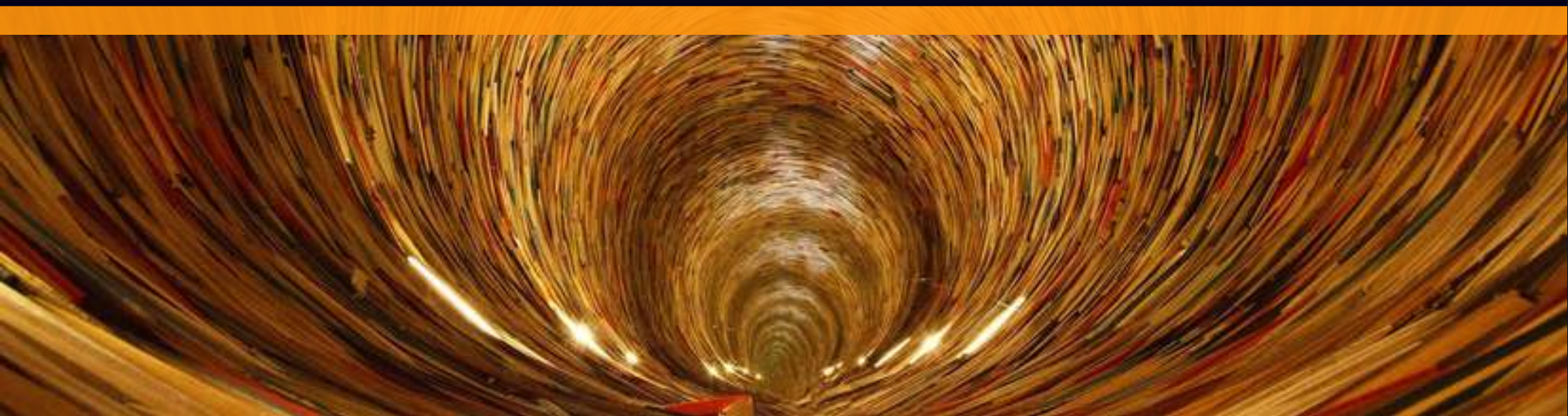


Neural Vector Representations beyond Words: Sentence and Document Embeddings

Gerard de Melo

<http://gerard.demelo.org>

Rutgers University



Gerard de Melo



**Assistant Professor
at Rutgers University**

Head of Deep Data Lab



Vector Representations

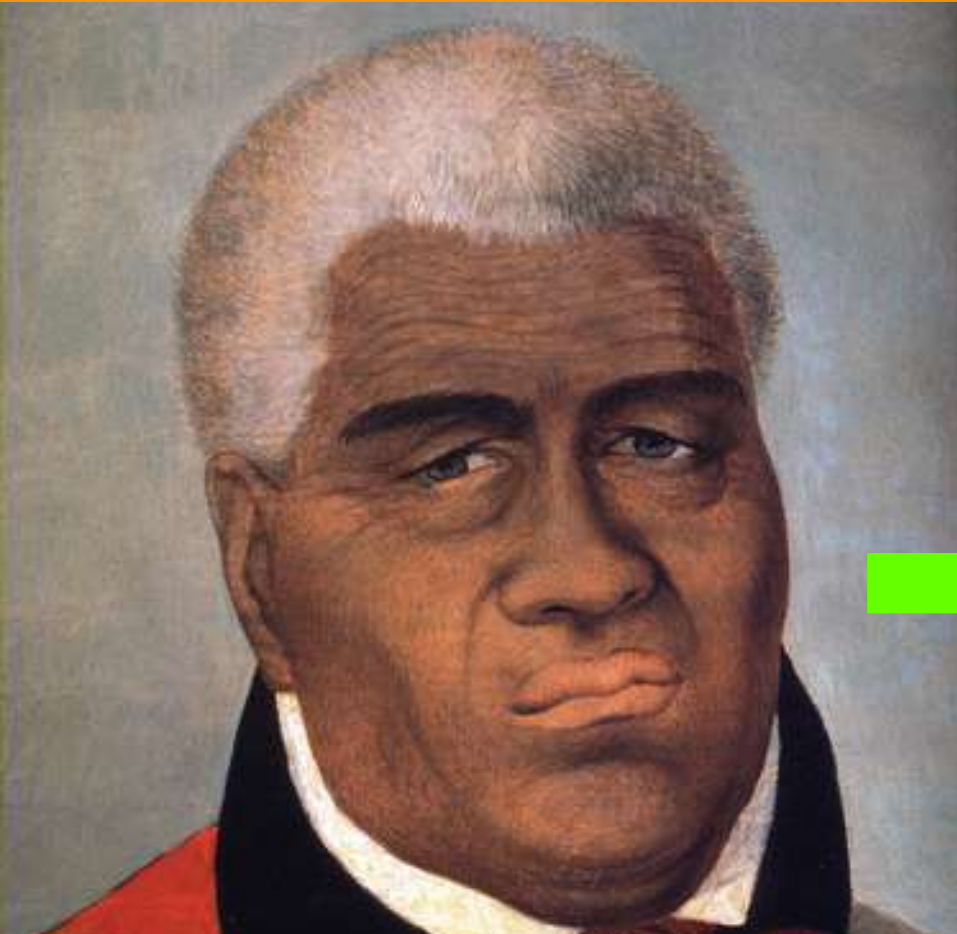


“ants”



\vec{x}_i
0.23
1.21
0.76
0.00
0.12

Vector Representations



“King
Kamehameha I”



\vec{x}_i
0.33
0.53
0.37
1.00
0.73

Vector Representations

“Octel said the purchase was expected.”



\vec{x}_i
0.33
0.53
0.37
1.00
0.73

Vector Representations



\vec{x}_i
0.33
0.53
0.37
1.00
0.73

Why Vectors Representations?



1. Facilitate Machine Learning

Incorrect

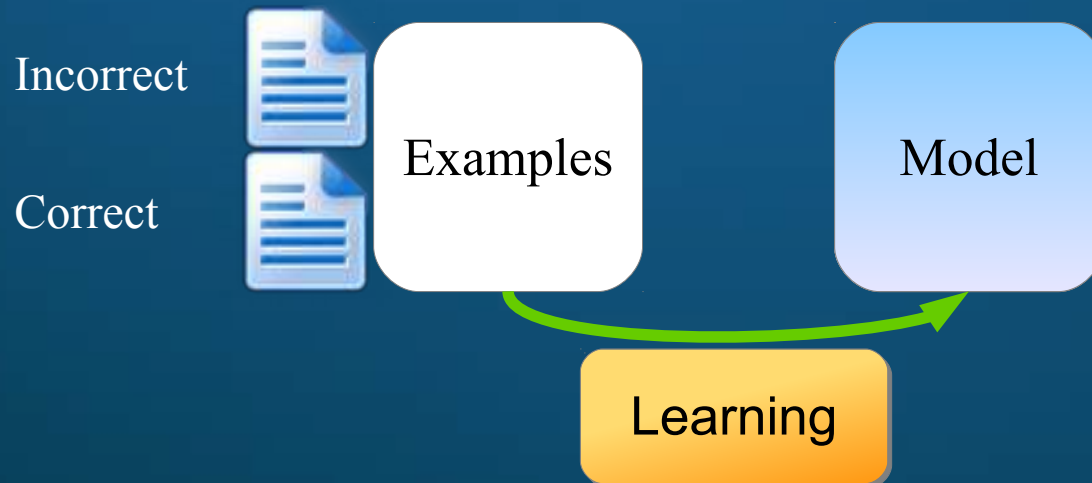


Correct

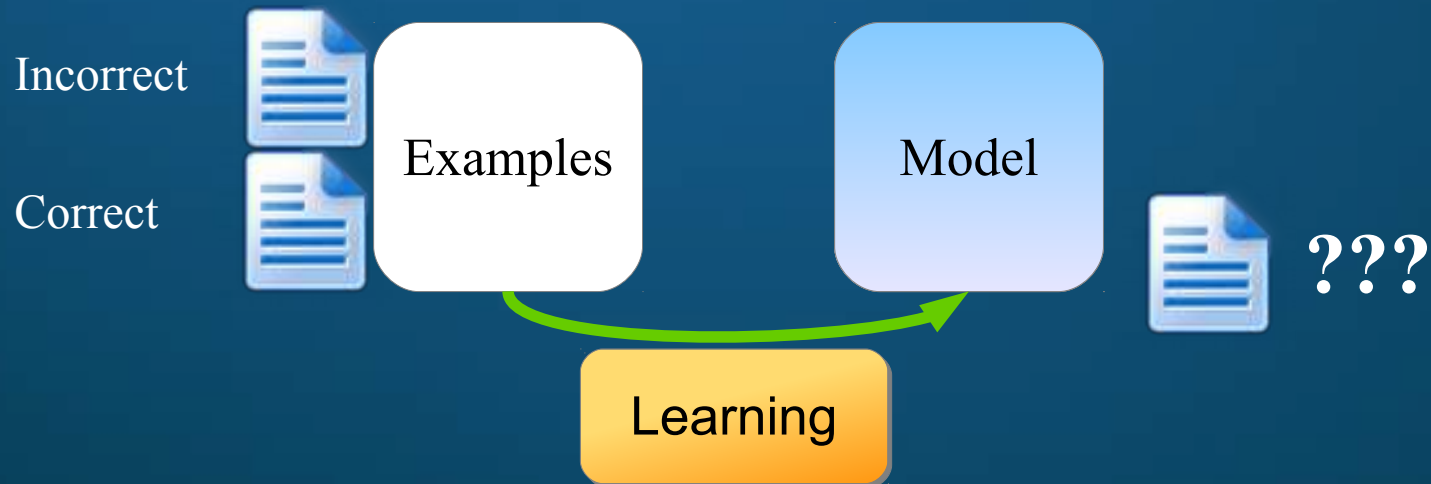


Examples

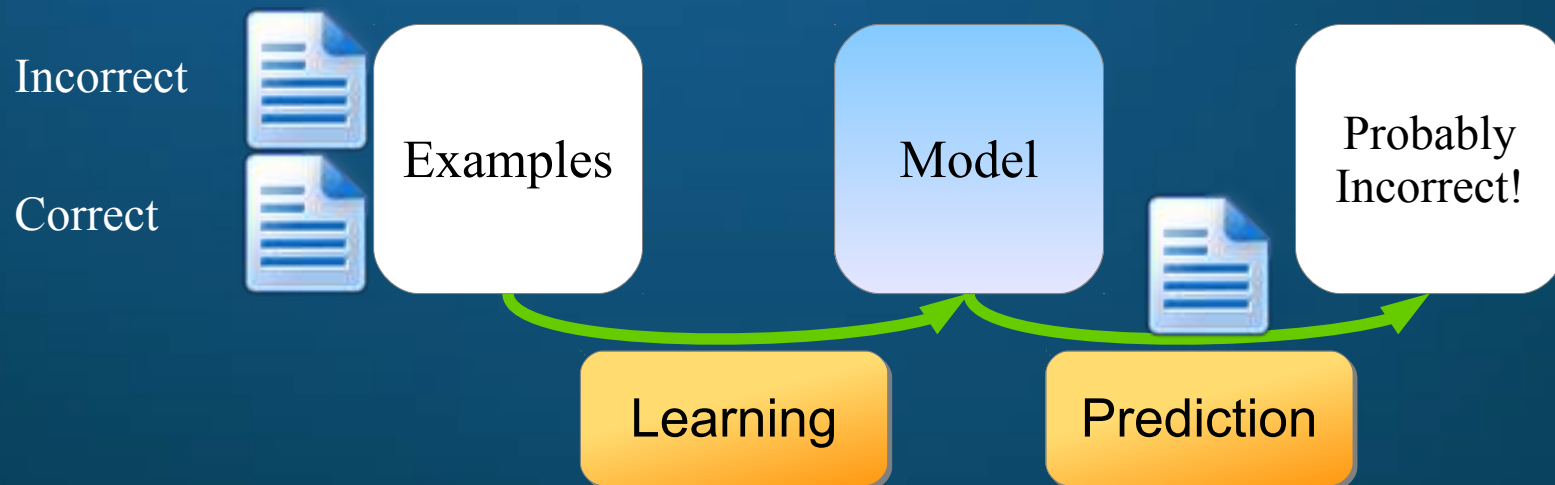
1. Facilitate Machine Learning



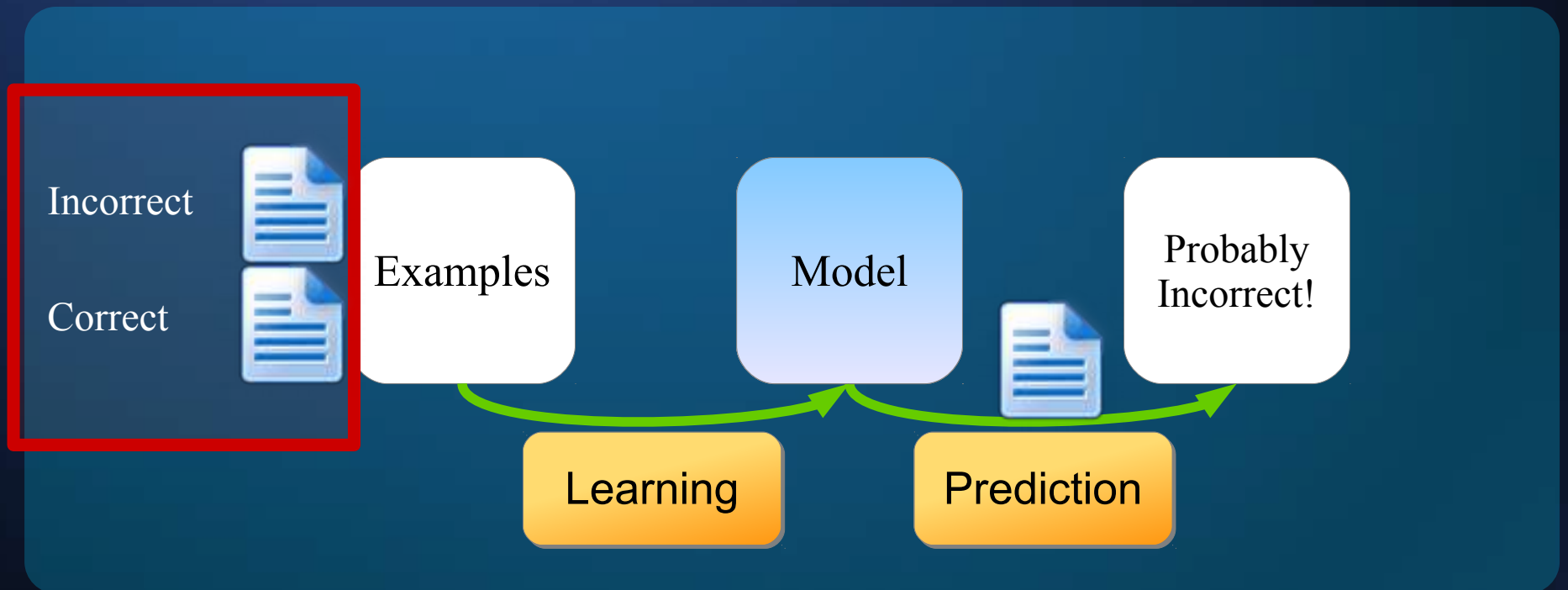
1. Facilitate Machine Learning



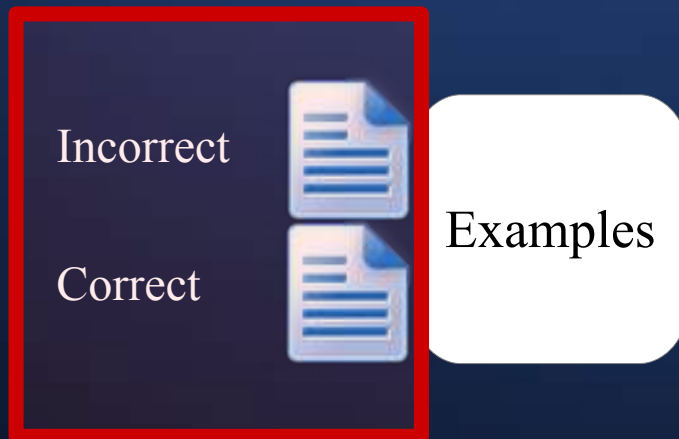
1. Facilitate Machine Learning



1. Facilitate Machine Learning



1. Facilitate Machine Learning



Problem:
Training Data usually
requires human work,
which is slow!

Example:
Penn Chinese TreeBank

2 years for
4000 sentences

New effort required
for each new **language**
and **domain**
(e.g. news vs. tweets
vs. biomedical text)



1. Facilitate Machine Learning



Raw Text

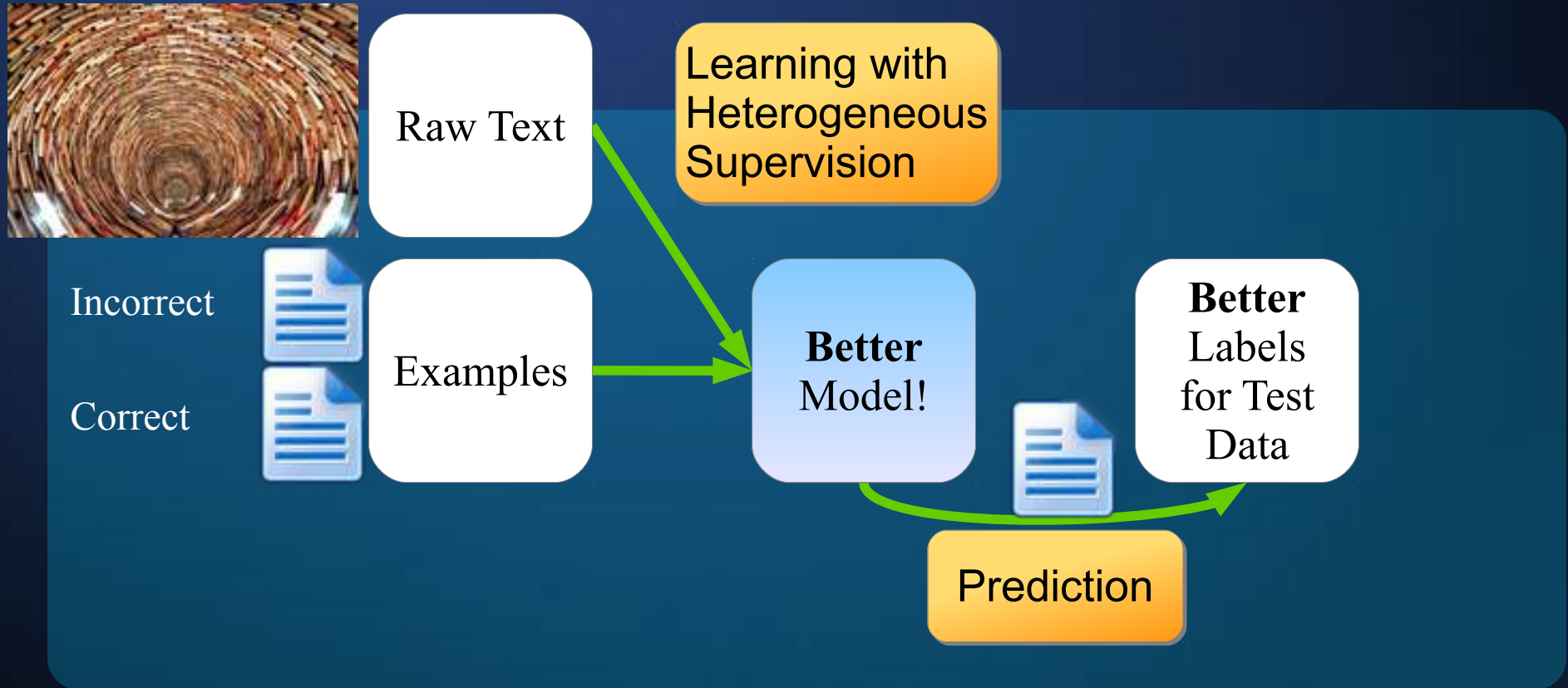
Model

Labels
for Test
Data

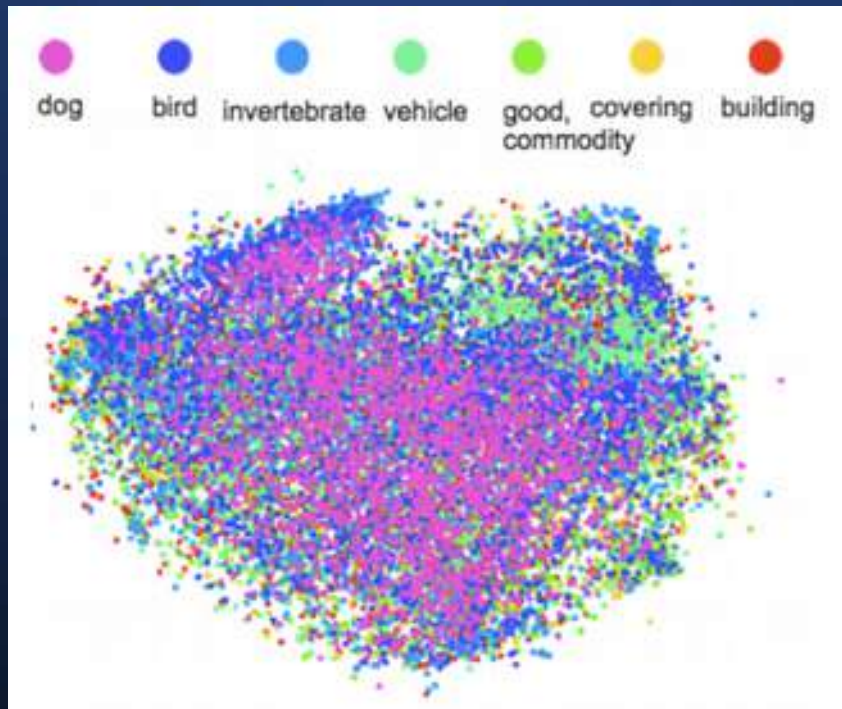
Prediction



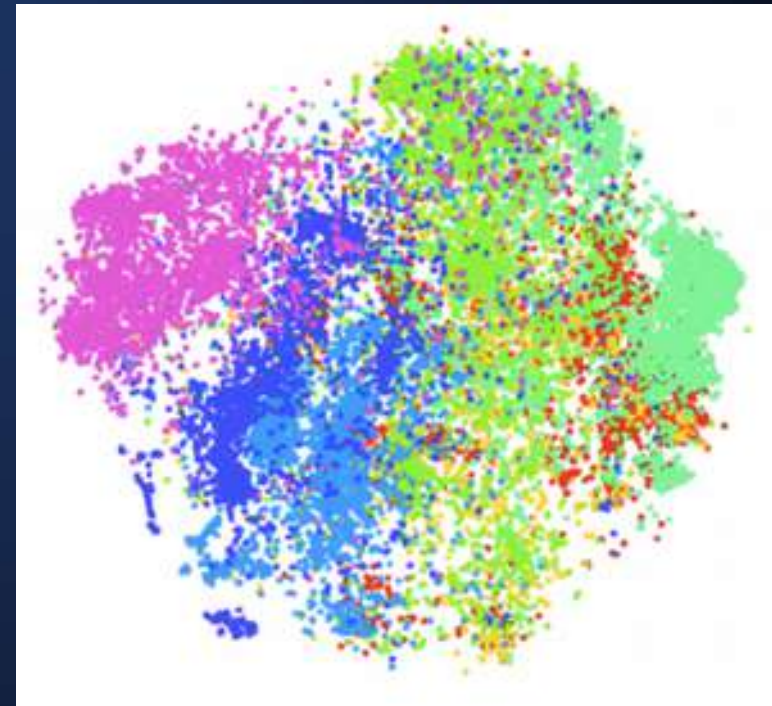
1. Facilitate Machine Learning



1. Facilitate Machine Learning

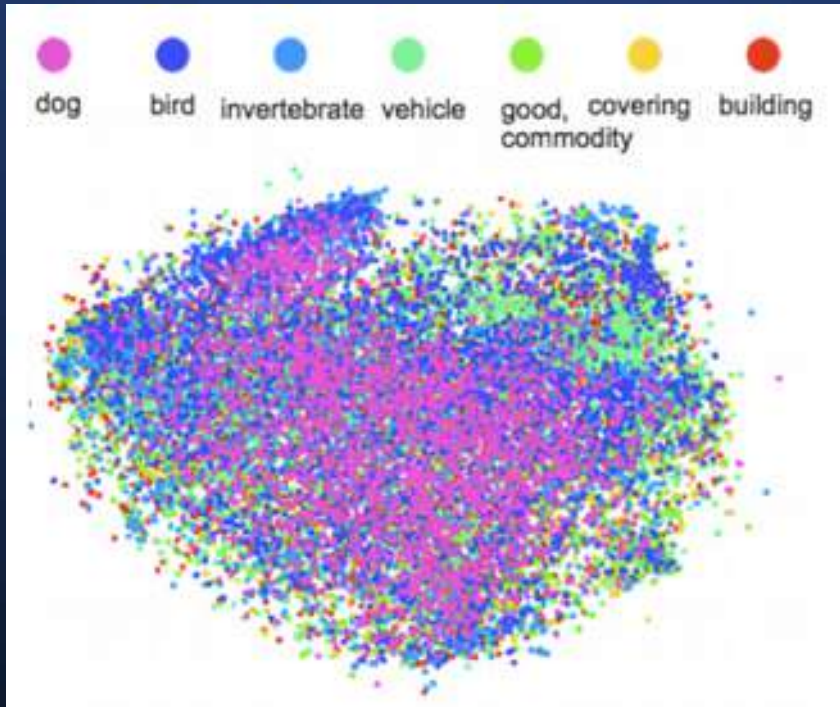


Vectors with
raw/low-level
data
(e.g. pixels)

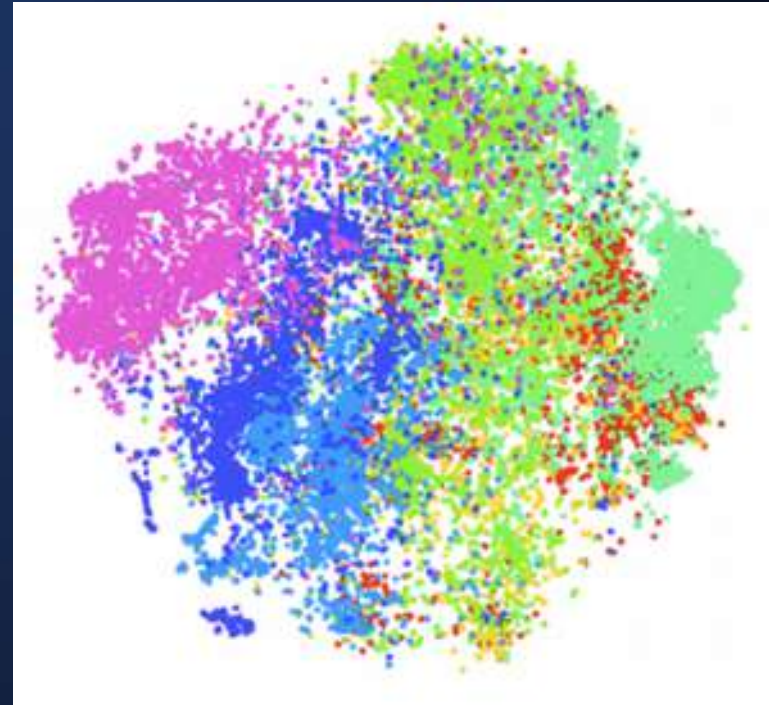


Vectors with
higher-level
representations
(DeCAF)

1. Facilitate Machine Learning



Vectors with
raw/low-level
data
(e.g. pixels)



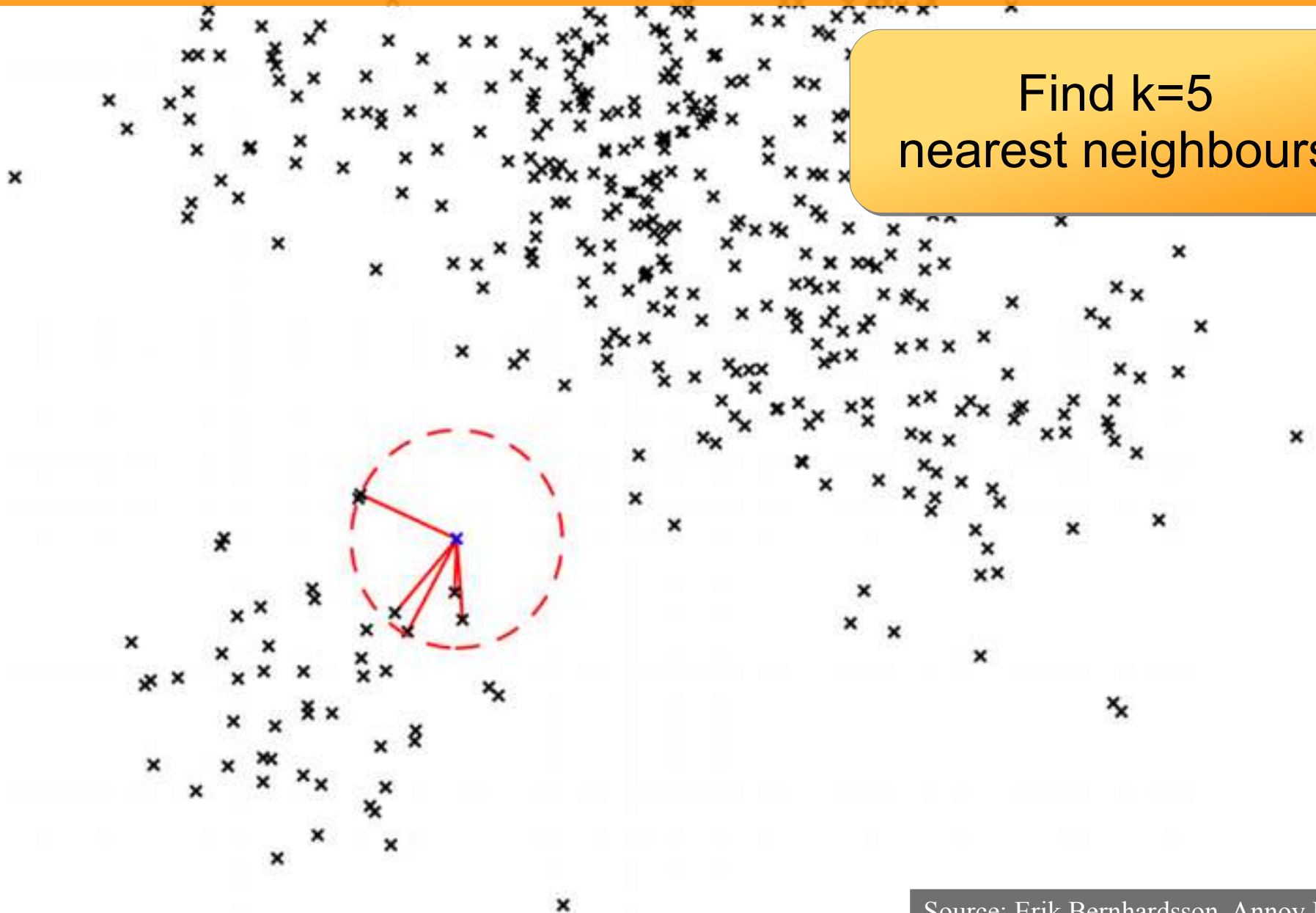
**With good representations,
machine learning becomes
much easier**

2. Find Nearest Neighbours



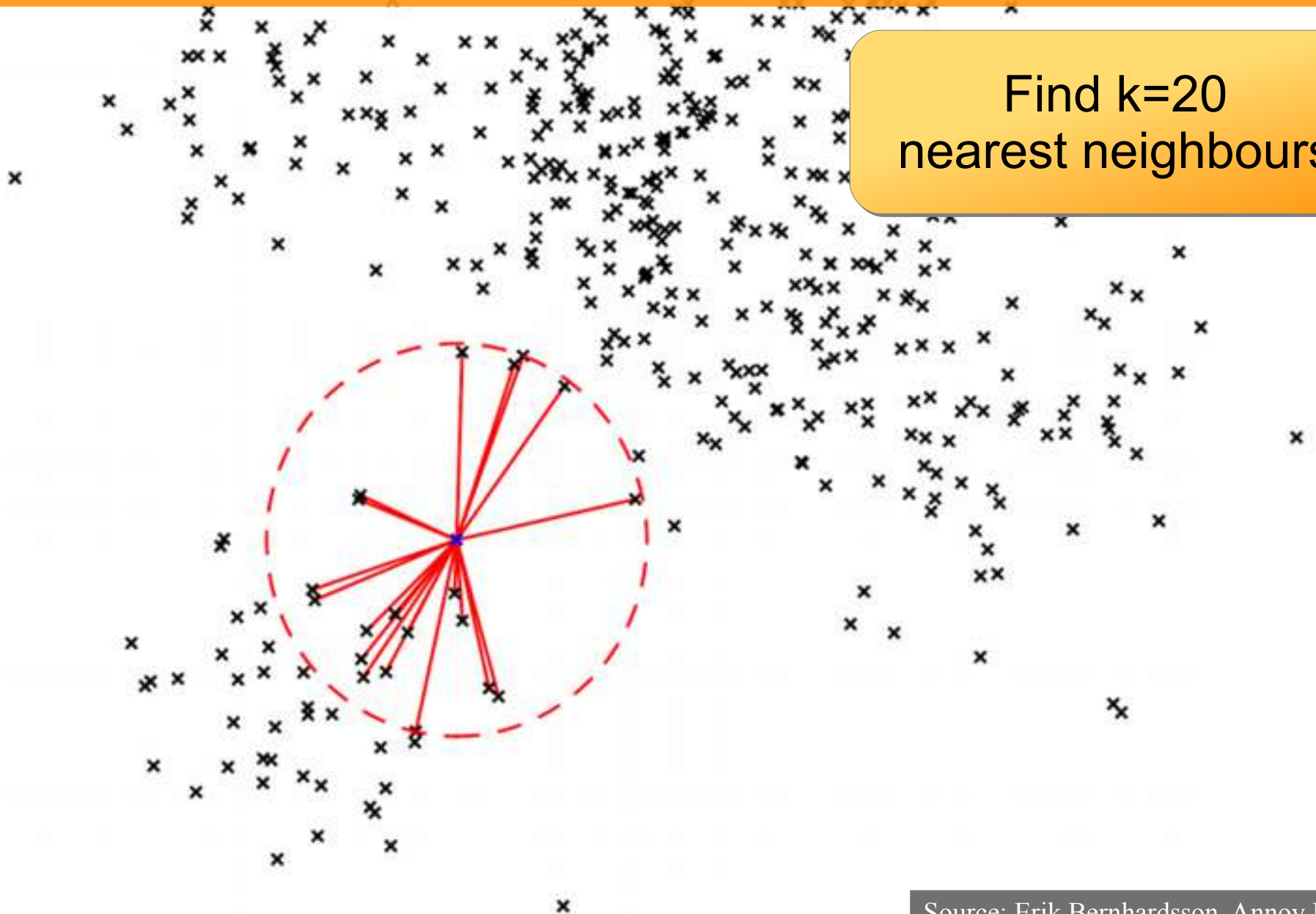
2. Find Nearest Neighbours

Find $k=5$
nearest neighbours



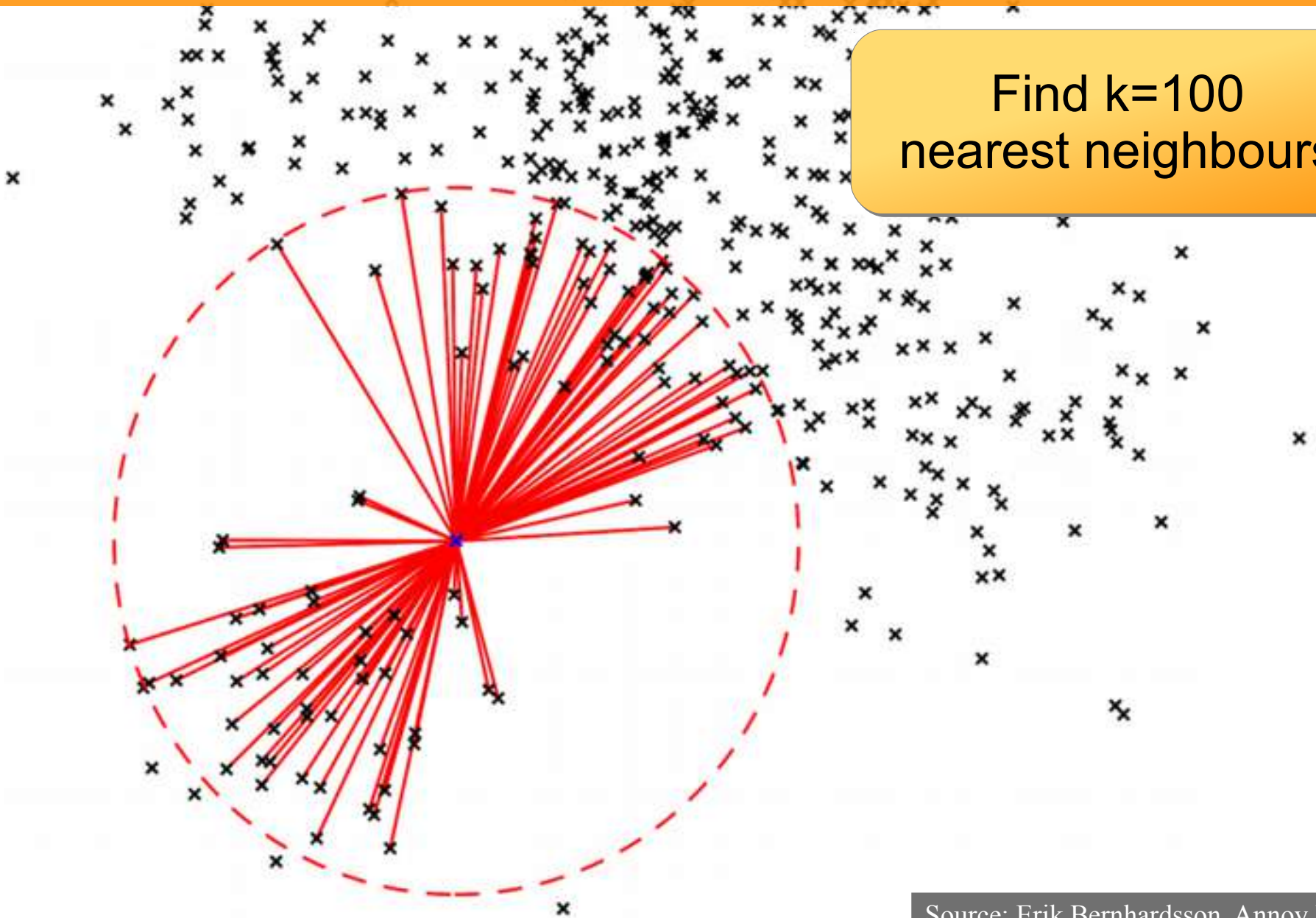
2. Find Nearest Neighbours

Find $k=20$
nearest neighbours



2. Find Nearest Neighbours

Find $k=100$
nearest neighbours



2. Find Nearest Neighbours

Enter word or sentence (EXIT to break): Chinese river

Word	Cosine distance
Yangtze_River	0.667376
Yangtze	0.644091
Qiantang_River	0.632979
Yangtze_tributary	0.623527
Xiangjiang_River	0.615482
Huangpu_River	0.604726
Hanjiang_River	0.598110
Yangtze_river	0.597621
Hongze_Lake	0.594108
Yangtse	0.593442

Outline

- **Word Representations**
- **Phrase Representations**
- **Sentence Representations**
- **Document Representations**
- **Applications and Outlook**

Outline

- **Word Representations**
- **Phrase Representations**
- **Sentence Representations**
- **Document Representations**
- **Applications and Outlook**

Word Representations

“dry” ≠ “arid”

Word Representations

“dry”

0
0
0
1
0
0
0
...
0
0

≠

“arid”

0
0
0
0
1
0
0
...
0
0

Even when you have
Bag-of-Words
representations

Word Representations

“dry”

0.00
0.23
0.03
0.31
0.01
0.03
0.91
...
0.31
0.50

≈

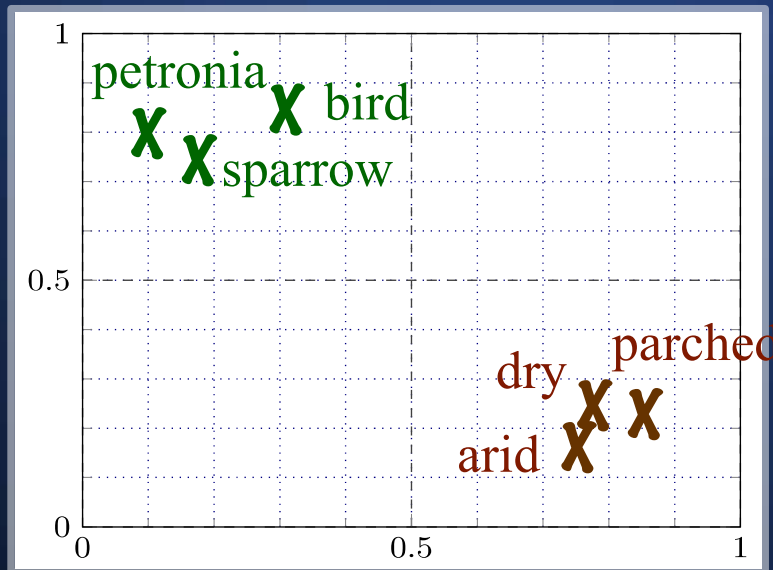
“arid”

0.01
0.23
0.01
0.29
0.00
0.03
0.92
...
0.30
0.51

**Distributed
Vector
Representation**

(e.g. with d=300
dimensions)

Word Representations



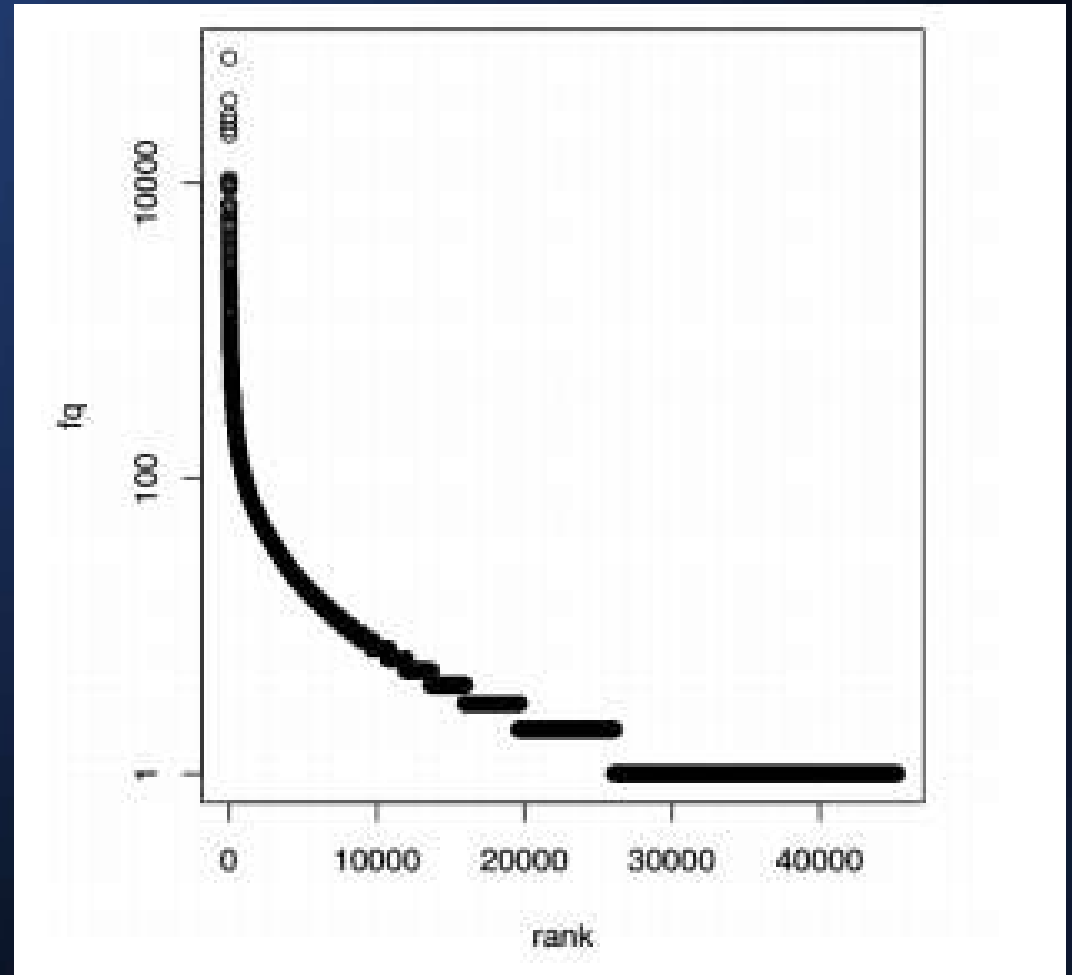
Why Generalization is Important

Data Sparsity:

Features may be rare.

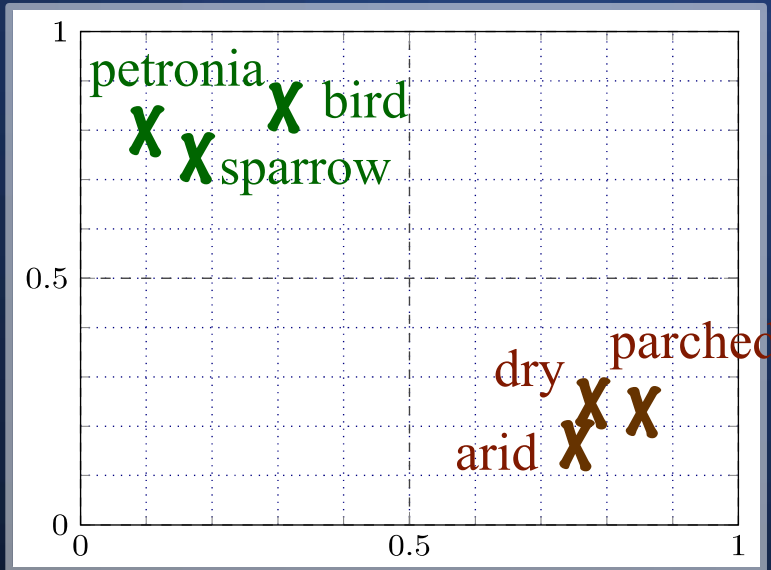
→ Missing in training data

E.g. most words have a low frequency (in the “long tail”)



Frequencies of Words in Brown Corpus

Word Representations



How can we obtain such vector representations?

Distributional Semantics

We found a little, hairy
wampimuk
sleeping behind the tree

cf. McDonald & Ramscar (2001)



Image: https://commons.wikimedia.org/wiki/File:Mahogany_glider.jpg

Distributional Semantics

He filled the **wampimuk**,
passed it around
and we all drunk some.

cf. McDonald & Ramscar (2001)

Term-Term Cooccurrence Matrix

dogs are animals
cats are animals
orchids are plants
roses are plants



	Animals	Are	Cats	Dogs	Orchids	Plants	Roses
Animals		X	X	X			
Are	X		X	X	X	X	X
Cats	X	X					
Dogs	X	X					
Orchids		X				X	
Plants		X			X		X
Roses		X				X	

Term-Term Cooccurrence Matrix

dogs are animals
cats are animals
orchids are plants
roses are plants



	Animals	Are	Cats	Dogs	Orchids	Plants	Roses
Animals		X	X	X			
Are	X		X	X	X	X	X
Cats	X	X					
Dogs	X	X					
Orchids		X				X	
Plants		X			X		X
Roses		X				X	

Term-Document Matrix

D1: dogs are animals
D2: cats are animals
D3: orchids are plants
D4: roses are plants



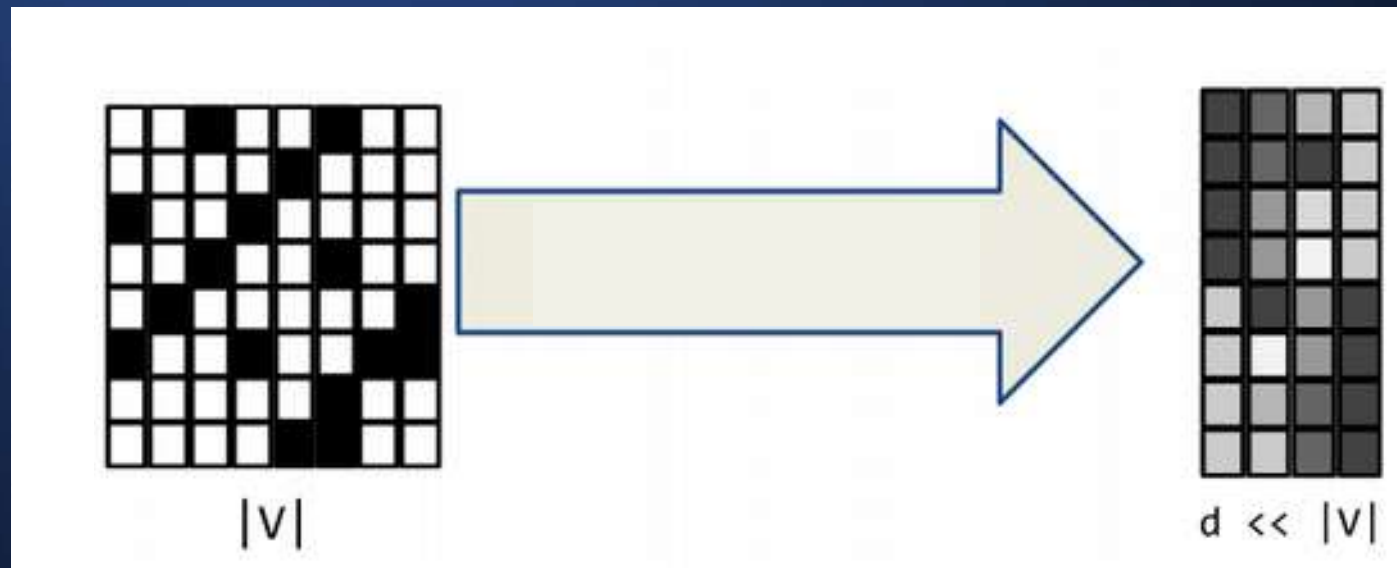
	D1	D2	D3	D4
Animals	X	X		
Are	X	X	X	X
Cats		X		
Dogs	X			
Orchids			X	
Plants			X	X
Roses				X

Large context scope.
Typically, more topic-oriented similarity.

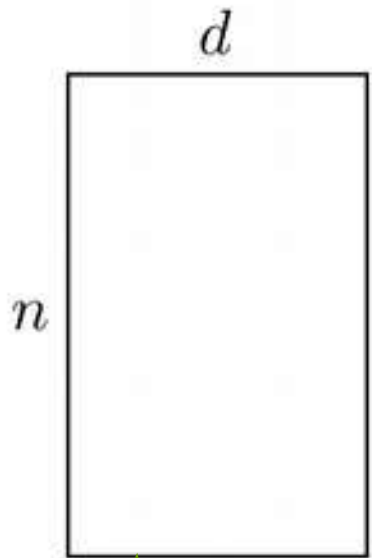
Image: Rafael Banchs

**Classic
methods:
Dimensionality
Reduction on
matrix (SVD)**

Learning Representations



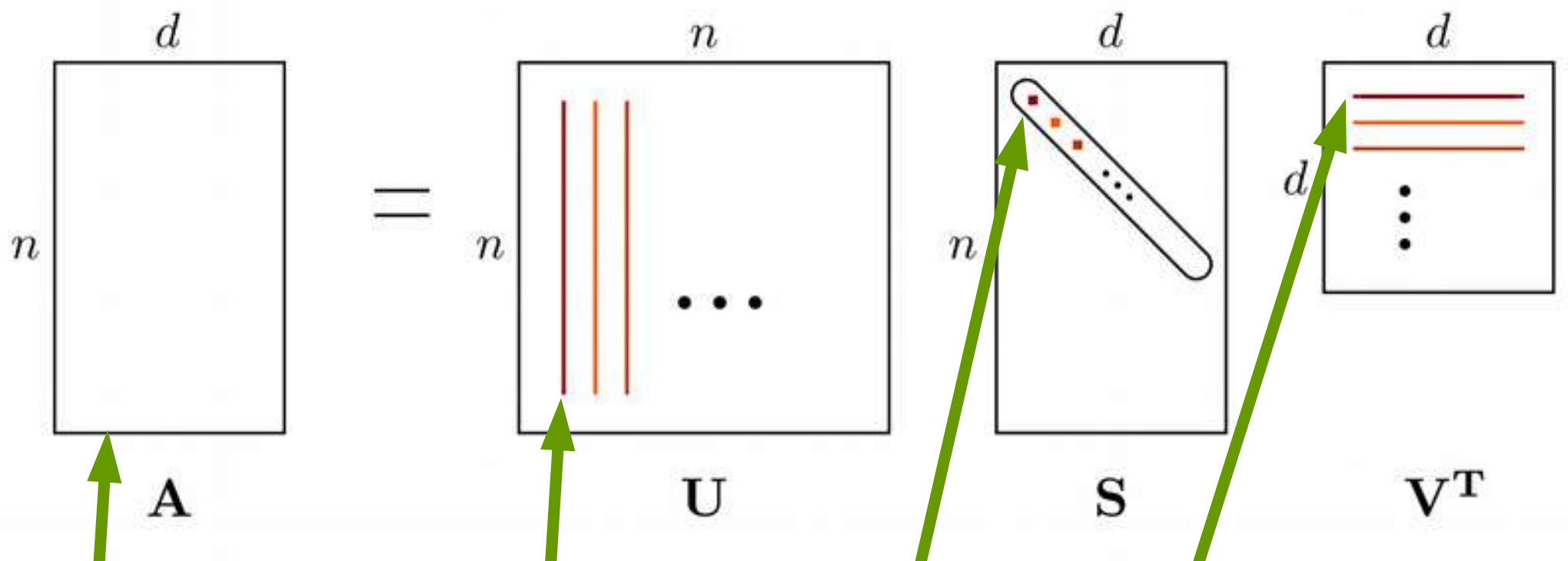
SVD



A

E.g.
Cooccurrence Matrix

SVD



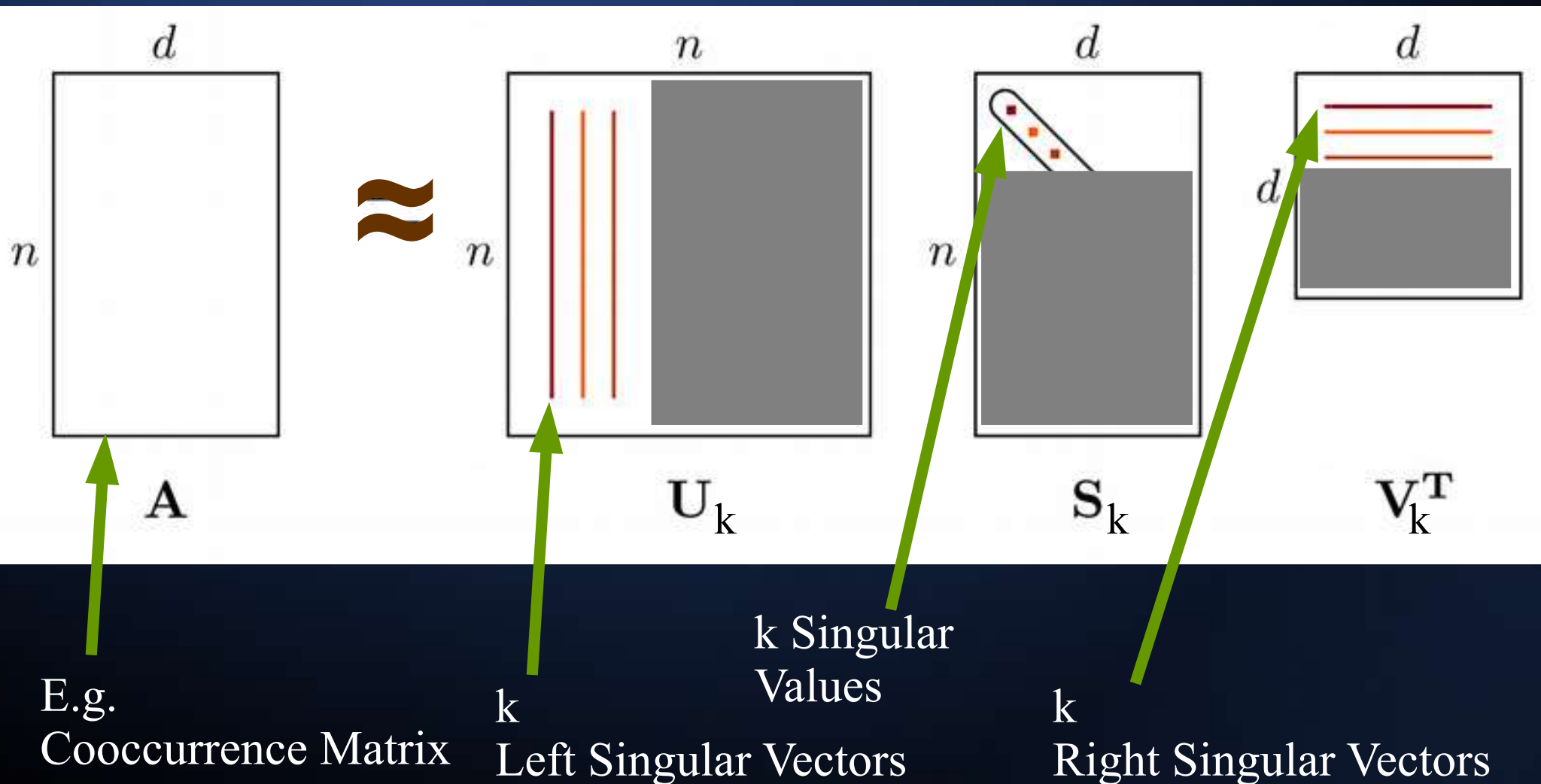
E.g.
Cooccurrence Matrix

Left Singular Vectors

Singular Values

Right Singular Vectors

Low-Rank Approximation via SVD



Language Models

bing

no one uses



no one uses **bing**

no one uses **the white crayon**

no one uses **myspace** anymore

no one uses **aol**

no one uses **twitter**

no one uses **icq**

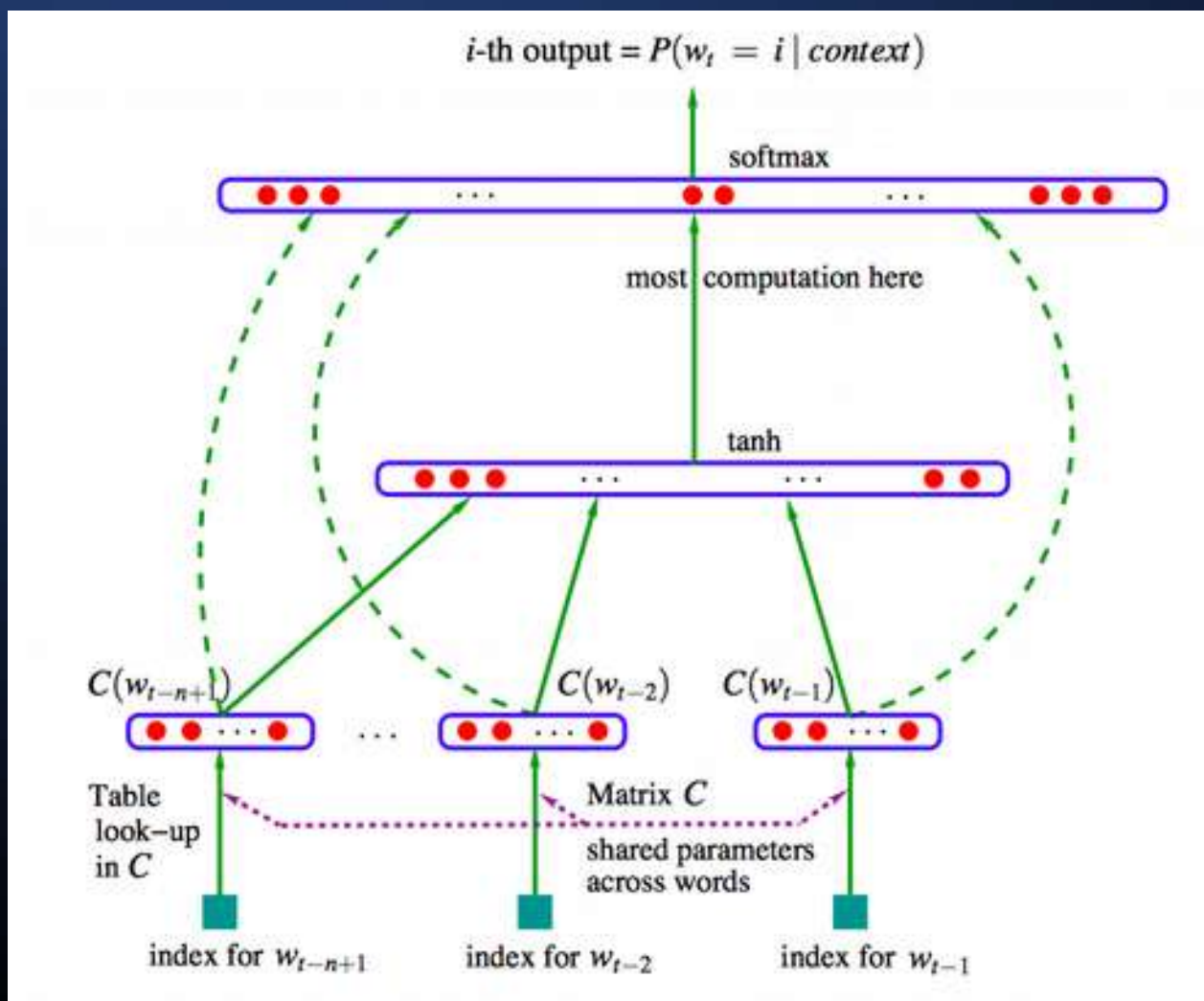
no one uses **mics** in battlefield 3

no one uses **highway** in simcity 4

[Manage search history](#)

Language Models for Neural Word Vectors

Bengio et al. (2003). A Neural Probabilistic Language Model

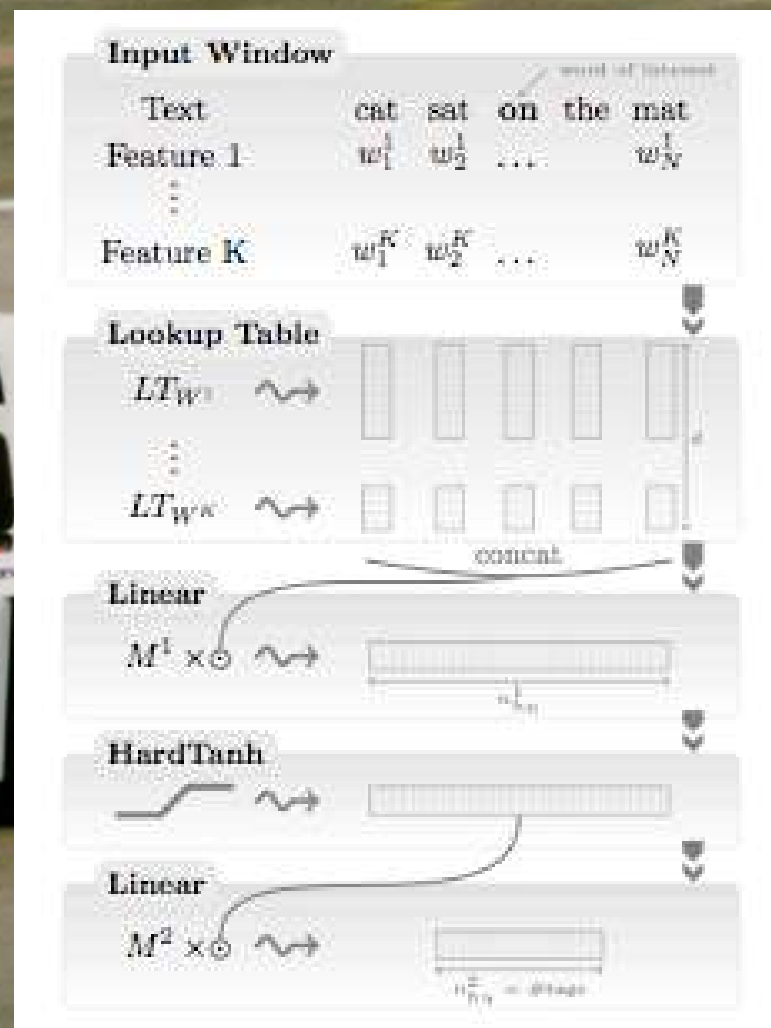


SENNA Embeddings (2008, 2011)

**Collobert et al.
(2008, 2011):**

**Natural Language Processing
(Almost) From Scratch**

**One architecture for
part-of-speech tagging,
chunking, and
named entity recognition**



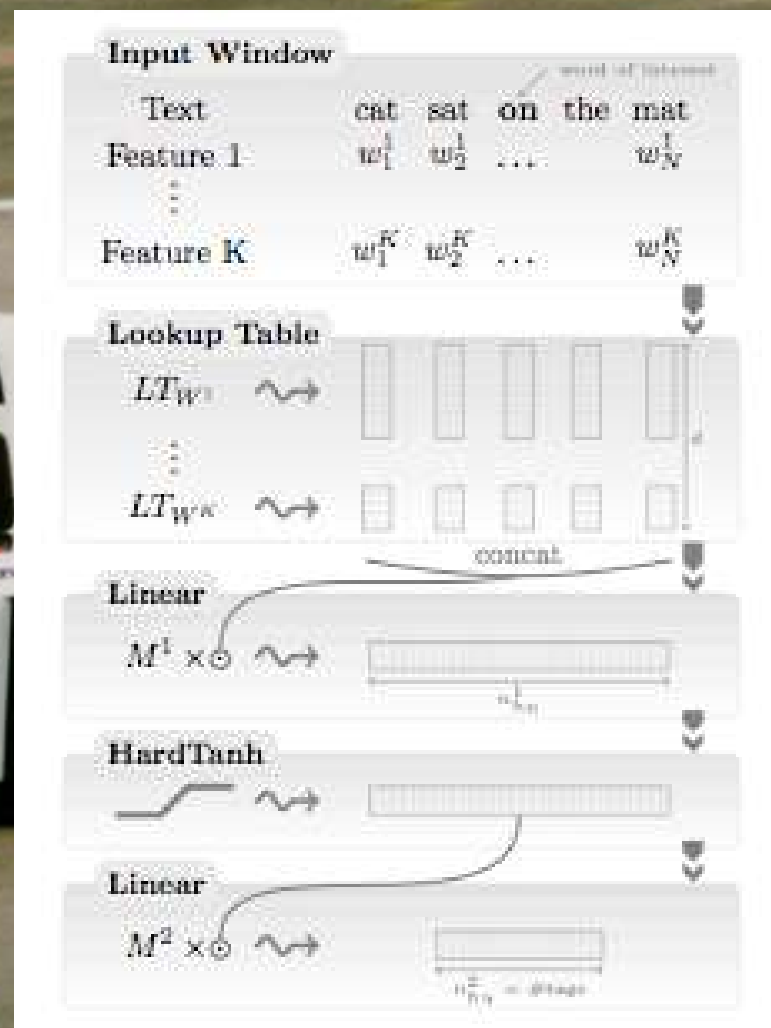
SENNA Embeddings (2008, 2011)

**Collobert et al.
(2008, 2011):**

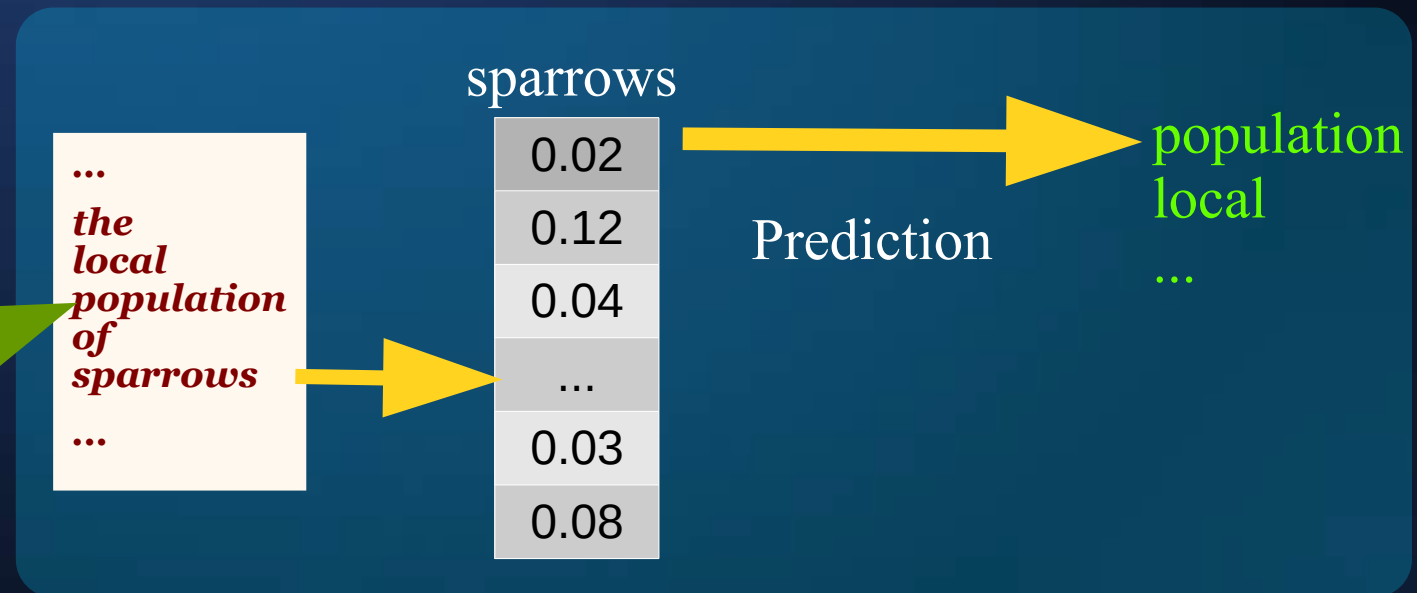
**Natural Language Processing
(Almost) From Scratch**

**Looked at both sides
of context to train
embeddings.**

**Training Time (2011):
4 weeks + 3 weeks**



Word Vector Representations: word2vec



word2vec Skip-Gram Model

Word Vector Representations: word2vec



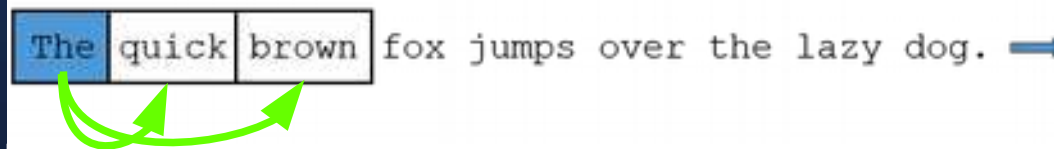
Word Vector Representations: word2vec



word2vec

Tool for computing continuous distributed representations of words.

Source Text



Training
Samples

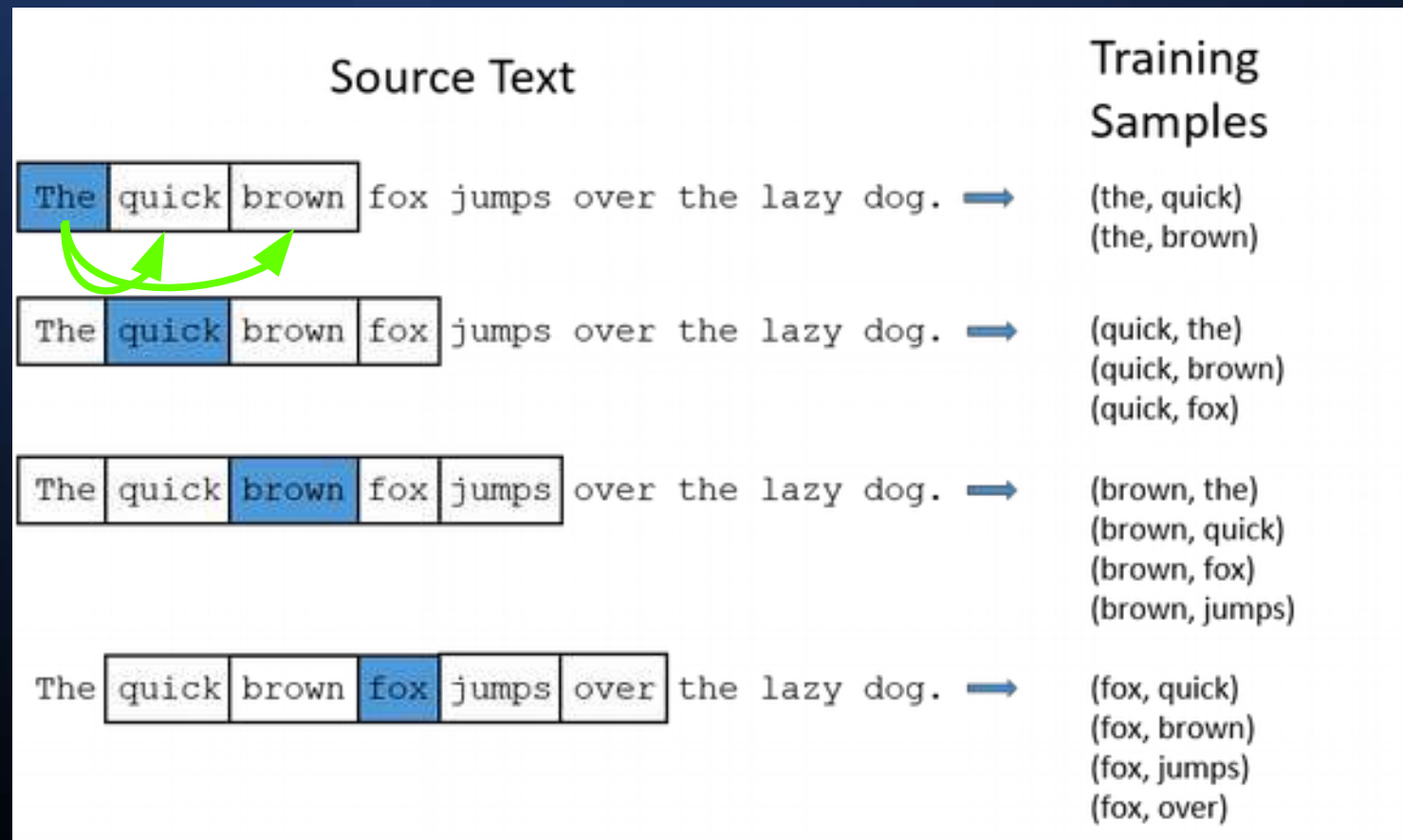
(the, quick)
(the, brown)

Word Vector Representations: word2vec

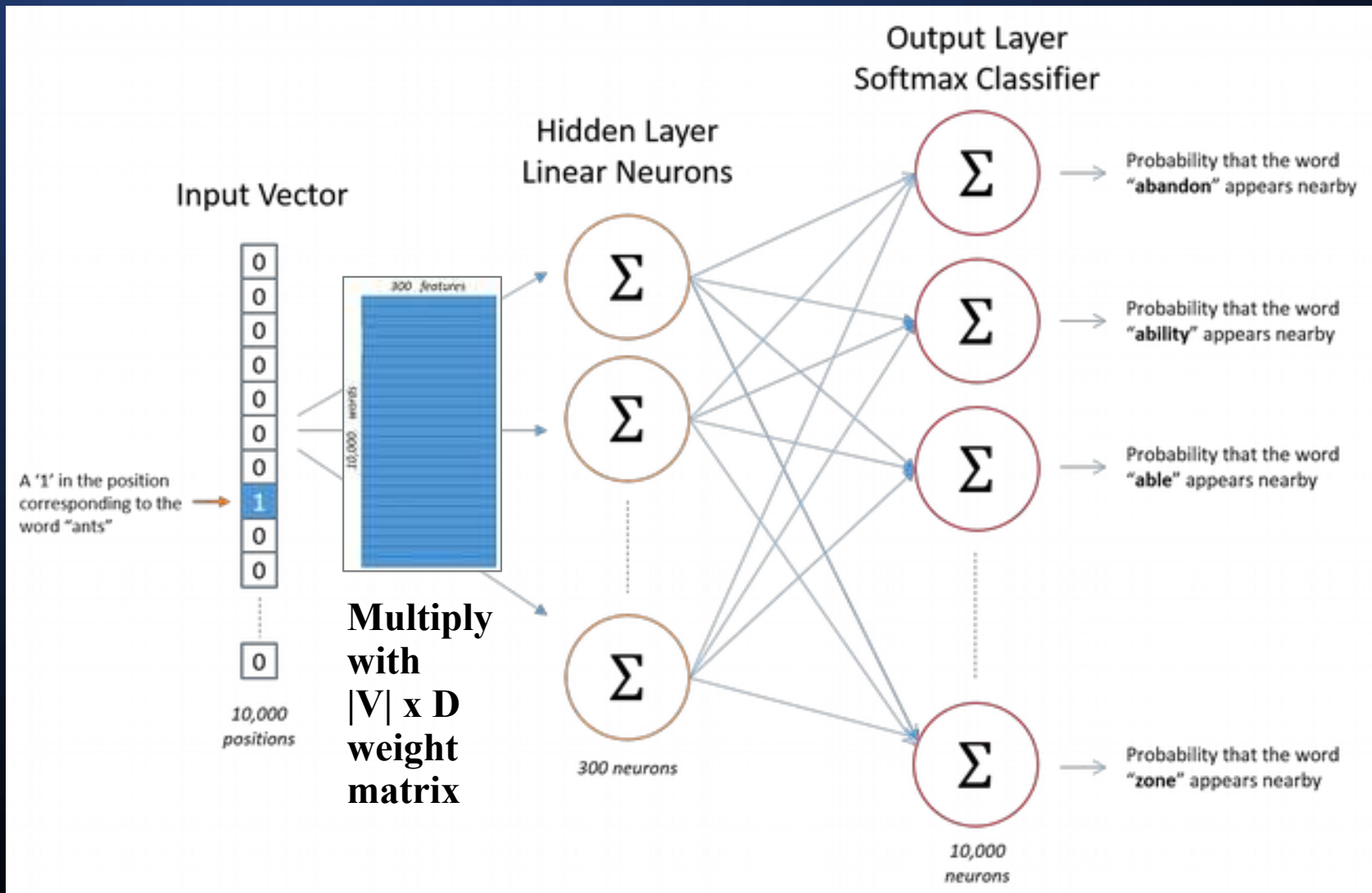


word2vec

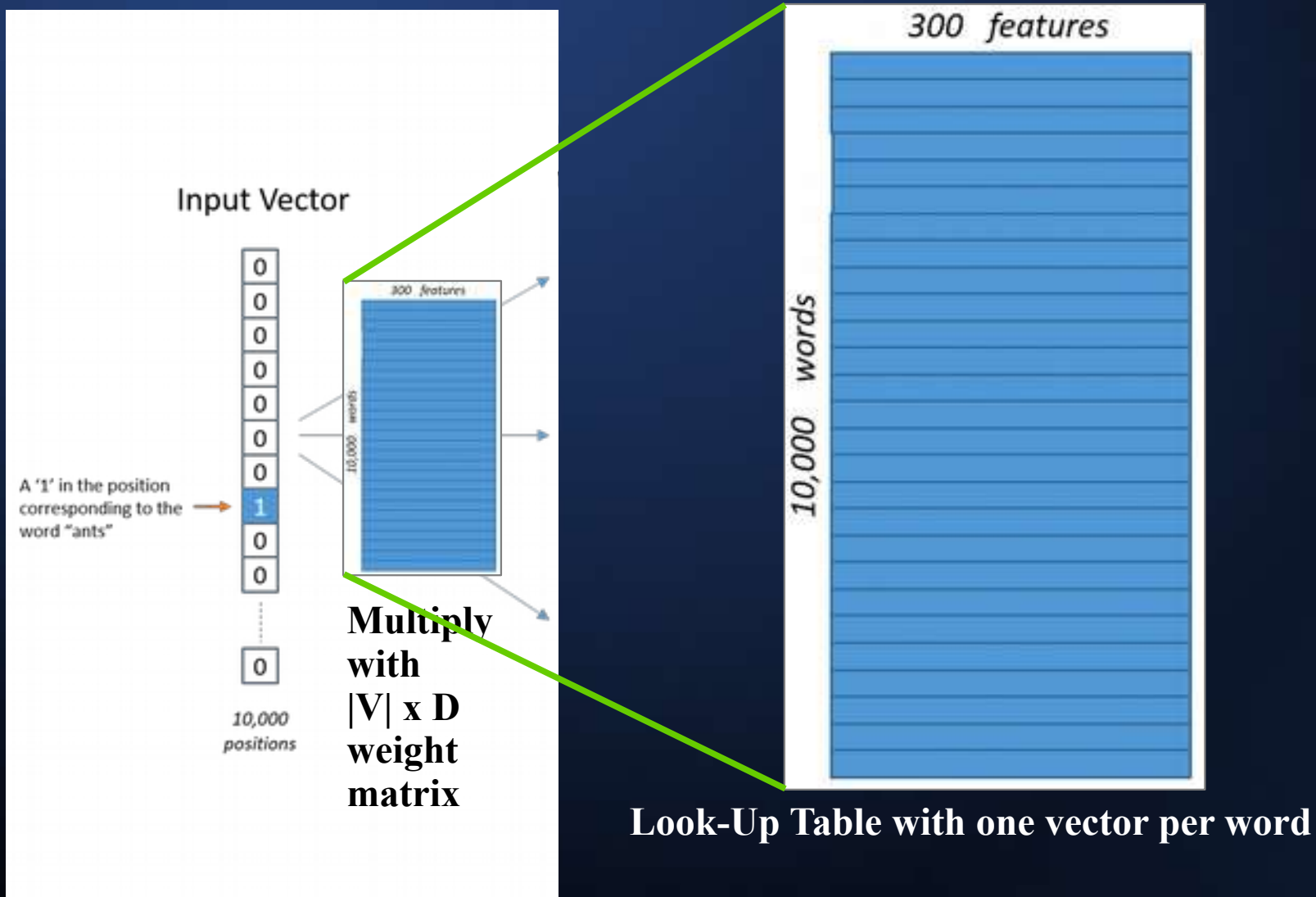
Tool for computing continuous distributed representations of words.



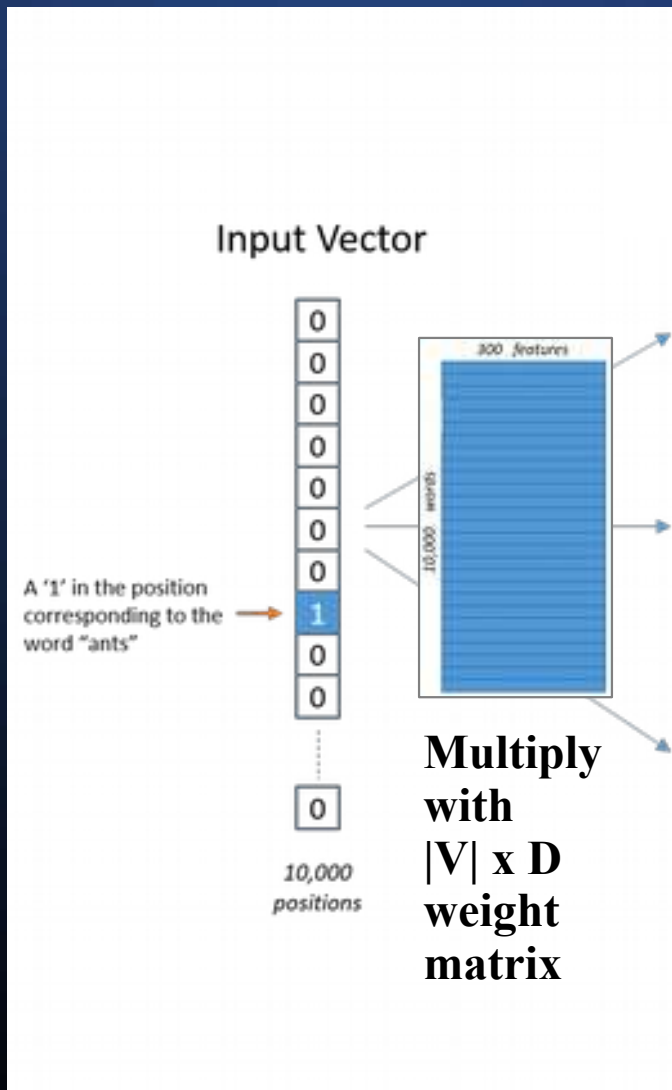
Word Vector Representations: word2vec



Word Vector Representations: word2vec



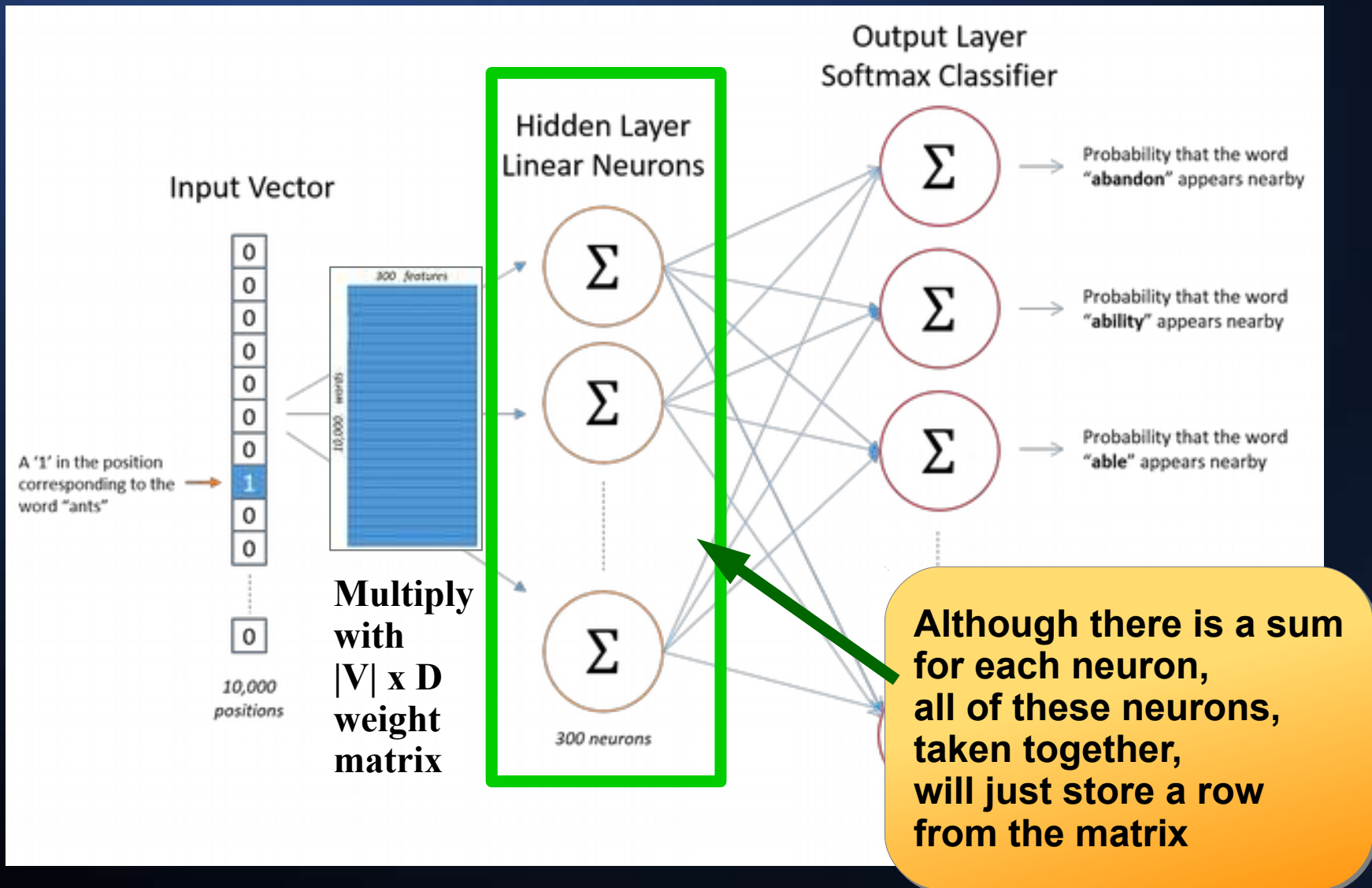
Word Vector Representations: word2vec



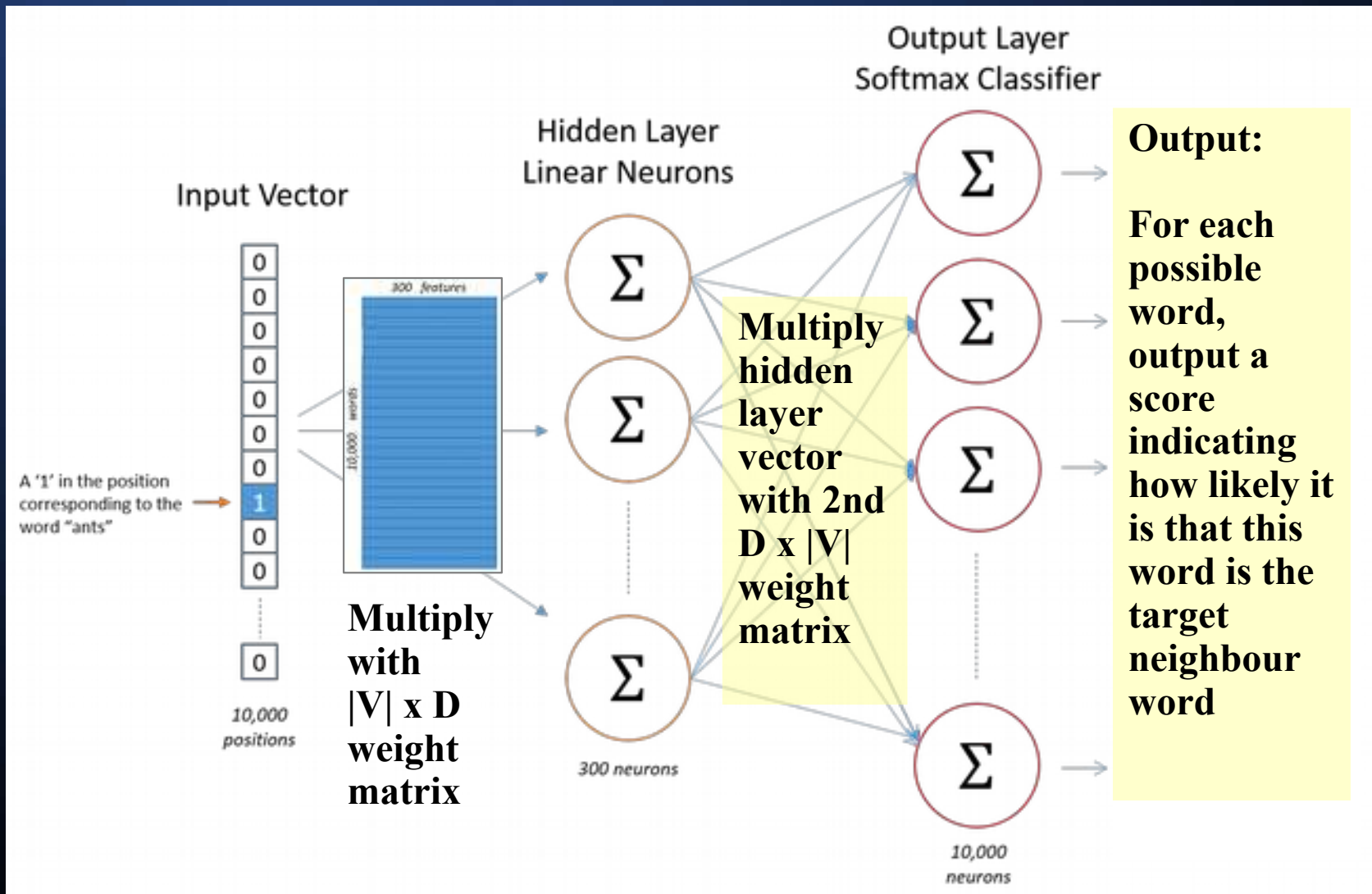
$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$

Multiplying a one-hot vector with this matrix will simply look-up a row in the matrix!

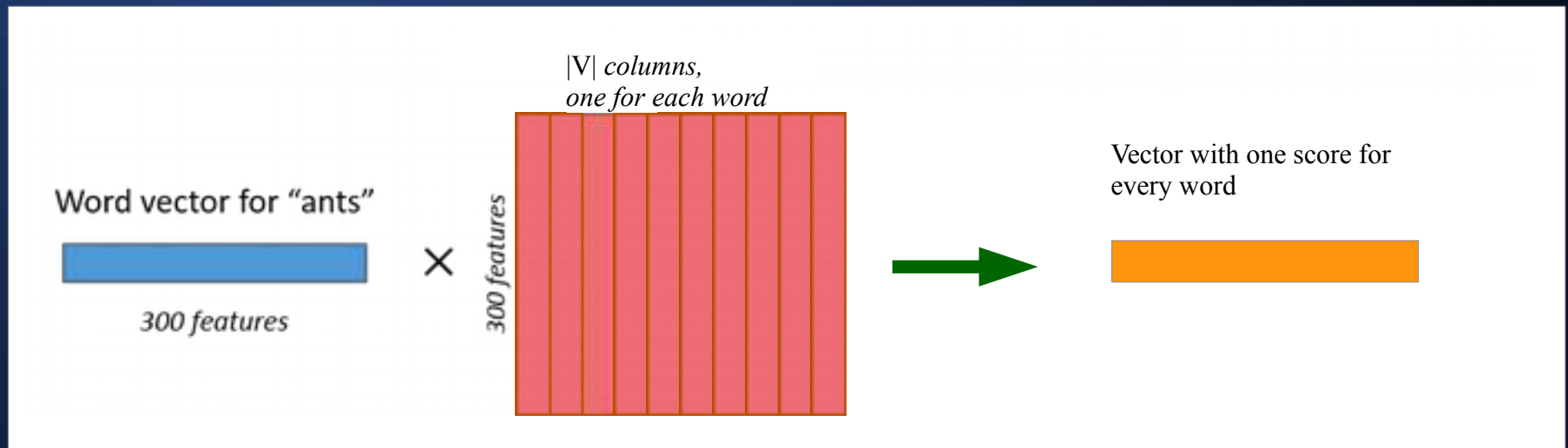
Word Vector Representations: word2vec



Word Vector Representations: word2vec



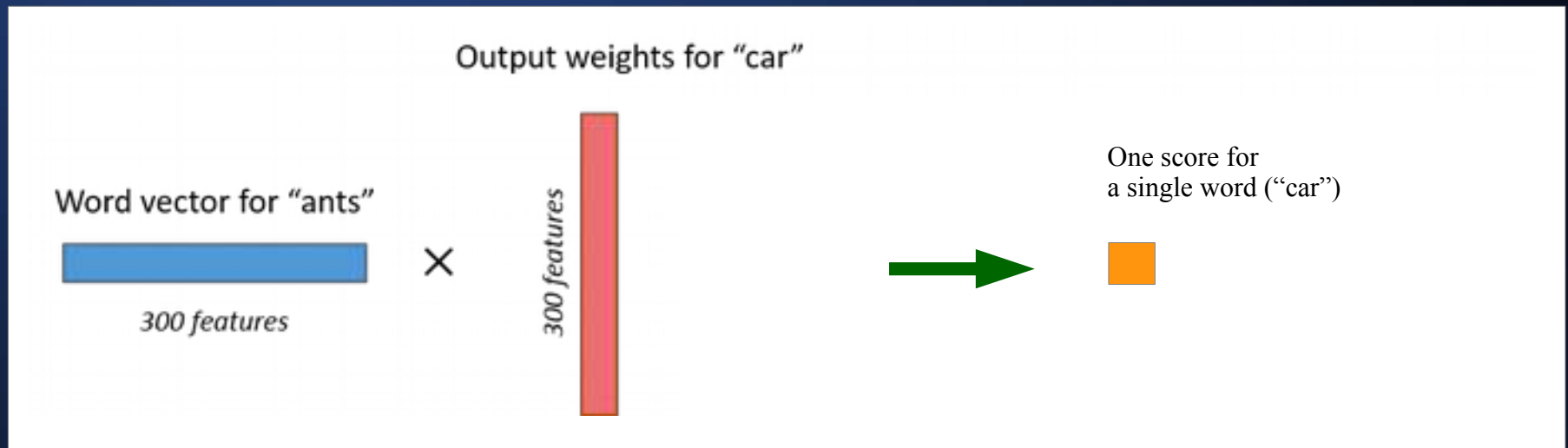
Word Vector Representations: word2vec



“The ant trail at the car ...”



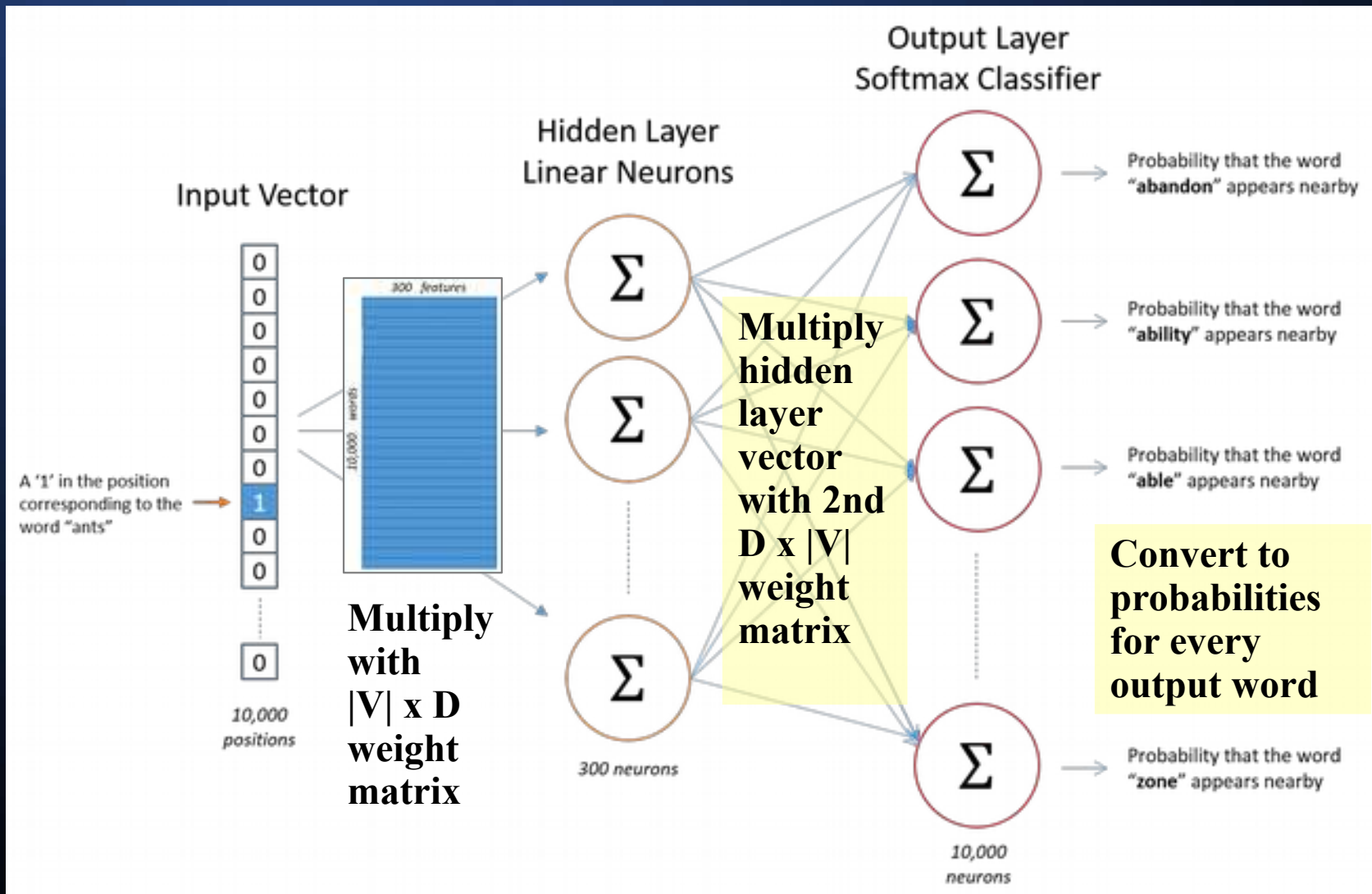
Word Vector Representations: word2vec



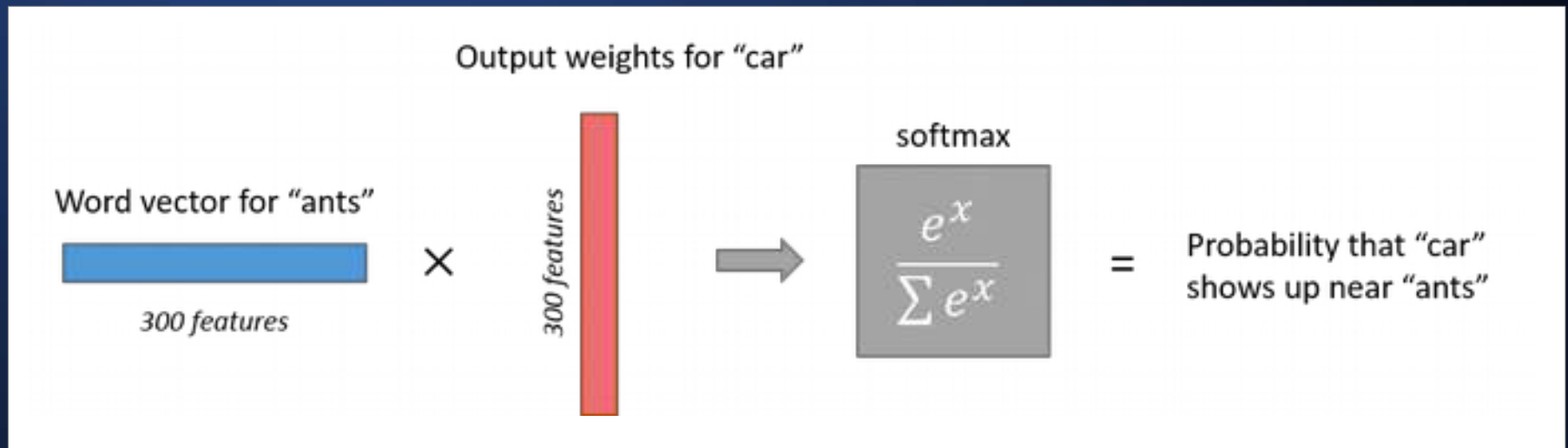
“The ant trail at the car ...”



Word Vector Representations: word2vec



Word Vector Representations: word2vec

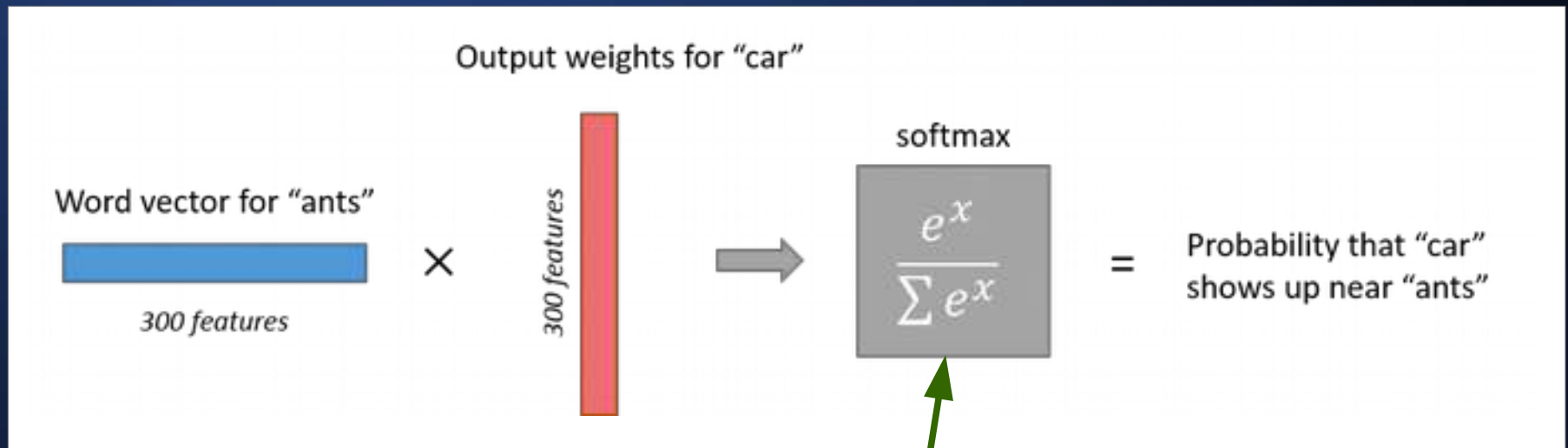


“The ant trail at the car ...”



Probability should be highest for the current target word we are considering.

Word Vector Representations: word2vec

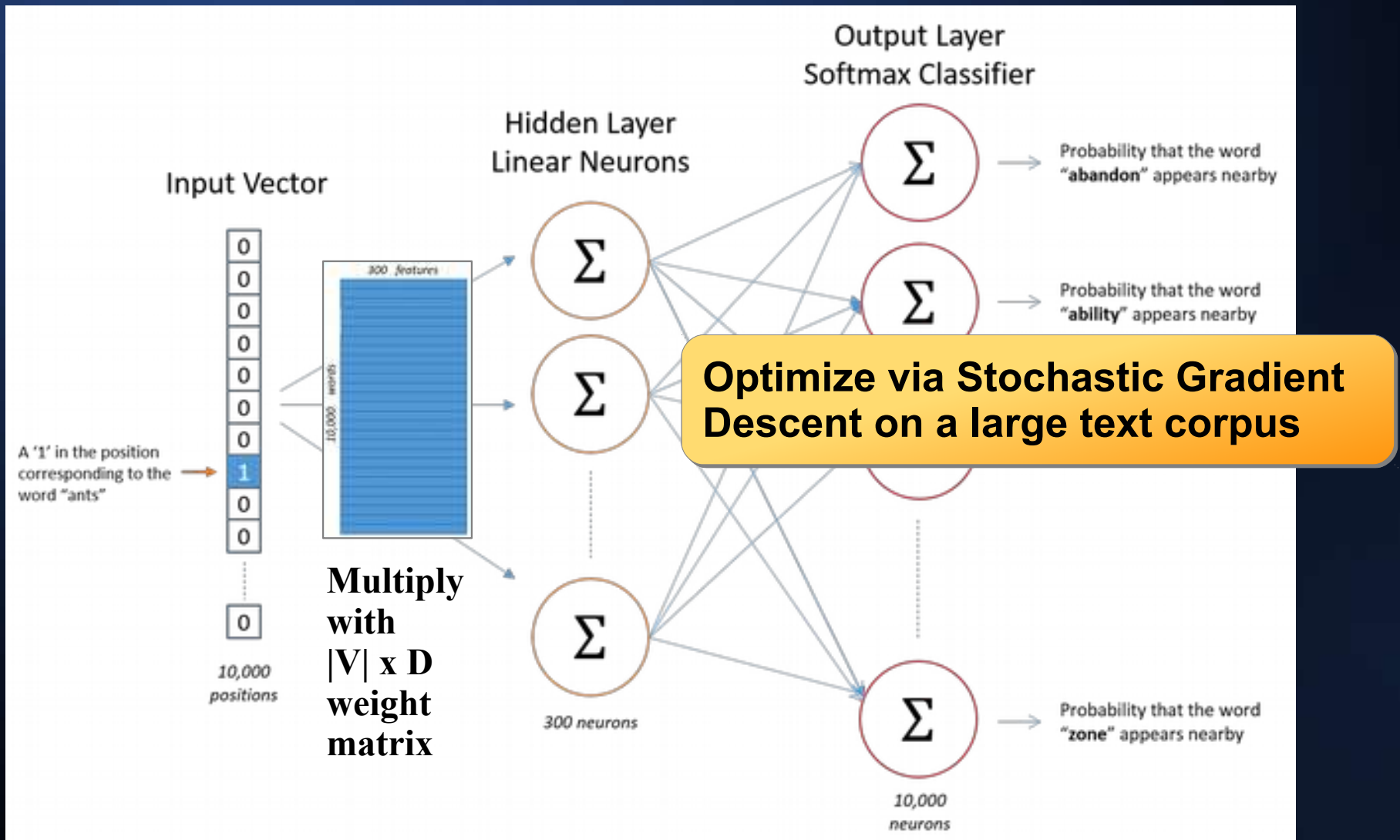


“The ant trail at the car ...”

Optimization:
Instead of normalizing over all alternative words, just consider k negative example words in addition to current positive example target word

Probability should be highest for the current target word we are considering.

Word Vector Representations: word2vec

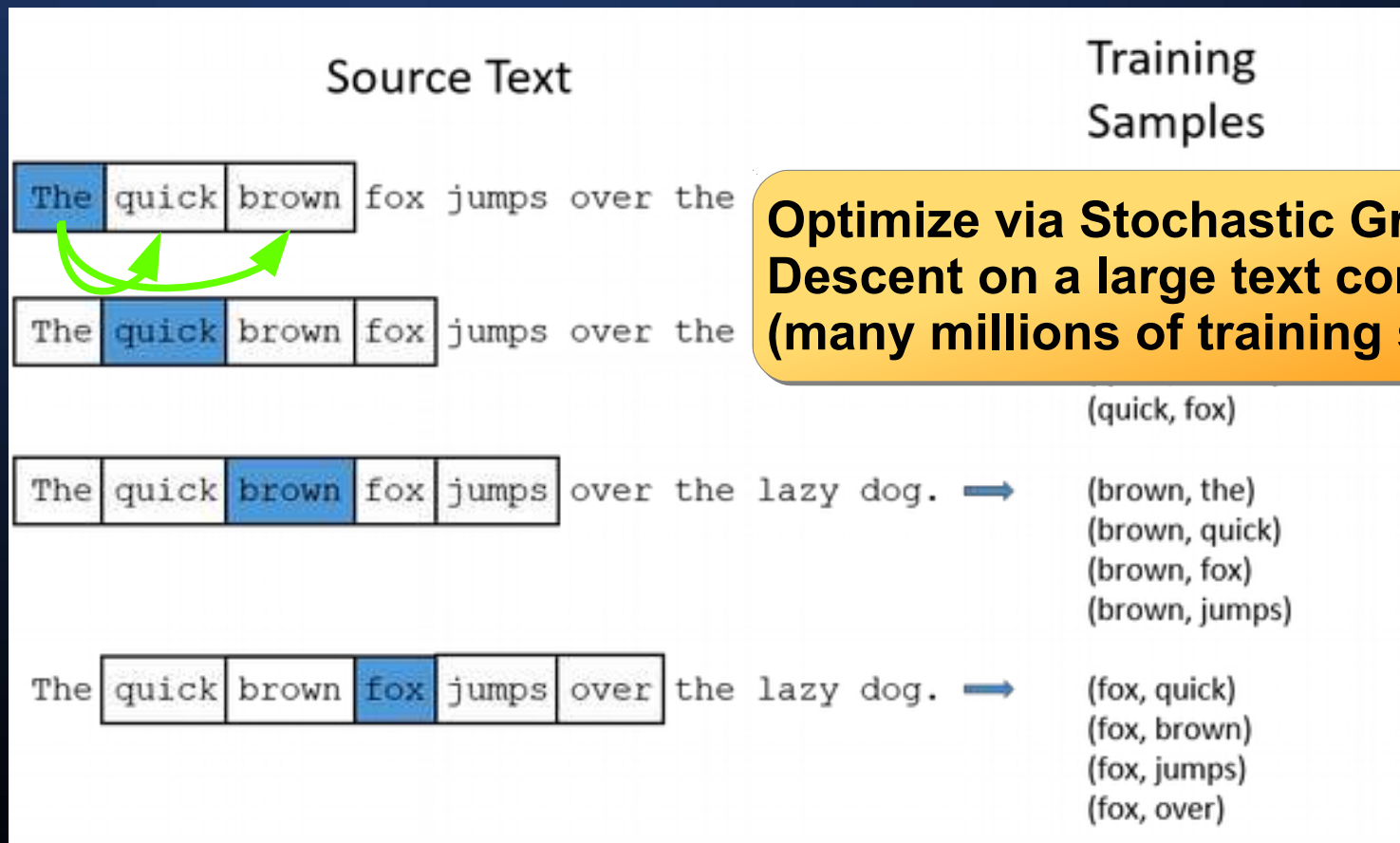


Word Vector Representations: word2vec



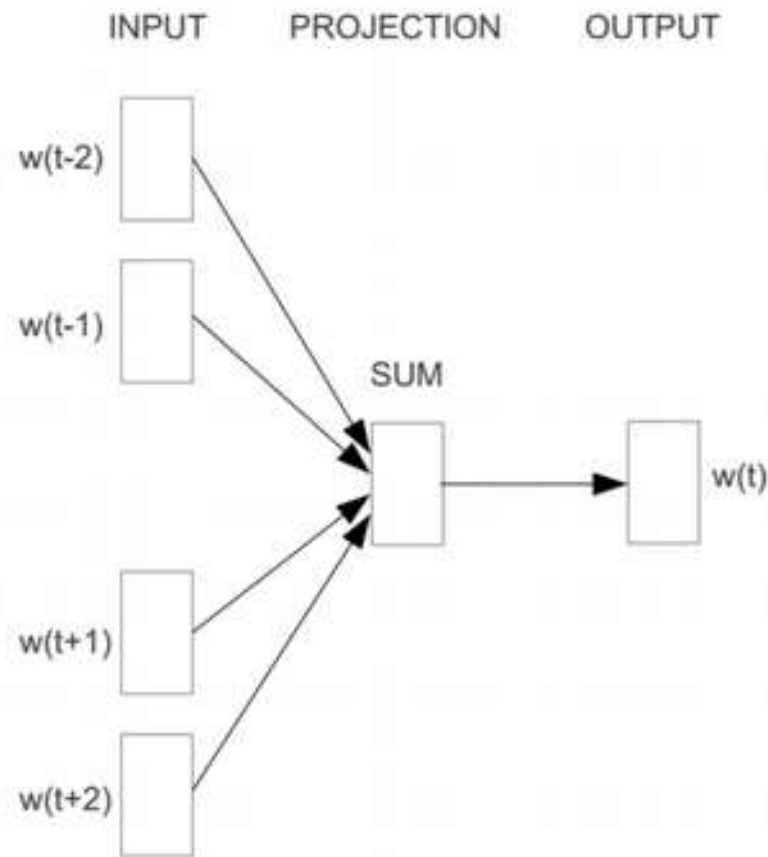
word2vec

Tool for computing continuous distributed representations of words.

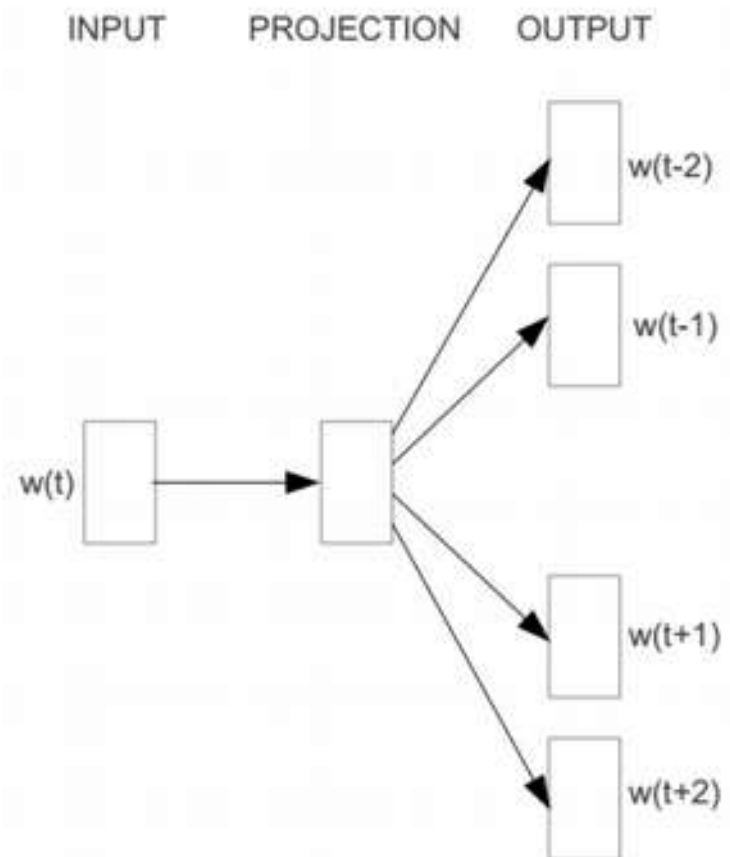


Optimize via Stochastic Gradient Descent on a large text corpus (many millions of training samples)

Word Vector Representations: word2vec



CBOW



Skip-gram

Mikolov et al. Distributed Representations of Words and Phrases and their Compositionality. Proc. NIPS 2013

Questions?

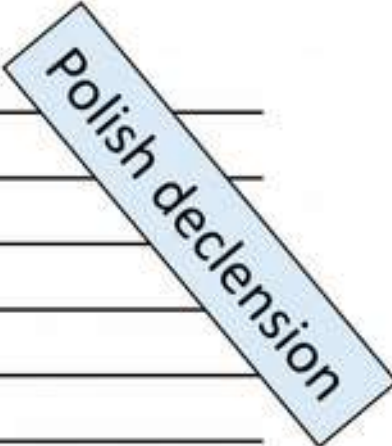


Rare Words

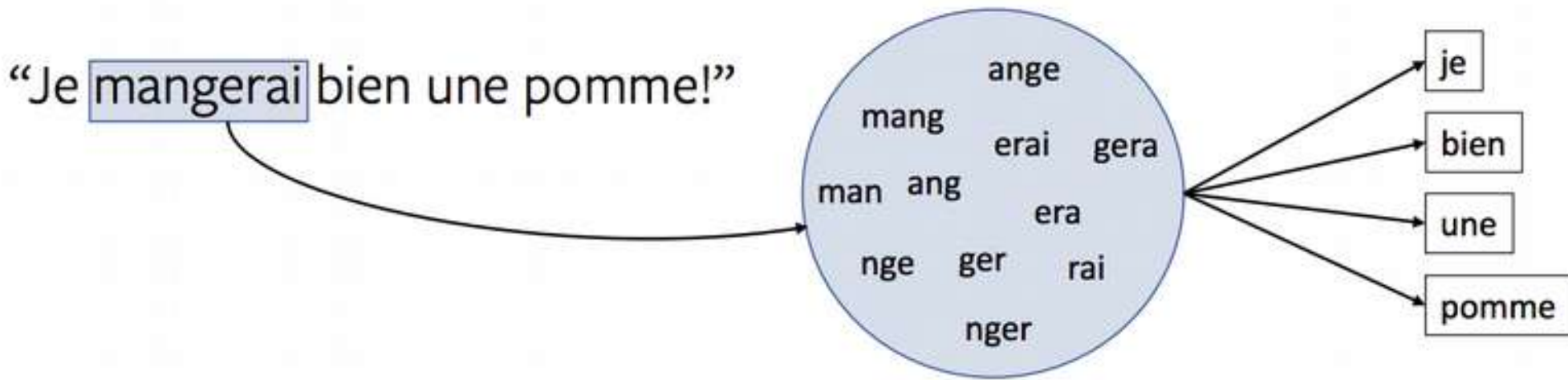


Case 1: Rare Forms

	Singular	Plural
Nominative	uniwersytet	uniwersytety
Genetive	uniwersytetu	uniwersytetów
Dative	uniwersytetowi	uniwersytetom
Accusative	uniwersytet	uniwersytety
Instrumental	uniwersytetem	uniwersytetami
Locative	uniwersytecie	uniwersytetach
Vocative	uniwersytecie	uniwersytety



Character N-Gram Approach as in fastText



$$h_w = \sum_{g \in w} x_g$$

mang erai ange
man ang gera
nge ger rai nger

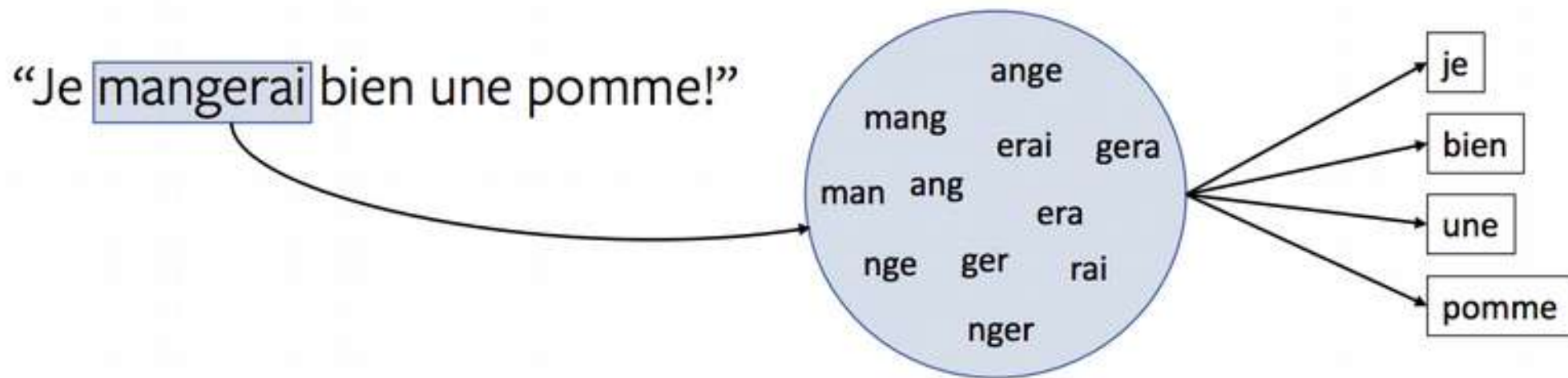
Character n-grams

+

mangerai

Word itself

Character N-Gram Approach as in fastText



$$h_w = \sum_{g \in w} x_g$$

The diagram shows the formula $h_w = \sum_{g \in w} x_g$ on the left. On the right, it shows the components of the sum: "Character n-grams" (listing the same n-grams as the circle) and "Word itself" (listing "mangerai"). A large orange X is drawn over the "Word itself" part, indicating it can be omitted.

**Can also omit word itself and
support out-of-vocabulary forms**

Character N-Gram Approach as in fastText

The logo for fastText, with 'fast' in red italicized font and 'Text' in blue bold font.

***fast*Text**

a library for efficient text classification
and word representation

**fastText is not the first work to
consider morphology, but one
of the easiest tools to use**

Case 2: Rare Words/Names

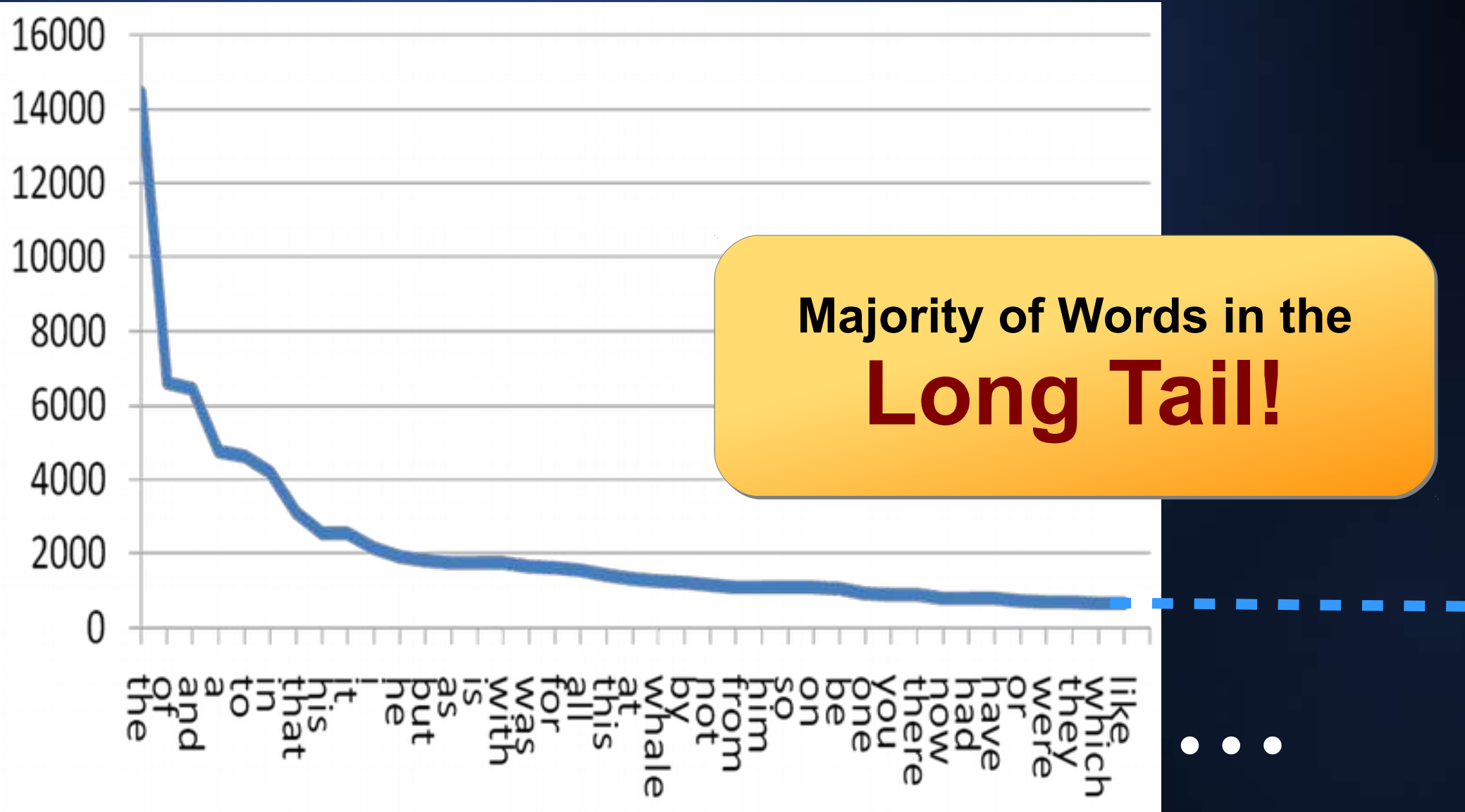
**Does anyone know what
“Mangalia”
is?**

Case 2: Rare Words/Names



**Mangalia:
Romanian port city**

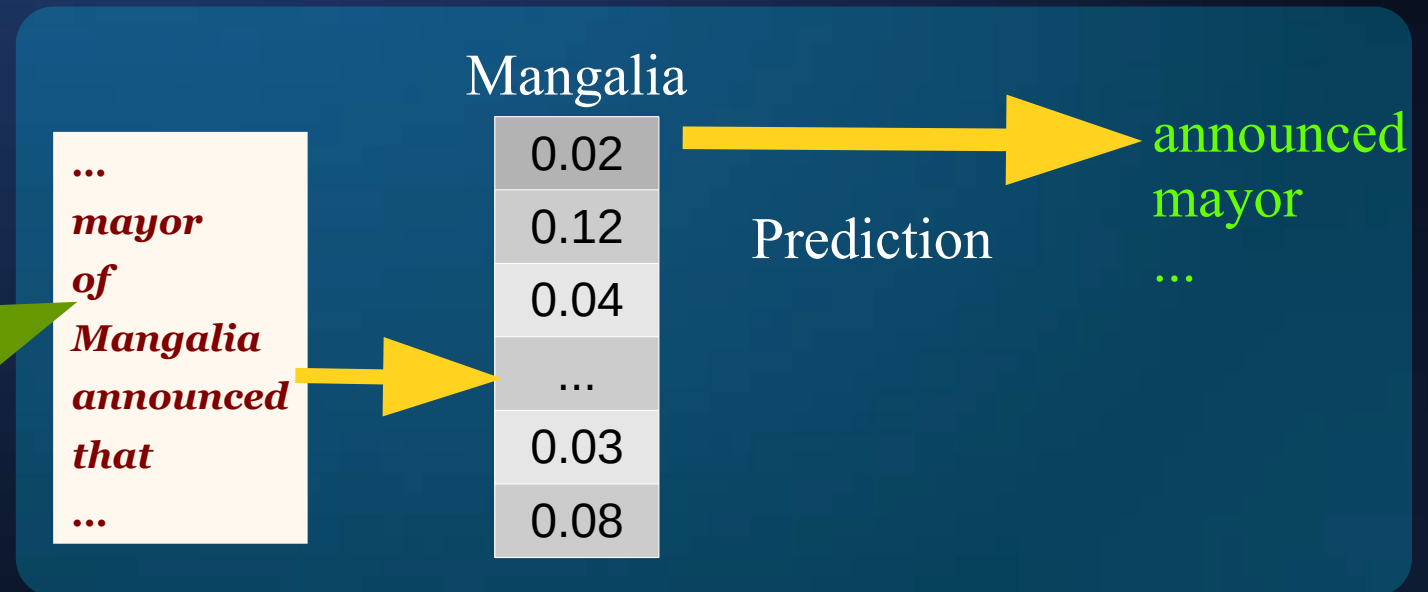
Zipf's Law



Out-of-Vocabulary Problem

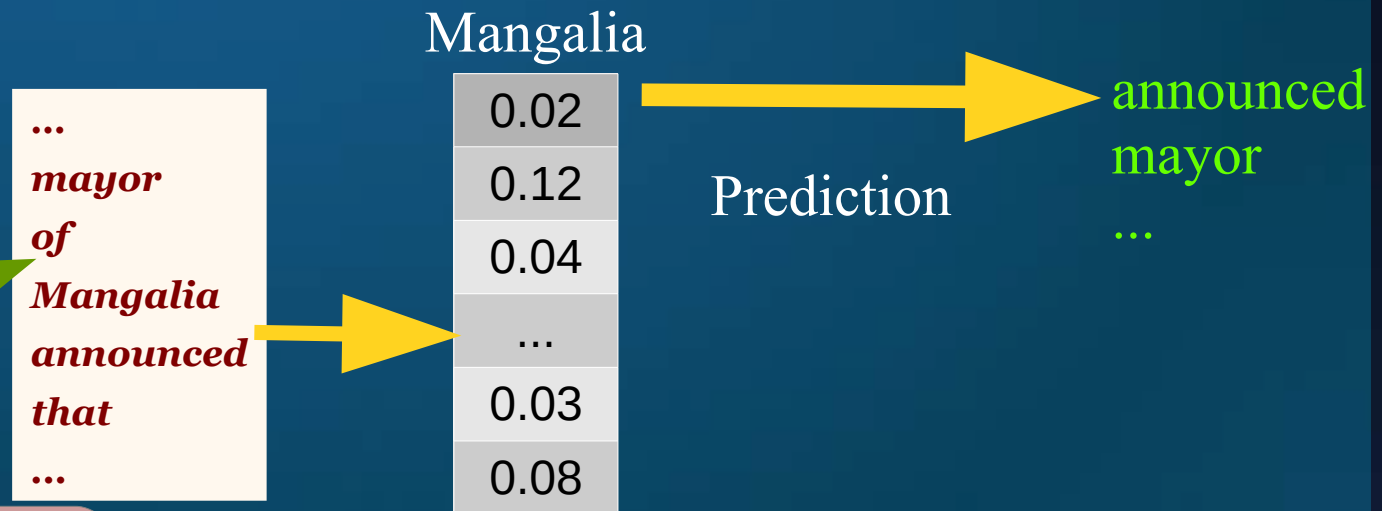
The mayor of *UNK* traveled to *UNK*.

Word Vector Representations: word2vec



word2vec Skip-Gram Model

Word Vector Representations: word2vec



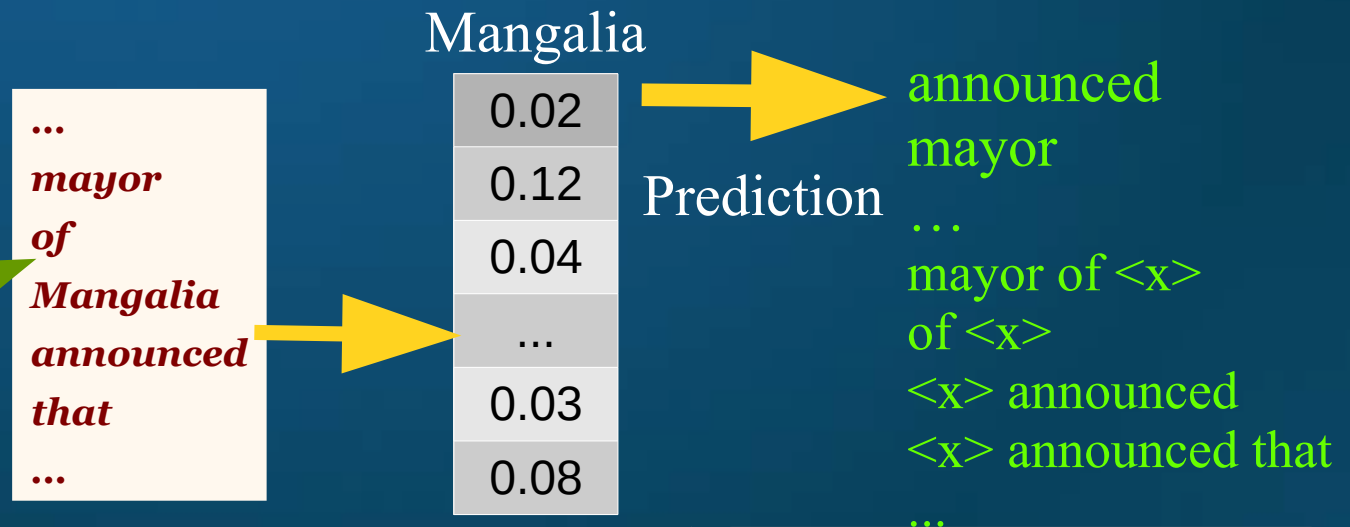
Problem:
Not very discriminative.
E.g. “*of mayor Laszkovic
who announced that*”

word2vec Skip-Gram Model

Context-Based Few-Shot Word Representation Learning

Idea 1

1. Predict n-grams (e.g. “New York”, “how much”)
2. Predict position-specific patterns

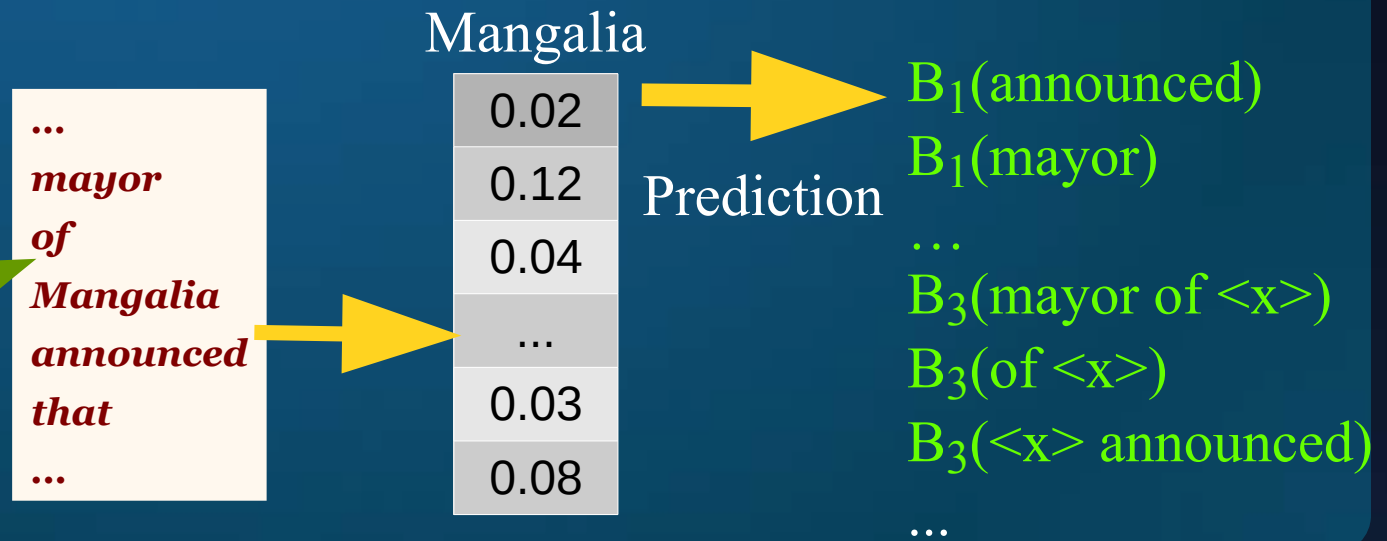


Problem:
Combinatorial
explosion in no. of
prediction targets!

Context-Based Few-Shot Word Representation Learning

Idea 2

Constrain # prediction targets by binning into feature buckets

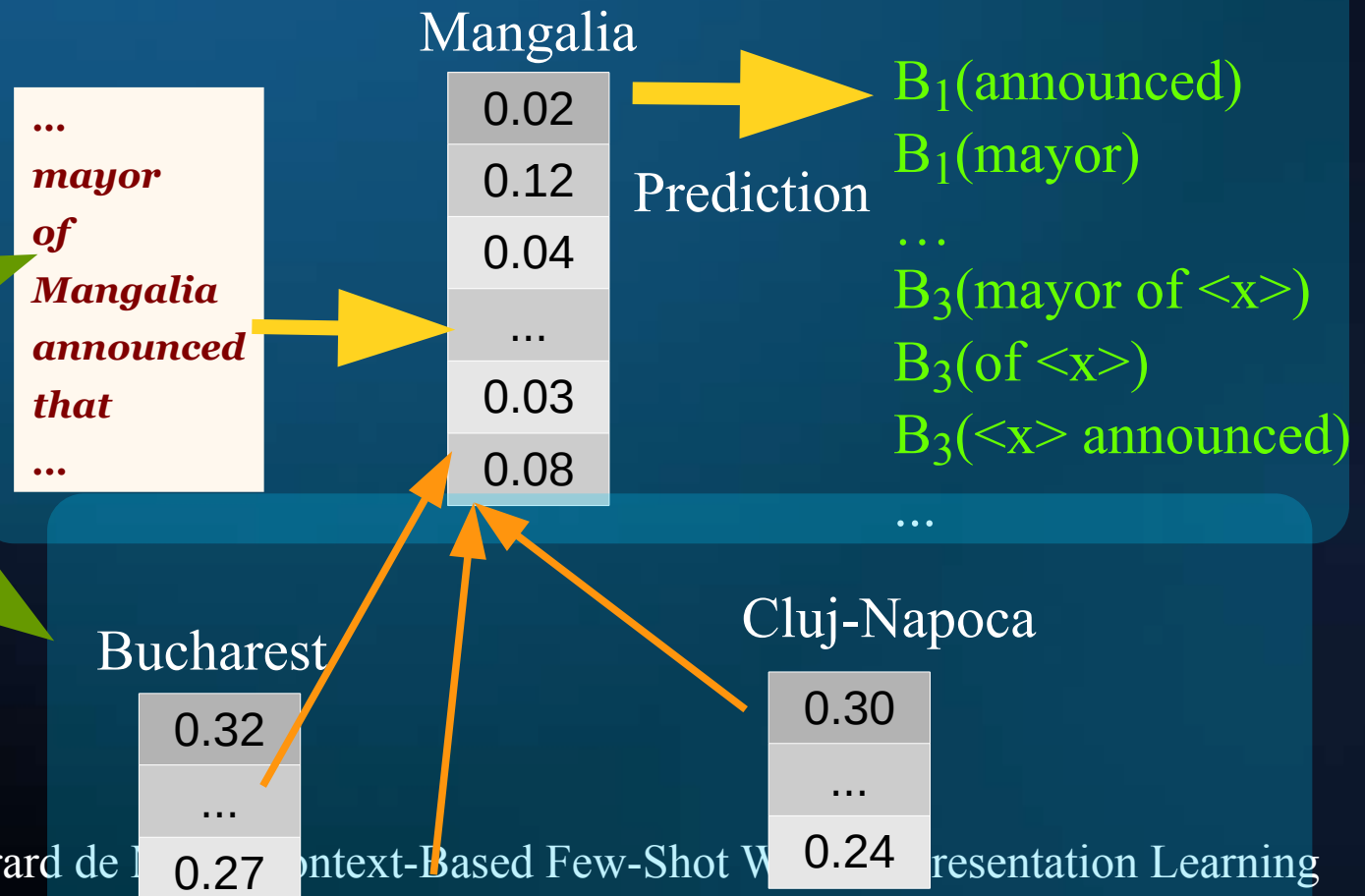


$$f_i = \{f_0 \mid h(f_0) \bmod |F| = i\}$$

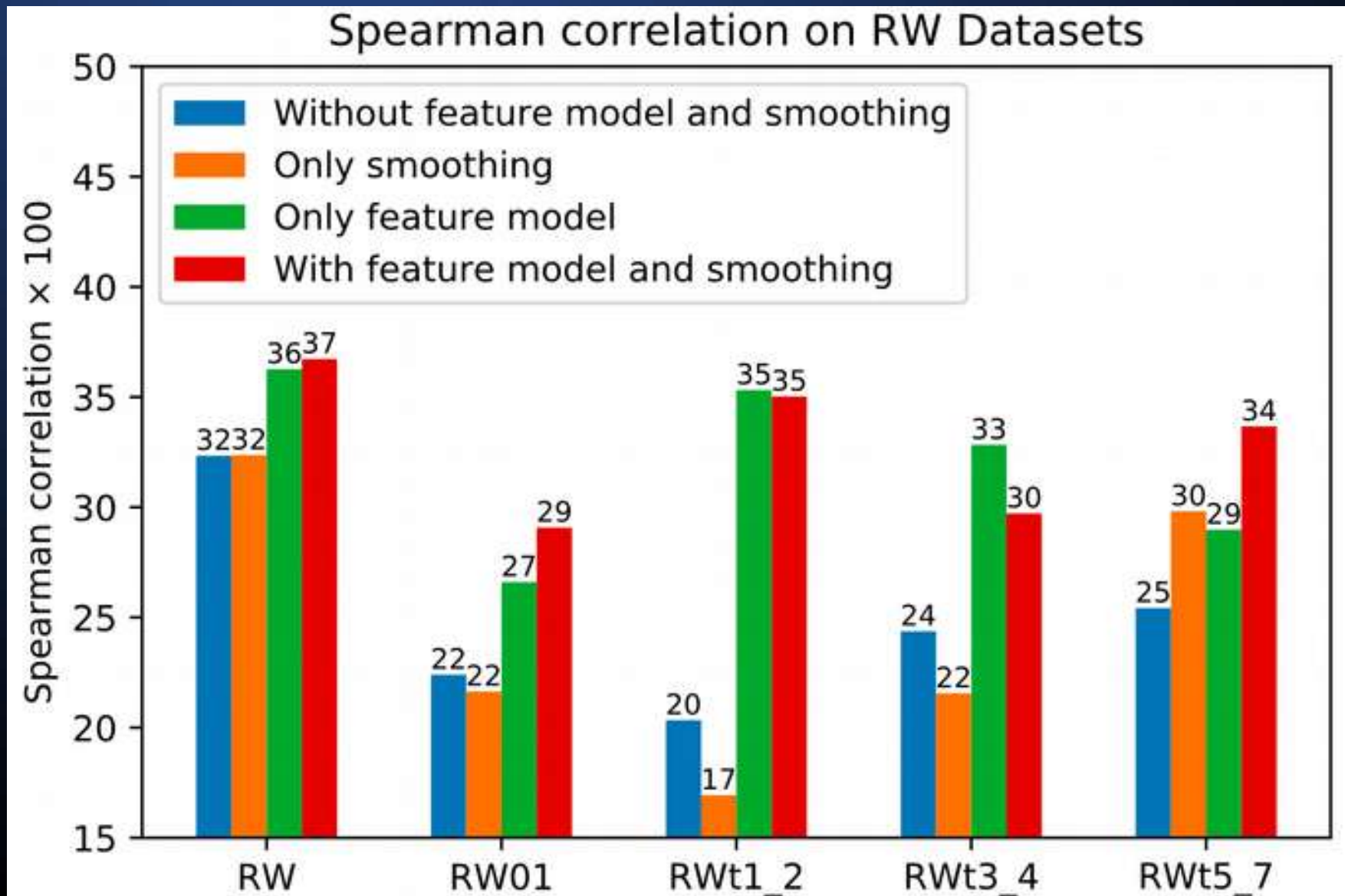
Context-Based Few-Shot Word Representation Learning

Idea 3

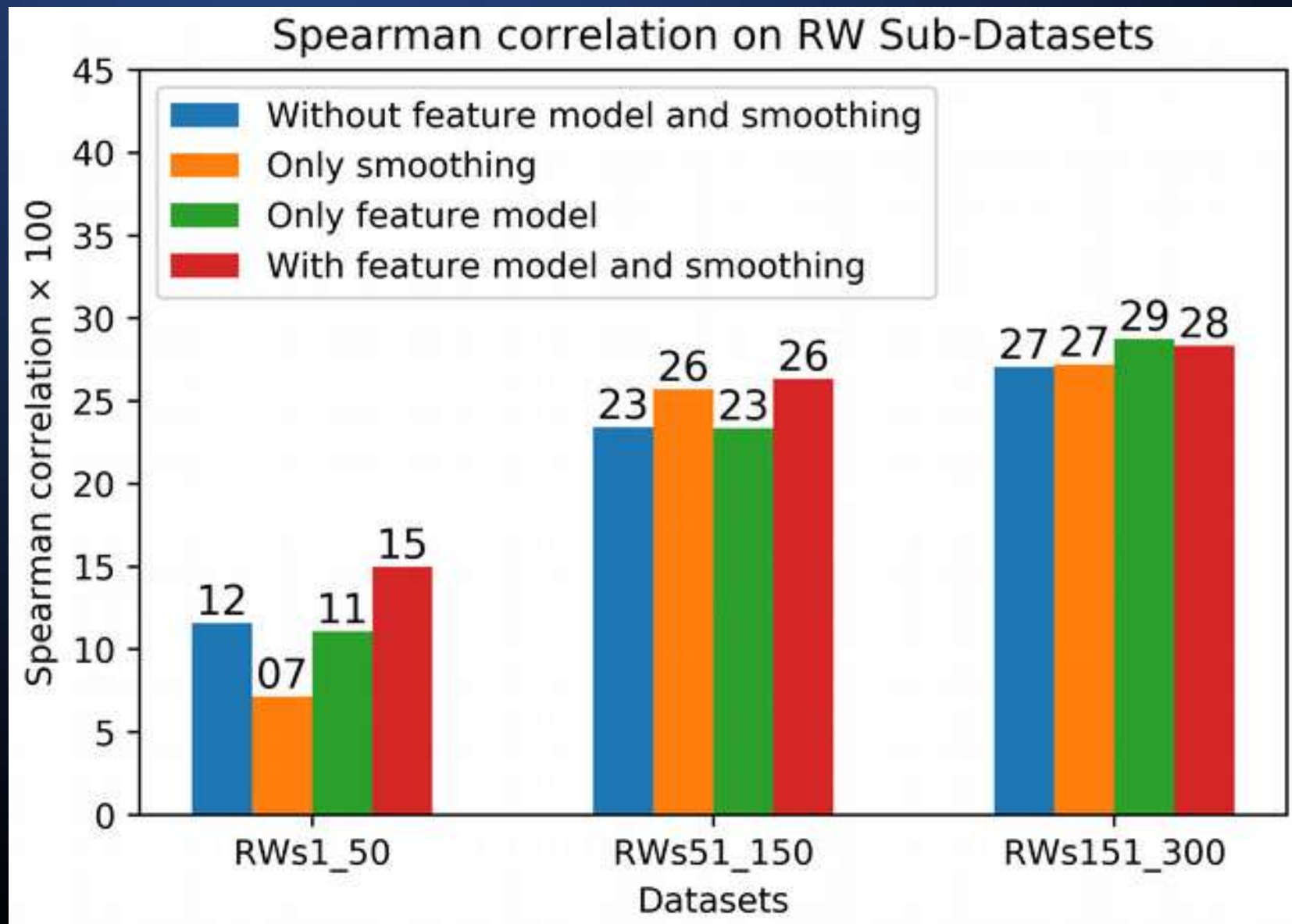
Smoothen vectors of rare words via vectors of more frequent words occurring in similar contexts



Context-Based Few-Shot Word Representation Learning



Context-Based Few-Shot Word Representation Learning



Questions?



Image: <https://www.flickr.com/photos/opensourceway/5556249000>

