

Attention and Transformers-I

Two Example Tasks

1. Machine Translation:

- Input: “He loved to eat”
- Output: “Er liebte zu essen” (German)

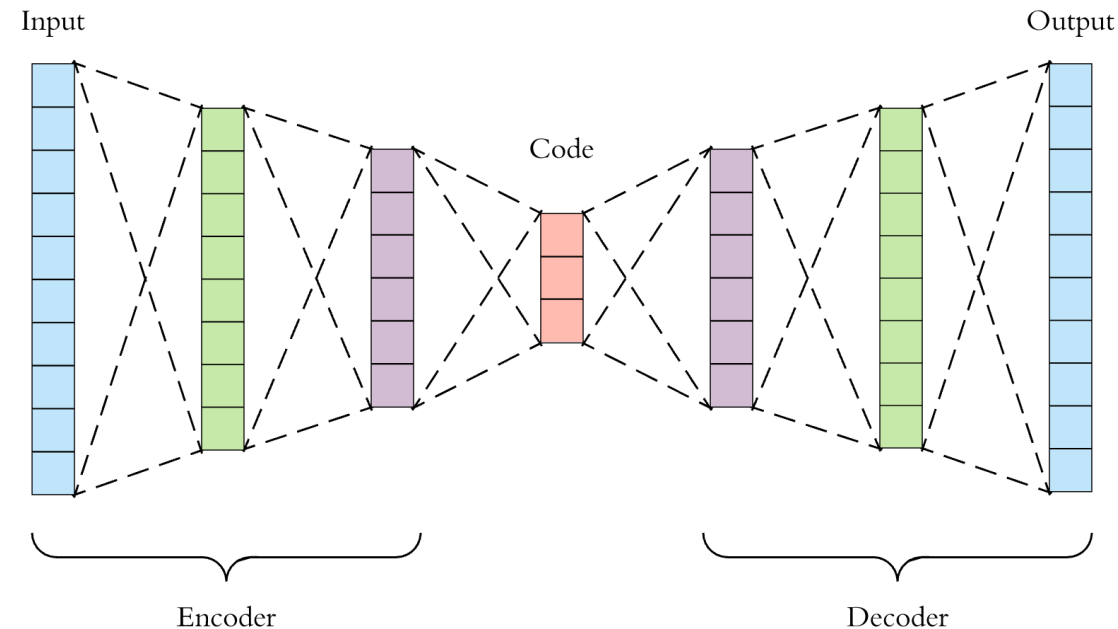
2. Image Captioning

- Input: Image
 - Output: “Cat Sat Outside”
- Two Possible Steps
 - Encode the input information
 - Decode into the target modality/language

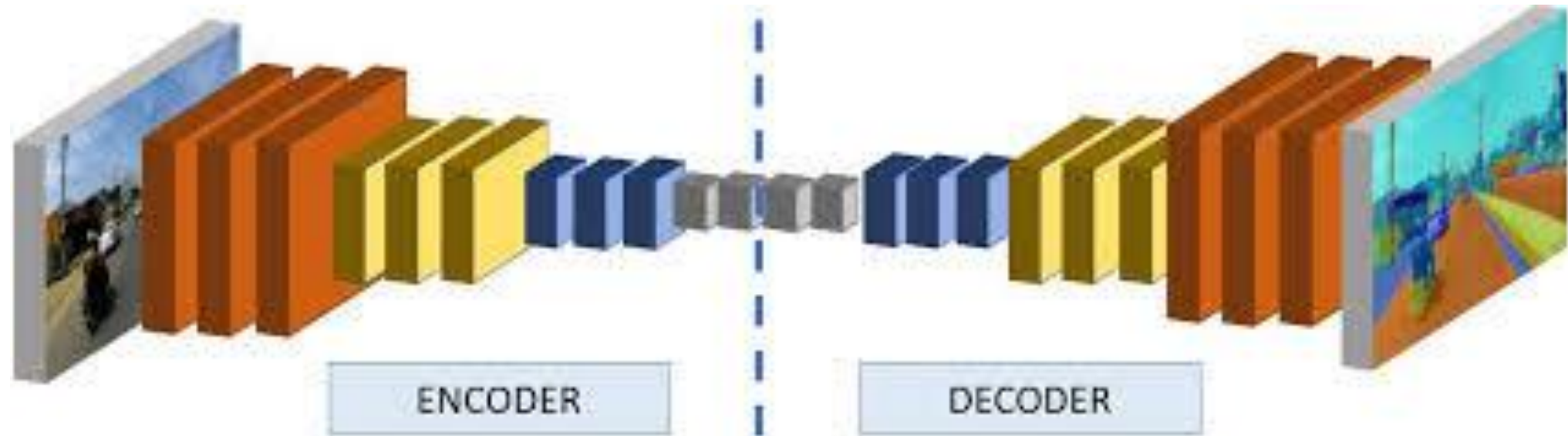


Encoder-Decoder: Auto encoders

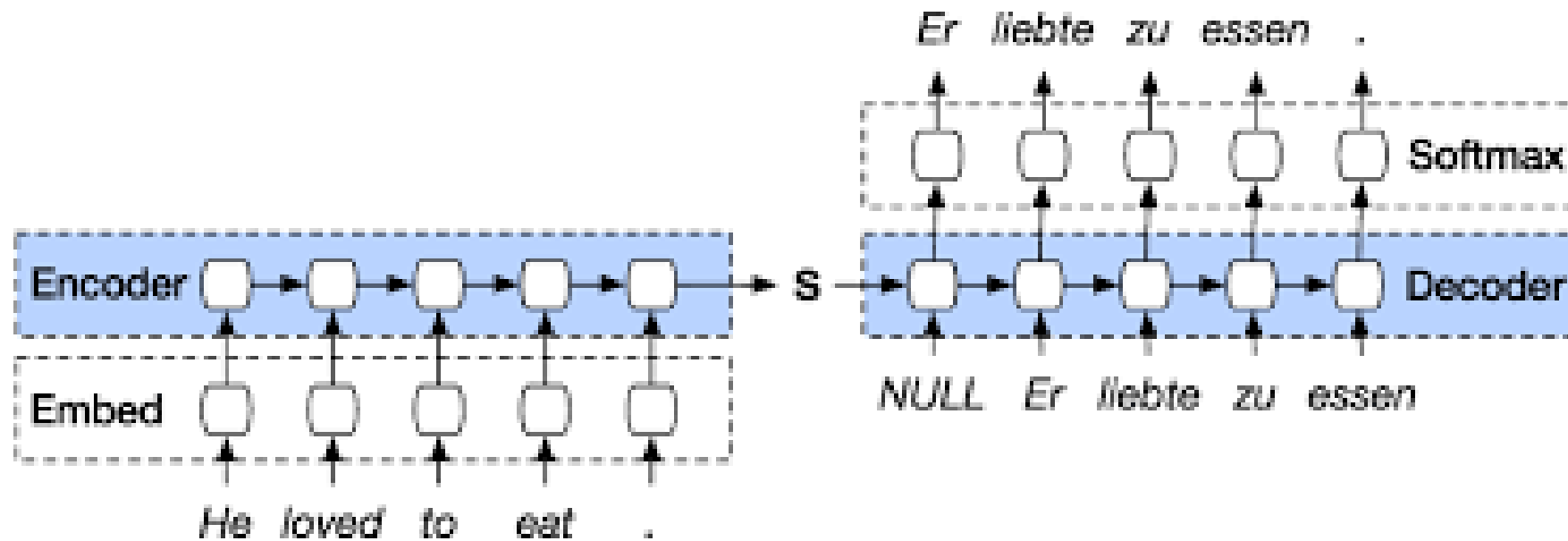
- Encoder Decoder architecture
 - Use for compression
 - Or throw away the decoder after training use the compressed code for supervised learning tasks (particularly useful in problems with small training set)



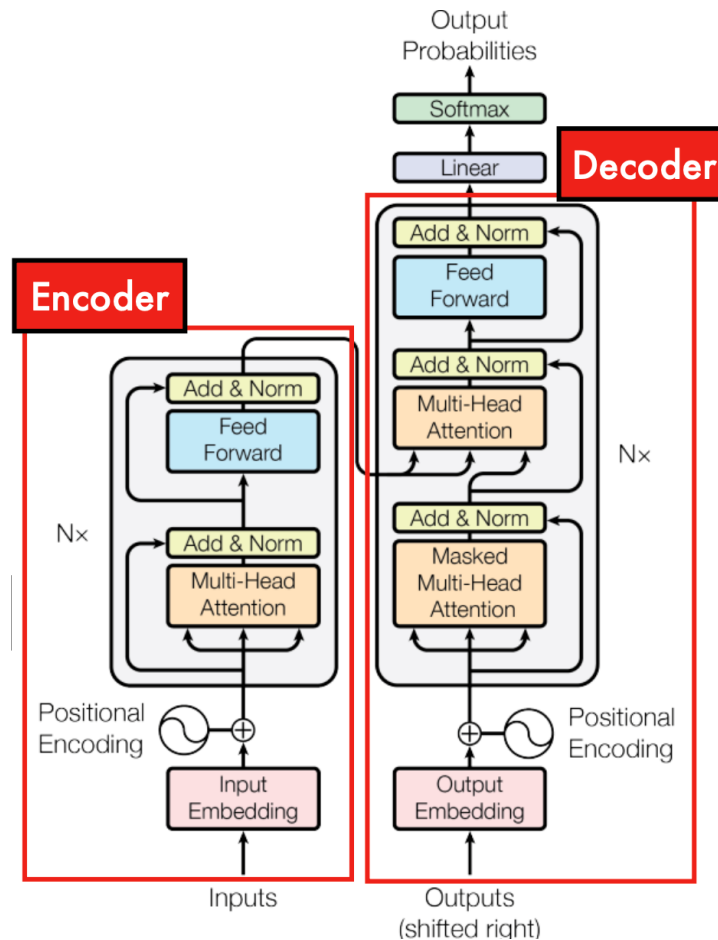
Encoder-Decoder



Encoder-Decoder Models: Seq2Seq



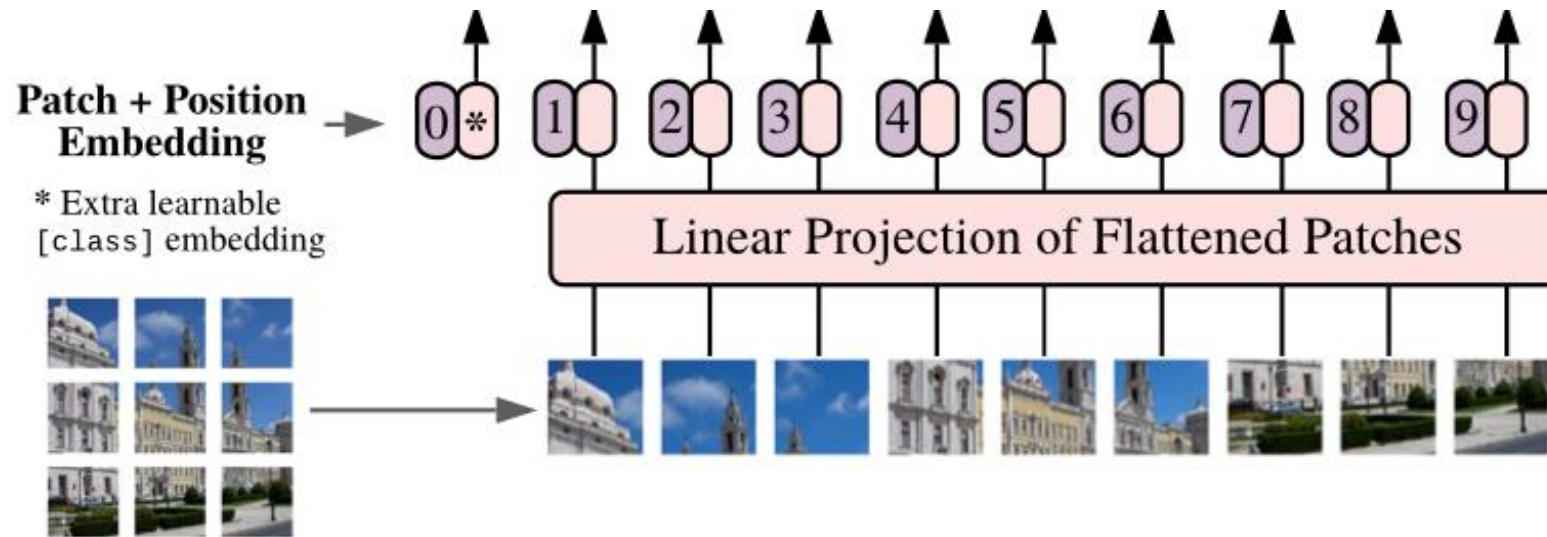
Transformer: Encoder and Decoder



- **Encoder:**
 - Learn useful representation from the input
 - Eg. BERT (encoder only)
- **Decoder:**
 - Decodes the learned representation along with any other inputs and predict the output
 - Eg. GPT-3 (decoder only)
- **Encoder-Decoder:**
 - Eg. Sequence to Sequence Tasks, like Machine Translation

Modelling Inter Dependency of Elements/Patches

1. **The animal** did not cross the road because **it** was too tired.
2. The animal did not cross **the road** because **it** was too wide.



Human Knowledge

Use of Data

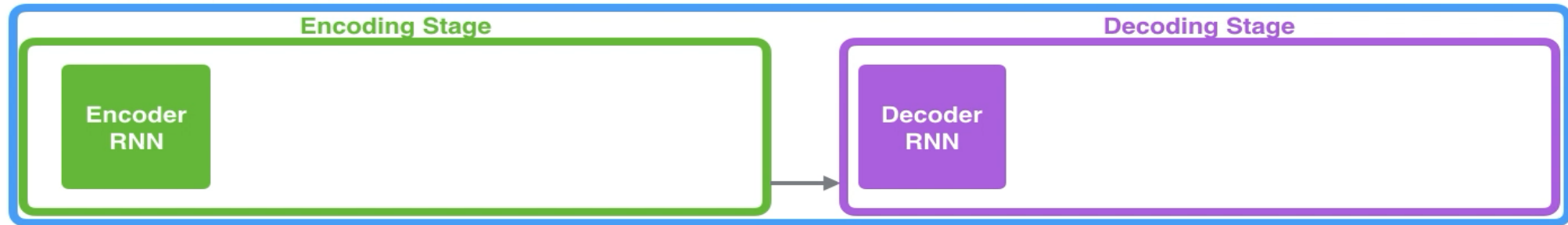
Transformers

CNNs

Hand Crafted

Seq2Seq

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL



Je

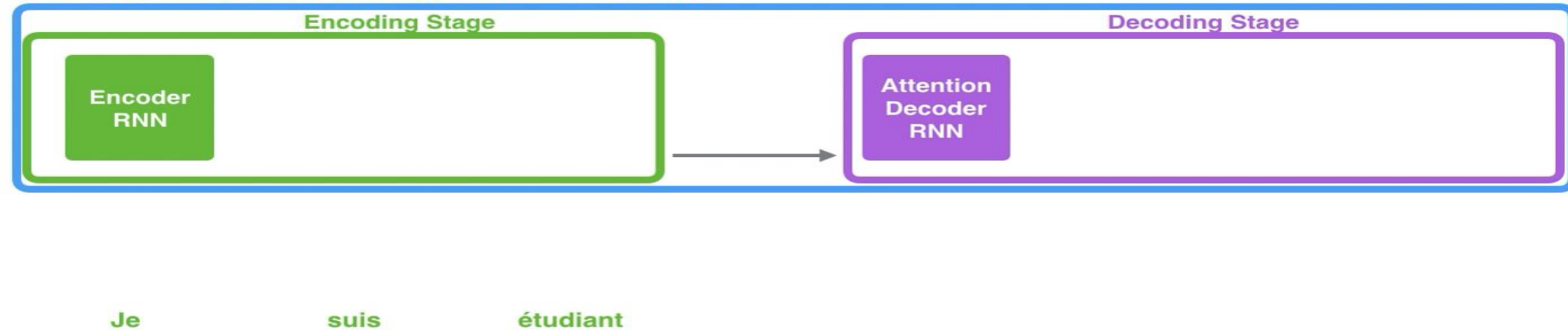
suis

étudiant

Seq2Seq with attention

Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



The encoder passes a lot more data/information/hidden states to the decoder.

Strategy : Seq2Seq with Context

Input: Sequence x_1, \dots, x_T

Output: Sequence y_1, \dots, y_T

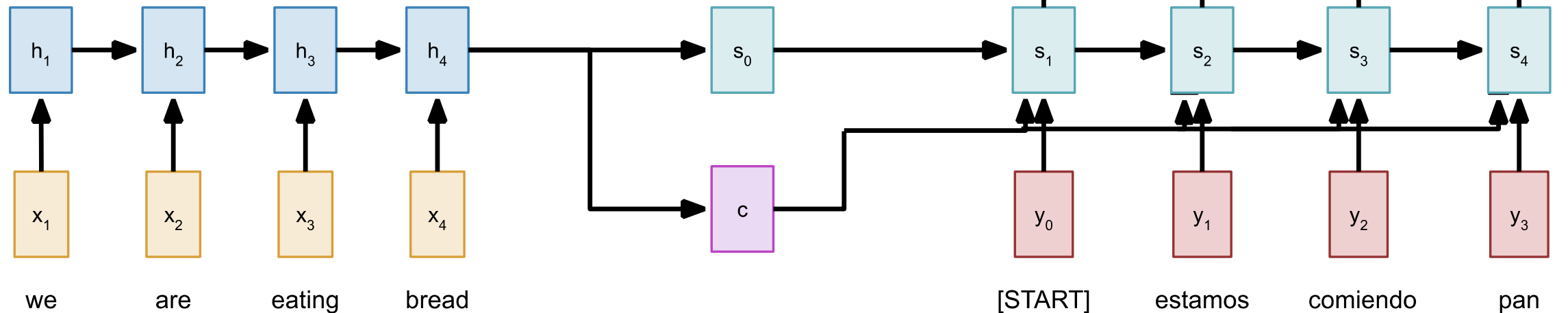
Decoder: $s_t = g_U(y_{t-1}, s_{t-1}, c)$

Encoder: $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:

Initial decoder state s_0

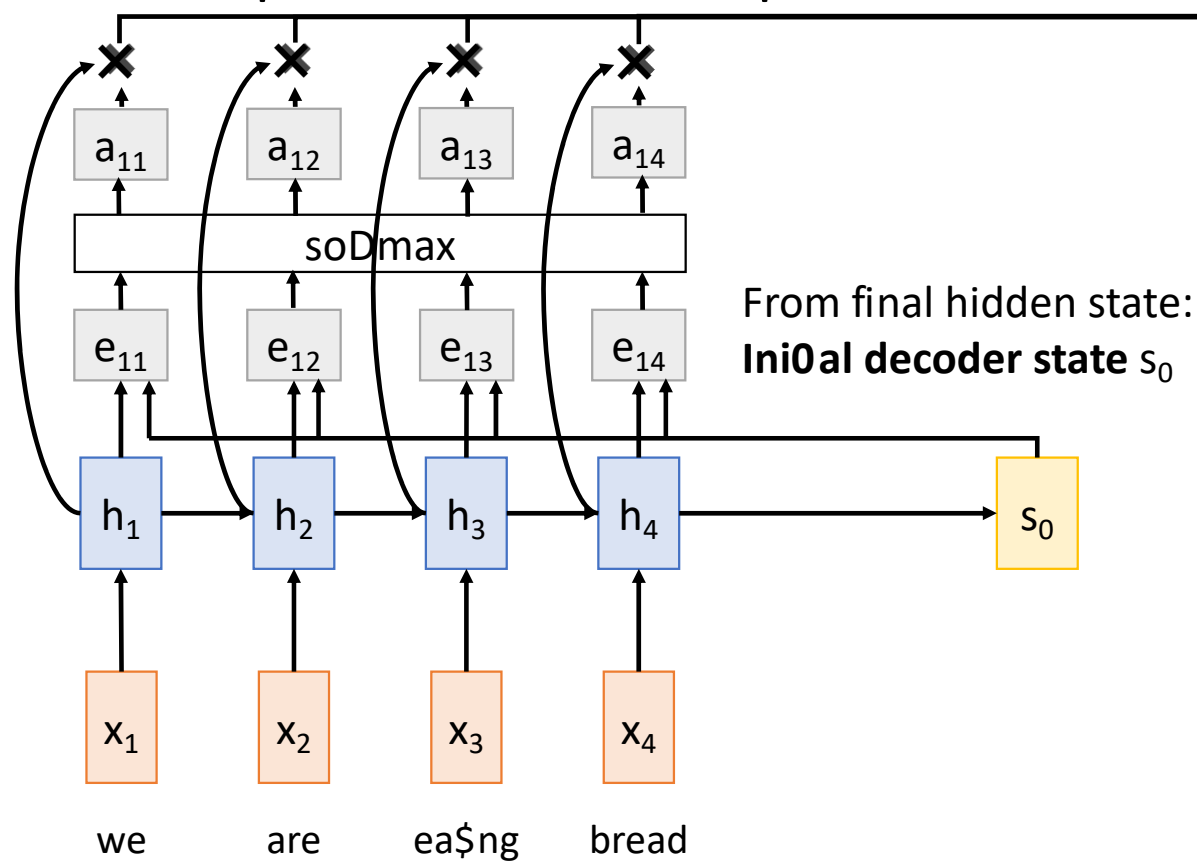
Context vector c (often $c=h_T$)



Problem

- Input sequence is summarized (bottlenecked) through a fixed size fixed vector.
 - Not good for long sequences/inputs
- **Solution:**
 - Provide more detailed context to the decoder.
 - Use new context vector at each time step.

Improved Strategy: Computation with “Attention”



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

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

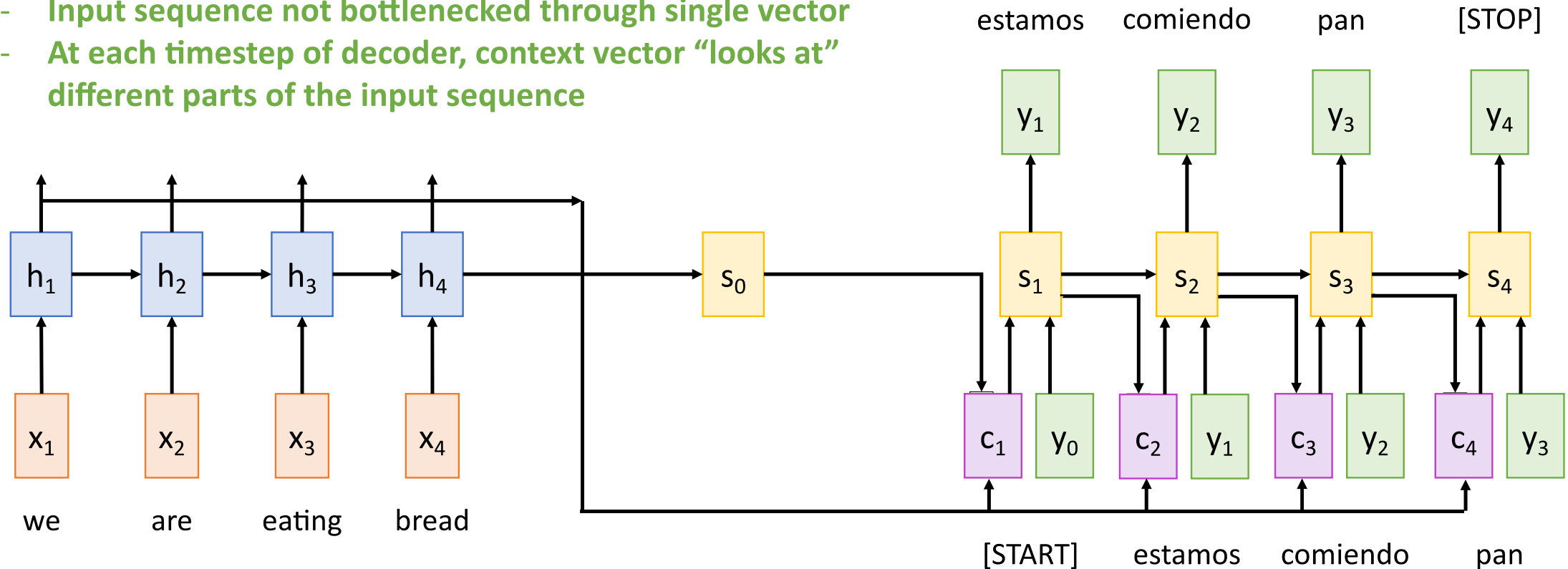
Use context vector in
 decoder: $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

**This is all differentiable! Do not
 supervise attention weights –
 backprop through everything**

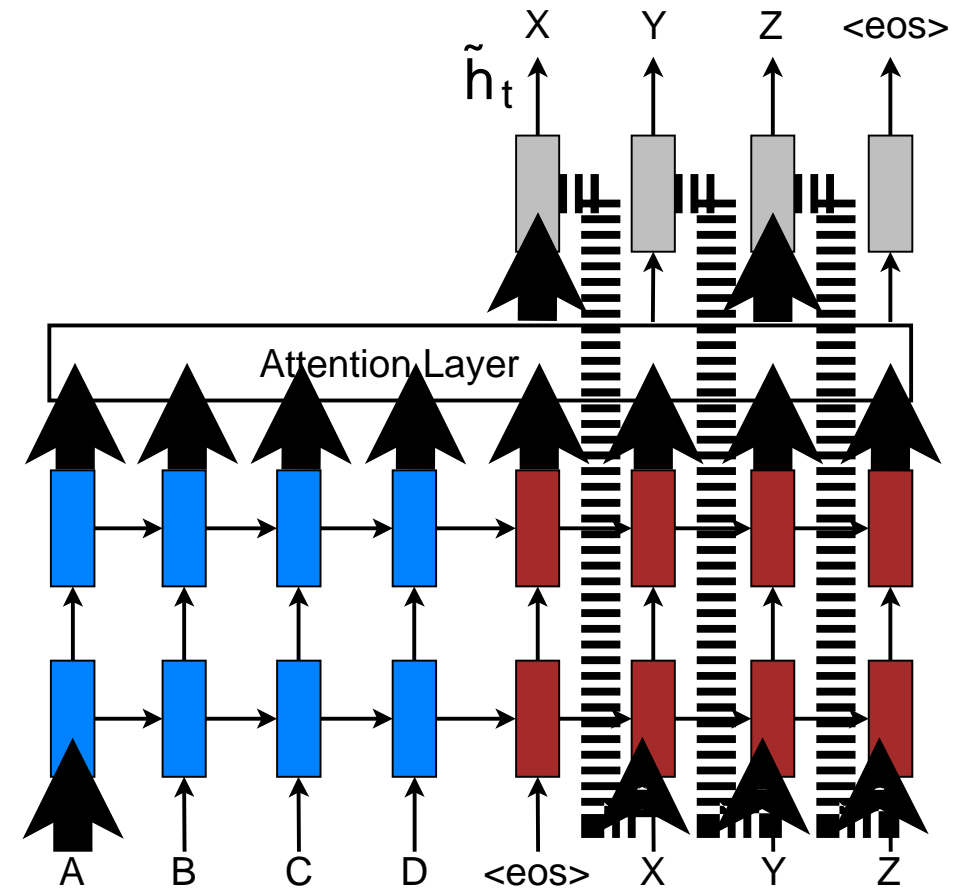
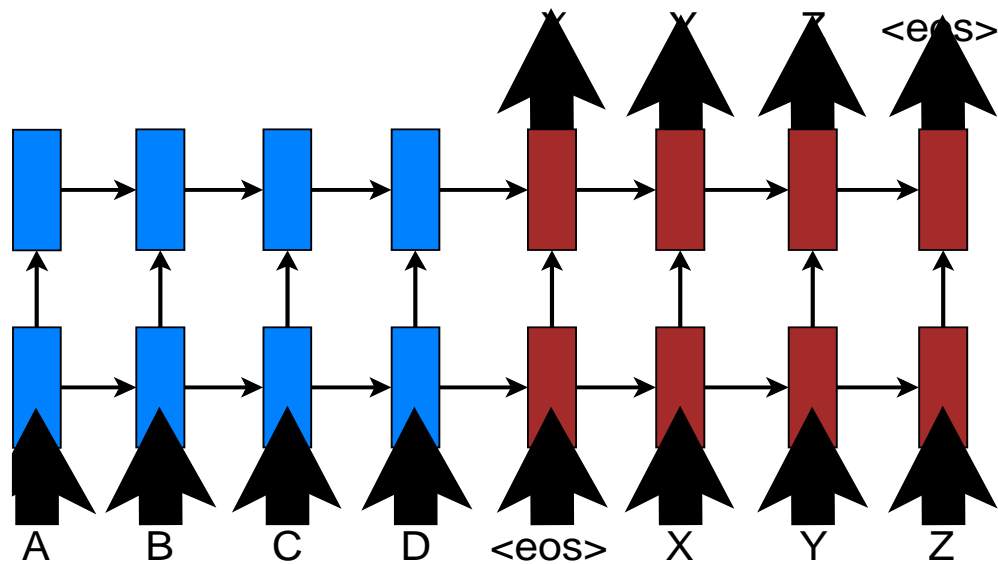
Attention in Translation (Seq2Seq)

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



Attention Layer



Visualization

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 $a_{t,i}$

Diagonal attention means words correspond in order

Attention figures out different word orders

Diagonal attention means words correspond in order

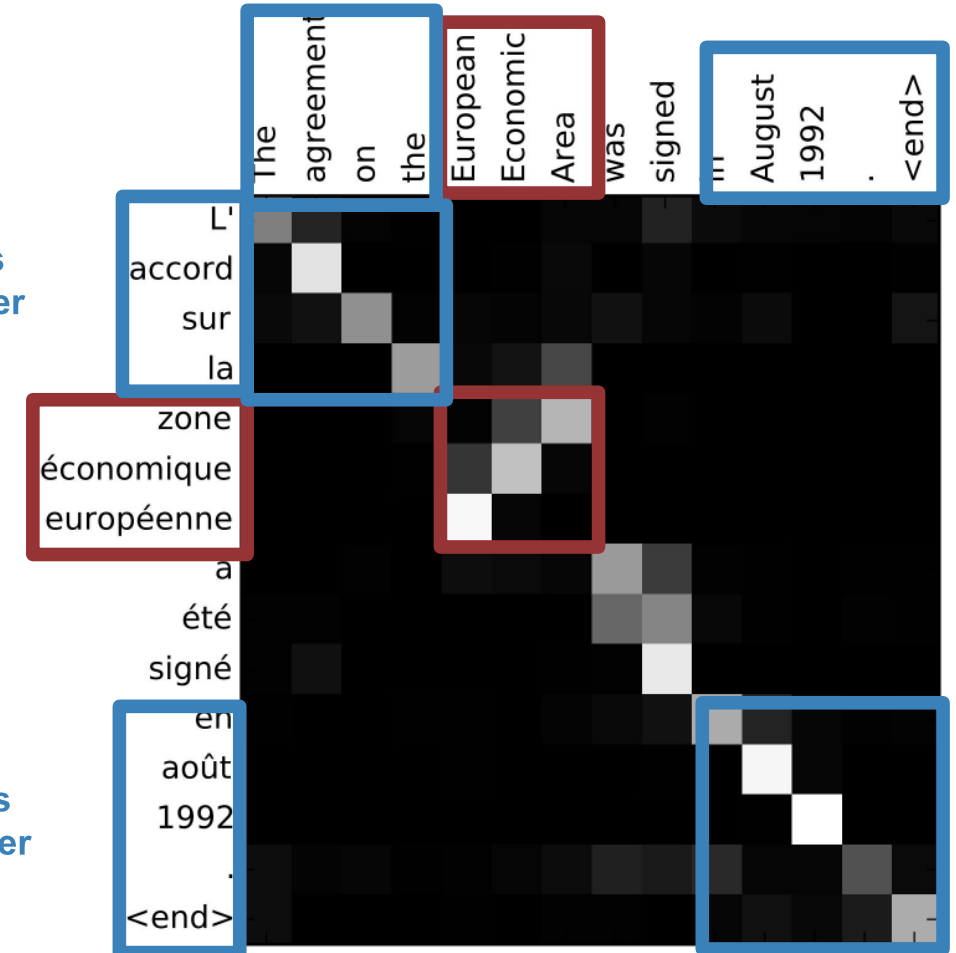


Image captioning with RNN and Attention

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

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

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



CNN

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

Alignment scores

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

Attention weights

$a_{2,1,1}$	$a_{2,1,2}$	$a_{2,1,3}$
$a_{2,2,1}$	$a_{2,2,2}$	$a_{2,2,3}$
$a_{2,3,1}$	$a_{2,3,2}$	$a_{2,3,3}$

softmax

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

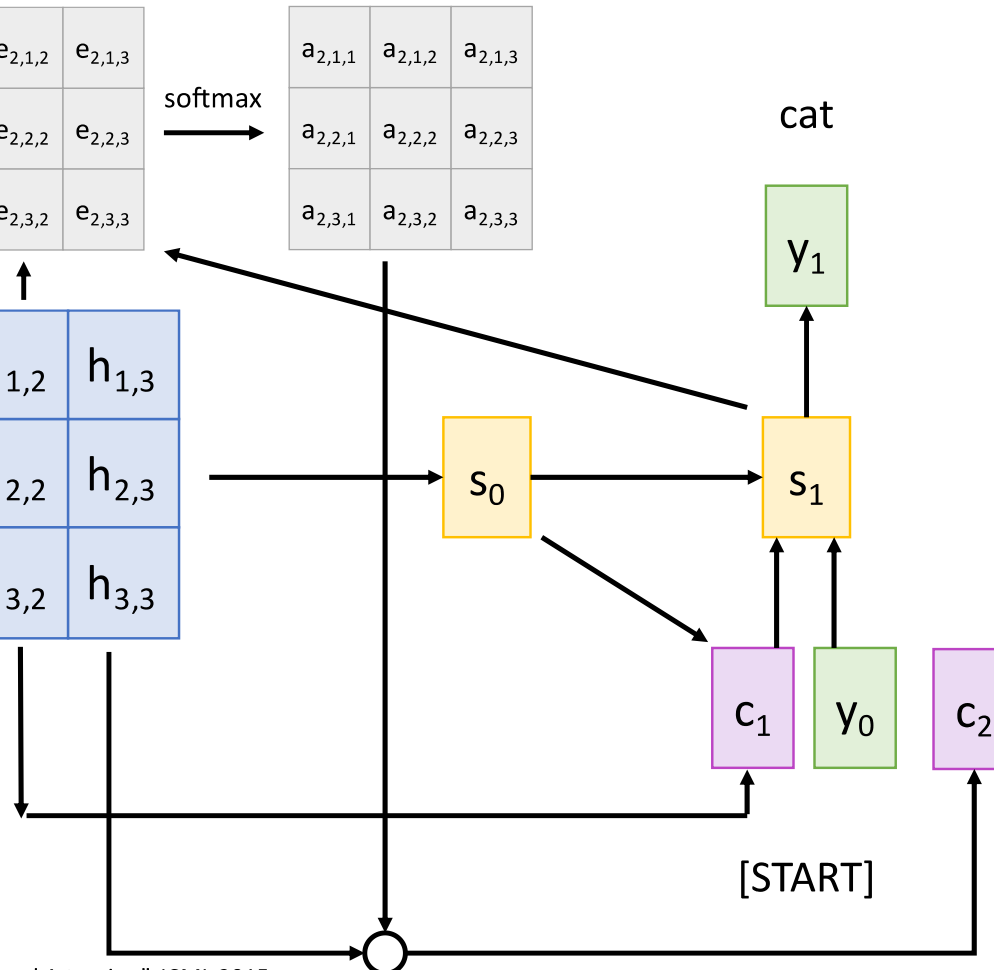


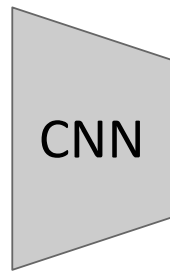
Image captioning with RNN and Attention

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

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

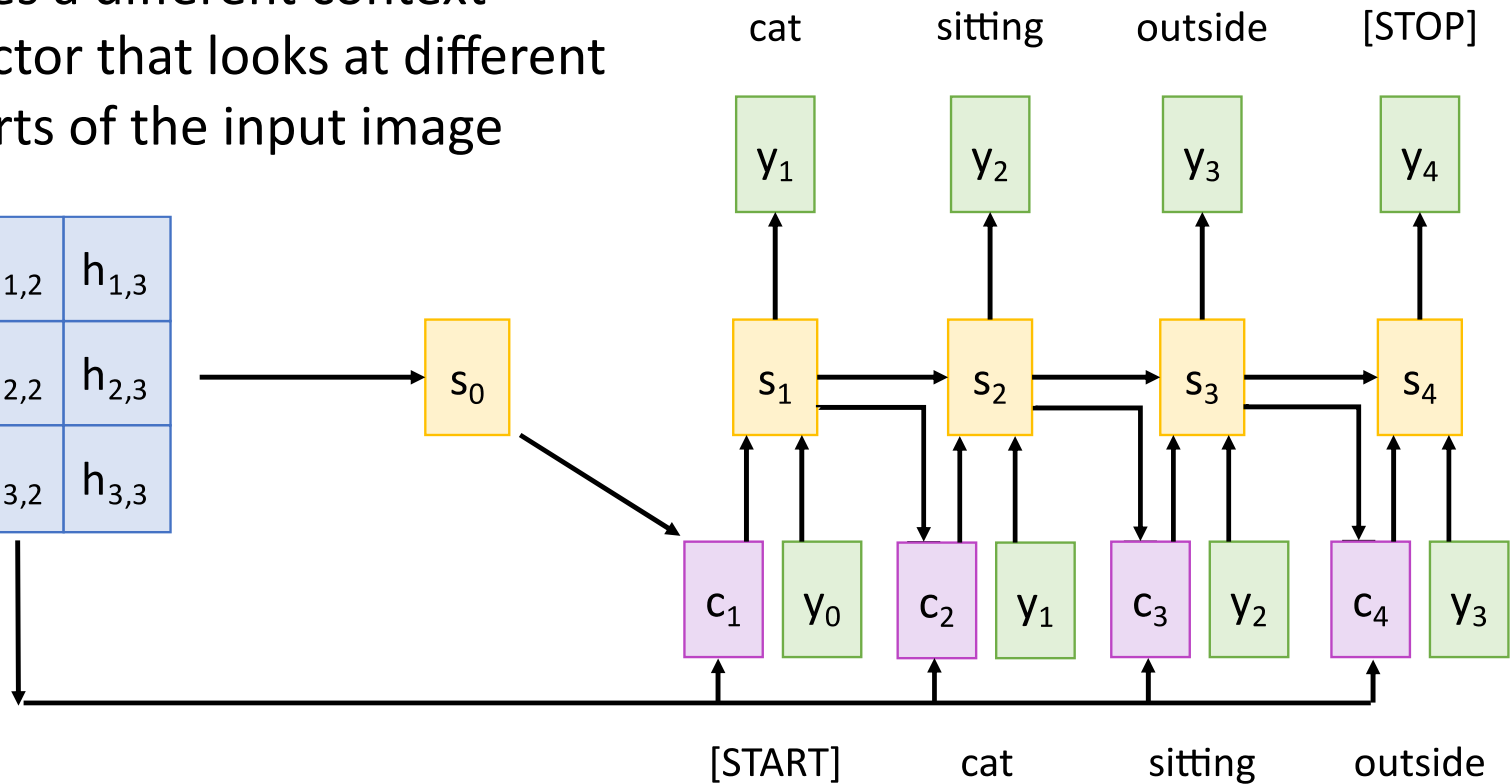
$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$

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



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

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



Blank Slide

Blank Slide

Query, Key and Value

Query, Key and Value

We project each embedding:

Queries

Keys

Values

Queries: "Here's what I'm looking for" $W^Q \in \mathbb{R}^{D \times d_k}$

Keys: "Here's what I have" $W^K \in \mathbb{R}^{D \times d_k}$

Values: "What gets communicated" $W^V \in \mathbb{R}^{D \times d_v}$

d_k is dimension of **queries** & **keys**, d_v is dimension of **values**

$$Q = XW^Q \in \mathbb{R}^{N \times d_k}$$

$$K = XW^K \in \mathbb{R}^{N \times d_k}$$

$$V = XW^V \in \mathbb{R}^{N \times d_v}$$

Query, Key and Value:

- Mimics the retrieval of a value v for a query q based on k in the database

q	K1	V1
	K2	V2
	K3	V3
	K4	V4
	K5	V5

A_i = Similarity (query, key- i)

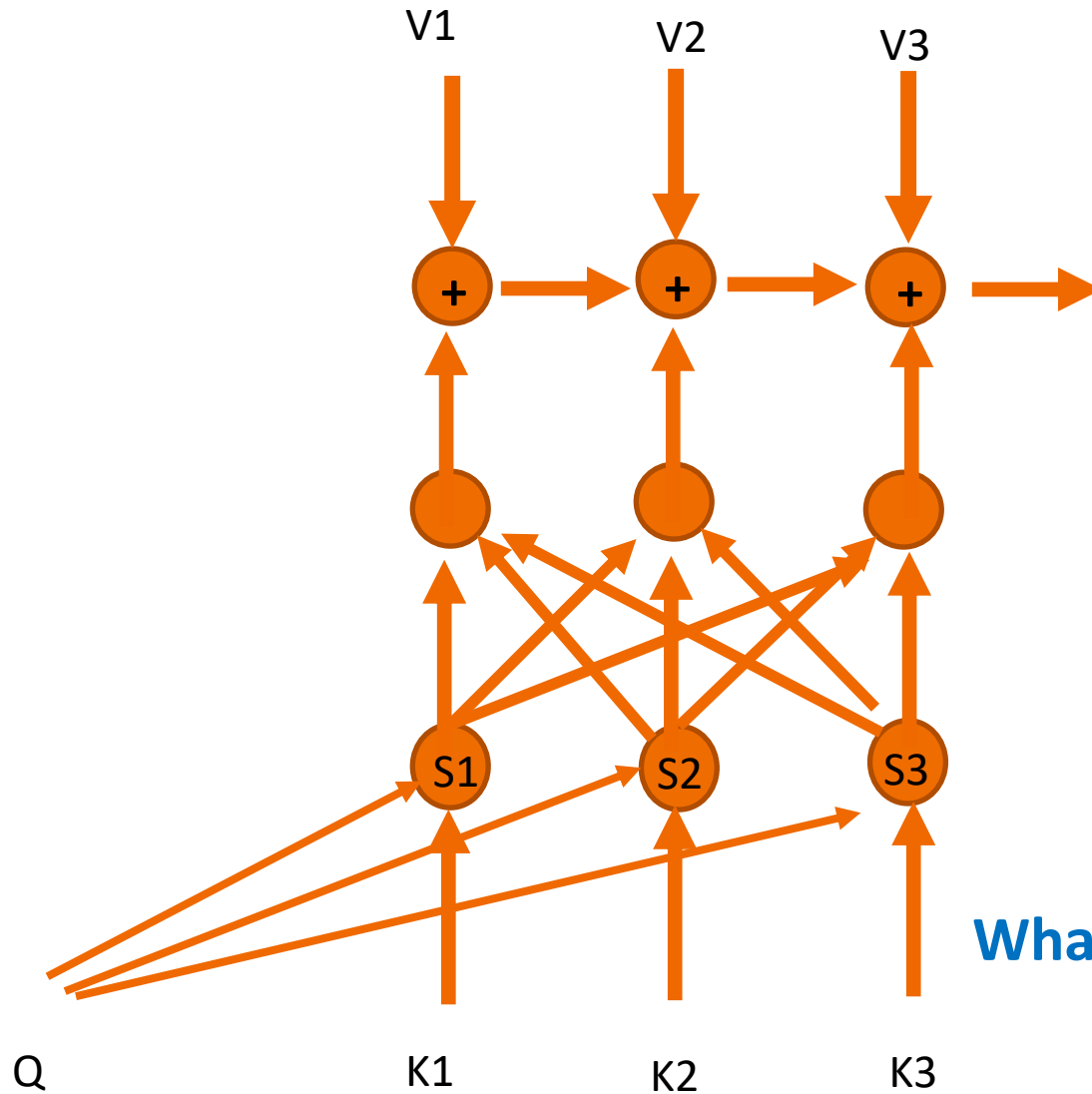
V = weighted sum of A_i and V_i

$$att(q, k, v) = \sum_i s(q, k_i) v_i$$

We compare "Query" with "Key" and "Value" is returned.

What is "Query", "Key" and "Value"?

Attention



What is “Query”, “Key” and “Value”?

Query, Key and Value

We project each embedding:

Queries

Keys

Values

Queries: "Here's what I'm looking for" $W^Q \in \mathbb{R}^{D \times d_k}$

Keys: "Here's what I have" $W^K \in \mathbb{R}^{D \times d_k}$

Values: "What gets communicated" $W^V \in \mathbb{R}^{D \times d_v}$

d_k is dimension of **queries** & **keys**, d_v is dimension of **values**

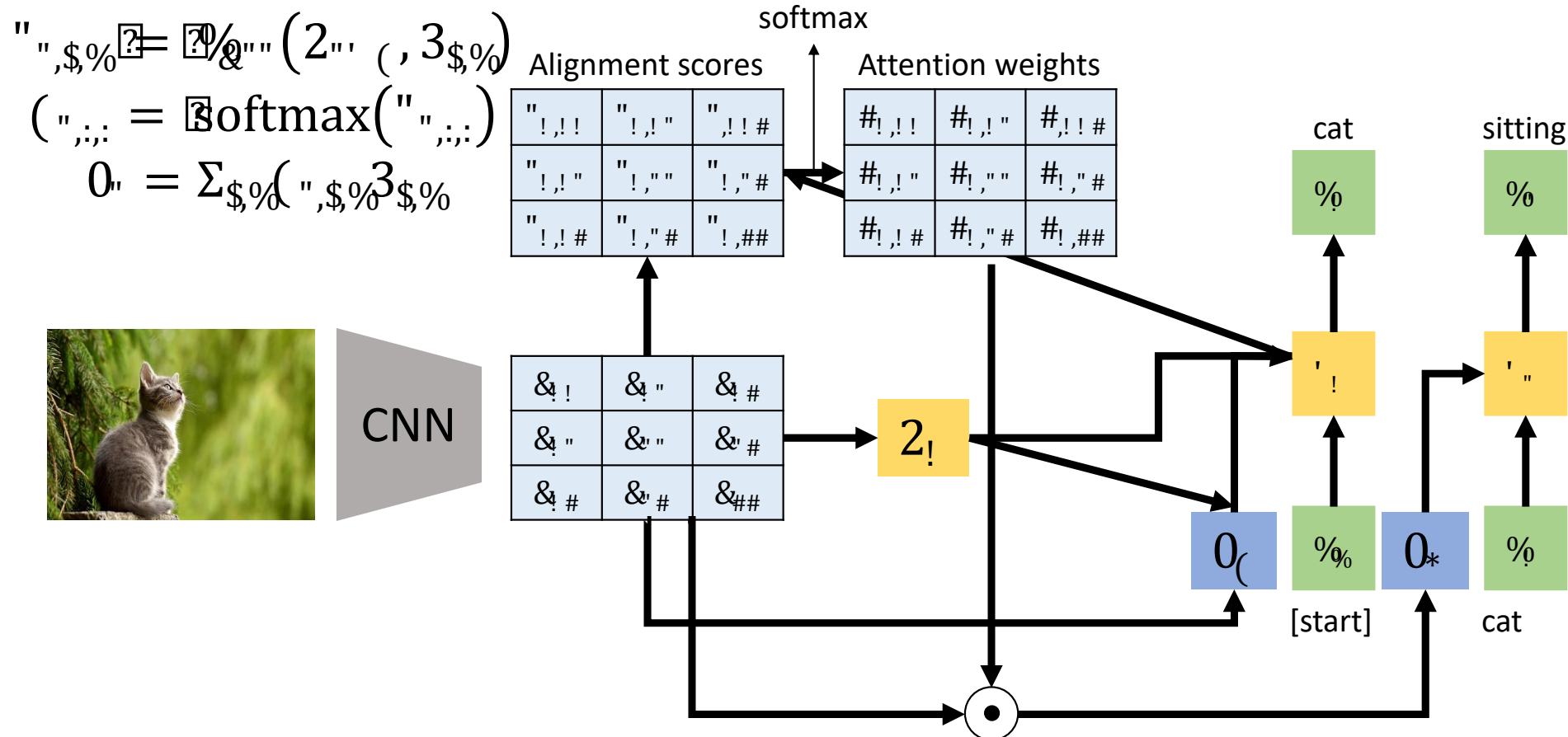
$$Q = XW^Q \in \mathbb{R}^{N \times d_k}$$

$$K = XW^K \in \mathbb{R}^{N \times d_k}$$

$$V = XW^V \in \mathbb{R}^{N \times d_v}$$

Towards Scalable and More Generic

Notation change: ! to ", and # to \$



Towards Scalable and More Generic

Inputs:

Query vector: Q (Shape: $N_Q \times D_Q$)

Input vectors: X (Shape: $N_X \times D_Q$)

Key matrix: W_K (Shape: $D_X \times D_Q$)

Value matrix: W_V (Shape: $D_X \times D_V$)

Computation:

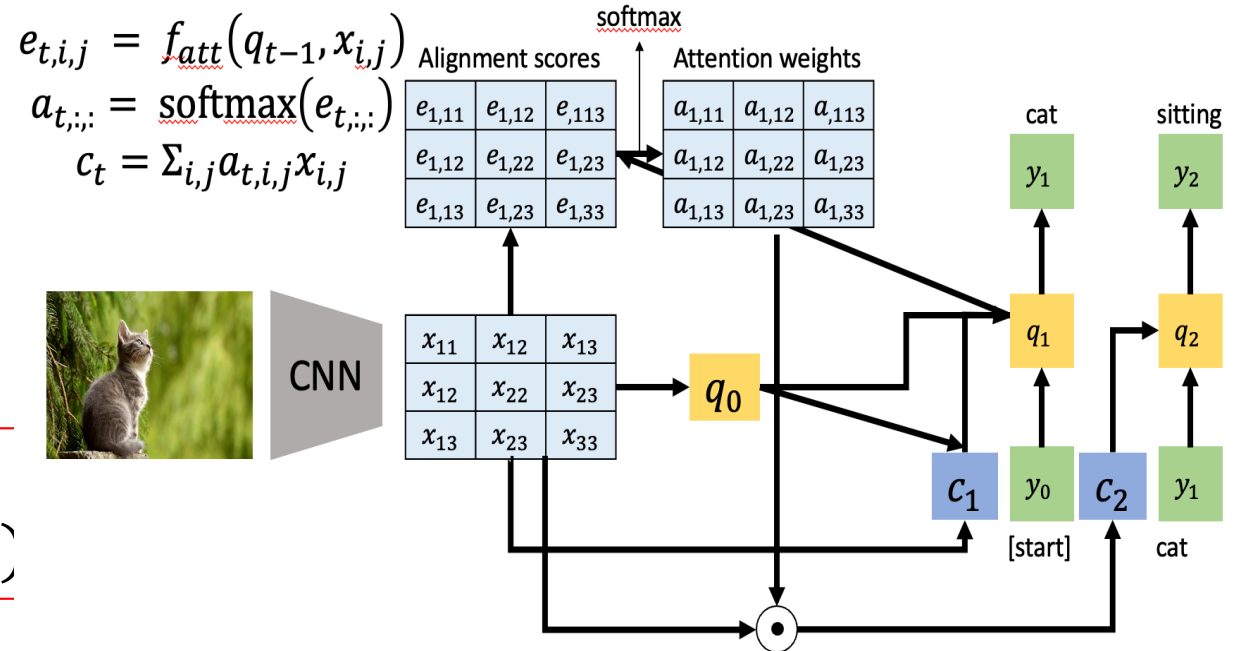
Key vectors: $K = XW_K$ (Shape: $N_X \times D_K$)

Value Vectors: $V = XW_V$ (Shape: $N_X \times D_V$)

Similarities: $E = \frac{QK^T}{\sqrt{D_K}}$ (Shape: $N_Q \times N_X$), $E_{0,0} = (Q_{0,0} \cdot K_{0,0}) / \sqrt{D_K}$

Attention weights: $A = \text{softmax}(E, \text{dim} = 1)$ (Shape: $N_Q \times N_X$)

Output vectors: $Y = AV$ (Shape: $N_Q \times D_V$) $Y_{0,0} = \sum_i A_{0,i} V_i$



Changes:

- Use **scaled** dot product for similarity
- Multiple query vectors
- Separate key and value

Blank Slide

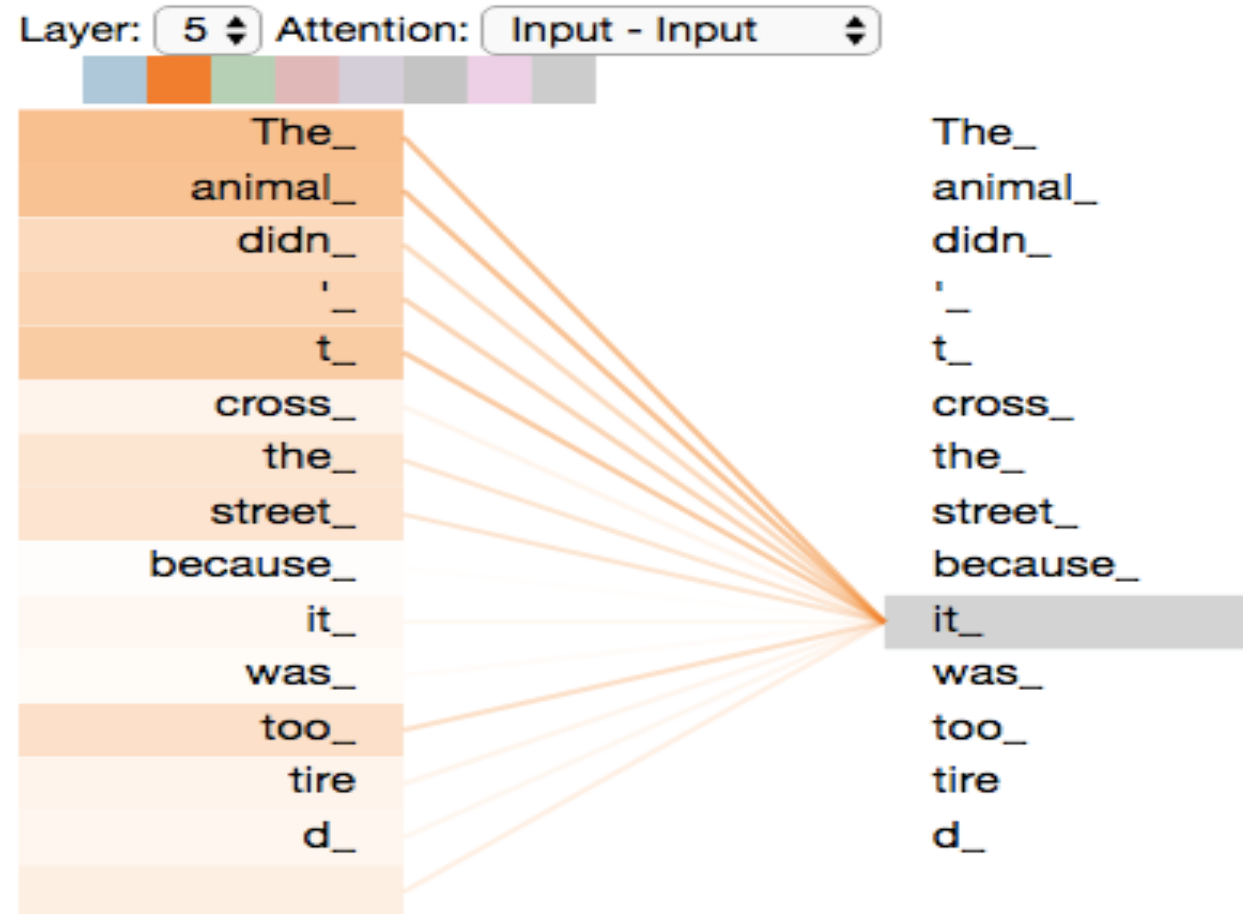
Blank Slide

Self Attention

Self Attention

- Self attention learns the relationship between elements in a sequence.
 - say between words in a sentence

Self Attention



1. **The animal** did not cross the road because **it** was too tired.
2. The animal did not cross **the road** because **it** was too wide.

Blank Slide

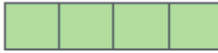
Self Attention: Step 1

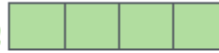
Input

Thinking

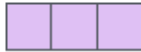
Machines

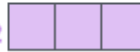
Embedding

X_1 

X_2 

Queries

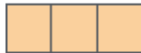
q_1 

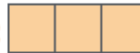
q_2 



W^Q

Keys

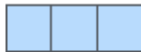
k_1 

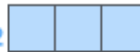
k_2 



W^K

Values

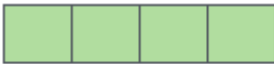
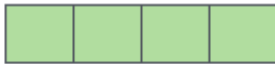


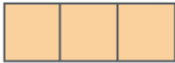


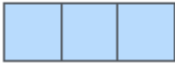
v_1 

v_2 



W^V

Self Attention: Step 2

Input	Thinking	Machines
Embedding	x_1 	x_2 
Queries	q_1 	q_2 
Keys	k_1 	k_2 
Values	v_1 	v_2 
Score	$q_1 \cdot k_1 = 112$	$q_1 \cdot k_2 = 96$
Divide by 8 ($\sqrt{d_k}$)	14	12
Softmax	0.88	0.12

Self Attention: Step 3

Input

Embedding

Queries

Keys

Values

Score

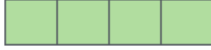
Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax
X
Value

Sum

Thinking

x_1 

q_1 

k_1 

v_1 

$q_1 \cdot k_1 = 112$

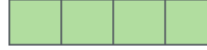
14

0.88

v_1 

z_1 

Machines

x_2 

q_2 

k_2 

v_2 

$q_1 \cdot k_2 = 96$

12

0.12

v_2 

z_2 

Self Attention: Matrix View



Self Attention

One query per input vector

Inputs:

Input vectors: X (Shape: $N_X \times D_Q$)

Key matrix: W_K (Shape: $D_X \times D_Q$)

Value matrix: W_V (Shape: $D_X \times D_V$)

Query matrix: W_Q (Shape: $D_Q \times D_Q$)

Computation:

Query Vectors $Q = XW_Q$

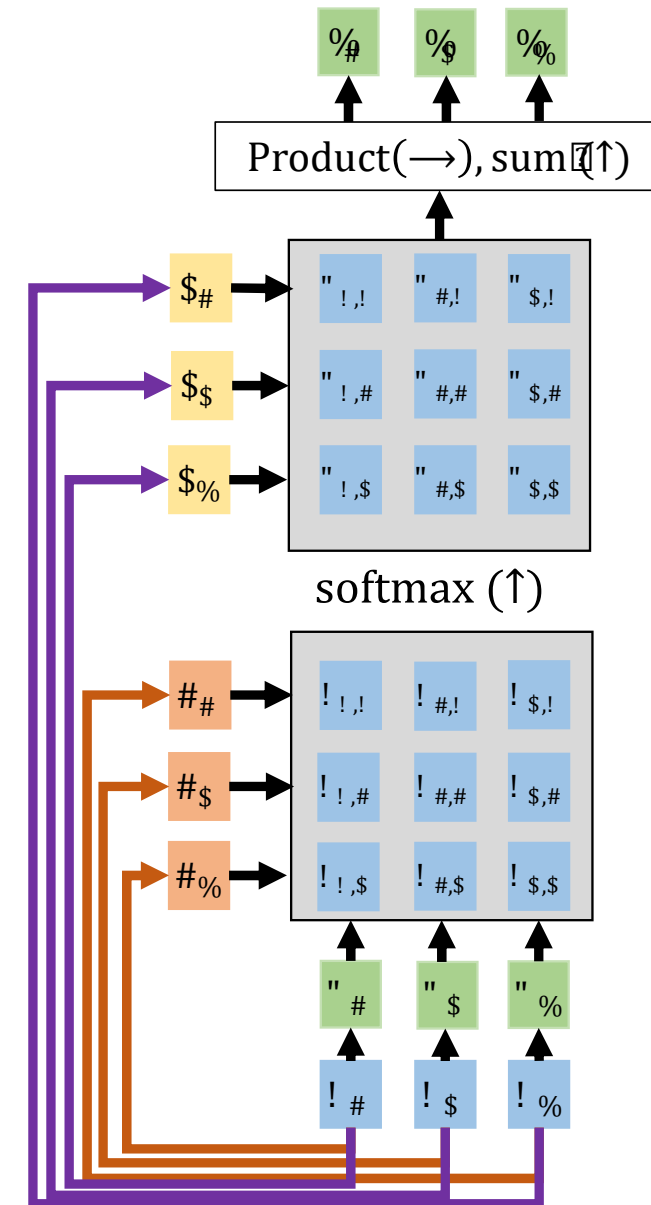
Key vectors: $K = XW_K$ (Shape: $N_n \times D_k$)

Value Vectors: $V = XW_V$ (Shape: $N_n \times D_v$)

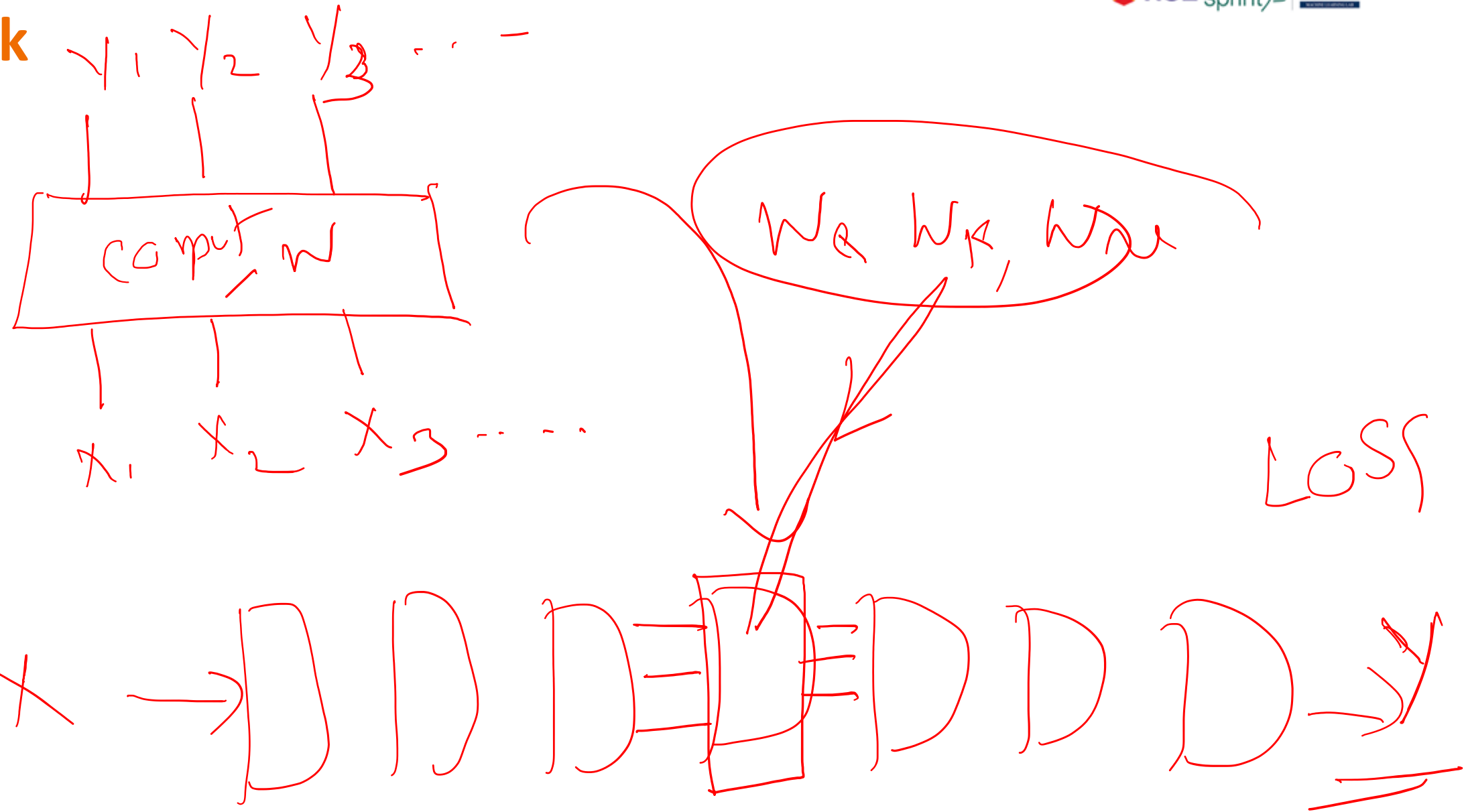
Similarities: $E = \frac{Q \cdot K^T}{\sqrt{D_k}}$ (Shape: $N_n \times N_n$) $E_{\%} = (Q_{\%} \cdot K^T) / \sqrt{D_k}$

Attention weights: $A = \text{softmax}(E, \text{dim} = 1)$ (Shape: $N_n \times N_n$)

Output vectors: $Y = AV$ (Shape: $N_n \times D_v$) $Y_{\%} = \sum A_{\%} V$



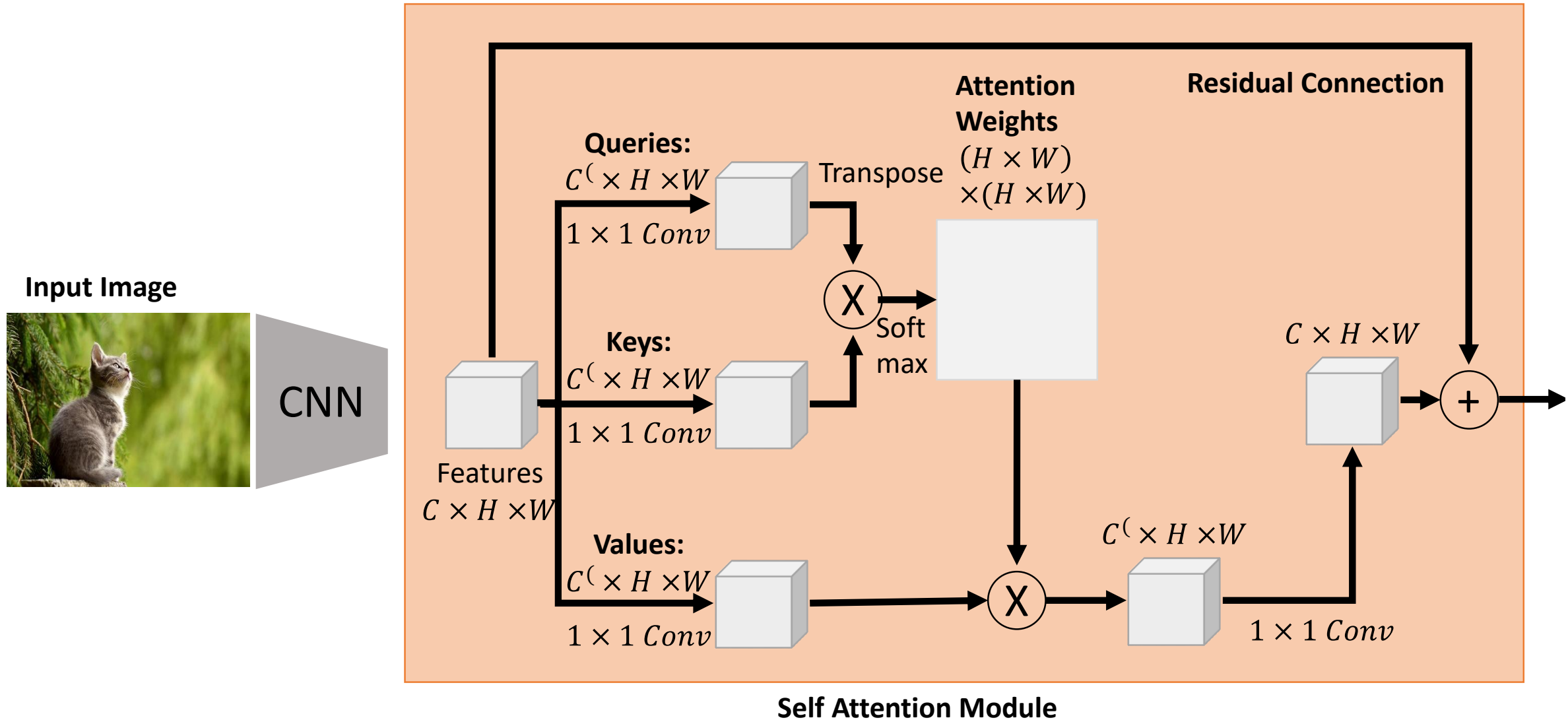
Blank



Self Attention Layer



CNN with Attention



Blank Slide

Convolution Layer Vs SA Layer



Self Attention: Comments

- Self attention learns the relationship between elements in a sequence.
 - say between words in a sentence
- Self Attention Vs Convolution
 - Filters are dynamically calculated instead of static filters
 - SA is invariant to changes in the input points
 - SA can operate on irregular inputs
- SA allows to learn global and local features
 - Hierarchical feature learning by cascading

Where we want to go?

The Transformer

Transformer Block:

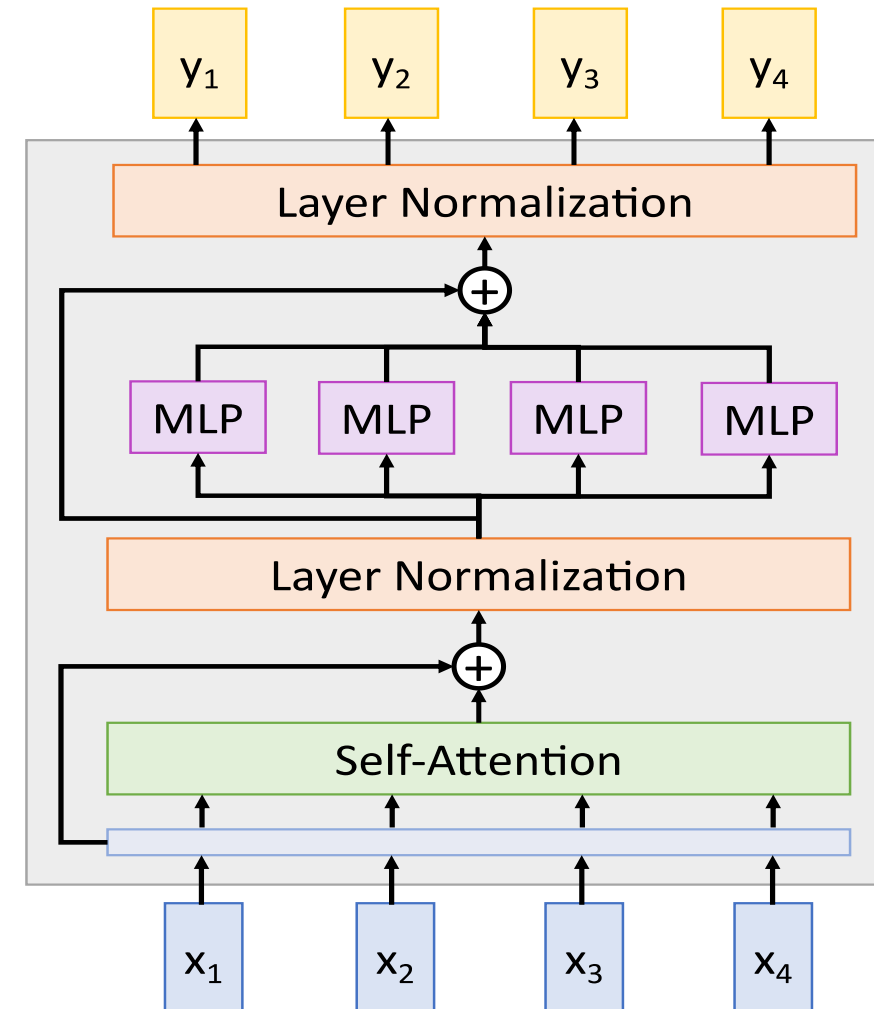
Input: Set of vectors x

Output: Set of vectors y

Self-attention is the only interaction between vectors!

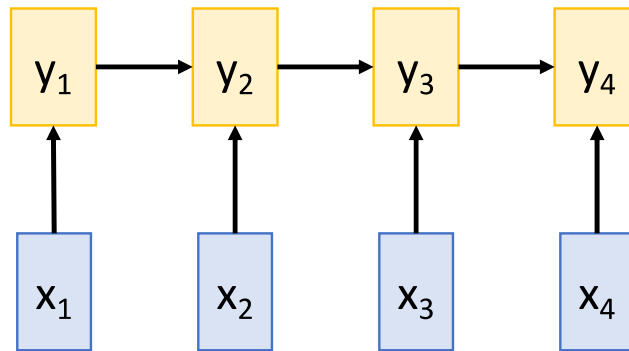
Layer norm and MLP work independently per vector

Highly scalable, highly parallelizable



Conceptual Comparison

Recurrent Neural Network

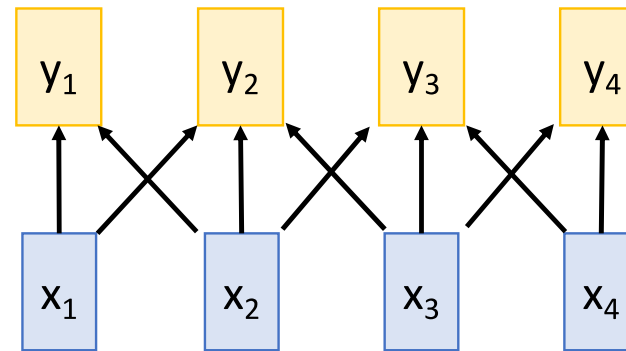


Works on **Ordered Sequences**

(+) **Good at long sequences:** After one RNN layer, h_T "sees" the whole sequence

(-) **Not parallelizable:** need to compute hidden states sequentially

1D Convolution

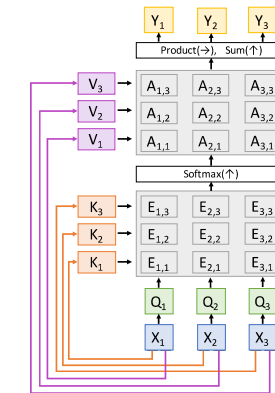


Works on **Multidimensional Grids**

(-) **Bad at long sequences:** Need to stack many conv layers for outputs to "see" the whole sequence

(+) **Highly parallel:** Each output can be computed in parallel

Self-Attention



Works on **Sets of Vectors**

(-) **Good at long sequences:** after one self-attention layer, each output "sees" all inputs!

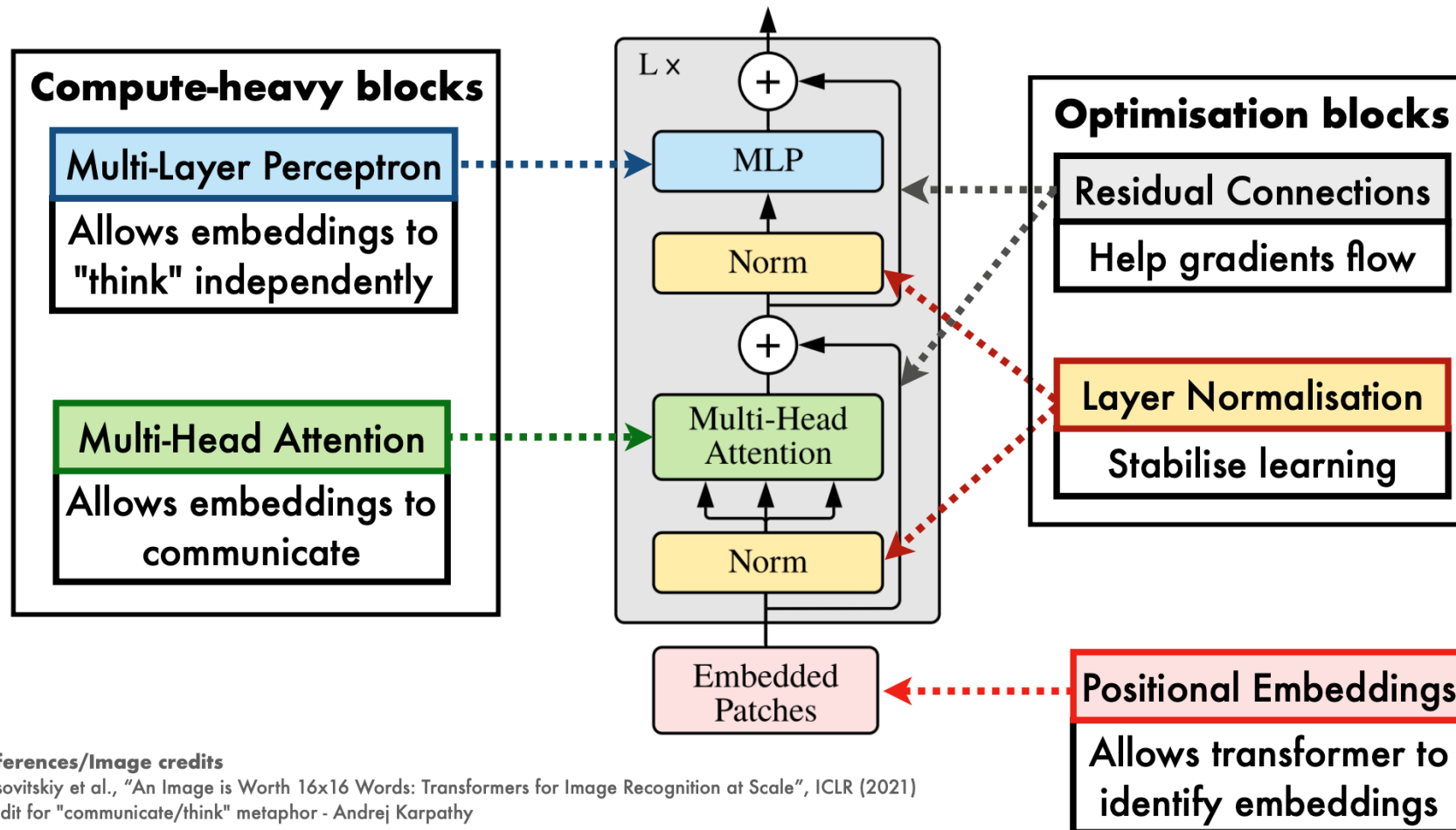
(+) **Highly parallel:** Each output can be computed in parallel

(-) **Very memory intensive**

Transformers

Transformer Encoder

Five key ideas



References/Image credits

Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR (2021)

Credit for "communicate/think" metaphor - Andrej Karpathy

Thanks!!

Questions?