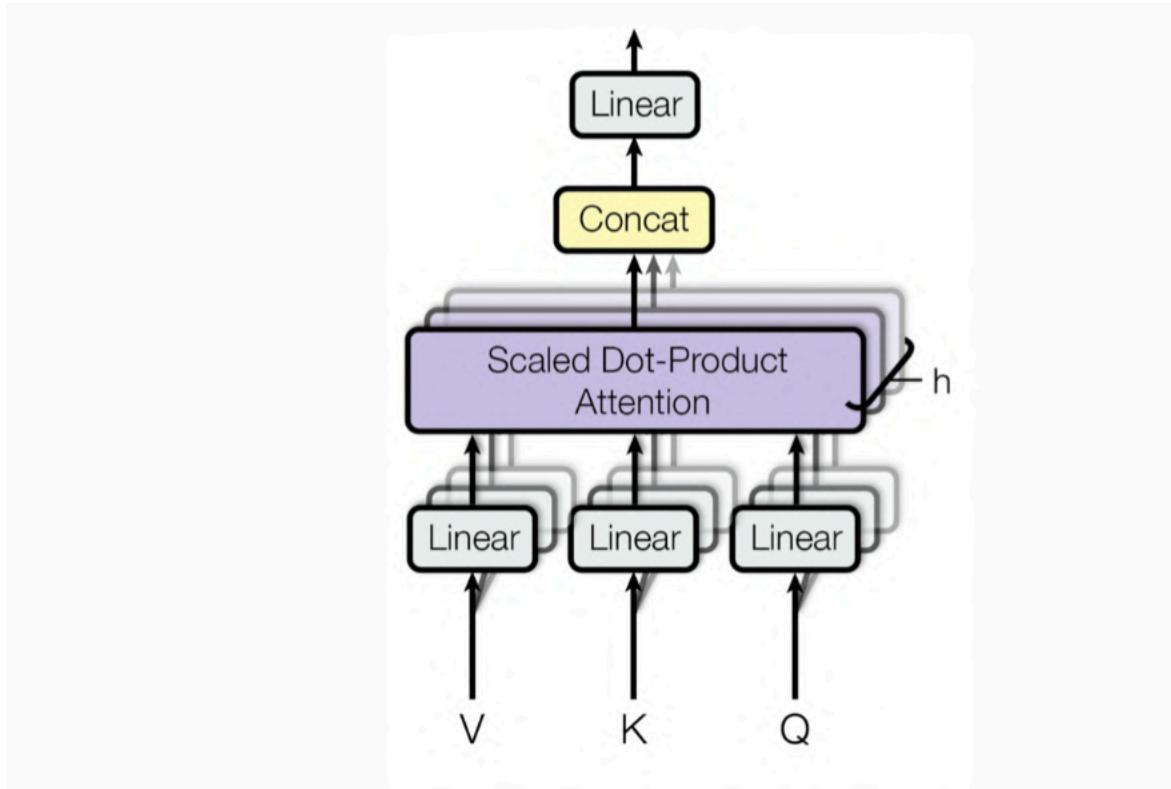


Attention and Transformers



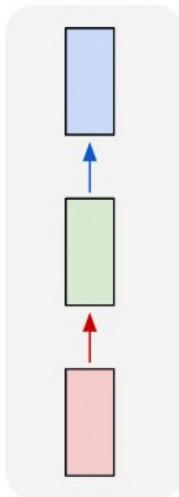
Outline

- General module of Attention
- Self-Attention
- Transformers
- Visual Transformer (ViT)

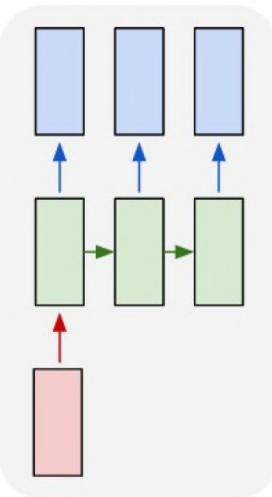
Last Week

Common RNN Structures

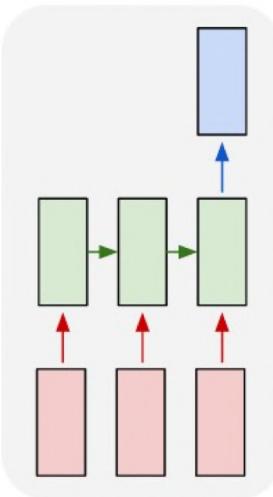
one to one



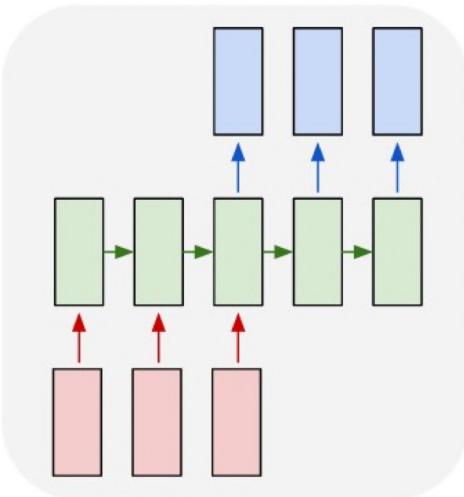
one to many



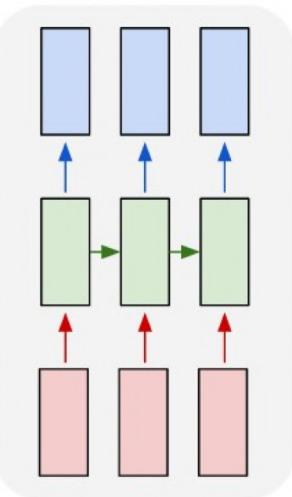
many to one



many to many

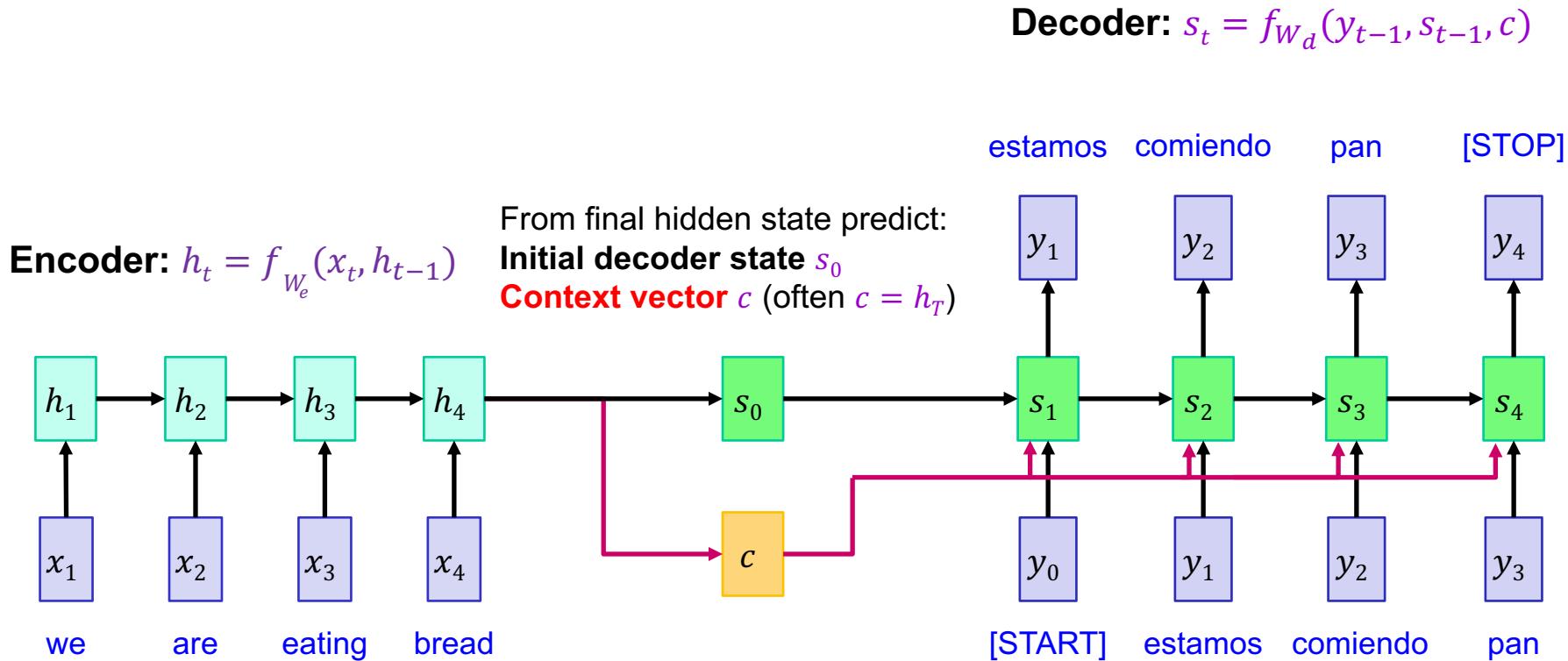


many to many



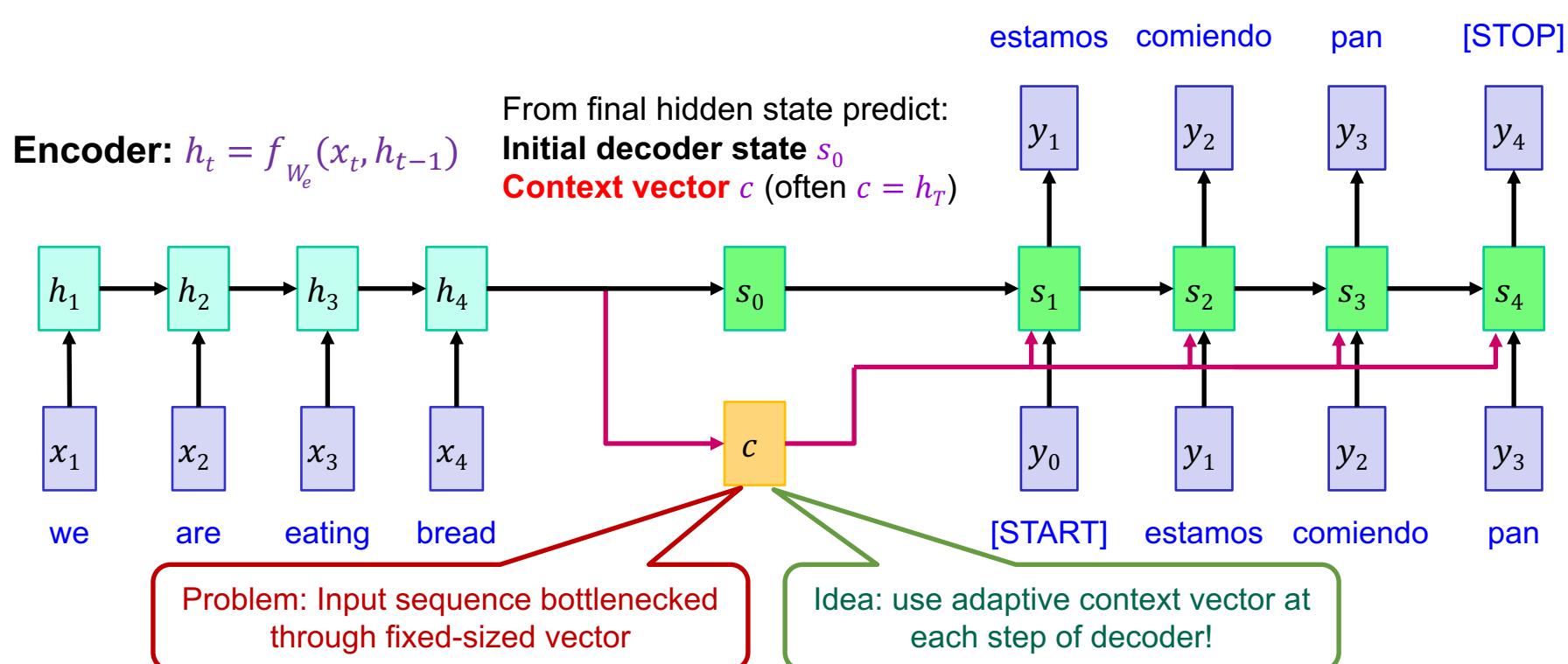
- One of the benefits of recurrent neural networks is the ability to **handle arbitrary length inputs and outputs**.
- This flexibility allows us to define a broad range of tasks.

Seq2seq: RNNs with a fixed context



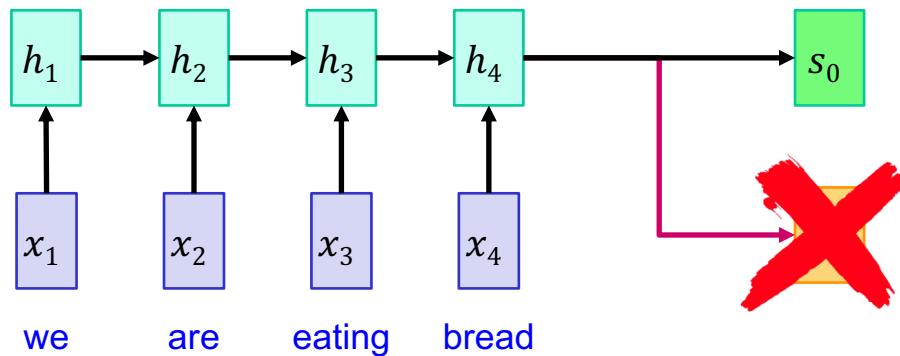
Seq2seq: RNNs with a fixed context

Decoder: $s_t = f_{W_d}(y_{t-1}, s_{t-1}, c)$

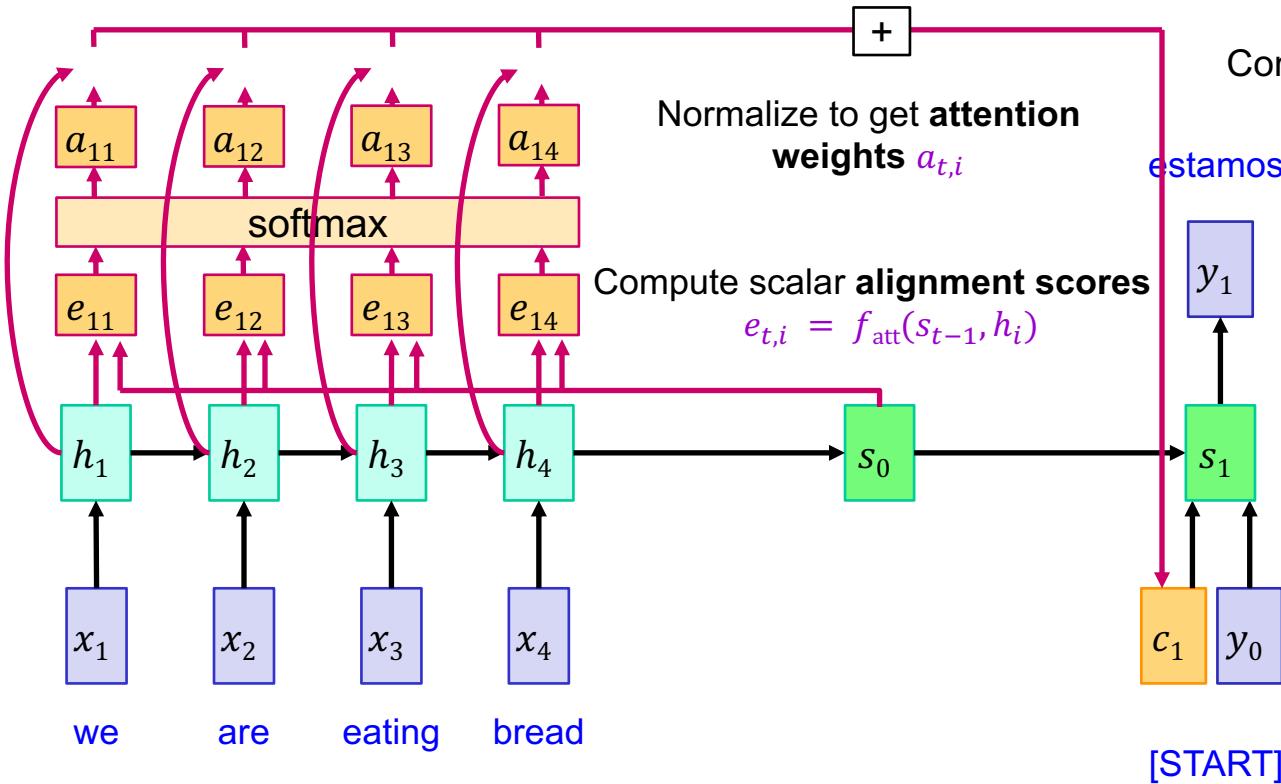


Seq2seq: RNNs with attention

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



Seq2seq: RNNs with attention



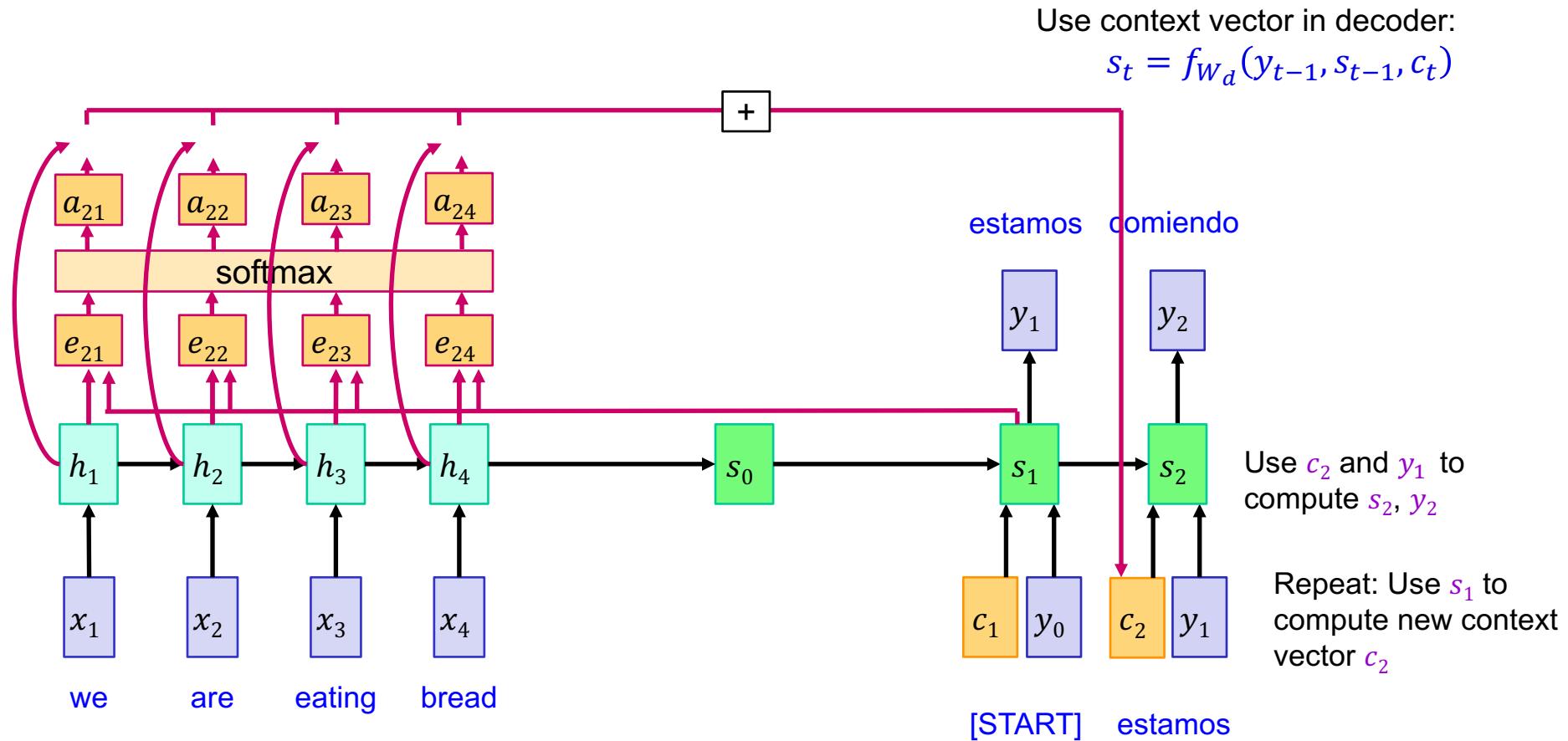
Use context vector in decoder:

$$s_t = f_{W_d}(y_{t-1}, s_{t-1}, c_t)$$

Compute **context vector** as

$$c_t = \sum_i a_{t,i} h_i$$

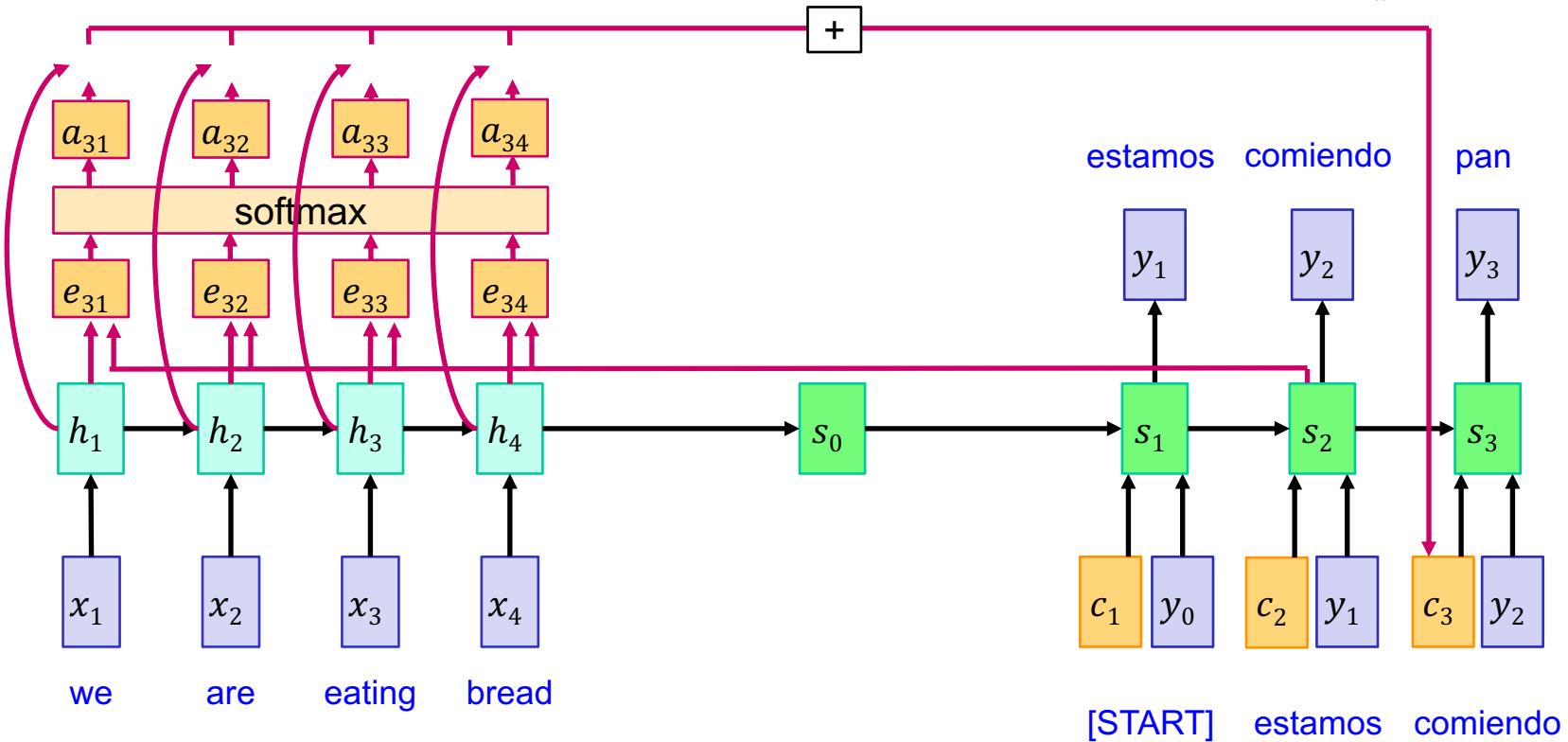
Seq2seq: RNNs with attention



Seq2seq: RNNs with attention

Use context vector in decoder:

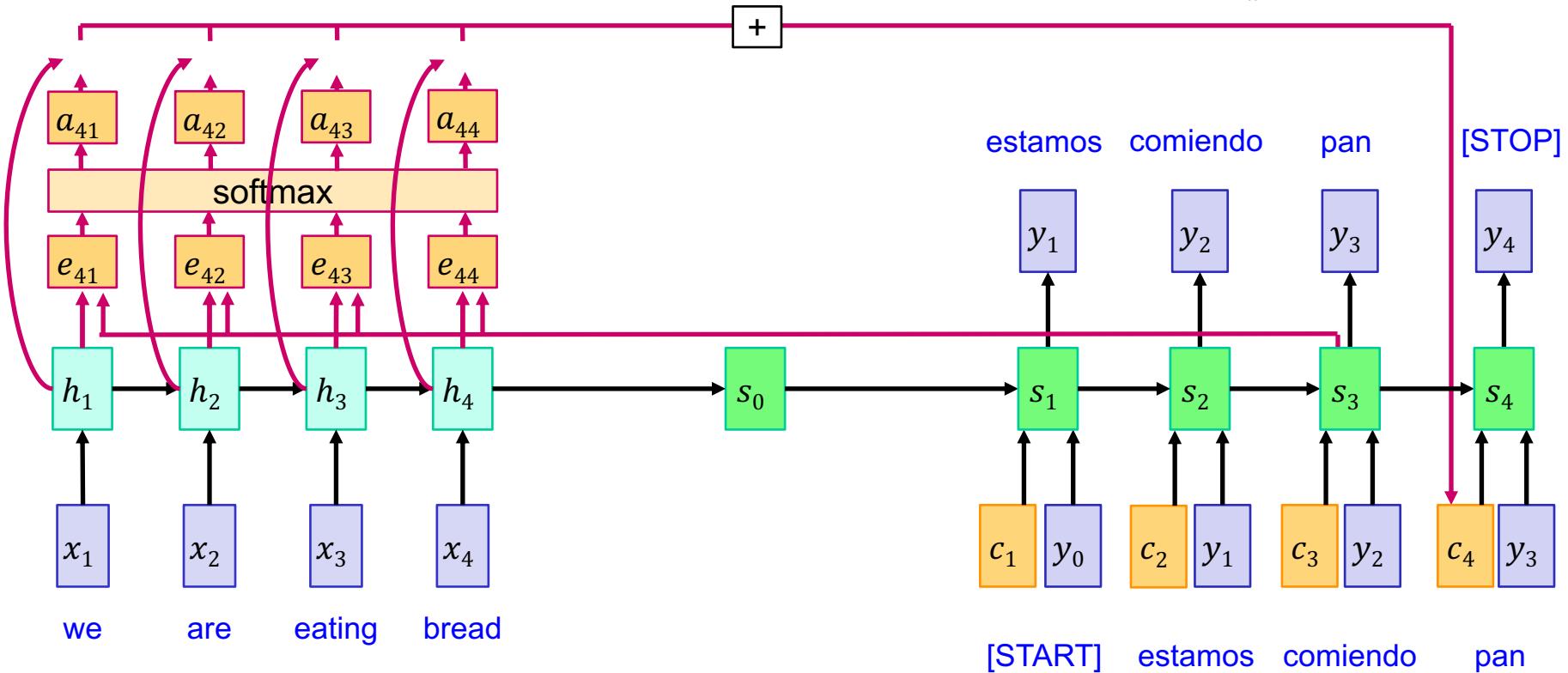
$$s_t = f_{W_d}(y_{t-1}, s_{t-1}, c_t)$$



Seq2seq: RNNs with attention

Use context vector in decoder:

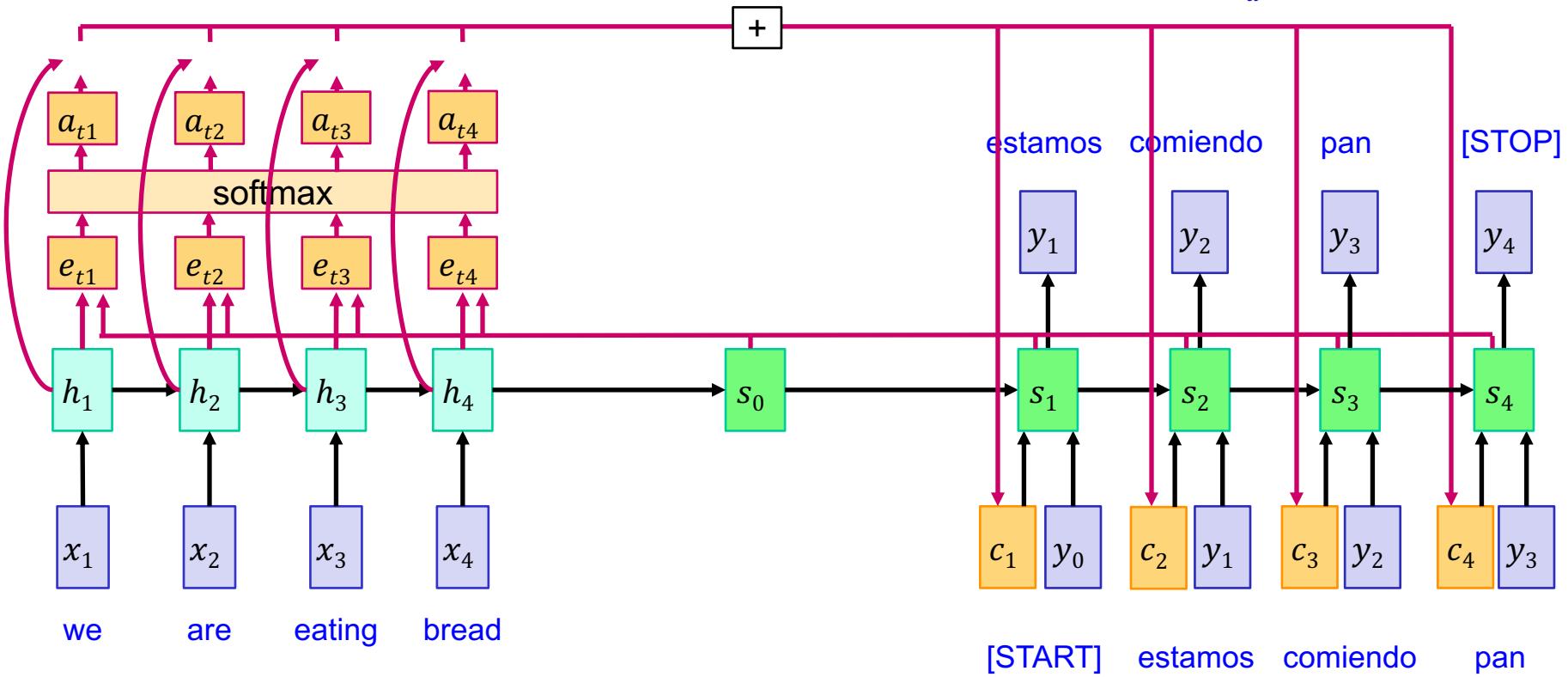
$$s_t = f_{W_d}(y_{t-1}, s_{t-1}, c_t)$$



Seq2seq: RNNs with attention

Use context vector in decoder:

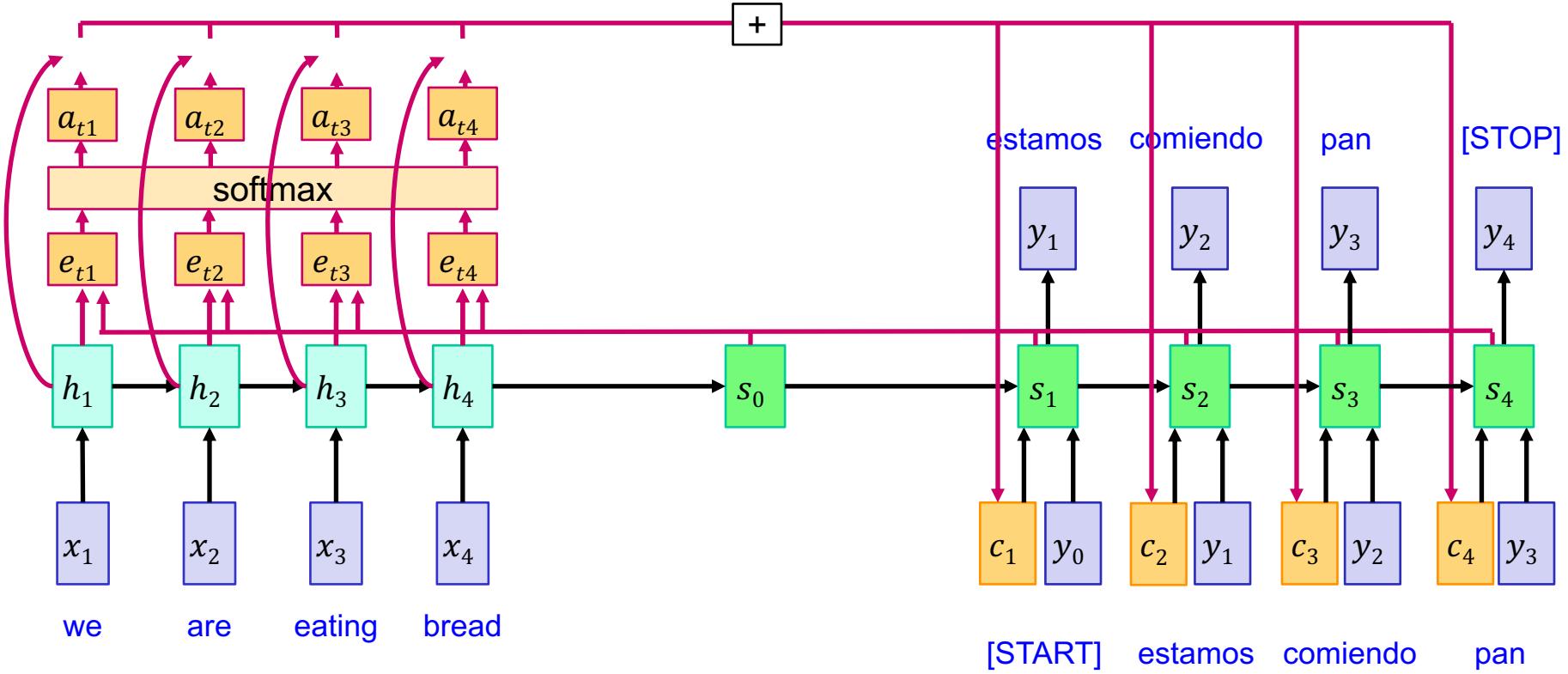
$$s_t = f_{W_d}(y_{t-1}, s_{t-1}, c_t)$$



Seq2seq: RNNs with attention

Use context vector in decoder:

$$s_t = f_{W_d}(y_{t-1}, s_{t-1}, c_t)$$



- Intuition: Context vector attends to the relevant part of the input sequence, e.g. “estamos”=“we are”, so maybe $a_{11} = a_{12} = 0.45$ and $a_{13} = a_{14} = 0.05$

Seq2seq: RNNs with attention

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

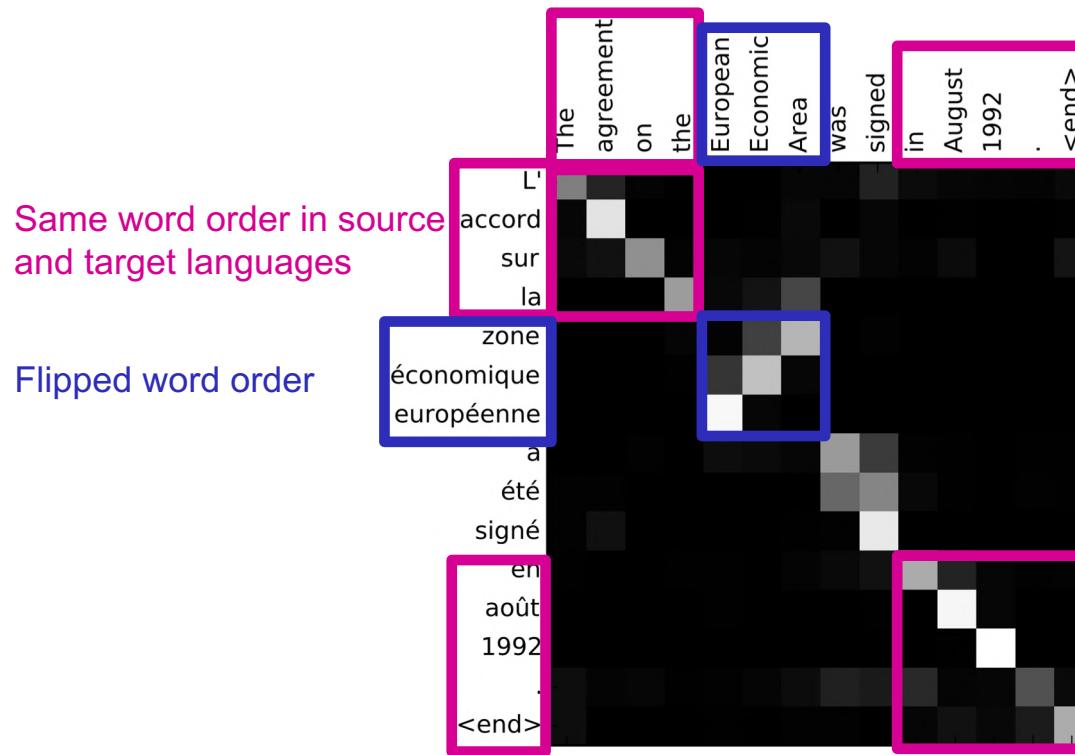


Image captioning with attention

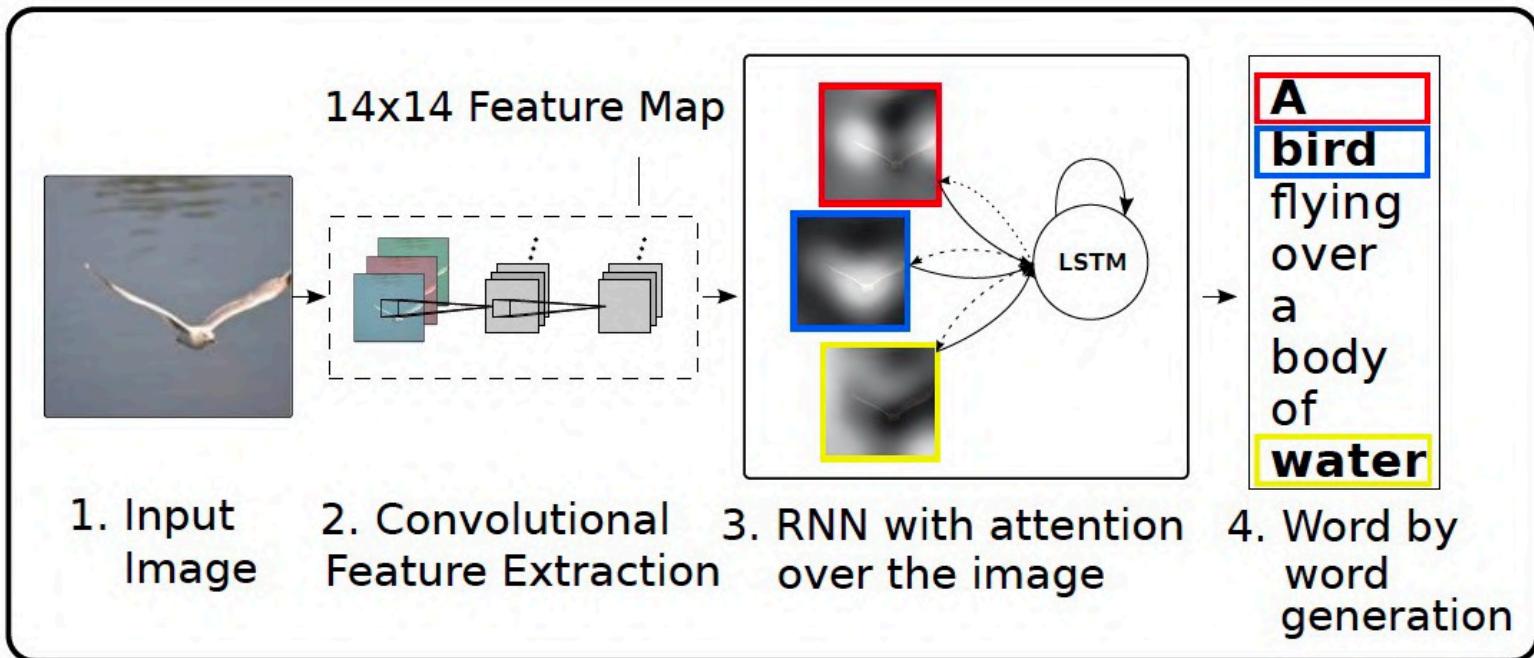


Image captioning with RNNs and attention

- 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

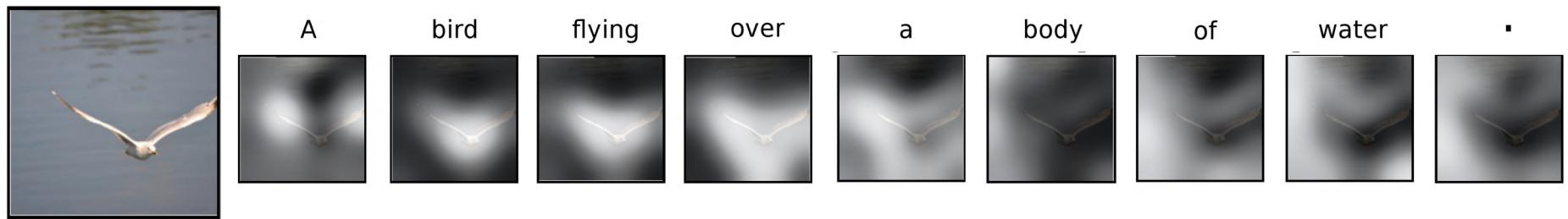
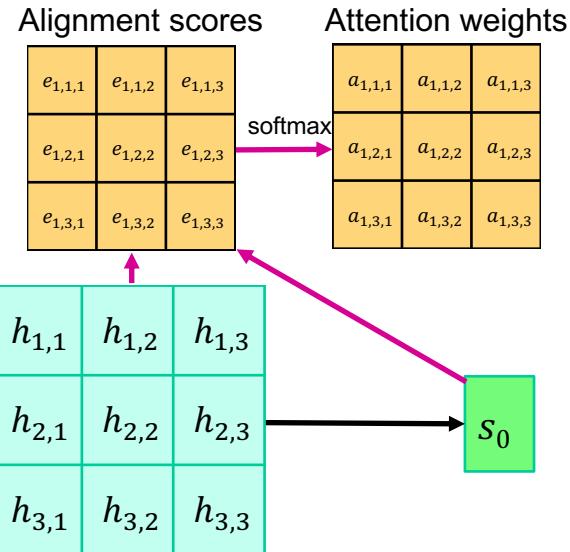


Image captioning with RNNs and attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$



Use a CNN to
compute a grid of
features for an image

Image captioning with RNNs and attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

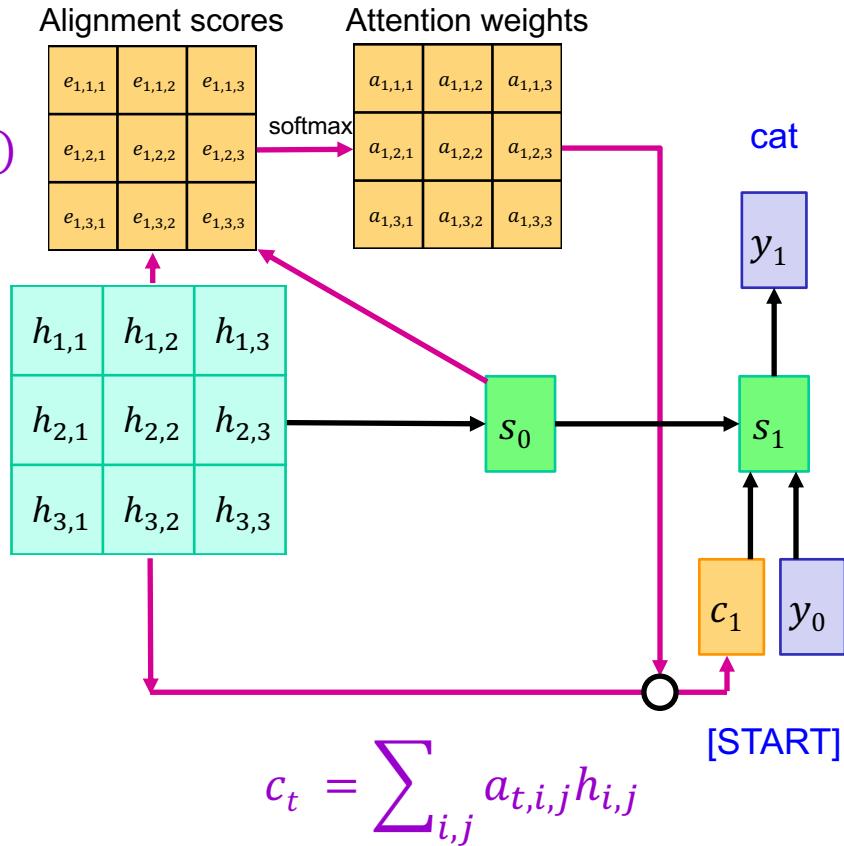


Image captioning with RNNs and attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

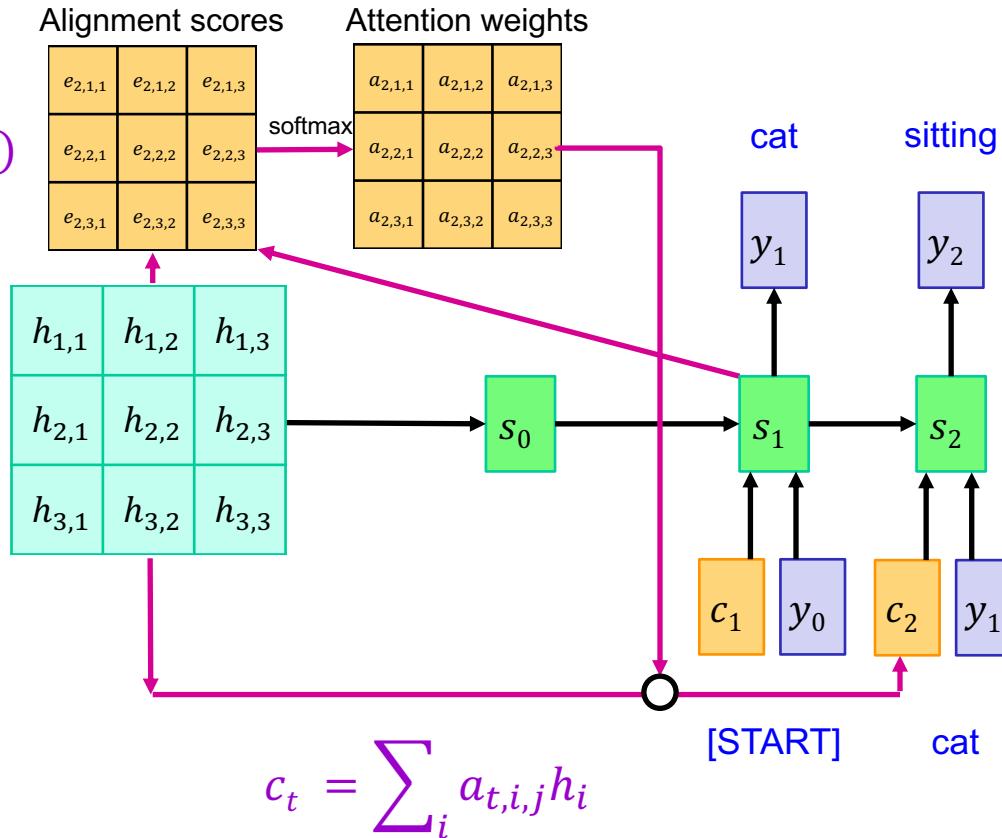
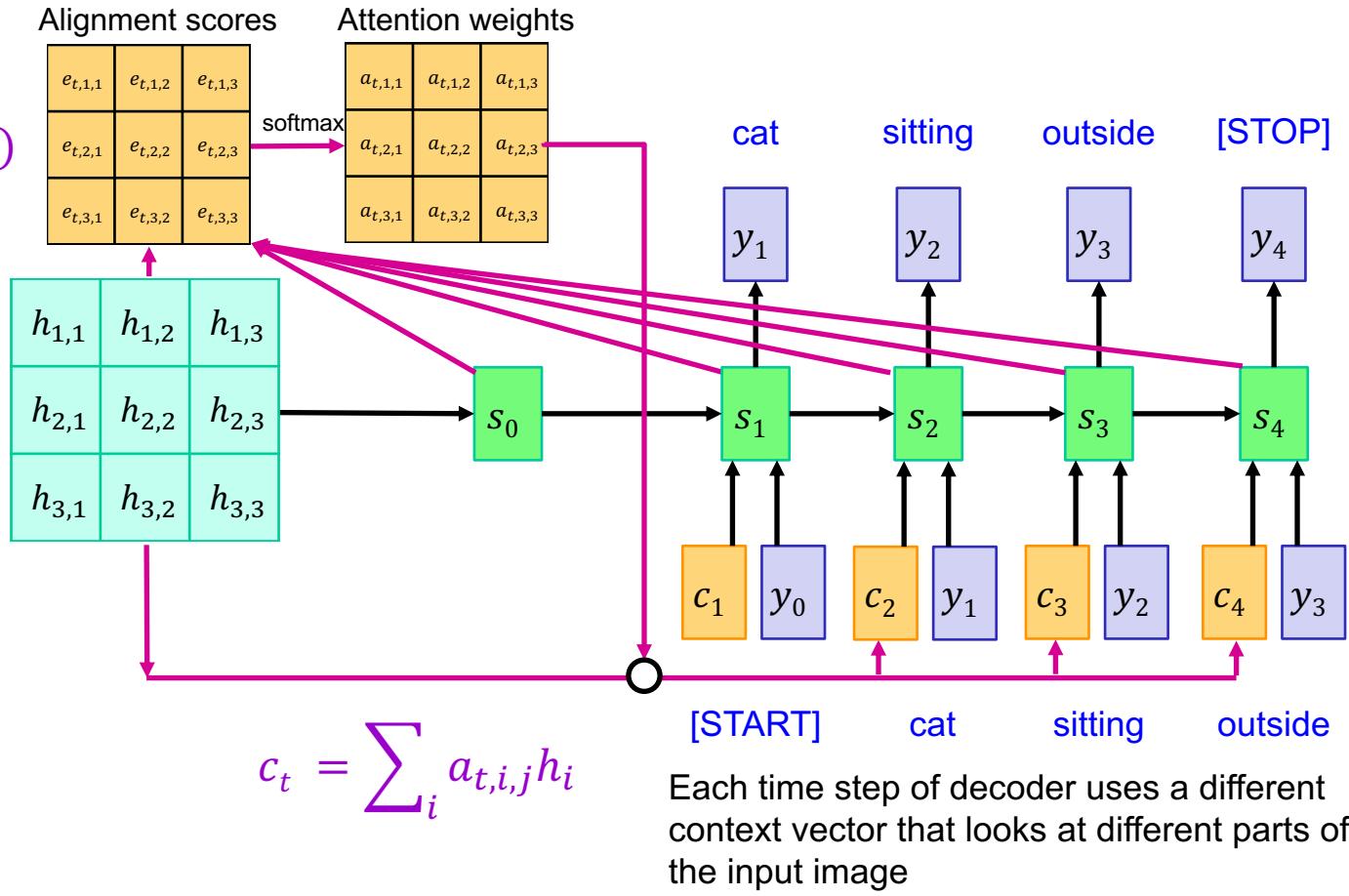


Image captioning with RNNs and attention

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$



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.

Example results

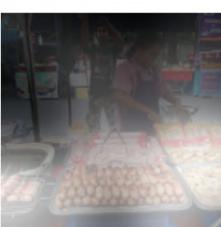
- Mistakes



A large white bird standing in a forest.

A woman holding a clock in her hand.

A man wearing a hat and a hat on a skateboard.



A person is standing on a beach with a surfboard.

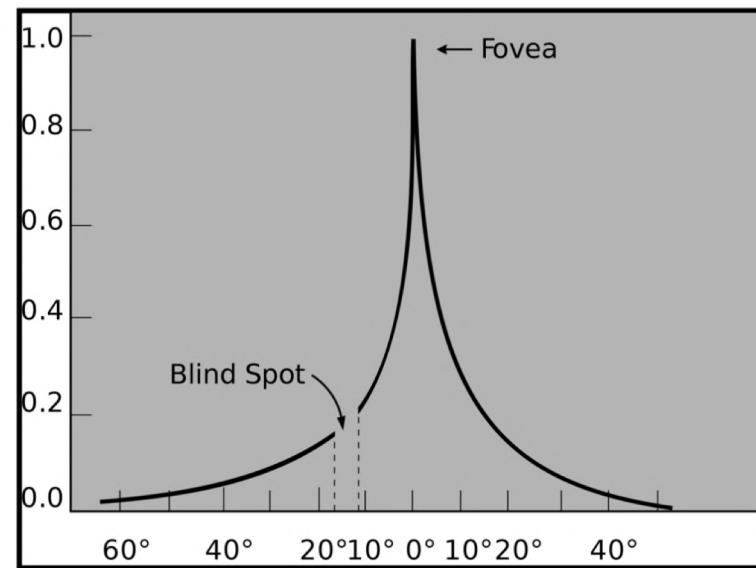
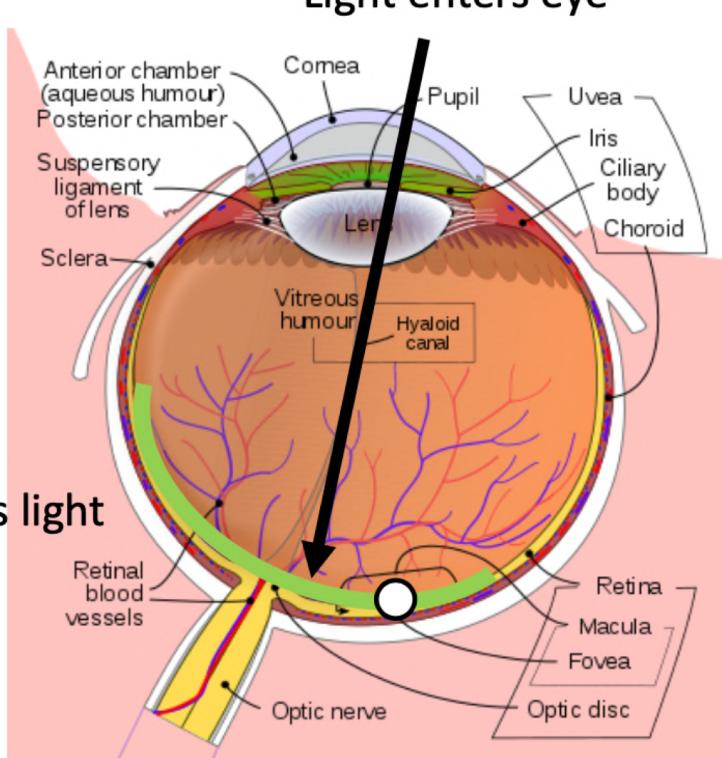
A woman is sitting at a table with a large pizza.

A man is talking on his cell phone while another man watches.

Motivation

The Fovea in Human Vision:

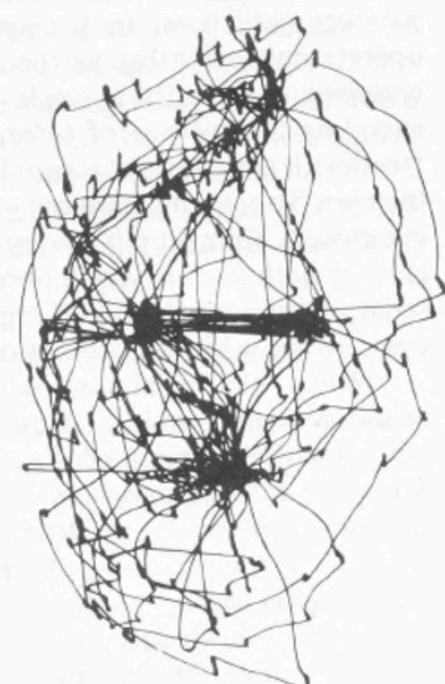
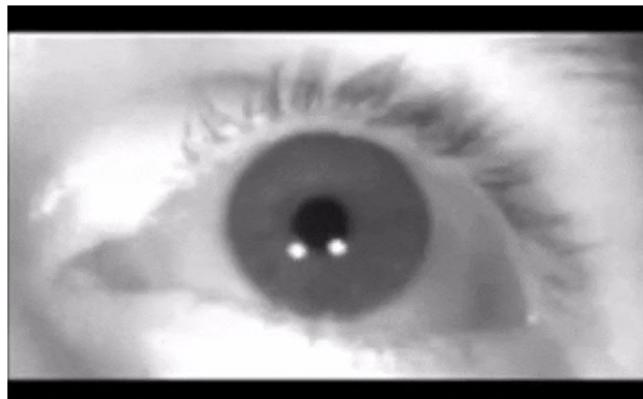
- The Fovea is a tiny region in the retina that can see in high acuity.



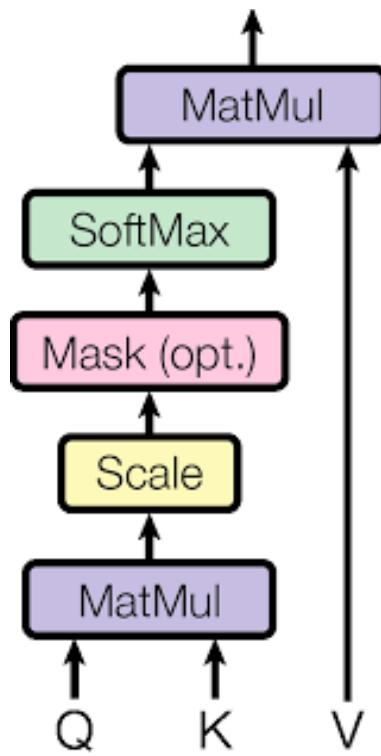
Motivation

The Fovea in Human Vision:

- The Fovea is a tiny region in the retina that can see in high acuity.
- Saccade: Human eyes are constantly moving to "fill the gaps".
- Attention weights at each time step is similar to the saccades in the human eye.



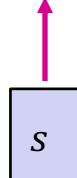
A General Attention Model



We saw Attention in image captioning

Features

$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$



Inputs:

Features: \mathbf{h} (*shape*: $H \times W \times D$)
Query: \mathbf{s} (*shape*: D)

Attention we saw in image captioning

Features

$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$

$e_{1,1}$	$e_{1,2}$	$e_{1,3}$
$e_{2,1}$	$e_{2,2}$	$e_{2,3}$
$e_{3,1}$	$e_{3,2}$	$e_{3,3}$

Alignment

s

Operations:

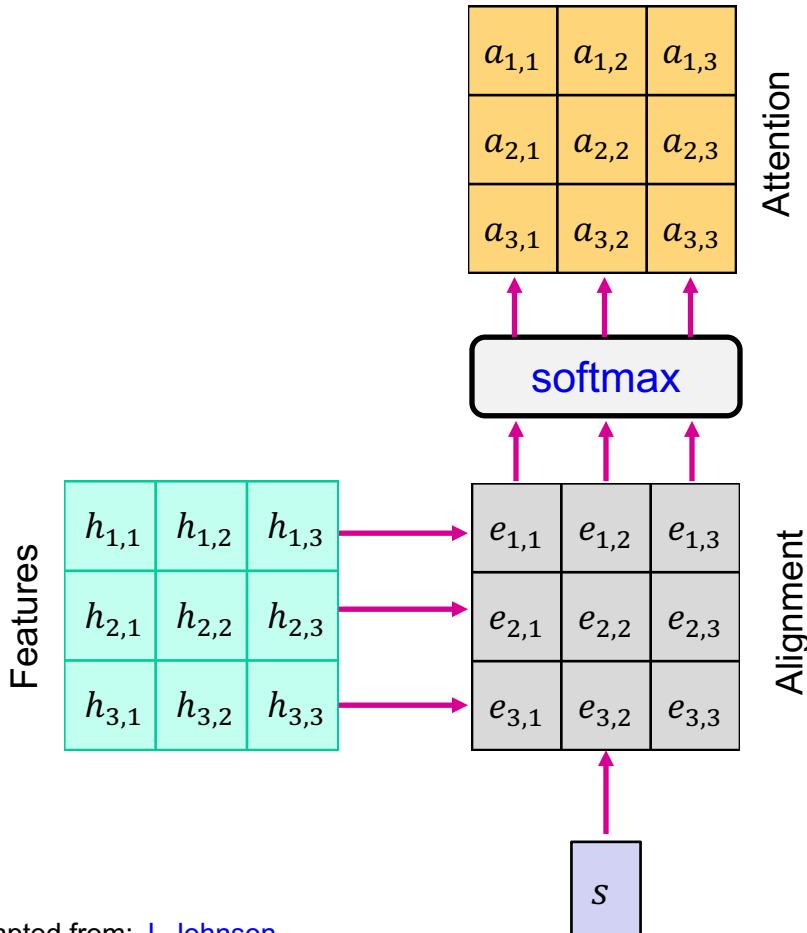
Alignments: $e_{i,j} = f_{att}(s, h_{i,j})$

Inputs:

Features: h (*shape*: $H \times W \times D$)

Query: s (*shape*: D)

Attention we saw in image captioning



Operations:

Alignments: $e_{i,j} = f_{att}(\mathbf{s}, \mathbf{h}_{i,j})$

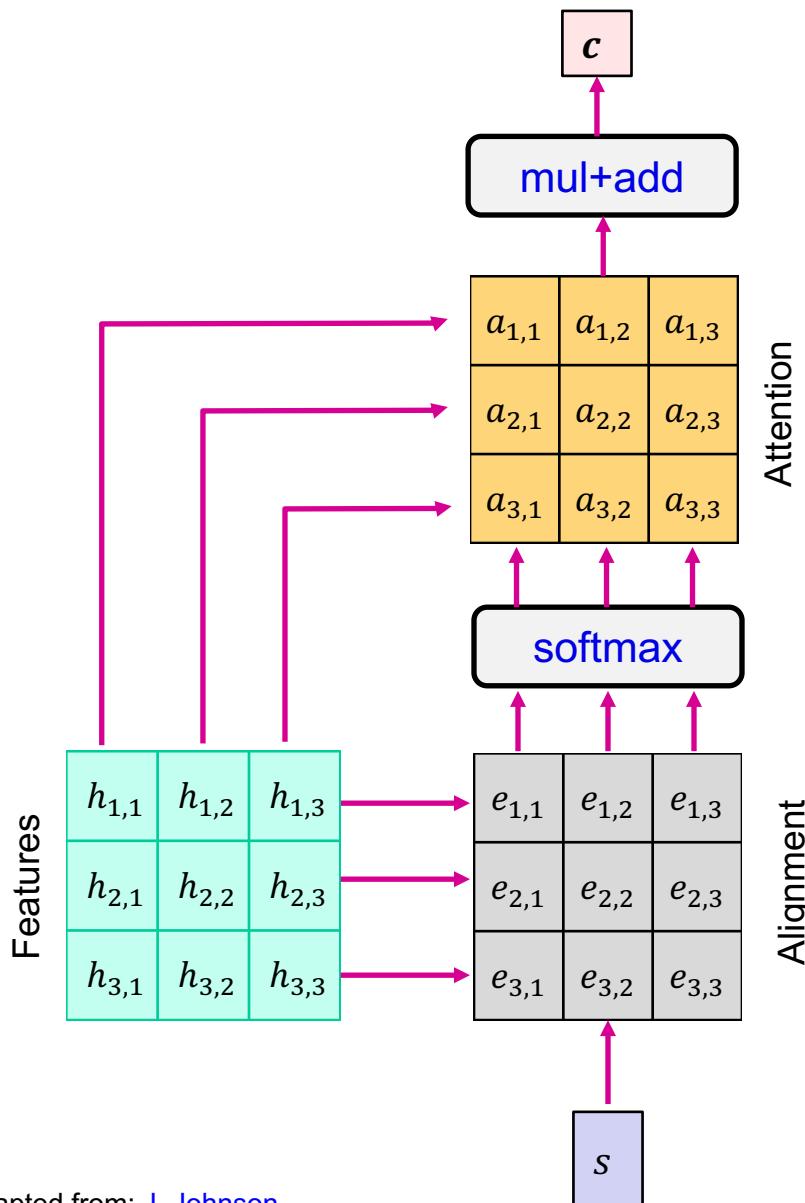
Attention: $\mathbf{a} = \text{softmax}(\mathbf{e})$

Inputs:

Features: \mathbf{h} (shape: $H \times W \times D$)

Query: \mathbf{s} (shape: D)

Attention we saw in image captioning



Outputs:

Context vector: c (shape: D)

Operations:

Alignments: $e_{i,j} = f_{att}(s, h_{i,j})$

Attention: $a = \text{softmax}(e)$

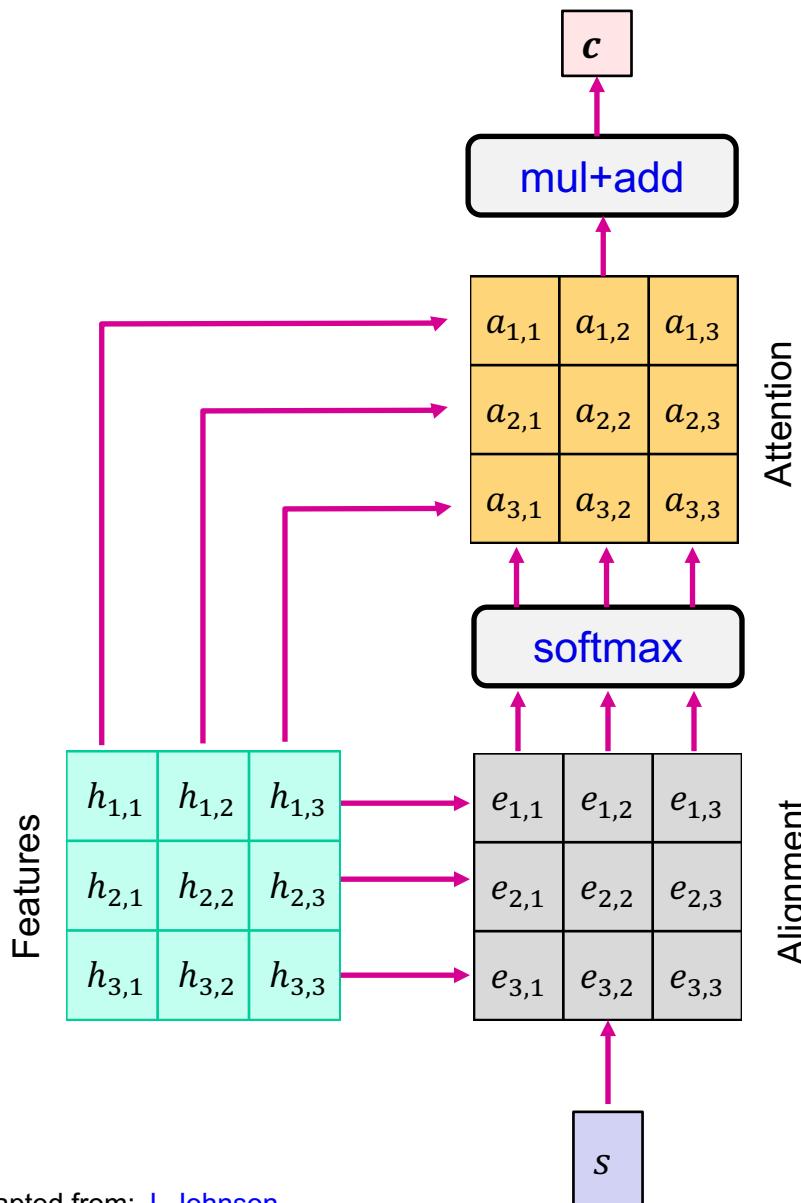
Output: $c = \sum_{i,j} a_{i,j} h_{i,j}$

Inputs:

Features: h (shape: $H \times W \times D$)

Query: s (shape: D)

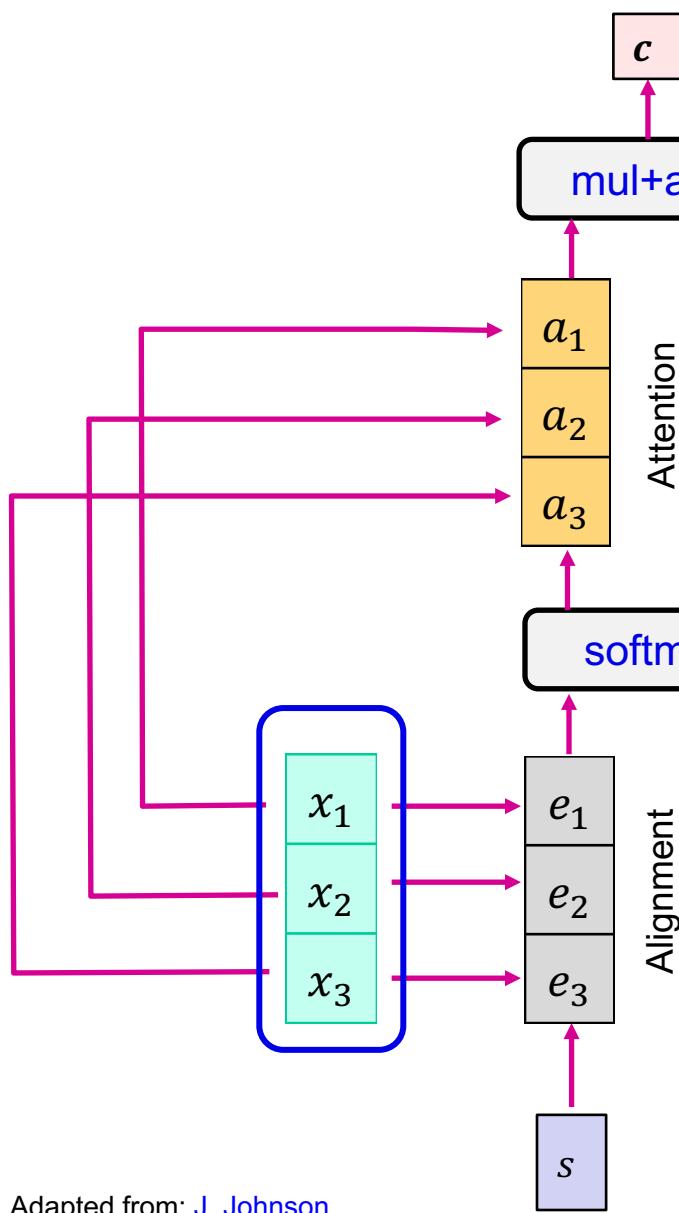
Attention we saw in image captioning



General attention Layer

- Attention operation is **permutation invariant**.
- Doesn't care about ordering of the features
- Stretch $H \times W = N$ feature vectors into N vectors x_i

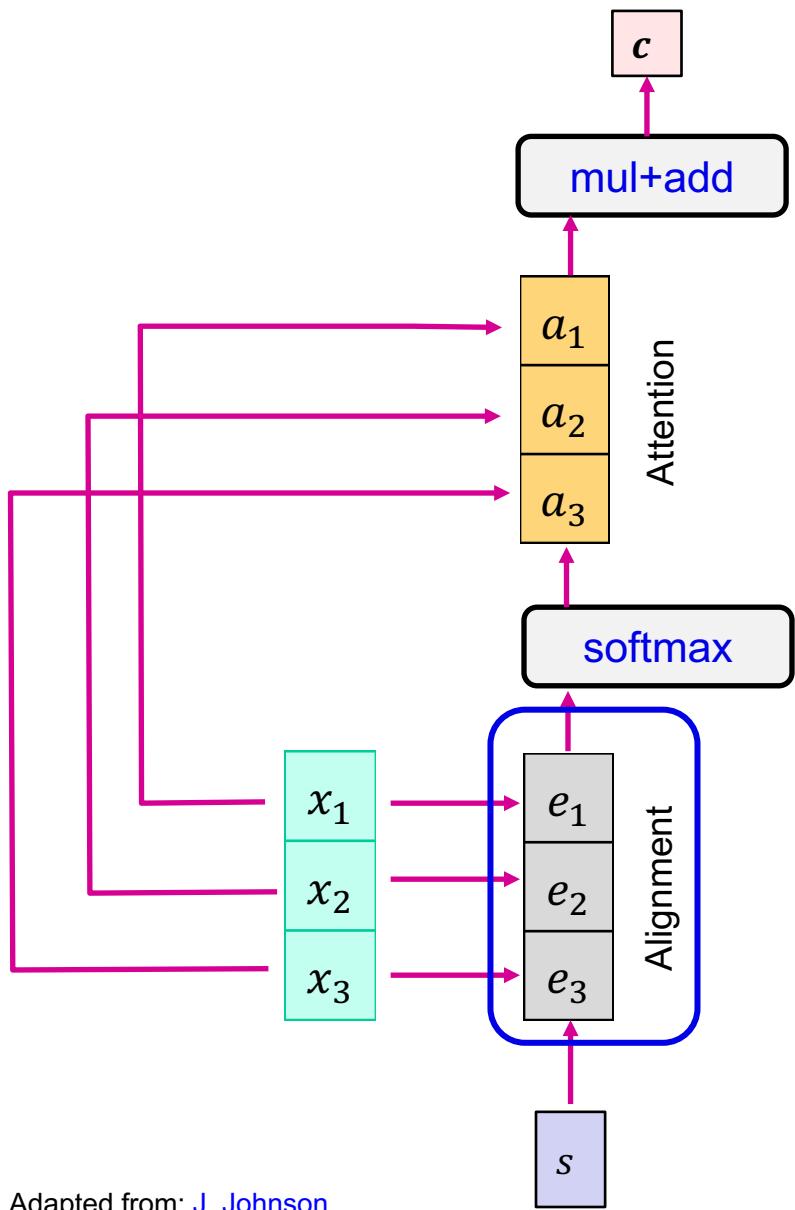
General attention layer



Inputs:

Input vectors: $\mathbf{X} = [x_1, \dots]^T$ (shape: $N \times D$)
Query: \mathbf{s} (shape: D)

General attention layer



Outputs:

Operations: Scaled inner product

Alignments: $e_{i,j} = f_{att}(s, h_{i,j})$

Alignments: $e_i = \mathbf{s} \cdot \mathbf{x}_i / \sqrt{D}$

Recall $\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}||\mathbf{b}| \cos(\text{angle})$

Suppose that \mathbf{a} and \mathbf{b} are constant vectors of dimension D .

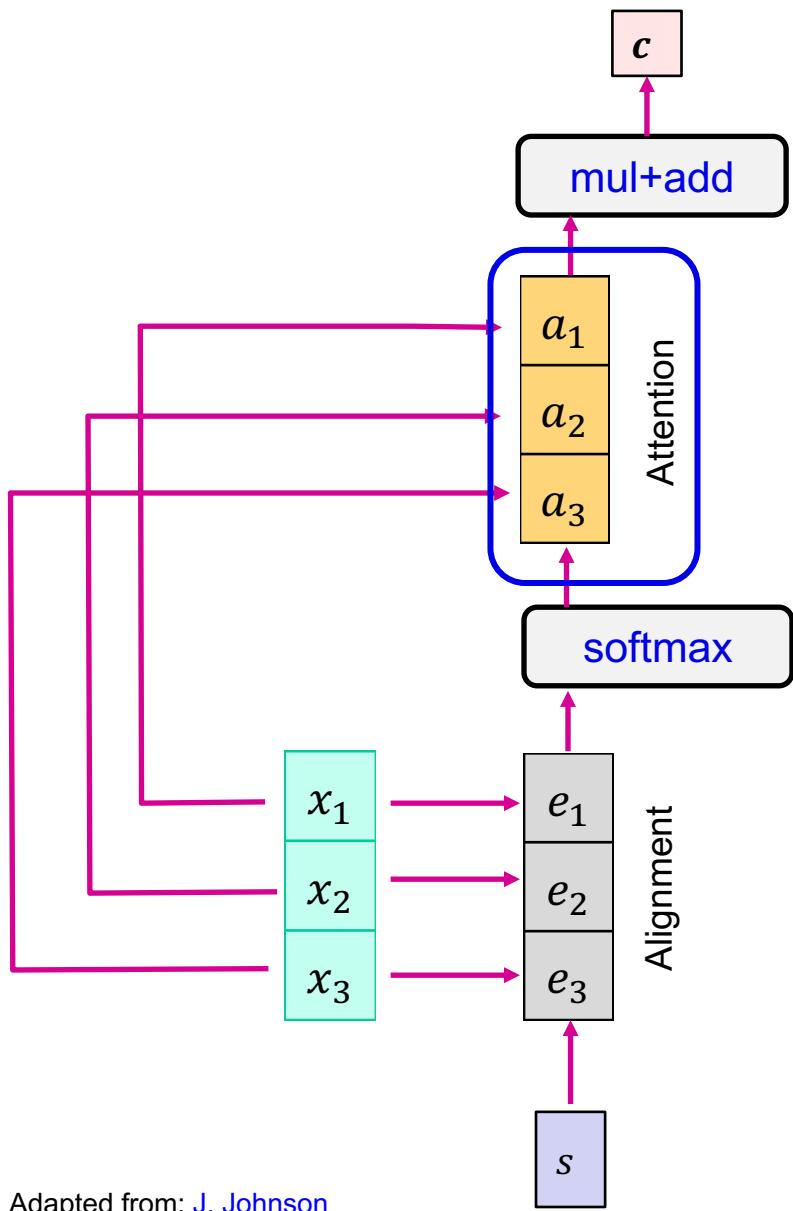
Then, $|\mathbf{a}| = (\sum_i a^2)^{1/2} = a\sqrt{D}$

Inputs:

Input vectors: $\mathbf{X} = [x_1, \dots]^T$ (shape: $N \times D$)

Query: \mathbf{s} (shape: D)

General attention layer



Outputs:

Operations:

Alignments: $e_{i,j} = f_{att}(s, h_{i,j})$

Alignments: $e_i = \mathbf{s} \cdot \mathbf{x}_i / \sqrt{D}$

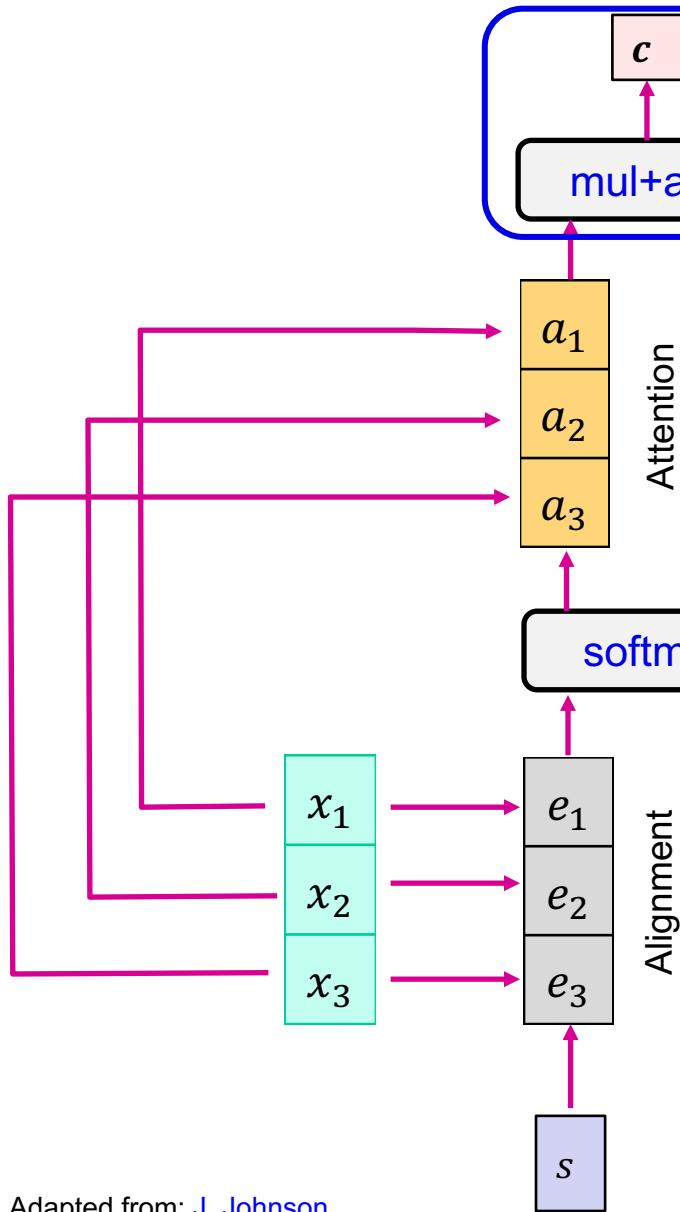
Attention: $\mathbf{a} = softmax(\mathbf{e})$ (shape: N)

Inputs:

Input vectors: $\mathbf{X} = [x_1, \dots]^T$ (shape: $N \times D$)

Query: \mathbf{s} (shape: D)

General attention layer



Outputs:

$$c = \sum_i a_i x_i \quad (\text{shape: } D)$$

Operations:

Alignments: $e_{i,j} = f_{att}(s, h_{i,j})$

Alignments: $e_i = s \cdot x_i / \sqrt{D}$

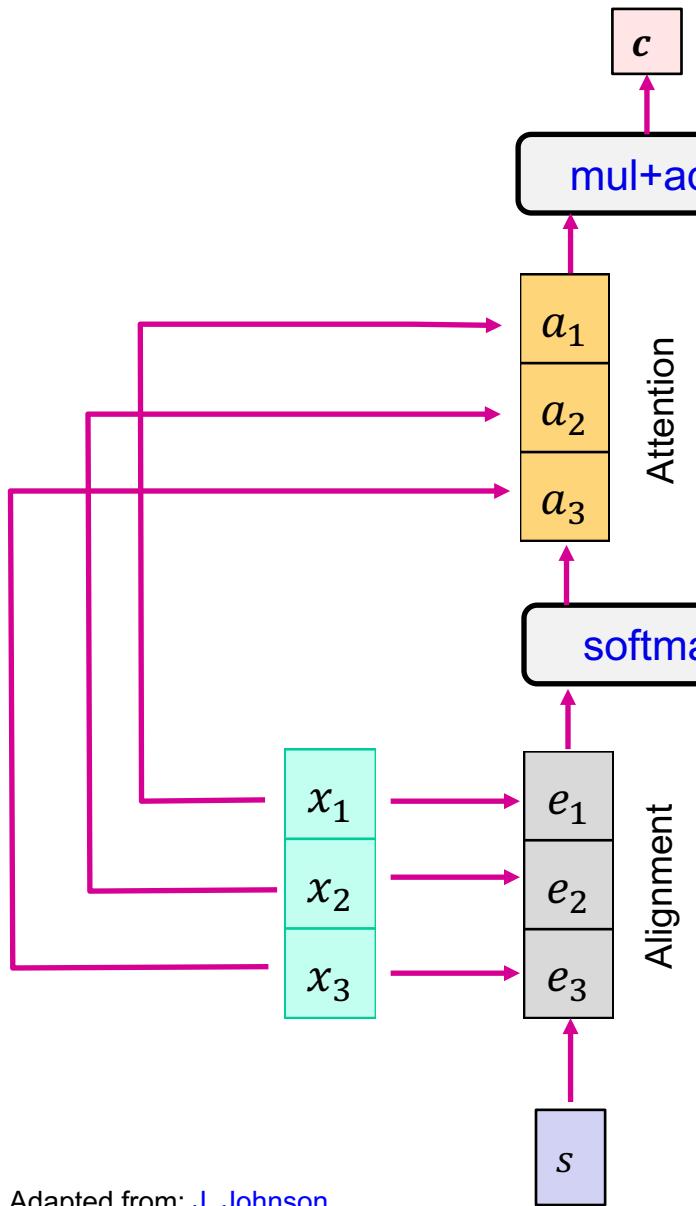
Attention: $a = \text{softmax}(e)$ (shape: N)

Inputs:

Input vectors: $\mathbf{X} = [x_1, \dots]^T$ (shape: $N \times D$)

Query: s (shape: D)

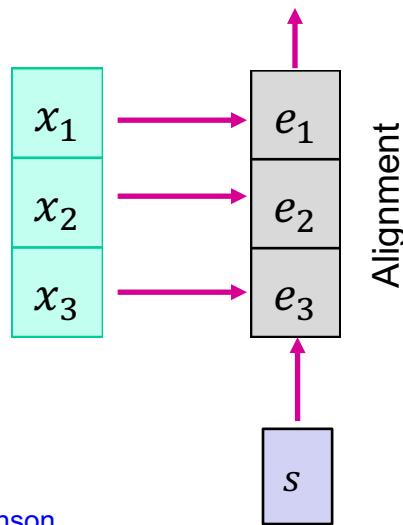
General attention layer: multiple queries



Multiple Queries

- We can extend the attention layer into multiple query vectors
- Each query creates a new output context vector

General attention layer: multiple queries

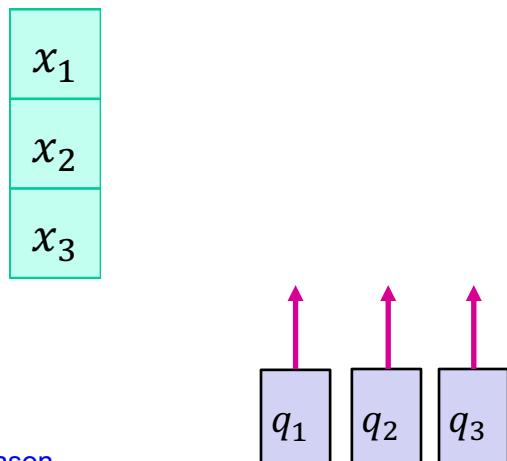


Inputs:

Input vectors: $\mathbf{X} = [x_1, \dots]^T$ (*shape*: $N \times D$)

Query: \mathbf{s} (*shape*: D)

General attention layer

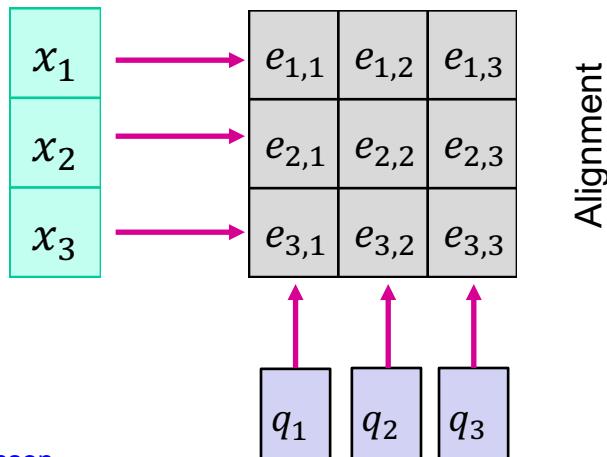


Inputs:

Input vectors: $\mathbf{X} = [x_1, \dots]^T$ (shape: $N \times D$)

Query: $\mathbf{Q} = [q_1, q_2, \dots]^T$ (shape: $M \times D$)

General attention layer



Operations:

Alignments: $e_{i,j} = \mathbf{x}_i \cdot \mathbf{q}_j / \sqrt{D}$

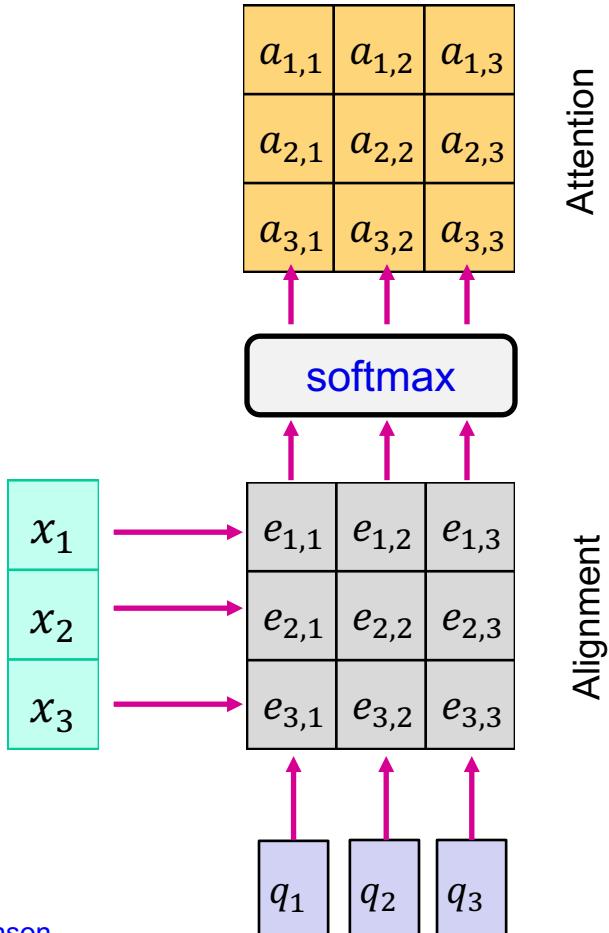
$E = \mathbf{X}\mathbf{Q}^T / \sqrt{D}$ (shape: $N \times M$)

Inputs:

Input vectors: $\mathbf{X} = [x_1, \dots]^T$ (shape: $N \times D$)

Query: $\mathbf{Q} = [q_1, q_2, \dots]^T$ (shape: $M \times D$)

General attention layer



Operations:

Alignments: $e_{i,j} = \mathbf{x}_i \cdot \mathbf{q}_j / \sqrt{D}$

$E = \mathbf{X}\mathbf{Q}^T / \sqrt{D}$ (shape: $N \times M$)

Attention: $a_{*,j} = \text{softmax}(e_{*,j})$

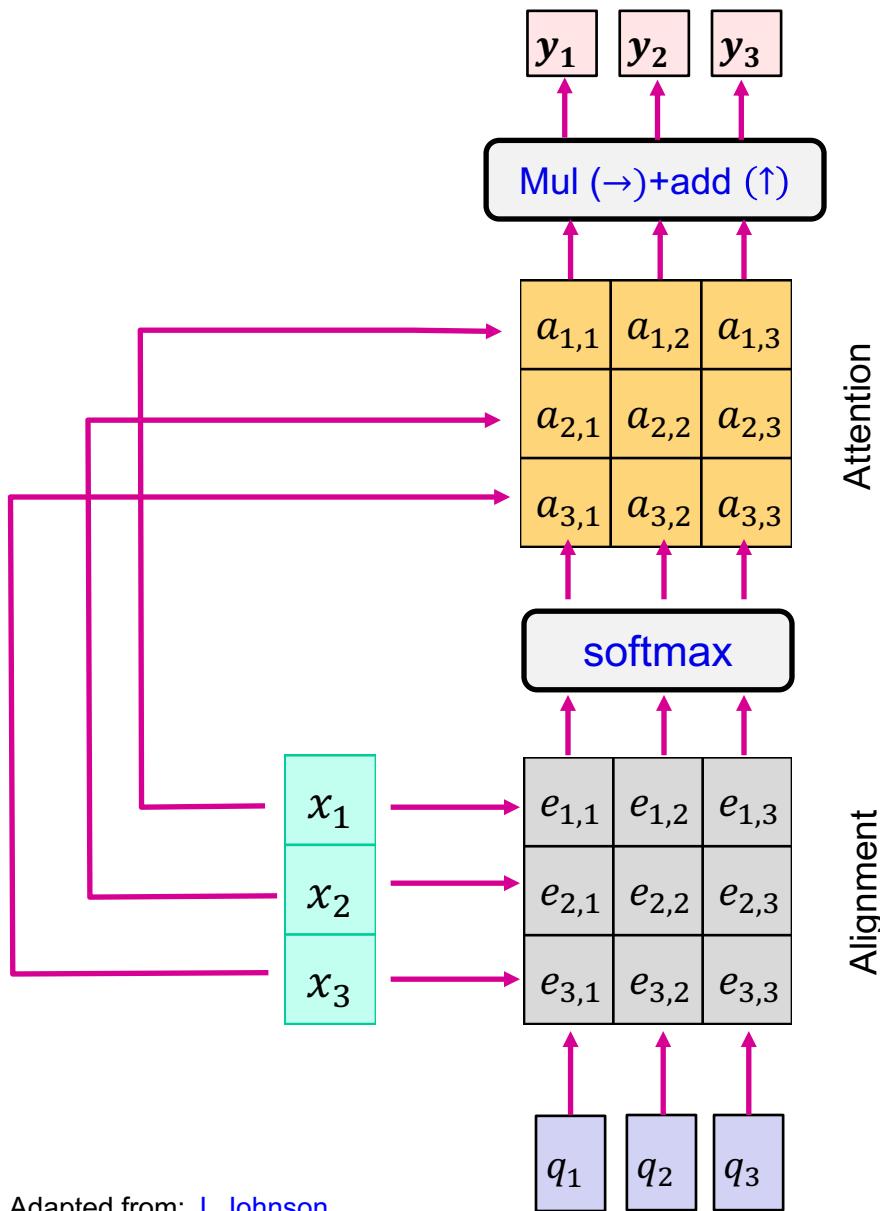
$A = \text{softmax}(E, \text{dim} = 1)$ ($N \times M$)

Inputs:

Input vectors: $\mathbf{X} = [x_1, \dots]^T$ (shape: $N \times D$)

Query: $\mathbf{Q} = [q_1, q_2, \dots]^T$ (shape: $M \times D$)

General attention layer



Outputs:

$$y_j = \sum_i a_{i,j} x_i ; \quad \mathbf{y} = A^T \mathbf{X} \quad (\text{shape: } M \times D)$$

Operations:

$$\text{Alignments: } e_{i,j} = \mathbf{x}_i \cdot \mathbf{q}_j / \sqrt{D}$$

$$E = \mathbf{X} \mathbf{Q}^T / \sqrt{D} \quad (\text{shape: } N \times M)$$

$$\text{Attention: } a_{*,j} = \text{softmax}(e_{*,j})$$

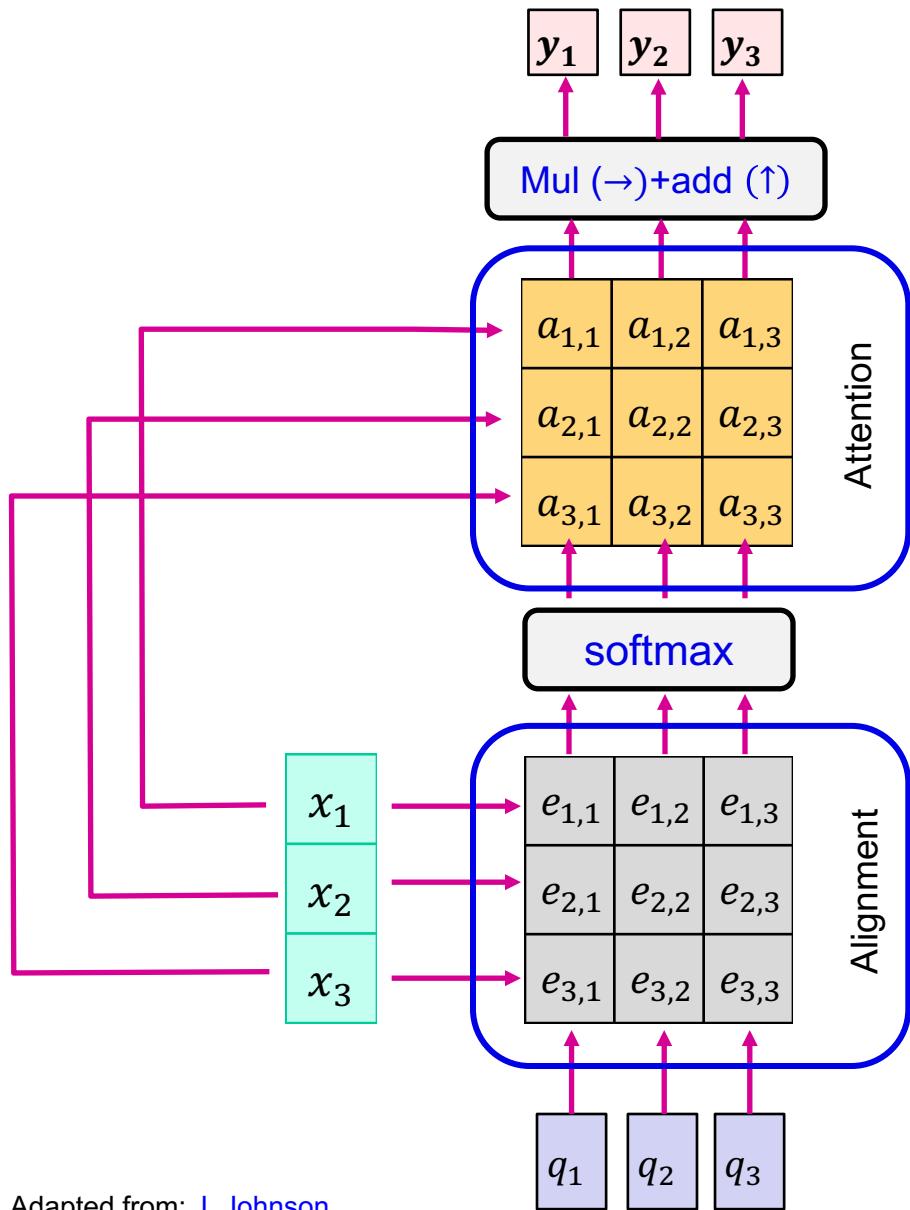
$$A = \text{softmax}(E, \text{dim} = 1) \quad (N \times M)$$

Inputs:

$$\text{Input vectors: } \mathbf{X} = [x_1, \dots]^T \quad (\text{shape: } N \times D)$$

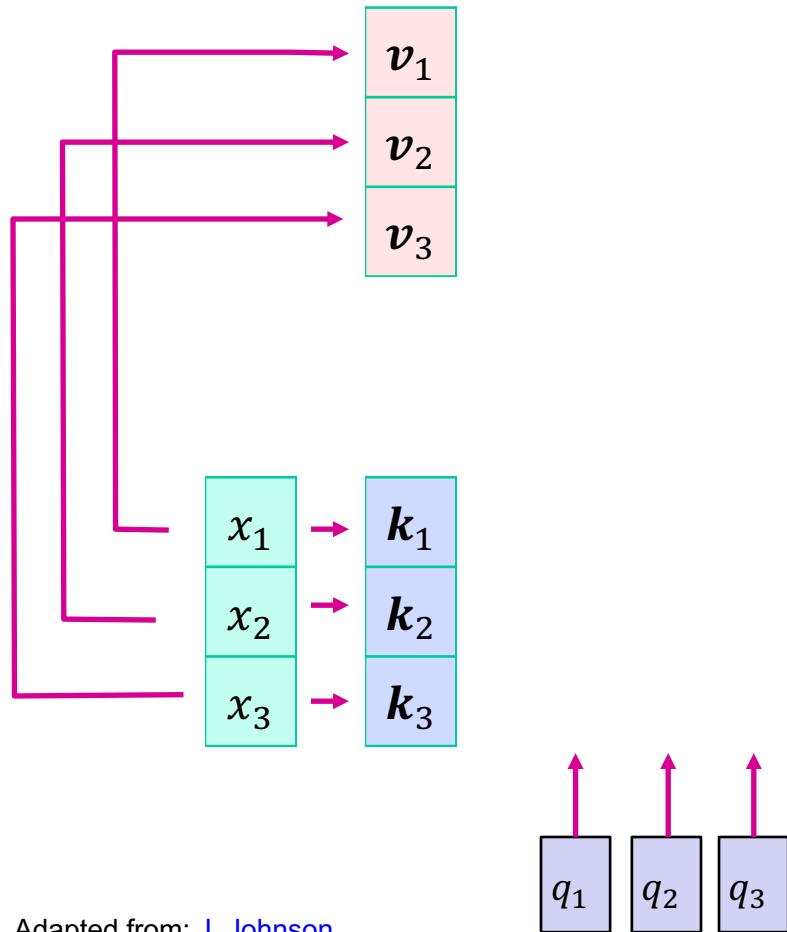
$$\text{Query: } \mathbf{Q} = [q_1, q_2, \dots]^T \quad (\text{shape: } M \times D)$$

General attention layer



- Notice, that the input vectors are used for both, the alignment as well as the attention calculations.
- In this scheme the attention model is fixed and not learned.
- We can add more expressivity to the layer by adding a different FC layer before each of the two steps.

General attention layer: adding keys/values



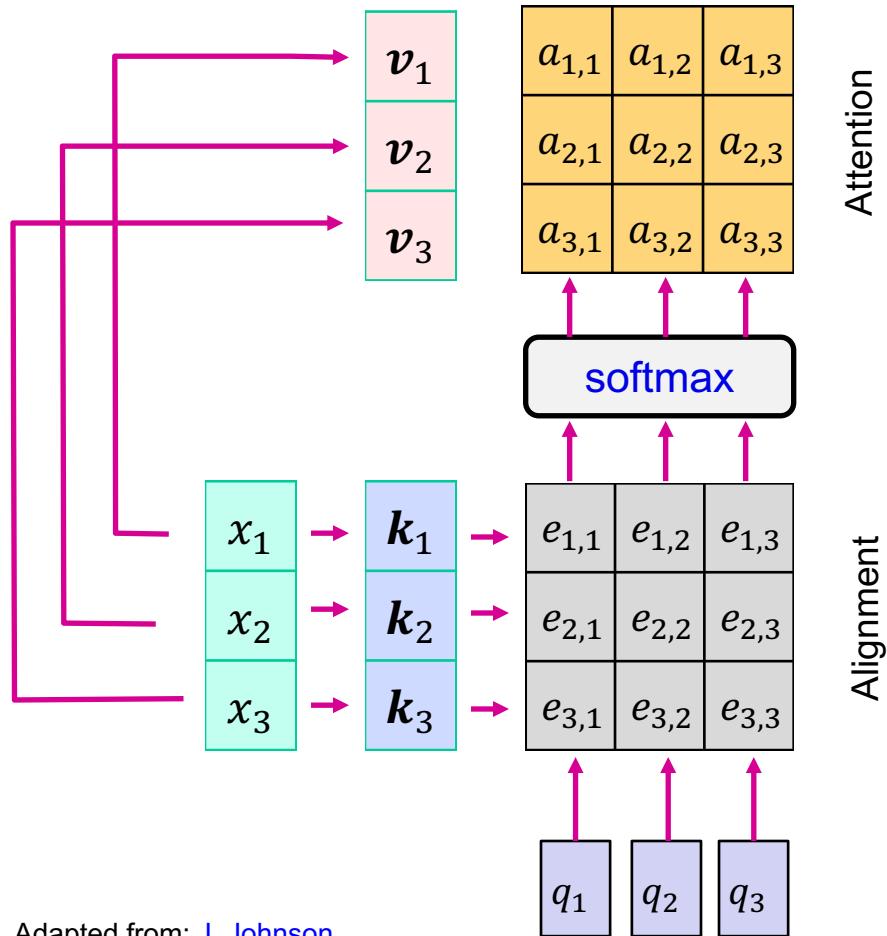
Operations:

Key vectors: $k_i = W_k x_i$ (shape: D_k)

Value vectors: $v_i = W_v x_i$ (shape: D_v)

Both W_k and W_v are learned in the training phase

General attention layer



Operations:

Key vectors: $\mathbf{k}_i = W_k \mathbf{x}_i$ (shape: D_k)

Value vectors: $\mathbf{v}_i = W_v \mathbf{x}_i$ (shape: D_v)

Alignments: $e_{i,j} = \mathbf{k}_i \cdot \mathbf{q}_j / \sqrt{D}$ ($N \times M$)

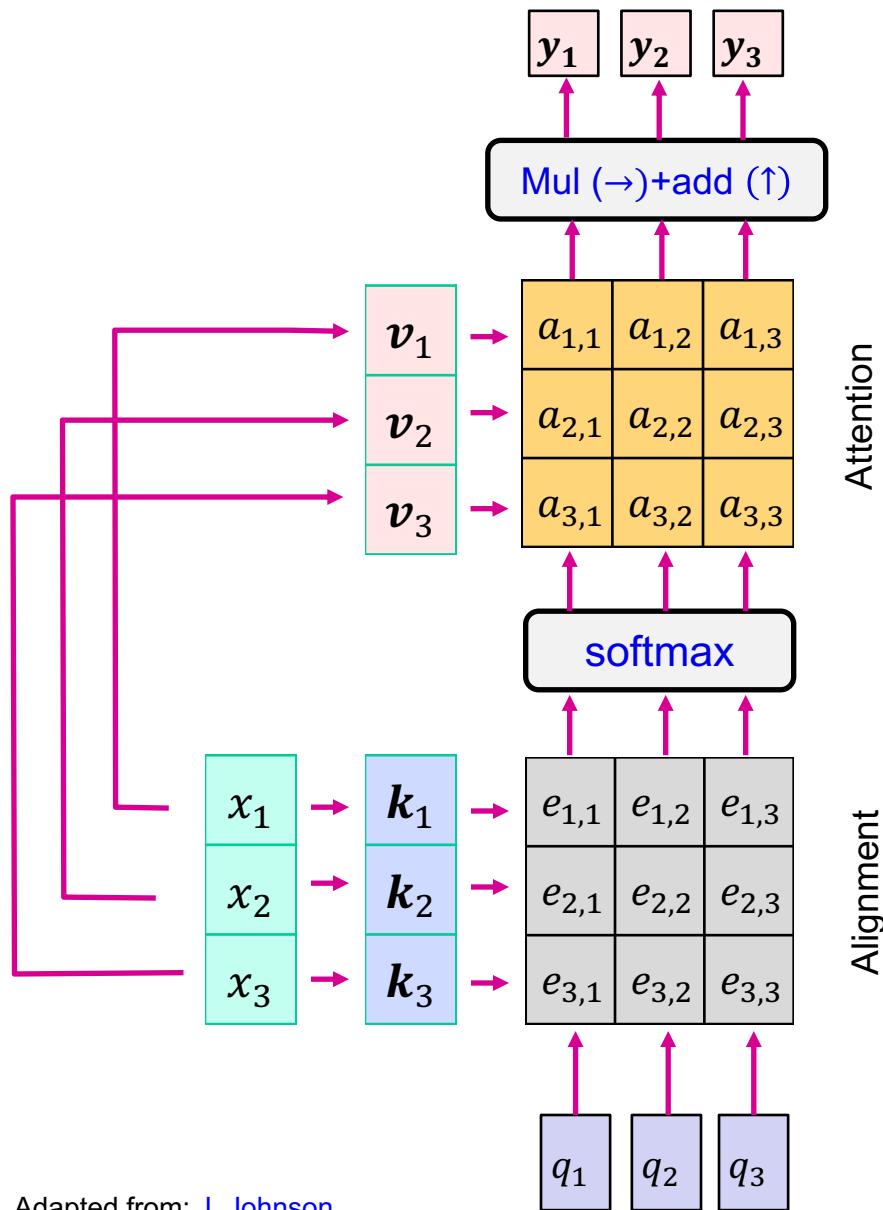
Attention: $a_{*,j} = \text{softmax}(e_{*,j})$ ($N \times M$)

Inputs:

Input vectors: $\mathbf{X} = [x_1, \dots]^T$ (shape: $N \times D$)

Query: $\mathbf{Q} = [q_1, q_2, \dots]^T$ (shape: $M \times D_k$)

General attention layer



Outputs:

$$y_j = \sum_i a_{i,j} v_i ; Y \quad (M \times D_v)$$

Operations:

Key vectors: $k_i = W_k x_i$ (shape: D_k)

Value vectors: $v_i = W_v x_i$ (shape: D_v)

Alignments: $e_{i,j} = k_i \cdot q_j / \sqrt{D}$ ($N \times M$)

Attention: $a_{*,j} = \text{softmax}(e_{*,j})$ ($N \times M$)

Inputs:

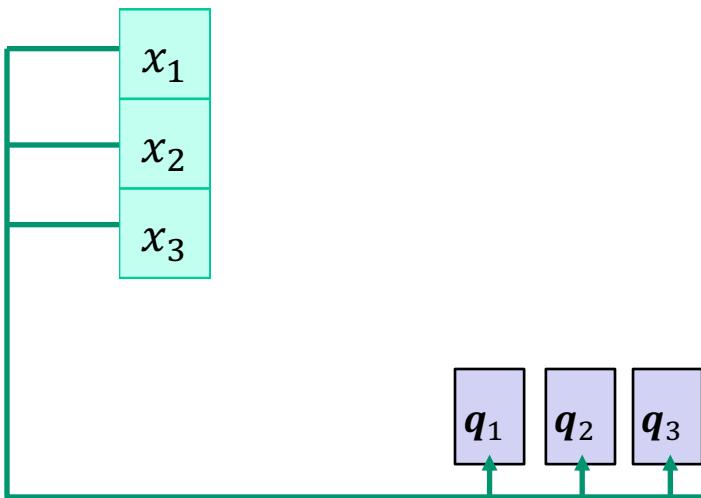
Input vectors: $X = [x_1, \dots]^T$ (shape: $N \times D$)

Query: $Q = [q_1, q_2, \dots]^T$ (shape: $M \times D_k$)

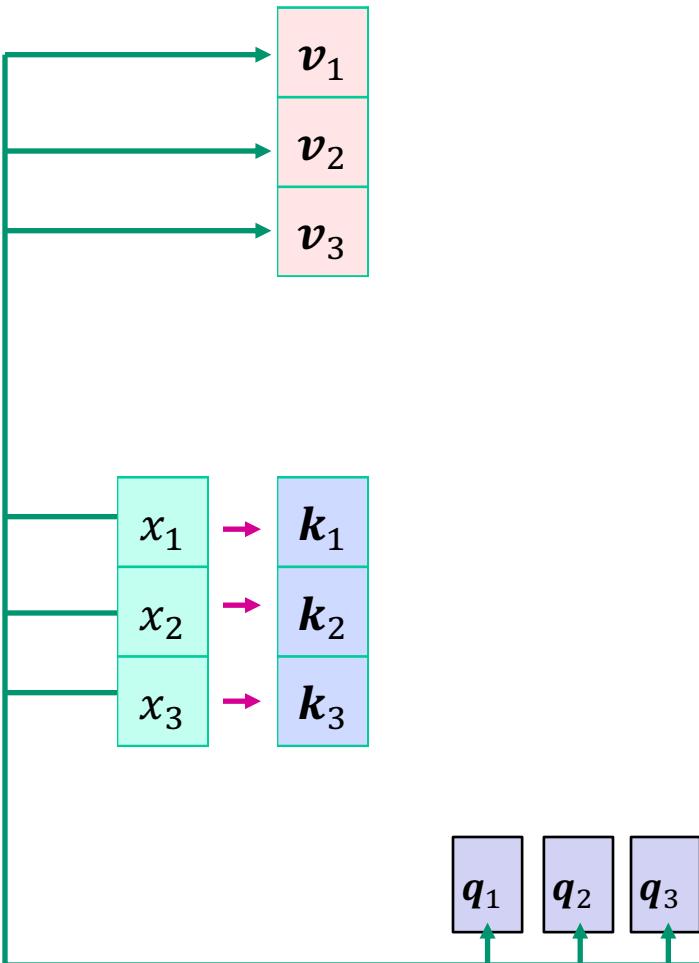
General attention layer: The self-attention

Self attention layer

- We calculate also the **query** vectors from the input vectors
- This defines a "**self-attention**" layer.
- Query vectors: $q = W_q x$



Self attention layer



Operations:

Query vectors: $\mathbf{q} = \mathbf{W}_q \mathbf{x}$ ($N \times D_k$)

Key vectors: $\mathbf{k} = \mathbf{W}_k \mathbf{x}$ ($N \times D_k$)

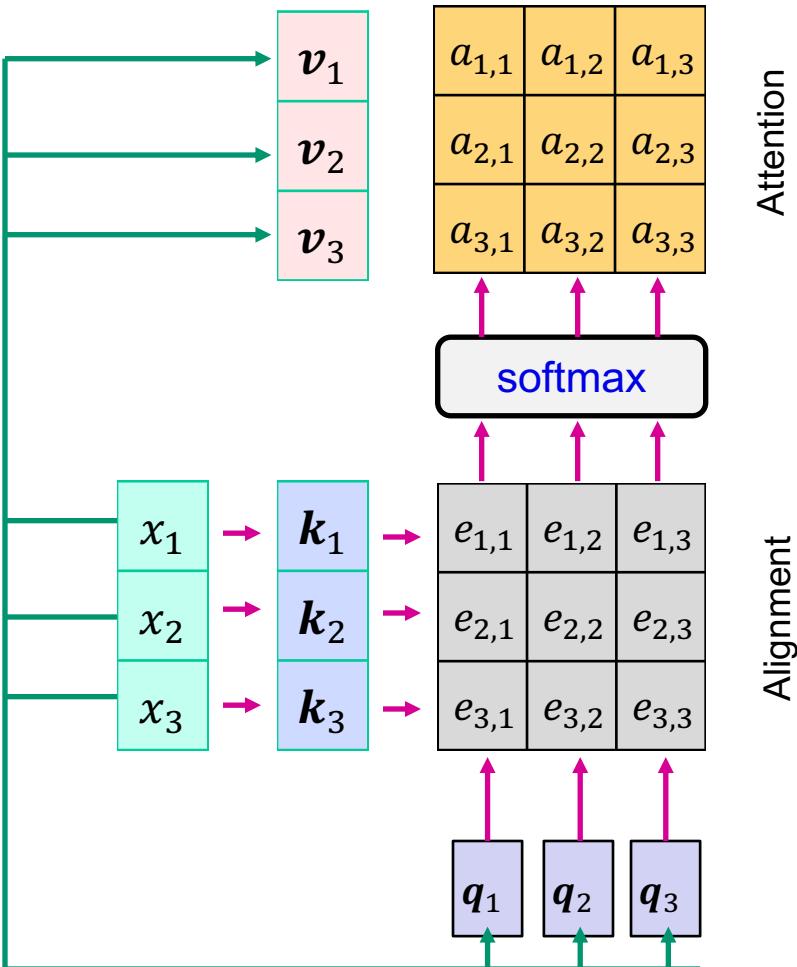
Value vectors: $\mathbf{v} = \mathbf{W}_v \mathbf{x}$ ($N \times D$)

All matrices \mathbf{W}_q and \mathbf{W}_k and \mathbf{W}_v are learned in the training phase

Inputs:

Input vectors: \mathbf{x} (*shape: $N \times D$*)

Self attention layer



Operations:

Query vectors: $q = W_q x$ ($N \times D_k$)

Key vectors: $k = W_k x$ ($N \times D_k$)

Value vectors: $v = W_v x$ ($N \times D$)

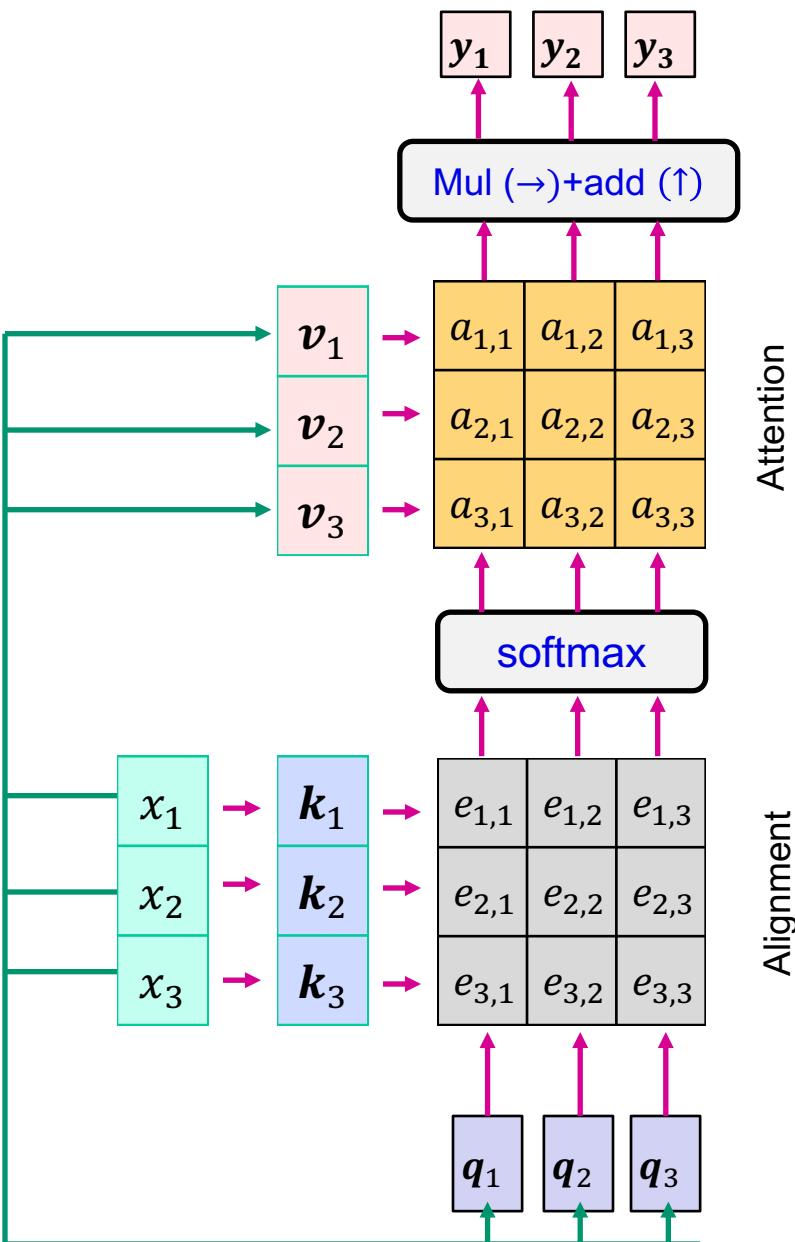
Alignments: $e_{i,j} = k_i \cdot q_j / \sqrt{D}$

Attention: $a_{*,j} = \text{softmax}(e_{*,j})$

Inputs:

Input vectors: x (*shape*: $N \times D$)

Self attention layer



Outputs:

Multiple Context vectors: \mathbf{y} ($N \times D$)

Operations:

Query vectors: $\mathbf{q} = W_q \mathbf{x}$ ($N \times D_k$)

Key vectors: $\mathbf{k} = W_k \mathbf{x}$ ($N \times D_k$)

Value vectors: $\mathbf{v} = W_v \mathbf{x}$ ($N \times D$)

Alignments: $e_{i,j} = \mathbf{k}_i \cdot \mathbf{q}_j / \sqrt{D}$

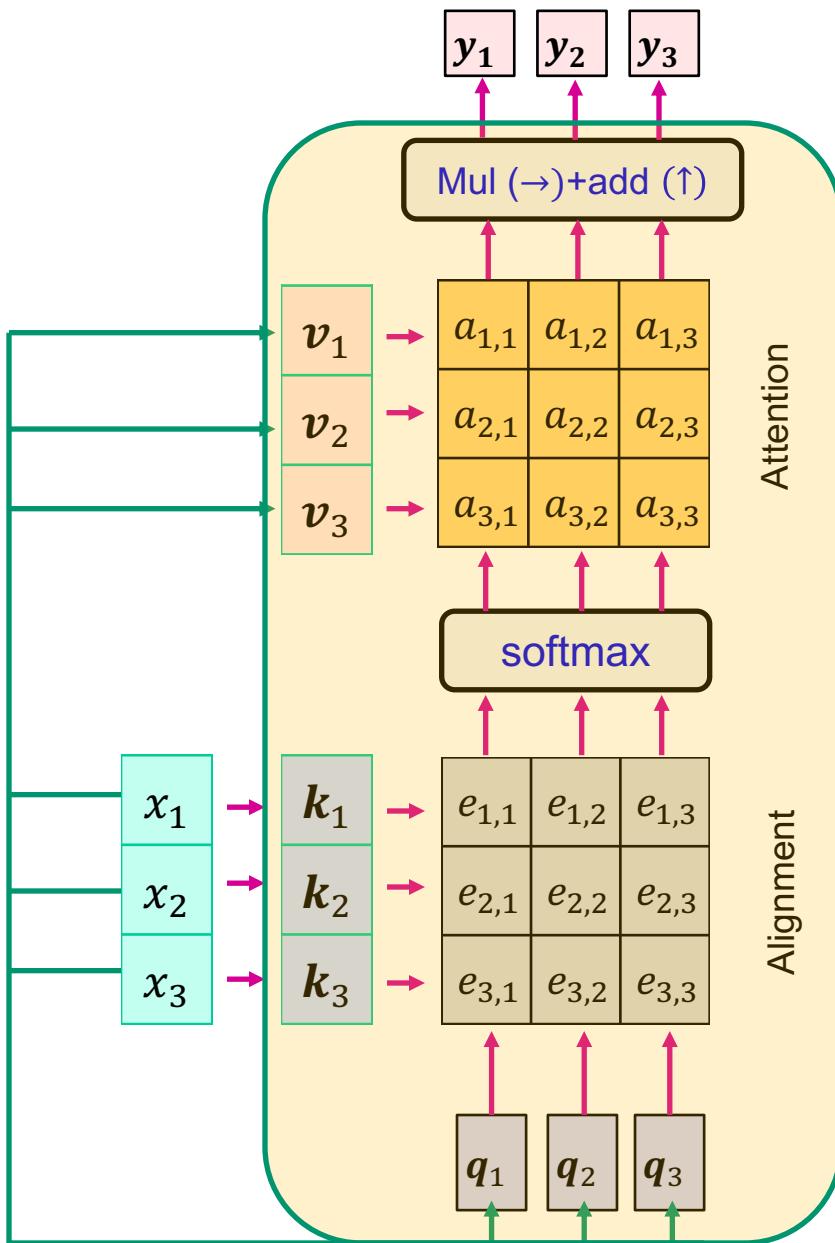
Attention: $a_{*,j} = \text{softmax}(e_{*,j})$

Output: $\mathbf{y}_j = \sum_i a_{i,j} \mathbf{v}_i$

Inputs:

Input vectors: \mathbf{x} (shape: $N \times D$)

Self attention layer



Input:

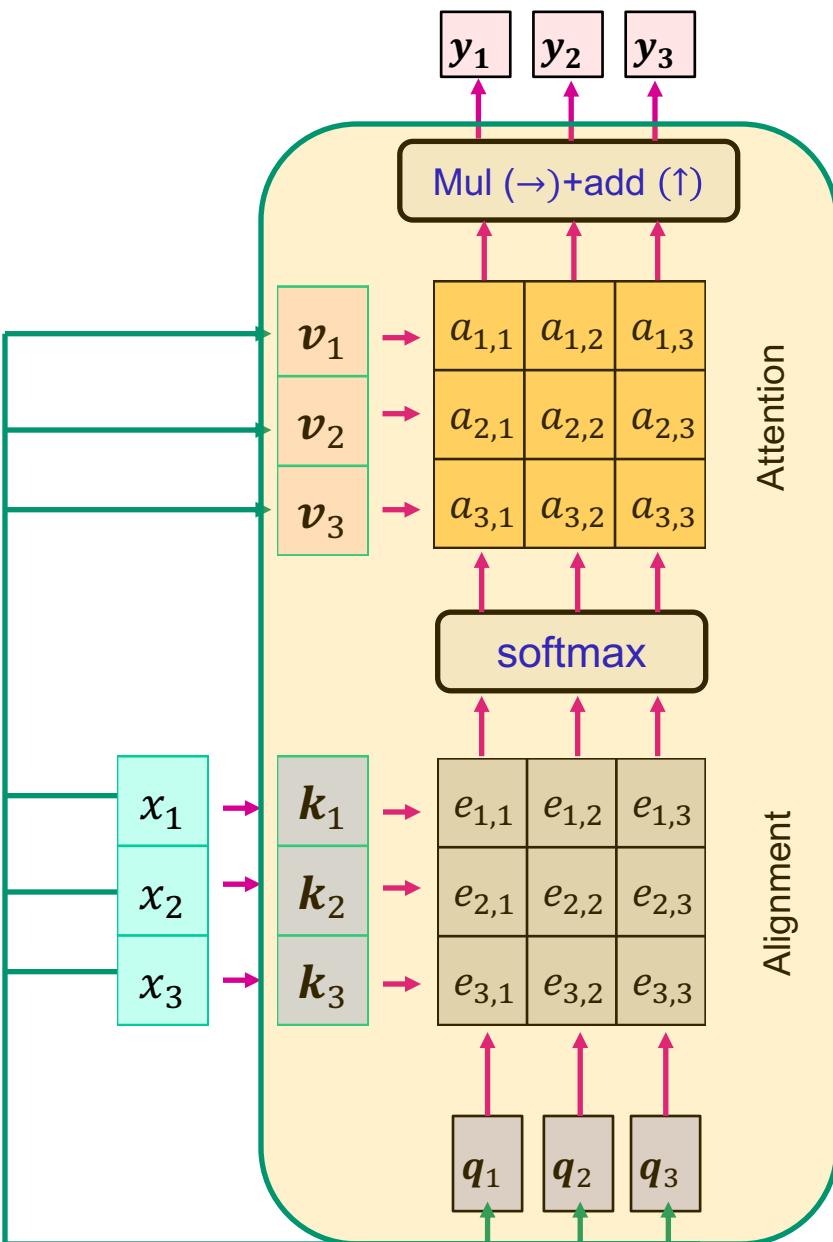
Input vectors: $N \times D$ inputs

output:

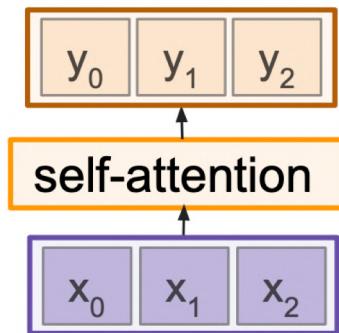
output vectors: $N \times D$ inputs

Learned: W_q, W_k, W_v

Self attention layer



Self-attention layer:
is used to capture context *within the sequence*.



Input:
Input vectors: $N \times D$ inputs

output:
output vectors: $N \times D$ inputs

Self attention layer

- 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”

General attention layer: Matrix form

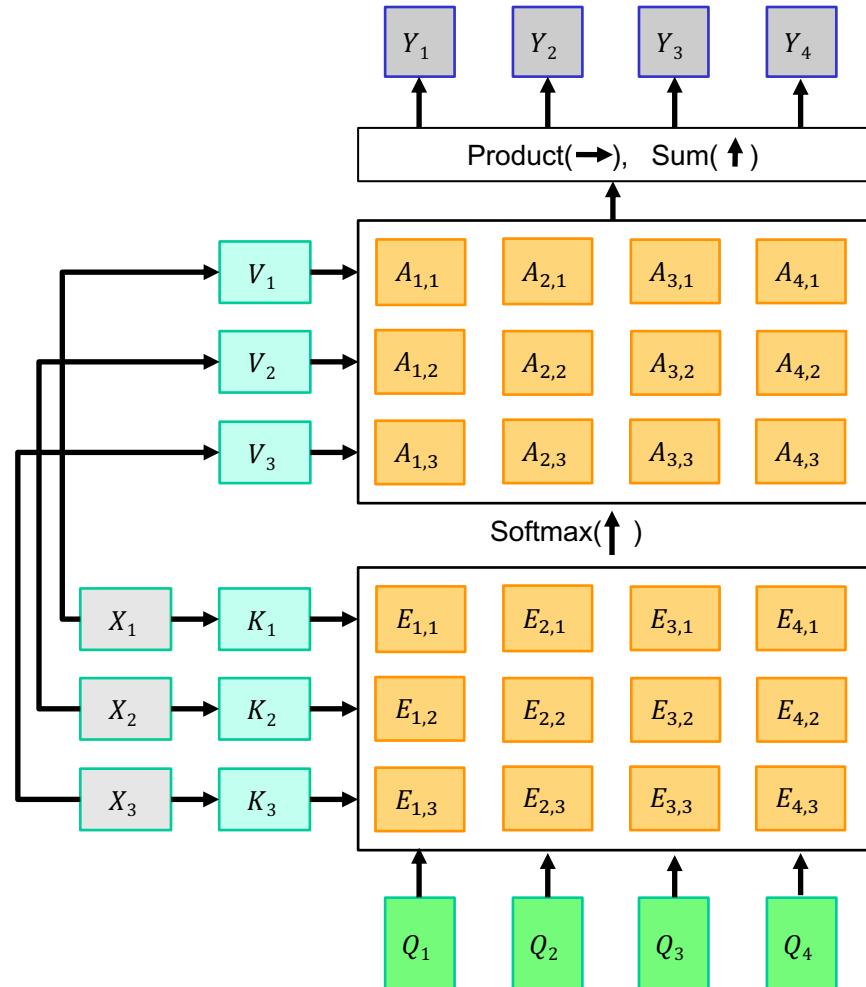
- Key vectors: $K = XW_K$ ($N \times D_k$)
- Value Vectors: $V = XW_V$ ($N \times D_v$)
- Query vectors ($M \times D_k$)
- Similarities: *scaled dot-product attention*

$$E_{i,j} = \frac{(Q_i \cdot K_j)}{\sqrt{D}} \text{ or } E = QK^T / \sqrt{D}$$

(D is the dimensionality of the keys)

- Attn. weights: $A = \text{softmax}(E, \text{dim} = 1)$
- Output vectors:

$$Y_i = \sum_j A_{i,j} V_j \text{ or } Y = AV \quad (M \times D_v)$$



Self attention layer: Matrix form

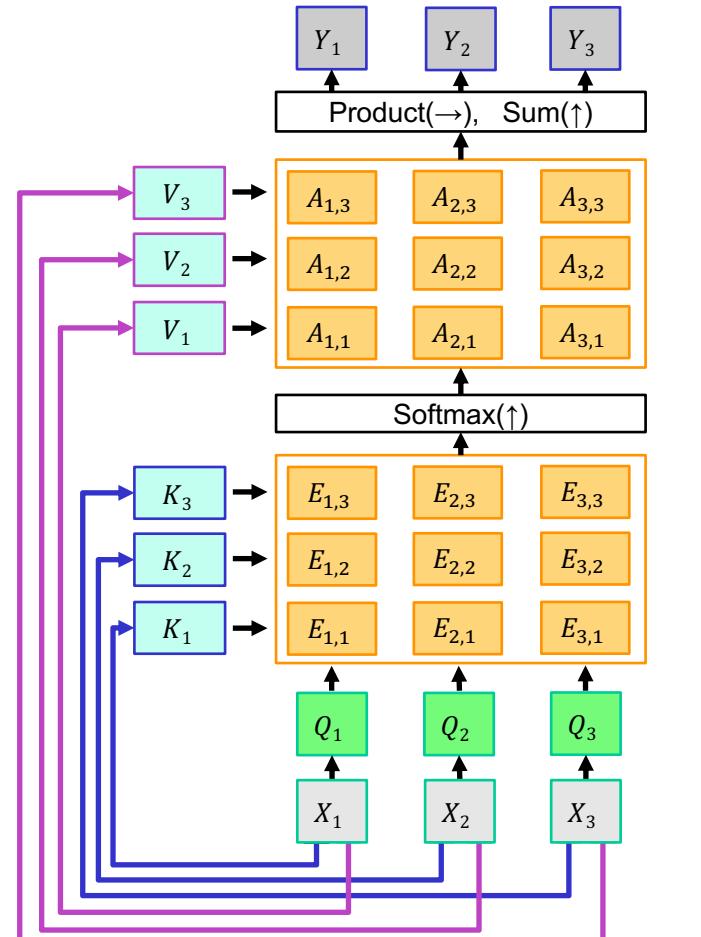
- Query vectors: $Q = XW_Q$
- Key vectors: $K = XW_K$
- Value vectors: $V = XW_V$
- Similarities: *scaled dot-product attention*

$$E_{i,j} = \frac{(Q_i \cdot K_j)}{\sqrt{D}} \text{ or } E = QK^T / \sqrt{D}$$

(D is the dimensionality of the keys)

- Attn. weights: $A = \text{softmax}(E, \text{dim} = 1)$
- Output vectors:

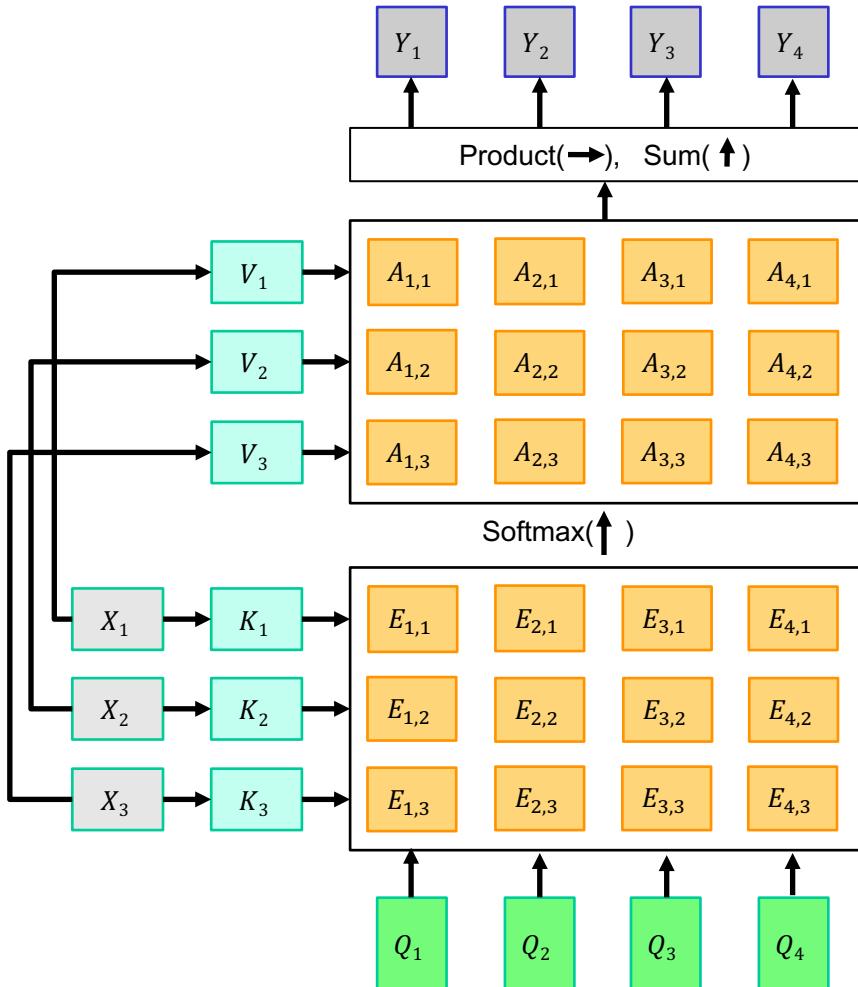
$$Y_i = \sum_j A_{i,j} V_j \text{ or } Y = AV \quad (\text{NxD})$$



One query per input vector

Attention layer:

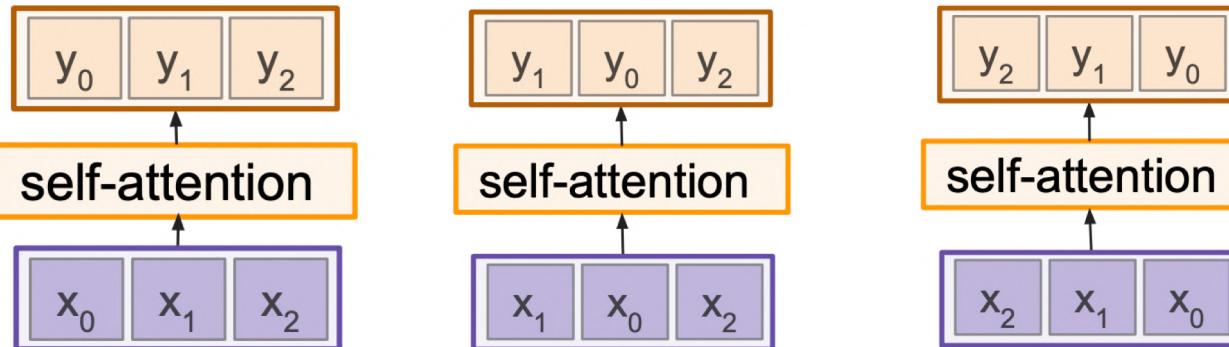
- How does permuting the order of the *queries* change the output?
- How does changing the order of the *keys/values* change the output?



Attention layer: Positional Encoding

Permutation invariant:

- The self attention layer attends over a set of inputs
- Attention operation is **permutation invariant**: it doesn't care about ordering of the inputs
- **Problem:** how can we encode ordered sequences like language or spatially ordered image features?



Attention layer: Positional Encoding

- **Positional encoding** provide the network information about the position of each input.
- Concatenate *positional encoding* p_j to each input vector x_j
- We use a function $\text{pos}: \mathbb{N} \rightarrow \mathbb{R}^d$ to encode the position j of the vector x_j into a d -dimensional vector, so $p_j = \text{pos}(j)$
- Requirements for $\text{pos}(j)$:
 - It should output a **unique** encoding for each position
 - **Distance** between any two time-steps should be consistent across sequences with different lengths.
 - Should generalize to **longer** sentences without any efforts.

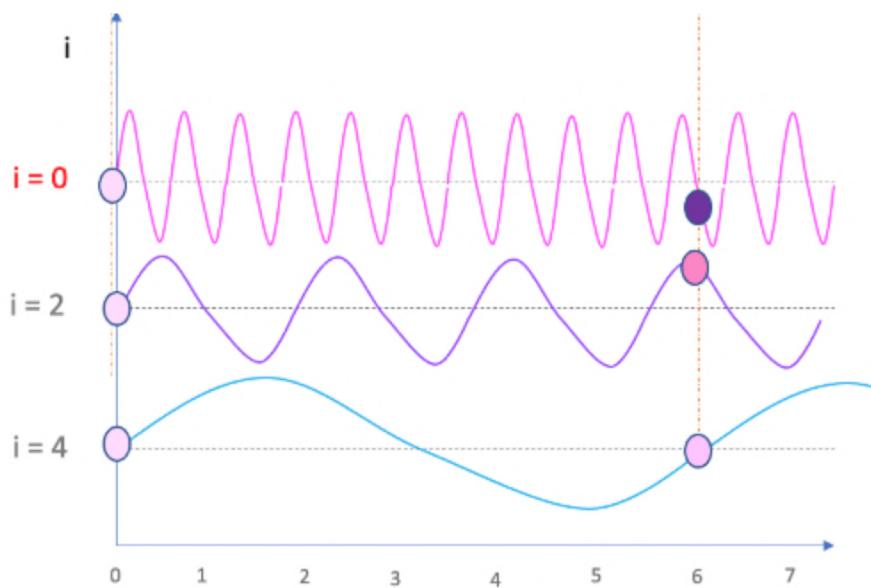
Attention layer: Positional Encoding

- Common use for $\text{pos}(j)$: sin+cos waves:

$$\begin{bmatrix} \sin(\omega_1 \cdot j) \\ \cos(\omega_1 \cdot j) \\ \sin(\omega_2 \cdot j) \\ \cos(\omega_2 \cdot j) \\ \vdots \\ \sin(\omega_{d/2} \cdot j) \\ \cos(\omega_{d/2} \cdot j) \end{bmatrix}$$

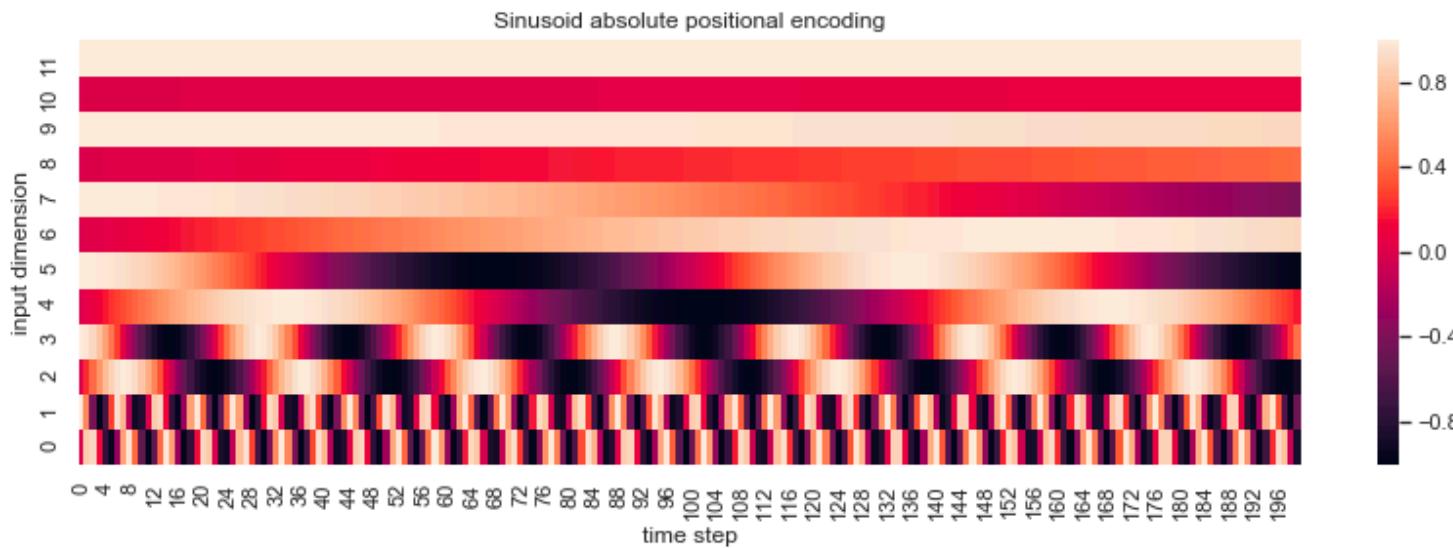
where $\omega_k = \frac{1}{10000^{2k/d}}$

Attention layer: Positional Encoding



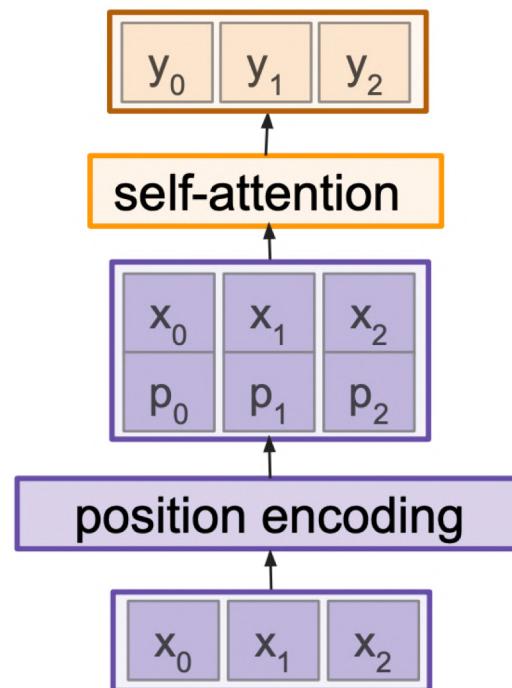
Intuition:

0 :	0	0	0	0	0	1	0	0	0
1 :	0	0	0	1	0	1	0	0	1
2 :	0	0	1	0	0	1	0	1	0
3 :	0	0	1	1	0	1	1	1	1
4 :	0	1	0	0	0	1	1	0	0
5 :	0	1	0	1	0	1	1	0	1
6 :	0	1	1	0	0	1	1	1	0
7 :	0	1	1	1	1	1	1	1	1



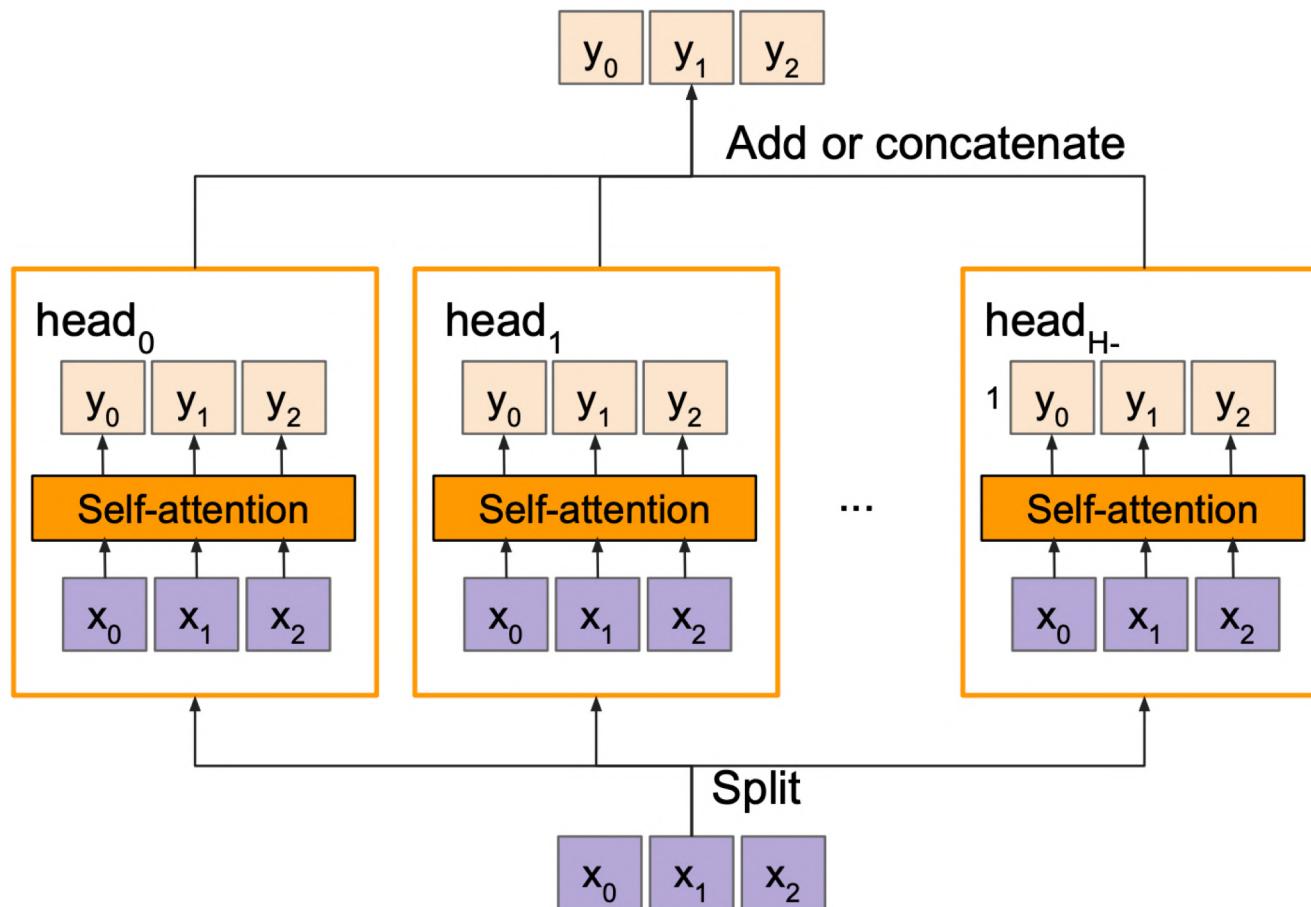
Attention layer: Positional Encoding

Self attention layer with positional encoding:

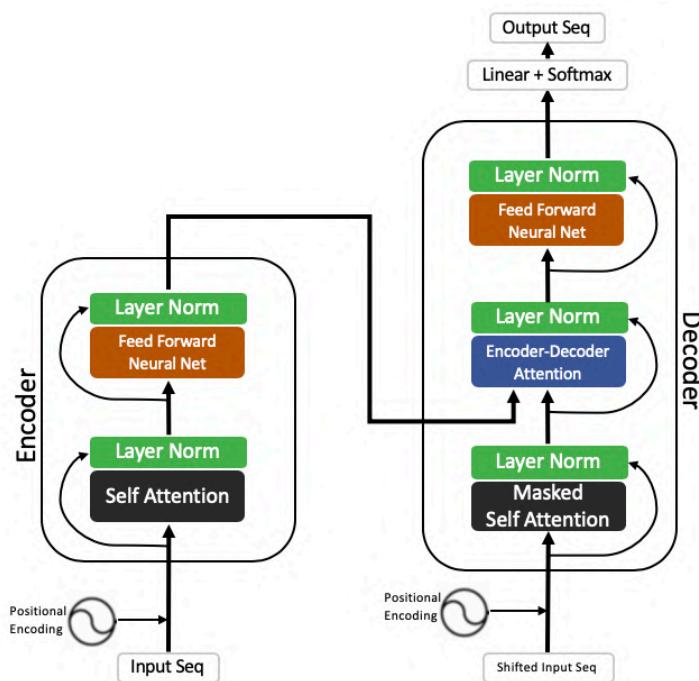


Multi-head self attention layer

Multi-head self attention layer:



Transformers

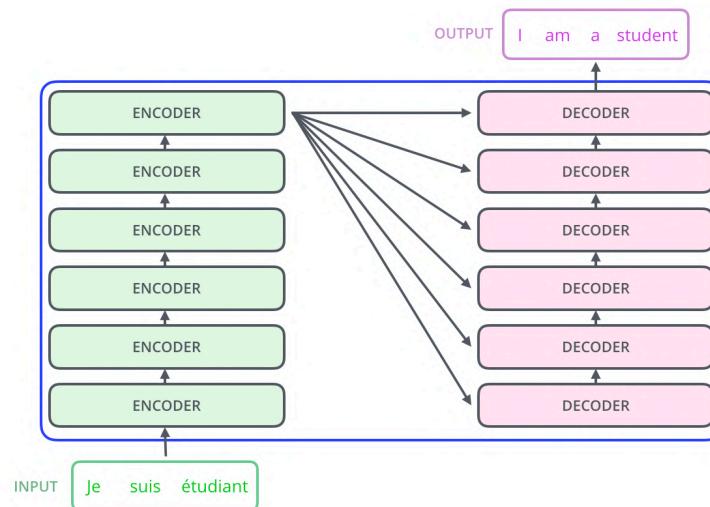


Basic transformer model (review)

- Sequence-to-sequence architecture using only point-wise processing and attention (no recurrent units or convolutions)
- More efficient and parallelizable than recurrent or convolutional architectures, faster to train, better accuracy.

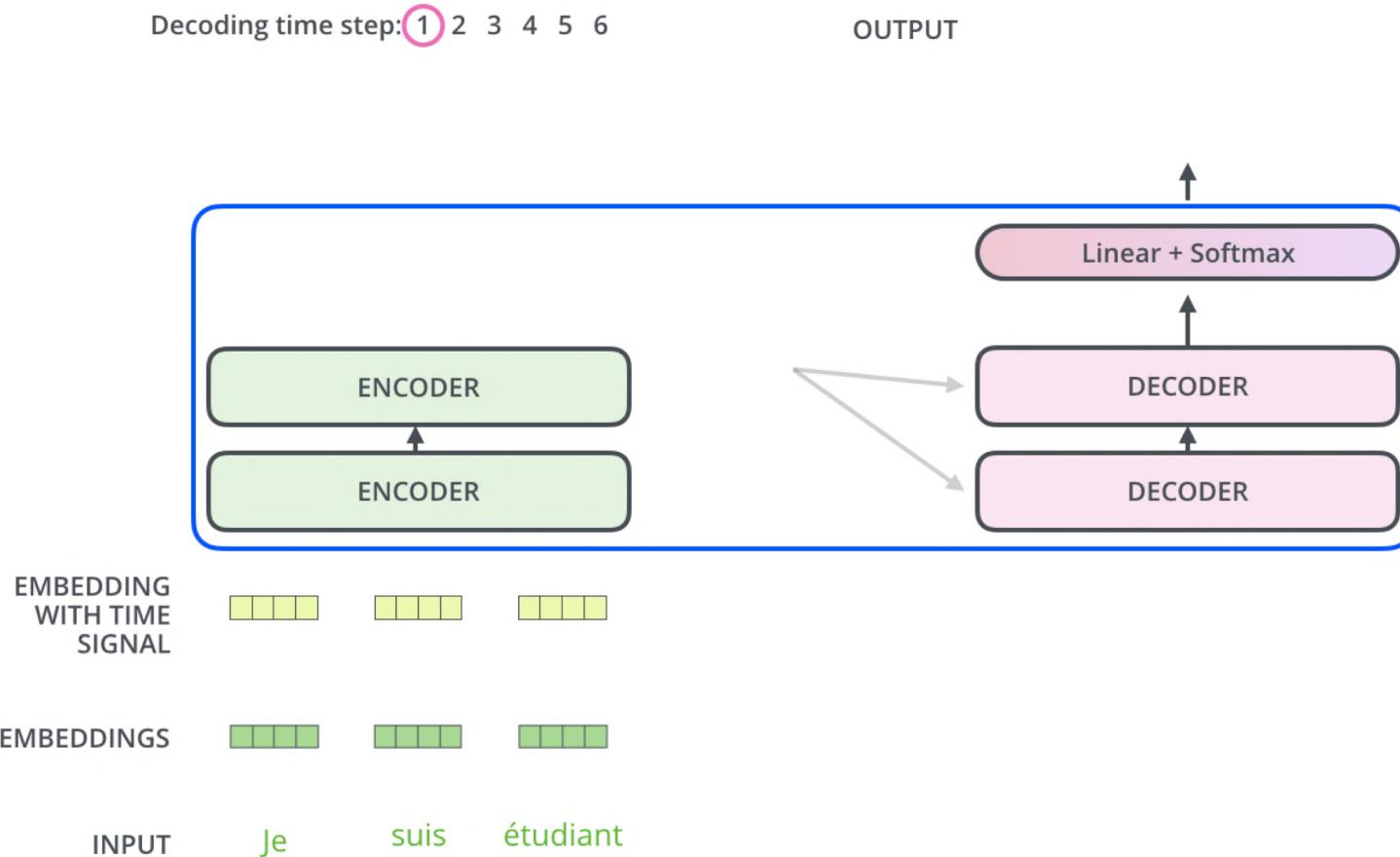
Encoder: receives entire input sequence and outputs encoded sequence of the same length

Decoder: predicts next token conditioned on encoder output and previously predicted tokens



Basic transformer model (review)

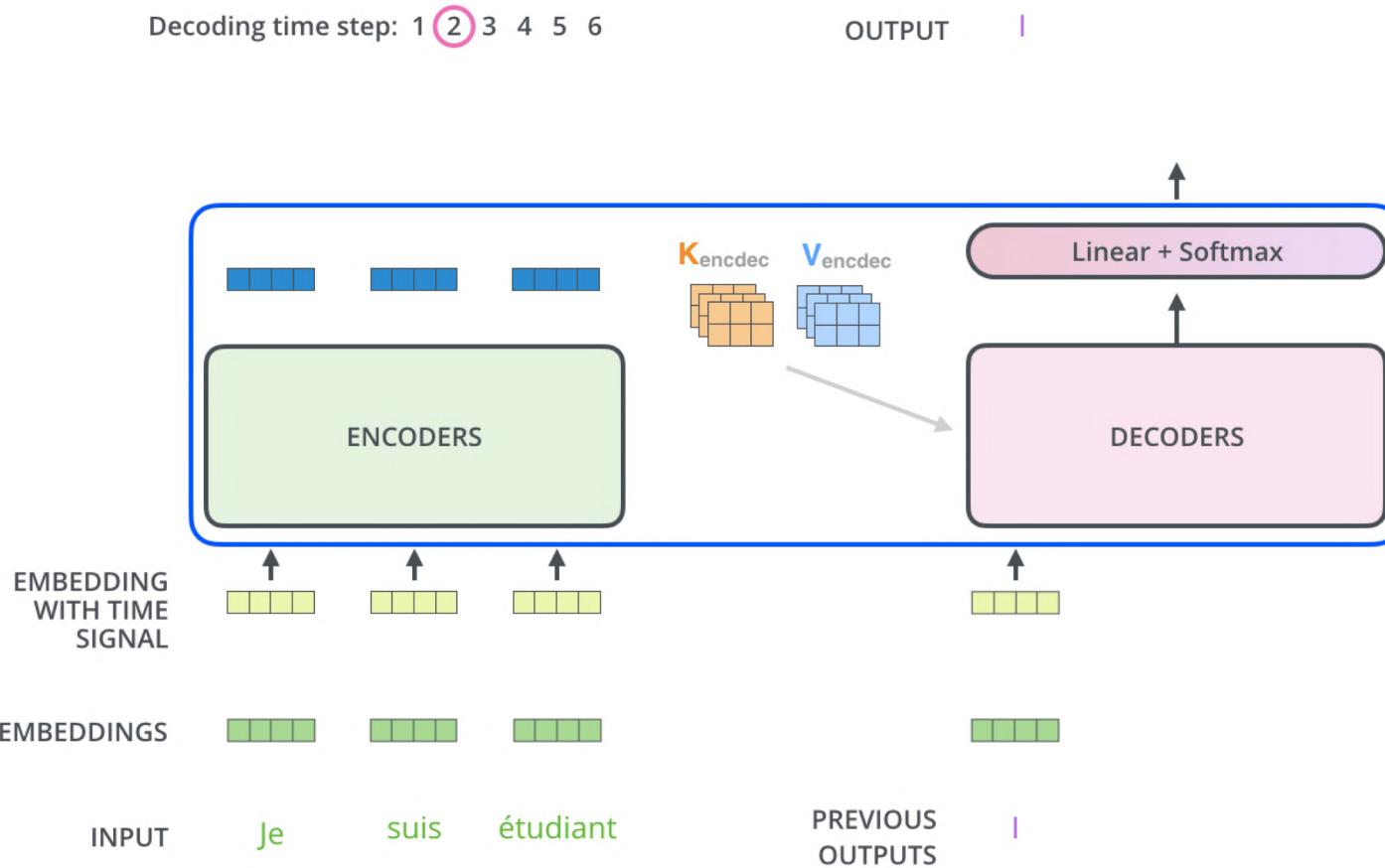
Illustrated transformer - 1st step:



From: [The Illustrated Transformer](#): Jay Alammar

Basic transformer model (review)

Illustrated transformer – next steps:



From: [The Illustrated Transformer](#): Jay Alammar

Image captioning using transformers

Input: image I

Output: $y = y_1, y_2, y_3 \dots, y_T$



Spatial features from
a pretrained CNN

$Z_{1,1}$	$Z_{1,2}$	$Z_{1,3}$
$Z_{2,1}$	$Z_{2,2}$	$Z_{2,3}$
$Z_{3,1}$	$Z_{3,2}$	$Z_{3,3}$

Features
 $H \times W \times D$

Image captioning using transformers

Input: image I

Output: $y = y_1, y_2, y_3 \dots, y_T$

Encoder: $c = T_w(z)$

where z is spatial CNN features and $T_w(z)$ is the transformer encoder

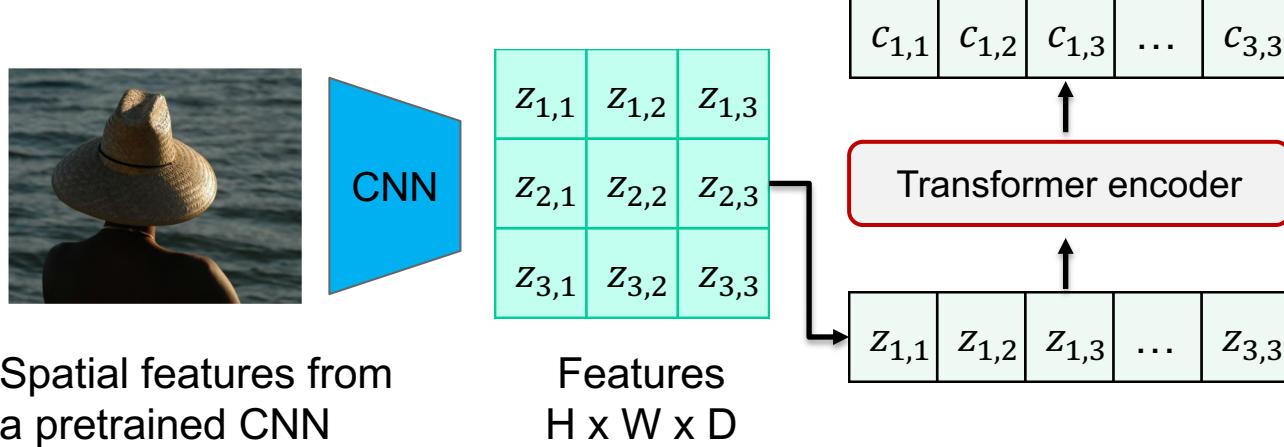


Image captioning using transformers

Input: image I

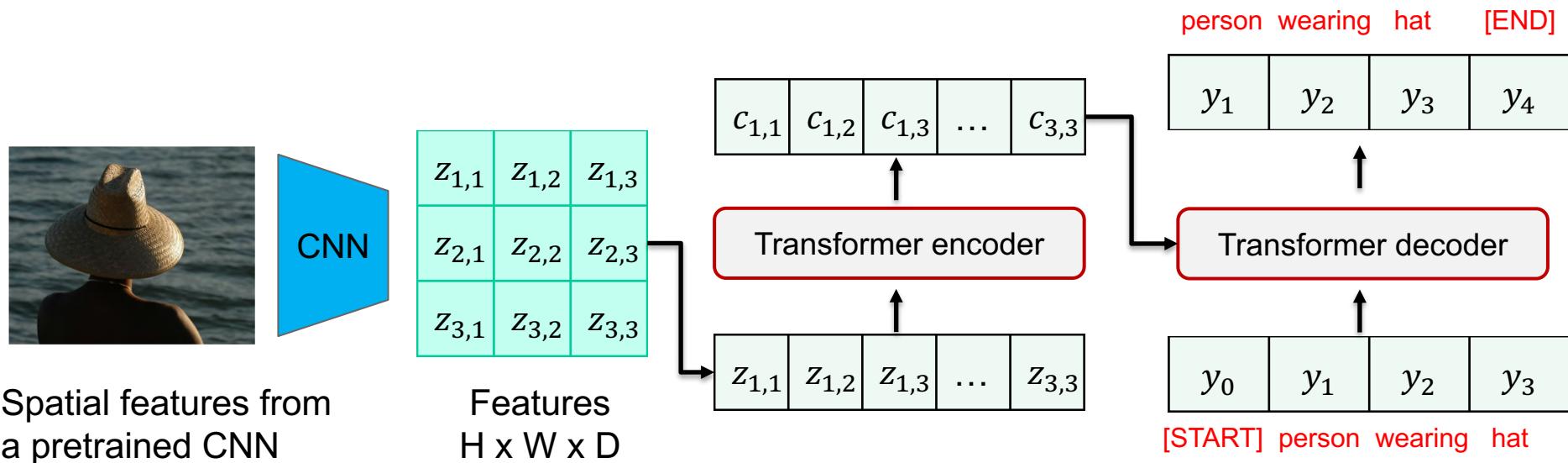
Output: $y = y_1, y_2, y_3 \dots, y_T$

Encoder: $c = T_w(z)$

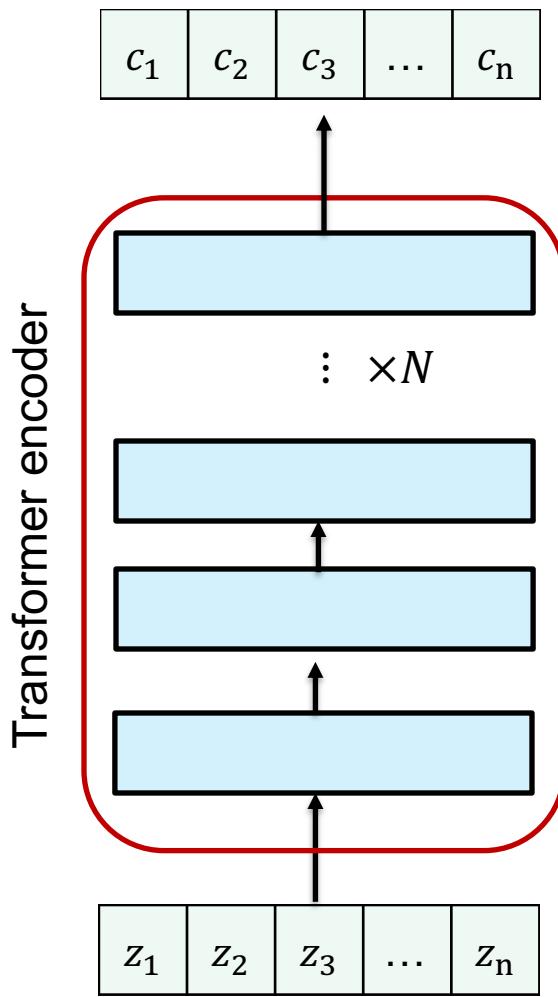
where z is spatial CNN features and $T_w(z)$ is the transformer encoder

Decoder: $y_t = T_D(y_{0:t-1}, c)$

where $T_D(\cdot)$ is the decoder transformer



The Transformer encoder block

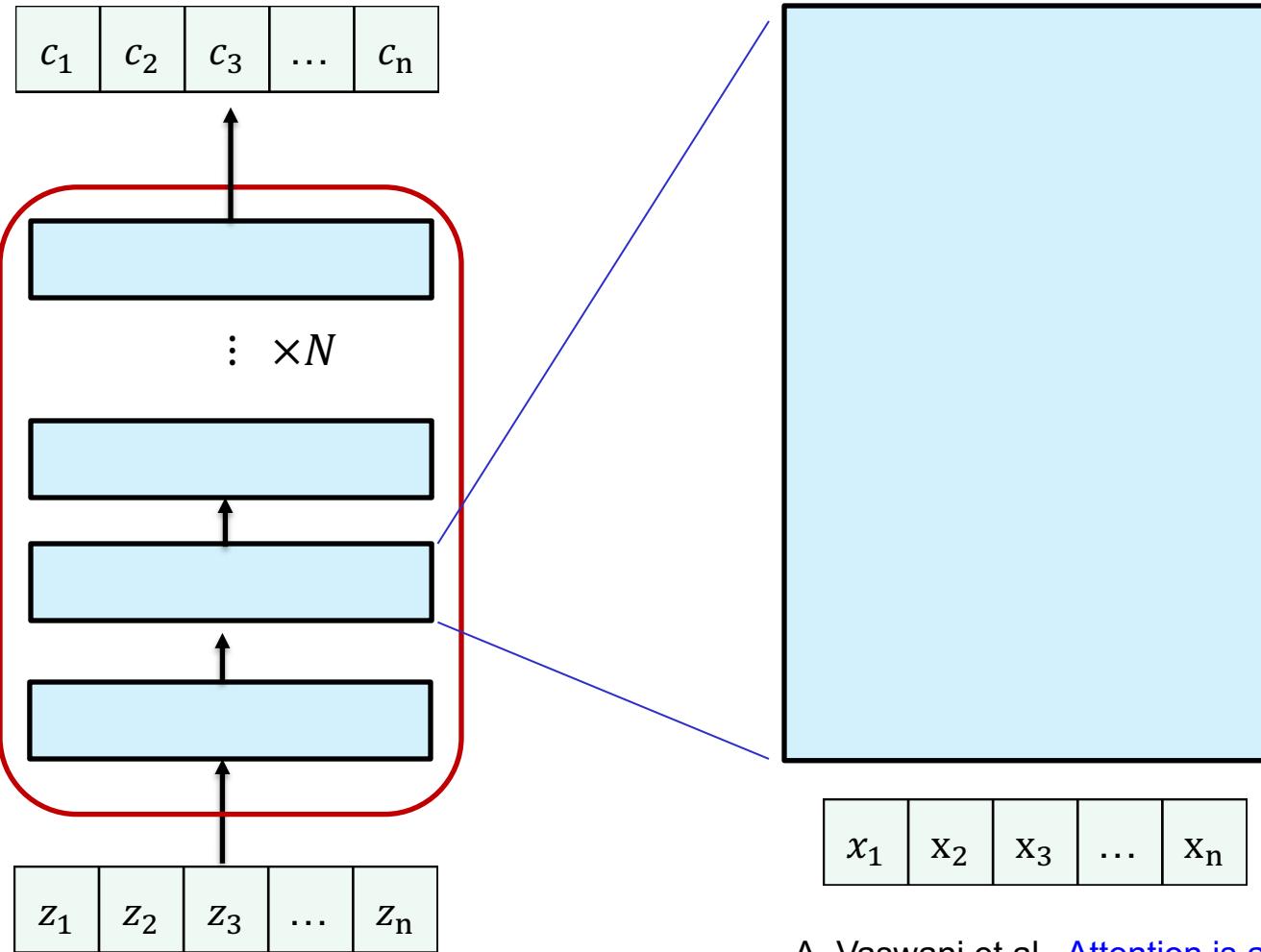


Made up of N encoder blocks

In Vaswani et al. $N=12$, $D_q=512$

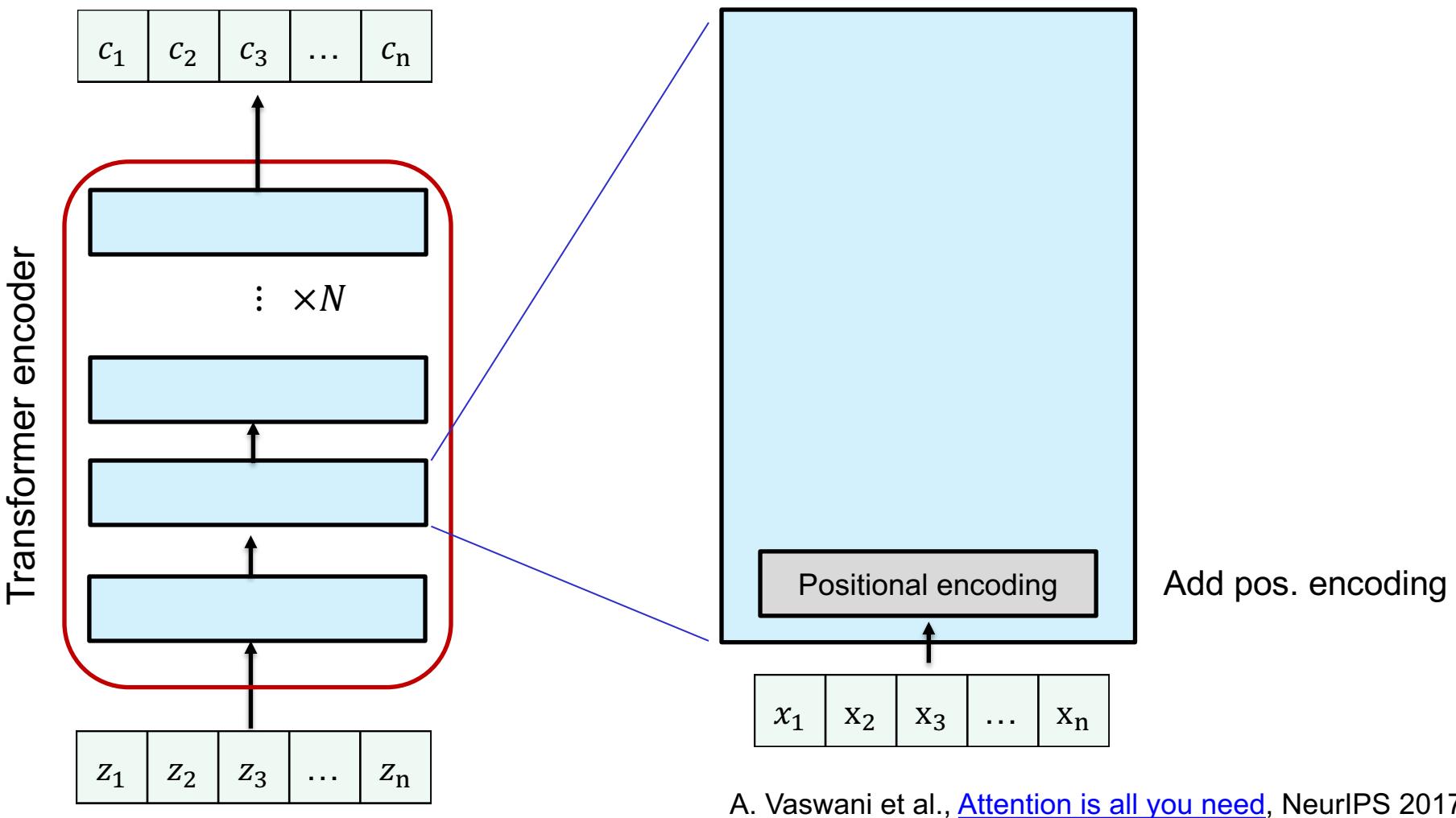
The Transformer encoder block

Transformer encoder

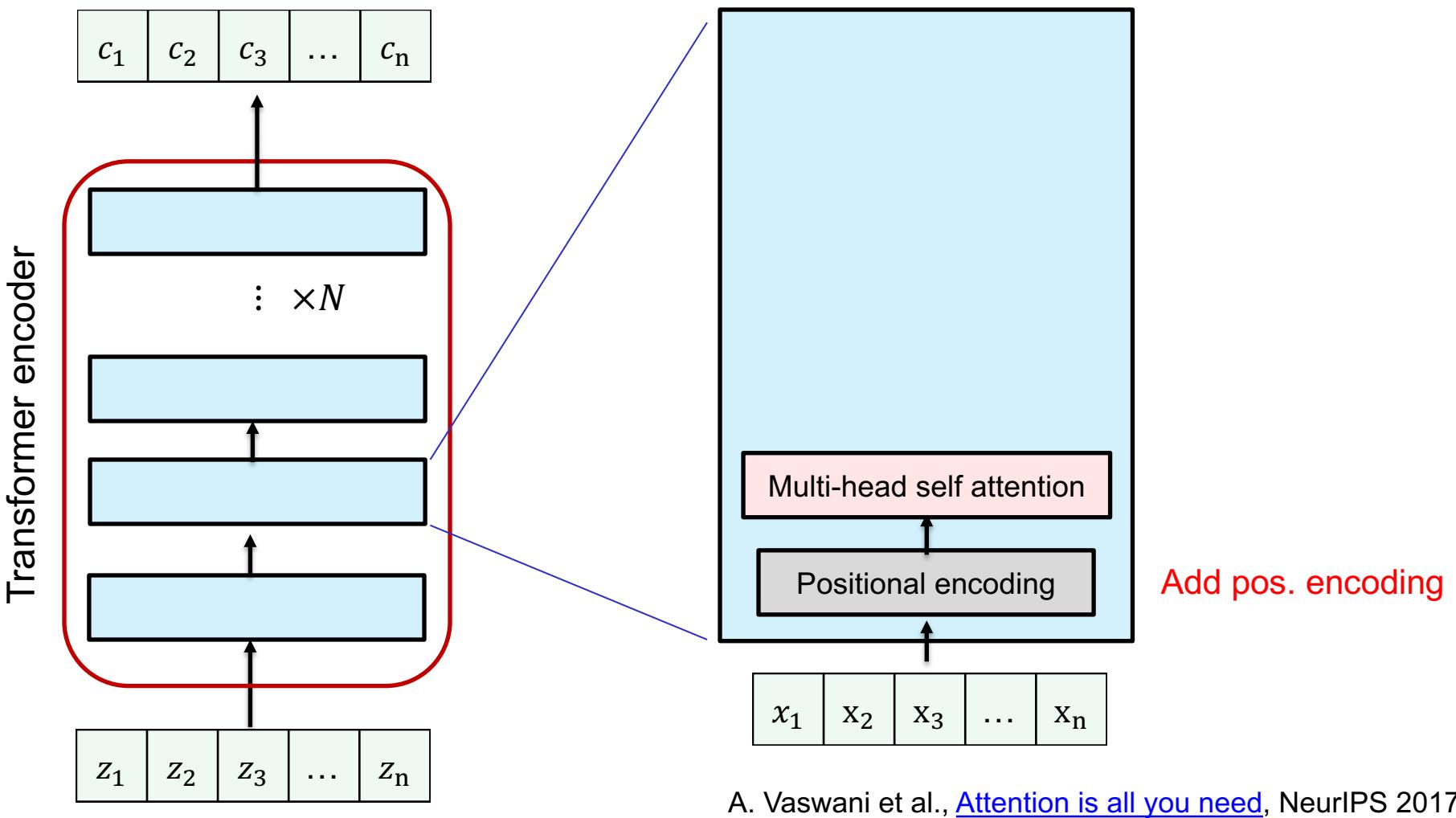


A. Vaswani et al., [Attention is all you need](#), NeurIPS 2017

The Transformer encoder block

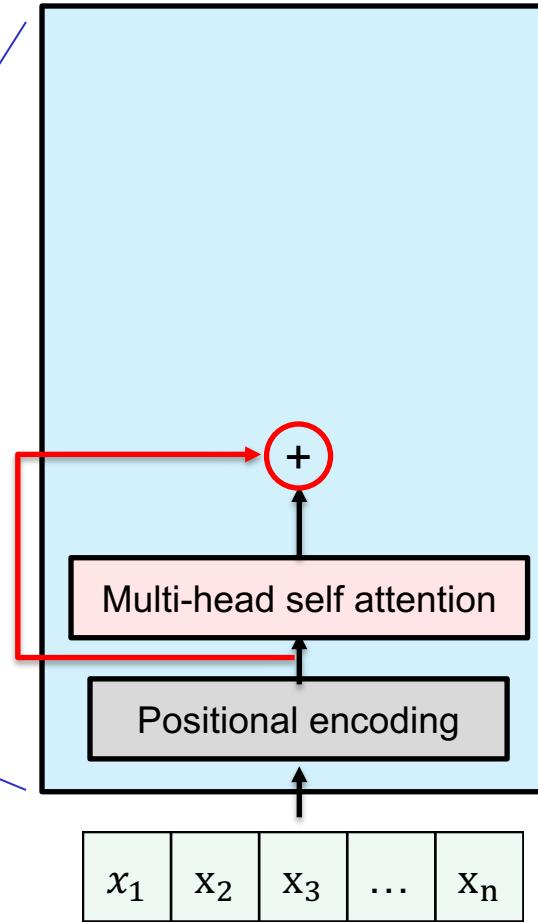
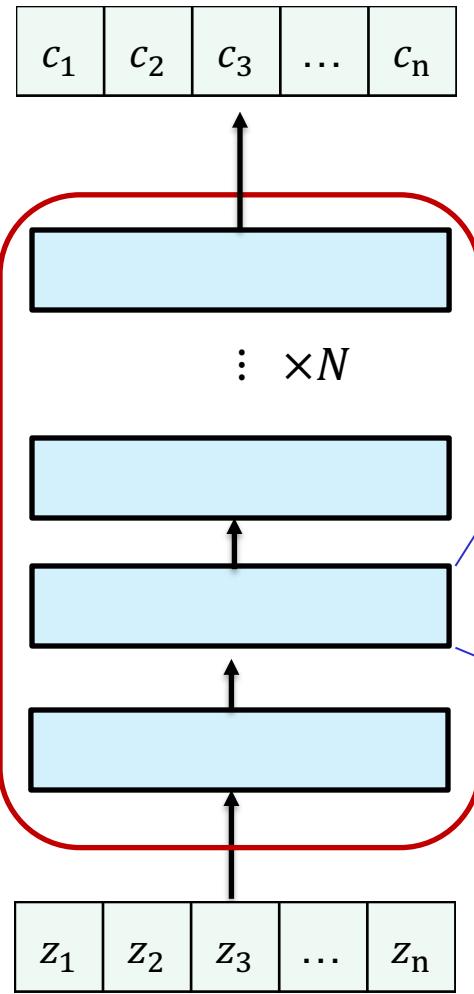


The Transformer encoder block



The Transformer encoder block

Transformer encoder

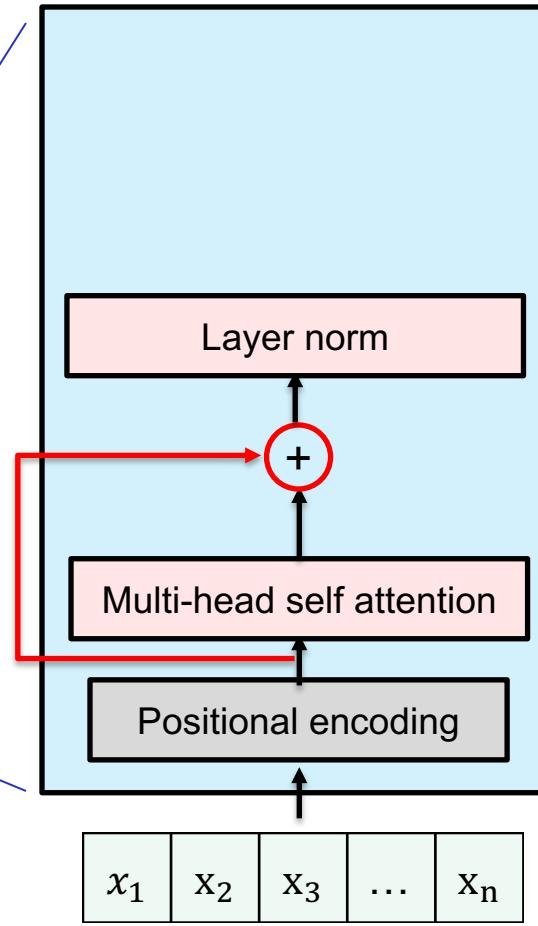
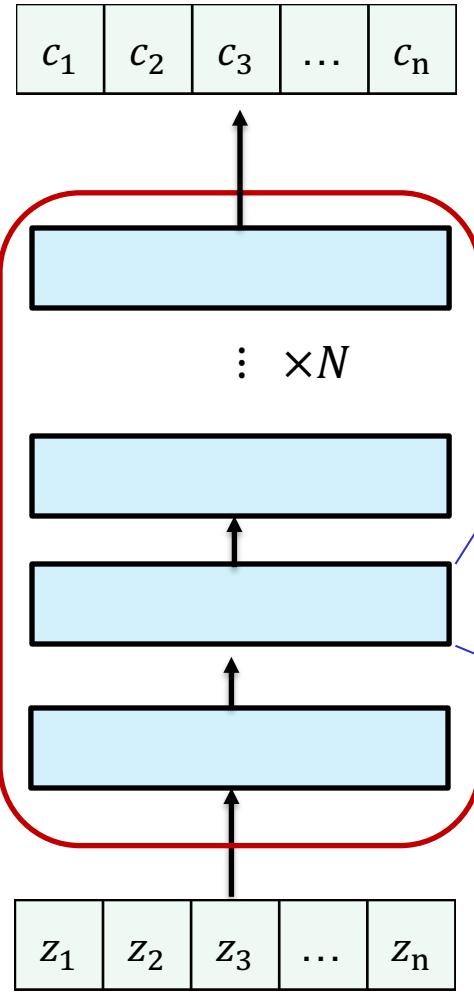


Residual connection

Add pos. encoding

The Transformer encoder block

Transformer encoder



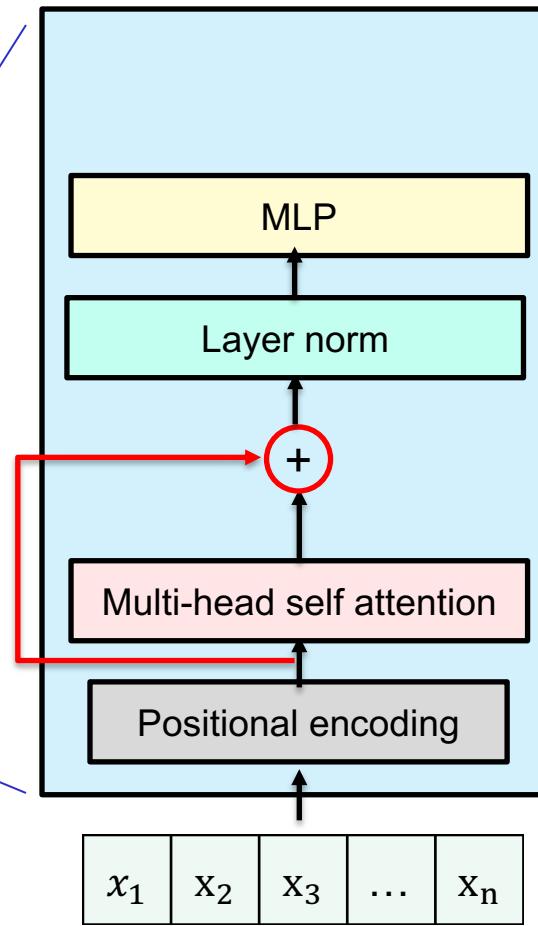
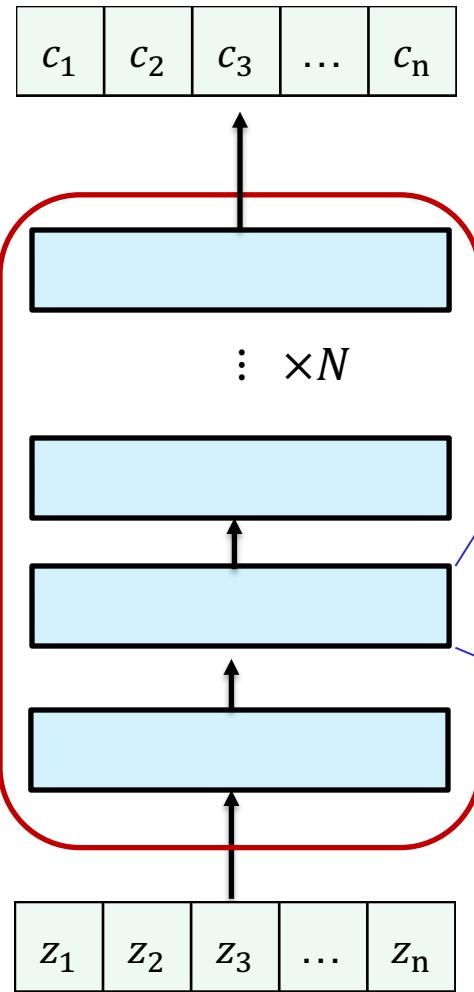
Layer norm over each vector individually

Residual connection

Add pos. encoding

The Transformer encoder block

Transformer encoder



MLP over each vector individually

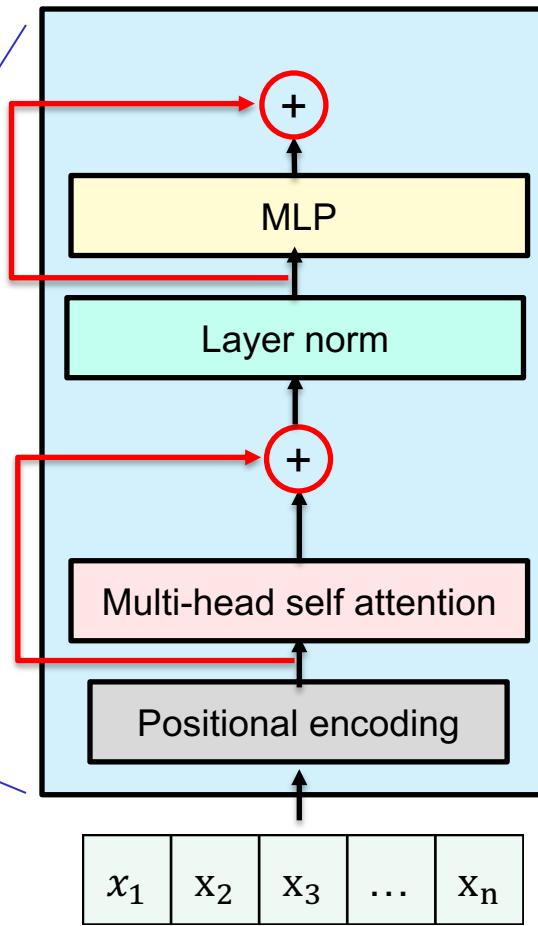
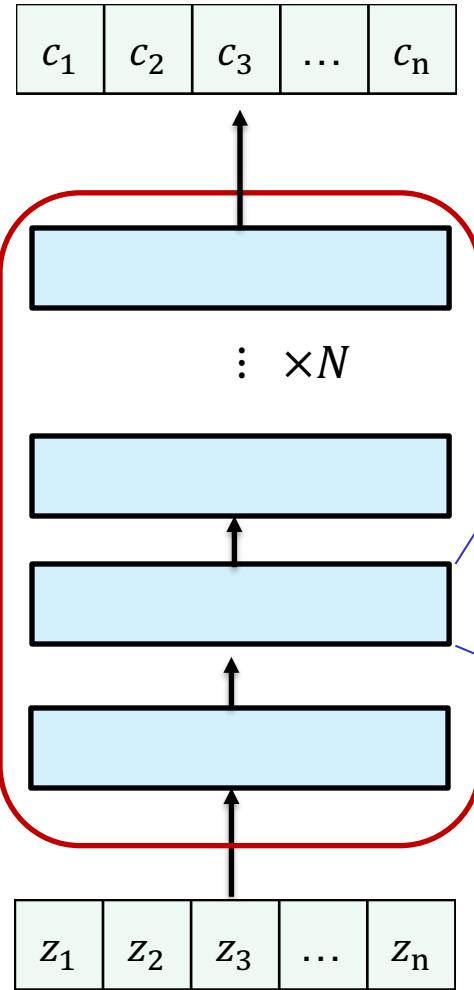
Layer norm over each vector individually

Residual connection

Add pos. encoding

The Transformer encoder block

Transformer encoder



Residual connection

MLP over each vector individually

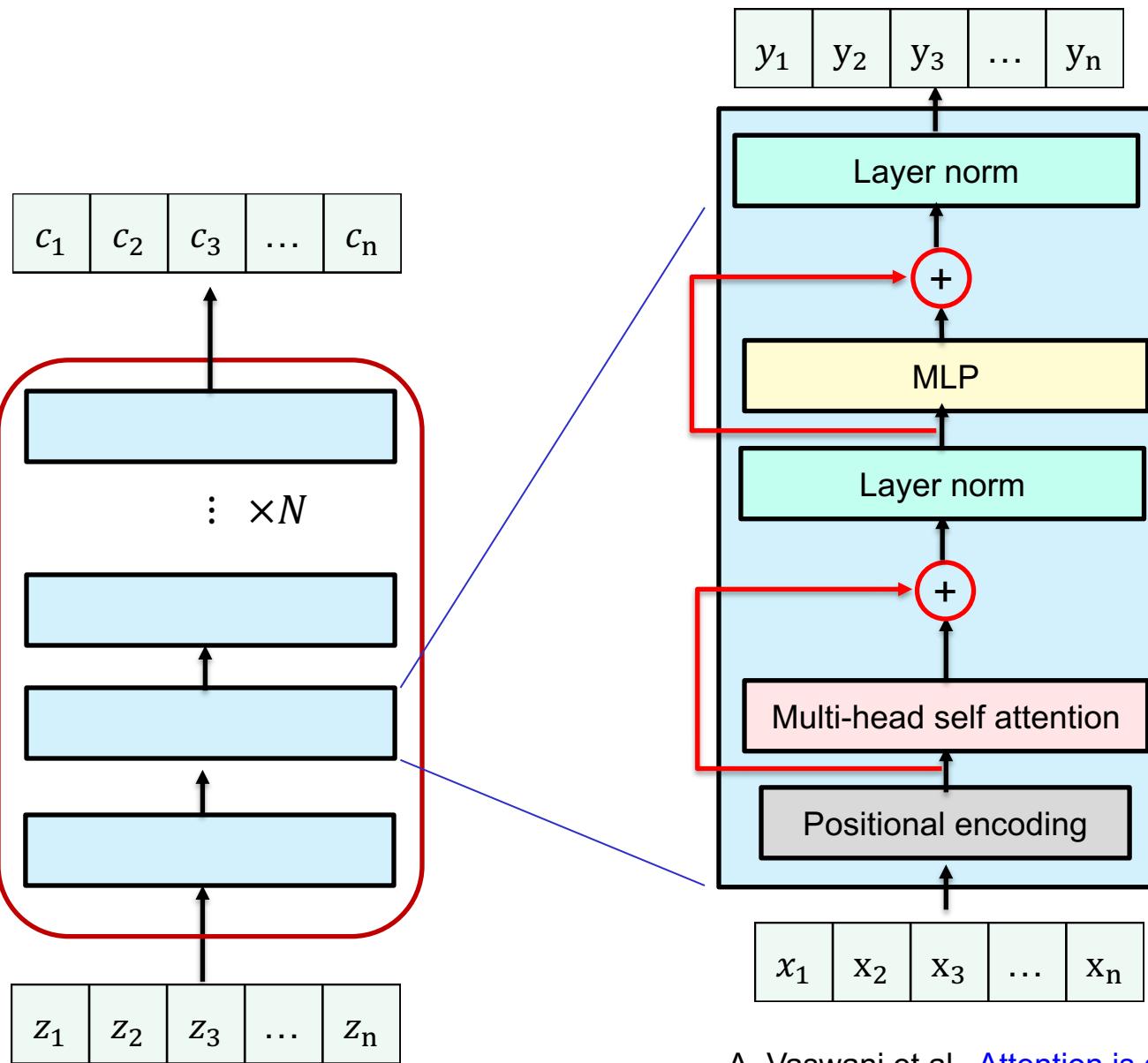
Layer norm over each vector individually

Residual connection

Add pos. encoding

The Transformer encoder block

Transformer encoder



Encoder Block:

Inputs: Set of vectors \mathbf{x}

Outputs: Set of vectors \mathbf{y}

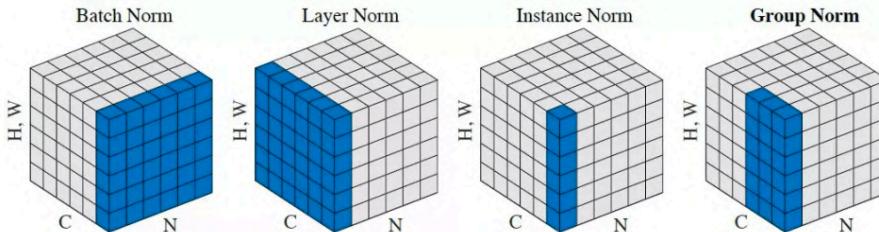
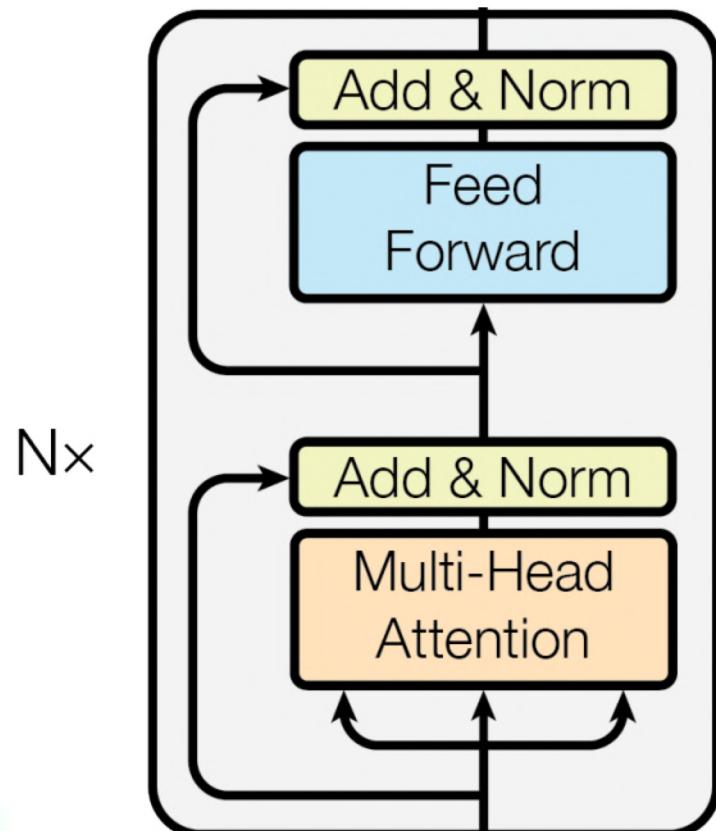
Self-attention is the only interaction between vectors.

Layer norm and MLP operate independently per vector.

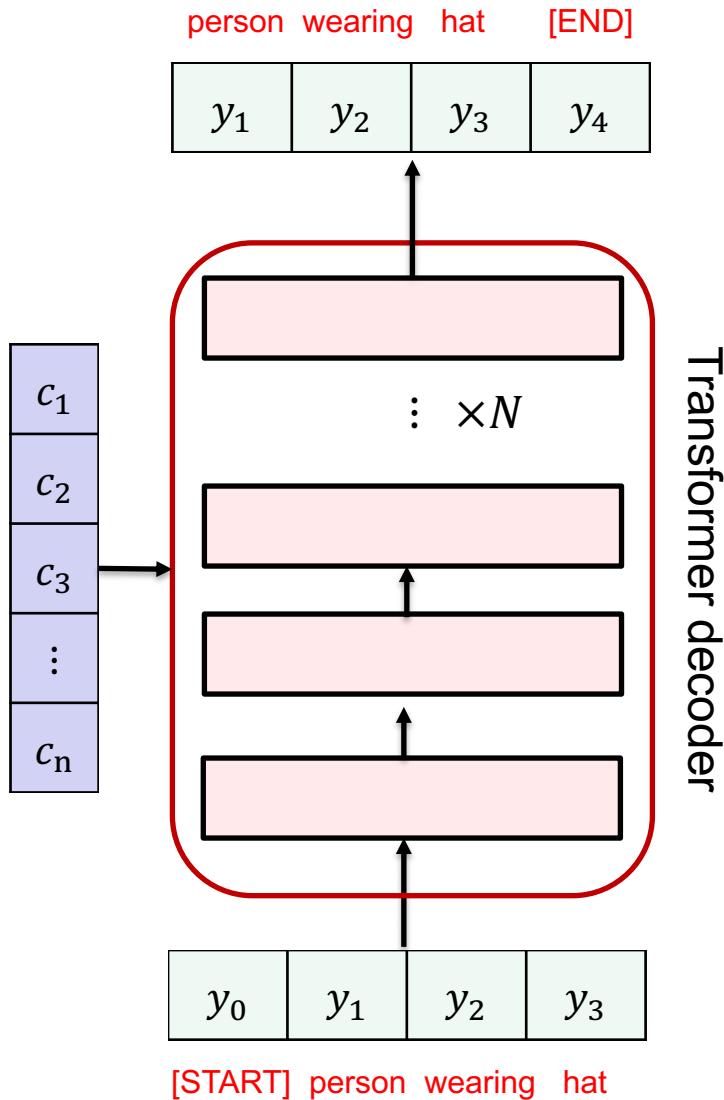
Highly scalable, highly parallelizable, but high memory usage.

The Transformer encoder block

- A **Transformer** is a sequence of transformer blocks
 - Vaswani et al.: N=12 blocks, embedding dimension = 512, 6 attention heads
 - **Add & Norm**: residual connection followed by layer normalization
 - **Feedforward**: two linear layers with ReLUs in between, applied independently to each vector
- Attention is the only interaction between inputs!

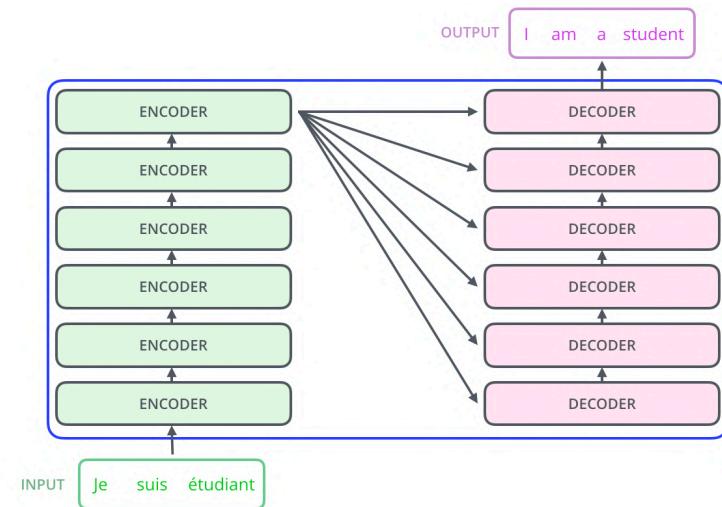


The Transformer decoder block

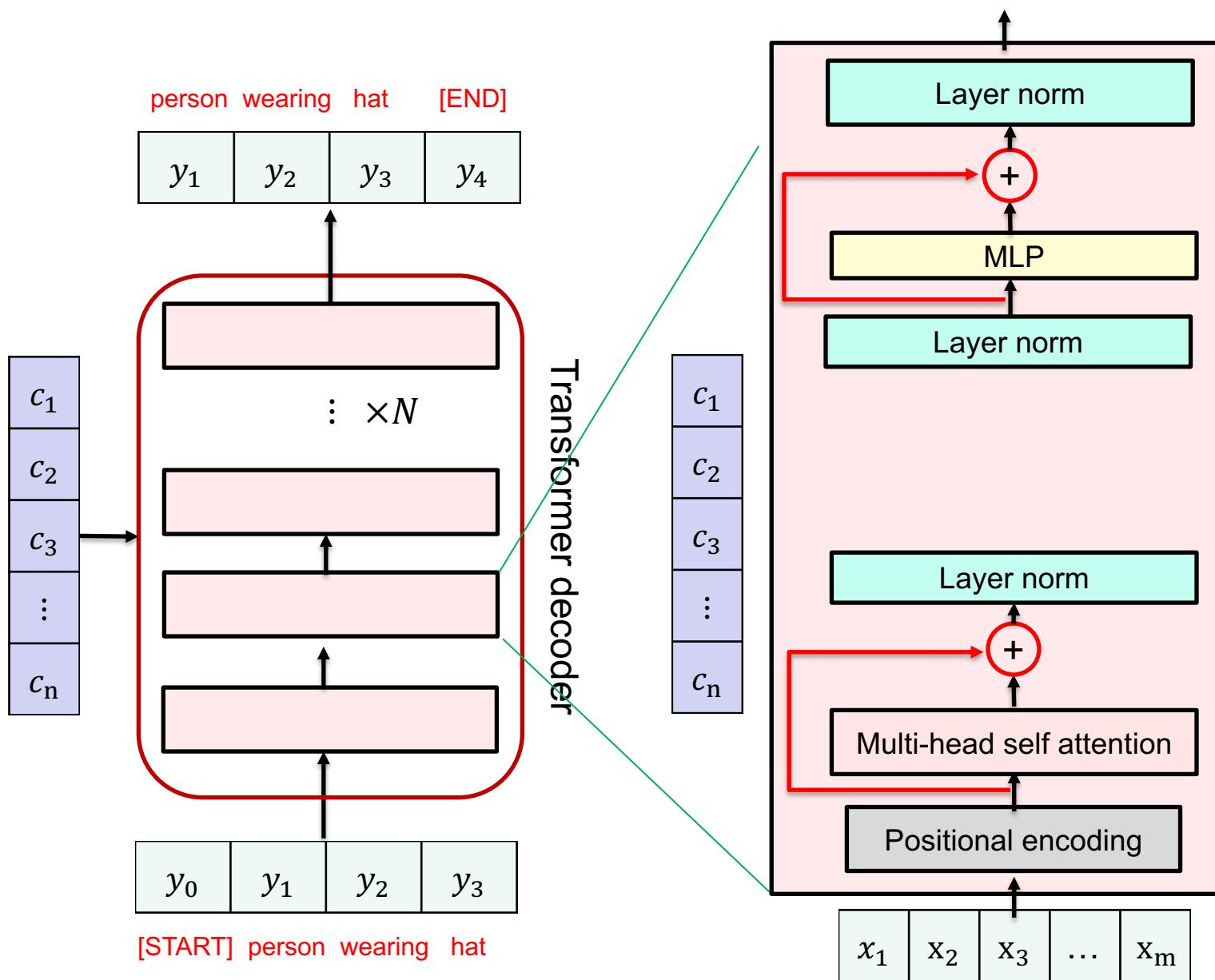


Made up of N decoder blocks

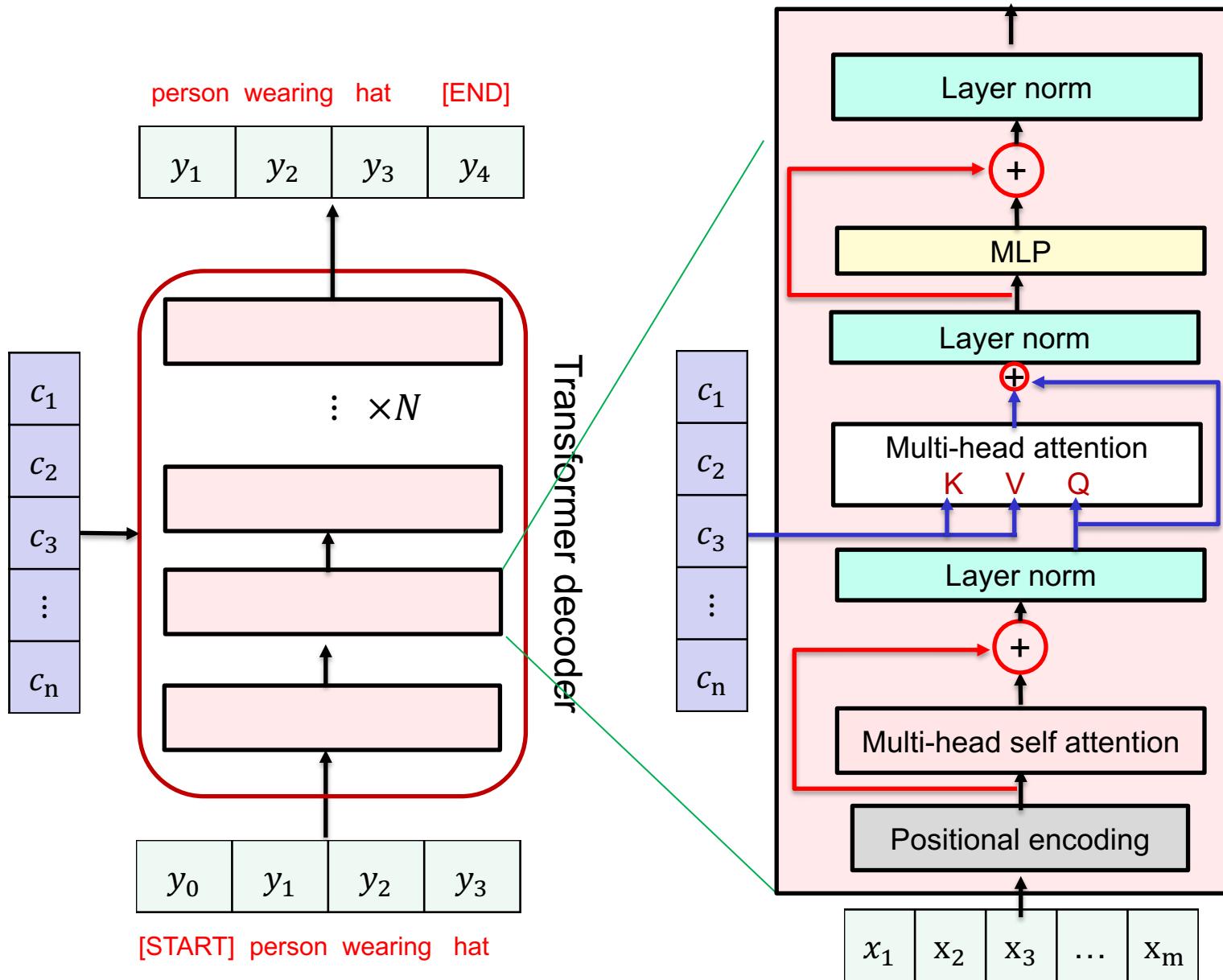
In Vaswani et al. $N=6$, $D_q=512$



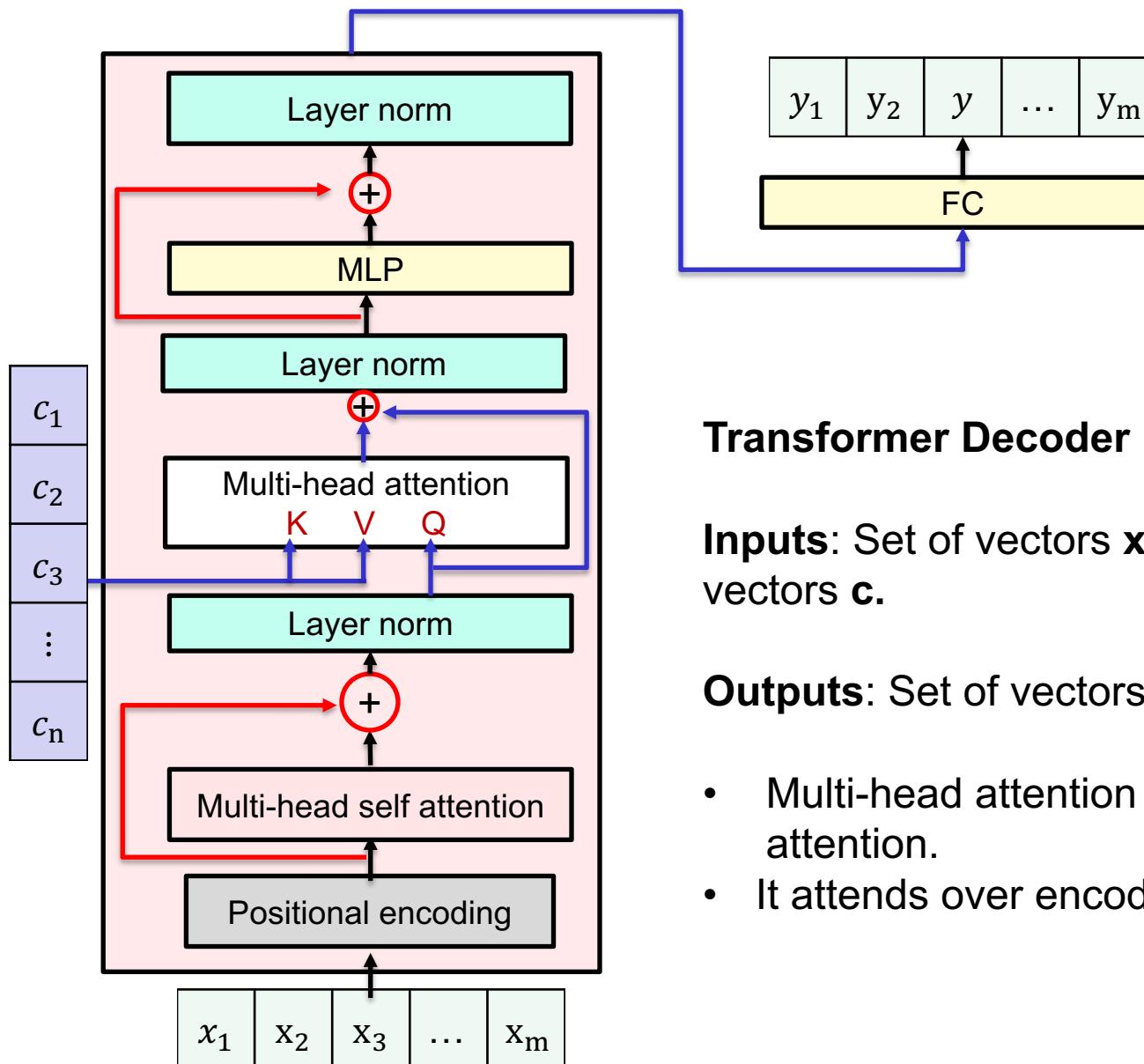
The Transformer decoder block



The Transformer decoder block



The Transformer decoder block



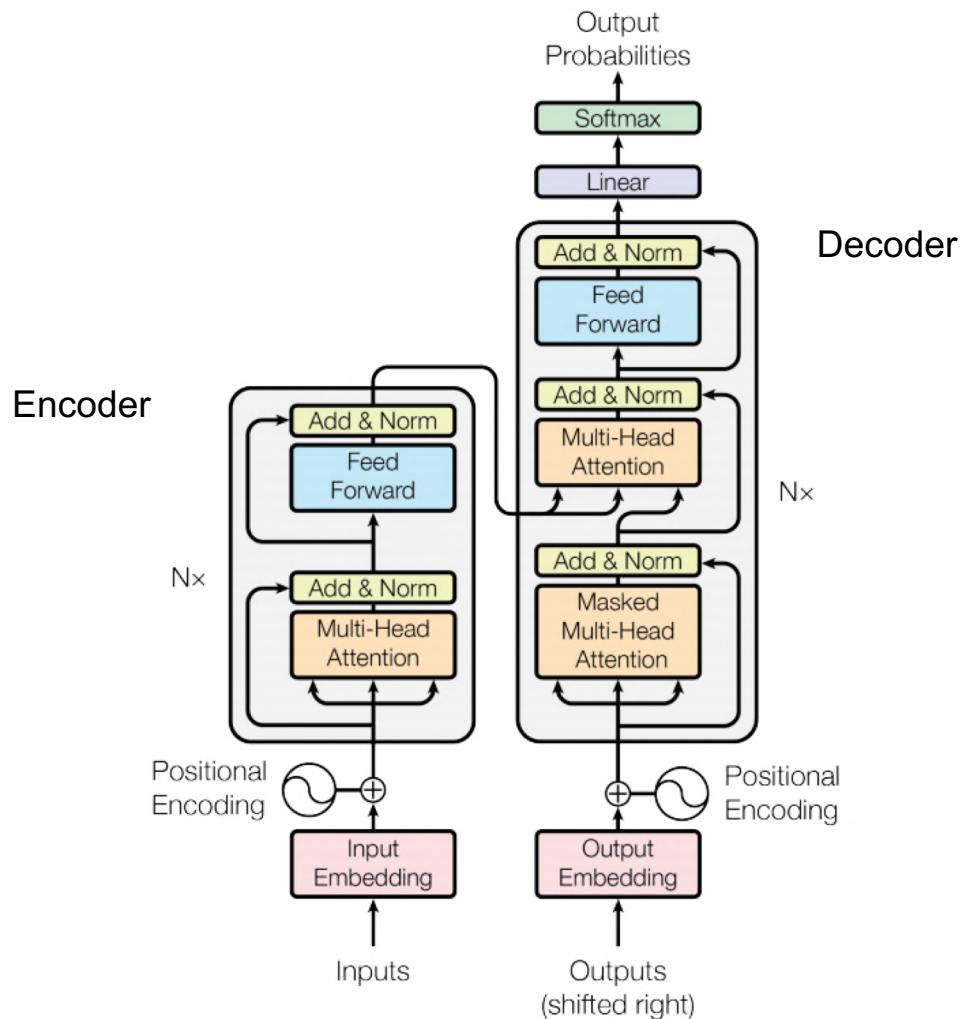
Transformer Decoder Block:

Inputs: Set of vectors \mathbf{x} and Set of context vectors \mathbf{c} .

Outputs: Set of vectors \mathbf{y} .

- Multi-head attention block is NOT self-attention.
- It attends over encoder outputs.

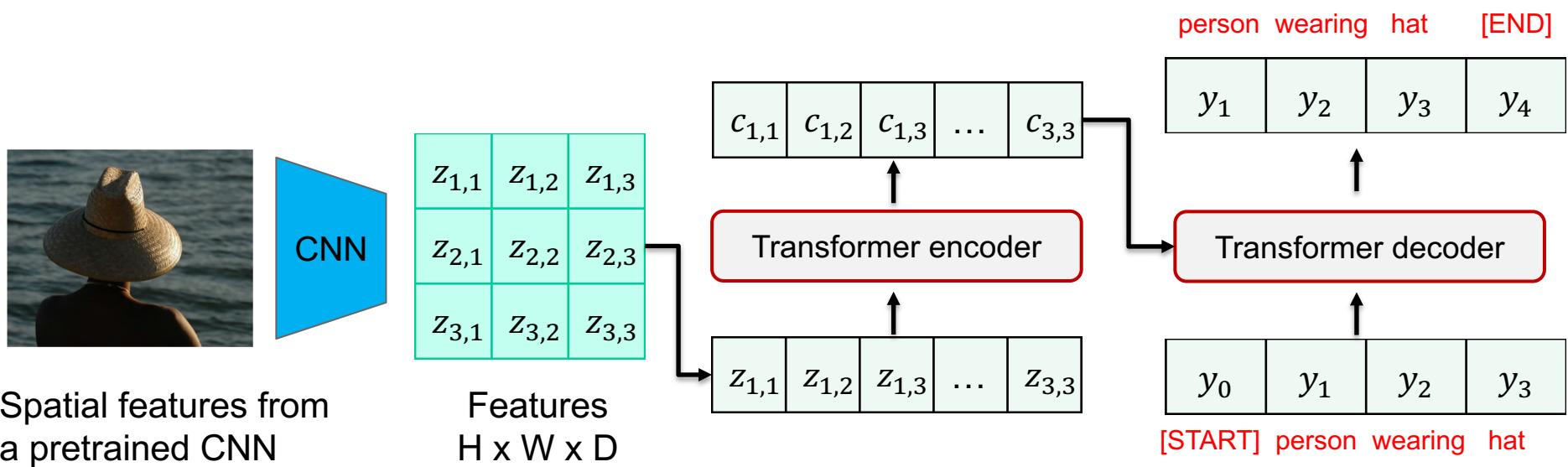
Transformer architecture: Details



Back to Image captioning

Input: image I

Output: $y = y_1, y_2, y_3 \dots, y_T$



Back to Image captioning

Perhaps we don't need convolution at all!

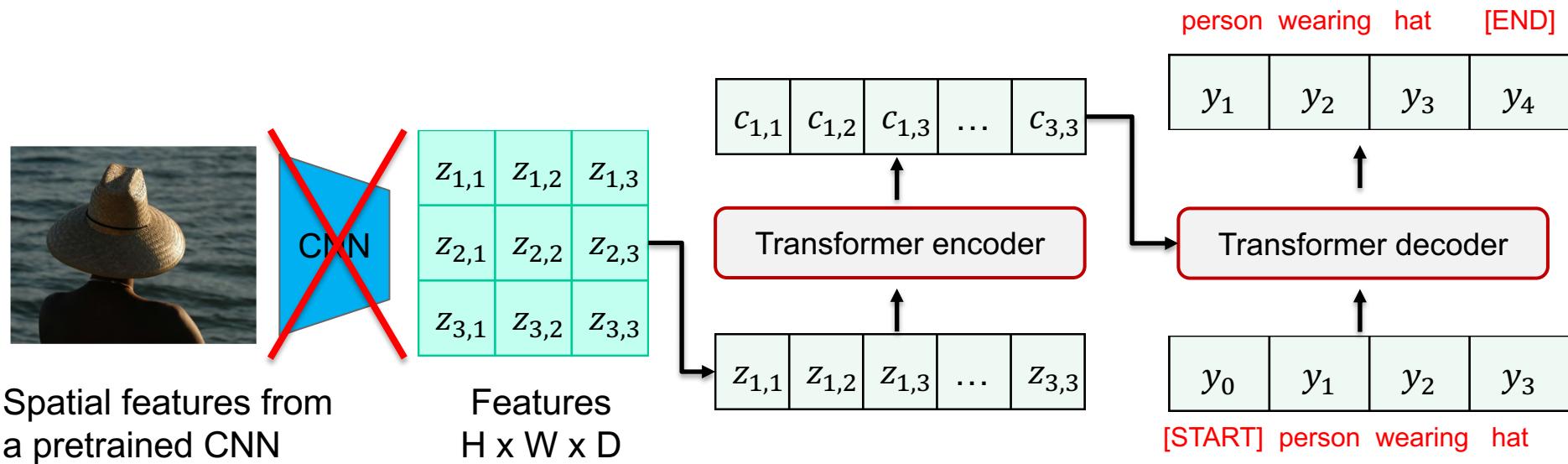
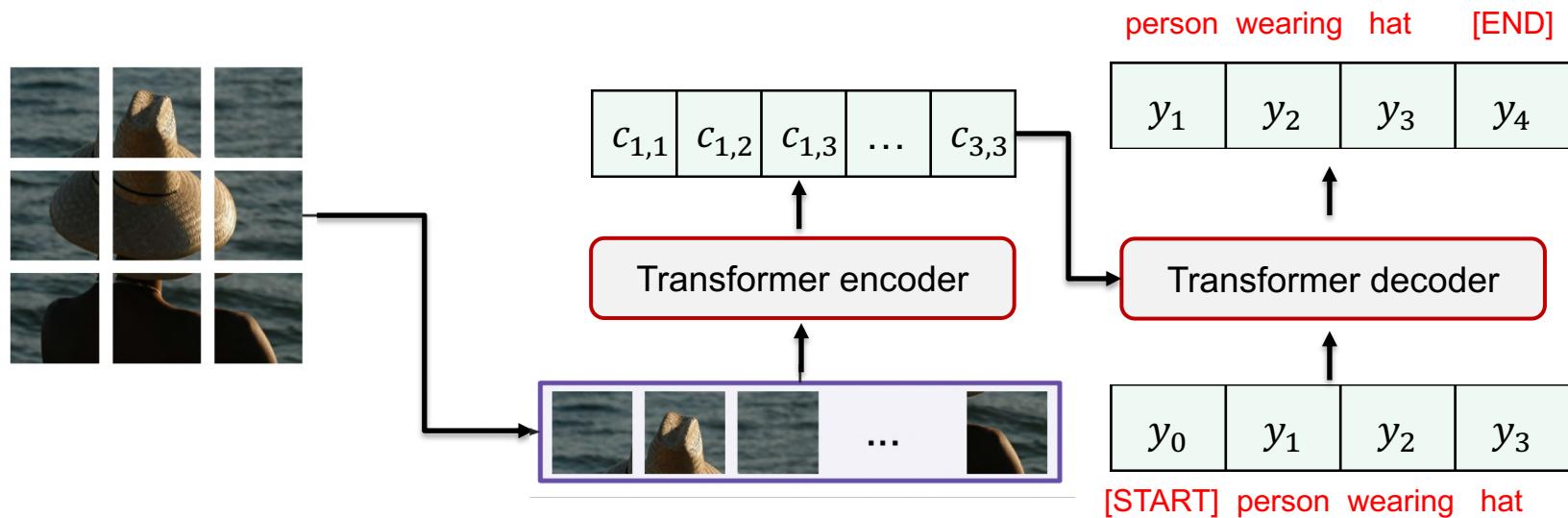


Image captioning using **only** transformers

Transformers directly from pixels to language.

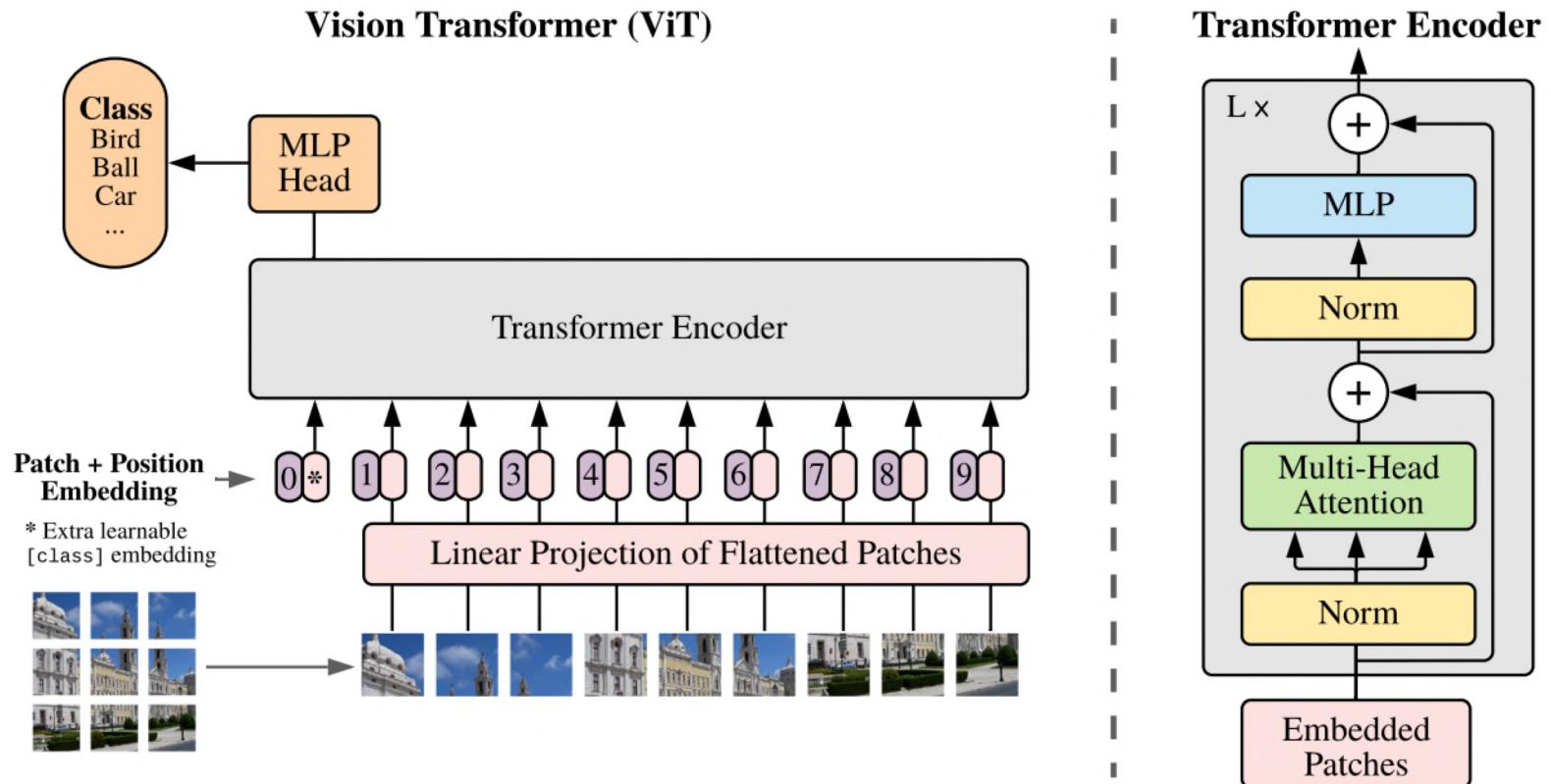
This is called a “**Vision Transformer**” (ViT).

Suggested in: A. Dosovitskiy et al. [An image is worth 16x16 words: Transformers for image recognition at scale](#). ICLR 2021



Vision Transformers for Image Classification

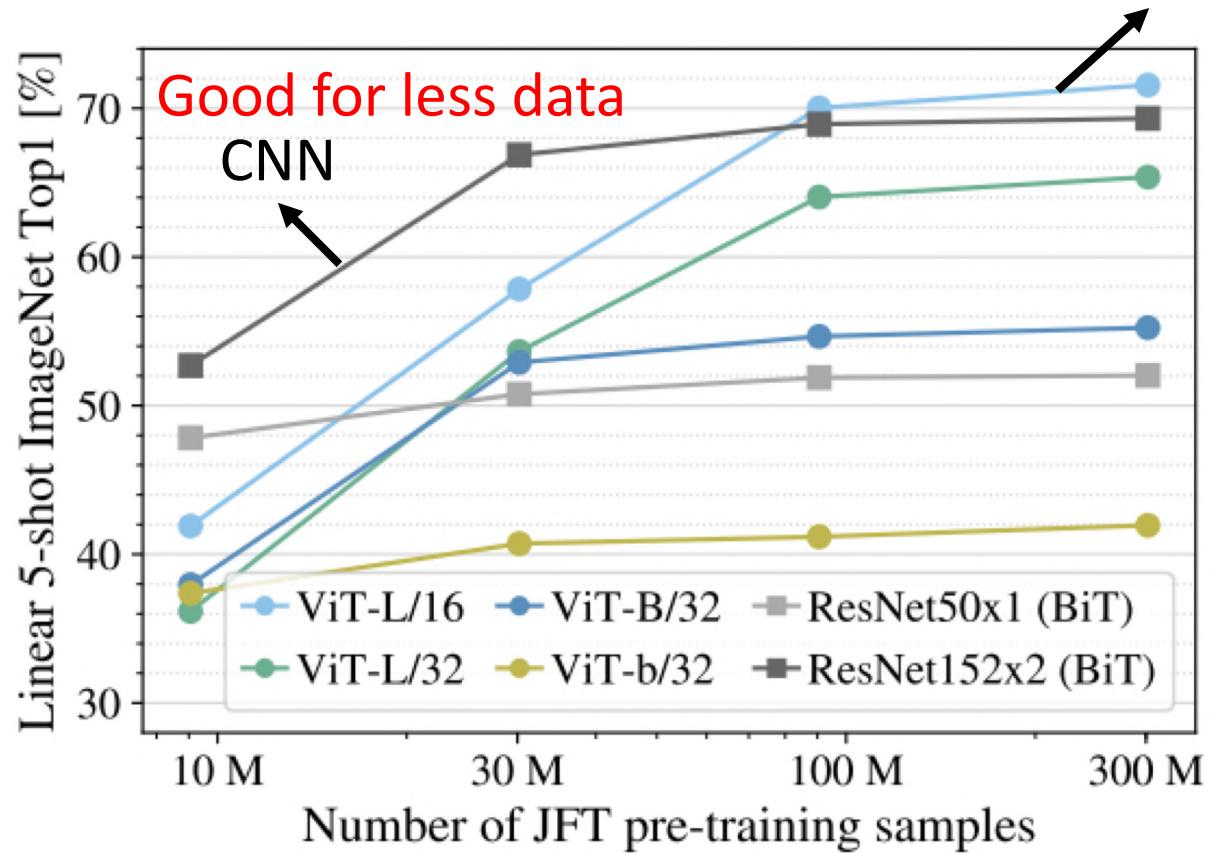
- Split an image into patches, feed linearly projected patches into standard transformer encoder



Self-attention v.s. CNN

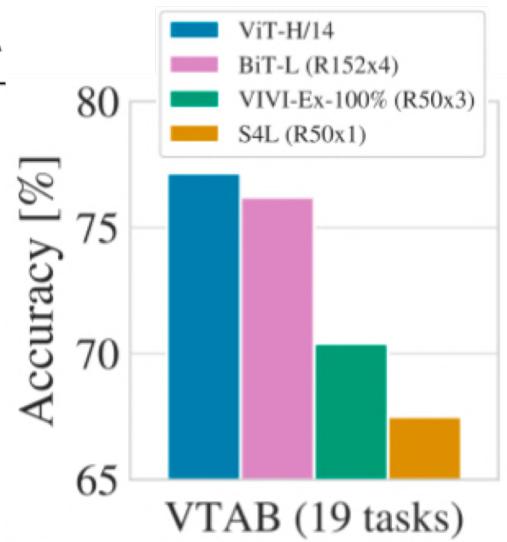
Good for more data

Self-attention



ViT Results

	ViT-H	Previous SOTA
ImageNet	88.55	88.5
ImageNet-Real	90.72	90.55
Cifar-10	99.50	99.37
Cifar-100	94.55	93.51
Pets	97.56	96.62
Flowers	99.68	99.63



Vision Transformer matches or outperforms state-of-the-art CNNs on popular benchmarks. Left: Popular image classification tasks (ImageNet, including new validation labels [Real](#), and [CIFAR](#), [Pets](#), and [Flowers](#)). Right: Average across 19 tasks in the VTAB classification suite.

Image transformer – Google

- Image generation and super-resolution with 32x32 output, attention restricted to local neighborhoods

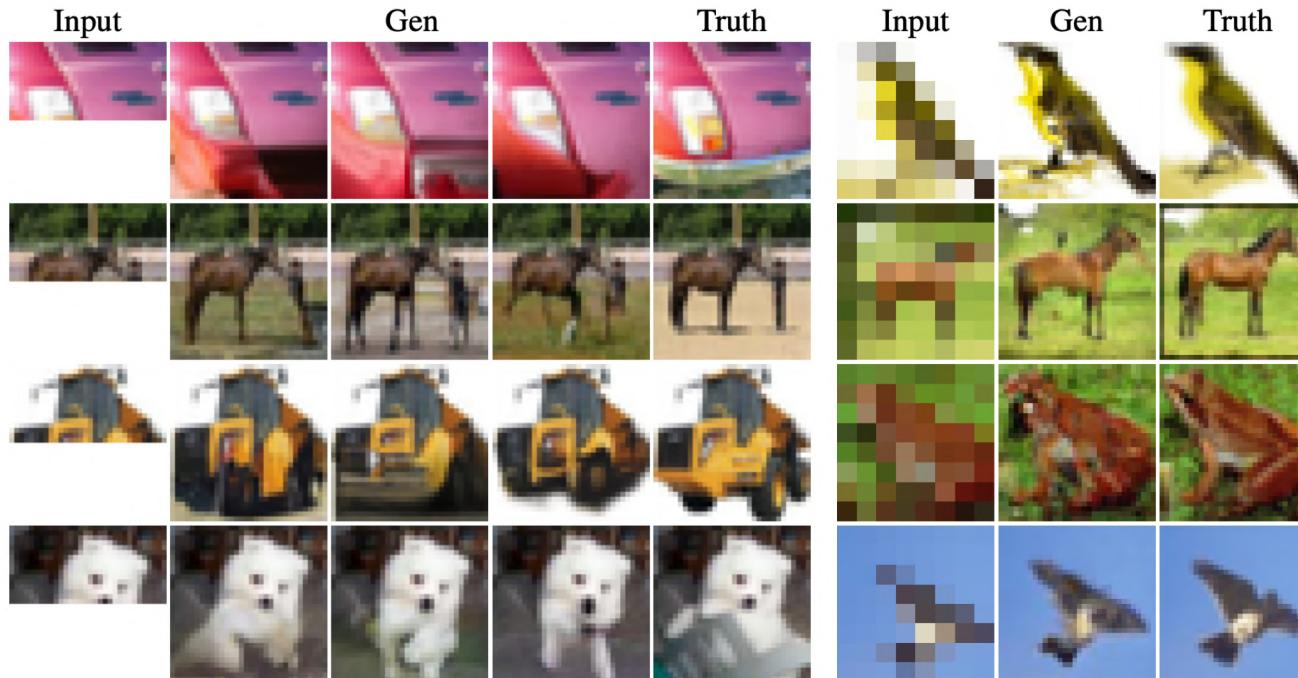
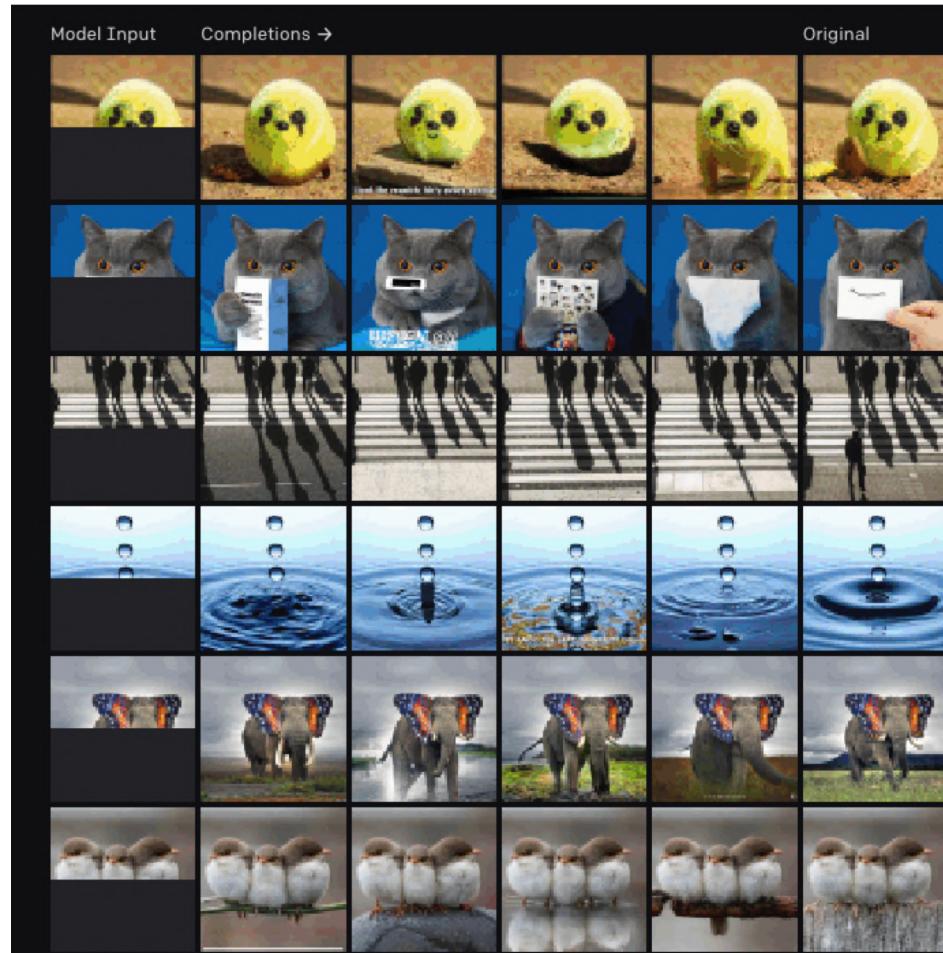


Table 2. On the left are image completions from our best conditional generation model, where we sample the second half. On the right are samples from our four-fold super-resolution model trained on CIFAR-10. Our images look realistic and plausible, show good diversity among the completion samples and observe the outputs carry surprising details for coarse inputs in super-resolution.

Image GPT – OpenAI



M. Chen et al., [Generative pretraining from pixels](#), ICML 2020

<https://openai.com/blog/image-gpt/>

Image transformer – Google

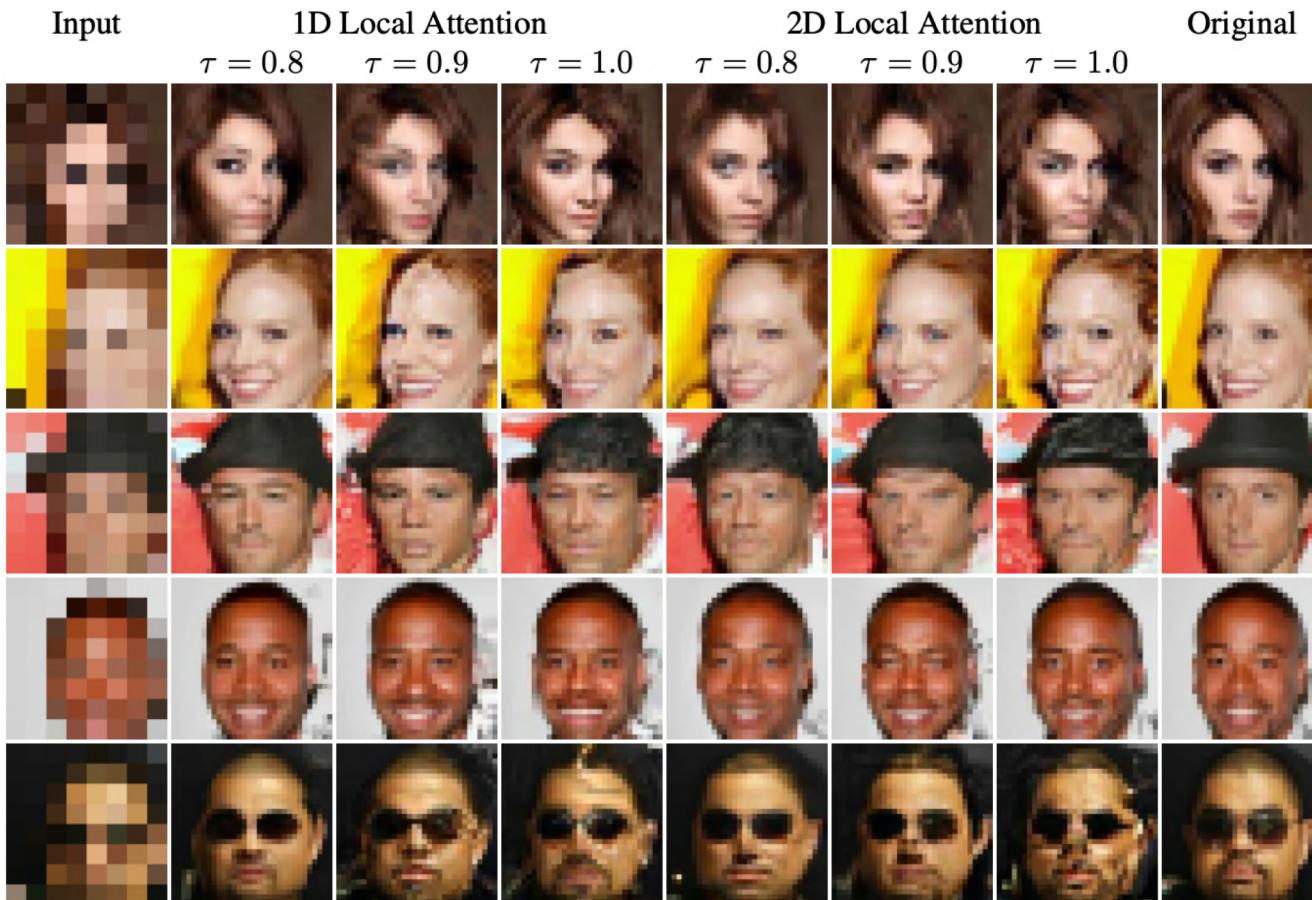


Table 6. Images from our 1D and 2D local attention super-resolution models trained on CelebA, sampled with different temperatures. 2D local attention with $\tau = 0.9$ scored highest in our human evaluation study.

CLIP for image editing

- Combine powerful recent image-text embedding technique ([CLIP](#)) with pre-trained StyleGAN for text-based image editing



“Emma Stone”

“Mohawk hairstyle”

“Without makeup”

“Cute cat”

“Lion”

“Gothic church”

O. Patashnik et al.. [StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery](#). arXiv 2021

Text-based image generation: DALL-E

- Learn a joint sequential transformer model that can be used to generate image based on text prompt



(a) a tapir made of accordion.
a tapir with the texture of an
accordion.

(b) an illustration of a baby
hedgehog in a christmas
sweater walking a dog

(c) a neon sign that reads
“backprop”. a neon sign that
reads “backprop”. backprop
neon sign

Text-based image generation: DALL-E

- Learn a joint sequential transformer model that can be used to generate image based on text prompt

a store front that has the word 'openai' written on it. . . .



Text-based image generation: DALL-E

- Learn a joint sequential transformer model that can be used to generate image based on text prompt

the exact same cat on the top as a sketch on the bottom



Text-based image generation: DALL-E

- Learn a joint sequential transformer model that can be used to generate image based on text prompt

an armchair in the shape of an avocado. . .



Summary

- Adding **attention** to RNNs allows them to "attend" to different parts of the input at every time step
- The **general attention layer** is a new type of layer that can be used to design new neural network architectures
- **Transformers** are a type of layer that uses **self-attention** and layer norm.
 - It is highly **scalable** and highly **parallelizable**
 - **Faster** training, **larger** models, **better** performance across vision and language tasks
 - They are quickly replacing RNNs, LSTMs, and may even replace convolutions.