

Transformer Model

Shusen Wang



Transformer Model

- **Original paper:** Vaswani et al. [Attention Is All You Need](#). In *NIPS*, 2017.

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

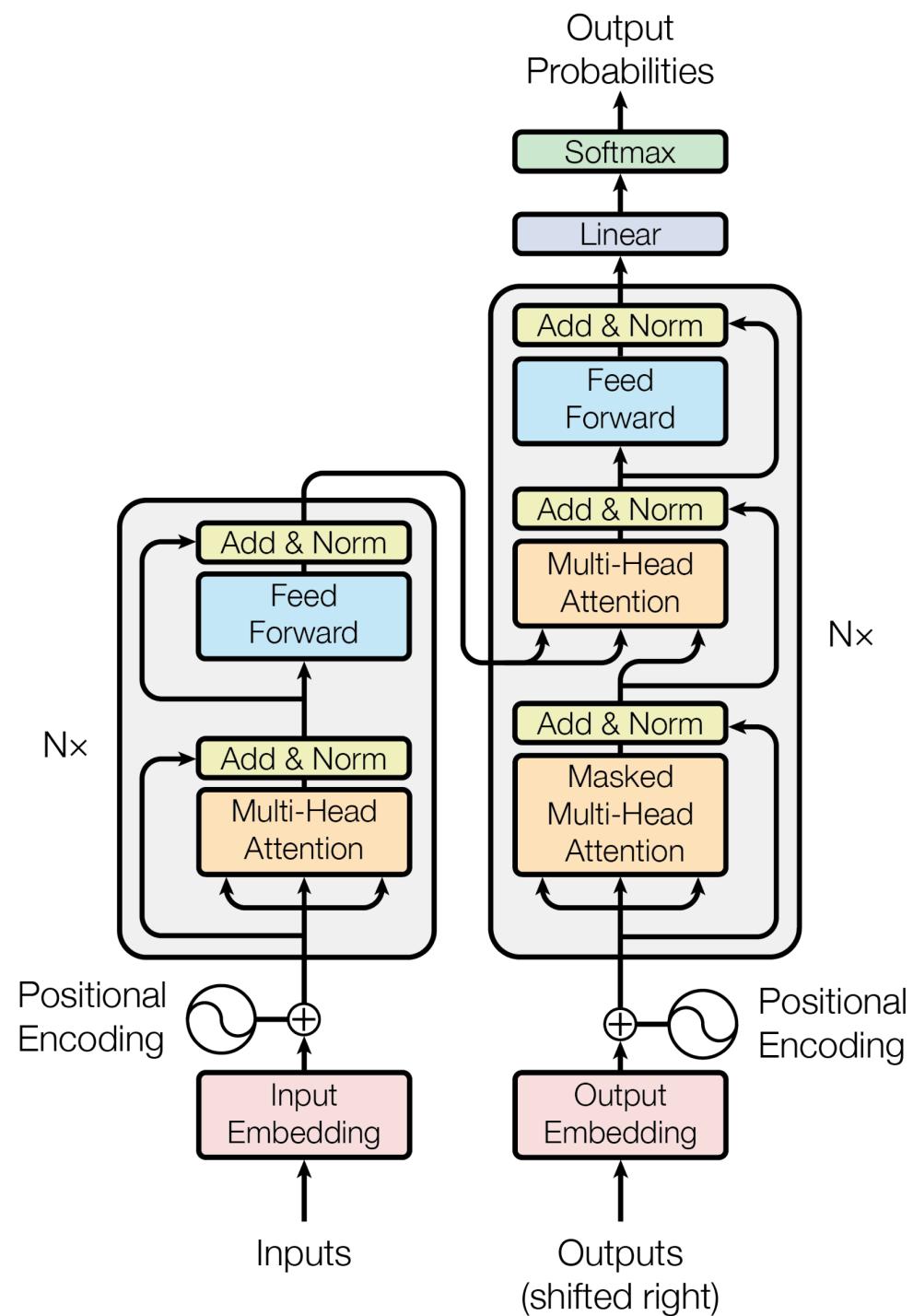
Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

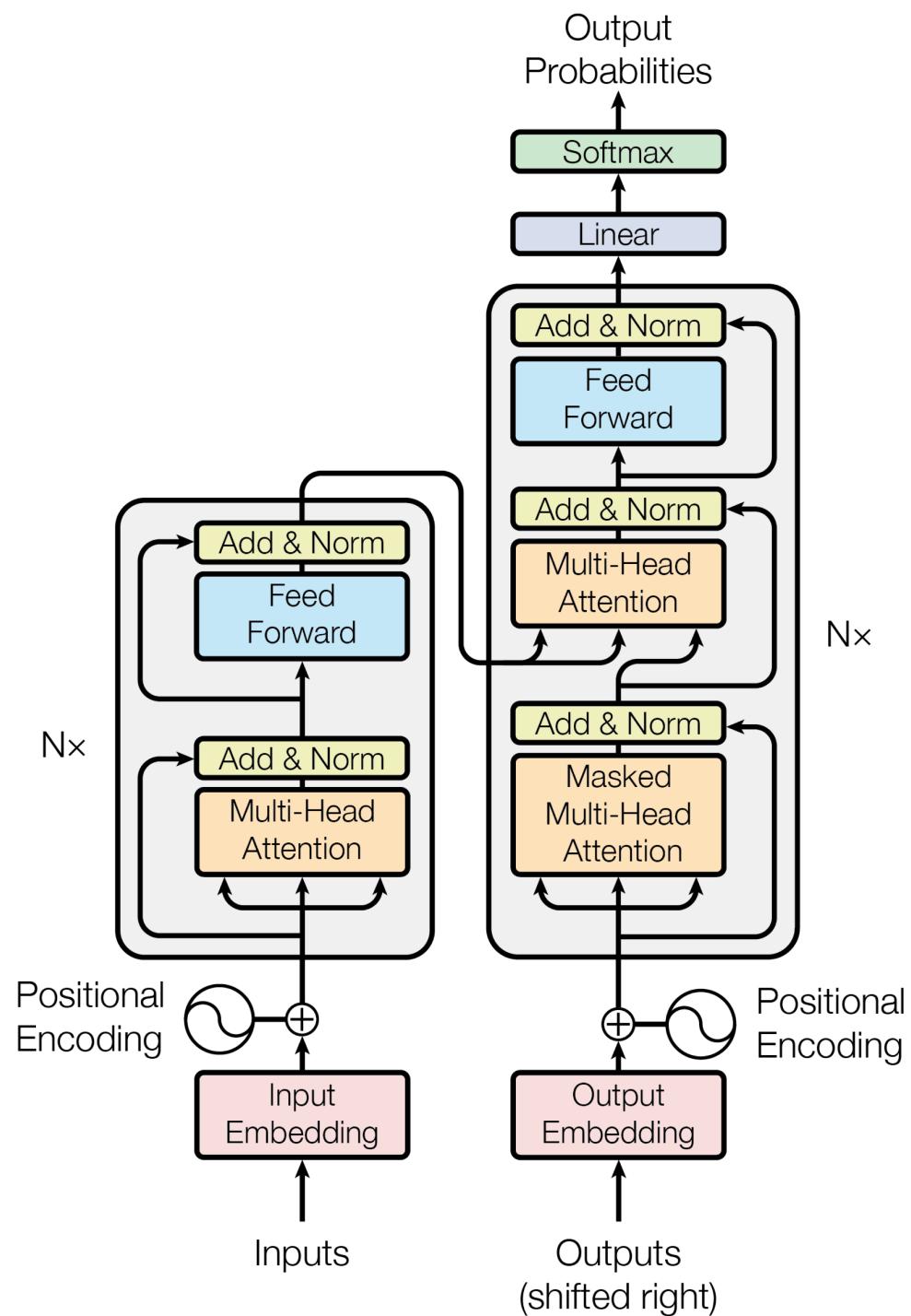
Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com



Transformer Model

- Transformer is a Seq2Seq model.
- Transformer is not RNN.
- Purely based attention and fully-connected layers.
- Much more computations than RNNs.
- Higher performance than RNNs on large datasets.

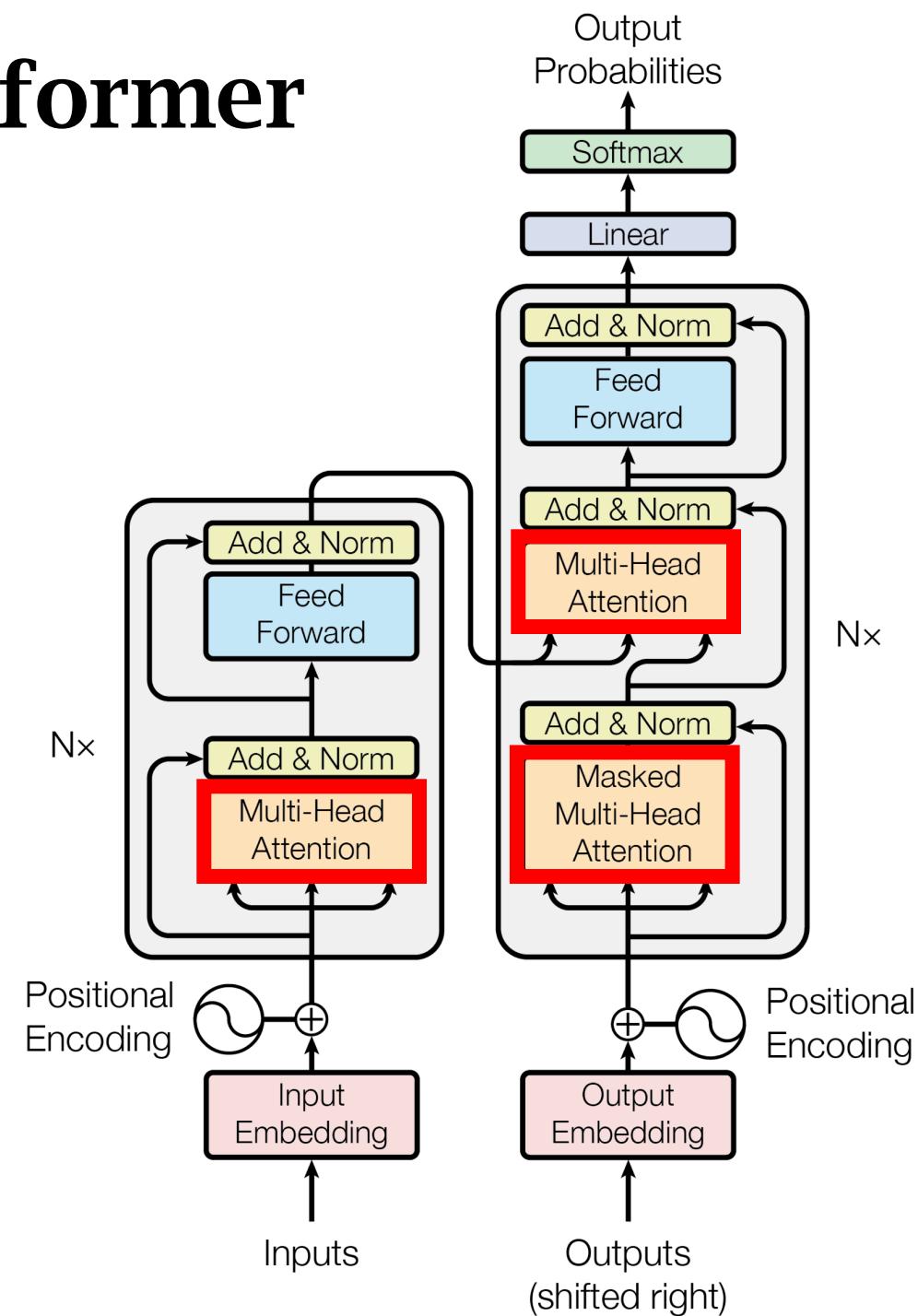


Attention beyond RNNs

Attention in Transformer

Multi-head attention:

- Multiple **single-head attentions**, each has its own parameter matrices.
- Concatenate the outputs.



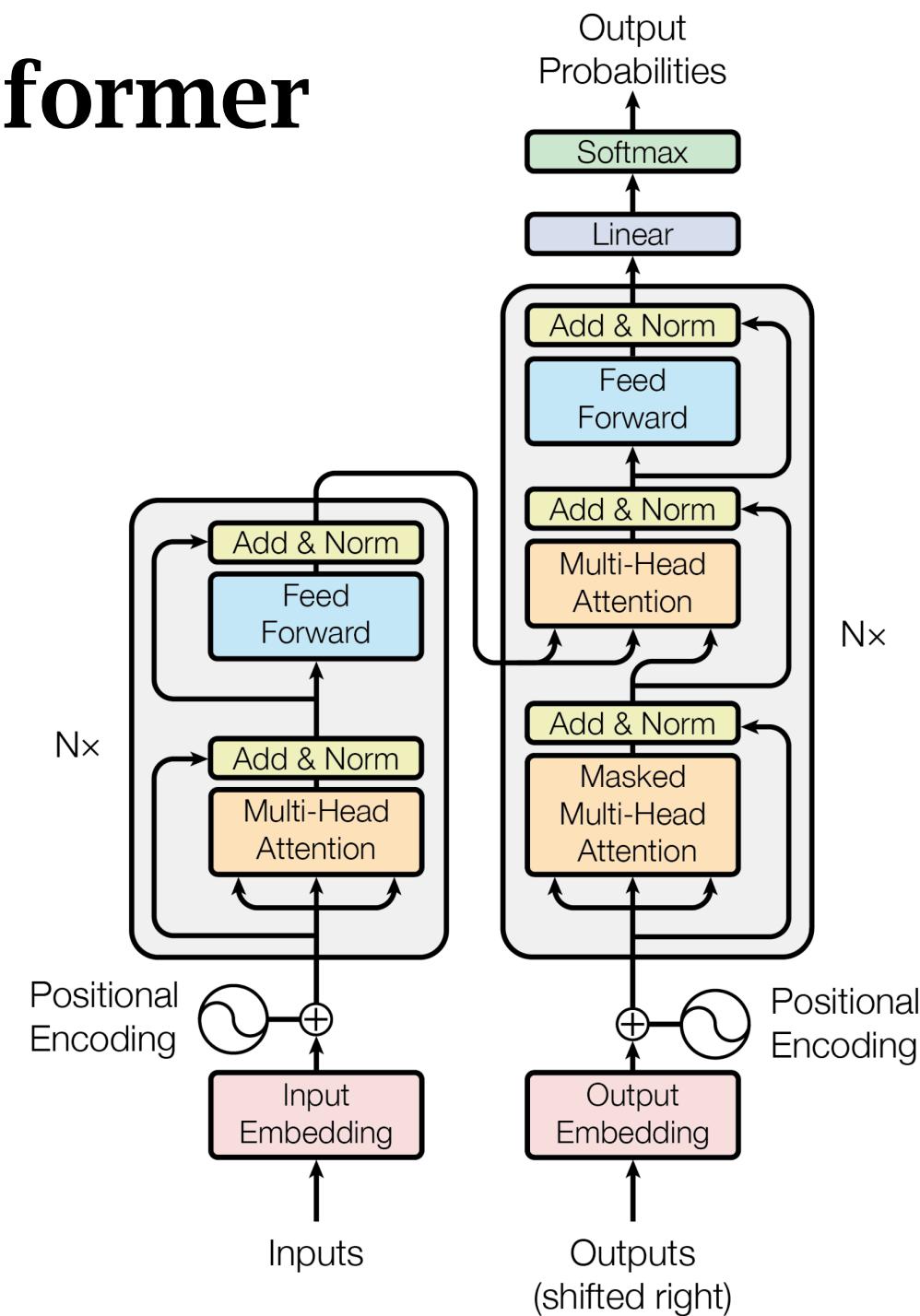
Attention in Transformer

Multi-head attention:

- Multiple single-head attentions, each has its own parameter matrices.
- Concatenate the outputs.

Single-head attention:

- $Z = \text{Attn}(Q, K, V)$.
-
- The diagram shows three input vectors labeled "query", "key", and "value". A red arrow points from "query" to the first input. A green arrow points from "key" to the second input. A purple arrow points from "value" to the third input. The three inputs are grouped together and connected to a central processing block labeled "Attn(Q, K, V)".



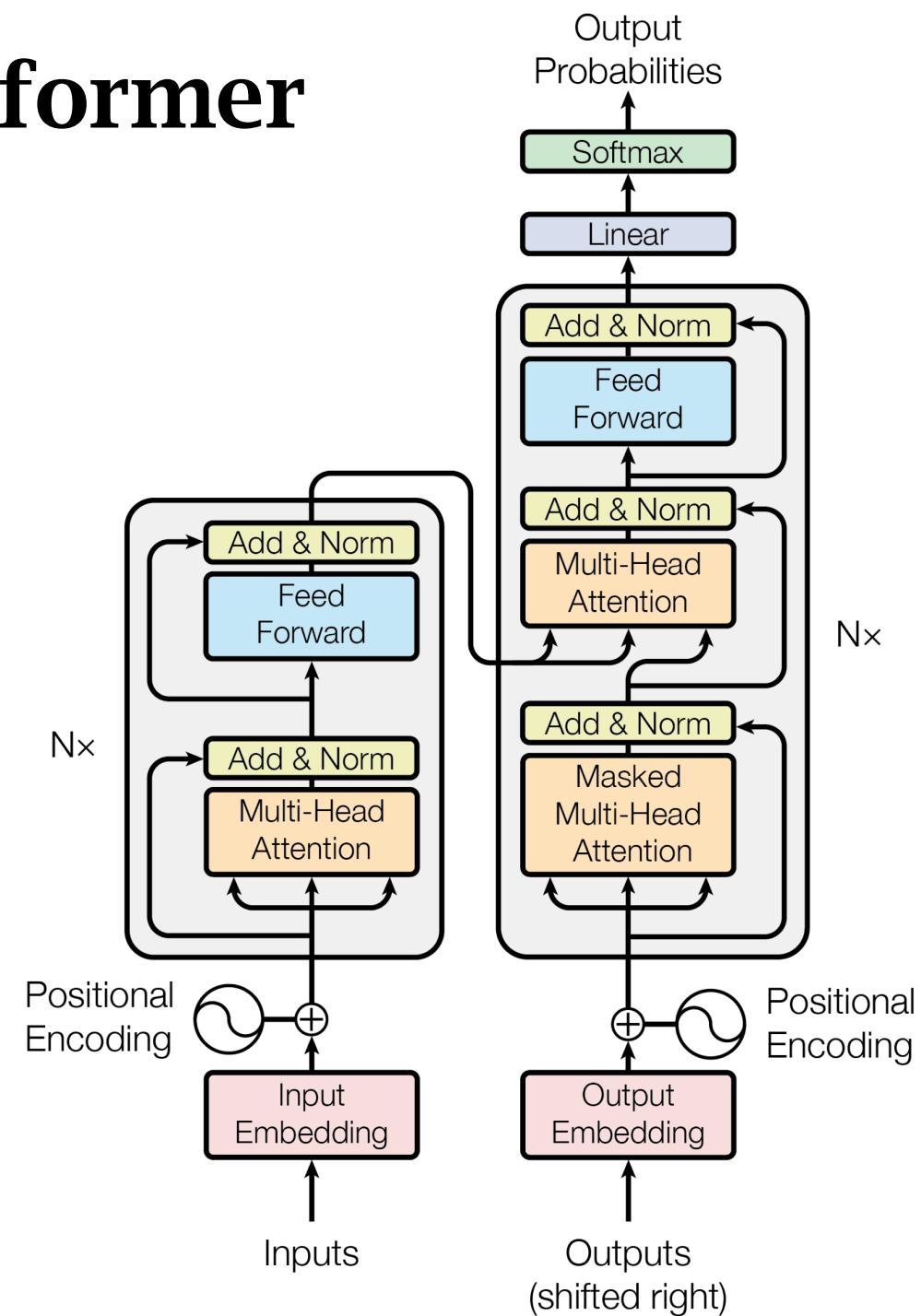
Attention in Transformer

Multi-head attention:

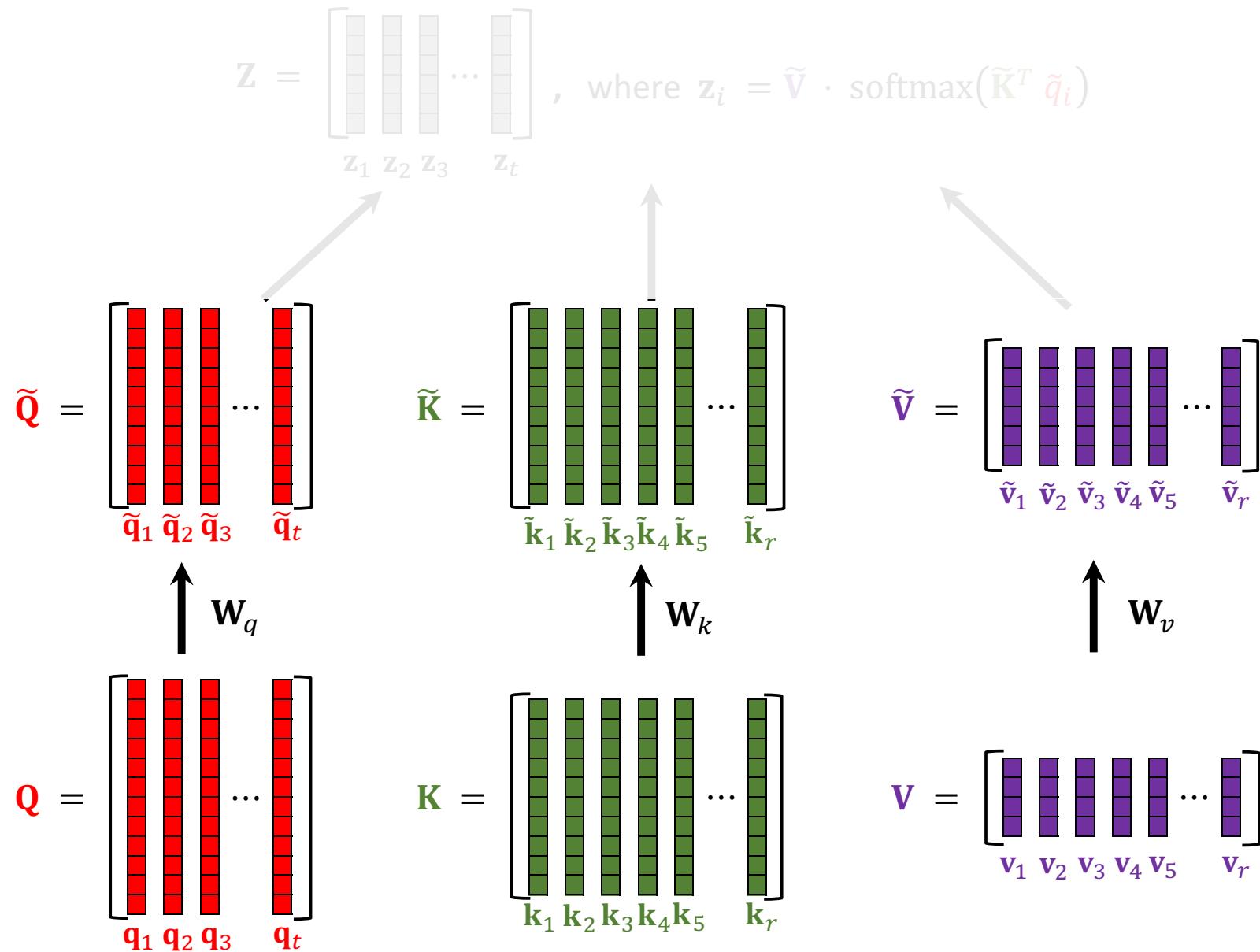
- Multiple single-head attentions, each has its own parameter matrices.
- Concatenate the outputs.

Single-head attention:

- $Z = \text{Attn}(Q, K, V)$.
- Q and Z have t columns.
- t : sequence length.
- K and V have r columns (r is arbitrary).



Single-Head Attention: $Z = \text{Attn}(Q, K, V)$.



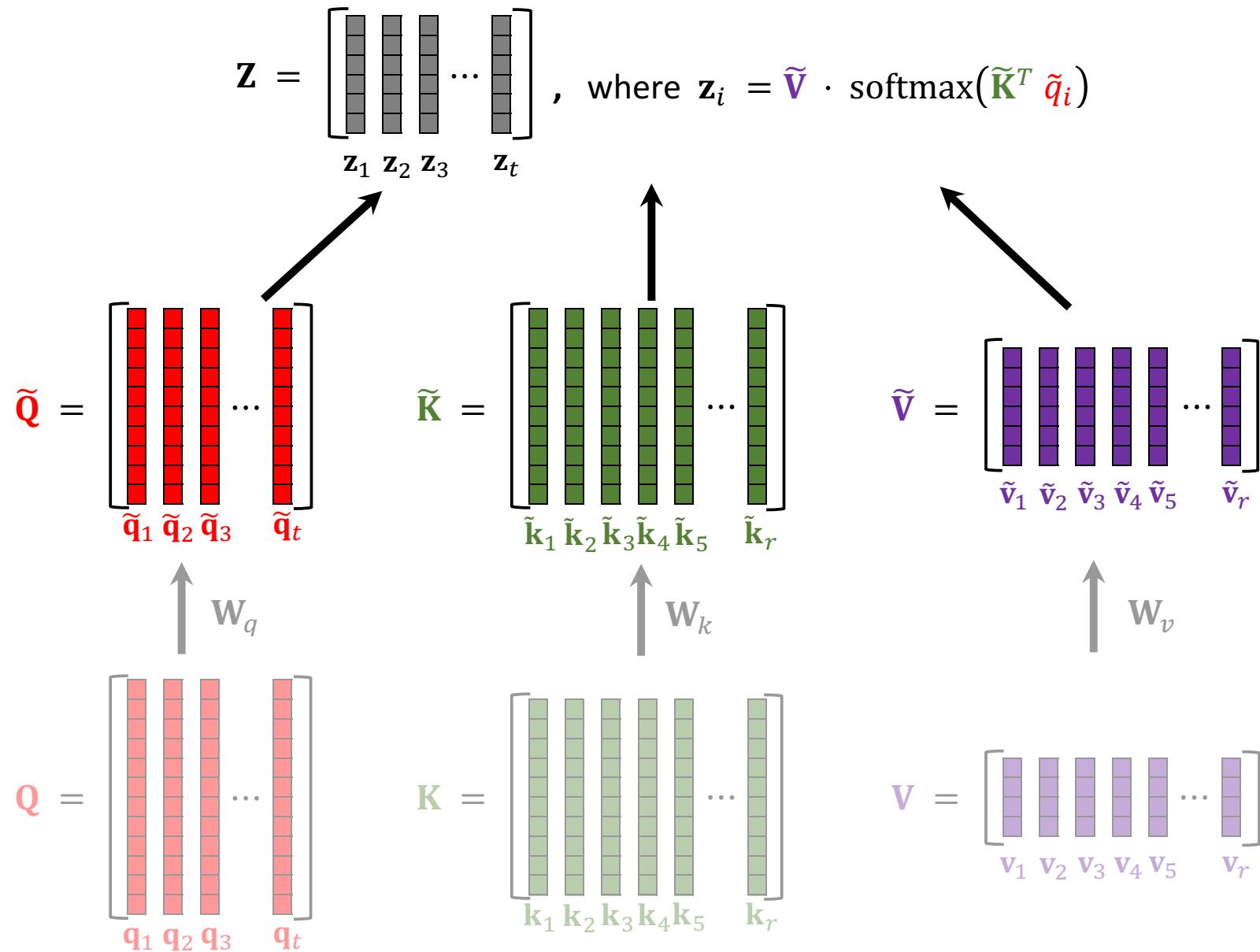
Linear maps:

- $\tilde{Q} = W_q Q$,
- $\tilde{K} = W_k K$,
- $\tilde{V} = W_v V$.

Trainable parameters:

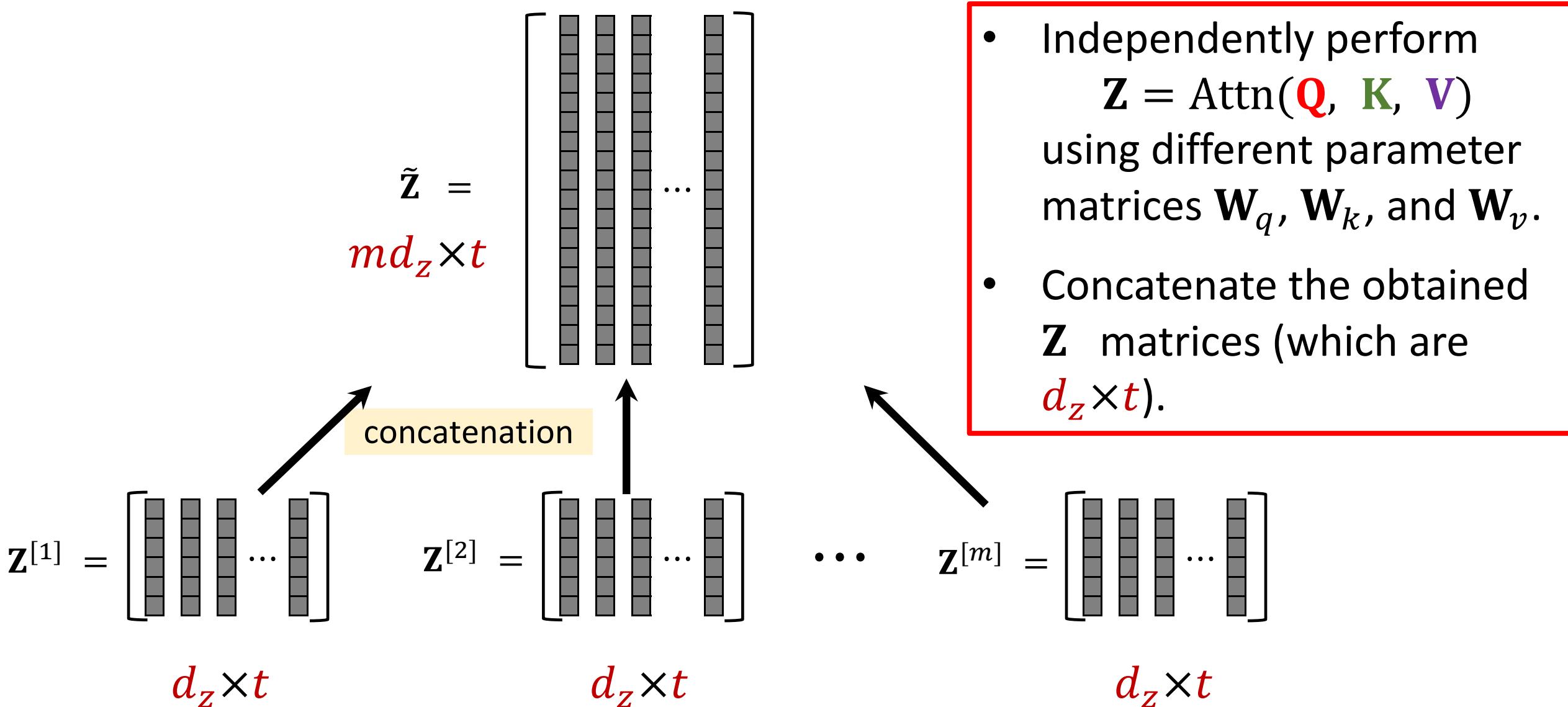
- W_q, W_k, W_v .

Single-Head Attention: $Z = \text{Attn}(Q, K, V)$.



- $\text{rows}(\tilde{Q}) = \text{rows}(\tilde{K})$
- $\text{cols}(\tilde{K}) = \text{cols}(\tilde{V})$
- $\text{rows}(Z) = \text{rows}(\tilde{V})$
- $\text{cols}(Z) = \text{cols}(\tilde{Q})$

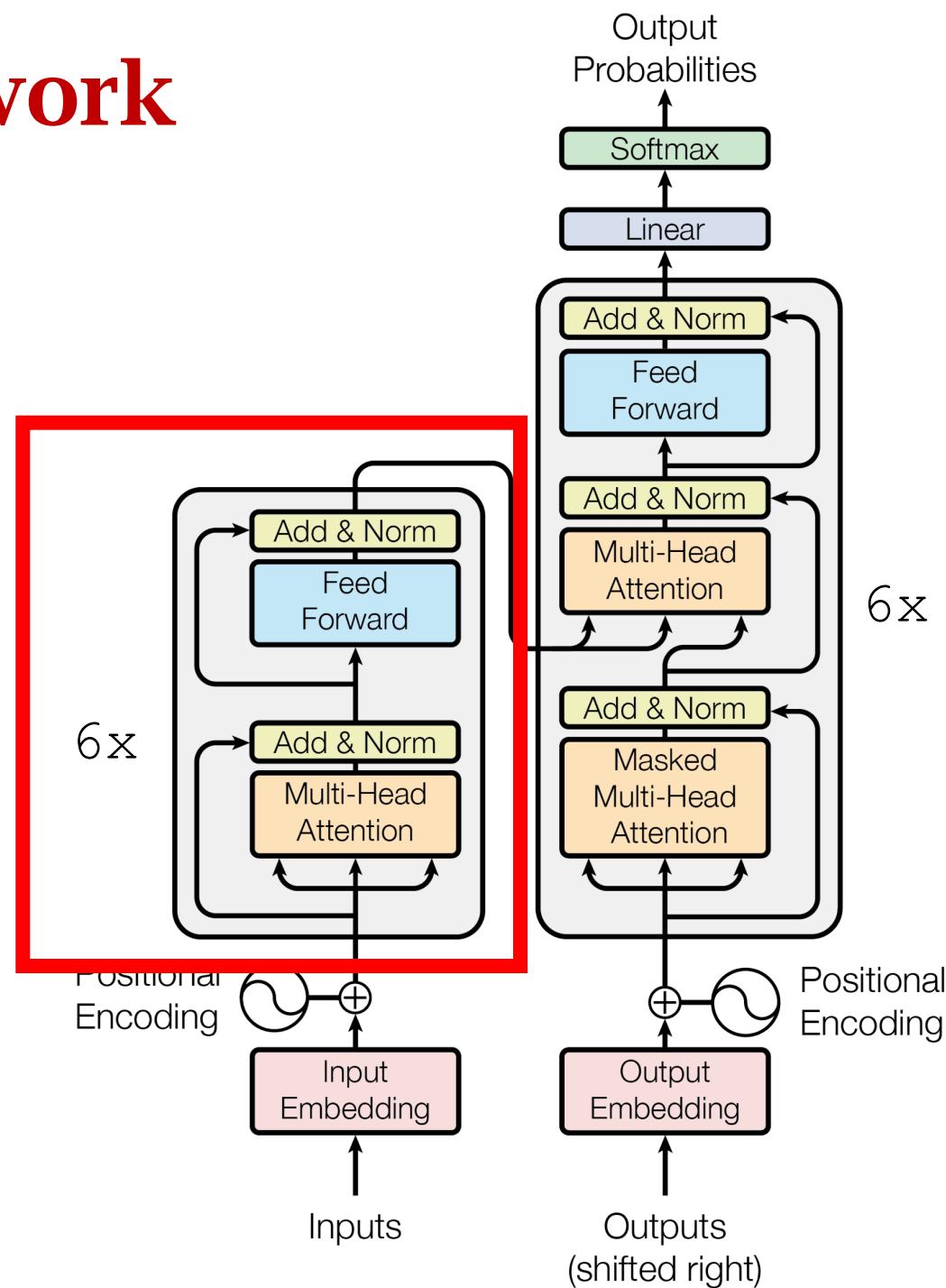
Multi-Head Attention



Encoder of Transformer

Encoder Network

- Encoder has 6 **blocks**.
- 1 Block = **Multi-head attention + Dense**.
- 6 is the result of hyper-parameter tuning; nothing magical about 6.
- Other tricks:
 - Skip connection.
 - Normalization.



Multi-Head Attention + Dense Layer

Multi-Head Attention

$$\tilde{\mathbf{z}} = [\tilde{\mathbf{z}}_1, \tilde{\mathbf{z}}_2, \tilde{\mathbf{z}}_3, \dots, \tilde{\mathbf{z}}_t] \in \mathbb{R}^{md_z \times t}$$

Concatenation

$$\mathbf{z}^{[1]} \in \mathbb{R}^{d_z \times t}$$

$$\mathbf{z}^{[2]} \in \mathbb{R}^{d_z \times t}$$

...

$$\mathbf{z}^{[m]} \in \mathbb{R}^{d_z \times t}$$

Attention w/ different
parameter matrices.

(Q,

K,

V)

Multi-Head Attention + Dense Layer

$$\tilde{\mathbf{U}} = [\tilde{\mathbf{u}}_1, \tilde{\mathbf{u}}_2, \tilde{\mathbf{u}}_3, \dots, \tilde{\mathbf{u}}_t] \in \mathbb{R}^{d_u \times t}$$

Dense layer is applied to every column independently.

$$\tilde{\mathbf{Z}} = [\tilde{\mathbf{z}}_1, \tilde{\mathbf{z}}_2, \tilde{\mathbf{z}}_3, \dots, \tilde{\mathbf{z}}_t] \in \mathbb{R}^{md_z \times t}$$

Concatenation

$$\mathbf{z}^{[1]} \in \mathbb{R}^{d_z \times t}$$

$$\mathbf{z}^{[2]} \in \mathbb{R}^{d_z \times t}$$

...

$$\mathbf{z}^{[m]} \in \mathbb{R}^{d_z \times t}$$

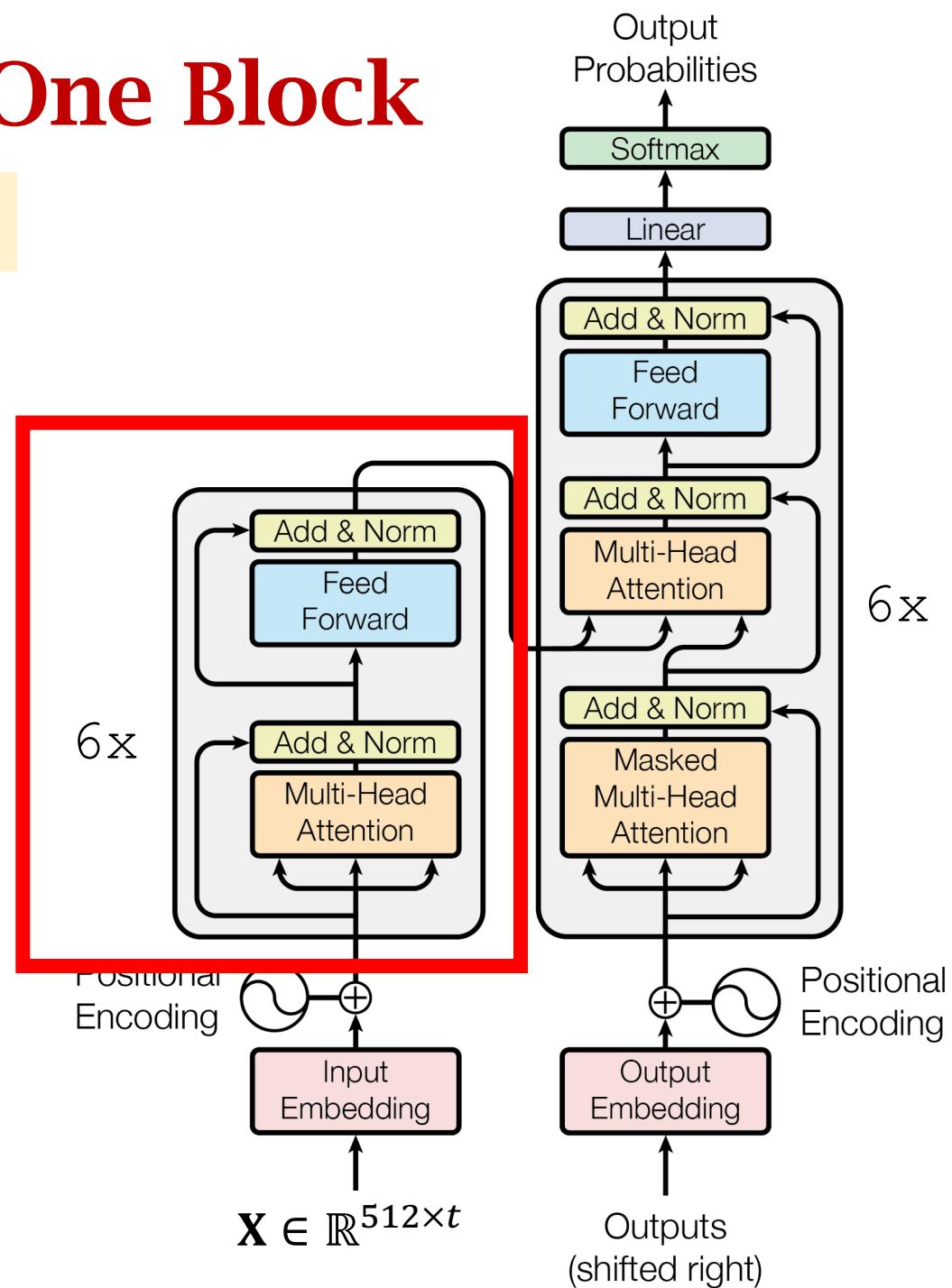
Attention w/ different parameter matrices.

$$(\mathbf{Q}, \mathbf{K}, \mathbf{V})$$

Encoder Network: One Block

Ignore skip connection and normalization.

- Input: $\mathbf{X} \in \mathbb{R}^{512 \times t}$; (t is the seq length.)
- Set $\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{X}$.

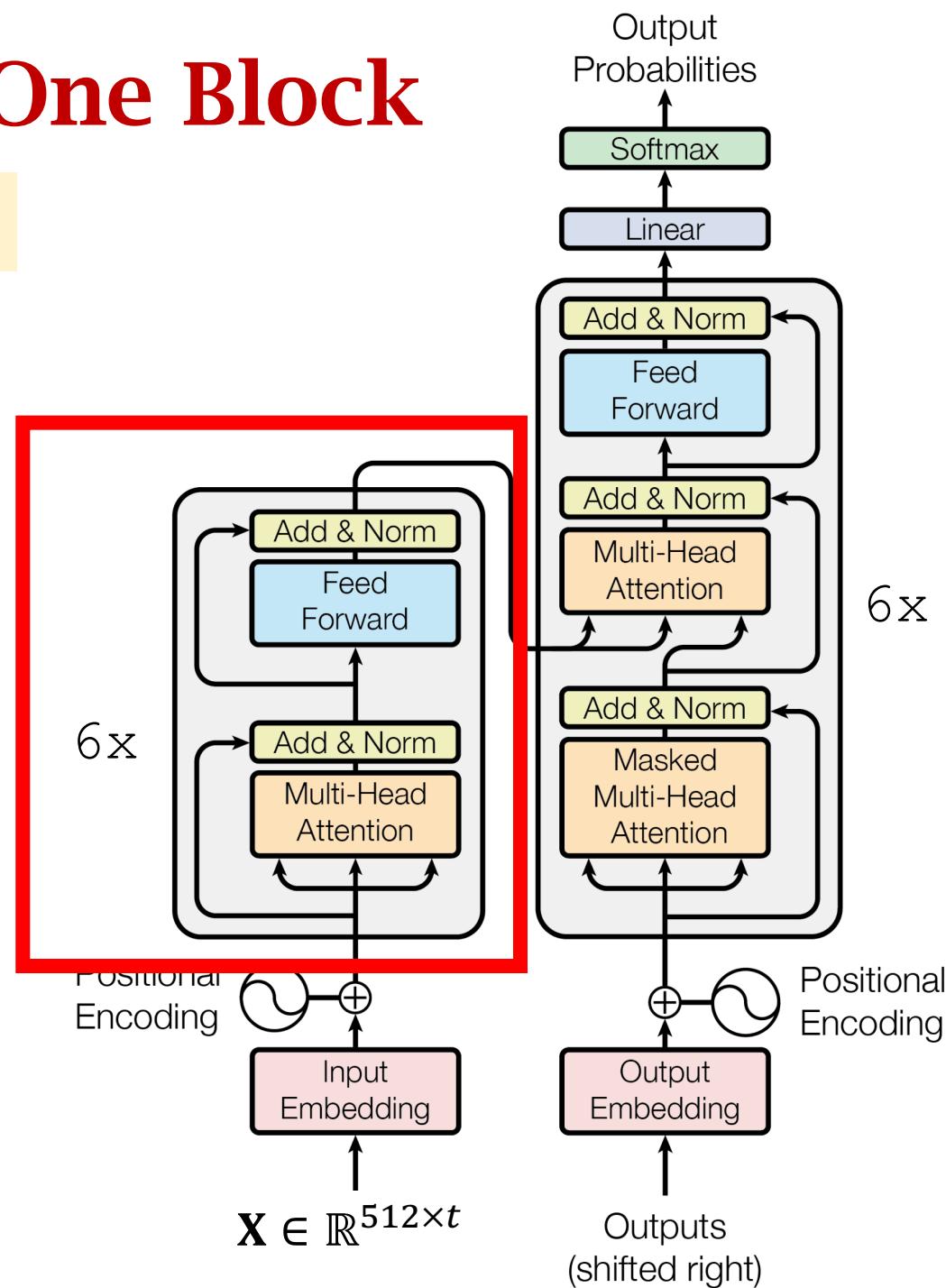


Encoder Network: One Block

Ignore skip connection and normalization.

- Input: $\mathbf{X} \in \mathbb{R}^{512 \times t}$; (t is the seq length.)
- Set $\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{X}$.
- Repeat $m = 8$ times:
$$\mathbf{z}^{[i]} = \text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \in \mathbb{R}^{64 \times t}.$$
- $\tilde{\mathbf{Z}} = \text{Concatenate}(\mathbf{z}^{[1]}, \dots, \mathbf{z}^{[m]}) \in \mathbb{R}^{512 \times t}.$

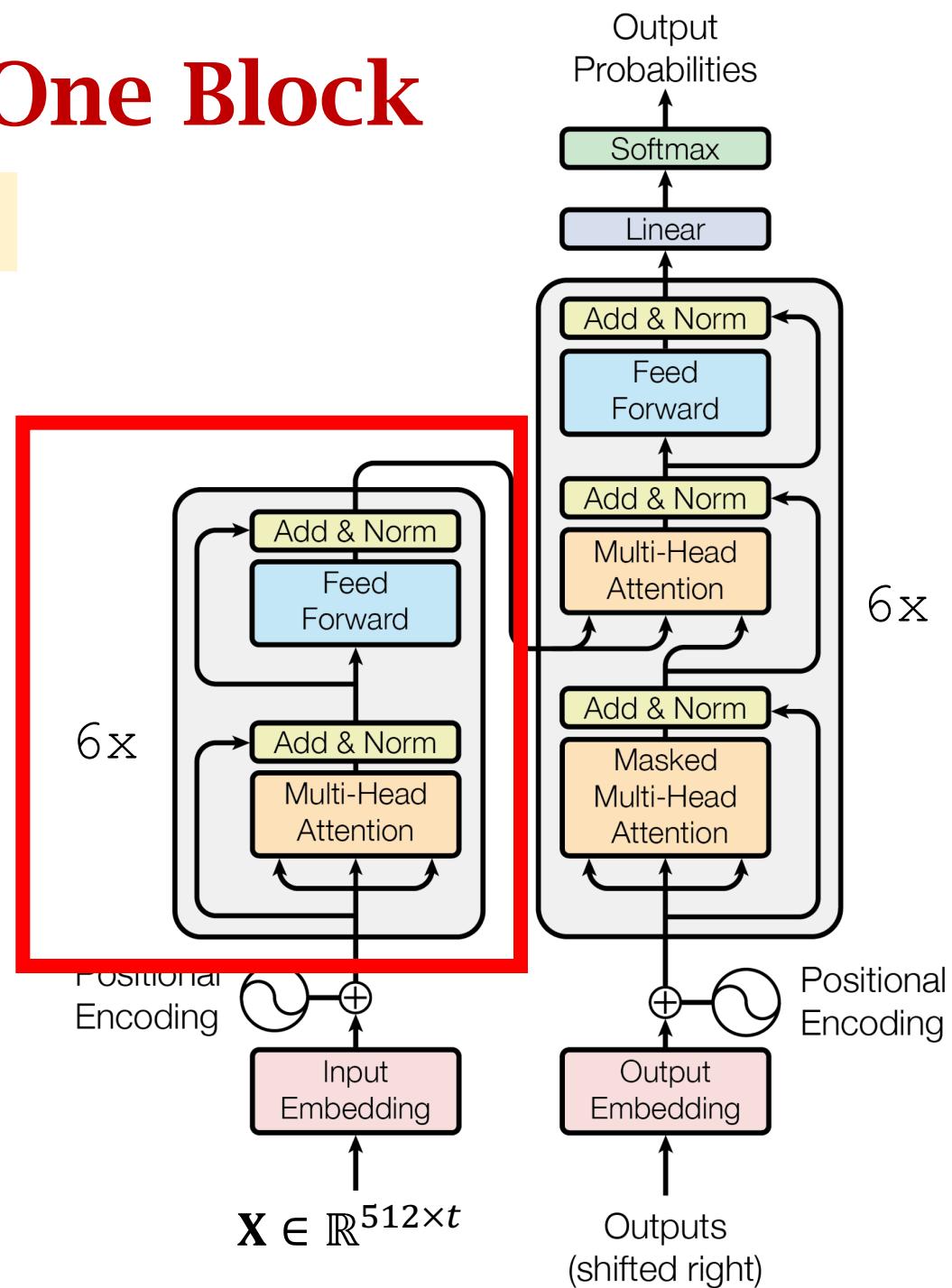
- Make sure the **input shape** and **output shape** are the same.
- Otherwise, skip connection cannot be applied.



Encoder Network: One Block

Ignore skip connection and normalization.

- Input: $\mathbf{X} \in \mathbb{R}^{512 \times t}$; (t is the seq length.)
- Set $\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{X}$.
- Repeat $m = 8$ times:
$$\mathbf{z}^{[i]} = \text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \in \mathbb{R}^{64 \times t}.$$
- $\tilde{\mathbf{Z}} = \text{Concatenate}(\mathbf{z}^{[1]}, \dots, \mathbf{z}^{[m]}) \in \mathbb{R}^{512 \times t}.$
- $\mathbf{U} = \text{DenseLayer}(\tilde{\mathbf{Z}}) \in \mathbb{R}^{512 \times t}.$
- Output: $\mathbf{U} \in \mathbb{R}^{512 \times t}$. (The same shape as \mathbf{X} .)

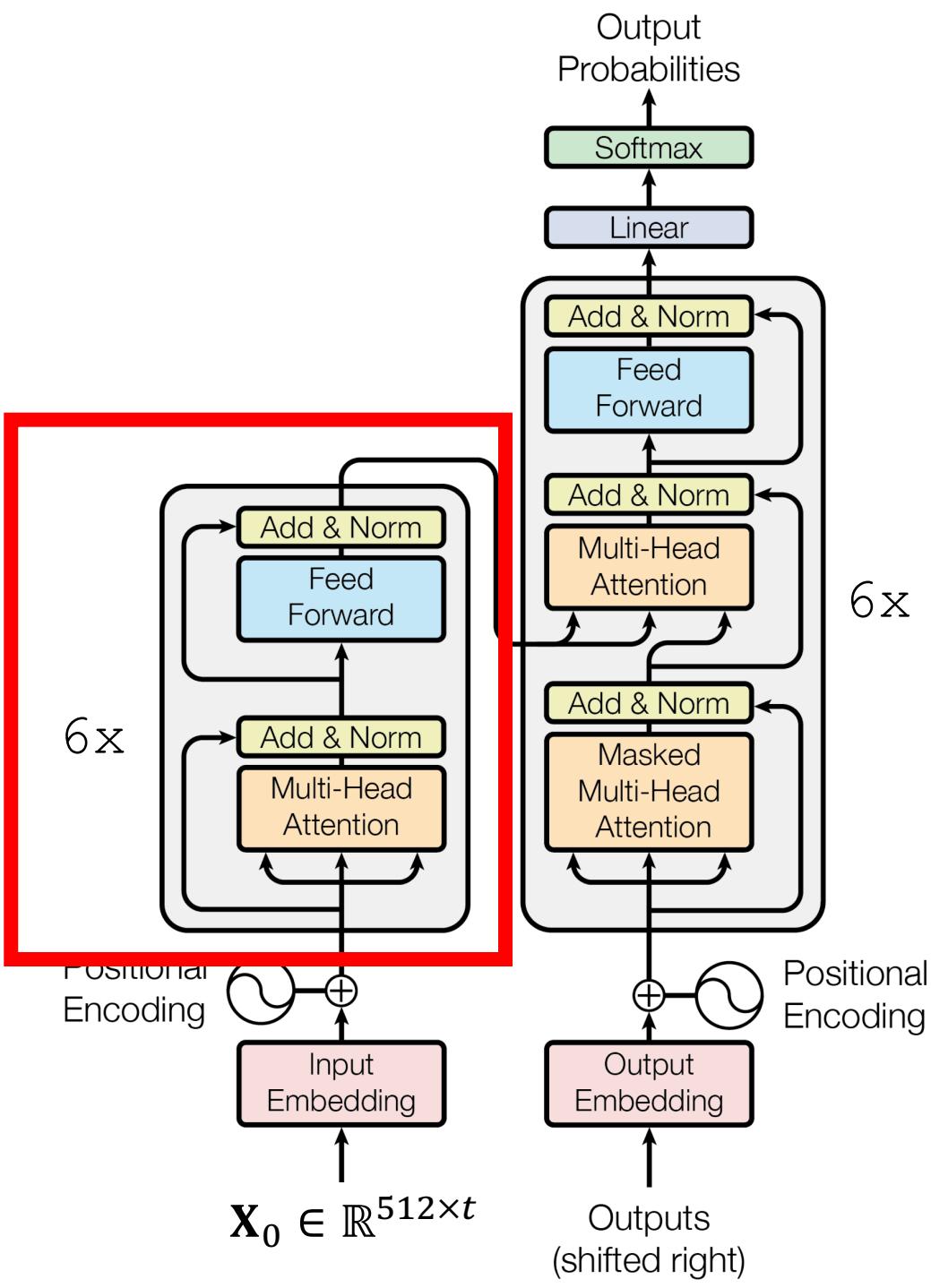


Encoder Network

$$\mathbf{X}_{(1)} \in \mathbb{R}^{512 \times t}$$

Block 1

$$\mathbf{X}_{(0)} \in \mathbb{R}^{512 \times t}$$



Encoder Network

Encoder

$$\mathbf{X}_{(6)} \in \mathbb{R}^{512 \times t}$$

Block 6

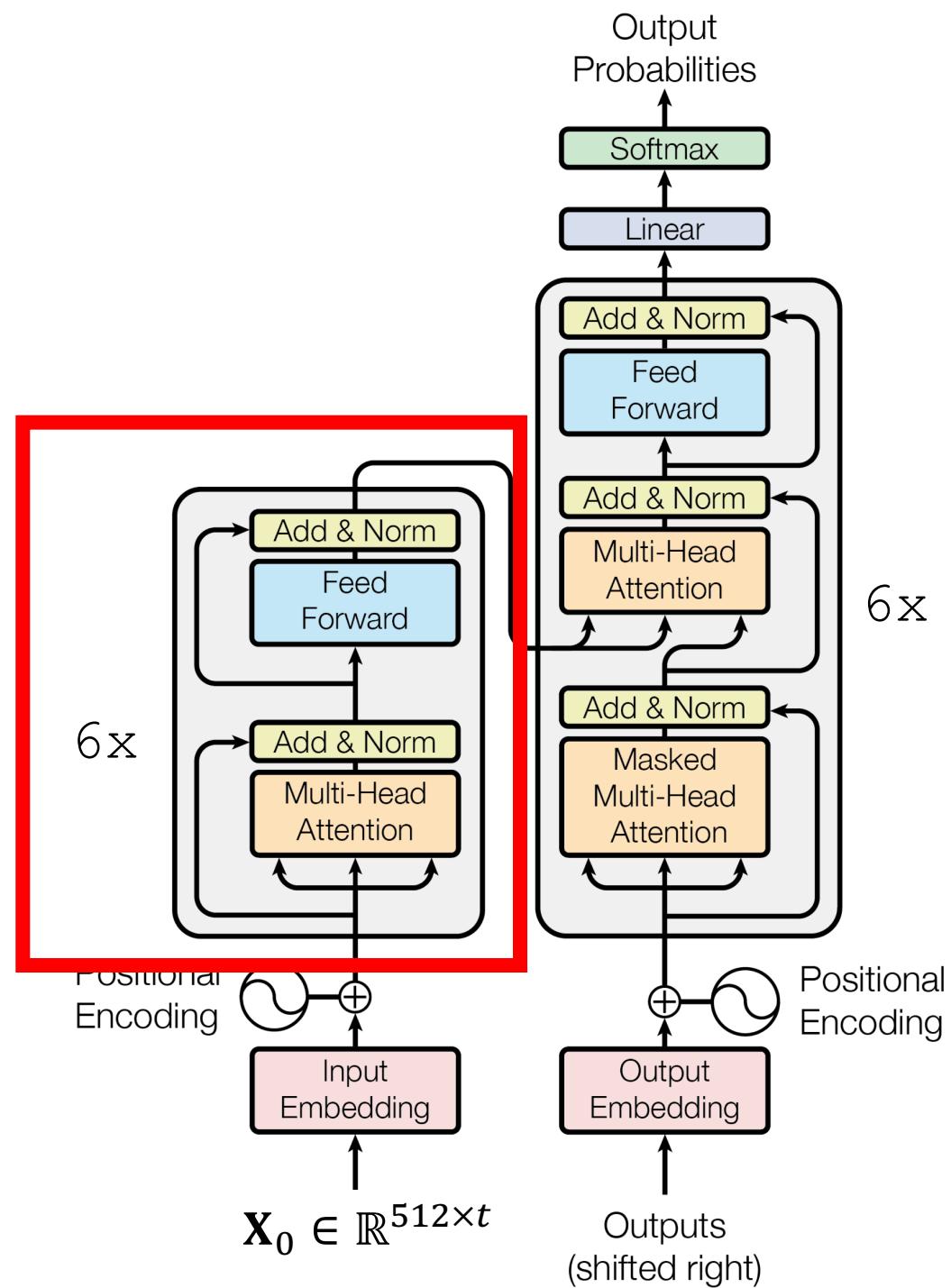


Block 2

$$\mathbf{X}_{(1)} \in \mathbb{R}^{512 \times t}$$

Block 1

$$\mathbf{X}_{(0)} \in \mathbb{R}^{512 \times t}$$



Decoder of Transformer

Decoder Network: One Block

Encoder

$$X_{(6)} \in \mathbb{R}^{512 \times t}$$

Block 6



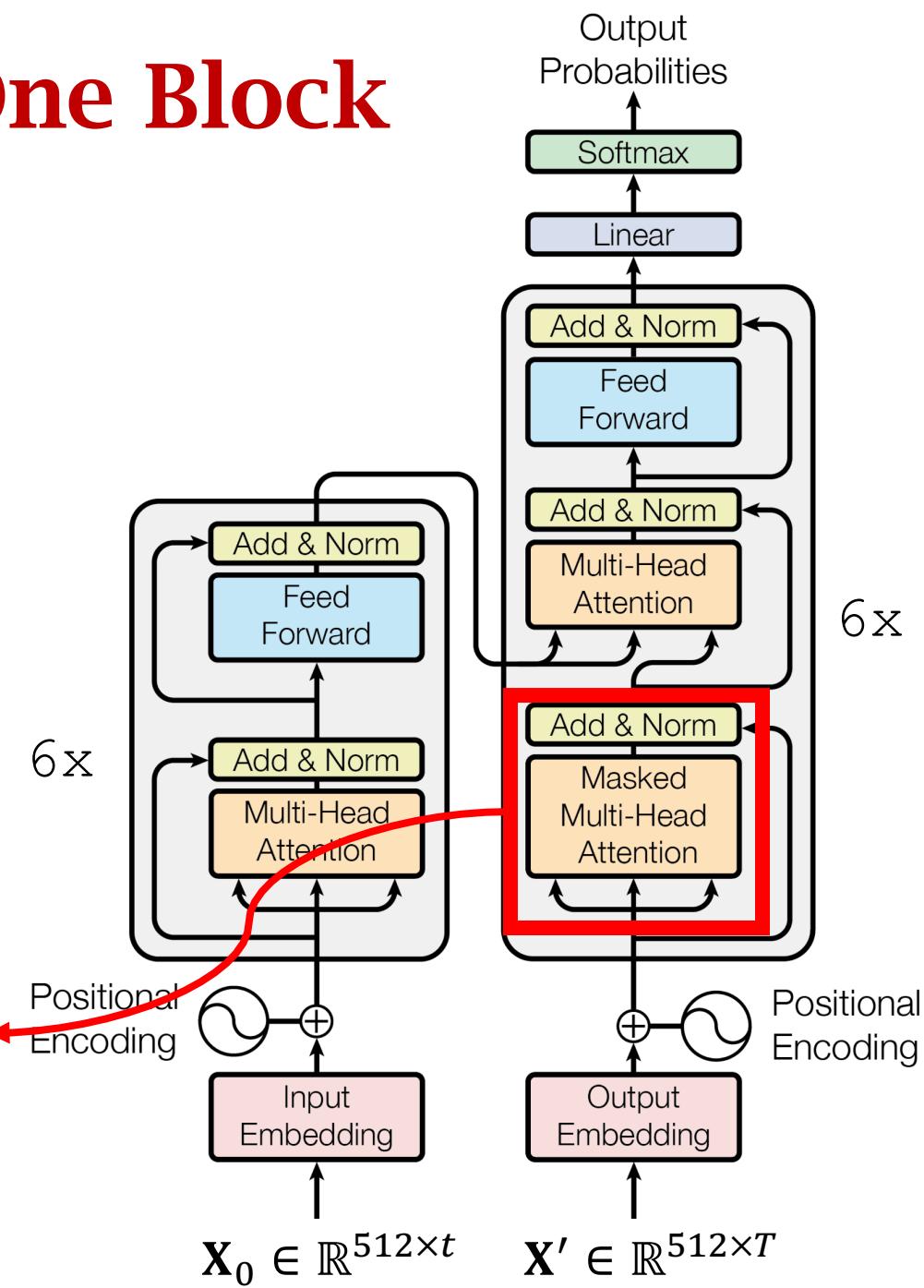
Block 2

$$X_{(1)} \in \mathbb{R}^{512 \times t}$$

Block 1

$$X_{(0)} \in \mathbb{R}^{512 \times t}$$

- Similar to encoder.
- Set $Q = K = V = X'$.



Decoder Network: One Block

Encoder

$$\mathbf{X}_{(6)} \in \mathbb{R}^{512 \times t}$$

Block 6

Block 2

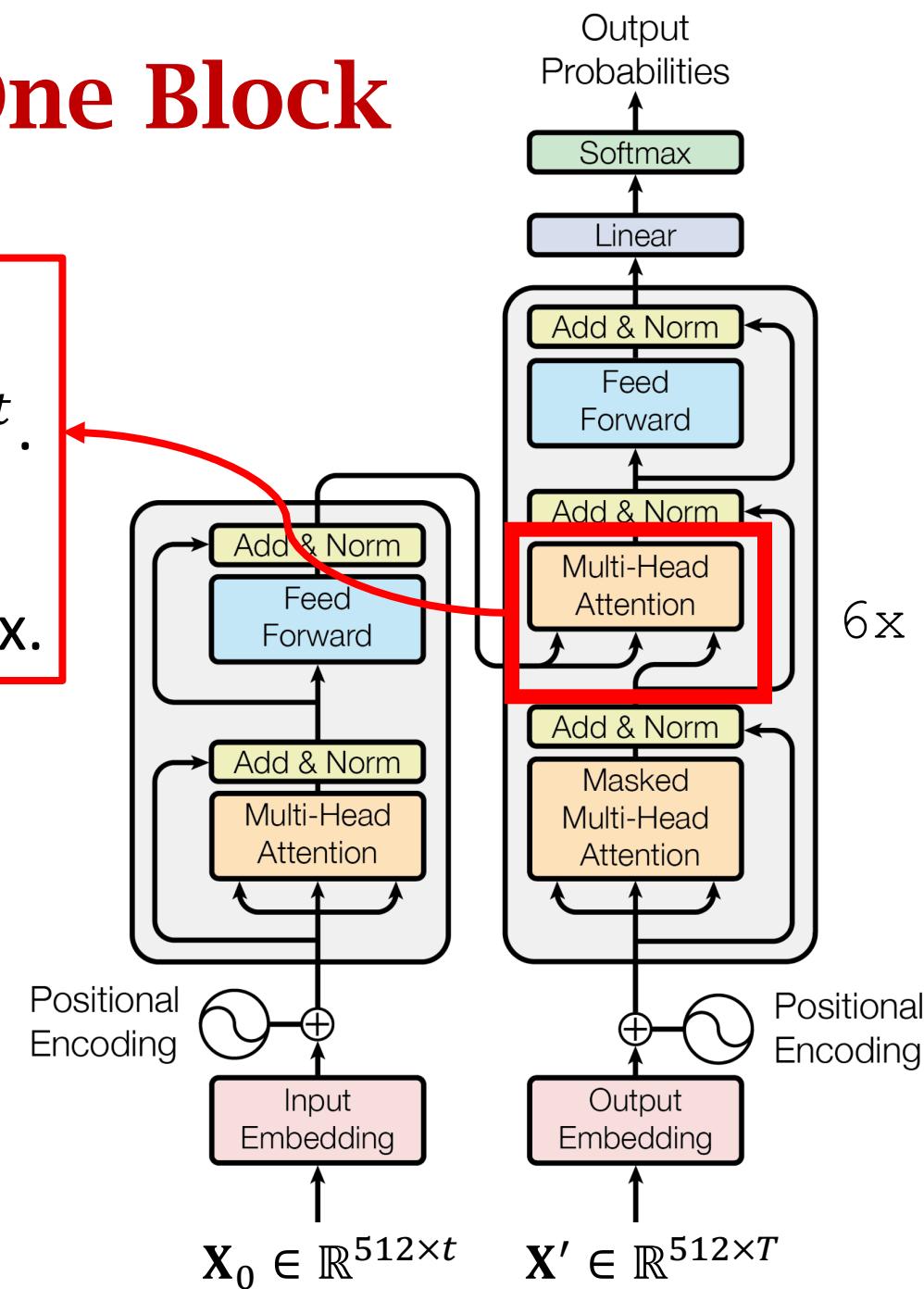
$$\mathbf{X}_{(1)} \in \mathbb{R}^{512 \times t}$$

Block 1

$$\mathbf{X}_{(0)} \in \mathbb{R}^{512 \times t}$$

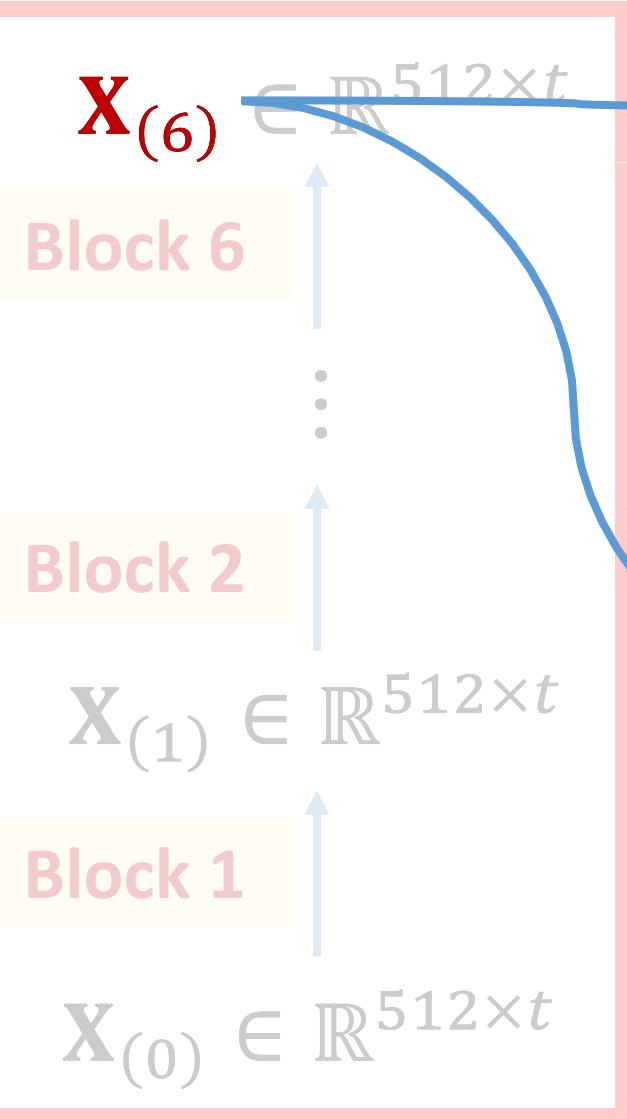
- Set $\mathbf{Q} = \mathbf{X}' \in \mathbb{R}^{512 \times T}$.
- $\mathbf{K} = \mathbf{V} = \mathbf{X}_{(6)} \in \mathbb{R}^{512 \times t}$.
- Multi-head attention outputs a $512 \times T$ matrix.

- Similar to encoder.
- Set $\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{X}'$.

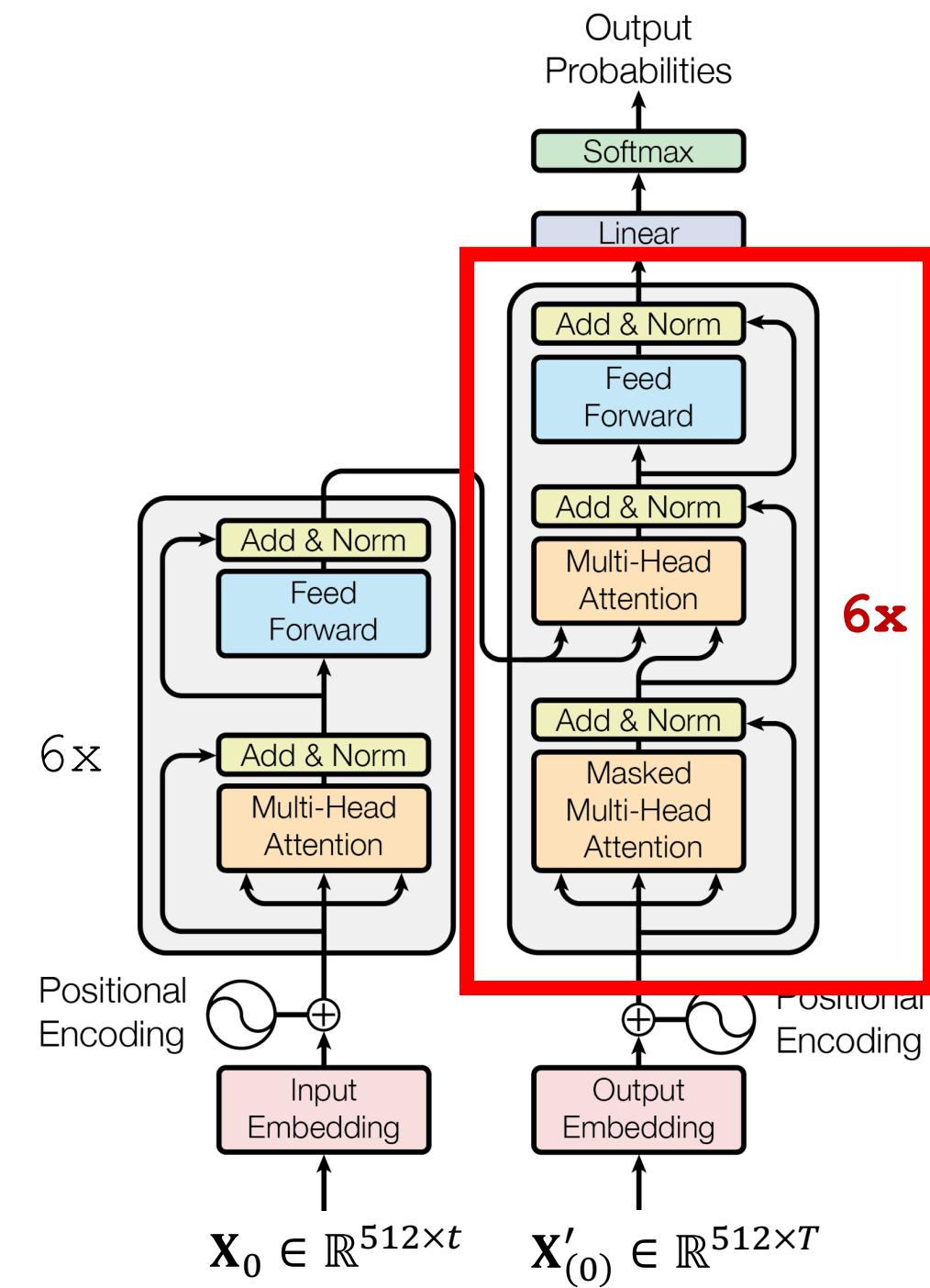
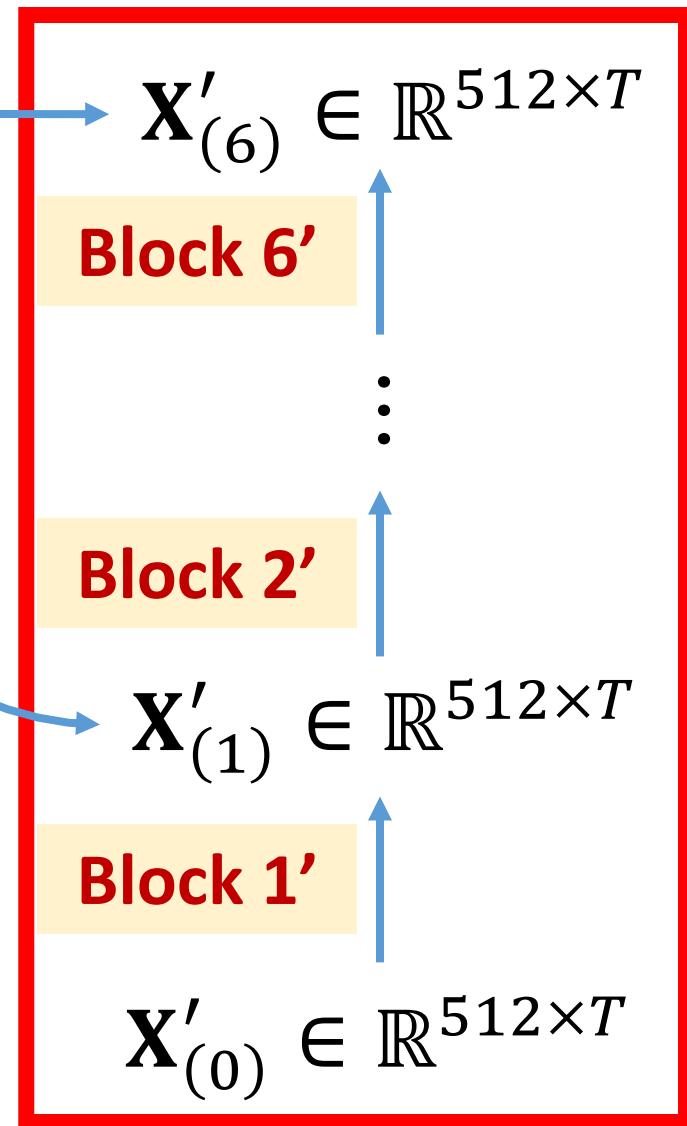


Decoder Network

Encoder

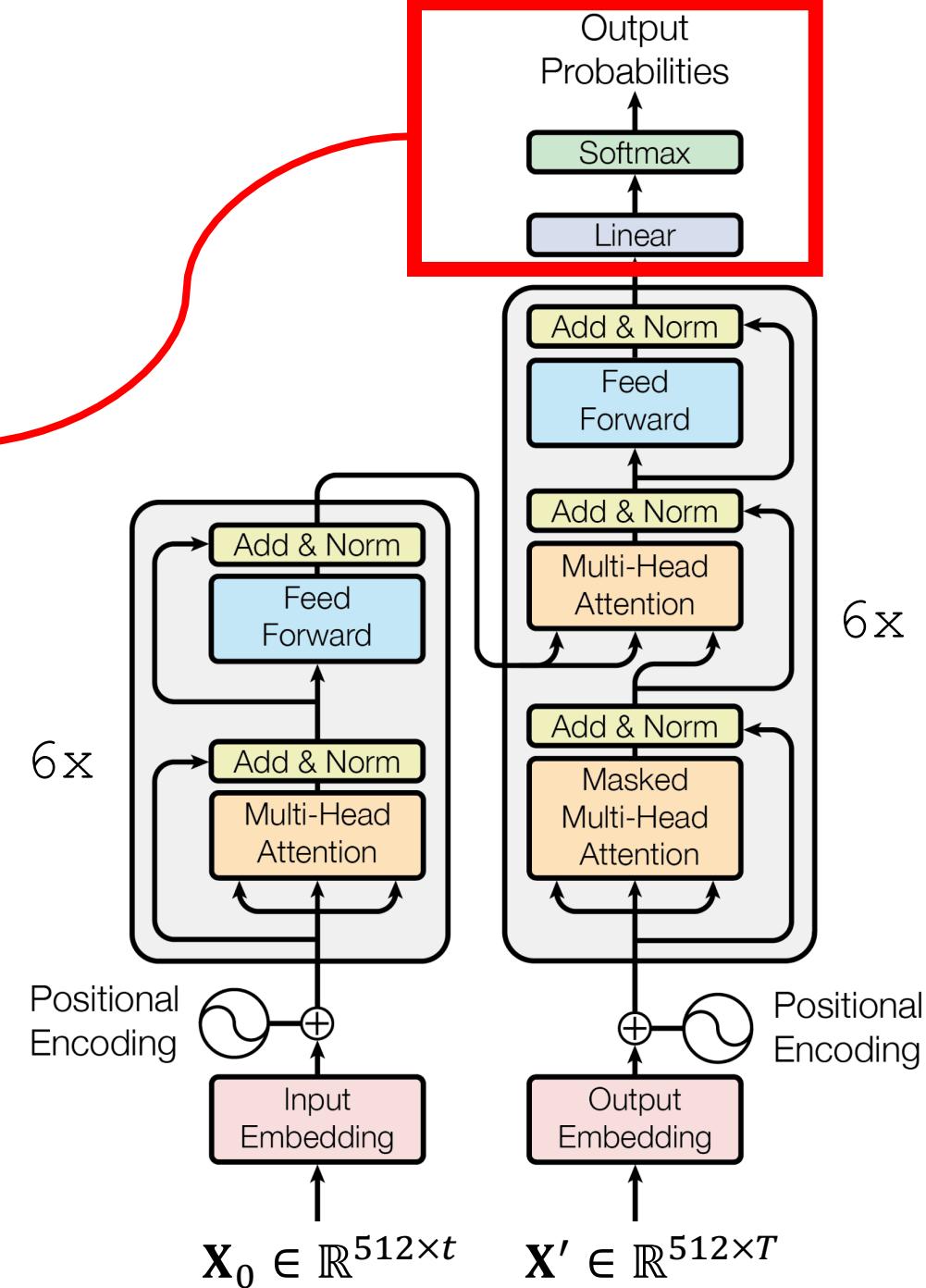


Decoder



Decoder Network

- Output a distribution over the vocabulary.
- Sample the next word according to the distribution.
- Append the new word's embedding to \mathbf{X}' .
- Run the decoder again, taking $\mathbf{X}' \in \mathbb{R}^{512 \times (T+1)}$ as input.



Summary

Summary

- Transformer model is **not RNN**.
 - Transformer is based on **attention** and **self-attention**.
 - **Upside:** Outperform all the state-of-the-art RNN models.
 - **Downside:** Much more expensive than RNN models.
-
- Read the original paper: Vaswani et al. [Attention Is All You Need](#). In *NIPS*, 2017.
 - Google “*transformer model explained*” and read the articles.

Key Concept: Multi-Head Attention

- Inputs: query \mathbf{Q} , key \mathbf{K} , and value \mathbf{V} .
- Linear maps: $\tilde{\mathbf{Q}} = \mathbf{W}_q \mathbf{Q}$, $\tilde{\mathbf{K}} = \mathbf{W}_k \mathbf{K}$, and $\tilde{\mathbf{V}} = \mathbf{W}_v \mathbf{V}$.
- Single-head attention:

$$\mathbf{z} = \text{attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \tilde{\mathbf{V}} \cdot \text{softmax}(\tilde{\mathbf{K}}^T \tilde{\mathbf{Q}}).$$

Key Concept: Multi-Head Attention

- Inputs: **query \mathbf{Q}** , **key \mathbf{K}** , and **value \mathbf{V}** .
- Linear maps: $\tilde{\mathbf{Q}} = \mathbf{W}_q \mathbf{Q}$, $\tilde{\mathbf{K}} = \mathbf{W}_k \mathbf{K}$, and $\tilde{\mathbf{V}} = \mathbf{W}_v \mathbf{V}$.

- Single-head attention:

$$\mathbf{z} = \text{attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \tilde{\mathbf{V}} \cdot \text{softmax}(\tilde{\mathbf{K}}^T \tilde{\mathbf{Q}}).$$

- Multi-head attention:

- Repeat $\text{attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$ using different parameters $\mathbf{W}_q, \mathbf{W}_v, \mathbf{W}_v$.
- Get $\mathbf{z}^{[1]}, \mathbf{z}^{[2]}, \dots, \mathbf{z}^{[m]} \in \mathbb{R}^{d_z \times t}$.
- Concatenate the m matrices to get $\tilde{\mathbf{Z}} \in \mathbb{R}^{md_z \times t}$.

Attention in the encoder:

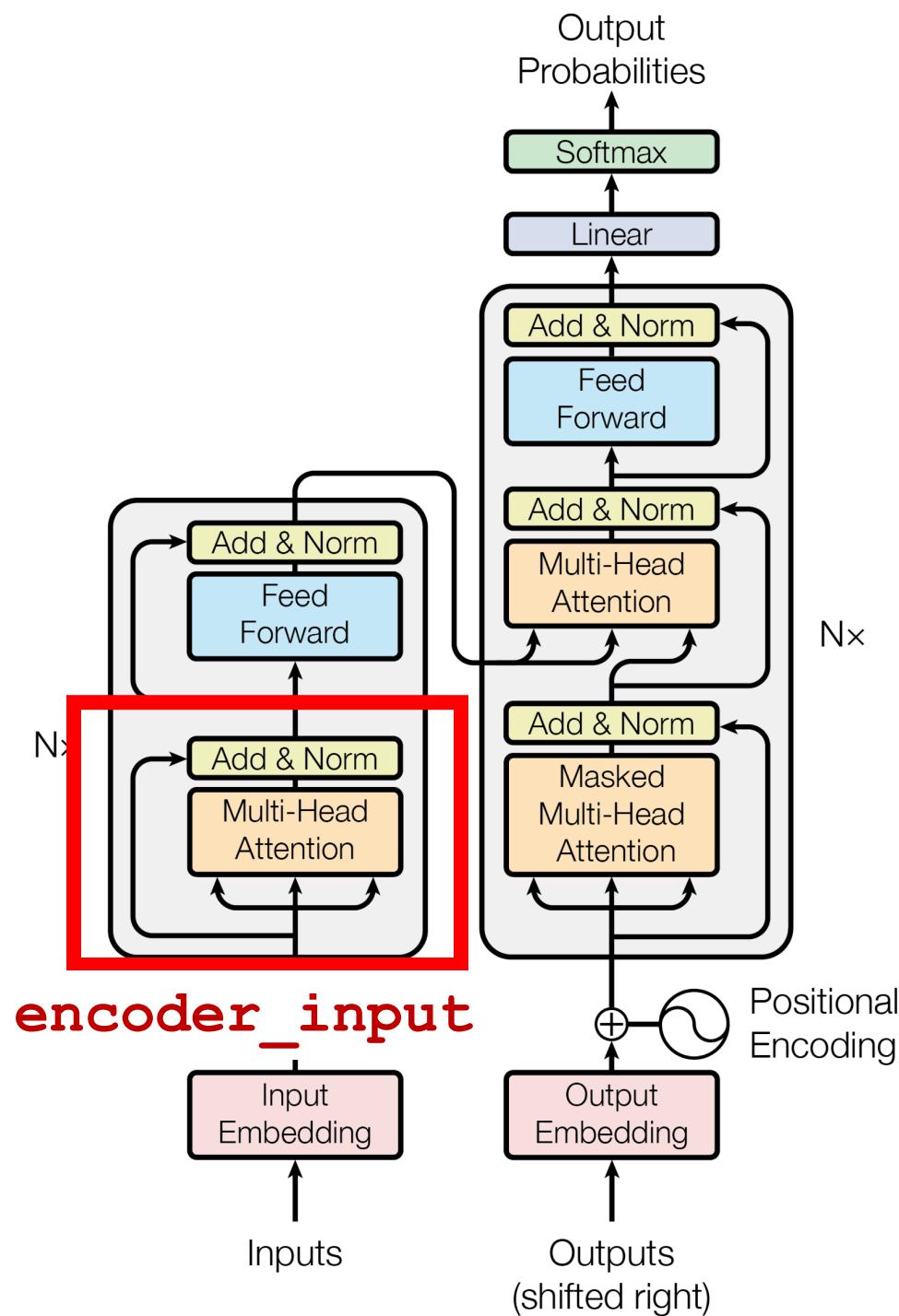
- $Q = K = V = \text{encoder_input}$.

1st attention in the decoder:

- $Q = K = V = \text{decoder_input}$.

2nd attention in the decoder

- $Q = \text{decoder_input}$
- $K = V = \text{encoder_output}$.



Attention in the encoder:

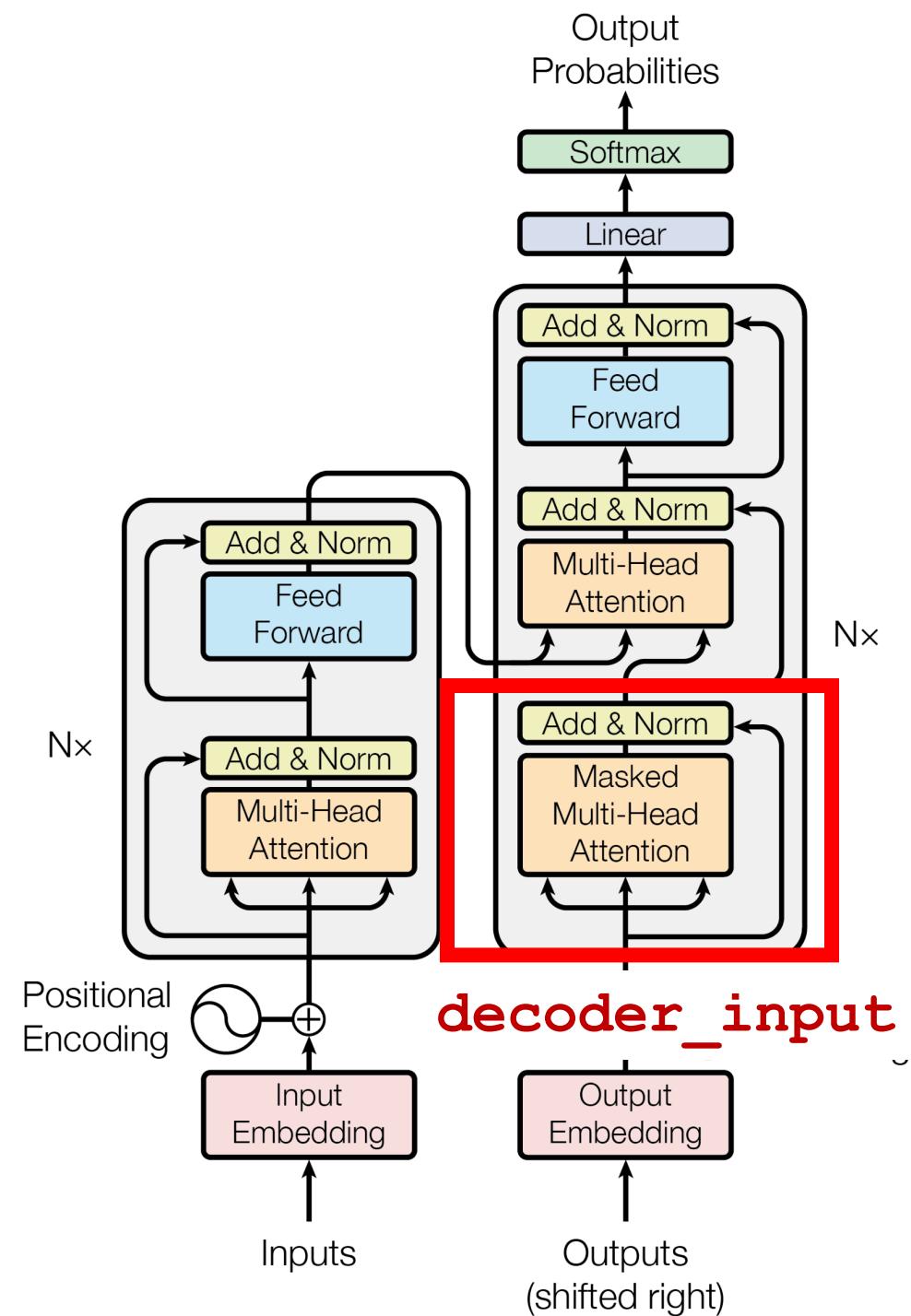
- $Q = K = V = \text{encoder_input}$.

1st attention in the decoder:

- $Q = K = V = \text{decoder_input}$.

2nd attention in the decoder

- $Q = \text{decoder_input}$
- $K = V = \text{encoder_output}$.



Attention in the encoder:

- $Q = K = V = \text{encoder_input}$.

1st attention in the decoder:

- $Q = K = V = \text{decoder_input}$.

2nd attention in the decoder

- $Q = \text{decoder_input}$
- $K = V = \text{encoder_output}$.

encoder_output

