

# Fundamentals of Large Language Models

---

Embeddings are all you need



# Steve Zielinski, CSEP

Currently consulting on neuromodulation device development.

Previously

Vice President, Product Development,  
Bioelectronic Therapies at Orchestra  
BioMed

BA University of St. Thomas  
BEE University of Minnesota  
MSDD University of St. Thomas

## The CSEP Study Guide for INCOSE SEH v4

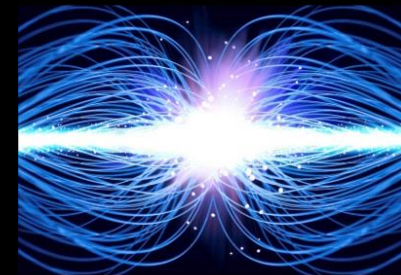
Steve Zielinski, CSEP  
and the staff of  
SystemsEngineeringPrep.com



Forward by  
William G. King Jr. IntPE, ESEP

## The Great Convergence

Merging  
Lean, Agile, and, Knowledge-based  
New Product Development



Steve Zielinski

# Goals

You should walk away with:

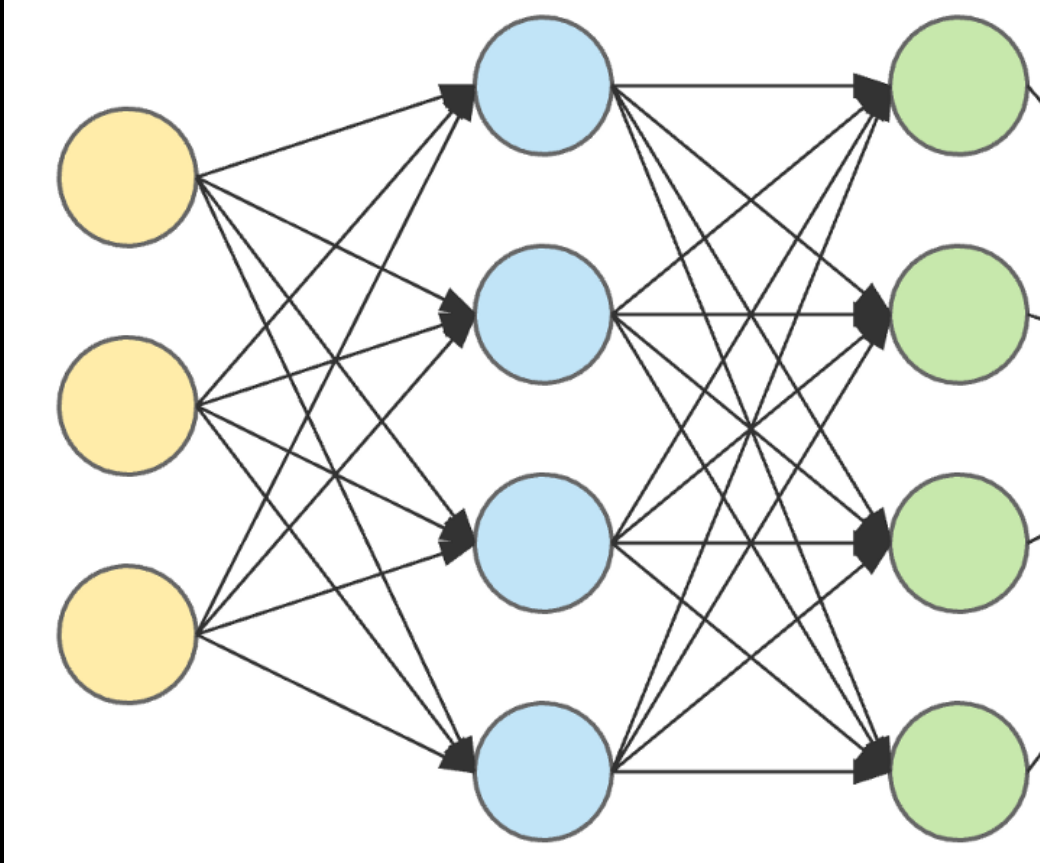
1. A deeper understanding of how ChatGPT works
2. Knowing what a Transformer is
3. A little scared/excited about the future.

# Natural Language Processing

# Starting Goal: Word Prediction

I want a glass of orange \_\_\_\_\_.

I want a glass of apple\_\_\_\_\_.



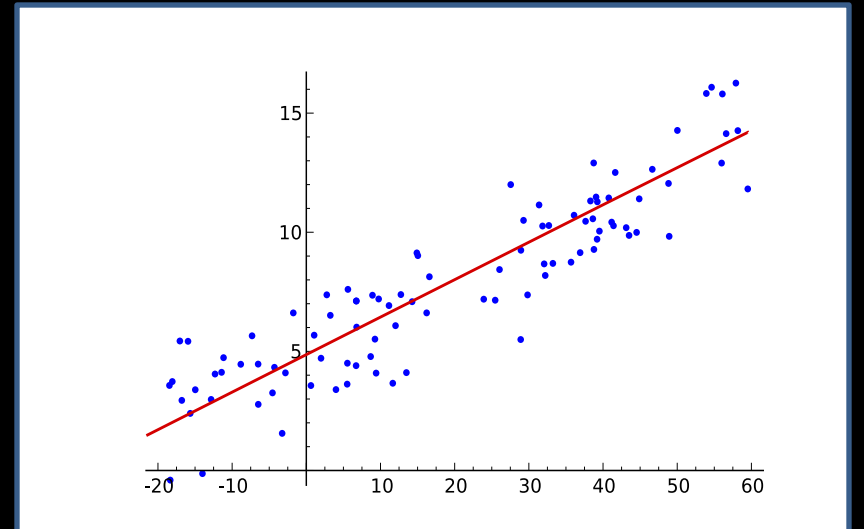
0.001

0.06

.5

.03

The NN is making a prediction about the most probable next word.



# The Big Problem:

Computers don't operate on words. They operate on numbers. How do you represent words as numbers?

Start with a vocabulary

$V = [a, aaron, \dots, zulu]$

These vocabularies can be anywhere from 10,000 to 100,000 (or more) words.

You'll frequently hear numbers in the 30,000 – 50,000 word range.



GPT-3 used 50,257 words

(Actually 50,257 tokens)

# 1-hot representation

man  
(5391)

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$$

woman  
(9853)

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

king  
(4914)

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

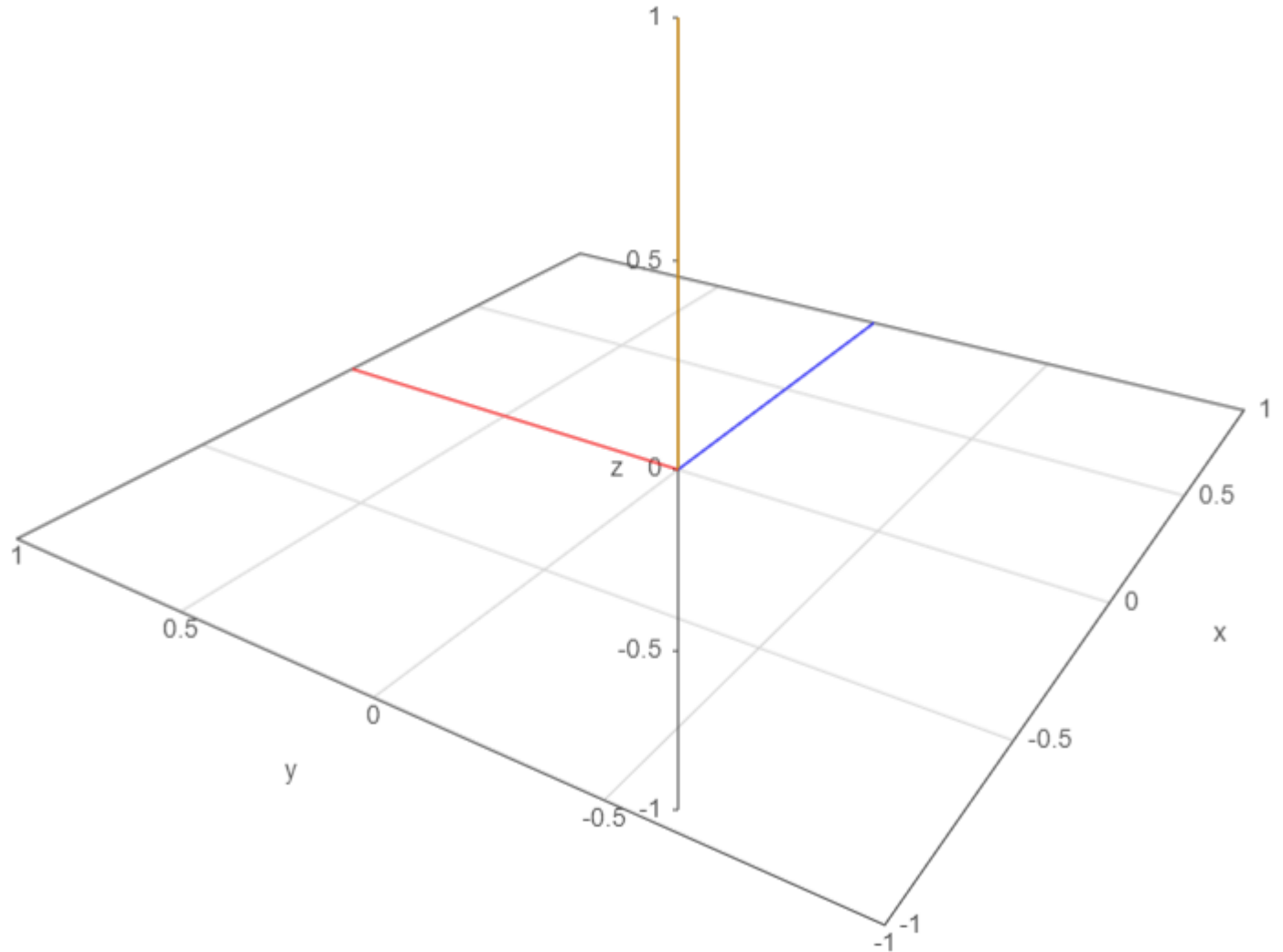
orange  
(6257)

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

# One-hot doesn't work very well

Every word is orthogonal to every other word.

The “distance” between each word is the same.



# Featurized representation: word embedding

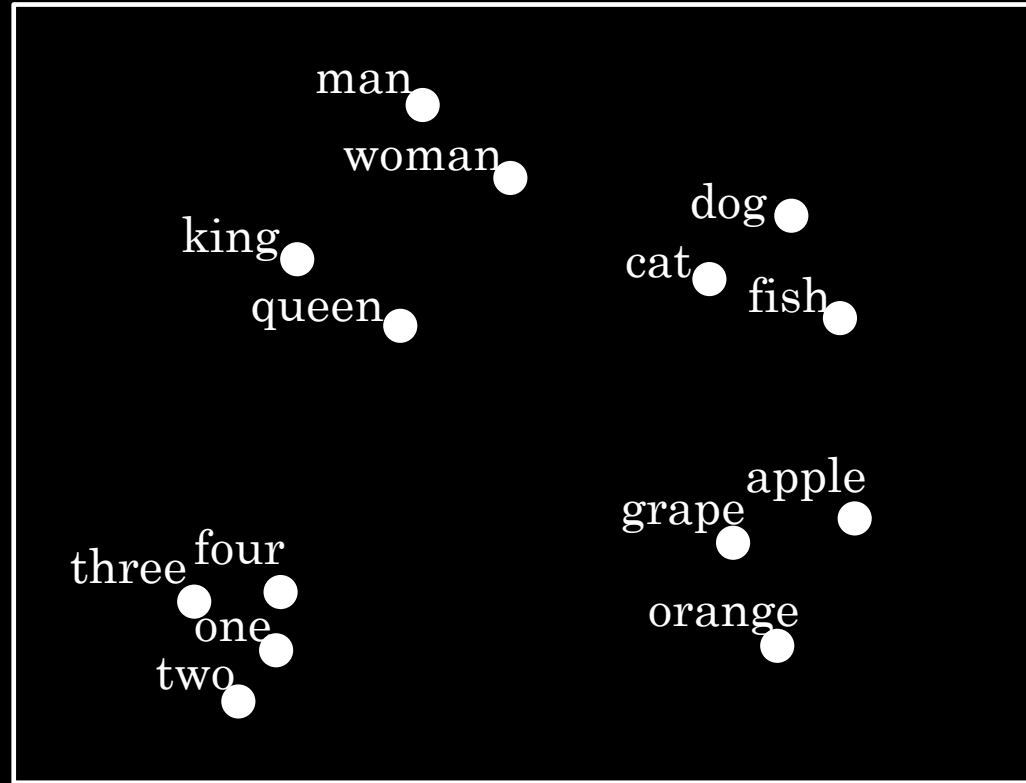
	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03		0.7	0.69	0.03	-0.02
Food			0.02	0.01	0.95	0.97

Embedding values are not assigned like this picture might imply.

They are learned from large amounts of text. Similar items congregate in the same area of space.

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03		0.7	0.69	0.03	-0.02
Food			0.02	0.01	0.95	0.97

# Visualizing word embeddings



GPT-3 uses word embeddings with  
12,288 characteristics

12,288-dimensional space





The goal of a transformer is to move the location of the starting embedding vector with information found in other parts of the sentence or document.

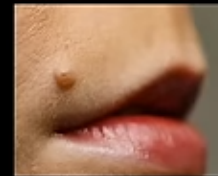
# Context changes meaning



American shrew **mole**

$6.02 \times 10^{23}$

One **mole** of carbon dioxide



Take a biopsy of the **mole**



American shrew mole

↓  
[ 6.0  
2.2  
3.9  
7.7  
6.1  
⋮  
6.3 ]

↓  
[ 0.4  
5.7  
5.0  
1.8  
9.7  
⋮  
5.4 ]

↓  
[ 5.8  
9.9  
2.5  
3.7  
9.1  
⋮  
2.1 ]

One mole of carbon dioxide

↓  
[ 5.2  
7.8  
2.5  
5.9  
9.8  
⋮  
2.7 ]

↓  
[ 5.8  
9.9  
2.5  
3.7  
9.1  
⋮  
2.1 ]

↓  
[ 5.8  
7.0  
4.0  
0.1  
4.3  
⋮  
4.5 ]

↓  
[ 7.6  
4.5  
5.7  
8.1  
5.6  
⋮  
4.8 ]

↓  
[ 9.9  
1.8  
6.1  
9.8  
9.1  
⋮  
0.4 ]

Take a biopsy of the mole

↓  
[ 4.9  
2.1  
4.7  
9.6  
8.0  
⋮  
2.2 ]

↓  
[ 3.5  
9.7  
3.6  
8.3  
0.8  
⋮  
8.9 ]

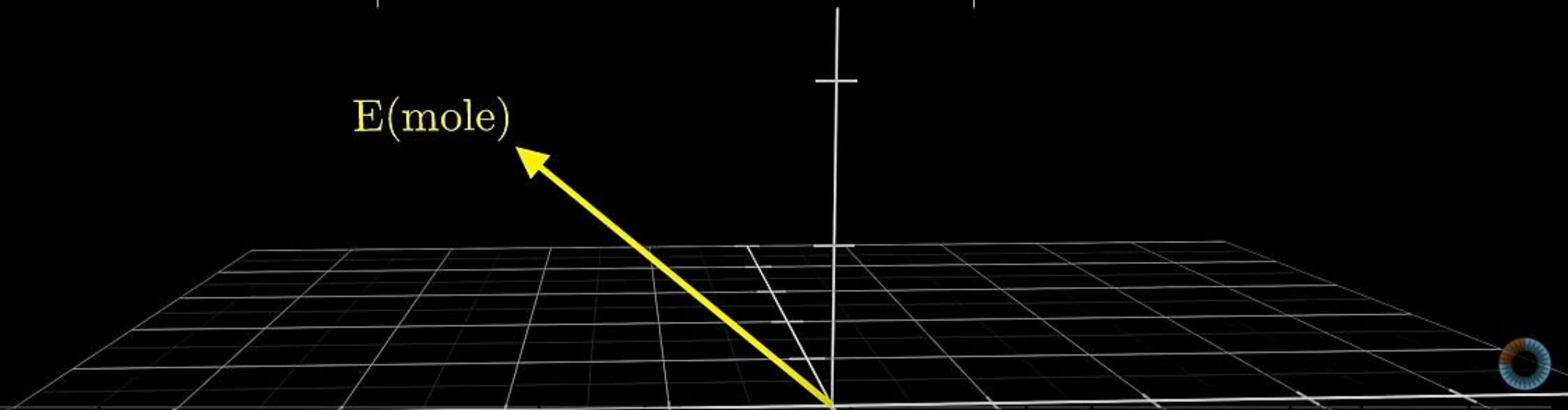
↓  
[ 1.7  
8.7  
3.4  
2.7  
4.7  
⋮  
2.3 ]

↓  
[ 5.8  
7.0  
4.0  
0.1  
4.3  
⋮  
4.5 ]

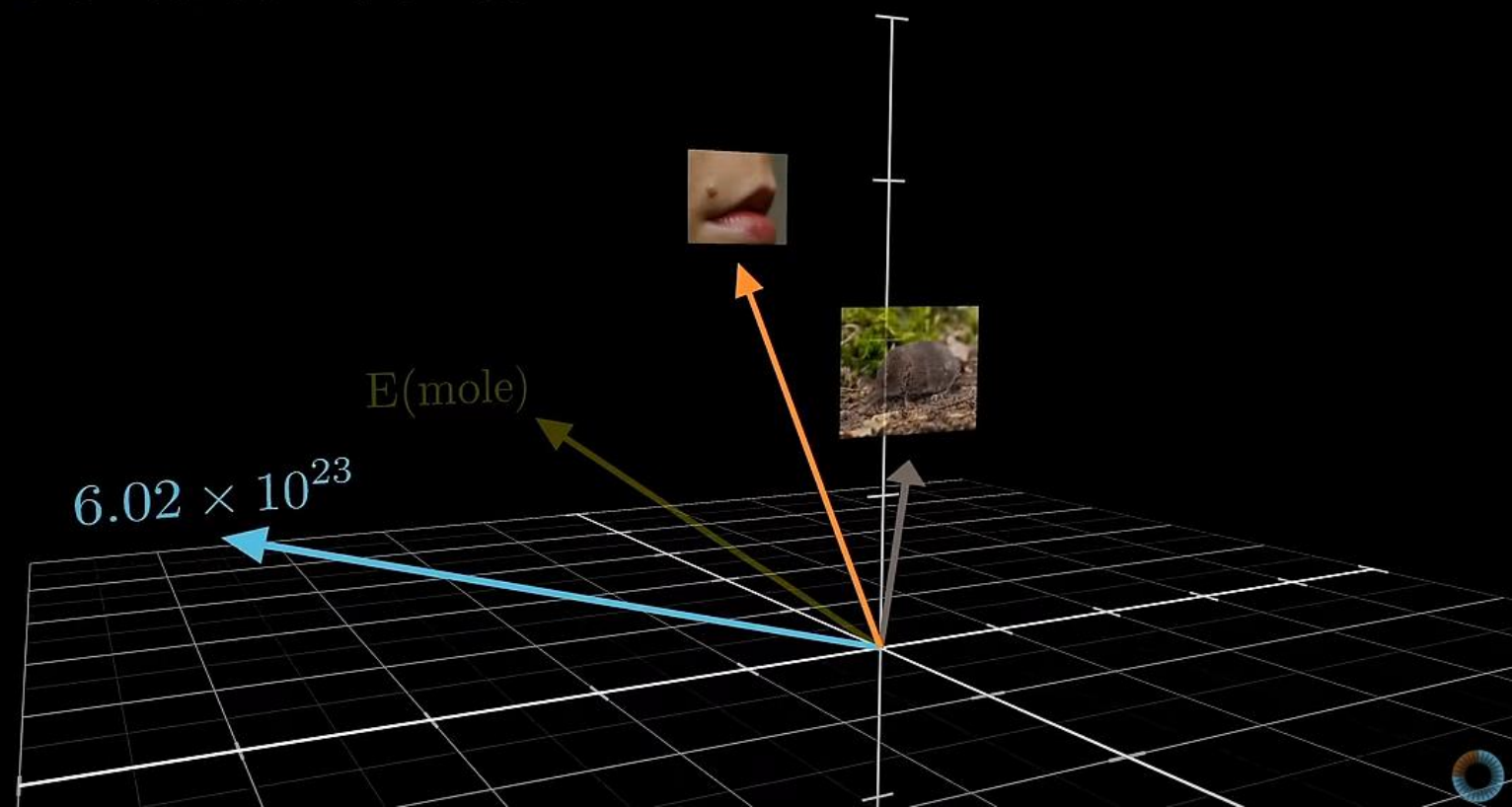
↓  
[ 2.3  
4.9  
6.4  
3.2  
4.4  
⋮  
6.5 ]

↓  
[ 5.8  
9.9  
2.5  
3.7  
9.1  
⋮  
2.1 ]

E(mole)



One mole of carbon dioxide





The attention block  
“moves” the vector  
for the word to  
different locations  
in the embedding  
space.



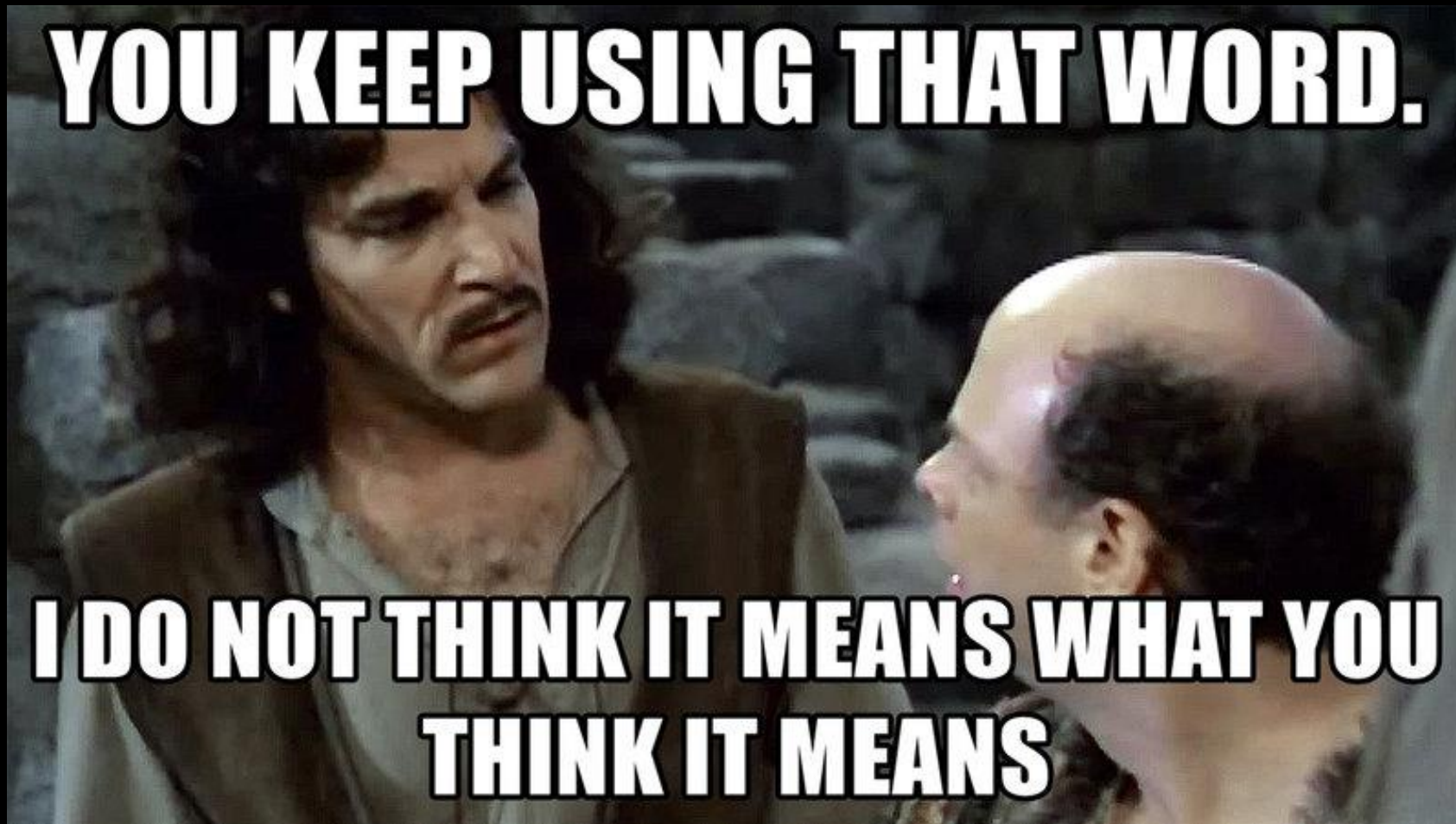
# Transformer

- Innovation #1

	One	mole	of	carbon	dioxide
0	0	0	-0.95	0.97	0.00
0.01	0.01	0.02	0.93	0.95	-0.01
0.03			0.7	0.69	0.03
			0.02	0.01	0.95
Position					



Position





# Queen

# and

# king

0.33	0.71	0.91	0.23	0.15
------	------	------	------	------

+

1	0	0	0	0
---	---	---	---	---

0.12	0.22	0.31	0.03	0.10
------	------	------	------	------

+

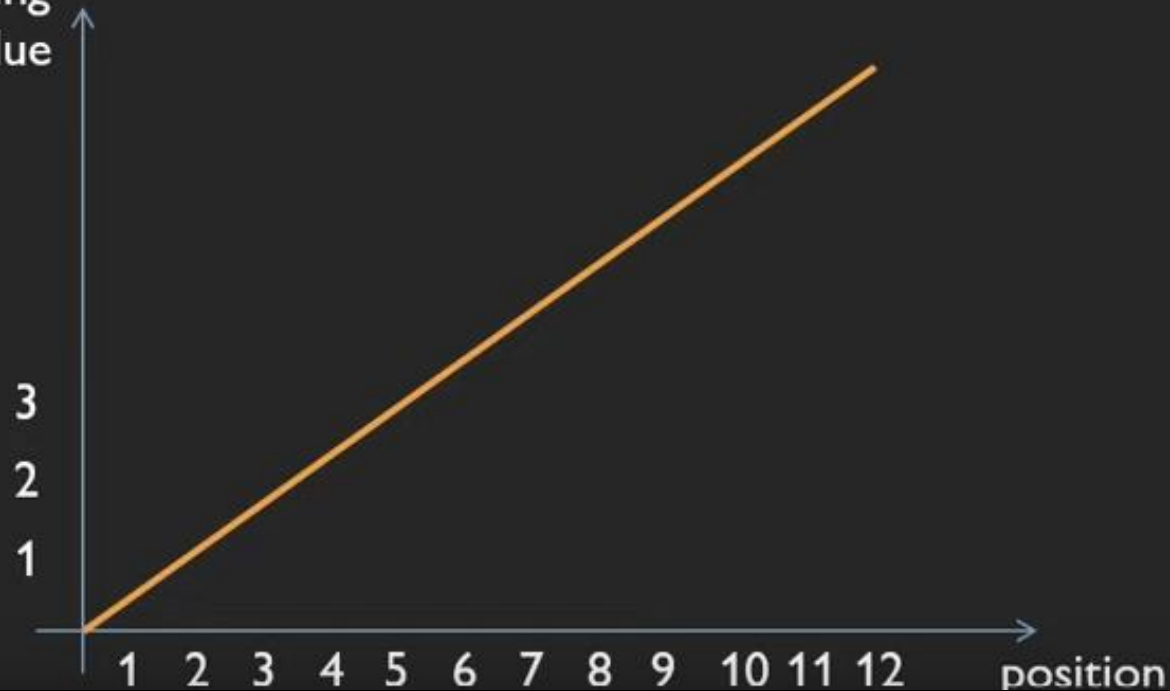
2	0	0	0	0
---	---	---	---	---

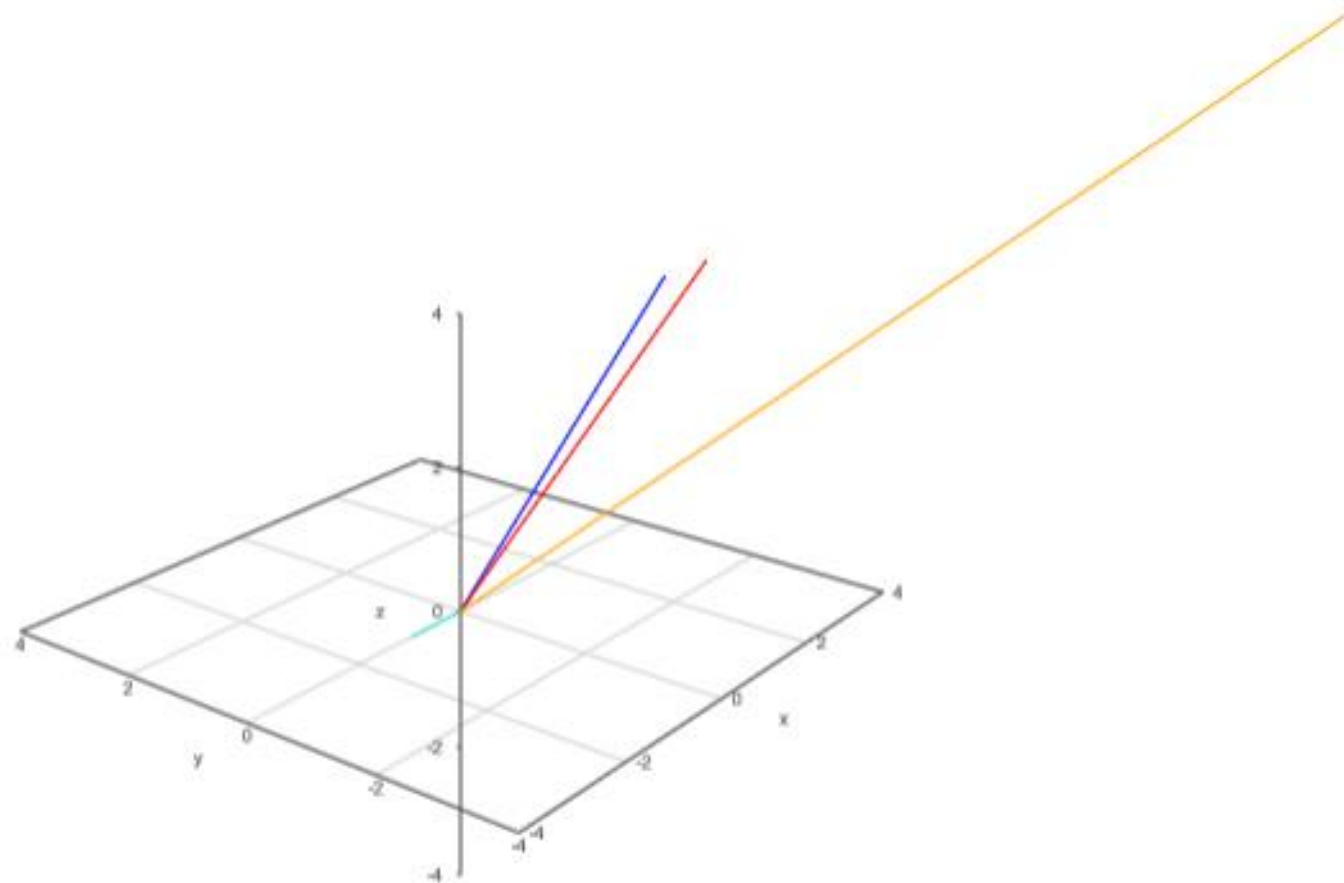
0.34	0.70	0.95	0.21	0.17
------	------	------	------	------

+

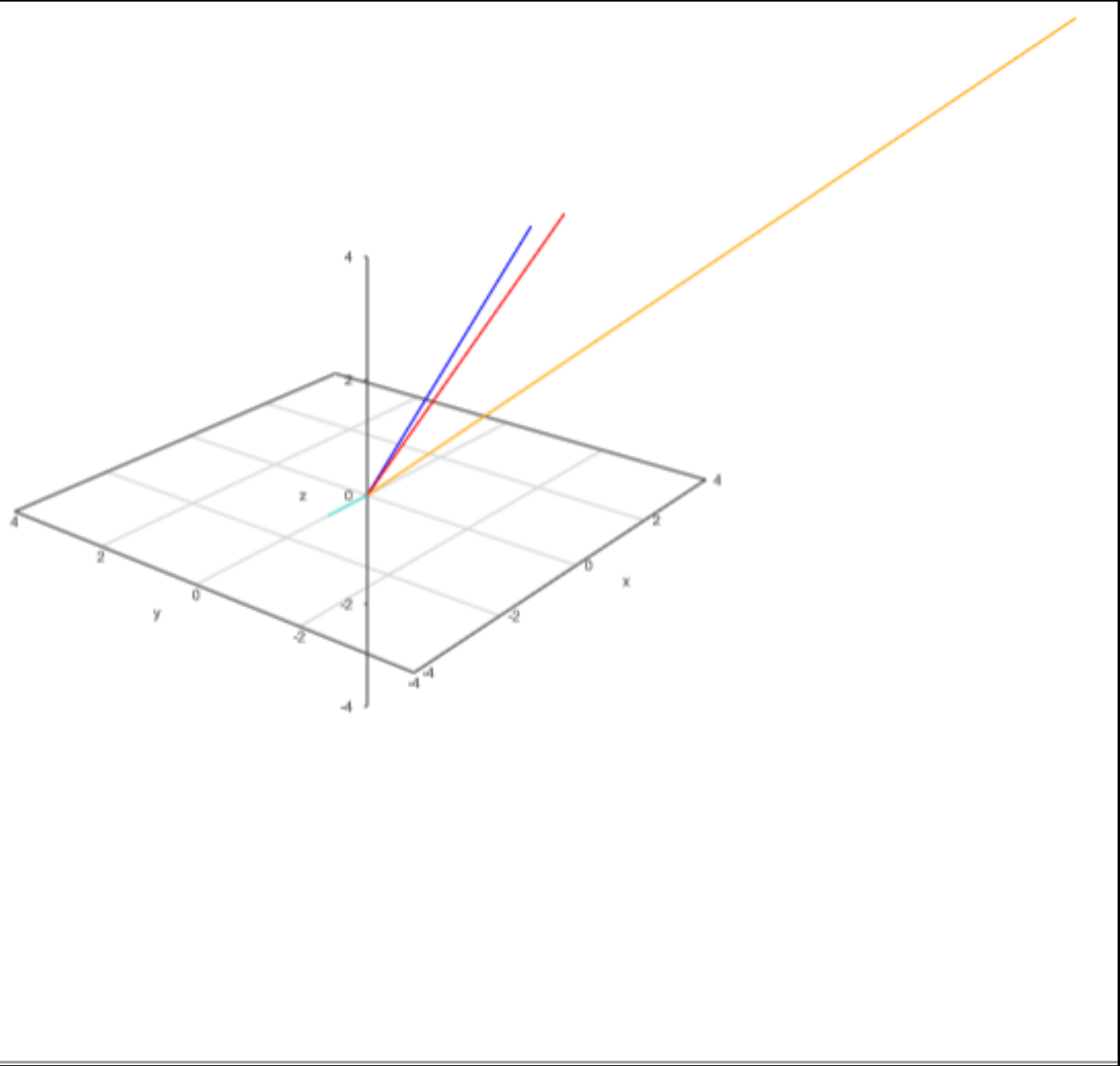
3	0	0	0	0
---	---	---	---	---

bedding  
value





The word  
position in the  
sentence is  
overwhelming  
the semantics.



$$PE(\text{position}, 2i) = \sin \left( \frac{\text{position}}{10000^{\frac{2i}{d_{model}}}} \right)$$

$$PE(\text{position}, 2i + 1) = \cos \left( \frac{\text{position}}{10000^{\frac{2i}{d_{model}}}} \right)$$

# Queen and king

0.33	0.71	0.91	0.23	0.15
------	------	------	------	------

+

0.2	0.7	0.1	0.99	0.01
-----	-----	-----	------	------

0.12	0.22	0.31	0.03	0.10
------	------	------	------	------

+

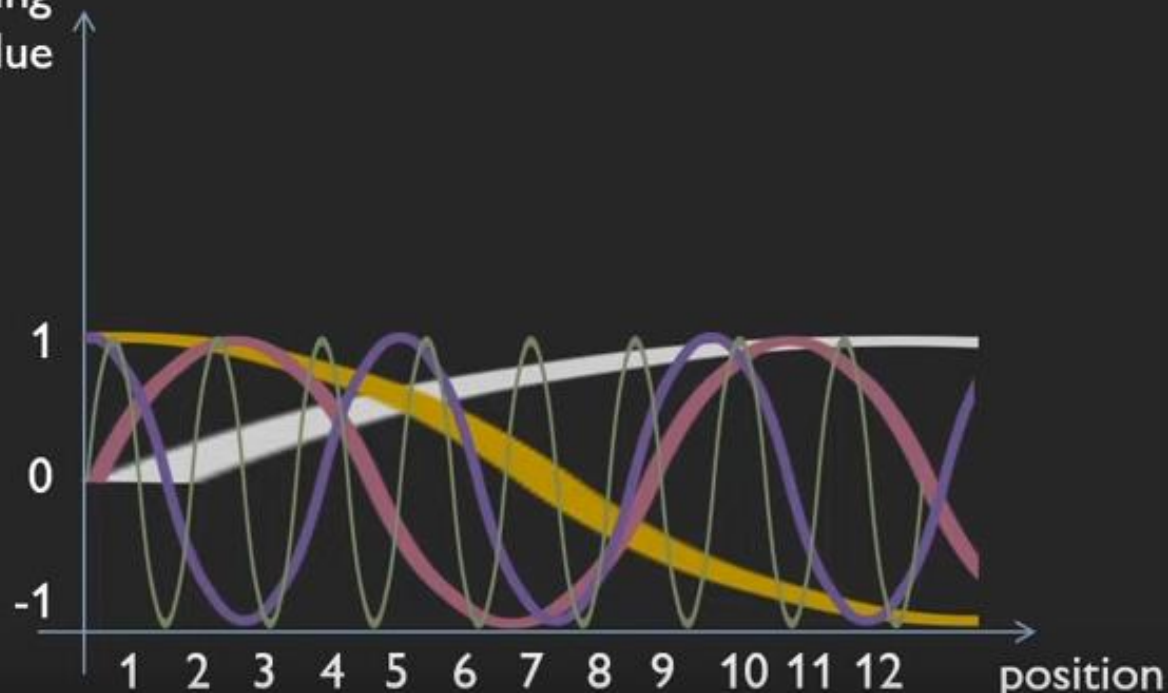
0.01	-0.8	0.9	0.80	0.02
------	------	-----	------	------

0.34	0.70	0.95	0.21	0.17
------	------	------	------	------

+

-0.9	-0.7	0.8	0.70	0.03
------	------	-----	------	------

embedding  
value



Queen

0.33	0.71	0.91	0.23	0.15
------	------	------	------	------

+

0.01	0.01	0.01	0.01	0.01
------	------	------	------	------



and

0.12	0.22	0.31	0.03	0.10
------	------	------	------	------

+

0.02	0.21	0.33	0.43	0.98
------	------	------	------	------



king

0.34	0.70	0.95	0.21	0.17
------	------	------	------	------

+

0.03	0.32	0.85	0.91	0.24
------	------	------	------	------



Prince

0.23	0.72	0.62	0.21	0.17
------	------	------	------	------

+

0.01	0.01	0.01	0.01	0.01
------	------	------	------	------



[PAD]

0.00	0.00	0.00	0.00	0.00
------	------	------	------	------

+

0.02	0.21	0.33	0.43	0.98
------	------	------	------	------



[PAD]

0.00	0.00	0.00	0.00	0.00
------	------	------	------	------

+

0.03	0.32	0.85	0.91	0.24
------	------	------	------	------



# Parallel computations

- Because word position is now embedded in the vector, I can operate on all the input words in parallel!
- This is why GPUs are conquering the world!

# Transformer Components

Attention : What part of the input should we focus?

	Focus		Attention Vectors
The	→	The big red dog	$[0.71 \quad 0.04 \quad 0.07 \quad 0.18]^T$
big	→	The big red dog	$[0.01 \quad 0.84 \quad 0.02 \quad 0.13]^T$
red	→	The big red dog	$[0.09 \quad 0.05 \quad 0.62 \quad 0.24]^T$
dog	→	The big red dog	$[0.03 \quad 0.03 \quad 0.03 \quad 0.91]^T$



The attention vectors establish a relationship between words.

**QUERIES, KEYS, AND VALUES (OH, MY!)**

Query: A word “asks”: “What is in front of me?”

dog  $\rightarrow$  The big red dog  $[0.03 \quad 0.03 \quad 0.03 \quad 0.91]^T$

$\downarrow$   
 $\vec{E}_3$

$\downarrow$   
 $\vec{E}_4$

$\downarrow W_Q$

$\vec{Q}_4$

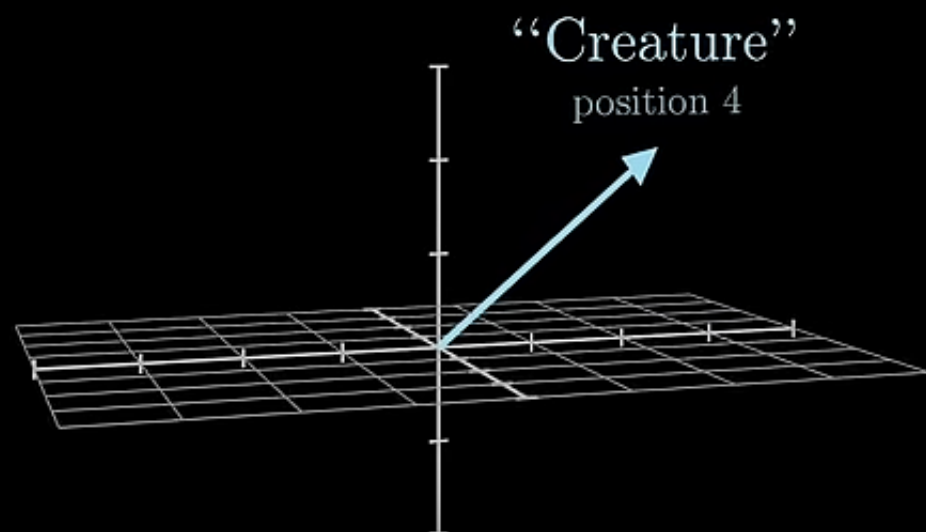
Any adjectives  
in front of me?

We need a vector “capable” of asking the question.

$$\overbrace{\begin{bmatrix} +7.5 & -3.2 & +9.1 & -5.3 & +8.9 & +8.7 & +5.9 & +2.6 & +7.4 & -4.1 & \cdots & +2.3 \\ -9.6 & -3.0 & -7.0 & +9.5 & -0.4 & -0.1 & +2.8 & -2.6 & -7.2 & +6.4 & \cdots & +0.2 \\ -5.5 & -8.0 & +7.2 & +9.4 & +9.1 & +8.0 & +5.4 & -3.3 & -8.3 & -1.8 & \cdots & -7.3 \\ -8.8 & +4.5 & -9.7 & +5.4 & -7.0 & -8.3 & -8.1 & +3.4 & -5.0 & -1.6 & \cdots & +7.1 \\ +4.5 & -4.5 & -7.3 & -8.8 & -3.9 & -4.7 & -0.9 & +3.6 & +3.9 & -4.3 & \cdots & -6.3 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ -9.0 & +5.9 & -8.4 & +0.4 & -3.8 & +1.5 & +9.1 & +2.9 & -9.2 & -1.4 & \cdots & +0.7 \end{bmatrix}}^{W_Q} \begin{bmatrix} 2.9 \\ 2.4 \\ 1.0 \\ 0.2 \\ 9.2 \\ 6.6 \\ 7.8 \\ 2.8 \\ 5.8 \\ 0.6 \\ \vdots \\ 9.7 \end{bmatrix} = \begin{bmatrix} +310.6 \\ -95.2 \\ -2.1 \\ -152.0 \\ -123.2 \\ \vdots \\ -12.7 \end{bmatrix}$$

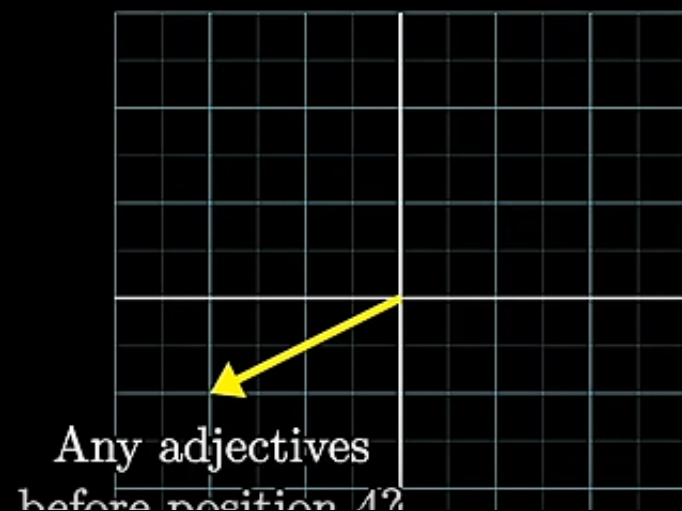
$\vec{E}_i$        $\vec{Q}_i$

Embedding space  
12,288-dimensional

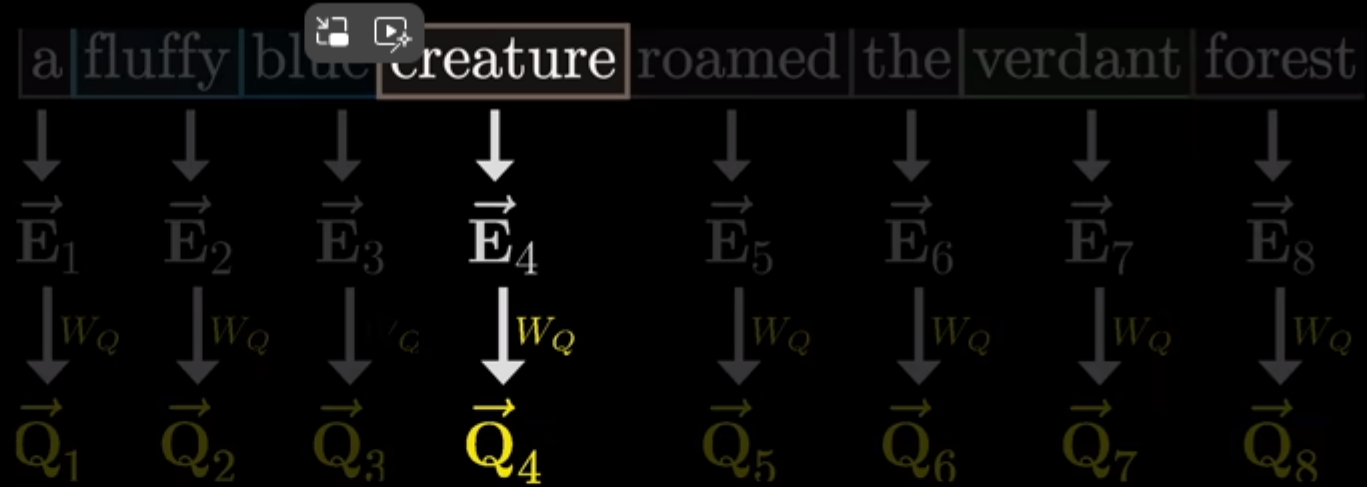


$W_Q$

Query/Key space  
128-dimensional



The Key matrix  
is used to  
answer to the  
query.



`a`  $\rightarrow$   $\vec{E}_1$   $\xrightarrow{W_k}$   $\vec{K}_1$

I'm an adjective!  
I'm there!

`fluffy`  $\rightarrow$   $\vec{E}_2$   $\xrightarrow{W_k}$   $\vec{K}_2$

Any adjectives  
in front of me?

`blue`  $\rightarrow$   $\vec{E}_3$   $\xrightarrow{W_k}$   $\vec{K}_3$

`creature`  $\rightarrow$   $\vec{E}_4$   $\xrightarrow{W_k}$   $\vec{K}_4$

I'm an adjective!  
I'm there!

`roamed`  $\rightarrow$   $\vec{E}_5$   $\xrightarrow{W_k}$   $\vec{K}_5$

$$\boxed{\text{a}} \rightarrow \vec{E}_1 \xrightarrow{W_k} \vec{K}_1$$

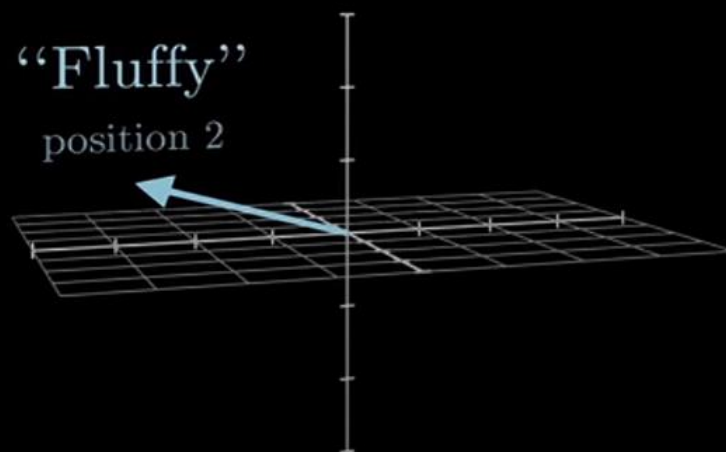
$$\boxed{\text{fluffy}} \rightarrow \vec{E}_2 \xrightarrow{W_k} \vec{K}_2$$

$$\boxed{\text{blue}} \rightarrow \vec{E}_3 \xrightarrow{W_k} \vec{K}_3$$

$$\boxed{\text{creature}} \rightarrow \vec{E}_4 \xrightarrow{W_k} \vec{K}_4$$

 $W_v$ 

Embedding space  
12,288-dimensional


 $W_K$ 

Query/Key space  
128-dimensional



Whoa! Hold on. Where did we get these Query and Key matrices?





- But, we haven't answered how “good” an answer the key is to the query.
- For that we use a dot product.
- This is referred to as the “similarity” between the words.

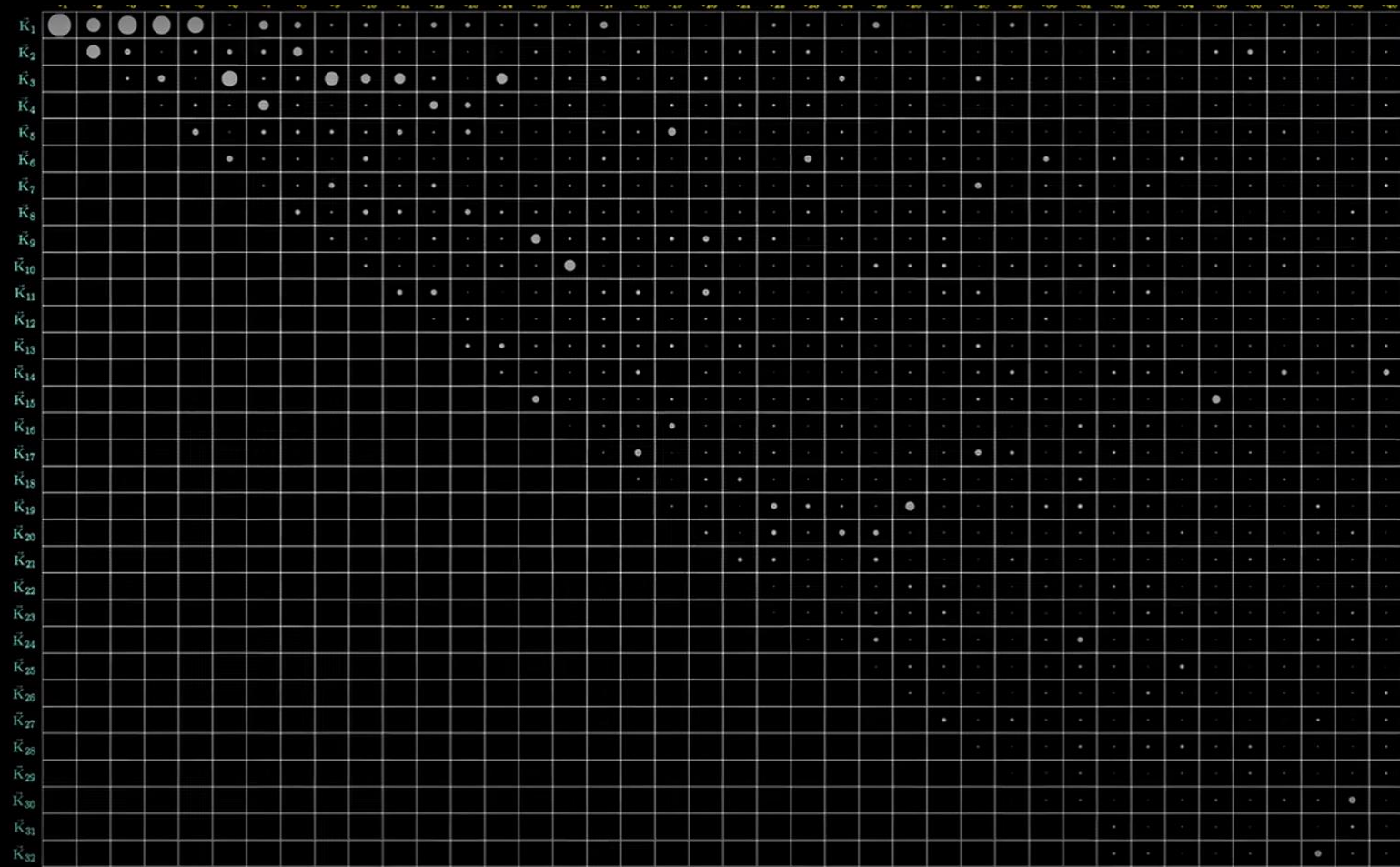
	a	fluffy	blue	creature	roamed	the	verdant	forest	
	$\downarrow$ $\vec{E}_1$ $\downarrow^{W_Q}$ $\vec{Q}_1$	$\downarrow$ $\vec{E}_2$ $\downarrow^{W_Q}$ $\vec{Q}_2$	$\downarrow$ $\vec{E}_3$ $\downarrow^{W_Q}$ $\vec{Q}_3$	$\downarrow$ $\vec{E}_4$ $\downarrow^{W_Q}$ $\vec{Q}_4$	$\downarrow$ $\vec{E}_5$ $\downarrow^{W_Q}$ $\vec{Q}_5$	$\downarrow$ $\vec{E}_6$ $\downarrow^{W_Q}$ $\vec{Q}_6$	$\downarrow$ $\vec{E}_7$ $\downarrow^{W_Q}$ $\vec{Q}_7$	$\downarrow$ $\vec{E}_8$ $\downarrow^{W_Q}$ $\vec{Q}_8$	
$\boxed{\text{a}} \rightarrow \vec{E}_1 \xrightarrow{W_k} \vec{K}_1$	$\vec{K}_1 \cdot \vec{Q}_1$	$\vec{K}_1 \cdot \vec{Q}_2$	$\vec{K}_1 \cdot \vec{Q}_3$	$\vec{K}_1 \cdot \vec{Q}_4$	$\vec{K}_1 \cdot \vec{Q}_5$	$\vec{K}_1 \cdot \vec{Q}_6$	$\vec{K}_1 \cdot \vec{Q}_7$	$\vec{K}_1 \cdot \vec{Q}_8$	
$\boxed{\text{fluffy}} \rightarrow \vec{E}_2 \xrightarrow{W_k} \vec{K}_2$	$\vec{K}_2 \cdot \vec{Q}_1$	$\vec{K}_2 \cdot \vec{Q}_2$	$\vec{K}_2 \cdot \vec{Q}_3$	$\vec{K}_2 \cdot \vec{Q}_4$	$\vec{K}_2 \cdot \vec{Q}_5$	$\vec{K}_2 \cdot \vec{Q}_6$	$\vec{K}_2 \cdot \vec{Q}_7$	$\vec{K}_2 \cdot \vec{Q}_8$	
$\boxed{\text{blue}} \rightarrow \vec{E}_3 \xrightarrow{W_k} \vec{K}_3$	$\vec{K}_3 \cdot \vec{Q}_1$	$\vec{K}_3 \cdot \vec{Q}_2$	$\vec{K}_3 \cdot \vec{Q}_3$	$\vec{K}_3 \cdot \vec{Q}_4$	$\vec{K}_3 \cdot \vec{Q}_5$	$\vec{K}_3 \cdot \vec{Q}_6$	$\vec{K}_3 \cdot \vec{Q}_7$	$\vec{K}_3 \cdot \vec{Q}_8$	
$\boxed{\text{creature}} \rightarrow \vec{E}_4 \xrightarrow{W_k} \vec{K}_4$	$\vec{K}_4 \cdot \vec{Q}_1$	$\vec{K}_4 \cdot \vec{Q}_2$	$\vec{K}_4 \cdot \vec{Q}_3$	$\vec{K}_4 \cdot \vec{Q}_4$	$\vec{K}_4 \cdot \vec{Q}_5$	$\vec{K}_4 \cdot \vec{Q}_6$	$\vec{K}_4 \cdot \vec{Q}_7$	$\vec{K}_4 \cdot \vec{Q}_8$	
$\boxed{\text{roamed}} \rightarrow \vec{E}_5 \xrightarrow{W_k} \vec{K}_5$	$\vec{K}_5 \cdot \vec{Q}_1$	$\vec{K}_5 \cdot \vec{Q}_2$	$\vec{K}_5 \cdot \vec{Q}_3$	$\vec{K}_5 \cdot \vec{Q}_4$	$\vec{K}_5 \cdot \vec{Q}_5$	$\vec{K}_5 \cdot \vec{Q}_6$	$\vec{K}_5 \cdot \vec{Q}_7$	$\vec{K}_5 \cdot \vec{Q}_8$	
$\boxed{\text{the}} \rightarrow \vec{E}_6 \xrightarrow{W_k} \vec{K}_6$	$\vec{K}_6 \cdot \vec{Q}_1$	$\vec{K}_6 \cdot \vec{Q}_2$	$\vec{K}_6 \cdot \vec{Q}_3$	$\vec{K}_6 \cdot \vec{Q}_4$	$\vec{K}_6 \cdot \vec{Q}_5$	$\vec{K}_6 \cdot \vec{Q}_6$	$\vec{K}_6 \cdot \vec{Q}_7$	$\vec{K}_6 \cdot \vec{Q}_8$	
$\boxed{\text{verdant}} \rightarrow \vec{E}_7 \xrightarrow{W_k} \vec{K}_7$	$\vec{K}_7 \cdot \vec{Q}_1$	$\vec{K}_7 \cdot \vec{Q}_2$	$\vec{K}_7 \cdot \vec{Q}_3$	$\vec{K}_7 \cdot \vec{Q}_4$	$\vec{K}_7 \cdot \vec{Q}_5$	$\vec{K}_7 \cdot \vec{Q}_6$	$\vec{K}_7 \cdot \vec{Q}_7$	$\vec{K}_7 \cdot \vec{Q}_8$	

	<div>a</div> <div><math>\downarrow</math> <math>\vec{E}_1</math> <math>\downarrow^{W_Q}</math> <math>\vec{Q}_1</math></div>	<div>fluffy</div> <div><math>\downarrow</math> <math>\vec{E}_2</math> <math>\downarrow^{W_Q}</math> <math>\vec{Q}_2</math></div>	<div>blue</div> <div><math>\downarrow</math> <math>\vec{E}_3</math> <math>\downarrow^{W_Q}</math> <math>\vec{Q}_3</math></div>	<div>creature</div> <div><math>\downarrow</math> <math>\vec{E}_4</math> <math>\downarrow^{W_Q}</math> <math>\vec{Q}_4</math></div>	<div>roamed</div> <div><math>\downarrow</math> <math>\vec{E}_5</math> <math>\downarrow^{W_Q}</math> <math>\vec{Q}_5</math></div>	<div>the</div> <div><math>\downarrow</math> <math>\vec{E}_6</math> <math>\downarrow^{W_Q}</math> <math>\vec{Q}_6</math></div>	<div>verdant</div> <div><math>\downarrow</math> <math>\vec{E}_7</math> <math>\downarrow^{W_Q}</math> <math>\vec{Q}_7</math></div>	<div>forest</div> <div><math>\downarrow</math> <math>\vec{E}_8</math> <math>\downarrow^{W_Q}</math> <math>\vec{Q}_8</math></div>	
<div>a</div> $\rightarrow \vec{E}_1 \xrightarrow{W_k} \vec{K}_1$	+0.7	-83.7	-24.7	-27.8	-5.2	-89.3	-45.2	-36.1	
<div>fluffy</div> $\rightarrow \vec{E}_2 \xrightarrow{W_k} \vec{K}_2$	-73.4	+2.9	-5.4	+93.0	-48.2	-87.3	-49.7	+7.8	
<div>blue</div> $\rightarrow \vec{E}_3 \xrightarrow{W_k} \vec{K}_3$	-53.4	-5.7	+1.8	+93.4	-55.6	-56.0	-26.1	-62.1	
<div>creature</div> $\rightarrow \vec{E}_4 \xrightarrow{W_k} \vec{K}_4$	-21.5	-29.7	-56.1	+4.9	-32.4	-92.3	-9.5	-28.1	
<div>roamed</div> $\rightarrow \vec{E}_5 \xrightarrow{W_k} \vec{K}_5$	-20.1	-40.9	-87.8	-55.4	+0.6	-64.7	-96.7	-18.9	
<div>the</div> $\rightarrow \vec{E}_6 \xrightarrow{W_k} \vec{K}_6$	-87.9	-33.3	-22.6	-31.4	+5.5	+0.6	-4.6	-96.8	
<div>verdant</div> $\rightarrow \vec{E}_7 \xrightarrow{W_k} \vec{K}_7$	-41.2	-55.5	-42.3	-59.8	-79.0	-97.9	+3.7	+93.8	

The the columns  
get normalized.

	<span>a</span> ↓ $\vec{E}_1$ ↓ $W_Q$ $\vec{Q}_1$	<span>fluffy</span> ↓ $\vec{E}_2$ ↓ $W_Q$ $\vec{Q}_2$	<span>blue</span> ↓ $\vec{E}_3$ ↓ $W_Q$ $\vec{Q}_3$	<span>creature</span> ↓ $\vec{E}_4$ ↓ $W_Q$ $\vec{Q}_4$	<span>roamed</span> ↓ $\vec{E}_5$ ↓ $W_Q$ $\vec{Q}_5$	<span>the</span> ↓ $\vec{E}_6$ ↓ $W_Q$ $\vec{Q}_6$	<span>verdant</span> ↓ $\vec{E}_7$ ↓ $W_Q$ $\vec{Q}_7$	<span>forest</span> ↓ $\vec{E}_8$ ↓ $W_Q$ $\vec{Q}_8$	
<span>a</span> → $\vec{E}_1$ → $W_k$ $\vec{K}_1$	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
<span>fluffy</span> → $\vec{E}_2$ → $W_k$ $\vec{K}_2$	0.00	1.00	0.00	0.42	0.00	0.00	0.00	0.00	
<span>blue</span> → $\vec{E}_3$ → $W_k$ $\vec{K}_3$	0.00	0.00	1.00	0.58	0.00	0.00	0.00	0.00	
<span>creature</span> → $\vec{E}_4$ → $W_k$ $\vec{K}_4$	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
<span>roamed</span> → $\vec{E}_5$ → $W_k$ $\vec{K}_5$	0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.00	
<span>the</span> → $\vec{E}_6$ → $W_k$ $\vec{K}_6$	0.00	0.00	0.00	0.00	0.99	1.00	0.00	0.00	
<span>verdant</span> → $\vec{E}_7$ → $W_k$ $\vec{K}_7$	0.00	0.00	0.00	0.00	0.00	0.00	1.00	1.00	

# Attention matrix size is $N^2$ of the context window!



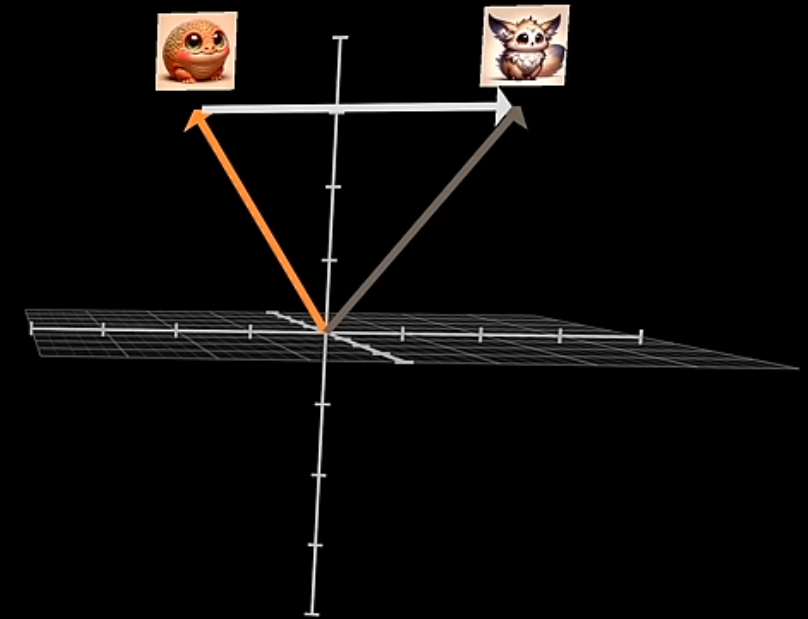
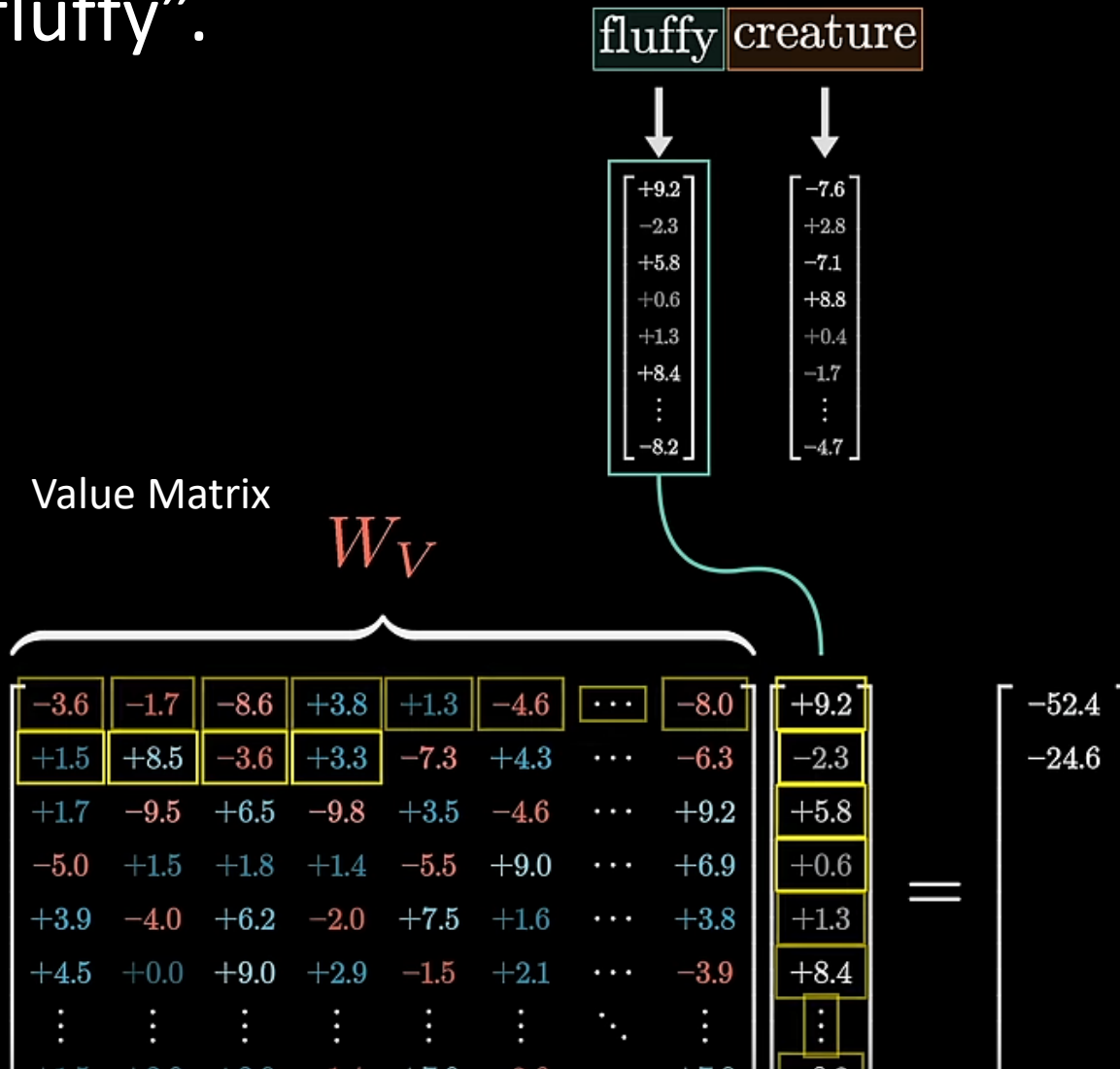


At this moment, I have a column of numbers telling me how important another word is to the word of interest (WOI).

	a	fluffy	blue	creature	roamed	the	verdant	forest	
	$\downarrow \vec{E}_1$ $\downarrow^{W_Q} \vec{Q}_1$	$\downarrow \vec{E}_2$ $\downarrow^{W_Q} \vec{Q}_2$	$\downarrow \vec{E}_3$ $\downarrow^{W_Q} \vec{Q}_3$	$\downarrow \vec{E}_4$ $\downarrow^{W_Q} \vec{Q}_4$	$\downarrow \vec{E}_5$ $\downarrow^{W_Q} \vec{Q}_5$	$\downarrow \vec{E}_6$ $\downarrow^{W_Q} \vec{Q}_6$	$\downarrow \vec{E}_7$ $\downarrow^{W_Q} \vec{Q}_7$	$\downarrow \vec{E}_8$ $\downarrow^{W_Q} \vec{Q}_8$	
$\boxed{\text{a}} \rightarrow \vec{E}_1 \xrightarrow{W_k} \vec{K}_1$	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
$\boxed{\text{fluffy}} \rightarrow \vec{E}_2 \xrightarrow{W_k} \vec{K}_2$	0.00	1.00	0.00	0.42	0.00	0.00	0.00	0.00	
$\boxed{\text{blue}} \rightarrow \vec{E}_3 \xrightarrow{W_k} \vec{K}_3$	0.00	0.00	1.00	0.58	0.00	0.00	0.00	0.00	
$\boxed{\text{creature}} \rightarrow \vec{E}_4 \xrightarrow{W_k} \vec{K}_4$	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	
$\boxed{\text{roamed}} \rightarrow \vec{E}_5 \xrightarrow{W_k} \vec{K}_5$	0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.00	
$\boxed{\text{the}} \rightarrow \vec{E}_6 \xrightarrow{W_k} \vec{K}_6$	0.00	0.00	0.00	0.00	0.99	1.00	0.00	0.00	
$\boxed{\text{verdant}} \rightarrow \vec{E}_7 \xrightarrow{W_k} \vec{K}_7$	0.00	0.00	0.00	0.00	0.00	0.00	1.00	1.00	

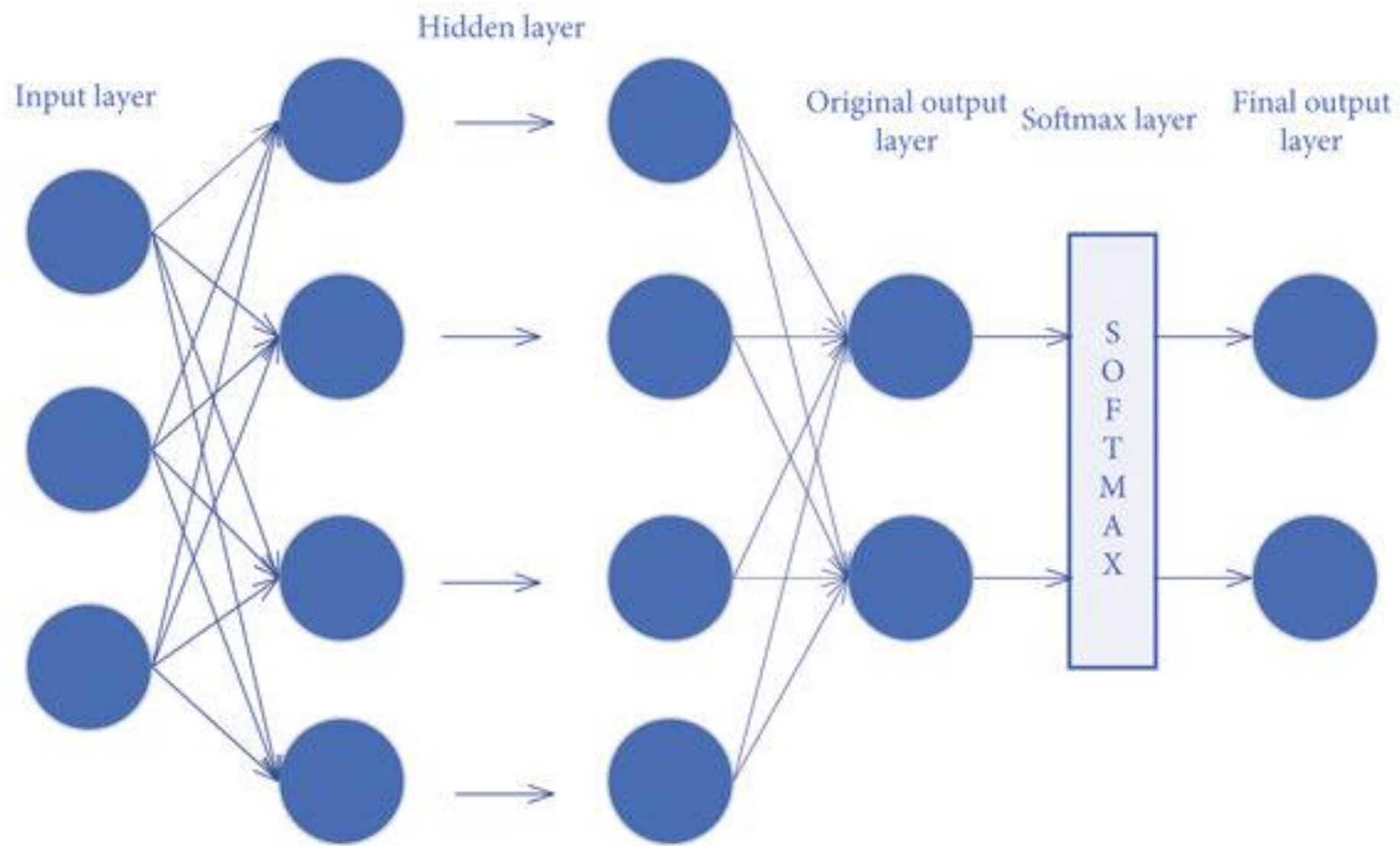
- What I don't know is which direction I should move the WOI.
- Enter the Value matrix.


The Value Matrix for “creature” is being updated with the column from “fluffy”.





- Having “transformed” a word to its new position, we now ask “What word is most likely to follow you?”
- You give the 12,288-element vector (“transformed” word embedding) to yet another neural network that has 50,257 outputs.

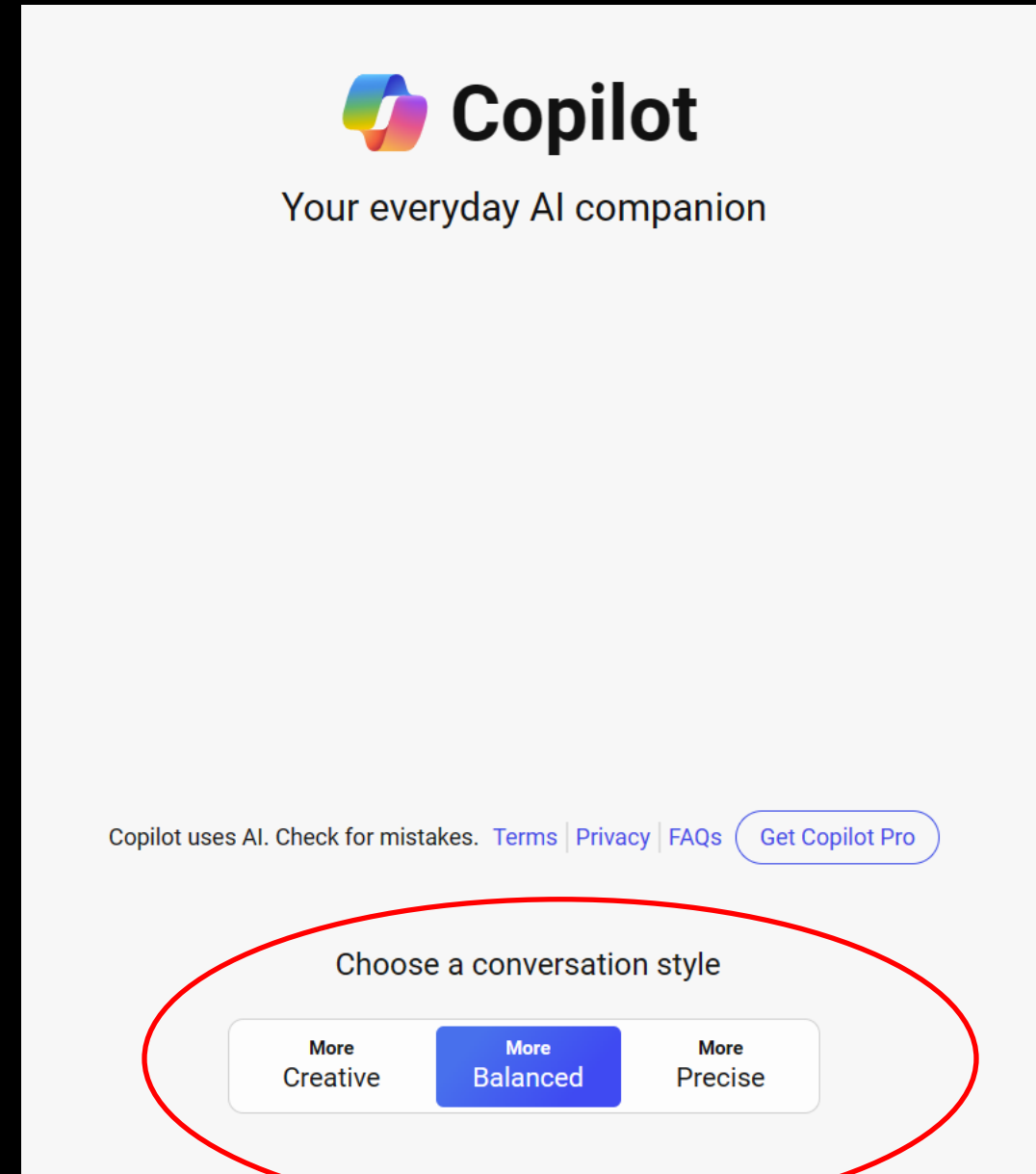




Each output is a number between 0..1 representing how probable it is that the vocabulary word is the next word.

# Temperature

- They don't always take the most probable next word.



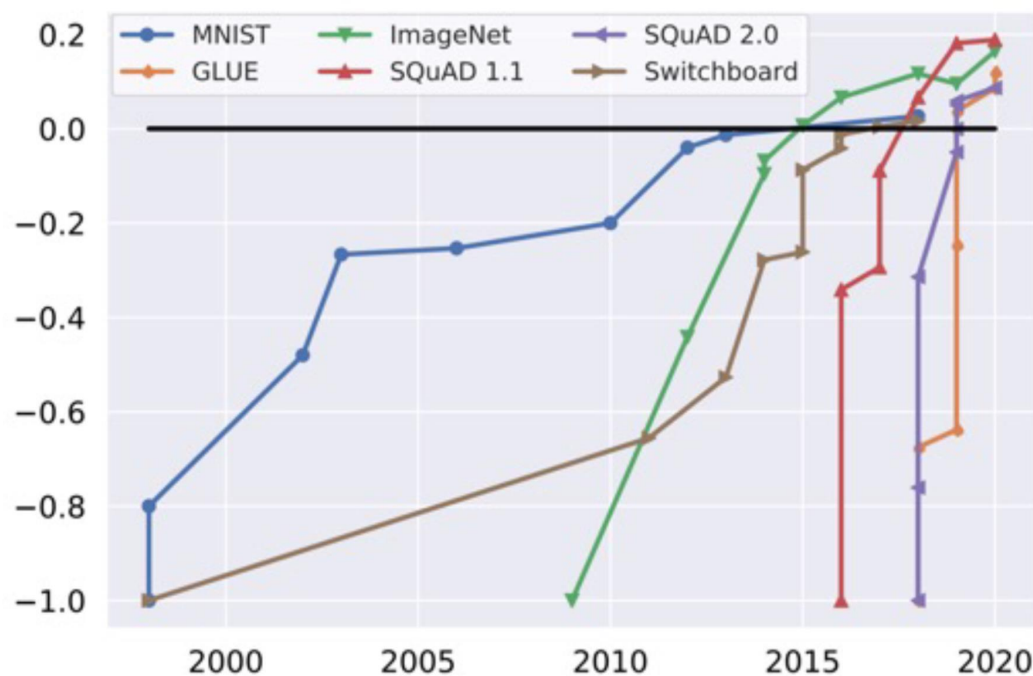






**Figure 1:** ML-benchmark saturation relative to human performance (black line) [Kielia et al., 2021].

**SQuAD2.0** combines the 100,000 questions in SQuAD1.1 with over 50,000 unanswerable questions written adversarially by crowdworkers to look similar to answerable ones. To do well on SQuAD2.0, systems must not only answer questions when possible, but also determine when no answer is supported by the paragraph and abstain from answering.

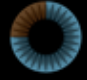


**Figure 1:** ML-benchmark saturation relative to human performance (black line) [Kiela et al., 2021].

The General Language Understanding Evaluation (GLUE) benchmark is a collection of resources for training, evaluating, and analyzing natural language understanding systems.

Switchboard: Recognition of conversational speech over telephone.

# Credits

- A number of images are credit to 3Blue1Brown.com
  - Their logo on the slides helps identify images from them. 
  - Please check out their YouTube channel.