

COURS Reconnaissance Visuelle par deep learning

<https://cord.isir.upmc.fr/teaching-multimedia/>

Matthieu Cord  
Sorbonne University  
Computer Science - ISIR

# Outline

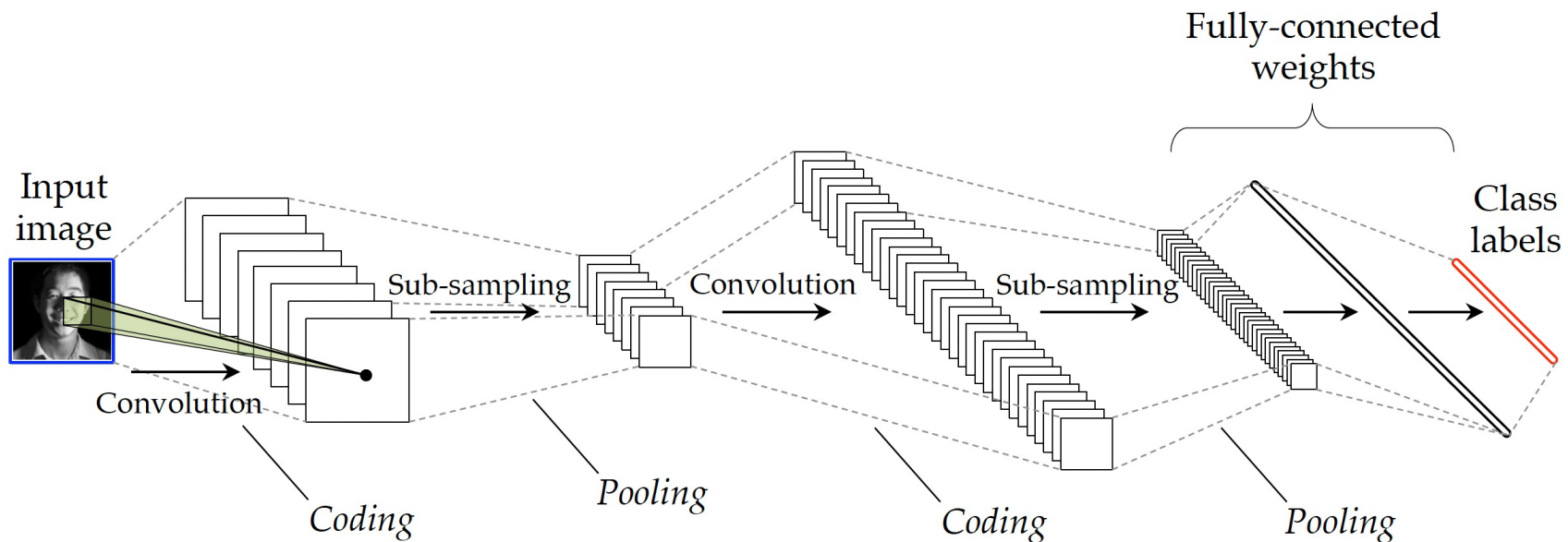
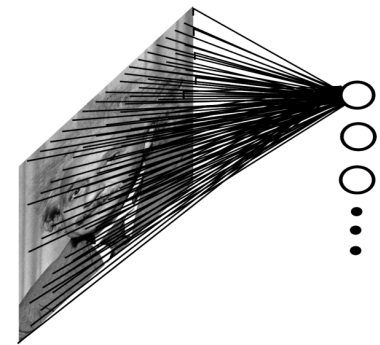
## 1. Attention and Vision Transformers

- NLP: Attention is all you need

# Attention process in ConvNets

In ConvNets, what information is shared between pixels (or features) in one block? => *2D spatial locality* (typically 3x3) => *attention is done locally*

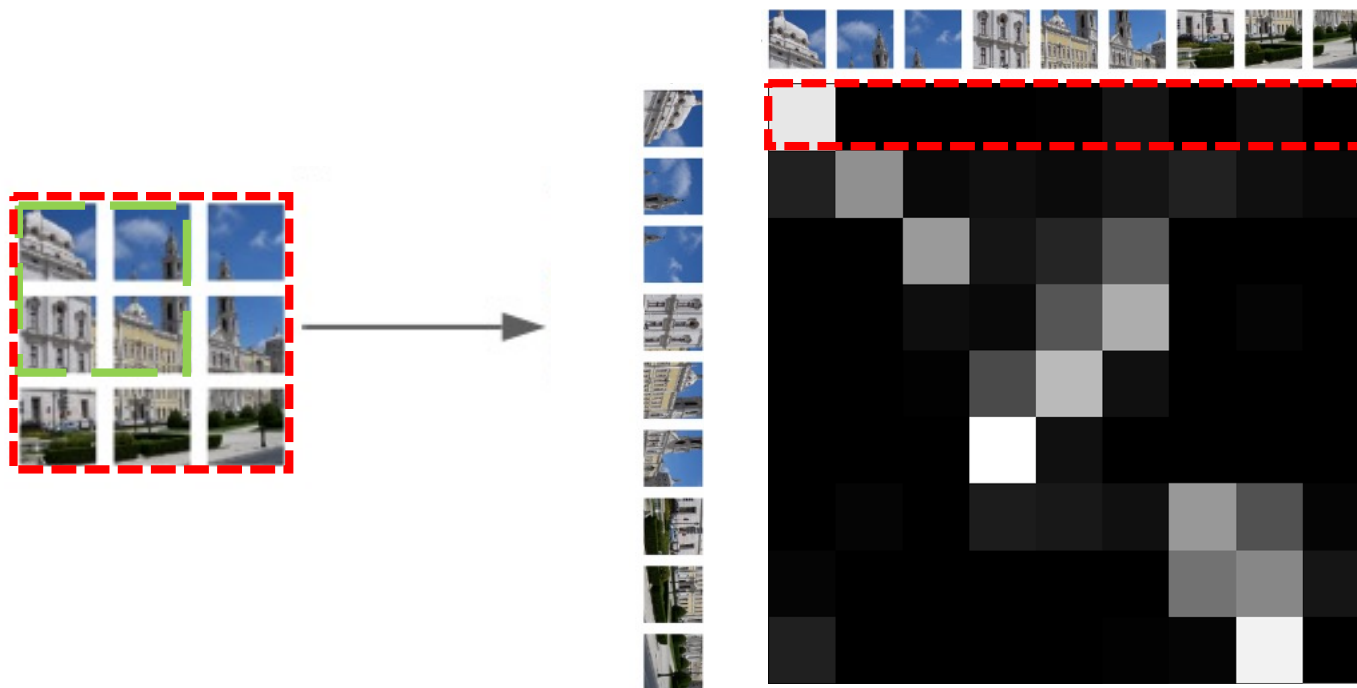
*Rq: less local after many layers*



# Global (Self) attention

How to build a deep architecture with ~~local~~ **global** attention inside?  
Meaning that one patch may interact with all others!

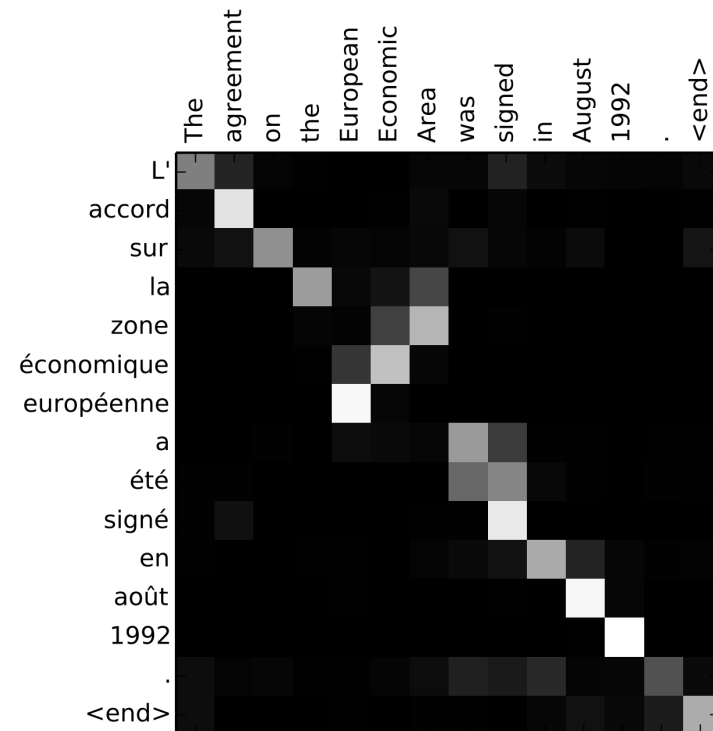
=> Different than convNet!



Let's see what they do in Natural Language Processing (NLP):

## Attention between words in Machine translation process:

1. Computing of weights
2. Use them to compute new features

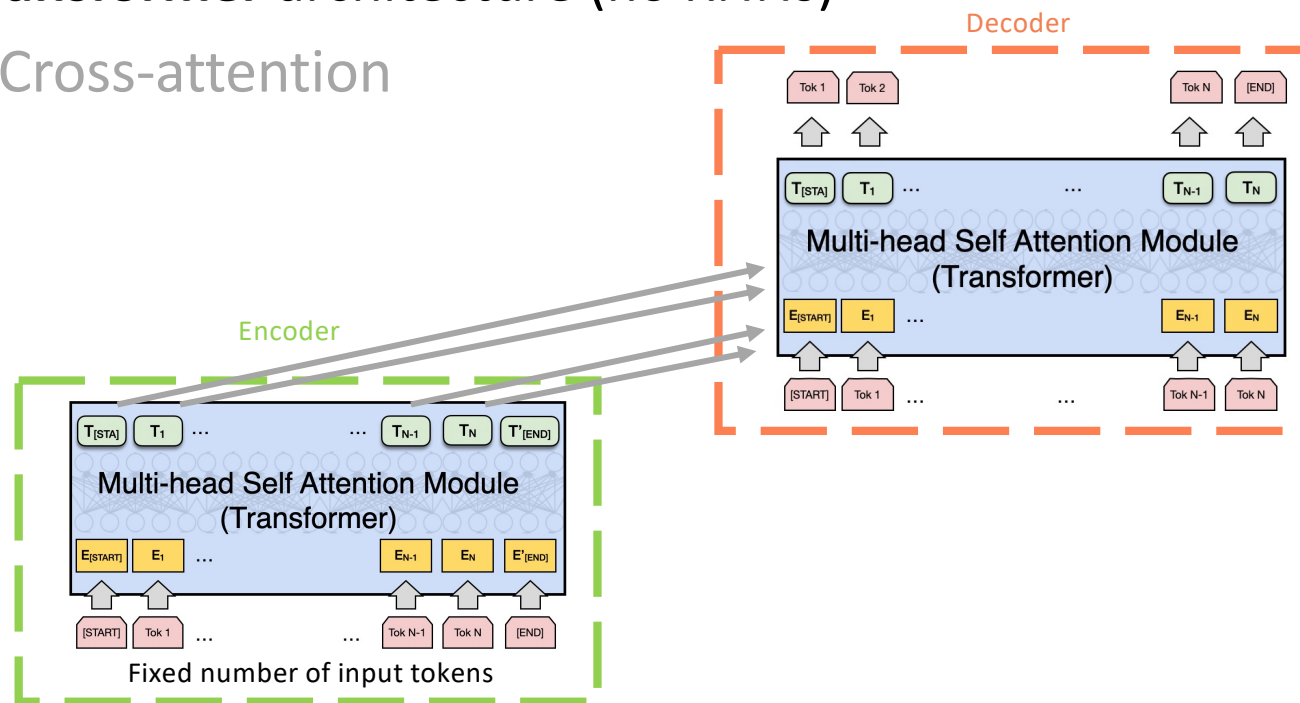


# Attention process in NLP

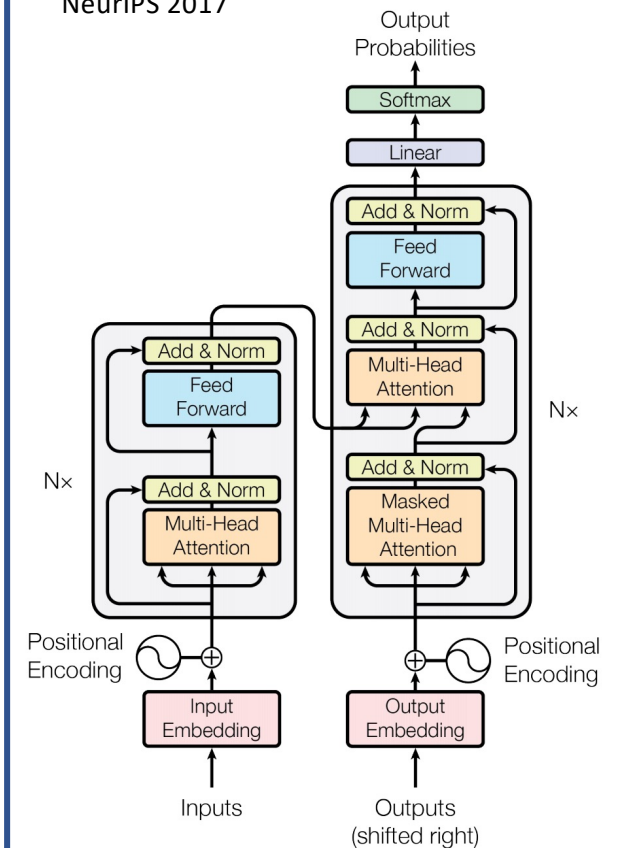
Basic language translation models: **Encoder/Decoder**

**Transformer** architecture (no RNNs)

- Cross-attention



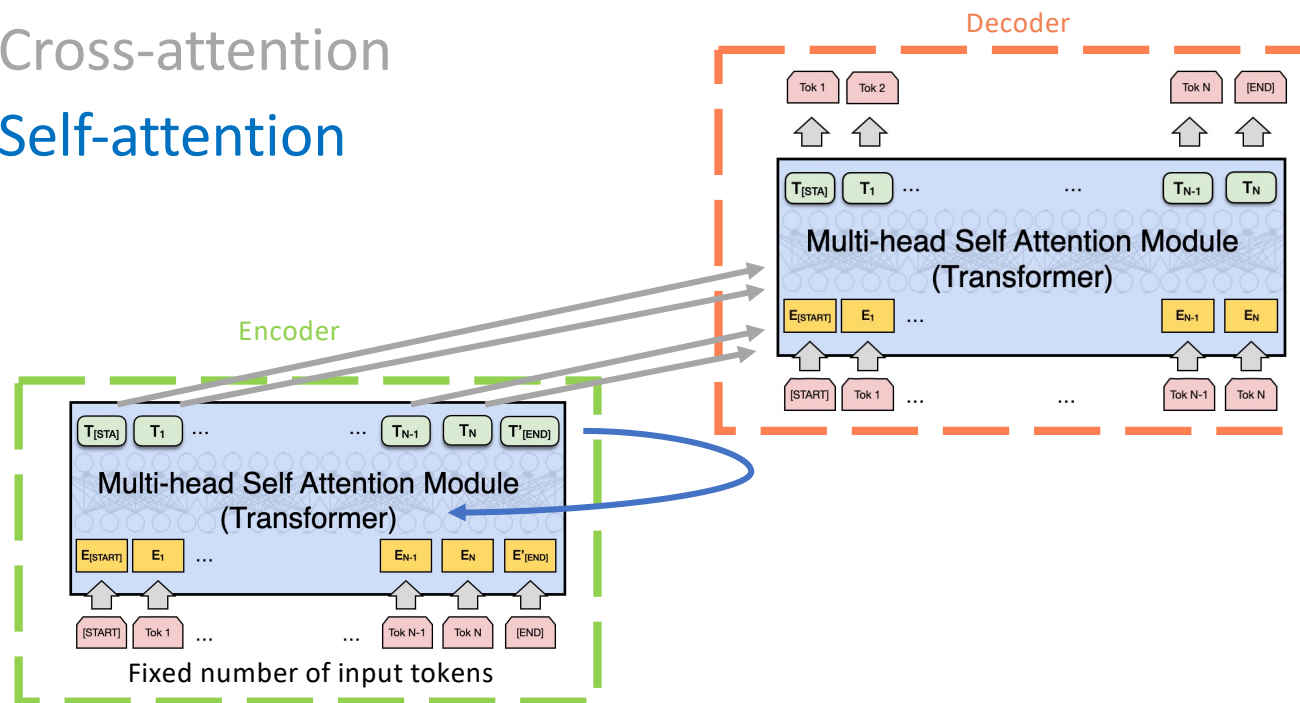
[Vaswani et al. Attention is all you need]  
<https://arxiv.org/abs/1706.03762>  
NeurIPS 2017



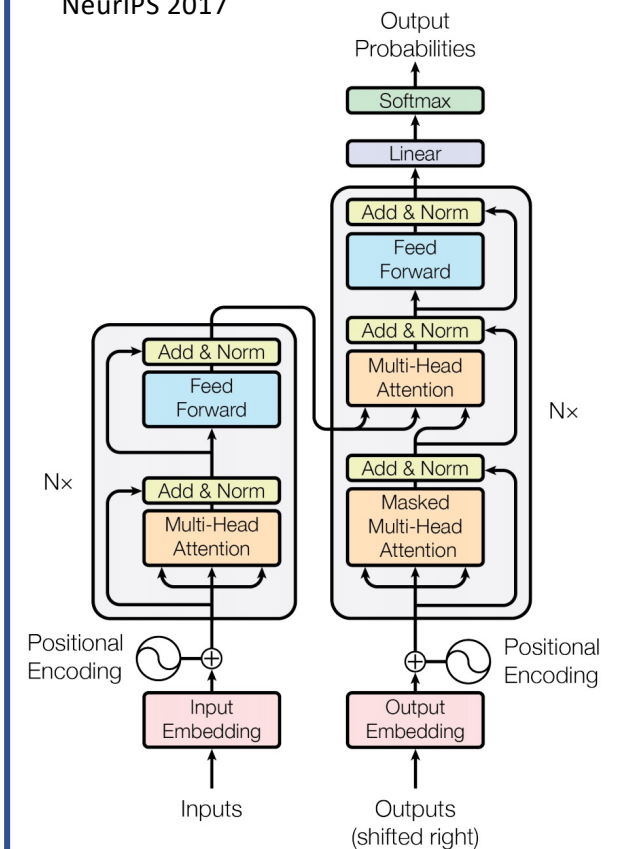
# Attention process in NLP

Basic language translation models: **Encoder/Decoder**  
**Transformer** architecture (no RNNs)

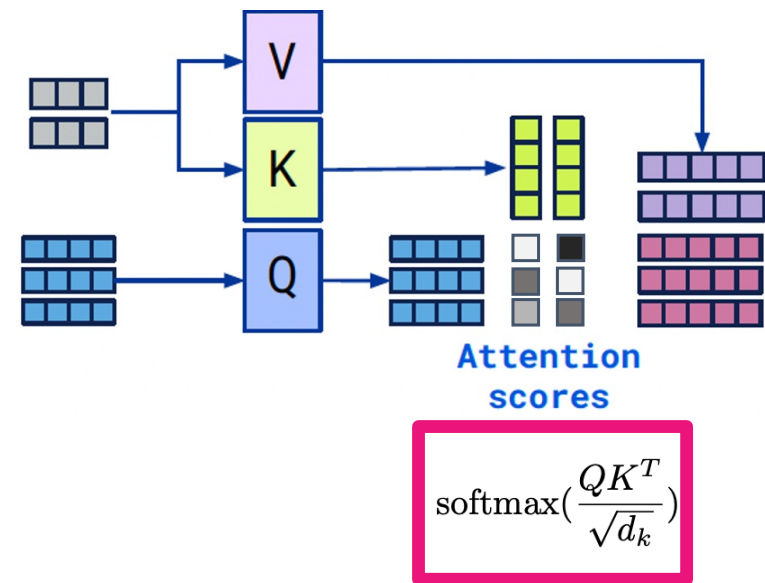
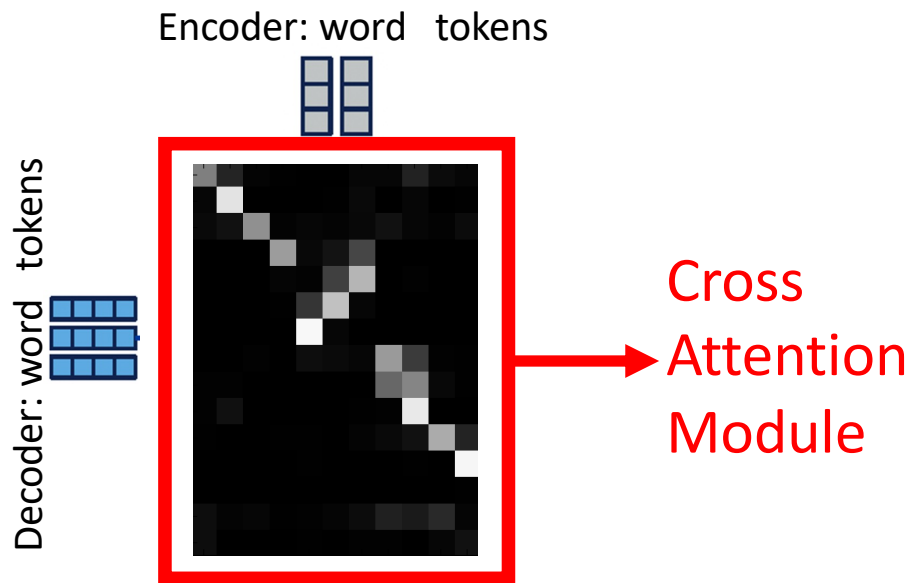
- Cross-attention
- Self-attention



[Vaswani et al. Attention is all you need]  
<https://arxiv.org/abs/1706.03762>  
NeurIPS 2017



# Attention process in NLP



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

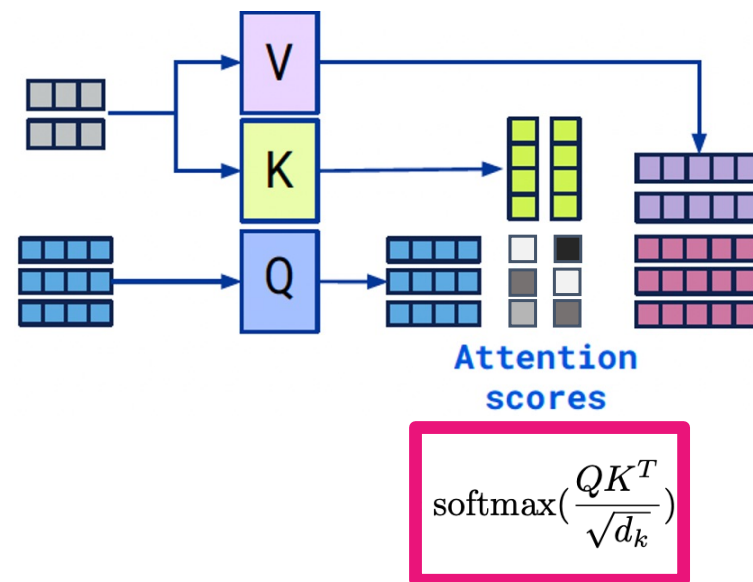
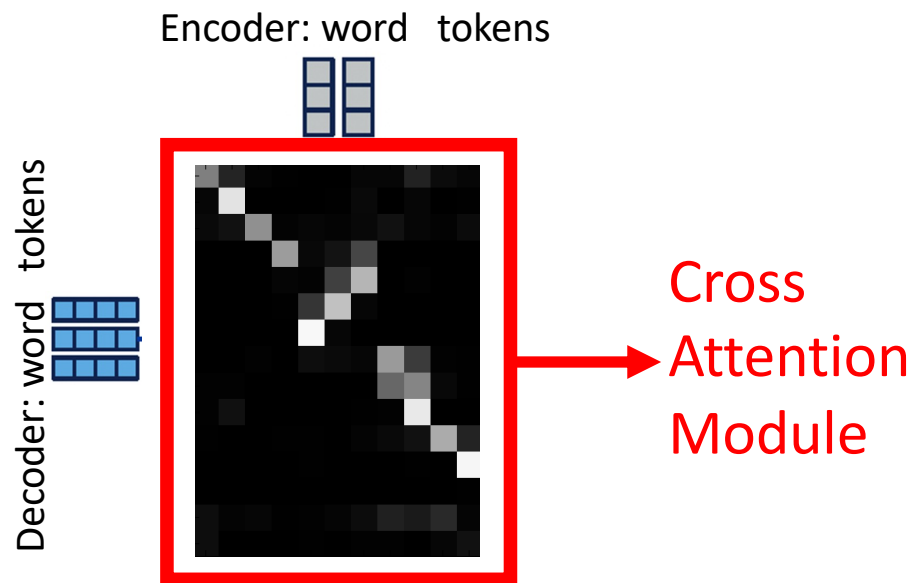


# Outline

## 1. Attention and Vision Transformers (ViT)

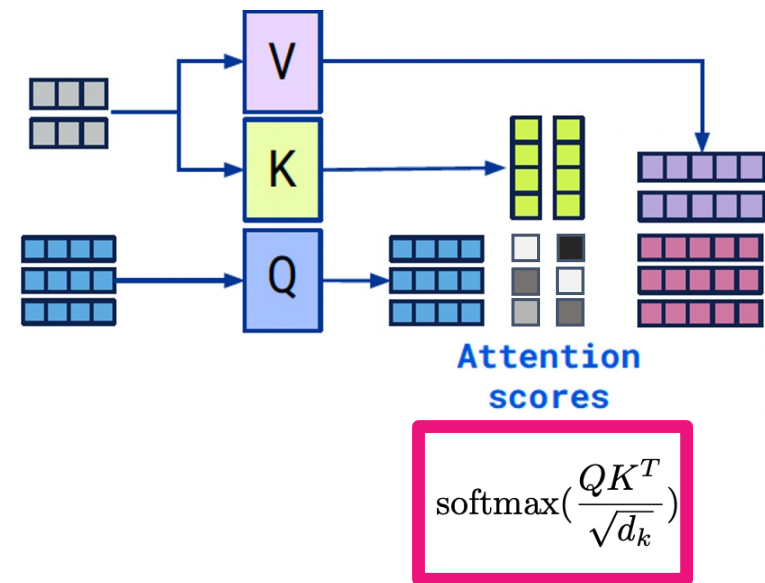
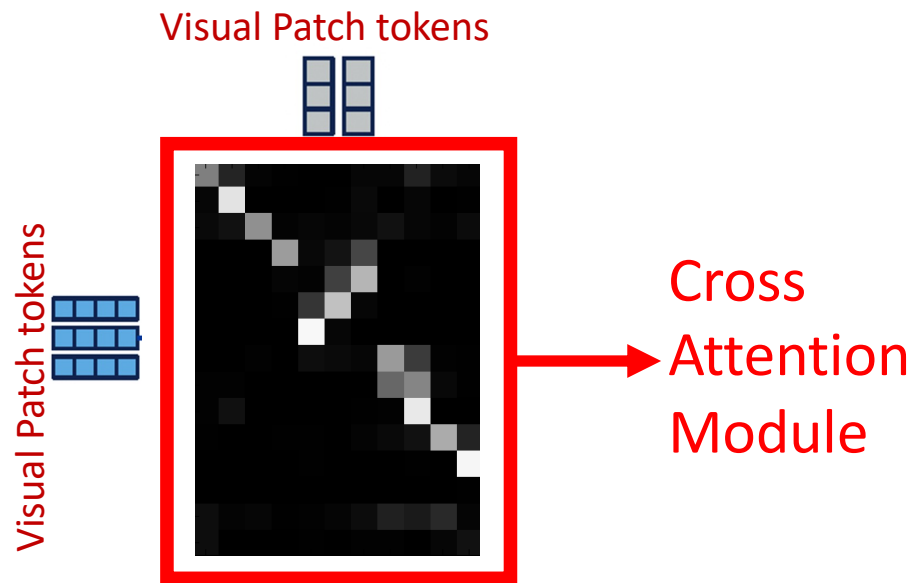
- NLP: Attention is all you need
- **Transformer for image classification**

# Attention process in NLP



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

# Attention process in Vision



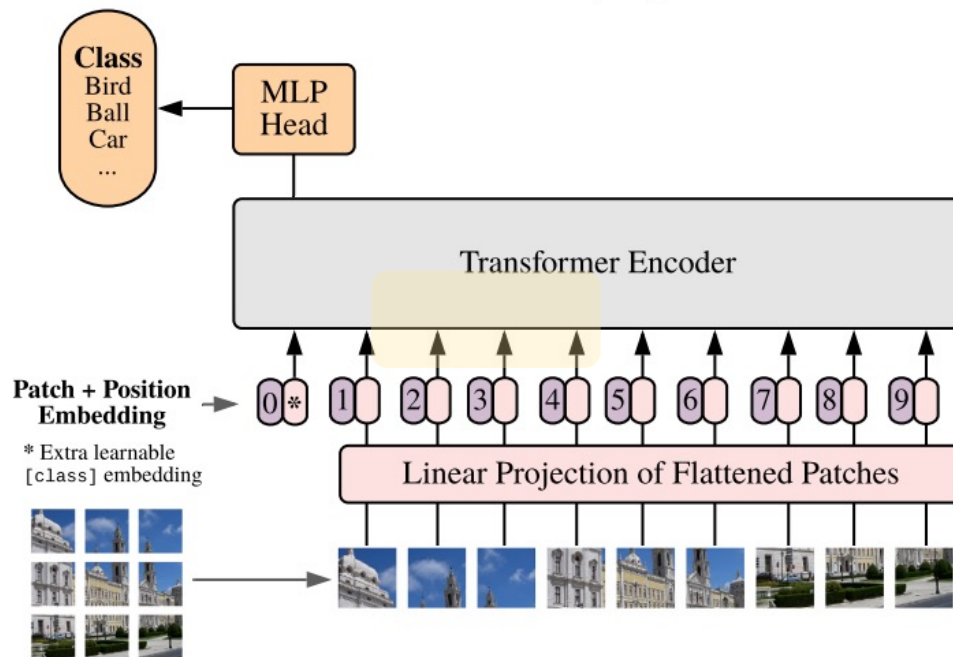
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Very similar except that Visual token is definitively less natural than word for NLP

# Attention process in Vision

Is it possible to mimic this attention-based architecture for vision processing?

Yes! **ViT** (Vision image Transformers) architecture

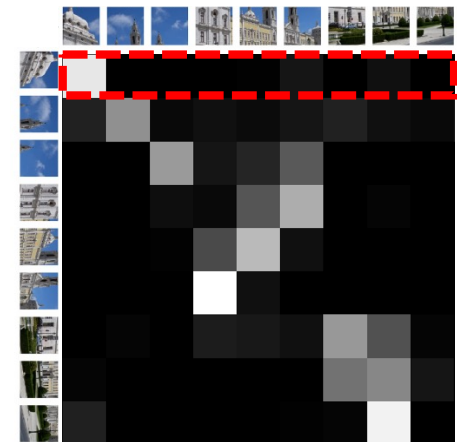


Published as a conference paper at ICLR 2021

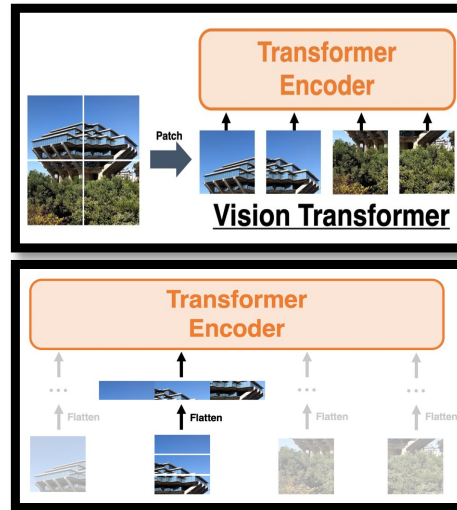
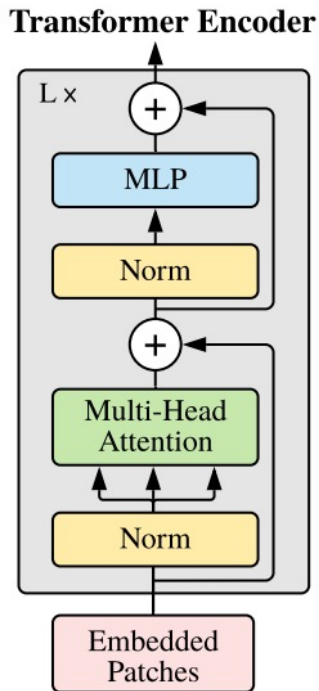
## AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy<sup>\*,†</sup>, Lucas Beyer<sup>\*</sup>, Alexander Kolesnikov<sup>\*</sup>, Dirk Weissenborn<sup>\*</sup>,  
Xiaohua Zhai<sup>\*</sup>, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby<sup>\*,†</sup>

<sup>\*</sup>equal technical contribution, <sup>†</sup>equal advising  
Google Research, Brain Team  
{adosovitskiy, neilhoulby}@google.com



# Attention process in Vision



$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}},$$

$$\mathbf{z}'_\ell = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1},$$

$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell,$$

$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0)$$

$$\mathbf{x} \in \mathbb{R}^{H \times W \times C}$$

$$\mathbf{x}_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$$

$$N = HW/P^2$$

$$\mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D}$$

CLS token

$$\ell = 1 \dots L$$

$$\ell = 1 \dots L$$

[class=CLS] token: a learnable embedding to the sequence of embedded patches

Layernorm (LN) before every block, and residual connections after every block

MSA: Multi Head Self Attention

MLP: two layers with a GELU non-linearity

Hybrid Architecture : Raw image patches --> Feature map of a CNN

