

# 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

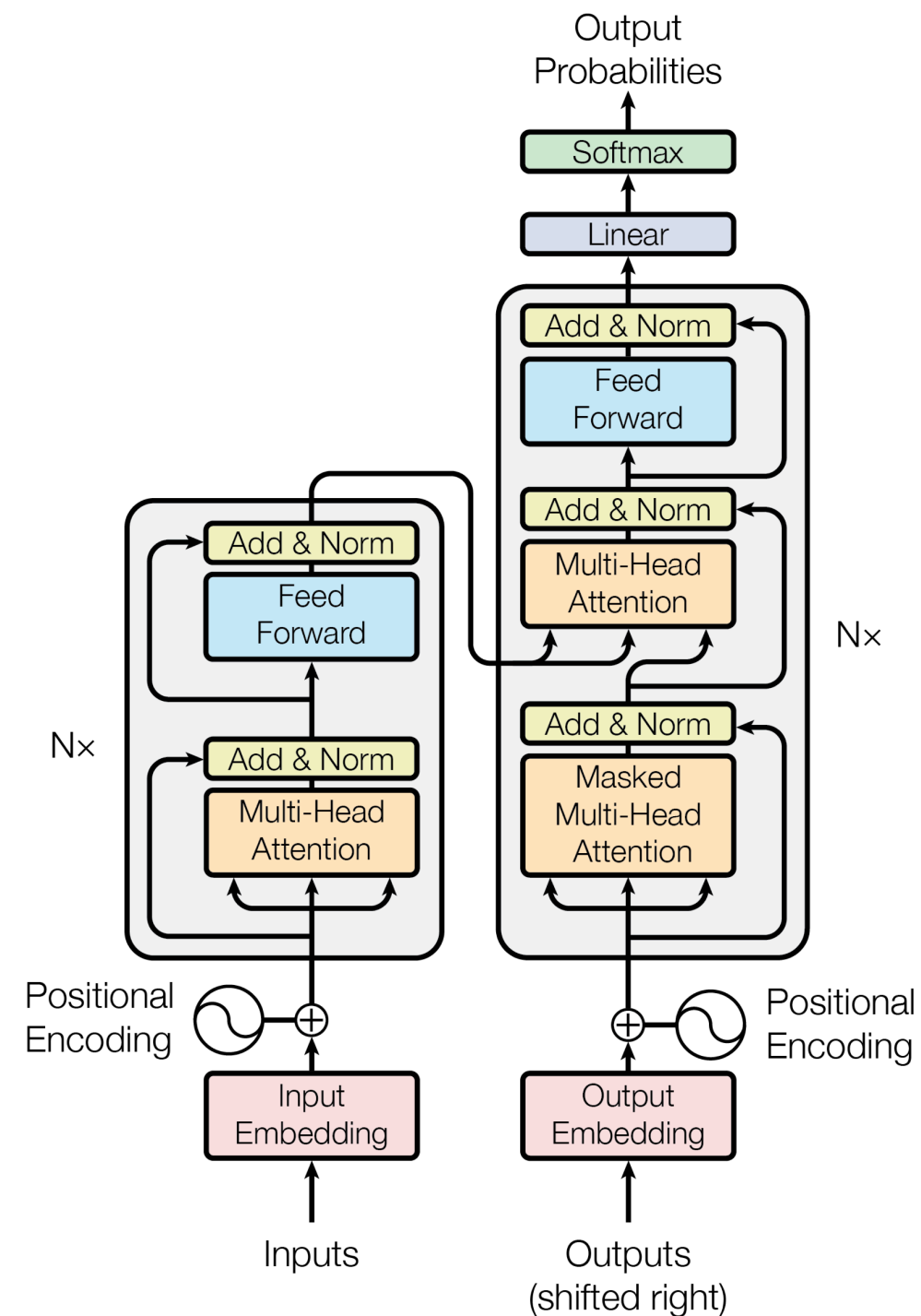
**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

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

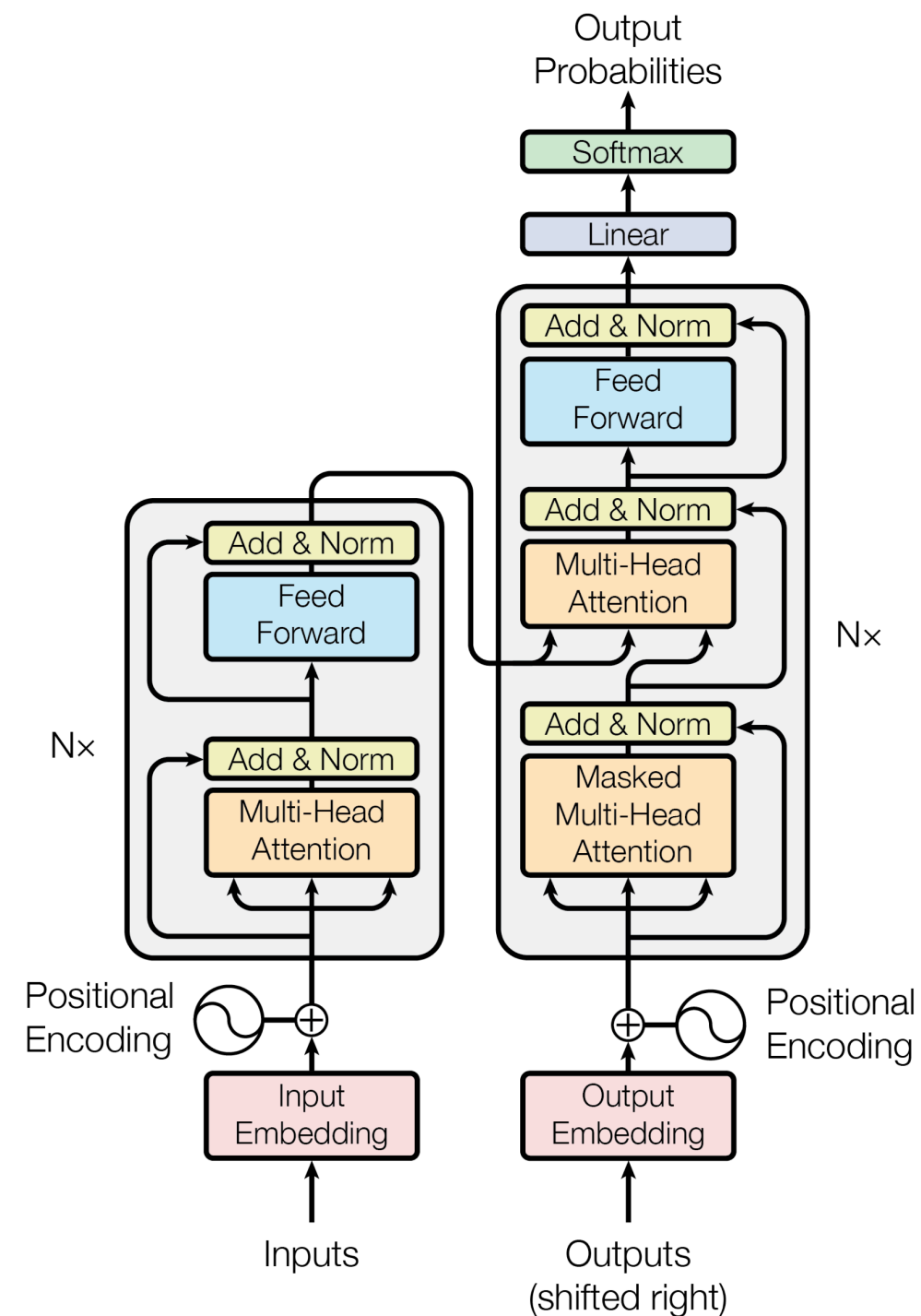
**Łukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com



# Transformer Model

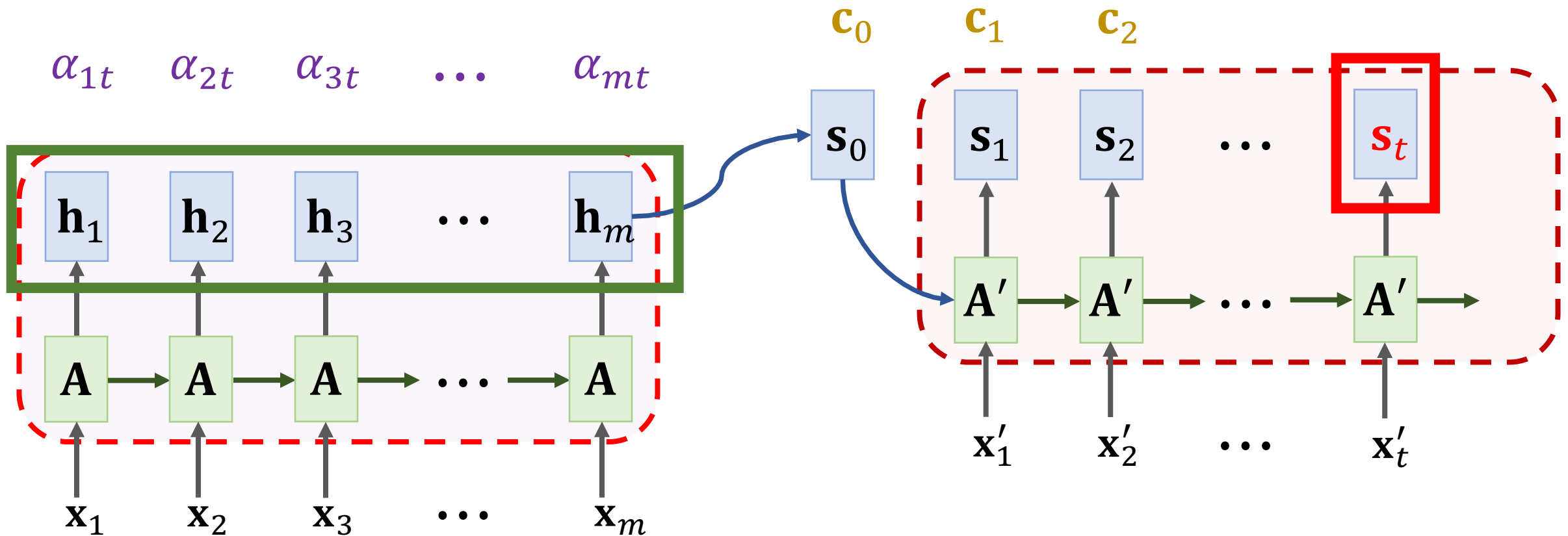
- Transformer is a Seq2Seq model.
- Transformer is not RNN.
- Purely based attention and fully-connected layers.
- Much more computations than RNNs.
- Higher performance than RNNs on large datasets.



# Revisit Attention

# Attention for RNN

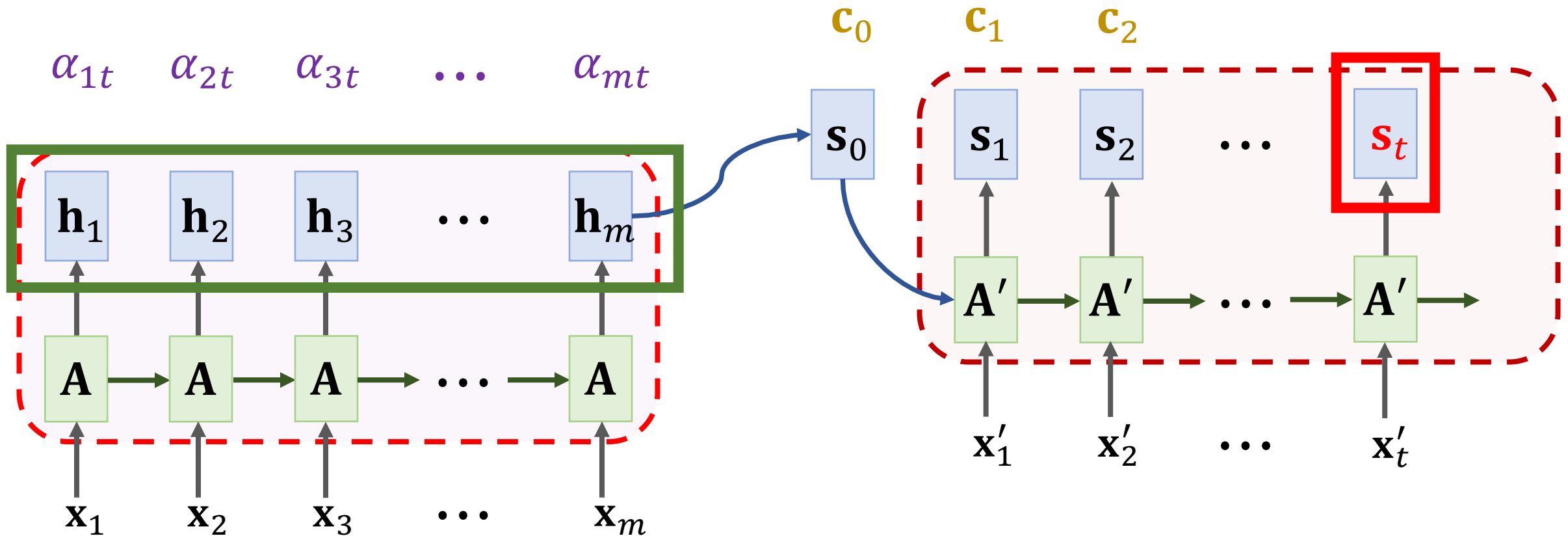
Weights:  $\alpha_{it} = \text{similarity}(\mathbf{h}_i, \mathbf{s}_t)$



# Attention for RNN

**Weights:**  $\alpha_{it} = \text{similarity}(\mathbf{h}_i, \mathbf{s}_t)$

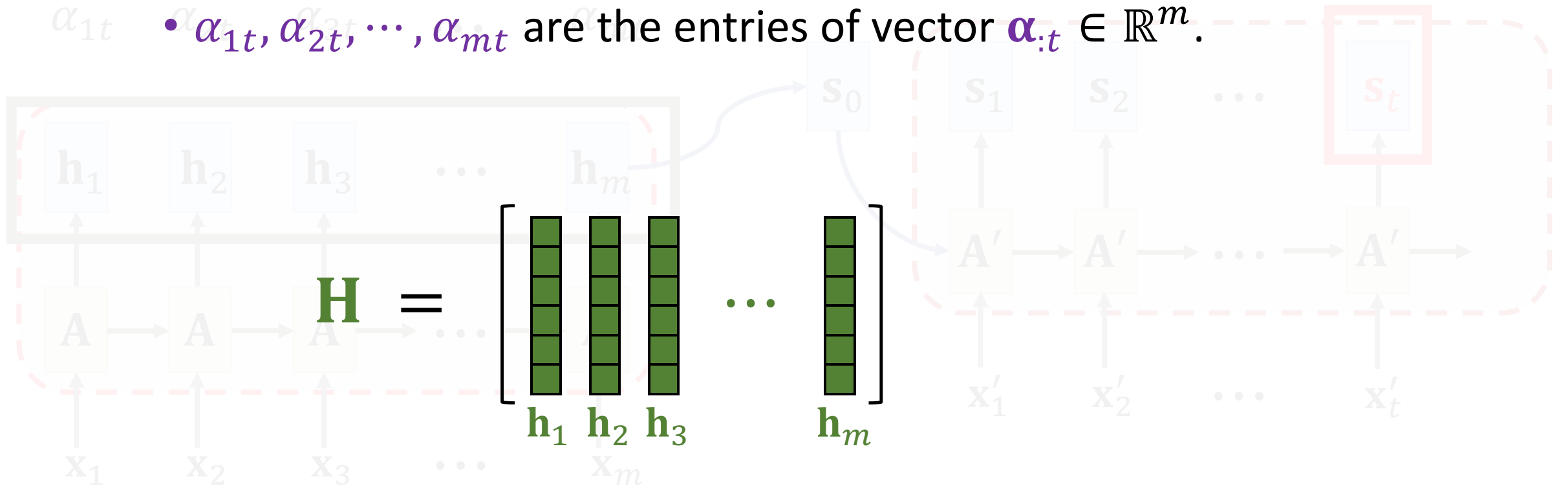
Here,  $\mathbf{h}_i$  is called “key”, and  $\mathbf{s}_t$  is called “query”.



# Attention for RNN

**Weights:**  $\alpha_{it} = \text{similarity}(\mathbf{h}_i, \mathbf{s}_t)$

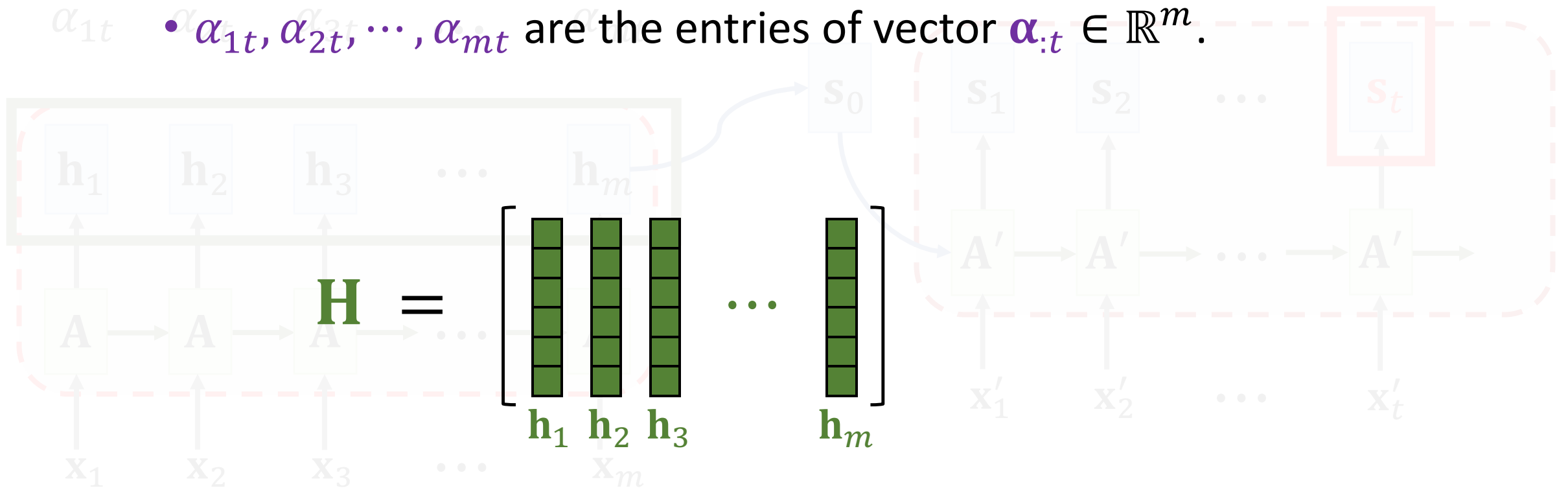
- Define  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_m] \in \mathbb{R}^{d \times m}$ .
- Compute weights:  $\boldsymbol{\alpha}_{:t} = \text{Softmax}\left((\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t)\right) \in \mathbb{R}^m$ .
- $\alpha_{1t}, \alpha_{2t}, \dots, \alpha_{mt}$  are the entries of vector  $\boldsymbol{\alpha}_{:t} \in \mathbb{R}^m$ .



# Attention for RNN

**Weights:**  $\alpha_{:t} = \text{Softmax} \left( (\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t) \right)$

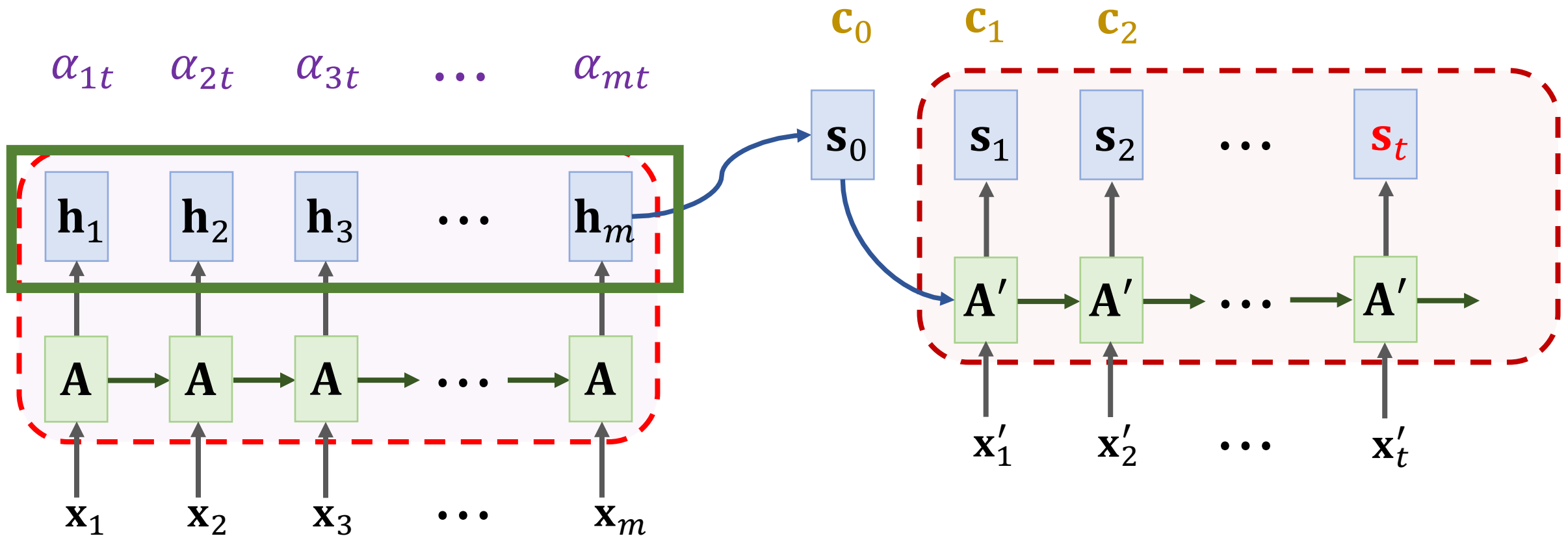
- Define  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_m] \in \mathbb{R}^{d \times m}$ .
- Compute weights:  $\alpha_{:t} = \text{Softmax} \left( (\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t) \right) \in \mathbb{R}^m$ .
- $\alpha_{1t}, \alpha_{2t}, \dots, \alpha_{mt}$  are the entries of vector  $\alpha_{:t} \in \mathbb{R}^m$ .





# Attention for RNN

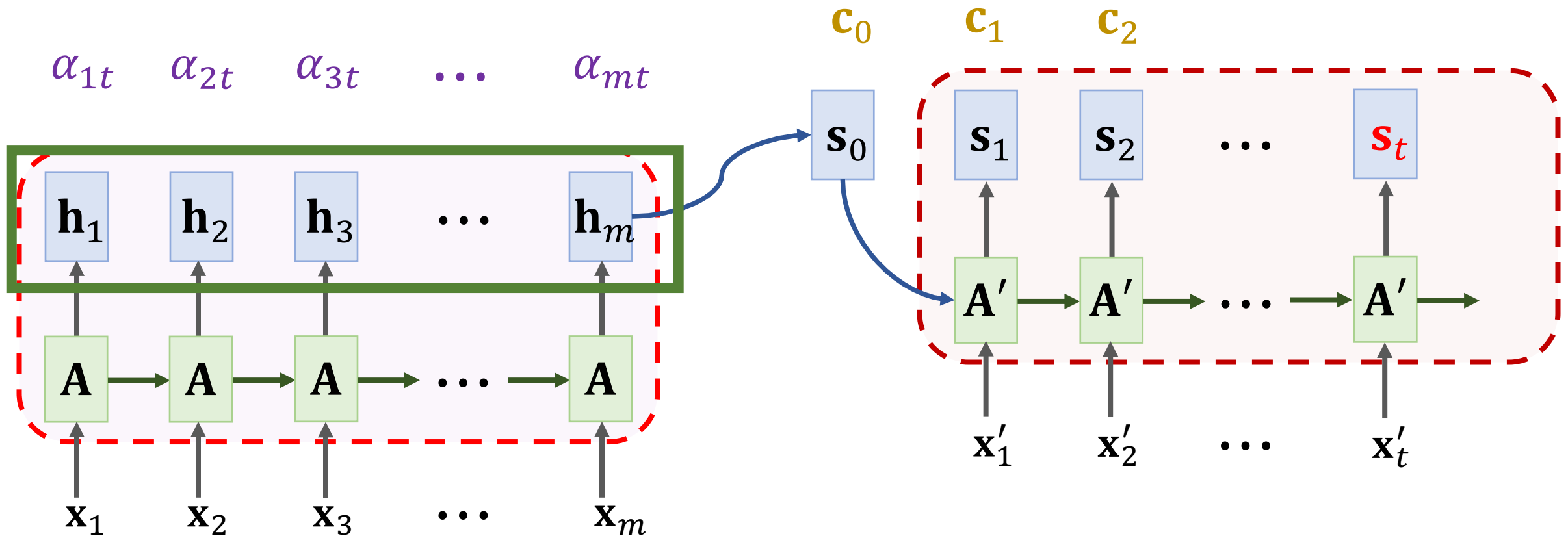
Weights:  $\alpha_{:,t} = \text{Softmax} \left( (\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t) \right)$



# Attention for RNN

**Weights:**  $\alpha_{:t} = \text{Softmax} \left( (\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t) \right)$

**Context vector:**  $\mathbf{c}_t = \alpha_{1t} \mathbf{h}_1 + \dots + \alpha_{mt} \mathbf{h}_m.$



# Attention for RNN

**Weights:**  $\alpha_{:t} = \text{Softmax} \left( (\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t) \right)$

**Context vector:**  $\mathbf{c}_t = \alpha_{1t} \mathbf{h}_1 + \dots + \alpha_{mt} \mathbf{h}_m.$

- Note that  $\mathbf{H} \alpha_{:t} = \alpha_{1t} \mathbf{h}_1 + \dots + \alpha_{mt} \mathbf{h}_m.$
- Thus  $\mathbf{c}_t = \alpha_{1t} \mathbf{h}_1 + \dots + \alpha_{mt} \mathbf{h}_m = \mathbf{H} \alpha_{:t}.$

$$\mathbf{H} = \begin{bmatrix} \begin{array}{c} \text{green box} \\ \text{green box} \\ \text{green box} \\ \text{green box} \\ \text{green box} \end{array} & \begin{array}{c} \text{green box} \\ \text{green box} \\ \text{green box} \\ \text{green box} \\ \text{green box} \end{array} & \begin{array}{c} \text{green box} \\ \text{green box} \\ \text{green box} \\ \text{green box} \\ \text{green box} \end{array} & \dots & \begin{array}{c} \text{green box} \\ \text{green box} \\ \text{green box} \\ \text{green box} \\ \text{green box} \end{array} \end{bmatrix}$$

$\mathbf{h}_1 \quad \mathbf{h}_2 \quad \mathbf{h}_3 \qquad \qquad \mathbf{h}_m$

# Attention for RNN

**Weights:**  $\alpha_{:t} = \text{Softmax} \left( (\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t) \right)$

**Context vector:**  $\mathbf{c}_t = \mathbf{H} \alpha_{:t}.$

# Attention for RNN

**Weights:**  $\alpha_{:t} = \text{Softmax} \left( (\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t) \right)$

**Context vector:**  $\mathbf{c}_t = \mathbf{W}_V \mathbf{H} \alpha_{:t}$ .

A different way to computing context vector.

# **Single-Head & Multi-Head Attention**

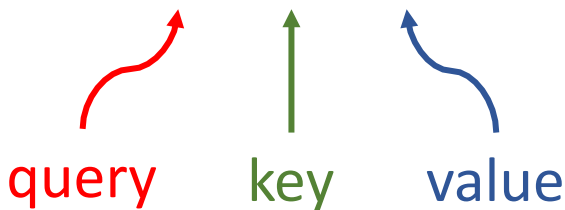
# Attention for RNN

**Context vector:**  $\mathbf{c}_t = (\mathbf{W}_V \mathbf{H}) \cdot \text{Softmax} \left( (\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t) \right).$

# Single-Head Attention

**Context vector:**  $\mathbf{c}_t = (\mathbf{W}_V \mathbf{H}) \cdot \text{Softmax} \left( (\mathbf{W}_K \mathbf{H})^T (\mathbf{W}_Q \mathbf{s}_t) \right).$

## Single-Head Attention:

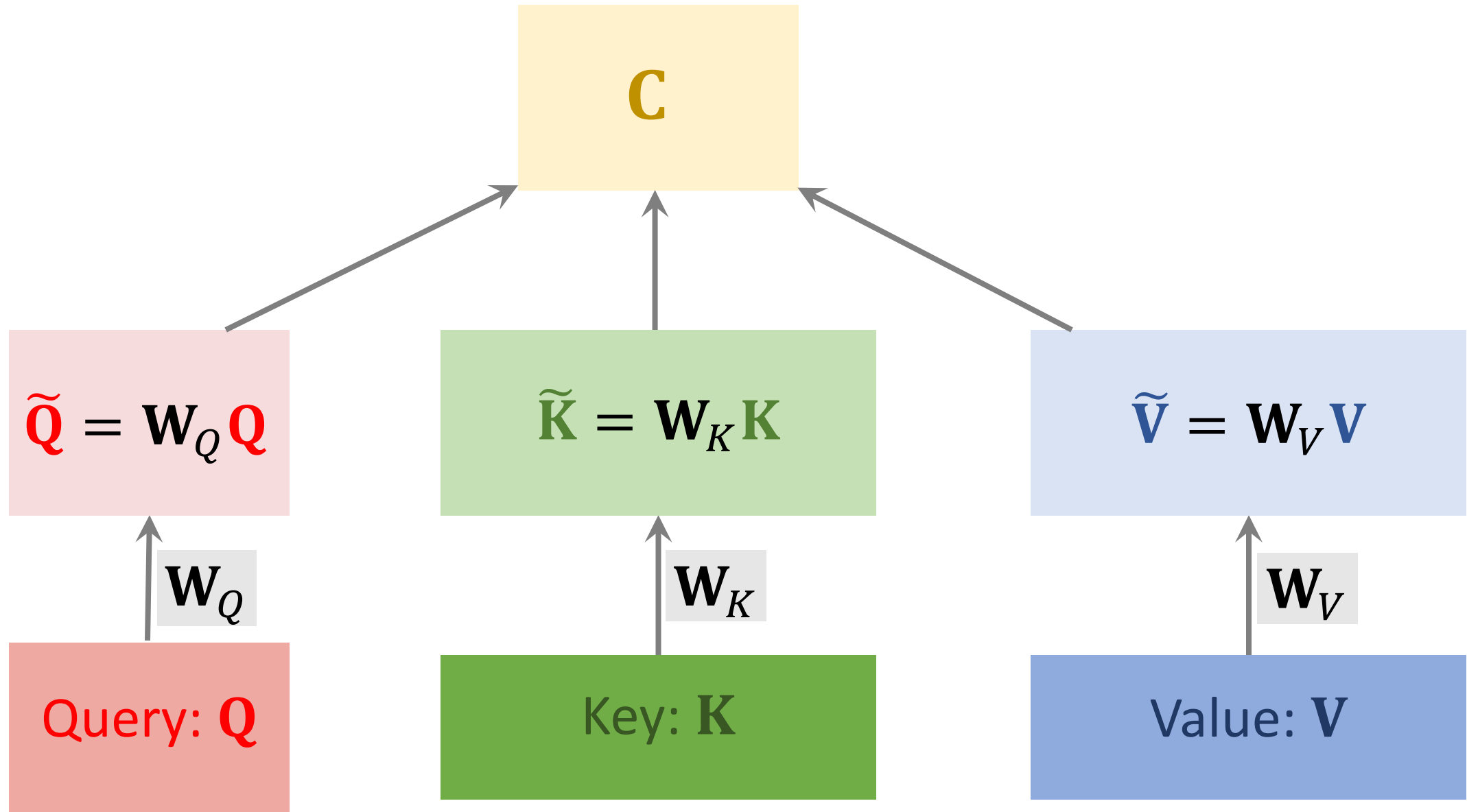
- $\mathbf{C} = \text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}).$   


- The  $t$ -th column of  $\mathbf{C}$  is:

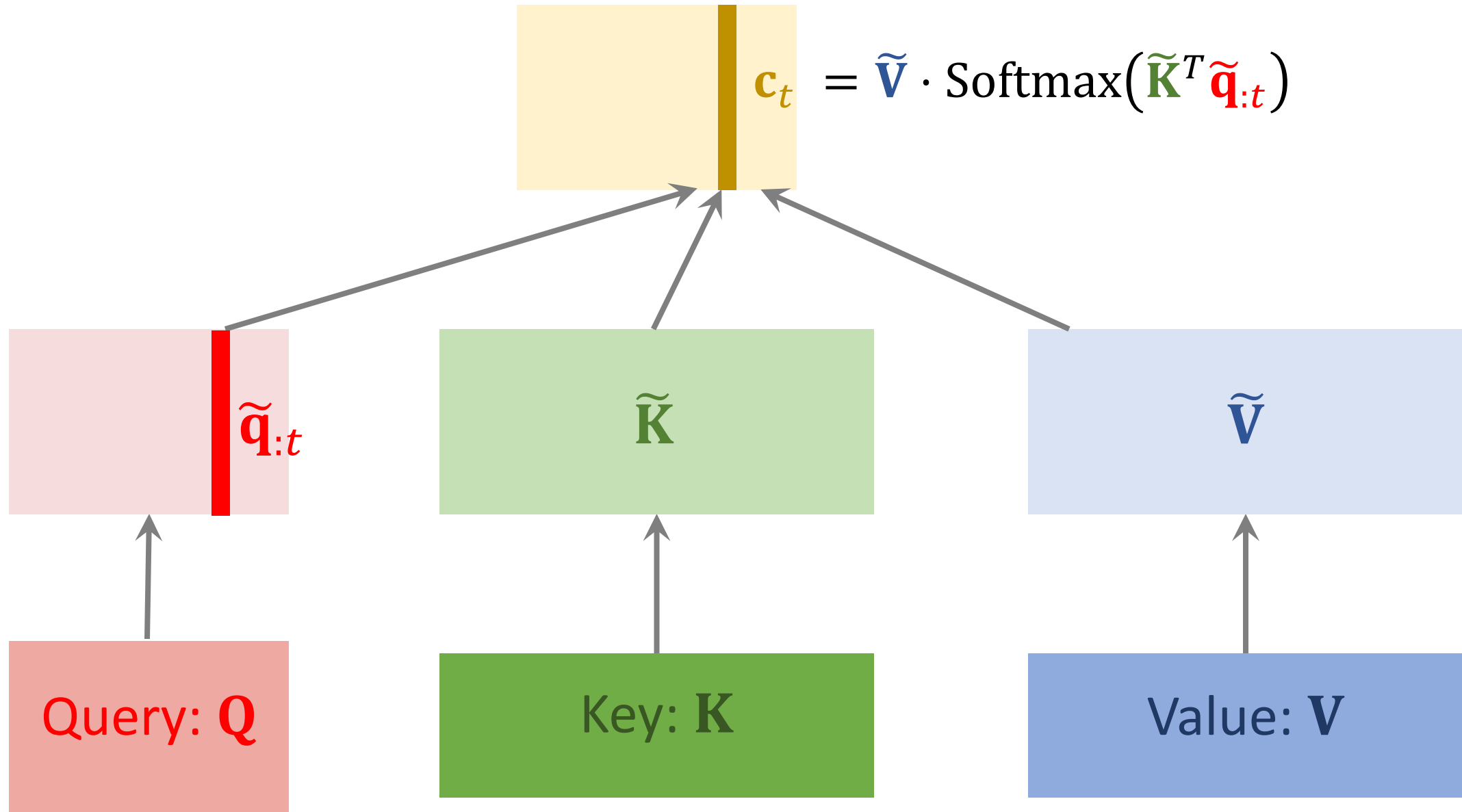
$$\mathbf{c}_t = (\mathbf{W}_V \mathbf{V}) \cdot \text{Softmax} \left( (\mathbf{W}_K \mathbf{K})^T (\mathbf{W}_Q \mathbf{q}_{:t}) \right).$$



# Single-Head Attention



# Single-Head Attention



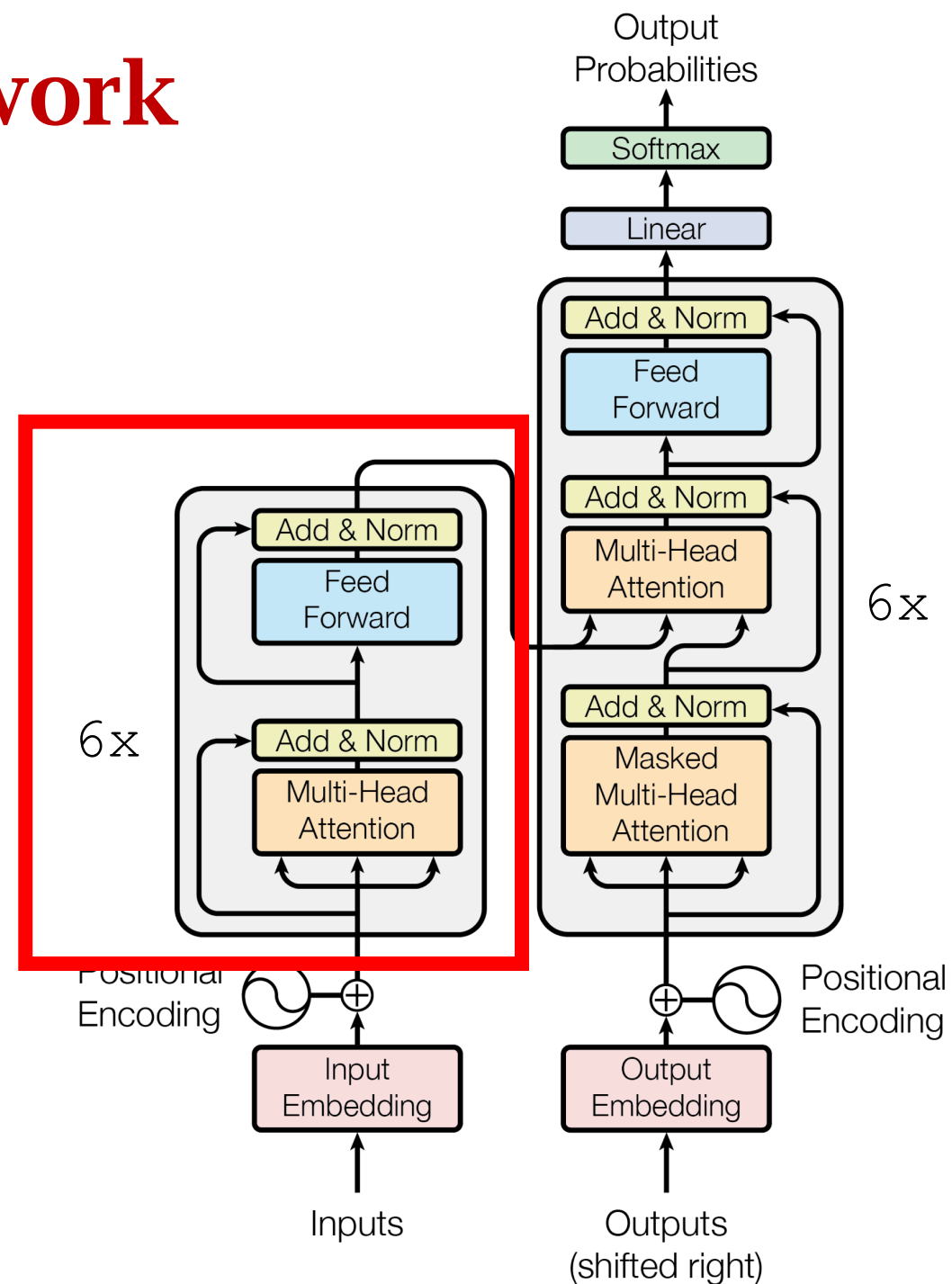
# Multi-Head Attention

- **Single-head attention:**  $\mathbf{C} = \text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$ .
  - Its  $t$ -th column is  $\mathbf{c}_t = (\mathbf{W}_V \mathbf{V}) \cdot \text{Softmax}\left((\mathbf{W}_K \mathbf{K})^T (\mathbf{W}_Q \mathbf{q}_{:t})\right)$ .
  - $\mathbf{W}_Q$ ,  $\mathbf{W}_K$ , and  $\mathbf{W}_V$  are trainable parameters.
  - Output shape:  $d_c \times t$ . (Matrix  $\mathbf{Q}$  has  $t$  columns.)
- **Multi-head attention:**
  - Using  $l$  single-head attentions (which do not share parameters.)
  - Totally  $3l$  parameters matrices  $\mathbf{W}$ .
  - Concatenating the output  $\mathbf{C}$  matrices of the single-head attentions.
  - Output shape:  $(ld_c) \times t$ .

# 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

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

Concatenation

$$\mathbf{C}^{[1]} \in \mathbb{R}^{d_c \times t}$$

$$\mathbf{C}^{[2]} \in \mathbb{R}^{d_c \times t}$$

...

$$\mathbf{C}^{[l]} \in \mathbb{R}^{d_c \times t}$$

Attention w/ different  
parameter matrices.

(**Q**, **K**, **V**)

# Multi-Head Attention + Dense Layer

$\tilde{\mathbf{C}}$  's number of columns,  $t$ , is determined by  $\mathbf{Q}$ .

## Multi-Head Attention

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

Concatenation

$$\mathbf{C}^{[1]} \in \mathbb{R}^{d_c \times t}$$

$$\mathbf{C}^{[2]} \in \mathbb{R}^{d_c \times t}$$

...

$$\mathbf{C}^{[l]} \in \mathbb{R}^{d_c \times t}$$

Attention w/ different  
parameter matrices.

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

# 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.

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

Concatenation

$$\mathbf{C}^{[1]} \in \mathbb{R}^{d_c \times t}$$

$$\mathbf{C}^{[2]} \in \mathbb{R}^{d_c \times t}$$

...

$$\mathbf{C}^{[l]} \in \mathbb{R}^{d_c \times t}$$

**Attention** w/ different parameter matrices.

(**Q**, **K**, **V**)



# Encoder Network: One Block

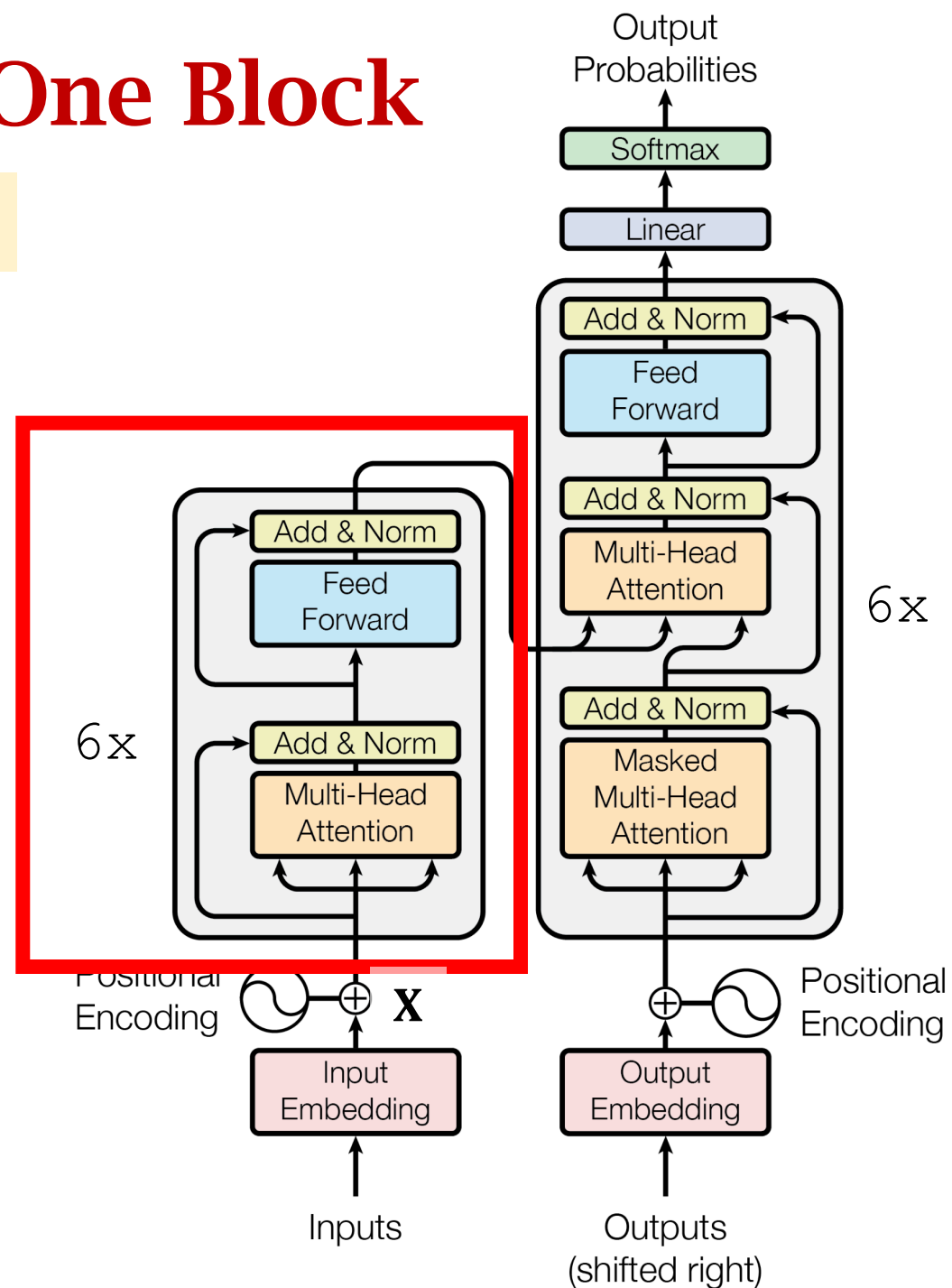
Ignore skip connection and normalization.

- Input:  $\mathbf{X} \in \mathbb{R}^{512 \times m}$ ; ( $m$  is the seq length.)

- Set  $\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{X}$ .

query      key      value

Similar to self-attention.

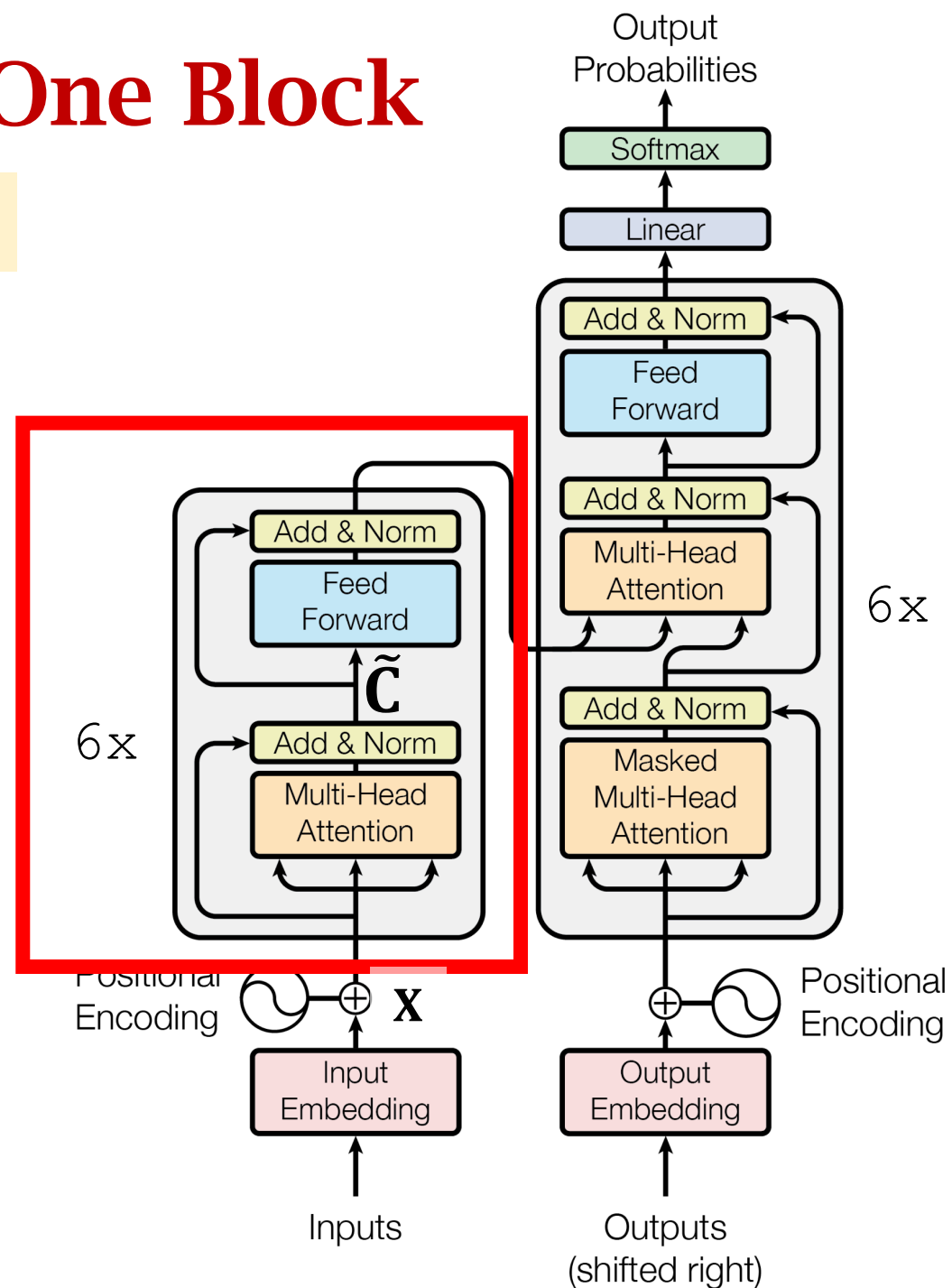


# Encoder Network: One Block

Ignore skip connection and normalization.

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

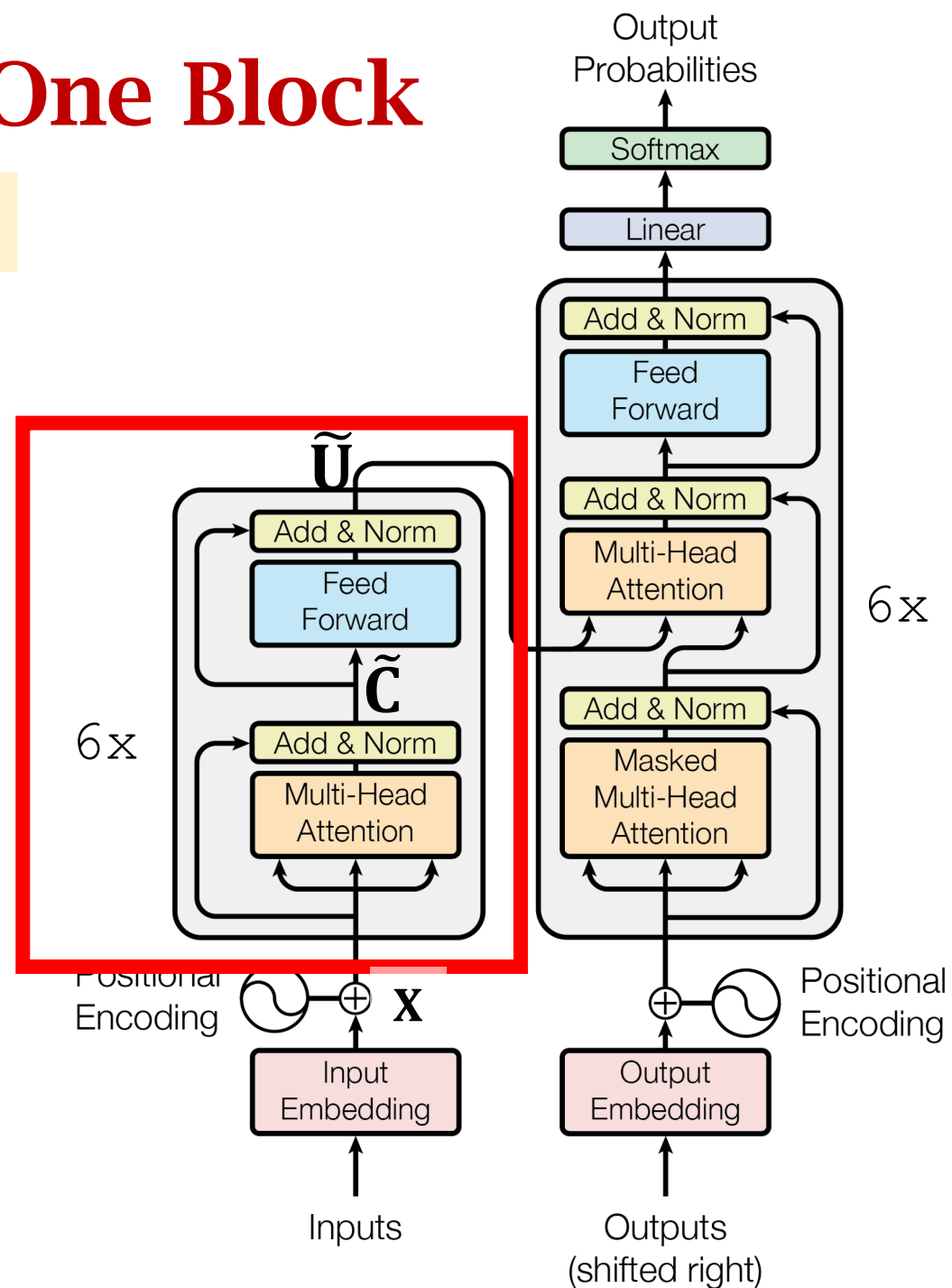
- Make sure the **input shape** and **output shape** are the same.
- Otherwise, skip connection cannot be applied.



# Encoder Network: One Block

Ignore skip connection and normalization.

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

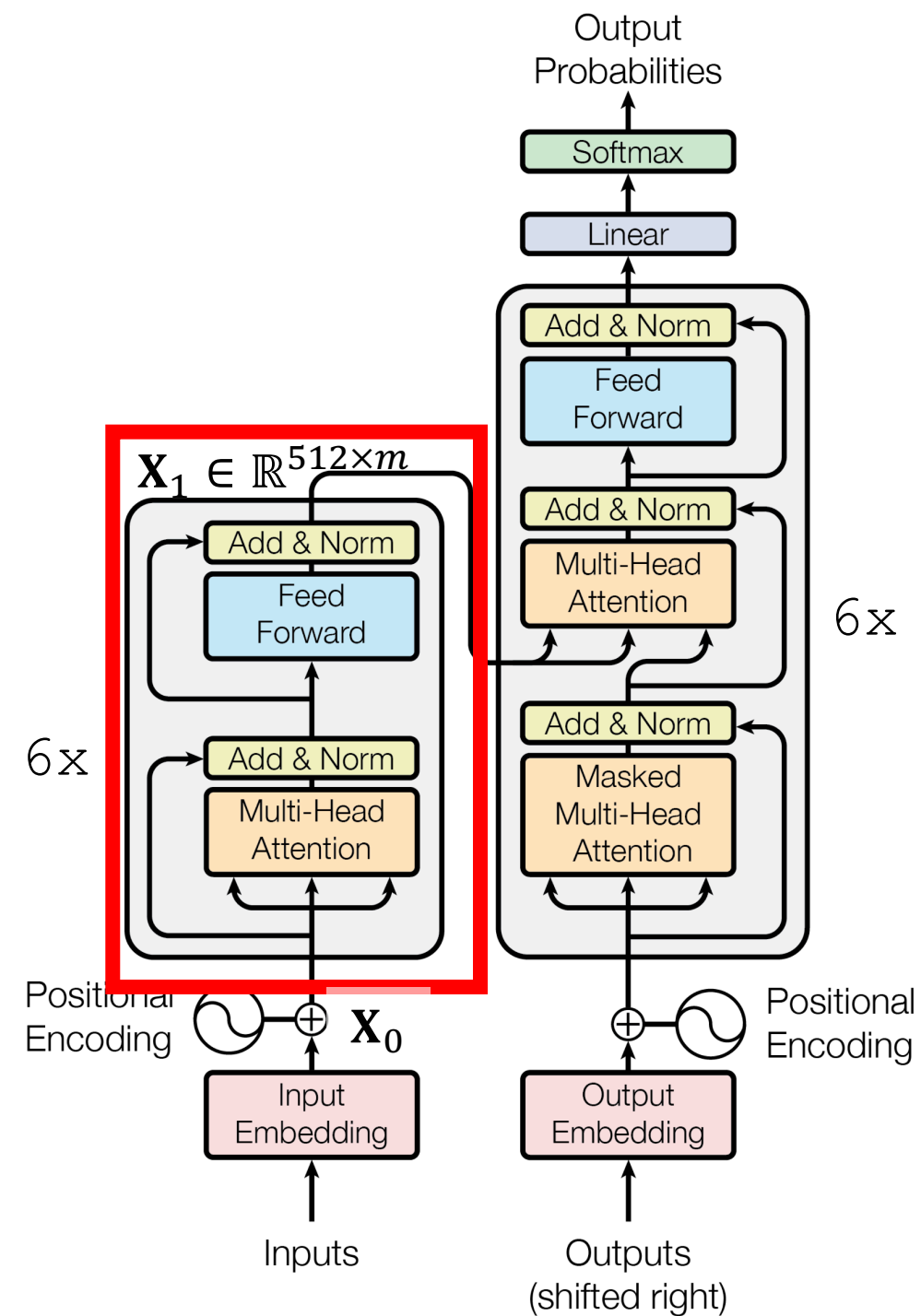


# Encoder Network

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

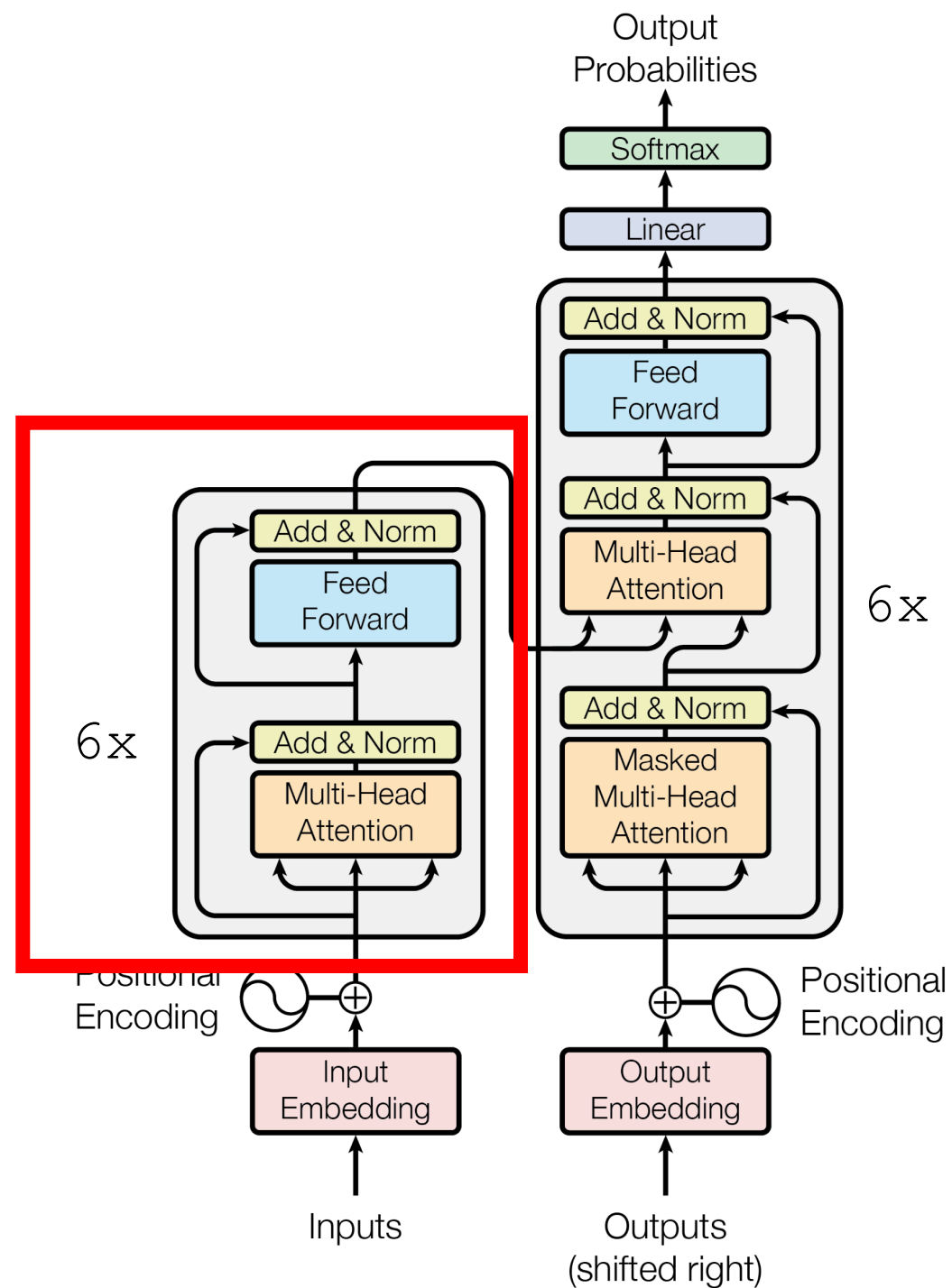
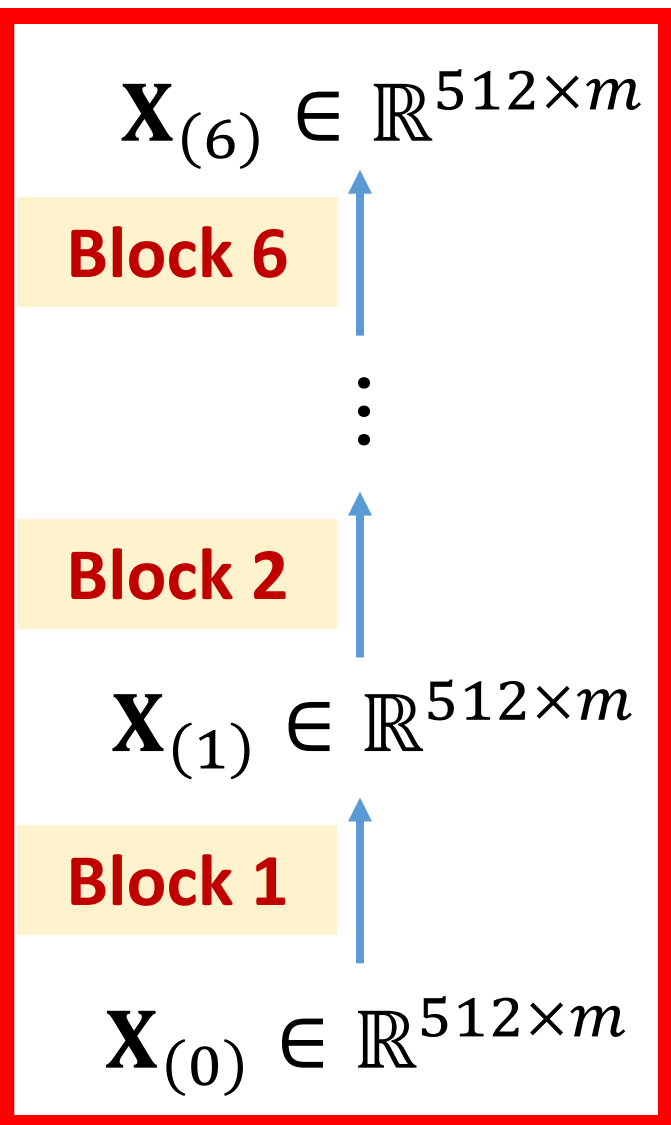
**Block 1**

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



# Encoder Network

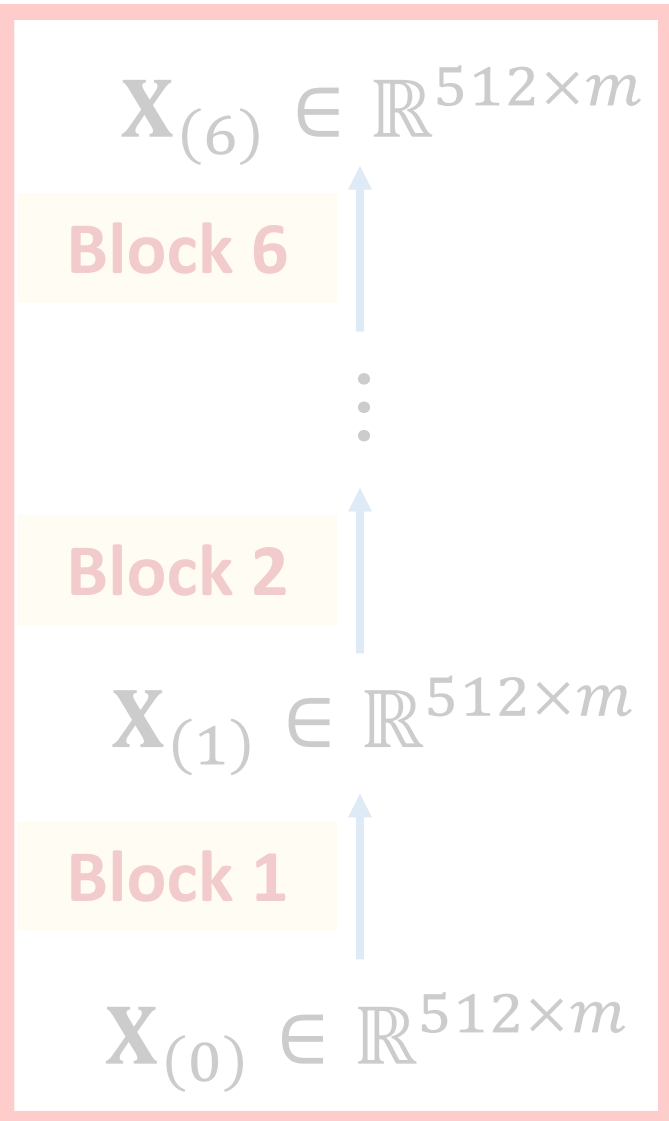
## Encoder



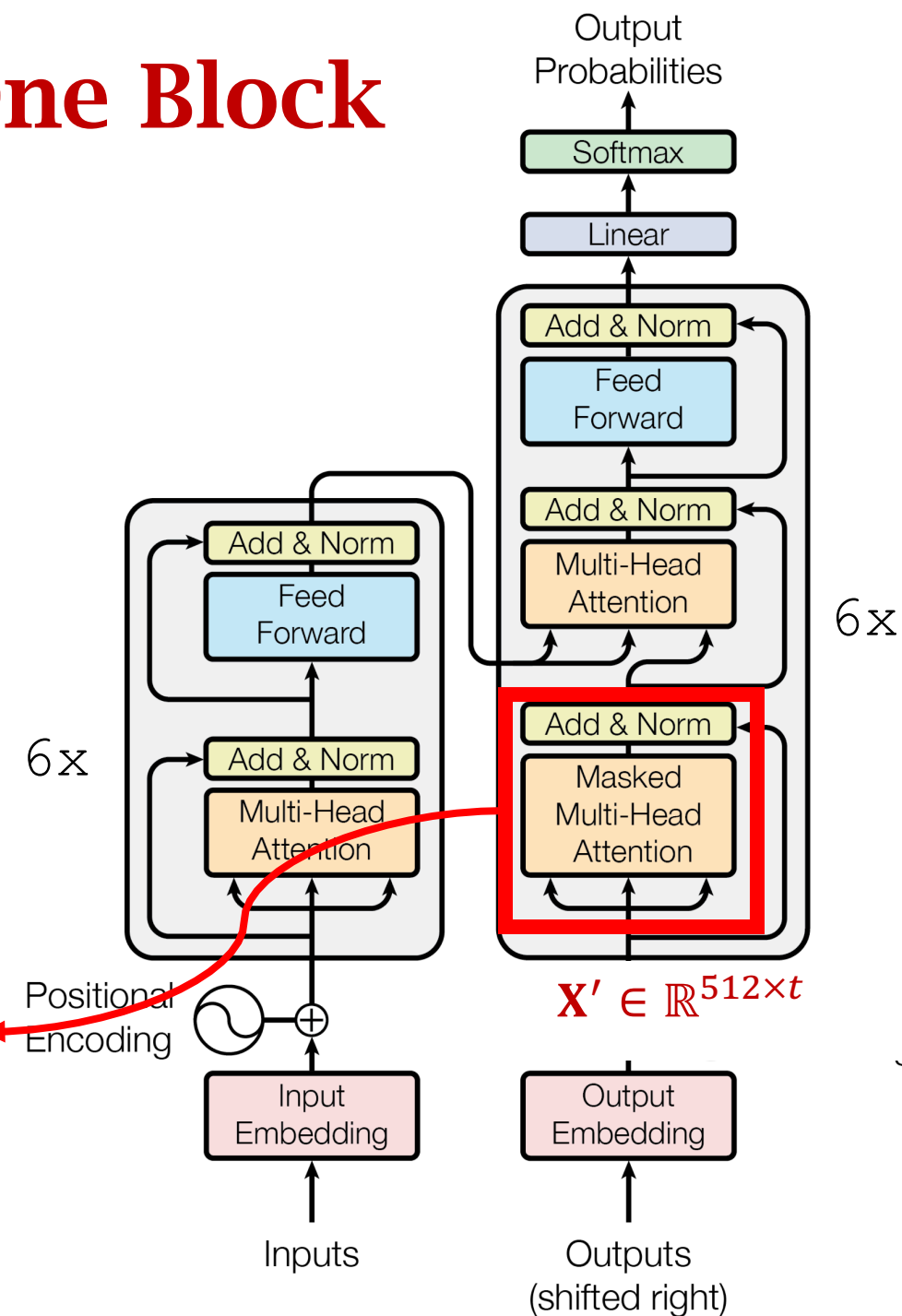
# Decoder of Transformer

# Decoder Network: One Block

## Encoder

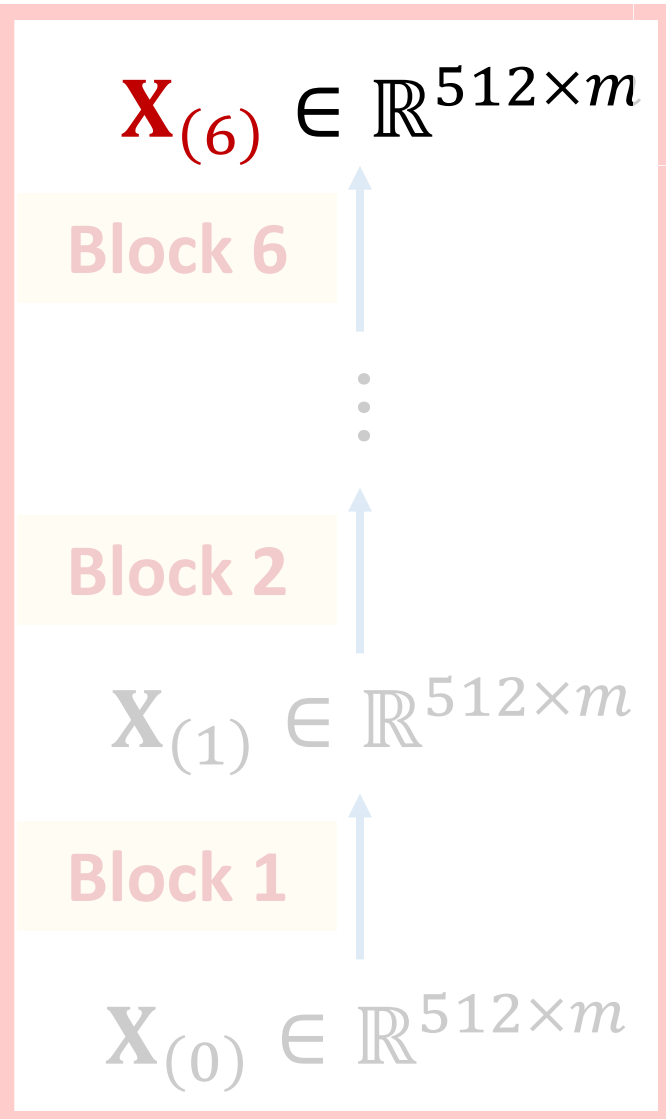


- Similar to encoder.
- Set  $\mathbf{Q} = \mathbf{K} = \mathbf{V} = \mathbf{X}'$ .



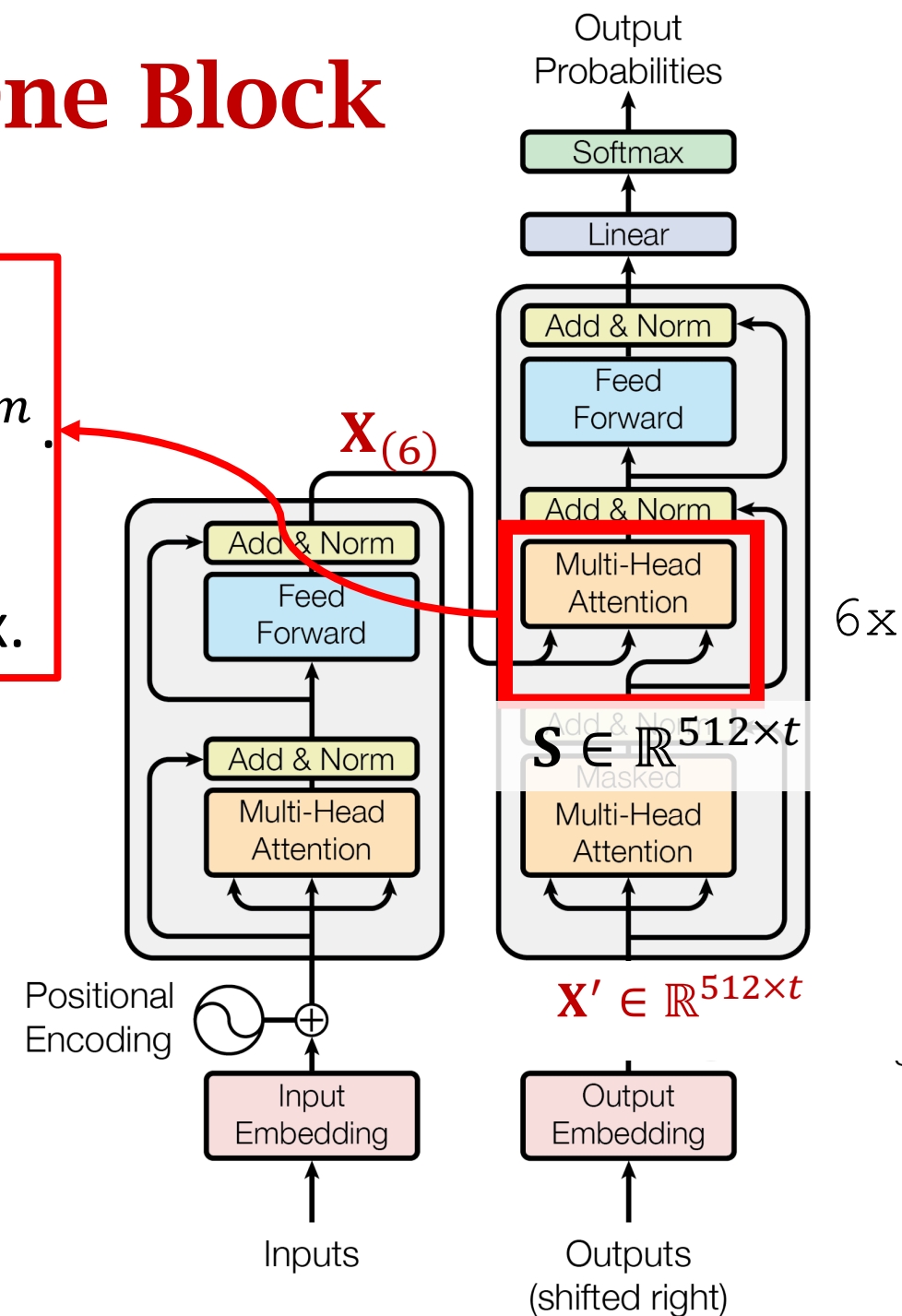
# Decoder Network: One Block

Encoder



- Set  $\mathbf{Q} = \mathbf{S} \in \mathbb{R}^{512 \times t}$ .
- $\mathbf{K} = \mathbf{V} = \mathbf{X}_{(6)} \in \mathbb{R}^{512 \times m}$ .
- Multi-head attention outputs a  $512 \times t$  matrix.

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

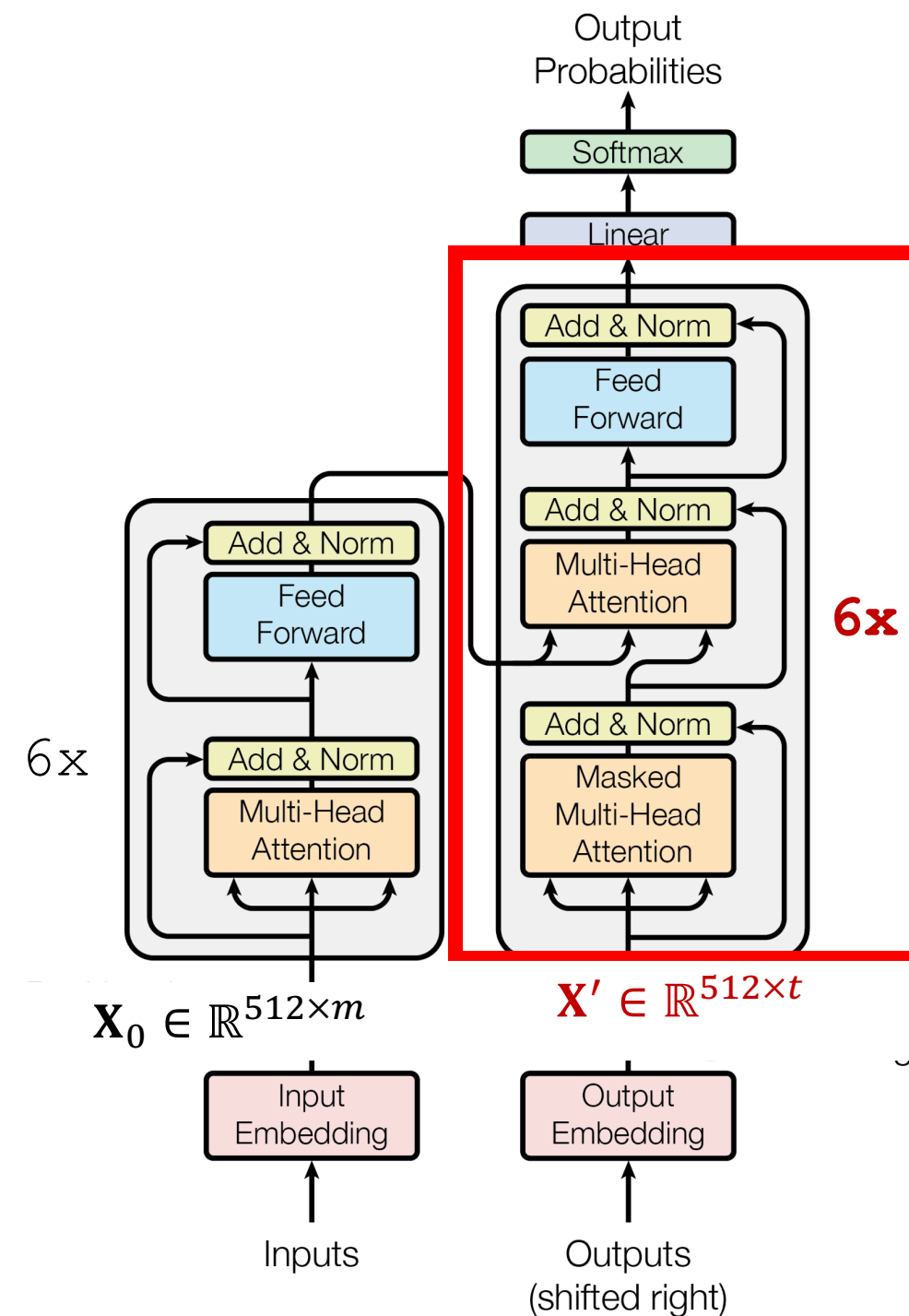
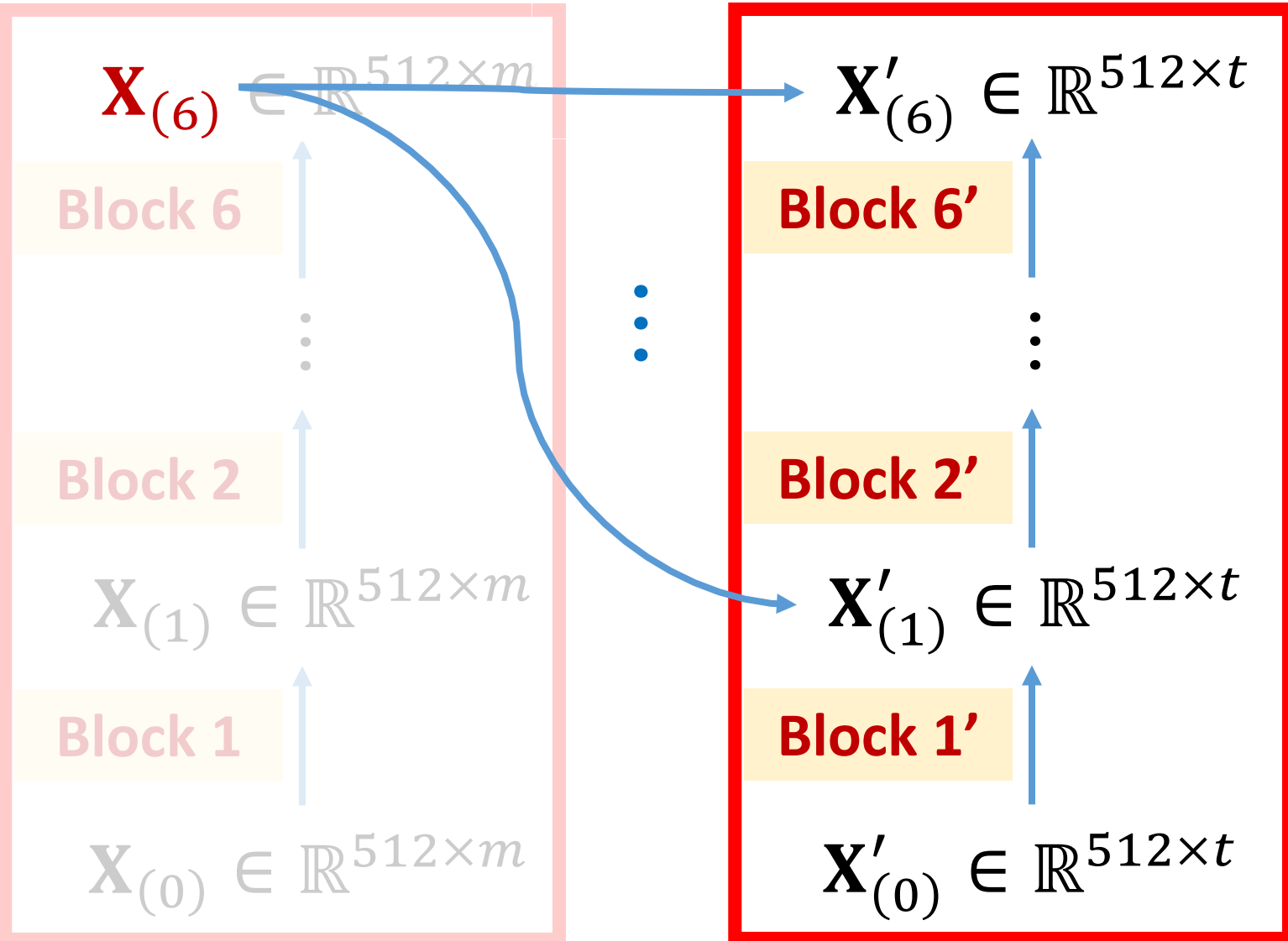




# Decoder Network

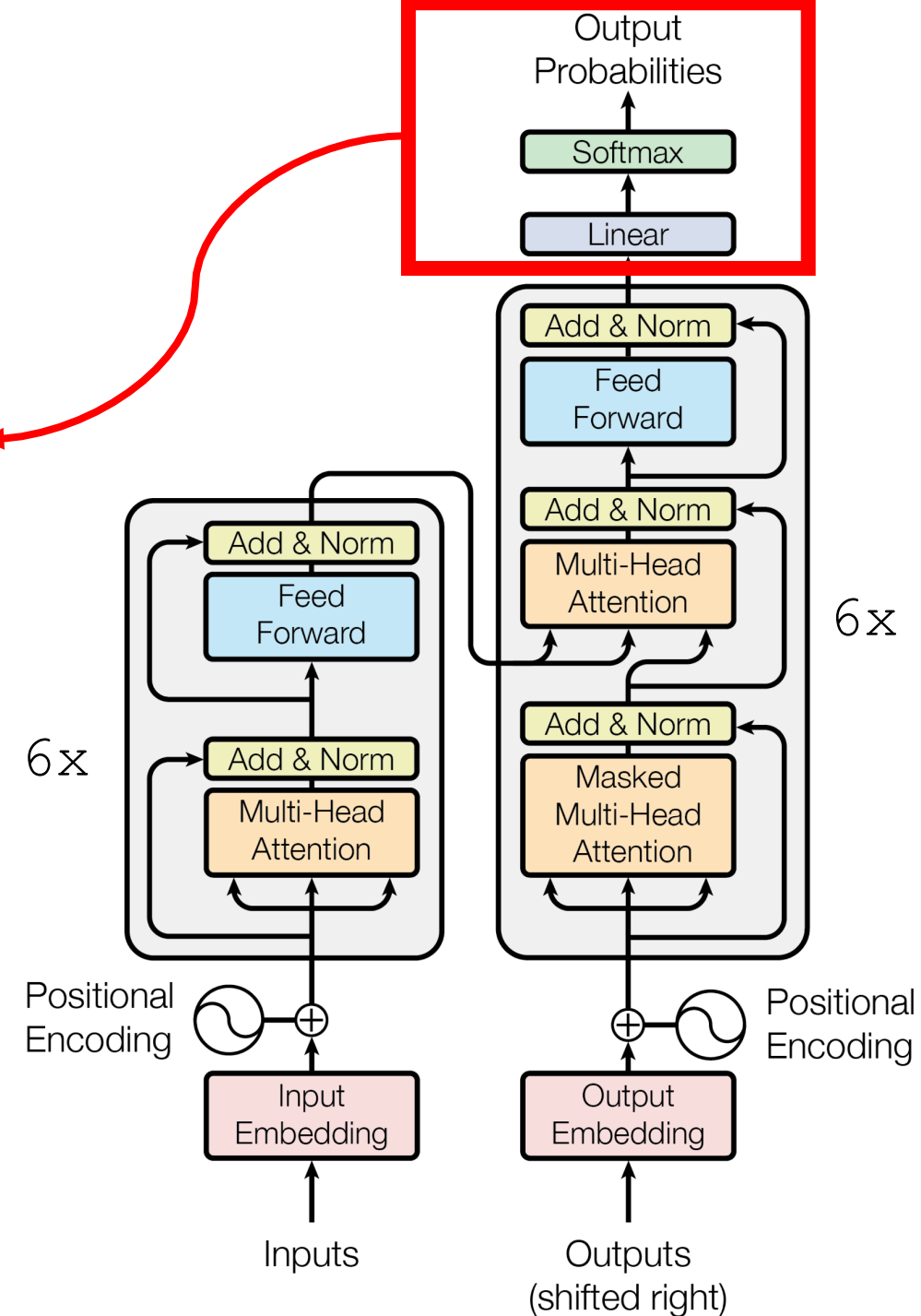
Encoder

Decoder



# Decoder Network

- Output a distribution over the vocabulary.
- Compare the distribution with the one-hot encode of the label.
- ➔ Loss, e.g., cross-entropy.
- ➔ Gradient.
- ➔ Update model parameters.



# 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:

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

# 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:

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

- Multi-head attention:
  - Repeat  $\text{attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$  using different parameters  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ .
  - Get  $\mathbf{C}^{[1]}, \mathbf{C}^{[2]}, \dots, \mathbf{C}^{[l]} \in \mathbb{R}^{d_c \times t}$ .
  - Concatenate the  $m$  matrices to get  $\tilde{\mathbf{C}} \in \mathbb{R}^{ld_c \times t}$ .

## Attention in the **encoder**:

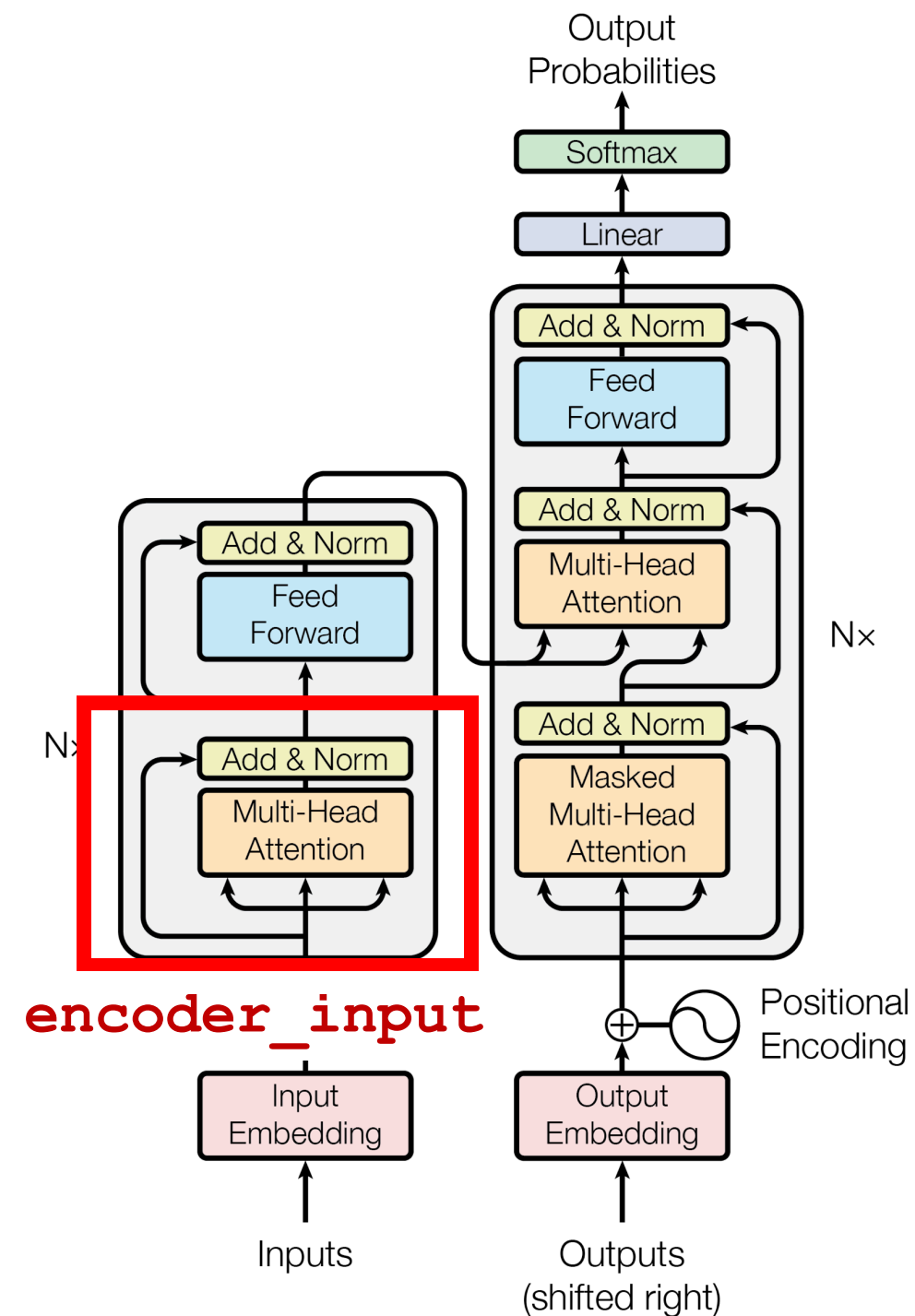
- $Q = K = V = \text{encoder\_input}$ .

## 1<sup>st</sup> attention in the **decoder**:

- $Q = K = V = \text{decoder\_input}$ .

## 2<sup>nd</sup> attention in the **decoder**

- $Q = \text{decoder\_input}$
- $K = V = \text{encoder\_output}$ .



### Attention in the **encoder**:

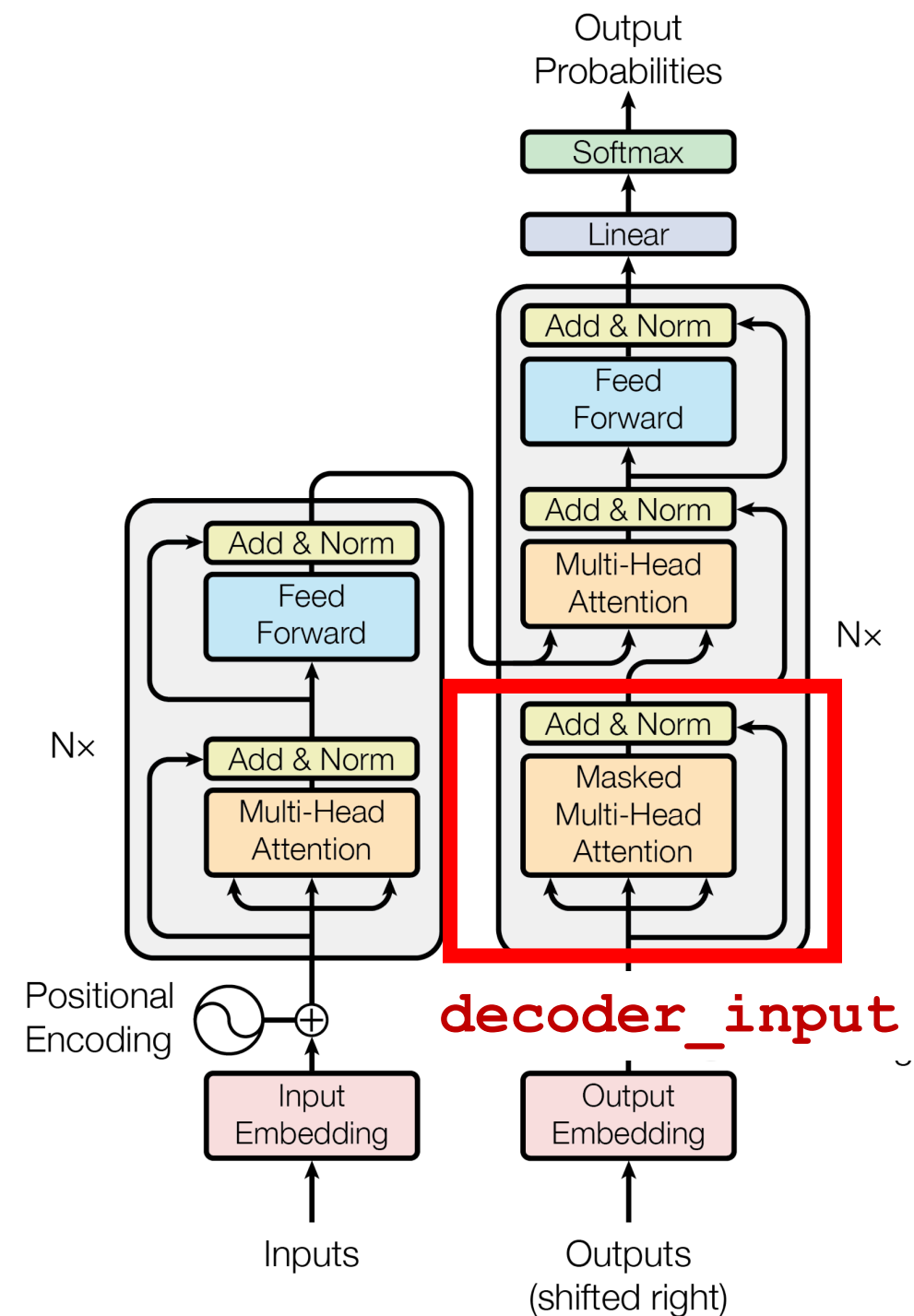
- $Q = K = V = \text{encoder\_input}.$

### **1<sup>st</sup>** attention in the **decoder**:

- $Q = K = V = \text{decoder\_input}.$

### **2<sup>nd</sup>** attention in the **decoder**

- $Q = \text{decoder\_input}$
- $K = V = \text{encoder\_output}.$





Attention in the **encoder**:

- $Q = K = V = \text{encoder\_input}.$

1<sup>st</sup> attention in the **decoder**:

- $Q = K = V = \text{decoder\_input}.$

2<sup>nd</sup> attention in the **decoder**

- $Q = \text{decoder\_input}$

- $K = V = \text{encoder\_output}.$

