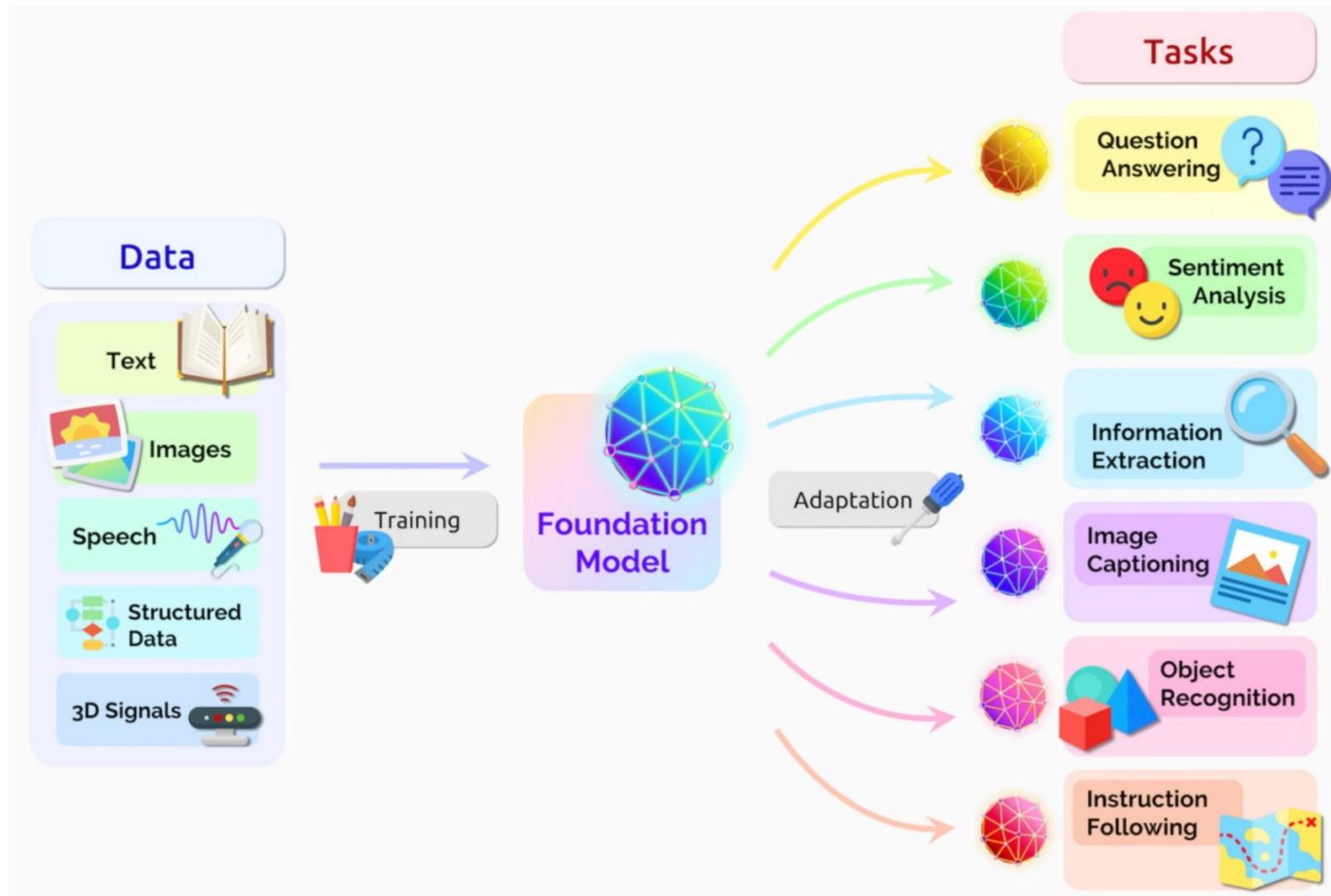


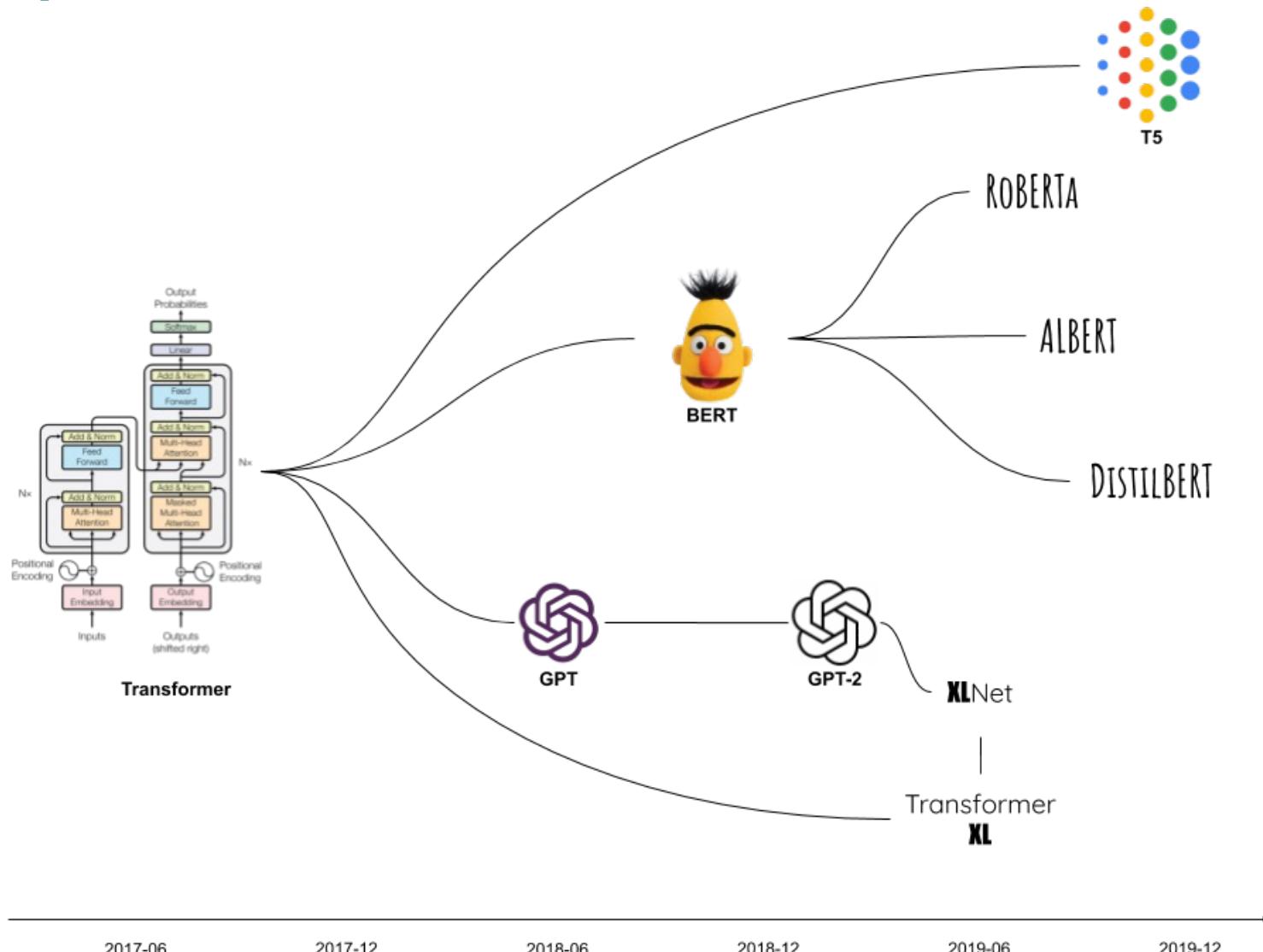
Computer Vision: Transformer

Siheng Chen 陈思衡

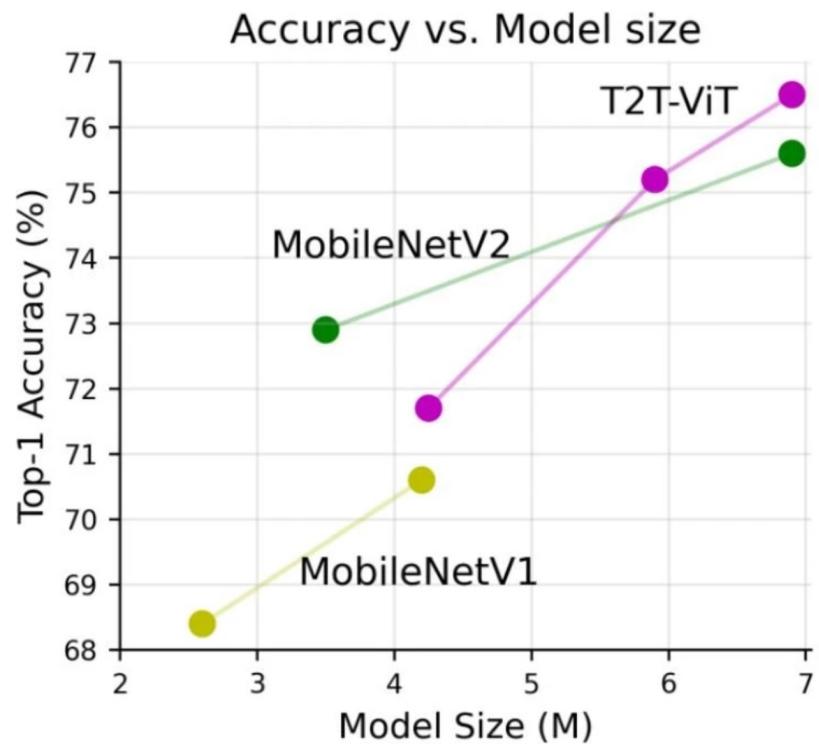
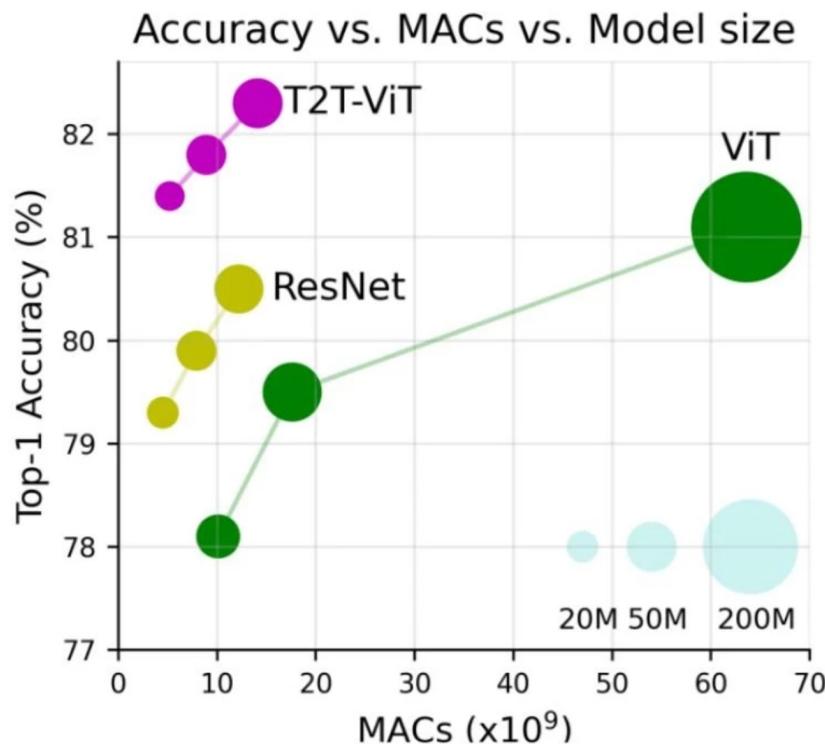
Why Transformer?



Why Transformer?

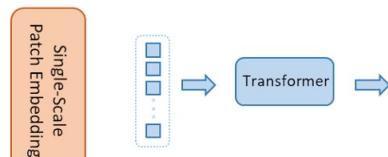
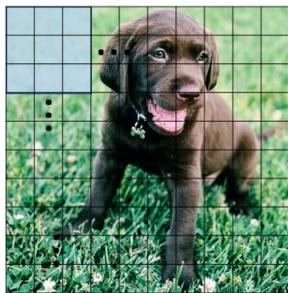


Why Transformer?



Why Transformer?

ViT-variants : Single-scale patch + Single-Path structure



Ours : Multi-scale Patches + Multi-Path structure

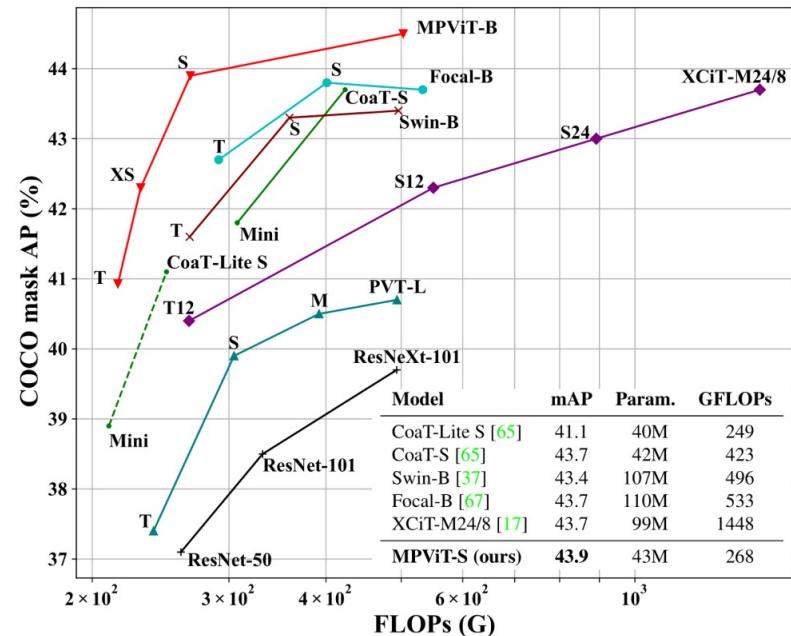
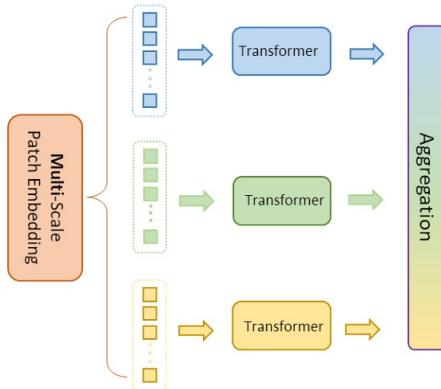
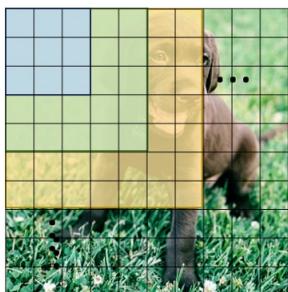
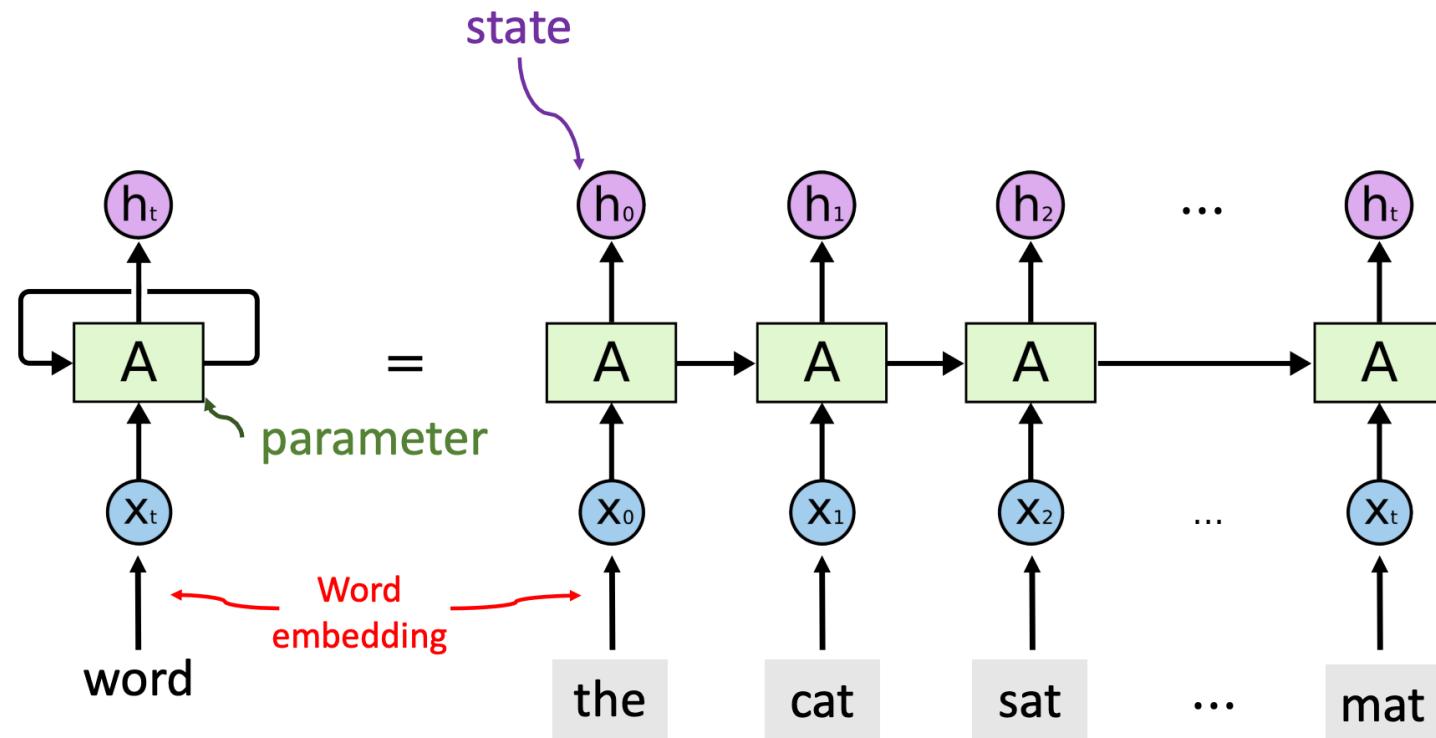


Figure 2. **FLOPs vs. COCO mask AP** on Mask R-CNN. MPViTs outperform state-of-the-art Vision Transformers while having fewer parameters and FLOPs. B, S, XS, and T at the end of the model names denote base, small, extra-small and tiny respectively. Complete results are in Table 3.

Recurrent NN

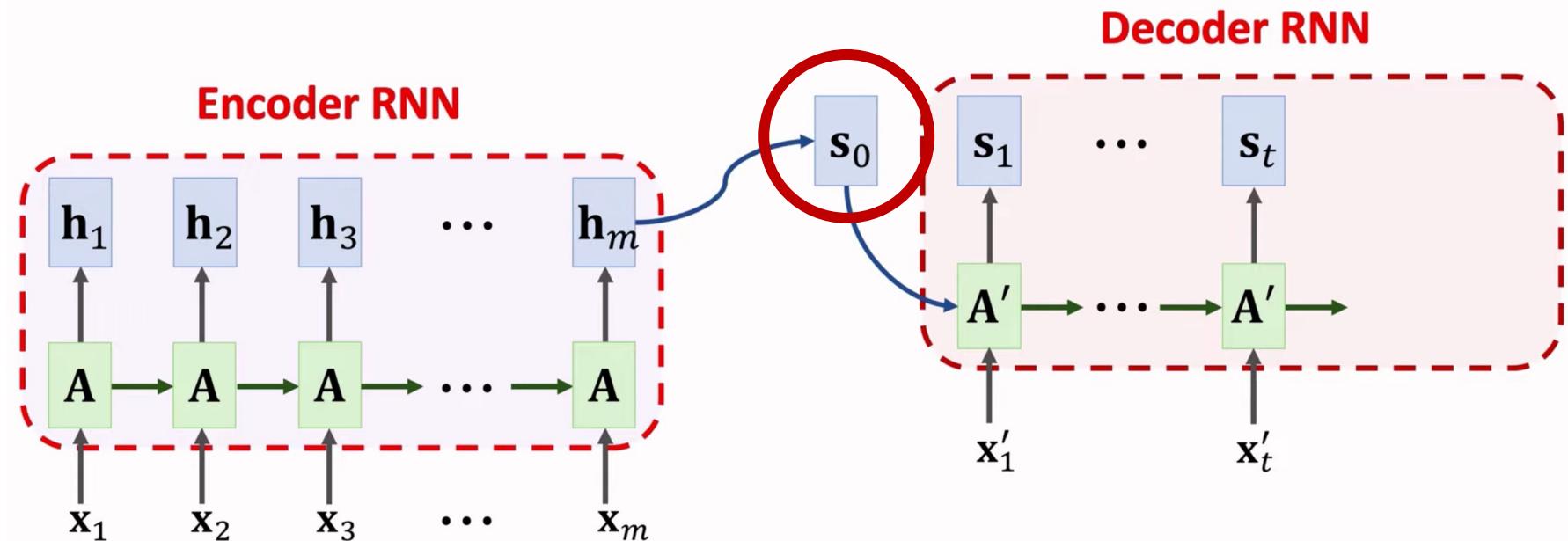


Roadmap

1. Recurrent neural network (RNN)
2. Seq2seq model
3. Attention for Seq2seq model
4. Self-attention for RNN
5. Attention without RNN: Attention layer / self-attention layer
6. Transformer
7. Vision Transformer

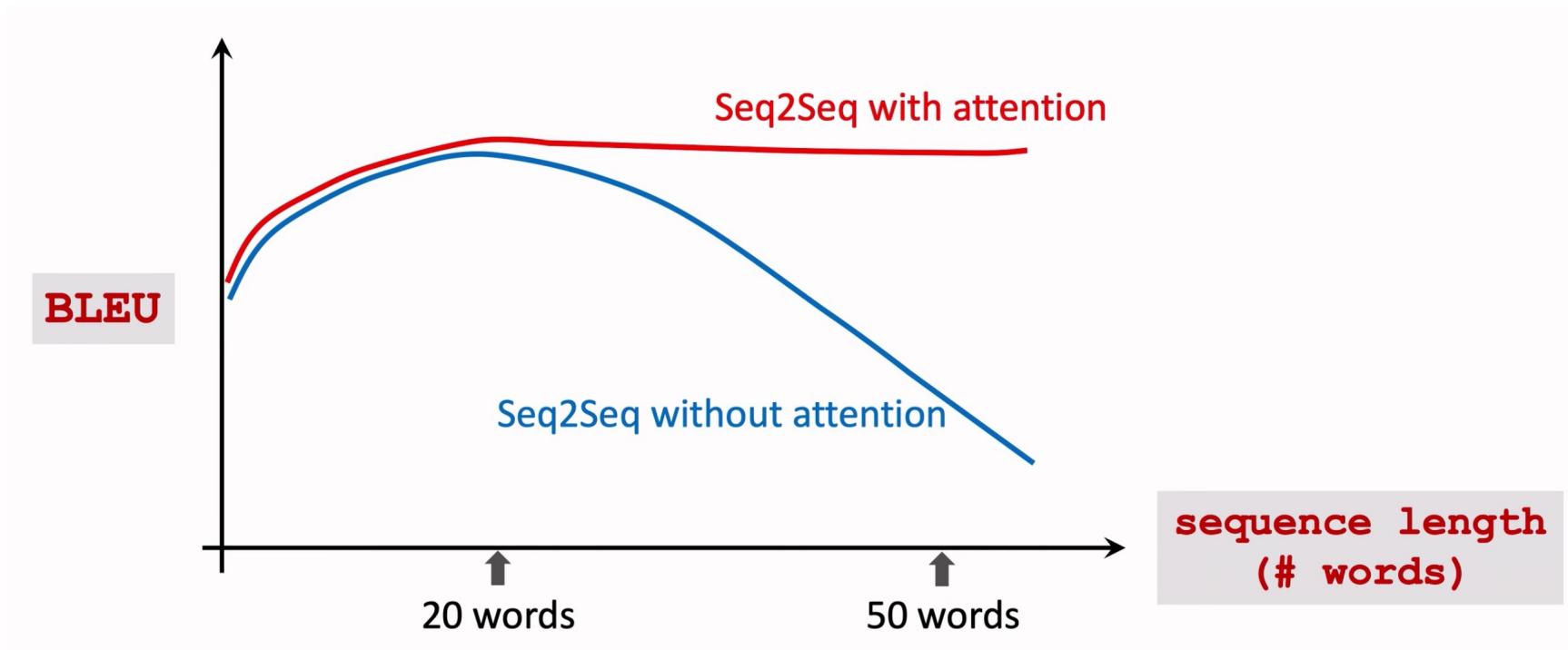
Seq2Seq model

Translation: from one sentence to another sentence



Problem: The final state is incapable of remembering a long sequence

Seq2Seq model



Problem: The final state is incapable of remembering a long sequence

Solution: Attention

Attention for Seq2Seq model

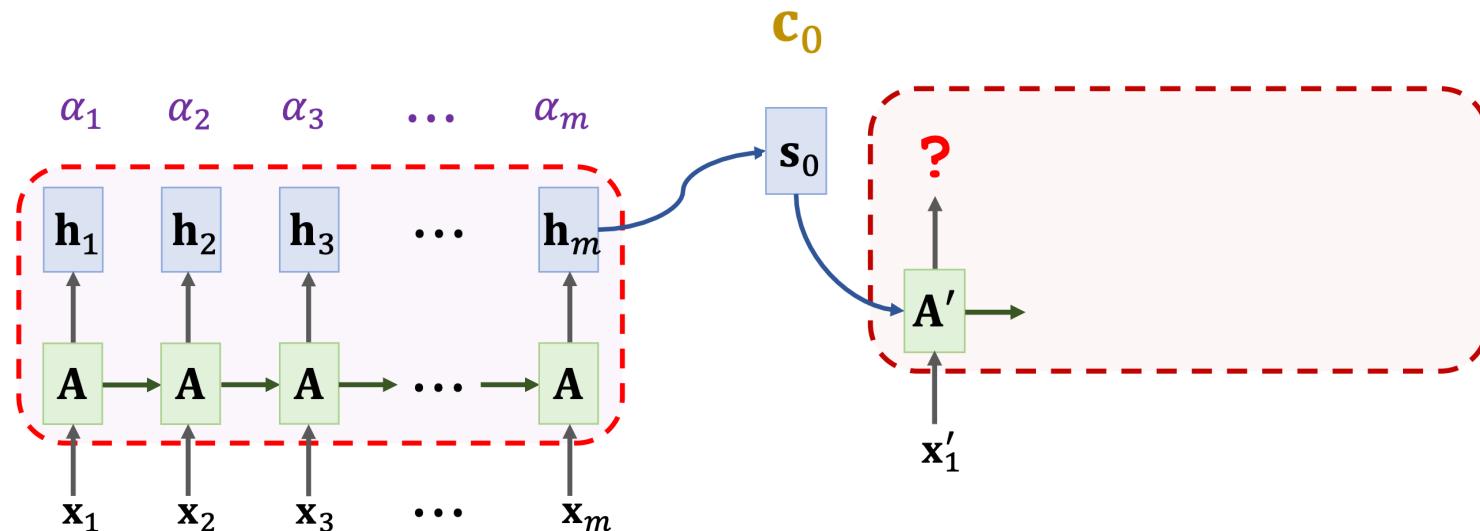
Weight: $\alpha_i = \text{align}(\mathbf{h}_i, \mathbf{s}_0)$.

SimpleRNN:

$$\mathbf{s}_1 = \tanh \left(\mathbf{A}' \cdot \begin{bmatrix} \mathbf{x}'_1 \\ \mathbf{s}_0 \end{bmatrix} + \mathbf{b} \right)$$

SimpleRNN + Attention:

$$\mathbf{s}_1 = \tanh \left(\mathbf{A}' \cdot \begin{bmatrix} \mathbf{x}'_1 \\ \mathbf{s}_0 \\ \mathbf{c}_0 \end{bmatrix} + \mathbf{b} \right)$$

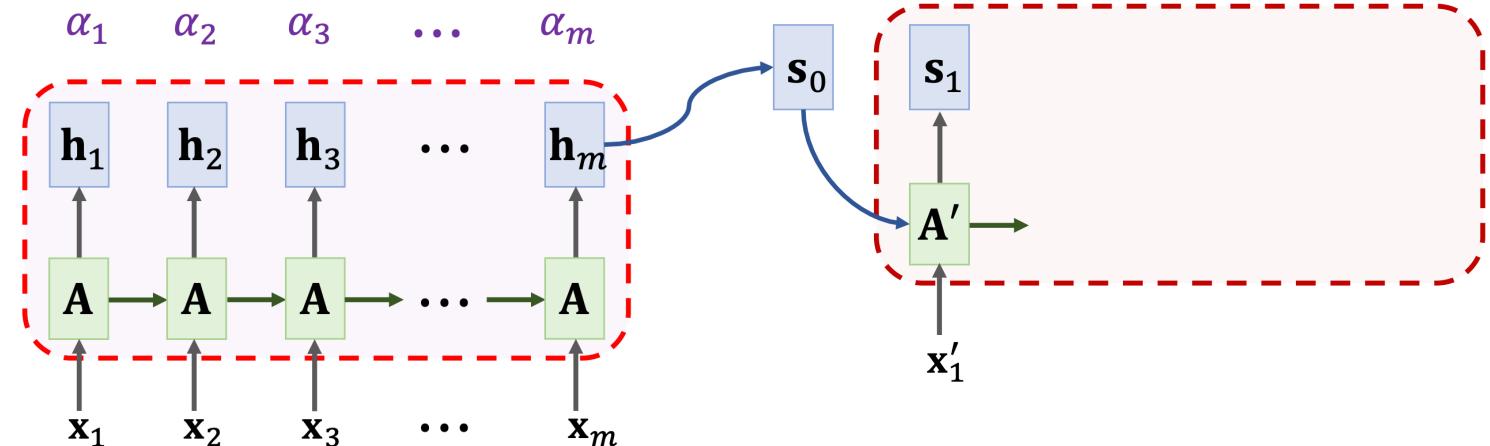


Attention for Seq2Seq model

Weight: $\alpha_i = \text{align}(\mathbf{h}_i, \mathbf{s}_1)$.

$$\mathbf{s}_2 = \tanh \left(\mathbf{A}' \cdot \begin{bmatrix} \mathbf{x}'_2 \\ \mathbf{s}_1 \\ \mathbf{c}_1 \end{bmatrix} + \mathbf{b} \right)$$

re-compute each time!



Attention for Seq2Seq model

Weight: $\alpha_i = \text{align}(\mathbf{h}_i, \mathbf{s}_0)$.

Attention weight computation

1. Linear maps:

- $\mathbf{k}_i = \mathbf{W}_K \cdot \mathbf{h}_i$, for $i = 1$ to m .
- $\mathbf{q}_0 = \mathbf{W}_Q \cdot \mathbf{s}_0$.

2. Inner product:

- $\tilde{\alpha}_i = \mathbf{k}_i^T \mathbf{q}_0$, for $i = 1$ to m .

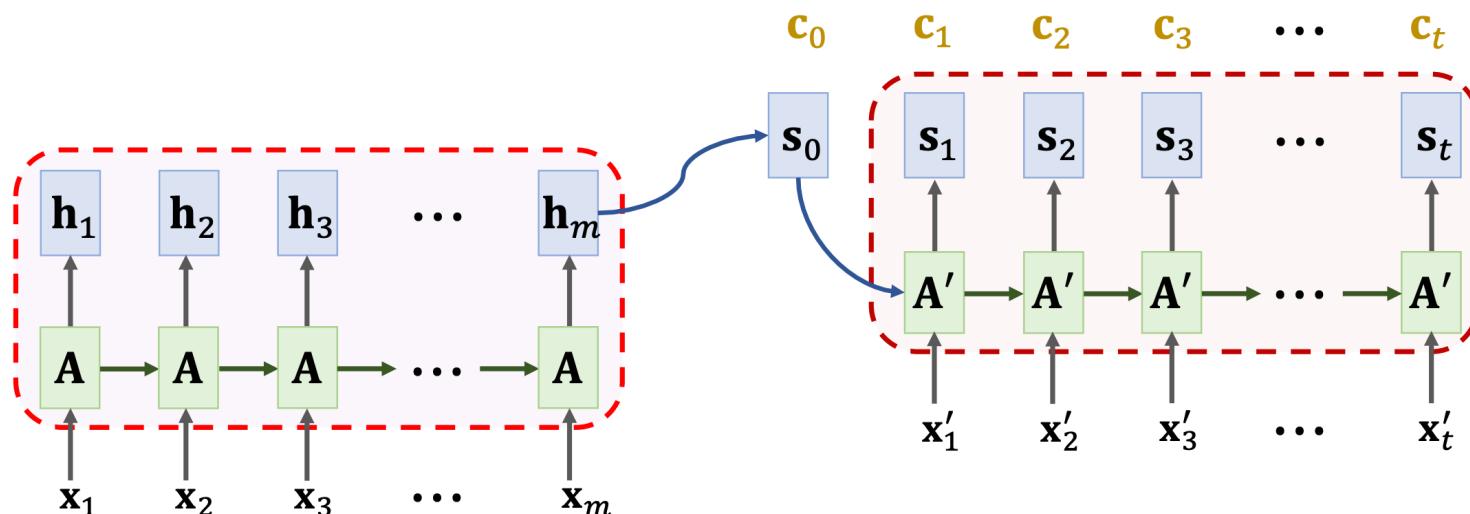
3. Normalization:

- $[\alpha_1, \dots, \alpha_m] = \text{Softmax}([\tilde{\alpha}_1, \dots, \tilde{\alpha}_m])$.

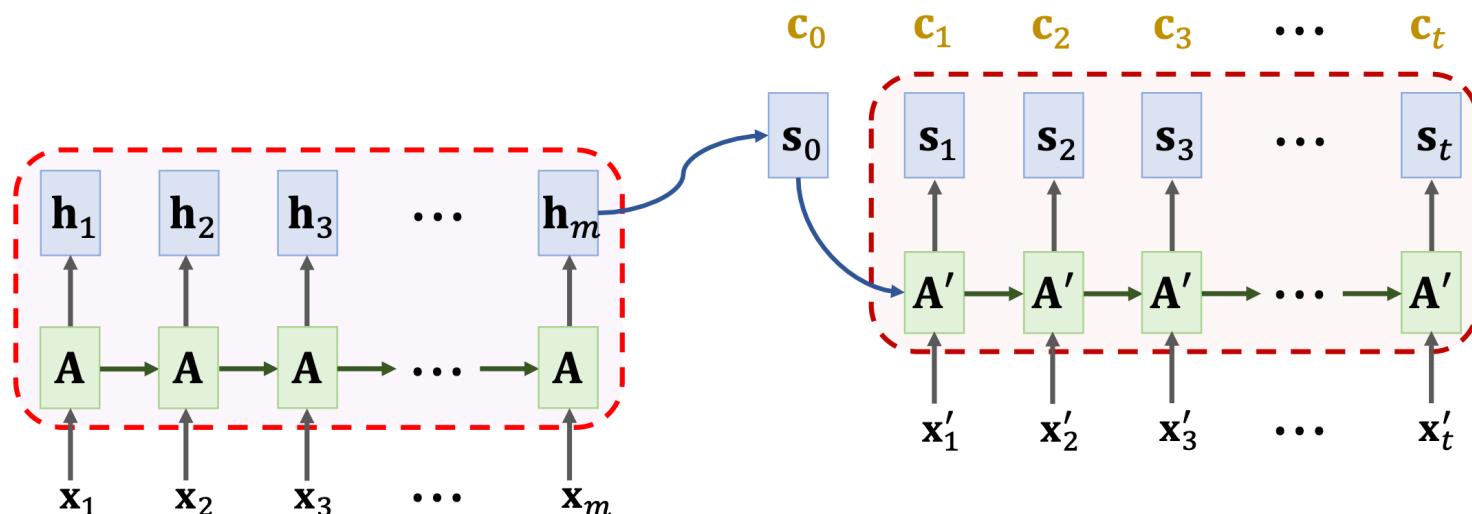
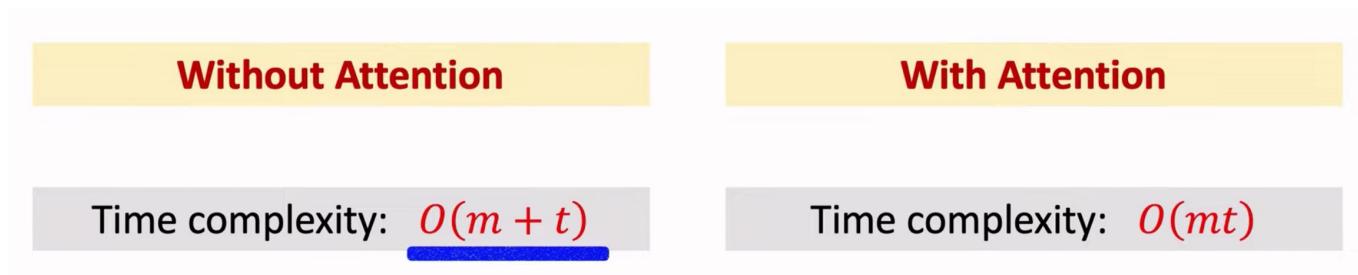
Attention for Seq2Seq model

Question: How many attention weights have been computed?

- To compute each c_i , we compute m weights
- The decode has t states, so there are totally mt weights



Attention for Seq2Seq model



Attention for Seq2Seq model

Advantages

- Standard Seq2Seq model: the decoder looks at only its current state
- Attention: decoder additionally looks at all the states of the encoder
- Attention: decoder knows where to focus

Disadvantages

- Much higher time complexity.

Self-Attention for RNN

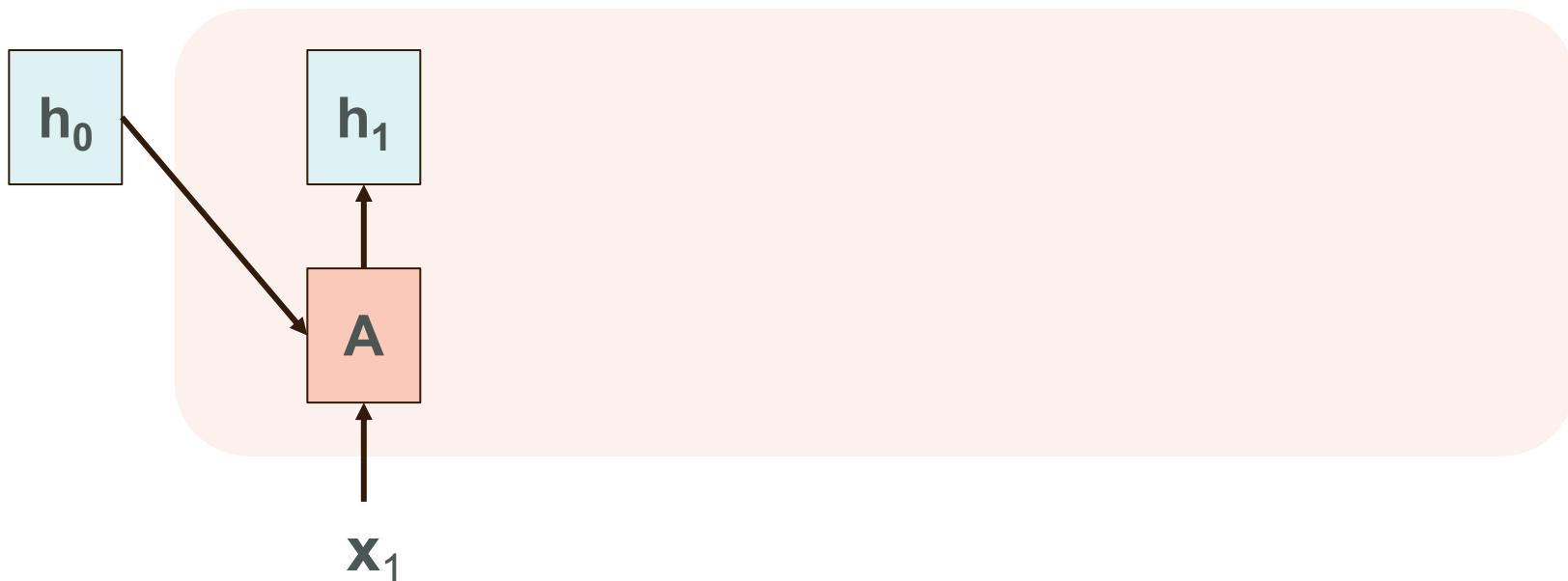
SimpleRNN:

$$\mathbf{h}_1 = \tanh(\mathbf{A} \cdot [\mathbf{x}_1 \mathbf{h}_0] + \mathbf{b})$$

SimpleRNN + Self-Attention:

$$\mathbf{h}_1 = \tanh(\mathbf{A} \cdot [\mathbf{x}_1 \mathbf{c}_1] + \mathbf{b})$$

$$\mathbf{c}_0 = 0 \quad \mathbf{c}_1 = \mathbf{h}_1$$



Self-Attention for RNN

SimpleRNN + Self-Attention:

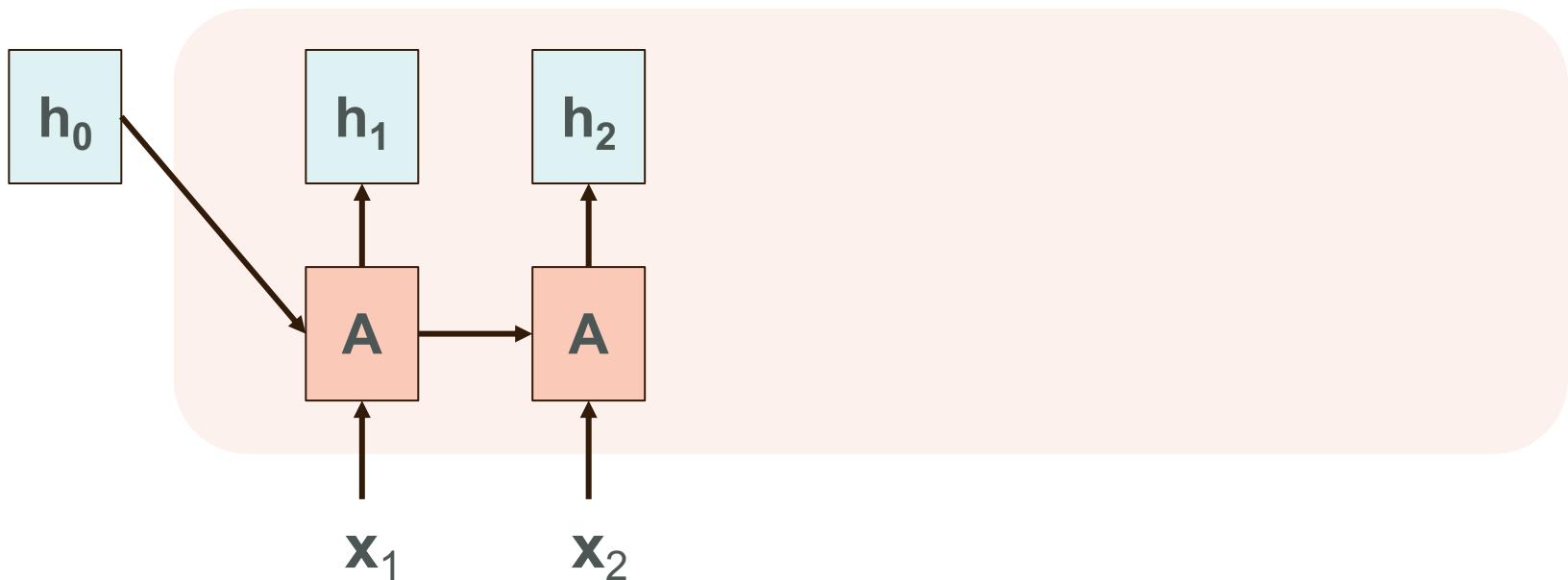
$$\mathbf{h}_1 = \tanh(\mathbf{A} \cdot [\mathbf{x}_1 | \mathbf{c}_0] + \mathbf{b})$$

Weights: $\alpha_i = \text{align}(\mathbf{h}_i, \mathbf{h}_2)$.

$$\mathbf{c}_0 = 0$$

$$\mathbf{c}_1 = \mathbf{h}_1$$

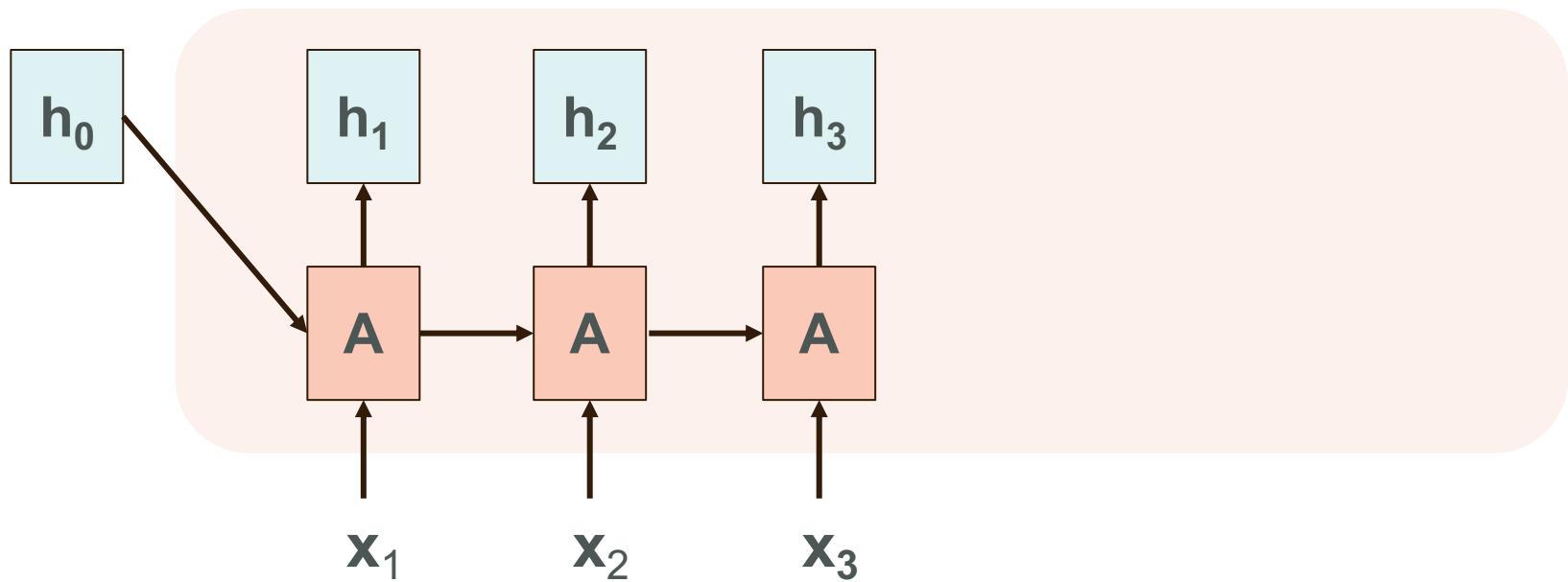
$$\mathbf{c}_2 = \alpha_1 \mathbf{h}_1 + \alpha_2 \mathbf{h}_2$$



Self-Attention for RNN

$$h_i \longrightarrow \alpha_1, \alpha_2, \dots, \alpha_i \longrightarrow c_i$$

$$c_0 = 0 \quad c_1 = h_1 \quad c_2 \quad c_3 = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3$$



Self-Attention for RNN

- With self-attention, RNN is less likely to forget
- Pay attention to the context relevant to the new input

The

The FBI

The FBI is

The FBI is chasing

The FBI is chasing a

The FBI is chasing a criminal

The FBI is chasing a criminal on

The FBI is chasing a criminal on the

The FBI is chasing a criminal on the run

The FBI is chasing a criminal on the run .

Attention without RNN

Encoder's inputs:

$x_1 \quad x_2 \quad \dots \quad x_m$

Decoder's inputs:

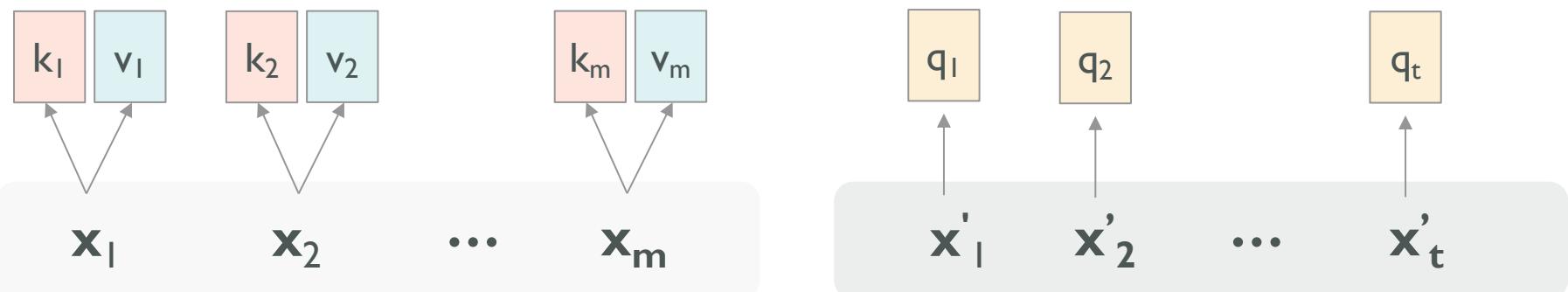
$x'_1 \quad x'_2 \quad \dots \quad x'_t$

Attention without RNN

Key: $k_i = W_K x_i$

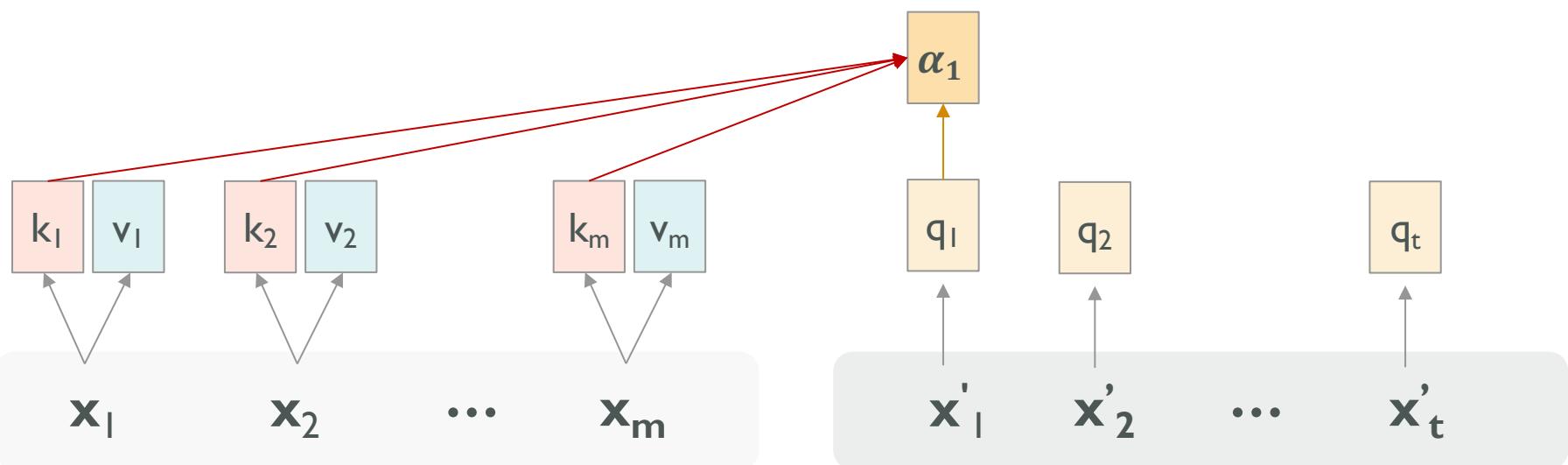
Value: $v_i = W_V x_i$

Query: $q_j = W_Q x'_j$



Attention without RNN

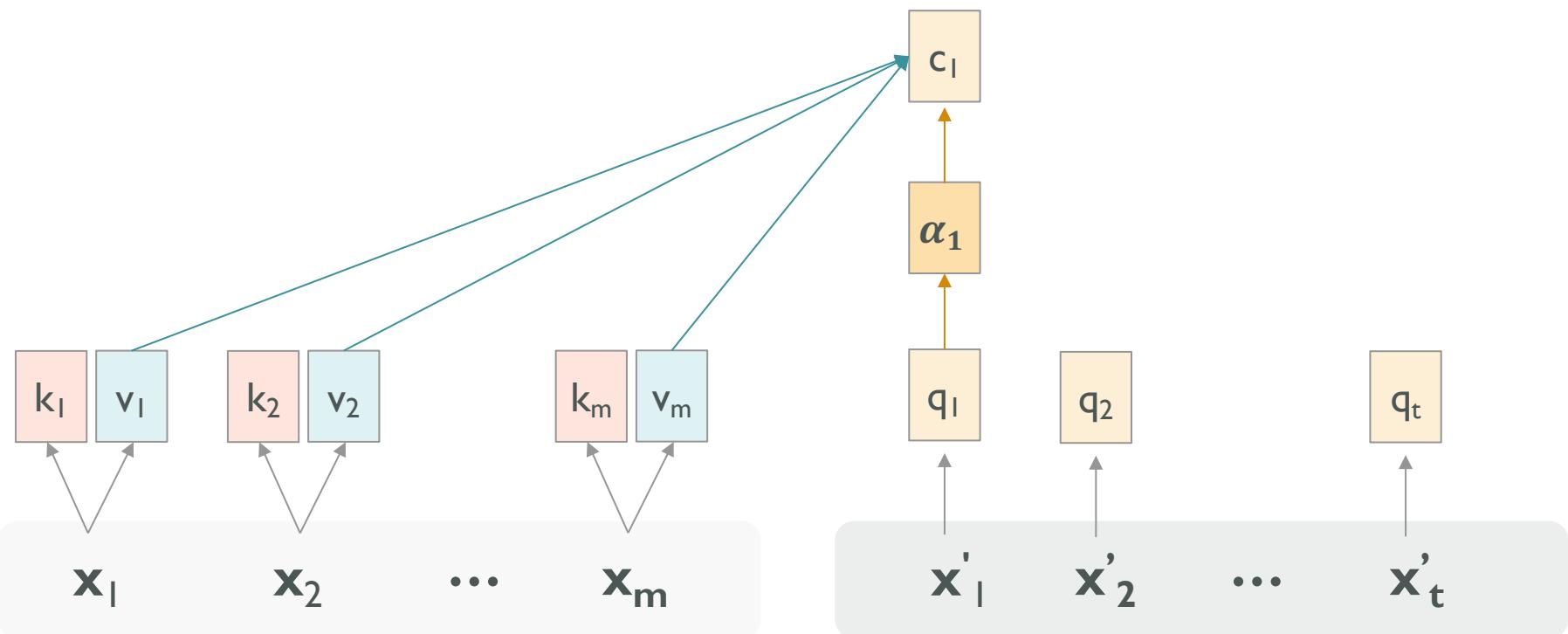
Compute attention weights: $\alpha_1 = \text{softmax} (K^T q_1)$



Attention without RNN

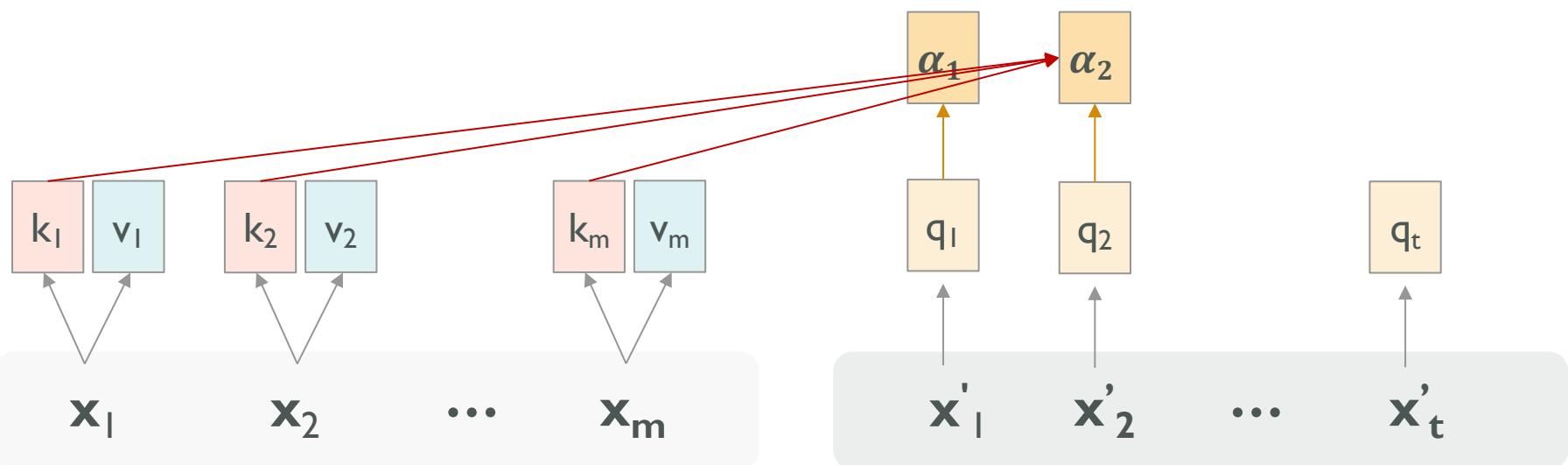
Compute attention weights: $\alpha_1 = \text{softmax} (K^T q_1)$

Compute context vector: $c_1 = V^T \alpha_1$



Attention without RNN

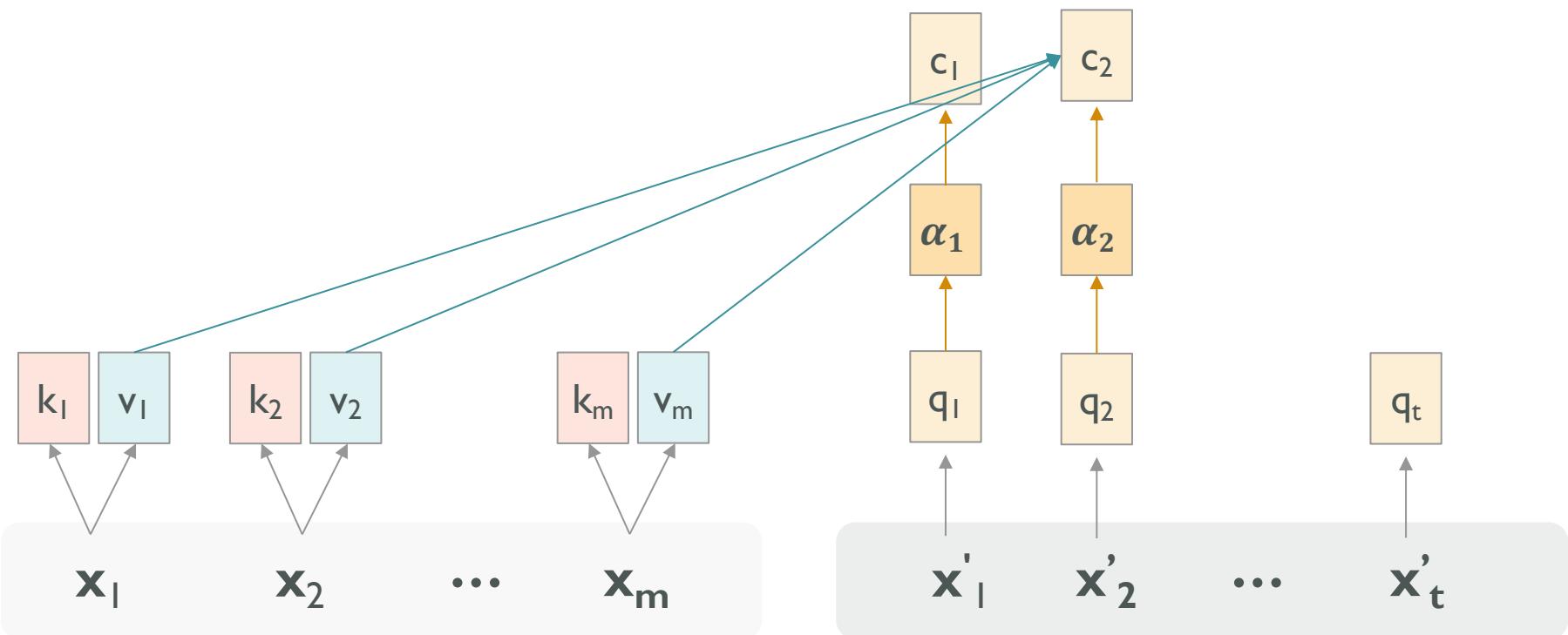
Compute attention weights: $\alpha_2 = \text{softmax} (K^T q_2)$



Attention without RNN

Compute attention weights: $\alpha_2 = \text{softmax} (K^T q_2)$

Compute context vector: $c_2 = V^T \alpha_2$

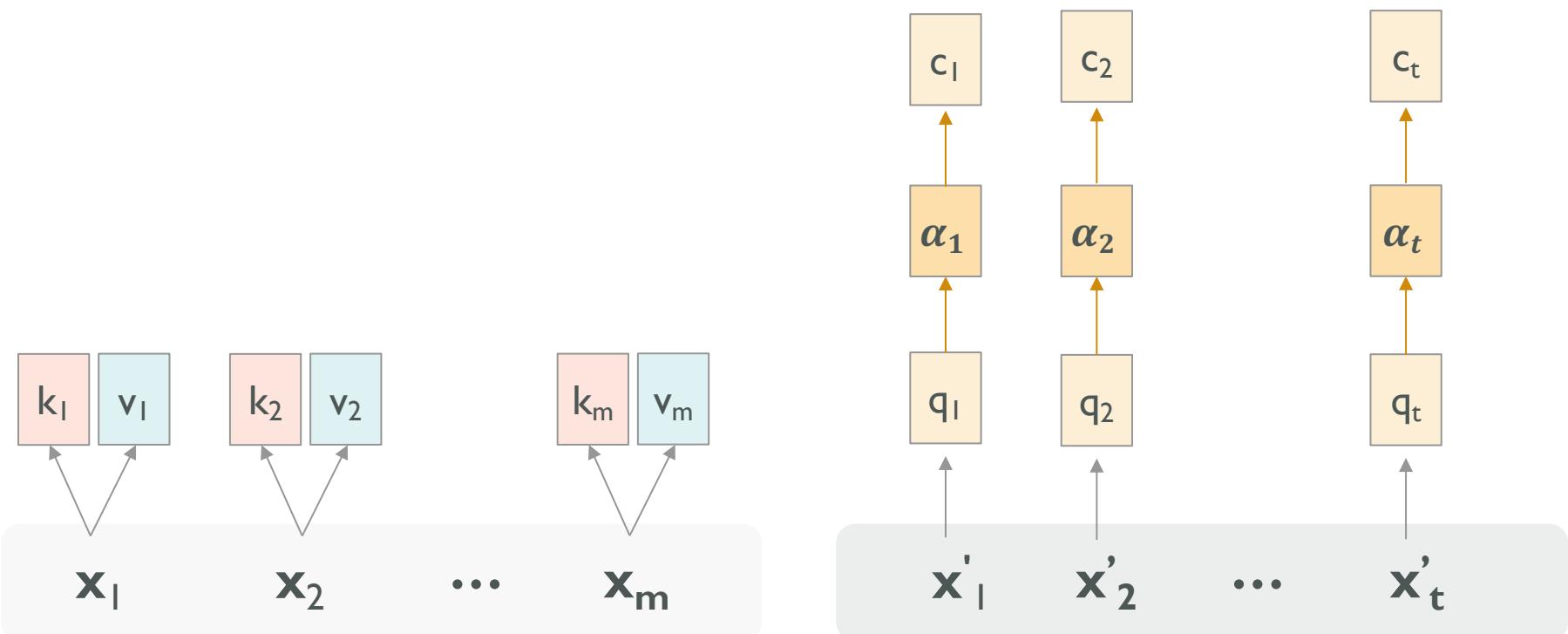


Attention layer

Compute attention weights: $\alpha_i = \text{softmax} (K^T q_i)$

Compute context vector: $c_i = V^T \alpha_i$

Output of attention: $C = [c_1, c_2, \dots, c_t]$



Attention layer

Attention layer: $C = attn(X, X')$

- Encoder's input: $X = [x_1, x_2, \dots, x_m]$
- Decoder's input: $X' = [x'_1, x'_2, \dots, x'_m]$

Attention Layer (Parameters: W_Q, W_K, W_V)

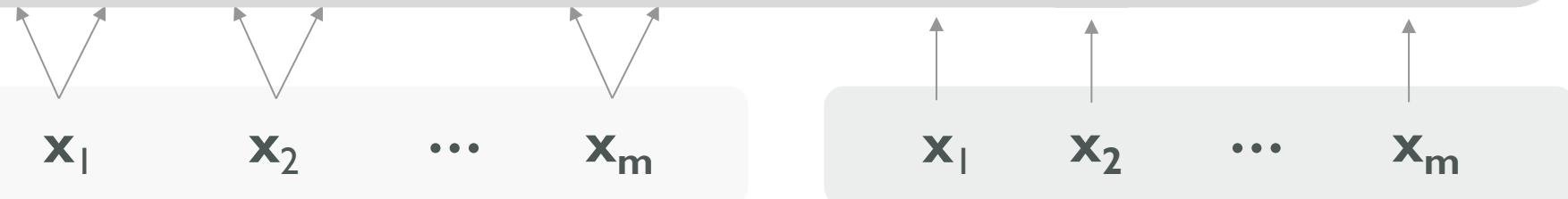


Self-attention layer

Self-attention layer: $C = attn(X, X)$

- Input: $X = [x_1, x_2, \dots, x_m]$

Self-Attention Layer (Parameters: W_Q, W_K, W_V)

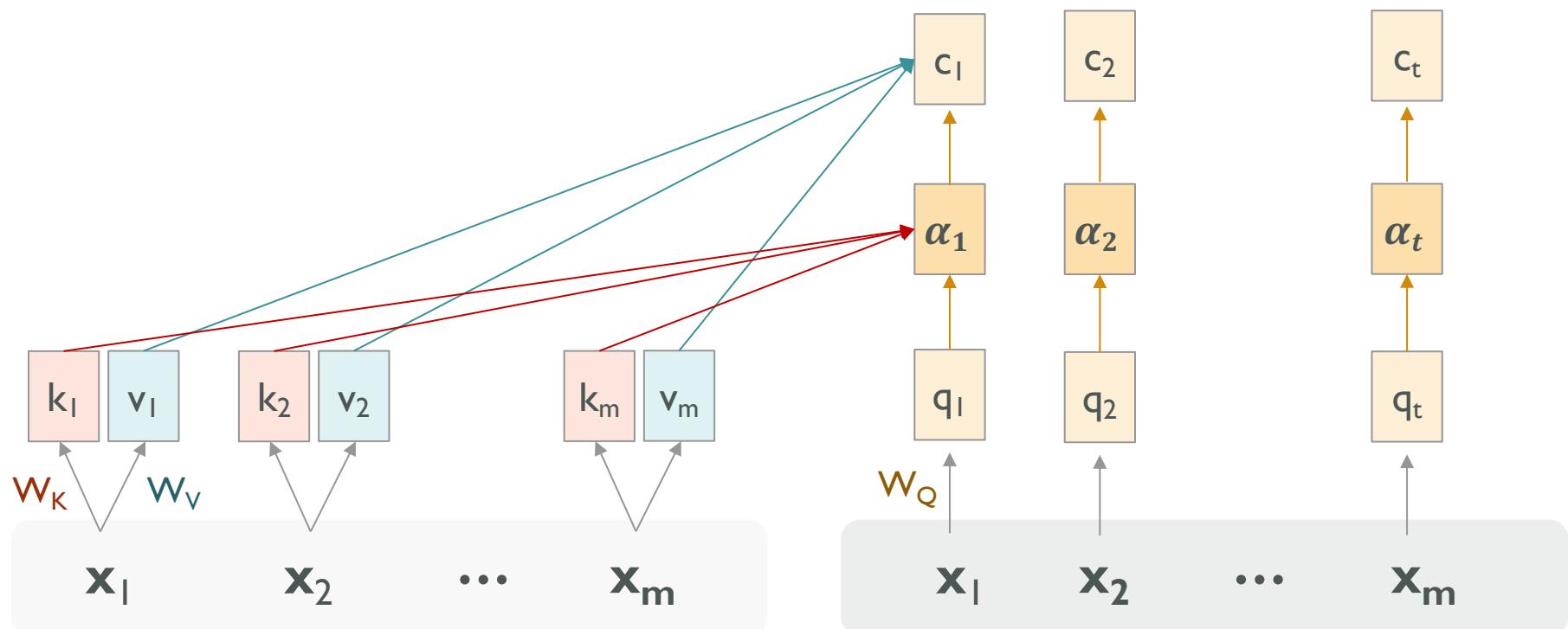


Self-attention layer

Compute attention weights: $\alpha_i = \text{softmax} (K^T q_i)$

Compute context vector: $c_i = V^T \alpha_i$

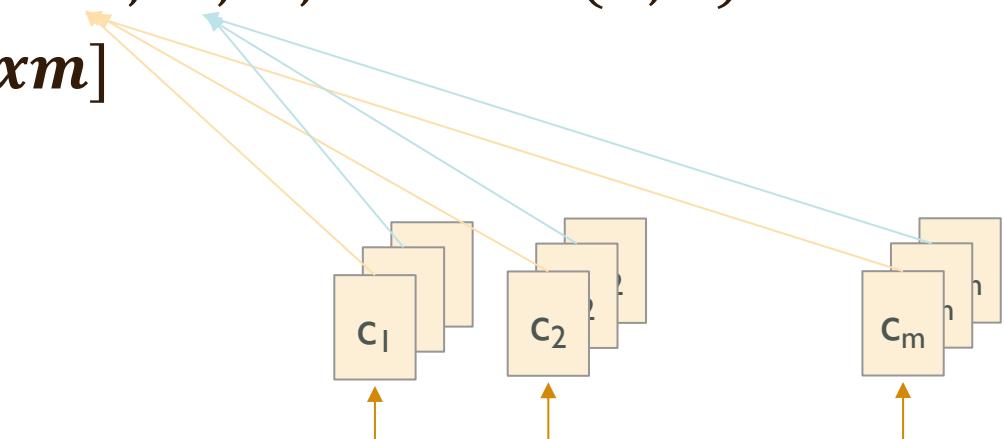
Output of attention: $C = [c_1, c_2, \dots, c_t]$



Multi-head attention layer

Multi-head attention layer: $C^1, C^2, \dots, Cl = mha(X, X)$

- Input: $X = [x_1, x_2, \dots, xm]$



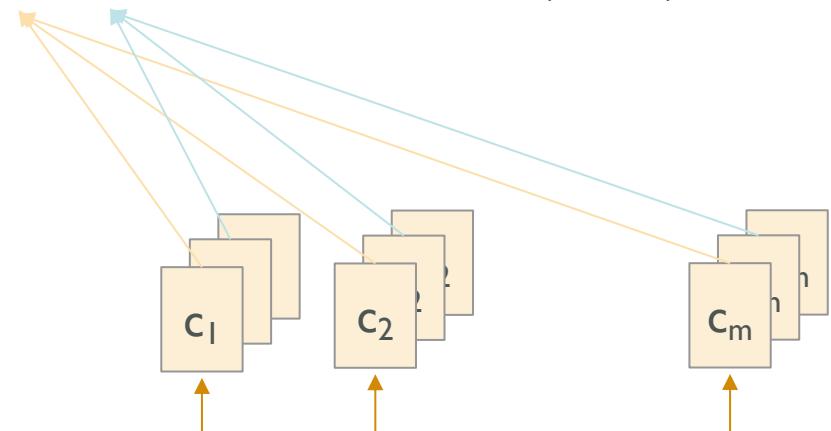
Multi-Head Attention Layer (Parameters: $\{W_{Qi}, W_{Ki}, W_{Vi}\}_i$)



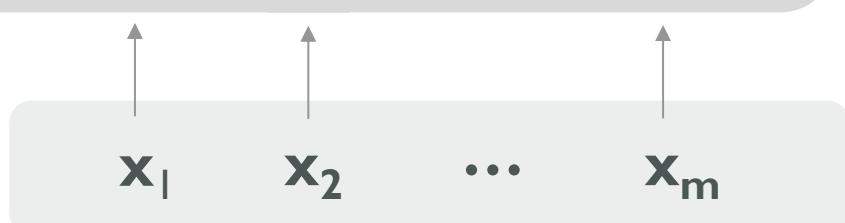
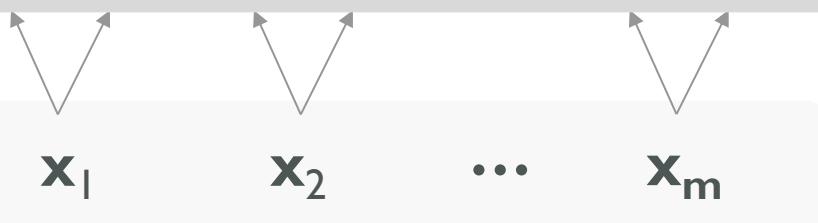
Multi-head self-attention layer

Multi-head self-attention layer: $C^1, C^2, \dots, C^l = mhsa(X, X)$

- Input: $X = [x_1, x_2, \dots, x_m]$



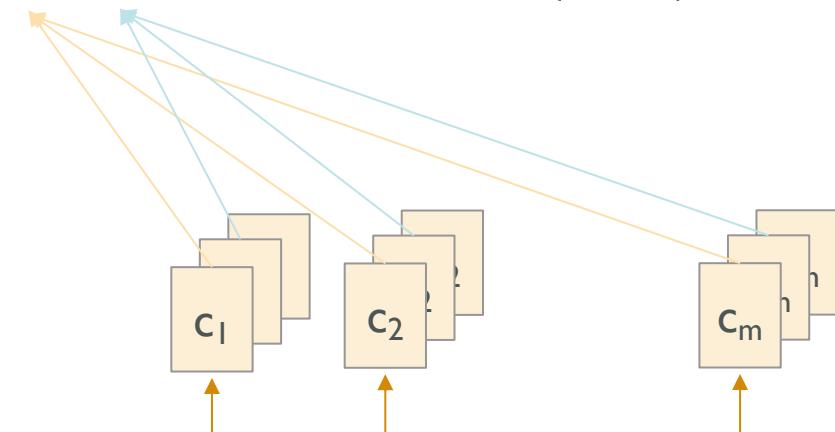
Multi-Head Self-Attention Layer (Parameters: $\{W_{Qi}, W_{Ki}, W_{Vi}\}_i$)



Multi-head self-attention layer

Multi-head self-attention layer: $C^1, C^2, \dots, C^l = mhsa(X, X)$

- Input: $X = [x_1, x_2, \dots, x_m]$



Multi-Head Self-Attention Layer (Parameters: $\{W_{Qi}, W_{Ki}, W_{Vi}\}_i$)



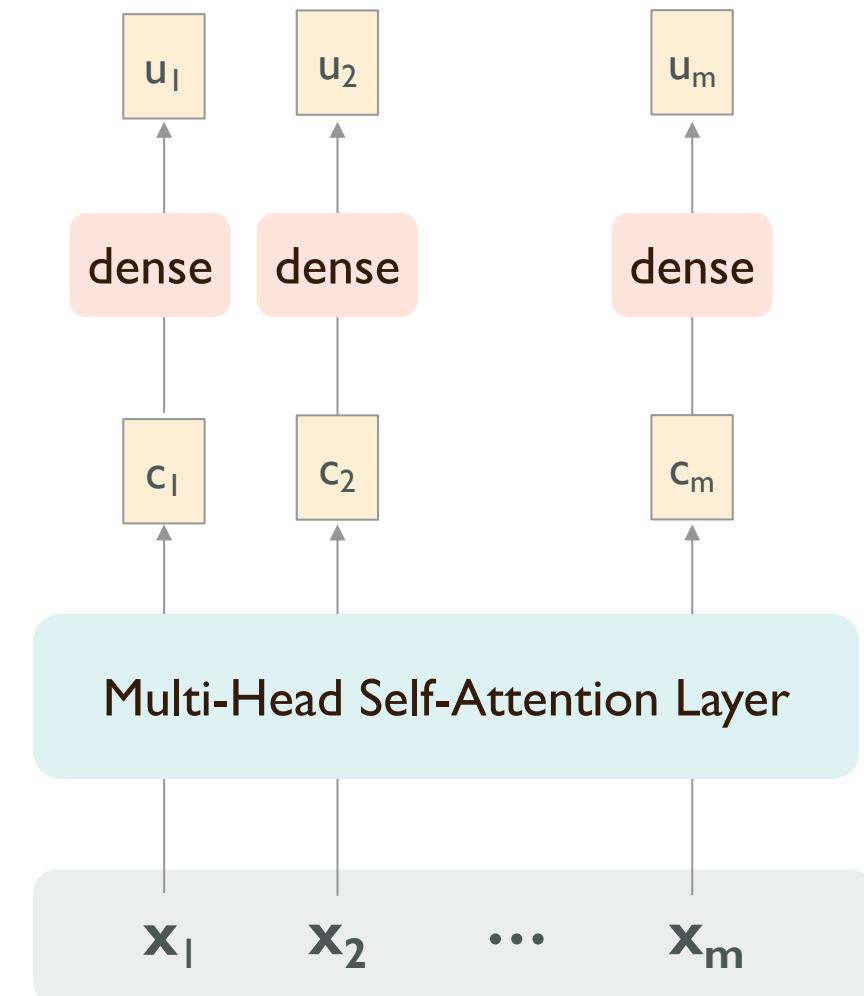
Transformer

One **block** of encoder

$512 \times m$ matrix

Dense layer
Self-attention layer

$X: 512 \times t$ matrix



Transformer

One **block** of encoder

$512 \times m$ matrix

Dense layer

Self-attention layer

$X: 512 \times t$ matrix

Transformer's encoder

Block 6

Block 2

Block 1

$X: 512 \times t$ matrix

Transformer

One **block** of decoder

$512 \times t$ matrix

Dense layer

Attention layer

Self-attention layer

$X: 512 \times m$ matrix

$X': 512 \times t$ matrix

Transformer's decoder

Block 6

Block 2

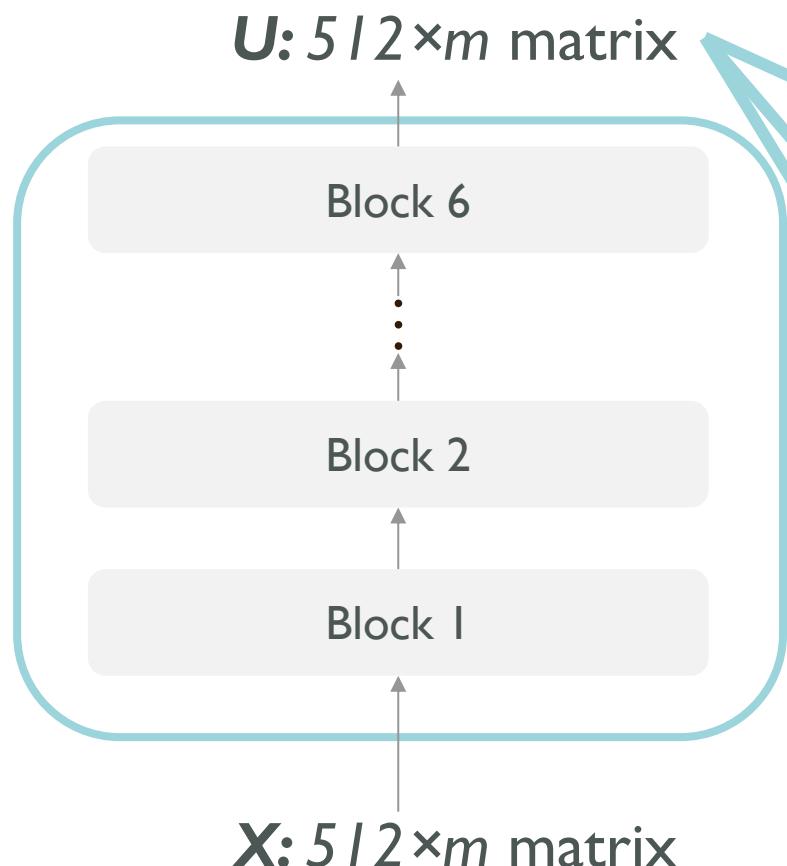
Block 1

from encoder

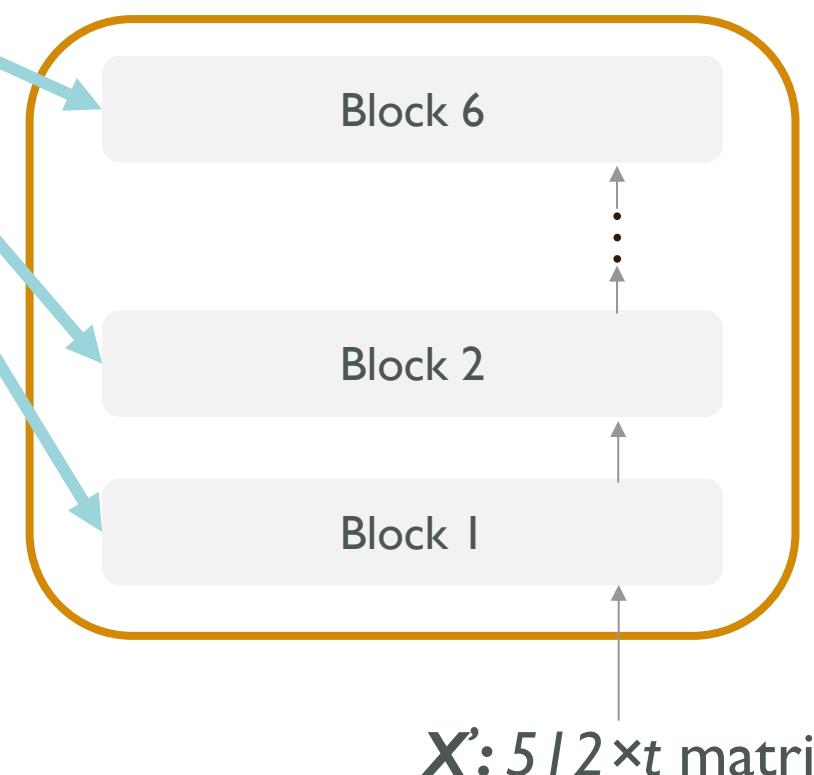
$X': 512 \times t$ matrix

Transformer

Transformer's encoder

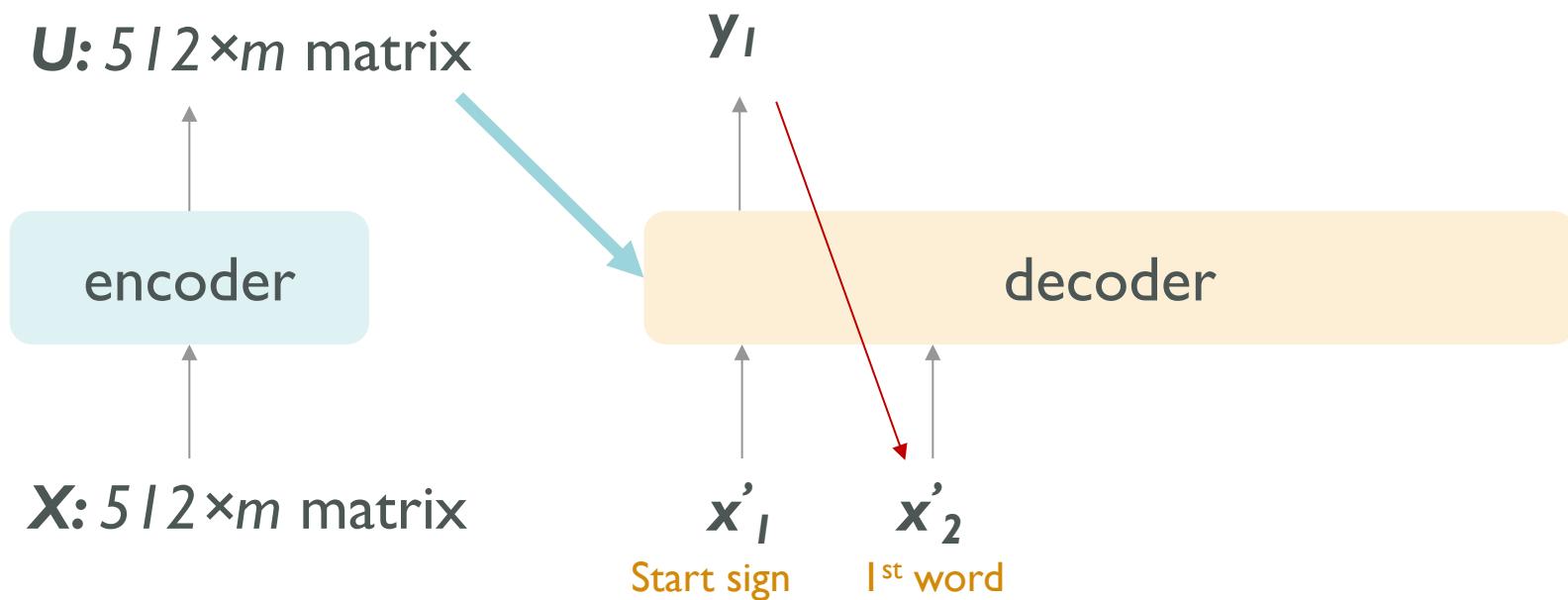


Transformer's decoder



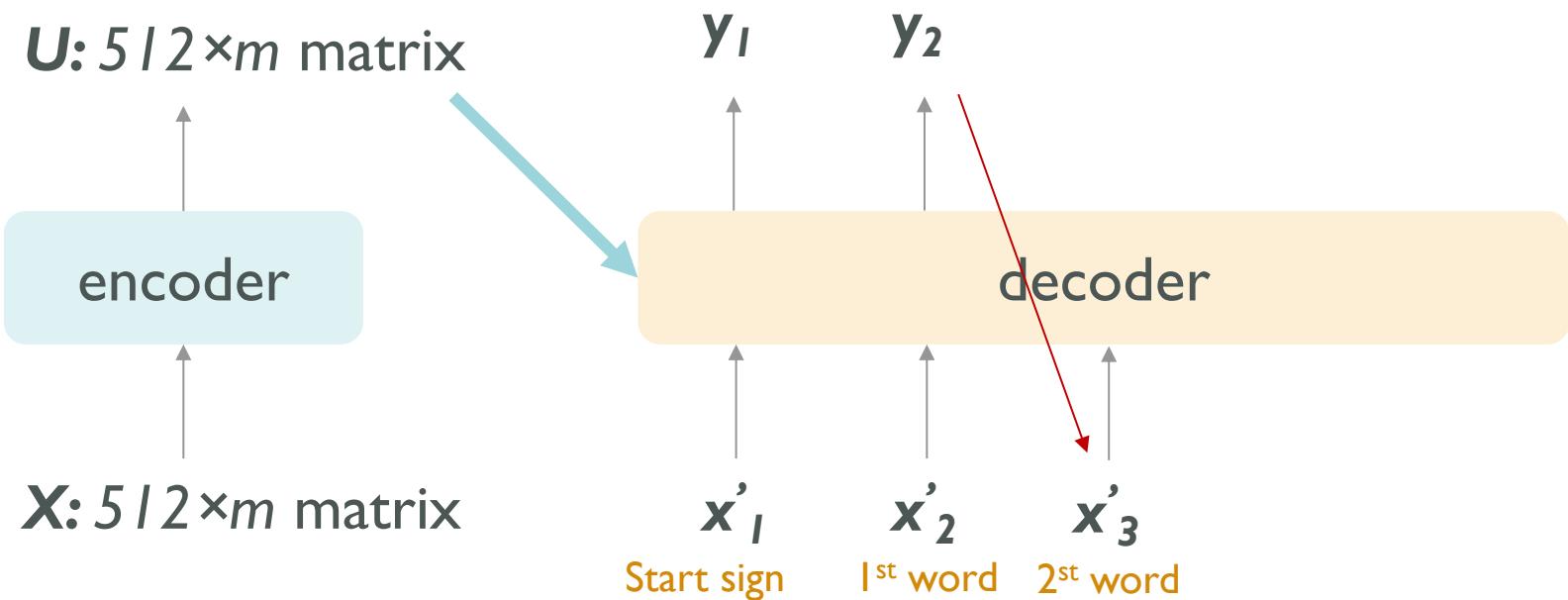
Transformer

Example: Translation



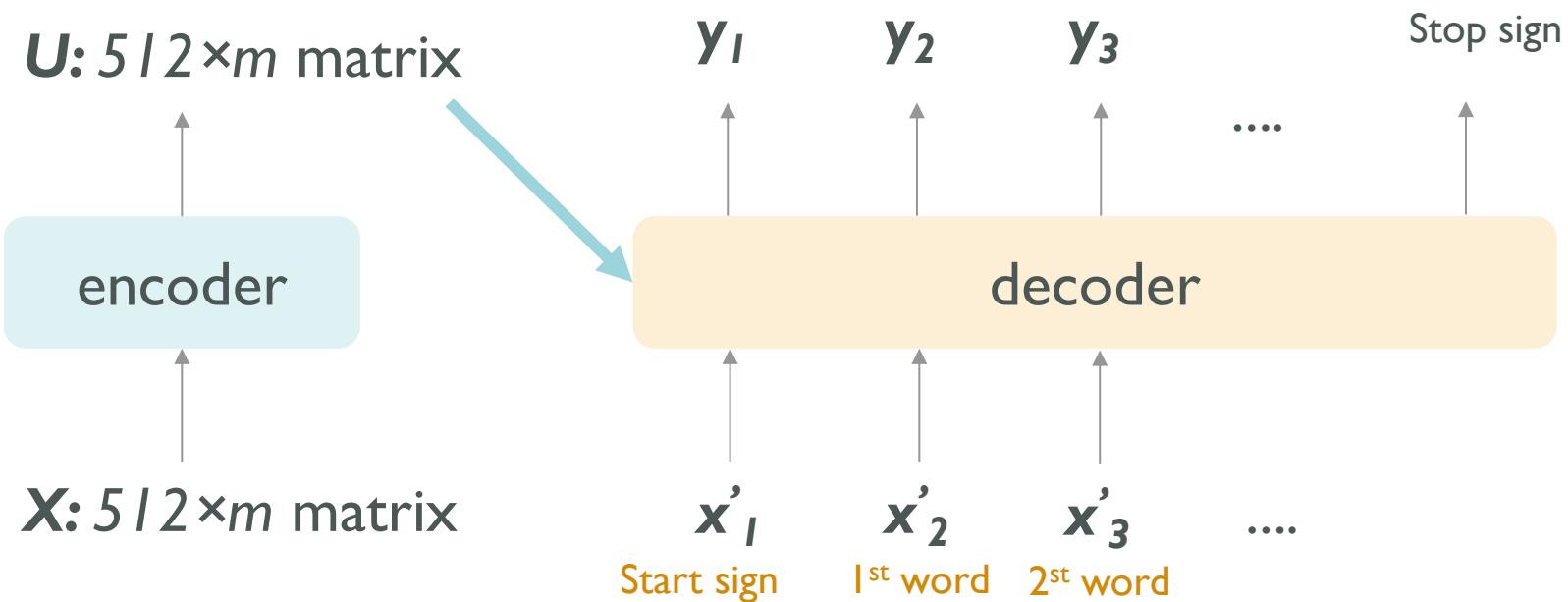
Transformer

Example: Translation



Transformer

Example: Translation



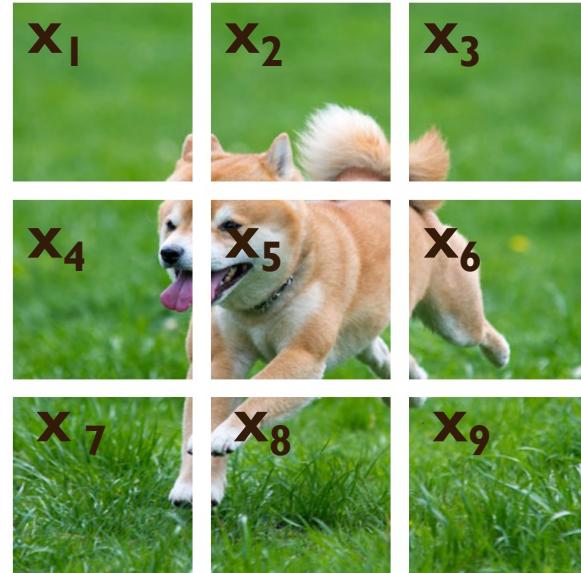
Transformer

- Transformer is Seq2Seq model; it has an encoder and a decoder.
- Transformer model is not RNN
- Transformer is based on attention and self-attention
- Transformer outperforms all the state-of-the-art RNN models.

Vision Transformer (ViT)

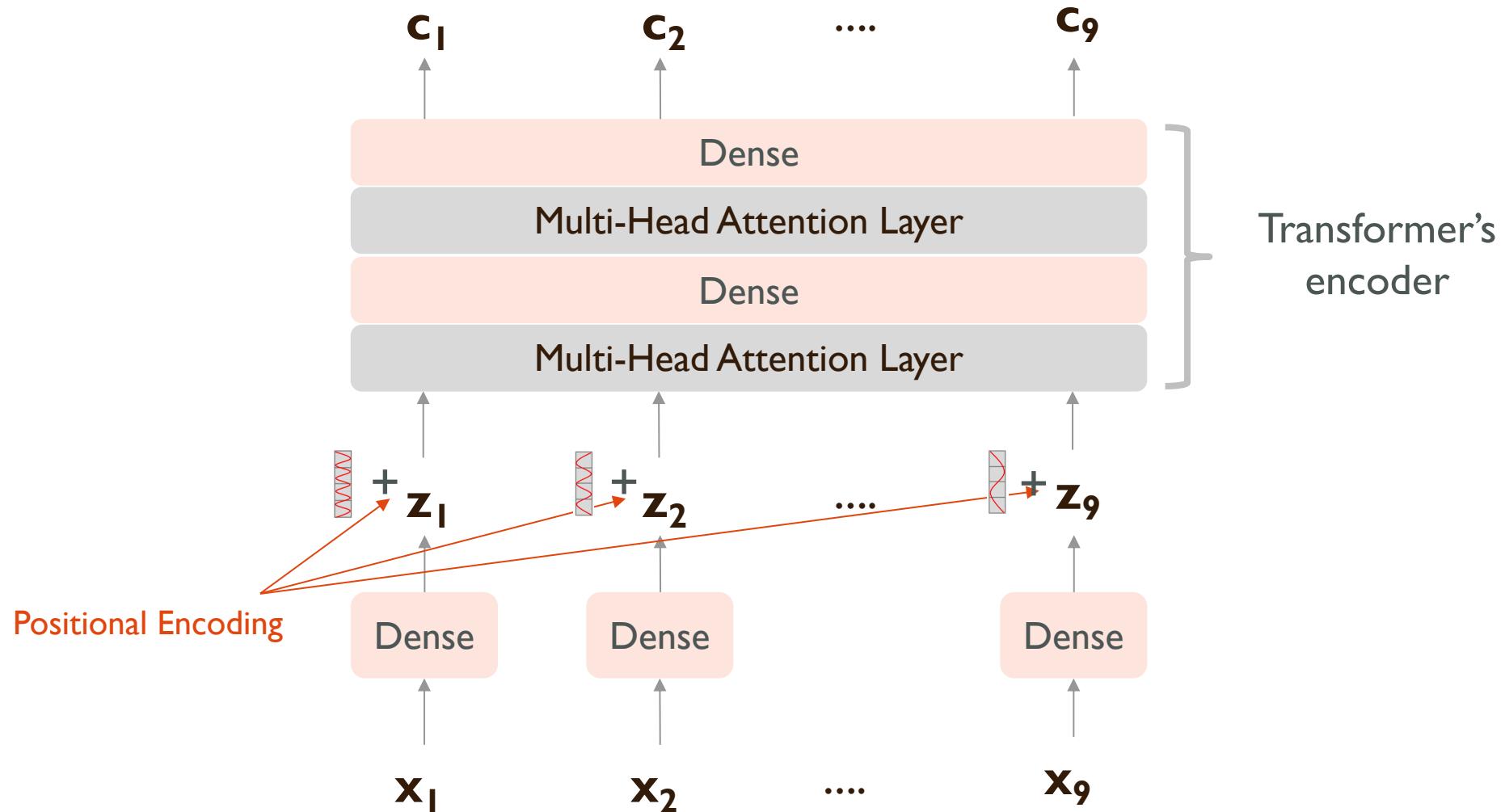
- CNNs, e.g., ResNet, were the best solutions to image classification
- Vision Transformer (ViT) [1] beats CNNs (by a small margin), if the dataset for pretraining is sufficiently large (at least 100 million images)
- ViT is based on Transformer (for NLP)

Vision Transformer (ViT)

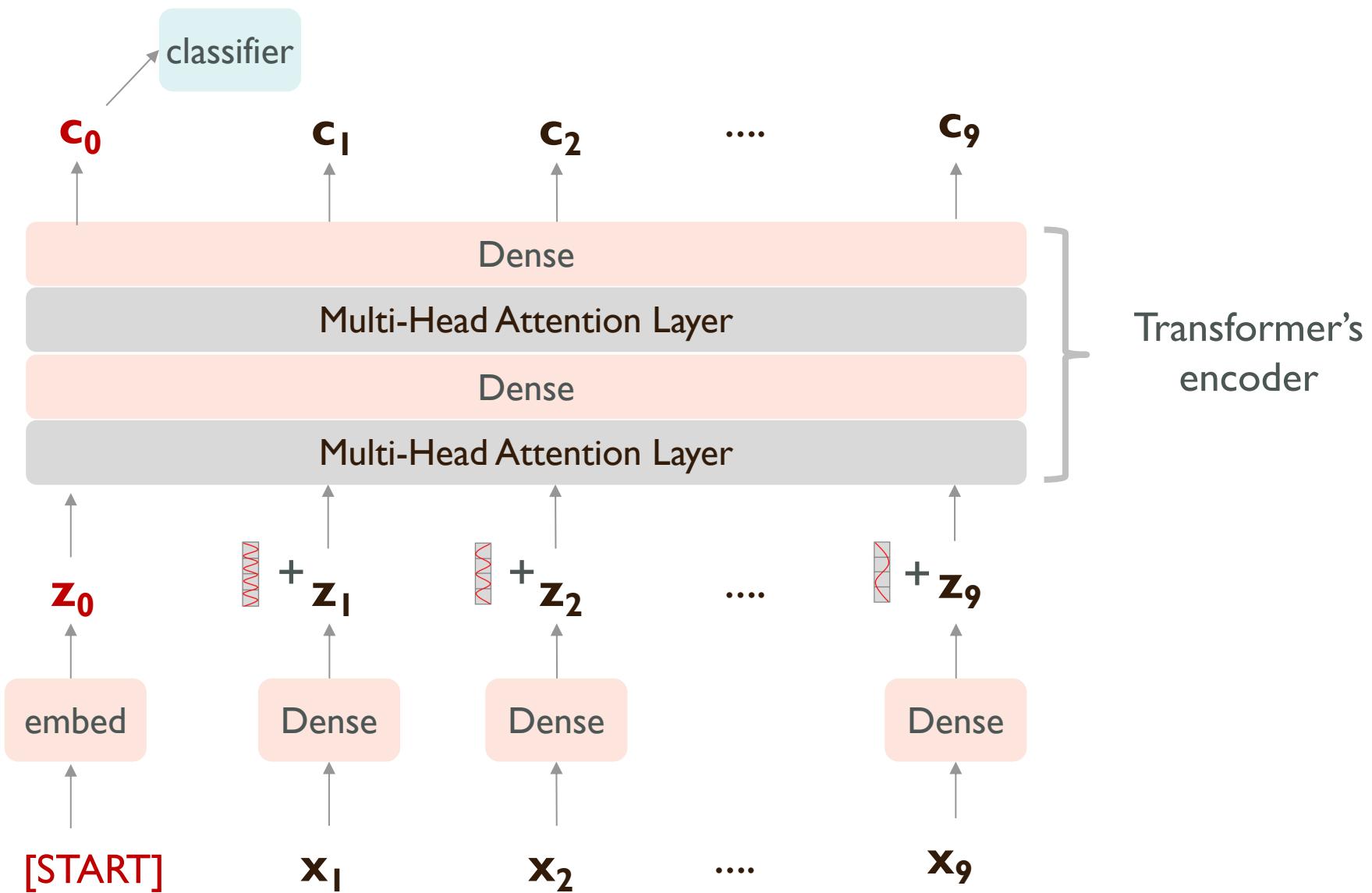


- Represent each image as several patches
- Treat each patch as one word
- Use Transformer's encoder to embed each patch

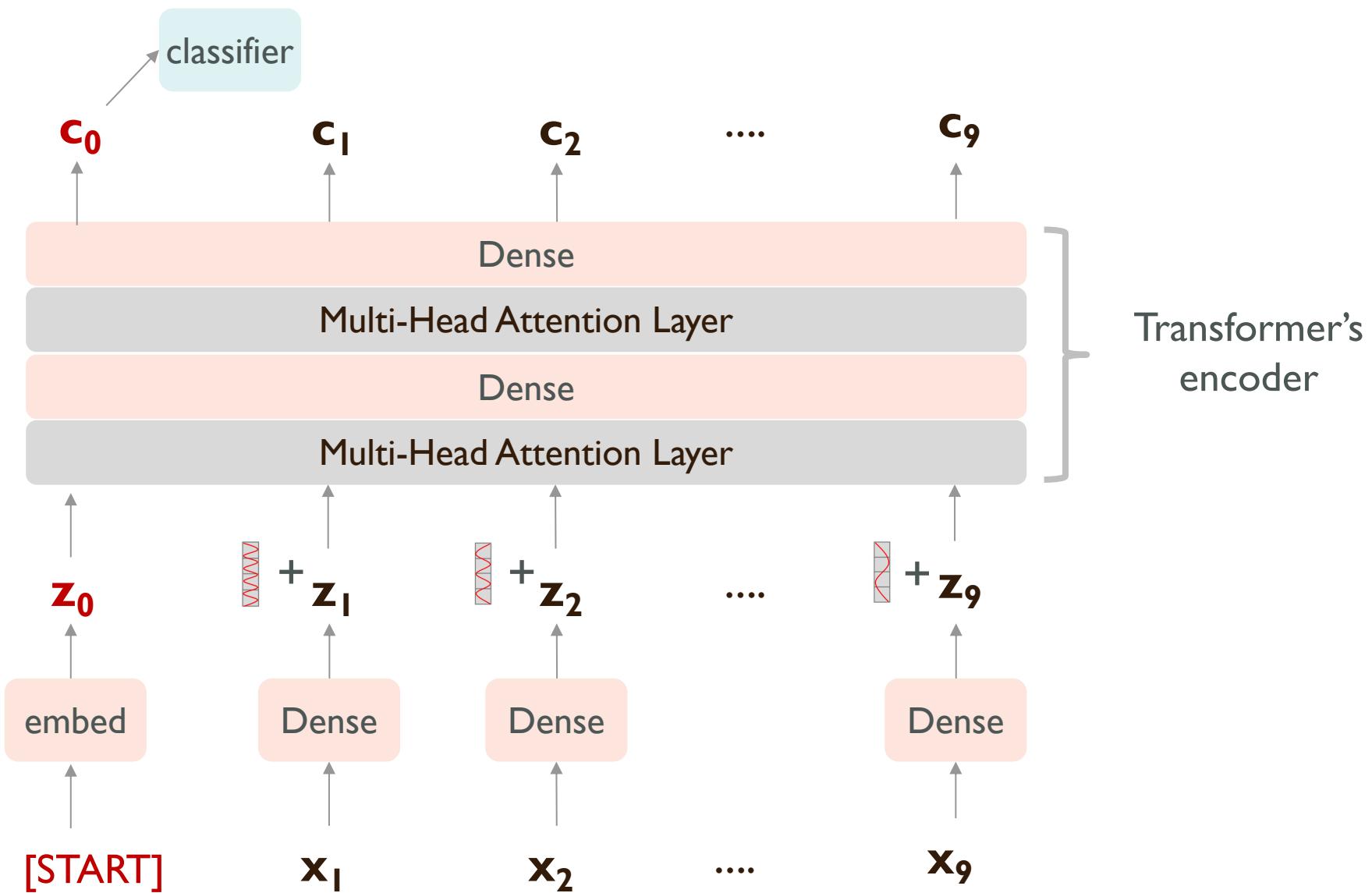
Vision Transformer (ViT)



Vision Transformer (ViT)



Vision Transformer (ViT)



Vision Transformer (ViT)

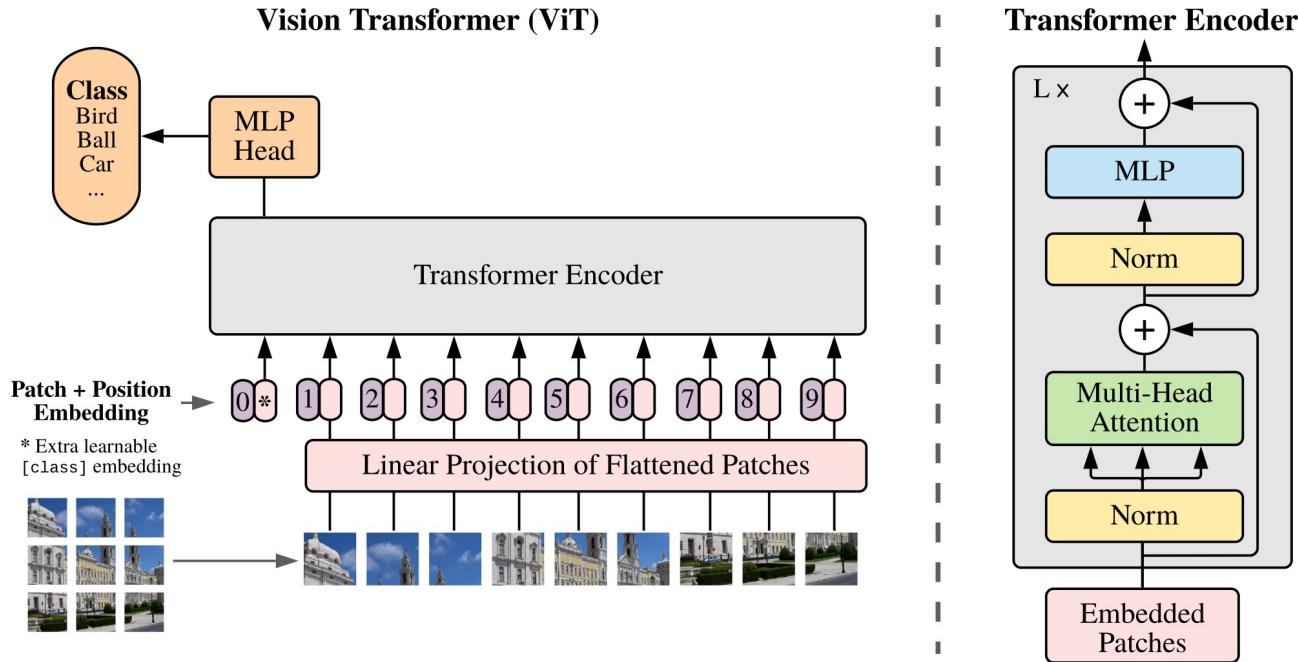
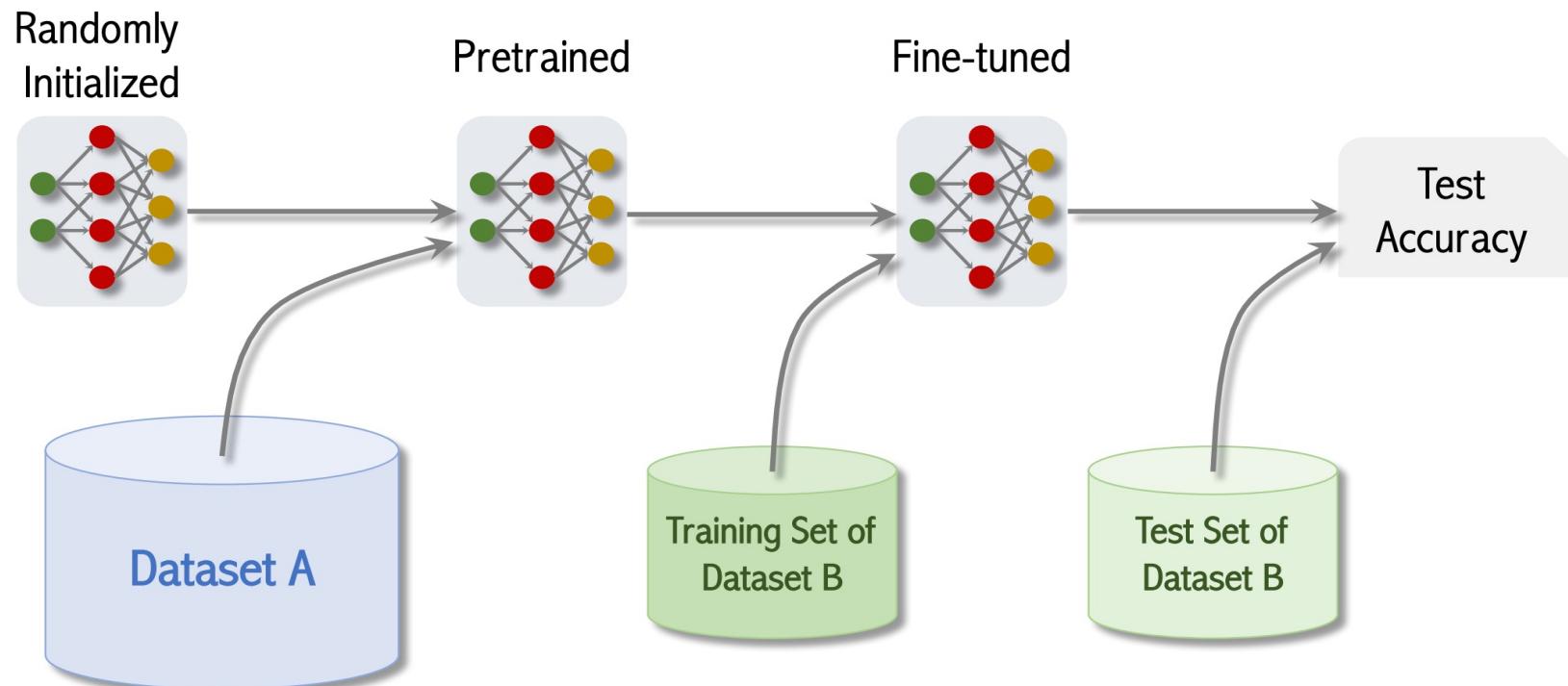
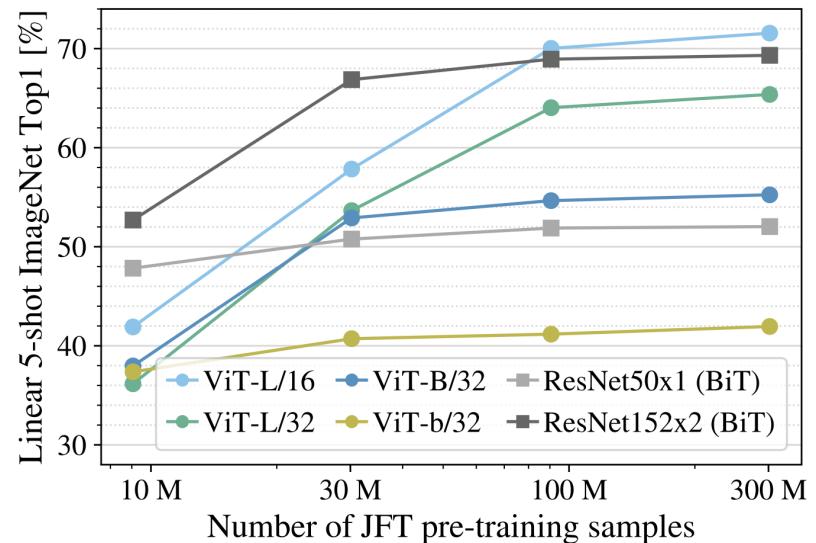
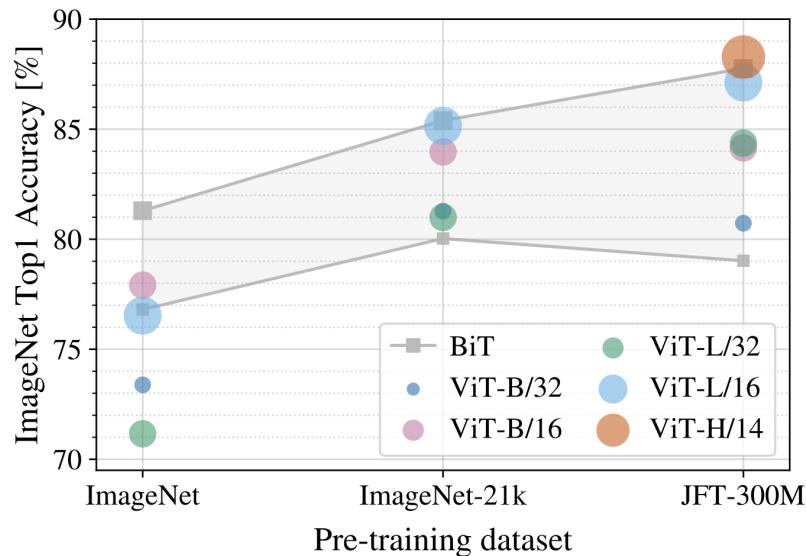


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

Vision Transformer (ViT)



Vision Transformer (ViT)



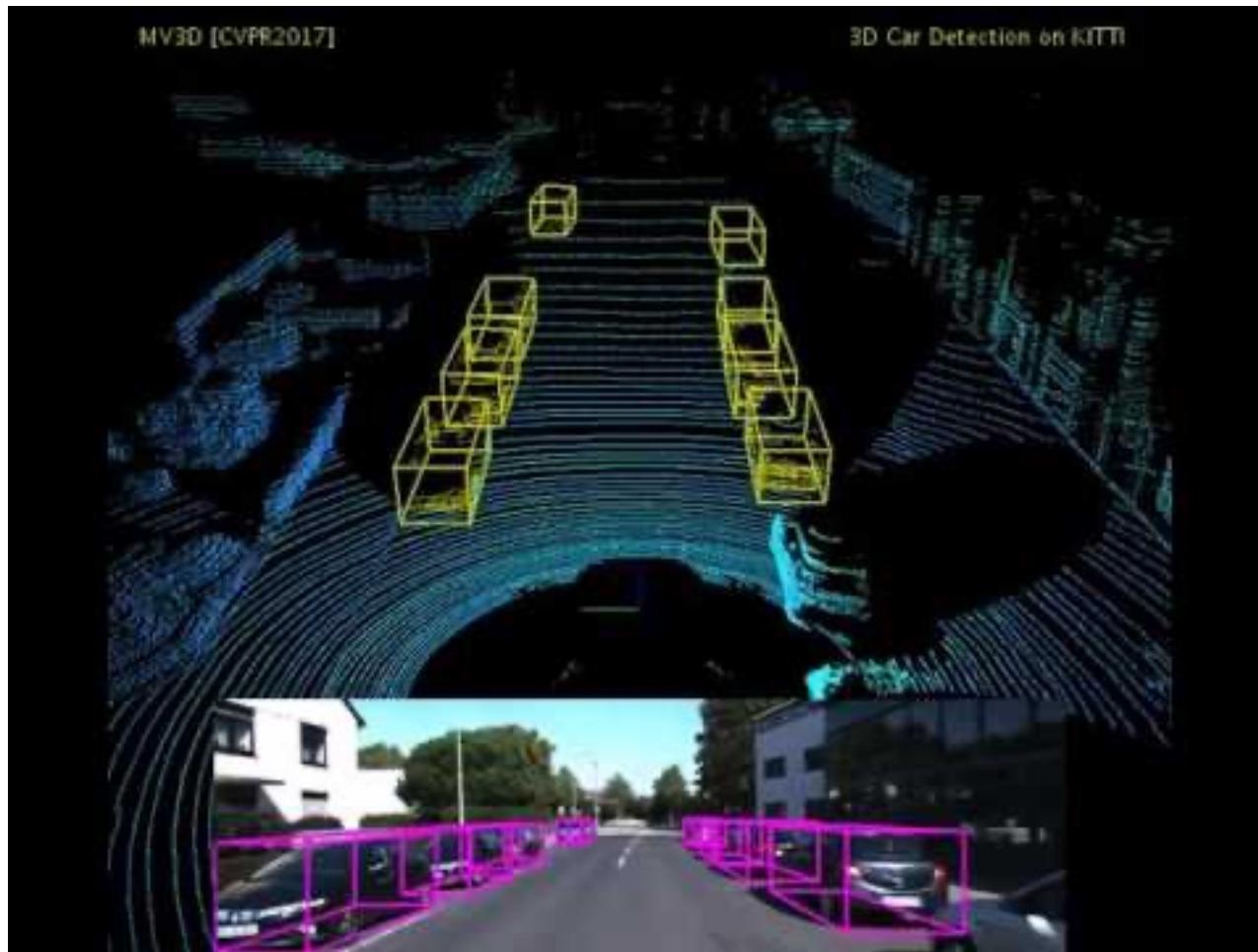
- Pretrained on ImageNet (small), ViT is slightly worse than ResNet
- Pretrained on ImageNet-21K (medium), ViT is comparable to ResNet
- Pretrained on JFT (large), ViT is slightly better than ResNet.

Vision Transformer (ViT)

Pretrain the model on Dataset A, fine-tune the model on Dataset B, and evaluate the model on Dataset B.

	# of Images	# of classes
ImageNet (Small)	1.3 Million	1 Thousand
ImageNet-21K (Medium)	14 Million	21 Thousand
JFT (Big)	300 Million	18 Thousand

Next lecture: Object detection



Thank you very much!

sihengc@sjtu.edu.cn