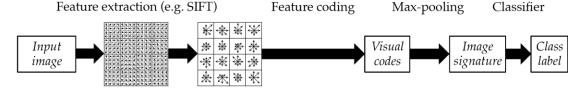
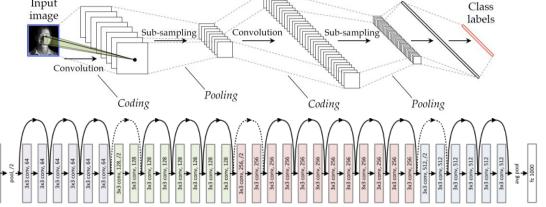
Image classification: where we are and what is missing?

Input

The 2000s: *BoWs + SVM*



The 2010s: Very *Large* ConvNets



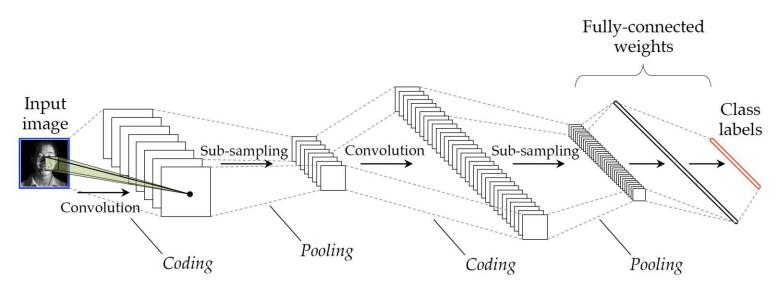
2020: The star: **ResNet**

Next step? What is missing?

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

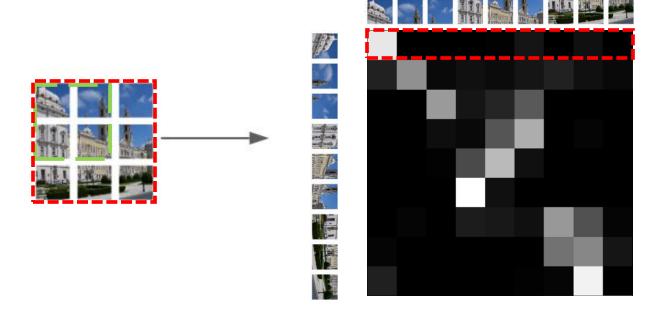
Rq: less local after many layers



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!



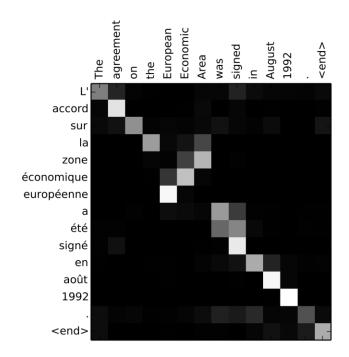
Vision Transformers

- 1. NLP: Attention is all you need
- 2. Transformer for Image Classification

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

Attention between words in Machine translation process:

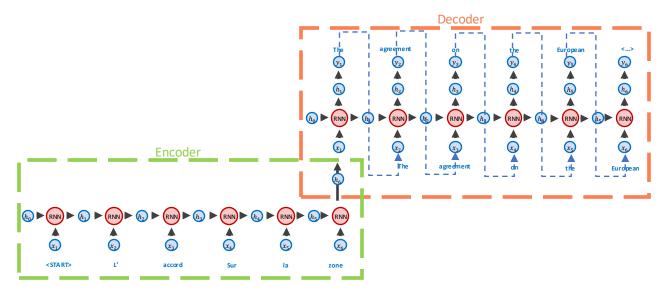
- 1. Computing of weights
- 2. Use them to compute new features



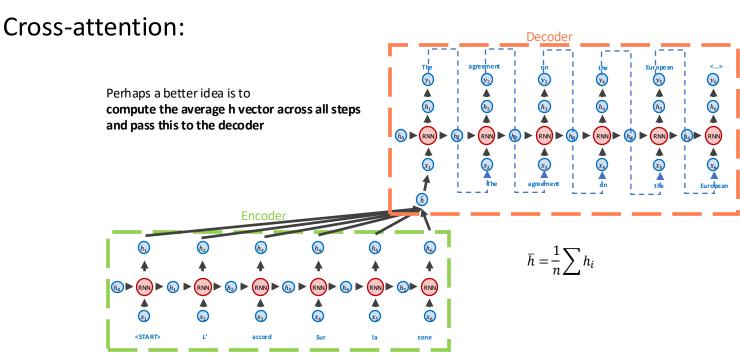
Basic language translation models: Encoder/Decoder

Ex.: Seq2Seq -- RNNs2RNNs

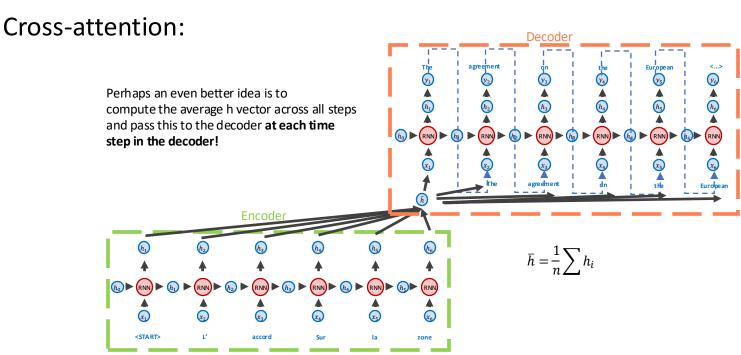
Cross-attention for language translation in at the end of Encoder



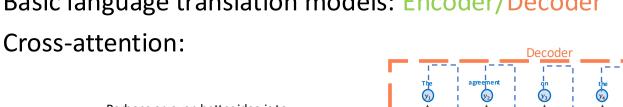
Basic language translation models: Encoder/Decoder

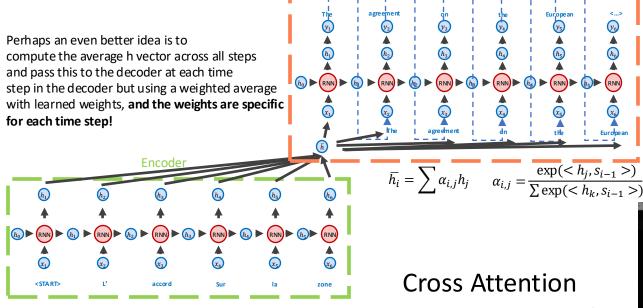


Basic language translation models: Encoder/Decoder



Basic language translation models: Encoder/Decoder

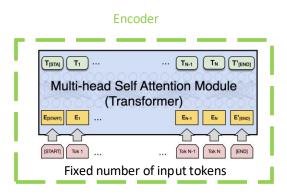


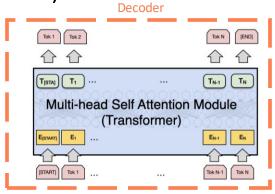


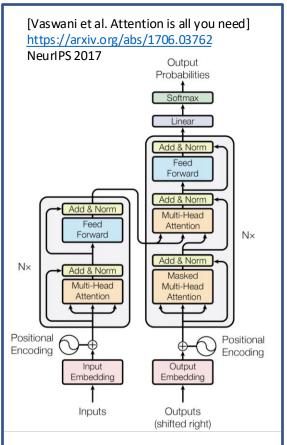
Encoder/ Decoder

Basic language translation models: Encoder/Decoder

Transformer architecture (no RNNs)



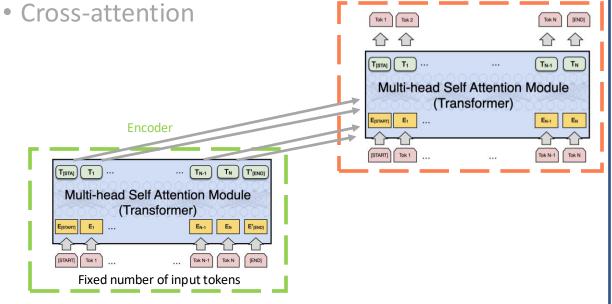


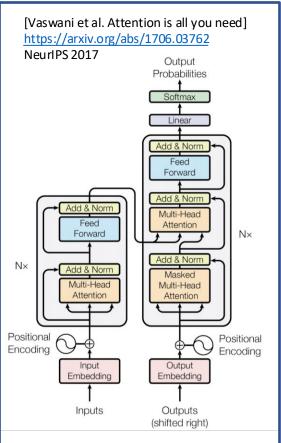


Basic language translation models: Encoder/Decoder

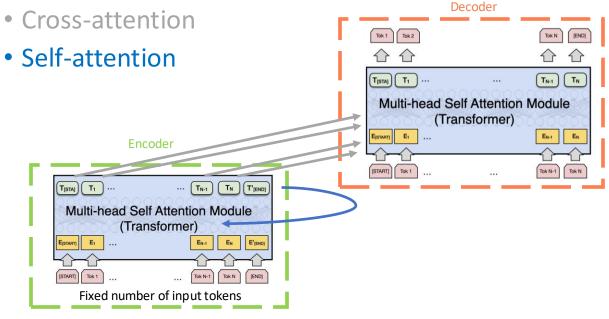
Transformer architecture (no RNNs)

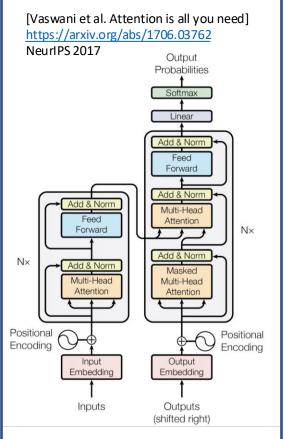
Decoder

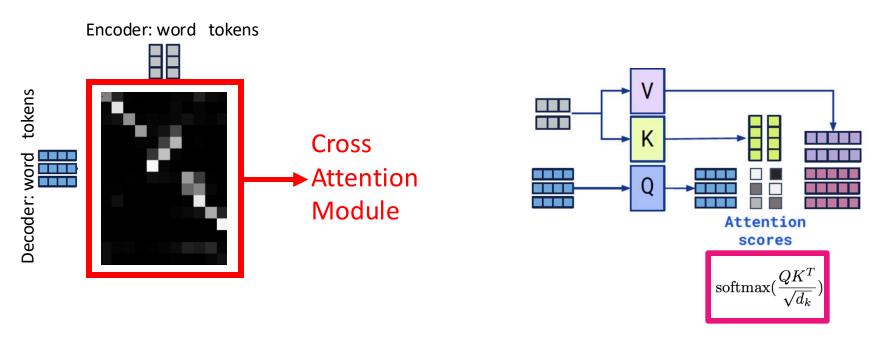




Basic language translation models: Encoder/Decoder Transformer architecture (no RNNs)



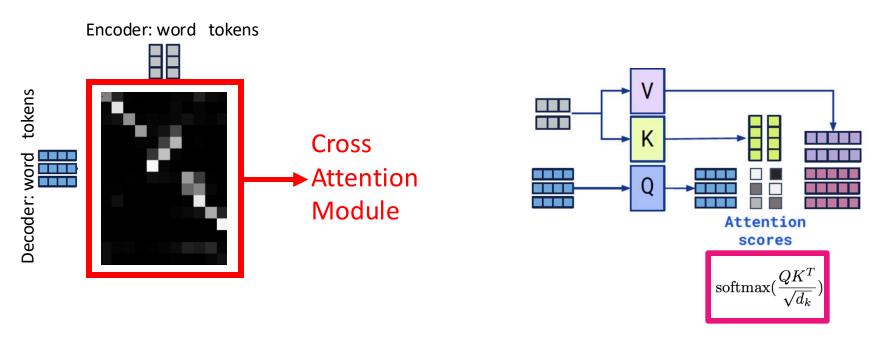




Attention
$$(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

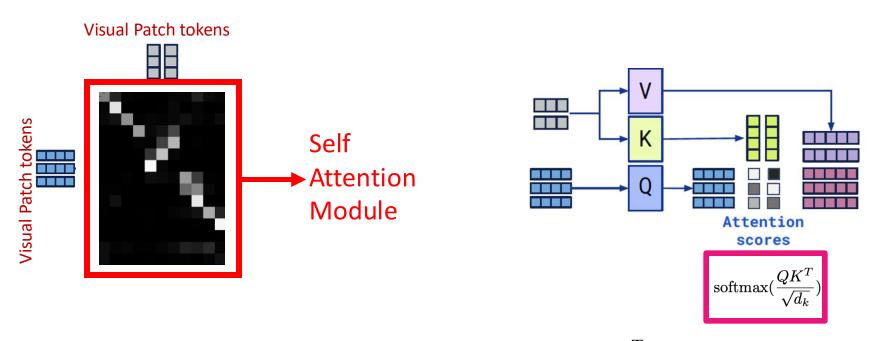
Vision Transformers

- 1. NLP: Attention is all you need
- 2. Transformer for Image Classification



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

Attention process in Vision



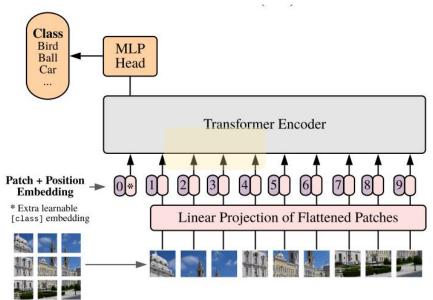
Attention
$$(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

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

Attention process in Vision

Is it possible to mimic this attentionbased 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

