



# Master in Computer Vision Barcelona

---

[\[http://pagines.uab.cat/mcv/\]](http://pagines.uab.cat/mcv/)

## Module 6 - Day 7 The Transformer in Vision

28th March 2023

# Acknowledgments



Xavier Giro-i-Nieto



[@DocXavi](https://twitter.com/DocXavi)



[xavier.giro@upc.edu](mailto:xavier.giro@upc.edu)

**Applied Scientist**  
Amazon Science Barcelona.

Associate Professor  
Universitat Politècnica de Catalunya

# Outline

1. Vision Transformer (ViT)
2. Beyond ViT

# The Transformer for Vision: ViT



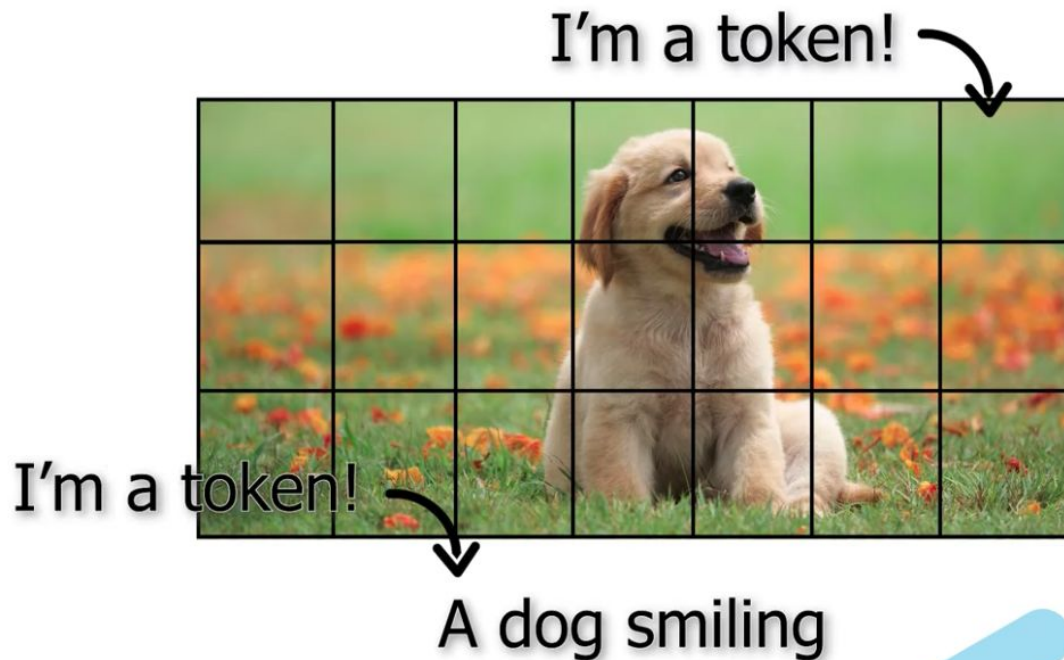
**#ViT** Dosovitskiy, Alexey, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani et al. ["An image is worth 16x16 words: Transformers for image recognition at scale."](#) ICLR 2021. [\[blog\]](#) [\[code\]](#) [\[video by Yannic Kilcher\]](#)

# Outline

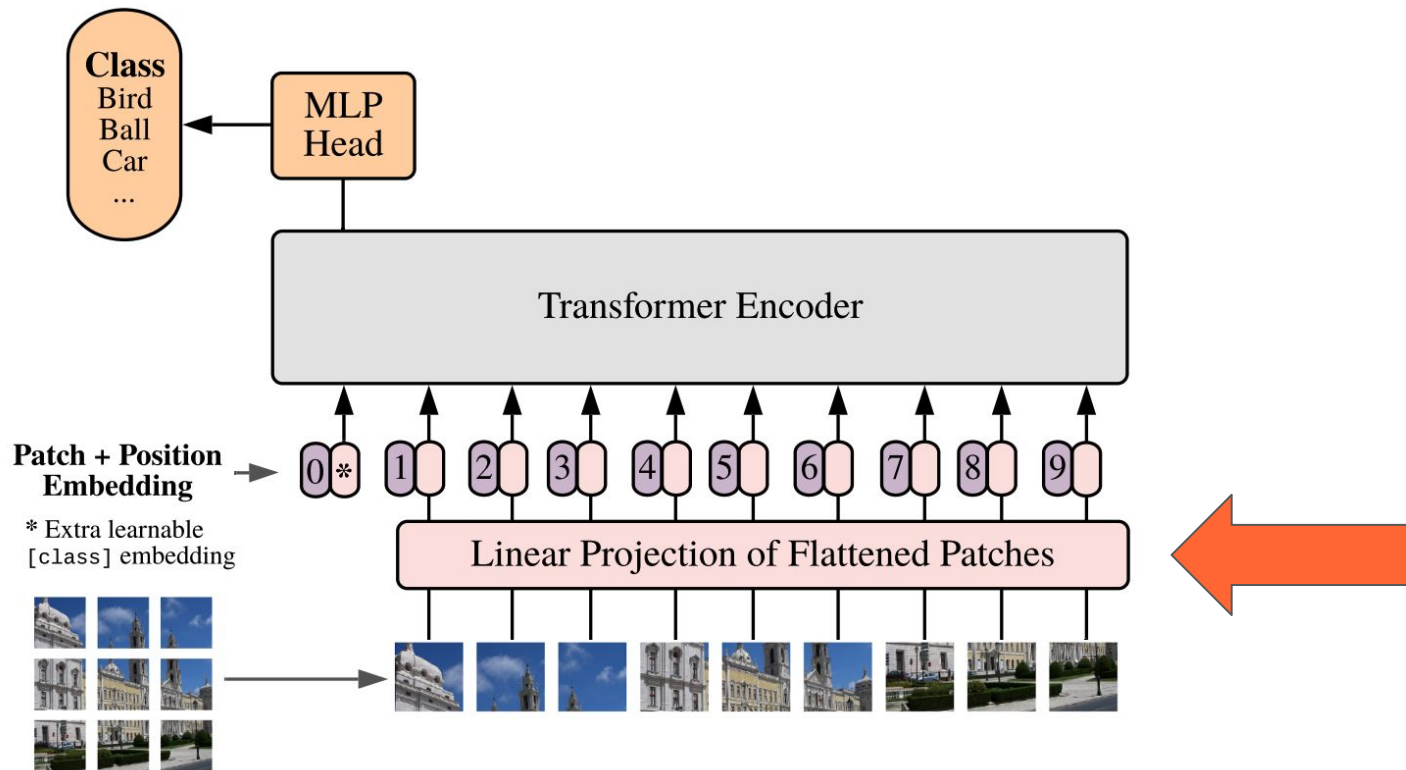
1. Vision Transformer (ViT)
  - a. Tokenization
  - b. Position embeddings**
  - c. Class embedding
  - d. Receptive field
  - e. Performance
2. Beyond ViT

# The Transformer for Vision

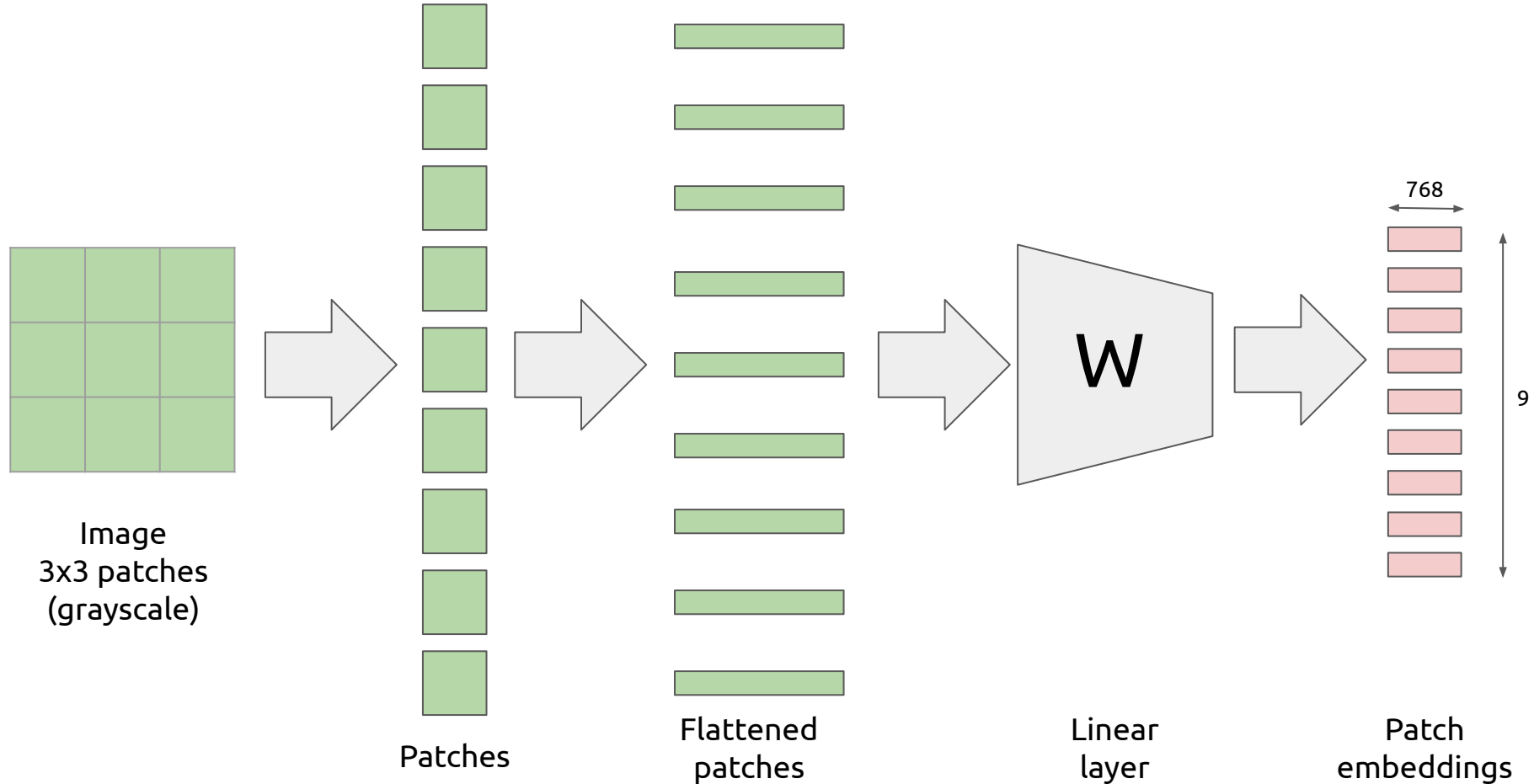
@whats\_ai



# The Transformer for Vision: ViT



# Linear projection of Flattened Patches

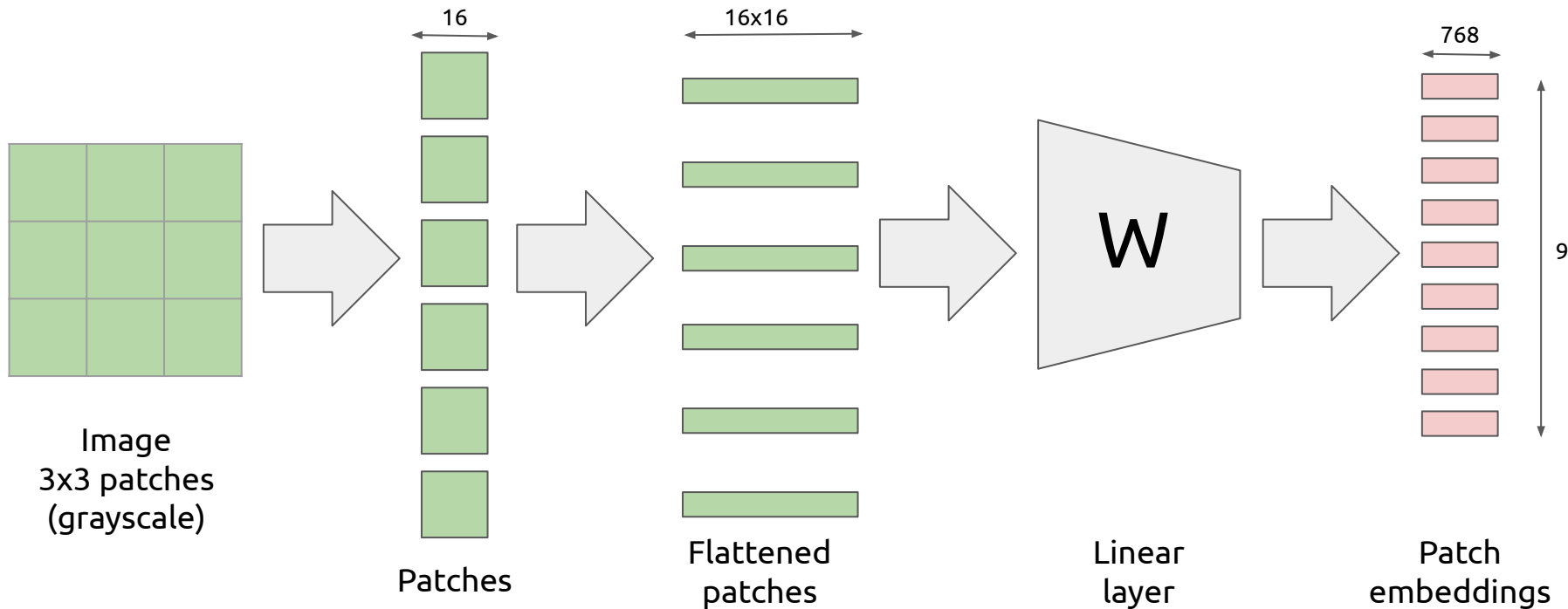




# The Transformer for Vision

...

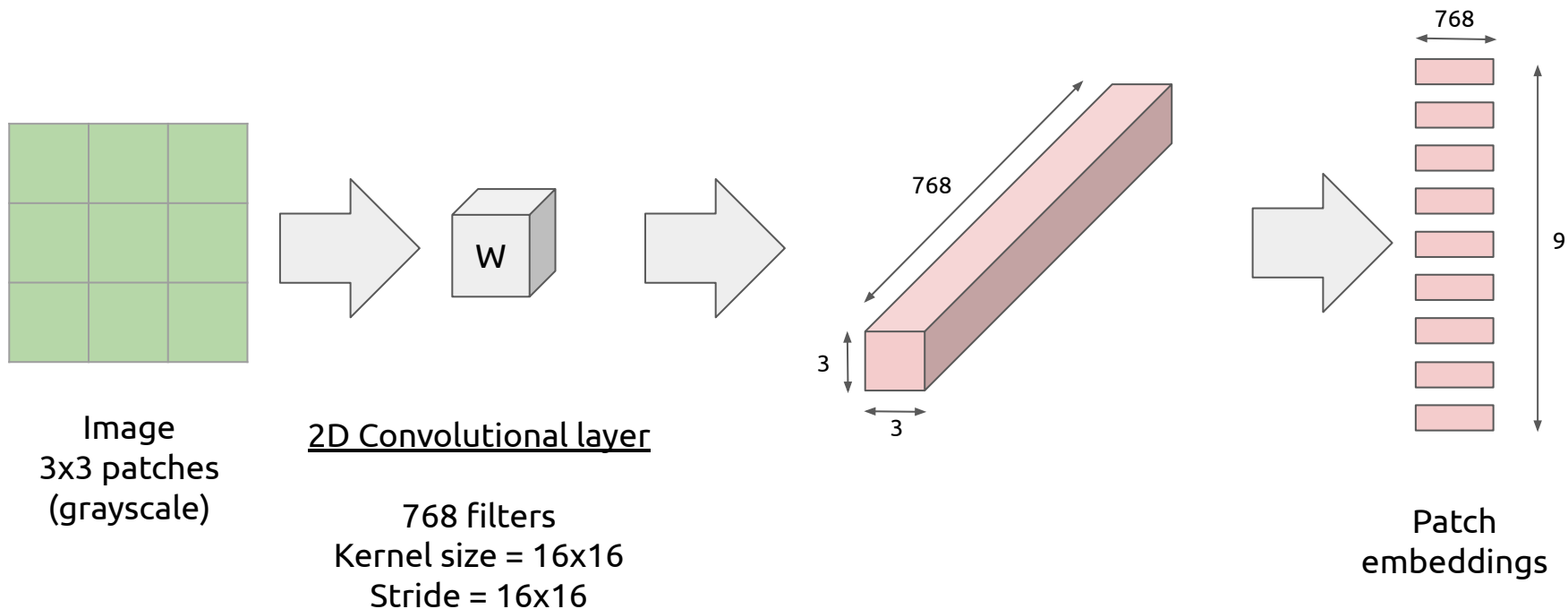
Consider the case of patches of  $16 \times 16$  pixels and their embedding size of  $D=768$ , as in ViT-Base. How could the linear layer be implemented with a convolutional layer ?



# The Transformer for Vision

...

Consider the case of patches of  $16 \times 16$  pixels and their embedding size of  $D=768$ , as in ViT-Base. How could the linear layer be implemented with a convolutional layer ?



# The Transformer for Vision



**Yann LeCun**

@ylecun

...

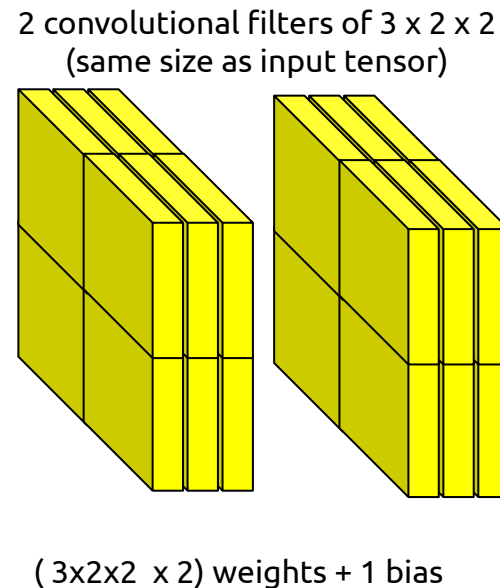
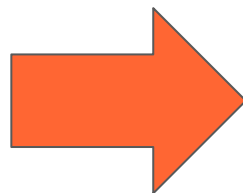
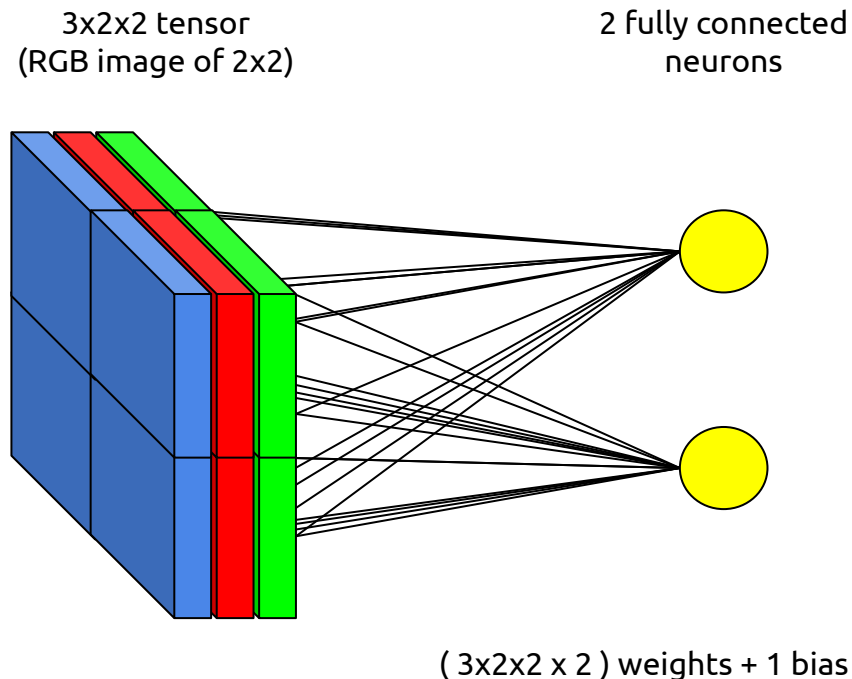
Wondering why the first layer of some recent DL architectures for vision are called  
"linear embedding of 16x16 non-overlapping patches"  
instead of  
"Convolutional layer with 16x16 kernels and 16x16 stride"  
???

[Tradueix el tuit](#)

11:32 p. m. · 6 de maig de 2021 · Twitter for Android

# The Transformer for Vision

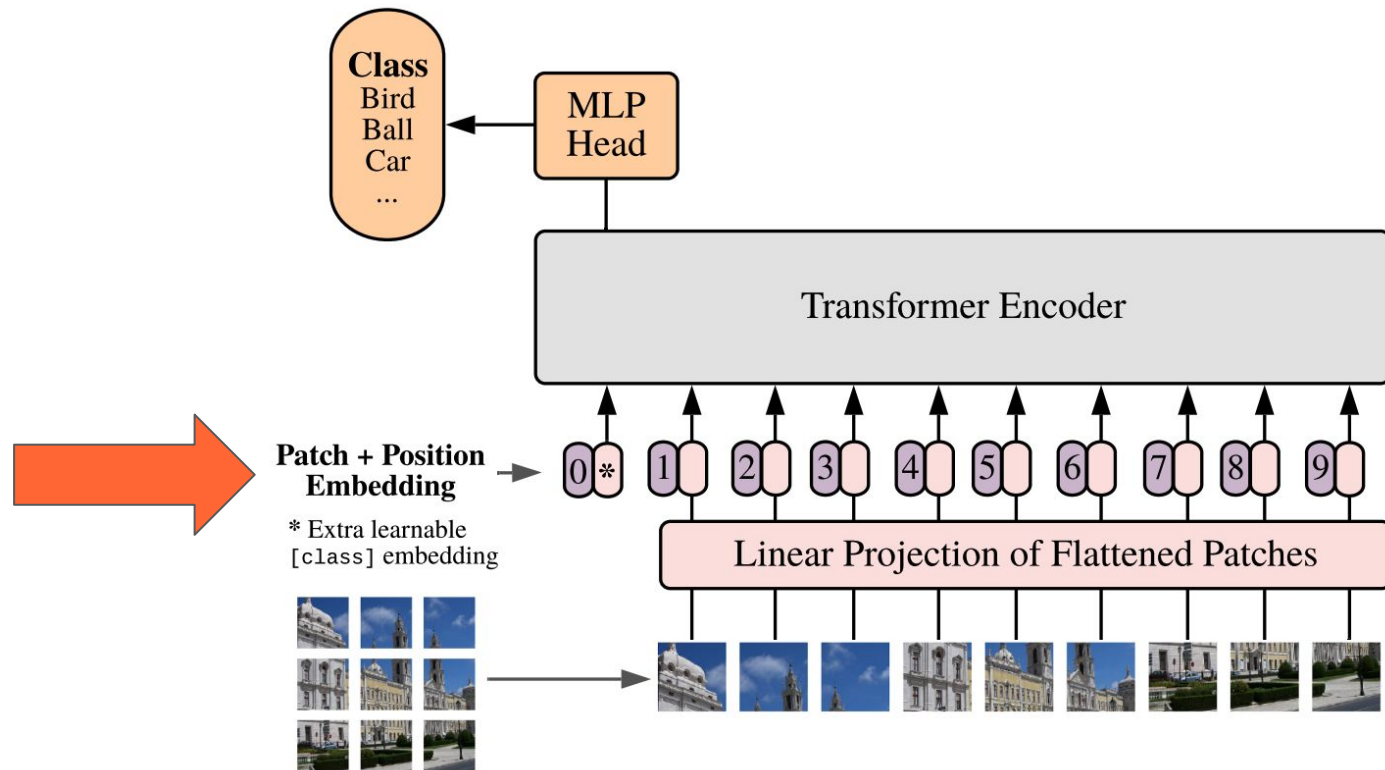
Observation: Fully connected neurons could be implemented as convolutional ones.



# Outline

1. Vision Transformer (ViT)
  - a. Tokenization
  - b. Position embeddings**
  - c. Class embedding
  - d. Receptive field
  - e. Performance

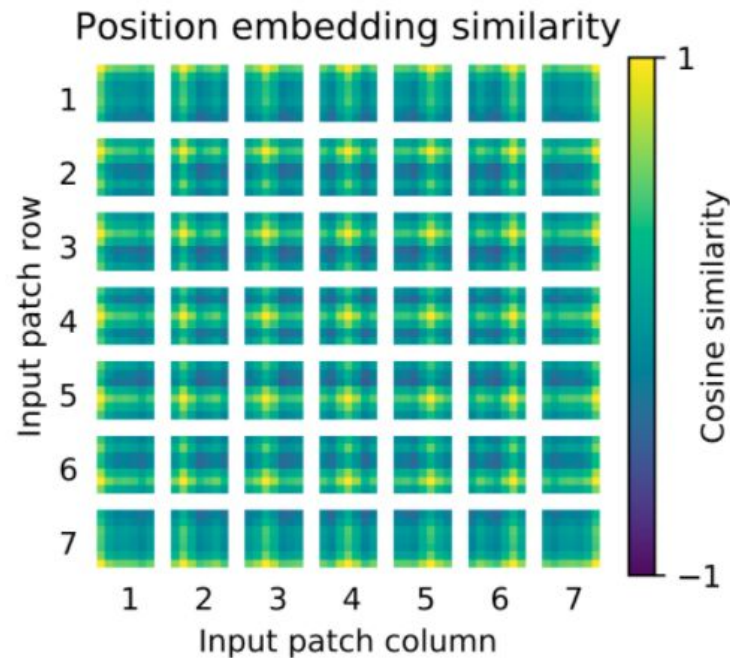
# Position Embeddings



# Position embeddings

The model **learns** to encode the relative position between patches.

Each position embedding is most similar to others in the same row and column, indicating that the model has recovered the grid structure of the original images.



# Outline

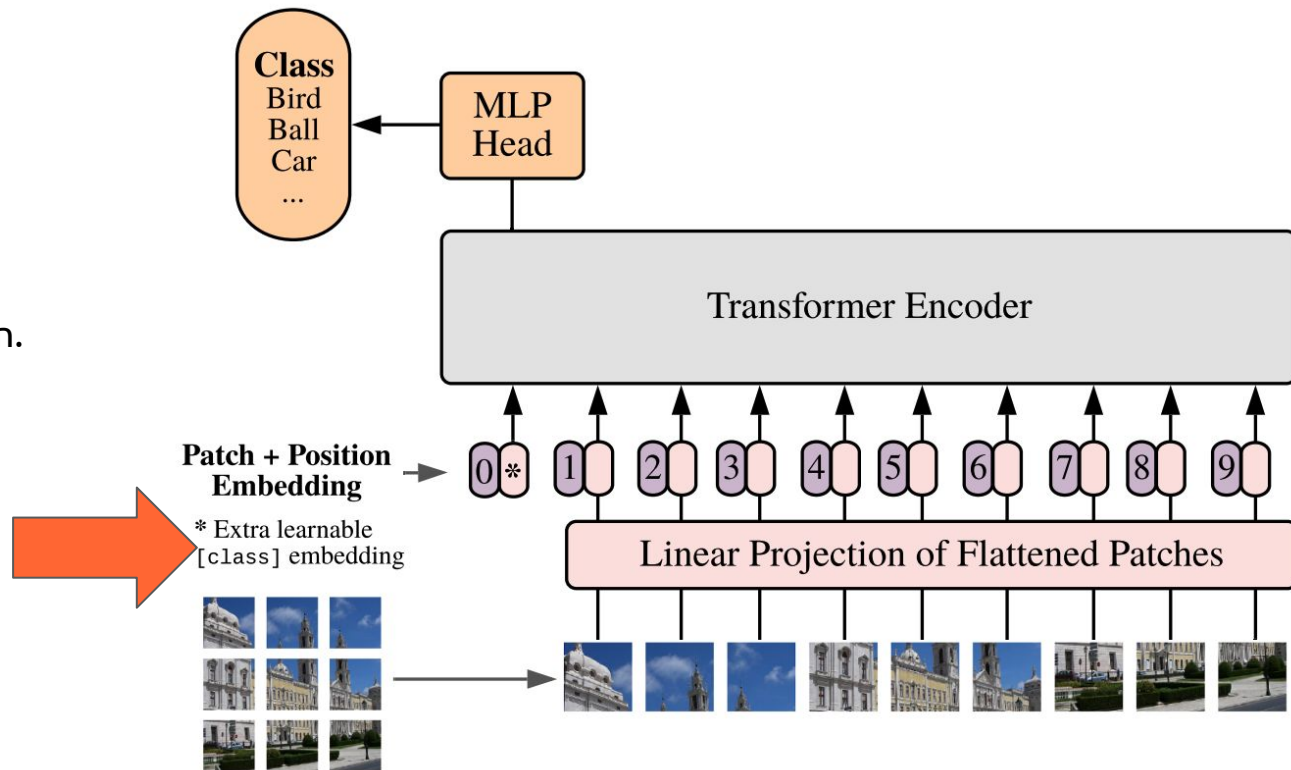
1. Vision Transformer (ViT)
  - a. Tokenization
  - b. Position embeddings
  - c. **Class embedding**
  - d. Receptive field
  - e. Performance



# Class embedding

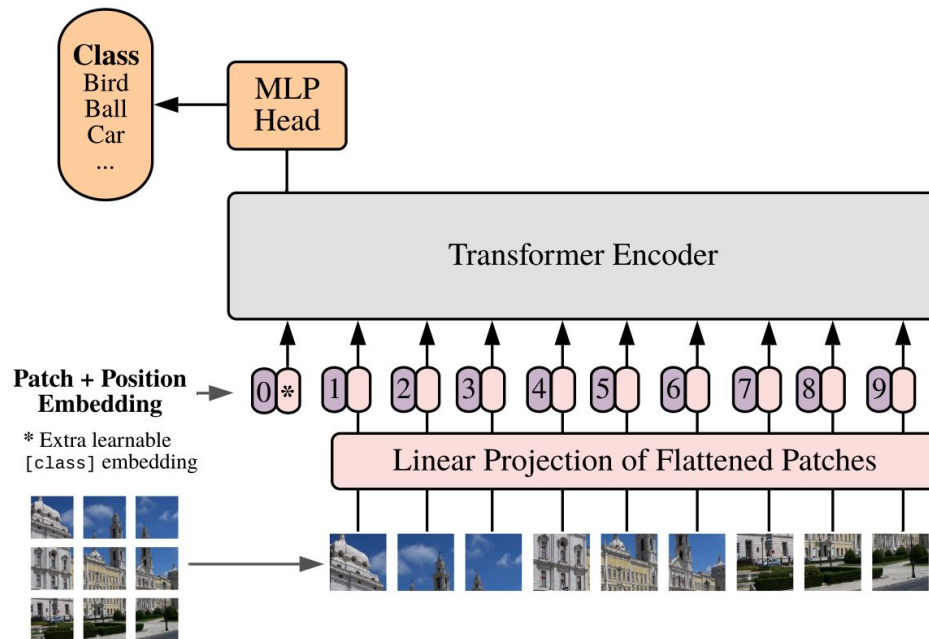
**[class]** is a special learnable embedding added in front of every input example.

It triggers the class prediction.



# Class embedding

Why does the ViT not have a decoder in its architecture ?

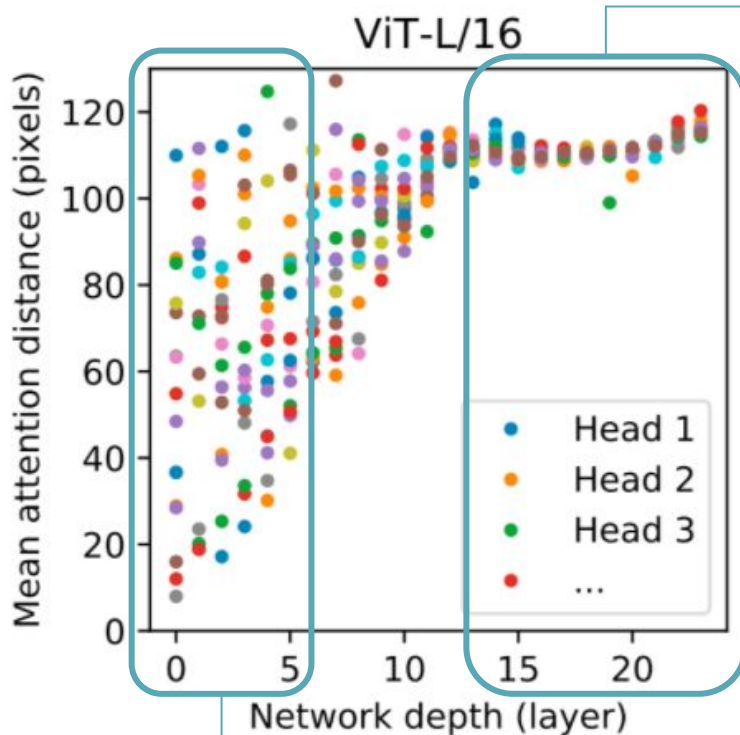


# Outline

1. Vision Transformer (ViT)
  - a. Tokenization
  - b. Position embeddings
  - c. Class embedding
  - d. Receptive field**
  - e. Performance

# Receptive field

Average spatial distance between one element attending to another for each transformer block:



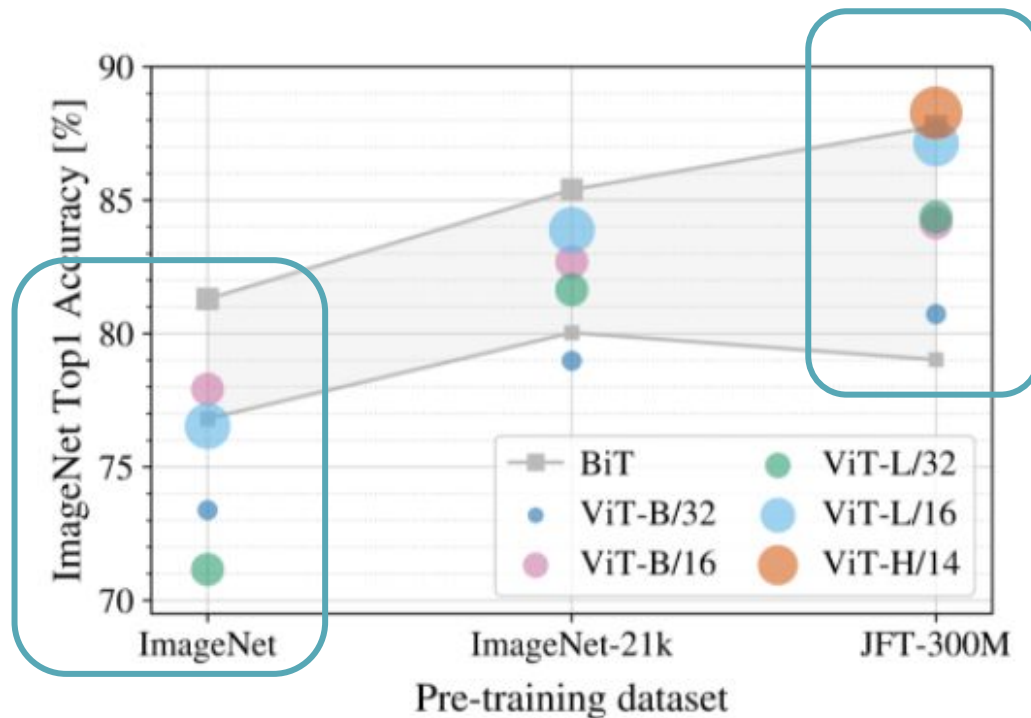
Both short & wide attention ranges in early layers (CNN can only learn short ranges)

Deeper layers attend all over the image.

# Outline

1. Vision Transformer (ViT)
  - a. Tokenization
  - b. Position embeddings
  - c. Class embedding
  - d. Receptive field
  - e. **Performance**

# Performance: Accuracy



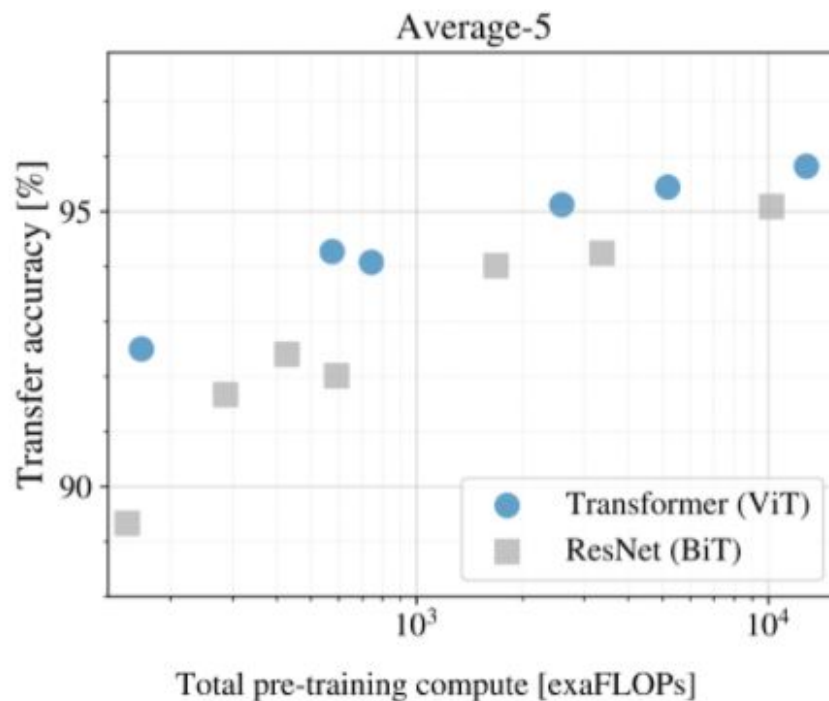
Worse performance than CNN (BiT) with ImageNet data only.

Slight improvement over CNN (BiT) when very large amounts of training data available.

**#BiT** Kolesnikov, Alexander, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. ["Big transfer \(bit\): General visual representation learning."](#) ECCV 2020.

**#ViT** Dosovitskiy, Alexey, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani et al. ["An image is worth 16x16 words: Transformers for image recognition at scale."](#) ICLR 2021. [\[blog\]](#) [\[code\]](#) [\[video by Yannic Kilcher\]](#)

# Performance: Computation



Requires less training computation than comparable CNN (BiT).

**#BiT** Kolesnikov, Alexander, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. ["Big transfer \(bit\): General visual representation learning."](#) ECCV 2020.

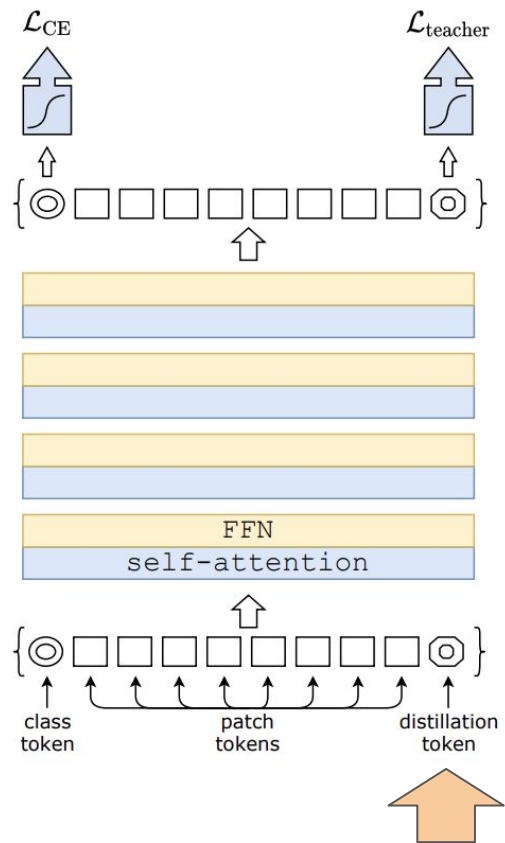
**#ViT** Dosovitskiy, Alexey, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani et al. ["An image is worth 16x16 words: Transformers for image recognition at scale."](#) ICLR 2021. [\[blog\]](#) [\[code\]](#) [\[video by Yannic Kilcher\]](#)

# Outline

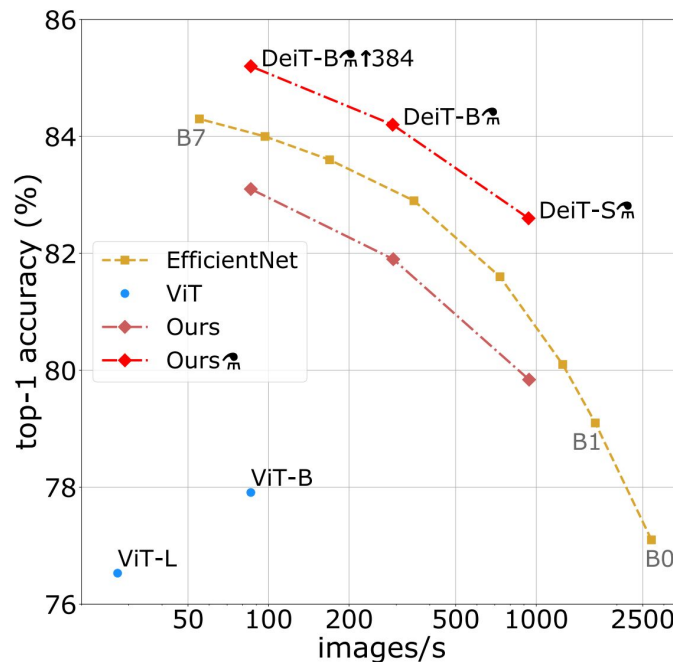
1. Vision Transformer (ViT)
  - a. Tokenization
  - b. Position embeddings
  - c. Class embedding
  - d. Receptive field
  - e. Performance
2. **Beyond ViT**



# Data-efficient Transformer (DeiT)



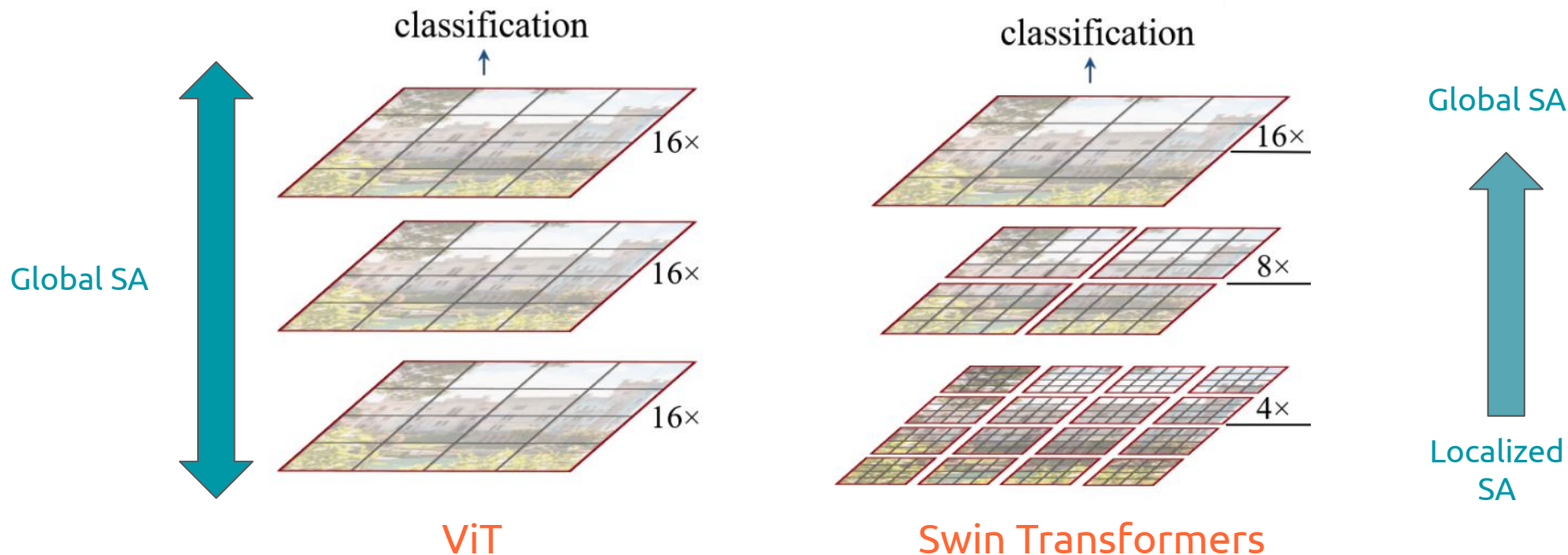
**Distillation token** that aims at predicting the label estimated by a teacher CNN. This allows introducing the convolutional bias in ViT.



**#DeiT** Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., & Jégou, H. [Training data-efficient image transformers & distillation through attention](#). ICML 2021.

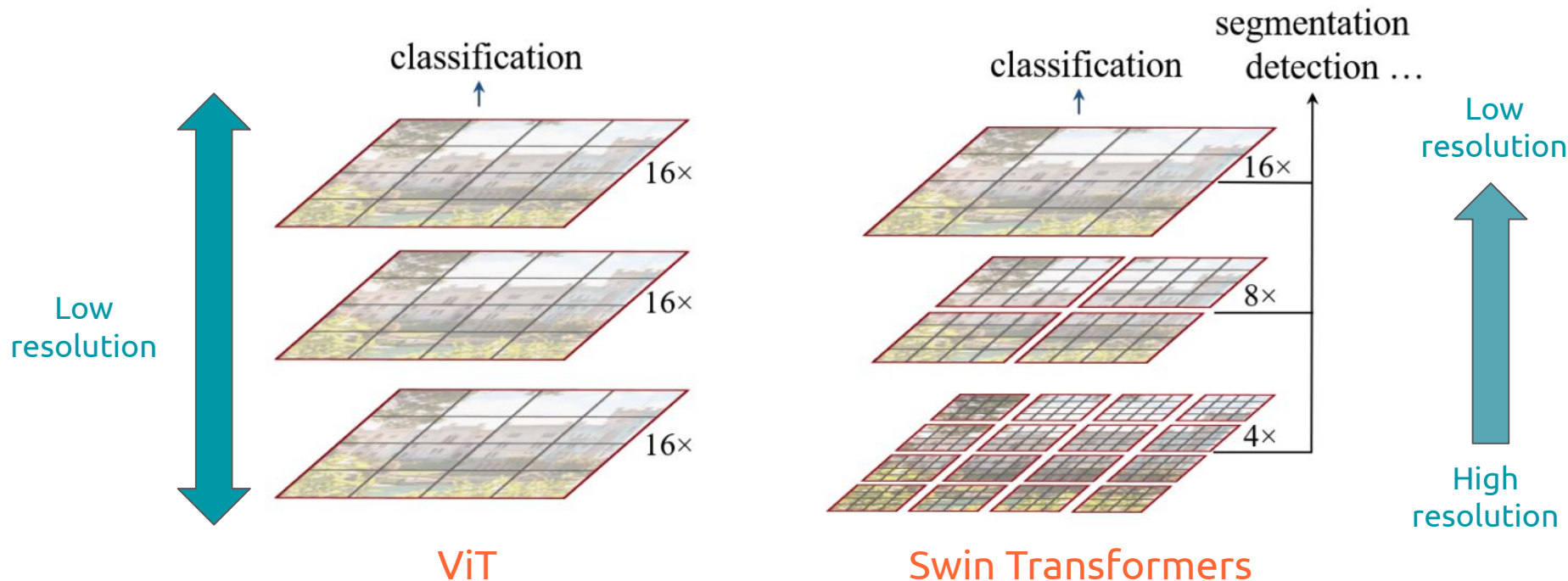
# Shifted WINDOW (SWIN) Self-Attention (SA)

Less computation by self-attending only in local windows (in grey).



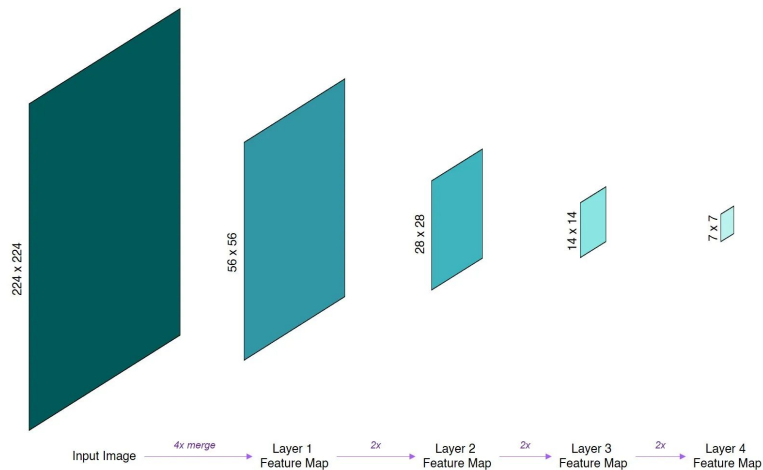
# Hierarchical ViT Backbone

Hierarchical features maps by merging **image patches (in red)** across layers.



# Hierarchical ViT Backbone

Hierarchical features maps across layers.

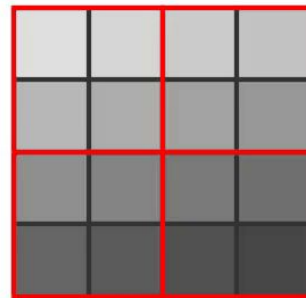


## Shifted Window MSA

Step 1: Shift window by a factor of  $M/2$ , where  $M$  = window size

Step 2: For efficient batch computation, move patches into empty slots to create a complete window.

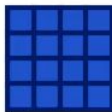
This is known as 'cyclic shift' in the paper.



# Hierarchical ViT Backbone

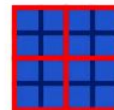
## Standard MSA

Attention for each patch is computed against all patches, resulting in quadratic complexity



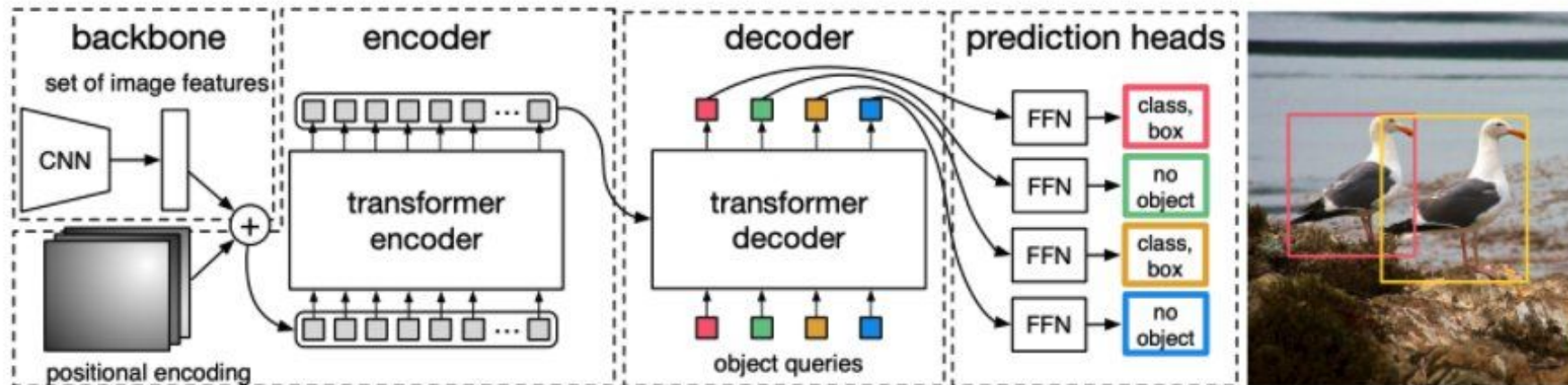
## Window-based MSA

Attention for each patch is only computed within its own window (drawn in red). Window size is 2x2 in this example.



# Object Detection

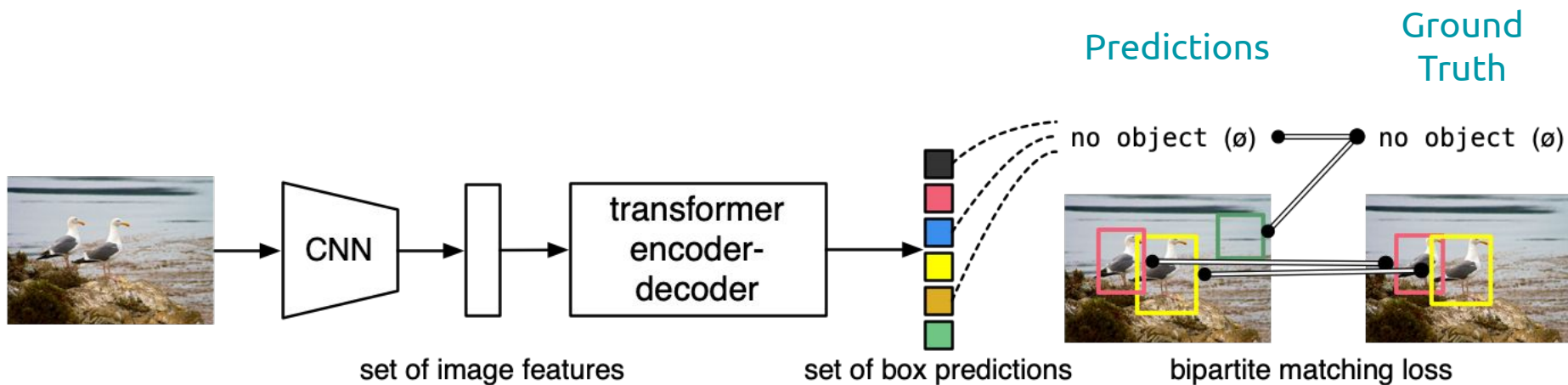
- Object detection formulated as a set prediction problem.
- DETR infers a fixed-size amount of predictions.
- Comparable performance to Faster R-CNN.



**#DETR** Carion, Nicolas, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. ["End-to-End Object Detection with Transformers."](#) ECCV 2020. [\[code\]](#) [\[colab\]](#)

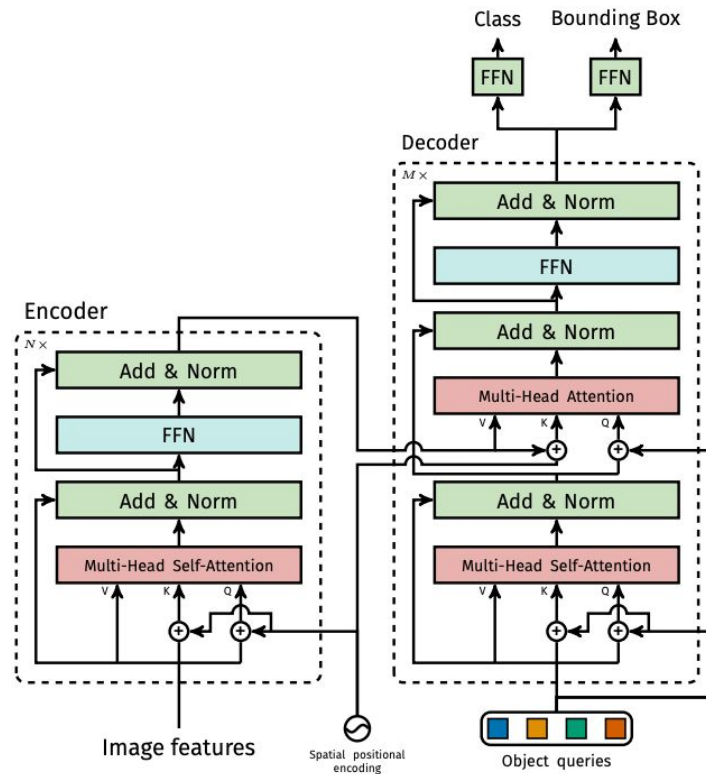
# Object Detection

- During training, bipartite matching uniquely assigns predictions with ground truth boxes.
- Prediction with no match should yield a “no object” ( $\emptyset$ ) class prediction.



# Object Detection

- Architecture of DETR.



**#DETR** Carion, Nicolas, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. ["End-to-End Object Detection with Transformers."](#) ECCV 2020. [\[code\]](#) [\[colab\]](#)



**Questions?**