

# Advanced Computer Vision

Tanveer Hussain

[htanveer3797@gmail.com](mailto:htanveer3797@gmail.com)



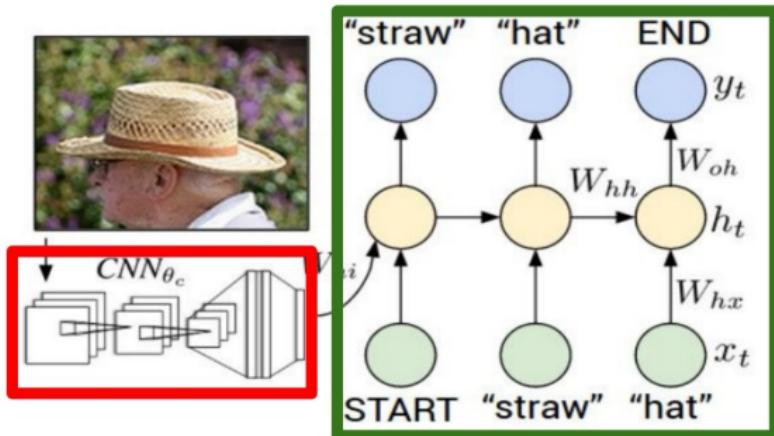
جامعة الملك عبد الله  
لعلوم والتكنولوجيا  
King Abdullah University of  
Science and Technology

KAUST Academy  
King Abdullah University of Science and Technology

Course designed by **Naeem Ullah Khan** ([naeemullah.khan@kaust.edu.sa](mailto:naeemullah.khan@kaust.edu.sa)), updated by Tanveer Hussain

1. Recurrent Neural Networks and Attention
2. General Attention and ViT
3. Video Classification
4. Generative AI (Part 1)
5. Generative AI (Part 2), CLIP model

## Recurrent Neural Network



## Convolutional Neural Network

# Image Captioning

- ▶ A computer Vision task in which we generate captions for images



A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



A dog is running in the grass with a frisbee



A white teddy bear sitting in the grass



Two people walking on the beach with surfboards



A tennis player in action on the court

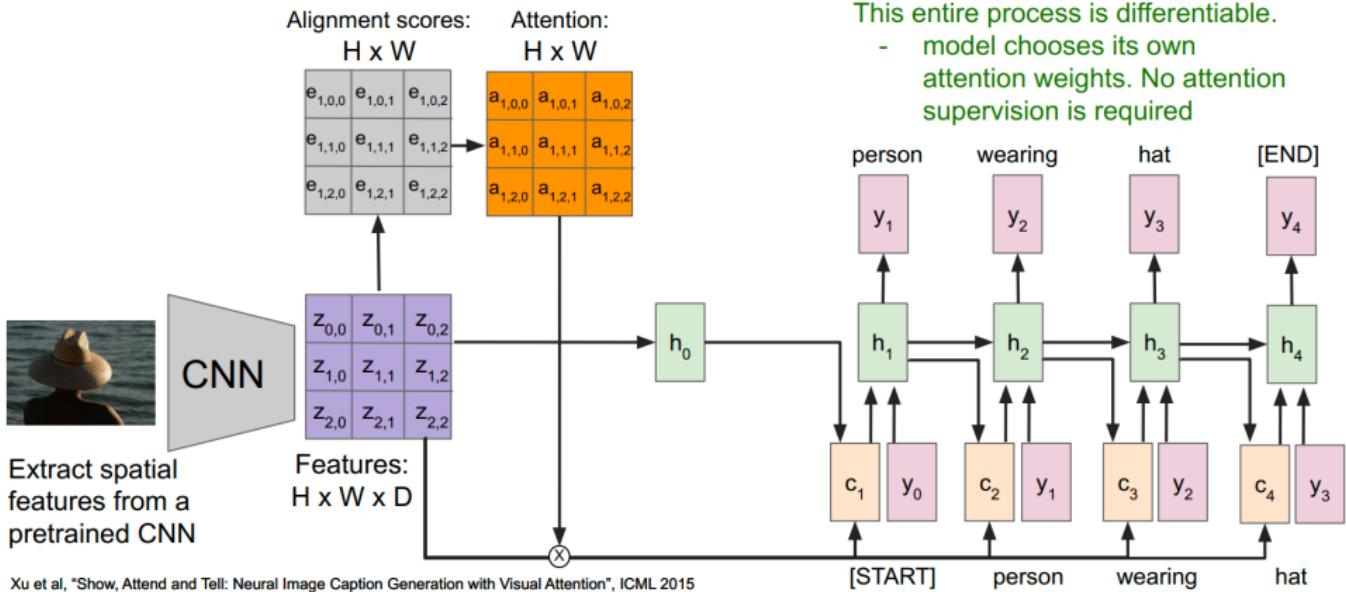


Two giraffes standing in a grassy field

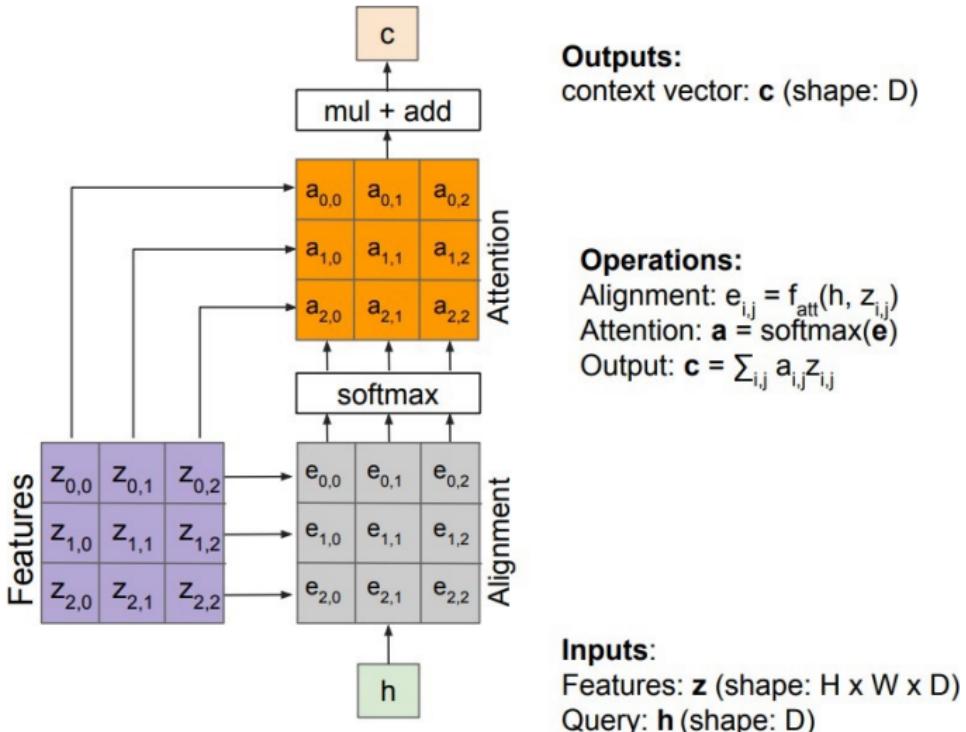


A man riding a dirt bike on a dirt track

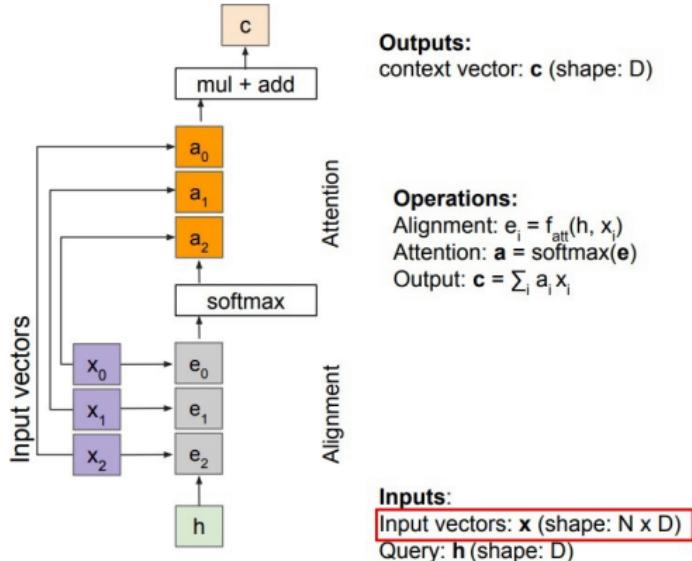
# Image Captioning with RNNs and Attention



# Attention we just saw in image captioning

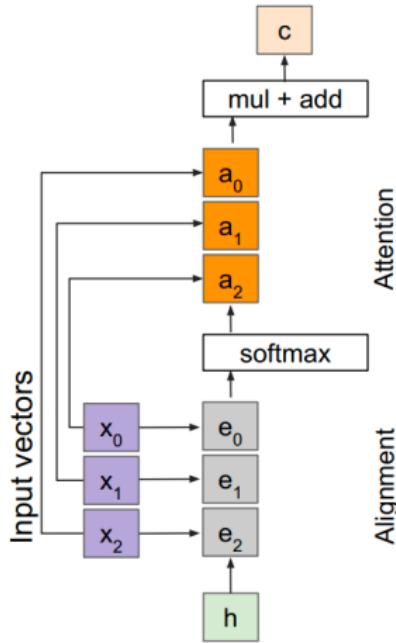


# General Attention Layer



- ▶ Attention operation is permutation invariant
- ▶ Doesn't care about ordering of the features
- ▶ Stretch  $H \times W = N$  into N vectors

# General Attention Layer



## Outputs:

context vector:  $c$  (shape: D)

## Operations:

Alignment:  $e_i = h \cdot x_i$

Attention:  $a = \text{softmax}(e)$

Output:  $c = \sum_i a_i x_i$

Change  $f_{\text{att}}(\cdot)$  to a simple dot product

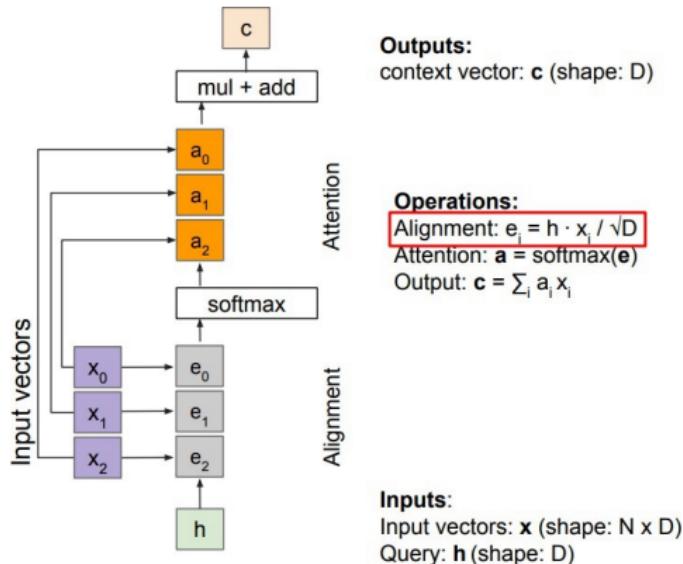
- only works well with key & value transformation trick (will mention in a few slides)

## Inputs:

Input vectors:  $x$  (shape: N x D)

Query:  $h$  (shape: D)

# General Attention Layer



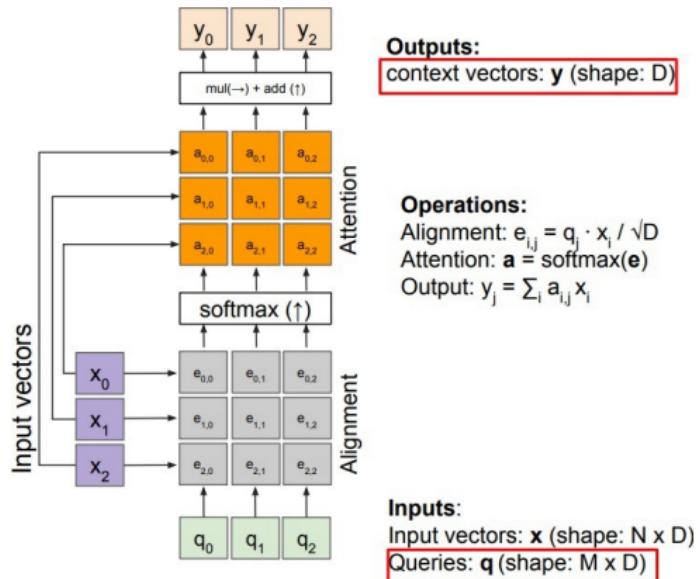
This scaling helps to:

- Normalize the logits to a more consistent range, reducing variance.
- Prevent extremely high values that lead to low-entropy softmax outputs.
- Ensures the softmax distribution remains balanced and informative.

Change  $f_{att}(\cdot)$  to a scaled simple dot product

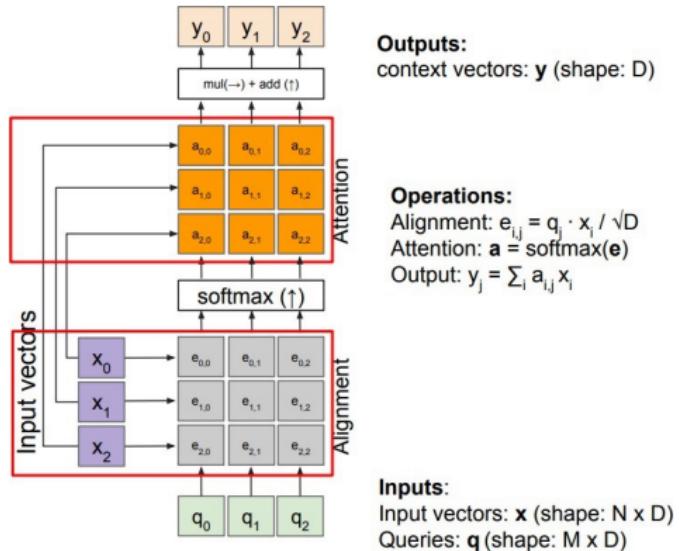
- ▶ Larger dimensions means more terms in the dot product sum.
- ▶ So, the variance of the logits is higher. Large magnitude vectors will produce much higher logits.
- ▶ So, the post-softmax distribution has lower-entropy, assuming logits are IID.
- ▶ Ultimately, these large magnitude vectors will cause softmax to peak and assign very little weight to all others
- ▶ Divide by  $\sqrt{D}$  to reduce effect of large magnitude vectors

# General Attention Layer



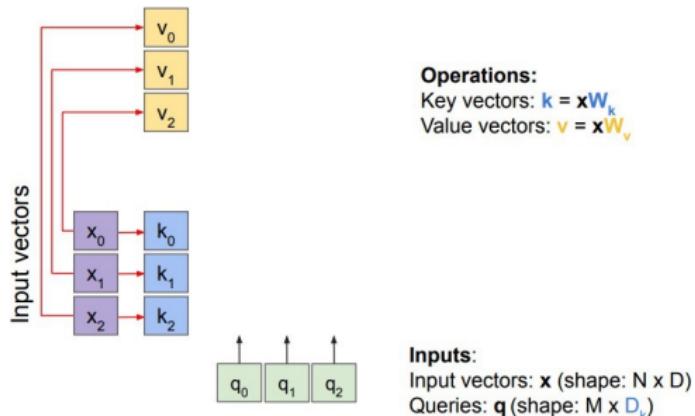
- ▶ We can have multiple Query Vectors.
- ▶ Each query creates a new output context vector

# General Attention Layer



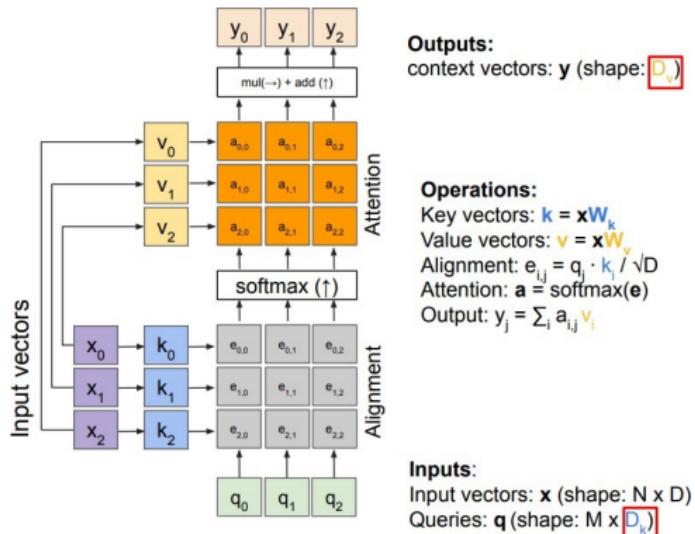
- ▶ Notice that the input vectors are used for both the alignment as well as the attention calculations.

# General Attention Layer



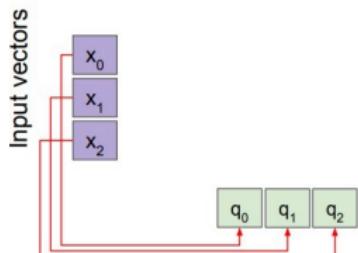
- ▶ Notice that the input vectors are used for both the alignment as well as the attention calculations.
- ▶ We can add more expressivity to the layer by adding a different FC layer before each of the two steps.

# General Attention Layer



- ▶ Notice that the input vectors are used for both the alignment as well as the attention calculations.
- ▶ We can add more expressivity to the layer by adding a different FC layer before each of the two steps.
- ▶ The input and output dimensions can now change depending on the key and value FC layers

# Self Attention Layer



## Operations:

Key vectors:  $k = xW_k$

Value vectors:  $v = xW_v$

Query vectors:  $q = xW_q$

Alignment:  $e_{ij} = q \cdot k_i / \sqrt{D}$

Attention:  $a = \text{softmax}(e)$

Output:  $y_j = \sum_i a_{ij} v_i$

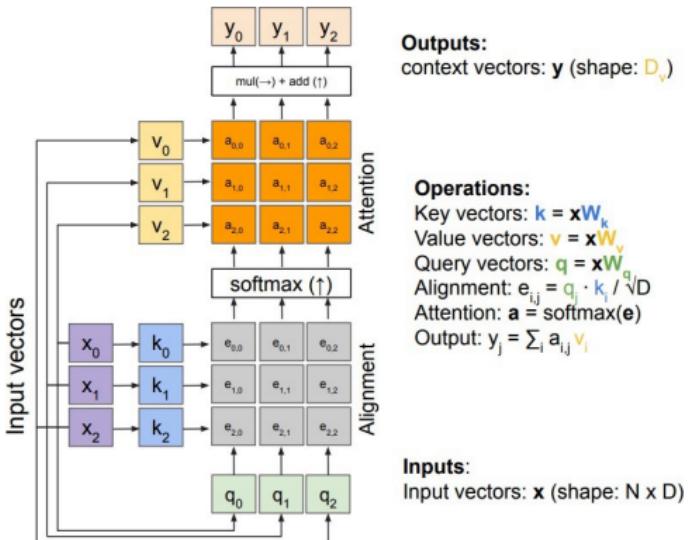
## Inputs:

Input vectors:  $x$  (shape:  $N \times D$ )

Queries:  $q$  (shape:  $M \times D_q$ )

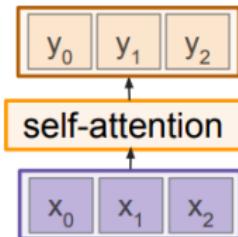
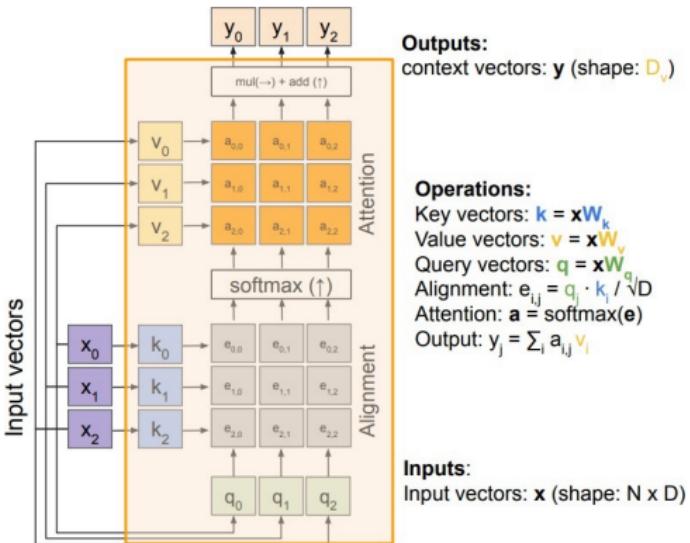
- Recall that the query vector was a function of the input vectors
- We can calculate the query vectors from the input vectors, therefore, defining a "self-attention" layer.
- No input query vectors anymore
- Instead, query vectors are calculated using a FC layer.

# Self Attention Layer

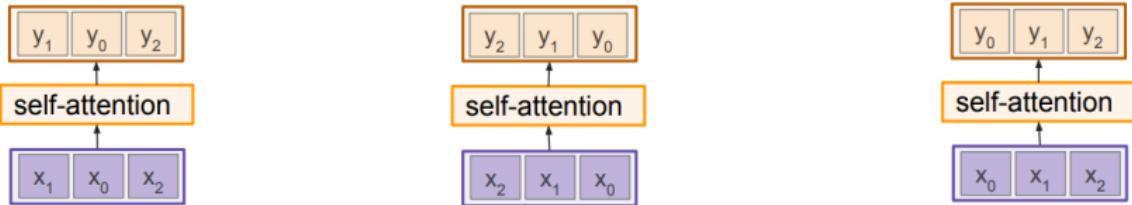


- ▶ Recall that the query vector was a function of the input vectors
- ▶ We can calculate the query vectors from the input vectors, therefore, defining a "self-attention" layer.
- ▶ No input query vectors anymore
- ▶ Instead, query vectors are calculated using a FC layer.

# Self Attention Layer



# Permutation Invariance

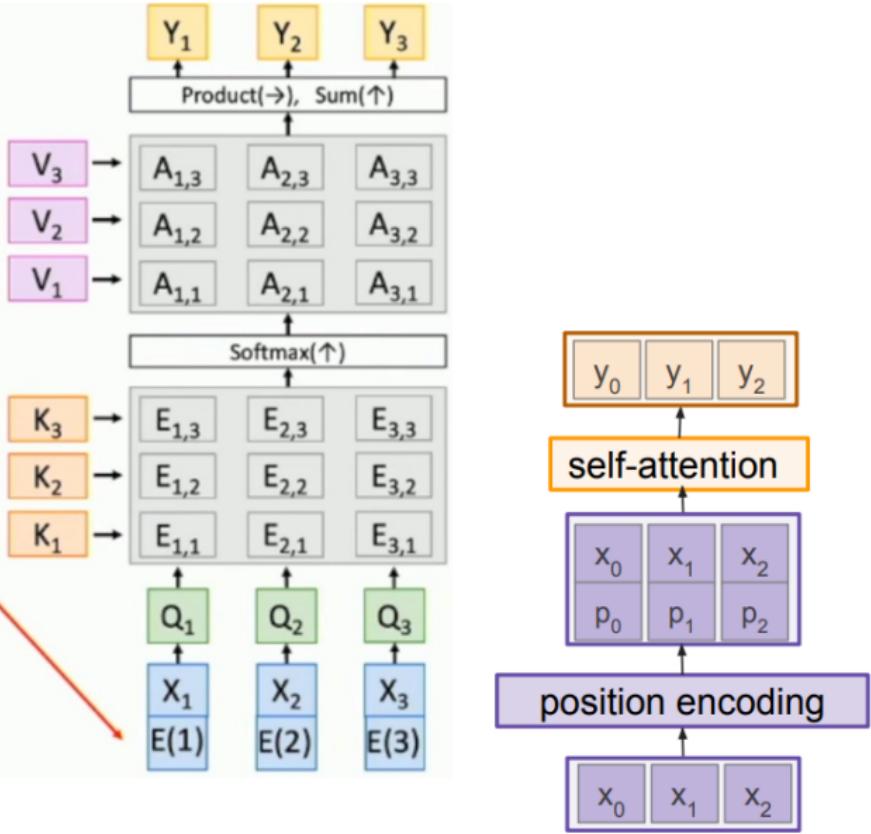


- ▶ Self-Attention Layer is Permutation equivariant
- ▶ It doesn't care about the orders of the inputs!
- ▶ **Problem:** How can we encode ordered sequences like language or spatially ordered image features?

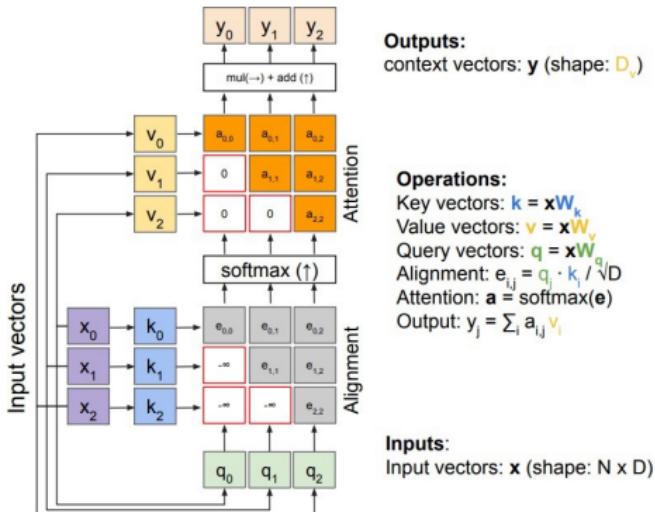
# Positional encoding

Self attention doesn't "know" the order of the vectors it is processing!

In order to make processing position-aware, concatenate input with **positional encoding**



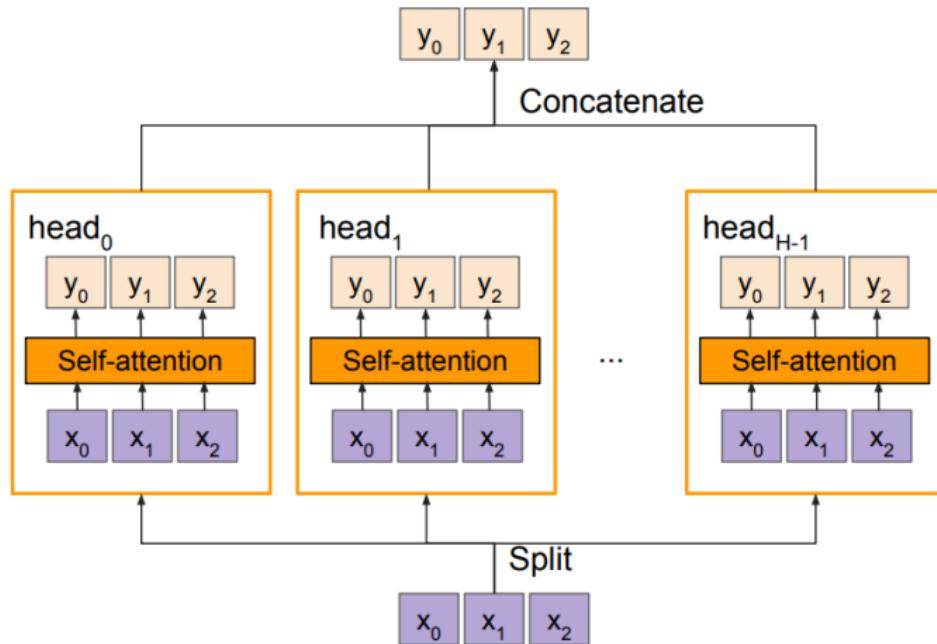
# Masked self-attention layer



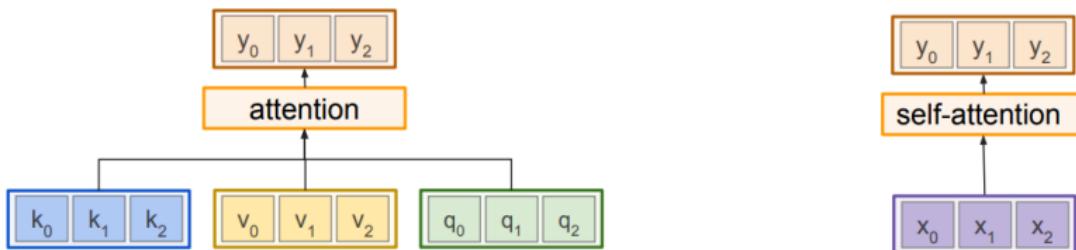
- ▶ Prevent vectors from looking at past/future vectors.
- ▶ Manually set alignment scores to  $-\infty$

# Multi-head self-attention layer

- ▶ Multiple self-attention heads in parallel



# General attention versus self-attention



# Example: CNN with Self-Attention

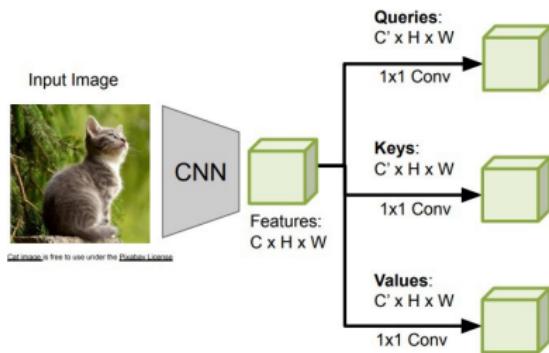
Input Image



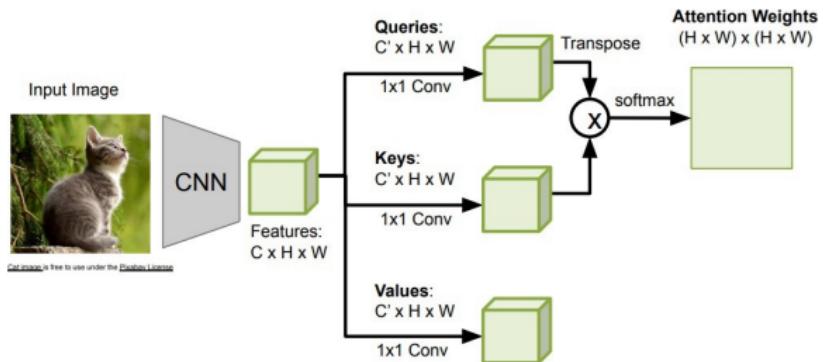
Features:  
 $C \times H \times W$

Cat image is free to use under the [Creative License](#)

# Example: CNN with Self-Attention

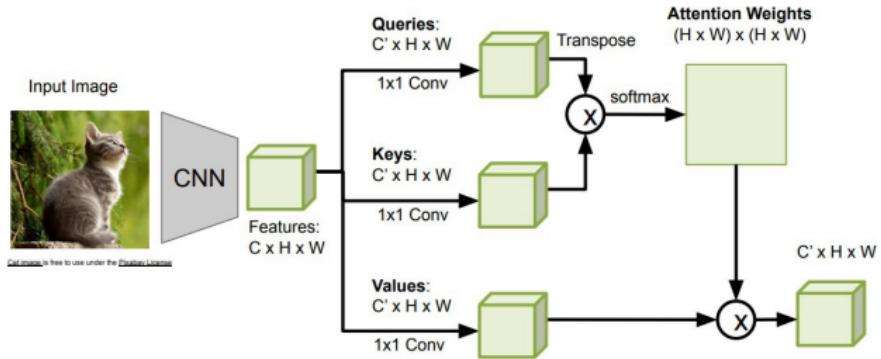


# Example: CNN with Self-Attention

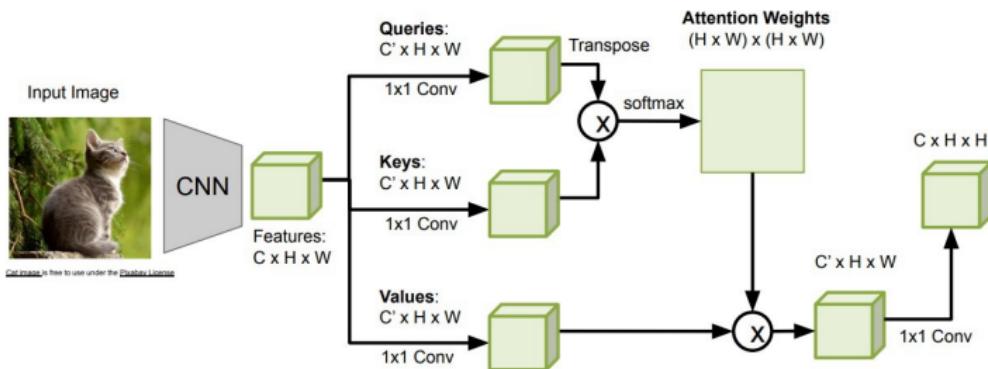


Cat image is free to use under the [Creative Commons License](#)

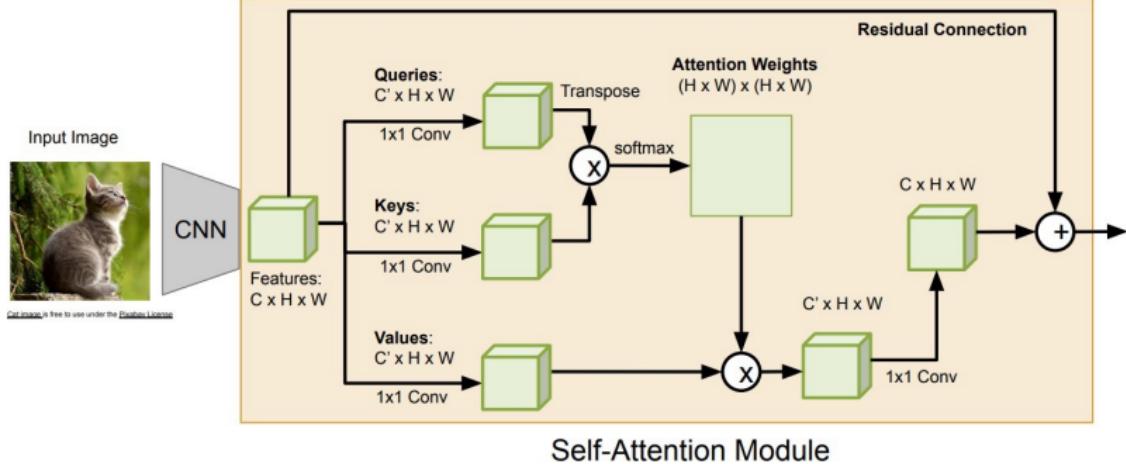
# Example: CNN with Self-Attention



# Example: CNN with Self-Attention



# Example: CNN with Self-Attention

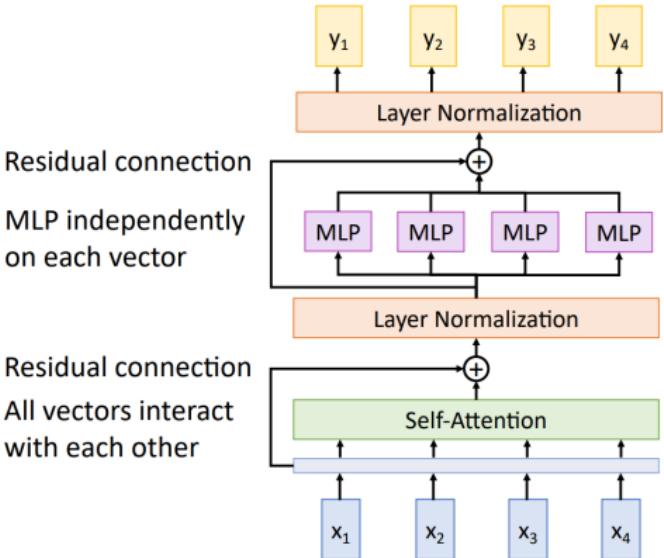


## Attention is all you need

Vaswani et al, NeurIPS 2017

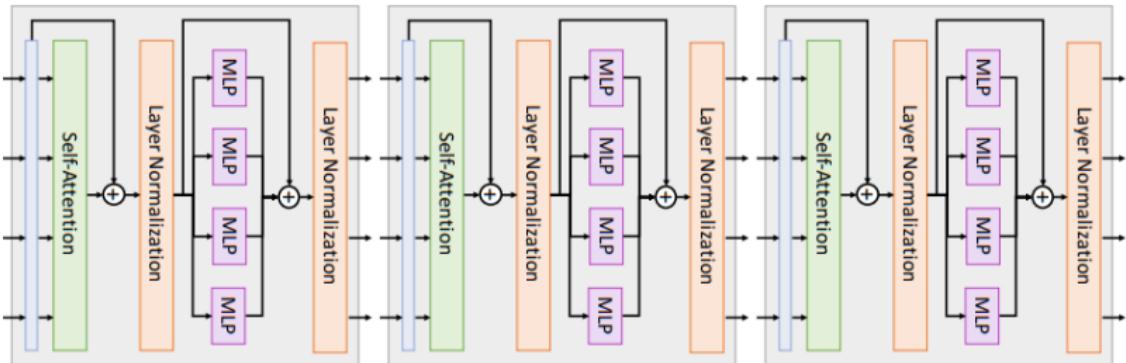
# Transformer Block

- ▶ **Input:** Set of vectors  $x$
- ▶ **Output:** Set of vectors  $y$
- ▶ Self-attention is the only interaction between vectors!
- ▶ Layer norm and MLP work independently per vector
- ▶ Highly scalable, highly parallelizable

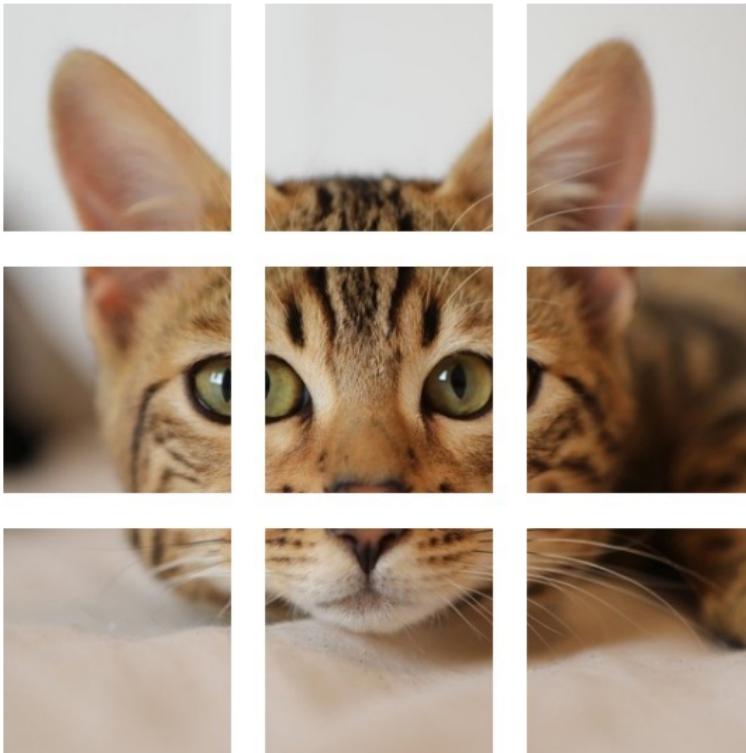


# Transformers

- ▶ A Transformer is a sequence of transformer blocks
- ▶ Vaswani et al used 12 blocks,  $DQ = 512$  and 6 attention heads



# Vision Transformers



# Vision Transformers

N input patches, each  
of shape 3x16x16



# Vision Transformers

Linear projection to  
D-dimensional vector

N input patches, each  
of shape  $3 \times 16 \times 16$

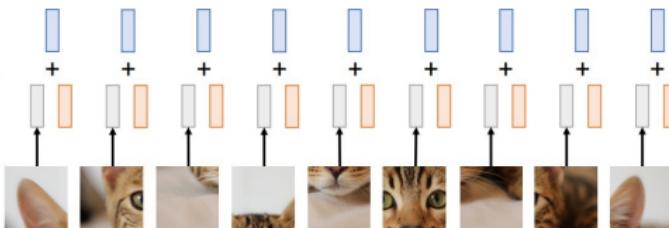


# Vision Transformers

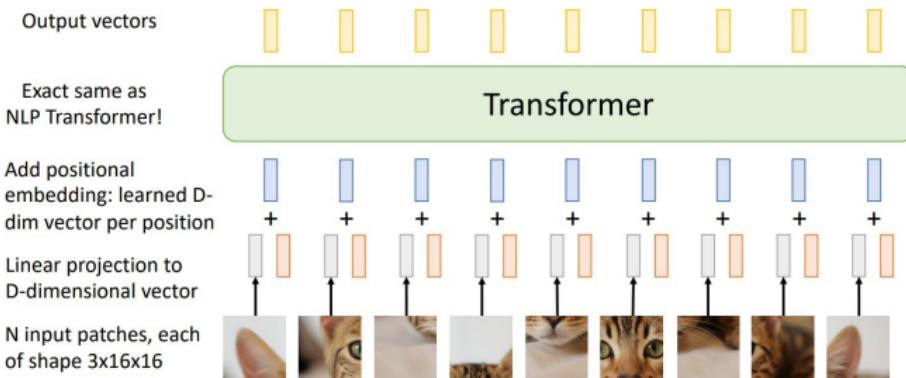
Add positional embedding: learned D-dim vector per position

Linear projection to D-dimensional vector

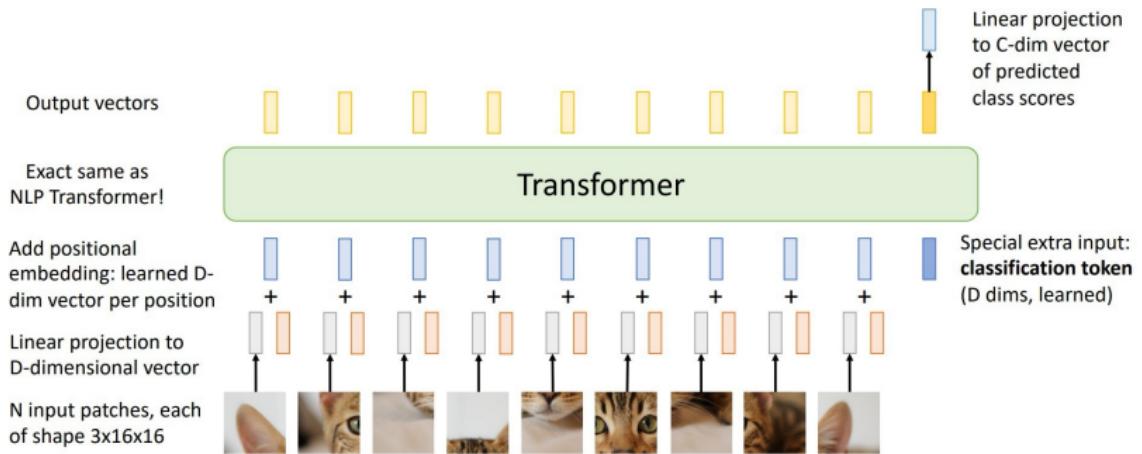
N input patches, each of shape 3x16x16



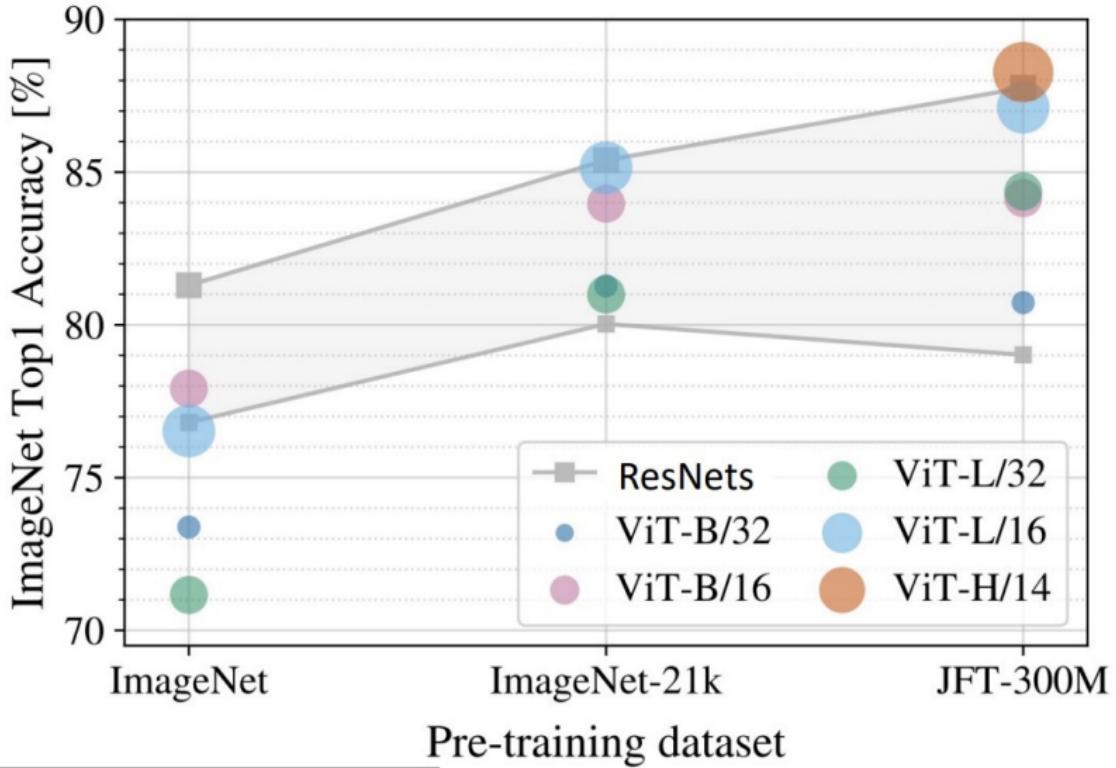
# Vision Transformers



# Vision Transformers



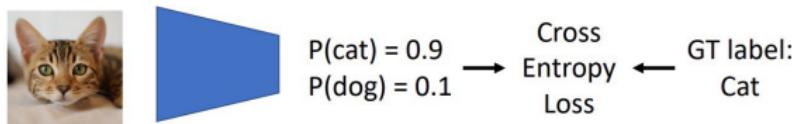
# Vision Transformers



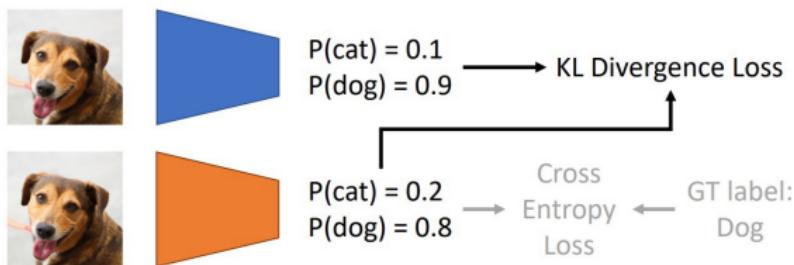
<sup>0</sup>Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

# Improving ViT: Distillation

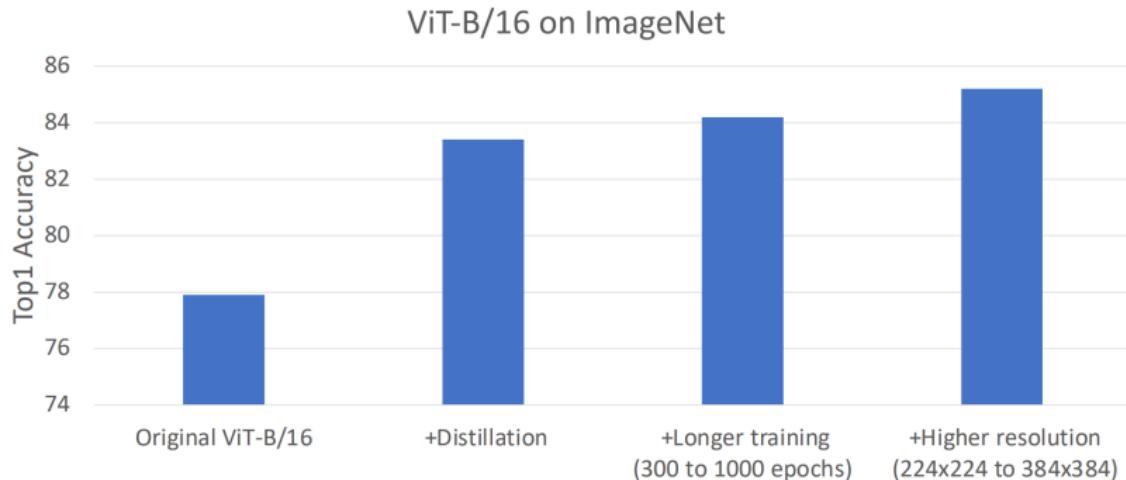
Step 1: Train a teacher CNN on ImageNet



Step 2: Train a student ViT to match ImageNet predictions from the **teacher CNN** (and match GT labels)

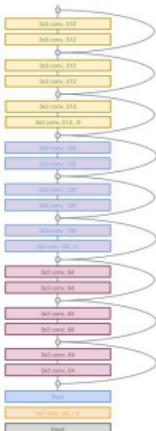


# Improving ViT: Distillation



<sup>0</sup>Touvron et al, “Training data-efficient image transformers distillation through attention”, ICML 2021

Stage 3:  
 $256 \times 14 \times 14$



Stage 2:  
 $128 \times 28 \times 28$

Stage 1:  
 $64 \times 56 \times 56$

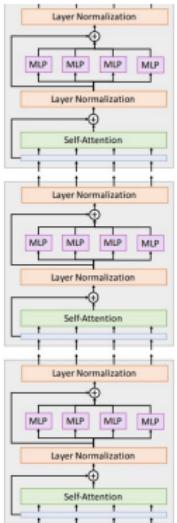
Input:  
 $3 \times 224 \times 224$

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

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

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



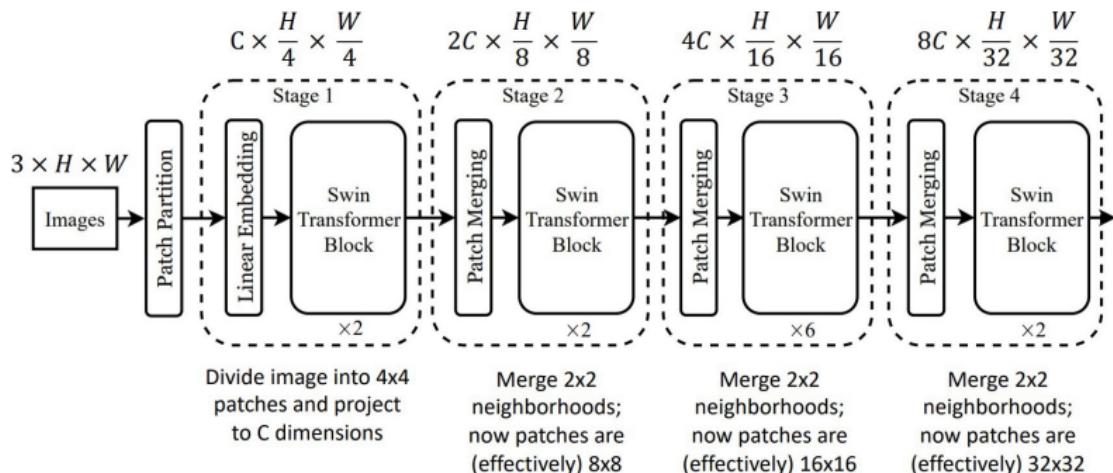
3<sup>rd</sup> block:  
 $768 \times 14 \times 14$

2<sup>nd</sup> block:  
 $768 \times 14 \times 14$

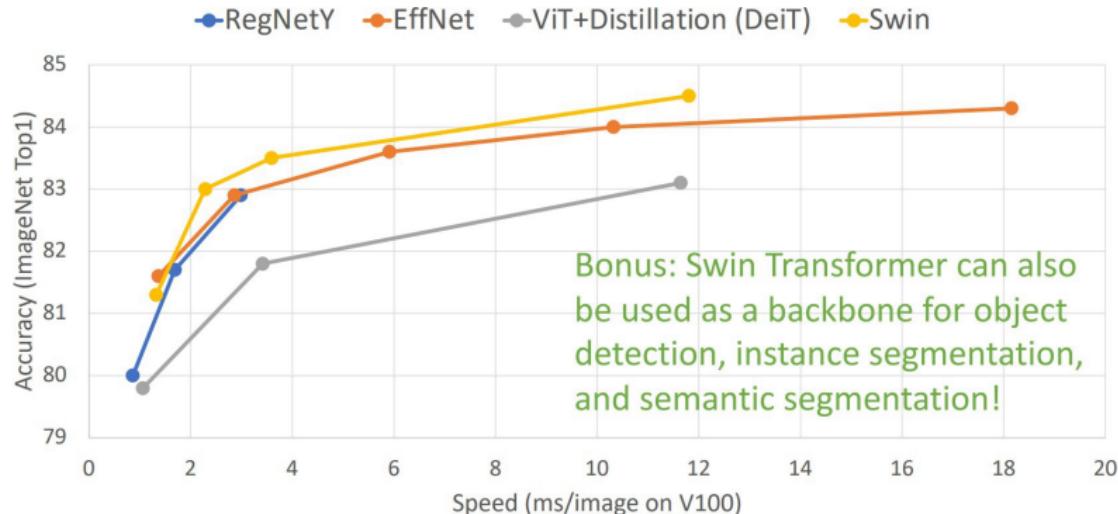
1<sup>st</sup> block:  
 $768 \times 14 \times 14$

Input:  
 $3 \times 224 \times 224$

# Hierarchical ViT: Swin Transformer



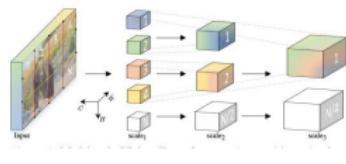
# Hierarchical ViT: Swin Transformer



<sup>0</sup>Liu et al, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows", CVPR 2021

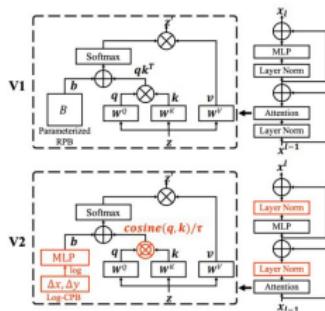
# Other Hierarchical Vision Transformers

MViT



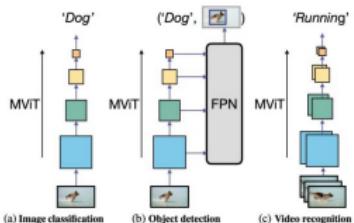
Fan et al, "Multiscale Vision  
Transformers", ICCV 2021

Swin-V2



Liu et al, "Swin Transformer V2: Scaling  
up Capacity and Resolution", CVPR 2022

Improved MViT

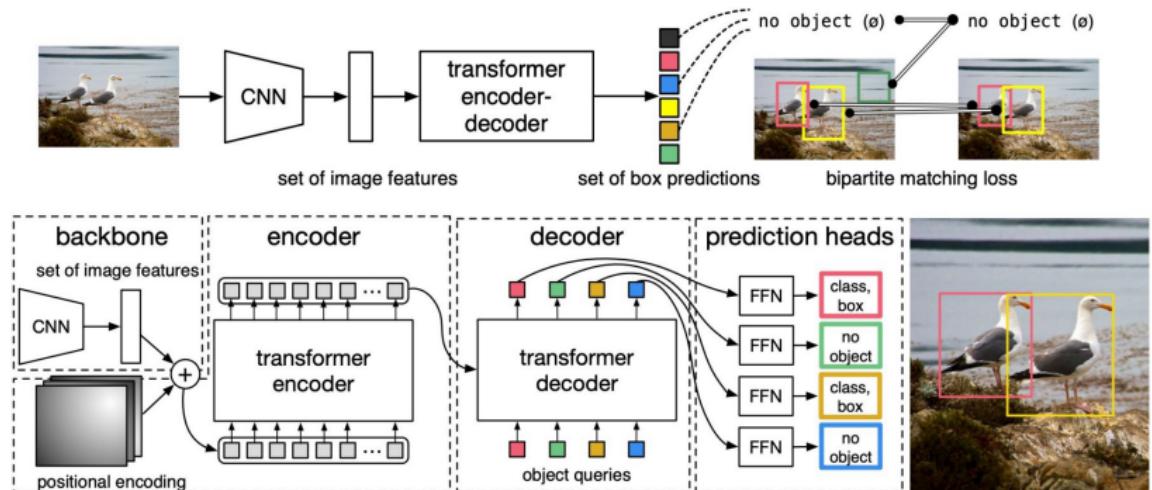


Li et al, "Improved Multiscale Vision Transformers  
for Classification and Detection", arXiv 2021

# Object 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



<sup>0</sup>Carion et al, "End-to-End Object Detection with Transformers", ECCV 2020

These slides have been adapted from

- ▶ Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: [Deep Learning for Computer Vision](#)
- ▶ Assaf Shocher, Shai Bagon, Meirav Galun & Tali Dekel, WAIC DL4CV [Deep Learning for Computer Vision: Fundamentals and Applications](#)
- ▶ Justin Johnson, UMich EECS 498.008/598.008: [Deep Learning for Computer Vision](#)