

# 344.063/163 KV Special Topic: Natural Language Processing with Deep Learning Transformers



Navid Rekab-Saz

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

**Institute of Computational Perception**

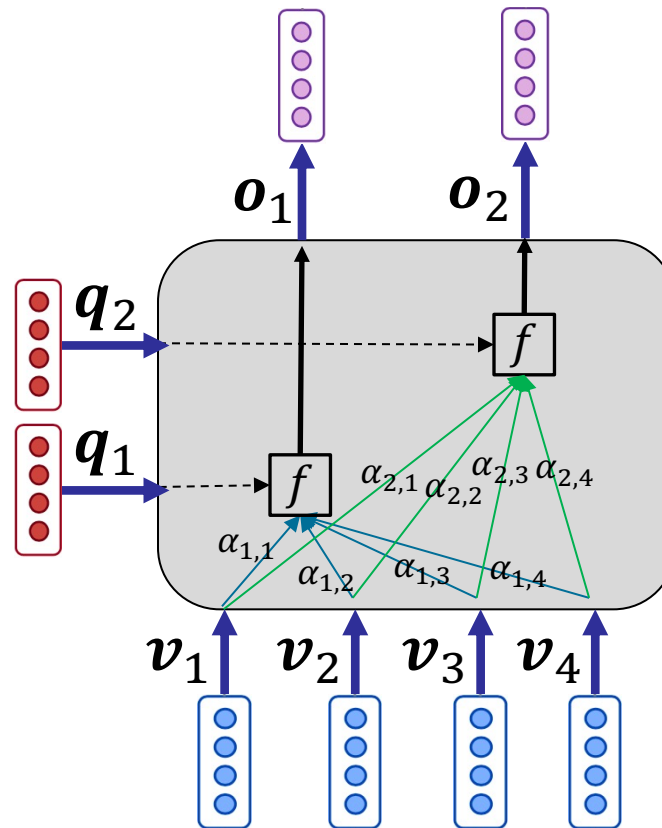
# Agenda

- 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 of query  $q_i$  on value  $v_j$

$\alpha_i$  is the vector of attentions of query  $q_i$  on value vectors  $V$

$\alpha_i$  is a probability distribution

$f$  is the attention function

## Attention Networks formulation – recap

- Given the query vector  $\mathbf{q}_i$ , an attention network assigns **attention**  $\alpha_{i,j}$  to each value vector  $\mathbf{v}_j$  using **attention function**  $f$ :

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

where  $\alpha_i$  forms a **probability distribution** over vector values:

$$\sum_{j=1}^{|V|} \alpha_{i,j} = 1$$

- The output regarding each query is the **weighted sum** of the value vectors using attentions as weights:

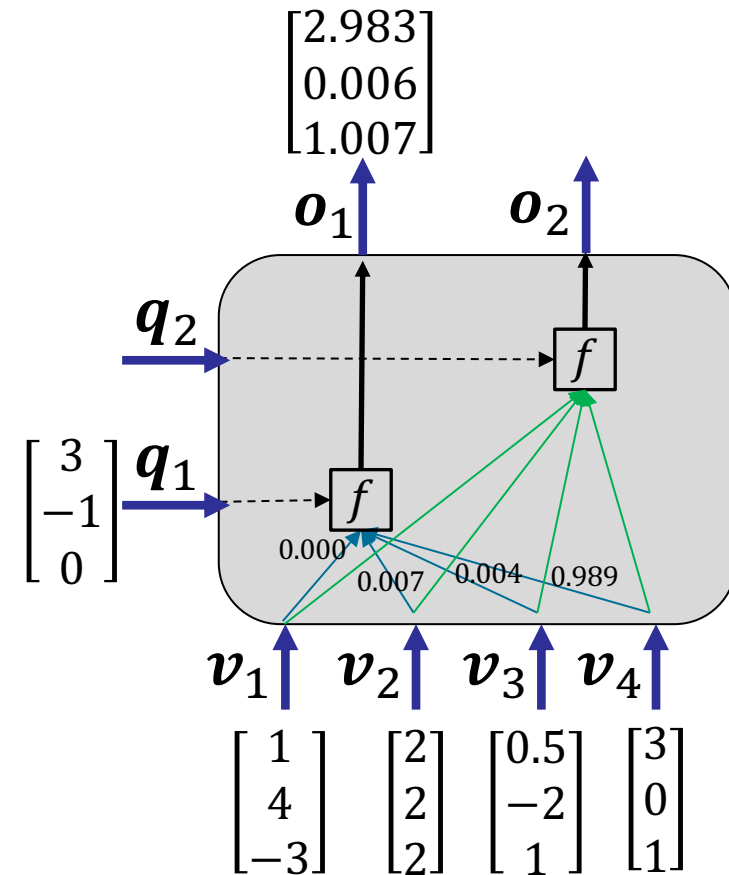
$$\mathbf{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \mathbf{v}_j$$

## Example – recap

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

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

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

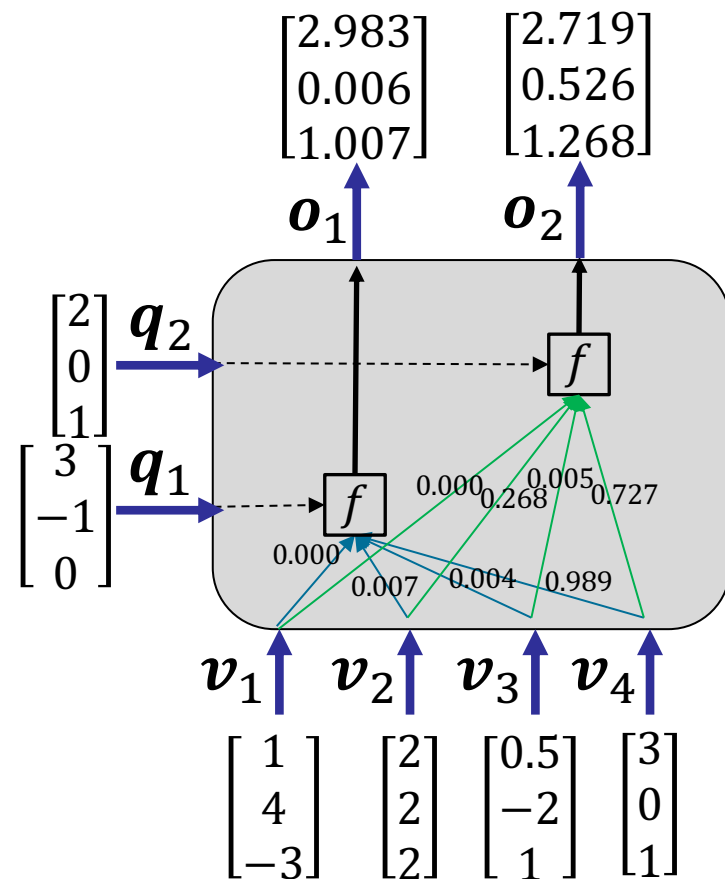


## Example – recap

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

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

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

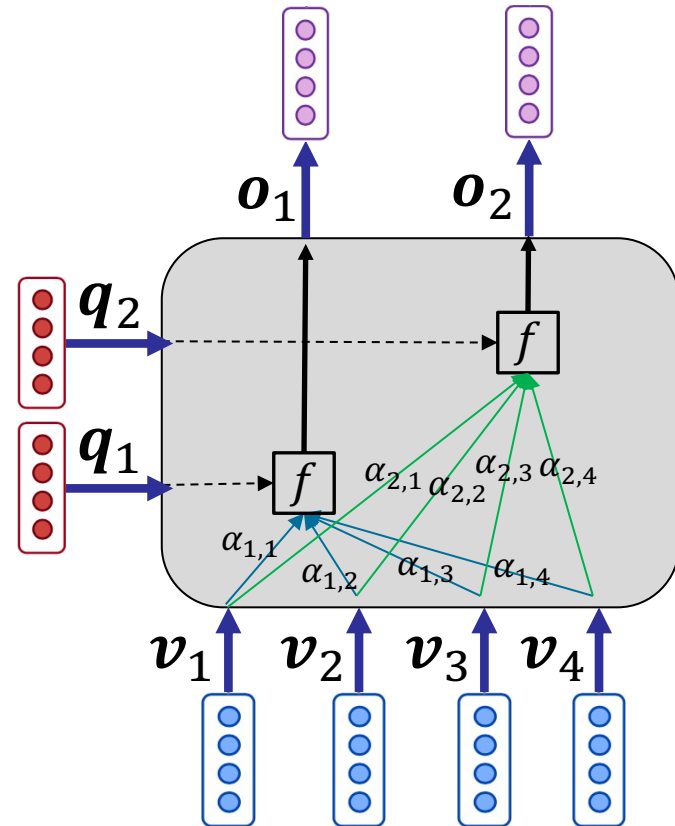


# Attention table

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

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

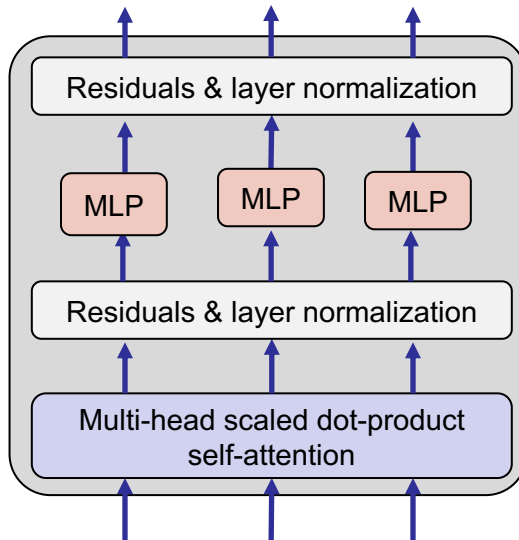




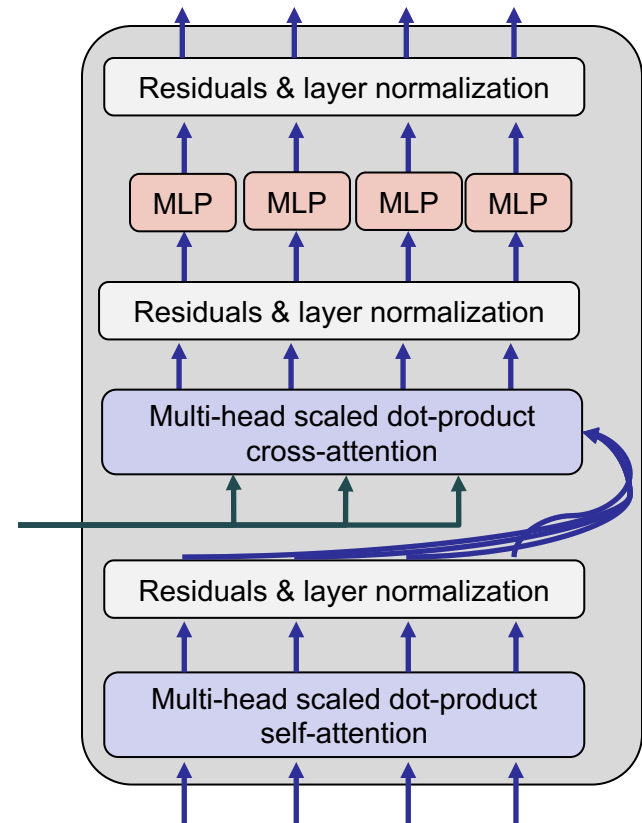
# Transformers

- 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

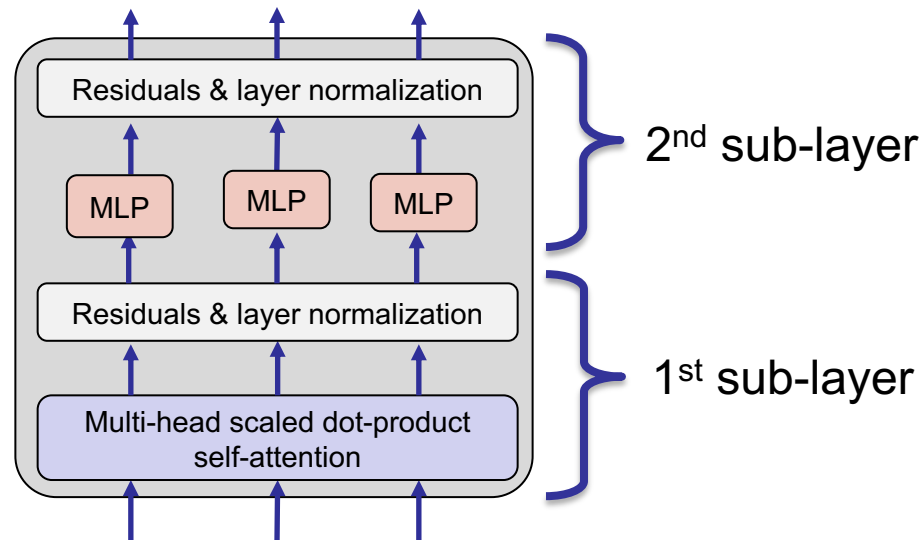


## Transformer decoder



# Transformer encoder

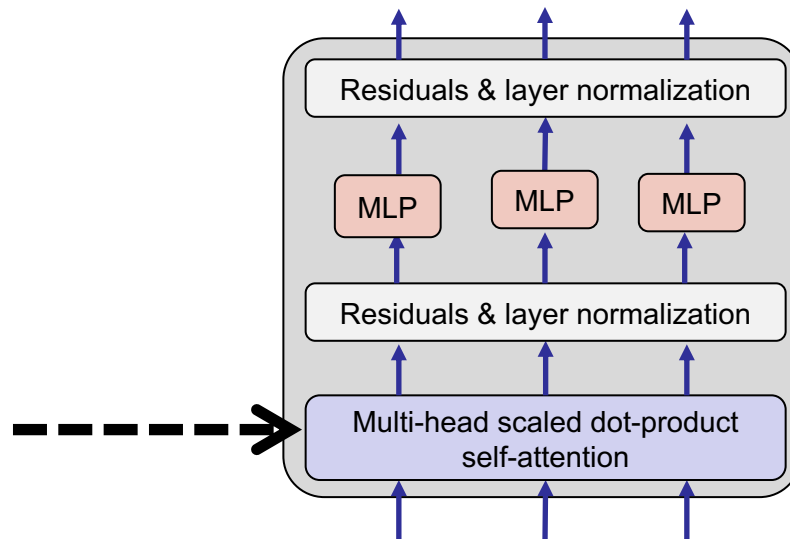
- Transformer encoder consists of two sub-layers:
  - 1<sup>st</sup> : Multi-head scaled dot-product self-attention
  - 2<sup>nd</sup> : Position-wise multi-layer perceptron (feed forward)
- Each sub-layer is followed by residual networks and layer normalization
  - drop-outs are applied after each computation



# 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



# Basic dot-product attention – recap

- First, non-normalized attention scores:

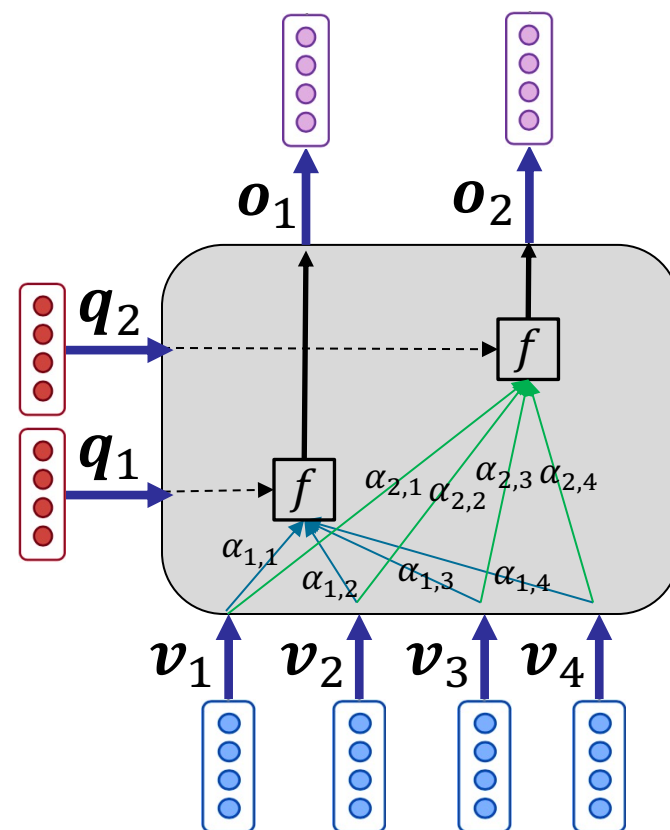
$$\tilde{\alpha}_{i,j} = \mathbf{q}_i \mathbf{v}_j^T$$

- In this variant  $d_q = d_v$
- There is no parameter to learn!

- Then, softmax over values:

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

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



# Scaled dot-product attention

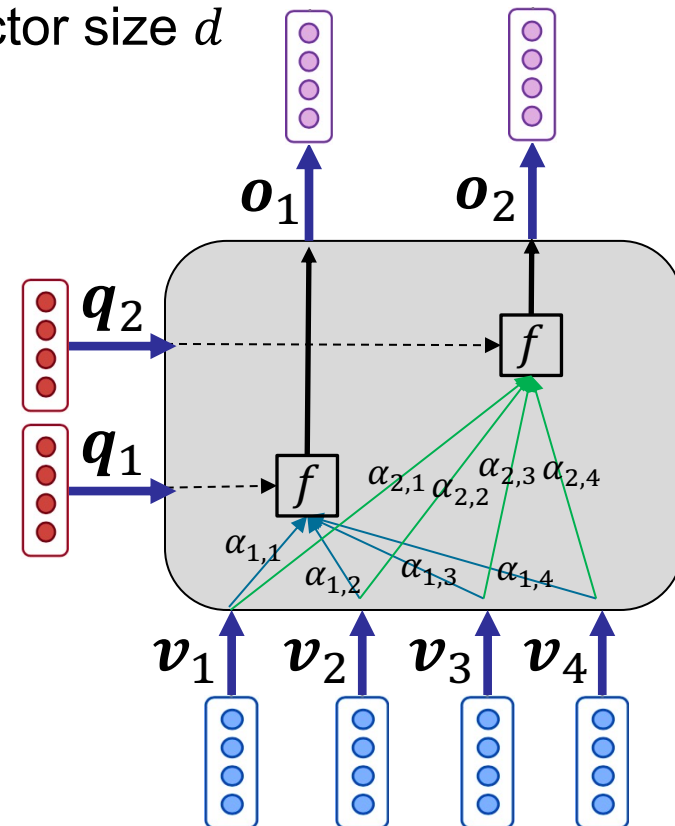
- Problem with basic dot-product attention:
  - As  $d$  gets large, the variance of  $\tilde{\alpha}_{i,j}$  increases ...
  - ... this makes softmax very peaked for some values of  $\tilde{\alpha}_i$  ...
  - ... and hence its gradient gets smaller
- One approach: normalize/scale  $\tilde{\alpha}_{i,j}$  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:  $\mathbf{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \mathbf{v}_j$

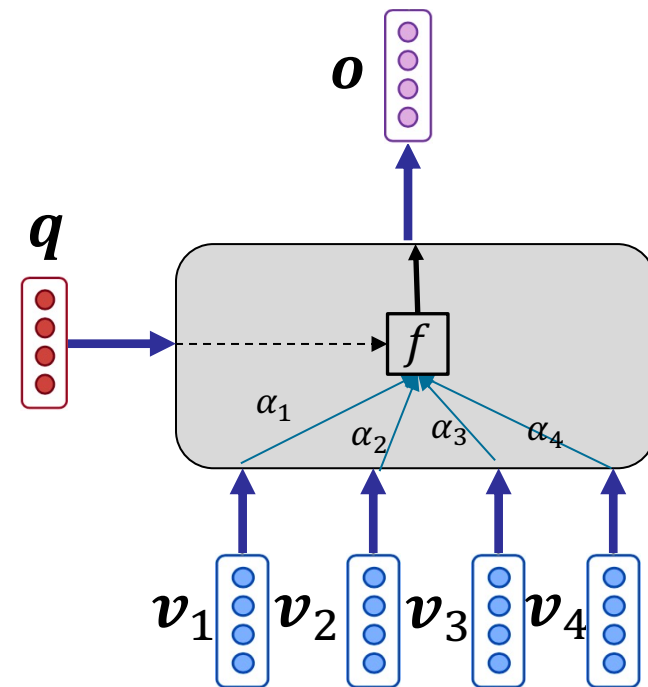


# From single-head to multi-head attention

- 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 to its object
- However, in a single-head attention network, all concepts are aggregated in one attention set
- Due to softmax, value vectors must compete for the attention of query vector  
→ **softmax bottleneck**



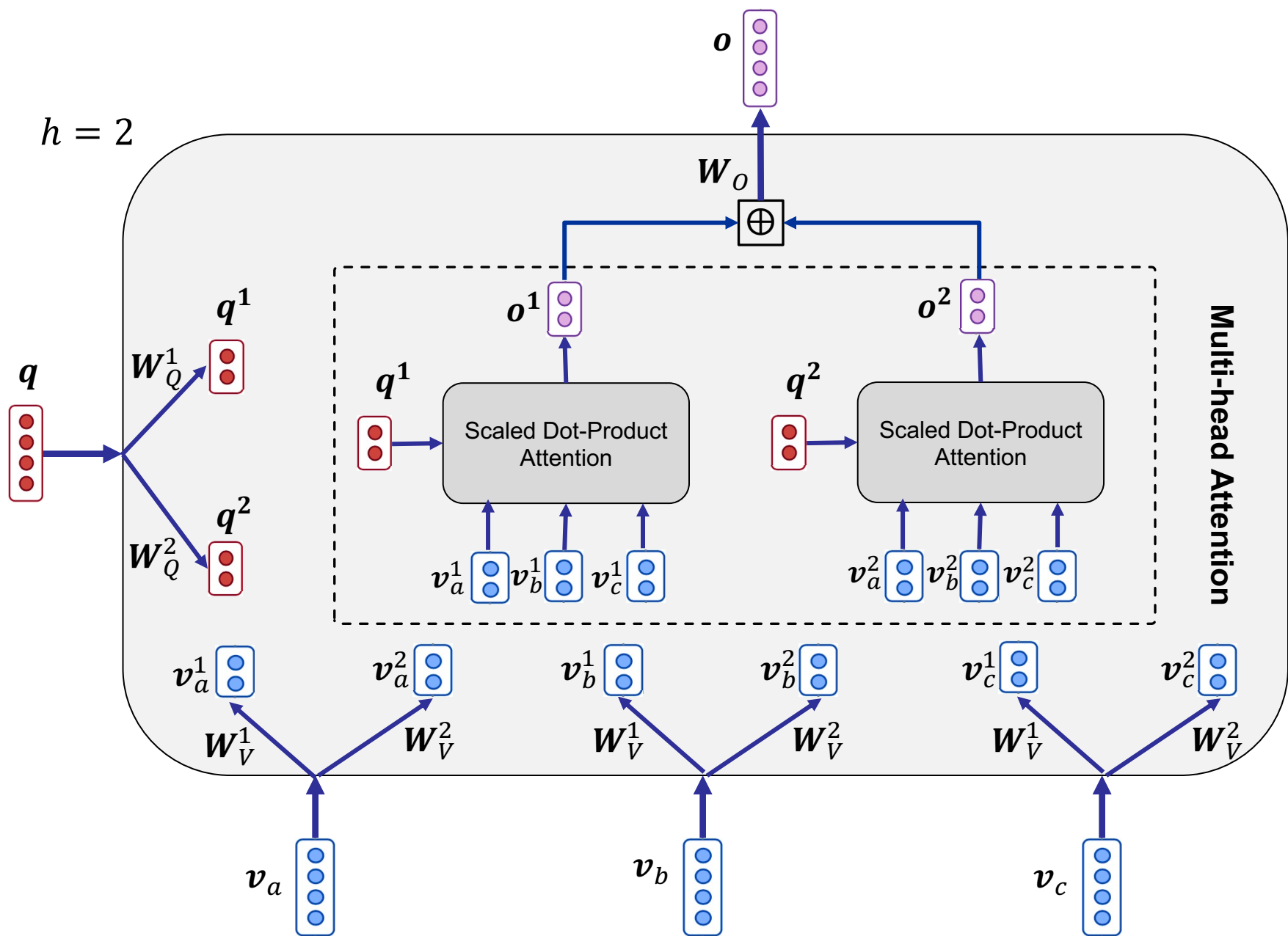
# Multi-head attention

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

## Multi-head attention:

1. Transfer each query/value vector to  $h$  subspaces (**heads**)
  2. In each subspace, apply a single-head attention network using the queries and values transferred to the subspace to achieve the output vectors of that subspace
  3. **Concatenate** the output vectors of all subspaces 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 relations

# 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  $q_i$ :

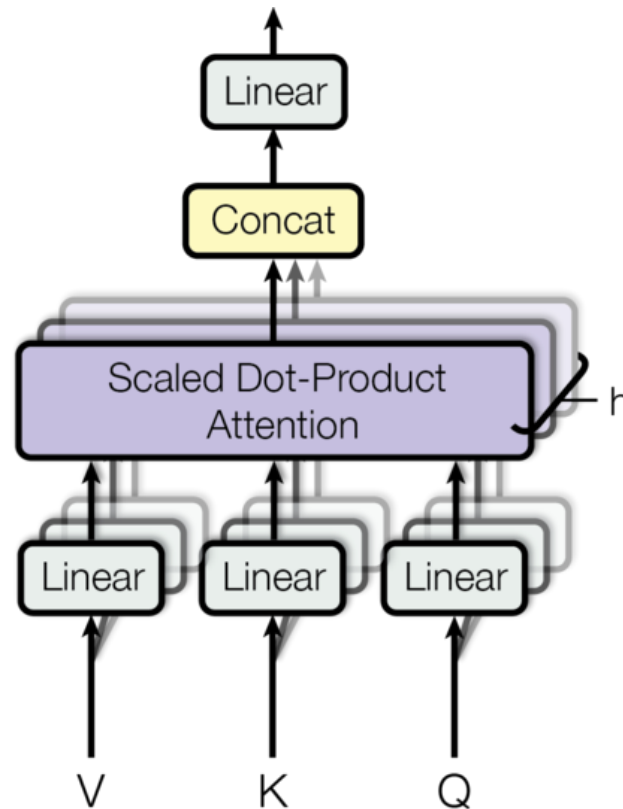
$$\boxed{\text{size: } d/h} \leftarrow o_i^1 = \text{ATT}(q_i^1, V^1) \quad \dots \quad o_i^h = \text{ATT}(q_i^h, V^h)$$

- 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

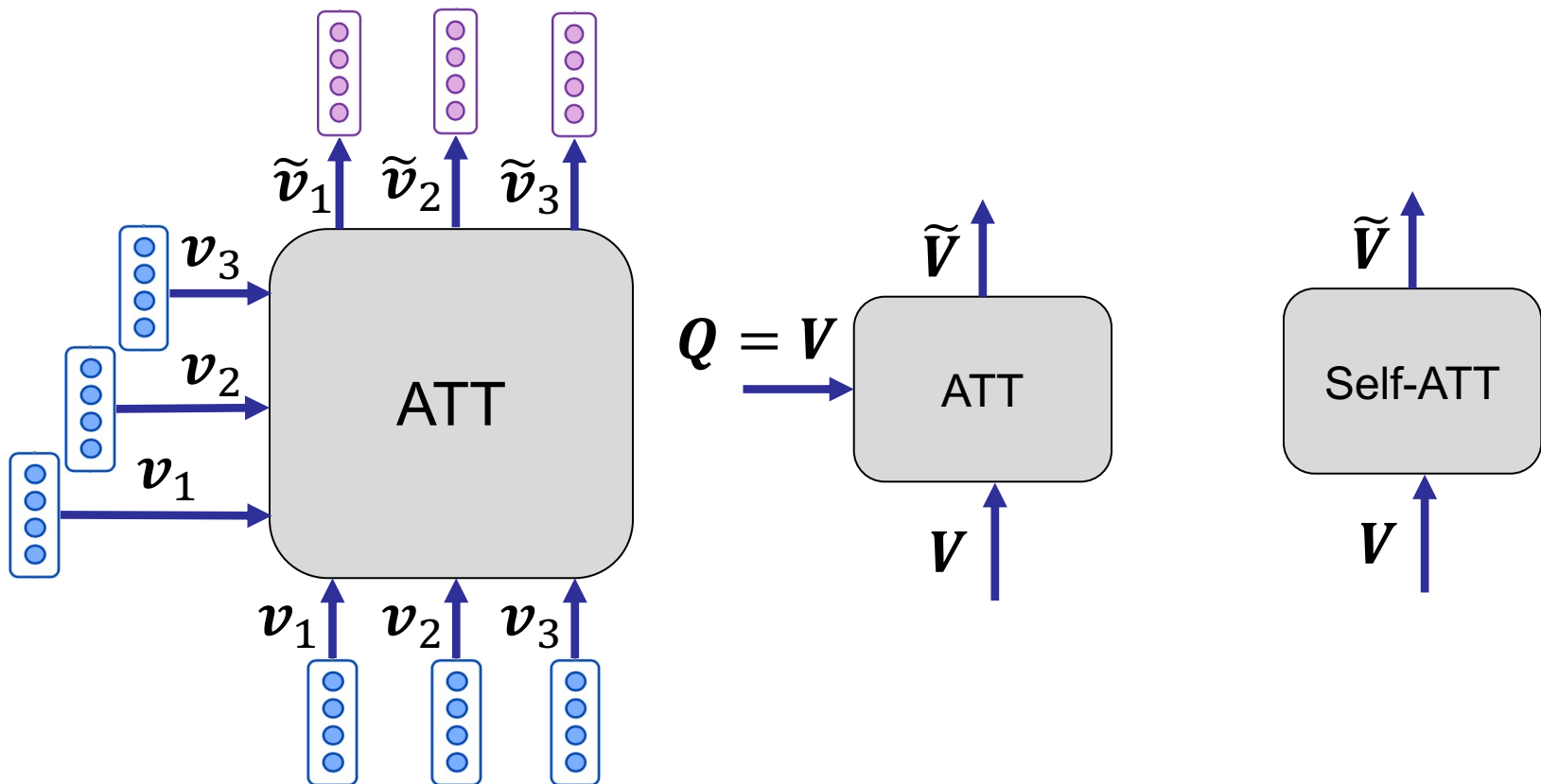
# Multi-head attention – graphic in original paper



- 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

# Self-attention – recap

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



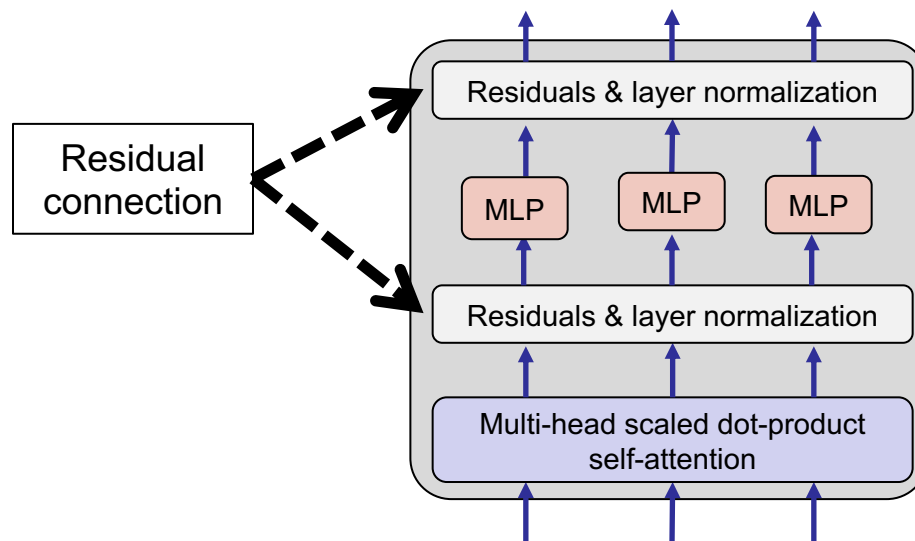
# Residuals

- Residual (short-cut) connection:

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

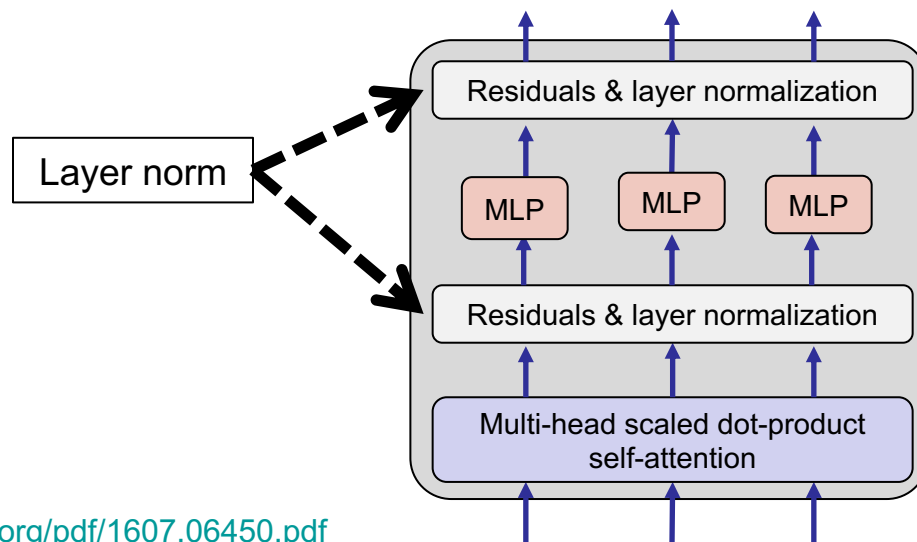
- Learn in detail:

- 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>



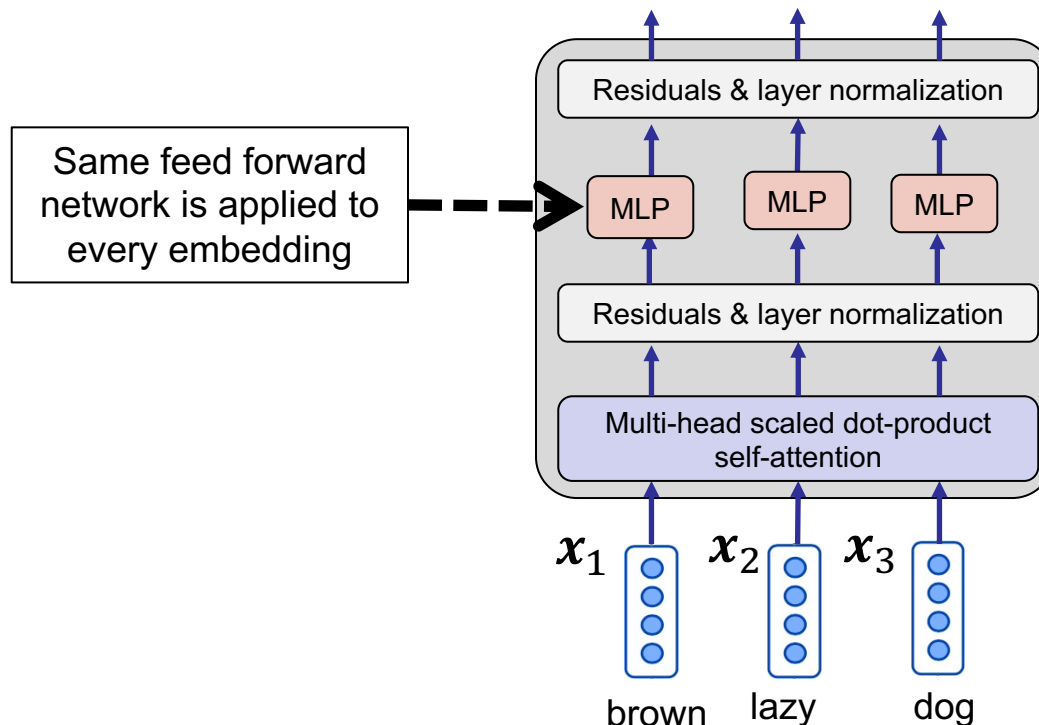
# Layer normalization

- 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



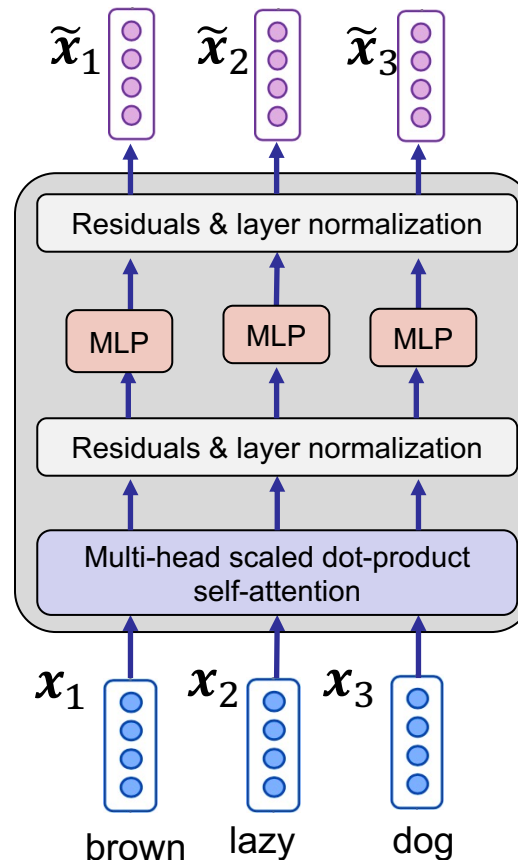
# Multi-layer perceptron on embedding

- In Transformers, a two-layer multi-layer perceptron (feed forward) neural network with ReLU is applied to each output embedding
  - The feed forward layer provides capacity for non-linear transformations 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 here more and study its PyTorch implementation:  
<http://nlp.seas.harvard.edu/2018/04/03/attention.html>



# Position embeddings

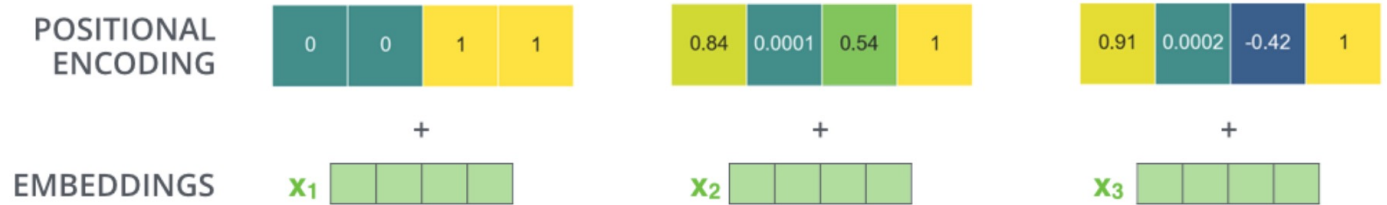
- Transformers are **agnostic** to the **position of tokens**
  - A context token in long-distance has the same effect as one in short-distance (no locality bias)
- However depending on the task, the positions of tokens in a sequence can be informative

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

- Create embeddings representing **positions** in a sequence, and **add** the corresponding position embedding to the token embedding
  - 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 word at different locations in a sequence will have different representations

# Position embeddings – examples

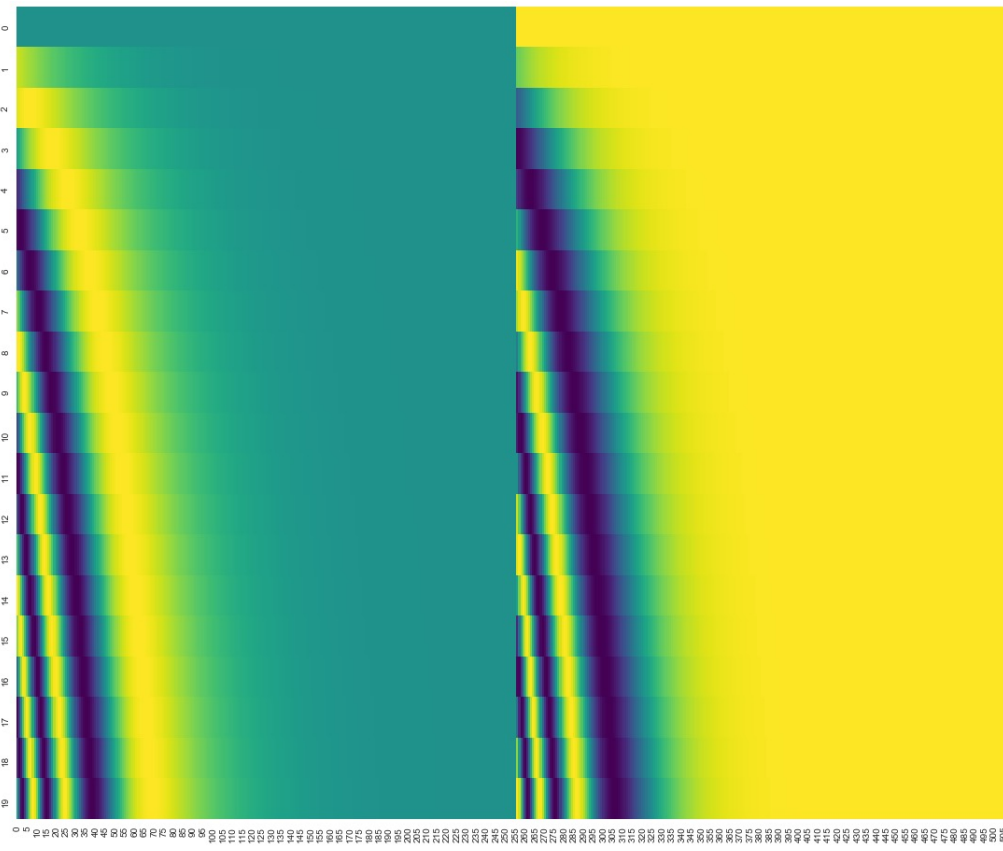
An example of embeddings with four dimensions:



Position embedding for location 0

Position embeddings

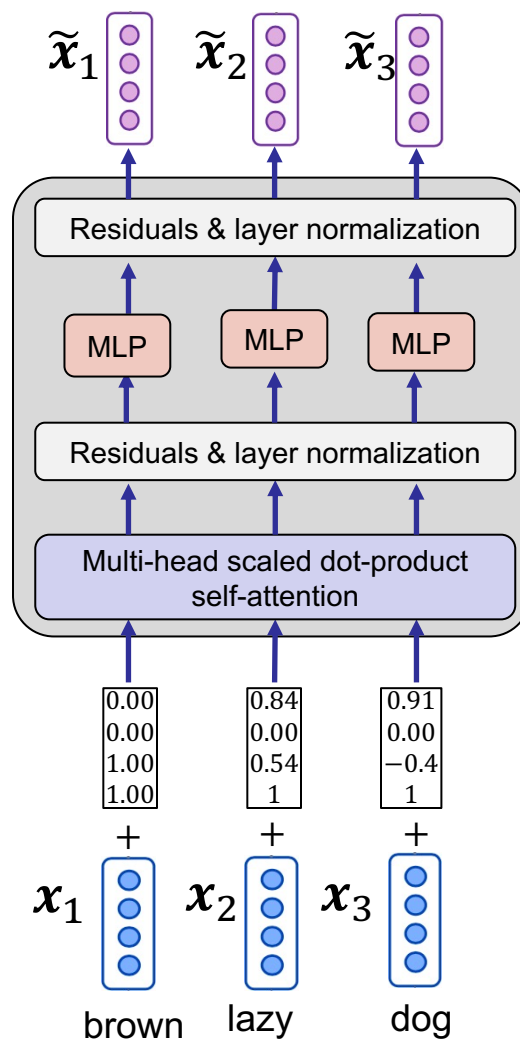
Position embedding for location 20



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

Dimensions (512)

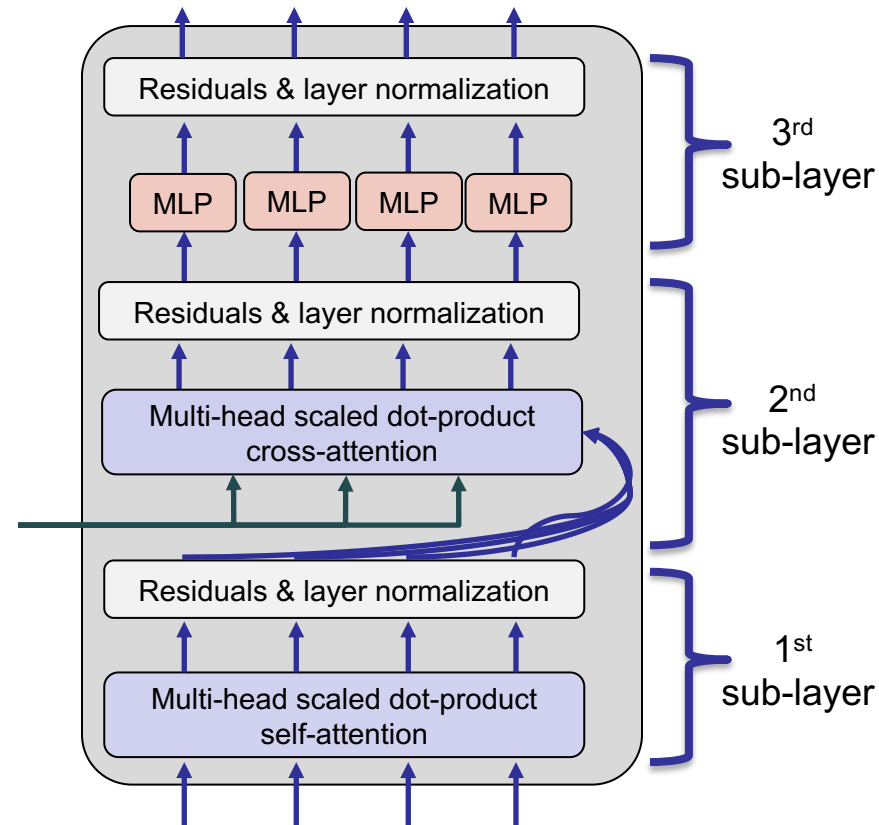
# Transformers with position embedding



# Transformer decoder

Transformer decoder consists of three sub-layers:

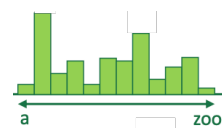
- 1<sup>st</sup> : **Masked** multi-head self-attention
- 2<sup>nd</sup> : Multi-head **cross attention**
  - Queries are the outputs of the previous sub-layer (contextualized embeddings)
  - Values come from outside (encoder)
  - Transformer decoder attends to the embeddings of an encoder (usually a Transformer encoder)
- 3<sup>rd</sup> : Position-wise multi-layer perceptron (feed forward)



# Agenda

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

# Seq2seq with Transformers



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

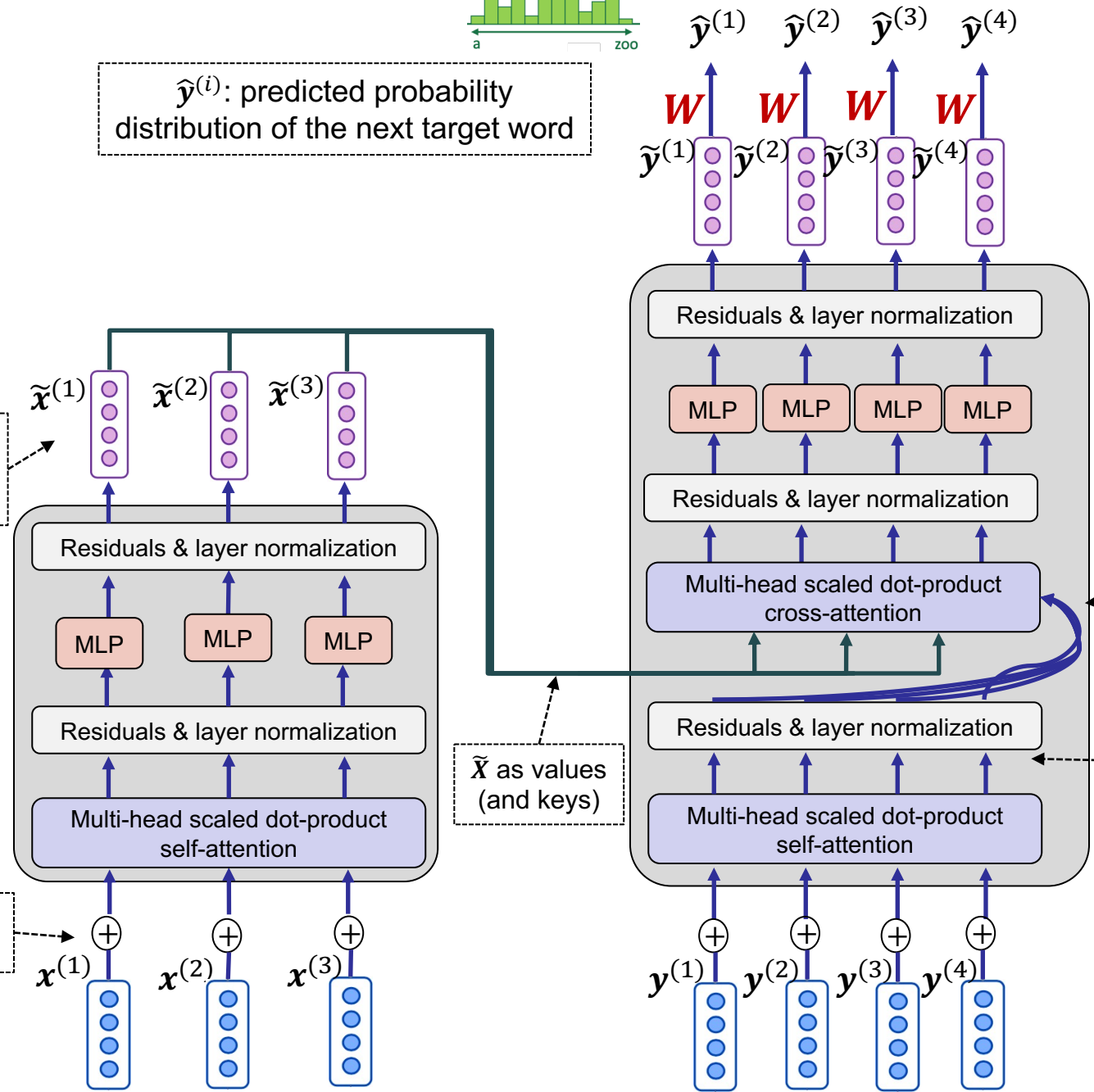
Contextualized input embeddings  $\tilde{X}$

Position embedding

$\tilde{X}$  as values (and keys)

target embeddings as queries

Contextualized target embeddings



# Seq2seq with Transformers – training

1. Data preparation: Source sequence  $X$  and target sequence  $Y$  are both started with  $\langle \text{bos} \rangle$  and ended with  $\langle \text{eos} \rangle$
2. Transformer encoder: outputs contextualized embeddings  $\tilde{X}$ 
  - $\tilde{X}$  is in size  $|X| \times d$ , where  $d$  is the embedding vector dimensions
3. Decoder self-attention: create contextualized embeddings  $\tilde{Y}$ 
  - $\tilde{Y}$  is in size  $|Y| \times d$
4. Decoder cross-attention:  $\tilde{Y}$  vectors attend to  $\tilde{X}$  vectors ( $\tilde{Y}$  queries  $\tilde{X}$  values)
  - Resulting matrix is in size  $|Y| \times d$
5. Prediction: Transformer decoder outputs are used to calculate  $\hat{Y}$  – the probability distributions of the next tokens
  - $\hat{Y}$  is in size  $|Y| \times |\mathbb{V}_d|$
  - E.g., vector  $\hat{y}^{(2)}$  predicts the probability distribution of the token at position 3
6. Loss: is calculated using Negative Log Likelihood of the actual next words
  - E.g., based on the probability value of token  $y^{(3)}$  in vector  $\hat{y}^{(2)}$

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

- Every token has access to what it is supposed to predict!

# Masking attentions

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

## Example

- Non-normalized self-attention table of decoder:

		attends to			
		$y^{(1)}$	$y^{(2)}$	$y^{(3)}$	$y^{(4)}$
Vectors of output (target) sequence	$y^{(1)}$	5	3	1	-4
	$y^{(2)}$	1	4	-2	3
	$y^{(3)}$	0	2	2	-3
	$y^{(4)}$	3	-1	1	4



### Non-normalized self-attention values

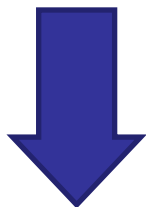
	$y^{(1)}$	$y^{(2)}$	$y^{(3)}$	$y^{(4)}$
$y^{(1)}$	5	3	1	-4
$y^{(2)}$	1	4	-2	3
$y^{(3)}$	0	-2	2	-3
$y^{(4)}$	3	-1	1	4

### attentions masks

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

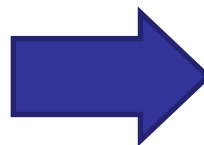
Applying masks to attention scores

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



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

softmax



### Final self-attention values

	$y^{(1)}$	$y^{(2)}$	$y^{(3)}$	$y^{(4)}$
$y^{(1)}$	1.00	0.00	0.00	0.00
$y^{(2)}$	0.04	0.96	0.00	0.00
$y^{(3)}$	0.11	0.01	0.86	0.00
$y^{(4)}$	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.

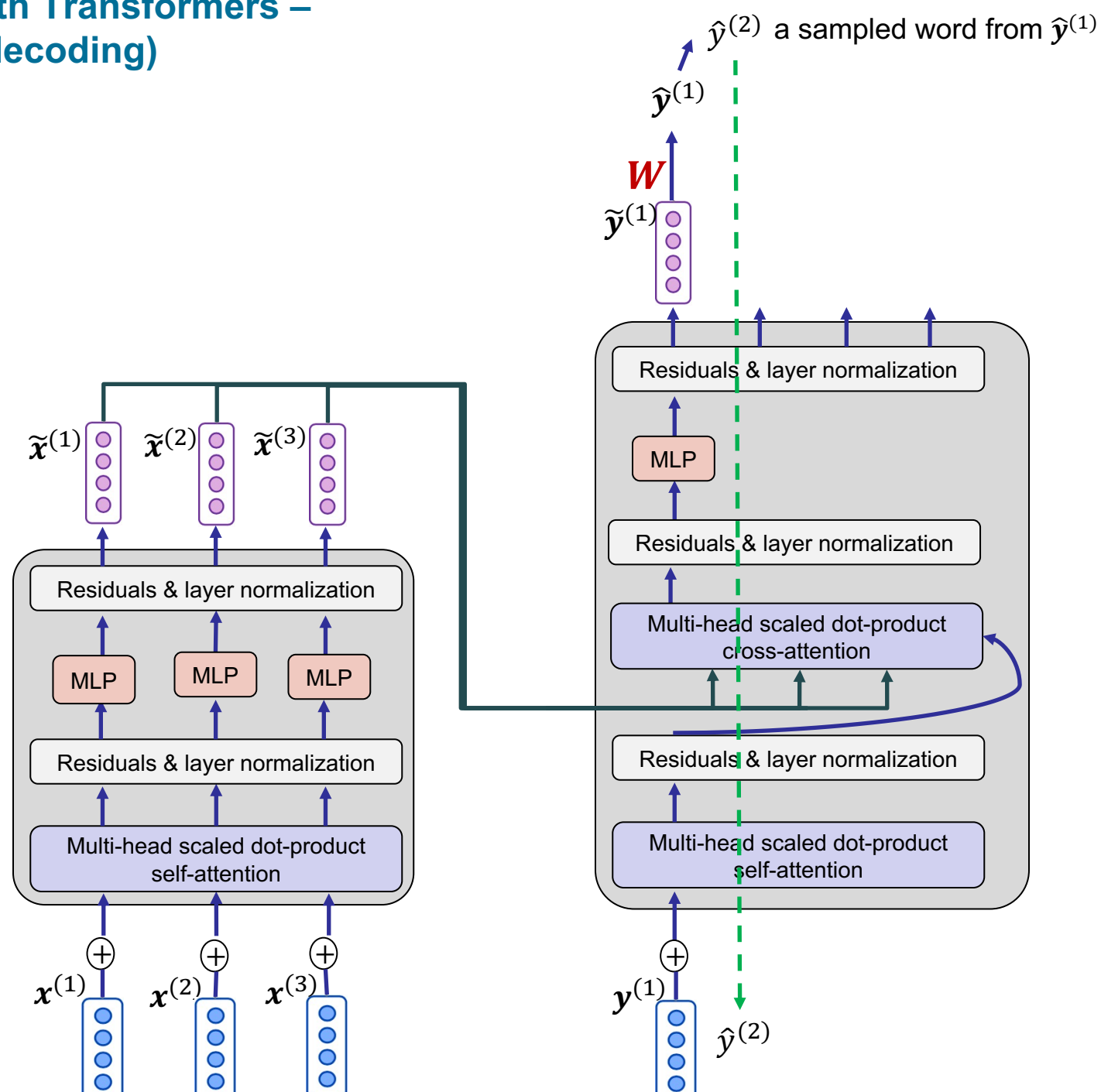
# Seq2seq with Transformers – training (completed!)

1. Data preparation: Source sequence  $X$  and target sequence  $Y$  are both started with  $\langle \text{bos} \rangle$  and ended with  $\langle \text{eos} \rangle$
2. Transformer encoder: outputs contextualized embeddings  $\tilde{X}$ 
  - $\tilde{X}$  is in size  $|X| \times d$ , where  $d$  is the embedding vector dimensions
3. Decoder self-attention: create contextualized embeddings  $\tilde{Y}$  **while masking future tokens**
  - $\tilde{Y}$  is in size  $|Y| \times d$
4. Decoder cross-attention:  $\tilde{Y}$  vectors attend to  $\tilde{X}$  vectors ( $\tilde{Y}$  queries  $\tilde{X}$  values)
  - Resulting matrix is in size  $|Y| \times d$
5. Prediction: Transformer decoder outputs are used to calculate  $\hat{Y}$  – the probability distributions of the next tokens
  - $\hat{Y}$  is in size  $|Y| \times |\mathbb{V}_d|$
  - E.g., vector  $\hat{y}^{(2)}$  predicts the probability distribution of the token at position 3
6. Loss: is calculated using Negative Log Likelihood of the actual next words
  - E.g., based on the probability value of token  $y^{(3)}$  in vector  $\hat{y}^{(2)}$

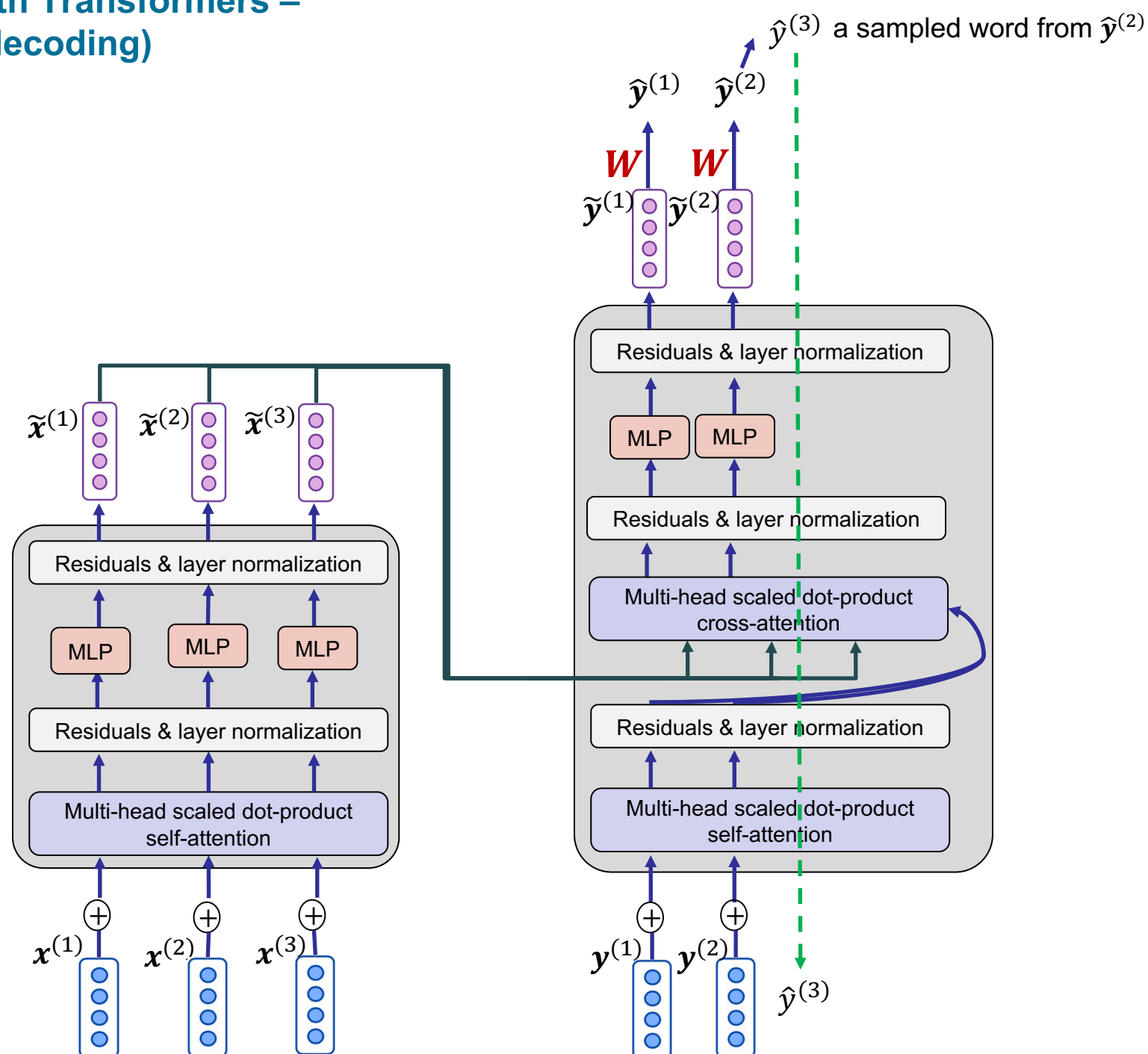
# Inference (decoding)

- In inference similar to training, the encoding of input sequence is done with a single computation (non-autoregressive)
- Decoding is however autoregressive (similar to seq2seq with RNNs):
  - Pass the 1<sup>st</sup> target token, generate the 2<sup>nd</sup> token
  - Pass the 1<sup>st</sup> and 2<sup>nd</sup> generated target tokens, generate the 3<sup>rd</sup> token
  - ...

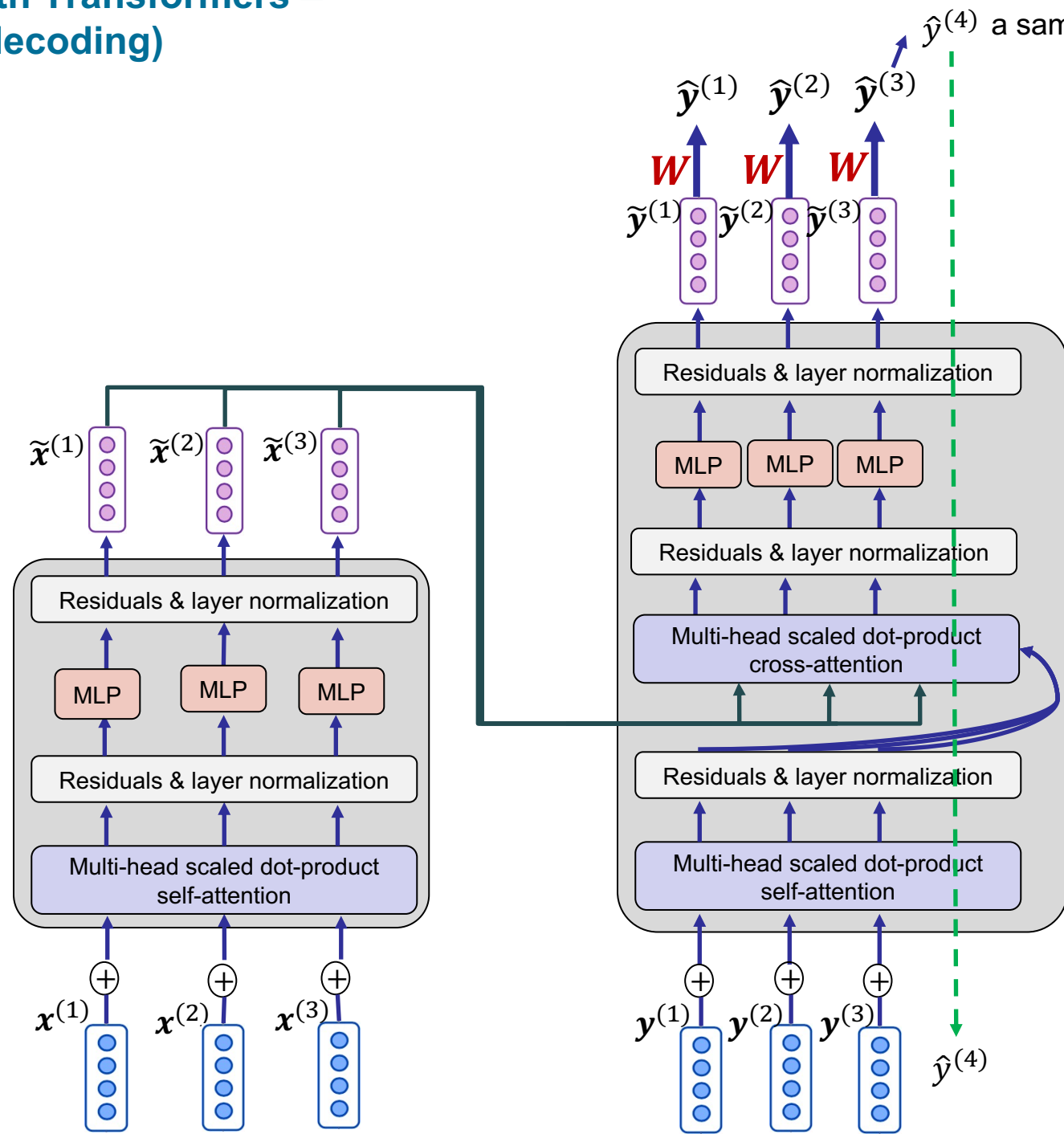
# Seq2seq with Transformers – inference (decoding)



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]
```

