# Mitigate the Gap: Investigating Approaches for Improving Cross-Modal Alignment in CLIP

**Sedigheh Eslami**
Hasso Plattner Institute
sedigheh.eslami@hpi.de

**Gerard de Melo**
Hasso Plattner Institute
gerard.demelo@hpi.de

## Abstract

Contrastive Language–Image Pre-training (CLIP) has manifested remarkable improvements in zero-shot classification and cross-modal vision-language tasks. Yet, from a geometrical point of view, the CLIP embedding space has been found to have a pronounced modality gap. This gap renders the embedding space overly sparse and disconnected, with different modalities being densely distributed in distinct subregions of the hypersphere. In this work, we aim at answering three main questions: 1. Does sharing the parameter space between the multi-modal encoders reduce the modality gap? 2. Can the gap be mitigated by pushing apart the uni-modal embeddings via intra-modality separation? 3. How do these gap reduction approaches affect the downstream performance? We design AlignCLIP, in order to answer these questions and through extensive experiments, we show that AlignCLIP achieves noticeable enhancements in the cross-modal alignment of the embeddings, and thereby, reduces the modality gap, while improving the performance across several zero-shot and fine-tuning downstream evaluations. The source code for reproducing our experiments is available at https://github.com/sarahESL/AlignCLIP.

## 1 Introduction

One of the most prominent and widely used pre-trained vision–language models is OpenAI's Contrastive Language–Image Pre-training (CLIP) model (Radford et al., 2021). CLIP is a dual-stream vision–language encoder trained for learning a shared representation space, in which image and text modalities can be jointly embedded. It has demonstrated exceptional zero-shot capabilities for image classification, multi-modal retrieval as well as robustness to natural distribution shifts.

Despite the outstanding performance of CLIP, recent work has shed light on a pronounced *modality gap* in the CLIP embedding space (Liang et al., 2022; Tyshchuk et al., 2023; Schrodi et al., 2024), leading to large distances between image and text embeddings. We illustrate these phenomena more comprehensively in Figure 1A, in which the DOSNES (Lu et al., 2019) projection of the CLIP-encoded image–text pairs from CC3M (Sharma et al., 2018) using the pre-trained ViT-B-32 backend is shown. As can be seen, each modality densely populates a separate small subregion of the CLIP's representation embedding space. On the contrary, with a meaningful and well-structured cross-modal embedding space, we envision having similar samples fairly closely aligned with each other, regardless of their modalities. An example of such an embedding space is shown in Figure 1B, on the right side, with similar images and texts closely located on the hypersphere. In contrast, the left side of Figure 1B shows an example of an unaligned embedding space, inspired by the CLIP visualization.

Through extensive experiments, Liang et al. (2022) and Tyshchuk et al. (2023) showed that the modality gap is caused by a combination of the model initialization and the contrastive loss optimization. Recently, Schrodi et al. (2024) showed that the driving factor behind the modality gap is the
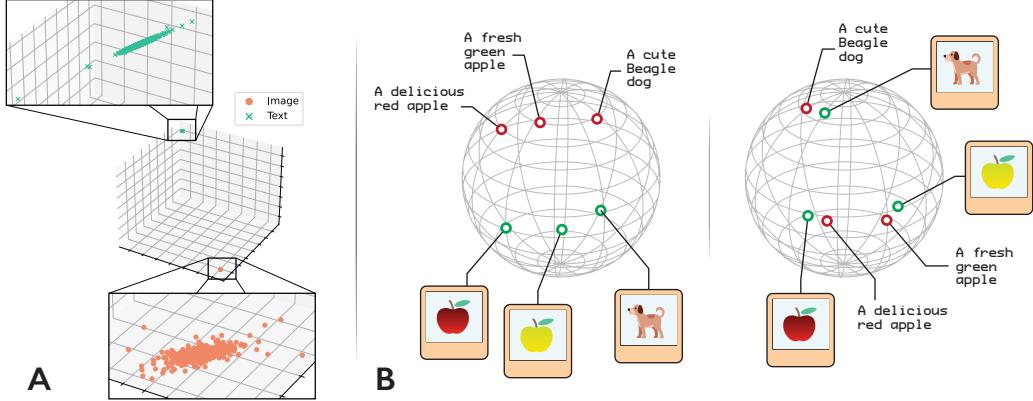
Figure 1: Modality gap and cross-modal alignment. (A) The average pairwise cosine similarity of the encoded image–text pairs is about 0.22, i.e., the average angle is about 78 degrees. (B) Schematic illustration of unaligned (left) versus aligned (right) embedding spaces.

information-imbalance between the two modalities, i.e., texts in the training datasets usually have less detailed information in comparison to their pairing images.

Attempting to reduce the modality gap, or equivalently, increase the cross-modal alignment in CLIP while enhancing the performance in downstream tasks, under the same amount of information-imbalance, we propose AlignCLIP. Ultimately, we answer three main questions: 1. With the same training dataset, i.e., under the same amount of information-imbalance between the two modalities, does sharing the parameter space between the multi-modal encoders reduce the modality gap? 2. Can the gap be further mitigated by pushing apart the uni-modal embeddings via intra-modality separation? 3. How do these gap reduction approaches affect the downstream performance?

We answer the first question by sharing the transformer encoder and the projection layer in the vision and language encoders, and observe that it already results in noticeable cross-modal alignment improvements as well as performance increase. As for the second question, we introduce the Intra-Modality Separation objective that encourages the embeddings within the visual modality to be expanded and pushed towards the language modality. As a result, each modality gets moved towards the embeddings from the opposite modality and thereby, the modality gap can get reduced. Our experimental results support that the two aforementioned refinements show substantial improvement of the cross-modal alignment while improving the zero-shot transfer performance across a wide variety of downstream tasks as well as linear probing for image classification and fine-tuning in cross-modal retrieval. As apposed to previous work that attempted at reducing the gap with naive isomorphic translations with respect to the distance of image-text pairs (Liang et al., 2022; Schrodi et al., 2024; Tyshchuk et al., 2023), which can hurt the distances of the unpaired samples, and therefore, distort the meaningful structure of the embedding space, AlignCLIP reduces the gap by modifications that improve the semantic structure of the latent modality space, as motivated by (Jiang et al., 2023).

## 2 Background, Notations and Concepts

Given a set of $N$ image–text pairs, we consider the CLIP image encoder to obtain the $l_2$-normalized vector $\vec{e}_v^i \in \mathbb{R}^d$ as a $d$-dimensional embedding vector for image $v_i$ and the CLIP text encoder to obtain the $l_2$-normalized text embedding $\vec{e}_t^i \in \mathbb{R}^d$ for the text sample $t_i$. We denote a batch of encoded image–text pairs by $E_v \in \mathbb{R}^{b \times d}$ for images and $E_t \in \mathbb{R}^{b \times d}$ for texts, where $b$ refers to the batch size.

**Background on CLIP**. The contrastive objective in CLIP is the average of the vision to language and the language to vision Info-NCE contrastive loss functions formulated as:

$$\mathcal{L}_{v \to l} = \frac{-1}{N} \sum_{i=1}^{N} \log \frac{\exp[(\vec{e}_v^i \cdot \vec{e}_t^i) \, / \, \tau]}{\sum_{j=1}^{N} \exp[(\vec{e}_v^i \cdot \vec{e}_t^j) \, / \, \tau]}, \quad \mathcal{L}_{l \to v} = \frac{-1}{N} \sum_{j=1}^{N} \log \frac{\exp[(\vec{e}_v^j \cdot \vec{e}_t^j) \, / \, \tau]}{\sum_{i=1}^{N} \exp[(\vec{e}_v^i \cdot \vec{e}_t^j) \, / \, \tau]},$$
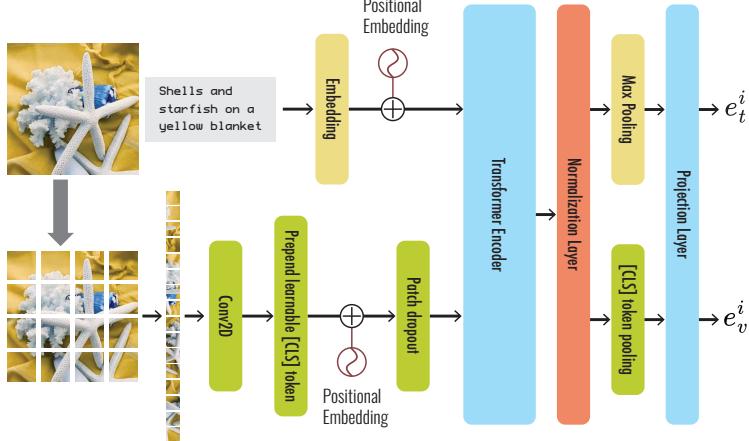(1)

Figure 2: Overview of sharing the transformer and projection layer in SharedCLIP.

respectively, where $\tau$ is the learnable temperature parameter. The overall CLIP loss is then:

$$\mathcal{L}_{\text{clip}} = \frac{1}{2} \left[ \mathcal{L}_{v \to l} + \mathcal{L}_{l \to v} \right]. \tag{2}$$

In practice, a symmetric cross-entropy loss is employed using the vision and language logits. The label $y \in \mathbb{R}$, which is the index of the image-text pair in the batch, represents the correspondence of the paired samples. For a batch of image–text pairs, $Y \in \mathbb{R}^b$ denotes the set of labels. The visual and textual logits, $\hat{y}_{\text{v}} \in \mathbb{R}^{b \times b}$ and $\hat{y}_{\text{t}} \in \mathbb{R}^{b \times b}$, are then calculated as:

$$\hat{y}_{\text{v}} = \exp(\tau) E_{\text{v}} E_{\text{t}}^{\mathsf{T}} \quad , \quad \hat{y}_{\text{t}} = \hat{y}_{\text{v}}^{\mathsf{T}}, \tag{3}$$

respectively. Ultimately, the overall CLIP loss is calculated using the cross-entropy loss, $H$, as:

$$\mathcal{L}_{\text{clip}} = \frac{1}{2} [H(\hat{y}_{\text{v}}, Y) + H(\hat{y}_{\text{t}}, Y)]. \tag{4}$$

**Cross-Modal Alignment Score.** In contrastive representation learning, the goal is to learn similar representations for positive pairs and distant representations for irrelevant negative samples. In cross-modal vision–language learning, often the paired image–texts form the positive pair distribution. Alignment entails mapping positive pairs to close embedding vectors such that a perfect alignment is achieved when $f(x_1) = f(x_2)$ for a given encoding function $f$ and a randomly drawn positive pair of samples $x_1, x_2$ (Wang and Isola, 2020). In this work, we are particularly interested in studying the alignment property for the CLIP embedding space as it represents the modality gap. We adopt the alignment measurement proposed by Goel et al. (2022) and define it as the average cosine similarity between the positive pairs in the CLIP embedding space, i.e., the paired image and text embeddings:

$$\text{alignment} = \frac{1}{N} \sum_{i=1}^{N} \vec{e}_{\text{v}}^i \cdot \vec{e}_{\text{t}}^i, \qquad \text{alignment} \in [-1, 1]. \tag{5}$$

Higher scores demonstrate better alignment, and, therefore, a decrease in the modality gap.

## 3 AlignCLIP

### 3.1 Sharing the Learnable Parameter Space in CLIP

For answering whether a shared parameter space reduces the gap, AlignCLIP employs a transformer-based encoder architecture (Vaswani et al., 2017), where the transformer is shared between the two vision and language modality encoders. We suspect that one of the main reasons that the modality gap exists in the original CLIP's embedding space is the fact that each modality has a separate disentangled set of parameters for optimization. As a result, even with the same architecture for encoding each of the modalities, the final values of the learned parameters associated to each modality

can diverge quite radically. This difference can lead to distinct functions that map each modality to completely different subregions of the embedding space.

Therefore, we seek to align the outputs of the vision and language encoding functions by sharing their parameter space to the extent possible. Specifically, we share the parameters of the transformer encoder as well as the projection layer between the vision and language modalities. We utilize a standard transformer encoder architecture (Dosovitskiy et al., 2020). An overview of our refined overall model architecture is given in Figure 2. Yellow components are designated for the language modality, green components for the visual modality, and the blue parts are shared between the two modalities. For encoding texts, the yellow and blue parts of the model actively contribute in the embedding calculations. Similarly, when encoding images, the green components as well as the blue ones are invoked for the computation. Following the original CLIP, we use max-pooling for text embeddings and $[CLS]$ token embedding for the image embedding. For simplicity, we refer to this architecture as SharedCLIP throughout the rest of the paper. Sharing parameters in CLIP has been previously investigated from the perspective of downstream performance (You et al., 2022). In this work, we rather show its effectiveness with respect to the cross-modal alignment property.

### 3.2 Intra-Modality Separation

As shown in Figure 1A, each modality resides in a distinct dense subregion of the CLIP embedding space. We hypothesize that this phenomenon is a direct result of the *cross-modal* contrastive objective in CLIP, which merely optimizes the relative distances of image embeddings and text embeddings. The cross-modal contrastive loss alone is not sufficient for imposing meaningful distances within the uni-modal embeddings, i.e., pairwise distances of text embeddings and pairwise distances of image embeddings. Therefore, we define an objective function that enforces reasonable distances within each modality by separating the uni-modal embeddings that are semantically dissimilar. In other words, we impose a semantically-regularized Intra-Modality Separation (IMSep) in addition to the CLIP's objective function.

IMSep is achieved by a vision to vision contrastive loss, where image–text pairs are the positive samples and any pairwise combination of image–image is considered as negative samples, respectively:

$$\mathcal{L}_{v \to v} = \frac{-1}{N} \sum_{i=1}^{N} \log \frac{\exp[(\vec{e}_{\mathrm{v}}^{i} \cdot \vec{e}_{\mathrm{t}}^{i}) \, / \, \tau]}{\sum_{\substack{j=1, \\ i \neq j}}^{N} \exp[(\vec{e}_{\mathrm{v}}^{i} \cdot \vec{e}_{\mathrm{v}}^{j}) \, / \, \tau]}. \tag{6}$$

In practice, given a batch of encoded image–text pairs, IMSep creates the vision-vision $\hat{y}_{\mathrm{vsep}}$ logits and then minimizes the cross-entropy loss over $Y$ and $\hat{y}_{\mathrm{vsep}}$. To this end, first the pairwise cosine similarities of the images in the batch are calculated by:

$$\mathcal{V} = E_{\mathrm{v}} E_{\mathrm{v}}^{\mathsf{T}}, \quad \text{where} \quad \mathcal{V} \in \mathbb{R}^{b \times b}. \tag{7}$$

One should notice that while enforcing intra-modality separation, by minimizing the denominator in Eq. 6, some samples might indeed be semantically similar to each other and therefore, must not be separated immensely. To this end, we regularize the intra-modality separation with respect to the pairwise semantic similarity of the samples. In order to calculate the semantic similarity within the image samples, we utilize their pairing texts as the semantic supervision signal, and invoke a pre-trained sentence encoder for encoding each text as $\vec{e}_{\mathrm{s}}^{i} \in \mathbb{R}^{d}$. We denote the corresponding batch of semantically encoded texts as $E_{\mathrm{s}} \in \mathbb{R}^{b \times d}$, and proceed to calculate the pairwise semantic similarity $S \in \mathbb{R}^{b \times b}$ and the distance $D \in \mathbb{R}^{b \times b}$ of the texts:

$$\mathcal{S} = \frac{E_{\mathrm{s}} E_{\mathrm{s}}^{\mathsf{T}}}{\|E_{\mathrm{s}}\|^{2}}, \qquad \mathcal{D} = 1 - \mathcal{S}, \tag{8}$$

respectively. We rely on $\mathcal{D}$ in order to re-scale $\mathcal{V}$. The goal of this re-scaling is to enforce a smaller dot product of the encoded images, if they are semantically similar, i.e., have a smaller distance in $D$. Conversely, the dot products in $\mathcal{V}$ obtain larger values when the text samples are semantically distant according to $D$. By this re-scaling, we seek to control the minimization of $(\vec{e}_{\mathrm{v}}^{i} \cdot \vec{e}_{\mathrm{v}}^{j})$ when the samples are semantically similar. This is because the re-scaling mechanism enforces reducing these in-modality dot product values, and therefore, the dot products do not get strongly minimized in the cross-entropy loss, since the values are already small. Re-scaling is performed by:

$$\mathcal{V}_{\mathcal{D}} = \mathcal{V} \odot \mathcal{D}, \quad \text{where} \quad \mathcal{V}_{\mathcal{D}} \in \mathbb{R}^{b \times b} \tag{9}$$
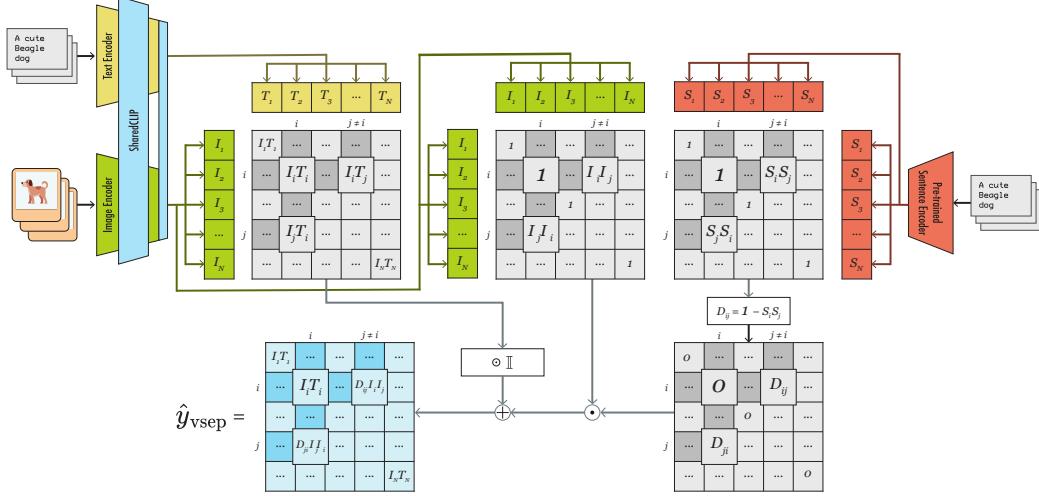
Figure 3: Schematic summary of the Intra-Modality Separation approach in AlignCLIP.

where $\odot$ denotes the element-wise product.

Afterwards, we calculate $\mathcal{M} = E_{\mathrm{v}} E_{\mathrm{t}}^{\intercal}$ and mask the non-diagonal values by $\mathrm{diag}(\mathcal{M}) = \mathbb{I} \odot \mathcal{M}$. We then obtain the vision–vision logits by:

$$\hat{y}_{\mathrm{vsep}} = \exp(\tau) \cdot [\mathrm{diag}(\mathcal{M}) + \mathcal{V}_{\mathcal{D}}], \quad \hat{y}_{\mathrm{vsep}} \in \mathbb{R}^{b \times b} \tag{10}$$

In Figure 3, our approach for obtaining IMSep is schematically summarized. Ultimately, we define the Intra-Modality Separation loss as:

$$\mathcal{L}_{\mathrm{IMsep}} = H(\hat{y}_{\mathrm{vsep}}, Y), \tag{11}$$

and adopt the core of the CLIP loss to represent the cross-modal separation:

$$\mathcal{L}_{\mathrm{CRsep}} = H(\hat{y}_{\mathrm{v}}, Y) + H(\hat{y}_{\mathrm{t}}, Y). \tag{12}$$

The final loss function optimized in AlignCLIP is:

$$\mathcal{L} = \alpha \mathcal{L}_{\mathrm{CRsep}} + \beta \mathcal{L}_{\mathrm{IMsep}} \tag{13}$$

Note that it is sufficient to define the Intra-Modality Separation function for only one of the modalities, image in our case, since the cross-modal objective defined by $\mathcal{L}_{\mathrm{CRsep}}$ already behaves as a supervision for the other modality, text in our case, and enforces the Intra-Modality Separation in that modality as well. To better support this statement, we extensively experiment the effects of adding the similar Intra-Modality Separation on text embeddings in Section A.2 and show that it is sufficient to impose the Intra-Modality Separation on the image embeddings.

## 4 Experiments

### 4.1 Training Dataset and Setup

We used the Conceptual Caption 12M (CC12M) dataset, which has been similarly used in previous work (Goel et al., 2022; Mu et al., 2022; Li et al., 2022; Changpinyo et al., 2021), for pre-training the models. We set our setup quite similar to that of the original CLIP with the ViT-B-16 backend, in order to ensure a fair comparison of CLIP with SharedCLIP and AlignCLIP. The pre-trained semantic encoder utilized in AlignCLIP for re-scaling image–image cosine similarities is the SBERT all-mpnet-base-v2 model. In order to fairly compare the effectiveness of each model, we trained all of them from scratch using the CC12M dataset and the OpenCLIP implementation (Cherti et al., 2023; Ilharco et al., 2021). Each model was trained for 6 days using an NVIDIA H100 GPU with batch size 512 for 30 epochs using AdamW optimization with a starting learning rate of $1 \times 10^{-3}$, cosine

| MODEL | CC3M | MSCOCO | IMAGENET-1K | CIFAR-100 | CIFAR-10 |
|-------|------|--------|-------------|-----------|----------|
| CLIP | 0.42 | 0.47 | 0.41 | 0.38 | 0.4 |
| SHAREDCLIP | 0.59 | 0.62 | 0.57 | 0.54 | 0.54 |
| ALIGNCLIP | **0.64** | **0.67** | **0.63** | **0.62** | **0.64** |

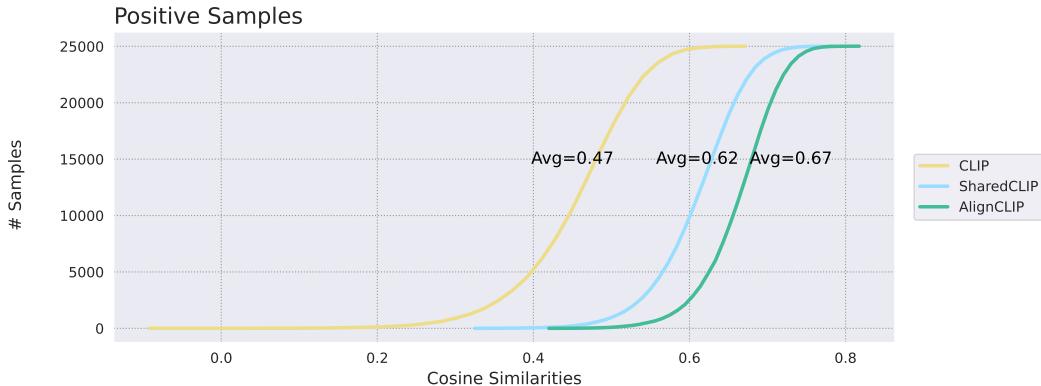Table 1: Comparison of the alignment score.



Figure 4: Cumulative distribution of pairwise cosine similarities of positive samples in MSCOCO.

scheduler, 10,000 warmup steps, a weight decay of 0.1, and an initial temperature value of 0.07. The output embedding dimension for both the vision and language modalities is set to 768. We used the checkpoint from the last epoch in our evaluations of downstream experiments. In AlignCLIP, we set $\alpha = 1$ and $\beta = \frac{1}{2}$. In Section A.1, more details about the training setup is provided.

Since our goal is to shed lights on a specific shortcoming of the original CLIP model, i.e., the modality gap problem, we compare our results to the original CLIP model in this paper. Further comparisons to other state-of-the-art models is not the focus of this study, as the goal is to investigate modifications that reduce the modality gap in CLIP without losing performance in downstream tasks, rather than achieving the state-of-the-art results.

## 4.2 Cross-Modal Alignment

We start by reporting and comparing the alignment scores when using CLIP, SharedCLIP, and AlignCLIP models on the validation sets from CC3M, MSCOCO as well as the ImageNet-1K, CIFAR-100, and CIFAR-10 test datasets. Table 1 summarizes the corresponding alignment scores. We observe that the original CLIP model has relatively low alignment scores, varying within $[0.38, 0.47]$, across all five datasets. These scores mean that the average angle between the paired image–text embeddings is a value between 61 and 68 degrees. In contrast, sharing the parameter space in SharedCLIP results in noticeable improvements of up to 0.17 in the alignment scores. As a result, the average angle between the paired image–text embeddings decreases to about 51 degrees. Furthermore, using AlignCLIP yields even better alignment scores, ranging from 0.62 to 0.67, and a decreased average angle of 47 degrees between the multi-modal paired samples. These observations confirm that AlignCLIP improves the cross-modal alignment in CLIP and thereby, reduces the modality gap.

We also plot the cumulative distribution of cosine similarities of the positive samples from MSCOCO validation dataset when encoded using the original CLIP, SharedCLIP, and AlignCLIP in Figure 4. We find that using SharedCLIP noticeably shifts the distribution of the cosine similarity of positive samples towards higher similarity values. A higher similarity of positive samples, i.e., image–text pairs, means achieving greater cross-modal alignment, and thus a lower modality gap. Furthermore, AlignCLIP shifts the distribution even more to the right, resulting in higher similarities of positive samples, better cross-modal alignment, and a reduction of the modality gap.

| MODEL | IMAGENET-1K | | CIFAR-100 | | CIFAR-10 | | FLOWERS-102 | | STANFORD CARS | |
|---|---|---|---|---|---|---|---|---|---|---|
| | TOP1 | TOP5 | TOP1 | TOP5 | TOP1 | TOP5 | TOP1 | TOP5 | TOP1 | TOP5 |
| CLIP | 31.4 | 58.7 | 28.1 | 55.9 | 61.5 | 95.6 | 18 | 39.1 | 11.6 | 36.5 |
| SHAREDCLIP | 32.1 | 59.7 | 26.4 | 54.7 | 56.9 | 95.2 | 18.2 | 38.9 | 10.7 | 35.4 |
| ALIGNCLIP | **32.8** | **60.6** | **36.5** | **66.4** | **69.3** | **97.8** | **18.8** | **40.3** | **11.8** | **38.1** |

Table 2: Accuracy scores for zero-shot image classification.

| MODEL | IMAGENET-1K | CIFAR-100 | CIFAR-10 | FLOWERS-102 | STANFORD CARS |
|---|---|---|---|---|---|
| CLIP | 50 | 62.6 | 85 | 71.5 | 42.2 |
| SHAREDCLIP | 51.2 | 63 | 85 | 74.4 | 40.5 |
| ALIGNCLIP | **51.5** | **67.4** | **87.2** | **76.8** | **45.6** |

Table 3: Accuracy scores for image classification with linear probing.

## 4.3 Classification

CLIP's pre-training objective for predicting whether a text is paired with an image has resulted in outstanding image classification performance when tested in a zero-shot setting as well as after linear probing (Radford et al., 2021). Therefore, we further assess how sharing the learnable parameters and the intra-modality loss affect the classification performance in these settings.

**Zero-Shot Image Classification.** We conduct the zero-shot classification experiments on ImageNet-1K (Russakovsky et al., 2015), CIFAR-100, CIFAR-10 (Krizhevsky et al., 2009), Flowers-102 (Nilsback and Zisserman, 2008), and Stanford Cars (Krause et al., 2013) with the combination of text prompts used by CLIP (Radford et al., 2021), e.g., "a photo of the {label}". The experimental results summarized in Table 2 show that SharedCLIP reduces the accuracy of the zero-shot classification on CIFAR-10 by about $5\%$ and $0.4\%$ when measuring Top-1 and Top-5 accuracy scores, respectively. Similarly, SharedCLIP's results on CIFAR-100 show about $2\%$ and $1\%$ decrease in accuracy. The trend of accuracy reduction is also observed on Flowers-102 and Stanford Cars datasets. However, on ImageNet-1K, SharedCLIP evinces up to $1\%$ improvement of accuracy. In contrast, AlignCLIP yields the best scores across all five datasets. It achieves $1.4\%$ and $2\%$ improvement of Top1 and Top5 accuracy, respectively, on ImageNet-1K when compared to the original CLIP. On CIFAR-10, AlignCLIP achieves up to $8\%$ and $2\%$ improvements in Top-1 and Top-5 accuracy scores. Similarly, using AlignCLIP for CIFAR-100 results in about $8\%$ and $11\%$ enhancement in Top-1 and Top-5 accuracy scores in comparison to the original CLIP. The trend of improvement is also observed on the Flowers-102 and Stanford Cars datasets. Our experiments thus evince that via sharing parameters and the additional intra-modality separation, AlignCLIP improves the cross-modal alignment of the embeddings while enhancing the performance on the downstream zero-shot image classification task.

**Linear Probing.** We further test the performance of SharedCLIP and AlignCLIP when performing linear probing for image classification and report the Top1 accuracy results in Table 3. For all datasets, we train the linear classifier layer with a batch size of 128, for 30 epochs, with AdamW optimizer, and a cosine scheduler with a starting learning rate of 5e-4. Table 3 shows the superiority of AlignCLIP in the task of image classification with linear probing across all 5 datasets.

## 4.4 Robustness to Natural Distribution Shift

In the zero-shot image classification task, CLIP has additionally shown impressive robustness to natural distribution shifts and promising generalizability to out-of-distribution images. Therefore, we expand our evaluations and investigate to what extent SharedCLIP and AlignCLIP change the performance with natural distribution shifts. We use the ImageNetV2, ImageNet-R, ImageNet-A, and ImageNetSketch datasets for these evaluations and report the corresponding results in Table 4,

| MODEL | IMAGENETV2 | | IMAGENET-R | | IMAGENET-A | | IMAGENETSKETCH | |
|---|---|---|---|---|---|---|---|---|
| | TOP1 | TOP5 | TOP1 | TOP5 | TOP1 | TOP5 | TOP1 | TOP5 |
| CLIP | 27.1 | 53.3 | 39.8 | 65.8 | 6.5 | 25.4 | 19.4 | 41.8 |
| SHAREDCLIP | 27.5 | 53.5 | 40.2 | **67.3** | 6.7 | 25.5 | 20.6 | **43.2** |
| ALIGNCLIP | **29.1** | **54.4** | **41.2** | 67.3 | **7** | **25.6** | **20.7** | **43.2** |

Table 4: Accuracy scores for zero-shot classification and natural distribution shift.

| MODEL | MSCOCO | | | | | | FLICKR30K | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | I → T | | | T → I | | | I → T | | | T → I | | |
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
| CLIP | 31.4 | 57 | 68.6 | 20.5 | 44.1 | 55.9 | 53.2 | 80.5 | 88.6 | 39.9 | 69 | 78.5 |
| SHAREDCLIP | 33.5 | 59.6 | 70.8 | 21.8 | **45.4** | **57.3** | **58.3** | **83.6** | **89.8** | 42.5 | 70 | **79.1** |
| ALIGNCLIP | **35.1** | **60.8** | **71.4** | **21.9** | **45.4** | 56.8 | 57.2 | 82.3 | **89.8** | **42.7** | **70.2** | **79.1** |

Table 5: Zero-shot cross-modal retrieval summarized with R@{1, 5, 10}.

in terms of Top-1 and Top-5 accuracy. It is first observed that SharedCLIP generally improves the classification accuracy in comparison to the CLIP model. Secondly, AlignCLIP achieves the best classification results across all datasets. In summary, when comparing to the CLIP model, AlignCLIP achieves about 2% and 1% improvement in Top1 and Top5 accuracy on the ImageNetV2 dataset. On ImageNet-R, AlignCLIP improves both Top1 and Top5 scores by about 2%. The positive general trend of the accuracy enhancement is also observed on the ImageNet-A and ImageNetSketch datasets. Based on these observations, we conclude that sharing parameters and applying intra-modality separation in SharedCLIP and AlignCLIP improves the robustness to natural distribution shifts when compared to the original CLIP model and at the same time, improve the modality gap.

### 4.5 Multi-Modal Retrieval

**Zero-Shot Transfer.** In addition to classification, we evaluate SharedCLIP and AlignCLIP in the application of zero-shot image-to-text and text-to-image retrieval using the validation splits from the MSCOCO (Lin et al., 2014) and Flickr30K (Plummer et al., 2015) datasets. The results of these evaluations are reported in Table 5. In all settings, the text prompt "a photo of the {caption}" is used. Our experiments show that both SharedCLIP and AlignCLIP improve the retrieval results measured by $R@\{1, 5, 10\}$ on both datasets when compared to the original CLIP model. In addition, AlignCLIP achieves the best overall results in comparison to SharedCLIP. When testing text-to-image retrieval on the MSCOCO dataset, SharedCLIP outperforms AlignCLIP at R@10. Similarly, for the image-to-text retrieval on the Flickr dataset, SharedCLIP achieves the best results. These experiments demonstrate that it is possible to reduce the modality gap in CLIP via parameter sharing while improving the downstream multi-modal retrieval tasks. Furthermore, the addition intra-modality separation improves the alignment noticeably while noticeably enhancing the retrieval performance.

**Fine-tuning Multi-Modal Retrieval.** Table 6 summarizes the result of multi-modal retrieval when each model is fine-tuned using the corresponding training set. We fine-tuned each model for 8 and 20 epochs on MSCOCO and Flickr, respectively. In both cases, the batch size was set to 128 and the AdamW optimizer with the learning rate of 5e-6 and a weight decay of 0.2 was used. Our results show that both SharedCLIP and AlignCLIP outperform the CLIP model in the fine-tuning scenario. Additionally, AlignCLIP generally achieves a slightly better performance in comparison to SharedCLIP.

### 4.6 Ablation Study

This section provides an ablation study on the effectiveness of the re-scaling mechanism proposed in Eq. 9. In order to see the impacts on the different types of downstream tasks, i.e., image classification,

| MODEL | MSCOCO | | | | | | FLICKR30K | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | I → T | | | T → I | | | I → T | | | T → I | | |
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
| CLIP | 39.6 | 67.5 | 78.3 | 26.7 | 53.4 | 65.7 | 64.8 | 87.8 | 93.9 | 47.4 | 75.9 | 84 |
| SHAREDCLIP | 40.7 | 69.2 | 79.6 | **27.9** | **55** | **66.7** | **66.5** | **89.1** | 94.1 | 48.9 | 76.4 | 84.3 |
| ALIGNCLIP | **41.7** | **69.3** | **80.1** | **27.9** | 54.8 | 66.2 | 66 | **89.1** | **94.5** | **49.4** | **76.7** | **84.4** |

Table 6: Fine-tuned cross-modal retrieval summarized with R@{1, 5, 10}.



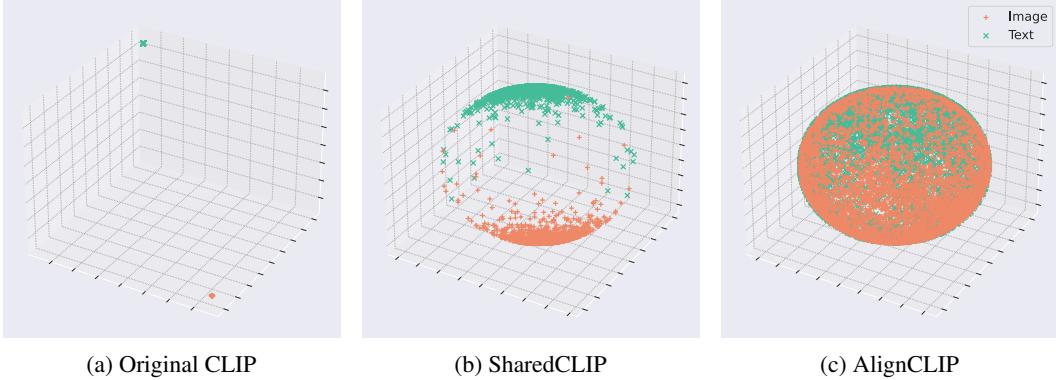(a) Original CLIP      (b) SharedCLIP      (c) AlignCLIP

Figure 5: DOSNES visualization of the multi-modal embeddings using CC3M

classification with distribution shift and multi-modal retrieval, we compare the performance of AlignCLIP with and without the re-scaling mechanism on ImageNet-1K, CIFAR-100, CIFAR-10, ImageNetV2, MSCOCO and Flickr30K datasets and summarize the results in Table 7. This study substantially concludes that the re-scaling mechanism is effective in controlling the separation of similar image samples in the batch and therefore, increasing the results in the downstream tasks.

| MODEL | IMAGENET-1K | CIFAR-100 | CIFAR-10 | IMAGENETV2 | MSCOCO | | FLICKR | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | TOP1 | TOP1 | TOP1 | TOP1 | I → T | T → I | I → T | T → I |
| ALIGNCLIP-W/O RESCALING | 32.8 | 34.7 | 64.2 | 27.9 | 34 | 59.7 | 56 | 42.5 |
| ALIGNCLIP | 32.8 | **36.5** | **69.3** | **29.1** | **35.1** | **60.8** | **57.2** | **42.7** |

Table 7: Ablation study on the re-scaling mechanism in AlignCLIP.

## 4.7 Results Analysis

**DOSNES Visualization.** For a more comprehensive comparison of the distribution of each modality, we visualize the DOSNES projection of the encoded image–texts from the MSCOCO validation set in Figure 5. As can be seen, the uni-modal embeddings in CLIP are densely located on opposite sides of the hypersphere. In contrast, the embeddings get spread out when using the SharedCLIP model. Finally, the embedding space of AlignCLIP achieves the best spread of the uni-modal embeddings and substantially greater alignment of image and text embeddings. Thus, the intra-modality separation leads to a better alignment and a substantial reduction of the modality gap. Furthermore, Figure 5c shows that AlignCLIP reduces the sparsity of the embeddings on the hypersphere.

**Comparison of Qualitative Examples.** In Table 8, examples from the MSCOCO validation dataset is provided where the images convey the same general semantics but one of the ground truth texts provides more detailed information. On the left side of Table 8, two images with their corresponding ground truth captions are provided. We use CLIP, SharedCLIP and AlignCLIP for encoding the images and texts, and provide the cosine similarities on the right side of the table. As can be seen,

| | | CLIP | | SHAREDCLIP | | ALIGNCLIP | |
|---|---|---|---|---|---|---|---|
| | | $T_1$ | $T_2$ | $T_1$ | $T_2$ | $T_1$ | $T_2$ |
| **1** [An image] **2** [An image] | $I_1$ | 0.37 | 0.42 | 0.53 | 0.55 | 0.61 | 0.58 |
| | $I_2$ | 0.39 | 0.52 | 0.61 | 0.64 | 0.63 | 0.71 |

**1** A bowl that has food inside of it

**2** An orange bowl filled with lots of noodles and beef

Table 8: Qualitative example of semantically similar samples from MSCOCO.

when querying the first image, the similarity of the first image and the second text using the CLIP embeddings is higher in comparison to the ground truth caption. Suggesting that the second text will get selected as the predicted caption when using CLIP. This flaw still appears when using SharedCLIP for encoding the images and texts. However, when using AlignCLIP, the cosine similarity of the first image and the first text is higher in comparison to the second text, meaning that when querying the first image, the ground truth caption, which is semantically more correct in comparison to the second text, successfully gets selected. This suggests that the semantic regularization of the Intra-Modality Separation in AlignCLIP, which is calculated using the semantics of the text samples, potentially contributes in improving the retrieval performance.

## 5   Related Work

**Modifications of CLIP.** Prior work has sought to improve several different aspects of CLIP. UniCLIP (Lee et al., 2022) proposed a multi-positive contrastive loss for training augmentation-aware feature embeddings. SLIP (Mu et al., 2022) added an additional image-based contrastive objective function using augmented images as positive samples to the original CLIP's loss function. Furthermore, CyCLIP (Goel et al., 2022) formalized the geometric consistency of the image and text embeddings in the CLIP embedding space and added a loss function for regularizing the cross-modal and in-modal similarity scores. The authors provide an analysis of the cross-modal alignment showing that the final CyCLIP model does not improve the alignment of image–text embeddings. In xCLIP (Zhou et al., 2023), a non-contrastive training regimen based on image–text pairs is adopted in order to mitigate the requirement of large batch sizes in contrastive learning. With a similar goal, SigLIP (Zhai et al., 2023) employed a pairwise sigmoid loss for scaling up the batch size. The recently published ReCLIP (Hu et al., 2024) had the goal of refining the fine-tuning of CLIP in domain adaptation settings, based on unlabeled target samples. While ReCLIP results in noticeable gains for the considered classification task, its effectiveness with regard to the cross-modal alignment property is not studied by the authors. In contrast to previous work, AlignCLIP studies the cross-modal alignment and modality gap phenomena in CLIP by investigating the effects of parameter-sharing as well as intra-modality separation.

**Modality Gap in CLIP.** The modality gap in CLIP's embedding space was first studied by Liang et al. (2022). Authors showed that the modality gap is caused by a combination of the model initialization and the contrastive loss optimization. Furthermore, they showed that CLIP's embedding space is very sparse such that the effective embedding space is an extremely narrow cone. In a similar study, Tyshchuk et al. (2023) measured the alignment of image and text embeddings in CLIP with a focus on the isotropic properties. More recently, Schrodi et al. (2024) showed that one of the key factors contributing in both the modality gap and the bias in CLIP is the information-imbalance between the two modalities, i.e., images often have much more detailed information of the scene in comparison to their corresponding text descriptions. Furthermore, Jiang et al. (2023) showed that even under the perfect alignment, the prediction error in downstream tasks cannot be smaller than the information gap that exists between the modalities. They further propose intra-modality as well as inter-modality regularization using augmented samples in order to improve the latent embeddings structures. In contrast to the previous work, we study the effectiveness of parameter sharing between the modality

encoders as well as semantically-regularized separation of the uni-modal embeddings on reducing the modality gap, under the same amount of information gap across modalities.

# 6 Conclusion

This work investigates the potential of reducing the modality gap in CLIP by sharing the learnable parameter space of the vision and language encoders as well as enforcing a contrastive intra-modality separation objective. We further examine the effects of modality gap reduction by the aforementioned refinements in the performance of downstream tasks. Through extensive experiments, we show that SharedCLIP and AlignCLIP improve various zero-shot as well as fine-tuning downstream applications when compared to CLIP, while substantially improving the cross-modal alignment, and therefore, reducing the gap. Our work shows that it is possible to mitigate the modality gap in CLIP via parameter sharing and intra-modality separation without losing downstream performance.

# Acknowledgments

# References

Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3558–3568, June 2021.

Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2829, June 2023.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.

Shashank Goel, Hritik Bansal, Sumit Bhatia, Ryan Rossi, Vishwa Vinay, and Aditya Grover. Cyclip: Cyclic contrastive language-image pretraining. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 6704–6719. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/2cd36d327f33d47b372d4711edd08de0-Paper-Conference.pdf.

Xuefeng Hu, Ke Zhang, Lu Xia, Albert Chen, Jiajia Luo, Yuyin Sun, Ken Wang, Nan Qiao, Xiao Zeng, Min Sun, Cheng-Hao Kuo, and Ram Nevatia. Reclip: Refine contrastive language image pre-training with source free domain adaptation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2994–3003, January 2024.

Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL https://doi.org/10.5281/zenodo.5143773. If you use this software, please cite it as below.

Qian Jiang, Changyou Chen, Han Zhao, Liqun Chen, Qing Ping, Son Dinh Tran, Yi Xu, Belinda Zeng, and Trishul Chilimbi. Understanding and constructing latent modality structures in multi-modal representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7661–7671, June 2023.

Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV) Workshops*, June 2013.

Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

Janghyeon Lee, Jongsuk Kim, Hyounguk Shon, Bumsoo Kim, Seung Hwan Kim, Honglak Lee, and Junmo Kim. Uniclip: Unified framework for contrastive language-image pre-training. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 1008–1019. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/072fd0525592b43da661e254bbaadc27-Paper-Conference.pdf.

Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=zq1iJkNk3uN.

Victor Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Y Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 17612–17625. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/702f4db7543a7432431df588d57bc7c9-Paper-Conference.pdf.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.

Yao Lu, Jukka Corander, and Zhirong Yang. Doubly stochastic neighbor embedding on spheres. *Pattern Recognition Letters*, 128:100–106, 2019. ISSN 0167-8655. doi: https://doi.org/10.1016/j.patrec.2019.08.026. URL https://www.sciencedirect.com/science/article/pii/S0167865518305099.

Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pre-training. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision – ECCV 2022*, pages 529–544, Cham, 2022. Springer Nature Switzerland. ISBN 978-3-031-19809-0.

Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages 722–729, 2008. doi: 10.1109/ICVGIP.2008.47.

Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/radford21a.html.

Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115 (3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.

Simon Schrodi, David T Hoffmann, Max Argus, Volker Fischer, and Thomas Brox. Two effects, one trigger: On the modality gap, object bias, and information imbalance in contrastive vision-language representation learning. In *ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2024. URL https://openreview.net/forum?id=7QwFMLzQHH.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In Iryna Gurevych and Yusuke Miyao, editors, *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1238. URL https://aclanthology.org/P18-1238.

Kirill Tyshchuk, Polina Karpikova, Andrew Spiridonov, Anastasiia Prutianova, Anton Razzhigaev, and Alexander Panchenko. On isotropy of multimodal embeddings. *Information*, 14(7), 2023. ISSN 2078-2489. doi: 10.3390/info14070392. URL https://www.mdpi.com/2078-2489/14/7/392.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.

Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 9929–9939. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/wang20k.html.

Haoxuan You, Luowei Zhou, Bin Xiao, Noel Codella, Yu Cheng, Ruochen Xu, Shih-Fu Chang, and Lu Yuan. Learning visual representation from modality-shared contrastive language-image pre-training. In *European Conference on Computer Vision*, pages 69–87, 2022.

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11975–11986, October 2023.

Jinghao Zhou, Li Dong, Zhe Gan, Lijuan Wang, and Furu Wei. Non-contrastive learning meets language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11028–11038, June 2023.

| MODEL | IMAGENET-1K | | CIFAR-100 | | CIFAR-10 | | FLOWERS-102 | | STANFORD CARS | |
|---|---|---|---|---|---|---|---|---|---|---|
| | TOP1 | TOP5 | TOP1 | TOP5 | TOP1 | TOP5 | TOP1 | TOP5 | TOP1 | TOP5 |
| ALIGNCLIP (TT) | 31.1 | 58.6 | 31.5 | 60.9 | 64.8 | 95.5 | 18.7 | 39.5 | 10 | 36.4 |
| ALIGNCLIP (II) | **32.8** | **60.6** | **36.5** | **66.4** | **69.3** | **97.8** | **18.8** | **40.3** | **11.8** | **38.1** |
| ALIGNCLIP (II-TT) | 32.4 | 60 | 31.2 | 61.8 | 66.4 | 96.9 | 18.7 | 39.9 | 11.8 | 37.1 |

Table 9: Accuracy scores for zero-shot image classification.

| MODEL | IMAGENET-1K | CIFAR-100 | CIFAR-10 | FLOWERS-102 | STANFORD CARS |
|---|---|---|---|---|---|
| ALIGNCLIP (TT) | 51.2 | 64.5 | 86.3 | 74.3 | 41.6 |
| ALIGNCLIP (II) | **51.5** | **67.4** | **87.2** | **76.8** | **45.6** |
| ALIGNCLIP (II-TT) | 50.3 | 63.4 | 85.8 | 72.2 | 42.2 |

Table 10: Accuracy scores for image classification with linear probing.

# A   Appendix / Supplemental Material

## A.1   Training Dataset and Setup

For pre-training, we used the Conceptual Caption 12M (CC12M) dataset (Changpinyo et al., 2021). In comparison to the data used for pre-training CLIP, i.e., 400M image–text pairs, CC12M is a much smaller dataset, with about 12M noisy image–text pairs. Therefore, our results are not directly comparable to the numbers reported in the original CLIP paper (Radford et al., 2021). Nonetheless, CC12M is one of the popularly used large-scale and publicly available datasets that enables pre-training vision–language models and analyzing the effectiveness of different training paradigms on benchmark evaluations.

We adopted a transformer encoder consisting of 12 layers and 12 heads in SharedCLIP and AlignCLIP. The same input pre-processing and augmentations employed in CLIP have been used for SharedCLIP as well as AlignCLIP, including random cropping of images to the size $224 \times 224$. The image patch size for encoding visual data is set to $16 \times 16$. When encoding texts, the maximum sequence length is set to 77 tokens and the vocabulary size for the embedding layer is set to 49,408. The output embedding dimension for both the vision and language modalities is set to 768. We chose our setup to be quite similar to that of the original CLIP with the ViT-B-16 backend, in order to ensure a fair comparison of CLIP and our proposed model. SBERT all-mpnet-base-v2 model has been utilized in AlignCLIP for re-scaling text–text and image–image cosine similarities is the .

In order to fairly compare the effectiveness of each model, we trained all of the models , i.e., the original CLIP with ViT-B-16 backend, SharedCLIP, and AlignCLIP, from scratch using the CC12M dataset and the OpenCLIP implementation (Cherti et al., 2023; Ilharco et al., 2021). For all models, we used AdamW optimization with a starting learning rate of $1 \times 10^{-3}$, cosine scheduler, 10,000 warmup steps, and a weight decay of $0.1$. The initial temperature value for all models were set to $0.07$. Each model was trained using an NVIDIA H100 GPU with batch size 512 for 30 epochs. We used the checkpoint from the last epoch in our evaluations of downstream experiments. In AlignCLIP, we set $\alpha = 1$ and $\beta = \frac{1}{2}$.

## A.2   Effects of text embedding separation

In Section 3.2, we proposed IMSep to enforce intra-modality separation among the image embeddings. A potential question that arises is that what happens if we enforce the separation amongst the text embeddings, or even, what happens of we enforce the separation inside the image embeddings as well as the text embeddings. To answer these questions, we trained versions of AlignCLIP without image-image separation and with only text-text separation, namely AlignCLIP-TT. In this case, we also used the re-scaling mechanism described in Section 3.2. Furthermore, we trained a version including both

image-image embeddings separation and text-text embeddings separation, i.e., AlignCLIP-II-TT. We followed the same experimental setup described in Section A.1 and compare the results for zero-shot image classification as well as linear probing in Table 9 and Table 10, respectively.

## B  Limitations

For the zero-shot classification experiments, we have tested the models using the ImageNet-1k, CIFAR-10, CIFAR-100, Stanford Cars, and Flowers-102 datasets. More classification analysis using other benchmarks datasets such as Food-100 and Hateful Memes has not been studied in this work. Furthermore, this study is limited to the English language and further analyses on multilingual and cross-lingual representation learning are necessary. Moreover, the re-scaling mechanism in the intra-modality separation loss is dependant on the choice of the pre-trained sentence encoder and has not been fully benchmarked in this work.

## C  Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning at a fundamental level with regard to cross-modal representation learning. We acknowledge that this line of work has a broad range of potential societal implications. For instance, vision–language models may exhibit harmful biases and stereotypes, particularly when trained on data crawled from the Web. Due to their incorporation in prominent generative AI models, vision–language models may also contribute toward the model's ability to produce images portraying trademarked characters or notable figures. These sorts of concerns need to be carefully considered before incorporating such models into real-world applications.