

Transformer 系列课程

授课老师：AI_distil



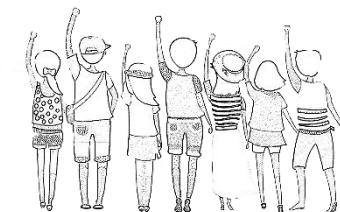
目录

CONTENTS

I 注意力机制

II Q、K、V矩阵

III transformer



TEACHER

PART ONE

注意力机制



$$1+5=\frac{1}{3}$$

A+
8/2

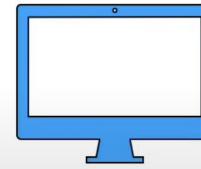
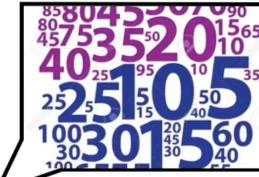


注意力机制



Andy

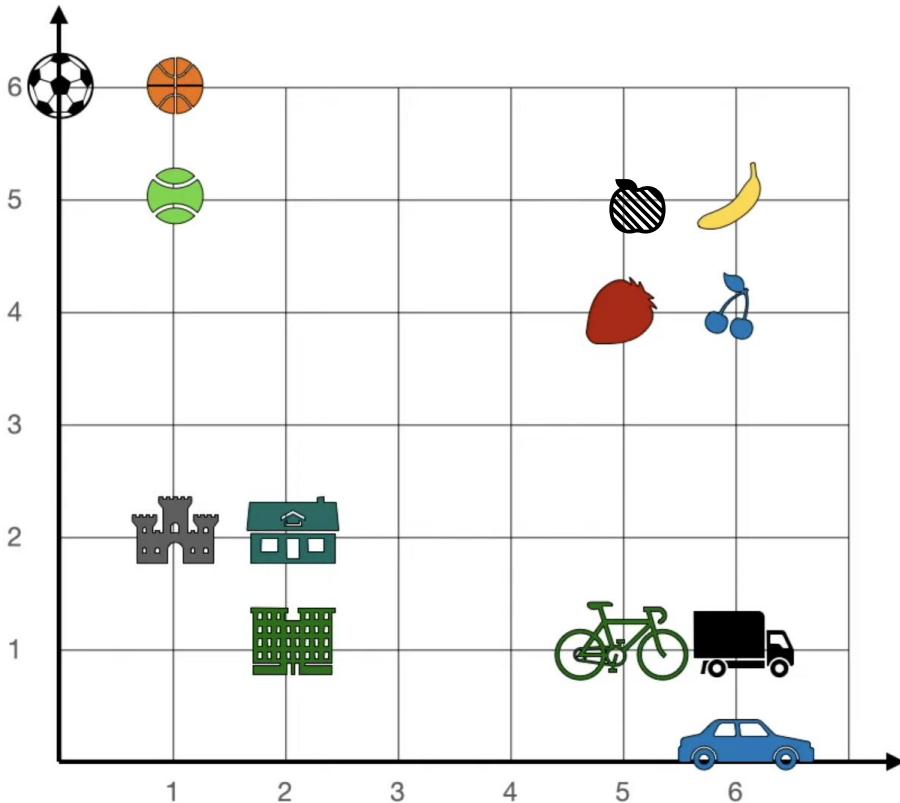
embeddings 即词嵌入，
是NLP领域最重要部分



文本变成数字, embeddings 越好模型就会越好

- sz/

embeddings



Andy - *szi* ?

相似的单词会被赋予相似的数字

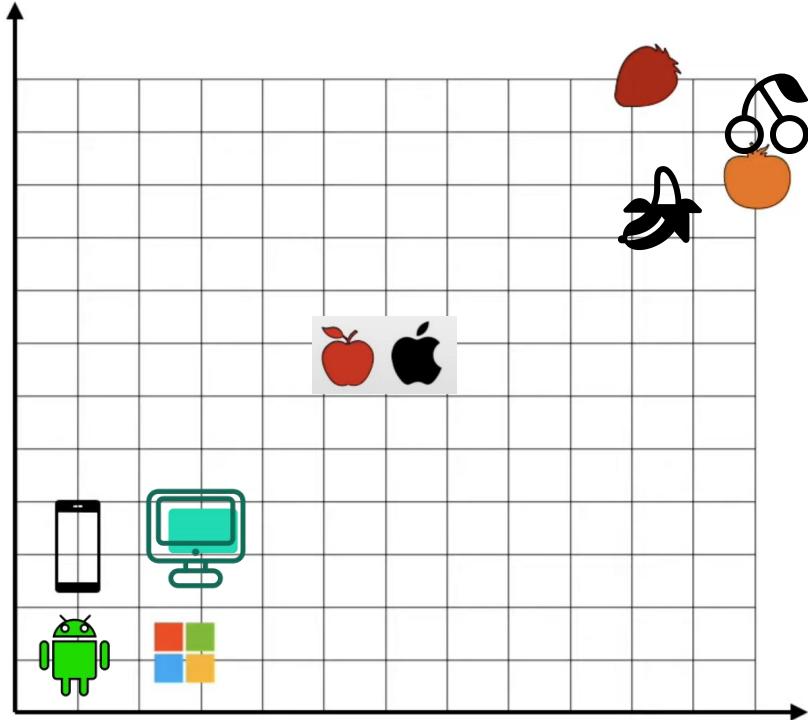
Word	Numbers	
Apple	5	5
Soccer	0	6
House	2	2
Car	6	0



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

4096

problems



Andy_szl

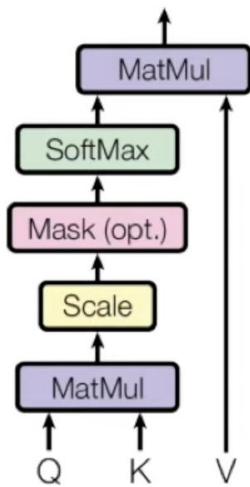
Apple ?

Andy_szl

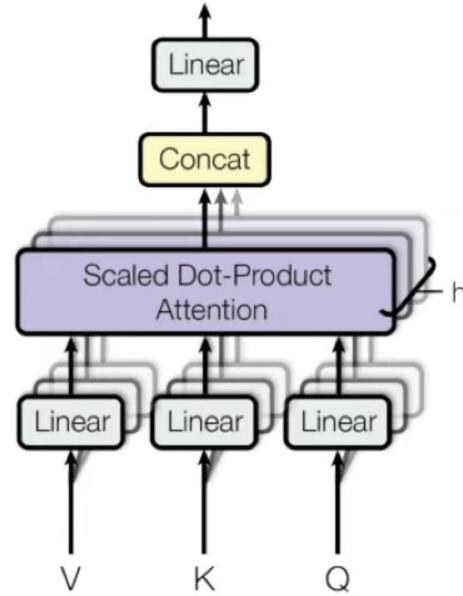


Attention

Scaled Dot-Product Attention



Multi-Head Attention

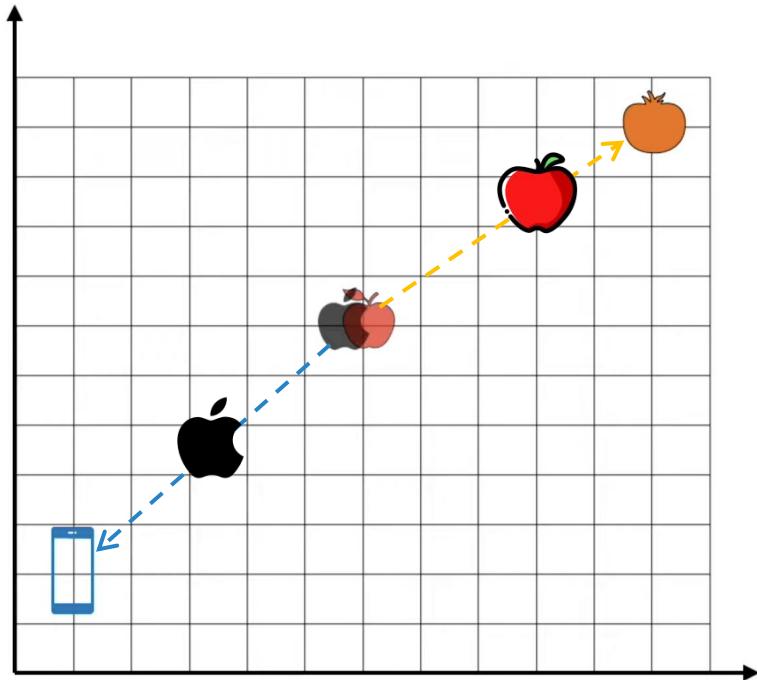


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

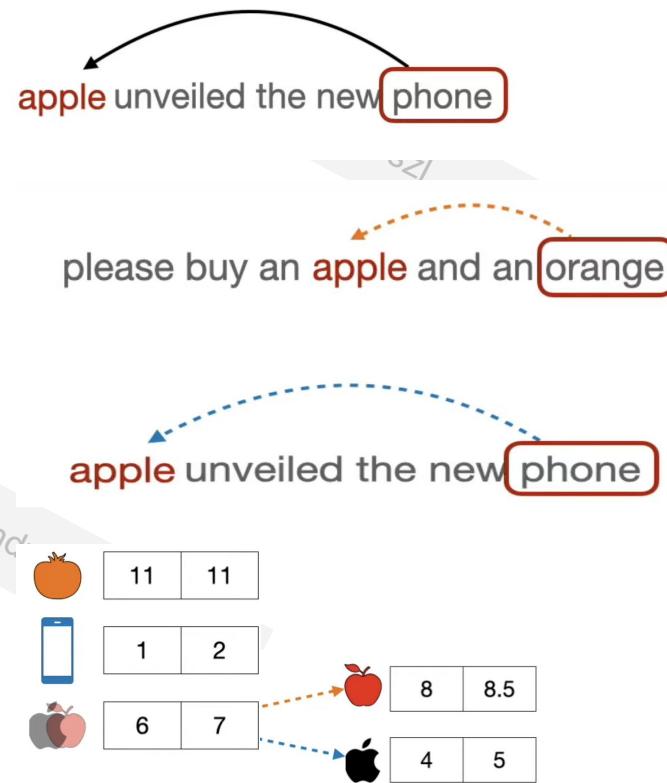
$$\begin{aligned} \text{MultiHead}(Q, K, V) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \\ \text{where } \text{head}_i &= \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \end{aligned}$$



Attention

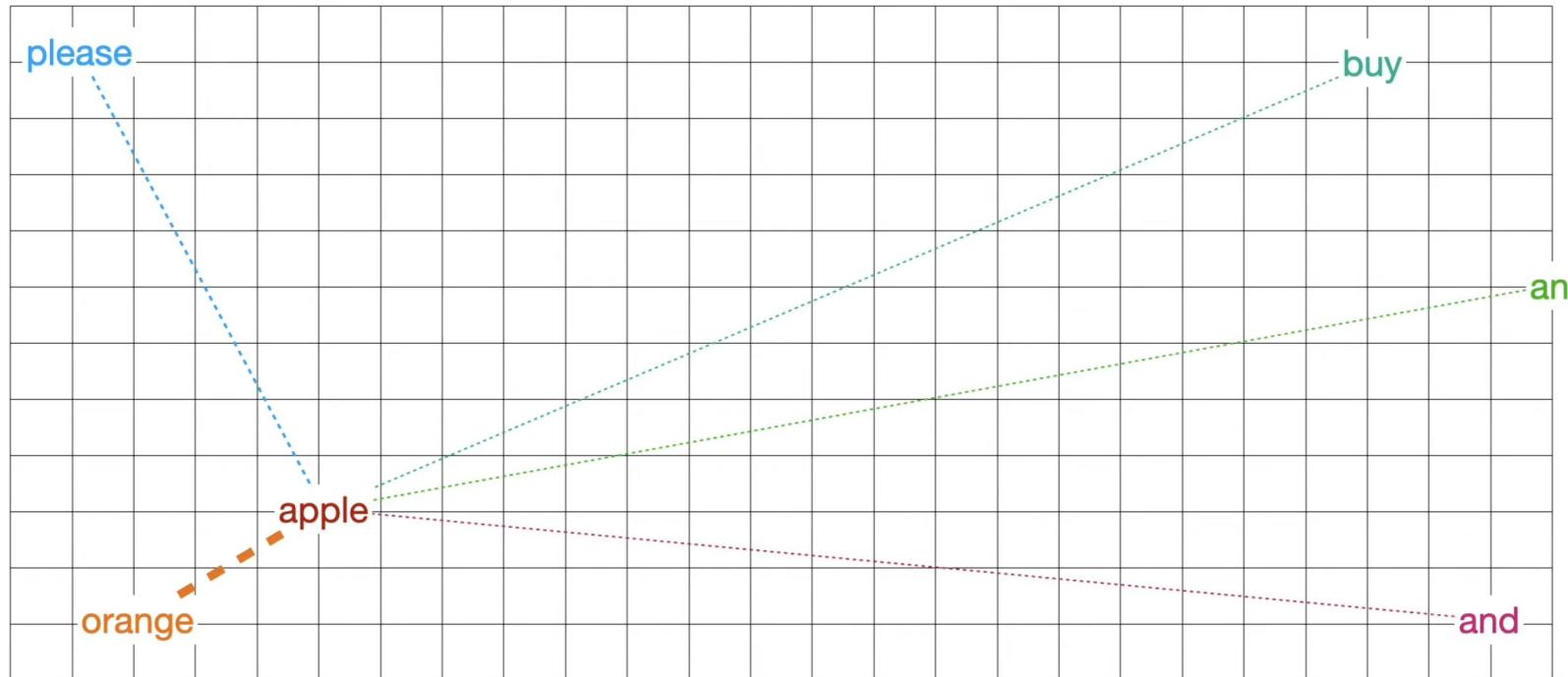


)

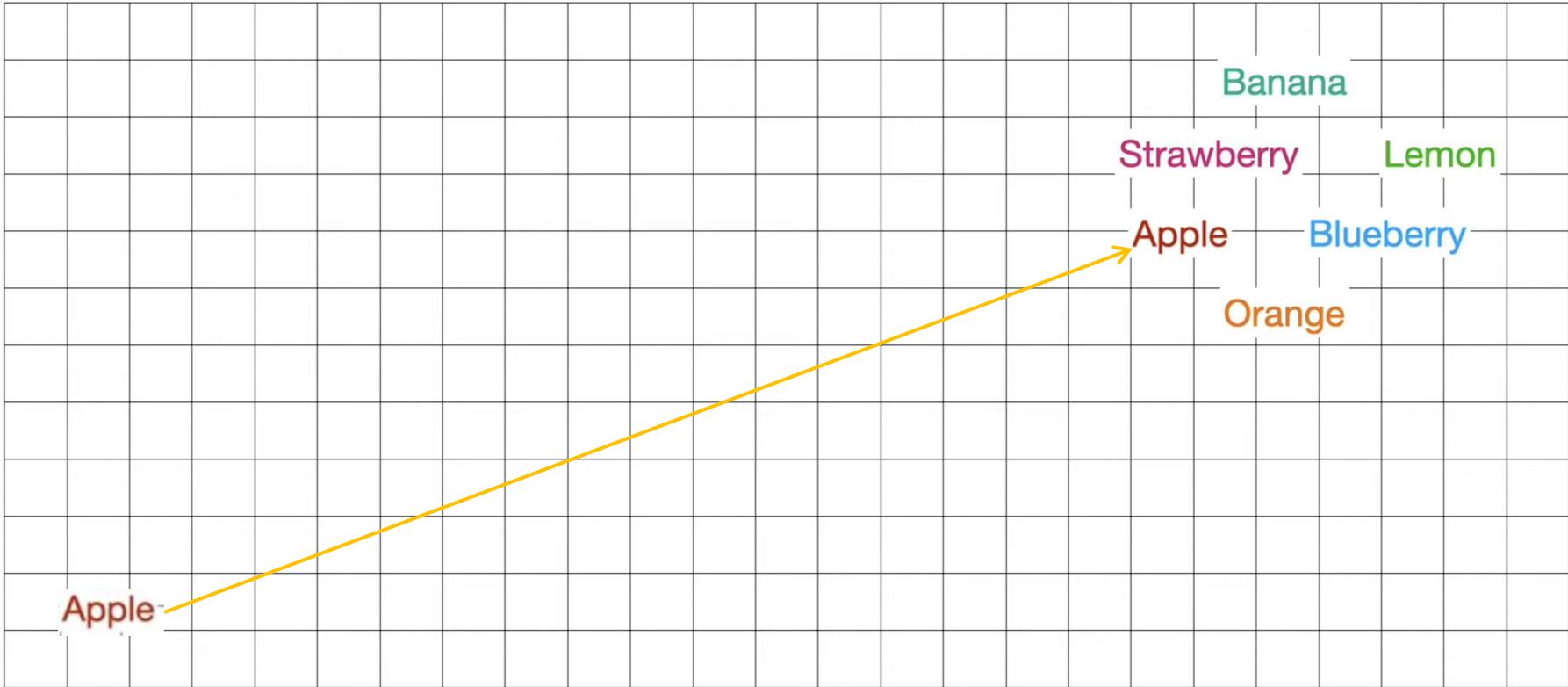


Attention

please buy an apple and an orange

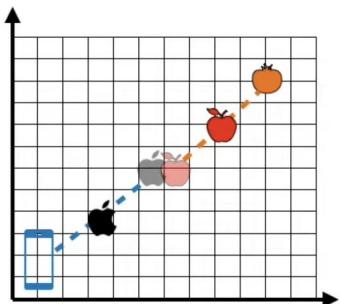
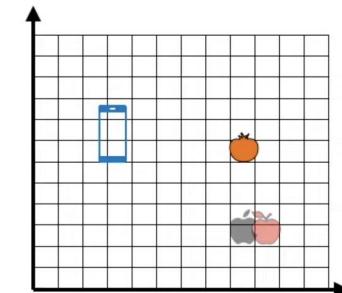
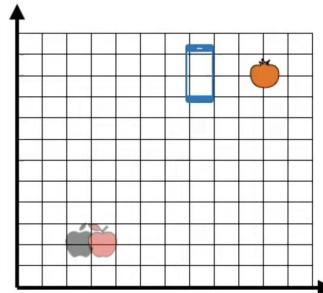
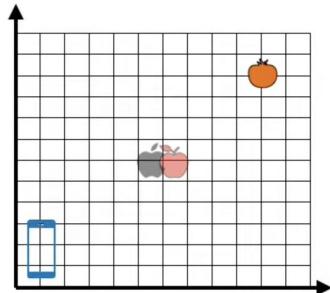


Attention

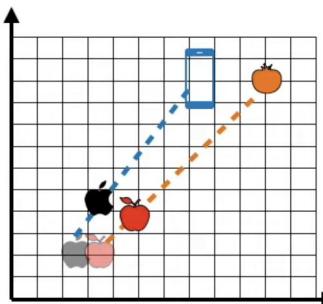




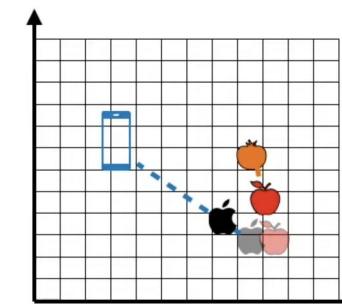
Multi-head attention



Good



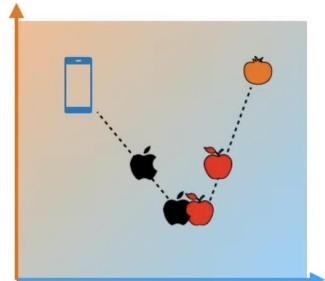
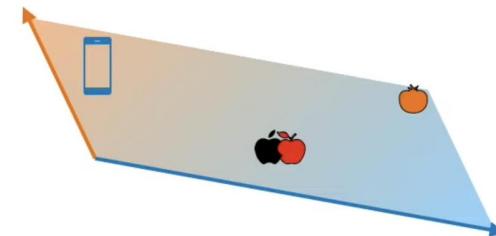
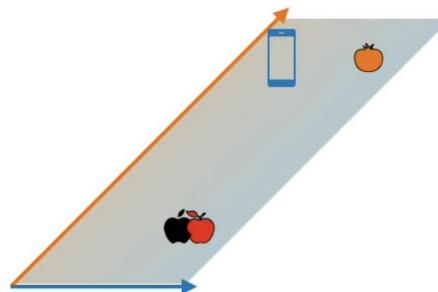
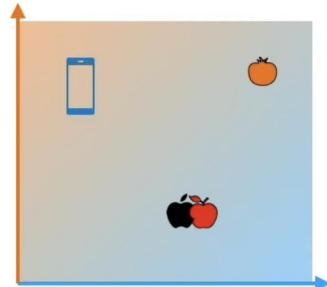
Bad



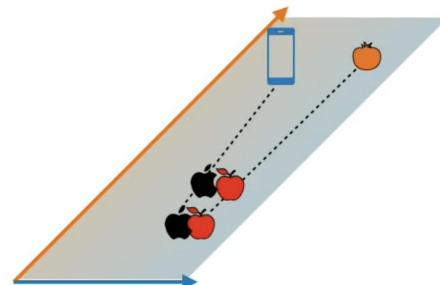
So-so



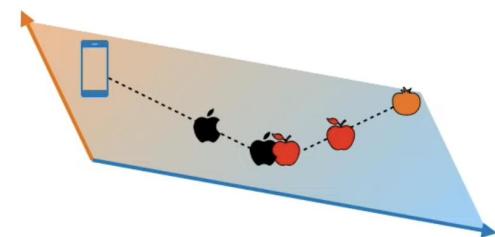
Linear transformations



Okay

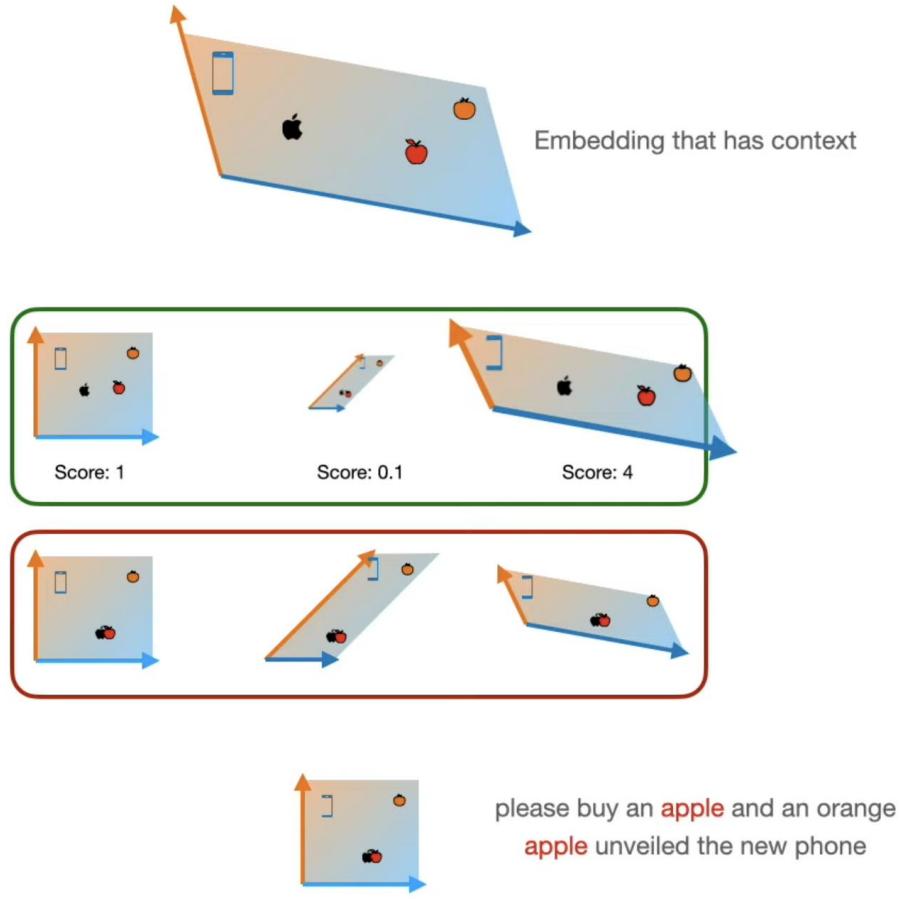
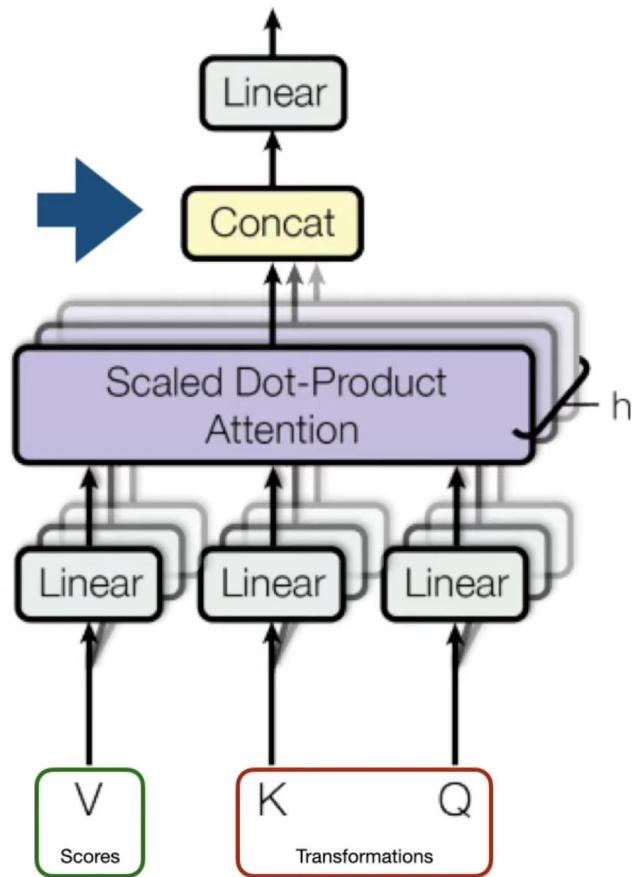


Bad



Good

Multi-Head Attention



TEACHER

PART two

Q、K、V矩阵

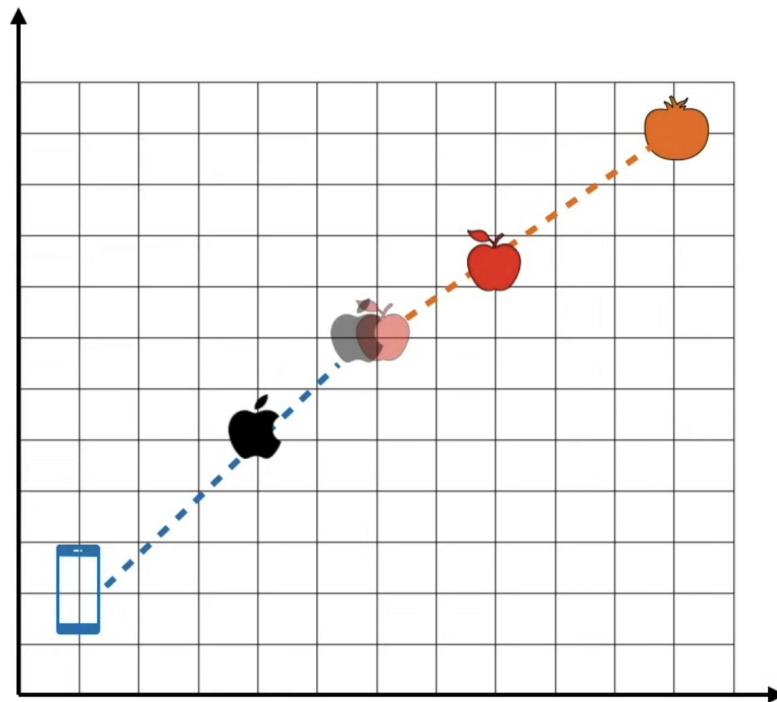
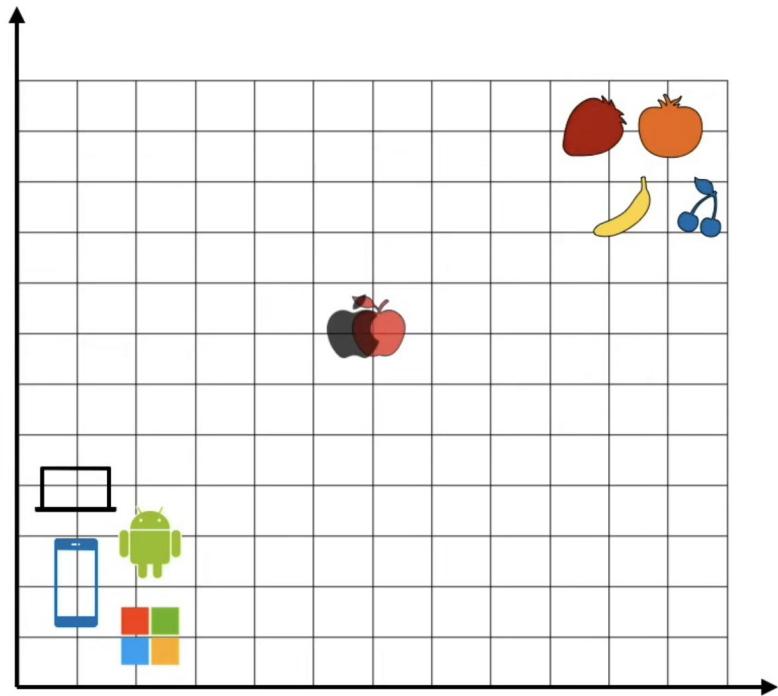


$$1+5=\frac{1}{3}$$

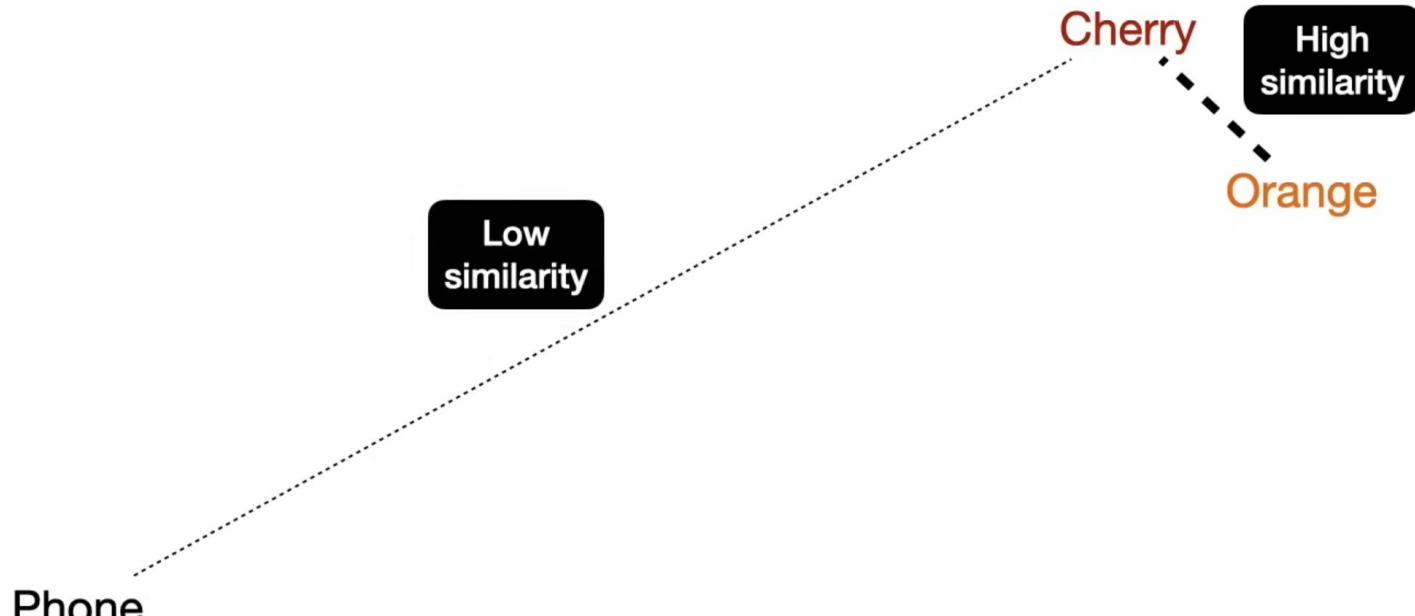
$$A+\frac{1}{2}$$



词义的相似性



similarity





1、点积



Andy - size

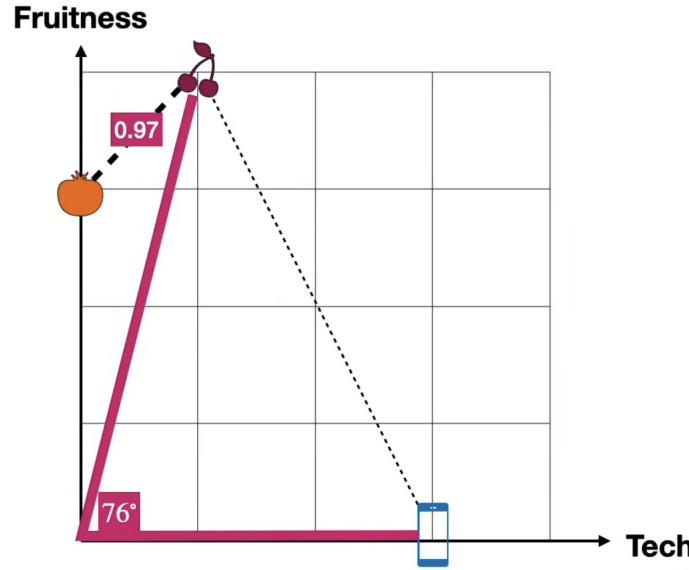
Tech	Fruitness
1	4
0	3
1	4
3	0
0	3
3	0

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

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

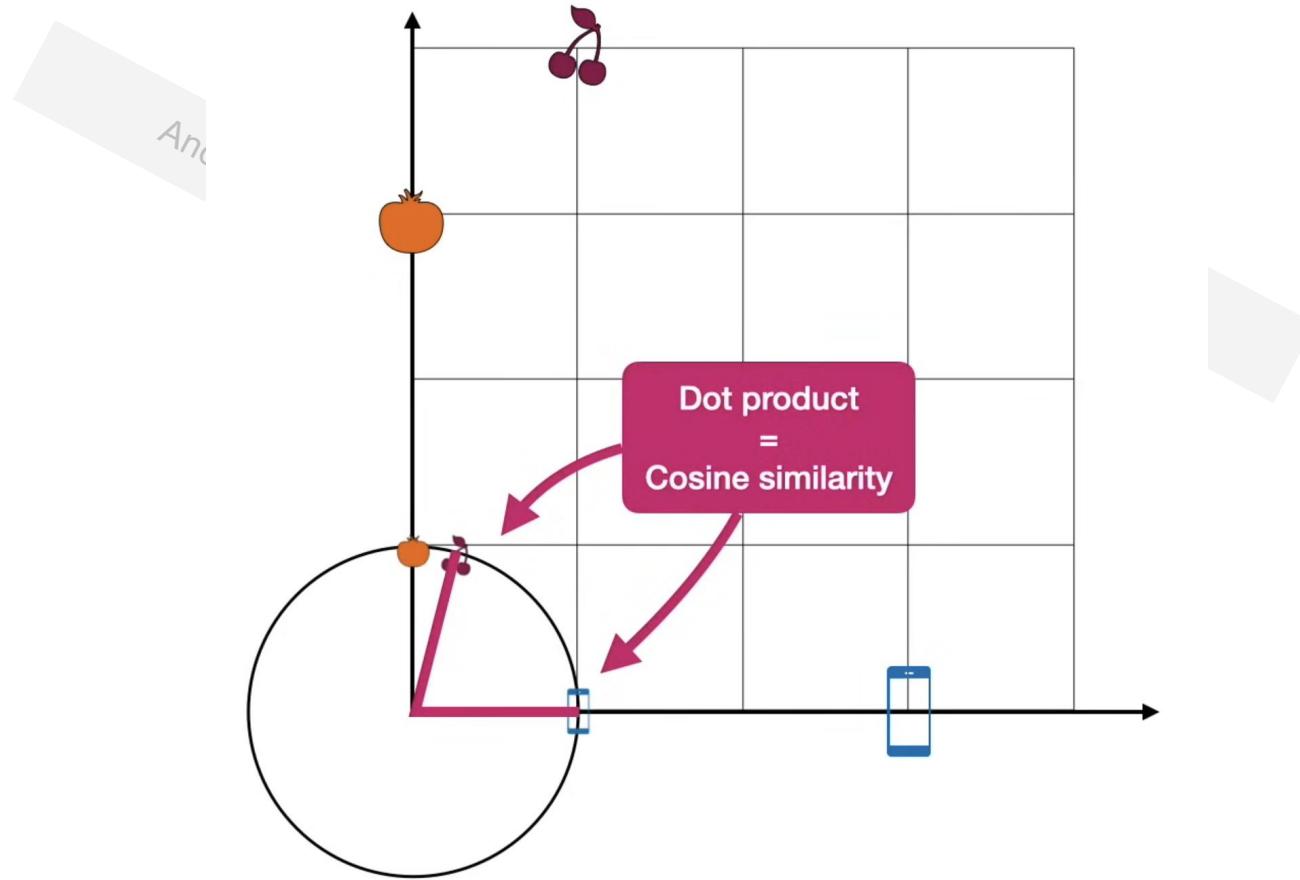
$0 \cdot 3 + 3 \cdot 0 = 0$

2、余弦相似度



- $\cos(14^\circ) = 0.97$
- $\cos(76^\circ) = 0.24$
- $\cos(90^\circ) = 0$

点积和余弦相似度





比例点积

Sim



1	4
---	---



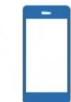
0	3
---	---

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

Sim



1	4
---	---



3	0
---	---

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

Sim



0	3
---	---

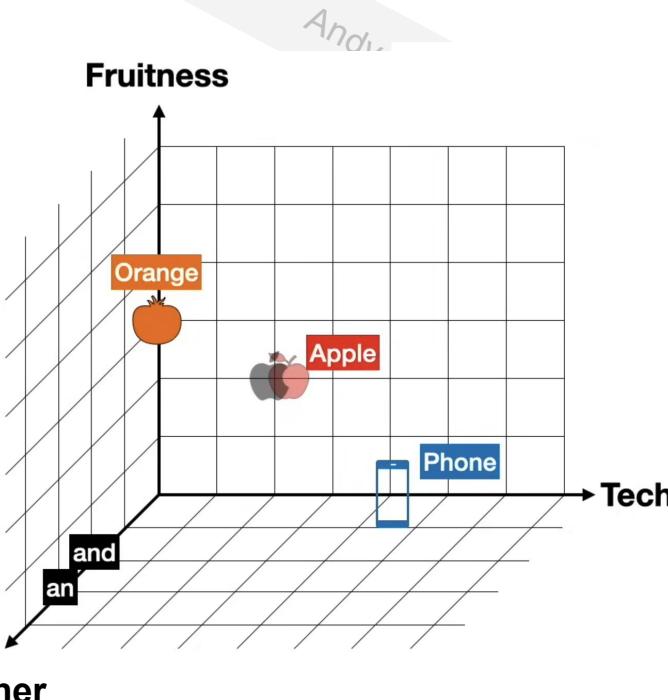


3	0
---	---

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

例子

>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

计算过程

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

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

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

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

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

计算过程

an apple phone

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

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

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

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

归一化

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



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



soft Max

$$x \rightarrow e^x$$

Andy
Sawyer

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

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



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

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

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

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

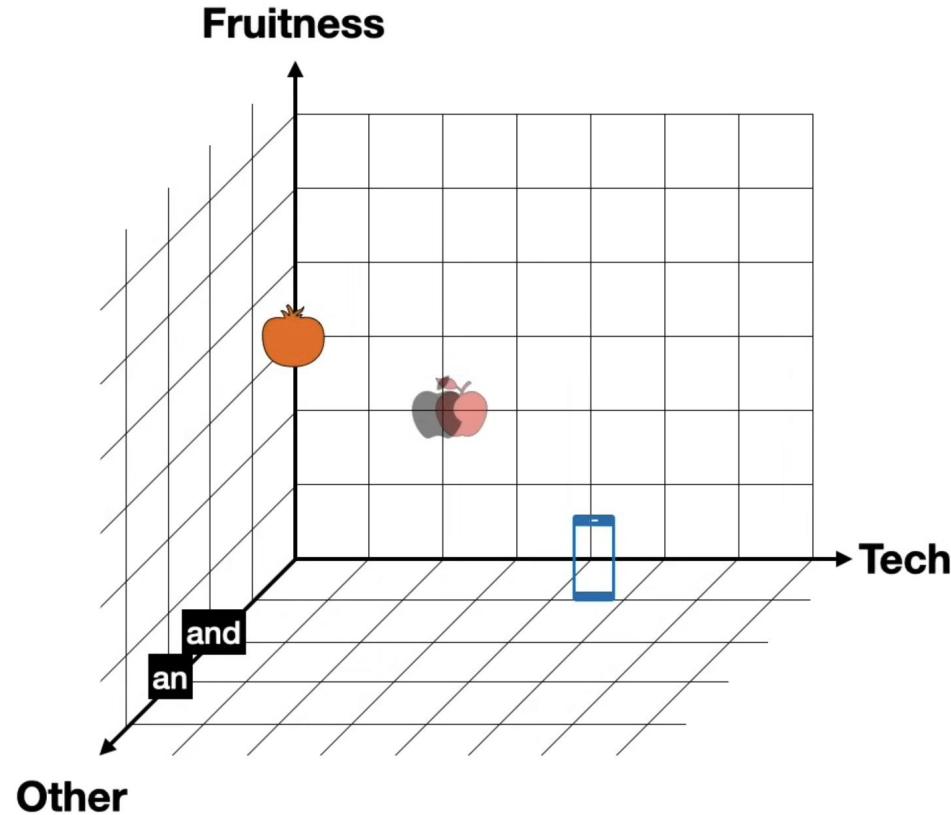
真实结果

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

Orange → 0.4 Orange + 0.3 Apple + 0.15 And + 0.15 An



移动过程



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

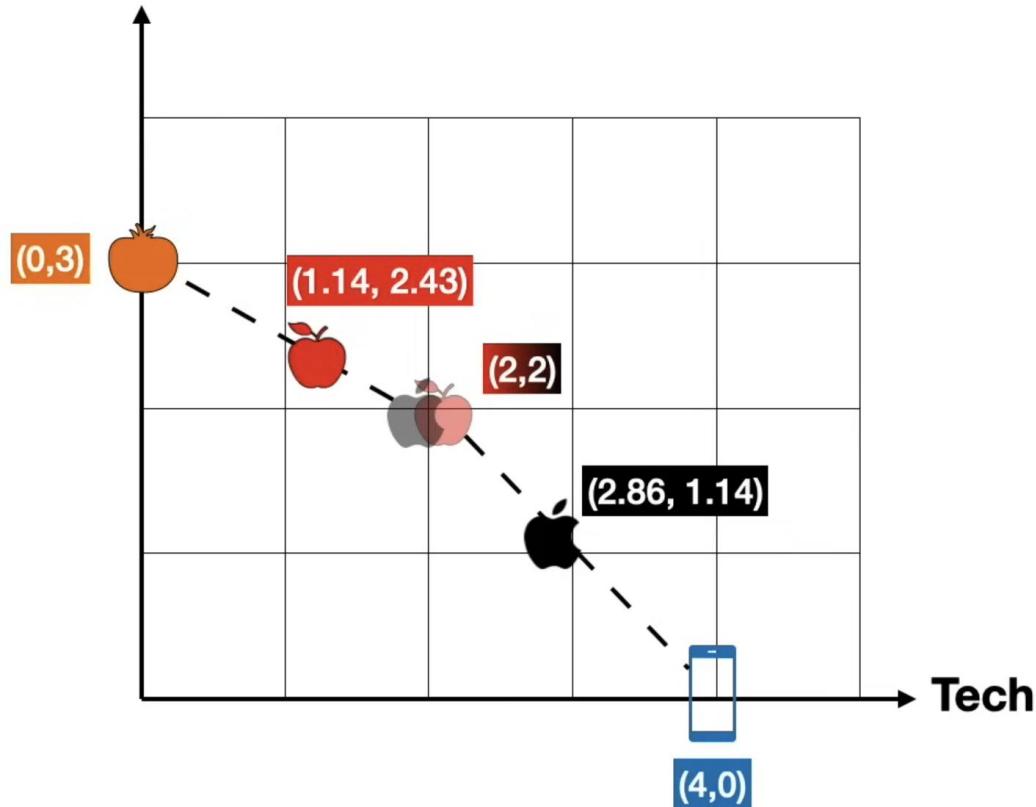
$$\text{Apple} \rightarrow 0.43 \text{ Orange} + 0.57 \text{ Apple}$$

an **apple** phone

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

移动过程

Fruitness



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

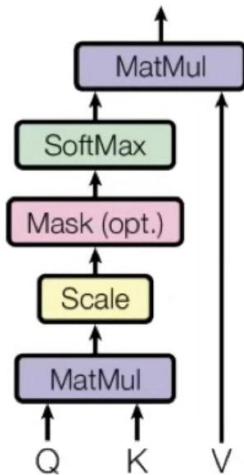
$$\text{Apple} \rightarrow 0.43 \text{ Orange} + 0.57 \text{ Apple}$$

an **apple** phone

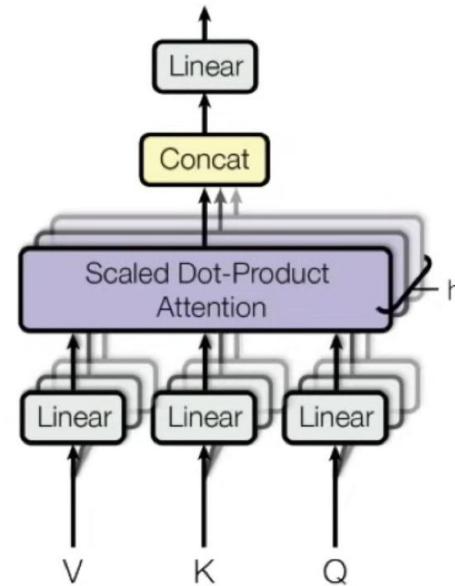
$$\text{Apple} \rightarrow 0.43 \text{ Phone} + 0.57 \text{ Apple}$$

Q, K, V

Scaled Dot-Product Attention



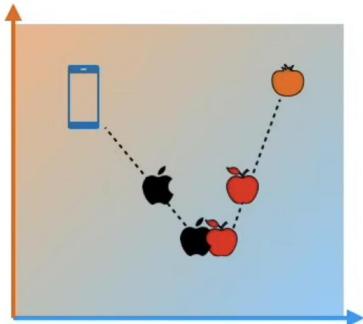
Multi-Head Attention



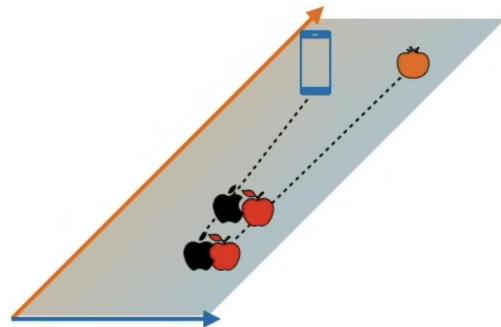
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

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

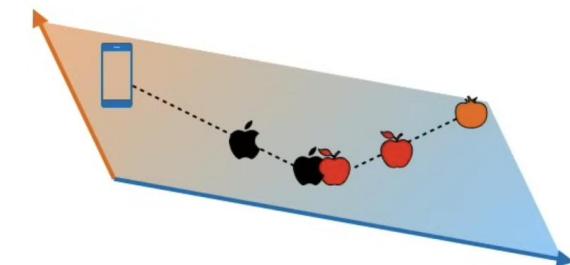
K, Q



Okay

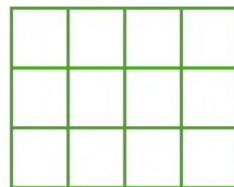


Bad

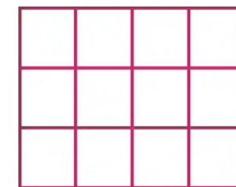


Good

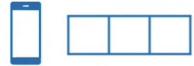
Keys



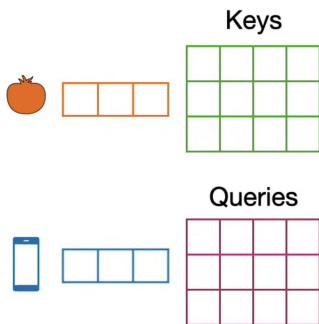
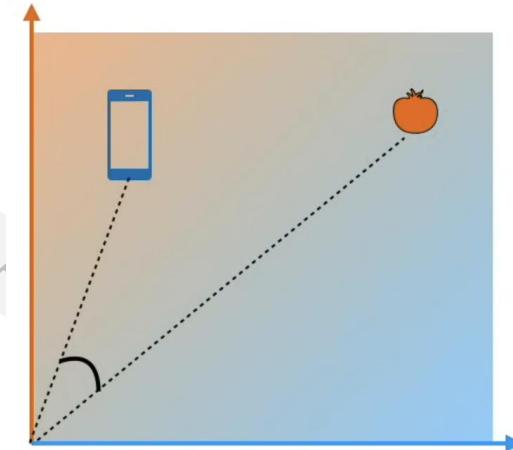
Queries



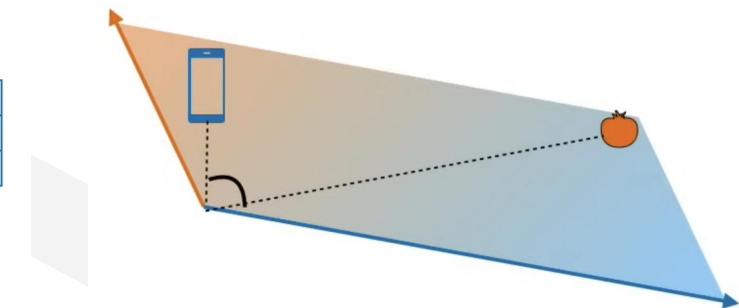
相似度



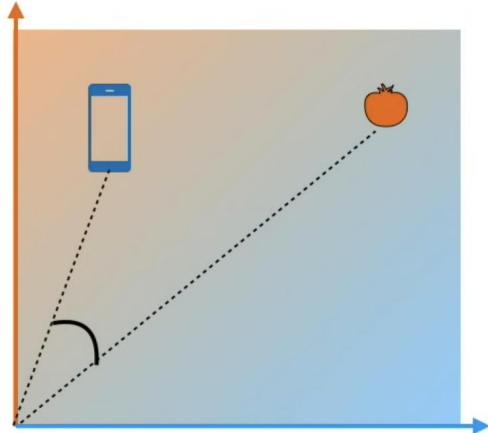
$$\text{Similarity}(\text{Tomato}, \text{Smartphone}) = \begin{matrix} \text{Tomato} \\ \text{Smartphone} \end{matrix} \cdot \begin{matrix} \text{Tomato} \\ \text{Smartphone} \end{matrix}$$



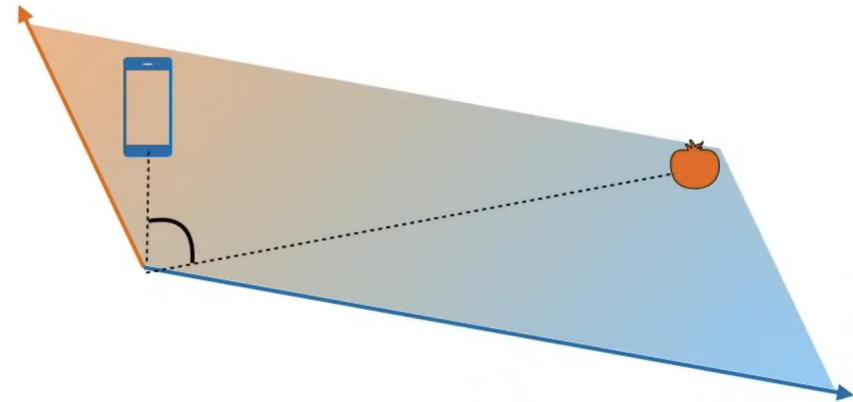
$$\text{Similarity}(\text{Tomato}, \text{Smartphone}) = \begin{matrix} \text{Tomato} \\ \text{Smartphone} \end{matrix} \cdot \begin{matrix} \text{Tomato} \\ \text{Smartphone} \end{matrix}$$



相似性的修改

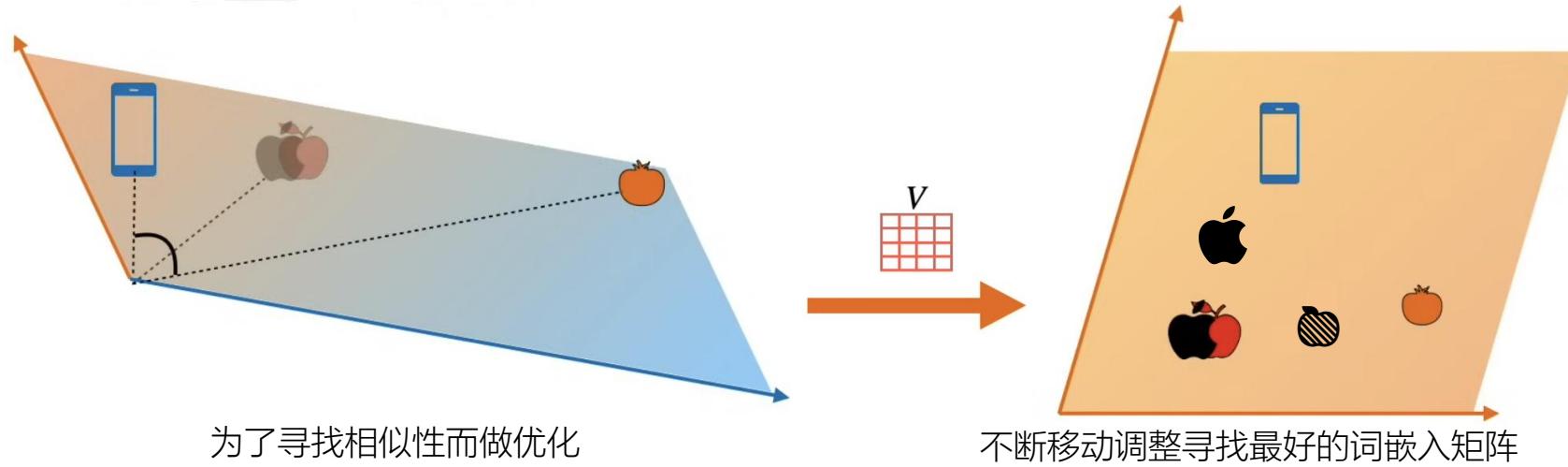


$$\text{Similarity}(\text{apple}, \text{phone}) = \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline\end{array} \quad \begin{array}{|c|c|}\hline & \\ \hline & \\ \hline\end{array}$$



$$\text{Similarity}(\text{apple}, \text{phone}) = \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline\end{array} \quad \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline\end{array} \quad \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline\end{array} \quad \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline\end{array} \quad \begin{array}{|c|c|c|}\hline & & \\ \hline & & \\ \hline & & \\ \hline\end{array}$$

值矩阵



Keys

Queries

$Q, K \Rightarrow \text{embedding}$:
颜色、大小、气味.....

Values



实现过程

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

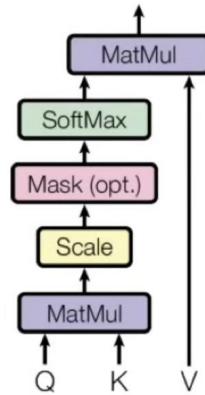
	Orange	Apple	And	An
Orange	v_{11}	v_{12}	v_{13}	v_{14}
Apple	v_{21}	v_{22}	v_{23}	v_{24}
And	v_{31}	v_{32}	v_{33}	v_{34}
An	v_{41}	v_{42}	v_{43}	v_{44}

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

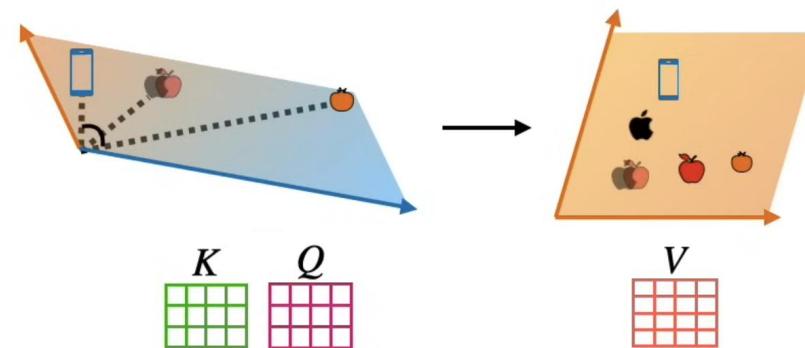
$$\begin{aligned}
 \text{apple} \longrightarrow & v_{21} \cdot \text{orange} \\
 & + v_{22} \cdot \text{apple} \\
 & + v_{23} \cdot \text{and} \\
 & + v_{24} \cdot \text{an}
 \end{aligned}$$

自注意力机制

Scaled Dot-Product Attention

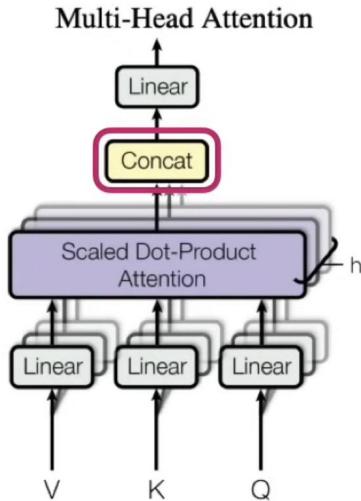


$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



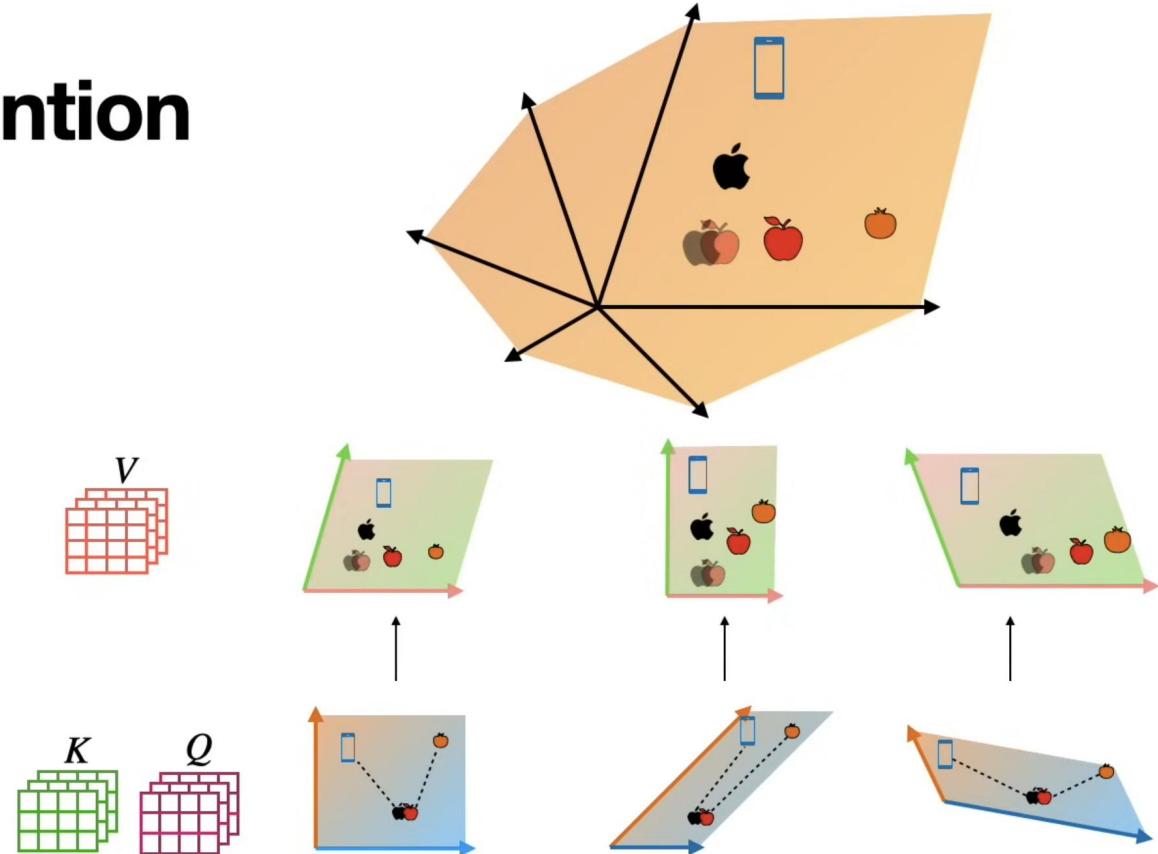
多头注意力

Multi-head attention



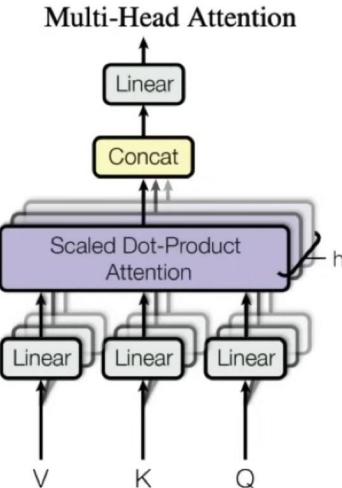
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

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



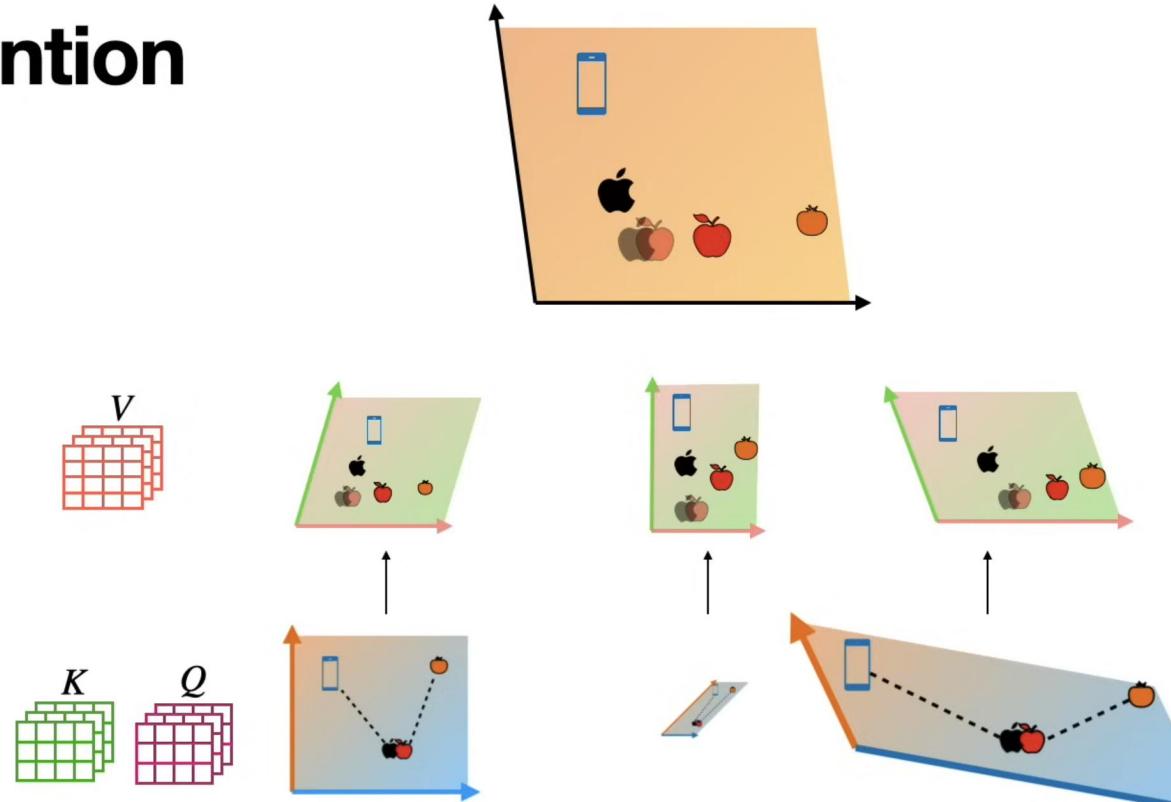
多头注意力

Multi-head attention



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

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



TEACHER

PART three

Transformer



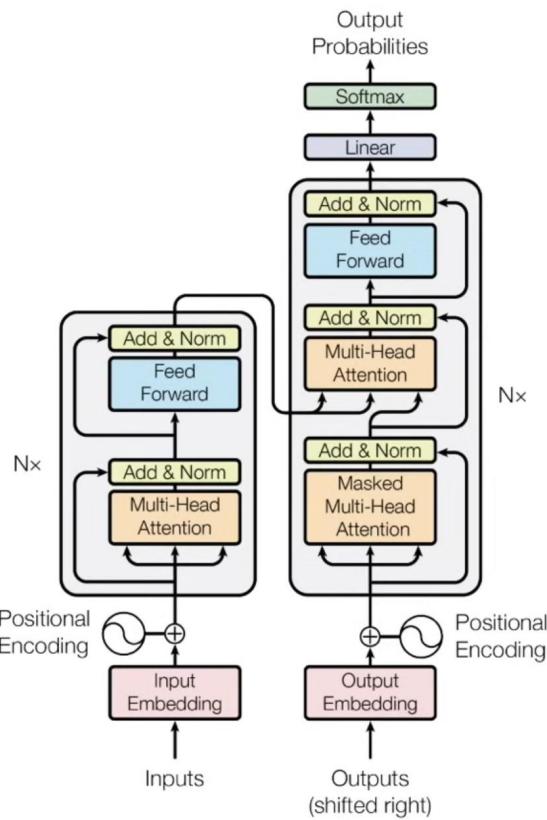
$$1 + 5 = \frac{1}{3}$$

$$A^+$$

$$\frac{9}{12}$$



transformer



 one by one

最近怎____

Andy - szl

怎样

怎会

怎么样

怎能

1-Gram

最近怎____

怎样

怎么样

3-Gram

.....

神经网络

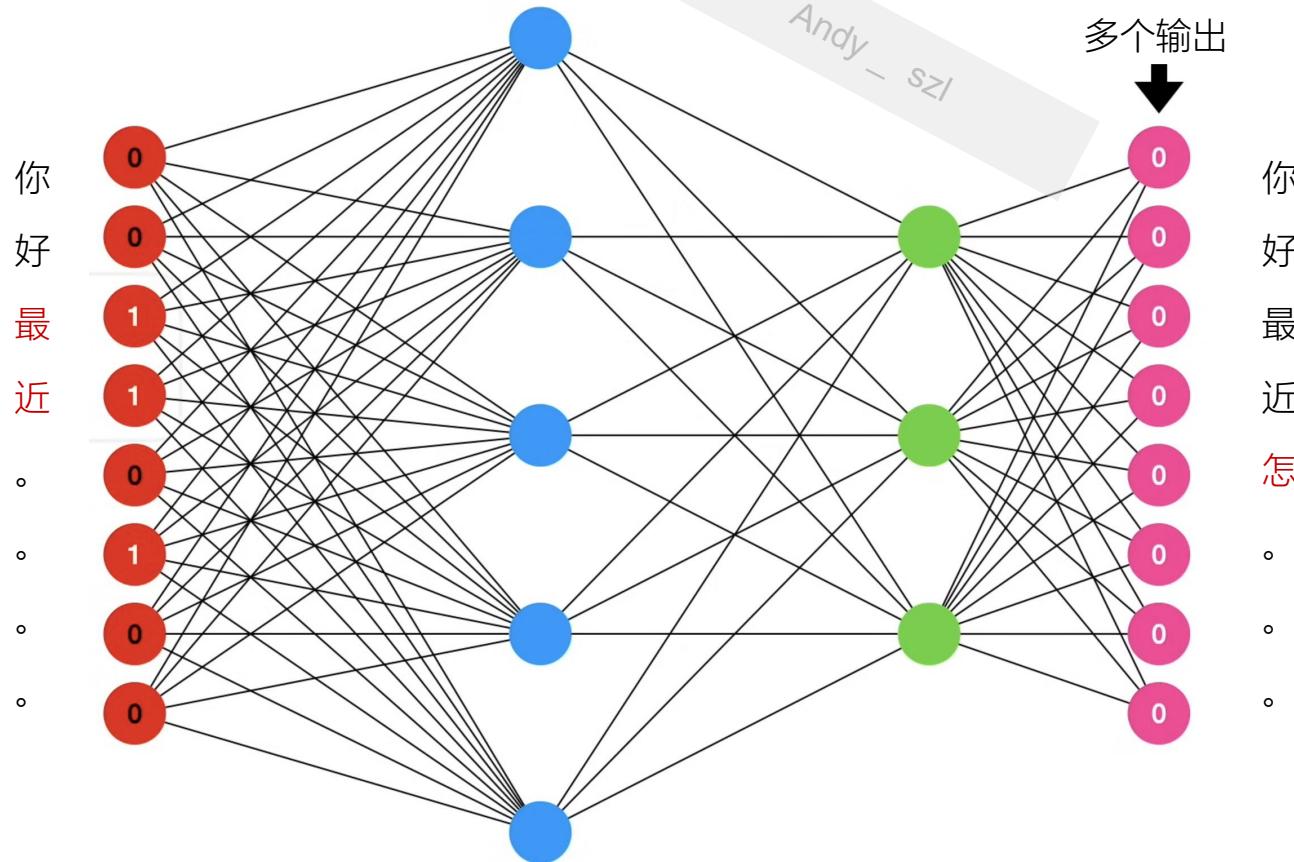


最近怎 → 神经网络 → 么

Andy - szl

Andy - szl

神经网络

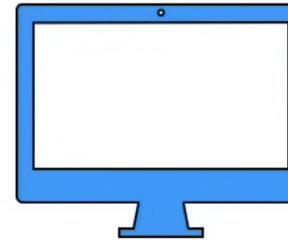
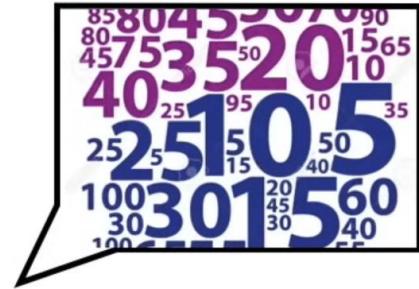
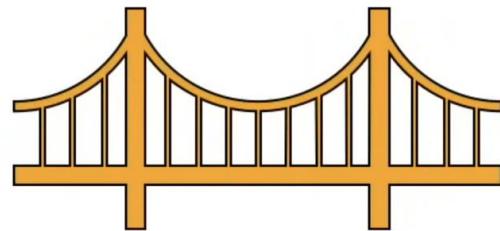


流程: tokenization

Write a story. → Write A story .

你最近怎么样 → 你 最近 怎么 样 ?

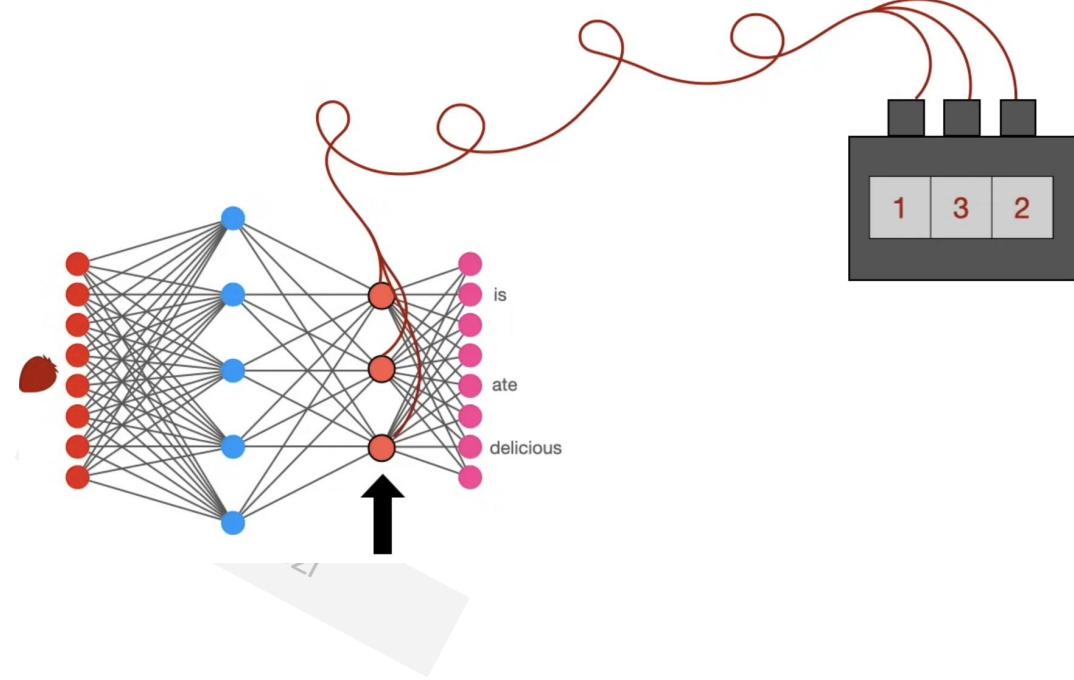
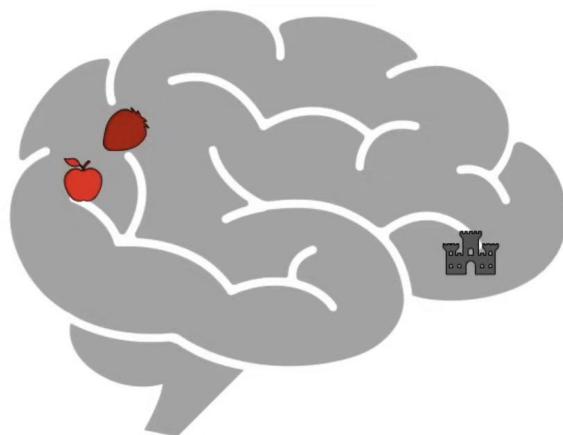
Embedding



相似性高的词会给予更加相近的向量

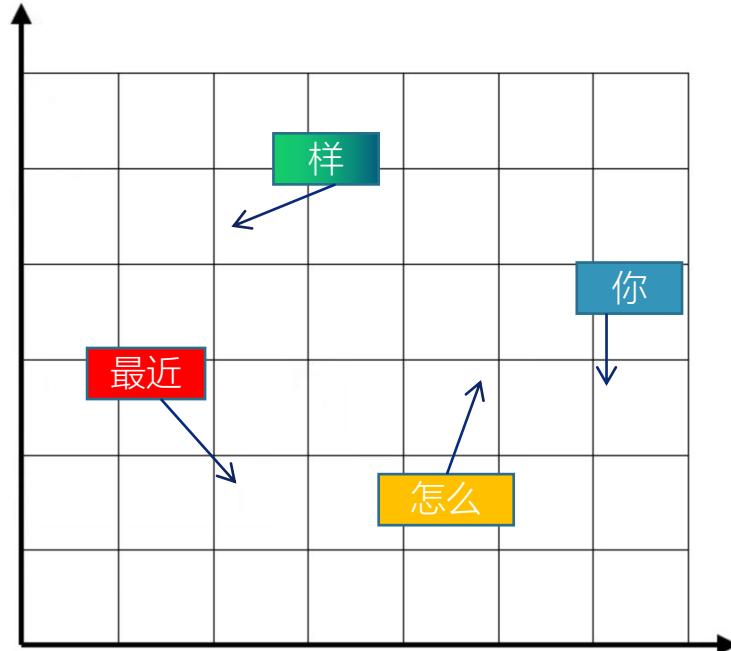
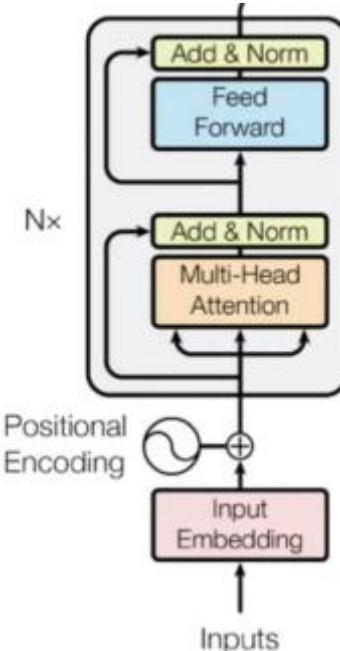
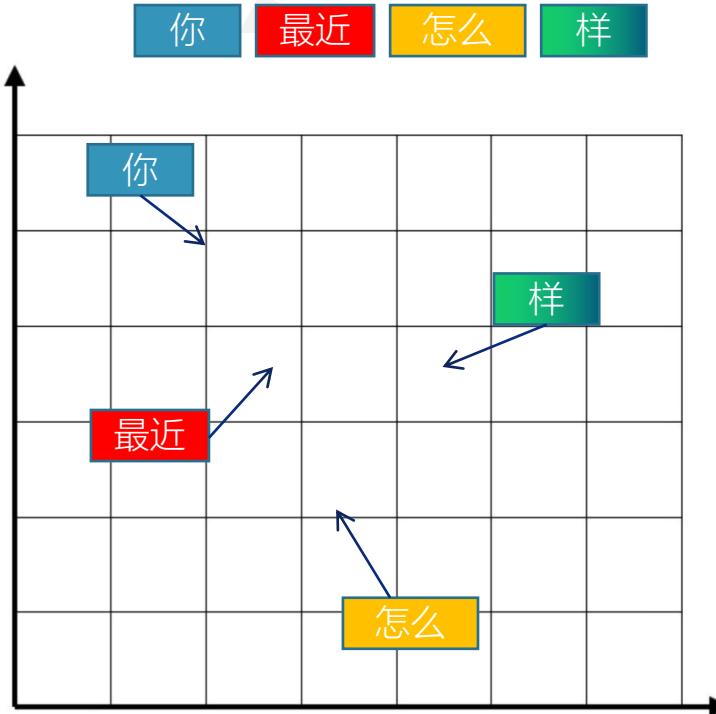


如何创建embedding

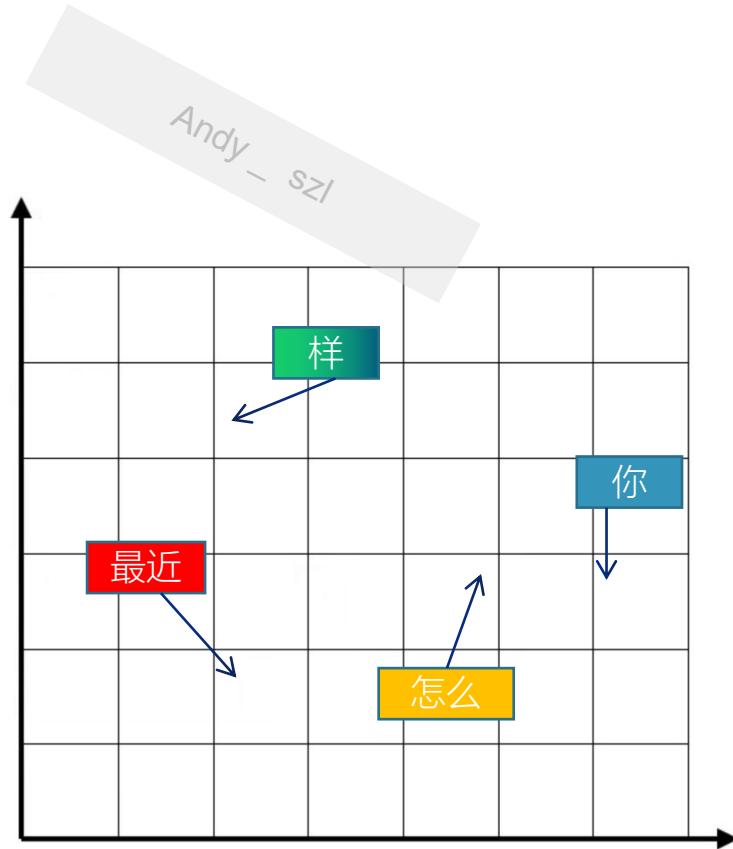
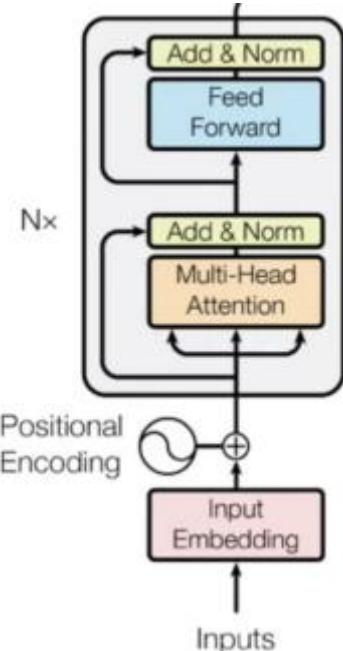
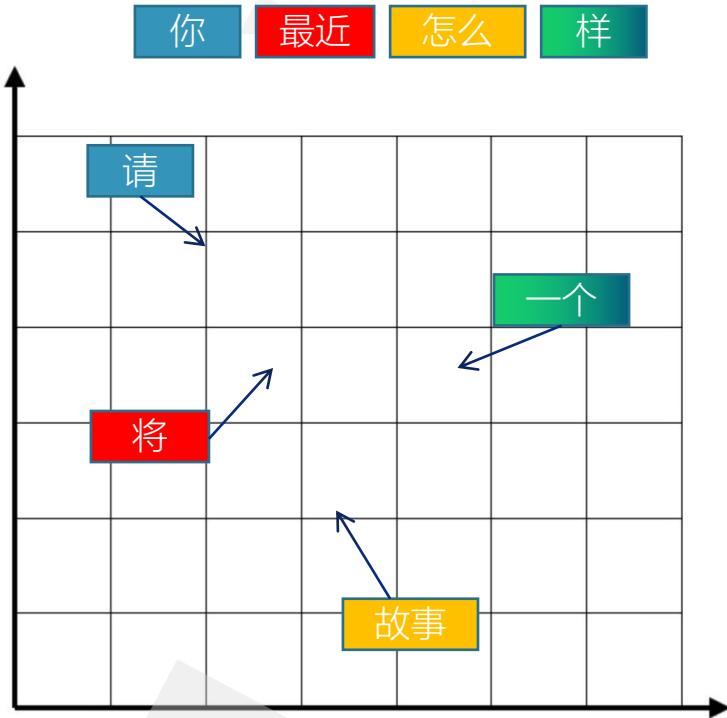




位置编码

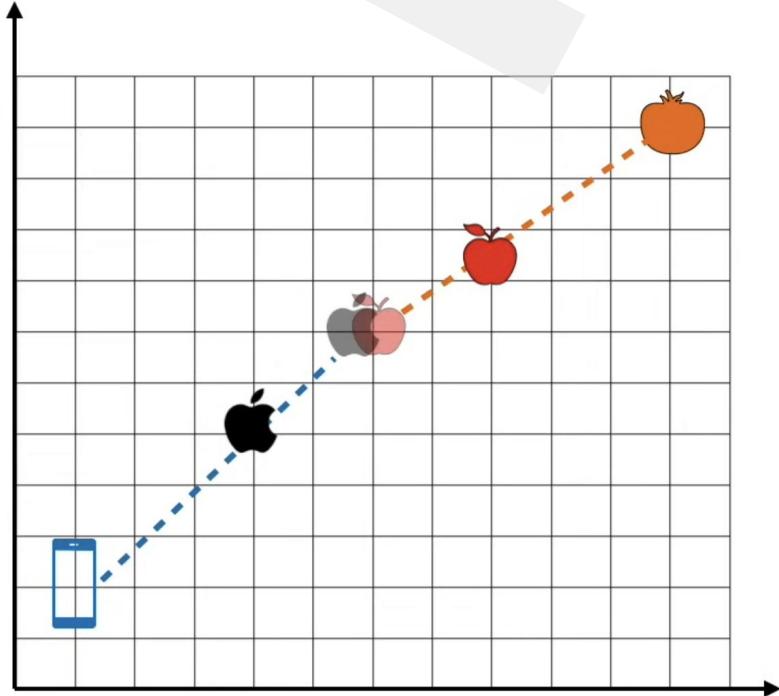
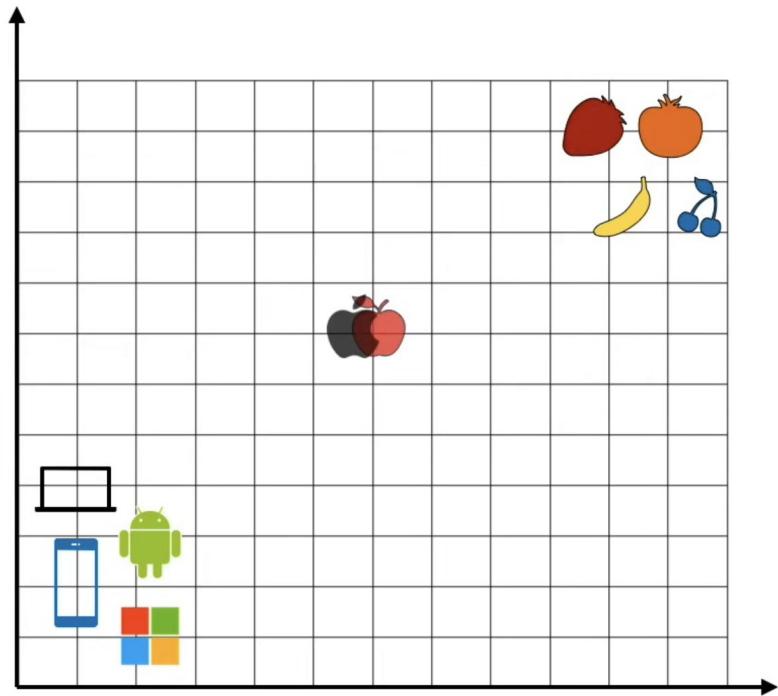


位置编码



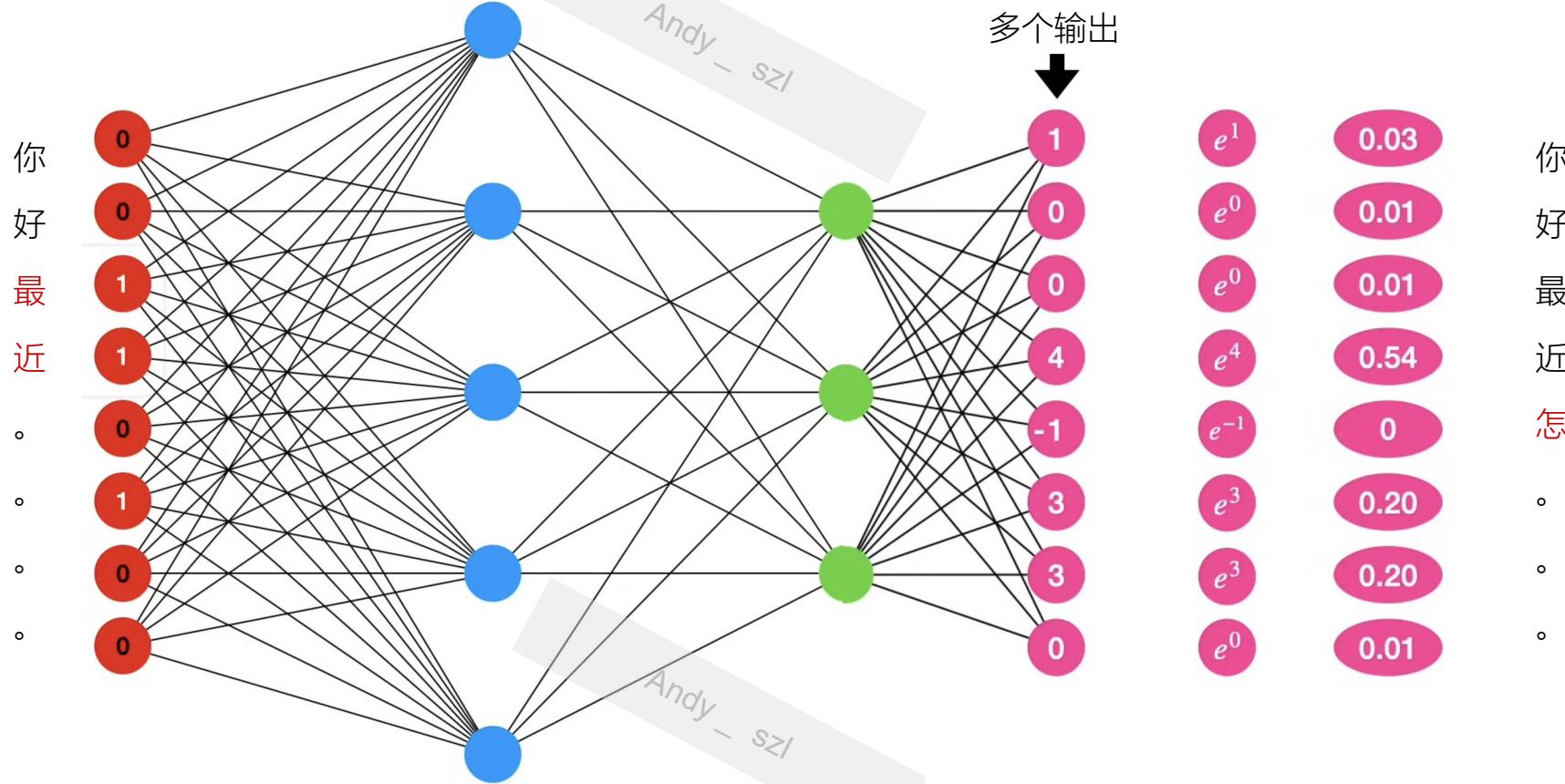


自注意力机制

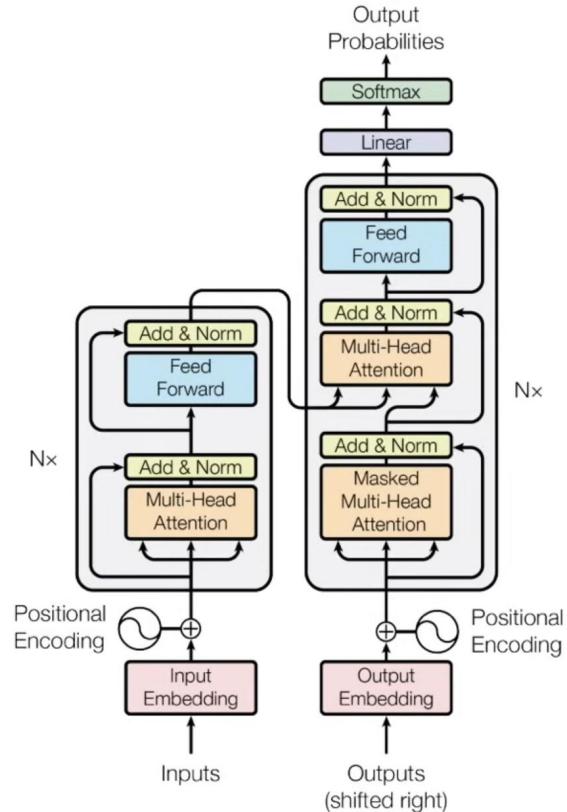




Soft Max



整个架构

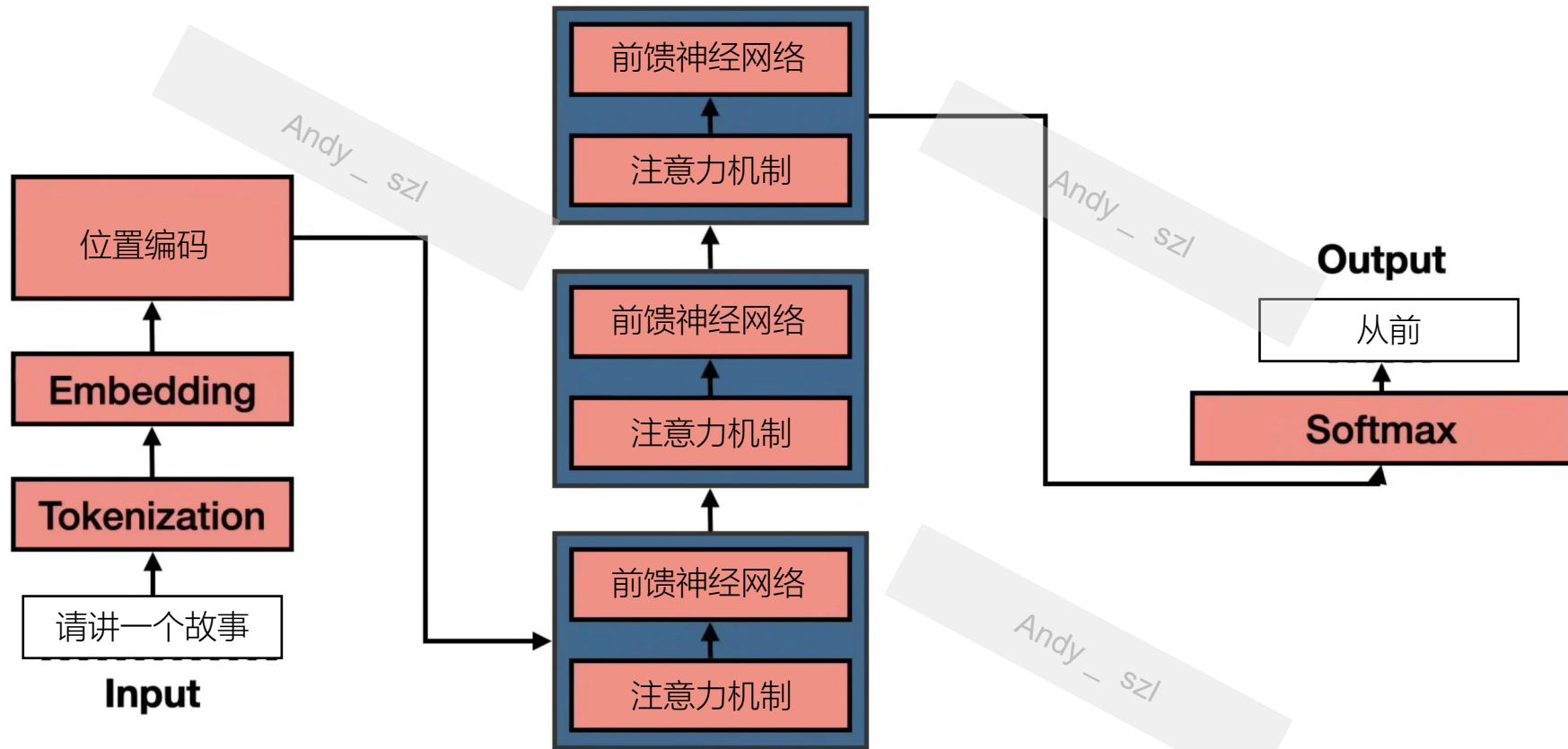


Andy

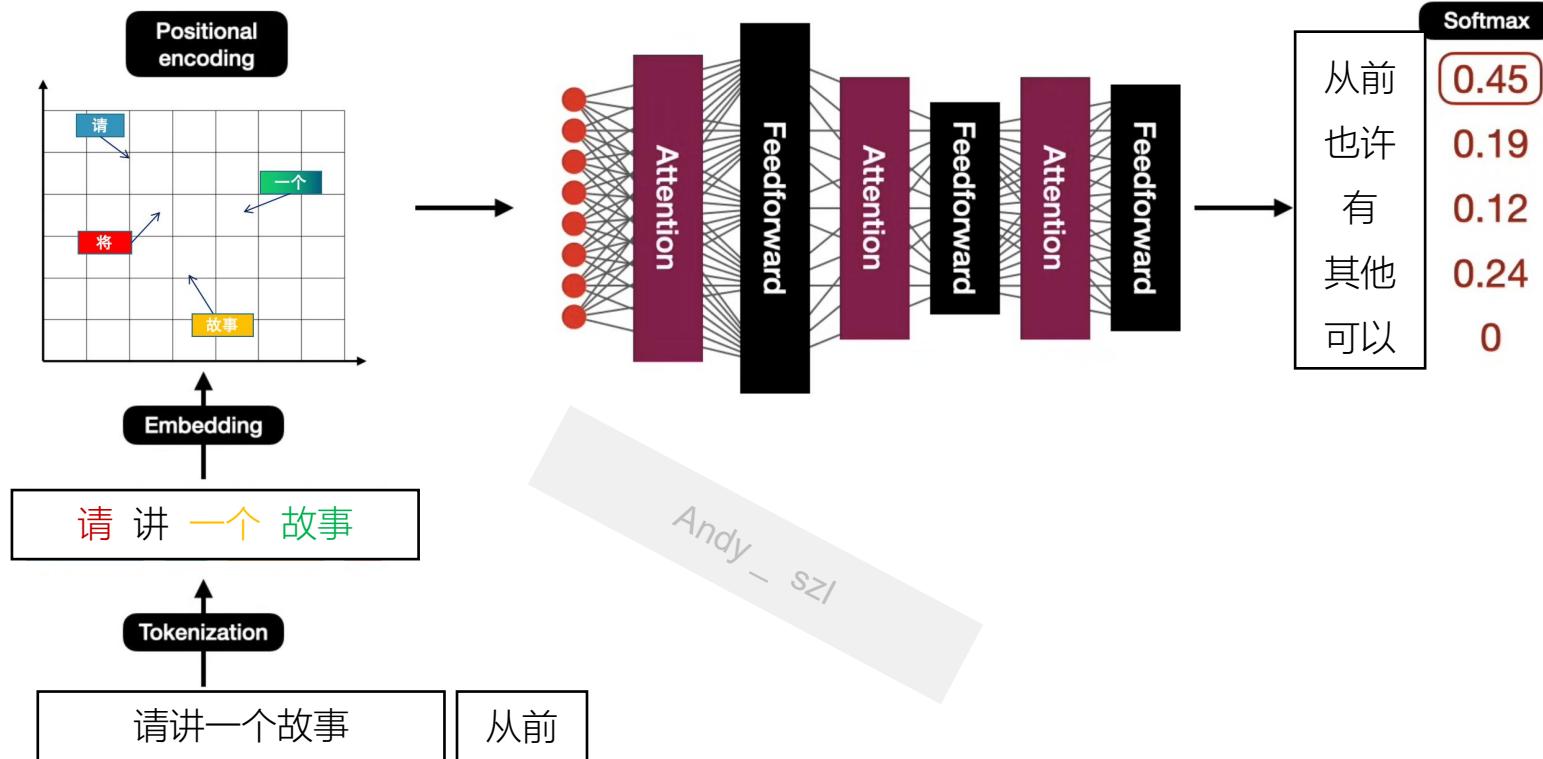
Andy - szl

27

整个架构



整个架构



位置编码

$$P(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$P(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$



Andy - szl

PART four

Andy - szl

编码器和解码器

Andy - szl

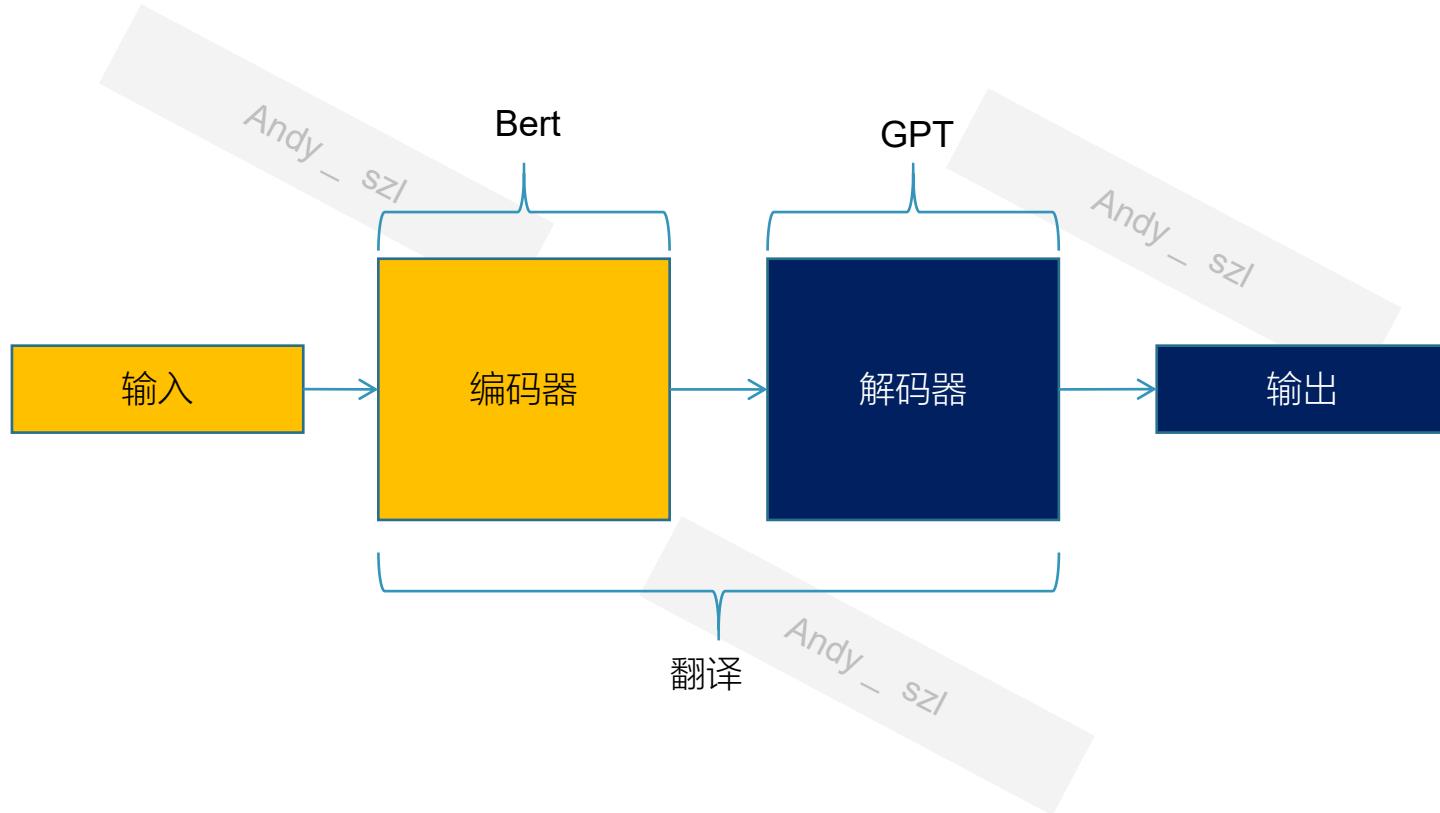


$$1+5=\frac{1}{3}$$

$$A+\frac{1}{2}$$



编码器和解码器



翻译示例

一加三

编码器

解码器

one plus three equals four



one
plus
three
equals
four

几个问题

- 1 解码器的输入是编码器的全部内容还是部分内容？
- 2 编码器和解码器使用的是同一套embedding还是两套？
- 3 会有使用一套embedding的模型吗
- 4 使用两套embedding时，两套embedding如何建立联系呢

Thank you

致謝

experiment

