

Prompt-Based Augmentation



Video

Swin Transformer



Audio

Wav2vec 2.0

Aww, man, then I won't get to hear Jonah lecture us.

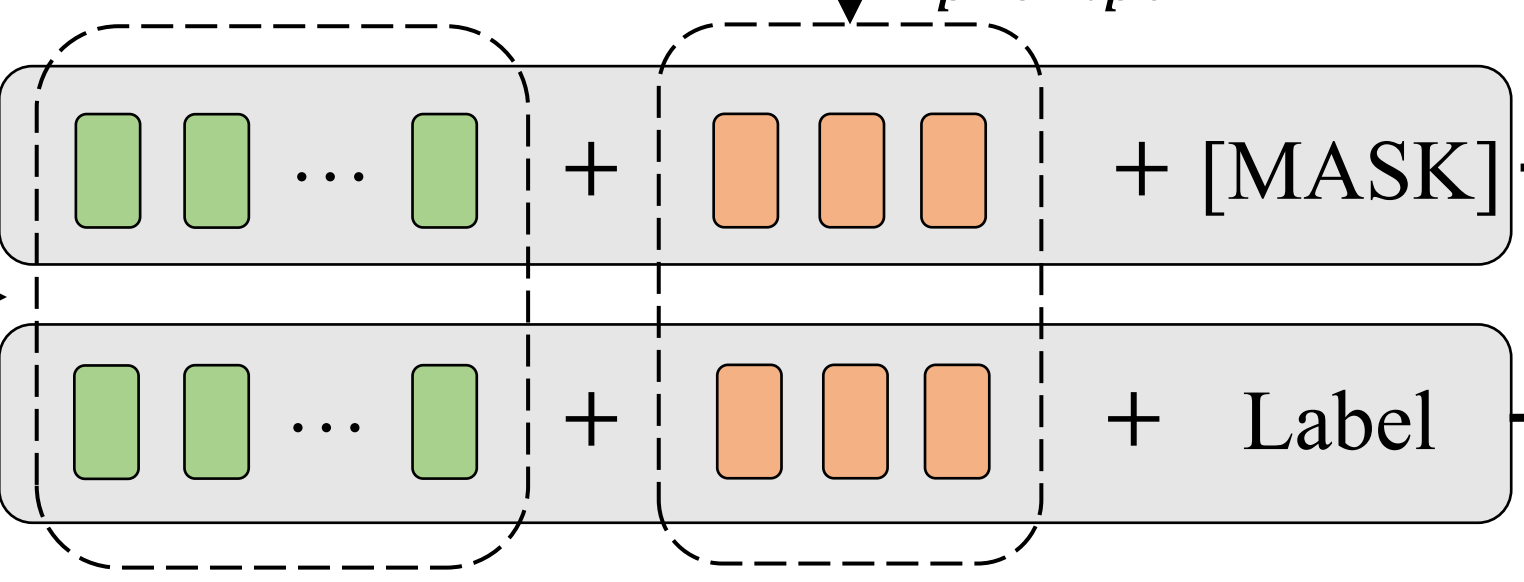
Text

BERT Embedding Layer

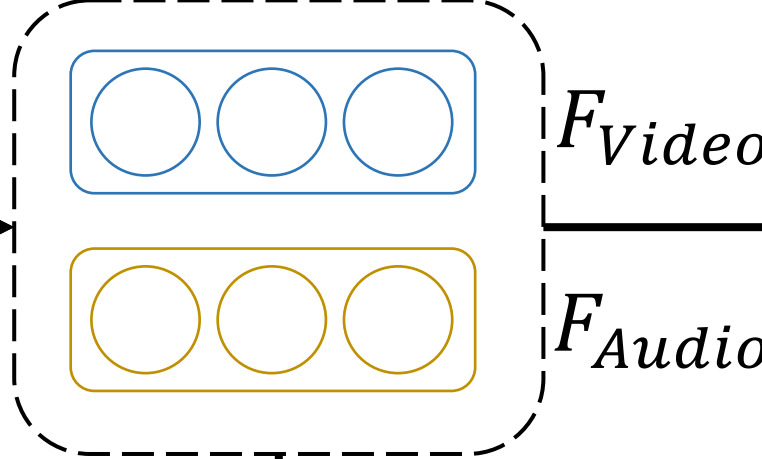
Z_{text}

Modality-Aware Prompting

Z_{prompt}



Augmented Sample Pair

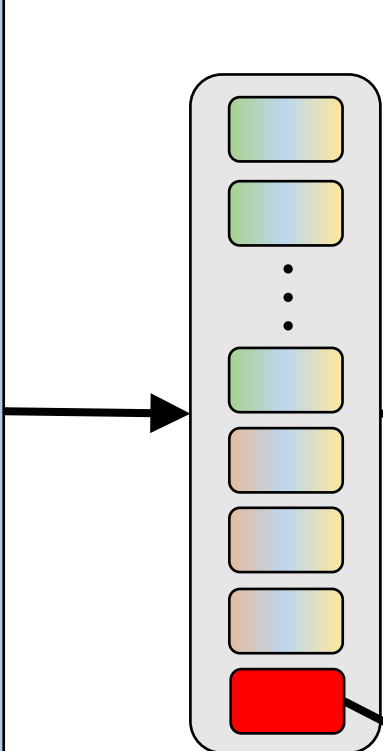


F_{video}

F_{audio}

Representation Learning

Multimodal Fusion Layer

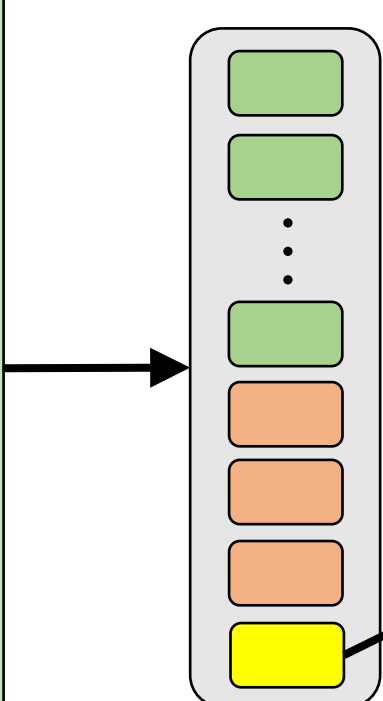


Normal Tokens

Label 1 Label 2 ... Label N

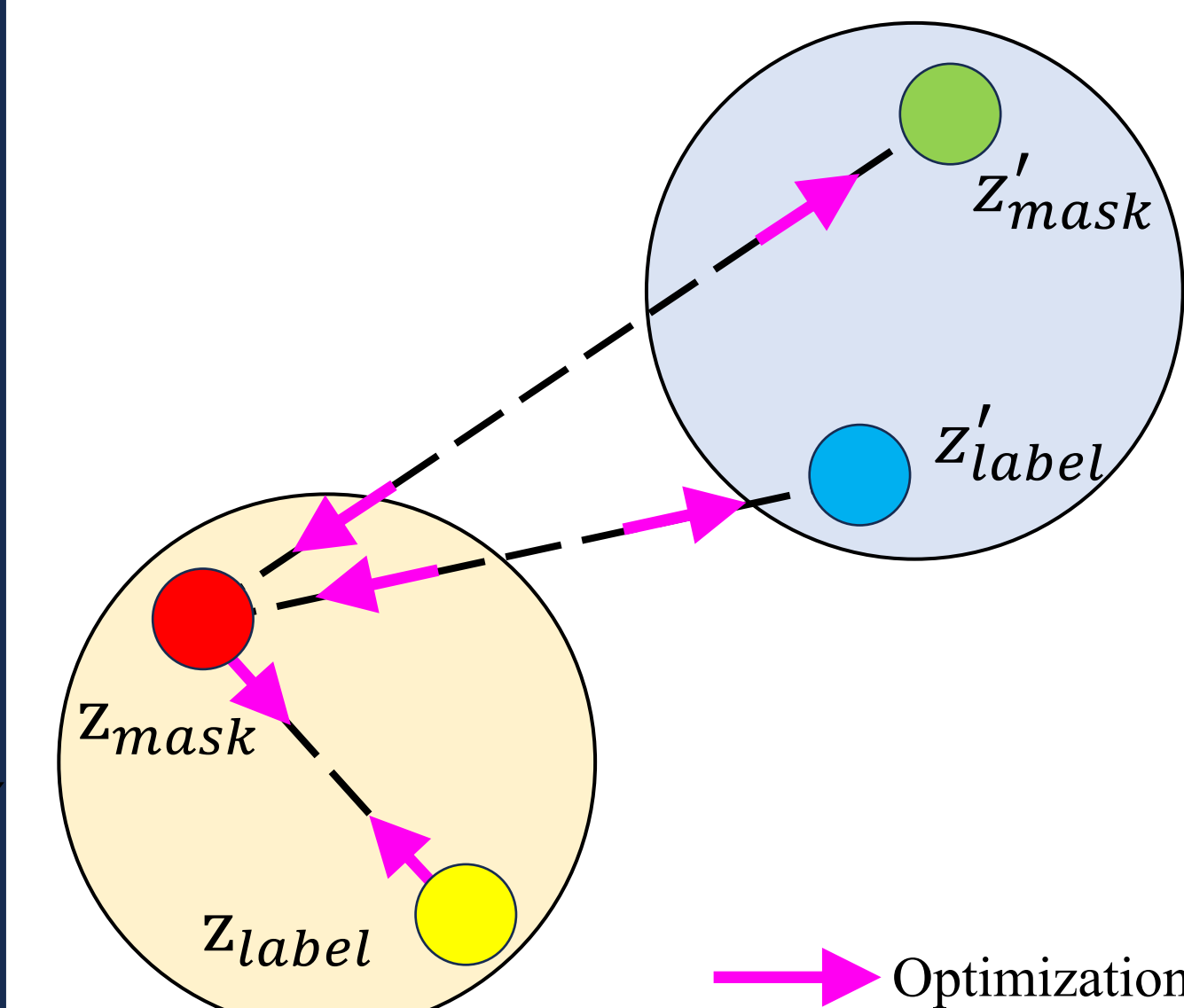
Cross Entropy Loss

BERT Encoder



Augmented Tokens

Token-Level Contrastive Learning



Optimization