

344.063 KV Special Topic:

# Natural Language Processing with Deep Learning Transformers



Navid Rekab-saz

[navid.rekabsaz@jku.at](mailto:navid.rekabsaz@jku.at)

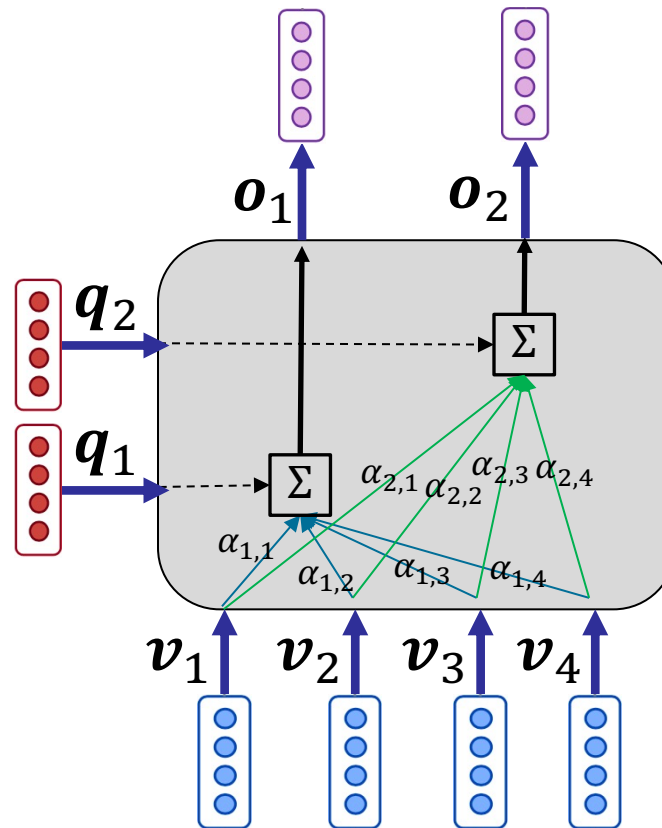
# Agenda

- Transformers
  - Transformer encoder
  - Transformer decoder
- seq2seq with Transformers

# Agenda

- **Transformers**
  - **Transformer encoder**
  - **Transformer decoder**
- seq2seq with Transformers

# Attentions! – recap



$\alpha_{i,j}$  is the attention score of query  $q_i$  on value  $v_j$

$\alpha_i$  is the vector of attentions of query  $q_i$  over value vectors  $V$  which forms a probability distribution

# Attention Networks – recap

- Given query vector  $\mathbf{q}_i$ , an attention network uses the **attention similarity function**  $f$  to assign a **non-normalized attention score**  $\tilde{\alpha}_{i,j}$  to value vector  $\mathbf{v}_j$ :

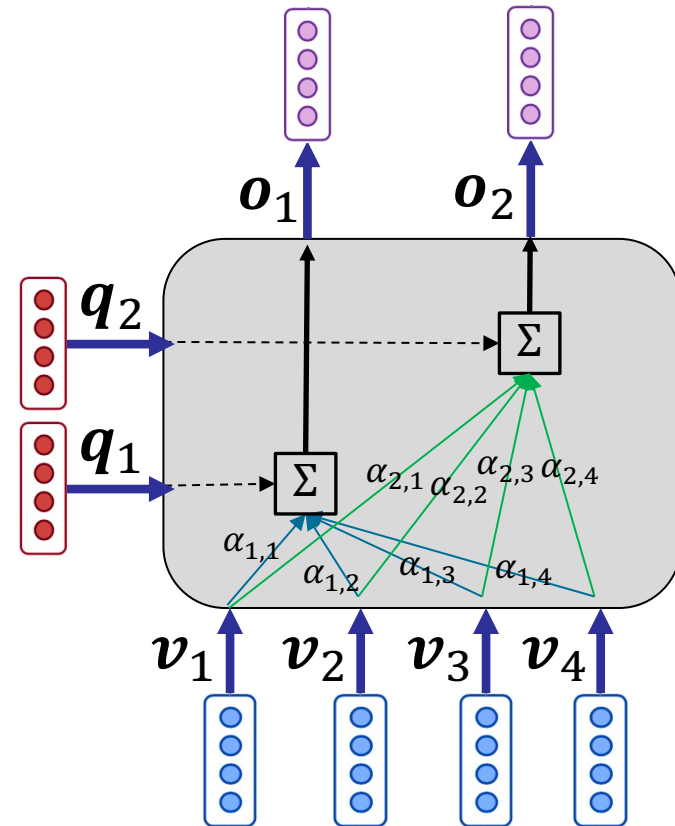
$$\tilde{\alpha}_{i,j} = f(\mathbf{q}_i, \mathbf{v}_j)$$

- Then, the attention scores over values are turned to a probability distribution using softmax:

$$\alpha_i = \text{softmax}(\tilde{\alpha}_i), \quad \sum_{j=1}^{|V|} \alpha_{i,j} = 1$$

- Finally, output vector  $\mathbf{o}_i$  regarding query  $\mathbf{q}_i$  is defined as the **sum** of the value vectors **weighted** by their corresponding attentions:

$$\mathbf{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \mathbf{v}_j$$

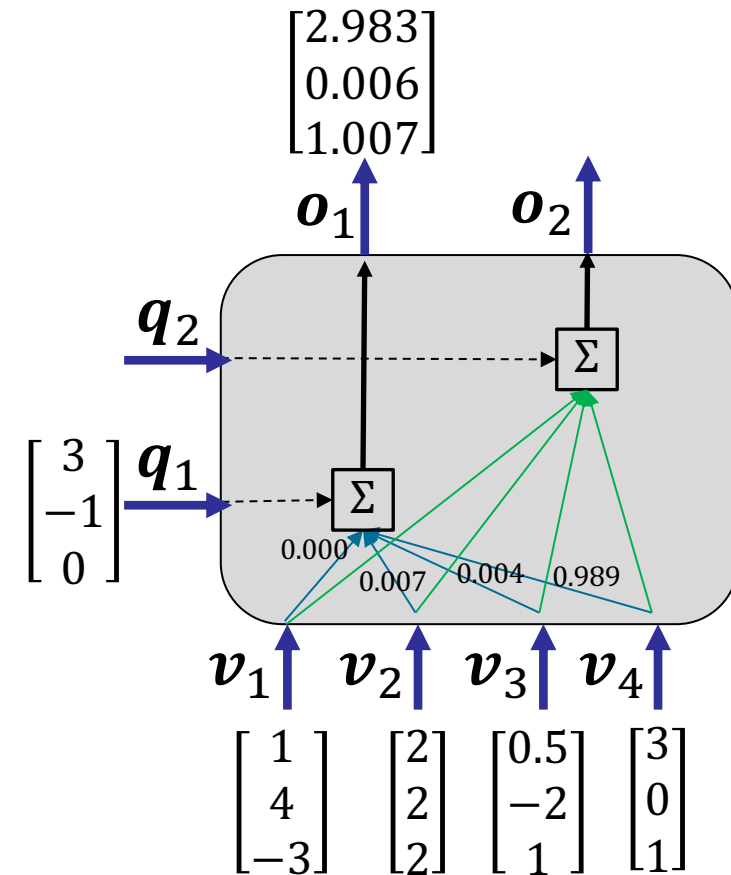


## Example – recap

$$\tilde{\alpha}_1 = \begin{bmatrix} \mathbf{q}_1 \mathbf{v}_1^T = -1 \\ \mathbf{q}_1 \mathbf{v}_2^T = 4 \\ \mathbf{q}_1 \mathbf{v}_3^T = 3.5 \\ \mathbf{q}_1 \mathbf{v}_4^T = 9 \end{bmatrix} \rightarrow \alpha_1 = \begin{bmatrix} 0.000 \\ 0.007 \\ 0.004 \\ 0.989 \end{bmatrix}$$

$$\mathbf{o}_1 = 0.000 \begin{bmatrix} 1 \\ 4 \\ -3 \end{bmatrix} + 0.007 \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} + 0.004 \begin{bmatrix} 0.5 \\ -2 \\ 1 \end{bmatrix} + 0.989 \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}$$

$$\mathbf{o}_1 = \begin{bmatrix} 2.983 \\ 0.006 \\ 1.007 \end{bmatrix}$$

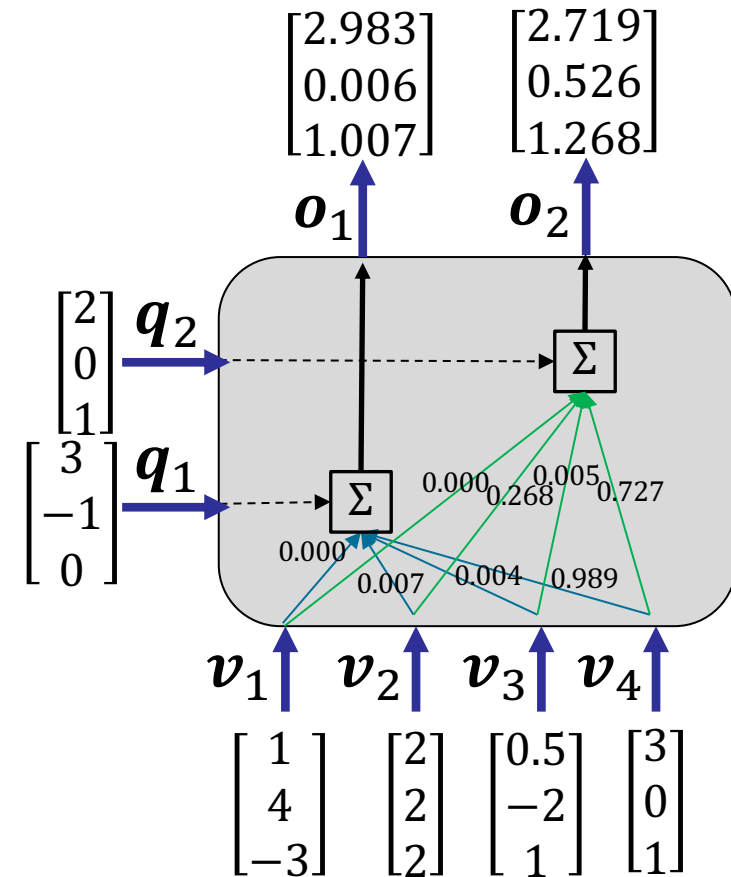


## Example – recap

$$\tilde{\alpha}_2 = \begin{bmatrix} \mathbf{q}_2 \mathbf{v}_1^T = -1 \\ \mathbf{q}_2 \mathbf{v}_2^T = 6 \\ \mathbf{q}_2 \mathbf{v}_3^T = 2 \\ \mathbf{q}_2 \mathbf{v}_4^T = 7 \end{bmatrix} \rightarrow \alpha_2 = \begin{bmatrix} 0.000 \\ 0.268 \\ 0.005 \\ 0.727 \end{bmatrix}$$

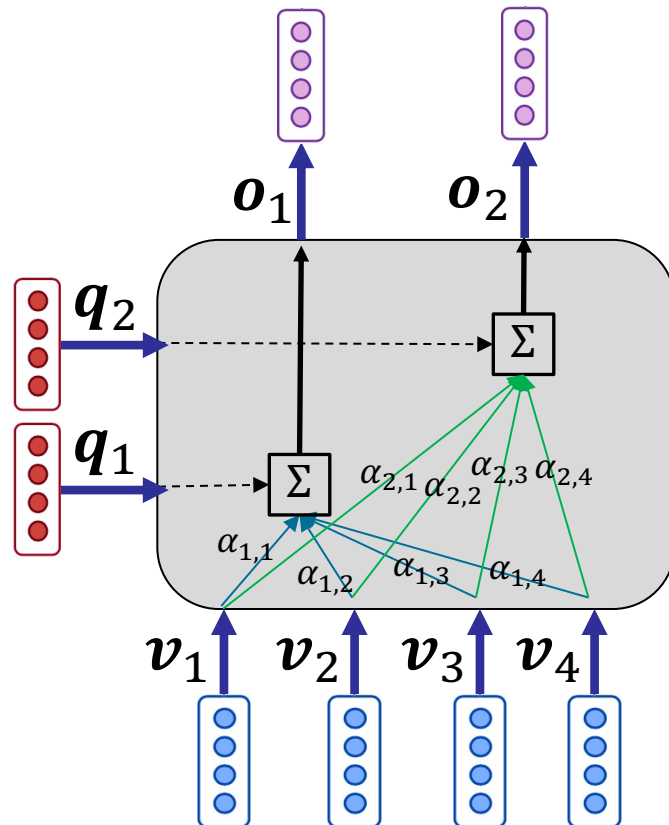
$$\mathbf{o}_2 = 0.000 \begin{bmatrix} 1 \\ 4 \\ -3 \end{bmatrix} + 0.268 \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} + 0.005 \begin{bmatrix} 0.5 \\ -2 \\ 1 \end{bmatrix} + 0.727 \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}$$

$$\mathbf{o}_2 = \begin{bmatrix} 2.719 \\ 0.526 \\ 1.268 \end{bmatrix}$$



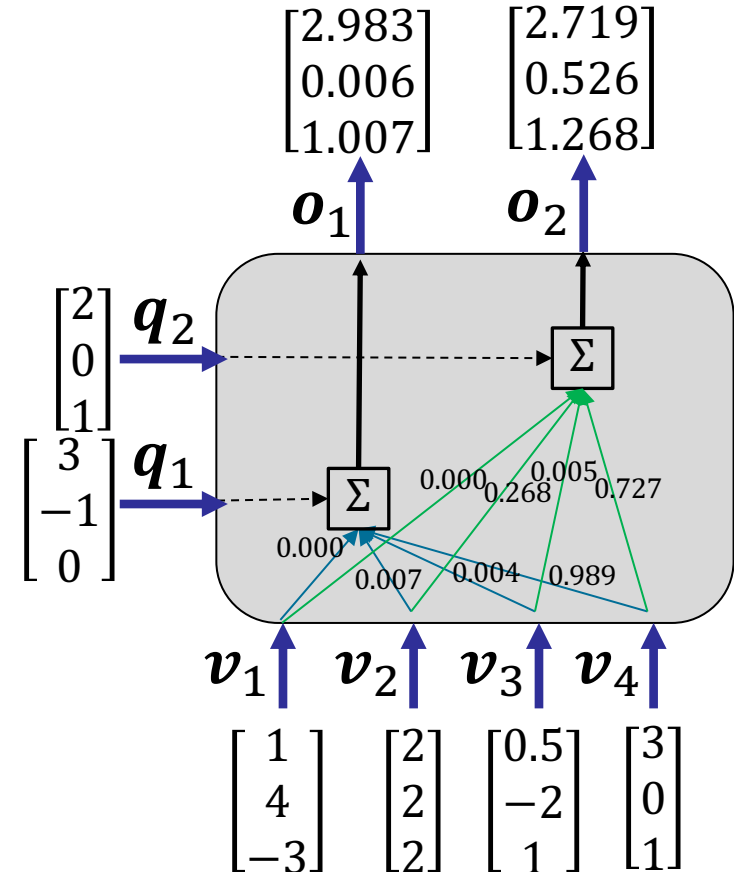
# Attention table

	$v_1$	$v_2$	$v_3$	$v_4$
$q_1$	$\alpha_{1,1}$	$\alpha_{1,2}$	$\alpha_{1,3}$	$\alpha_{1,4}$
$q_2$	$\alpha_{2,1}$	$\alpha_{2,2}$	$\alpha_{2,3}$	$\alpha_{2,4}$



In the example:

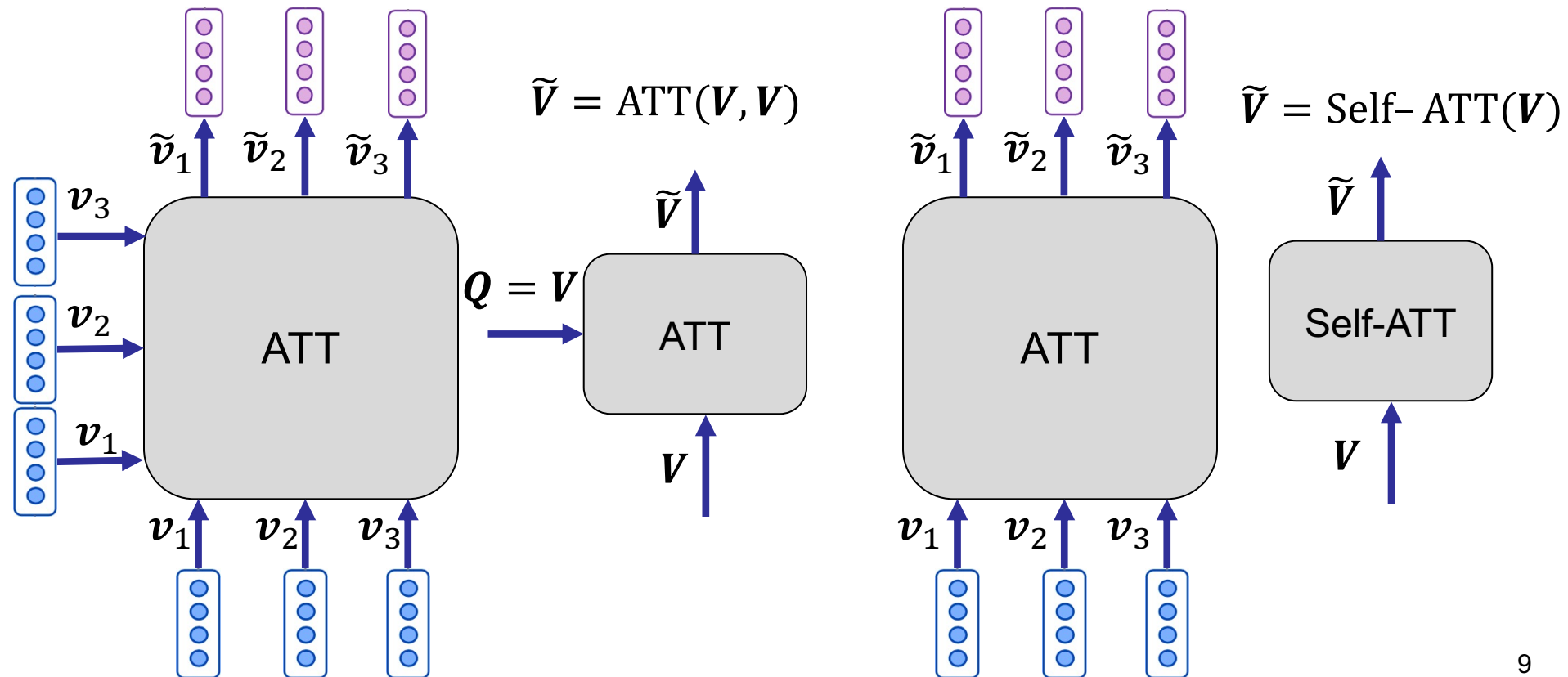
	$v_1$	$v_2$	$v_3$	$v_4$
$q_1$	0.000	0.007	0.004	0.989
$q_2$	0.000	0.268	0.005	0.727





# Self-attention

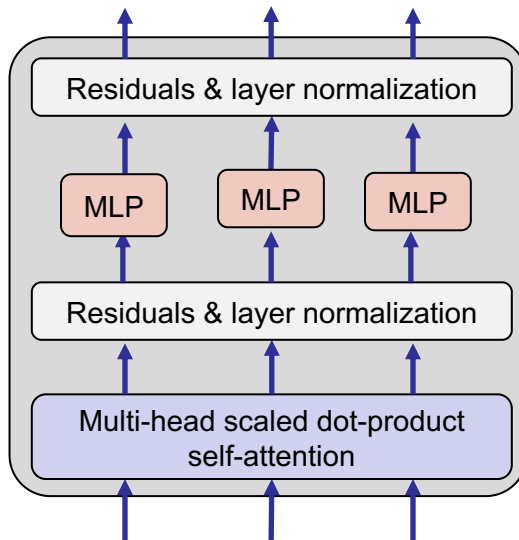
- Self-attention is when the values are also given as the queries:  $Q = V$
- Self-attention **encodes** a sequence  $V$  to a **contextualized sequence**  $\tilde{V}$ 
  - In self-attention, each input vector  $v_i$  attends to all other input vectors  $V$ , and outputs  $\tilde{v}_i$  as a composition of input vectors
  - Output vector  $\tilde{v}_i$  is the **contextual embedding** of the input vector  $v_i$



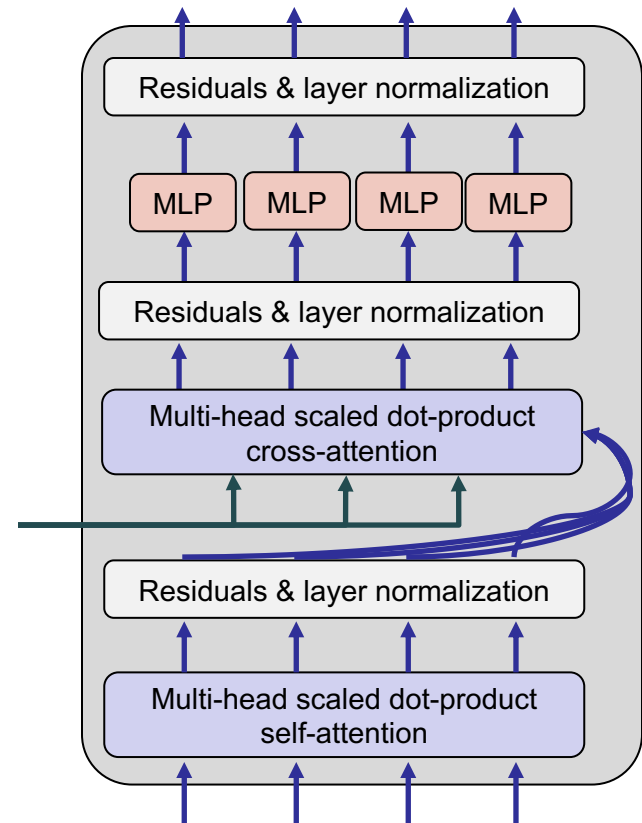
# Transformers

- Attention network with DL best practices!
  - Originally introduced in the context of machine translation and is now widely adopted for [sequence encoding](#) and [decoding](#)

## Transformer Encoder

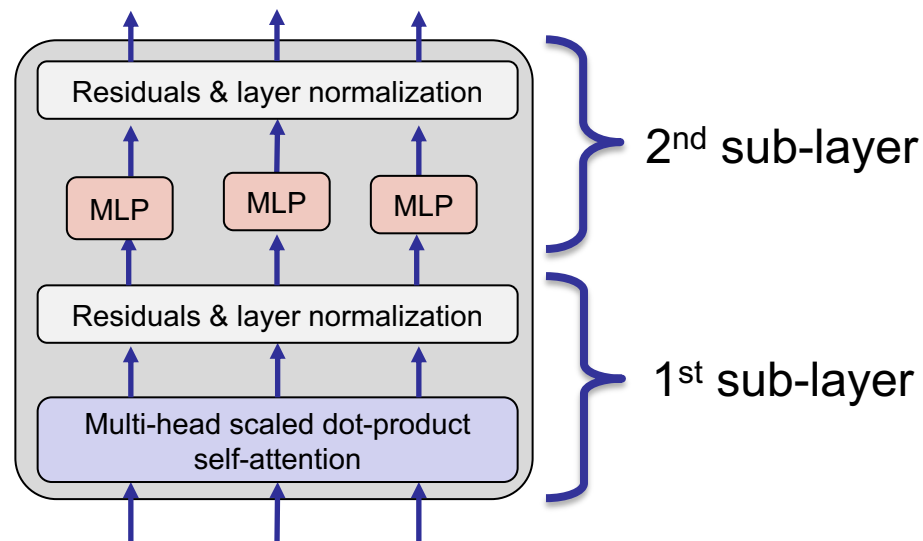


## Transformer Decoder



# Transformer Encoder

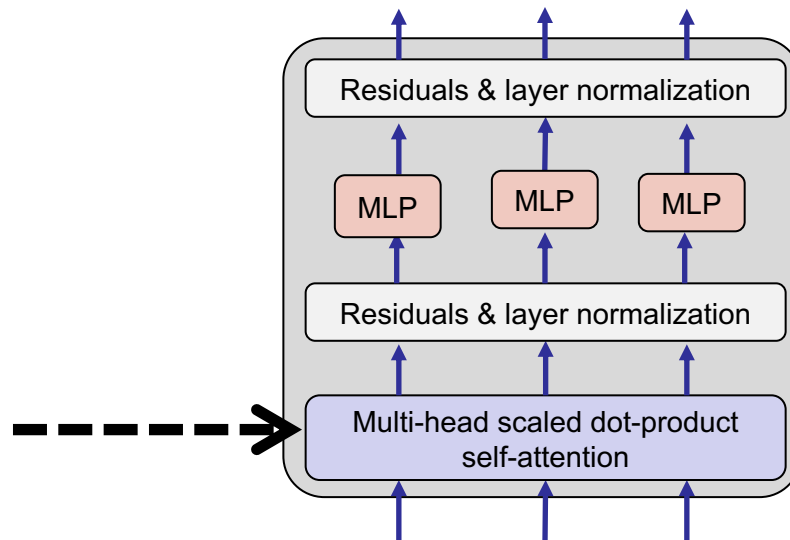
- Transformer Encoder consists of two sub-layers:
  - 1<sup>st</sup> : Multi-head scaled dot-product self-attention
  - 2<sup>nd</sup> : Position-wise multi-layer perceptron (feed forward)
- Each sub-layer is followed by residual networks and layer normalization
  - Drop-outs are applied after each computation



# Transformer Encoder

Let's start from multi-head scaled dot-product self-attention:

1. Scaled dot-product attention
2. Multi-head attention
3. self-attention

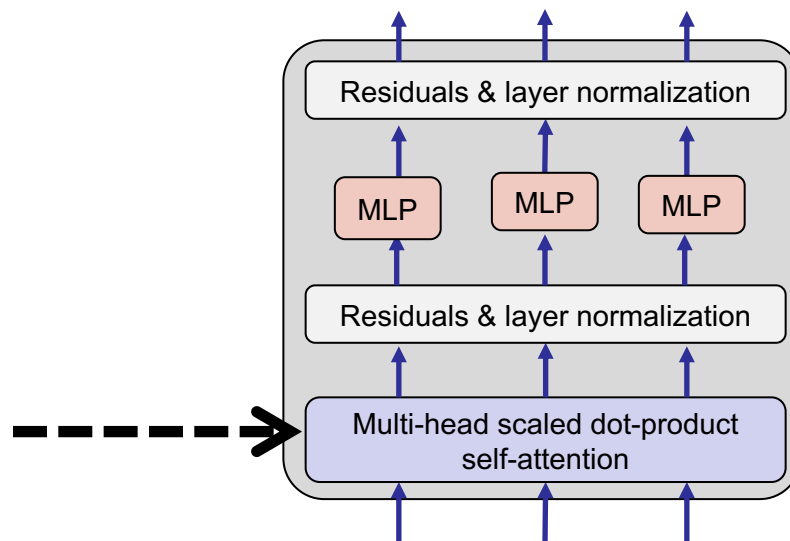


**Transformer Encoder**

# Transformer Encoder

Let's start from multi-head scaled dot-product self-attention:

1. **Scaled dot-product attention**
2. Multi-head attention
3. self-attention



**Transformer Encoder**

# Basic dot-product attention – recap

- Non-normalized attention scores:

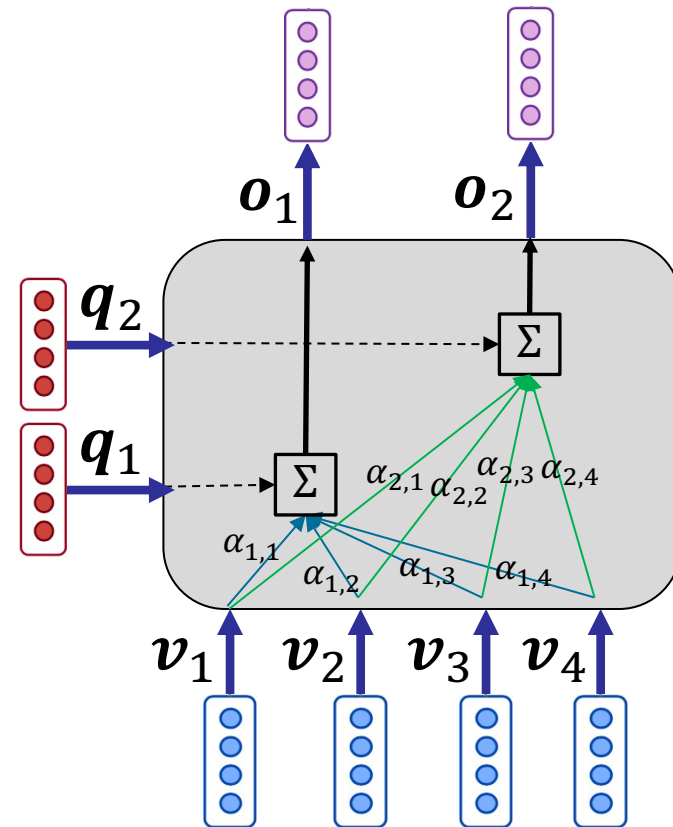
$$\tilde{\alpha}_{i,j} = f(\mathbf{q}_i, \mathbf{v}_j)$$

$$\tilde{\alpha}_{i,j} = \mathbf{q}_i \mathbf{v}_j^T$$

- In this case,  $d_q = d_v$
  - Attention network has no parameter to learn!
- Softmax over value vectors:

$$\alpha_i = \text{softmax}(\tilde{\alpha}_i)$$

- Output (weighted sum):  $\mathbf{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \mathbf{v}_j$



# Scaled dot-product attention

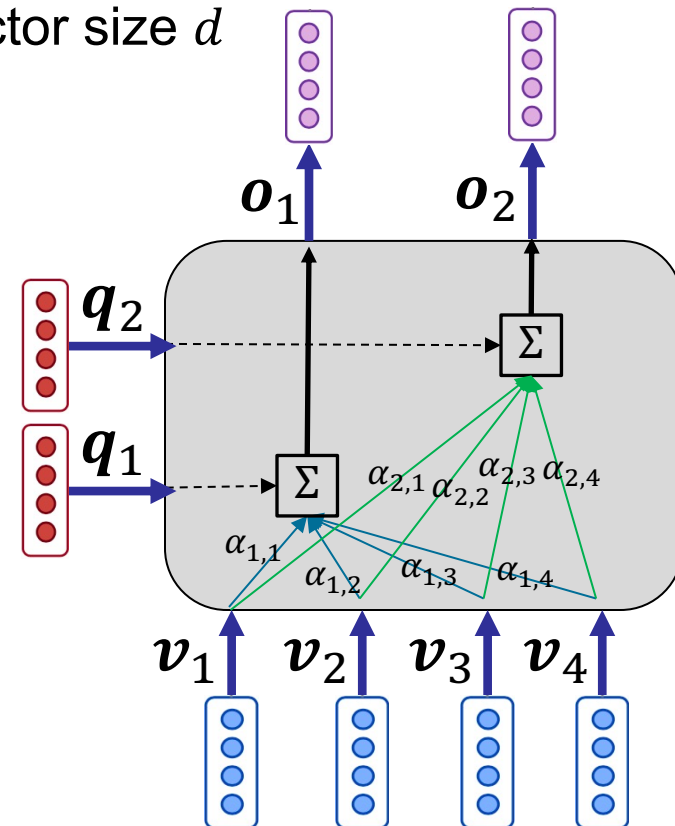
- Problem with basic dot-product attention:
  - As  $d$  gets large, the variance of  $\tilde{\alpha}_{i,j}$  increases ...
  - ... this makes softmax very peaked for some values of  $\tilde{\alpha}_i$  ...
  - ... and hence its gradient gets smaller
- One approach: normalize/scale  $\tilde{\alpha}_{i,j}$  by vector size  $d$

## Scaled dot-product attention

- Non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \frac{\mathbf{q}_i \mathbf{v}_j^T}{\sqrt{d}}$$

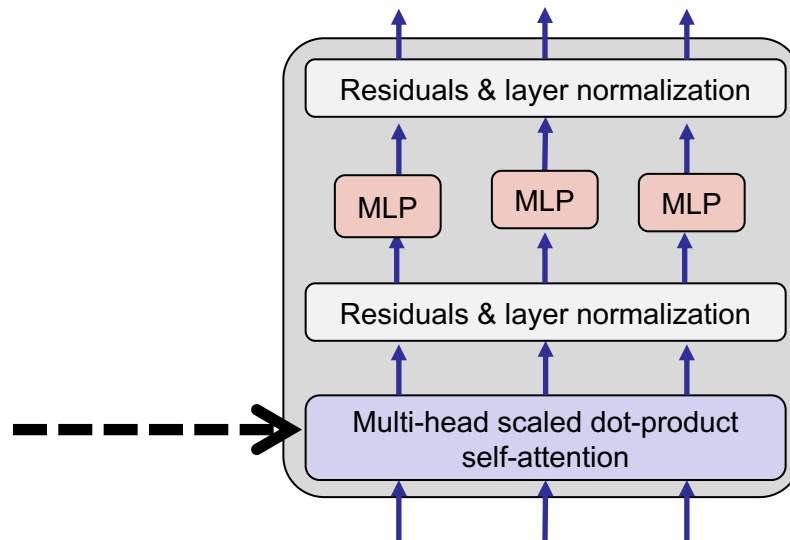
- Softmax over values:  $\alpha_i = \text{softmax}(\tilde{\alpha}_i)$
- Output:  $\mathbf{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \mathbf{v}_j$



# Transformer Encoder

Let's start from multi-head scaled dot-product self-attention:

1. Scaled dot-product attention
- 2. Multi-head attention**
3. self-attention



**Transformer Encoder**

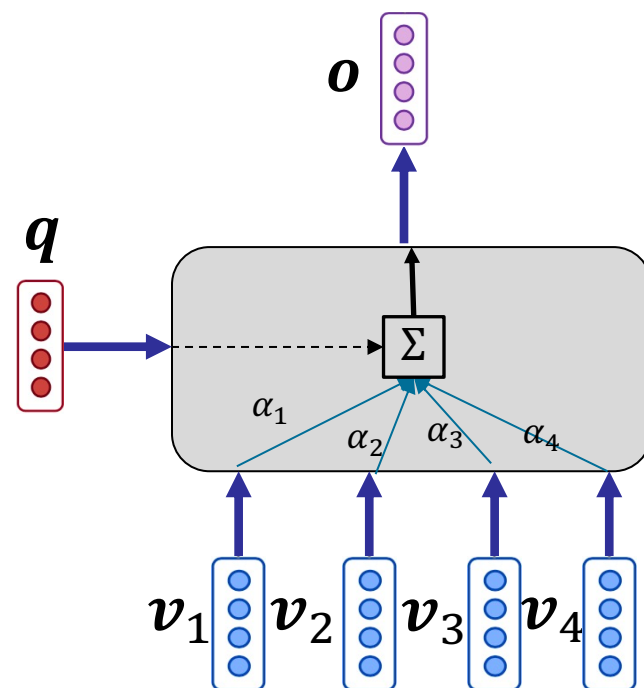


# Softmax bottleneck!

- Softmax is applied to non-normalized attention vectors
  - Recall: softmax makes the **maximum value** much higher than the other

$$\mathbf{z} = [1 \quad 2 \quad 5 \quad 6] \rightarrow \text{softmax}(\mathbf{z}) = [0.004 \quad 0.013 \quad 0.264 \quad 0.717]$$

- Common in language, a word may be related to several other words in a sequence, each through a **specific concept**
  - Like the relations of a verb to its subject and object
- However, normal (single-head) attention network aggregates all concepts in one set
- In this case, due to softmax, value vectors must compete for the attention of query vector  $\rightarrow$  **softmax bottleneck**



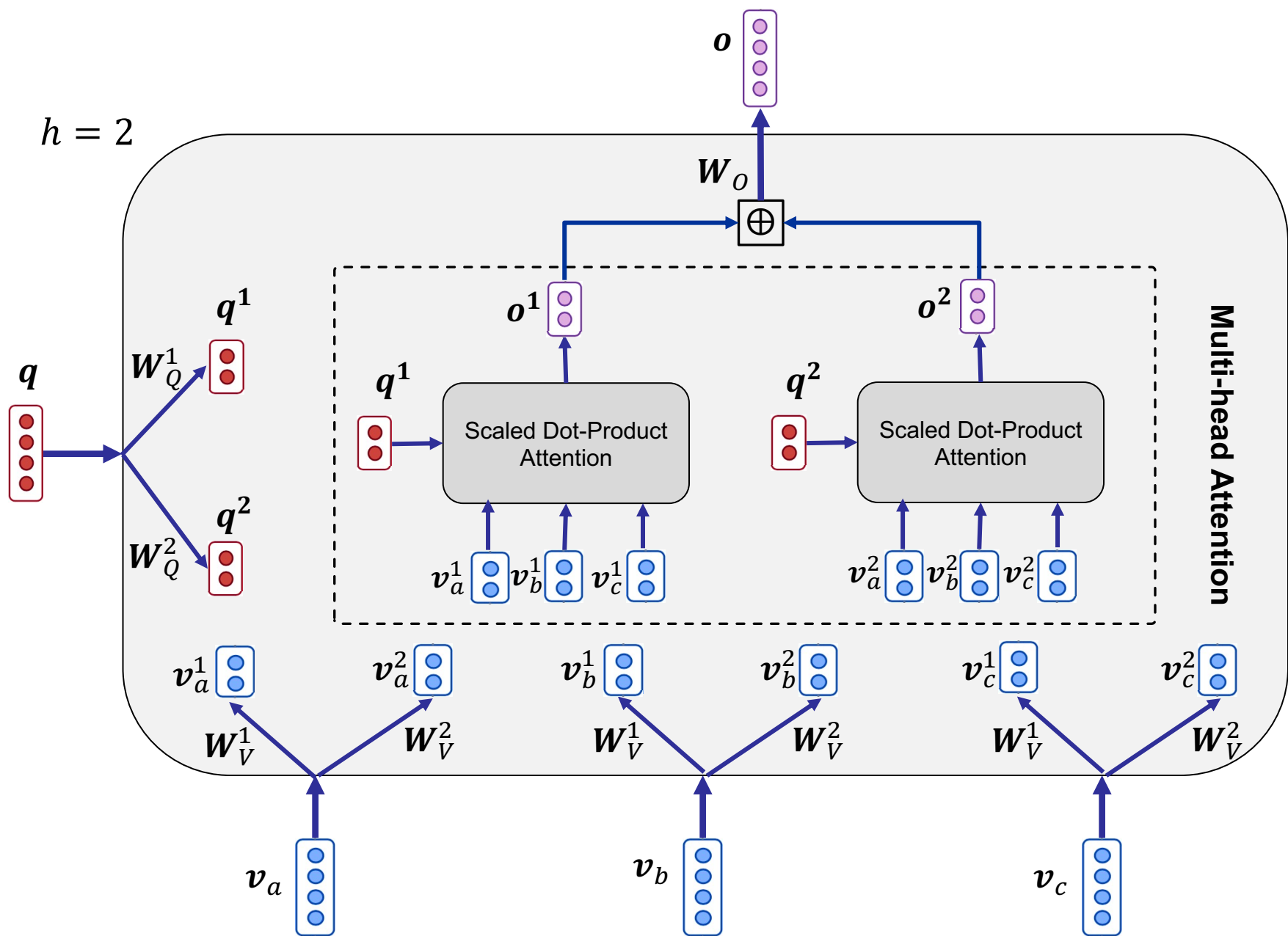
# Multi-head attention

- Multi-head attention approaches *softmax bottleneck* by calculating **multiple sets of attentions** between a query and values

## Multi-head attention:

1. Transfer each query/value vector to  $h$  query/value subspaces, each called a **head**
  2. In each subspace, apply a normal (single-head) attention network using the queries and values transferred to the subspace to achieve the output vectors of that head
  3. **Concatenate** the output vectors of all heads in respect to a query to achieve the **final output** of the query
- In multi-head attention, **each head** (and each subspace) can specialize on capturing a **specific kind** of relation

# Multi-head attention



# Multi-head attention – formulation

- Transfer every query  $q_i$  to  $h$  vectors, each with size  $d/h$ :

$$\boxed{\text{size: } d/h} \leftarrow q_i^1 = q_i \mathbf{W}_Q^1 \quad \dots \quad q_i^h = q_i \mathbf{W}_Q^h \rightarrow \boxed{\text{Matrix size: } d \times d/h}$$

- Transfer every value  $v_j$  to  $h$  vectors, each with size  $d/h$ :

$$\boxed{\text{size: } d/h} \leftarrow v_j^1 = v_j \mathbf{W}_V^1 \quad \dots \quad v_j^h = v_j \mathbf{W}_V^h \rightarrow \boxed{\text{Matrix size: } d \times d/h}$$

- Calculate outputs of subspaces corresponding to  $q_i$ :

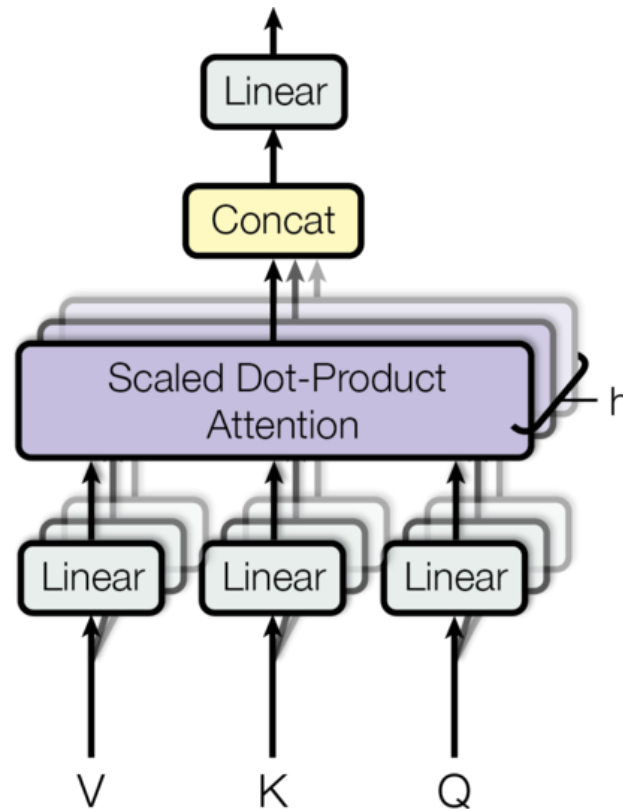
$$\boxed{\text{size: } d/h} \leftarrow o_i^1 = \text{ATT}(q_i^1, V^1) \quad \dots \quad o_i^h = \text{ATT}(q_i^h, V^h)$$

- Concatenate outputs of subspaces for  $q_i$  as its final output:

$$\boxed{\text{size: } d} \leftarrow o_i = \mathbf{W}_O [o_i^1; \dots; o_i^h]$$

Size:  $d \times d$   
This matrix linearly combines  
the dimensions of the  
concatenated vectors

# Multi-head attention – graphic in original paper

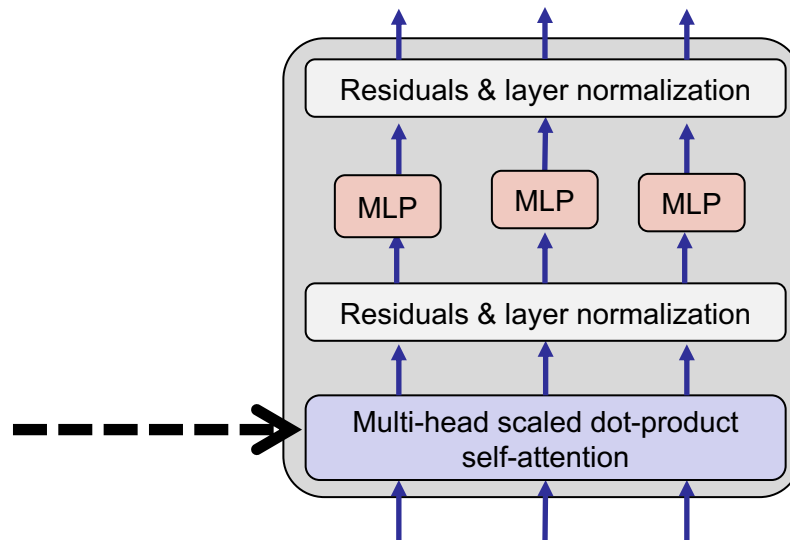


- Default number of heads in Transformers:  $h = 8$
- Recall: Attentions (and Transformers) in fact have three inputs (not two), namely queries, keys, and values.
  - Keys are used to calculate attentions
  - Values are used to produce outputs

# Transformer Encoder

Let's start from multi-head scaled dot-product self-attention:

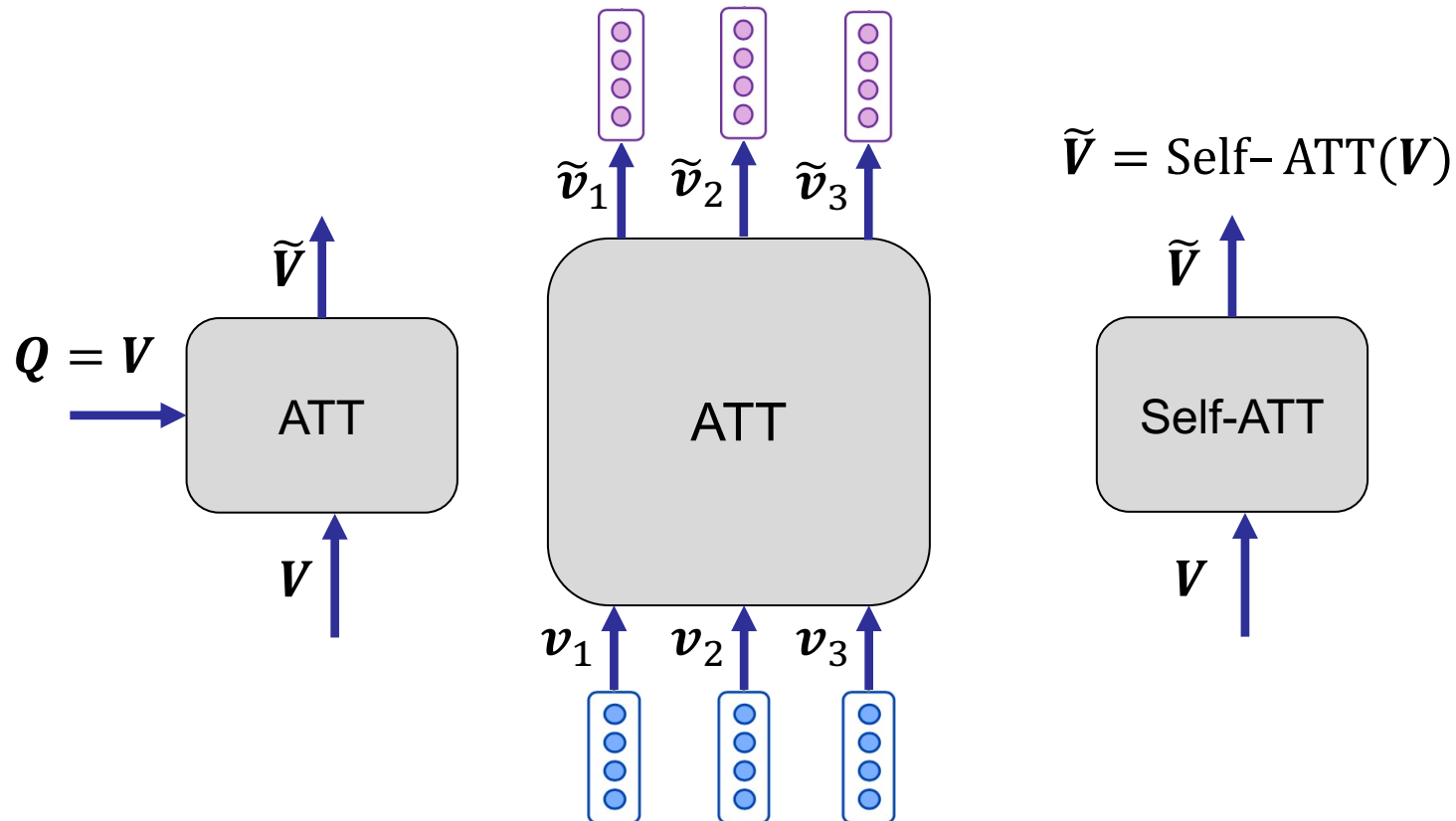
1. Scaled dot-product attention
2. Multi-head attention
3. **Self-attention**



**Transformer Encoder**

# Self-attention (recap)

- Values are the same as queries
- Each output vector is the **contextual embedding** of the corresponding input vector
  - $\tilde{v}_i$  is the contextual embedding of  $v_i$

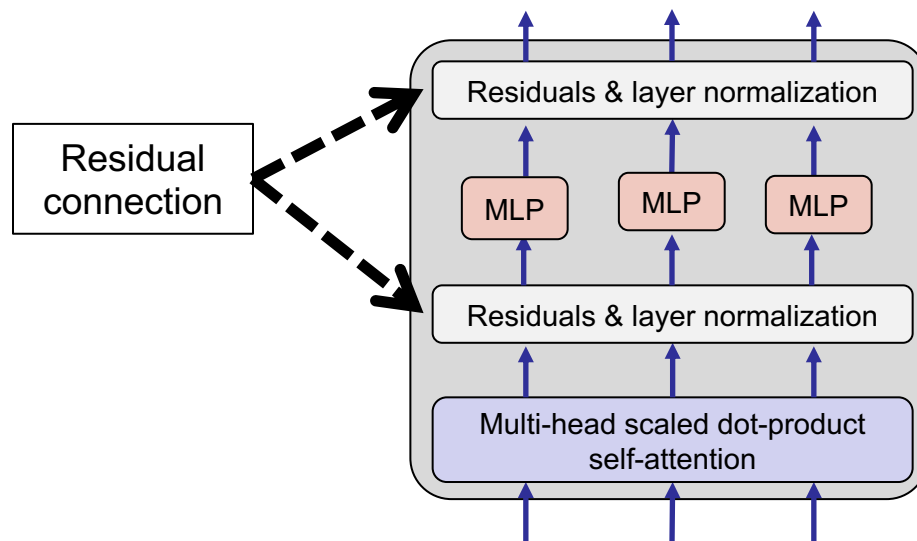


# Residuals

- Residual (short-cut) connection:

$$\text{output} = f(x) + x$$

- Learn in detail:
  - He, Kaiming; Zhang, Xiangyu; Ren, Shaoqing; Sun, Jian (2016). "Deep Residual Learning for Image Recognition" . In proc. of CVPR
  - Srivastava, Rupesh Kumar; Greff, Klaus; Schmidhuber, Jürgen (2015). "Highway Networks". <https://arxiv.org/pdf/1505.00387.pdf>

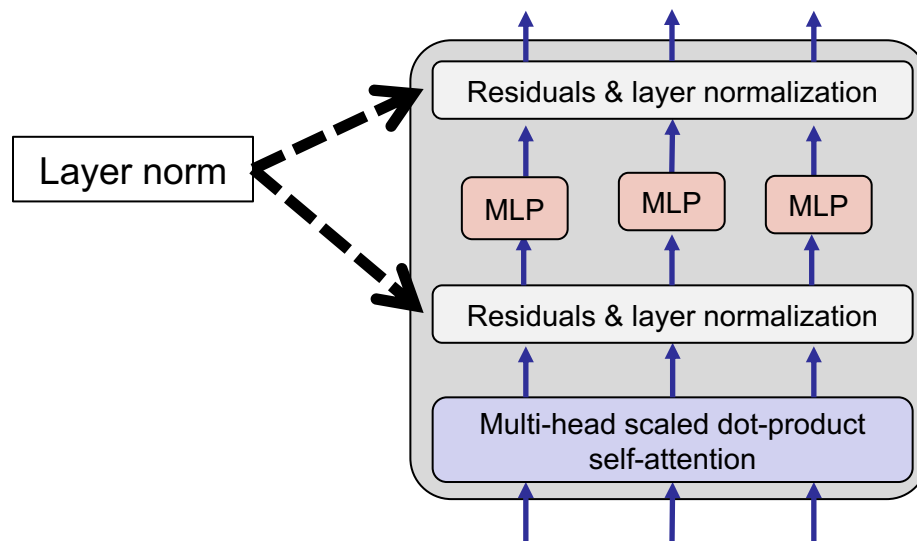


**Transformer Encoder**



# Layer normalization

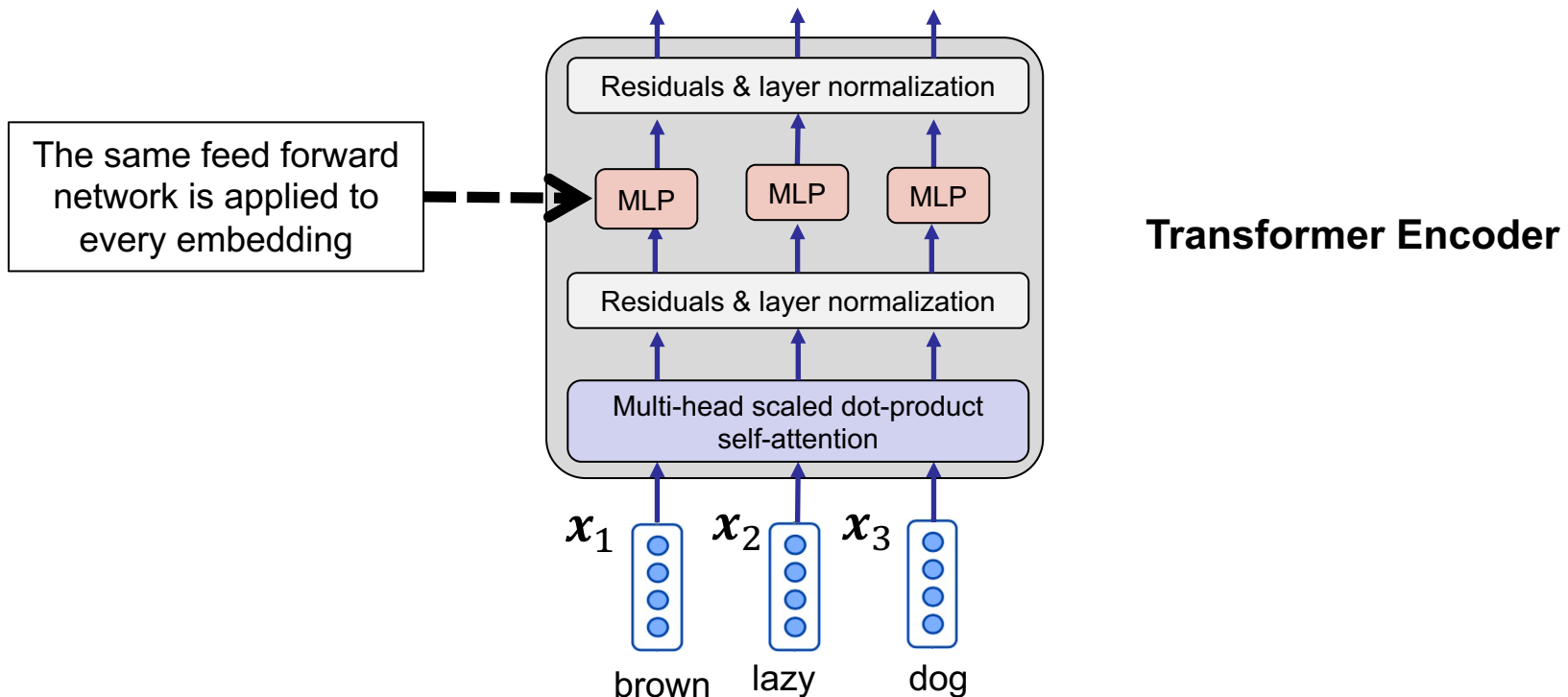
- Layer normalization changes the activations of each vector to have mean 0 and variance 1 ...
  - ... and learns two parameters per layer to shift the mean and variance



**Transformer Encoder**

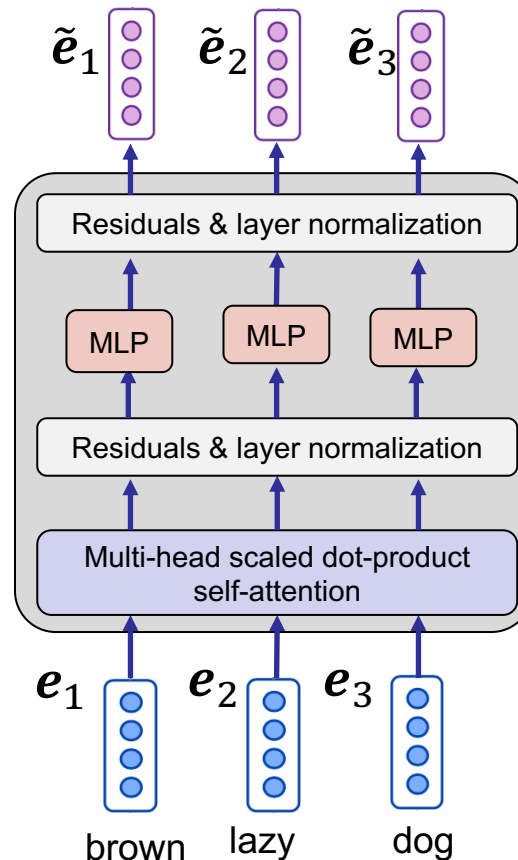
# Multi-layer perceptron on embedding

- A two-layer multi-layer perceptron (with ReLU) is applied to each output embedding
  - This layer provides the capacity for a non-linear transformation over each (contextualized) embedding



# Transformer Encoder – all together

- Transformer Encoder receive input embeddings and outputs the corresponding contextualized embeddings
  - Processing all inputs happen at the same time → non auto-regressive



# Transformer Encoder – summary

- A self-attention model using
  - multi-head scaled dot-product attention
  - followed by the same feed-forward layer applied to each embedding
  - all packed with residuals, layer norms, and dropouts

## Transformers as in attentions ...

- do not have **locality (position) bias**
  - A long-distance context has “equal opportunity”
- process all the input together with a **single computation** per each layer
  - Friendly with parallel computations in GPU

Learn more and study the PyTorch implementation: <http://nlp.seas.harvard.edu/2018/04/03/attention.html>

# Position embeddings

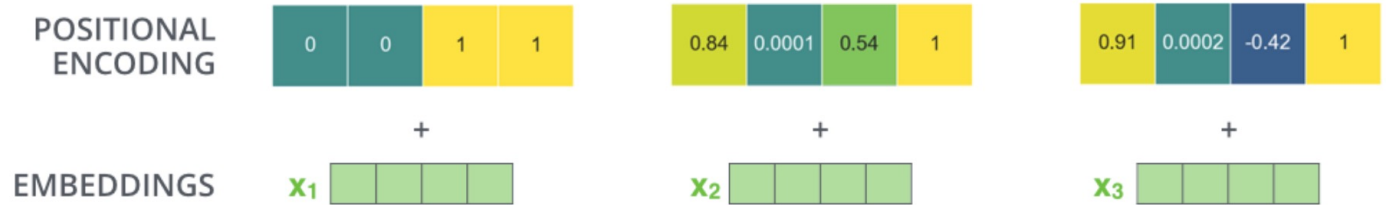
- Transformers are **agnostic** to the **position of tokens**
  - A context token in long-distance has the same effect as the one in short-distance (no *locality bias*)
- However, the positions of tokens in a sequence might be informative and important in some tasks

## **Position embeddings** – a common approach in Transformers:

- Create embeddings representing **positions** in a sequence, and **add** the values of the position embeddings to the token embeddings at corresponding positions
  - Position embedding is usually created using a sine/cosine function
    - It can also be learned end-to-end with the model parameters
  - Using position embeddings, the same token at different positions of a sequence will have different final representations

# Position embeddings – examples

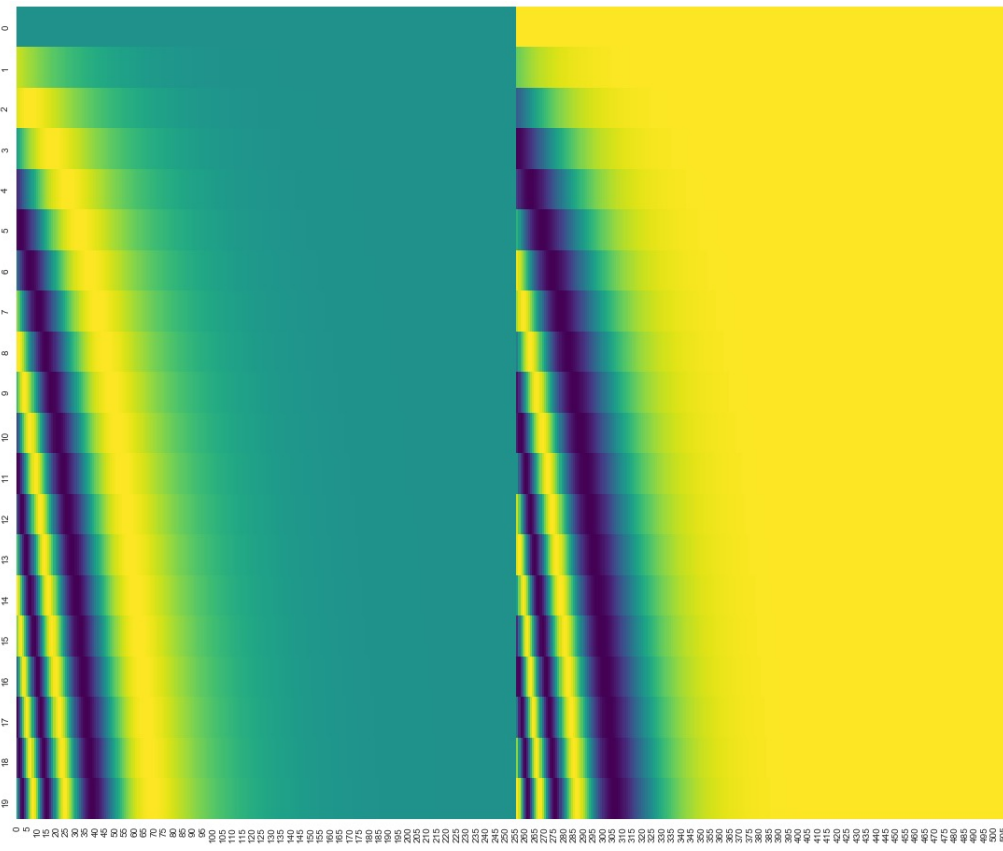
An example of embeddings with four dimensions:



Position embedding for location 0

Position embeddings

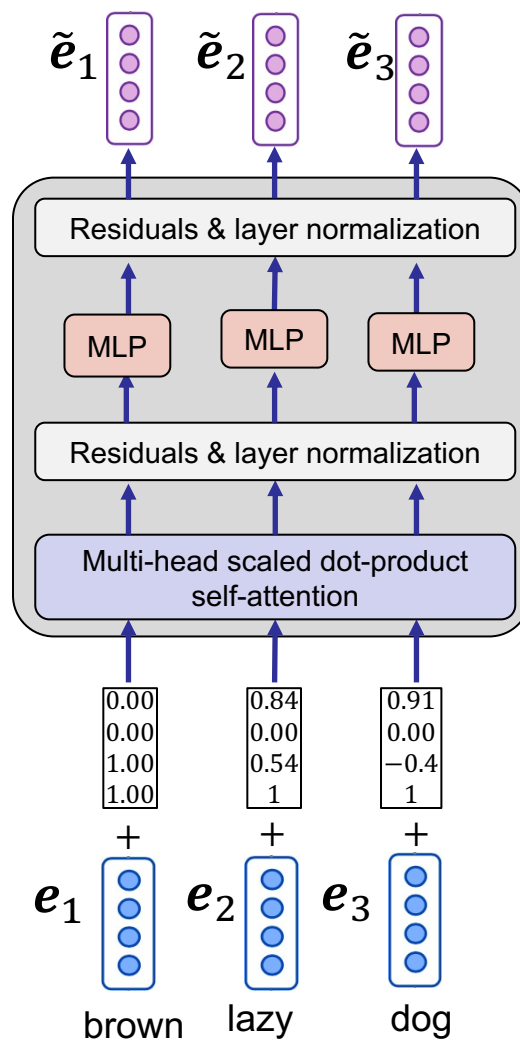
Position embedding for location 20



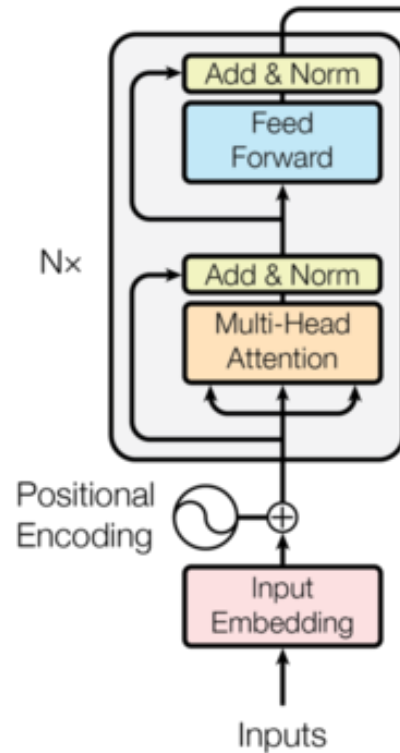
Values from -1 (dark) to +1 (light)

Dimensions (512)

# Transformer Encoder with position embedding



# Transformer Encoder with position embedding

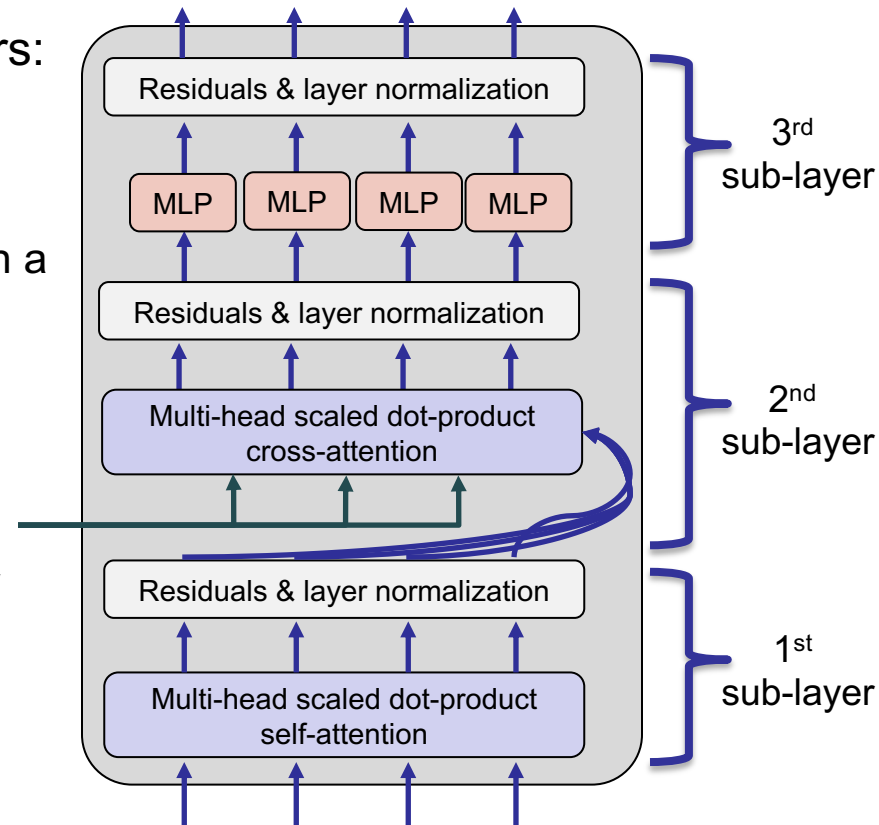




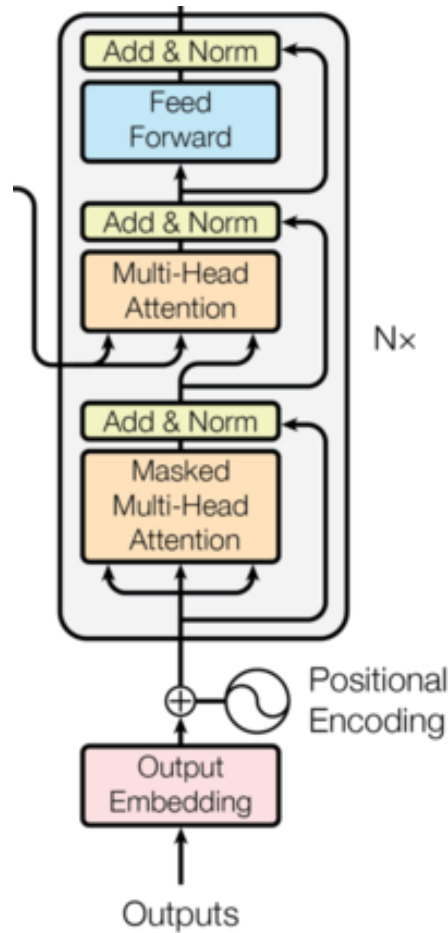
# Transformer Decoder

Transformer Decoder consists of three sub-layers:

- 1<sup>st</sup> : **Masked** multi-head self-attention
  - Exactly like Transformer Encoder but also with a *masking* functionality
- 2<sup>nd</sup> : Multi-head **cross attention**
  - Values are given from outside
    - Like from the outputs of a Transformer Encoder
  - Queries are the outputs of the 1<sup>st</sup> sub-layer
- 3<sup>rd</sup> : Position-wise multi-layer perceptron
  - Exactly like Transformer Encoder



# Transformer Decoder with position embedding



# Agenda

- Transformers
  - Transformer encoder
  - Transformer decoder
- **seq2seq with Transformers**

## Sequence-to-sequence modeling – recap

- Given the source sequence  $X = \{x^{(1)}, x^{(2)}, \dots, x^{(L)}\}, \dots$
- generate the target sequence  $Y = \{y^{(1)}, y^{(2)}, \dots, y^{(T)}\}$
- A seq2seq model estimates the **conditional probability**:

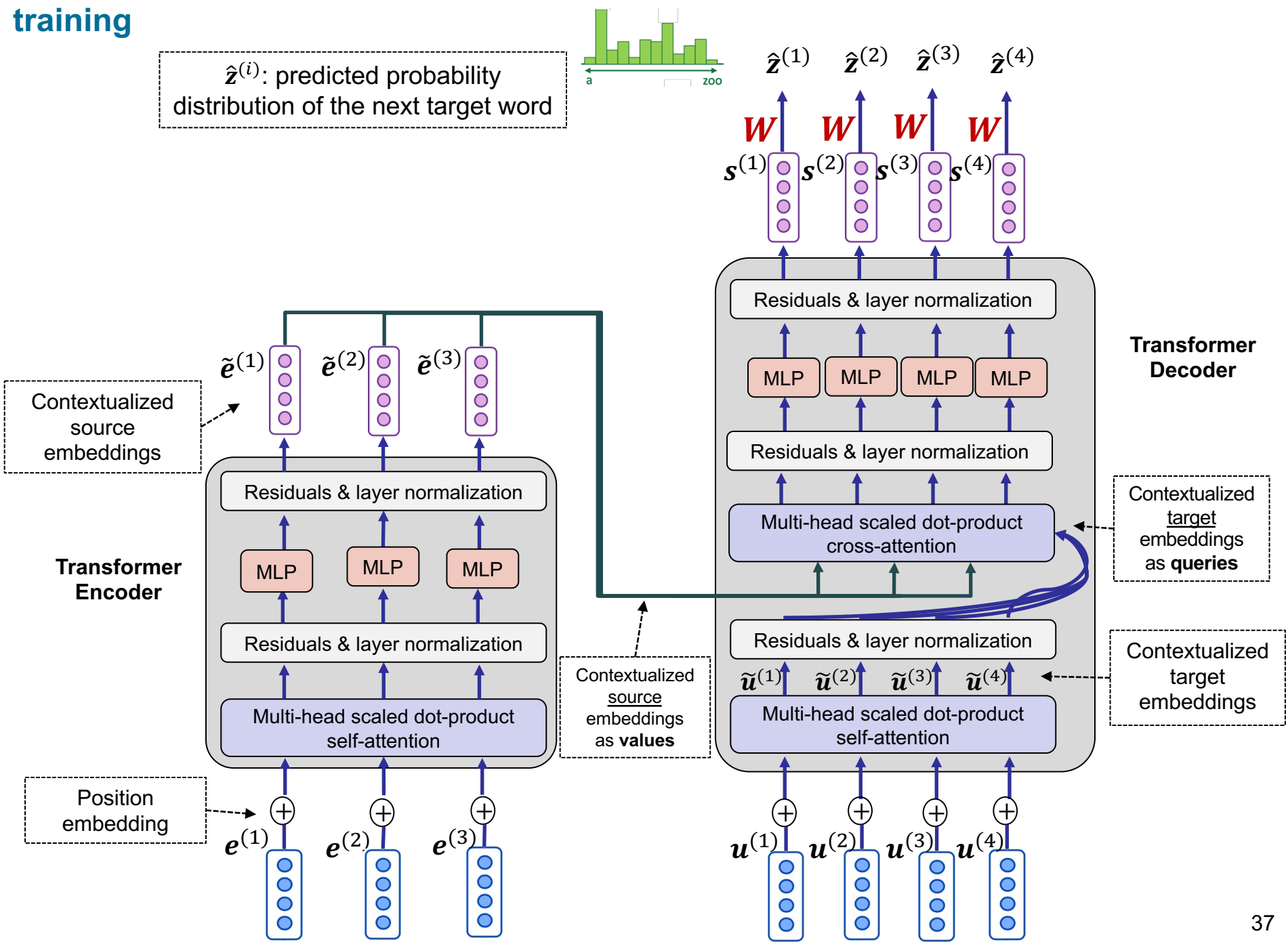
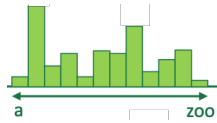
$$P(Y|X)$$

- and at inference time, it generates a new sequence  $Y^*$  such that:

$$Y^* = \operatorname{argmax}_Y P(Y|X)$$

# Seq2seq with Transformers – training

$\hat{\mathbf{z}}^{(i)}$ : predicted probability distribution of the next target word

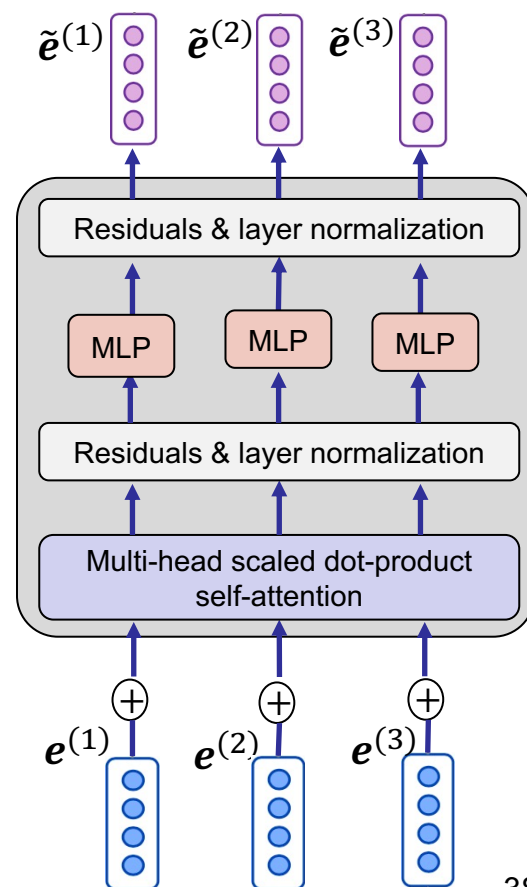


# Seq2seq with Transformers – training

- Two sets of vocabularies
  - $\mathbb{V}_e$  is the set of vocabularies for **source sequences**
  - $\mathbb{V}_d$  is the set of vocabularies for **target sequences**
- Source sequence  $X$  and target sequence  $Y$ 
  - Both are typically started/ended with  $\langle \text{bos} \rangle / \langle \text{eos} \rangle$

## Encoder

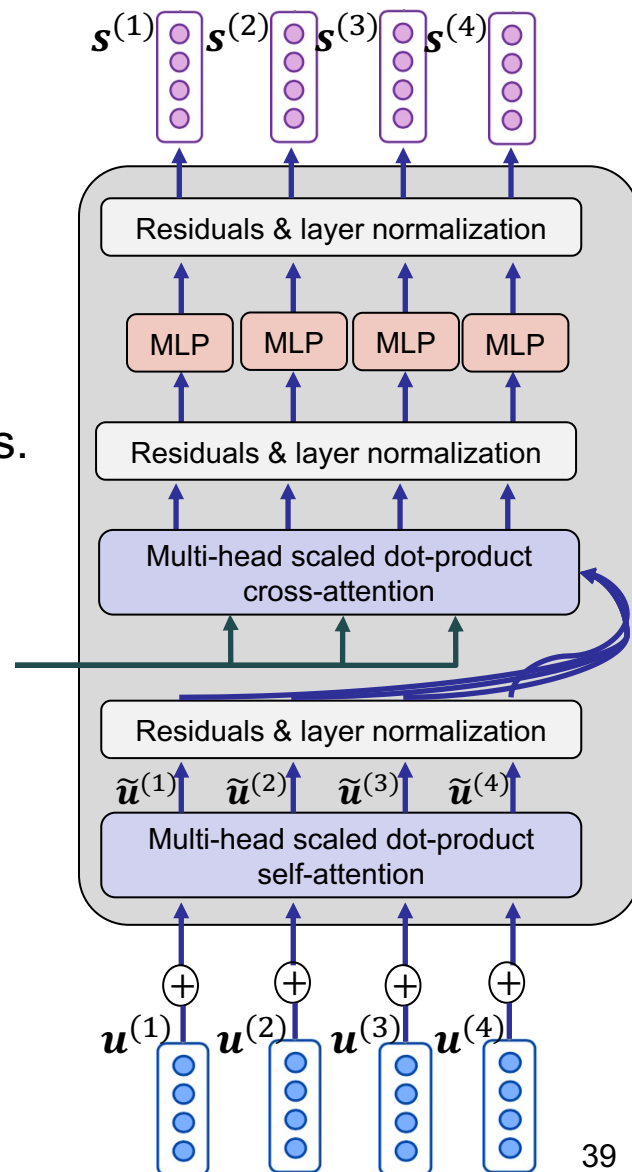
- Transformer encoder
  - passes source embeddings  $[e^{(1)}, \dots, e^{(L)}]$  and creates contextualized source embeddings:  $[\tilde{e}^{(1)}, \dots, \tilde{e}^{(L)}]$



# Seq2seq with Transformers – training

## Decoder

- Transformer Decoder **self-attention** layer
  - passes target embeddings  $[u^{(1)}, \dots, u^{(T)}]$  and creates contextualized target embeddings:  $[\tilde{u}^{(1)}, \dots, \tilde{u}^{(T)}]$
- Transformer Decoder **cross-attention** layer
  - applies attention with  $[\tilde{u}^{(1)}, \dots, \tilde{u}^{(T)}]$  as queries. and  $[\tilde{e}^{(1)}, \dots, \tilde{e}^{(L)}]$  as values (and keys)
- Transformer Decoder output
  - A set of vectors  $[s^{(1)}, \dots, s^{(T)}]$



**Incomplete version!**

# Seq2seq with Transformers – training

## Decoder (cont.)

- Decoder output prediction
  - uses  $[s^{(1)}, \dots, s^{(T)}]$  to calculate  $[\hat{z}^{(1)}, \dots, \hat{z}^{(T)}]$ , the vectors of the predicted probability distribution at the next position:

$$\hat{z}^{(t)} = \text{softmax}(\mathbf{W}[s^{(t)}] + \mathbf{b}) \in \mathbb{R}^{|\mathbb{V}_d|}$$

- Training loss for each position  $t$ 
  - NLL of the predicted probability of the next target word  $y^{(t+1)}$

$$\mathcal{L}^{(t)} = -\log \hat{z}_{y^{(t+1)}}^{(t)}$$

- Overall loss is the average of loss values over the target sequence:

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T \mathcal{L}^{(t)}$$



# Let's revisit the decoder!

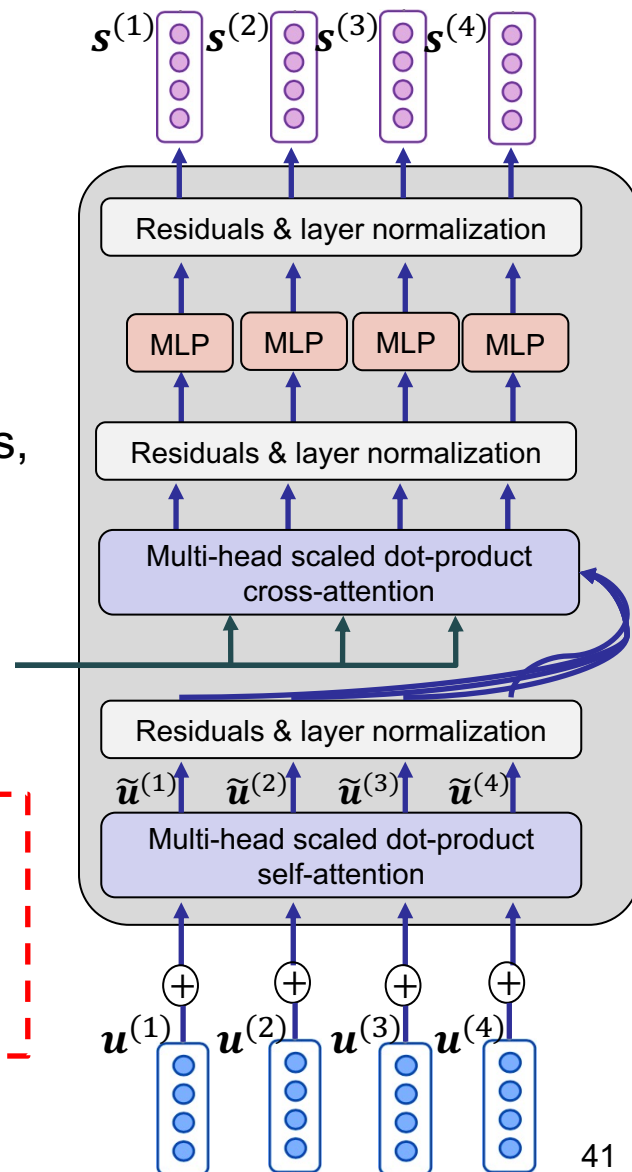
## Decoder

- Transformer Decoder **self-attention** layer
  - passes target embeddings  $[u^{(1)}, \dots, u^{(T)}]$  and creates contextualized target embeddings:  $[\tilde{u}^{(1)}, \dots, \tilde{u}^{(T)}]$
- Transformer Decoder **cross-attention** layer
  - applies attention with  $[\tilde{u}^{(1)}, \dots, \tilde{u}^{(T)}]$  as queries, and  $[\tilde{e}^{(1)}, \dots, \tilde{e}^{(L)}]$  as values (and keys)
- Transformer Decoder output
  - A set of vectors  $[s^{(1)}, \dots, s^{(T)}]$

**Problem:** in **self-attention** part, every token looks at all other tokens, namely the previous ones but also the next tokens!

- Every token has access to what it suppose to predict!

**Incomplete version!**



# Masking attentions

- In seq2seq with Transformers, we mask the attentions to every **future token** according to the self-attentions table of the Transformer Decoder

## Example

- **Non-normalized self-attention** scores of Transformer Decoder:

attends to ...  
other target embeddings

$\mathbf{u}^{(1)}$   $\mathbf{u}^{(2)}$   $\mathbf{u}^{(3)}$   $\mathbf{u}^{(4)}$

Each target embedding

$\mathbf{u}^{(1)}$	5	3	1	-4
$\mathbf{u}^{(2)}$	1	4	-2	3
$\mathbf{u}^{(3)}$	0	2	2	-3
$\mathbf{u}^{(4)}$	3	-1	1	4

## Non-normalized self-attention scores

$u^{(1)} \quad u^{(2)} \quad u^{(3)} \quad u^{(4)}$

$u^{(1)}$	5	3	1	-4
$u^{(2)}$	1	4	-2	3
$u^{(3)}$	0	-2	2	-3
$u^{(4)}$	3	-1	1	4

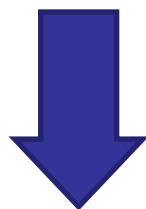
## attentions masks

$u^{(1)} \quad u^{(2)} \quad u^{(3)} \quad u^{(4)}$

$u^{(1)}$	1	0	0	0
$u^{(2)}$	1	1	0	0
$u^{(3)}$	1	1	1	0
$u^{(4)}$	1	1	1	1

Applying masks to attention scores

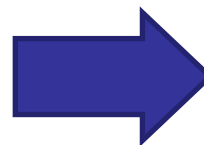
- adds  $-\infty$  for every mask value 0
- adds 0 for every mask value 1



$u^{(1)} \quad u^{(2)} \quad u^{(3)} \quad u^{(4)}$

$u^{(1)}$	5	$-\infty$	$-\infty$	$-\infty$
$u^{(2)}$	1	4	$-\infty$	$-\infty$
$u^{(3)}$	0	-2	2	$-\infty$
$u^{(4)}$	3	-1	1	4

softmax



## Final self-attention scores

$u^{(1)} \quad u^{(2)} \quad u^{(3)} \quad u^{(4)}$

$u^{(1)}$	1.00	0.00	0.00	0.00
$u^{(2)}$	0.04	0.96	0.00	0.00
$u^{(3)}$	0.11	0.01	0.86	0.00
$u^{(4)}$	0.25	0.01	0.34	0.70

👉 In Transformers, there are  $h$  times of such attention matrices. The same masking is applied to each of them.

# Seq2seq with Transformers – training

## Decoder

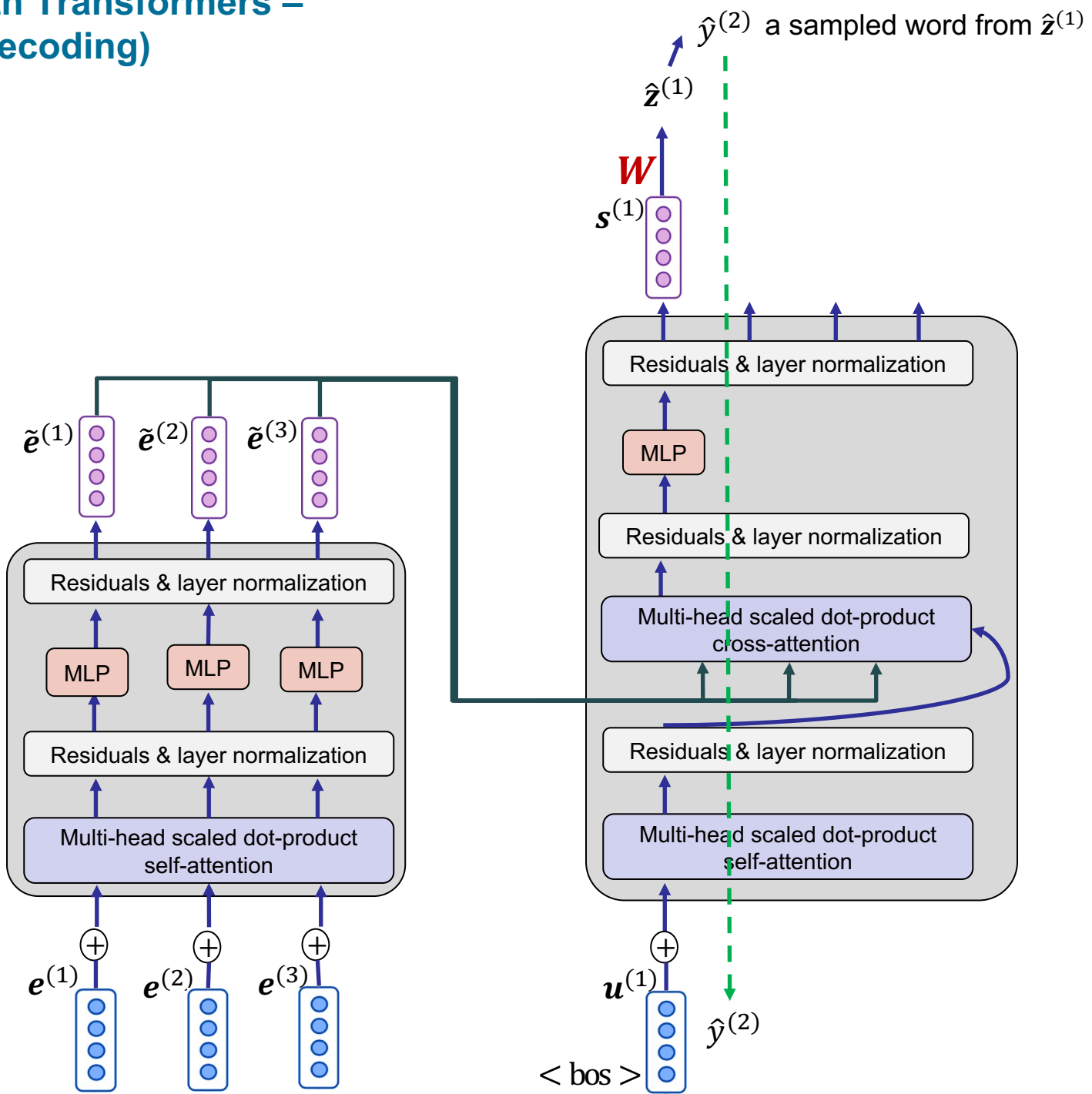
- Transformer Decoder self-attention layer
  - passes target embeddings  $[\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(T)}]$  and creates contextualized target embeddings:  $[\tilde{\mathbf{u}}^{(1)}, \dots, \tilde{\mathbf{u}}^{(T)}]$  **while masking future tokens**
- Transformer Decoder cross-attention layer
  - applies attention with  $[\tilde{\mathbf{u}}^{(1)}, \dots, \tilde{\mathbf{u}}^{(T)}]$  as queries and  $[\tilde{\mathbf{e}}^{(1)}, \dots, \tilde{\mathbf{e}}^{(L)}]$  as values (and keys)
- Transformer Decoder output
  - A set of vectors  $[\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(T)}]$

**Complete version!**

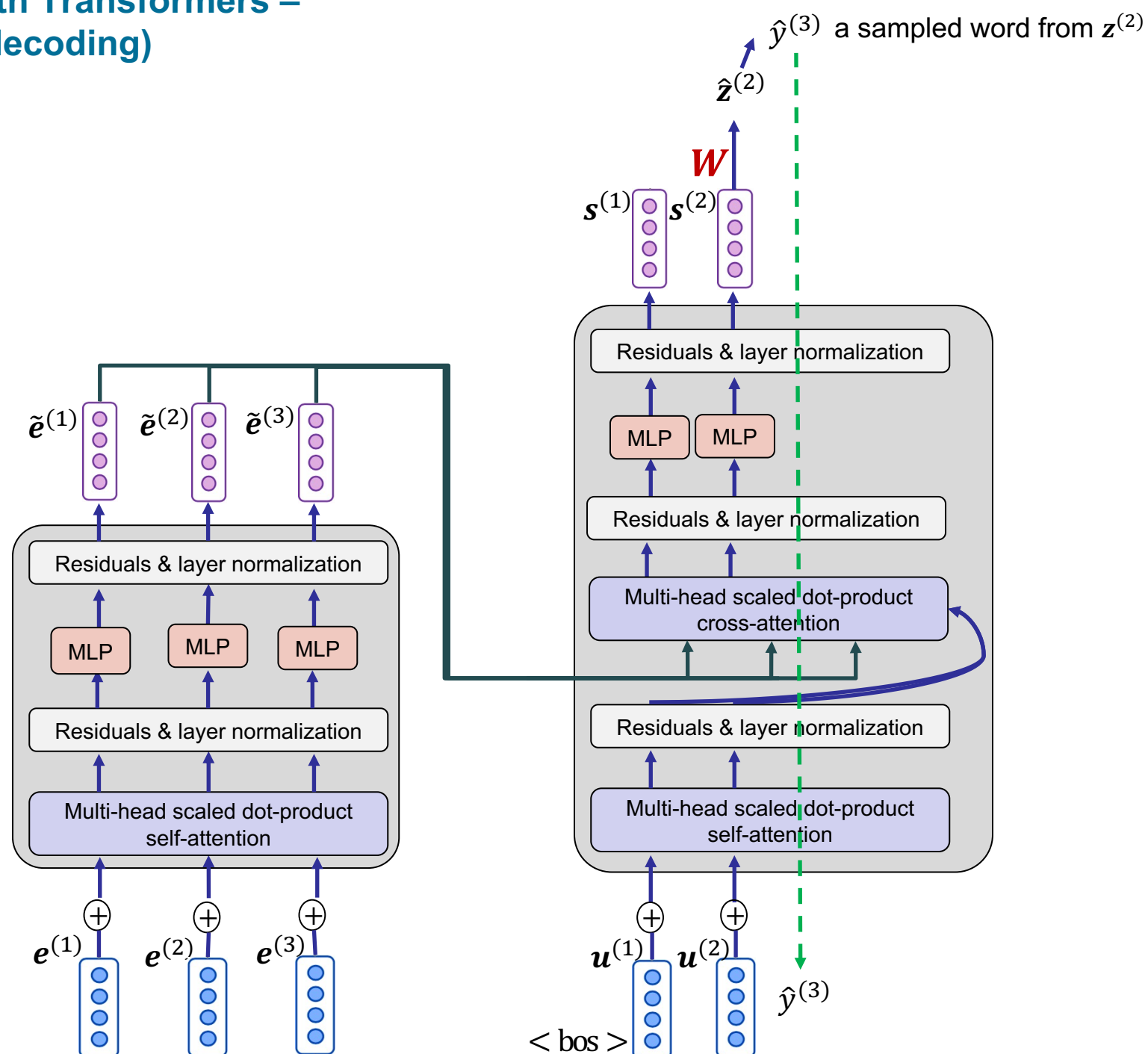
# Inference (decoding)

- During inference, as in training, the encoding of input sequence is done with a single computation (non-autoregressive)
- However, as in seq2seq with RNNs, decoding of seq2seq with Transformers is done in autoregressive fashion (one token after each other):
  - Pass the 1<sup>st</sup> target token (< bos >), generate the 2<sup>nd</sup> token
  - Pass the 1<sup>st</sup> token + the 2<sup>nd</sup> generated target tokens, generate the 3<sup>rd</sup> token
  - Pass the 1<sup>st</sup> token + the 2<sup>nd</sup> and 3<sup>rd</sup> generated target tokens, generate the 4<sup>th</sup> token
  - ...

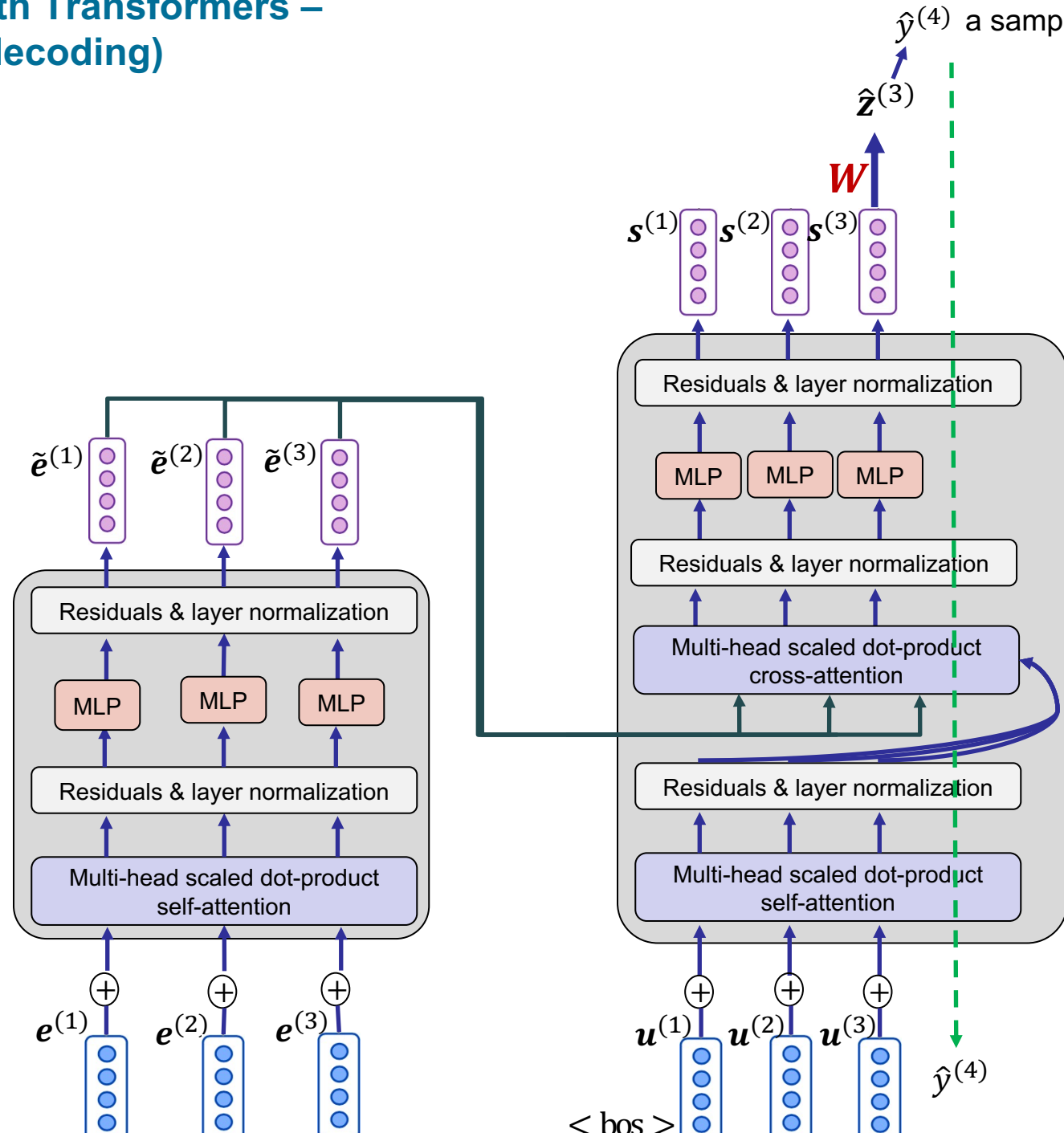
# Seq2seq with Transformers – inference (decoding)



# Seq2seq with Transformers – inference (decoding)



# Seq2seq with Transformers – inference (decoding)





# Seq2seq with Transformers – code

- Each Transformer encoder/decoder is a block. You can stack them several times and make the network deep!

```
CLASS torch.nn.TransformerEncoder(encoder_layer, num_layers, norm=None) [SOURCE]
```

```
CLASS torch.nn.TransformerEncoderLayer(d_model, nhead,  
dim_feedforward=2048, dropout=0.1, activation='relu') [SOURCE]
```

```
CLASS torch.nn.TransformerDecoder(decoder_layer, num_layers, norm=None) [SOURCE]
```

```
CLASS torch.nn.TransformerDecoderLayer(d_model, nhead,  
dim_feedforward=2048, dropout=0.1, activation='relu') [SOURCE]
```

```
forward(tgt, memory, tgt_mask=None, memory_mask=None,  
tgt_key_padding_mask=None, memory_key_padding_mask=None) [SOURCE]
```

