

Vision Transformers

(for image classification)



Image Classification

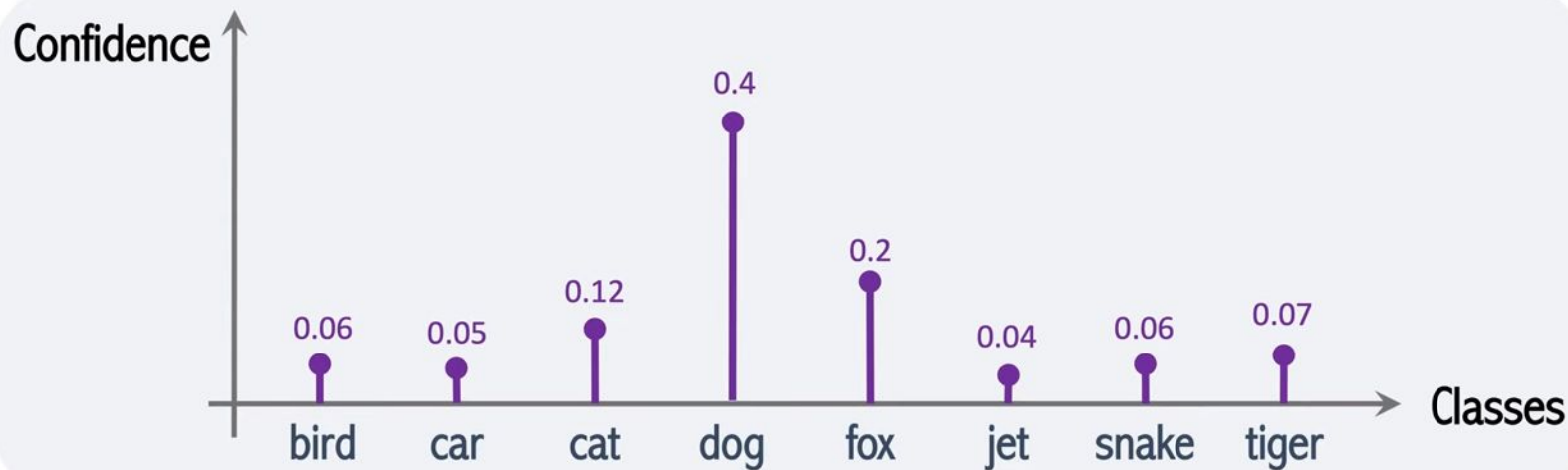
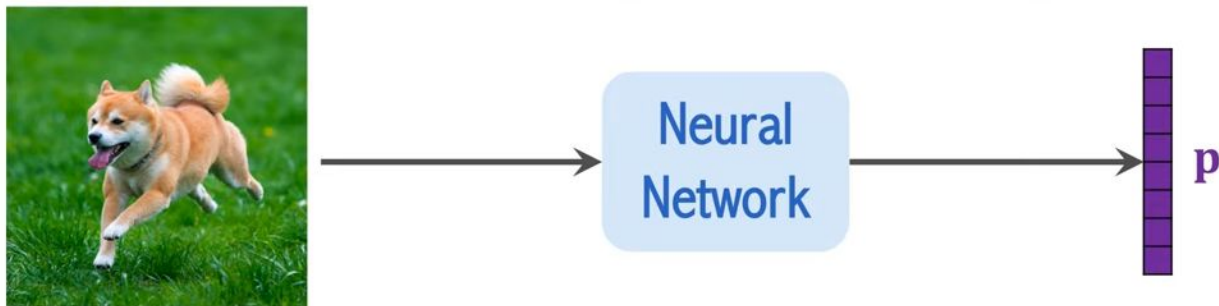
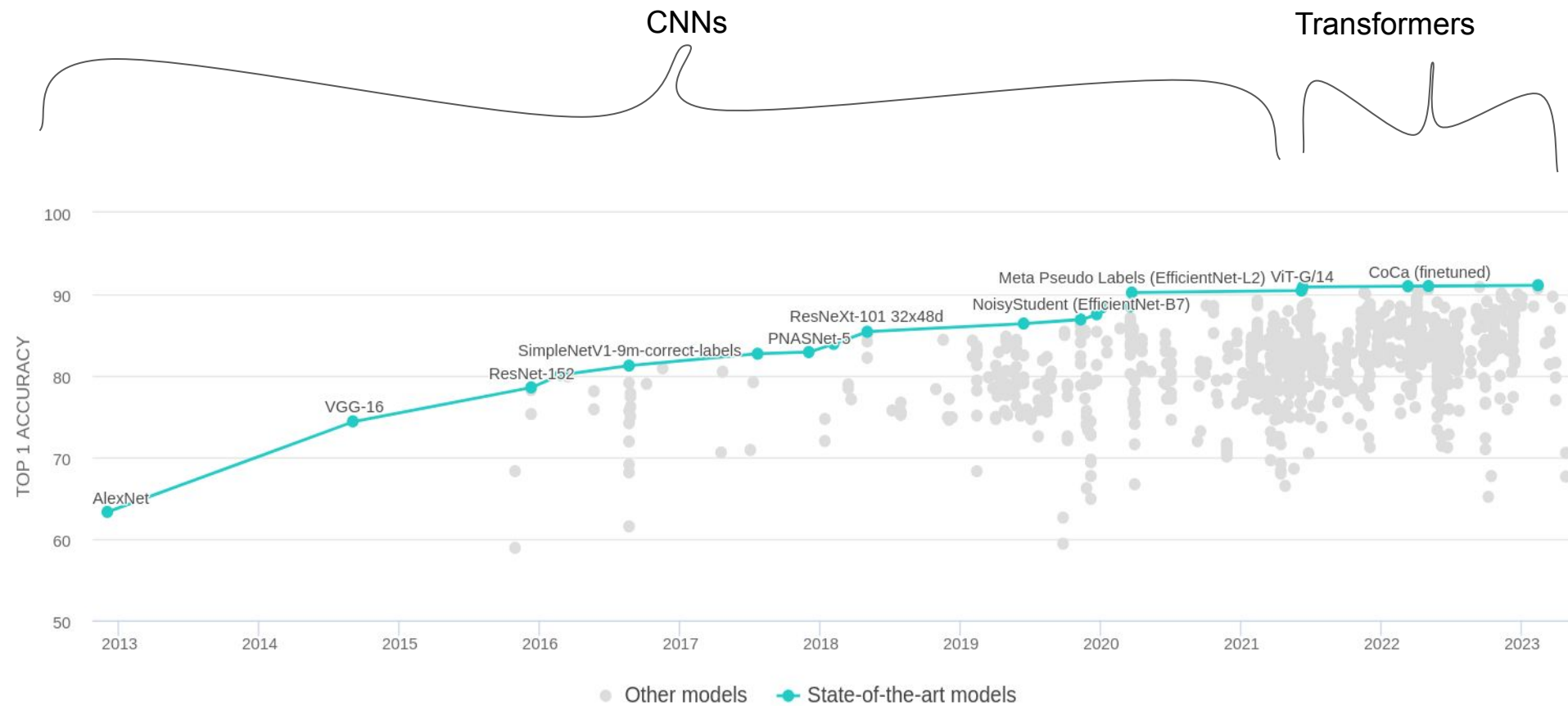
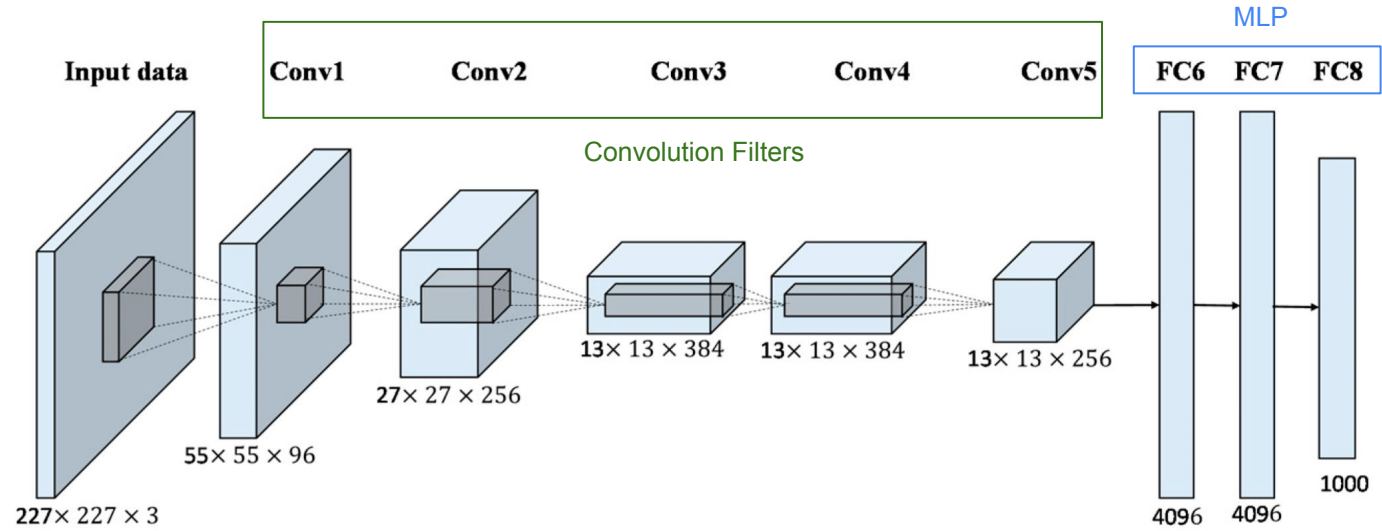


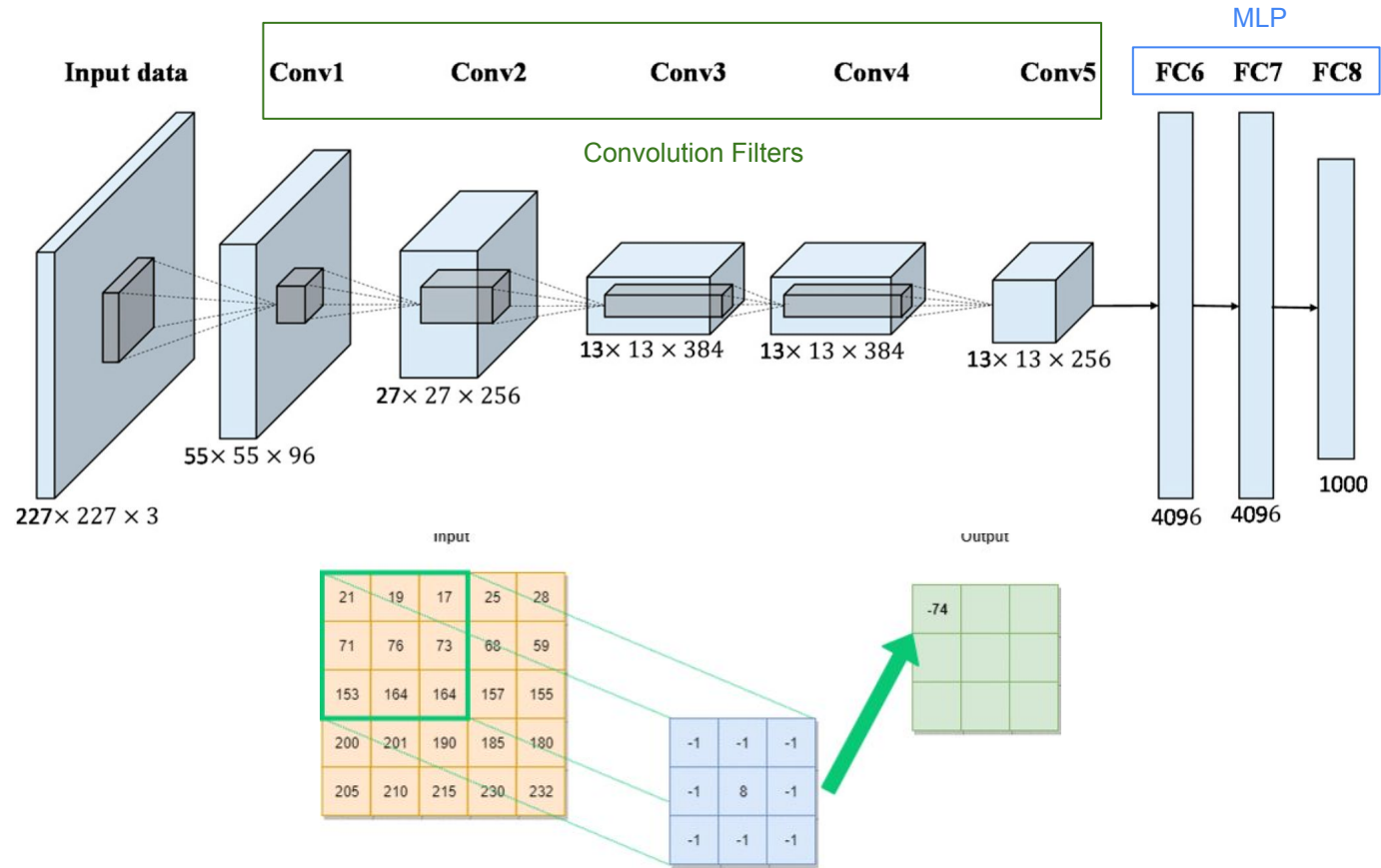
Image Classification SotA



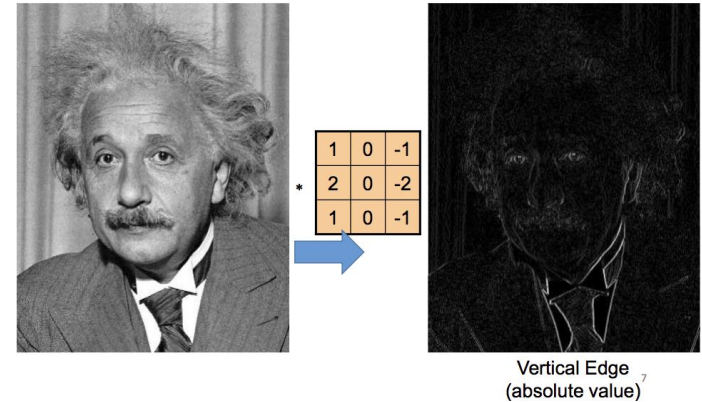
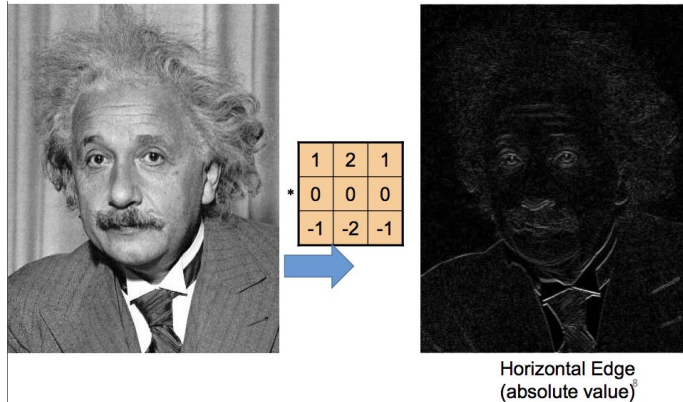
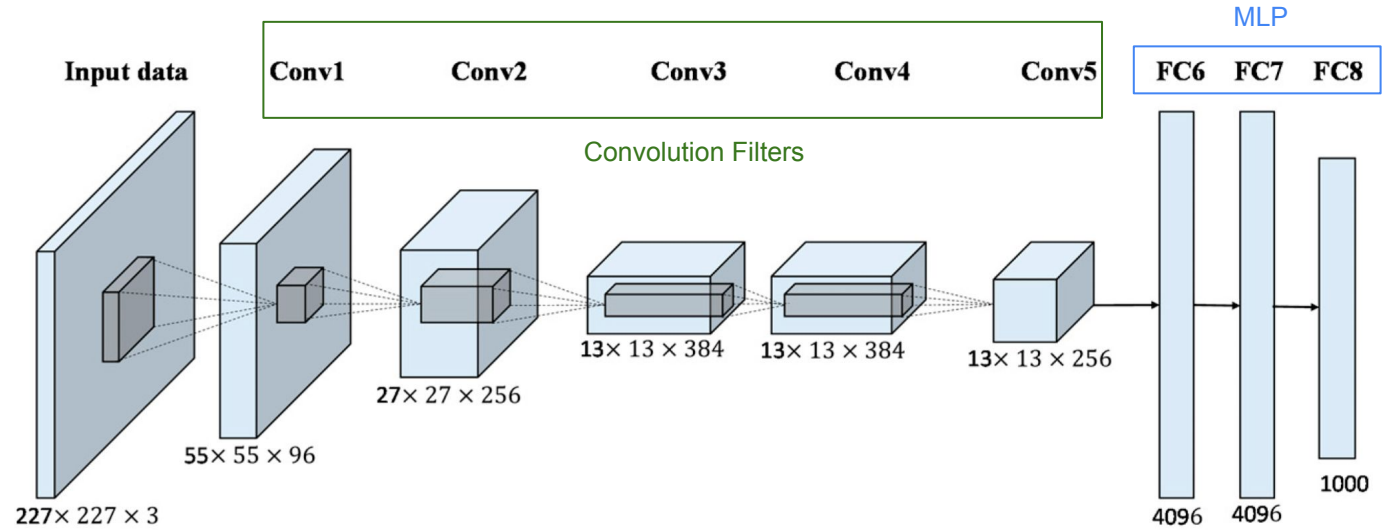
CNN



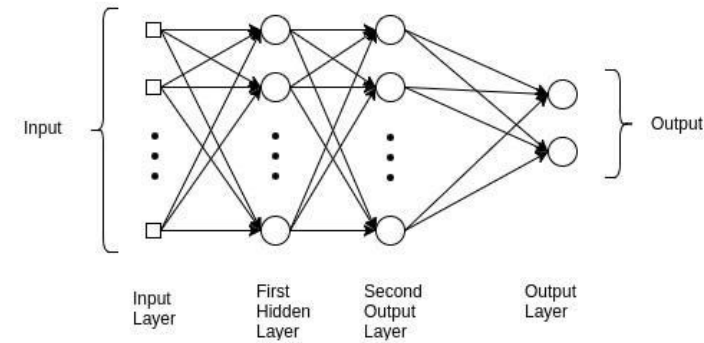
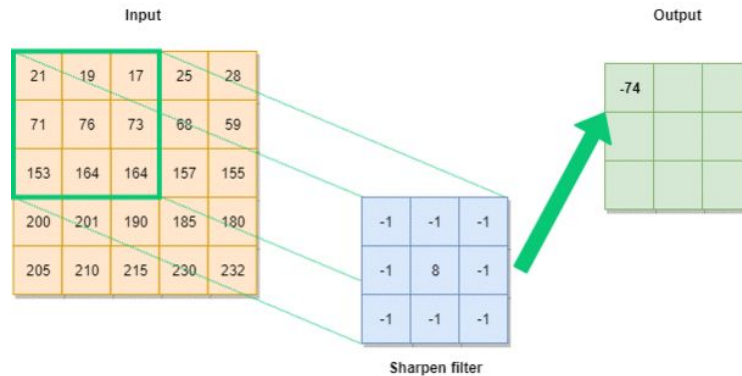
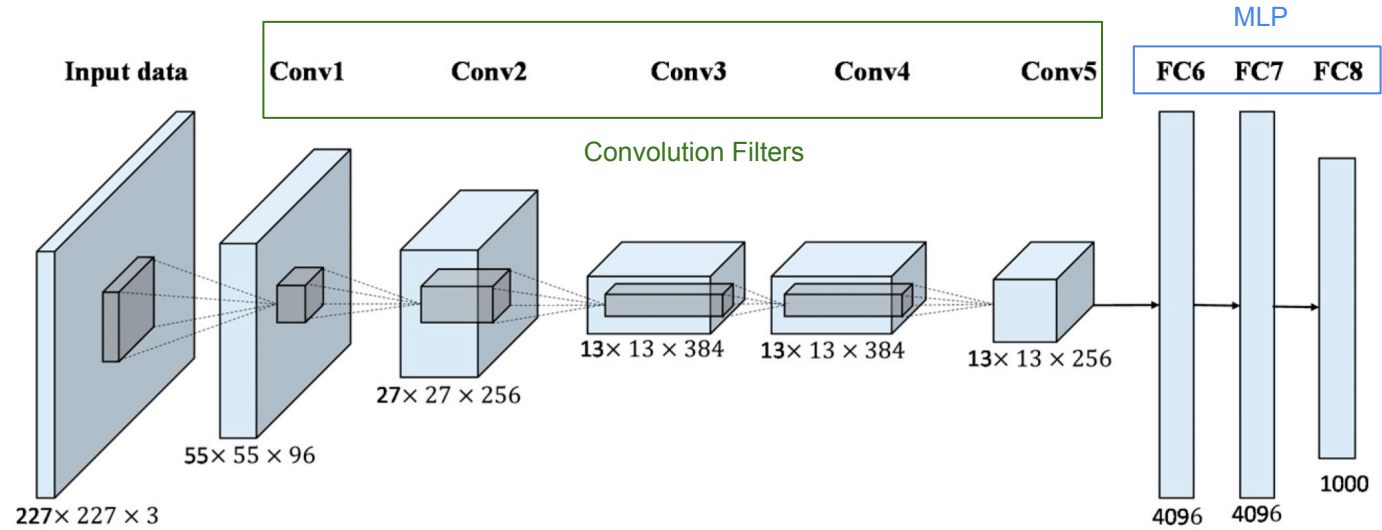
CNN



CNN



CNN

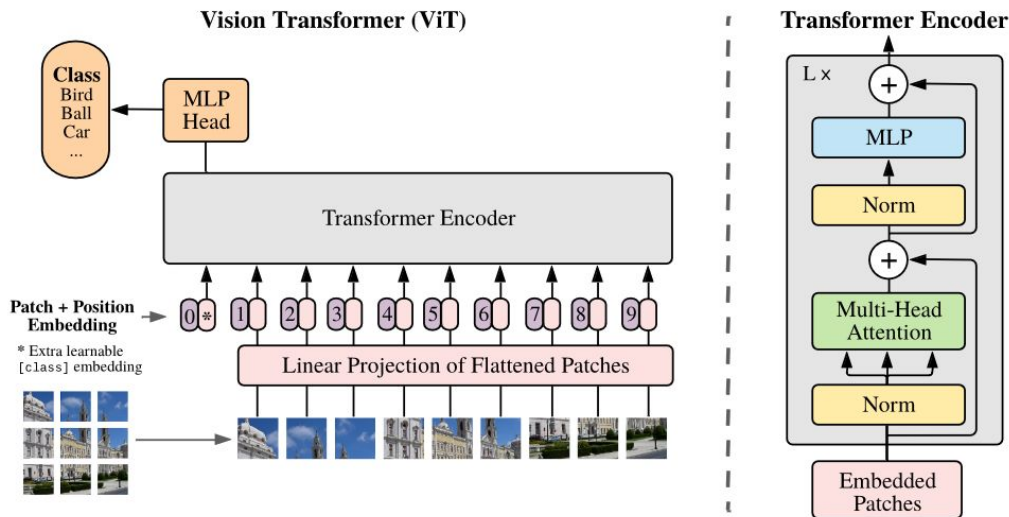


AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

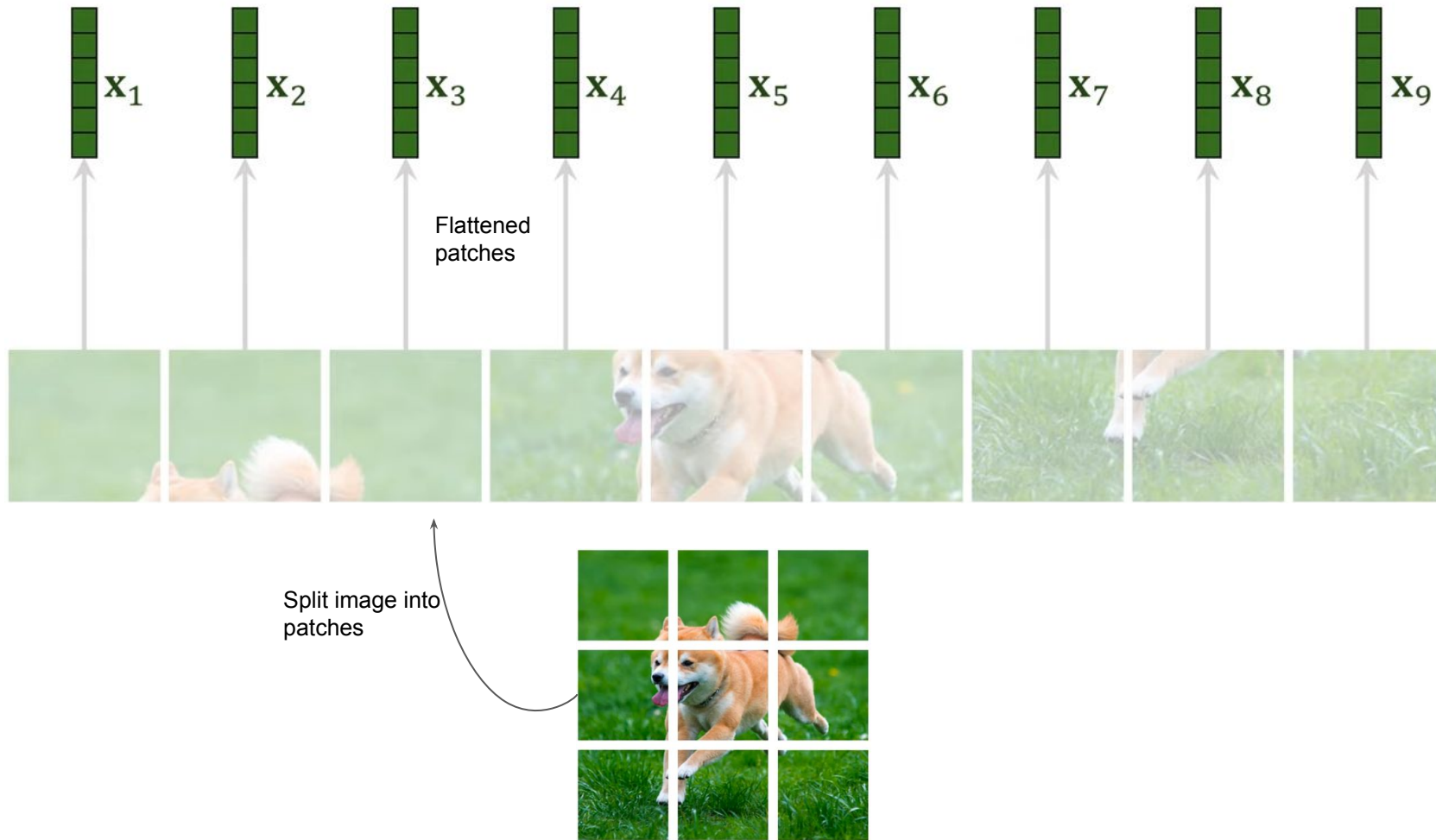
Alexey Dosovitskiy^{*,†}, Lucas Beyer^{*}, Alexander Kolesnikov^{*}, Dirk Weissenborn^{*},
Xiaohua Zhai^{*}, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby^{*,†}

^{*}equal technical contribution, [†]equal advising

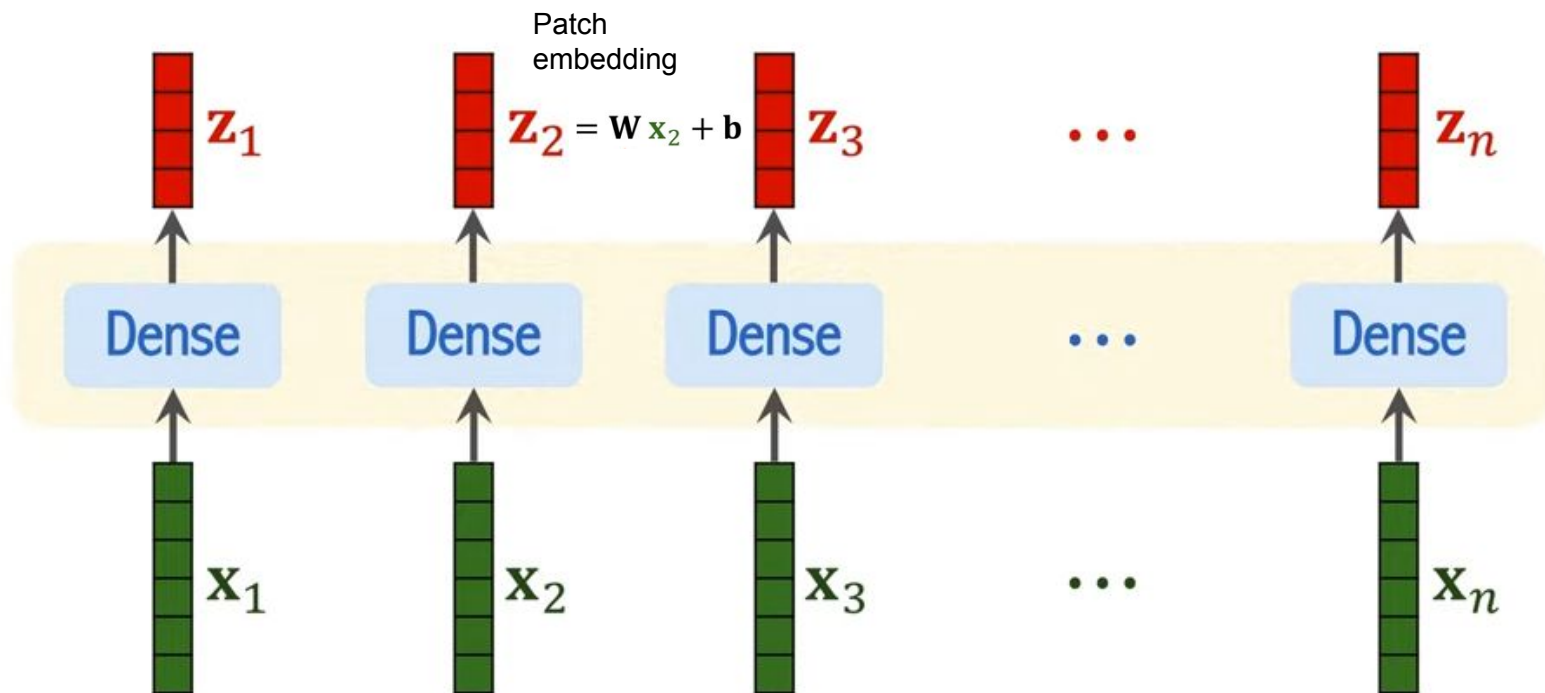
Google Research, Brain Team



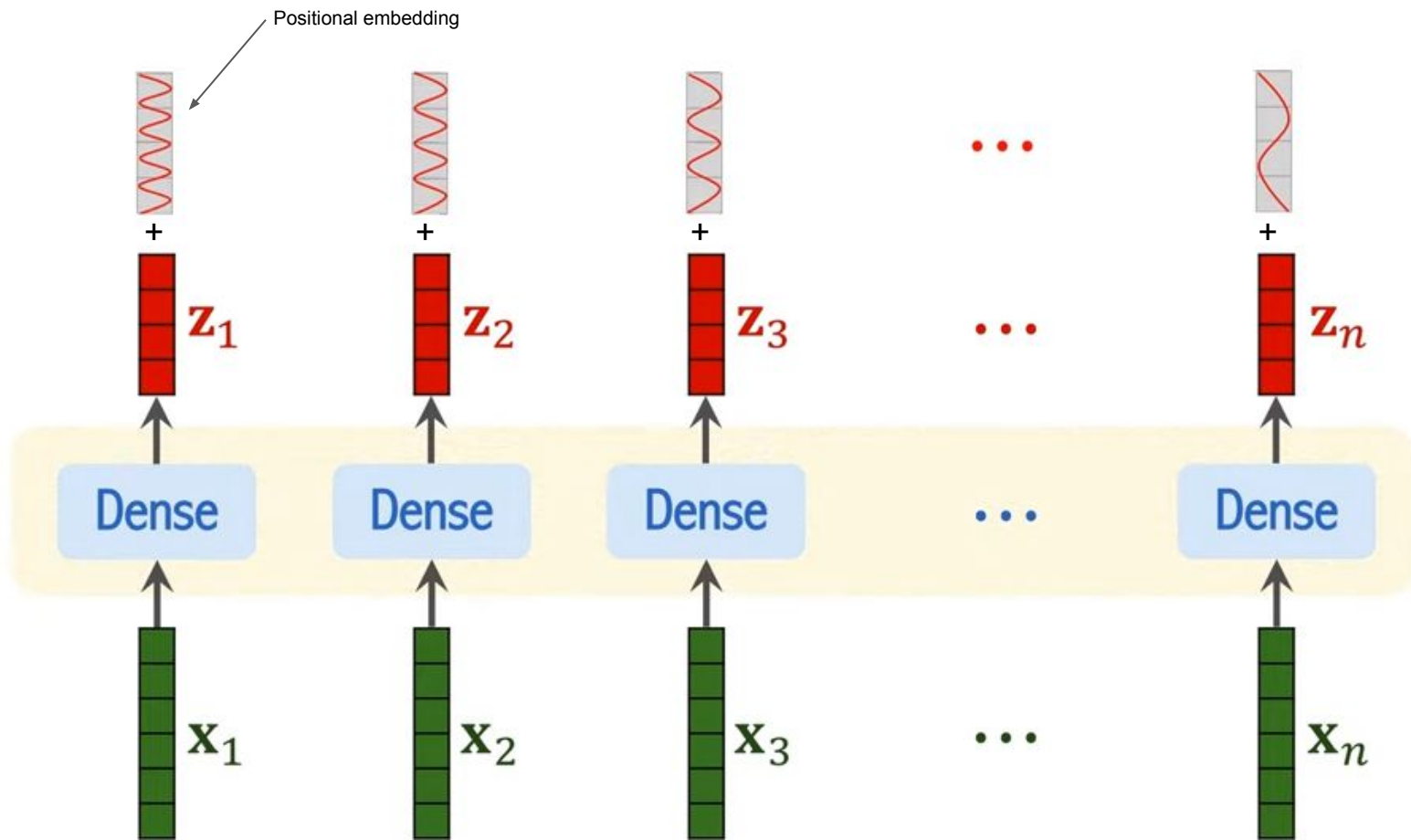
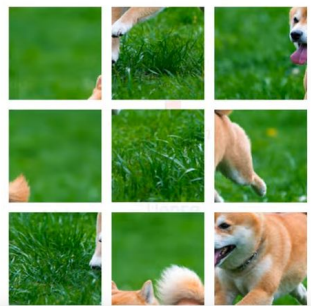
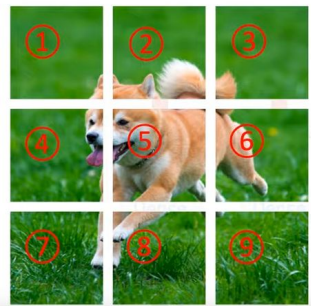
Embedding



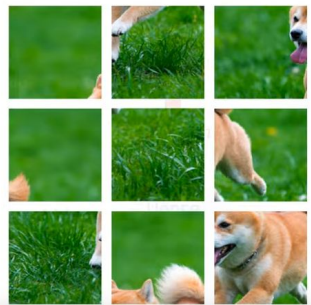
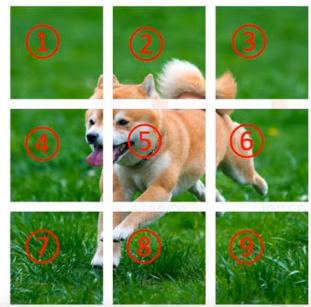
Embedding



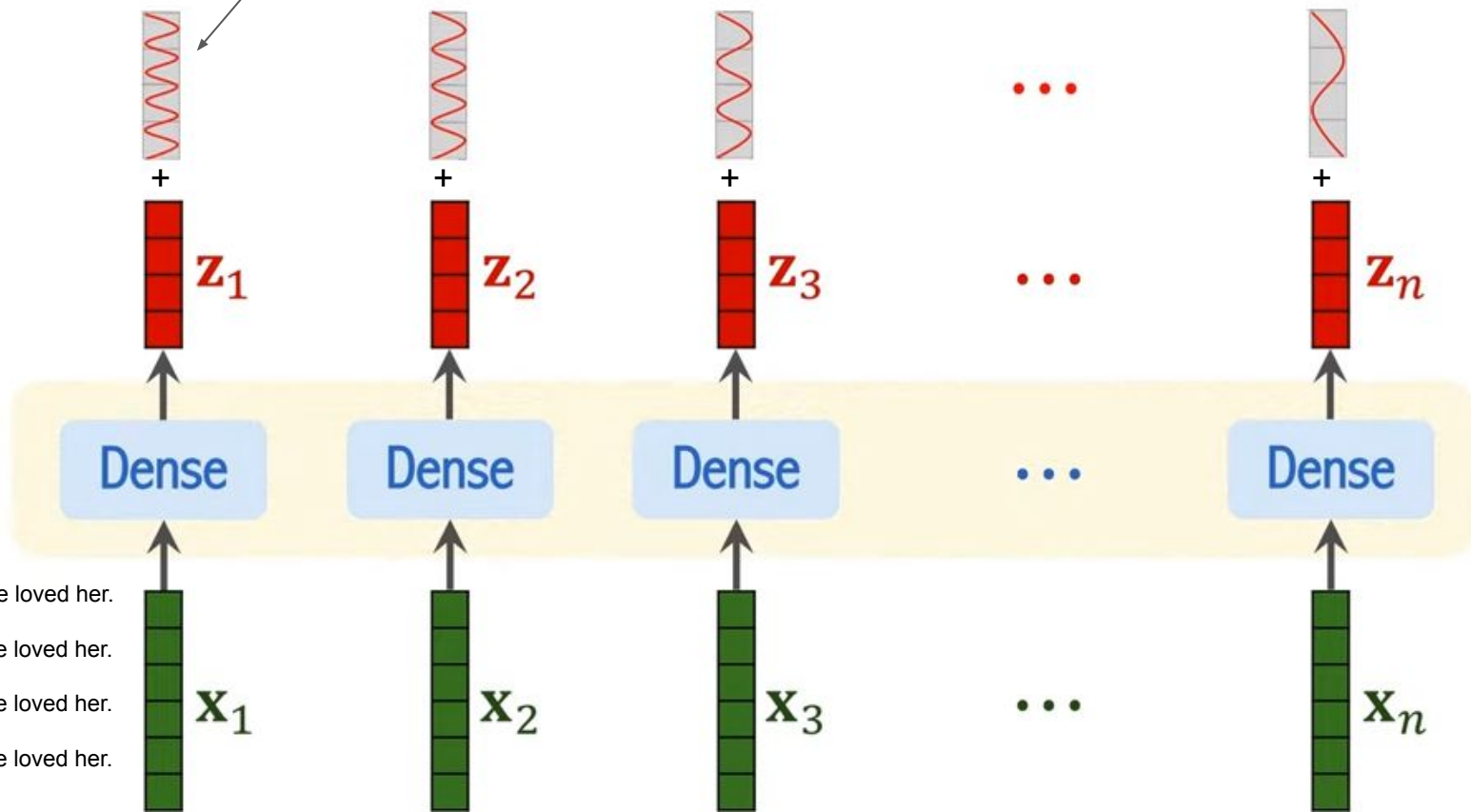
Embedding



Embedding



Positional embedding



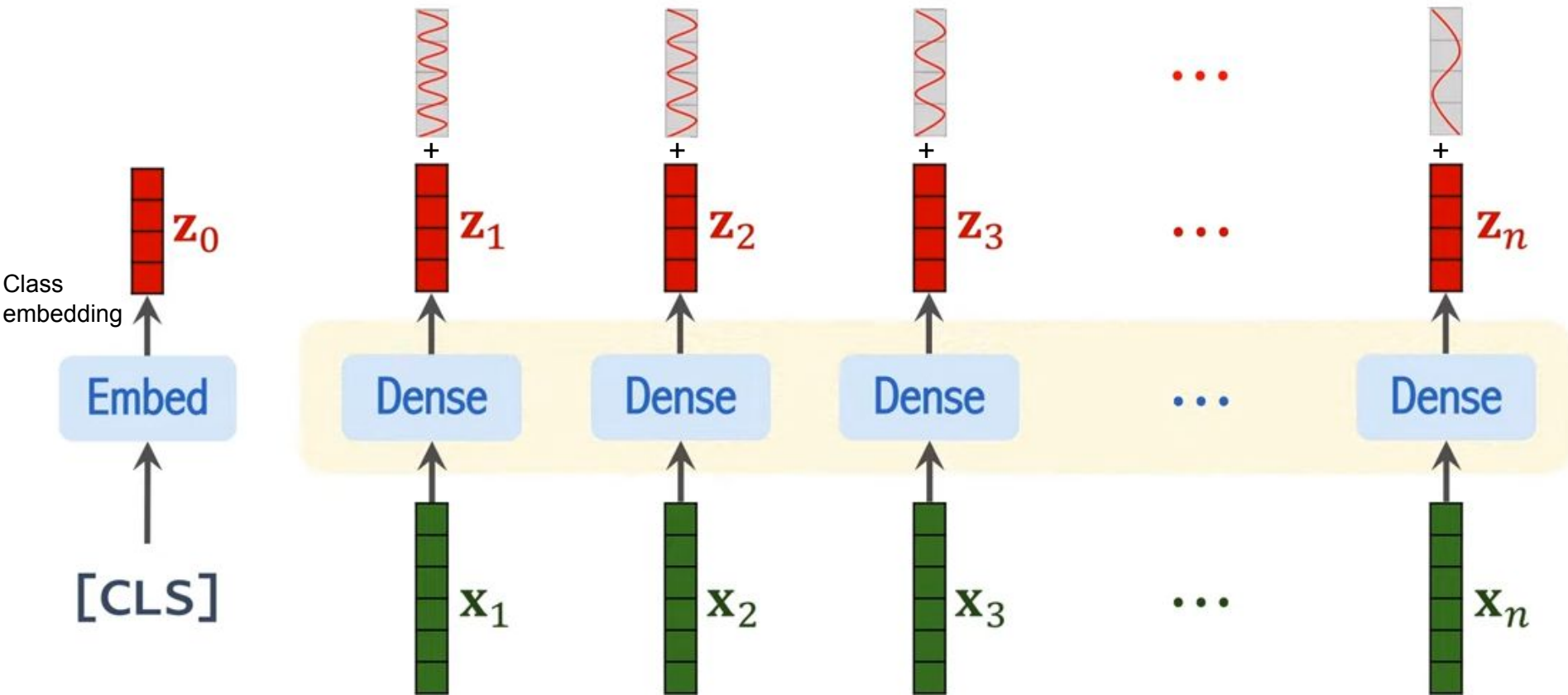
Only he told his mistress that he loved her.

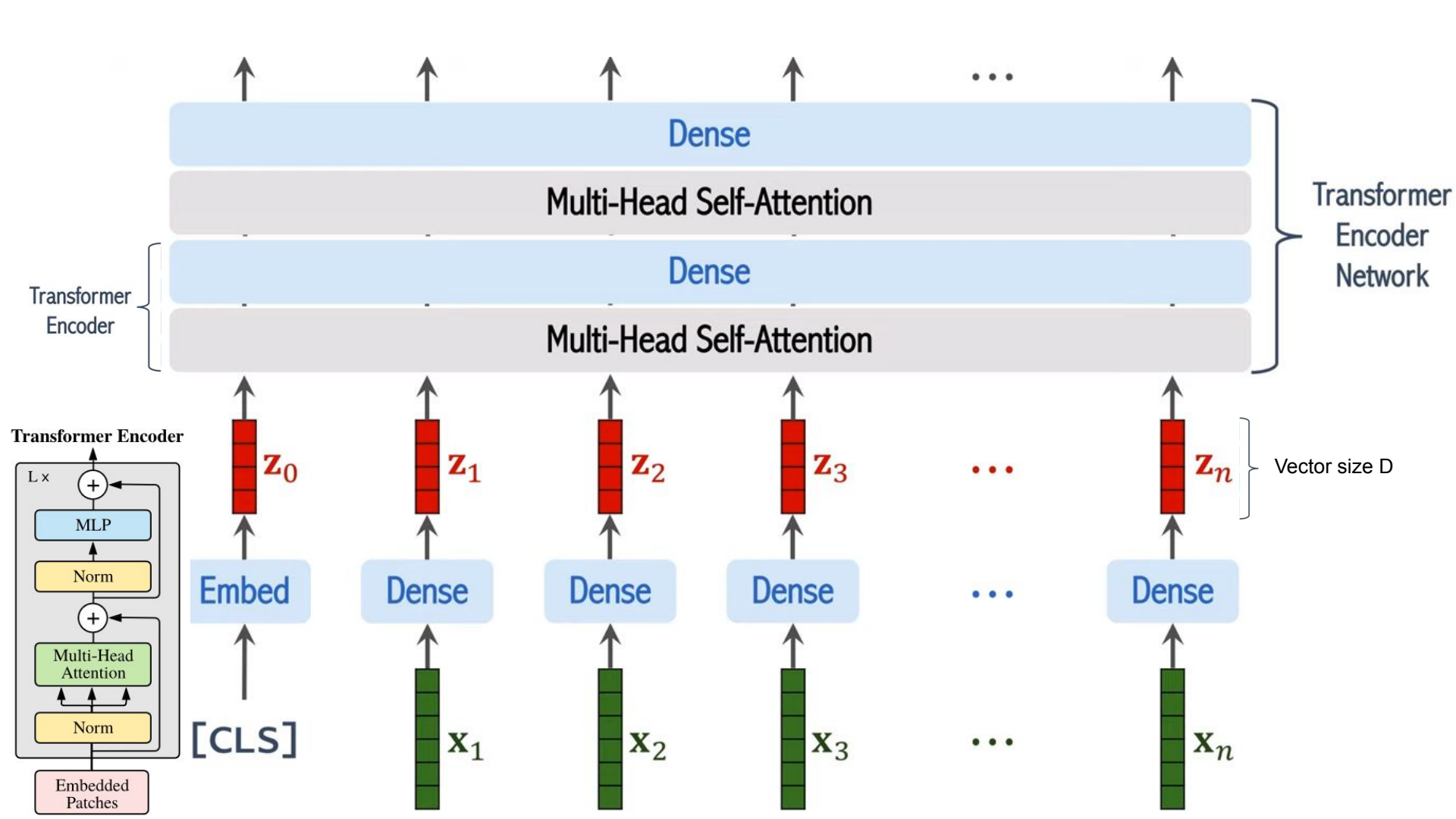
He told only his mistress that he loved her.

He told his only mistress that he loved her.

He told his mistress that only he loved her.

Embedding





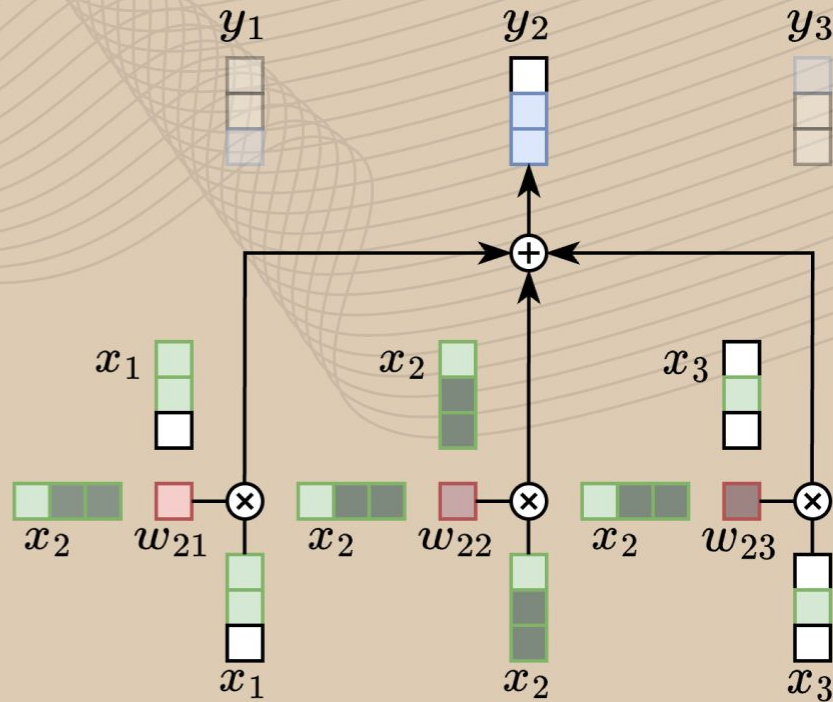
¿Qué es el Self-Attention?

- Es la operación fundamental en cualquier arquitectura de los Transformer
- Es una operación entre secuencias de vectores
- Es la única operación que propaga información entre vectores

$$y_i = \sum_j w_{ij} x_j$$

$$w'_{ij} = x_i^T x_j$$

$$w_{ij} = \frac{e^{w'_{ij}}}{\sum_j e^{w'_{ij}}} \quad (Softmax)$$



Input

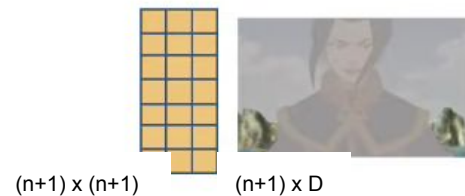


Attention

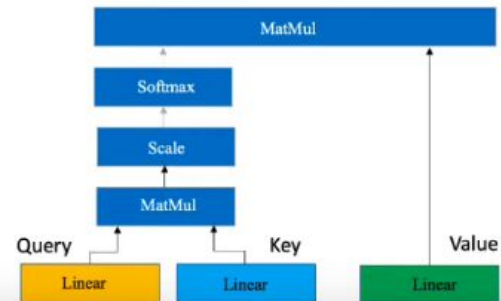
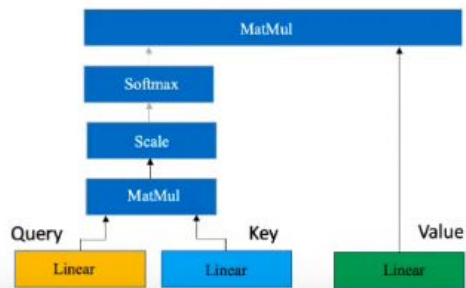
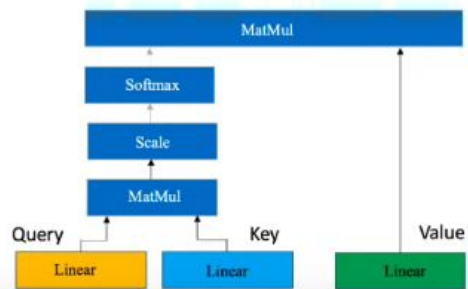
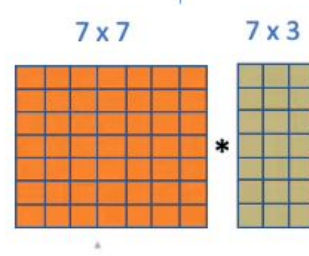
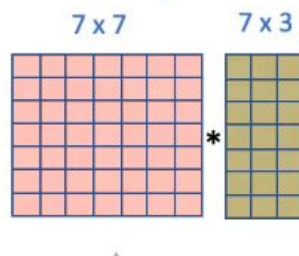
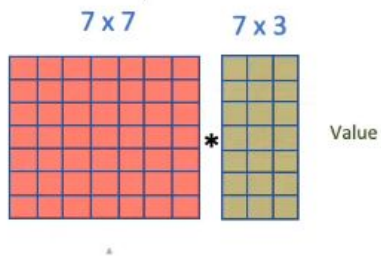


Multi-head self-attention

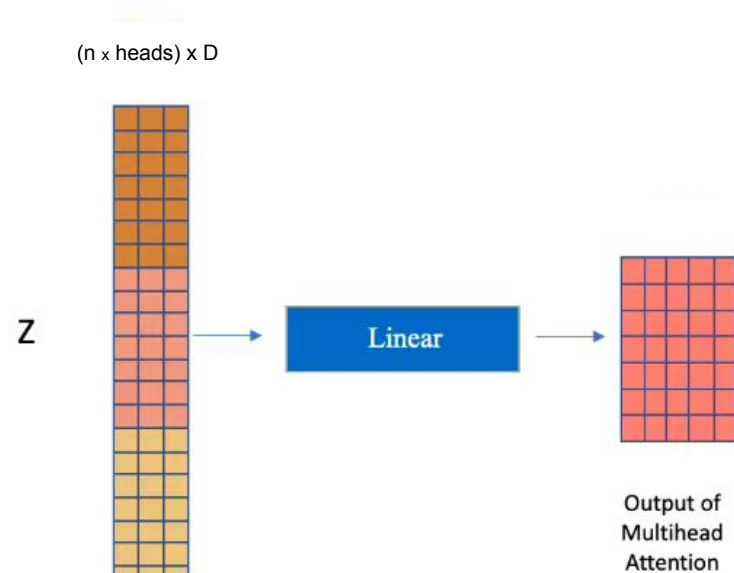
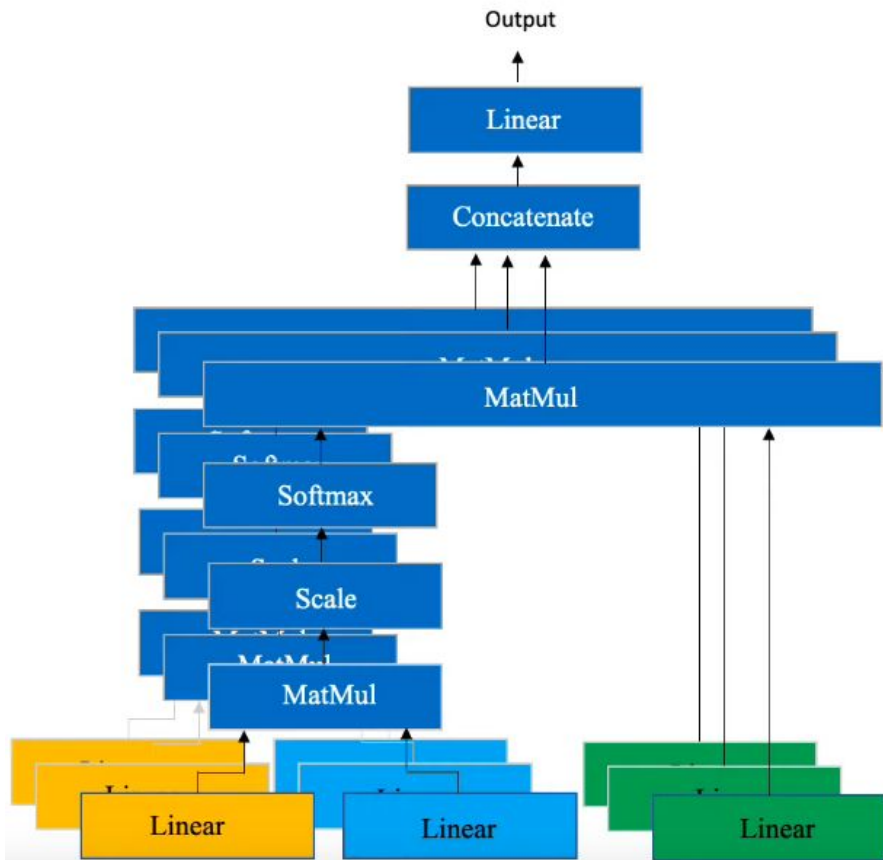
Filtered Value

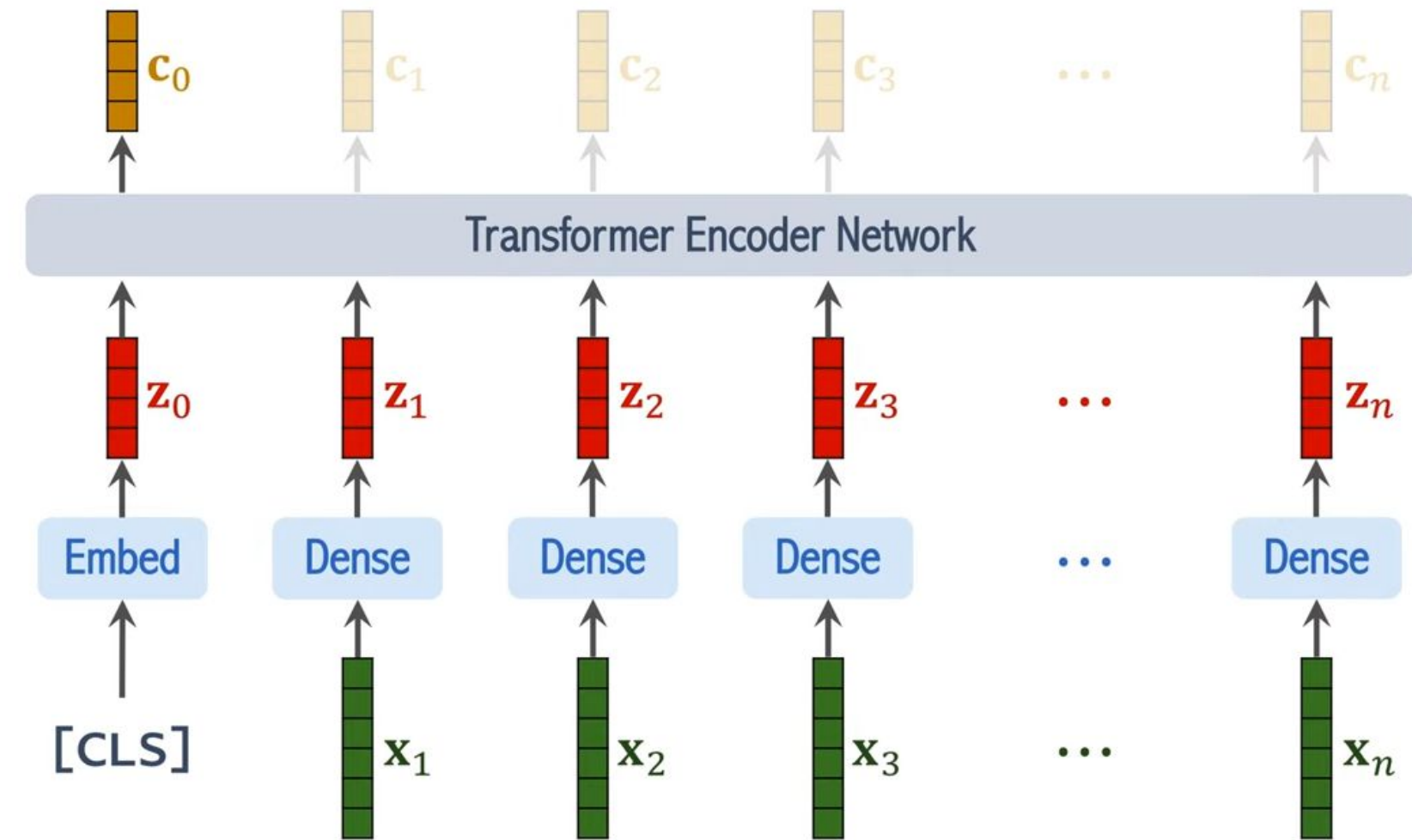


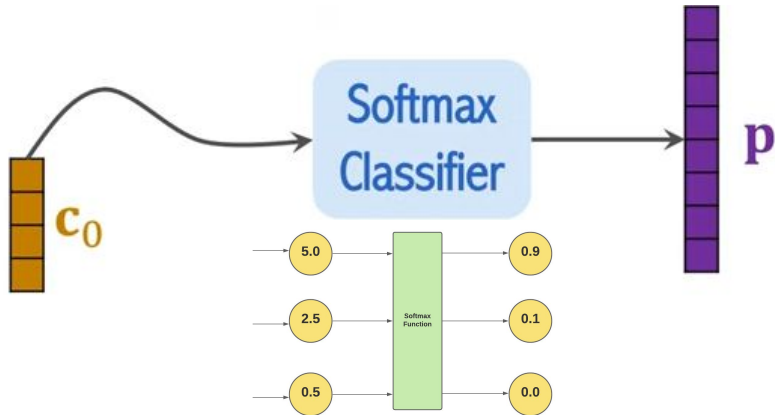
Attention Filter



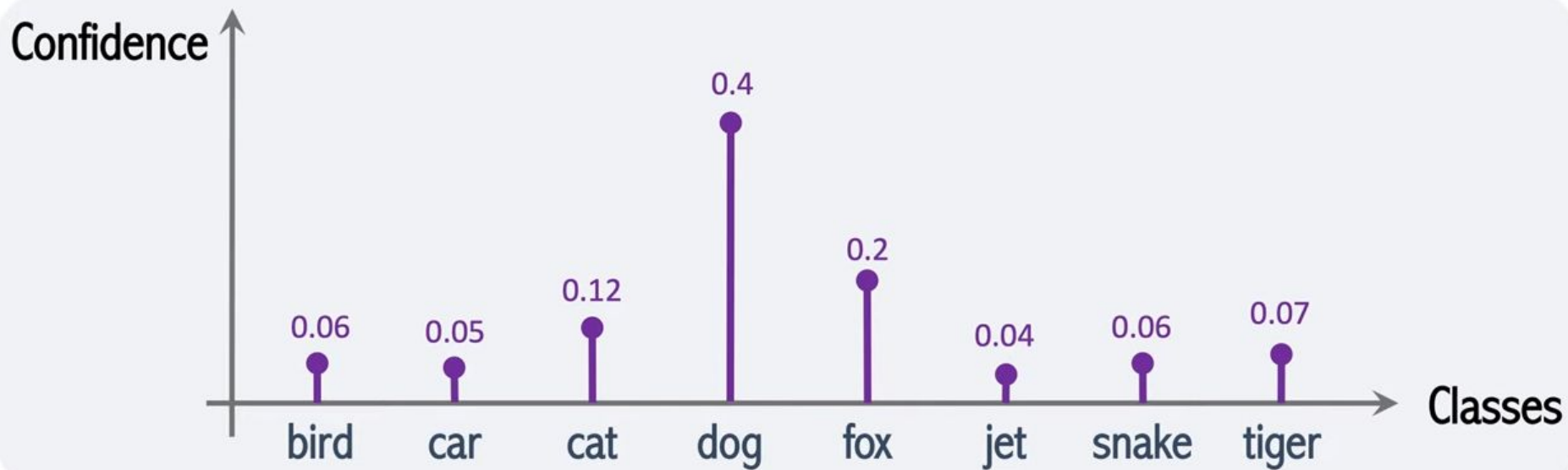
Multi-head self-attention

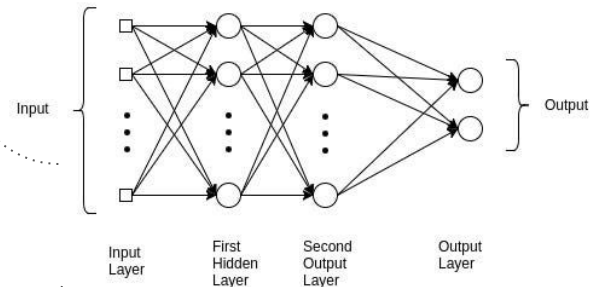
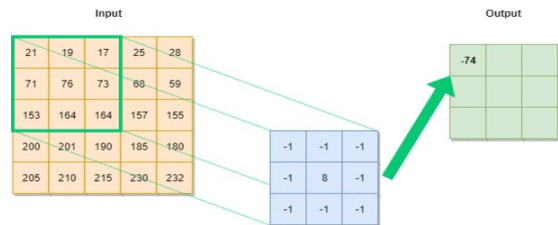
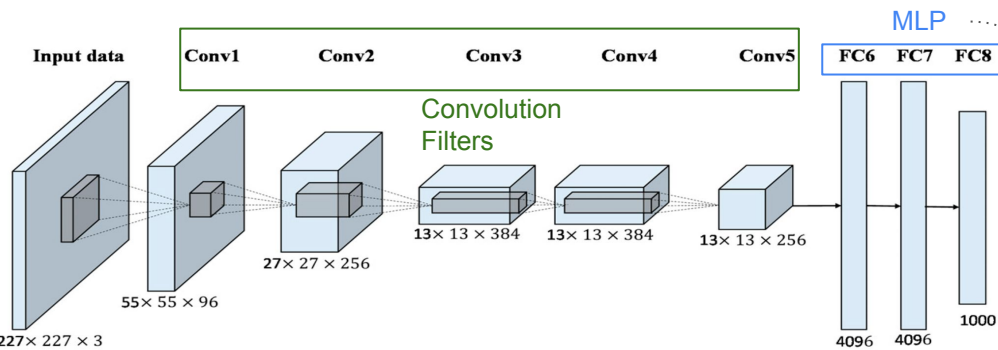




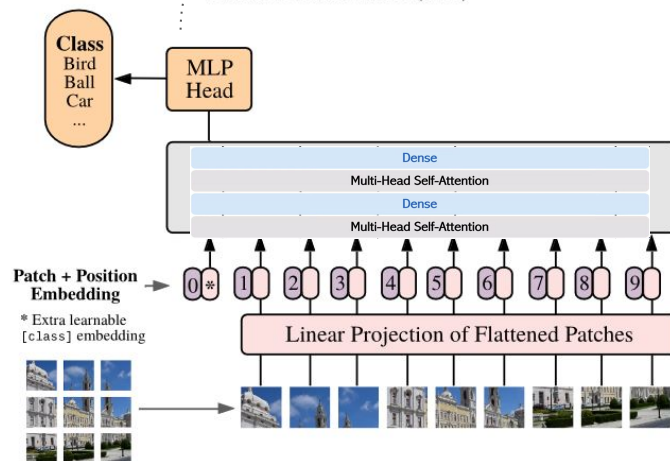


$$\text{Loss} = - \sum_{i=1}^{\text{output size}} y_i \cdot \log \hat{y}_i$$

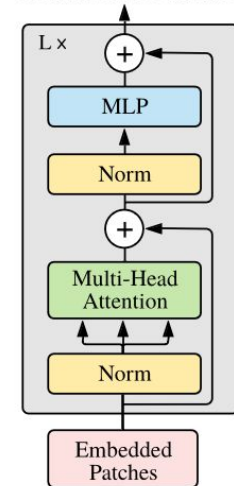


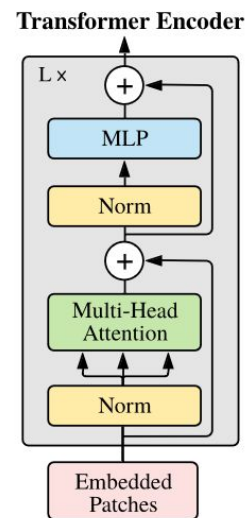
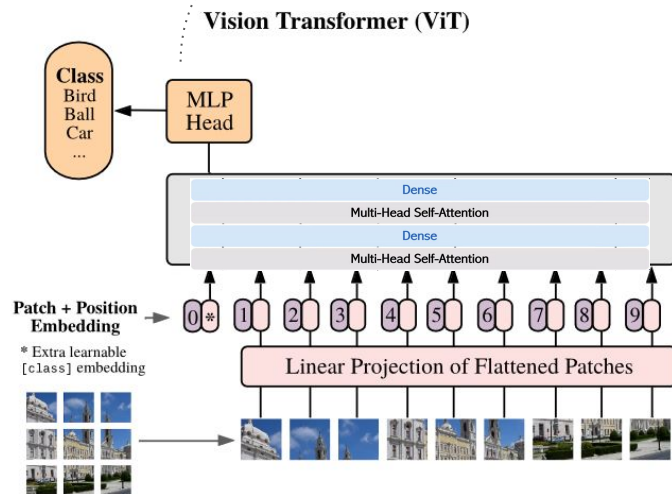
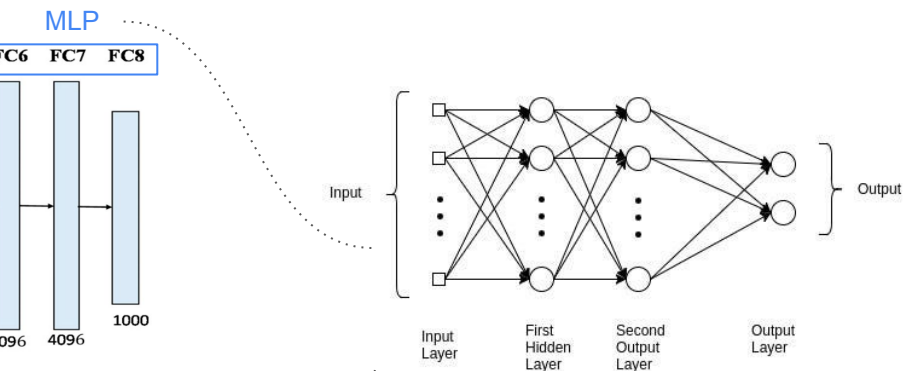
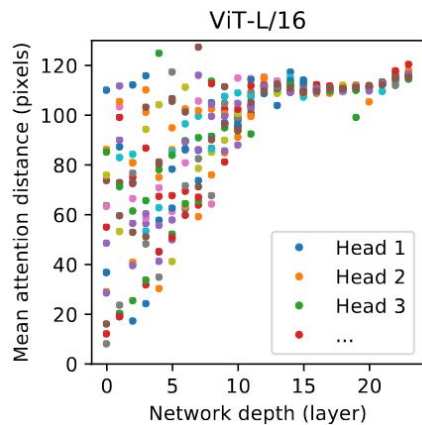
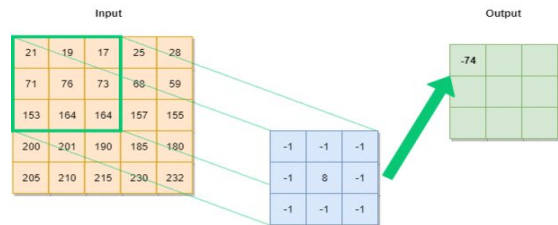
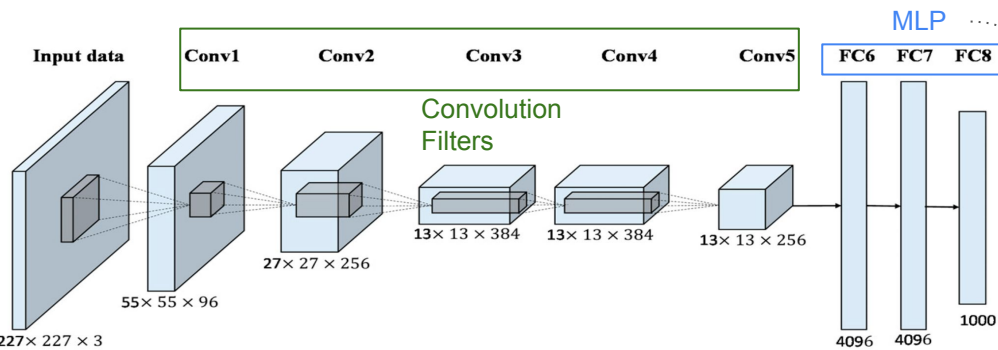


Vision Transformer (ViT)

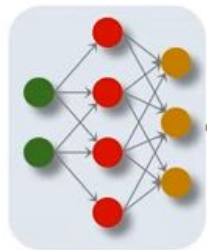


Transformer Encoder

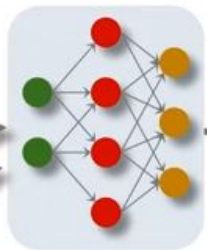




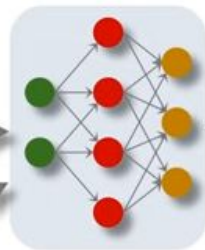
Randomly
Initialized



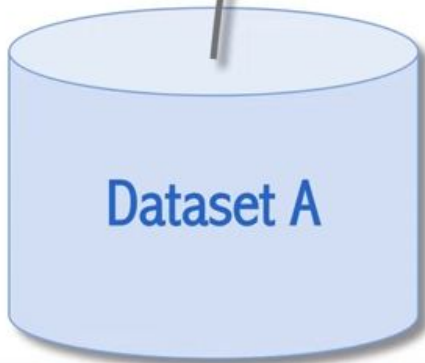
Pretrained



Fine-tuned



Test
Accuracy



Training Set of
Dataset B



Test Set of
Dataset B



	# of Images	# of Classes
ImageNet (Small)	1.3 Million	1 Thousand
ImageNet-21K (Medium)	14 Million	21 Thousand
JFT (Big)	300 Million	18 Thousand

- Pretrain the model on Dataset A, fine-tune the model on Dataset B, and evaluate the model on Dataset B.
- Pretrained on ImageNet (small), ViT is slightly worse than ResNet.
- Pretrained on ImageNet-21K (medium), ViT is comparable to ResNet.
- Pretrained on JFT (large), ViT is slightly better than ResNet.

