$$\boldsymbol{H}^{w}_{L-F} = \text{Encoder}_{w}(\boldsymbol{H}^{w}_{0})$$
$$\boldsymbol{H}^{v}_{L-F} = \text{Encoder}_{v}(\boldsymbol{H}^{v}_{0})$$
$$\boldsymbol{H}^{w}_{L} = \text{Encoder}_{s}(\boldsymbol{H}^{w}_{L-F})$$
$$\boldsymbol{H}^{v}_{L} = \text{Encoder}_{s}(\boldsymbol{H}^{v}_{L-F}) \tag{1}$$
$$\boldsymbol{H}^{w}_{l} = \left[\boldsymbol{w}^{[\text{T\_CLS}]}_{l}, \boldsymbol{w}^{1}_{l}, ..., \boldsymbol{w}^{M}_{l}, \boldsymbol{w}^{[\text{T\_SEP}]}_{l}\right]$$
$$\boldsymbol{H}^{v}_{l} = \left[\boldsymbol{v}^{[\text{I\_CLS}]}_{l}, \boldsymbol{v}^{1}_{l}, ..., \boldsymbol{v}^{N}_{l}\right]$$

- with $l \in \{1, ..., L - F, ..., L\}$

- we define $\boldsymbol{H}^{w}_{L}$ as the final output of the student model for the caption, and $\boldsymbol{H}^{v}_{L}$ as the final output of the student model for the image, with $\boldsymbol{H}^{w}_{L} \in \mathbb{R}^{(M+2)\times D}$ and $\boldsymbol{H}^{v}_{L} \in \mathbb{R}^{(N+1)\times D}$

**Image-Text Matching with Feature Fusion**