

UniCLIP: Unified Framework for Contrastive Language-Image Pretraining

[Abstract]

UniCLIP integrates the contrastive loss of both inter-domain pairs and intra-domain pairs into a single universal space. The discrepancies that occur when integrating contrastive loss between different domains are resolved by the three key components of UniCLIP:

- (1) augmentation-aware feature embedding
- (2) MP-NCE loss
- (3) domain dependent similarity measure

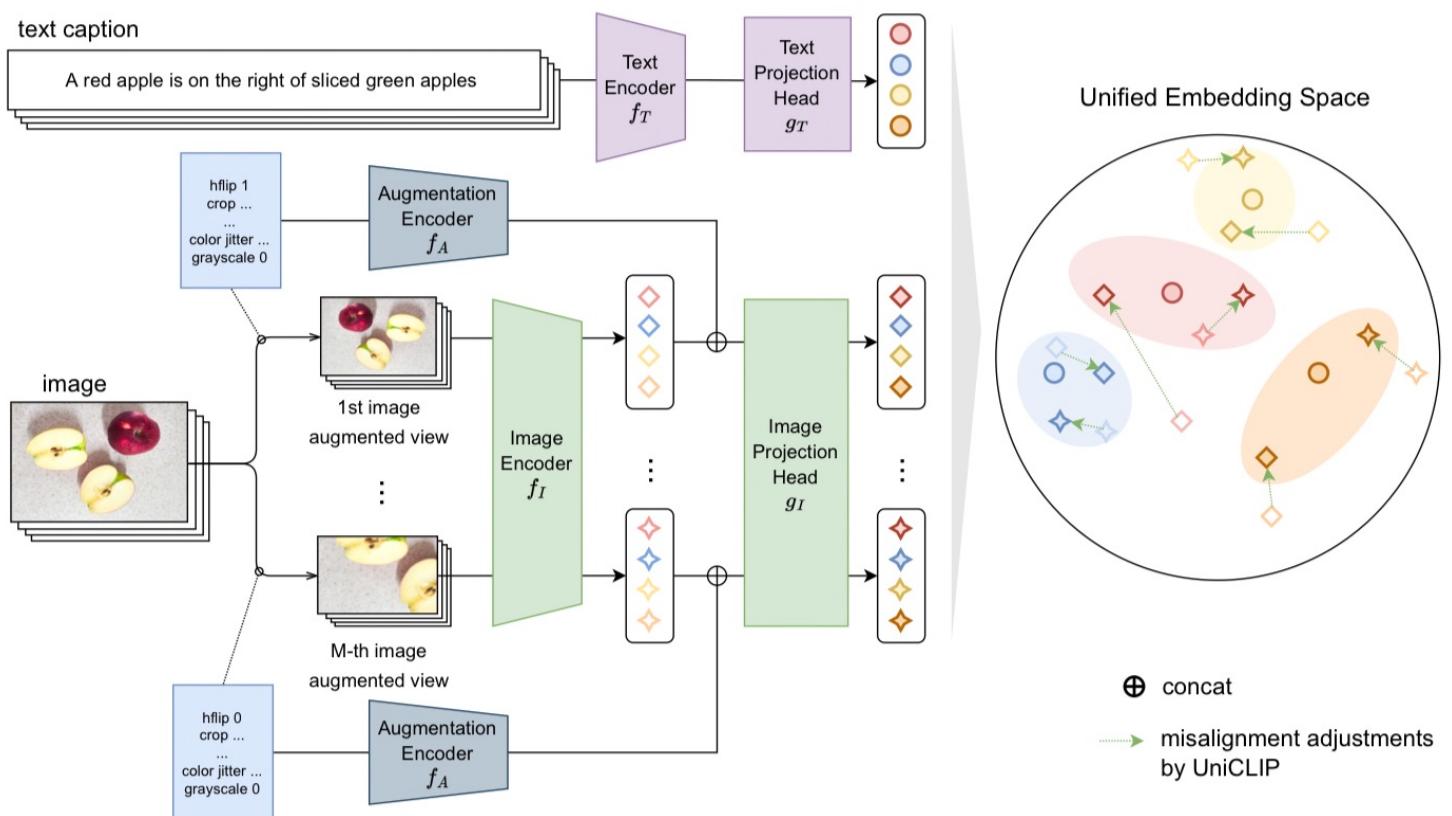
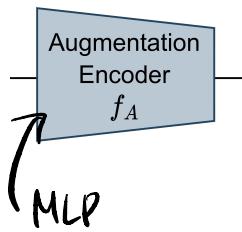


Figure 2: Overview of the UniCLIP framework.

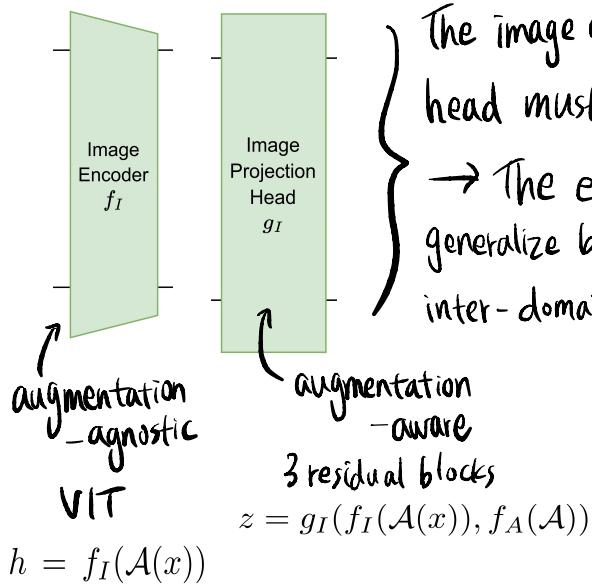
[Methods]

■ Augmentation Encoder



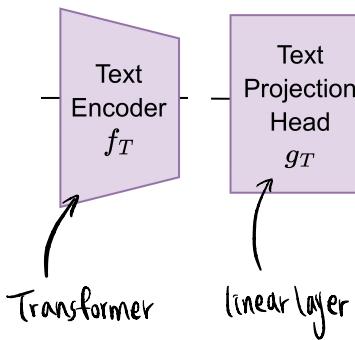
To use as an input to the network, augmentation instruction A is described as a real vector containing information about how much each basic transformation is applied to data

■ Image Encoder & Image Projection Head



The image encoder must be augmentation-agnostic and image projection head must be augmentation-aware.
 → The encoder can fully enjoy the benefits of data augmentation and generalize better, while the projection heads is still able to correct inter-domain misalignments caused by the augmentations.

■ Text Encoder & Text Projection Head



tokenized text x : tokenize raw text by byte pair encoding
 → wrap with a start token and an end token

$$h = f_T(x)$$

$$z = g_T(f_T(x))$$

Contrastive loss functions for multiple positive pairs

$$P_i = \{j | (z_i, z_j) \text{ is a positive pair and } j \neq i\}$$

$$N_i = \{j | (z_i, z_j) \text{ is a negative pair}\}$$

MIL-NCE Loss MIL-NCE loss [23] for the i -th embedding is defined by

$$\mathcal{L}_i^{\text{MIL-NCE}} = -\log \frac{\sum_{p \in P_i} s_{i,p}}{\sum_{p \in P_i} s_{i,p} + \sum_{n \in N_i} s_{i,n}}.$$

maximize the sum of all positive pair similarity scores

and minimize the sum of all negative pair similarity scores

For some $q \in P_i$,

$$\frac{\partial \mathcal{L}_i^{\text{MIL-NCE}}}{\partial s_{i,q}} = -\frac{\sum_{n \in N_i} s_{i,n}}{\left(\sum_{p \in P_i} s_{i,p}\right) \left(\sum_{p \in P_i} s_{i,p} + \sum_{n \in N_i} s_{i,n}\right)}$$

the gradient will vanish to zero when $\sum_{p \in P_i} s_{i,p}$ is already large

even if the positive score $s_{i,q}$ is small (= hard positive pair)

→ hard positive pair cannot receive enough gradient

→ easy positive pairs hinder the training of hard positive pairs

SupCon Loss SupCon loss [17] for the i -th embedding is described by

$$\mathcal{L}_i^{\text{SupCon}} = \mathbb{E}_{p \in P_i} \left[-\log \frac{s_{i,p}}{\sum_{p' \in P_i} s_{i,p'} + \sum_{n \in N_i} s_{i,n}} \right]$$

each positive pair $s_{i,p}$ is compared with the negative pairs

$$\frac{\partial \mathcal{L}_i^{\text{SupCon}}}{\partial s_{i,q}} = \frac{s_{i,q} - \frac{1}{|P_i|} \left(\sum_{p \in P_i} s_{i,p} + \sum_{n \in N_i} s_{i,n} \right)}{s_{i,q} \left(\sum_{p \in P_i} s_{i,p} + \sum_{n \in N_i} s_{i,n} \right)}$$