

## 1. Recurrent Neural Networks (RNNs) and Their Trade-offs

- **Advantages:**
  - Can process sequences of arbitrary length.
  - Shares weights across time, making it parameter-efficient.
- **Disadvantages:**
  - Sequential processing leads to slow computation.
  - Struggles with long-term dependencies due to vanishing gradients.

## 2. seq2seq Task

### 2.1 only RNNs

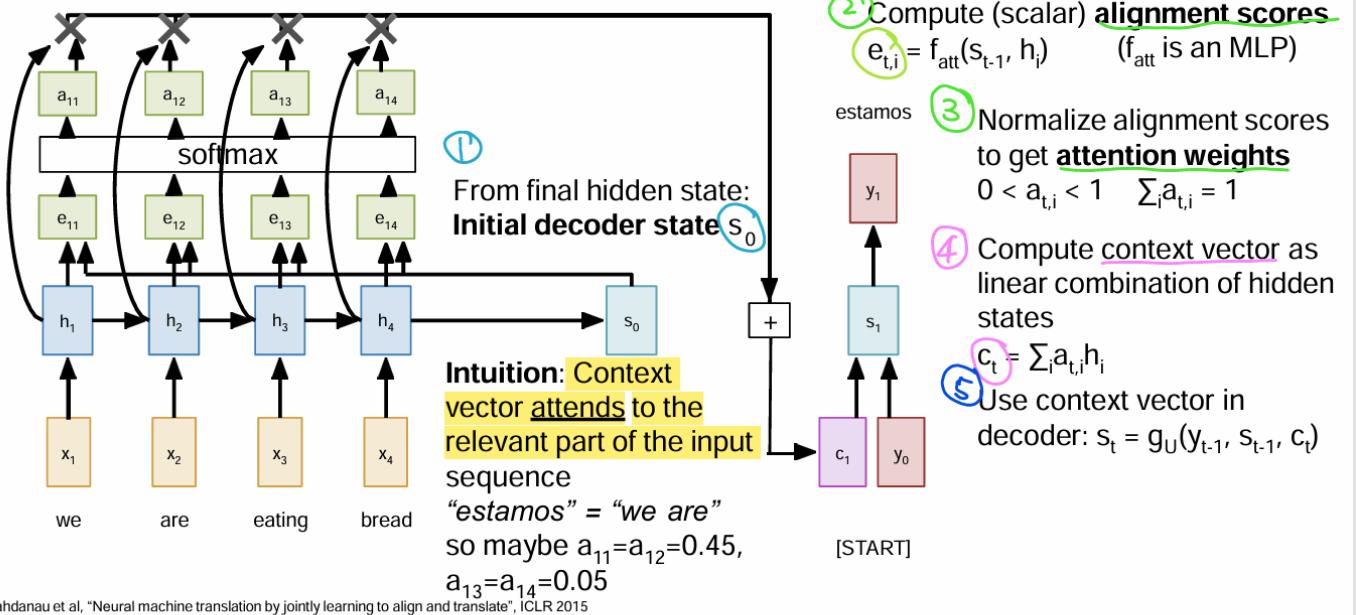
Problem: The encoder compresses the entire input sequence into a single fixed-size vector (the final hidden state of the encoder). This vector serves as the sole context for the decoder to generate the output sequence. For long input sequences, this fixed-size vector often fails to capture all relevant information, leading to loss of critical details.  
Decoders struggle to generate accurate outputs for longer or more complex sequences.

### 2.2 RNN + Attention

- Attention addresses the bottleneck of fixed-size context vectors by allowing the decoder to focus on specific parts of the input sequence dynamically.
- **Key Features:**
  - Alignment Scores

- Context Vector

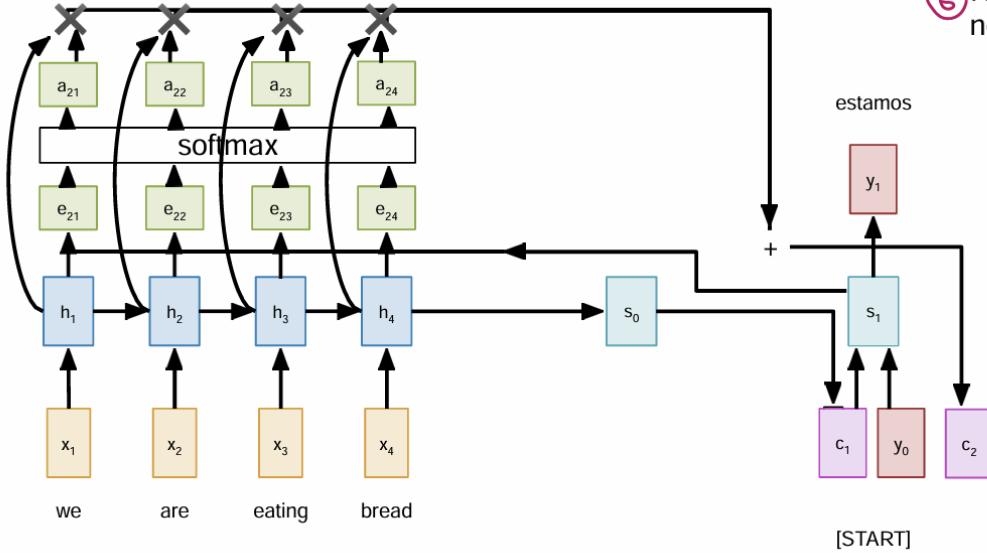
## Sequence to Sequence with RNNs and Attention



30

## Sequence to Sequence with RNNs and Attention

6 Repeat: Use  $s_1$  to compute new context vector  $c_2$



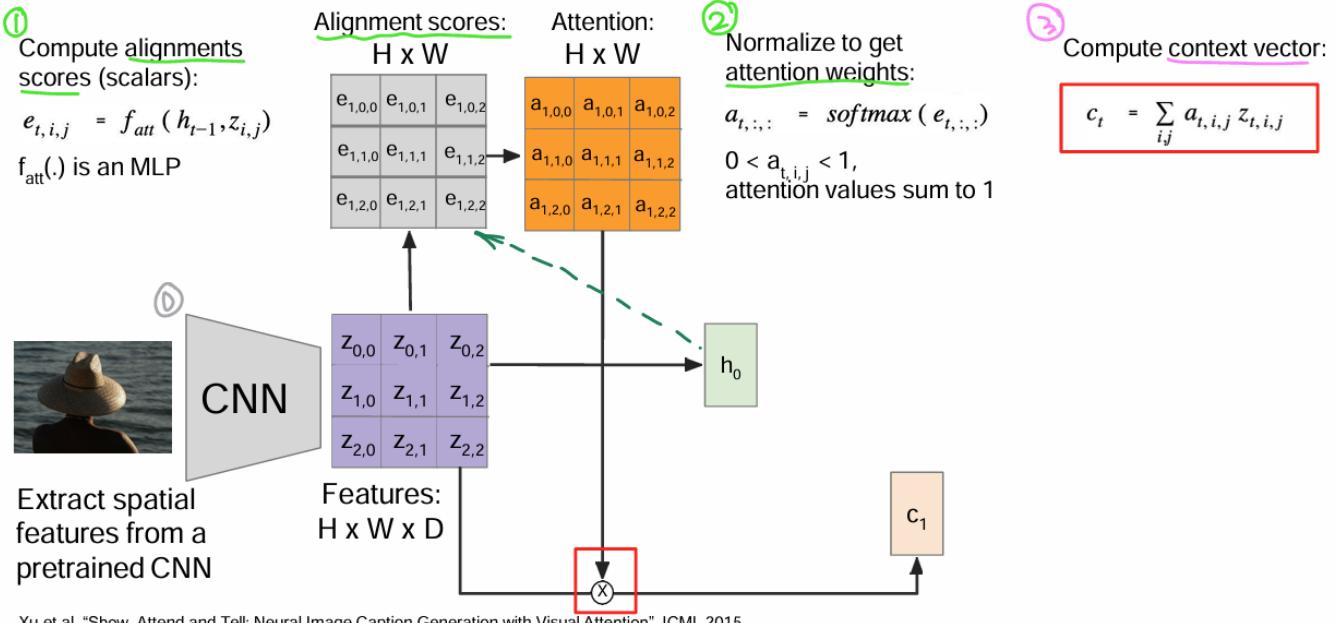
- different context vector in each timestep of decoder
- Input sequence not bottlenecked through single vector
- At each timestep of decoder, context vector "looks at" different parts of the input sequence

### 3. Image Captioning Task

#### 3.1 RNN + Attention

New context vector at every time step. Each context vector will attend to different image regions.

# Image Captioning with RNNs and Attention

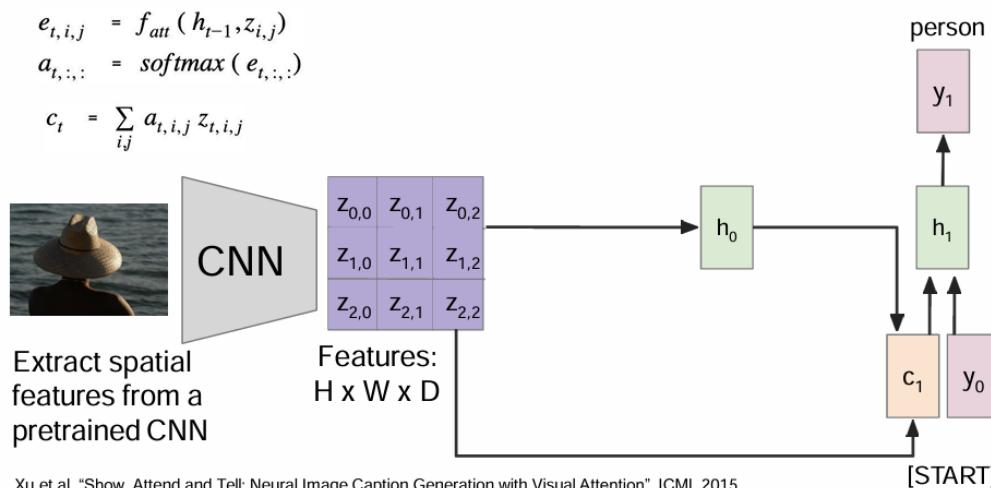


38

# Image Captioning with RNNs and Attention

Each timestep of decoder uses a different context vector that looks at different parts of the input image

**Decoder:**  $y_t = g_V(y_{t-1}, h_{t-1}, c_t)$   
New context vector at every time step



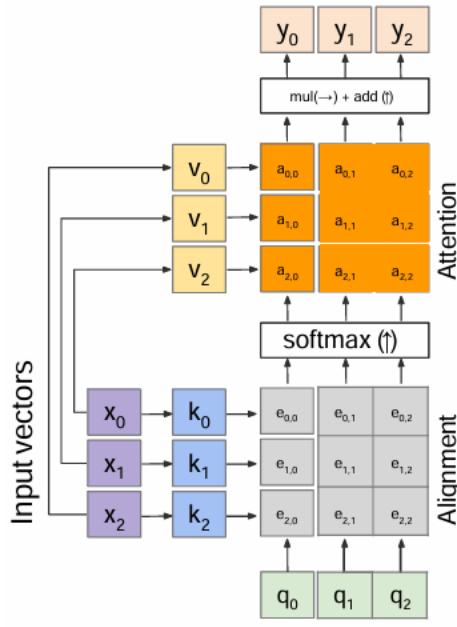
This entire process is differentiable.  
Model chooses its own attention weights, no attention supervision is required.

## 4. General att. V.S. Self att.

Aspect	General Attention Layer	Self-Attention Layer
Purpose, focuses on relationship	between two sets of vectors (e.g., encoder-decoder).	within the same sequence.
Inputs	Query (decoder), Keys and Values (encoder).	Query, Key, and Value are all derived from

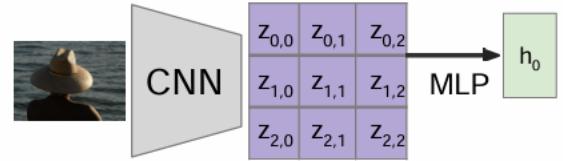
Aspect	General Attention Layer	Self-Attention Layer
Context Vector	information from the external input sequence.	relationships among tokens in the same sequence
Usage	Used in encoder-decoder attention for seq2seq tasks.	Core of self-contained models like Transformer
Applications	Machine translation, image captioning, multimodal tasks.	Language modeling, text generation, image captioning
Complexity	Depends on the sizes of the query and key-value sets.	Scales quadratically with sequence length

## General attention layer



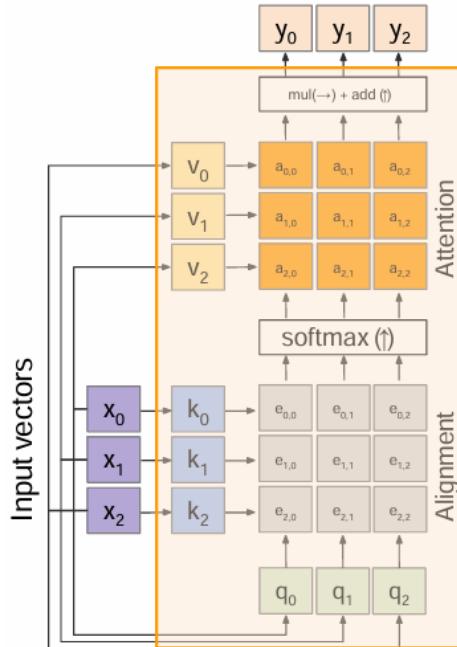
Recall that the query vector was a function of the input vectors

**Encoder:**  $h_0 = f_w(\mathbf{z})$   
where  $\mathbf{z}$  is spatial CNN features  
 $f_w(\cdot)$  is an MLP

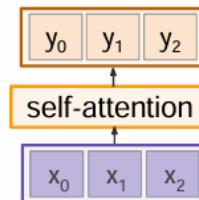


Inputs:  
Input vectors:  $\mathbf{x}$  (shape:  $N \times D$ )  
Queries:  $\mathbf{q}$  (shape:  $M \times D_k$ )

## Self attention layer - attends over sets of inputs



**Self-attention**



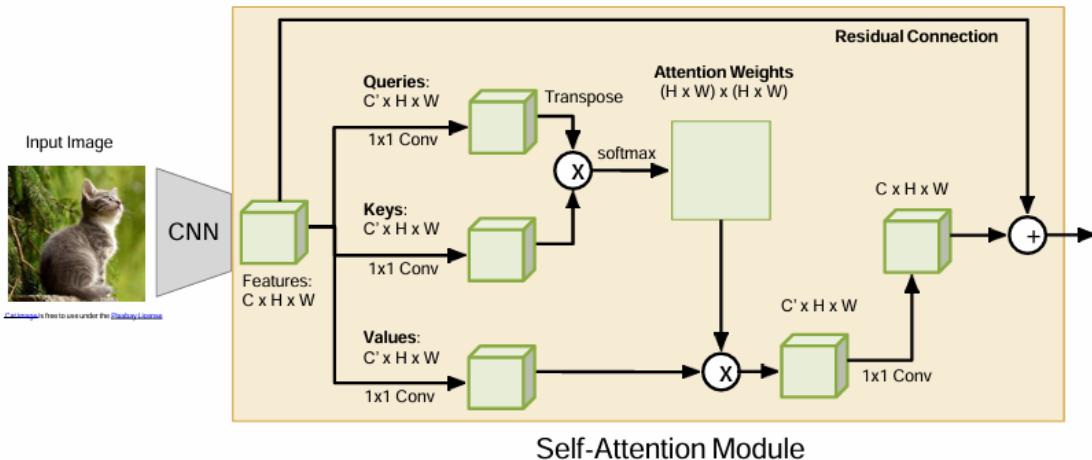
Inputs:  
Input vectors:  $\mathbf{x}$  (shape:  $N \times D$ )

Self-attention layer doesn't care about the orders of the inputs!

- Problem: How to encode ordered sequences like language and spatially ordered img features?

- Solution: Positional encoding

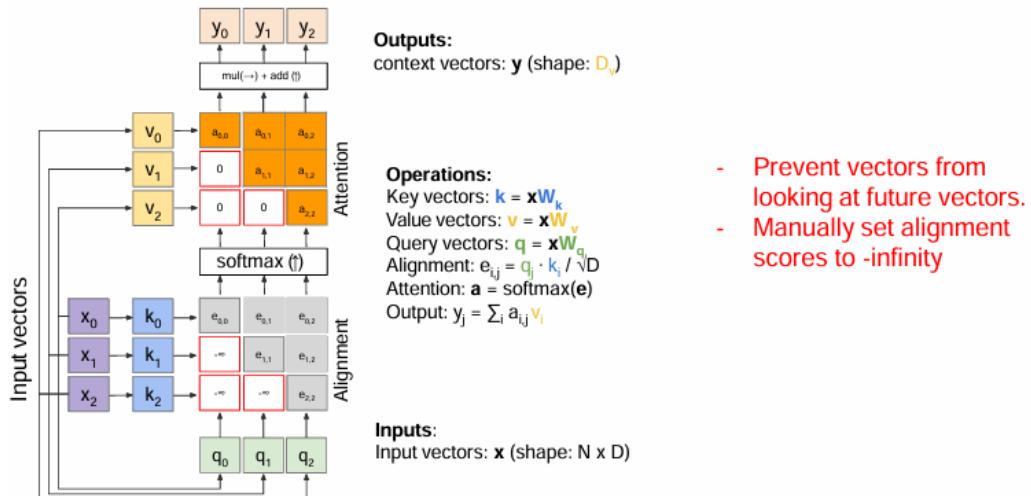
## Example: CNN with Self-Attention



## 5. Transformers

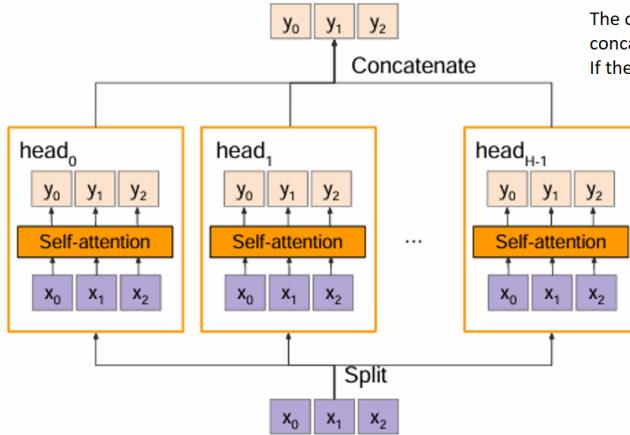
- Introduced in "Attention Is All You Need" (Vaswani et al., 2017).
- Transformers are a type of layer that uses self-attention and layer norm.
  - It is highly scalable and highly parallelizable
  - Faster training, larger models, better performance across vision and language tasks
- Key Innovations: multi-head att self-att, positional encoding, masked att.

### Masked self-attention layer



# Multi-head self-attention layer

- Multiple self-attention heads in parallel



Each attention head produces an output of shape  $(N, dk)$  where:

- N: Sequence length (number of tokens).
- dk: Dimension of the head's output.

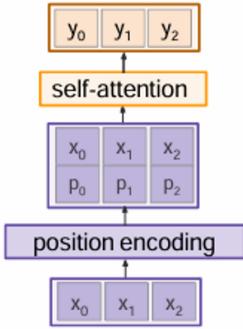
The outputs of all heads (e.g.,  $\text{head}_0, \text{head}_1, \dots, \text{head}_{H-1}$ ) are concatenated along the **feature dimension**.

If there are H heads, the concatenated output has shape  $(N, H \cdot dk)$

Options for *pos()*:

1. Learn a lookup table:
  - o Learn parameters to use for *pos(t)* for  $t \in [0, T]$
  - o Lookup table contains  $T \times d$  parameters.
2. Design a fixed function with the desiderata

# Positional encoding



Concatenate **special positional encoding  $p$**  to **each input vector  $x_i$** .

We use a function  $\text{pos}: N \rightarrow \mathbb{R}^d$  to process the position  $j$  of the vector into a  $d$ -dimensional vector

So,  $p_j = \text{pos}(j)$

$$p(t) = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_d$$

Intuition:

0 :	0 0 0 0	8 :	1 0 0 0
1 :	0 0 0 1	9 :	1 0 0 1
2 :	0 0 1 0	10 :	1 0 1 0
3 :	0 0 1 1	11 :	1 0 1 1
4 :	0 1 0 0	12 :	1 1 0 0
5 :	0 1 0 1	13 :	1 1 0 1
6 :	0 1 1 0	14 :	1 1 1 0
7 :	0 1 1 1	15 :	1 1 1 1

Vaswani et al, "Attention is all you need", NeurIPS 2017

things desired

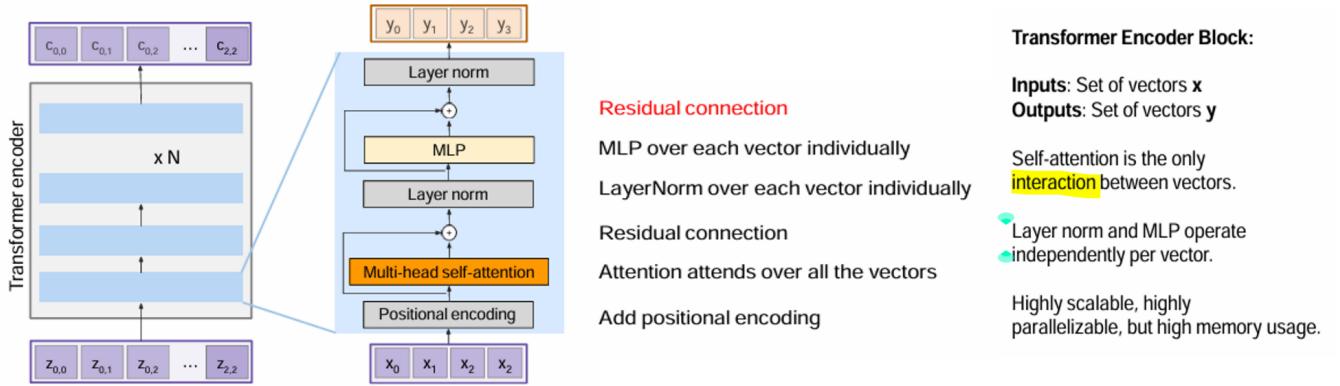
Desiderata of *pos()*:

1. It should output a **unique** encoding for each **time-step** (word's position in a sentence)
2. **Distance** between any two time-steps should be **consistent** across sentences with different lengths.
3. Our model should generalize to **longer** sentences without any efforts. Its values should be bounded.
4. It must be **deterministic**.

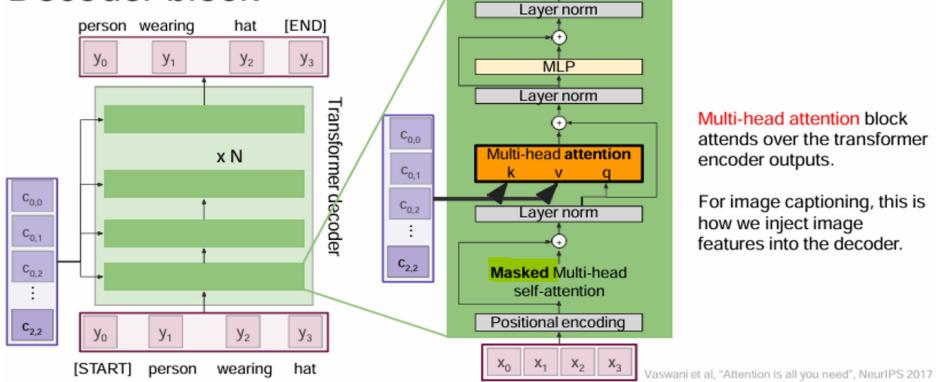
Vaswani et al, "Attention is all you need", NeurIPS 2017

The output or behavior can be exactly predicted based on the input and the defined rules.

## The Transformer encoder block



## The Transformer Decoder block



## Comparing RNNs to Transformer

### RNNs

- (+) LSTMs work reasonably **well** for long sequences.
- (-) Expects an **ordered** sequences of inputs
- (-) Sequential computation: subsequent hidden states can only be computed after the previous ones are done.

### Transformer:

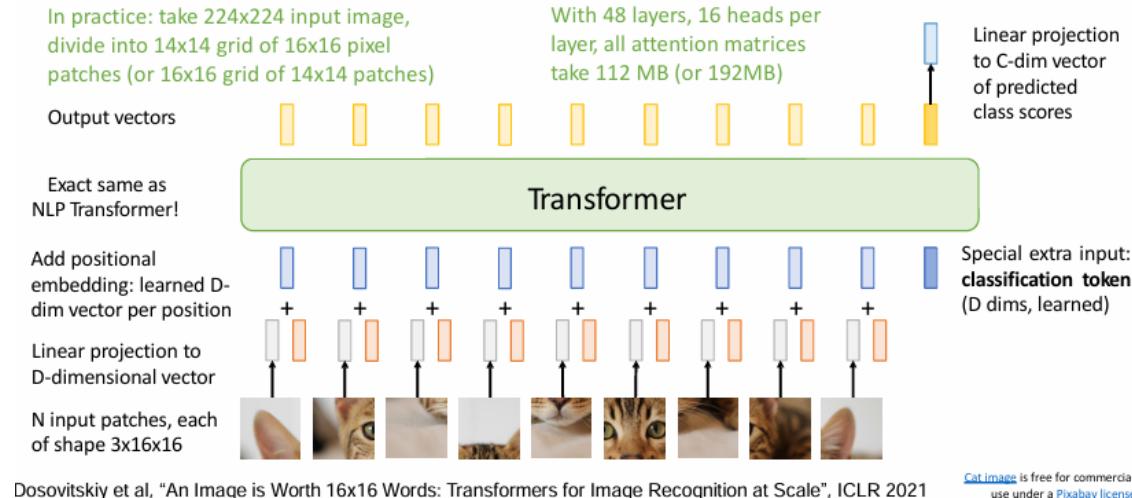
- (+) Good at **long** sequences. Each attention calculation looks at all inputs.
- (+) Can operate over **unordered** sets or **ordered** sequences with positional encodings.
- (+) **Parallel** computation: All alignment and attention scores for all inputs can be done in parallel.
- (-) Requires a lot of **memory**:  $N \times M$  alignment and attention scalers need to be calculated and stored for a single self-attention head. (but GPUs are getting bigger and better)

## 6. Vision Transformers (ViTs)

- Adapt Transformers for image processing.
- Split images into patches (e.g.,  $(16 \times 16)$ ) and treat each patch as a sequence token.

- Pretrained on large-scale image datasets, demonstrating competitive performance with CNNs.

## Vision Transformer (ViT)



## ViT vs CNN

Stage 3:  
256 x 14 x 14

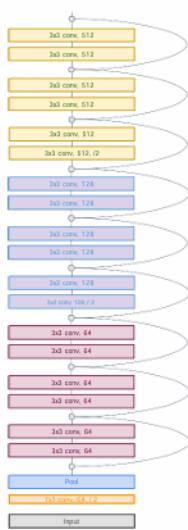
In most CNNs (including ResNets), **decrease resolution** and **increase channels** as you go **deeper** in the network (Hierarchical architecture)

Useful since objects in images can occur at various **scales**

Stage 2:  
128 x 28 x 28

Stage 1:  
64 x 56 x 56

Input:  
3 x 224 x 224



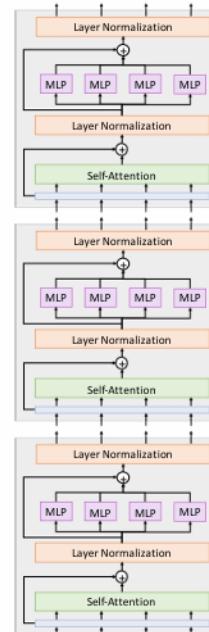
Can we build a **hierarchical ViT model?**

In a ViT, all blocks have **same resolution** and **number of channels** (**Isotropic** architecture)  
各向同性

3rd block:  
768 x 14 x 14

2nd block:  
768 x 14 x 14

1st block:  
768 x 14 x 14



Input:  
3 x 224 x 224

### 6.1 Hierarchical ViT:

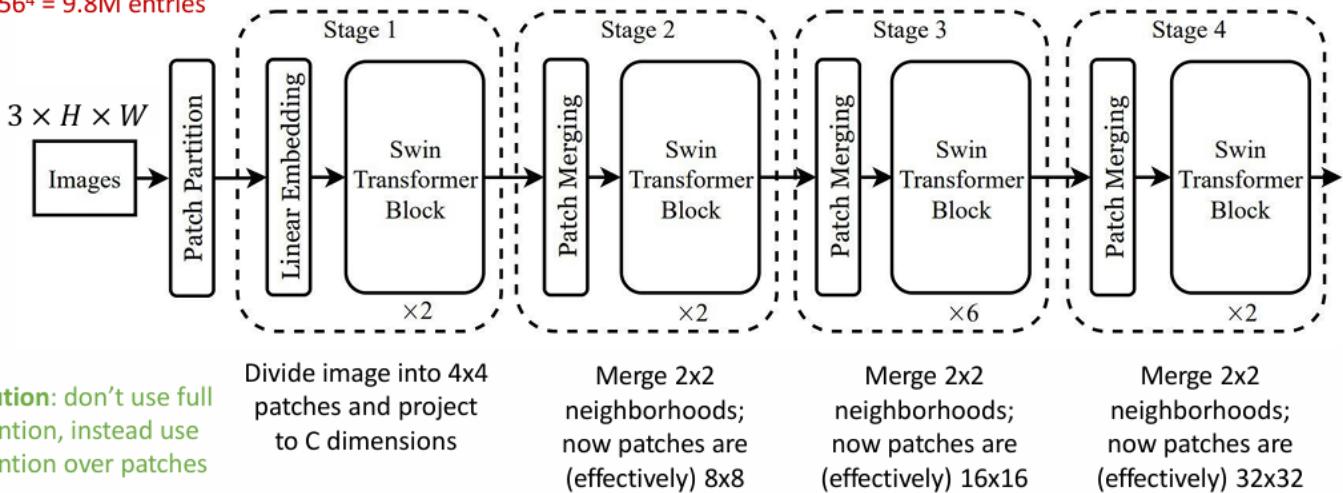
MViT, Improved MViT, Swin (Swin Transformer: shifted window attention)

Swin:

# Hierarchical ViT: Swin Transformer

**Problem:** 224x224 image with 56x56 grid of 4x4 patches: attention matrix has  $56^4 = 9.8M$  entries

$$C \times \frac{H}{4} \times \frac{W}{4} \quad 2C \times \frac{H}{8} \times \frac{W}{8} \quad 4C \times \frac{H}{16} \times \frac{W}{16} \quad 8C \times \frac{H}{32} \times \frac{W}{32}$$



Liu et al, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows", CVPR 2021

**Problem:** tokens only interact with other tokens within the same window; no communication across windows



With  $H \times W$  grid of tokens, each attention matrix is  $H^2W^2$  – quadratic in image size

Rather than allowing each token to attend to all other tokens, instead divide into windows of  $M \times M$  tokens (here  $M=4$ ); only compute attention within each window

Total size of all attention matrices is now:  
 $M^4(H/M)(W/M) = M^2HW$

Linear in image size for fixed  $M$ !  
 Swin uses  $M=7$  throughout the network

al, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows", CVPR 2021

**Solution:** Alternate between normal windows and shifted windows in successive Transformer blocks

**Detail: Relative Positional Bias**

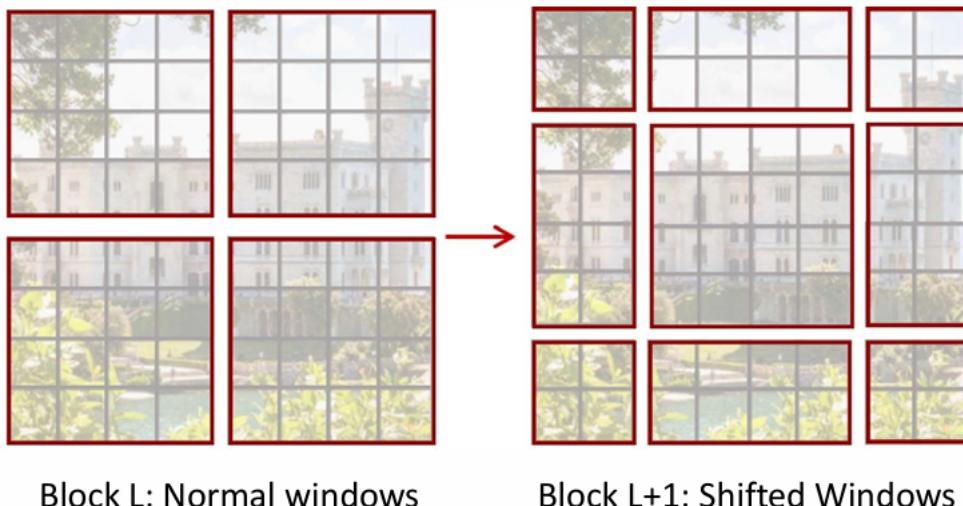
ViT adds positional embedding to input tokens, encodes *absolute position* of each token in the image

Swin does not use **positional embeddings**, instead encodes **relative position** between patches when computing attention:

Attention with relative bias:

$$A = \text{Softmax} \left( \frac{QK^T}{\sqrt{D}} + B \right) V$$

$Q, K, V: M^2 \times D$  (Query, Key, Value)  
 $B: M^2 \times M^2$  (learned biases)



u et al, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows", CVPR 2021

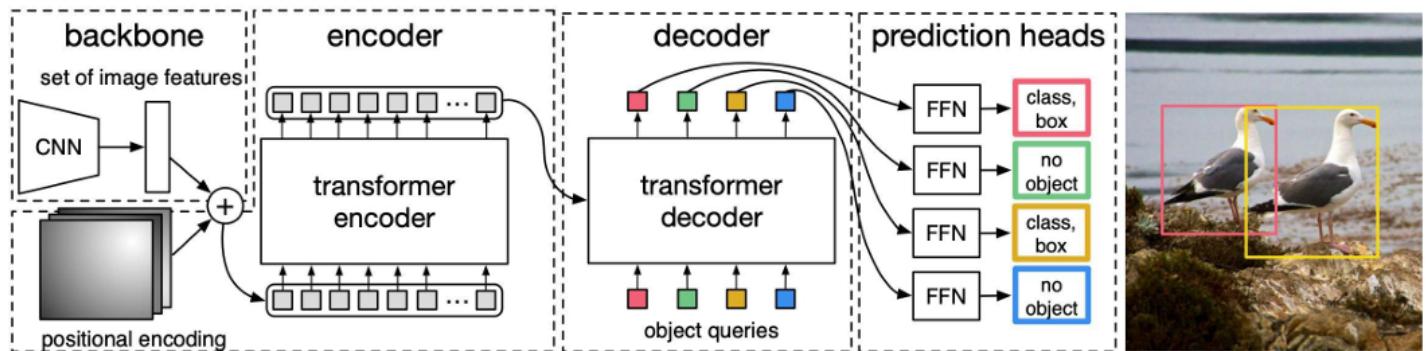
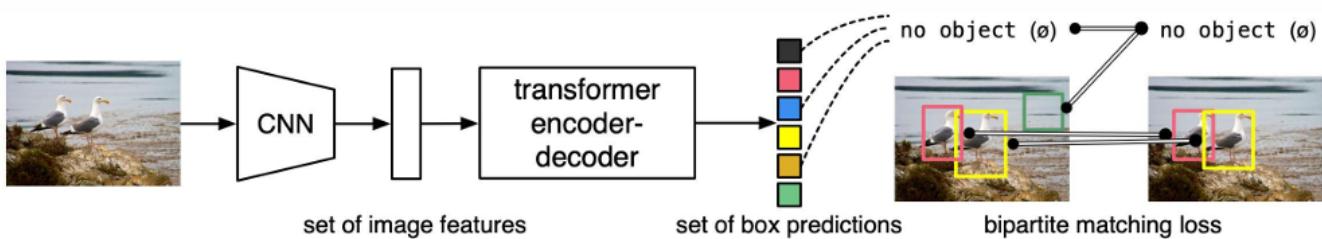
## 7. Obj detection with Transformers: DETR

Simple object detection pipeline: directly output a set of boxes from a Transformer

No anchors, no regression of box transforms

Match predicted boxes to GT boxes with bipartite matching; train to regress box coordinates

# Object Detection with Transformers: DETR



ion et al, "End-to-End Object Detection with Transformers", ECCV 2020

Architecture of DETR:

### 1. Backbone (Feature Extraction):

A CNN (e.g., ResNet) is used as a backbone to extract feature maps from the input image.

### 2. Transformer Encoder:

Processes the flattened feature maps and encodes the global context of the image.

### 3. Transformer Decoder:

Takes learned positional embeddings (object queries) and interacts with the encoder outputs to predict objects.

### 4. Prediction Heads:

- Each decoder output predicts:
  - Class label: The category of the object.
  - Bounding box: The coordinates of the object's location, represented as normalized values.

+:

1. End-to-End Training: Removes the need for complex heuristics like anchors, proposals, or NMS.
2. Global context: transformer captures global relationships in the image
3. flexibility: easily extended to other tasks like panoptic (全景式) segmentation.

-:

1. Slow Convergence
2. High Computational Cost
3. Localization Precision: Bounding box regression can sometimes be less precise compared to traditional methods.

## References

[https://cs231n.stanford.edu/slides/2022/lecture\\_10\\_ruohan.pdf](https://cs231n.stanford.edu/slides/2022/lecture_10_ruohan.pdf)

[https://cs231n.stanford.edu/slides/2022/lecture\\_11\\_ruohan.pdf](https://cs231n.stanford.edu/slides/2022/lecture_11_ruohan.pdf)