

Exercise Class

Transformer Applications

Nguyen Quoc Thai

CONTENT

(1) – Sequence-to-Sequence

(2) – Transformer

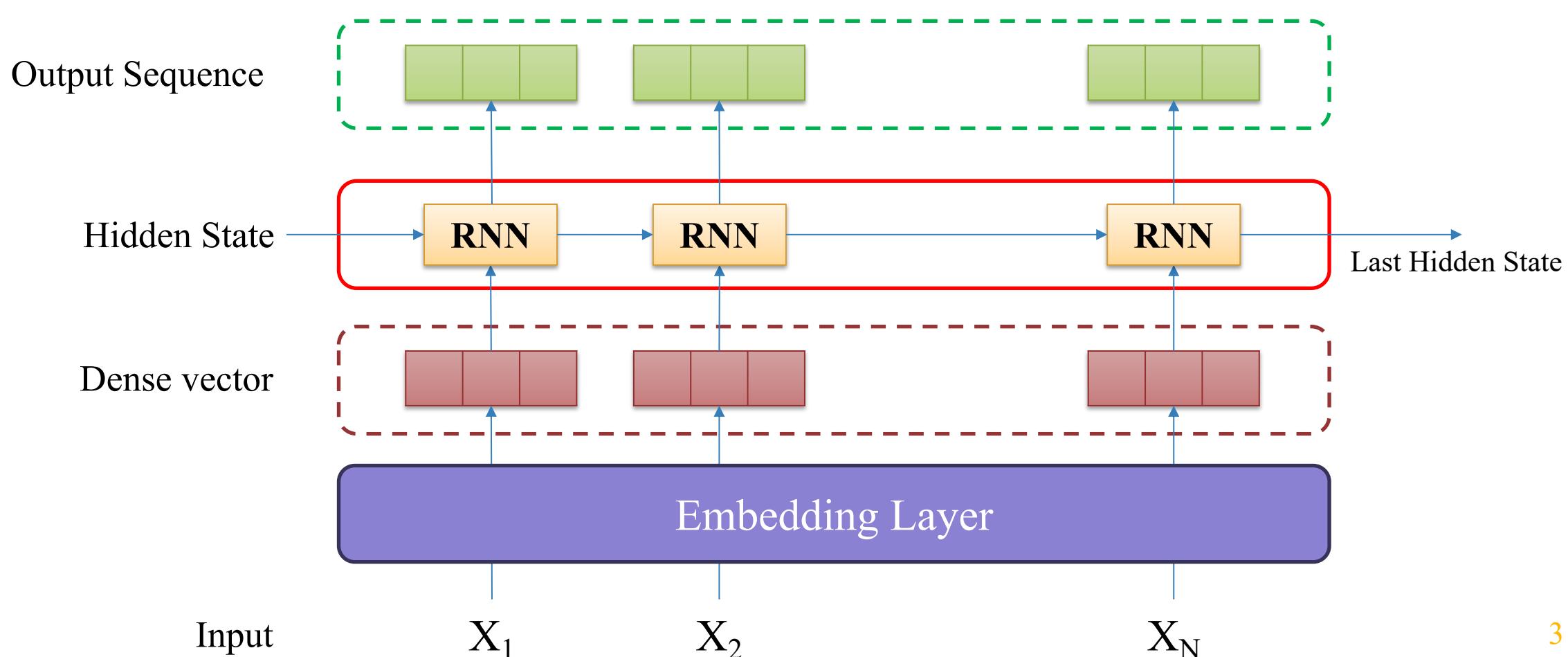
(3) – BERT

(4) – Vision Transformer

1 – Sequence-to-Sequence



RNNs Model

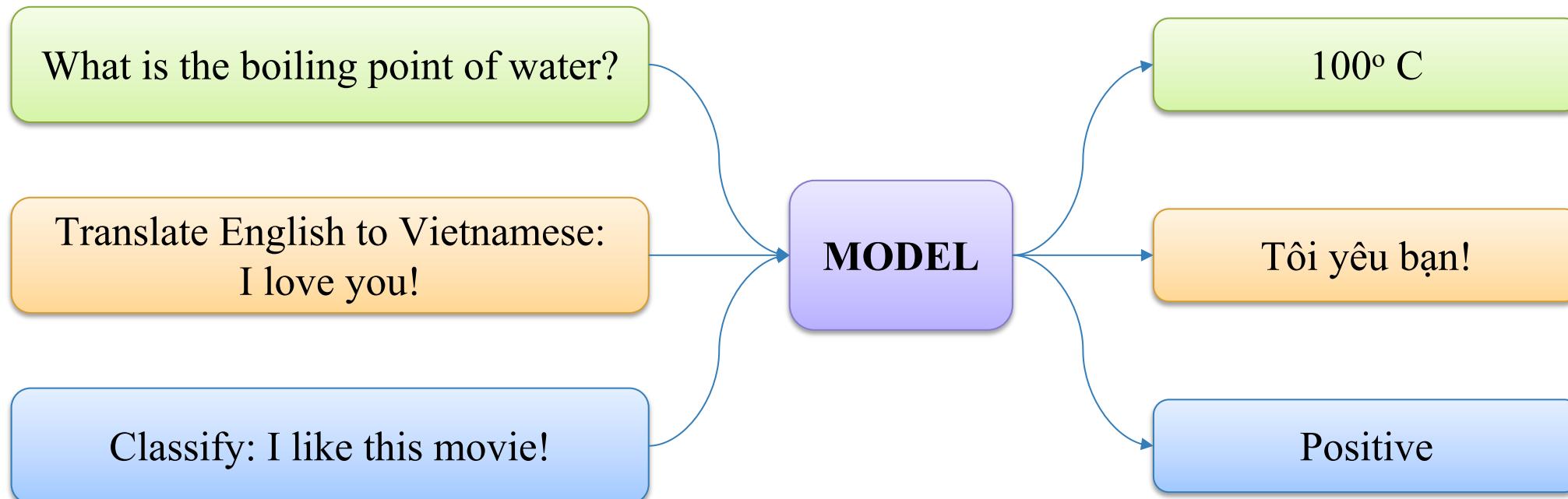


1 – Sequence-to-Sequence

!

Text Generation

- ❖ Take a sequence as input, predict a sequence as output

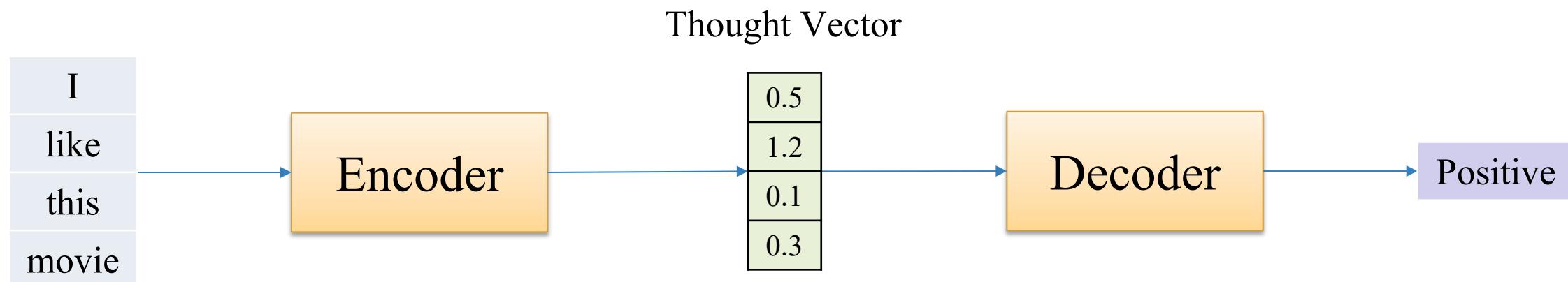


1 – Sequence-to-Sequence

!

Sequence-to-Sequence Architecture

- ❖ Take a sequence as input, predict a sequence as output
- ❖ Encoder: encoding the inputs into state (thought vector)
- ❖ Decoder: the state is passed into the decoder to generate the output



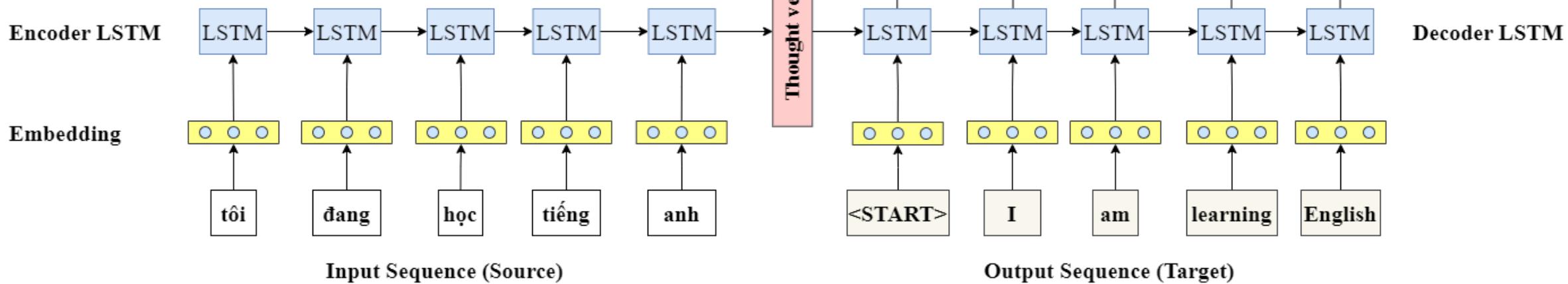
1 – Sequence-to-Sequence

!

Sequence-to-Sequence Architecture

❖ Training - Teacher Forcing

$$J = \frac{1}{T} \sum_{i=1}^T J_i$$

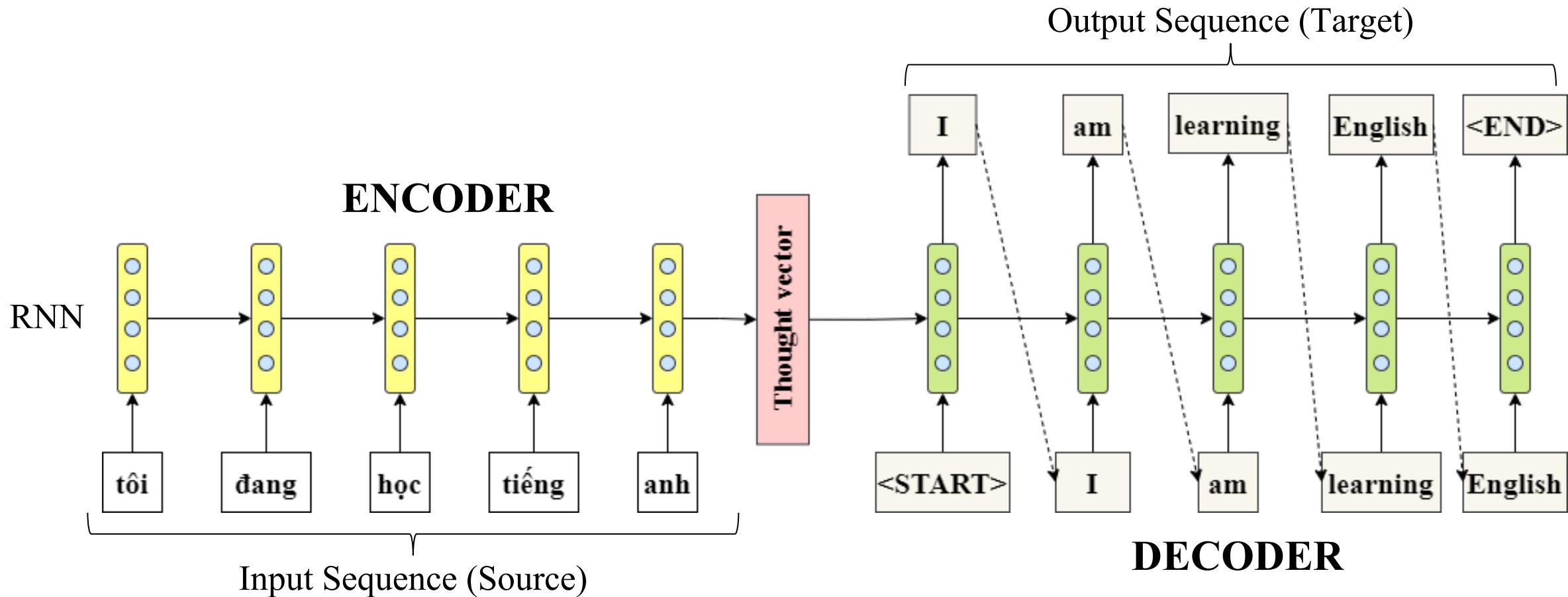


1 – Sequence-to-Sequence

!

Sequence-to-Sequence Architecture

- ❖ Inference – Decoding – Greedy Search

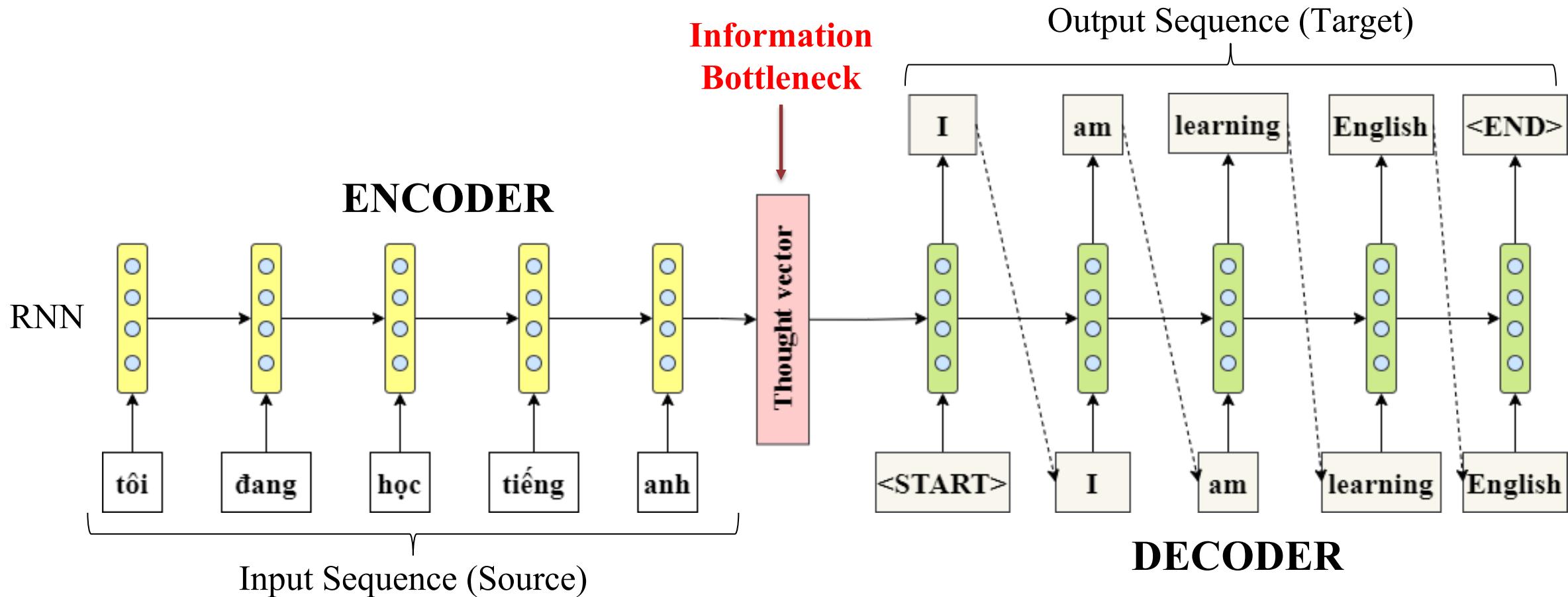


1 – Sequence-to-Sequence

!

Information Bottleneck

- Thought vector: capture all information of input sentence

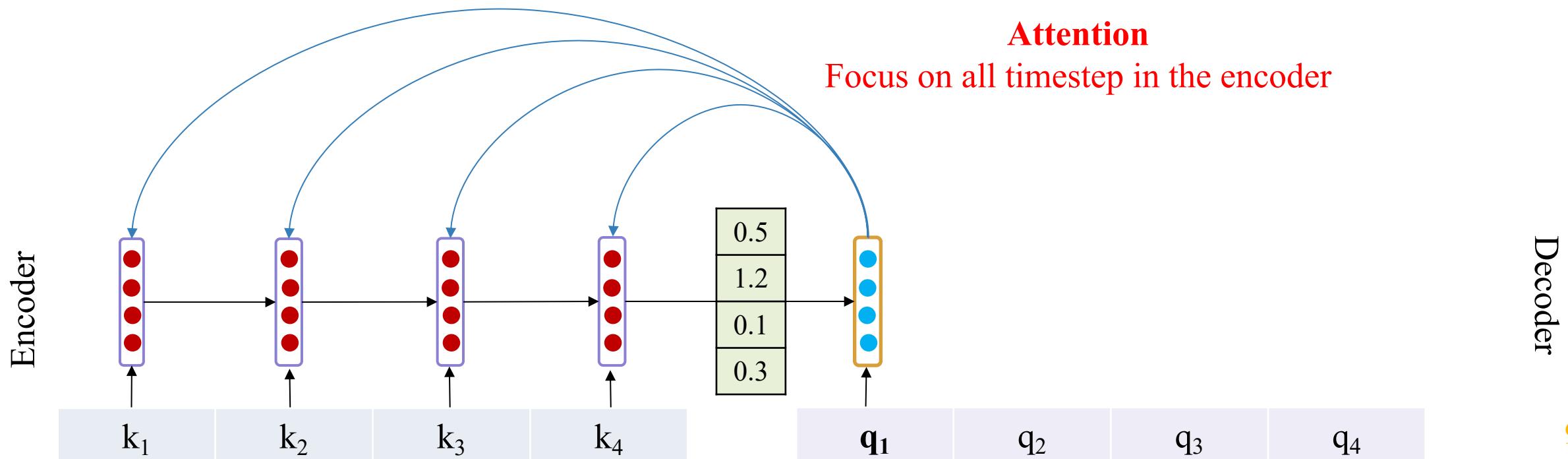


1 – Sequence-to-Sequence



Information Bottleneck

- Thought vector: capture all information of input sentence

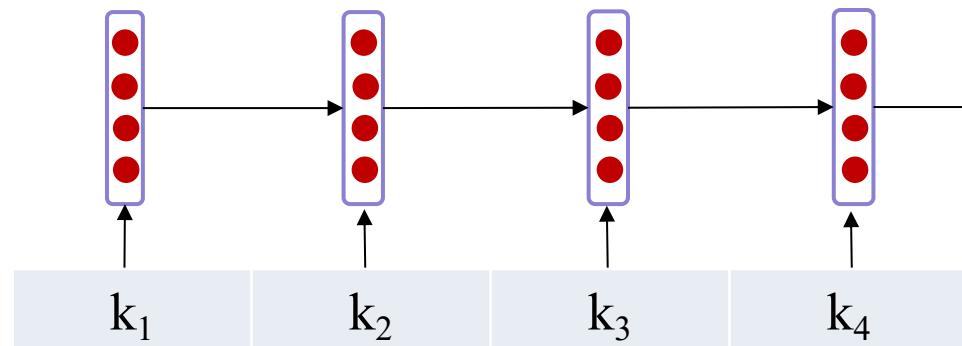


1 – Sequence-to-Sequence

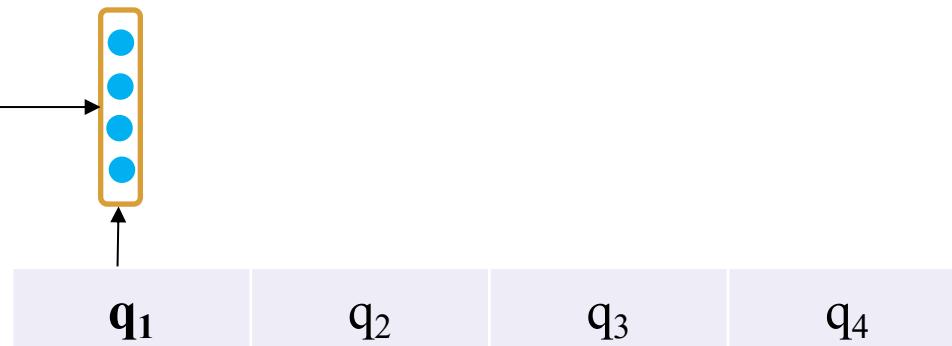
!

Scaled Dot-Product Attention

Encoder



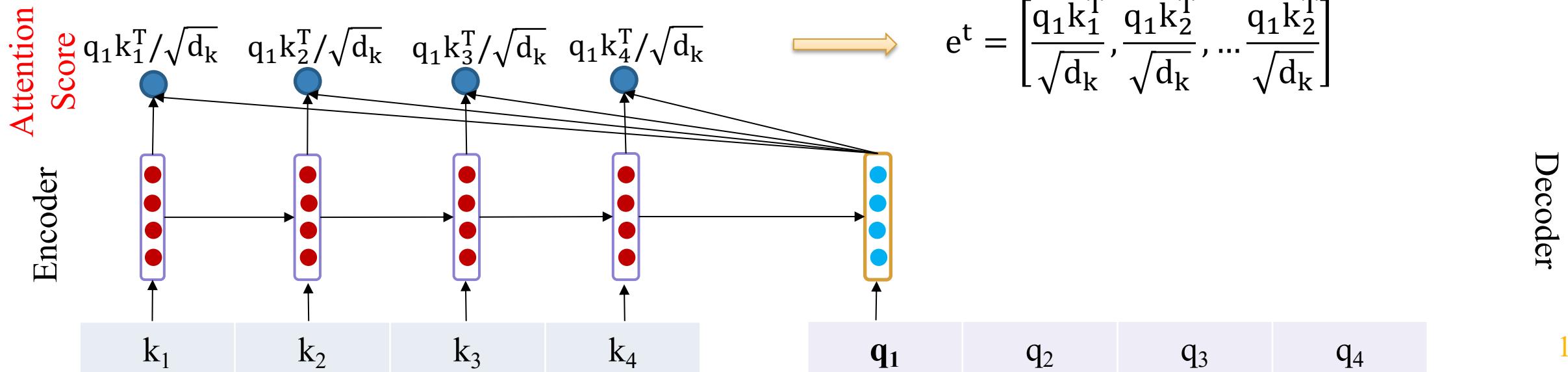
Decoder



1 – Sequence-to-Sequence

!

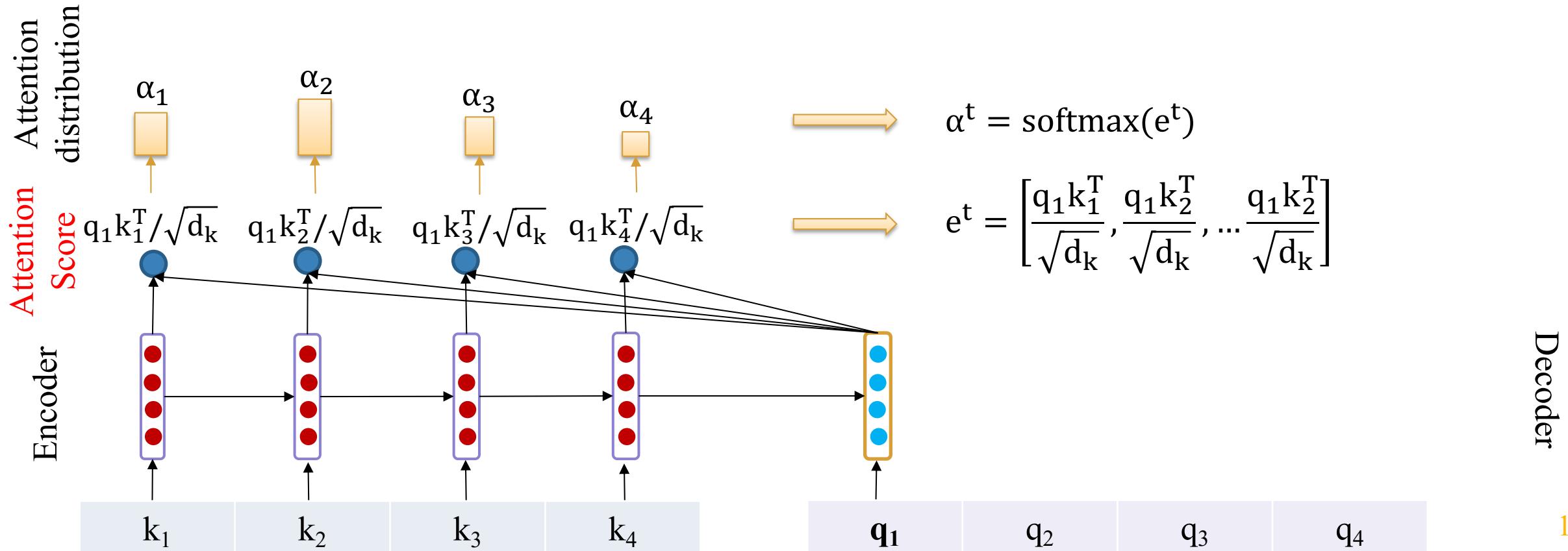
Scaled Dot-Product Attention



1 – Sequence-to-Sequence

!

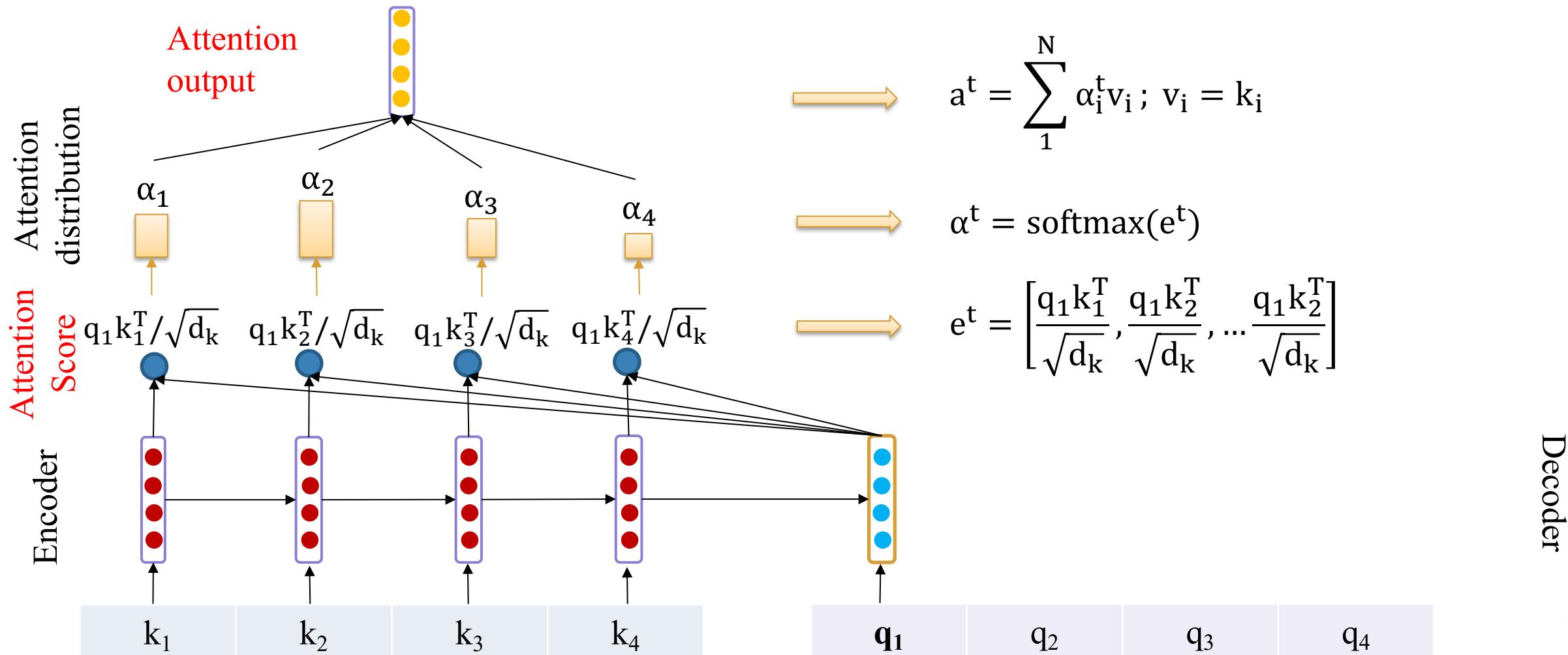
Scaled Dot-Product Attention



1 – Sequence-to-Sequence

!

Scaled Dot-Product Attention



1 – Sequence-to-Sequence

!

Scaled Dot-Product Attention

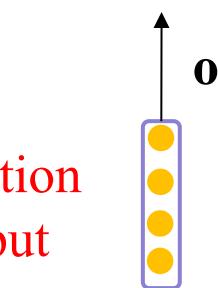
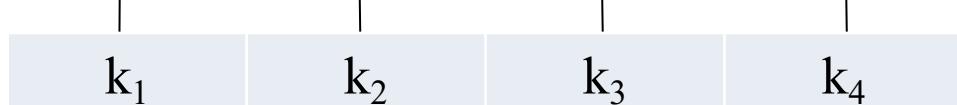
$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V; \quad K = V$$

Loss - Prediction

Attention
output

Scaled Dot-Product Attention

Encoder



Decoder

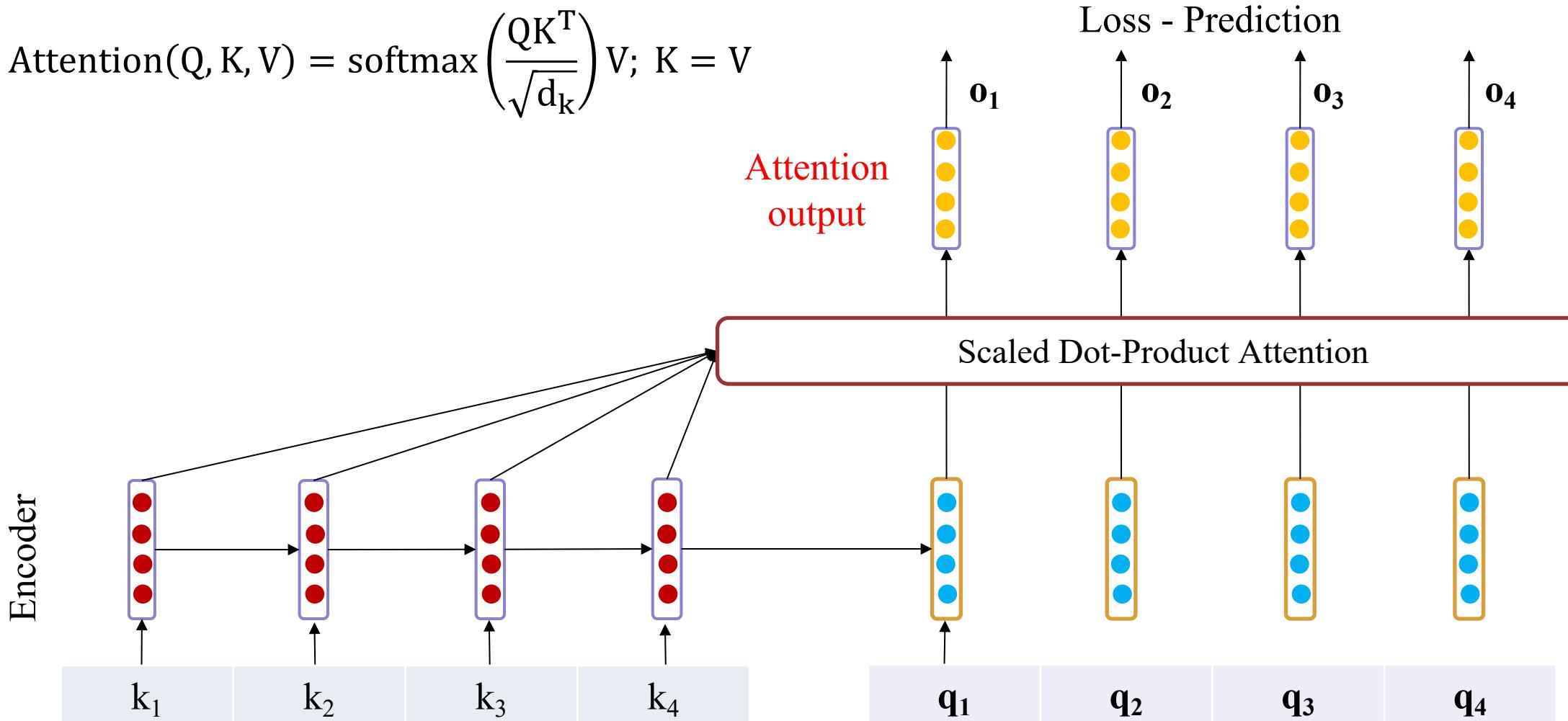


1 – Sequence-to-Sequence



Scaled Dot-Product Attention

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



1 – Sequence-to-Sequence



Attention Variants

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

Multiplicative Attention

$$\text{softmax}(QWK^T)V$$

Additive Attention

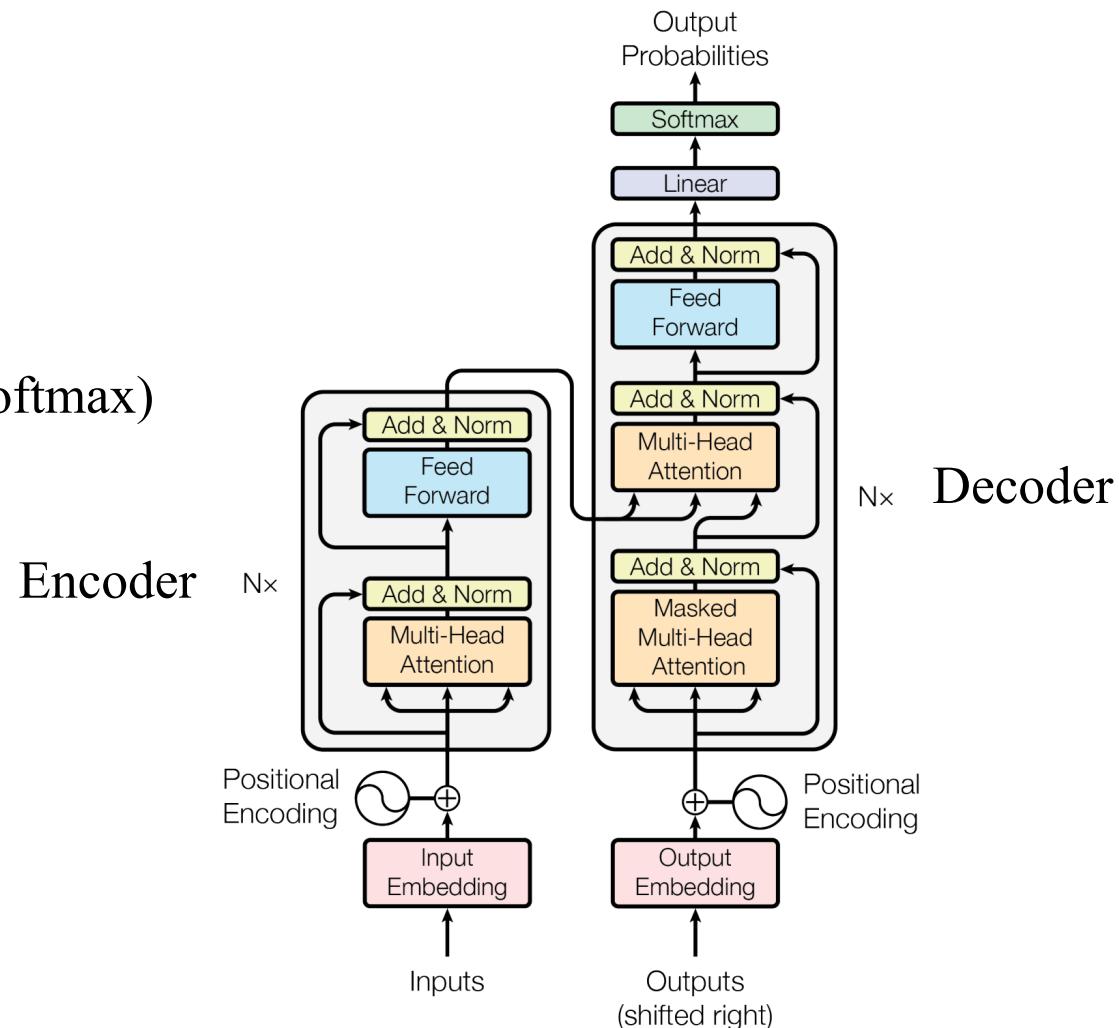
$$\tanh(W_1Q + W_2K)V$$

2 – Transformer



Transformer

- ❖ Architecture:
 - Input Embedding
 - Positional Encoding
 - N Encoder Layer
 - N Decoder Layer
 - Language Model Head (Projection + Softmax)
- ❖ Core technique: attention
- ❖ Loss function: cross-entropy

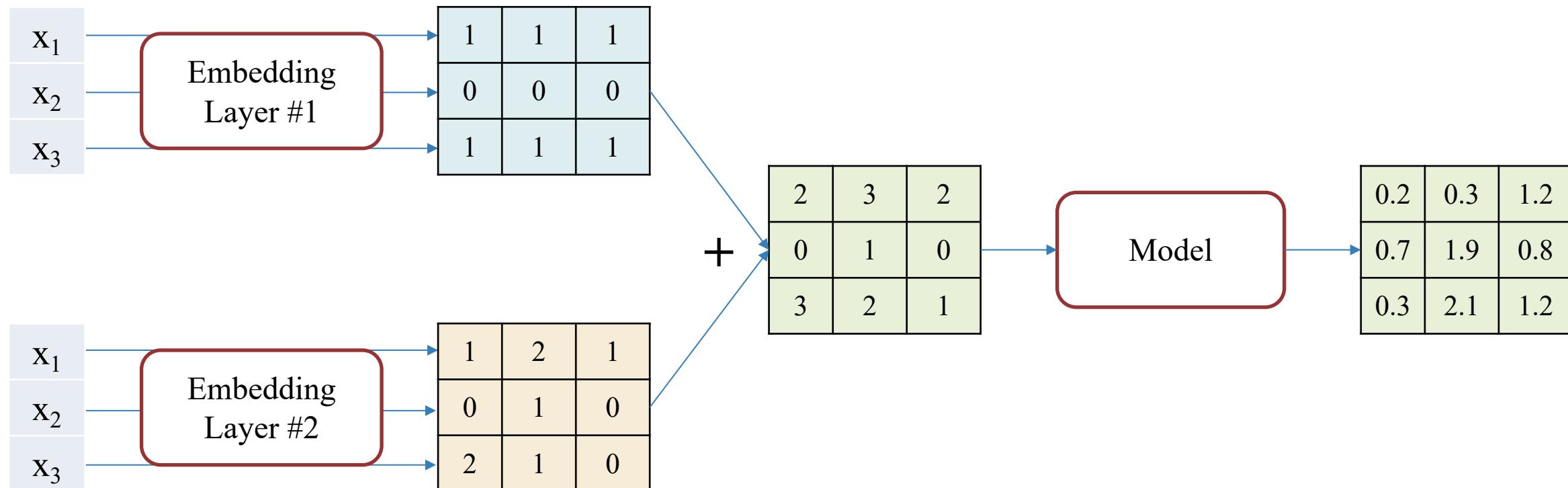


2 – Transformer

!

Input Embedding - 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)

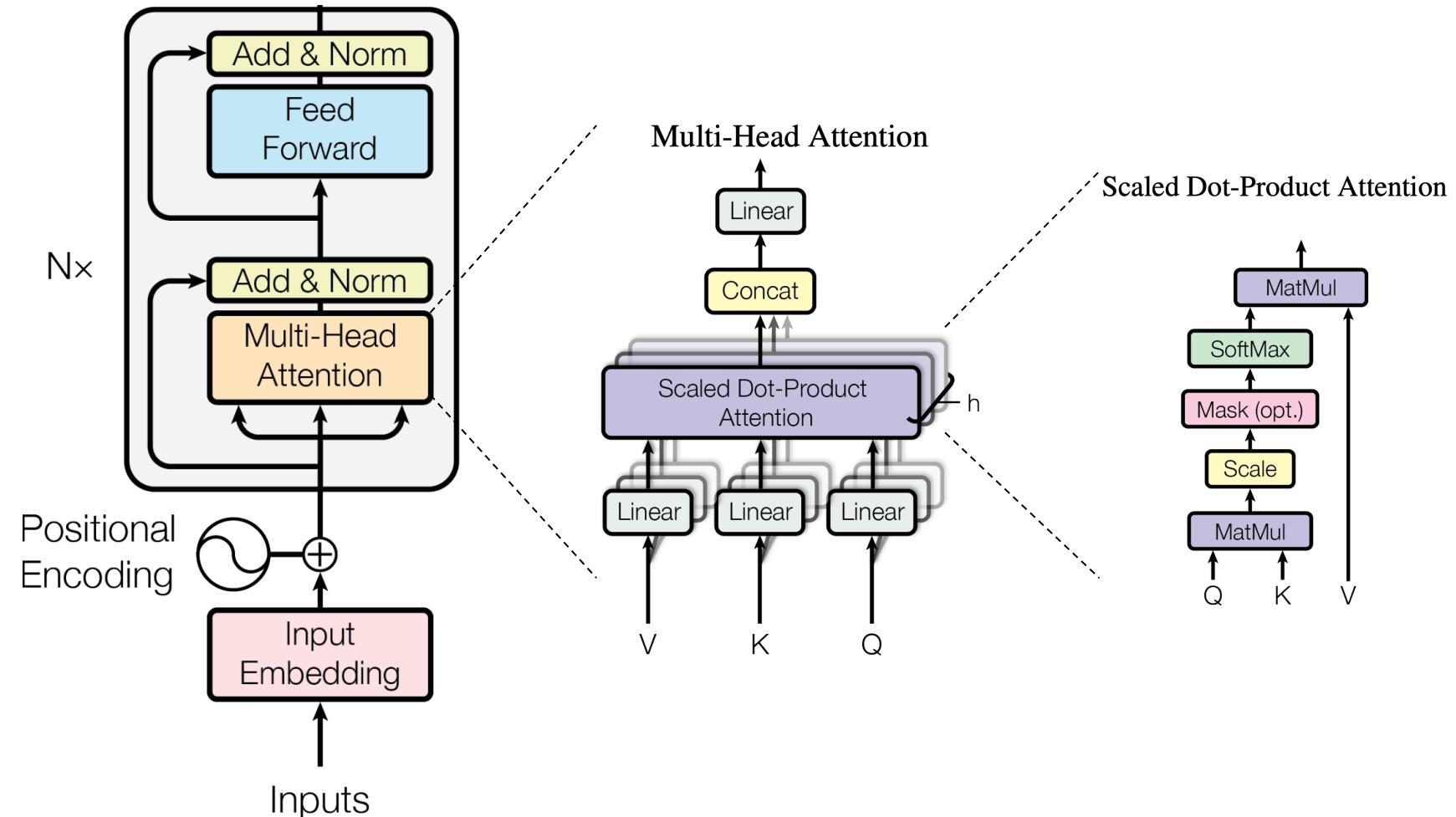


2 – Transformer



Transformer-Encoder

- ❖ Multi-Head Attention
- ❖ Feed Forward
- ❖ Add & Norm

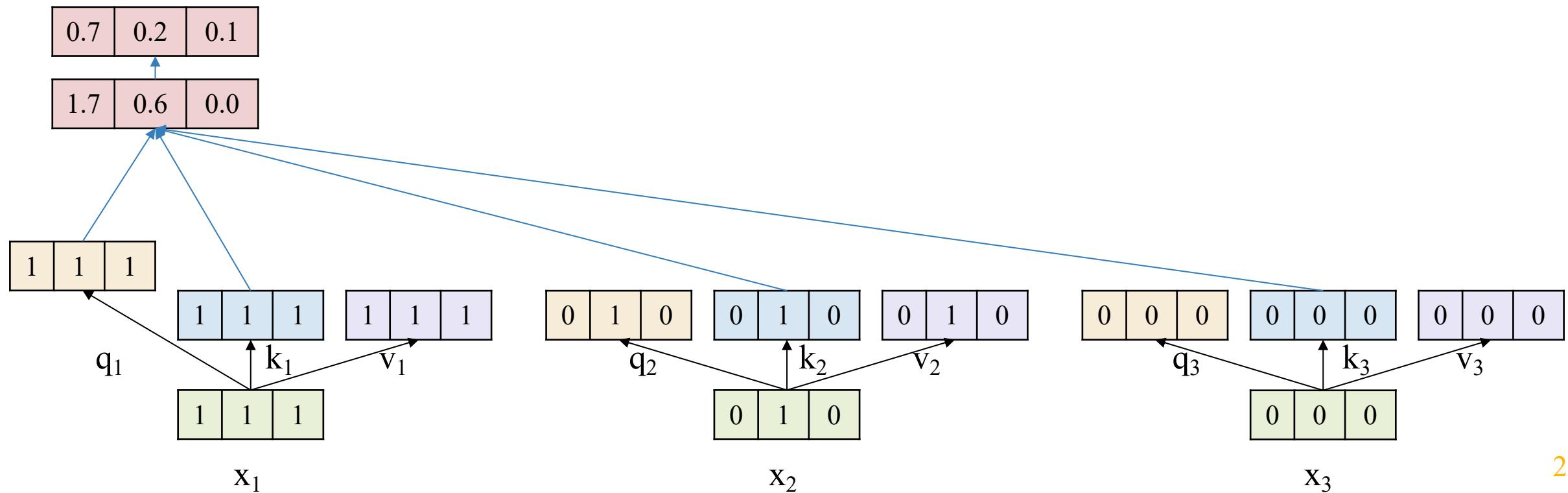


2 – Transformer

!

Transformer-Encoder

❖ Self-Attention

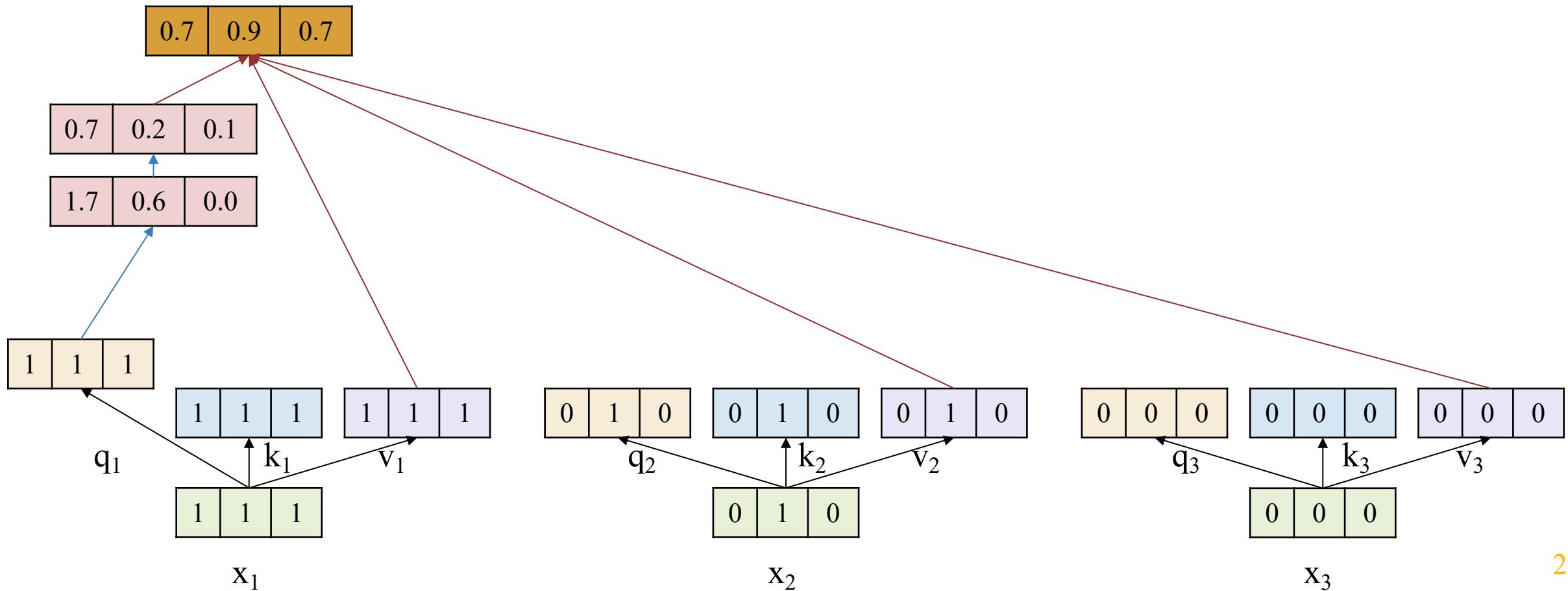


2 – Transformer



Transformer-Encoder

❖ Self-Attention

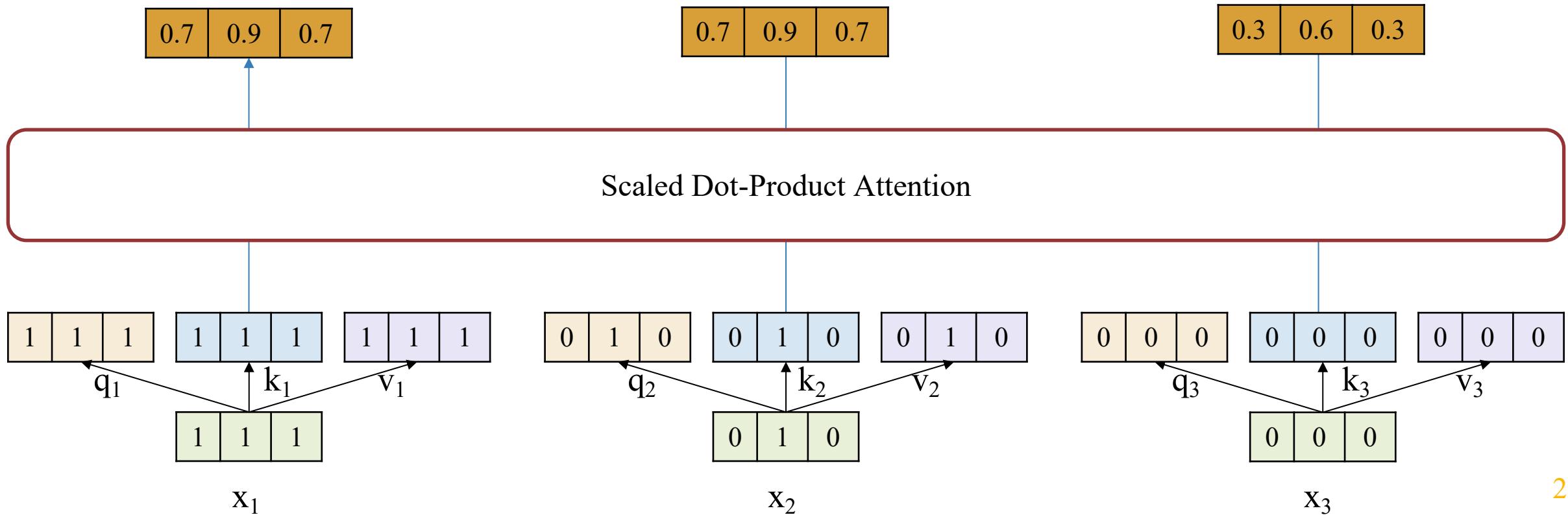


2 – Transformer



Transformer-Encoder

❖ Self-Attention

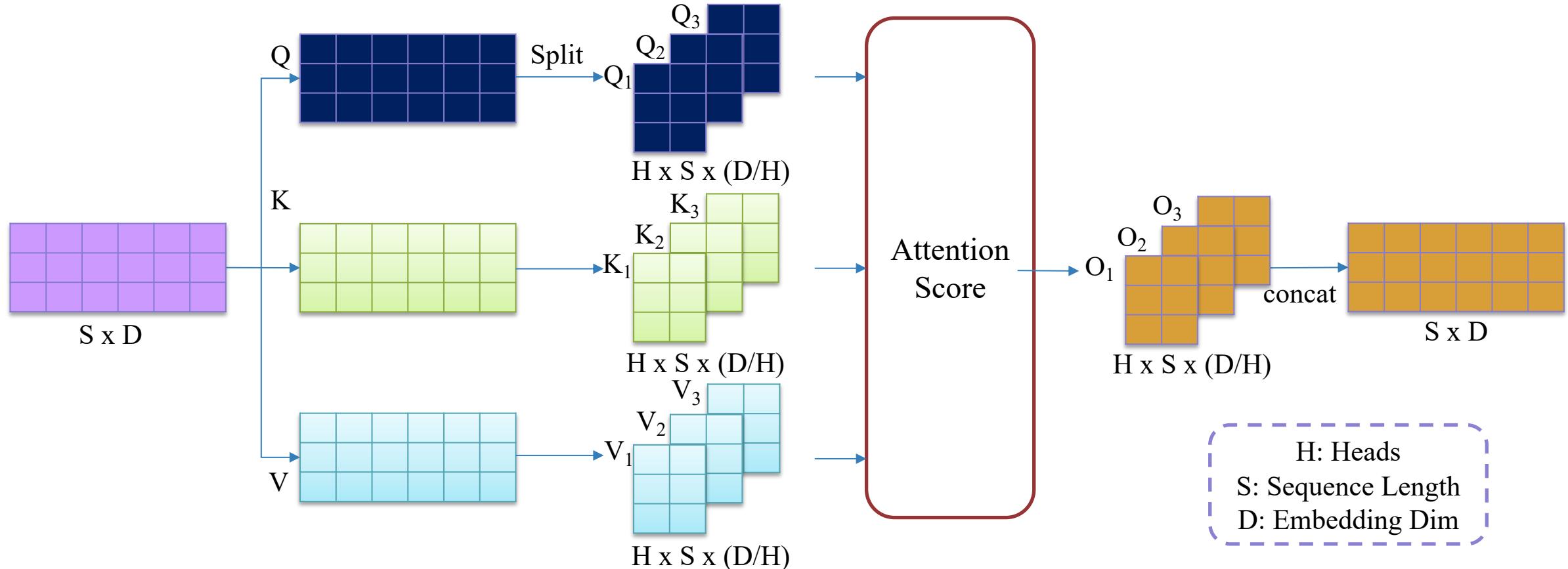


2 – Transformer



Transformer-Encoder

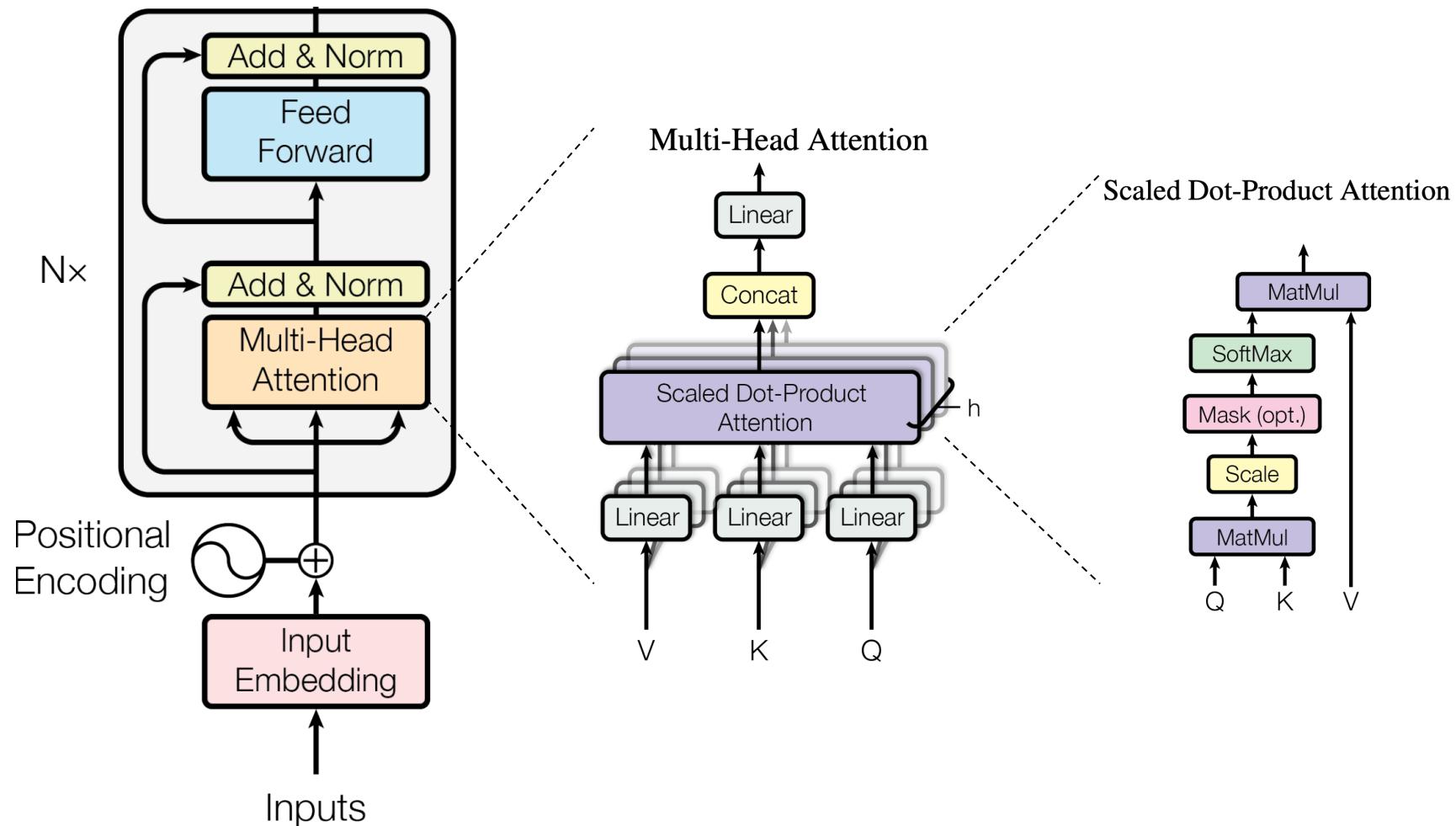
- ❖ **Multi-Head Attention**
- ❖ Split into the multiple attention heads (process independently) => self-attention => concat



2 – Transformer



Transformer-Encoder



2 – Transformer



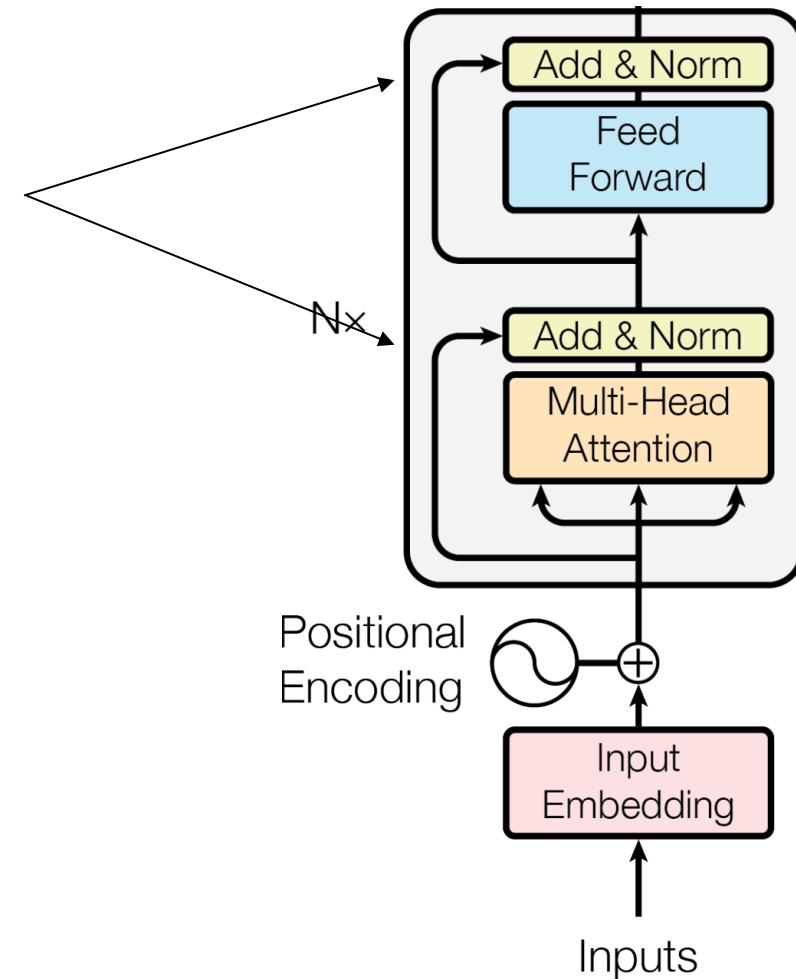
Transformer-Encoder

❖ Layer Normalization

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

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

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

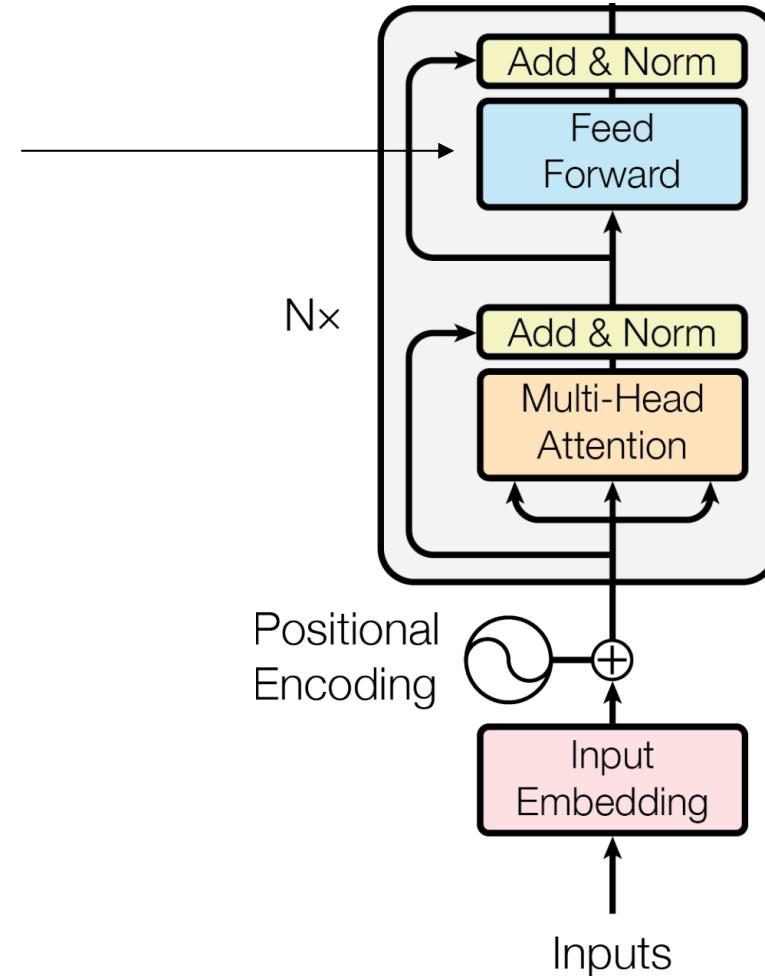
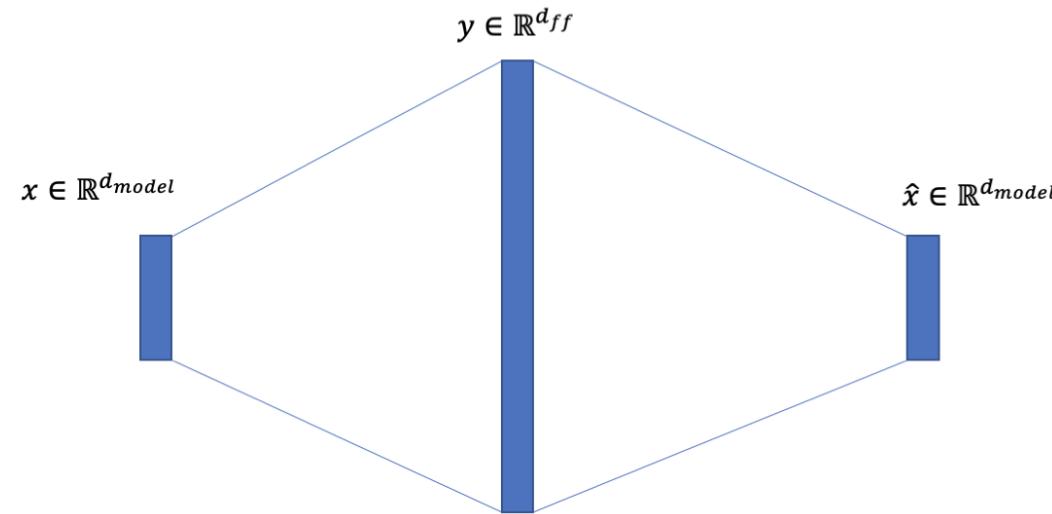


2 – Transformer



Transformer-Encoder

- ❖ Feed Forward
- ❖ 2 FC Layer



2 – Transformer



Transformer-Encoder – Demo

```
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
```

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

```
embedded.shape
```

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

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

```
encoded.shape
```

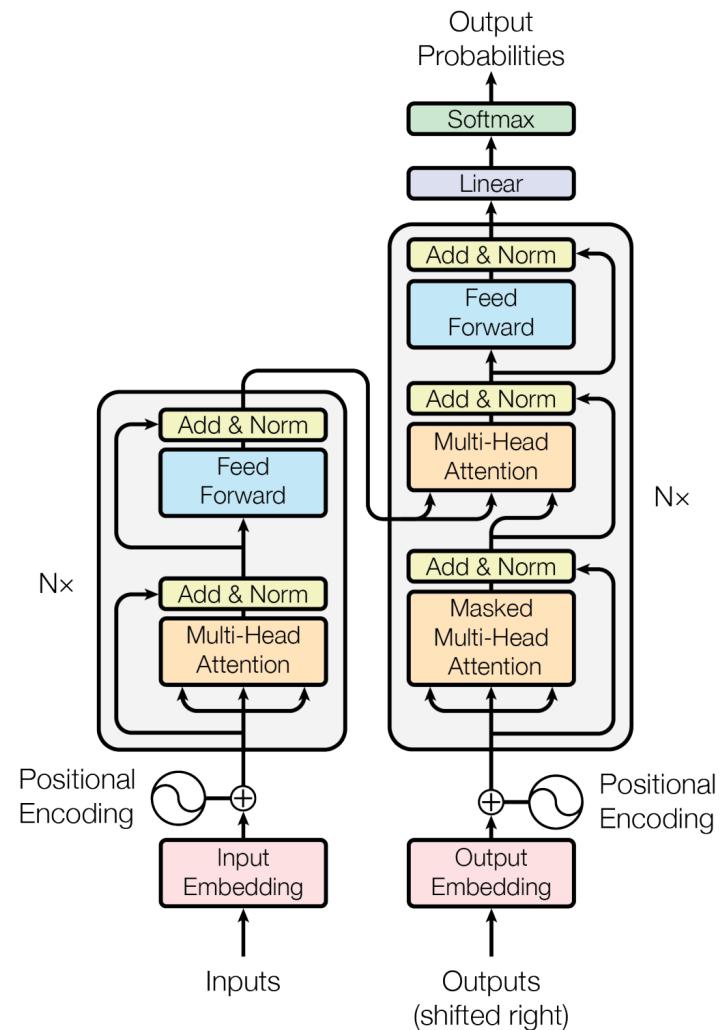
```
torch.Size([32, 50, 200])
```

2 – Transformer



Transformer-Decoder

- ❖ Masked Multi-Head Attention
- ❖ Multi-Head Attention
- ❖ Layer Normalization
- ❖ Feed Forward



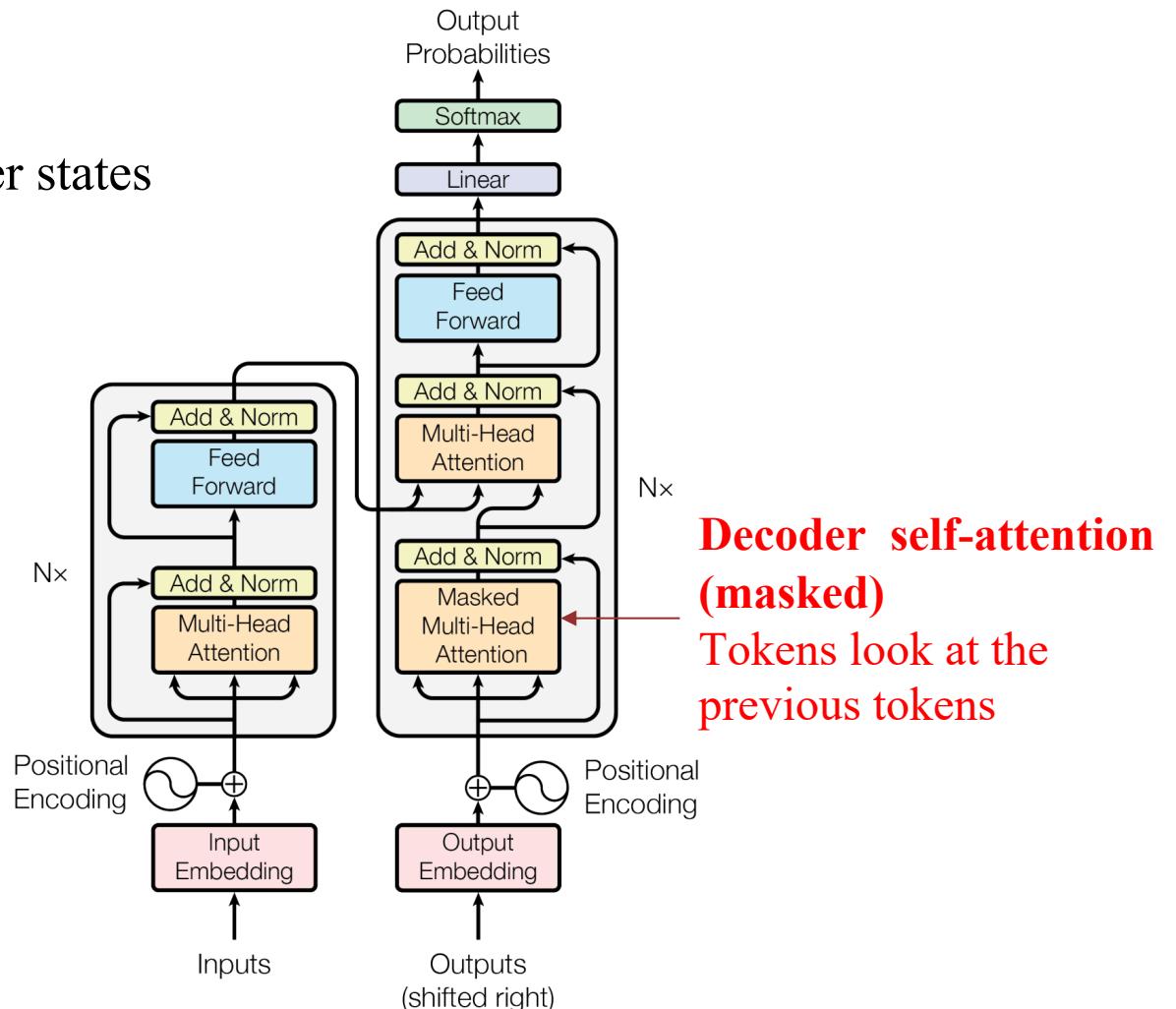
2 – Transformer



Transformer-Decoder

- ❖ **Masked Multi-Head Attention**
- ❖ Predict a next token based on the history token
- ❖ Queries, keys, values are computed from decoder states

Not Seen					
I	am	learn	#ing	nlp	
I	1.2	2.3	2.5	1.2	3.5
am	2.3	1.5	0.5	0.6	1.2
learn	0.5	1.4	1.6	0.3	4.8
#ing	0.6	1.8	2.4	0.3	1.2
nlp	2.1	2.3	0.2	2.0	2.5



2 – Transformer



Transformer-Decoder

- ❖ **Masked Multi-Head Attention**
- ❖ Predict a next token based on the history token
- ❖ Queries, keys, values are computed from decoder states

Not Seen				
I	am	learn	#ing	nlp
I	1.2	2.3	2.5	1.2
am	2.3	1.5	0.5	0.6
learn	0.5	1.4	1.6	0.3
#ing	0.6	1.8	2.4	0.3
nlp	2.1	2.3	0.2	2.0

1.2	-inf	-inf	-inf	-inf
2.3	1.5	-inf	-inf	-inf
0.5	1.4	1.6	-inf	-inf
0.6	1.8	2.4	0.3	-inf
2.1	2.3	0.2	2.0	2.5

Softmax

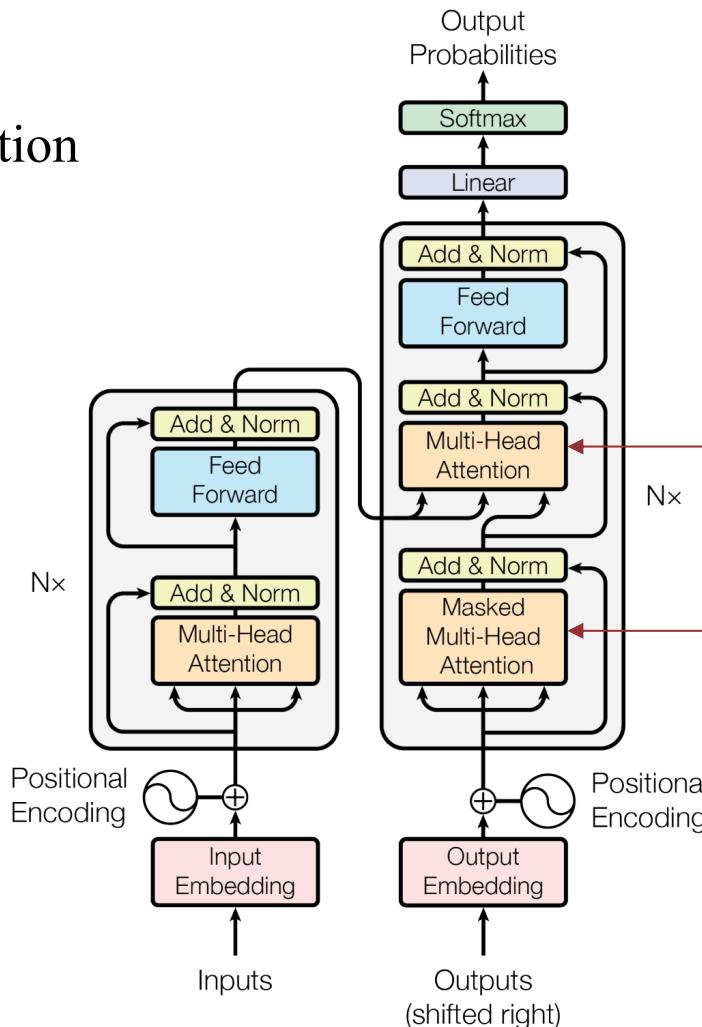
1.0	0.0	0.0	0.0	0.0
0.69	0.31	0.0	0.0	0.0
0.15	0.38	0.46	0.0	0.0
0.6	1.8	2.4	0.3	0.0
2.1	2.3	0.2	2.0	2.5

2 – Transformer



Transformer-Decoder

- ❖ Multi-Head Attention
- ❖ Sequence-to-sequence with attention
- ❖ Queries – from decoder states
- ❖ Keys, values from encoder states



Decoder – encoder attention
Target token look at the source

Decoder self-attention (masked)
Tokens look at the previous tokens

2 – Transformer



Transformer-Decoder – Demo

```
class TransformerDecoderBlock(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.cross_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.layernorm_3 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
        self.dropout_1 = nn.Dropout(p=dropout)
        self.dropout_2 = nn.Dropout(p=dropout)
        self.dropout_3 = nn.Dropout(p=dropout)
```

```
    def forward(self, x, enc_output, src_mask, tgt_mask):
        attn_output, _ = self.attn(x, x, x, attn_mask=tgt_mask)
        attn_output = self.dropout_1(attn_output)
        out_1 = self.layernorm_1(x + attn_output)

        attn_output, _ = self.cross_attn(
            out_1, enc_output, enc_output, attn_mask=src_mask
        )
        attn_output = self.dropout_2(attn_output)
        out_2 = self.layernorm_2(out_1 + attn_output)

        ffn_output = self.ffn(out_2)
        ffn_output = self.dropout_2(ffn_output)
        out_3 = self.layernorm_2(out_2 + ffn_output)
        return out_3
```

2 – Transformer



Transformer – Demo

```
class Transformer(nn.Module):
    def __init__(self,
                 src_vocab_size, tgt_vocab_size,
                 embed_dim, max_length, num_layers, num_heads, ff_dim,
                 dropout=0.1, device='cpu'):
        super().__init__()
        self.device = device
        self.encoder = TransformerEncoder(
            src_vocab_size, embed_dim, max_length, num_layers, num_heads, ff_dim)
        self.decoder = TransformerDecoder(
            tgt_vocab_size, embed_dim, max_length, num_layers, num_heads, ff_dim)
        self.fc = nn.Linear(embed_dim, tgt_vocab_size)
```

```
def generate_mask(self, src, tgt):
    src_seq_len = src.shape[1]
    tgt_seq_len = tgt.shape[1]

    src_mask = torch.zeros(
        (src_seq_len, src_seq_len),
        device=self.device).type(torch.bool)

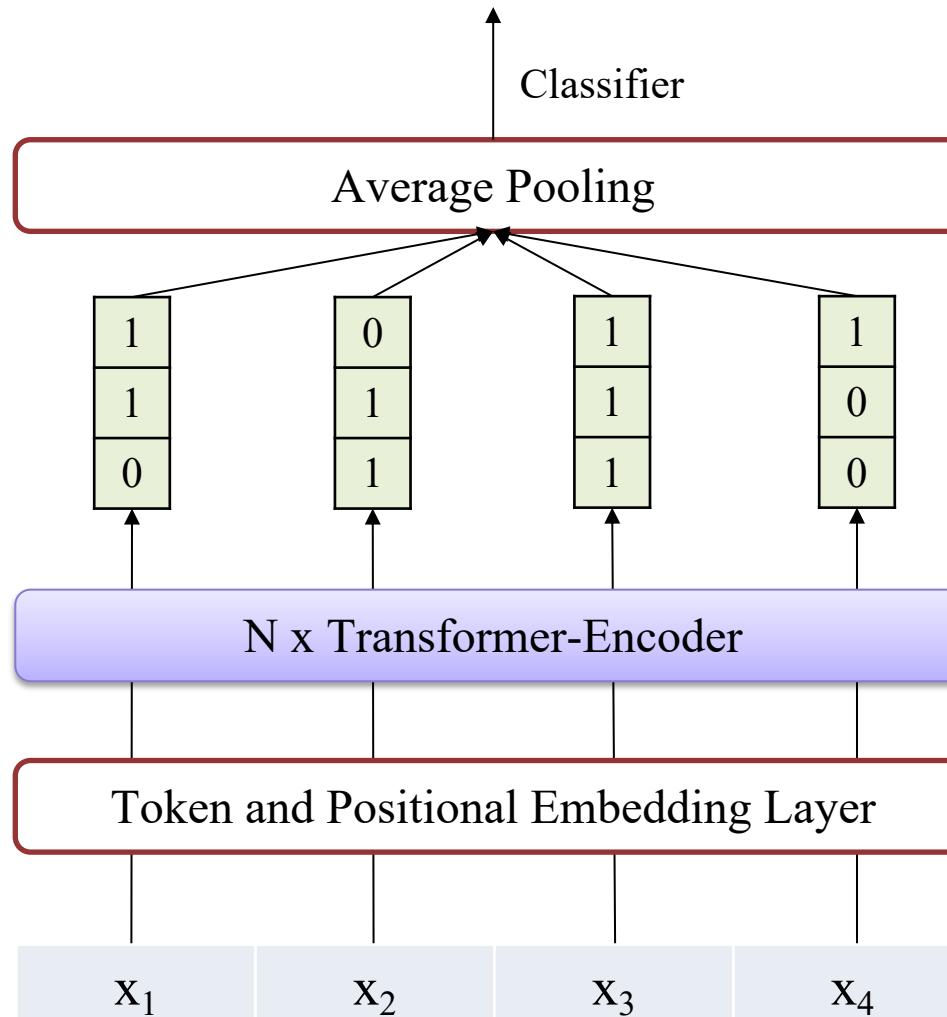
    tgt_mask = (torch.triu(torch.ones(
        (tgt_seq_len, tgt_seq_len),
        device=self.device))
        == 1).transpose(0, 1)
    tgt_mask = tgt_mask.float().masked_fill(
        tgt_mask == 0, float('-inf')).masked_fill(tgt_mask == 1, float(0.0))
    return src_mask, tgt_mask

def forward(self, src, tgt):
    src_mask, tgt_mask = self.generate_mask(src, tgt)
    enc_output = self.encoder(src)
    dec_output = self.decoder(tgt, enc_output, src_mask, tgt_mask)
    output = self.fc(dec_output)
    return output
```

2 – Transformer

!

Text Classification using Transformer-Encoder – Demo



2 – Transformer



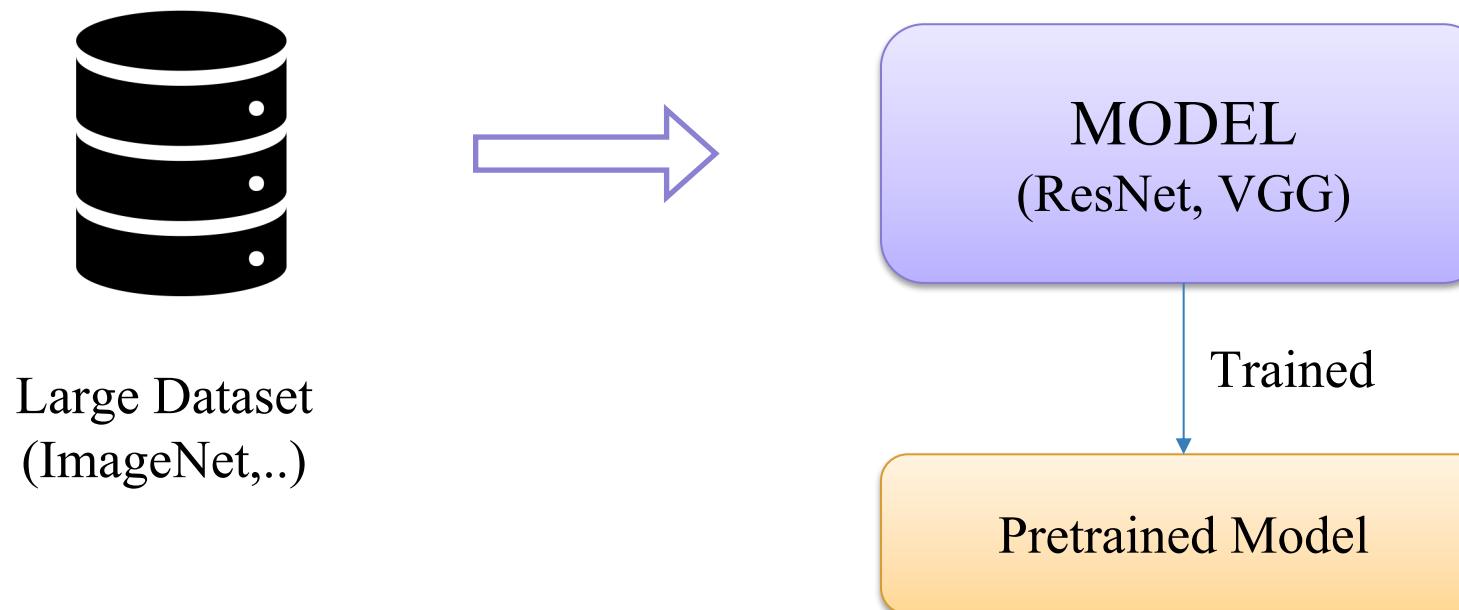
Text Classification using Transformer-Encoder

```
class TransformerEncoderCls(nn.Module):
    def __init__(self,
                 vocab_size, max_length, num_layers, embed_dim, num_heads, ff_dim,
                 dropout=0.1, device='cpu'):
        super().__init__()
        self.encoder = TransformerEncoder(
            vocab_size, embed_dim, max_length, num_layers, num_heads, ff_dim, dropout, device)
        self.pooling = nn.AvgPool1d(kernel_size=max_length)
        self.fc1 = nn.Linear(in_features=embed_dim, out_features=20)
        self.fc2 = nn.Linear(in_features=20, out_features=2)
        self.dropout = nn.Dropout(p=dropout)
        self.relu = nn.ReLU()
    def forward(self, x):
        output = self.encoder(x)
        output = self.pooling(output.permute(0, 2, 1)).squeeze()
        output = self.dropout(output)
        output = self.fc1(output)
        output = self.dropout(output)
        output = self.fc2(output)
        return output
```



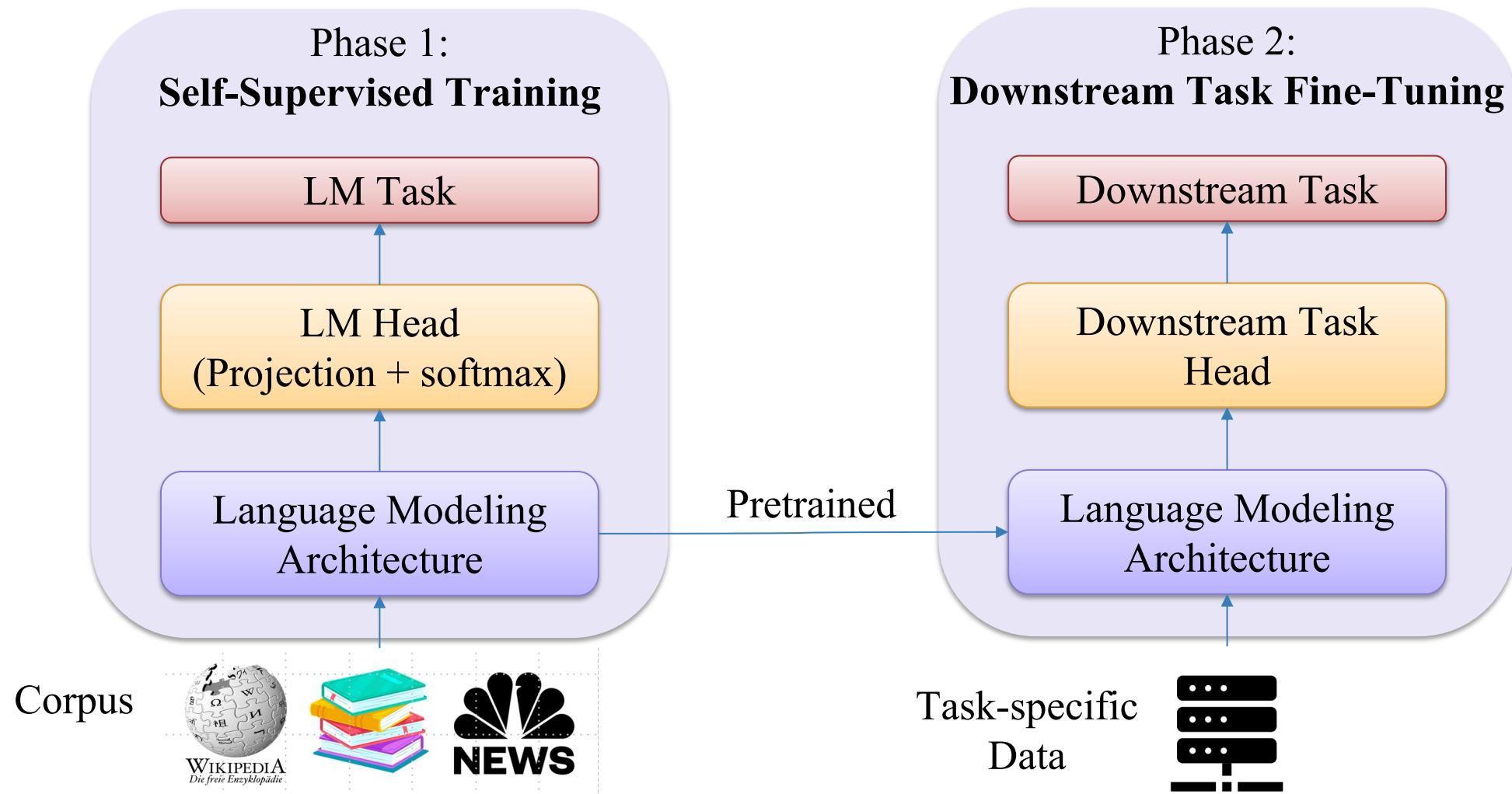
Pretrained Model for Image

- ❖ **ImageNet**
- ❖ Training: 1,281,167 images. Validation: 50,000 images. Testing: 100,000 images
- ❖ Object classes: 1,000





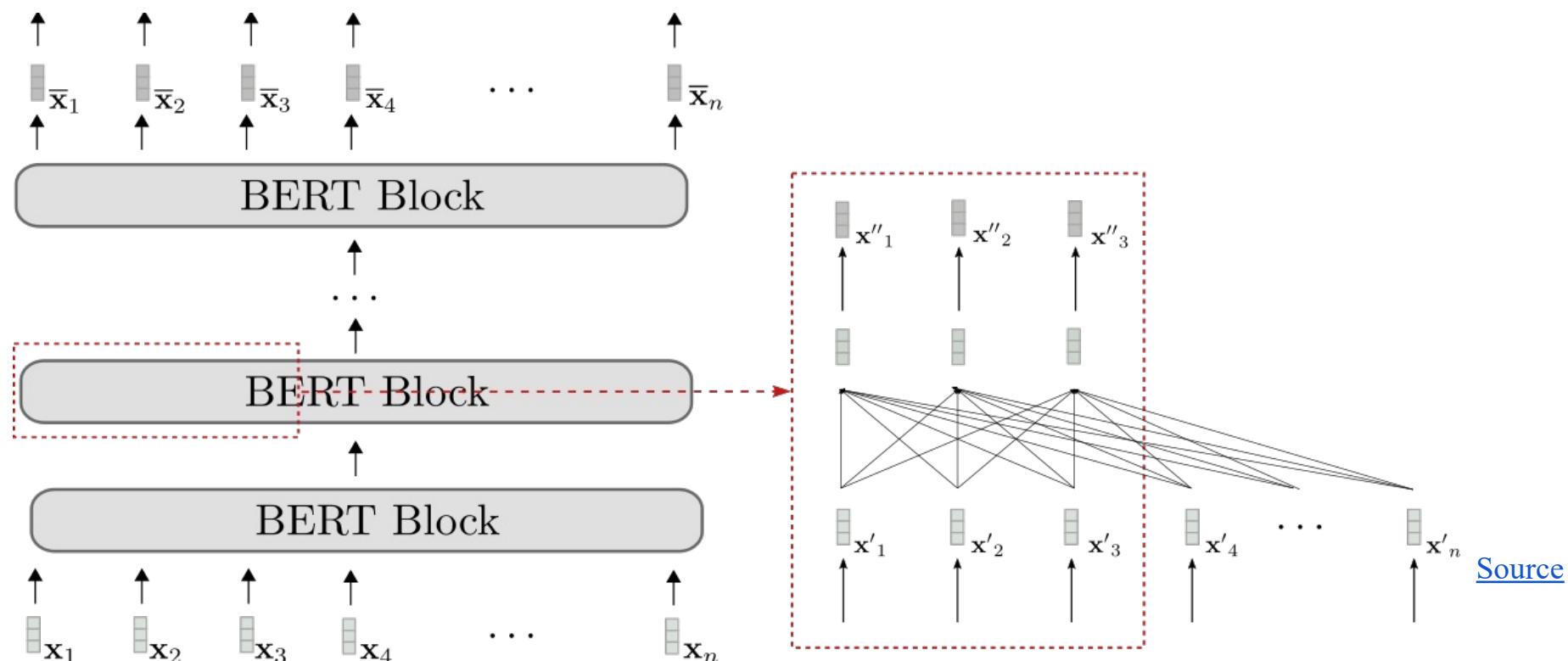
Pretrained Model for Text





BERT

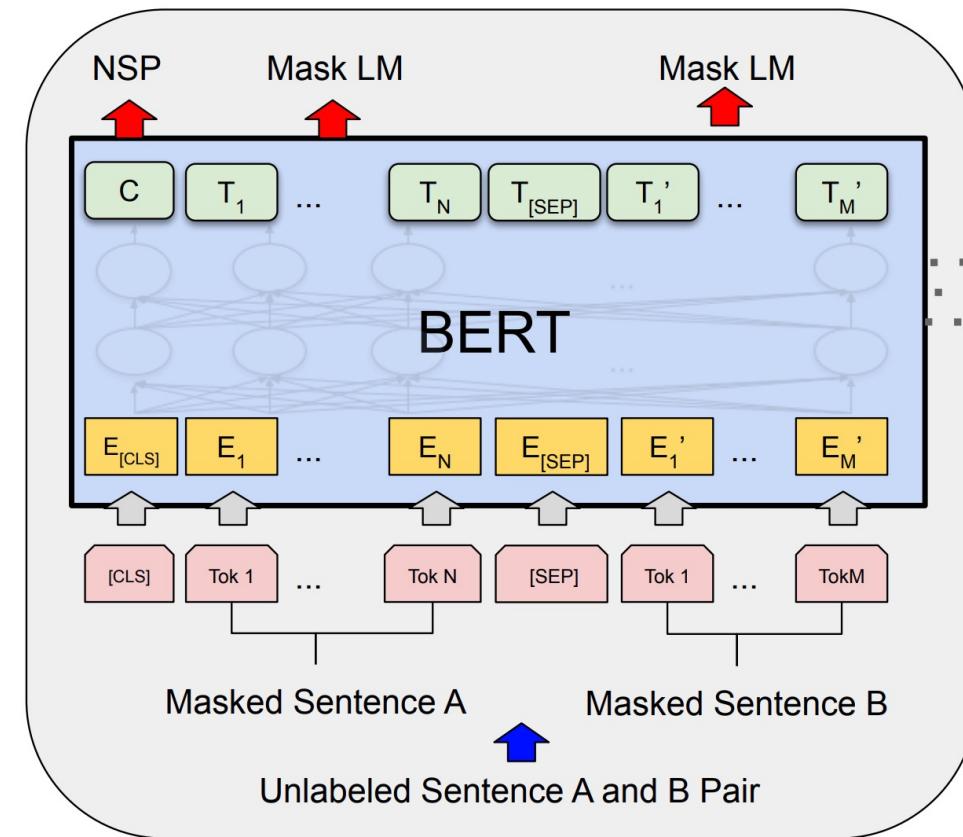
- ❖ BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding
- ❖ BERT: An encoder-only model
- ❖ Maps an input sequence to a contextualized sequence: $f_{\theta_{BERT}}: \mathbf{X}_{1:n} \rightarrow \overline{\mathbf{X}}_{1:n}$





Objective Functions

- ❖ Masked LM (15% token):
 - 80%: replace with [MASK]
 - 10%: replace with a random word
 - 10%: keep unchanged
- ❖ Next Sentence Prediction (NSP)
 - Classification Task
 - 2 Labels: IsNext and NotNext
 - Use [SEP] [CLS] token



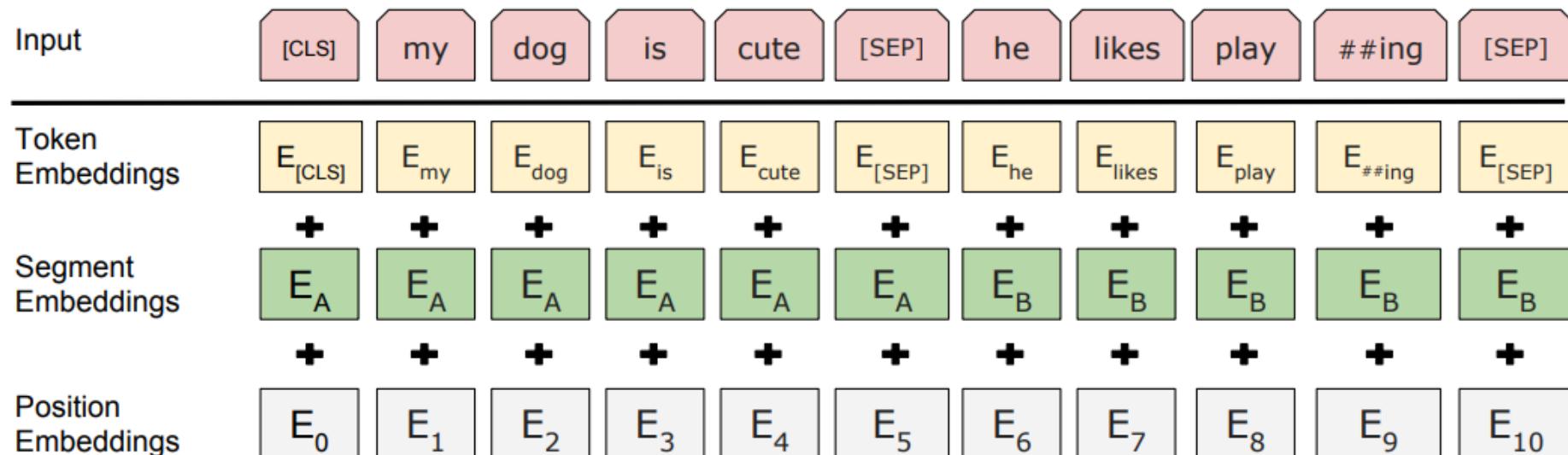
Pre-training

[Source](#)



Input Representation

❖ BERT Input Representation





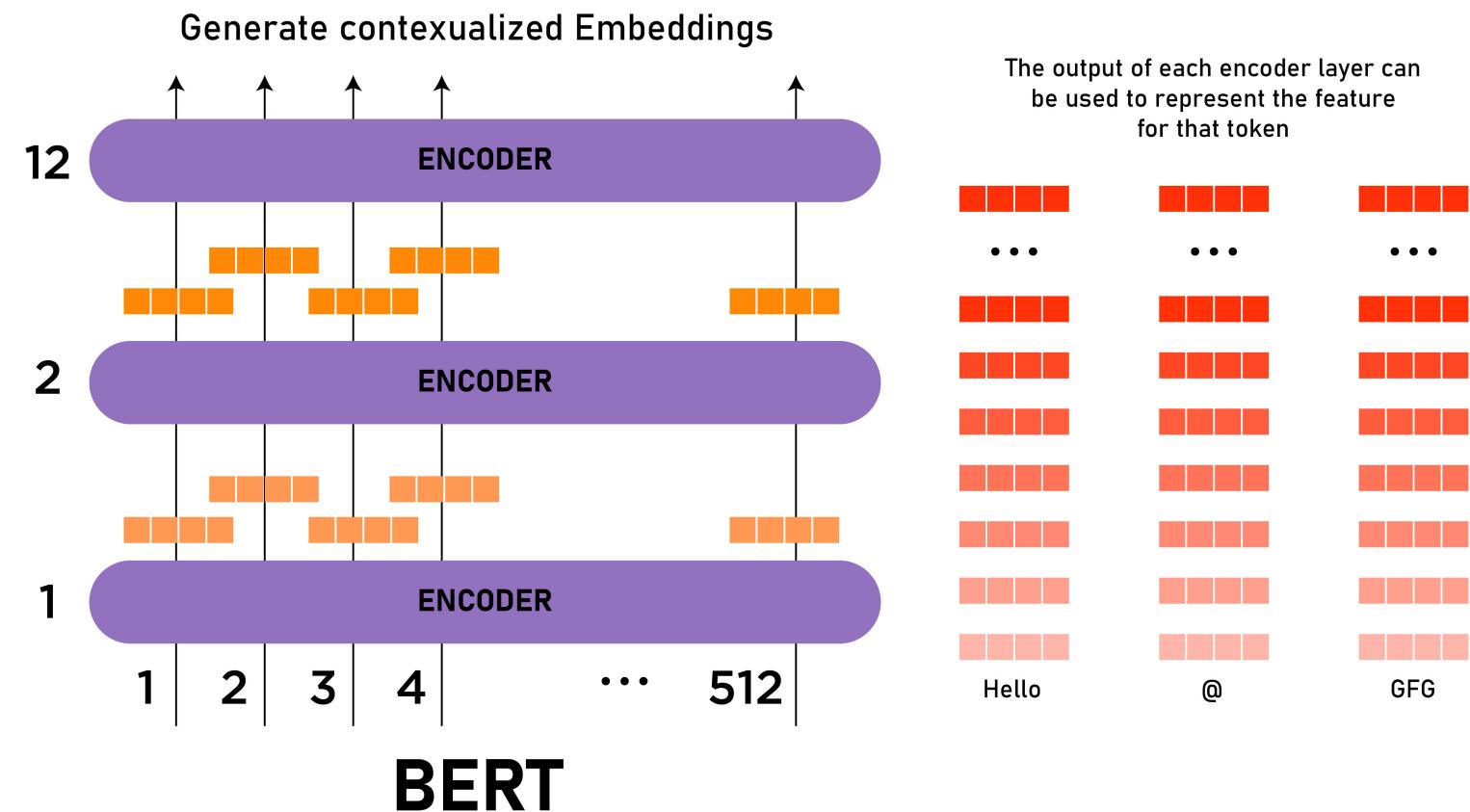
BERT Training

❖ BERT base:

L=12, H=768, A=12, P=110M

❖ BERT large:

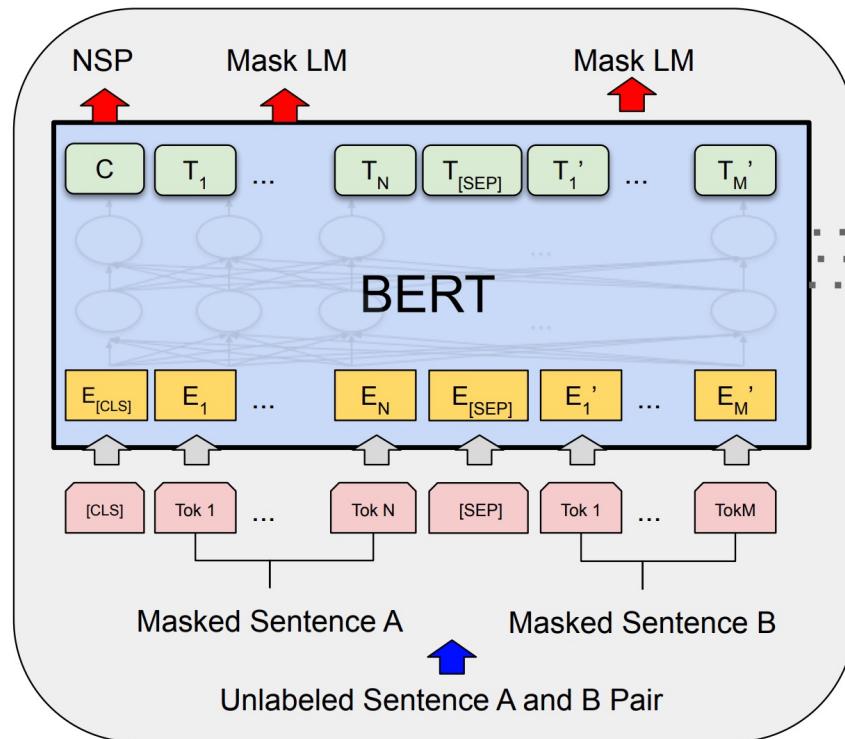
L=24, H=1024, A=16, P=340M



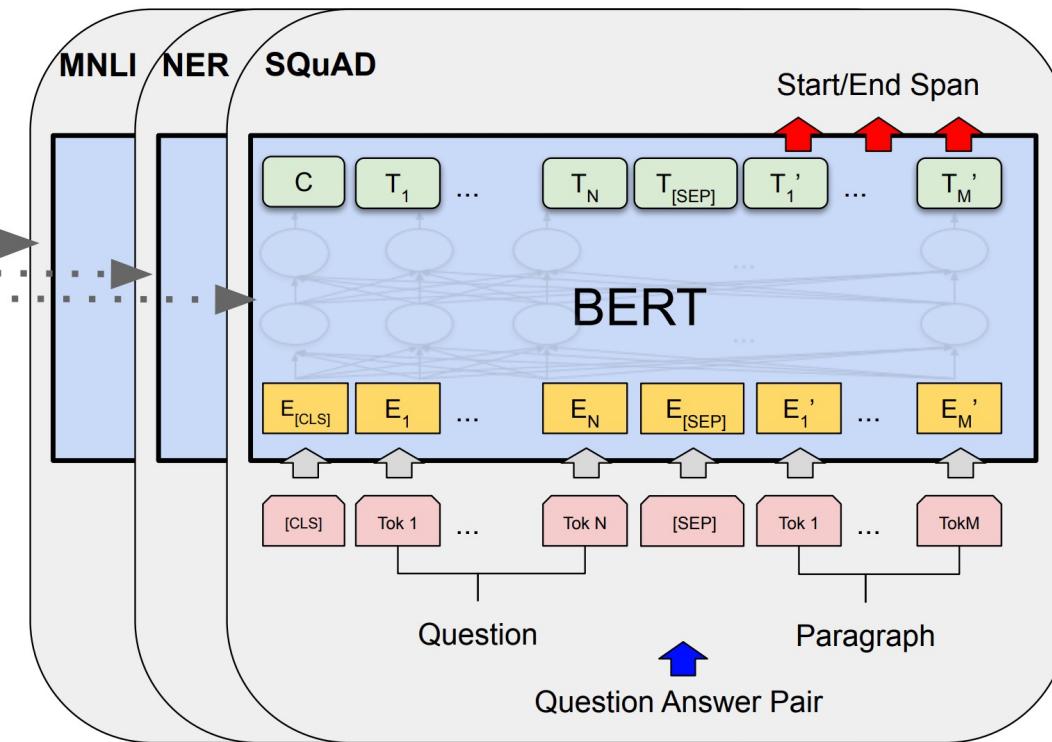
3 – BERT



Fine-Tuning



Pre-training

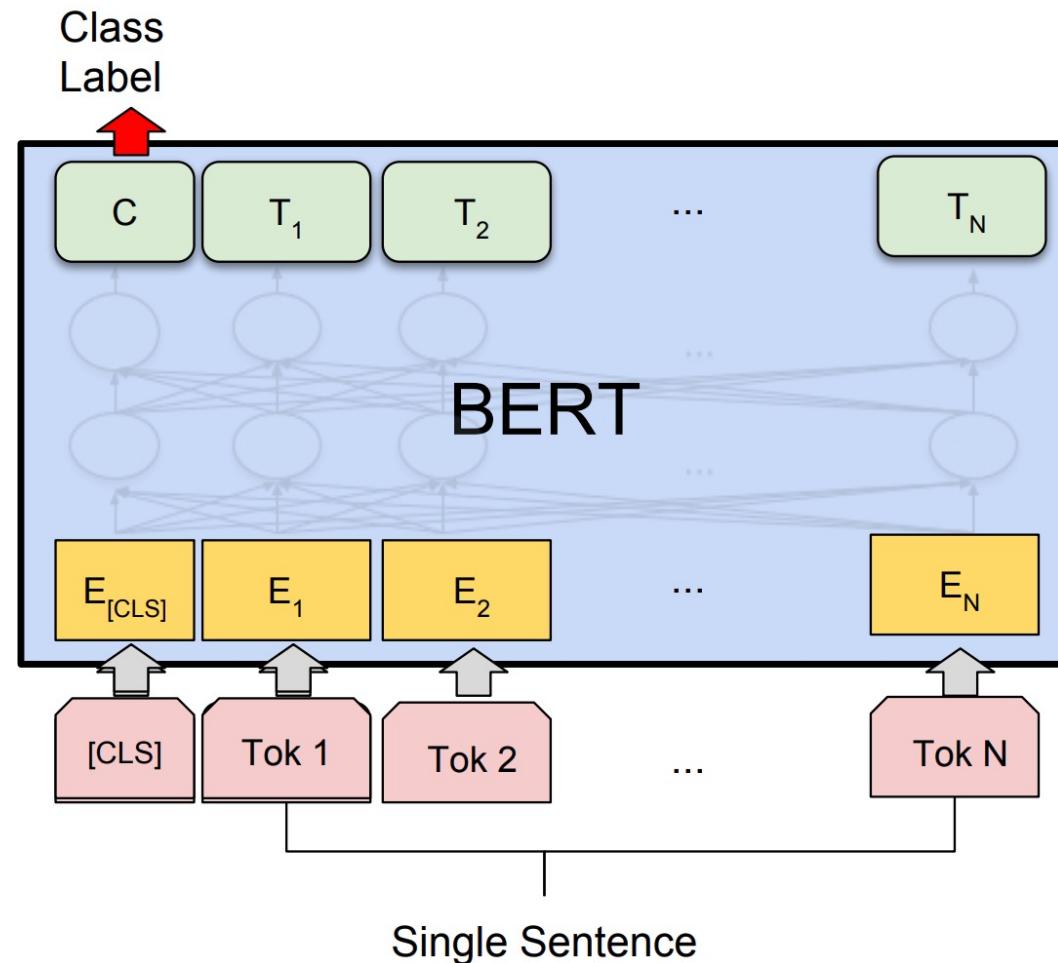


Fine-Tuning

Source



Text Classification using BERT – Demo





Text Classification using BERT

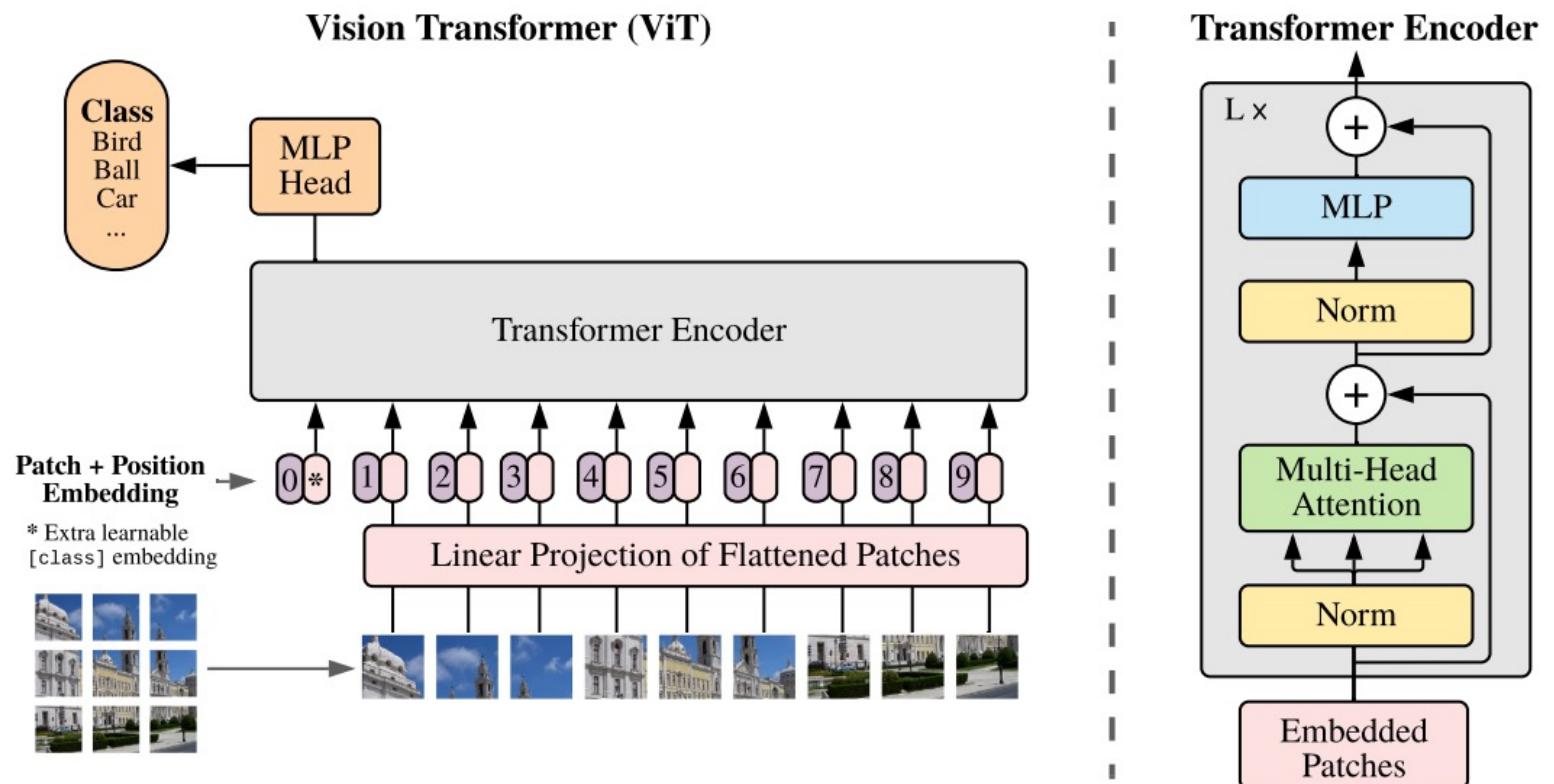
- ❖ NTC-SCV Dataset

Model	Accuracy
Transformer-Encoder	82.55
BERT	85.52
RNN	80.92
TextCNN	88.78

4 – Vision Transformer



Architecture



4 – Vision Transformer

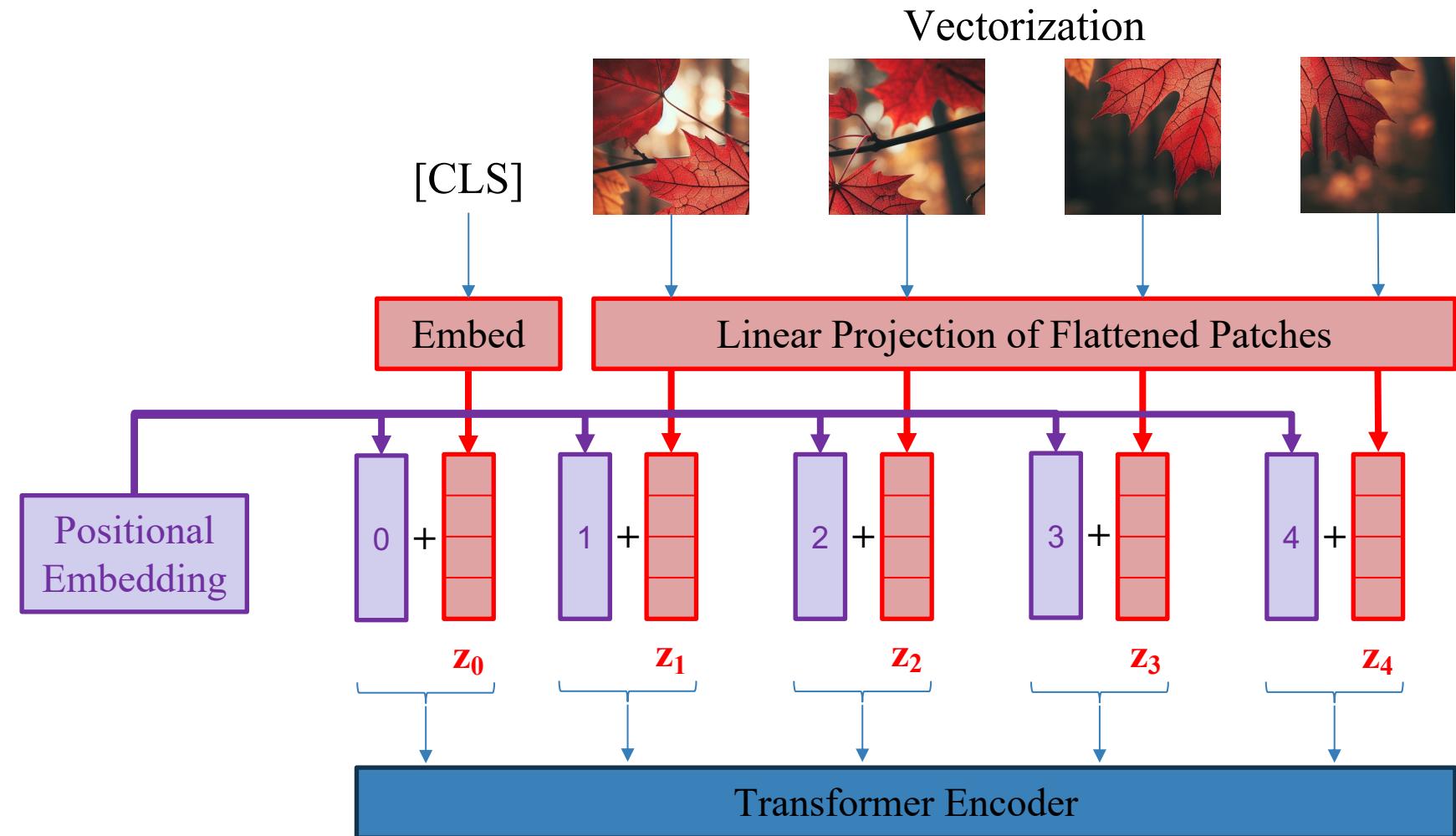


Embedded Patches

Patches



Original



4 – Vision Transformer

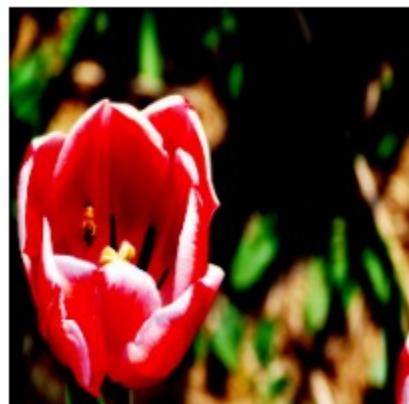
!

Flower Dataset

dandelion



tulips



roses



sunflowers



sunflowers



4 – Vision Transformer



Implementation (Scratch) - Linear Projection – Demo

```
1  class TransformerEncoder(nn.Module):
2      def __init__(self, embed_dim, num_heads, ff_dim, dropout=0.1):
3          super().__init__()
4          self.attn = nn.MultiheadAttention(
5              embed_dim=embed_dim,
6              num_heads=num_heads,
7              batch_first=True
8          )
9          self.ffn = nn.Sequential(
10             nn.Linear(in_features=embed_dim, out_features=ff_dim, bias=True),
11             nn.ReLU(),
12             nn.Linear(in_features=ff_dim, out_features=embed_dim, bias=True)
13         )
14          self.layernorm_1 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
15          self.layernorm_2 = nn.LayerNorm(normalized_shape=embed_dim, eps=1e-6)
16          self.dropout_1 = nn.Dropout(p=dropout)
17          self.dropout_2 = nn.Dropout(p=dropout)
18
19      def forward(self, query, key, value):
20          attn_output, _ = self.attn(query, key, value)
21          attn_output = self.dropout_1(attn_output)
22          out_1 = self.layernorm_1(query + attn_output)
23          ffn_output = self.ffn(out_1)
24          ffn_output = self.dropout_2(ffn_output)
25          out_2 = self.layernorm_2(out_1 + ffn_output)
26          return out_2
```

4 – Vision Transformer

!

Implementation (Scratch) - Linear Projection – Demo

```
1  class PatchPositionEmbedding(nn.Module):
2      def __init__(self, image_size=224, embed_dim=512, patch_size=16, device='cpu'):
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          scale = embed_dim ** -0.5
6          self.positional_embedding = nn.Parameter(scale * torch.randn((image_size // patch_size) ** 2, embed_dim))
7          self.device = device
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(self.device)
15
16         return x
```

add positional
embedding to
patch
embedding

initialize
positional
embedding

4 – Vision Transformer

!

Implementation (Scratch) - Linear Projection – Demo

```
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.embed_layer = PatchPositionEmbedding(
8              image_size=image_size, embed_dim=embed_dim, patch_size=patch_size, device=device
9          )
10         self.transformer_layer = TransformerEncoder(
11             embed_dim, num_heads, ff_dim, dropout
12         )
13         self.fc1 = nn.Linear(in_features=embed_dim, out_features=20)
14         self.fc2 = nn.Linear(in_features=20, out_features=num_classes)
15         self.dropout = nn.Dropout(p=dropout)
16         self.relu = nn.ReLU()
17     def forward(self, x):
18         output = self.embed_layer(x)
19         output = self.transformer_layer(output, output, output)
20         output = output[:, 0, :]
21         output = self.dropout(output)
22         output = self.fc1(output)
23         output = self.dropout(output)
24         output = self.fc2(output)
25     return output
```

4 – Vision Transformer

!

Implementation (Pretrained) – Demo

```
from transformers import ViTForImageClassification

id2label = {id:label for id, label in enumerate(classes)}
label2id = {label:id for id,label in id2label.items()}

model = ViTForImageClassification.from_pretrained(
    'google/vit-base-patch16-224-in21k',
    num_labels=num_classes,
    id2label=id2label,
    label2id=label2id
)
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
model.to(device)
```

4 – Vision Transformer

!

Pretrained - google/vit-base-patch16-224-in21k

Dataset	ImageNet-21k
Num of Images	14 millions
Num of Classes	21,000
Resized	224 x 224
Normalized Mean	(0.5, 0.5, 0.5)
Normalized Standard Deviation	(0.5, 0.5, 0.5)
Hardware	TPUv3 (8 cores)
Batch Size	4096

Thanks!

Any questions?