

Lecture 7

Attention Mechanism

Rui Dai

✉ rui.dai@inria.fr

About Me

Rui Dai

- Home page: <https://dairui01.github.io/>
- Ph.D. student at INRIA, STARS team.
- Research topic: “Action detection using Deep Learning”.



Outline

- Introduction to Attention Mechanism
- Attention Modules
- Self-Attention
- Transformer

Section 1

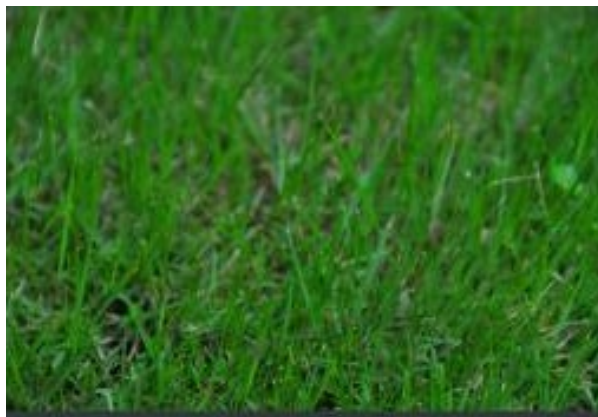
Introduction

What is Attention?

In psychology, attention is the cognitive process of selectively concentrating on one or a few things while ignoring others.

What is Attention?

In psychology, attention is the cognitive process of selectively concentrating on one or a few things while ignoring others.



What is Attention?

In psychology, attention is the cognitive process of selectively concentrating on one or a few things while ignoring others.



What is Attention?

In psychology, attention is the cognitive process of selectively concentrating on one or a few things while ignoring others.



What is Attention?

In psychology, attention is the cognitive process of selectively concentrating on one or a few things while ignoring others.



What is Attention?

In psychology, attention is the cognitive process of selectively concentrating on one or a few things while ignoring others.



Introduction



What object is it?

Introduction



What object is it?

What is Attention?

With Attention Mechanism:

Selectively concentrating on a few relevant things, while ignoring others in deep neural networks.



What object is it?

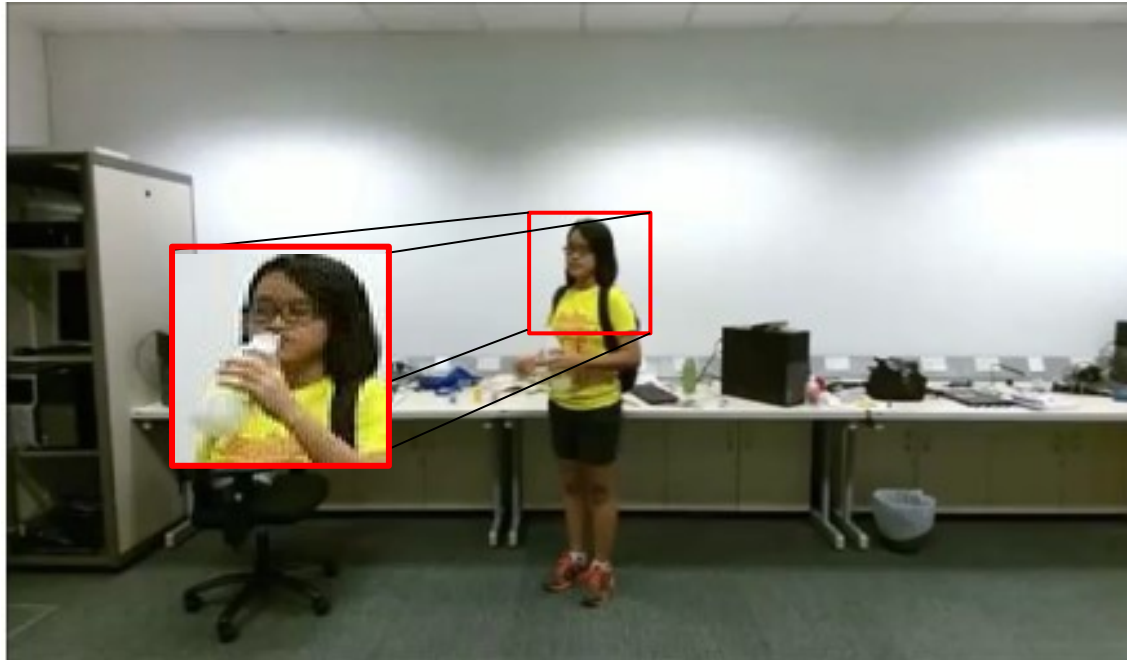
Example in Action Recognition

The girl is drinking water from a bottle



Do we really need the whole video to infer that?

Example in Action Recognition

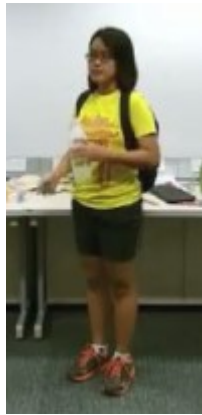


➤ Isn't this enough for an inference?

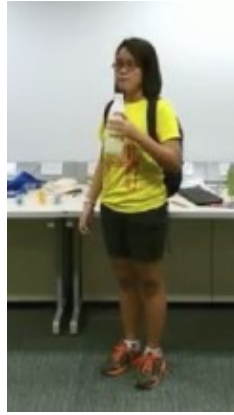
Focus in the **Spatial** space is required!

Example in Action Recognition

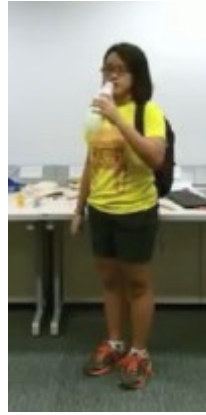
Time -1



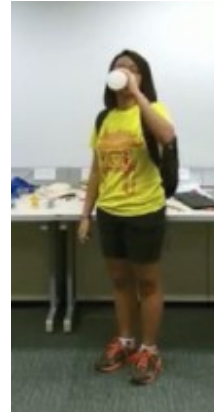
Time -2



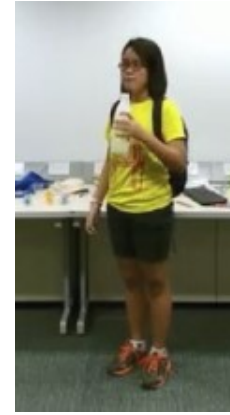
Time -3



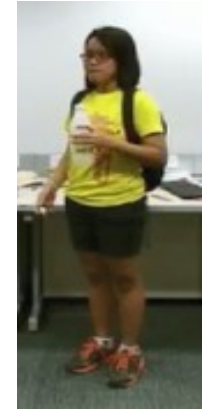
Time -4



Time -5

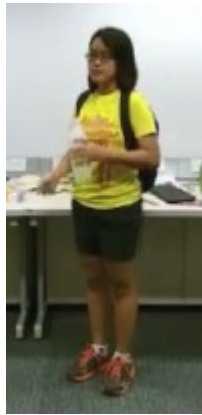


Time 6

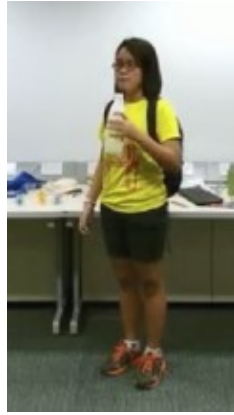


Example in Action Recognition

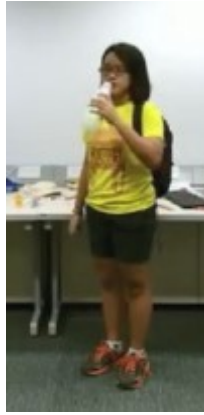
Time -1



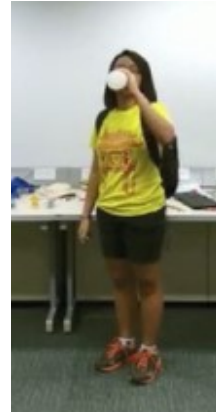
Time -2



Time -3



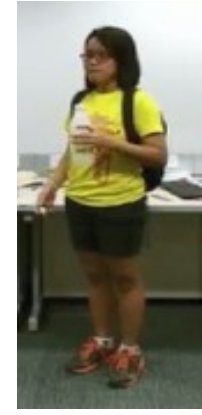
Time -4



Time -5



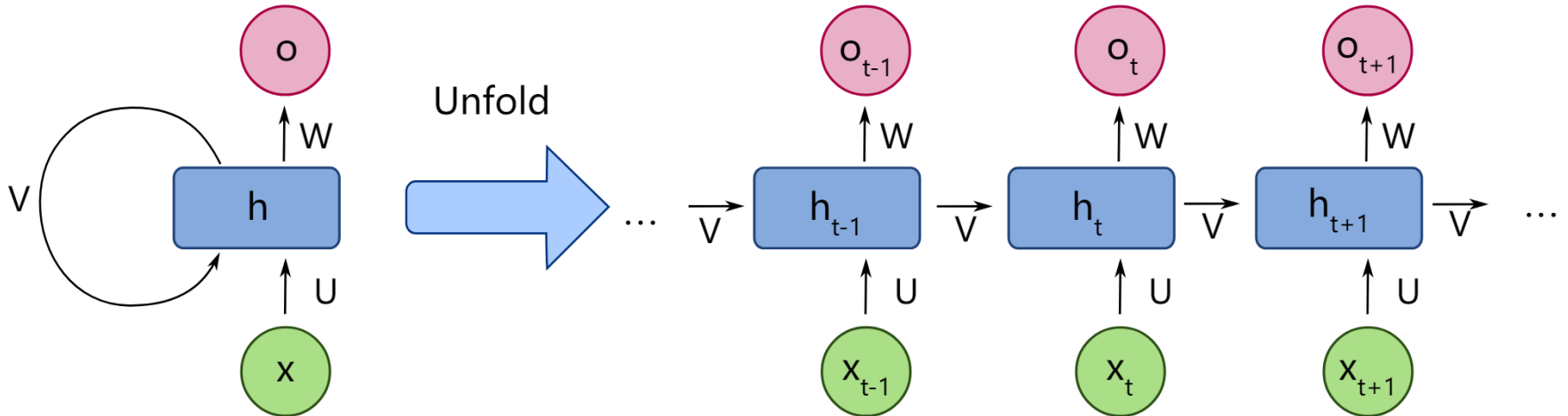
Time 6



Section 2

Attention Modules

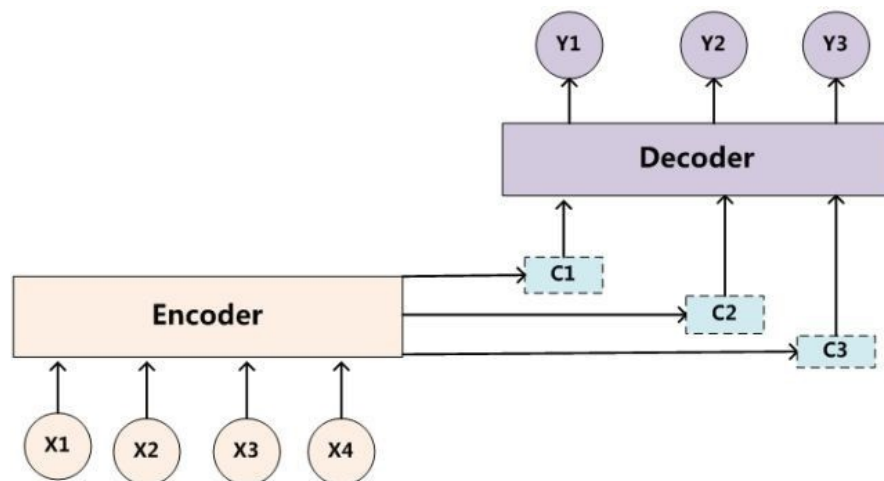
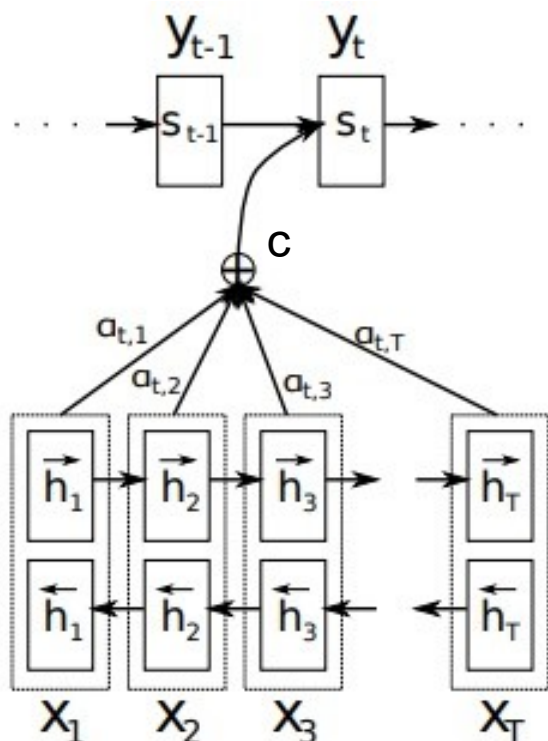
RNN



Drawbacks:

Can not learn long-term information.

Attention in RNN

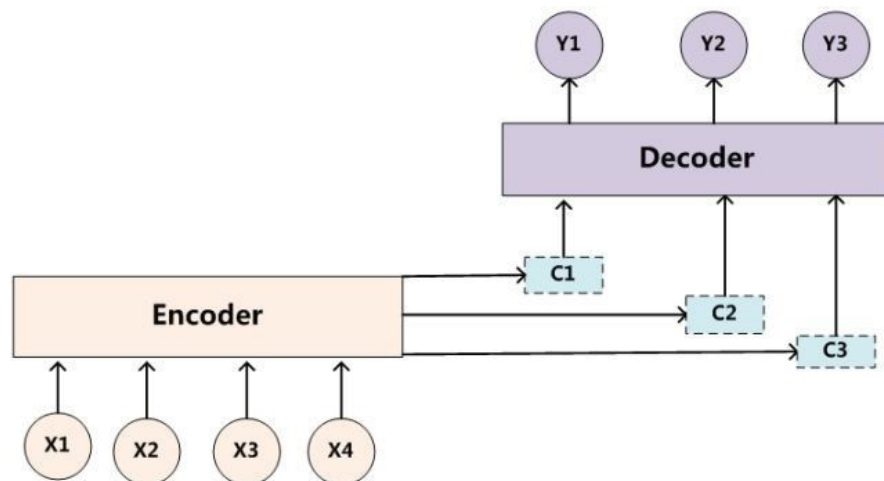
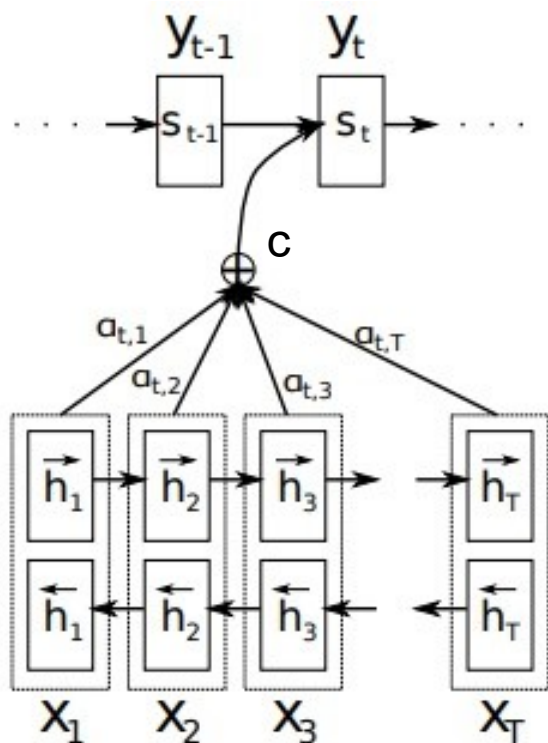


$$e_{ij} = a(s_{i-1}, h_j) \\ = v_a^T \tanh(W_a s_{i-1} + U_a h_j)$$

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^T \exp(e_{ik})}$$

a compute the correlation between the target s and input h

Attention in RNN



$$c_i = \sum_{j=1}^T a_{ij} h_j$$

$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$

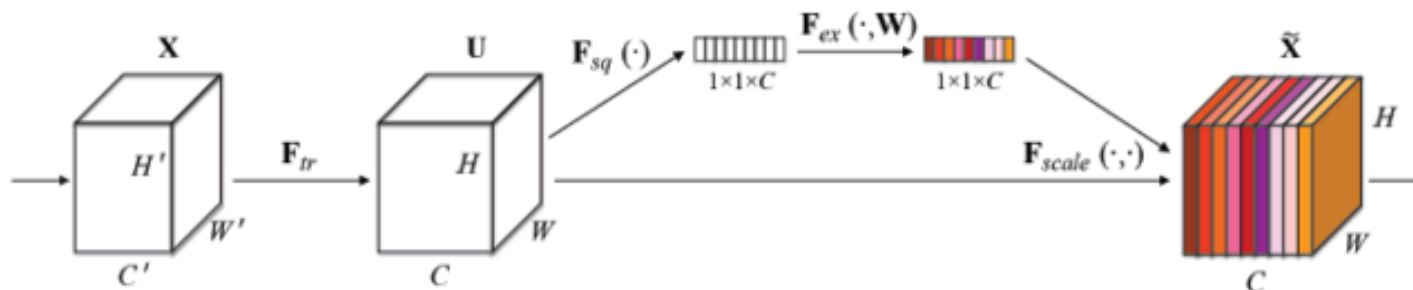
$$y_i = g(y_{i-1}, s_i, c_i)$$

Squeeze-and-Excitation Attention

Attention on Channels

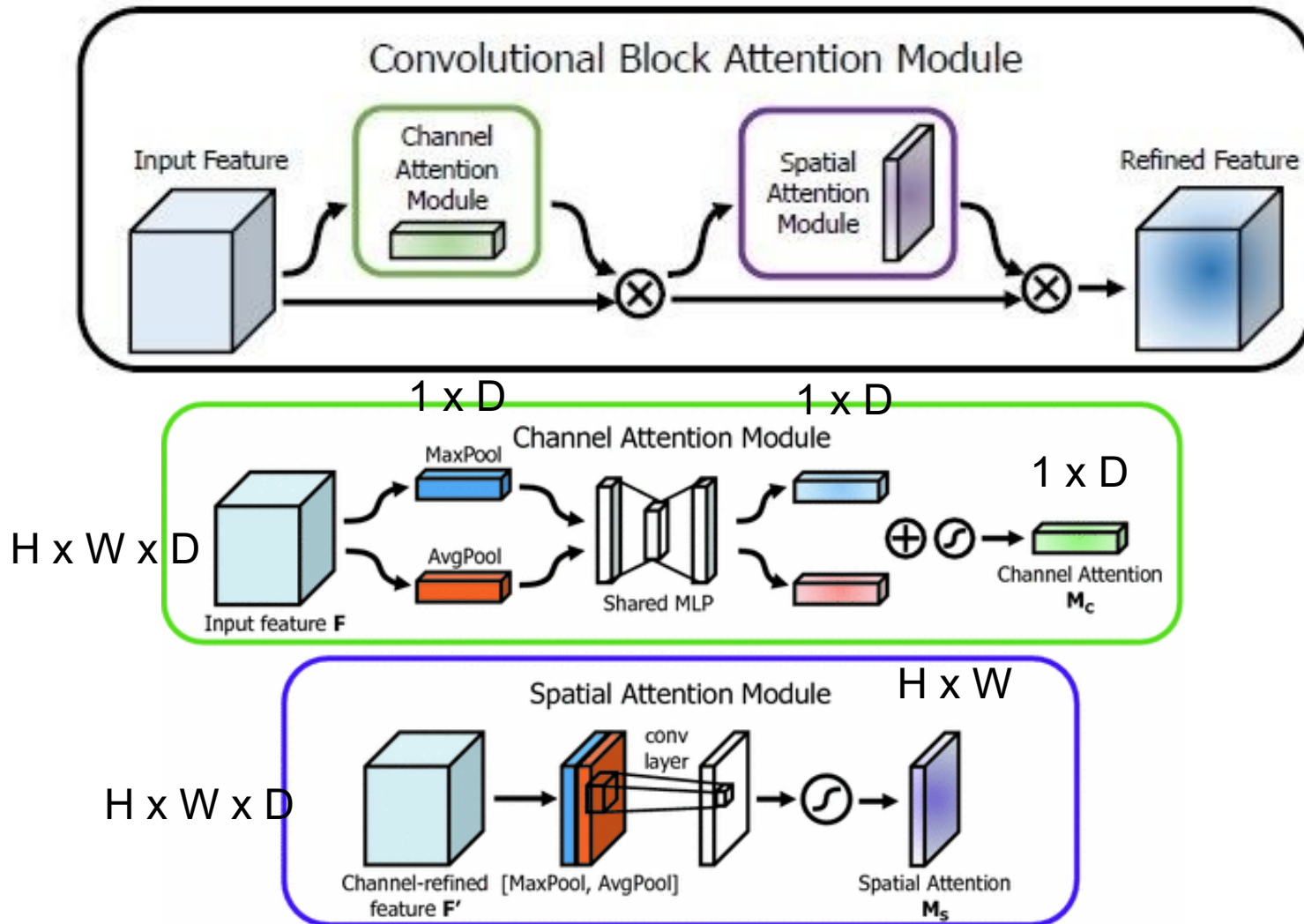
Squeeze: Average Pool

Excitation: FC + Sigmoid

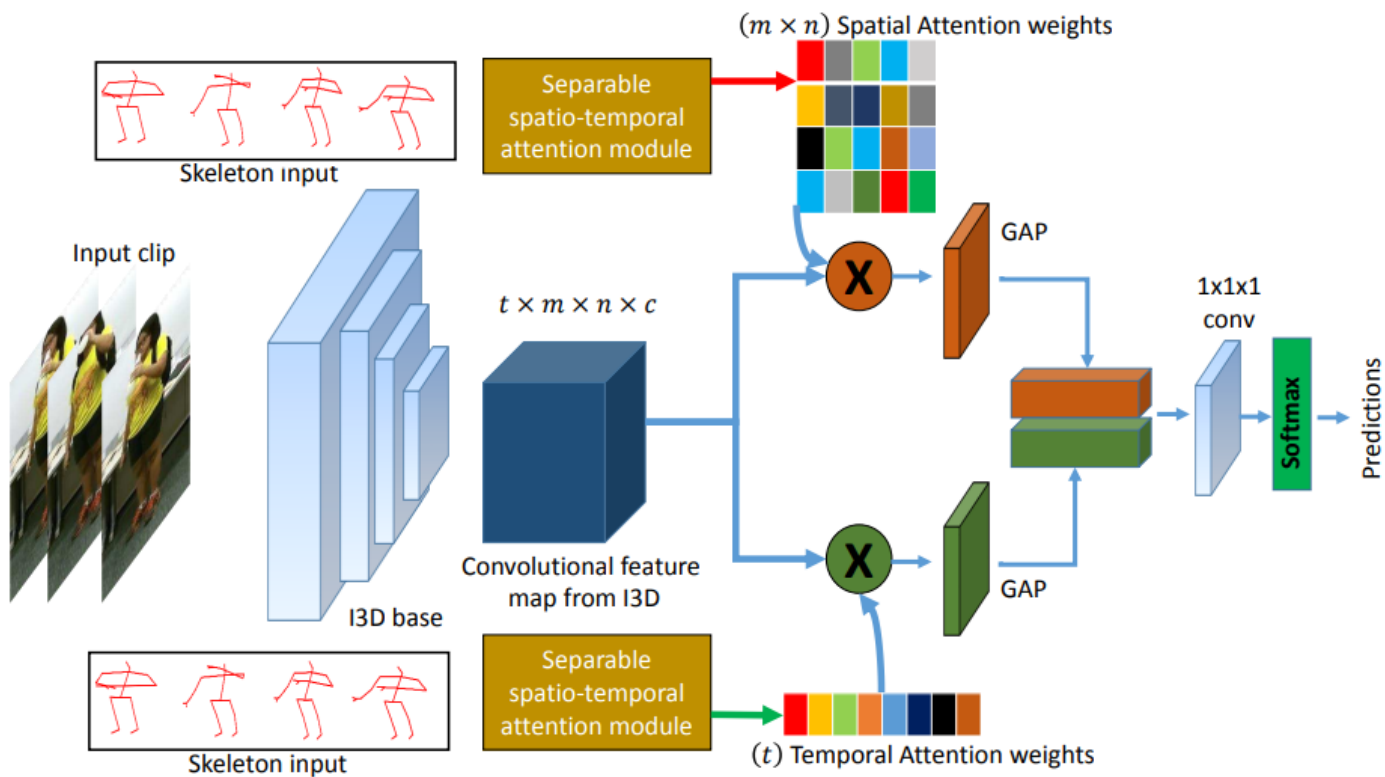


Convolutional Block Attention Module

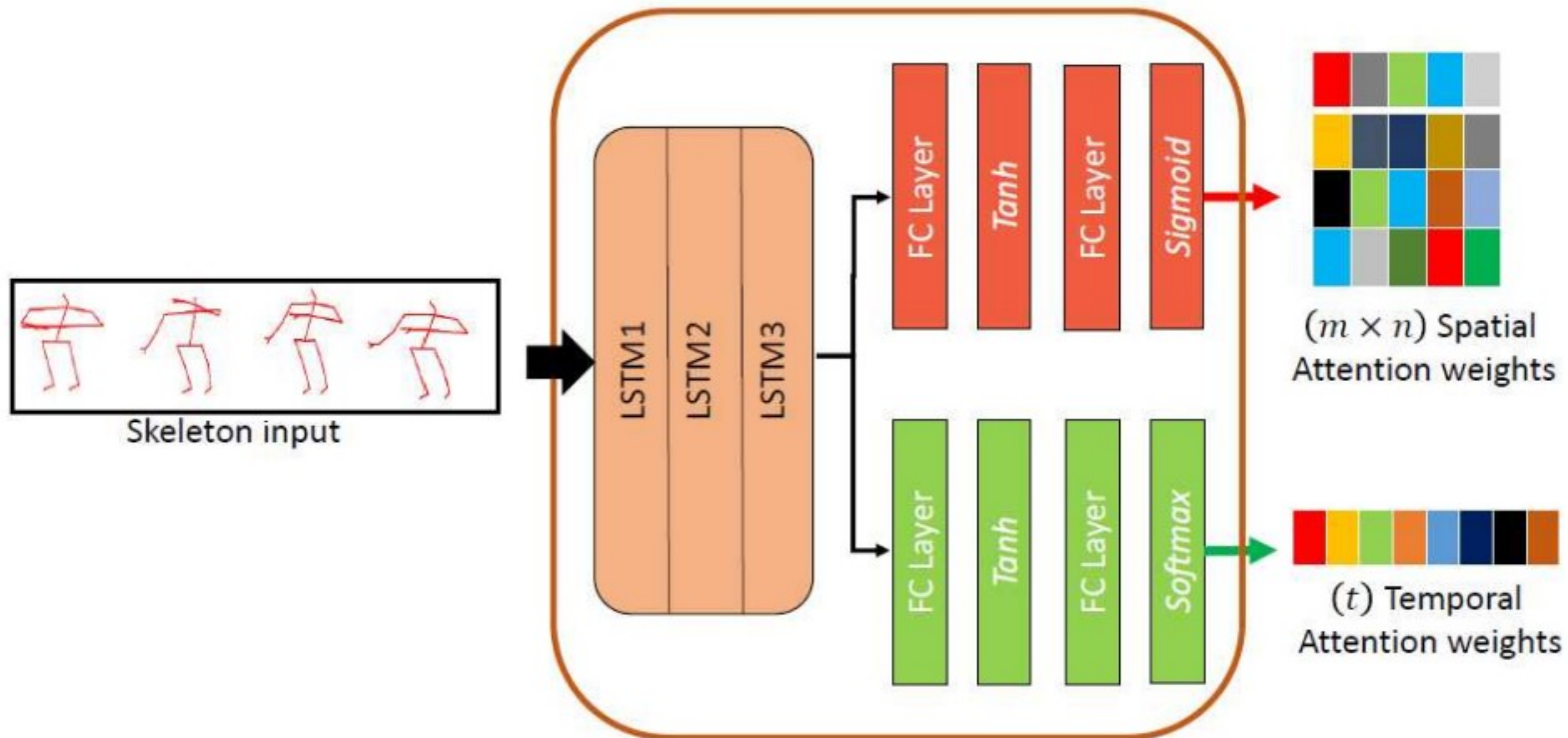
Channel + Spatial



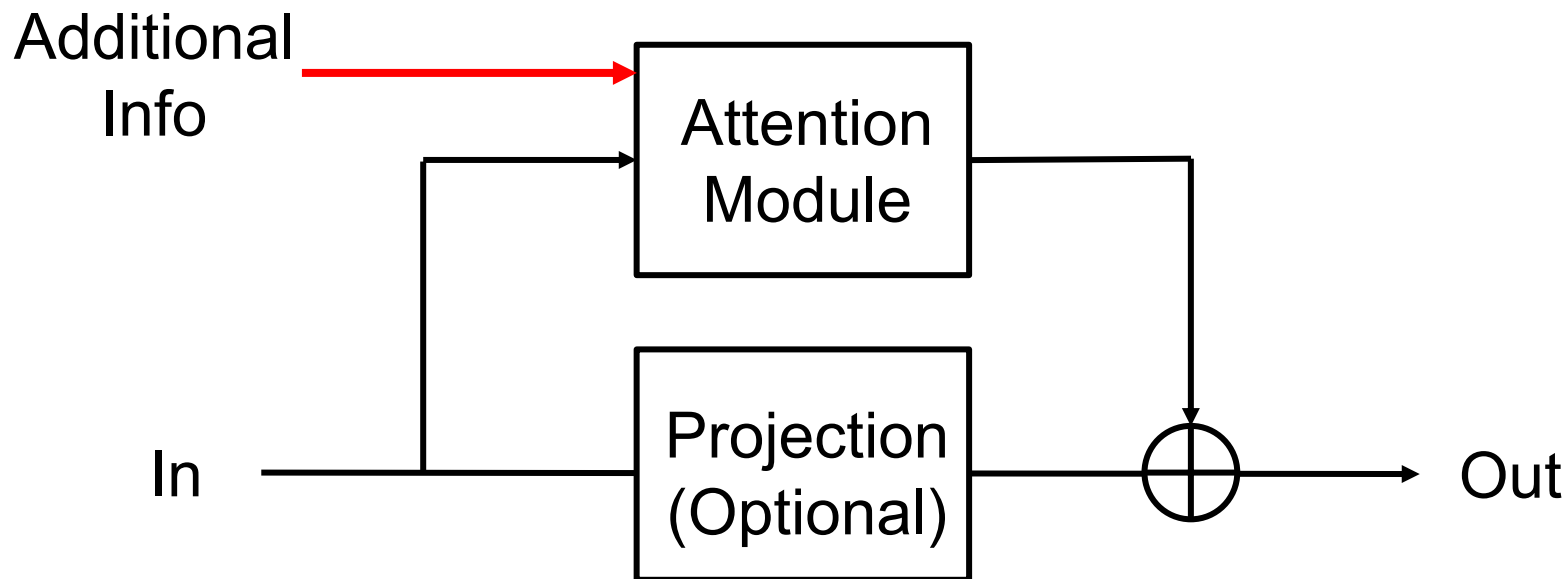
Spatial-Temporal Attention Network



Spatial-Temporal Attention Network



Attention Module Structure

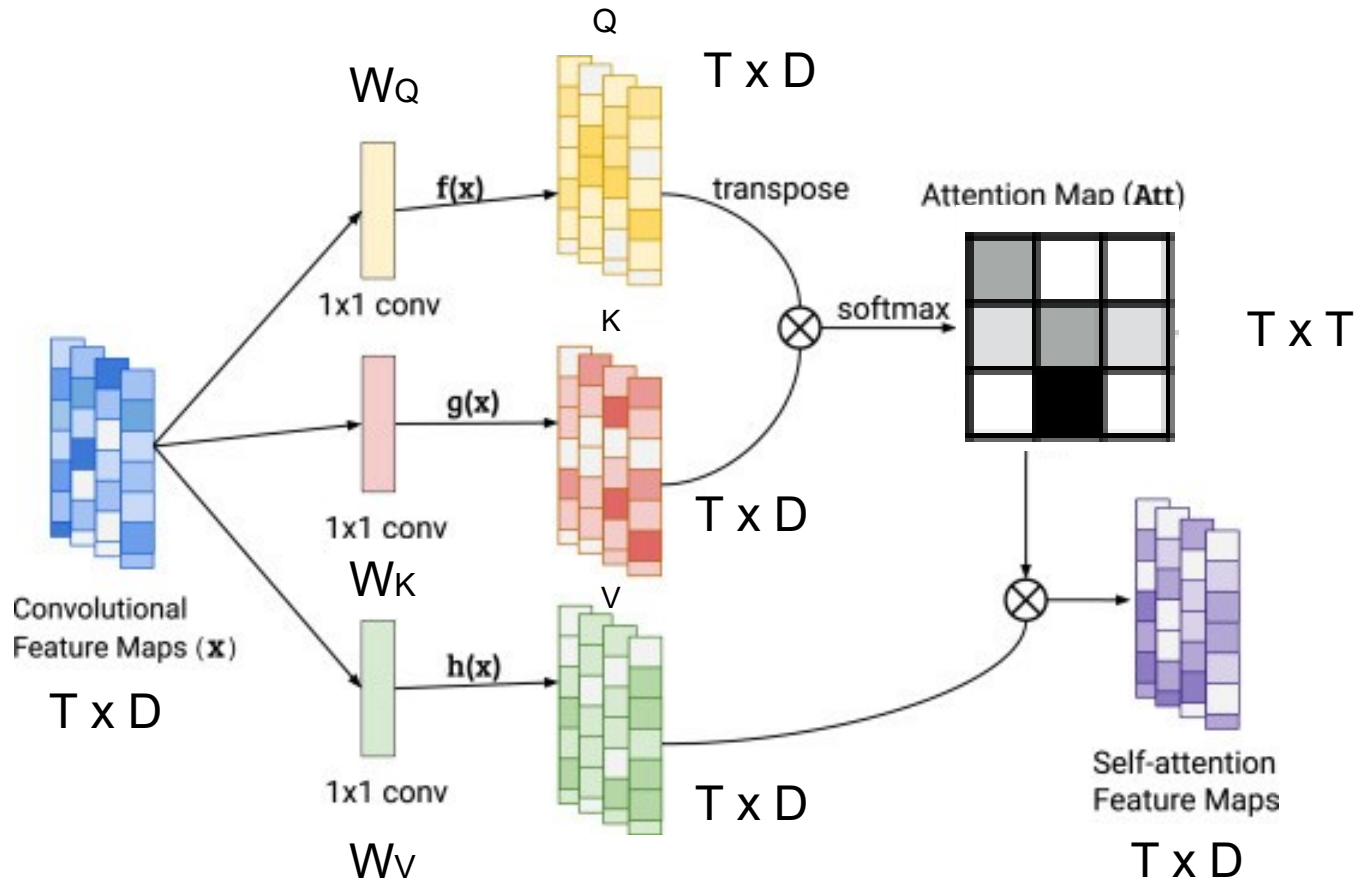


Section 3

Self-Attention

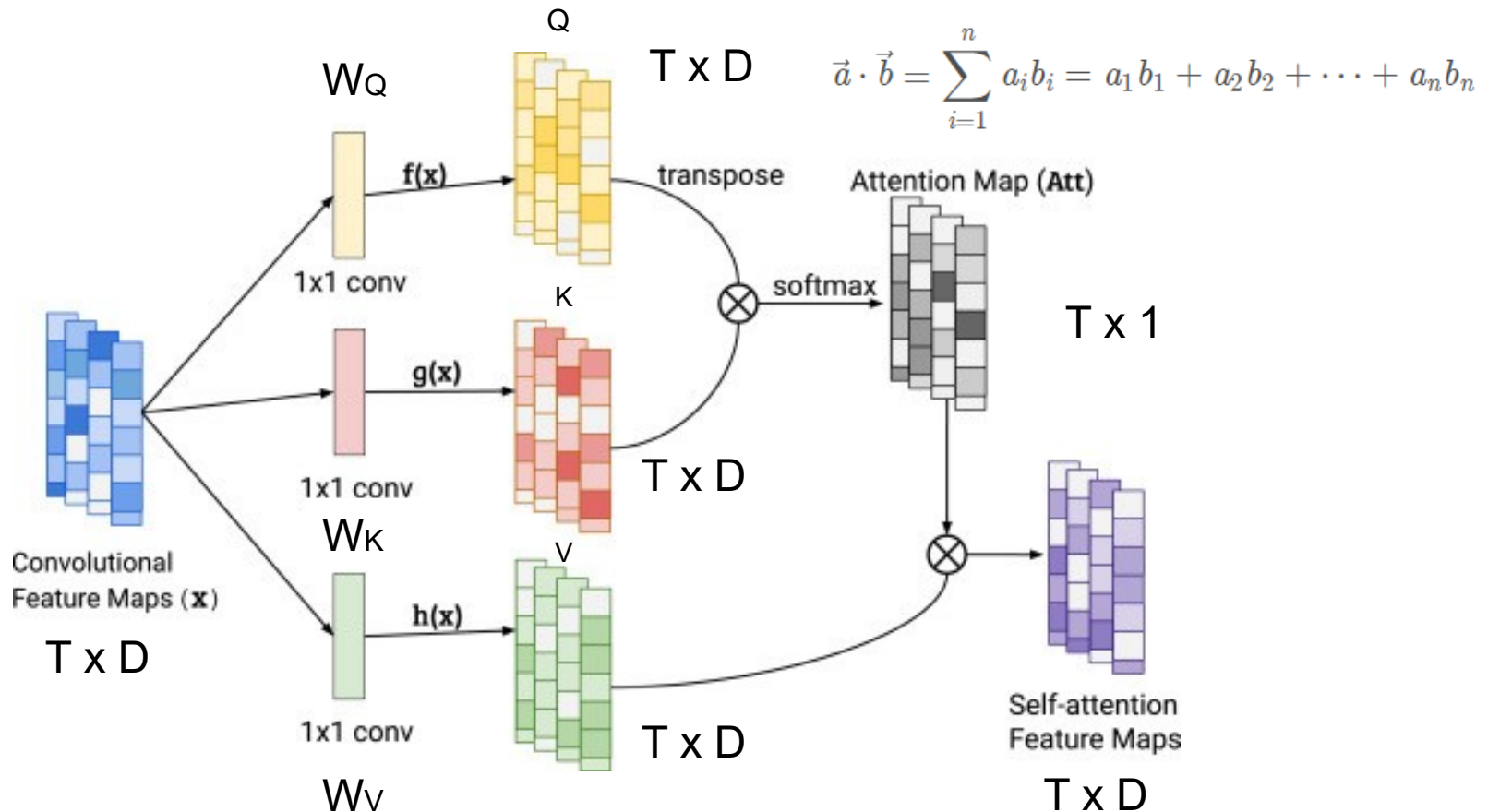
Self-Attention (1D)

Compute the Correlation in a sequence



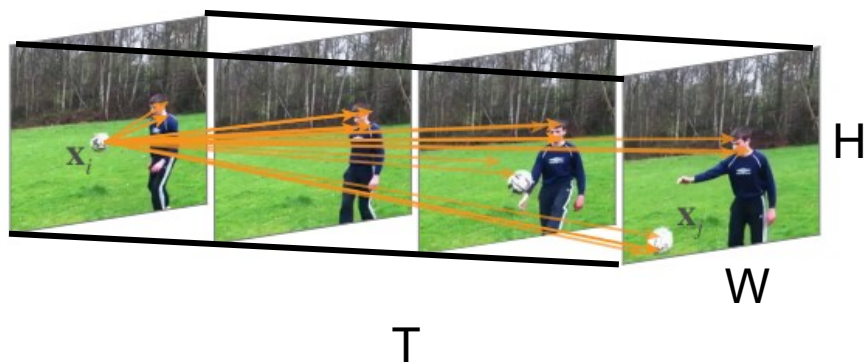
Self-Attention (1D)

Compute the Correlation in a sequence

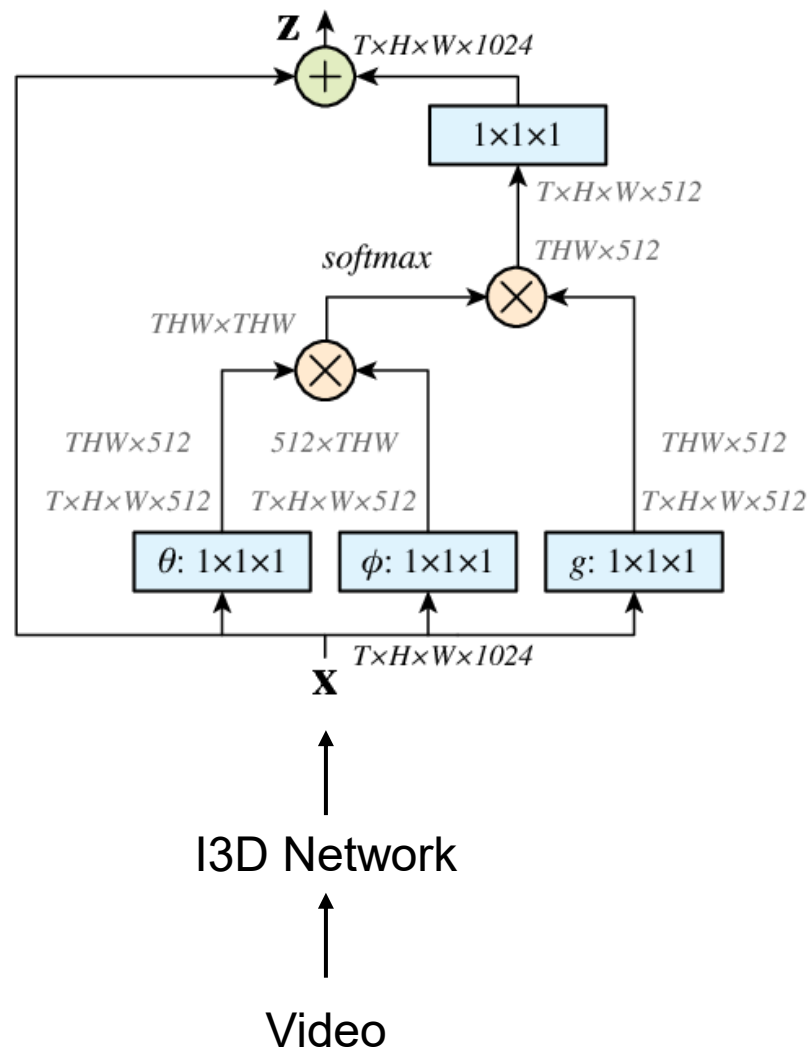


Dot product version

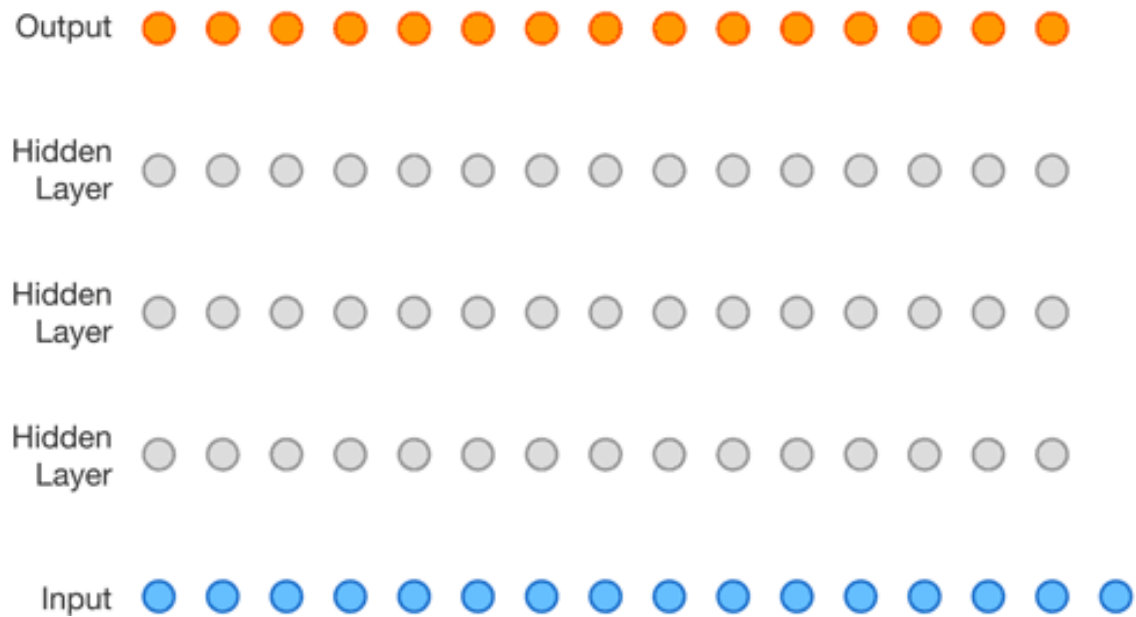
Application- Non local block (3D)



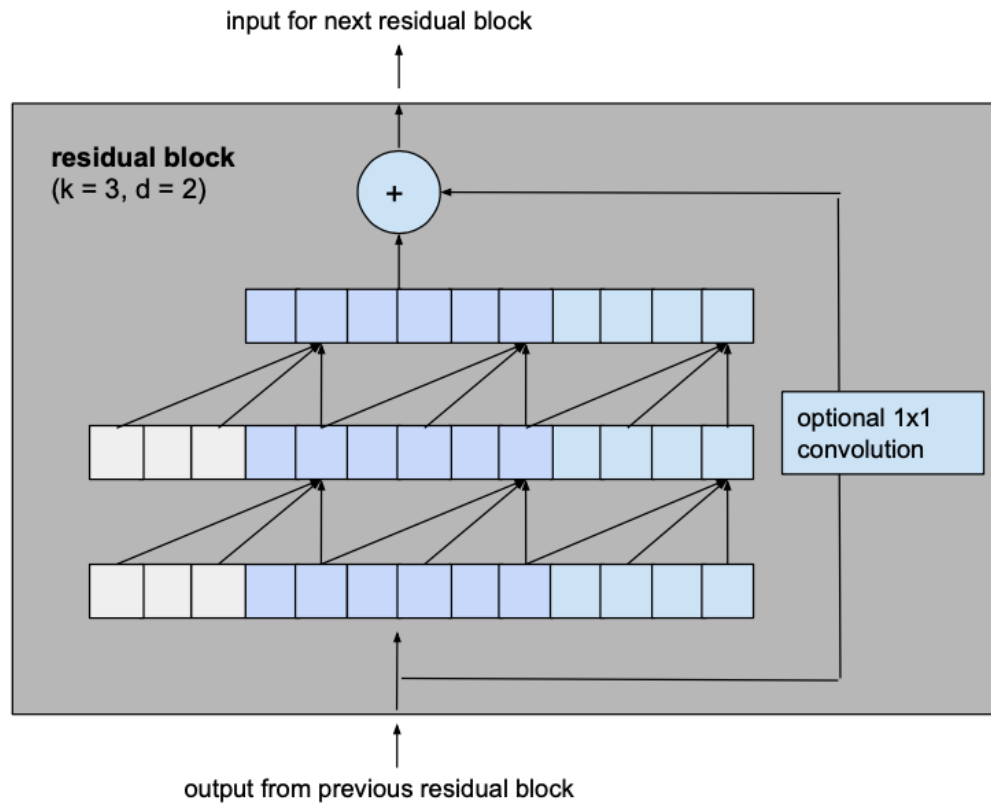
Statistic the Global information
as complementary to the local embeddings



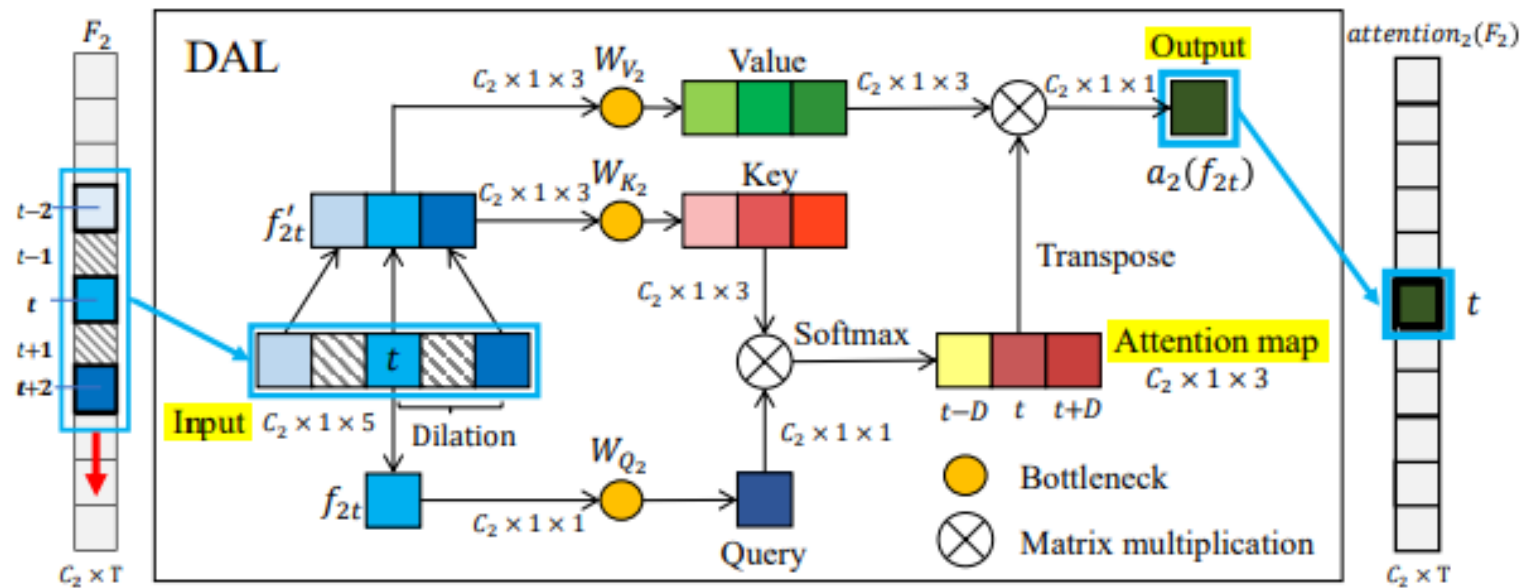
1 Dimensional Convolution



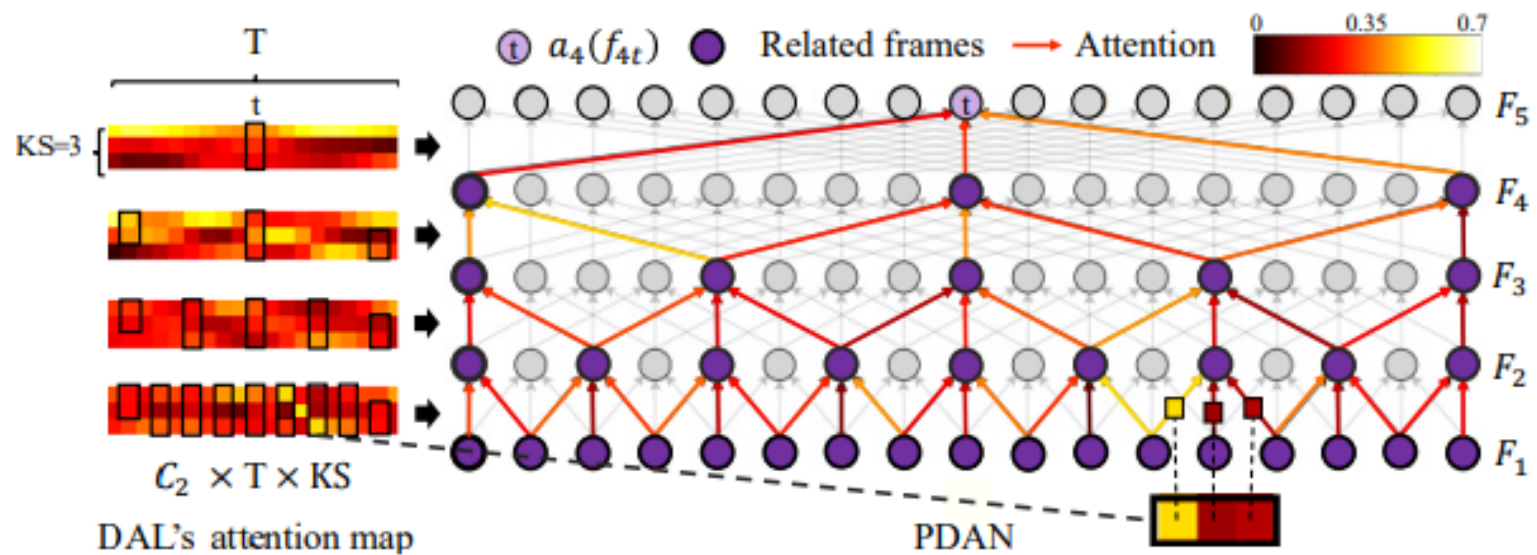
1 Dimensional Convolution



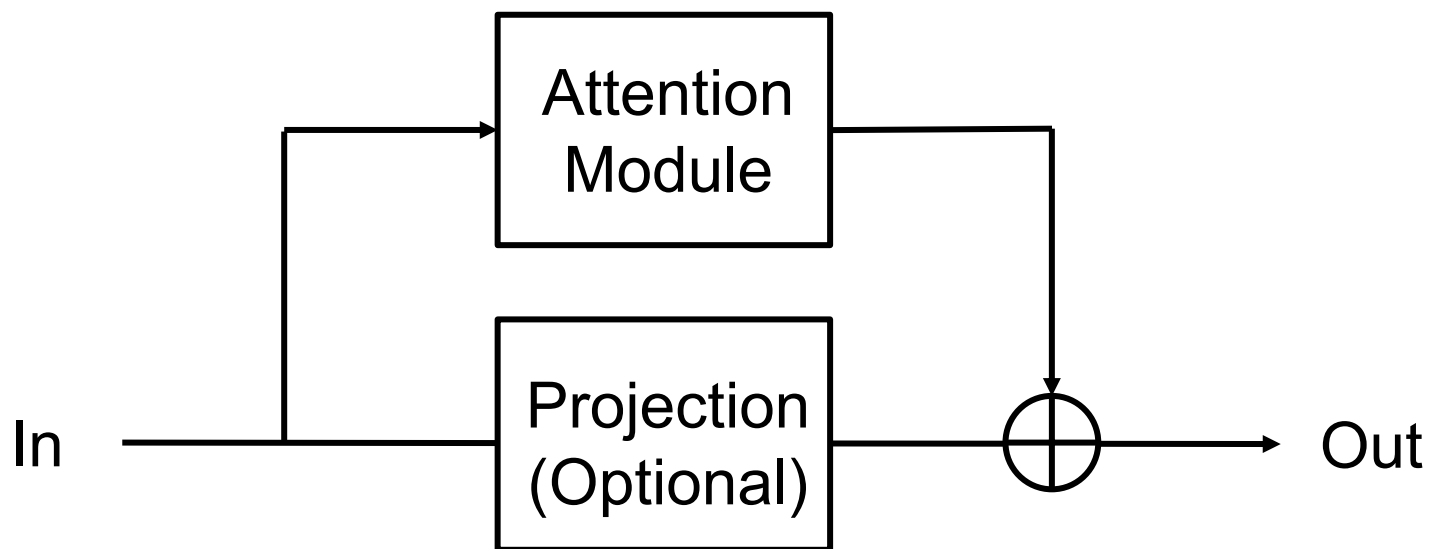
Application- PDAN (1D)



Application- PDAN (1D)



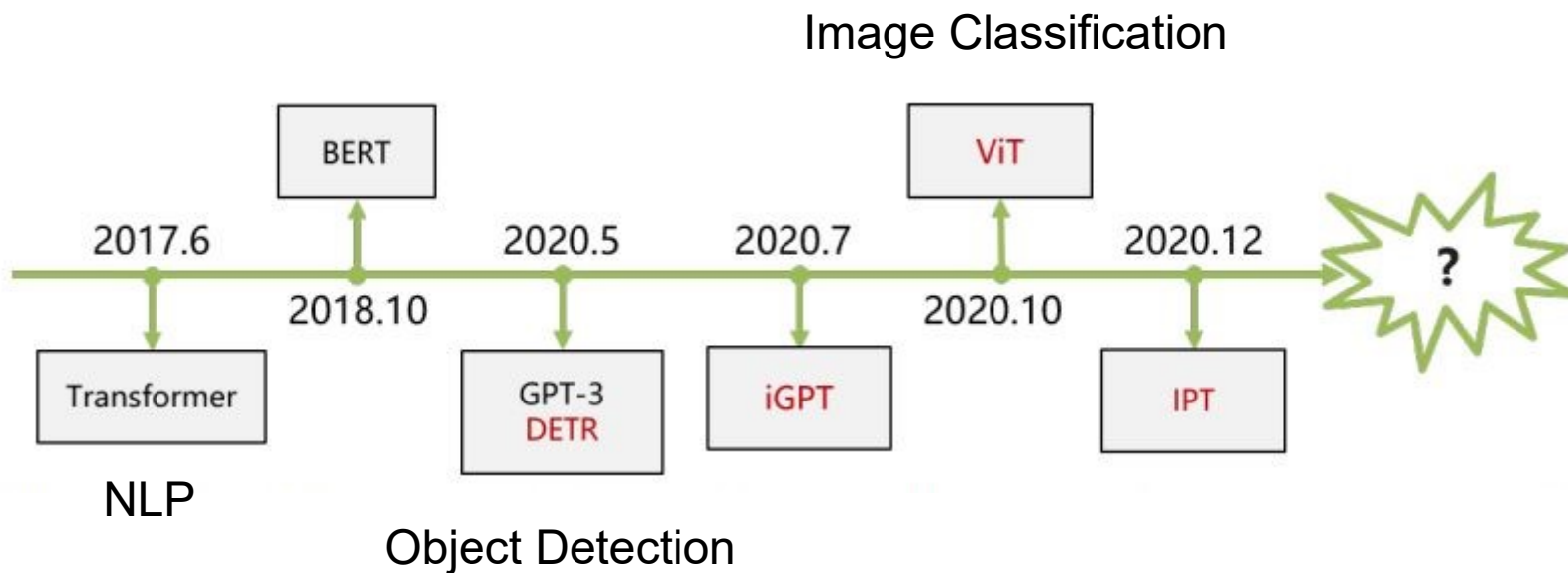
Attention Mechanism Structure

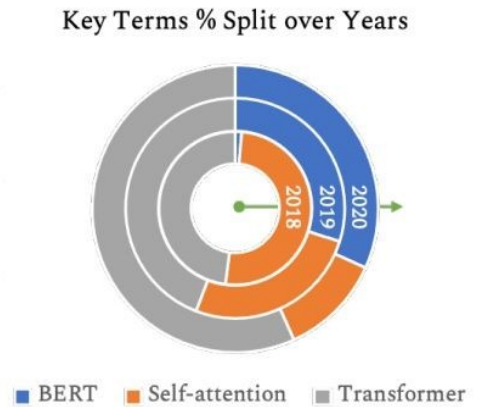
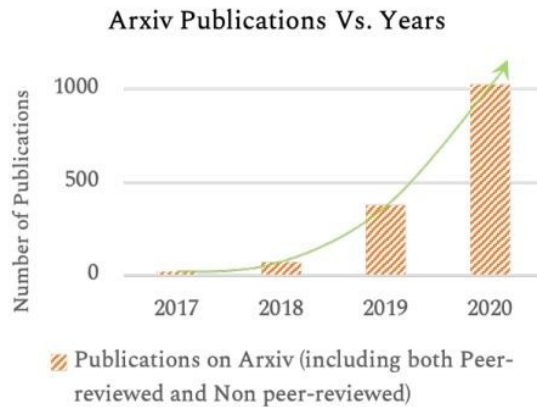
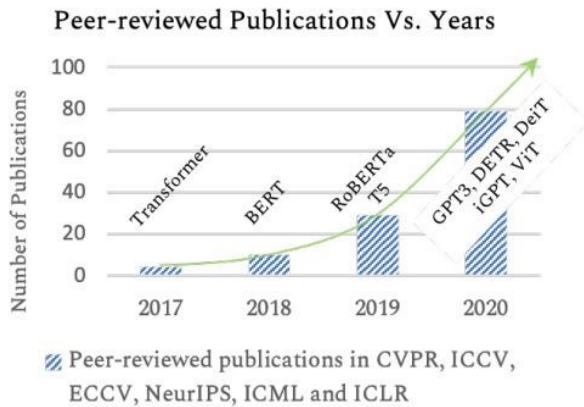


Section 4

Transformer

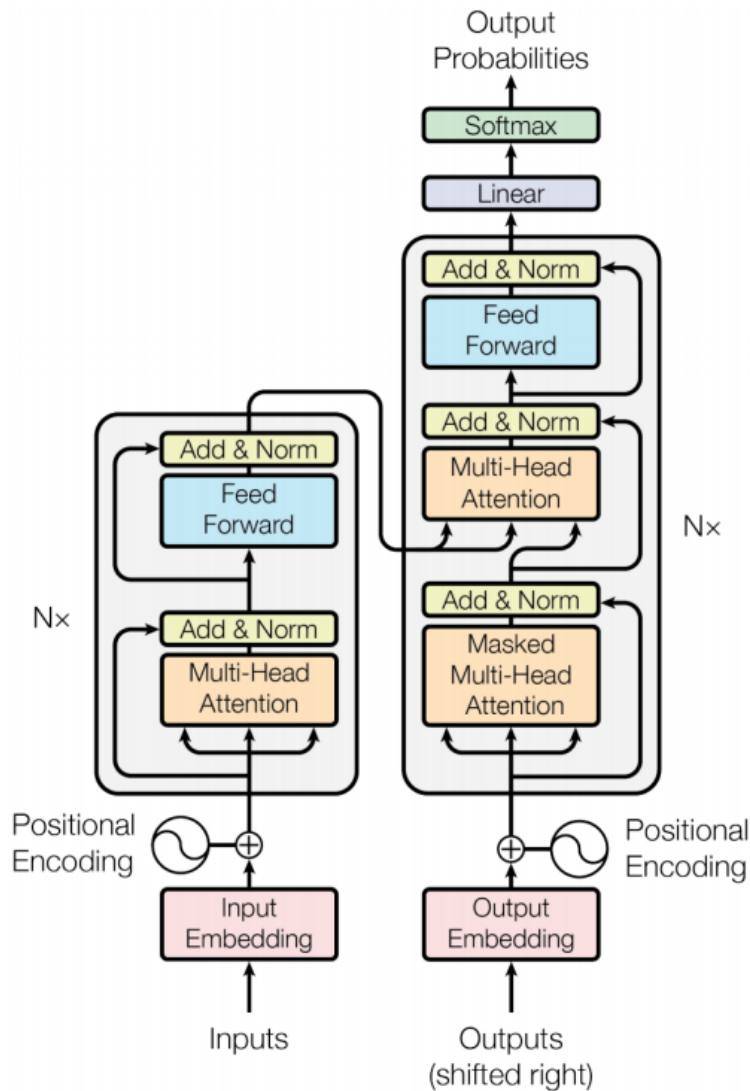
History





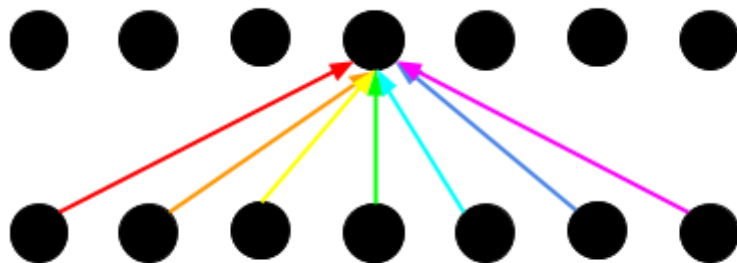
Transformer

- Positional Encoding
- Multi-head attention
- Feed Forward
- Outputs

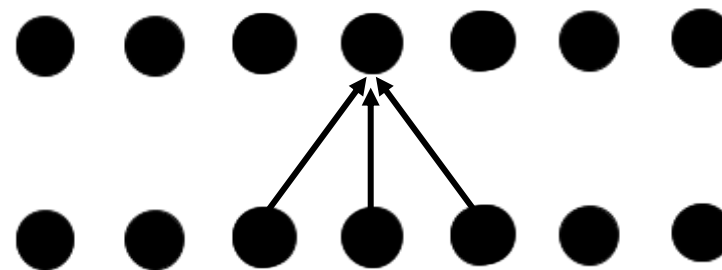


Transformer and Convolution

Transformer layer



Convolution layer



Pros:

Global relation

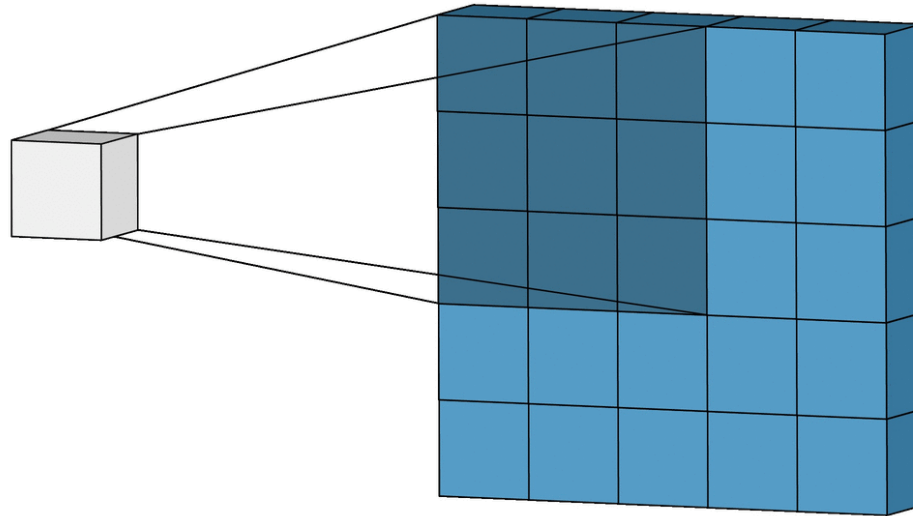
Attention enhanced

Cons:

More Flops

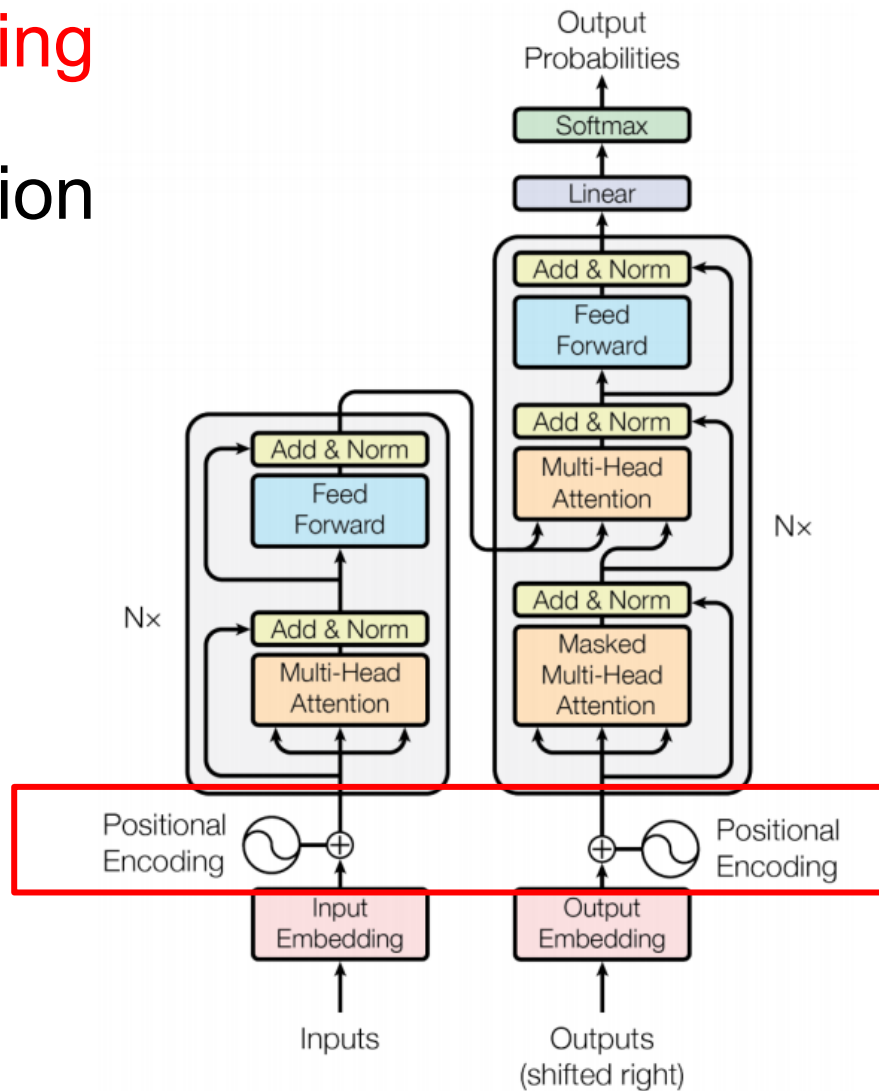
Loss location info

Missing location info



Transformer

- Positional Encoding
- Multi-head attention
- Feed Forward
- Outputs



Positional Encoding

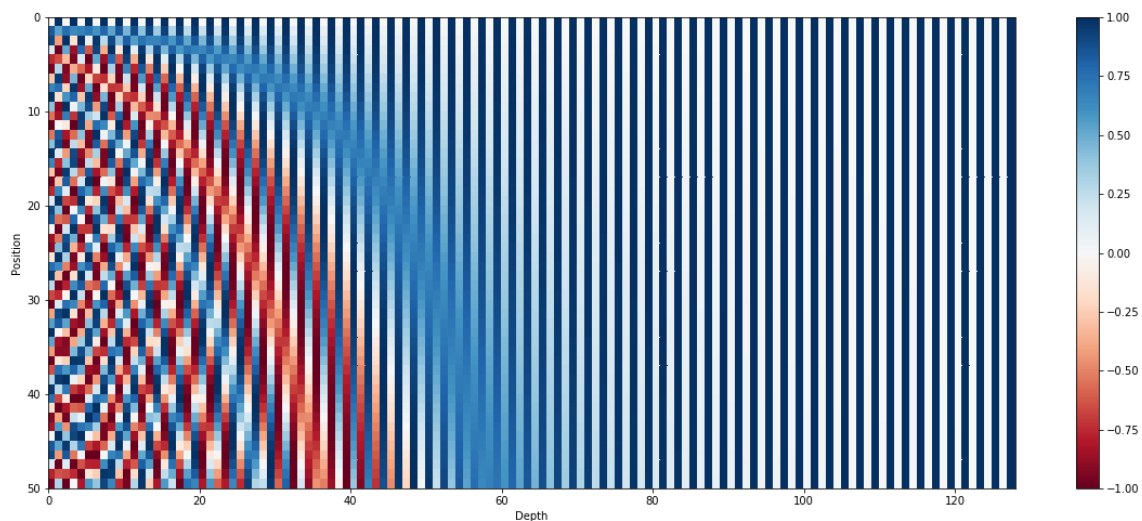
Indicating the positional information

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

Positional Encoding

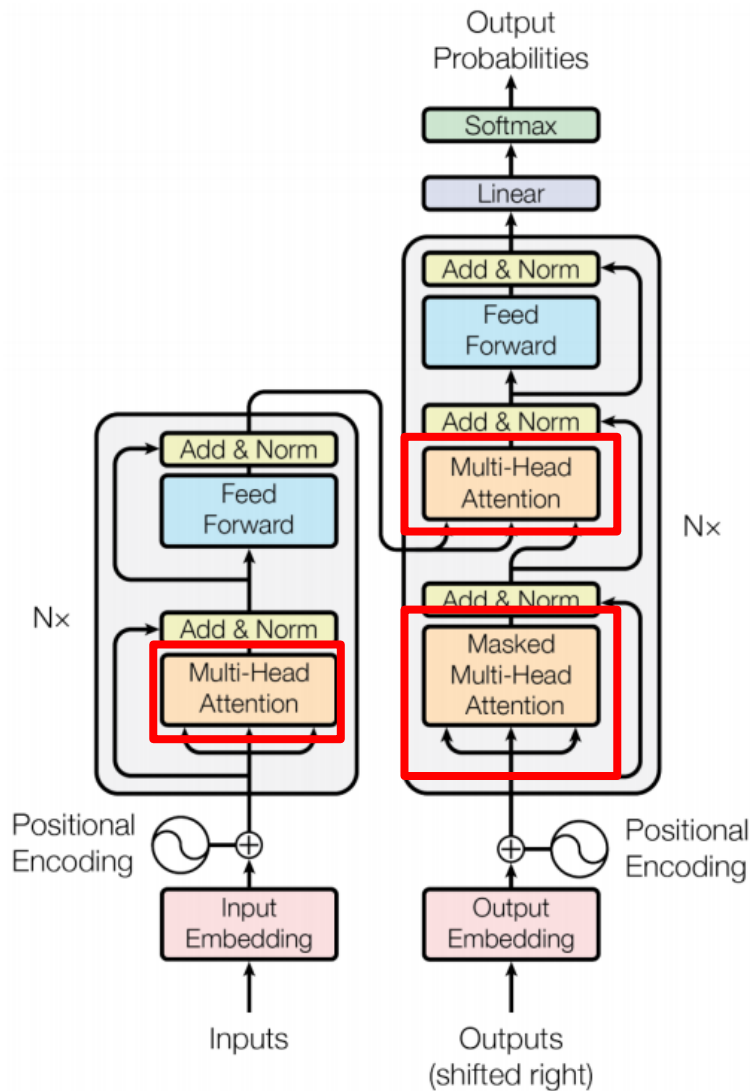
$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases} \quad \omega_k = \frac{1}{10000^{2k/d}}$$

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

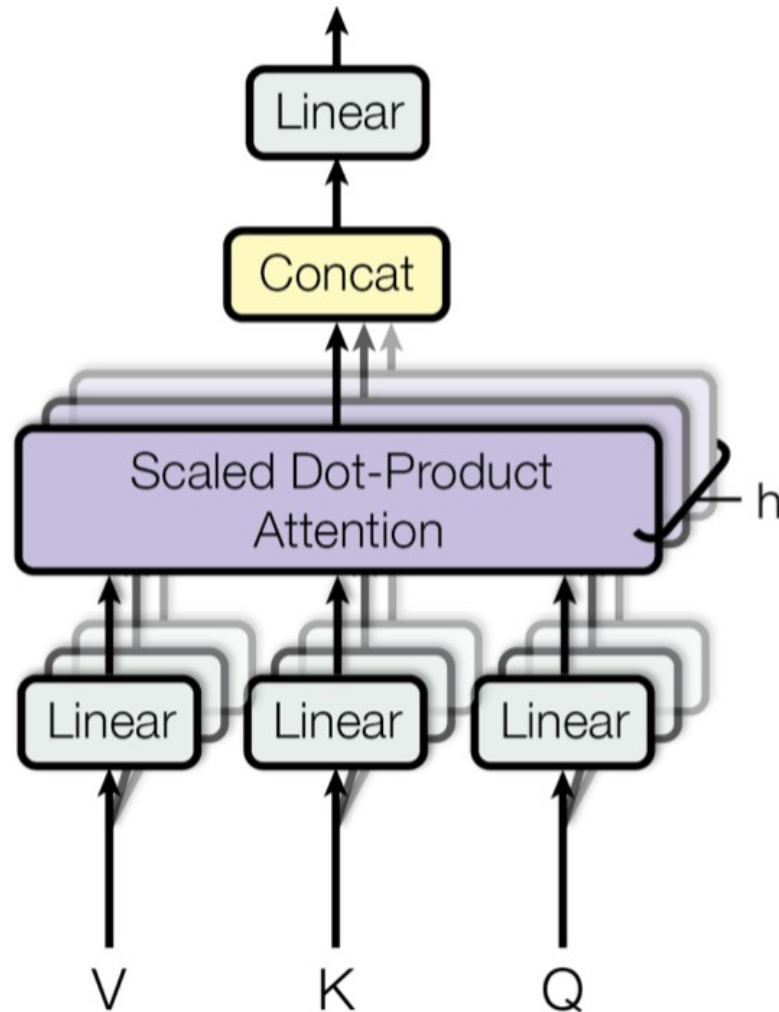


Transformer

- Positional Encoding
- Multi-head attention
- Feed Forward
- Outputs

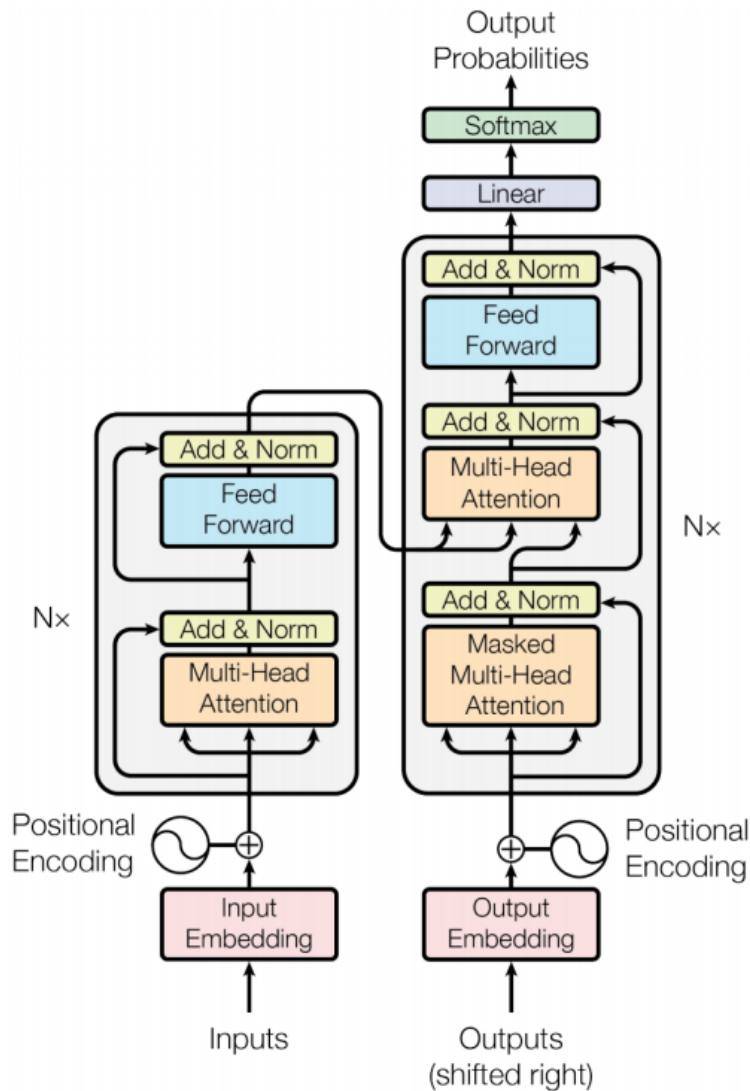


Multi-head Attention



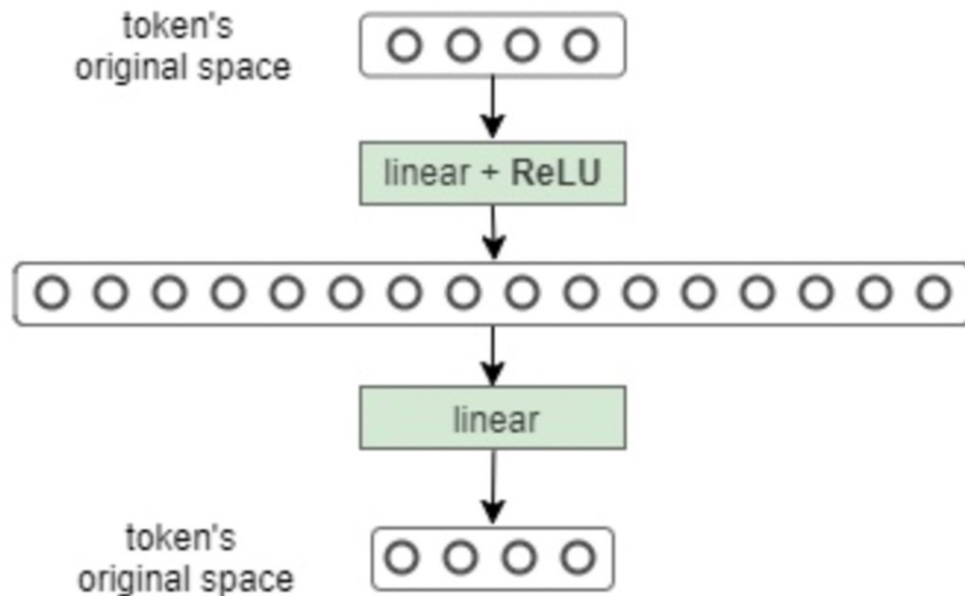
Transformer

- Positional Encoding
- Multi-head attention
- Feed Forward
- Outputs



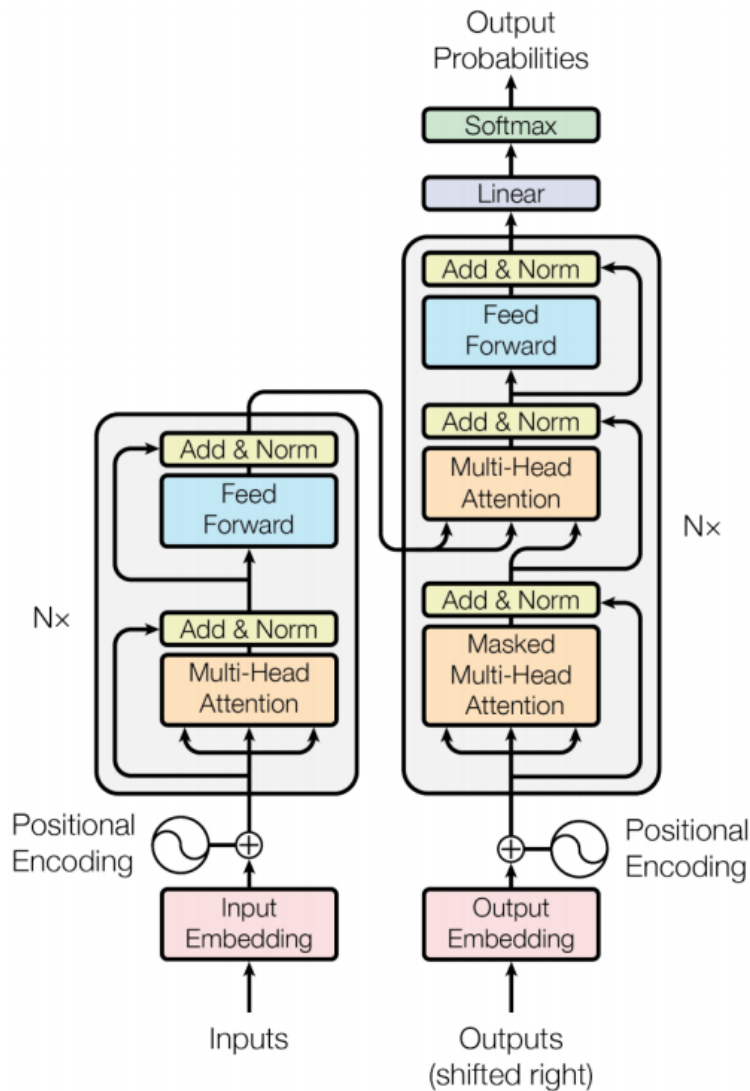
Transformer

Project the feature into a larger space, for extracting the features easier (Similar to SVM)

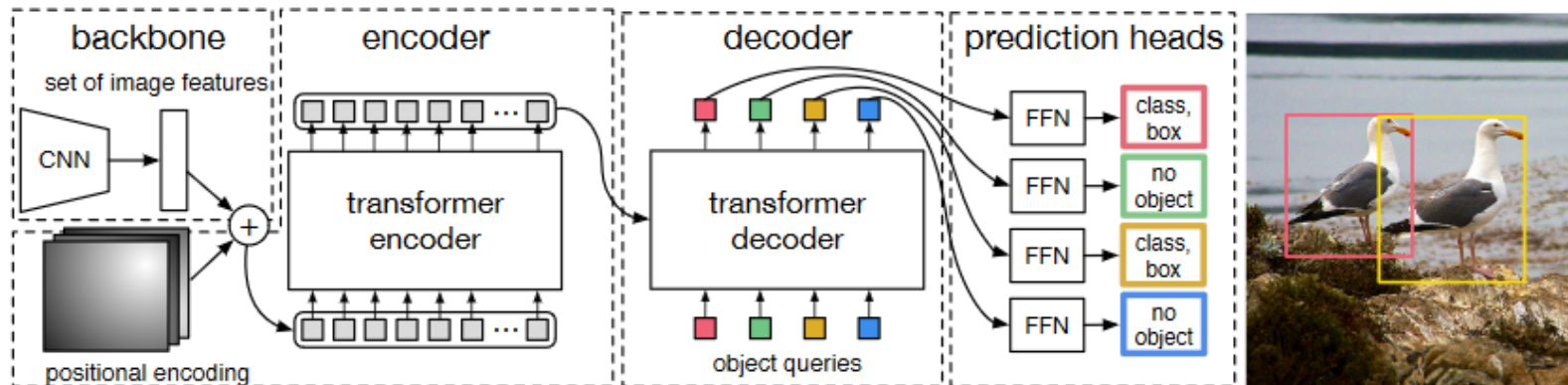


Transformer

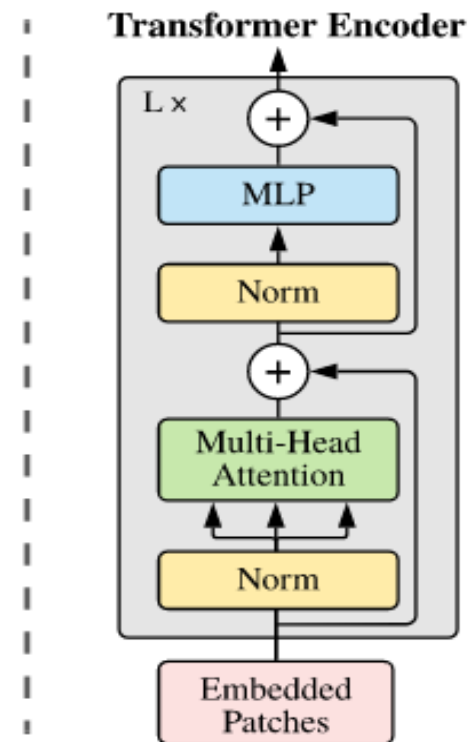
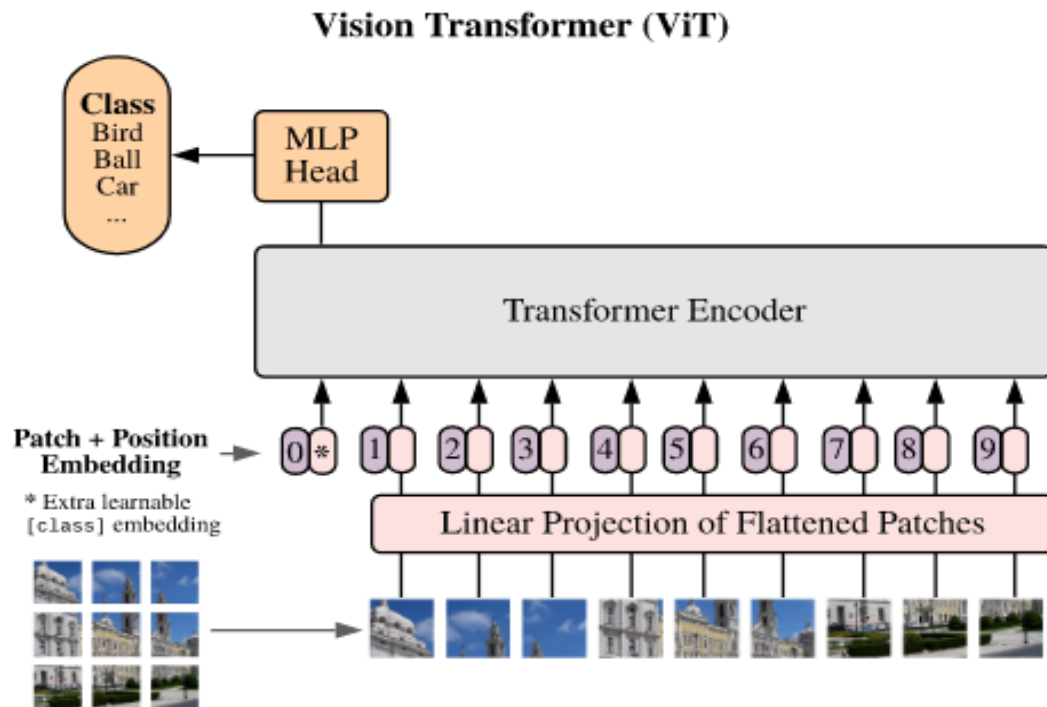
- Positional Encoding
- Multi-head attention
- Feed Forward
- **Outputs**



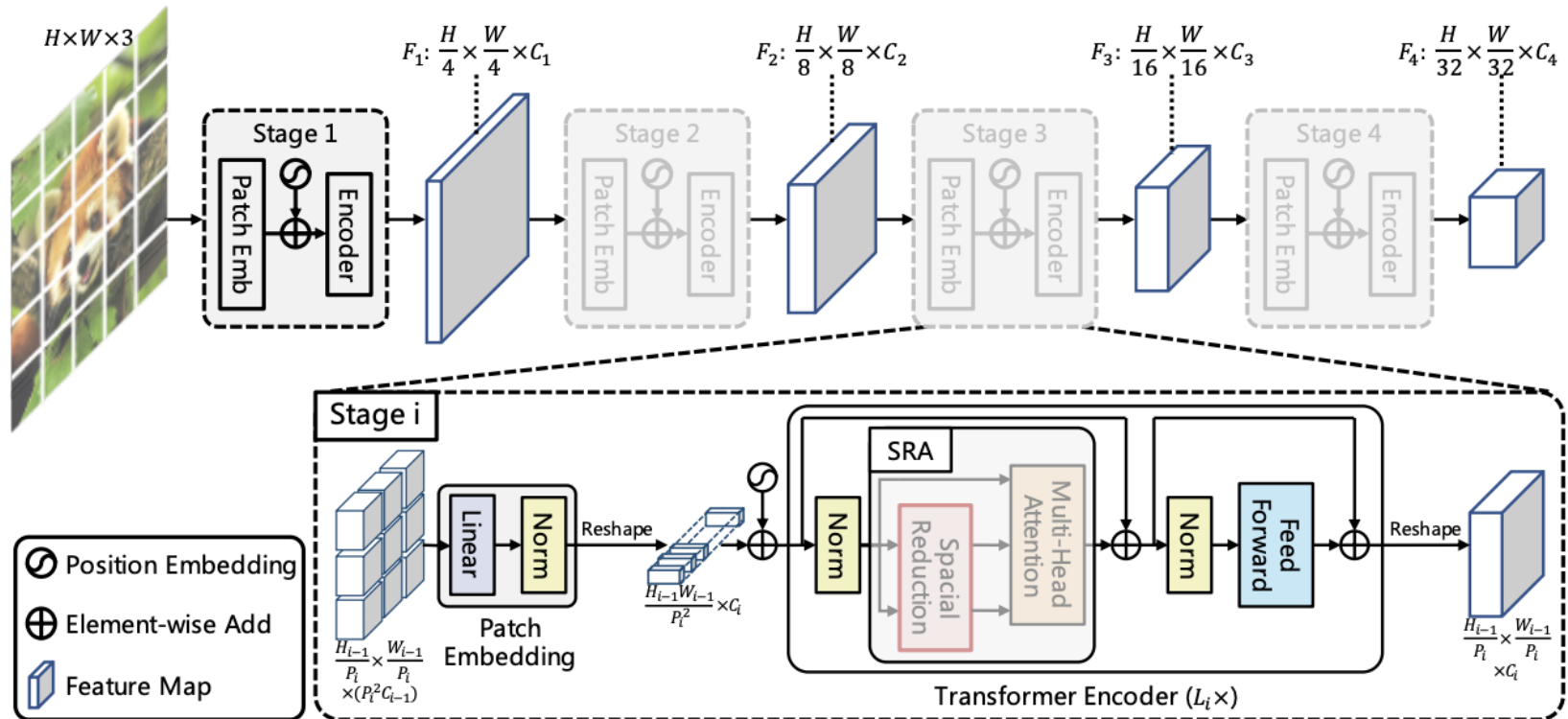
DETR



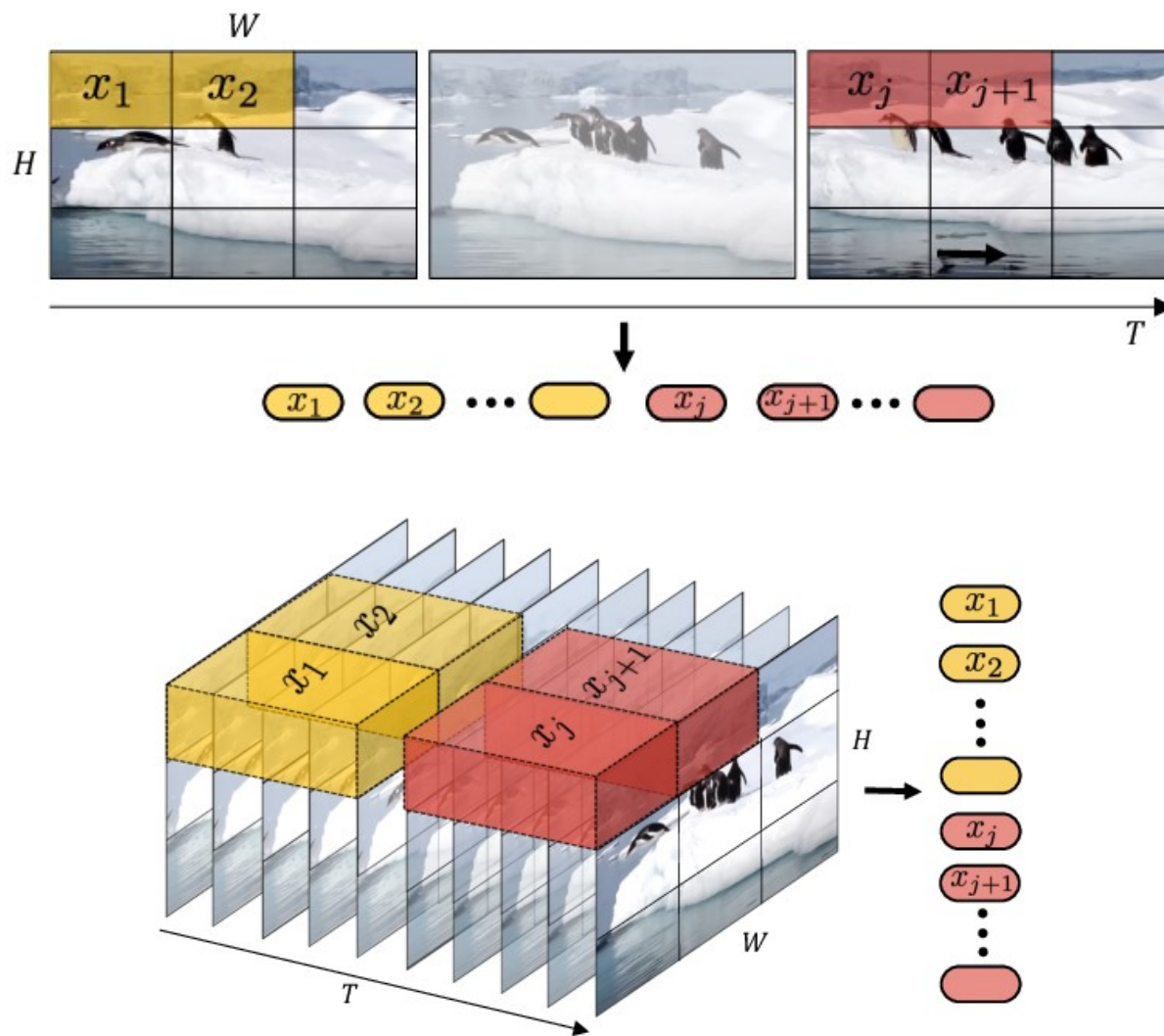
Vision Transformer



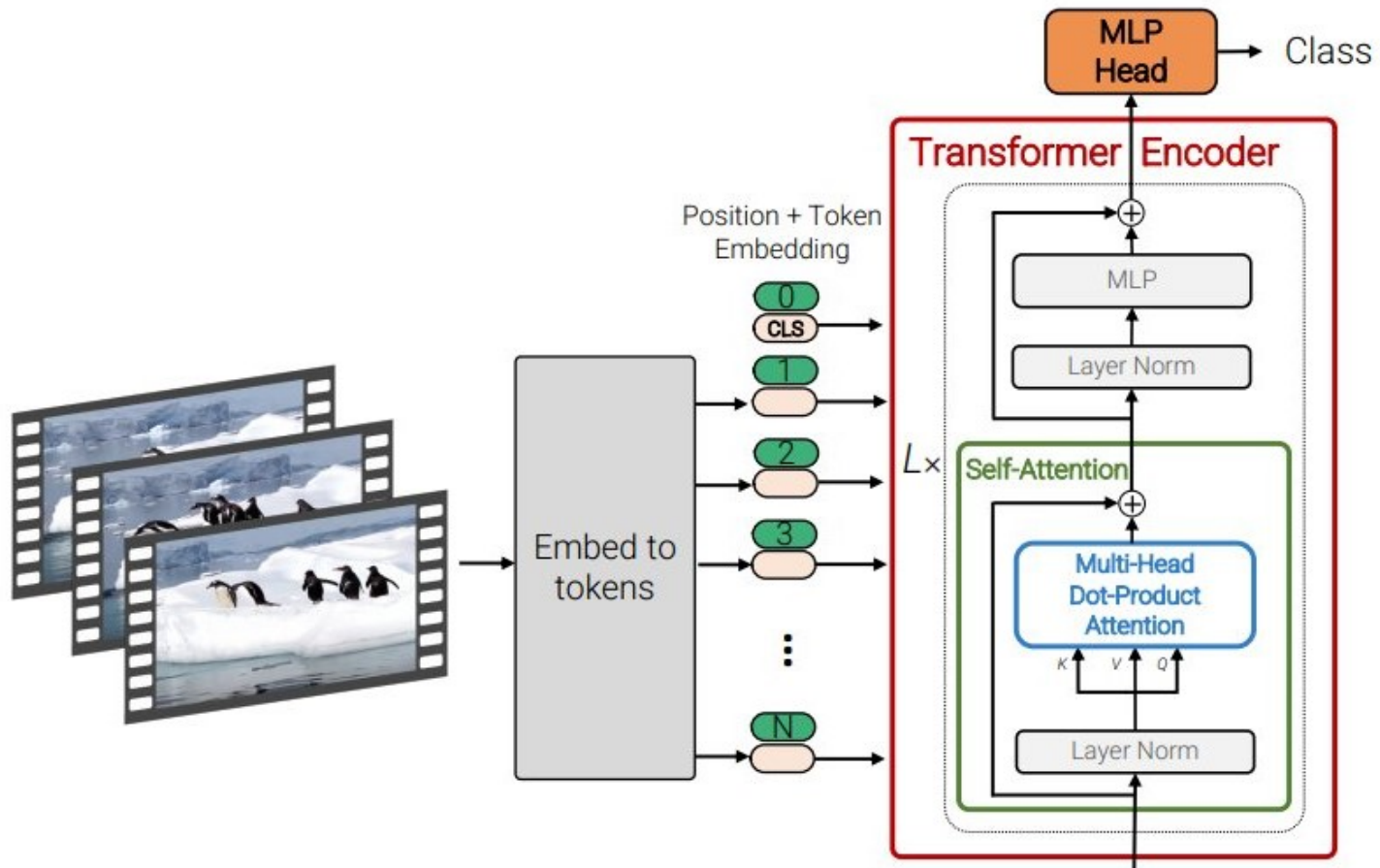
Pyramid Vision Transformer



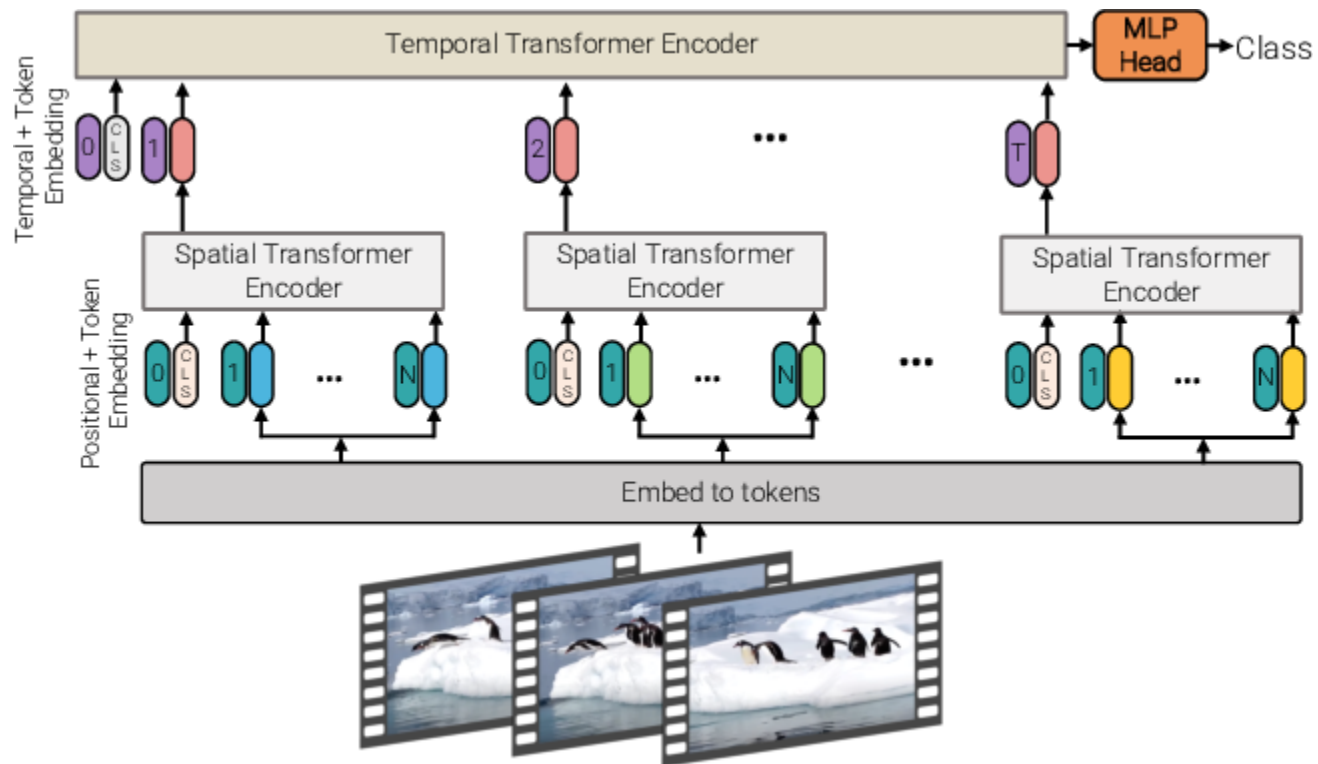
ViVit - patch



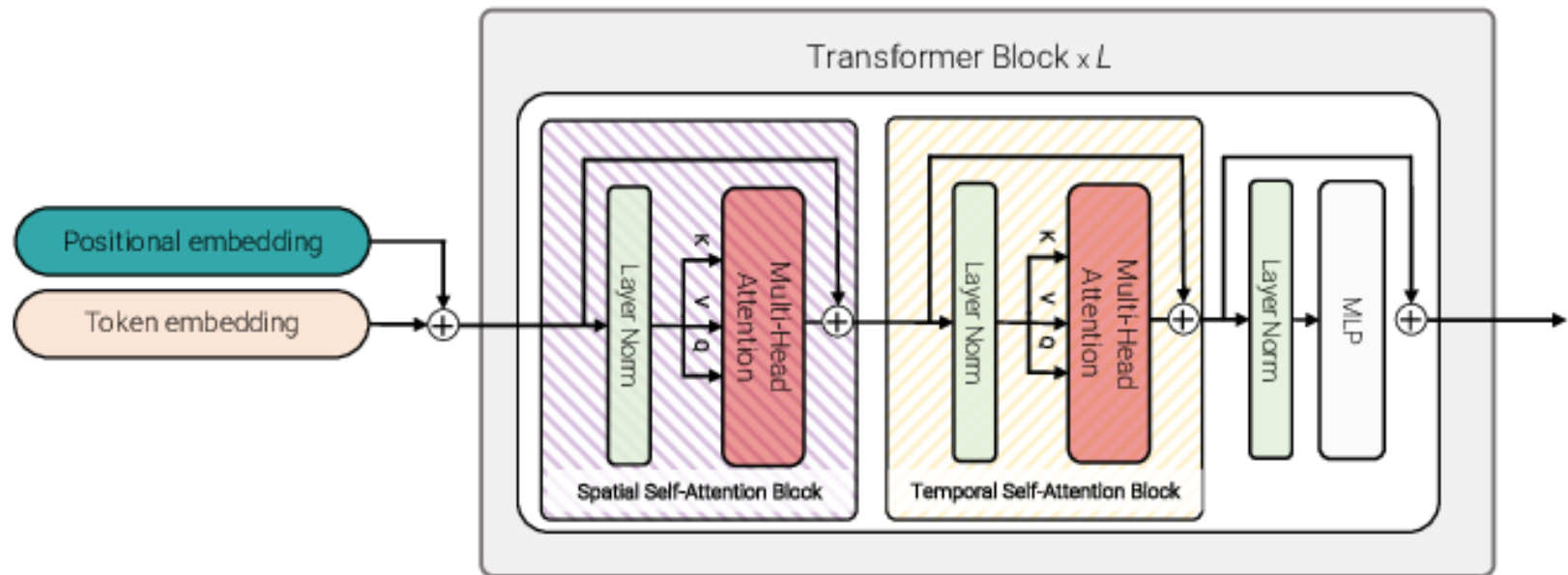
ViViT –M1



ViViT – M2



ViViT – M3



Travaux Pratiques

Practice

- Positional encoding:

https://colab.research.google.com/drive/1ibERwNZ_QcDXqb52Ac3f61v1wBfvqvpn?usp=sharing

- Self-attention:

<https://colab.research.google.com/drive/19M9W5fx6yx7LZ275ccNy7aQKcM3cVwSi?usp=sharing>

- Transformer for Text task:

<https://colab.research.google.com/drive/18LTQ5FgDJKSQiU1f0nAZ-K3eyKA5ZPH4?usp=sharing>

Thanks!

E-mail: ruい.dai@inria.fr