

$$I^{(ms)} \in R^{4 \times H \times W} \quad I^{(pan)} \in R^{1 \times 4H \times 4W}$$

$$F^{(ms)} = \text{Encoder}(ms)(I^{(ms)}) \in R^{512 \times 16 \times 16}$$

$$F^{(pan)} = \text{Encoder}(pan)(I^{(pan)}) \in R^{512 \times 16 \times 16}$$

$$\text{超像素掩码} \{S_k\}_{k=1}^{65}$$

$$z_k^{(pan)} = \text{Pool}(F^{(pan)} \odot S_k) \in R^{512}$$

$$z_k^{(ms)} = \text{Pool}(F^{(ms)} \odot S_k) \in R^{512}$$

$$\text{设两个模态超像素特征为} \{z_i^{(pan)}\}_{i=1}^K, \{z_j^{(ms)}\}_{j=1}^K$$

$$\text{Sim}(z_i^{(pan)}, z_j^{(ms)}) = \frac{z_i^{(pan)} \cdot z_j^{(ms)}}{\|z_i^{(pan)}\| \cdot \|z_j^{(ms)}\|}$$

$$\bar{s}_i^{(pan \rightarrow ms)} = \frac{1}{65} \sum_j \text{Sim}(z_i^{(pan)}, z_j^{(ms)})$$

$$\bar{s}_j^{(ms \rightarrow pan)} = \frac{1}{65} \sum_i \text{Sim}(z_i^{(pan)}, z_j^{(ms)})$$

$$\mathcal{M}^{pan} = \{i \mid \bar{s}_i^{(pan \rightarrow ms)} \geq T^{(pan)}\}$$

$$\mathcal{M}^{ms} = \{j \mid \bar{s}_j^{(ms \rightarrow pan)} \geq T^{(ms)}\}$$

$$M = \{M^{pan} \cap M^{ms}\}$$

$$M = M^{ms} \odot M^{pan}$$

其中 $M^{ms}, M^{pan} \in \{0, 1\}^N$

$$S_{ms}^+, S_{ms}^-, S_{pan}^+, S_{pan}^- \in R^N$$

$$\tilde{S}_{ms}^+ = S_{ms}^+ \odot M, \tilde{S}_{ms}^- = S_{ms}^- \odot M$$

$$\tilde{S}_{pan}^+ = S_{pan}^+ \odot M, \tilde{S}_{pan}^- = S_{pan}^- \odot M$$

$$L_{ms}^{(i)} = \max(0, \text{margin} + \tilde{S}_{ms}^-(i) - \tilde{S}_{ms}^+(i))$$

$$L_{pan}^{(i)} = \max(0, \text{margin} + \tilde{S}_{pan}^-(i) - \tilde{S}_{pan}^+(i))$$

$$L_{final} = \frac{1}{2 \cdot \max(10, \sum_{i=1}^N M(i))} \left(\sum_{i=1}^N M(i) \cdot L_{ms}^{(i)} + \sum_{i=1}^N M(i) \cdot L_{pan}^{(i)} \right)$$

\Downarrow $N_v \geq 10$ \Uparrow

$$L = 0$$

$$N_v < 10$$

$$N_v = \sum_i M(i)$$

$$F_q = att_ms \cdot q_ms + att_pan \cdot q_pan$$

$$F_p = att_ms \cdot p_ms + att_pan \cdot p_pan$$