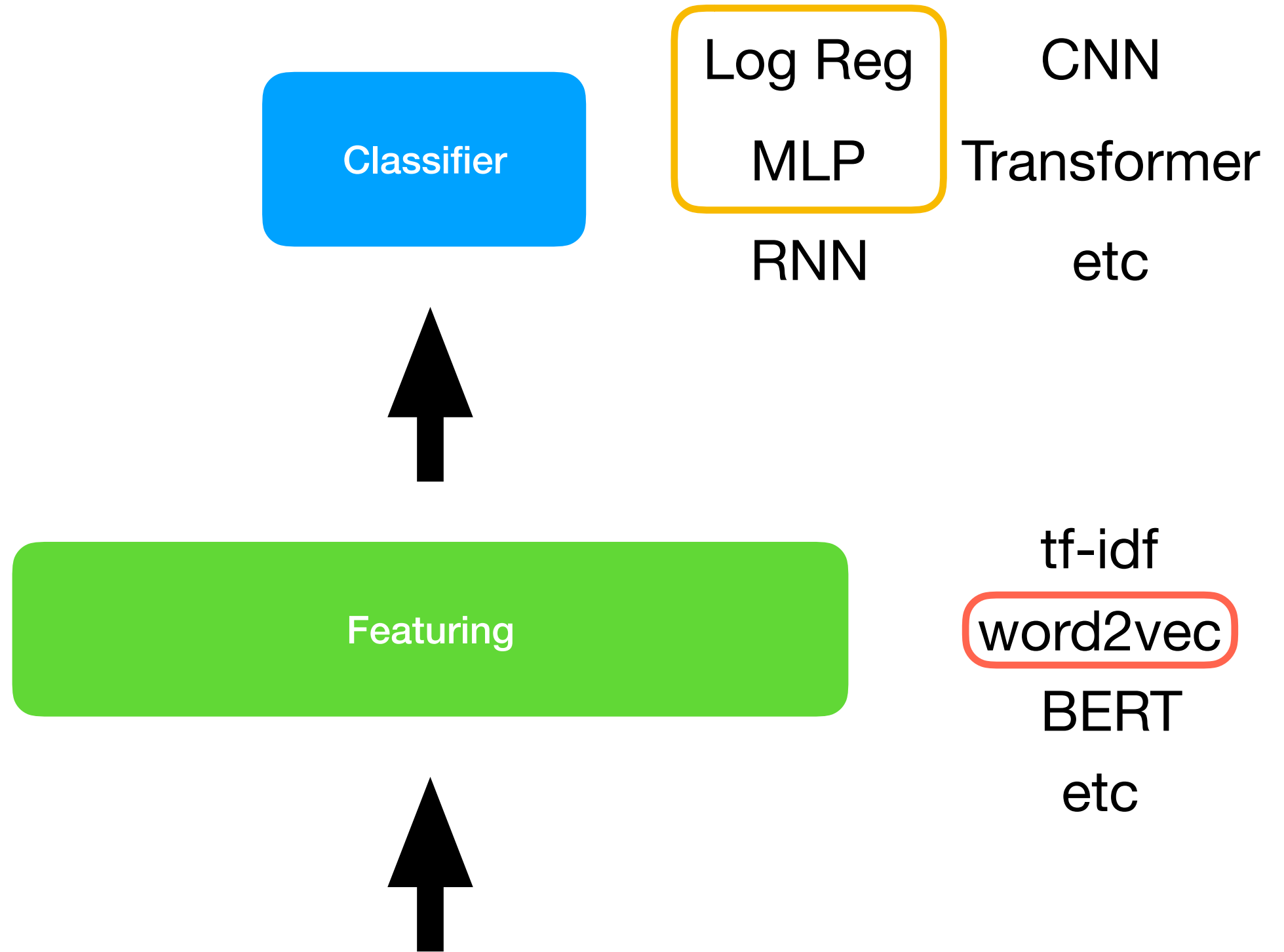


Word Embeddings

Materials MS DL Course of Boris Zubarev @bobazooba

Training Pipeline



The iPhone X is the huge leap forward

One Hot Encoding

Bag of words

```
motel = [0 0 0 0 0 0 0 0 0 0 1 0 0 0 0]  
hotel = [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0]
```

Orthogonal vectors

Dimension = len(vocabulary)

Similarity

Dot Product

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}, \quad [0, 1]$$

Vector Norms

TF-IDF

$$w_{x,y} = tf_{x,y} \times \log \left(\frac{N}{df_x} \right)$$

TF-IDF

Term x within document y

$tf_{x,y}$ = frequency of x in y

df_x = number of documents containing x

N = total number of documents

TF-IDF

$$w_{x,y} = tf_{x,y} \times \log \left(\frac{N}{df_x} \right)$$

TF-IDF

Term x within document y

$tf_{x,y}$ = frequency of x in y

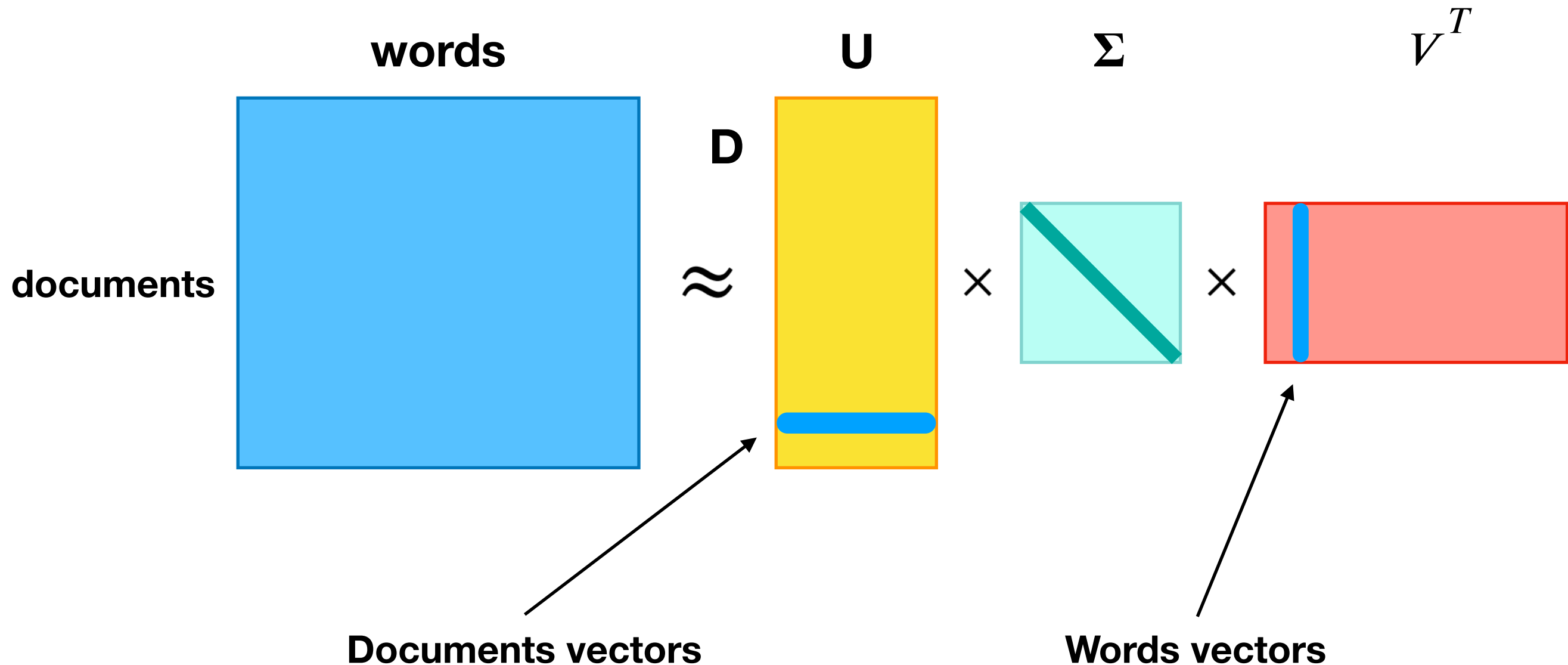
df_x = number of documents containing x

N = total number of documents

text1	0	0	0	0	0.47	0	0.23	0
text2	0	0.68	0	0	0.32	0	0	0
text3	0.11	0	0.19	0	0	0	0	0

Co-occurrence Matrix

$$X \approx \hat{X} = U \Sigma V^T$$



Computational expensive

Co-occurrence Vectors

«You shall know a word by the company it keeps» — Firth, 1957

Corpus sentences

He also found five fish swimming in murky water in an old **bathtub**.

We do abhor dust and dirt, and stains on the bathtub, and any kind of filth.

Above At the far end of the garden room a **bathtub** has been planted with herbs for the winter.

They had been drinking Cisco, a fruity, wine-based fluid that smells and tastes like a mixture of cough syrup and **bathtub** gin.

Science finds that a surface tension on the water can draw the boats together, like toy boats in a **bathtub**.

In fact, the godfather of gloom comes up with a plot that takes in Windsor Davies (the ghost of sitcoms past), a **bathtub** and a big box of concentrated jelly.

'I'll tell him,' said the Dean from the bathroom above the sound of bathwater falling from a great height into the ample Edwardian **bathtub**.

Co-occurrence counts

the	12
a	9
of	7
and	6
in	5
...	...
like	2
water	2
boat	2
from	2
stain	1
toy	1
god-father	1
Cisco	1
...	...

vector

$$\begin{pmatrix} 12 \\ 9 \\ 7 \\ 6 \\ 5 \\ \vdots \\ 2 \\ 2 \\ 2 \\ 2 \\ 1 \\ 1 \\ 1 \\ 1 \\ \vdots \end{pmatrix}$$

Dimensionality reduction

small vector

[illegible]

Co-occurrence Matrix

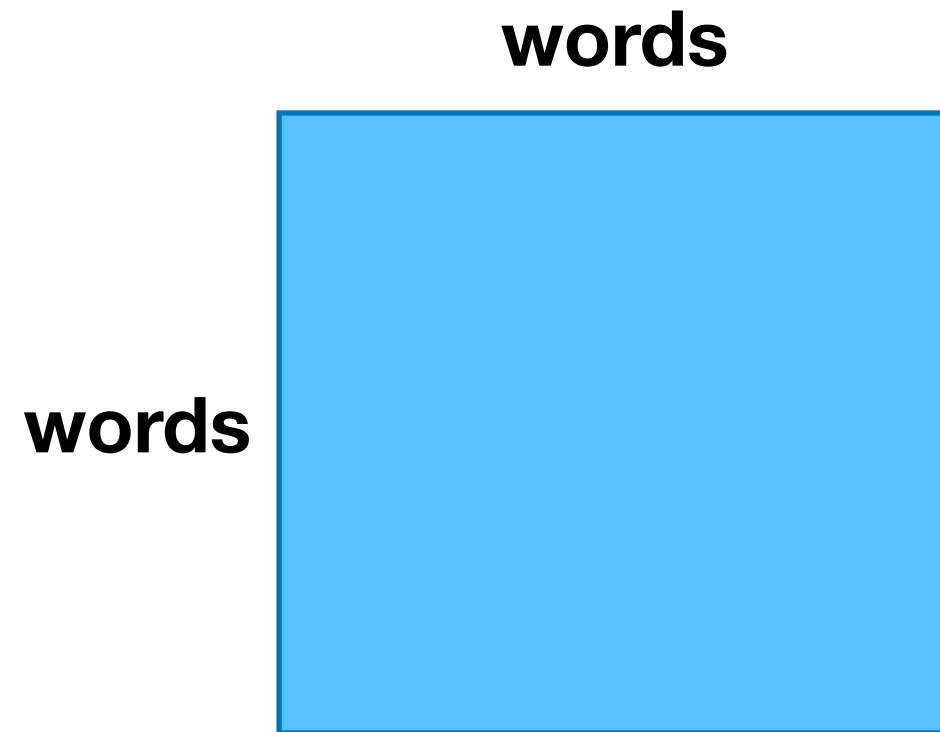
words



words

№	Словосочетание	Документы	Частота
1	и не	22732	201352
2	и в	27048	193983
3	потому что	14926	117401
4	я не	10675	113767
5	у меня	9734	97102
6	может быть	16086	96065
7	то что	17195	95251
8	что он	11786	92743
9	не было	13196	92729
10	в том	21604	89842

Co-occurrence Matrix

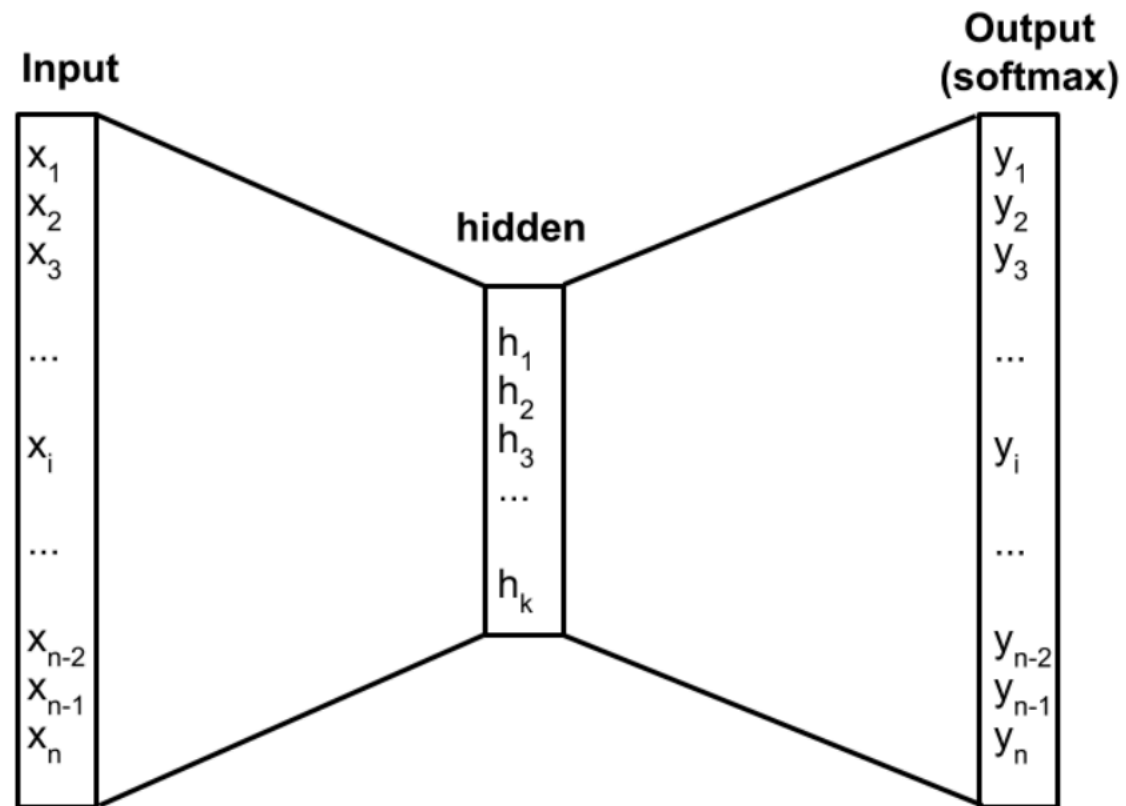


$$\text{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)}$$

$$\text{ppmi} = \max(\text{pmi}, 0)$$

Word2Vec

Word2Vec

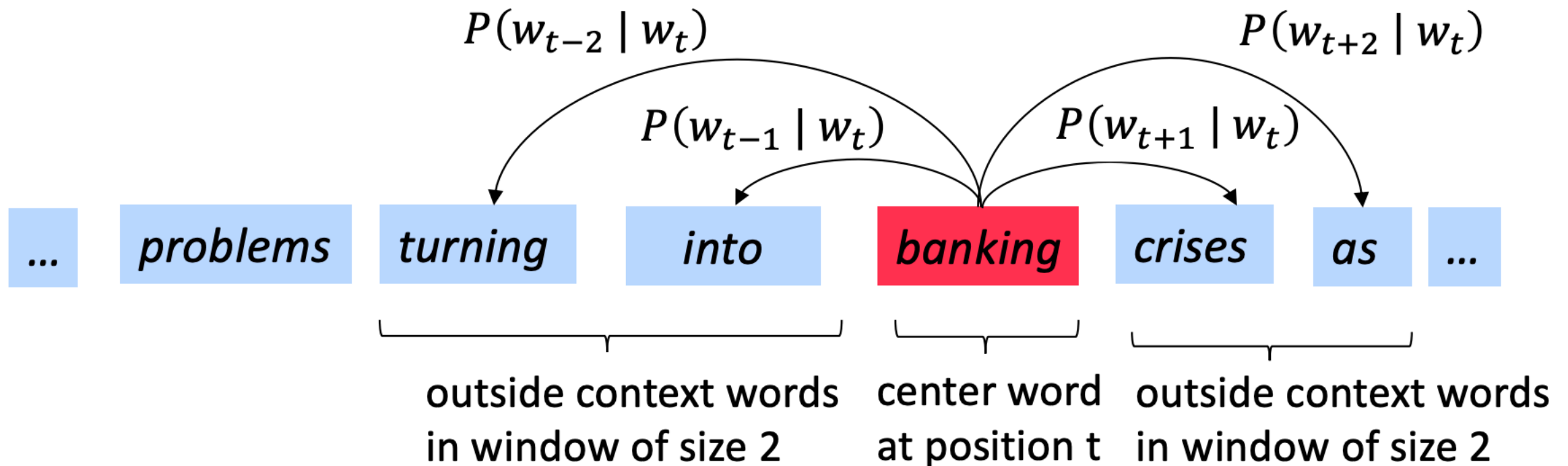
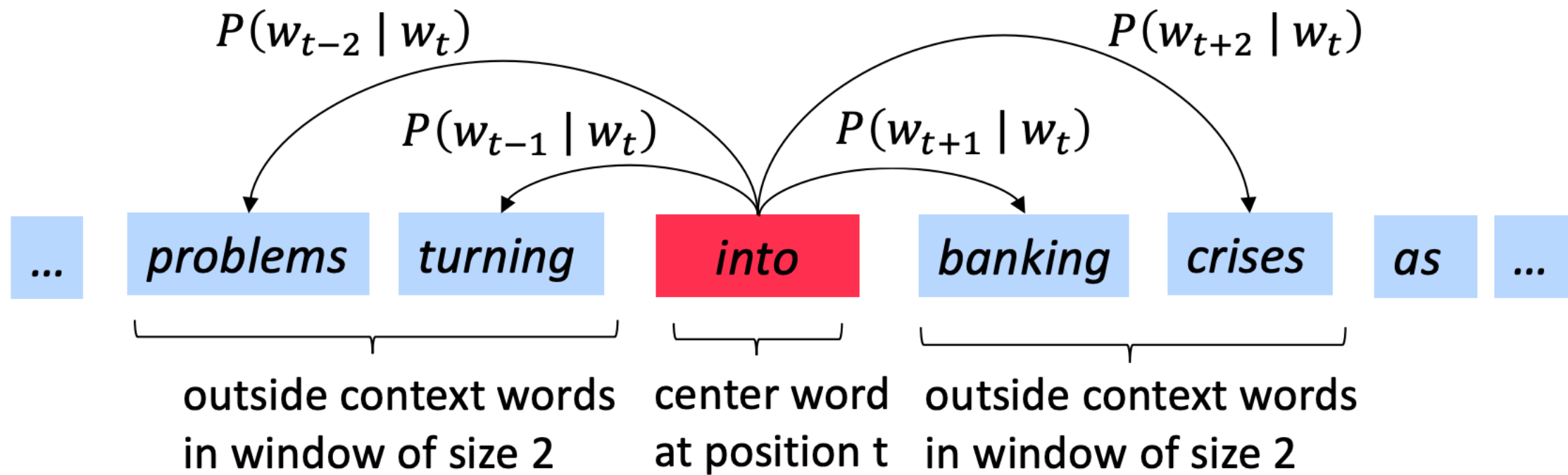


$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

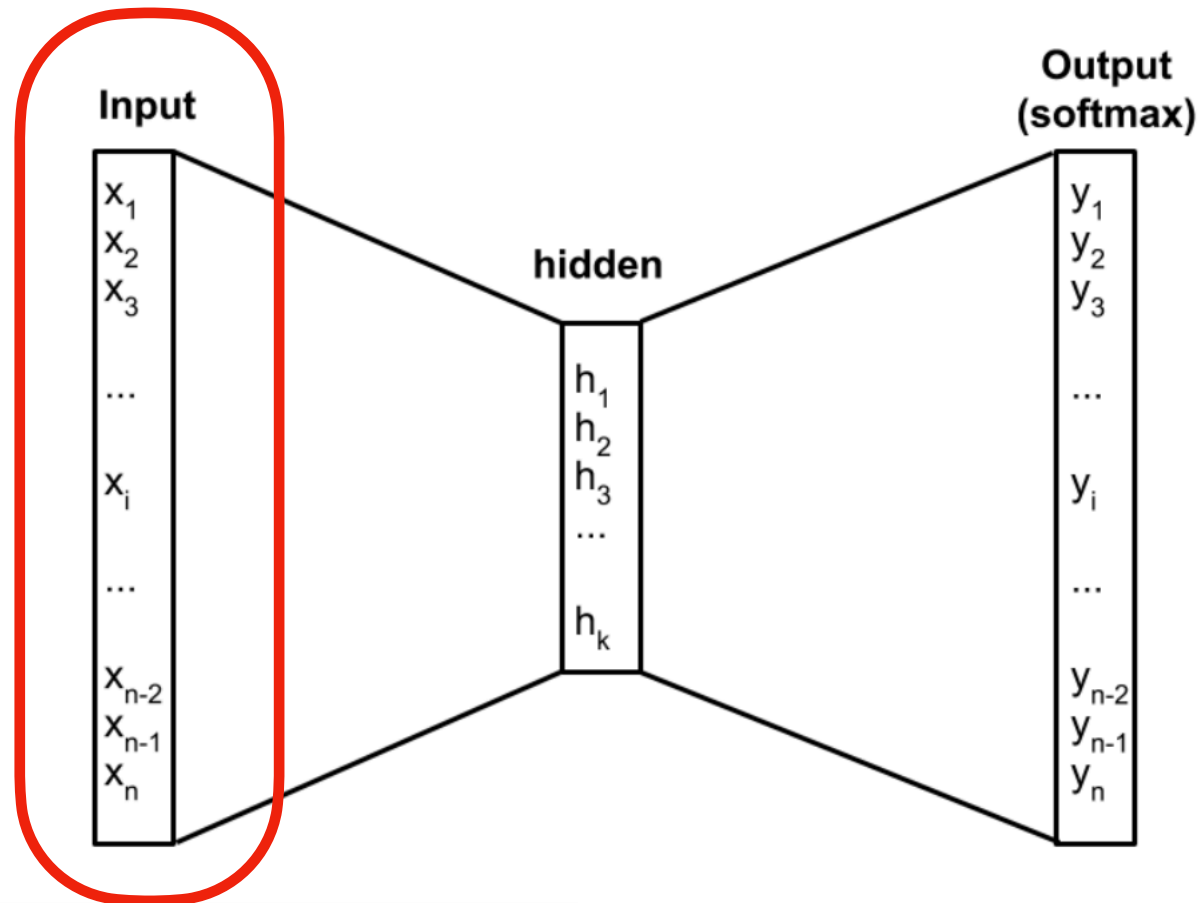
Word2Vec

Source Text	Training Samples						
<table><tr><td>The</td><td>quick</td><td>brown</td></tr></table> fox jumps over the lazy dog. ➡	The	quick	brown	(the, quick) (the, brown)			
The	quick	brown					
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td></tr></table> jumps over the lazy dog. ➡	The	quick	brown	fox	(quick, the) (quick, brown) (quick, fox)		
The	quick	brown	fox				
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td></tr></table> over the lazy dog. ➡	The	quick	brown	fox	jumps	(brown, the) (brown, quick) (brown, fox) (brown, jumps)	
The	quick	brown	fox	jumps			
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td><td>over</td></tr></table> the lazy dog. ➡	The	quick	brown	fox	jumps	over	(fox, quick) (fox, brown) (fox, jumps) (fox, over)
The	quick	brown	fox	jumps	over		

Word2Vec

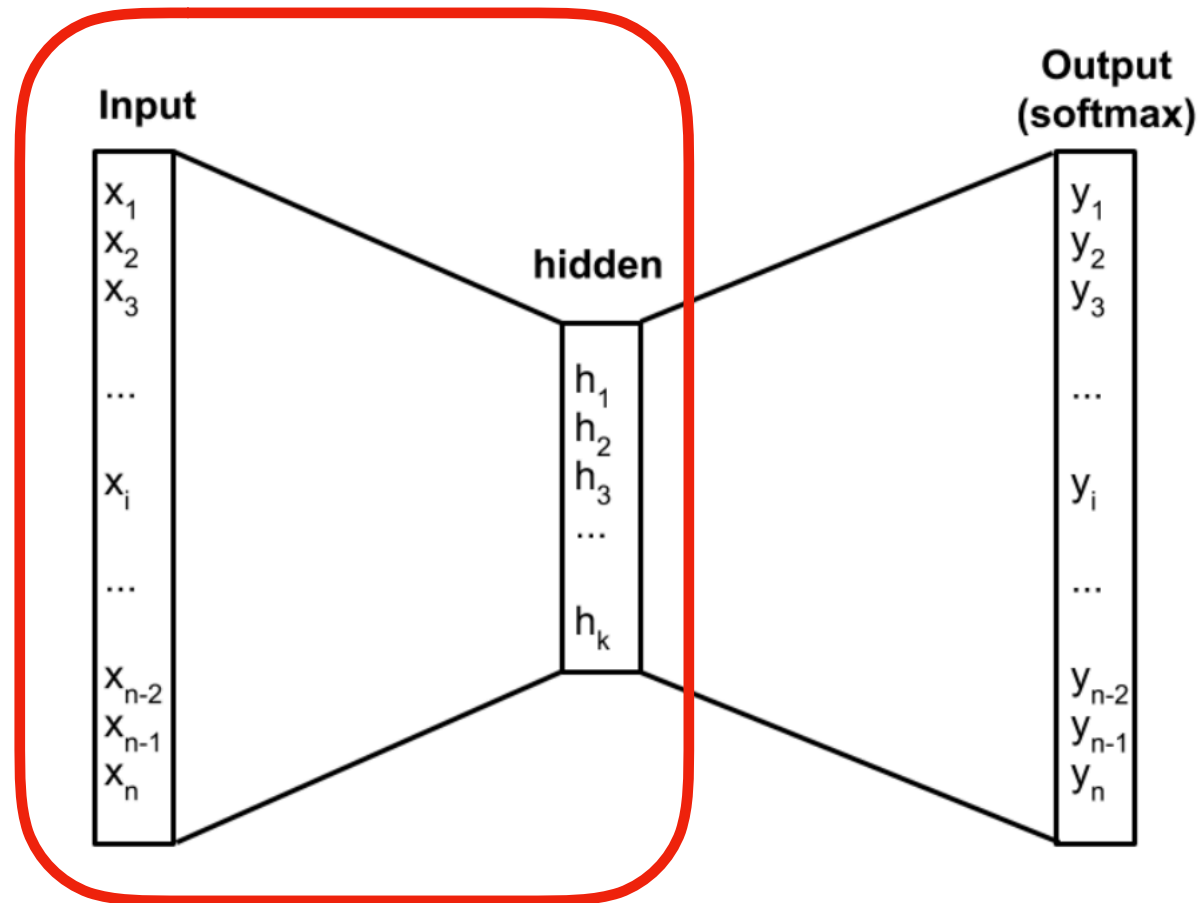


Word2Vec



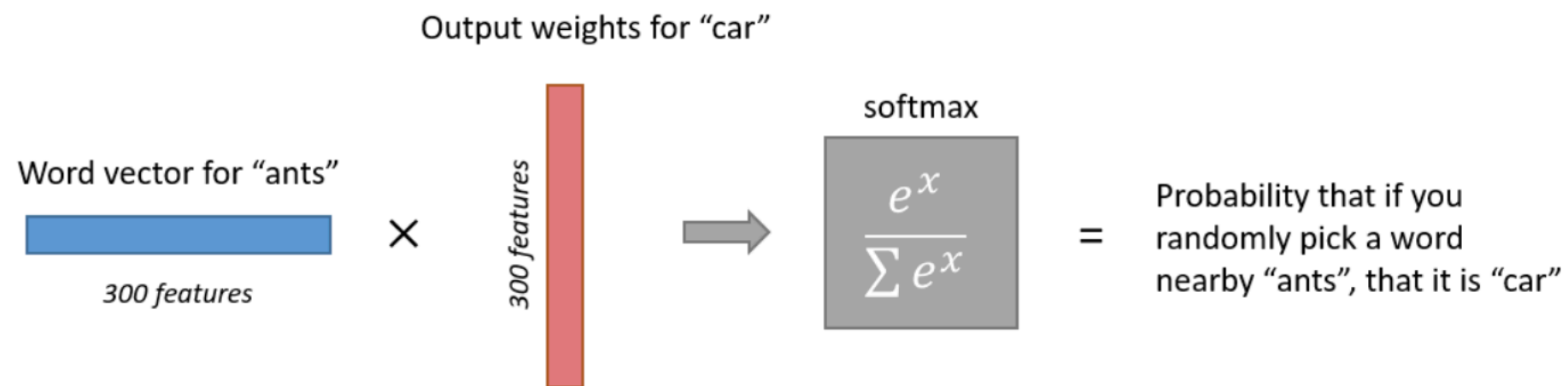
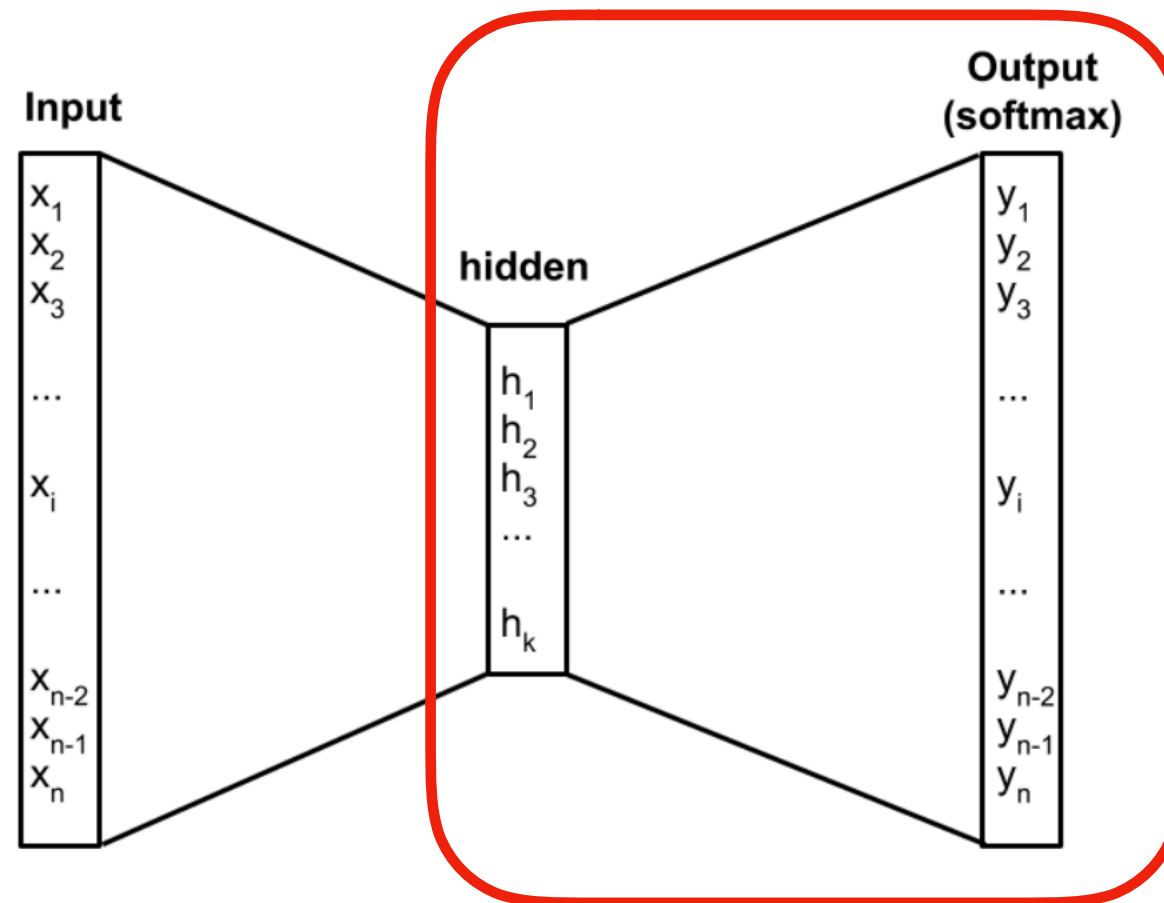
$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$

Word2Vec

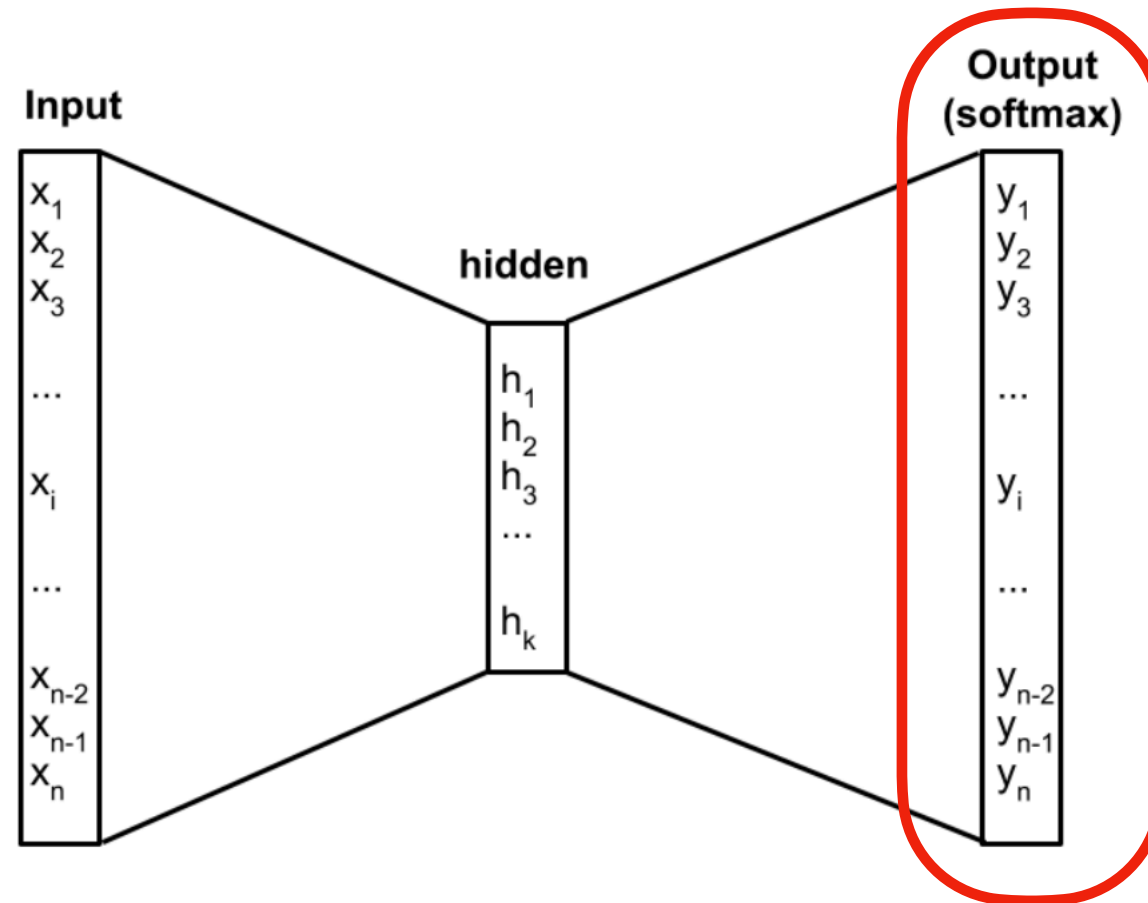


$$[0 \quad 0 \quad 0 \quad \mathbf{1} \quad 0] \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ \mathbf{10} & \mathbf{12} & \mathbf{19} \\ 11 & 18 & 25 \end{bmatrix} = [10 \quad 12 \quad 19]$$

Word2Vec

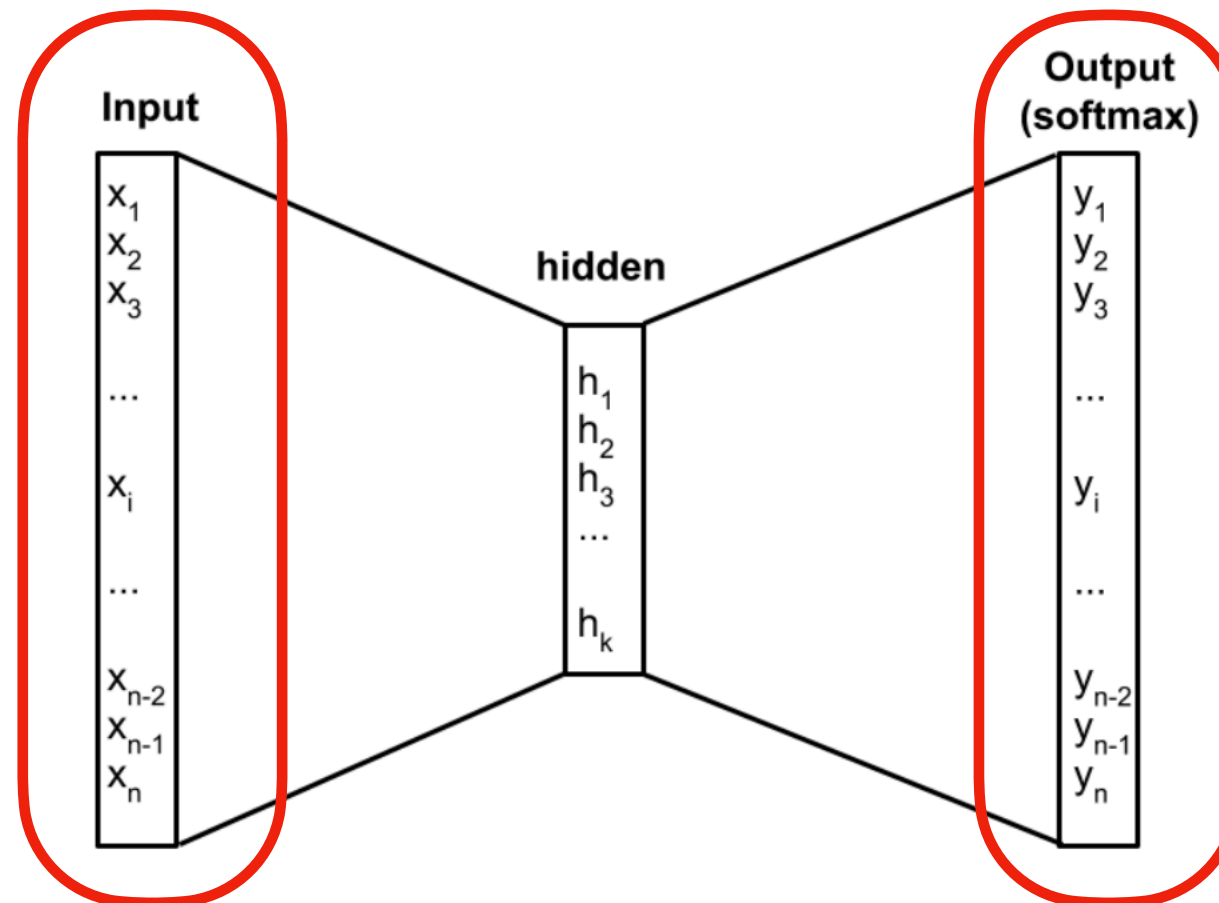


Word2Vec



$$p(w_O|w_I) = \frac{\exp \left(v'_{w_O}{}^\top v_{w_I} \right)}{\sum_{w=1}^W \exp \left(v'_w{}^\top v_{w_I} \right)}$$

Word2Vec



$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV} \quad \begin{matrix} \text{Same words} \\ \\ \\ \text{Different vectors} \end{matrix} \quad \theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$

Word2Vec

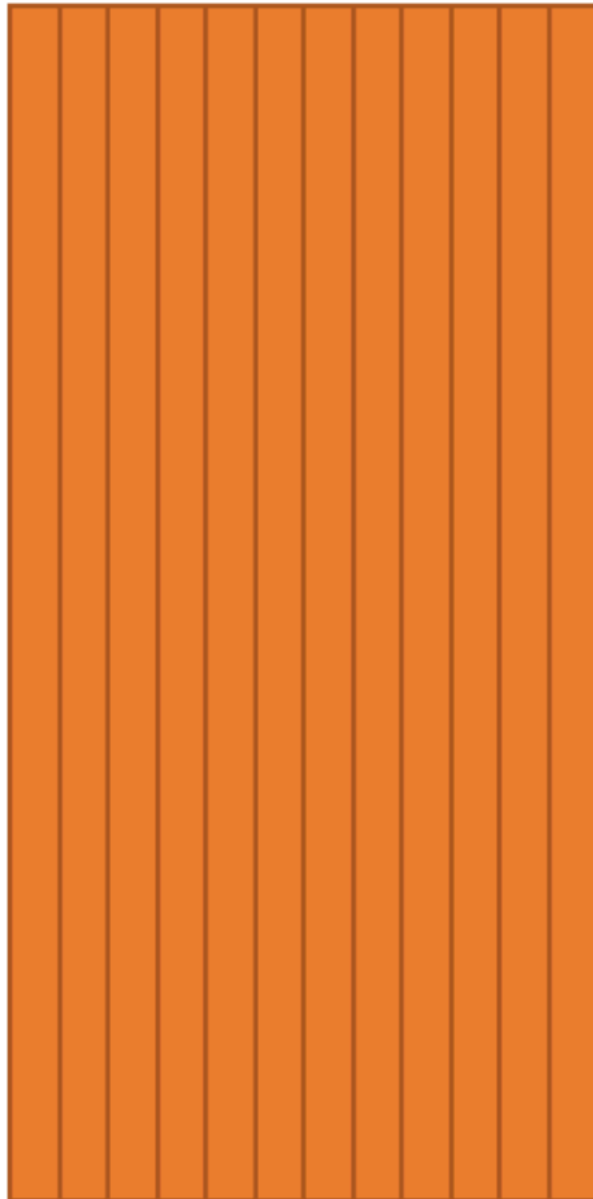
Hidden Layer
Weight Matrix



*Word Vector
Lookup Table!*

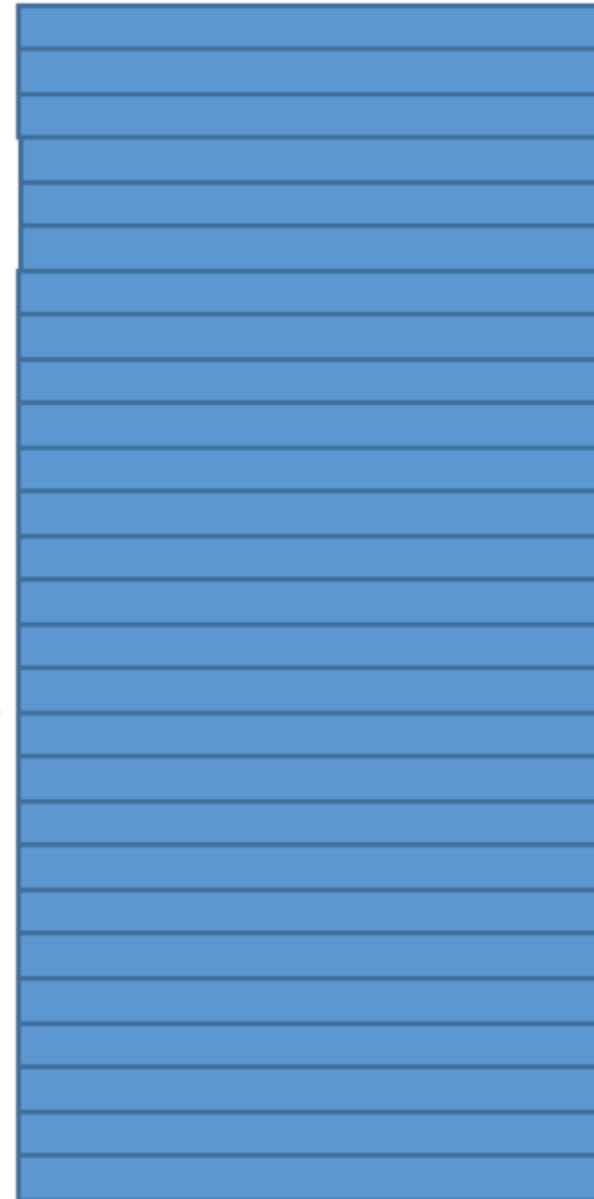
300 neurons

10,000 words

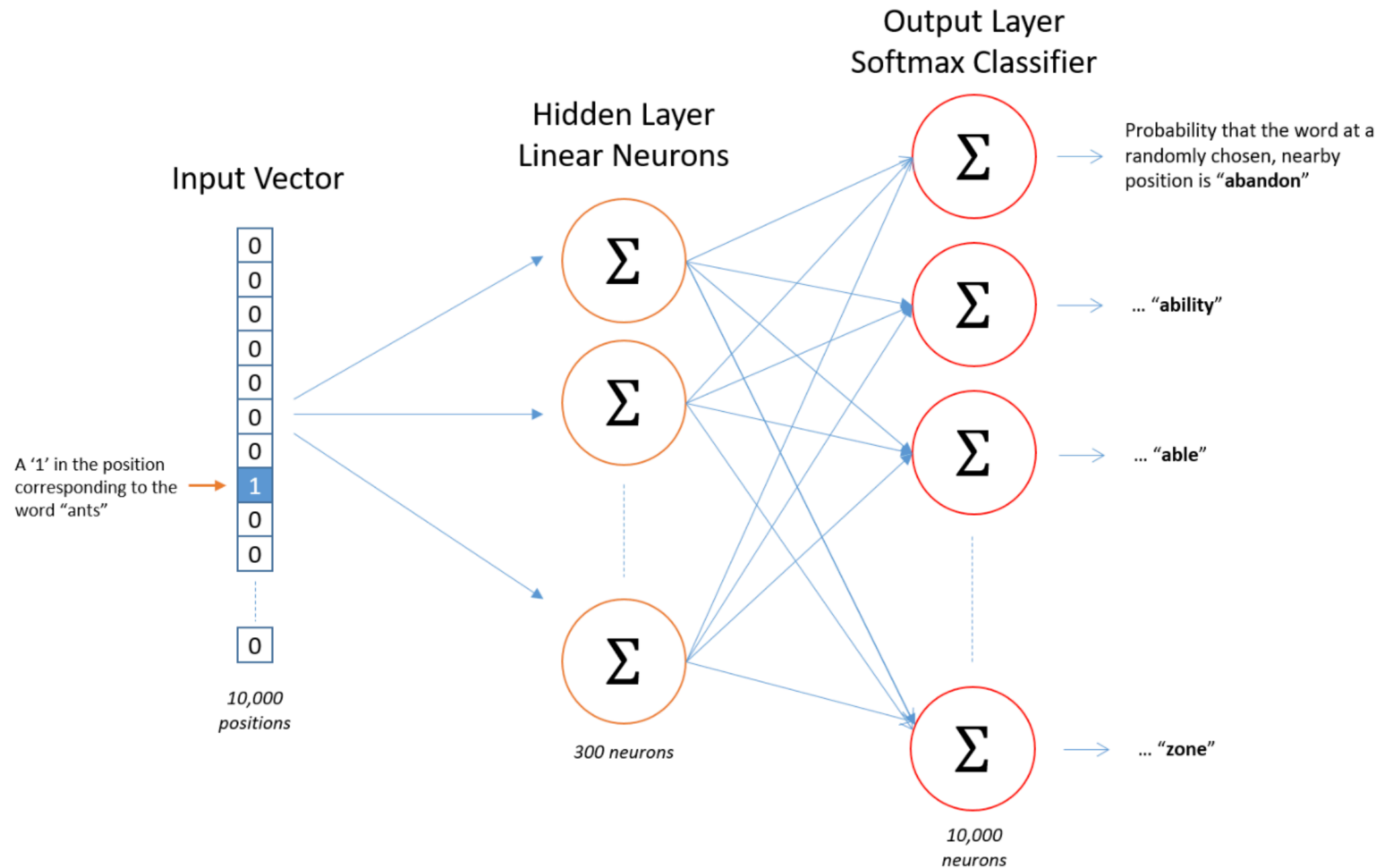


300 features

10,000 words



Word2Vec

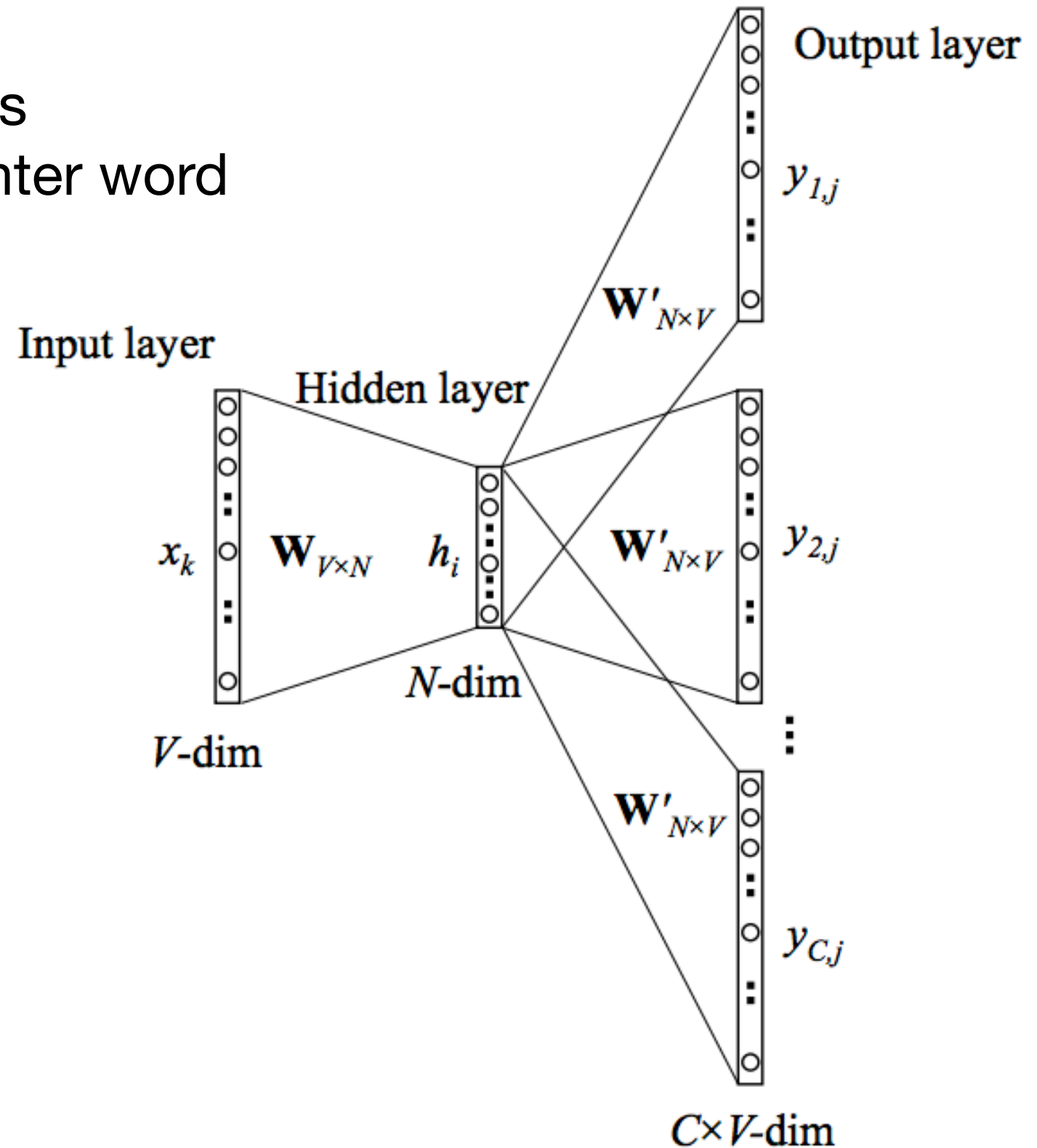


Word2Vec

Skipgrams

Predict context ("outside") words
(position independent) given center word

Better for rare words

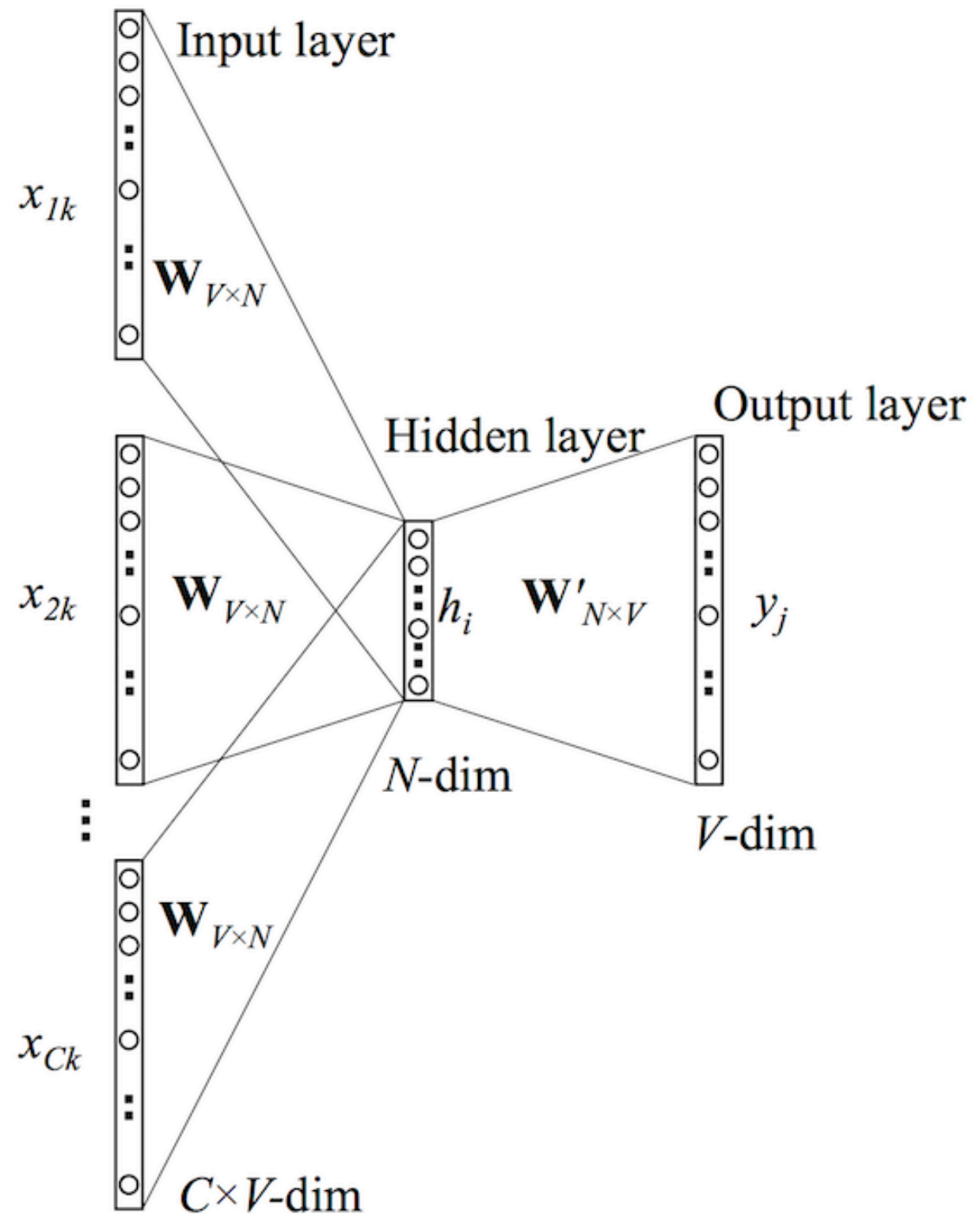


Word2Vec

CBOW

Predict center word from
(bag of) context words

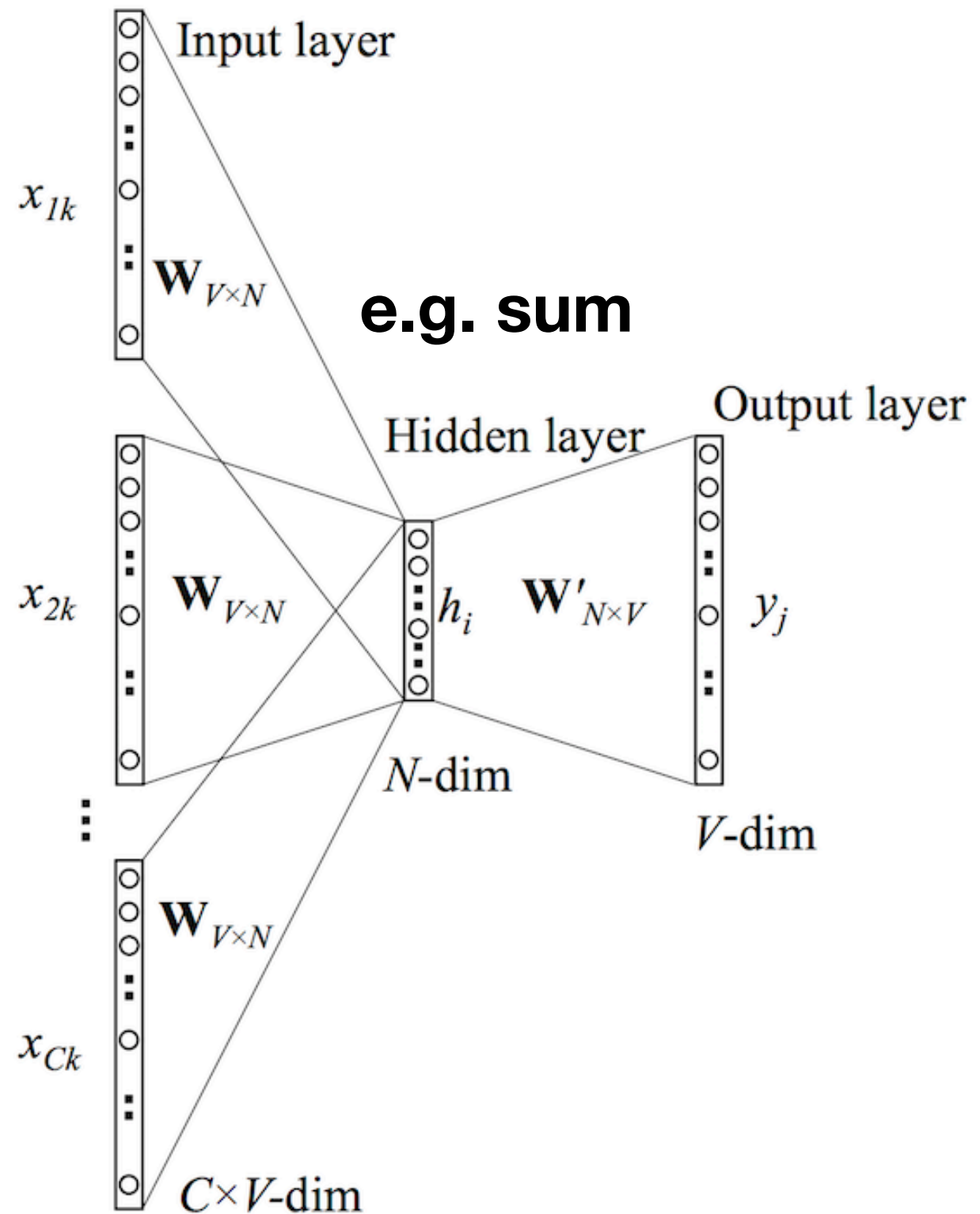
Faster



Word2Vec

CBOW

Predict center word from
(bag of) context words



Visualization

<https://projector.tensorflow.org/>

- BERT Embedding Projector


Visualization



Частотность слова

☒ Высокая ☒ Средняя ☐ Низкая

НКРЯ и Wikipedia

1. **англия** PROPN 0.58 
2. **европа** PROPN 0.54 
3. **великобритания** PROPN 0.52 
4. **страна** NOUN 0.48 
5. **франция** PROPN 0.47 

Visualization

Some vector close to queen

$$\text{word2vec}(\text{king}) - \text{word2vec}(\text{man}) + \text{word2vec}(\text{woman}) = \text{word2vec}(\text{queen})$$

Word2Vec

Fasttext

OOV

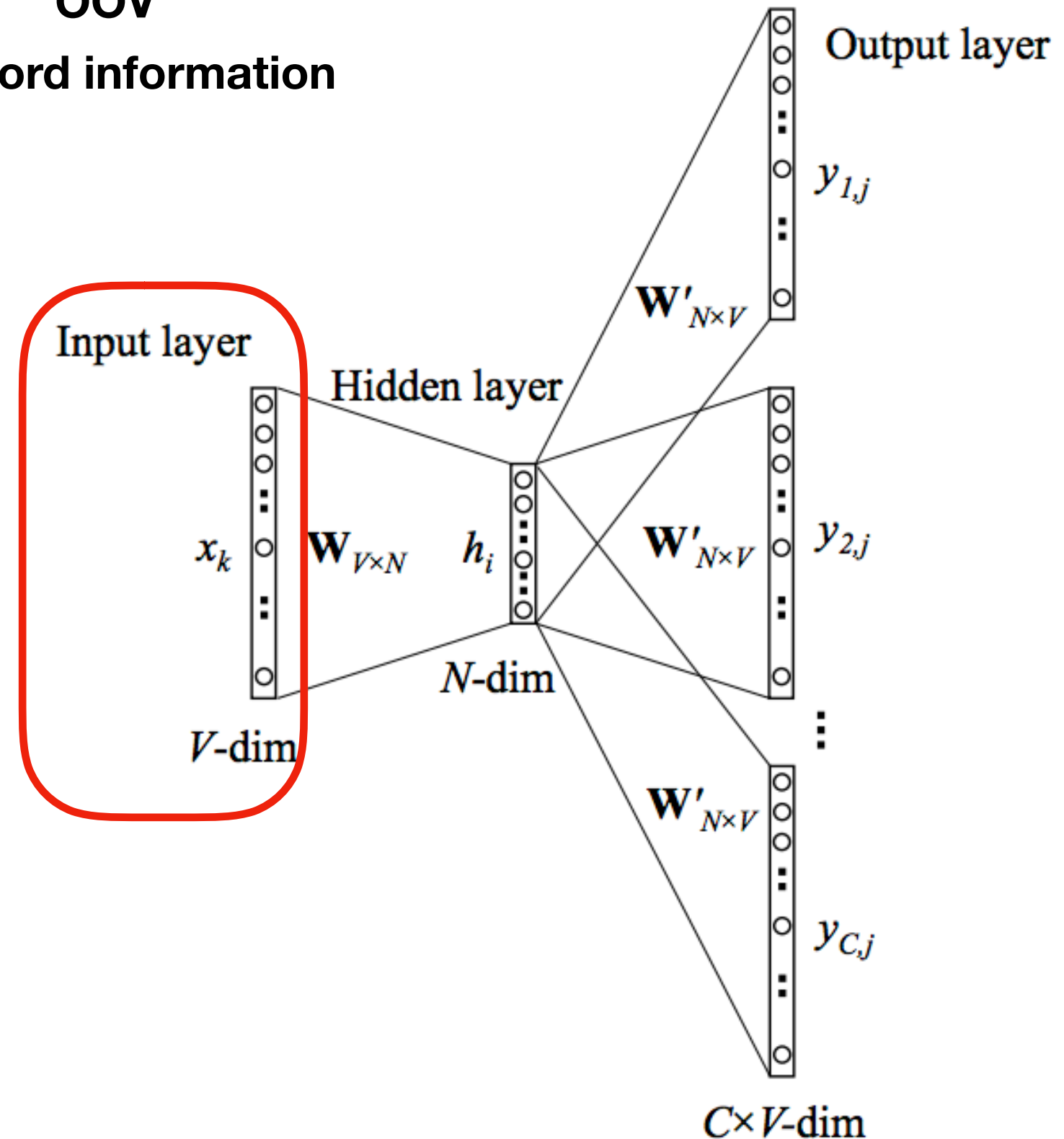
Subword information

where

=

<wh + whe + her + ere + re>

3 — 6 char n-gram length



Word2Vec

$$p(w_O|w_I) = \frac{\exp \left(v'_{w_O}{}^\top v_{w_I} \right)}{\sum_{w=1}^W \exp \left(v'_w{}^\top v_{w_I} \right)}$$

Word2Vec

$$p(w_O | w_I) = \frac{\exp \left(v'_{w_O}{}^\top v_{w_I} \right)}{\sum_{w=1}^W \exp \left(v'_w{}^\top v_{w_I} \right)}$$

Computational expensive



Word2Vec

- Hierarchical softmax
- Naive softmax
- Subset of vocabulary
- Negative sampling
- Binary classification

$$p(w_O | w_I) = \frac{\exp(v'_{w_O}{}^\top v_{w_I})}{\sum_{w=1}^W \exp(v'_w{}^\top v_{w_I})}$$

Word2Vec

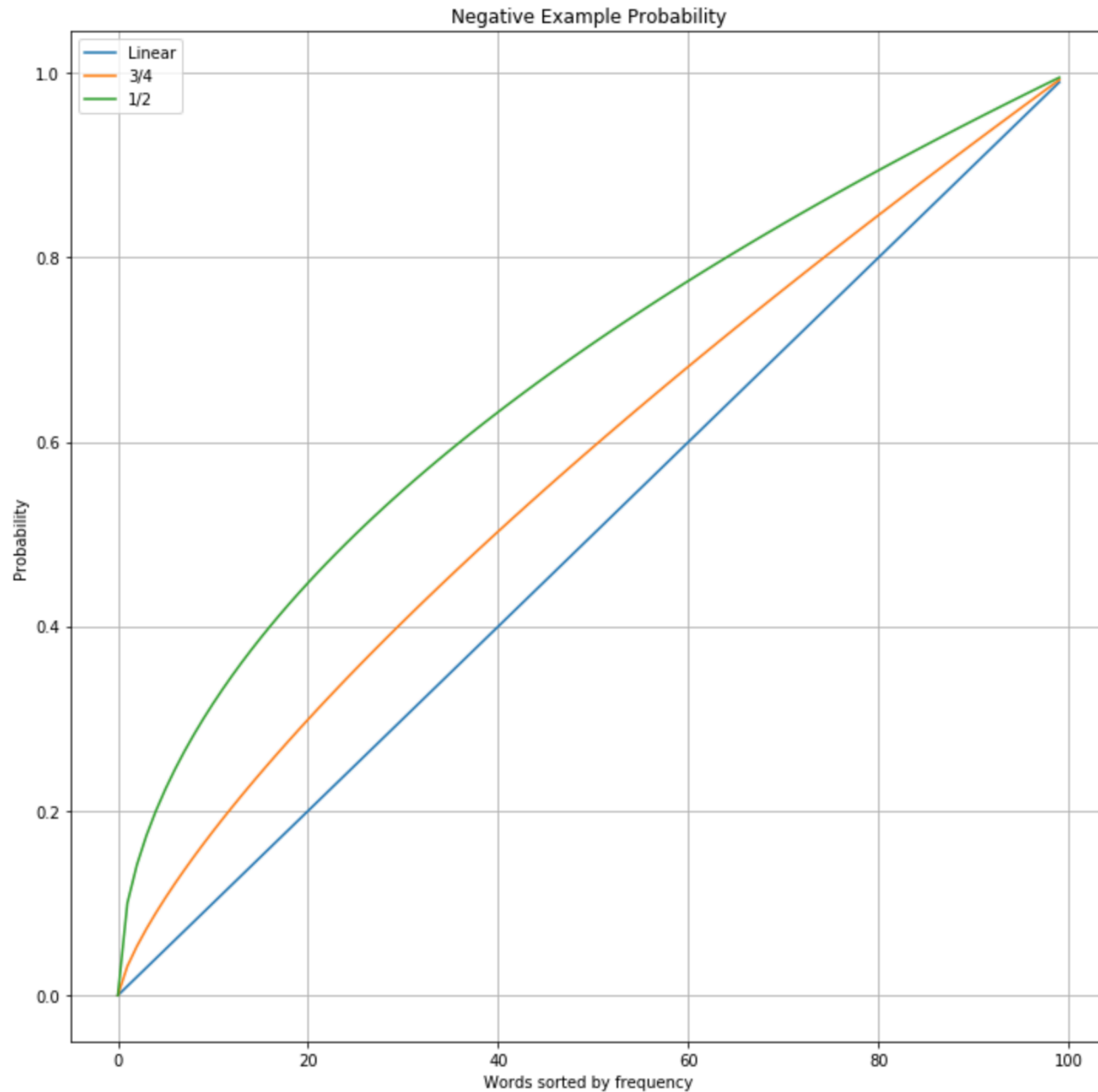
Negative sampling

$$J_{neg-sample}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$$

Sampling negatives

P(word) 3/4 e.g. K = 5
Increase rare word probability

Word2Vec



Word2Vec

Negative sampling

$$J_{neg-sample}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$$

Sampling negatives

$$P(\text{word})^{3/4}$$

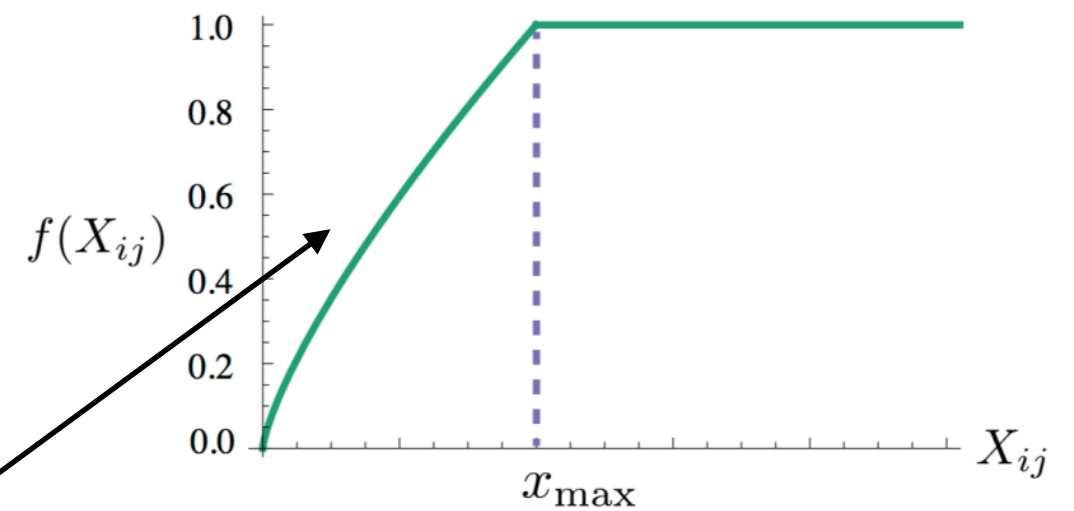
Subsampling frequent words

$$P(w_i) = \frac{10^{-3}}{p_i} \left(\sqrt{10^3 p_i} + 1 \right)$$

Removing pairs

GloVe

$$J(\theta) = \frac{1}{2} \sum_{i,j=1}^W f(P_{ij}) (u_i^T v_j - \log P_{ij})^2$$



Less influence of rare words

GloVe

nearest neighbors of
frog

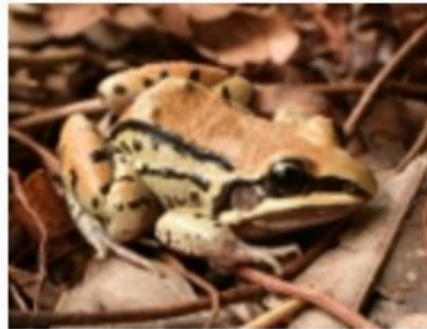
Litoria

Leptodactylidae

Rana

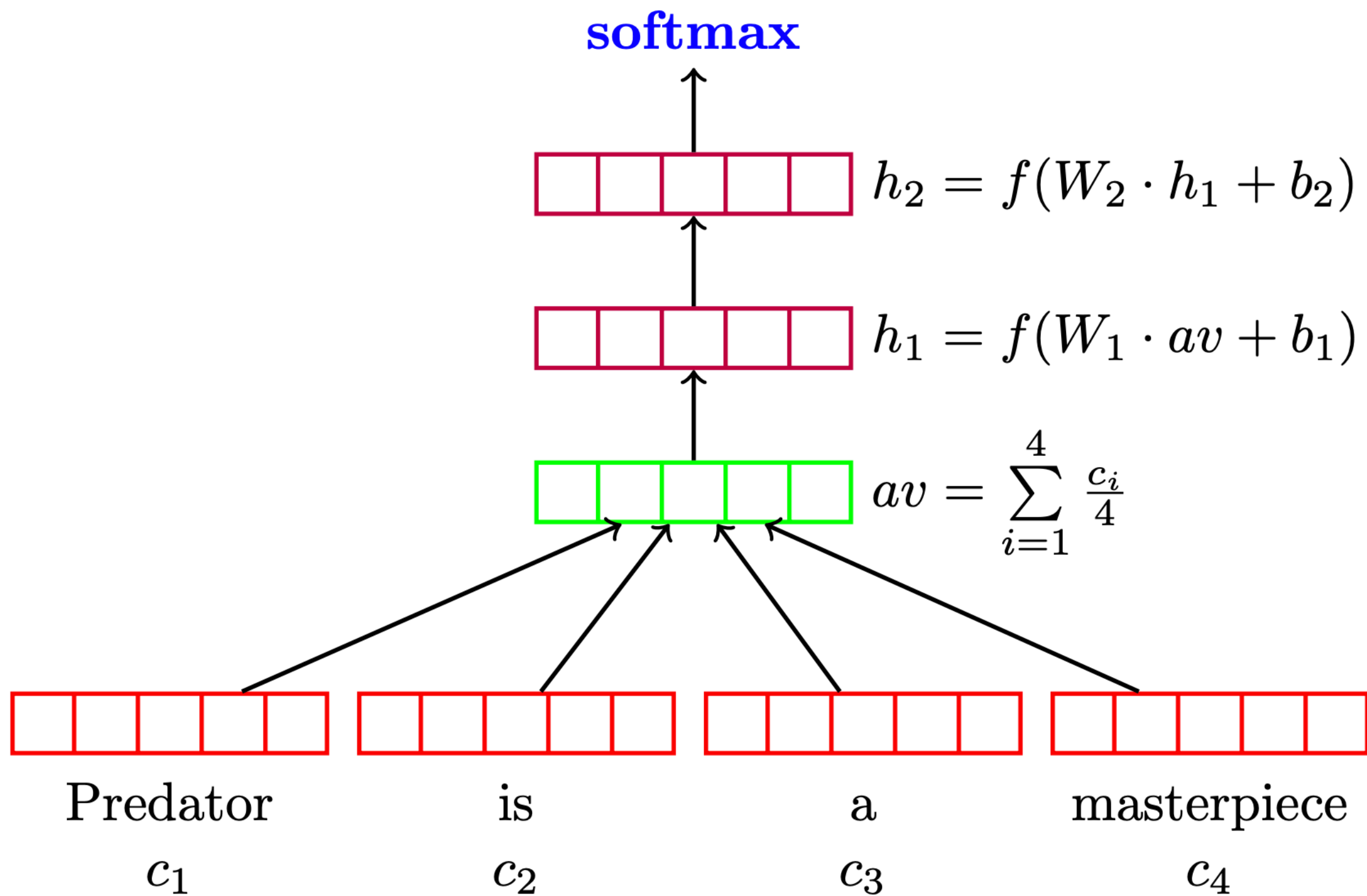
Eleutherodactylus

Pictures



How to choose embeddings?

DAN



DAN

