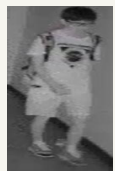


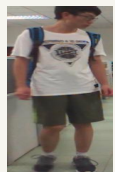
Text Generation



[filter caption]



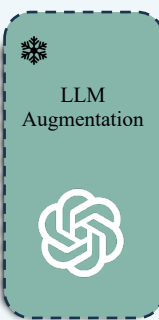
[full caption]



Incremental Fine-tuning



IR Feature



Text Feature

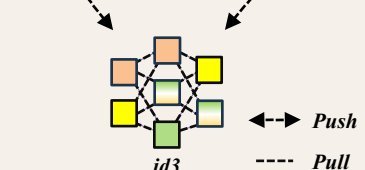
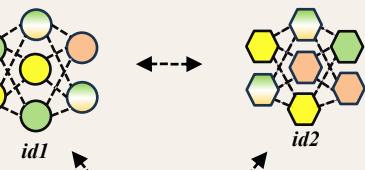
IR Enhancement
With Info Fusion

Fusion Feature

RGB Feature

Modality Joint Learning

$$L = L_{id} + L_{wrt}$$

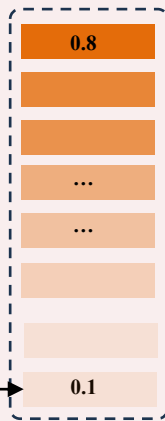


Push (solid arrow), Pull (dashed arrow)

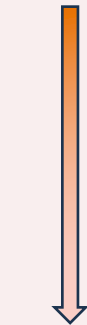


Modality Ensemble Retrieving

Ranking list



High Similarity



Low Similarity

Voting Similarity

