

# Deep Learning for Natural Language Processing (NLP)

# Representation of words as vectors

# One-Hot Encoding

write all words → give them numbers.

Document ID	Text
D1	Dog bites man.
D2	Man bites dog.
D3	Dog eats meat.
D4	Man eats food.

Our toy corpus having six unique words: dog = 1, bites = 2, man = 3, meat = 4, food = 5, eats = 6.

# One-Hot Encoding

Document ID	Text
D1	Dog bites man.
D2	Man bites dog.
D3	Dog eats meat.
D4	Man eats food.

Our toy corpus having six unique words: dog = 1, bites = 2, man = 3, meat = 4 , food = 5, eats = 6.

D3 is represented as

[ [1 0 0 0 0 0], [0 0 0 0 0 1], [0 0 0 1 0 0]].

# Bag of Words (BoW) → Ek me hi like do.

Document ID	Text
D1	Dog bites man.
D2	Man bites dog.
D3	Dog eats meat.
D4	Man eats food.

Our toy corpus having six unique words: dog = 1, bites = 2, man = 3, meat = 4, food = 5, eats = 6.

# Bag of Words (BoW)

Document ID	Text
D1	Dog bites man.
D2	Man bites dog.
D3	Dog eats meat.
D4	Man eats food.

Our toy corpus having six unique words: dog = 1, bites = 2, man = 3, meat = 4, food = 5, eats = 6.

D3 is represented as [1 0 0 1 0 1]

2 - Grams

## Bigram (Bag of 2-Grams)

Document ID	Text
D1	<u>Dog</u> <u>bites</u> man.
D2	<u>Man</u> <u>bites</u> dog.
D3	<u>Dog</u> <u>eats</u> meat.
D4	<u>Man</u> <u>eats</u> food.

{dog bites,  
bites man,  
man bites,  
bites dog,  
dog eats,  
eats meat,  
man eats,  
eats food}.

- We need eight dimensions to represent a word in document.

# Bigram (Bag of 2-Grams)

Document ID	Text
D1	Dog bites man.
D2	Man bites dog.
D3	Dog eats meat.
D4	Man eats food.

{dog bites,  
bites man,  
man bites,  
bites dog,  
dog eats,  
eats meat,  
man eats,  
eats food}.

- We need eight dimensions to represent a word in document.
- Bigram representation for D1 is as follows: D1 : [1,1,0,0,0,0,0,0]



# Bigram (Bag of 2-Grams)

Document ID	Text
D1	Dog bites man.
D2	Man bites dog.
D3	Dog eats meat.
D4	Man eats food.

{dog bites,  
bites man,  
man bites,  
bites dog,  
dog eats,  
eats meat,  
man eats,  
eats food}.

- We need eight dimensions to represent a word in document.
- Bigram representation for D1 is as follows: D1 : [1,1,0,0,0,0,0,0]

Bag of N-Grams (BoN): Break text into chunks of n contiguous words (or tokens).

## Bigram (Bag of 2-Grams)

Document ID	Text
D1	Dog bites man.
D2	Man bites dog.
D3	Dog eats meat.
D4	Man eats food.

{dog bites,  
bites man,  
man bites,  
bites dog,  
dog eats,  
eats meat,  
man eats,  
eats food}.

- We need eight dimensions to represent a word in document.
- Bigram representation for D1 is as follows: D1 : [1,1,0,0,0,0,0,0]



**Limitation: BoW or BoN representations does not have semantic meaning.**

**Bag of N-Grams (BoN): Break text into chunks of n contiguous words (or tokens).**

# Word Embedding

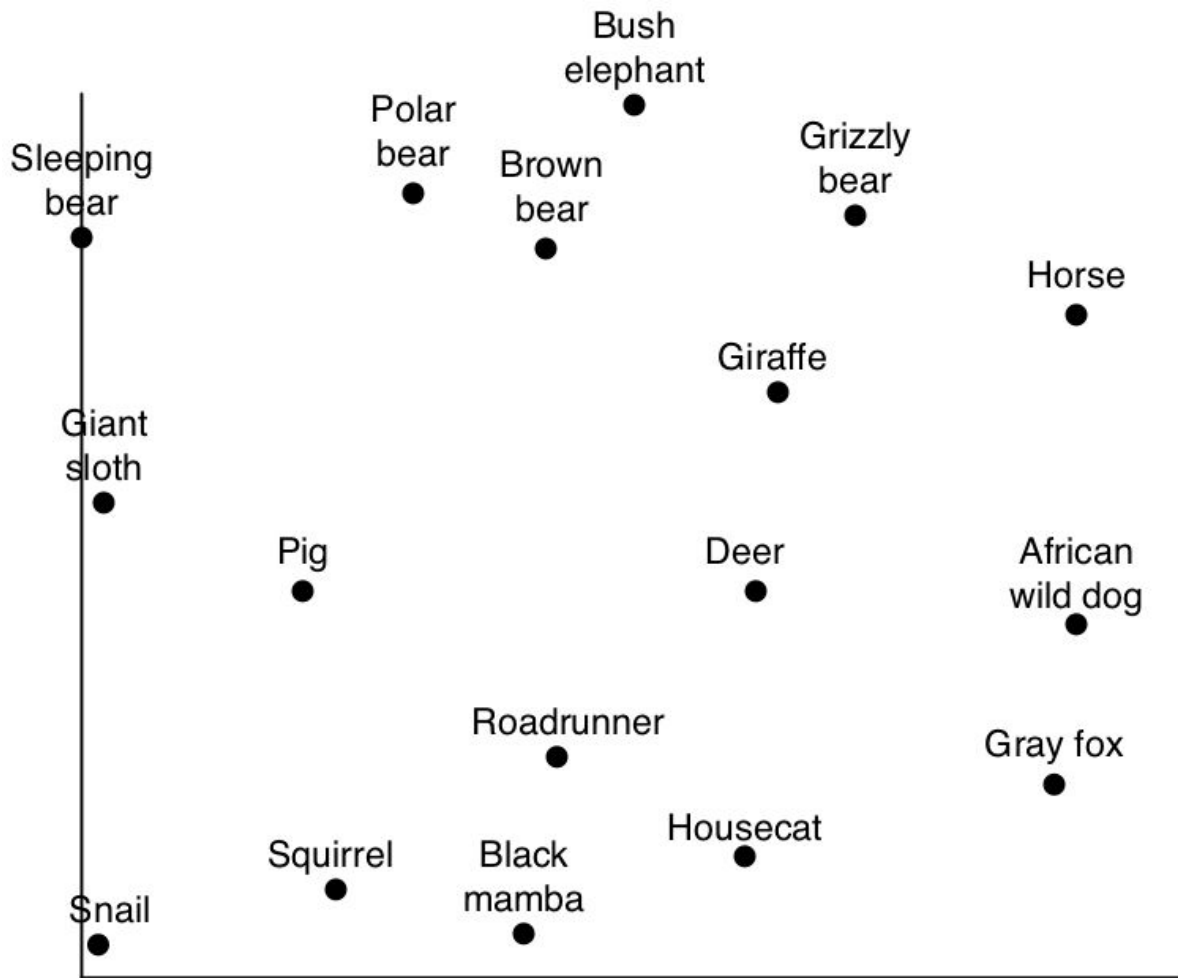
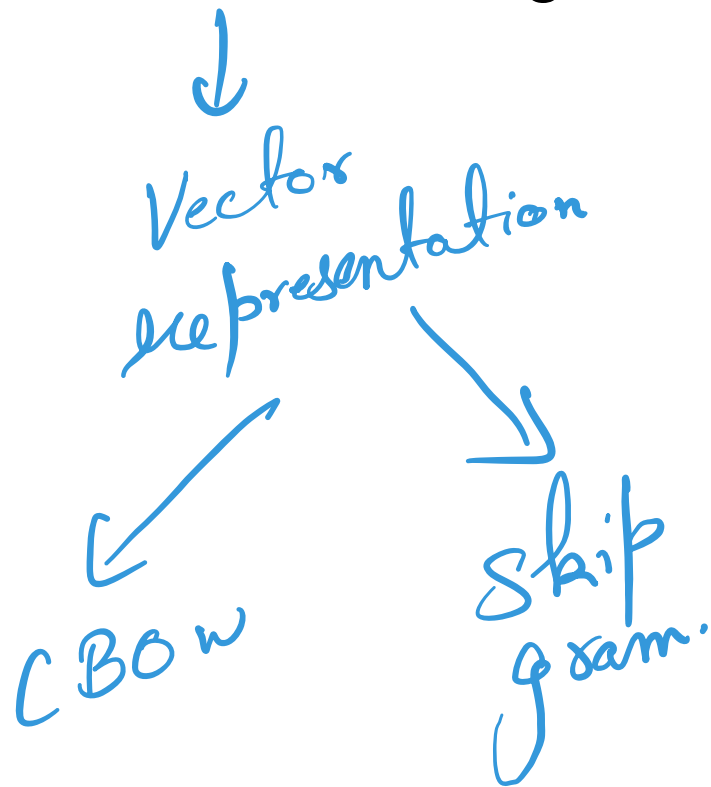


Figure 20-1: A collection of animals, organized roughly by land speed horizontally and adult weight vertically, though those axis labels aren't shown (data from Reisner 2020)

# Word Embeddings (dense vector)



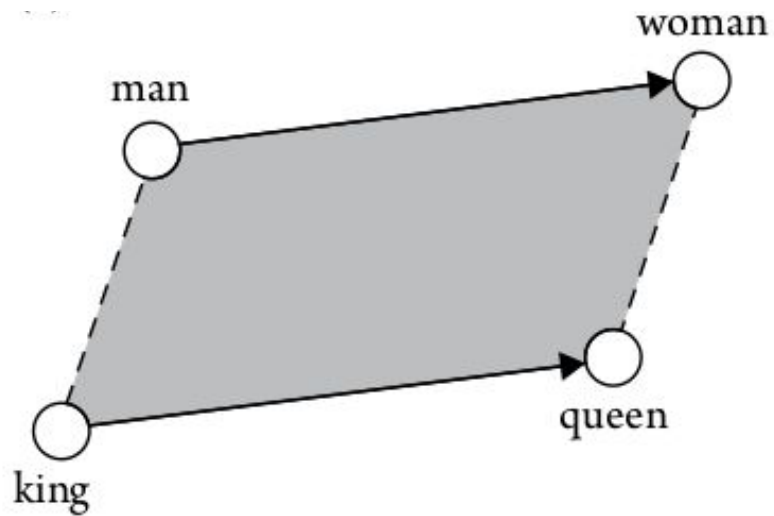
Horse

Horse

=

$$\begin{bmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{bmatrix}$$

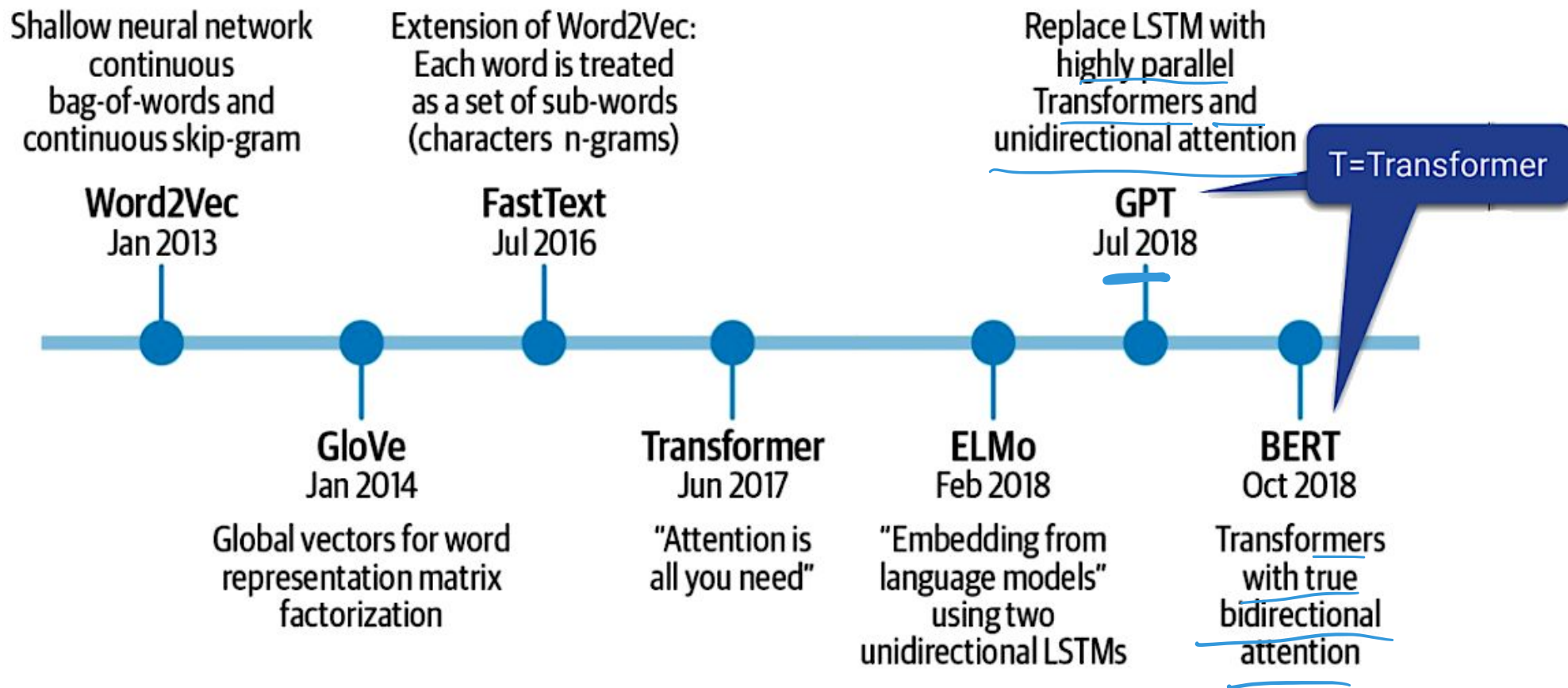
# Word analogies as parallelograms



Man is to king as  
woman is to queen.

$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{king}} - \overrightarrow{\text{queen}}$$

# Progress over the years



# Word2vec →

## Input:

- Text corpora: Wikipedia, Twitter, Common Crawl.
- $V$ : a predefined vocabulary
- $d$ : dimension of word vectors

## Output:

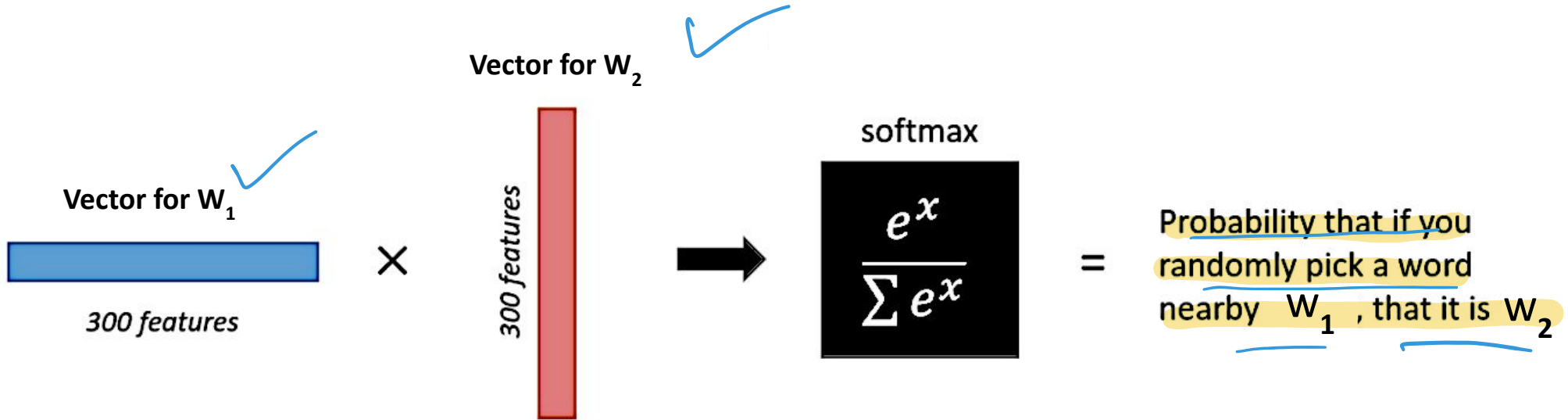
$$f : V \rightarrow \mathbb{R}^d$$

← Train a classifier on prediction task

**Idea:** Train a classifier on a prediction task:  $w_1$  likely to show up near  $w_2$ ?

Use running text as implicitly supervised training data!

# Word2vec (intuitive idea)





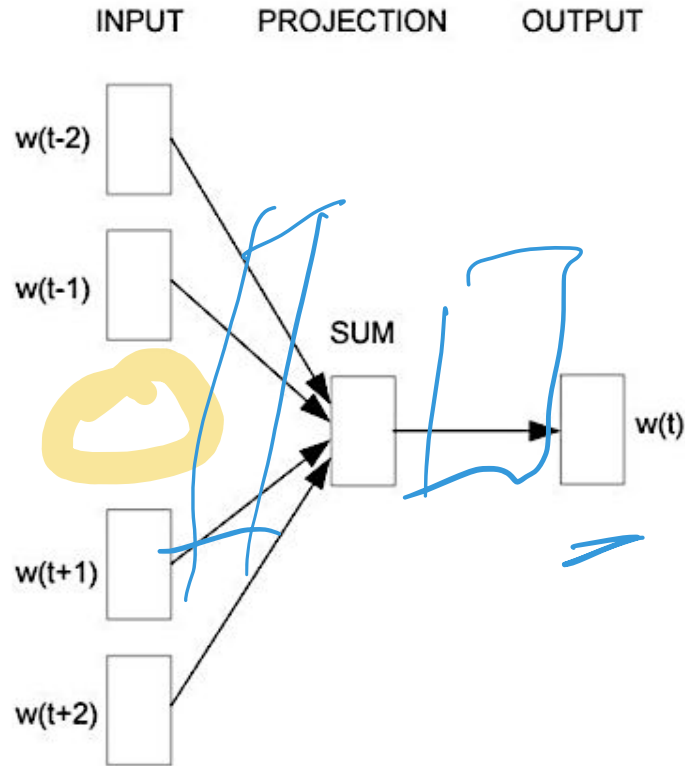
# Main Idea of word2vec

- The word2vec model is a neural network-based approach for generating word embeddings.
- The model is trained on a large corpus of text data using either a skip-gram or a continuous bag-of-words (CBOW) architecture.
- During training, the model learns to associate words in the embedding space that appear in similar contexts.

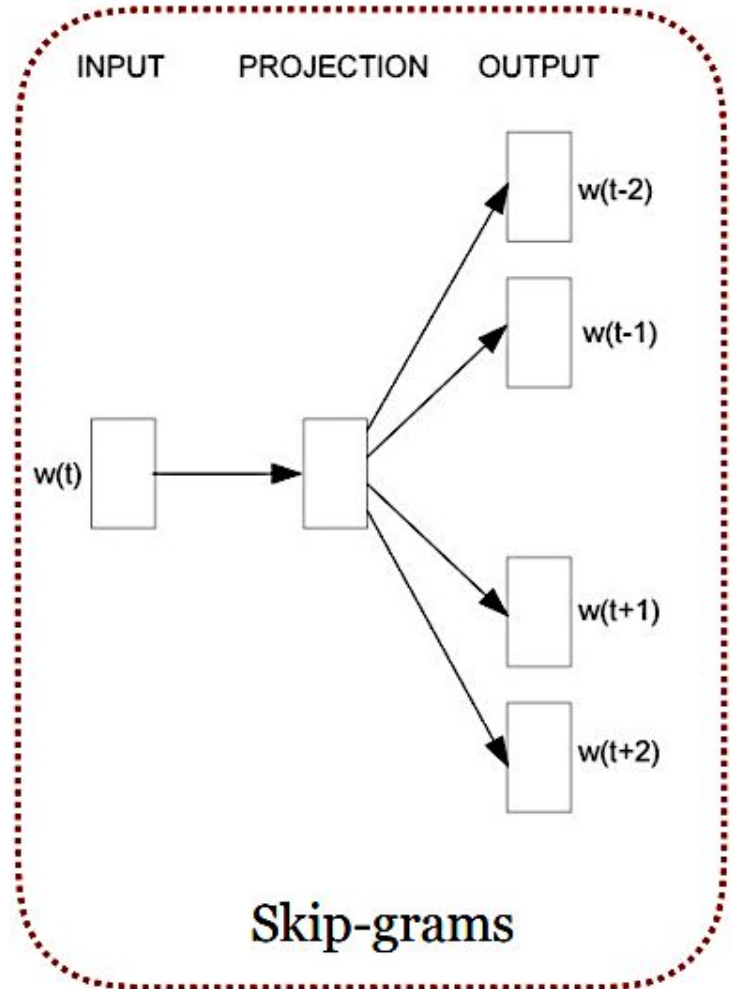
# Word2vec

*Context window = 2*

*$w_t$*



Continuous Bag of Words (CBOW)

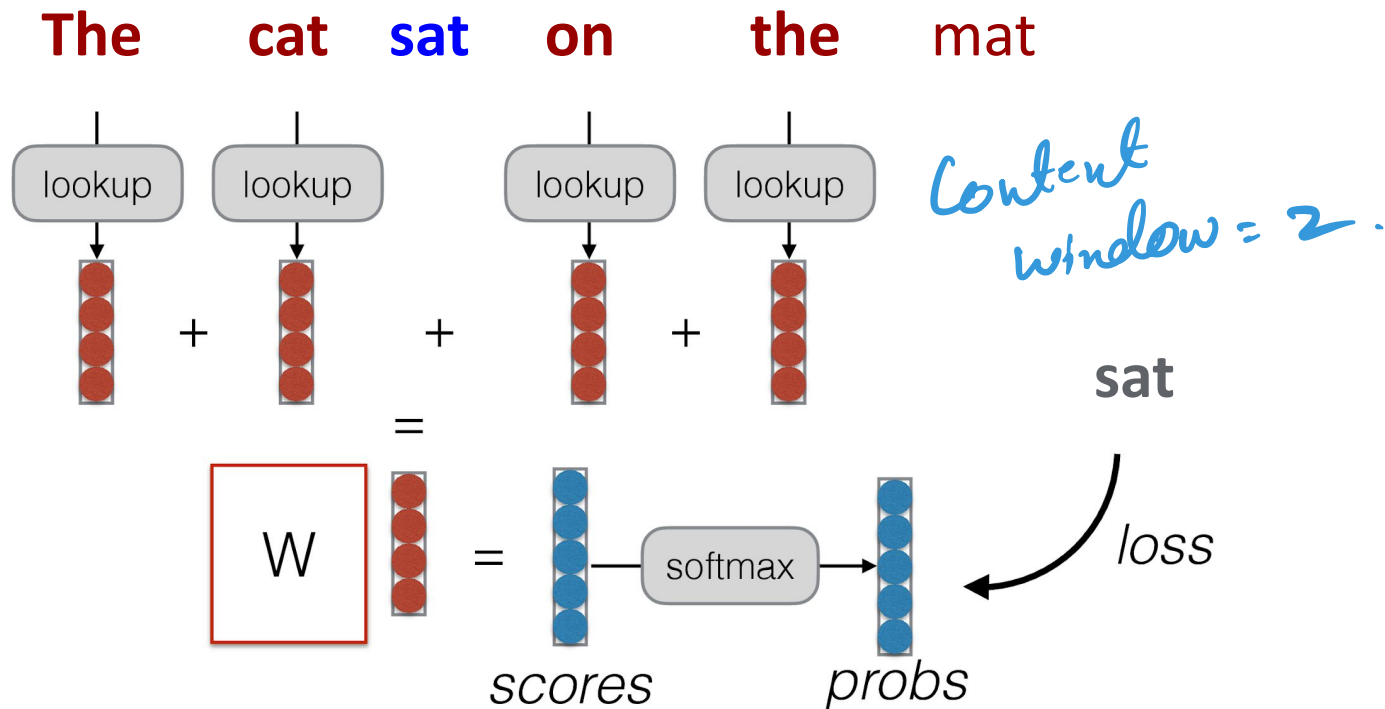


Skip-grams

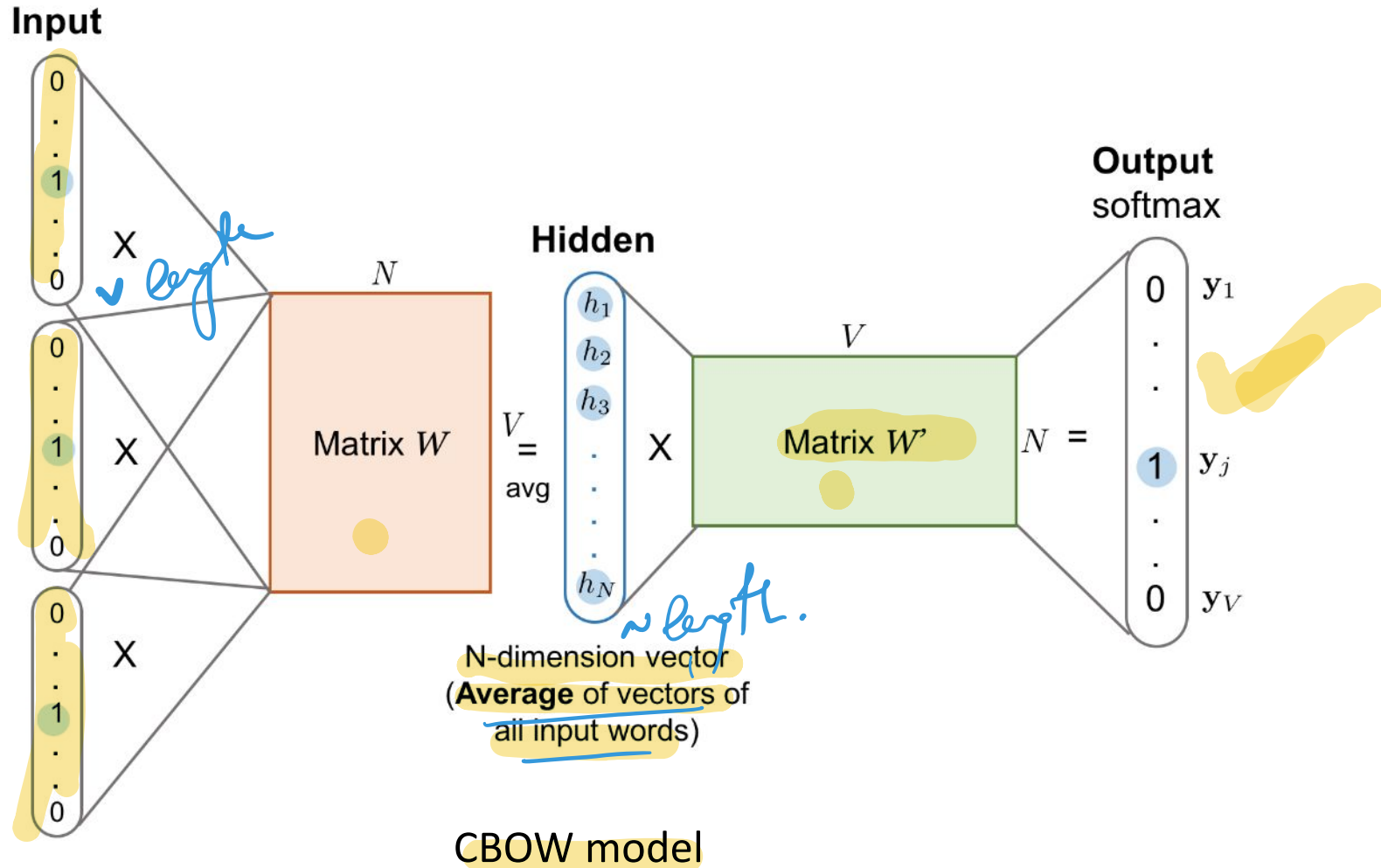
# Continuous Bag of Words Model (CBOW)

# Continuous Bag of Words Model (CBOW)

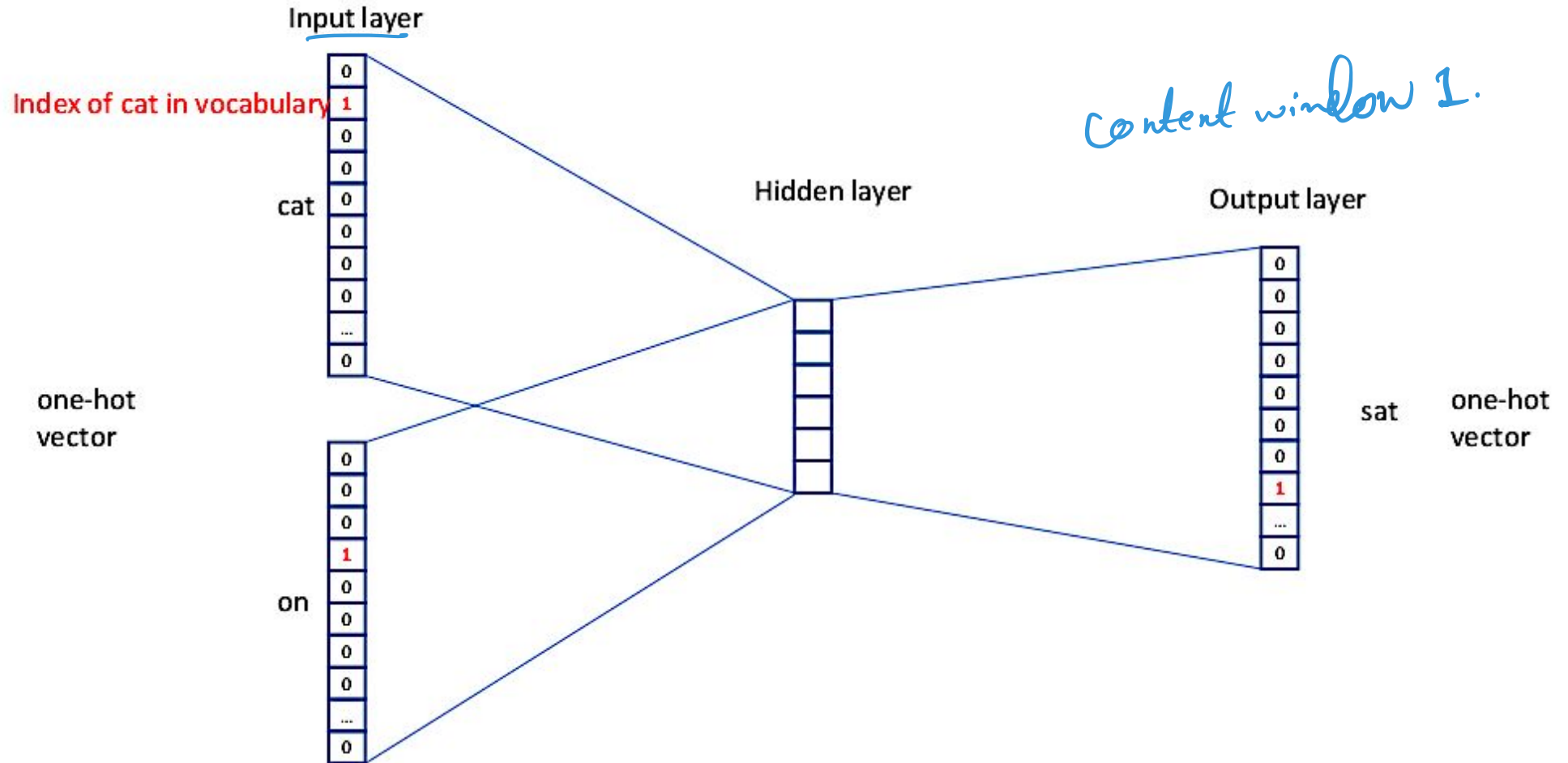
- Generate/predict the center word by considering the context words.



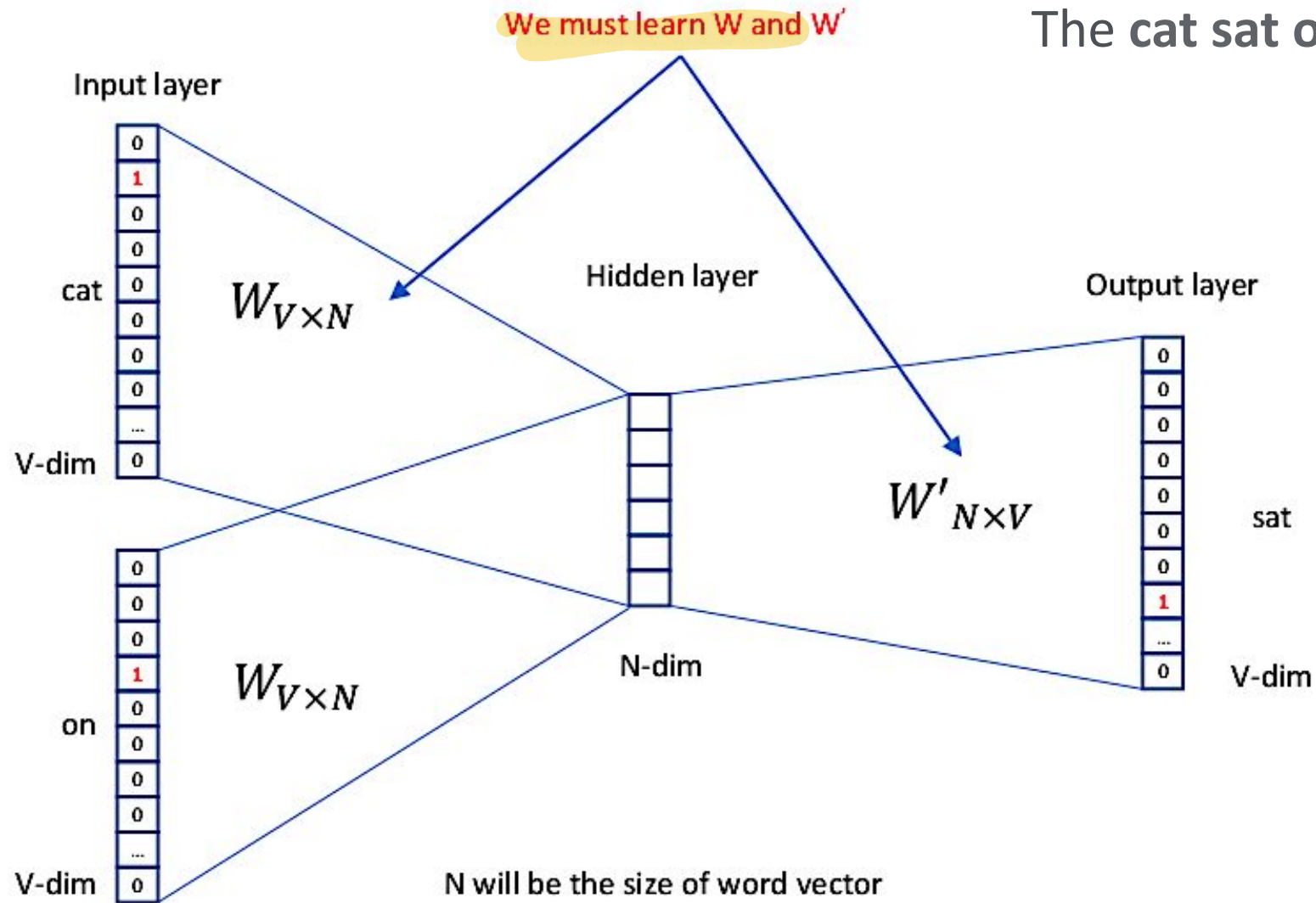
# Continuous Bag of Words Model (CBOW)

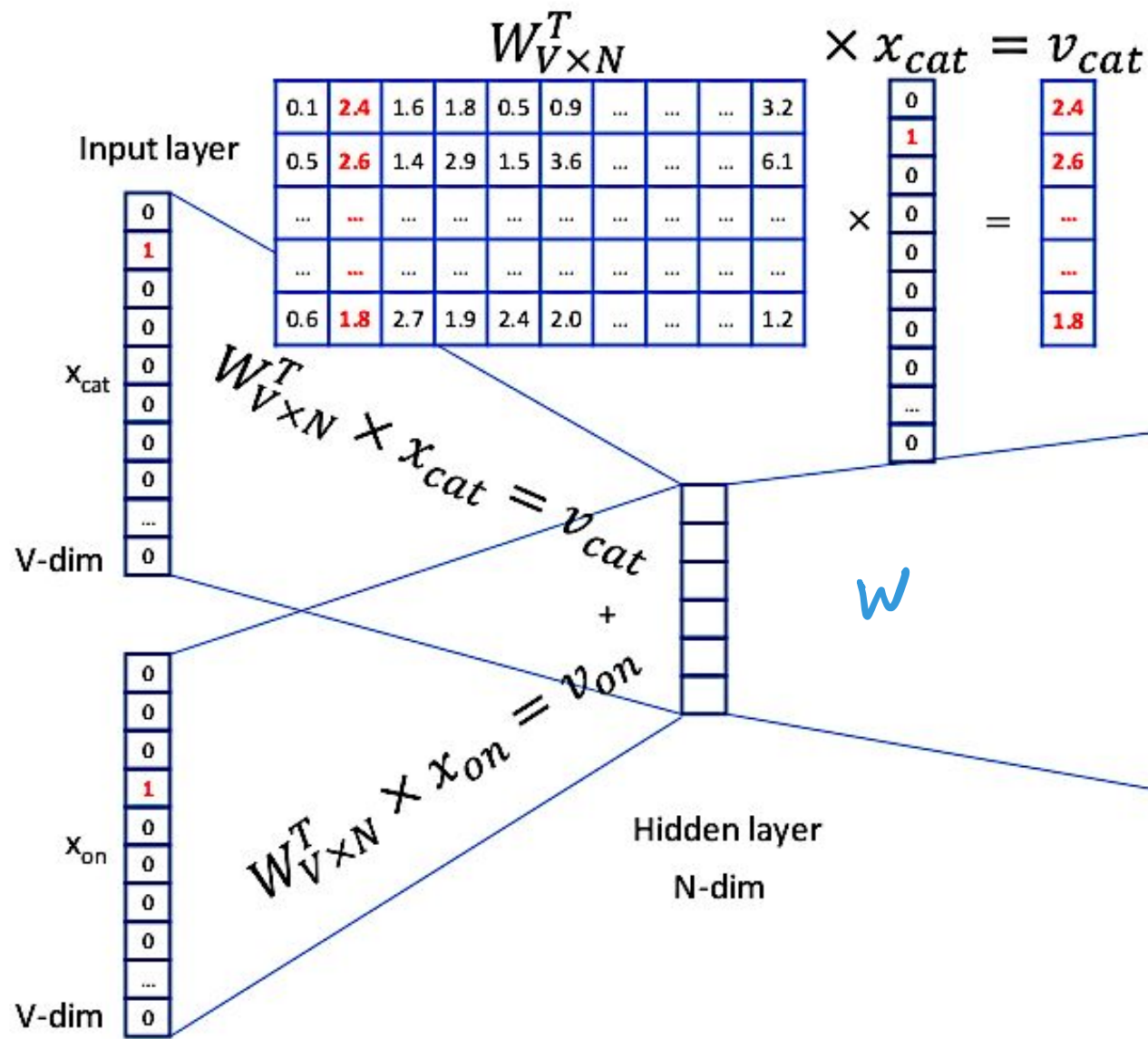


The **cat** sat on the mat

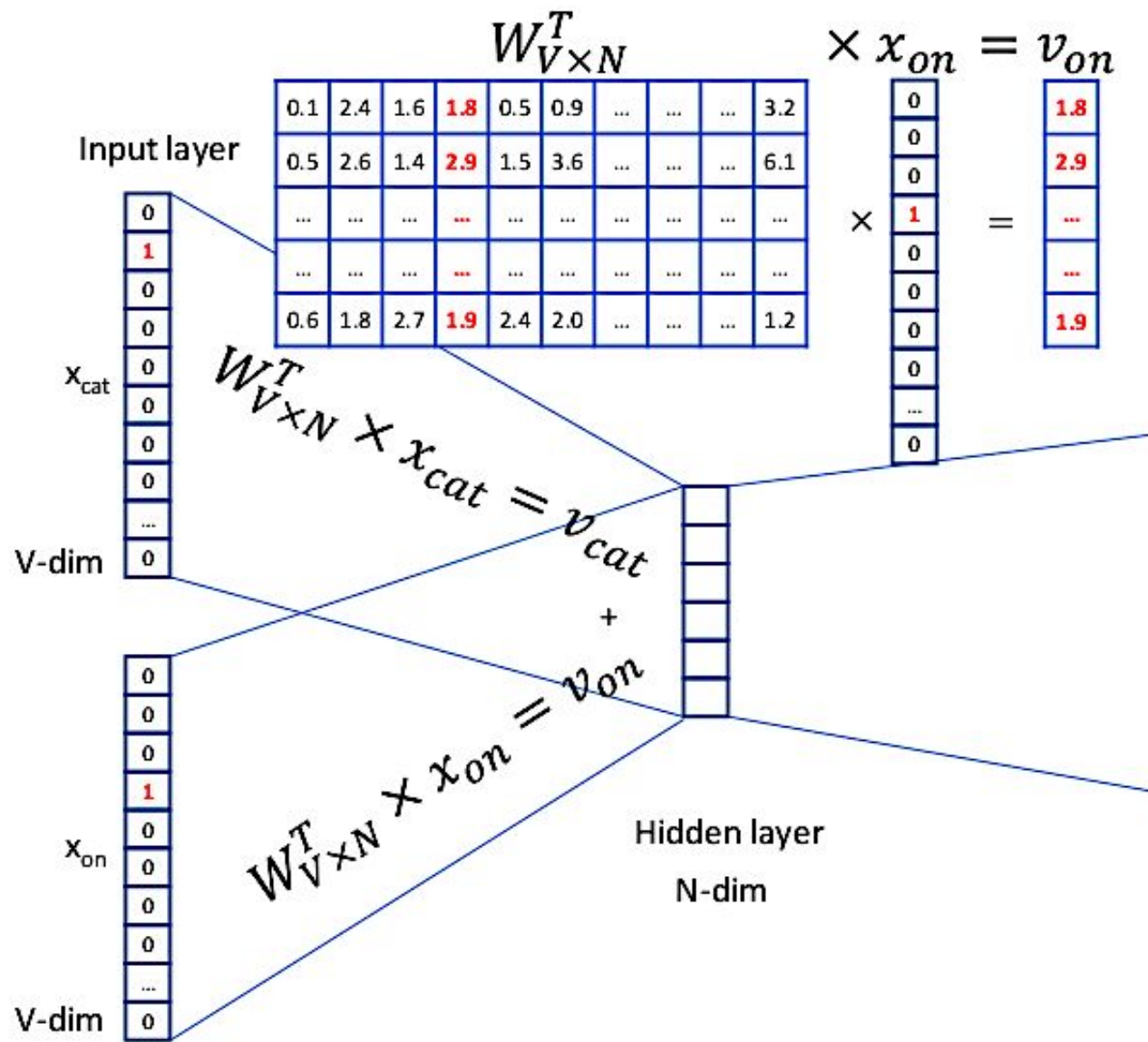


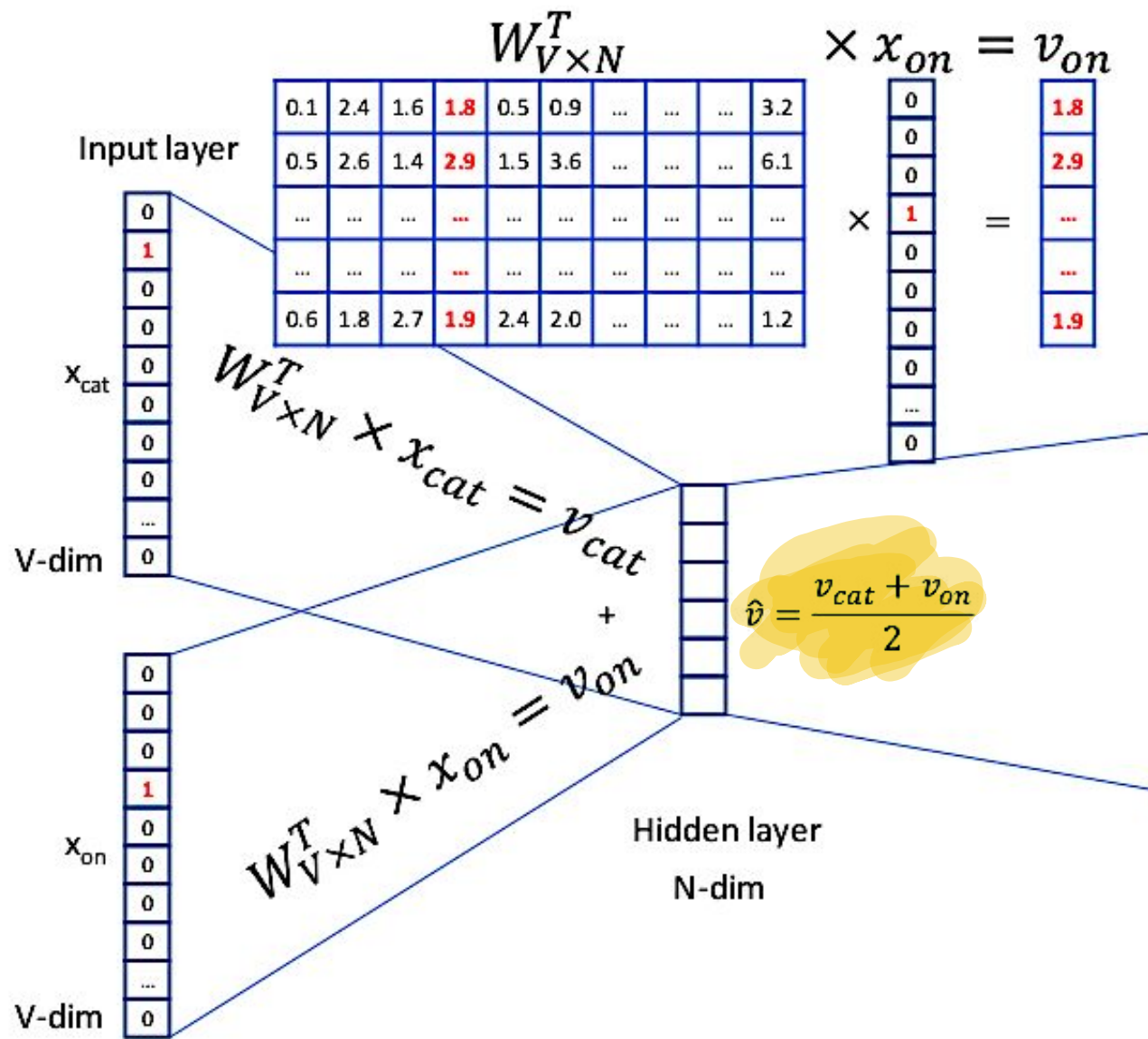
The cat sat on the mat

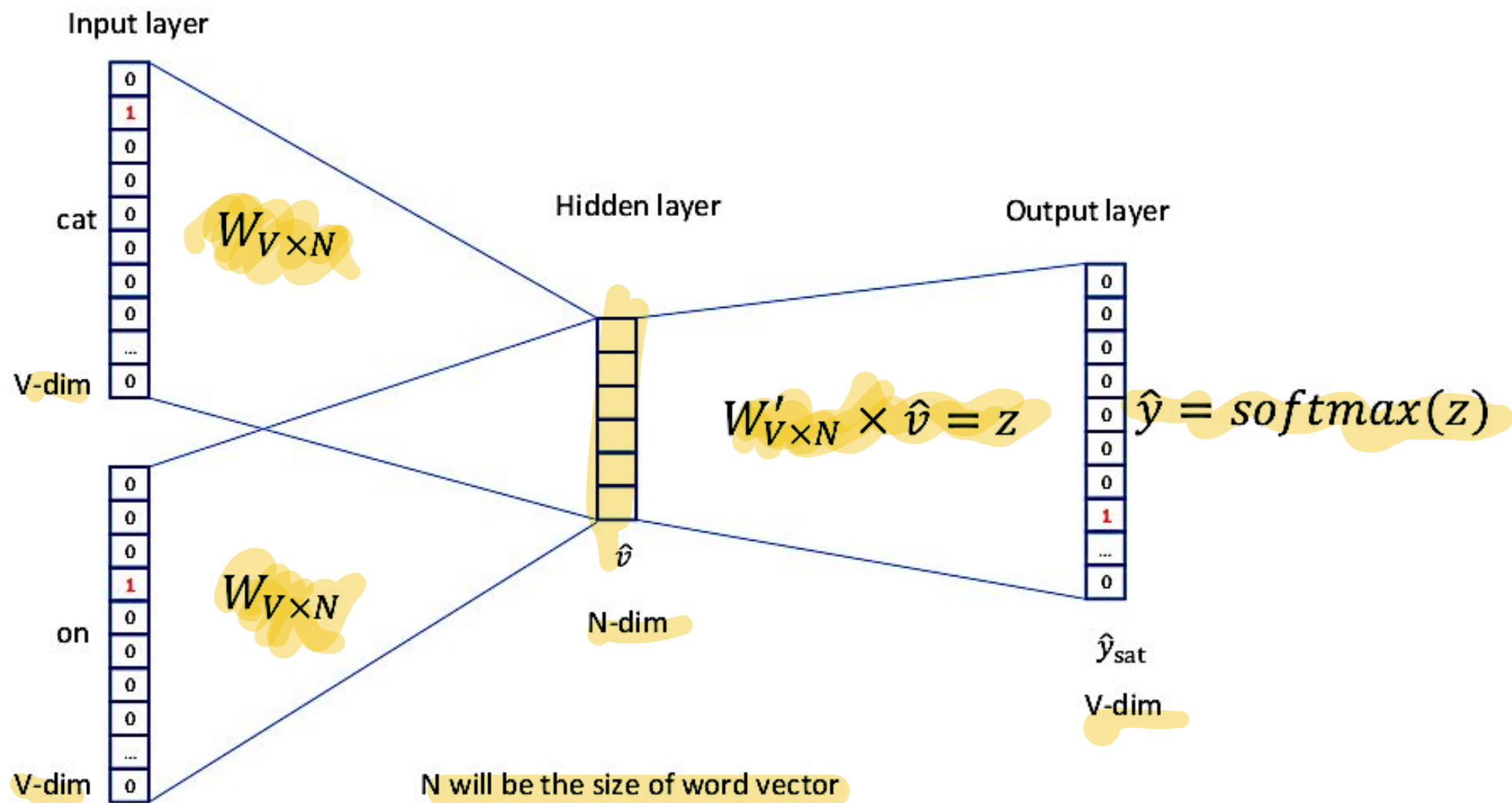


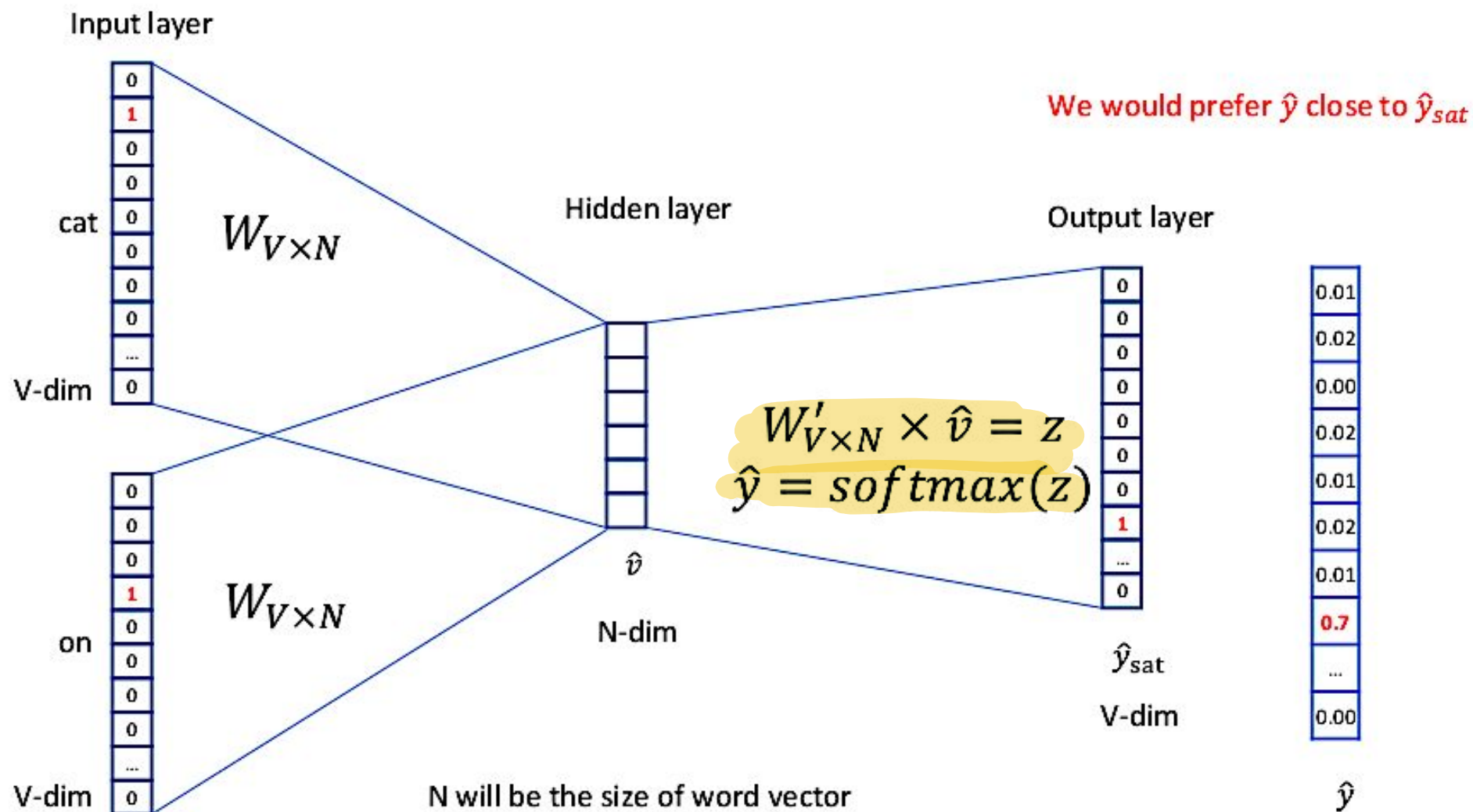


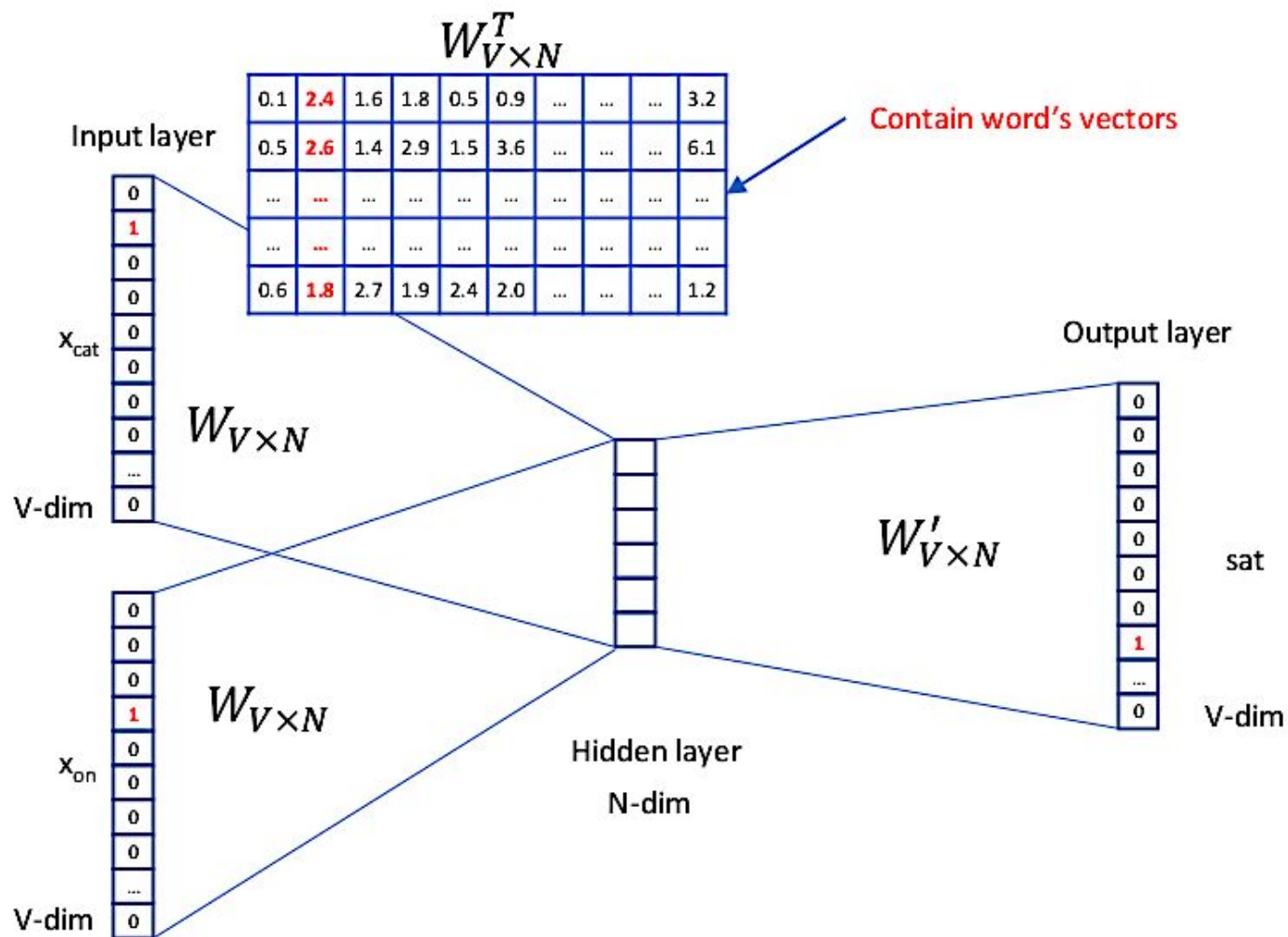










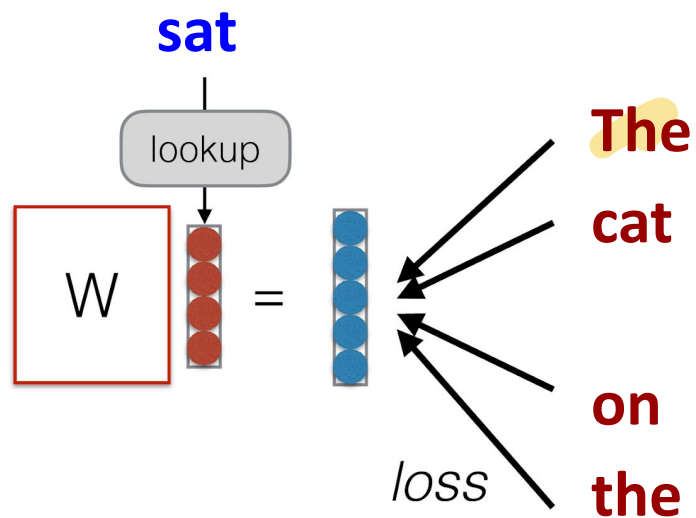


We can consider either  $W$  or  $W'$  as the word's representation. Or even take the average.

# Skip-Gram Task

# Skip-gram

- Predict each word in the context given the word



## Skip-gram training data example

*data of CBOW will be the reverse.*  
The man who passes the sentence should swing the sword

Input

The

man

who

passes

*context window = 2*

*Context words*

*["man" "who"]*

*The who passes.*

*The man passes*



# Skip-gram training data example

The man who passes the sentence should swing the sword

Sliding window (size = 5)	Target word	Context
[The man who]	the	man, who
[The man who passes]	man	the, who, passes
[The man who passes the]	who	the, man, passes, the
[man who passes the sentence]	passes	man, who, the, sentence
...	...	...
[sentence should swing the sword]	swing	sentence, should, the, sword
[should swing the sword]	the	should, swing, sword
[swing the sword]	sword	swing, the

# Skip-gram

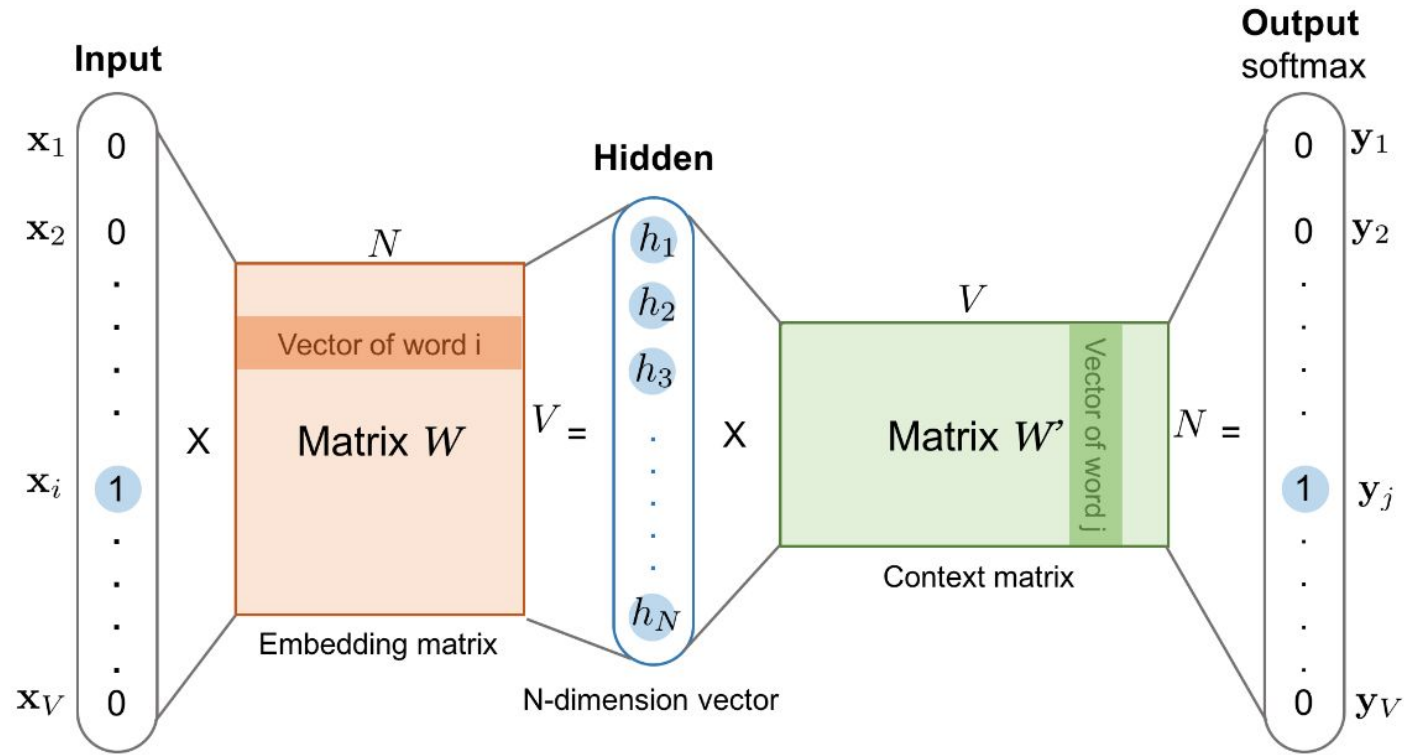
The man who passes the sentence should swing the sword

- Each context-target pair is treated as a new observation in the data.
- For example, the target word “swing” in the above case produces four training samples:
  - (“swing”, “sentence”), (“swing”, “should”), (“swing”, “the”), and (“swing”, “sword”).

swing

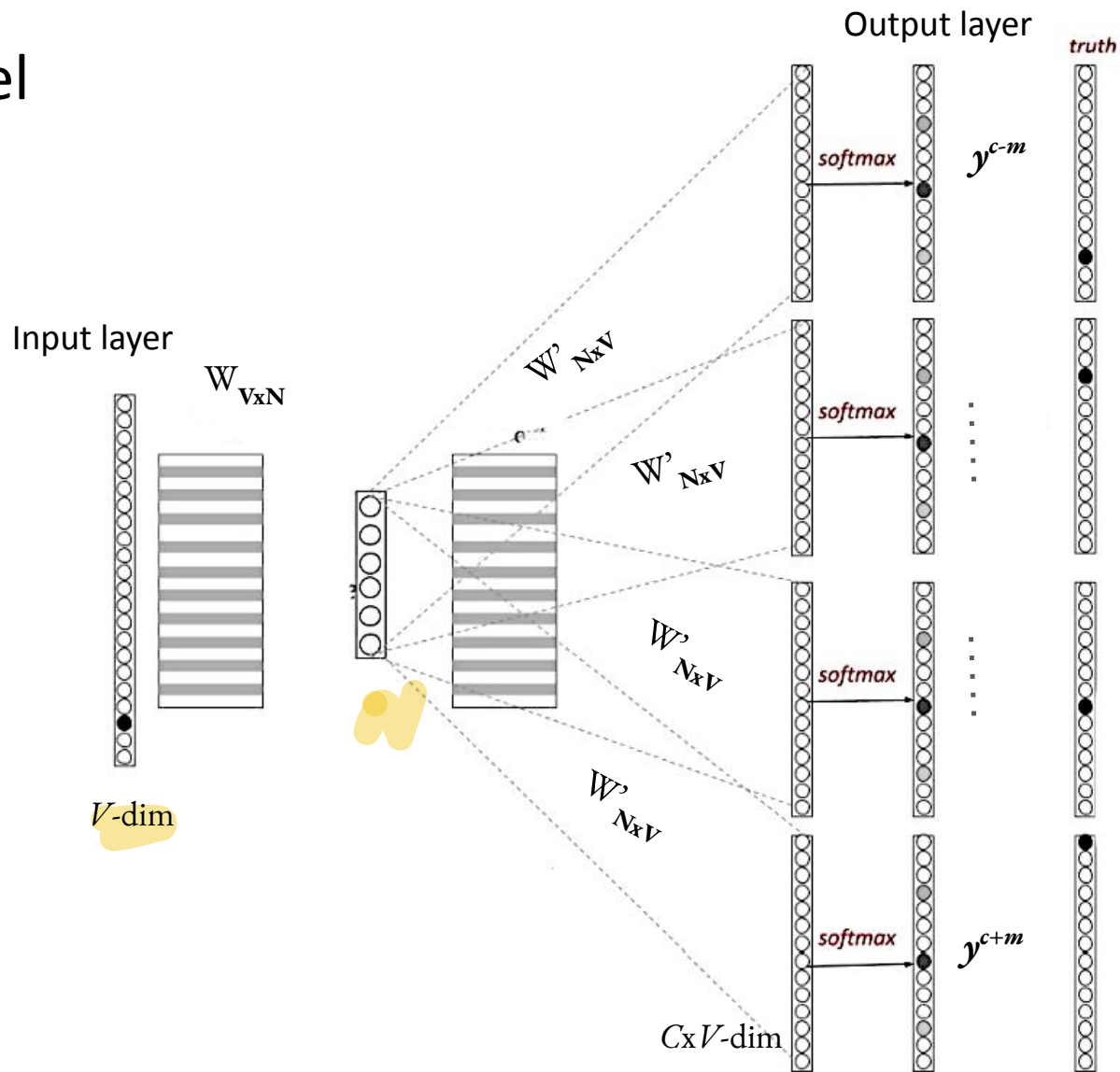
sentence, should, the, sword

# Skip-gram



The skip-gram model. Both the input vector  $x$  and the output  $y$  are one-hot encoded word representations. The hidden layer is the word embedding of size  $N$ .

# Skip-gram Model



# FastText

$$s(w, c)$$

↑  
context

$$p(w_t | w_c) = \frac{e^{s(w_c, w_t)}}{\sum}$$

$$\sum_a z_g v_c$$

# FastText

Represents a word  
by sum of its  $n$ -grams.

- It enhances the Word2Vec algorithm
- Limitation of Word2Vec:
  - Out of Vocabulary(OOV) Words
    - **Tensor flow** TensorFlow
  - Morphology
    - go, goes, going, gone
- FastText overcome these limitations
- FastText allows to compute word representation for words not in training data

FastText and n-grams → represent words as a sum of its character n-grams.

- FastText is the modification to the skip-gram method Sub-word generation using generation of character n-grams of length 3 to 6.
- E.g., *where* with  $n=3$ : <wh, whe, her, ere, re>, <where>
- Represent words as sum of its character n-grams
- Grammatical variations still share most of n-grams.
- Compound nouns are easy to model, e.g., Noun + Noun: lunchtime

goes  
↓  
<go, goe, goes, es>

# FastText

- Let us consider that we are given a scoring function  $s(w, c)$  which maps pairs of (word, context) to scores.



# FastText

- Let us consider that we are given a scoring function  $s(w, c)$  which maps pairs of (word, context) to scores.
- Suppose that you are given a dictionary of  $n$  grams of size  $G$ .

# FastText

- Let us consider that we are given a scoring function  $s(w, c)$  which maps pairs of (word, context) to scores.
- Suppose that you are given a dictionary of  $n$  grams of size  $G$ .
- Given a word  $w$ , let us denote by  $\mathcal{G}_w \subset \{1, \dots, G\}$  the set of  $n$ -grams appearing in  $w$ .

# FastText

- Let us consider that we are given a scoring function  $s(w, c)$  which maps pairs of (word, context) to scores.
- Suppose that you are given a dictionary of  $n$  grams of size  $G$ .
- Given a word  $w$ , let us denote by  $\mathcal{G}_w \subset \{1, \dots, G\}$  the set of  $n$ -grams appearing in  $w$ .
- We associate a vector representation  $\mathbf{z}_g$  to each  $n$ -gram  $g$ .

# FastText

- Let us consider that we are given a scoring function  $s(w, c)$  which maps pairs of (word, context) to scores.
- Suppose that you are given a dictionary of  $n$  grams of size  $G$ .
- Given a word  $w$ , let us denote by  $\mathcal{G}_w \subset \{1, \dots, G\}$  the set of  $n$ -grams appearing in  $w$ .
- We associate a vector representation  $\mathbf{z}_g$  to each  $n$ -gram  $g$ .
- We represent a word by the sum of the vector representations of its  $n$ -grams. We thus obtain the scoring function:

$$s(w, c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^\top \mathbf{v}_c. \quad (1)$$

# FastText

- Let us consider that we are given a scoring function  $s(w, c)$  which maps pairs of (word, context) to scores.
- Suppose that you are given a dictionary of  $n$  grams of size  $G$ .
- Given a word  $w$ , let us denote by  $\mathcal{G}_w \subset \{1, \dots, G\}$  the set of  $n$ -grams appearing in  $w$ .
- We associate a vector representation  $\mathbf{z}_g$  to each  $n$ -gram  $g$ .
- We represent a word by the sum of the vector representations of its  $n$ -grams. We thus obtain the scoring function:

$$s(w, c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^\top \mathbf{v}_c. \quad (1)$$

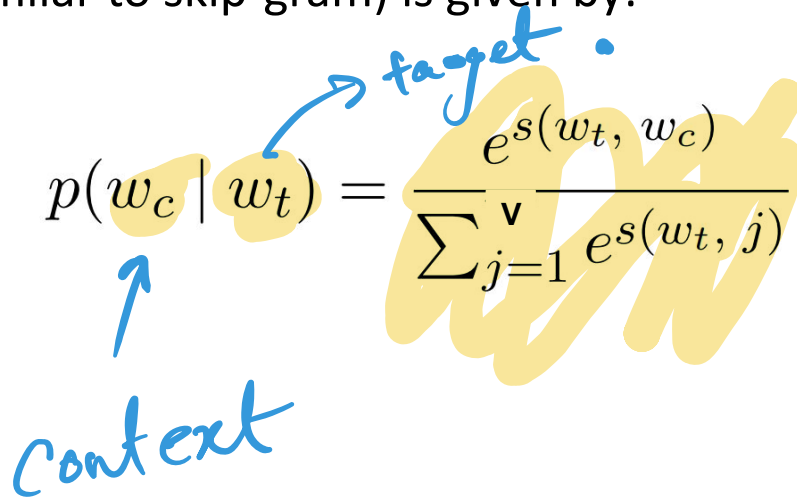
- This simple model allows sharing the representations across words, thus allowing to learn reliable representation for rare words.

# FastText

The scoring function is described as:

$$s(w, c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^\top \mathbf{v}_c.$$

Now, the final model (similar to skip-gram) is given by:



The diagram shows the equation  $p(w_c | w_t) = \frac{e^{s(w_t, w_c)}}{\sum_{j=1}^v e^{s(w_t, j)}}$  with handwritten blue annotations. An arrow points from the word  $w_c$  to the word "context". Another arrow points from the word  $w_t$  to the word "target". The entire equation is highlighted with a large yellow brushstroke.

$$p(w_c | w_t) = \frac{e^{s(w_t, w_c)}}{\sum_{j=1}^v e^{s(w_t, j)}}$$

# References

- Deep Learning: A Visual Approach Book by Andrew Glassner.
- Bolukbasi, Tolga, et al. "Man is to computer programmer as woman is to homemaker? debiasing word embeddings." Advances in neural information processing systems (2016).
- Practical Natural Language Processing: A Comprehensive Guide to Building Real-World NLP Systems by Anuj Gupta.
- Bojanowski, Piotr, Edouard Grave, Armand Joulin, and Tomas Mikolov. "Enriching word vectors with subword information." Transactions of the association for computational linguistics 5 (2017): 135-146.
- <https://lilianweng.github.io/posts/2017-10-15-word-embedding/>