

Prev: BOW }

TFIDF }

Task ↗
Agenda: Words → Vectors
d-dim
(capture Semantics)
e.g: good; great;
best

Plan: Let us build/invent techniques
together
(by building on ideas/concepts we
already know)

Efficient Estimation of Word Representations in Vector Space

GloVe: Global Vectors for Word Representation

How can we solve for this problem?

- Can we match keywords from user queries that are present in abstract?
 - If we do keyword matching, will we be able to understand the user's intent? For e.g: 'origin' and 'discovery'
 - Should we consider the context of the words, then?

Let us build a search engine using Word Embeddings

Datset

Dataset Corpus

Dataset

$d_1 \rightarrow w_3 \ w_6 \ w_5 \ \tilde{w_3} \ \dots$

$d_2 \rightarrow w_1 \ w_{12} \ w_6 \ w_3 \ \dots$

\vdots

$d_n \rightarrow \dots$

$w_i \rightarrow \delta\text{-dim}$
~~rep~~
(semantic)

Semantically sensible

v_1
good

v_2
great

$v_1 \& v_2 \rightarrow \underline{\underline{dim}}$

$$\text{Cosine-Sim}(v_1, v_2) \rightarrow 1$$

good , terrible
 $\downarrow v_3$

$$\text{Cos-Sim}(v, v_3) : \underline{\underline{low}}$$

Task:

$w_i \rightarrow v_i \in \mathbb{R}^d$ using corpus

Brain storm:

1.

Clustering (?)

→ where are the vectors to cluster?

2.

Time-series

layer ... RNN | LSTM | GRU

prediction



SARIMAX

(wigs to scalars?)

3.

PCA → [data matrix?]

HINT: Combine 1 & 3

MF

MF
Recap:

Data Matrix

$$A_{n \times m} = B_{n \times d} \cdot C_{d \times m}$$

Users $u_i \rightarrow$ []
Movies [] \downarrow

$u_i \left[\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} \right]_{n \times d}$

d [] m_j

Q: what is my data matrix in NLP task

doc-term occurrence

	w ₁	w ₂	...	w _N
d ₁	1	2	-	-
d ₂	0	2	-	-
:	:			
d _n				

Corpus, n-docs
≠

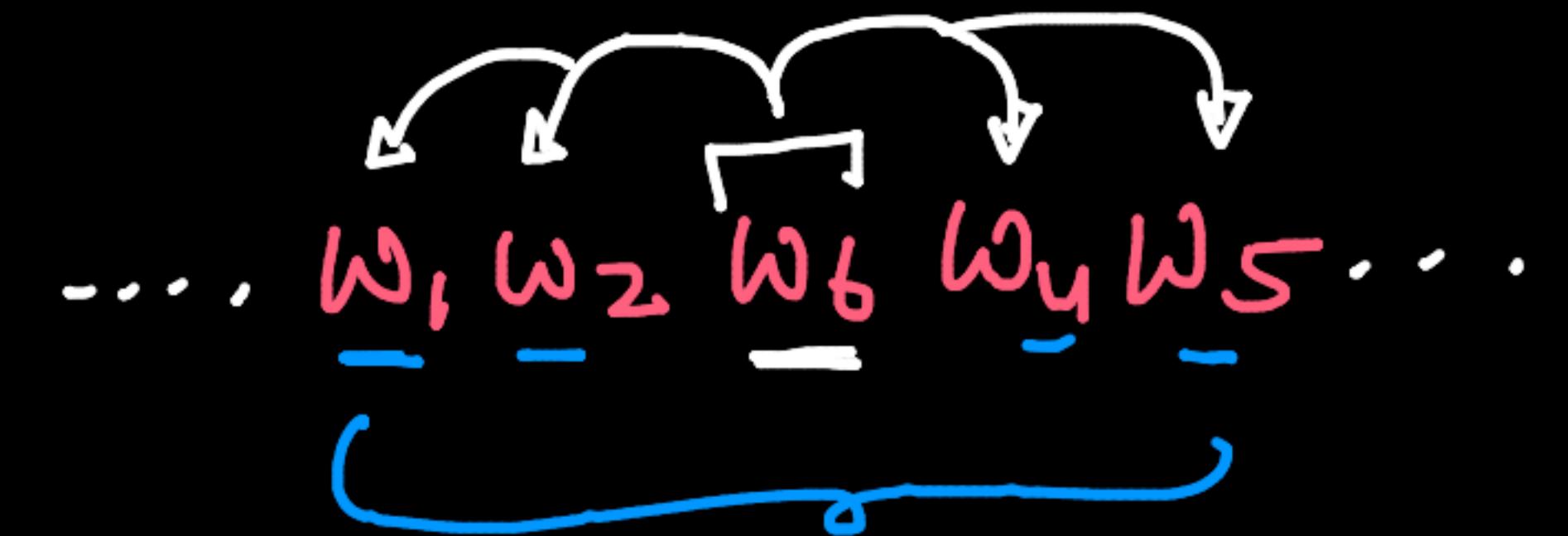
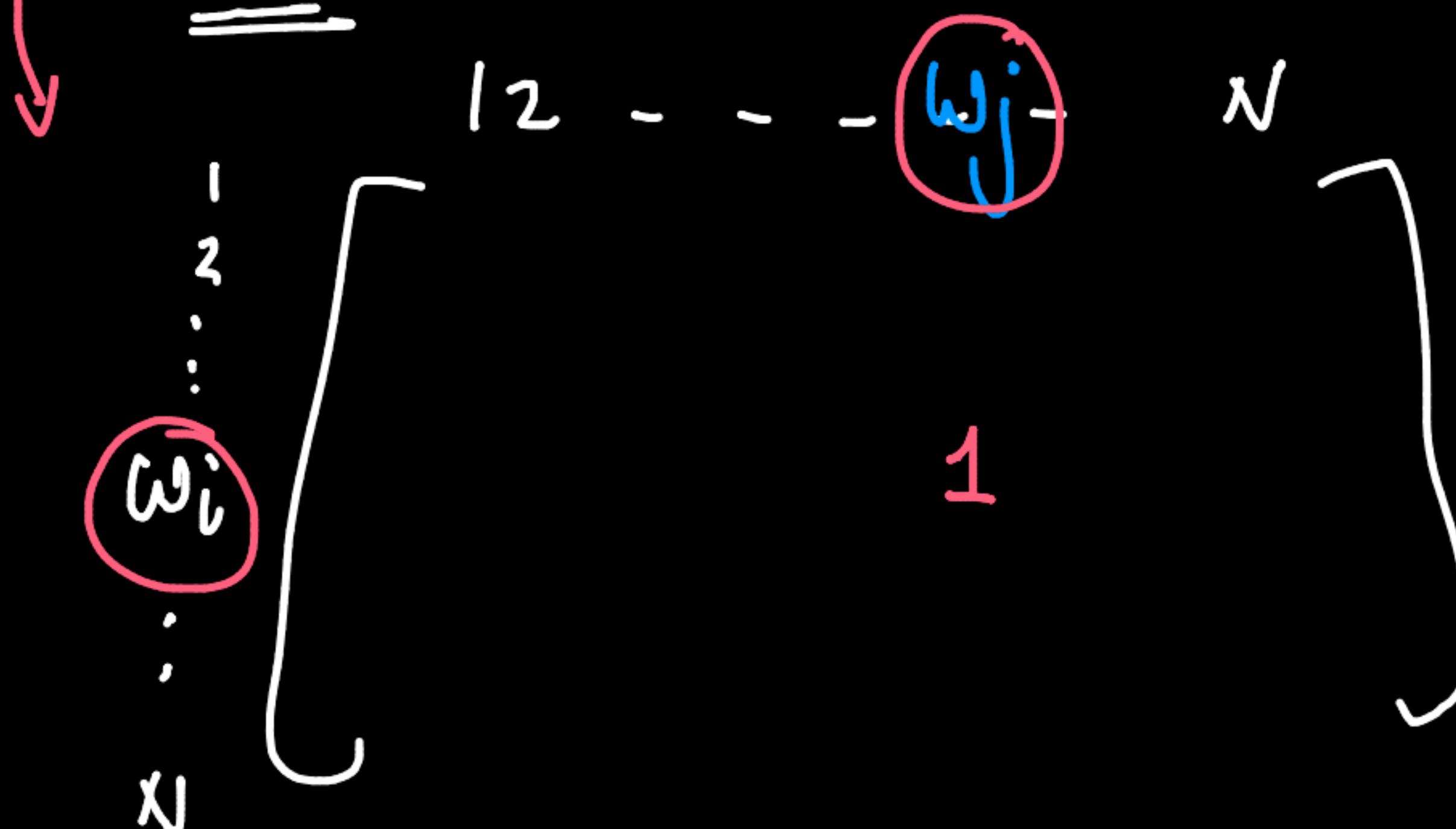


Voc. $\bigcirc N \bigcirc$

→ Problem?
→ very sparse
→ N is large,
→ sequence is long

Data Matrix:

Co-Occurrence - Matrix



window of size 5
(context)

↓
partially taking
Case of the
seq. / context
info :-

$\mathcal{D}W_{n \times N}$

\checkmark
V. Sparse $\longrightarrow M\hat{F}$

$C_{N \times N}$

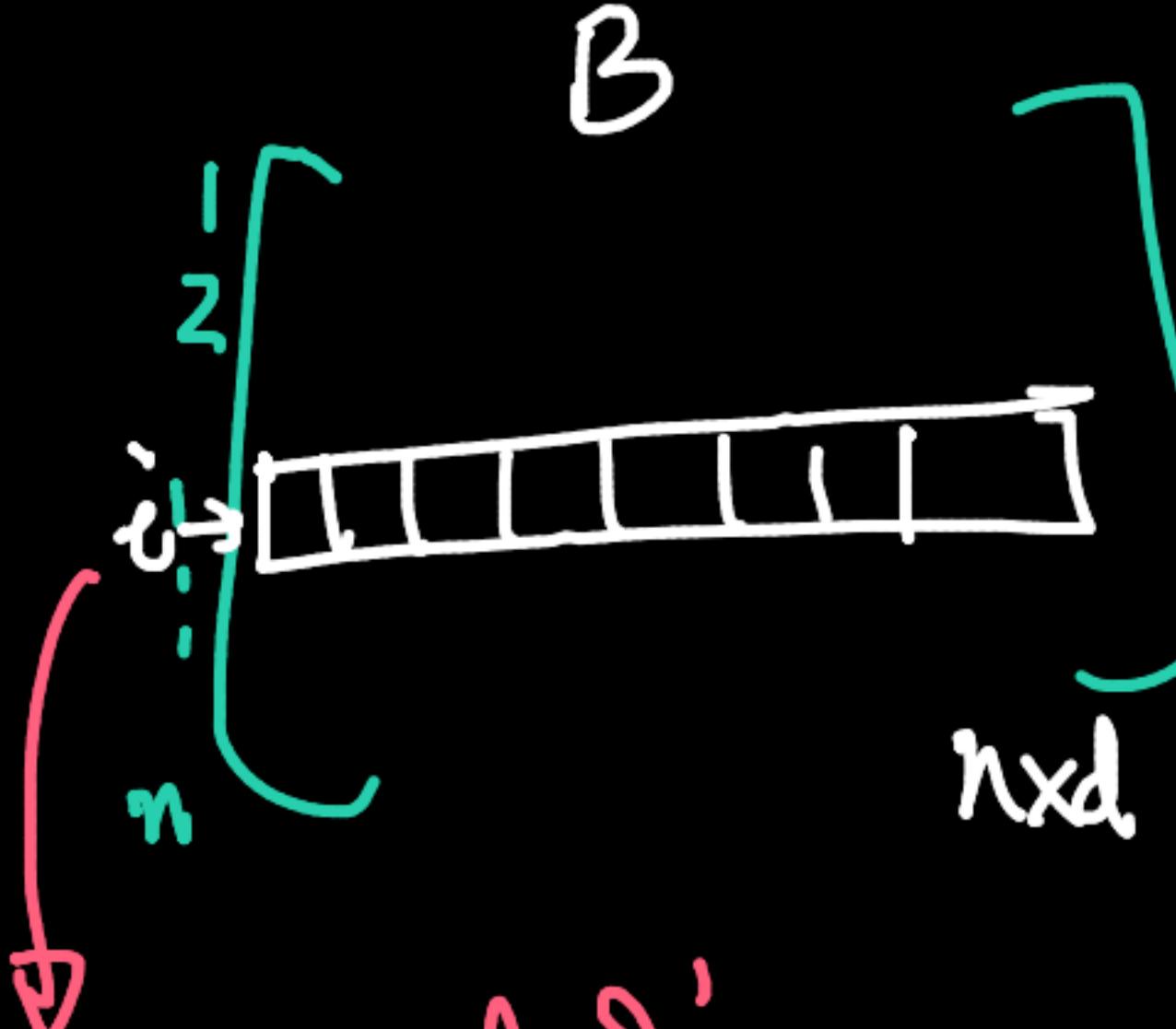
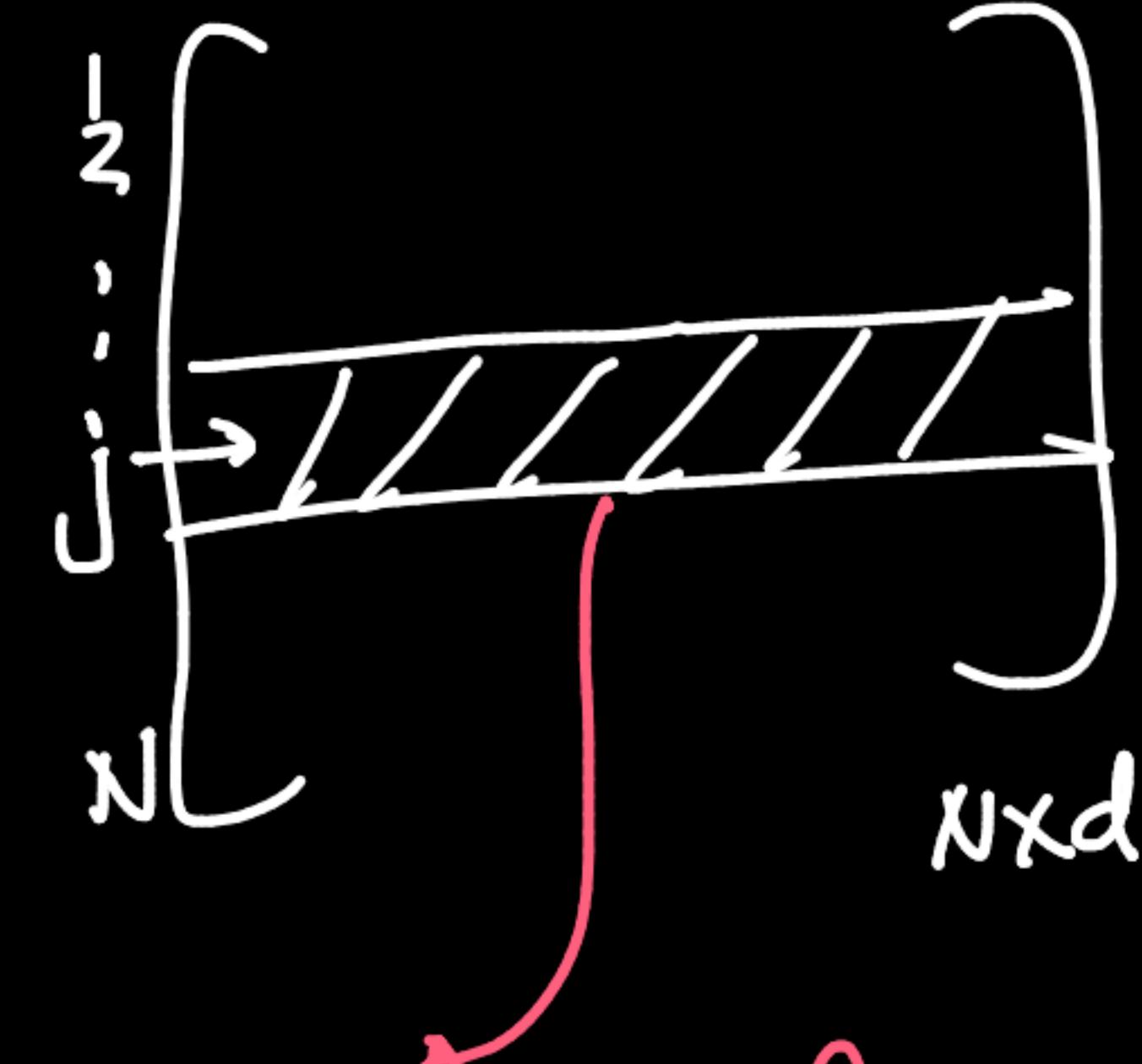


$$\tilde{A}_{n \times N} = B_{n \times d} \cdot C^T_{d \times N}$$

loss:

$$\sum_{i,j} (A_{ij} - B_i^T C_j)^2 + \dots$$

A_{ij} ≠ 0

Vec. rep of \vec{a}_i Vec. rep of \vec{w}_j 

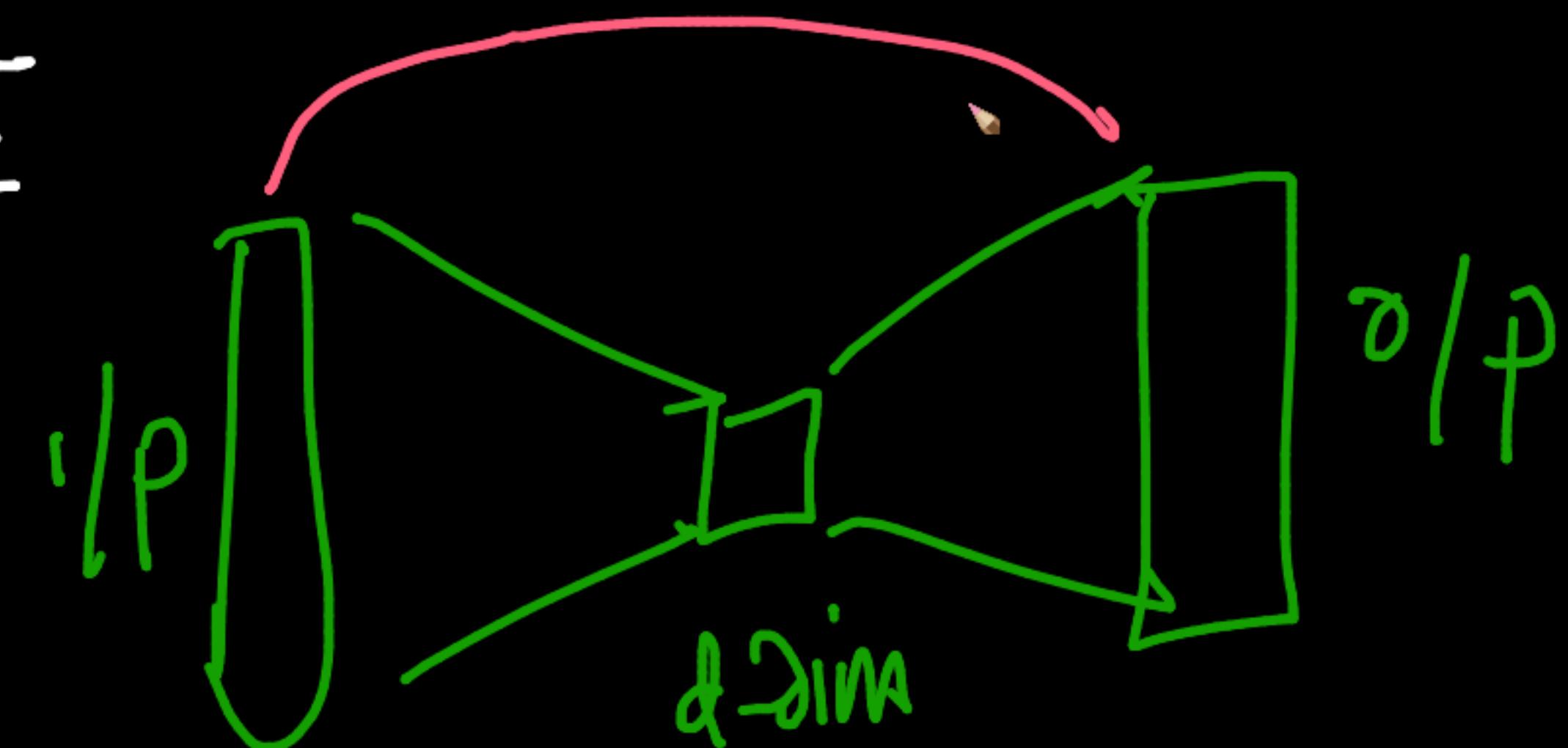
MF /SVD / - - -

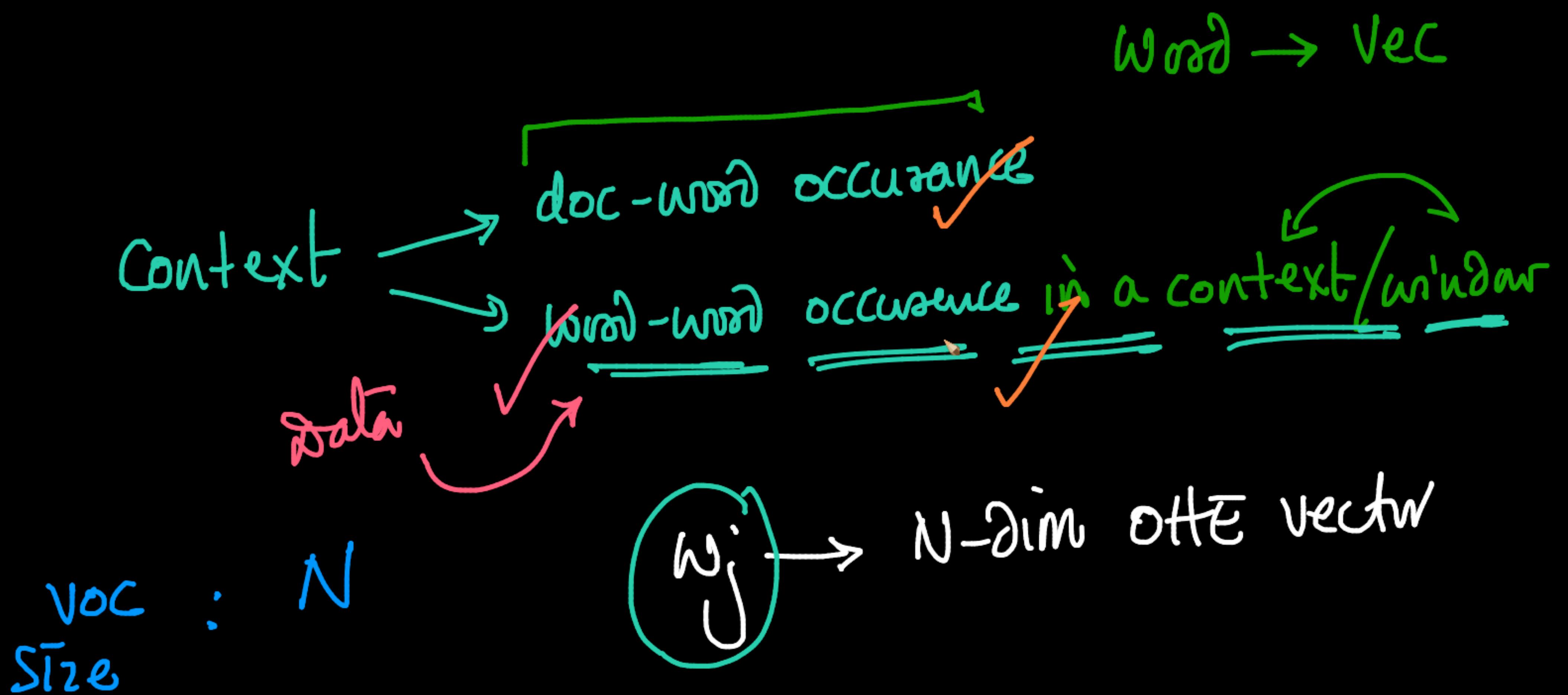


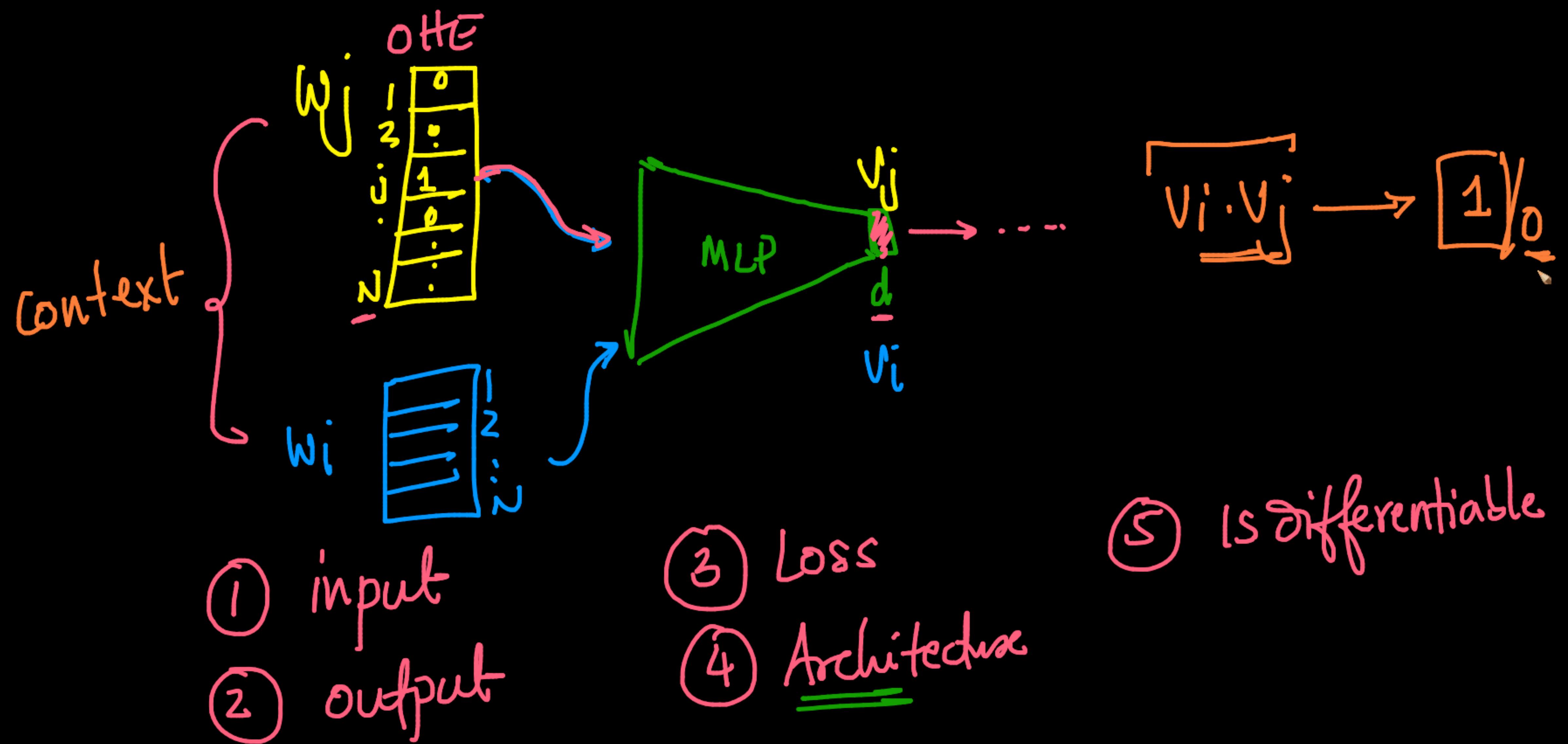
Mf → Classical ML

? → deep-learning (NN)

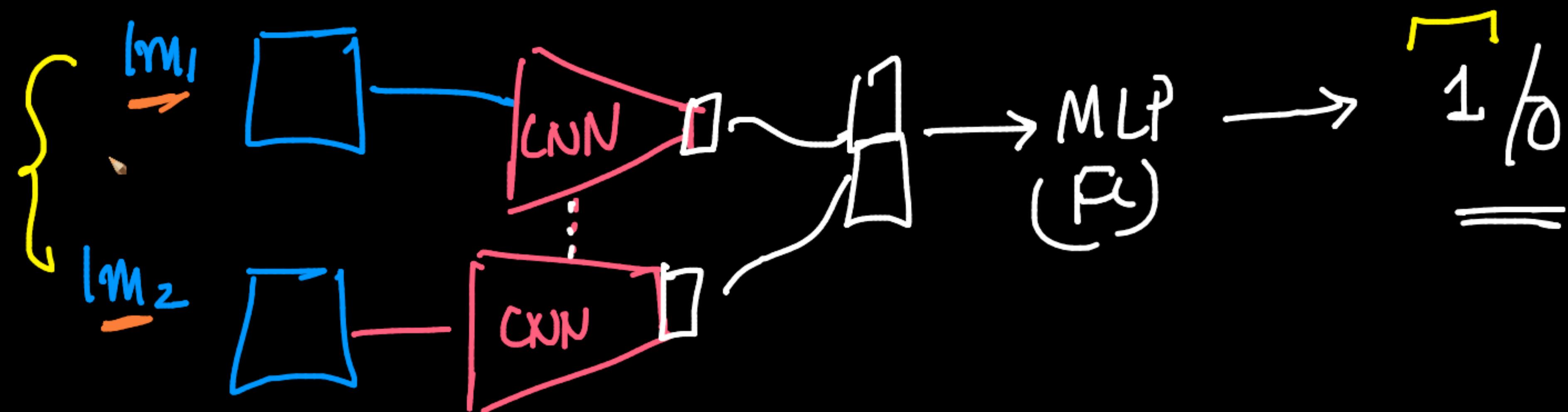
✓ Auto-enc



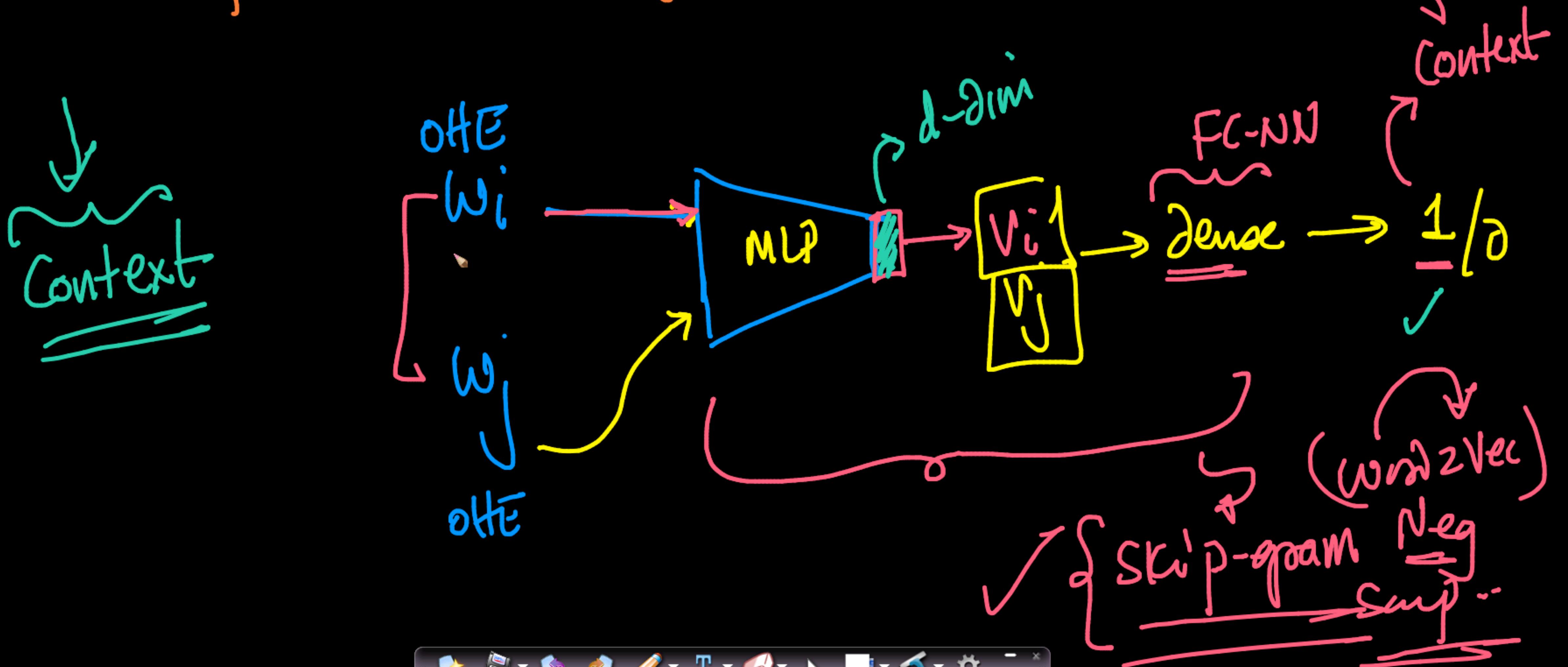




Siamese n/w:



Corpus $\rightarrow w_i \cdot w_j =$ cooccurrent Matrix $(1/d)$



History:

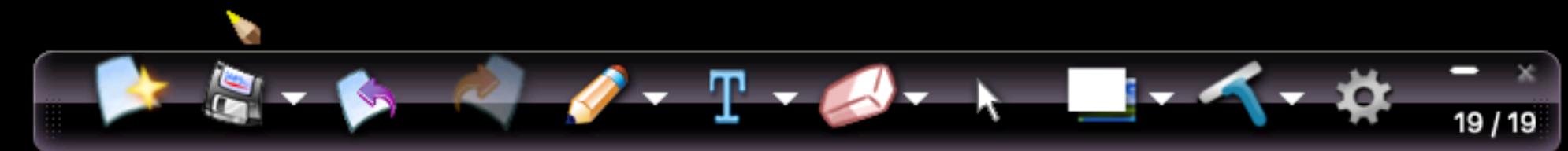
(2012-13)

Word2Vec

→ CBOW ✓

↓
→ skip gram ✓

→ skipgram + Negative sampling ✓



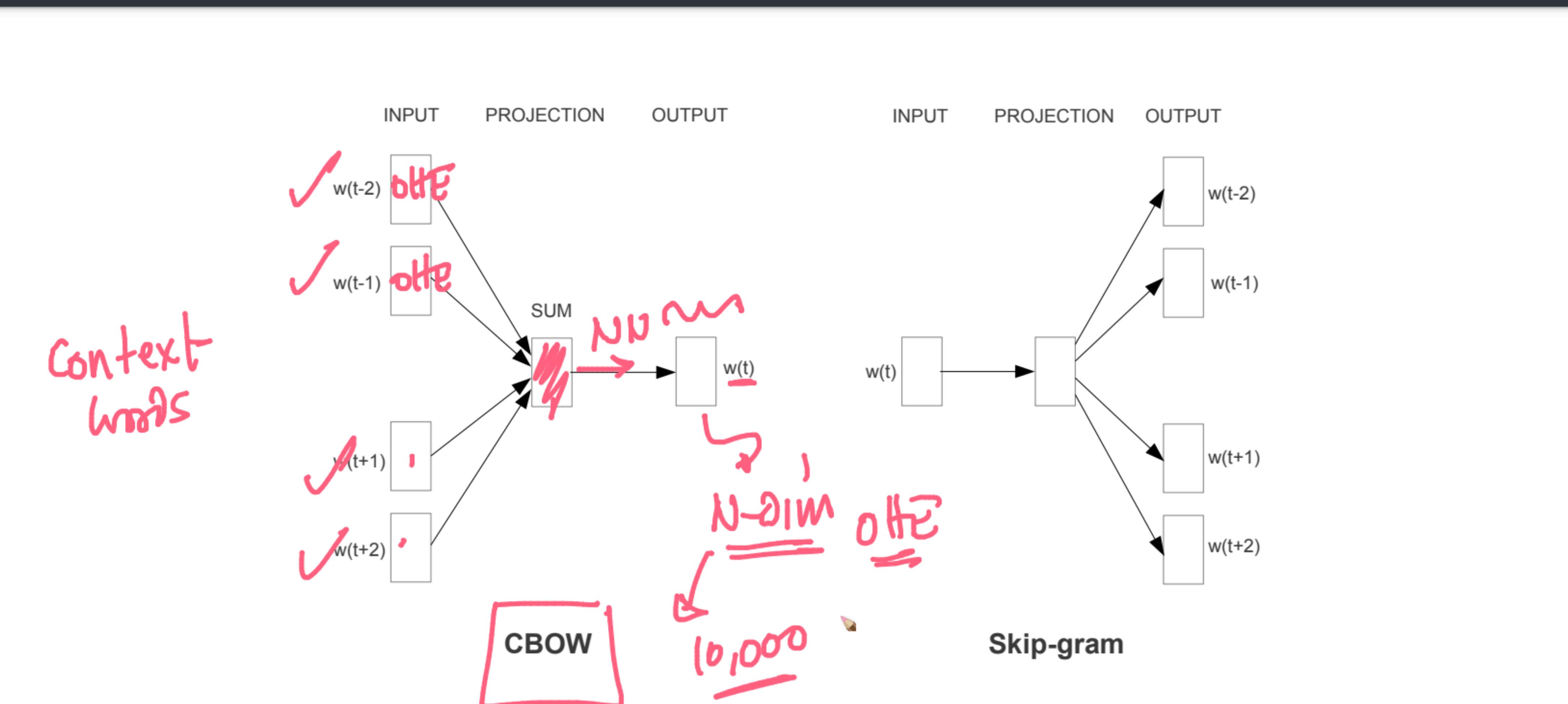


Figure 1: New model architectures. The CBOW architecture predicts the current word based on the context, and the Skip-gram predicts surrounding words given the current word.

WordEmbeddings_(Word2Vec).pdf 1301.3781.pdf word2vec models.word2vec – Word2vec glove.pdf

colab.research.google.com/drive/1fG62NwfWU4PRQ2HpN6kTExL6D-rjVze#scrollTo=eca3715c

+ Code + Text Connect |  

• Finally, the network parameters W_1 and W_2 are adjusted using gradient descent.

Is CBoW used in practice?

- When the context window > 1, CBoW **overfits on frequent words** as the input is an average of one-hot encoded vectors of the context word.
- CBoW does not produce good representations for **rare words** in the corpus.

We can over-come these dis-advantages of CBoW by using an alternative architecture of Word2Vec called the Skip-Gram.

What is Skip-Gram?

- Skip-Gram is simply an **inversion of the CBoW architecture**.
- Takes in **input as center/context word and predict its surrounding words**.

Architecture of Skip Gram



21 / 21

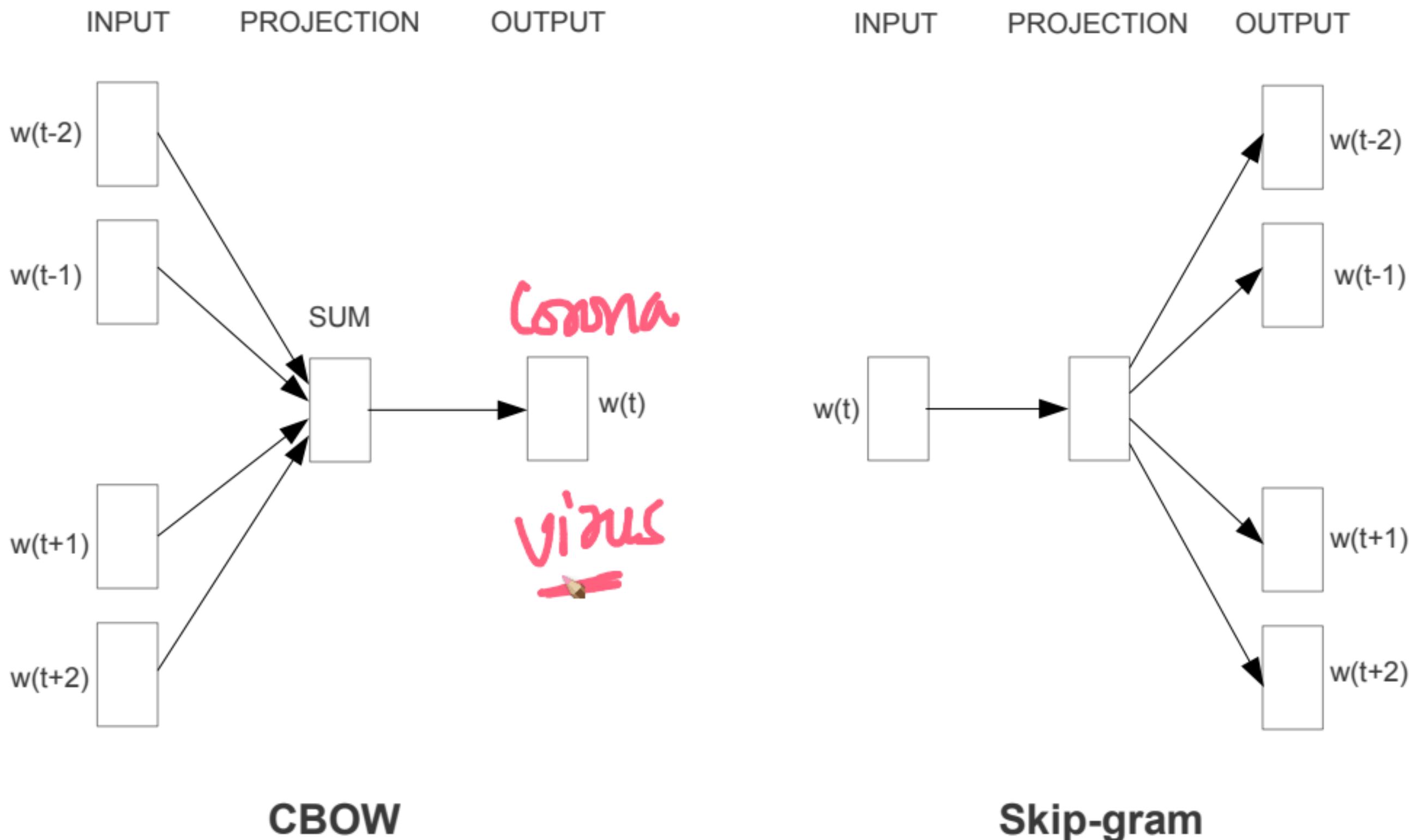
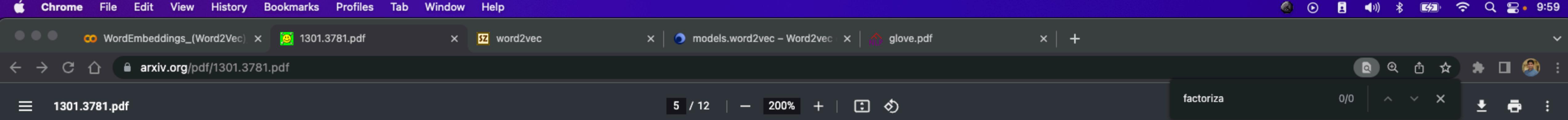


Figure 1: New model architectures. The CBOW architecture predicts the current word based on the context, and the Skip-gram predicts surrounding words given the current word.

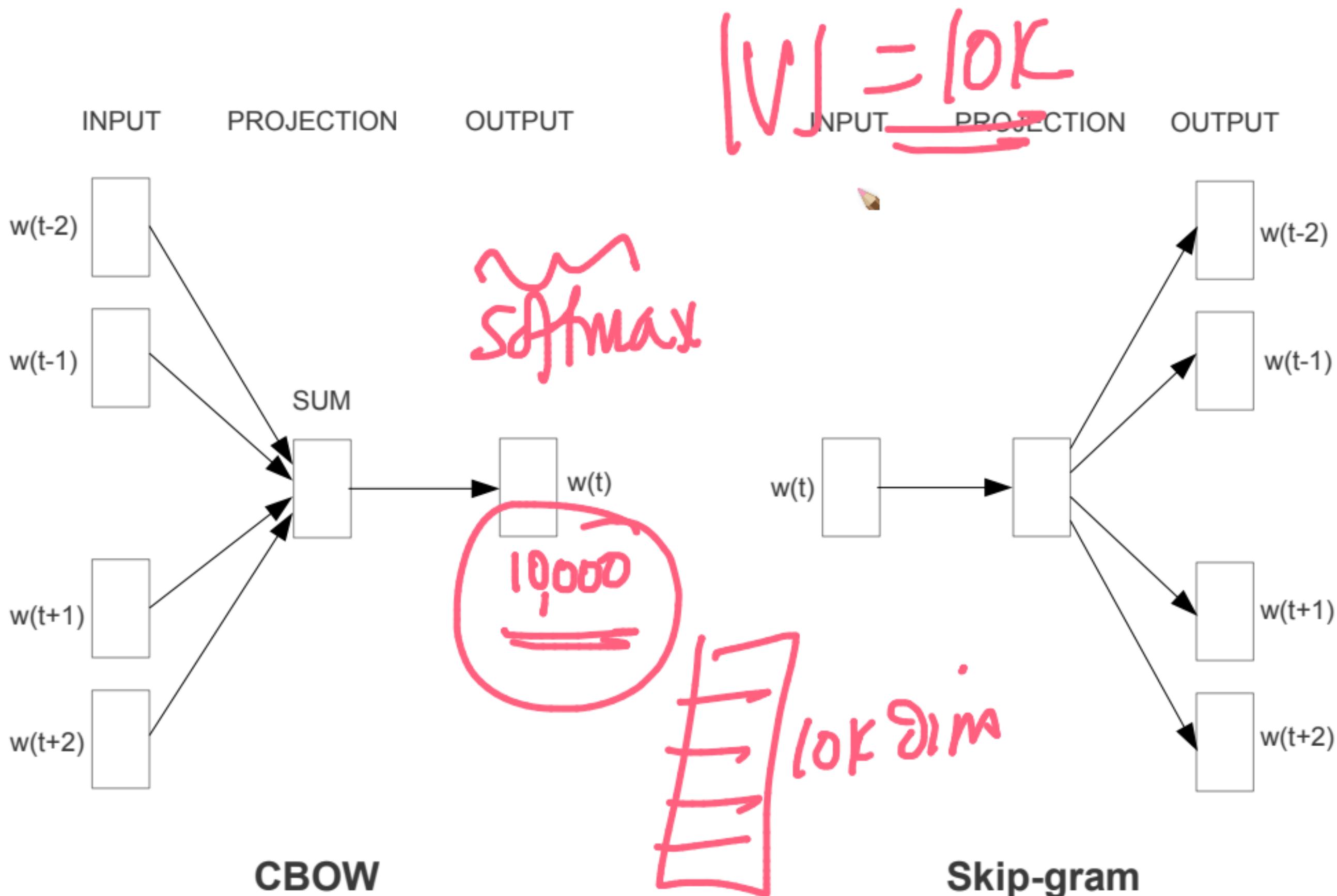
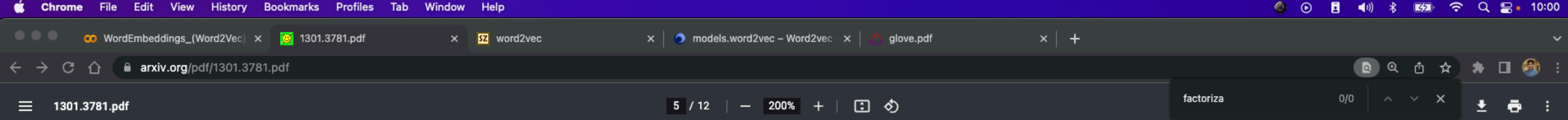
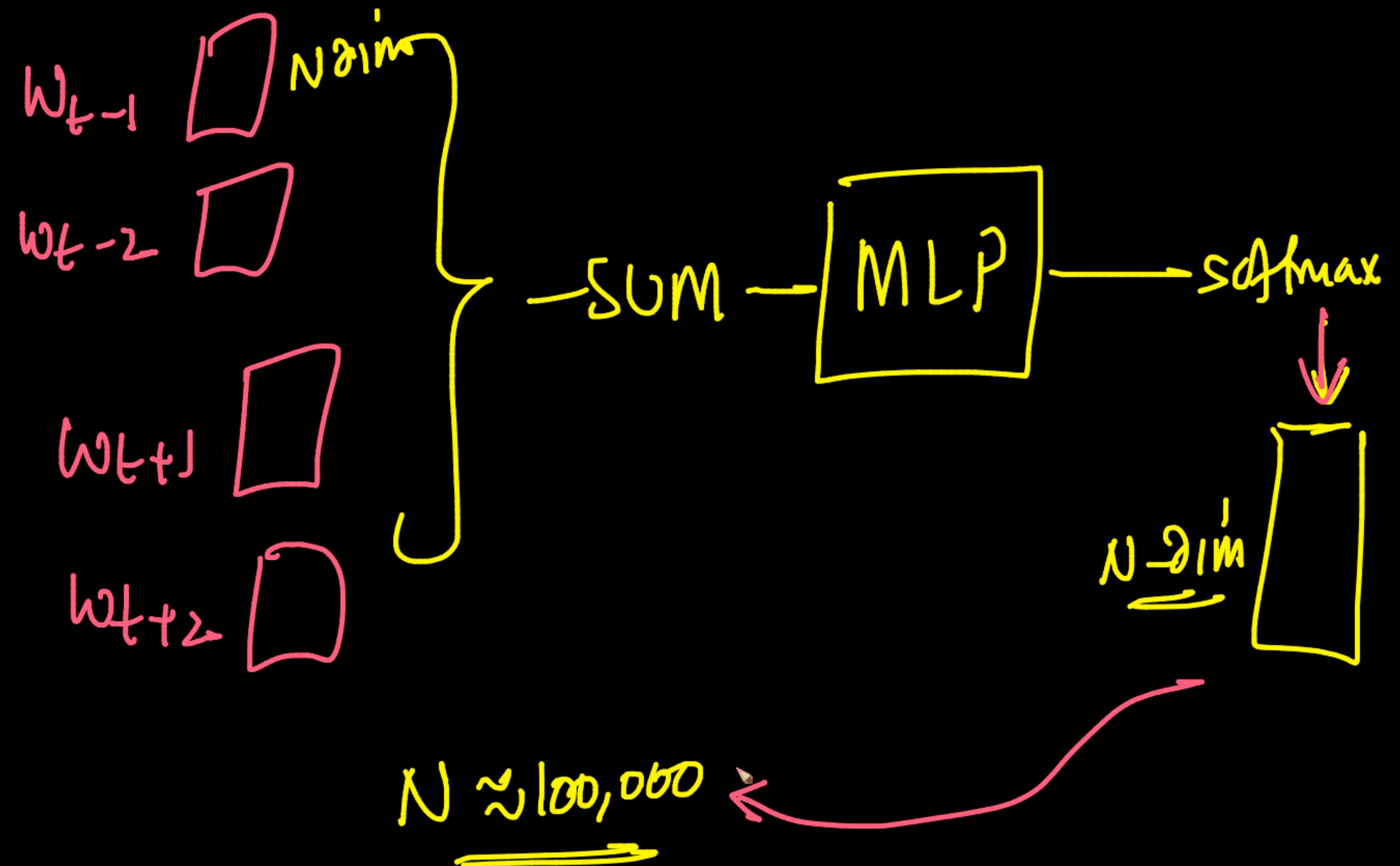
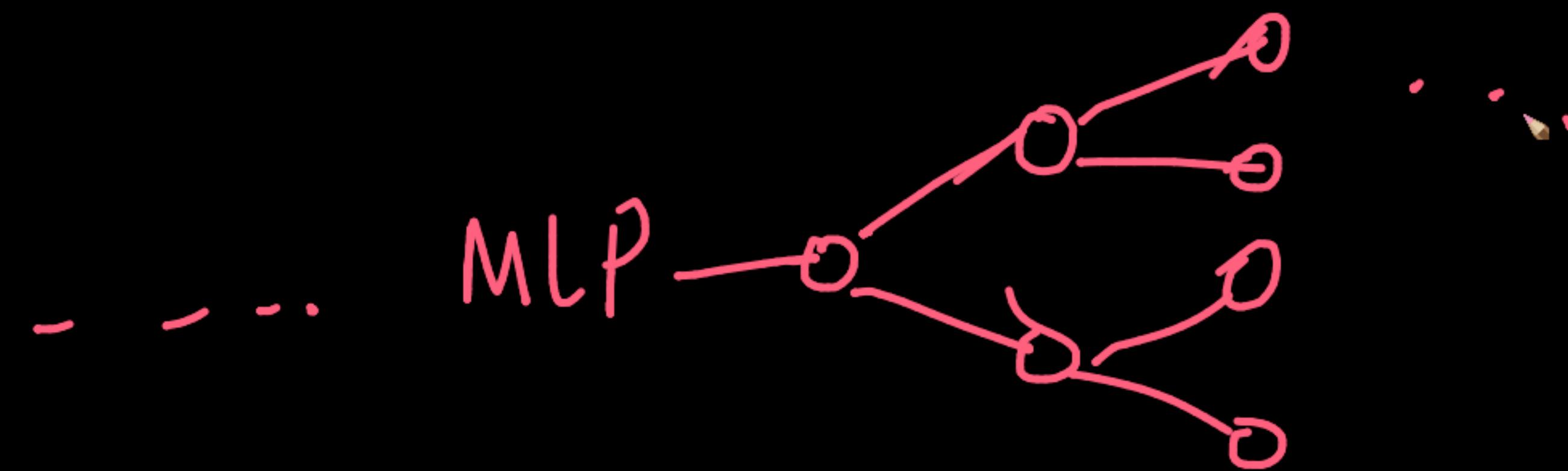


Figure 1: New model architectures. The CBOW architecture predicts the current word based on the context, and the Skip-gram predicts surrounding words given the current word.





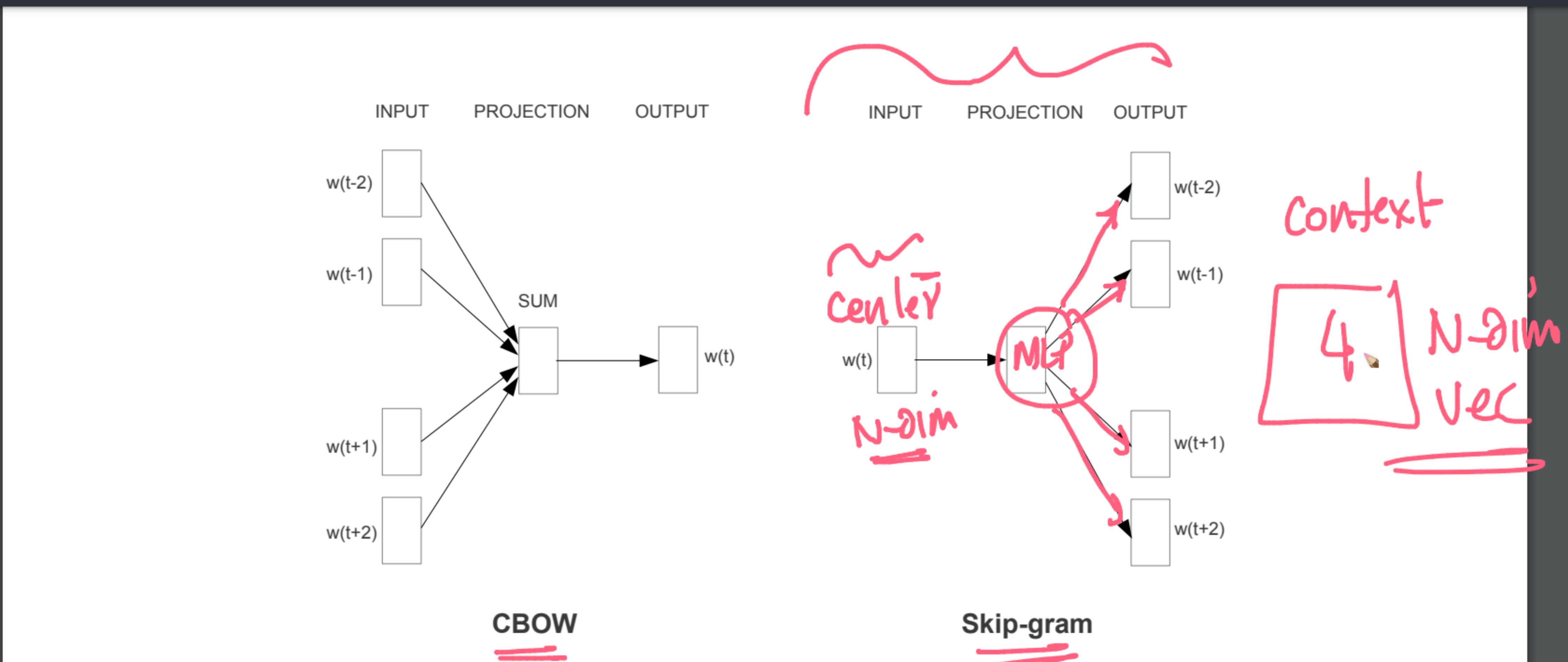
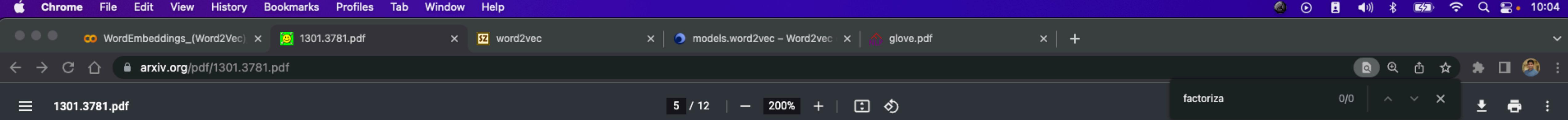


Figure 1: New model architectures. The CBOW architecture predicts the current word based on the context, and the Skip-gram predicts surrounding words given the current word.

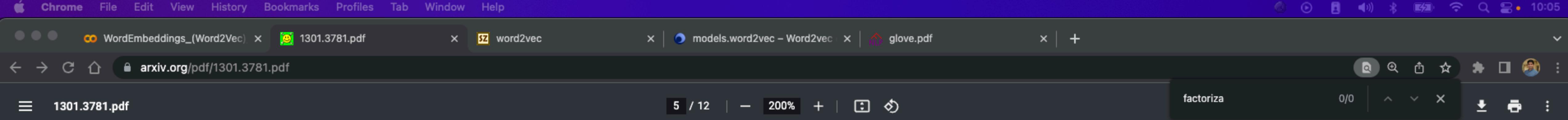


Figure 1: New model architectures. The CBOW architecture predicts the current word based on the context, and the Skip-gram predicts surrounding words given the current word.

R words from the context window has to do $R \times 2$

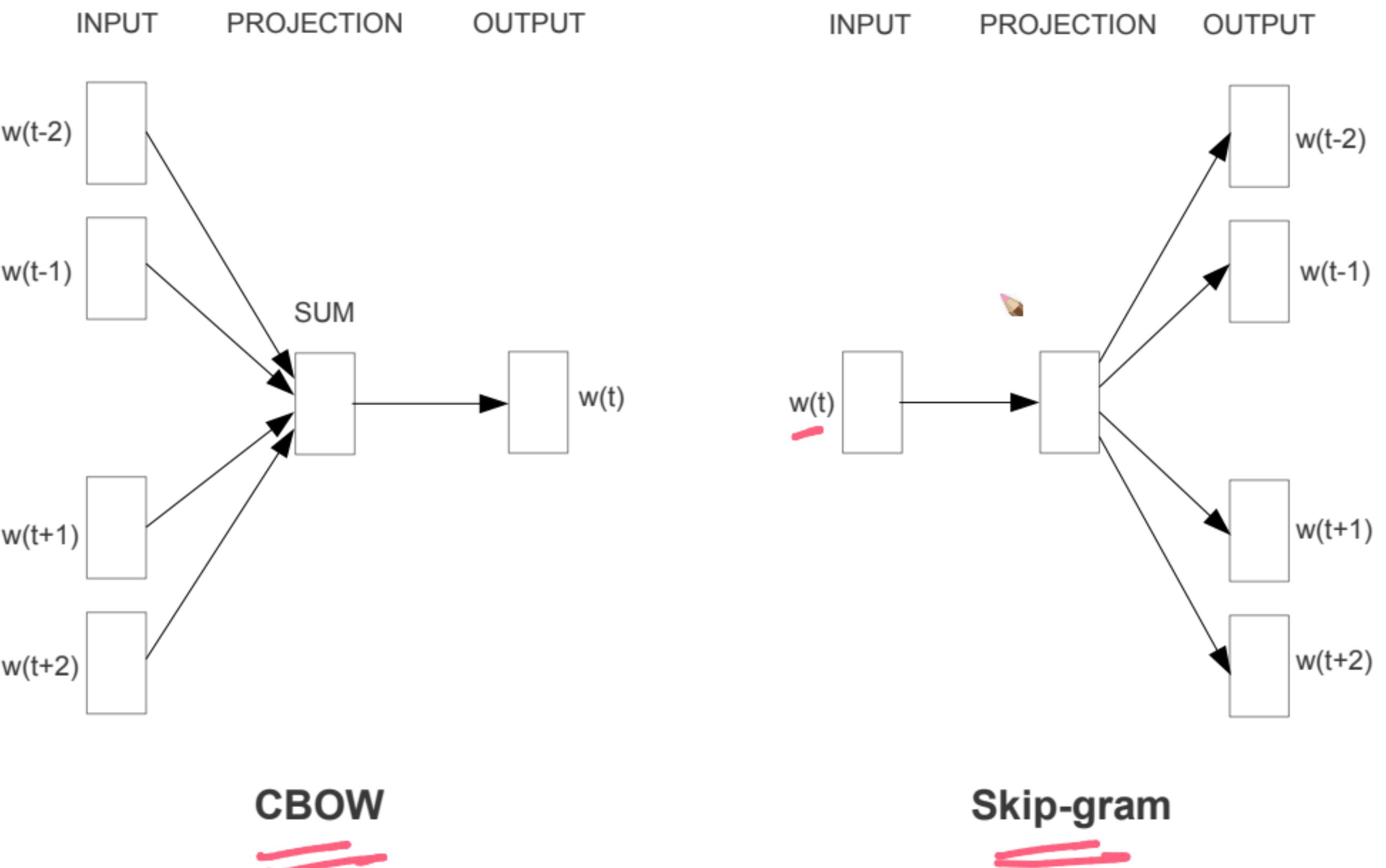
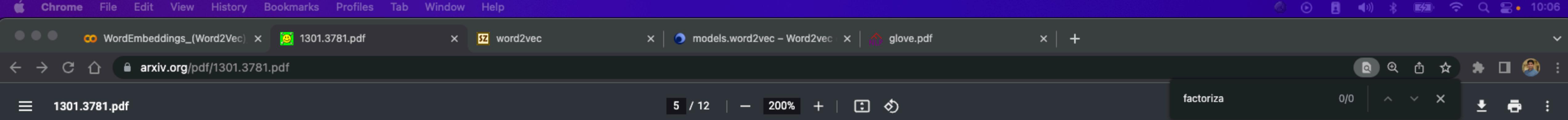


Figure 1: New model architectures. The CBOW architecture predicts the current word based on the context, and the Skip-gram predicts surrounding words given the current word.

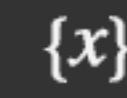


+ Code + Text

Connect

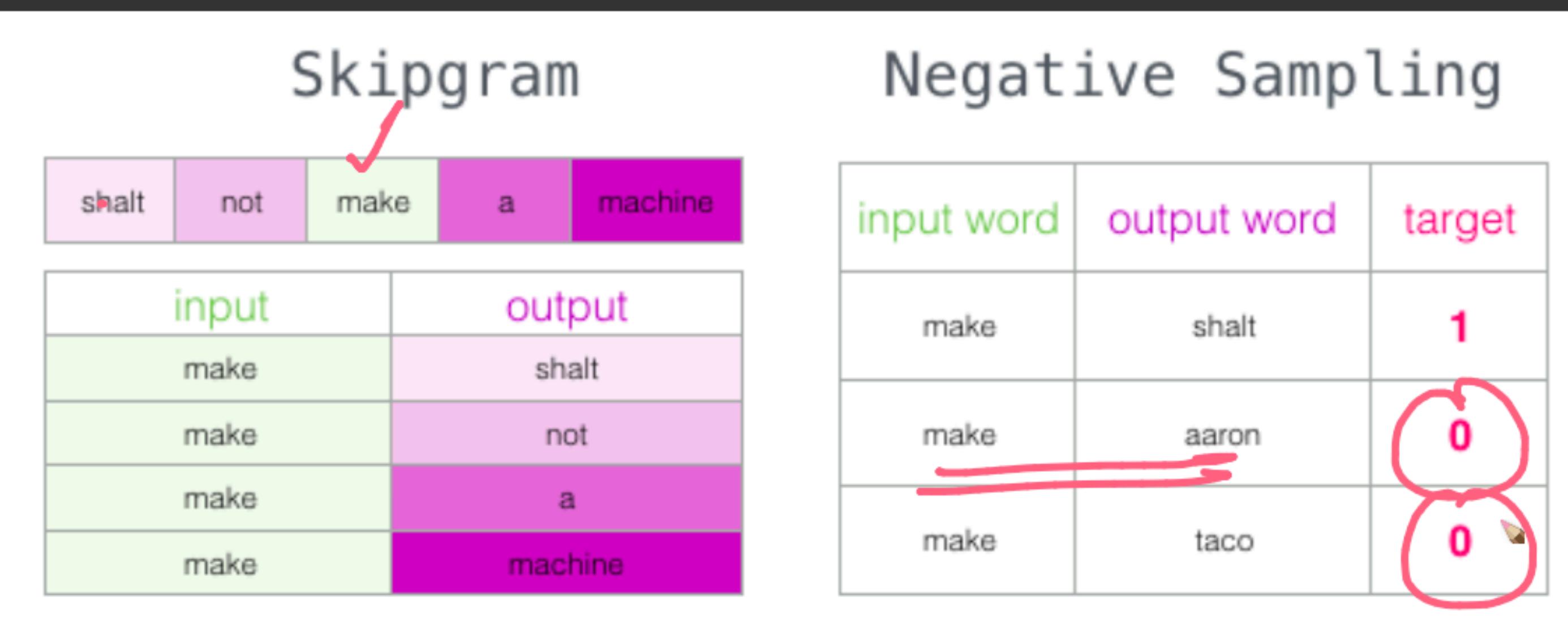


- Given the objective function, the model parameters **W** and **W'** are fine-tuned using gradient descent via backpropagation.
- Similar to CBoW, the **weight matrix W_1** is considered as the **word vector**.



Negative Sampling for Skipgram training:

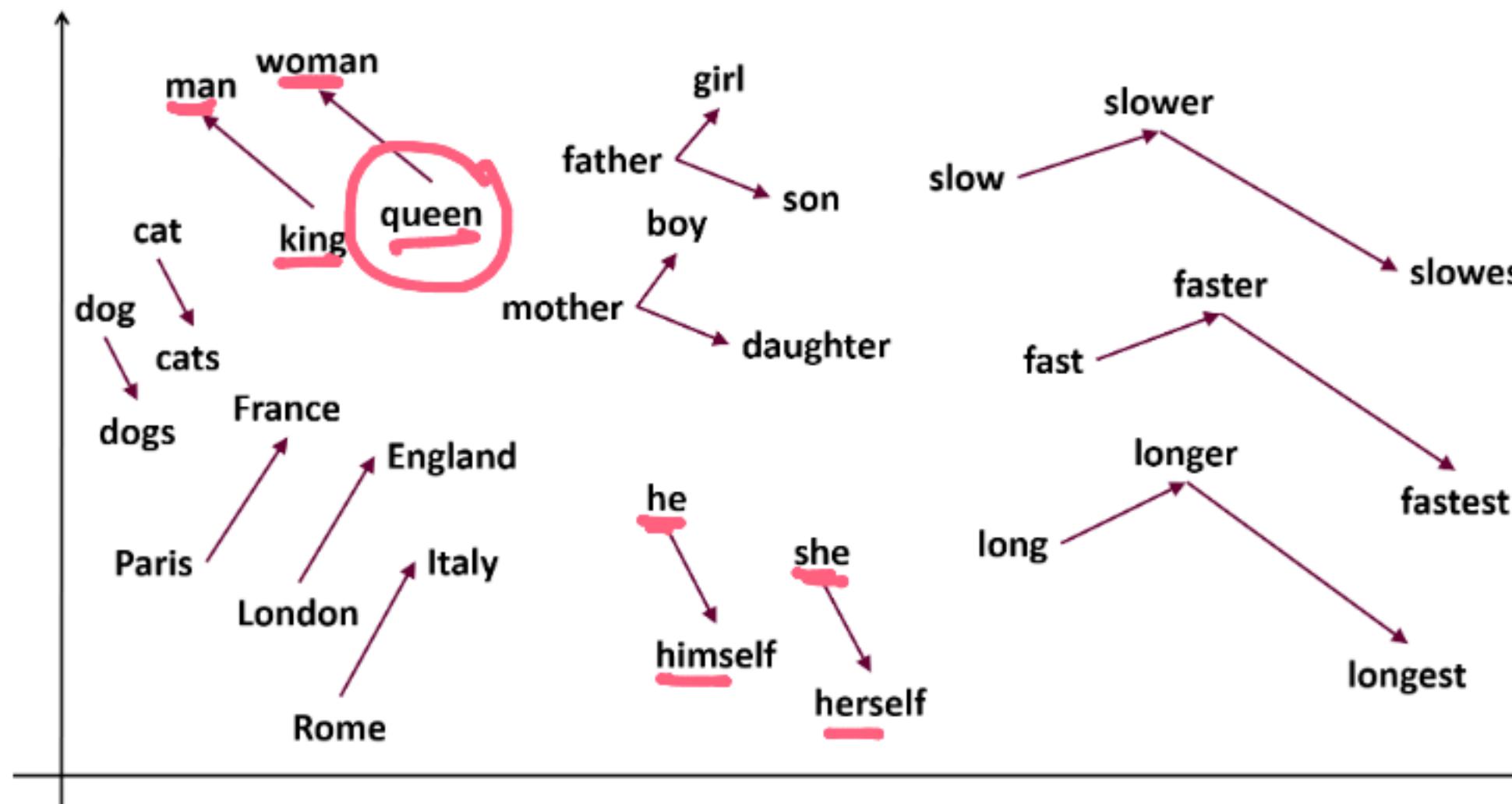
Refer: <https://jalammar.github.io/illustrated-word2vec/>



Now that we've understood how Word2Vec can help generate meaningful and contextual word

Vking - Vman
+ Vwoman

$$= V_{\text{queen}}$$



$w_j \rightarrow \underline{\underline{d-dim}}$

NLP (Natural Language Processing) is a fast developing field of research in recent years, especially by Google, which depends on NLP technologies for managing its vast repositories of text contents.

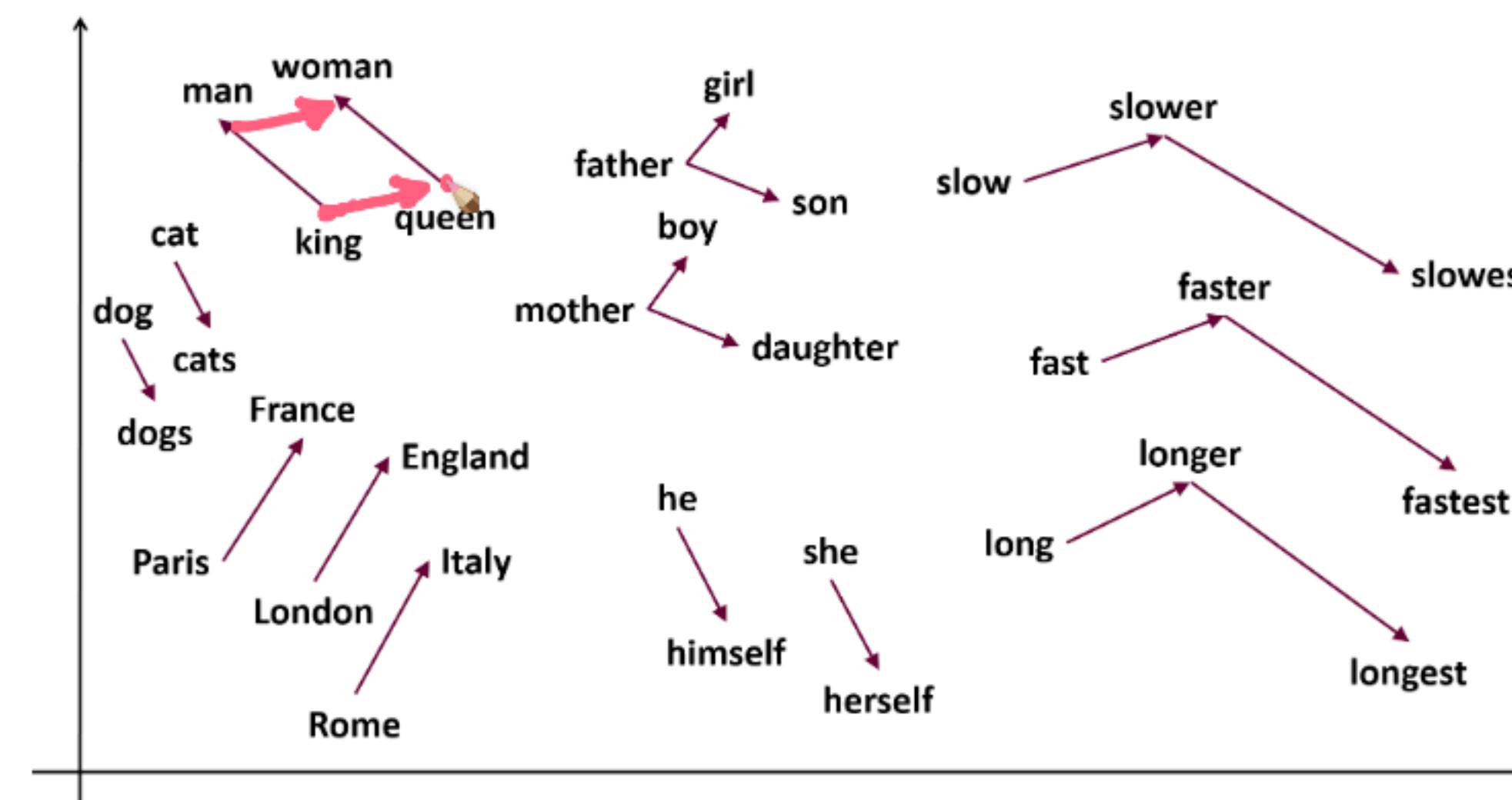
In this study unit we will lay a simple introduction to this field through the use of the excellent **gensim** Python package of [Radim Rehurek's](#) and his excellent **word2vec Tutorial**:
<https://rare-technologies.com/word2vec-tutorial>.

We have followed Radim's code with some supplements and more examples, and adapted it to an **IMDB** movie reviews dataset from Cornell university:

<https://www.cs.cornell.edu/people/pabo/movie-review-data>

You may download this dataset more conveniently from here:
<http://www.samyzaf.com/ML/nlp/aclImdb.zip>

NLP with gensim (word2vec)



NLP (Natural Language Processing) is a fast developing field of research in recent years, especially by Google, which depends on NLP technologies for managing its vast repositories of text contents.

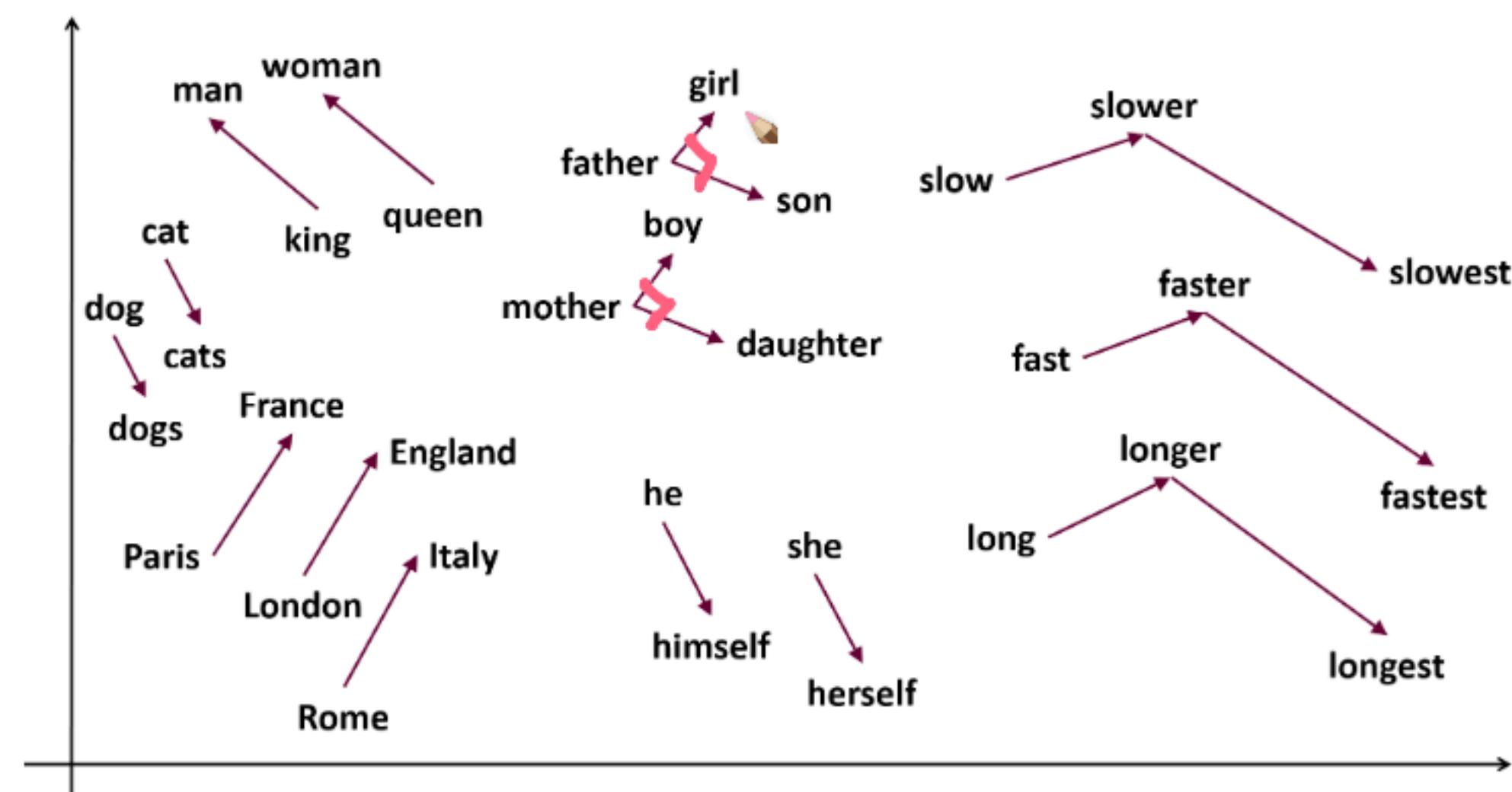
In this study unit we will lay a simple introduction to this field through the use of the excellent **gensim** Python package of [Radim Rehurek's](#) and his excellent **word2vec Tutorial**:
<https://rare-technologies.com/word2vec-tutorial>.

We have followed Radim's code with some supplements and more examples, and adapted it to an **IMDB** movie reviews dataset from Cornell university:

<https://www.cs.cornell.edu/people/pabo/movie-review-data>

You may download this dataset more conveniently from here:
<http://www.samyzaf.com/ML/nlp/aclImdb.zip>

NLP with gensim (word2vec)



NLP (Natural Language Processing) is a fast developing field of research in recent years, especially by Google, which depends on NLP technologies for managing its vast repositories of text contents.

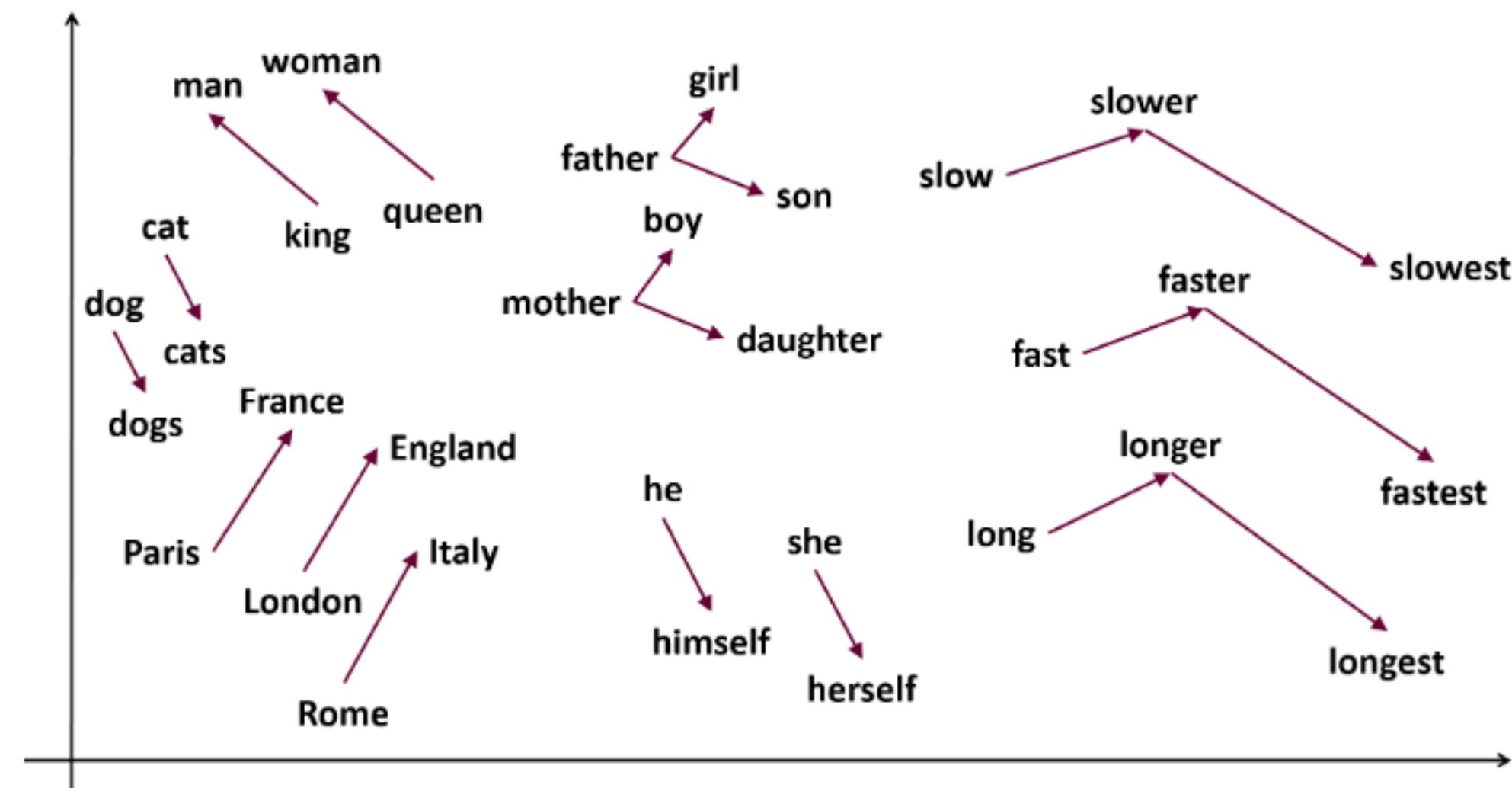
In this study unit we will lay a simple introduction to this field through the use of the excellent **gensim** Python package of [Radim Rehurek's](#) and his excellent **word2vec Tutorial**:
<https://rare-technologies.com/word2vec-tutorial>.

We have followed Radim's code with some supplements and more examples, and adapted it to an **IMDB** movie reviews dataset from Cornell university:

<https://www.cs.cornell.edu/people/pabo/movie-review-data>

You may download this dataset more conveniently from here:
<http://www.samyzaf.com/ML/nlp/aclImdb.zip>

NLP with gensim (word2vec)



2013 - ...

2018 -

NLP (Natural Language Processing) is a fast developing field of research in recent years, especially by Google, which depends on NLP technologies for managing its vast repositories of text contents.

In this study unit we will lay a simple introduction to this field through the use of the excellent **gensim** Python package of [Radim Rehurek's](#) and his excellent **word2vec Tutorial**:
<https://rare-technologies.com/word2vec-tutorial>.

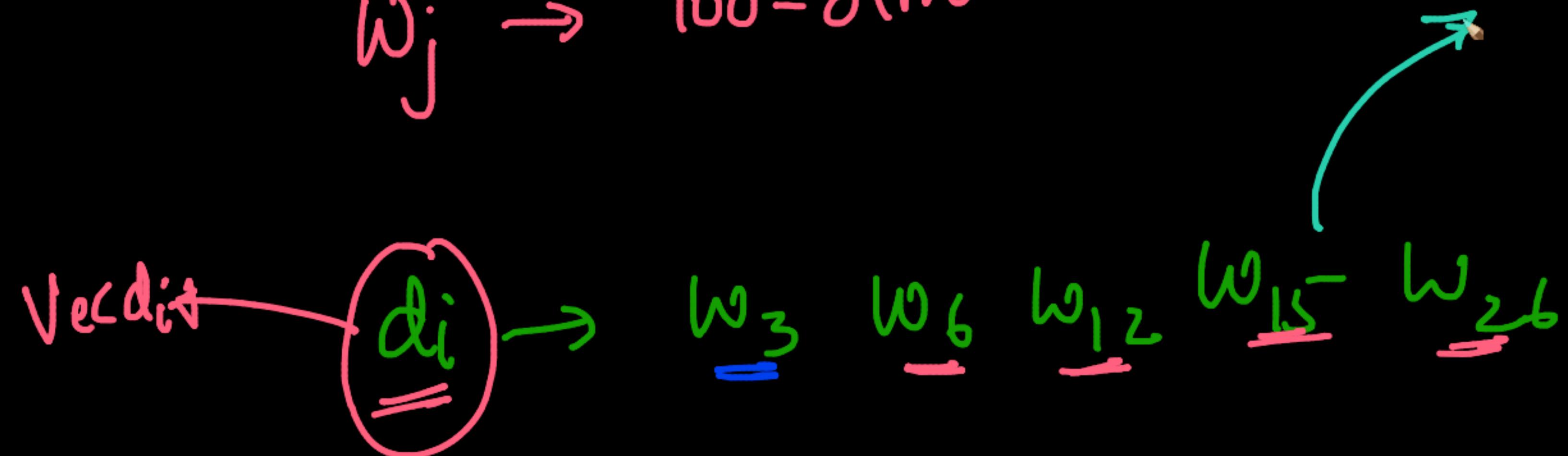
We have followed Radim's code with some supplements and more examples, and adapted it to an **IMDB** movie reviews dataset from Cornell university:

<https://www.cs.cornell.edu/people/pabo/movie-review-data>

You may download this dataset more conveniently from here:

<http://www.samyzaf.com/ML/nlp/aclImdb.zip>

(Q)

 $w_j \rightarrow (00\text{-dim})$ 

Task: Query: $w_6 w_4 w_5$ → Vec_Q

↓
docs Most similar

DOC-VEC

① avg

✓ ② centroid③ weighted avg \rightarrow TFIDF

Renewer

$$\{ d_1: w_3 \quad w_6 \quad w_4 = df \\ \downarrow \quad \downarrow \quad \downarrow \\ \text{TFIDF}_3 + \text{TFIDF}_6 + \quad v_4 \cdot \bar{f}(\partial f_4) \}$$

WordEmbeddings_(Word2Vec).pdf 1301.3781.pdf word2vec models.word2vec – Word2vec Embedding projector - visualization glove.pdf

colab.research.google.com/drive/1fG62NwfWU4PRQ2HpN6kTExL6D-rjVze#scrollTo=be5e94a6

+ Code + Text

```
for item in cosine_list[:num]:  
    [ ]     papers_list.append((item[0], item[1], item[2]))  
return papers_list
```

{x}

▼ Method to use the base function and retrieve top matching documents

□

```
[ ] def query(query, top_matches=10):  
    model_to_use = model  
    df_covid_to_use = df_covid  
    return rank_docs(model_to_use, query, df_covid_to_use, top_matches)  
  
[ ] query('origin of corona virus')
```

```
[ ('Role of Nonstructural Proteins in the Pathogenesis of SARS-CoV-2',  
  '10.1002/jmv.25858',  
  2.381886614774109),  
  ('Re-emergence of a genetic outlier strain of equine arteritis virus: Impact on phylogeny',  
   '10.1016/j.virusres.2014.12.009',  
   2.3744738115509545),  
  ('Emerging and Neglected Viruses of Zoonotic Importance in Croatia',  
   '10.3390/pathogens10010073',  
   2.370784610584278),  
  ('The global emergence of severe acute respiratory syndrome coronavirus 2 in human',  
   '10.1007/s13337-020-00613-y',  
   2.369730359573394),
```

<>

≡

>-

+ Code + Text

```
banyangvirus in Northern bats from Germany',  
[ ] '10.1038/s41598-020-58466-w',  
2.357634998523873),  
('Update of the current knowledge on genetics, evolution, immunopathogenesis, and transmission for coronavirus  
disease 19 (COVID-19)',  
'10.7150/ijbs.48812',  
2.3548538950775884)]
```

Conclusion

- Continuous text representations are used to capture syntactic and semantic similarity that, discrete text representations don't capture.
- SVD based methods suffer from scaling and in practice, iterative method like Word2Vec is used.
- Word2Vec models the data as multi-class classification problem and in doing so, learns the word representations.
- Two architectures of Word2Vec - (i) CBoW and (ii) Skip-Gram.
- Training methods available - (i) Negative Sampling and (ii) Hierarchical Softmax
- Gensim library is used for building Word2Vec from scratch.
- Finally, we built a search engine based on semantic similarity between query and abstracts of the publications.

Post Read

+ Code + Text

Connect |  

```
banyangvirus in Northern bats from Germany',  
[ ] '10.1038/s41598-020-58466-w',  
2.357634998523873),  
('Update of the current knowledge on genetics, evolution, immunopathogenesis, and transmission for coronavirus  
disease 19 (COVID-19)',  
'10.7150/ijbs.48812',  
2.3548538950775884)]
```

Conclusion

- **Continuous text representations** are used to **capture syntactic and semantic similarity** that, **discrete text representations** don't capture.
- **SVD** based methods **suffer from scaling** and in practice, iterative method like **Word2Vec** is used.
- **Word2Vec** models the data as multi-class classification problem and in doing so, **learns the word representations**.
- Two architectures of Word2Vec - (i) **CBoW** and (ii) **Skip-Gram**.
- Training methods available - (i) **Negative Sampling** and (ii) **Hierarchical Softmax**
- ✓ **Gensim** library is used for building Word2Vec from scratch.
- Finally, we built a **search engine** based on **semantic similarity between query and abstracts** of the publications.

Post Read

✓ Word2Vec (Google) (2013) → Very small-context / window
= = =

→ ✓ Glove (Stanford) (2014) → classical MF (doc-wm) + W2Vec
= = = AnxN

→ fasttext (FB) (2016-17)
= = =
Non-vocabulary words
facebook

With over 50,000 publications and research papers on the corona-virus family till date, it has become difficult to search across and get useful insights for medical practitioners.

As a *DataScientist* at Google, you are tasked with solving this problem with the help of Machine Learning.

Efficient Estimation of Word Representations in Vector Space

GloVe: Global Vectors for Word Representation

How can we solve for this problem?

- Can we match keywords from user queries that are present in abstract?
 - If we do keyword matching, will we be able to understand the user's intent? For e.g: 'origin' and 'discovery'
 - Should we consider the context of the words, then?

Let us build a search engine using Word Embeddings

Dataset