# Extra Class

# Introduction to Transformer

Nguyen Quoc Thai

# CONTENT

**AI VIET NAM**
@aivietnam.edu.vn

!  **RNNs Model**

Output Sequence

Hidden State → **RNN** → **RNN** → **RNN** → Last Hidden State

Dense vector

Embedding Layer

Input $X_1$ $X_2$ $X_N$

**AI VIET NAM**
@aivietnam.edu.vn

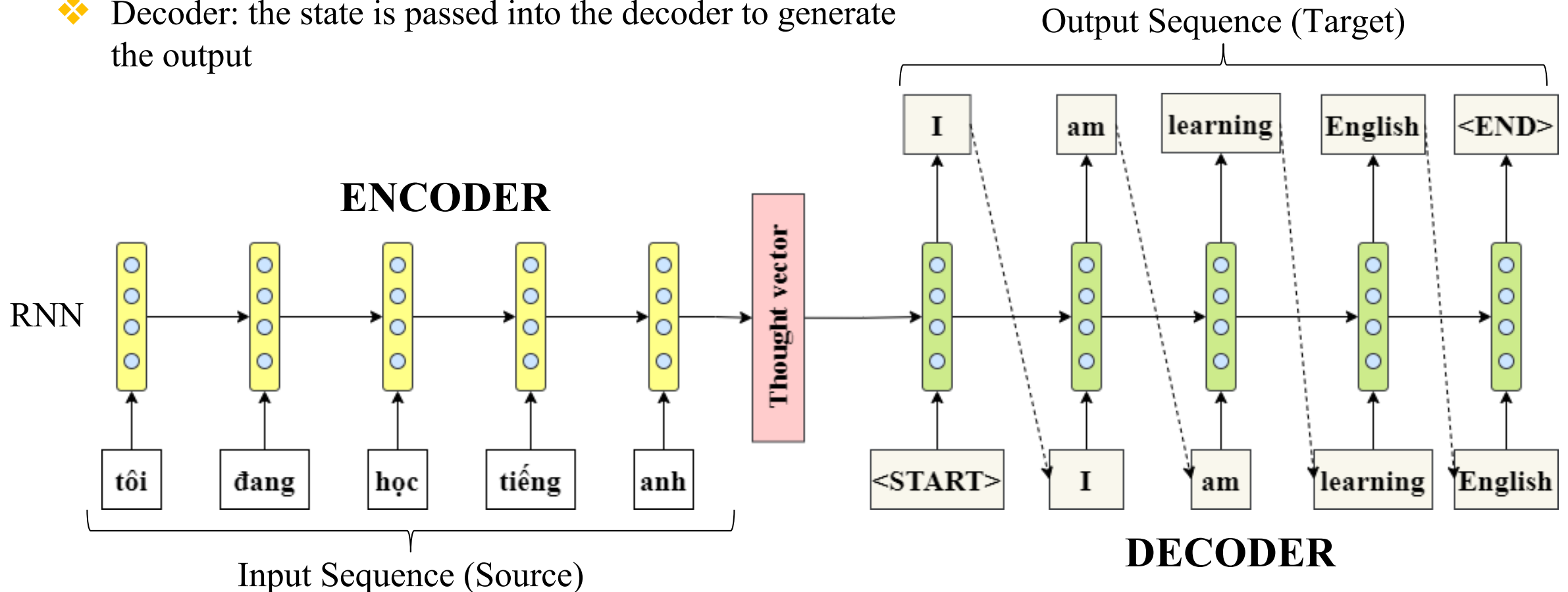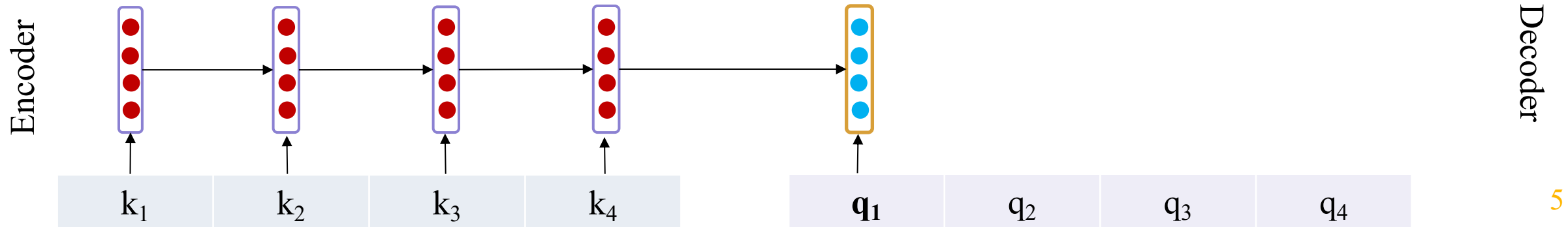!  **Sequence-to-Sequence Architecture**

❖ Encoder: encoding the inputs into state (thought vector)
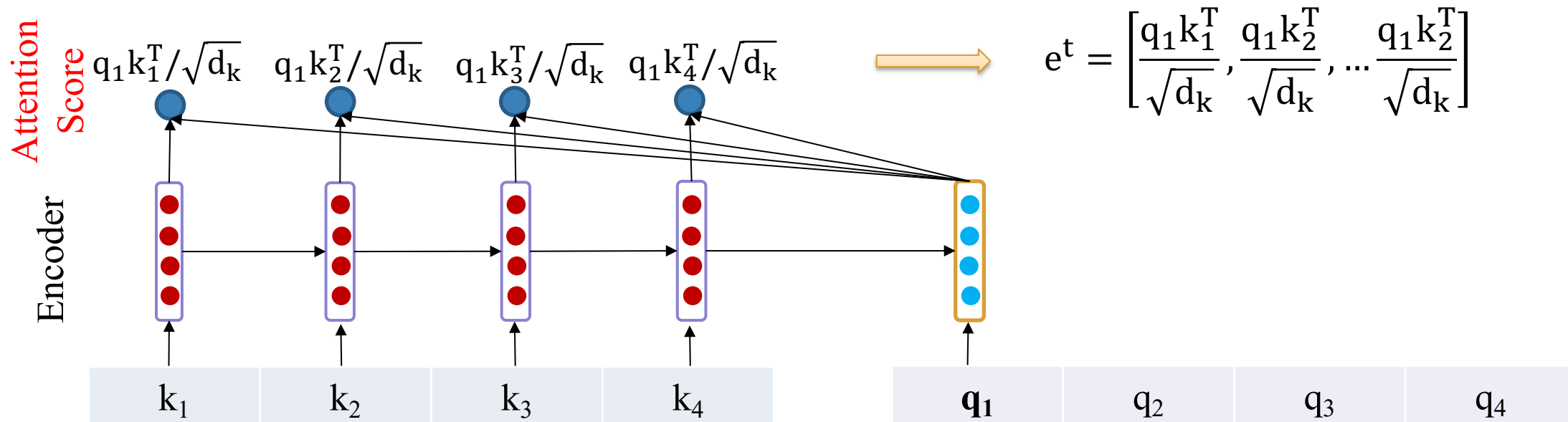❖ Decoder: the state is passed into the decoder to generate the output



RNN

ENCODER

DECODER

Thought vector

Input Sequence (Source)

Output Sequence (Target)

## Scaled Dot-Product Attention

**AI VIET NAM**
@aivietnam.edu.vn

!

**Scaled Dot-Product Attention**



$q_1 k_1^T/\sqrt{d_k}$    $q_1 k_2^T/\sqrt{d_k}$    $q_1 k_3^T/\sqrt{d_k}$    $q_1 k_4^T/\sqrt{d_k}$

$$e^t = \left[ \frac{q_1 k_1^T}{\sqrt{d_k}}, \frac{q_1 k_2^T}{\sqrt{d_k}}, \dots \frac{q_1 k_2^T}{\sqrt{d_k}} \right]$$

Attention Score

Encoder

Decoder

$k_1$    $k_2$    $k_3$    $k_4$    $q_1$    $q_2$    $q_3$    $q_4$

6

**AI VIET NAM**
@aivietnam.edu.vn

**!** Scaled Dot-Product Attention



Attention distribution

$\alpha_1$   $\alpha_2$   $\alpha_3$   $\alpha_4$

$\alpha^t = \text{softmax}(e^t)$

Attention Score

$q_1 k_1^T / \sqrt{d_k}$   $q_1 k_2^T / \sqrt{d_k}$   $q_1 k_3^T / \sqrt{d_k}$   $q_1 k_4^T / \sqrt{d_k}$

$$e^t = \left[ \frac{q_1 k_1^T}{\sqrt{d_k}}, \frac{q_1 k_2^T}{\sqrt{d_k}}, \dots \frac{q_1 k_2^T}{\sqrt{d_k}} \right]$$

Encoder

Decoder

$k_1$   $k_2$   $k_3$   $k_4$   $\mathbf{q_1}$   $q_2$   $q_3$   $q_4$

7

AI VIET NAM
@aivietnam.edu.vn

**Scaled Dot-Product Attention**



Attention output

Attention distribution

$\alpha_1$    $\alpha_2$    $\alpha_3$    $\alpha_4$

Attention Score

$q_1 k_1^T / \sqrt{d_k}$    $q_1 k_2^T / \sqrt{d_k}$    $q_1 k_3^T / \sqrt{d_k}$    $q_1 k_4^T / \sqrt{d_k}$

Encoder

Decoder

$$a^t = \sum_1^N \alpha_i^t v_i \; ; \; v_i = k_i$$

$$\alpha^t = \text{softmax}(e^t)$$

$$e^t = \left[ \frac{q_1 k_1^T}{\sqrt{d_k}}, \frac{q_1 k_2^T}{\sqrt{d_k}}, \dots \frac{q_1 k_2^T}{\sqrt{d_k}} \right]$$

$k_1$    $k_2$    $k_3$    $k_4$      $q_1$    $q_2$    $q_3$    $q_4$

8

**AI VIET NAM**
@aivietnam.edu.vn

## Scaled Dot-Product Attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^{T}}{\sqrt{d_k}}\right)V; \ K = V$$
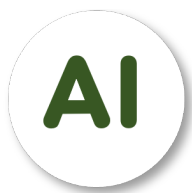
**AI VIET NAM**
@aivietnam.edu.vn

!

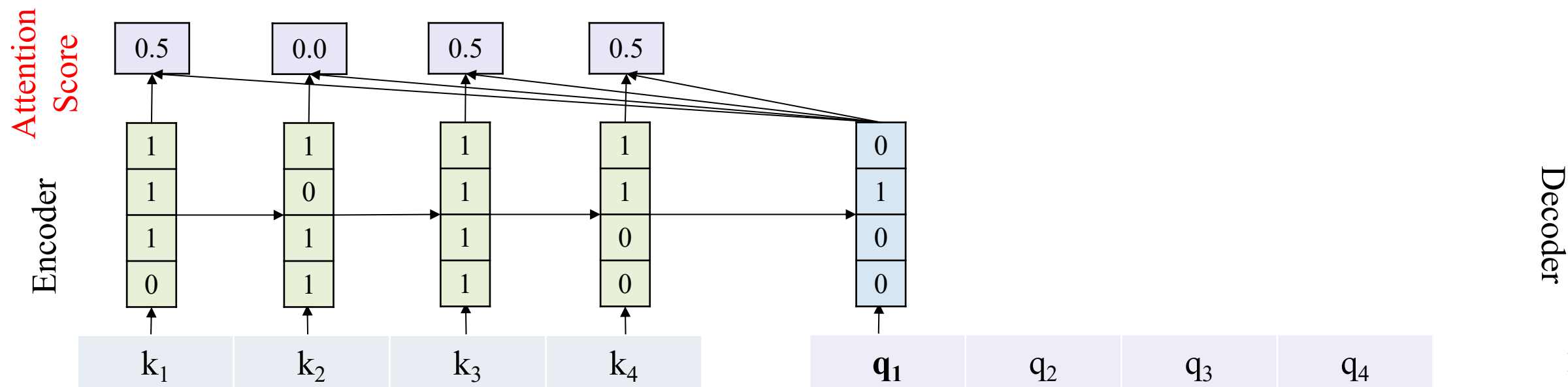**Scaled Dot-Product Attention - Example**

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V; \ K = V$$



Attention distribution
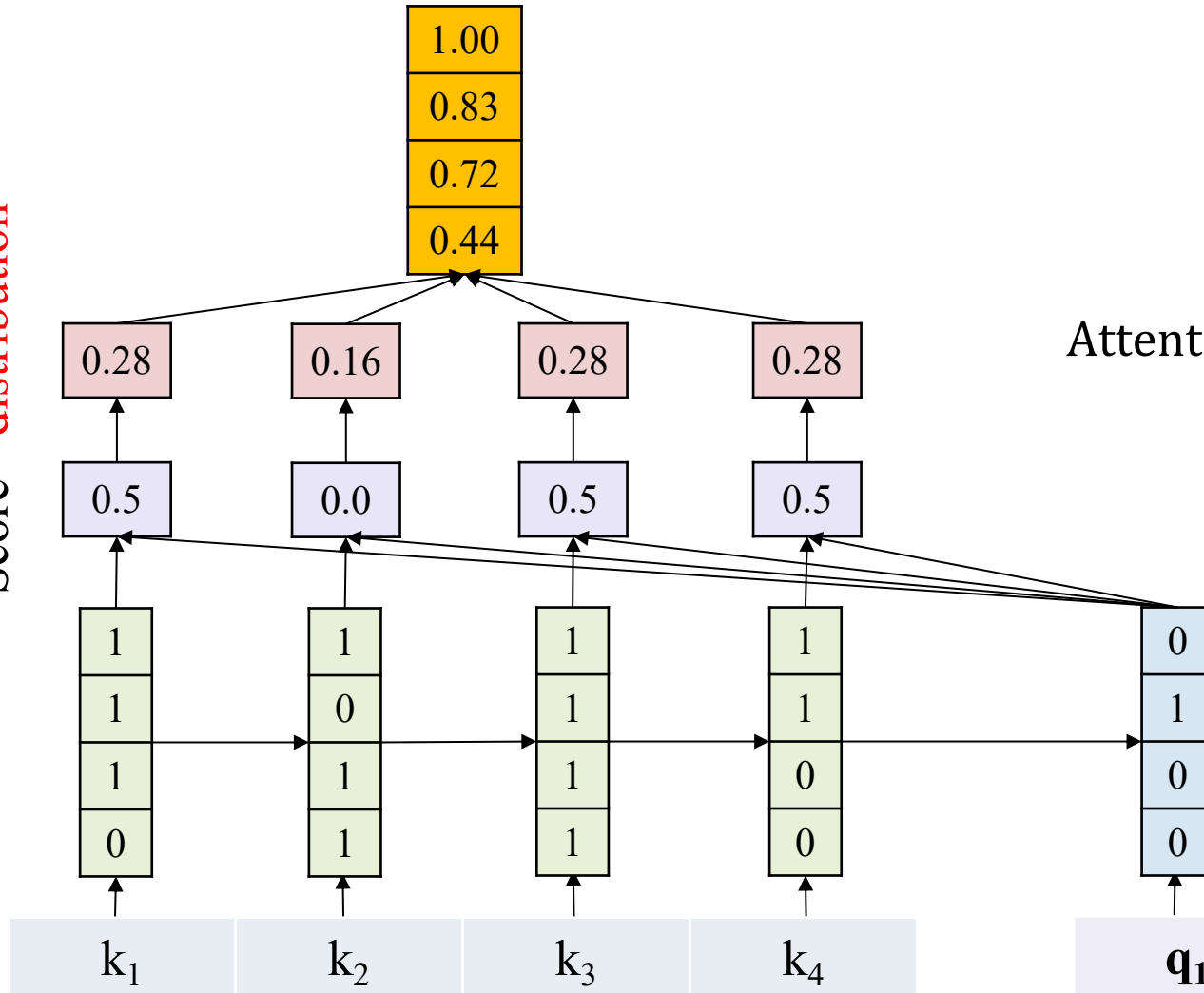
Attention Score

Encoder

Decoder

# 1 – Attention
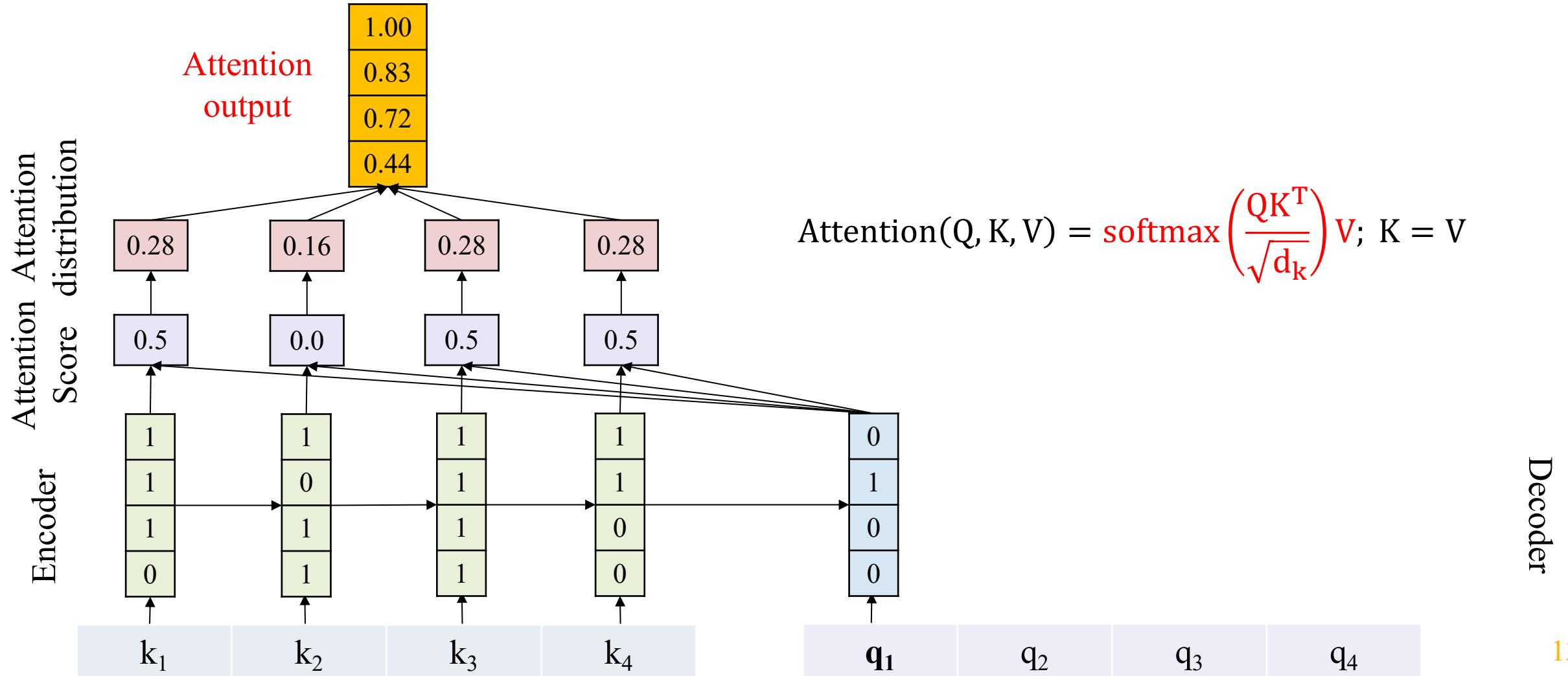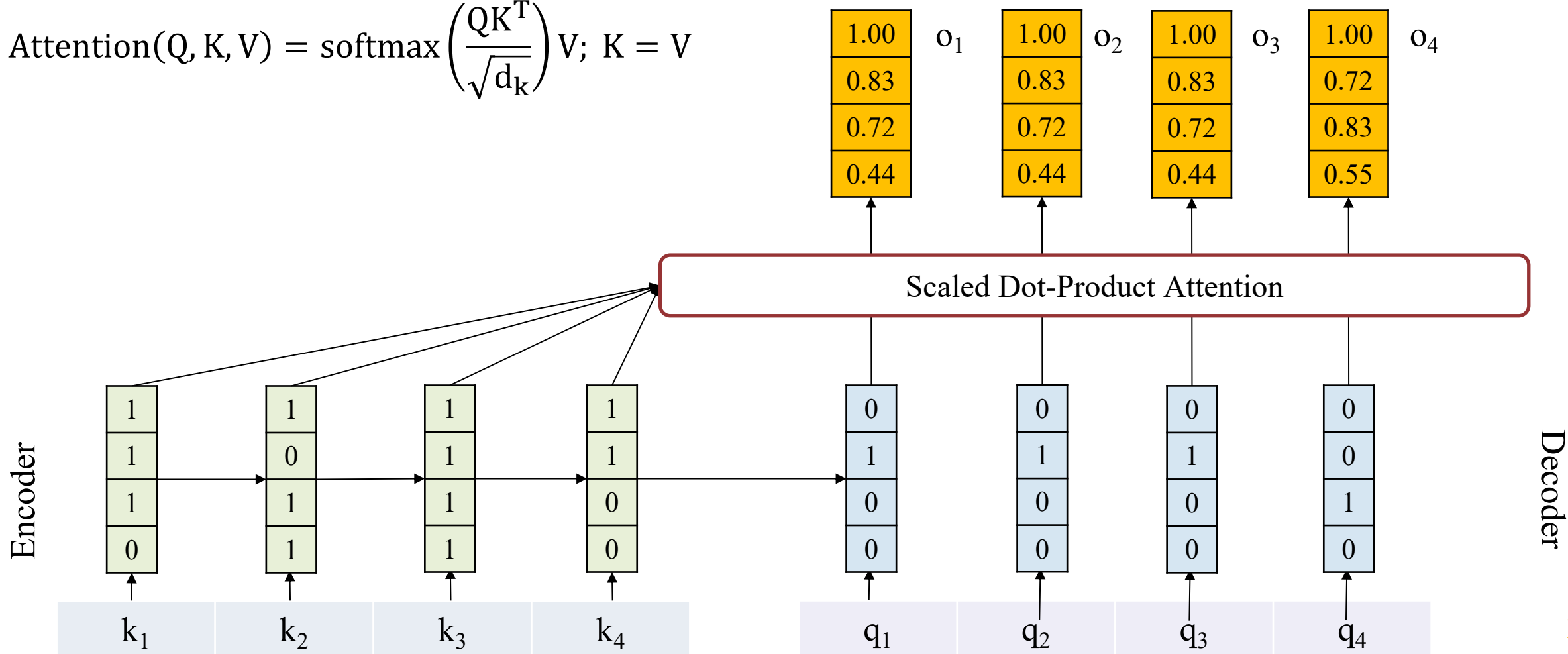
**AI VIET NAM**

! **Scaled Dot-Product Attention - Example**

Attention output

1.00
0.83
0.72
0.44

Attention distribution

0.28 | 0.16 | 0.28 | 0.28

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V; \ K = V$$

Attention Score

0.5 | 0.0 | 0.5 | 0.5

Encoder

| 1 | | 1 | | 1 | | 1 |
| 1 | | 0 | | 1 | | 1 |
| 1 | | 1 | | 1 | | 0 |
| 0 | | 1 | | 1 | | 0 |

Decoder

0
1
0
0

$k_1$ | $k_2$ | $k_3$ | $k_4$ | $\mathbf{q_1}$ | $q_2$ | $q_3$ | $q_4$

13

# 1 – Attention

AI VIET NAM
@aivietnam.edu.vn

**Scaled Dot-Product Attention - Example**

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V; \ K = V$$

**AI VIET NAM**
@aivietnam.edu.vn

**!**  **Scaled Dot-Product Attention - Demo**

```python
query = torch.randint(
    high=2,
    size=(1, 4, 4), # batch_size x seq_len x embedding_dim
    dtype=torch.float32
)
query
```

```
tensor([[[0., 1., 0., 0.],
         [0., 1., 0., 0.],
         [0., 1., 0., 0.],
         [0., 0., 1., 0.]]])
```

```python
key = torch.randint(
    high=2,
    size=(1, 4, 4),
    dtype=torch.float32
)
key
```

```
tensor([[[1., 1., 1., 0.],
         [1., 0., 1., 1.],
         [1., 1., 1., 1.],
         [1., 1., 0., 0.]]])
```

```python
value = key
value
```

```
tensor([[[1., 1., 1., 0.],
         [1., 0., 1., 1.],
         [1., 1., 1., 1.],
         [1., 1., 0., 0.]]])
```

```python
attentionn_weight = F.scaled_dot_product_attention(
    query=query,
    key=key,
    value=value
)
attentionn_weight
```
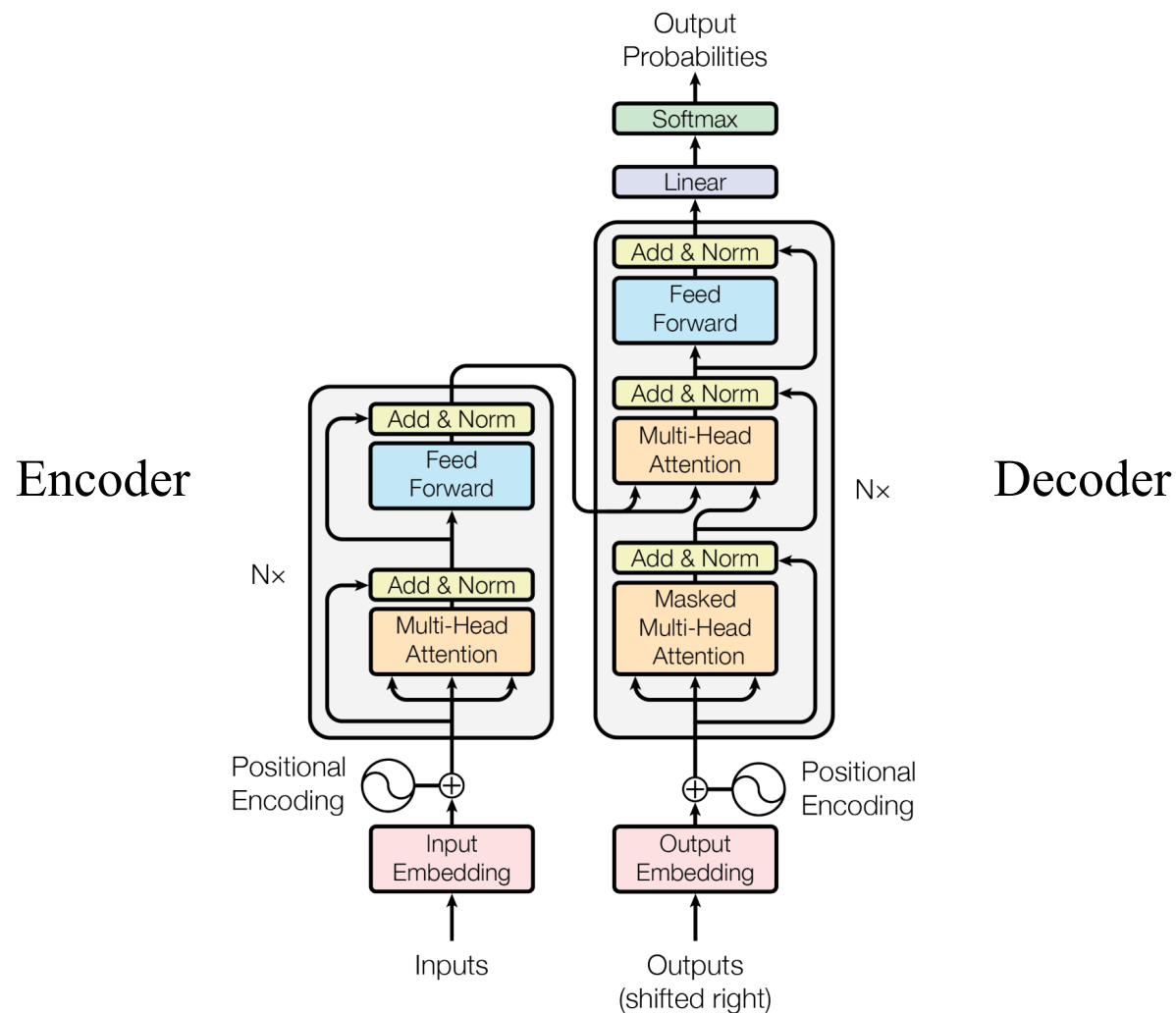
```
tensor([[[1.0000, 0.8318, 0.7227, 0.4455],
         [1.0000, 0.8318, 0.7227, 0.4455],
         [1.0000, 0.8318, 0.7227, 0.4455],
         [1.0000, 0.7227, 0.8318, 0.5545]]])
```

15

!

**Transformer**

❖ Architecture:
    N Encoder Layer
    N Decoder Layer
❖ Core technique: attention
❖ Loss function: cross-entropy

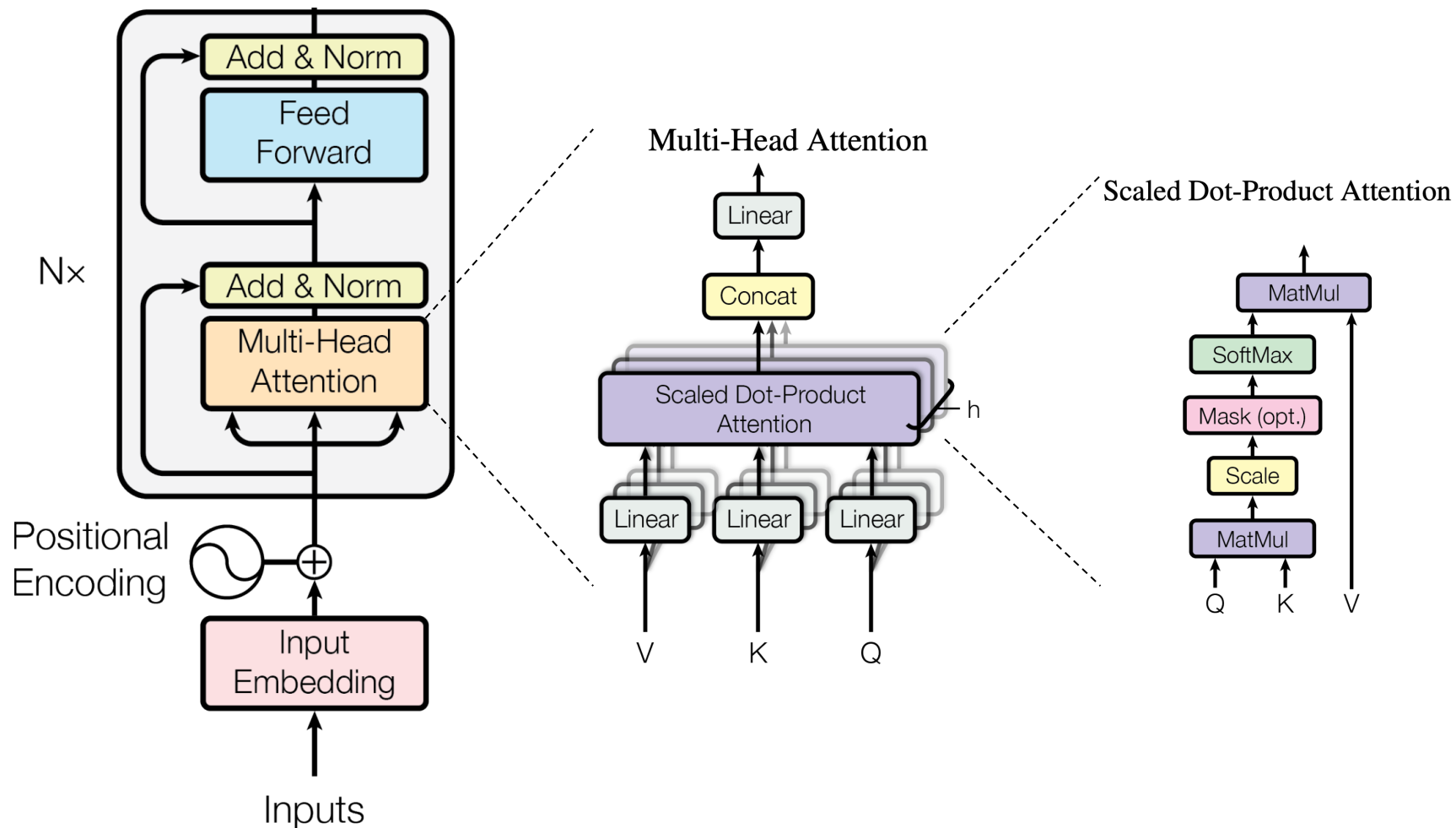Encoder                             Decoder



16

**AI VIET NAM**
@aivietnam.edu.vn

**!**

## Transformer-Encoder

- ❖ Input Embedding
- ❖ Positional Encoding
- ❖ Multi-Head Attention
- ❖ Feed Forward
- ❖ Add & Norm



17

**AI VIET NAM**
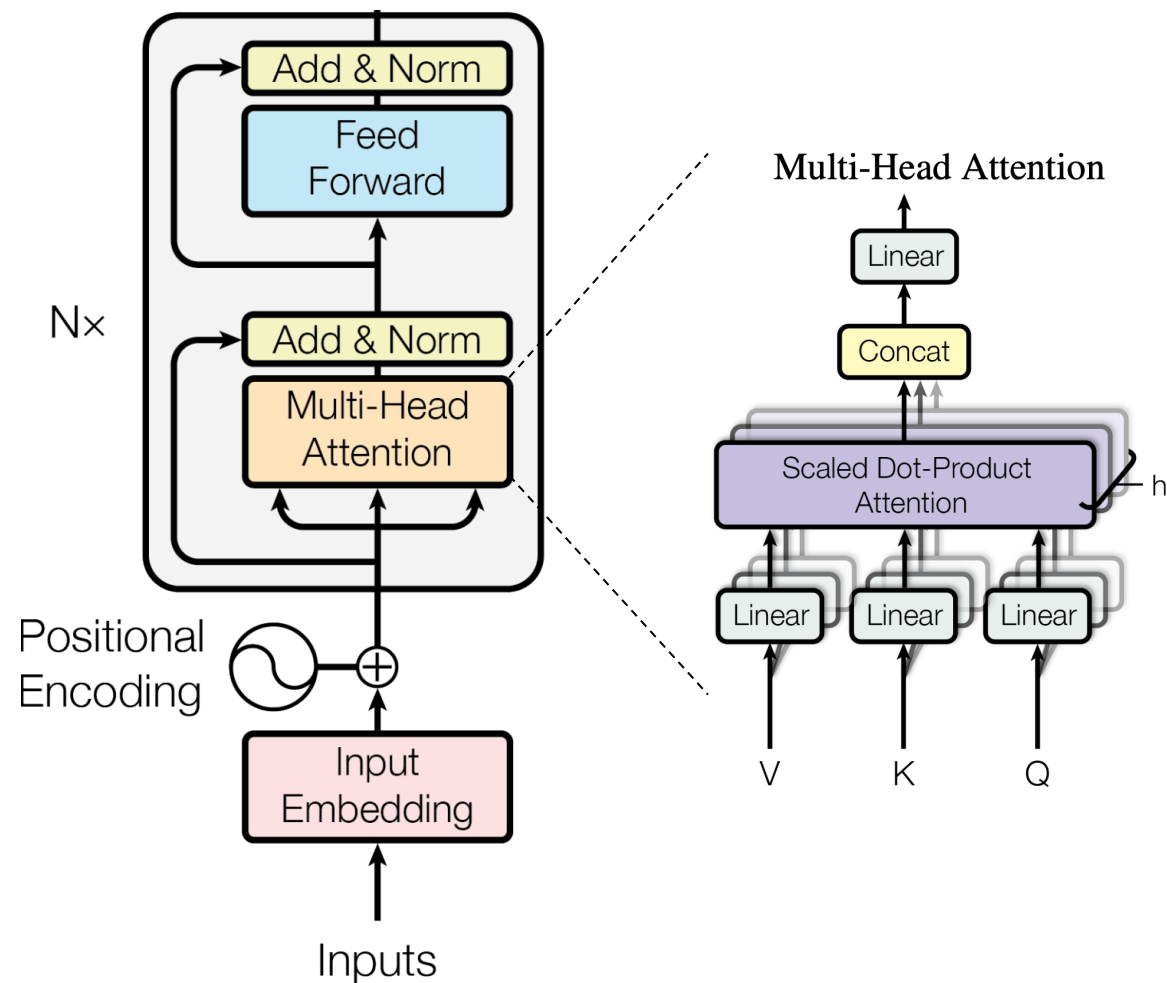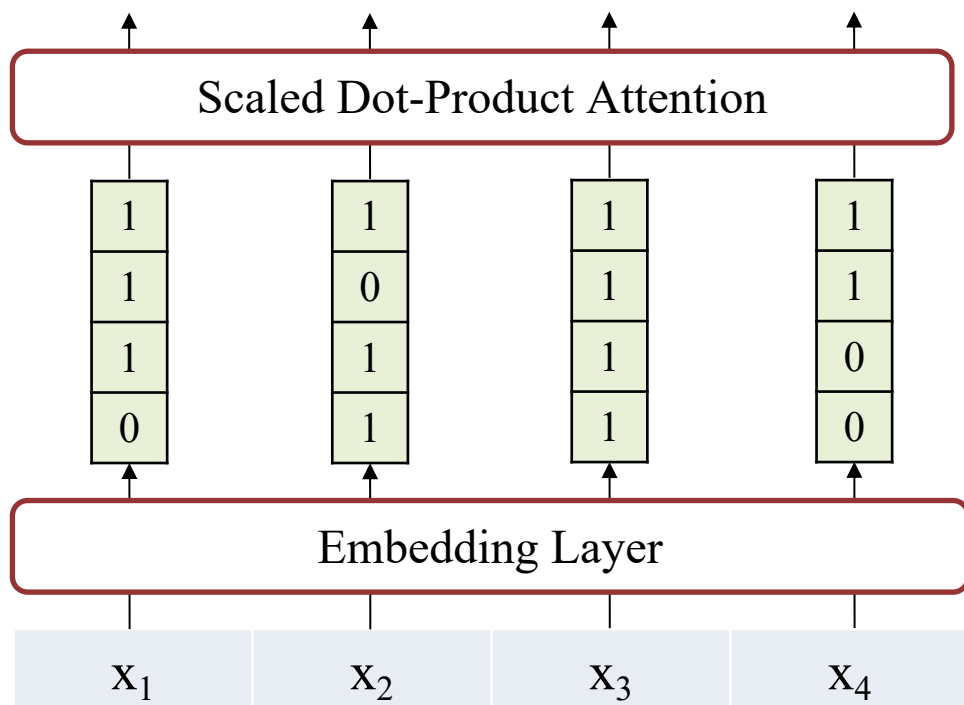@aivietnam.edu.vn

!

## Input Embedding

❖ Input Embedding: Embedding Layer

**QUERY – KEY – VALUE ?**

**AI VIET NAM**
@aivietnam.edu.vn

**!**

**Self-Attention**

**QUERY = KEY = VALUE = EMBEDDED**

**AI VIET NAM**
@aivietnam.edu.vn

**Self-Attention**

AI VIET NAM
@aivietnam.edu.vn

**Self-Attention**



22

**AI VIET NAM**
@aivietnam.edu.vn

**!** Self-Attention

**Self-Attention**

**! Self-Attention**

❖ To learn the relationship between word in the sentence

**Ignore the order of words in the sentence ?**

**Positional Encoding**

❖ The position of a token in a sentence as unique representation – each position is mapped to a vector
❖ Methods: Sinusoid; Learned positional embedding (as learned input embedding)

| $x_1$ | Embedding Layer #1 |
| $x_2$ | |
| $x_3$ | |

| 1 | 1 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| 1 | 1 | 1 |

| $x_1$ | Embedding Layer #2 |
| $x_2$ | |
| $x_3$ | |

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

**+**

| 2 | 3 | 2 |
|---|---|---|
| 0 | 1 | 0 |
| 3 | 2 | 1 |

Self-Attention

| 0.2 | 0.3 | 1.2 |
|-----|-----|-----|
| 0.7 | 1.9 | 0.8 |
| 0.3 | 2.1 | 1.2 |

28

AI VIET NAM
@aivietnam.edu.vn

**!**

## Positional Encoding – Demo

```python
class TokenAndPositionEmbedding(nn.Module):
    def __init__(self, vocab_size, embed_dim, max_length, device='cpu'):
        super().__init__()
        self.device = device
        self.word_emb = nn.Embedding(
            num_embeddings=vocab_size,
            embedding_dim=embed_dim
        )
        self.pos_emb = nn.Embedding(
            num_embeddings=max_length,
            embedding_dim=embed_dim
        )

    def forward(self, x):
        N, seq_len = x.size()
        positions = torch.arange(0, seq_len).expand(N, seq_len).to(self.device)
        output1 = self.word_emb(x)
        output2 = self.pos_emb(positions)
        output = output1 + output2
        return output
```

```python
vocab_size = 10000
embed_dim = 200
max_length = 50
embedding = TokenAndPositionEmbedding(
    vocab_size,
    embed_dim,
    max_length
)
```

```python
batch_size = 32

input = torch.randint(
    high=2,
    size=(batch_size, max_length),
    dtype=torch.int64
)
```

```python
embedded = embedding(input)
```

```python
embedded.shape
```
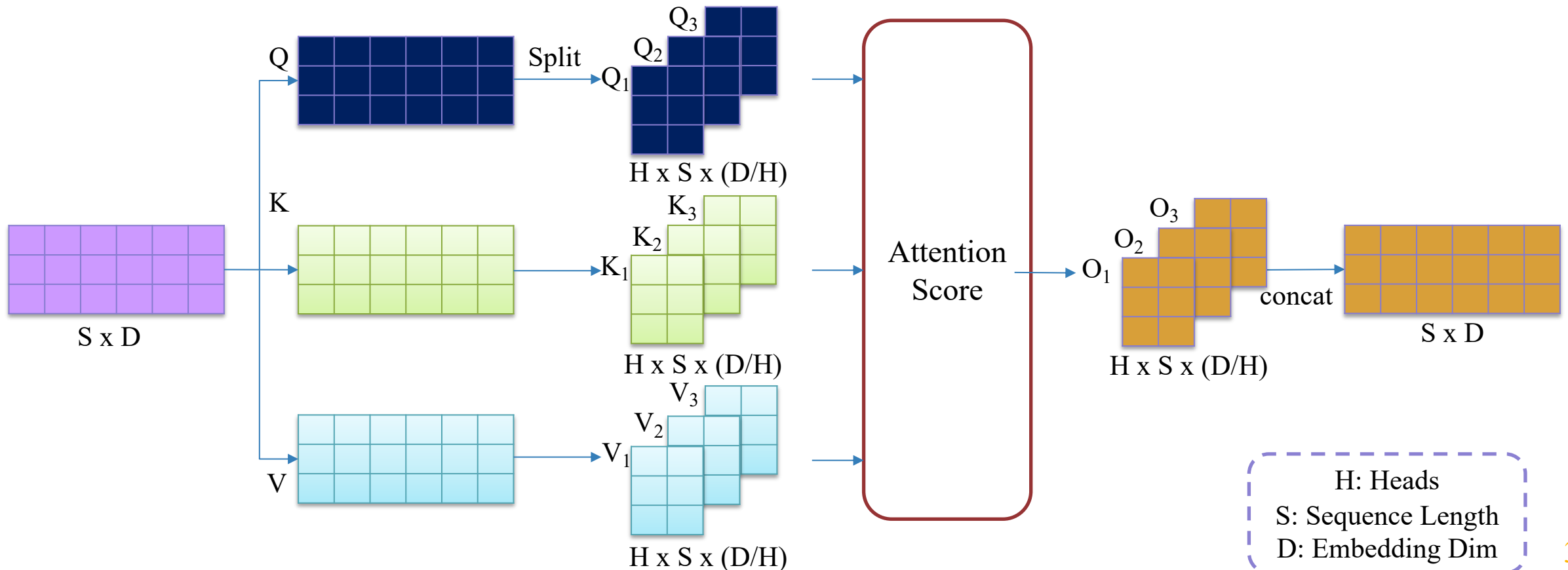
```
torch.Size([32, 50, 200])
```

29

AI VIET NAM
@aivietnam.edu.vn

**Multi-Head Attention**

❖ Split into the multiple attention heads (process independently) => self-attention => concat



H: Heads
S: Sequence Length
D: Embedding Dim

30

## Multi-Head Attention – Demo

```python
batch_size = 1
seq_len = 50
embedding_dim = 200

input = torch.randint(
    high=2,
    size=(batch_size, seq_len, embedding_dim),
    dtype=torch.float32
)
input
```

```
tensor([[[0., 1., 1.,  ..., 0., 1., 1.],
         [0., 1., 0.,  ..., 0., 0., 0.],
         [1., 0., 1.,  ..., 1., 1., 1.],
         ...,
         [0., 0., 0.,  ..., 1., 0., 0.],
         [1., 0., 1.,  ..., 0., 1., 1.],
         [0., 1., 1.,  ..., 1., 1., 1.]]])
```

```python
embedding_dim = 200
num_heads = 5

att_layer = nn.MultiheadAttention(
    embed_dim=embedding_dim,
    num_heads=num_heads,
    batch_first=True
)
```

```python
attn_output, attn_output_weights = att_layer(
    query=input,
    key=input,
    value=input
)
```

```python
attn_output.shape
```
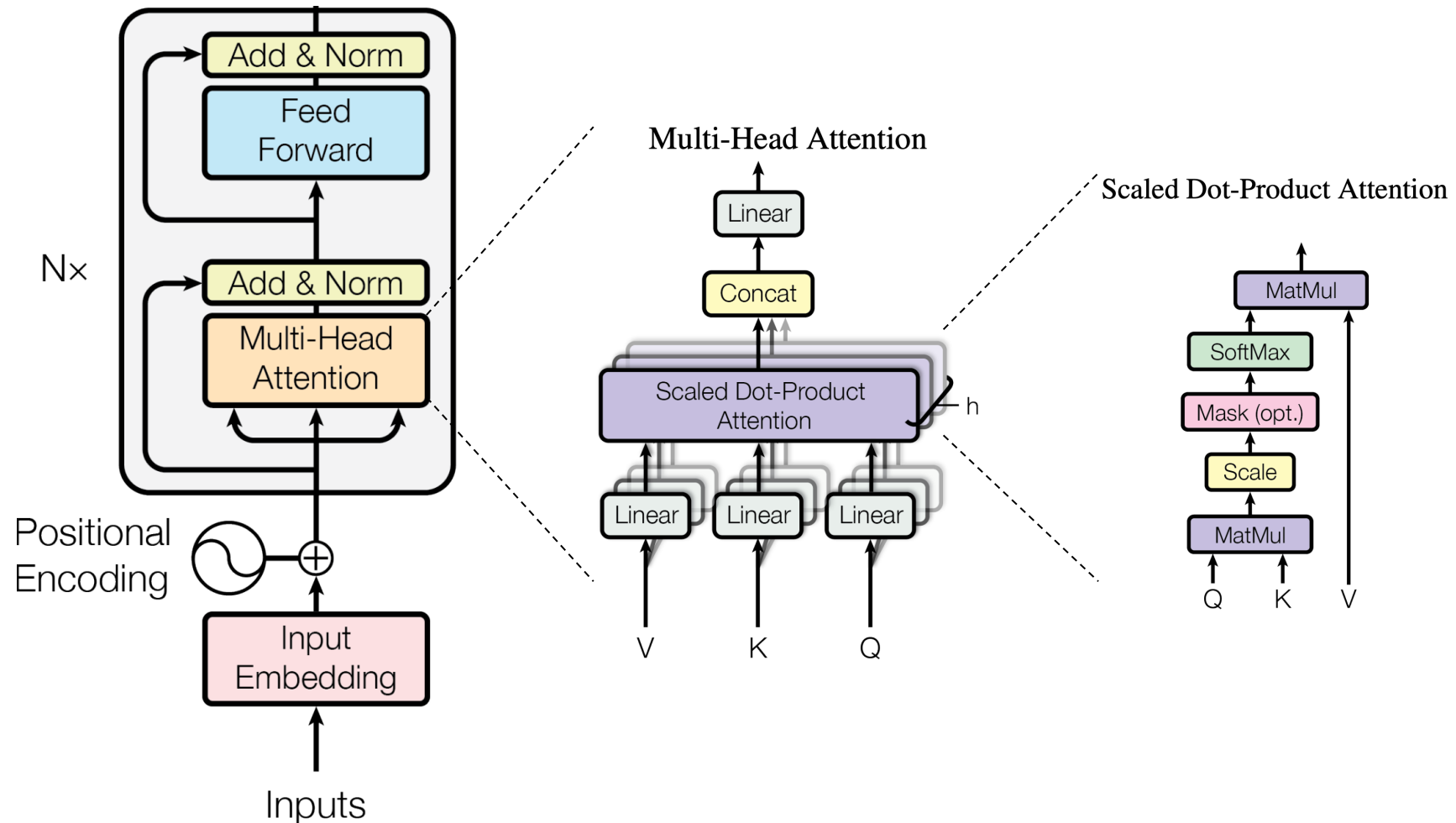
```
torch.Size([1, 50, 200])
```

```python
attn_output_weights.shape
```

```
torch.Size([1, 50, 50])
```

31

**AI VIET NAM**
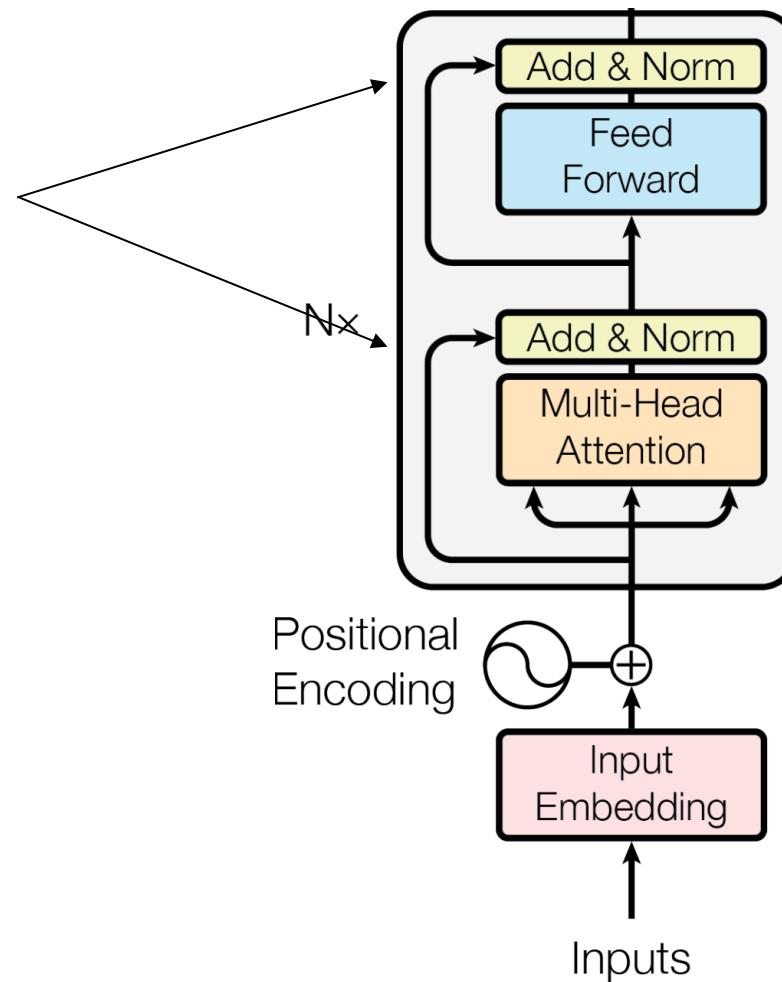@aivietnam.edu.vn

**Transformer-Encoder**

! 

## Layer Normalization

$$\mu_i = \frac{1}{m} \sum_{j=1}^{m} x_{ij}$$

$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^{m} (x_{ij} - \mu_i)^2$$

$$\hat{x}_{ij} = \frac{(x_{ij} - \mu_i)}{\sqrt{\sigma_i^2 + \epsilon}}$$

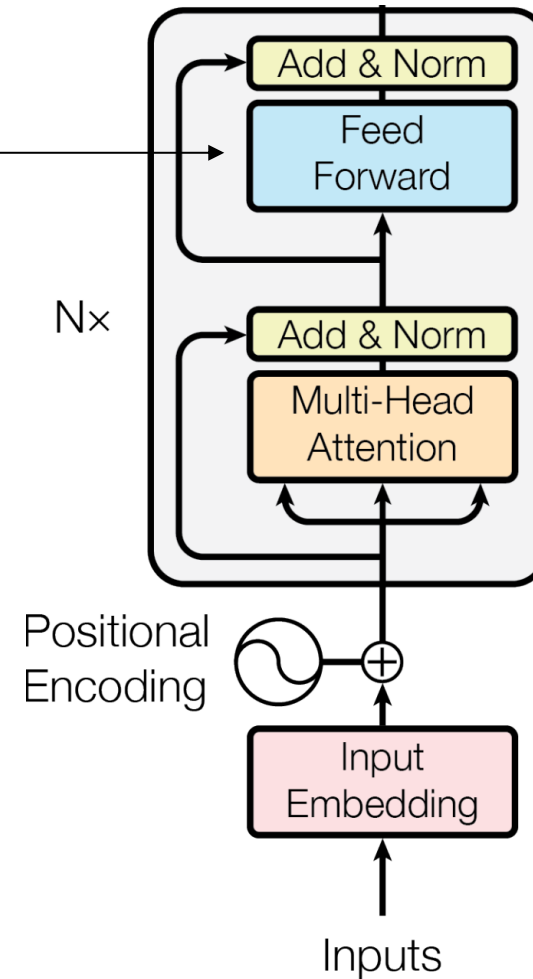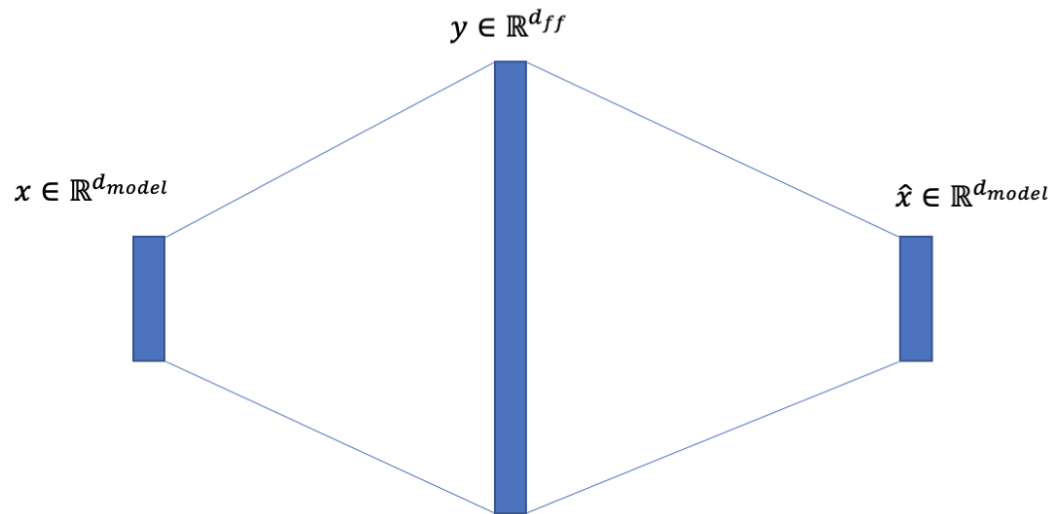> **Feed Forward**

❖ 2 FC Layer



$$y \in \mathbb{R}^{d_{ff}}$$

$$x \in \mathbb{R}^{d_{model}}$$

$$\hat{x} \in \mathbb{R}^{d_{model}}$$

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

34

! **Transformer-Encoder – Demo**

```python
class TransformerEncoder(nn.Module):
    def __init__(self, embed_dim, num_heads, ff_dim, dropout=0.1):
        super().__init__()
        self.attn = nn.MultiheadAttention(
            embed_dim=embed_dim,
            num_heads=num_heads,
            batch_first=True
        )
        self.ffn = nn.Sequential(
            nn.Linear(in_features=embed_dim, out_features=ff_dim, bias=True),
            nn.ReLU(),
            nn.Linear(in_features=ff_dim, out_features=embed_dim, bias=True)
        )
        self.layernorm_1 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
        self.layernorm_2 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
        self.dropout_1 = nn.Dropout(p=dropout)
        self.dropout_2 = nn.Dropout(p=dropout)

    def forward(self, query, key, value):
        attn_output, _ = self.attn(query, key, value)
        attn_output = self.dropout_1(attn_output)
        out_1 = self.layernorm_1(query + attn_output)
        ffn_output = self.ffn(out_1)
        ffn_output = self.dropout_2(ffn_output)
        out_2 = self.layernorm_2(out_1 + ffn_output)
        return out_2
```

```python
encoder_layer = TransformerEncoder(
    embed_dim=200,
    num_heads=5,
    ff_dim=1024
)
```

```python
embedded.shape
```

```
torch.Size([32, 50, 200])
```

```python
encoded = encoder_layer(embedded, embedded, embedded)
```

```python
encoded.shape
```

```
torch.Size([32, 50, 200])
```

**AI VIET NAM**
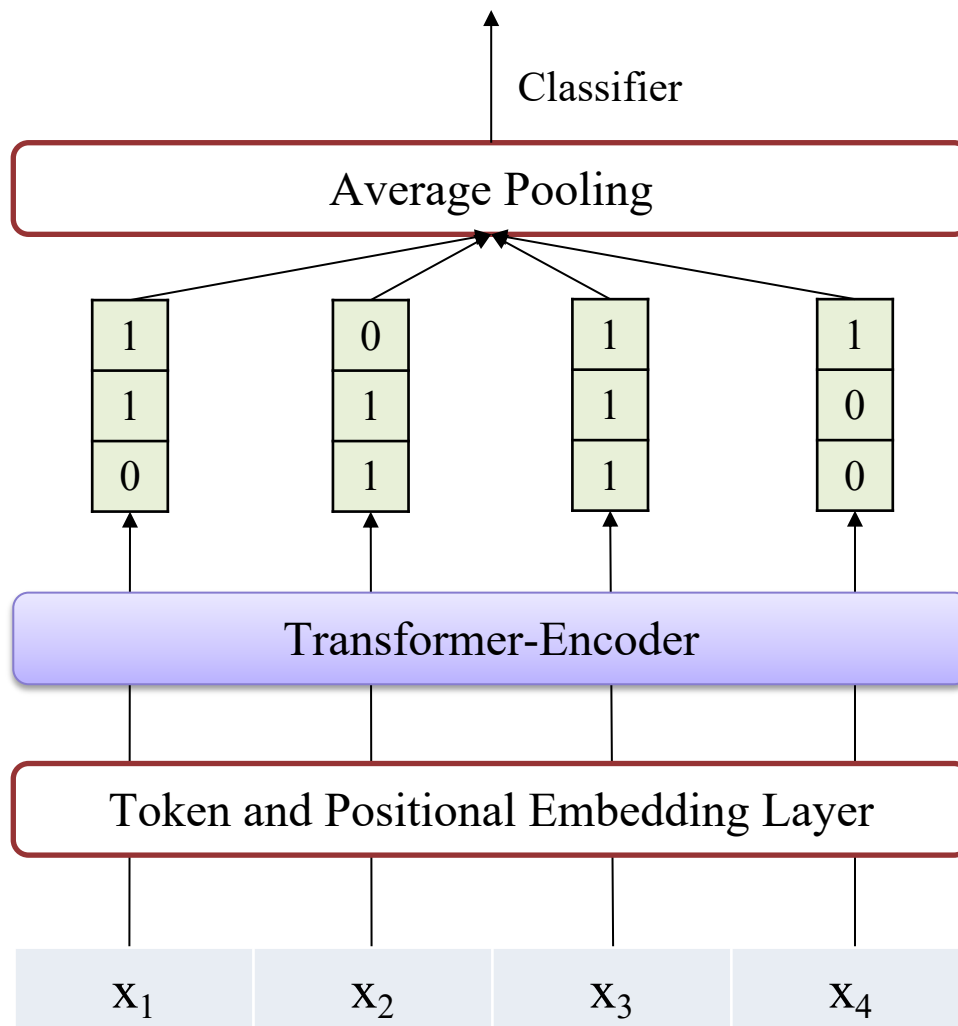@aivietnam.edu.vn

**!**

**NTC-SCV Dataset**

❖ **Sentiment Analysis**

| Positive Example | Negative Example |
|---|---|
| Mình được 1 cô bạn giới_thiệu đến đây , tìm địa_chỉ khá dễ . Menu nước uống chất khỏi nói . Mình muốn cũng đc 8 loại nước ở đây , món nào cũng ngon và bổ_dưỡng cả . | Quán chế_biến đồ_ăn lâu , Cá_Sapa nướng ướp rất dở , sò Lông ko tươi , nước_chấm ko ngon\n Tóm_lại sẽ ko bao_giờ ghé nữa , ăn_dở mà uổng tiền |
| Mỗi lần thèm trà sữa là làm 1 ly . Quán dễ kiếm , không_gian lại rộng_rãi . Nhân_viên thì dễ_thương gần_gũi . Nói_chung thèm trà sữa là mình ghé Quán ở đây vì gần nhà . | Quán này thấy khá nhiều người bảo mình nên mình đã đi ăn thử , nhưng thực_sự ăn xong thấy không được như mong_đợi lắm . |

**AI VIET NAM**
@aivietnam.edu.vn

! **Modeling**



Classifier

Average Pooling

| 1 | 0 | 1 | 1 |
| 1 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 |

Transformer-Encoder

Token and Positional Embedding Layer

| $x_1$ | $x_2$ | $x_3$ | $x_4$ |

37

**Modeling – Demo**

```python
class TransformerEncoder(nn.Module):
    def __init__(self, embed_dim, num_heads, ff_dim, dropout=0.1):
        super().__init__()
        self.attn = nn.MultiheadAttention(
            embed_dim=embed_dim,
            num_heads=num_heads,
            batch_first=True
        )
        self.ffn = nn.Sequential(
            nn.Linear(in_features=embed_dim, out_features=ff_dim, bias=True),
            nn.ReLU(),
            nn.Linear(in_features=ff_dim, out_features=embed_dim, bias=True)
        )
        self.layernorm_1 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
        self.layernorm_2 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
        self.dropout_1 = nn.Dropout(p=dropout)
        self.dropout_2 = nn.Dropout(p=dropout)

    def forward(self, query, key, value):
        attn_output, _ = self.attn(query, key, value)
        attn_output = self.dropout_1(attn_output)
        out_1 = self.layernorm_1(query + attn_output)
        ffn_output = self.ffn(out_1)
        ffn_output = self.dropout_2(ffn_output)
        out_2 = self.layernorm_2(out_1 + ffn_output)
        return out_2
```
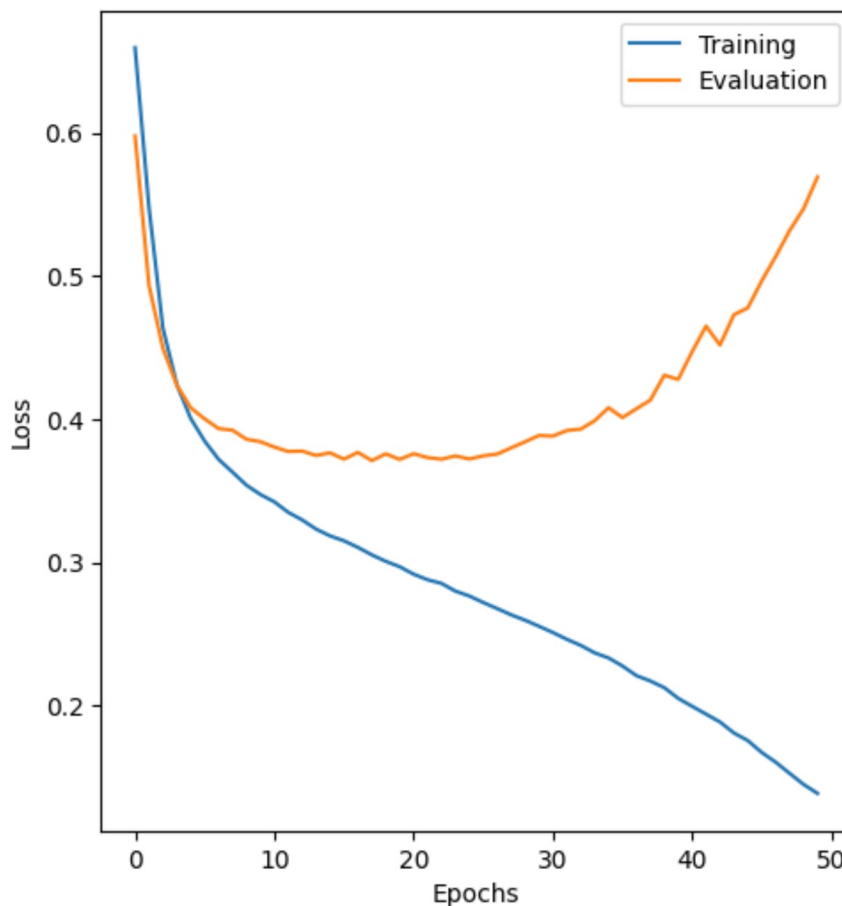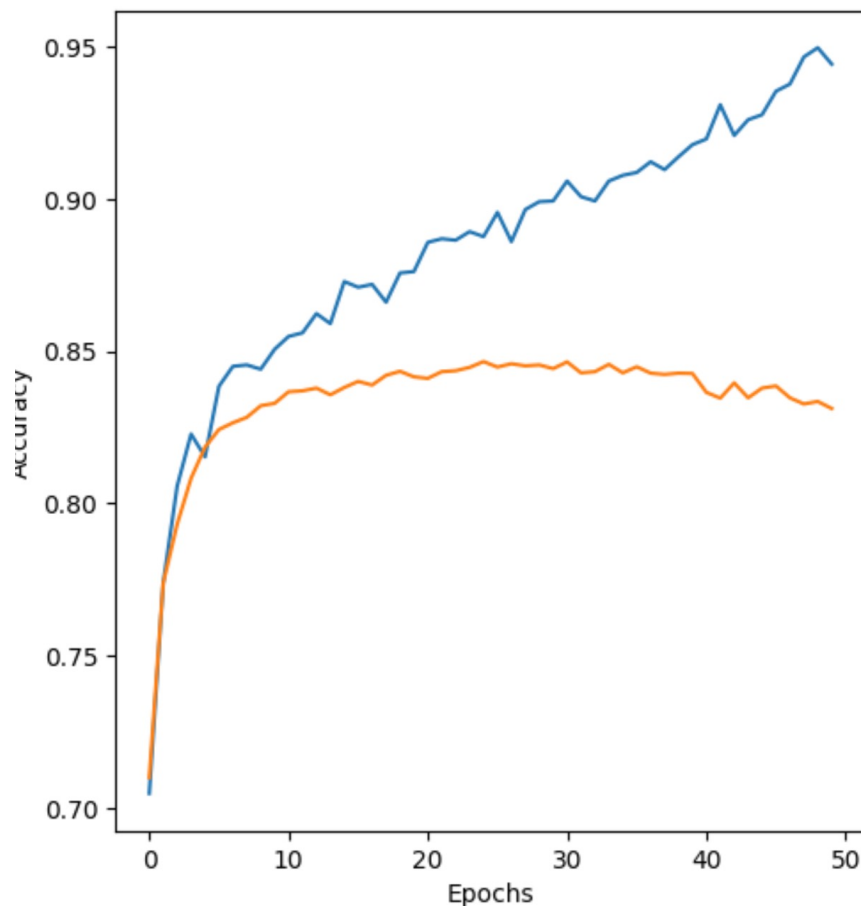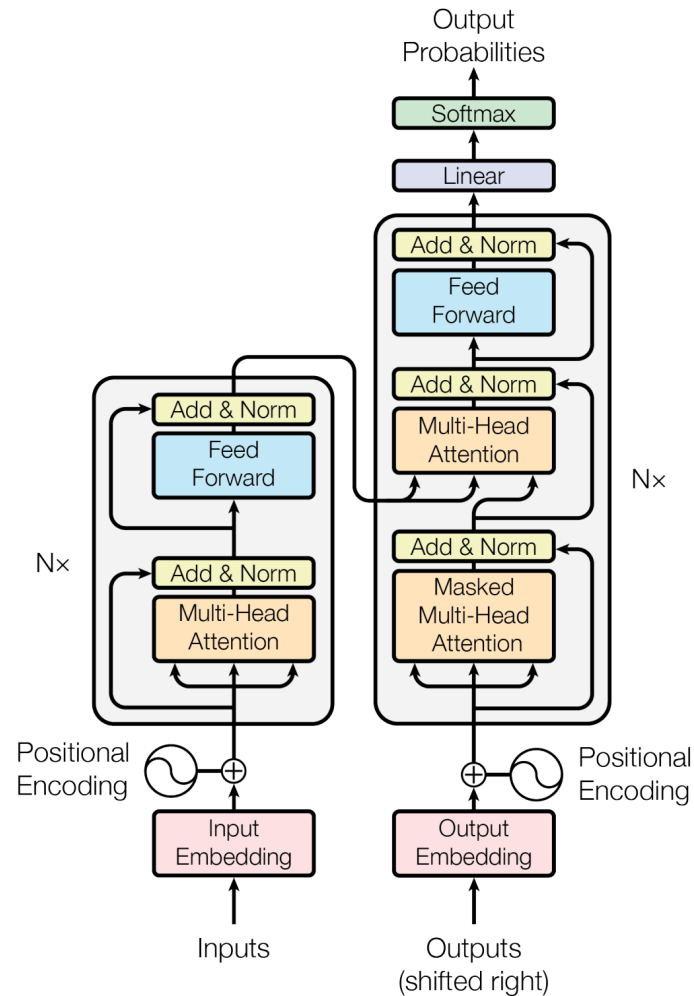
38

**AI VIET NAM**
@aivietnam.edu.vn

**!**

**Training**

❖ Testing: 83.66%

**AI VIET NAM**
@aivietnam.edu.vn

## ViT



Transformers are so successful in NLP,
Can we use them for images?

40

**AI VIET NAM**
@aivietnam.edu.vn

! **From text to image**



"I Love AI VN"

Average Pooling

| 1 | 0 | 1 | 1 |
| 1 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 |

Transformer-Encoder

Token and Positional Embedding Layer

Tokenization

| I | Love | AI | VN |

**!** From text to image

Can we tokenize an image?

**From text to image**

Can we tokenize an image?

**From text to image**

Can we tokenize an image?



$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6 \qquad x_{n-3} \quad x_{n-2} \quad x_{n-1} \quad x_n$

Flattening

AI VIET NAM
@aivietnam.edu.vn

## ViT Architecture



Vision Transformer (ViT)

Class
Bird
Ball
Car
...

MLP
Head

Transformer Encoder

Patch + Position Embedding

0 * 1 2 3 4 5 6 7 8 9

* Extra learnable [class] embedding

Linear Projection of Flattened Patches

Transformer Encoder

L ×

MLP

Norm

Multi-Head Attention

Norm

Embedded Patches

45

**AI VIET NAM**
@aivietnam.edu.vn

!

**Patch embedding**



**Vision Transformer (ViT)**

**Transformer Encoder**

46

**AI VIET NAM**
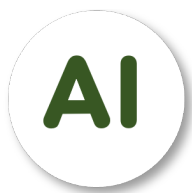@aivietnam.edu.vn

!

**Patch embedding**

```python
1 class PatchEmbedding(nn.Module):
2     def __init__(self, embed_dim=512, patch_size=16, image_size=224):
3         super().__init__()
4         self.conv1 = nn.Conv2d(in_channels=3, out_channels=embed_dim, kernel_size=patch_size, stride=patch_size, bias=False)
5
6     def forward(self, x):
7         x = self.conv1(x) # shape = [*, width, grid, grid]
8         x = x.reshape(x.shape[0], x.shape[1], -1)  # shape = [*, width, grid ** 2]
9         x = x.permute(0, 2, 1)  # shape = [*, grid ** 2, width]
10        return x
```

```python
1 patch_embedding = PatchEmbedding()
2 x = torch.randn(1, 3, 224, 224)
3
4 out = patch_embedding(x)
5 print(out.shape)
```

```
torch.Size([1, 196, 512])
```

47

**AI VIET NAM**
@aivietnam.edu.vn

! **Patch embedding**

Patch size: 4

**!** **Positional embedding**



Embed + add

49

**AI VIET NAM**
@aivietnam.edu.vn

!

## Positional embedding

```python
1 class PatchPositionEmbedding(nn.Module):
2     def __init__(self, embed_dim=512, patch_size=16, image_size=224):
3         super().__init__()
4         self.conv1 = nn.Conv2d(in_channels=3, out_channels=embed_dim, kernel_size=patch_size, stride=patch_size, bias=False)
5
6         scale = embed_dim ** -0.5
7         self.positional_embedding = nn.Parameter(scale * torch.randn((image_size // patch_size) ** 2, embed_dim))
8
9     def forward(self, x):
10        x = self.conv1(x) # shape = [*, width, grid, grid]
11        x = x.reshape(x.shape[0], x.shape[1], -1)  # shape = [*, width, grid ** 2]
12        x = x.permute(0, 2, 1)  # shape = [*, grid ** 2, width]
13
14        x = x + self.positional_embedding.to(x.dtype)
15        return x
```

```python
[36]  1 patchpos_embedding = PatchPositionEmbedding()
      2 x = torch.randn(1, 3, 224, 224)
      3
      4 out = patchpos_embedding(x)
      5 print(out.shape)

     torch.Size([1, 196, 512])
```

50

**AI VIET NAM**
@aivietnam.edu.vn

! **Positional embedding**

| | |
|---|---|
| 14 | 10 |

| | |
|---|---|
| 14 | 12 |

+

| | |
|---|---|
| 18 | 17 |

| | |
|---|---|
| 11 | 9 |

| | |
|---|---|
| 0 | 0 |

| | |
|---|---|
| 1 | 1 |

| | |
|---|---|
| 1 | 1 |

| | |
|---|---|
| 0 | 0 |

=

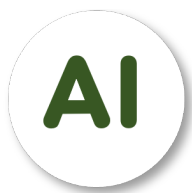| | |
|---|---|
| 14 | 10 |

| | |
|---|---|
| 15 | 13 |

| | |
|---|---|
| 19 | 18 |

| | |
|---|---|
| 11 | 9 |

Positional
Embedding

**AI VIET NAM**
@aivietnam.edu.vn

**!**

**[CLS] Token**

AI VIET NAM
@aivietnam.edu.vn

**[CLS] Token – Why?**

Alternatives?
- Global Average Pooling
- Max Pooling
- …

53

AI VIET NAM
@aivietnam.edu.vn

! 

**[CLS] Token – Demo**

```python
1  class PatchPositionEmbedding(nn.Module):
2      def __init__(self, embed_dim=512, patch_size=16, image_size=224):
3          super().__init__()
4          self.conv1 = nn.Conv2d(in_channels=3, out_channels=embed_dim, kernel_size=patch_size, stride=patch_size, bias=False)
5
6          scale = embed_dim ** -0.5
7          self.class_embedding = nn.Parameter(scale * torch.randn(embed_dim))
8          self.positional_embedding = nn.Parameter(scale * torch.randn((image_size // patch_size) ** 2 + 1, embed_dim))
9
10     def forward(self, x):
11         x = self.conv1(x) # shape = [*, width, grid, grid]
12         x = x.reshape(x.shape[0], x.shape[1], -1)  # shape = [*, width, grid ** 2]
13         x = x.permute(0, 2, 1)  # shape = [*, grid ** 2, width]
14         # expanding the CLS embedding
15         cls_embs = self.class_embedding.to(x.dtype) + torch.zeros(x.shape[0], 1, x.shape[-1], dtype=x.dtype, device=x.device)
16         x = torch.cat([cls_embs, x], dim=1)  # shape = [*, grid ** 2 + 1, width]
17
18         x = x + self.positional_embedding.to(x.dtype)
19         return x
```

54

**Modeling**

Change Token
Embedding with
Patch Embedding
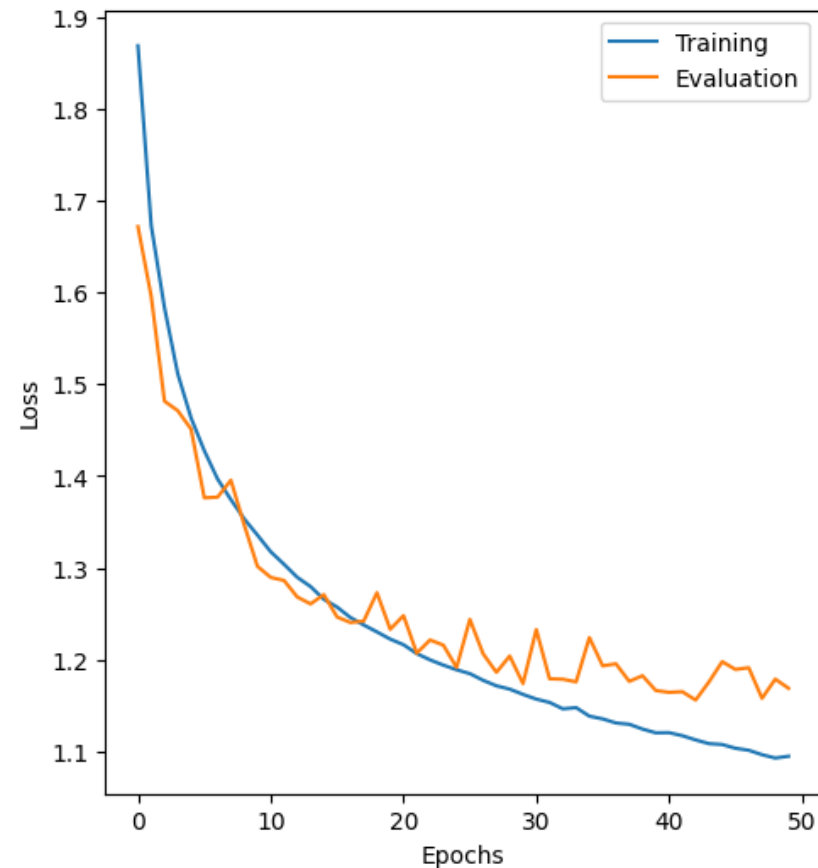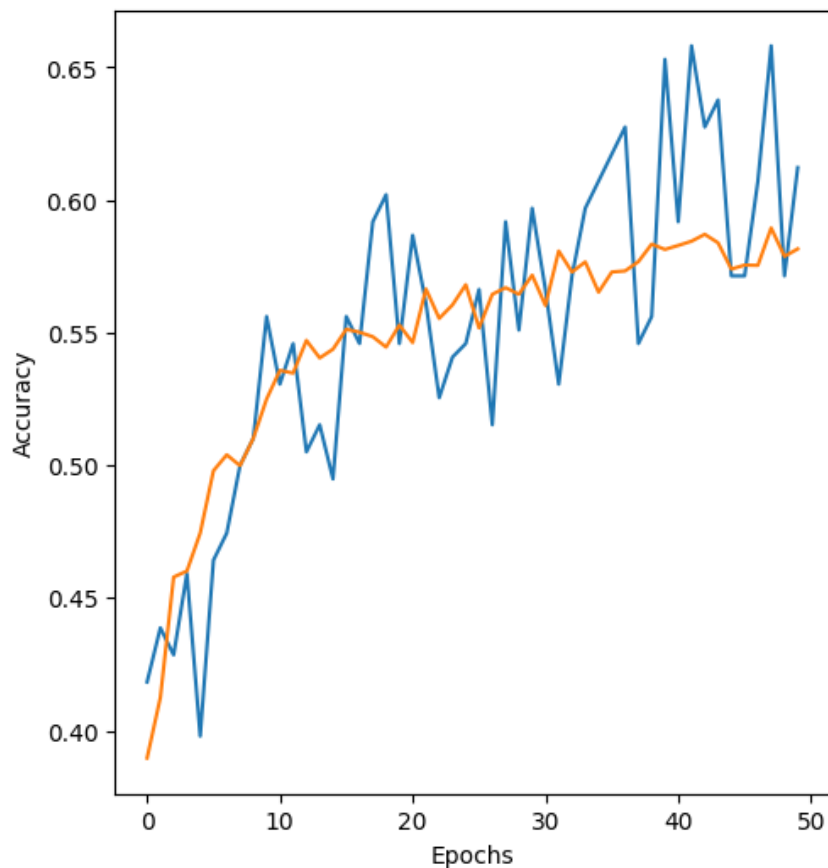
[CLS] token instead
of pooling (can still
use pooling)

```python
1  class VisionTransformerCls(nn.Module):
2      def __init__(self,
3                   image_size, embed_dim, num_heads, ff_dim,
4                   dropout=0.1, device='cpu', num_classes = 10, patch_size=16
5      ):
6          super().__init__()
7          self.embd_layer = PatchPositionEmbedding(
8              image_size=image_size, embed_dim=embed_dim, patch_size=patch_size
9          )
10         self.transformer_layer = TransformerEncoder(
11             embed_dim, num_heads, ff_dim, dropout
12         )
13         # self.pooling = nn.AvgPool1d(kernel_size=max_length)
14         self.fc1 = nn.Linear(in_features=embed_dim, out_features=20)
15         self.fc2 = nn.Linear(in_features=20, out_features=num_classes)
16         self.dropout = nn.Dropout(p=dropout)
17         self.relu = nn.ReLU()
18     def forward(self, x):
19         output = self.embd_layer(x)
20         output = self.transformer_layer(output, output, output)
21         output = output[:, 0, :]
22         output = self.dropout(output)
23         output = self.fc1(output)
24         output = self.dropout(output)
25         output = self.fc2(output)
26         return output
```

**!**

## Training

# Thanks!

## Any questions?