

COURS RDFIA deep Image

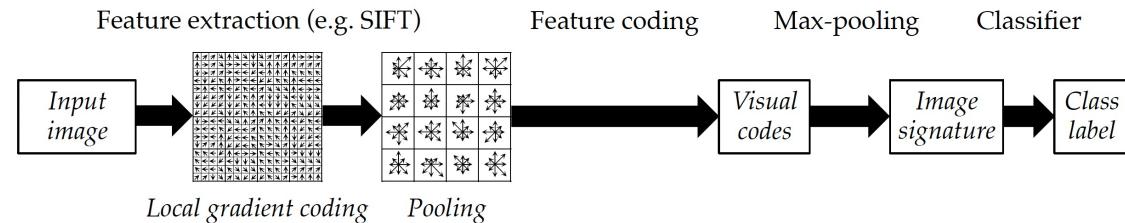
Matthieu Cord
Sorbonne University

Course Outline

- 1. Computer Vision and ML basics:** Visual (local) feature detection and description, Bag of Word Image representation, Linear classification (SVM)
- 2. Introduction to Neural Networks (NNs)**
- 3. Machine Learning basics (2):** Risk, Classification, Datasets, benchmarks and evaluation
- 4. Neural Nets for Image Classification**
- 5. Vision Transformers**

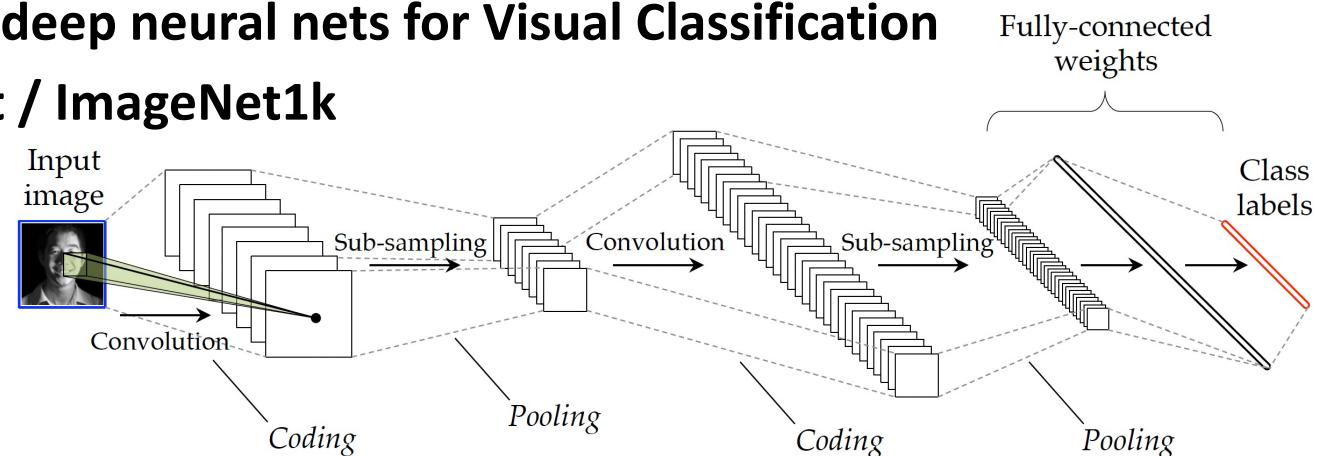
Context: Image classification Before/After ImageNet (2009)

The 2000s: *BoWs image modeling + SVMs* for Visual Classification



The 2010s: *Large deep neural nets for Visual Classification*

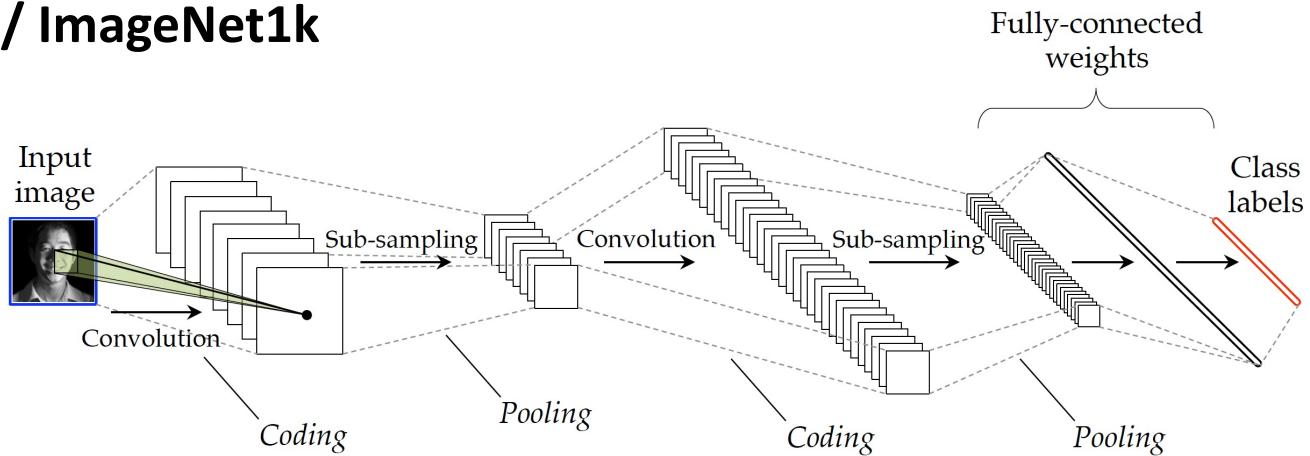
The star: **ConvNet / ImageNet1k**



Context: Image classification After ImageNet (2009)

The 2010s: *Large deep neural nets for Visual Classification*

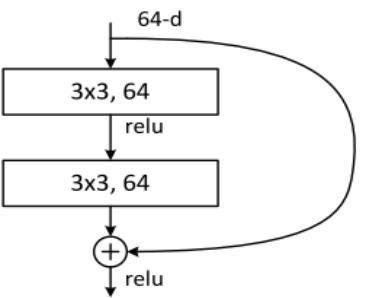
The star: **ConvNet / ImageNet1k**



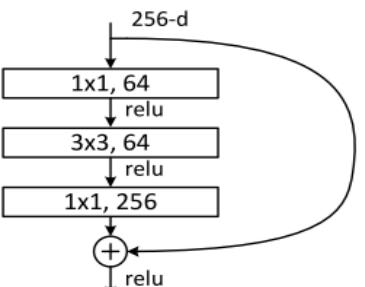
AlexNet 2012

- Same model as LeCun'98 but:
 - Bigger model (8 layers)
 - More data (10^6 vs 10^3 images)
 - GPU implementation (50x speedup over CPU)
 - Better regularization (DropOut)

Post-2012 revolution: ResNet Architecture

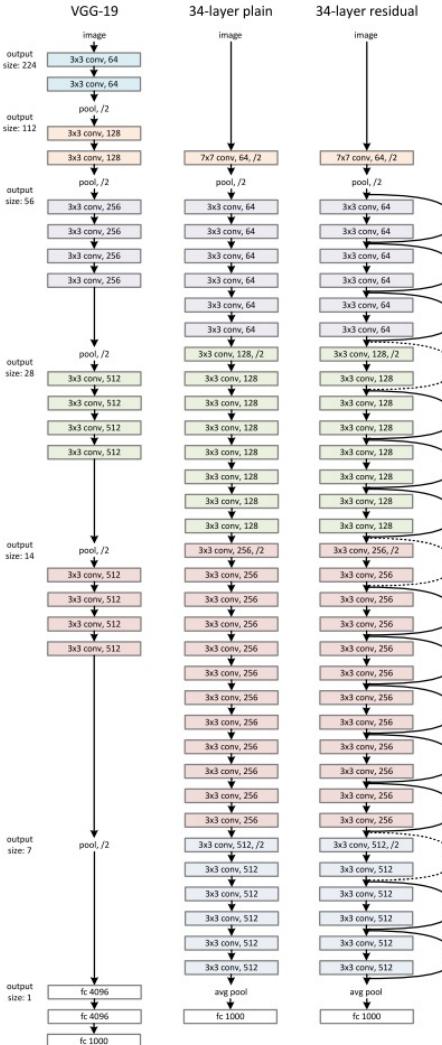


A naïve residual block



“bottleneck” residual block
(for ResNet-50/101/152)

ConvNet Configuration			
B	C	D	E
13 weight layers	16 weight layers	16 weight layers	19 weight layers
Input (224 × 224 RGB image)			
conv3-64	conv3-64	conv3-64	conv3-64
conv3-64	conv3-64	conv3-64	conv3-64
maxpool			
conv3-128	conv3-128	conv3-128	conv3-128
conv3-128	conv3-128	conv3-128	conv3-128
maxpool			
conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256	conv3-256
conv1-256	conv3-256	conv3-256	conv3-256
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
conv1-512	conv3-512	conv3-512	conv3-512
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
conv1-512	conv3-512	conv3-512	conv3-512
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			



Context: Beyond ImageNet?

The 2000s: *BoWs image modeling + SVMs* for Visual Classification

The 2010s: *Large deep neural nets* for Visual Classification

What is expected for the 2020s?

“*Attention is all you need*”: **Transformers** for Vision !?

And **datasets? Internet...**

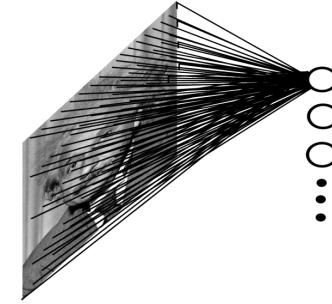
[Vaswani et al., Attention is all you need, NeurIPS 2017]

Outline

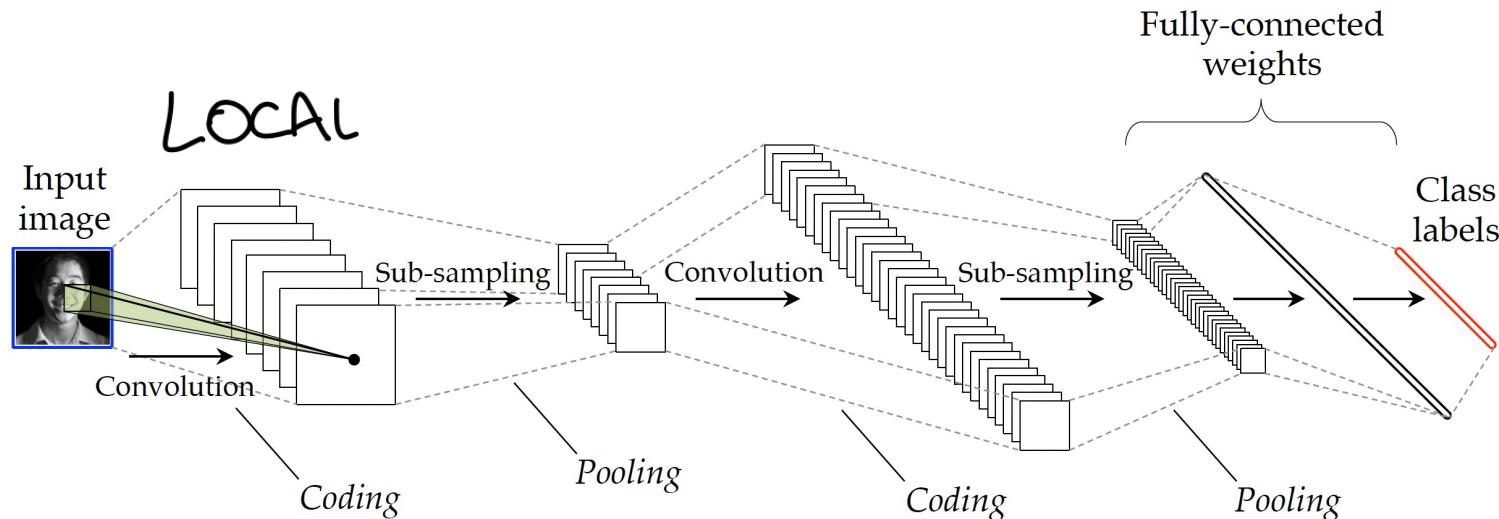
1. Computer Vision and ML basics: Visual (local) feature detection and description, Bag of Word Image representation, Linear classification (SVM)
2. Introduction to Neural Networks (NNs)
3. Machine Learning basics (2): Risk, Classification, Datasets, benchmarks and evaluation
4. Neural Nets for Image Classification
5. 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*



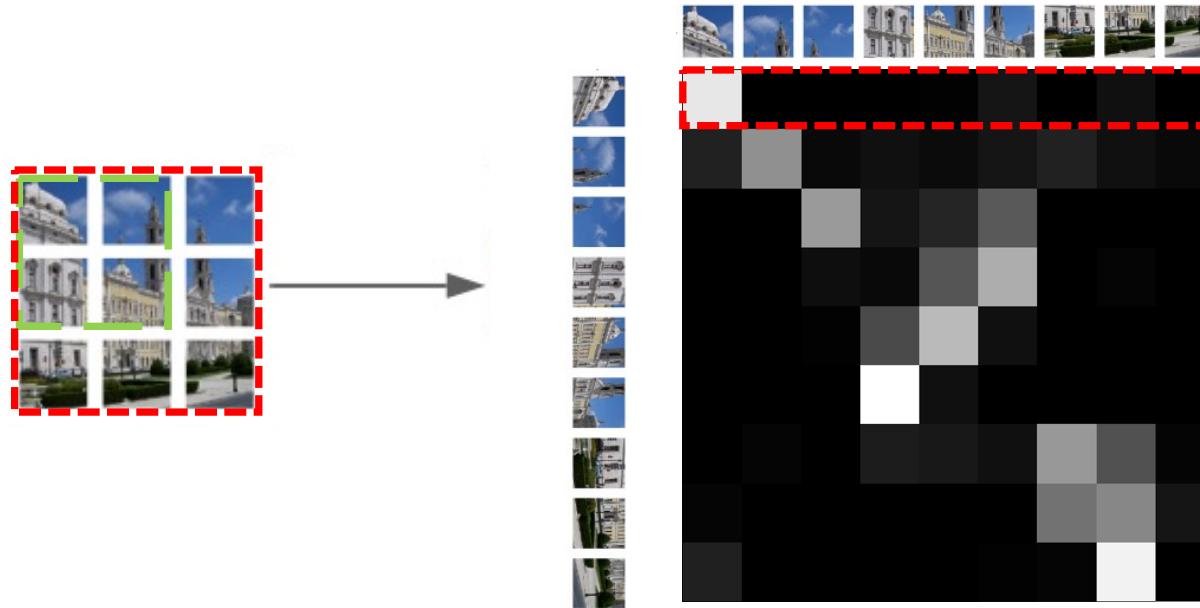
Rq: less local after many layers



Global (Self) attention

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

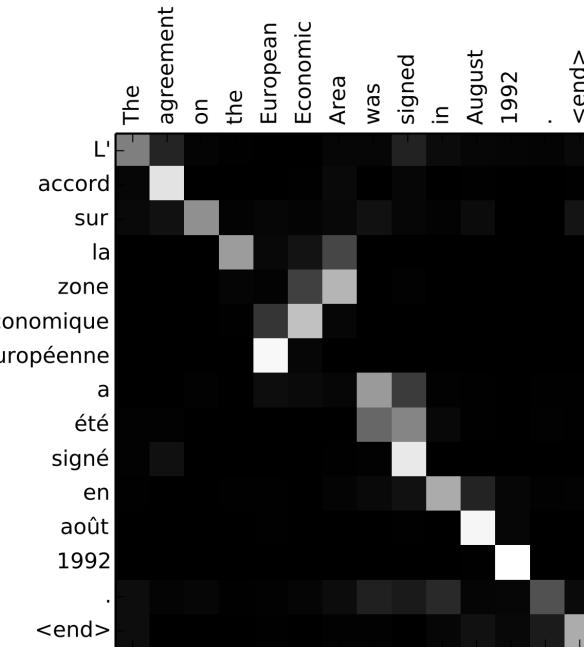
=> Different than convNet!



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

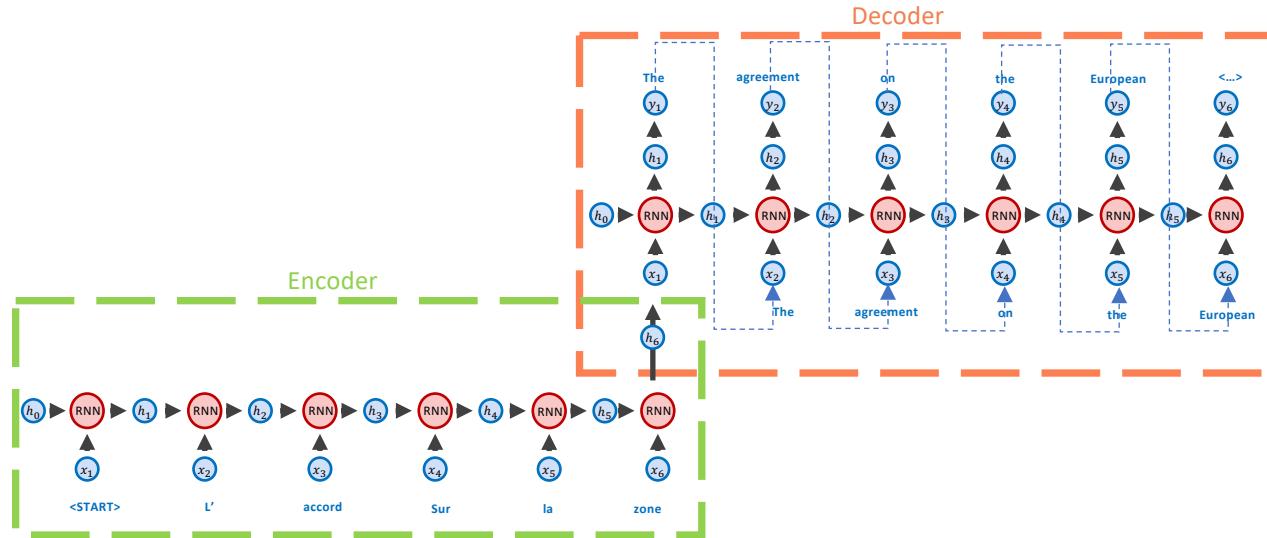


Attention process in NLP

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

Ex.: Seq2Seq -- RNNs2RNNs

Cross-attention for language translation in at the end of Encoder

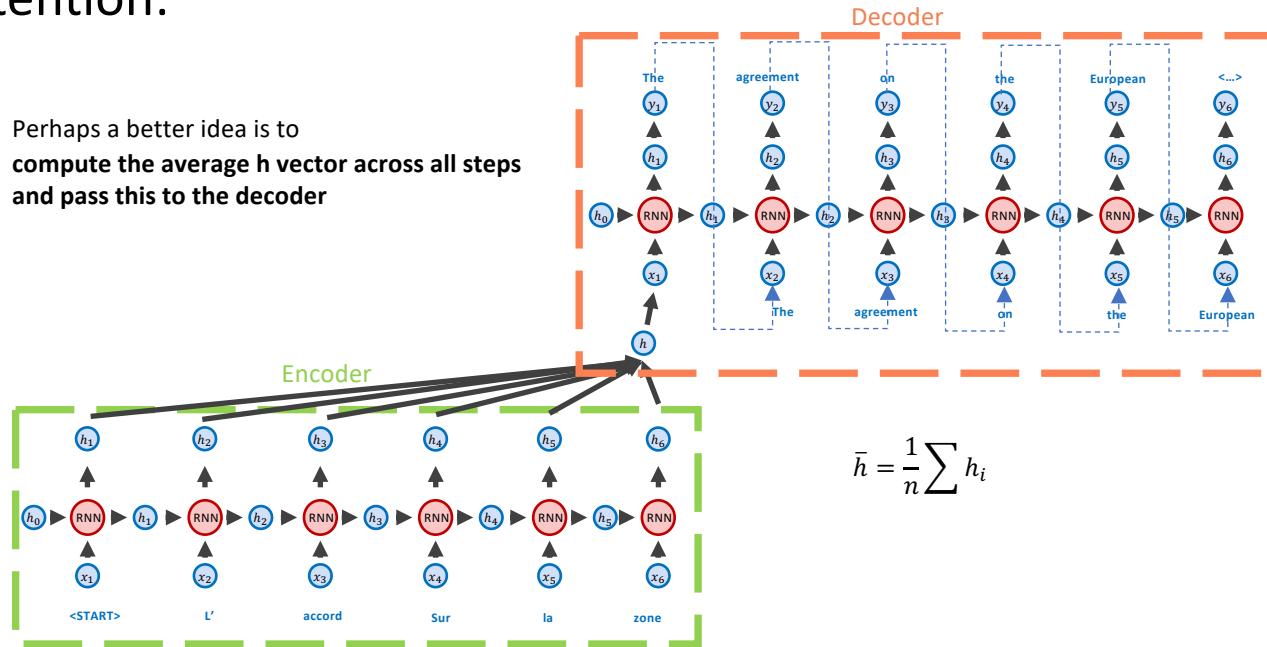


Attention process in NLP

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

Cross-attention:

Perhaps a better idea is to
compute the average h vector across all steps
and pass this to the decoder

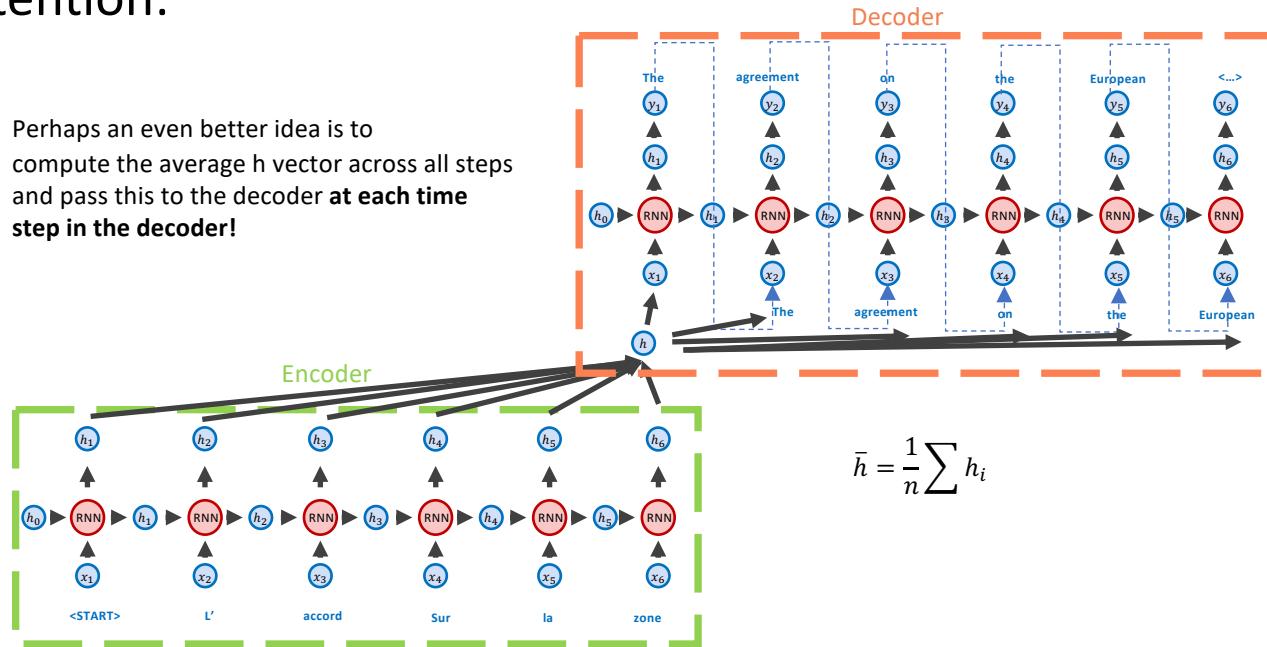


Attention process in NLP

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

Cross-attention:

Perhaps an even better idea is to compute the average h vector across all steps and pass this to the decoder at each time step in the decoder!

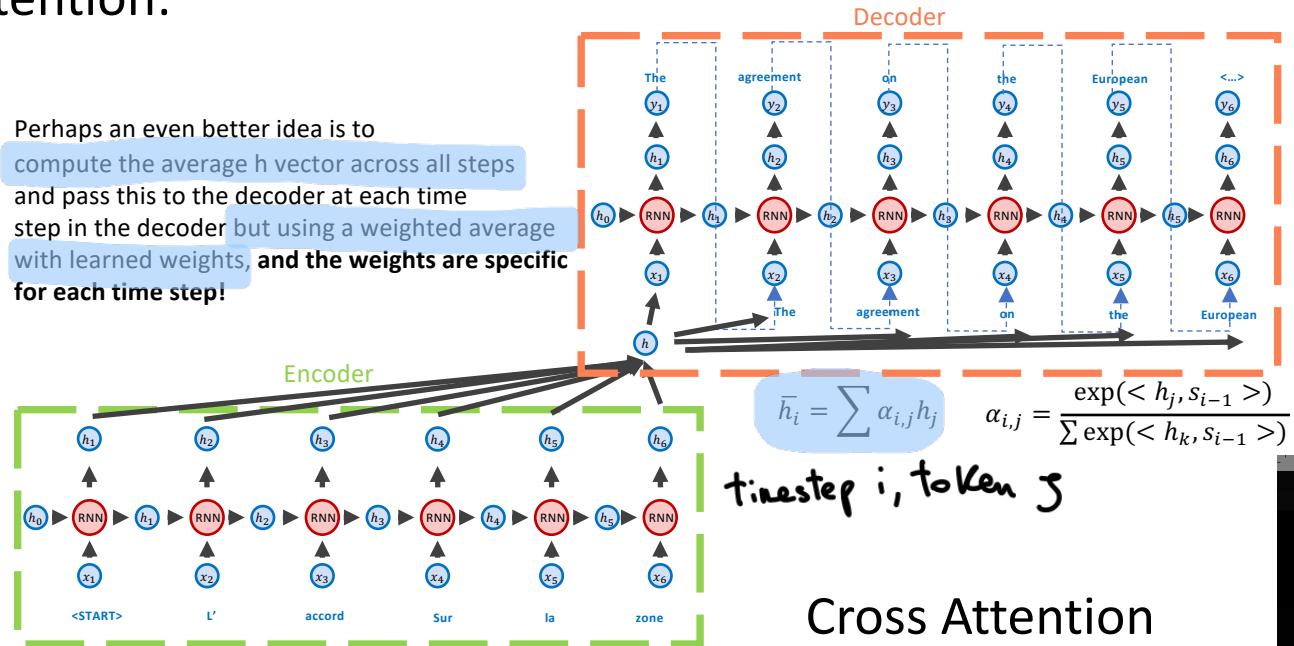


Attention process in NLP

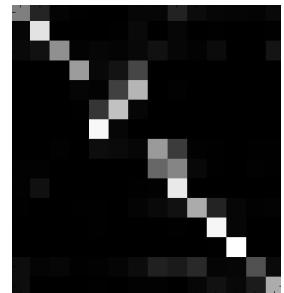
Basic language translation models: **Encoder/Decoder**

Cross-attention:

Perhaps an even better idea is to compute the average h vector across all steps and pass this to the decoder at each time step in the decoder but using a weighted average with learned weights, and the weights are specific for each time step!



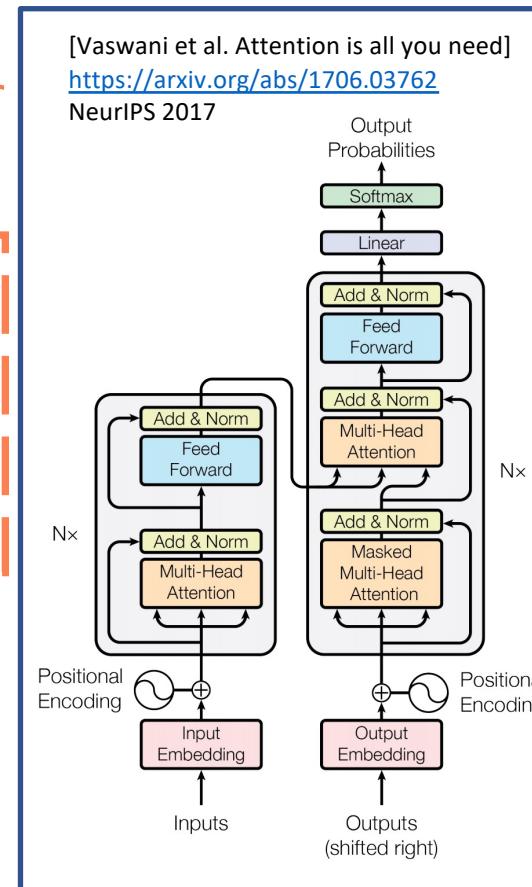
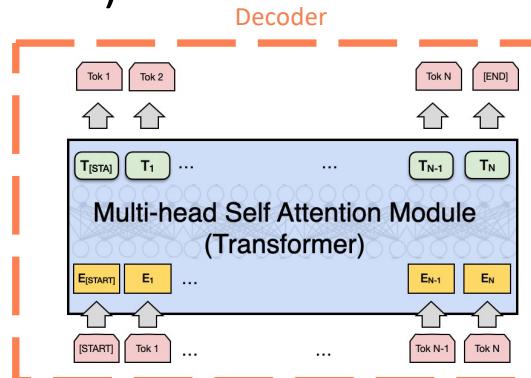
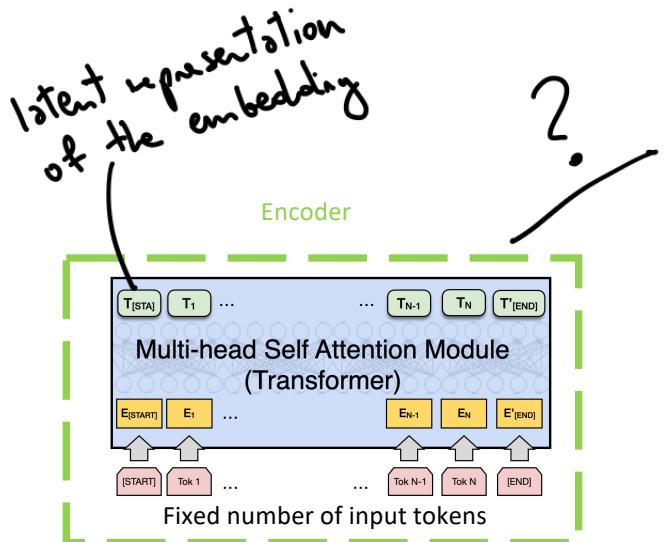
Cross Attention
Encoder/ Decoder



Attention process in NLP

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

Transformer architecture (no RNNs)

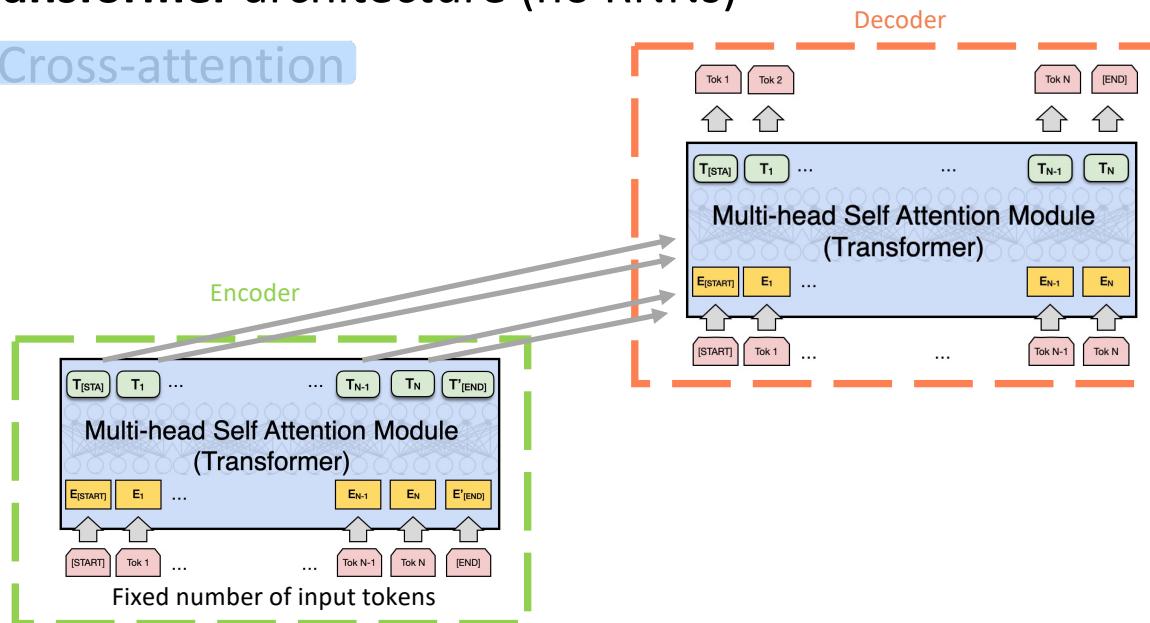


Attention process in NLP

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

Transformer architecture (no RNNs)

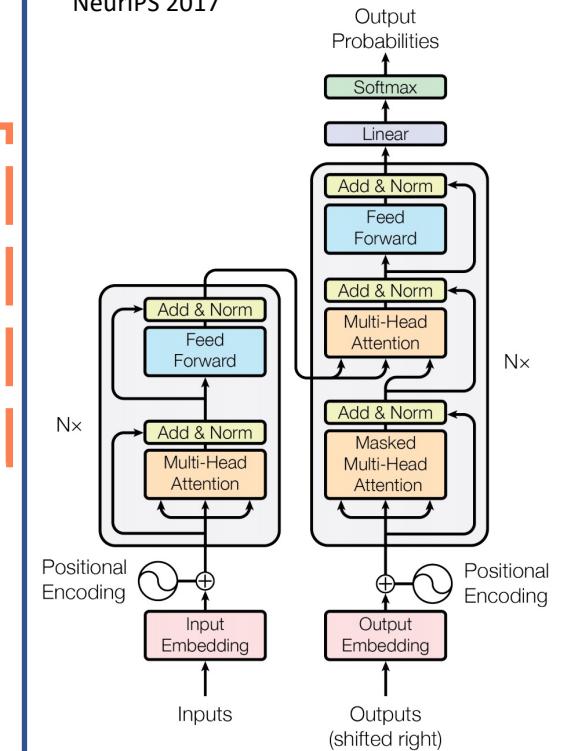
- **Cross-attention**



[Vaswani et al. Attention is all you need]

<https://arxiv.org/abs/1706.03762>

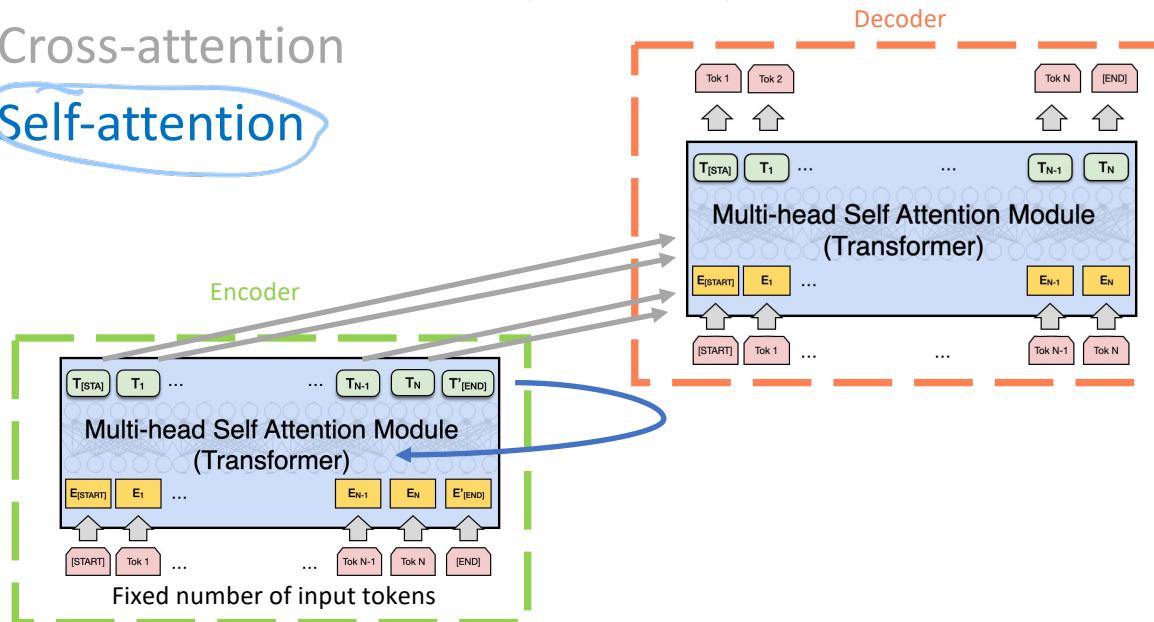
NeurIPS 2017



Attention process in NLP

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

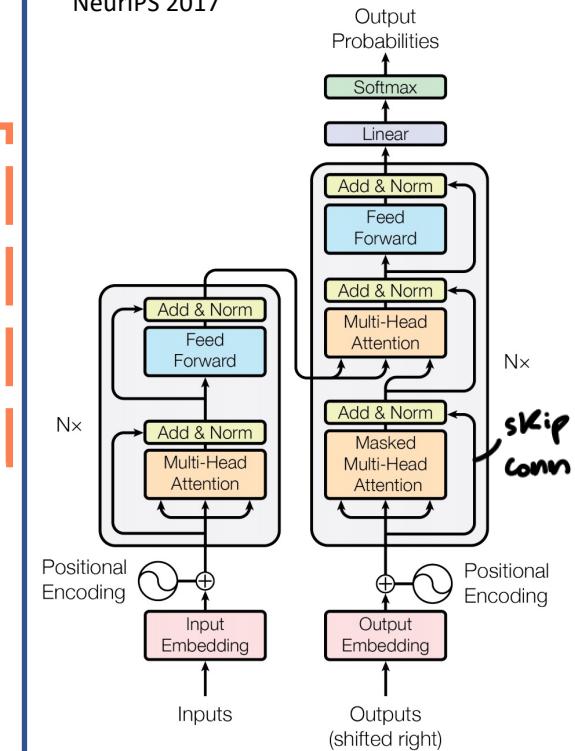
- Cross-attention
- **Self-attention**



[Vaswani et al. Attention is all you need]

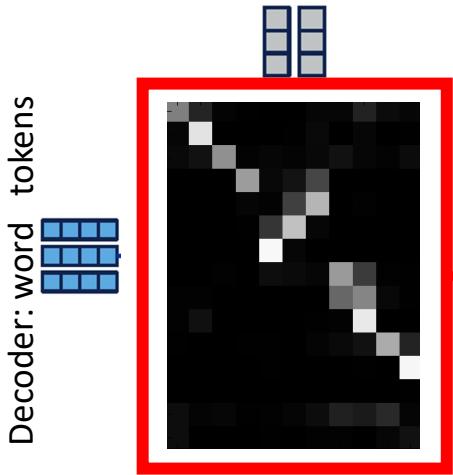
<https://arxiv.org/abs/1706.03762>

NeurIPS 2017



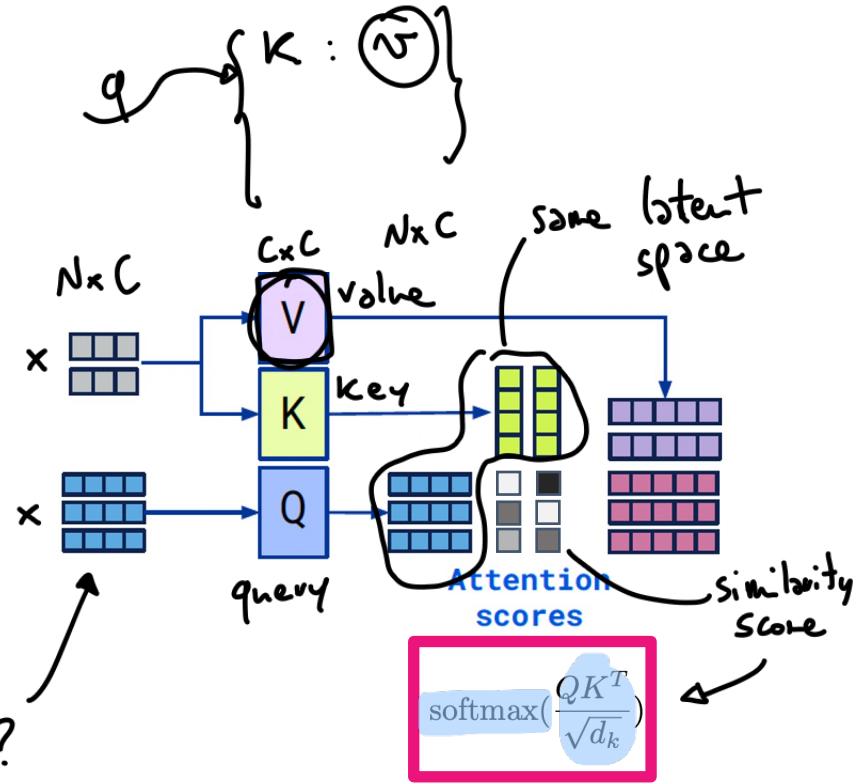
Attention process in NLP

Encoder: word tokens



Cross
Attention
Module

What are the
best tokens coming
from the decoder?



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

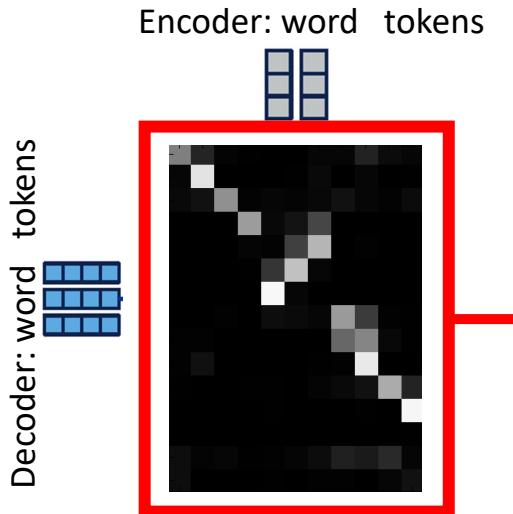
Outline

1. Computer Vision and ML basics: Visual (local) feature detection and description, Bag of Word Image representation, Linear classification (SVM)
 2. Introduction to Neural Networks (NNs)
 3. Machine Learning basics (2): Risk, Classification, Datasets, benchmarks and evaluation
 4. Neural Nets for Image Classification
- ## 5. Vision Transformers

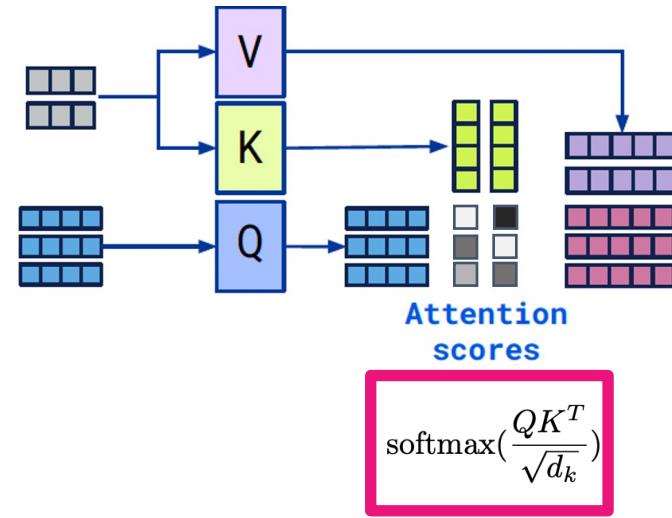
NLP: Attention is all you need

Transformer for image classification

Attention process in NLP

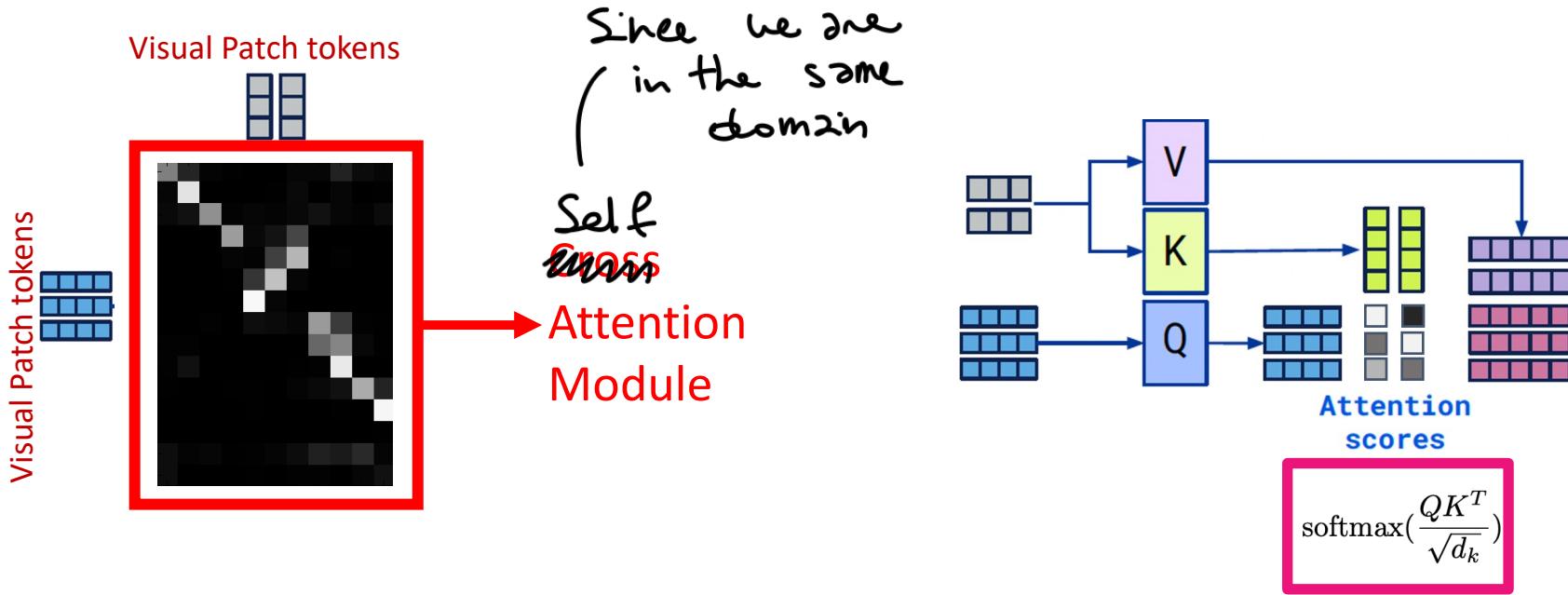


Cross
Attention
Module



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

Attention process in Vision



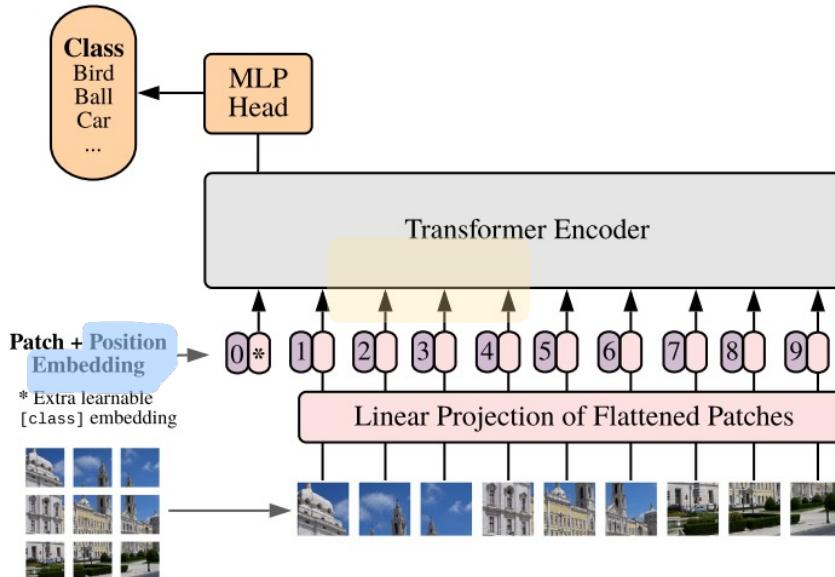
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Very similar except that Visual token is definitely 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

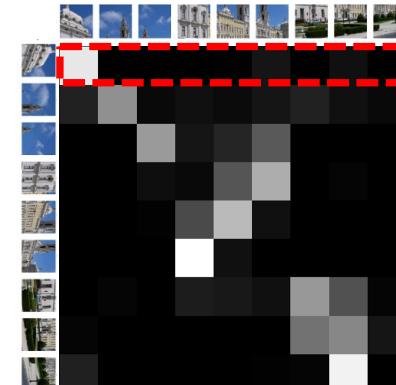
AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

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

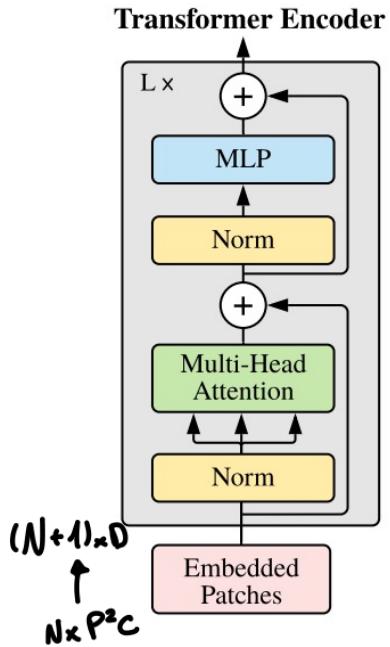
*equal technical contribution, †equal advising

Google Research, Brain Team

{adosovitskiy, neilhoulsby}@google.com



Attention process in Vision



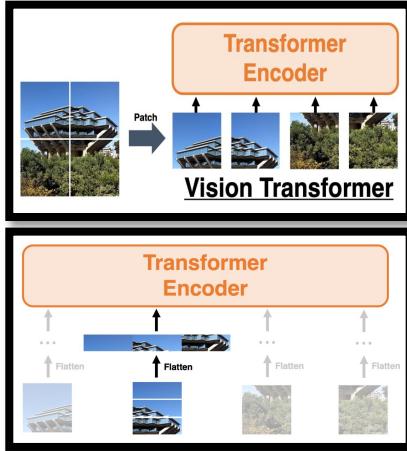
[class=CLS] token: a learnable embedding to the sequence of embedded patches

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

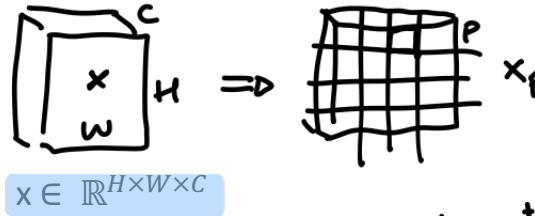
MSA: Multi Head Self Attention

MLP: two layers with a GELU non-linearity

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



each patch is
 $P \times P \times C$



single patch

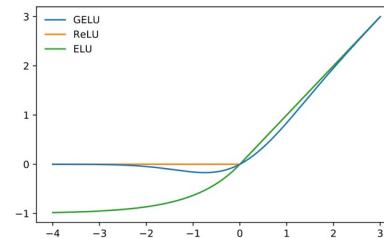


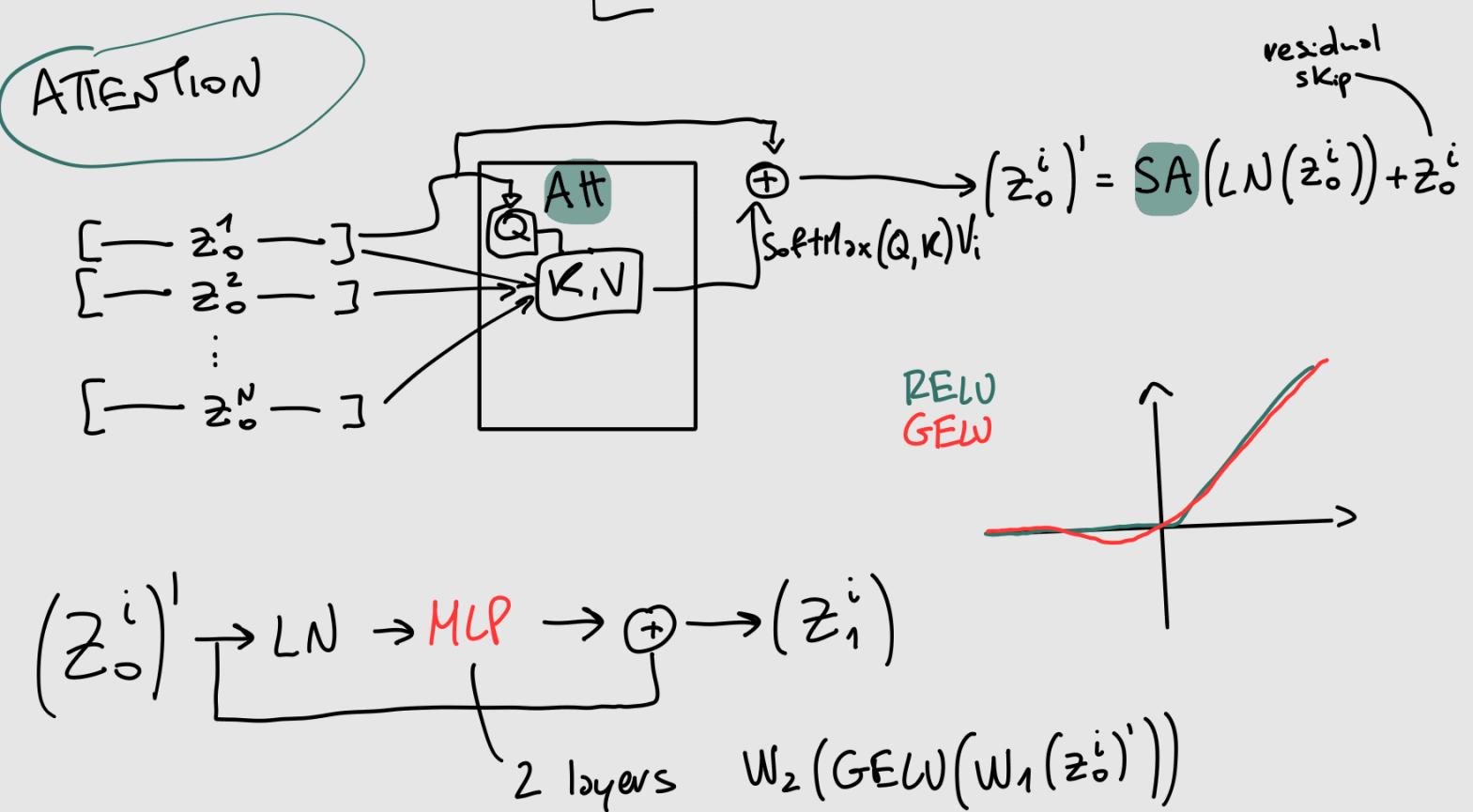
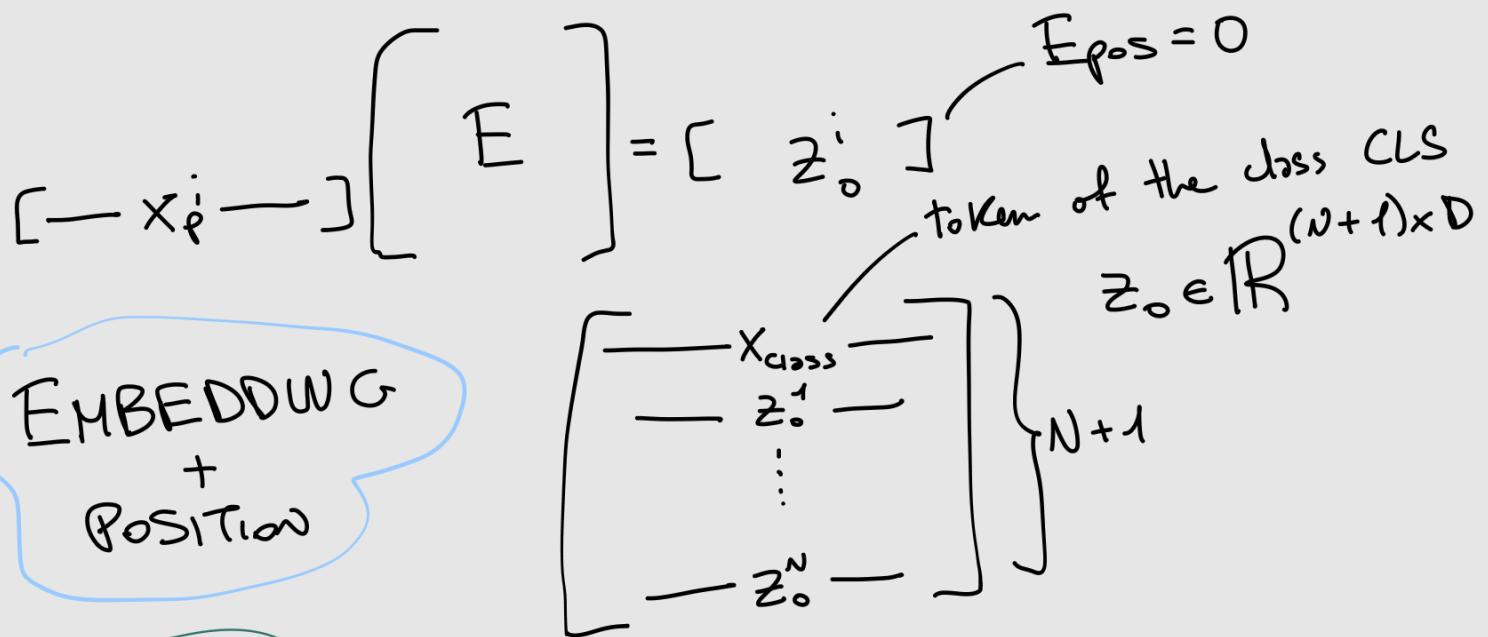
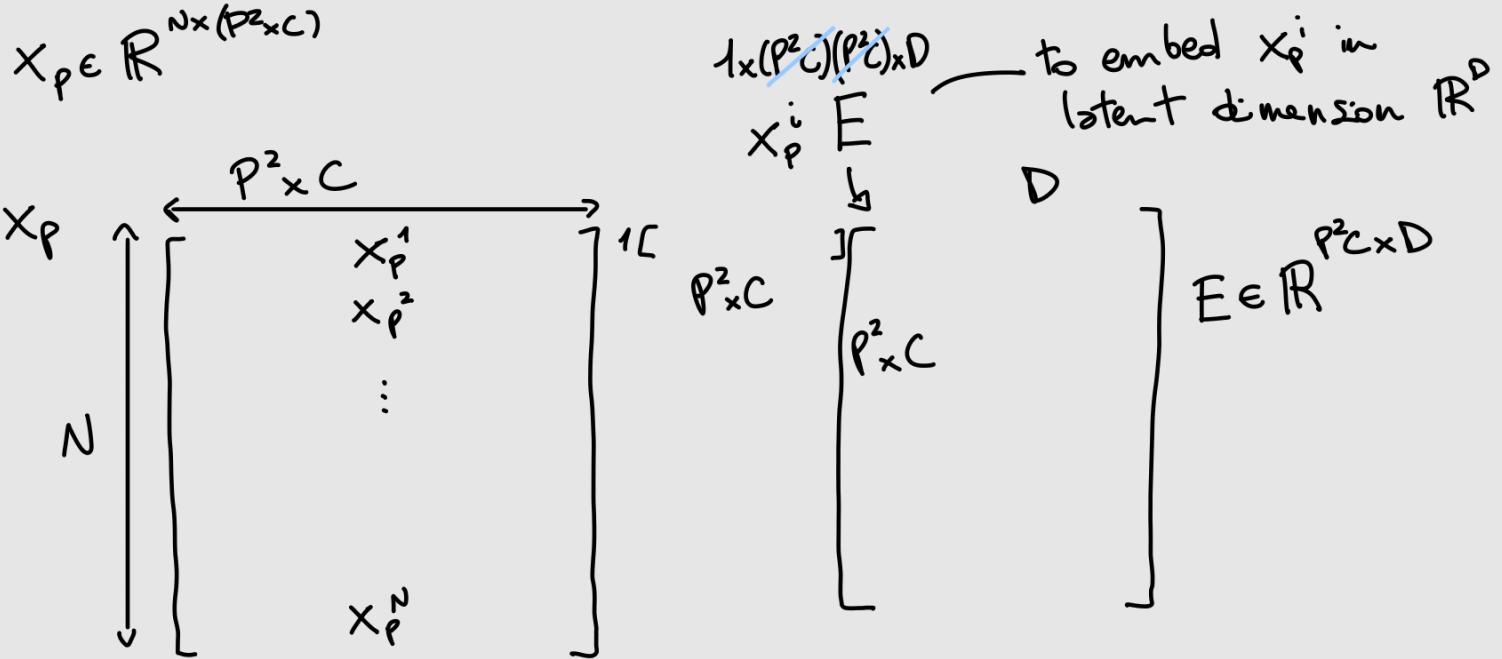
CLS token

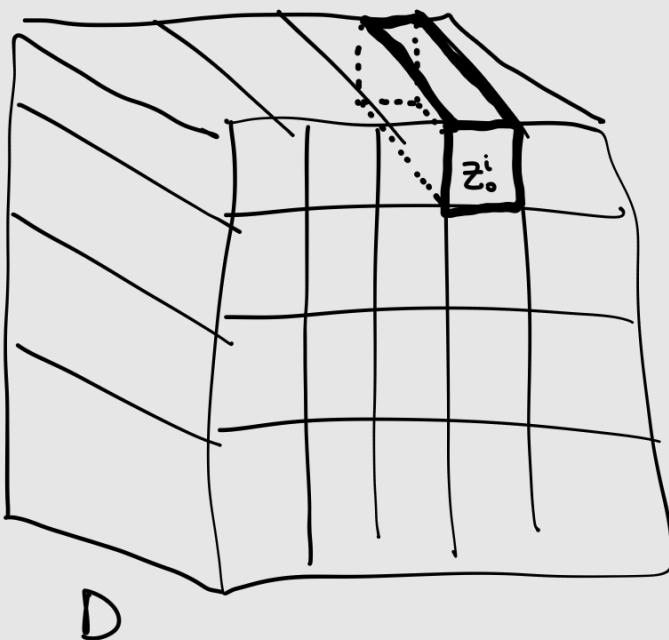
$$\mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{pos} \in \mathbb{R}^{(N+1) \times D}$$

$$\ell = 1 \dots L$$

$$\ell = 1 \dots L$$







MLP

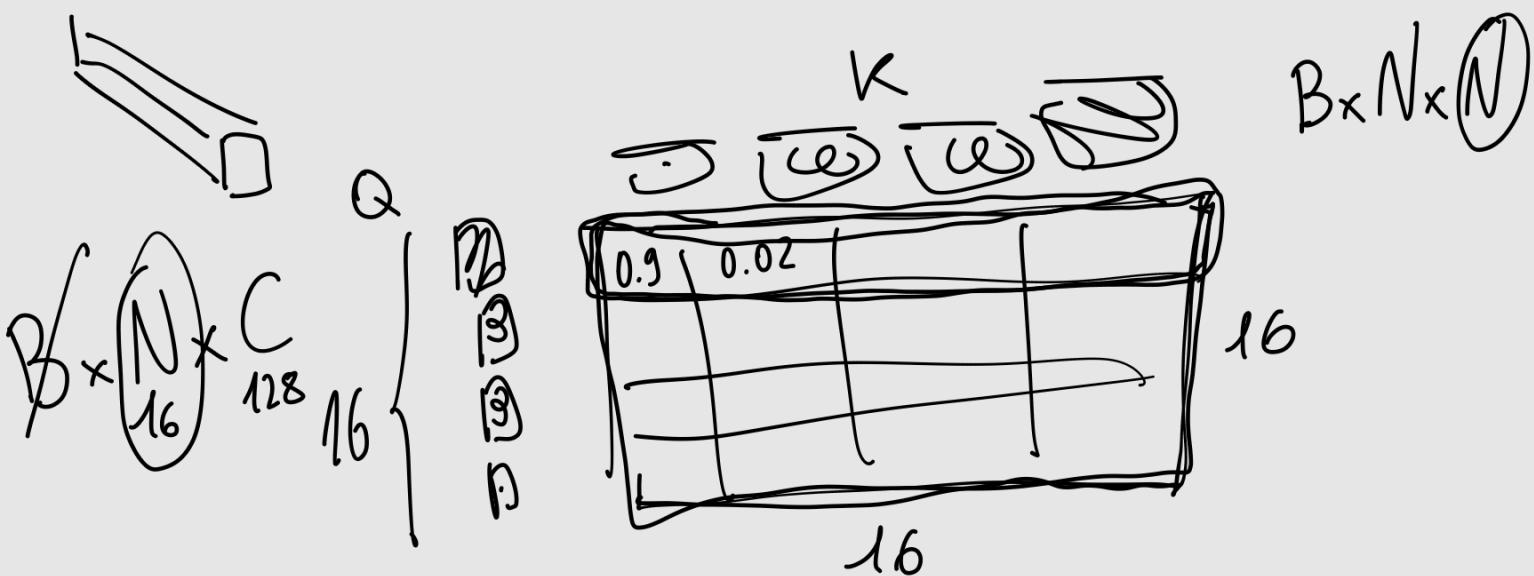
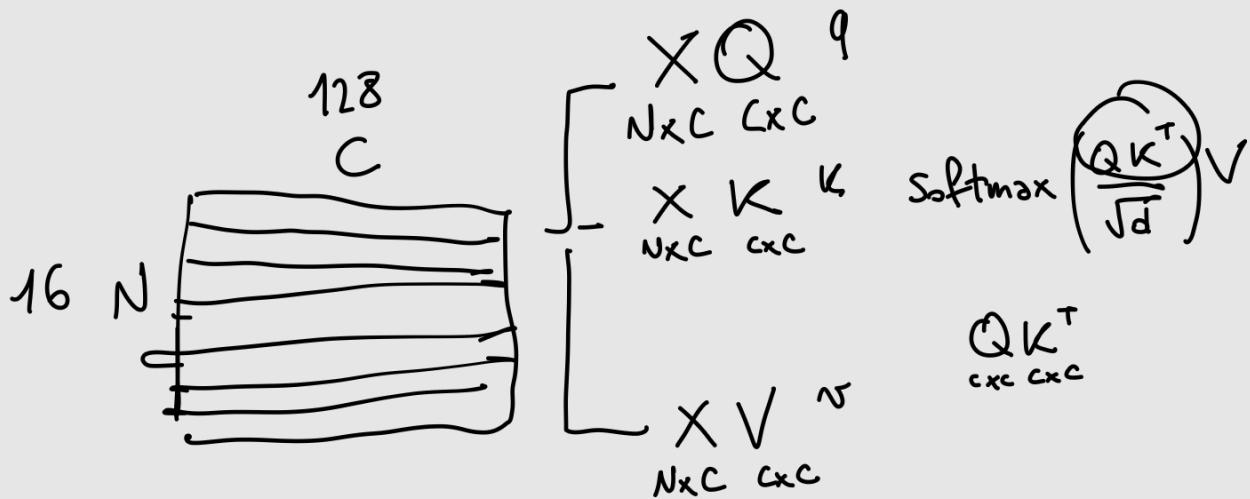
$$\begin{bmatrix} W_1 \\ W_{1J} \end{bmatrix} \begin{bmatrix} z_o^i \end{bmatrix}$$

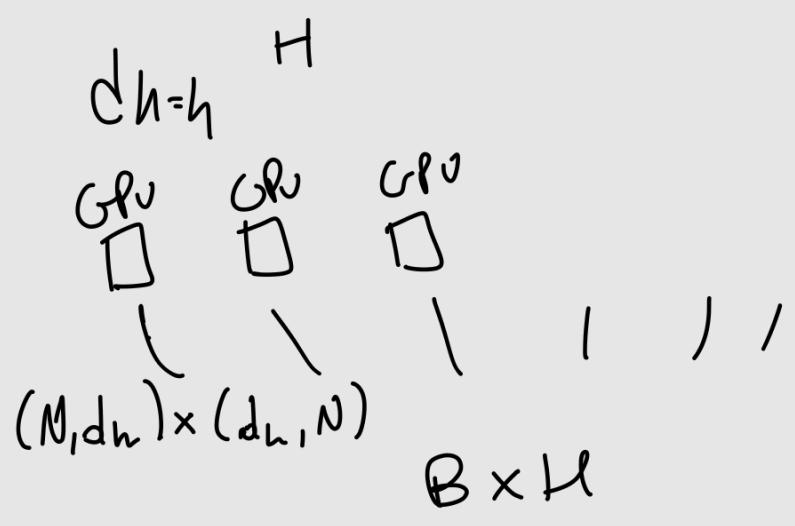
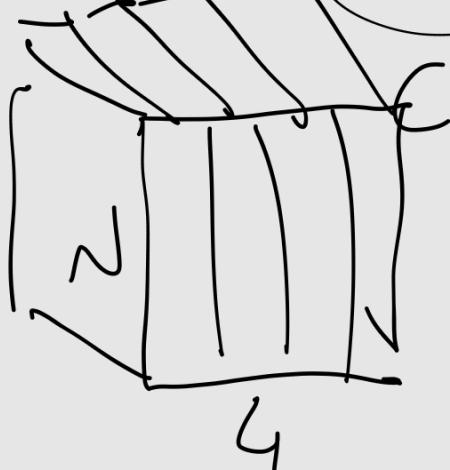
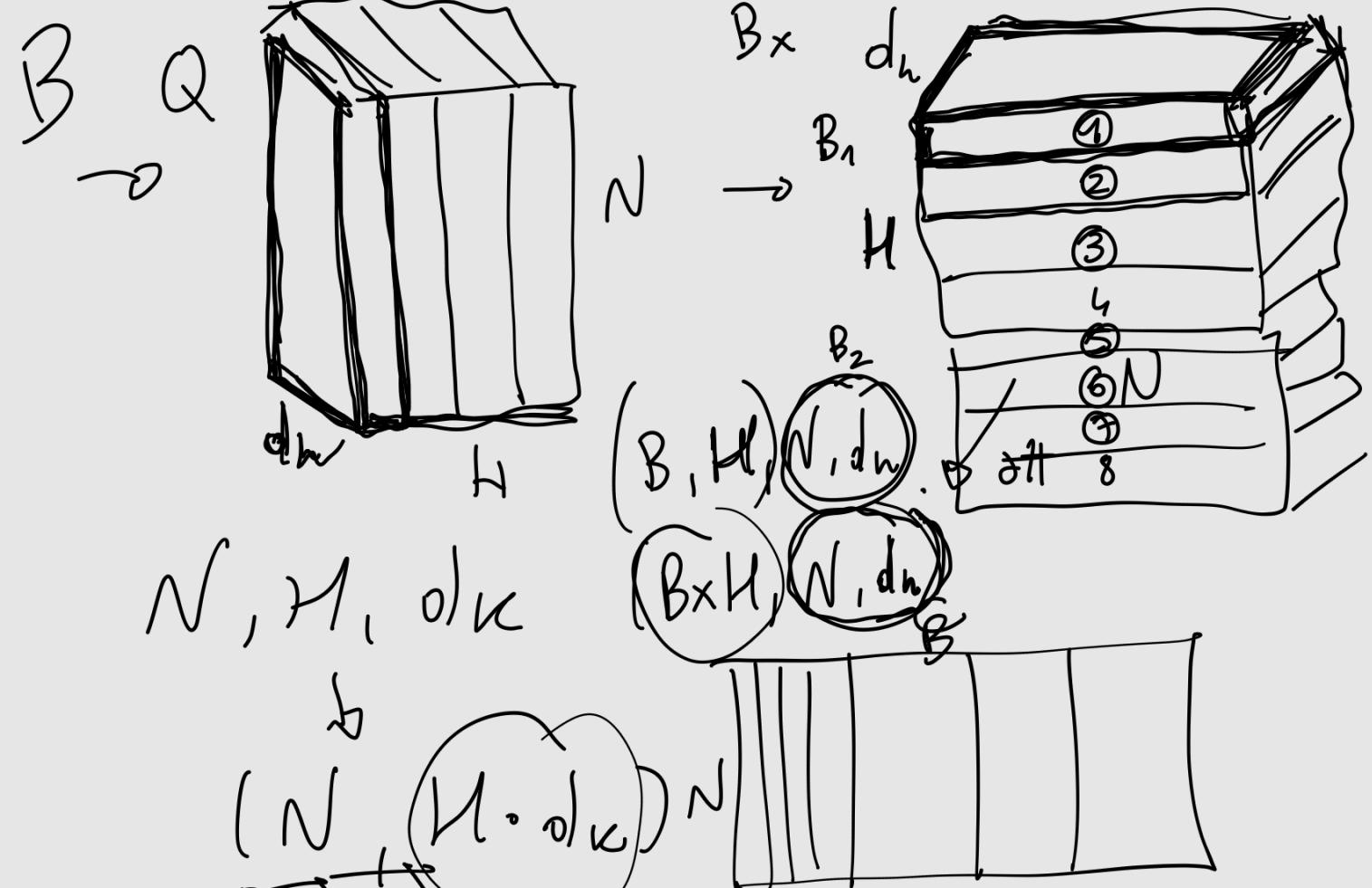
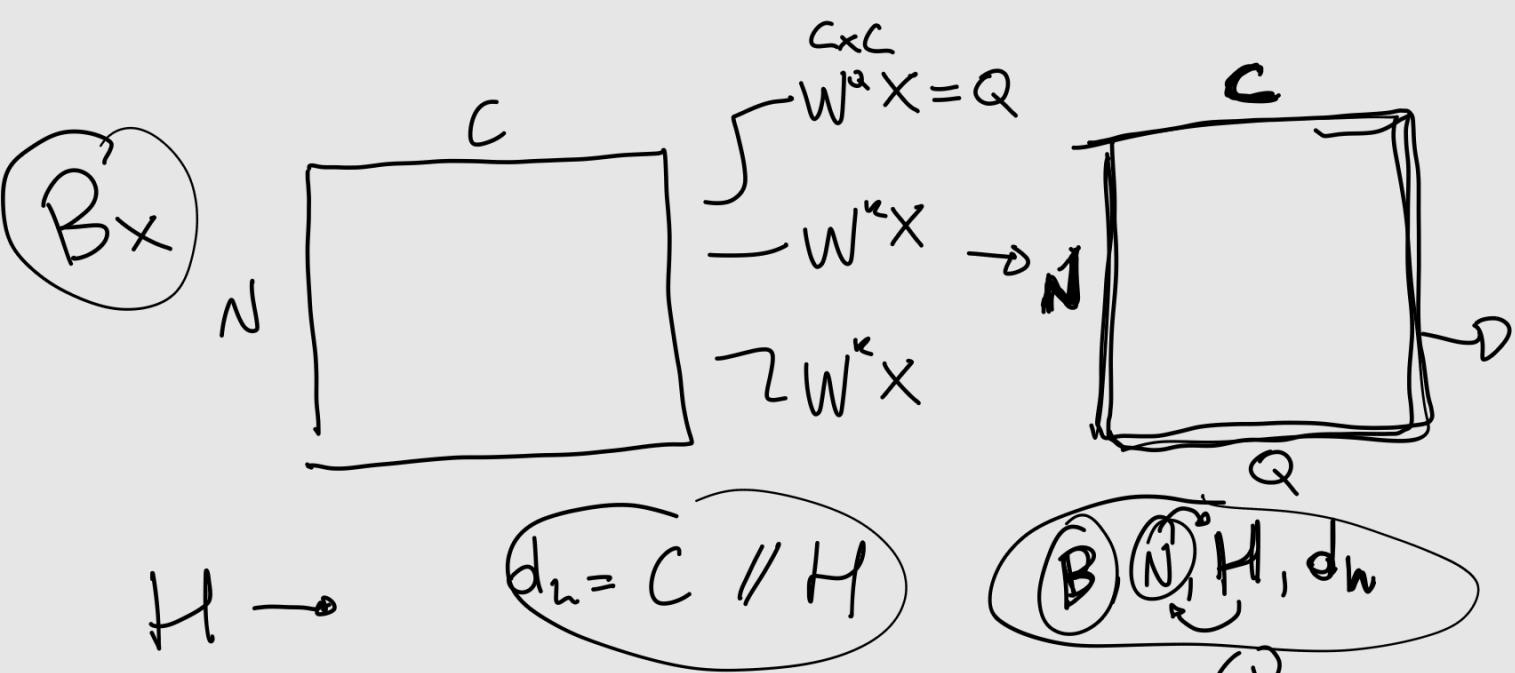
W_{1J} is
 1×1

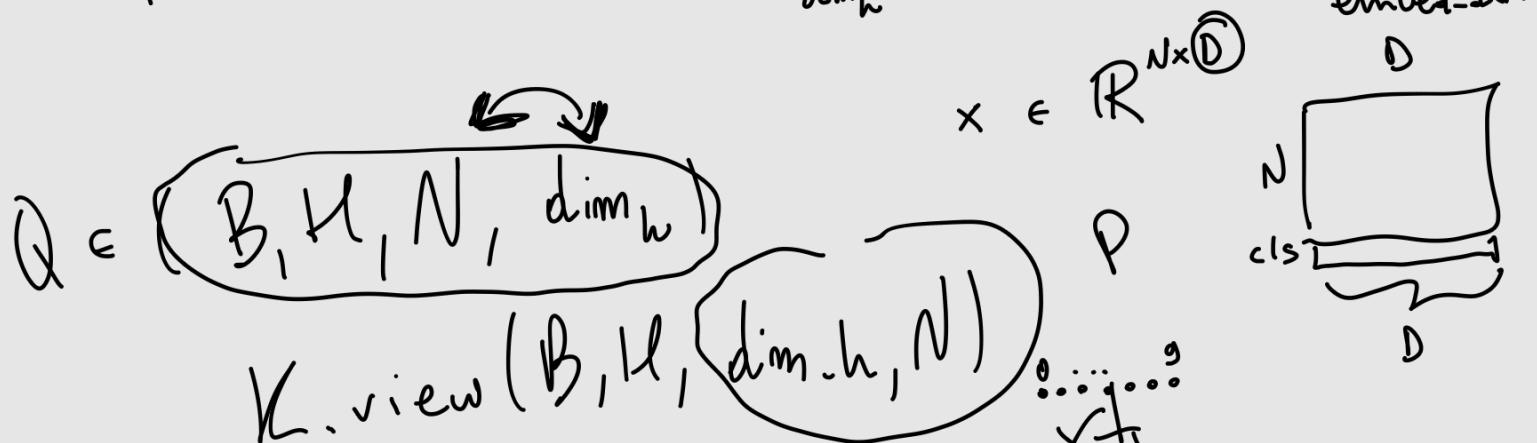
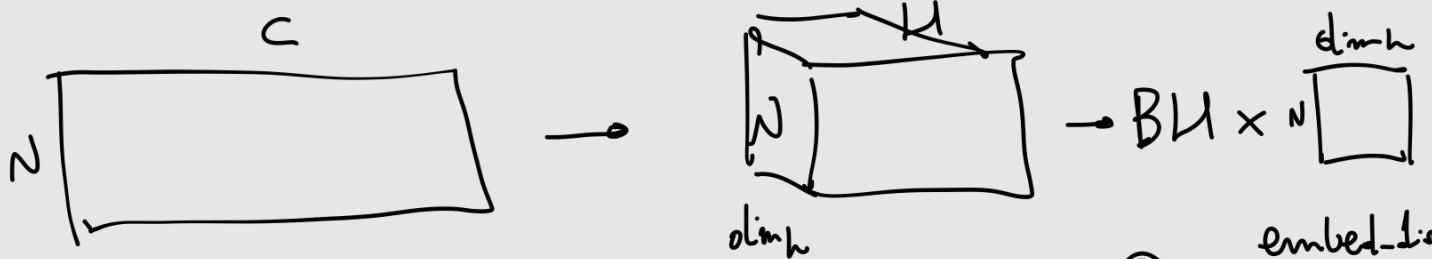
1×1 conv with D filters

$W_1 z_o^i$ has the same dimensions and it can be used as a saliency map to interpret the model

z_o

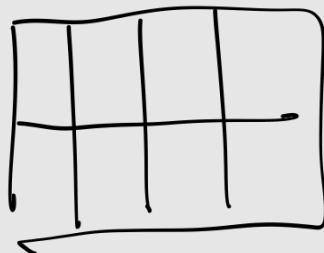
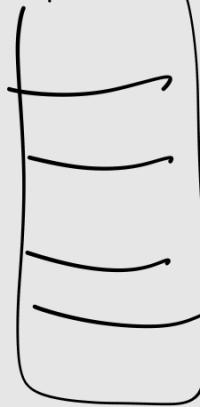






$x = B, H, N, \text{dim_h}$

$x = B, \overset{N}{\circlearrowleft}, C$

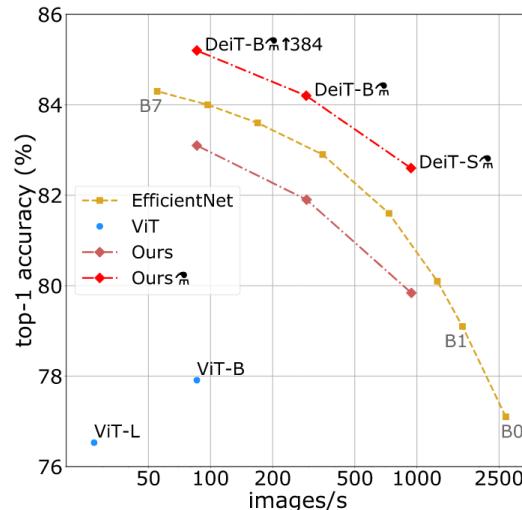


Attention process in Vision

Experiments with ViT (and variants DeiT, CaiT) transformers for image classification

State-of-the-art performance on ImageNet1k classification!

From ViT paper, **many tricks/discussions to simplify learning** in DeiT, CaiT, ...



Published as a conference paper at ICML 2021

Training data-efficient image transformers & distillation through attention

Hugo Touvron^{1,2} Matthieu Cord^{1,2} Matthijs Douze¹
Francisco Massa¹ Alexandre Sablayrolles¹ Hervé Jégou¹