
Attention and Transformers

Saehwa Kim

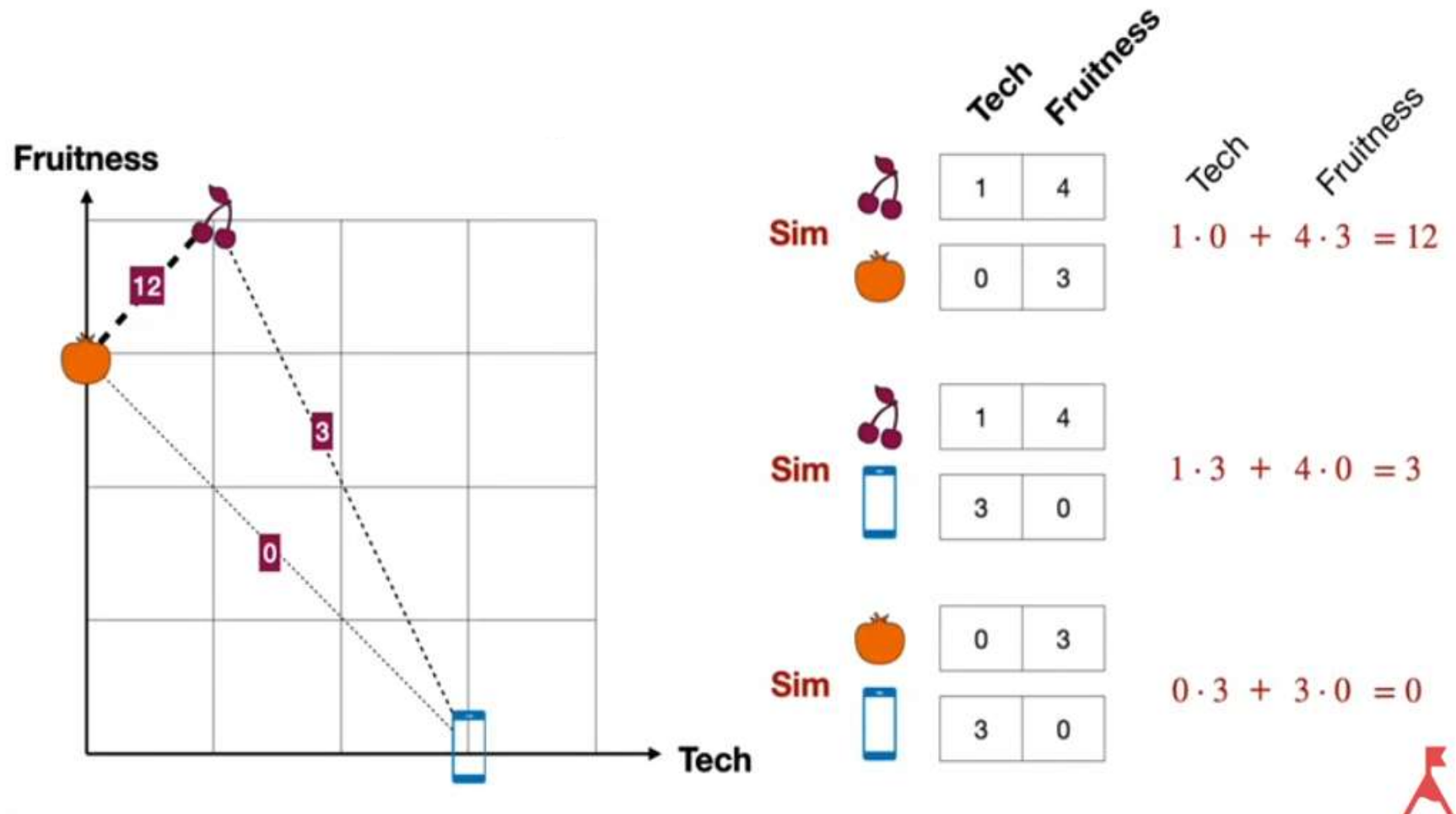
Information and Communications Engineering
Hankuk University of Foreign Studies



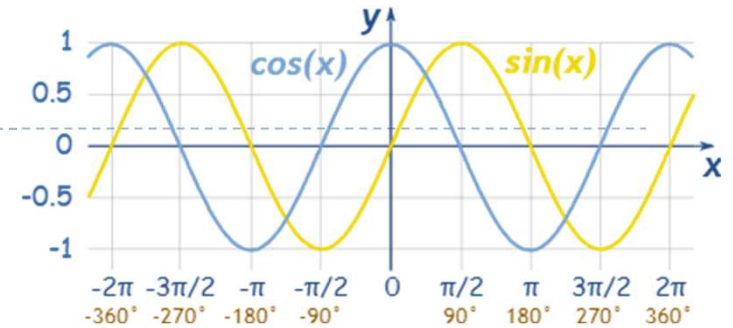
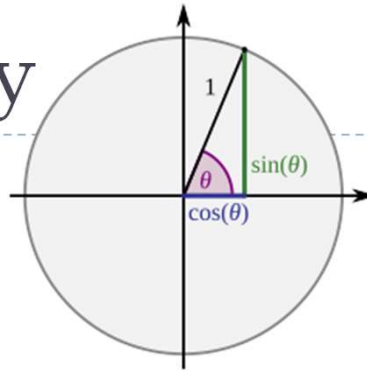
Outline

- ▶ Dot Product
- ▶ Attention Scores
- ▶ New Word Embedding with Attention Scores
- ▶ Attention in Transformers
- ▶ Transformer Architecture
- ▶ GPT and Transformers

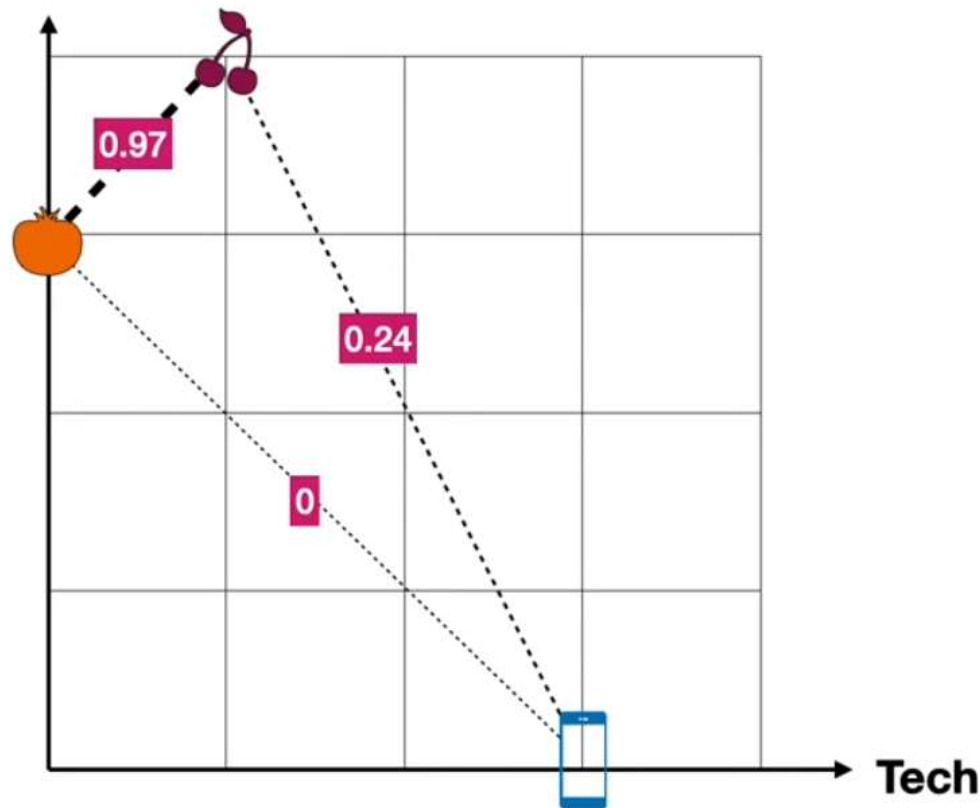
Dot Product



Cosine Similarity



Fruitness

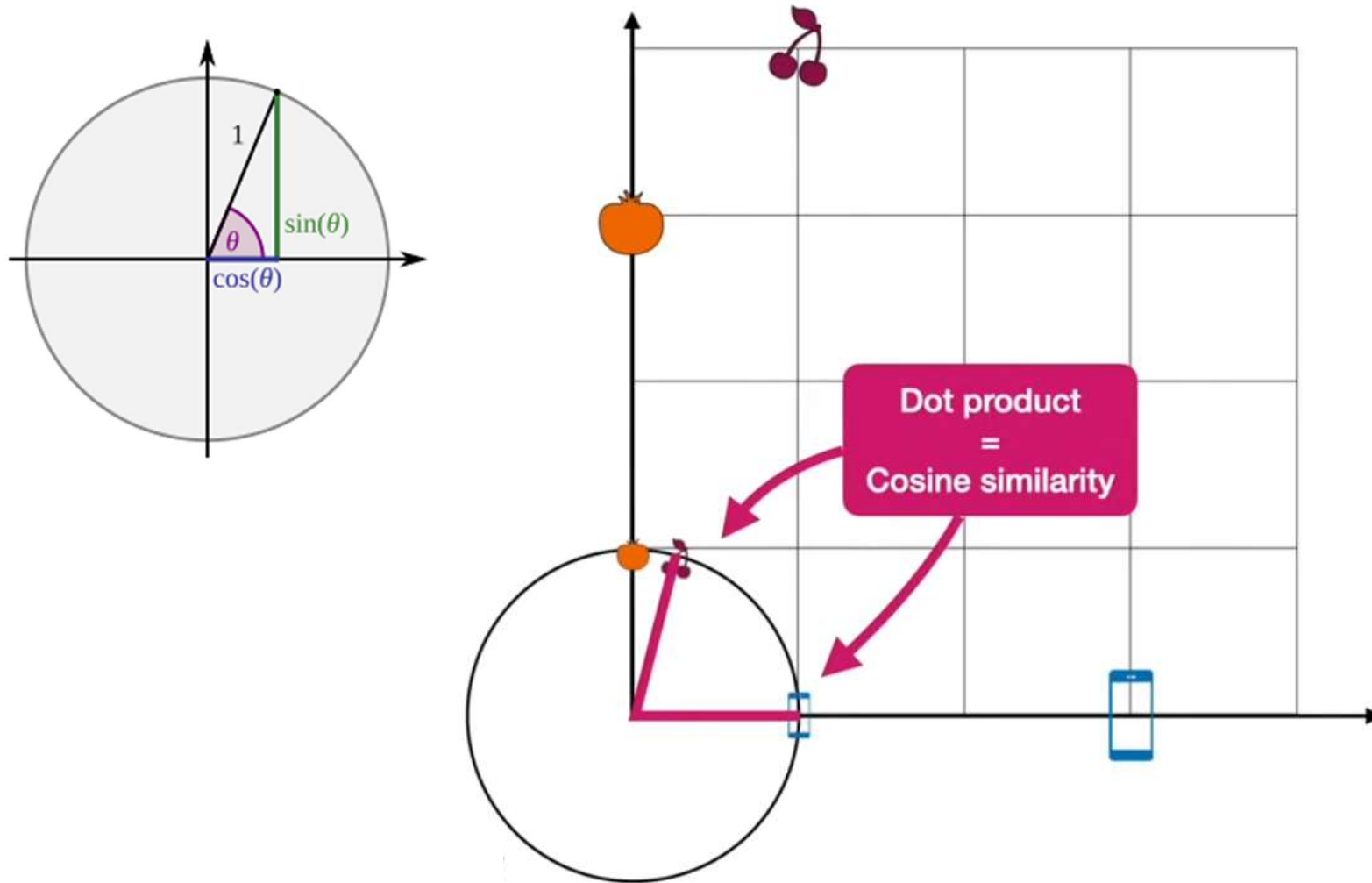


Sim 
 $\cos(14^\circ) = 0.97$

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

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

Dot Product \approx Cosine Similarity (1/2)

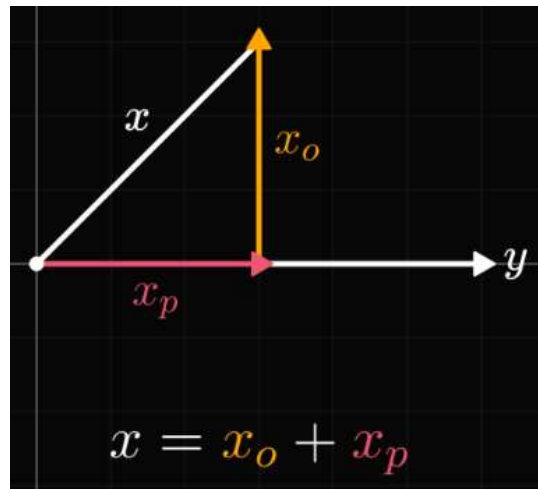


Dot Product \approx Cosine Similarity (2/2)

$$\langle x, y \rangle = \sum_{i=1}^n x_i y_i$$

$$x = (x_1, \dots, x_n)$$

$$y = (y_1, \dots, y_n)$$



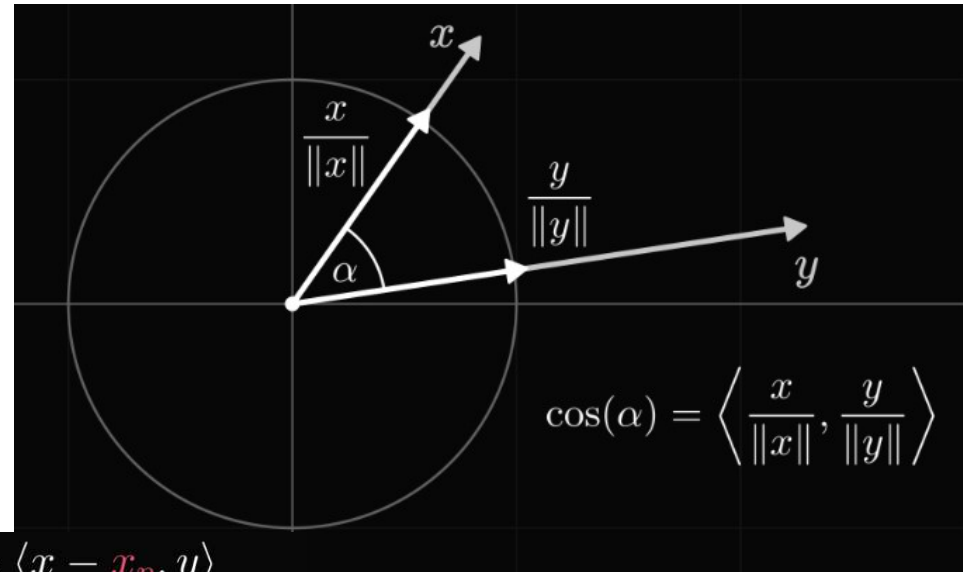
$$x_p = cy, \quad (c \in \mathbb{R})$$

$$\langle x_o, y \rangle = 0$$

$$\begin{aligned} \langle x_o, y \rangle &= \langle x - x_p, y \rangle \\ &= \langle x - cy, y \rangle \\ &= \langle x, y \rangle - c \langle y, y \rangle \\ &= 0 \end{aligned}$$

$$c = \frac{\langle x, y \rangle}{\langle y, y \rangle} = \frac{\langle x, y \rangle}{\|y\|^2}$$

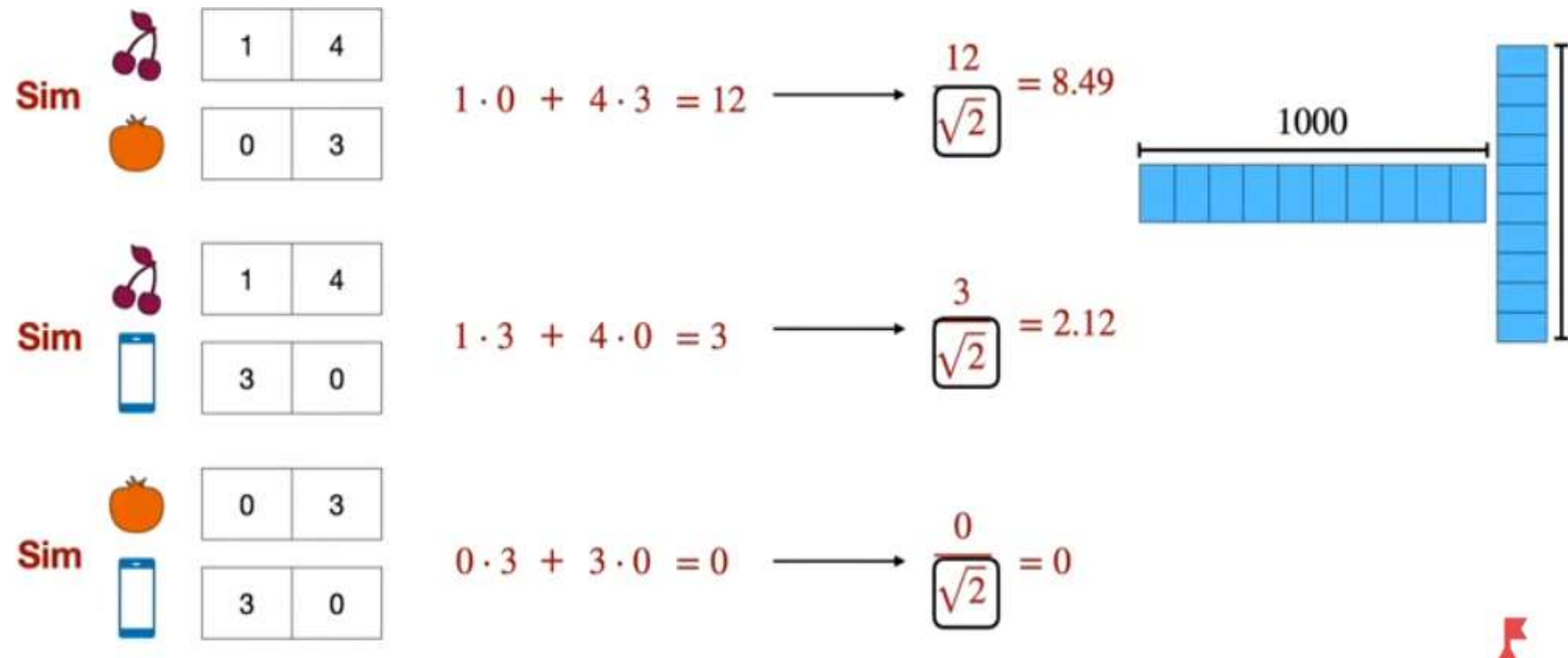
\uparrow
 $\langle y, y \rangle = \|y\|^2$
 (magnitude of y)



$$\begin{aligned} \cos(\alpha) &= \left\langle \frac{x}{\|x\|}, \frac{y}{\|y\|} \right\rangle \\ &= \frac{\langle x, y \rangle}{\|x\| \|y\|} \end{aligned}$$

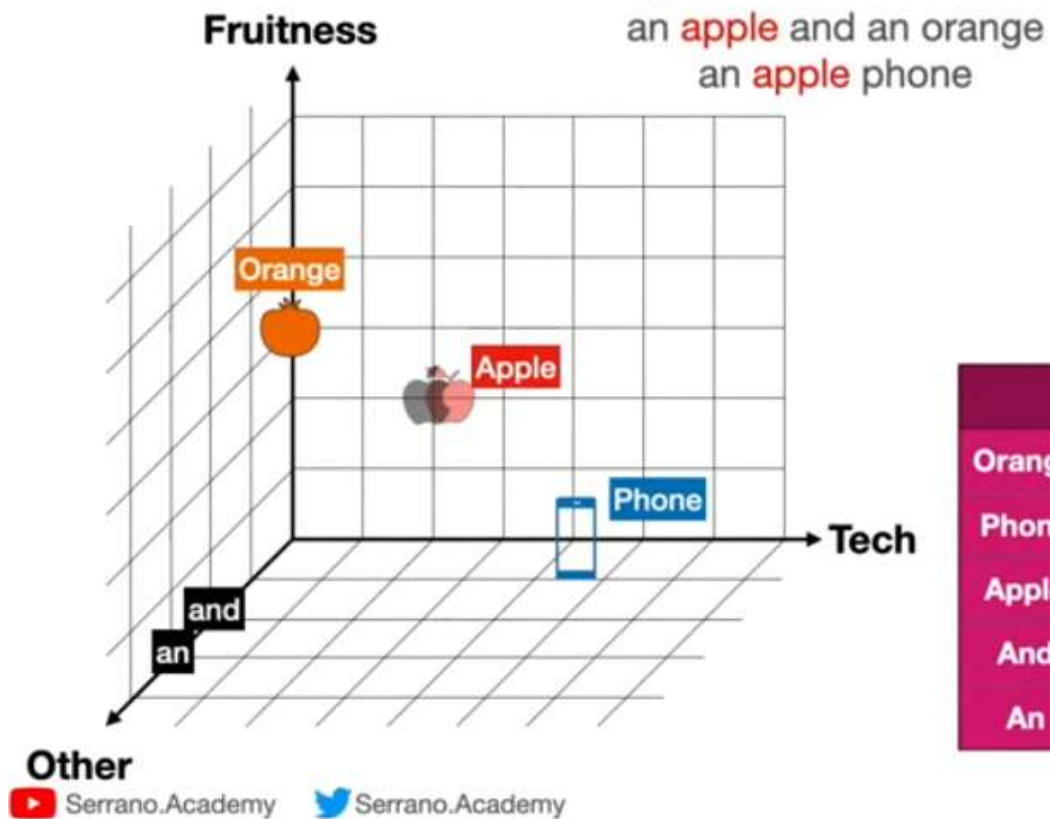
Scaled Dot Product

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



Attention Scores

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	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



New Word Embedding with Attention Scores (1/4)

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

Orange $\rightarrow 1$ Orange + 0.71 Apple

Apple $\rightarrow 0.71$ Orange + 1 Apple

And $\rightarrow 1$ And + 1 An

An $\rightarrow 1$ An + 1 And

an **apple** phone

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

Phone $\rightarrow 1$ Phone + 0.71 Apple

Apple $\rightarrow 0.71$ Phone + 1 Apple

An $\rightarrow 1$ An

New Word Embedding with Attention Scores (2/4)

► Scaling

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

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

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

$$\text{And} \rightarrow 0.5 \text{ And} + 0.5 \text{ An}$$

$$\text{An} \rightarrow 0.5 \text{ An} + 0.5 \text{ And}$$

an **apple** phone

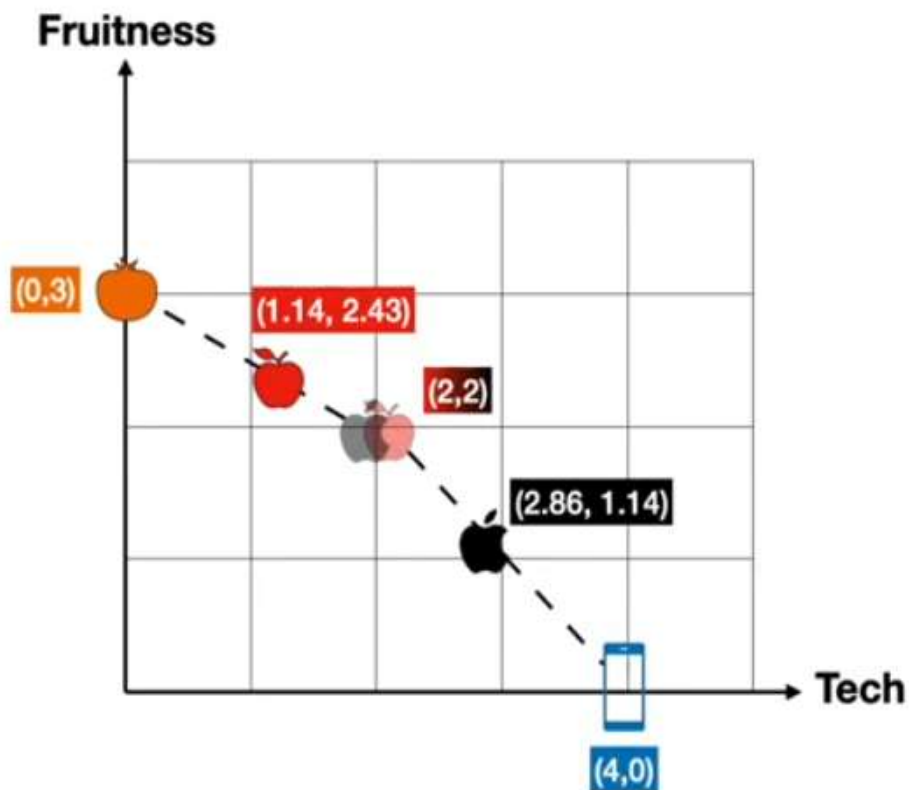
	Phone	Apple	An
Phone	1	0.71	0
Apple	0.71	1	0
An	0	0	1

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

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

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

New Word Embedding with Attention Scores (3/4)

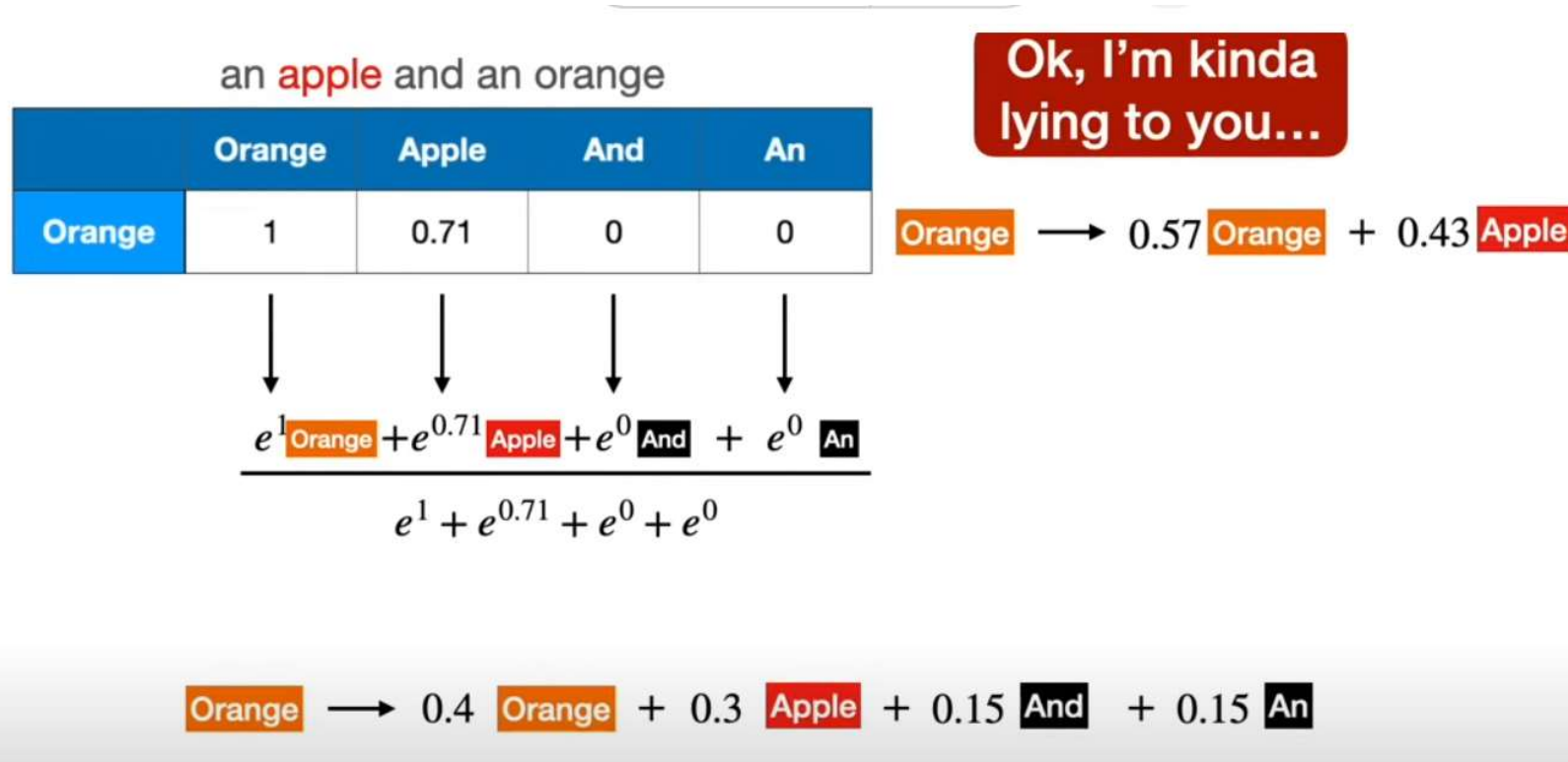


an **apple** and an orange
Apple \rightarrow 0.43 **Orange** + 0.57 **Apple**

an **apple** phone
Apple \rightarrow 0.43 **Phone** + 0.57 **Apple**

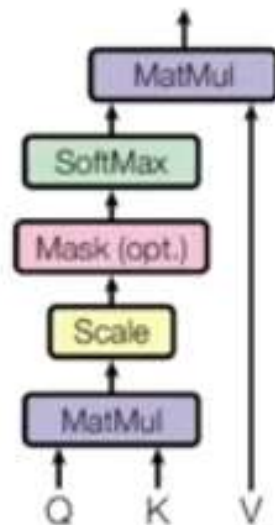
New Word Embedding with Attention Scores (4/4)

► Softmax normalization



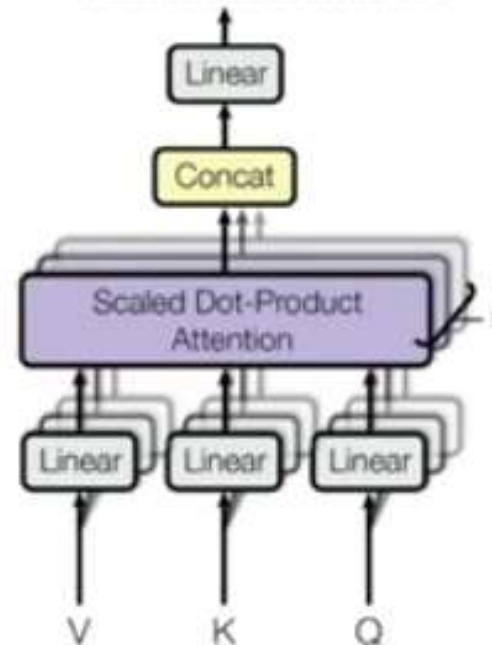
Attention in Transformers

Scaled Dot-Product Attention



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

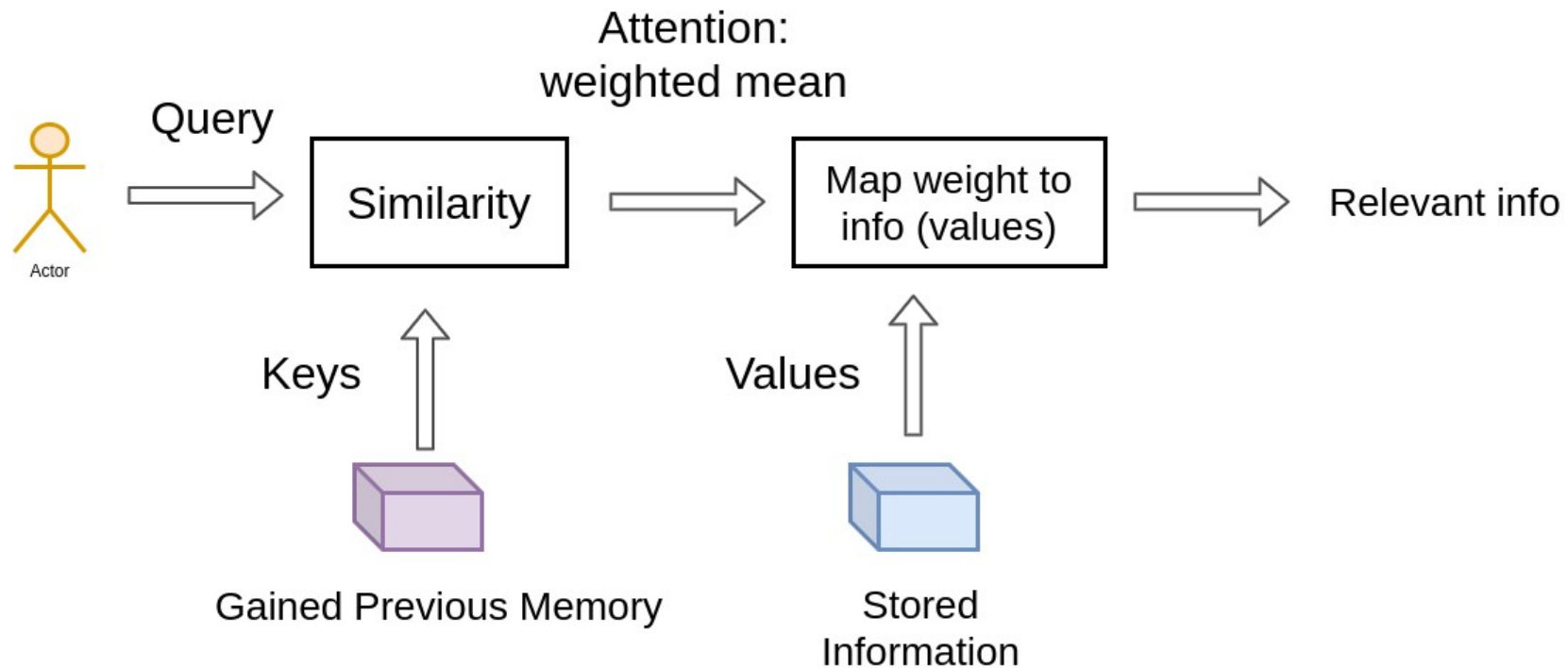
Multi-Head Attention



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

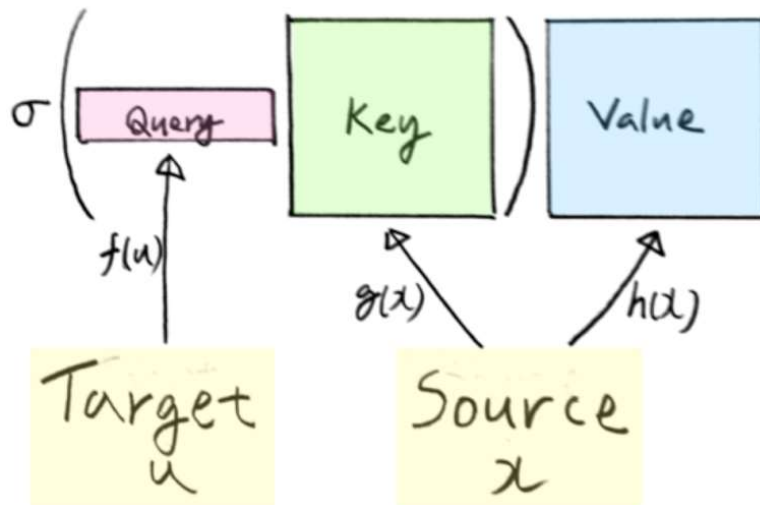
where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

Notion of Query, Key, and Values

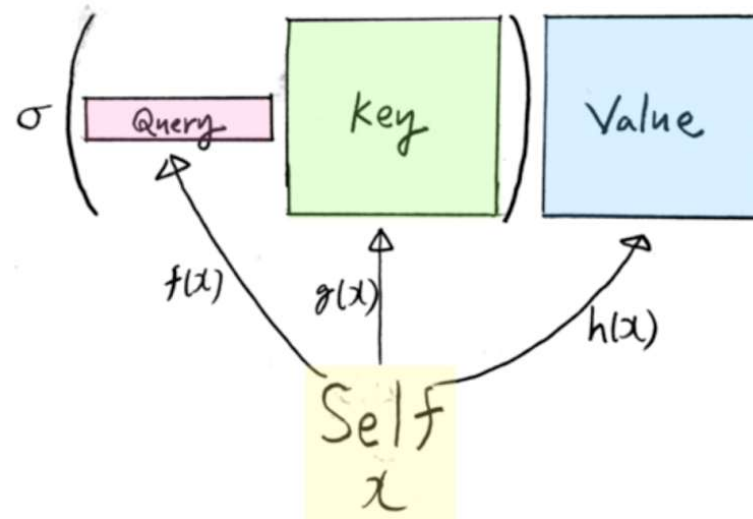


Source-Target Attention vs. Self-Attention

(Source-Target-Attention)

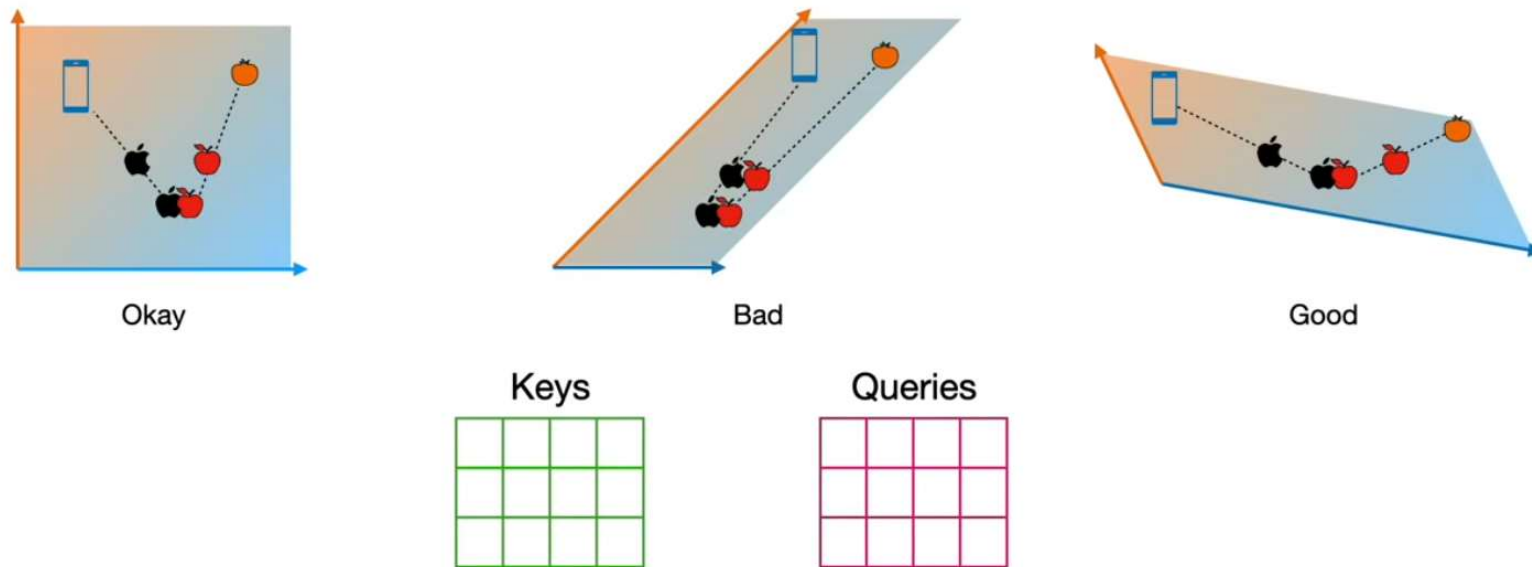


(Self-Attention)



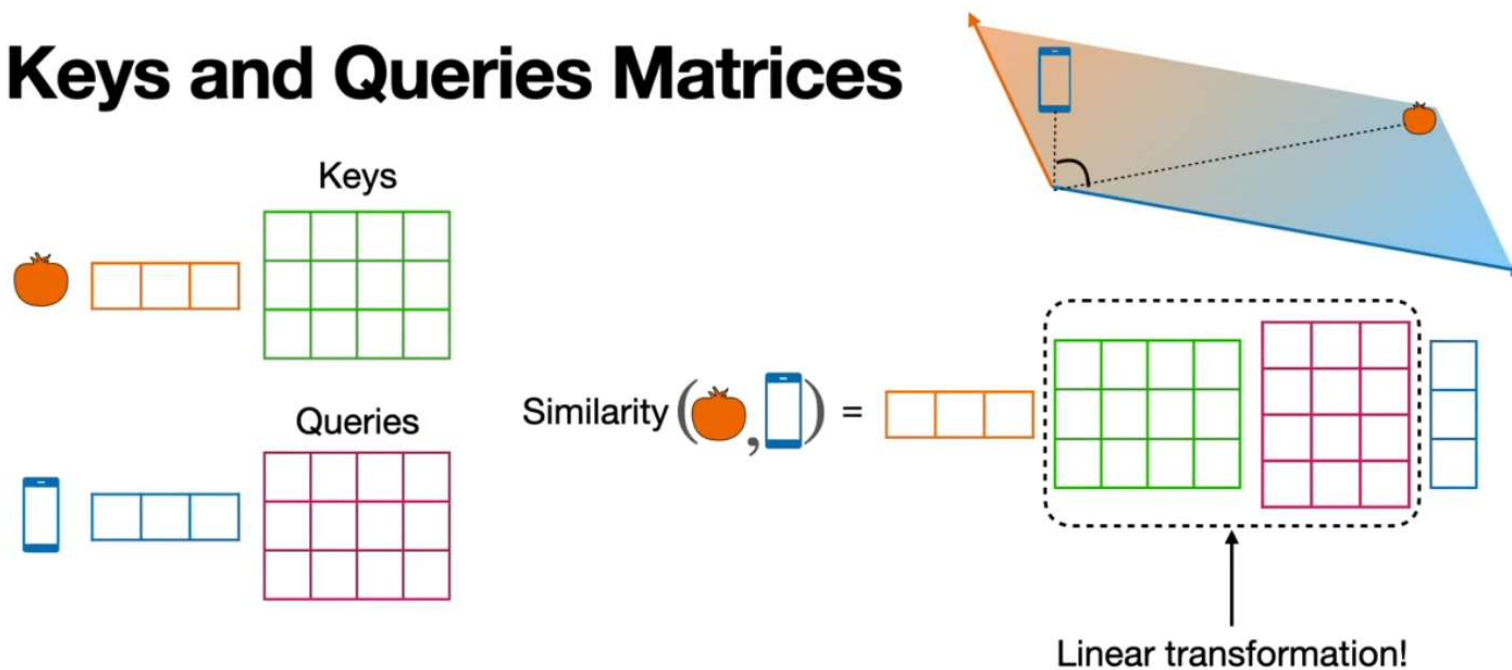
Weights for Keys and Queries

Get new embeddings from existing ones

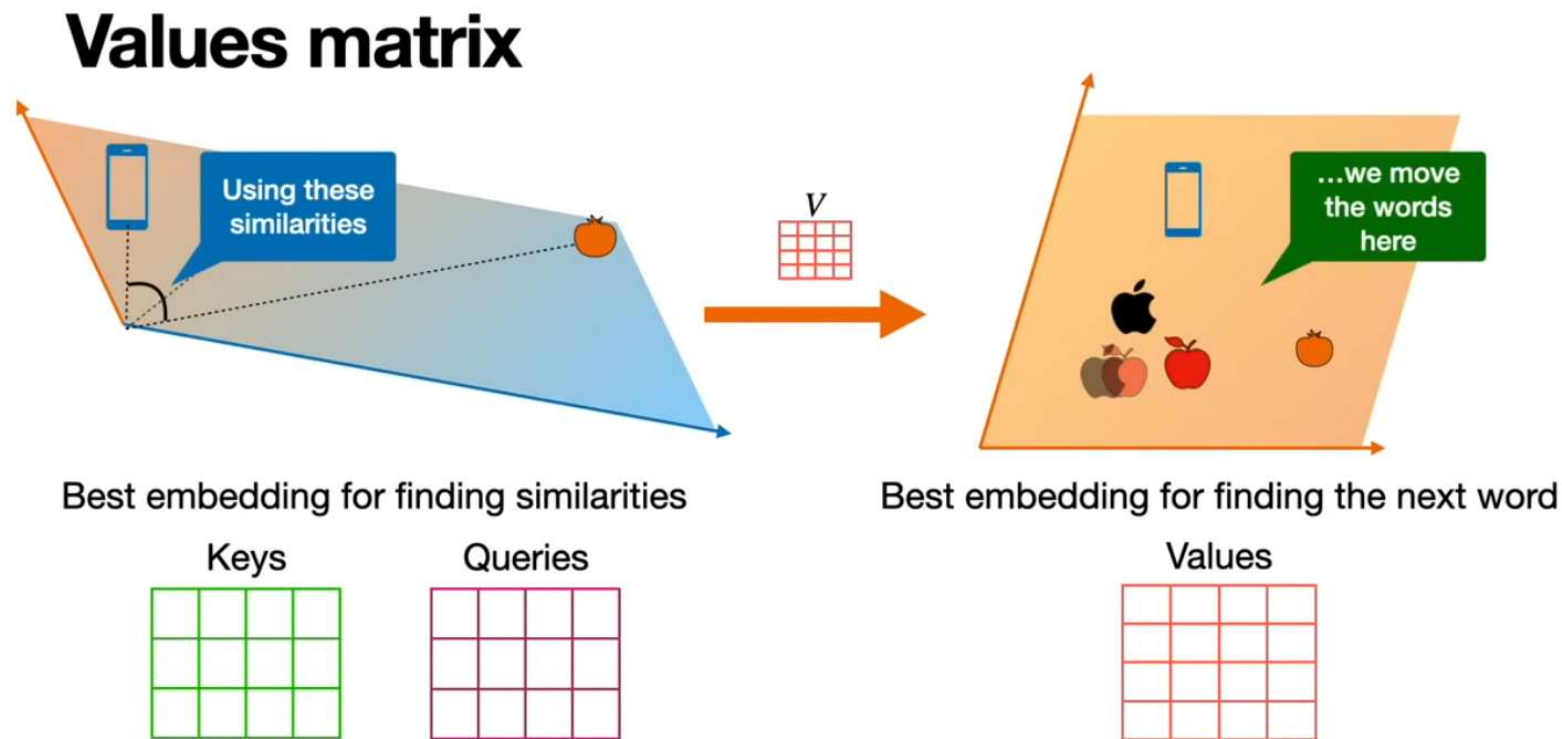


Getting Similarity of Queries and Keys

Keys and Queries Matrices

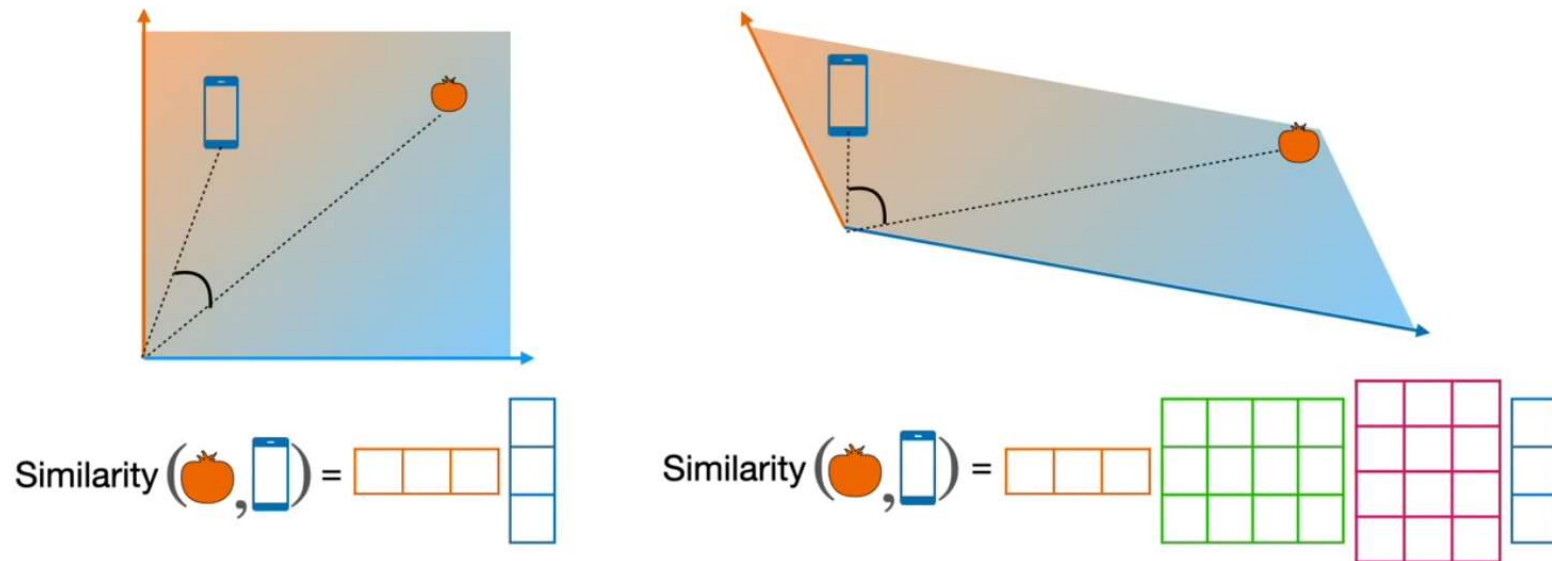


Applying Value Matrix



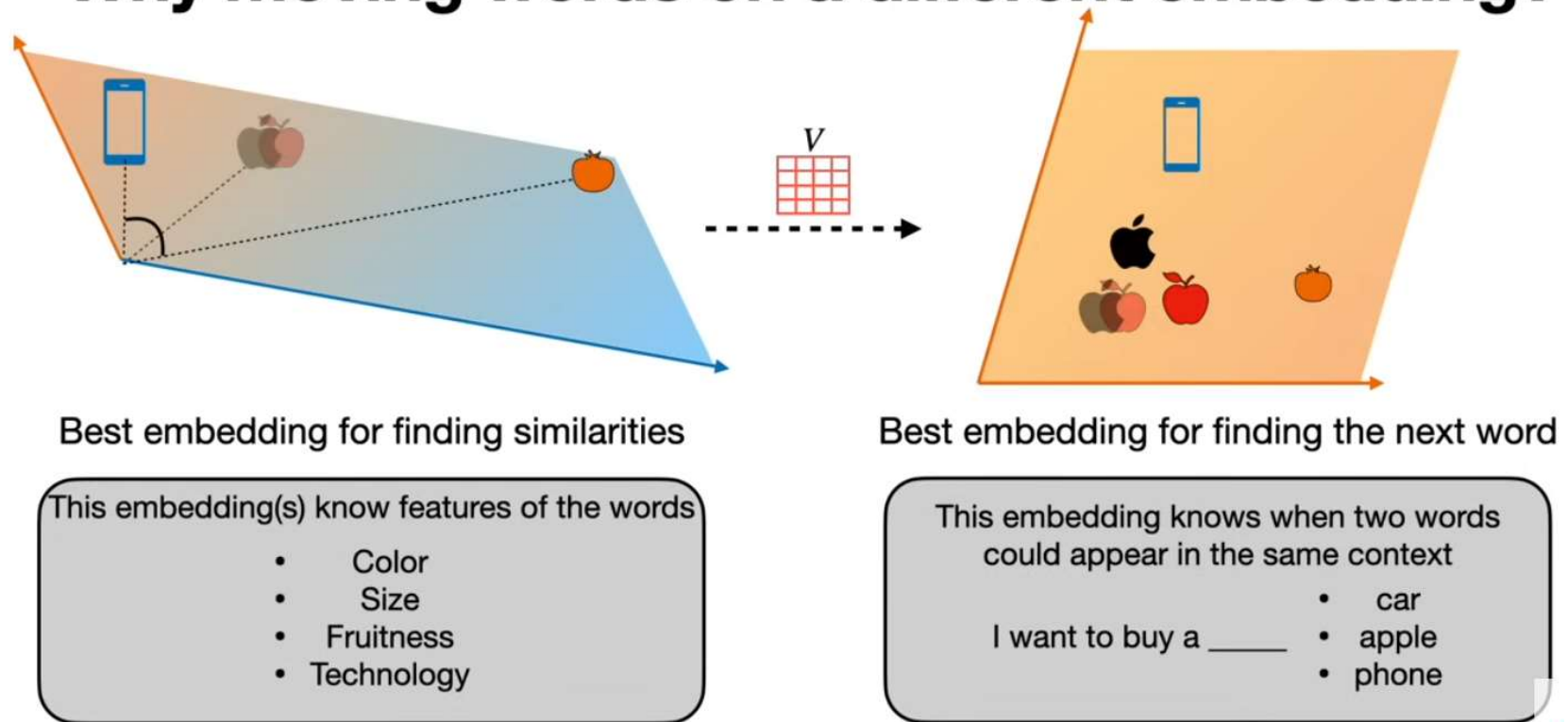
Effects of Linear Transformations

Similarity on a transformed embedding



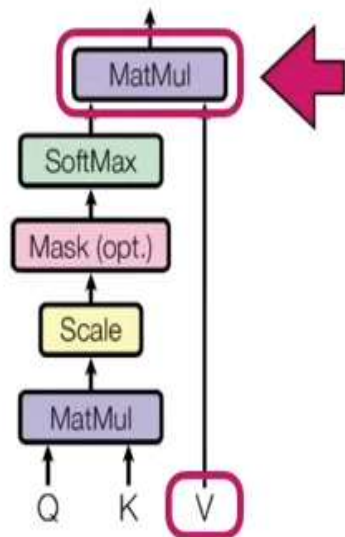
Effects of Value Matrix

Why moving words on a different embedding?

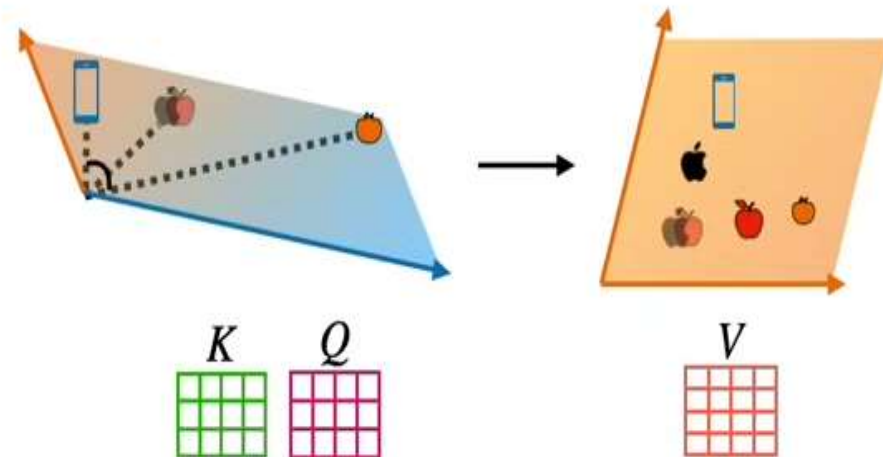


Scaled Dot-Product Attention

Scaled Dot-Product Attention

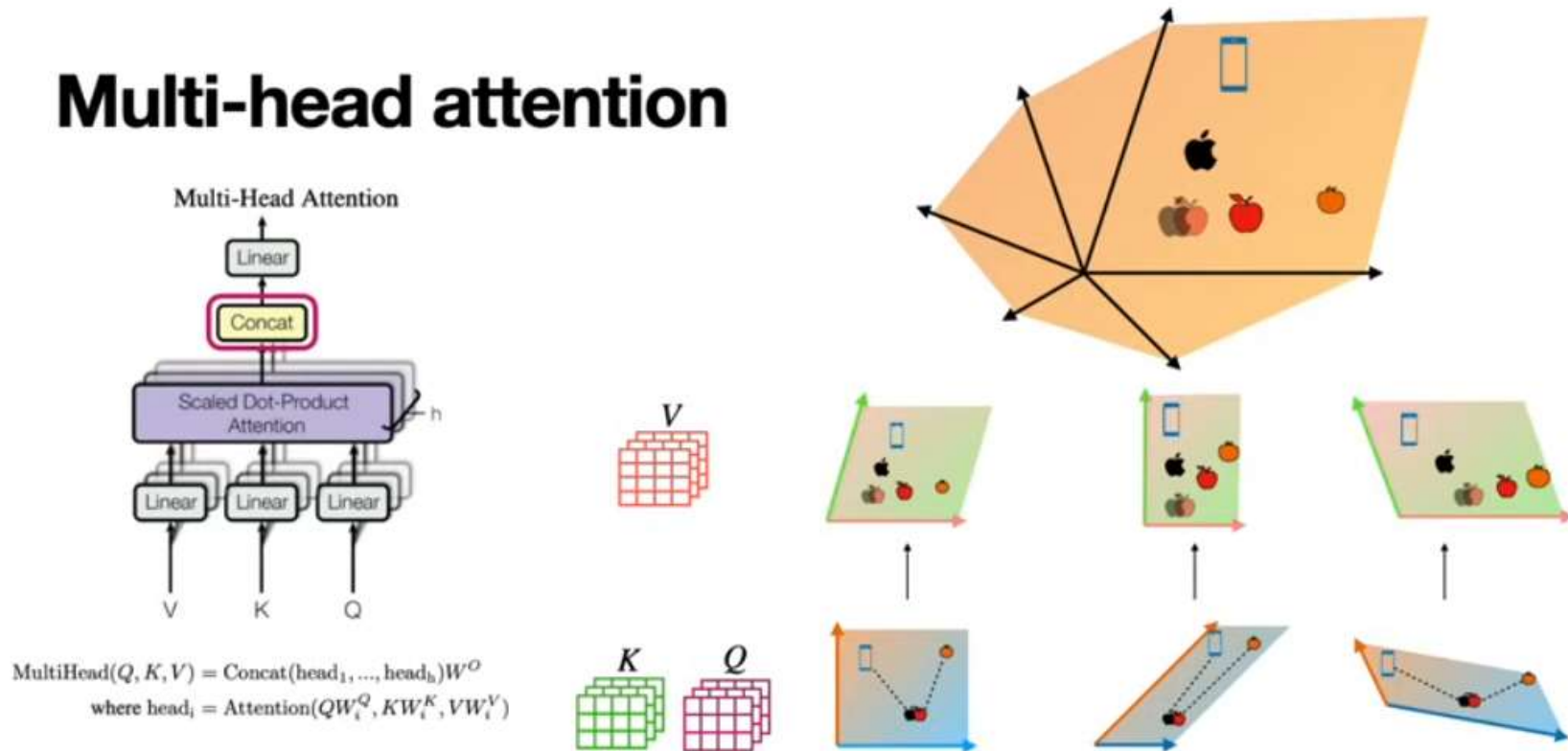


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$



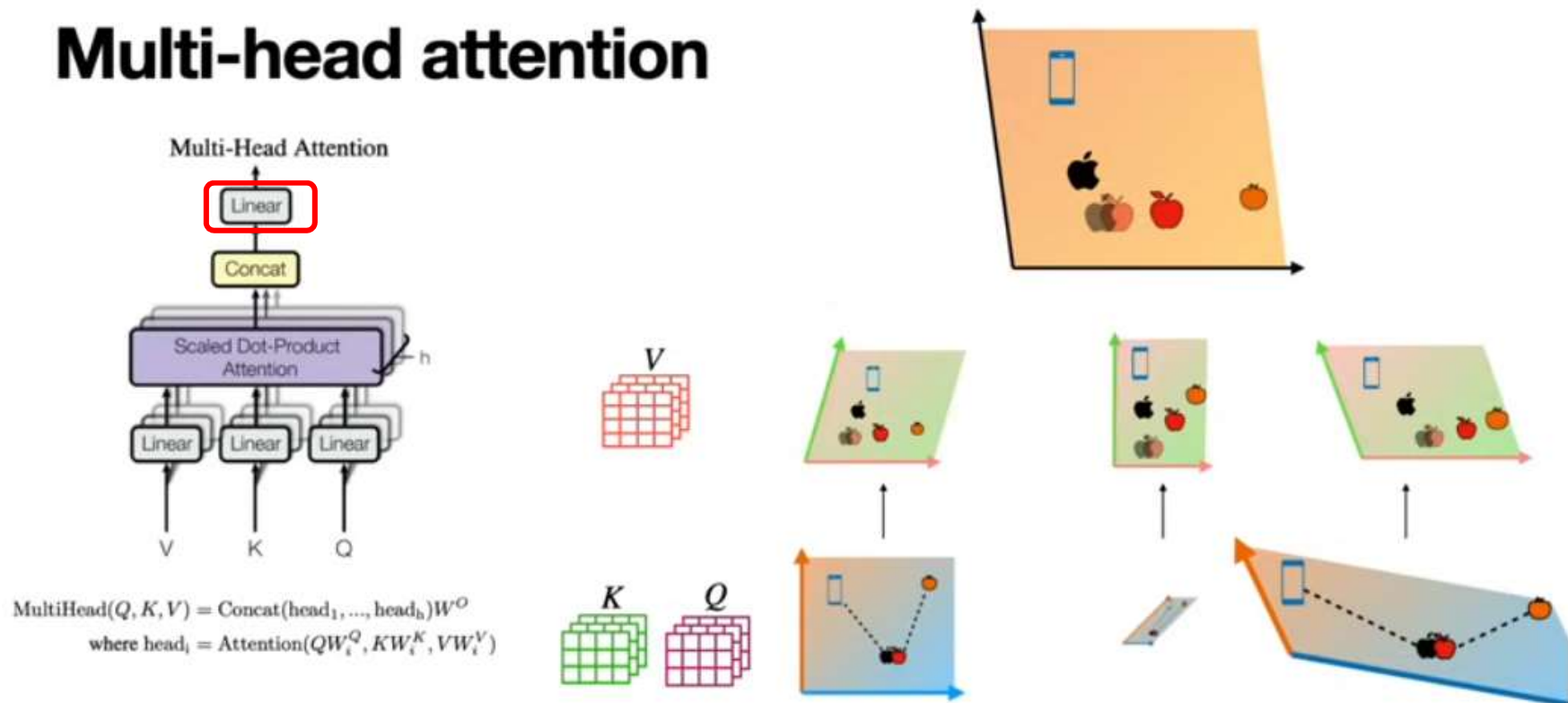
Multi-head Attention (1/3)

Multi-head attention



Multi-head Attention (2/3)

Multi-head attention



Multi-head Attention (3/3)

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



Transformer Architecture (1/2)

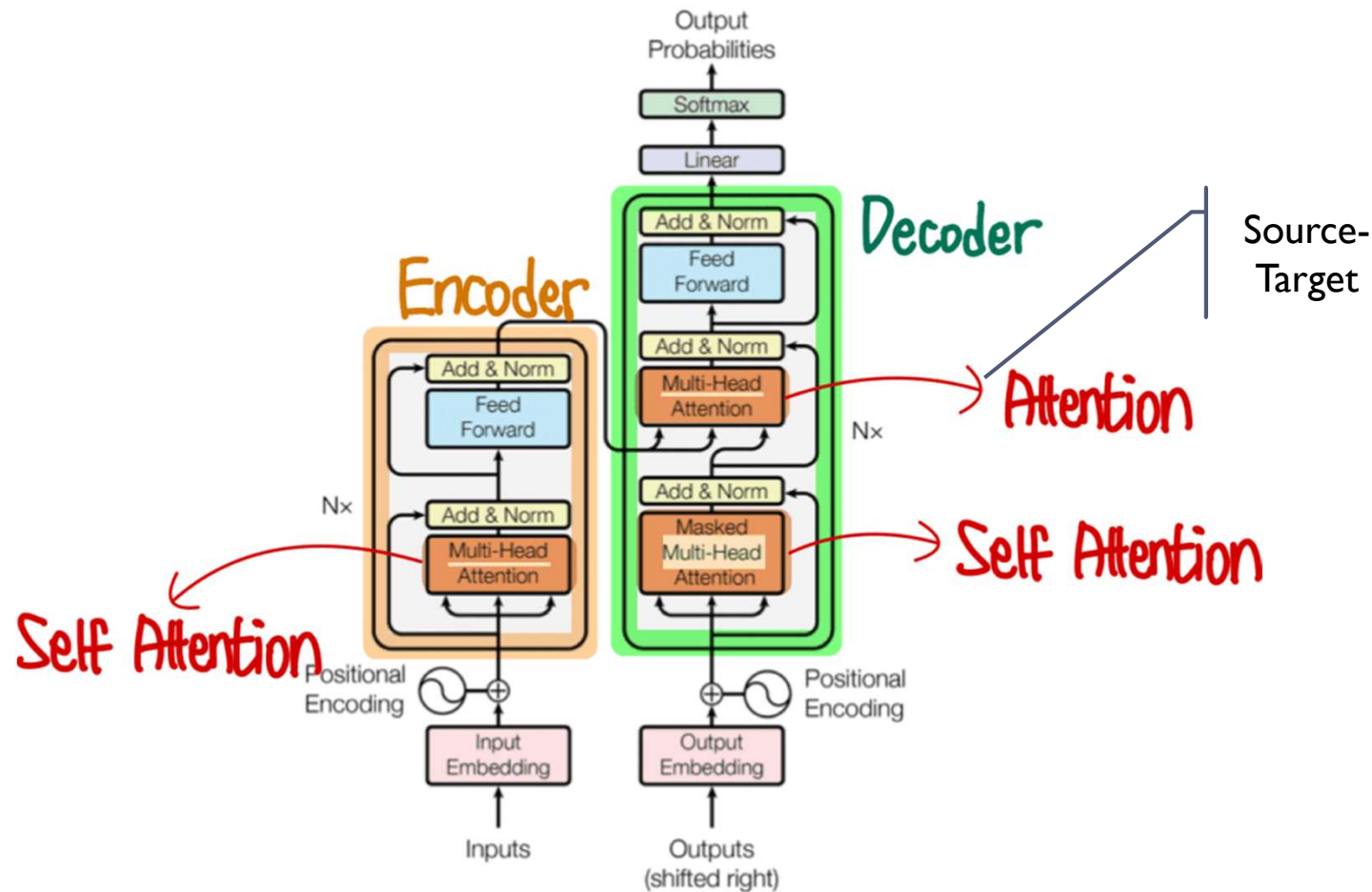
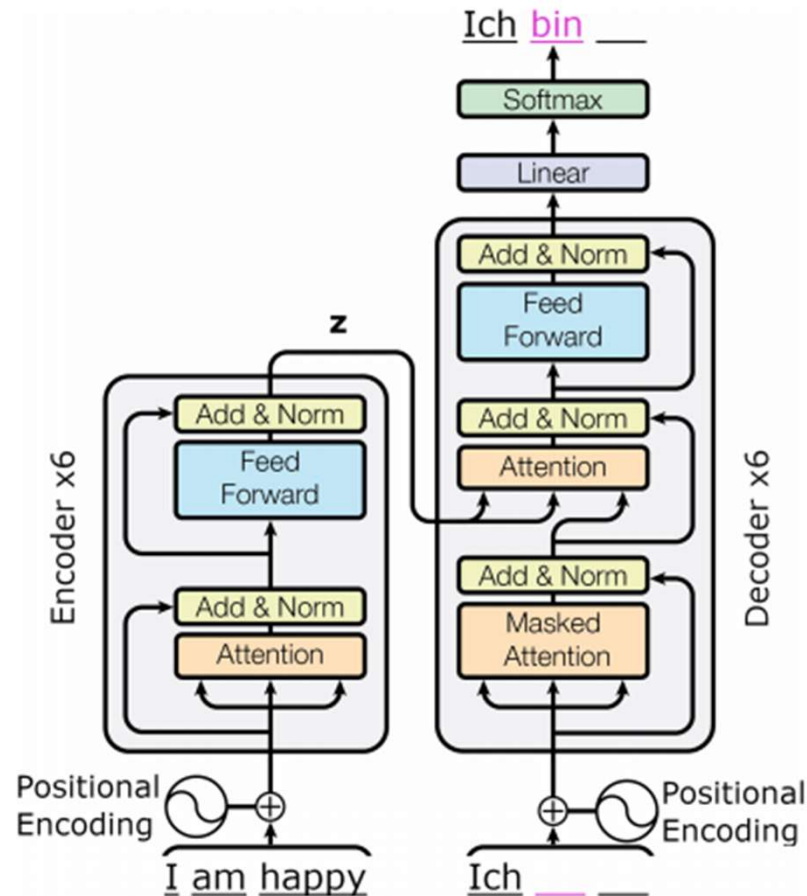


Figure 1: The Transformer - model architecture.

Transformer Architecture (2/2)



Chat-GPT and Transformers

- ▶ Chat-GPT adopted decoder-only transformer architecture
- ▶ Key Value caching in inference

