

# Attention and Transformer

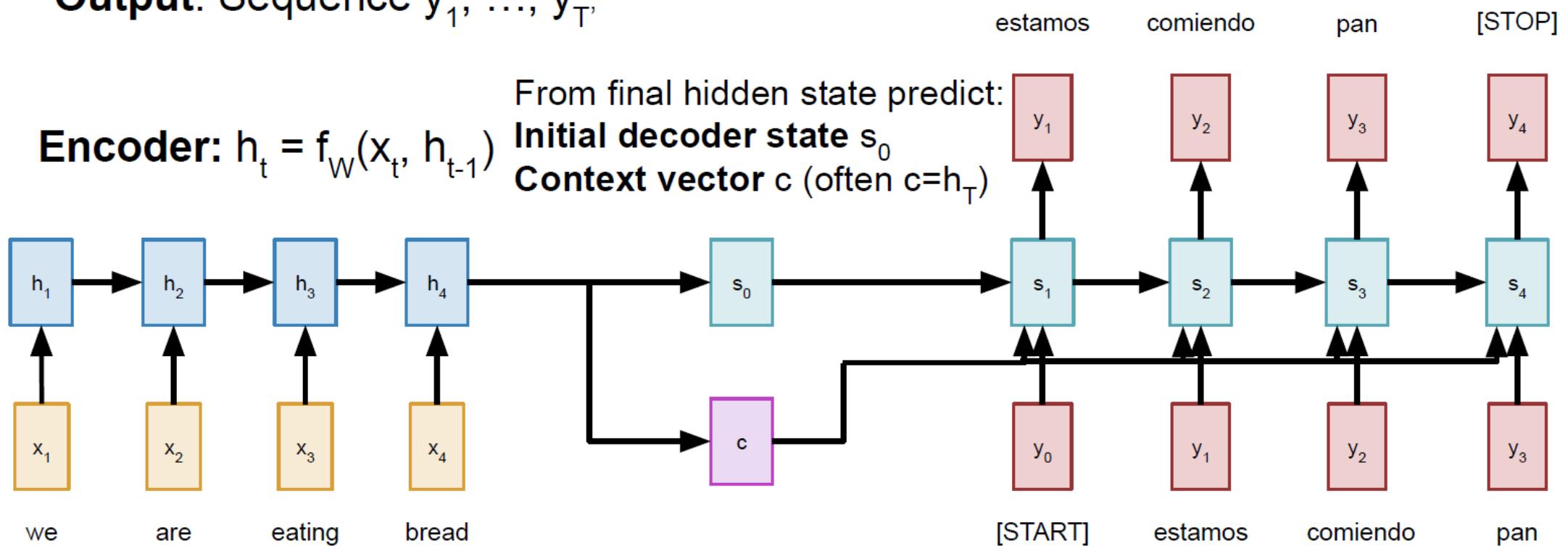
Slides refer to Prof. Fei-Fei Li's one at Stanford

# Sequence to Sequence with RNNs

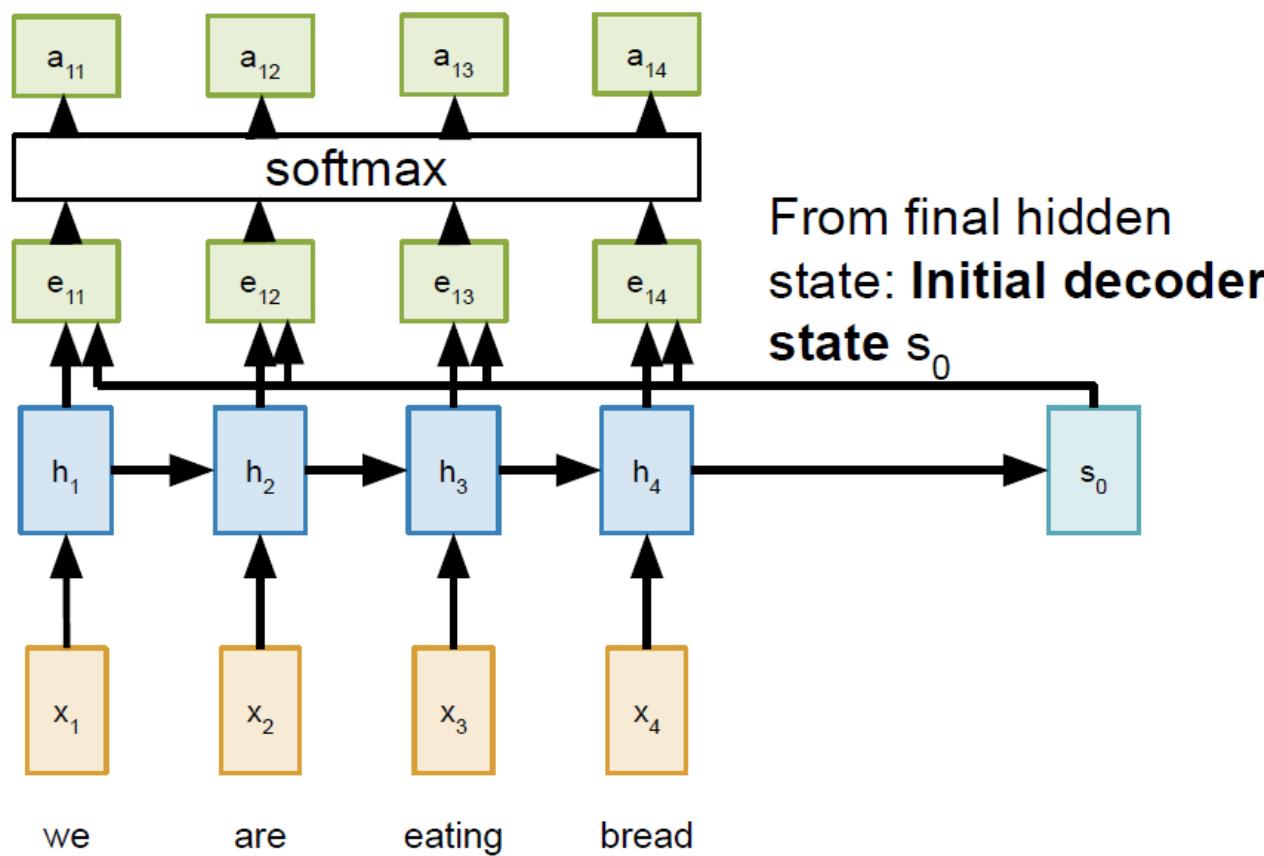
**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

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



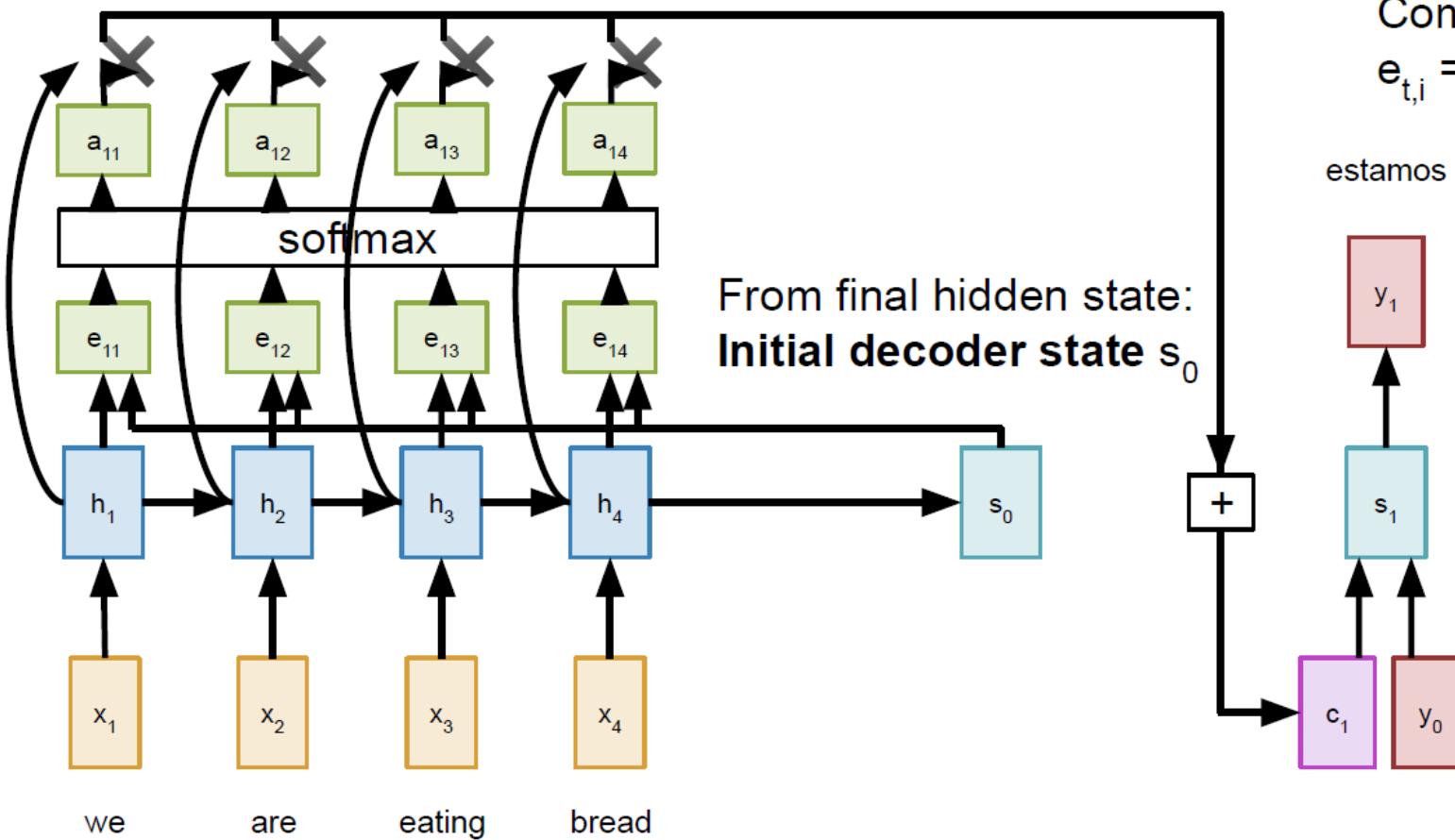
# Sequence to Sequence with RNNs and Attention



Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{att}(s_{t-1}, h_i)$  ( $f_{att}$  is an MLP)

Normalize alignment scores  
to get **attention weights**  
 $0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$

# Sequence to Sequence with RNNs and Attention



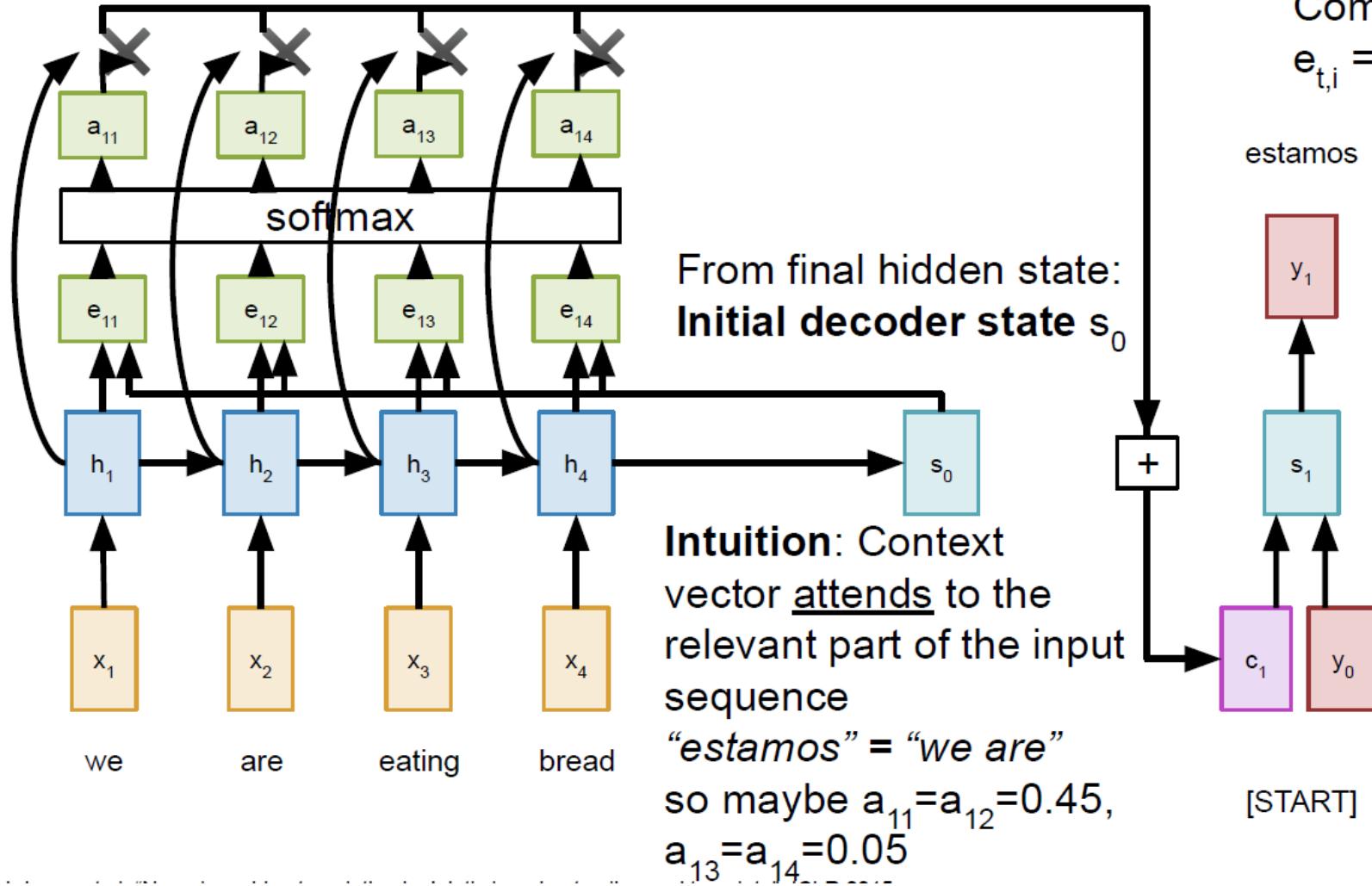
Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{att}(s_{t-1}, h_i)$  ( $f_{att}$  is an MLP)

estamos

Normalize alignment scores  
to get **attention weights**  
 $0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$

Compute context vector as  
linear combination of hidden  
states  
 $c_t = \sum_i a_{t,i} h_i$

# Sequence to Sequence with RNNs and Attention



Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$  ( $f_{\text{att}}$  is an MLP)

estamos

Normalize alignment scores to get **attention weights**  
 $0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$

Compute context vector as linear combination of hidden states

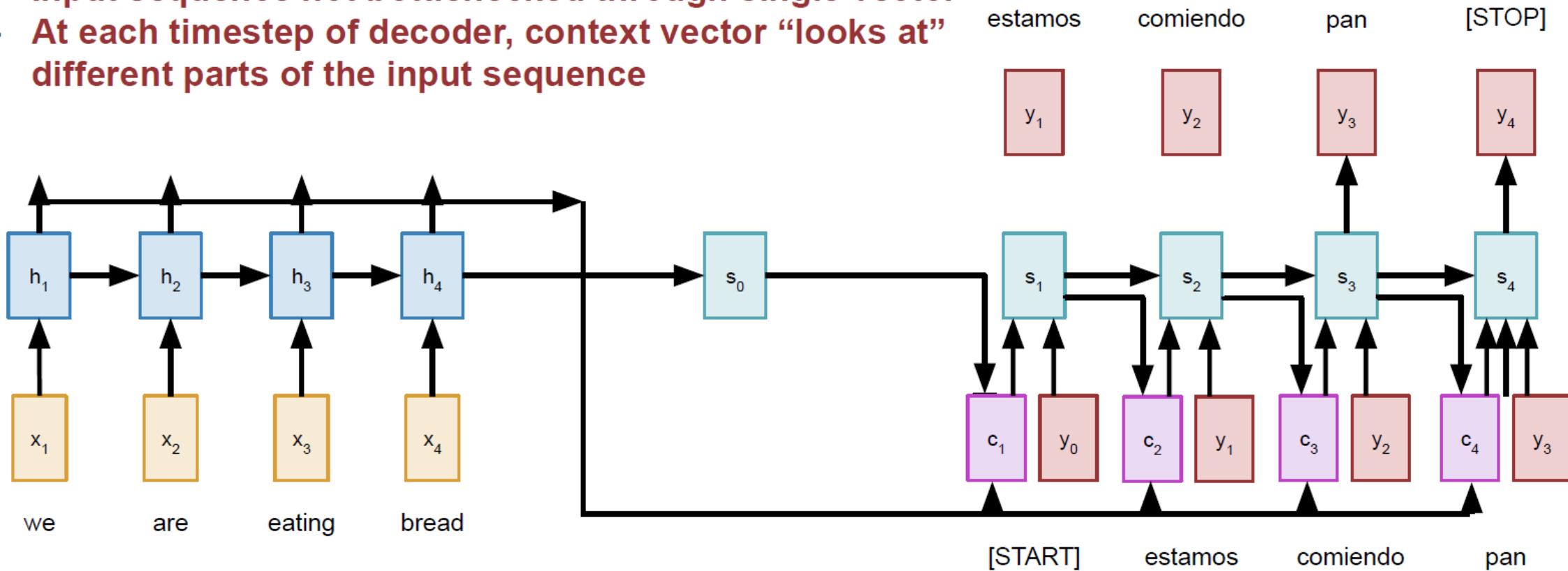
$$c_t = \sum_i a_{t,i} h_i$$

Use context vector in decoder:  $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

# Sequence to Sequence with RNNs and Attention

Use a different context vector in each timestep of decoder

- Input sequence not bottlenecked through single vector
- At each timestep of decoder, context vector “looks at” different parts of the input sequence



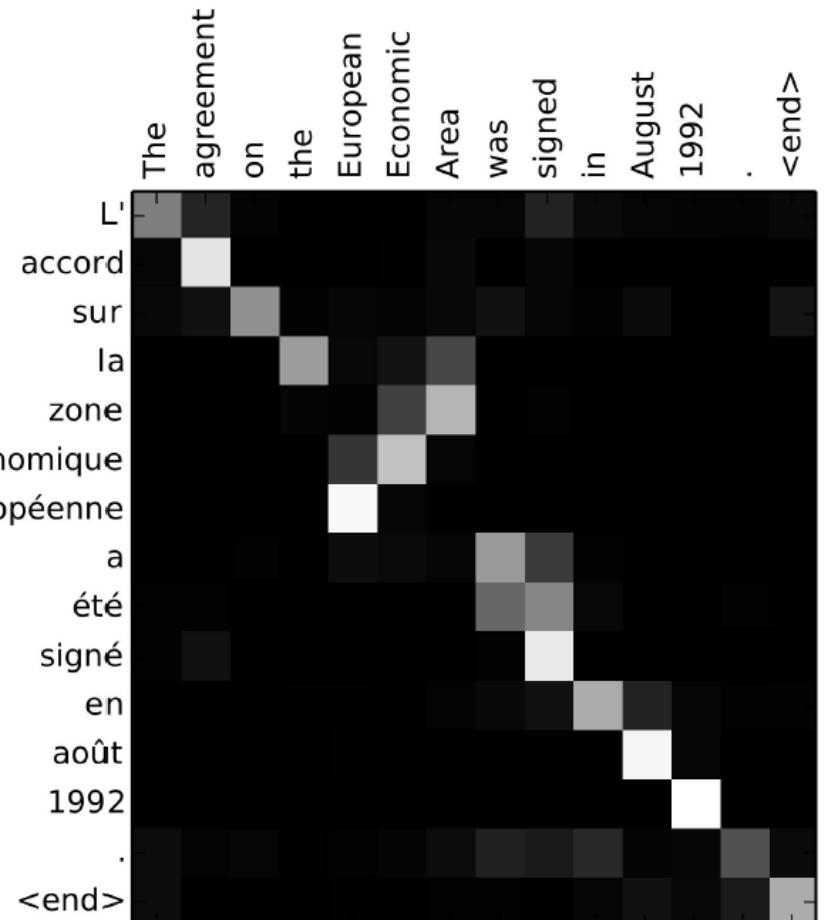
# Sequence to Sequence with RNNs and Attention

**Example:** English to French translation

**Input:** “The agreement on the European Economic Area was signed in August 1992.”

**Output:** “L'accord sur la zone économique européenne a été signé en août 1992.”

Visualize attention weights  $a_{t,i}$



# Sequence to Sequence with RNNs and Attention

**Example:** English to French translation

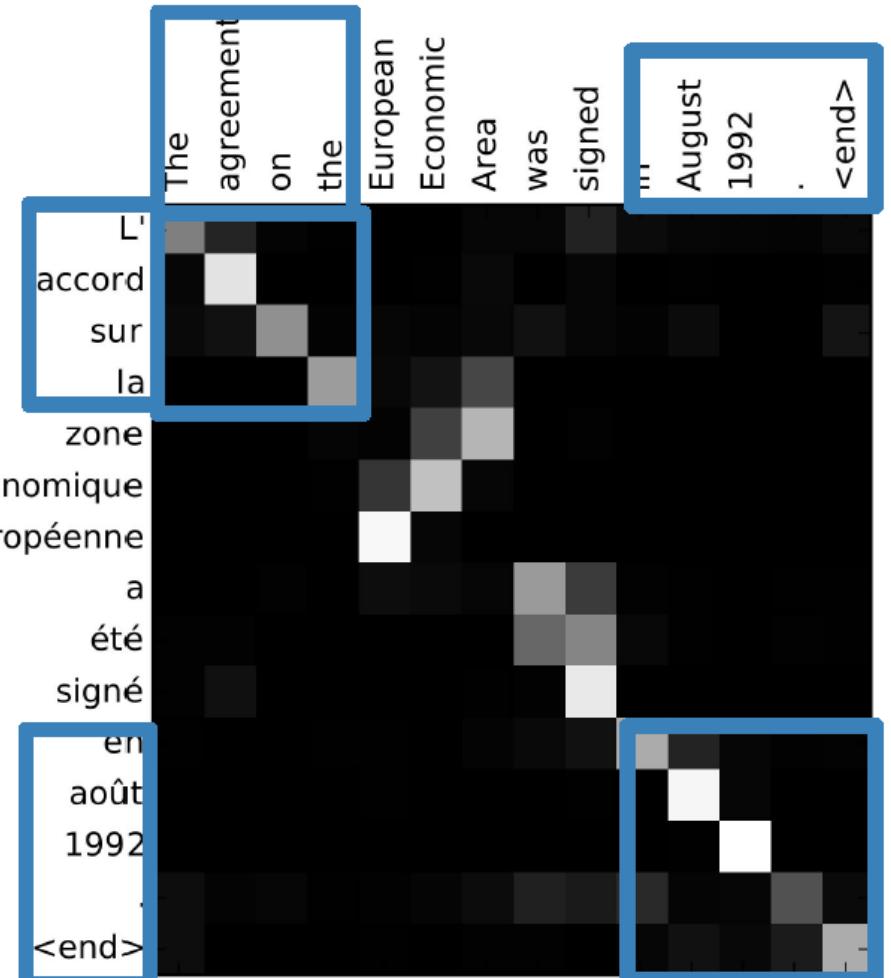
**Input:** “**The agreement on the European Economic Area was signed in August 1992.**”

**Output:** “**L'accord sur la zone économique européenne a été signé en août 1992.**”

Diagonal attention means words correspond in order

Diagonal attention means words correspond in order

Visualize attention weights  $a_{t,i}$



# Sequence to Sequence with RNNs and Attention

**Example:** English to French translation

**Input:** “**The agreement on the European Economic Area** was signed **in August 1992**.”

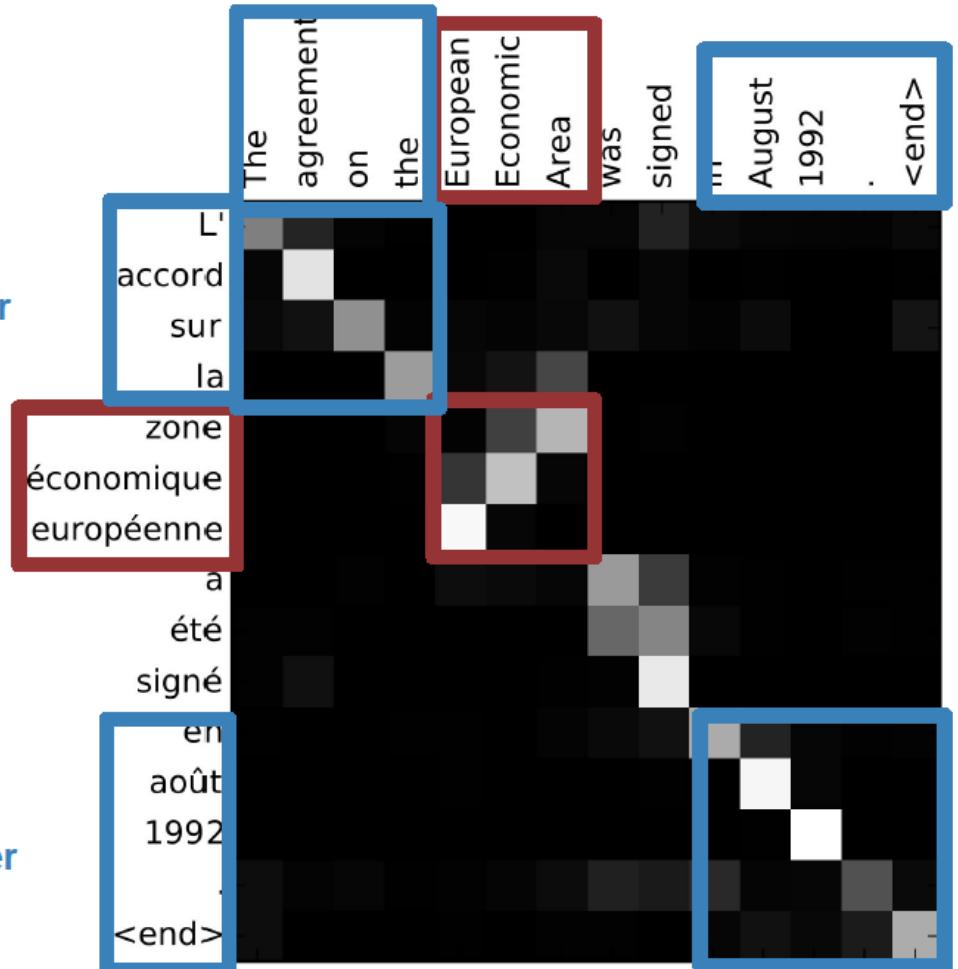
**Output:** “**L'accord sur la zone économique européenne** a été signé **en août 1992**.”

Diagonal attention means words correspond in order

Attention figures out different word orders

Diagonal attention means words correspond in order

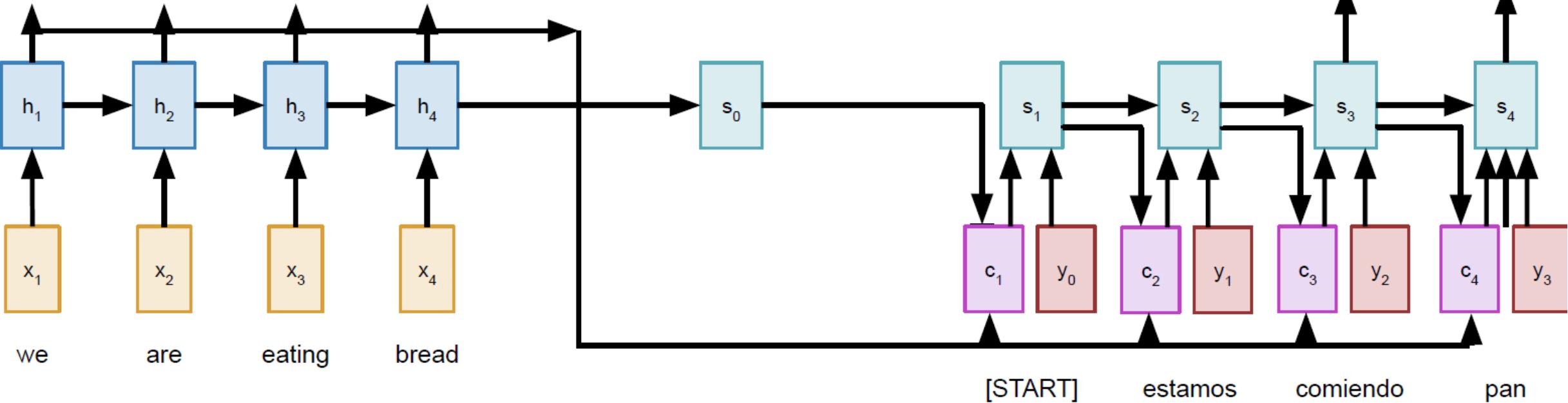
Visualize attention weights  $a_{t,i}$



# Sequence to Sequence with RNNs and Attention

The decoder doesn't use the fact that  $h_i$  form an ordered sequence – it just treats them as an unordered set  $\{h_i\}$

Can use similar architecture given any set of input hidden vectors  $\{h_i\}$ !



# Image Captioning using spatial features

**Input:** Image I

**Output:** Sequence  $y = y_1, y_2, \dots, y_T$

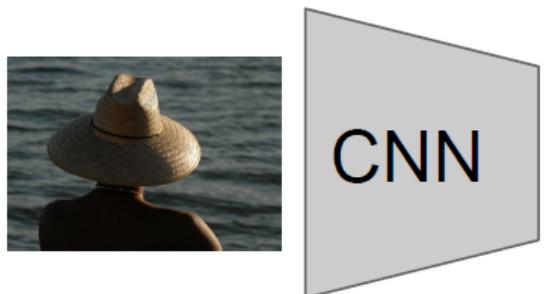
**Decoder:**  $y_t = g_v(y_{t-1}, h_{t-1}, c)$

where context vector  $c$  is often  $c = h_0$

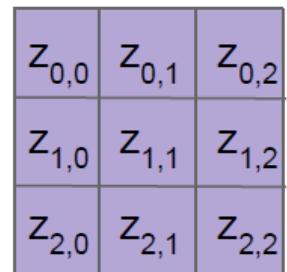
**Encoder:**  $h_0 = f_w(z)$

where  $z$  is spatial CNN features

$f_w(\cdot)$  is an MLP

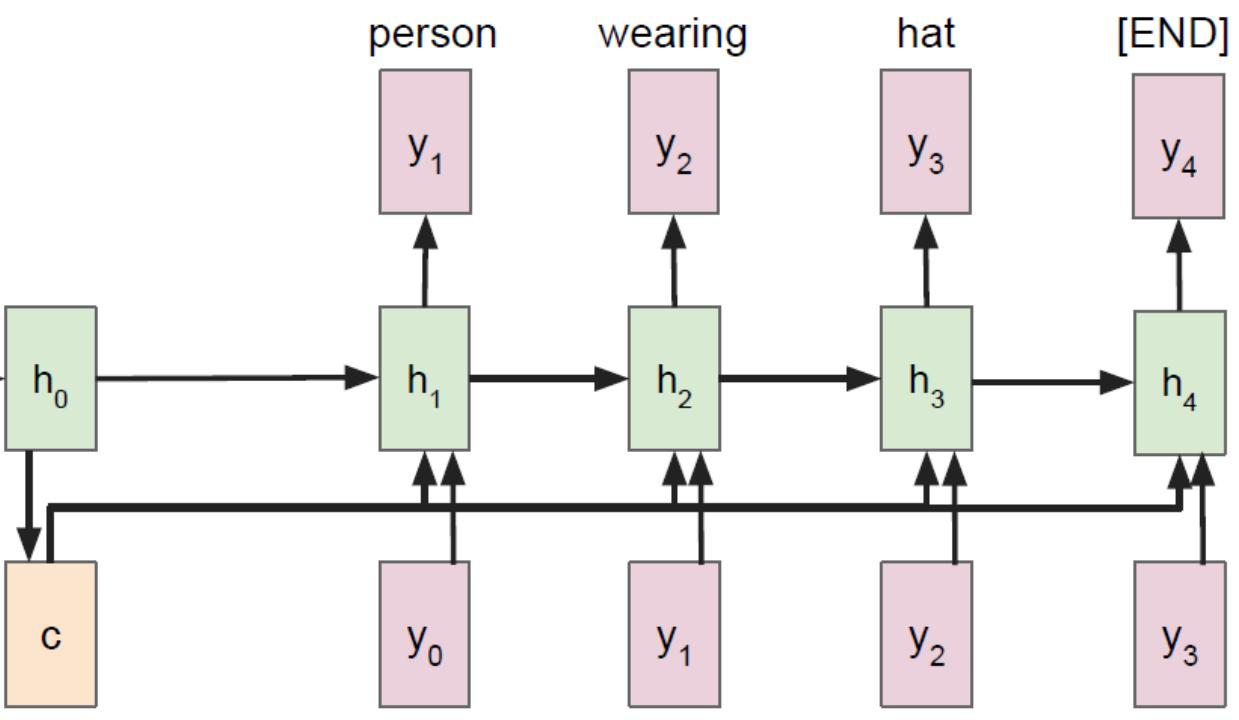


Extract spatial  
features from a  
pretrained CNN



Features:  
 $H \times W \times D$

MLP

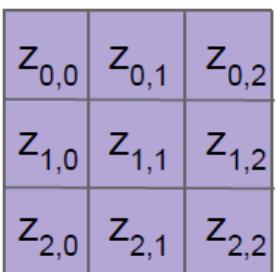
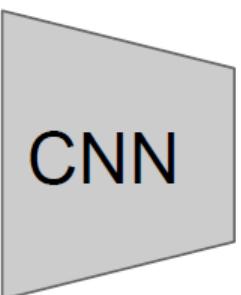


# Image Captioning using spatial features

**Problem: Input is "bottlenecked" through c**

- Model needs to encode everything it wants to say within c

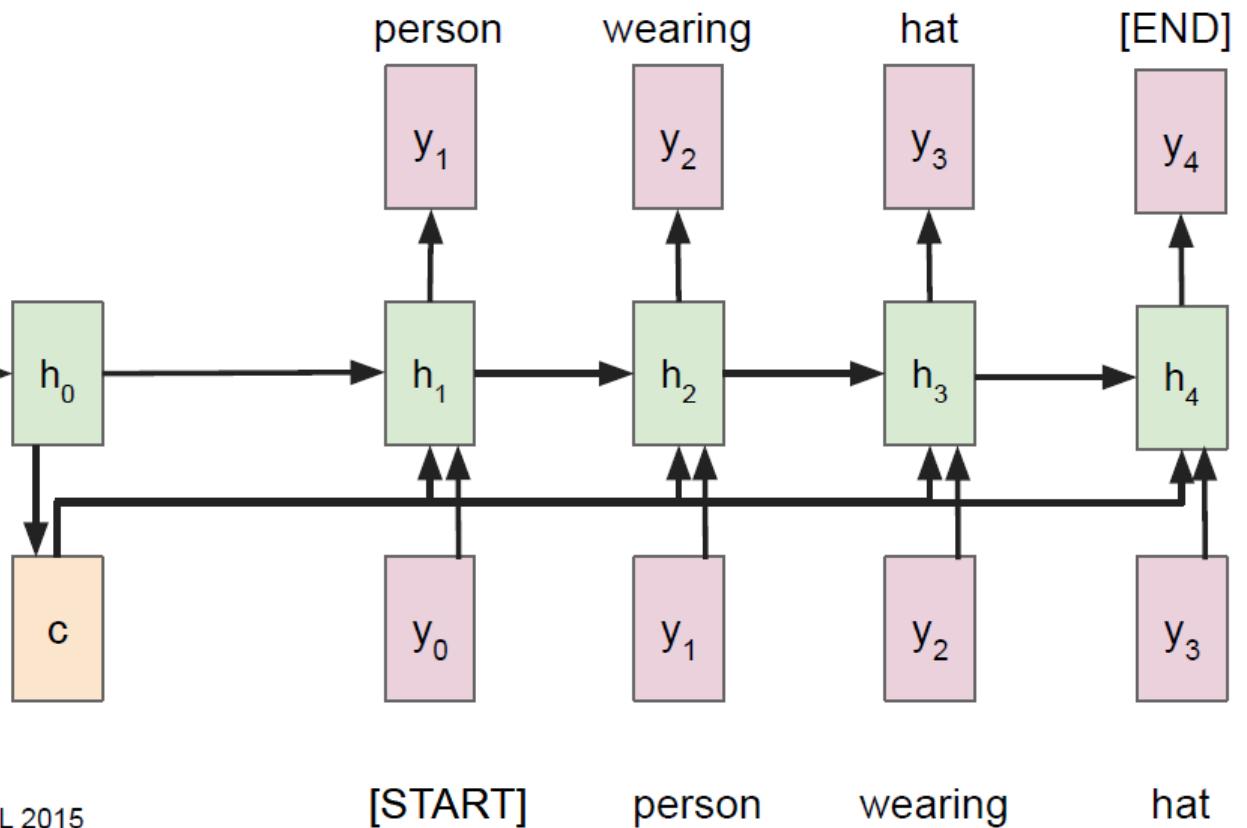
This is a problem if we want to generate really long descriptions? 100s of words long



Extract spatial features from a pretrained CNN

Features:  
 $H \times W \times D$

MLP

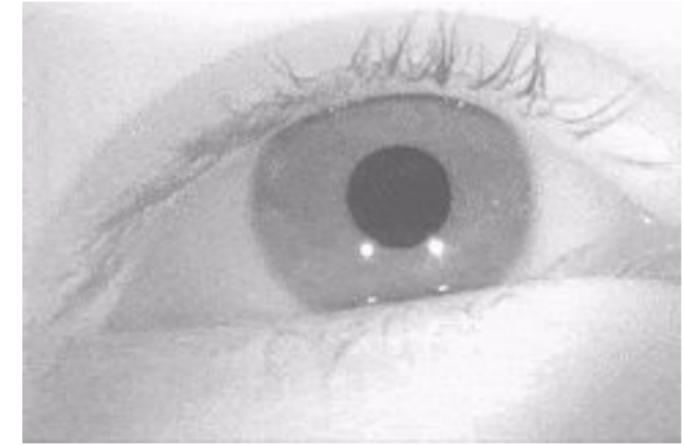


# Image Captioning with RNNs and Attention

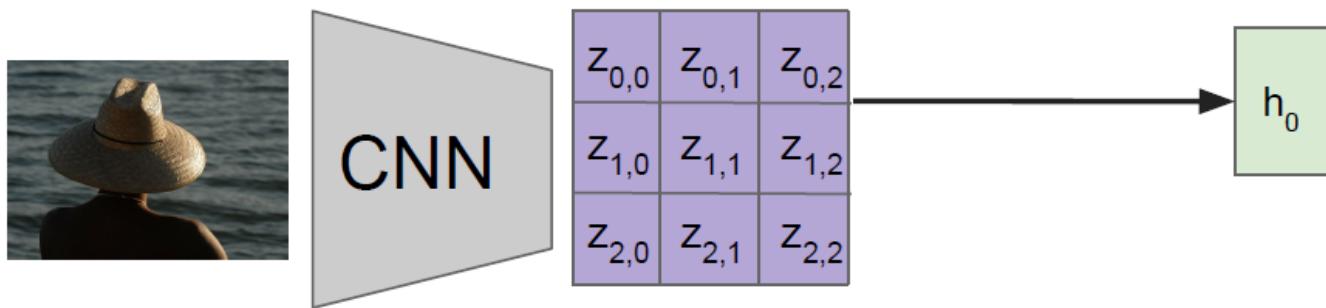
[gif source](#)

Attention idea: New context vector at every time step.

Each context vector will attend to different image regions



Attention Saccades in humans



Extract spatial  
features from a  
pretrained CNN

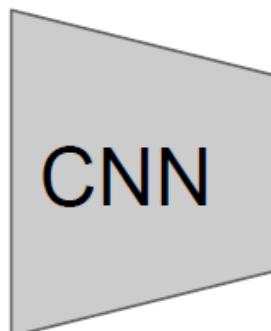
Features:  
 $H \times W \times D$

# Image Captioning with RNNs and Attention

Compute alignments scores (scalars):

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

$f_{att}(\cdot)$  is an MLP



Extract spatial features from a pretrained CNN

Alignment scores:

$H \times W$

$e_{1,0,0}$	$e_{1,0,1}$	$e_{1,0,2}$
$e_{1,1,0}$	$e_{1,1,1}$	$e_{1,1,2}$
$e_{1,2,0}$	$e_{1,2,1}$	$e_{1,2,2}$

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

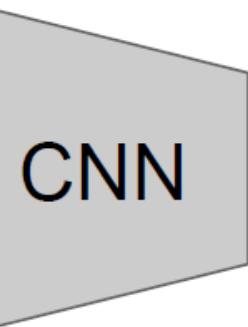
$h_0$

# Image Captioning with RNNs and Attention

Compute alignments scores (scalars):

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

$f_{att}(\cdot)$  is an MLP



Extract spatial features from a pretrained CNN

Alignment scores:

$H \times W$

$e_{1,0,0}$	$e_{1,0,1}$	$e_{1,0,2}$
$e_{1,1,0}$	$e_{1,1,1}$	$e_{1,1,2}$
$e_{1,2,0}$	$e_{1,2,1}$	$e_{1,2,2}$

Attention:

$H \times W$

$a_{1,0,0}$	$a_{1,0,1}$	$a_{1,0,2}$
$a_{1,1,0}$	$a_{1,1,1}$	$a_{1,1,2}$
$a_{1,2,0}$	$a_{1,2,1}$	$a_{1,2,2}$

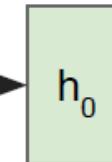
Normalize to get attention weights:

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$0 < a_{t,i,j} < 1$ ,  
attention values sum to 1

Features:  
 $H \times W \times D$

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

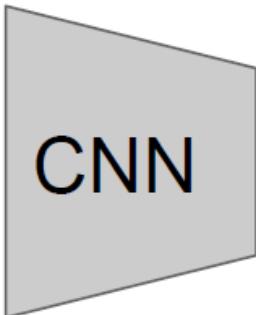


# Image Captioning with RNNs and Attention

Compute alignments scores (scalars):

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

$f_{att}(\cdot)$  is an MLP



Extract spatial features from a pretrained CNN

Alignment scores:

$$H \times W$$

$e_{1,0,0}$	$e_{1,0,1}$	$e_{1,0,2}$
$e_{1,1,0}$	$e_{1,1,1}$	$e_{1,1,2}$
$e_{1,2,0}$	$e_{1,2,1}$	$e_{1,2,2}$

Attention:

$$H \times W$$

$a_{1,0,0}$	$a_{1,0,1}$	$a_{1,0,2}$
$a_{1,1,0}$	$a_{1,1,1}$	$a_{1,1,2}$
$a_{1,2,0}$	$a_{1,2,1}$	$a_{1,2,2}$

Normalize to get attention weights:

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$0 < a_{t,i,j} < 1$ ,  
attention values sum to 1

Compute context vector:

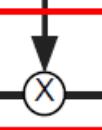
$$c_t = \sum_{i,j} a_{t,i,j} z_{t,i,j}$$

Features:  
 $H \times W \times D$

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

$$h_0$$

$$c_1$$



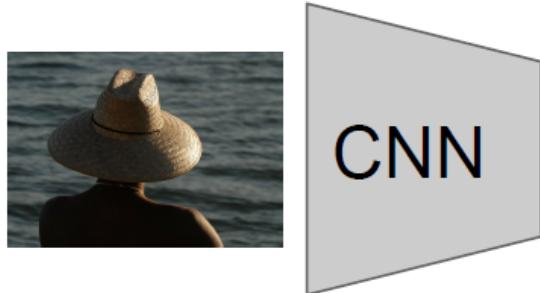
# Image Captioning with RNNs and Attention

Each timestep of decoder uses a different context vector that looks at different parts of the input image

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

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} z_{t,i,j}$$

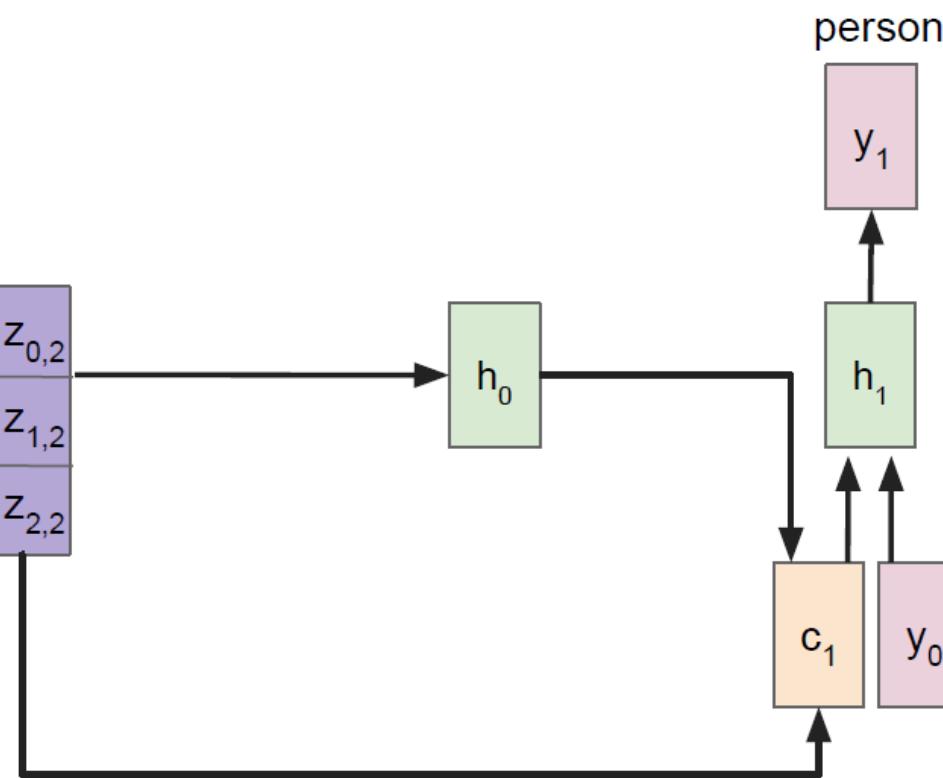


Extract spatial  
features from a  
pretrained CNN

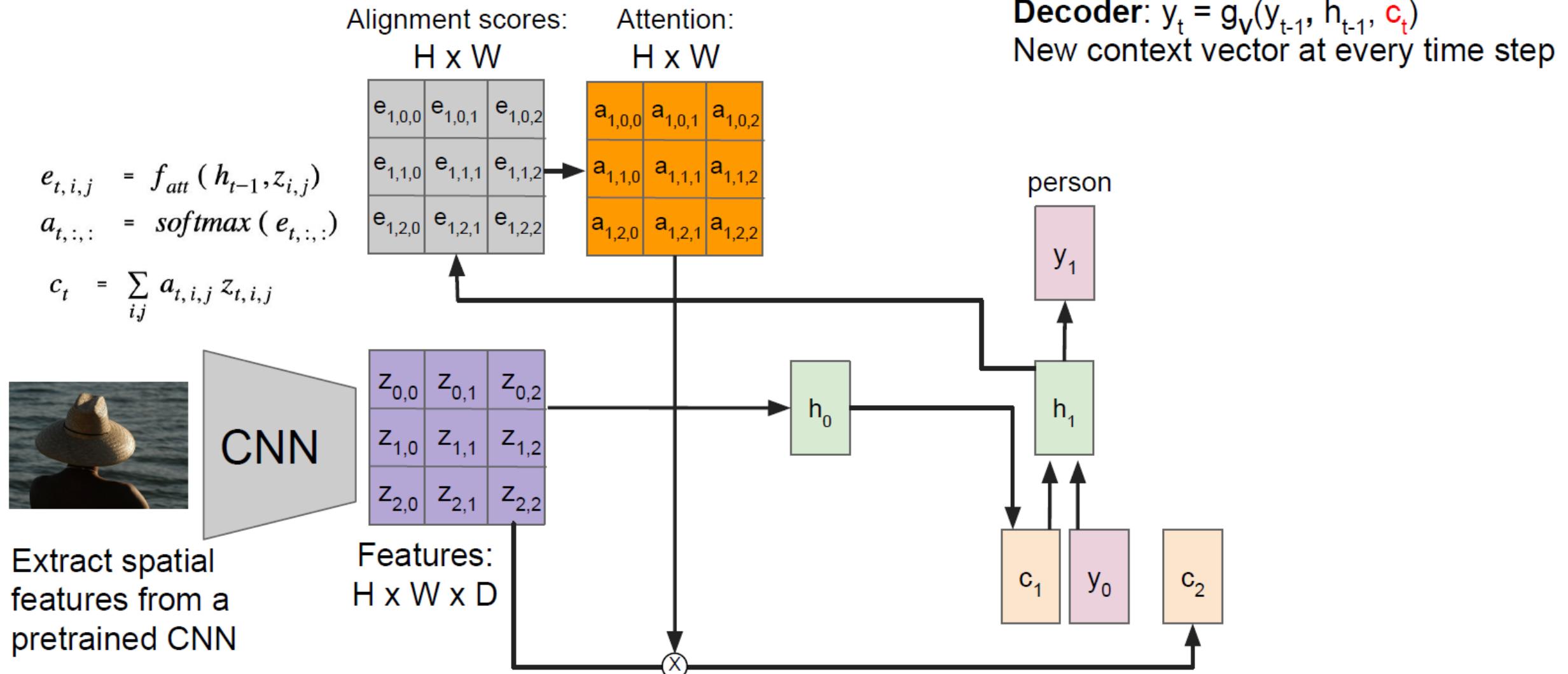
$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

**Decoder:**  $y_t = g_v(y_{t-1}, h_{t-1}, c_t)$   
New context vector at every time step



# Image Captioning with RNNs and Attention



# Image Captioning with RNNs and Attention

Each timestep of decoder uses a different context vector that looks at different parts of the input image

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

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} z_{t,i,j}$$

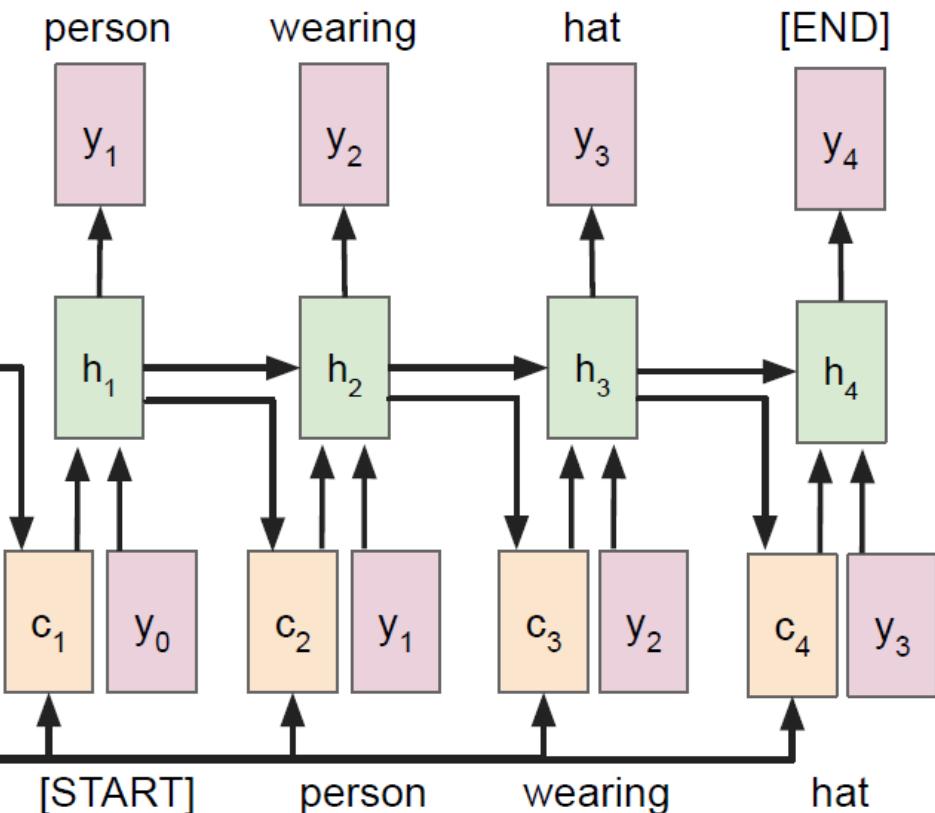


Extract spatial features from a pretrained CNN

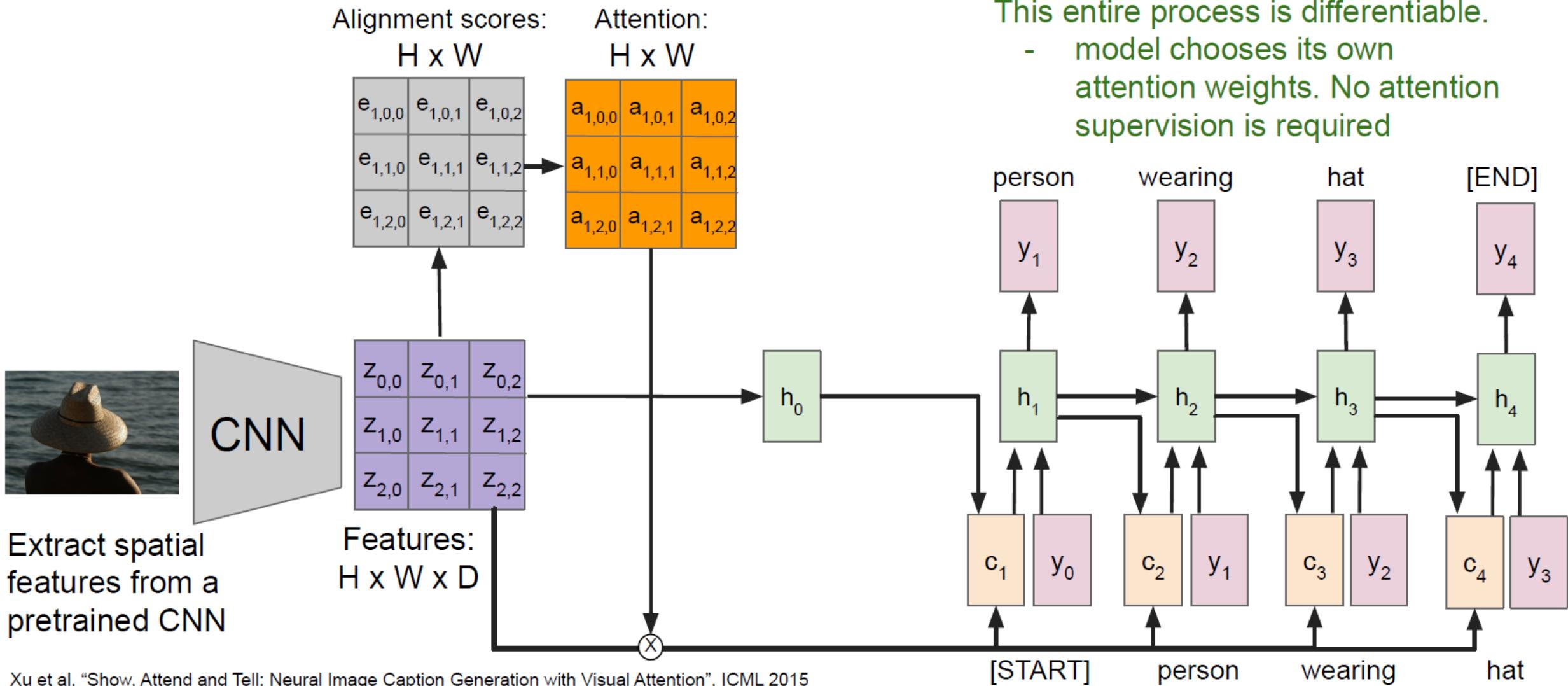
$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

**Decoder:**  $y_t = g_V(y_{t-1}, h_{t-1}, c_t)$   
New context vector at every time step



# Image Captioning with RNNs and Attention

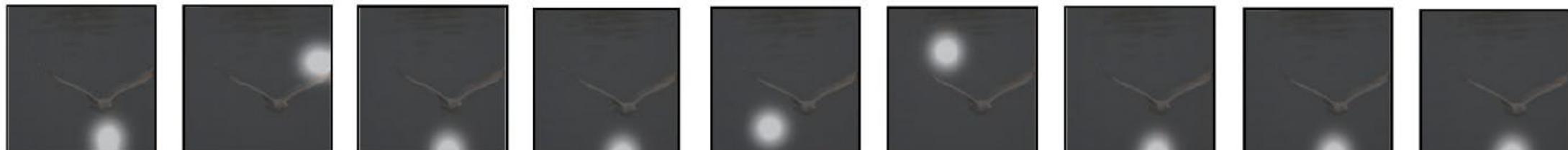


# Image Captioning with Attention

Soft attention



Hard attention  
(requires  
reinforcement  
learning)



A

bird

flying

over

a

body

of

water

.

# Image Captioning with Attention



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.



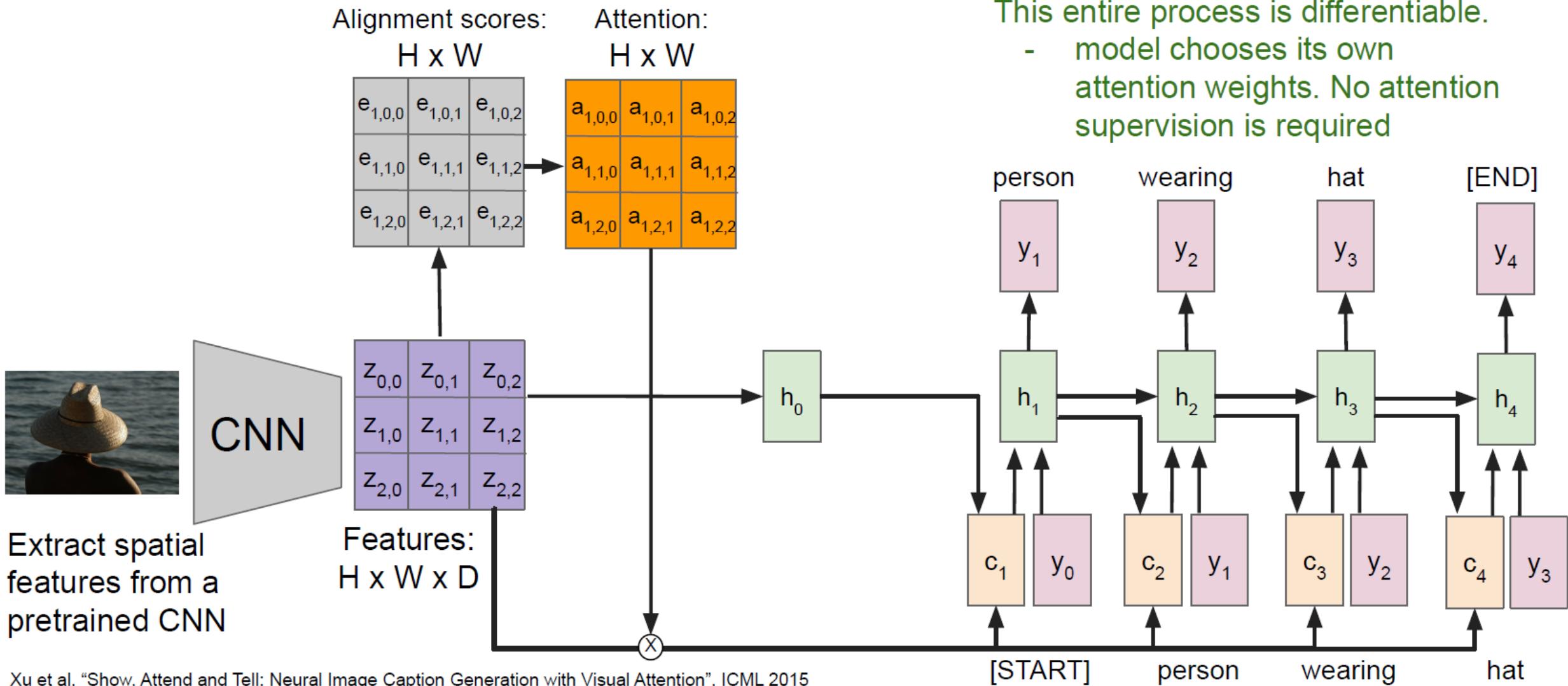
A group of people sitting on a boat in the water.



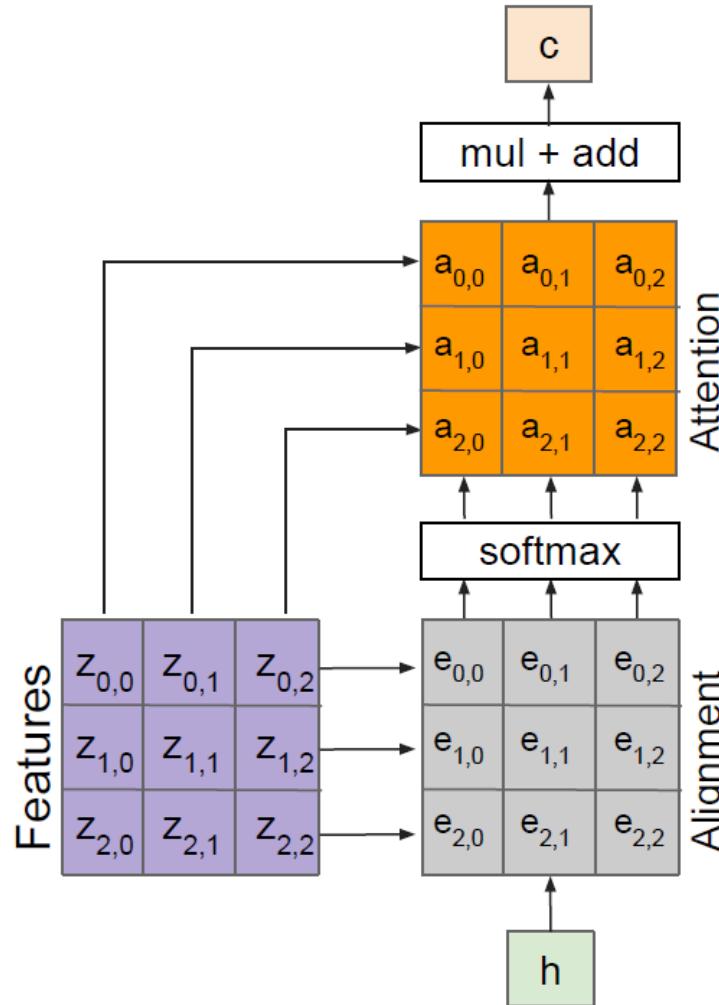
A giraffe standing in a forest with trees in the background.



# Image Captioning with RNNs and Attention



# Attention we just saw in image captioning



## Outputs:

context vector:  $\mathbf{c}$  (shape: D)

## Operations:

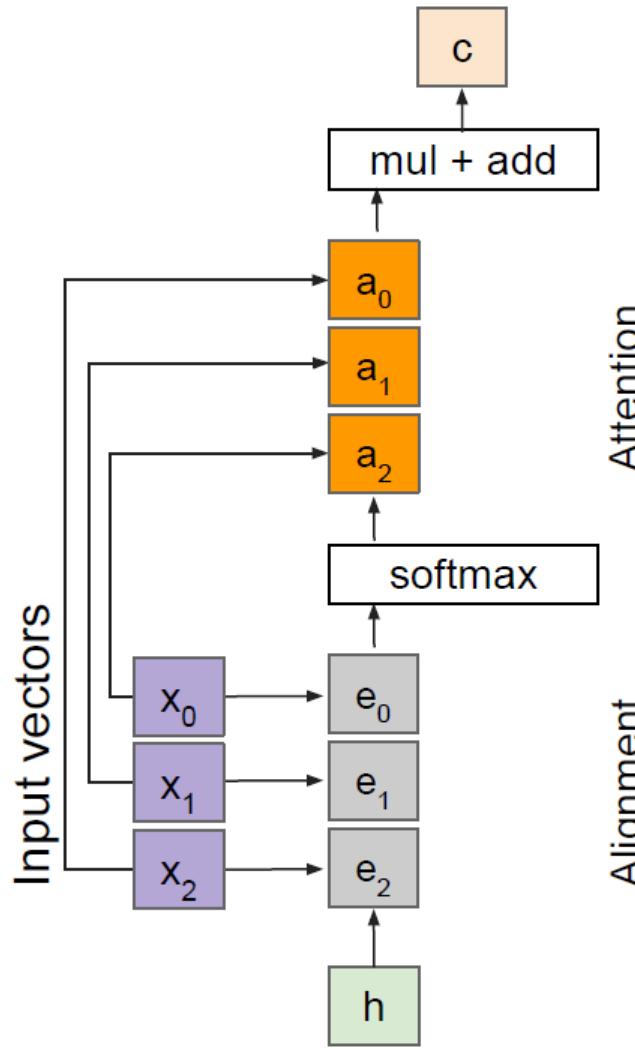
Alignment:  $e_{i,j} = f_{att}(h, z_{i,j})$   
Attention:  $\mathbf{a} = \text{softmax}(\mathbf{e})$

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

## Inputs:

Features:  $\mathbf{z}$  (shape: H x W x D)  
Query:  $\mathbf{h}$  (shape: D)

# General attention layer



## Outputs:

context vector:  $\mathbf{c}$  (shape: D)

## Operations:

Alignment:  $e_i = f_{att}(h, x_i)$

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

Output:  $\mathbf{c} = \sum_i a_i x_i$

Attention operation is **permutation invariant**.

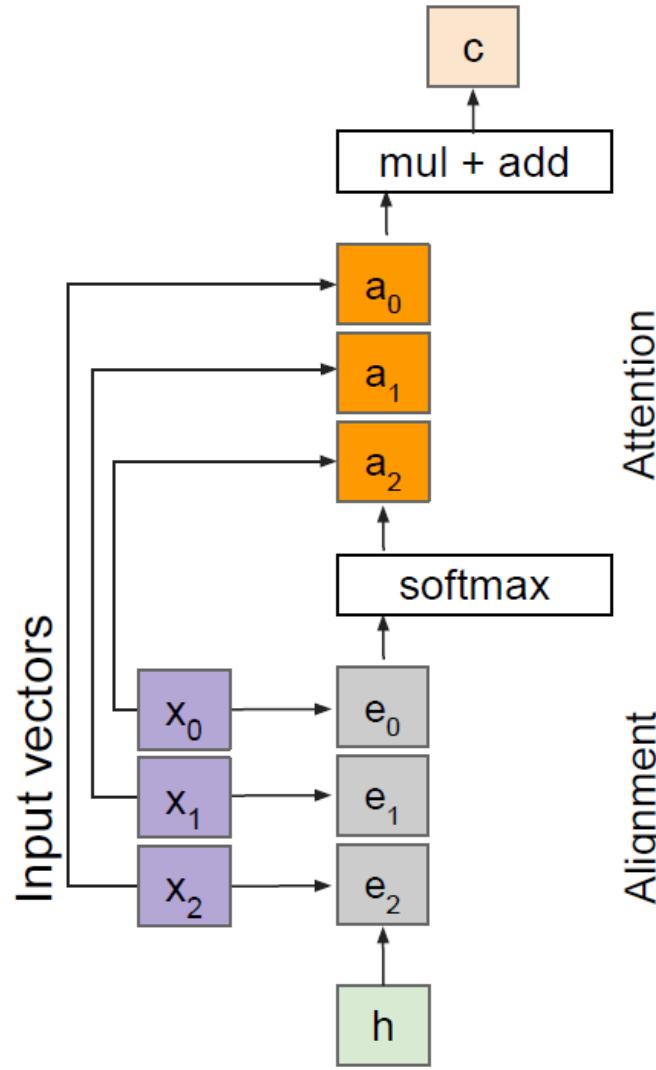
- Doesn't care about ordering of the features
- Stretch  $H \times W = N$  into N vectors

## Inputs:

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

Query:  $\mathbf{h}$  (shape: D)

# General attention layer



## Outputs:

context vector:  $\mathbf{c}$  (shape: D)

## Operations:

$$\text{Alignment: } \mathbf{e}_i = \mathbf{h} \cdot \mathbf{x}_i$$

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

$$\text{Output: } \mathbf{c} = \sum_i \mathbf{a}_i \mathbf{x}_i$$

Change  $f_{\text{att}}(\cdot)$  to a simple dot product

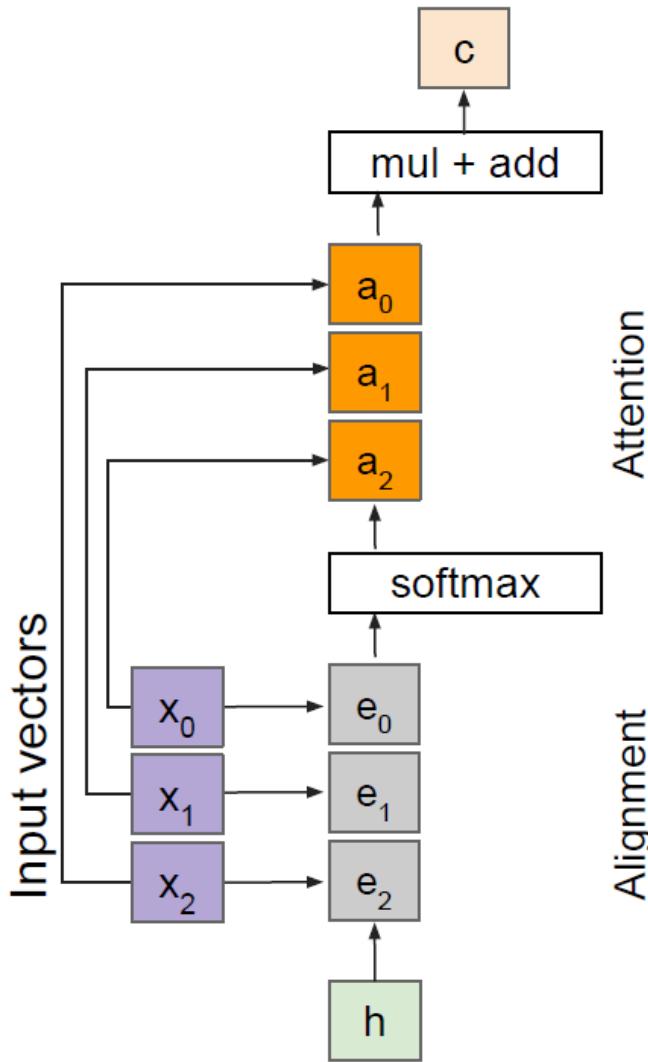
- only works well with key & value transformation trick (will mention in a few slides)

## Inputs:

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

Query:  $\mathbf{h}$  (shape: D)

# General attention layer



## Outputs:

context vector:  $\mathbf{c}$  (shape: D)

## Operations:

Alignment:  $e_i = \mathbf{h} \cdot \mathbf{x}_i / \sqrt{D}$

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

Output:  $\mathbf{c} = \sum_i a_i \mathbf{x}_i$

## Inputs:

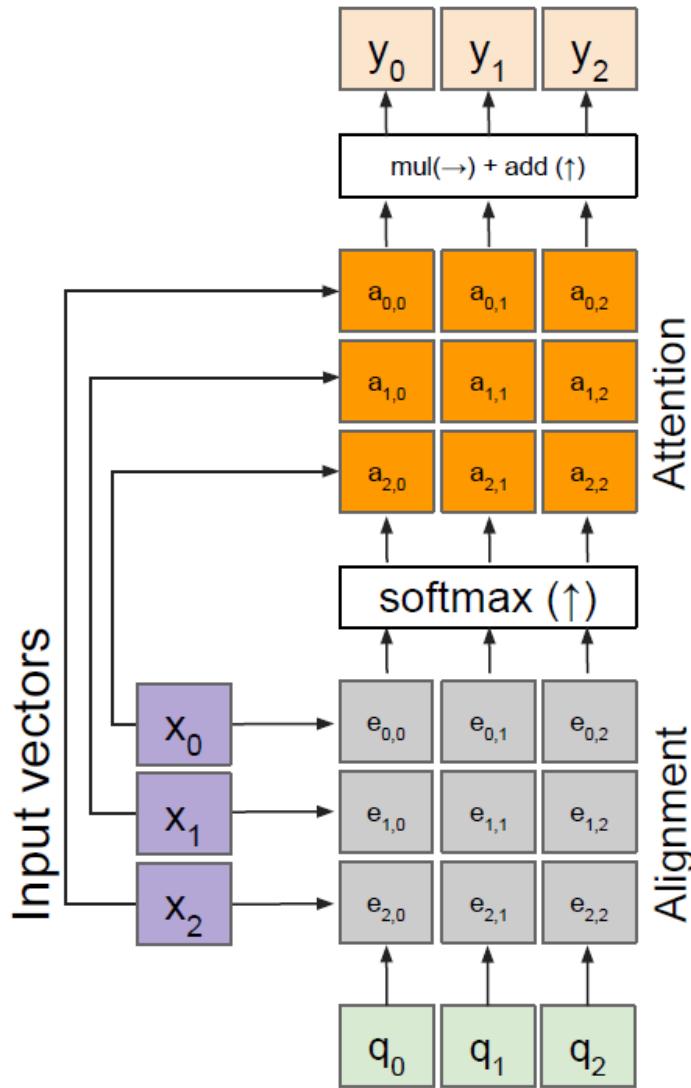
Input vectors:  $\mathbf{x}$  (shape: N x D)

Query:  $\mathbf{h}$  (shape: D)

Change  $f_{\text{att}}(\cdot)$  to a **scaled** simple dot product

- Larger dimensions means more terms in the dot product sum.
- So, the variance of the logits is higher. Large magnitude vectors will produce much higher logits.
- So, the post-softmax distribution has lower-entropy, assuming logits are IID.
- Ultimately, these large magnitude vectors will cause softmax to peak and assign very little weight to all others
- Divide by  $\sqrt{D}$  to reduce effect of large magnitude vectors

# General attention layer



## Outputs:

context vectors:  $\mathbf{y}$  (shape: D)

## Multiple query vectors

- each query creates a new output context vector

## Operations:

Alignment:  $e_{i,j} = q_j \cdot x_i / \sqrt{D}$

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

Output:  $y_j = \sum_i a_{i,j} x_i$

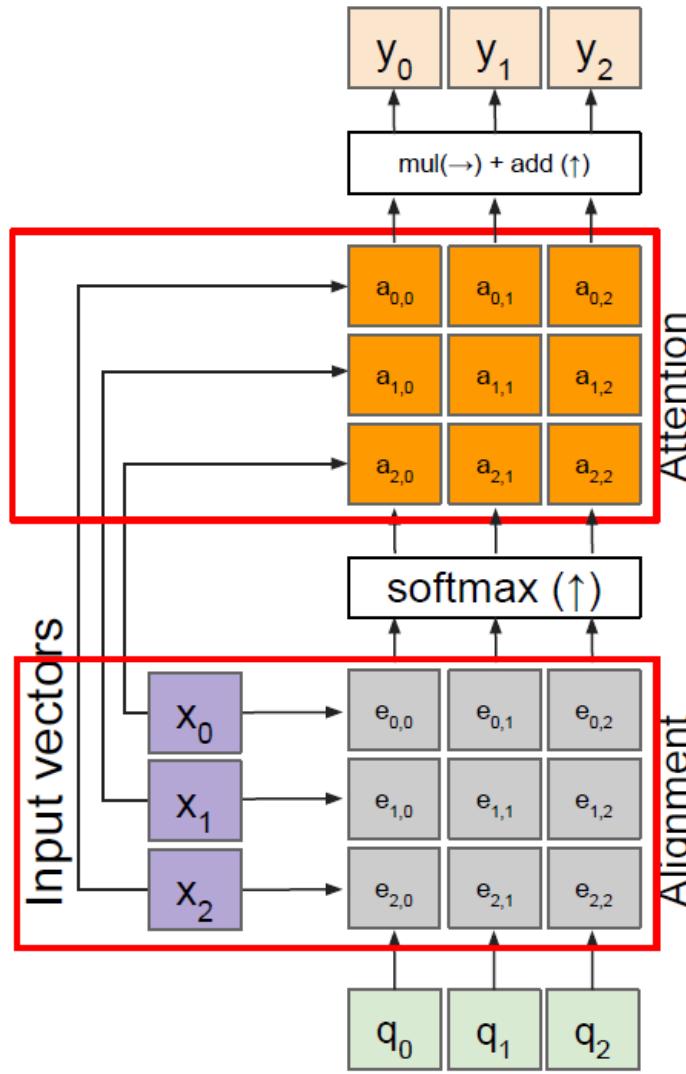
## Multiple query vectors

## Inputs:

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

Queries:  $\mathbf{q}$  (shape: M x D)

# General attention layer



**Outputs:**

context vectors:  $\mathbf{y}$  (shape: D)

**Operations:**

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

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

Output:  $y_j = \sum_i a_{i,j} \mathbf{x}_i$

Notice that the input vectors are used for both the alignment as well as the attention calculations.

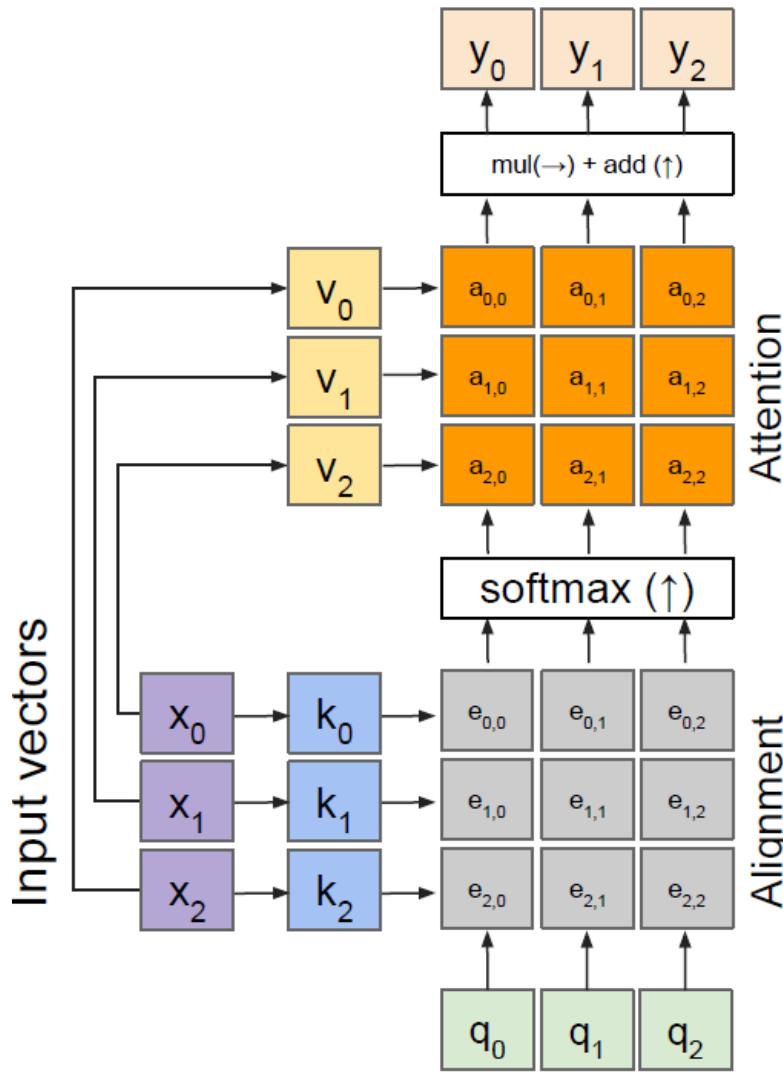
- We can add more expressivity to the layer by adding a different FC layer before each of the two steps.

**Inputs:**

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

Queries:  $\mathbf{q}$  (shape: M x D)

# General attention layer



## Outputs:

context vectors:  $\mathbf{y}$  (shape:  $D_y$ )

The input and output dimensions can now change depending on the key and value FC layers

## Operations:

Key vectors:  $\mathbf{k} = \mathbf{x}W_k$

Value vectors:  $\mathbf{v} = \mathbf{x}W_v$

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

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

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

Notice that the input vectors are used for both the alignment as well as the attention calculations.

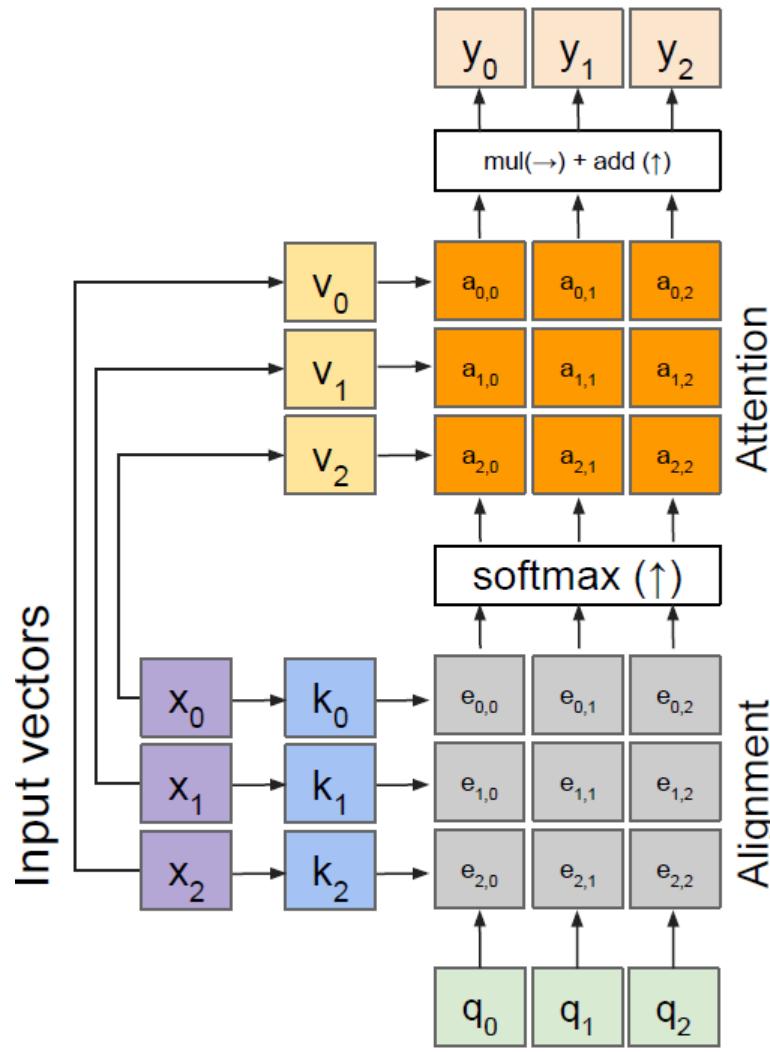
- We can add more expressivity to the layer by adding a different FC layer before each of the two steps.

## Inputs:

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

Queries:  $\mathbf{q}$  (shape:  $M \times D_k$ )

# General attention layer



## Outputs:

context vectors:  $\mathbf{y}$  (shape:  $D_v$ )

## Operations:

Key vectors:  $\mathbf{k} = \mathbf{x}W_k$

Value vectors:  $\mathbf{v} = \mathbf{x}W_v$

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

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

Output:  $y_j = \sum_i a_{i,j} v_i$

## Inputs:

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

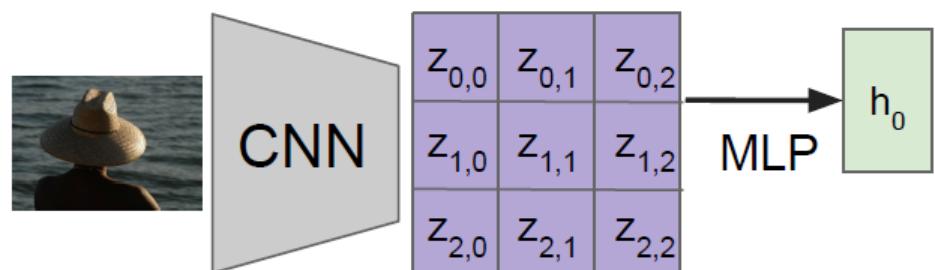
Queries:  $\mathbf{q}$  (shape:  $M \times D_k$ )

Recall that the query vector was a function of the input vectors

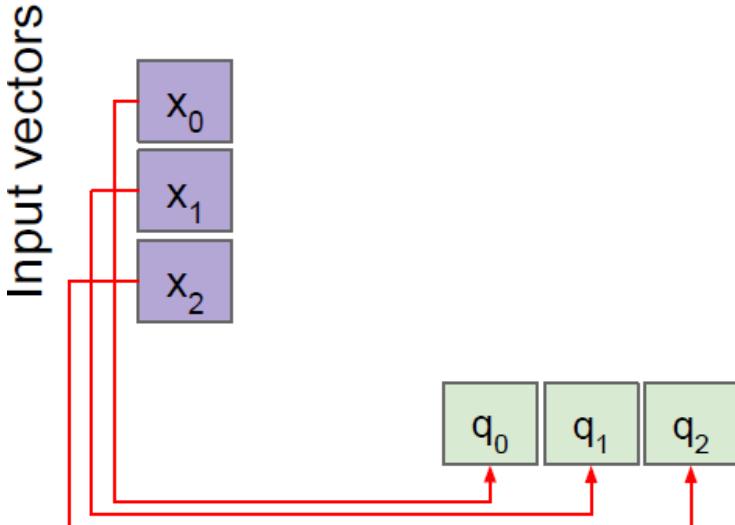
**Encoder:**  $h_0 = f_w(\mathbf{z})$

where  $\mathbf{z}$  is spatial CNN features

$f_w(\cdot)$  is an MLP



# Self attention layer



## Operations:

Key vectors:  $\mathbf{k} = \mathbf{x}W_k$

Value vectors:  $\mathbf{v} = \mathbf{x}W_v$

Query vectors:  $\mathbf{q} = \mathbf{x}W_q$

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

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

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

We can calculate the query vectors from the input vectors, therefore, defining a "self-attention" layer.

Instead, query vectors are calculated using a FC layer.

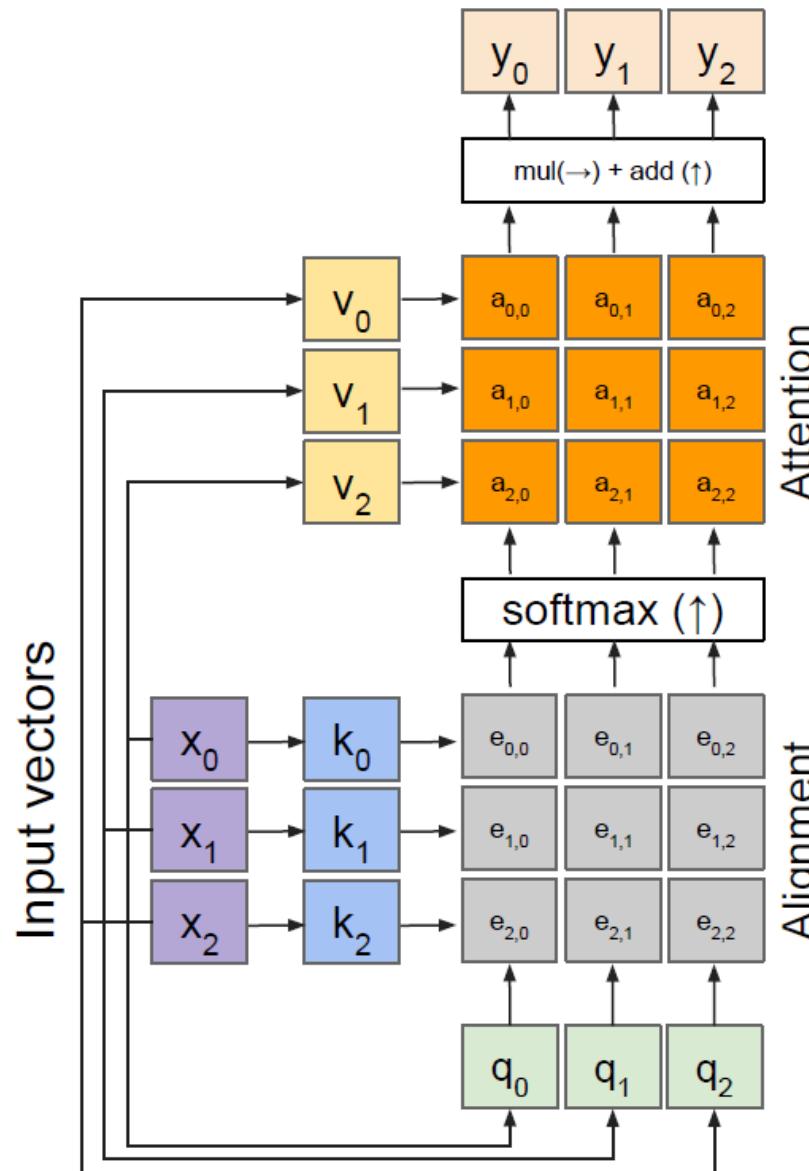
## Inputs:

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

Queries:  $\mathbf{q}$  (shape:  $M \times D_k$ )

No input query vectors anymore

# Self attention layer



## Outputs:

context vectors:  $\mathbf{y}$  (shape:  $D_v$ )

## Operations:

Key vectors:  $\mathbf{k} = \mathbf{x}\mathbf{W}_k$

Value vectors:  $\mathbf{v} = \mathbf{x}\mathbf{W}_v$

Query vectors:  $\mathbf{q} = \mathbf{x}\mathbf{W}_q$

Alignment:  $\mathbf{e}_{i,j} = \mathbf{q}_j \cdot \mathbf{k}_i / \sqrt{D}$

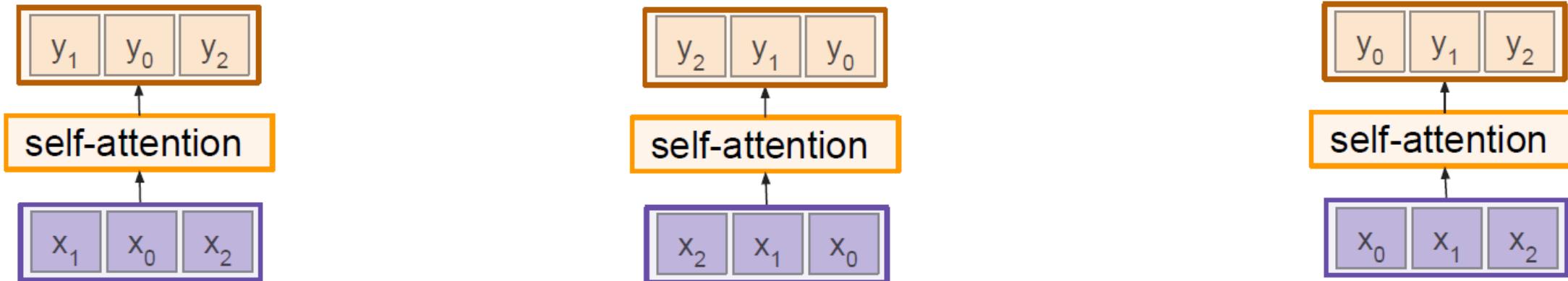
Attention:  $\mathbf{a} = \text{softmax}(\mathbf{e})$

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

## Inputs:

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

# Self attention layer - attends over sets of inputs

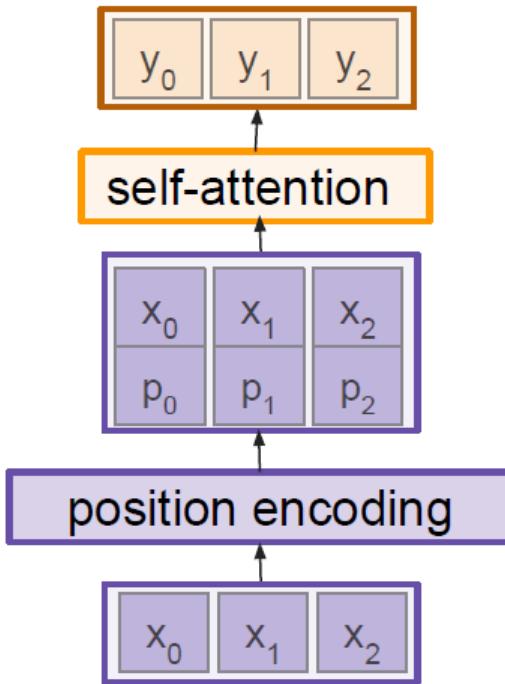


Permutation equivariant

Self-attention layer doesn't care about the orders of the inputs!

**Problem:** How can we encode ordered sequences like language or spatially ordered image features?

# Positional encoding



Concatenate/add special positional encoding  $p_j$  to each input vector  $x_j$

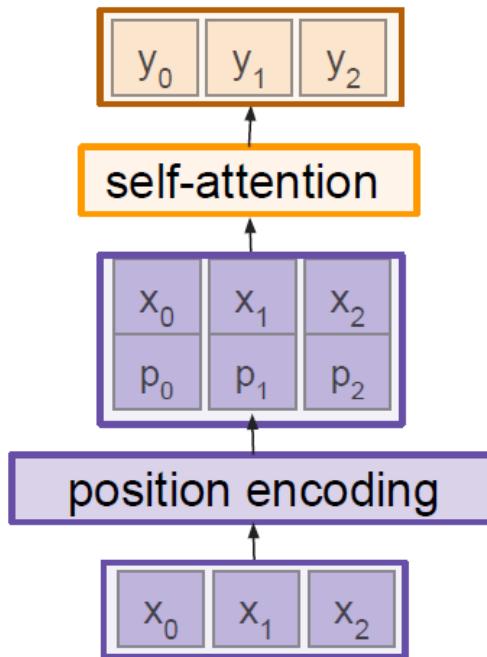
We use a function  $pos: N \rightarrow \mathbb{R}^d$  to process the position  $j$  of the vector into a  $d$ -dimensional vector

$$\text{So, } p_j = pos(j)$$

Desiderata of  $pos(\cdot)$  :

1. It should output a **unique** encoding for each time-step (word's position in a sentence)
2. **Distance** between any two time-steps should be consistent across sentences with different lengths.
3. Our model should generalize to **longer** sentences without any efforts. Its values should be bounded.
4. It must be **deterministic**.

# Positional encoding



Concatenate special positional encoding  $p_j$  to each input vector  $x_j$

We use a function  $pos: N \rightarrow \mathbb{R}^d$  to process the position  $j$  of the vector into a  $d$ -dimensional vector

So,  $p_j = pos(j)$

Options for  $pos(\cdot)$

1. Learn a lookup table:
  - o Learn parameters to use for  $pos(t)$  for  $t \in [0, T)$
  - o Lookup table contains  $T \times d$  parameters.
2. Design a fixed function with the desiderata

$$p(t) = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_d$$

Intuition:

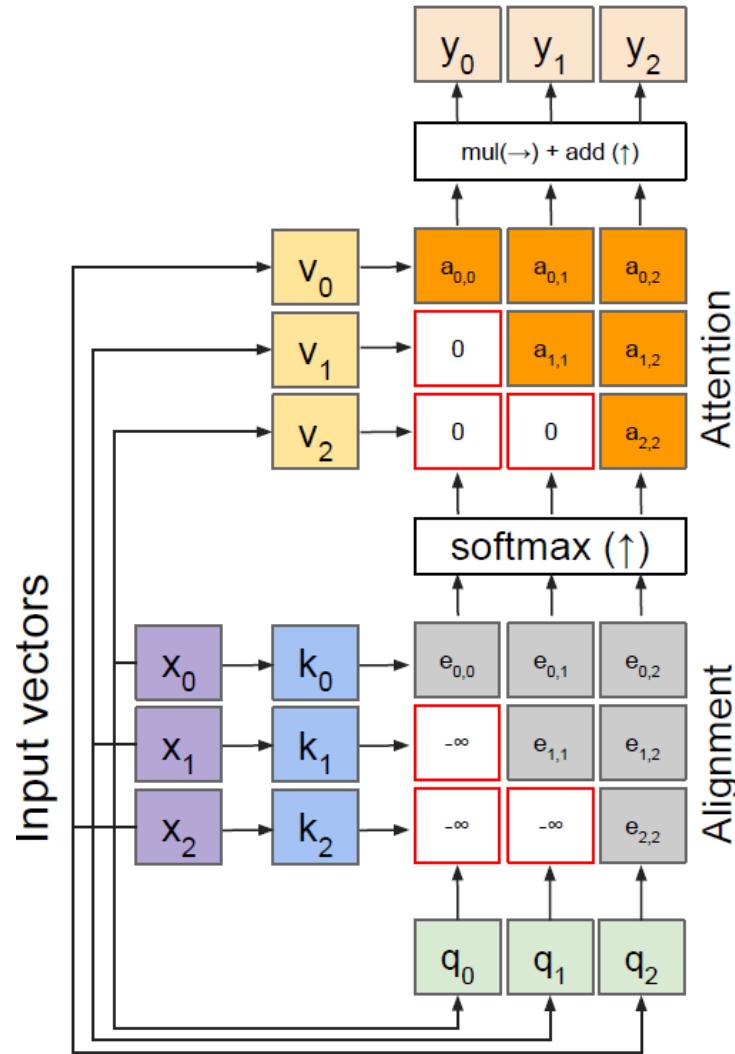
0 :	0	0	0	0	8 :	1	0	0	0
1 :	0	0	0	1	9 :	1	0	0	1
2 :	0	0	1	0	10 :	1	0	1	0
3 :	0	0	1	1	11 :	1	0	1	1
4 :	0	1	0	0	12 :	1	1	0	0
5 :	0	1	0	1	13 :	1	1	0	1
6 :	0	1	1	0	14 :	1	1	1	0
7 :	0	1	1	1	15 :	1	1	1	1

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

image source

Vaswani et al, "Attention is all you need", NeurIPS 2017

# Masked self-attention layer



## Outputs:

context vectors:  $\mathbf{y}$  (shape:  $D_v$ )

## Operations:

Key vectors:  $\mathbf{k} = \mathbf{x}W_k$

Value vectors:  $\mathbf{v} = \mathbf{x}W_v$

Query vectors:  $\mathbf{q} = \mathbf{x}W_q$

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

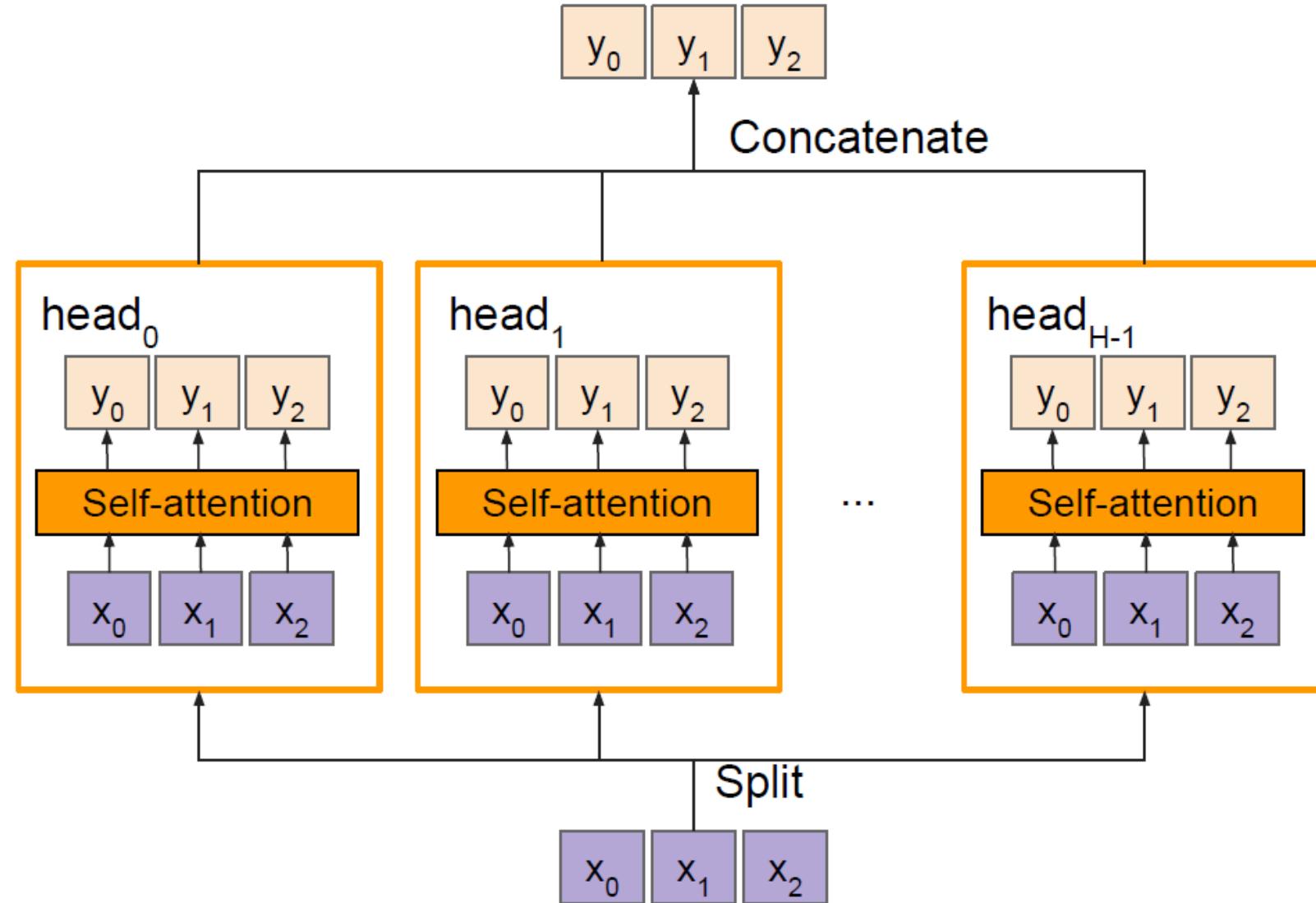
Attention:  $\mathbf{a} = \text{softmax}(\mathbf{e})$

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

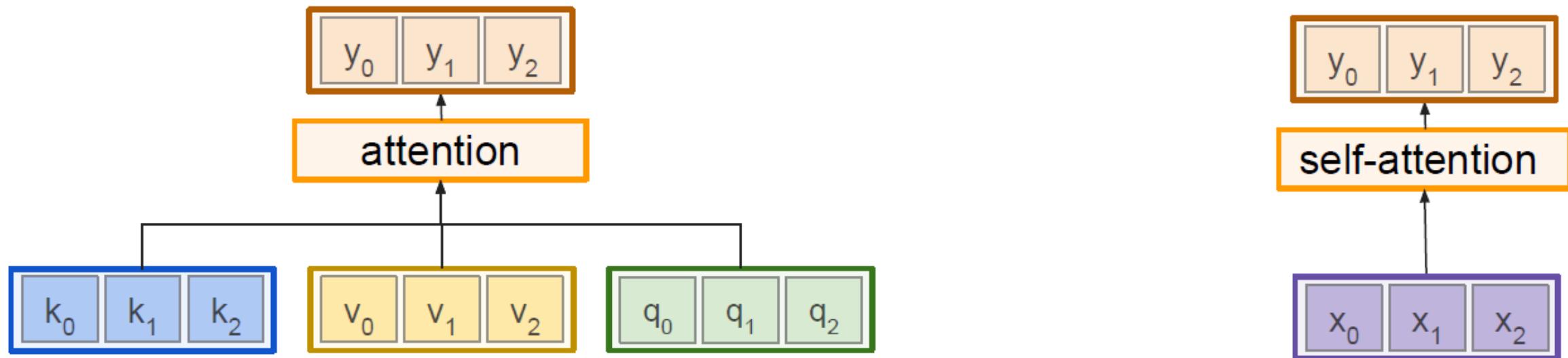
- Prevent vectors from looking at future vectors.
- Manually set alignment scores to -infinity

# Multi-head self attention layer

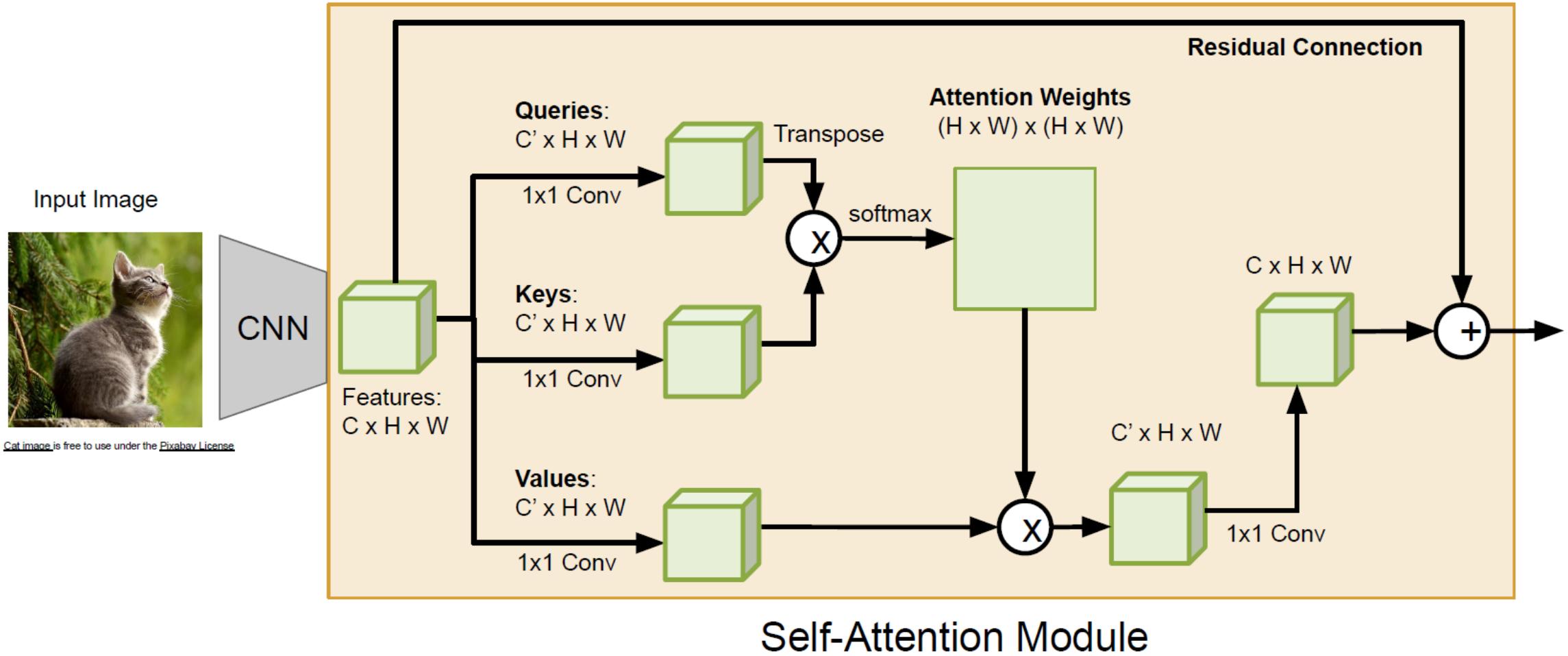
- Multiple self-attention heads in parallel



# General attention versus self-attention



# Example: CNN with Self-Attention



# Comparing RNNs to Transformer

## RNNs

- (+) LSTMs work reasonably well for long sequences.
- (-) Expects an ordered sequences of inputs
- (-) Sequential computation: subsequent hidden states can only be computed after the previous ones are done.

## Transformer:

- (+) Good at long sequences. Each attention calculation looks at all inputs.
- (+) Can operate over unordered sets or ordered sequences with positional encodings.
- (+) Parallel computation: All alignment and attention scores for all inputs can be done in parallel.
- (-) Requires a lot of memory:  $N \times M$  alignment and attention scalers need to be calculated and stored for a single self-attention head. (but GPUs are getting bigger and better)

# Image Captioning using Transformers

**Input:** Image  $I$

**Output:** Sequence  $y = y_1, y_2, \dots, y_T$

**Encoder:**  $c = T_w(z)$

where  $z$  is spatial CNN features

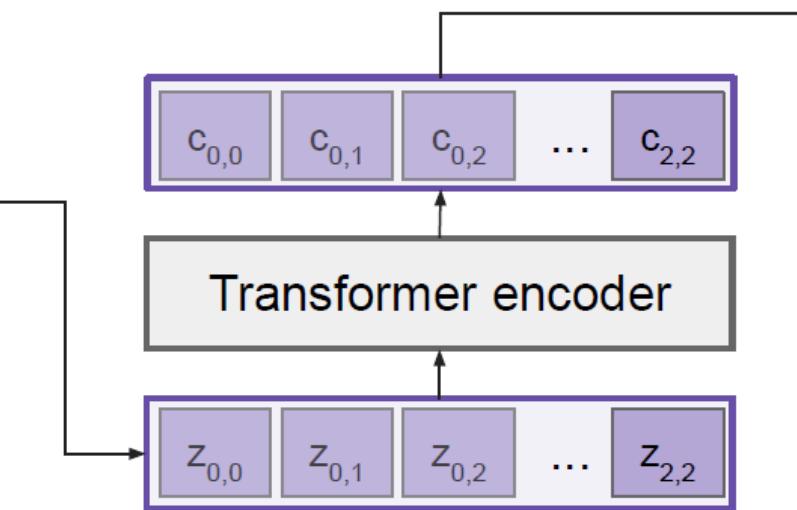
$T_w(\cdot)$  is the transformer encoder



Extract spatial  
features from a  
pretrained CNN

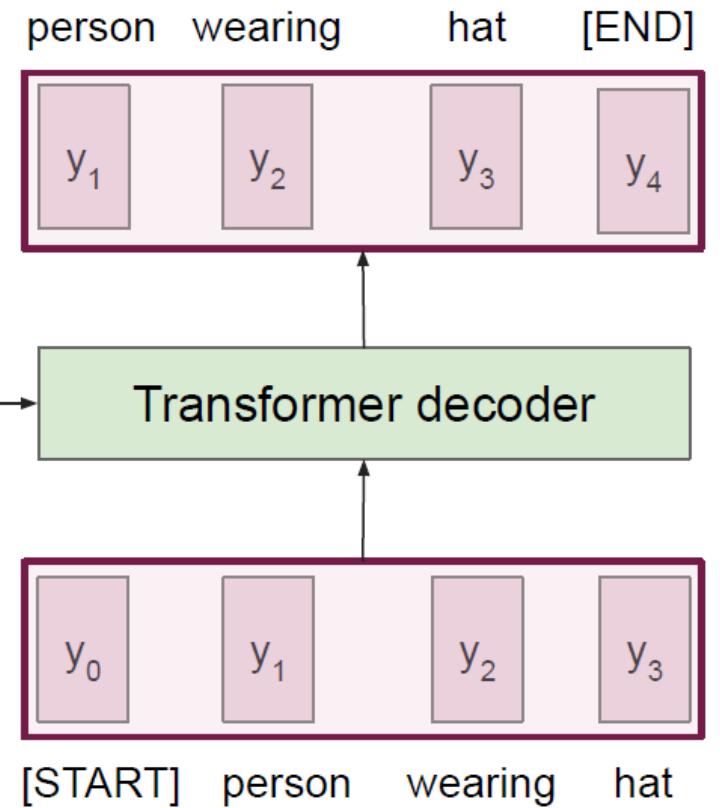
$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

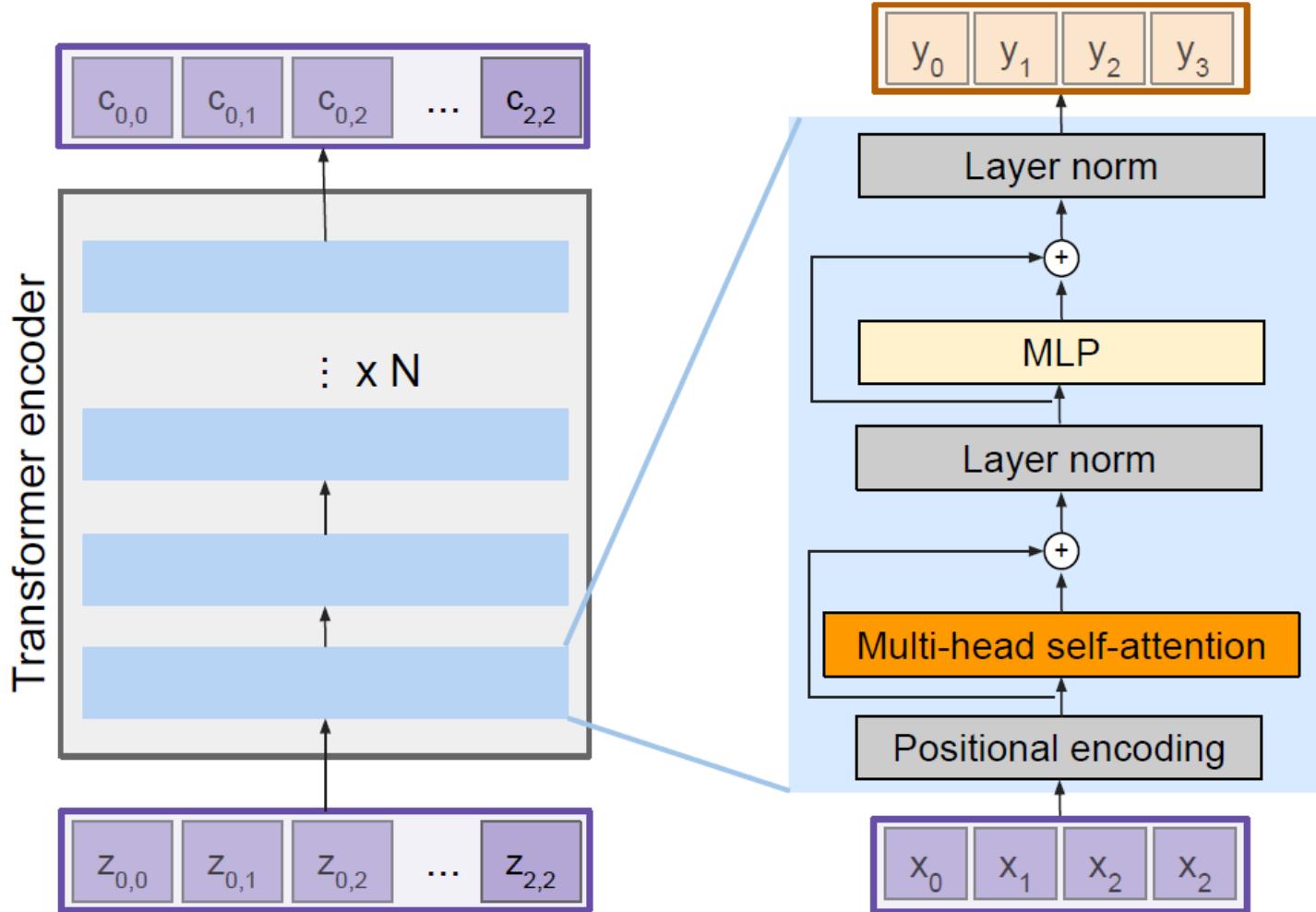


**Decoder:**  $y_t = T_D(y_{0:t-1}, c)$

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



# The Transformer encoder block



**Transformer Encoder Block:**

**Inputs:** Set of vectors  $x$

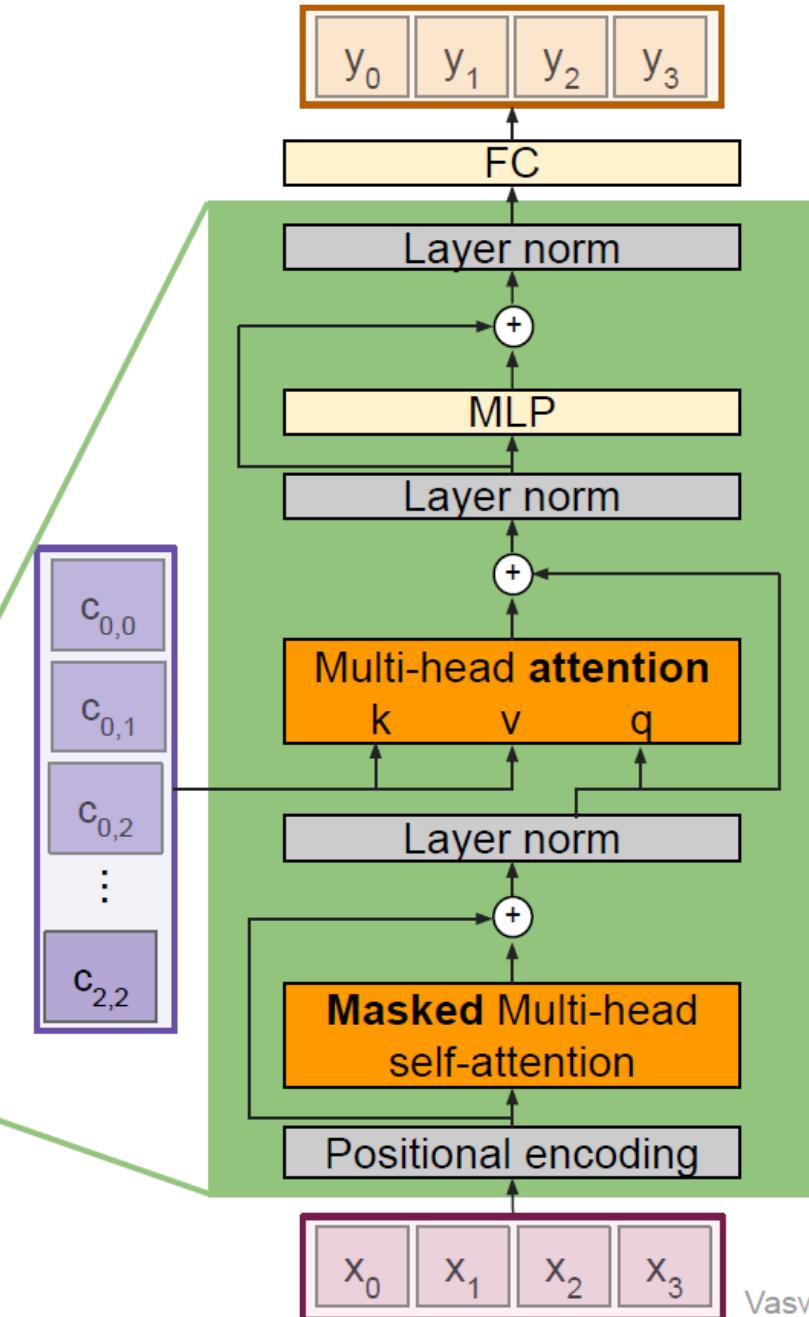
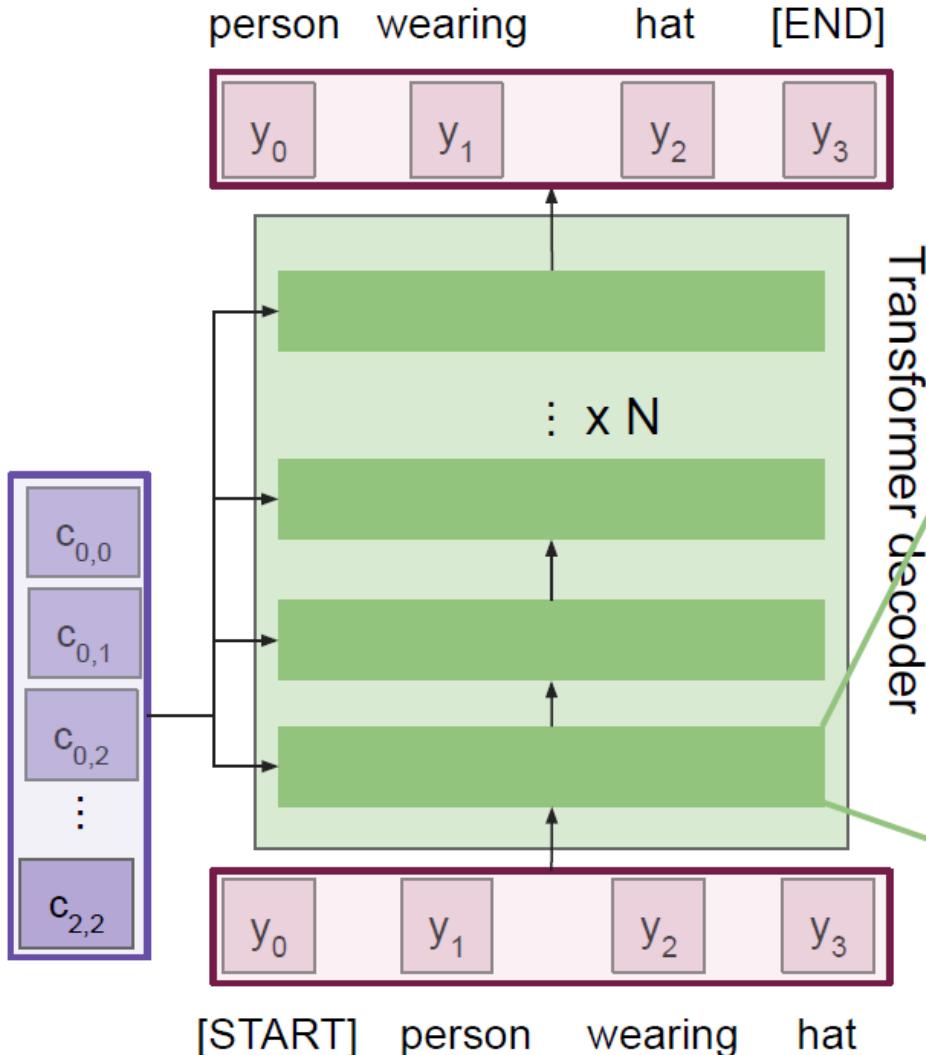
**Outputs:** Set of vectors  $y$

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 Decoder block



Transformer Decoder Block:

**Inputs:** Set of vectors  $x$  and Set of context vectors  $c$ .

**Outputs:** Set of vectors  $y$ .

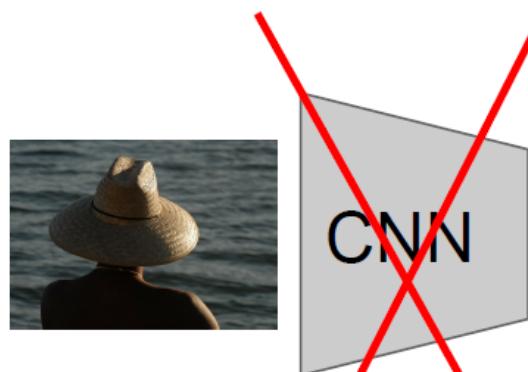
Masked Self-attention only interacts with past inputs.

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

Highly scalable, highly parallelizable, but high memory usage.

# Image Captioning using transformers

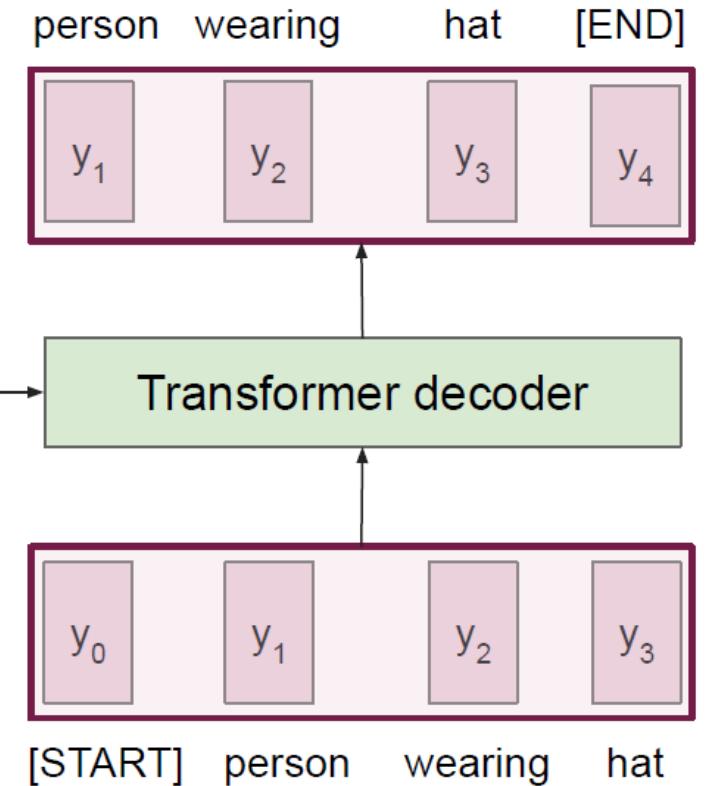
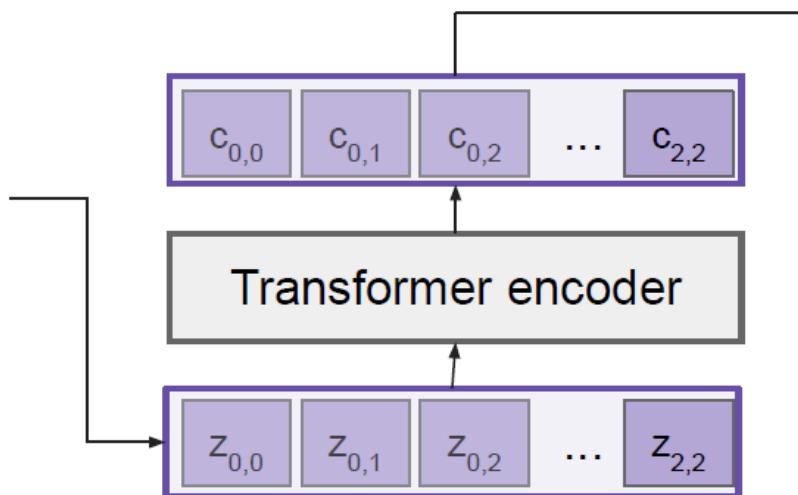
- Perhaps we don't need convolutions at all?



Extract spatial  
features from a  
pretrained CNN

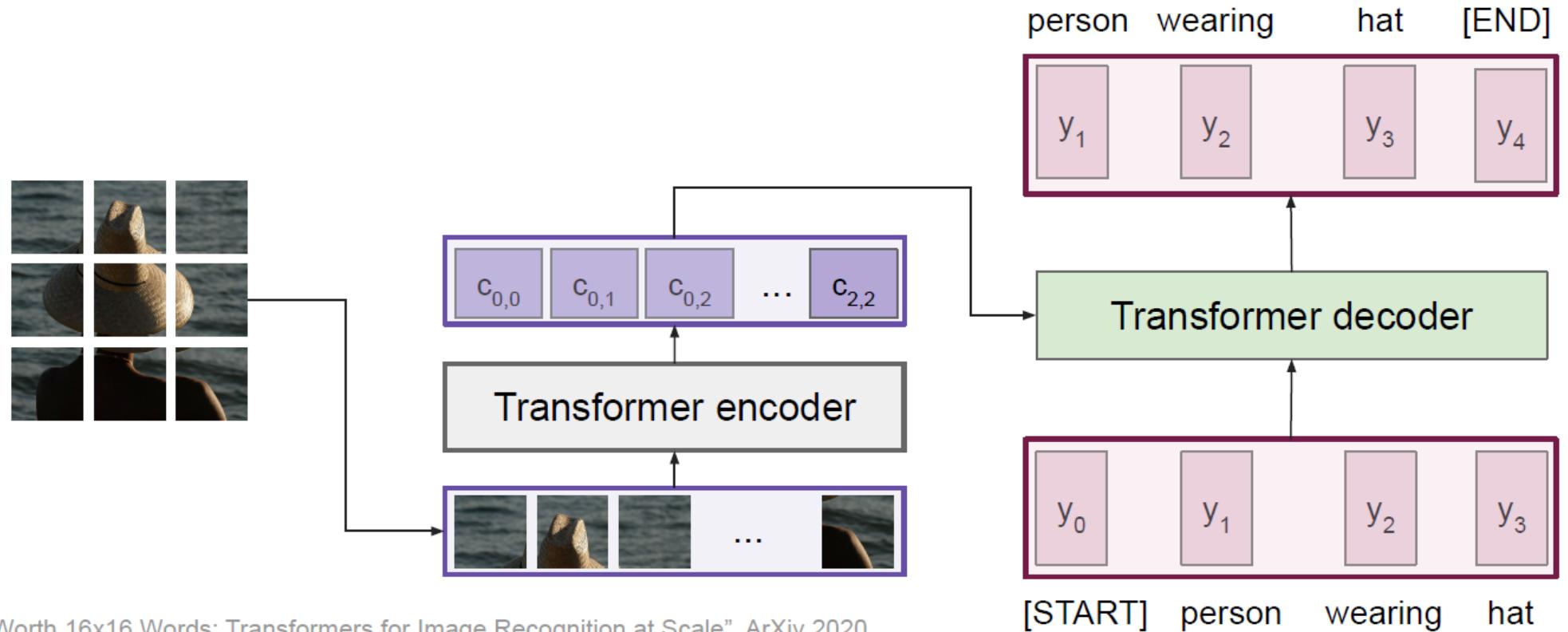
$Z_{0,0}$	$Z_{0,1}$	$Z_{0,2}$
$Z_{1,0}$	$Z_{1,1}$	$Z_{1,2}$
$Z_{2,0}$	$Z_{2,1}$	$Z_{2,2}$

Features:  
 $H \times W \times D$



# Image Captioning using ONLY transformers

- Transformers from pixels to language



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ArXiv 2020  
[Colab link](#) to an implementation of vision transformers