

# Advanced Computer Vision

Naeemullah Khan  
[naeemullah.khan@kaust.edu.sa](mailto:naeemullah.khan@kaust.edu.sa)

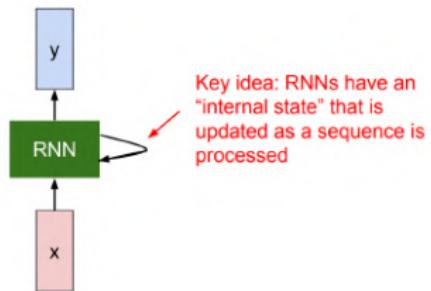


جامعة الملك عبد الله  
للغات والتكنولوجيا  
King Abdullah University of  
Science and Technology

KAUST Academy  
King Abdullah University of Science and Technology

June 15, 2023

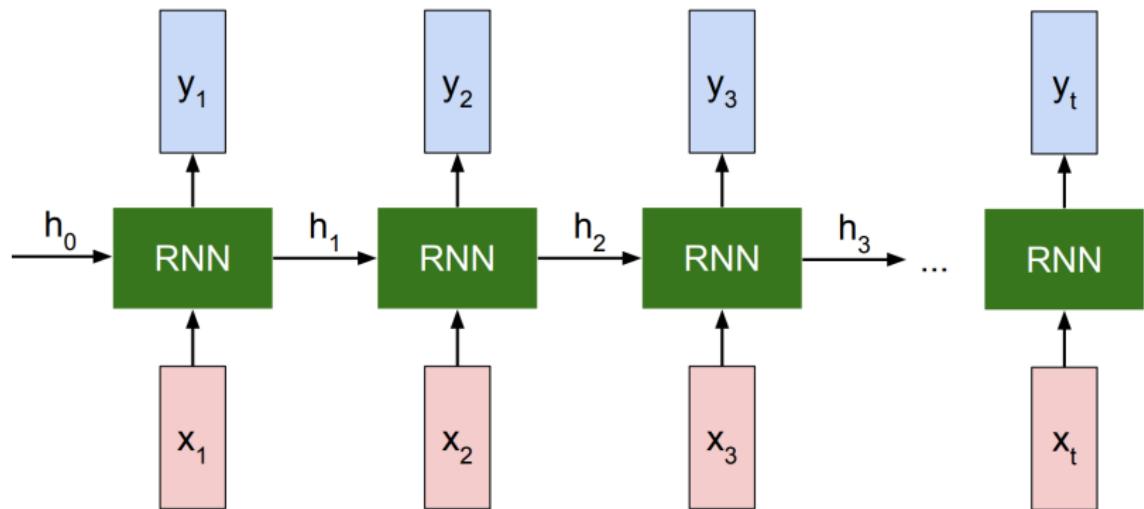
# Recurrent Neural Networks (RNN)



$$h_t = f_W(h_{t-1}, x_t)$$

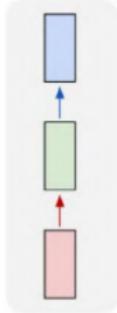
new state                  old state      input vector at some time step  
some function with parameters W

# Recurrent Neural Networks (RNN)

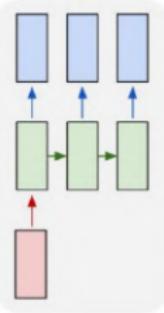


# Recurrent Neural Networks (RNN)

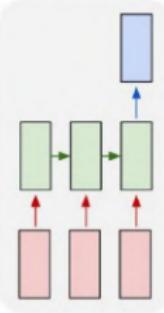
one to one



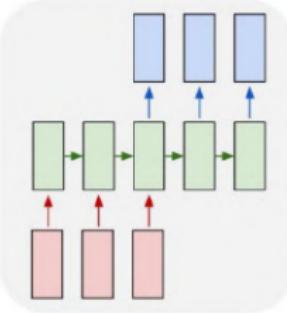
one to many



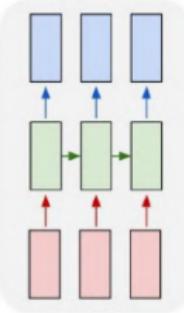
many to one



many to many



many to many



# Recurrent Neural Networks (RNN)

## ► RNN Advantages:

- Can process any length input
- Computation for step  $t$  can (in theory) use information from many steps back
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

# Recurrent Neural Networks (RNN)

## ► RNN Advantages:

- Can process any length input
- Computation for step  $t$  can (in theory) use information from many steps back
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

## ► RNN Disadvantages:

- Recurrent computation is slow
- In practice, difficult to access information from many steps back

# 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

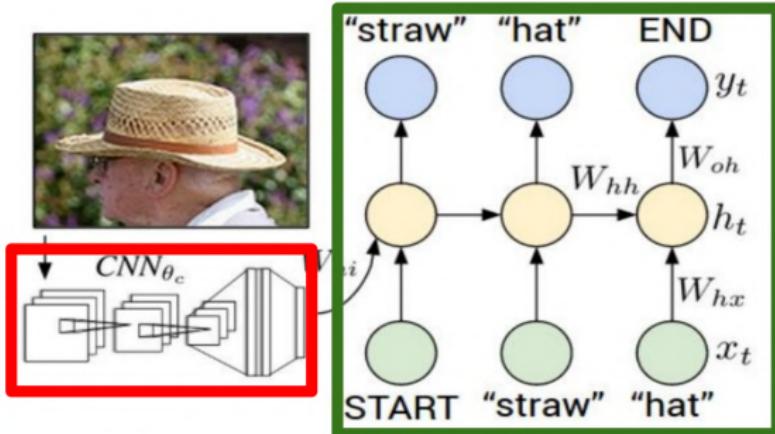


Two giraffes standing in a grassy field



A man riding a dirt bike on a dirt track

## Recurrent Neural Network



## Convolutional Neural Network

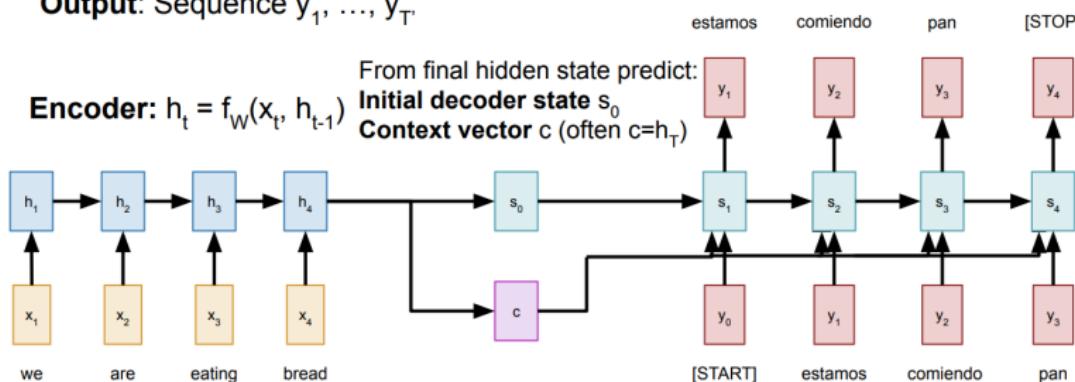
# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$   
**Output:** Sequence  $y_1, \dots, y_T$

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:  
**Initial decoder state**  $s_0$   
**Context vector**  $c$  (often  $c=h_T$ )



Sutskever et al., "Sequence to sequence learning with neural networks"

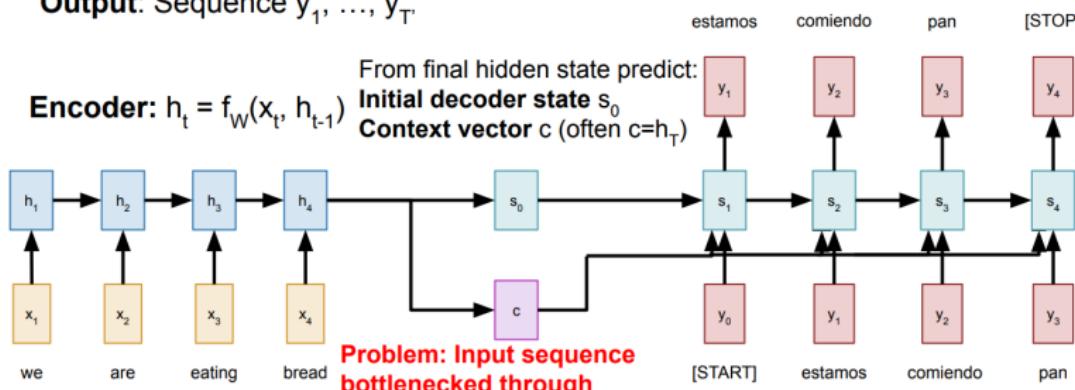
# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$   
**Output:** Sequence  $y_1, \dots, y_{T'}$

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:  
**Initial decoder state**  $s_0$   
**Context vector**  $c$  (often  $c=h_T$ )



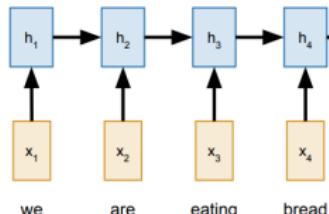
Problem: Input sequence  
bottlenecked through  
fixed-sized vector. What if  
 $T=1000$ ?

Sutskever et al., "Sequence to sequence learning with neural networks", NIPS 2014

# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$   
**Output:** Sequence  $y_1, \dots, y_T$

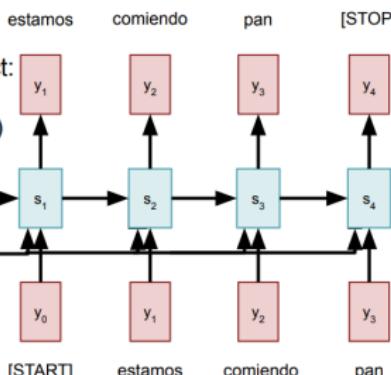
**Encoder:**  $h_t = f_W(x_t, h_{t-1})$



From final hidden state predict:  
**Initial decoder state**  $s_0$   
**Context vector**  $c$  (often  $c=h_T$ )

**Problem:** Input sequence  
bottlenecked through  
fixed-sized vector. What if  
 $T=1000$ ?

**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$



**Idea:** use new context vector  
at each step of decoder!

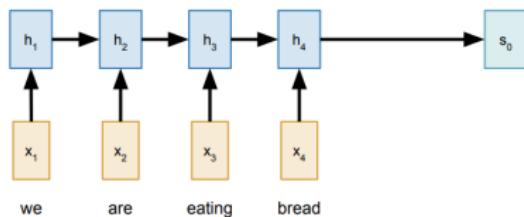
Sutskever et al., "Sequence to sequence learning with neural networks", NIPS 2014

# Sequence to Sequence with RNNs and Attention

**Input:** Sequence  $x_1, \dots, x_T$

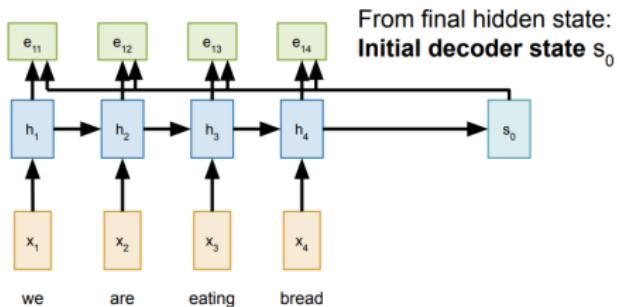
**Output:** Sequence  $y_1, \dots, y_T$

**Encoder:**  $h_t = f_W(x_t, h_{t-1})$  From final hidden state:  
**Initial decoder state**  $s_0$

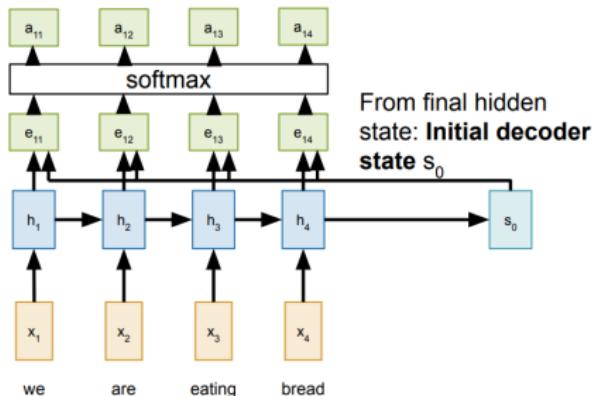


# Sequence to Sequence with RNNs and Attention

Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{att}(s_{t-1}, h_i)$       ( $f_{att}$  is an MLP)



# Sequence to Sequence with RNNs and Attention

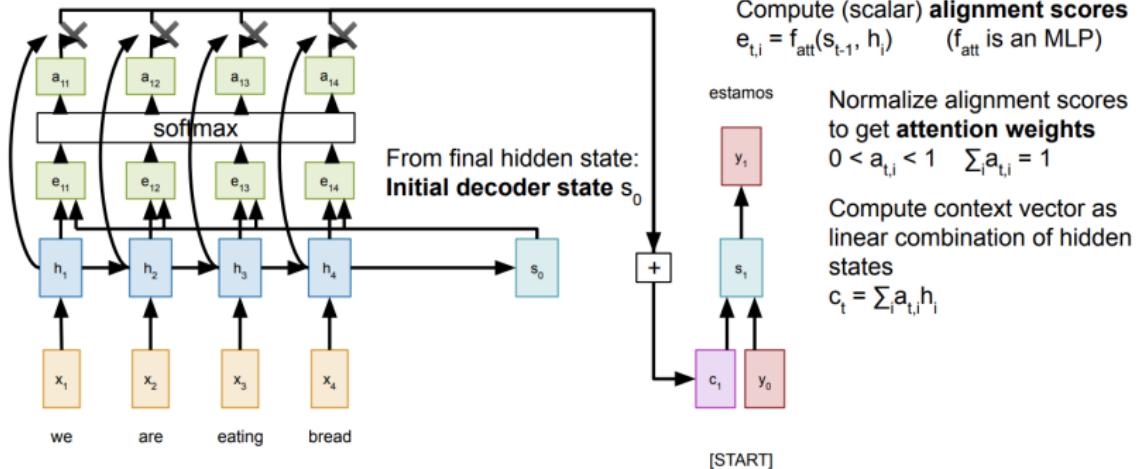


From final hidden state:  
**Initial decoder state  $s_0$**

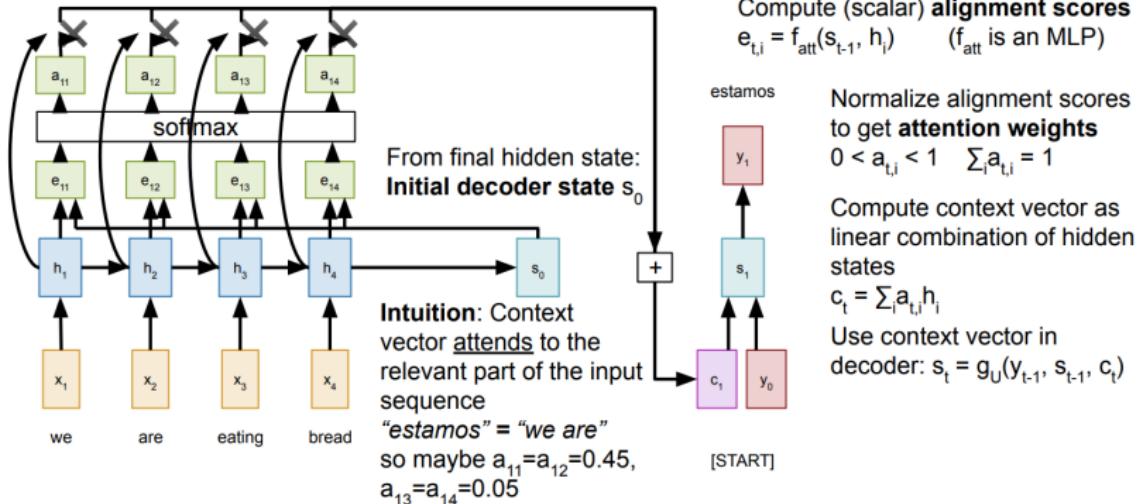
Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{att}(s_{t-1}, h_i)$  ( $f_{att}$  is an MLP)

Normalize alignment scores  
to get **attention weights**  
 $0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$

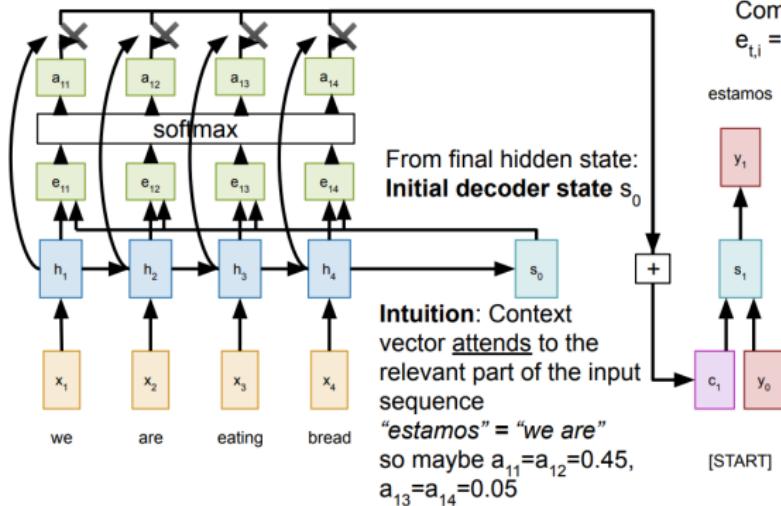
# Sequence to Sequence with RNNs and Attention



# Sequence to Sequence with RNNs and Attention



# Sequence to Sequence with RNNs and Attention



Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{att}(s_{t-1}, h_i)$  ( $f_{att}$  is an MLP)

estamos

Normalize alignment scores to get **attention weights**  
 $0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$

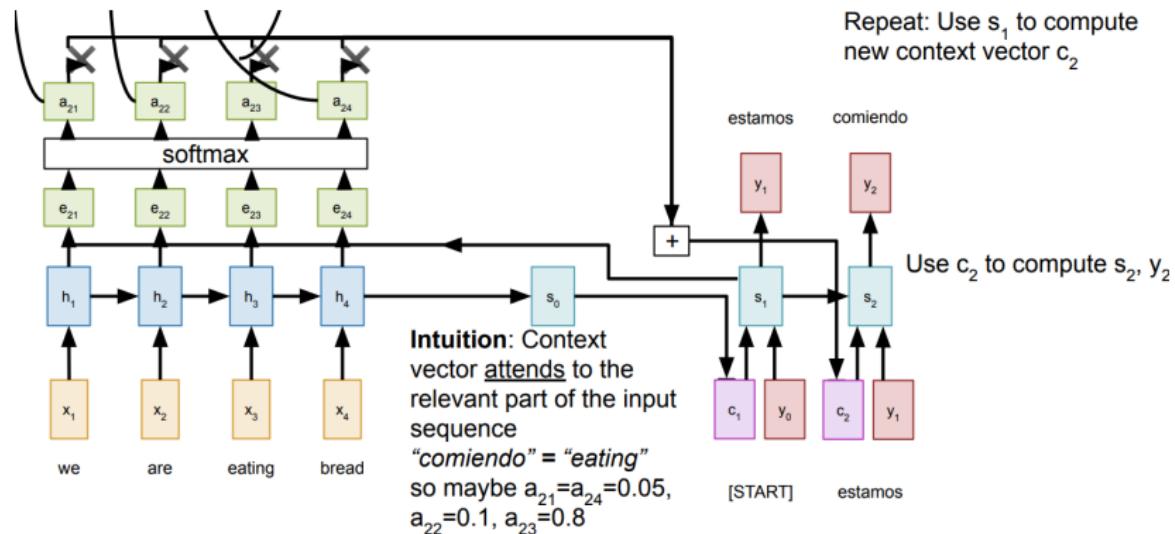
Compute context vector as linear combination of hidden states

$$c_t = \sum_i a_{t,i} h_i$$

Use context vector in decoder:  $s_t = g_0(y_{t-1}, s_{t-1}, c_t)$

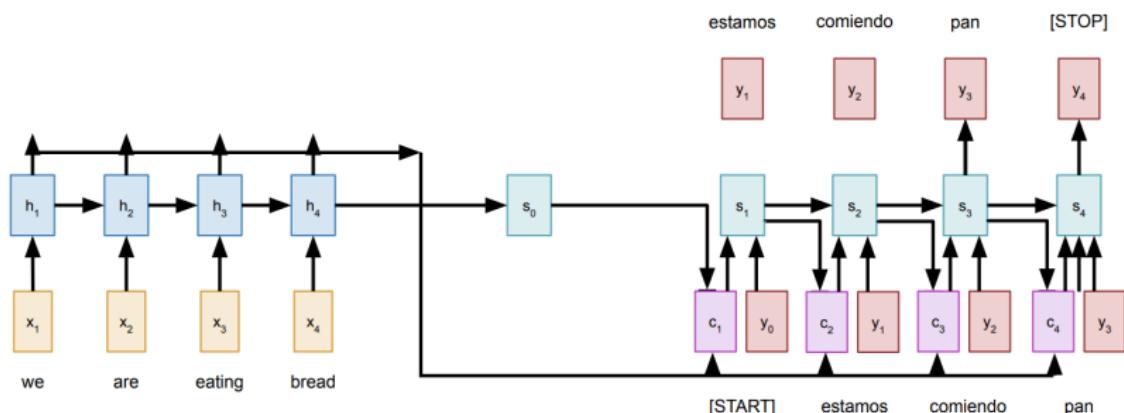
This is all differentiable! No supervision on attention weights – backprop through everything

# Sequence to Sequence with RNNs and Attention



# Sequence to Sequence with RNNs and Attention

- ▶ Use a 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



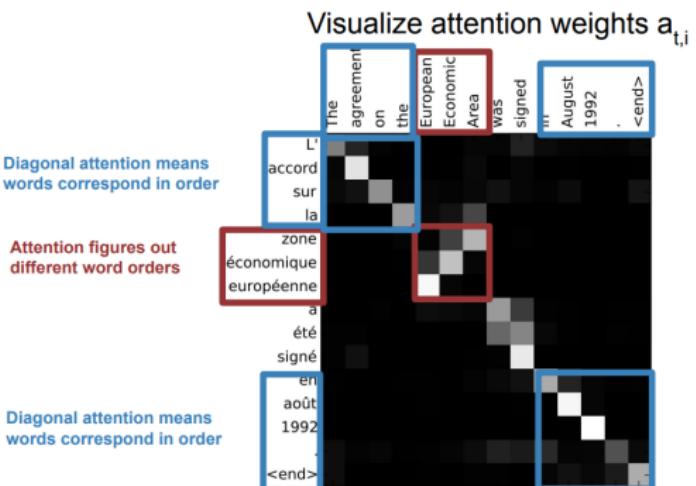
# Sequence to Sequence with RNNs and Attention

**Example:** English to French translation

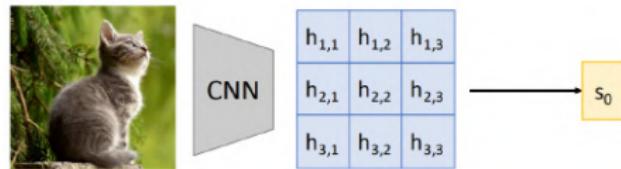
**Input:** “The agreement on the European Economic Area was signed in August 1992.”

**Output:** “L'accord sur la zone économique européenne a été signé en août 1992.”

Bahdanau et al. “Neural machine translation by jointly learning to align and translate”, ICLR 2015

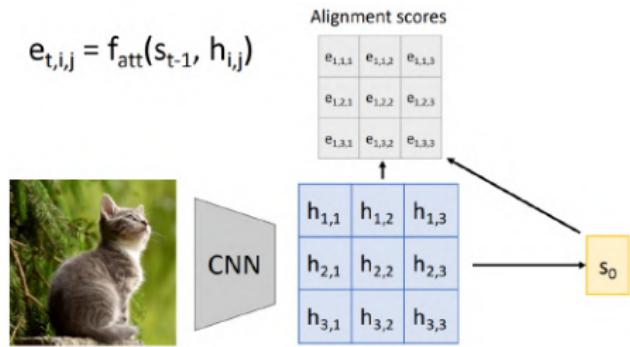


# Image Captioning with RNNs and Attention



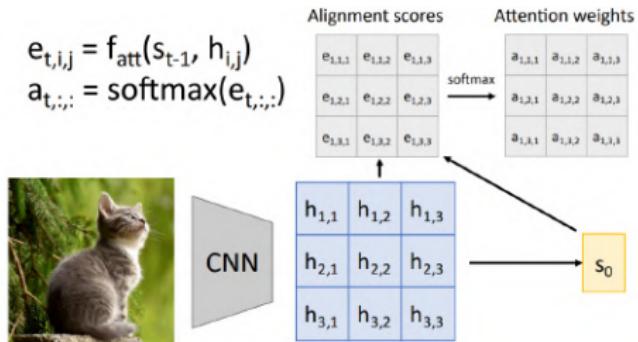
Use a CNN to compute a  
grid of features for an image

# Image Captioning with RNNs and Attention



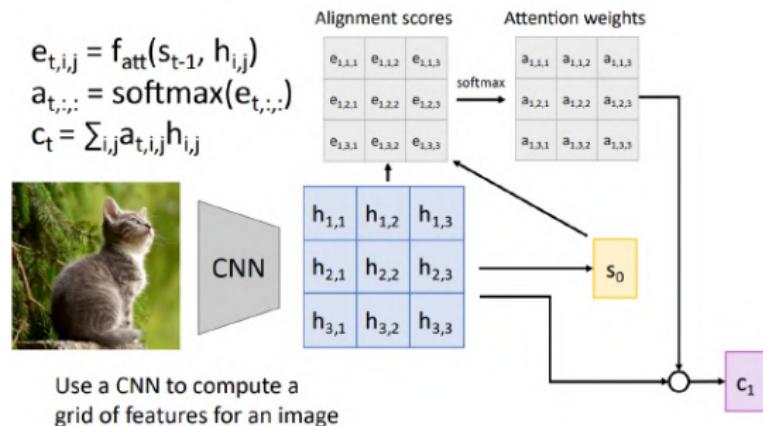
Use a CNN to compute a grid of features for an image

# Image Captioning with RNNs and Attention

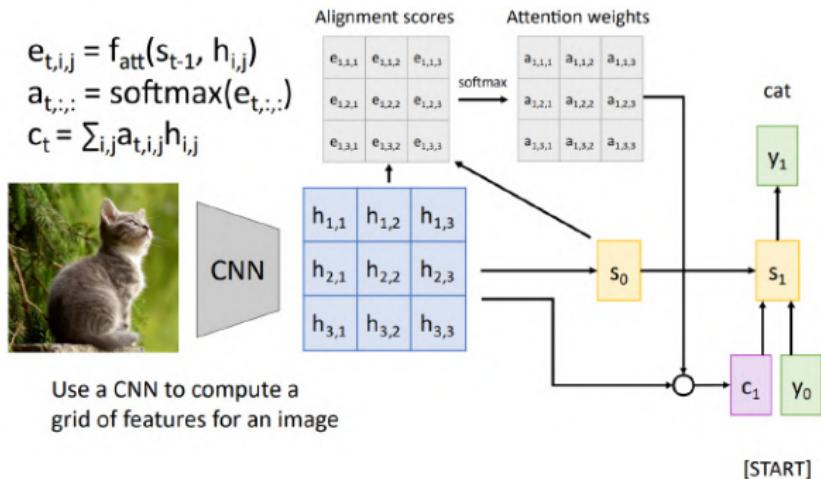


Use a CNN to compute a grid of features for an image

# Image Captioning with RNNs and Attention



# Image Captioning with RNNs and Attention



# Image Captioning with RNNs and Attention

$$e_{t,i,j} = f_{att}(s_{t-1}, h_{i,j})$$

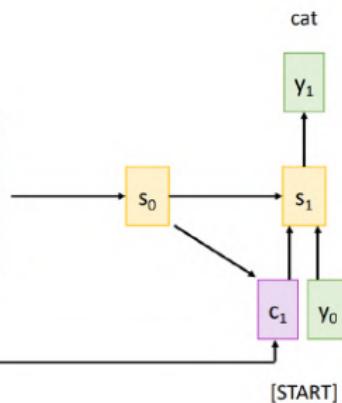
$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} h_{i,j}$$

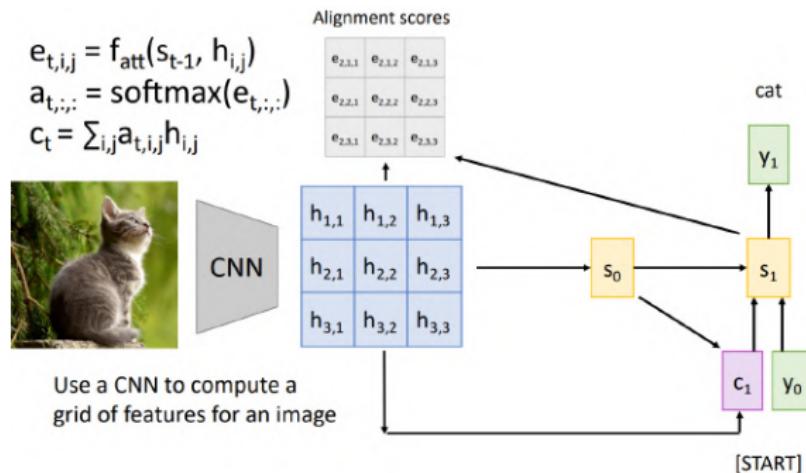


$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$

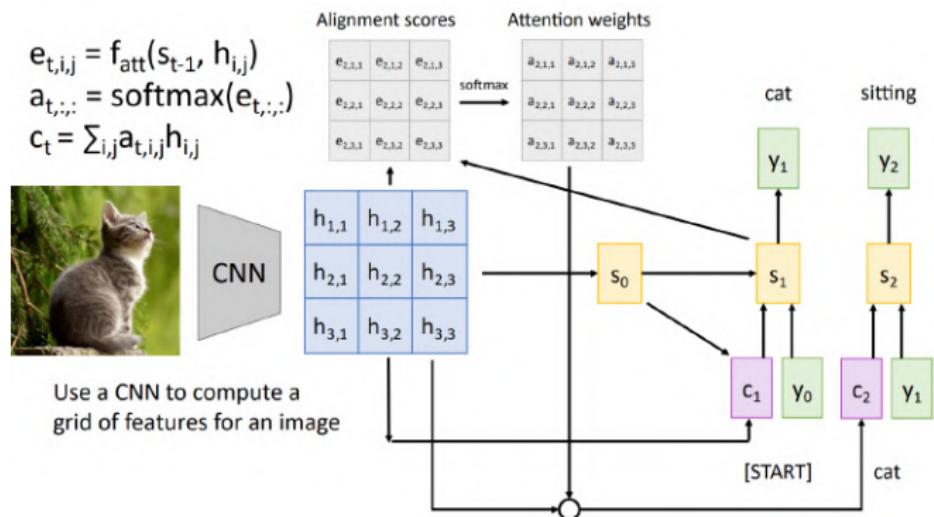
Use a CNN to compute a grid of features for an image



# Image Captioning with RNNs and Attention



# Image Captioning with RNNs and Attention



# Image Captioning with RNNs and Attention

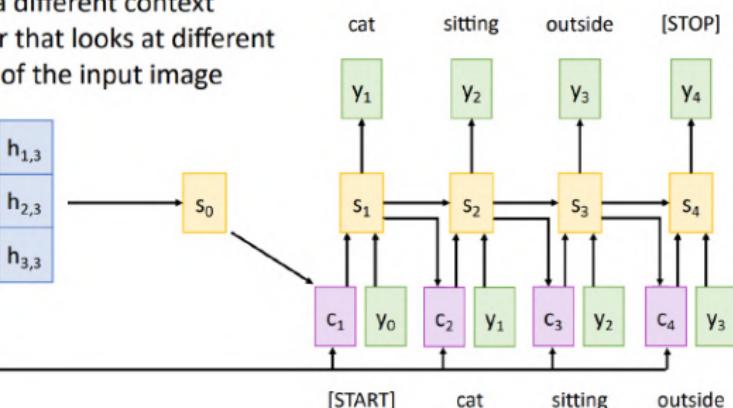
$$\begin{aligned} e_{t,i,j} &= f_{\text{att}}(s_{t-1}, h_{i,j}) \\ a_{t,:,:} &= \text{softmax}(e_{t,:,:}) \\ c_t &= \sum_{i,j} a_{t,i,j} h_{i,j} \end{aligned}$$

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



Use a CNN to compute a grid of features for an image

$h_{1,1}$	$h_{1,2}$	$h_{1,3}$
$h_{2,1}$	$h_{2,2}$	$h_{2,3}$
$h_{3,1}$	$h_{3,2}$	$h_{3,3}$



# Image Captioning with RNNs and Attention



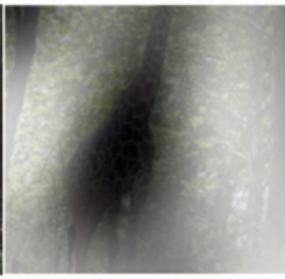
A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



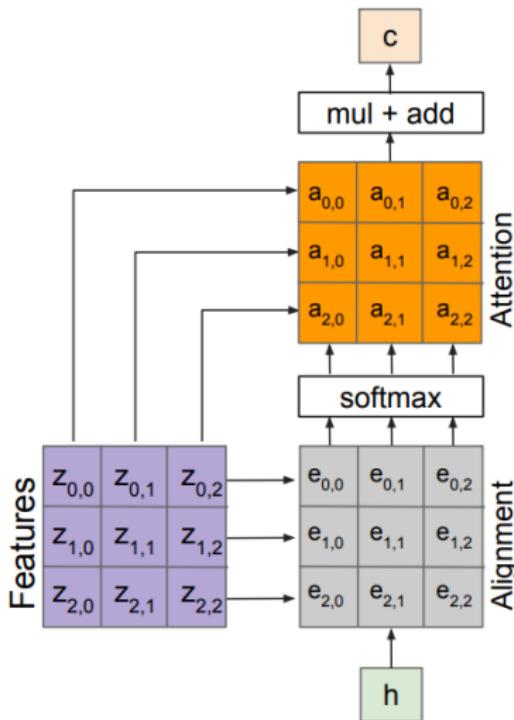
A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

<sup>0</sup>Xu et al, "Show, Ask, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Attention we just saw in image captioning



## Outputs:

context vector:  $\mathbf{c}$  (shape: D)

## Operations:

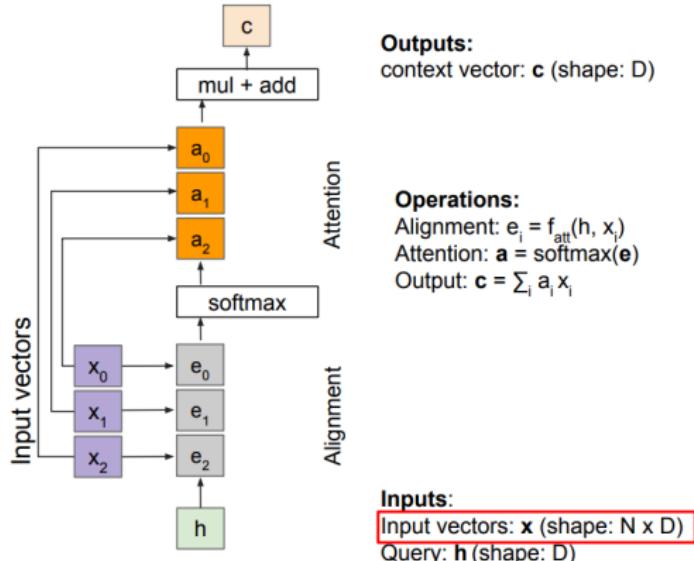
Alignment:  $e_{i,j} = f_{att}(\mathbf{h}, z_{i,j})$   
Attention:  $\mathbf{a} = \text{softmax}(\mathbf{e})$   
Output:  $\mathbf{c} = \sum_{i,j} a_{i,j} z_{i,j}$

## Inputs:

Features:  $\mathbf{z}$  (shape: H x W x D)

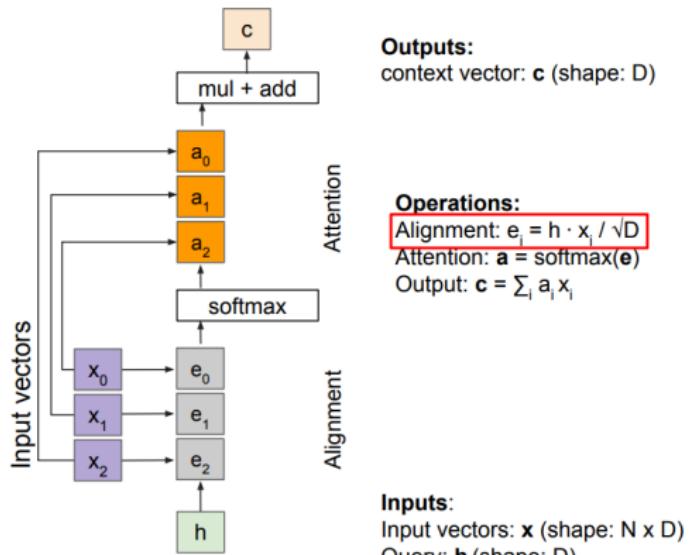
Query:  $\mathbf{h}$  (shape: D)

# General Attention Layer



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

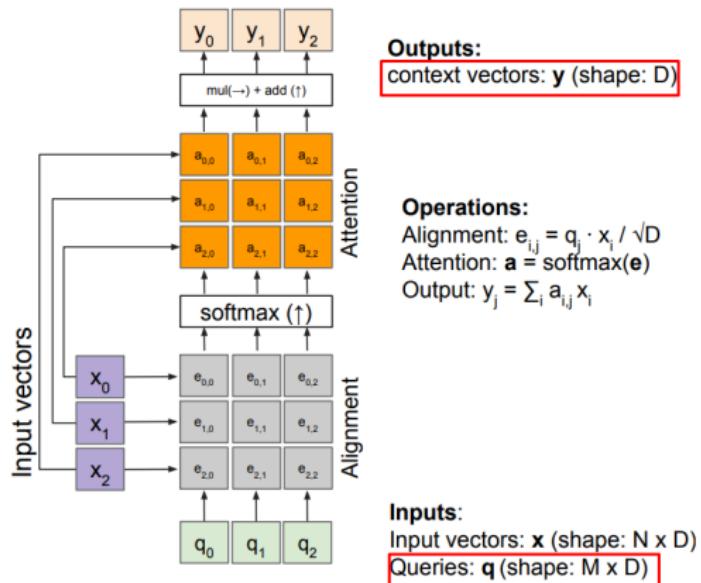
# General Attention Layer



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

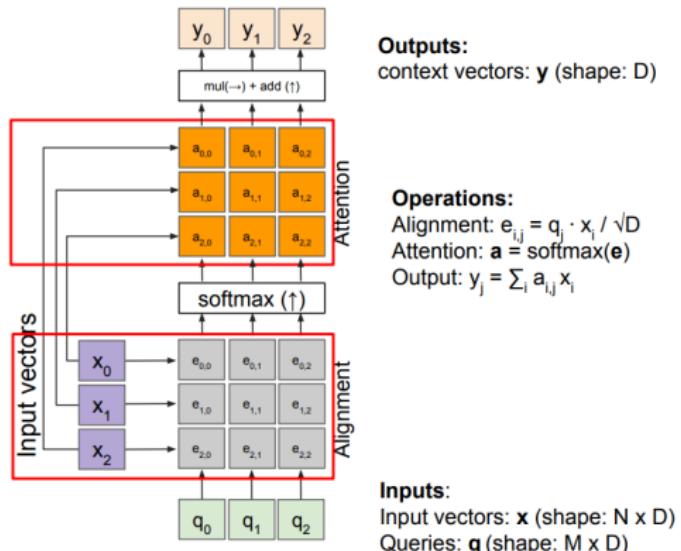
- ▶ Larger dimensions means more terms in the dot product sum.
- ▶ So, the variance of the logits is higher. Large magnitude vectors will produce much higher logits.
- ▶ So, the post-softmax distribution has lower-entropy, assuming logits are IID.
- ▶ Ultimately, these large magnitude vectors will cause 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



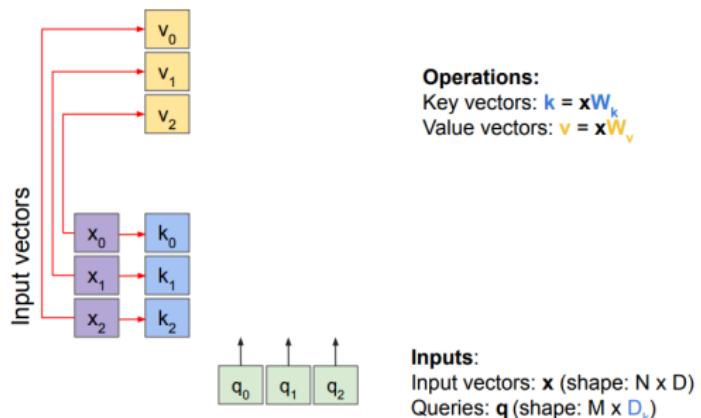
- ▶ We can have multiple Query Vectors.
- ▶ Each query creates a new output context vector

# General Attention Layer



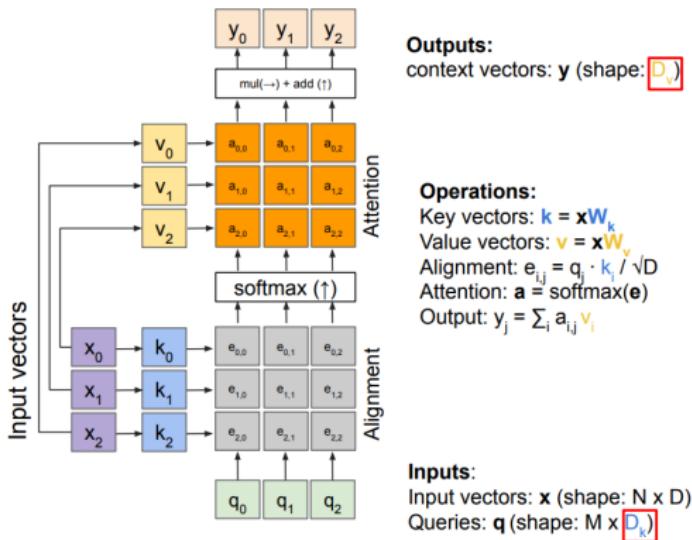
- ▶ Notice that the input vectors are used for both the alignment as well as the attention calculations.

# General Attention Layer



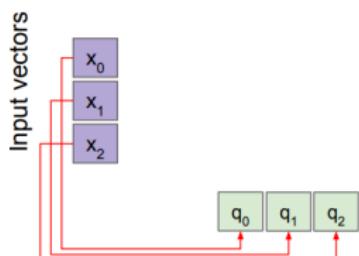
- ▶ Notice that the input vectors are used for both the alignment as well as the attention calculations.
- ▶ We can add more expressivity to the layer by adding a different FC layer before each of the two steps.

# 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:  $e_{ij} = q_i \cdot k_j / \sqrt{D}$

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

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

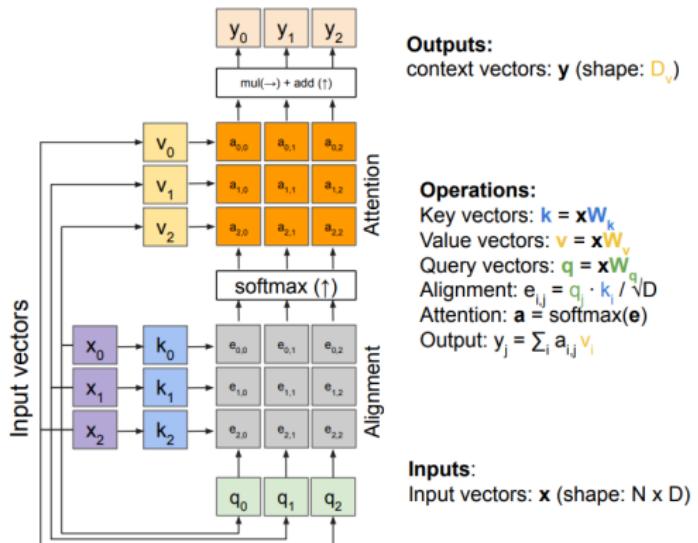
## Inputs:

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

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

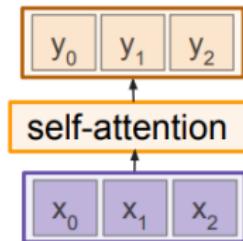
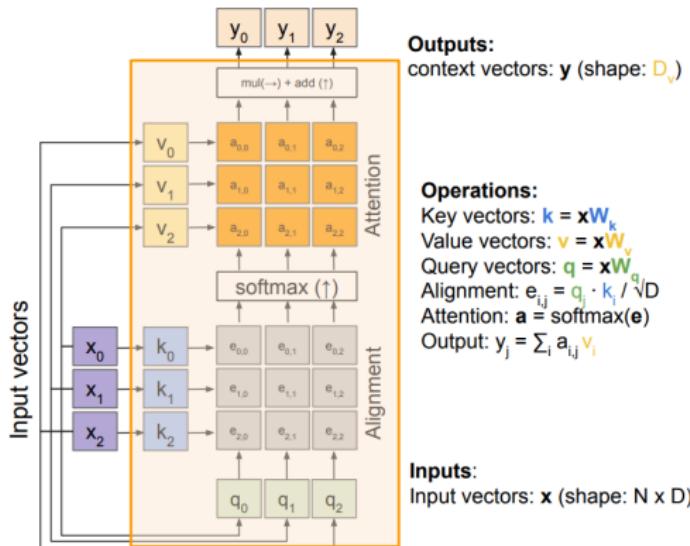
- ▶ Recall that the query vector was a function of the input vectors
- ▶ We can calculate the query vectors from the input vectors, therefore, defining a "self-attention" layer.
- ▶ No input query vectors anymore
- ▶ Instead, query vectors are calculated using a FC layer.

# Self Attention Layer

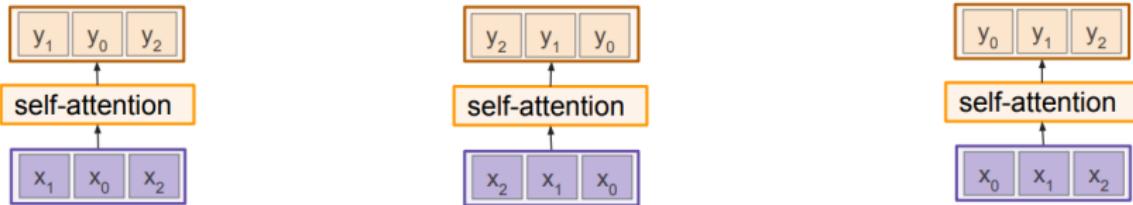


- ▶ Recall that the query vector was a function of the input vectors
- ▶ We can calculate the query vectors from the input vectors, therefore, defining a "self-attention" layer.
- ▶ No input query vectors anymore
- ▶ Instead, query vectors are calculated using a FC layer.

# Self Attention Layer



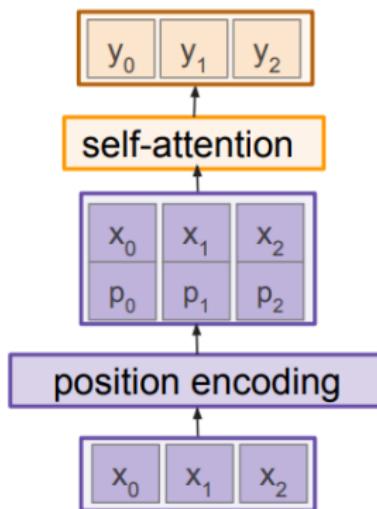
# Permutation Invariance



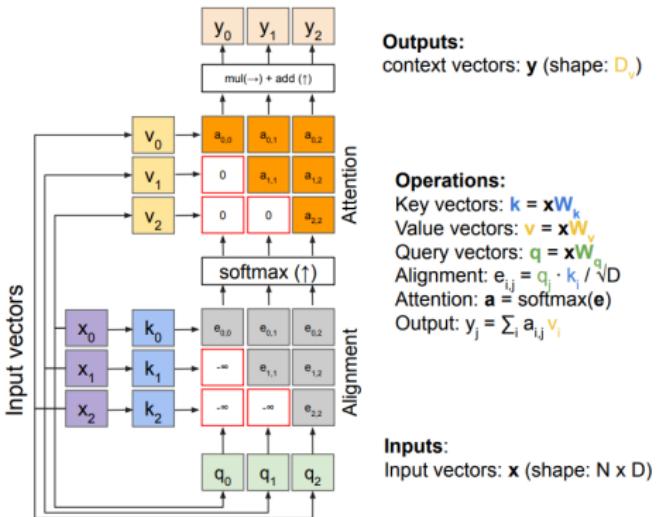
- ▶ Self-Attention Layer is Permutation equivariant
- ▶ It doesn't care about the orders of the inputs!
- ▶ **Problem:** How can we encode ordered sequences like language or spatially ordered image features?

# Positional encoding

- ▶ Concatenate/add special positional encoding  $p_j$  to each input vector  $x_j$
- ▶ 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,  $p_j = \text{pos}(j)$



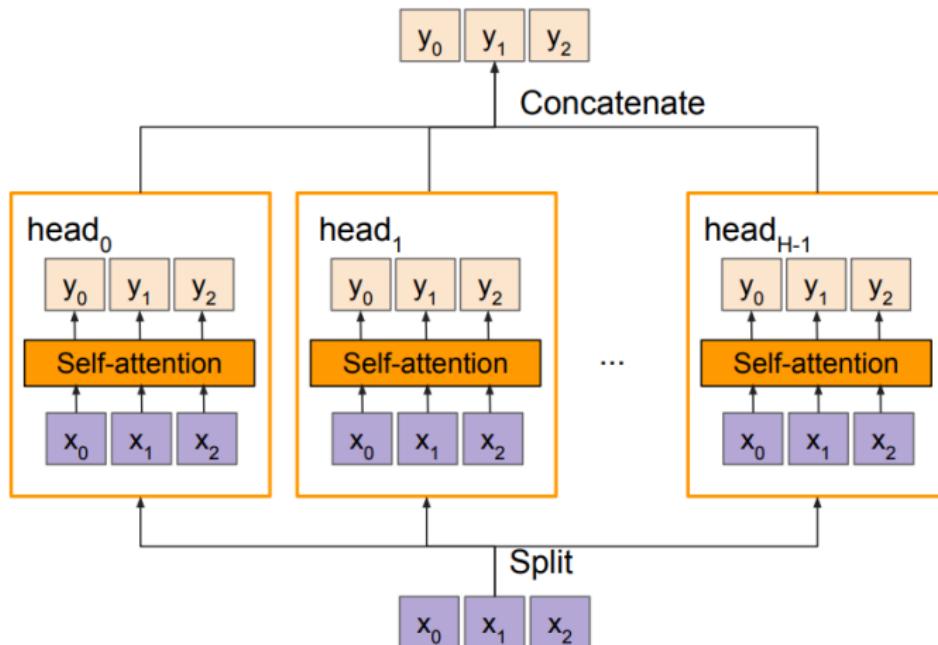
# Masked self-attention layer



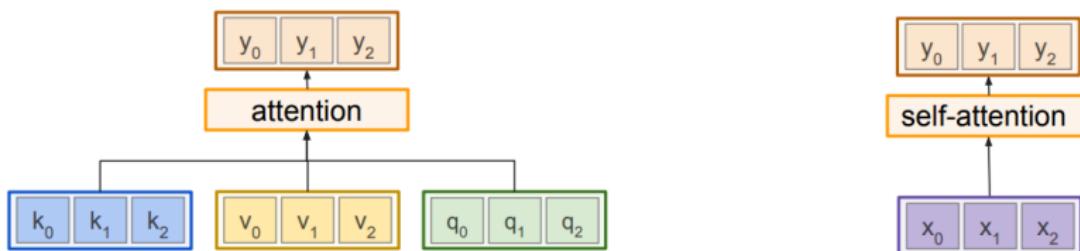
- ▶ Prevent vectors from looking at future vectors.
- ▶ Manually set alignment scores to  $-\infty$

# Multi-head self-attention layer

- ▶ Multiple self-attention heads in parallel



# General attention versus self-attention



# Example: CNN with Self-Attention

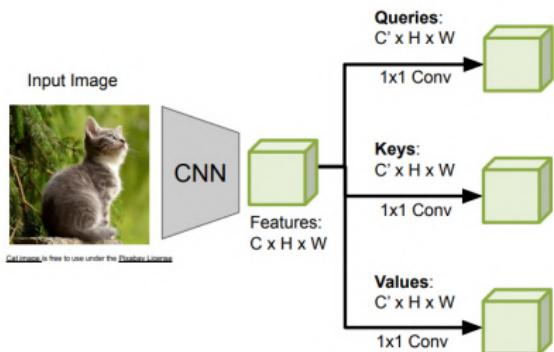
Input Image



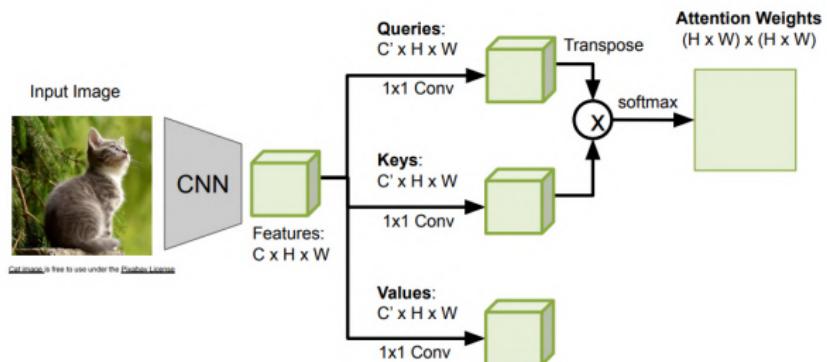
Features:  
 $C \times H \times W$

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

# Example: CNN with Self-Attention

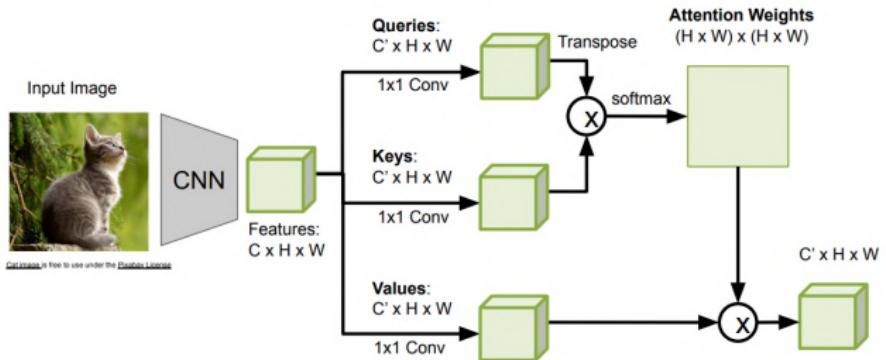


# Example: CNN with Self-Attention

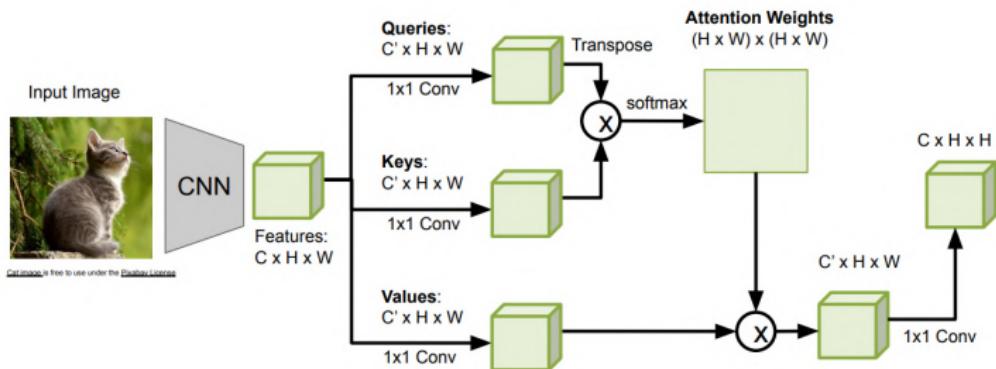


Cat image is free to use under the [Creative Commons](#)

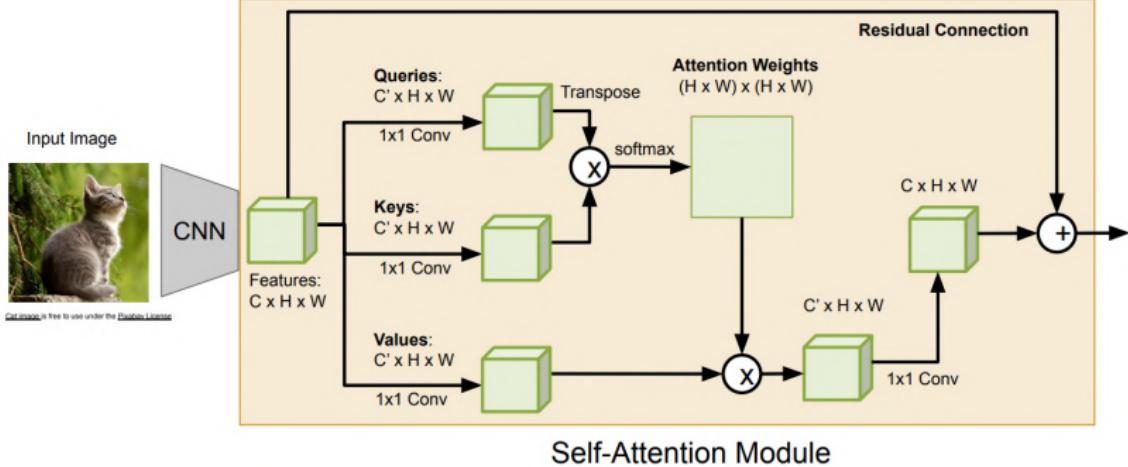
# Example: CNN with Self-Attention



# Example: CNN with Self-Attention



# Example: CNN with Self-Attention

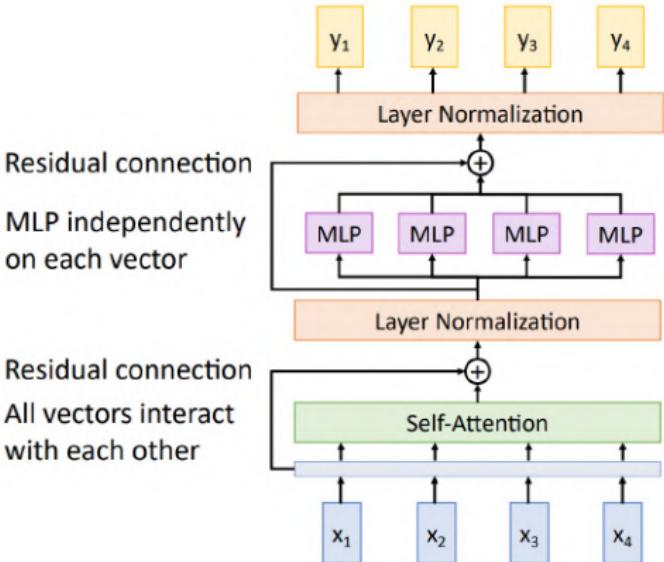


# Attention is all you need

Vaswani et al, NeurIPS 2017

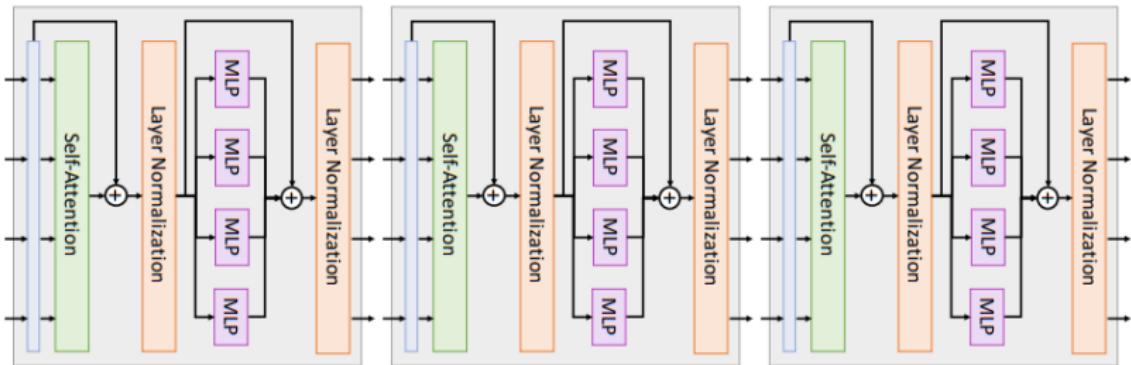
# Transformer Block

- ▶ **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

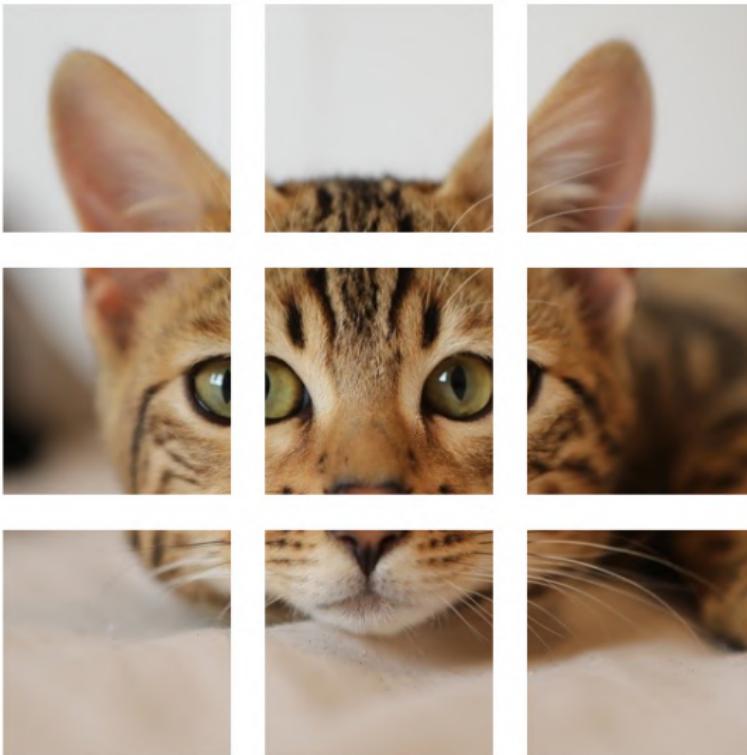


# Transformers

- ▶ A Transformer is a sequence of transformer blocks
- ▶ Vaswani et al used 12 blocks,  $D_Q = 512$  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  $3 \times 16 \times 16$

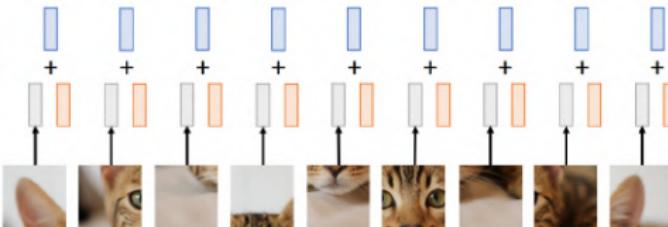


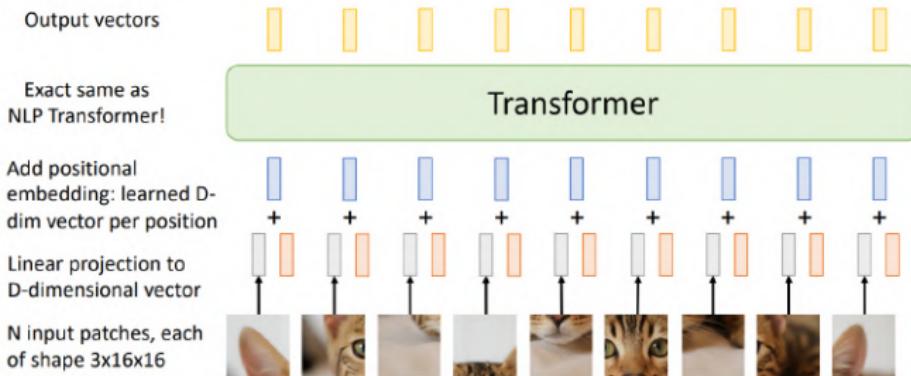
# Vision Transformers

Add positional embedding: learned D-dim vector per position

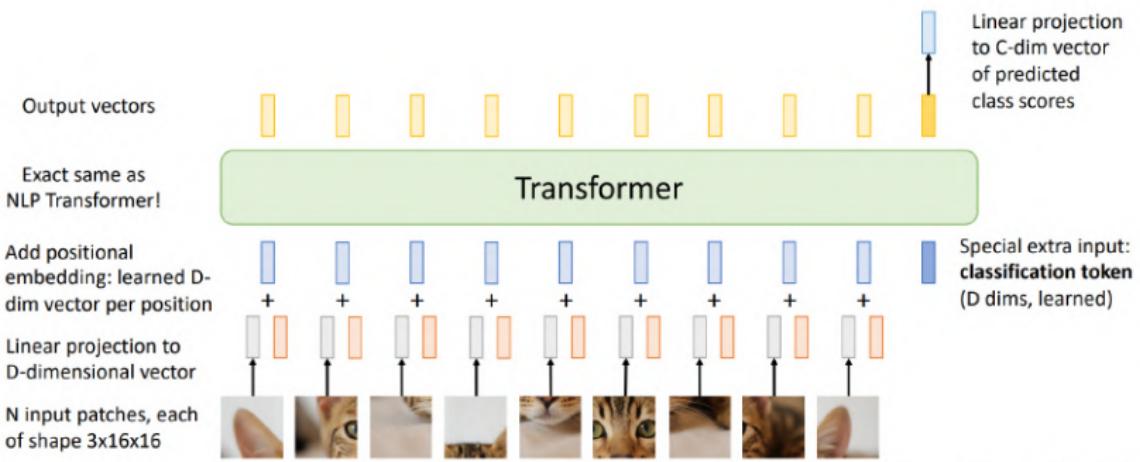
Linear projection to D-dimensional vector

N input patches, each of shape 3x16x16

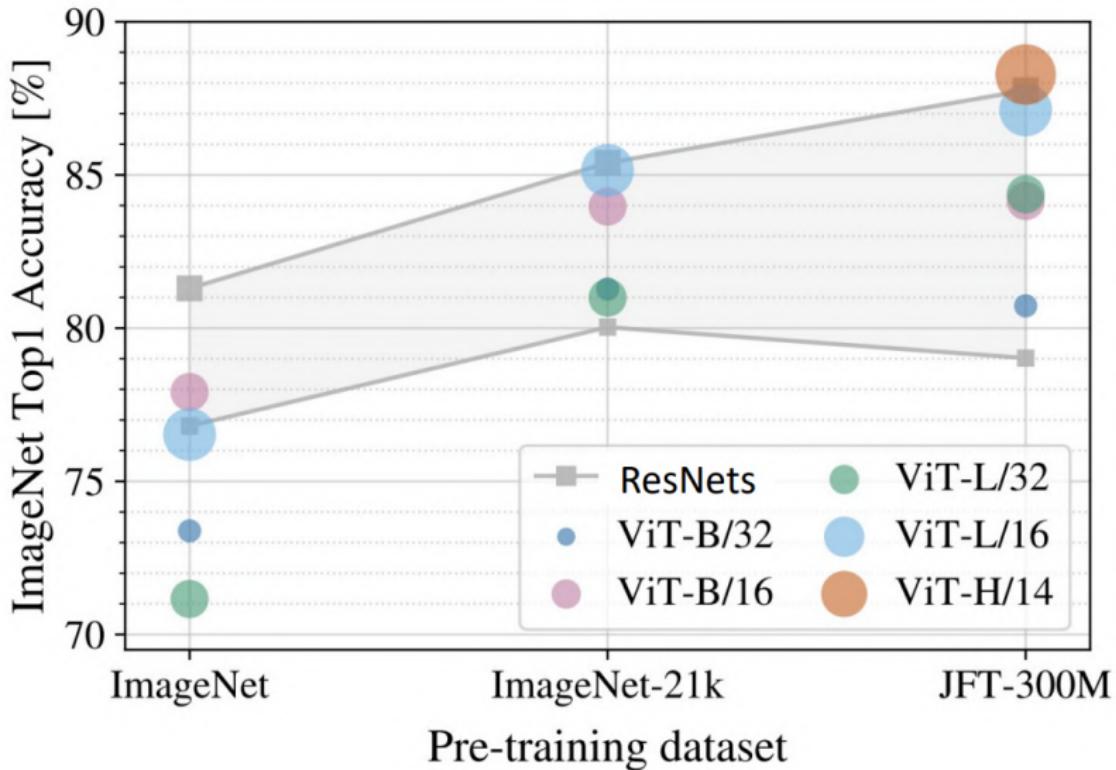




# Vision Transformers



# Vision Transformers

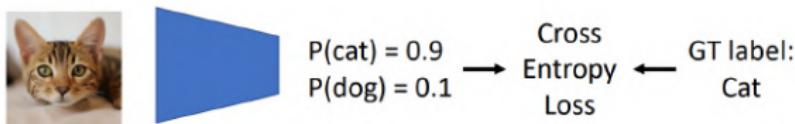


<sup>0</sup>Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

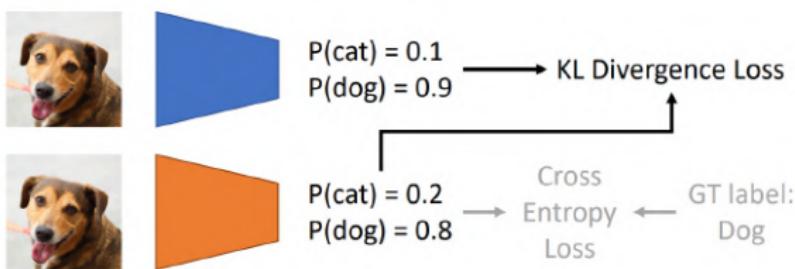
# Frame Title

# Improving ViT: Distillation

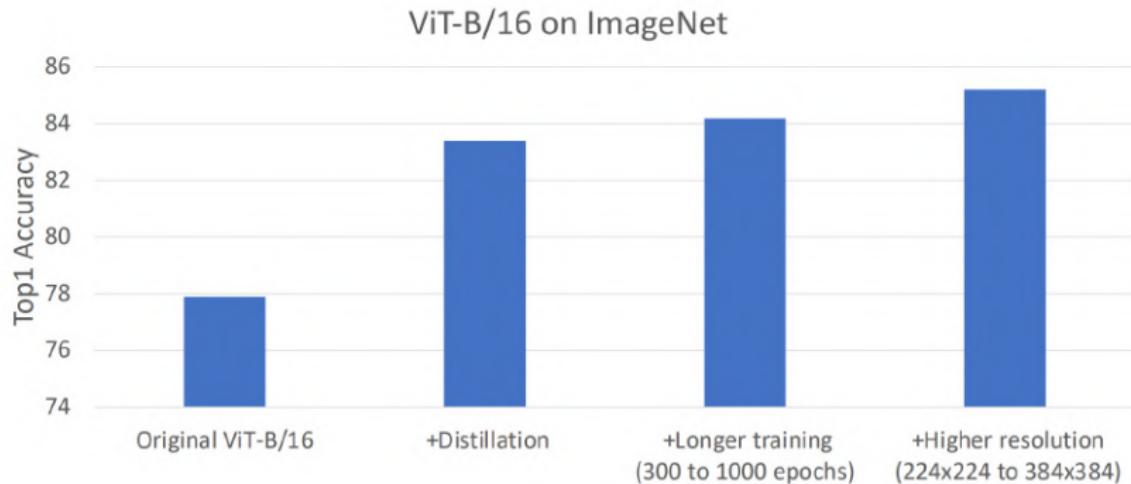
Step 1: Train a teacher CNN on ImageNet



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

## CNN



1. Maintain 2D structure logic



2. Shift equivariant



3. Consider only local correlations



4. Hierarchically growing field of view



5. Hierarchically progressing complexity



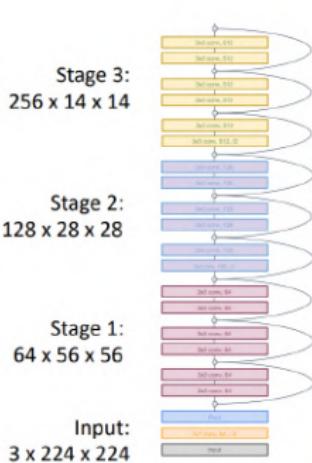
6. Reasonable amount of params



7. Global representation

## ViT



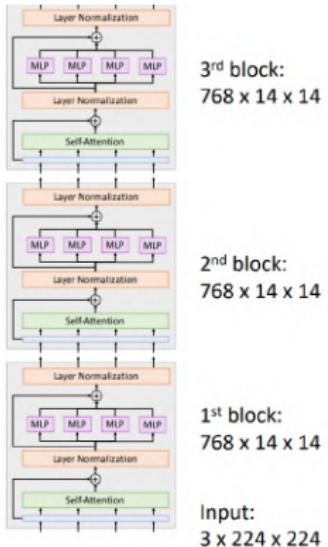


In most CNNs (including ResNets), **decrease** resolution and **increase** channels as you go deeper in the network (Hierarchical architecture)

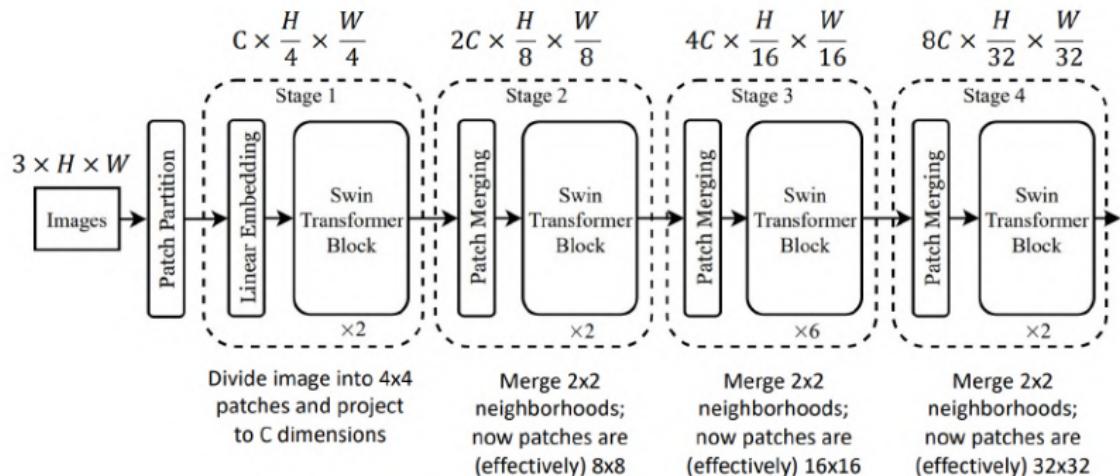
Useful since objects in images can occur at various scales

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

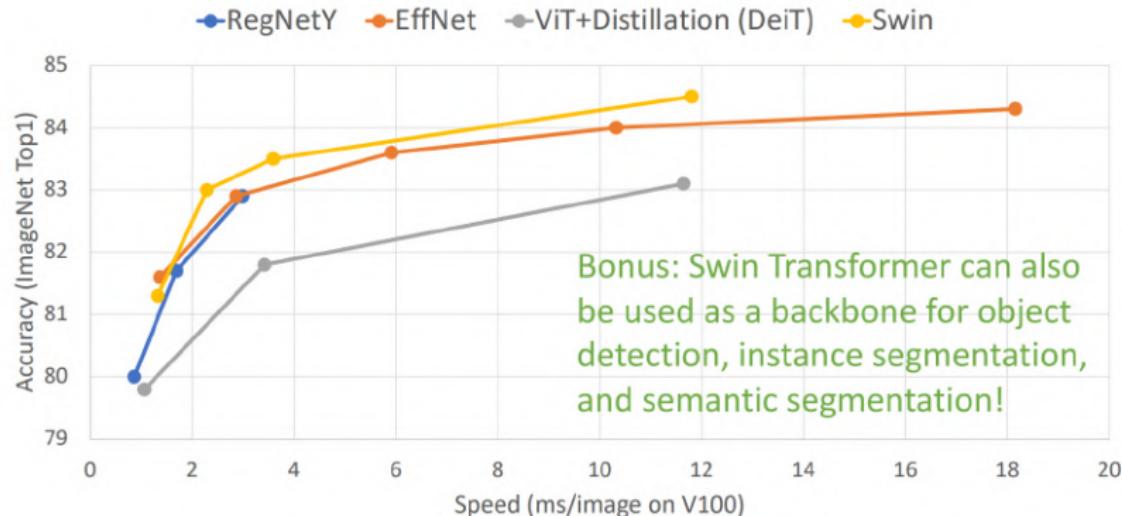
Can we build a **hierarchical** ViT model?



# Hierarchical ViT: Swin Transformer



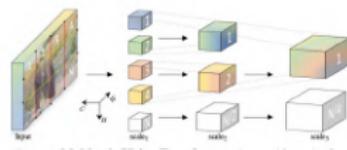
# Hierarchical ViT: Swin Transformer



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

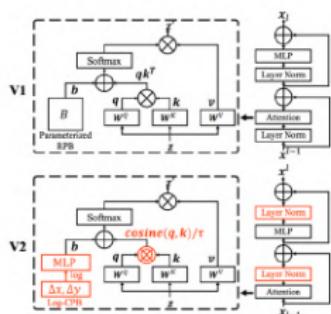
# Other Hierarchical Vision Transformers

MViT



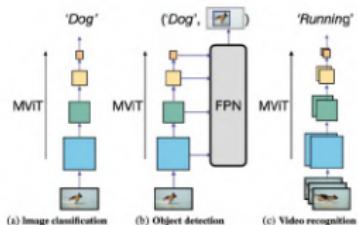
Fan et al, "Multiscale Vision  
Transformers", ICCV 2021

Swin-V2



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

Improved MViT

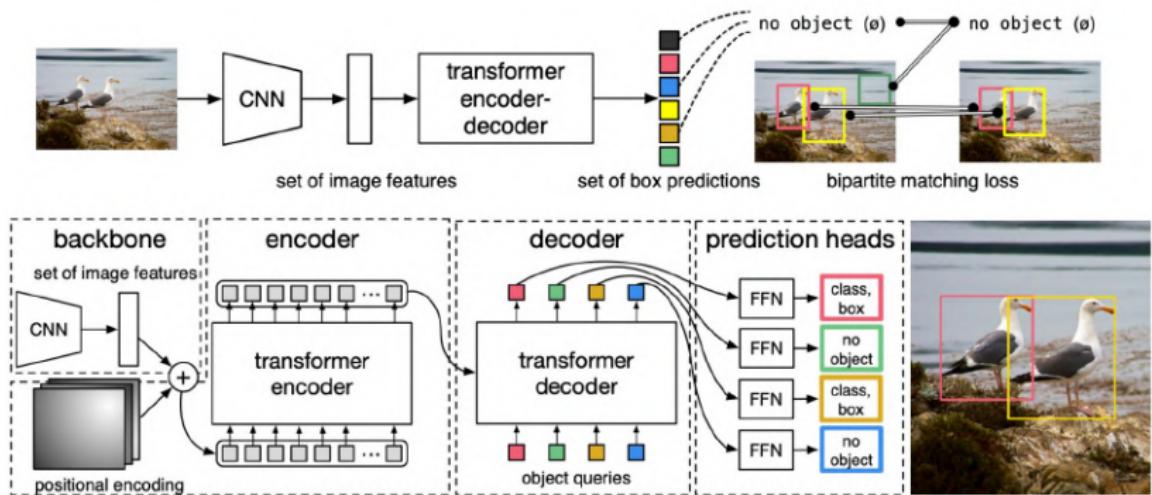


Li et al, "Improved Multiscale Vision Transformers  
for Classification and Detection", arXiv 2021

# 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



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