

AI VIETNAM
All-in-One Course
(TA Session)

Introduction to Text Classification and POS Tagging



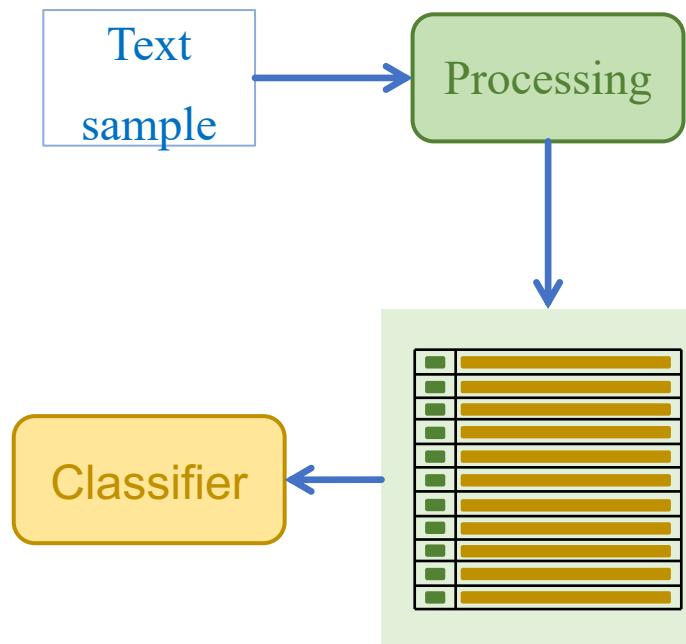
AI VIET NAM
[@aivietnam.edu.vn](http://aivietnam.edu.vn)

**Thuan Duong – TA
Nguyen Nha – TA
Thu Van – sTA**

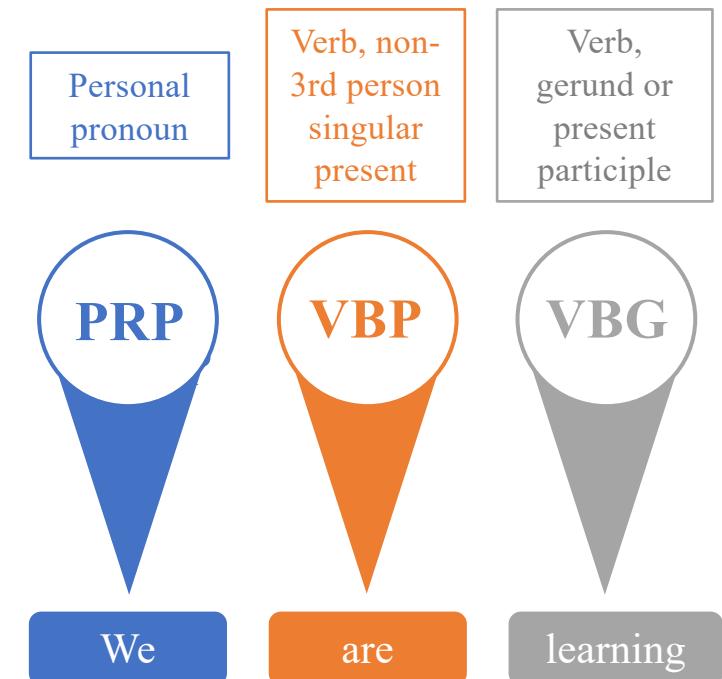
Objectives



❖ Text Classification



❖ POS Tagging



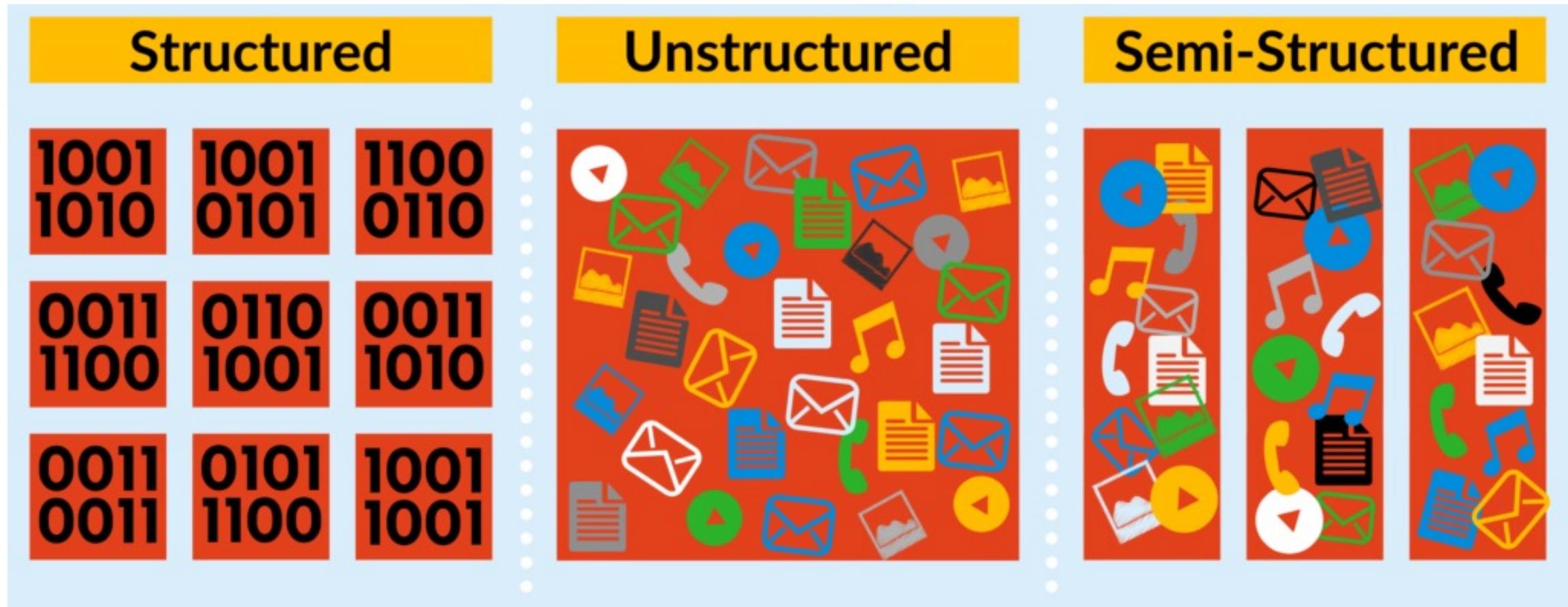
Outline

- Introduction to NLP
- Text Classification
- Part-of-Speech Tagging
- Question

Introduction to NLP

Introduction to NLP

❖ Data Types



Introduction to NLP

❖ Structured Data

Data organized in a predefined format (rows and columns)

order_id	name	order_date	total_amount	name	total_spent
1	Alice Johnson	2024-02-01	150	Bob Smith	300
2	Alice Johnson	2024-02-05	200	Alice Johnson	350
3	Bob Smith	2024-02-03	300	Charlie Brown	50
4	Charlie Brown	2024-02-06	50		



Tables in database

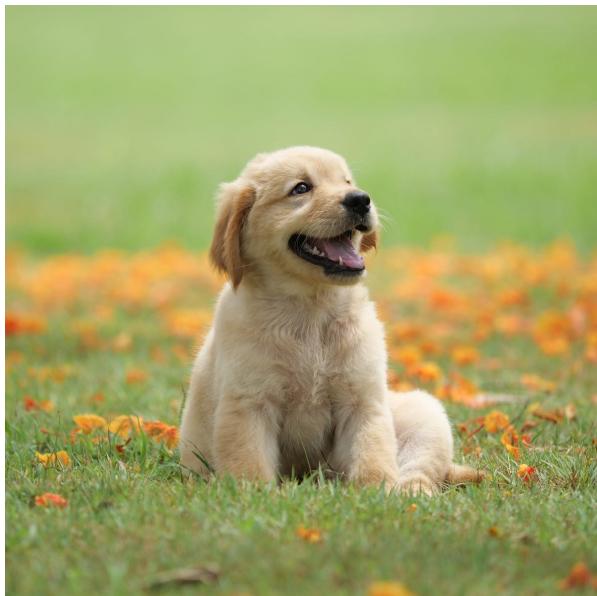
Excel sheets

Pandas DataFrame

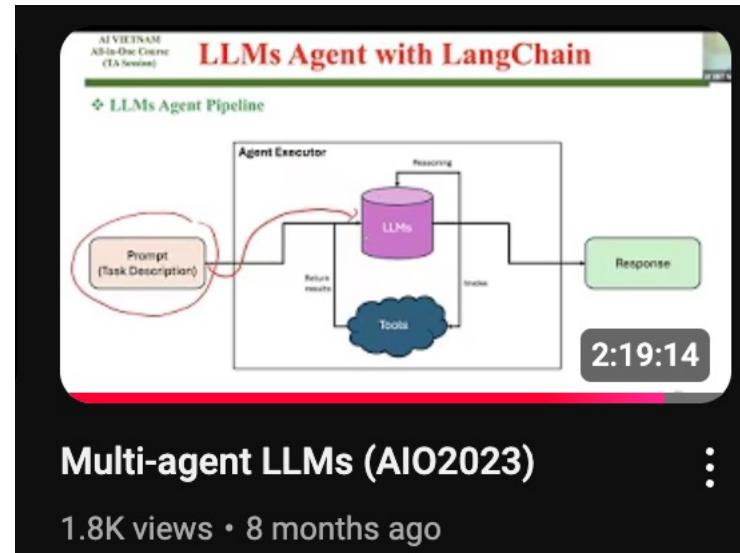
Introduction to NLP

❖ Unstructured Data

Data that does not have a predefined structure or format.



Image



Video

Congratulations!

Your Certificate is Ready

Dear Duong Nguyen Thuan,

Congratulations! You've successfully completed Deep Neural Networks. On behalf of AI Viet Nam we are pleased to issue your official Deep Neural Networks Certificate.

Free-Text

Introduction to NLP

❖ Semi-Structured Data

Semi-structured data is a type of data that does not have a rigid structure like relational databases but still contains some organizational properties.

```
{  
  "_id": "123",  
  "user": "john_doe",  
  "email": "john@example.com",  
  "preferences": {  
    "theme": "dark",  
    "notifications": true  
  }  
}
```

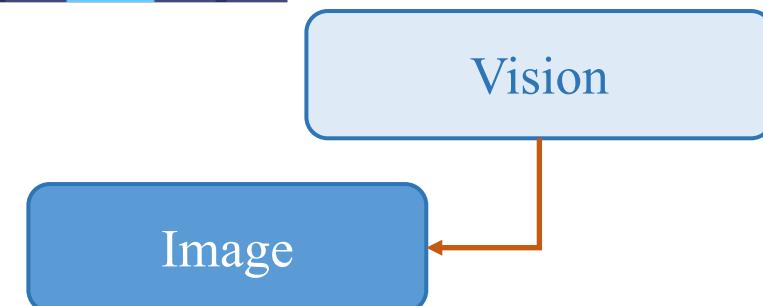
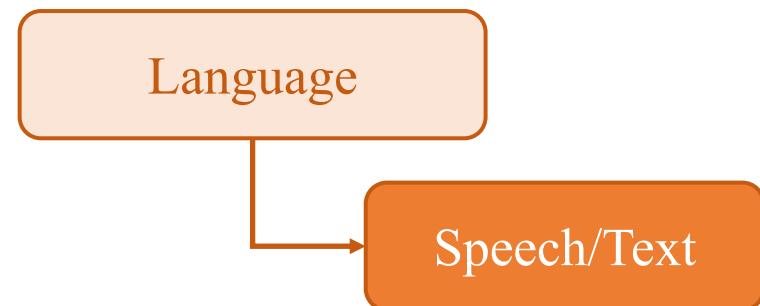
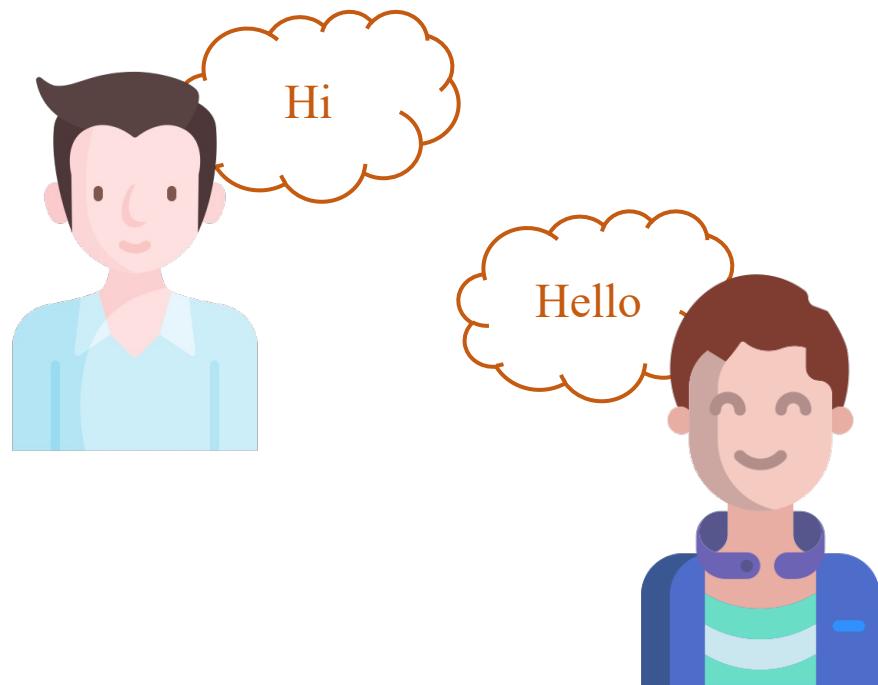
Name, Age, Skills
Alice, 25, "Python, Machine Learning"
Bob, 30, "Java, DevOps, Cloud"

CSV

NoSQL Database

Introduction to NLP

❖ Human vs Computer



Introduction to NLP

❖ Human vs Computer



Collaborative coding with AI tools is allowed

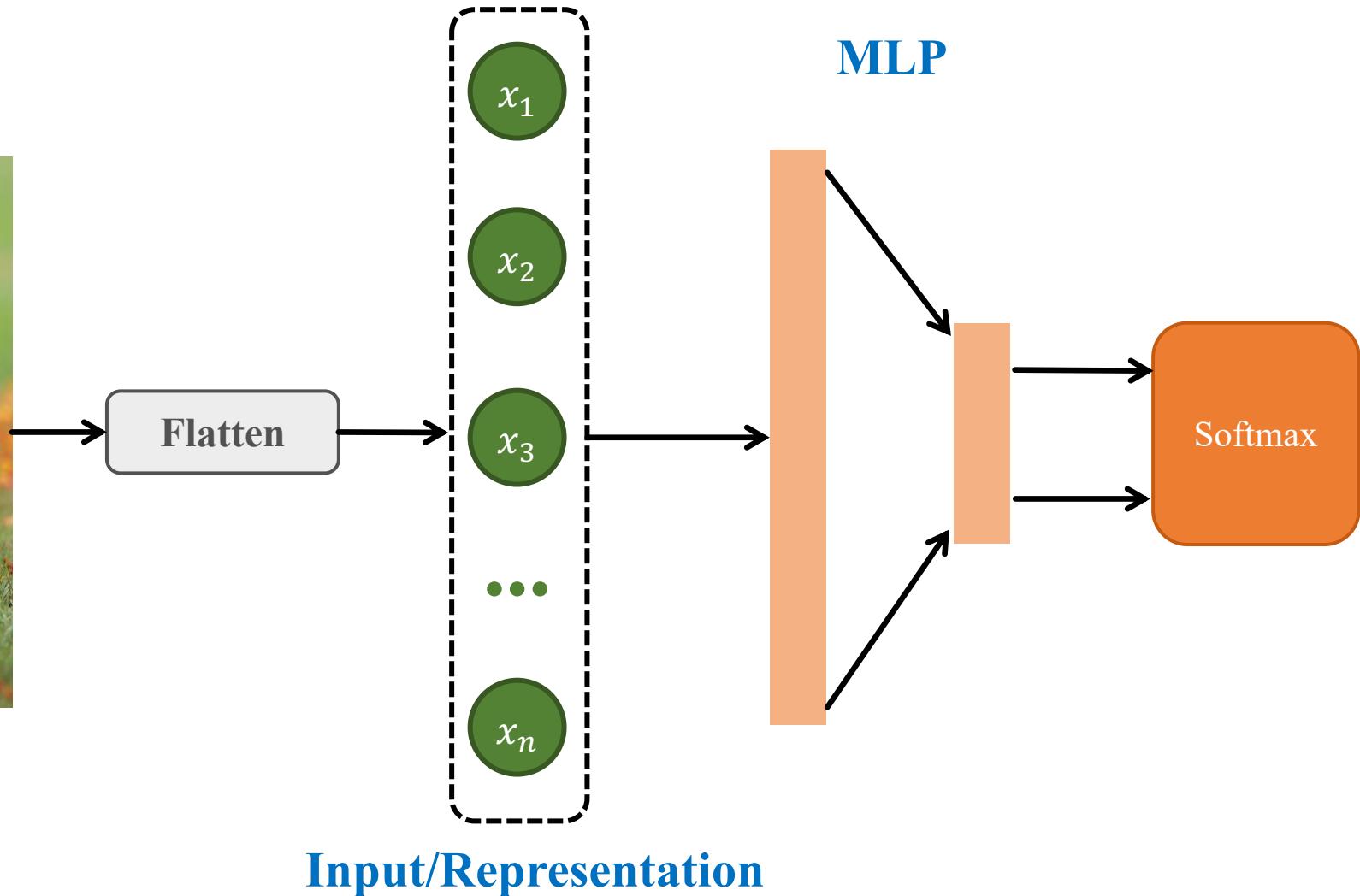
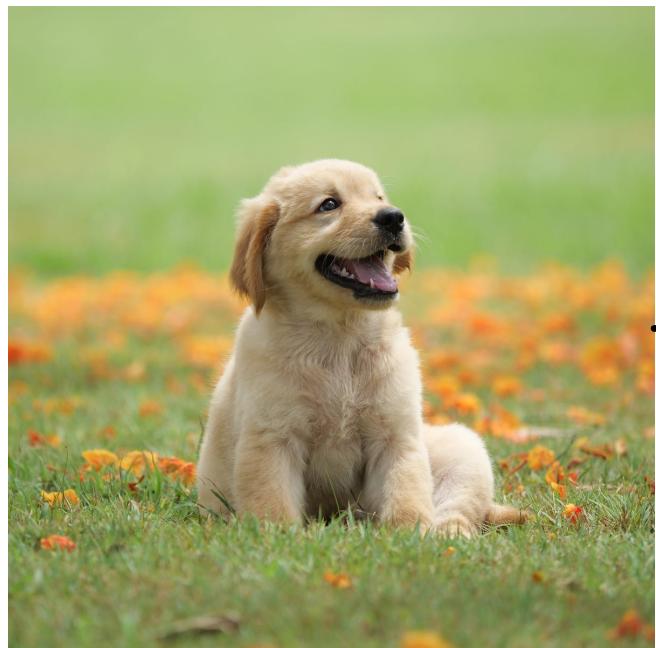
Speech/Text

Image



Introduction to NLP

❖ Vision in Computer

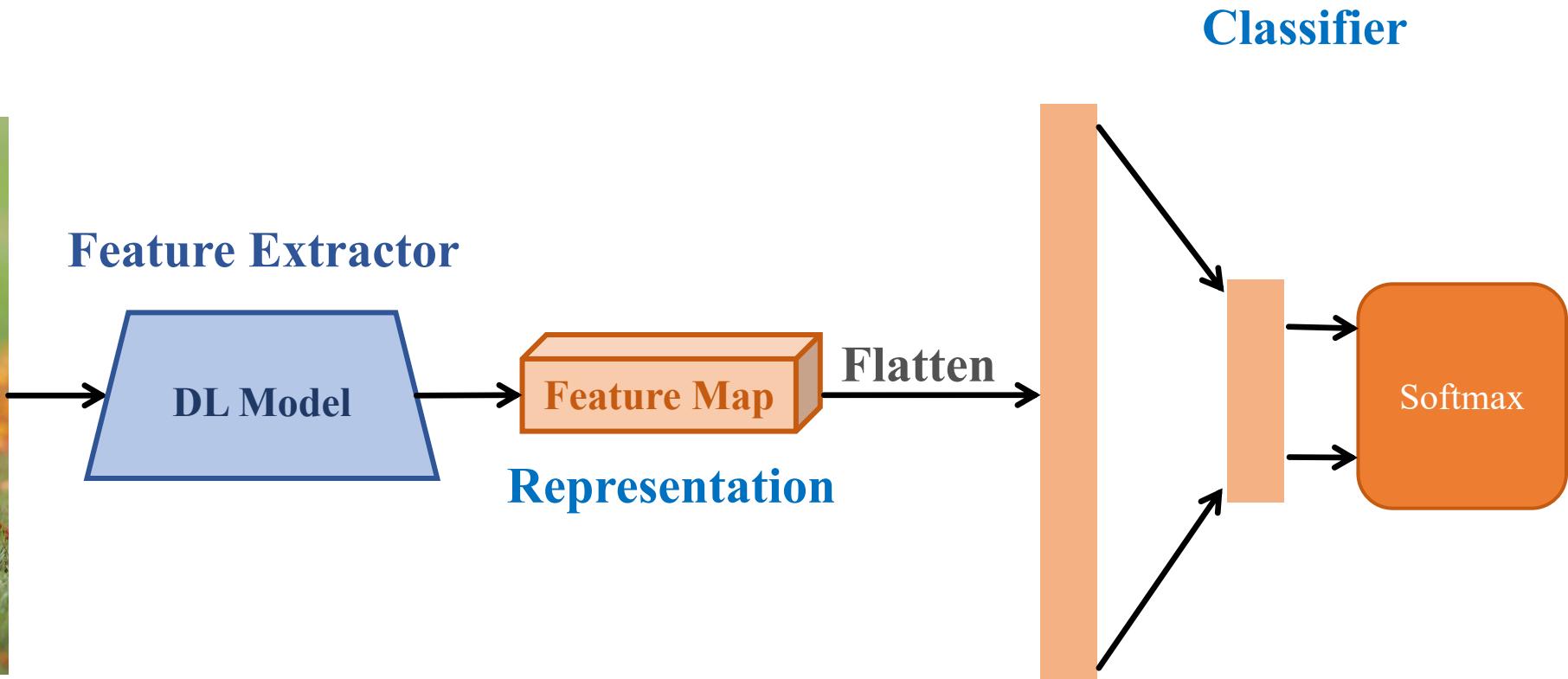


Introduction to NLP

❖ Vision in Computer

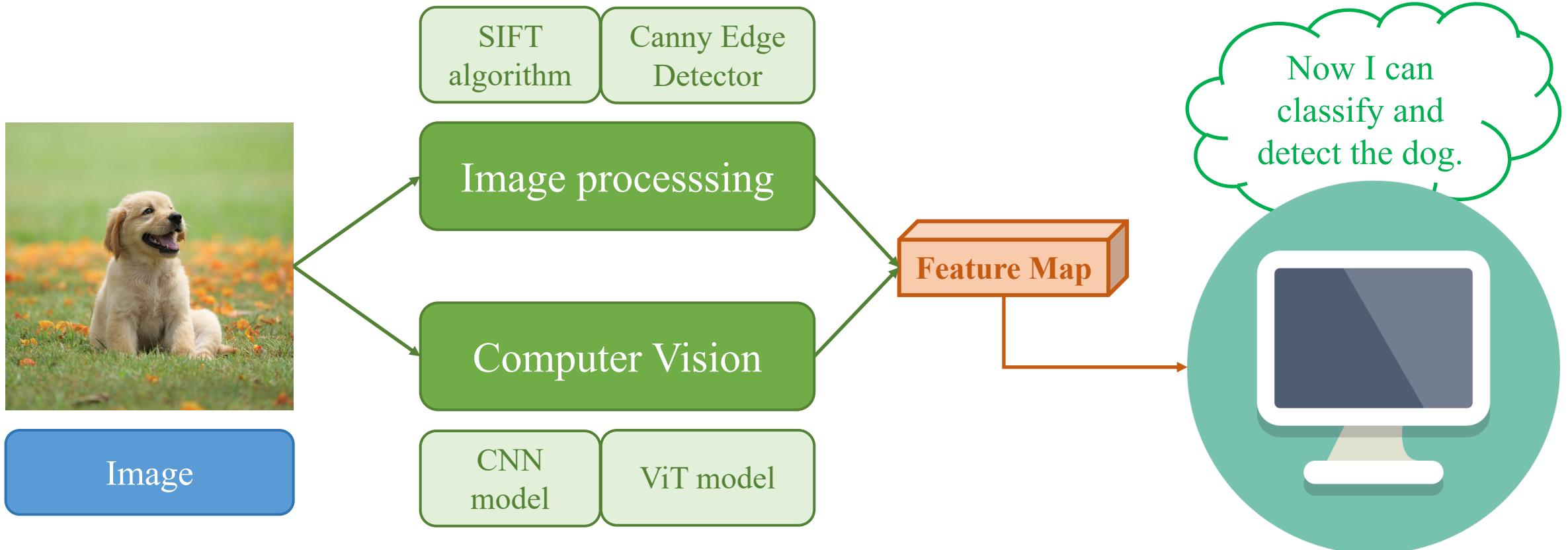


Input



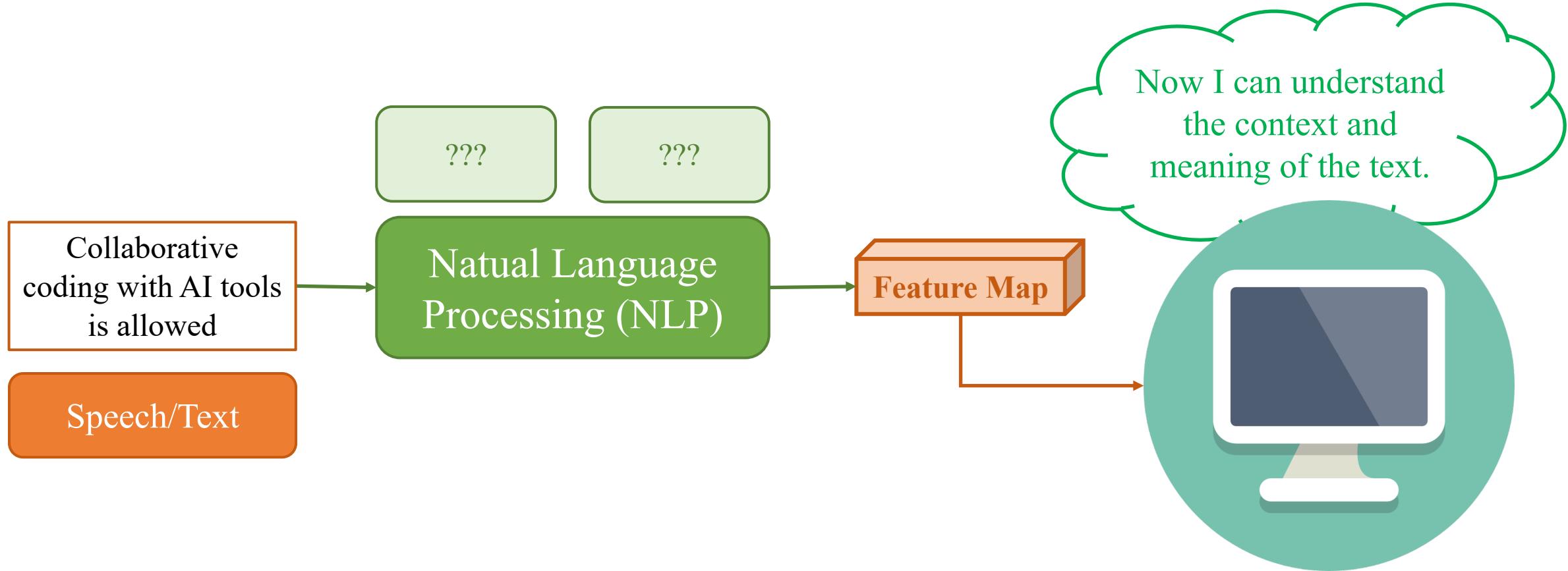
Introduction to NLP

❖ Vision in Computer



Introduction to NLP

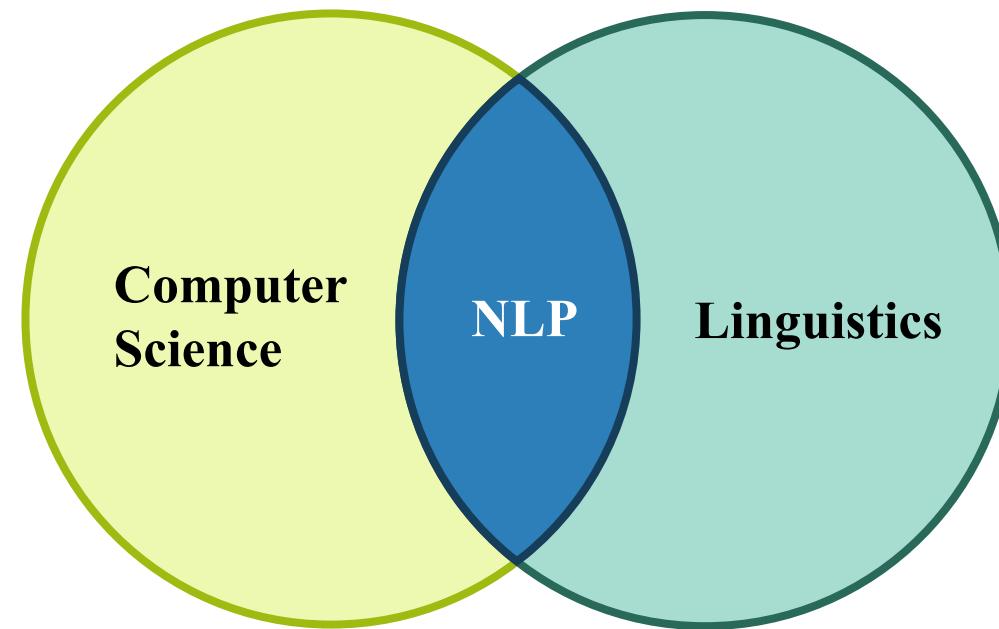
❖ Natural Language in Computer



Introduction to NLP

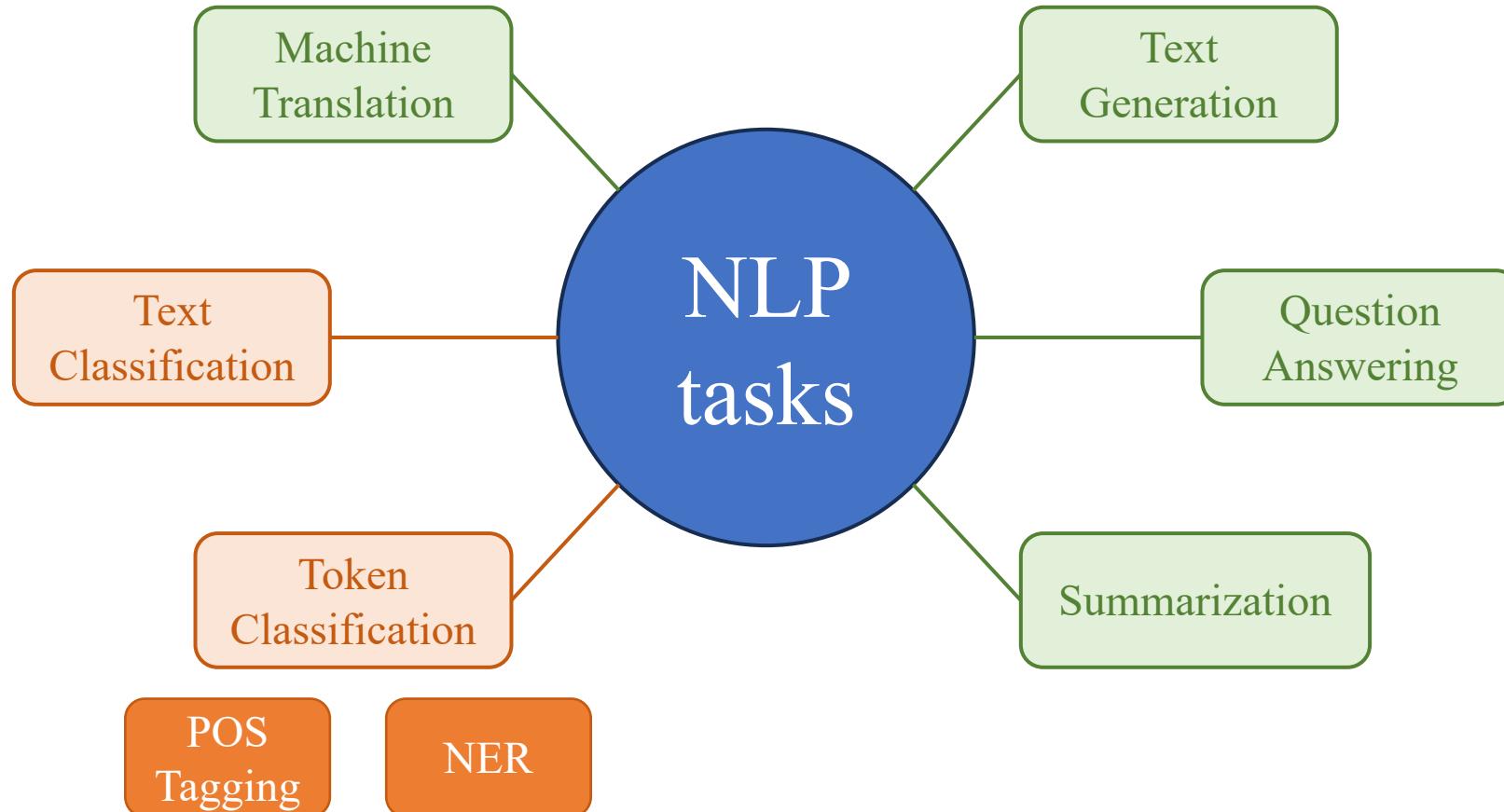
❖ Getting Started

Natural language processing (NLP) is a subfield of computer science and especially artificial intelligence. It is primarily concerned with providing computers with the ability to process data encoded in natural language and is thus closely related to information retrieval, knowledge representation and computational linguistics, a subfield of linguistics.



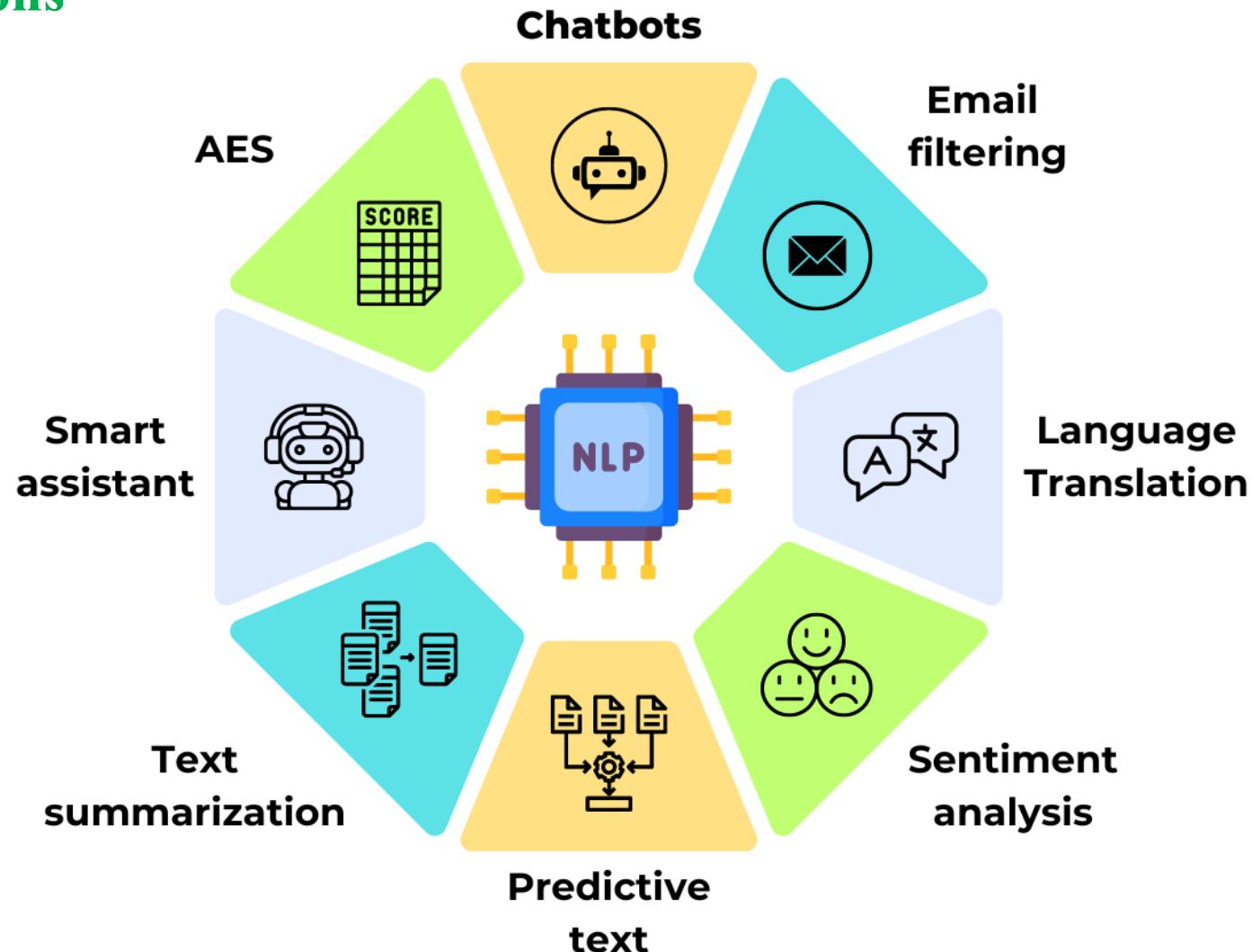
Introduction to NLP

❖ NLP tasks



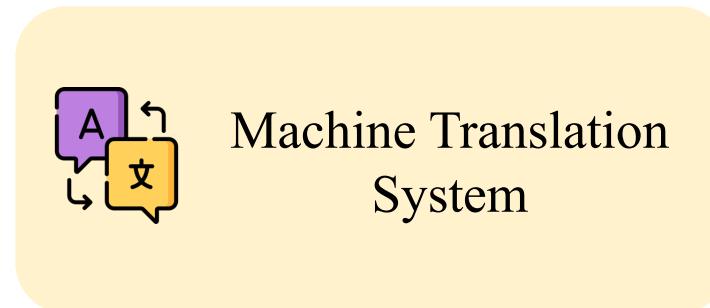
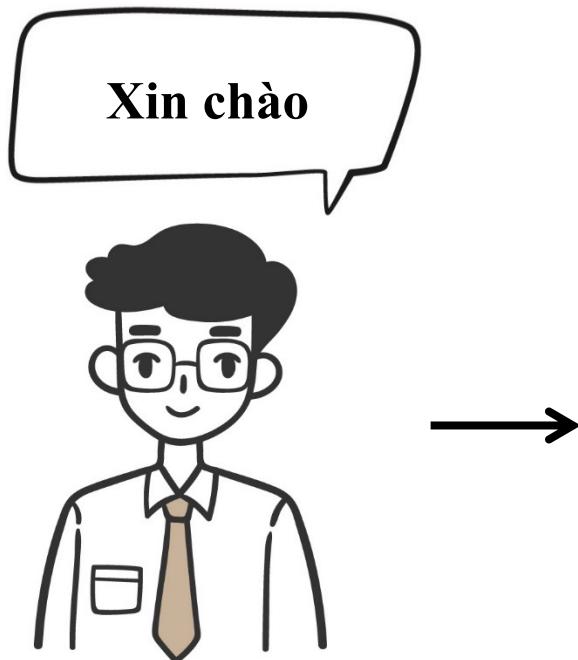
Introduction to NLP

❖ NLP Applications

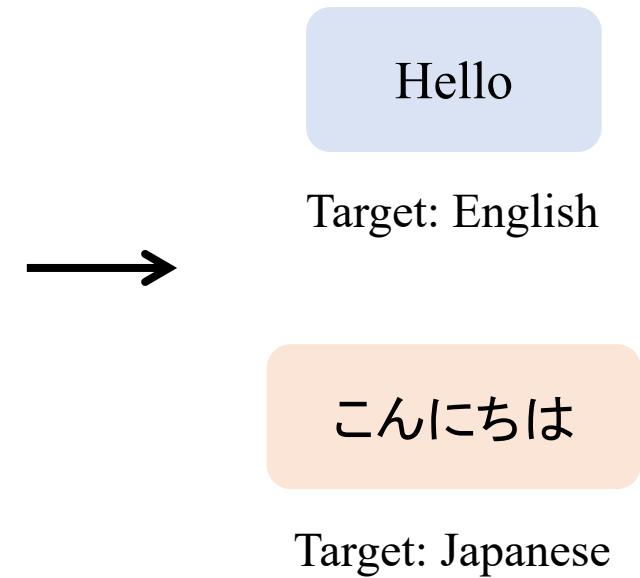


Introduction to NLP

❖ Translation

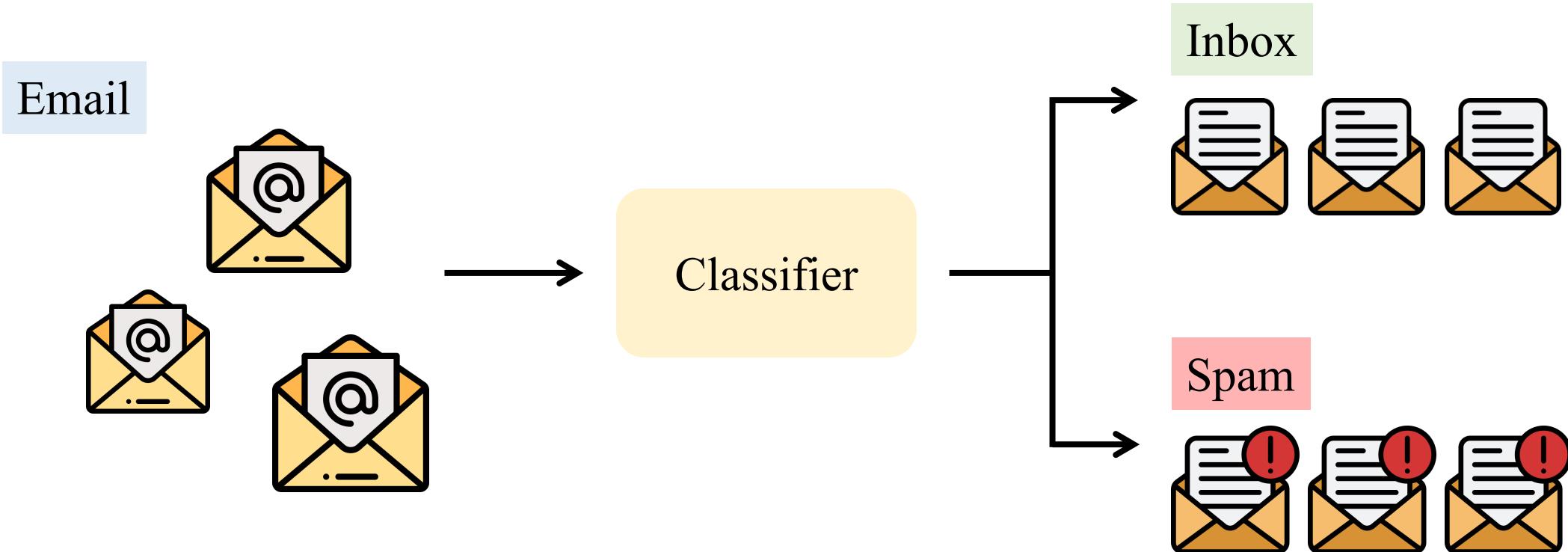


Source: Vietnamese



Introduction to NLP

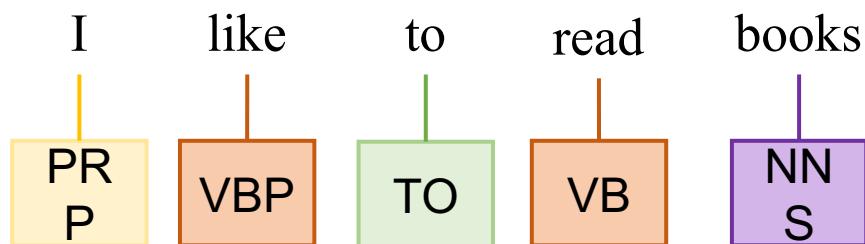
❖ Email filtering



Introduction to NLP

❖ Information Extraction

Part of Speech Tagging (POS Tagging)



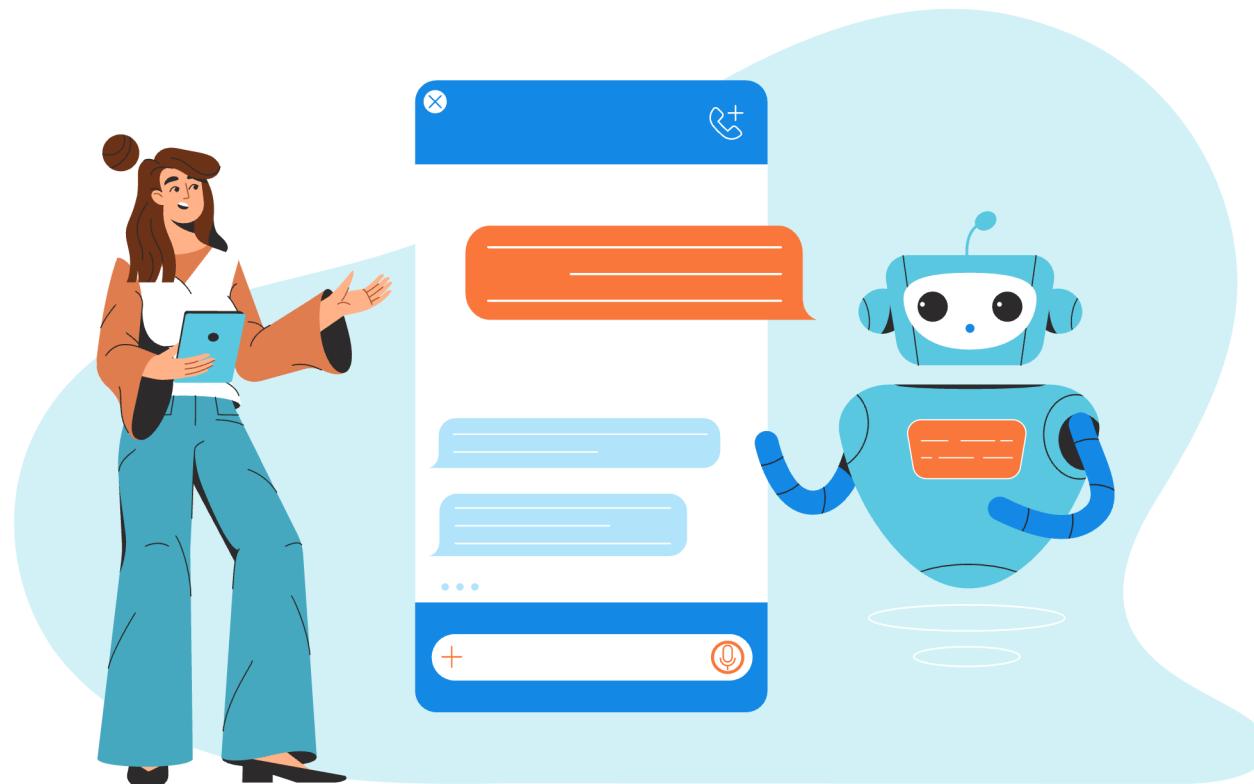
Named Entity Recognition (NER)

The diagram shows a text snippet with named entities highlighted in colored boxes. Above the text is a legend with categories: Person (p), Loc (l), Org (o), Event (e), Date (d), and Other (z). The text is: "Barack Hussein Obama II (born August 4, 1961) is an American attorney and politician who served as the 44th President of the United States from January 20, 2009, to January 20, 2017. A member of the Democratic Party, he was the first African American to serve as president. He was previously a United States Senator from Illinois and a member of the Illinois State Senate." The entities identified are: Barack Hussein Obama II (Person), August 4, 1961 (Date), American (Other), the United States (Loc), January 20, 2009 (Date), January 20, 2017 (Date), Democratic Party (Other), African American (Other), Illinois (Loc), and Illinois State Senate (Other).

Introduction to NLP

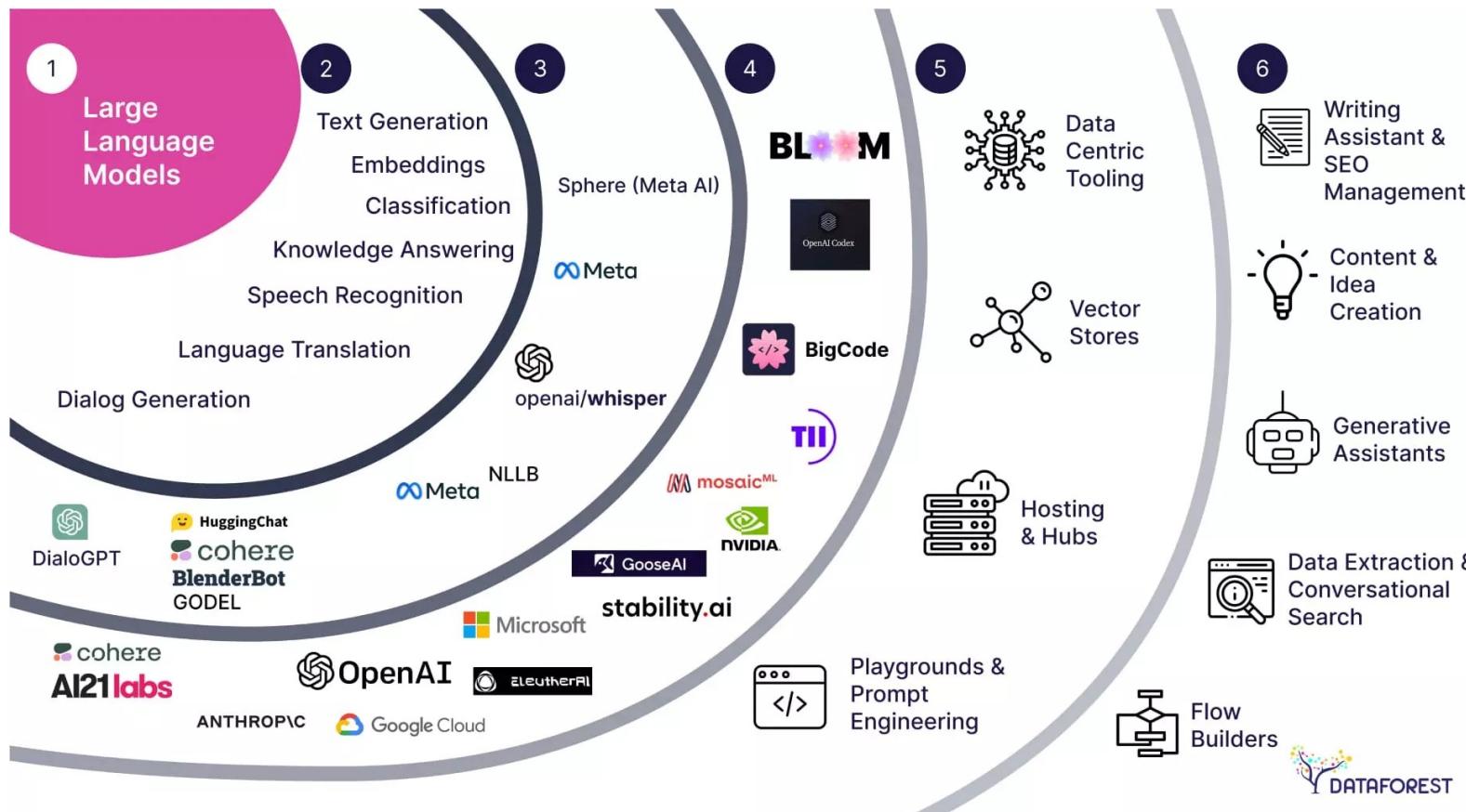
❖ Chatbot

Chatbot is a software application or web interface that is designed to mimic human conversation through text or voice interactions.



Introduction to NLP

❖ Large Language Models (LLMs)



A **large language model (LLM)** is a type of machine learning model designed for natural language processing tasks such as language generation. LLMs are language models with many parameters and are trained with self-supervised learning on a vast amount of text.

Introduction to NLP

❖ Large Language Models (LLMs)



GPT-4



LLaMA



Qwen



BERT



Claude



Introduction to NLP

❖ LLM: ChatGPT

Bao lâu nữa sẽ tới tết Nguyên Đán 2026?



Tết Nguyên Đán 2026 sẽ rơi vào ngày **17 tháng 2 năm 2026 (Thứ Ba)**.

Hôm nay là **4 tháng 2 năm 2025**, vậy còn khoảng **1 năm 13 ngày** nữa là đến Tết Nguyên Đán 2026. 🎉



Xử lý ngôn ngữ tự nhiên là gì?



Xử lý ngôn ngữ tự nhiên (Natural Language Processing - NLP) là gì?

Xử lý ngôn ngữ tự nhiên (NLP) là một lĩnh vực trong trí tuệ nhân tạo (AI) tập trung vào việc giúp máy tính hiểu, diễn giải, và tạo ra ngôn ngữ con người một cách tự động. NLP kết hợp các kỹ thuật từ **khoa học máy tính, ngôn ngữ học, và học máy** để xử lý văn bản hoặc giọng nói theo cách tự nhiên nhất có thể.

Introduction to NLP

❖ LLM: DeepSeek

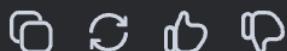
Bao lâu nữa sẽ tới Tết Nguyên Đán 2026?



Để tính khoảng thời gian từ hiện tại đến Tết Nguyên Đán 2026, chúng ta cần biết ngày hiện tại. Tuy nhiên, nếu giả sử hôm nay là ngày 10 tháng 10 năm 2023, thì Tết Nguyên Đán 2026 sẽ rơi vào ngày 29 tháng 1 năm 2026.

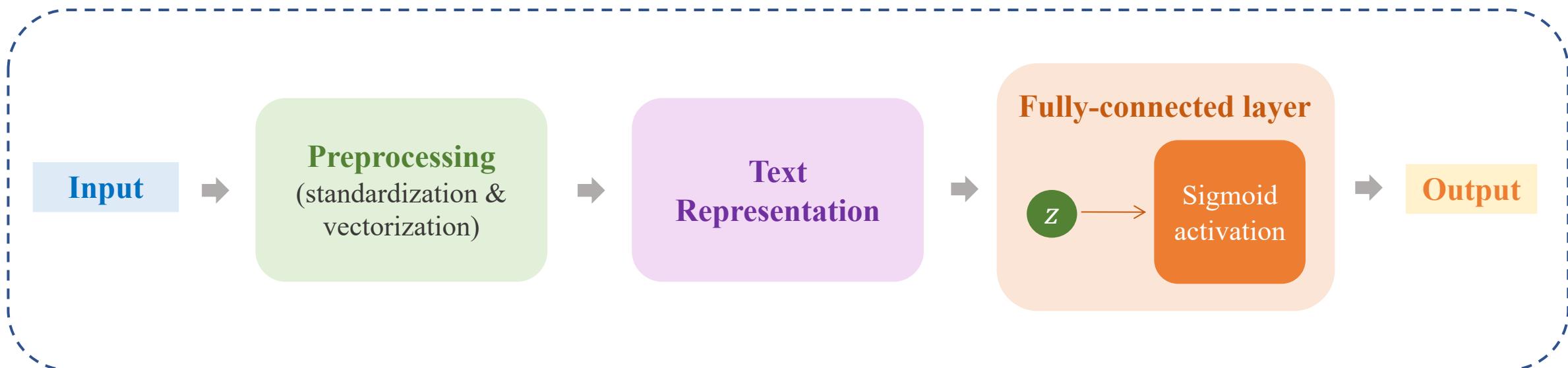
Vậy, từ ngày 10 tháng 10 năm 2023 đến ngày 29 tháng 1 năm 2026 là khoảng **2 năm, 3 tháng và 19 ngày**.

Nếu bạn muốn biết chính xác hơn, bạn có thể tính toán dựa trên ngày hiện tại của bạn.



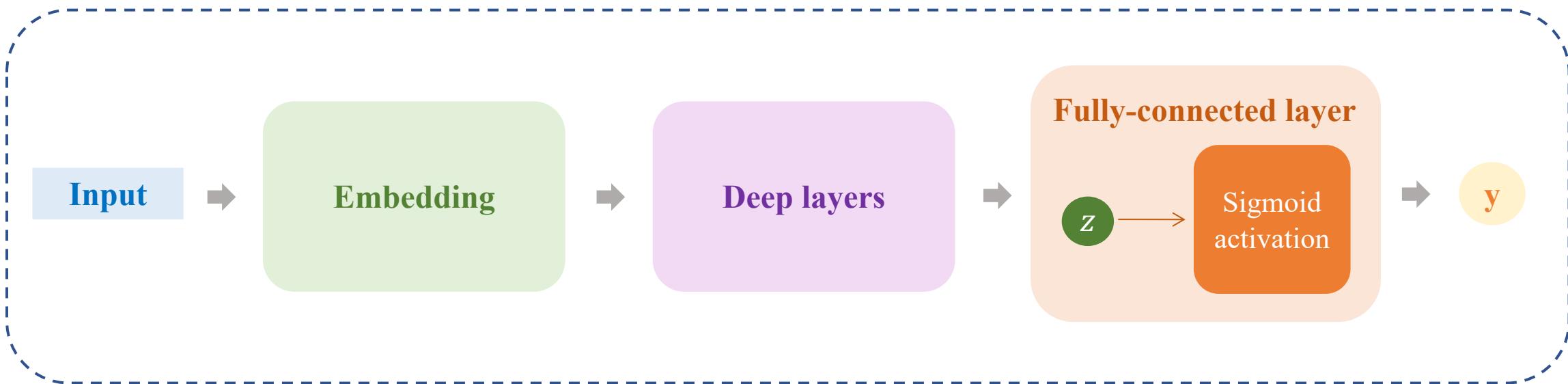
Introduction to NLP

❖ Basic NLP pipeline



Text Representation

❖ Sentence Classification



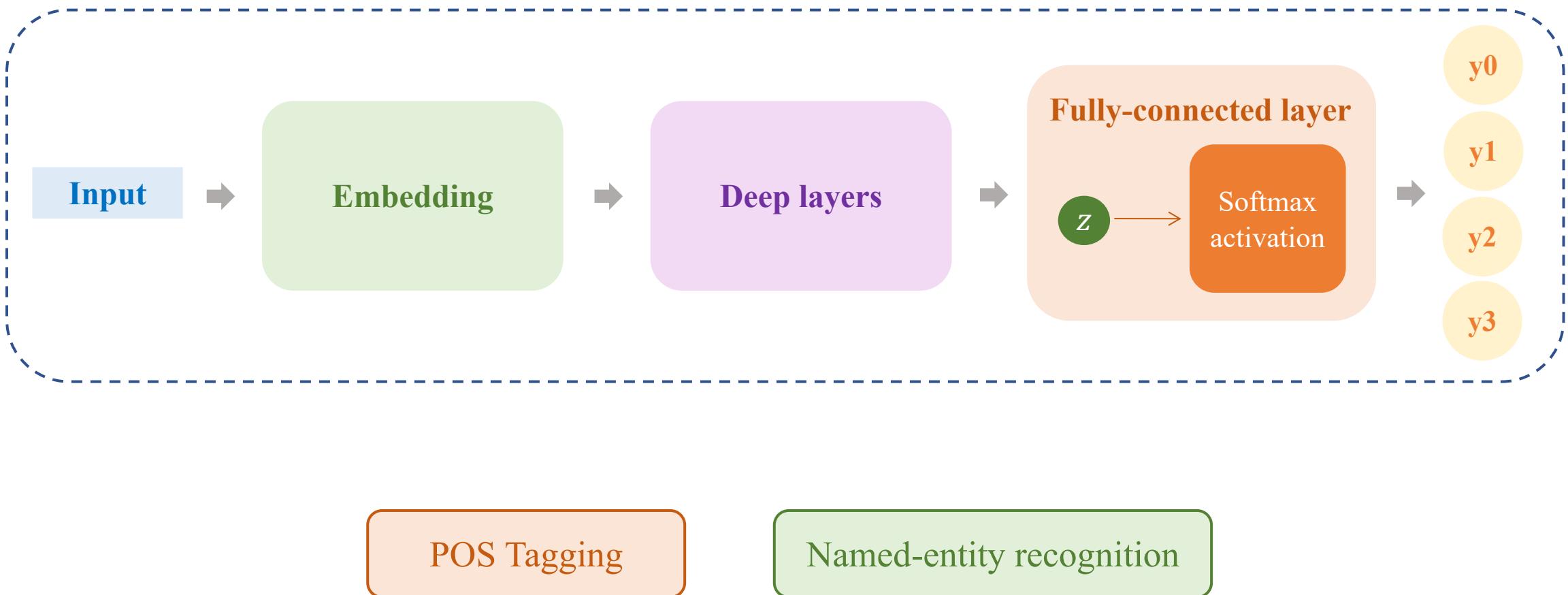
Text Classification

Sentiment Analysis

Spam/ Not Spam

Text Representation

❖ Token Classification



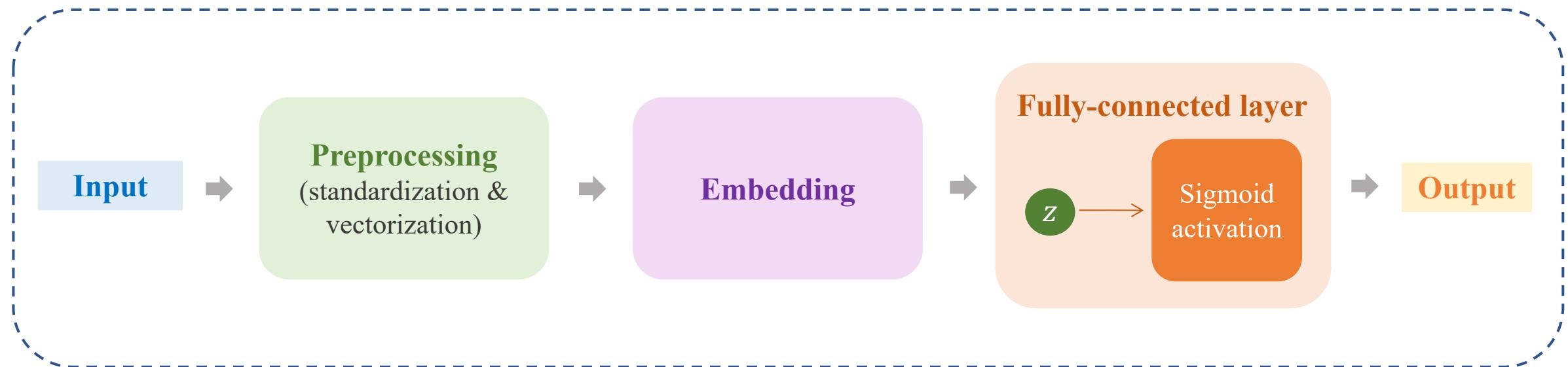
Text Classification

Text Classification

- Preprocessing
- Embedding
- Text Classification

Text Classification

❖ Simple approach



Preprocessing

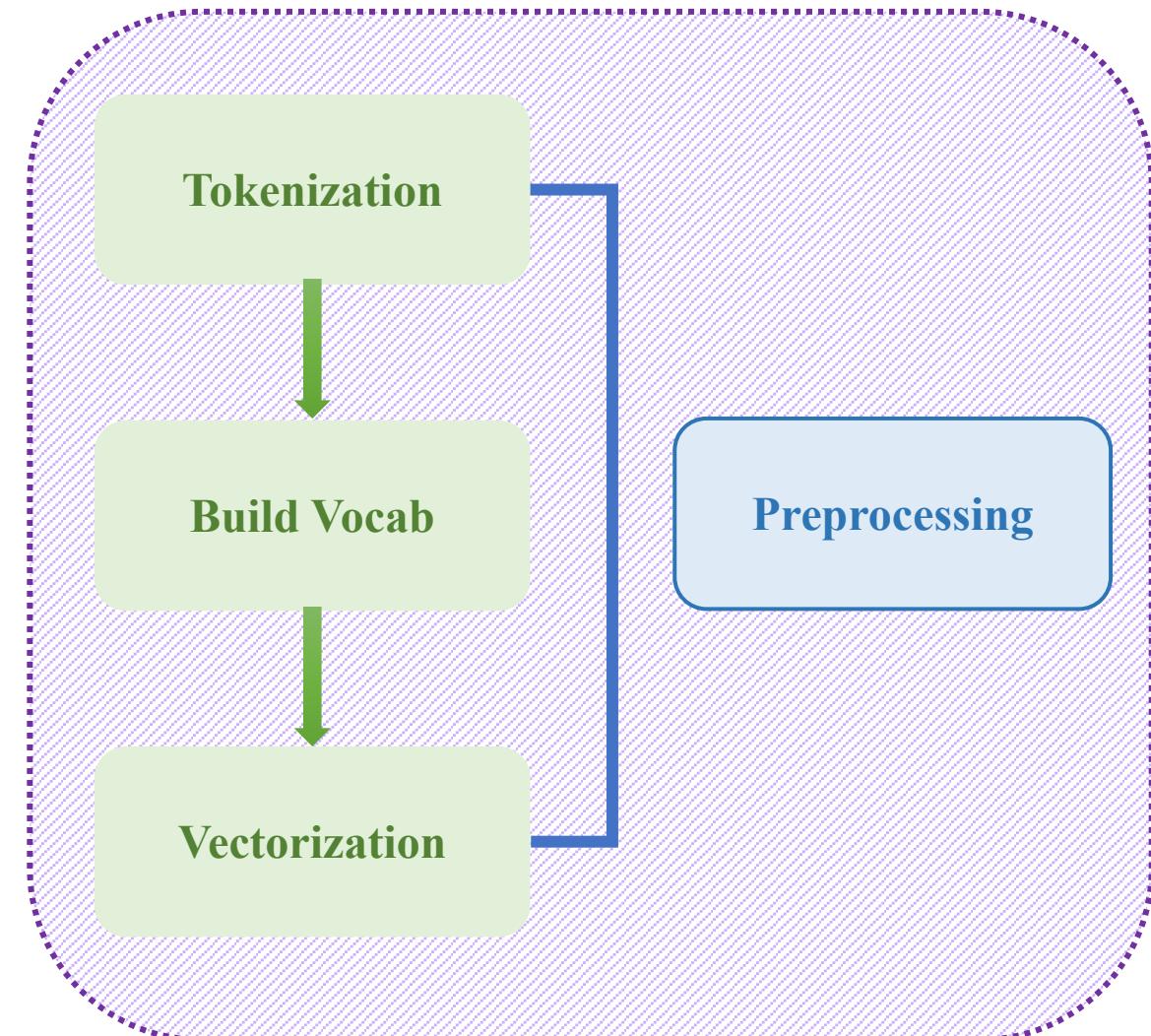
❖ Text Corpus

Deepseek sell off stock market

Deepseek release all their
technical documents

```
corpus = [  
    "Deepseek sell off stock market",  
    "Deepseek release all their technical documents"  
]
```

```
# 0: negative - 1: positive  
labels = [0, 1]
```



Tokenization

Example corpus

sample1: ‘Deepseek sell off stock market’

sample2: ‘Deepseek release all their
technical documents’

Deepseek sell off stock market

Deepseek release all their
technical documents

↓ Tokenize

Deepseek | sell | off | stock

market

Deepseek | release | all | their

technical

documents

```
from torchtext.data.utils import get_tokenizer
# version 0.17.0

sample1 = corpus[0]
sample2 = corpus[1]

#Define tokenizer function
tokenizer = get_tokenizer('basic_english')
sample1_tokens = tokenizer(sample1)
sample2_tokens = tokenizer(sample2)

print(sample1_tokens)
print(sample2_tokens)

['deepseek', 'sell', 'off', 'stock', 'market']
['deepseek', 'release', 'all', 'their', 'technical', 'documents']
```

Building Vocab

```
from torchtext.data.utils import get_tokenizer
from torchtext.vocab import build_vocab_from_iterator

# Define tokenizer function
tokenizer = get_tokenizer('basic_english')

# Create a function to yield list of tokens
def yield_tokens(examples):
    for text in examples:
        yield tokenizer(text)

# Create vocabulary
vocab = build_vocab_from_iterator(yield_tokens(corpus),
                                    max_tokens=vocab_size,
                                    specials=["<unk>", "<pad>"])
vocab.set_default_index(vocab["<unk>"])

    vocab.get_stoi()
{'release': 7,
 'off': 6,
 'market': 5,
 'documents': 4,
 'deepseek': 2,
 '<pad>': 1,
 'all': 3,
 '<unk>': 0}
```

0	<unk>
1	<pad>
2	deepseek
3	all
4	documents
5	market
6	off
7	release

Vocabulary

Vectorization

0	<unk>
1	<pad>
2	deepseek
3	all
4	documents
5	market
6	off
7	release

Vocabulary

```
sample1_ids = [vocab[token] for token in sample1_tokens]  
sample2_ids = [vocab[token] for token in sample2_tokens]
```

```
print(sample1_ids)  
print(sample2_ids)
```

[2, 0, 6, 0, 5]
[2, 7, 3, 0, 0, 4]

Deepseek sell off stock market



Vectorize

2 0 6 0 5

Problem?

Deepseek release all their technical documents

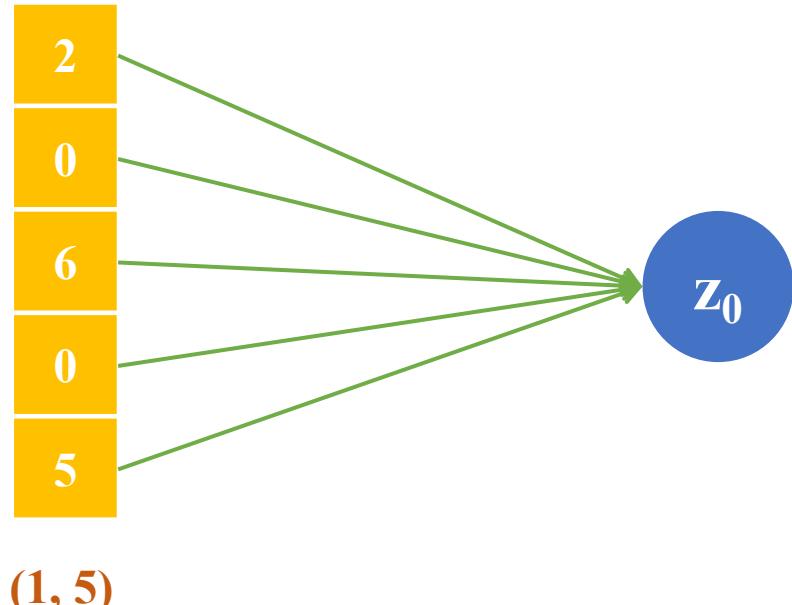


Vectorize

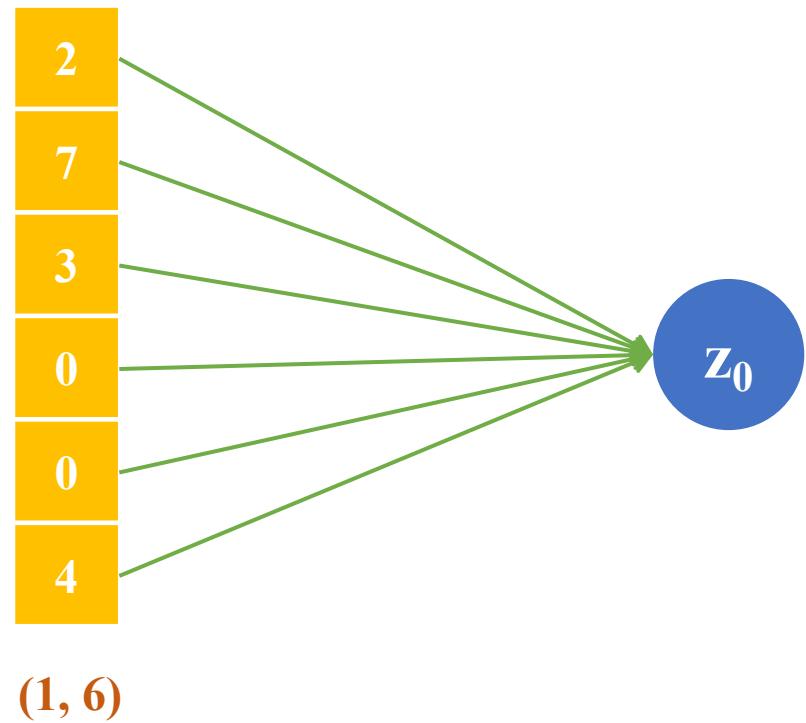
2 7 3 0 0 4

Vectorization

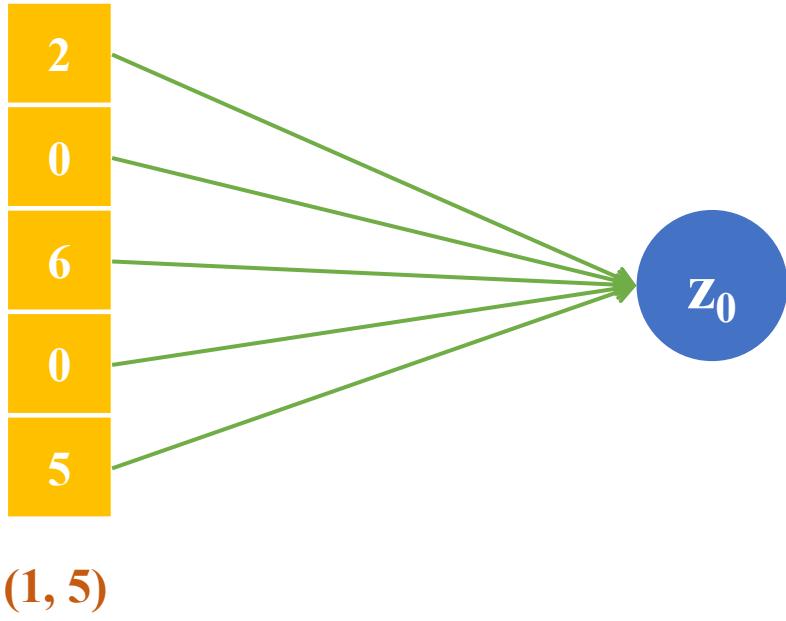
$$Z = X \cdot W^T + b$$



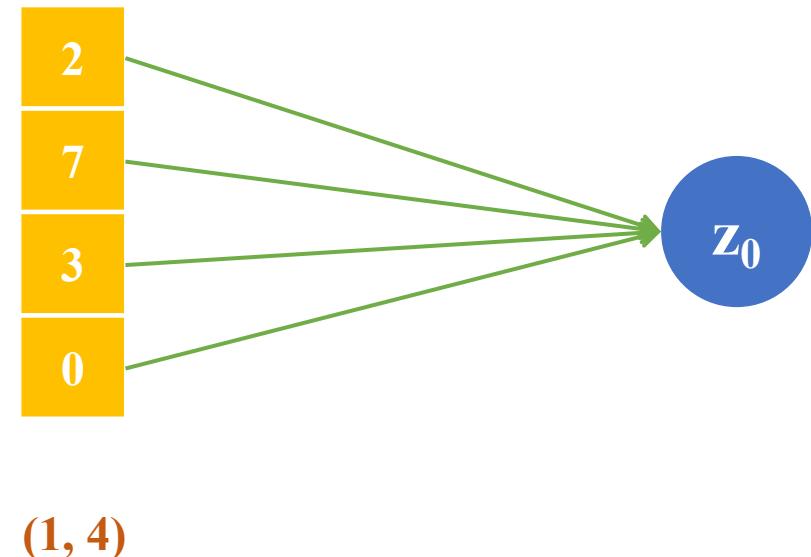
Shape W ?



Vectorization



Shape W ?



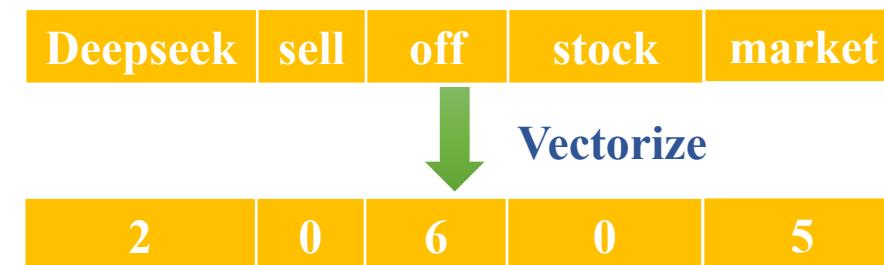
Vectorization

0	<unk>
1	<pad>
2	deepseek
3	all
4	documents
5	market
6	off
7	release

Vocabulary

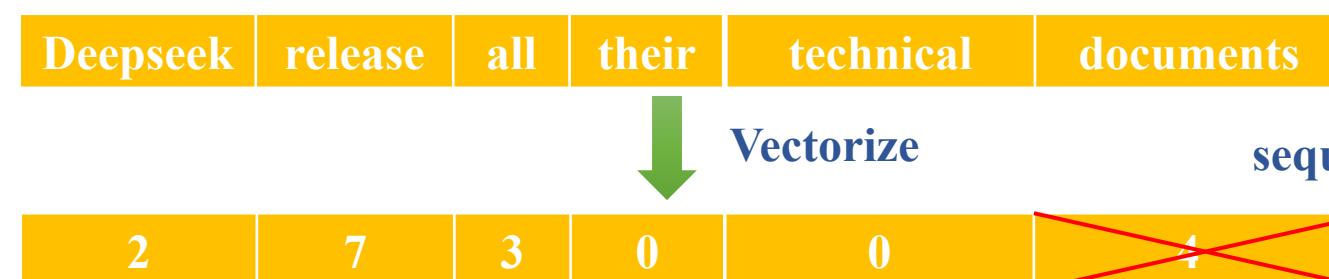
```
# Tokenize and numericalize your samples
def vectorize(text, vocab, sequence_length):
    tokens = tokenizer(text)
    token_ids = [vocab[token] for token in tokens][:sequence_length]
    token_ids = token_ids + [vocab["<pad>"]] * (sequence_length - len(tokens))
    return torch.tensor(token_ids, dtype=torch.long)

# Vectorize the samples
corpus_ids = []
for sentence in corpus:
    corpus_ids.append(vectorize(sentence, vocab, sequence_length))
```



```
for v in corpus_ids:
    print(v)
```

```
tensor([2, 0, 6, 0, 5])
tensor([2, 7, 3, 0, 0])
```



sequence_length = 5

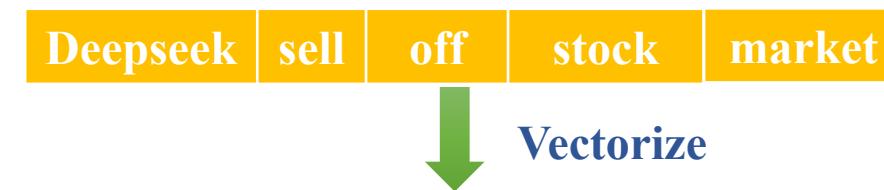
Vectorization

0	<unk>
1	<pad>
2	deepseek
3	all
4	documents
5	market
6	off
7	release

Vocabulary

```
# Tokenize and numericalize your samples
def vectorize(text, vocab, sequence_length):
    tokens = tokenizer(text)
    token_ids = [vocab[token] for token in tokens][:sequence_length]
    token_ids = token_ids + [vocab["<pad>"]] * (sequence_length - len(tokens))
    return torch.tensor(token_ids, dtype=torch.long)

# Vectorize the samples
corpus_ids = []
for sentence in corpus:
    corpus_ids.append(vectorize(sentence, vocab, sequence_length))
```



```
for v in corpus_ids:
    print(v)
```

```
tensor([2, 0, 6, 0, 5])
tensor([2, 7, 3, 0, 0])
```



sequence_length = 5

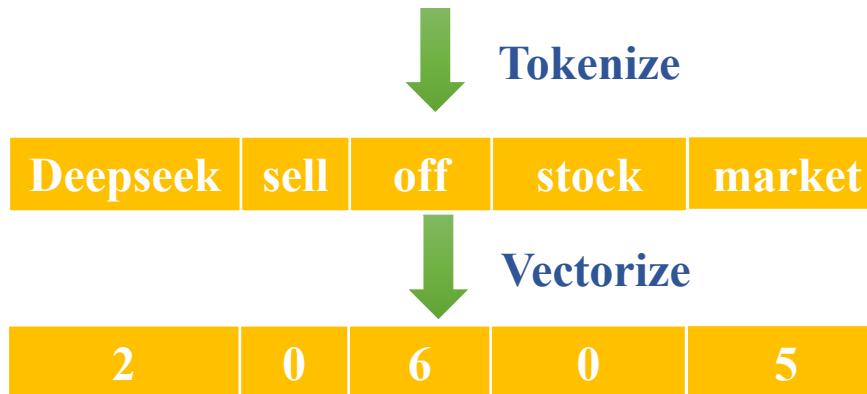
Preprocessing

Example corpus

sample1: 'Deepseek sell off stock market'

sample2: 'Deepseek release all their
technical documents'

Deepseek sell off stock market



Vocabulary

0	1	2	3	4	5	6	7
<unk>	<pad>	deepseek	all	documents	market	off	release

Deepseel release all their technical documents



0	1	2	3	4	5	6	7
<unk>	<pad>	deepseek	all	documents	market	off	release

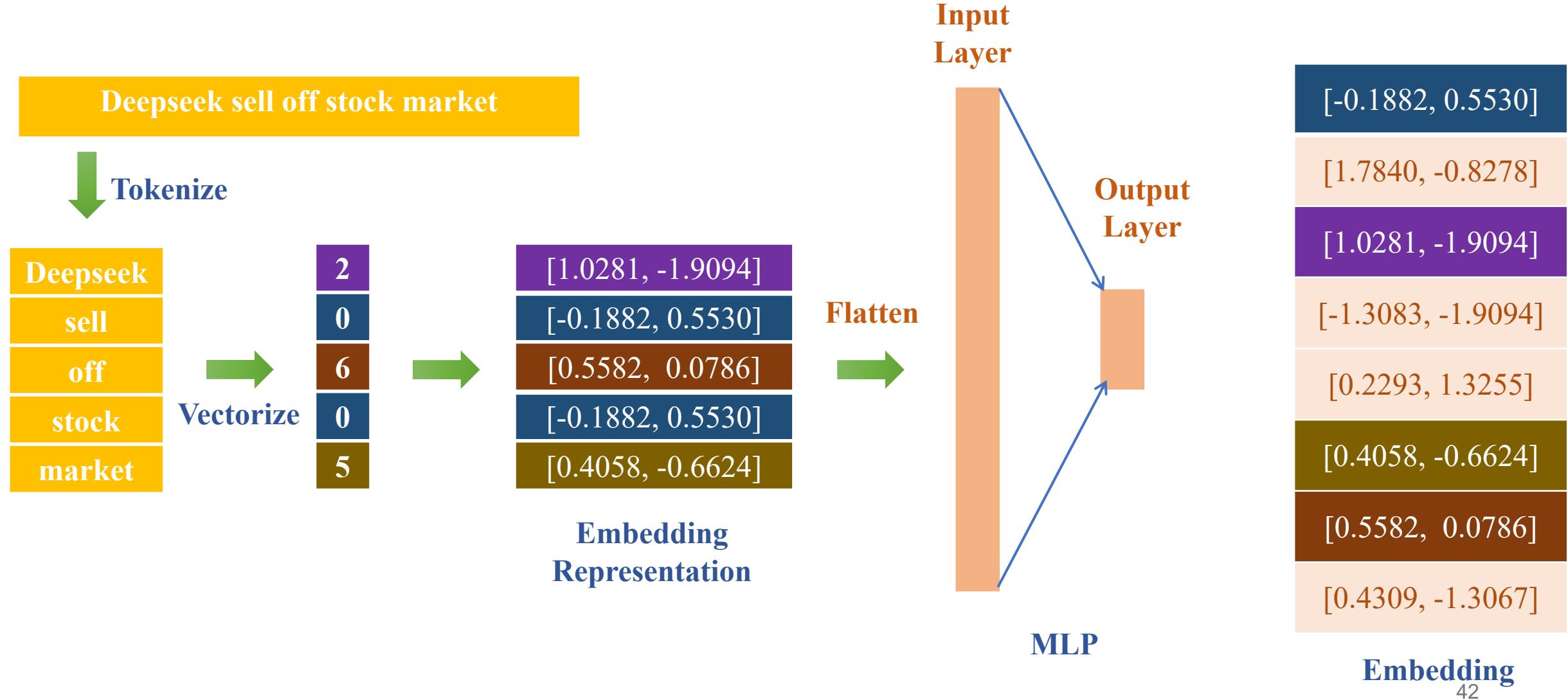
Embedding

Index	Token	Embedding_dims = 2
0	<unk>	 [-0.1882, 0.5530]
1	<pad>	 [1.7840, -0.8278]
2	deepseek	 [1.0281, -1.9094]
3	all	 [-1.3083, -1.9094]
4	documents	 [0.2293, 1.3255]
5	market	 [0.4058, -0.6624]
6	off	 [0.5582, 0.0786]
7	release	 [0.4309, -1.3067]

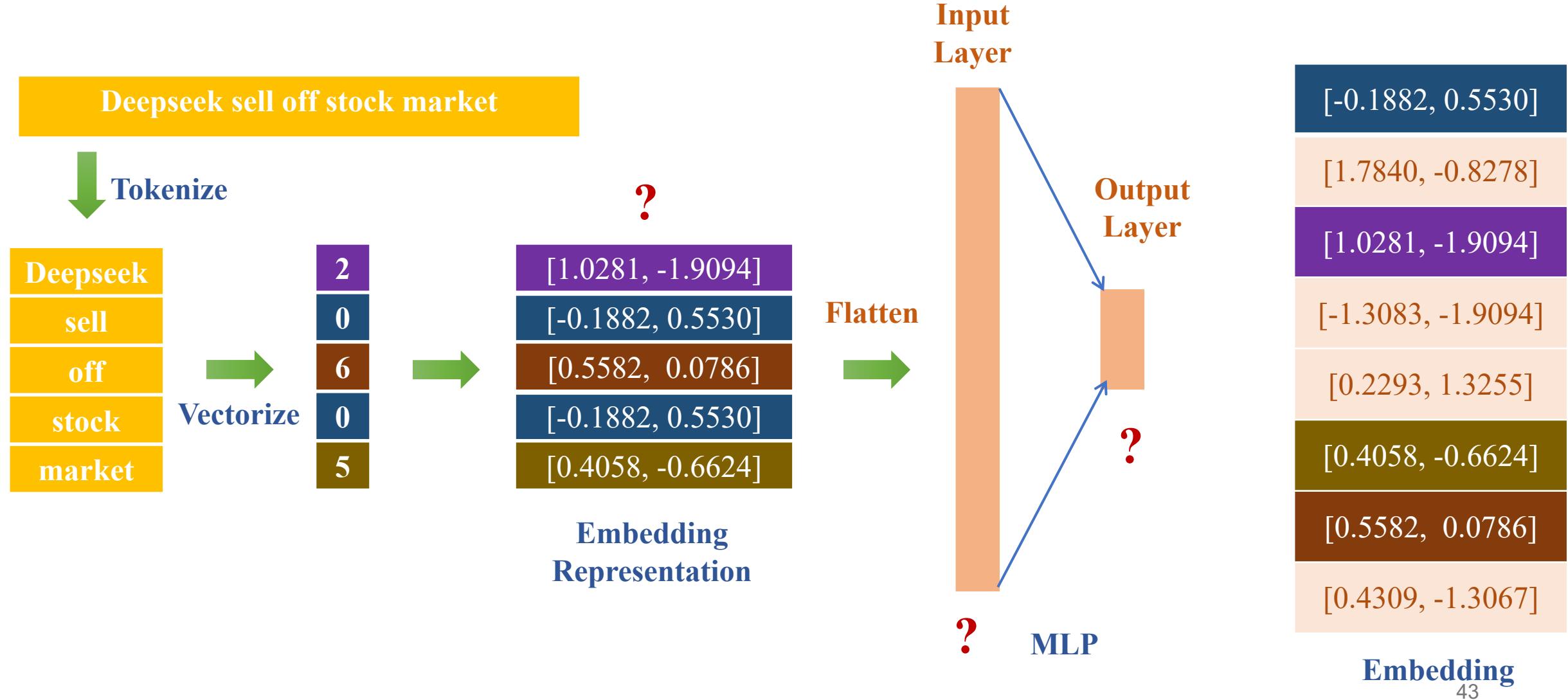
```
vocab_size = 8
embedding_dim = 2
embedding = nn.Embedding(vocab_size, embedding_dim)
custom_weights = torch.tensor( [-0.1882,  0.5530],
                               [ 1.7840, -0.8278],
                               [ 1.0281, -1.9094],
                               [-1.3083, -0.0987],
                               [ 0.2293,  1.3255],
                               [ 0.4058, -0.6624],
                               [ 0.5582,  0.0786],
                               [ 0.4309, -1.3067]).float()
embedding.weight = nn.Parameter(custom_weights)
```

Initialization

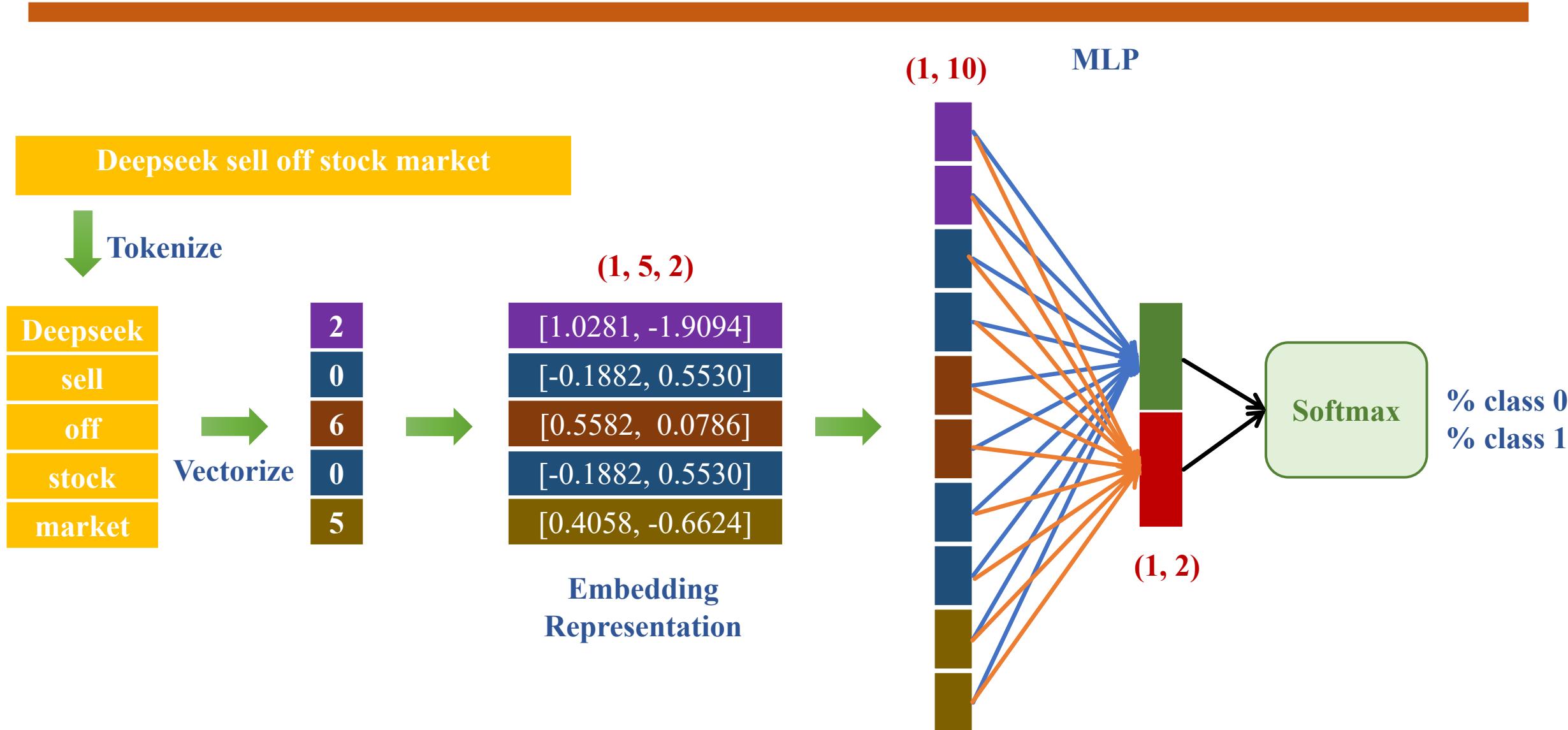
Text Classification



Text Classification



Text Classification



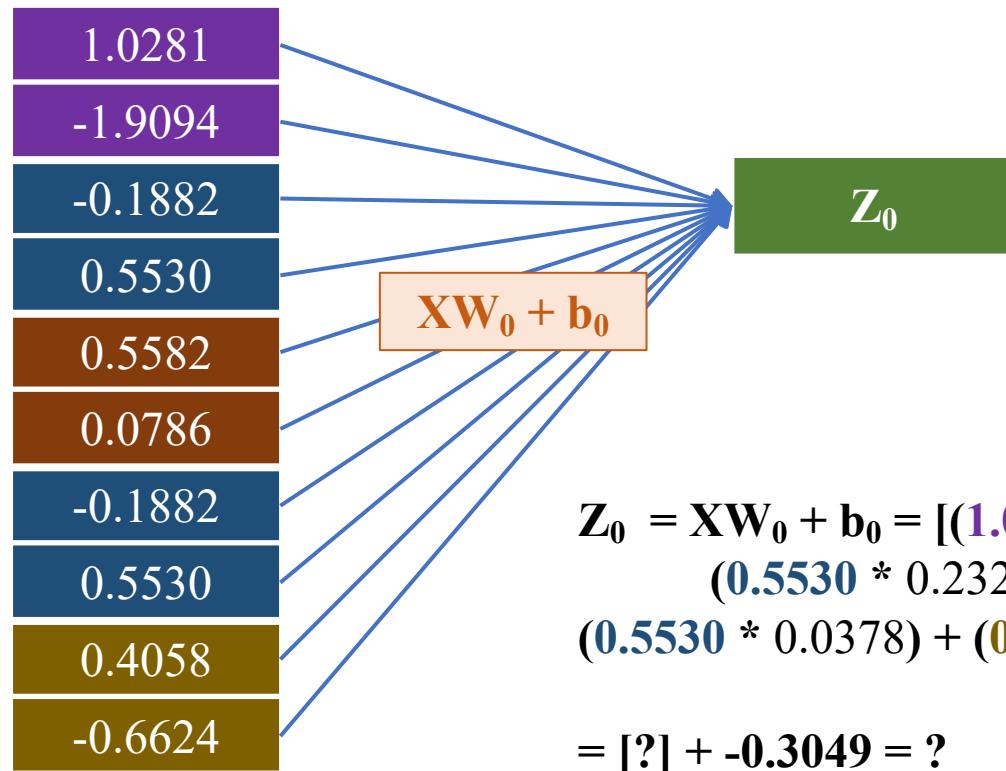
FC Forward

$$W = \begin{bmatrix} [0.2108, -0.0074, 0.2760, 0.2325, -0.0518, -0.1876, 0.0194, 0.0378, 0.0210, 0.2982] \\ [0.0284, 0.2968, -0.0260, 0.1251, -0.0282, 0.0175, -0.1817, 0.2483, 0.2338, 0.2985] \end{bmatrix}$$

$$w_0 = \begin{bmatrix} [-0.3049, b_0] \\ 0.1028 \end{bmatrix}$$

$$w_1 = \begin{bmatrix} b_1 \end{bmatrix}$$

X



$$\begin{aligned} Z_0 &= XW_0 + b_0 = [(1.0281 * 0.2108) + (-1.9094 * -0.0074) + (-0.1882 * 0.2760) + \\ &\quad (0.5530 * 0.2325) + (0.5582 * -0.0518) + (0.0786 * -0.1876) + (-0.1882 * 0.0194) + \\ &\quad (0.5530 * 0.0378) + (0.4058 * 0.0210) + (-0.6624 * 0.2982)] + -0.3049 \\ &= [?] + -0.3049 = ? \end{aligned}$$

FC Forward

$$W = [0.2108, -0.0074, 0.2760, 0.2325, -0.0518, -0.1876, 0.0194, 0.0378, 0.0210, 0.2982] \quad w_0 \\ [0.0284, 0.2968, -0.0260, 0.1251, -0.0282, 0.0175, -0.1817, 0.2483, 0.2338, 0.2985] \quad w_1$$

$$b = [-0.3049, b_0 \\ 0.1028] \quad b_1$$

X

1.0281
-1.9094
-0.1882
0.5530
0.5582
0.0786
-0.1882
0.5530
0.4058
-0.6624

$$X \cdot W_0 = \begin{matrix} 0.2108 \\ -0.0074 \\ 0.2760 \\ 0.2325 \\ -0.0518 \\ -0.1876 \\ 0.0194 \\ 0.0378 \\ 0.0210 \\ 0.2982 \end{matrix}$$

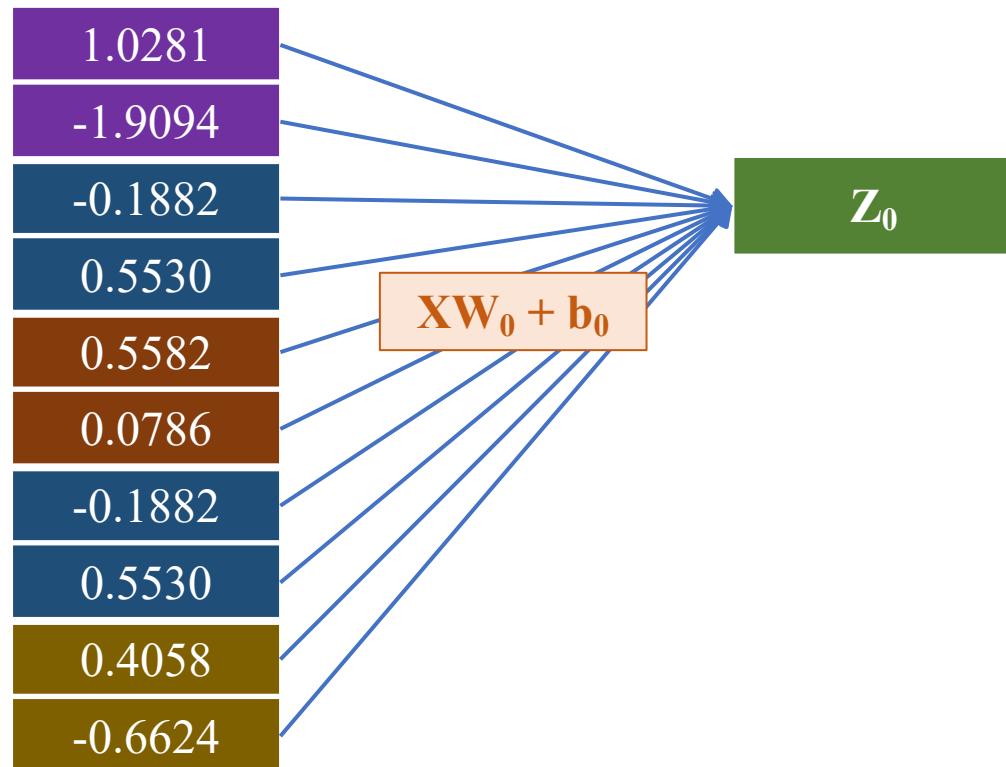
$$= \begin{matrix} 0.2167 \\ 0.0141 \\ -0.0519 \\ 0.1286 \\ -0.0289 \\ -0.0147 \\ -0.0037 \\ 0.0209 \\ 0.0085 \\ -0.1975 \end{matrix}$$

$\text{Sum} = 0.0921 + -0.3049 = Z_0$
 $XW_0 + b_0 = Z_0$

FC Forward

$$\mathbf{W} = \begin{bmatrix} [0.2108, -0.0074, 0.2760, 0.2325, -0.0518, -0.1876, 0.0194, 0.0378, 0.0210, 0.2982] \\ [0.0284, 0.2968, -0.0260, 0.1251, -0.0282, 0.0175, -0.1817, 0.2483, 0.2338, 0.2985] \end{bmatrix}$$
$$w_0 \quad w_1$$
$$\mathbf{b} = \begin{bmatrix} [-0.3049, b_0] \\ 0.1028] \end{bmatrix} \quad b_1$$

X



$$Z_0 = \mathbf{X}\mathbf{W}_0 + \mathbf{b}_0 = 0.9210 + -0.3049 = -0.2128$$

FC Forward

$$W = \begin{bmatrix} [0.2108, -0.0074, 0.2760, 0.2325, -0.0518, -0.1876, 0.0194, 0.0378, 0.0210, 0.2982] \\ [0.0284, 0.2968, -0.0260, 0.1251, -0.0282, 0.0175, -0.1817, 0.2483, 0.2338, 0.2985] \end{bmatrix}$$

$$w_0 \quad w_1$$

$$b = \begin{bmatrix} [-0.3049, b_0] \\ [0.1028] b_1 \end{bmatrix}$$

X

1.0281
-1.9094
-0.1882
0.5530
0.5582
0.0786
-0.1882
0.5530
0.4058
-0.6624

$$Z_1 = XW_1 + b_1 = [(1.0281 * 0.0284) + (-1.9094 * 0.2968) + (-0.1882 * -0.0260) + (0.5530 * 0.1251) + (0.5582 * -0.0282) + (0.0786 * 0.0175) + (-0.1882 * -0.1817) + (0.5530 * 0.2483) + (0.4058 * 0.2338) + (-0.6624 * 0.2985)] + 0.1028$$

$$= [?] + 0.1028 = ?$$

$XW_1 + b_1$

Z_1

FC Forward

$$W = \begin{bmatrix} [0.2108, -0.0074, 0.2760, 0.2325, -0.0518, -0.1876, 0.0194, 0.0378, 0.0210, 0.2982] & w_0 \\ [0.0284, 0.2968, -0.0260, 0.1251, -0.0282, 0.0175, -0.1817, 0.2483, 0.2338, 0.2985] & w_1 \end{bmatrix}$$

$$b = \begin{bmatrix} [-0.3049, b_0] \\ 0.1028] b_1 \end{bmatrix}$$

X

$$X = \begin{bmatrix} 1.0281 \\ -1.9094 \\ -0.1882 \\ 0.5530 \\ 0.5582 \\ 0.0786 \\ -0.1882 \\ 0.5530 \\ 0.4058 \\ -0.6624 \end{bmatrix}$$

$$W_1 = \begin{bmatrix} 0.0284 \\ 0.2968 \\ -0.0260 \\ 0.1251 \\ -0.0282 \\ 0.0175 \\ -0.1817 \\ 0.2483 \\ 0.2338 \\ 0.2985 \end{bmatrix}$$

=

$$\text{Sum} = -0.4091 + 0.1028 = Z_1$$

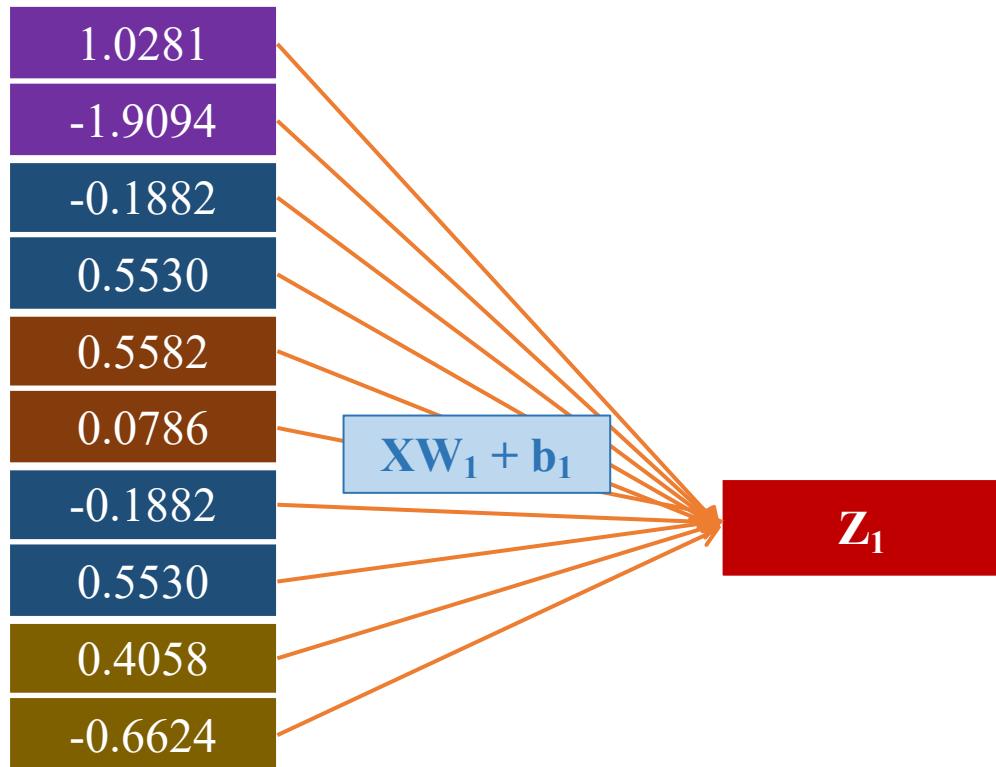
XW_1

b_1

FC Forward

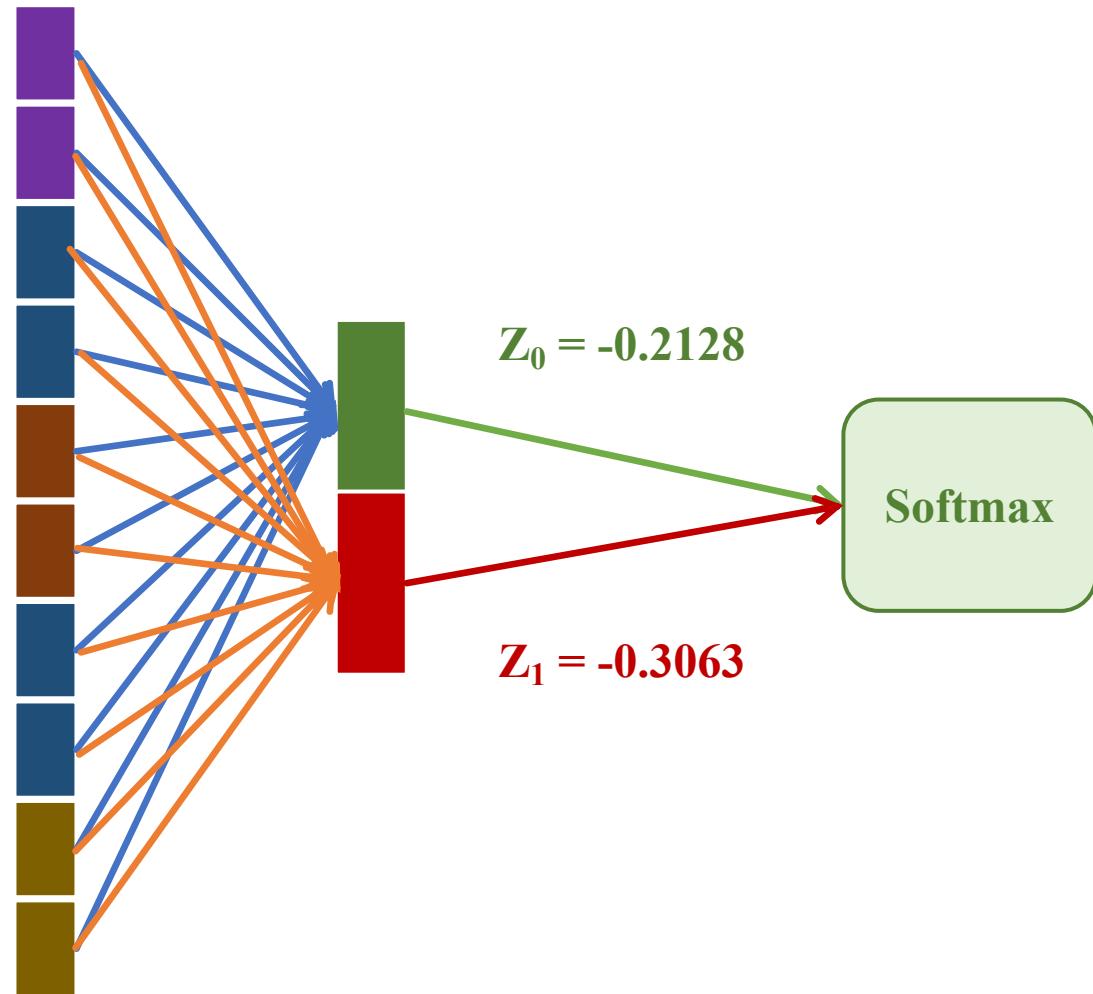
$$W = \begin{bmatrix} [0.2108, -0.0074, 0.2760, 0.2325, -0.0518, -0.1876, 0.0194, 0.0378, 0.0210, 0.2982] \\ [0.0284, 0.2968, -0.0260, 0.1251, -0.0282, 0.0175, -0.1817, 0.2483, 0.2338, 0.2985] \end{bmatrix}$$
$$w_0$$
$$w_1$$
$$b = [-0.3049, 0.1028]$$
$$b_0$$
$$b_1$$

X

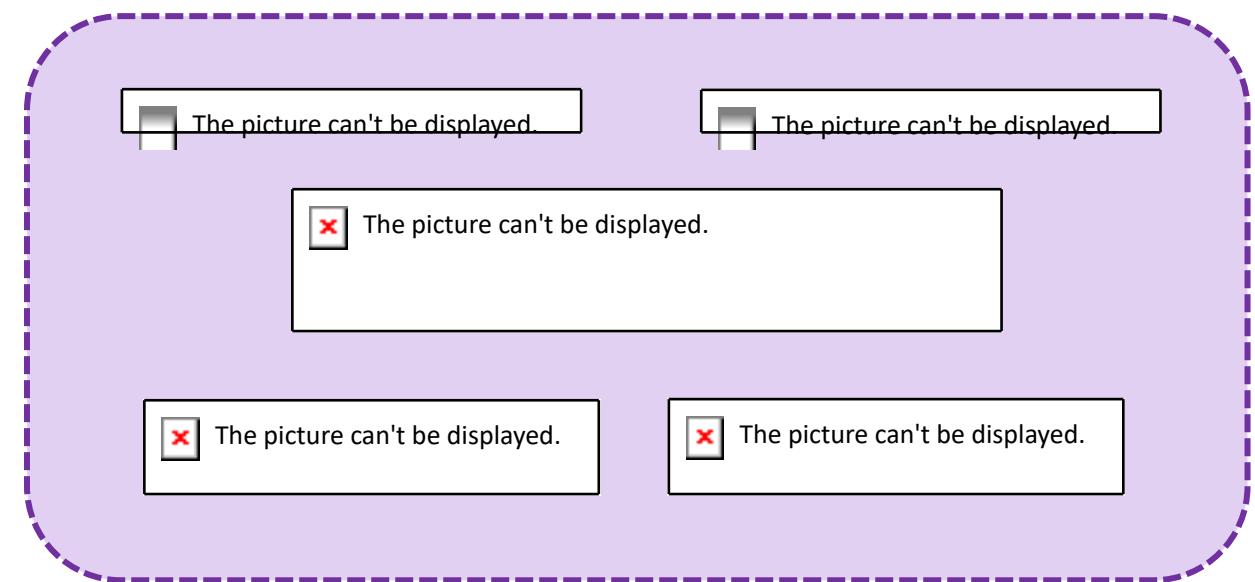


$$Z_1 = XW_1 + b_1 = -0.4091 + 0.1028 = -0.3063$$

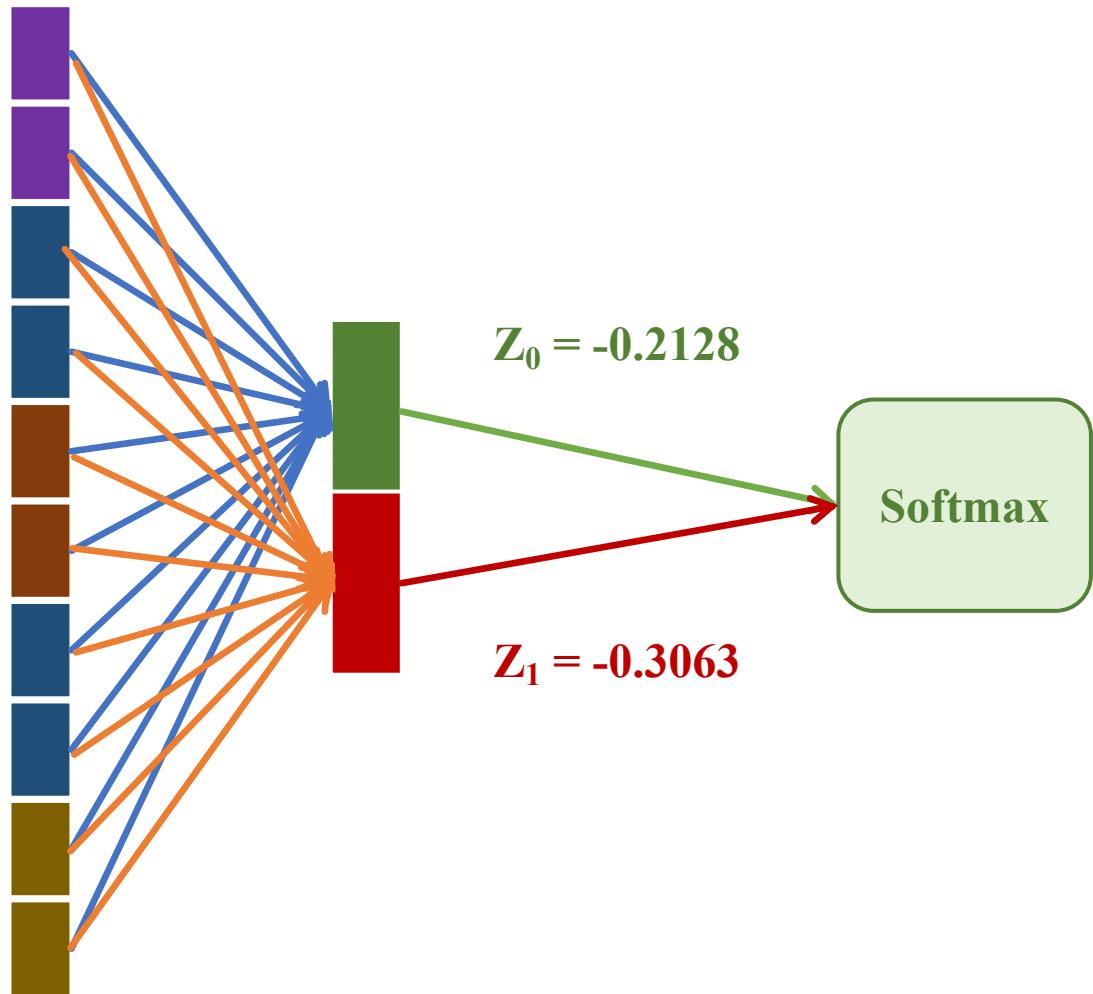
FC Forward



$$S(Z)_i = \frac{e^{Z_i}}{\sum_{j=1}^k e^{Z_j}}$$



FC Forward



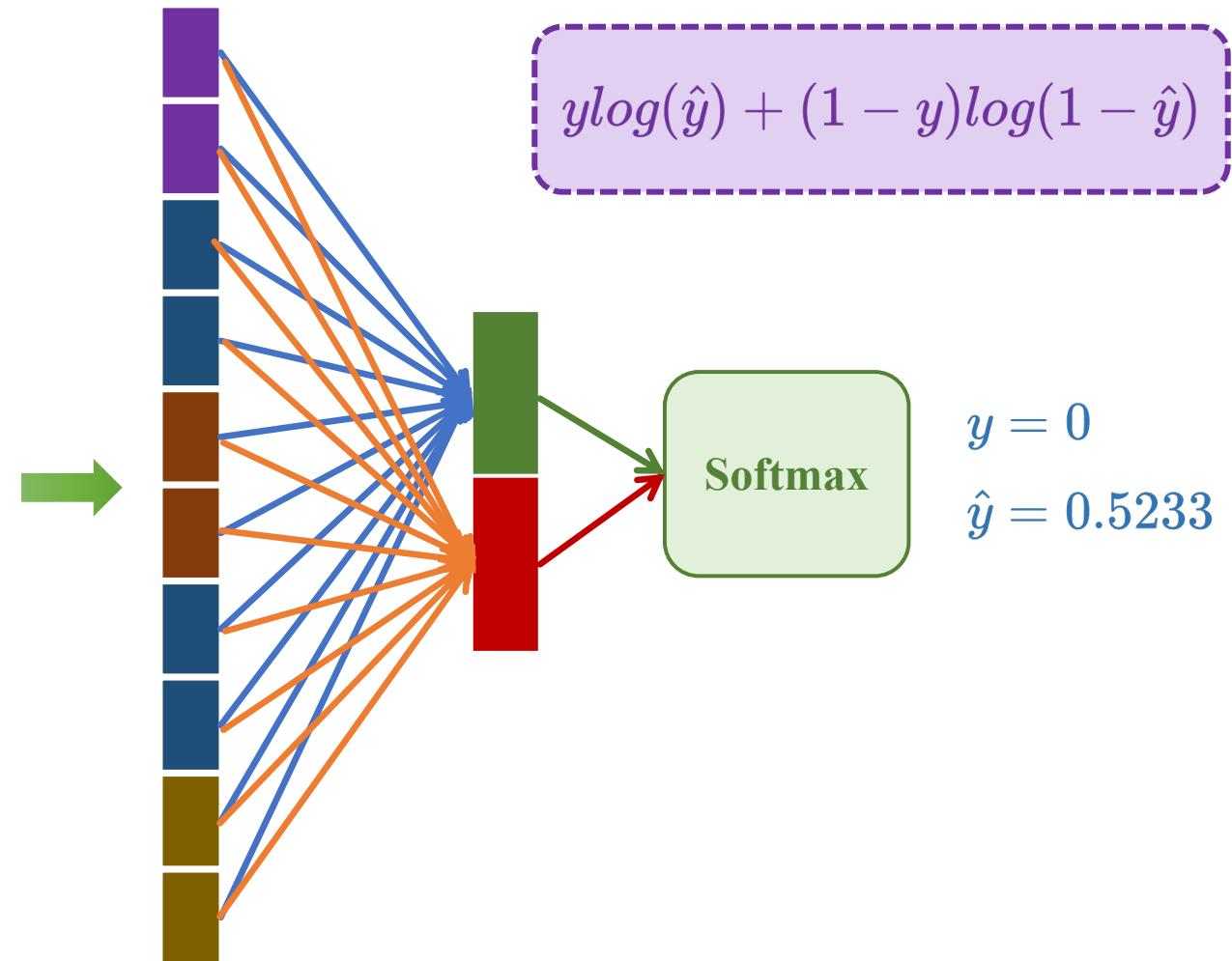
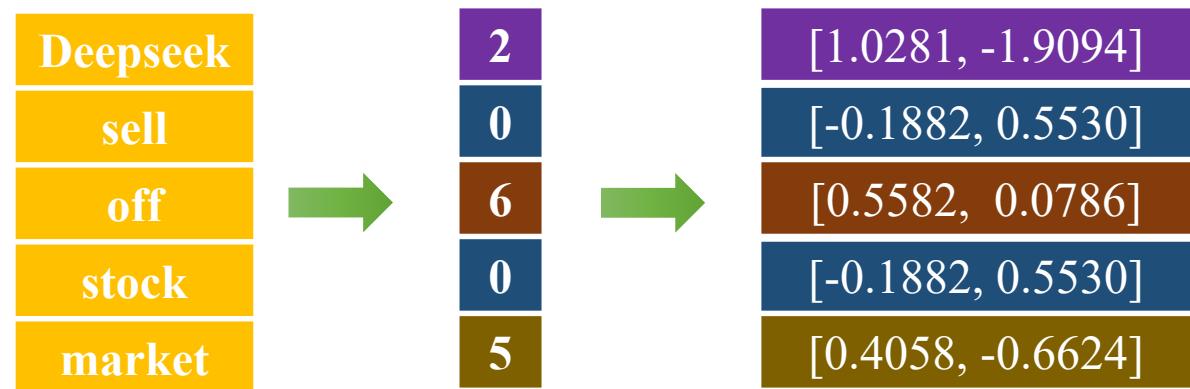
$$S(Z)_i = \frac{e^{Z_i}}{\sum_{j=1}^k e^{Z_j}}$$

$$e^{Z_0} = e^{-0.2128} = 0.8083 \quad e^{Z_1} = e^{-0.3063} = 0.7362$$

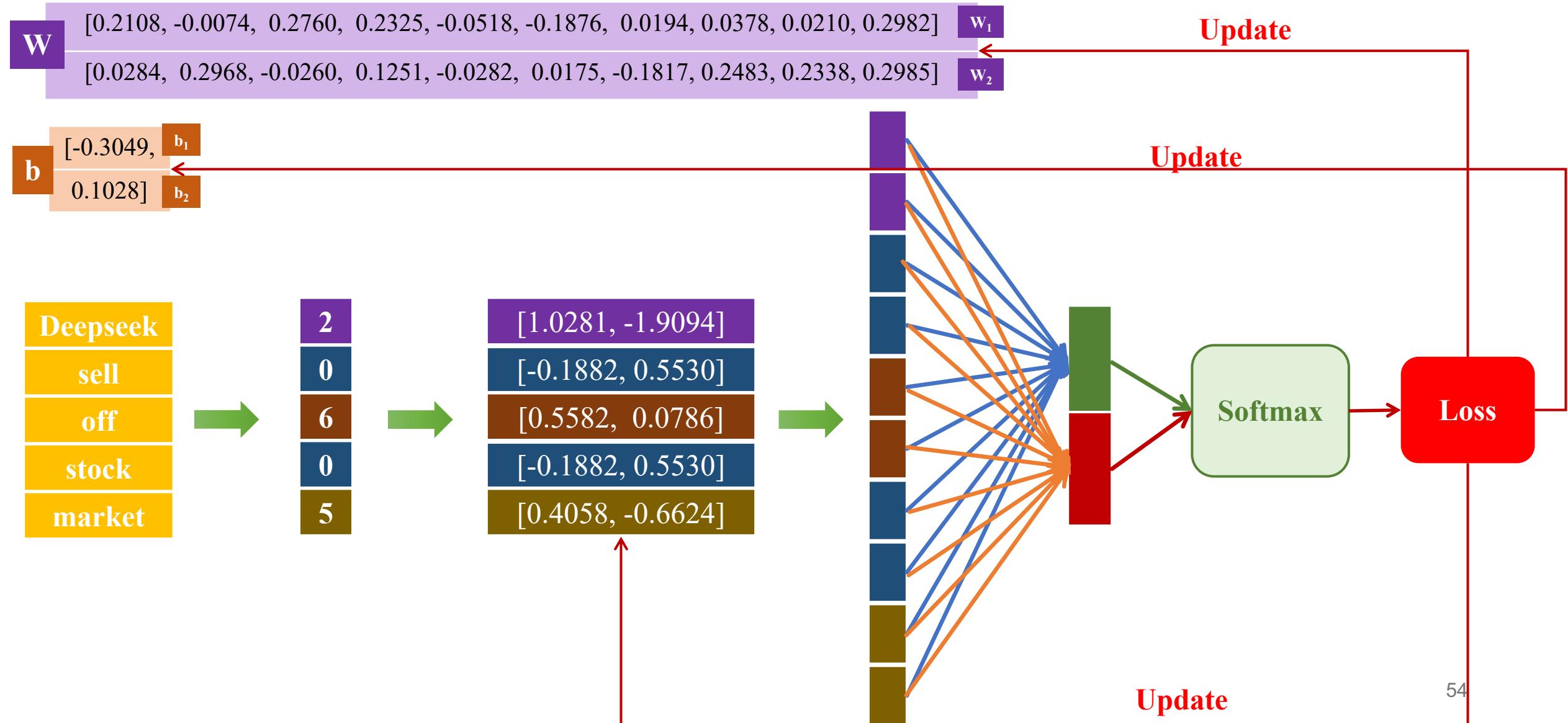
$$\sum_{j=1}^k e^{Z_j} = 0.8083 + 0.7362 = 1.5445$$

$$S_0 = \frac{0.8083}{1.5445} = 0.5233 \quad S_1 = \frac{0.7361}{1.5445} = 0.4766$$

FCN Forward



Text Classification



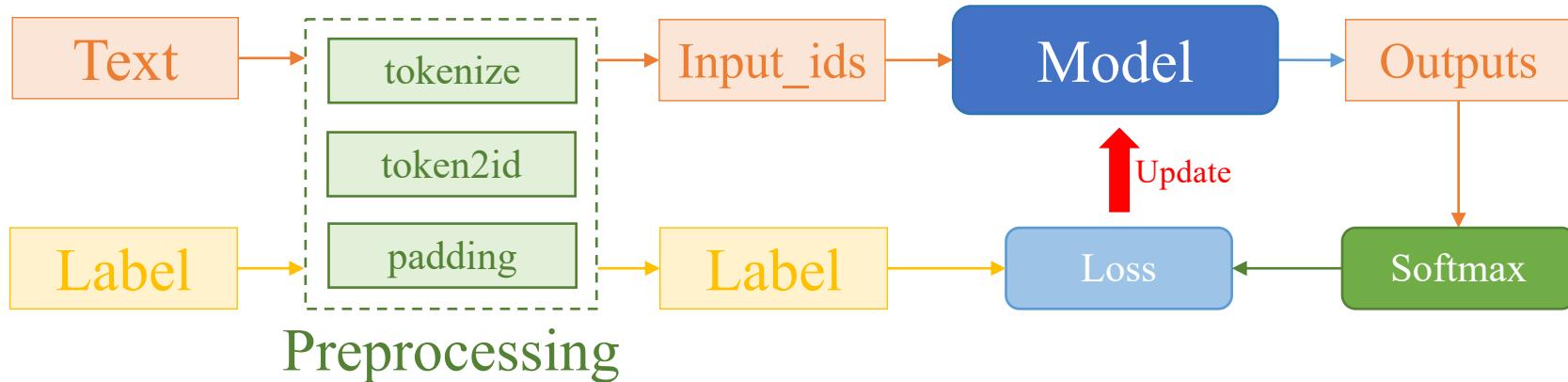
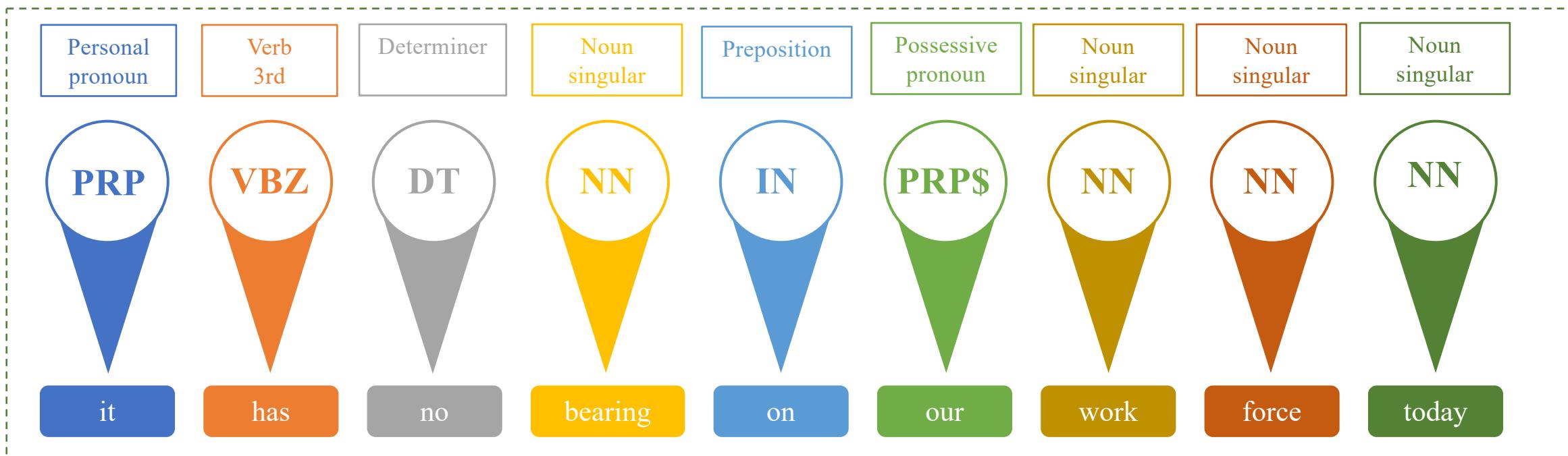
Part-of-Speech Tagging

Code

Part-of-Speech Tagging

- Preprocessing
- Embedding
- POS Tagging

Part-of-Speech Tagging



Preprocessing

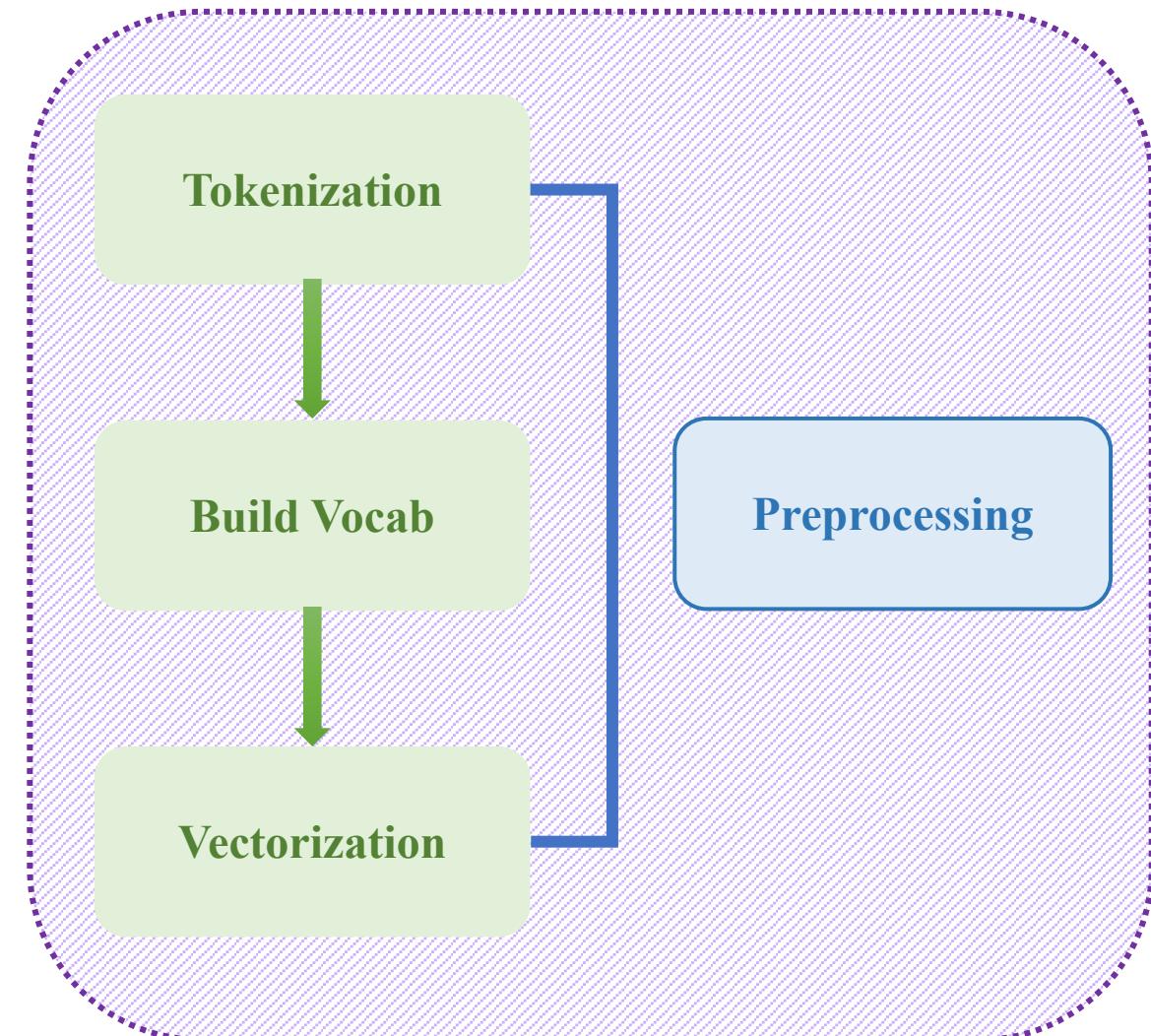
❖ Text Corpus

AI is growing fast

Models are large

```
corpus = [  
    "AI is growing fast",  
    "Models are large"  
]
```

```
# 0: noun/pronoun - 1: verb - others - 2  
labels = [[0, 1, 1, 2],  
          [0, 1, 2]]
```

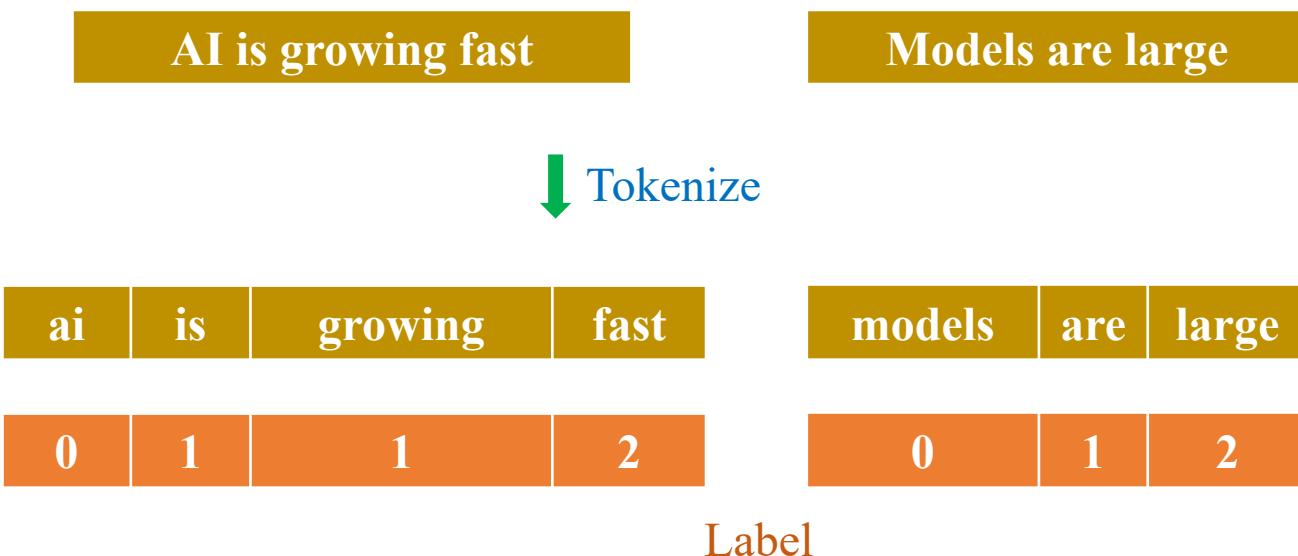


Tokenization

Example corpus

sample1: ‘AI is growing fast’

sample2: ‘Models are large’



```
from torchtext.data.utils import get_tokenizer

sample1 = corpus[0]
sample2 = corpus[1]

# Define tokenizer
tokenizer = get_tokenizer('basic_english')

# Tokenize sample 1
sample1_tokens = tokenizer(sample1)
print(sample1_tokens)

# Tokenize sample 2
sample2_tokens = tokenizer(sample2)
print(sample2_tokens)

['ai', 'is', 'growing', 'fast']
['models', 'are', 'large']
```

Building Vocab

```
from torchtext.data.utils import get_tokenizer
from torchtext.vocab import build_vocab_from_iterator

# Define tokenizer
tokenizer = get_tokenizer('basic_english')

# Create a function to yield list of tokens
def yield_tokens(examples):
    for text in examples:
        yield tokenizer(text)

# Create vocabulary
vocab = build_vocab_from_iterator(yield_tokens(corpus),
                                    max_tokens=vocab_size,
                                    specials=["<unk>", "<pad>"])
vocab.set_default_index(vocab["<unk>"])
```

vocab.get_stoi()

```
{'is': 6,
 'growing': 5,
 'models': 8,
 'are': 3,
 'ai': 2,
 'fast': 4,
 '<pad>': 1,
 'large': 7,
 '<unk>': 0}
```

0	<unk>
1	<pad>
2	ai
3	are
4	fast
5	growing
6	is
7	large
8	models

Vocabulary

Vectorization

0	<unk>
1	<pad>
2	ai
3	are
4	fast
5	growing
6	is
7	large
8	models

Vocabulary

```
sample1_ids = [vocab[token] for token in sample1_tokens]  
sample2_ids = [vocab[token] for token in sample2_tokens]
```

```
print(sample1_ids)  
print(sample2_ids)
```

[2, 6, 5, 4]
[8, 3, 7]

ai | is | growing | fast



Vectorize

Input: 2 | 6 | 5 | 4

Label: 0 | 1 | 1 | 2

models | are | large



Vectorize

Input: 8 | 3 | 7

Label: 0 | 1 | 2

Vectorization

0	<unk>
1	<pad>
2	ai
3	are
4	fast
5	growing
6	is
7	large
8	models

Vocabulary

```
sample1_ids = [vocab[token] for token in sample1_tokens]  
sample2_ids = [vocab[token] for token in sample2_tokens]
```

```
print(sample1_ids)  
print(sample2_ids)
```

[2, 6, 5, 4]
[8, 3, 7]

ai | is | growing | fast



Vectorize

Input:

2 | 6 | 5 | 4

len = 4

Not equal length

Label:

0 | 1 | 1 | 2

models | are | large



Vectorize

Input:

8 | 3 | 7

len = 3

Label:

0 | 1 | 2

Vectorization

0	<unk>
1	<pad>
2	ai
3	are
4	fast
5	growing
6	is
7	large
8	models

Vocabulary

```
sample1_ids = [vocab[token] for token in sample1_tokens]
sample2_ids = [vocab[token] for token in sample2_tokens]
```

```
print(sample1_ids)
print(sample2_ids)
```

[2, 6, 5, 4]
[8, 3, 7]

sequence_length = 4

ai | is | growing | fast



Vectorize

Input:

2	6	5	4
---	---	---	---

len = 4

Not equal length

Label:

0	1	1	2
---	---	---	---

models | are | large



Vectorize

Input:

8	3	7
---	---	---

Label:

0	1	2
---	---	---

len = 3

Add <pad> token

Vectorization

0	<unk>
1	<pad>
2	ai
3	are
4	fast
5	growing
6	is
7	large
8	models

```
sample1_ids = [vocab[token] for token in sample1_tokens]  
sample2_ids = [vocab[token] for token in sample2_tokens]
```

```
print(sample1_ids)  
print(sample2_ids)
```

[2, 6, 5, 4]
[8, 3, 7]

sequence_length = 4

ai | is | growing | fast



Vectorize

Input: 2 | 6 | 5 | 4

Label: 0 | 1 | 1 | 2

models | are | large | <pad>



Vectorize

Input: 8 | 3 | 7 | 1

Vectorization

0	<unk>
1	<pad>
2	ai
3	are
4	fast
5	growing
6	is
7	large
8	models

Vocabulary

```
sample1_ids = [vocab[token] for token in sample1_tokens]
sample2_ids = [vocab[token] for token in sample2_tokens]
```

```
print(sample1_ids)
print(sample2_ids)
```

[2, 6, 5, 4]
[8, 3, 7]

sequence_length = 4

ai | is | growing | fast



Input:

2	6	5	4
---	---	---	---

Label:

0	1	1	2
---	---	---	---

models | are | large | <pad>



Input:

8	3	7	1
---	---	---	---

Label:

0	1	2	?
---	---	---	---

Index	Label
0	(Pro)noun
1	Verb
2	Others

Vectorization

0	<unk>
1	<pad>
2	ai
3	are
4	fast
5	growing
6	is
7	large
8	models

ai | is | growing | fast

↓ Vectorize

Input: 2 | 6 | 5 | 4

Label: 0 | 1 | 1 | 2

sequence_length = 4

models | are | large | <pad>

↓ Vectorize

Input: 8 | 3 | 7 | 1

Label: 0 | 1 | 2 | 3

Index	Label
0	(Pro)noun
1	Verb
2	Others
3	<pad>

0	<unk>
1	<pad>
2	ai
3	are
4	fast
5	growing
6	is
7	large
8	models

Vocabulary

```
# Tokenize and numericalize your samples
def vectorize(text, vocab, sequence_length, sequence_label):
    tokens = tokenizer(text)

    token_ids = [vocab[token] for token in tokens]
    token_ids = token_ids + [vocab["<pad>"]] * (sequence_length - len(tokens))
    sequence_label = sequence_label + [3] * (sequence_length - len(tokens))

    return torch.tensor(token_ids, dtype=torch.long), torch.tensor(sequence_label, dtype=torch.long)

# Vectorize the samples
sentence_vecs = []
label_vecs = []
for sentence, labels in zip(corpus, labels):
    sentence_vec, labels_vec = vectorize(sentence, vocab, sequence_length, labels)
    sentence_vecs.append(sentence_vec)
    label_vecs.append(labels_vec)
```

ai | is | growing | fast

sequence_length = 4

Vectorize

Input: 2 | 6 | 5 | 4

Label: 0 | 1 | 1 | 2

models | are | large | <pad>

Vectorize

Input: 8 | 3 | 7 | 1

Label: 0 | 1 | 2 | 3

Index	Label
0	(Pro)noun
1	Verb
2	Others
3	<pad>

Vectorization

0	<unk>
1	<pad>
2	ai
3	are
4	fast
5	growing
6	is
7	large
8	models

Vocabulary

AI is growing fast Models are large

↓ Tokenize

ai | is | growing | fast models | are | large

↓ Vectorize

Input: 2 | 6 | 5 | 4 8 | 3 | 7 | 1

Label: 0 | 1 | 1 | 2 0 | 1 | 2 | 3

→ Token classification problem

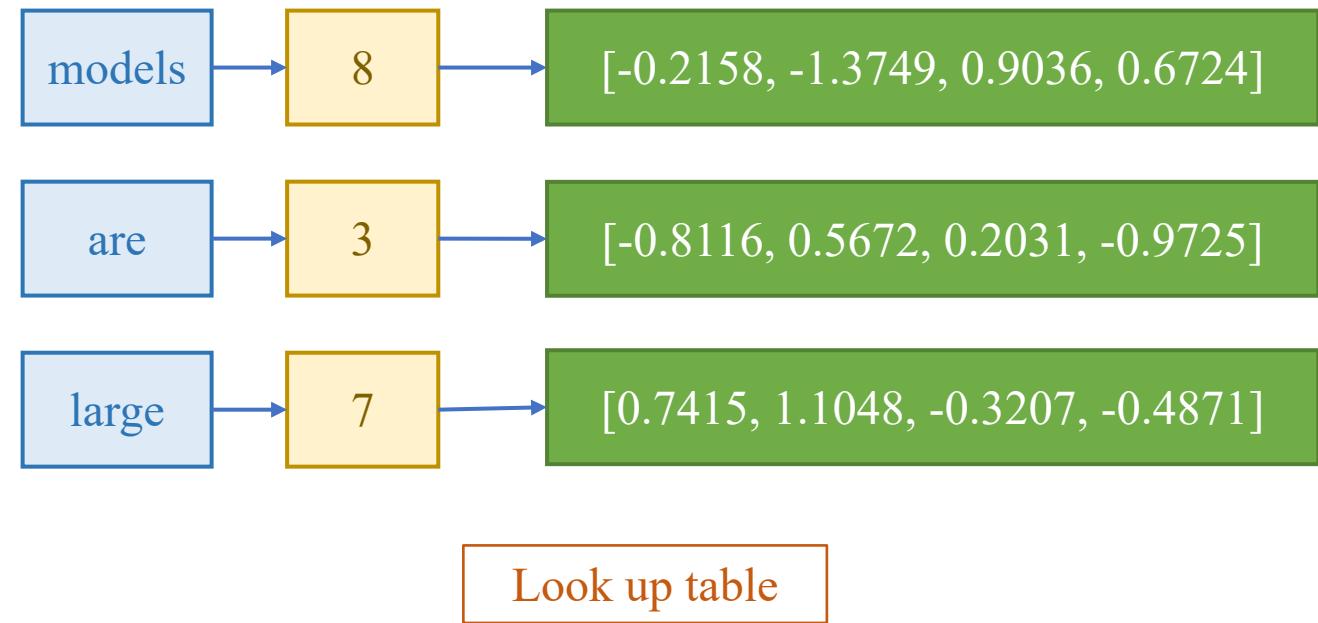
Part-of-Speech Tagging

❖ Embedding

index	vocab	Embedding vector
0	<unk>	[0.7321, -1.2503, 0.8154, -0.4987]
1	<pad>	[-0.2165, 1.7823, -0.6892, 0.9476]
2	ai	[1.1043, 0.3289, -1.5284, 0.6520]
3	are	[-0.8116, 0.5672, 0.2031, -0.9725]
4	fast	[0.4528, -0.1476, 1.2930, 0.7842]
5	growing	[-1.3459, 0.6791, -0.8853, 1.1207]
6	is	[0.9823, -0.5504, 0.4389, -1.2396]
7	large	[0.7415, 1.1048, -0.3207, -0.4871]
8	models	[-0.2158, -1.3749, 0.9036, 0.6724]

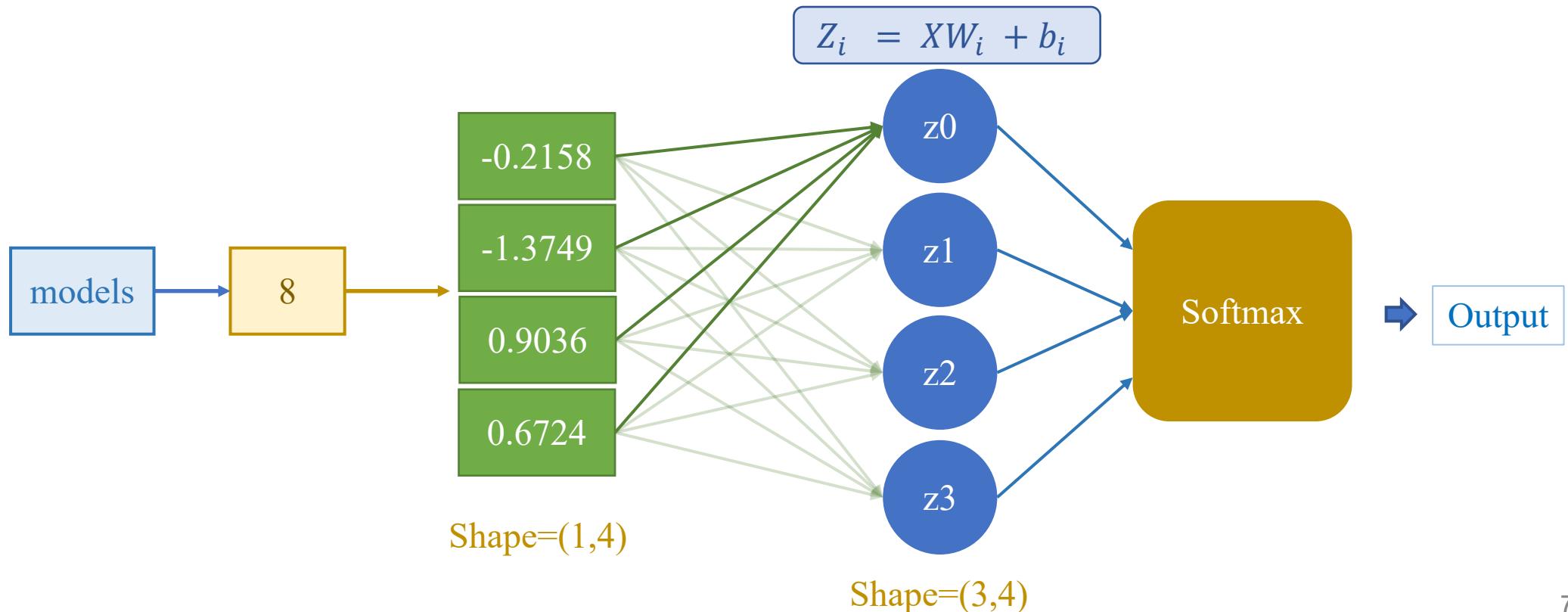
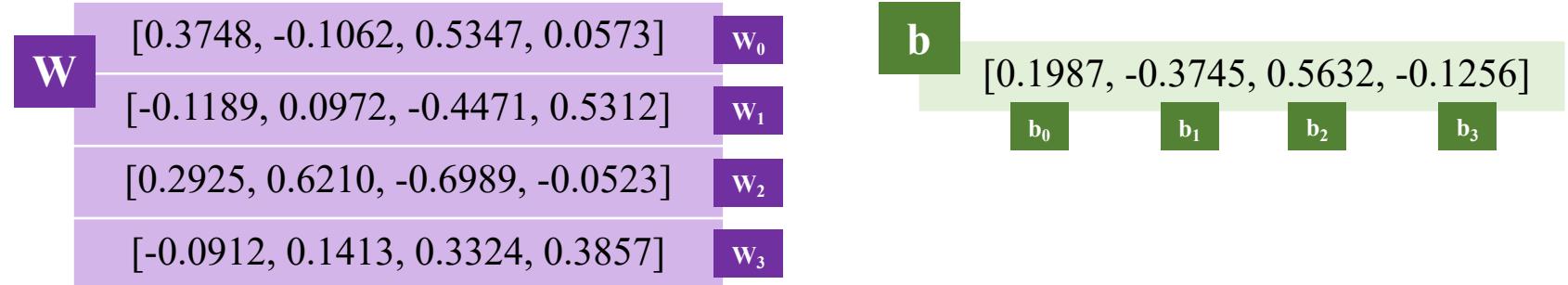


From Word2Vec models or
learnable parameters



Part-of-Speech Tagging

❖ A word input

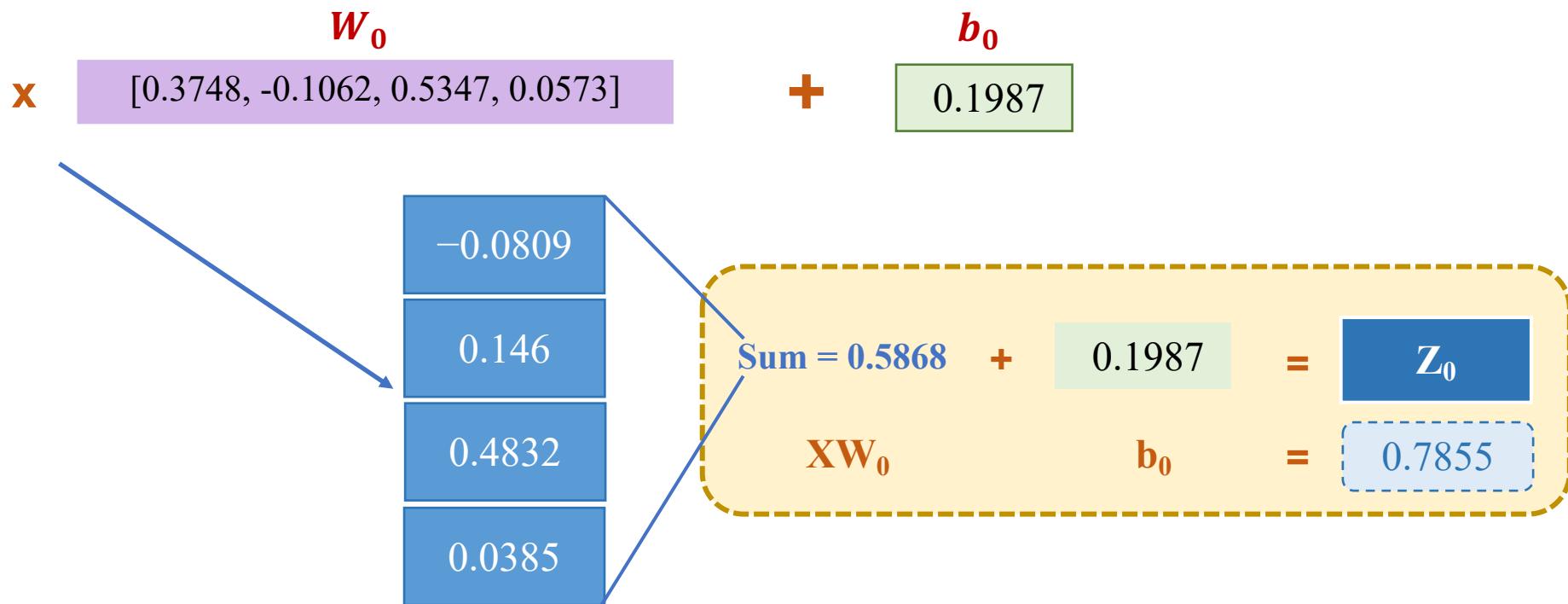
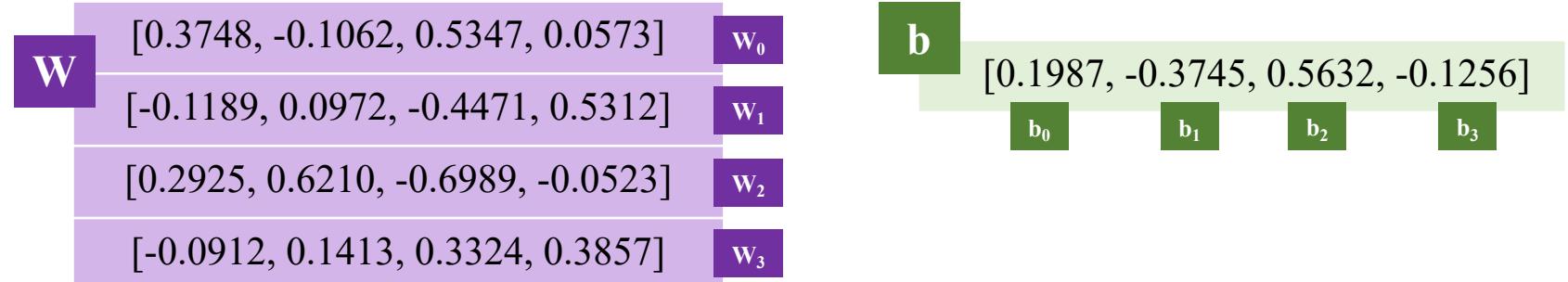


Part-of-Speech Tagging

❖ A word input

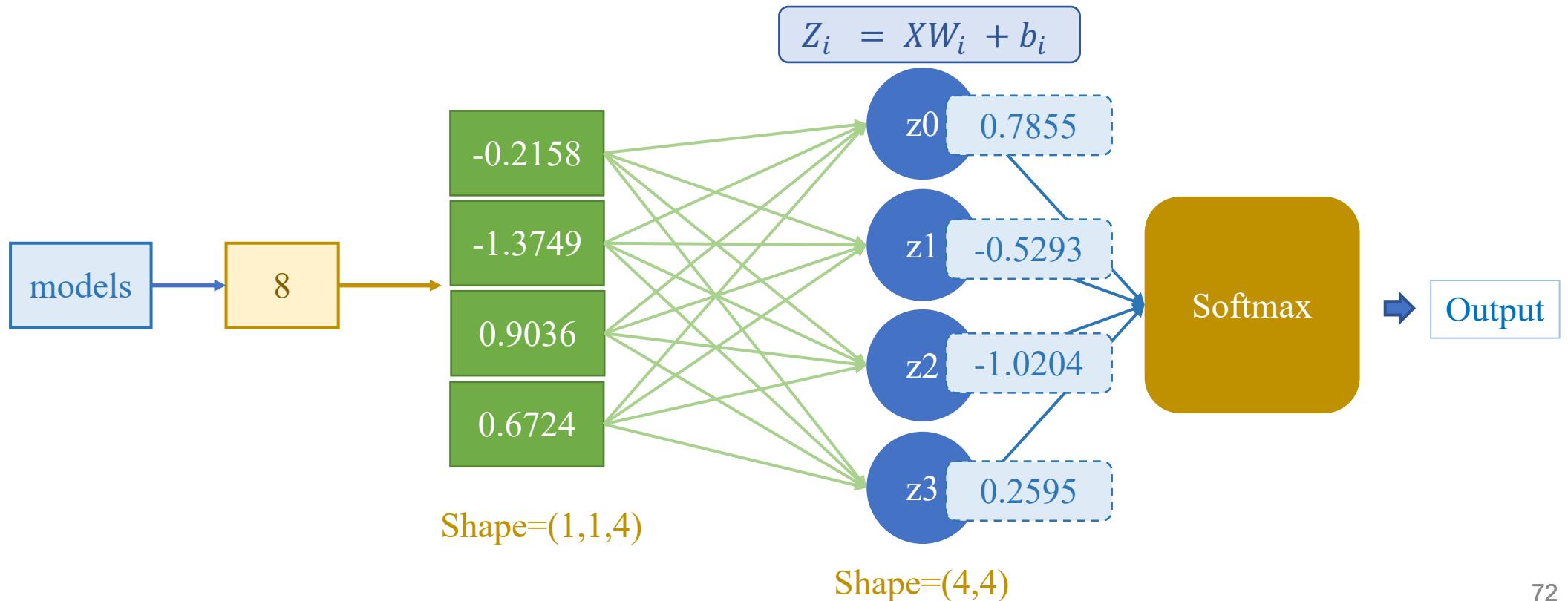
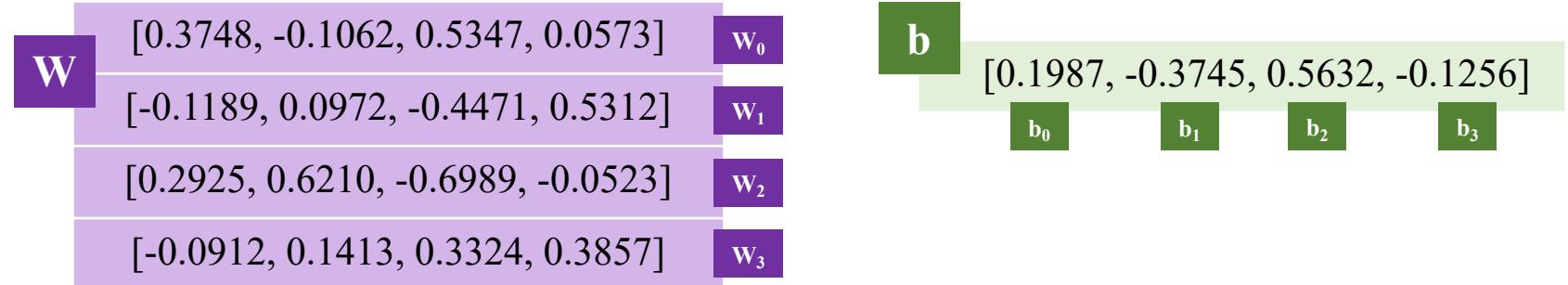
$$\boxed{Z_0 = XW_0 + b_0}$$

X	-0.2158
	-1.3749
	0.9036
	0.6724



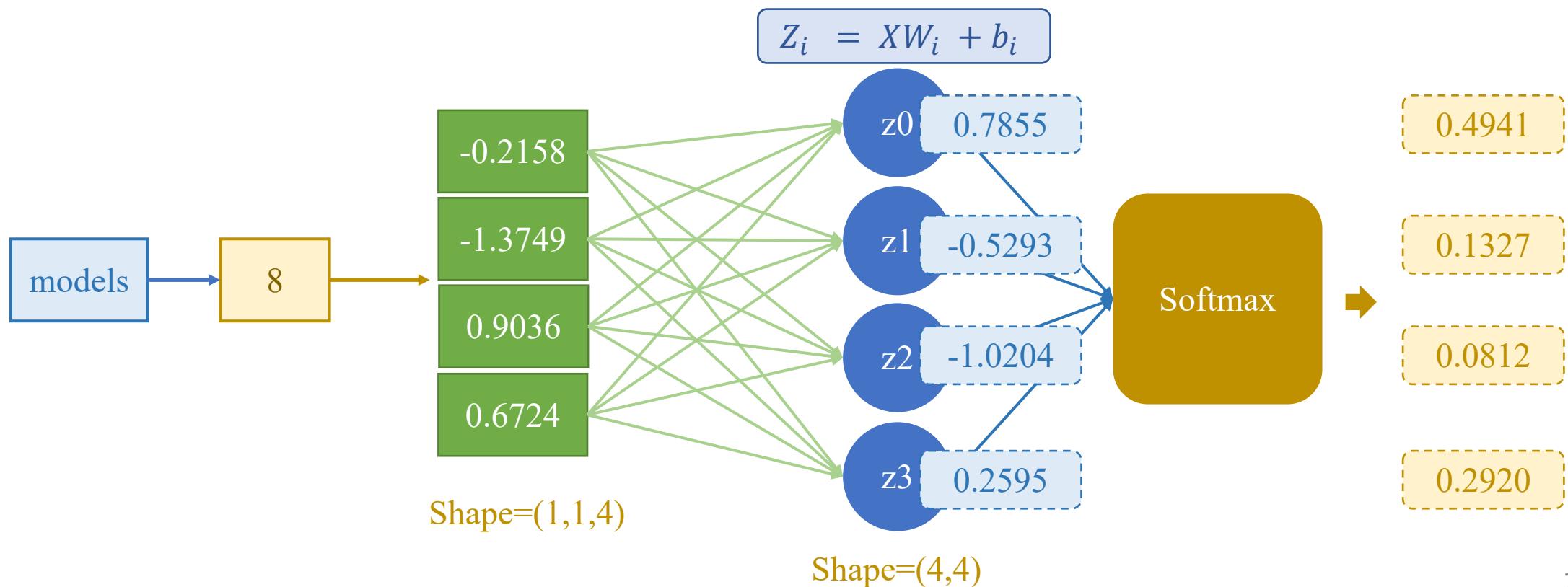
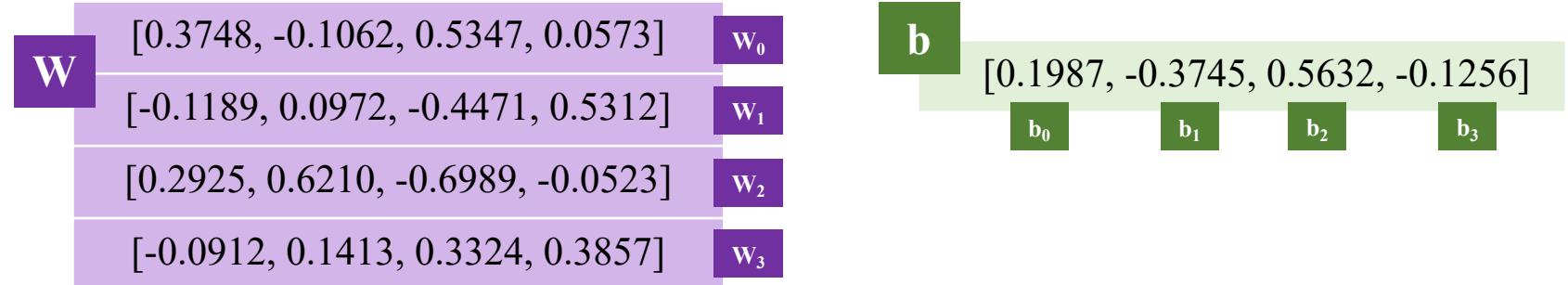
Part-of-Speech Tagging

❖ A word input



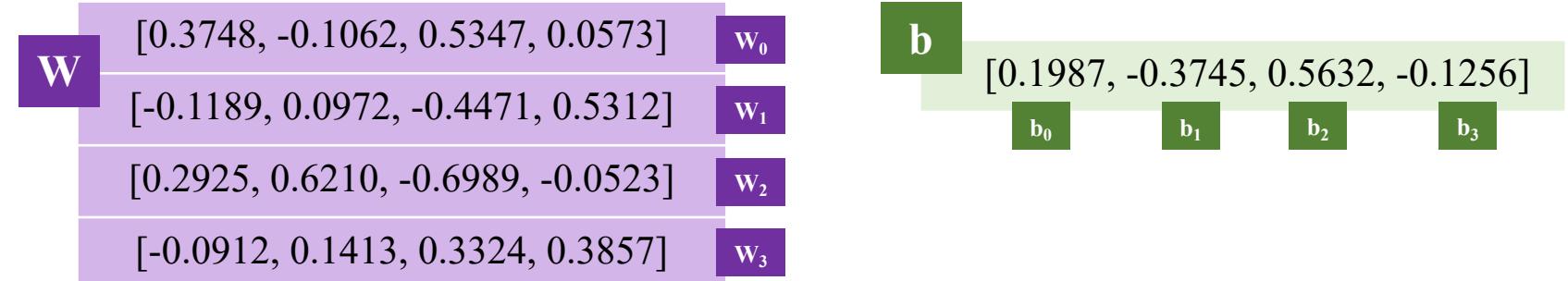
Part-of-Speech Tagging

❖ A word input

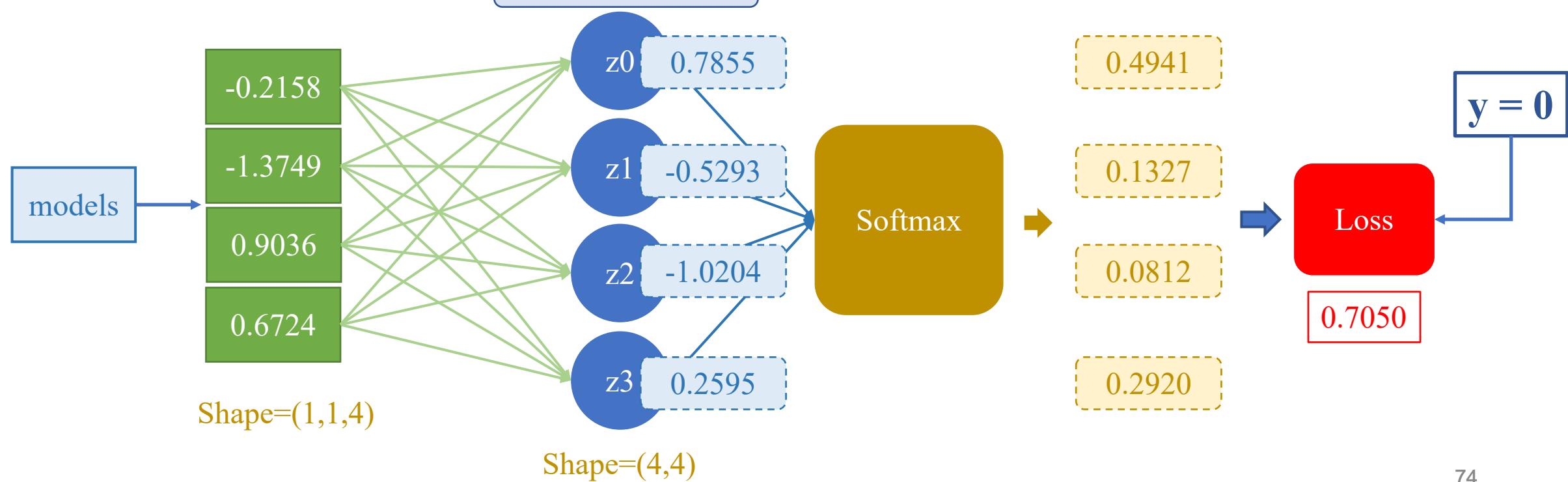


Part-of-Speech Tagging

❖ A word input

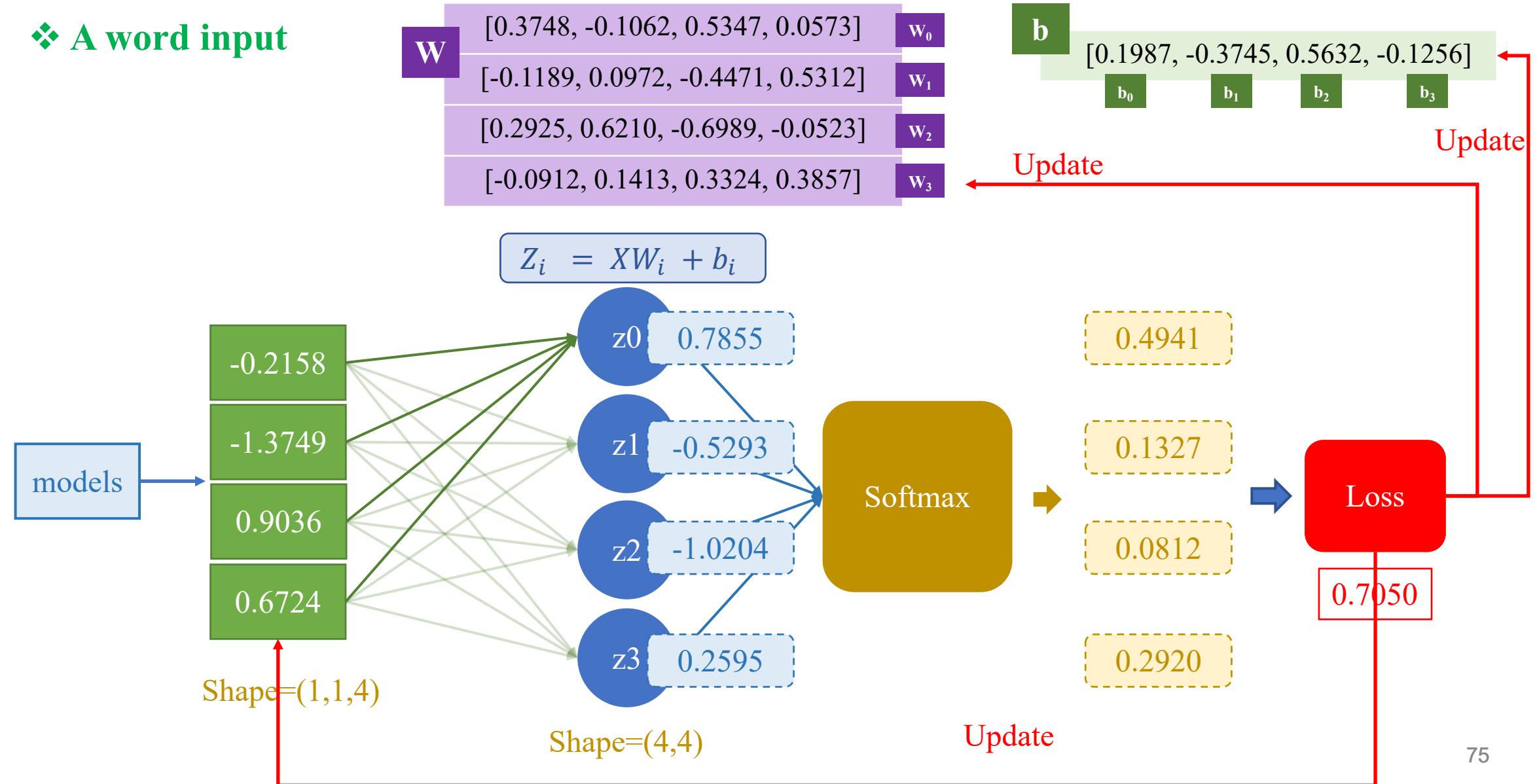


$$Z_i = XW_i + b_i$$



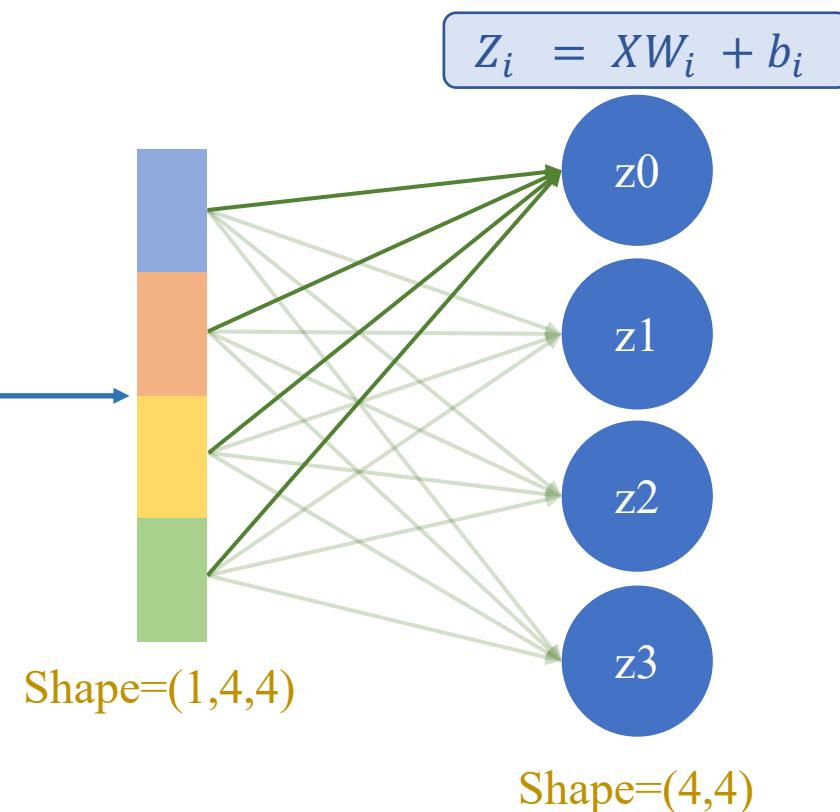
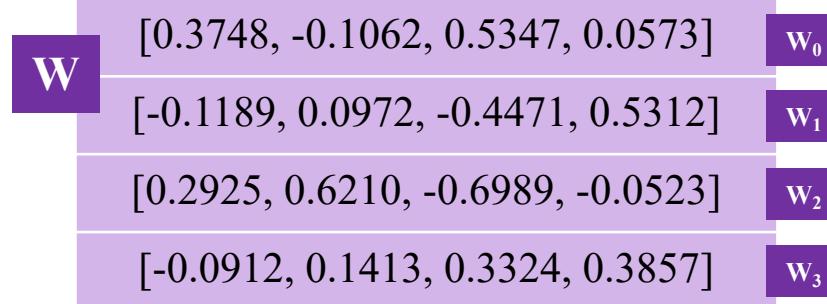
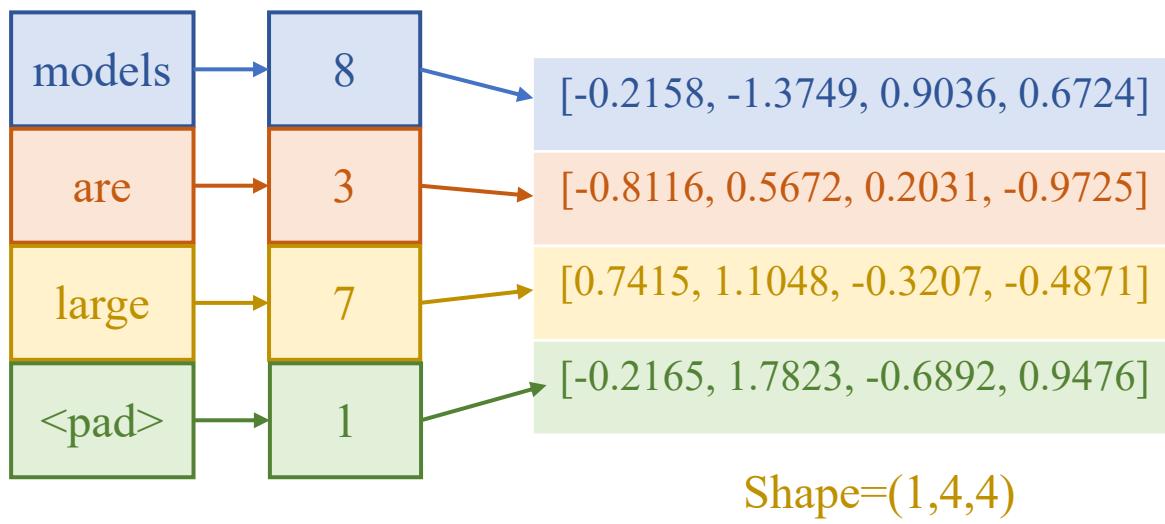
Part-of-Speech Tagging

❖ A word input



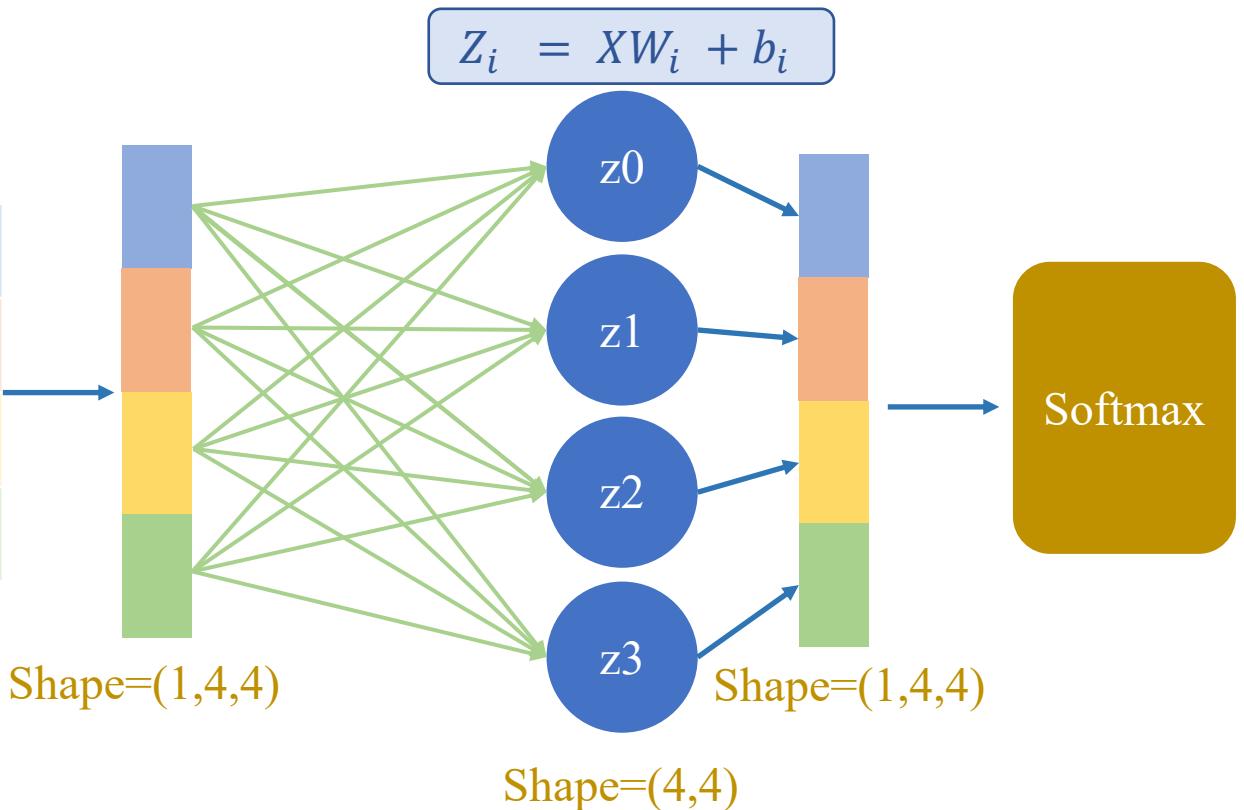
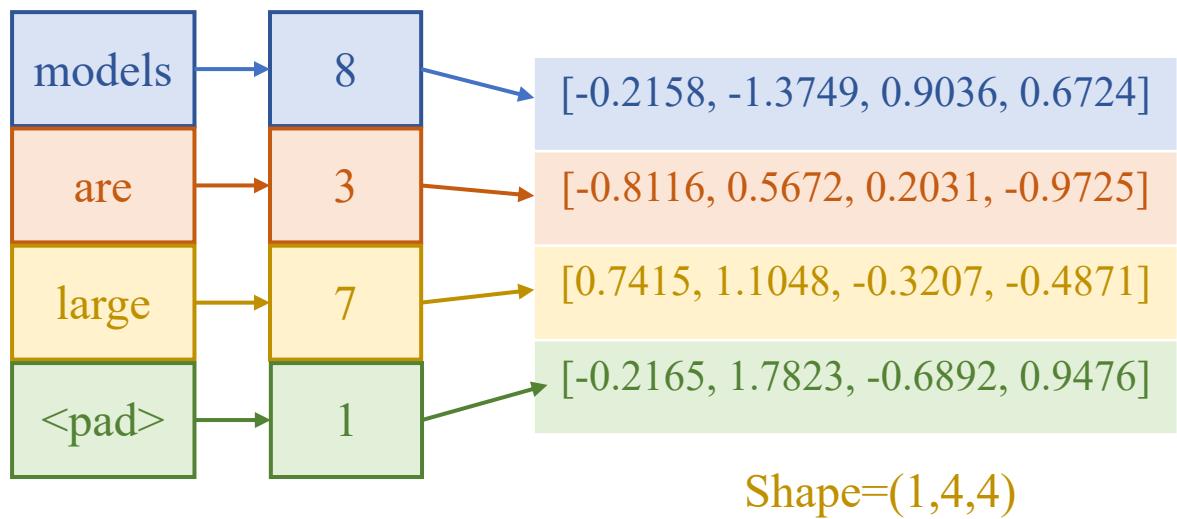
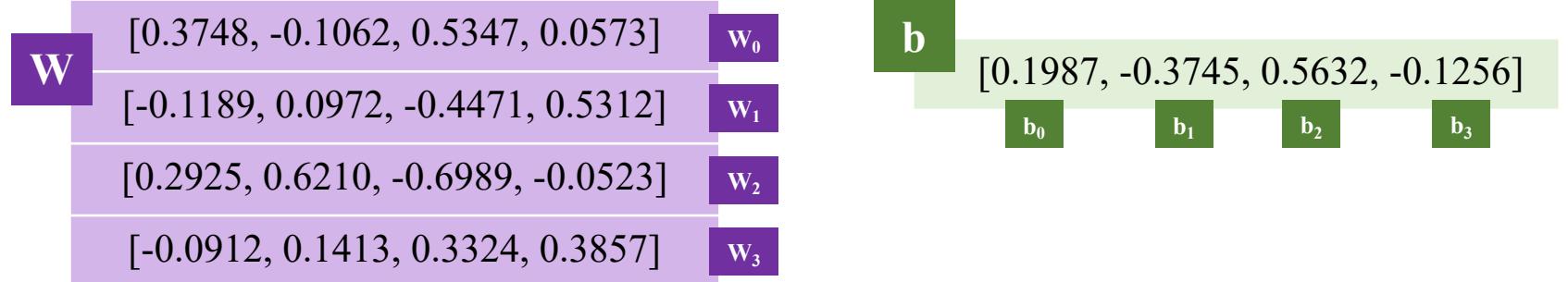
Part-of-Speech Tagging

❖ A word input



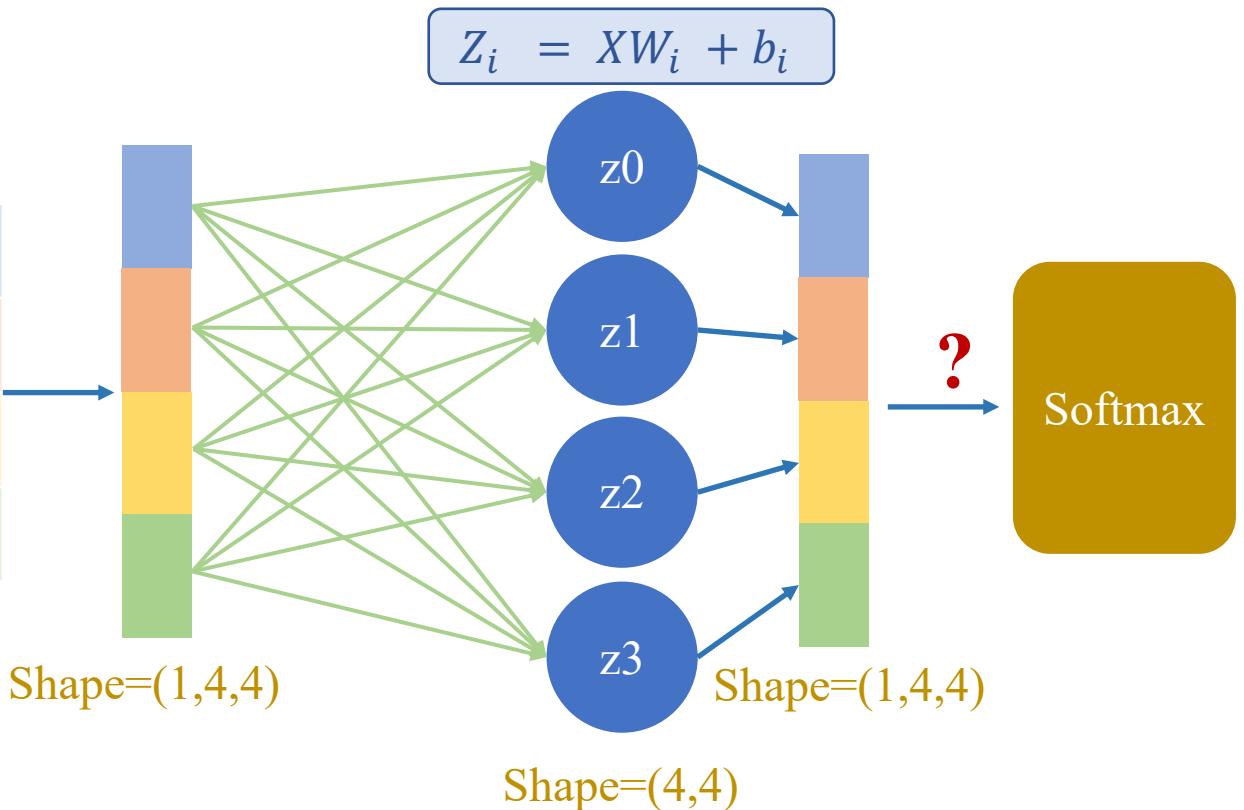
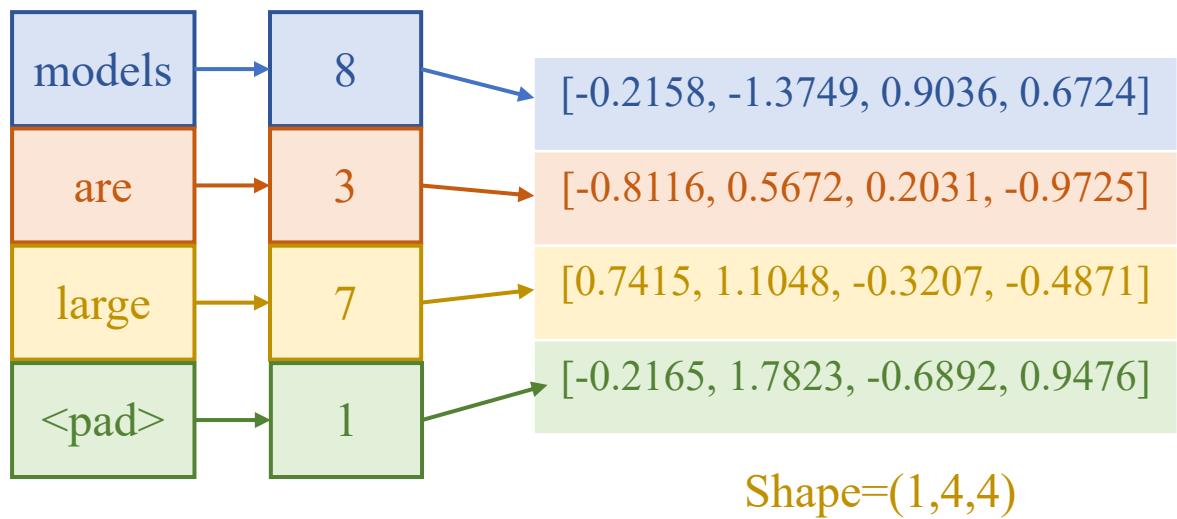
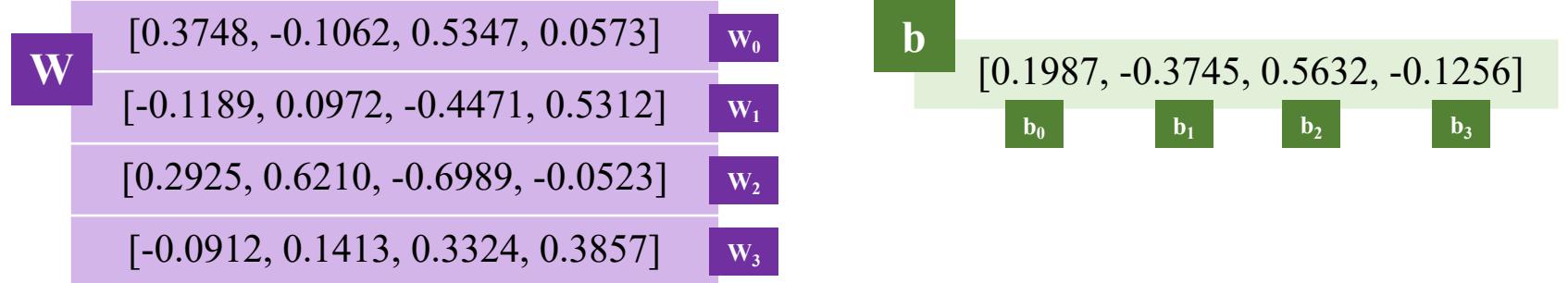
Part-of-Speech Tagging

❖ A word input



Part-of-Speech Tagging

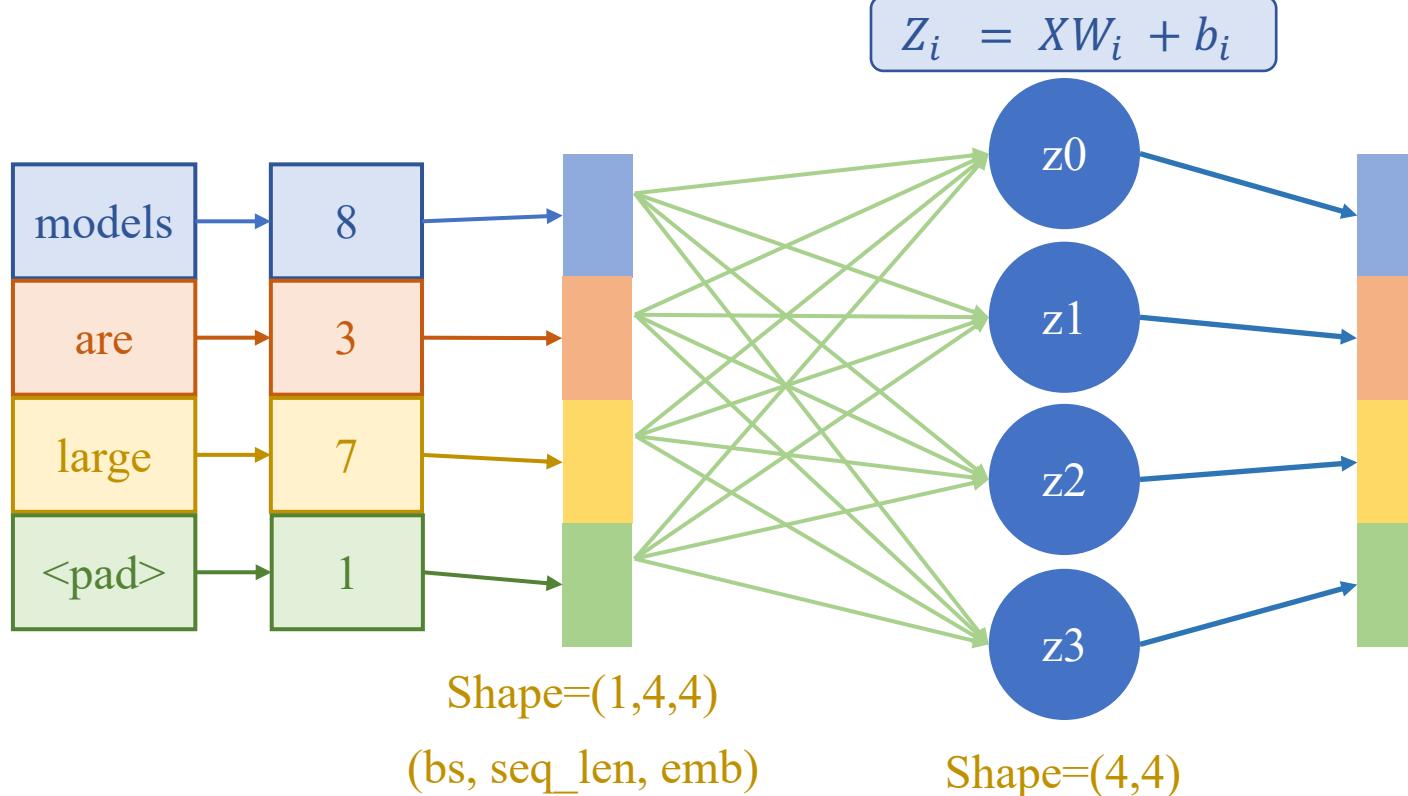
❖ A word input



Part-of-Speech Tagging

❖ A word input

nn.Linear

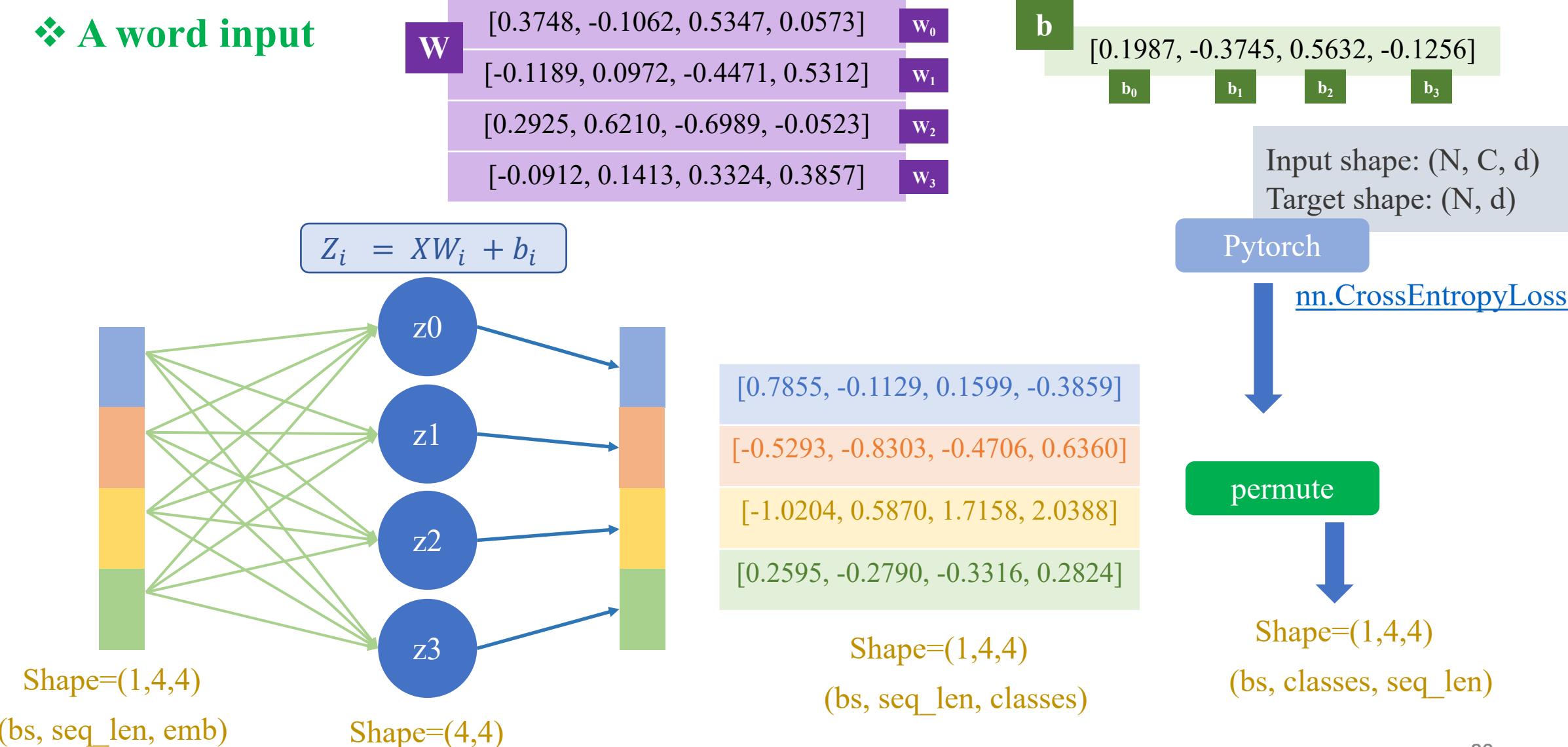


[0.7855, -0.1129, 0.1599, -0.3859]
[-0.5293, -0.8303, -0.4706, 0.6360]
[-1.0204, 0.5870, 1.7158, 2.0388]
[0.2595, -0.2790, -0.3316, 0.2824]

Shape = (1,4,4)
(bs, seq_len, classes)

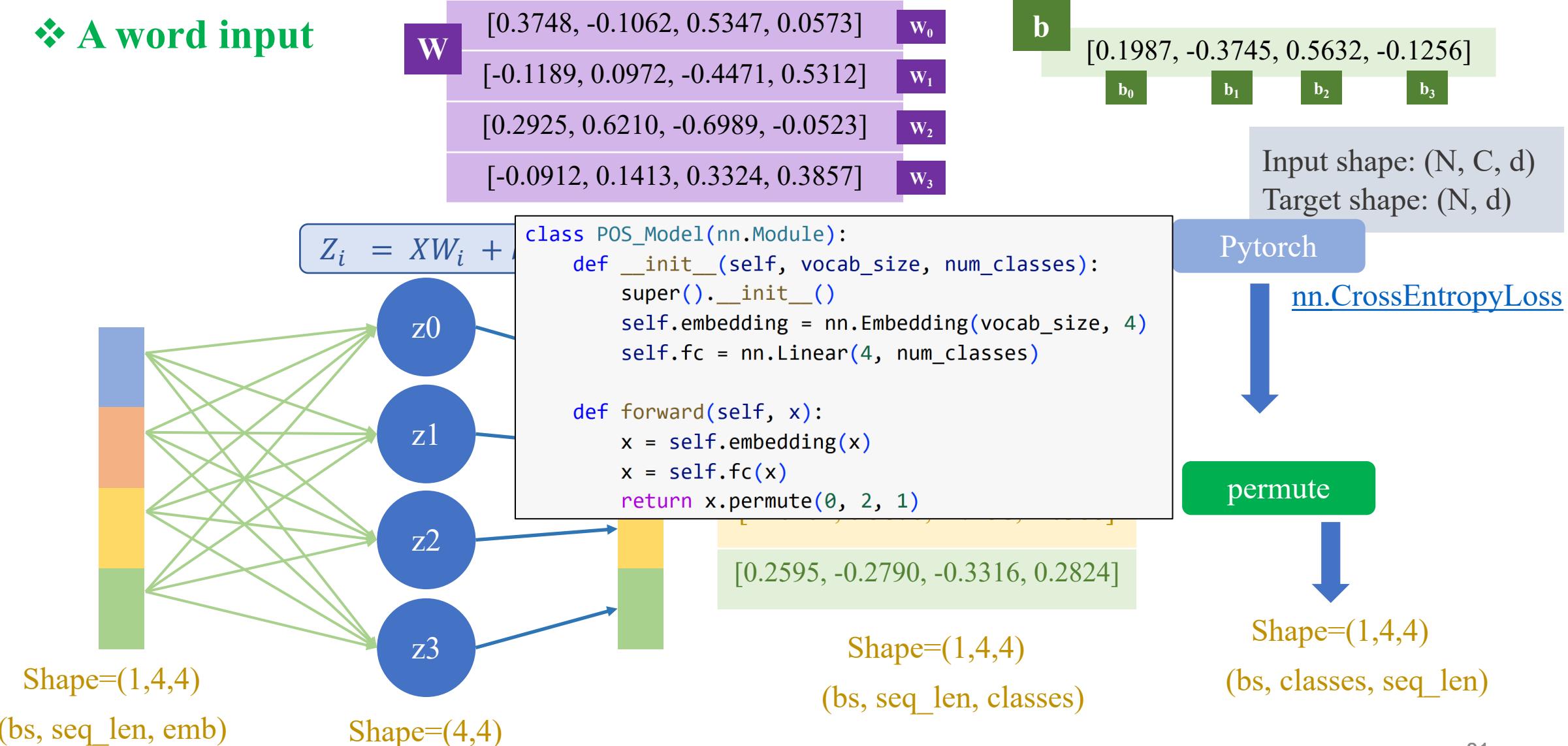
Part-of-Speech Tagging

❖ A word input



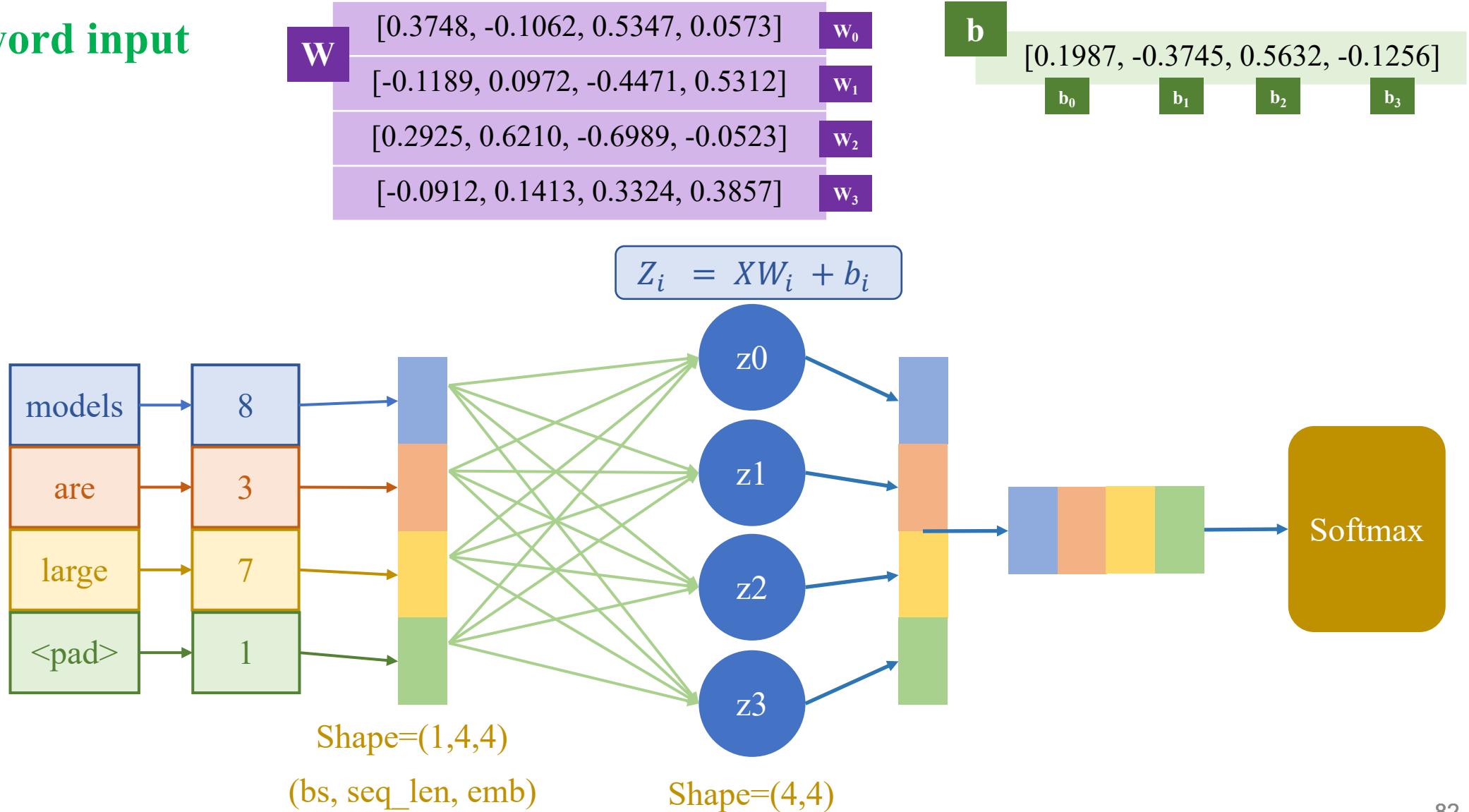
Part-of-Speech Tagging

❖ A word input



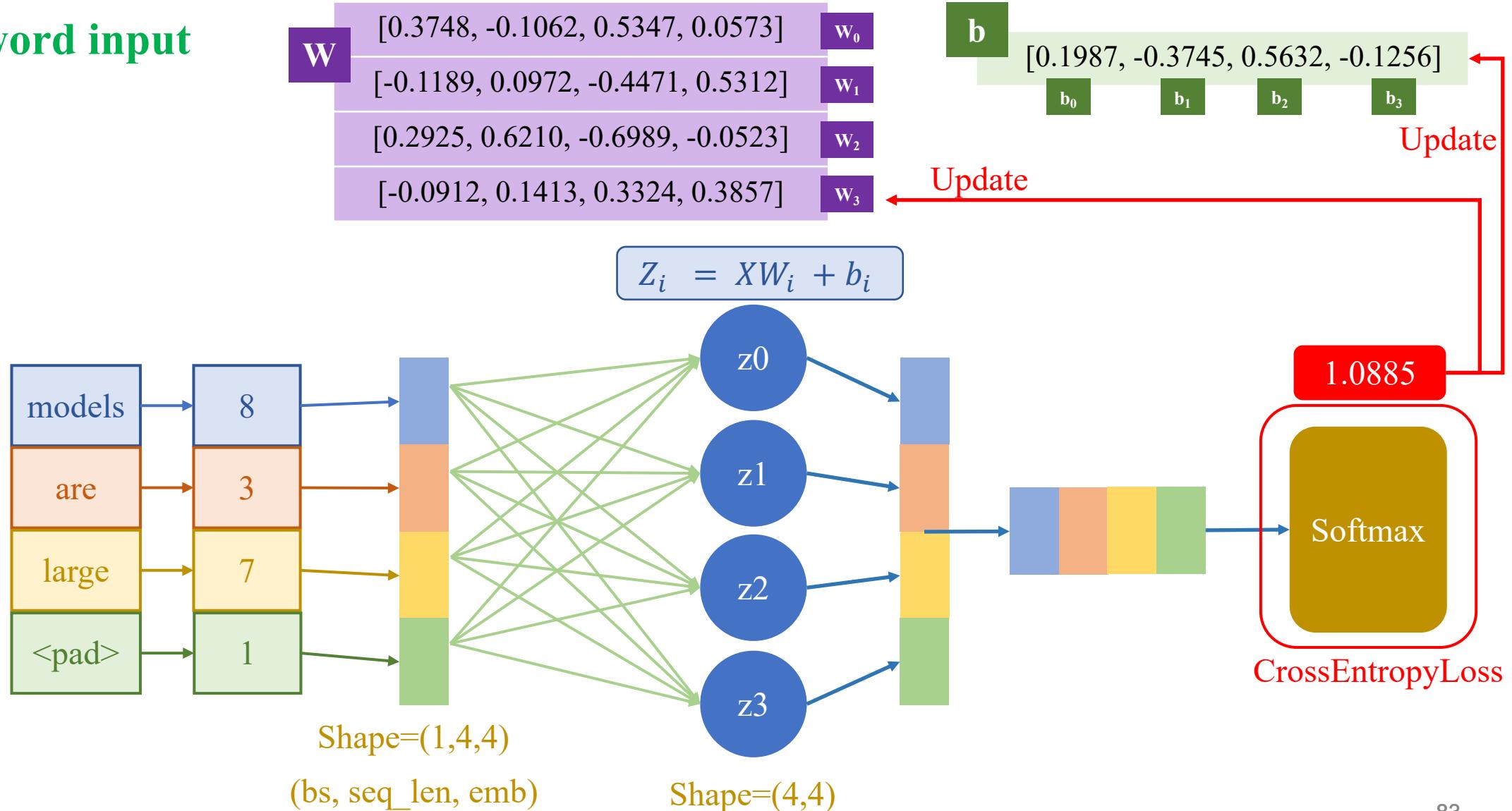
Part-of-Speech Tagging

❖ A word input



Part-of-Speech Tagging

❖ A word input



Question

