

# Deep Learning (M1)

## Attention mechanisms

Denis Coquenet



## Knowledge

- Principle of attention mechanisms
- Transformer architecture
- Transformer attention formula

## Skills and know-how

- Integrate transformer attention in encoder / decoder
- Determine size and meaning of intermediate representations in attention mechanisms
- Analyze attention maps

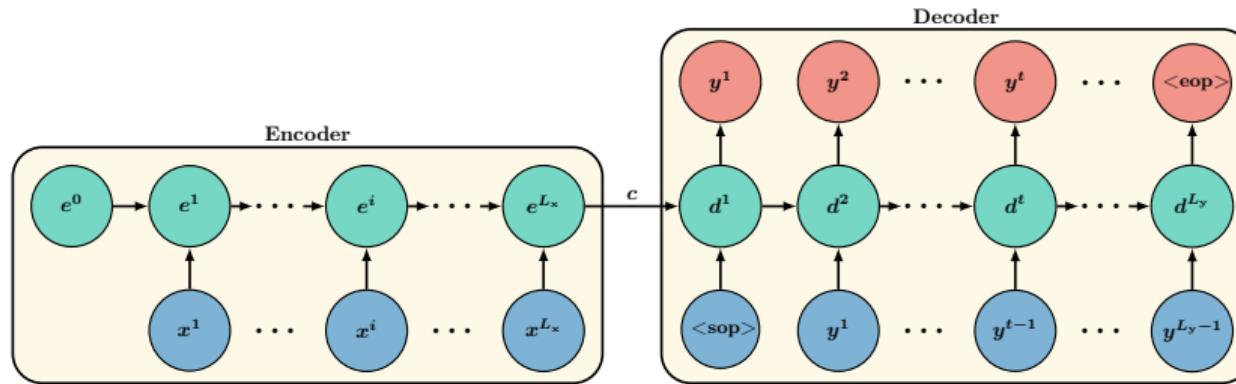
## Limitations of RNNs for large sequences

For feature extraction (encoder):

- Difficult to preserve information from start to end of the sequence

For seq2seq:

- The size of the context vector  $c$  is fixed, no matter the sequence length (compression, information loss)



1 Attention mechanisms

2 Transformer

3 Dealing with increasing need for data

4 Transformer's everywhere

## Idea

From a static context vector:

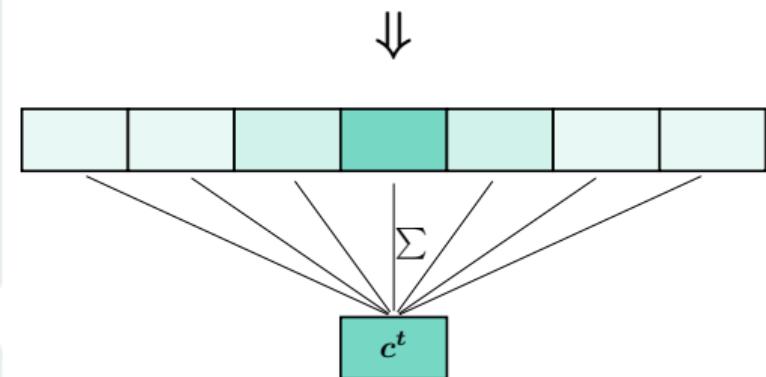
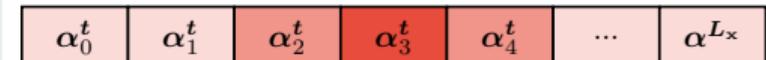
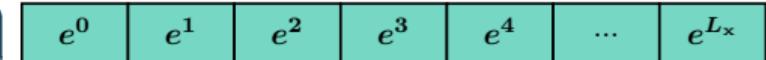
$$c = f(x)$$

to a dynamic context vector:

$$c^t = \sum_{i=1}^{L_x} \alpha_i^t e^i$$

with:  $\sum_i \alpha_i = 1$ .

$\alpha_i$  are called attention weights



## Encoder

Forward pass performed only once leading to  $\{e^1, \dots, e^{L_x}\}$

## Decoder

Iterative process, at iteration  $t$ :

- ① Compute score/energy for each encoded features  $s_i^t = a(\dots)$
- ② Normalize these scores through softmax, leading to attention weights:

$$\alpha_i^t = \frac{e^{s_i^t}}{\sum_j e^{s_j^t}}$$

- ③ Compute the context vector:  $c^t = \sum_i \alpha_i^t e^i$
- ④ Make prediction for iteration  $t$ :  $\hat{y}^t = b(c^t)$

## How to compute the scores?

### Content-based attention

The scores are based on features and previous context:  $s_i^t = a(\mathbf{c}^{t-1}, \mathbf{e}^i)$

### Location-based attention

The scores are based on previous attention weights and previous context:

$$s_i^t = a(\mathbf{c}^{t-1}, \boldsymbol{\alpha}^{t-1})$$

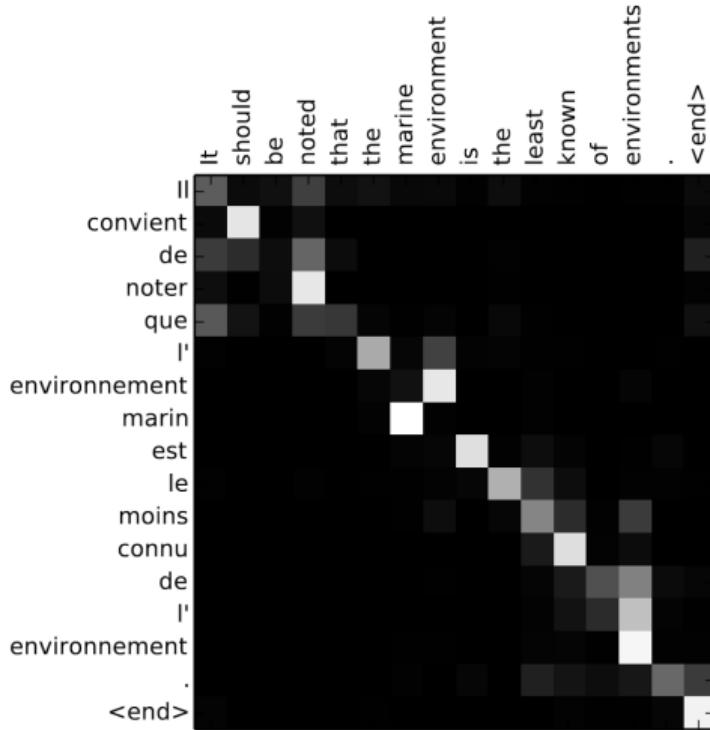
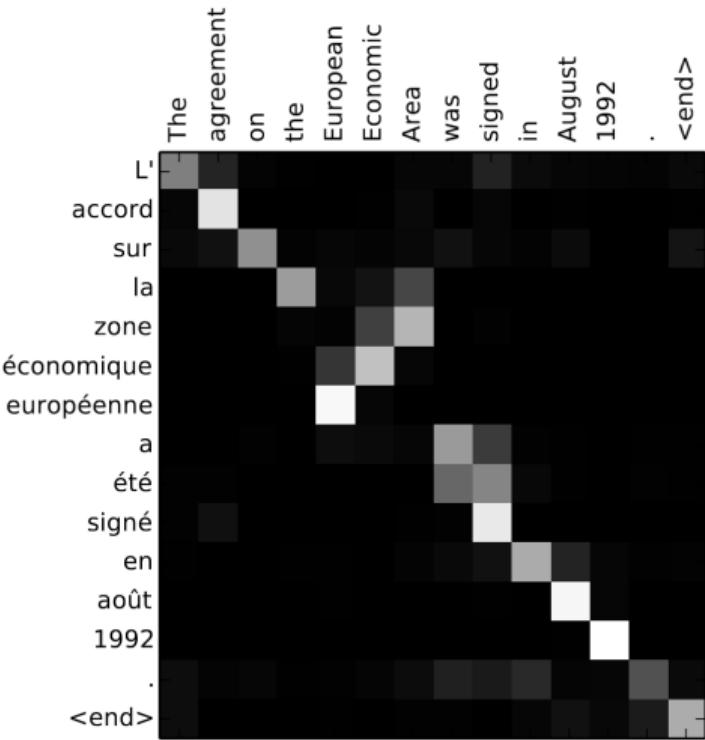
### Hybrid attention

Generally, both approaches are mixed together:  $s_i^t = a(\mathbf{c}^{t-1}, \boldsymbol{\alpha}^{t-1}, \mathbf{e}^i)$

### Example

$$s_i^t = \mathbf{v}^T \tanh(\mathbf{W}_c \mathbf{c}^{t-1} + \mathbf{W}_\alpha \boldsymbol{\alpha}^{t-1} + \mathbf{W}_e \mathbf{e}^i)$$

# Attention maps for encoder-decoder attention



► Size: target sequence length × source sequence length

Images from [1]

### Goal

Embed a token as a weighted sum of the other tokens of the sequence

Input: a set of token embeddings  $\{e^1, \dots, e^{L_x}\}$

Output: a new set of token embeddings  $\{o^1, \dots, o^{L_x}\}$

For each token  $e^i$

- ① Compute a score with respect to each token:  $s_j^i = a(e^i, e^j)$
- ② Normalize the scores:

$$\alpha_j^i = \frac{e^{s_j^i}}{\sum_k e^{s_k^i}}$$

- ③ Compute the context vector:  $c^i = \sum_j \alpha_j^i e^j$
  - ④ Compute the new token representation:  $o^i = b(c^i)$
- All  $o^i$  can be computed in parallel

Let's assume an input sequence  $\mathbf{e} = [\mathbf{e}^1, \dots, \mathbf{e}^{L_x}]$  with  $\mathbf{e}^i \in \mathbb{R}^C$  and  $\mathbf{W} \in \mathbb{R}^{C_o \times C}$ .

- What is the size of  $s^i$ ,  $\alpha^i$ ,  $c^i$  and  $o^i$ ?

$$\mathbf{s}_j^i = \sum_k \mathbf{e}_k^i \mathbf{e}_k^j$$

$$\alpha_j^i = \frac{e^{\mathbf{s}_j^i}}{\sum_k e^{\mathbf{s}_k^i}}$$

$$\mathbf{c}^i = \sum_j \alpha_j^i \mathbf{e}^j$$

$$\mathbf{o}^i = \mathbf{W} \mathbf{c}^i$$

Let's assume an input sequence  $\mathbf{e} = [\mathbf{e}^1, \dots, \mathbf{e}^{L_x}]$   
with  $\mathbf{e}^i \in \mathbb{R}^C$

- There are  $L_x$  tokens
- For a given token  $i$ , one score per token  $j$ :  
 $s^i \in \mathbb{R}^{L_x}$
- For a given token, we want one attention  
weight for each token:  $\alpha^i \in \mathbb{R}^{L_x}$
- Context vector is a weighted sum of vectors of  
size  $C$ :  $\mathbf{c}^i \in \mathbb{R}^C$
- Fully connected layer applied to the context  
vector:  $\mathbf{o}^i \in \mathbb{R}^{C_o}$

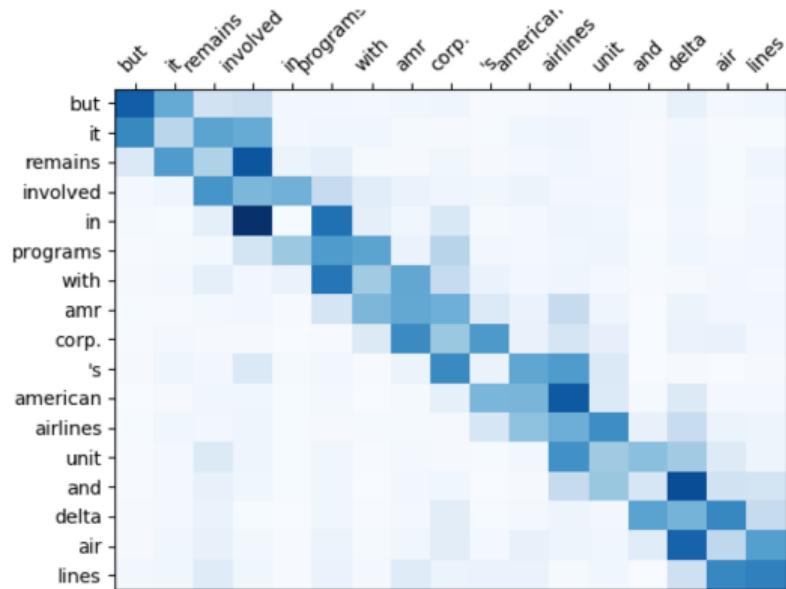
$$\mathbf{s}_j^i = \sum_k \mathbf{e}_k^i \mathbf{e}_k^j$$

$$\alpha_j^i = \frac{e^{\mathbf{s}_j^i}}{\sum_k e^{\mathbf{s}_k^i}}$$

$$\mathbf{c}^i = \sum_j \alpha_j^i \mathbf{e}^j$$

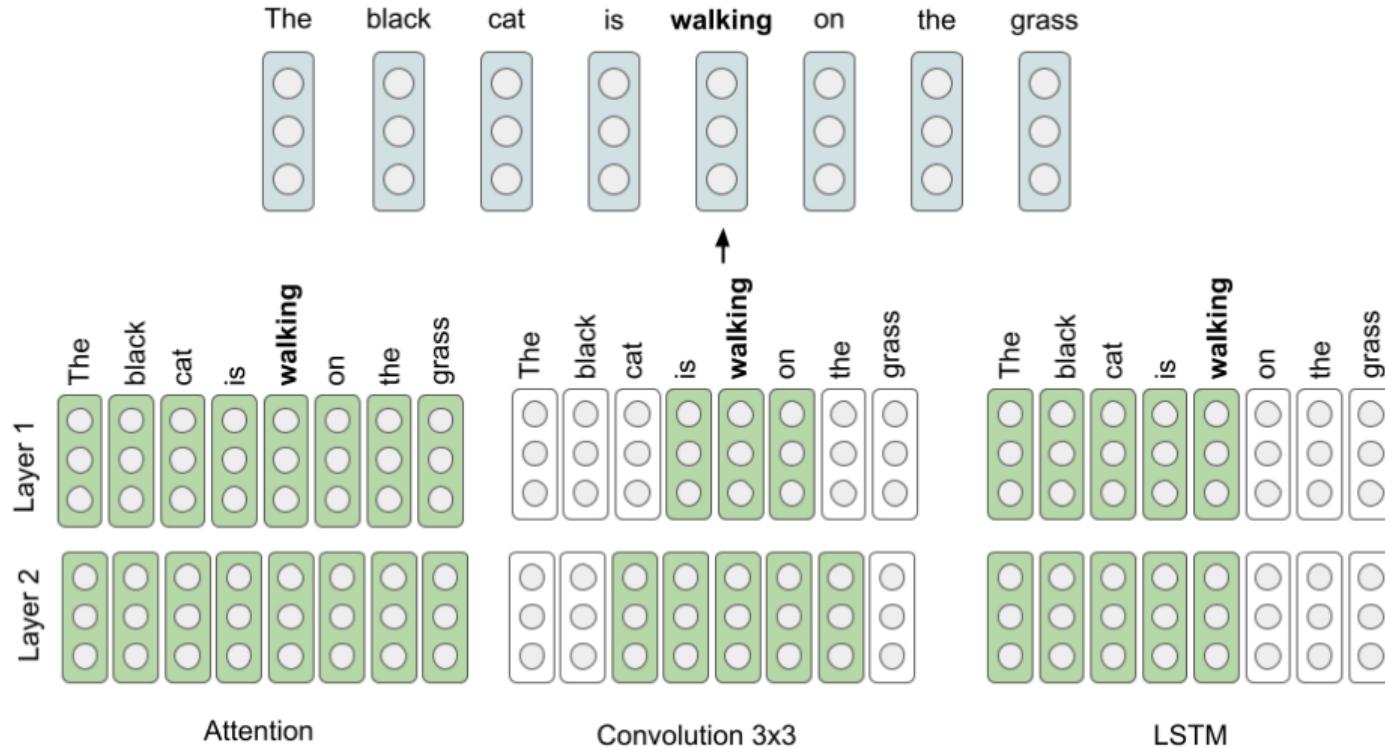
$$\mathbf{o}^i = \mathbf{W} \mathbf{c}^i$$

## Attention map for self-attention



- Size: source sequence length  $\times$  source sequence length

## Context modeling comparison



## 1 Attention mechanisms

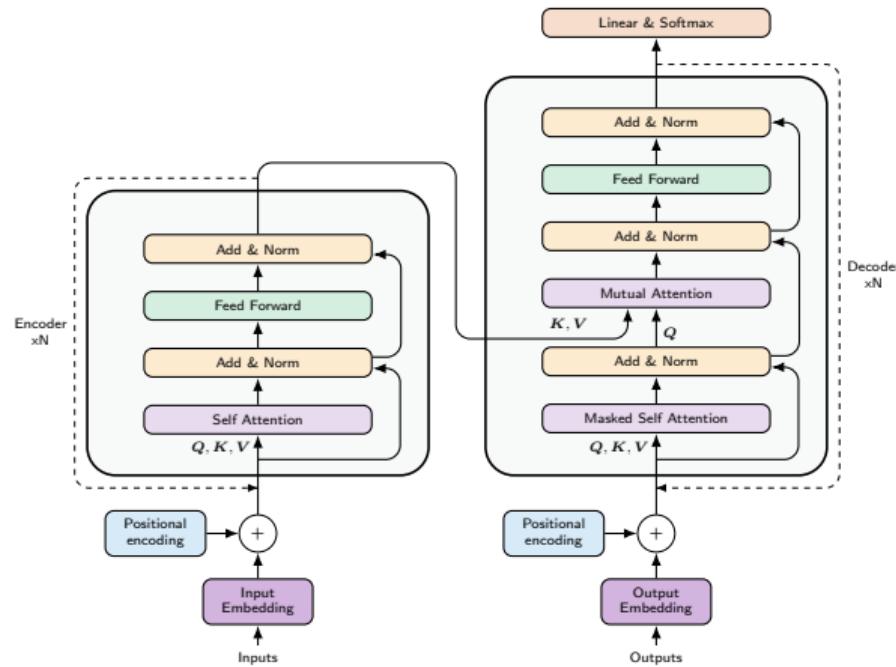
## 2 Transformer

- The Transformer architecture
- Query-key-value paradigm
- The Vision Transformer
- Swin

## 3 Dealing with increasing need for data

## 4 Transformer's everywhere

# Attention is all you need: Transformer (2017) [2]



Two main blocks:

- Transformer encoder
- Transformer decoder

One-shot encoding

$$\mathbf{f} = \text{encoder}(\mathbf{x} + \text{PE})$$

Iterative decoding

$$\mathbf{c}^t = \text{decoder}(\mathbf{f}, \{\hat{\mathbf{y}}^0, \dots, \hat{\mathbf{y}}^{t-1}\} + \text{PE})$$

$$\hat{\mathbf{y}}^t = \arg \max(\mathbf{W} \mathbf{c}^t)$$

► Mainly relies on the Query-Key-Value attention paradigm

## Concept

Goal: retrieve the useful information among a collection of data

Inputs: a *query* and a set of data

## Approach

- Compute a similarity score between the query and all the data
  - Generate an answer based on the most relevant data
- Two representations of the data are used: *keys* for comparison with the query, and *values* for the answer

## Analogy with library search engine

Query: keywords; Keys: book titles; Values: book contents

## Issue

Gradient propagation for hard attention (select only the most relevant item)

► Soft attention (weighted sum of items)

## Scaled Dot-Product Attention

$q \in \mathbb{R}^d$ : a query vector

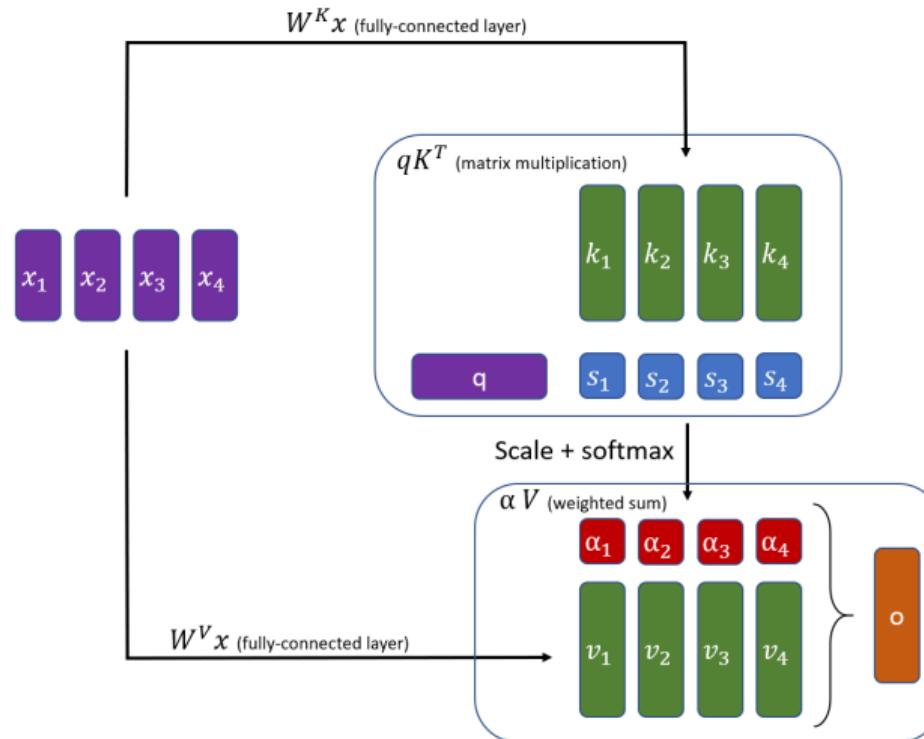
$K \in \mathbb{R}^{L \times d}$ : a matrix gathering key vectors for each available item

$V \in \mathbb{R}^{L \times d}$ : a matrix gathering values vectors for each available item

$$\text{Attention}(q, K, V) = \underbrace{\text{softmax}\left(\frac{qK^T}{\sqrt{d}}\right)}_{\alpha} V = c \quad (\in \mathbb{R}^d)$$

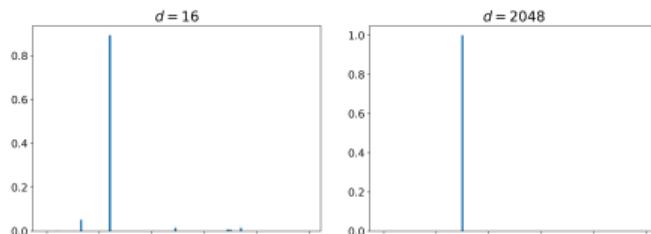
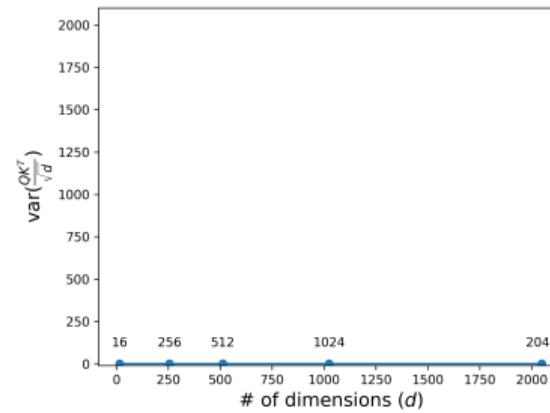
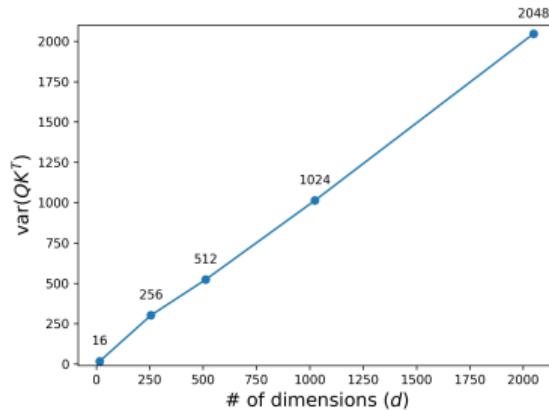
► Dot product as distance function

# Illustration of the computations

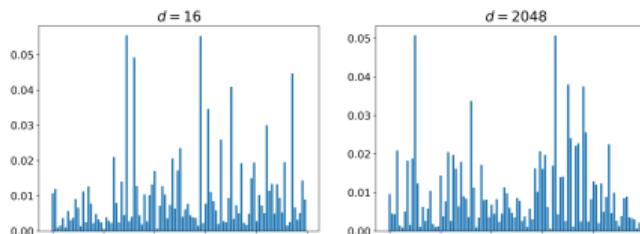


Why scaling by  $\frac{1}{\sqrt{d}}$ ?

For 1 query and 10,000 keys (from normal distribution):



Softmax over the first 100 keys



Softmax over the first 100 keys

- Queries, Keys and Values are from the same source

Query-key-value for feature extraction

$\mathbf{X}$ : the latent representation of an input sequence.

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^Q (\in \mathbb{R}^{L_x \times d})$$

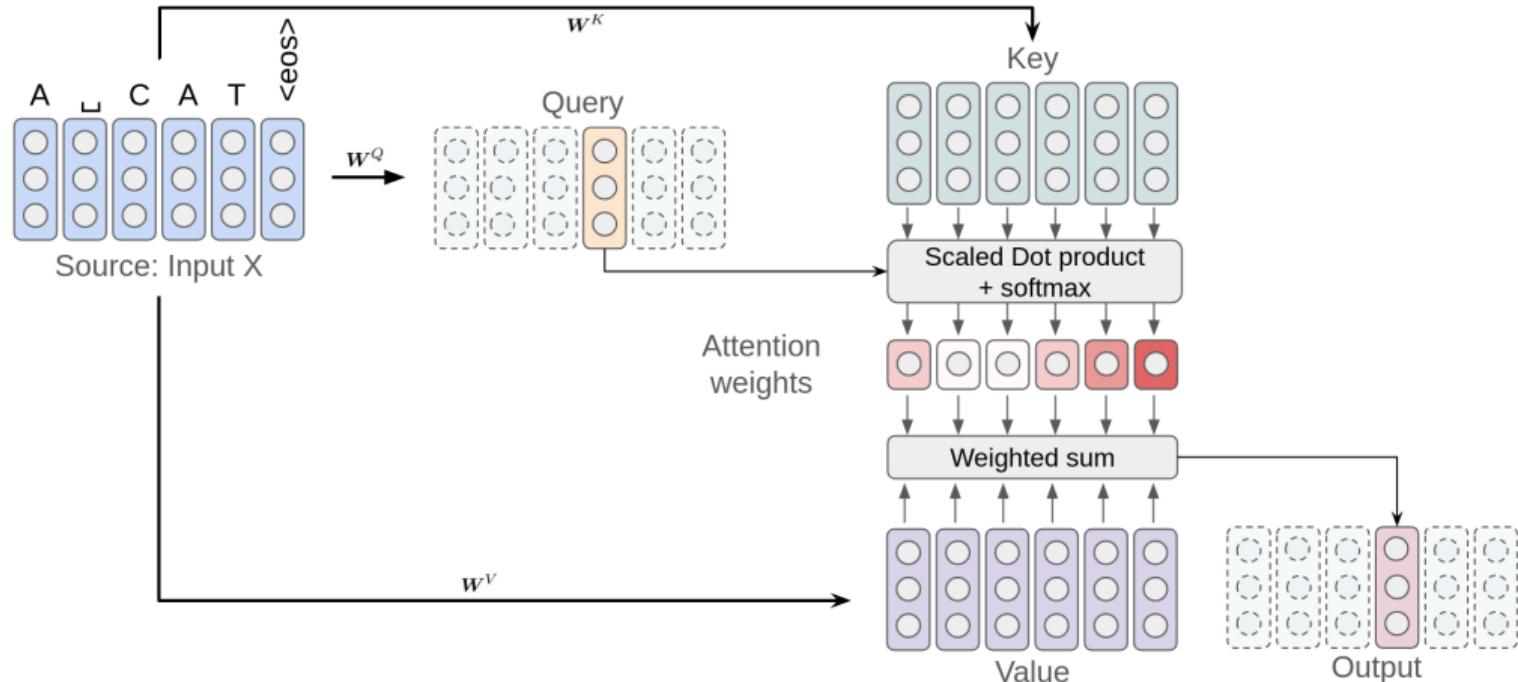
$$\mathbf{K} = \mathbf{X}\mathbf{W}^K (\in \mathbb{R}^{L_x \times d})$$

$$\mathbf{V} = \mathbf{X}\mathbf{W}^V (\in \mathbb{R}^{L_x \times d})$$

$$\mathbf{Y} = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} \right) \mathbf{V}$$

$\mathbf{W}^Q$ ,  $\mathbf{W}^K$  and  $\mathbf{W}^V$  are matrices of trainable weights (fully-connected layers)

## Self-attention



- Output can be computed in parallel through matrix multiplications
- Output length = Query length (= Key length = Value length)

## Multi-head self-attention (MSA)

$h$  self-attention in parallel

$$\mathbf{Q}_1 = \mathbf{X} \mathbf{W}_1^Q \quad \dots \quad \mathbf{Q}_h = \mathbf{X} \mathbf{W}_h^Q$$

$$\mathbf{K}_1 = \mathbf{X} \mathbf{W}_1^K \quad \dots \quad \mathbf{K}_h = \mathbf{X} \mathbf{W}_h^K$$

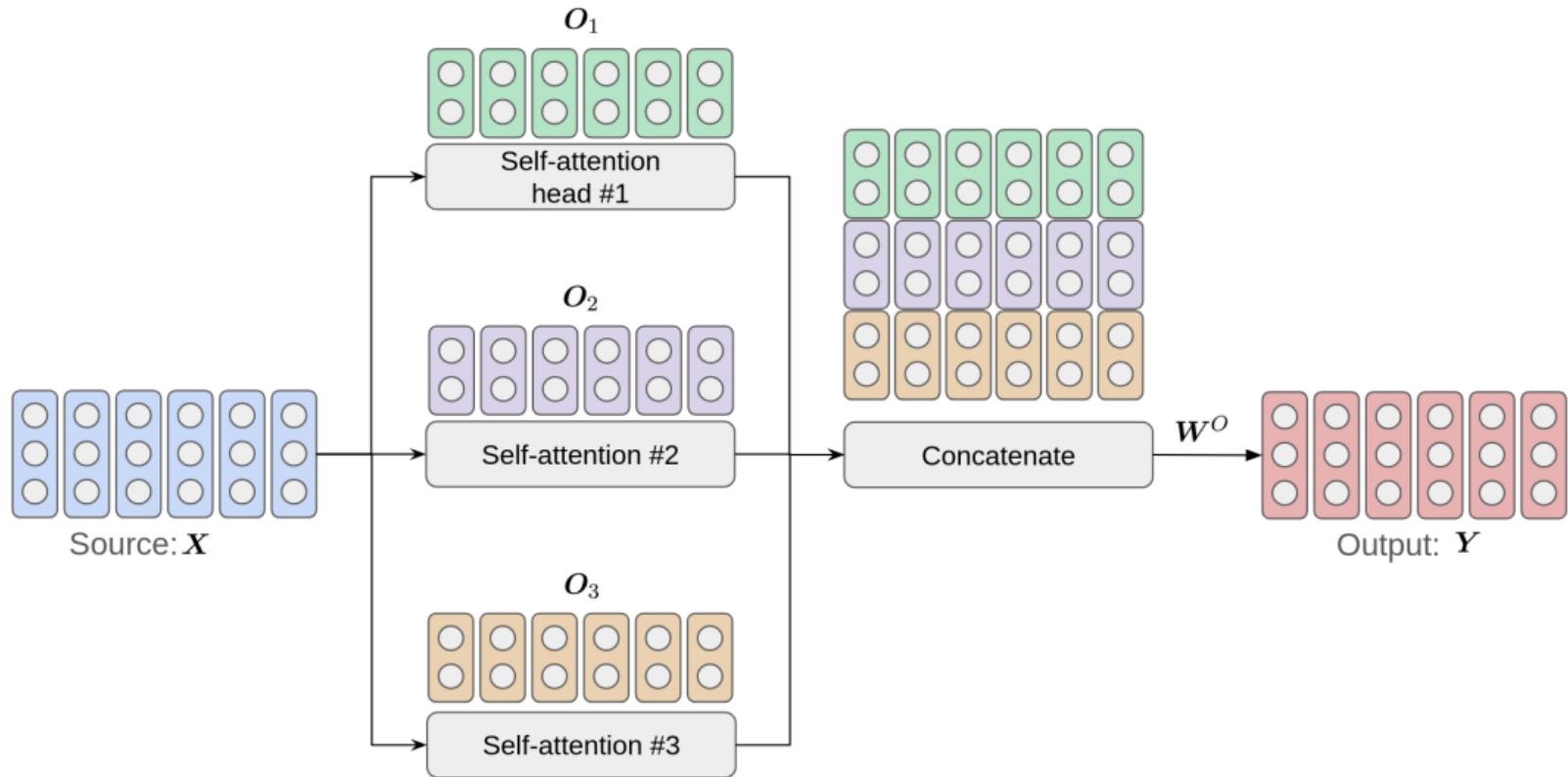
$$\mathbf{V}_1 = \mathbf{X} \mathbf{W}_1^V \quad \dots \quad \mathbf{V}_h = \mathbf{X} \mathbf{W}_h^V$$

$$\mathbf{O}_1 = \text{softmax} \left( \frac{\mathbf{Q}_1 \mathbf{K}_1^T}{\sqrt{d}} \right) \mathbf{V}_1 \quad \dots \quad \mathbf{O}_h = \text{softmax} \left( \frac{\mathbf{Q}_h \mathbf{K}_h^T}{\sqrt{d}} \right) \mathbf{V}_h$$

$$\mathbf{Y} = \text{concat}(\mathbf{O}_1, \dots, \mathbf{O}_h) \mathbf{W}^O$$

Heads can be seen as the equivalent of convolutional kernels

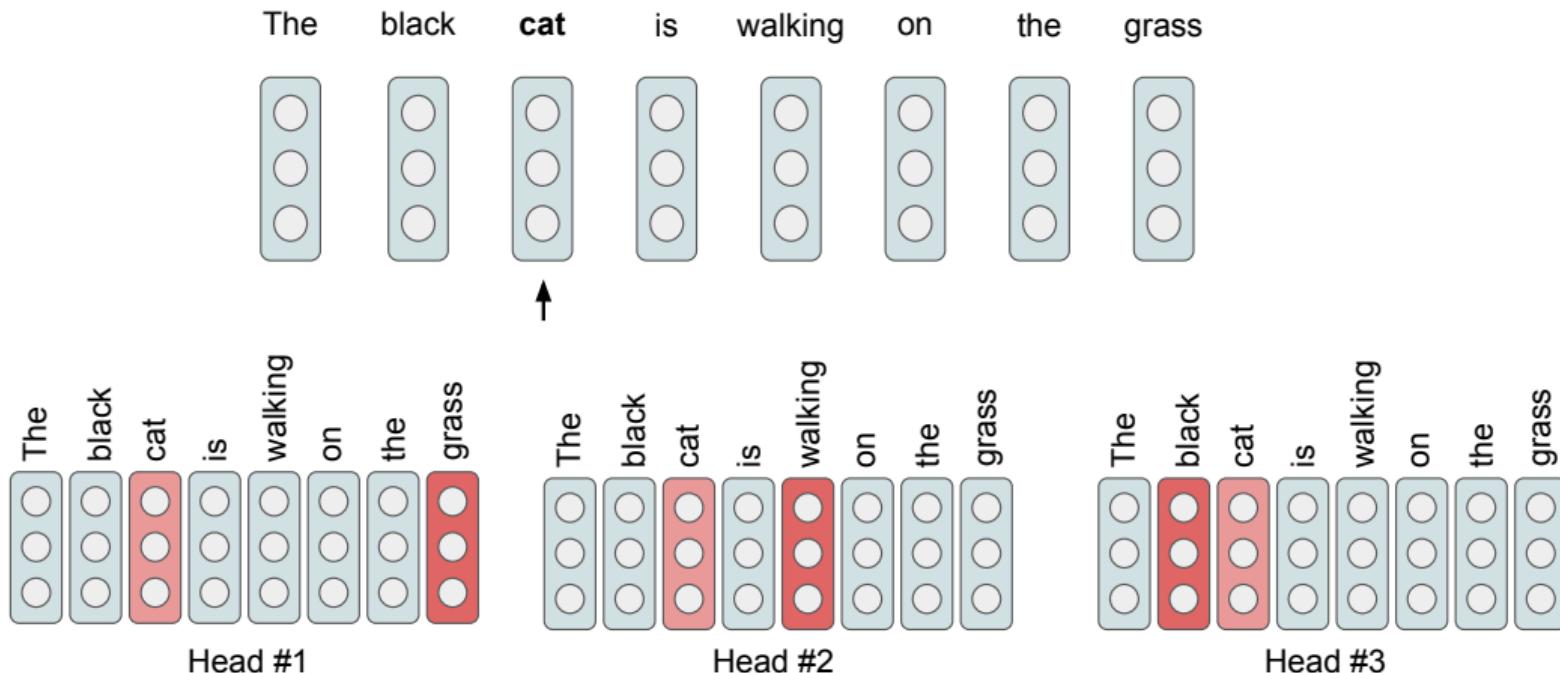
# Multi-head Self-Attention (MSA)



# Multi-head Self-Attention (MSA)

## Intuition

Each head is specialized for a specific query (e.g., action, location, description)



### Issue with Multi-head Self Attention (MSA)

Self-attention is equivariant to input permutation:

$$\text{MSA}([\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3]) = [\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3] \Rightarrow \text{MSA}([\mathbf{x}_2, \mathbf{x}_1, \mathbf{x}_3]) = [\mathbf{y}_2, \mathbf{y}_1, \mathbf{y}_3]$$

The position information is lost whereas it is crucial

- A word/sentence is a sequence of characters, not a set of characters

### Goal

Inject positional information to preserve the sequentiality

### Solution

Additive positional encoding:

- Input of transformer encoder:  $\mathbf{I} = \mathbf{X} + \text{PE}$

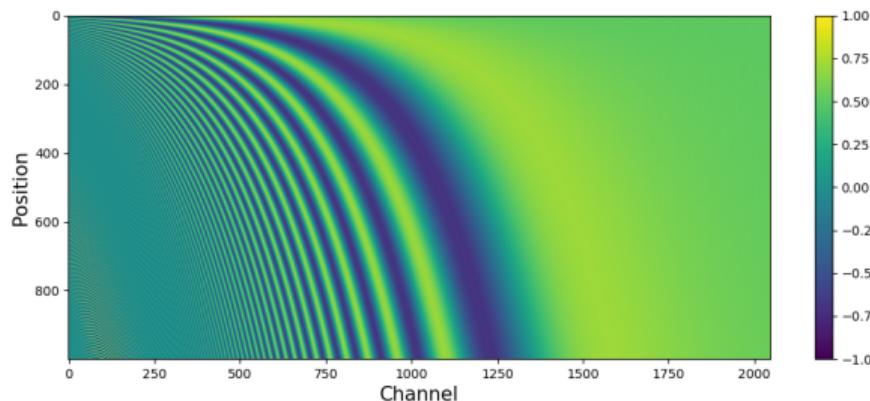
A combination of sines and cosines

$\text{PE}(p, k)$ :  $k^{\text{th}}$  value of positional embedding vector for position  $p$ .

$$\text{PE}(p, 2k) = \sin(w_k \cdot p)$$

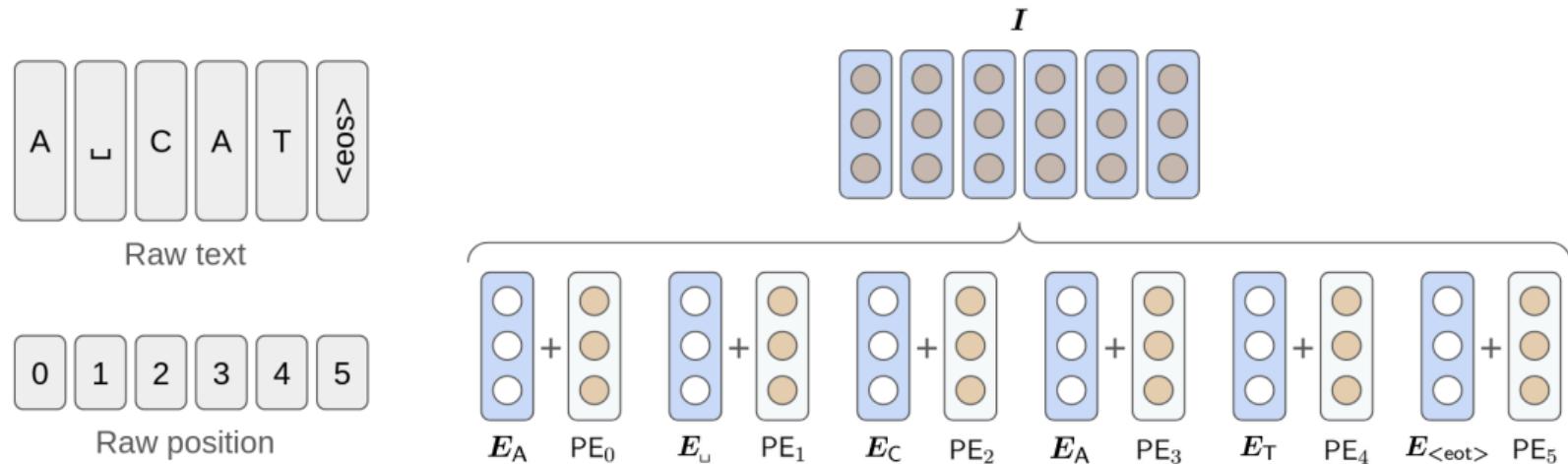
$$\text{PE}(p, 2k + 1) = \cos(w_k \cdot p)$$

$\forall k \in [0, d/2]$ , with  $w_k = 1/10000^{2k/d}$

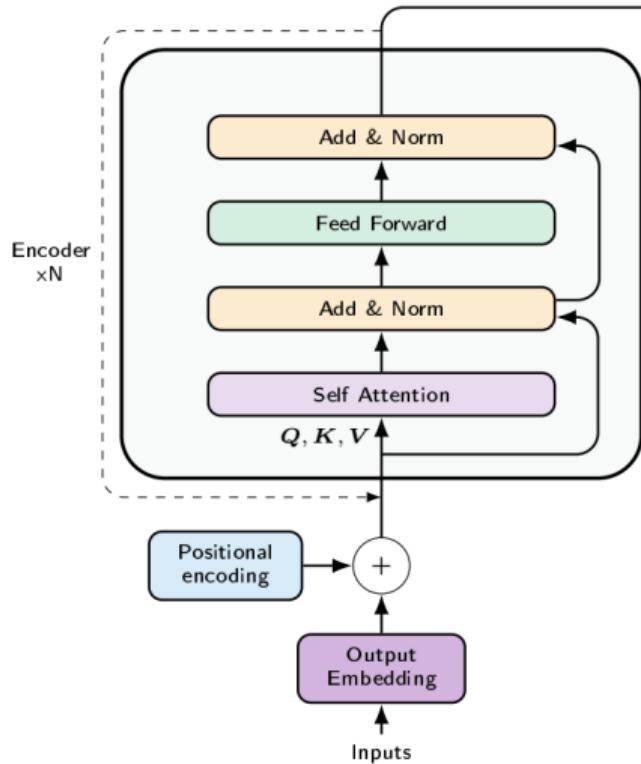


## Advantages

- Does not depend on the input length
- Encoding unique for each position
- Periodicity of sine/cosine could help to adapt to unseen position
- Bonus: no additional trainable parameters



## Transformer encoder



A stack of 6 transformer encoder layers

- Multi-head Self Attention
  - Global context modeling
- Add: residual connections
  - Multi-scale representation
  - Reinforce identity
- Norm: layer normalization
  - Stabilize training (range of values)
- Feed Forward: 2 fully-connected layers
  - Local projection

Input sequence of  $L_x$  token embeddings:  $\mathbf{X} \in \mathbb{R}^{L \times d}$

A simplified version of the attention layer, keeping the same dimension:

$$\mathbf{Y} = (\mathbf{XW}^Q)(\mathbf{XW}^K)^T(\mathbf{XW}^V)\mathbf{W}^O$$

$\mathbf{W}^Q \in \mathbb{R}^{d \times d}$ ,  $\mathbf{W}^K \in \mathbb{R}^{d \times d}$ ,  $\mathbf{W}^V \in \mathbb{R}^{d \times d}$ ,  $\mathbf{W}^O \in \mathbb{R}^{d \times d} \rightarrow 4d^2$  parameters

$\mathbf{XW}^i$ : from  $L \times d$  to  $L \times d$ ;  $2d^2L$  FLOPs ( $\times 3$ )

$\mathbf{S} = (\mathbf{XW}^Q)(\mathbf{XW}^K)^T$ : output  $L \times L$ ;  $2dL^2$  FLOPs

$\mathbf{O} = \mathbf{S}(\mathbf{XW}^V)$ : output  $L \times d$ ;  $2dL^2$  FLOPs

$\mathbf{Y} = \mathbf{OW}^O$ : output  $L \times d$ ;  $2d^2L$  FLOPs

►  $8d^2L + 4dL^2$  FLOPs

► Quadratic complexity with respect to the input length

## Cross attention (also known as mutual attention)

- Queries from target domain, and keys and values from source domain

### General case

Inputs:  $\mathbf{X}$ : the features from the encoded input sequence (source domain),  
 $\tilde{\mathbf{z}}$ : the query sequence (e.g., a question)  
 $\mathbf{z}$ : the query vector ( $\mathbf{z} = g(\tilde{\mathbf{z}})$ ).

$$\mathbf{q} = \mathbf{z}\mathbf{W}^Q (\in \mathbb{R}^d)$$

$$\mathbf{K} = \mathbf{X}\mathbf{W}^K (\in \mathbb{R}^{L_x \times d})$$

$$\mathbf{V} = \mathbf{X}\mathbf{W}^V (\in \mathbb{R}^{L_x \times d})$$

$$\mathbf{c} = \text{Attention}(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \underbrace{\text{softmax} \left( \frac{\mathbf{q}\mathbf{K}^T}{\sqrt{d}} \right)}_{\alpha} \mathbf{V}$$

Prediction:  $\hat{\mathbf{y}} = h(\mathbf{c})$ , with  $g$  and  $h$  parametric functions

## Cross attention (also known as mutual attention)

- Queries from target domain, and keys and values from source domain

### The iterative decoding process case

Inputs:  $\mathbf{X}$ : the features from the encoded input sequence (source domain),  
 $\tilde{\mathbf{z}} = \hat{\mathbf{y}}^{0:t-1}$ : an output sequence (target domain)  
 $\mathbf{z}^t$ : a query vector at iteration  $t$  ( $\mathbf{z}^t = g(\hat{\mathbf{y}}^0, \dots, \hat{\mathbf{y}}^{t-1})$ ).

$$\mathbf{q}^t = \mathbf{z}^t \mathbf{W}^Q (\in \mathbb{R}^d)$$

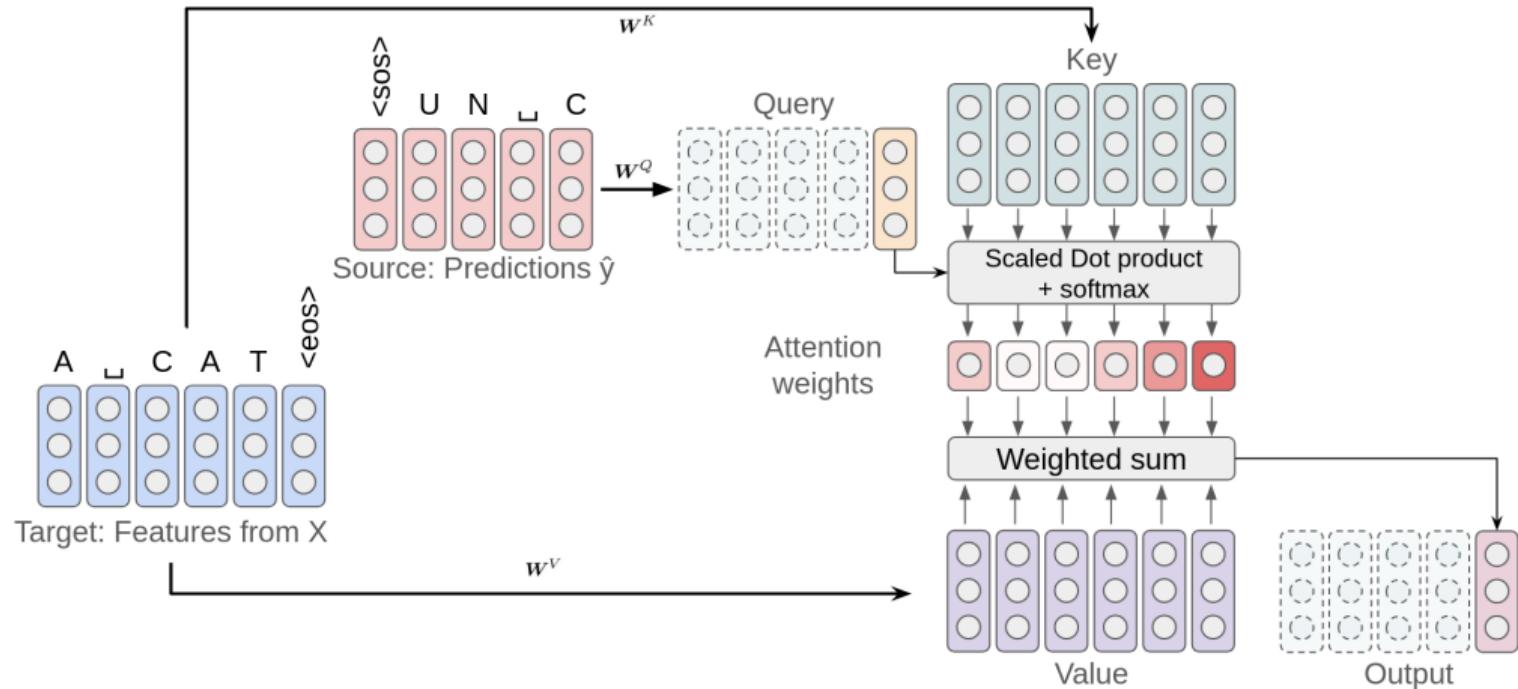
$$\mathbf{K} = \mathbf{X} \mathbf{W}^K (\in \mathbb{R}^{L_x \times d})$$

$$\mathbf{V} = \mathbf{X} \mathbf{W}^V (\in \mathbb{R}^{L_x \times d})$$

$$\mathbf{c}^t = \text{Attention}(\mathbf{q}^t, \mathbf{K}, \mathbf{V}) = \underbrace{\text{softmax}\left(\frac{\mathbf{q}^t \mathbf{K}^T}{\sqrt{d}}\right)}_{\alpha^t} \mathbf{V}$$

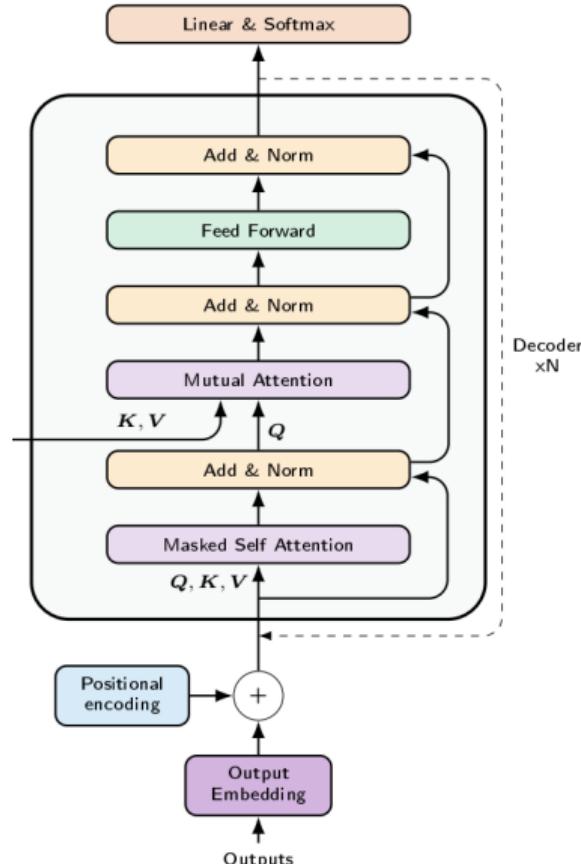
New prediction:  $\hat{\mathbf{y}}^t = h(\mathbf{c}^t)$ , with  $g$  and  $h$  parametric functions

## Cross attention attention



► Output length = Query length ( $\neq$  Key length = Value length)

# Transformer decoder

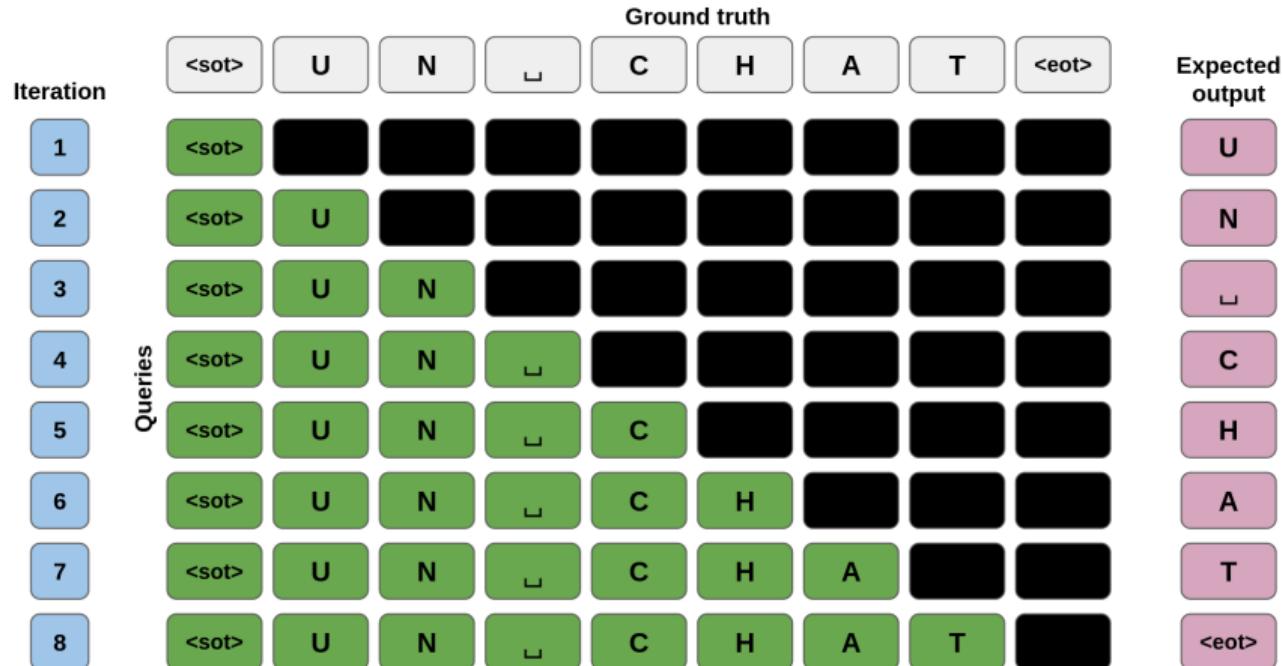


A stack of 6 transformer decoder layers

- Multi-head Masked Self Attention
  - Global query modeling
- Add: residual connections
  - Multi-scale representation
  - Reinforce identity
- Norm: layer normalization
  - Stabilize training (range of values)
- Multi-head mutual attention
  - Look for relevant information from source
- Feed Forward: 2 fully-connected layers
  - Local projection

## Masked self attention and teacher forcing

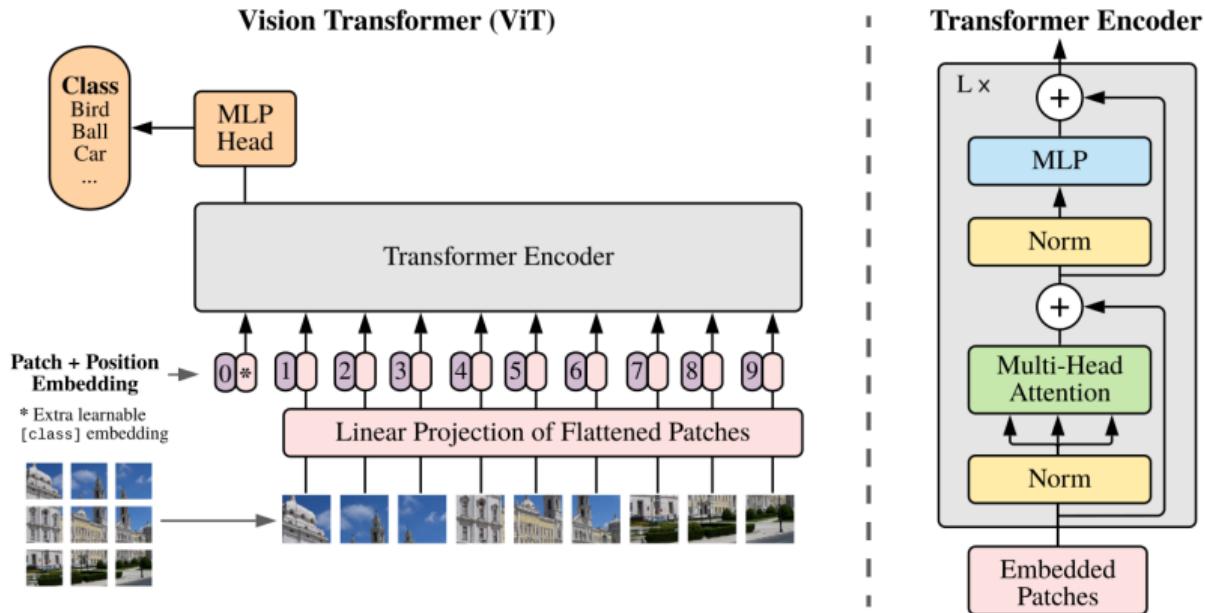
- Parallelize decoding process at training time using ground truth



## Conclusion

- A modular approach: expressivity can be enhanced by adding more encoder/decoder layers
  - But slower training/inference and more required data
- Self attention: a way to model global context without sequential computations (at each layer, each token can attend to all the other tokens)
  - But a quadratic complexity
- Parameter-free additive positional encoding
  - But the relation between positions must be learned through training
- Defined for Natural Language Processing: SOTA architecture
  - How to adapt it for computer vision?

# Transformer for image classification: Vision Transformer (2021) [3]

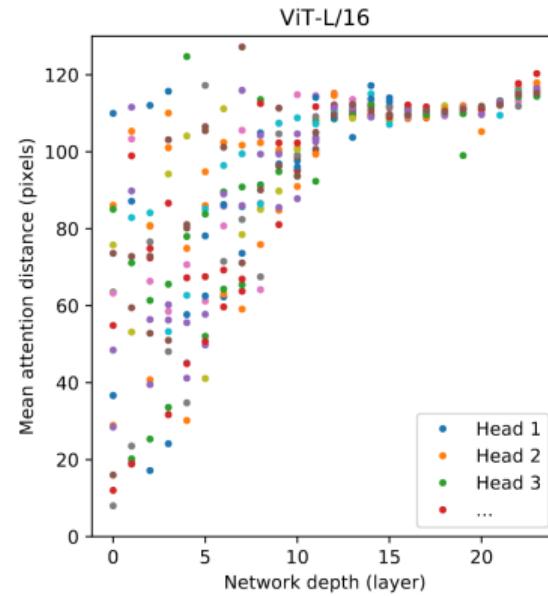


Input image:  $224 \times 224$ , patch size:  $16 \times 16 \Rightarrow 196$  patches  
Introduction of a trainable class token used for prediction

Input      Attention



Attention map for the class token



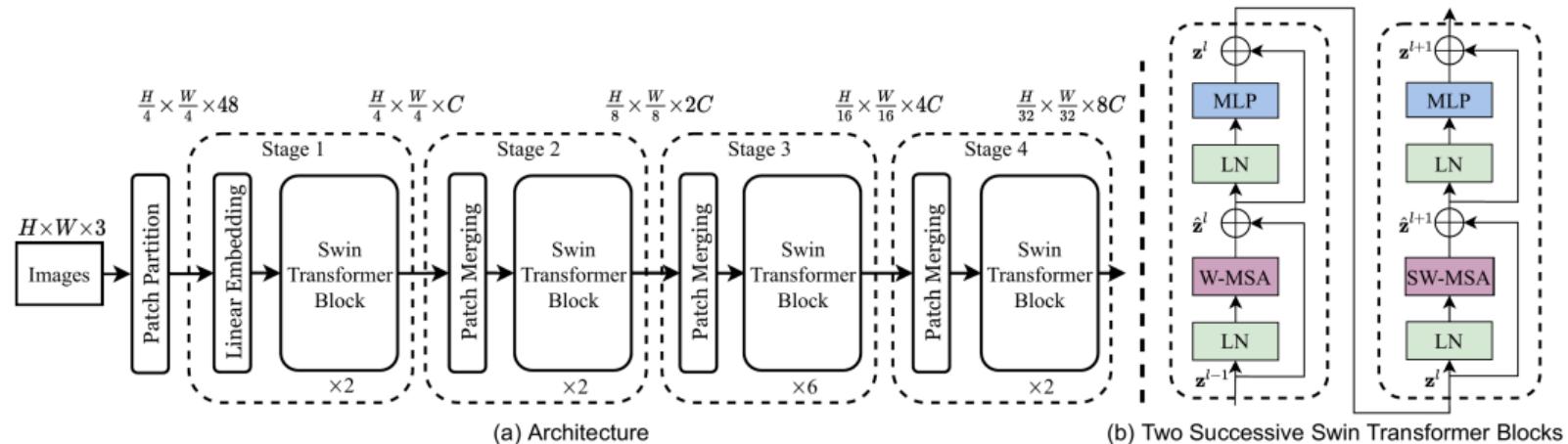
Architecture	Top-1 accuracy (%) on ImageNet	# params (M)
ResNet-50	79.26	26
ResNet-101	80.13	45
ResNet-152	<b>80.62</b>	60
ViT-B/16	77.91	86
ViT-B/32	73.38	86
ViT-L/16	76.53	307
ViT-L/32	71.16	307

Attention-based architectures require more data to exploit their full potential:  
► Spatial relations between patches must be learned from scratch

Performance on ImageNet (fine-tuning included)

Architecture	Top-1 accuracy (%) on ImageNet	
	Pre-trained on ImageNet	Pre-trained on JFT-300M
ResNet-50	79.26	✗
ResNet-101	80.13	✗
ResNet-152	80.62	✗
ViT-B/16	77.91	84.15
ViT-B/32	73.38	80.73
ViT-L/16	76.53	87.12
ViT-L/32	71.16	84.37
ViT-H/14	✗	<b>88.04</b>

## ► Hierarchical Vision Transformer using Shifted windows (Swin)

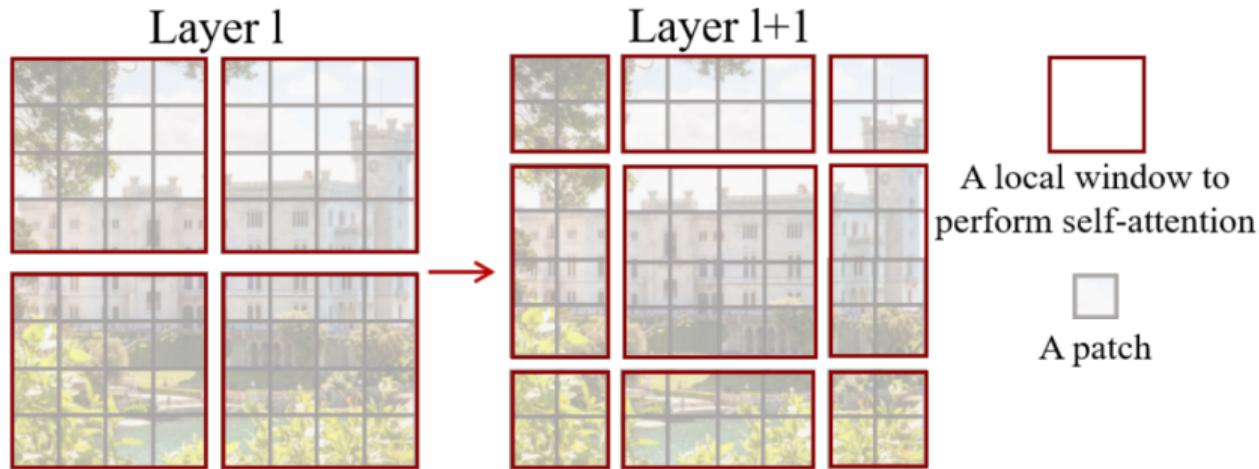


Patch size:  $4 \times 4$

Patch merging:  $2 \times 2$  adjacent patches are merged together

(S)W-MSA: (Shifted) Window Multi Self-Attention

► Idea: perform attention on patch windows (locally), in parallel



### Modeling global context with shifted windows

Windows are shifted from one layer to another to propagate the information from one window to its neighbours

Reminder: simplified self-attention

Number of FLOPs:  $8d^2L + 4dL^2$

With  $d$  the number of dimensions and  $L$  the sequence length

Compute the number of FLOPs for a latent representation of an image of size ( $H = 1024, W = 1024, C = 256$ )

- for a traditional self-attention layer
  - for the shifted windows approach (window size:  $8 \times 8$ , first step only)
- Remark: the number of dimensions is preserved (from 256 to 256)

Reminder: simplified self-attention

Number of FLOPs:  $8d^2L + 4dL^2$

With  $d$  the number of dimensions and  $L$  the sequence length

For a traditional self-attention layer

Flatten:  $L = HW = 1024 \times 1024 = 1,048,576$

►  $8 \times 256^2 \times 1,048,576 + 4 \times 256 \times 1,048,576^2 = 1.1 \times 10^{15}$  FLOPs

For a shifted-window self-attention layer

Window size:  $8 \times 8$

Flatten:  $L = HW = 8 \times 8 = 64$

►  $8 \times 256^2 \times 64 + 4 \times 256 \times 64^2 = 37.7 \times 10^6$  FLOPs per window

Number of windows:  $\frac{1024}{8} \times \frac{1024}{8} = 16,384$

►  $16,384 \times 37.7 \times 10^6 = 0.6 \times 10^9$  FLOPs

method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
R-101x3 [38]	$384^2$	388M	204.6G	-	84.4
R-152x4 [38]	$480^2$	937M	840.5G	-	85.4
ViT-B/16 [20]	$384^2$	86M	55.4G	85.9	84.0
ViT-L/16 [20]	$384^2$	307M	190.7G	27.3	85.2
Swin-B	$224^2$	88M	15.4G	278.1	85.2
Swin-B	$384^2$	88M	47.0G	84.7	86.4
Swin-L	$384^2$	197M	103.9G	42.1	87.3

Pre-training on ImageNet-22k

- 1 Attention mechanisms
- 2 Transformer
- 3 Dealing with increasing need for data
- 4 Transformer's everywhere

### Create artificial examples

- Data augmentation
- Synthetic data

### Benefit from other annotated datasets

- Transfer learning
- Fine-tuning
- Multi-tasking

### Benefit from unlabeled data

- Semi-supervised learning
- **Self-supervised learning**

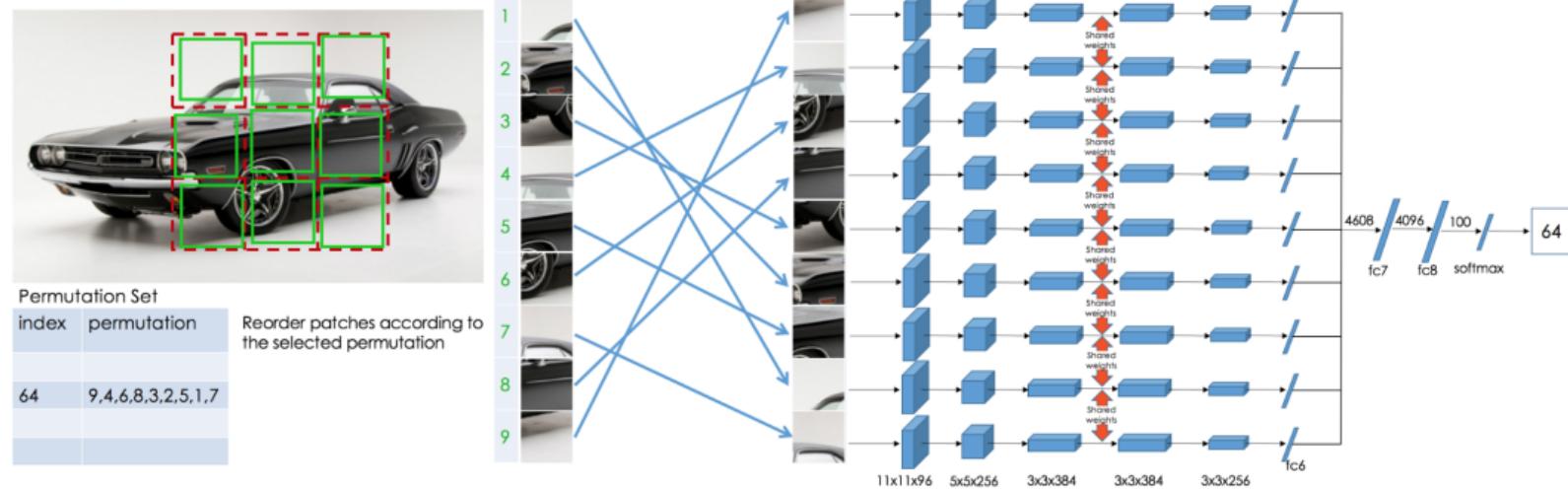
The goal is to solve a pretext task

Some examples:

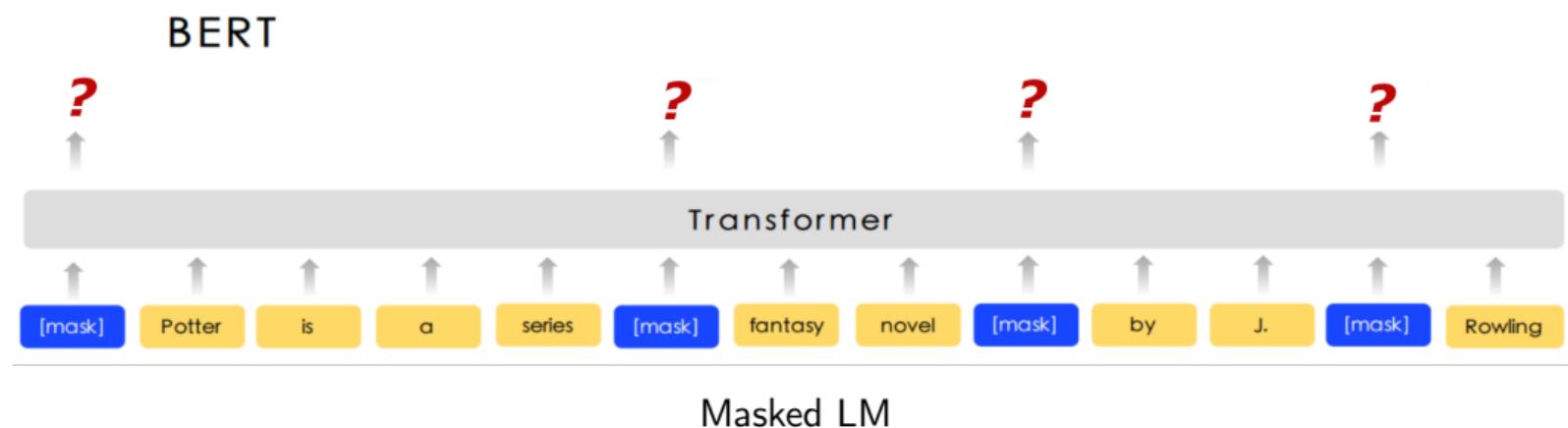
- Solve Jigsaw puzzles
- Reconstruct the input image:  
(Variational) Auto-encoders, Masked auto-encoders
- Contrastive learning: SimCLR, CLIP  
Learning from positive and negative samples

# Jigsaw puzzle (2016) [5]

► Find the right permutation: classification formulation

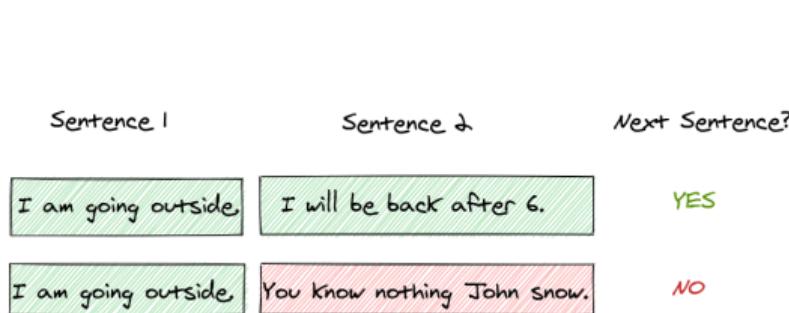
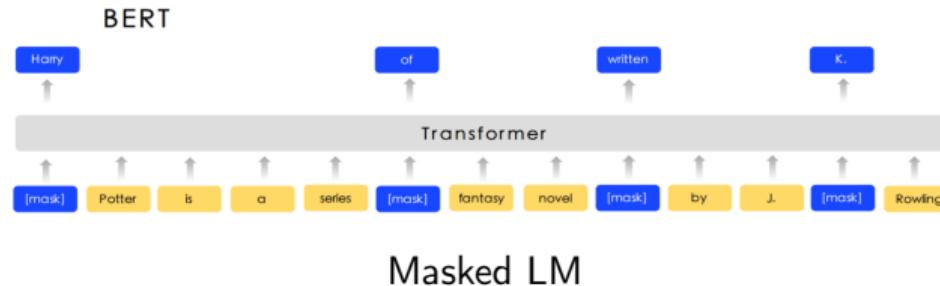


- Two self-supervised strategies

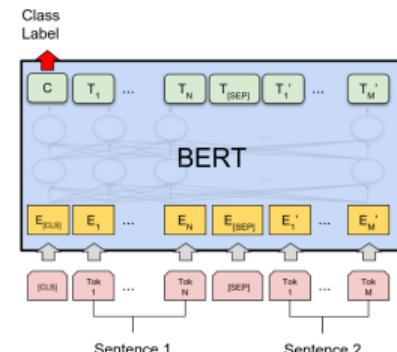


# BERT pre-training (2019) [6]

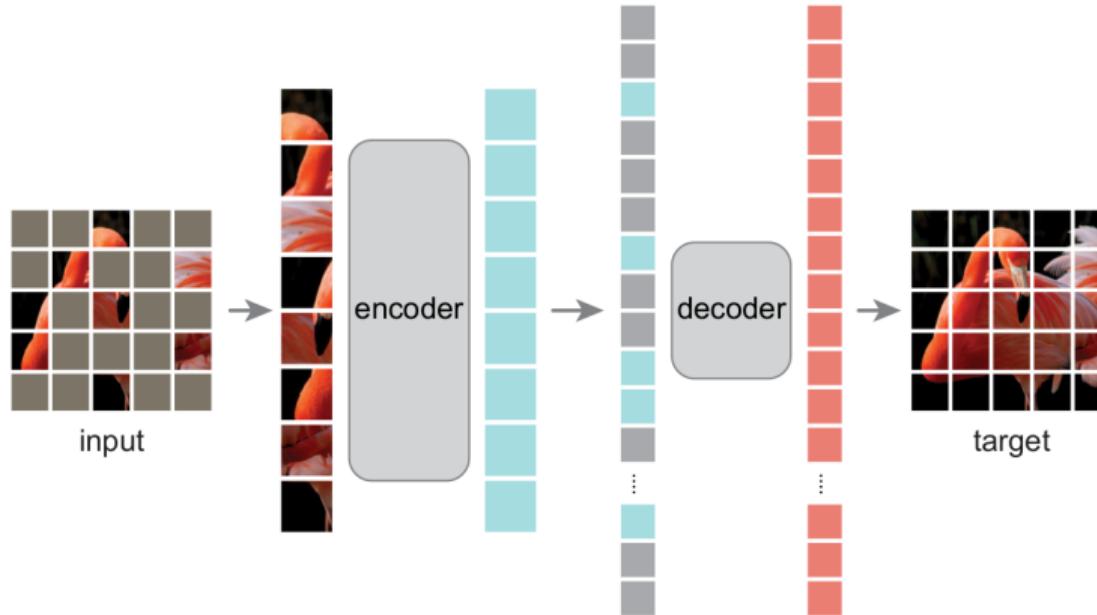
## ► Two self-supervised strategies



Next Sentence Prediction



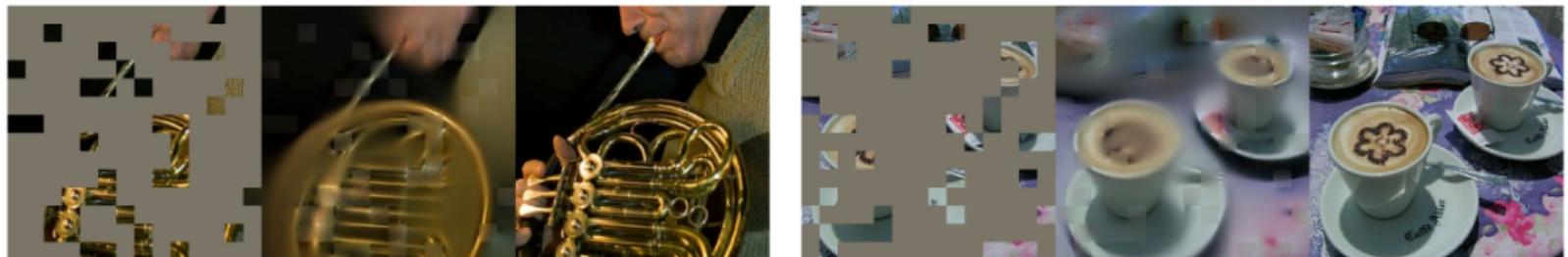
- Reconstruct the missing parts (pixel-level MSE loss)



- Vision Transformer architecture

- Classification performance (Top-1 accuracy) with ViT architectures, pre-trained and evaluated on ImageNet

Approach	ViT-B	ViT-L	ViT-H	ViT-H <sub>488</sub>
Supervised	82.3	82.6	83.1	X
MAE	83.6	85.9	86.9	<b>87.8</b>



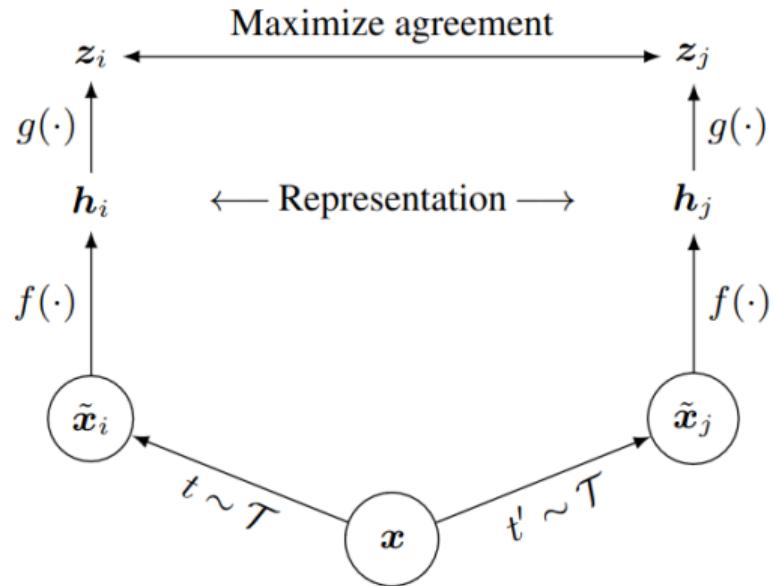
Left: masked input, Middle: reconstructed image, Right: ground truth

## Contrastive learning: SimCLR (2020) [8]

- Idea: generate two views for each input (= positive pair). Bring latent representations of positives closer and keep latent representations of negatives away
- Siamese networks: same architecture with same weights used on different inputs while comparing outputs

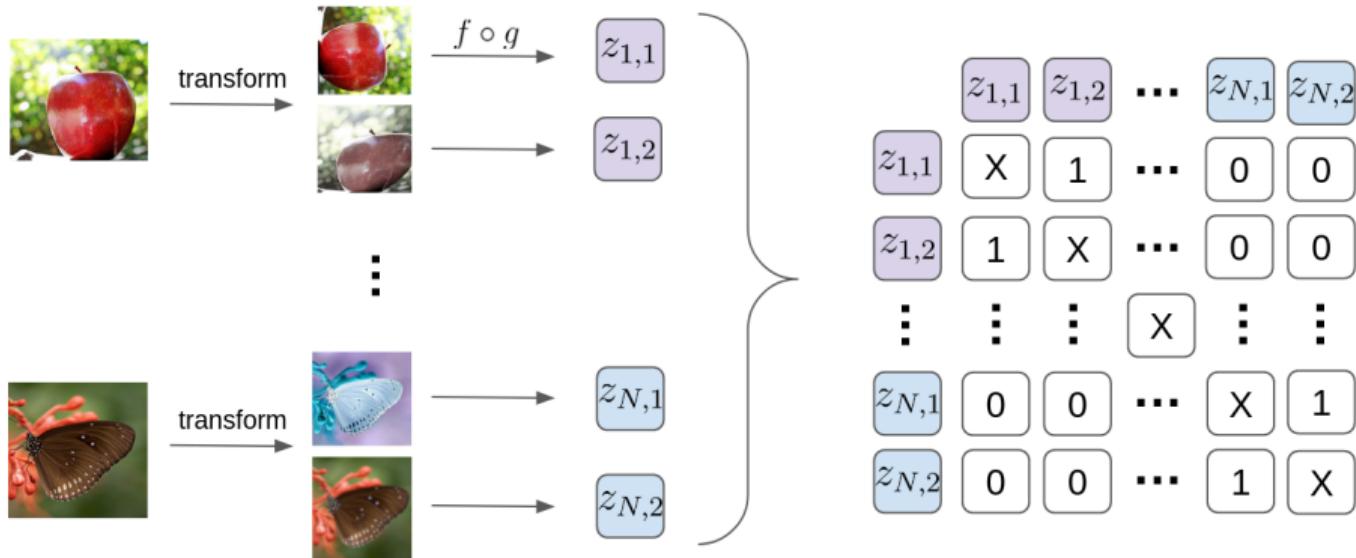
**Algorithm 1** SimCLR's main learning algorithm.

```
input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .  
for sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  do  
  for all  $k \in \{1, \dots, N\}$  do  
    draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$   
    # the first augmentation  
     $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$   
     $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation  
     $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection  
    # the second augmentation  
     $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$   
     $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation  
     $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection  
  end for  
  for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do  
     $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity  
  end for  
  define  $\ell(i, j)$  as  $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$   
   $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$   
  update networks  $f$  and  $g$  to minimize  $\mathcal{L}$   
end for  
return encoder network  $f(\cdot)$ , and throw away  $g(\cdot)$ 
```



## Contrastive learning: SimCLR (2020) [8]

- Find the matching pairs
- Cosine similarity as distance metric



- Efficiency depends on mini-batch size

## ► Contrastive Language-Image Pre-training (CLIP)

## Idea

Jointly train two encoders to avoid human annotation effort:

- An image encoder  $f$ : ResNet/ViT
- A text encoder  $g$ : transformer

► Both encoders are trained to generate a fixed-length latent representation from text/image input sharing the same feature space

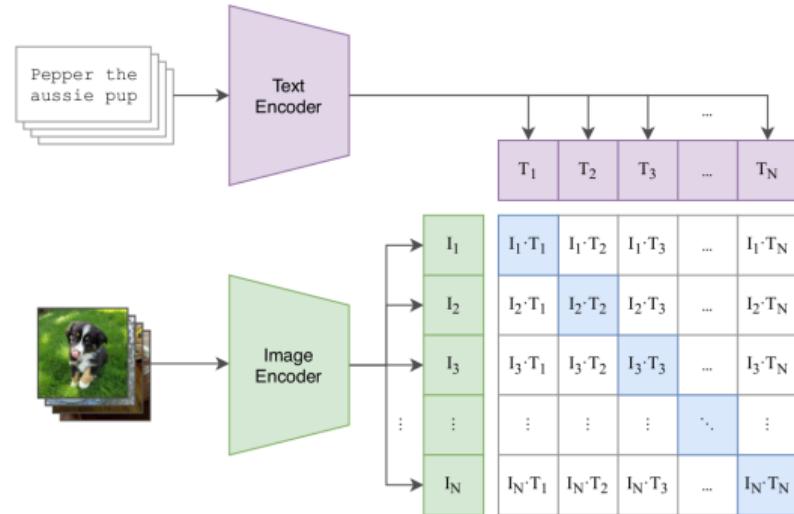
## How?

Contrastive learning between images and captions

► 400 million pairs (image, text) collected from the web

$x$ : image  
 $y$ : caption

Image encoder:  $f(x) = I$   
 Text encoder:  $g(y) = T$



## Symmetric loss

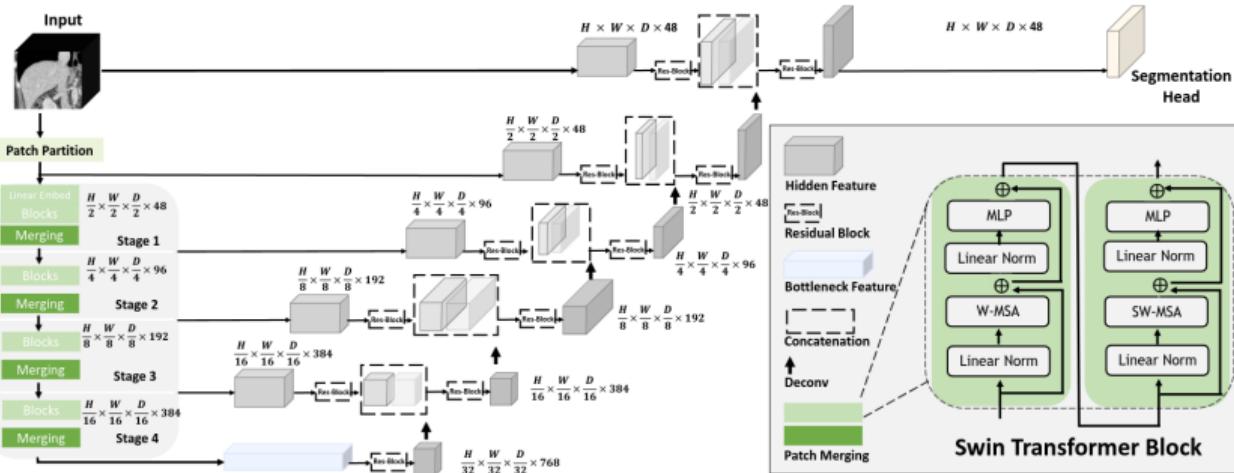
$$\mathcal{L} = \frac{1}{2} \left( \underbrace{\sum_{i=1}^N \mathcal{L}_{\text{CE}}(\hat{y}_i^I, y_i^I)}_{\text{image } i \text{ compared to all texts}} + \underbrace{\sum_{t=1}^N \mathcal{L}_{\text{CE}}(\hat{y}_t^T, y_t^T)}_{\text{text } t \text{ compared to all images}} \right)$$

$$\hat{y}_{i,j}^I = I_i \cdot T_j$$

$$\hat{y}_{t,j}^T = I_j \cdot T_t$$

- 1 Attention mechanisms
- 2 Transformer
- 3 Dealing with increasing need for data
- 4 Transformer's everywhere

## ► Medical image segmentation

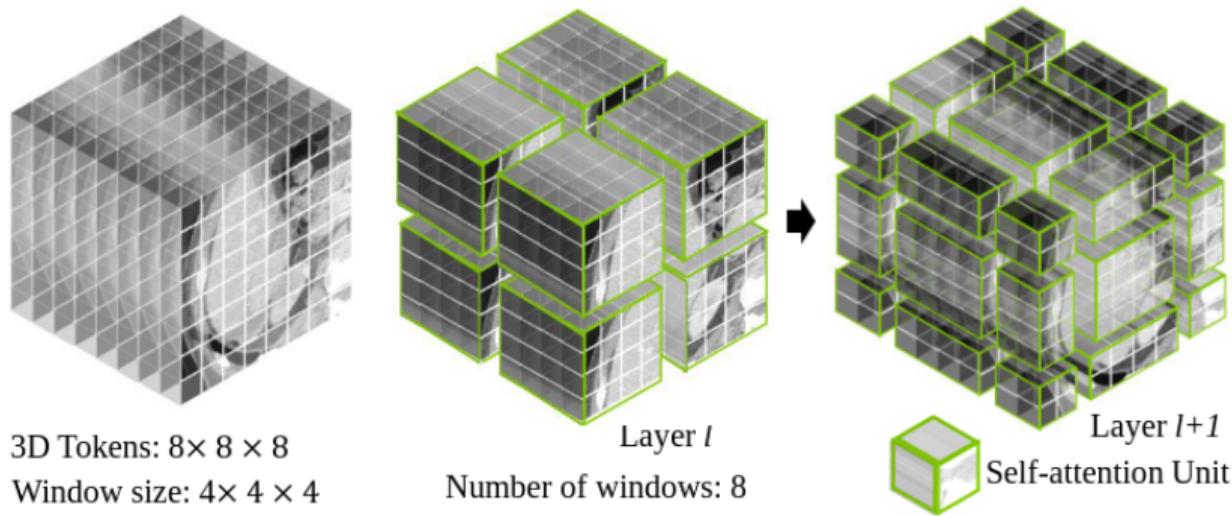


Goal: reducing information compression when patching

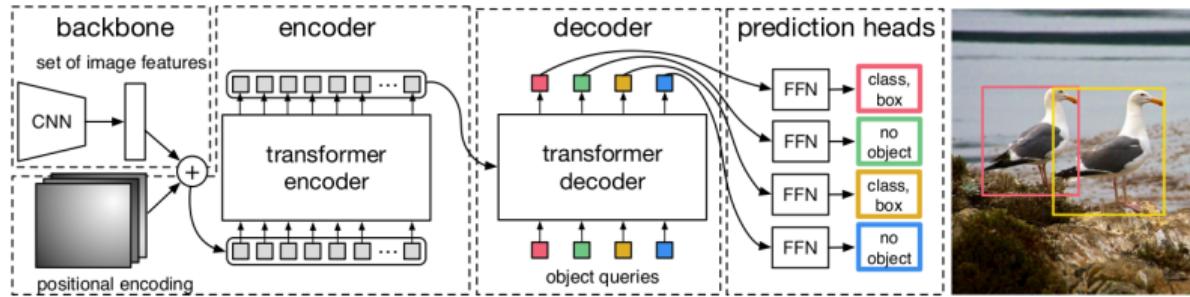
Size of input sub-volume:  $128 \times 128 \times 128 \times 4$

Patch size:  $2 \times 2 \times 2$

Input dimension for Swin:  $(64 \times 64 \times 64) \times 48 \rightarrow 262,144 \times 48$

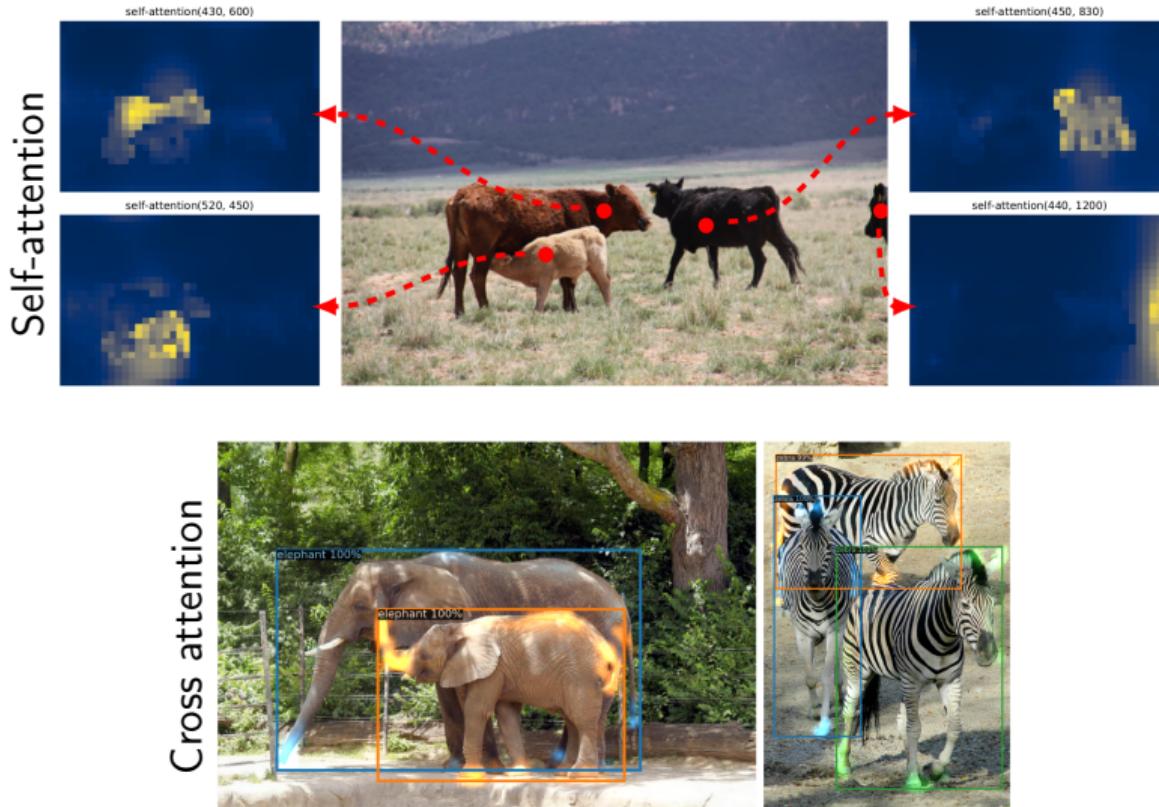


## ► Object detection

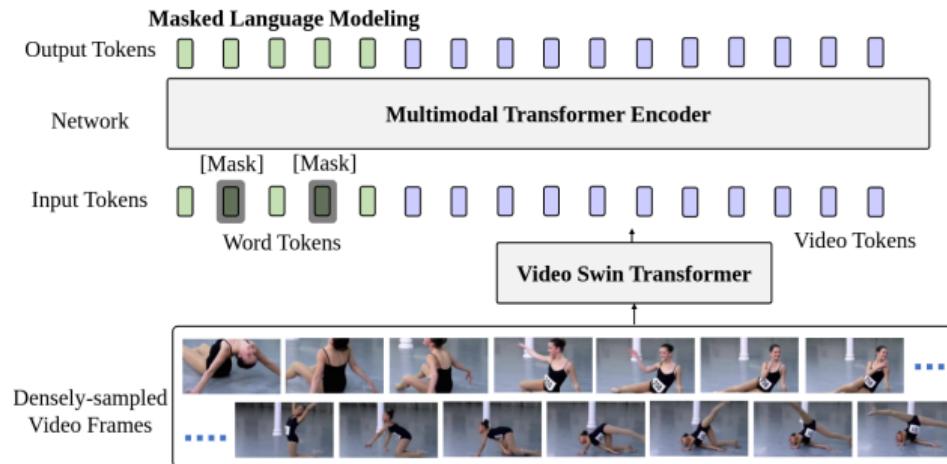


Set of  $N$  learned object queries (must be high enough)

All predictions in parallel



## ► Video captioning

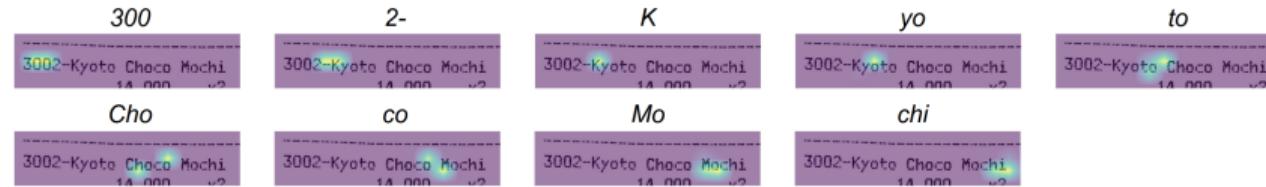
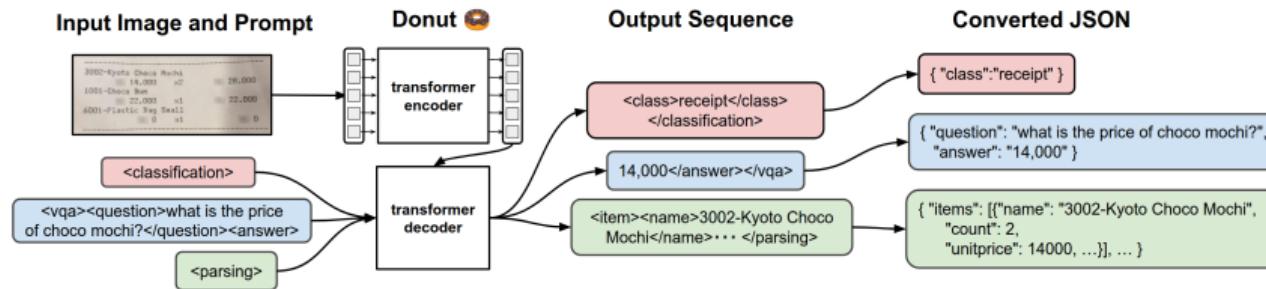


Generated caption: A dog is eating a watermelon  
GT1: A dog is eating a watermelon  
GT2: A dog eats watermelon

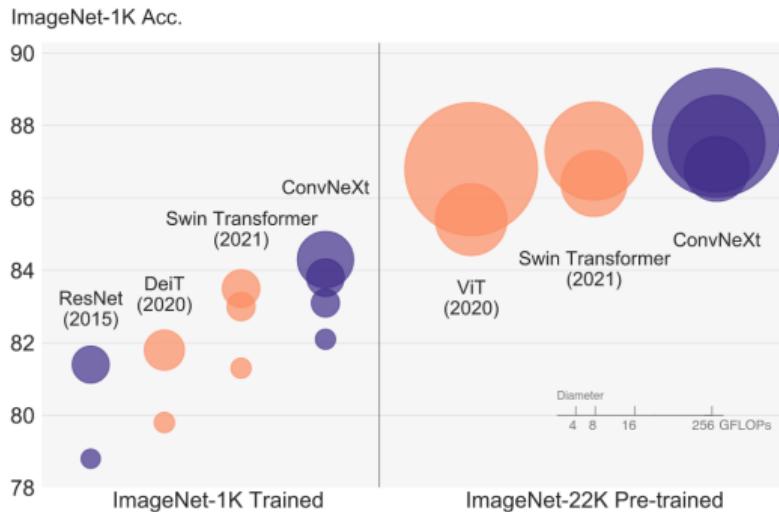


Generated caption: A boy is jumping on a trampoline  
GT1: A boy jumps on a trampoline  
GT2: Boy jumping on trampoline

## ► Document analysis and recognition



## ► ConvNext (2022) [14]



CNN with ViT/Swin patterns:

- ReLU → GELU
- Batch Norm → Layer Norm
- Less activations
- Larger kernel size
- etc

## Attention

- + A powerful tool for feature extraction and encoder/decoder paradigm
- Requires more data

## Transformer

- + State of the art
- + The same architecture for many modalities (text, image, audio)
- + Parallelizable decoding at training time (teacher forcing)
- Prohibitive for high resolution data (quadratic complexity)

## References |

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. "Neural Machine Translation by Jointly Learning to Align and Translate". In: *3rd International Conference on Learning Representations*. 2015.
- [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is All you Need". In: *Advances in Neural Information Processing Systems 30 (NIPS)*. 2017, pp. 5998–6008.
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale". In: *International Conference on Learning Representations*. 2021.
- [4] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows". In: *2021 IEEE/CVF International Conference on Computer Vision*. 2021, pp. 9992–10002.
- [5] Mehdi Noroozi and Paolo Favaro. "Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles". In: *14th European Conference on Computer Vision*. Vol. 9910. 2016, pp. 69–84.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding". In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Ed. by Jill Burstein, Christy Doran, and Thamar Solorio. Association for Computational Linguistics, 2019, pp. 4171–4186.

- [7] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. "Masked Autoencoders Are Scalable Vision Learners". In: *Conference on Computer Vision and Pattern Recognition*. 2022, pp. 15979–15988.
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. "A Simple Framework for Contrastive Learning of Visual Representations". In: *Proceedings of the 37th International Conference on Machine Learning*. Vol. 119. Proceedings of Machine Learning Research. 2020, pp. 1597–1607.
- [9] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. "Learning Transferable Visual Models From Natural Language Supervision". In: *Proceedings of the 38th International Conference on Machine Learning*. Vol. 139. Proceedings of Machine Learning Research. 2021, pp. 8748–8763.
- [10] Yucheng Tang, Dong Yang, Wenqi Li, Holger R. Roth, Bennett A. Landman, Daguang Xu, Vishwesh Nath, and Ali Hatamizadeh. "Self-Supervised Pre-Training of Swin Transformers for 3D Medical Image Analysis". In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 20698–20708.
- [11] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. "End-to-End Object Detection with Transformers". In: *16th European Conference on Computer Vision*. Vol. 12346. 2020, pp. 213–229.

## References III

- [12] Kevin Lin, Linjie Li, Chung-Ching Lin, Faisal Ahmed, Zhe Gan, Zicheng Liu, Yumao Lu, and Lijuan Wang. "SwinBERT: End-to-End Transformers with Sparse Attention for Video Captioning". In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 2022, pp. 17928–17937.
- [13] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. "OCR-Free Document Understanding Transformer". In: *17th European Conference on Computer Vision*. Vol. 13688. 2022, pp. 498–517.
- [14] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. "A ConvNet for the 2020s". In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2022, pp. 11966–11976.