

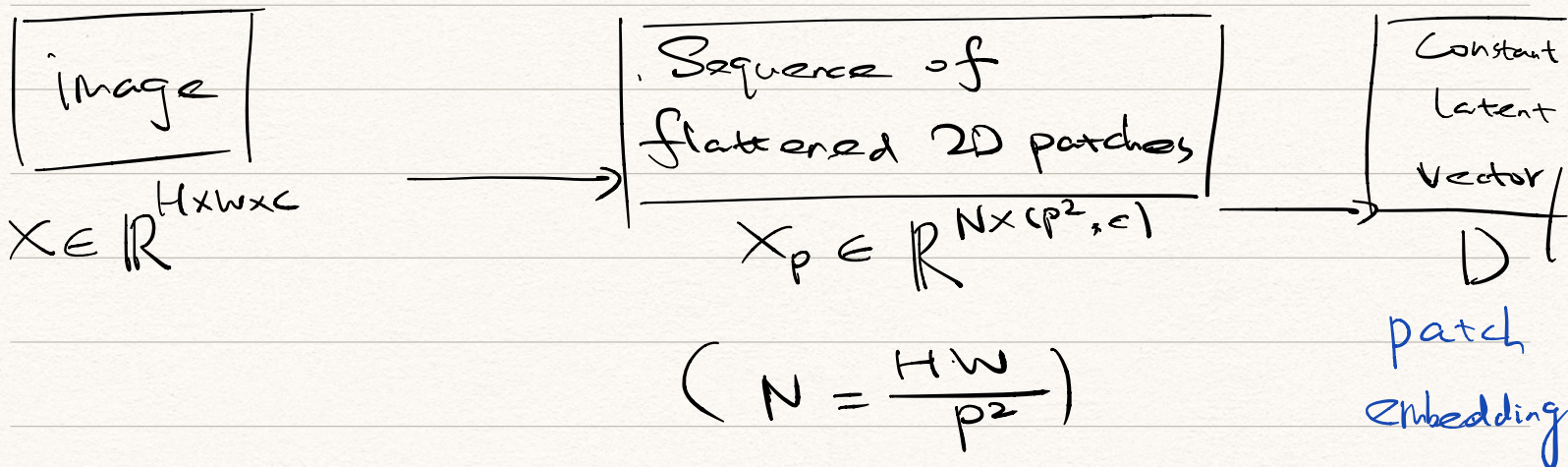
Transformer

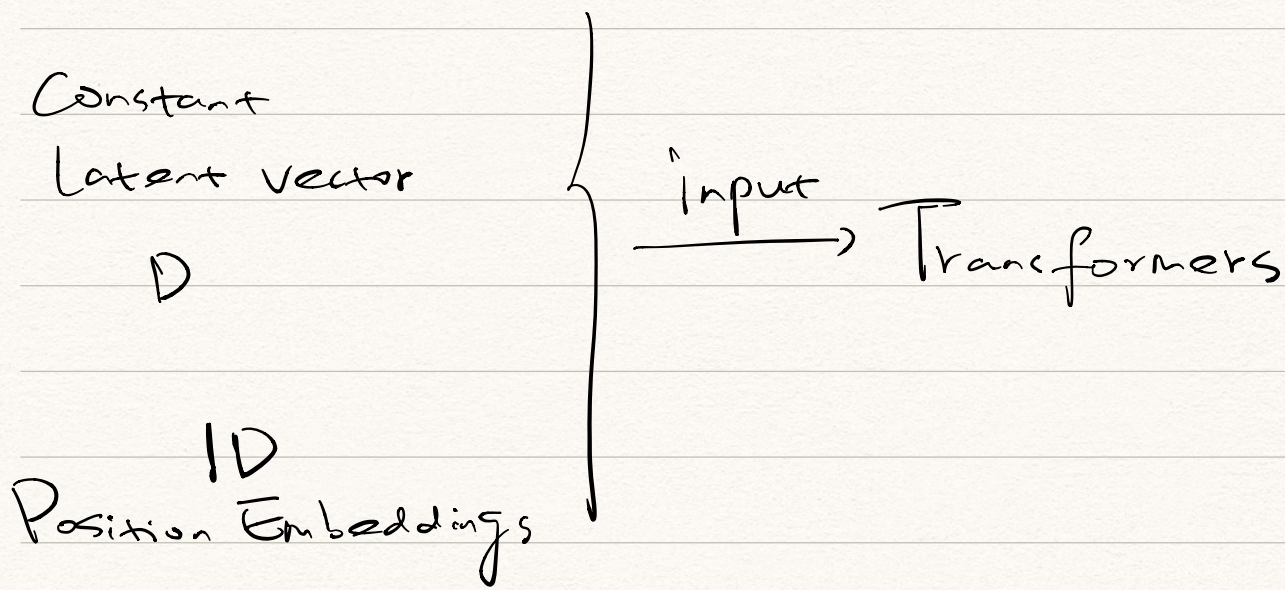
Pros: {

- computational efficiency
- scalability
- can handle arbitrary sequence lengths

Cons: lack some of the inductive biases inherent to CNNs, such as translation equivariance and locality, and therefore do not generalize well when trained on insufficient amounts of data.

Vision Transformer:





MSA: Multiheaded Self-attention

MLP: GELU non-linearity

$$z_0 = [X_{\text{class}} ; X_p^1 E ; X_p^2 E ; X_p^3 E ; \dots ; X_p^N E] + \bar{E}_{\text{pos}}$$

$$z'_1 = \text{MSA}(\text{LN}(z_{L-1})) + z_{L-1}$$

$$z_L = \text{MLP}(\text{LN}(z'_1)) + z'_1$$

$$y = \text{LN}(z_L)$$

$$\bar{E} \in \mathbb{R}^{(P^2 \times C) \times D} ; \bar{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D} ; l = 1, \dots, L$$

Attach a zero-initialized $D \times k$ feed forward layer

Fine-tuning Accuracies: capture the performance of each model after fine-tuning it on the respective dataset,

Few-shot accuracies: obtained by solving a regularized linear regression problem that maps the (frozen) representation of a subset of training images to $\{-1, 1\}^k$ target vectors.

(for fast on-the-fly evaluation where fine-tuning would be too costly)