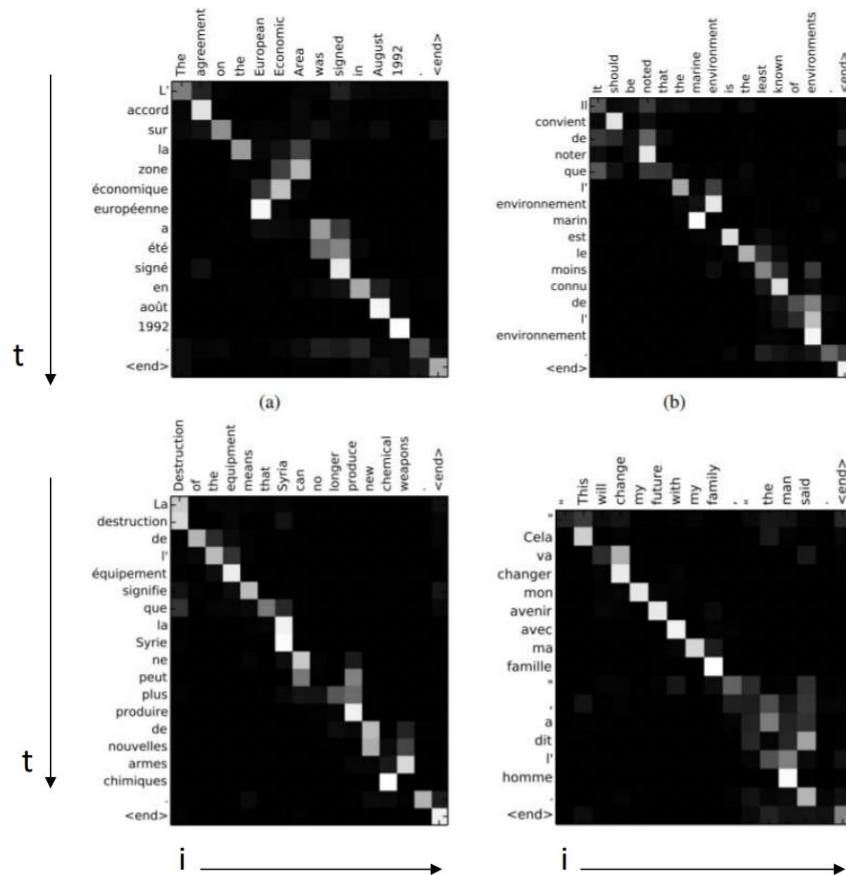


# Vision Transformers

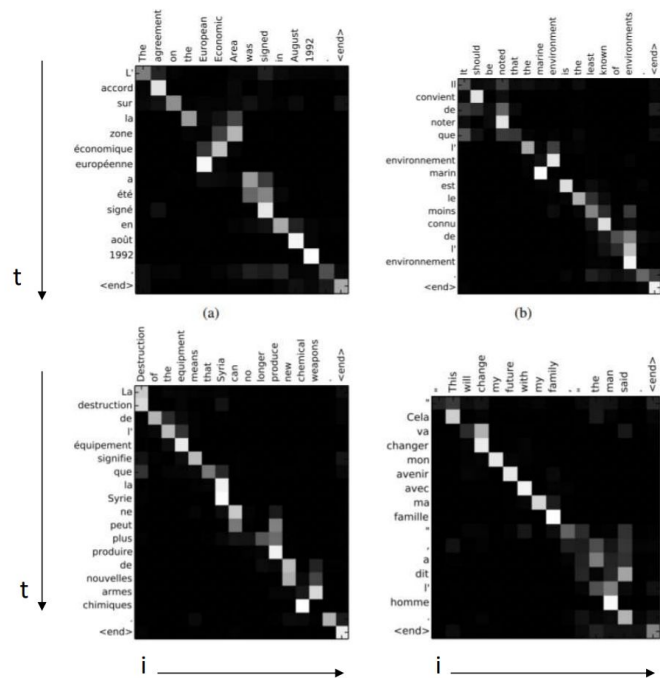
*An image is worth 16x16 words*

# Motivation

- Revolutionary paper which found a way to adapt conventional transformer architecture to computer vision tasks
- Previously transformers were predominantly used for NLP related tasks
- 'Pure transformer' → no convolutions



# Significance



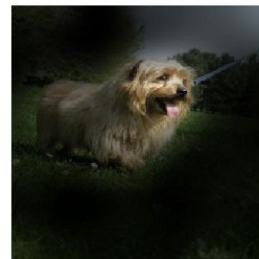
*How?*



Input



Attention



# Adapting conventional transformers

- Transformers were traditionally used where Input is a sequence of tokens.
- The attention mechanism is  $O(n^2)$  in the number of tokens
- If you were to naively apply this to images treating each pixel as a token the computational cost makes attention infeasible.
- Therefore the naive approach of each pixel attending to every other pixel wouldn't work



# Solution

1. Split an image into fixed size patches
2. Linearly embed each of them
3. Add position embeddings
4. Feed the resulting sequence into a standard transformer encoder
5. Extra learnable 'classification token'

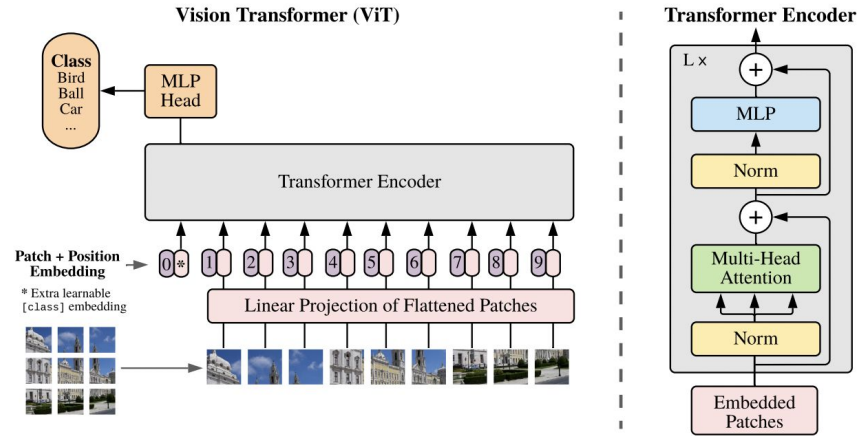
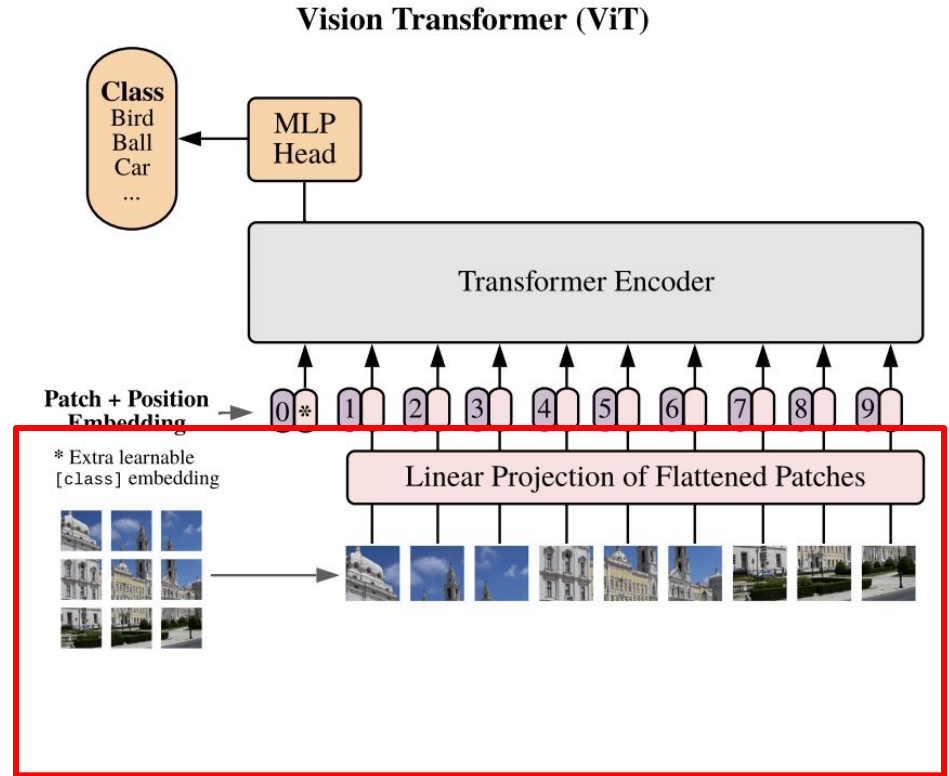


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

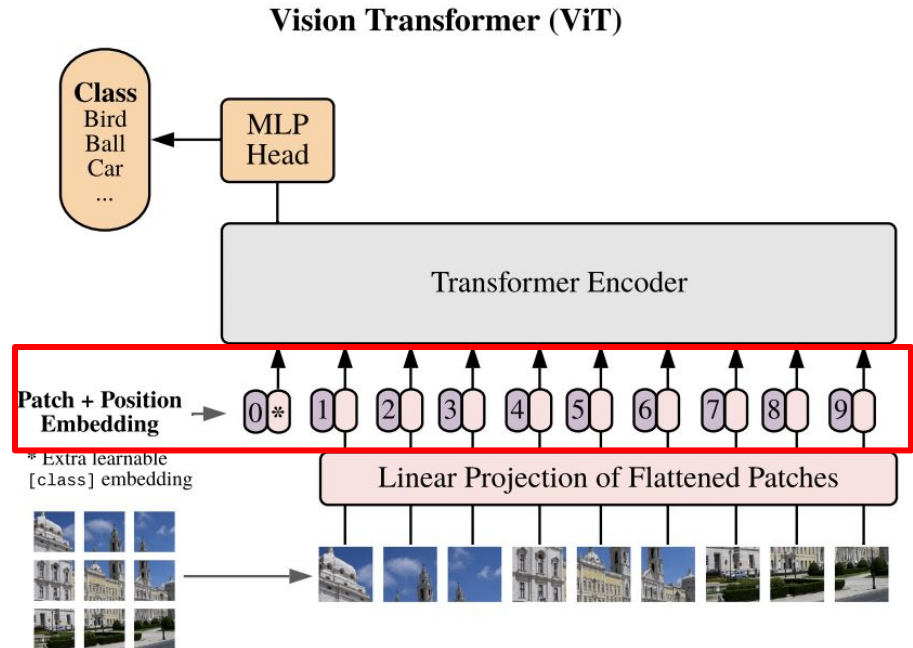
# Solution

- Split each image into 16x16 patches
- Flatten each patch into a vector, pass through a Linear layer to generate an embedding
- Steps 1 & 2 of the overall process.



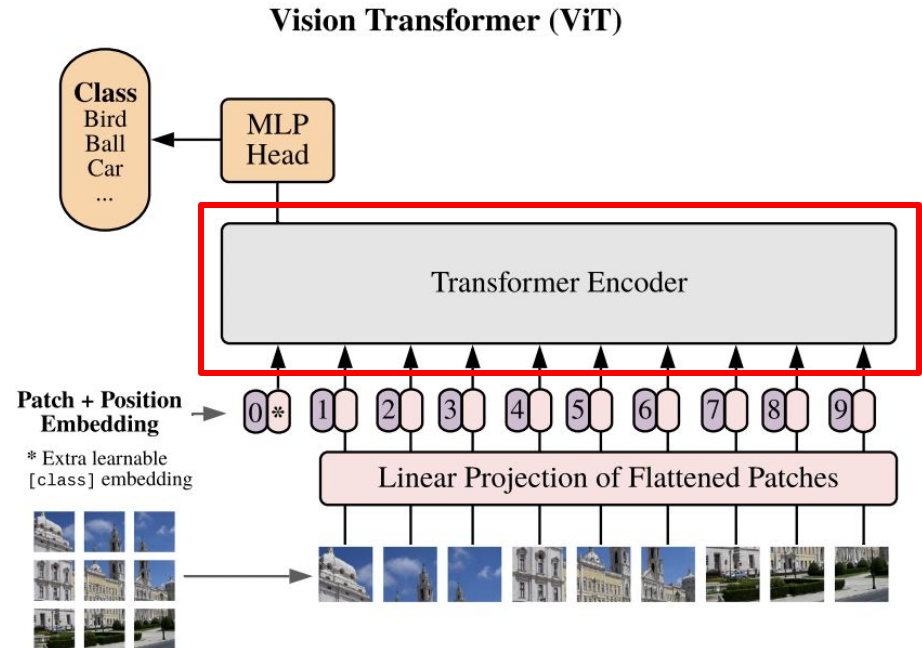
# Solution

- Combine linear embedding with position embedding



# Solution

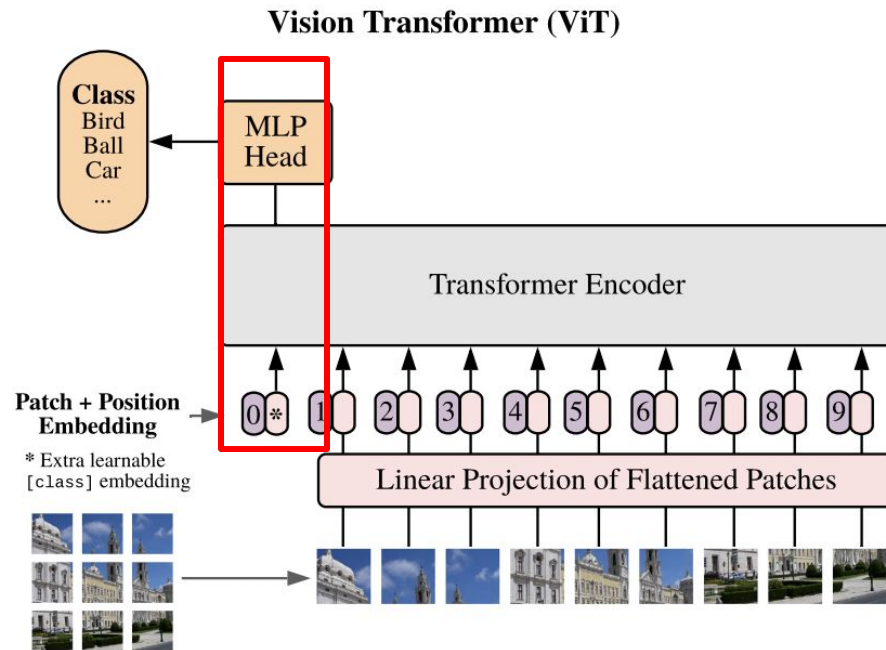
- Standard transformer encoder (BERT)





# Solution

- Use Extra learnable 'classification token' to perform classification task
- Learns a general embedding of the whole image which is then used in classification



## Did it work?

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	<b>88.55</b> $\pm 0.04$	87.76 $\pm 0.03$	85.30 $\pm 0.02$	87.54 $\pm 0.02$	88.4/88.5*
ImageNet Real	<b>90.72</b> $\pm 0.05$	90.54 $\pm 0.03$	88.62 $\pm 0.05$	90.54	90.55
CIFAR-10	<b>99.50</b> $\pm 0.06$	99.42 $\pm 0.03$	99.15 $\pm 0.03$	99.37 $\pm 0.06$	—
CIFAR-100	<b>94.55</b> $\pm 0.04$	93.90 $\pm 0.05$	93.25 $\pm 0.05$	93.51 $\pm 0.08$	—
Oxford-IIIT Pets	<b>97.56</b> $\pm 0.03$	97.32 $\pm 0.11$	94.67 $\pm 0.15$	96.62 $\pm 0.23$	—
Oxford Flowers-102	99.68 $\pm 0.02$	<b>99.74</b> $\pm 0.00$	99.61 $\pm 0.02$	99.63 $\pm 0.03$	—
VTAB (19 tasks)	<b>77.63</b> $\pm 0.23$	76.28 $\pm 0.46$	72.72 $\pm 0.21$	76.29 $\pm 1.70$	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

## Transformers or CNNs?

model	image size	FLOPs	throughput (image / s)	IN-1K / 22K trained, 1K acc.
○ Swin-T	224 <sup>2</sup>	4.5G	1325.6	81.3 / –
● ConvNeXt-T	224 <sup>2</sup>	4.5G	<b>1943.5</b> (+47%)	<b>82.1</b> / –
○ Swin-S	224 <sup>2</sup>	8.7G	857.3	83.0 / –
● ConvNeXt-S	224 <sup>2</sup>	8.7G	<b>1275.3</b> (+49%)	<b>83.1</b> / –
○ Swin-B	224 <sup>2</sup>	15.4G	662.8	83.5 / 85.2
● ConvNeXt-B	224 <sup>2</sup>	15.4G	<b>969.0</b> (+46%)	<b>83.8</b> / <b>85.8</b>
○ Swin-B	384 <sup>2</sup>	47.1G	242.5	84.5 / 86.4
● ConvNeXt-B	384 <sup>2</sup>	45.0G	<b>336.6</b> (+39%)	<b>85.1</b> / <b>86.8</b>
○ Swin-L	224 <sup>2</sup>	34.5G	435.9	– / 86.3
● ConvNeXt-L	224 <sup>2</sup>	34.4G	<b>611.5</b> (+40%)	<b>84.3</b> / <b>86.6</b>
○ Swin-L	384 <sup>2</sup>	103.9G	157.9	– / 87.3
● ConvNeXt-L	384 <sup>2</sup>	101.0G	<b>211.4</b> (+34%)	85.5 / <b>87.5</b>
● ConvNeXt-XL	224 <sup>2</sup>	60.9G	<b>424.4</b>	– / <b>87.0</b>
● ConvNeXt-XL	384 <sup>2</sup>	179.0G	<b>147.4</b>	– / <b>87.8</b>