

Data Efficient Image Transformer (DeiT)

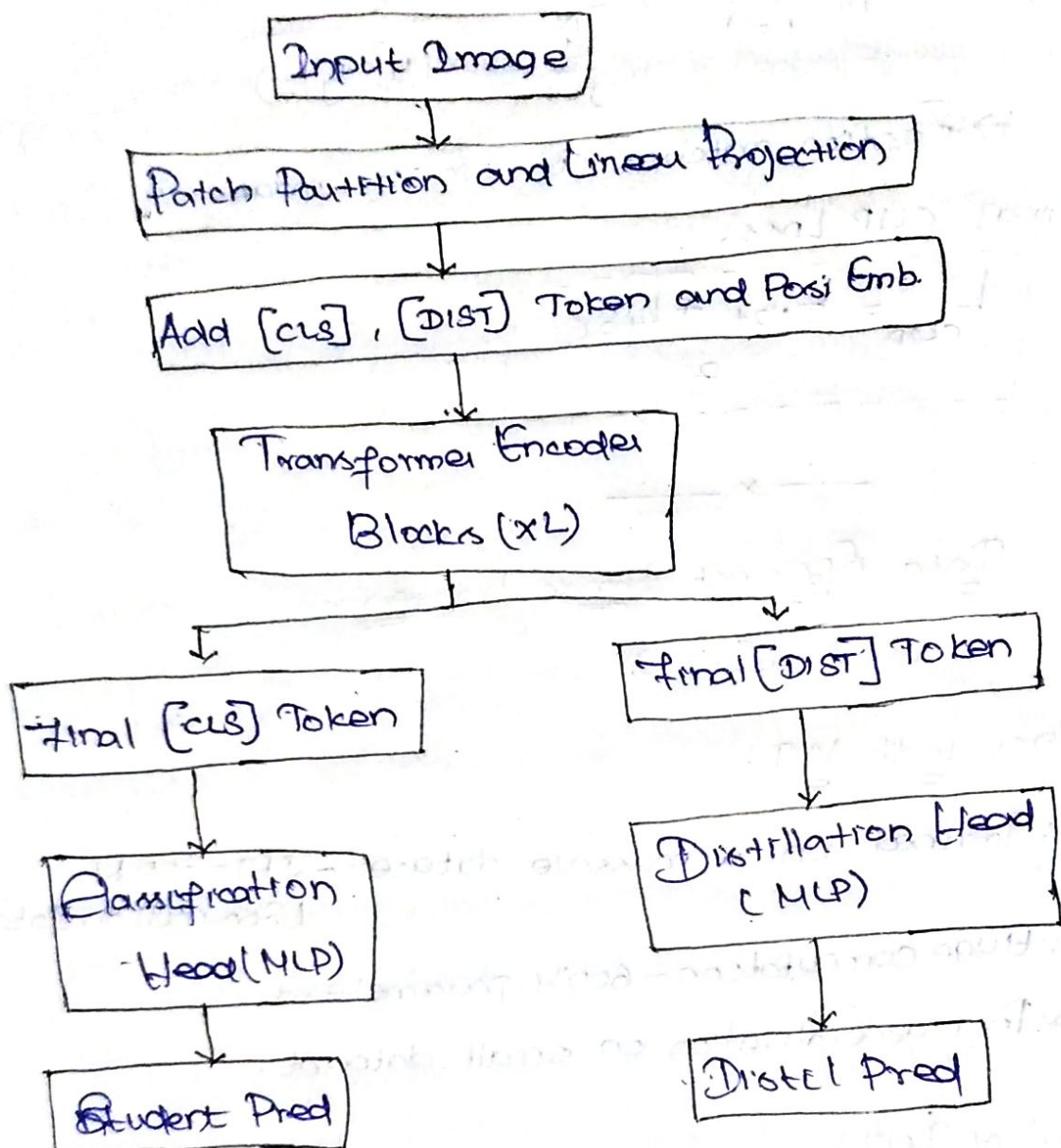
Problems with ViT:

- * Trained on a massive dataset - JFT-300M
(300 million dataset)
- * Huge computation - 600M parameters
- * Poor generalization on small dataset.

Goal of DeiT:

- * Train a Vision Transformer efficiently using limited data and reasonable compute, without sacrificing accuracy.
- use ImageNet-1K dataset (~1.2M Images)
- use smaller models (~84M params)
- achieve performance comparable to SOTA CNNs.

Block Diagram:



Loss vs Teacher Output.

Loss vs Ground Truth

⇒ The teacher output comes from a separate, pre-trained teacher model (often a CNN like ResNet)

*Dist = Vct + one additional token.

- CLS token learns from ground truth.
- DIST token learns from teacher.

Classification Loss: (CLS Token)

$$L_{CE} = - \sum_i y_i \log p_i^{CLS}$$

→ ground truth
→ CLS Token.

Distillation Loss:

$$L_{KD} = T^2 \sum p_c^{\text{teacher}} \log \left(\frac{p_c^{\text{teacher}}}{p_c^{\text{dist}}} \right)$$

* T is used to control the softmax sharpness before computing the loss.

* Extra T^2 is introduced to cancel $(1/T^2)$ while computing backprop.

Final Loss:

$$L = \alpha L_{CE} + (1-\alpha) L_{KD} \quad \text{Typically } \alpha = 0.5$$