

COURS RDFIA deep Image

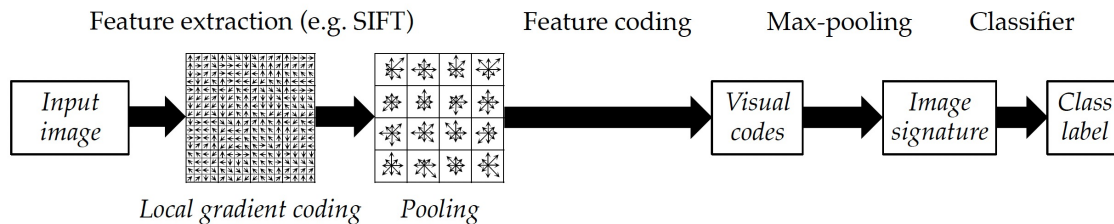
Matthieu Cord
Sorbonne University

Course Outline – Week timeline

1. **Computer Vision basics:** Visual (local) feature detection and description, Bag of Word Image representation
2. **Supervised learning:** Introduction to Neural Networks (NNs)
3. **Machine Learning basics:** Risk, Classification, Datasets, benchmarks and evaluation, Linear classification (SVM)
4. **Convolutional Nets** for visual classification
5. **Large deep convnets and Vision Transformers**
6. Beyond ImageNet: FCNs and Segmentation
7. Transfer Learning and domain adaptation
8. Generative models with (conditional) GANs
9. Vision-Language models
10. Control
11. Explainable AI and applications
- 12/14. Bayesian deep learning

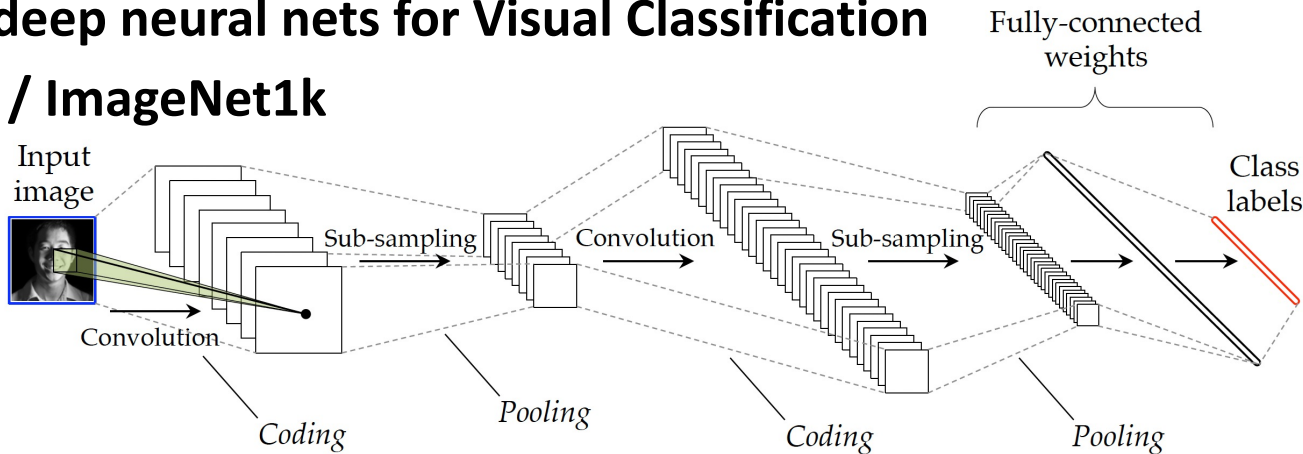
Context: Image classification **Before/After** ImageNet (2009)

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



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

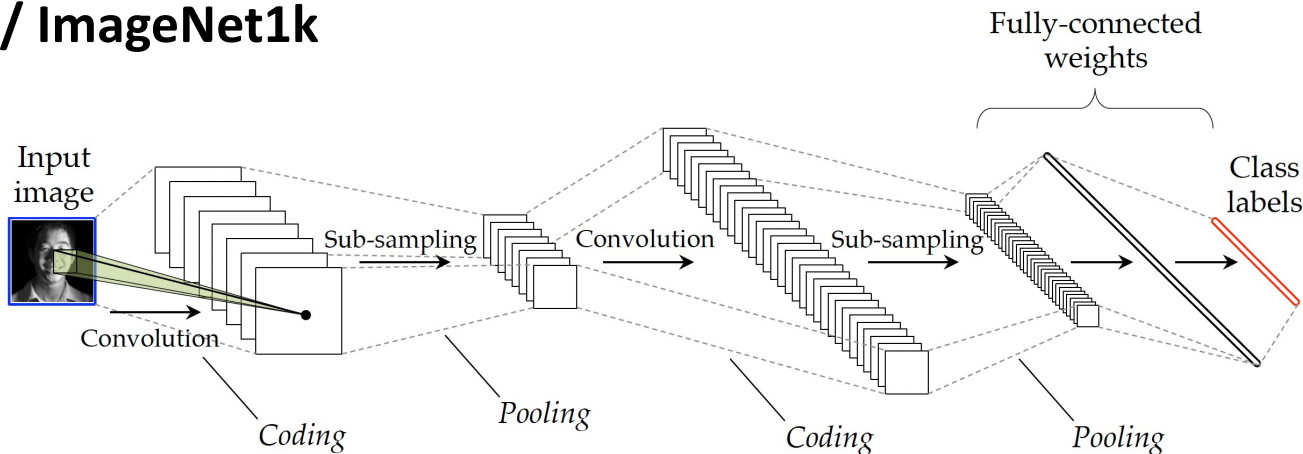
The star: **ConvNet / ImageNet1k**



Context: Image classification **After** ImageNet (2009)

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

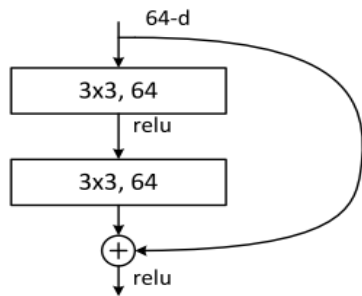
The star: **ConvNet** / ImageNet1k



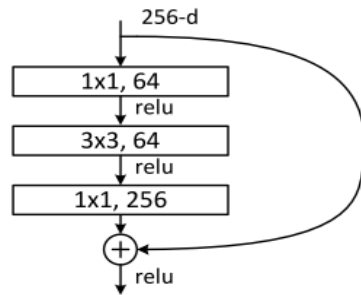
AlexNet 2012

- 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

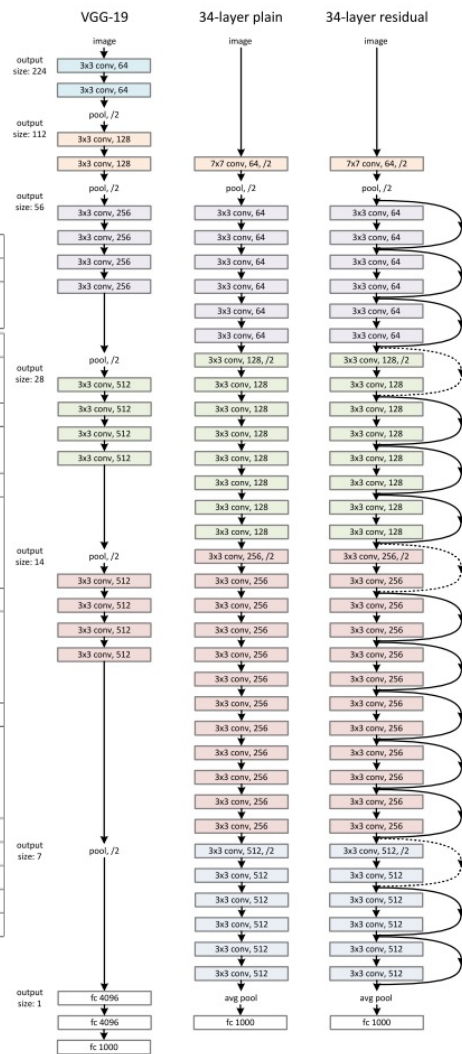


A naïve residual block



“**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
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
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
maxpool			
FC-4096	FC-4096	FC-4096	FC-4096
FC-4096	FC-4096	FC-4096	FC-4096
FC-1000	FC-1000	FC-1000	FC-1000
soft-max	soft-max	soft-max	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

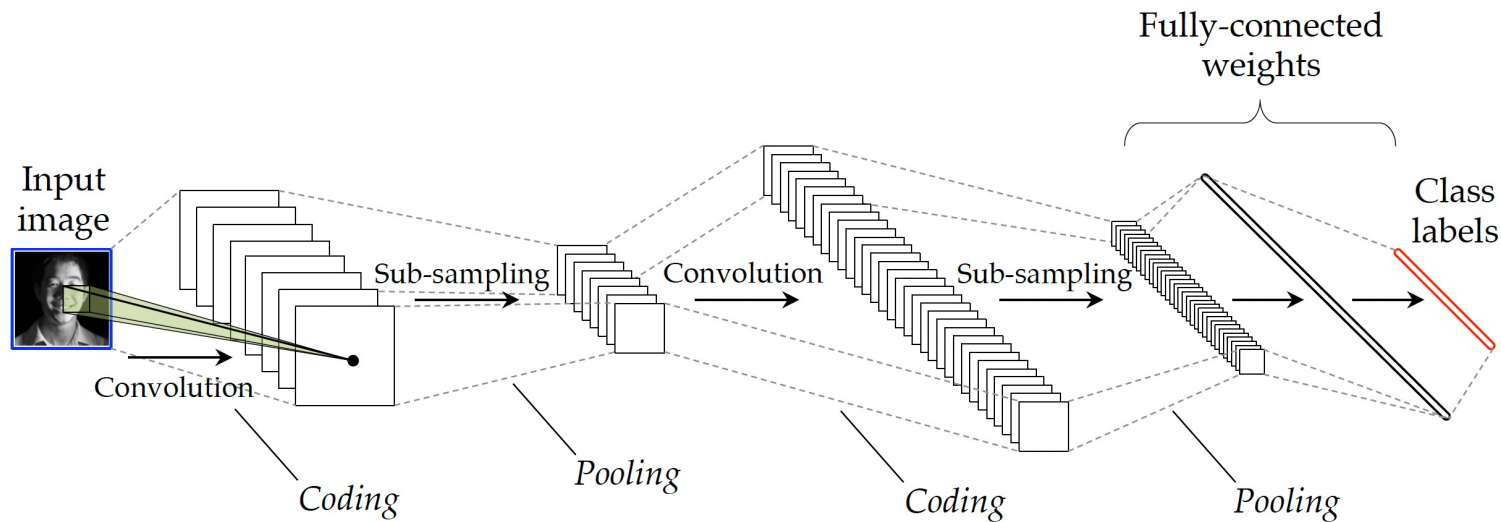
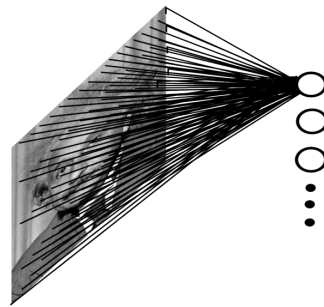
1. Attention and Vision Transformers

- 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*

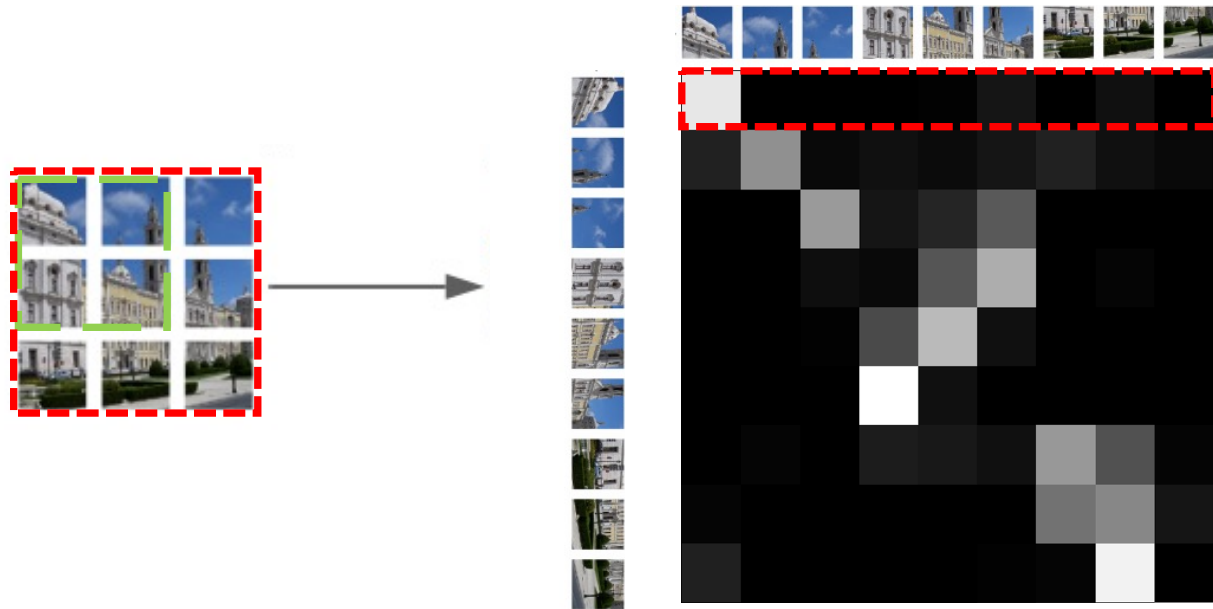
Rq: less local after many layers



Global (Self) attention

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

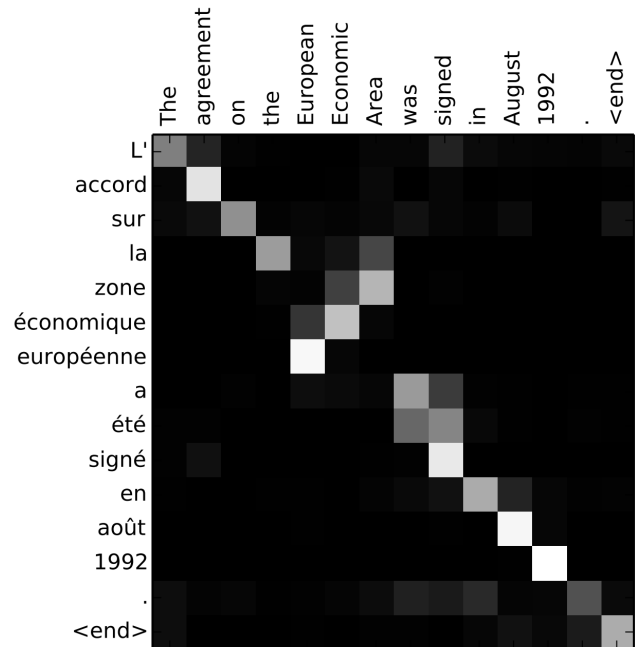
=> Different than convNet!



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

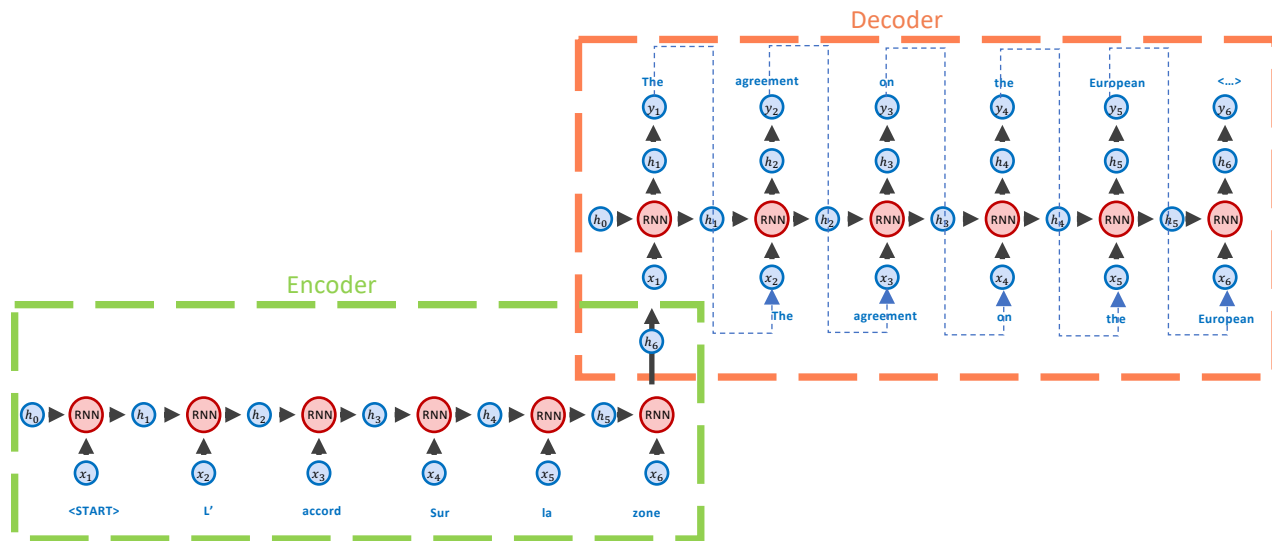


Attention process in NLP

Basic language translation models: Encoder/Decoder

Ex.: Seq2Seq -- RNNs2RNNs

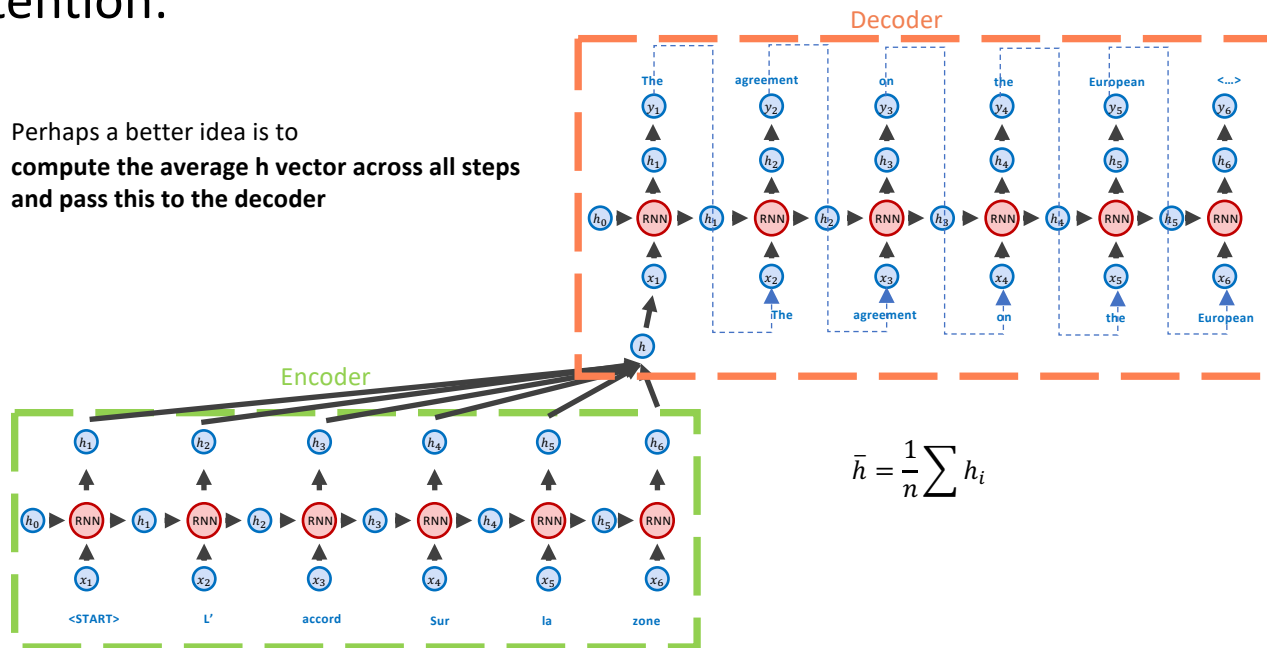
Cross-attention for language translation in at the end of Encoder



Attention process in NLP

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

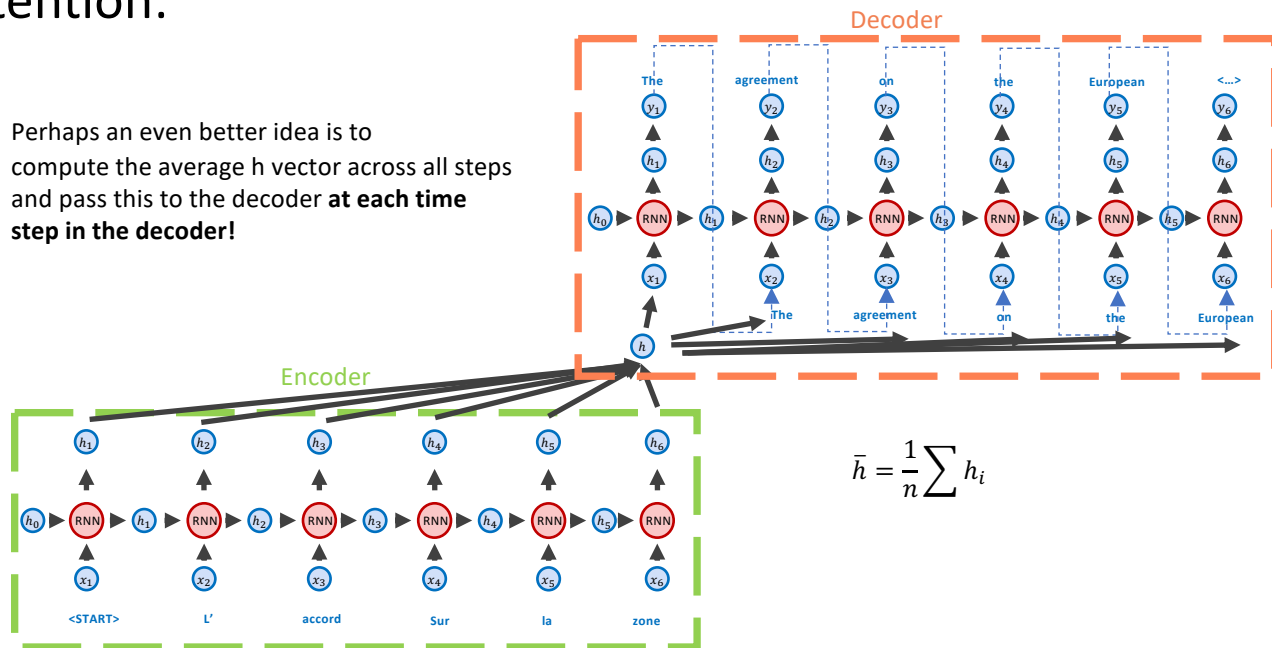
Cross-attention:



Attention process in NLP

Basic language translation models: Encoder/Decoder

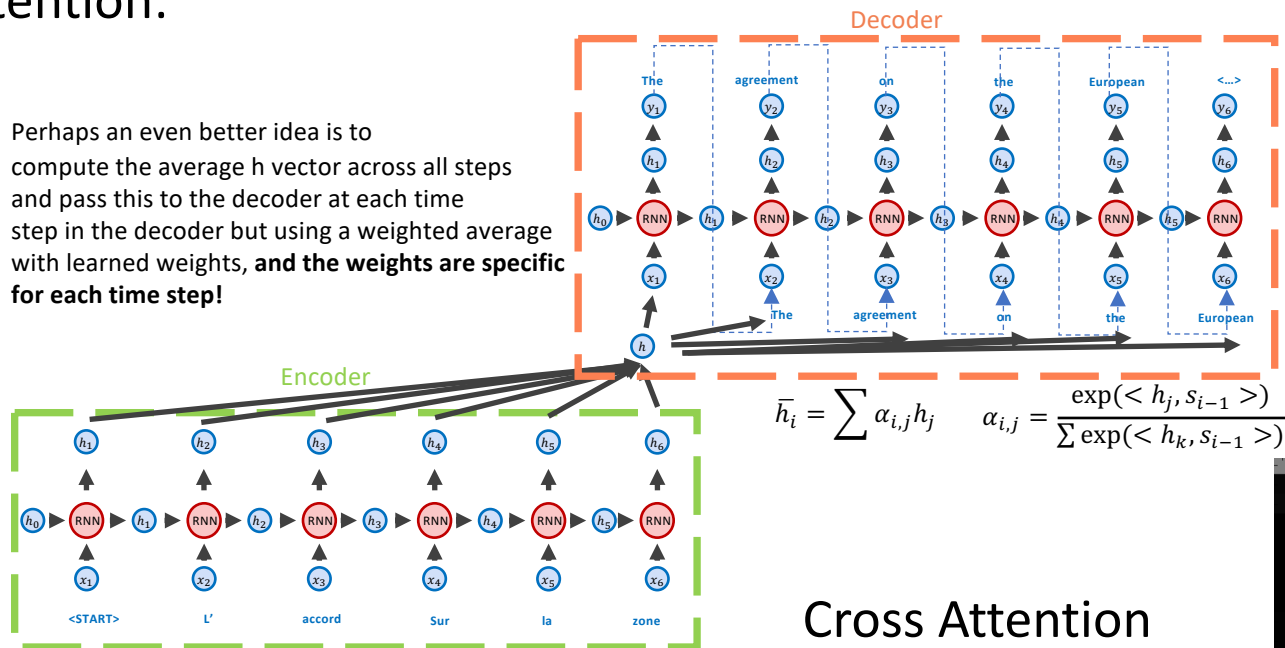
Cross-attention:



Attention process in NLP

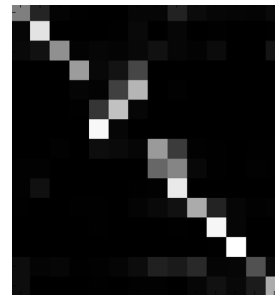
Basic language translation models: Encoder/Decoder

Cross-attention:



Cross Attention

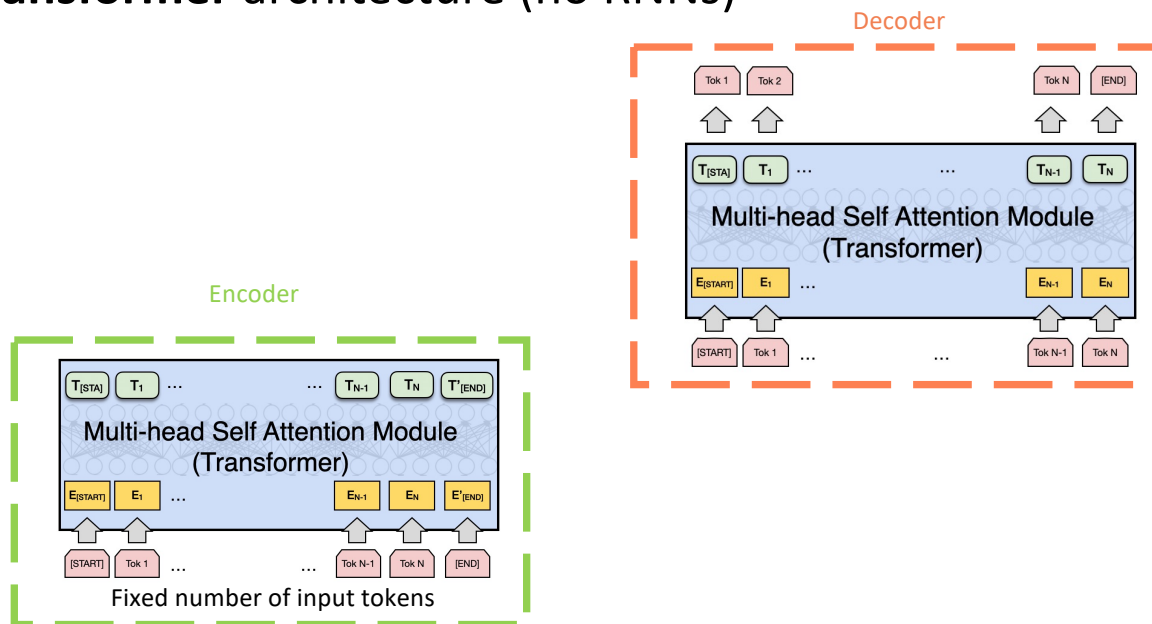
Encoder/ Decoder



Attention process in NLP

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

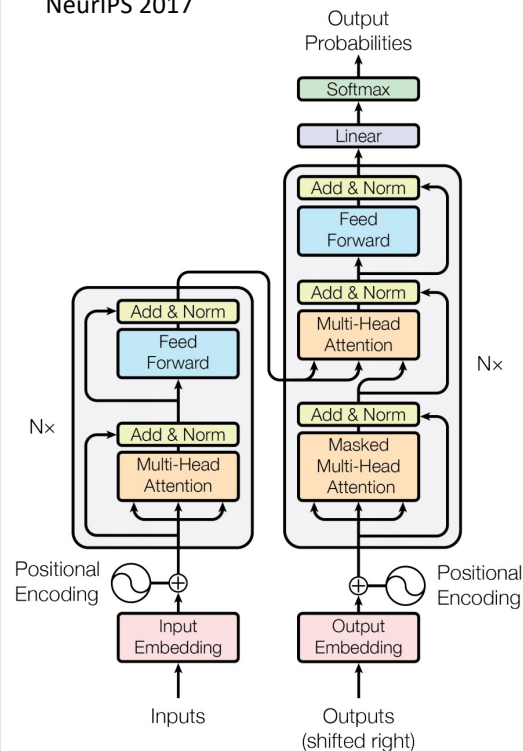
Transformer architecture (no RNNs)



[Vaswani et al. Attention is all you need]

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

NeurIPS 2017

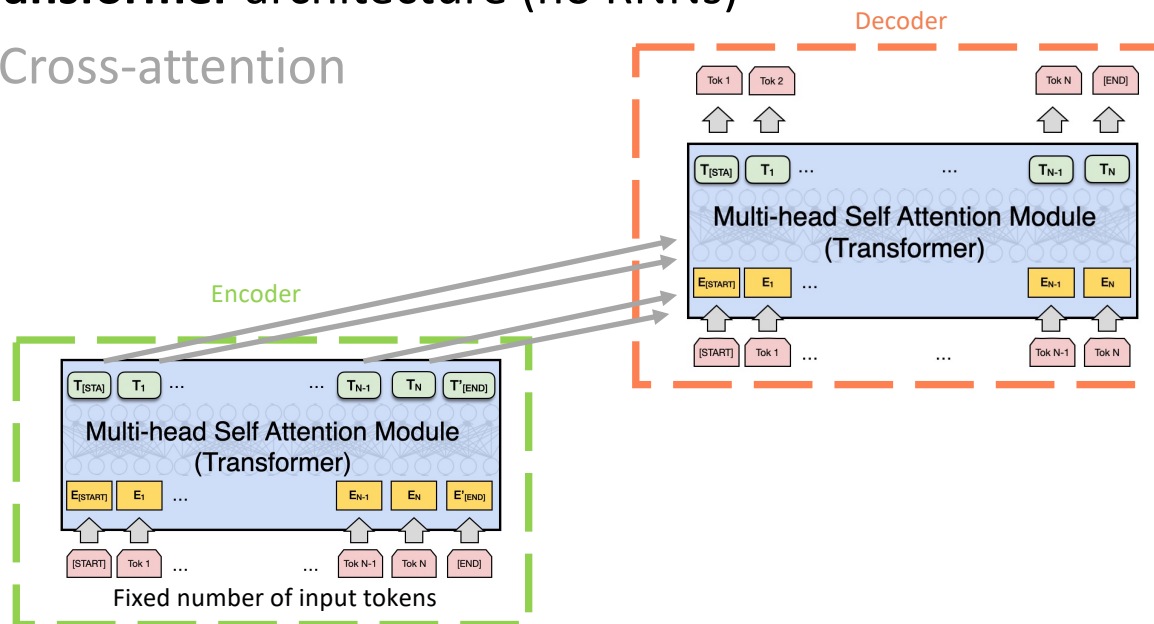


Attention process in NLP

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

Transformer architecture (no RNNs)

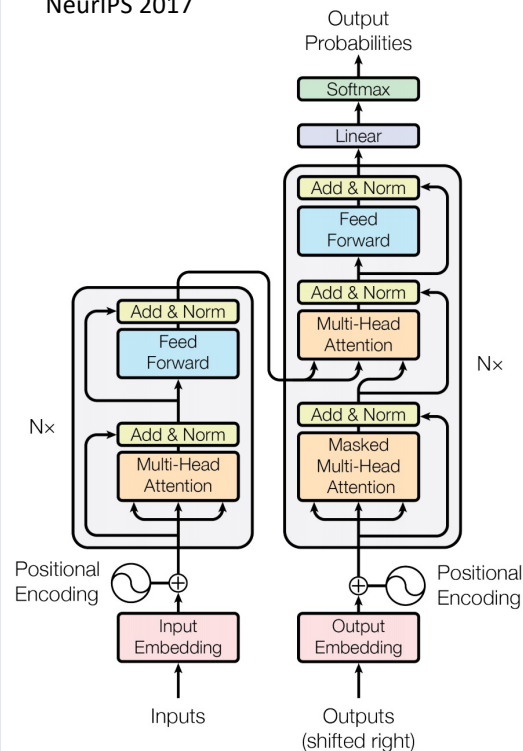
- Cross-attention



[Vaswani et al. Attention is all you need]

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

NeurIPS 2017

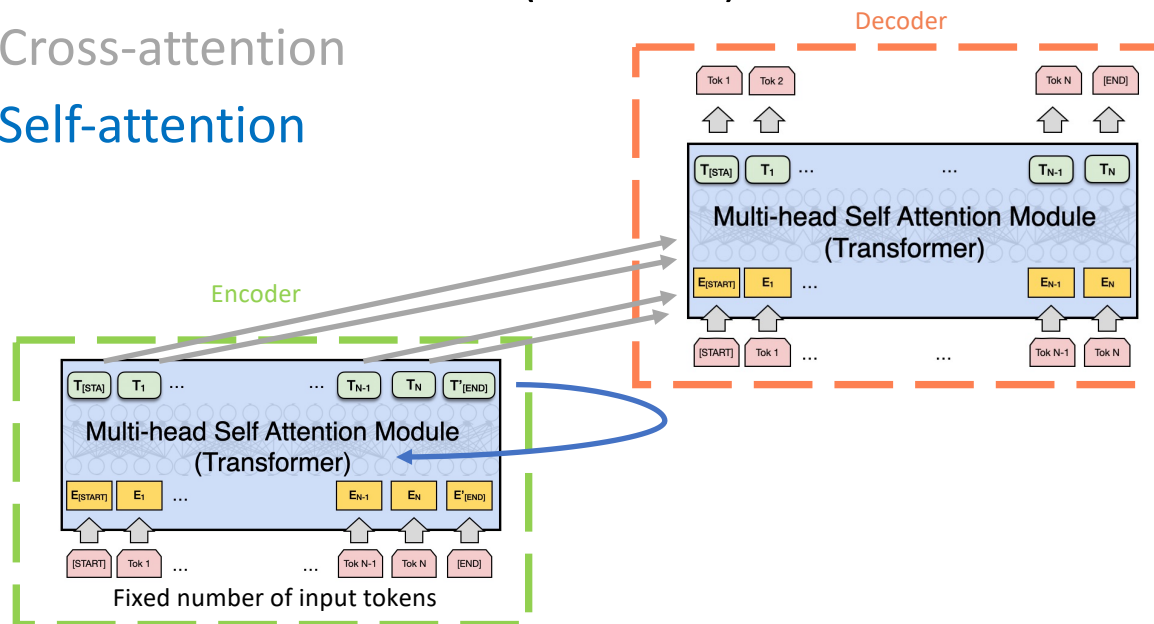


Attention process in NLP

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

Transformer architecture (no RNNs)

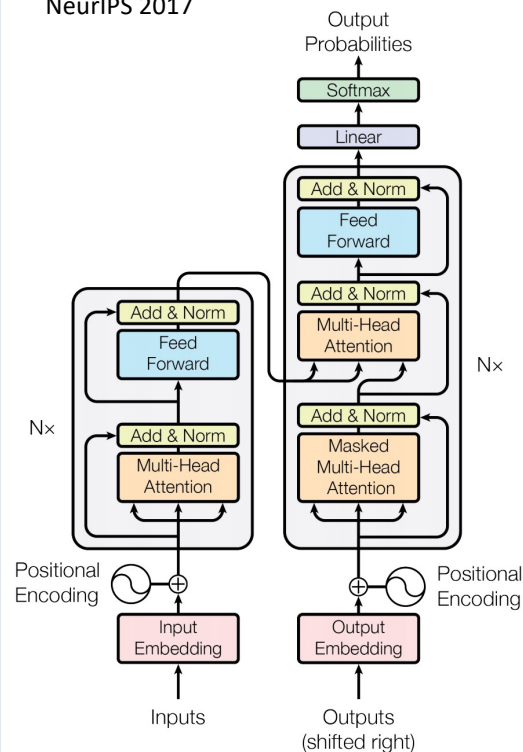
- Cross-attention
- Self-attention



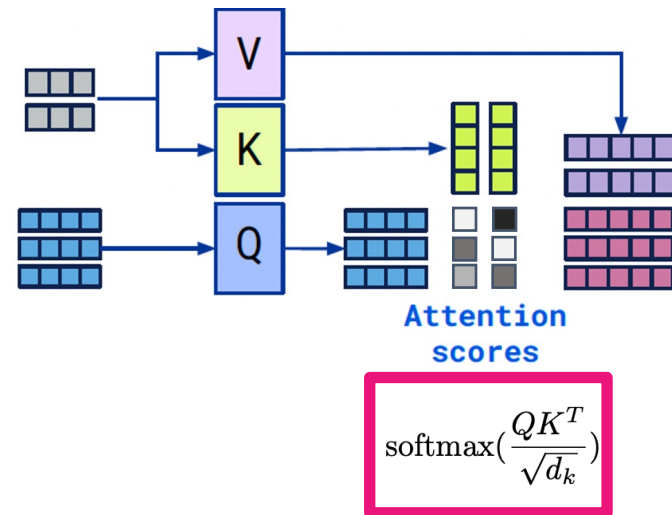
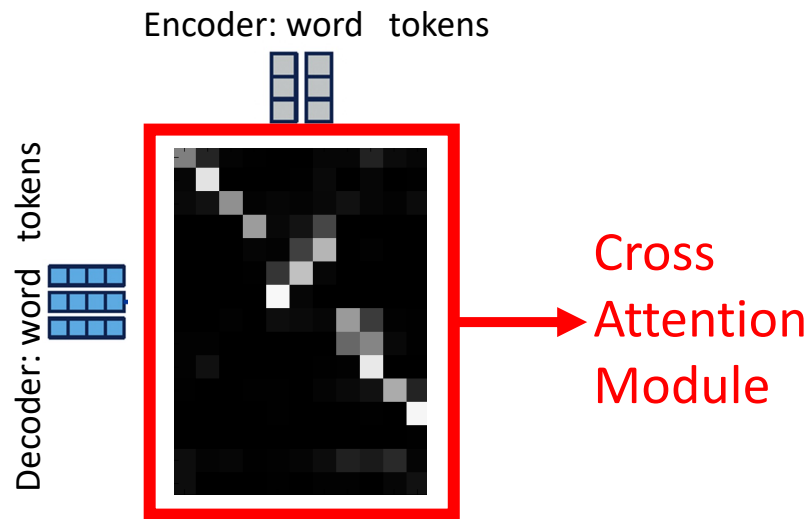
[Vaswani et al. Attention is all you need]

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

NeurIPS 2017



Attention process in NLP



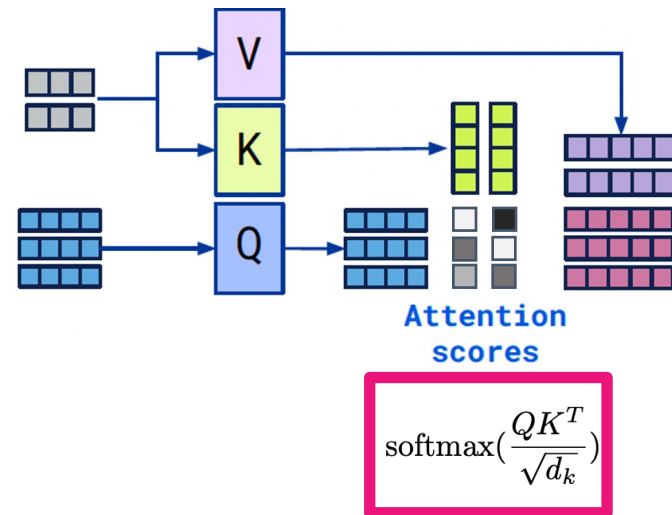
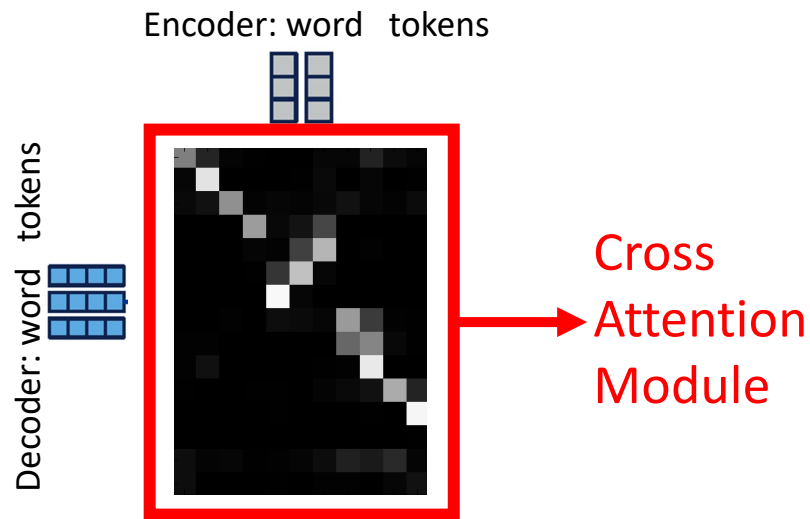
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Outline

1. Attention and Vision Transformers (ViT)

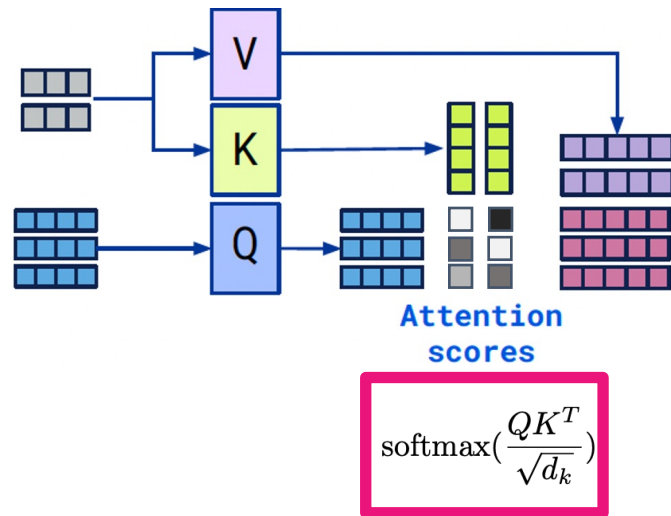
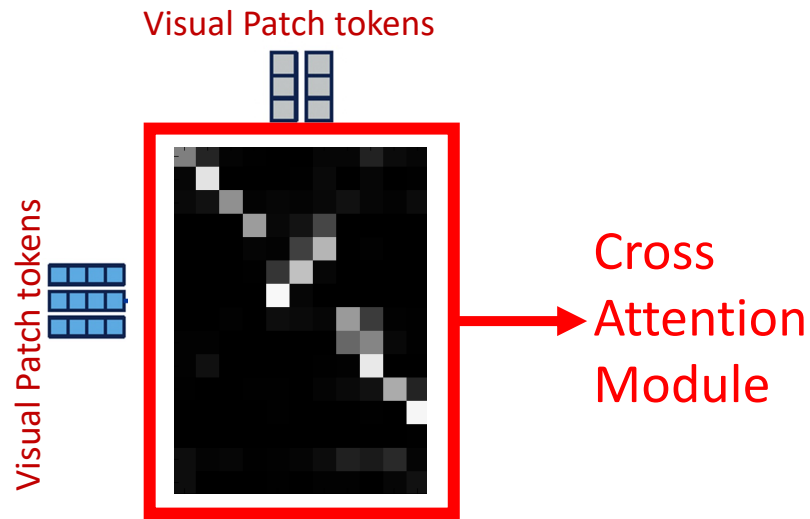
- NLP: Attention is all you need
- **Transformer for image classification**

Attention process in NLP



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

Attention process in Vision



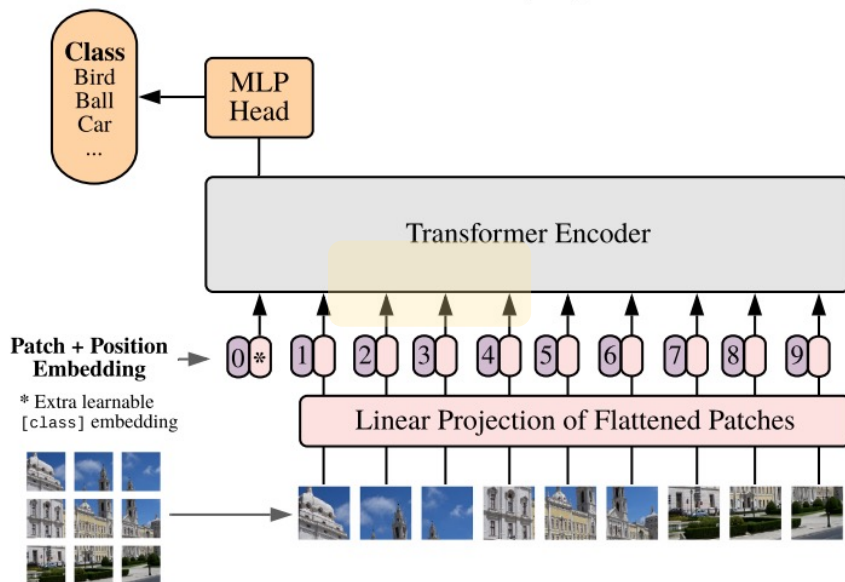
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

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

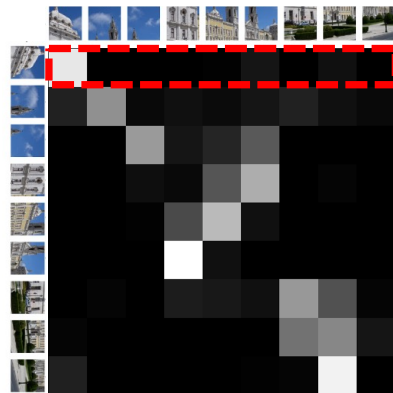
AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy^{*,†}, 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

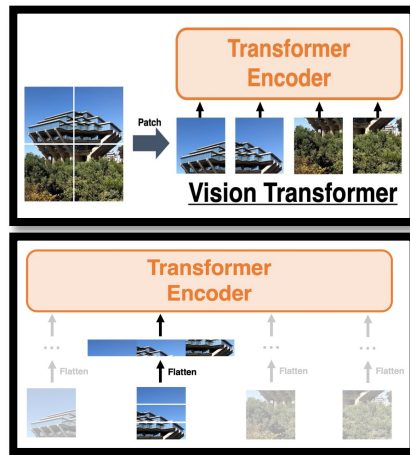
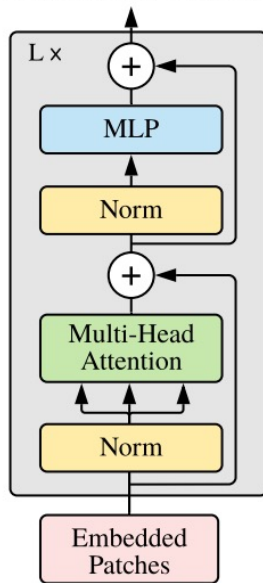
Google Research, Brain Team

{adosovitskiy, neilhoulby}@google.com



Attention process in Vision

Transformer Encoder



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

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

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

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

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

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

$$N = HW/P^2$$

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

CLS token

$$\ell = 1 \dots L$$

$$\ell = 1 \dots L$$

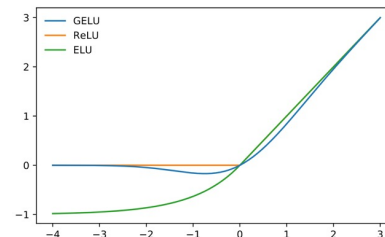
[class=CLS] token: a learnable embedding to the sequence of embedded patches

Layer norm (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

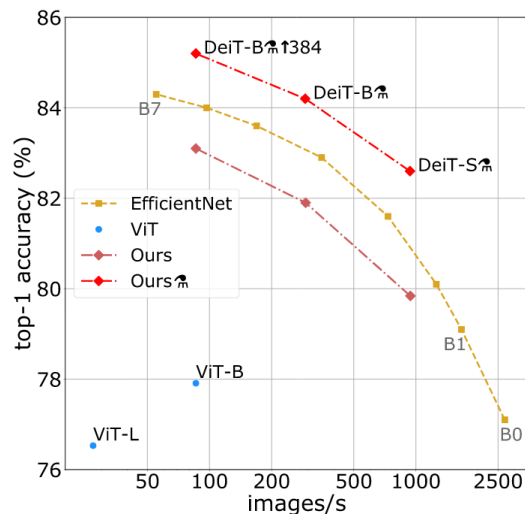


Attention process in Vision

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

Attention process in Vision

How to choose the image splitting?

Pb: quadratic complexity with the nb of patches

Many kinds of hybrid architectures with convnets/transformers

Ex: Swin Transformers

Published as a conference paper at ICCV 2021

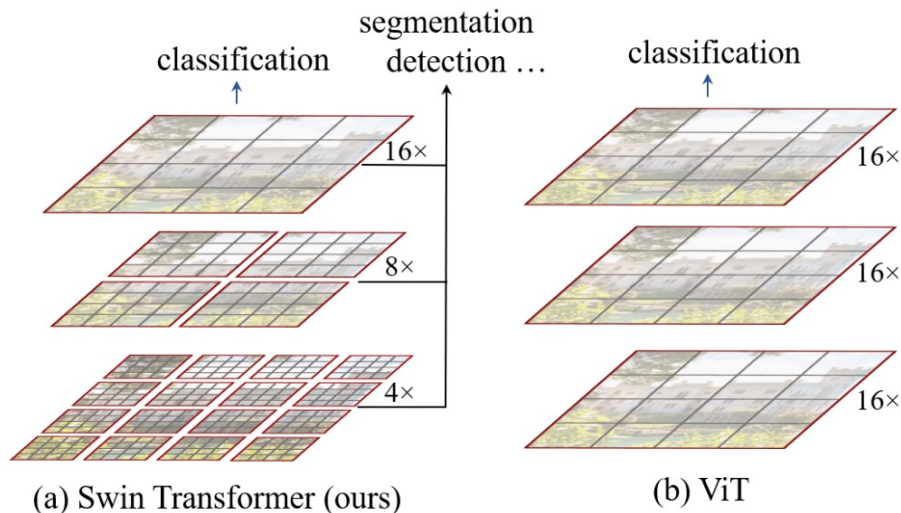
Swin Transformer: Hierarchical Vision Transformer using Shifted Windows

Ze Liu^{1,2†*} Yutong Lin^{1,3†*} Yue Cao^{1*} Han Hu^{1*‡} Yixuan Wei^{1,4†}
Zheng Zhang¹ Stephen Lin¹ Baining Guo¹

¹Microsoft Research Asia ²University of Science and Technology of China

³Xian Jiaotong University ⁴Tsinghua University

{v-zeliu1, v-yutlin, yuecao, hanhu, v-yixwe, zhez, stevelin, bainguo}@microsoft.com



Outline

1. Attention and Vision Transformers (ViT)

- NLP: Attention is all you need
- Transformer Encoder ViT with Self Attention for image classification

