

Transformer: Part I (Encoder)

CS7150, Spring 2025

Prof. Huaizu Jiang

Northeastern University

Announcement

- Fixing an error in Task 2.7 of PA2
 - Check the announcement on Canvas
- Grades and feedback of the project proposal will be released by the end of Friday

Recap

Recall: Batch Normalization for ConvNets

Batch Normalization for
fully-connected networks

$$\begin{array}{l} x : N \times D \\ \text{Normalize} \downarrow \\ \mu, \sigma : 1 \times D \\ \gamma, \beta : 1 \times D \\ y = \frac{(x - \mu)}{\sigma} \gamma + \beta \end{array}$$

Batch Normalization for
convolutional networks
(Spatial Batchnorm, BatchNorm2D)

$$\begin{array}{l} x : N \times C \times H \times W \\ \text{Normalize} \downarrow \quad \downarrow \quad \downarrow \\ \mu, \sigma : 1 \times C \times 1 \times 1 \\ \gamma, \beta : 1 \times C \times 1 \times 1 \\ y = \frac{(x - \mu)}{\sigma} \gamma + \beta \end{array}$$

How to use BatchNorm in an RNN?

Naïve way of applying BN in RNN

Why it is not a good idea?

1. Variable lengths in the input.
2. The recurrent nature of RNN.

Batch Normalization for
recurrent networks

$$\begin{array}{c} x : N \times L \times C \\ \quad \downarrow \quad \downarrow \\ \mu, \sigma : 1 \times 1 \times C \\ \gamma, \beta : 1 \times 1 \times C \\ y = \frac{(x - \mu)}{\sigma} \gamma + \beta \end{array}$$

Layer Normalization

Computing of μ, σ is independent of the batch and length dimension.

But the learnable parameters γ, β are still for each channel.

Largely adopted in Transformer.

Layer Normalization for **recurrent** networks

$$x : N \times L \times C$$

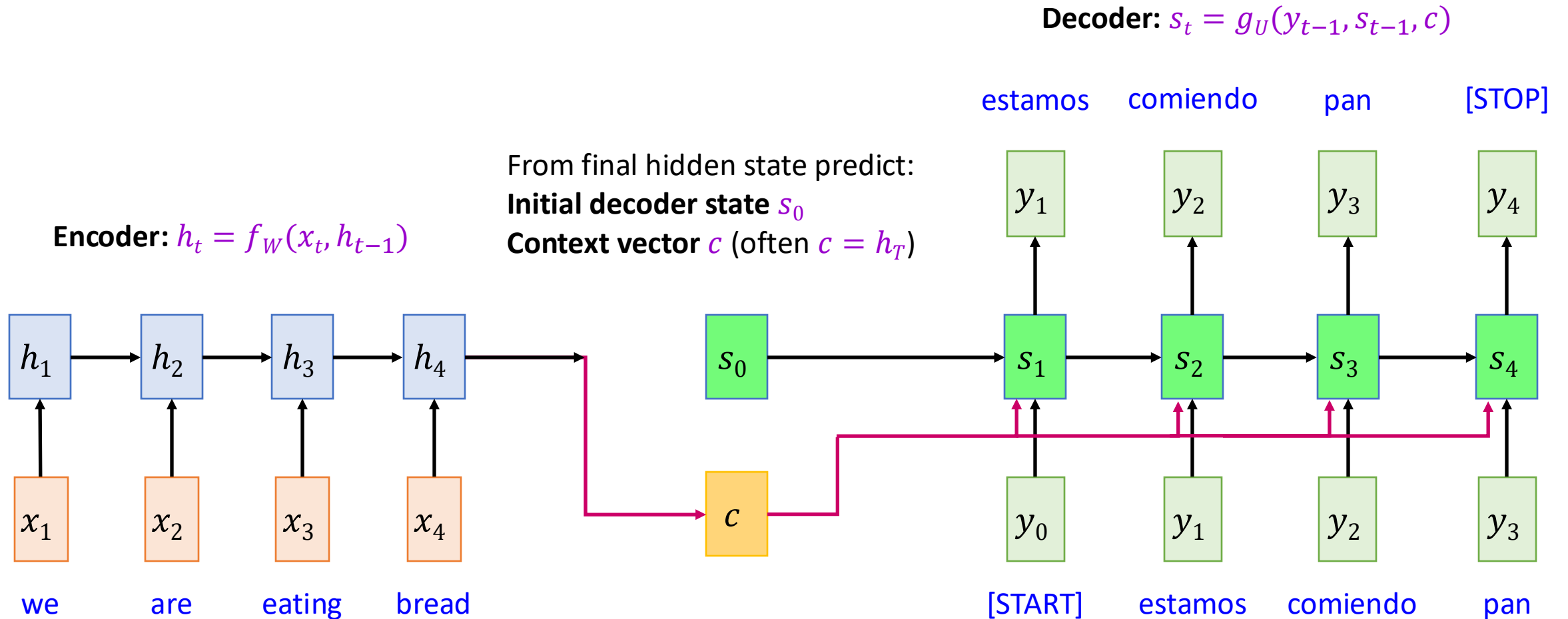


$$\mu, \sigma : N \times L \times 1$$

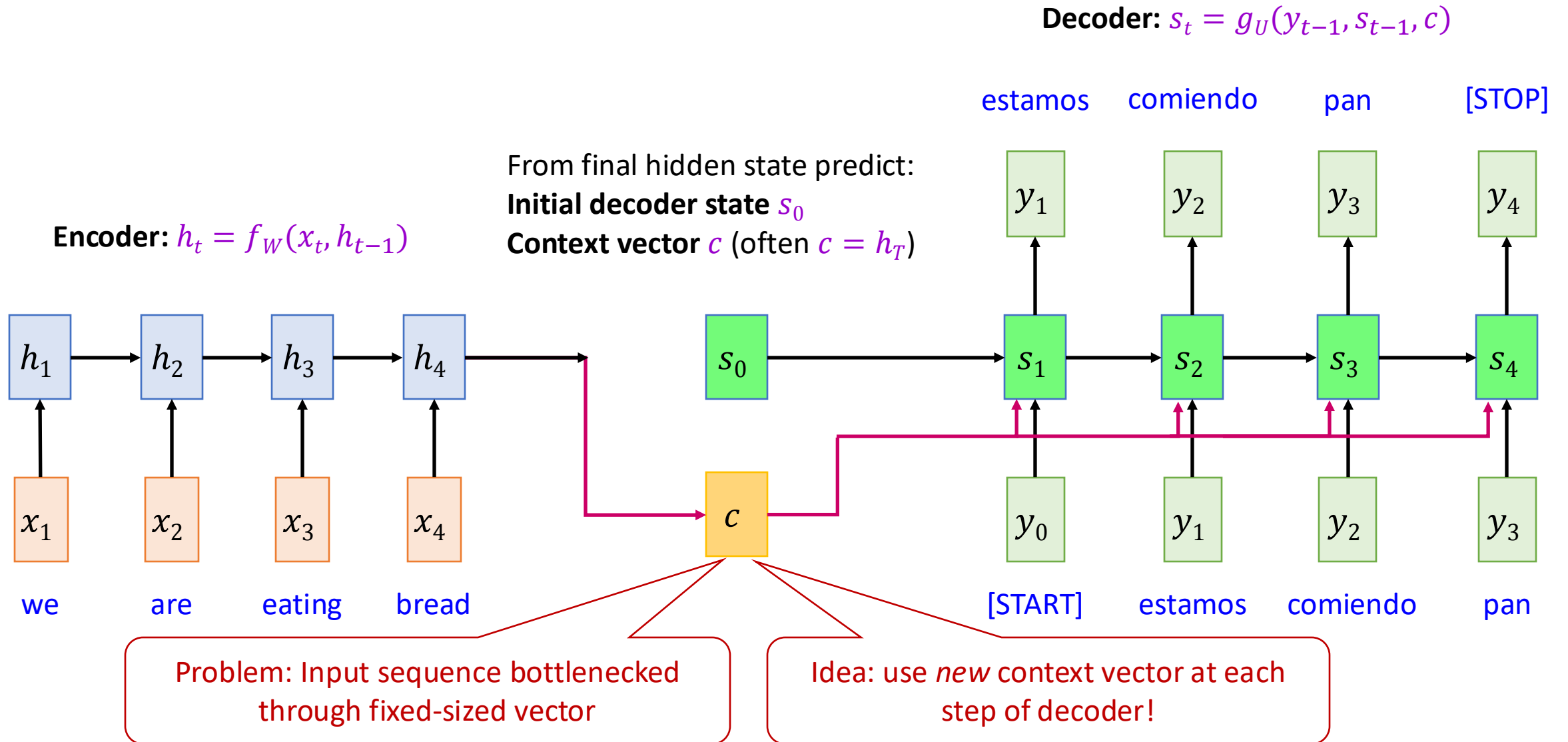
$$\gamma, \beta : 1 \times 1 \times C$$

$$y = \frac{(x - \mu)}{\sigma} \gamma + \beta$$

Sequence-to-sequence with RNNs

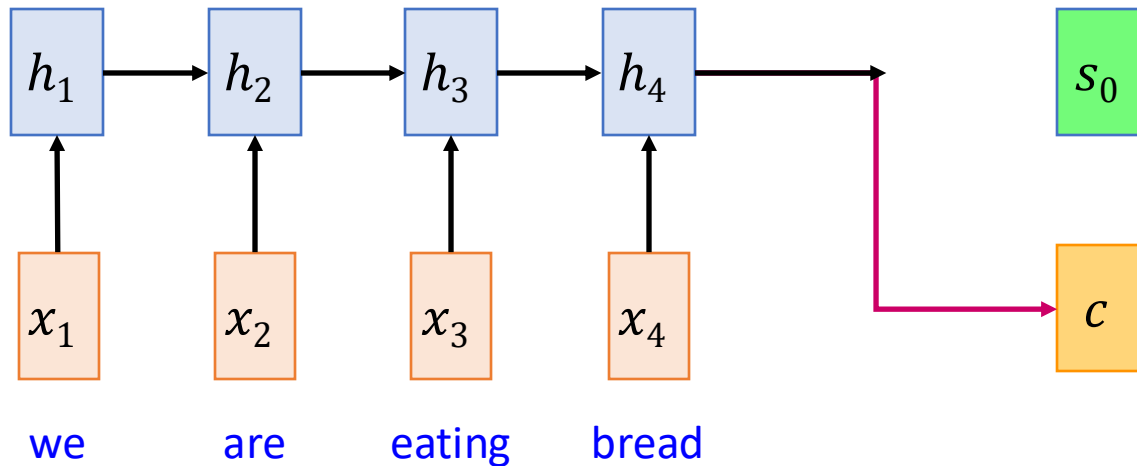


Sequence-to-sequence with RNNs



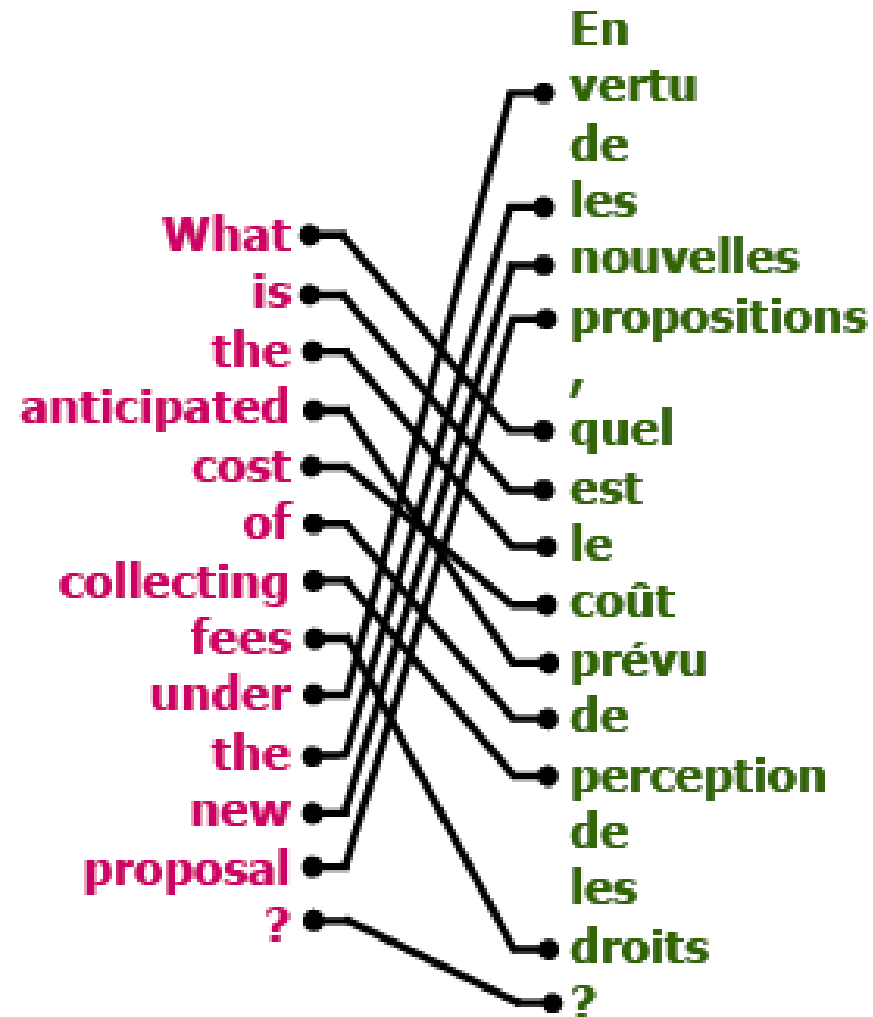
Sequence-to-sequence with RNNs and attention

- At each timestep of decoder, context vector “looks at” different parts of the input sequence

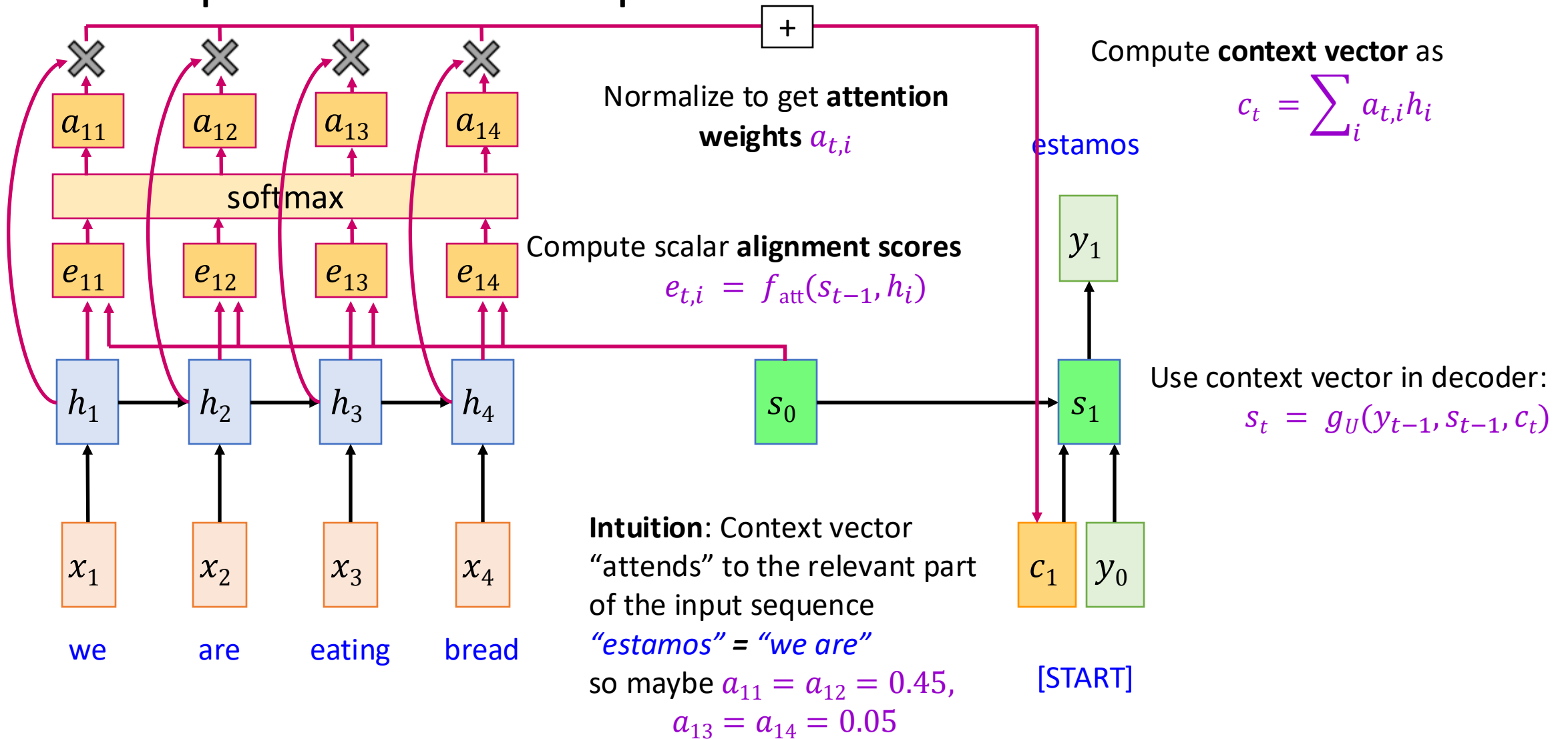


Sequence-to-sequence with RNNs and attention

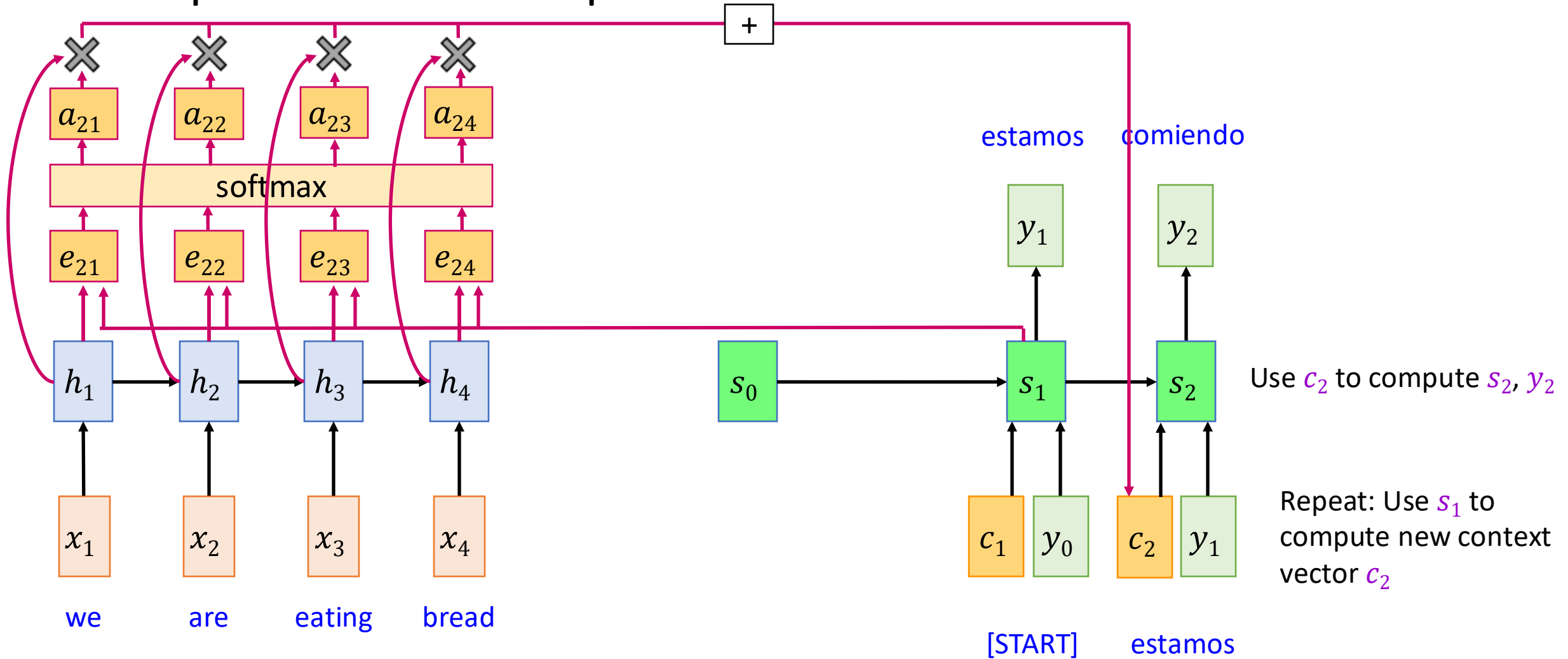
- Intuition: translation requires *alignment*



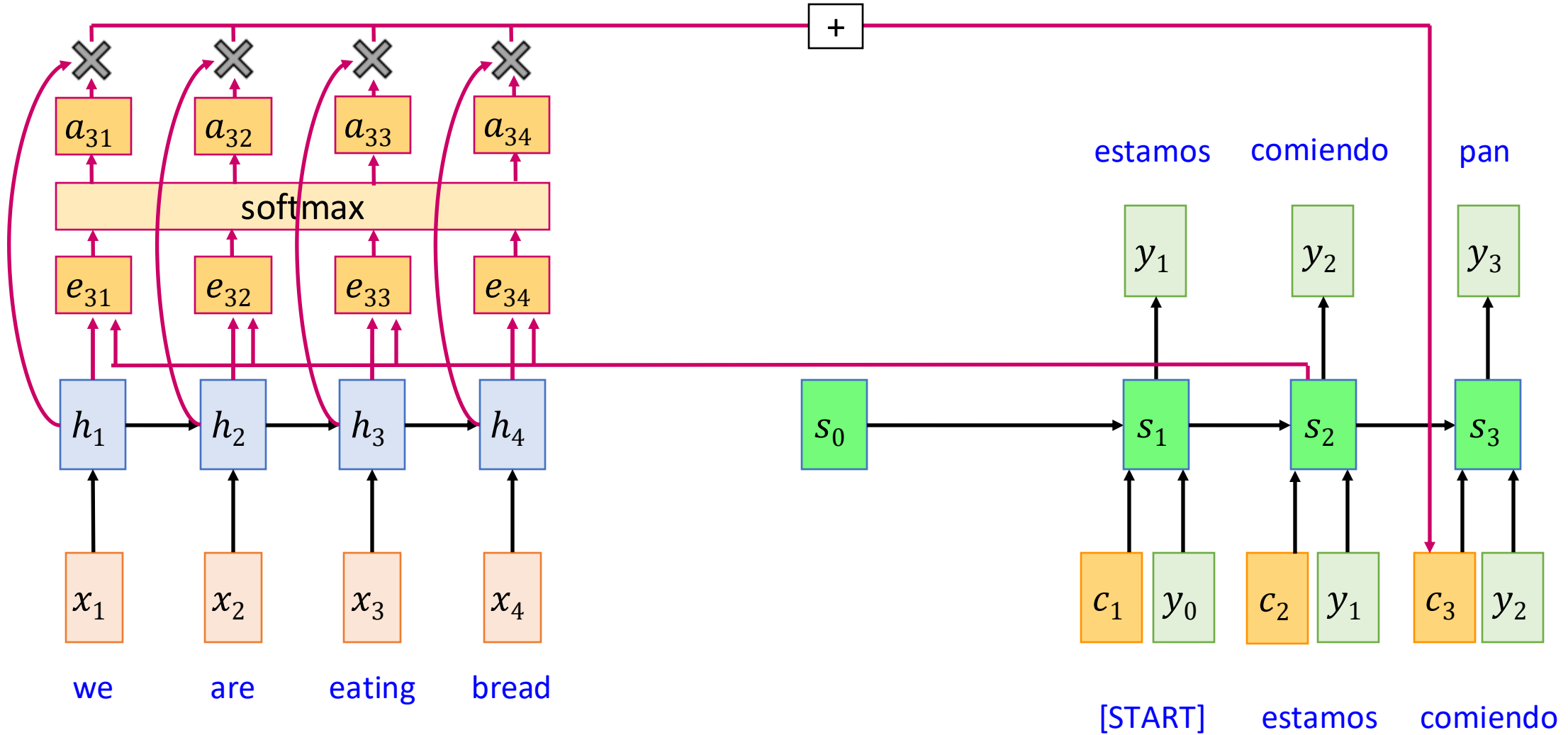
Sequence-to-sequence with RNNs and attention



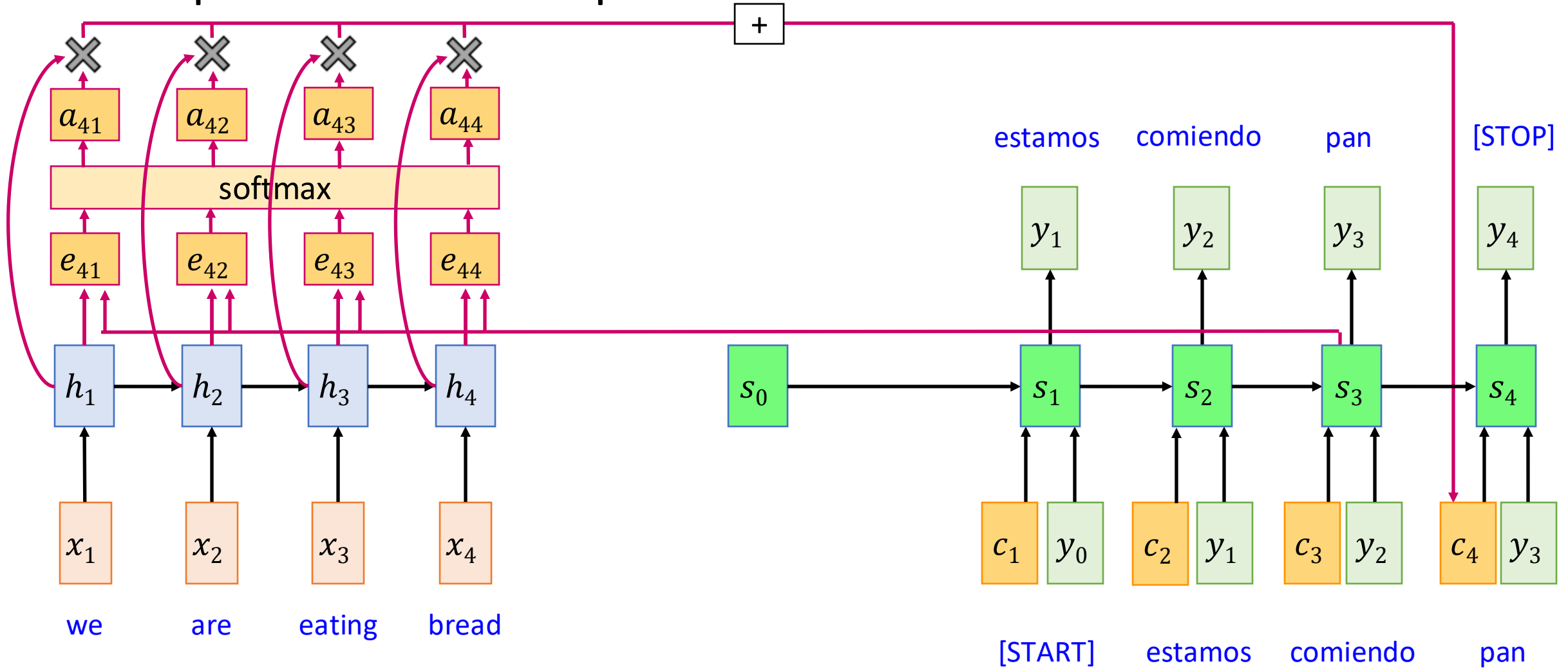
Sequence-to-sequence with RNNs and attention



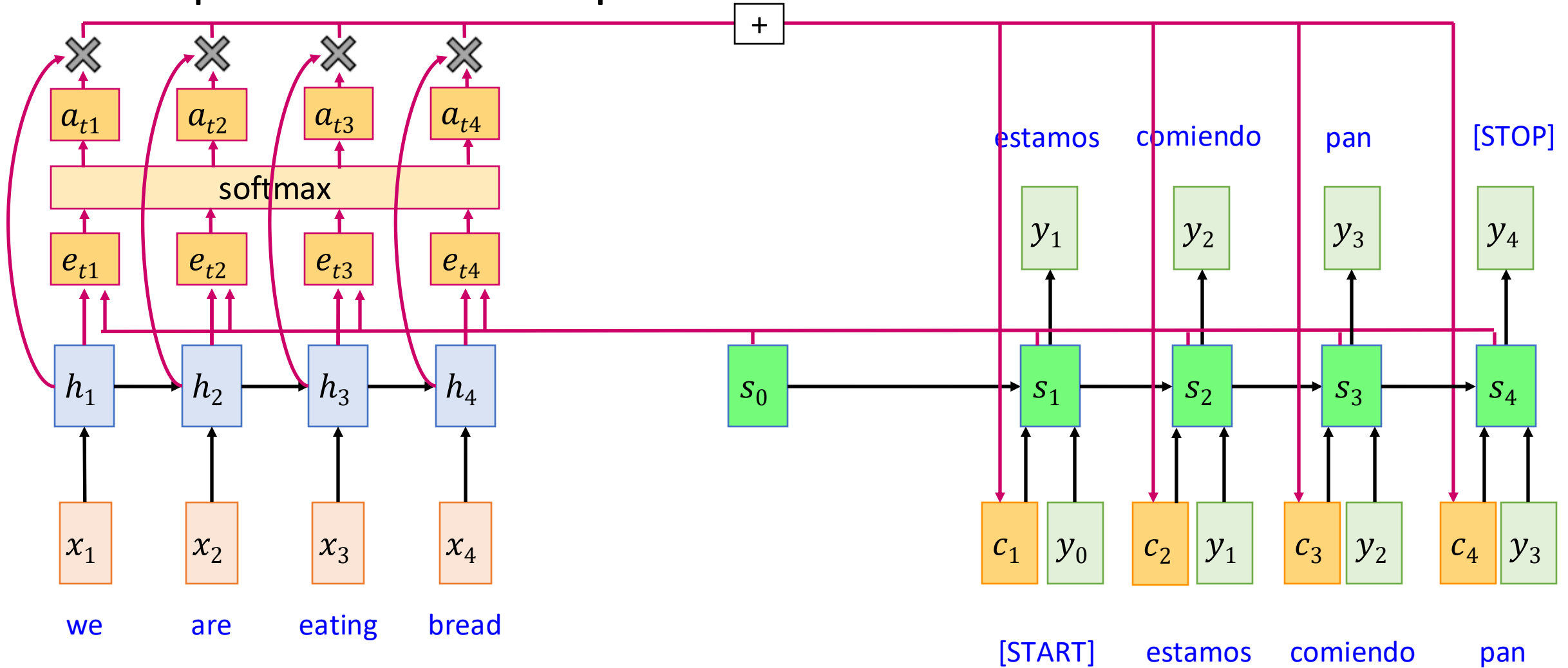
Sequence-to-sequence with RNNs and attention



Sequence-to-sequence with RNNs and attention

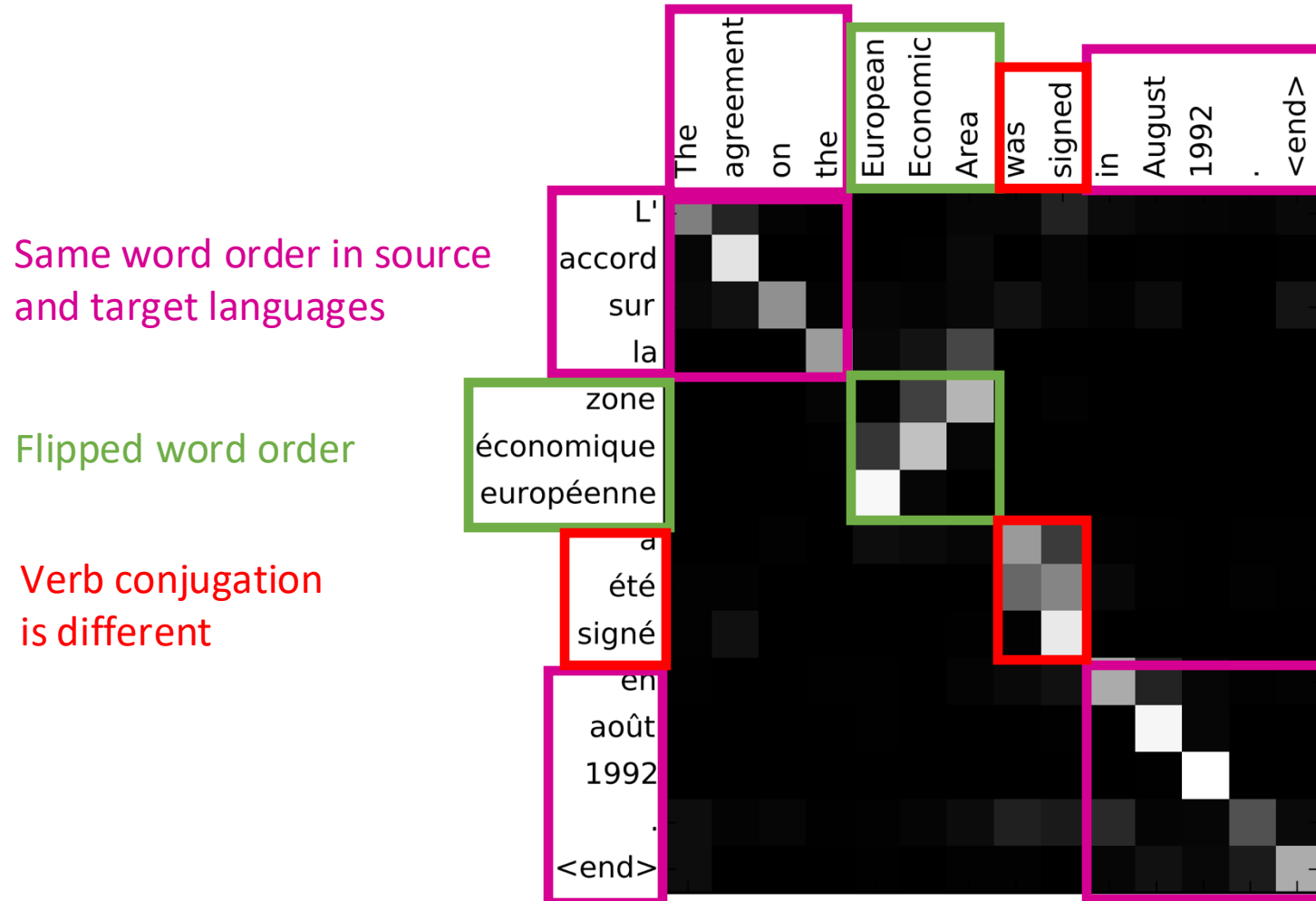


Sequence-to-sequence with RNNs and attention

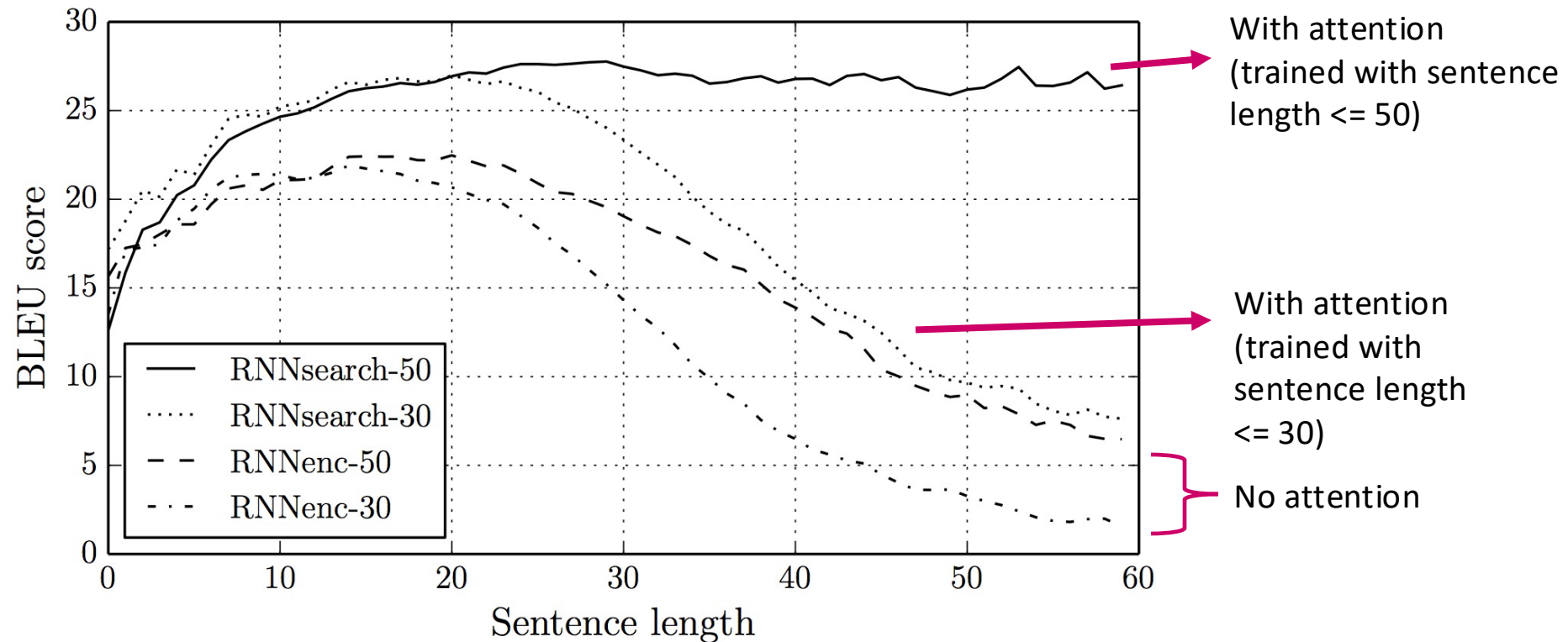


Sequence-to-sequence with RNNs and attention

- Visualizing attention weights (English source, French target):



Quantitative evaluation



D. Bahdanau, K. Cho, Y. Bengio, [Neural Machine Translation by Jointly Learning to Align and Translate](#), ICLR 2015

Attention in a Nutshell

Image captioning with RNNs and attention

- Idea: pay attention to different parts of the image when generating different words
- Automatically learn this *grounding* of words to image regions without direct supervision

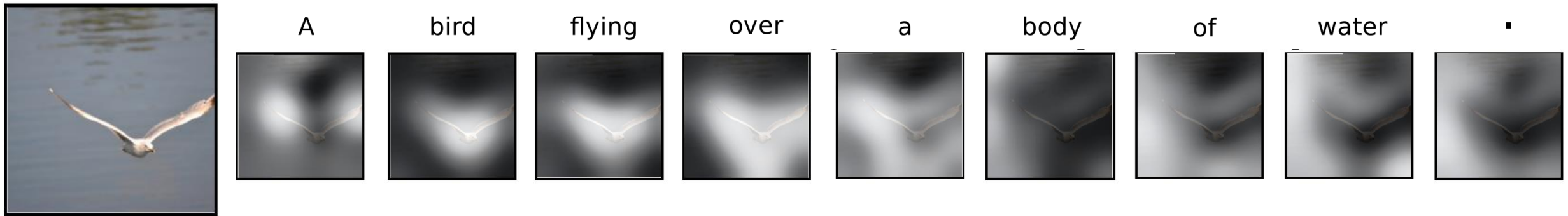
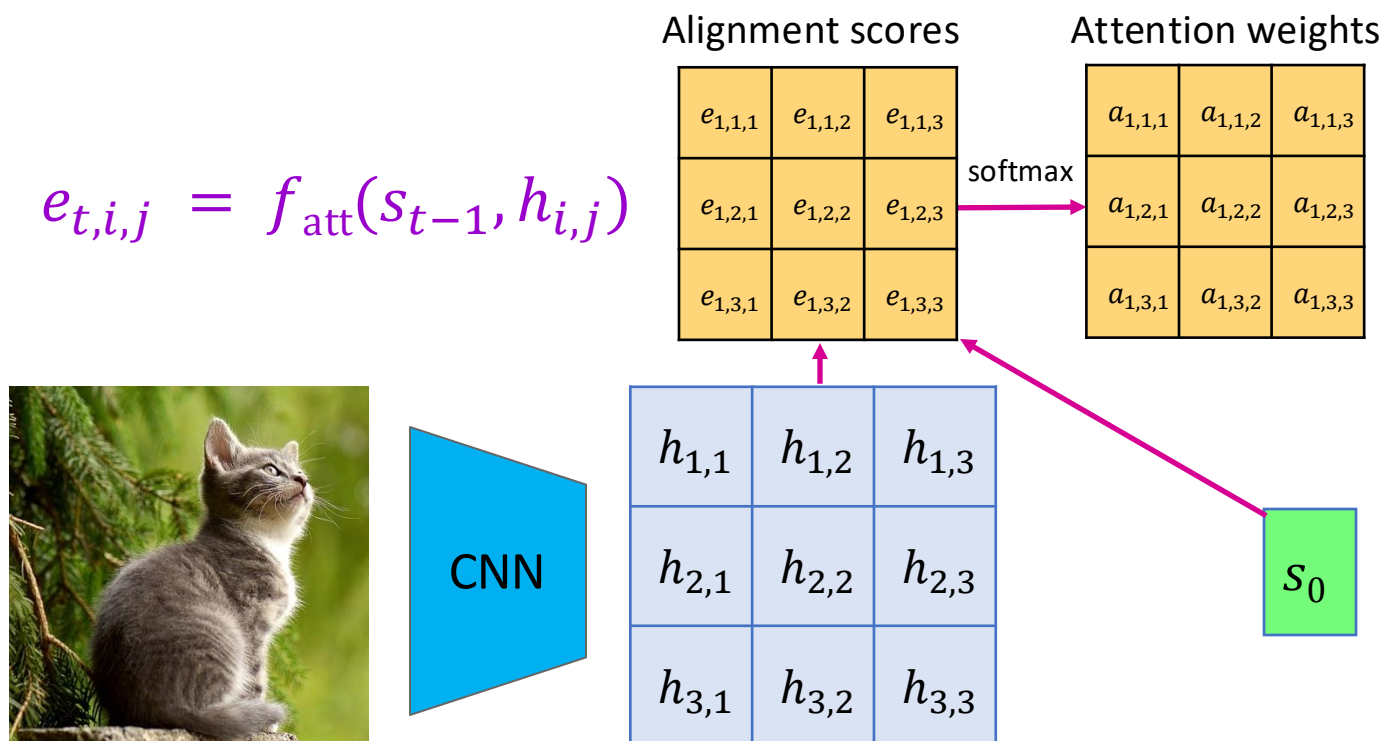


Image captioning with RNNs and attention



Use CNN to extract a grid of features

Image captioning with RNNs and attention

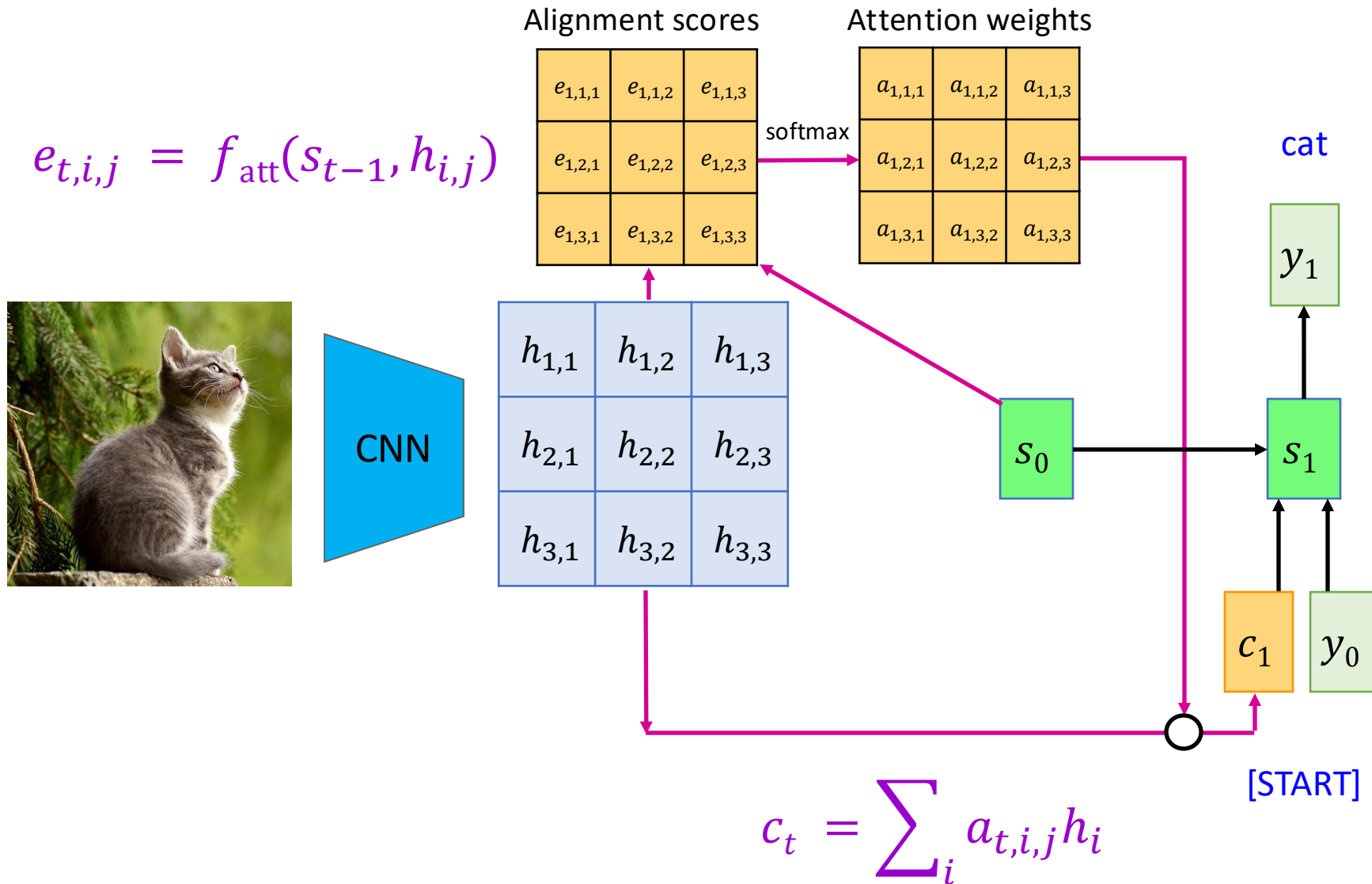


Image captioning with RNNs and attention

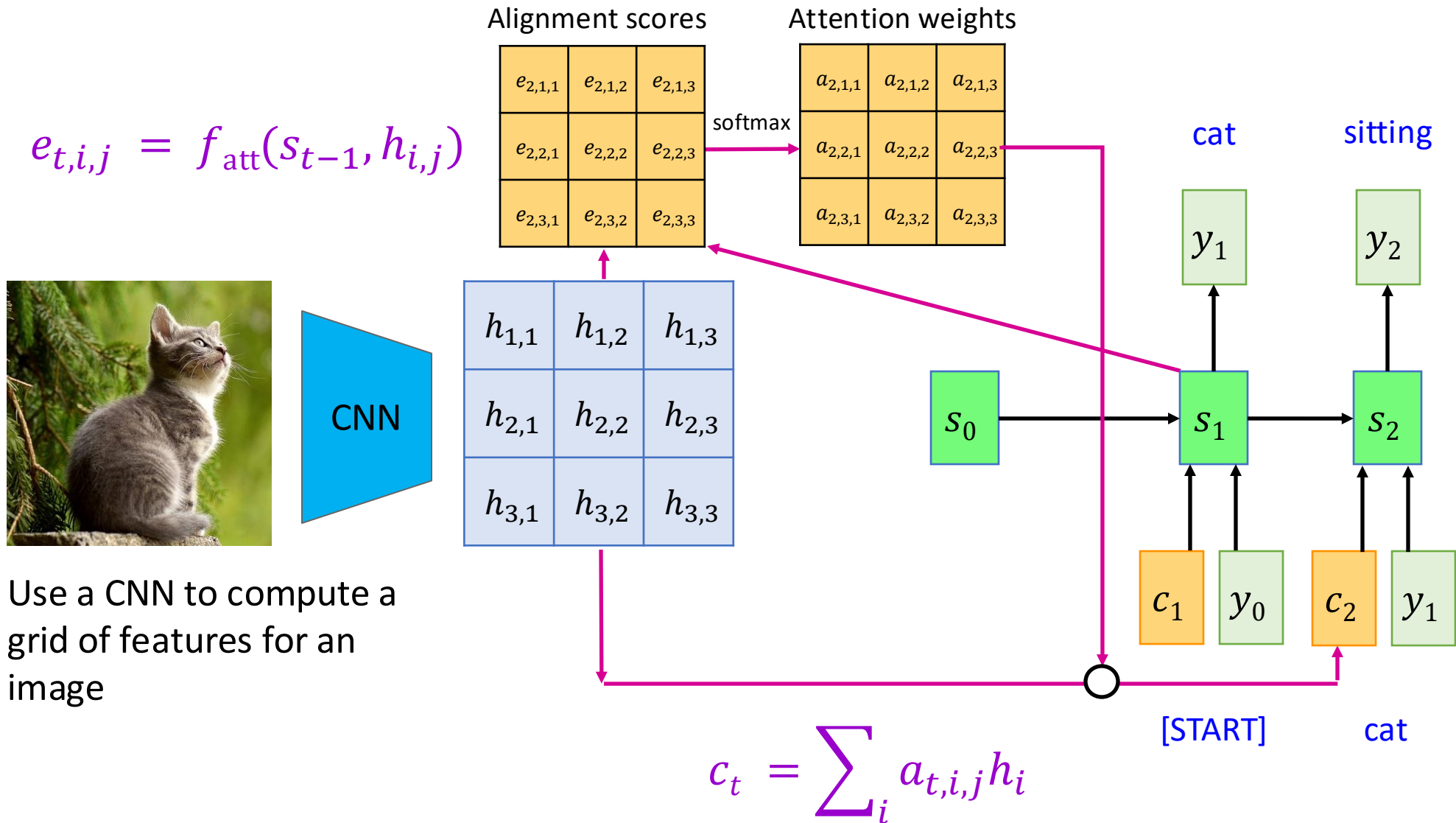
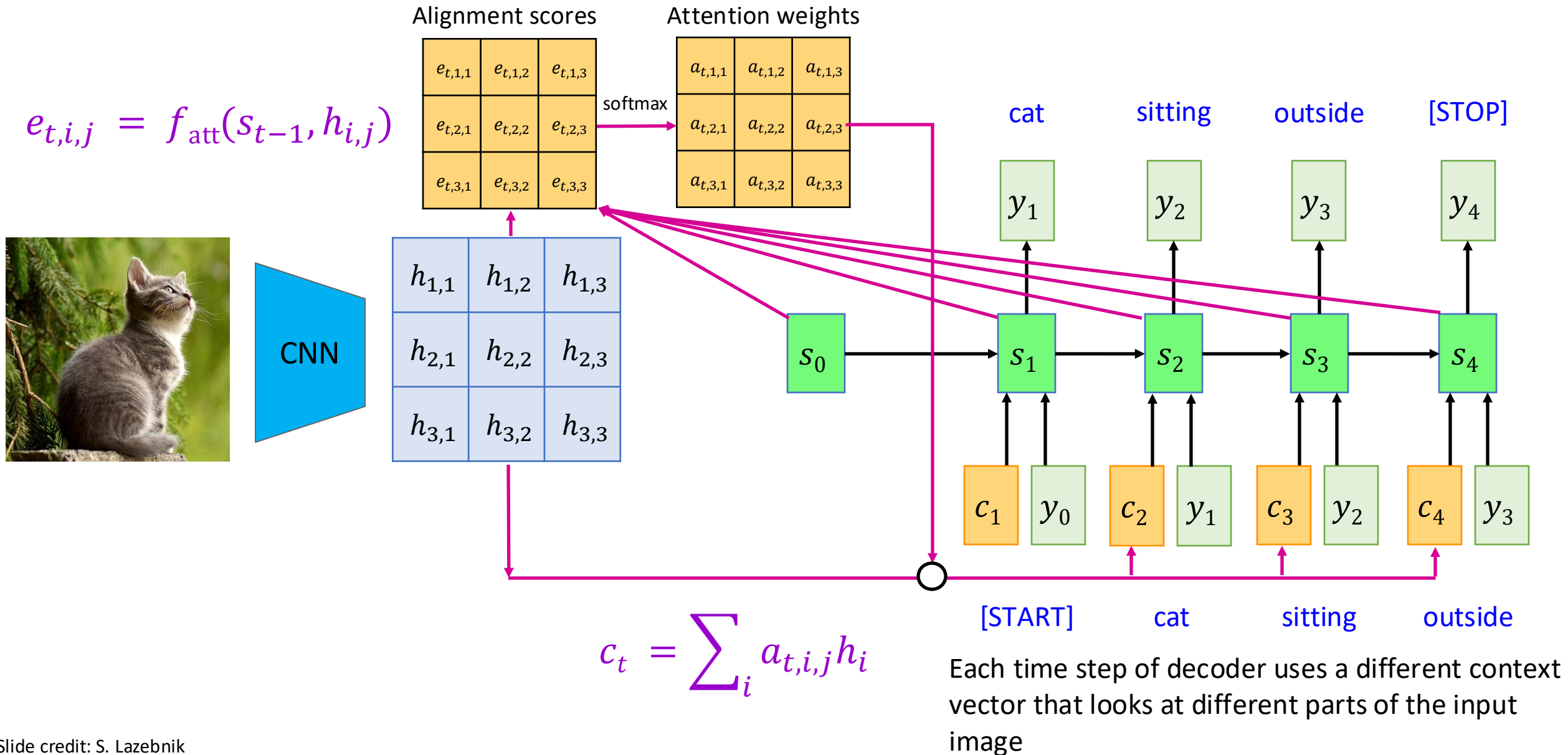


Image captioning with RNNs and attention



Example results

- Good captions



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



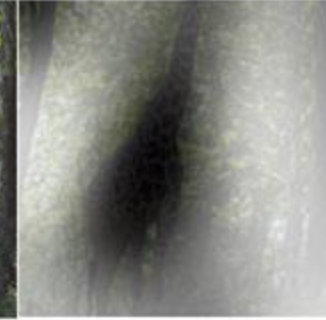
A little girl sitting on a bed with a teddy bear.



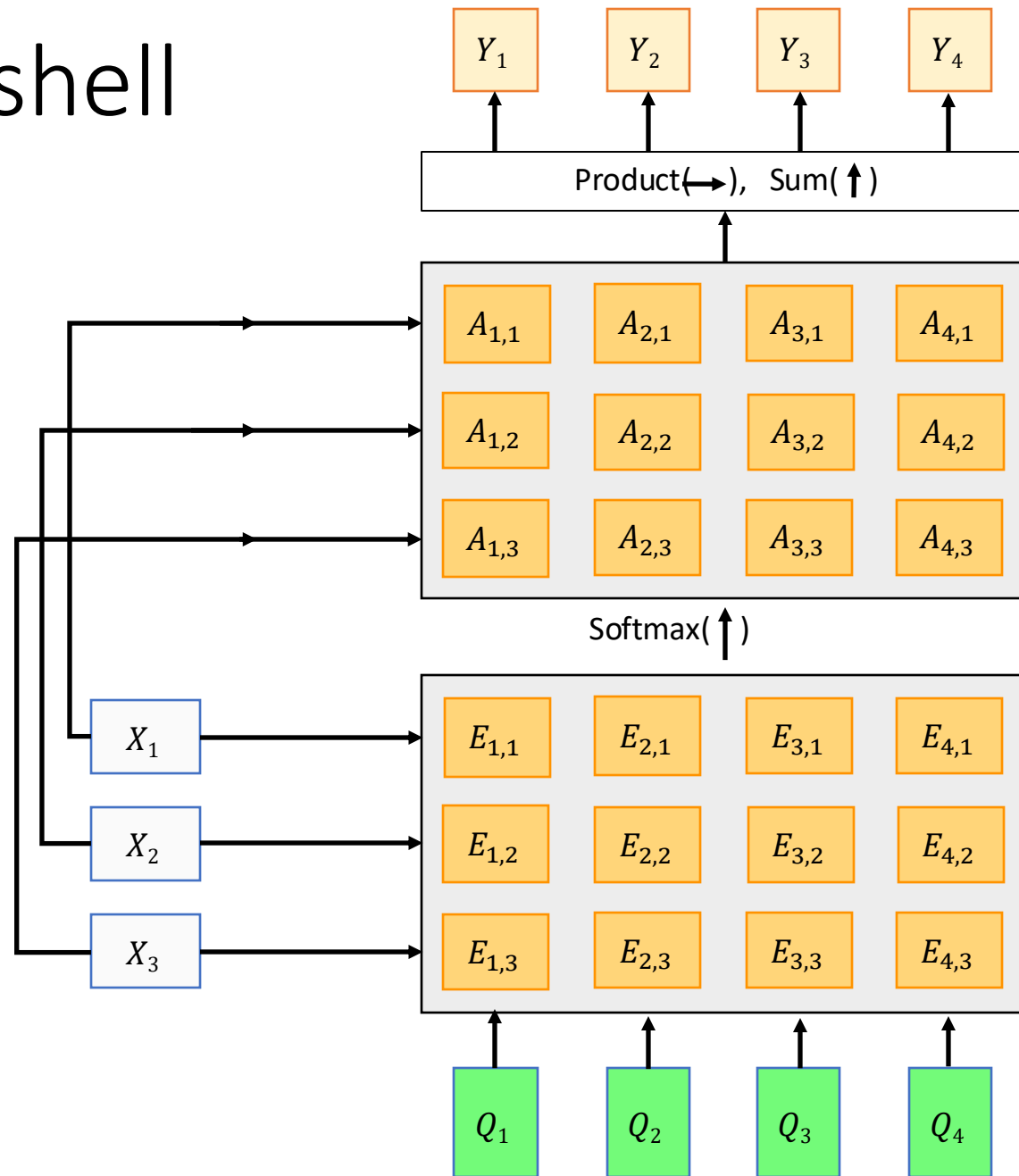
A group of people sitting on a boat in the water.



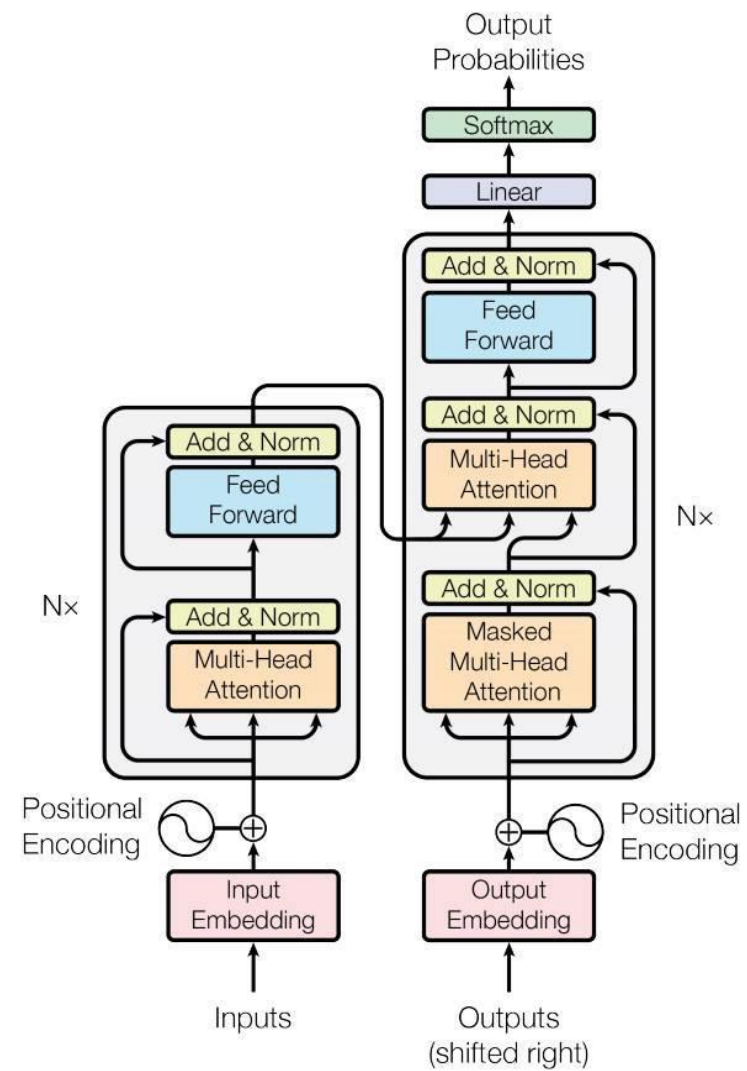
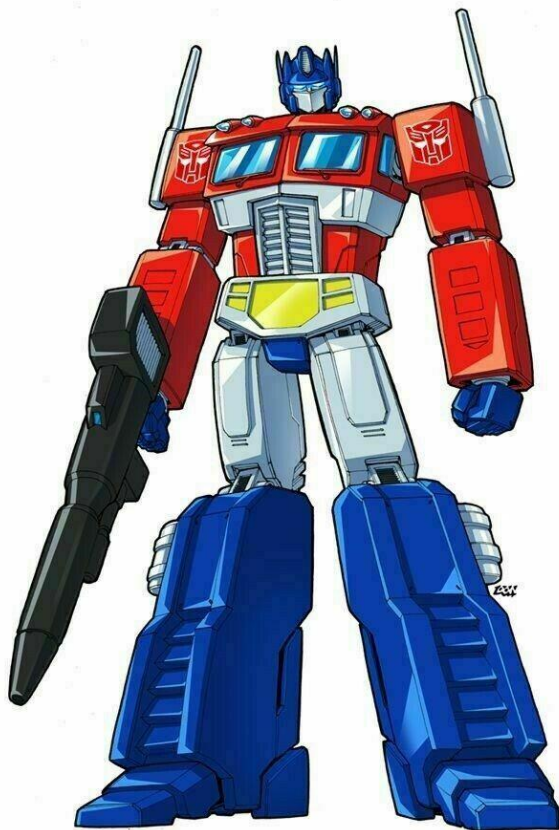
A giraffe standing in a forest with trees in the background.



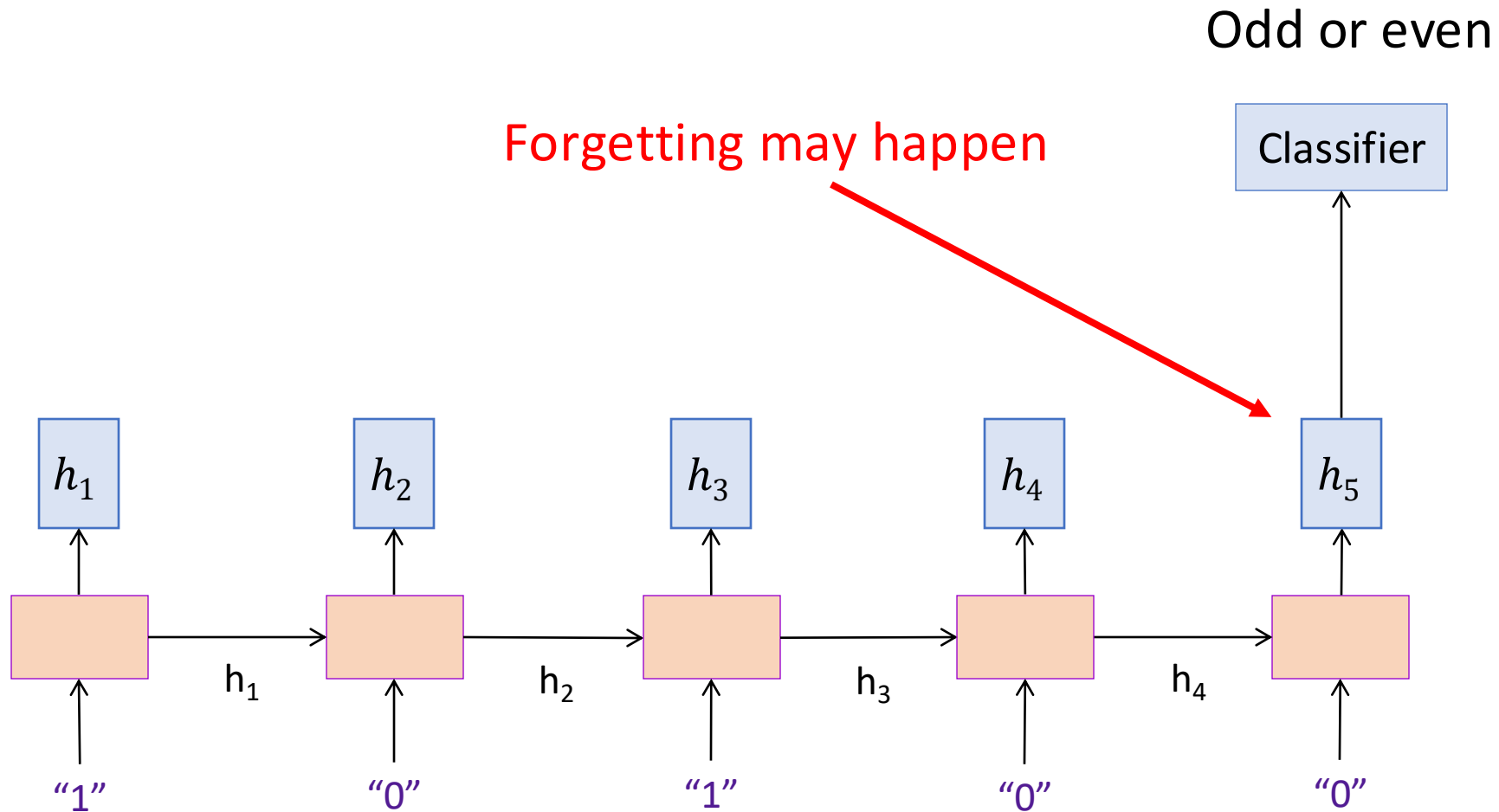
Attention in a nutshell



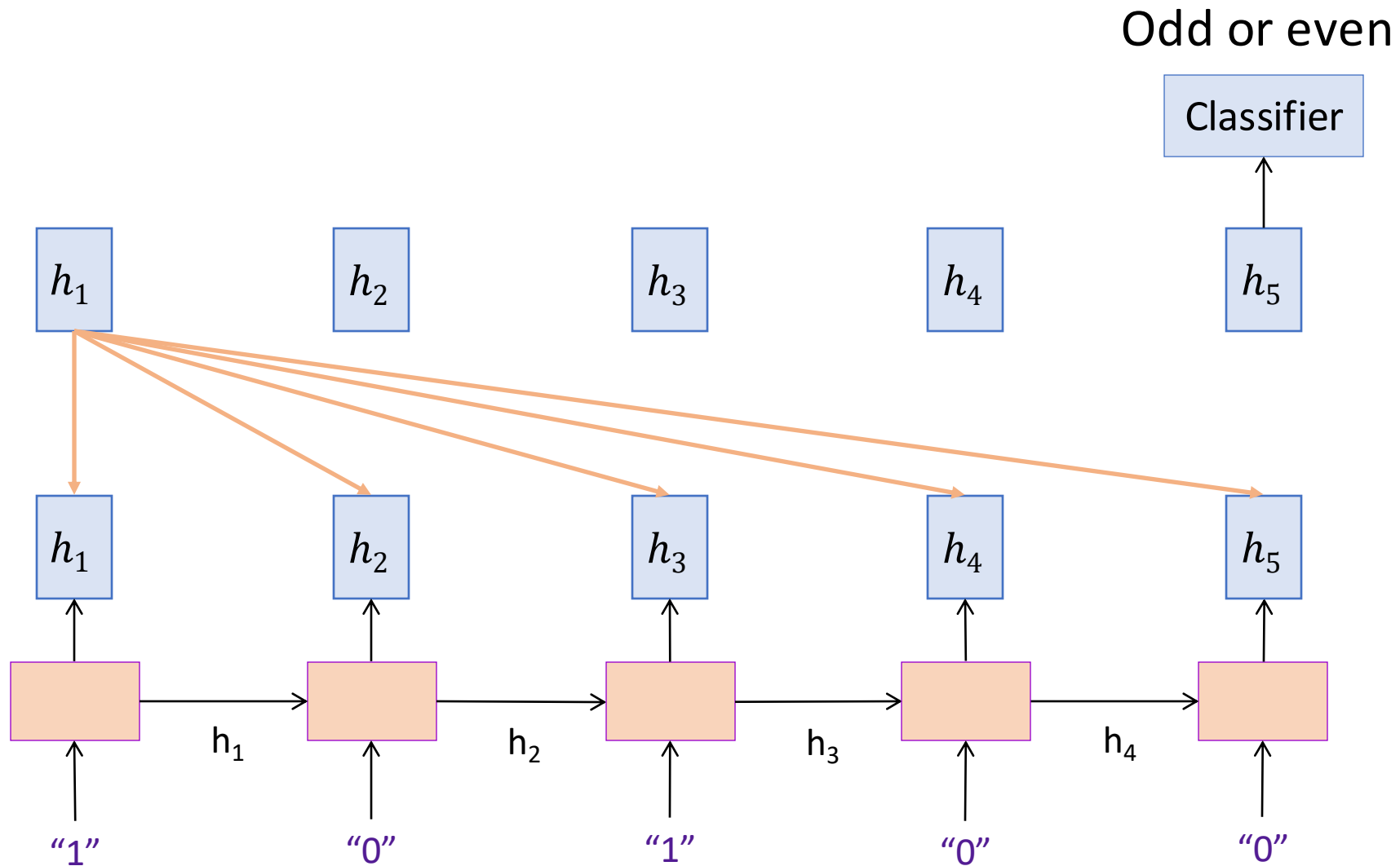
Transformer



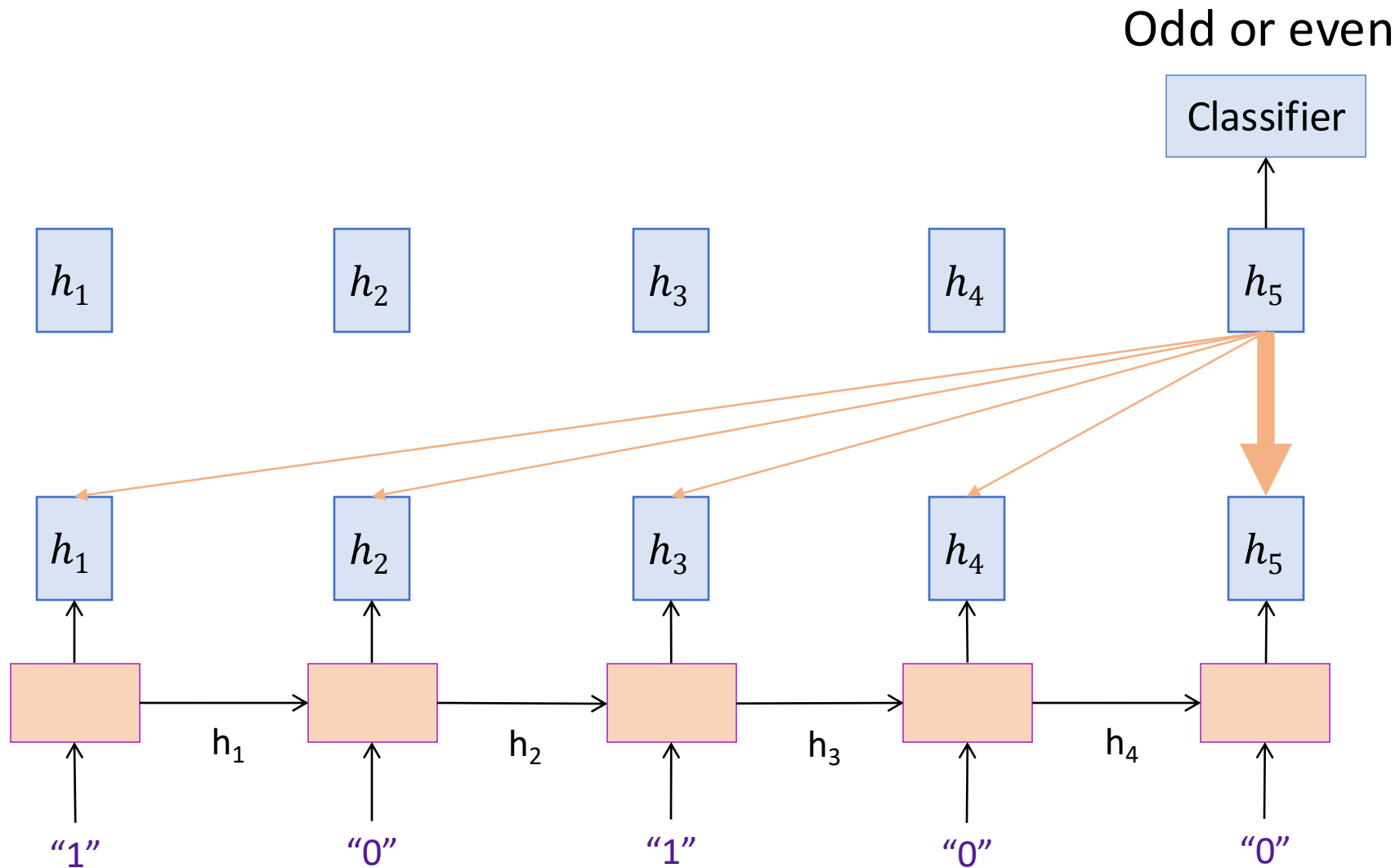
Intuition of Transformer: Self-attention



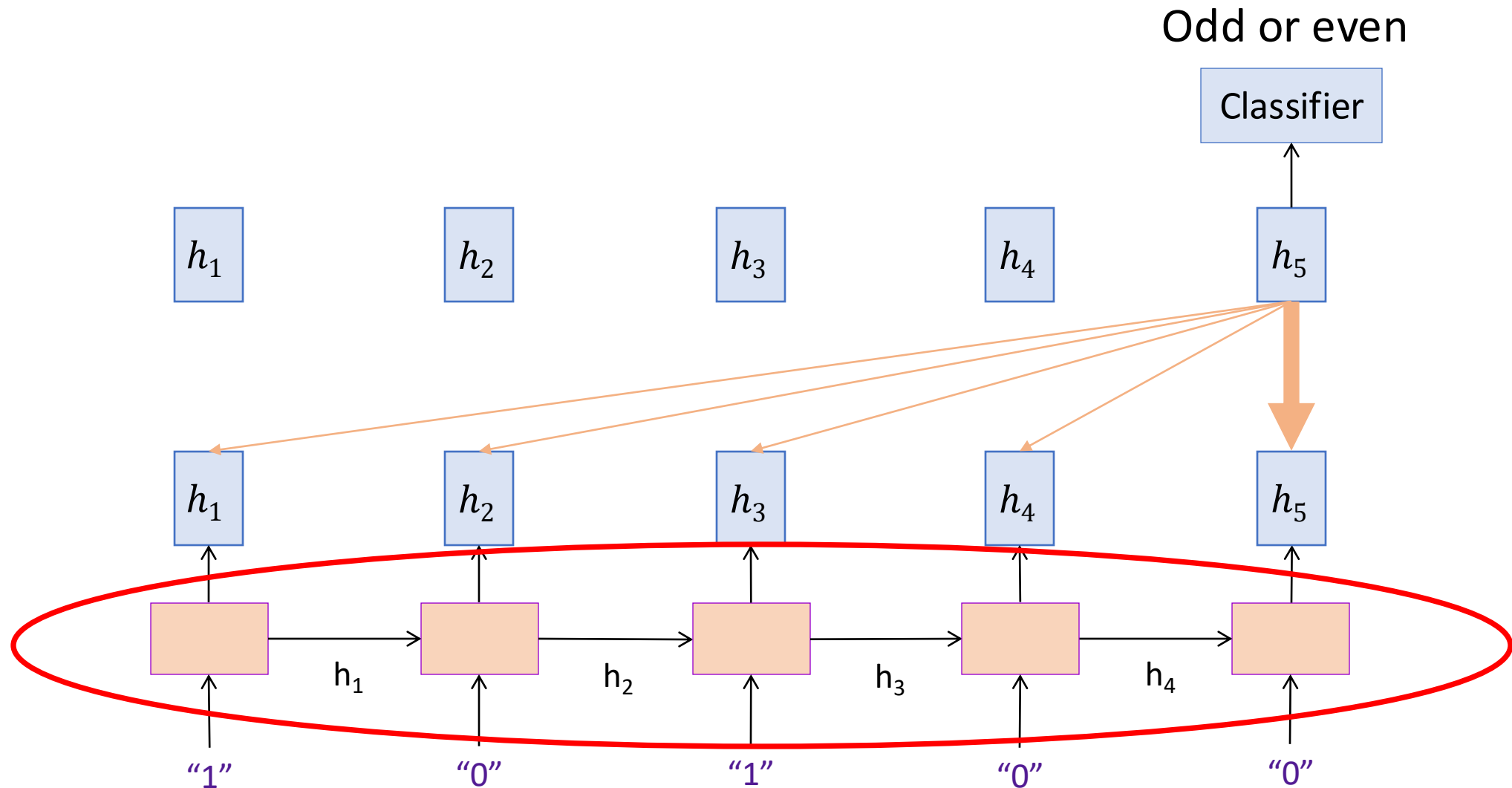
Intuition of Transformer: Self-attention



Intuition of Transformer: Self-attention

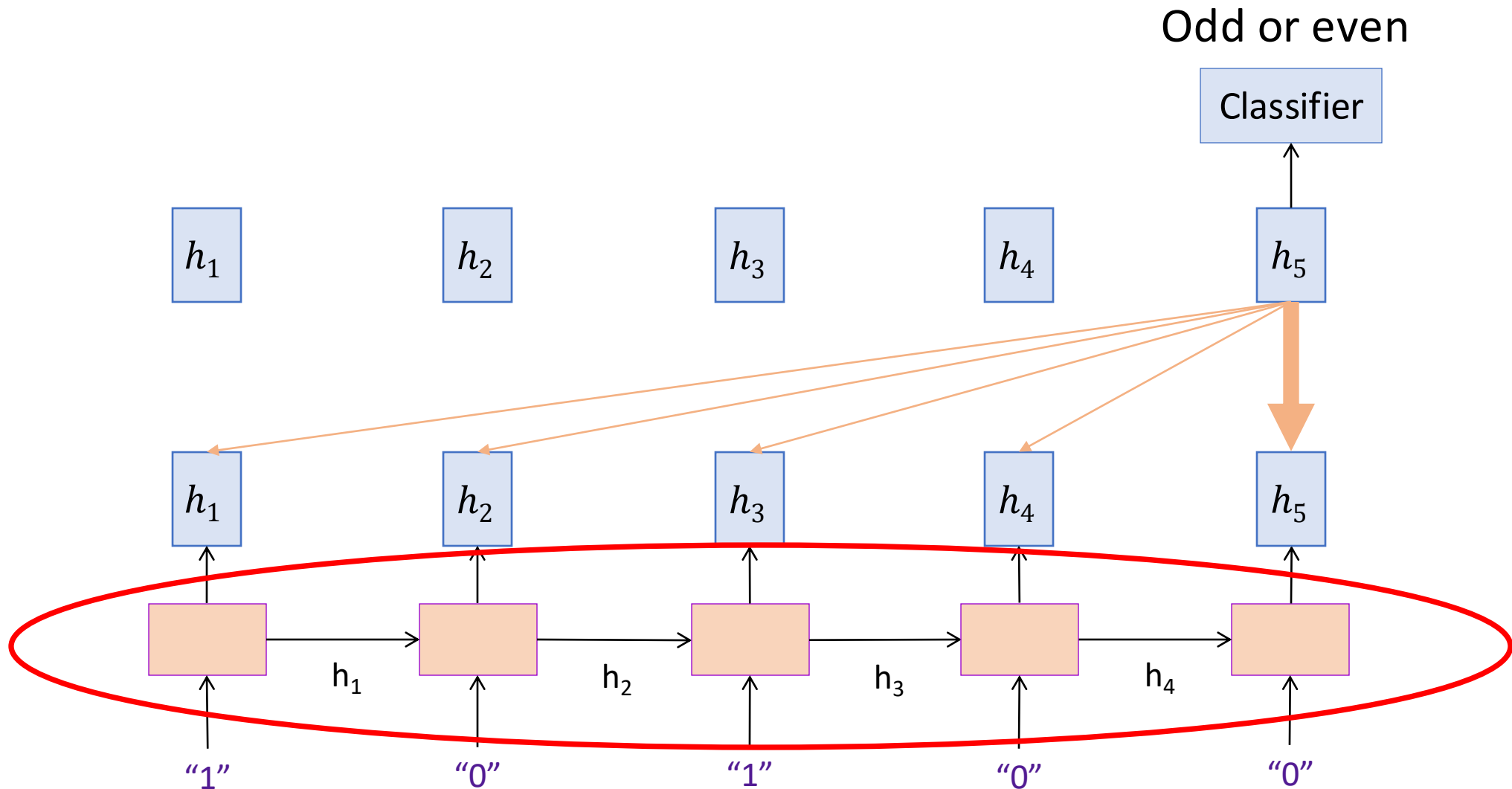


Intuition of Transformer: Self-attention



Do we still need the Recurrent Neural Network to encode the input?

Intuition of Transformer: Self-attention



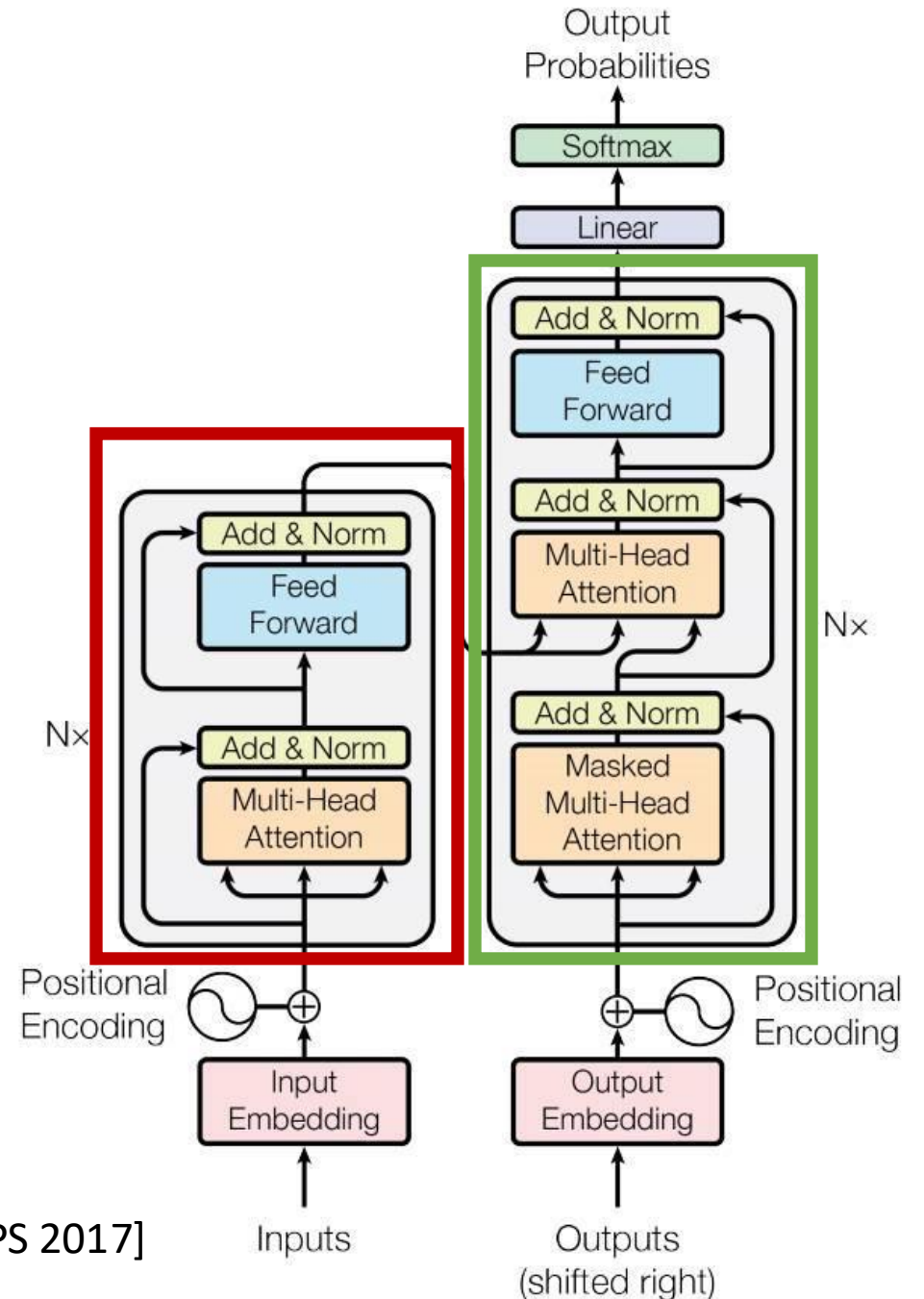
Do we still need the Recurrent Neural Network to encode the input?

Today's Class

Transformer Encoder

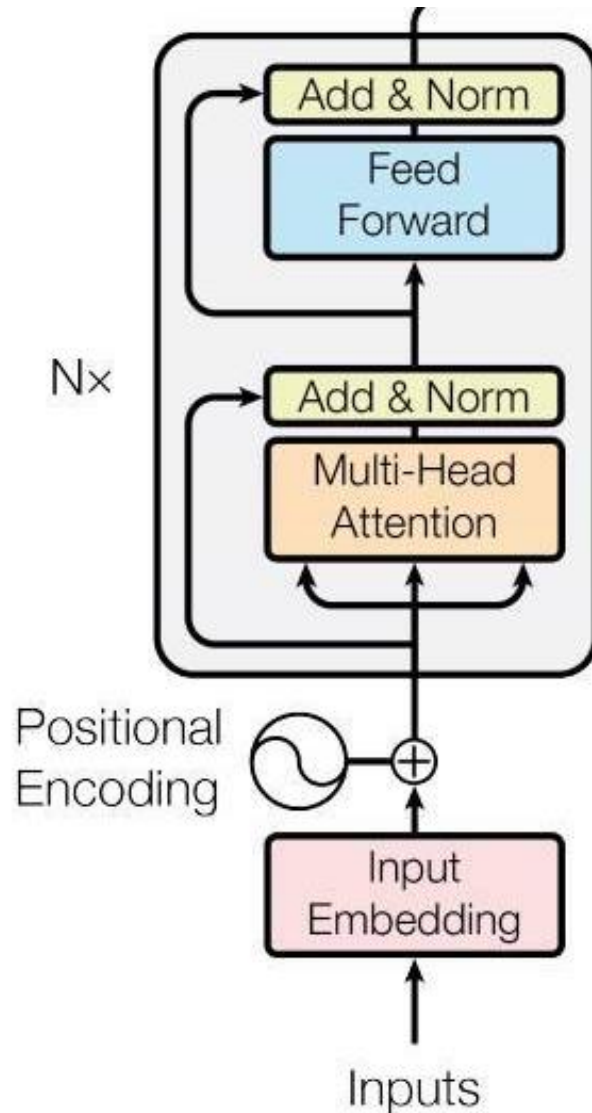
Overview of Transformer

- **Encoder** (left part)
 - multi-head self attention
 - Position embedding
 - Feedforward network
 - Non-linearity and normalization in-between
- **Decoder** (right part), **optional**
 - Multi-head cross attention
 - And others in the encoder



[Vaswani et al., [Attention is all you need](#). NeurIPS 2017]

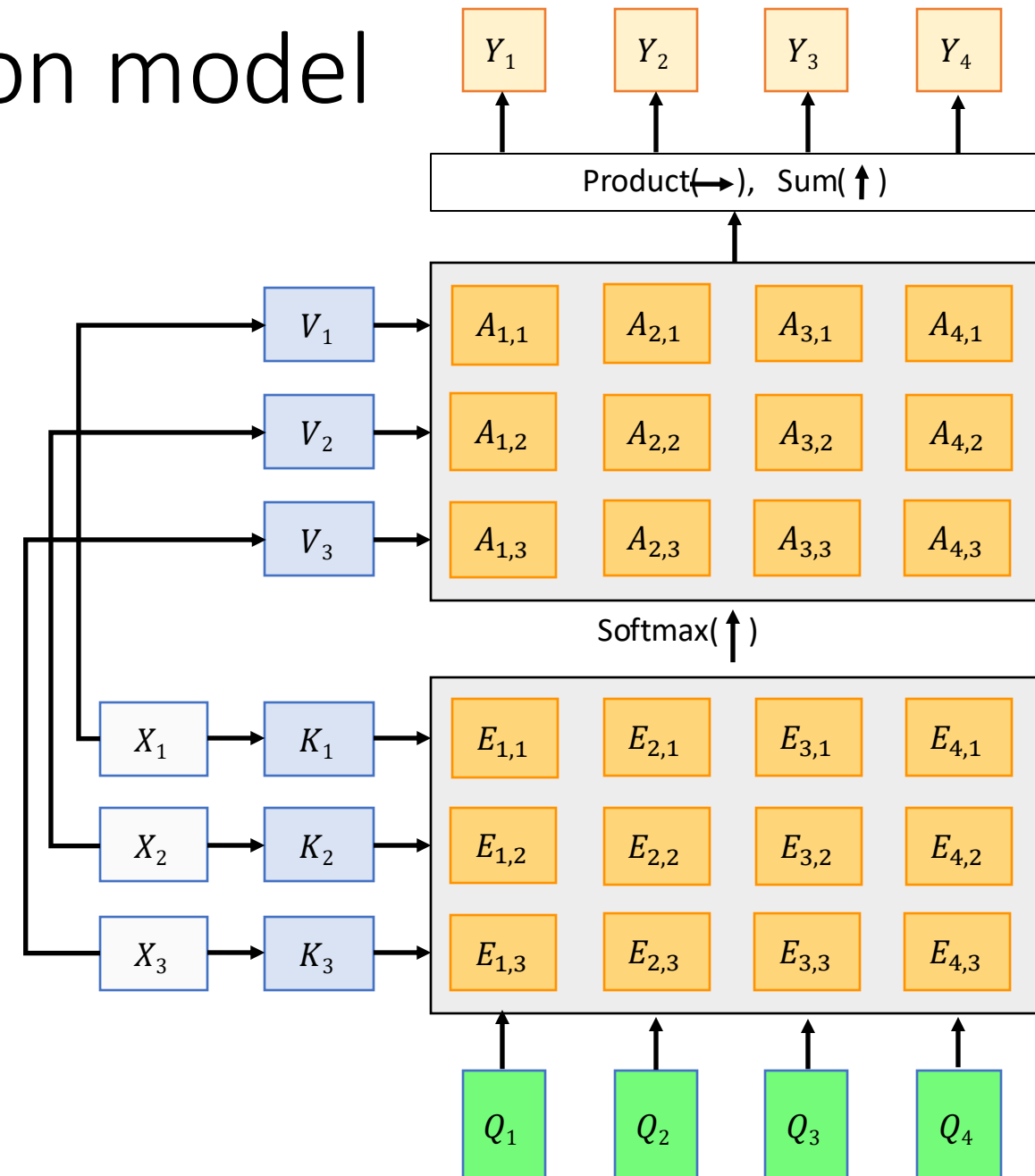
Overview of the Transformer Encoder (Cell)



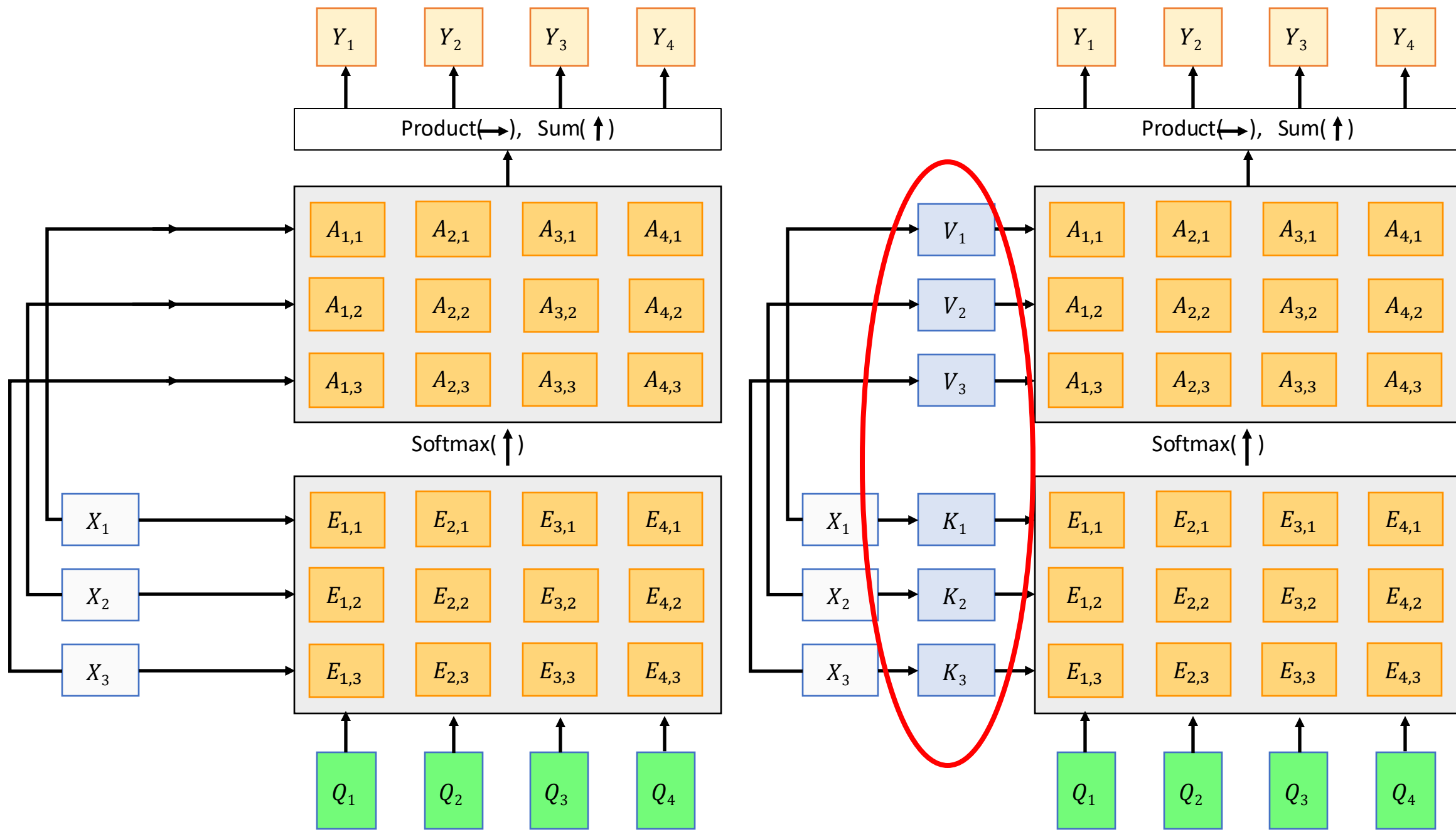
- Multi-head (self-)attention
- Positional encoding
- Feedforward network
- Residual connections
- Regularization tricks

Key-Value-Query attention model

- Key vectors: $K = XW_K$
- Value Vectors: $V = XW_V$
- Query vectors
- Similarities: *scaled dot-product attention*
 - $E_{i,j} = \frac{(Q_i \cdot K_j)}{\sqrt{D}}$ or $E = QK^T / \sqrt{D}$
 - (D is the dimensionality of the keys)
- Attn. weights: $A = \text{softmax}(E, \text{dim} = 1)$
- Output vectors:
 - $Y_i = \sum_j A_{i,j} V_j$ or $Y = AV$



Different Attention Mechanisms



Self-attention

- Used to capture context *within the sequence*

The animal didn't cross the street because **it** was too tired .

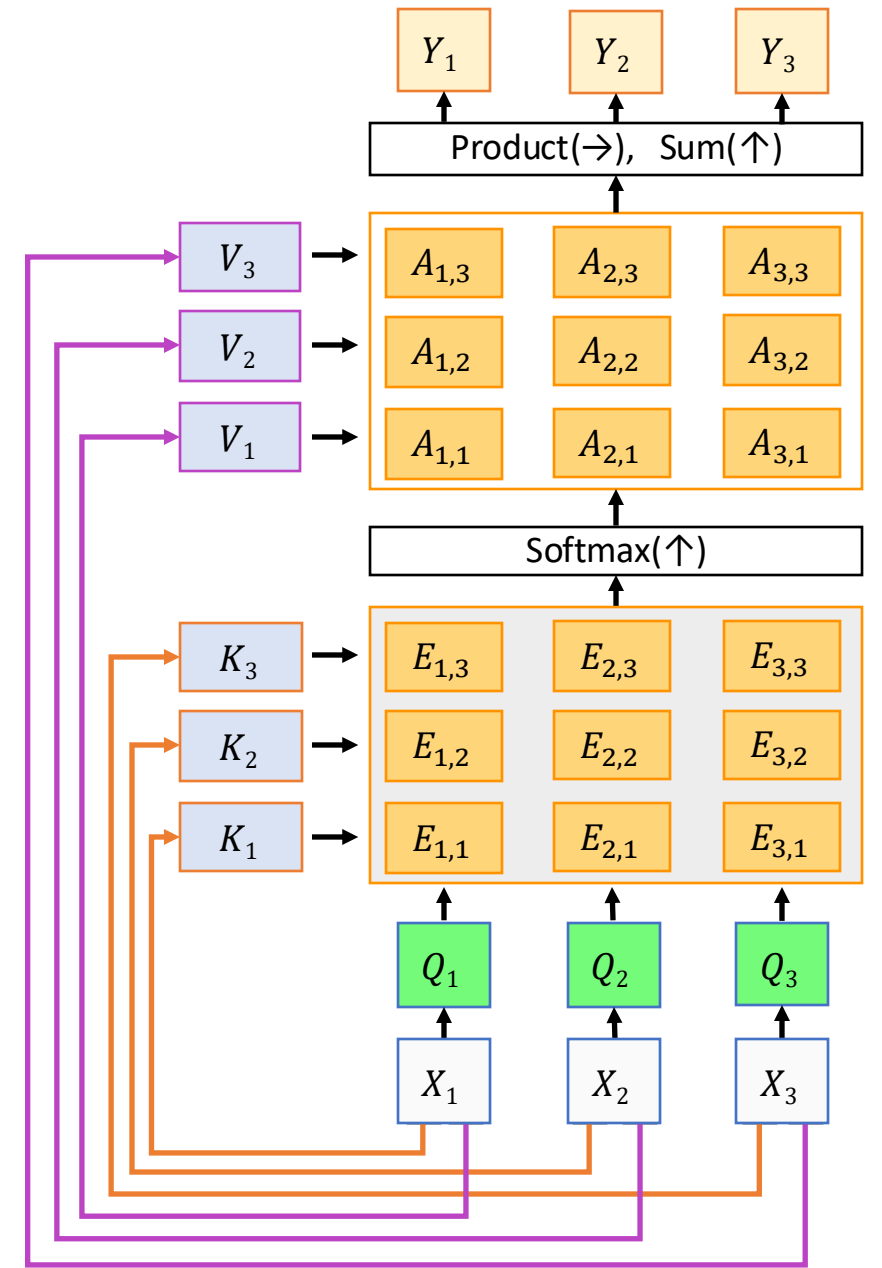
As we are encoding “it”, we should focus on “the animal”

The animal didn't cross the street because **it** was too wide .

As we are encoding “it”, we should focus on “the street”

Self-attention layer

- Query vectors: $Q = XW_Q$
- Key vectors: $K = XW_K$
- Value vectors: $V = XW_V$
- Similarities: *scaled dot-product attention*
 - $E_{i,j} = \frac{(Q_i \cdot K_j)}{\sqrt{D}}$ or $E = QK^T / \sqrt{D}$
 - (D is the dimensionality of the keys)
- Attn. weights: $A = \text{softmax}(E, \text{dim} = 1)$
- Output vectors:
 - $Y_i = \sum_j A_{i,j} V_j$ or $Y = AV$



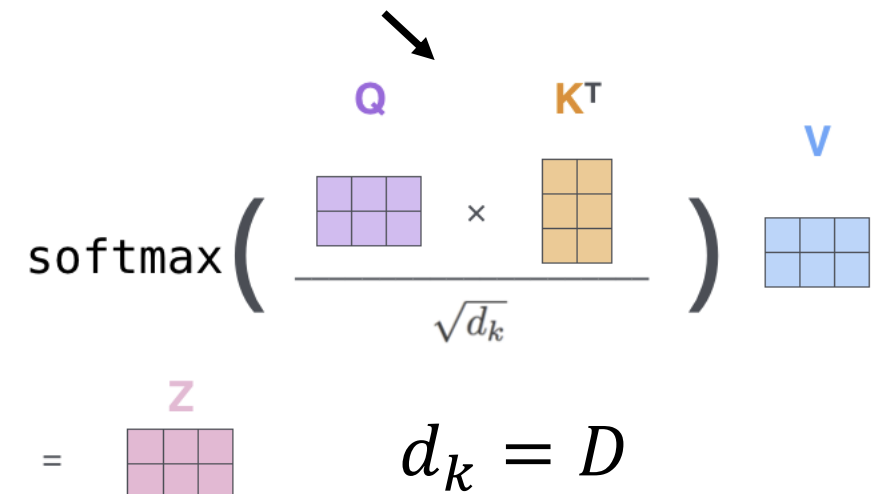
One query per input vector

Matrix Calculation of Self-Attention

$$X \in R^{L \times C}, W^Q \in R^{C \times D}, Q \in R^{L \times D}$$



Dimension: $R^{L \times L}$



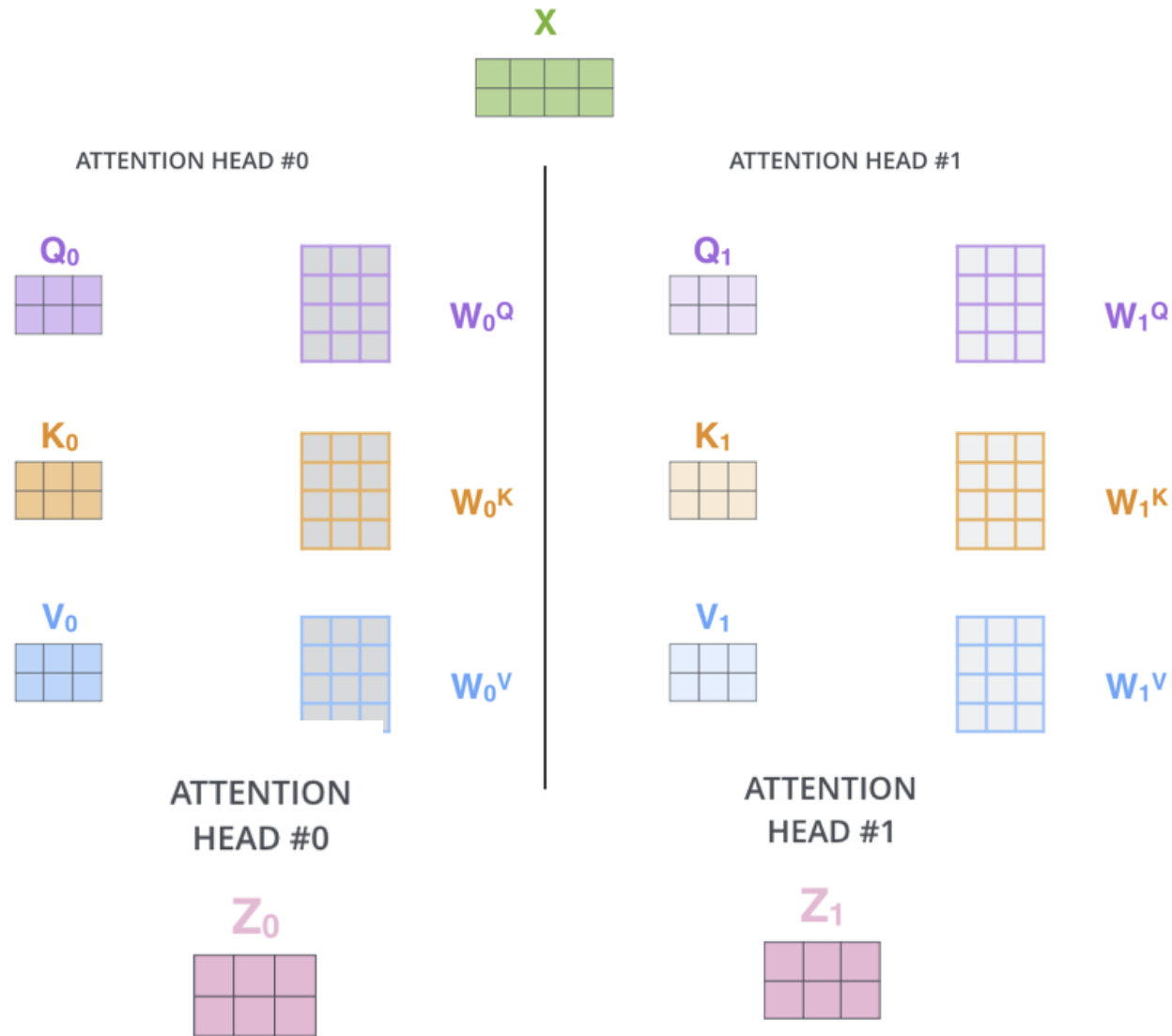
$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right) \times V$

$d_k = D$

Z

In the implementation, we need to process batched data, where $X \in R^{B \times L \times C}$, **check `torch.matmul` for batched matrix multiplication.**

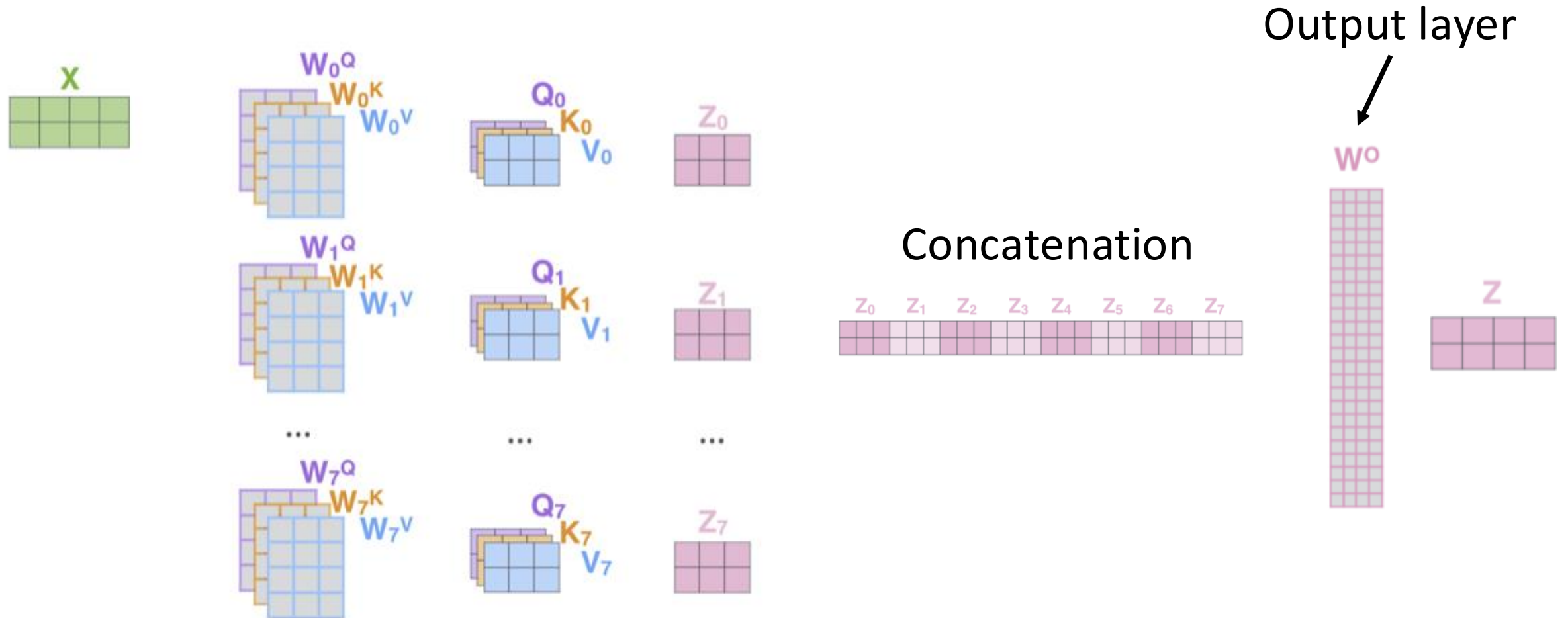
Multi-Head Self-Attention



Intuition:

Multiple attention heads can capture more information in the input

Multi-Head Self-Attention



Can we finish the entire operation without any for loops?

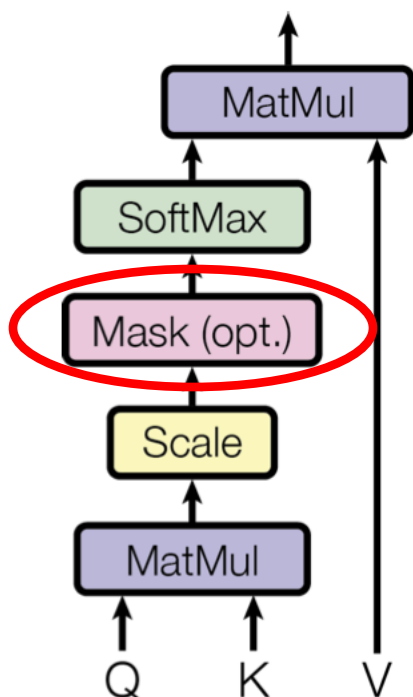
Finish Multi-head Attention without For Loops in PA3

- Linear transformations of Query, Key, and Value
 - Use a big (concatenated) transformation matrix for each of them, respectively
- How about the attention?
 - What are the shapes of Q, K, and V?
 - With and without the multi-head attention
 - The matrix multiplication is defined over a single attention
 - Solution:
 - Reshape the tensors so that the channel dimension is only about a single head
 - Merge the number of heads to the batch dimension

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{matrix} \text{3x3} \end{matrix} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{matrix} \text{3x3} \end{matrix} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{matrix} \text{3x3} \end{matrix} \end{matrix}$$

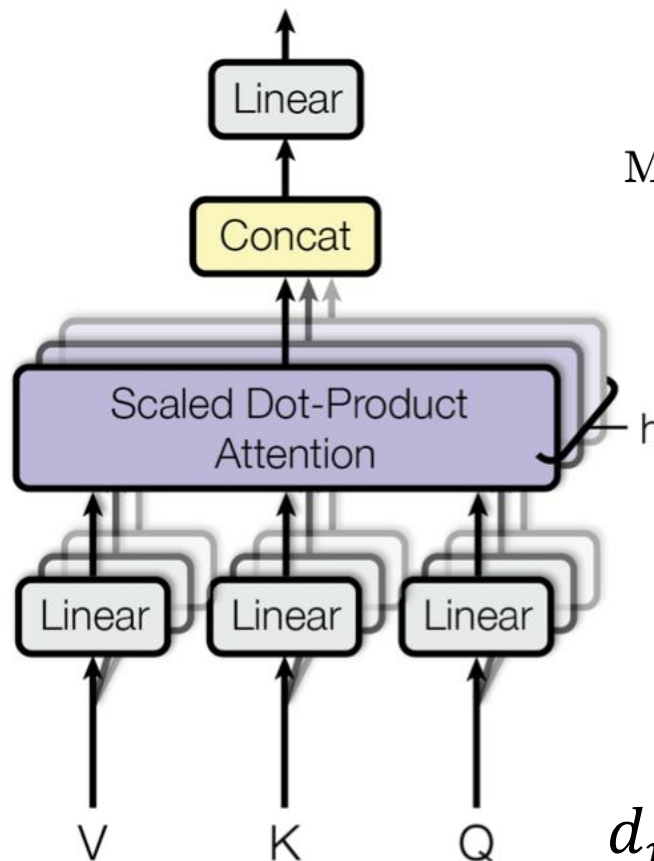
Multi-Head Self-Attention Summary

Scaled Dot-Product Attention



Will be clear when we
talk about the decoder

Multi-Head Attention



$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}, W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$$

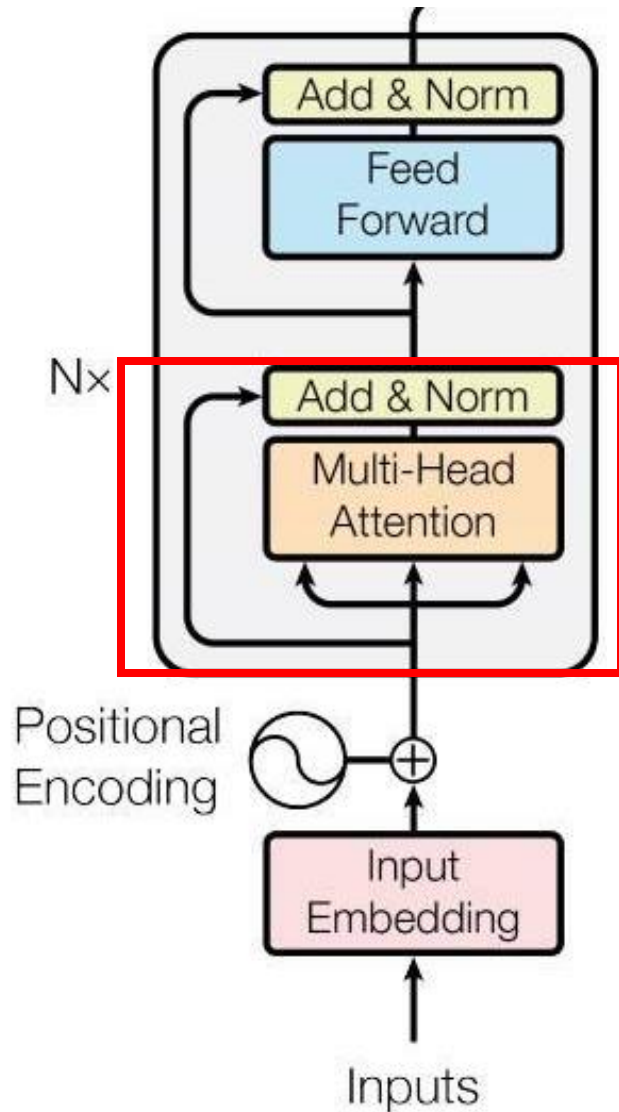
$$W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$$

$$W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$$

$$d_k = d_v = d_{\text{model}}/h$$

d_{model} : dimension of input embeddings
 h : number of attention heads

Residual Connections, Dropout, and Normalization



$$\begin{aligned} X_{att} &= \text{MHA}(X) \\ X_{att} &= \text{Dropout}(X_{att}, p = 0.1) \\ X &= \text{LayerNorm}(X + X_{att}) \end{aligned}$$

Recall: Layer Normalization

1. Variable lengths in the input.
2. The recurrent nature of RNN.

Computing of μ, σ is independent of the batch and length dimension.

But the learnable parameters γ, β are still for each channel.

Largely adopted in Transformer.

Layer Normalization for **recurrent** networks

$$x : N \times L \times C$$



$$\mu, \sigma : N \times L \times 1$$

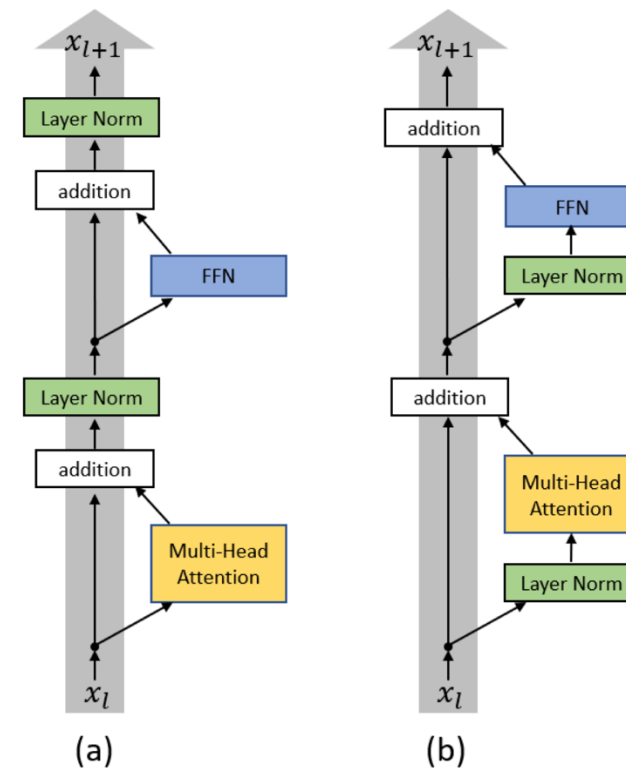
$$\gamma, \beta : 1 \times 1 \times C$$

$$y = \frac{(x - \mu)}{\sigma} \gamma + \beta$$

On the LayerNorm in Transformer

<https://arxiv.org/pdf/2002.04745>

Implement Post-Norm in PA3



Post-LN Transformer

$$\begin{aligned}
 x_{l,i}^{post,1} &= \text{MultiHeadAtt}(x_{l,i}^{post}, [x_{l,1}^{post}, \dots, x_{l,n}^{post}]) \\
 x_{l,i}^{post,2} &= x_{l,i}^{post} + x_{l,i}^{post,1} \\
 x_{l,i}^{post,3} &= \text{LayerNorm}(x_{l,i}^{post,2}) \\
 x_{l,i}^{post,4} &= \text{ReLU}(x_{l,i}^{post,3} W^{1,l} + b^{1,l}) W^{2,l} + b^{2,l} \\
 x_{l,i}^{post,5} &= x_{l,i}^{post,3} + x_{l,i}^{post,4} \\
 x_{l+1,i}^{post} &= \text{LayerNorm}(x_{l,i}^{post,5})
 \end{aligned}$$

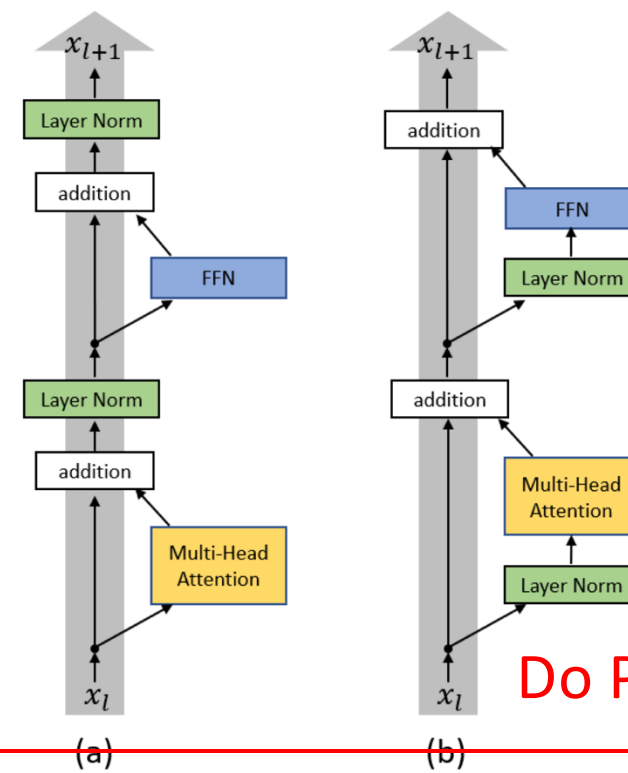
Pre-LN Transformer

$$\begin{aligned}
 x_{l,i}^{pre,1} &= \text{LayerNorm}(x_{l,i}^{pre}) \\
 x_{l,i}^{pre,2} &= \text{MultiHeadAtt}(x_{l,i}^{pre,1}, [x_{l,1}^{pre,1}, \dots, x_{l,n}^{pre,1}]) \\
 x_{l,i}^{pre,3} &= x_{l,i}^{pre} + x_{l,i}^{pre,2} \\
 x_{l,i}^{pre,4} &= \text{LayerNorm}(x_{l,i}^{pre,3}) \\
 x_{l,i}^{pre,5} &= \text{ReLU}(x_{l,i}^{pre,4} W^{1,l} + b^{1,l}) W^{2,l} + b^{2,l} \\
 x_{l+1,i}^{pre} &= x_{l,i}^{pre,5} + x_{l,i}^{pre,3}
 \end{aligned}$$

$$\text{Final LayerNorm: } x_{Final,i}^{pre} \leftarrow \text{LayerNorm}(x_{L+1,i}^{pre})$$

On the LayerNorm in Transformer

<https://arxiv.org/pdf/2002.04745>



Do Pre-Norm in practice

Post-LN Transformer

$$\begin{aligned}
 x_{l,i}^{post,1} &= \text{MultiHeadAtt}(x_{l,i}^{post}, [x_{l,1}^{post}, \dots, x_{l,n}^{post}]) \\
 x_{l,i}^{post,2} &= x_{l,i}^{post} + x_{l,i}^{post,1} \\
 x_{l,i}^{post,3} &= \text{LayerNorm}(x_{l,i}^{post,2}) \\
 x_{l,i}^{post,4} &= \text{ReLU}(x_{l,i}^{post,3} W^{1,l} + b^{1,l}) W^{2,l} + b^{2,l} \\
 x_{l,i}^{post,5} &= x_{l,i}^{post,3} + x_{l,i}^{post,4} \\
 x_{l+1,i}^{post} &= \text{LayerNorm}(x_{l,i}^{post,5})
 \end{aligned}$$

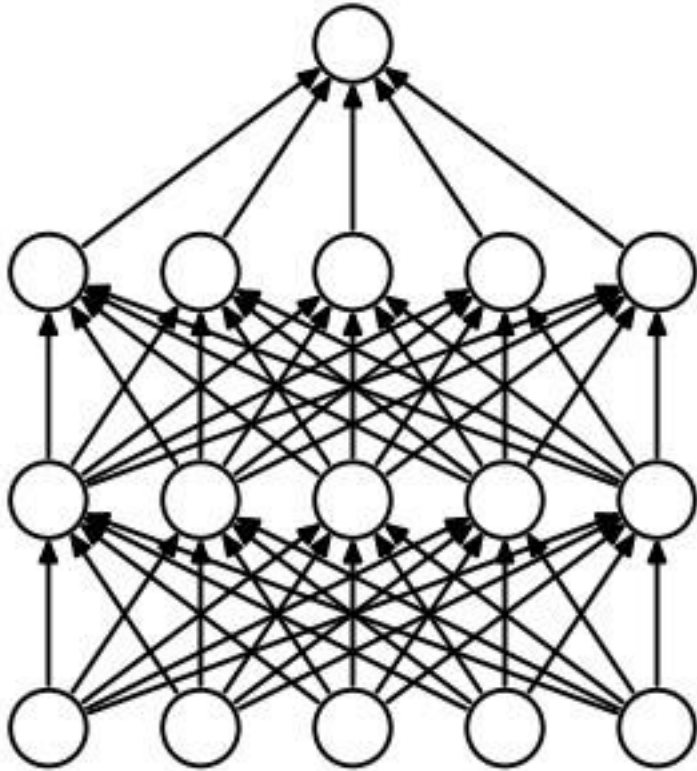
Pre-LN Transformer

$$\begin{aligned}
 x_{l,i}^{pre,1} &= \text{LayerNorm}(x_{l,i}^{pre}) \\
 x_{l,i}^{pre,2} &= \text{MultiHeadAtt}(x_{l,i}^{pre,1}, [x_{l,1}^{pre,1}, \dots, x_{l,n}^{pre,1}]) \\
 x_{l,i}^{pre,3} &= x_{l,i}^{pre} + x_{l,i}^{pre,2} \\
 x_{l,i}^{pre,4} &= \text{LayerNorm}(x_{l,i}^{pre,3}) \\
 x_{l,i}^{pre,5} &= \text{ReLU}(x_{l,i}^{pre,4} W^{1,l} + b^{1,l}) W^{2,l} + b^{2,l} \\
 x_{l+1,i}^{pre} &= x_{l,i}^{pre,5} + x_{l,i}^{pre,3}
 \end{aligned}$$

$$\text{Final LayerNorm: } x_{Final,i}^{pre} \leftarrow \text{LayerNorm}(x_{L+1,i}^{pre})$$

Regularization: **Dropout**

“randomly set some neurons to zero in the forward pass”

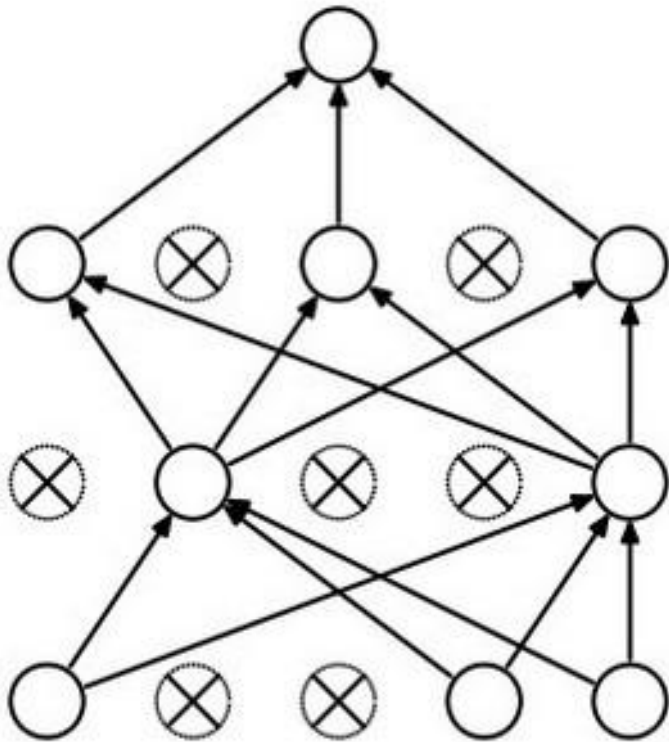


(a) Standard Neural Net

Randomly set the
output value of some
neurons to be 0

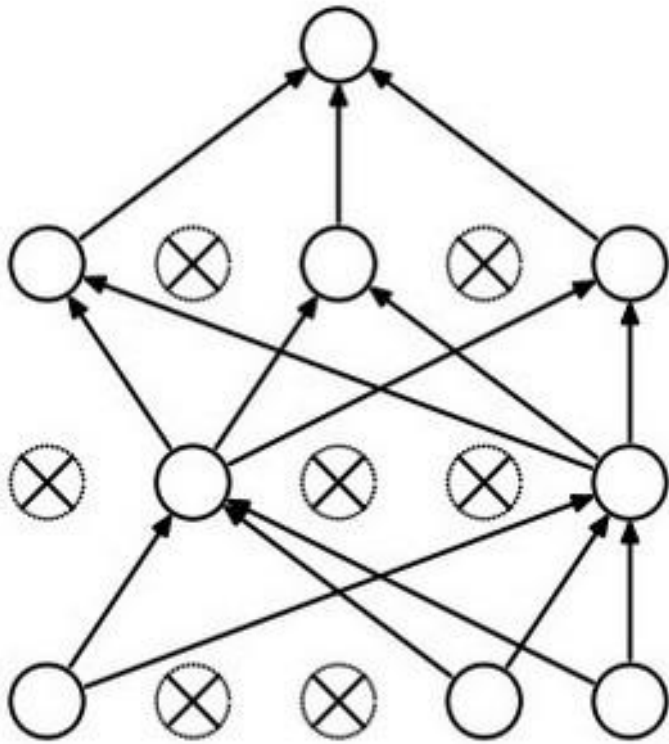
Waaaaait a second...

How could this possibly be a good idea?



Waaaaait a second...

How could this possibly be a good idea?



Forces the network to have a redundant representation.



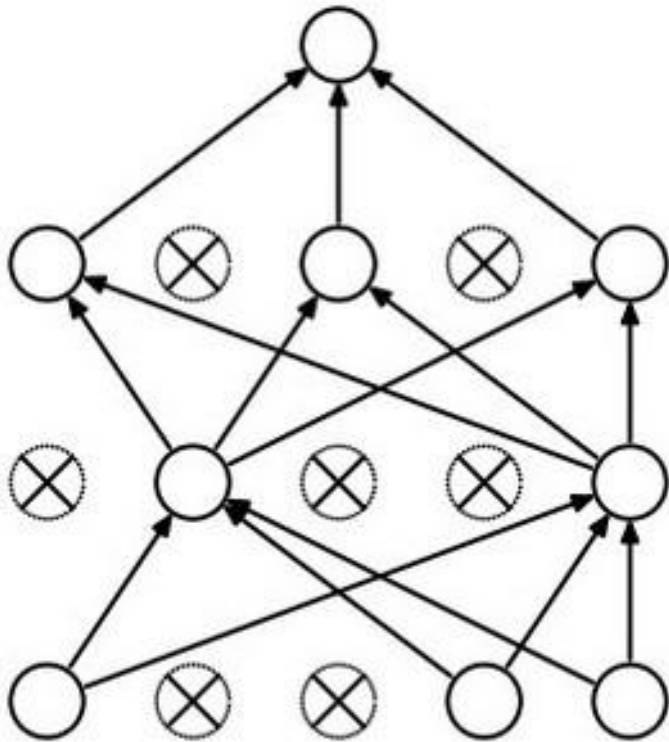
Training with occlusions?



Hiding because you
know it's "Vet" day !

Waaaaait a second...

How could this possibly be a good idea?



Another interpretation:

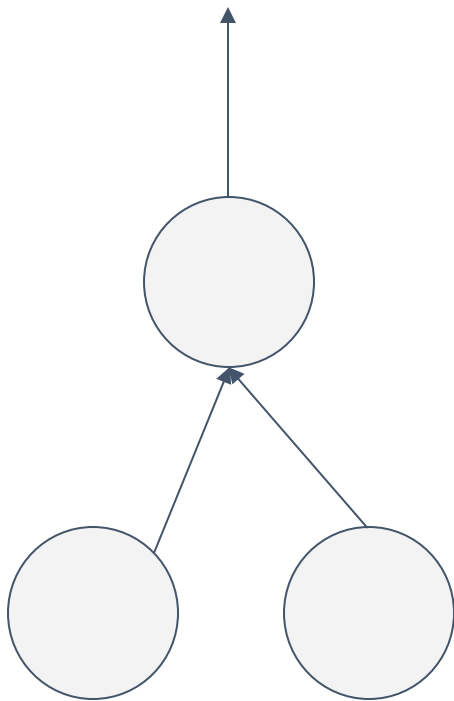
Dropout is training a large ensemble of models (that share parameters).

Each binary mask is one model, gets trained on only \sim one datapoint.

Dropout: Test Time

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).



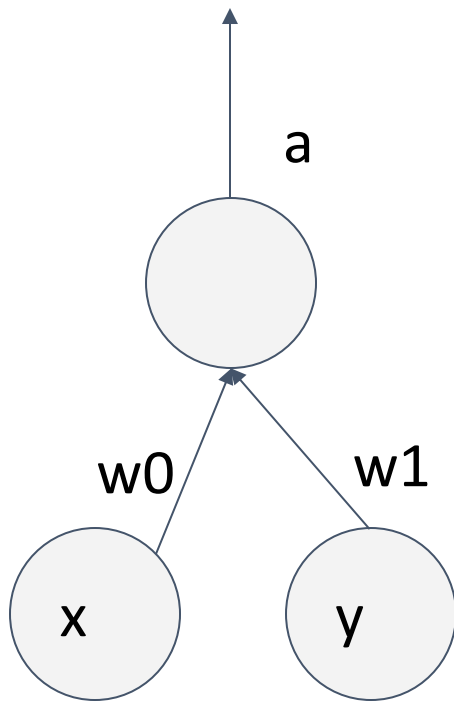
Q: Suppose that with all inputs present at test time the output of this neuron is a .

What would its output be during training time, in expectation? (e.g. if $p = 0.5$)

Dropout: Test Time

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).



during test: $a = w_0 * x + w_1 * y$

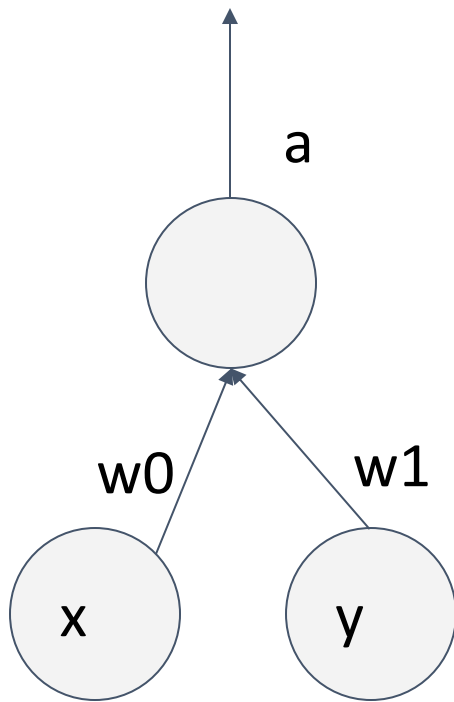
during train:

$$\begin{aligned} E[a] &= \frac{1}{4} * (w_0 * 0 + w_1 * 0 + \\ &\quad w_0 * 0 + w_1 * y + \\ &\quad w_0 * x + w_1 * 0 + \\ &\quad w_0 * x + w_1 * y) \\ &= \frac{1}{4} * (2 w_0 * x + 2 w_1 * y) \\ &= \frac{1}{2} * (w_0 * x + w_1 * y) \end{aligned}$$

At test time....

Can in fact do this with a single forward pass! (approximately)

Leave all input neurons turned on (no dropout).



during test: $a = w_0 * x + w_1 * y$

during train:

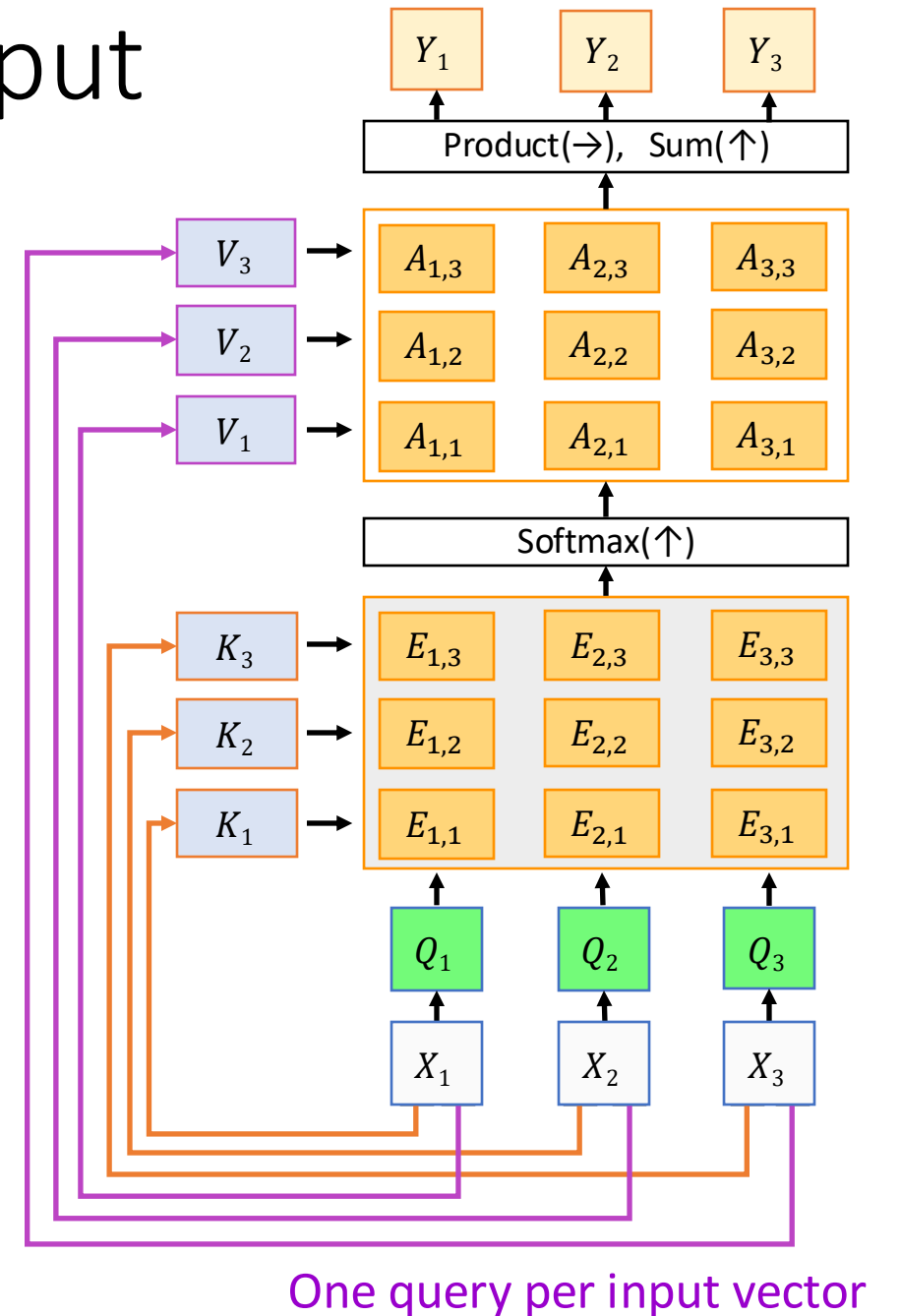
$$\begin{aligned} E[a] &= \frac{1}{4} * (w_0 * 0 + w_1 * 0 + \\ &\quad w_0 * 0 + w_1 * y + \\ &\quad w_0 * x + w_1 * 0 + \\ &\quad w_0 * x + w_1 * y) \\ &= \frac{1}{4} * (2 w_0 * x + 2 w_1 * y) \\ &= \frac{1}{2} * (w_0 * x + w_1 * y) \end{aligned}$$

With $p=0.5$, using all inputs in the forward pass would inflate the activations by 2x from what the network was “used to” during training!

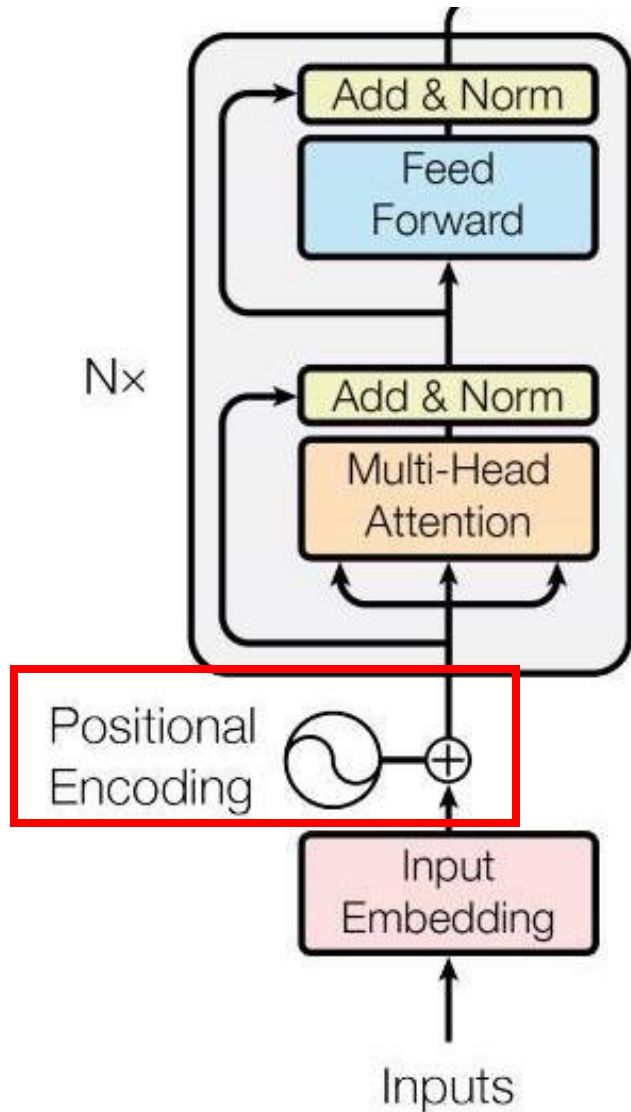
=> Have to compensate by scaling the activations back down by $\frac{1}{2}$

Capturing the Order of the Input

- Do the embeddings of each token/input change if we randomly permute the input?
- Is it good or bad?



Augmenting the MSA with Positional Encoding



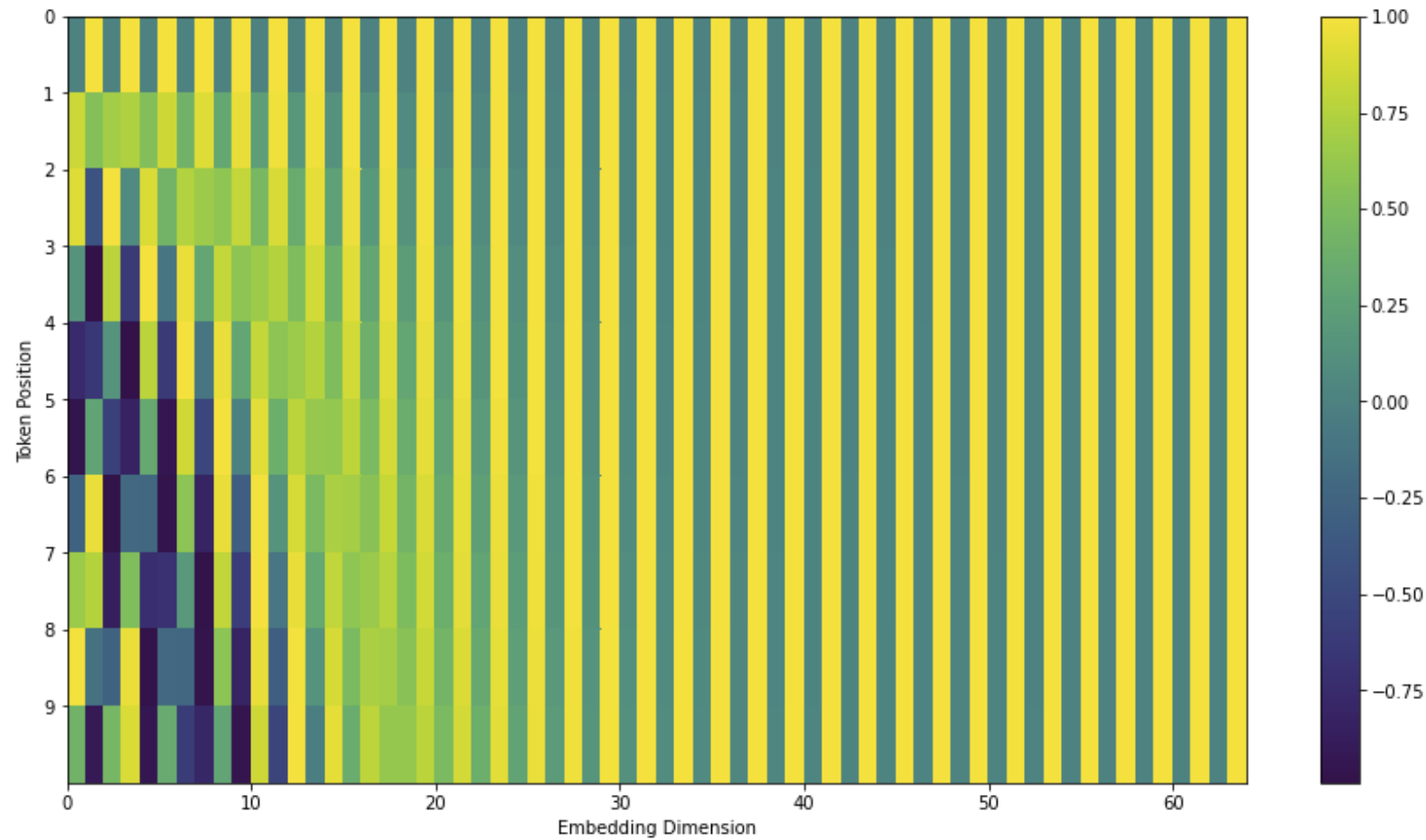
$$X = X + PE(X)$$

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

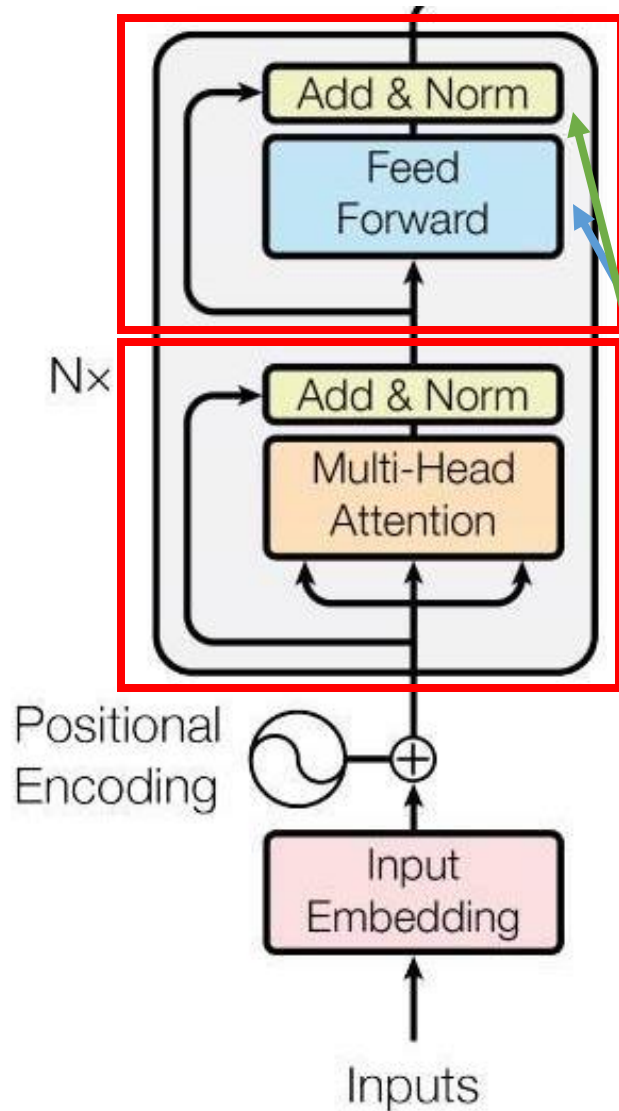
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

- PE has the same dimension with the input embeddings of each token.
- pos : Position of each token in the input, $[0, L-1]$
- i : index of the embeddings, $[0, d_{model} - 1]$
- What does 10000 mean here?
 - The maximum input length

What Does Positional Encoding Look Like?



Feedforward Network

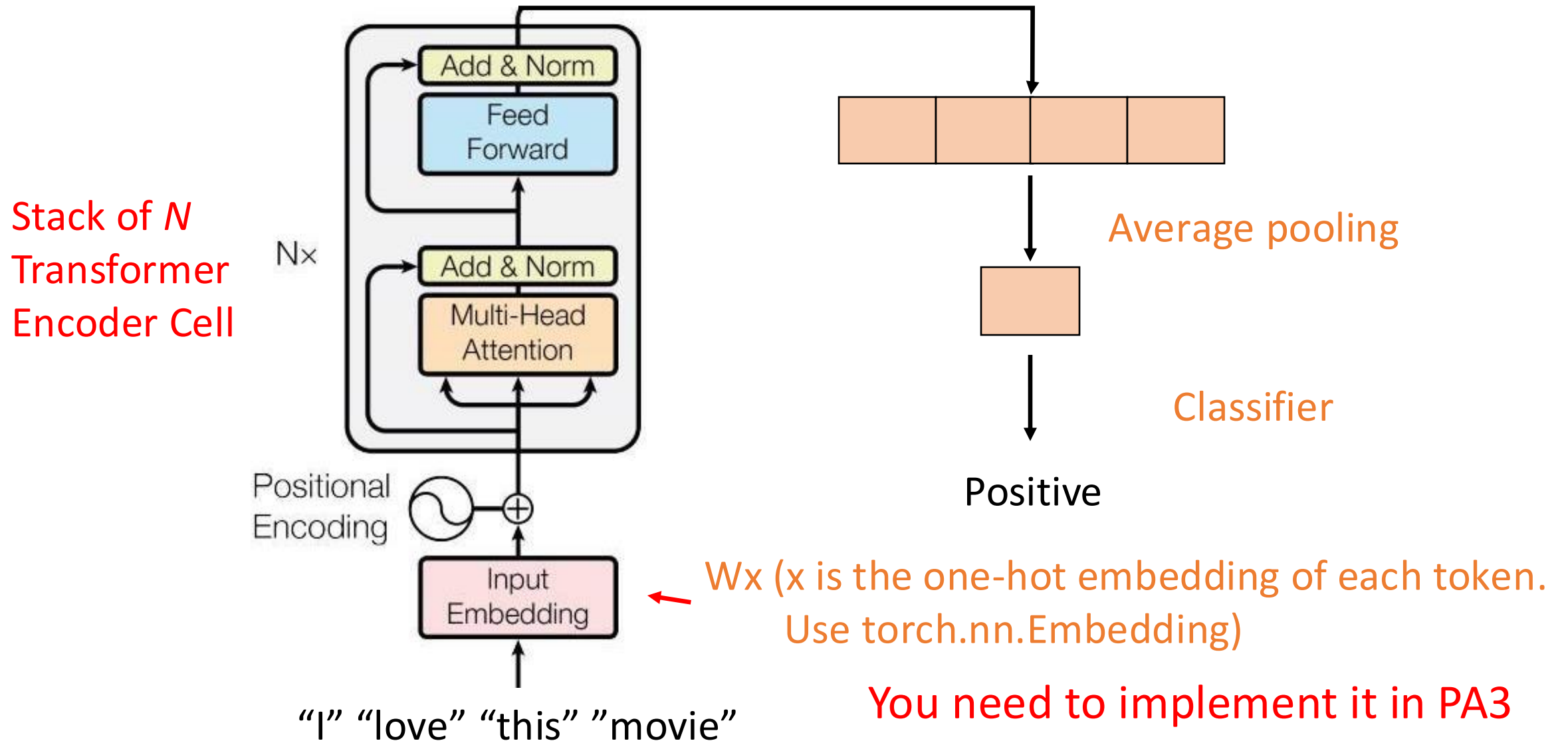


Are there any non-linearities in the lower part?

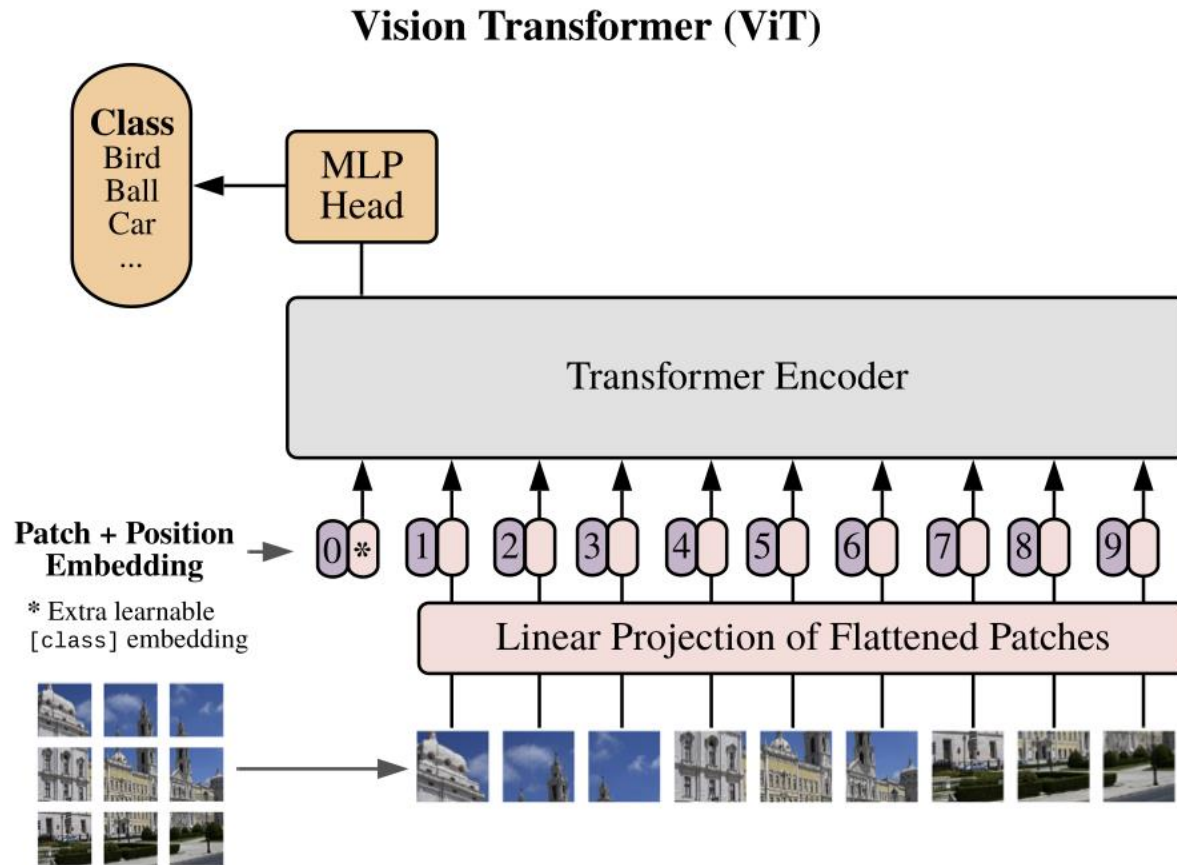
A two-layer fully-connected network (one hidden layer).

Residual connections, dropout, and normalization work in a similar way to the multi-head self-attention module

Transformer Encoder for Text Classification



Transformer Encoder for Image Classification

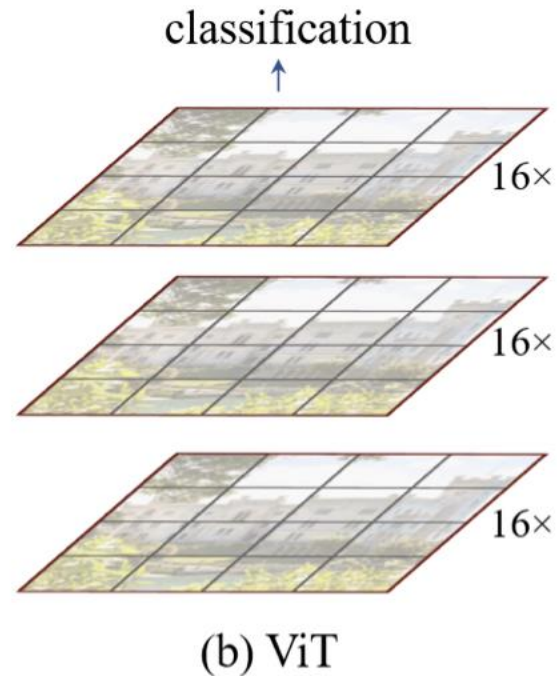


How to represent each patch?
Flattened RGB pixels

Extra credits in PA3:

- Use the special **[CLASS]** token to aggregate the image's information
- Alternatively, we get embeddings of each patch and do average pooling as in the previous example
- The details of the Transformer encoder is slightly different
- Positional encoding: 1D is sufficient.

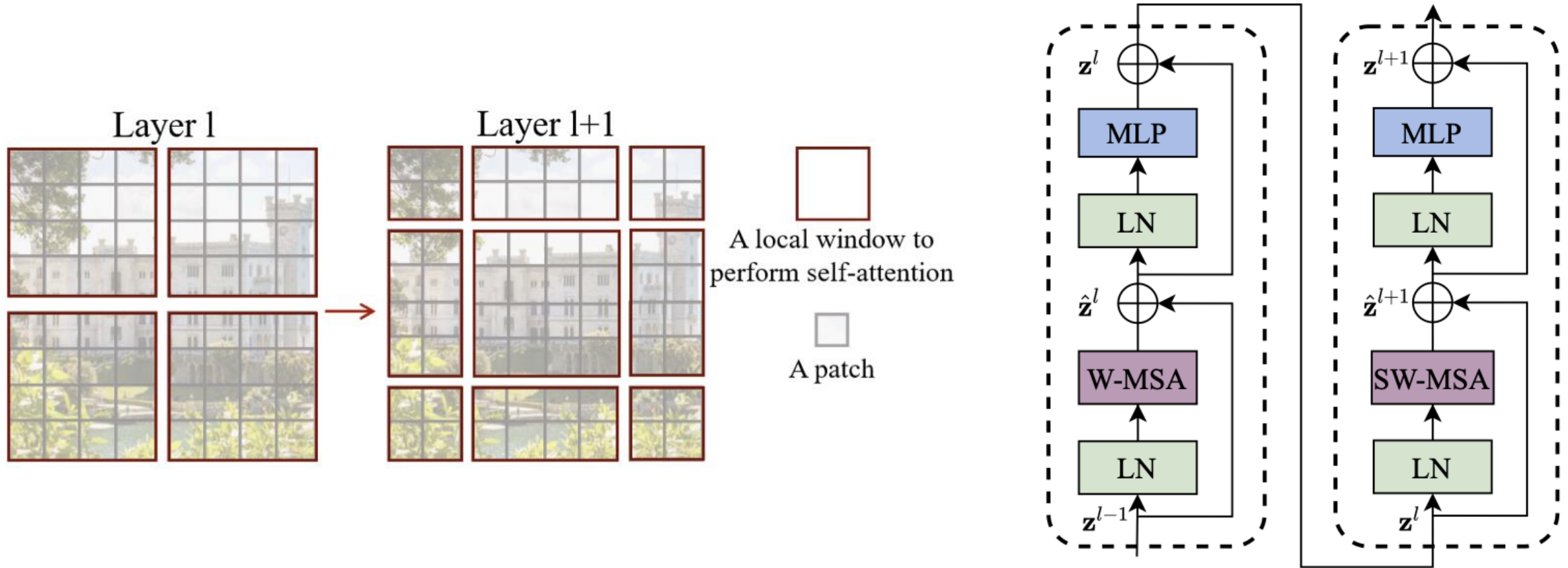
SwinTransformer



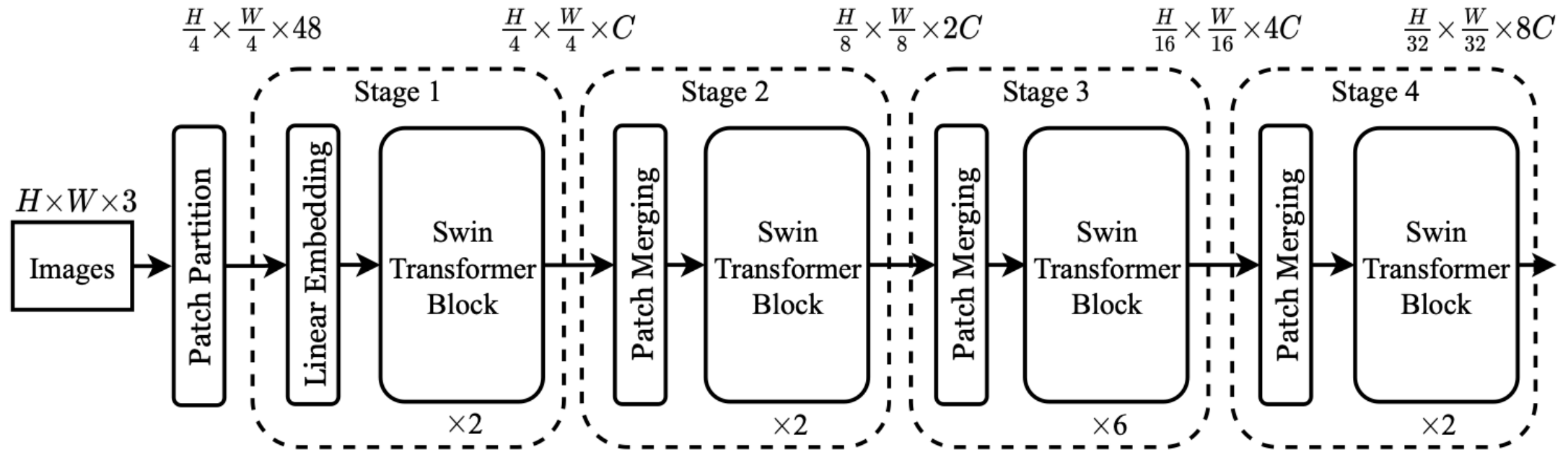
Lesson we learned from CNNs:

- We process the input at different hierarchies (resolutions) for efficiency purpose.
- We process the input at different hierarchies (resolutions) for different receptive field sizes.
- We need to fuse features from different hierarchies (resolutions) for dense prediction tasks

Shifted Window-based Efficient Attention



Hierarchical Representations in SwinTransformer



Transformer vs CNN

(a) Various frameworks							
Method	Backbone	AP ^{box}	AP ₅₀ ^{box}	AP ₇₅ ^{box}	#param.	FLOPs	FPS
Cascade	R-50	46.3	64.3	50.5	82M	739G	18.0
Mask R-CNN	Swin-T	50.5	69.3	54.9	86M	745G	15.3
ATSS	R-50	43.5	61.9	47.0	32M	205G	28.3
	Swin-T	47.2	66.5	51.3	36M	215G	22.3
RepPointsV2	R-50	46.5	64.6	50.3	42M	274G	13.6
	Swin-T	50.0	68.5	54.2	45M	283G	12.0
Sparse R-CNN	R-50	44.5	63.4	48.2	106M	166G	21.0
	Swin-T	47.9	67.3	52.3	110M	172G	18.4

(b) Various backbones w. Cascade Mask R-CNN									
	AP ^{box}	AP ₅₀ ^{box}	AP ₇₅ ^{box}	AP ^{mask}	AP ₅₀ ^{mask}	AP ₇₅ ^{mask}	param	FLOPs	FPS
DeiT-S [†]	48.0	67.2	51.7	41.4	64.2	44.3	80M	889G	10.4
R50	46.3	64.3	50.5	40.1	61.7	43.4	82M	739G	18.0
Swin-T	50.5	69.3	54.9	43.7	66.6	47.1	86M	745G	15.3
X101-32	48.1	66.5	52.4	41.6	63.9	45.2	101M	819G	12.8
Swin-S	51.8	70.4	56.3	44.7	67.9	48.5	107M	838G	12.0
X101-64	48.3	66.4	52.3	41.7	64.0	45.1	140M	972G	10.4
Swin-B	51.9	70.9	56.5	45.0	68.4	48.7	145M	982G	11.6

Next Class

- Transformer Decoder