

The paper "DynamicViT: Efficient Vision Transformers with Dynamic Token Sparsification" proposes an attention masking strategy to implement differentiable token pruning in vision transformers. Here's a mathematical explanation of their attention masking implementation:

1. Standard Attention Calculation:

The standard attention matrix A is computed as:

$$A = \text{Softmax}(QK^T / C)$$

2. Attention Masking Implementation:

To prune tokens while maintaining differentiability, they:

a) First compute the unnormalized attention scores:

$$P = QK^T / C \in \mathbb{R}^{N \times N}$$

b) Create a binary mask G based on the pruning decisions \hat{D} :

$$G_{\{ij\}} = \begin{cases} 1, & \text{if } i = j \text{ (self-connection)} \\ \hat{D}_j, & \text{if } i \neq j \end{cases}$$

This means:

- A token always attends to itself (self-loop)
- A token only attends to other tokens if they're not pruned ($\hat{D}_j = 1$)

c) Compute the masked attention matrix:

$$\tilde{A}_{\{ij\}} = \exp(P_{\{ij\}}) G_{\{ij\}} / \sum_{k=1}^N \exp(P_{\{ik\}}) G_{\{ik\}}$$

Key Properties:

1. Differentiable: The pruning decision \hat{D} is incorporated through multiplication, maintaining differentiability
2. Hardware Friendly: Maintains fixed tensor shapes ($N \times N$) during training
3. Exact Pruning: When $\hat{D}_j = 0$, token j has no effect on other tokens (except itself)
4. Self-Loops: Preserve gradient flow to pruned tokens

This implementation allows end-to-end training of both the transformer and the pruning prediction module while achieving actual speedups during inference by completely removing pruned tokens.