
CYCLIP: Cyclic Contrastive Language-Image Pretraining

Shashank Goel*
UCLA
shashankgoel@ucla.edu

Hritik Bansal*
UCLA
hbansal@ucla.edu

Sumit Bhatia
MDSR Lab, Adobe Systems
sumit.bhatia@adobe.com

Ryan A. Rossi
Adobe Research
ryrossi@adobe.com

Vishwa Vinay
Adobe Research
vinay@adobe.com

Aditya Grover
UCLA
adityag@cs.ucla.edu

Abstract

Recent advances in contrastive representation learning over paired image-text data have led to models such as CLIP [46] that achieve state-of-the-art performance for zero-shot classification and distributional robustness. Such models typically require joint reasoning in the image and text representation spaces for downstream inference tasks. Contrary to prior beliefs, we demonstrate that the image and text representations learned via a standard contrastive objective are not interchangeable and can lead to inconsistent downstream predictions. To mitigate this issue, we formalize *consistency* and propose CYCLIP, a framework for contrastive representation learning that explicitly optimizes for the learned representations to be *geometrically consistent* in the image and text space. In particular, we show that consistent representations can be learned by explicitly symmetrizing (a) the similarity between the two mismatched image-text pairs (cross-modal consistency); and (b) the similarity between the image-image pair and the text-text pair (in-modal consistency). Empirically, we show that the improved consistency in CYCLIP translates to significant gains over CLIP, with gains ranging from 10% – 24% for zero-shot classification accuracy on standard benchmarks (CIFAR-10, CIFAR-100, ImageNet1K) and 10% – 27% for robustness to various natural distribution shifts. The code is available at <https://github.com/goel-shashank/CyCLIP>.

1 Introduction

The ability to learn general-purpose representations from diverse data modalities is a long-standing goal of artificial intelligence (AI) [4, 32]. In this regard, recent instantiations such as CLIP [46], ALIGN [29], and BASIC [42] have scaled up vision-language contrastive pretraining to jointly learn image and text embeddings, by exploiting an enormous amount of paired image-text data on the web. Post pretraining, these embeddings exhibit impressive zero-shot classification performance [13] and robustness to natural distribution shifts [49, 58, 24, 26]. Recently, these embeddings have been extended to text-guided generation of natural images [47, 12, 38, 48] and transferred to modalities such as 3-D shapes [51] by emphasizing the interchangeability of the image and text embeddings.

In the context of vision-language pretraining, the standard contrastive learning objective aims to maximize the similarity between matched image-text pairs (“positives”) against all the mismatched image-text pairs (“negatives”) [45, 7, 41, 22]. While such an objective aligns the true image-text pairs, it poses no constraints on the overall geometry of all data pairs, including the mismatched

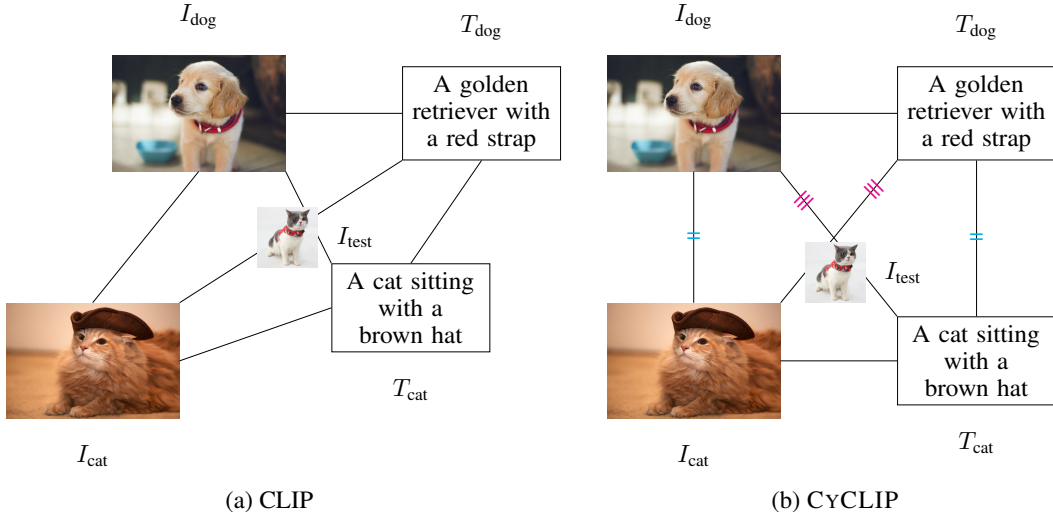


Figure 1: An illustration of the planar geometry of the learned representations of image-text pairs by (a) CLIP and (b) CYCLIP. The edges indicate the distance between the representations i.e., $d(e_1, e_2) = 1 - \langle e_1, e_2 \rangle$, where $\langle \cdot, \cdot \rangle$ is the inner product. CYCLIP is cyclic consistent between image-text pairs as the in-modal distances, $d(T_{\text{cat}}, T_{\text{dog}}) \sim d(I_{\text{cat}}, I_{\text{dog}})$, and the cross-modal distances, $d(T_{\text{cat}}, I_{\text{dog}}) \sim d(I_{\text{cat}}, T_{\text{dog}})$, are similar to each other unlike CLIP. Due to explicit consistency constraints, the test image of a cat is classified as a cat in the image as well as the text space.

pairs and pairs within the same modality. In Figure 1 (a), we illustrate this effect where matched image-text pairs, $(I_{\text{dog}}, T_{\text{dog}})$ and $(I_{\text{cat}}, T_{\text{cat}})$, get close to each other but the overall geometry of pairwise distances can be highly irregular (see e.g., $(I_{\text{dog}}, T_{\text{cat}})$ and $(I_{\text{cat}}, T_{\text{dog}})$). If we use such representations for downstream inference, such irregularities can translate into inconsistent reasoning in the image and text spaces. For example, CLIP designs proxy captions for class labels and uses the most similar class caption to perform zero-shot classification for images; using the default captions in Figure 1 (a), this would imply that a test image I_{test} gets classified as a dog in the image space even when a simple nearest neighbor classifier in the text space would correctly infer the label to be a cat.

To mitigate these challenges, we propose **Cyclic Contrastive Language-Image Pretraining (CYCLIP)**, a framework that imposes additional geometric structure on the learned representations. Specifically, given two image-text pairs, we augment the contrastive learning objective with two symmetrization terms. The first term provides for in-modal consistency by encouraging the distance between the two image embeddings to be close to the distance between the corresponding text embeddings. The second term for the cross-modal consistency that encourages the distance between the image and text embedding from the first and second pairs respectively to be close to the distance between the text and image embeddings from the first and second pairs respectively. As shown in Figure 1 (b), if representations of any two image-text pairs, $(I_{\text{dog}}, T_{\text{dog}})$ and $(I_{\text{cat}}, T_{\text{cat}})$ exactly satisfy both forms of cyclic consistency, then we can guarantee that any test image I_{test} respects the ordering of distances in both image and text spaces (i.e., if $d(I_{\text{test}}, I_{\text{dog}}) > d(I_{\text{test}}, I_{\text{cat}})$, then $d(I_{\text{test}}, T_{\text{dog}}) > d(I_{\text{test}}, T_{\text{cat}})$).

Empirically, we demonstrate that the improved consistency in CYCLIP translates to improvements over CLIP. In all cases, we pre-train our models on the Conceptual Captions 3M dataset[53]. On zero-shot classification, we observe that CYCLIP improves over CLIP by 10.2% on ImageNet1K, 10.6% on CIFAR-10 and 23.9% on CIFAR-100 respectively. Further, CYCLIP outperforms CLIP with an average relative gain of +17% on ImageNet natural distribution shift benchmarks. We further analyze the improved performance of CYCLIP and find that the additional geometric structure in the representation space better captures the coarse and fine-grained concept hierarchies of datasets.

Our contributions are as follows:

1. We analyze contrastive learning for representation learning jointly over image and text modalities. We identify a critical shortcoming in the geometry of the learned representation space that can lead to inconsistent predictions in image and text domains.

2. We propose CYCLIP, a simple and effective framework for contrastive representation learning with two additional cycle consistency constraints for mitigating the above issue.
3. We demonstrate that CYCLIP achieves significant empirical improvements over CLIP on zero-shot classification and robustness benchmarks. We further explain these improvements by analyzing the impact of consistency on the hierarchical structure of datasets.

2 Cycle Consistent Representation Learning

2.1 Preliminaries

We are interested in using text supervision to learn general-purpose visual representations that can be generalized to downstream predictive tasks. To this end, there have been several recent advances in language-image pretraining concerning model architectures, training objectives, and sources of supervision. Our work is most closely related to Contrastive Language-Image Pretraining (CLIP) [46] which combines many such advances in a highly scalable and generalizable learning framework.

CLIP is trained on millions of images with their captions scraped from the web. Formally, we consider a dataset $S \subset \mathcal{I} \times \mathcal{T}$ consisting of pairs (I_j, T_j) where I_j is a raw image and T_j is a text caption. We use \mathcal{I} and \mathcal{T} to denote the domain of images and text, respectively. The CLIP architecture consists of 3 components: (i) an image encoder network, $f_I : \mathcal{I} \mapsto \mathbb{R}^d$, to encode the raw image into an embedding vector of dimension d , (ii) a text encoder network, $f_T : \mathcal{T} \mapsto \mathbb{R}^d$, to encode the raw text into an embedding vector of dimension d , (iii) a contrastive objective that pulls the embeddings of paired image-caption pairs together while pushing apart embeddings of unmatched pairs.

Formally, during training, consider a batch of N image-captions pairs, $\{I_j, T_j\}_{j=1}^N$, where I_j and T_j represent the raw image and text pair, respectively. The image embedding $I_j^e \in \mathbb{R}^d$ and text embedding $T_j^e \in \mathbb{R}^d$ are obtained by passing I_j and T_j through the image encoder f_I and text encoder f_T , respectively; i.e. $I_j^e = f_I(I_j)$ and $T_j^e = f_T(T_j)$. Further, we assume they are normalized to have unit ℓ_2 -norm. The contrastive objective in CLIP aims to align the image and text representations by minimizing the loss function $\mathcal{L}_{\text{CLIP}}$ shown below:

$$\mathcal{L}_{\text{CLIP}} = -\frac{1}{2N} \sum_{j=1}^N \log \underbrace{\left[\frac{\exp(\langle I_j^e, T_j^e \rangle / \tau)}{\sum_{k=1}^N \exp(\langle I_j^e, T_k^e \rangle / \tau)} \right]}_{\text{Contrasting images with the texts}} - \frac{1}{2N} \sum_{k=1}^N \log \underbrace{\left[\frac{\exp(\langle I_k^e, T_k^e \rangle / \tau)}{\sum_{j=1}^N \exp(\langle I_j^e, T_k^e \rangle / \tau)} \right]}_{\text{Contrasting texts with the images}} \quad (1)$$

where $\langle \cdot, \cdot \rangle$ represents the inner product, and τ is a trainable temperature parameter. CLIP and its variants can be used to perform zero-shot image classification, i.e., classifying test images into categories not seen at training time. We first transform each category into a suitable caption (e.g., the airplane category in CIFAR-10 can be expressed as ‘a photo of an airplane’). Then, the similarity of the test image to each caption is computed (e.g., cosine distance), and the model predicts the category for which the image-caption similarity is the highest.

2.2 Inconsistent Representation Learning in CLIP

As illustrated in Figure 1 (a), the standard contrastive objective in CLIP can learn image-text representations such that the predicted labels for the test image are different in the image and text spaces. Here, we reason about such inconsistencies more formally in the context of downstream classification. As discussed above, we can predict a label in the text embedding space (zero-shot setting) by selecting the label that is closest to the test image (P_T). Additionally, for classification in the image embedding space, if we had access to a labeled training set, then one natural way to infer the predicted label (P_I^k) of a test image I_{test} is by taking a majority vote from the true labels

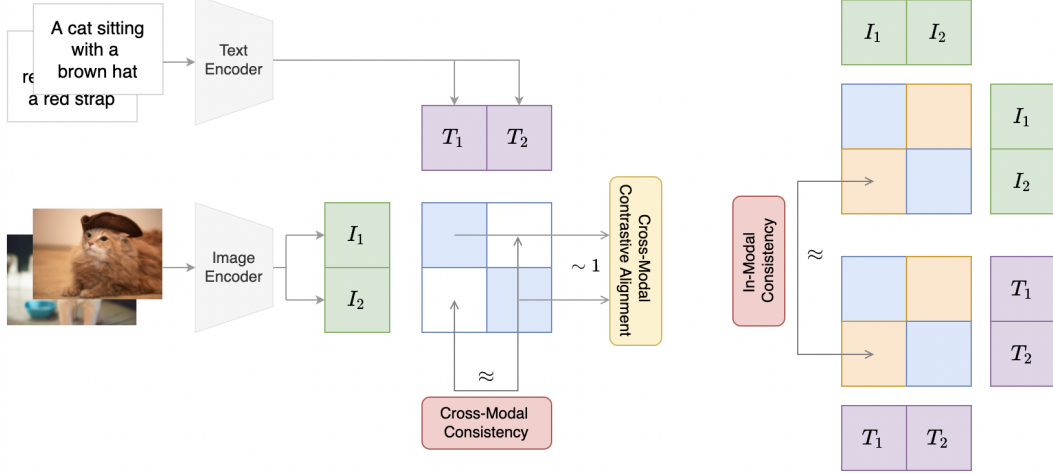


Figure 2: Illustrative overview for CYCLIP ($N = 2$). It consists of 3 major components: (a) cross-modal contrastive alignment, (b) cross-modal consistency, and (c) in-modal consistency. Only (a) is present in CLIP, whereas our proposed regularizers in (b) and (c) mitigate inconsistency.

associated with the k -nearest training images. Formally, we define a consistency score that measures the synchrony between the predicted labels in the image and text spaces as:

$$\text{Consistency Score}_k = \frac{1}{N} \sum_{j=1}^N \mathbb{1} [P_I^k(I_j) = P_T(I_j)] \quad (2)$$

where N is the number of test images. In our experiments (discussed in detail in §3), we found the CLIP’s consistency score ($k = 1$) to be 44%, 16%, and 16% on the standard benchmarks CIFAR-10, CIFAR-100, and ImageNet1K, respectively, showing a very high degree of disagreement in the image and text spaces. In the following section, we describe our approach to alleviate the inconsistent inference problem and quantitatively show that our solution improves the consistency score in §4.1.

2.3 Cycle Consistent Representation Learning via CYCLIP

We showed that the visual representations learned by CLIP could be inconsistent when used for inference in the image and text spaces. To mitigate this problem, we propose CYCLIP, a learning framework that builds upon CLIP by augmenting the contrastive loss in Eq. 1 with additional geometric consistency regularizers. The intuition follows directly from Figure 1 (b), where we showed that inconsistency in the image and text spaces could be eliminated if we symmetrize the similarity between the two mismatched image-text pairs and the similarity between the image-image pair and the text-text pair. We formalize this intuition with two consistency regularizers.

(1) The **cross-modal consistency** regularizer reduces the gap in the similarity scores between the embeddings of all the mismatched image-text pairs in a batch, two at a time:

$$\mathcal{L}_{\text{C-Cyclic}} = \frac{1}{N} \sum_{j=1}^N \sum_{k=1}^N (\langle I_j^e, T_k^e \rangle - \langle I_k^e, T_j^e \rangle)^2. \quad (3)$$

(2) The **in-modal consistency** regularizer reduces the gap in the similarity scores between the embeddings of all combinations of image pairs and their corresponding text pairs in a batch:

$$\mathcal{L}_{\text{I-Cyclic}} = \frac{1}{N} \sum_{j=1}^N \sum_{k=1}^N (\langle I_j^e, I_k^e \rangle - \langle T_k^e, T_j^e \rangle)^2. \quad (4)$$

Hence, our overall loss for CYCLIP is given as:

$$\mathcal{L}_{\text{CYCLIP}} = \mathcal{L}_{\text{CLIP}} + \lambda_1 \mathcal{L}_{\text{I-Cyclic}} + \lambda_2 \mathcal{L}_{\text{C-Cyclic}} \quad (5)$$

where $\lambda_1 > 0$ and $\lambda_2 > 0$ are hyperparameters controlling the importance of the in-modal and cross-modal cyclic consistency regularizers relative to the contrastive loss in CLIP. We can also characterize the effect of the regularizers in terms of symmetrizing the in-modal and cross-modal similarity matrices, as illustrated in Figure 2. Note that the optimal solution to the contrastive loss formulation would push the similarity between the normalized embeddings of the matched pairs towards 1 while forcing all other pairs of similarities to 0, thereby also symmetrizing the cross-modal similarity matrix and minimizing the cross-modal consistency loss. However, this idealized scenario does not occur in practice, and we find that explicit regularization via cycle-consistency in CYCLIP facilitates improved learning, as we show in our experiments.

3 Experiments

Setup: We use Conceptual Captions 3M [53] (CC3M) image-caption pairs as the source of multi-modal pretraining data for all our models. Note while this dataset is smaller than the custom dataset (400 million pairs) used in the original work on CLIP [46], it is suitable for our available data and compute and has been used for benchmark evaluations in many subsequent works on language-image pretraining [5, 33, 37, 57]. Following prior work [46], our CLIP models use ResNet-50 as the image encoder and a transformer architecture as the text encoder. Further, we train our models from scratch for 64 epochs on 4 V100 GPUs with a batch size of 128 and an initial learning rate of 0.0005 with cosine scheduling and 10000 warmup steps. The dimension of the image and text embeddings is 1024. For CYCLIP, we use $\lambda_1 = 0.25$ and $\lambda_2 = 0.25$ across all our experiments.

3.1 Zero-Shot Transfer

We compare the zero-shot performance of CLIP and CYCLIP on standard image classification datasets: CIFAR-10, CIFAR-100 [31], and ImageNet1K [50]. We follow the evaluation strategy suggested by [46] for zero-shot classification using prompt engineering. For each dataset, we use the names of the classes to form a set of natural sentences such as ‘a photo of the {class name}’, ‘a sketch of the {class name}’ and more. These are passed through the text encoder to get a set of text embeddings for that class. This set of text embeddings are ℓ_2 -normalized, averaged, and further ℓ_2 -normalized to obtain a single text embedding for that class. For a given image, the image embedding is obtained as described in §2. The class whose text embedding (as described above) is closest to the test image is taken to be the predicted label. The zero-shot performance of the models is presented in Table 1.

Table 1: Zero-shot TopK classification accuracy (%) where $K \in \{1, 3, 5\}$

	CIFAR-10			CIFAR-100			ImageNet1K		
	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5
CLIP	46.54	78.22	91.16	18.69	34.72	43.97	20.03	33.04	39.35
CYCLIP	51.45	79.57	91.80	23.15	41.46	50.66	22.08	35.98	42.30
%GAIN	+10.6	+1.7	+0.7	+23.9	+19.4	+15.2	+10.2	+8.9	+7.5

We observe that the CYCLIP outperforms CLIP across all the datasets and on all TopK metrics, with gains in the range of 10% - 24% for $K=1$. Our results on zero-shot transfer indicate the usefulness of having geometrical consistency for improved downstream performance of CLIP.

3.2 Robustness to Natural Distribution Shifts

One of the major successes of CLIP was its state-of-the-art performance on the natural distribution shift benchmarks. These benchmarks include images depicting sketches, cartoons, adversaries generated using attacks on trained ImageNet models. In Table 2, we evaluate the zero-shot classification accuracy of CYCLIP on four natural distribution shift benchmarks for the ImageNet dataset: ImageNetV2 [49], ImageNetSketch [58], ImageNet-A [27], and ImageNet-R [25].

For most of the distribution shift benchmarks, both CLIP and CYCLIP undergo a significant reduction in their zero-shot performance compared to the original ImageNet1K dataset (last three columns in

Table 2: Zeroshot Classification on Natural Distribution Shifts (%)

	ImageNetV2			ImageNetSketch			ImageNet-A			ImageNet-R		
	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5
CLIP	16.91	29.28	34.99	10.37	19.15	24.20	4.23	11.35	16.88	24.32	39.69	47.20
CYCLIP	19.22	32.29	38.41	12.26	22.56	28.17	5.35	13.53	19.51	26.79	42.31	50.03
%GAIN	+13.7	+10.3	+9.8	+18.2	+17.8	+16.4	+26.5	+19.2	+15.6	+10.2	+6.6	+6.0

Table 1). However, we observe that CYCLIP outperforms CLIP on all of the datasets considered in this experiment by a significant margin of improvement (10 - 27%). This result indicates that having cyclic consistency in the learned representations preserves the robustness on the traditional datasets.

3.3 Linear Probing

While the primary focus of CLIP and CYCLIP is zero-shot generalization, we can also assess if the benefits of our cyclic consistency constraints in mitigating inconsistency can be recovered with extra in-domain and in-modality supervision i.e., in the presence of in-distribution training samples from in-domain visual datasets. To this end, we conduct an additional experiments on linear probing where we fit a linear classifier on the representations learned by the visual encoder (ResNet-50) of CLIP and CYCLIP on a range of image classification datasets.

Table 3: Transfer CLIP and CYCLIP to 14 downstream visual datasets using linear probing. Our CYCLIP performs marginally better on 9 out of 14 datasets. For training ImageNet1K, we use a random subset of 50K images from its original training dataset.

	Caltech101	CIFAR-10	CIFAR-100	DTD	Aircraft	Flowers102	Food101	GTSRB	ImageNet1K	OxfordPets	SST2	StanfordCars	STL10	SVHN	Average
CLIP	79.80	78.26	54.85	59.02	28.00	83.50	54.44	69.72	35.93	57.66	53.82	20.00	89.23	47.28	57.96
CYCLIP	80.33	76.98	55.74	63.44	27.86	82.96	54.96	71.70	37.12	56.82	53.74	22.14	90.10	48.01	58.71

We present our results in Table 3. We find that both CLIP and CYCLIP can recover most of the performance lost due to inconsistency when provided extra in-domain and in-modality supervision, with CYCLIP marginally outperforming the CLIP on 9 out of 14 visual datasets.

4 Analysis

Previously, we demonstrated the gains of CYCLIP over CLIP on downstream tasks that involve joint reasoning over the image and text spaces. In the current section, we wish to better understand the relative behavior of the two models on a set of challenging tasks.

4.1 Consistency in Image and Text Spaces

We begin by quantitatively measuring the inconsistency problem illustrated in Figure 1. That is, we wish to evaluate to what extent are the predictions in the image-text space (zero-shot) consistent with the ones made purely within the image space, as measured by our consistency metric in Eq. 2.

Table 4 presents our results over standard benchmarks (CIFAR-10, CIFAR-100, ImageNet1K). The consistency score is calculated over 10K, 10K, and 50K testing images of the CIFAR-10, CIFAR-100 and ImageNet dataset respectively. We use 50K samples from the training set of each dataset for k-Nearest Neighbor prediction. CYCLIP is more consistent than CLIP across all the datasets as we explicitly symmetrize the cross-modal and in-modal distances. Hence, the representations learned by CYCLIP can be better used interchangeably than CLIP.

Table 4: Consistency score (%) trend for CLIP and CYCLIP across standard benchmarks . Top-k consistency score implies the fraction of times, the zero-shot predicted label in the text space is identical to the k-Nearest Neighbor predicted label in the image space (using the training dataset).

	CIFAR-10				CIFAR-100				ImageNet1K			
	Top1	Top3	Top5	Top10	Top1	Top3	Top5	Top10	Top1	Top3	Top5	Top10
CLIP	44.60	46.04	47.06	48.45	16.21	17.28	18.42	19.36	16.34	17.42	18.58	19.78
CYCLIP	48.81	50.89	52.30	53.71	20.43	21.96	23.18	24.31	19.20	20.31	21.95	23.94
%GAIN	+8.6	+9.5	+10.0	+9.8	+20.7	+21.3	+20.5	+20.4	+14.9	+14.2	+15.4	+17.4

4.2 Fine-grained and Coarse-grained Performance

In §3.1, we observed that CYCLIP outperforms CLIP on zero-shot transfer across various datasets. We perform an error analysis investigating both models’ coarse and fine-grained classification performance to understand the transfer phenomena better. Given a hierarchical class structure dataset, coarse-grained classification differentiates between high-level (parent) classes, i.e., zero-shot classification into aquatic mammals and fish. The fine-grained classification task focuses on differentiating low-level (child) classes, i.e., zero-shot classification into a dolphin, otter, and seal (subclasses of aquatic mammals). We perform this analysis on the CIFAR-100, ImageNet1K, ImageNetV2, ImageNetSketch, ImageNet-A, and ImageNet-R datasets.

Formally, we consider a test set of N image-subclass-superclass triplets, $\{I_j, C_j, P_j\}_{j=1}^N$, where I_j, C_j, P_j represent the image, the subclass (child) and superclass (parent) respectively. The image embedding $I_j^e \in \mathbb{R}^d$ is obtained as described in §2, and the subclass embedding $C_j^e \in \mathbb{R}^d$ and superclass embedding $P_j^e \in \mathbb{R}^d$ are obtained as described in §3.1. Let the total number of superclasses and subclasses in the dataset be n_p and n_c , respectively. Further, let F be a unique mapping from a subclass to the superclass, and G denote the inverse mapping from a superclass to the set of subclasses i.e. $\forall P \in \{1, \dots, n_p\}, G(P) = \{C : F(C) = P \text{ and } C \in \{1, \dots, n_c\}\}$. Under this setup, the fine-grained and coarse-grained accuracies are defined as:

$$\text{Fine-grained Accuracy} = \frac{1}{N} \sum_{j=1}^N \mathbb{1} \left[\operatorname{argmax}_{C \in G(P_j)} \langle I_j^e, C \rangle = C_j \right] \quad (6)$$

$$\text{Coarse-grained Accuracy} = \frac{1}{N} \sum_{j=1}^N \mathbb{1} \left[\operatorname{argmax}_{C \in \{1, \dots, n_c\}} \langle I_j^e, C \rangle \in G(P_j) \right] \quad (7)$$

In Figure 3 we visualize how CLIP and CYCLIP compare with each other on the above metrics. The difference between the zero-shot performance of CYCLIP and CLIP is much more significant for coarse-grained classification than fine-grained classification across all the datasets. This observation indicates that concept-level knowledge is better captured in CYCLIP compared to CLIP. The drastic difference in the coarse-grained performance of CYCLIP and CLIP may be attributed to the rigid separation that the default cross-entropy loss in CLIP enforces between the positive pairs and negative pairs, which might degrade performance when some pairs in the negative batch belong to a similar entity. However, CYCLIP does not suffer from this problem as much because it poses cycle constraints on the overall geometry of all the data pairs rather than forcing a rigid separation.

4.3 Alignment and Uniformity on the Unit Hypersphere

[59] argues that contrastive learning directly optimizes for (a) alignment (closeness) of the representations of the positive pairs and (b) uniformity (coverage) of the representation space on the unit hypersphere. We extend these properties for multimodal contrastive representation learning as:

$$\text{Alignment} = \frac{1}{N} \sum_{j=1}^N \langle I_j^e, T_j^e \rangle \quad \text{Uniformity} = \log \left(\frac{1}{N(N-1)} \sum_{j=1}^N \sum_{k=1, j \neq k}^N e^{-\langle I_j^e, T_k^e \rangle} \right) \quad (8)$$

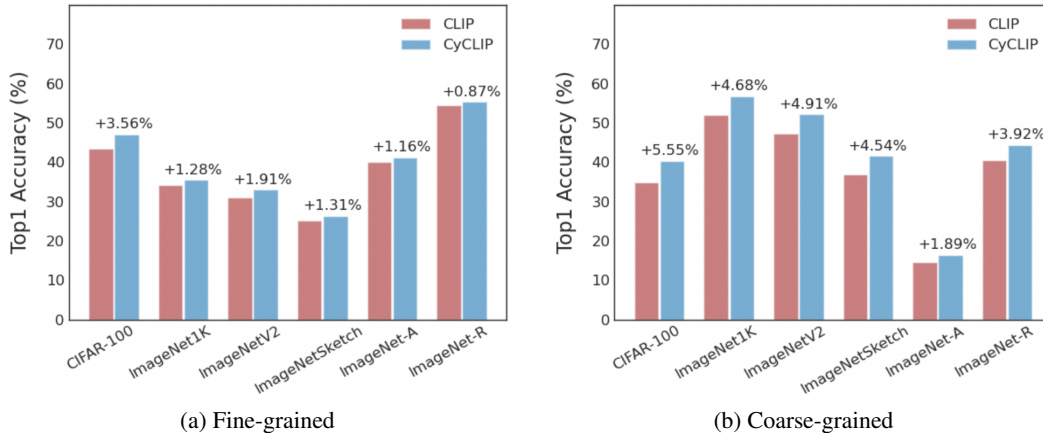


Figure 3: The gap between the performances of CLIP and CyCLIP is much larger in coarse-grained scenario highlighting better entity-level knowledge representation in CYCLIP.

We desire our models to achieve high alignment and uniformity scores so that the image-text representations are close for the matched pairs and better spread over the unit hypersphere for different categories. We analyze the effect of cross-modal and in-modal consistency on the alignment and uniformity of the shared representations. For this, we train two ablated versions of CYCLIP, 1) C-CYCLIP with only cross-modal consistency component i.e. $\lambda_1 = 0, \lambda_2 = 0.5$, and 2) I-CYCLIP with only in-modal consistency component i.e. $\lambda_1 = 0.5, \lambda_2 = 0$ (in Eq. 5). We design proxy captions for classes as discussed in §3.1 to act as text embeddings. We present the results in Table 5.

Table 5: Alignment and Uniformity values for CLIP and Cyclic CLIP models. We abbreviate Alignment by A, Uniformity by U, and Zero-shot Top1 classification accuracy (%) by ZS-Top1.

Model	CIFAR-10			CIFAR-100			ImageNet1K		
	A	U	ZS-Top1	A	U	ZS-Top1	A	U	ZS-Top1
CLIP	0.36	-0.27	46.54	0.36	-0.25	18.69	0.39	-0.18	20.03
CYCLIP	0.36	-0.34	51.45	0.37	-0.33	23.15	0.38	-0.32	22.08
I-CYCLIP	0.60	-0.57	50.97	0.60	-0.57	22.35	0.61	-0.55	21.21
C-CYCLIP	0.05	-0.02	55.52	0.06	-0.02	25.49	0.07	-0.02	21.73

We observe that I-CYCLIP learns representations that are better aligned in the representation space; however, they do not cover the hypersphere uniformly. The representations learned by C-CYCLIP are more uniformly spread but poorly aligned compared to I-CYCLIP. In this light, the components of CYCLIP can be seen to encourage a balance of good alignment and uniformity. Further, we find that CLIP is more uniform than CYCLIP in all datasets, but contrary to prior beliefs, this does not translate to improved downstream performance. C-CYCLIP has the best downstream zero-shot performance for CIFAR-10 and CIFAR-100 despite its poor alignment score. Further, all 3 variants of CYCLIP outperform CLIP on all 3 datasets, with CYCLIP performing the best on ImageNet1K.

4.4 Image-Text Retrieval

We evaluate the effectiveness of the proposed method on the cross-modal (image to text and text to image) retrieval downstream task in the zero-shot as well as fine-tuned settings. We consider the standard benchmark datasets: Flickr30K [43] and MSCOCO [8]. We assess our models on the test set of Flickr30K (1K) and MSCOCO (5K) obtained from the well-known Karpathy [30] split. Both the datasets contains 5 paired captions per image that makes text retrieval per image more easier than image retrieval per caption. We confirm the same in our results below. We perform fine-tuning on the Karpathy’s training split with the batch size of 48. We fine-tune on Flickr30K for 10 epochs and MSCOCO for 5 epochs. All the other hyperparameters are identical to that of pre-training.

Table 6: Zero-shot and fine-tuned cross-modal image-text retrieval (text-to-image and image-to-text) results of CLIP and CYCLIP on Flickr30K and MSCOCO datasets.

		Flickr30K (1K)						MSCOCO (5K)					
		Text Retrieval			Image Retrieval			Text Retrieval			Image Retrieval		
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
Zero-shot	CLIP	88.2	93.9	95.8	29.9	57.2	68.0	82.1	85.6	87.8	8.4	19.5	26.6
	CYCLIP	88.1	93.7	95.9	30.9	57.8	69.1	82.1	85.6	87.7	8.6	20.0	27.0
Fine-tuned	CLIP	91.9	97.0	98.0	46.3	74.7	83.6	83.2	87.6	90.0	10.6	23.9	31.3
	CYCLIP	92.3	97.0	98.4	47.3	76.6	85.4	83.2	87.8	90.3	11.4	25.8	33.4

Table 6 presents our cross-modal image-text retrieval results for CLIP and CYCLIP. In the zero-shot setting, we find that CYCLIP marginally outperforms CLIP on the image retrieval task on both the datasets. The relatively lower performance of both CLIP and CYCLIP in the zero-shot setting may be attributed to the more complicated nature of the two datasets where the models are expected to find similarities between the image and text at multiple resolutions as opposed to image classification where there is mostly single object to be matched with a simpler caption. It is not clear as to what distinctions in the raw input and text space are reflected in the embedding space too. Hence, we perform fine-tuning on both the datasets to better inform our models of the downstream datasets. In the fine-tuning setting, we find that the performance of both the models increases across both the datasets. However, we observe clear benefits of the soft consistent regularization on the image retrieval results for both the datasets.

4.5 CYCLIP preserves the Effective Robustness of CLIP

[36] shows that there is a strong correlation between the in-distribution and out-of-distribution generalization of the models trained on ImageNet1K, as illustrated by the linear fit (red) in Figure 4. Ideally, any model that does not undergo distribution shift would fall on the $y = x$ trendline (black). For other models, the deviations of the models from this ideal fit indicate their effective robustness. Previously, [45] showed that the zero-shot CLIP classifier trained on 400M image-text pairs improves effective robustness significantly compared to prior approaches to robustness. Subsequently, [28] demonstrated that CLIP models trained at small scales also exhibit high effective robustness that allows them to be used as a proxy to study the robustness properties of CLIP.

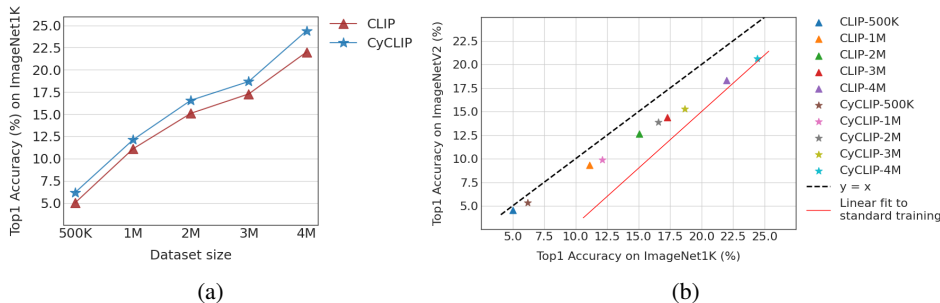


Figure 4: Effect of varying the training dataset size on (a) Classification accuracy on ImageNet1K and (b) Effective Robustness on ImageNetV2.

We evaluate the effect of cyclic consistency on effective robustness. We trained 4 CLIP and CYCLIP models, varying the training dataset sizes from 500K to 4M image-text pairs from the CC3M + CC12M datasets. In Figure 4, (a) we observe that for all training data sizes, CYCLIP shows a significant improvement over CLIP, showcasing its effectiveness in a diverse set of data regimes. Further, Figure 4 (b) shows that CYCLIP lies way above the baseline trend and preserves the effective robustness of CLIP.

5 Related Work

Our work fits into the broader theme of unsupervised pretraining with multiple modalities and has been successfully applied for learning representations of modalities such as images, text, and speech [2, 15, 1, 60, 44, 34]. Similar to the unimodal setting, two predominant approaches for multimodal pretraining are contrastive and generative, as described below.

Contrastive Representation Learning: Contrastive learning was originally proposed for self-supervised representation learning in the unimodal context where the embeddings of a sample are brought closer to an augmented version of the sample. In contrast, the embeddings are pushed away for other samples, and their augmentation [11, 52, 40, 56, 21, 7, 16, 41, 67, 23, 18]. [64] and [3] impose additional constraints to remove redundancies and prevent dimensional collapse in the visual representations. Recently, contrastive learning has also been used to learn robust representations of the multimodal data [63, 47]. Many works use additional losses to imbibe extra supervisory multimodal knowledge during the training process [55, 66, 65, 14, 35]. In this work, we focus on having cyclic consistency in addition to the contrastive loss to learn more robust image-text representations.

Contrastive Language-Image Pretraining: CLIP [46], ALIGN [29] and BASIC [42] have enjoyed great success in extending contrastive learning to paired image-text data, with impressive zero-shot classification and robustness performance. These works have been further extended recently to include visual self-supervision [37], additional nearest neighbor supervision [33], and utilization of unpaired data [57]. Our work complements much of this literature as it identifies consistency regularizers that can be augmented to the learning objective of the above works.

Generative Representation Learning: Generative models have been applied for learning representations of multimodal data [61, 54]. In particular, [68, 62, 10] proposed a notion of cyclic consistency for learning from unpaired multimodal data using GANs [17], which was extended later to normalizing flows [19, 20]. While these works focus on regularizing a generative mapping between modalities, our notion of cycle consistency applies to embeddings learned via a contrastive framework.

6 Conclusion

We presented CYCLIP, a framework for cycle consistent multimodal representation learning for image and text modality. The main benefits of CYCLIP stem from including cross-modal consistency and in-modal consistency regularizers to prevent inconsistent inference in the image and text spaces. Empirically, we show that CYCLIP performs much better than CLIP on zero-shot classification and is more robust on benchmarks for distributional robustness. We also showed that the representations learned by CYCLIP are more consistent than CLIP and better capture concept-level knowledge, as evidenced by our analysis of fine-grained and coarse-grained accuracies.

We believe this work can motivate further studies on understanding the geometry of the representation spaces learned via the contrastive objective applied to paired multimodal data and, in particular, identify conditions and regularization strategies under which the learned representations are synergistic across the various modalities for downstream applications.

One important future direction and a current limitation is scaling CYCLIP to larger datasets. While we do not possess the resources for this study, it is imperative to study the extent to which the benefits of cycle consistency remain at the scale on which the original CLIP was trained (400M image-text pairs). Finally, for real-world deployment of CLIP and their variants, such as CYCLIP, we need to be cautious about amplifying societal biases as these models are trained on large uncurated datasets scraped from the web [9]. Additionally, it is easy to add malicious data to the web, which poses a severe security threat [5]. Alleviating such harms is an important and active area of research.

Acknowledgements

This research is supported by an Adobe Data Science Research Award for Aditya Grover. We would like to thank the IDRE’s Research Technology group for the GPU computing resources on the UCLA Hoffman2 Cluster. We also want to thank Tung Duc Nguyen, Satvik Mashkaria, Siddharth Krishnamoorthy, Varuni Sarwal, and Ashima Suvana for their helpful suggestions.

References

- [1] Y. Aytar, C. Vondrick, and A. Torralba. See, hear, and read: Deep aligned representations. *arXiv preprint arXiv:1706.00932*, 2017.
- [2] T. Baltrušaitis, C. Ahuja, and L.-P. Morency. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443, 2018.
- [3] A. Bardes, J. Ponce, and Y. LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- [4] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- [5] N. Carlini and A. Terzis. Poisoning and backdooring contrastive learning. *arXiv preprint arXiv:2106.09667*, 2021.
- [6] S. Changpinyo, P. Sharma, N. Ding, and R. 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*, pages 3558–3568, 2021.
- [7] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [8] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- [9] J. Cho, A. Zala, and M. Bansal. Dall-eval: Probing the reasoning skills and social biases of text-to-image generative transformers. *arXiv preprint arXiv:2202.04053*, 2022.
- [10] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8789–8797, 2018.
- [11] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, volume 1, pages 539–546. IEEE, 2005.
- [12] K. Crowson, S. R. Biderman, D. Kornis, D. Stander, E. Hallahan, L. Castriato, and E. Raff. Vqgan-clip: Open domain image generation and editing with natural language guidance. *ArXiv*, abs/2204.08583, 2022.
- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [14] K. Desai and J. Johnson. Virtex: Learning visual representations from textual annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11162–11173, 2021.
- [15] J. Duan, L. Chen, S. Tran, J. Yang, Y. Xu, B. Zeng, C. Tao, and T. Chilimbi. Multi-modal alignment using representation codebook. *arXiv:2203.00048*, 2022.
- [16] T. Gao, X. Yao, and D. Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021.
- [17] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.

- [18] J.-B. Grill, F. Strub, F. Alché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*, 33: 21271–21284, 2020.
- [19] A. Grover, M. Dhar, and S. Ermon. Flow-gan: Combining maximum likelihood and adversarial learning in generative models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [20] A. Grover, C. Chute, R. Shu, Z. Cao, and S. Ermon. Alignflow: Cycle consistent learning from multiple domains via normalizing flows. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 4028–4035, 2020.
- [21] M. Gutmann and A. Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 297–304. JMLR Workshop and Conference Proceedings, 2010.
- [22] R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006.
- [23] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020.
- [24] D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli, M. Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021.
- [25] D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. L. Zhu, S. Parajuli, M. Guo, D. X. Song, J. Steinhardt, and J. Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8320–8329, 2021.
- [26] D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15262–15271, 2021.
- [27] D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. X. Song. Natural adversarial examples. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15257–15266, 2021.
- [28] G. Ilharco, M. Wortsman, R. Wightman, C. Gordon, N. Carlini, R. Taori, A. Dave, V. Shankar, H. Namkoong, J. Miller, H. Hajishirzi, A. Farhadi, and L. Schmidt. Openclip. *Zenodo*, July 2021. doi: 10.5281/zenodo.5143773. URL <https://doi.org/10.5281/zenodo.5143773>. If you use this software, please cite it as below.
- [29] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021.
- [30] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137, 2015.
- [31] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. *Citeseer*, 2009.
- [32] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.

- [33] Y. Li, F. Liang, L. Zhao, Y. Cui, W. Ouyang, J. Shao, F. Yu, and J. Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. *arXiv:2110.05208*, 2021.
- [34] K. Lu, A. Grover, P. Abbeel, and I. Mordatch. Pretrained transformers as universal computation engines. *arXiv preprint arXiv:2103.05247*, 2021.
- [35] S. Mai, Y. Zeng, S. Zheng, and H. Hu. Hybrid contrastive learning of tri-modal representation for multimodal sentiment analysis. *arXiv:2109.01797*, 2021.
- [36] J. P. Miller, R. Taori, A. Raghunathan, S. Sagawa, P. W. Koh, V. Shankar, P. Liang, Y. Carmon, and L. Schmidt. Accuracy on the line: on the strong correlation between out-of-distribution and in-distribution generalization. In *International Conference on Machine Learning*, pages 7721–7735. PMLR, 2021.
- [37] N. Mu, A. Kirillov, D. Wagner, and S. Xie. Slip: Self-supervision meets language-image pre-training. *arXiv:2112.12750*, 2021.
- [38] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv:2112.10741*, 2021.
- [39] C. G. Northcutt, A. Athalye, and J. Mueller. Pervasive label errors in test sets destabilize machine learning benchmarks. *ArXiv*, 2021. doi: 10.48550/ARXIV.2103.14749. URL <https://arxiv.org/abs/2103.14749>.
- [40] H. Oh Song, Y. Xiang, S. Jegelka, and S. Savarese. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4004–4012, 2016.
- [41] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [42] H. Pham, Z. Dai, G. Ghiasi, H. Liu, A. W. Yu, M.-T. Luong, M. Tan, and Q. V. Le. Combined scaling for zero-shot transfer learning. *arXiv preprint arXiv:2111.10050*, 2021.
- [43] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pages 2641–2649, 2015.
- [44] A. Quattoni, M. Collins, and T. Darrell. Learning visual representations using images with captions. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007.
- [45] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [46] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [47] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [48] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [49] B. Recht, R. Roelofs, L. Schmidt, and V. Shankar. Do imagenet classifiers generalize to imagenet? In *International Conference on Machine Learning*, pages 5389–5400. PMLR, 2019.
- [50] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.

- [51] A. Sanghi, H. Chu, J. G. Lambourne, Y. Wang, C.-Y. Cheng, and M. Fumero. Clip-forge: Towards zero-shot text-to-shape generation. *arXiv preprint arXiv:2110.02624*, 2021.
- [52] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [53] P. Sharma, N. Ding, S. Goodman, and R. Soiccut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018.
- [54] Y. Shi, B. Paige, P. Torr, et al. Variational mixture-of-experts autoencoders for multi-modal deep generative models. *Advances in Neural Information Processing Systems*, 32, 2019.
- [55] A. Singh, R. Hu, V. Goswami, G. Couairon, W. Galuba, M. Rohrbach, and D. Kiela. Flava: A foundational language and vision alignment model. *arXiv preprint arXiv:2112.04482*, 2021.
- [56] K. Sohn. Improved deep metric learning with multi-class n-pair loss objective. *Advances in neural information processing systems*, 29, 2016.
- [57] A. Tejankar, B. Wu, S. Xie, M. Khabsa, H. Pirsiavash, and H. Firooz. A fistful of words: Learning transferable visual models from bag-of-words supervision. *arXiv:2112.13884*, 2021.
- [58] H. Wang, S. Ge, Z. Lipton, and E. P. Xing. Learning robust global representations by penalizing local predictive power. *Advances in Neural Information Processing Systems*, 32, 2019.
- [59] T. Wang and P. Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *ICML*, 2020.
- [60] W. Wang, H. Bao, L. Dong, and F. Wei. Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. *arXiv preprint arXiv:2111.02358*, 2021.
- [61] M. Wu and N. Goodman. Multimodal generative models for scalable weakly-supervised learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- [62] Z. Yi, H. Zhang, P. Tan, and M. Gong. Dualgan: Unsupervised dual learning for image-to-image translation. In *Proceedings of the IEEE international conference on computer vision*, pages 2849–2857, 2017.
- [63] X. Yuan, Z. Lin, J. Kuen, J. Zhang, Y. Wang, M. Maire, A. Kale, and B. Faieta. Multimodal contrastive training for visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6995–7004, 2021.
- [64] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320. PMLR, 2021.
- [65] R. Zellers, X. Lu, J. Hessel, Y. Yu, J. S. Park, J. Cao, A. Farhadi, and Y. Choi. Merlot: Multimodal neural script knowledge models. *Advances in Neural Information Processing Systems*, 34, 2021.
- [66] P. Zhang, X. Li, X. Hu, J. Yang, L. Zhang, L. Wang, Y. Choi, and J. Gao. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5579–5588, 2021.
- [67] Y. Zhang, H. Jiang, Y. Miura, C. D. Manning, and C. P. Langlotz. Contrastive learning of medical visual representations from paired images and text. *arXiv preprint arXiv:2010.00747*, 2020.
- [68] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

A Additional Results

In addition to CYCLIP described in §2, we train two more instantiations of it by keeping either of the two consistency regularizers active in the loss objective (Eq. 5). The instantiation trained by setting $\lambda_1 = 0$ and $\lambda_2 = 0.5$ is termed as C-CYCLIP as only cross-modal consistency regularizer term is added to the loss objective. Similarly, we get I-CYCLIP where only in-modal consistency regularizer is added to the loss by setting $\lambda_1 = 0.5$ and $\lambda_2 = 0$. We evaluate C-CYCLIP and I-CYCLIP on most of the experiments discussed in the main text to understand their zero-shot transfer ability on standard datasets and robustness to natural distribution shifts.

A.1 Zero-shot Transfer

Table 7 presents our results of the zero-shot transfer experiment described in §3.1. We find that CYCLIP outperforms its sub-variants and the CLIP model on the ImageNet1K dataset. Interestingly, we observe that the C-CYCLIP model performs the best amongst all models on CIFAR-10 and CIFAR-100. Further, we notice that all the versions of CYCLIP (the last three rows of Table 7) are better than CLIP across all the datasets. This indicates that jointly improving the geometry of the learned image and text representations using cyclic consistency regularizers is practical for improved zero-shot transfer performance on the image classification task.

Table 7: Zero-shot TopK classification accuracy (%) where $K \in \{1, 3, 5\}$

	CIFAR-10			CIFAR-100			ImageNet1K		
	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5
CLIP	46.54	78.22	91.16	18.69	34.72	43.97	20.03	33.04	39.35
CYCLIP	51.45	79.57	91.80	23.15	41.46	50.66	22.08	35.98	42.30
C-CYCLIP	55.52	80.93	90.60	25.49	44.82	54.10	21.74	35.48	41.96
I-CYCLIP	50.97	79.51	90.75	22.35	39.25	48.45	21.22	34.72	41.05

A.2 Natural Distribution Shifts

We evaluate the performance of the ablated CYCLIP models on natural distribution shift benchmarks to evaluate their distributional robustness and present our results in Table 8. We find that all the CYCLIP models (the last three rows) outperform CLIP by a large margin across all the datasets. Interestingly, we observe that C-CYCLIP performs the best on three of the four natural distribution shift datasets (last three columns). Among the CYCLIP models, I-CYCLIP performs the worst across all the datasets. This indicates that having just in-modal consistency is not enough to preserve the distributional robustness on the traditional datasets.

Table 8: Zero-shot classification on Natural Distribution Shifts (%)

	ImageNetV2			ImageNetSketch			ImageNet-A			ImageNet-R		
	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5
CLIP	16.91	29.28	34.99	10.37	19.15	24.20	4.23	11.35	16.88	24.32	39.69	47.20
CYCLIP	19.22	32.29	38.41	12.26	22.56	28.17	5.35	13.53	19.51	26.79	42.31	50.03
C-CYCLIP	18.65	31.81	38.29	12.77	22.75	28.26	5.59	14.65	21.35	27.99	43.66	50.99
I-CYCLIP	18.40	30.63	36.82	10.87	20.17	25.49	4.75	11.84	16.92	24.55	38.71	45.68

A.3 Linear Probe

We perform linear probing to assess the performance of CLIP and CYCLIP models in the presence of extra in-domain and extra in-modality supervision, i.e., with access to the training datasets for the visual classification task. We first get the image embeddings of the training data points from our trained visual encoder (ResNet-50), followed by training a learnable classifier (linear mapping).

Further, we train our models from scratch for 32 epochs on a single RTX2080Ti GPU with a batch size of 16 and an initial learning rate of 0.005 with cosine scheduling. We use a weight decay of 0.01 for the non-bias parameters of the linear layer.

We experiment with 14 visual classification datasets, details of which are present in Table 9. We present our linear probing results in Table 10. In particular, we observe that the CYCLIP models marginally outperform the CLIP model on all the datasets except Flowers102 and OxfordIIITPet.

Table 9: Training data size, Testing data size, and the number of classes for different datasets used for Linear Probe evaluation.

Dataset	Classes	Train size	Test size
Caltech101	102	3060	6084
CIFAR-10	10	50000	10000
CIFAR-100	100	10000	10000
DTD	47	3760	1880
FGVCAircraft	100	6667	3333
Flowers102	102	2040	6149
Food101	101	75750	25250
GTSRB	43	26640	12630
ImageNet1K	1000	50000	50000
OxfordIIITPet	37	3680	3669
RenderedSST2	2	6920	1821
StanfordCars	196	8144	8041
STL10	10	5000	8000
SVHN	10	73257	26032

Table 10: Transfer CLIP and CYCLIP to 14 downstream visual datasets using linear probing. For training ImageNet1K, we use a random subset of 50K images from its original training dataset. Results have been averaged over 5 different seeds used for training the linear classifier.

	Caltech101	CIFAR-10	CIFAR-100	DTD	FGVCAircraft	Flowers102	Food101	GTSRB	ImageNet1K	OxfordIIITPet	RenderedSST2	StanfordCars	STL10	SVHN	Average
CLIP	79.80	78.26	54.85	59.02	28.00	83.50	54.44	69.72	35.93	57.66	53.82	20.00	89.23	47.28	57.96
CYCLIP	80.33	76.98	55.74	63.44	27.86	82.96	54.96	71.70	37.12	56.82	53.74	22.14	90.10	48.01	58.71
C-CYCLIP	80.83	79.15	56.15	61.13	28.25	82.14	56.68	70.33	37.81	57.31	56.33	22.05	90.05	46.81	58.93
I-CYCLIP	80.73	77.77	55.48	61.87	27.46	81.85	54.00	68.26	36.86	57.44	53.60	20.83	89.95	46.89	58.07

A.4 Consistency in Image and Text Spaces

Table 11 presents the results for the consistency score defined in Equation 2. We find that C-CYCLIP is the most consistent model despite missing the in-modal consistency regularizer term. Additionally, we observe that all the CYCLIP models (last three rows) are more consistent than the CLIP model across all the datasets. Hence, the explicit geometrization of the representations helps in improving the consistency of the image and text spaces.

A.5 Fine-grained and Coarse-grained Performance

In §4.2, we described a novel fine-grained and coarse-grained analysis to investigate the zero-shot transfer phenomena better. We have 1000, 200, and 200 subclasses distributed across 67, 59, and 49 superclasses for the ImageNet1K/V2/Sketch, ImageNet-A, and ImageNet-R datasets, respectively.

Table 11: Consistency score (%) trend for CLIP and CyCLIP across standard benchmarks . Top-k consistency score implies the fraction of times, the zero-shot predicted label in the text space is identical to the k-Nearest Neighbor predicted label in the image space (using the training dataset).

	CIFAR-10				CIFAR-100				ImageNet1K			
	Top1	Top3	Top5	Top10	Top1	Top3	Top5	Top10	Top1	Top3	Top5	Top10
CLIP	44.60	46.04	47.06	48.45	16.21	17.28	18.42	19.36	16.34	17.42	18.58	19.78
CyCLIP	48.81	50.89	52.30	53.71	20.43	21.96	23.18	24.31	19.20	20.31	21.95	23.94
C-CyCLIP	53.47	56.56	57.84	59.27	22.72	24.03	25.85	26.95	20.02	21.30	23.06	24.75
I-CyCLIP	48.34	49.98	51.27	52.55	20.32	21.32	22.56	23.65	18.82	19.81	21.33	22.92

Figure 5 illustrates the distribution of the subclasses across the superclasses for these ImageNet benchmarks. For the CIFAR-100 dataset, we have 100 subclasses uniformly distributed across 20 superclasses, i.e., 5 subclasses per superclass².

Figure 6 illustrates the performance of CLIP and CyCLIP models across the traditional datasets and their natural distribution shifted variants. In tandem with our findings in §4.2, we observe that even ablated CyCLIP models, I-CyCLIP and C-CyCLIP outperform the CLIP model on the coarse-grained and fine-grained classification metric. Additionally, the margin of improvement is larger in the case of coarse-grained analysis than the fine-grained one. This highlights that all the CyCLIP variants are better at capturing entity-level knowledge than CLIP in their joint representations. We also find that CyCLIP captures more entity-level knowledge than C-CyCLIP in ImageNet1K and ImageNetV2, datasets containing natural images belonging to large number of categories.

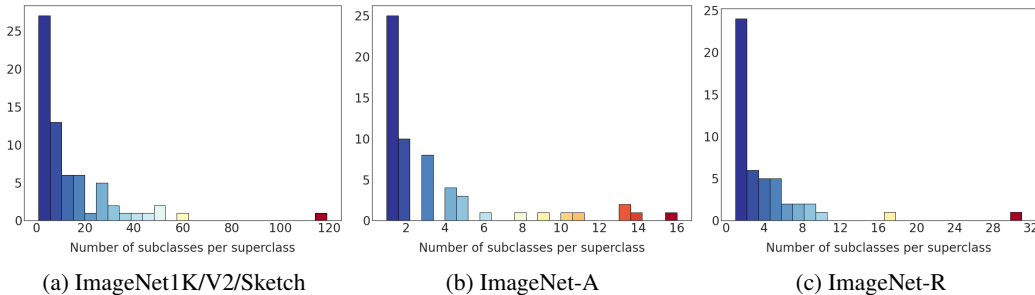


Figure 5: Distribution of superclasses across different number of subclasses. ImageNet1K/V2/Sketch and ImageNet-A/R have 1000 and 200 subclasses respectively.

A.6 Zero-shot Transfer: Additional Datasets

We present the results for zero-shot transfer downstream task, described in §3.1 in Table 12 across additional image classification datasets. The list of additional datasets is almost identical to the one used for linear probing §10. We include an additional dataset, *SUN397*, with 108754 images split across 397 classes. Further, we do not include some visual datasets as their labels were not amenable to constructing a simple proxy caption. For instance, *GTSRB* is a visual dataset with German traffic signs where a label is decided by the associated color, shape, and sign ID. Likewise, the images in the *FGVCAircraft* dataset are described by multiple attributes, and those in the *RenderedSST2* dataset are described by positive/negative sentiments, hence not relevant in this setting. Furthermore, we present our results on a cleaner test set of CIFAR-10, CIFAR100, and ImageNet1K [39]. We find that CyCLIP outperforms CLIP on the Top-1 classification accuracy with an average gain of 10.6% across all the datasets.

²<https://www.cs.toronto.edu/~kriz/cifar.html>

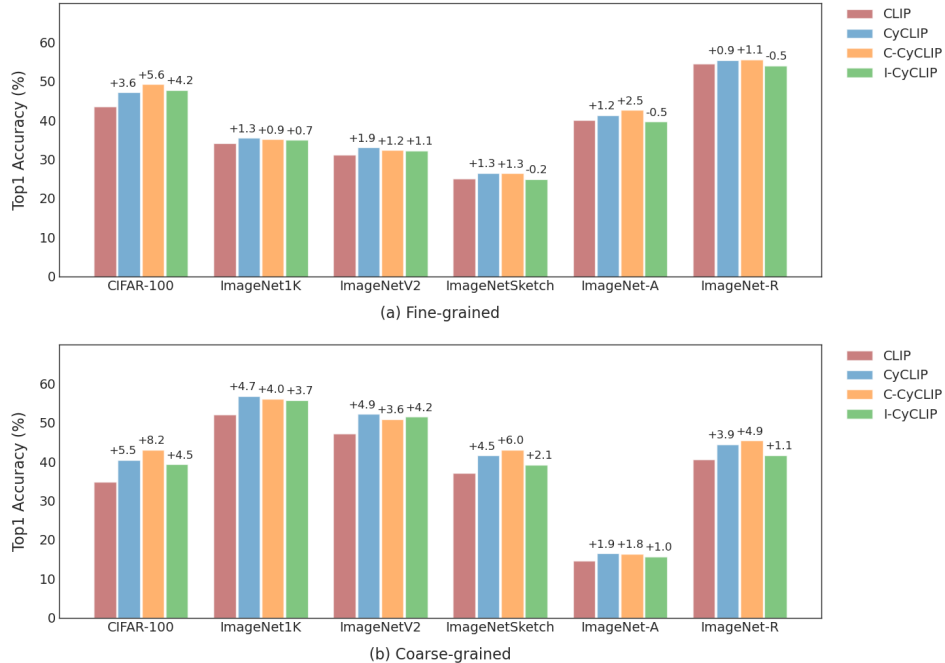


Figure 6: The gap between the performances of CLIP and our CyCLIP variants is much larger in coarse-grained scenario highlighting better entity-level knowledge representation in CyCLIPs.

Table 12: Zero-shot Top1 classification accuracy (%) across a battery of visual datasets. For the datasets marked with (*) we use their error-free labels instead of the original labels in their test sets.

	Caltech101	CIFAR-10*	CIFAR-100*	DTD	Flowers102	Food101	ImageNet1K*	OxfordIIITPets	StanfordCars	STL10	SVHN	SUN397	Average
CLIP	52.6	46.3	18.4	14.7	13.7	13.3	20.2	2.2	1.2	83.5	7.1	39.6	26.1
CyCLIP	60.7	50.8	22.0	19.2	10.9	14.0	22.5	2.3	1.6	86.1	13.3	42.7	28.8
Gain (%)	15.4	9.8	19.1	30.0	-20.4	5.0	11.8	7.3	34.4	3.1	86.3	7.7	10.6

B Additional Discussion

For the models pre-trained on the web-crawled data, one might argue against using the in-modal and cross-modal regularization as the image-text paired data is noisy where the caption may not fully describe the image at all resolutions. However, in an unsupervised setting, it is not pronounced as to what extent these distinctions should be reflected in the embedded representations.

As an example, we might think of the potential harm in bringing the embeddings of the two seemingly different captions (‘A person playing with a dog on a beach’, ‘A person playing with a cat in a room’) describing similar-looking images (with ‘dog, person, beach, ball, sky’, with ‘cat, person, room’) close using the in-modal consistency constraint. However, if the downstream task is to classify cats vs dogs, then the presence of any other objects in the images is a spurious correlation to be ignored. In such cases, we would desire the in-modal distance between text and image embeddings to be similar for consistent predictions in both modalities. Hence, the effectiveness of any information captured in a modality depends on the downstream task.

Moreover, one could make a similar counter-argument about the text and image modalities being inherently different against CLIP’s contrastive loss. It weighs each image and text pair equally, even

though some text captions might be more descriptive about the images than others and therefore should be assigned a higher weight. From a practical standpoint, our evidence suggests that soft consistency regularization in the form of additional loss terms, as in Equation 5, can be generally helpful for the downstream tasks and domain settings of interest. The relative gains can indeed be different across tasks and domains. For example, in Figure 3, while our regularizers lead to gains in fine and coarse-grained classification, we noted that the relative gains are much higher for coarse-grained settings due to the improved consistency.

C Pretraining and Implementation details

C.1 Dataset

Conceptual Captions 3M (CC3M) [53] is an open-source dataset consisting of approximately 3.3M image-caption pairs scraped from the web. The dataset (including train/validation split) is made available by Google³. Due to the broken image URLs, we use a subset of 2,631,703 image-caption pairs for pre-training. Further, Conceptual 12M (CC12M) [6] is a noisy extension to the CC3M dataset and contains approximately 12M image-caption pairs, covering a more diverse set of visual concepts. For the dataset ablation experiments in Figure 4 (a), with a data size of 3M pairs, we extend our CC3M dataset using a random subset of image-caption pairs (368,297 pairs) from CC12M⁴.

C.2 Model architecture

Our models use the same architecture as the original CLIP model presented in [46] with a ResNet-50 image encoder (38,316,896 parameters) and a transformer-based text encoder with a projection layer (63,690,240 parameters) to match the image embedding dimension of 1024. We use a weight decay for all the parameters during training, except for batch/layer norm, bias, and logit scale parameters.

C.3 Hyperparameter settings

The hyperparameters for the best-performing configuration were inspired from the original paper on CLIP [46] and mainly taken from mlfoundations⁵ repository. However, in our experiments, we use a batch size of 128 due to limited computation resources. To avoid hurting the performance of our models, we train our models for 64 epochs. For the cyclic consistency hyperparameters (λ_1, λ_2), we used a combination of zoom grid search and manual tuning on a small subset of 480K image-text pairs with training on 16 epochs and optimizing the contrastive loss on the Conceptual Captions validation set of 13,156 examples. We found that CyCLIP loss is not very sensitive in the non-zero parameter space between 0 and 1. Therefore, we choose a more simplistic setup with both the values as 0.25. We trained the CLIP and CyCLIP models on Azure with 4x Tesla-V100-SXM2-16GB GPUs and 24x CPUs for approximately 84 hours each. The hyperparameter settings are shown in table 13.

Table 13: Hyperparameters used for training the CLIP models

Hyperparameter	Value
Embedding dimension	1024
Logit scale range	0 to 4.6052
Epochs	64
Batch size	128
Learning rate	0.0005
Adam beta1	0.9
Adam beta2	0.99
Adam weight decay	0.1
Scheduler	Cosine
Learning rate warmup steps	10000

³<https://ai.google.com/research/ConceptualCaptions/download/>

⁴<https://github.com/google-research-datasets/conceptual-12m>

⁵https://github.com/mlfoundations/open_clip