

Based on the DynamicViT paper, the loss function that constrains the pruning ratio is implemented as an MSE loss between the target pruning ratio and the actual pruning ratio at each stage. The mathematical implementation is:

$$L_{\text{ratio}} = (1/BS) * \sum_{s=1}^S (\frac{1}{N} \sum_{i=1}^N D_{b,s,i})^2$$

Where:

- B is batch size
- S is number of pruning stages (typically 3)
- $\frac{1}{N} \sum_{i=1}^N D_{b,s,i}$  is ratio at stage s (e.g.,  $\frac{1}{N} \sum_{i=1}^N D_{b,s,i}$ )
- N is number of tokens
- $D_{b,s,i}$  is binary decision mask for sample b at stage s
- The term  $(1/N) \sum_{i=1}^N D_{b,s,i}$  calculates actual keeping ratio

This loss is combined with other losses (classification, distillation, KL divergence) with a weighting factor  $\mu$  (set to 2 in their experiments).

The full training objective is:

$$L = L_{\text{cls}} + \mu L_{\text{distill}} + L_{\text{ratio}}$$

This MSE loss ensures the model maintains the desired pruning ratio at each stage while allowing some flexibility for input-dependent pruning decisions.