

Text Generation

Quang-Vinh Dinh
Ph.D. in Computer Science

Outline

- **Motivation to Text Generation**
- **Simple Models**
- **Using RNNs**
- **Examples**
- **Using Masked Encoder (~Decoder)**

Self-Supervision Using Image Data

Noisy images



Model



Clean images



Grayscale images



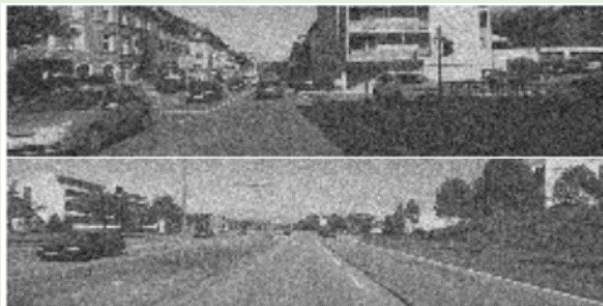
Model



Color images



Noisy and grayscale images



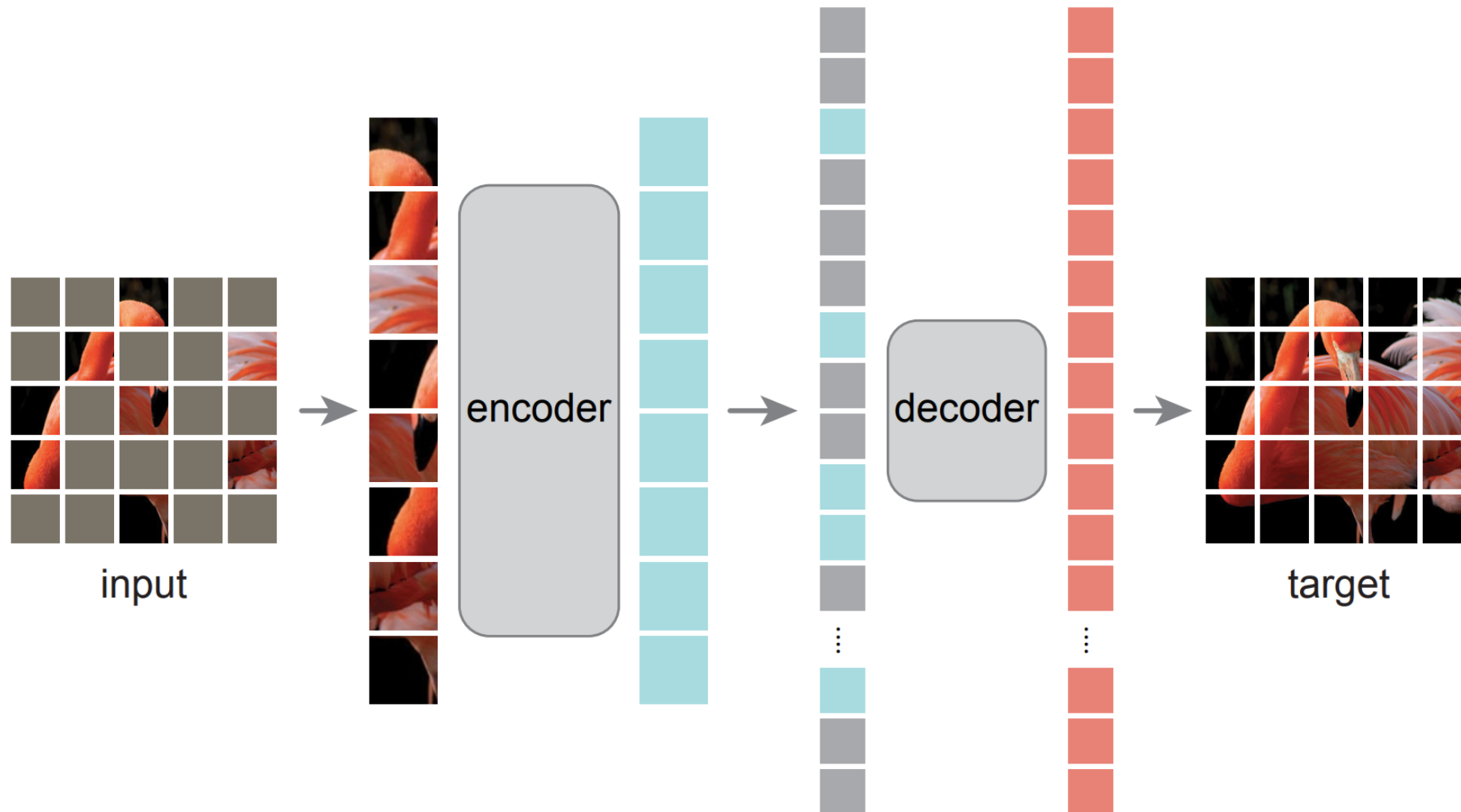
Model



Clean and color images



Self-Supervision Using Image Data



Self-Supervision Using Text Data

❖ How?

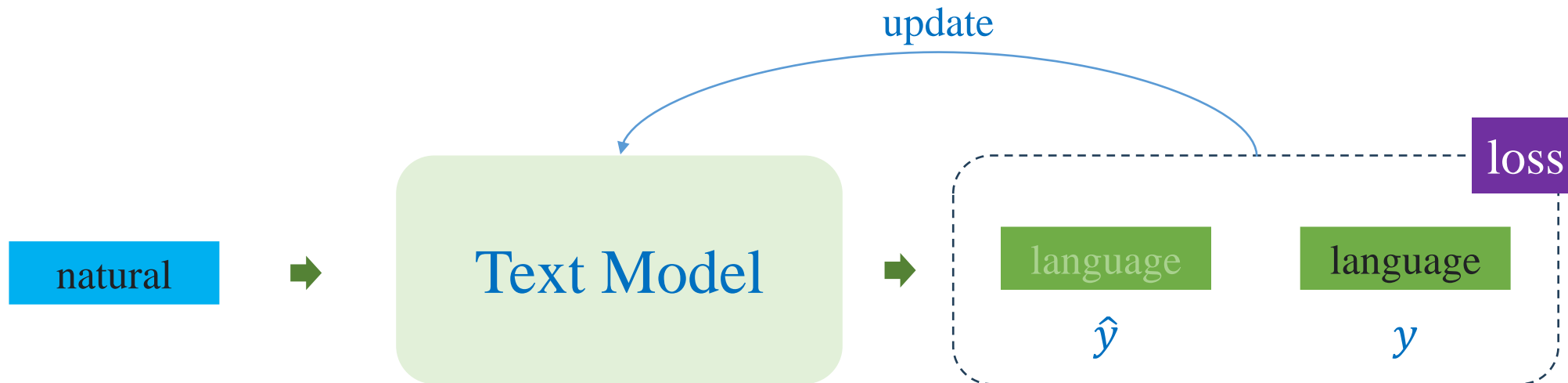
natural language processing is a branch of artificial intelligence

natural

language

x

y



Text Generation

❖ Applications

❖ From the viewpoint of users

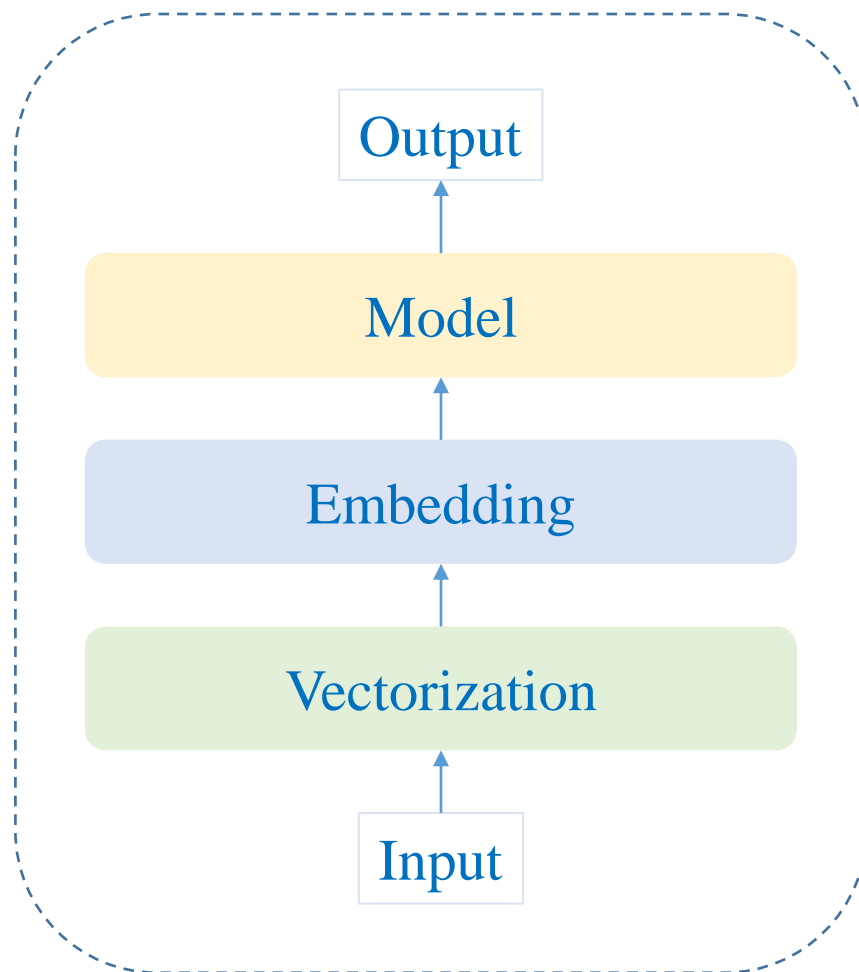
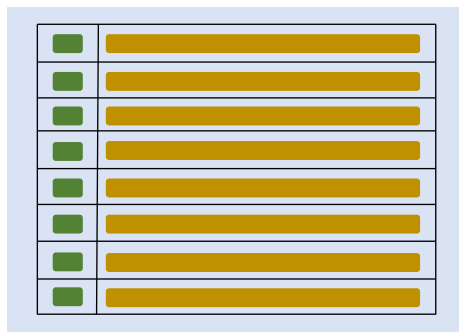
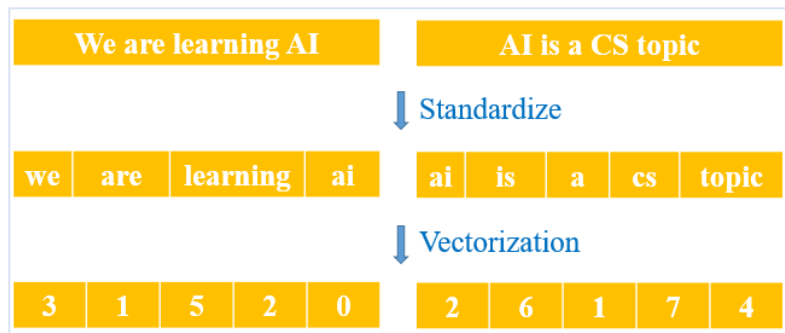
Inference

```
1 generate_text("thúy kiều", 4, generator_model)
```

'thúy kiều sắc sảo khôn ngoan'

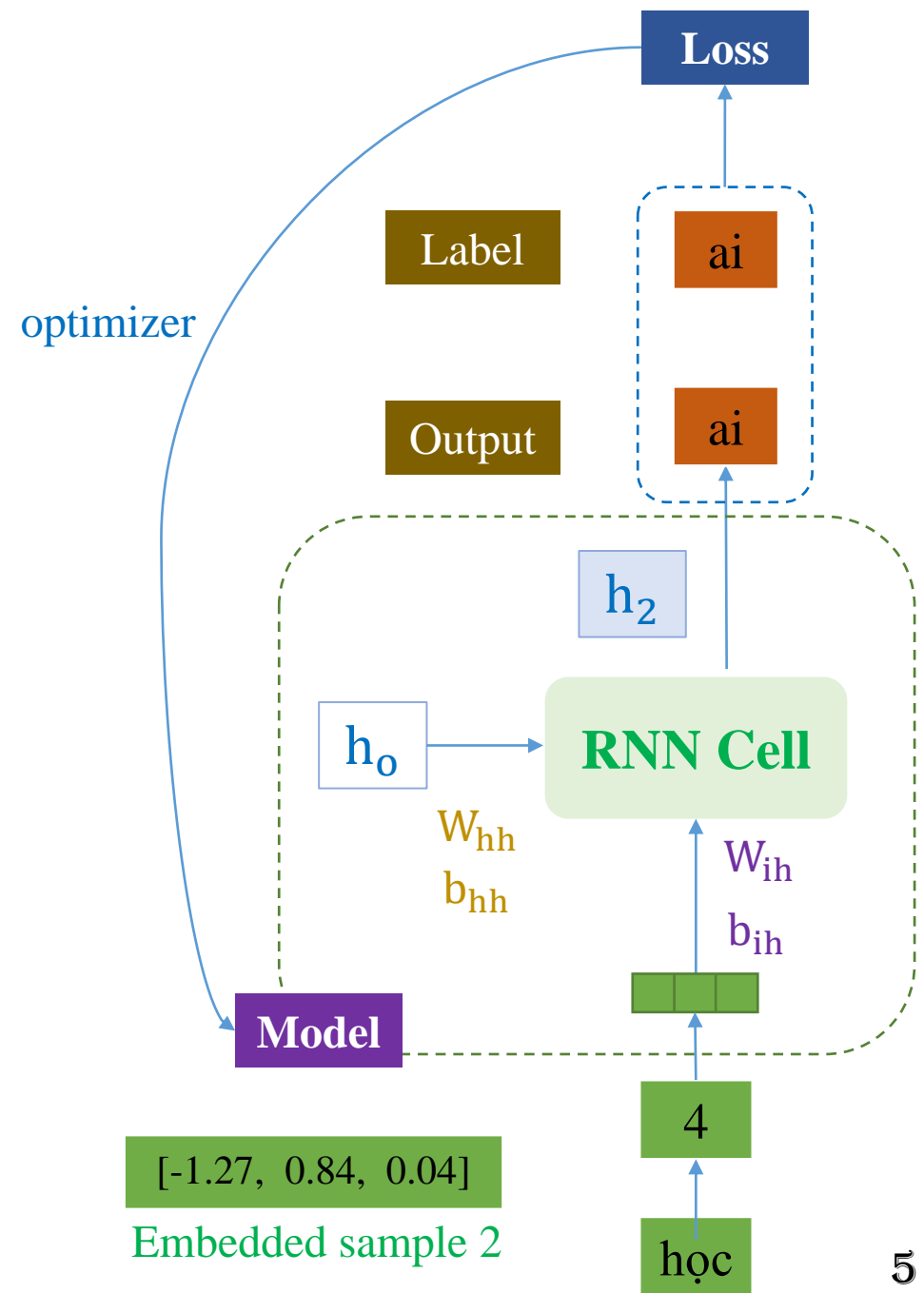
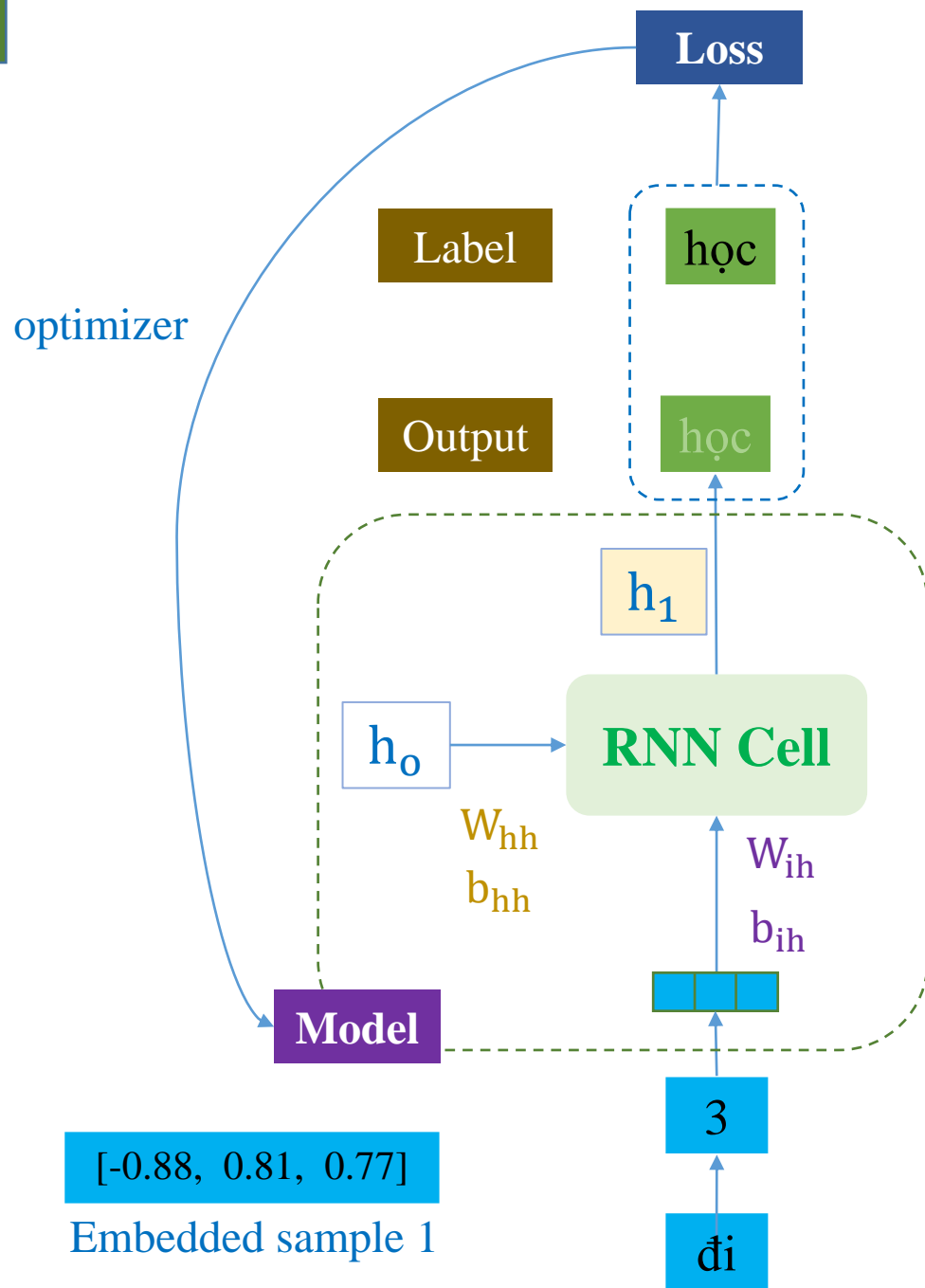
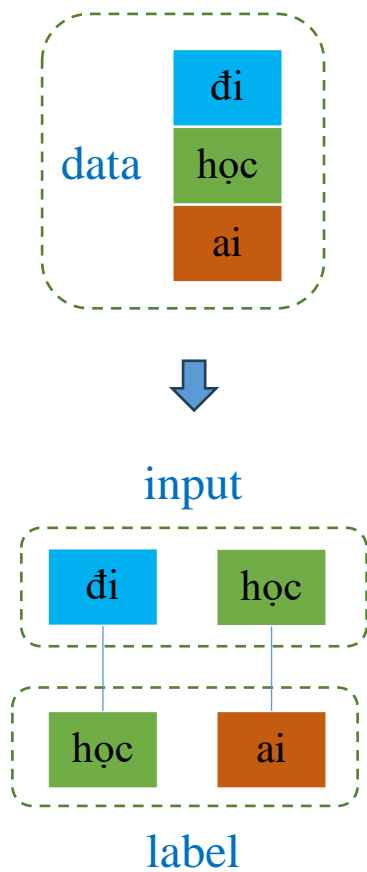
```
1 generate_text("kim trọng và thúy kiều", 6, generator_model)
```

'kim trọng và thúy kiều nhân gian sao tan tành ngay'



Simple Model

Simple Model

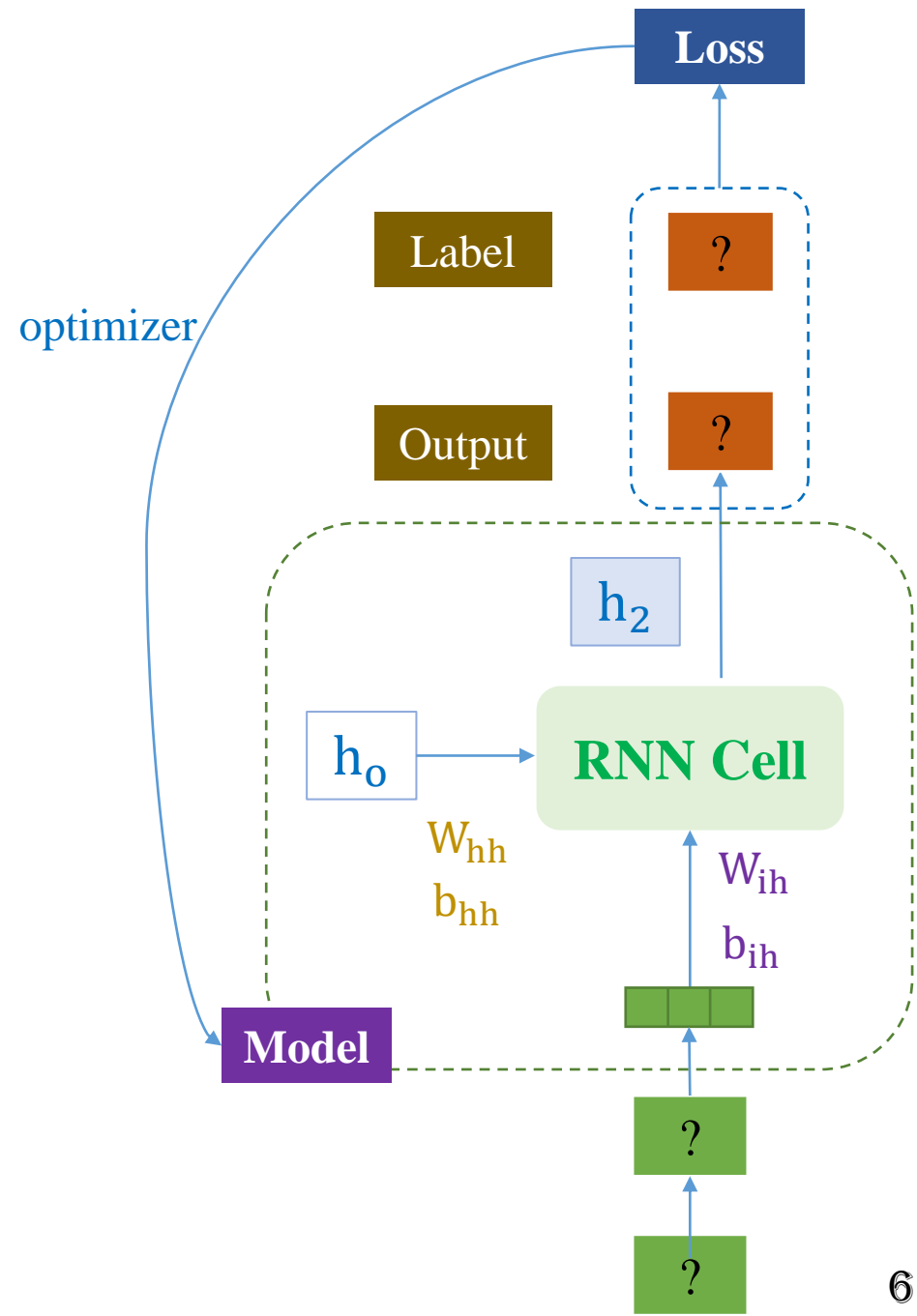
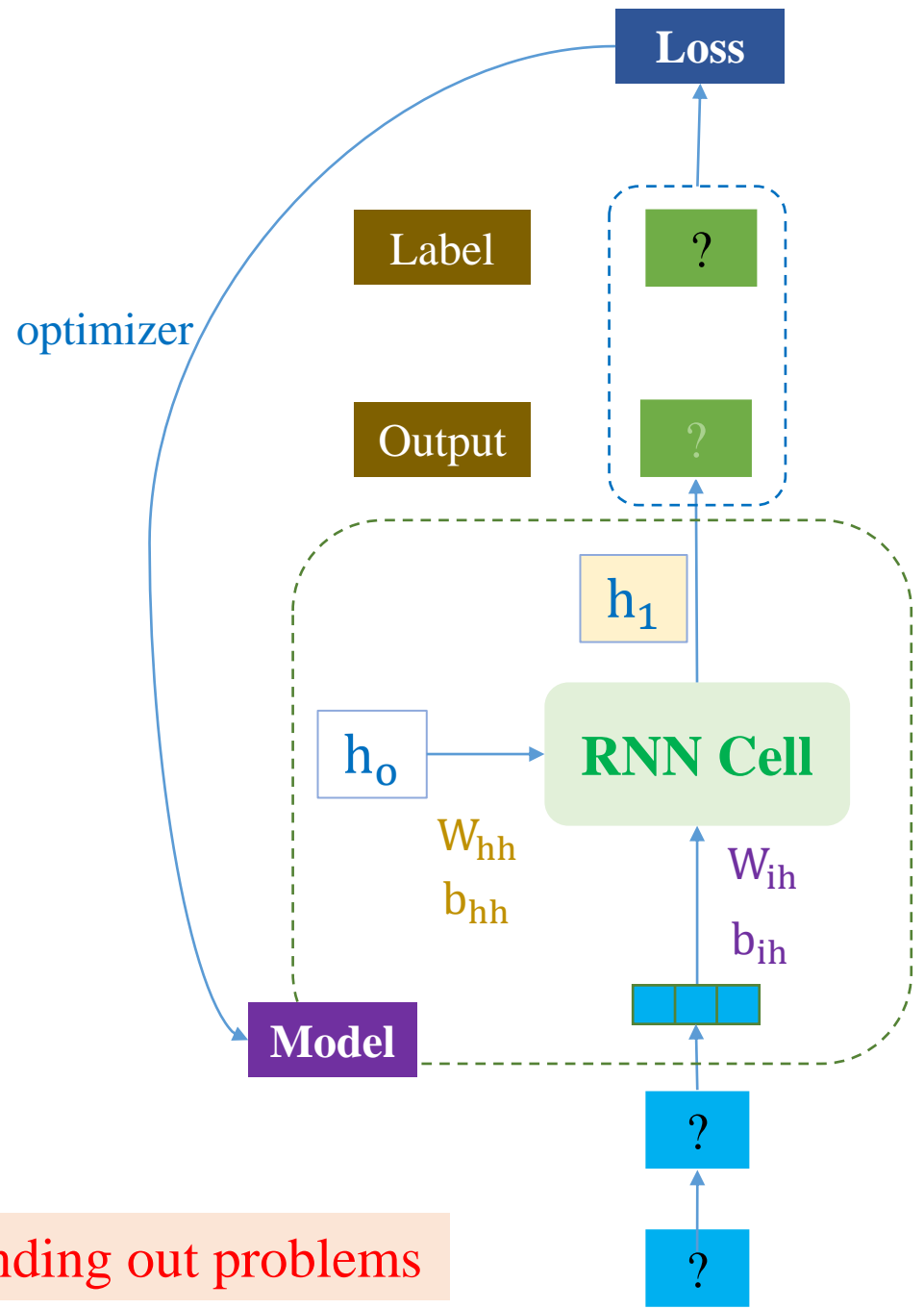


Another Example

data

có
làm
mới
có
ăn

input & label



Discussion and Finding out problems

Text Generation

❖ Input is a set of tokens

```
1 generate_text("thúy kiều", 4, generator_model)
```

'thúy kiều sắc sảo khôn ngoan'

```
1 generate_text("kim trọng và thúy kiều", 6, generator_model)
```

'kim trọng và thúy kiều nhân gian sao tan tành ngay'

data = 'trăm năm trong cõi người ta'

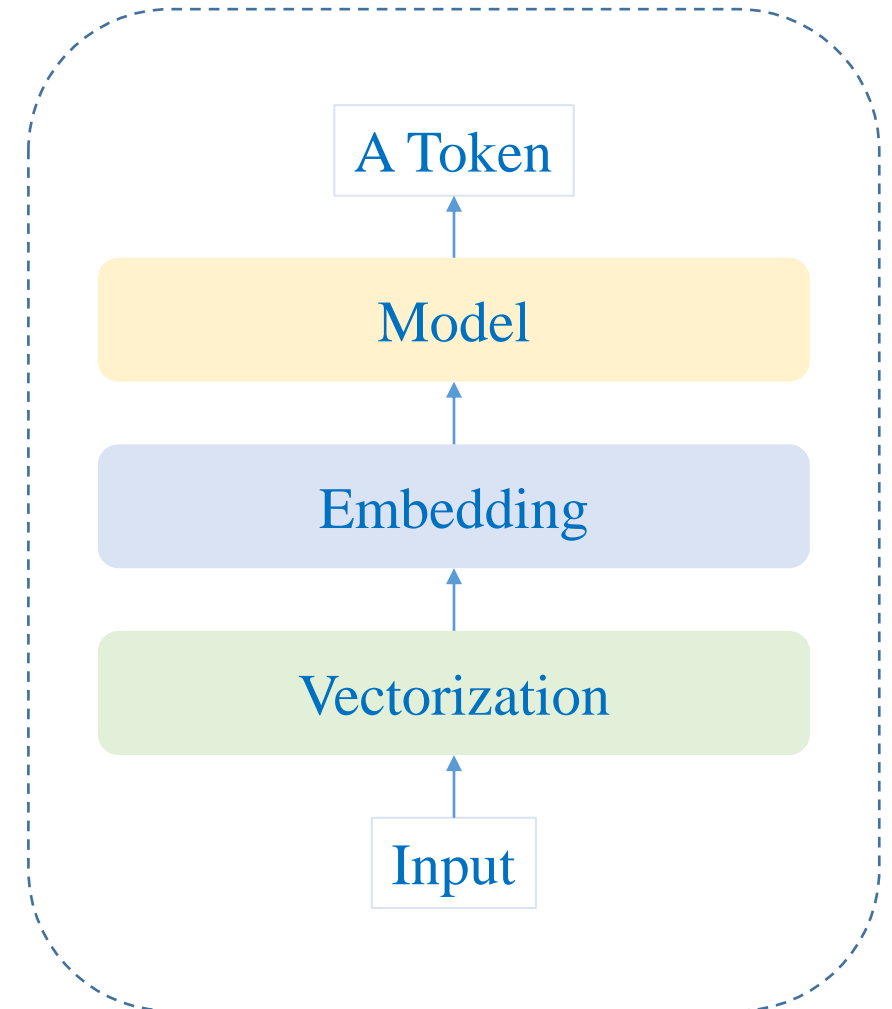
Input data: ['trăm'] => output data: năm

Input data: ['trăm', 'năm'] => output data: trong

Input data: ['trăm', 'năm', 'trong'] => output data: cõi

Input data: ['trăm', 'năm', 'trong', 'cõi'] => output data: người

Input data: ['trăm', 'năm', 'trong', 'cõi', 'người'] => output data: ta

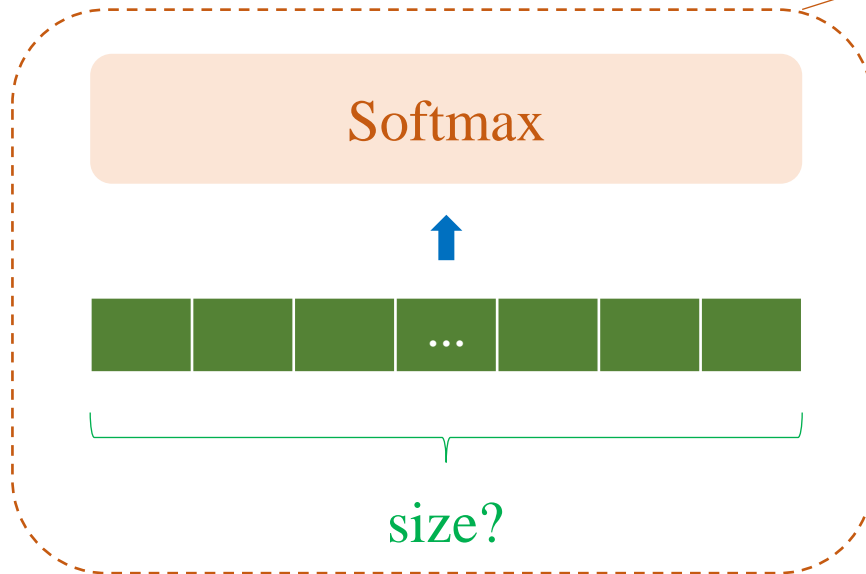


Text Generation

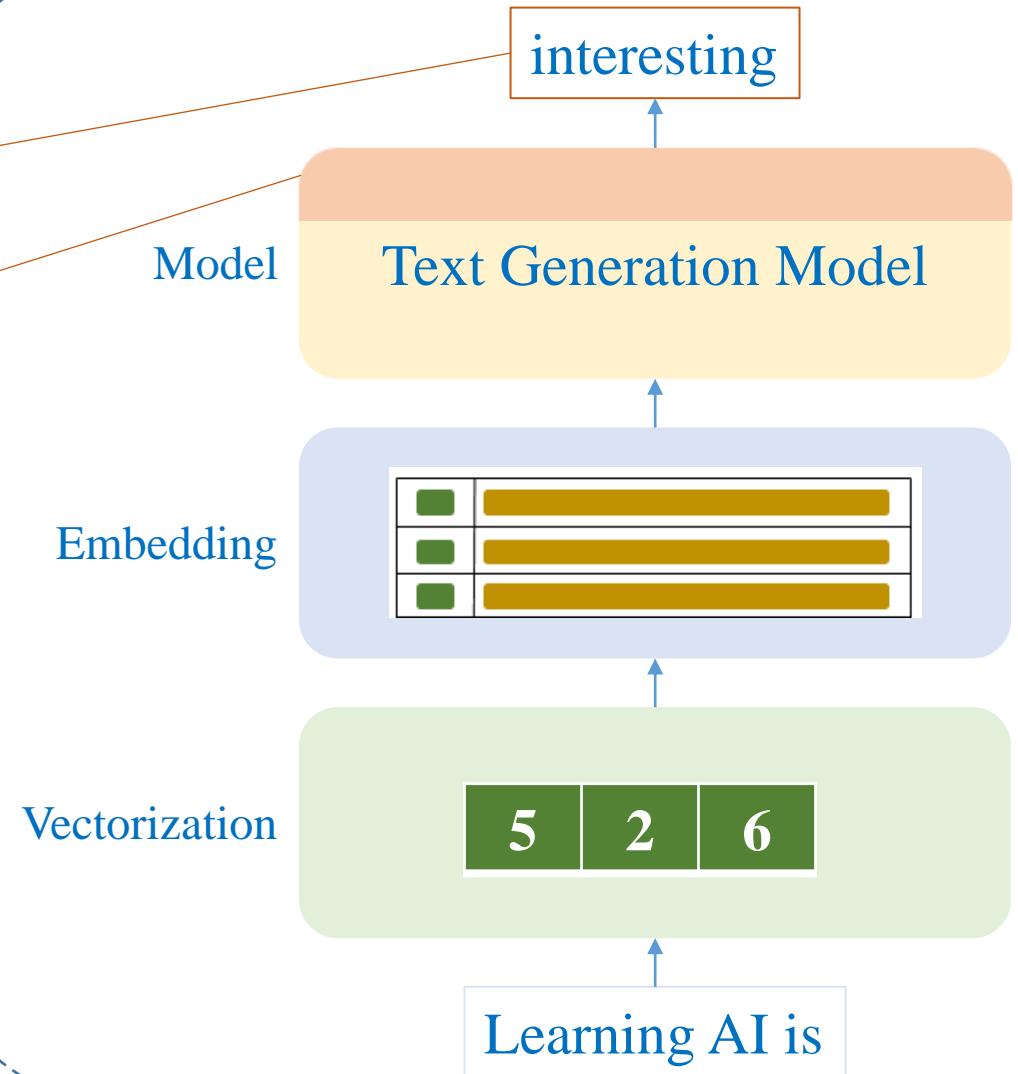
Another Example

Learning AI is interesting
AI is a CS Topic

Where is “interesting” from?

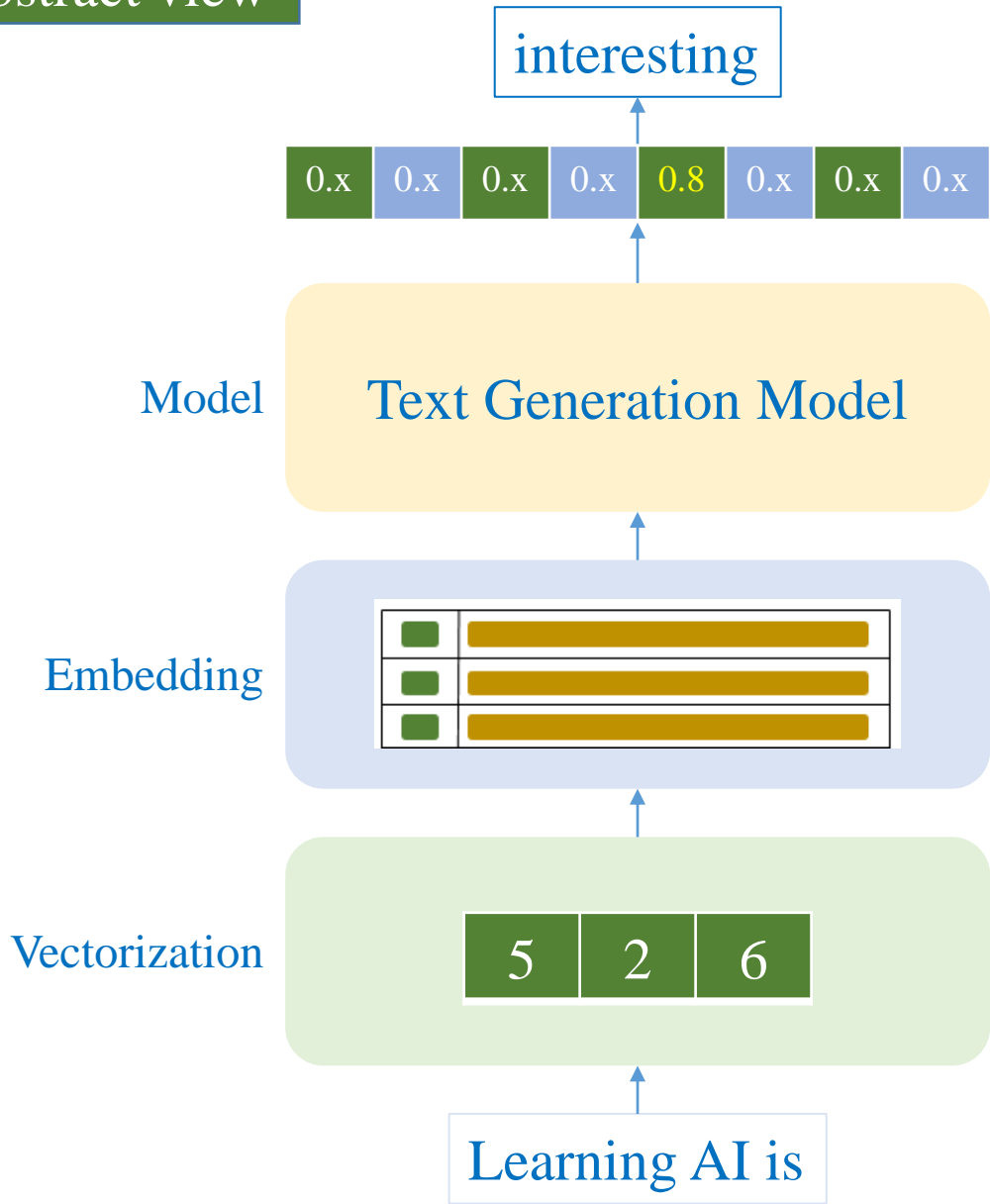


size = vocab_size

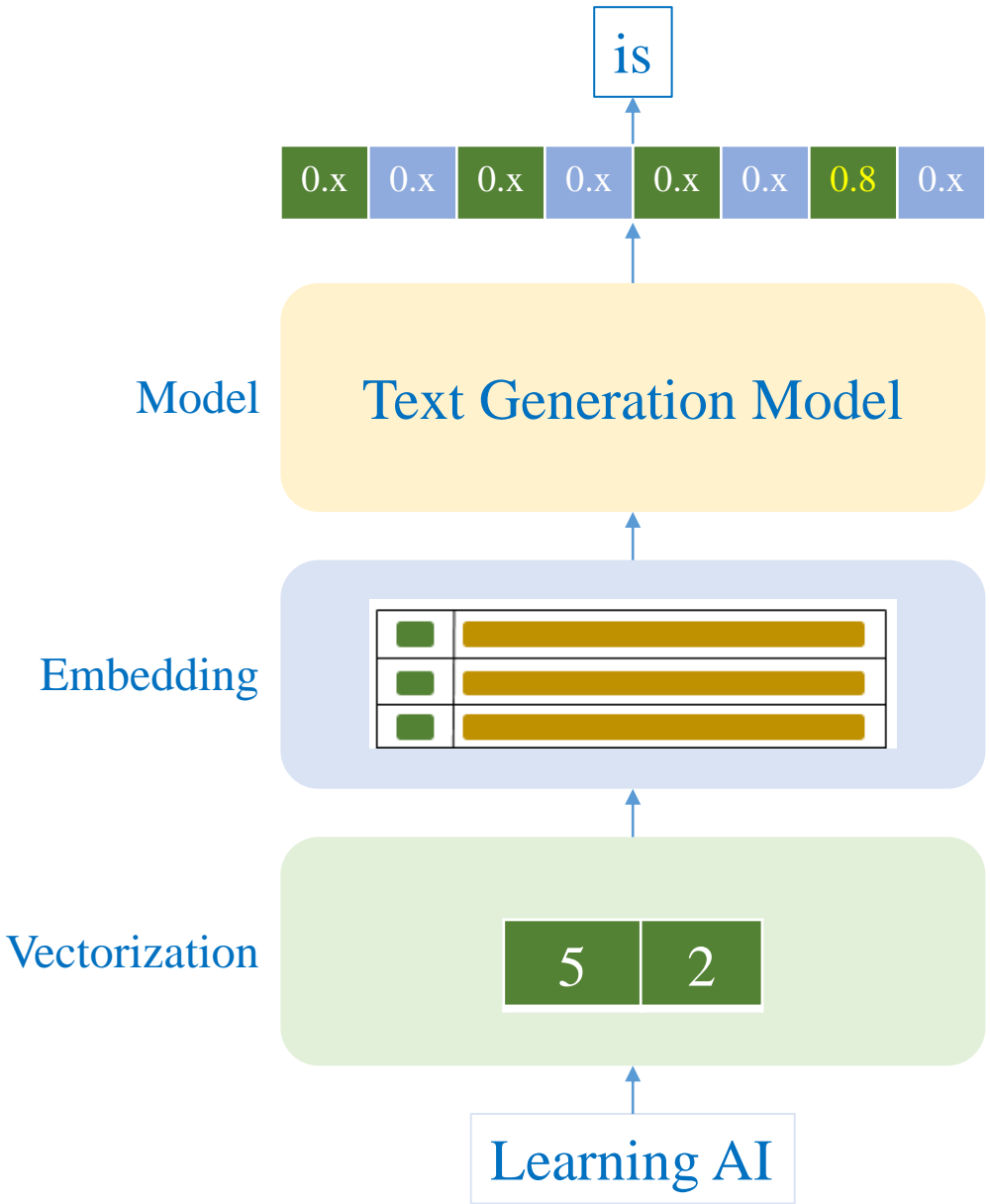


Text Generation

Abstract view



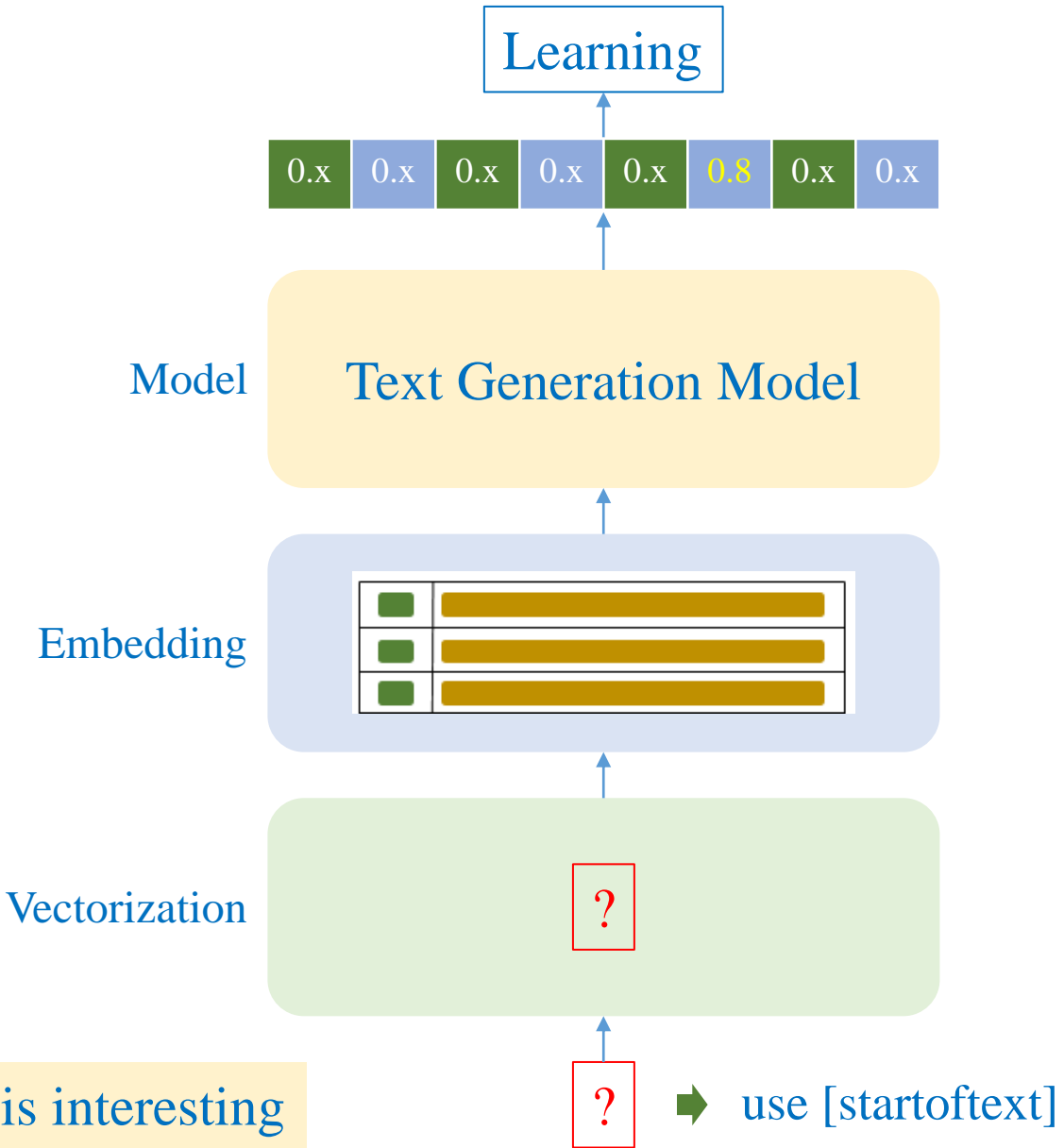
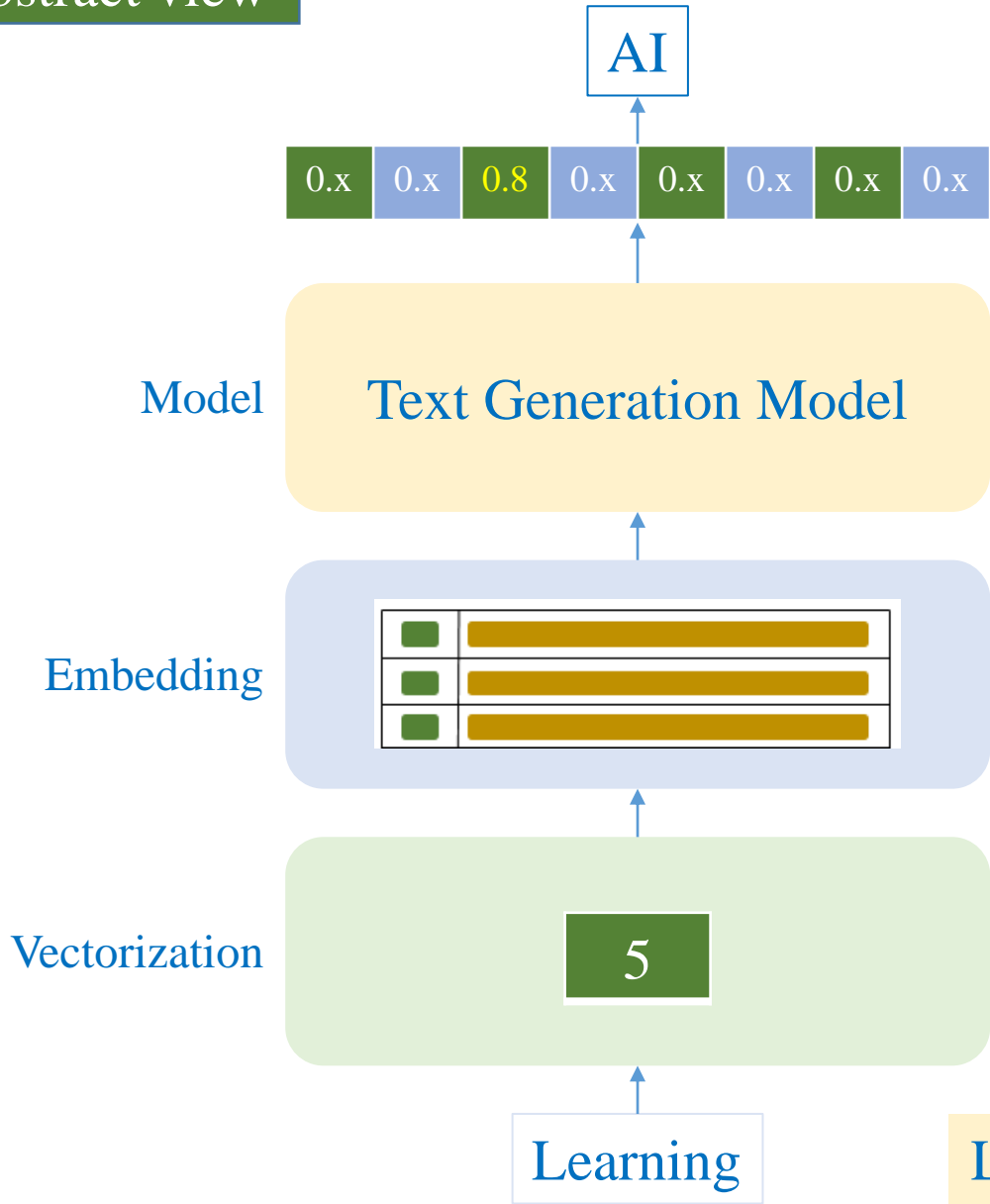
index	0	1	2	3	4	5	6	7
word	pad	[UNK]	ai	we	interesting	learning	is	cs



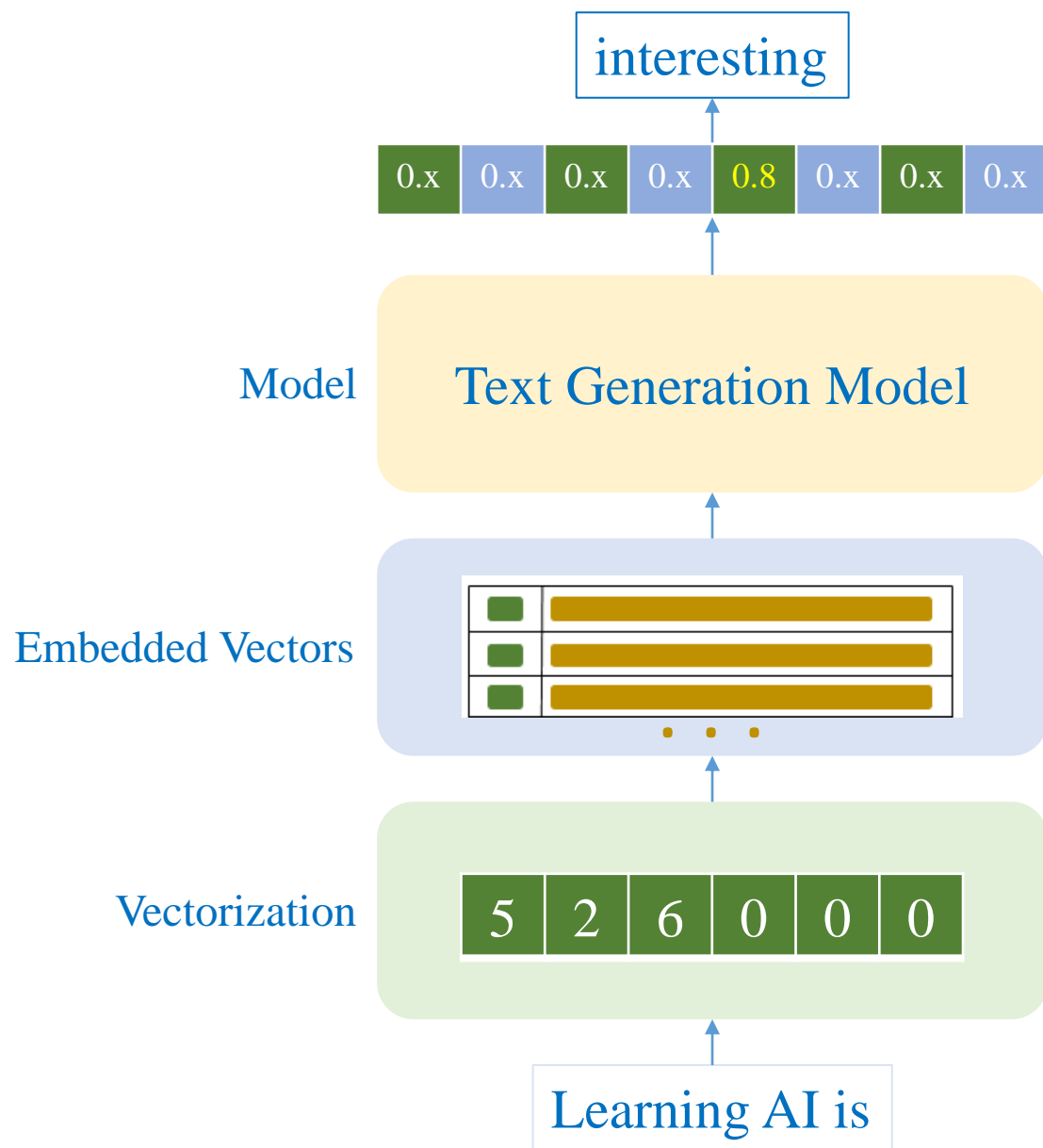
Text Generation

Abstract view

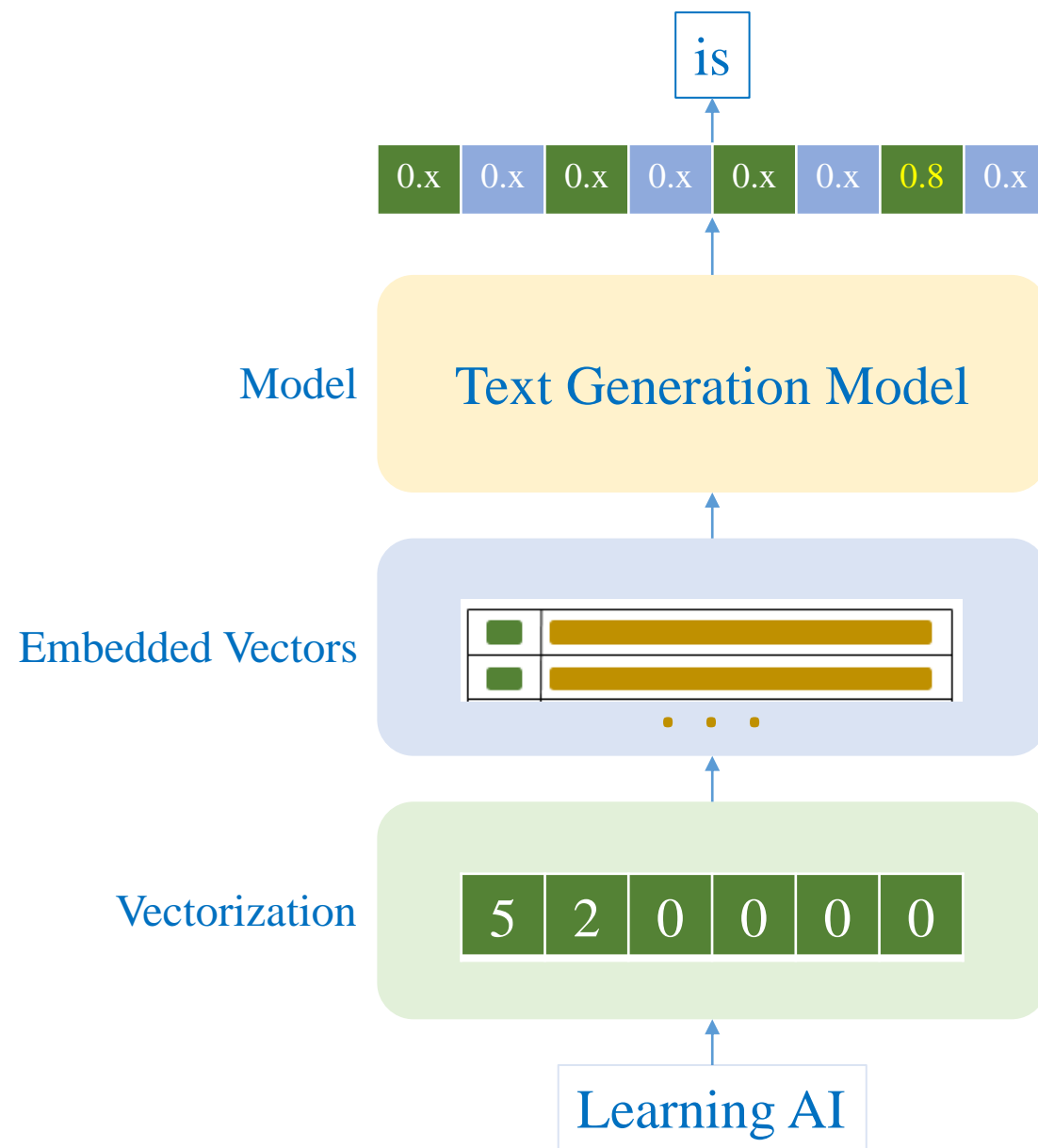
index	0	1	2	3	4	5	6	7
word	pad	[UNK]	ai	we	interesting	learning	is	cs



sequence_length = 6



index	0	1	2	3	4	5	6	7
word	pad	[UNK]	ai	we	interesting	learning	is	cs



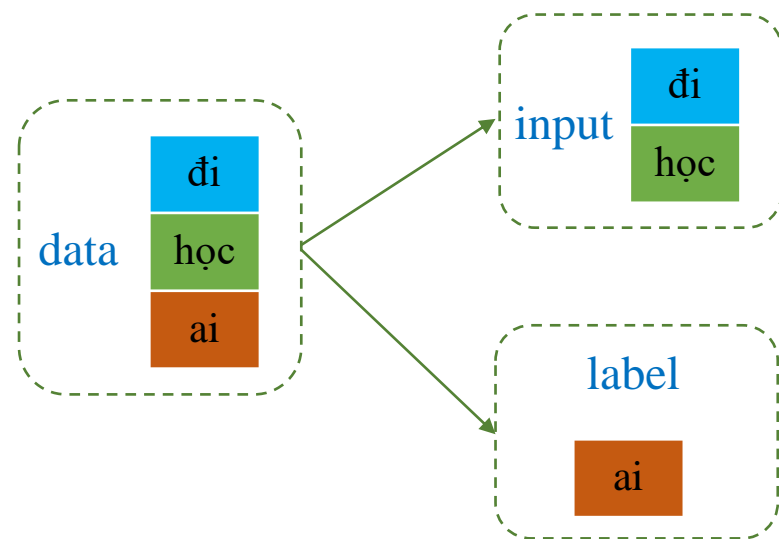
Outline

- **Motivation to Text Generation**
- **Simple Models**
- **Using RNNs**
- **Examples**
- **Using Masked Encoder (~Decoder)**

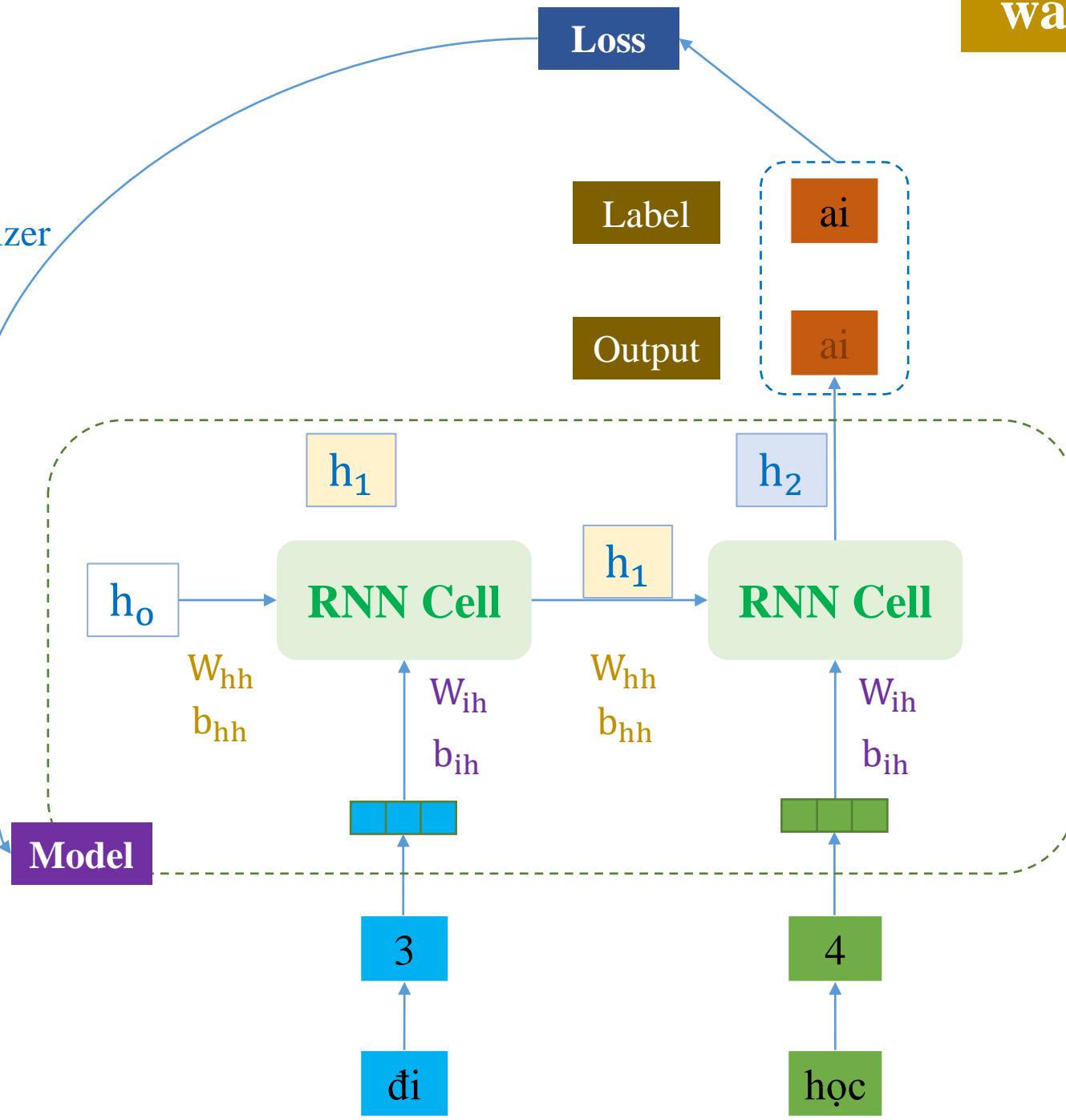
Implementation Using RNN

way 1

Extract context from input



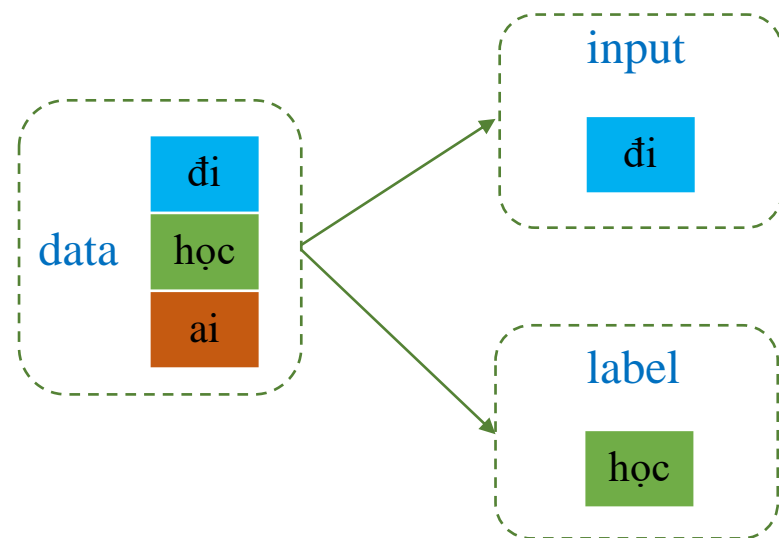
optimizer



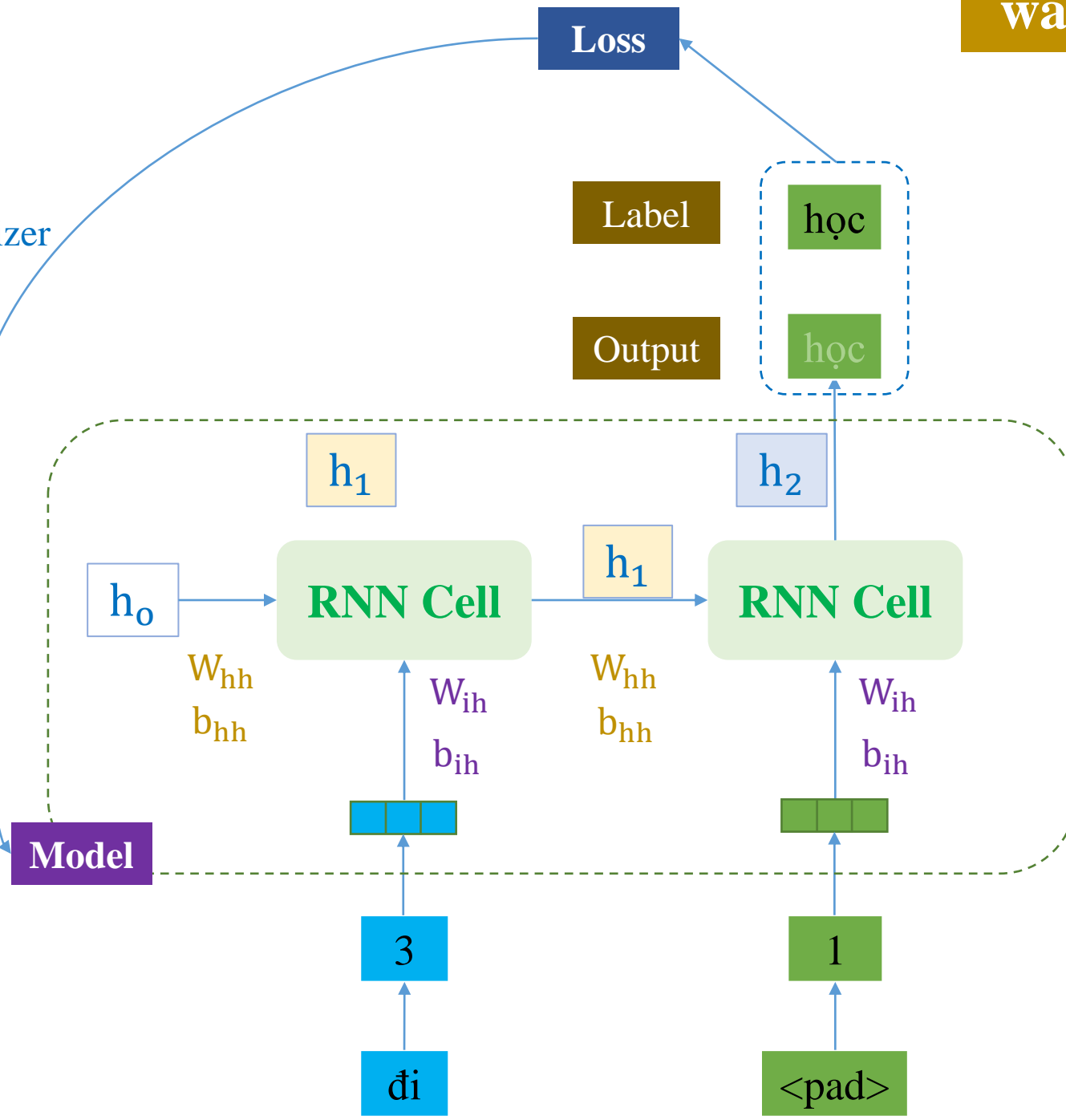
Implementation Using RNN

way 1

Extract context from input



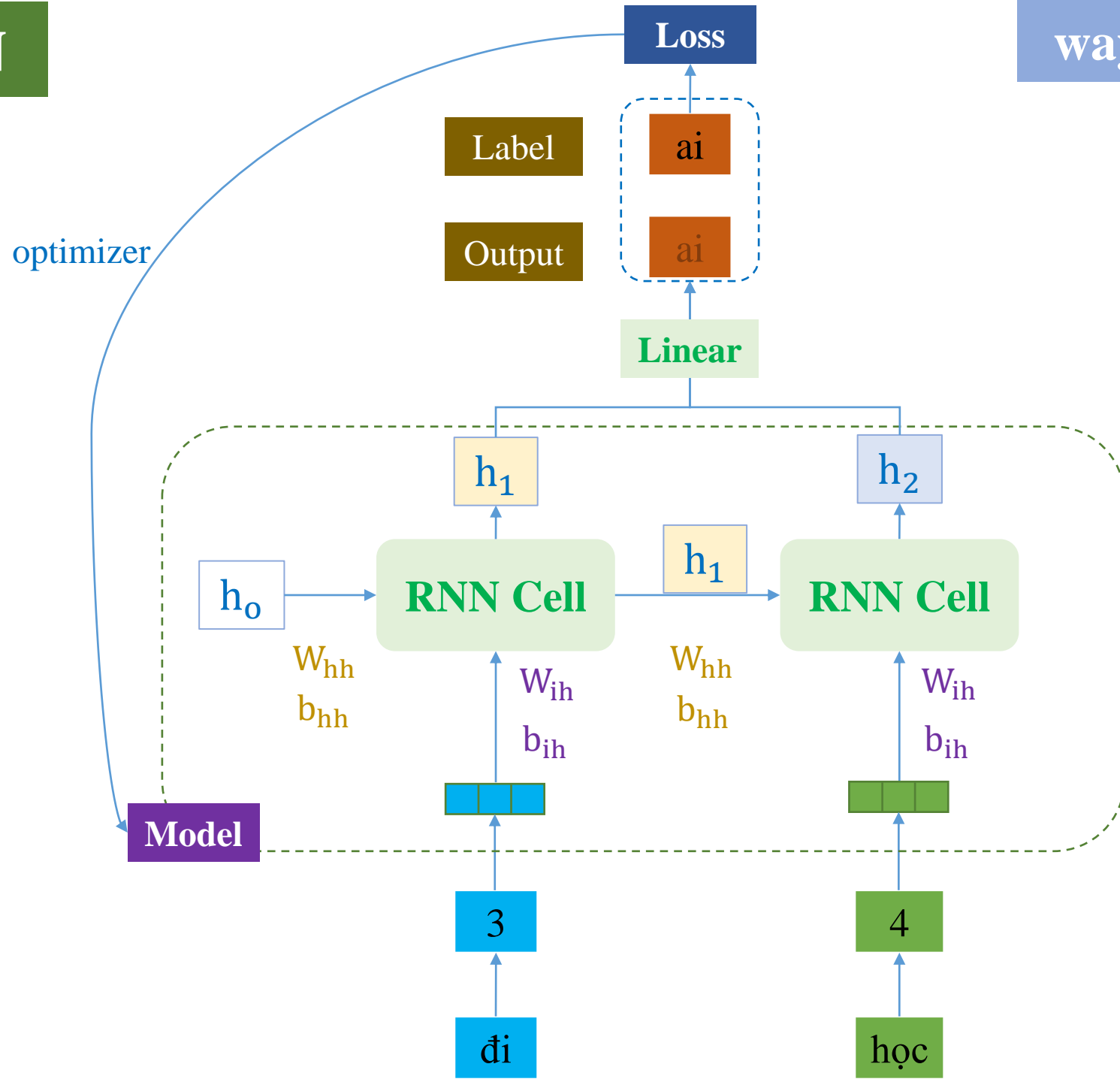
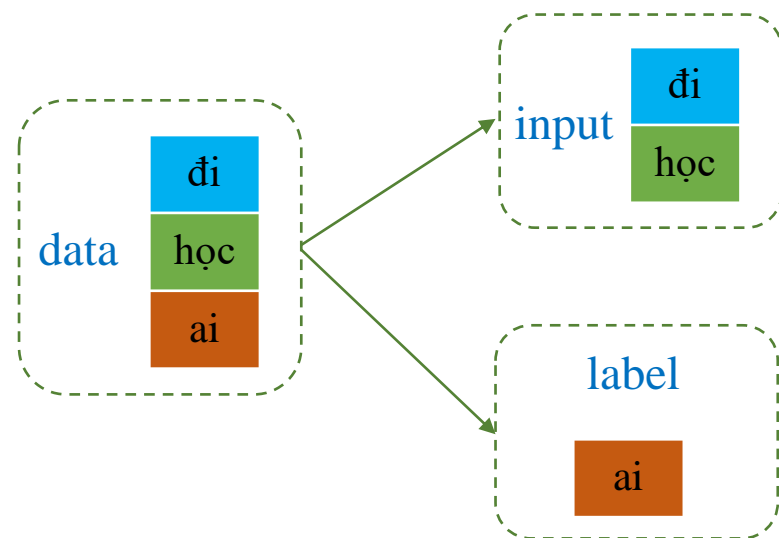
optimizer



Implementation Using RNN

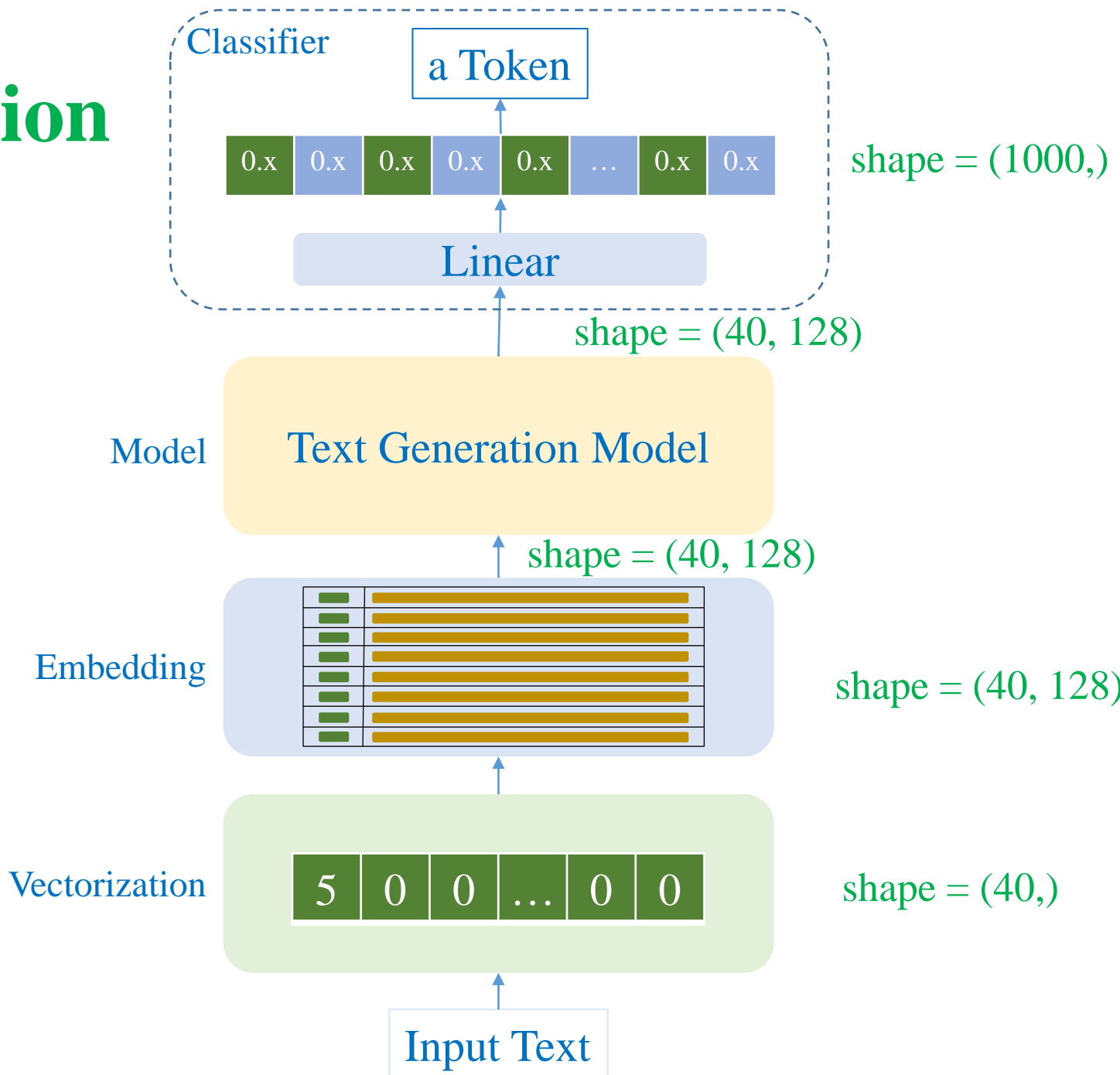
way 2

Using all the features



Text Generation

❖ Practice

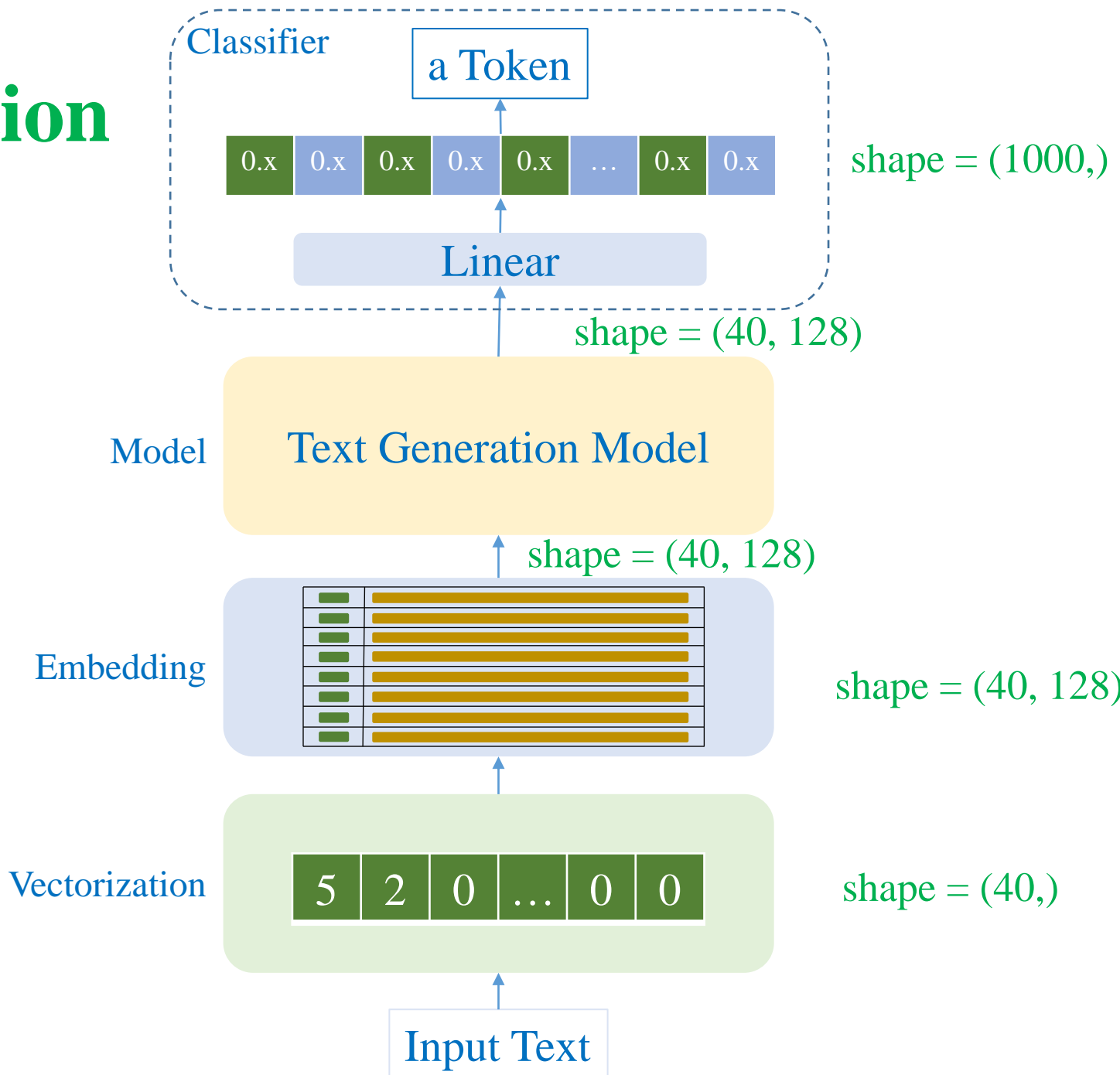


sequence_length = 40
embed_dim = 128
vocab_size = 1000

‘Learning AI is interesting’
 x = ‘Learning’
 y = ‘AI’

Text Generation

❖ Practice

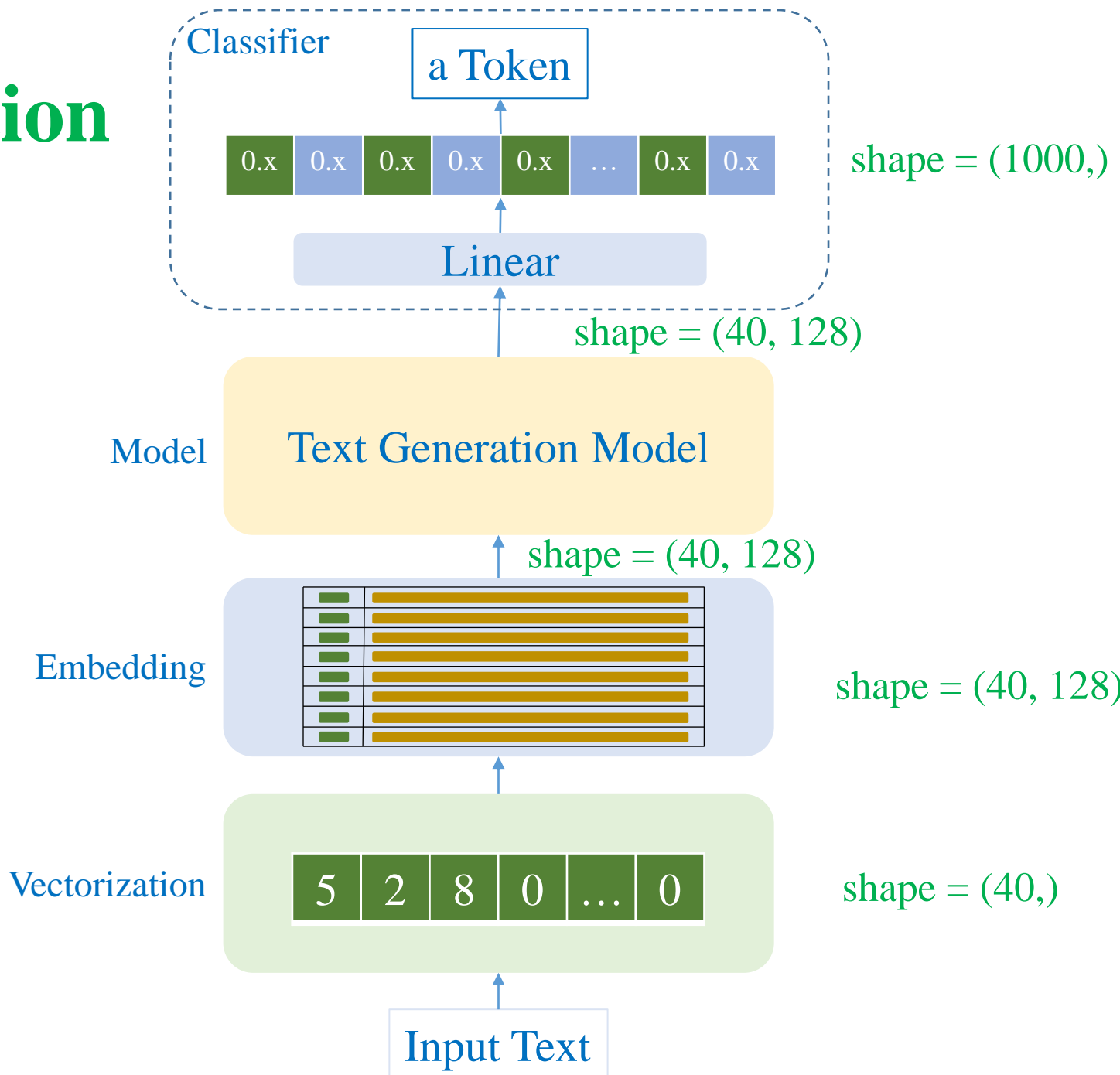


sequence_length = 40
embed_dim = 128
vocab_size = 1000

‘Learning AI is interesting’
 x = ‘Learning AI’
 y = ‘is’

Text Generation

❖ Practice

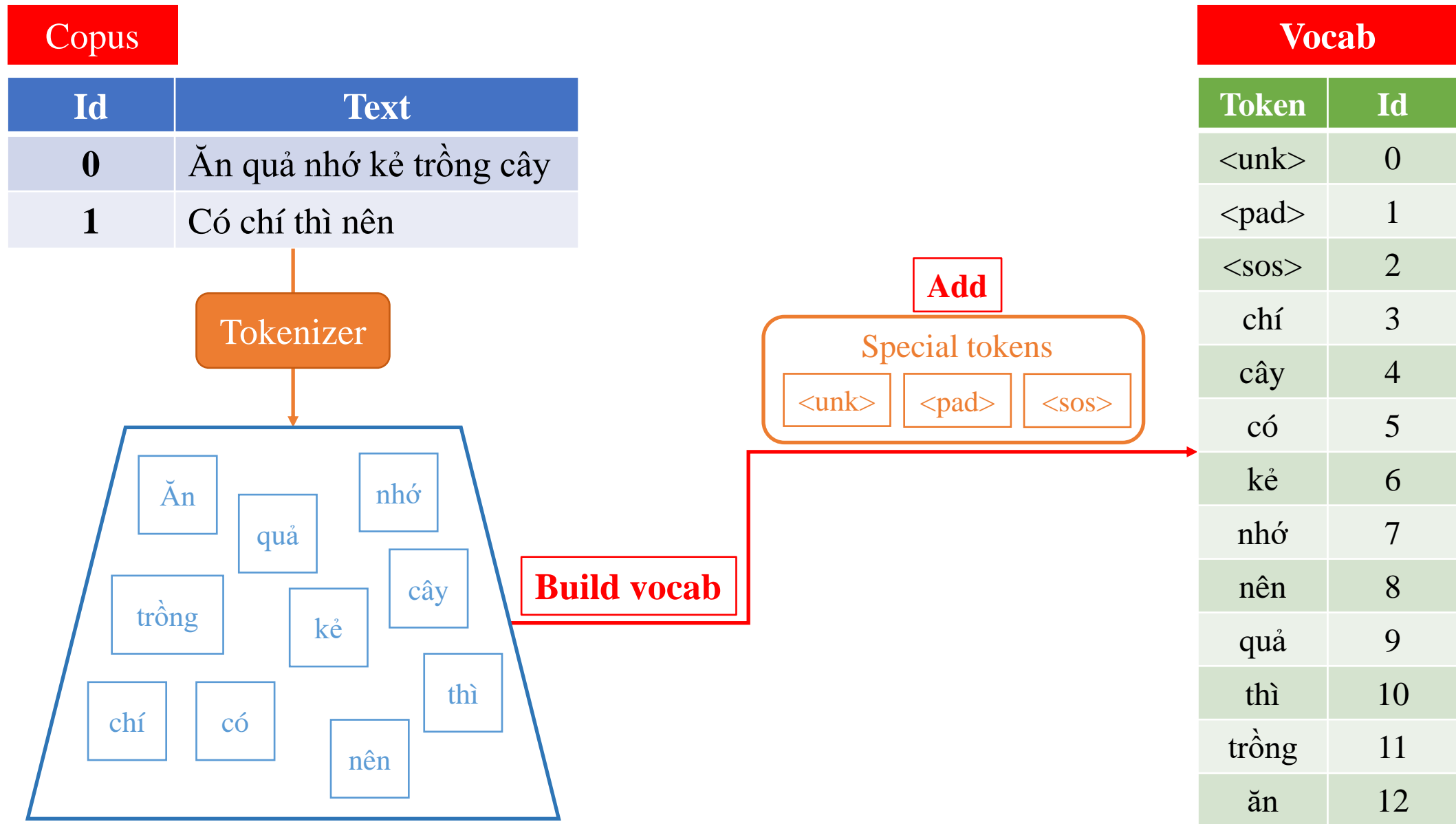


sequence_length = 40
embed_dim = 128
vocab_size = 1000

‘Learning AI is interesting’
 x = ‘Learning AI is’
 y = ‘interesting’

Loss function?

Example



Example

Id	Text
0	Ăn quả nhớ kẻ trồng cây
1	Có chí thì nên



Input tokens	Target token
<sos>	Ăn
<sos> Ăn	quả
<sos> Ăn quả	nhớ
<sos> Ăn quả nhớ	kẻ
<sos> Ăn quả nhớ kẻ	trồng
<sos> Ăn quả nhớ kẻ trồng	cây
<sos>	Có
<sos> Có	chí
<sos> Có chí	thì
<sos> Có chí thì	nên

Token	Id	Token	Id
<unk>	0	kẻ	6
<pad>	1	nhớ	7
<sos>	2	nên	8
chí	3	quả	9
cây	4	thì	10
có	5	trồng	11
		ăn	12

Vocab

padding

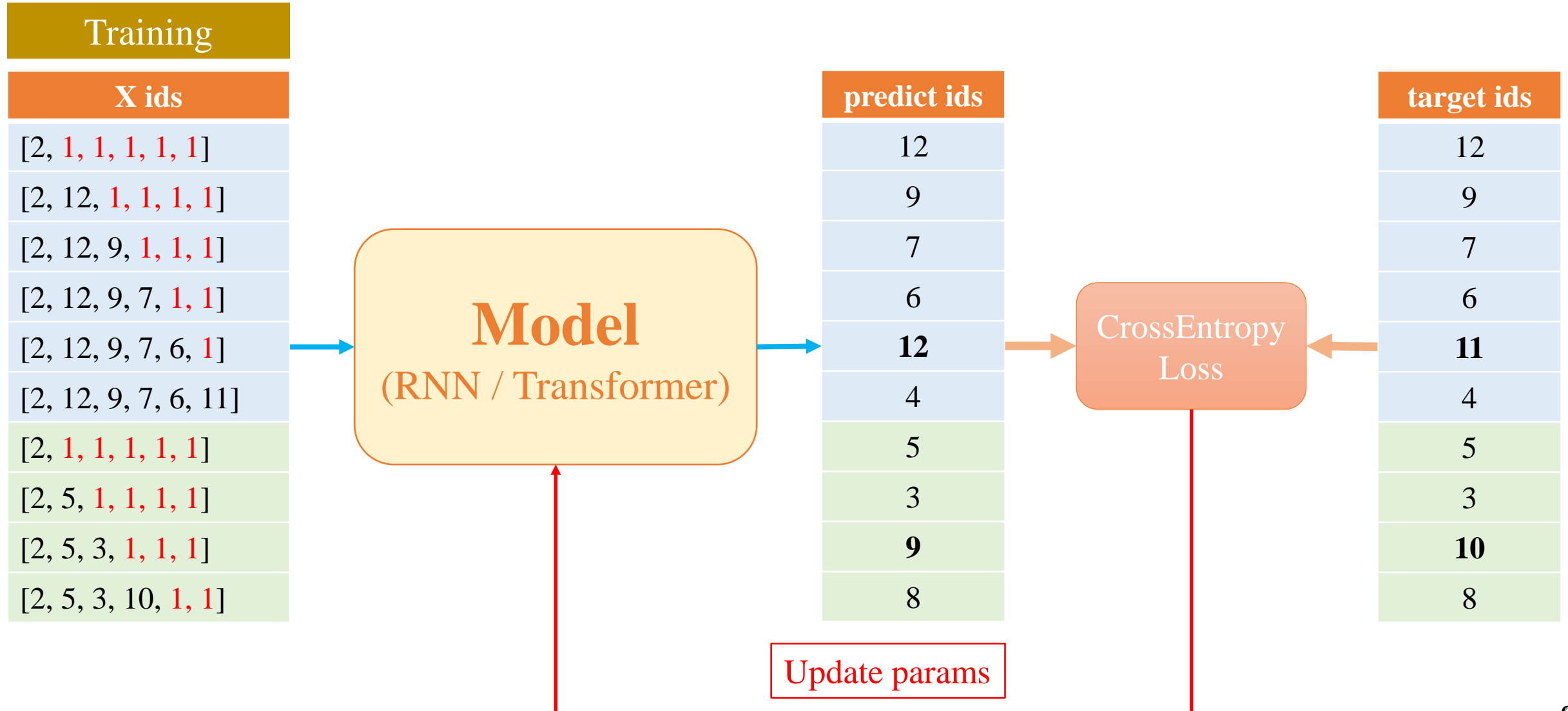
Input ids	Target ids
[2, 1, 1, 1, 1, 1]	12
[2, 12, 1, 1, 1, 1]	9
[2, 12, 9, 1, 1, 1]	7
[2, 12, 9, 7, 1, 1]	6
[2, 12, 9, 7, 6, 1]	11
[2, 12, 9, 7, 6, 11]	4
[2, 1, 1, 1, 1, 1]	5
[2, 5, 1, 1, 1, 1]	3
[2, 5, 3, 1, 1, 1]	10
[2, 5, 3, 10, 1, 1]	8

Training data

Next token prediction dataset

sequence_length = 6

Example



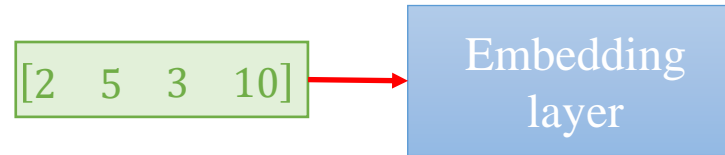
sequence_length = 4

hidden_dim = 4

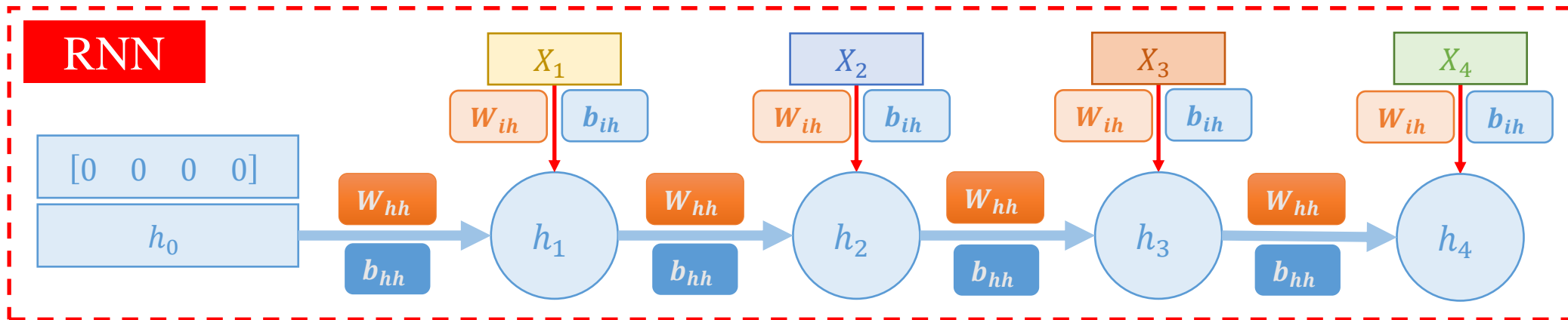
Feed Forward

Input tokens	Target token
<sos> Có chí thì	nên

Input ids	Target ids
[2, 5, 3, 10]	[8]



$X_1 = [-0.7521 \quad 1.6487 \quad -0.3925 \quad -1.4036]$
$X_2 = [-0.7581 \quad 1.0783 \quad 0.8008 \quad 1.6806]$
$X_3 = [-0.7279 \quad -0.5594 \quad -0.7688 \quad 0.7624]$
$X_4 = [-0.8371 \quad -0.9224 \quad 1.8113 \quad 0.1606]$



W_{hh}			
0.3035	-0.1187	0.2860	-0.3885
-0.2523	0.1524	0.1057	-0.1275
0.2980	0.3399	-0.3626	-0.2669
0.4578	-0.1687	-0.1773	-0.4838

W_{ih}			
-0.2982	-0.2982	0.4497	0.1666
0.4811	-0.4126	-0.4959	-0.3912
-0.3363	0.2025	0.1790	0.4155
-0.2582	-0.3409	0.2653	-0.2021

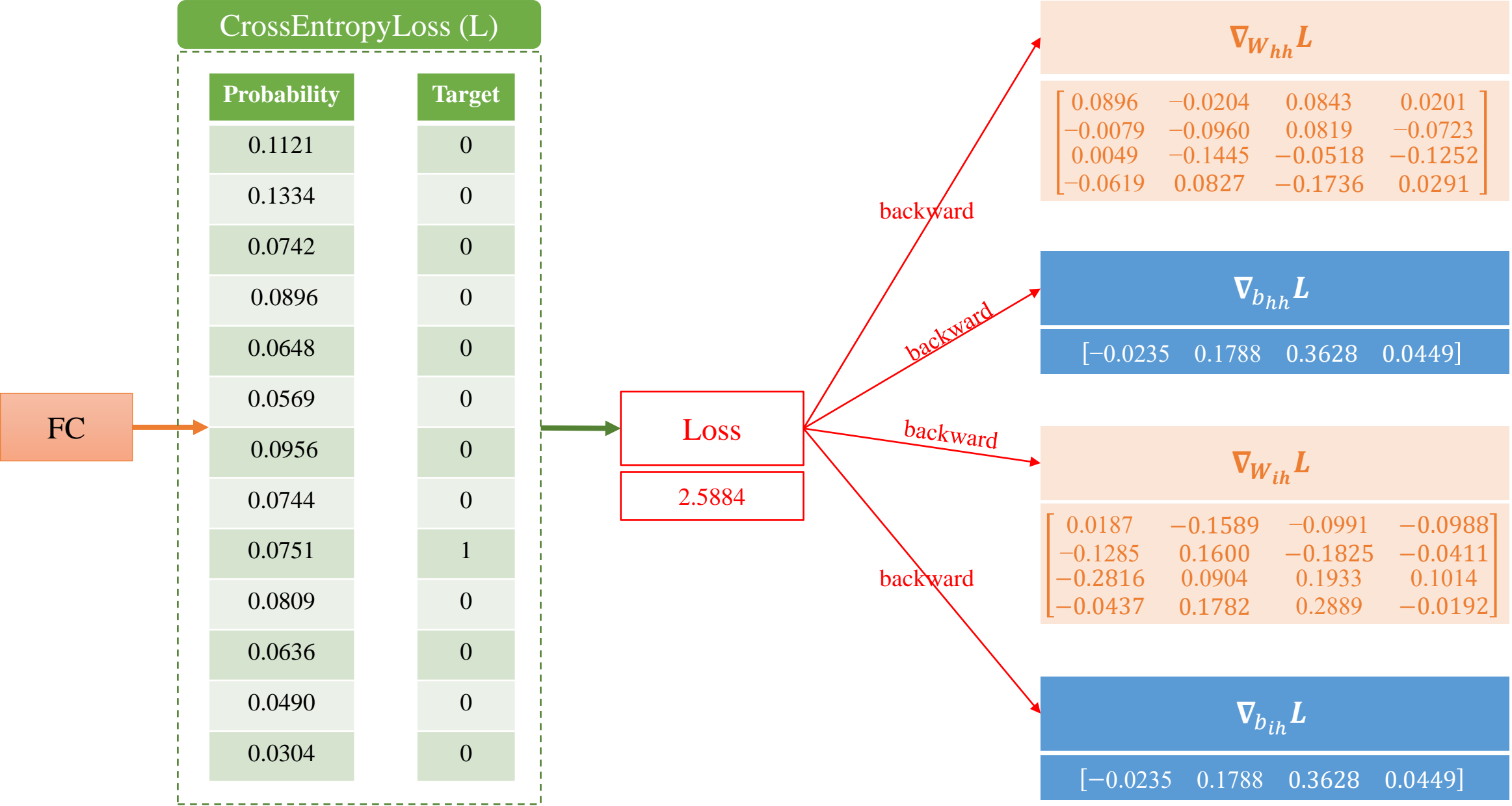
b_{hh}			
0.0117	-0.3415	-0.4242	-0.2753

b_{ih}			
-0.2863	0.1249	-0.0660	-0.3629

$h_1 = [-0.7409 \quad -0.4739 \quad -0.5055 \quad -0.6786]$
$h_2 = [0.2170 \quad -0.9590 \quad 0.6681 \quad -0.6519]$
$h_3 = [0.4735 \quad -0.2915 \quad -0.4692 \quad -0.1583]$
$h_4 = [0.8327 \quad -0.8839 \quad 0.2449 \quad 0.6446]$



Back-Propagation



Update Parameters
SGD($lr = 0.1$)

$$W_{hh} = W_{hh} - lr \times dW_{hh}$$

$$\begin{bmatrix} 0.2945 & -0.1166 & 0.2776 & -0.3905 \\ -0.2515 & 0.1620 & 0.0975 & -0.1203 \\ 0.2975 & 0.3544 & -0.3574 & -0.2544 \\ 0.4640 & -0.1770 & -0.1599 & -0.4867 \end{bmatrix}$$

$$b_{hh} = b_{hh} - lr \times db_{hh}$$

$$[0.0141 \quad -0.3594 \quad -0.4605 \quad -0.2798]$$

$$W_{ih} = W_{ih} - lr \times dW_{ih}$$

$$\begin{bmatrix} -0.3001 & -0.2823 & 0.4596 & 0.1765 \\ 0.4940 & -0.4286 & -0.4777 & -0.3871 \\ -0.3082 & 0.1935 & 0.1597 & 0.4053 \\ -0.2538 & -0.3587 & 0.2364 & -0.2002 \end{bmatrix}$$

$$b_{ih} = b_{ih} - lr \times db_{ih}$$

$$[-0.2840 \quad 0.1070 \quad -0.1023 \quad -0.3674]$$

Forward again

CrossEntropyLoss (L)

Probability	Target
0.1044	0
0.1163	0
0.0659	0
0.0865	0
0.0575	0
0.0521	0
0.0894	0
0.0673	0
0.1506	1
0.0754	0
0.0578	0
0.0488	0
0.0280	0

Loss

1.8929

Outline

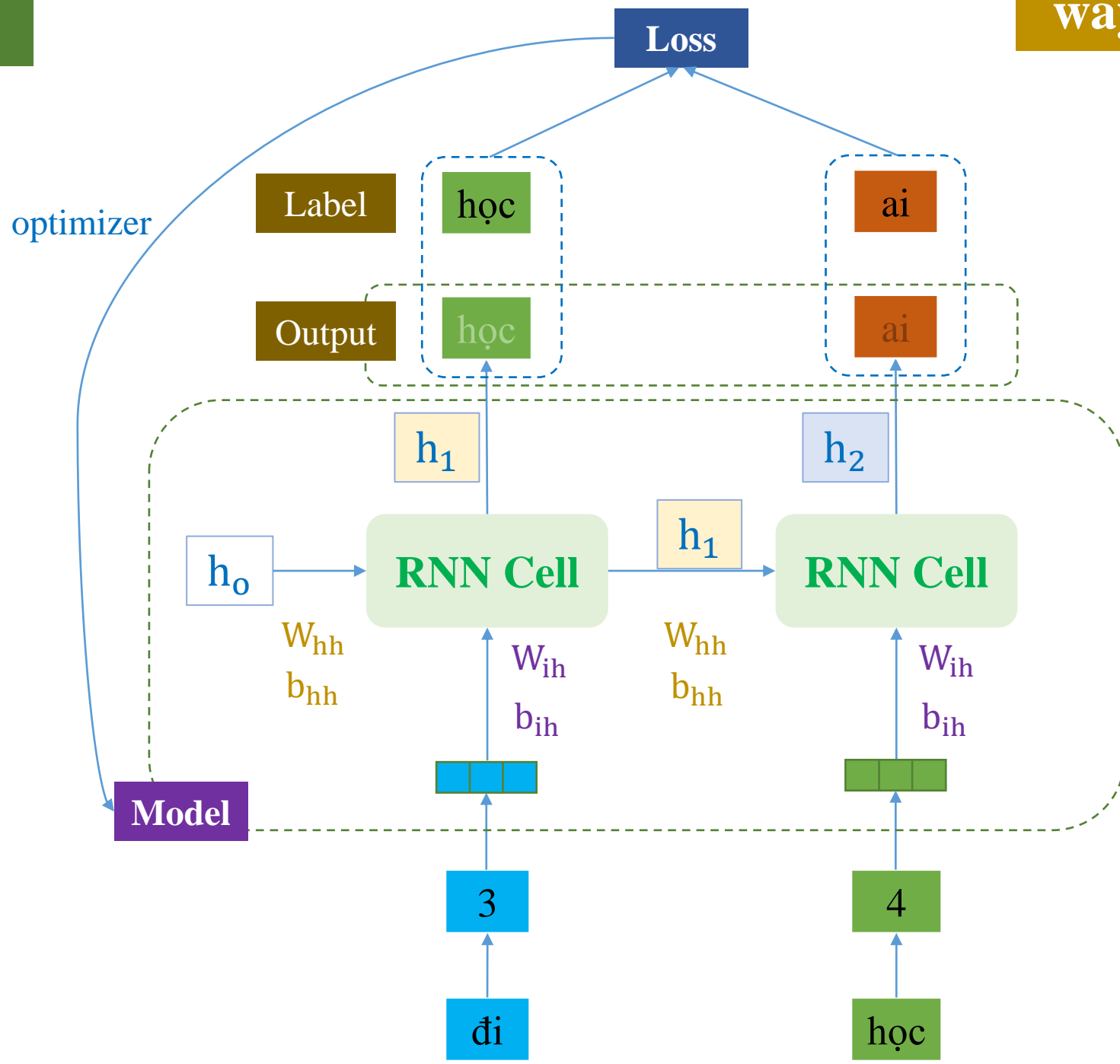
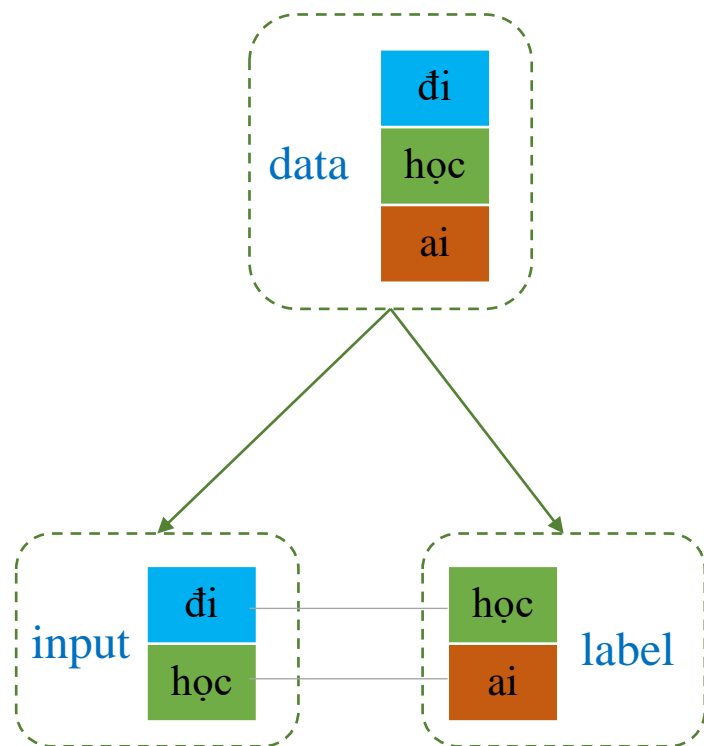
- **Motivation to Text Generation**
- **Simple Models**
- **Using RNNs**
- **Examples**
- **Using Masked Encoder (~Decoder)**

Implementation Using RNN

way 3

RNNs encode input sequentially

*Linear can be added (where?)



Implementation Using RNN

RNNs encode input sequentially

way 3

hidden_dim = vocab_size = 4

[-0.88, 0.81, 0.77]
[-1.27, 0.84, 0.04]

Embedded Input

W_{ih}

[0.25, 0.05, 0.17]
[-0.11, 0.35, 0.22]
[0.44, -0.22, -0.42]
[0.38, 0.36, -0.01]

b_{ih}

[-0.47, 0.13, 0.46, -0.15]

W_{hh}

[0.38, -0.06, 0.37, -0.19]
[-0.16, 0.42, 0.08, -0.02]
[-0.36, -0.01, 0.17, 0.39]
[-0.43, 0.09, 0.33, 0.03]

b_{hh}

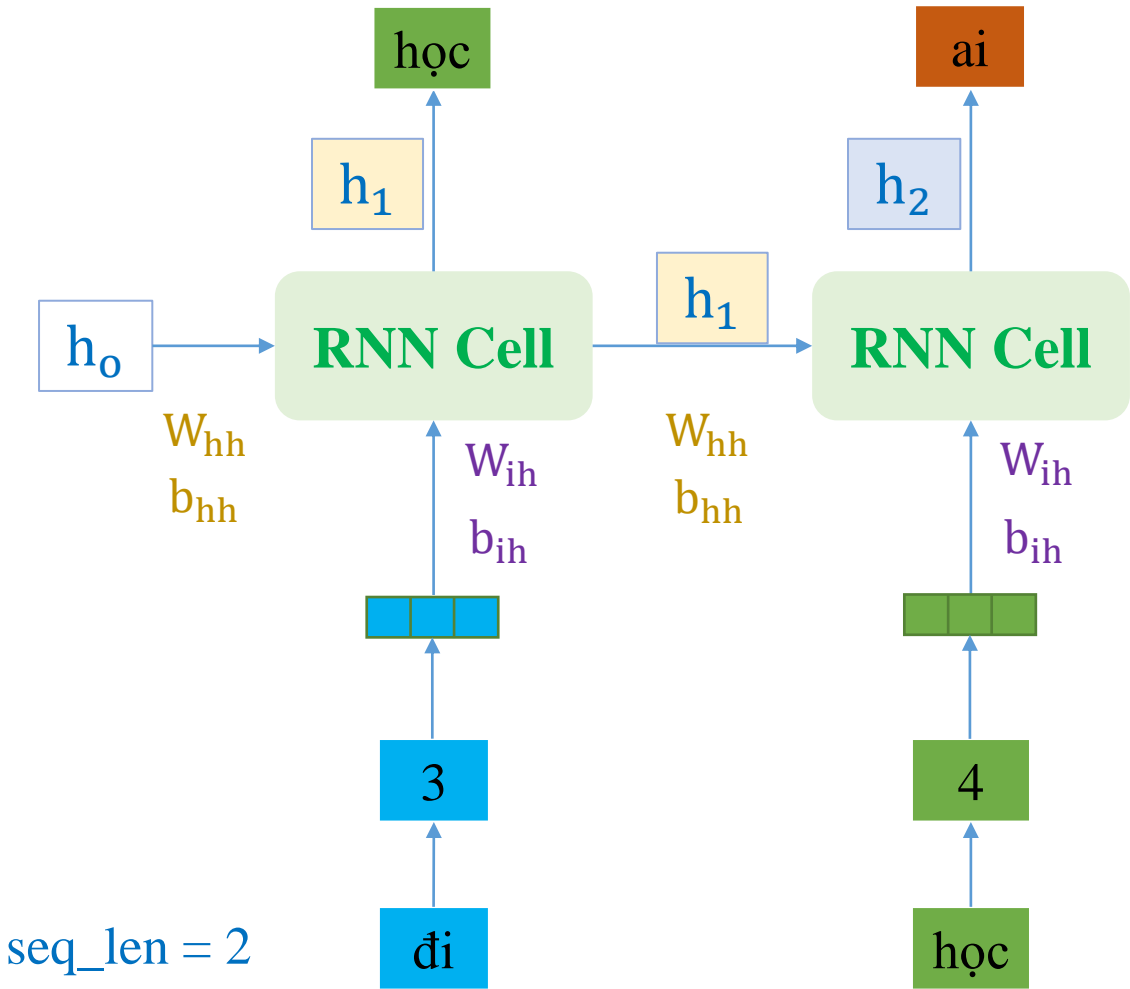
[0.43, 0.16, 0.34, -0.39]

Output

[-0.08, 0.68, -0.08, -0.53]
[-0.30, 0.77, -0.15, -0.58]

Hidden

[-0.30, 0.77, -0.15, -0.58]

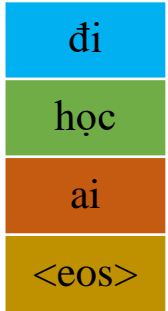


seq_len = 2

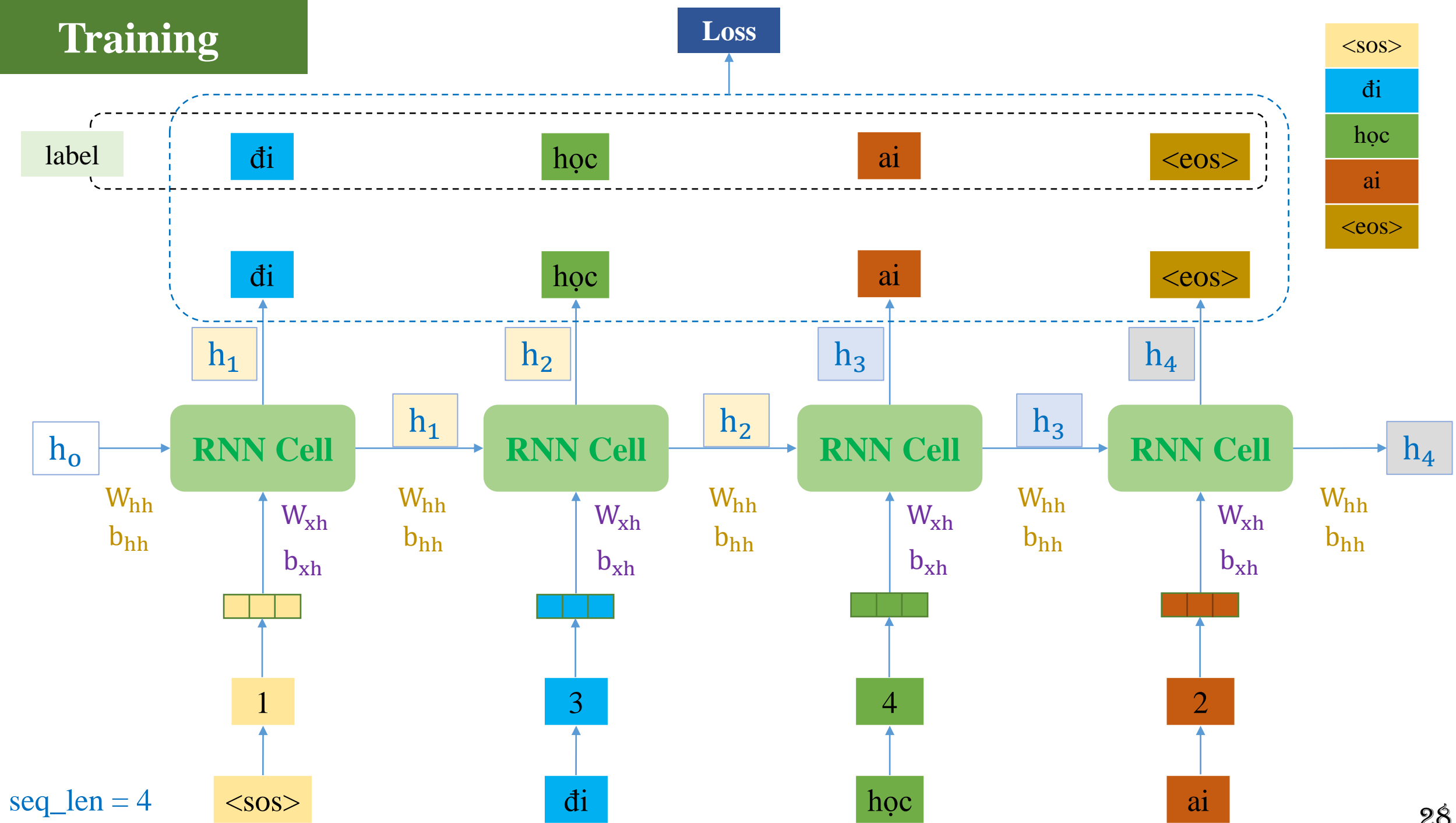
*Linear can be added

Training

Add more special tokens

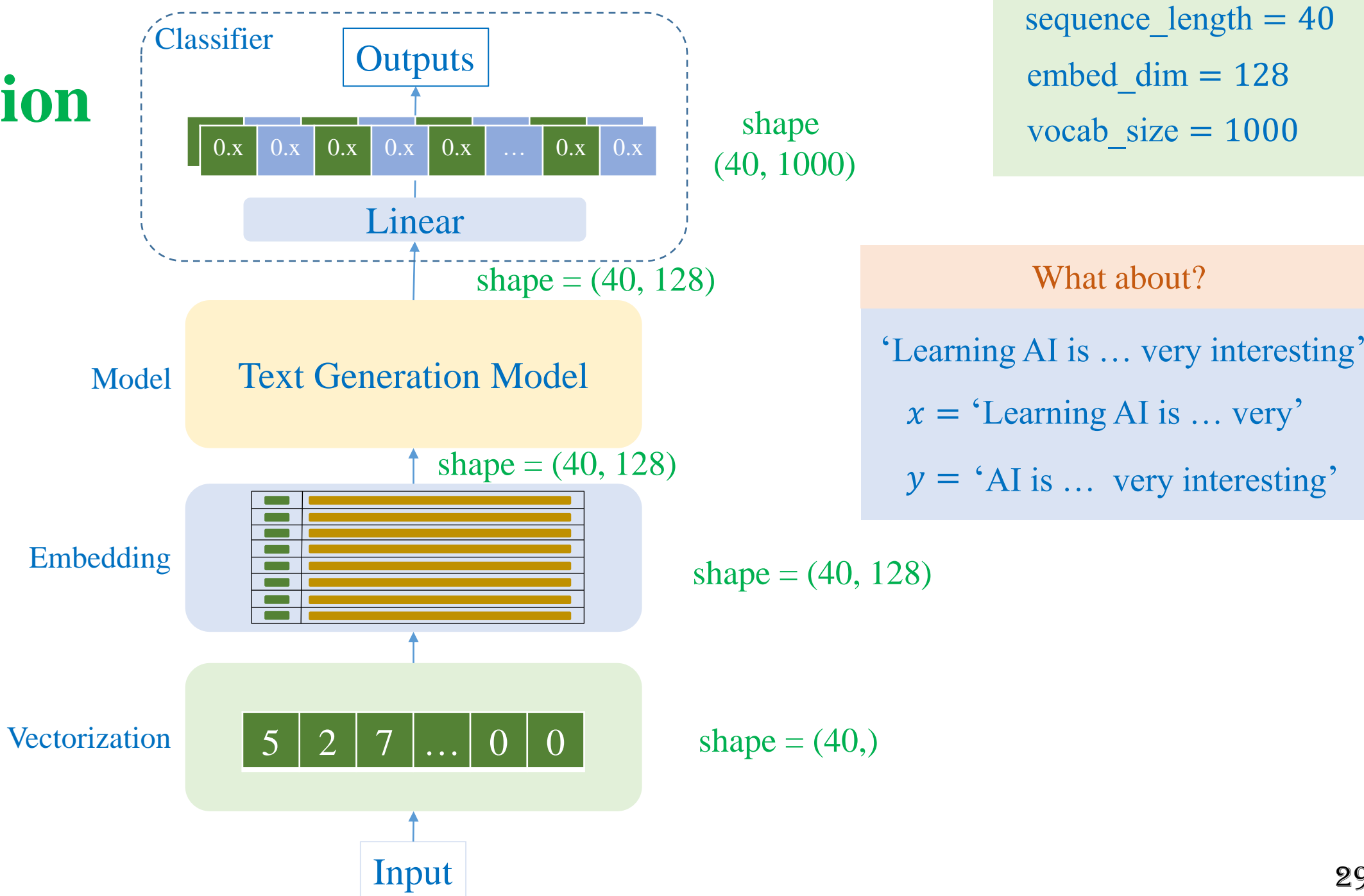


Training



Text Generation

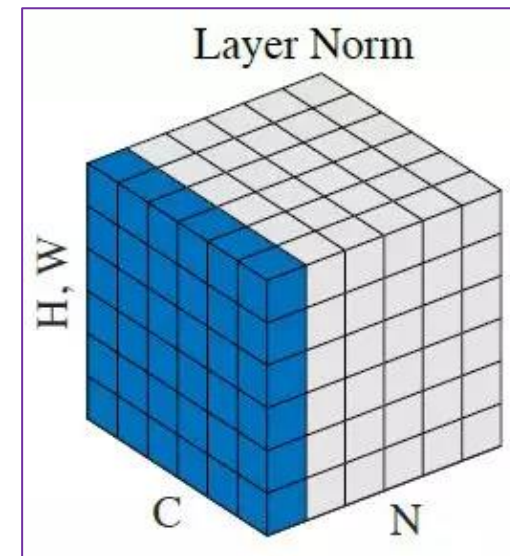
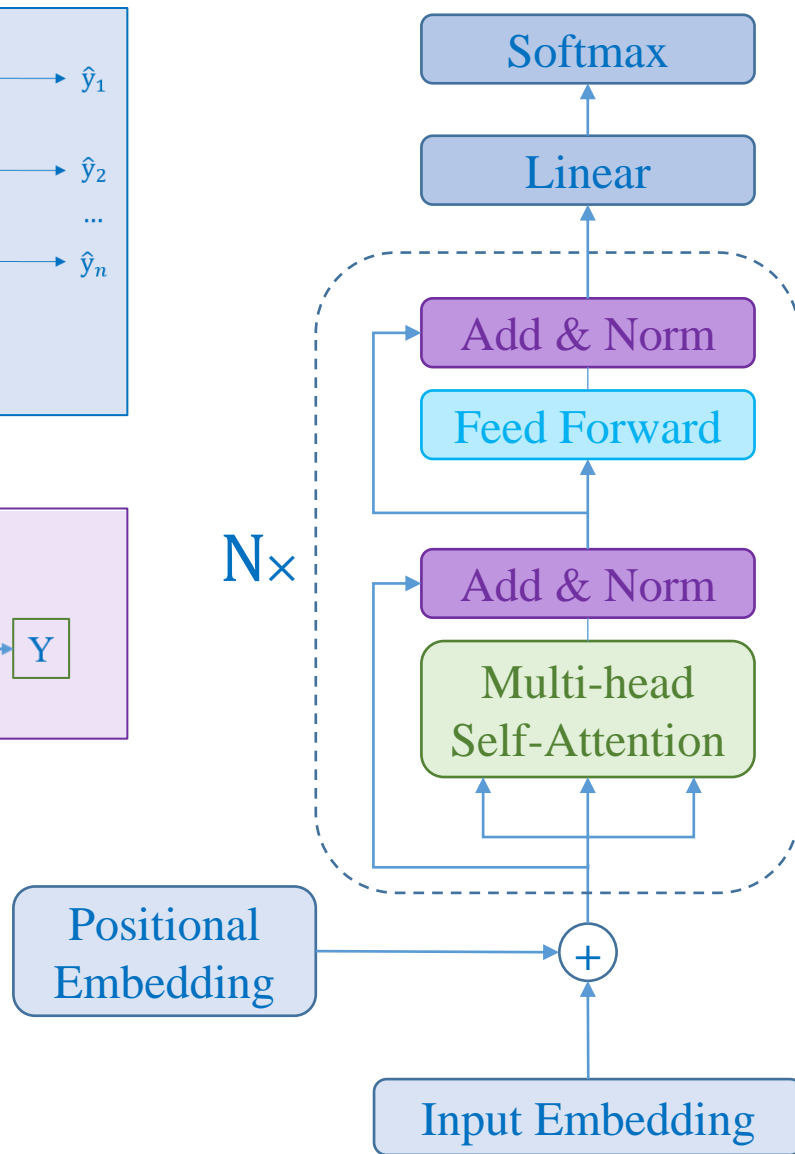
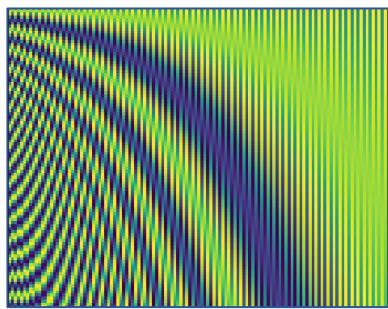
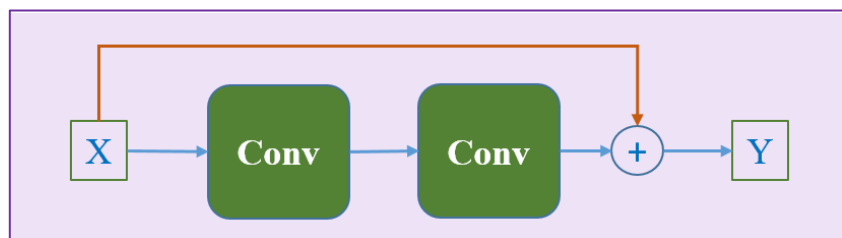
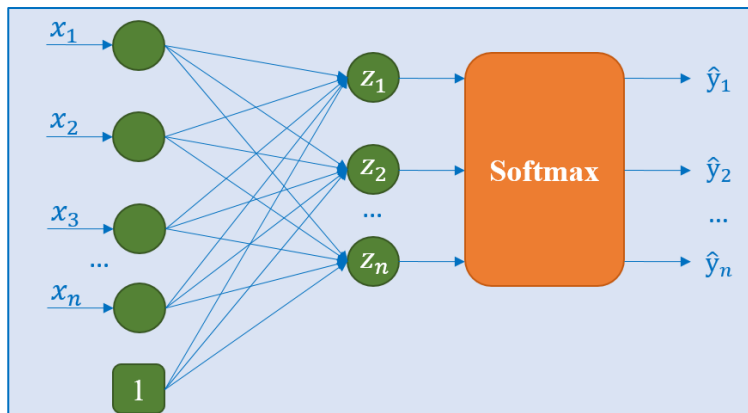
❖ Practice



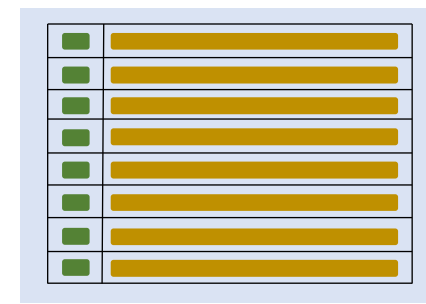
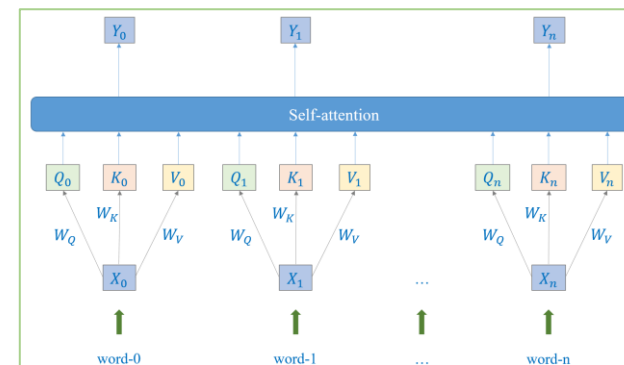
Outline

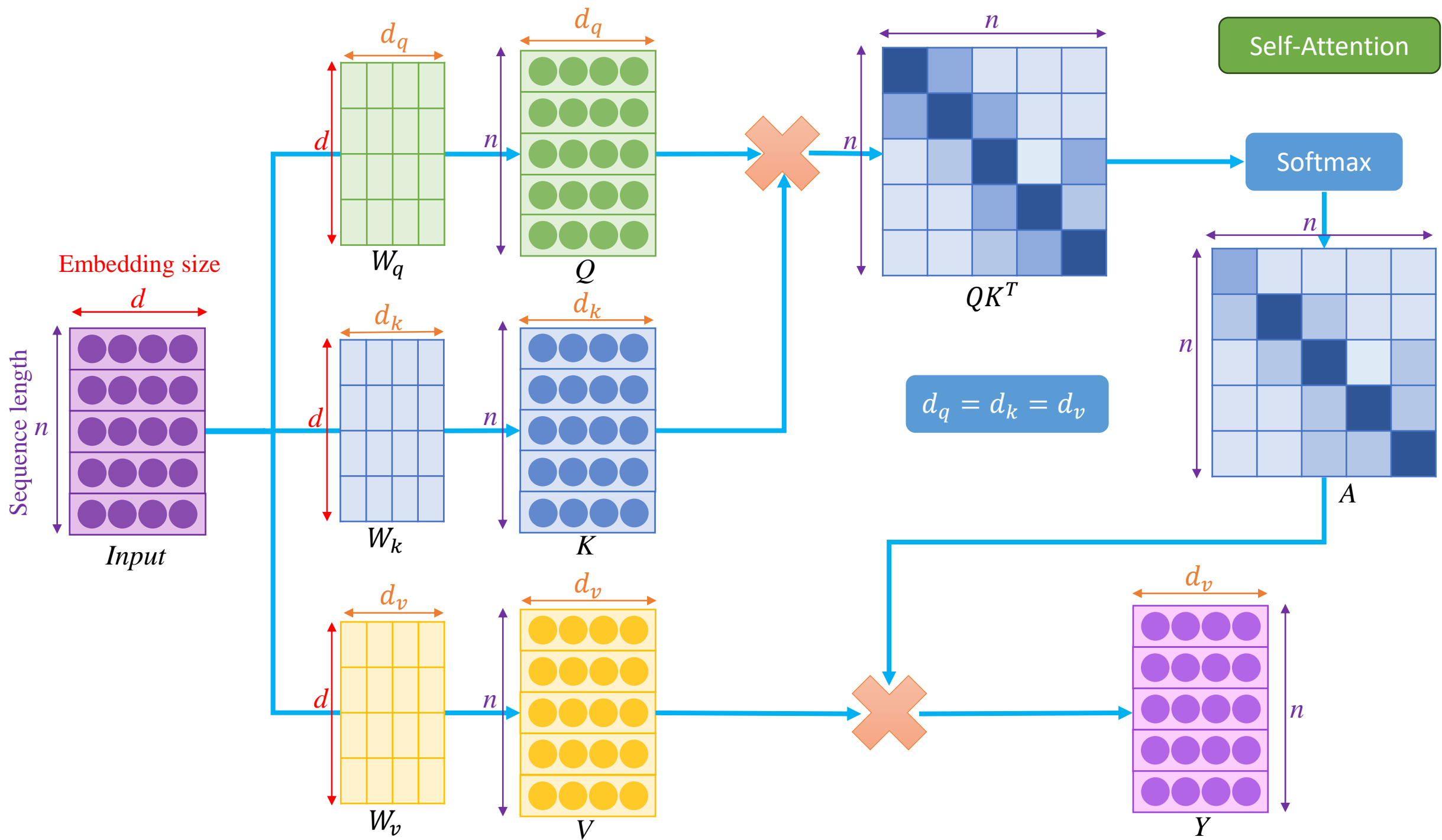
- **Motivation to Text Generation**
- **Simple Models**
- **Using RNNs**
- **Examples**
- **Using Masked Encoder (~Decoder)**

Transformer Encoder

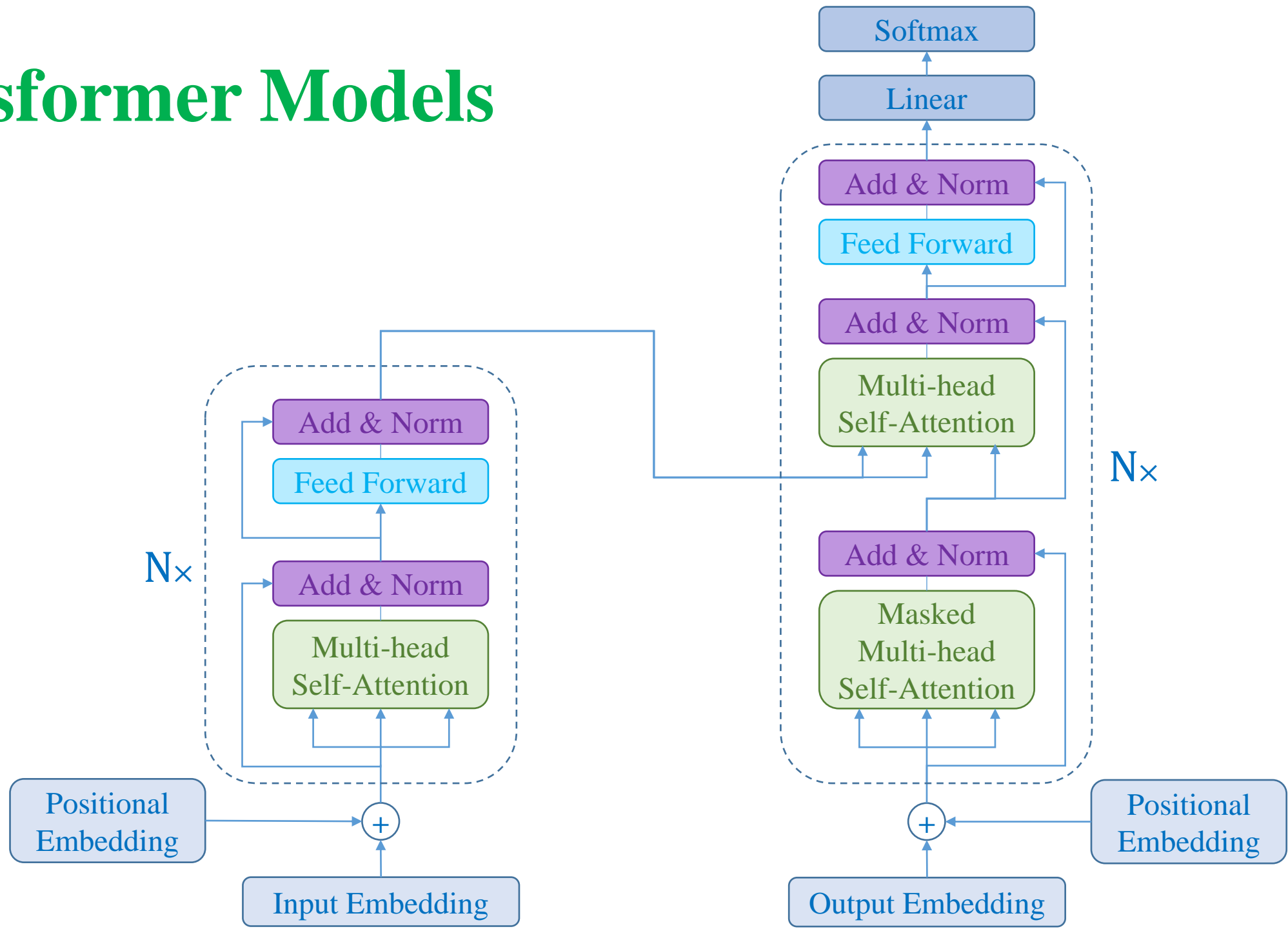


<https://arxiv.org/pdf/1803.08494.pdf>



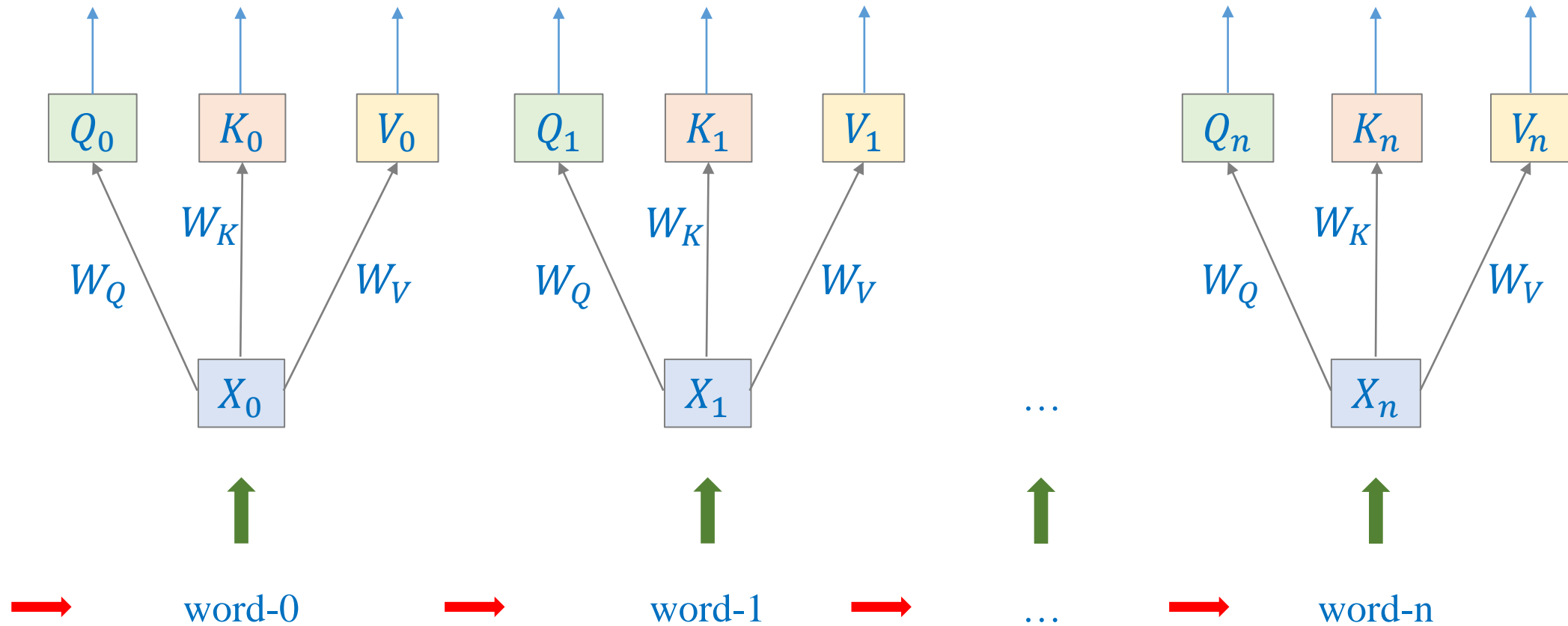


Transformer Models

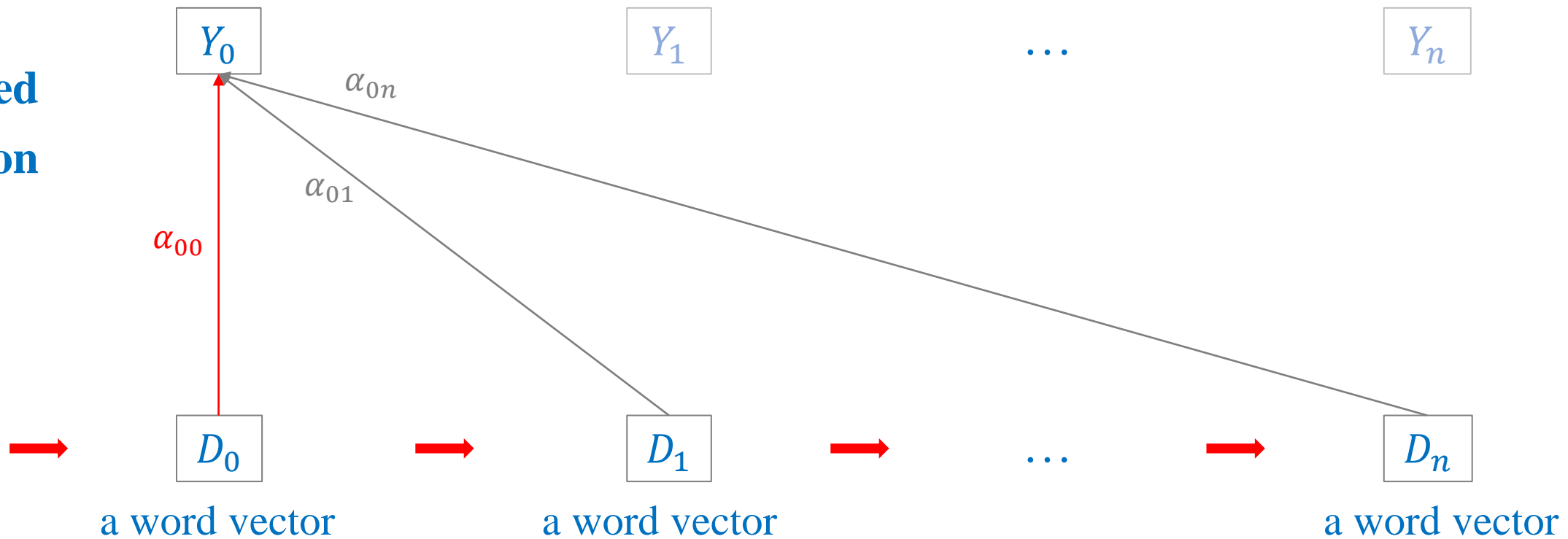


Transformer

❖ Masked self-attention



Masked self-attention



$$Y_0 = \alpha_{00}D_0 + \alpha_{01}D_1 + \dots + \alpha_{0n}D_n$$

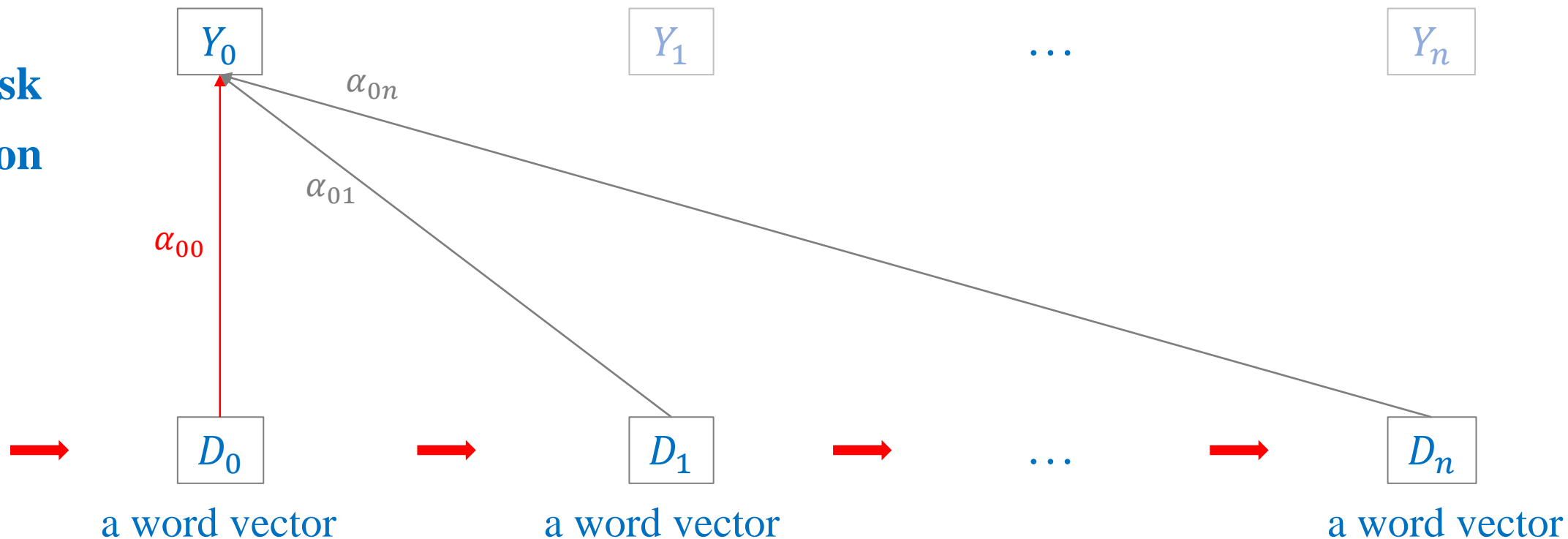


$$Y_0 = \alpha_{00}D_0 + 0 \times D_1 + \dots + 0 \times D_n$$

$$\alpha_0 = \text{softmax}\left(\frac{D_0 D^T}{\sqrt{d}}\right) = \begin{bmatrix} \alpha_{00} \\ \alpha_{01} \\ \dots \\ \alpha_{0n} \end{bmatrix} \quad \text{How to obtain kind of } \rightarrow \begin{bmatrix} \alpha_{00} \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

$$\alpha_0 = \text{softmax}\left(\frac{D_0 D^T}{\sqrt{d}}\right) * \begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix} = \begin{bmatrix} \alpha_{00} \\ 0 \\ \dots \\ 0 \end{bmatrix} \quad ?$$

Mask self-attention



$$Y_0 = \alpha_{00}D_0 + \alpha_{01}D_1 + \dots + \alpha_{0n}D_n$$

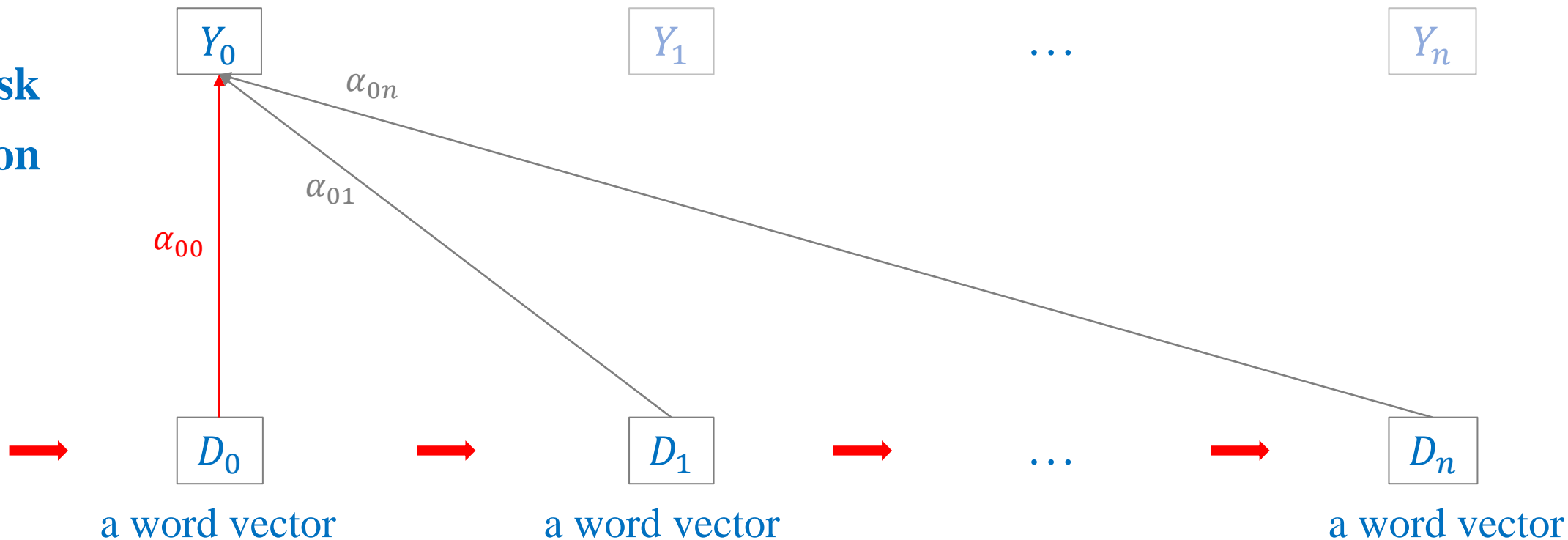


$$Y_0 = \alpha_{00}D_0 + 0 \times D_1 + \dots + 0 \times D_n$$

$$\alpha_0 = \text{softmax}\left(\frac{D_0 D^T}{\sqrt{d}}\right) = \begin{bmatrix} \alpha_{00} \\ \alpha_{01} \\ \dots \\ \alpha_{0n} \end{bmatrix} \quad \text{How to obtain kind of } \rightarrow \begin{bmatrix} \alpha_{00} \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

$$\alpha_0 = \text{softmax}\left(\frac{D_0 D^T}{\sqrt{d}} * \begin{bmatrix} 1 \\ 0 \\ \dots \\ 0 \end{bmatrix}\right) = \begin{bmatrix} \alpha_{00} \\ 0 \\ \dots \\ 0 \end{bmatrix} \quad ?$$

Mask self-attention



$$Y_0 = \alpha_{00}D_0 + \alpha_{01}D_1 + \dots + \alpha_{0n}D_n$$

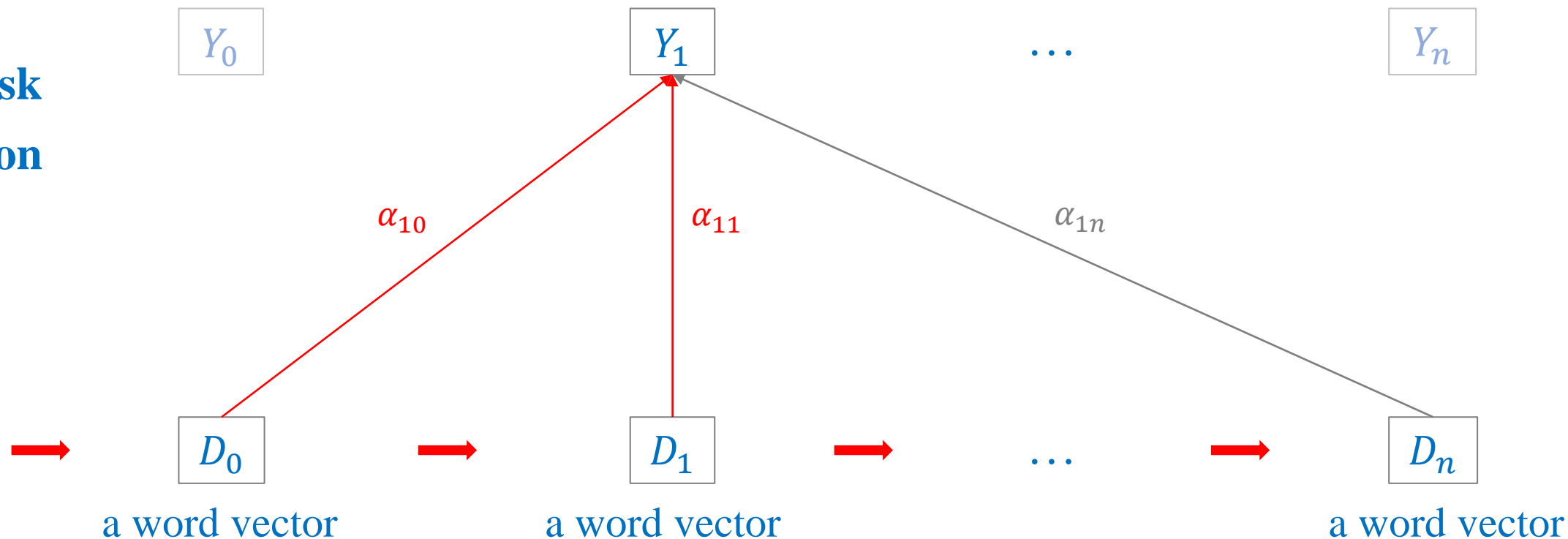


$$Y_0 = \alpha_{00}D_0 + 0 \times D_1 + \dots + 0 \times D_n$$

$$\alpha_0 = \text{softmax}\left(\frac{D_0 D^T}{\sqrt{m}}\right) = \begin{bmatrix} \alpha_{00} \\ \alpha_{01} \\ \dots \\ \alpha_{0n} \end{bmatrix} \quad \text{How to obtain kind of } \rightarrow \begin{bmatrix} \alpha_{00} \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

$$\alpha_0 = \text{softmax}\left(\frac{D_0 D^T}{\sqrt{d}} + \begin{bmatrix} 0 \\ -\infty \\ \dots \\ -\infty \end{bmatrix}\right) = \begin{bmatrix} \alpha_{00} \\ 0 \\ \dots \\ 0 \end{bmatrix} \quad ?$$

**Mask
self-attention**



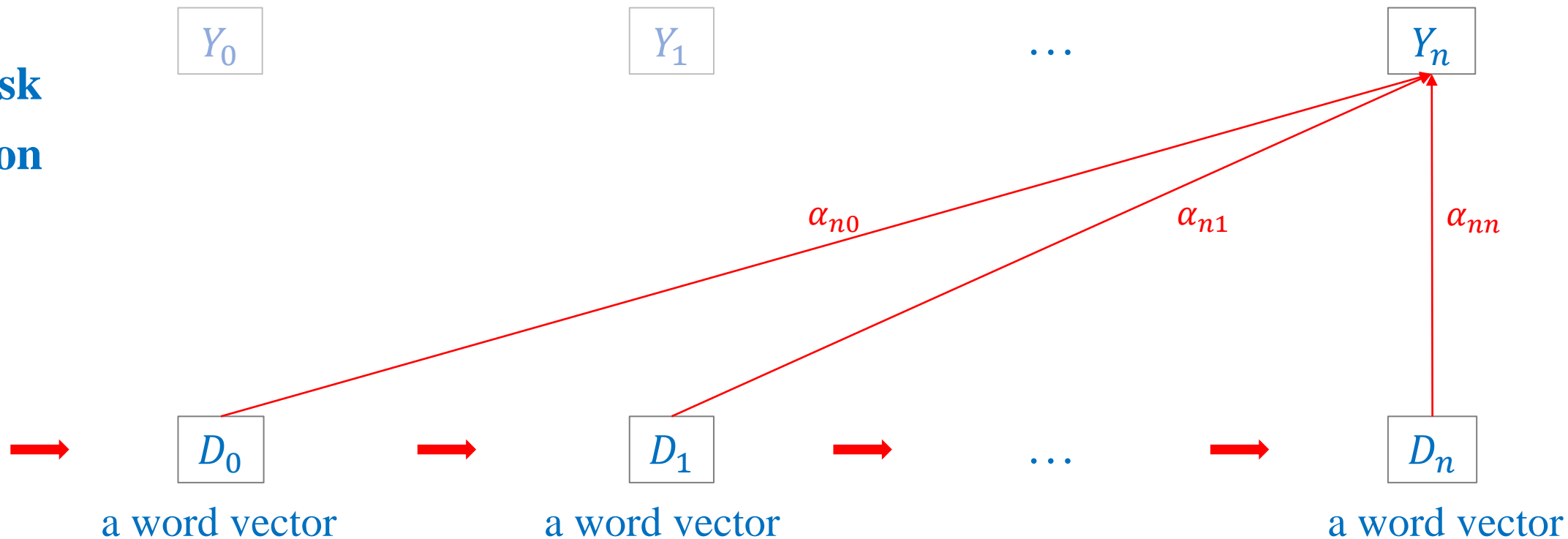
$$Y_1 = \alpha_{10}D_0 + \alpha_{11}D_1 + \cdots + \alpha_{1n}D_n$$



$$Y_1 = \alpha_{10}D_0 + \alpha_{11}D_1 + \cdots + 0 \times D_n$$

$$\alpha_1 = \text{softmax} \left(\frac{D_1 D^T}{\sqrt{d}} + \begin{bmatrix} 0 \\ 0 \\ -\infty \\ \vdots \\ -\infty \end{bmatrix} \right) = \begin{bmatrix} \alpha_{10} \\ \alpha_{11} \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Mask self-attention

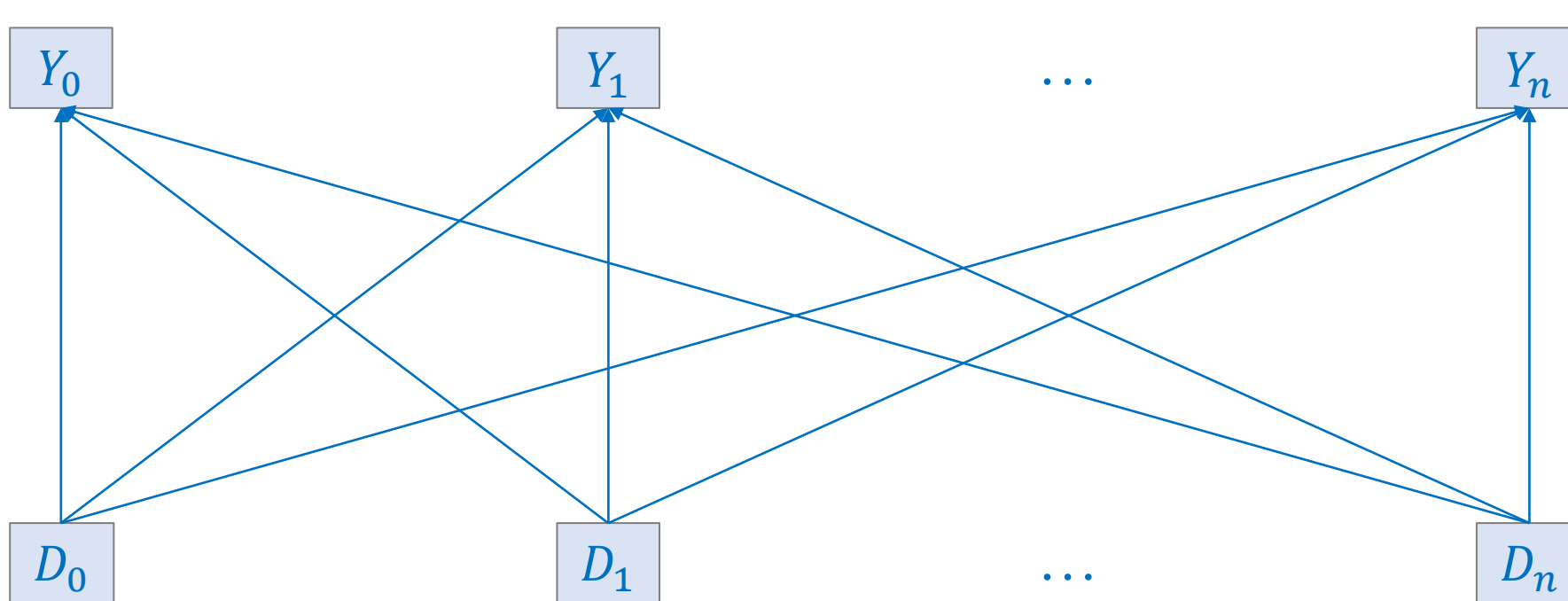


$$Y_n = \alpha_{n0}D_0 + \alpha_{n1}D_1 + \dots + \alpha_{nn}D_n$$

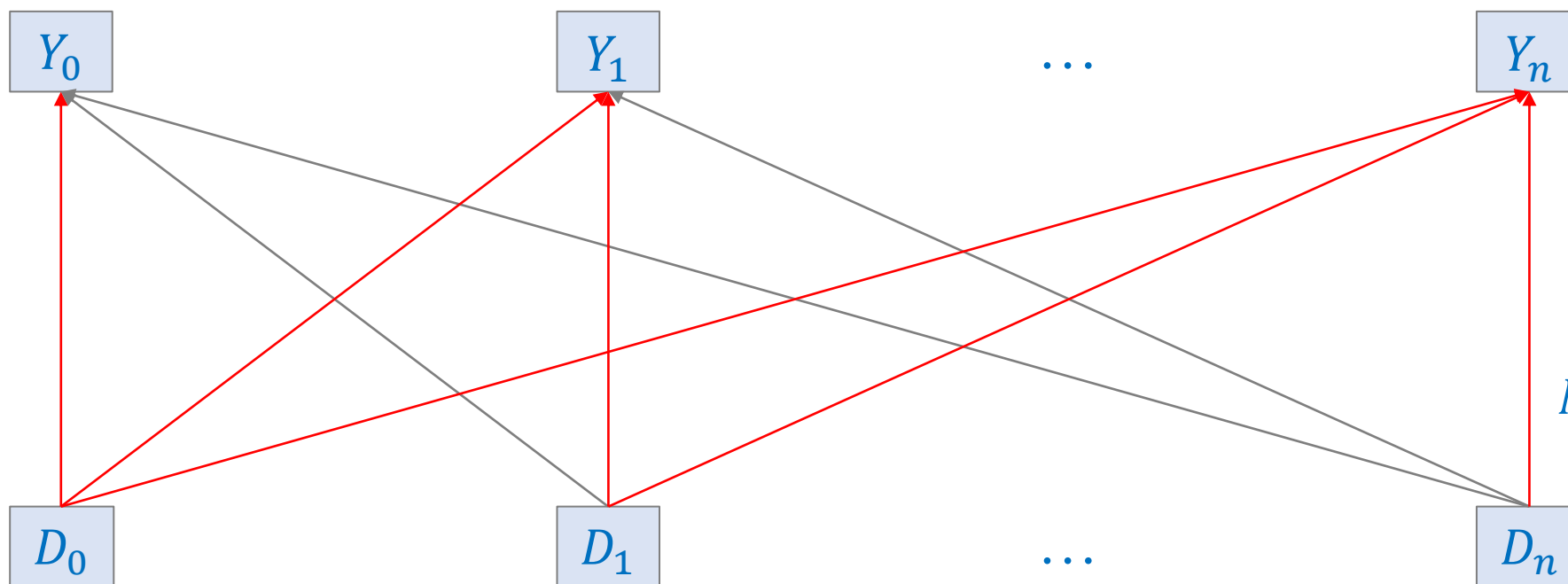


?

$$\alpha_n = \text{softmax} \left(\frac{D_n D^T}{\sqrt{d}} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ \dots \\ 0 \end{bmatrix} \right) = \begin{bmatrix} \alpha_{n0} \\ \alpha_{n1} \\ \alpha_{n2} \\ \dots \\ \alpha_{nn} \end{bmatrix}$$



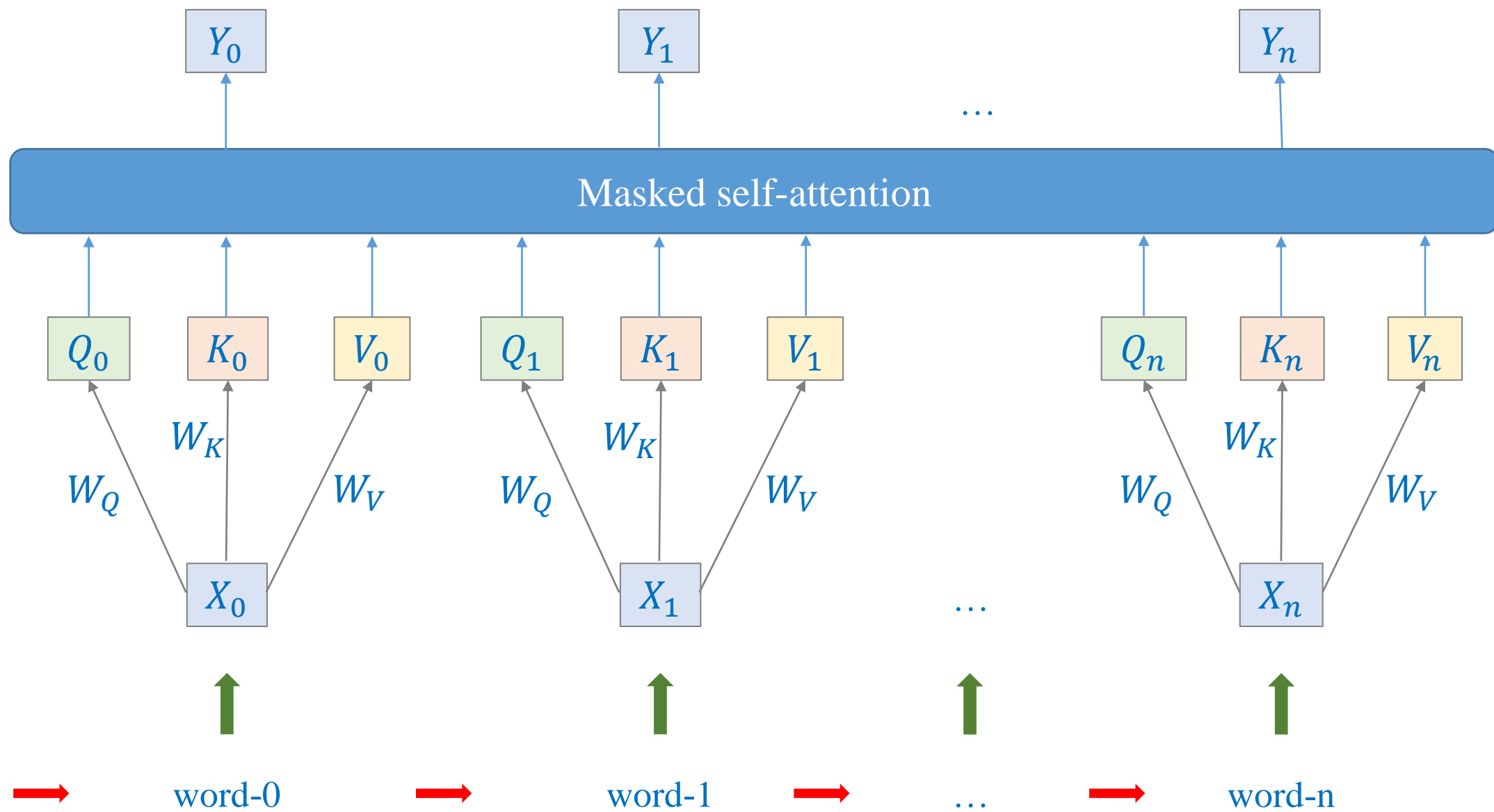
$$\alpha = \text{sigmoid}\left(\frac{DD^T}{\sqrt{d}}\right)$$



$$\alpha = \text{sigmoid}\left(\frac{DD^T}{\sqrt{d}} + M\right)$$

$$M = \begin{bmatrix} [0 & -\infty & -\infty & \dots & -\infty] \\ [0 & 0 & -\infty & \dots & -\infty] \\ [0 & 0 & \dots & 0 & \dots & 0] \end{bmatrix}$$

$$Y = \text{sigmoid} \left(\frac{QK^T}{\sqrt{d}} + M \right) V \quad M = \begin{bmatrix} [0 & -\infty & -\infty \dots -\infty] \\ [0 & 0 & -\infty \dots -\infty] \\ [0 & 0 & \dots 0 \dots 0] \end{bmatrix}$$



```

x = torch.tensor([[[[-0.1, 0.1, 0.3], [ 0.4, -1.1, -0.3]]]])
layer = nn.MultiheadAttention(embed_dim=3, num_heads=1, batch_first=True)

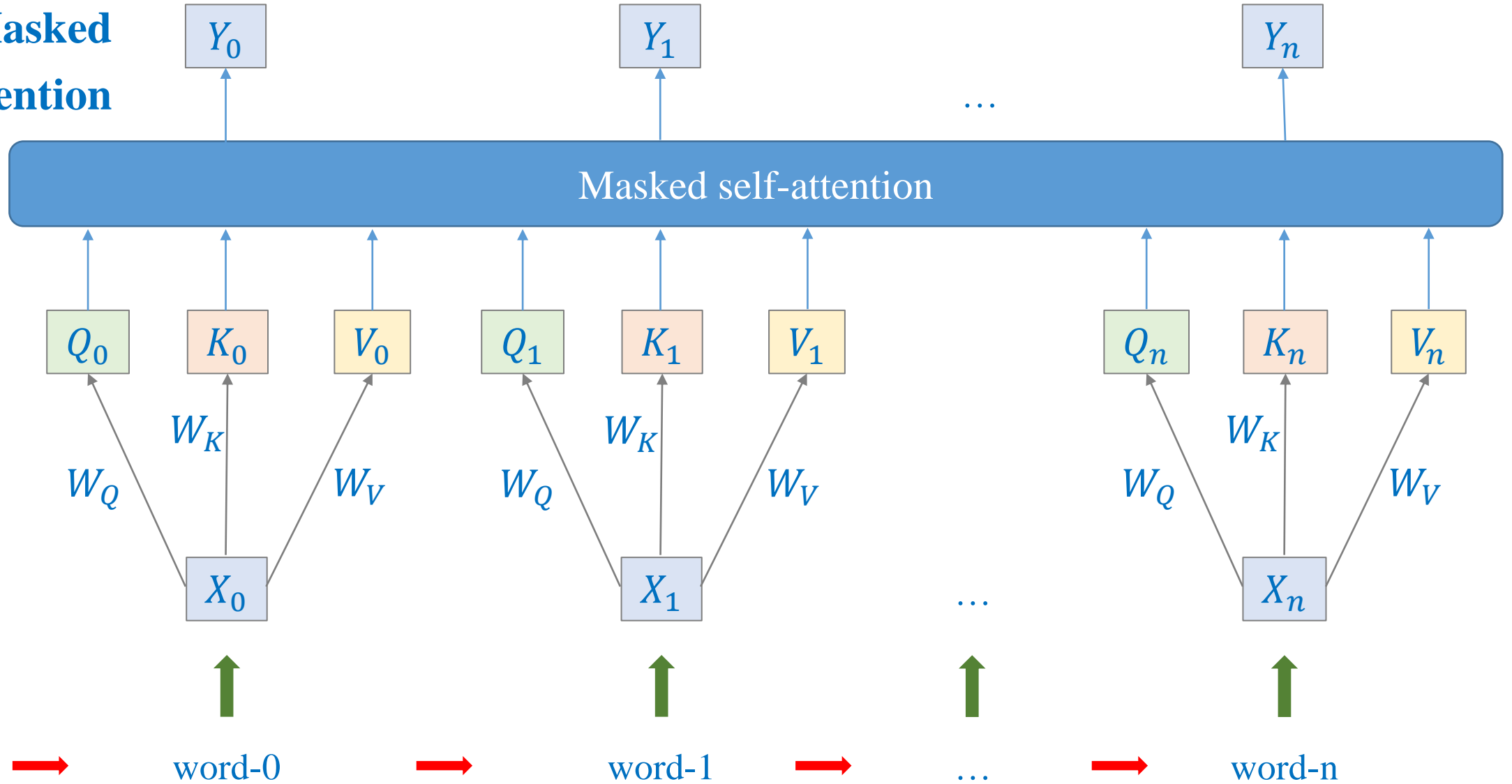
mask = torch.triu(torch.ones(1, 2, 2), diagonal=1).bool()
output_tensor, attn_output_weights = layer(query=x, key=x, value=x, attn_mask=mask)

```

$$A = \text{sigmoid}\left(\frac{QK^T}{\sqrt{d}} + M\right) V$$

$$Y = AW_O$$

Masked Self-Attention



Masked Multi-head Attention

head = 1

$$W_Q = \begin{bmatrix} -0.35 & 0.51 & 0.50 \\ 0.36 & -0.47 & -0.29 \\ -0.51 & -0.14 & -0.56 \end{bmatrix}$$

$$W_K = \begin{bmatrix} -0.49 & -0.68 & 0.18 \\ -0.44 & -0.46 & 0.18 \\ 0.07 & -0.10 & 0.44 \end{bmatrix}$$

$$W_V = \begin{bmatrix} -0.41 & 0.39 & -0.65 \\ -0.40 & -0.07 & -0.34 \\ -0.55 & -0.13 & -0.29 \end{bmatrix}$$

$$W_O = \begin{bmatrix} -0.36 & -0.08 & 0.32 \\ 0.27 & 0.05 & 0.15 \\ -0.05 & -0.28 & 0.05 \end{bmatrix}$$

$$X = \begin{bmatrix} -0.1 & 0.1 & 0.3 \\ 0.4 & -1.1 & -0.3 \end{bmatrix}$$

$$Q = XW_Q = \begin{bmatrix} -0.1 & 0.1 & 0.3 \\ 0.4 & -1.1 & -0.3 \end{bmatrix} \begin{bmatrix} -0.35 & 0.51 & 0.50 \\ 0.36 & -0.47 & -0.29 \\ -0.51 & -0.14 & -0.56 \end{bmatrix}$$

$$= \begin{bmatrix} -0.08 & -0.14 & -0.24 \\ -0.39 & 0.77 & 0.69 \end{bmatrix}$$

$$K = XW_K = \begin{bmatrix} -0.1 & 0.1 & 0.3 \\ 0.4 & -1.1 & -0.3 \end{bmatrix} \begin{bmatrix} -0.49 & -0.68 & 0.18 \\ -0.44 & -0.46 & 0.18 \\ 0.07 & -0.10 & 0.44 \end{bmatrix}$$

$$= \begin{bmatrix} 0.02 & -0.01 & 0.13 \\ 0.27 & 0.27 & -0.26 \end{bmatrix}$$

$$V = XW_V = \begin{bmatrix} -0.1 & 0.1 & 0.3 \\ 0.4 & -1.1 & -0.3 \end{bmatrix} \begin{bmatrix} -0.41 & 0.39 & -0.65 \\ -0.40 & -0.07 & -0.34 \\ -0.55 & -0.13 & -0.29 \end{bmatrix}$$

$$= \begin{bmatrix} -0.16 & -0.08 & -0.05 \\ -0.02 & -0.02 & 0.05 \end{bmatrix}$$

$$M = \begin{bmatrix} 0 & -\infty \\ 0 & 0 \end{bmatrix}$$

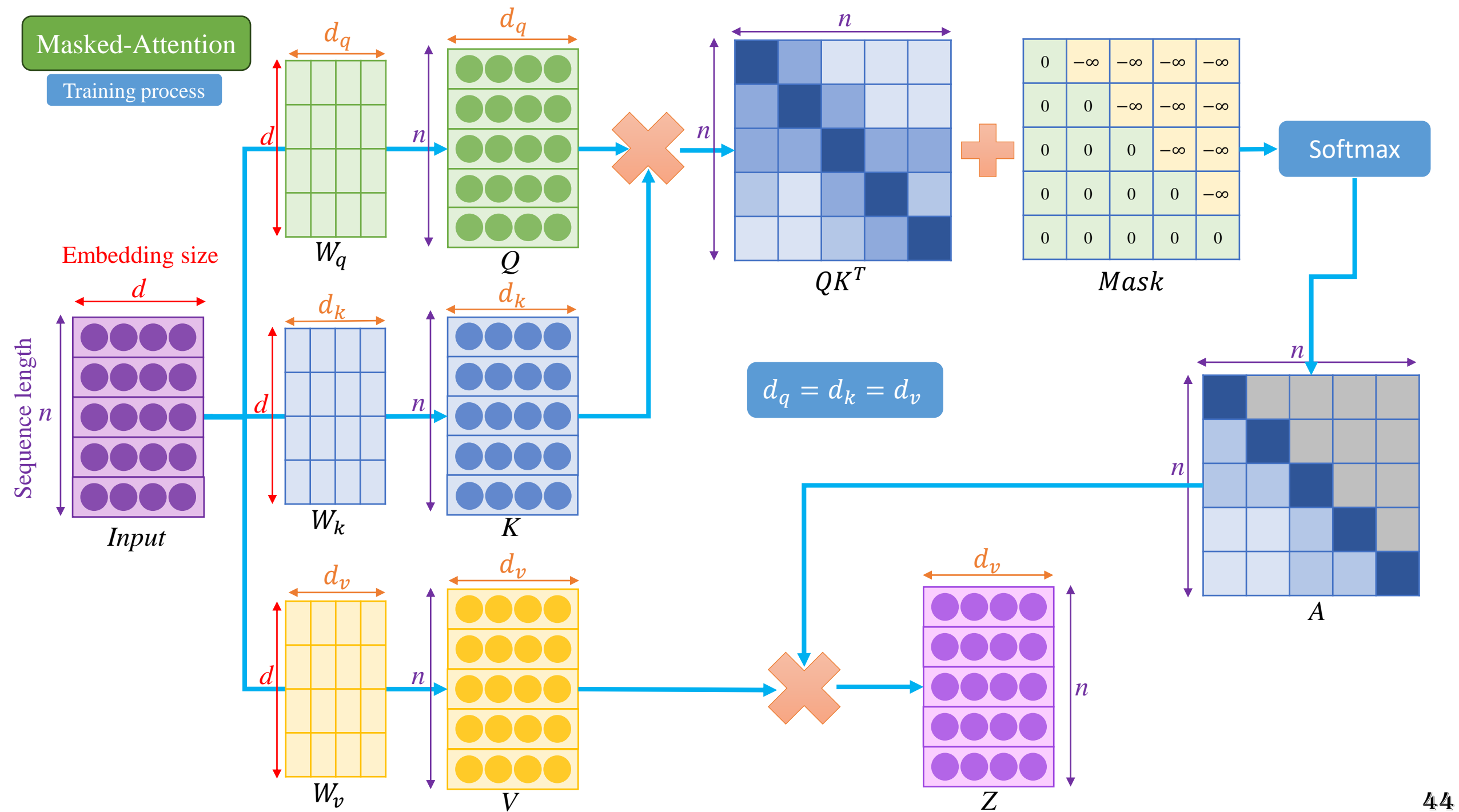
Masked Multi-head Attention

❖ Example

approximately

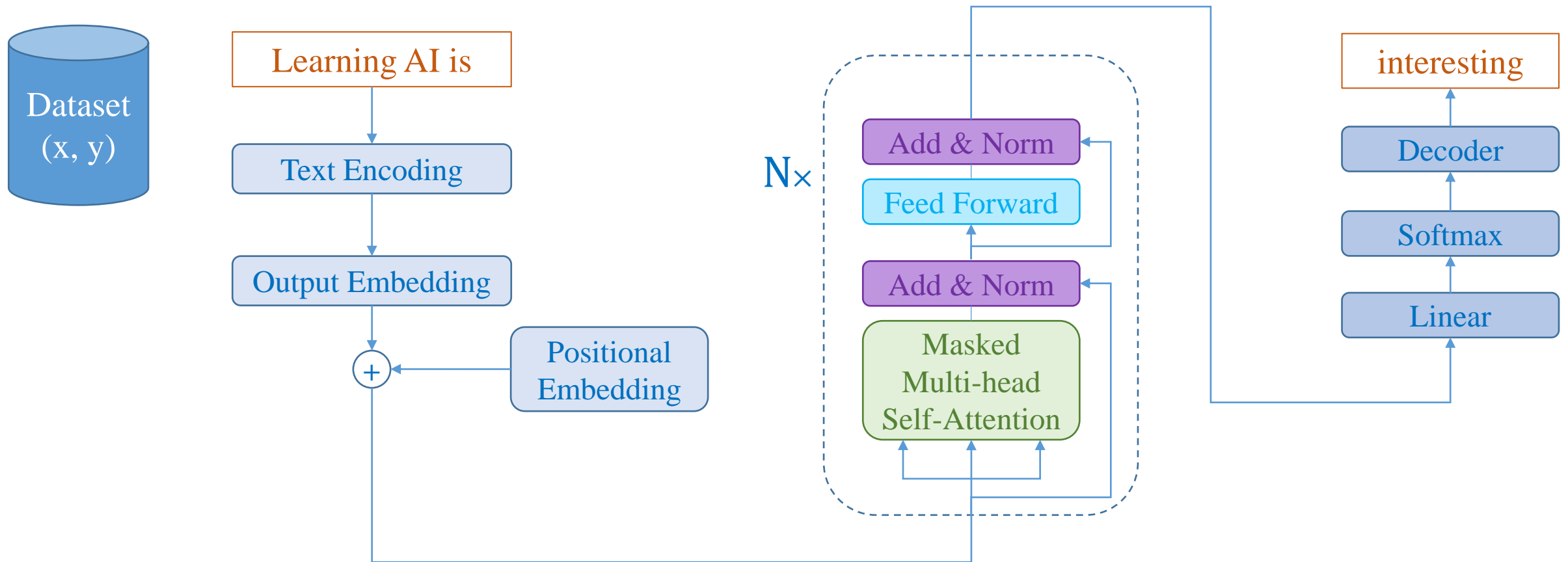
$$\begin{aligned}
 A &= \text{sigmoid} \left(\frac{QK^T}{\sqrt{d}} + M \right) V \\
 &= \text{sigmoid} \left(\begin{bmatrix} -0.08 & -0.14 & -0.24 \\ -0.39 & 0.77 & 0.69 \end{bmatrix} \begin{bmatrix} 0.02 & 0.27 \\ -0.01 & 0.27 \\ 0.13 & -0.26 \end{bmatrix} \frac{1}{\sqrt{d}} + \begin{bmatrix} 0 & -\infty \\ 0 & 0 \end{bmatrix} \right) \begin{bmatrix} -0.16 & -0.08 & -0.05 \\ -0.02 & -0.02 & 0.05 \end{bmatrix} \\
 &= \text{sigmoid} \left(\begin{bmatrix} -0.019 & 0.002 \\ 0.043 & -0.046 \end{bmatrix} + \begin{bmatrix} 0 & -\infty \\ 0 & 0 \end{bmatrix} \right) \begin{bmatrix} -0.16 & -0.08 & -0.05 \\ -0.02 & -0.02 & 0.05 \end{bmatrix} \\
 &= \begin{bmatrix} 1.0 & 0.0 \\ 0.52 & 0.48 \end{bmatrix} \begin{bmatrix} -0.16 & -0.08 & -0.05 \\ -0.02 & -0.02 & 0.05 \end{bmatrix} = \begin{bmatrix} -0.16 & -0.08 & -0.05 \\ 0.12 & 0.08 & 0.06 \end{bmatrix}
 \end{aligned}$$

$$Y = AW_O = \begin{bmatrix} -0.16 & -0.08 & -0.05 \\ 0.12 & 0.08 & 0.06 \end{bmatrix} \begin{bmatrix} -0.36 & -0.08 & 0.32 \\ 0.27 & 0.05 & 0.15 \\ -0.05 & -0.28 & 0.05 \end{bmatrix} = \begin{bmatrix} 0.03 & 0.02 & -0.06 \\ -0.02 & -0.02 & 0.05 \end{bmatrix}$$



Text Generation

❖ Architecture

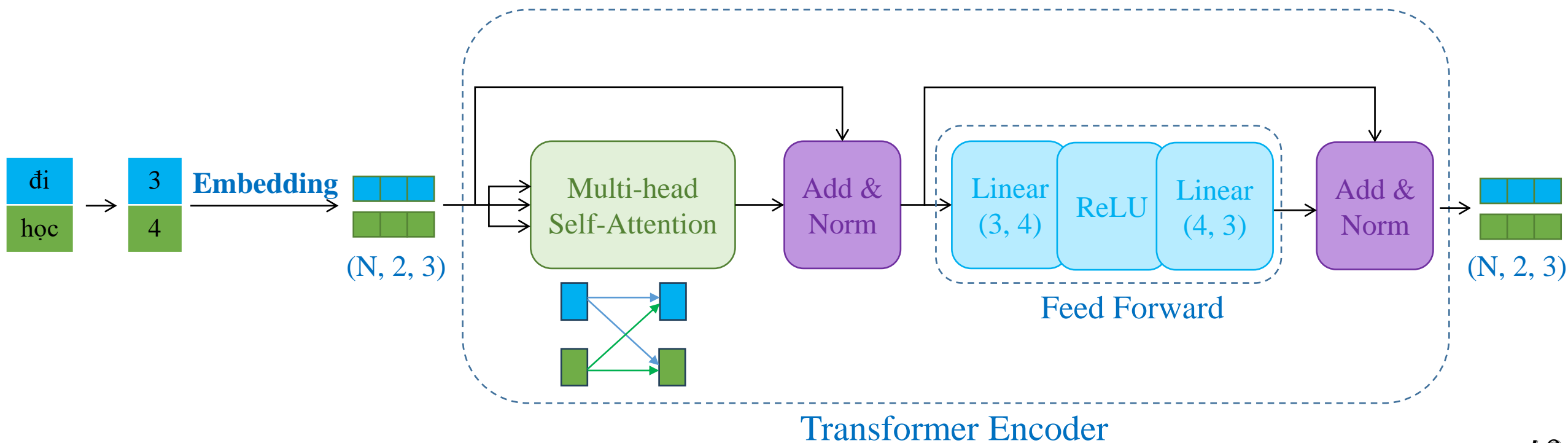


Encoder in PyTorch

Encoder in Pytorch

index	word
0	[UNK]
1	[pad]
2	ai
3	đi
4	học
...	...

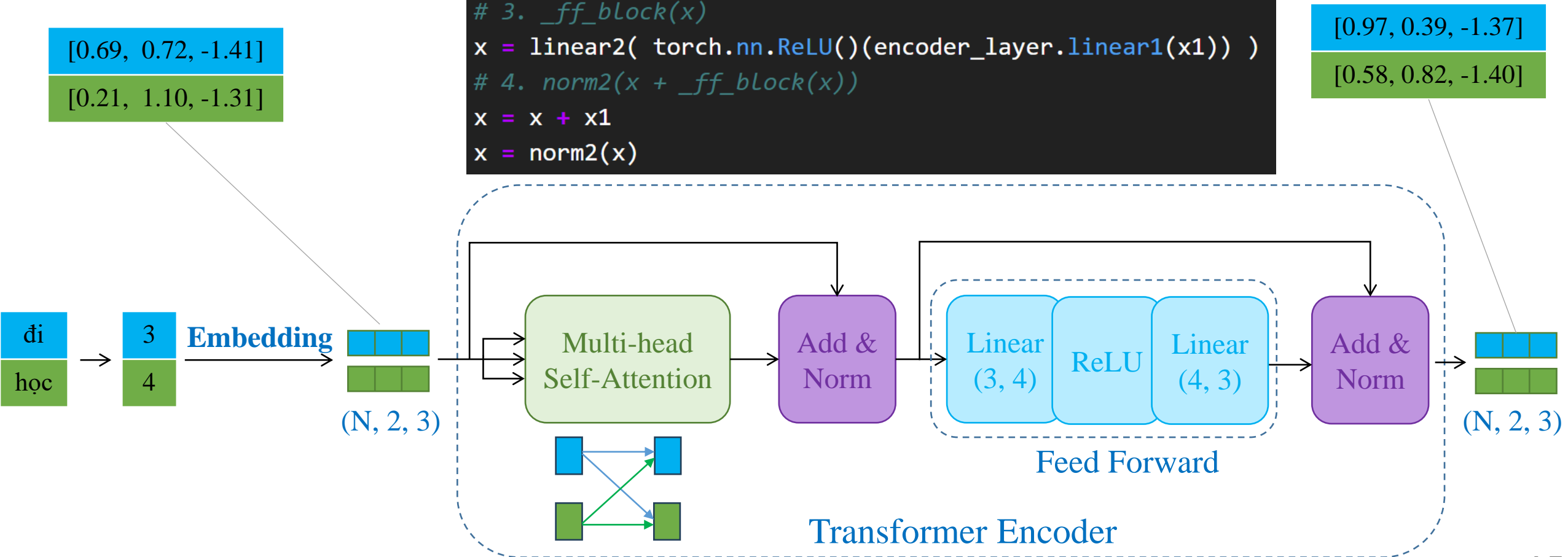
index	Embedding
0	$[-0.188, \dots, 0.7013]$
1	$[1.7840\dots, 1.3586]$
2	$[1.0281, \dots, 0.4211]$
3	$[-1.308, \dots, -0.3680]$
4	$[0.2293, \dots, 2.0501]$
...	...



Encoder in Pytorch

```
# input after thought embedding
embed_input = torch.Tensor([[[ 0.69,  0.72, -1.41],
                               [ 0.21,  1.10, -1.31]]])

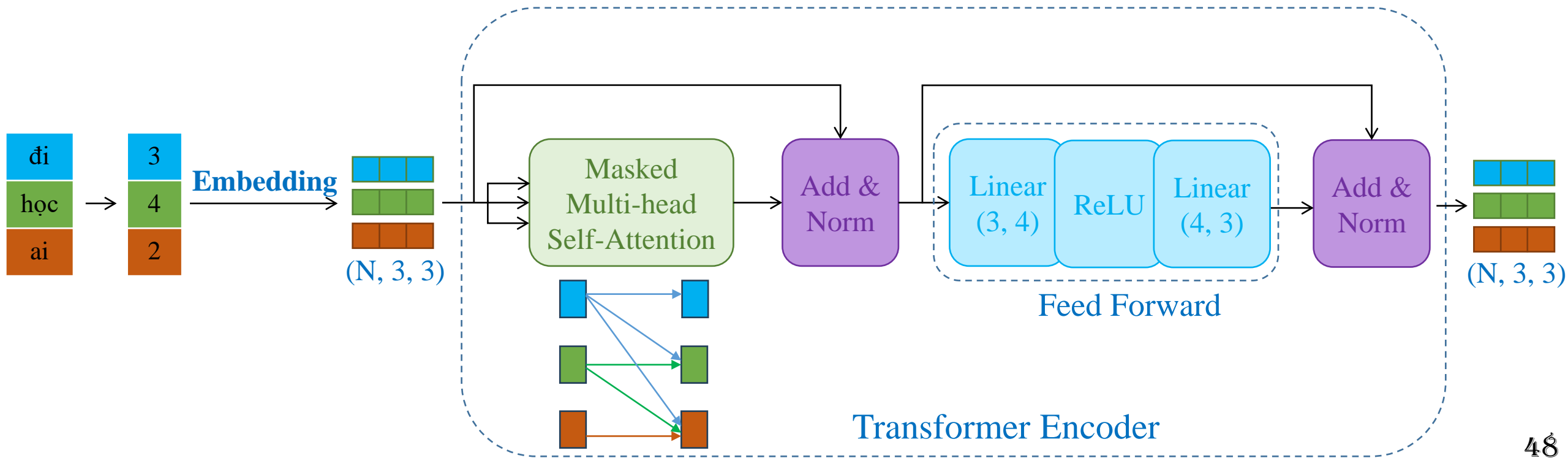
# 1. self_attn
x = self_attn(embed_input, embed_input, embed_input)[0]
# 2. norm1(x + self._sa_block(x))
x = src + x
x1 = norm1(x)
# 3. _ff_block(x)
x = linear2( torch.nn.ReLU()(encoder_layer.linear1(x1)) )
# 4. norm2(x + _ff_block(x))
x = x + x1
x = norm2(x)
```



Masked Encoder in PyTorch

Masked Encoder in Pytorch

index	word	index	Embedding
0	[UNK]	0	[-0.188, ..., 0.7013]
1	[pad]	1	[1.7840..., 1.3586]
2	ai	2	[1.0281, ..., 0.4211]
3	đi	3	[-1.308, ..., -0.3680]
4	học	4	[0.2293, ..., 2.0501]
...



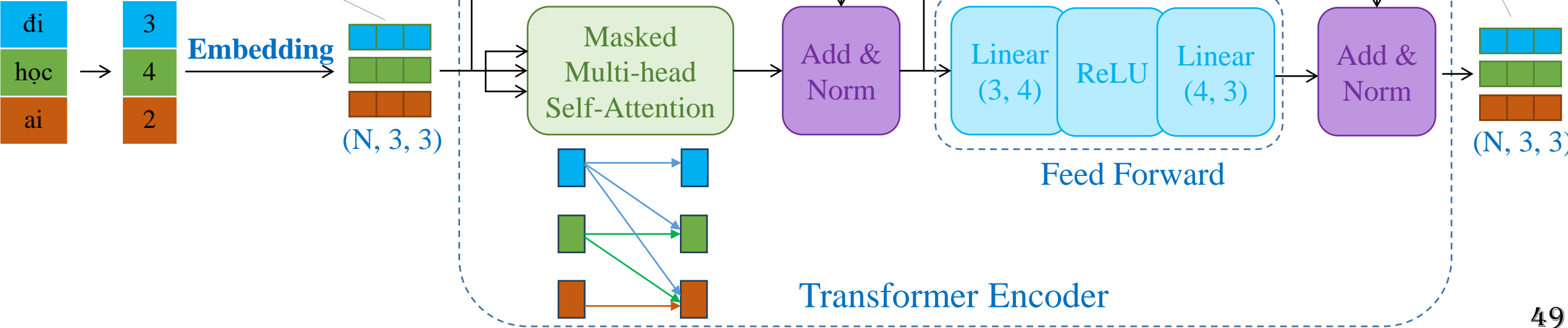
Masked Encoder in Pytorch

[0.69, 0.72, -1.41]
[0.21, 1.10, -1.31]
[-0.88, 0.60, -0.31]

```
# input after thought embedding
embed_input = torch.Tensor([[[ 0.69,  0.72, -1.41],
                               [ 0.21,  1.10, -1.31]]])

# 1. masked self_attn
mask = torch.triu(torch.ones(seq_len, seq_len), diagonal=1).bool()
x = self_attn(embed_input, embed_input, embed_input,attn_mask=mask)[0]
# 2. norm1(x + self._sa_block(x))
x = src + x
x1 = norm1(x)
# 3. _ff_block(x)
x = linear2( torch.nn.ReLU()(encoder_layer.linear1(x1)) )
# 4. norm2(x + _ff_block(x))
x = x + x1
x = norm2(x)
```

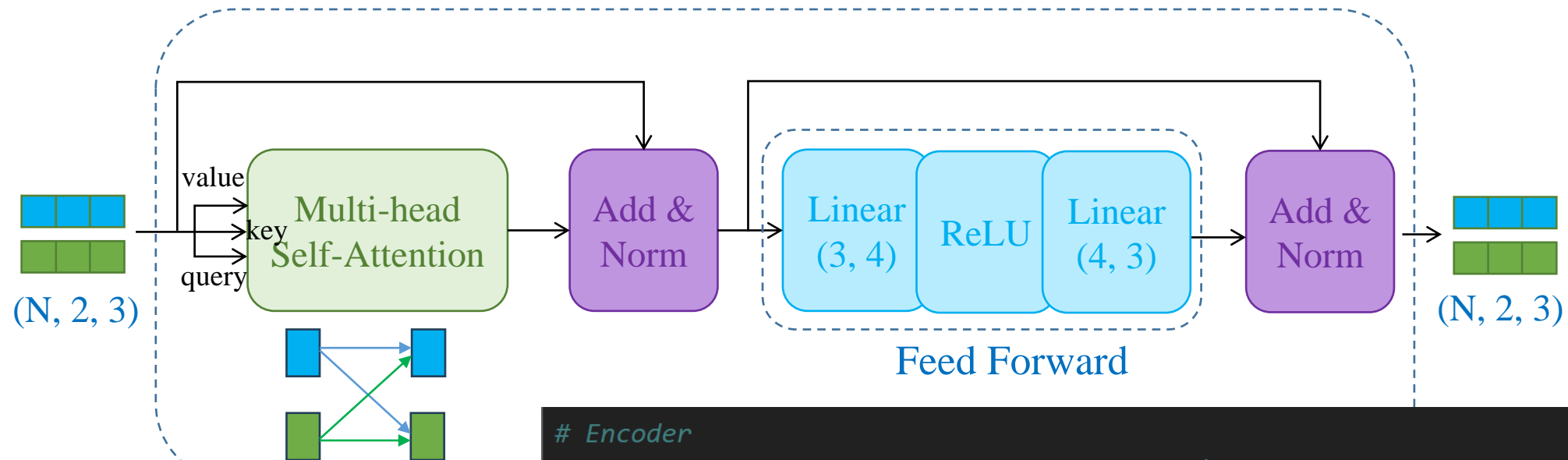
[0.97, 0.39, -1.37]
[0.58, 0.82, -1.40]
[-0.85, 1.40, -0.54]



Transformer in PyTorch

Transformer in PyTorch

❖ Transformer Encoder

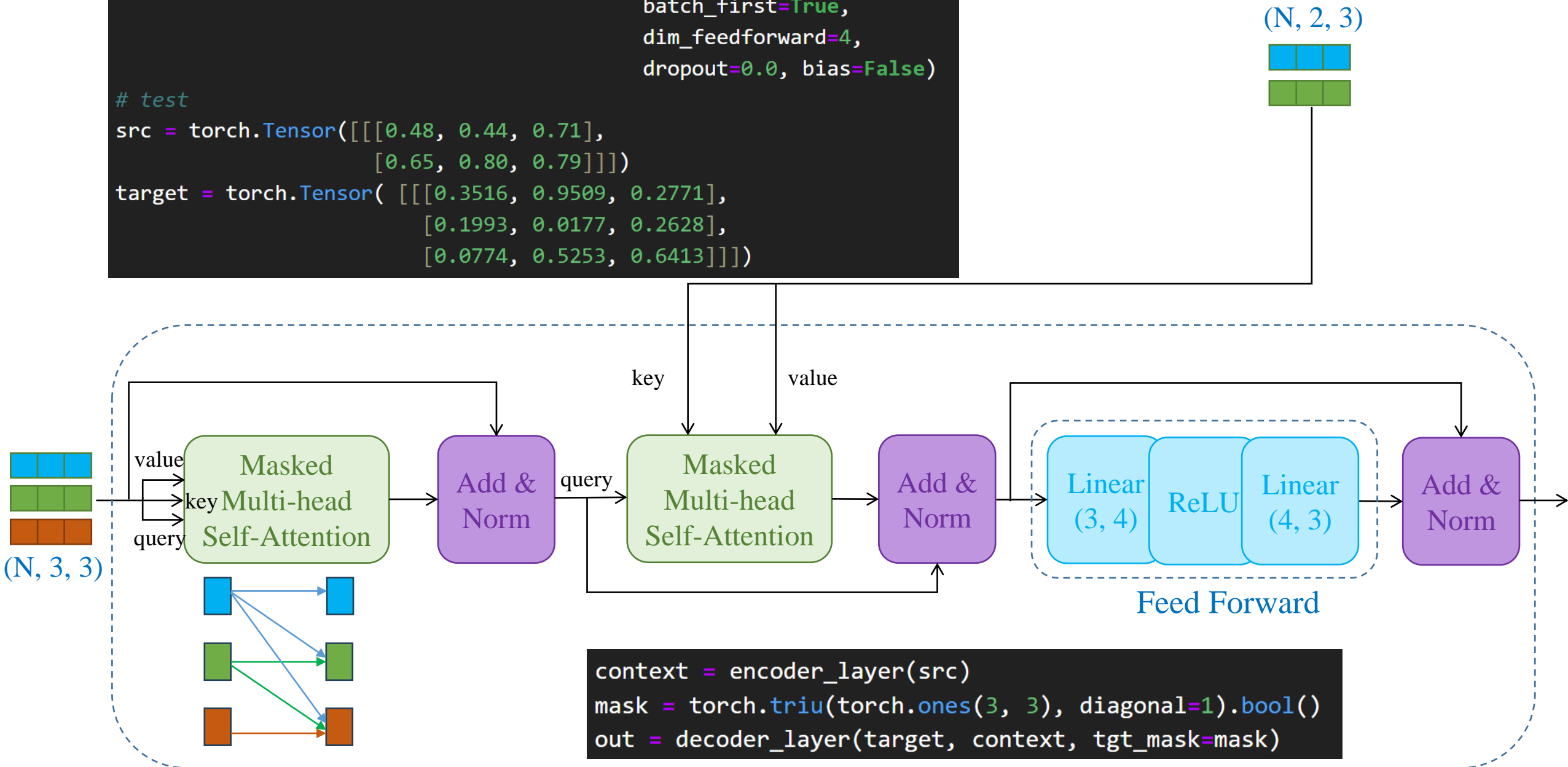


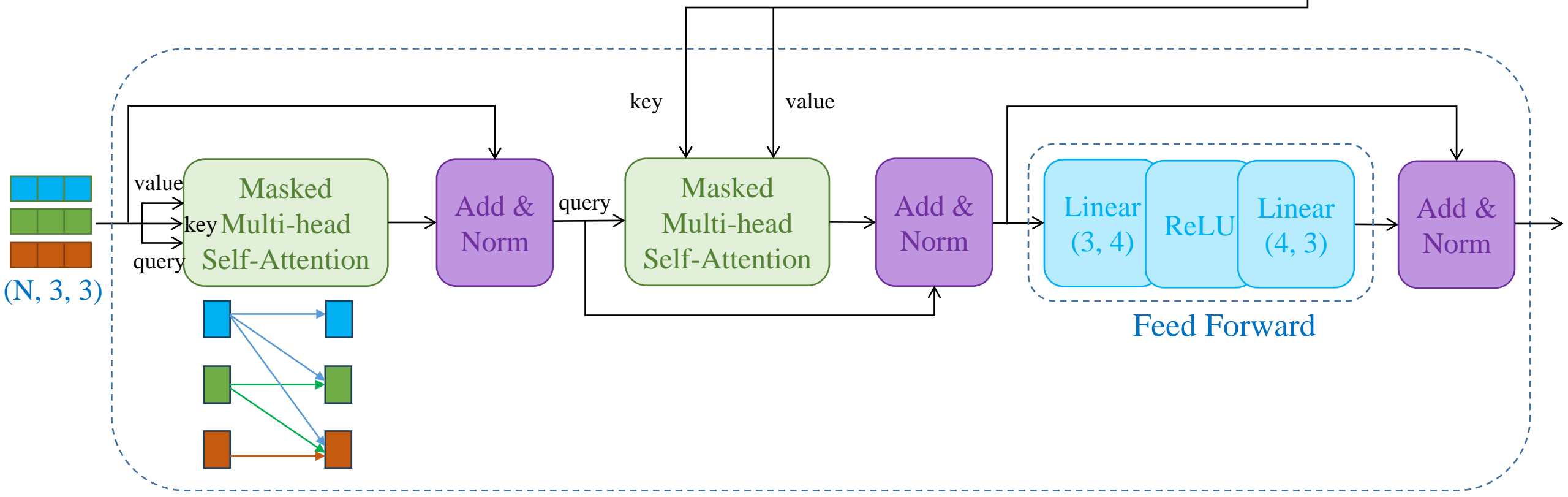
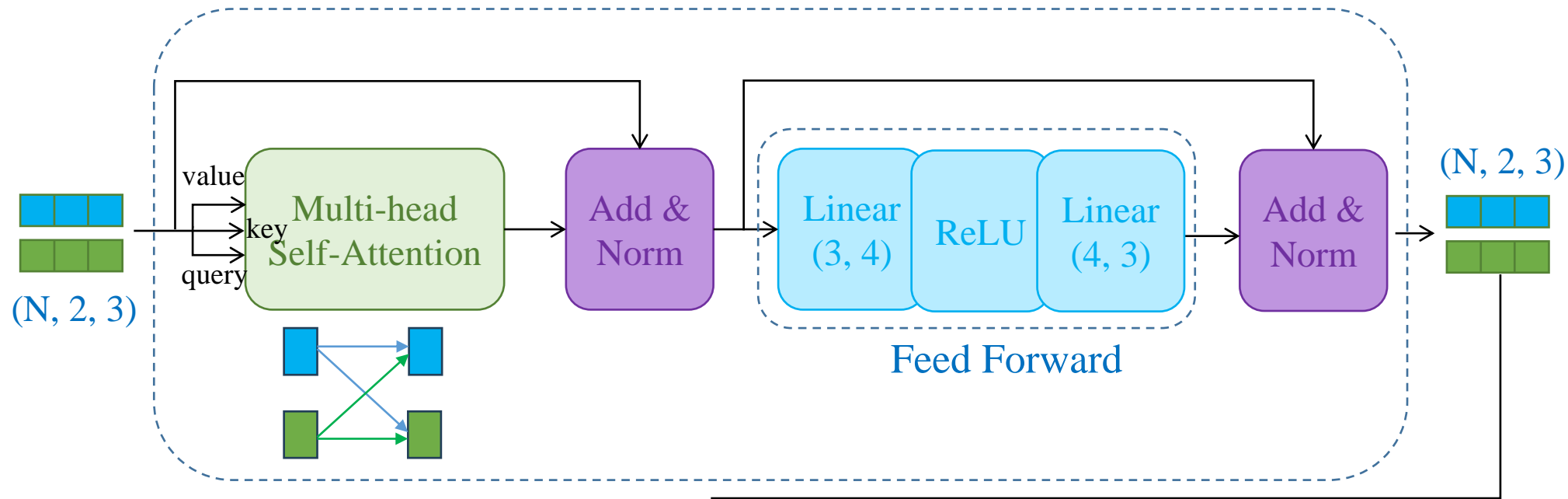
```
# Encoder
encoder_layer = nn.TransformerEncoderLayer(d_model=3, nhead=1,
                                           batch_first=True,
                                           dim_feedforward=4,
                                           dropout=0.0, bias=False)

# test
src = torch.Tensor([[[0.48, 0.44, 0.71],
                    [0.65, 0.80, 0.79]]])
context = encoder_layer(src)
```

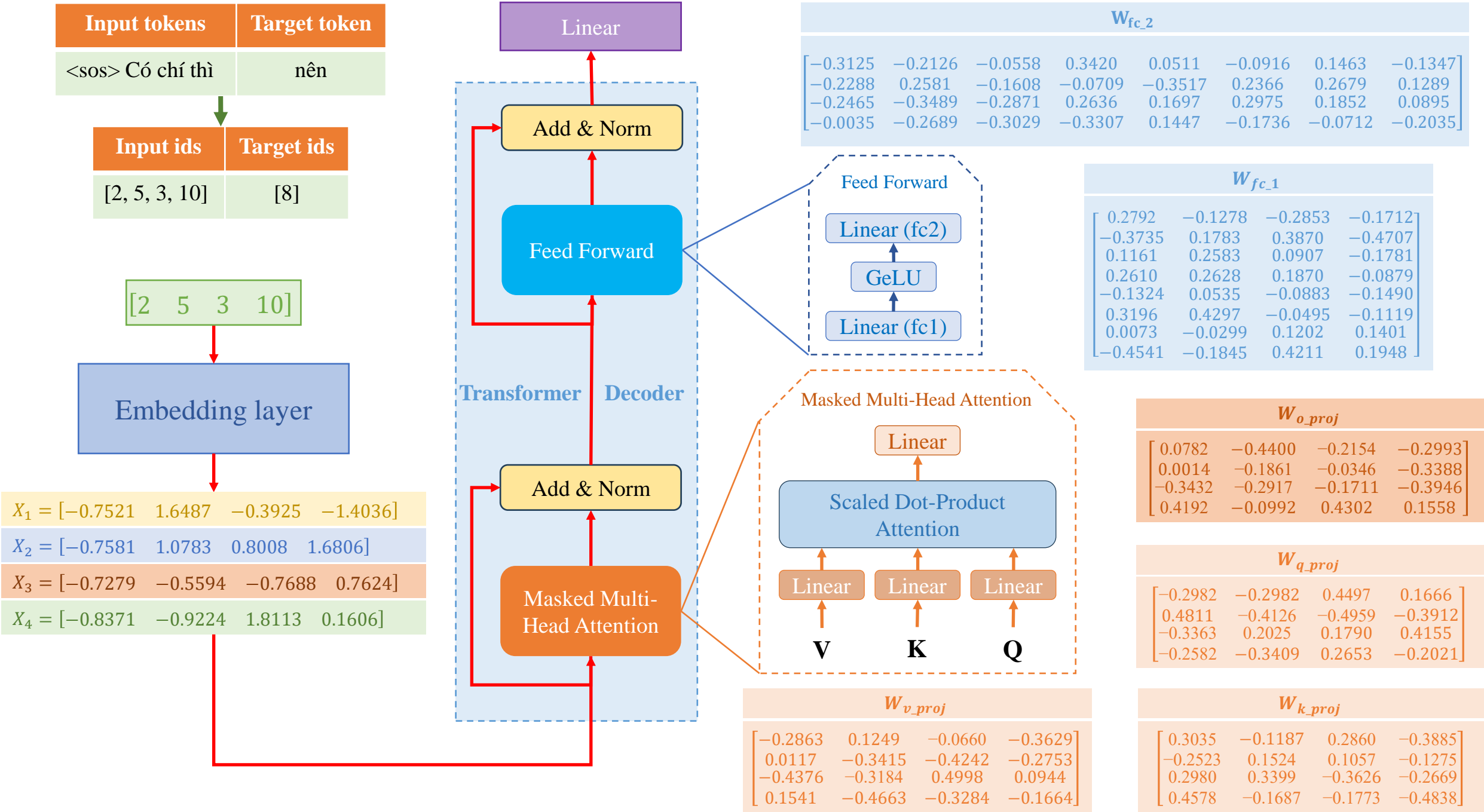
```
# Layers
encoder_layer = nn.TransformerEncoderLayer(...)
decoder_layer = nn.TransformerDecoderLayer(d_model=3, nhead=1,
                                           batch_first=True,
                                           dim_feedforward=4,
                                           dropout=0.0, bias=False)

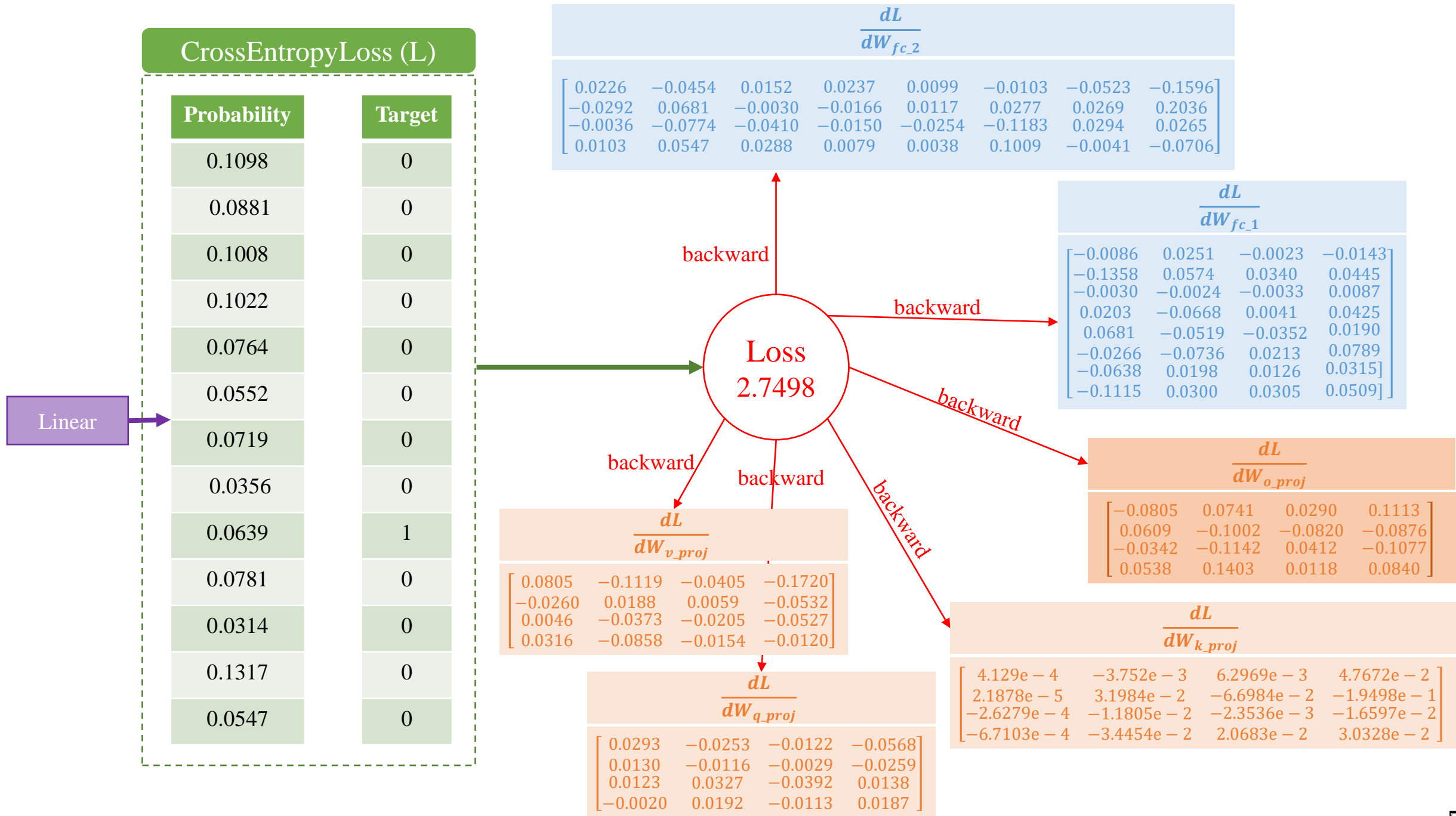
# test
src = torch.Tensor([[[0.48, 0.44, 0.71],
                     [0.65, 0.80, 0.79]]])
target = torch.Tensor([[[0.3516, 0.9509, 0.2771],
                        [0.1993, 0.0177, 0.2628],
                        [0.0774, 0.5253, 0.6413]]])
```





Example





Update Parameters
SGD($lr = 0.1$)

$$W_{fc_2} = W_{fc_2} - lr \times dW_{fc_2}$$

-0.3148	-0.2081	-0.0573	0.3396	0.0502	-0.0905	-0.1515	-0.1187
-0.2259	0.2512	-0.1605	-0.0692	-0.3529	0.2339	0.2652	0.1085
-0.2462	-0.3412	-0.2830	0.2651	0.1723	0.3093	0.1823	0.0868
-0.0045	-0.2743	-0.3058	-0.3315	0.1443	-0.1837	-0.0707	-0.1964

$$W_{fc_1} = W_{fc_1} - lr \times dW_{fc_1}$$

0.2800	-0.1303	-0.2851	-0.1698
-0.3600	0.1726	0.3836	-0.4752
0.1164	0.2585	0.0910	-0.1789
0.2589	0.2694	0.1866	-0.0921
-0.1392	0.0587	-0.0848	-0.1509
0.3223	0.4371	-0.0516	-0.1198
0.0137	-0.0318	0.1189	0.1370
-0.4430	-0.1875	0.4180	0.1897

$$W_{o_proj} = W_{o_proj} - lr \times dW_{o_proj}$$

0.0862	-0.4474	-0.2183	-0.3105
-0.0047	-0.1760	-0.0264	-0.3301
-0.3398	-0.2803	-0.1753	-0.3839
0.4139	-0.1133	0.4290	0.1474

$$W_{q_proj} = W_{q_proj} - lr \times dW_{q_proj}$$

-0.3011	-0.2957	0.4509	0.1723
0.4798	-0.4115	-0.4956	-0.3886
-0.3376	0.1992	0.1830	0.4141
-0.2580	-0.3428	0.2664	-0.2040

$$W_{k_proj} = W_{k_proj} - lr \times dW_{k_proj}$$

0.3034	-0.1183	0.2854	-0.3933
-0.2523	0.1492	0.1124	-0.1080
0.2981	0.3411	-0.3624	-0.2653
0.4579	-0.1653	-0.1793	-0.4868

$$W_{v_proj} = W_{v_proj} - lr \times dW_{v_proj}$$

-0.2944	0.1361	-0.0619	-0.3457
0.0143	-0.3434	-0.4248	-0.2700
-0.4381	-0.3146	0.5019	0.0997
0.1509	-0.4578	-0.3268	-0.1652

Forward again

CrossEntropyLoss (L)

Probability	Target
0.0503	0
0.0227	0
0.0684	0
0.0322	0
0.0225	0
0.0170	0
0.0445	0
0.0075	0
0.6174	1
0.0274	0
0.0108	0
0.0522	0
0.0270	0

Loss

0.1699



Precision and Recall of Words

Predict/Candidate/Output:

Tôi học NLP của AI VIET NAM

Reference:

Tôi đang học lớp AI của AI VIET NAM

Precision 1-gram

$$\frac{\text{correct}}{\text{candidate_length}} = \frac{6}{7}$$

Recall 1-gram

$$\frac{\text{correct}}{\text{reference_length}} = \frac{6}{9}$$

F1-score 1-gram

$$\frac{\text{precision} \times \text{recall}}{(\text{precision} + \text{recall})/2} = 0.75$$

BLEU Score

❖ BLEU score

Precision 1-gram

$$\frac{\text{correct}}{\text{candidate_length}} = \frac{6}{7}$$

Recall 1-gram

$$\frac{\text{correct}}{\text{reference_length}} = \frac{6}{9}$$

N-gram overlap between machine translation candidate and reference translation

Compute precision for n-grams of size 1 to 4

With 4-gram and add brevity penalty
(for too short translations):

$$\text{BLEU} = \min\left(1, \frac{\text{candidate_length}}{\text{reference_length}}\right) \left(\prod_{i=1}^4 \text{Precision}_i\right)^{1/4}$$

Precision and Recall of Words

Predict/Candidate/Output:

Tôi học NLP của AI VIET NAM

Reference:

Tôi đang học lớp CV và NLP của AI

Precision	1-gram	2-gram	3-gram	4-gram
	6/7	3/6	2/5	1/4

Multiple reference: N-grams may match in any of the reference and closest reference length used

Brevity penalty = 7/9

BLEU = 0.35

$$\text{BLEU} = \min \left(1, \frac{\text{candidate_length}}{\text{reference_length}} \right) \left(\prod_{i=1}^4 \text{Precision}_i \right)^{1/4}$$