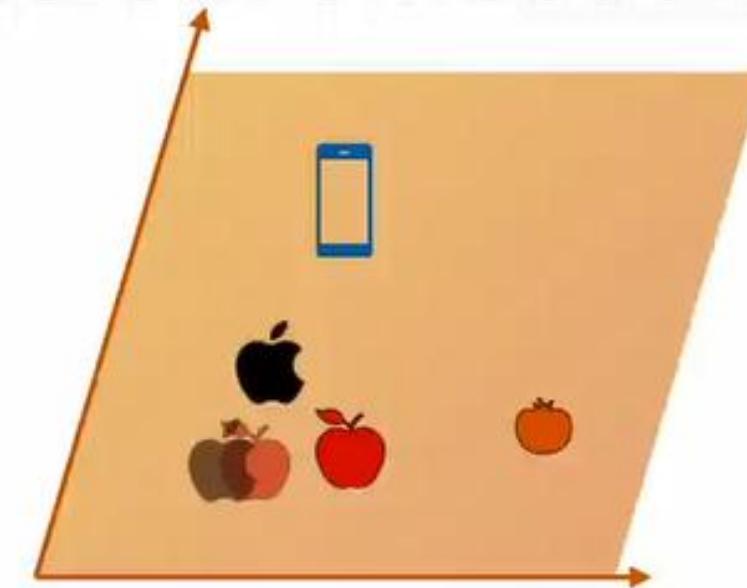
Why moving words on a different embedding?



Best embedding for finding similarities

This embedding(s) know features of the words

- Color
- Size
- Fruitness
- Technology



Best embedding for finding the next word

This embedding knows when two words could appear in the same context

I want to buy a _____

- car
- apple
- phone

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}

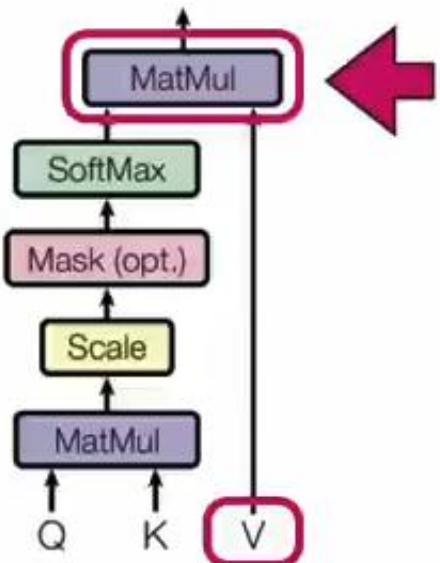
$$\begin{aligned} \text{apple} \longrightarrow & 0.3 \cdot \text{orange} \\ & + 0.4 \cdot \text{apple} \\ & + 0.15 \cdot \text{and} \\ & + 0.15 \cdot \text{an} \end{aligned}$$

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

MULTi-HEAD ATTENTiON

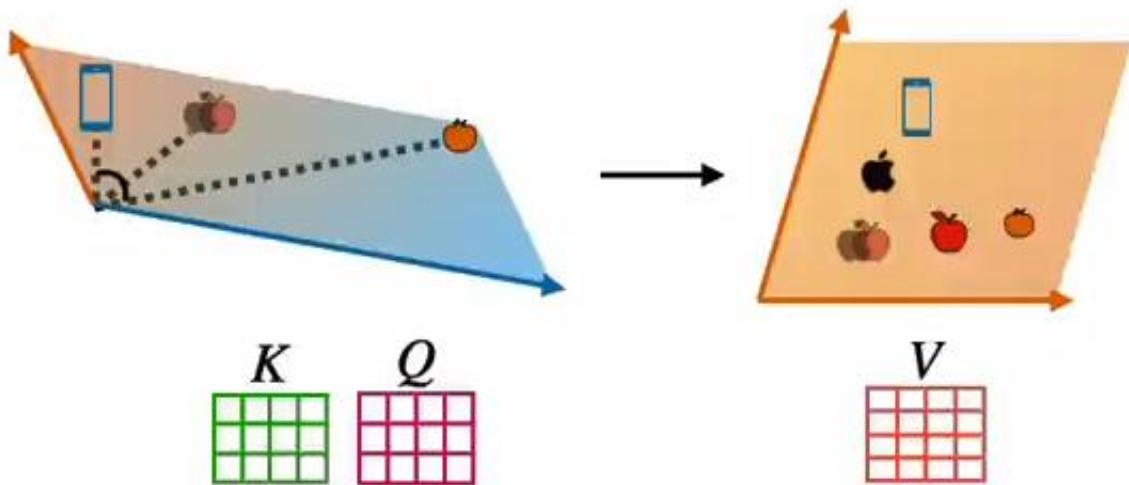
Self-attention

Scaled Dot-Product Attention

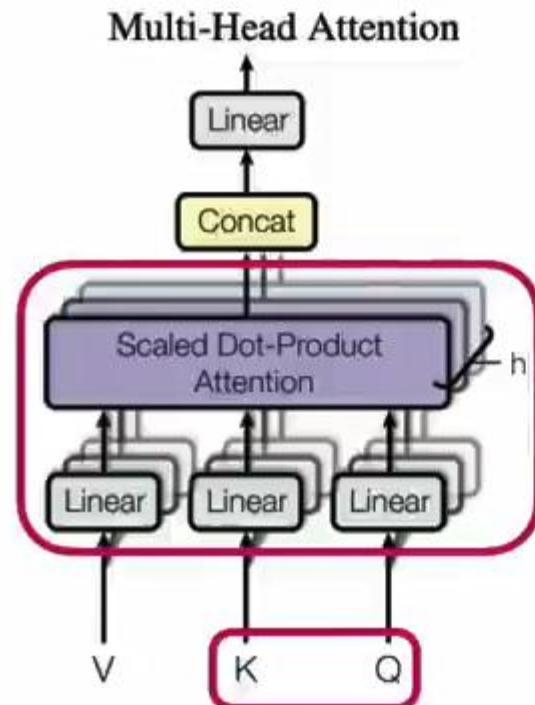


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

WE WILL USE THE VALUE MATRIX TO TRANSFORM THIS EMBEDDINGS IN A BETTER ONE. SO WE USE THE LINEAR TRANSFORMATION TO FIND THE EMBEDDING ON THE RIGHT AND MOVE THE WORDS IN A BETTER WAY.



Multi-head attention

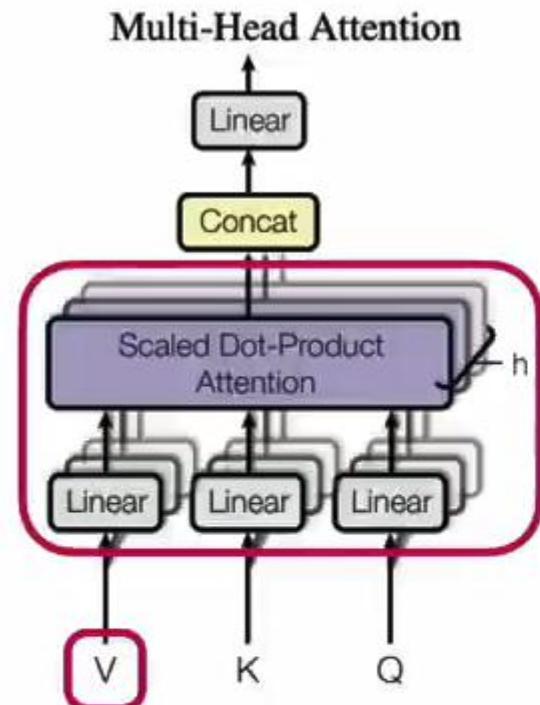


$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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



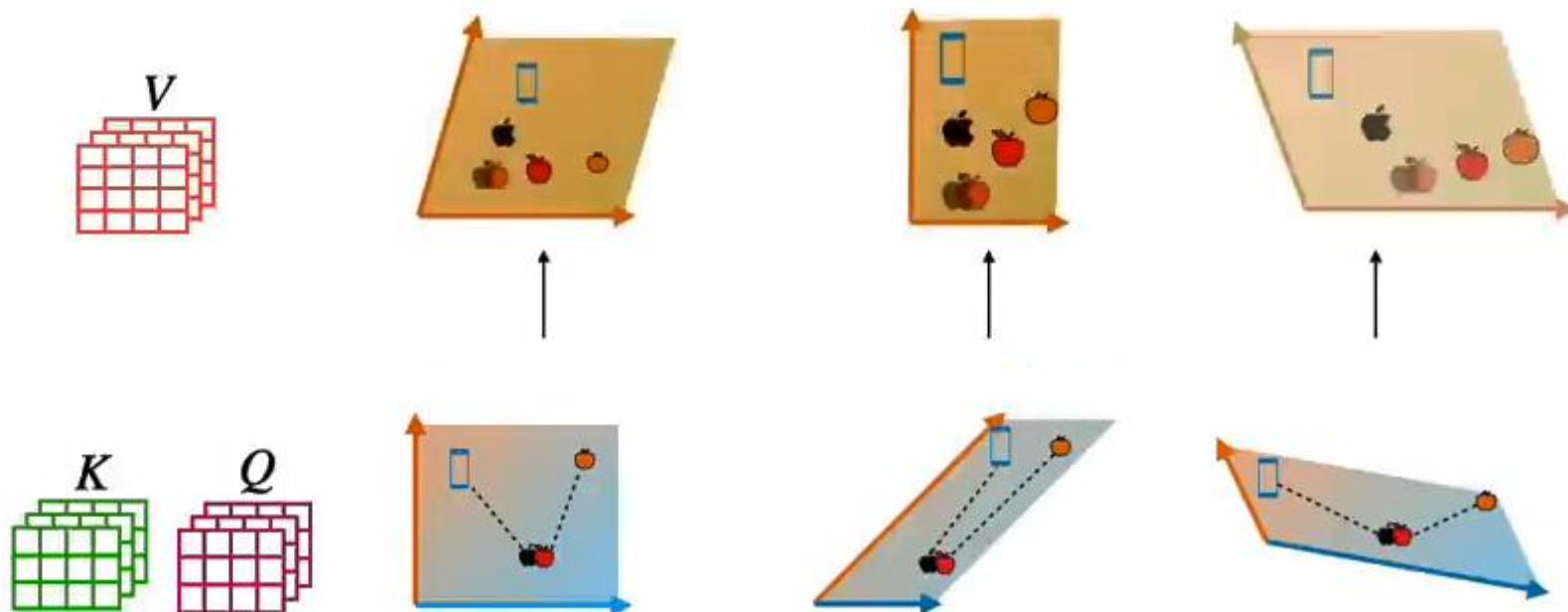
Multi-head attention



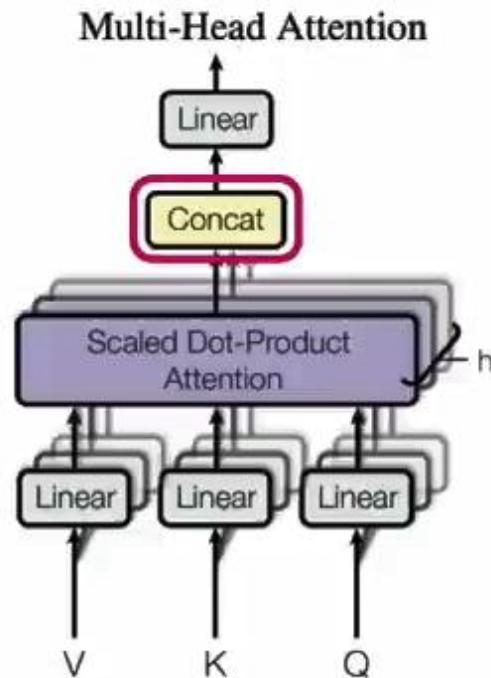
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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

As we said before, this V matrix transform these embedding where we calculate the similarities into embeddings where we can move the words around

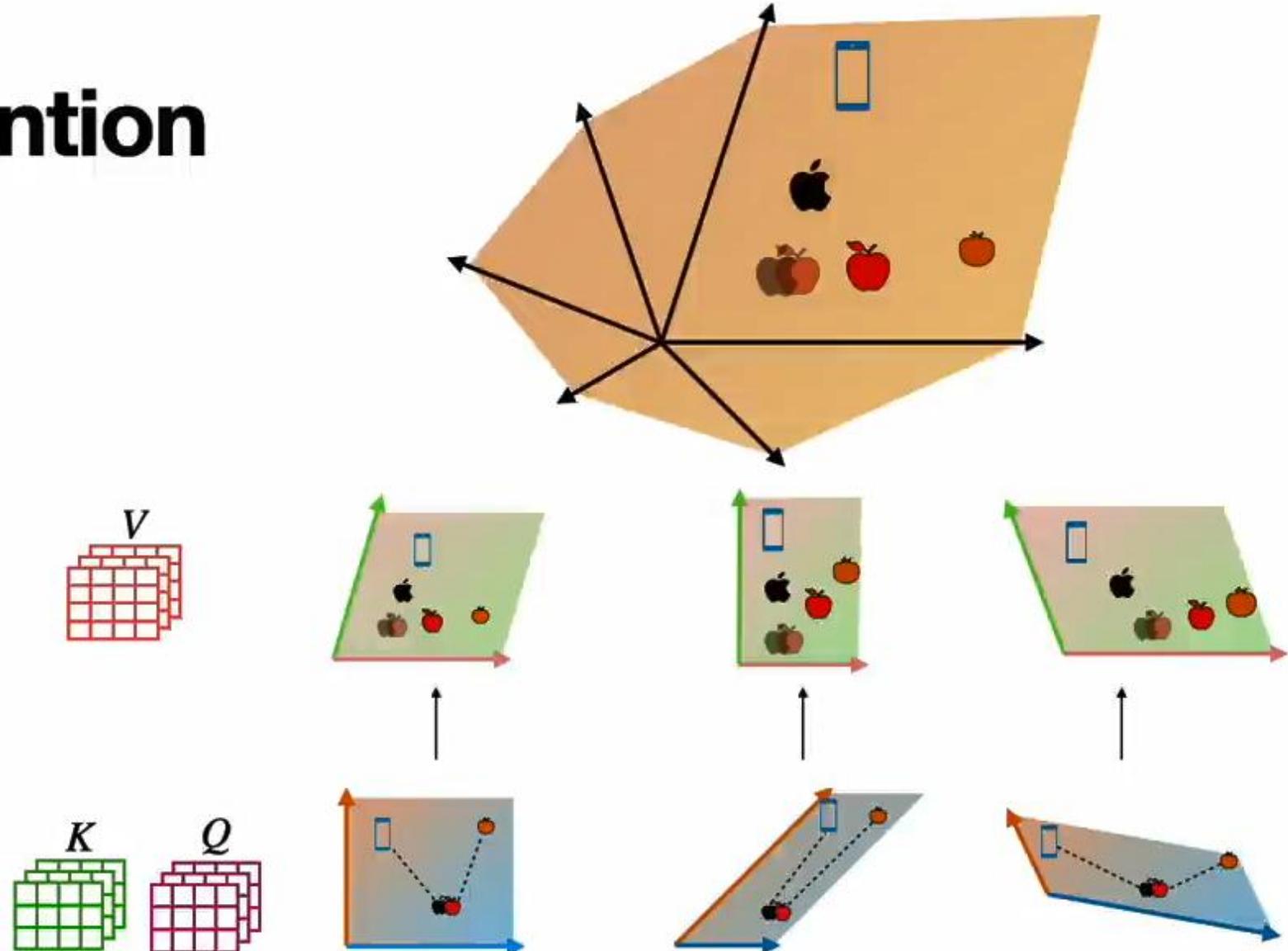


Multi-head attention



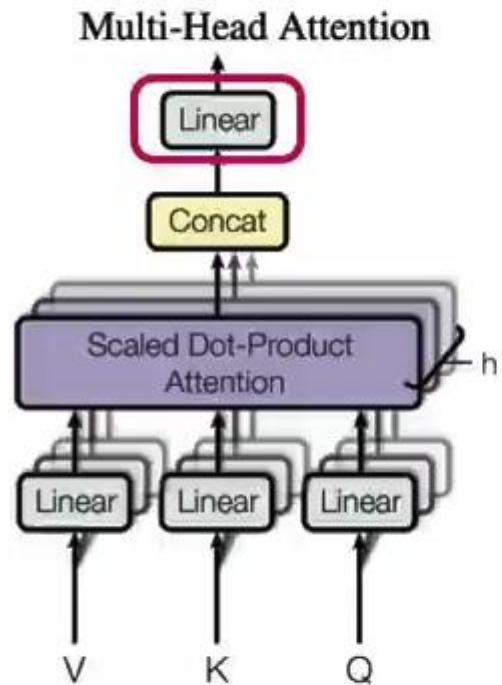
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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



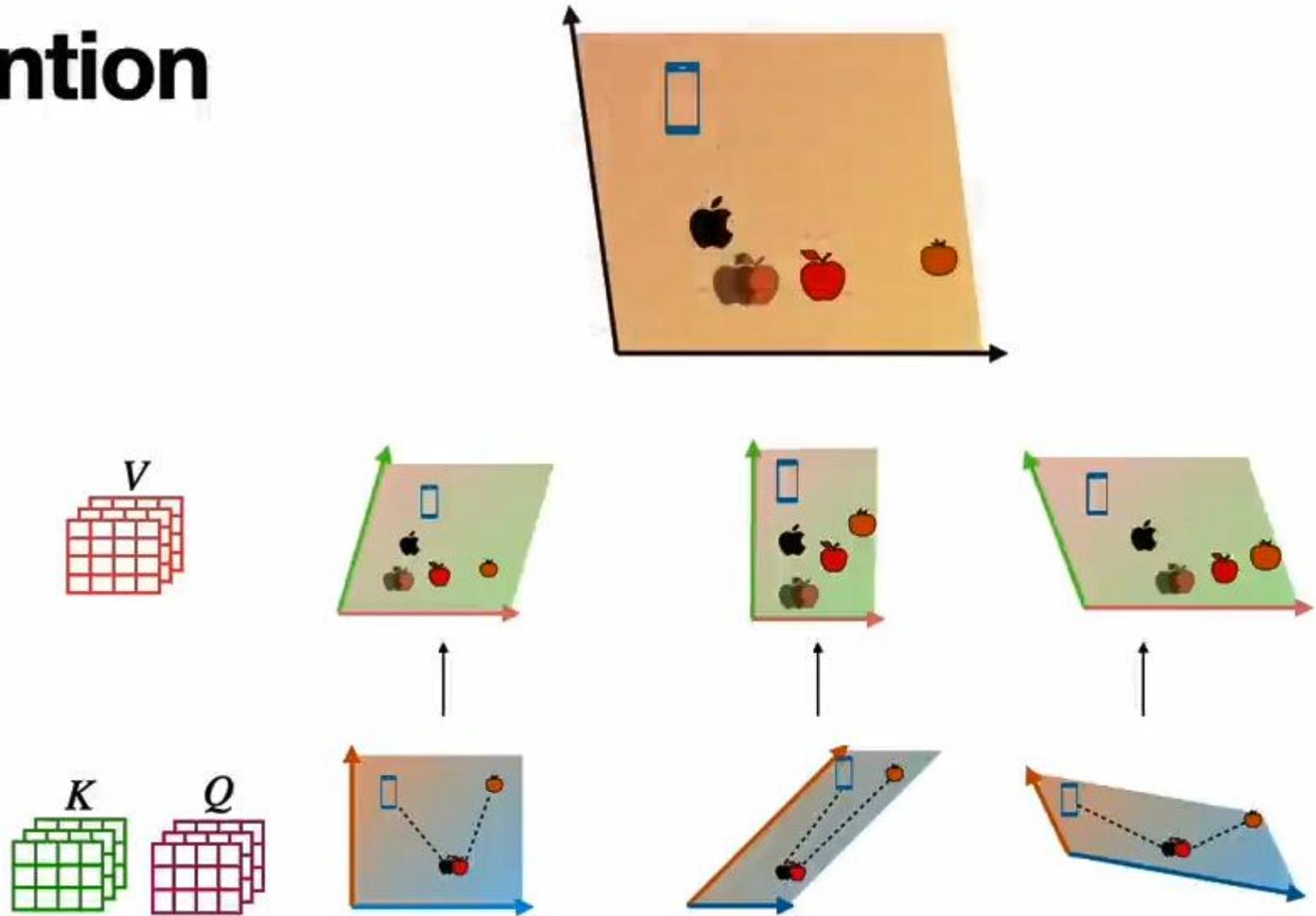
The concat the result, and I obtain an embeddings of 6 dimension. To reconduct this into lower dimension I will apply this Linear step (rectangular matrix)

Multi-head attention

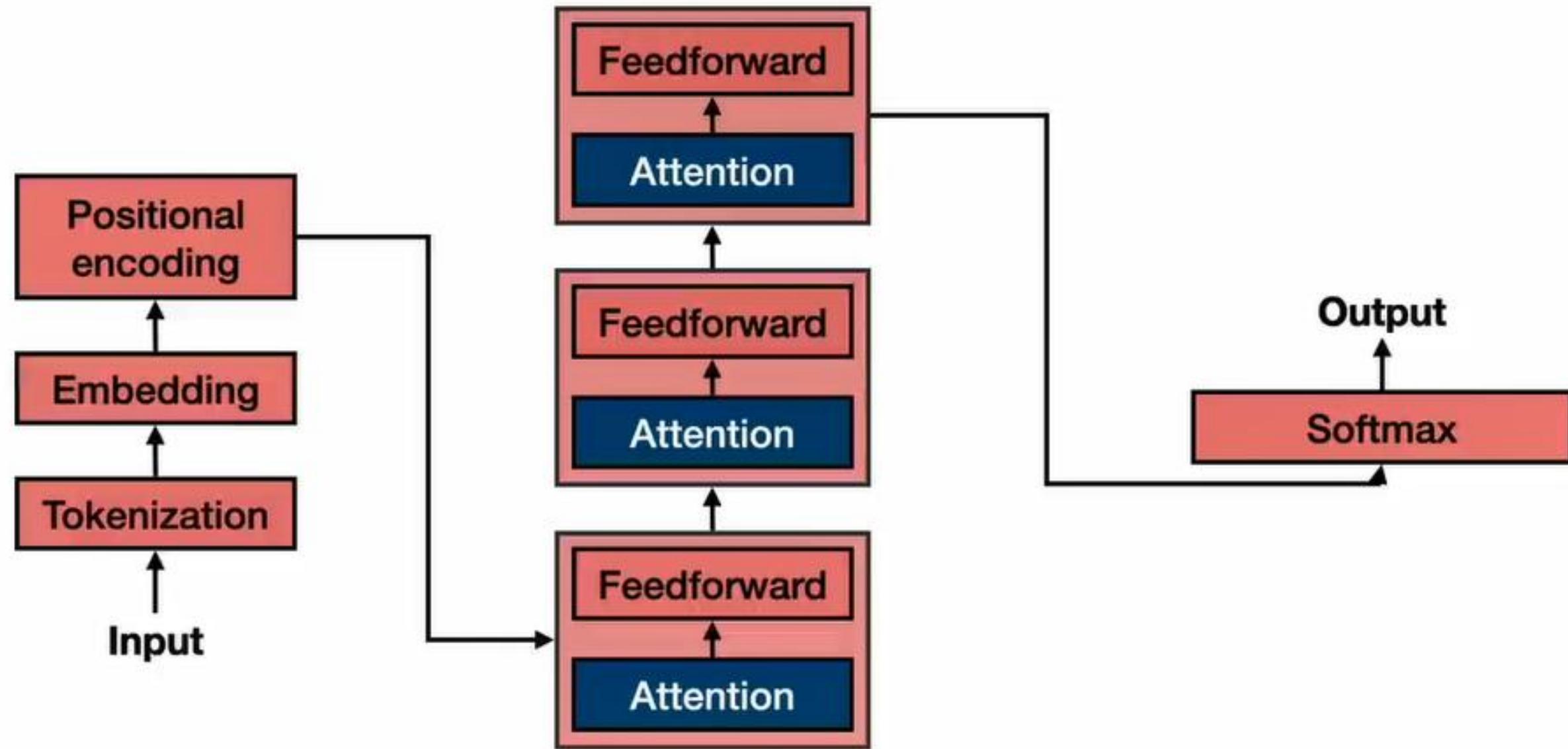


$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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



Weights get trained with the transformer model





That's all folks!