

Vision Transformer [ViT]

Applied AI Course.com

# AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

**Alexey Dosovitskiy<sup>\*,†</sup>, Lucas Beyer<sup>\*</sup>, Alexander Kolesnikov<sup>\*</sup>, Dirk Weissenborn<sup>\*</sup>,  
Xiaohua Zhai<sup>\*</sup>, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby<sup>\*,†</sup>**

<sup>\*</sup>equal technical contribution, <sup>†</sup>equal advising

Google Research, Brain Team

{adosovitskiy, neilhoulby}@google.com

<https://arxiv.org/pdf/2010.11929.pdf>

<https://ai.googleblog.com/2020/12/transformers-for-image-recognition-at.html>

<https://openreview.net/forum?id=YicbFdNTTy>

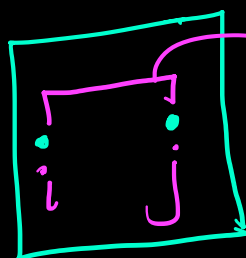
no convolution

ViT → Transformer based [NLP]  
comparable to SOTA CNNs  
lower compute-power

Major challenge for non-conv Transformers.

→ Attention

640x640 pixels



one attention weight  
for every pair of pixels

↓  
640x640  $c_2$  attention

~83.88B

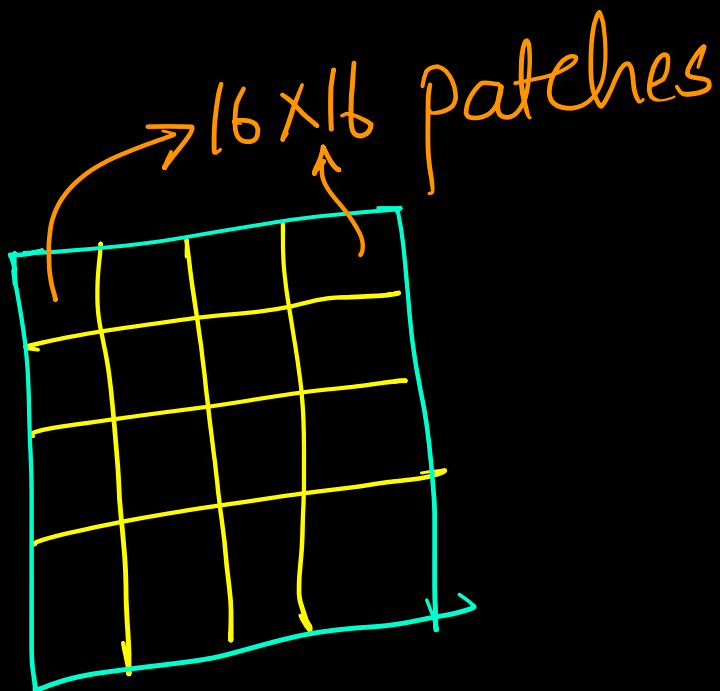


Attention  
(NLP)

← 512 →  
tokens

512  
 $C_2$

Key-idea



$W_1$   $W_2$   $W_3$  ...  $W_n$   $\rightarrow$  NLP

A sequence of weight matrices  $W_1, W_2, W_3, \dots, W_n$  is shown, each represented by a small square. An arrow points from the sequence to the text 'NLP'.

16 16 16 ... 16  $\rightarrow$  CV

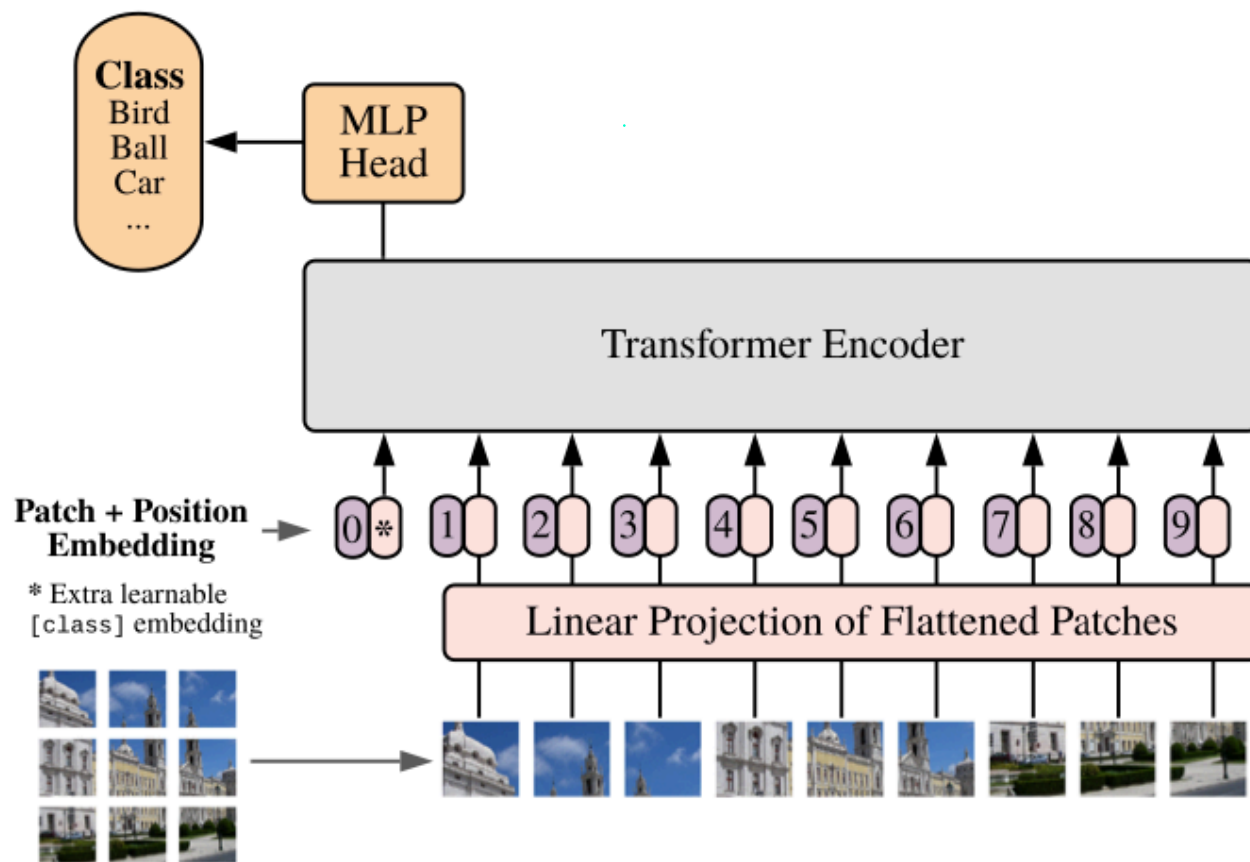
A sequence of patches is shown, each represented by a small square with '16' written above and to the left of it. An arrow points from the sequence to the text 'CV'.

patch

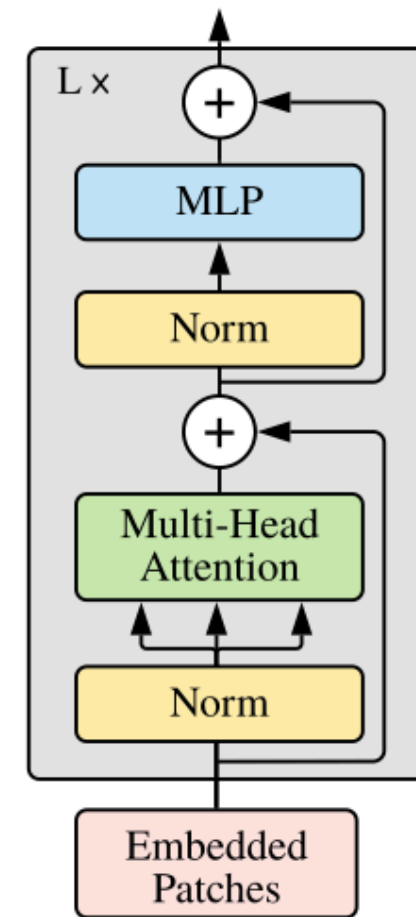
$640 \times 640 \rightarrow 40 \times 40$  patches (16x16 size)

Attention  $\rightarrow 40 \times 40 C_2 = 1.279$  Million  
 $\downarrow$   
more manageable

## Vision Transformer (ViT)

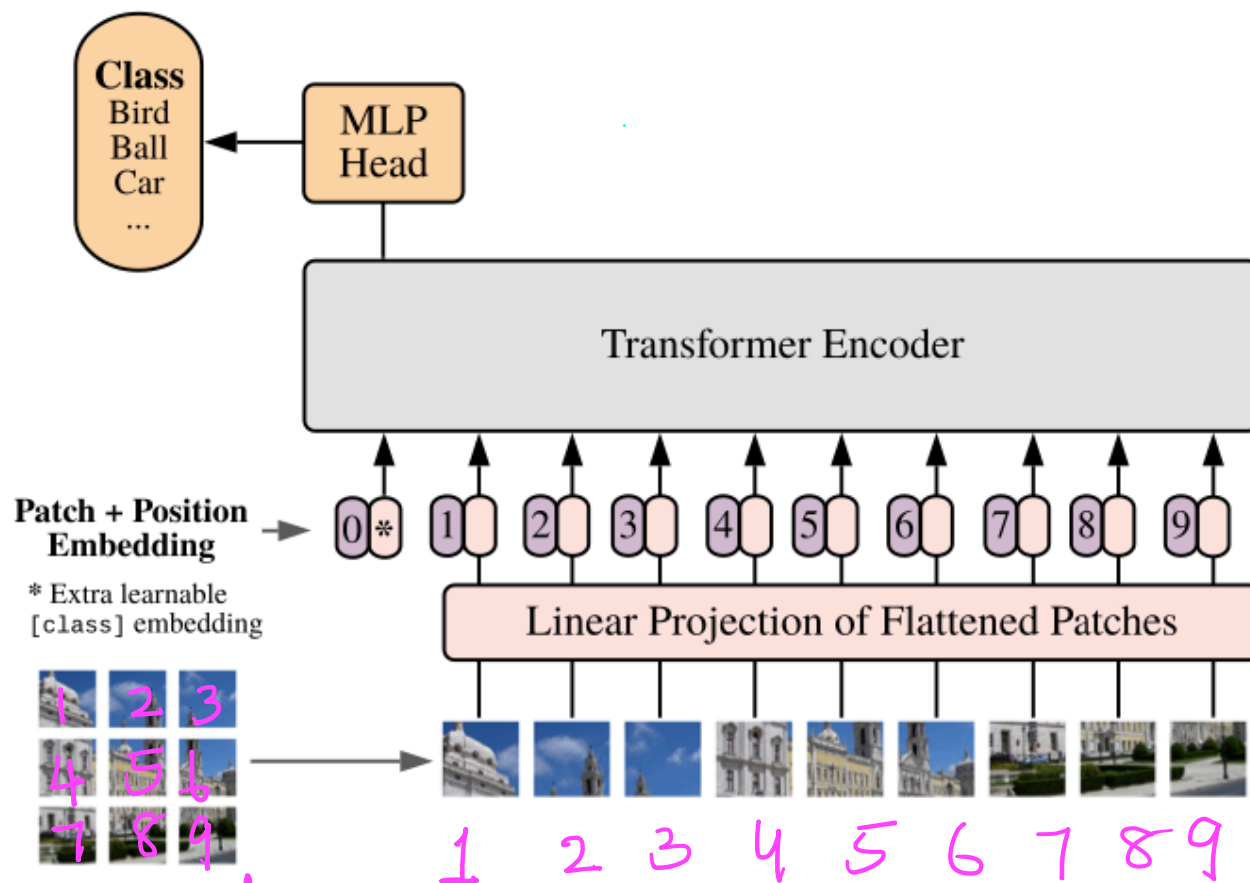


## Transformer Encoder



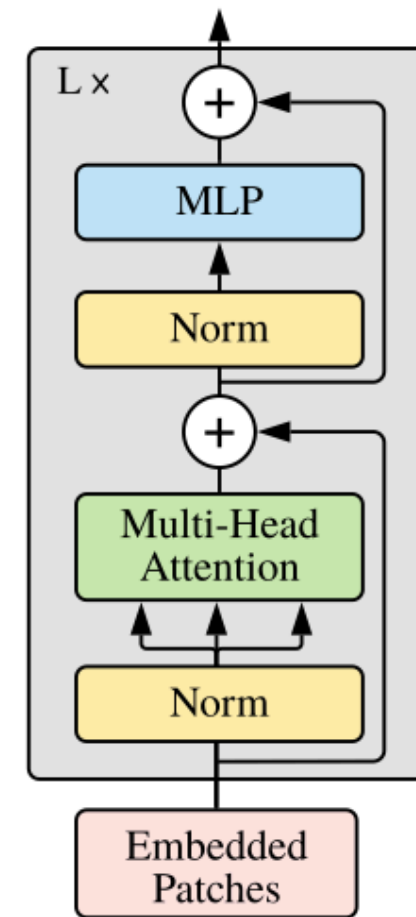


## Vision Transformer (ViT)

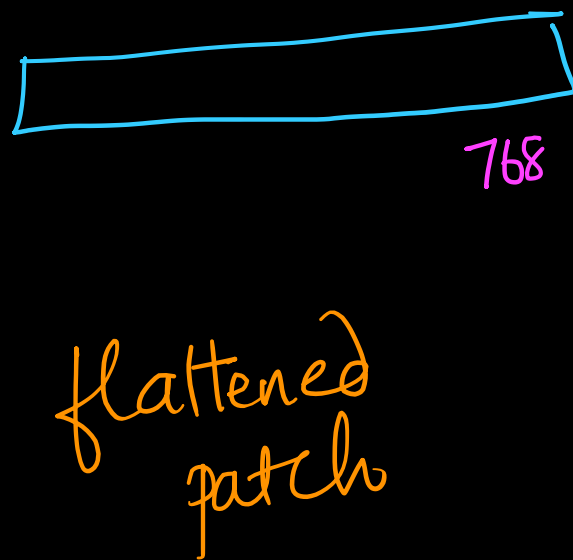
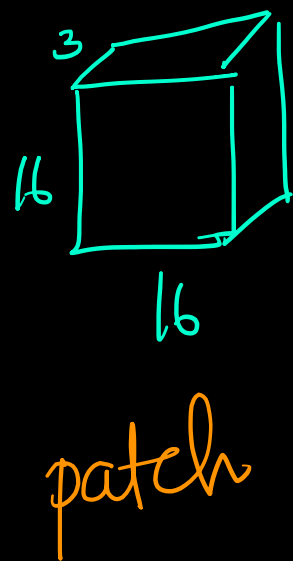


16x16 patches

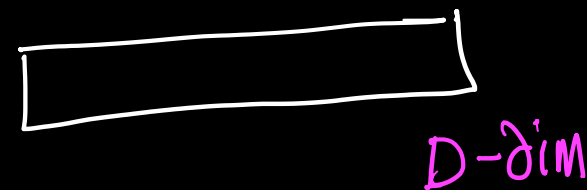
## Transformer Encoder



# Linear Projection of flattened patches

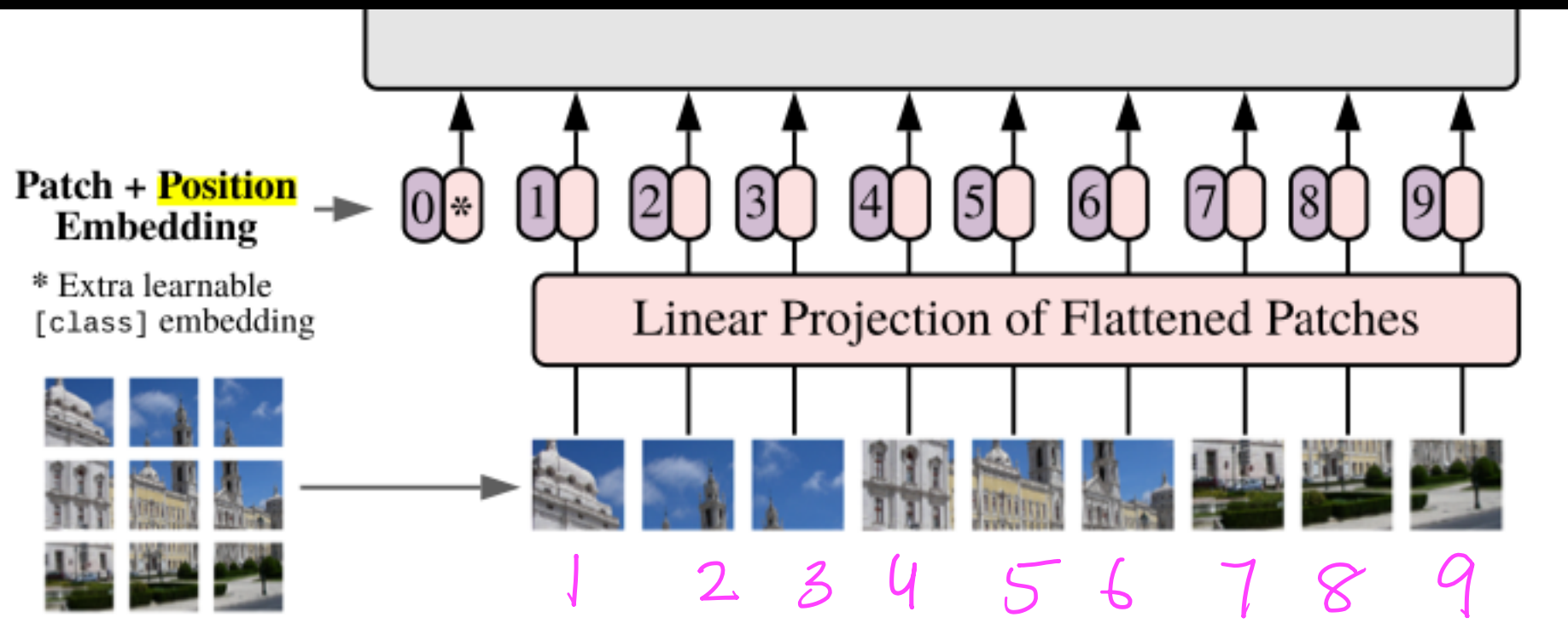


Linear  
Projection



patch-embedding  
(similar to word  
embedding)

# Positional Embedding (PE)



1-Dim PE

sin ↙      ↘ learnable

Pos. Emb.	Default/Stem	Every Layer	Every Layer-Shared
No Pos. Emb.	0.61382	N/A	N/A
1-D Pos. Emb.	0.64206	0.63964	0.64292
2-D Pos. Emb.	0.64001	0.64046	0.64022
Rel. Pos. Emb.	0.64032	N/A	N/A

Table 8: Results of the ablation study on positional embeddings with ViT-B/16 model evaluated on ImageNet 5-shot linear.

→ 1-D pos-Embed @ input is good enough

Learned PE captures  
patch positions

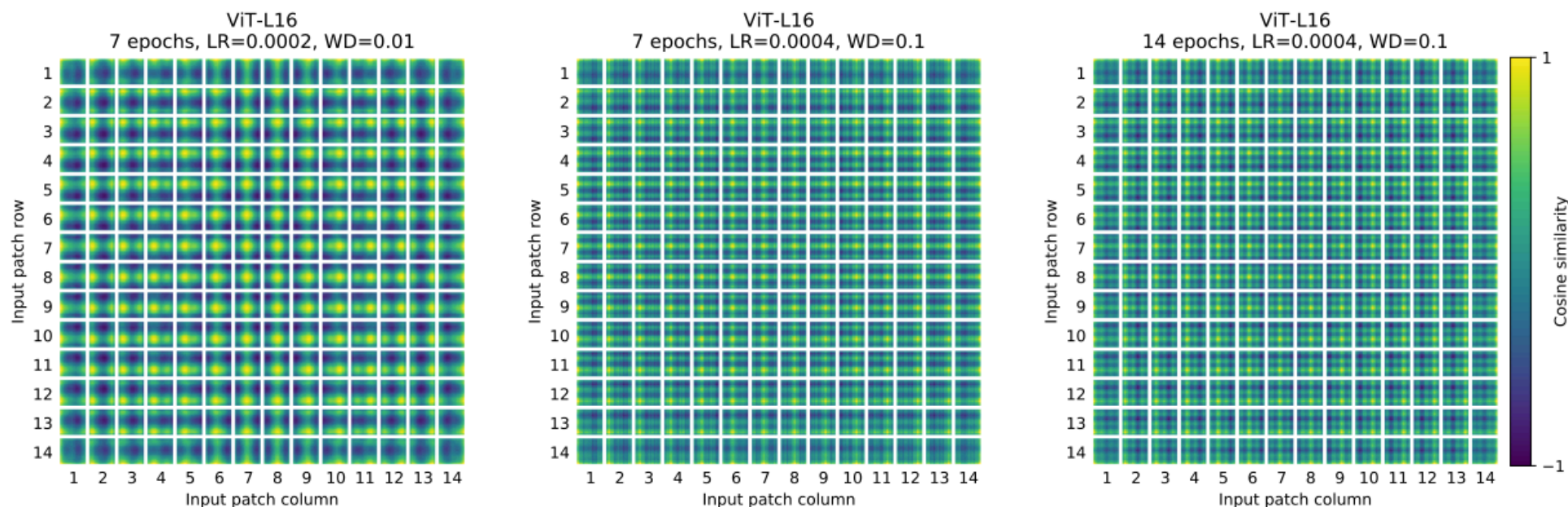
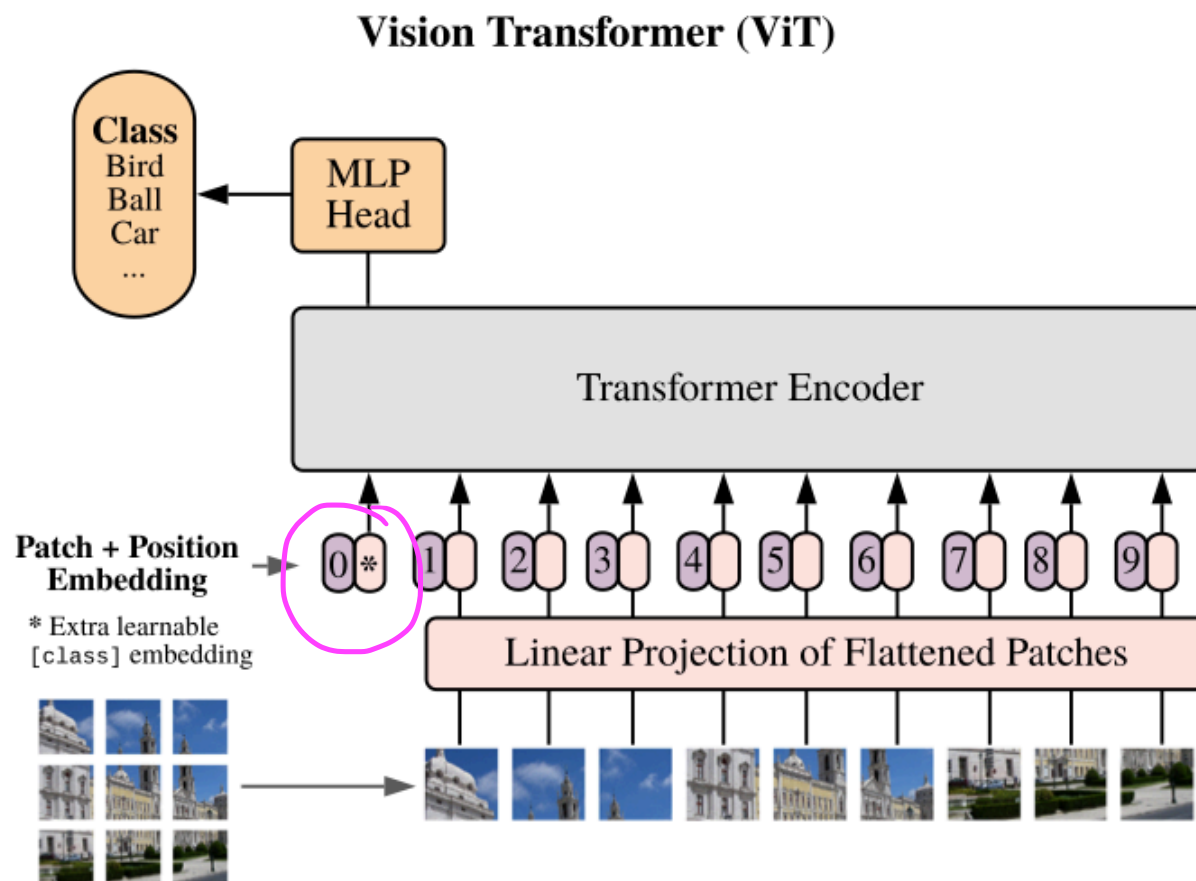


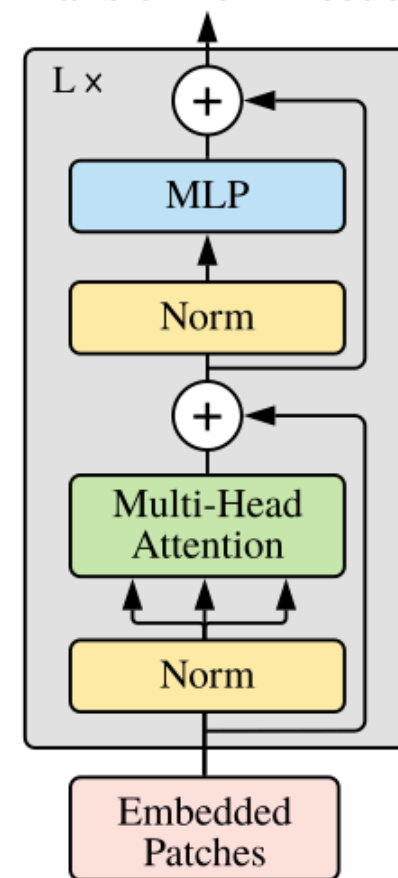
Figure 9: Position embeddings of models trained with different hyperparameters.

224x224 sized images as input

BERT [CLS] like  
O[\*]

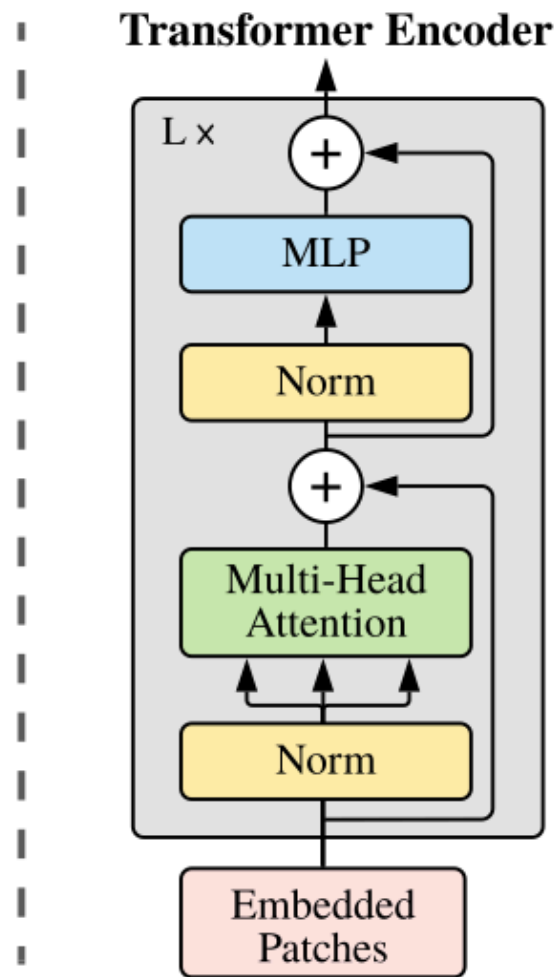
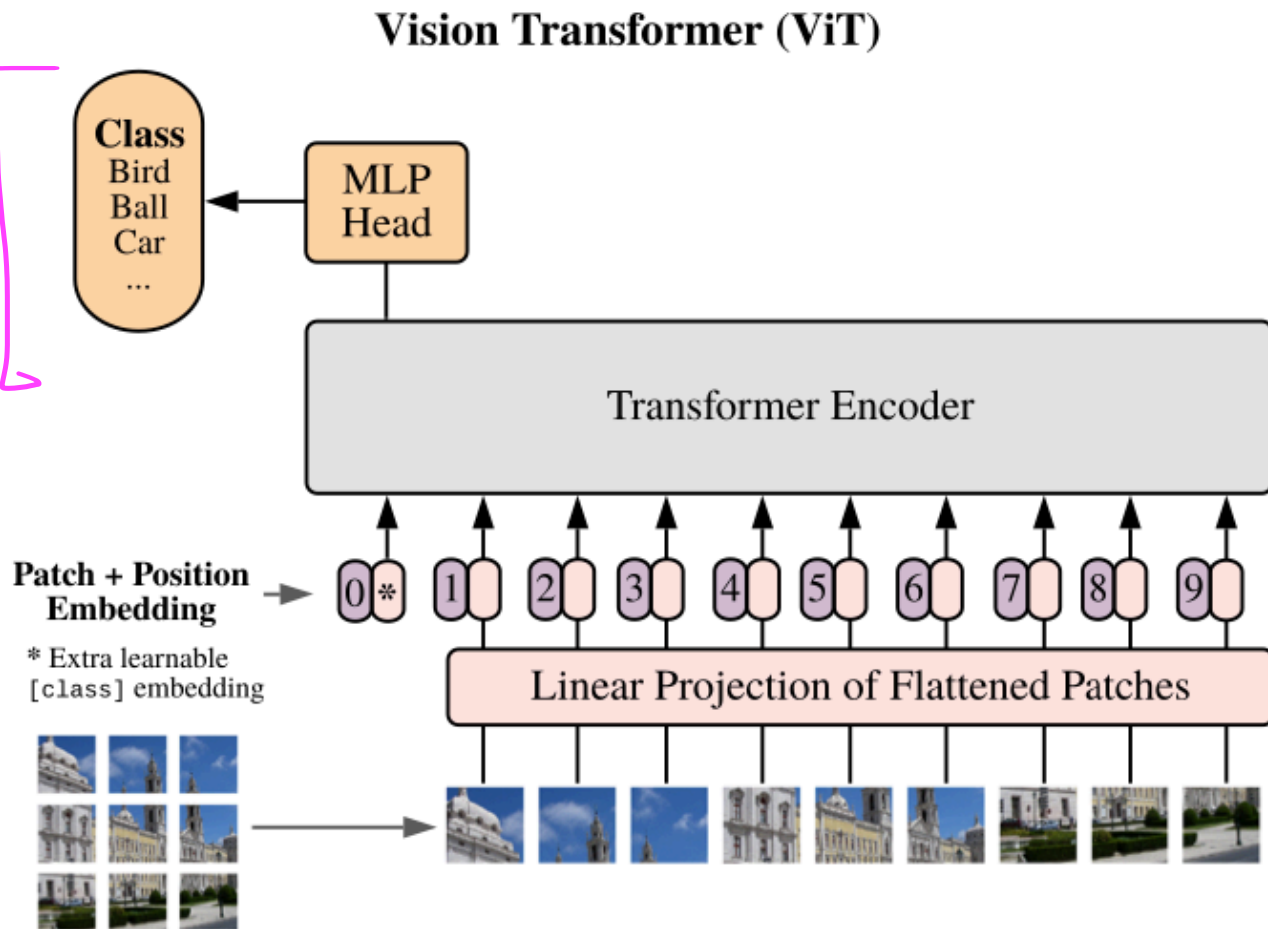


Transformer Encoder



Encoder-only Transformer  
[BERT]

classifi-  
cation





Model	Layers	Hidden size $D$	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

Table 1: Details of Vision Transformer model variants.



MLP's in Transformer  
Encoder

Smotherer ReLU

$$\sim 0.5x(1 + \tanh[\text{SQRT}(2/\pi)(x + 0.044715x^3)])$$

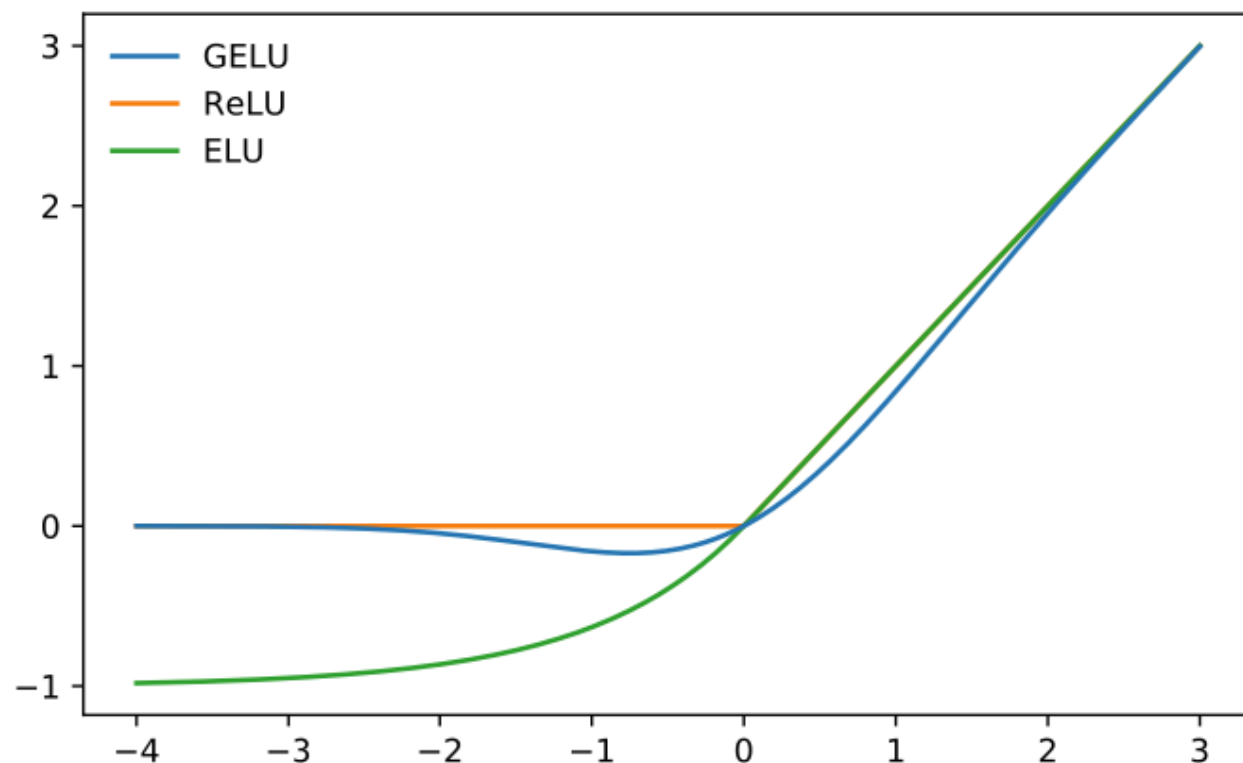
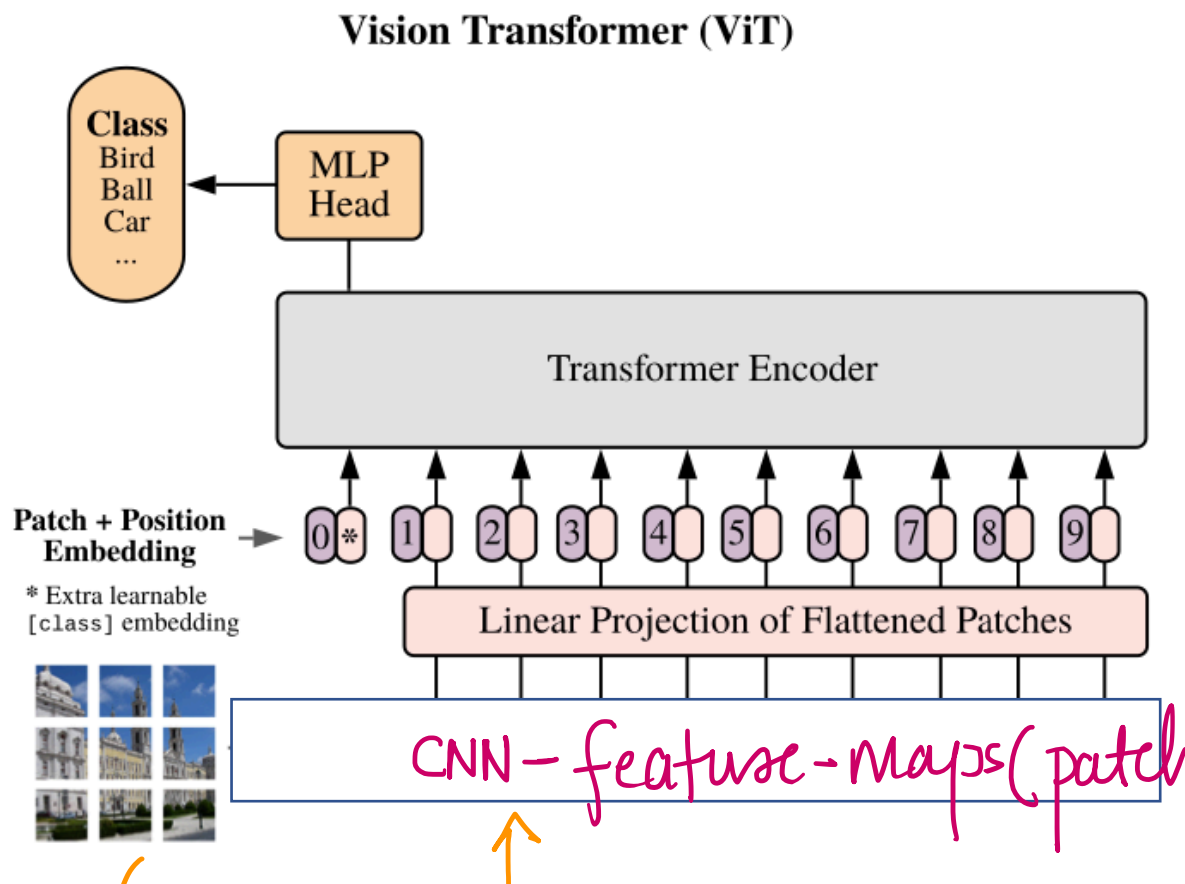


Figure 1: The GELU ( $\mu = 0, \sigma = 1$ ), ReLU, and ELU ( $\alpha = 1$ ).

<https://arxiv.org/pdf/1606.08415.pdf>

# Hybrid Architecture



CNN - feature - maps (patch-wise)

CNN

patches

# Experimental Results

ViT on ImageNet  $\rightarrow$  77.9% top-1 accuracy

SOTA  
[ResNet-based]  
CNN  $\rightarrow$  85.8%

$\rightarrow$  overfitting

ImageNet-

ImageNet-21K  
[1M, 21K classes]

JFT  
[300MM, 18K classes]

BiT  
[CNN-based]

✓

~

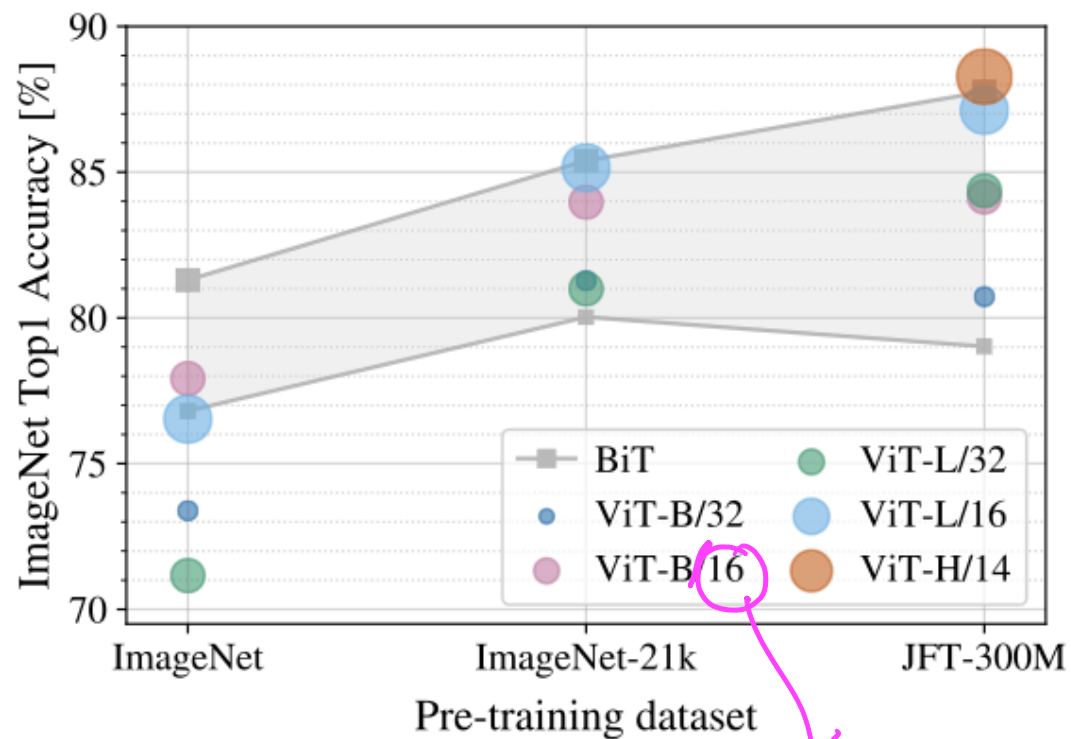
-

ViT

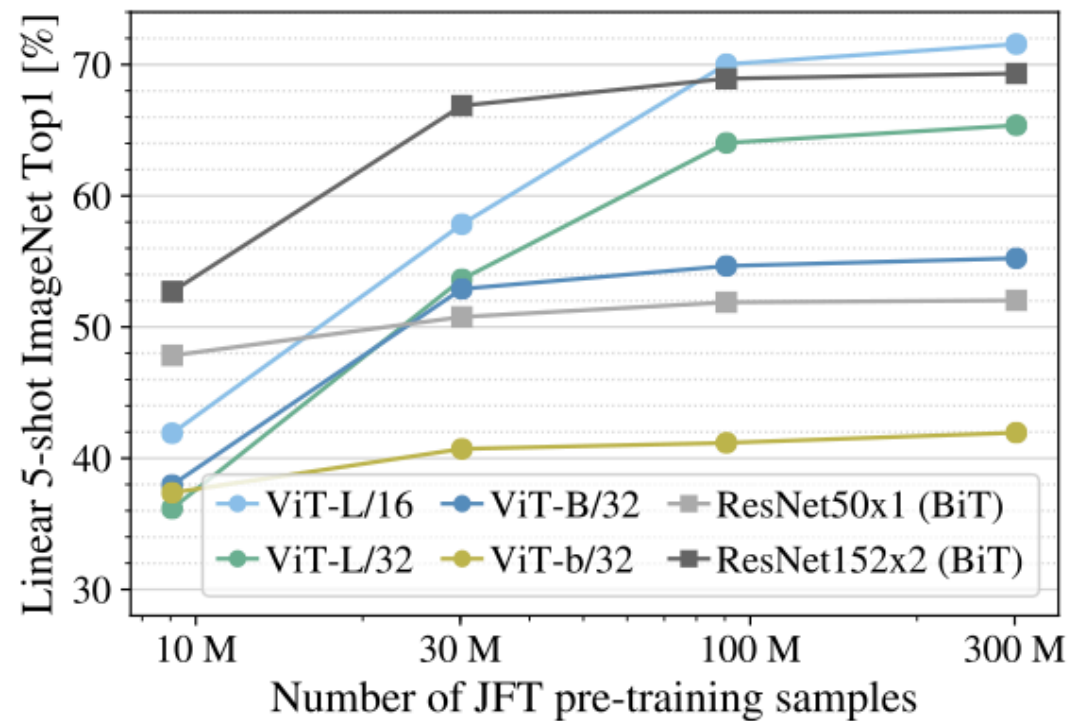
-

~

✓



16x16 patches



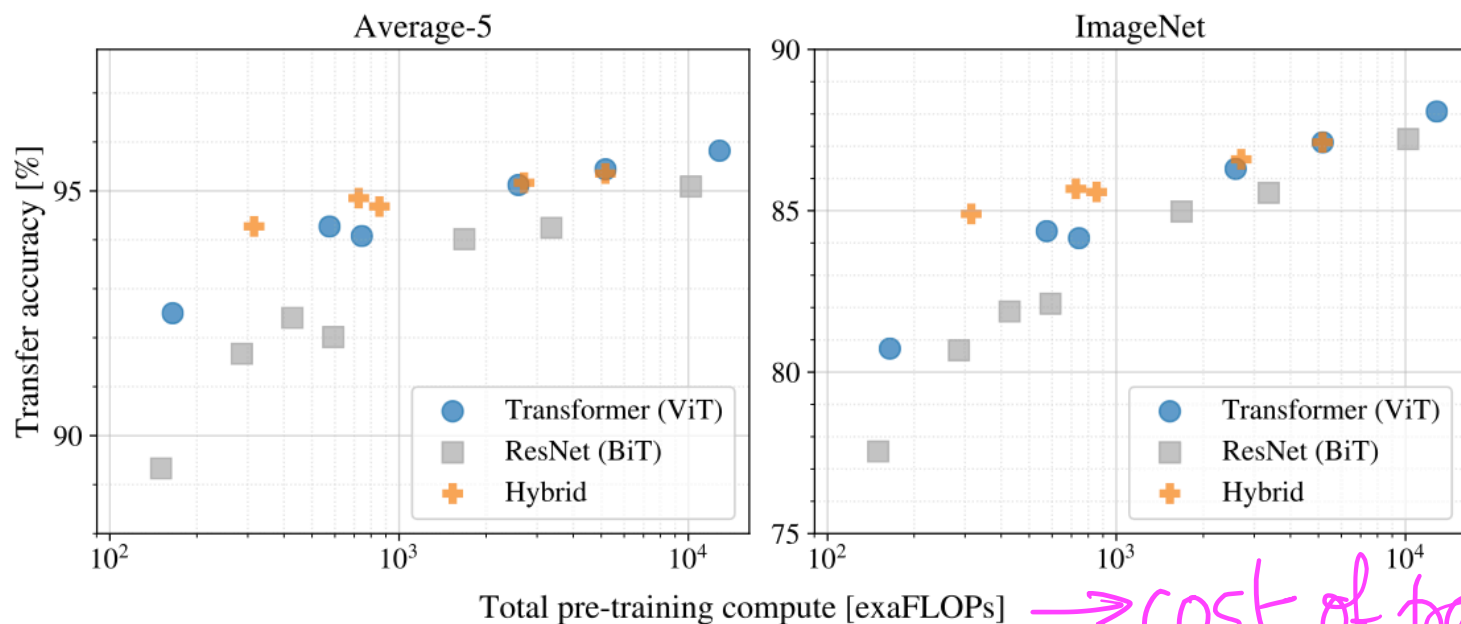


Figure 5: Performance versus cost for different architectures: Vision Transformers, ResNets, and hybrids. Vision Transformers generally outperform ResNets with the same computational budget. Hybrids improve upon pure Transformers for smaller model sizes, but the gap vanishes for larger models.

code

[https://github.com/google-research/vision\\_transformer](https://github.com/google-research/vision_transformer)

Model :

[https://github.com/google-research/vision\\_transformer/blob/master/vit\\_jax/models.py](https://github.com/google-research/vision_transformer/blob/master/vit_jax/models.py)

jax = autograd  
+ XLA (speed up ML)



using a pretrained Model:

[https://colab.research.google.com/github/google-research/vision\\_transformer/blob/master/vit\\_jax.ipynb](https://colab.research.google.com/github/google-research/vision_transformer/blob/master/vit_jax.ipynb)



## Conclusion

→ first steps towards non-conv models in CV

→ ~18-24 months for more models/variations

Q & A