

Algorithm 1: Attention Suppression Loss

Input: Attention matrix \mathbf{A}_{r2all} ; Index sets \mathcal{I}_p and \mathcal{I}_{np}

Parameter: Threshold factor γ ; top- k value k

Output: Suppression loss \mathcal{L}_{sup}

- 1: Obtain \mathbf{A}_{prompt} and $\mathbf{A}_{non-prompt}$ from \mathbf{A}_{r2all} using index sets \mathcal{I}_p and \mathcal{I}_{np}
- 2: $\mathbf{T}_k \leftarrow \text{TopK}(\mathbf{A}_{prompt}, k)$
- 3: $\mu_{non} \leftarrow \text{Mean}(\mathbf{A}_{non-prompt})$
- 4: **return** $\mathcal{L}_{sup} = \text{Mean}(\text{ReLU}(\mathbf{T}_k / \mu_{non} - \gamma))$

This loss can be implemented in a single line of Python:

```
loss = ReLU(TopK(A_prompt, k) / Mean(A_non_prompt) - factor).mean()
```
