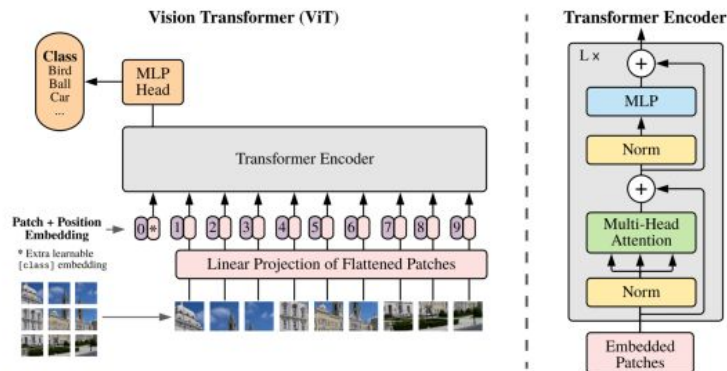
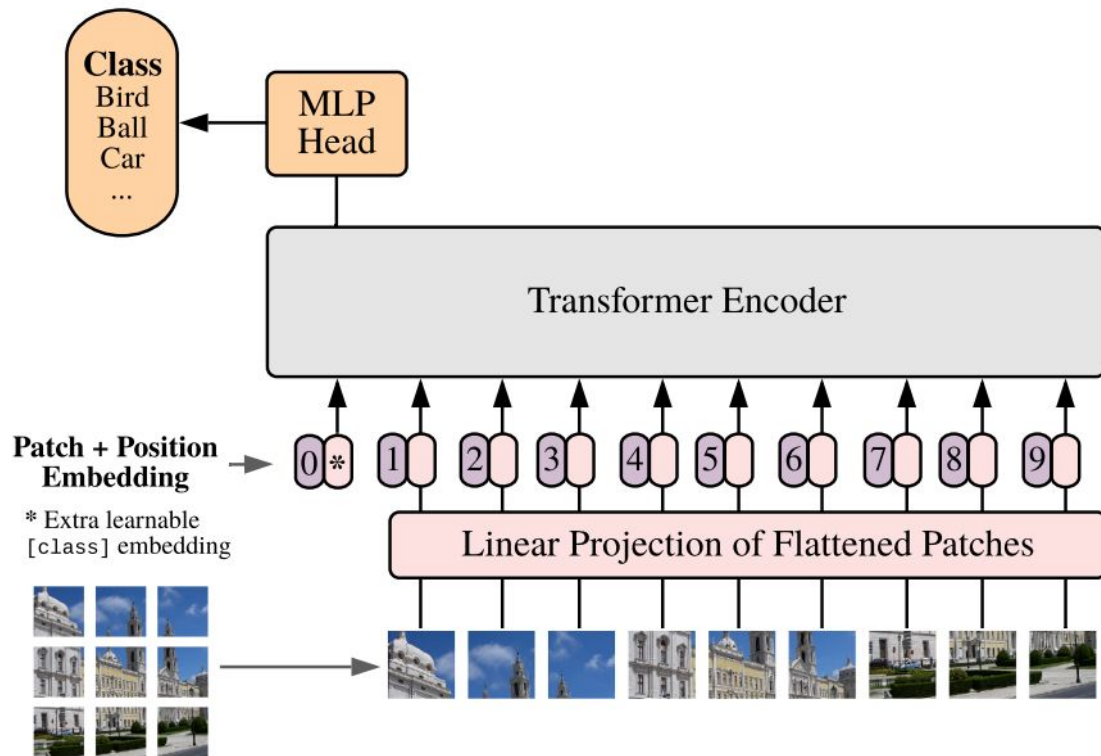


VISION TRANSFORMERS



What is Vision Transformer?



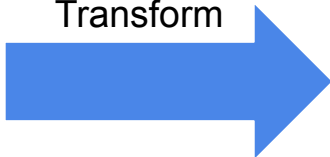
Input Image Processing

Input Image



Size: 512 x 512 x 3

Transform



Patches



Size: 32 x 32 x 3
Num Patches: 256

Input Image Processing

THEORY

Input Image: $H \times W \times C$.

Patch Size: $P_h \times P_w$

Number of patches (N): $(H \times W) / (P_h \times P_w)$

Transformed Input: $(N, P_h \times P_w \times C)$

$H \times W$ = Image height x width

C = Image channels

$P_h \times P_w$ = Patch height x width

N = Number of patches

EXAMPLE

Input Image: $512 \times 512 \times 3$

Patch Size: 32×32

Number of patches (N):

$$= (512 \times 512) / (32 \times 32)$$

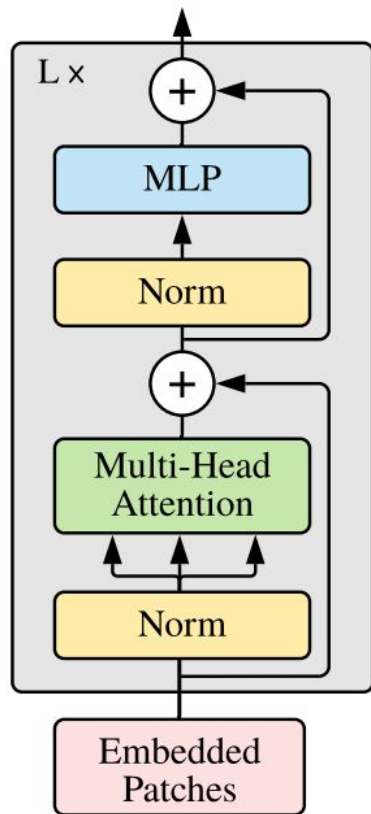
$$= 256$$

Transformed Input:

$$= (256, 32 \times 32 \times 3)$$

$$= (256, 3072)$$

Transformer Encoder



Norm: Layer Normalization.

MLP: Uses GELU activation function.

ViT Variants

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