# From Text Generation to Machine Translation

**Quang-Vinh Dinh**
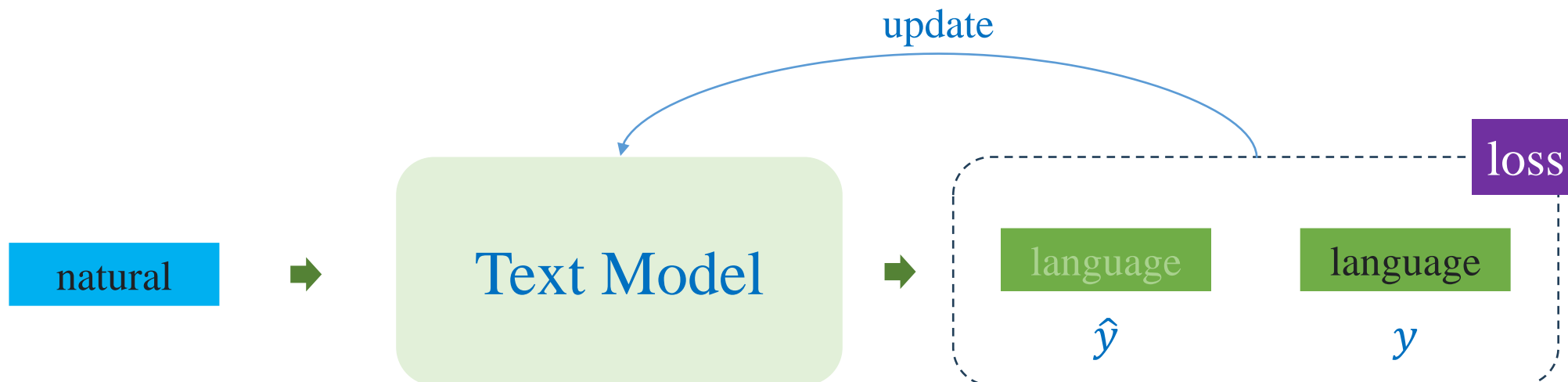**Ph.D. in Computer Science**

# Outline

- ➤ **Text Generation Using Transformer**
- ➤ **An improved Approach to Text Generation**
- ➤ **Machine Translation Using RNN**
- ➤ **Machine Translation Using Transformer**

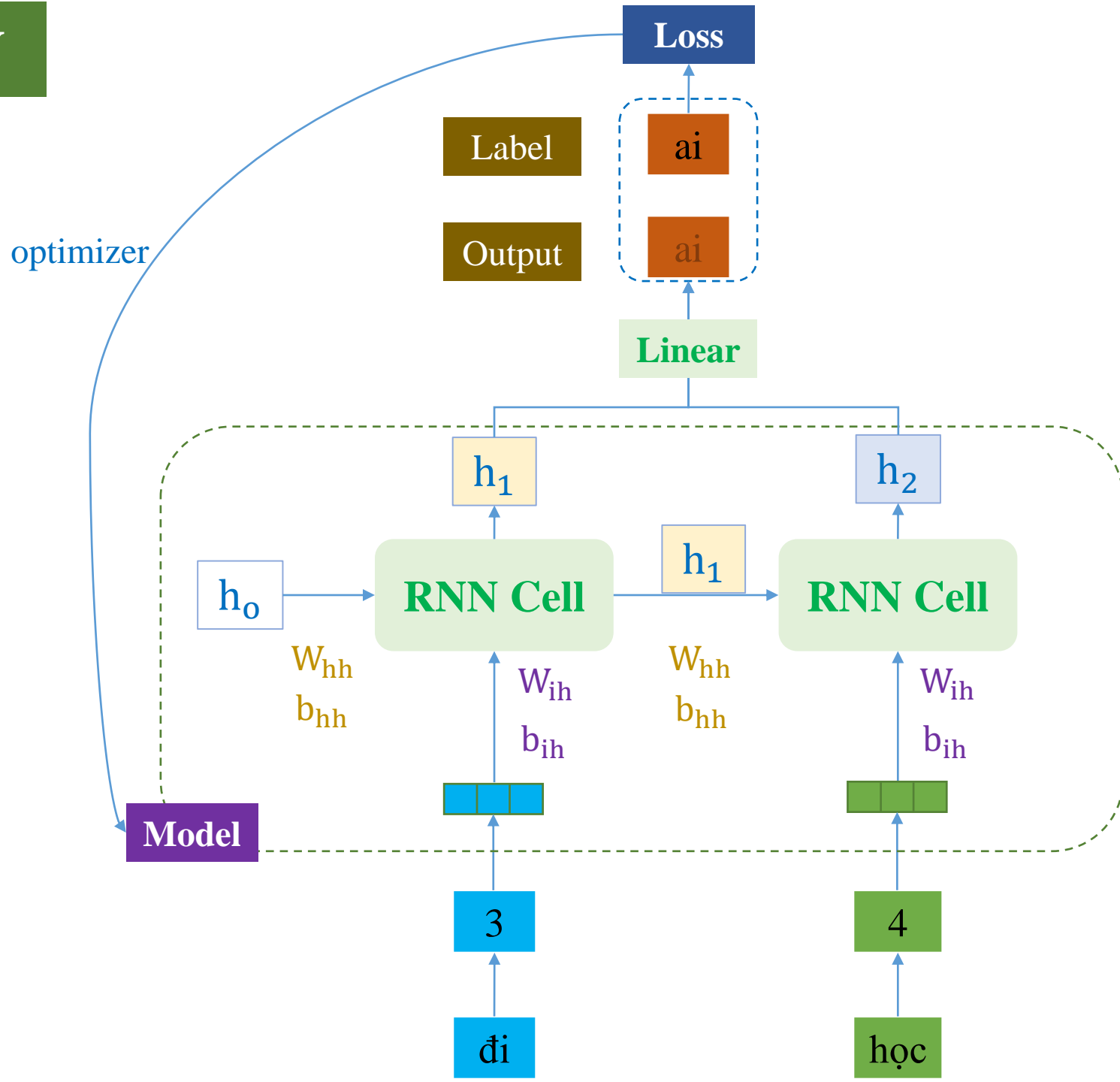# Self-Supervision Using Text Data

❖ **Encode the sequential relationship**
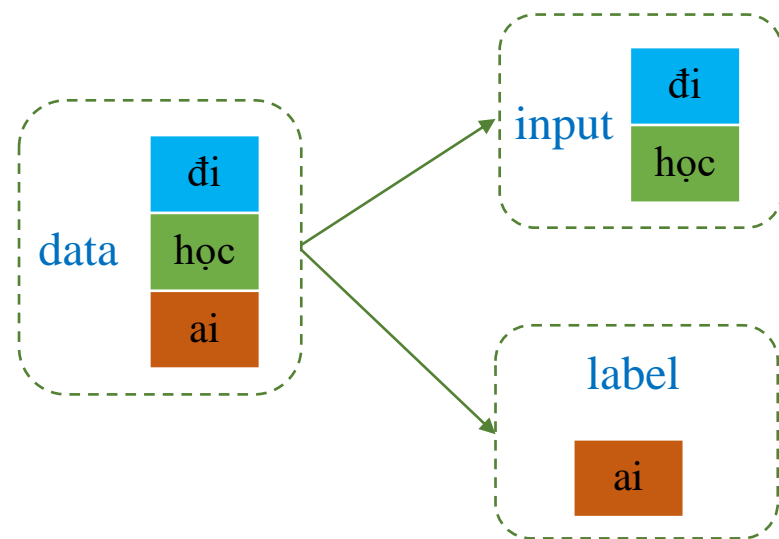
natural language processing is a branch of artificial intelligence

natural     language

X        y

update

natural → Text Model → language language

$\hat{y}$       $y$

loss

# Implementation Using RNN

Using all the features

# Example

## Copus

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

```python
tokenizer = get_tokenizer('basic_english')
def yield_tokens(examples):
    for text in examples:
        yield tokenizer(text)


# Create vocabulary
vocab_size = 13
build_vocab_from_iterator(yield_tokens(corpus),
        max_tokens=vocab_size,
        specials=["<unk>", "<pad>", "<sos>"])
```

**vocab_size = 13**

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

Tokenizer

Ăn quả nhớ trồng kẻ cây chí có nên thì

**Build vocab**

Special tokens: <unk> <pad> <sos>

**Add**

**Vocab**

| 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 |

**Next token prediction dataset**

```python
# create the next-token-prediction dataset
corpus = [
    "ăn quả nhớ kẻ trồng cây",
    "có chí thì nên"
]


data_x = []
data_y = []
for vector in corpus:
    vector = vector.split()

    for i in range(len(vector)):
        data_x.append(['<sos>'] + vector[:i])
        data_y.append(vector[i])
```

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

**vocab_size = 13**

4

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

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

| 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 |

**Next token prediction dataset**

**padding**

**Vocab**

sequence_length = 6

| data_x_ids | data_y_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**

4

# 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 |

**Next token prediction dataset**

```python
data_x_ids = []
data_y_ids = []
def vectorize(x, y, vocab, sequence_length):
    x_ids = [vocab[token] for token in x][:sequence_length]
    x_ids = x_ids + [vocab["<pad>"]]*(sequence_length - len(x))
    return x_ids, vocab[y]
for x, y in zip(data_x, data_y):
    x_ids, y_ids = vectorize(x, y, vocab, sequence_length)
    data_x_ids.append(x_ids)
    data_y_ids.append(y_ids)
```

**Training data**

sequence_length = 6

**Vocab**

**padding**

| data_x_ids | data_y_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 |

4

**Example**

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

vocab_size = 13

sequence_length = 6

**data_x_ids**

[2, 1, 1, 1, 1, 1]

**Model**
(RNN / Transformer)

**predictions**

[…, …, … … …, … ]

```python
class TG_Model(nn.Module):
    def __init__(self, vocab_size, sequence_length):
        super().__init__()
        self.recurrent = nn.RNN(4, 4, batch_first=True)
        self.linear = nn.Linear(sequence_length*4, vocab_size)

    def forward(self, x):
        x,_ = self.recurrent(x)  # [n, sequence_length, 4]
        x = nn.Flatten()(x)      # [n, 24]
        x = self.linear(x)       # [n, 13]
        return x

model = TG_Model(vocab_size, sequence_length)
outputs = model(data_x_ids)
```
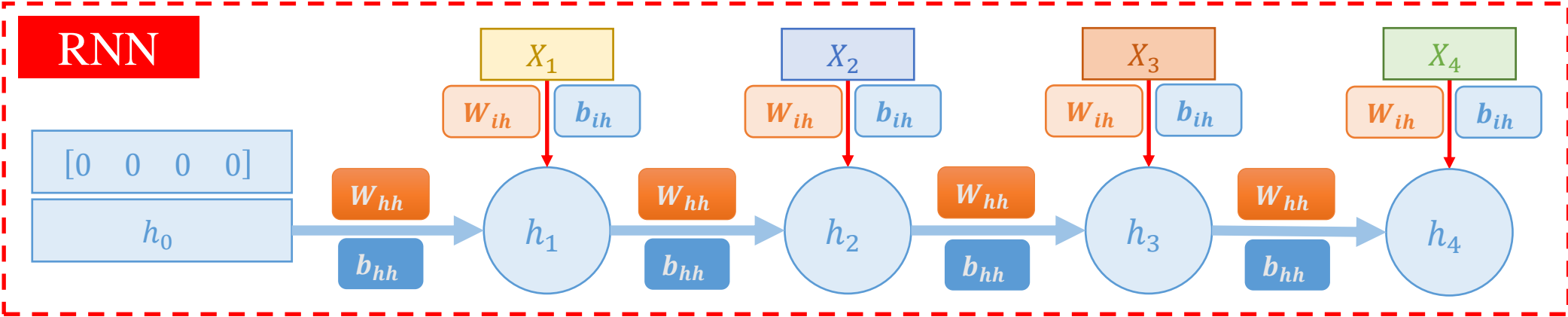
?

# Example

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

**vocab_size = 13**

sequence_length = 6

| data_x_ids |
|------------|
| [2, 1, 1, 1, 1, 1] |

**Model**
(RNN / Transformer)

| predictions |
|-------------|
| [..., ..., ... ... ..., ... ] |

```python
class TG_Model(nn.Module):
    def __init__(self, vocab_size, sequence_length):
        super().__init__()
        self.embedding = nn.Embedding(vocab_size, 4)
        self.recurrent = nn.RNN(4, 4, batch_first=True)
        self.linear = nn.Linear(sequence_length*4, vocab_size)


    def forward(self, x):
        x = self.embedding(x)      # [n, sequence_length, 4]
        x,_ = self.recurrent(x)    # [n, sequence_length, 4]
        x = nn.Flatten()(x)        # [n, 24]
        x = self.linear(x)         # [n, 13]
        return x


model = TG_Model(vocab_size, sequence_length)
outputs = model(data_x_ids)
```

sequence_length = 4

hidden_dim = 4

| Input tokens | Target token |
|---|---|
| <sos> Có chí thì | nên |

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

[2   5   3   10] → Embedding layer →

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

**RNN**

$X_1$  $W_{ih}$  $b_{ih}$
$X_2$  $W_{ih}$  $b_{ih}$
$X_3$  $W_{ih}$  $b_{ih}$
$X_4$  $W_{ih}$  $b_{ih}$

[0   0   0   0]

$h_0$  $W_{hh}$  $b_{hh}$  $h_1$  $W_{hh}$  $b_{hh}$  $h_2$  $W_{hh}$  $b_{hh}$  $h_3$  $W_{hh}$  $b_{hh}$  $h_4$

| $W_{hh}$ |
|---|
| $\begin{bmatrix} 0.3035 & -0.1187 & 0.2860 & -0.3885 \\ -0.2523 & 0.1524 & 0.1057 & -0.1275 \\ 0.2980 & 0.3399 & -0.3626 & -0.2669 \\ 0.4578 & -0.1687 & -0.1773 & -0.4838 \end{bmatrix}$ |

| $W_{ih}$ |
|---|
| $\begin{bmatrix} -0.2982 & -0.2982 & 0.4497 & 0.1666 \\ 0.4811 & -0.4126 & -0.4959 & -0.3912 \\ -0.3363 & 0.2025 & 0.1790 & 0.4155 \\ -0.2582 & -0.3409 & 0.2653 & -0.2021 \end{bmatrix}$ |

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

| $b_{hh}$ |
|---|
| [0.0117   -0.3415   -0.4242   -0.2753 ] |

| $b_{ih}$ |
|---|
| [-0.2863   0.1249   -0.0660   -0.3629] |

Flatten → Linear (16, 13)

6

# Back-Propagation

## CrossEntropyLoss (L)

| Probability | Target |
|:---:|:---:|
| 0.1121 | 0 |
| 0.1334 | 0 |
| 0.0742 | 0 |
| 0.0896 | 0 |
| 0.0648 | 0 |
| 0.0569 | 0 |
| 0.0956 | 0 |
| 0.0744 | 0 |
| 0.0751 | 1 |
| 0.0809 | 0 |
| 0.0636 | 0 |
| 0.0490 | 0 |
| 0.0304 | 0 |

FC

backward

Loss

2.5884

backward

backward

backward

$\nabla_{W_{hh}}L$

$$\begin{bmatrix} 0.0896 & -0.0204 & 0.0843 & 0.0201 \\ -0.0079 & -0.0960 & 0.0819 & -0.0723 \\ 0.0049 & -0.1445 & -0.0518 & -0.1252 \\ -0.0619 & 0.0827 & -0.1736 & 0.0291 \end{bmatrix}$$

$\nabla_{b_{hh}}L$

$$\begin{bmatrix} -0.0235 & 0.1788 & 0.3628 & 0.0449 \end{bmatrix}$$

$\nabla_{W_{ih}}L$

$$\begin{bmatrix} 0.0187 & -0.1589 & -0.0991 & -0.0988 \\ -0.1285 & 0.1600 & -0.1825 & -0.0411 \\ -0.2816 & 0.0904 & 0.1933 & 0.1014 \\ -0.0437 & 0.1782 & 0.2889 & -0.0192 \end{bmatrix}$$

$\nabla_{b_{ih}}L$

$$\begin{bmatrix} -0.0235 & 0.1788 & 0.3628 & 0.0449 \end{bmatrix}$$

# Example

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

**vocab_size = 13**

**data_x_ids**

[2, 1, 1, 1, 1, 1]

sequence_length = 6

**Model**
**(Transformer)**

**predictions**

[…, …, … … …, … ]

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-head Attention

Add & Norm

Masked Multi-head Attention

N×

Add & Norm

Feed Forward

Add & Norm

Multi-head Attention

N×

Positional Embedding

Input Embedding

Positional Embedding

Output Embedding

```
layer = nn.TransformerEncoderLayer(d_model=3,
                                   nhead=1,
                                   batch_first=True)

# feed forward
src = torch.Tensor([[[ 0.69,  0.72, -1.41],
                     [ 0.21,  1.10, -1.31]]])

out = layer(src)
```

# Encoder in Pytorch

| index | word |
|-------|------|
| 0 | [UNK] |
| 1 | [pad] |
| 2 | ai |
| 3 | đi |
| 4 | học |
| … | … |

| index | Embedding |
|-------|-----------|
| 0 | [-0.188, …,  0.7013] |
| 1 | [1.7840…,  1.3586] |
| 2 | [1.0281, …,  0.4211] |
| 3 | [-1.308,  …, -0.3680] |
| 4 | [0.2293, …,  2.0501] |
| … | … |



Transformer Encoder

15

```
layer = nn.TransformerEncoderLayer(d_model=3,
                                   nhead=1,
                                   batch_first=True)

# feed forward
src = torch.Tensor([[[ 0.69,  0.72, -1.41],
                     [ 0.21,  1.10, -1.31]]])

out = layer(src)
```

[0.69,  0.72, -1.41]

[0.21,  1.10, -1.31]

[0.97, 0.39, -1.37]

[0.58, 0.82, -1.40]



đi

học

3

4

**Embedding**

(N, 2, 3)

Multi-head
Self-Attention

Add &
Norm

Linear
(3, 4)

ReLU

Linear
(4, 3)

Add &
Norm

(N, 2, 3)

Feed Forward

Transformer Encoder

16

```python
layer = nn.TransformerEncoderLayer(d_model=3, nhead=1,
                                   batch_first=True)

# feed forward
src = torch.Tensor([[[ 0.69,  0.72, -1.41],
                     [ 0.21,  1.10, -1.31],
                     [-0.88,  0.60, -0.31]]])
mask = torch.triu(torch.ones(3, 3), diagonal=1).bool()
out = layer(src, src_mask=mask)
```

# Masked Encoder in Pytorch

| index | word  |
|-------|-------|
| 0     | [UNK] |
| 1     | [pad] |
| 2     | ai    |
| 3     | đi    |
| 4     | học   |
| …     | …     |

| index | Embedding              |
|-------|------------------------|
| 0     | [-0.188, …,  0.7013]   |
| 1     | [1.7840…,  1.3586]     |
| 2     | [1.0281, …,  0.4211]   |
| 3     | [-1.308,  …, -0.3680]  |
| 4     | [0.2293,  …,  2.0501]  |
| …     | …                      |

# Masked Encoder in Pytorch

```python
layer = nn.TransformerEncoderLayer(d_model=3, nhead=1,
                                    batch_first=True)

# feed forward
src = torch.Tensor([[[ 0.69,  0.72, -1.41],
                     [ 0.21,  1.10, -1.31],
                     [-0.88,  0.60, -0.31]]])
mask = torch.triu(torch.ones(3, 3), diagonal=1).bool()
out = layer(src, src_mask=mask)
```

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

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



đi → 3
học → 4
ai → 2

**Embedding**

(N, 3, 3)

Masked Multi-head Self-Attention

Add & Norm

Feed Forward
Linear (3, 4)  ReLU  Linear (4, 3)

Add & Norm

(N, 3, 3)

Transformer Encoder

18

# Masked Multi-head Attention

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

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

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

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

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

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

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

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

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

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

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

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

12

# Masked Multi-head Attention

❖ **Example**

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

$$= softmax\left(\begin{bmatrix} -0.08 & -0.14 & -0.24 \\ -0.39 & 0.77 & 0.69 \end{bmatrix}\begin{bmatrix} 0.02 & 0.27 \\ -0.01 & 0.27 \\ 0.13 & -0.26 \end{bmatrix}\frac{1}{\sqrt{d}} + \begin{bmatrix} 0 & -\infty \\ 0 & 0 \end{bmatrix}\right)\begin{bmatrix} -0.16 & -0.08 & -0.05 \\ -0.02 & -0.02 & 0.05 \end{bmatrix}$$

$$= softmax\left(\begin{bmatrix} -0.019 & 0.002 \\ 0.043 & -0.046 \end{bmatrix} + \begin{bmatrix} 0 & -\infty \\ 0 & 0 \end{bmatrix}\right)\begin{bmatrix} -0.16 & -0.08 & -0.05 \\ -0.02 & -0.02 & 0.05 \end{bmatrix}$$

$$= \begin{bmatrix} 1.0 & 0.0 \\ 0.52 & 0.48 \end{bmatrix}\begin{bmatrix} -0.16 & -0.08 & -0.05 \\ -0.02 & -0.02 & 0.05 \end{bmatrix} = \begin{bmatrix} -0.16 & -0.08 & -0.05 \\ 0.12 & 0.08 & 0.06 \end{bmatrix}$$

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

```python
nn.TransformerEncoderLayer(d_model=3,
                           nhead=1,
                           batch_first=True)
nn.TransformerDecoderLayer(d_model=3,
                           nhead=1,
                           batch_first=True)

# feed forward
src = torch.Tensor([[[0.48, 0.44, 0.71],
                     [0.65, 0.80, 0.79]]])

tgt = torch.Tensor([[[0.3516, 0.9509, 0.2771],
                     [0.1993, 0.0177, 0.2628],
                     [0.0774, 0.5253, 0.6413],
                     [0.6749, 0.5501, 0.1641]]])

context = encoder_layer(src)
mask = torch.triu(torch.ones(4, 4),
                  diagonal=1).bool()
out = decoder_layer(tgt, context,
                    tgt_mask=mask)
```
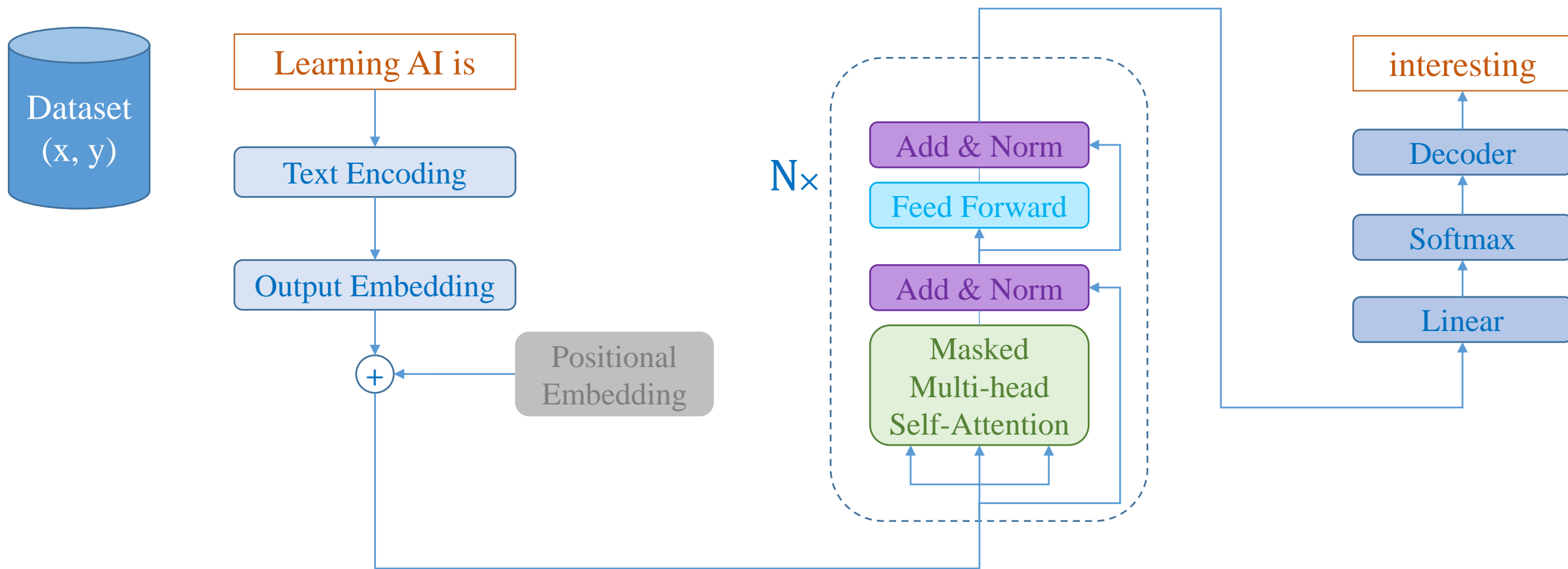
# Transformer in PyTorch

❖ **Transformer Encoder**



```python
# 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)
```

# Text Generation

❖ **Architecture**

# Example

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

**vocab_size = 13**

**data_x_ids**

[2, 1, 1, 1, 1, 1]

sequence_length = 6

**Model**
**(Transformer)**

**predictions**

[…, …, … … …, … ]

```python
class TG_Model(nn.Module):
    def __init__(self, vocab_size, embed_dim, num_heads, sequence_length):
        super().__init__()
        self.embedding = nn.Embedding(vocab_size, embed_dim)
        self.mask = torch.triu(torch.ones(sequence_length, sequence_length),
                               diagonal=1).bool()
        self.transformer =  nn.TransformerEncoderLayer(d_model=embed_dim,
                                                       nhead=num_heads,
                                                       batch_first=True,
                                                       dim_feedforward=4)

        self.linear = nn.Linear(sequence_length*embed_dim, vocab_size)


    def forward(self, x):
        x = self.embedding(x)                        # [n, seq_len, embed_dim]
        x = self.transformer(x, src_mask=self.mask)  # [n, seq_len, embed_dim]
        x = nn.Flatten()(x)                          # [n, seq_len*embed_dim]
        x = self.linear(x)                           # [n, vocab_size]
        return x


model = TG_Model(vocab_size, 8, 2, sequence_length)
```
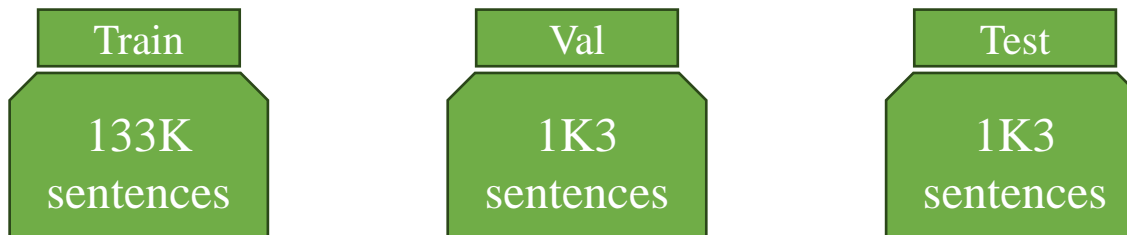
# Outline

- ➢ **Text Generation Using Transformer**
- ➢ **An improved Approach to Text Generation**
- ➢ **Machine Translation Using RNN**
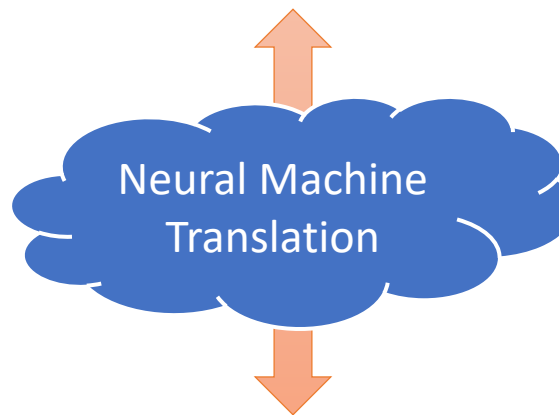- ➢ **Machine Translation Using Transformer**

An improved approach

*Linear can be added (where?)

# An improved approach

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]

học                    ai

$h_1$                  $h_2$

$h_o$  →  **RNN Cell**  —$h_1$→  **RNN Cell**

$W_{hh}$        $W_{ih}$        $W_{hh}$        $W_{ih}$
$b_{hh}$        $b_{ih}$        $b_{hh}$        $b_{ih}$

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

3                      4

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

seq_len = 2     đi        học        *Linear can be added

# Outline

- ➢ **Text Generation Using Transformer**
- ➢ **An improved Approach to Text Generation**
- ➢ **Machine Translation Using RNN**
- ➢ **Machine Translation Using Transformer**

(N, 2, 3)

value
key
query

Multi-head Self-Attention

Add & Norm

Linear (3, 4)  ReLU  Linear (4, 3)

Feed Forward

Add & Norm

(N, 2, 3)

**Source sequences**

good morning <eos>

ai books <eos>

key    value

value
key
query

Masked Multi-head Self-Attention

Add & Norm

query

Masked Multi-head Self-Attention

Add & Norm

Linear (3, 4)  ReLU  Linear (4, 3)

Feed Forward

Add & Norm

(N, 3, 3)

**Target sequences**

<sos> chào buổi sáng <eos>

<sos> sách ai <eos> <pad>