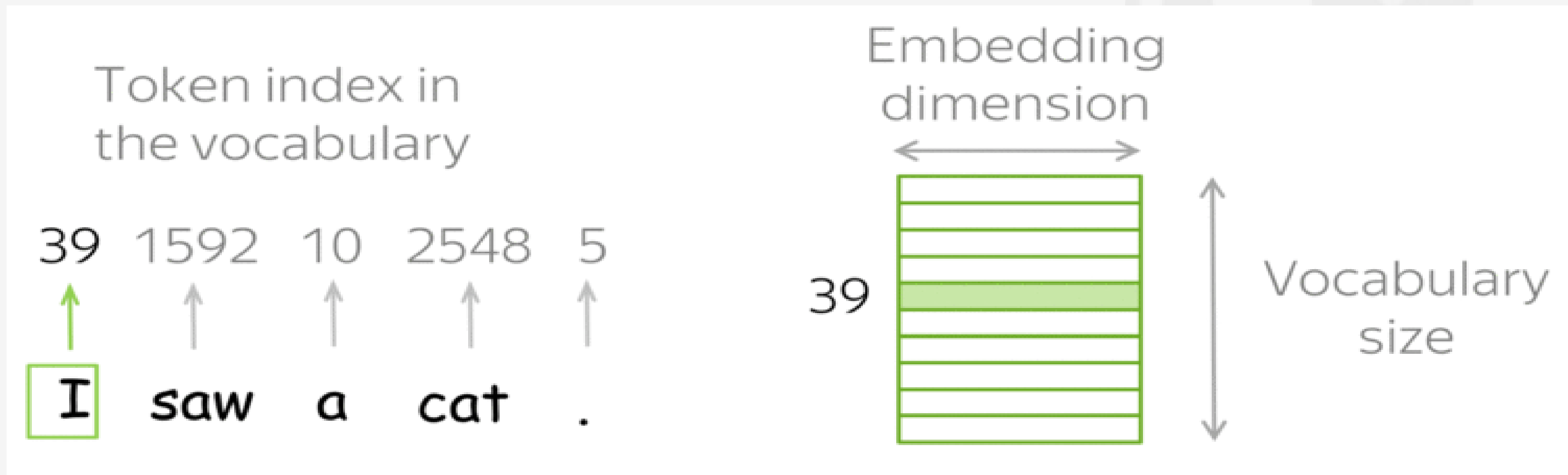


NLP WORKSHOPS WITH LLM INSIGHT

INTRODUCTION TO EMBEDDINGS



Mathematical Representation of words

WHY TEXT EMBEDDING

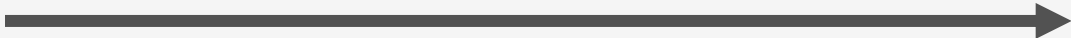
“so do all who
live to see such
times.”

0	0	8	0
2	25	0	7
13	0	12	9
15	32	10	0



ML MODEL

“so do all who
live to see such
times.”



0	0	8	0
2	25	0	7
13	0	12	9
15	32	10	0

- ❑ One-hot encoding
- ❑ Count based representation
- ❑ Embeddings



ONE-HOT ENCODING

Vocabulary = $\left(\begin{array}{cc} \text{I} & \text{Apple} \\ \text{You} & \text{Cat} \\ \text{Bag} & \text{Dog} \end{array} \right)$

I	You	Bag	Apple	Cat	Dog
0	1	2	3	4	5

ONE-HOT ENCODING

Vocabulary = $\left(\begin{array}{l} \text{I} \\ \text{You} \\ \text{Bag} \end{array} \quad \begin{array}{l} \text{Apple} \\ \text{Cat} \\ \text{Dog} \end{array} \right)$

I	You	Bag	Apple	Cat	Dog
0	0	1	0	0	0
0	1	2	3	4	5

BAG OF WORDS

A wizzard is never late.

A	He	Is	When	Never	To	Means	...

He arrives precisely when he means to.

--	--	--	--	--	--	--	--

BAG OF WORDS

A wizard is
never late.

A wizard	wizard is	is never	never late	he arrives	...
1	1	1	0	1	...

He arrives
precisely when
he means to.

0	0	0	0	1	...
---	---	---	---	---	-----

TF-IDF

$$\text{TF} = \frac{\text{number of times this words occurs}}{\text{number of words in the documents}}$$

$$\text{IDF} = \log \frac{\text{number of documents}}{\text{number of documents where this word occurred}}$$

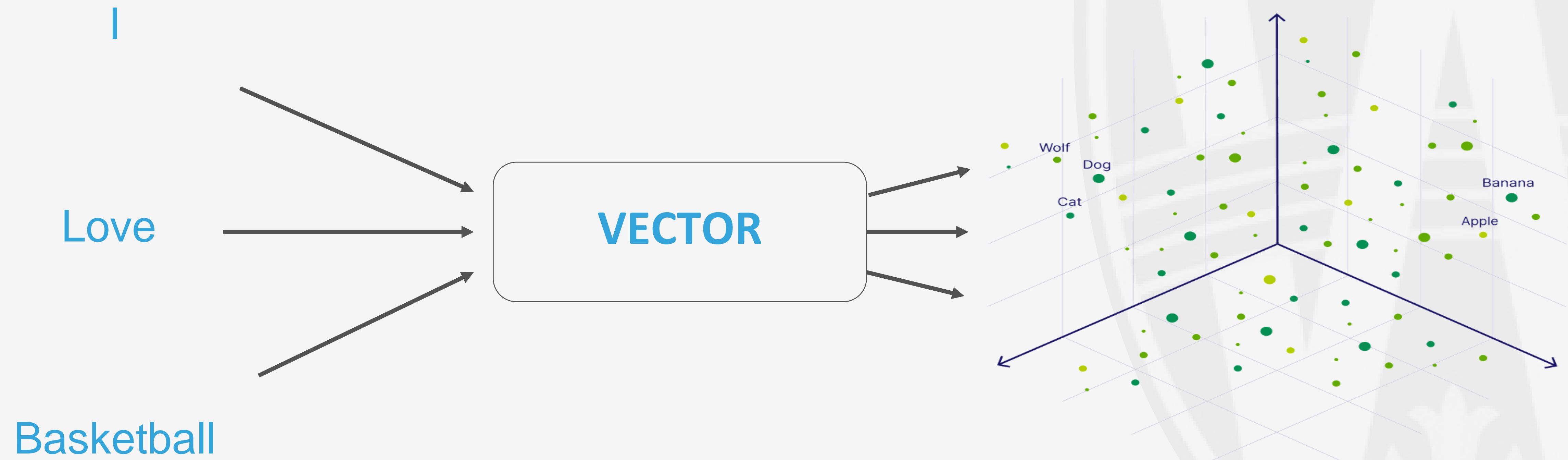
$$\text{TFIDF} = \text{TF} * \text{IDF} + 1$$

A **wizard** is never late, nor is he **early**.

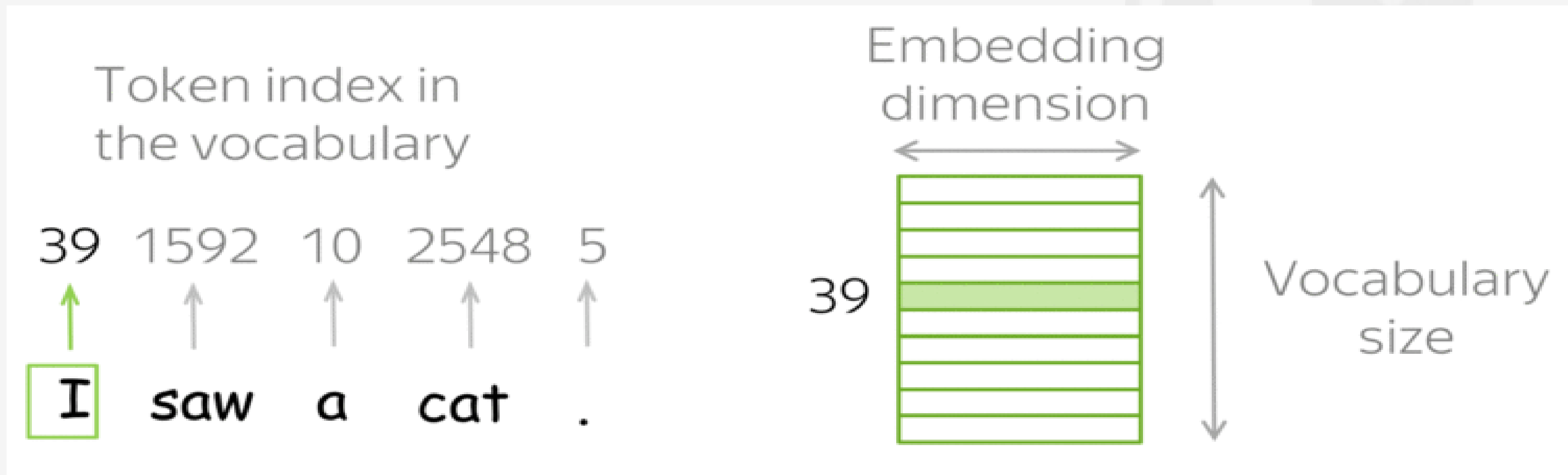
- ❑ No context
- ❑ Unknown words
- ❑ Sparse vectors



EMBEDDINGS



WHAT IS DENSE VECTOR



WHAT IS DENSE VECTOR

Tea



Breakfast

Coffee



Enjoy

Drink

WHAT IS DENSE VECTOR

Tea



Breakfast

Drink

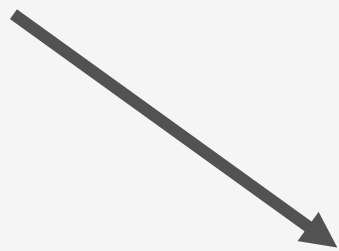
Pea



Enjoy

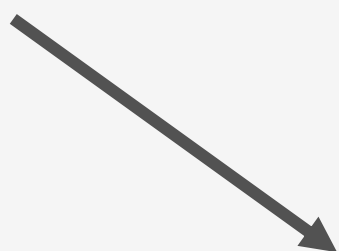
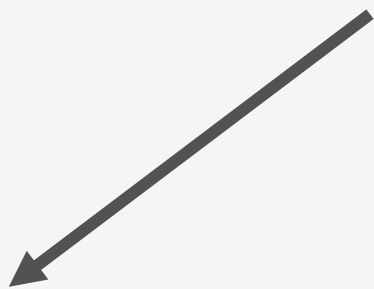
WHAT IS DENSE VECTOR

Tea



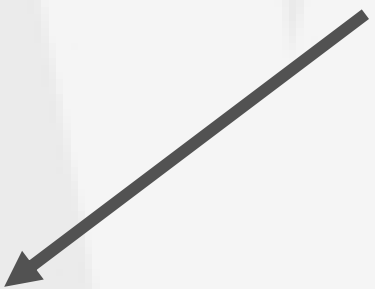
Distance
0.3

Pea



Distance
0.7

Coffee



WHAT IS EMBEDDINGS SPACE

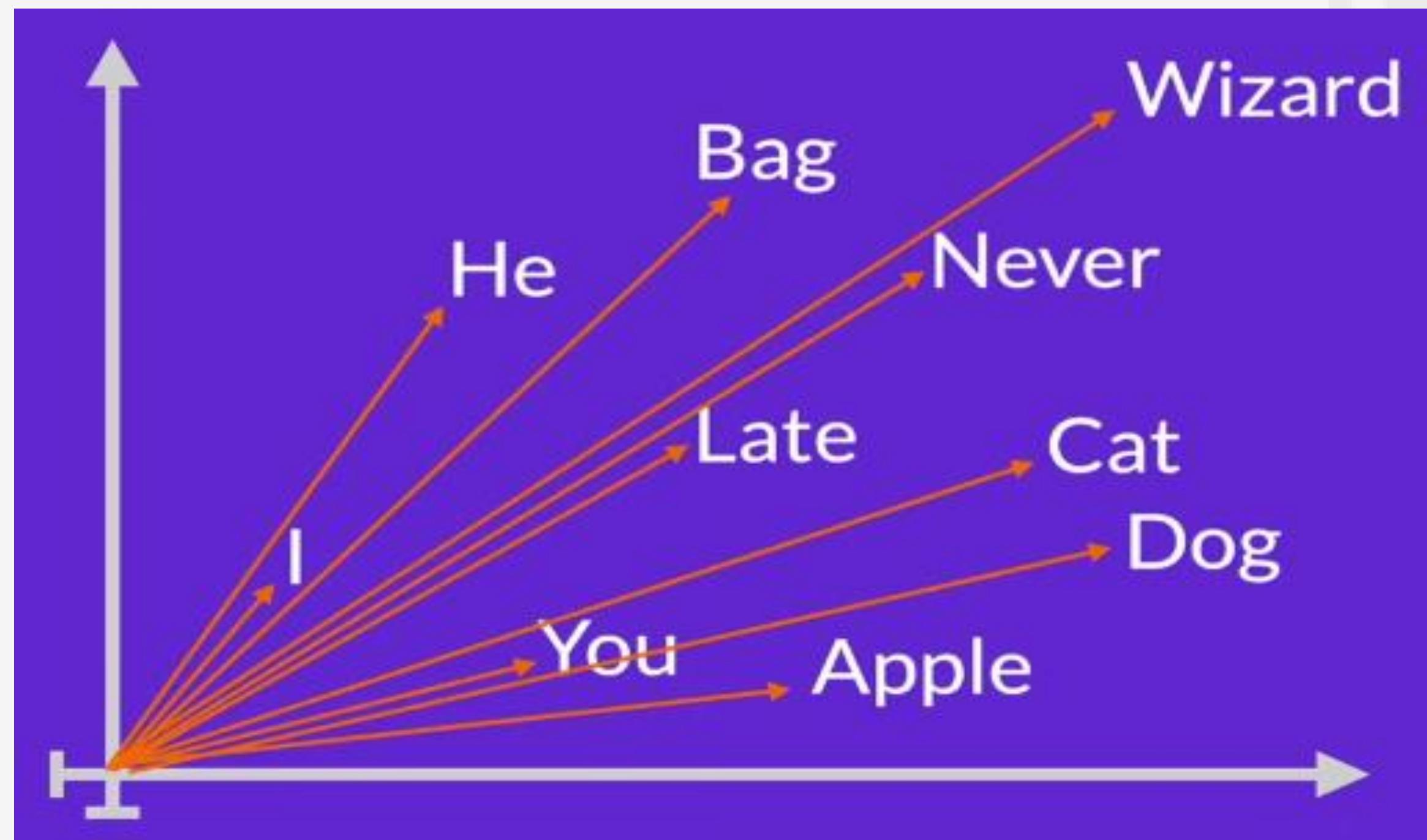
I	You	He	Late	Never	Bag	Wizard	Apple	Cat	Dog
2	11	3	7	13	5	18	8	20	23



Similarity
 $18 - 8 = 10$

WHAT IS EMBEDDINGS SPACE

I	You	He	Late	Never	Bag	Wizard	Apple	Cat	Dog
[2,5]	[11,3]	[3,8]	[7,5]	[13,10]	[5,8]	[18,20]	[8,4]	[20,2]	[23,1]

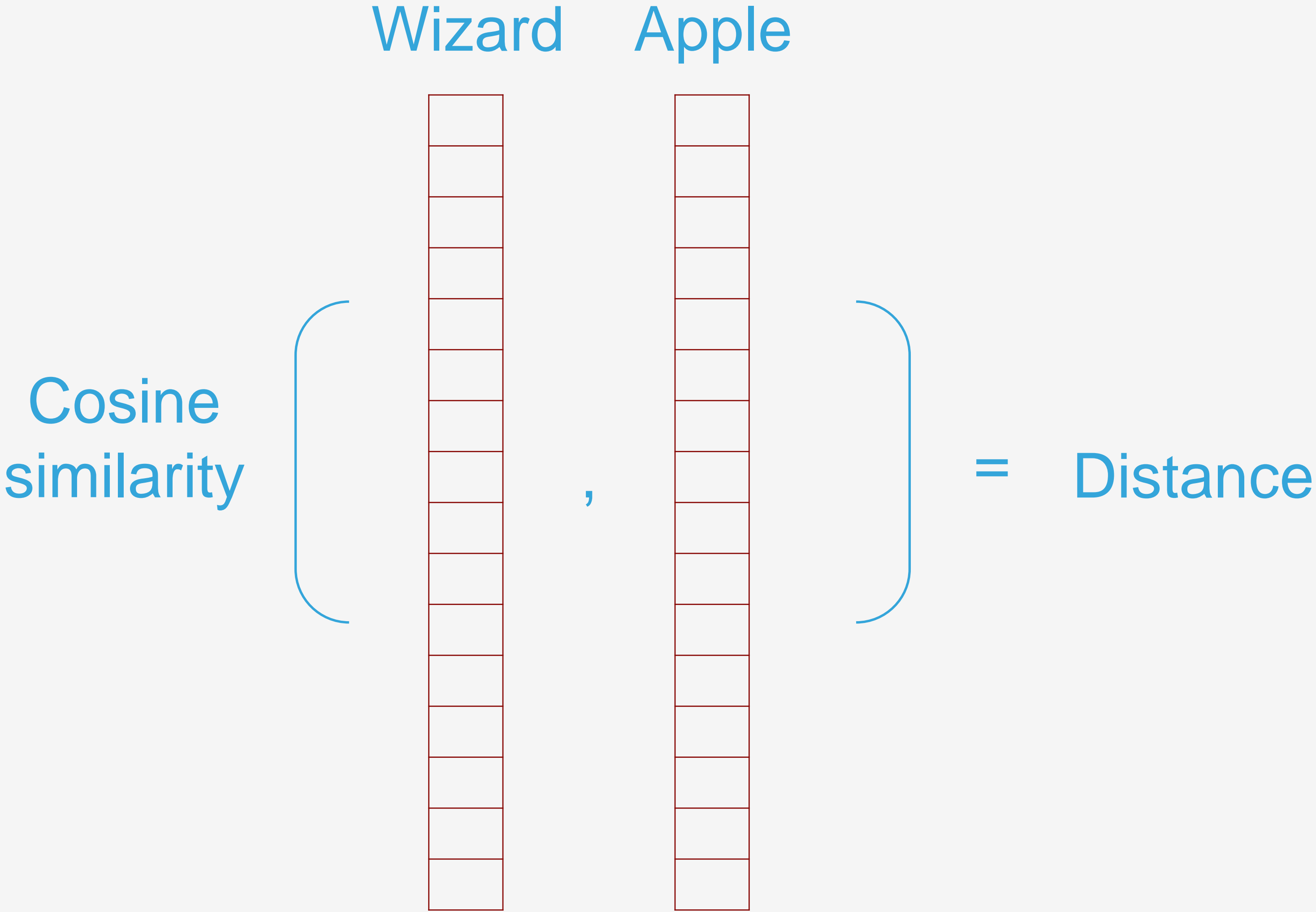


WHAT IS EMBEDDINGS SPACE

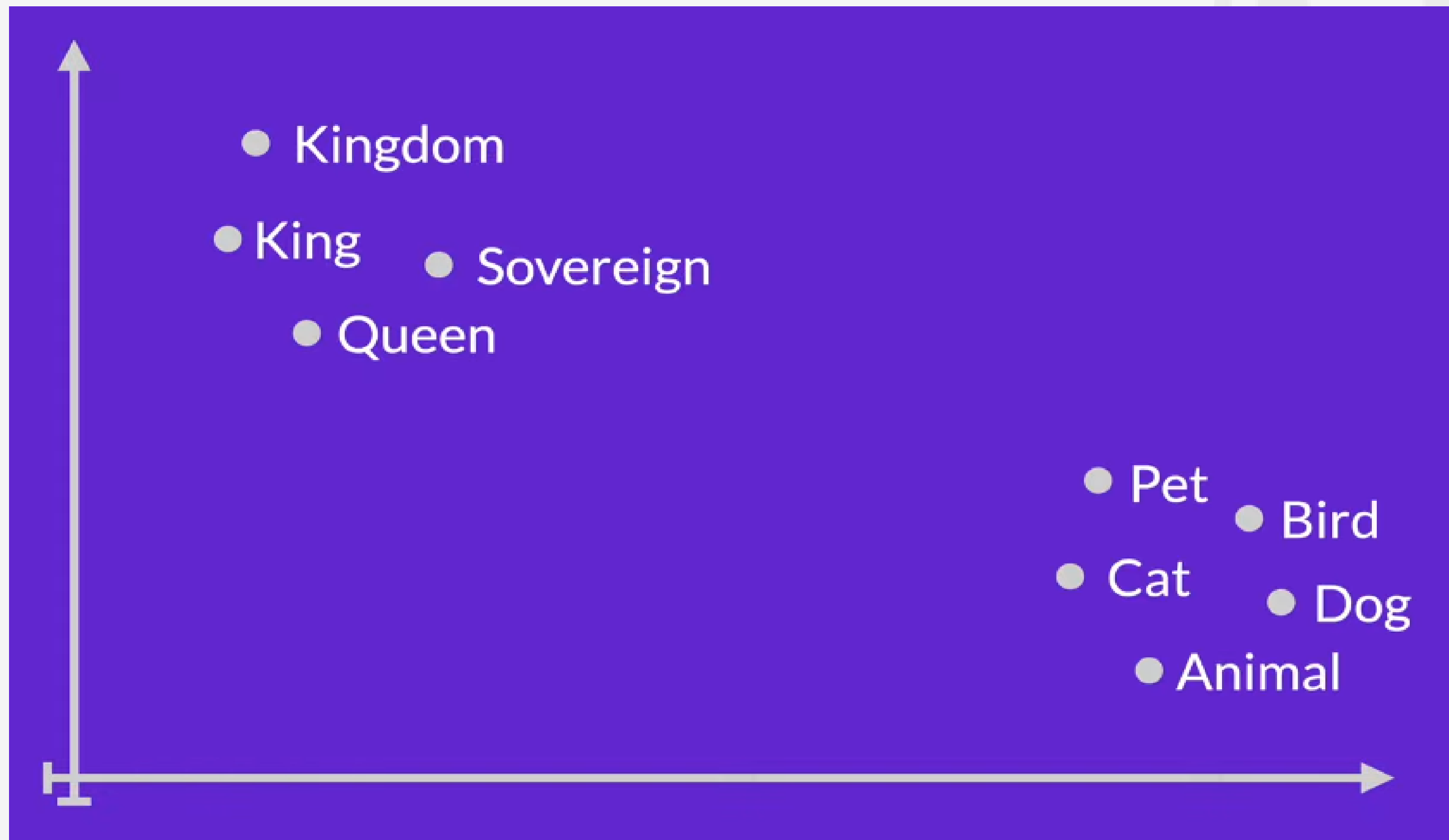
I You He Late Never Bag Wizard Apple Cat Dog

[illegible]

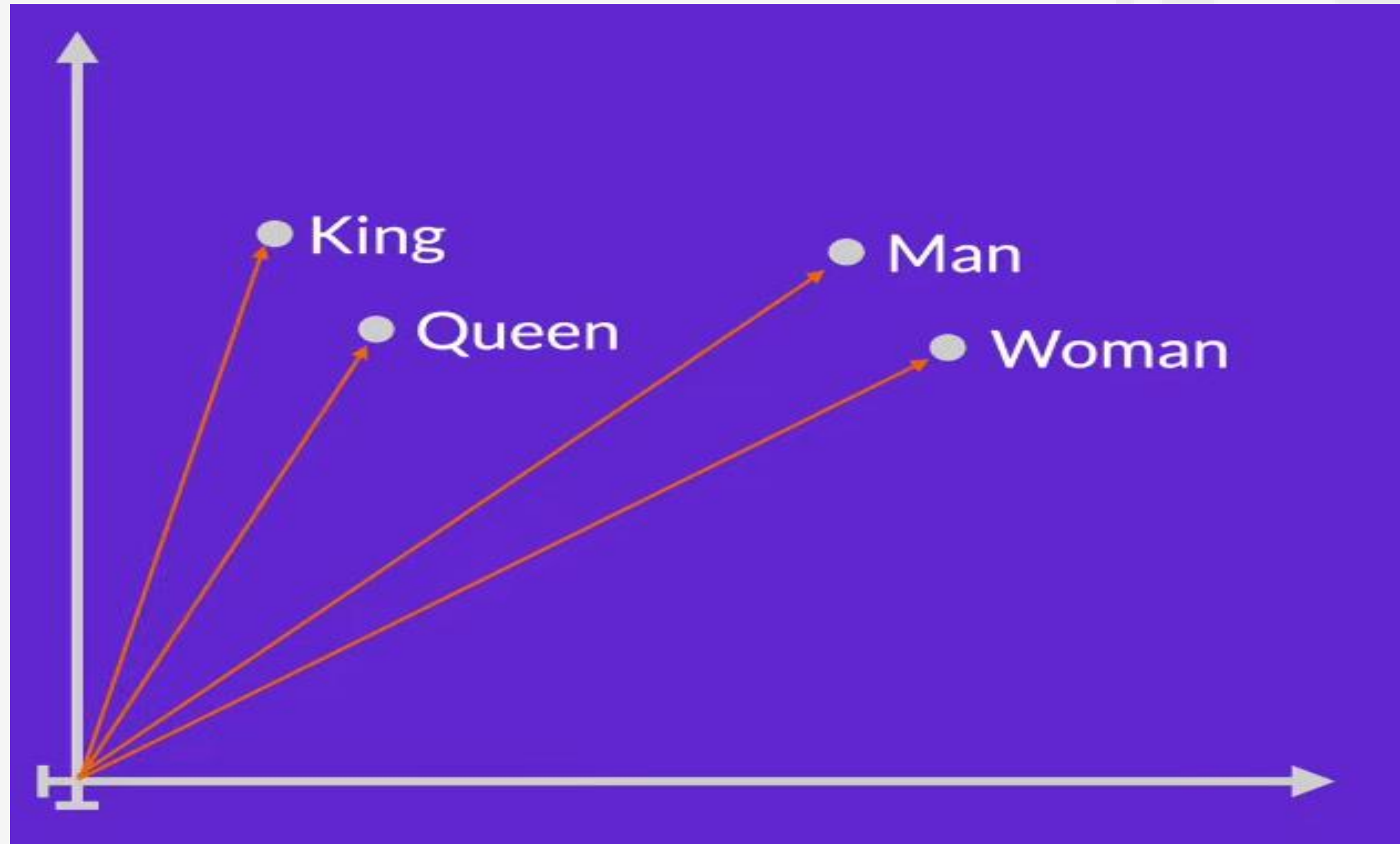
WHAT IS EMBEDDINGS SPACE



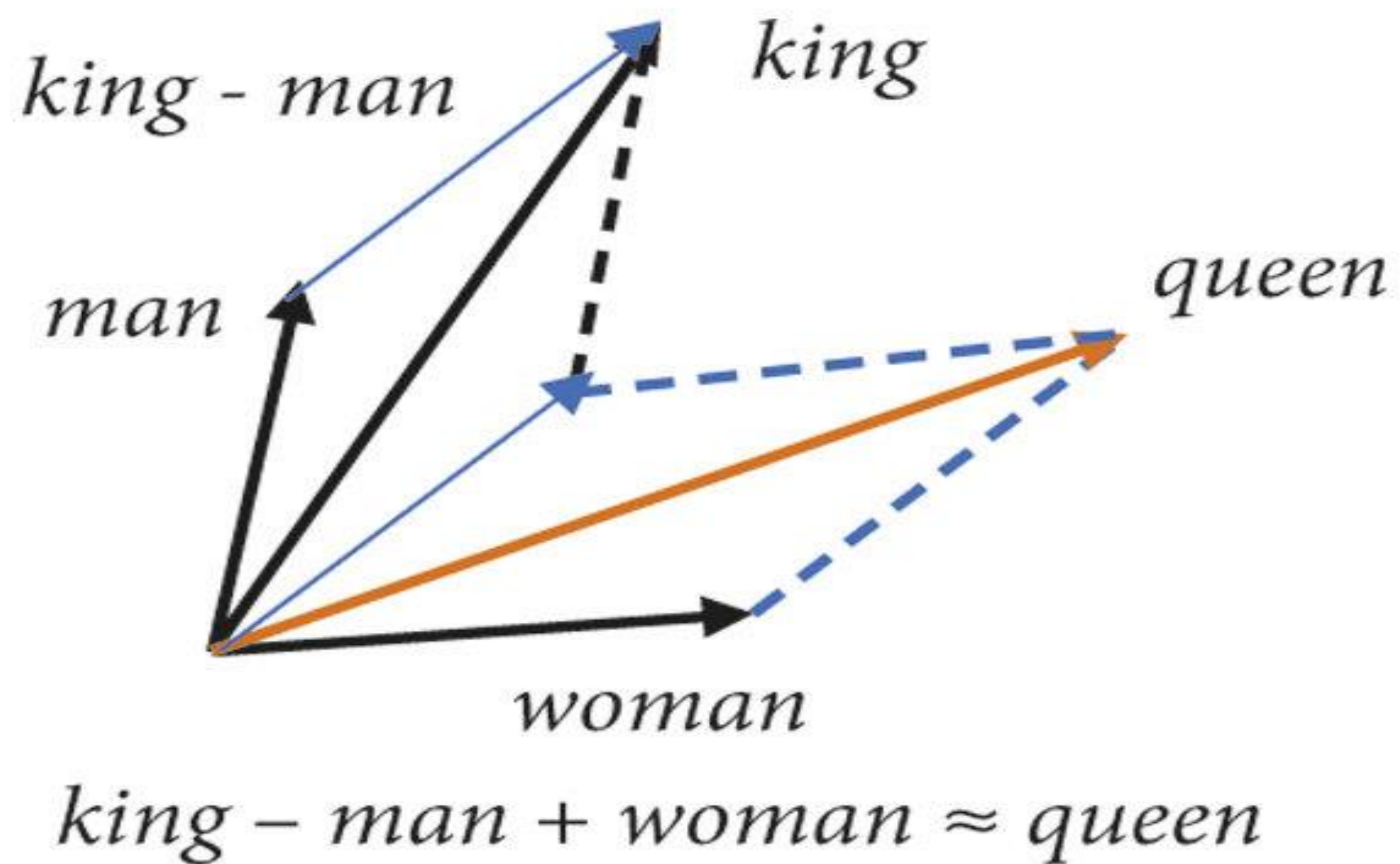
WHAT IS EMBEDDINGS SPACE



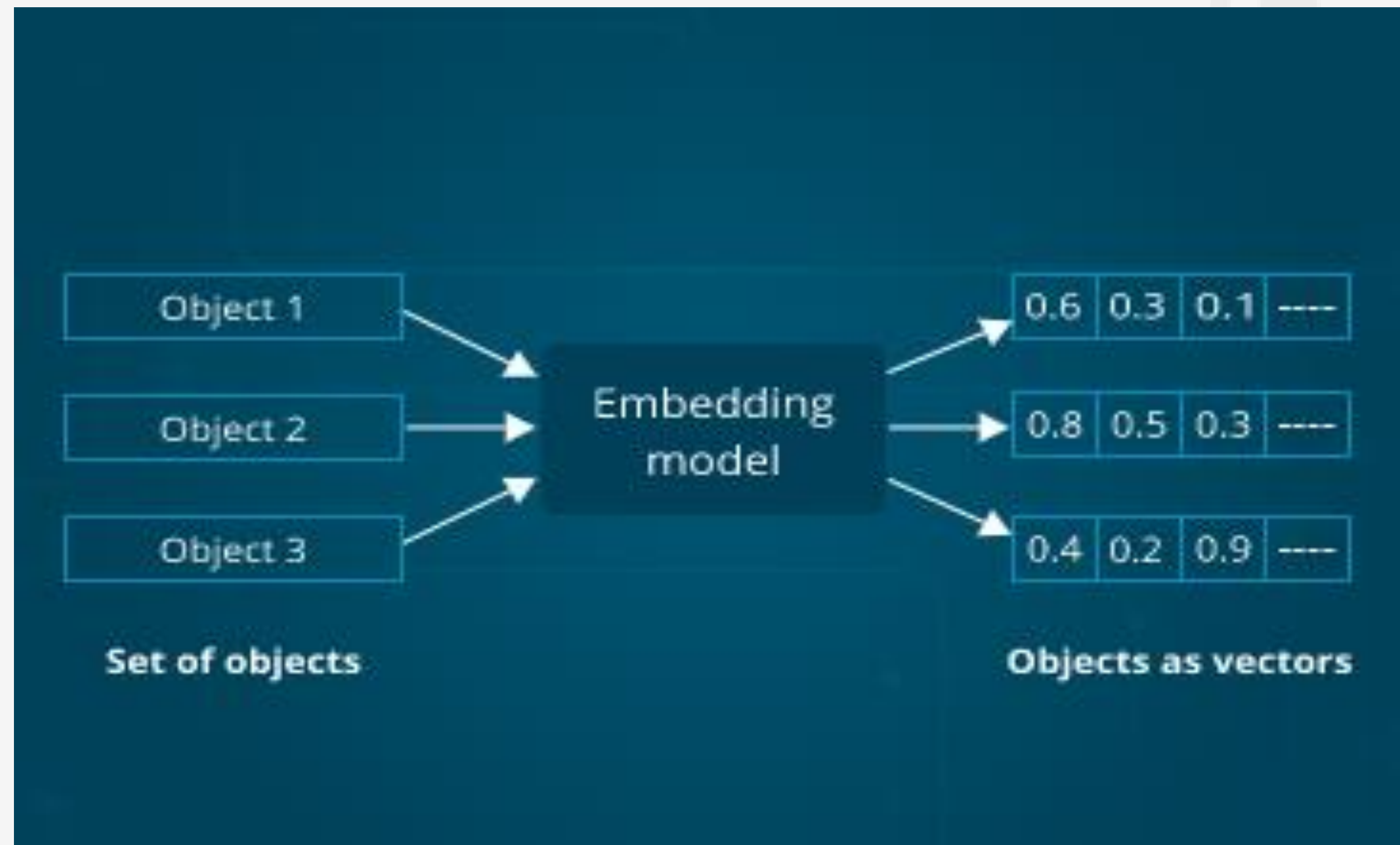
WHAT IS EMBEDDINGS SPACE



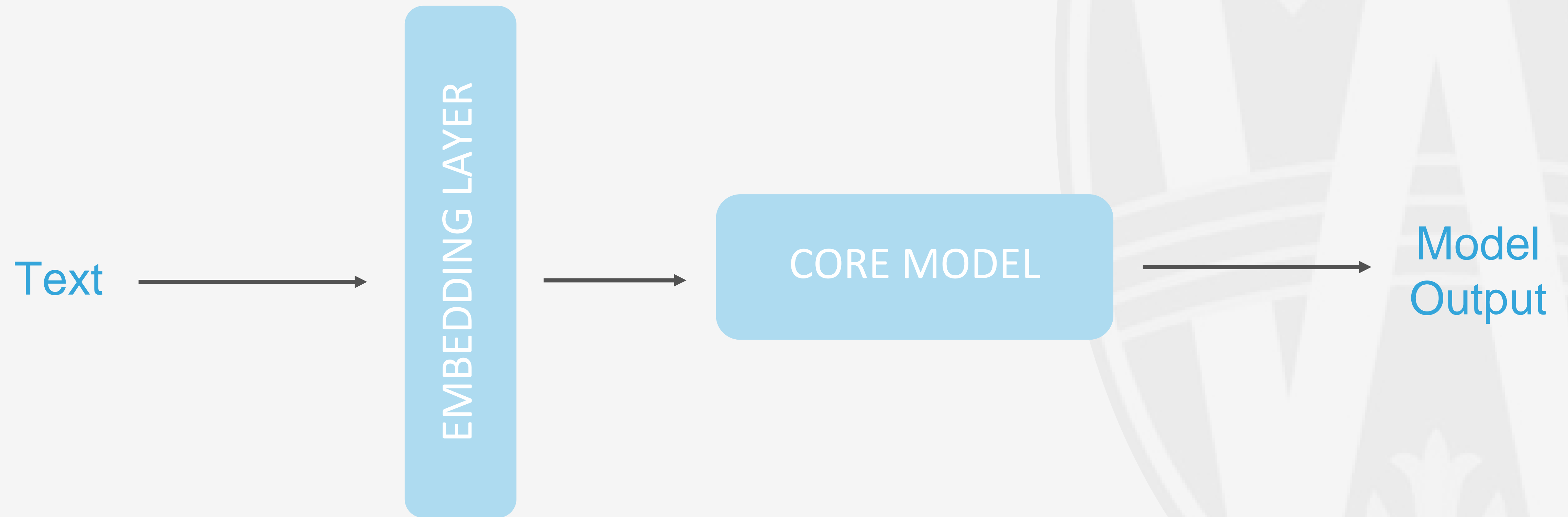
WHAT IS EMBEDDINGS SPACE



HOW ARE WORD EMBEDDINGS MADE



HOW ARE WORD EMBEDDINGS MADE



HOW ARE WORD EMBEDDINGS MADE

- **Pros:**

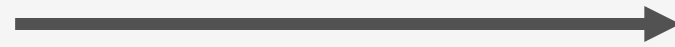
- + specialized embedding model

- **Cons:**

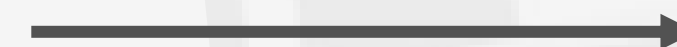
- Huge training set

WORD2VEC

One-hot
Encoded word

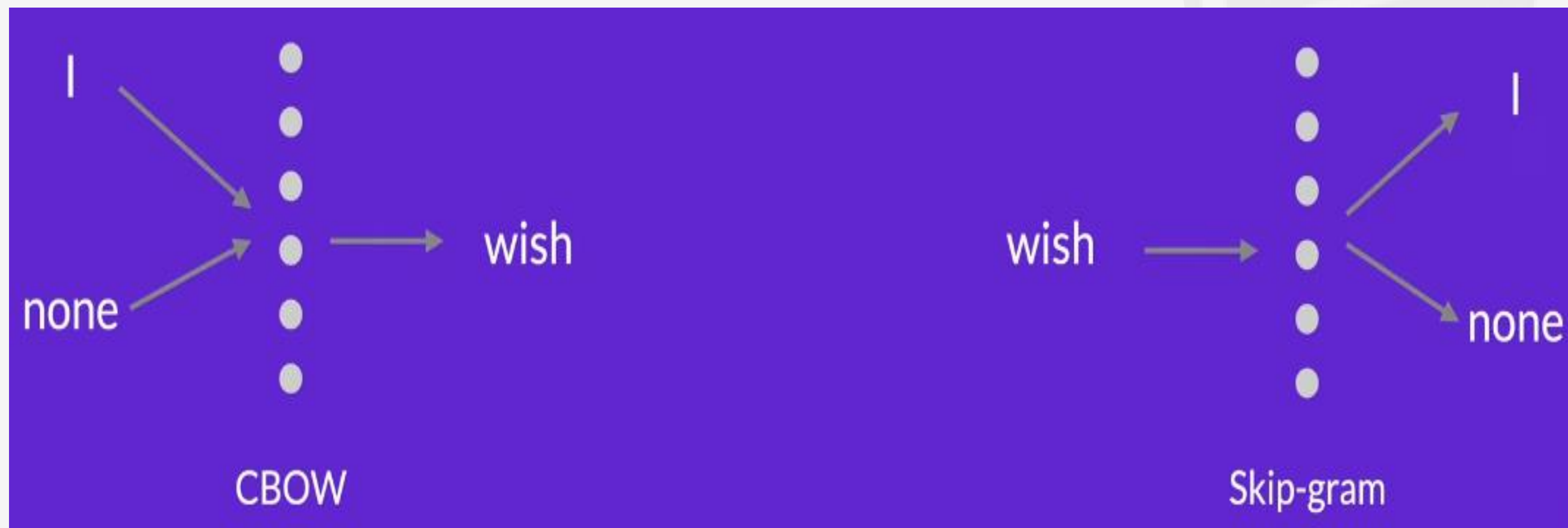


WORD2VEC



Word
Embedding

... I wish none of this has happened. So do all who live to see such times. But that is not up to us to decide. All we have to decide is what to do with the time that is given to us...



Local Dependencies

+

Global Context

Word2vec

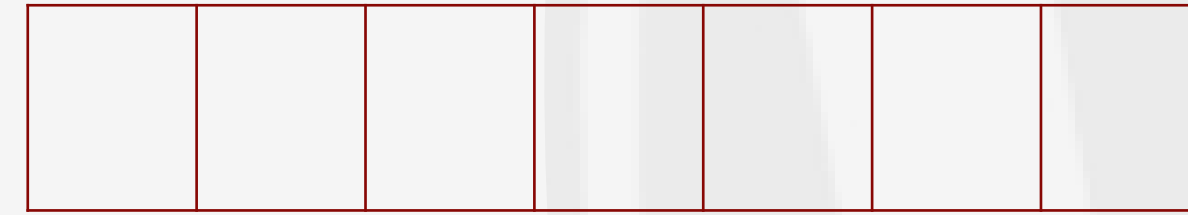
	Tea	Coffee	Pea	Apple	Monkey	App
Tea	1	6	0	2	1	0
Coffee	6	1	0	1	1	1
Pea	0	0	1	0	2	0
Apple	2	1	0	1	2	0
Monkey	1	1	2	2	1	1
App	0	1	0	0	1	1

GLOBAL VECTORS (GLOVES)

Tea



Coffee



$$\text{Tea} \bullet \text{Coffee} = \log \left(\begin{array}{l} \text{Likelihood of tea and coffee} \\ \text{being in the same sentence} \end{array} \right)$$

RPM Feedback - Zahra

