

Using Mamba pretrained weight for Vision Encoder

Using ALBEF pretrained weight for Text Encoder

[CLS] Token

Mask Token

RA Loss

Relation-Aware Loss

SA Loss

Sensitivity-Aware Loss

