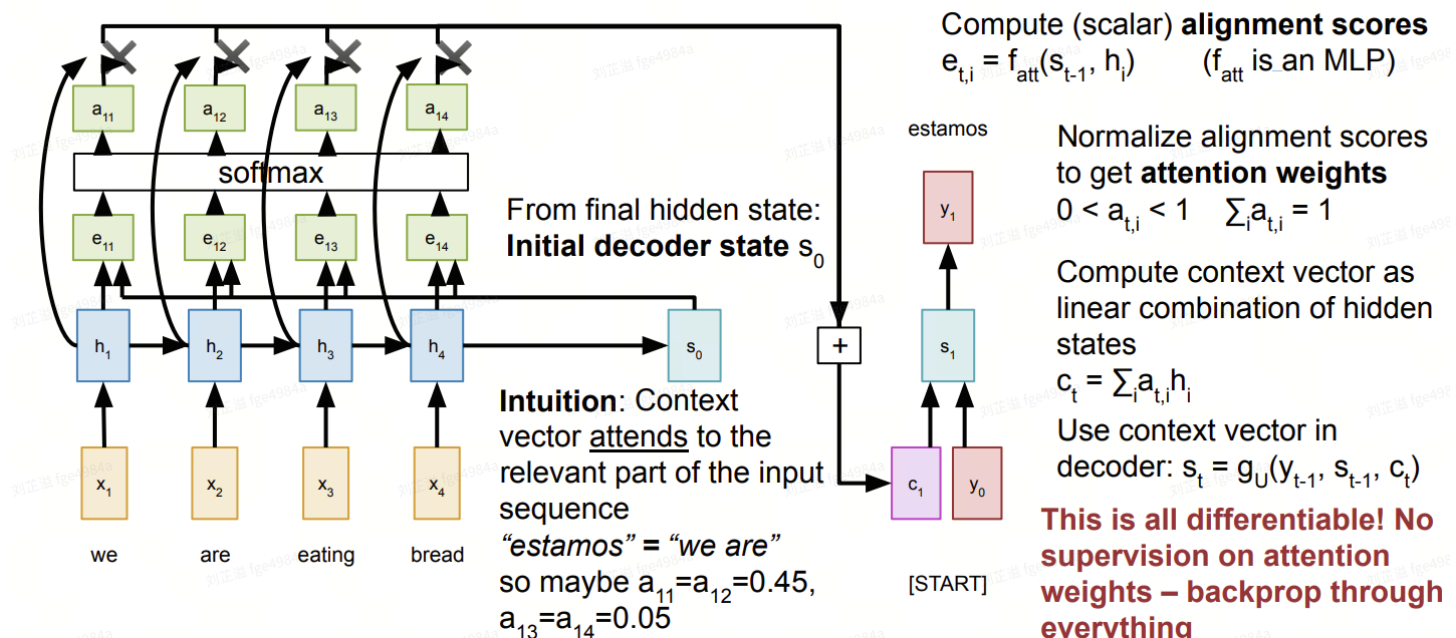


Attention and Transformers

Attention

- encoder中设置上下文向量
- decoder中每个时间步的上下文向量都不同

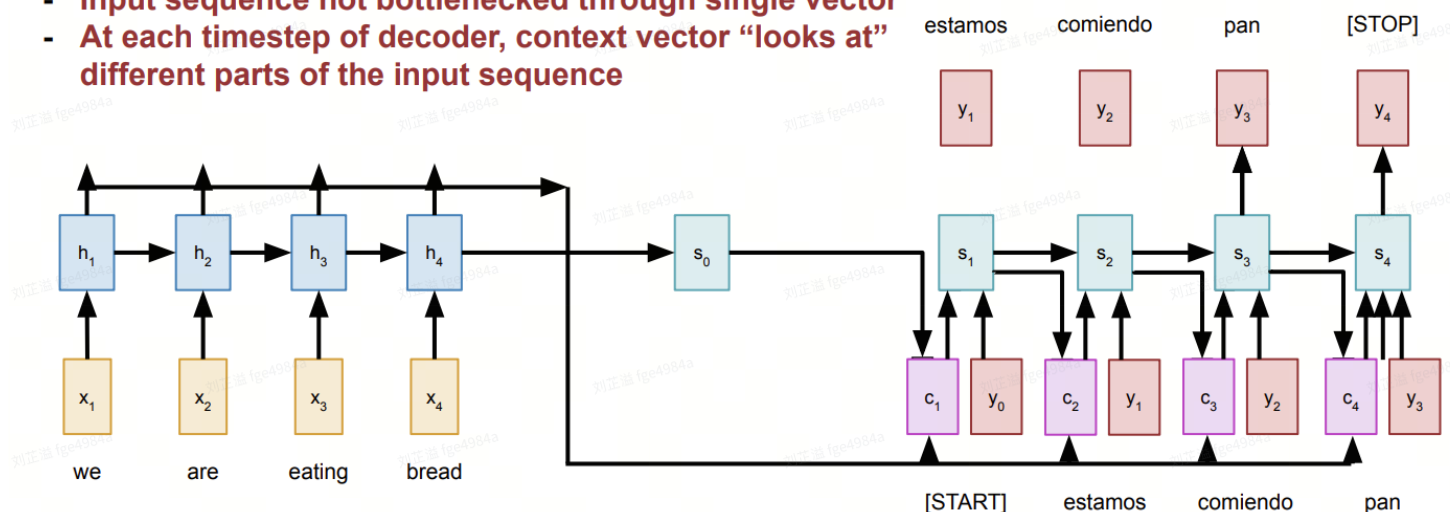
Sequence to Sequence with RNNs and Attention



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



CNN + Attention

CNN + RNN

- Input is "bottlenecked" through c ; because model needs to encode everything with c

Image Captioning using spatial features

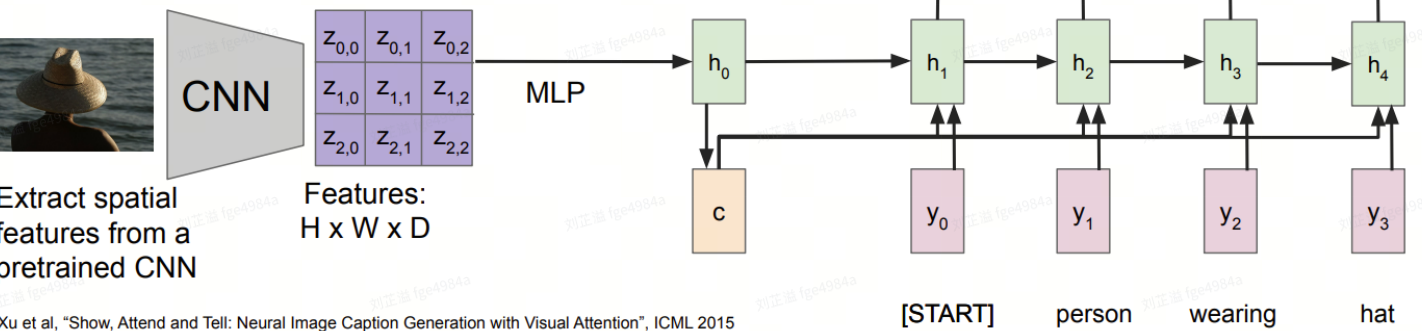
Input: Image I

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

Encoder: $h_0 = f_w(z)$

where z is spatial CNN features

$f_w(\cdot)$ is an MLP



CNN + Attention

- encoder每个步的上下文变量都不同，会重新根据注意力机制计算
- 模型自己选择注意力权重，不需要注意力监督

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

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:

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

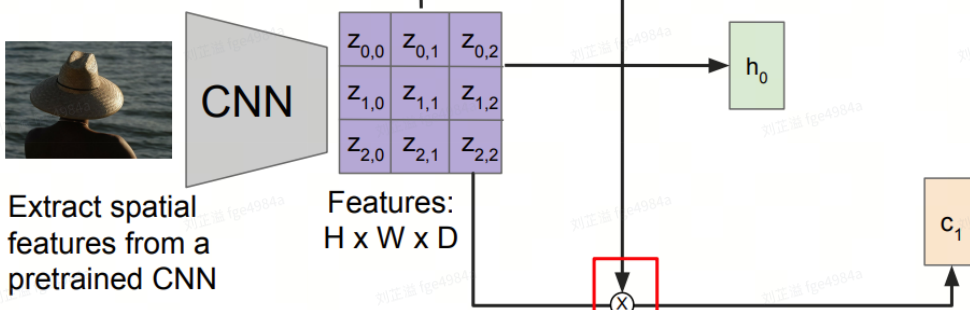
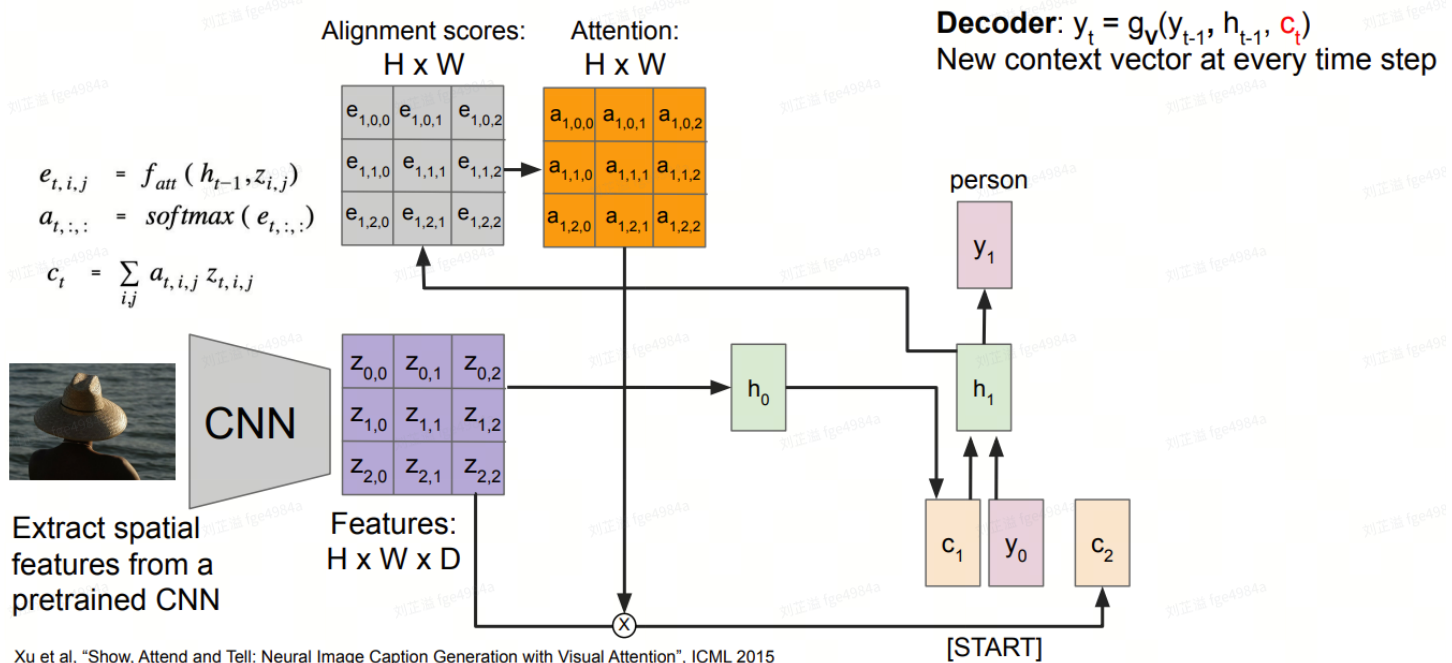


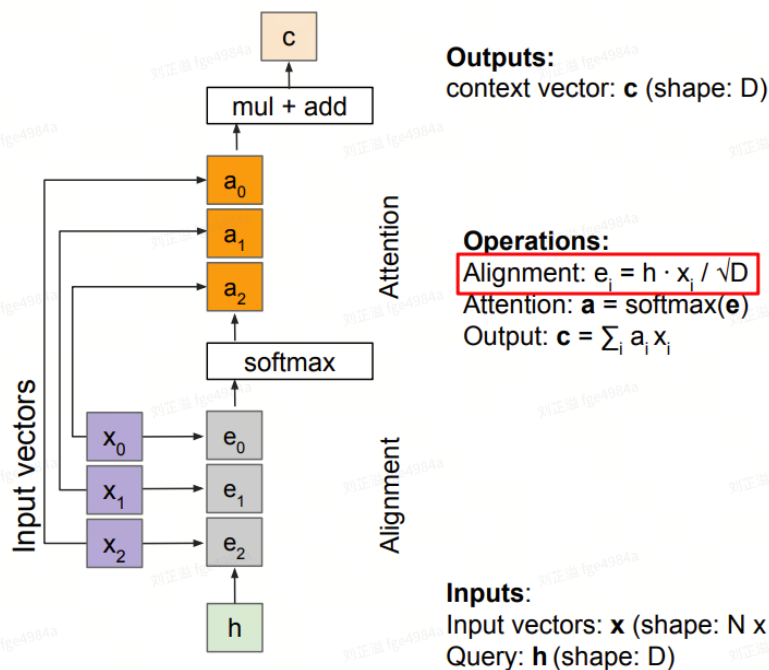
Image Captioning with RNNs and Attention



Attention Layer

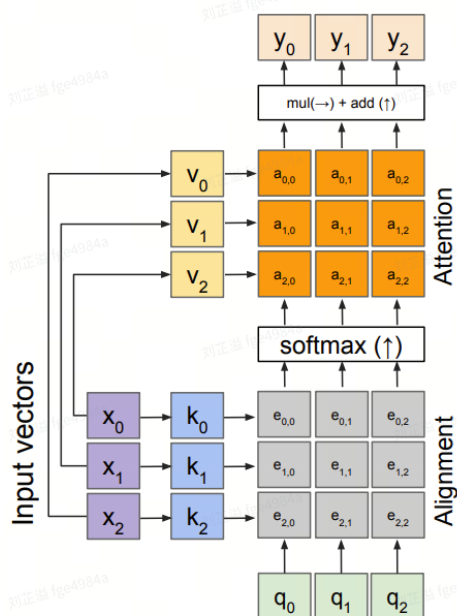
- 通过除以 $D^{*0.5}$ ，消除大幅度向量对整体向量的影响

General attention layer



- 增加不同的全连接层增加该层的表现力

General attention layer



Outputs:

context vectors: y (shape: D_v)

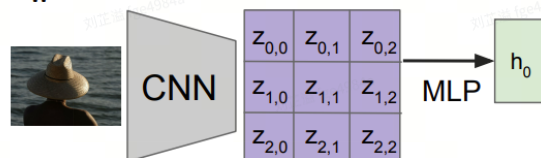
Operations:

Key vectors: $k = xW_k$
 Value vectors: $v = xW_v$
 Alignment: $e_{i,j} = q_i \cdot k_j / \sqrt{D}$
 Attention: $a = \text{softmax}(e)$
 Output: $y_j = \sum_i a_{i,j} v_i$

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

Encoder: $h_0 = f_w(z)$

where z is spatial CNN features
 $f_w(\cdot)$ is an MLP



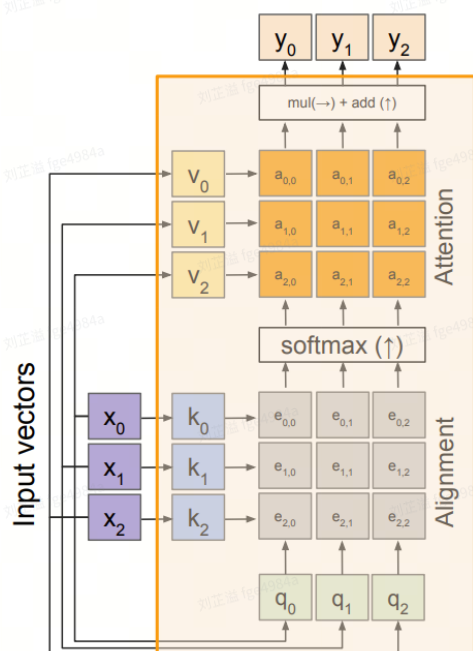
Inputs:

Input vectors: x (shape: $N \times D$)
 Queries: q (shape: $M \times D_q$)

Self Attention Layer

- 不关心输入的排列顺序，需要增加positional encoding

Self attention layer - attends over sets of inputs



Outputs:

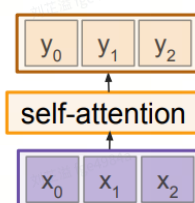
context vectors: y (shape: D_v)

Operations:

Key vectors: $k = xW_k$
 Value vectors: $v = xW_v$
 Query vectors: $q = xW_q$
 Alignment: $e_{i,j} = q_i \cdot k_j / \sqrt{D}$
 Attention: $a = \text{softmax}(e)$
 Output: $y_j = \sum_i a_{i,j} v_i$

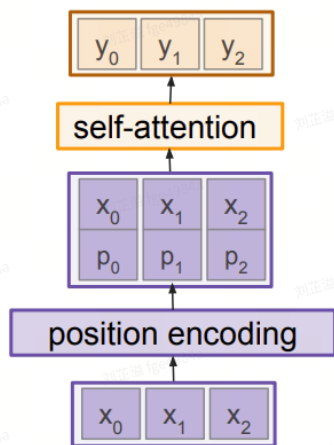
Inputs:

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



Positional Encoding

Positional encoding



Concatenate special positional encoding p_j to each input vector x_j

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

So, $p_i = pos(j)$

Options for $pos(.)$

1. Learn a lookup table:
 - Learn parameters to use for $pos(t)$ for $t \in [0, T)$
 - Lookup table contains $T \times d$ parameters.

Desiderata of $pos(.)$:

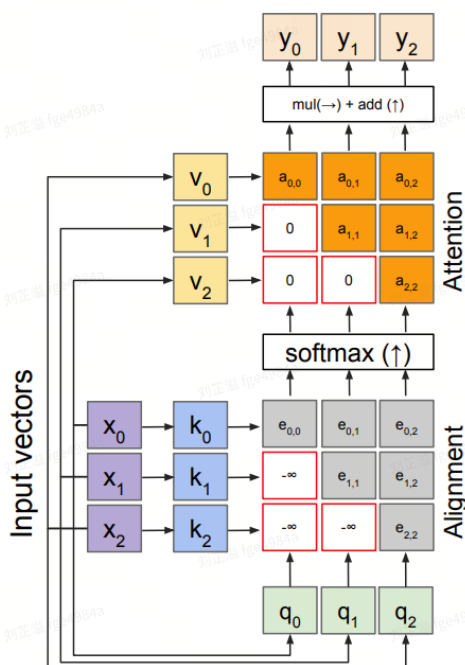
1. It should output a **unique** encoding for each time-step (word's position in a sentence)
2. **Distance** between any two time-steps should be consistent across sentences with different lengths.
3. Our model should generalize to **longer** sentences without any efforts. Its values should be bounded.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

Vaswani et al. "Attention is all you need"

Masked self-attention layer

- 防止向量查看未来向量

Masked self-attention layer



Outputs:
context vectors: y (shape: D_v)

Operations:
Key vectors: $k = xW_k$
Value vectors: $v = xW_v$
Query vectors: $q = xW_q$
Alignment: $e_{ij} = q_j \cdot k_i / \sqrt{D}$
Attention: $a = \text{softmax}(e)$
Output: $y_j = \sum_i a_{ij} v_i$

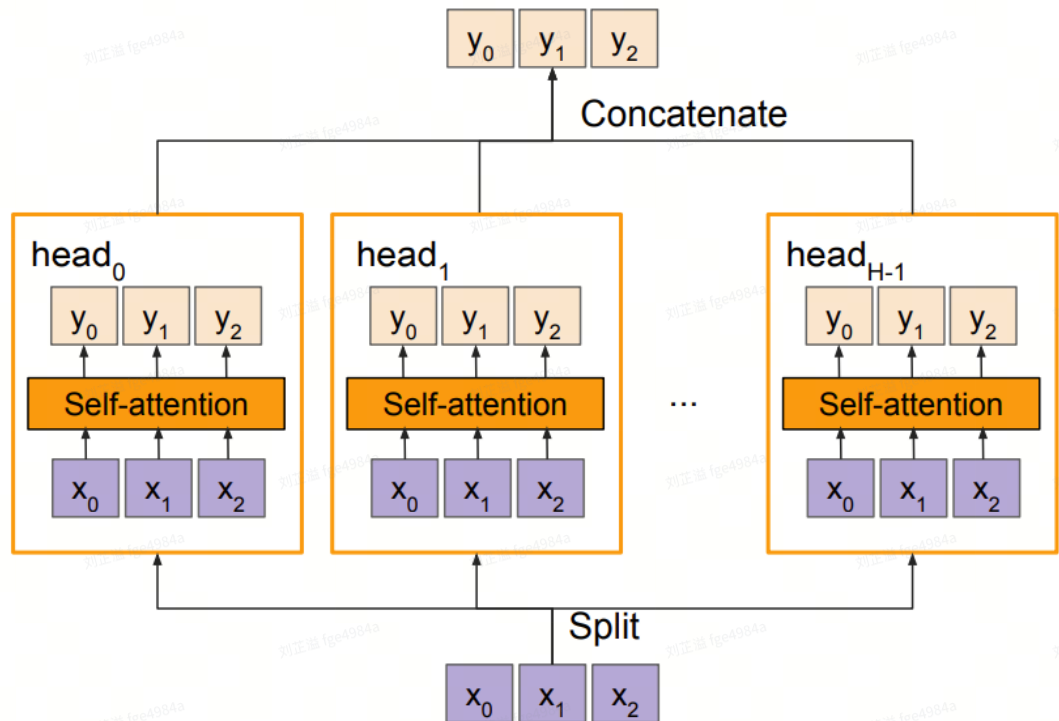
Inputs:
Input vectors: x (shape: $N \times D$)

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

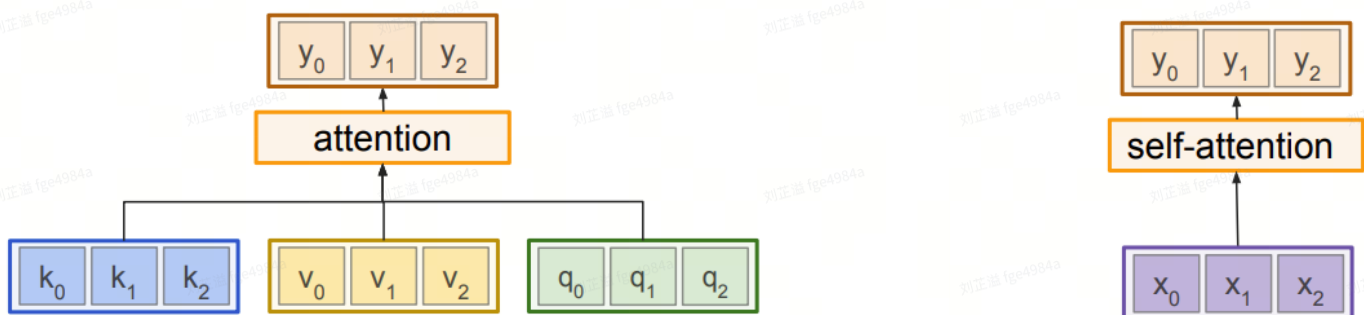
Multi-head self-attention layer

Multi-head self-attention layer

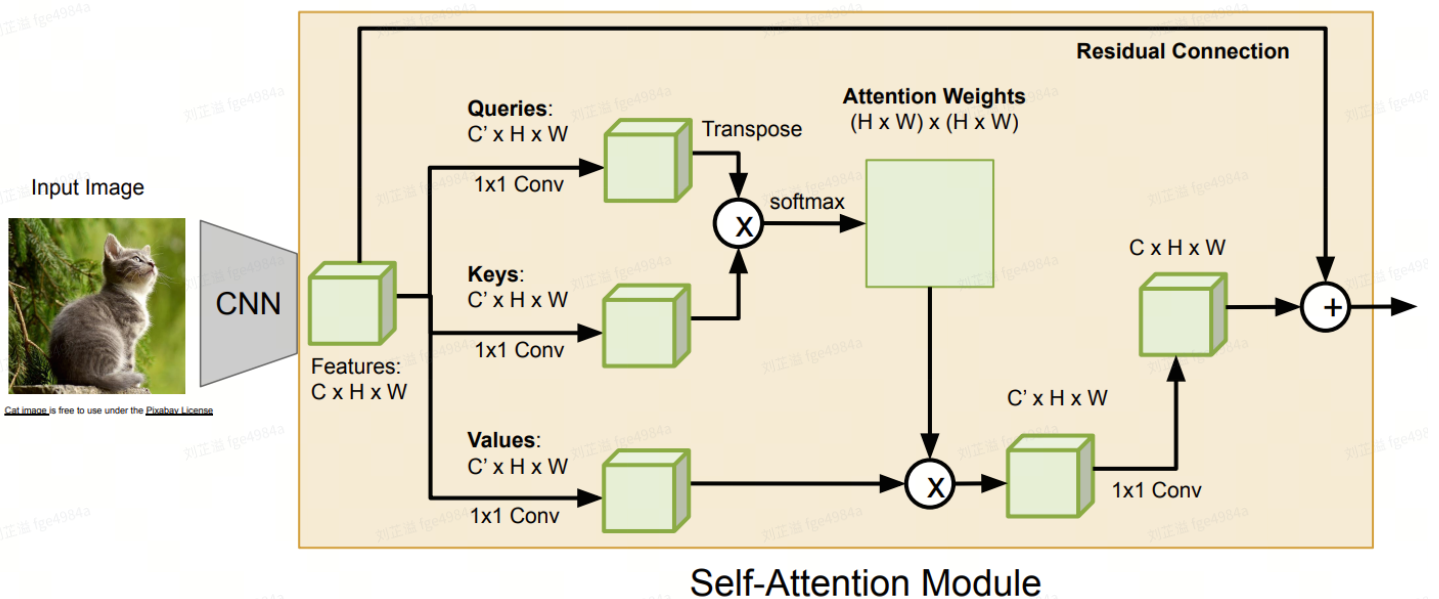
- Multiple self-attention heads in parallel



General attention versus self-attention



CNN with Self-Attention



Self-Attention Module

RNNs VS. Transformer

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 scalars need to be calculated and stored for a single self-attention head. (but GPUs are getting bigger and better)

Transformer

Image Captioning using Transformers

Input: Image I

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

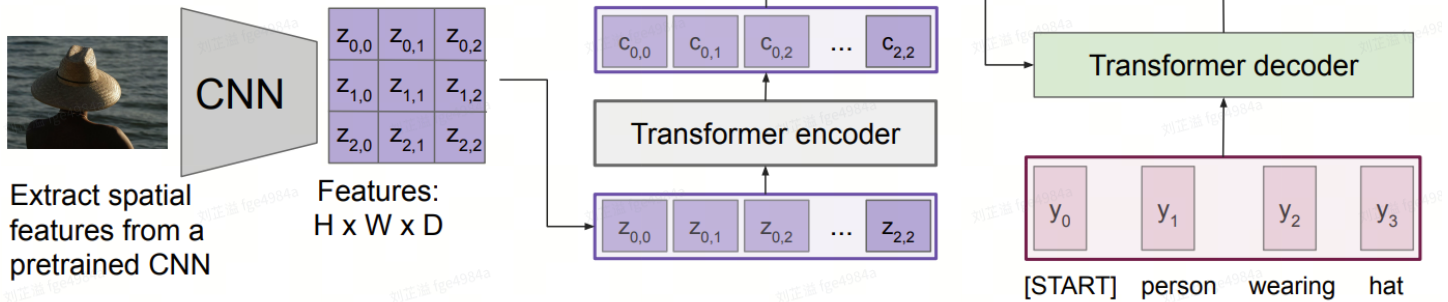
Decoder: $y_t = T_D(y_{0:t-1}, c)$

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

Encoder: $c = T_w(z)$

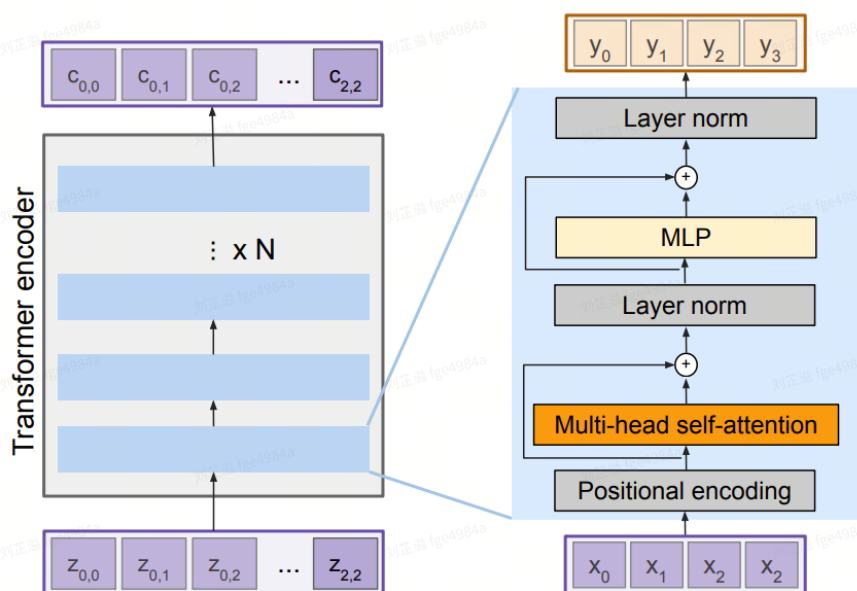
where z is spatial CNN features

$T_w(\cdot)$ is the transformer encoder



Transformer encoder block

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.

Transformer decoder block

The Transformer Decoder block

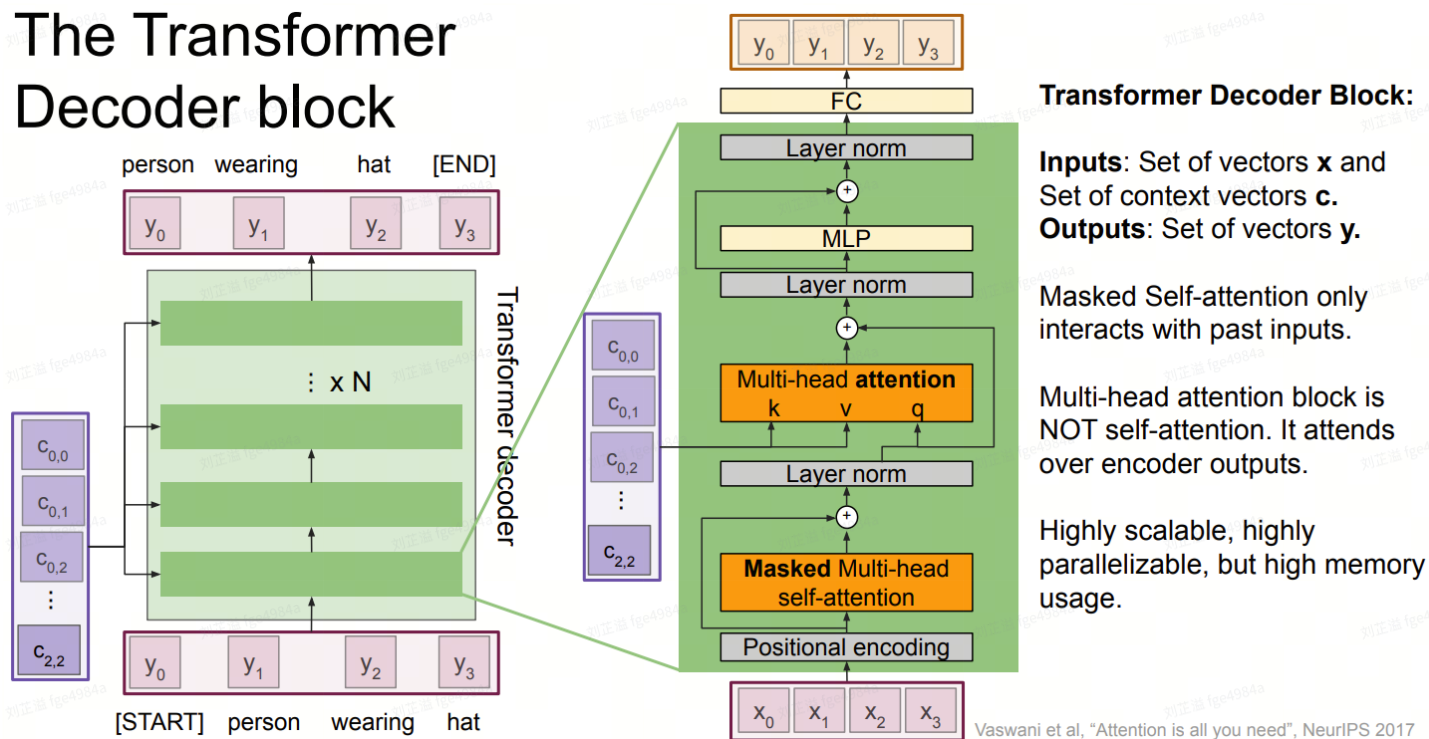


Image Captioning using ONLY Transformers

Image Captioning using **ONLY** transformers

- Transformers from pixels to language

