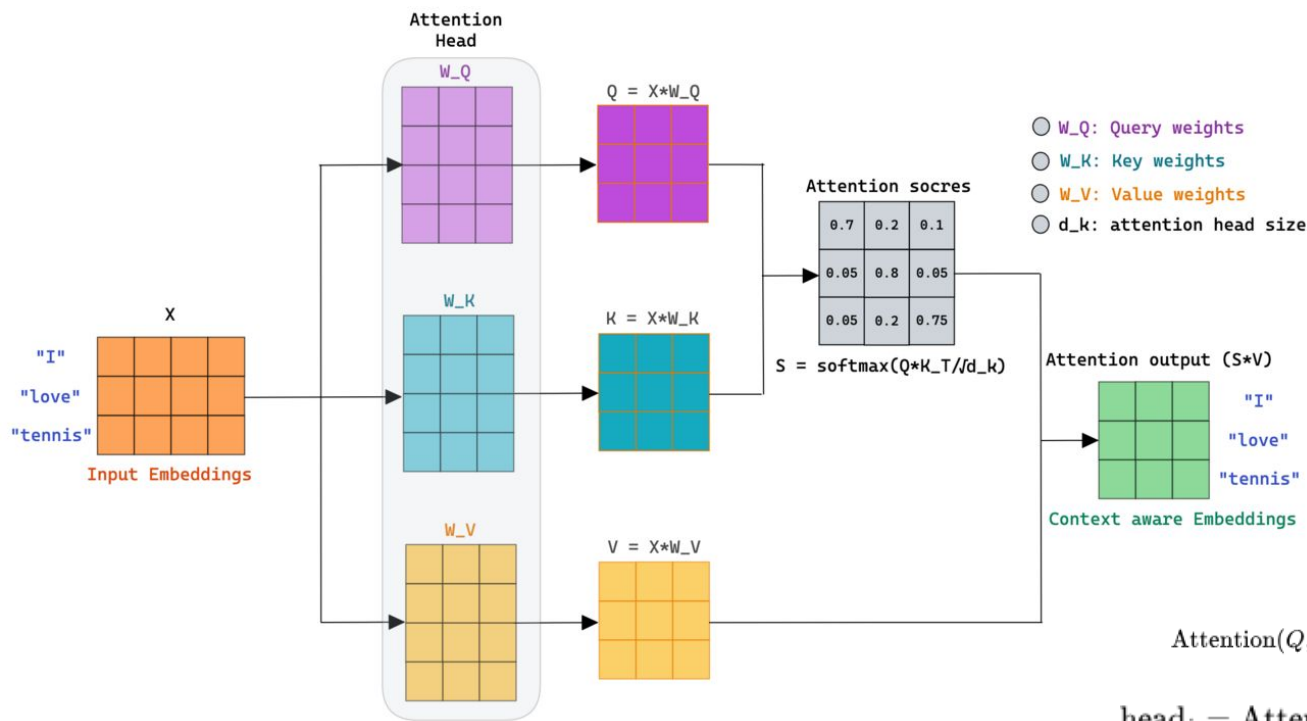


CV Transformers

Attention is All you Need

Attention is All you Need

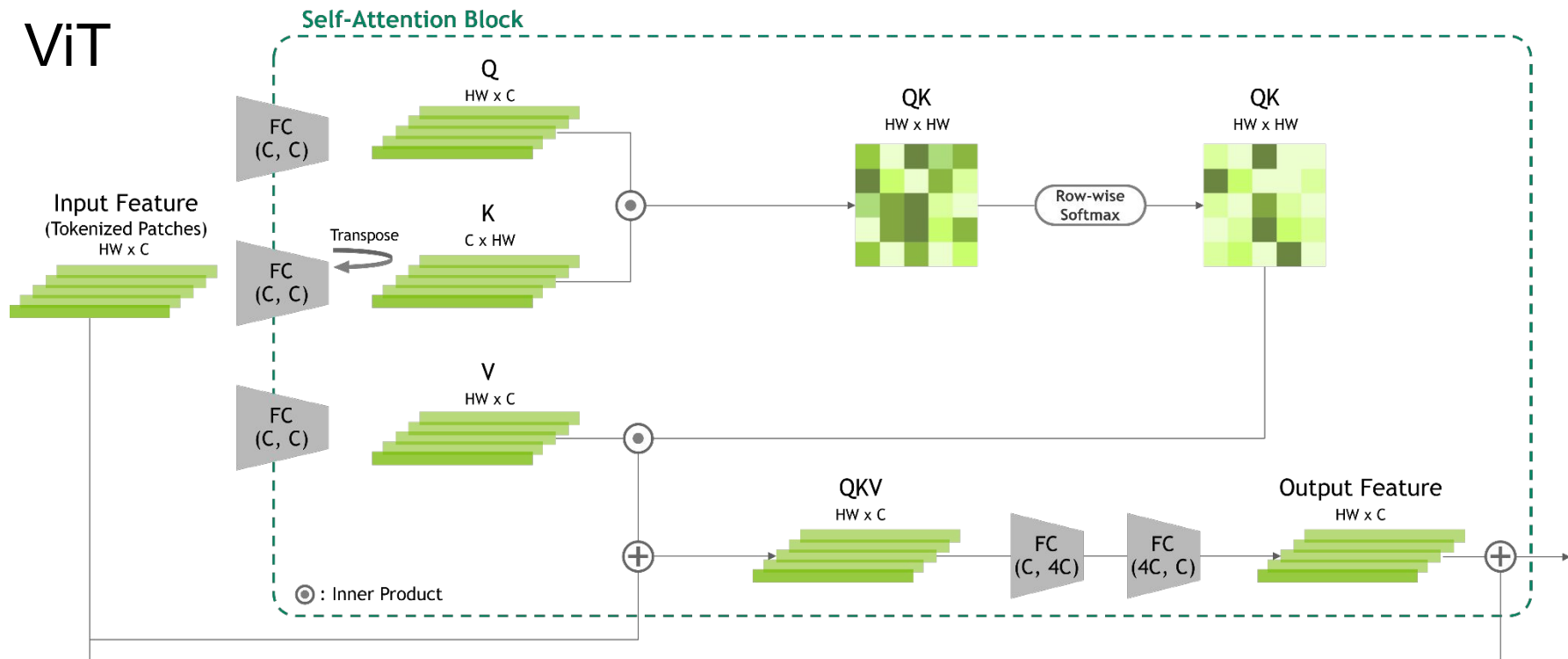


$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V,$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V).$$

$$\text{MultiHeadAttention}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O,$$

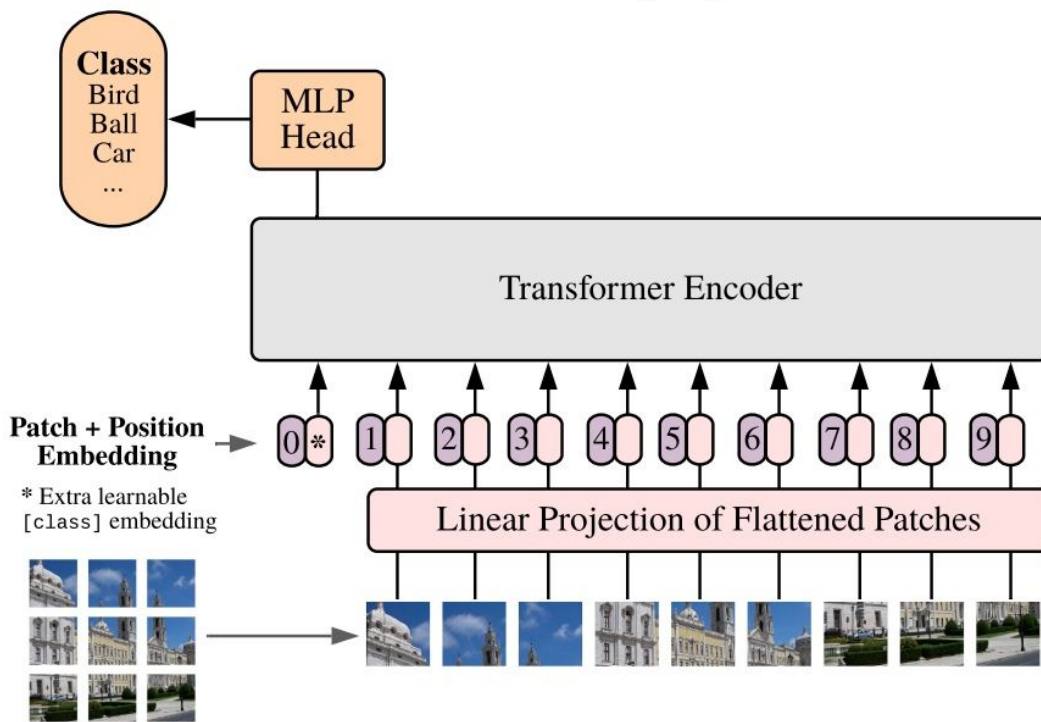
ViT



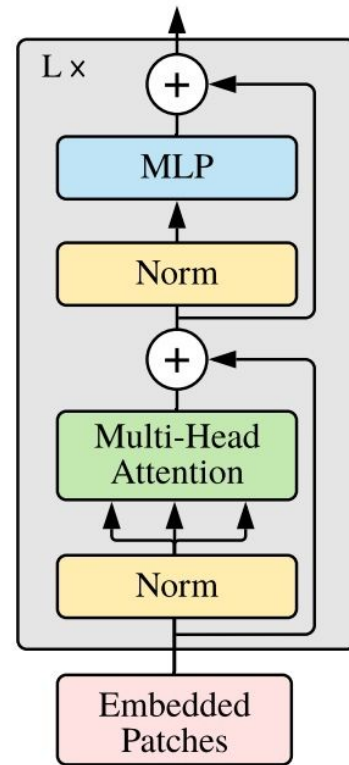
[link](#)

ViT

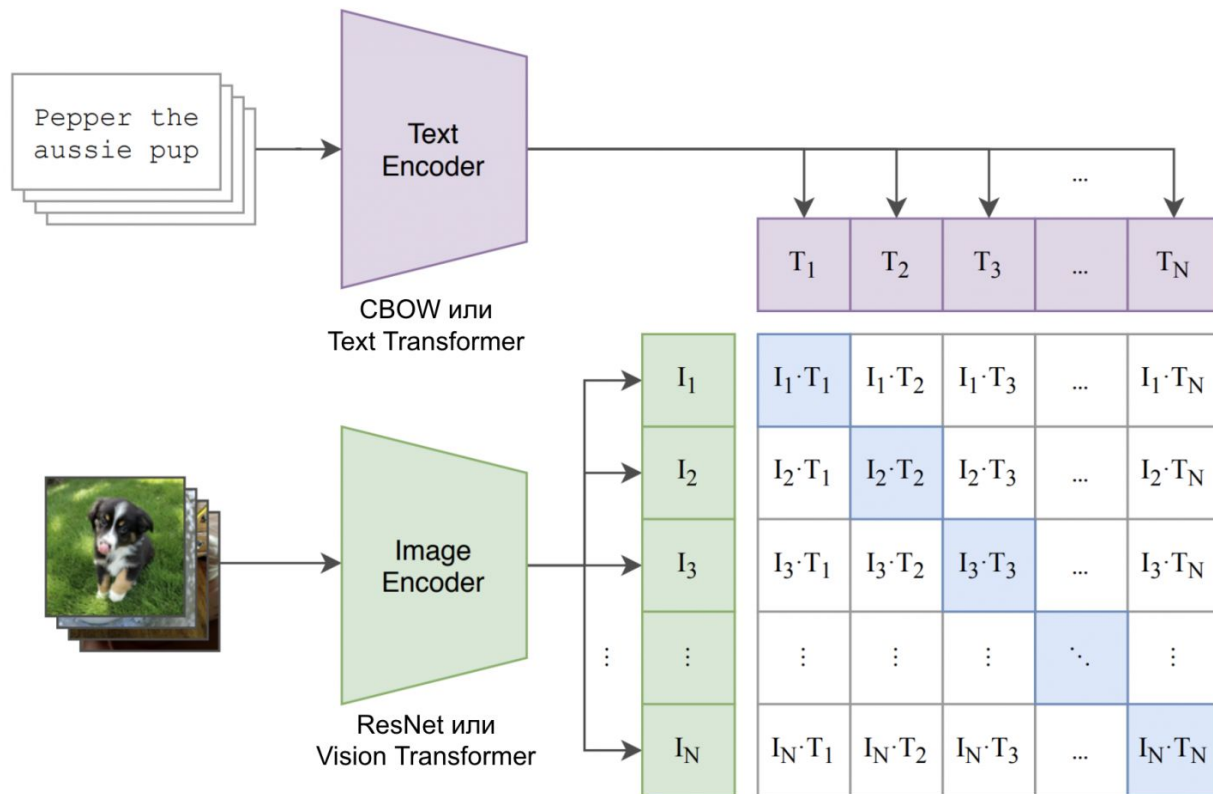
Vision Transformer (ViT)



Transformer Encoder

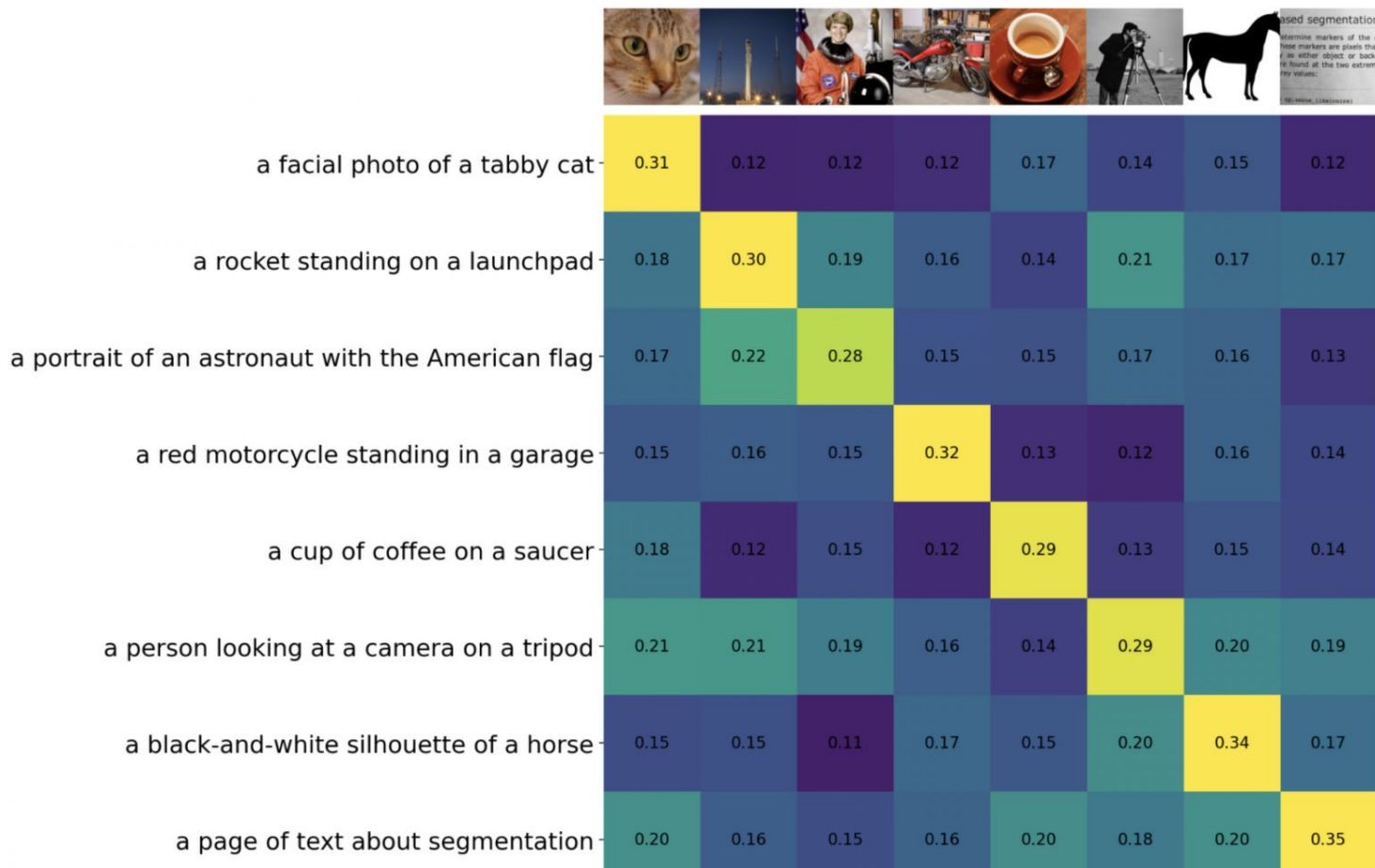


CLIP

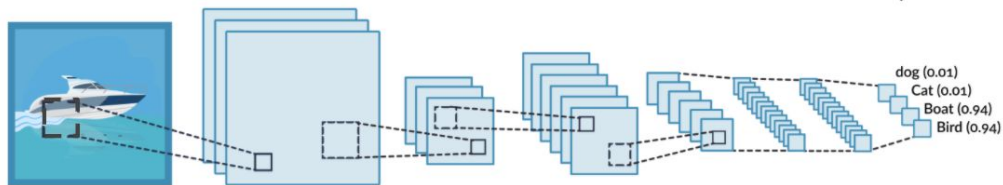


CLIP

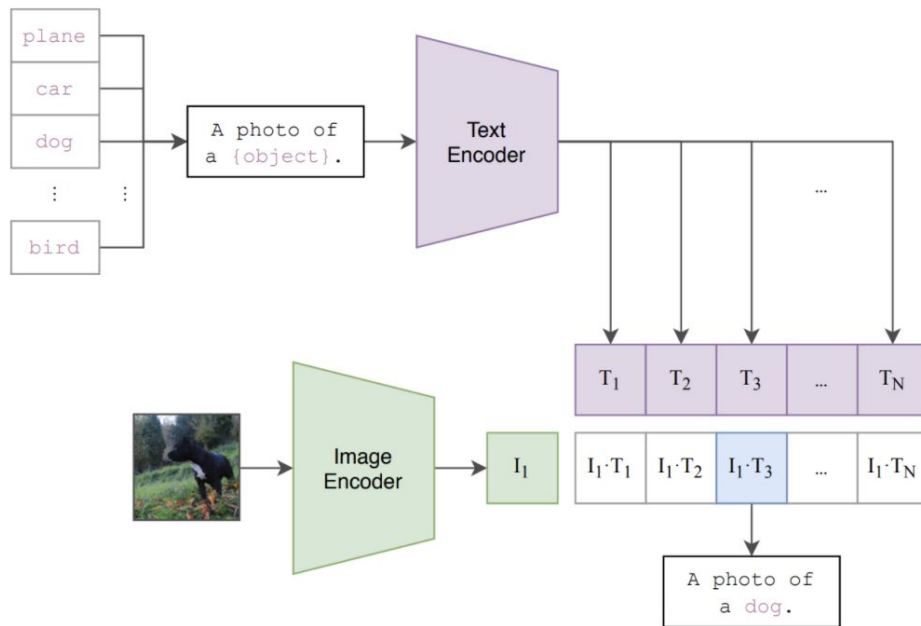
Cosine similarity between text and image features



CLIP

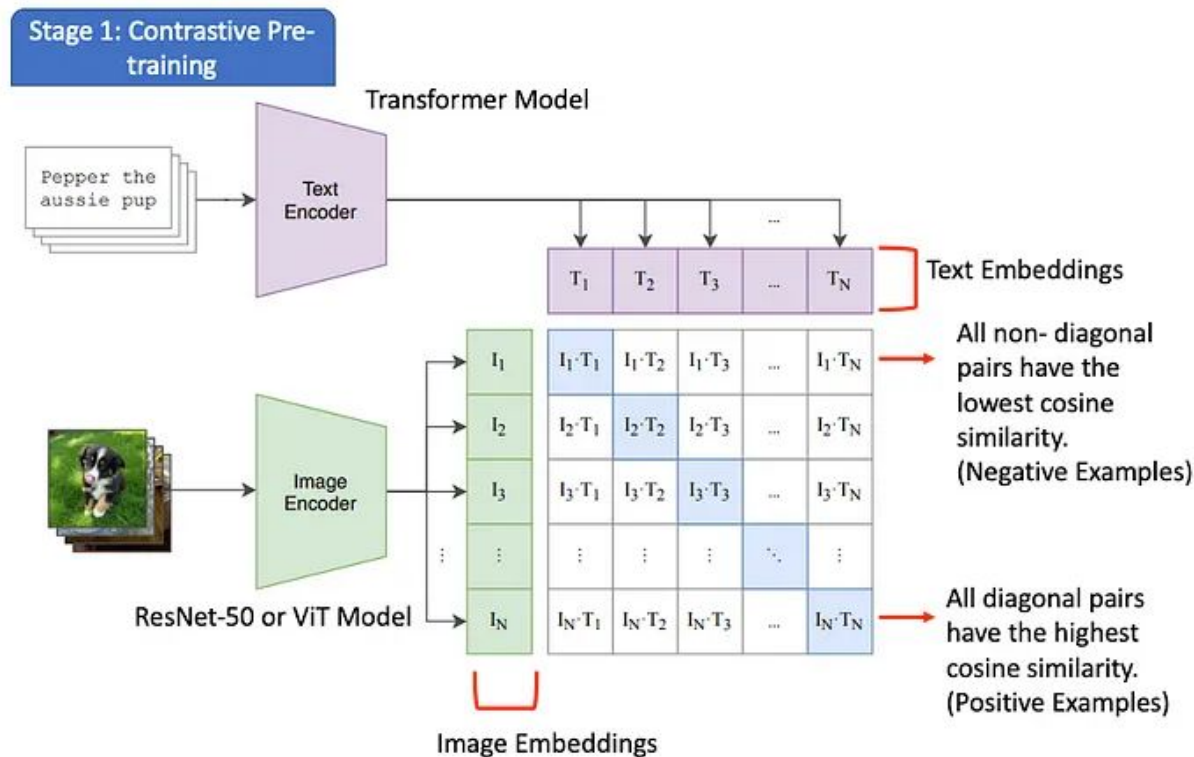


Классификация изображений при помощи сверточной нейросети в рамках “классической” парадигмы



Классификация изображений через схожесть репрезентаций сверточной сети и трансформера в рамках “гибридной” парадигмы

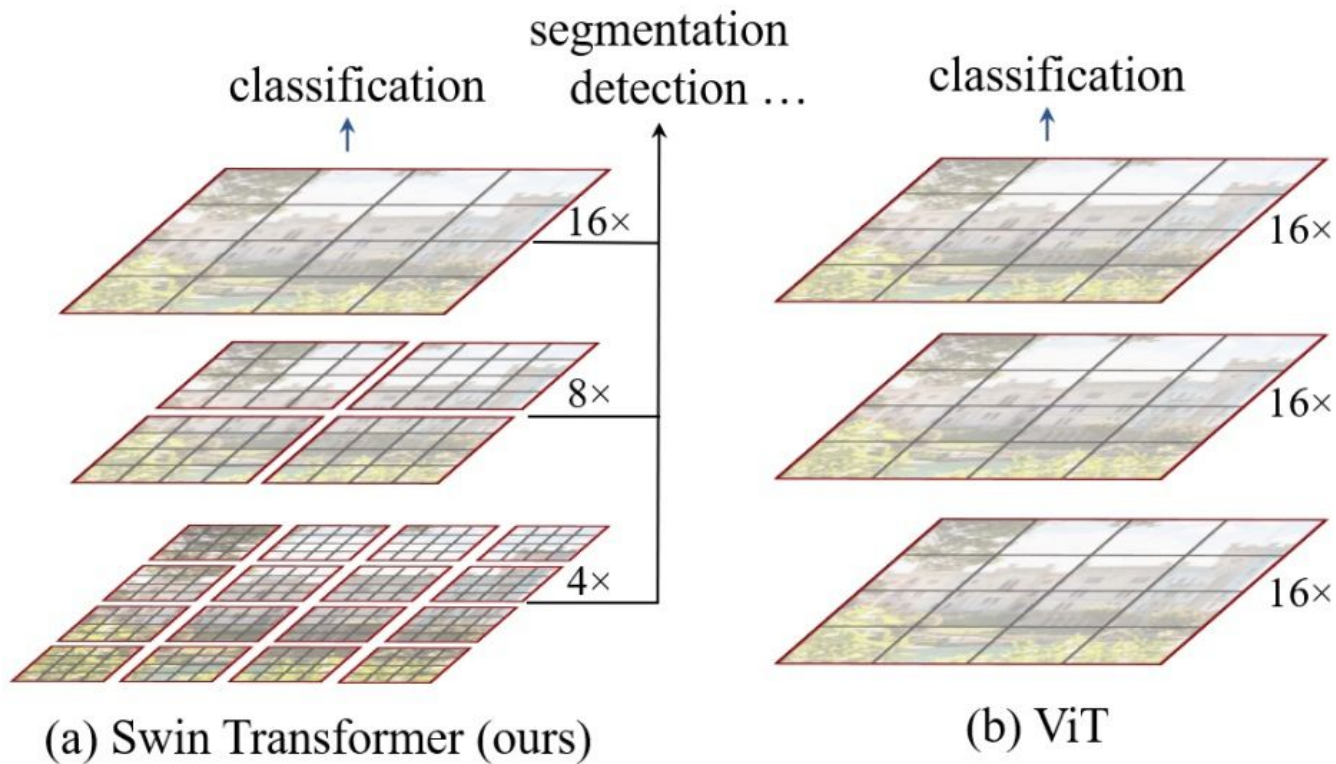
CLIP. Contrastive representation learning



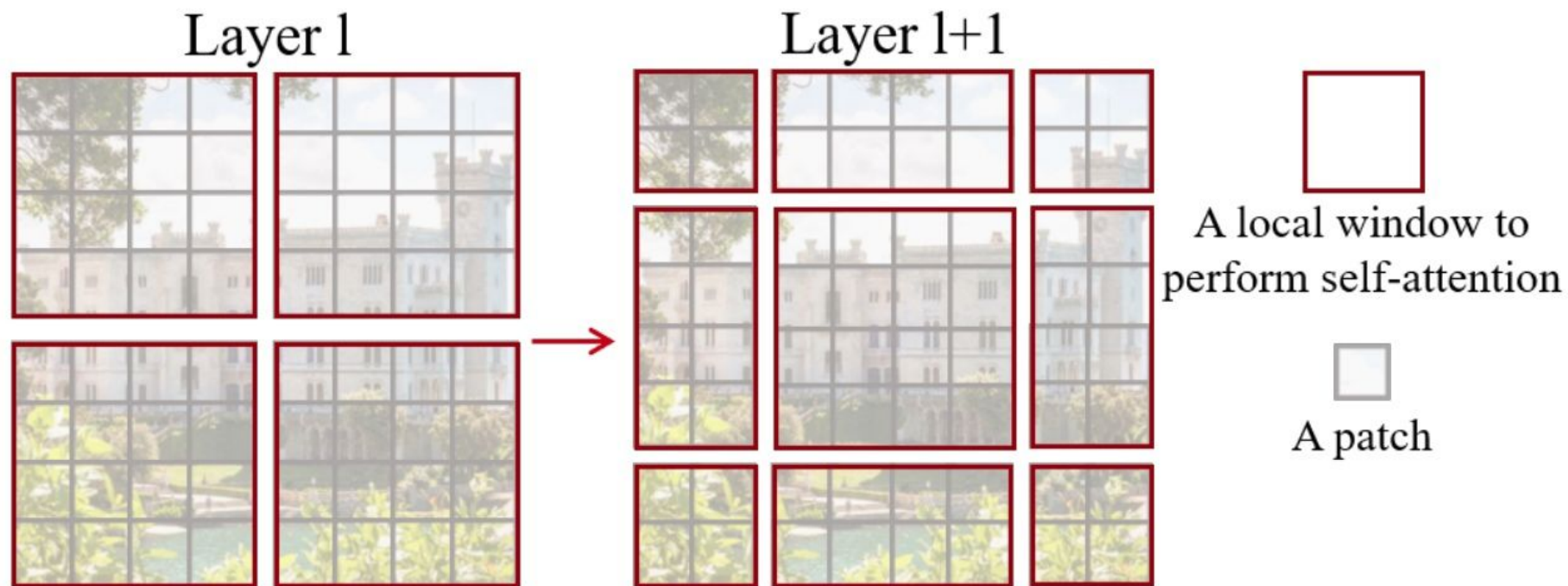
CLIP. Results

	Dataset Examples	ImageNet ResNet101	Zero-Shot CLIP	Δ Score
ImageNet		76.2	76.2	0%
ImageNetV2		64.3	70.1	+5.8%
ImageNet-R		37.7	88.9	+51.2%
ObjectNet		32.6	72.3	+39.7%
ImageNet Sketch		25.2	60.2	+35.0%
ImageNet-A		2.7	77.1	+74.4%

Swin



Swin



Swin

