# Transformers I

CS 4391 Introduction Computer Vision

Professor Yu Xiang

The University of Texas at Dallas

# Machine Translation

- Translate a phrase from one language to anther
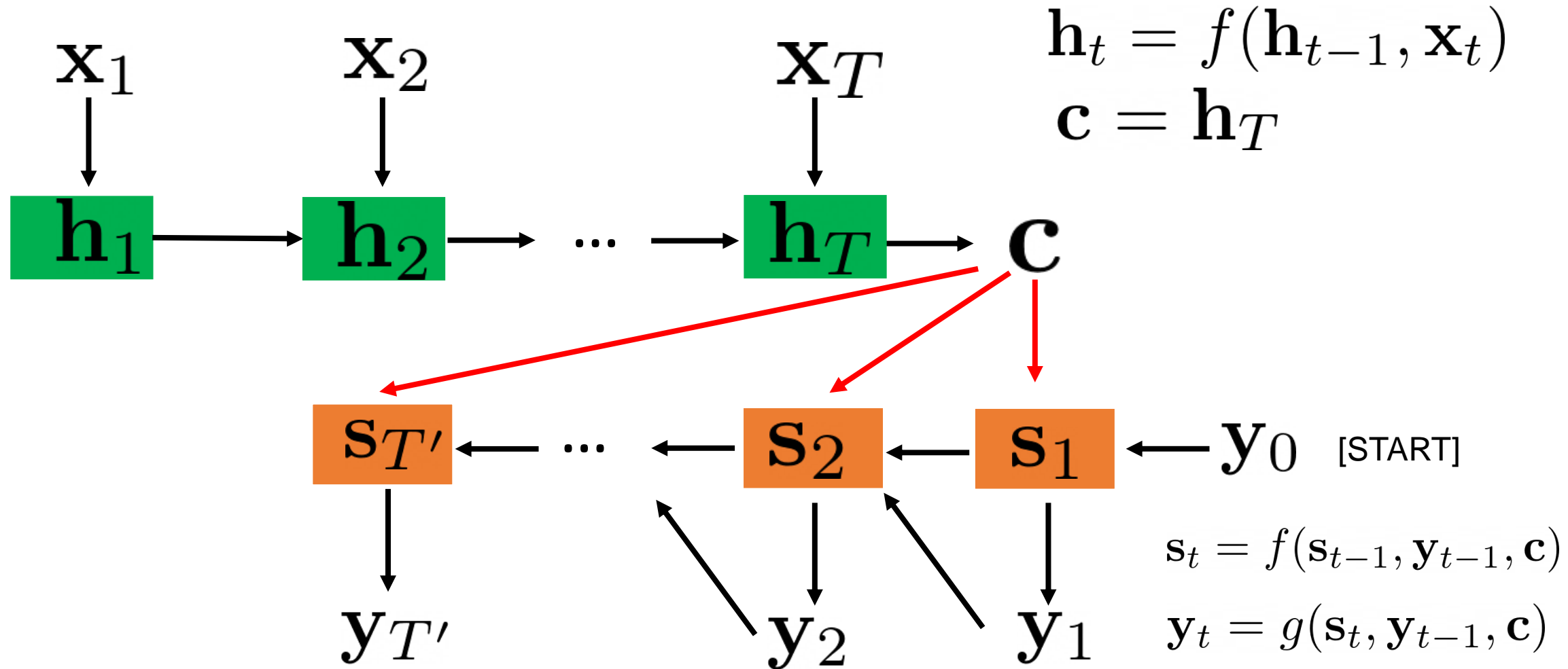  - E.g., English phrase to French phrase



Google Translation

13 words

15 words

# RNN Encoder-Decoder



$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t)$$
$$\mathbf{c} = \mathbf{h}_T$$

[START]

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}, \mathbf{y}_{t-1}, \mathbf{c})$$
$$\mathbf{y}_t = g(\mathbf{s}_t, \mathbf{y}_{t-1}, \mathbf{c})$$
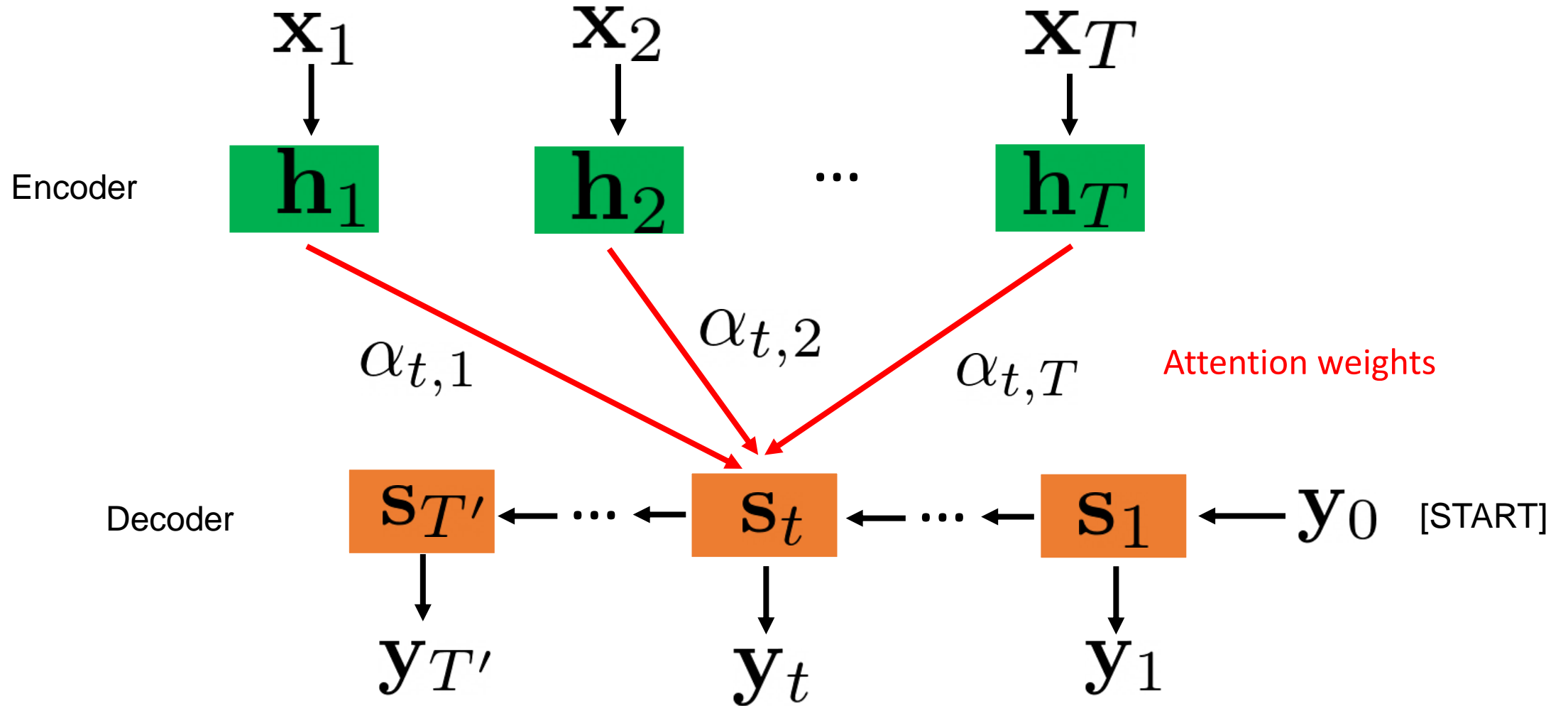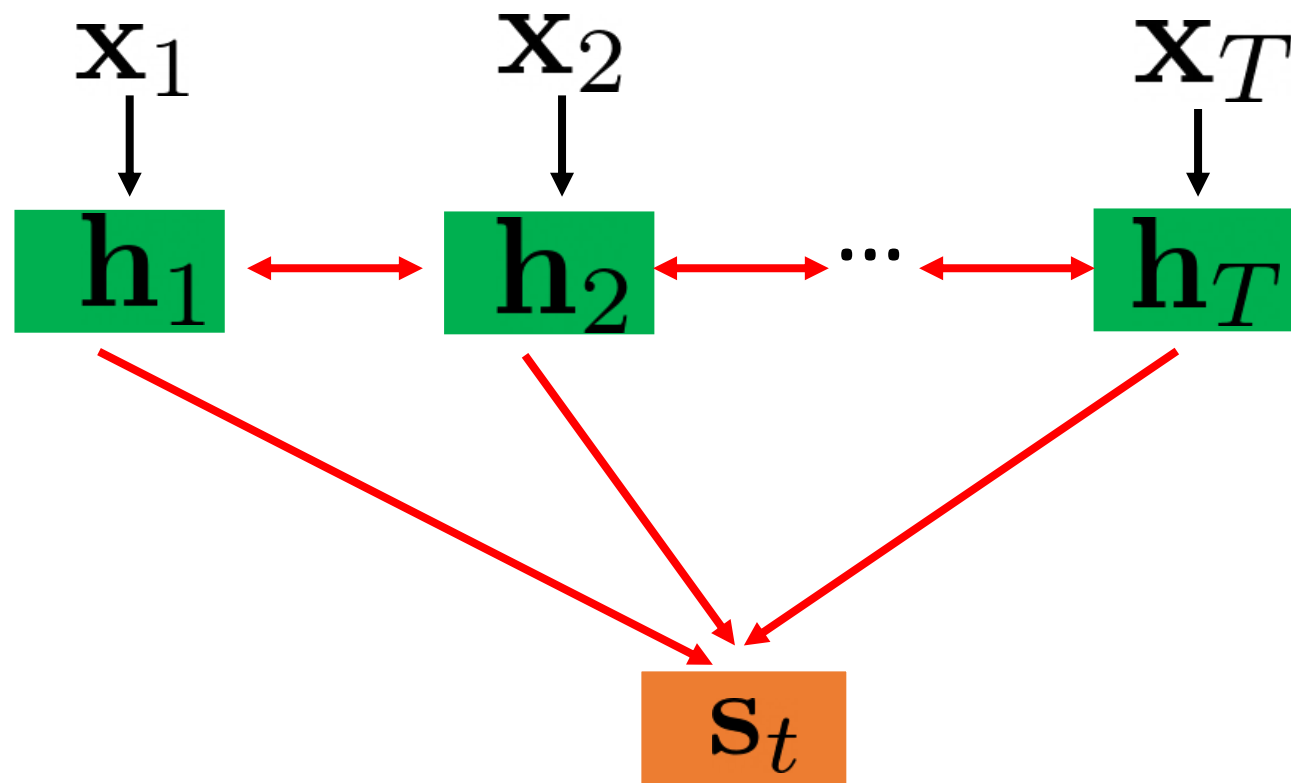
Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. Cho et al., EMNLP'14

# Transformer: Encoder-Decoder with Attention



Encoder

$\mathbf{x}_1$  $\mathbf{x}_2$  $\mathbf{x}_T$

$\mathbf{h}_1$  $\mathbf{h}_2$  ...  $\mathbf{h}_T$

$\alpha_{t,1}$  $\alpha_{t,2}$  $\alpha_{t,T}$  Attention weights

Decoder

$\mathbf{s}_{T'} \leftarrow \cdots \leftarrow \mathbf{s}_t \leftarrow \cdots \leftarrow \mathbf{s}_1 \leftarrow \mathbf{y}_0$  [START]

$\mathbf{y}_{T'}$  $\mathbf{y}_t$  $\mathbf{y}_1$

# Transformer: Attention

- Two types of attentions



- Self-attention



- Cross-attention

$\mathbf{x}_1 \quad\quad \mathbf{x}_2 \quad\quad\quad\quad \mathbf{x}_T$

$\mathbf{h}_1 \longleftrightarrow \mathbf{h}_2 \longleftrightarrow \cdots \longleftrightarrow \mathbf{h}_T$

$\mathbf{s}_t$

# Transformer: Attention

- Input
  - (key, value) pairs (think about python dictionary)
  - A query

- Output
  - Compare the query to all the keys to compute weights
  - Weighted sum of the values

Attention is all you need. Vaswani et al., NeurIPS'17

# Transformer: Attention

- Scaled Dot-Product Attention
  - Keys $\quad K : m \times d_k$
  - Values $\quad V : m \times d_v$
  - n queries $Q : n \times d_k$

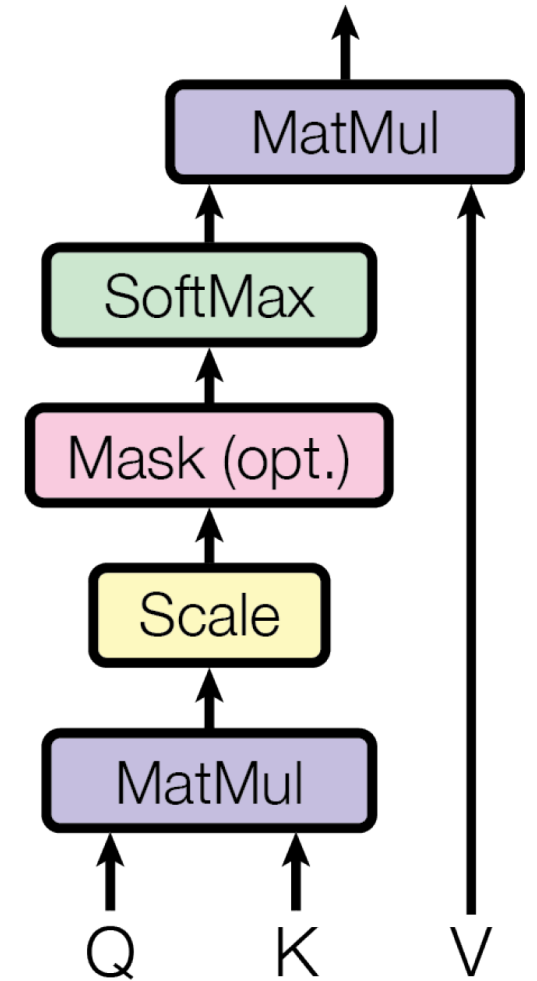$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Softmax function

$$\sigma(\mathbf{y})_i = \frac{e^{y_i}}{\sum_j^m e^{y_i}}$$

$$n \times d_v$$

weights



Attention is all you need. Vaswani et al., NeurIPS'17

# Transformer: Attention
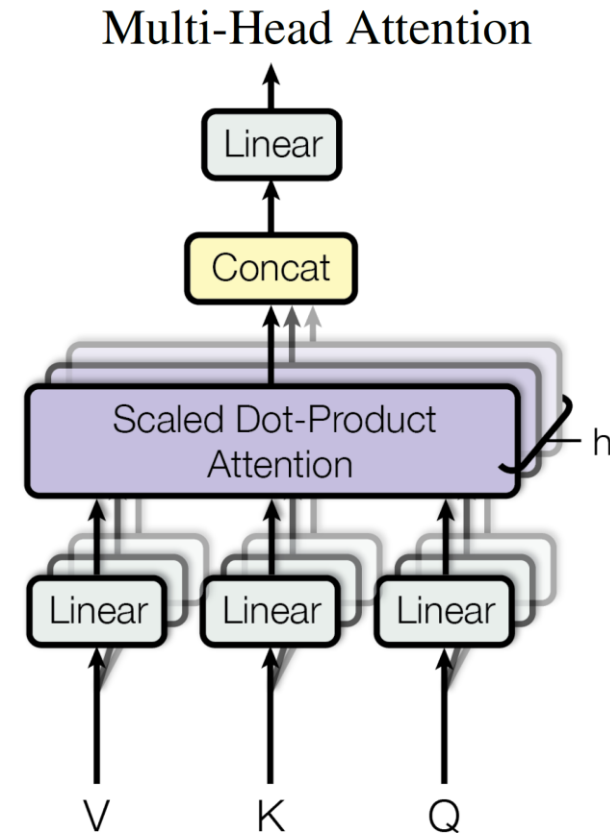
- Multi-Head Attention
  - Suppose the latent vector is with dimension $d_{\text{model}}$


Multi-Head Attention

$$m \times d_{\text{model}} \qquad d_{\text{model}} \times d_k$$

$$\text{head}_{\text{i}} = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad \textcolor{red}{\text{Projection}}$$

$$n \times d_v$$

$$n \times d_{\text{model}} \quad d_{\text{model}} \times d_k \qquad m \times d_{\text{model}} \quad d_{\text{model}} \times d_v$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_{\text{h}})W^O$$

$$n \times d_{\text{model}} \qquad\qquad n \times hd_v \qquad\qquad hd_v \times d_{\text{model}}$$

Attention is all you need. Vaswani et al., NeurIPS'17

# Transformer: Encoder

Encoder output

- Self-attention (repeat N times)
  - Keys, values and queries are all the same
  - n input tokens $\quad n \times d_{\text{model}}$

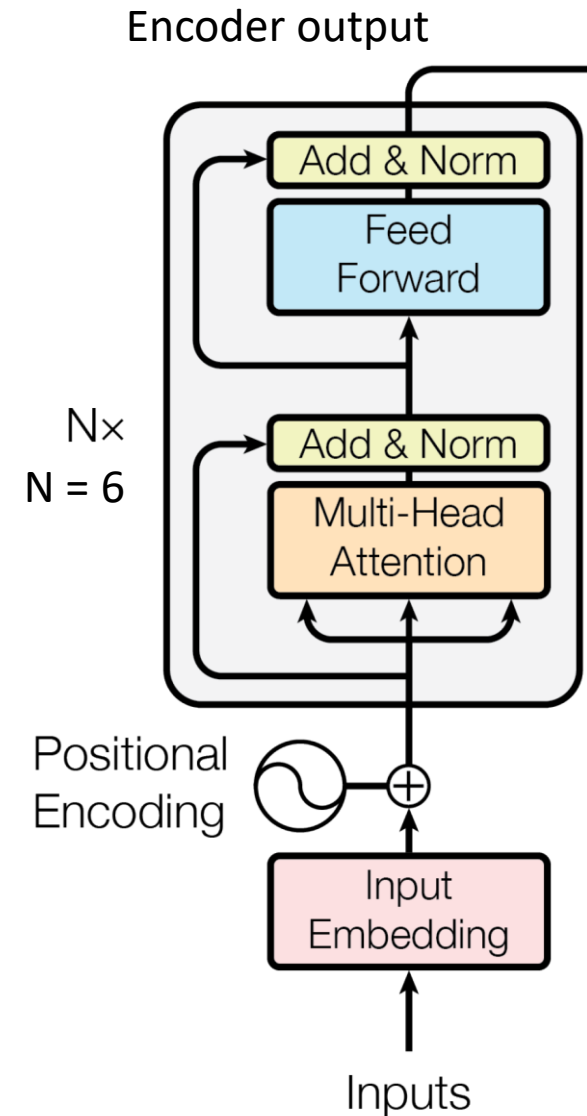$$\text{MultiHead}(Q, K, V)$$

  - Residual connection

$$\text{LayerNorm}(x + \text{Sublayer}(x))$$

  - Layer normalization

$$\mu^l = \frac{1}{H} \sum_{i=1}^{H} a_i^l \qquad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^{H} \left(a_i^l - \mu^l\right)^2} \qquad \frac{a^l - \mu^l}{\sigma^l}$$

Attention is all you need. Vaswani et al., NeurIPS'17

Nx
N = 6

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Transformer: Encoder
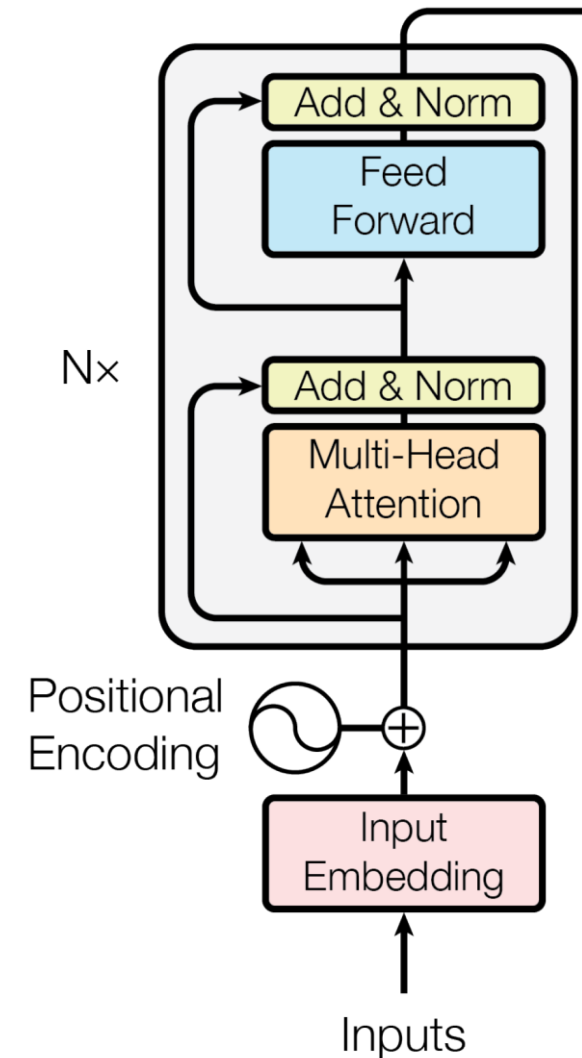
- Feed Forward Network

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

- Positional encoding
  - Make use the order of the sequence
  - With dimension $d_{\text{model}}$ for each input

$$PE_{(pos, 2i)} = sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos, 2i+1)} = cos(pos/10000^{2i/d_{\text{model}}})$$
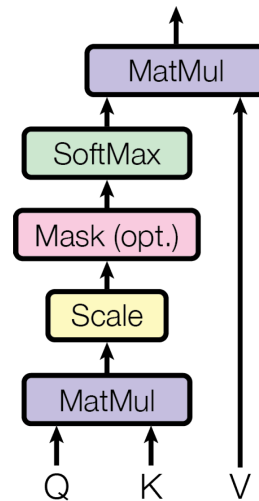
Attention is all you need. Vaswani et al., NeurIPS'17
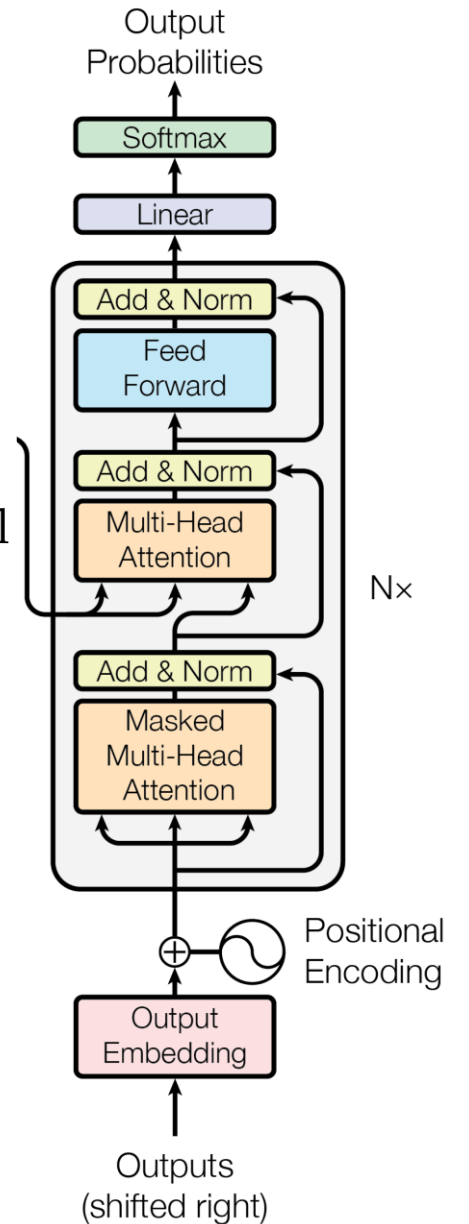
# Transformer: Decoder

- Output embedding

[START]

$$\mathbf{y}_0 \; \mathbf{y}_1 \cdots \mathbf{y}_{t-1} \; \mathbf{y}_t \; \mathbf{y}_{t+1} \cdots \mathbf{y}_{T'}$$

Shifted right by one position and insert the start token

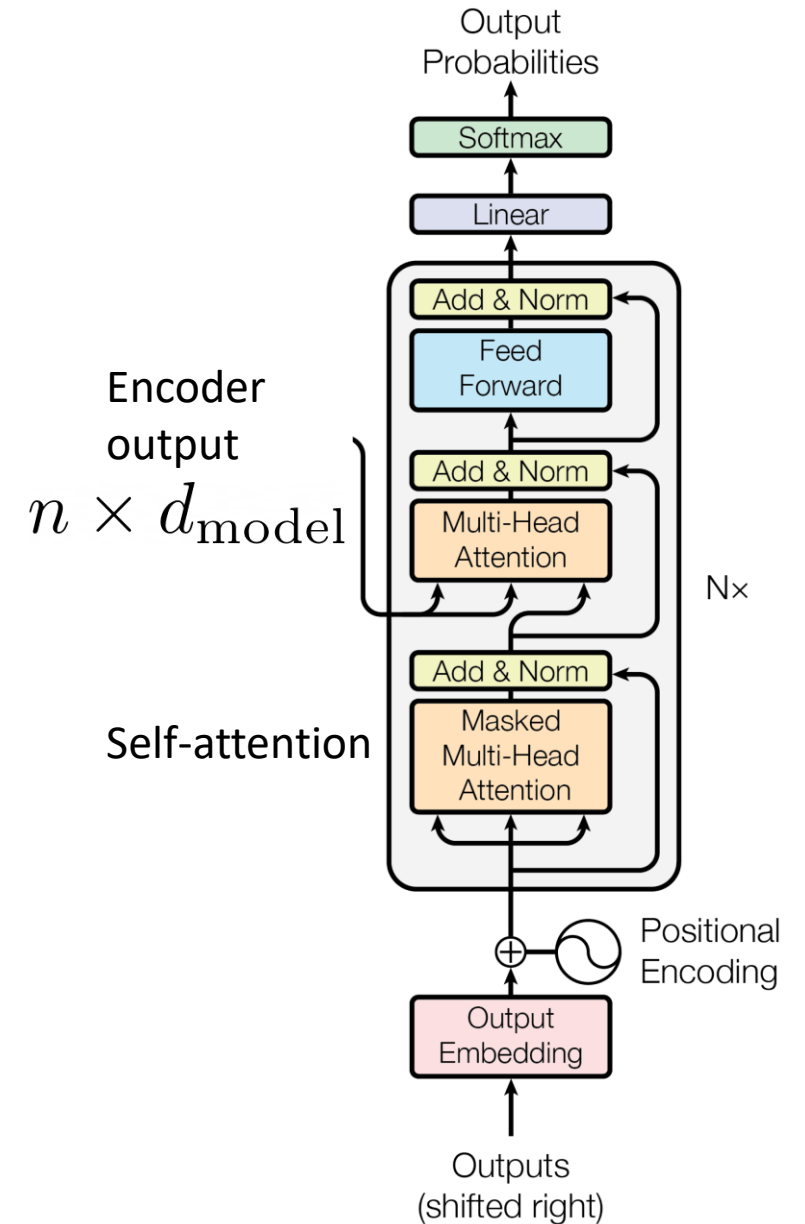Mask out current and future outputs during training (setting to $-\infty$ )

Encoder output $n \times d_{\mathrm{model}}$



Attention is all you need. Vaswani et al., NeurIPS'17
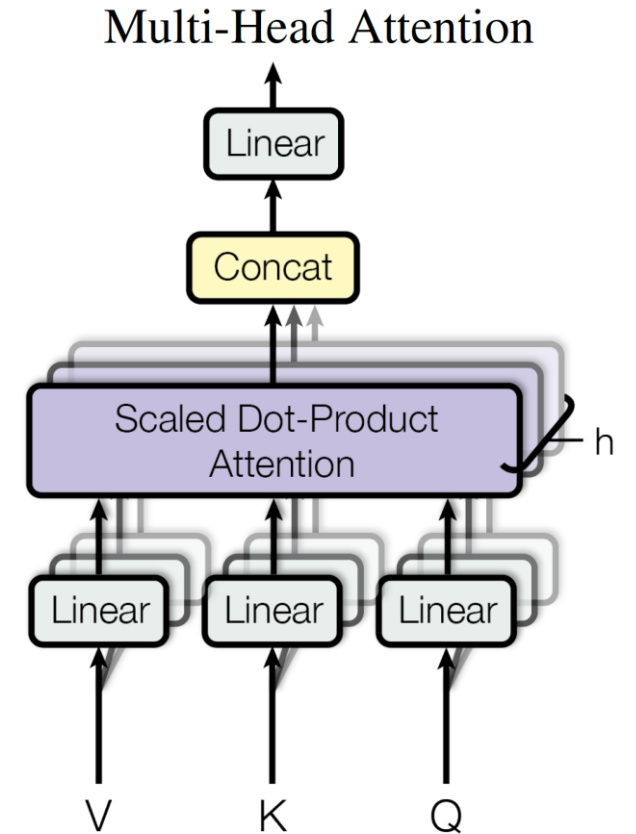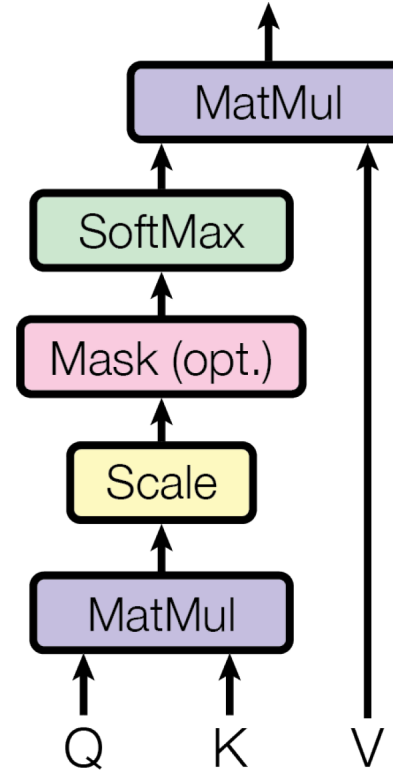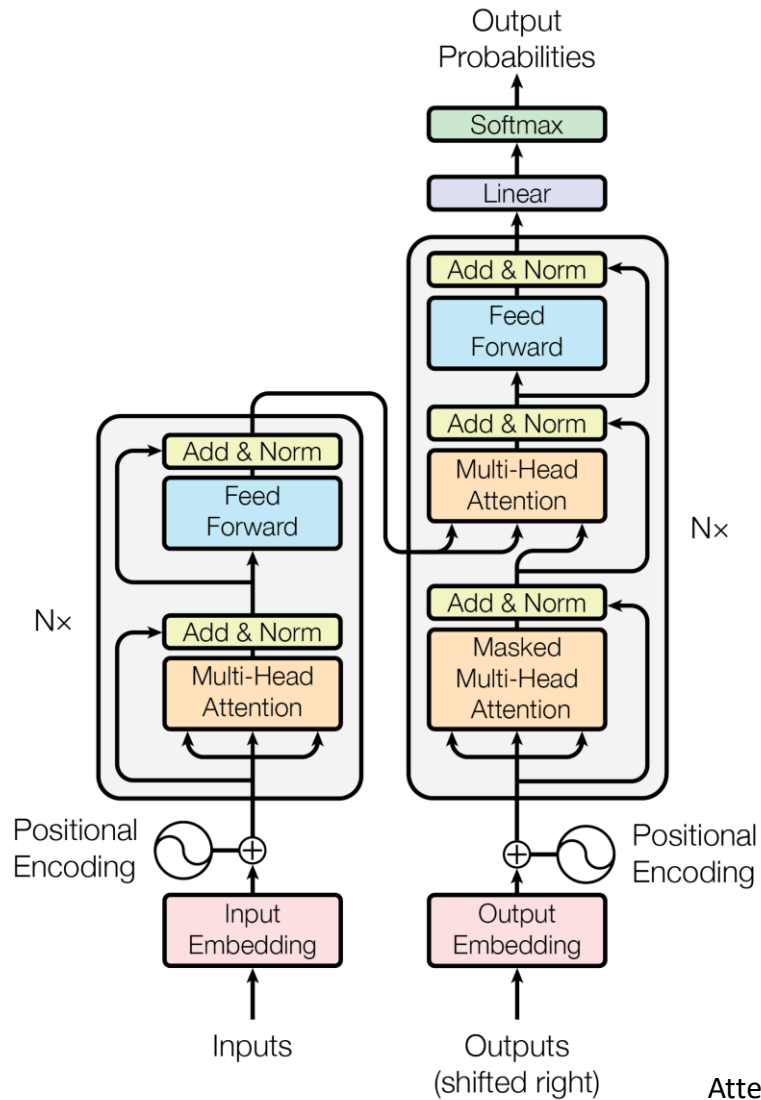
# Transformer: Decoder

- Encoder-decoder attention
  - (Key, value): encoder output
  - Queries: decoder output
  - Every position in the decoder attends to all positions in the input sequence

- Softmax
  - Predicts next-token probabilities



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Encoder output

$n \times d_{\mathrm{model}}$

Add & Norm

Multi-Head Attention

N×

Add & Norm

Self-attention

Masked Multi-Head Attention

Positional Encoding
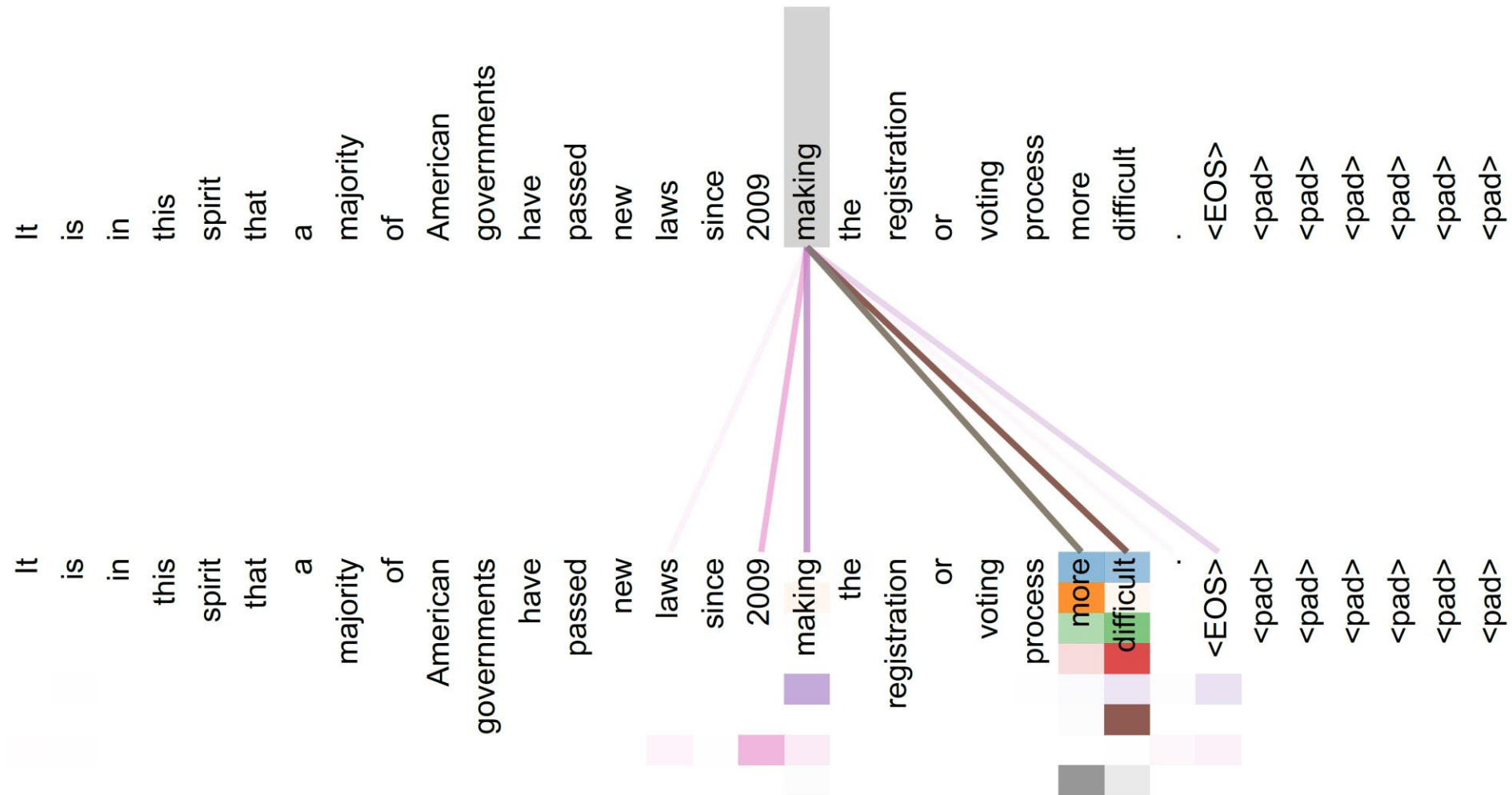
Output Embedding

Outputs (shifted right)

Attention is all you need. Vaswani et al., NeurIPS'17

# Transformer



Attention is all you need. Vaswani et al., NeurIPS'17

# Transformer: Attention Visualization



Attention is all you need. Vaswani et al., NeurIPS'17

# Summary

- Transformers
  - Can capture long-distance dependencies (global attention)

  - Computationally efficient, more parallelizable

# Further Reading

- Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation https://arxiv.org/abs/1406.1078

- Neural Machine Translation by Jointly Learning to Align and Translate https://arxiv.org/abs/1409.0473

- Transformer: Attention is all you need https://arxiv.org/abs/1706.03762