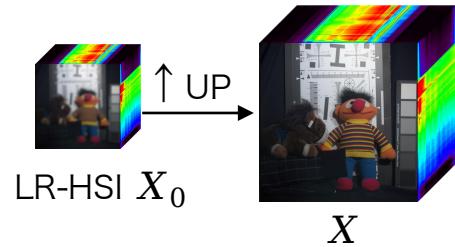
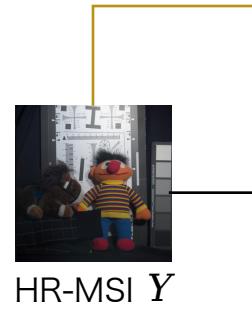
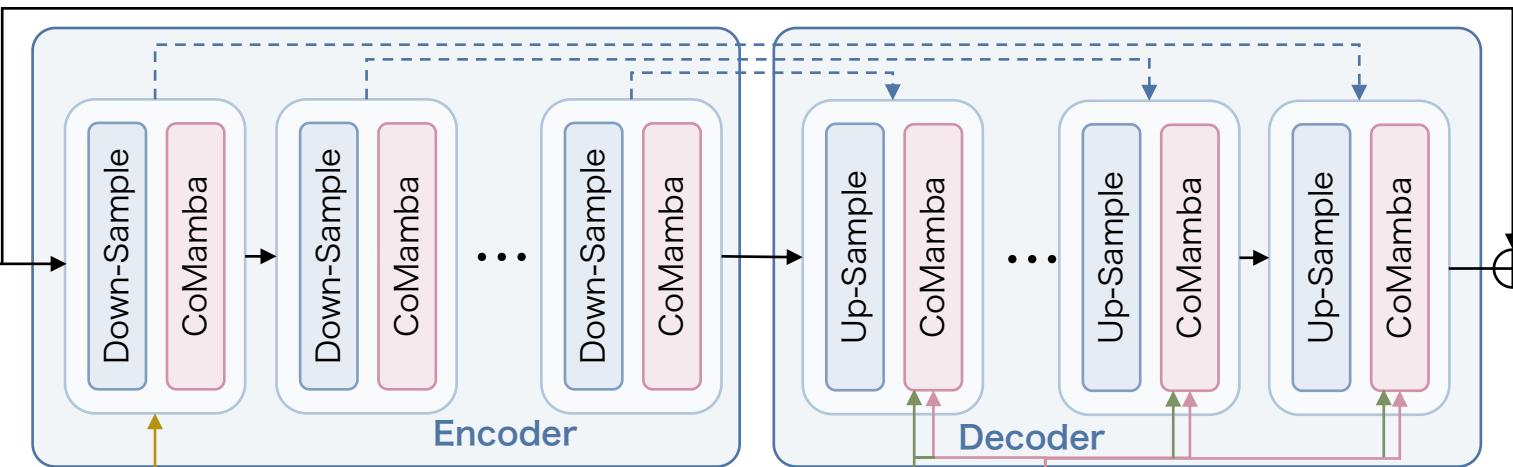


Addition
 Multiplication
 Sigmoid Function



X



A brown horse and a striped man stuffed toy are displayed with a math whiteboard and color chart.

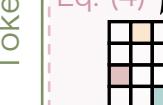
CLIP Image Encoder
Adapter

CLIP Text Encoder
Adapter

MSI Subspace
Maximum Similarity \mathcal{L}_{sim}

Text Subspace

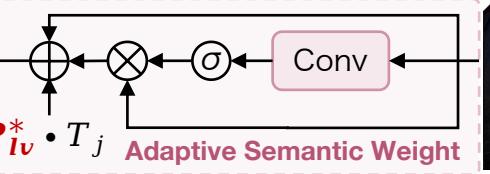
Text Token



$P_{lv}^* \cdot T_j$ Adaptive Semantic Weight

$$Eq. (4)$$

\mathcal{W}



T

$$P_{lv}^* = \arg \min_{P_{lv}} \sum_{i=1}^{N_v} \sum_{j=1}^{N_l} C_{lv}(i, j) \cdot P_{lv}(i, j) + \beta (\nabla P_{lv}(i, j))^2$$

cost matrix

$$\mathcal{I}(X, T)$$

mutual information

$$N_l$$

Cross-modal Optimal Transport (COT)

key details change

