# RANKCLIP: Ranking-Consistent Language-Image Pretraining

Yiming Zhang[*1,2]     Zhuokai Zhao[*3]
Zhaorun Chen[3]   Zhili Feng[4]   Zenghui Ding[1]   Yining Sun[1,2]
[1]HFIPS, Chinese Academy of Sciences   [2]University of Science and Technology of China
[3]University of Chicago   [4]Carnegie Mellon University

## Abstract

*Self-supervised contrastive learning models, such as CLIP, have set new benchmarks for vision-language models in many downstream tasks. However, their dependency on rigid one-to-one mappings overlooks the complex and often multifaceted relationships between and within texts and images. To this end, we introduce RANKCLIP, a novel pretraining method that extends beyond the rigid one-to-one matching framework of CLIP and its variants. By extending the traditional pair-wise loss to list-wise, and leveraging both in-modal and cross-modal ranking consistency, RANKCLIP improves the alignment process, enabling it to capture the nuanced many-to-many relationships between and within each modality. Through comprehensive experiments, we demonstrate the effectiveness of RANKCLIP in various downstream tasks, notably achieving significant gains in zero-shot classifications over state-of-the-art methods, underscoring the importance of this enhanced learning process. Code and model checkpoints are available at https://github.com/Jam1ezhang/RankCLIP.*

## 1. Introduction

In the realm of computer vision (CV) [55], natural language processing (NLP) [7], and multimodal deep learning [4, 25, 60], the alignment between visual and textual modalities [5, 49] has emerged as a cornerstone for downstream applications, ranging from image captioning [18] to zero-shot classification [42]. Contrastive Language-Image Pretraining (CLIP) [43] marks a significant advancement in this field, demonstrating incredible performance from training on large amounts of text-image pairs to create self-supervised models that understand [6, 22, 23] and generate [11, 44] descriptions of visual contents. Following the success of this contrastive learning paradigm, many recent works have been developed upon the original CLIP. More specifically, these enhancements focus on optimizing data efficiency through intrinsic supervision [28], as well
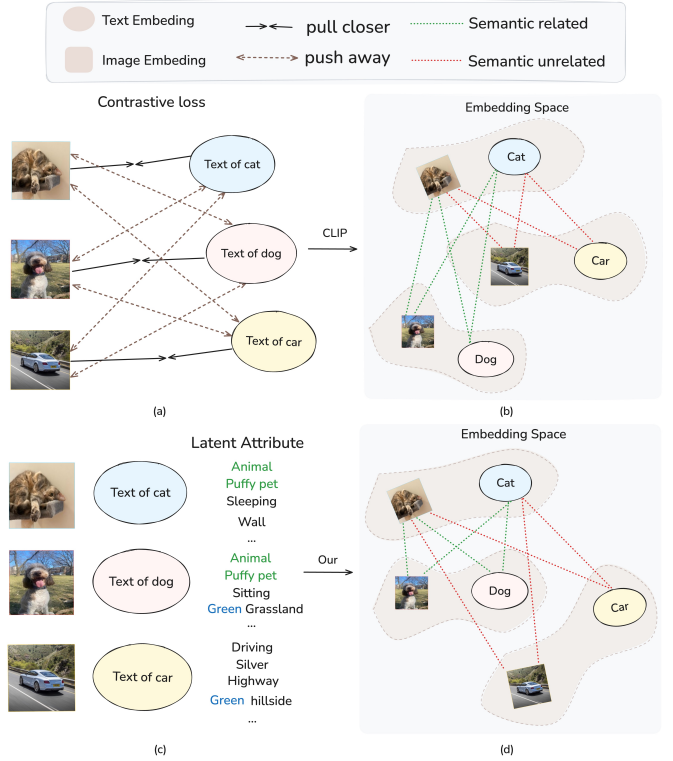


Figure 1. Comparison of learning outcomes between CLIP and RANKCLIP using three text-image pairs: dog (red), cat (blue), and car (yellow). (a) Contrastive loss treats all unmatched relationships equally, failing to distinguish latent similar attributes between dog and cat versus airplane. RANKCLIP addresses this issue by leveraging the shared attributes in (c) during training, improving the final trained embedding distribution from (b) to (d).

as improving downstream performance via cross-modal late interaction [59], hierarchical feature alignment [15], geometric consistency regularization [19], additional learning [36], adaptive loss [58], hierarchy-aware attentions [17], and softer cross-modal alignment [16].

Despite the improvements, these methods often have reliance on strict *pairwise, cross-modal, and one-to-one* mappings between images and texts, overlooking the actual *many-to-many* relationships that exist both *cross-modal* and

---

*These authors contribute equally. Correspondence to Zenghui Ding: dingzenghui@iim.ac.cn

*in-modal* in real-world data [8]. For example, as shown in Fig. 1, while pretrained models like CLIP can correctly classify `dog`, `cat` and `airplane`, they do not necessarily learn that `dog` and `cat` are more close to each other than `dog` and `airplane`, in terms of both in-modal (`dog` text is more similar to `cat` text than to `airplane` text) and cross-modal (`dog` text is more matched to `cat` image than to `airplane` image) similarities. Since it is rooted from the current contrastive loss that only the correct pairs are optimized while the rest of the unmatched pairs are treated the same, a large amount of information not used and unknown to the model during and after the training process.

Recognizing the complex *many-to-many* relationships as well as the rich information contained within both *in-modal* and *cross-modal* data, we introduce **Rank**ing-**C**onsistent **L**anguage-**I**mage **P**retraining, (**RANKCLIP**), which employs *ranking consistency* to learn and optimize similarity levels both between (cross-modal) and within (in-modal) the text-image pairs. The concept of ranking consistency stems from the simple observations that similar texts often correlate with similar images, as seen with the `dog`, `cat` and `airplane` example in Fig. 1. It effectively captures secondary similarity relationships among unmatched pairs, enabling the model to learn *more efficiently for free* compared to relying solely on matched pairs. Ranking consistency is conveniently modeled as an additional loss term to the traditional contrastive loss, requiring no extra external modules. It acts as a plug-and-play improvement for many existing methods, including those focusing on data-efficiency [28], potentially boosting performance in both efficiency and effectiveness.

The main contributions of this paper are: 1) RANKCLIP, a novel contrastive pretraining method that uses ranking consistency to exploit the many-to-many relationships within data, thereby enhancing performance in downstream tasks such as zero-shot classification and retrieval accuracy; and 2) through comprehensive experiments conducted on multiple datasets, we demonstrate the superior effectiveness of RANKCLIP in improving pretraining model performance without requiring any additional data or computational resources.

## 2. Related Work

Vision-language pretraining has witnessed significant advancements over the past years [3, 14, 31]. Models such as CLIP [43], ALIGN [26] and FLAVA [49] have pioneered the contrastive learning paradigm applied with text-image pairs, showcasing remarkable performance and robustness in downstream tasks. Many follow-up works, mostly built upon CLIP, have been proposed since then. Li et al. [28] introduced DeCLIP, improving zero-shot performance through intrinsic supervision. FILIP [59] advances CLIP's alignment between image patches and text with a cross-modal interaction mechanism. Gao et al. [15] developed PyramidCLIP, using hierarchical feature alignment to boost model efficiency and performance. Additionally, SLIP [36] merges self-supervised learning with CLIP pre-training for improved visual representation and accuracy. Goel et al. [19] introduced CyCLIP, augmenting CLIP with geometric consistency regularizers to enhance robustness and performance under varied conditions.

Recently, Yang et al. [58] introduced ALIP, an adaptive pre-training model that enhances language-image alignment using raw text and synthetic captions with dynamic adjustments. HiCLIP [17] refines CLIP by adding hierarchy-aware attentions to uncover semantic hierarchies in images and texts. EqSim [57] incorporates equivariance loss into vision-language models, significantly improving sensitivity to semantic changes in image-text pairs. Additionally, SoftCLIP [16] softens CLIP's one-to-one constraint, enabling more flexible cross-modal alignment through fine-grained adjustments.

Compared with existing approaches, RANKCLIP sets itself apart by fully leveraging the *many-to-many* relationships within each batch of text-image pairs, promoting learning from both matched and unmatched pairs with varying similarities by integrating in-modal and cross-modal *list-wise ranking consistencies* into the contrastive training objective. Crucially, RANKCLIP diverges from existing models' pair-wise training objective by adopting a global, list-wise optimization approach. In other words, it considers the rankings of all images and texts collectively within each batch, rather than focusing on pairwise similarities as seen in other methods.

## 3. RANKCLIP

RANKCLIP efficiently leverages the many-to-many relationships in real-world data by focusing on both matched and unmatched pairs. As in Fig. 2, it not only identifies if an image-text pair matches but also assesses their relative semantic similarities to other images and texts of both modalities in the dataset through self-supervised ranking consistency. Uniquely, RANKCLIP employs a list-wise loss for training batches, distinguishing it from other methods that solely rely on pair-wise relationships, as discussed in §2.

### 3.1. Ranking Model Formulation

RANKCLIP leverages the Plackett-Luce (PL) ranking model [20, 32, 41] to estimate the probability distribution over rankings for every image-text pair $(V_i, T_j)$, so that the consistency in their relative ordering with respect to a reference ranking can be measured. Specifically, for a given data pair, whether it is in-modal (image-image, text-text), or cross-modal (image-text), we calculate its in- or cross-modal cosine similarity $S_{ij}$ to serve as the score when measuring the alignment of its ranking with respect to another
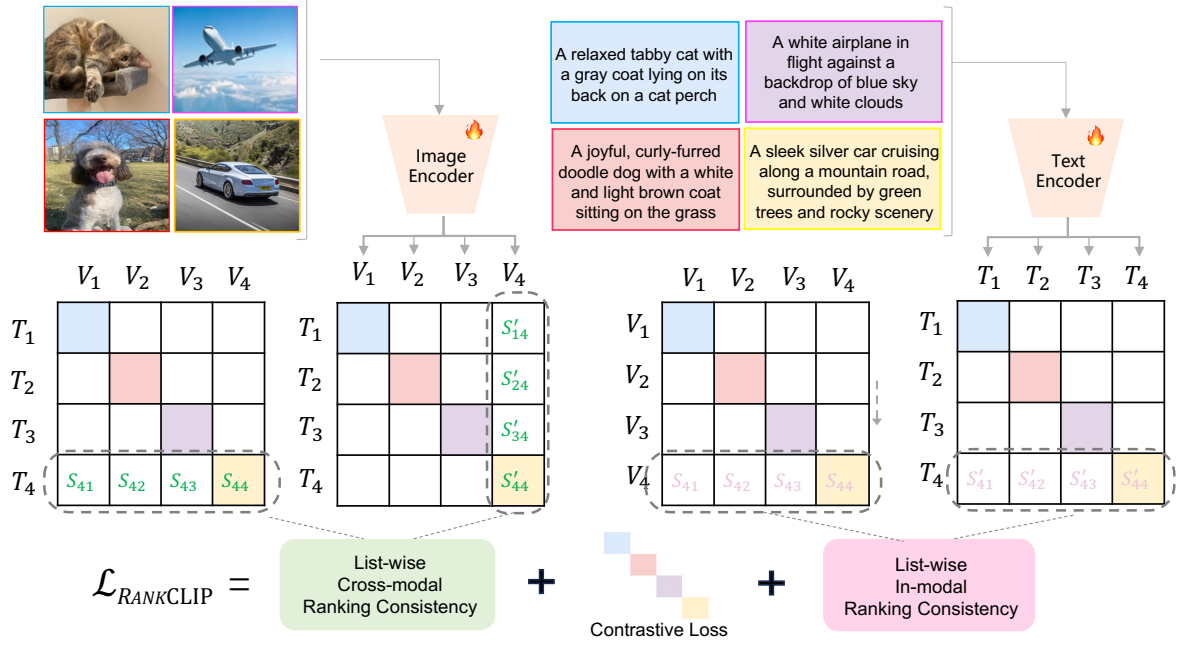
Figure 2. Overview of RANKCLIP. Unlike conventional contrastive loss, which includes only the middle term, RANKCLIP introduces both cross-modal and in-modal consistency terms by minimizing a self-supervised, list-wise ranking loss. Paired images and texts are indicated by matching contour line colors. $V$, $T$, and $S$ represent image embeddings, text embeddings, and similarity scores, respectively.

reference ranking $y_{\text{ref}}$.

Following [41], we first sort the reference ranking in a descending order to construct the optimal ranking $y^*$, and assume that the ego ranking $y$ is sampled from $y^*$. The probability that item $d$ with score $S_{ij}$ is ranked $k^{\text{th}}$ in the ego ranking $y$ from a set of items $\mathcal{D}$ is the score of $e^{S_{ij}}$ divided by the sum of scores for the items that have not been placed yet:

$$\pi(d \mid y_{1:k-1}, y_{\text{ref}}, \mathcal{D}) = \frac{e^{S_{ij}}}{\sum_{d' \in \mathcal{D} \setminus y_{1:k-1}} e^{S'_{ij}}}, \quad (1)$$

where $y_{1:k-1} = [y_1, y_2, ..., y_{k-1}]$ denotes the set of items ranked before $d$. Consequently, the probability of the entire ranking $y$ is the product of all individual placement probabilities:

$$\mathcal{P}(y, y_{\text{ref}}) = \prod_{k=1}^{K} \pi(y_k \mid y_{1:k-1}, \mathbf{y}_{\text{ref}}, \mathcal{D}). \quad (2)$$

RANKCLIP's objective is to maximize the consistency log-likelihood of the list ranking in one modality towards the reference ranking (from the same/in-modal and different/cross-modal data), which conveniently aligns with minimizing the negative log-likelihood loss:

$$\mathcal{L}_{\text{PL}} = -\log \mathcal{P}(y, y_{\text{ref}}) \quad (3)$$

### 3.2. Cross-modal Consistency Ranking

As illustrated by the green box in Fig. 2, RANKCLIP utilizes secondary relationships between unmatched visual and textual representations by constructing a list-wise rank loss. This approach ensures that the semantic similarity rankings between one image and multiple texts align with those between one corresponding text and multiple images. For example, as shown in Fig. 1, from the dog perspective, the semantic distance between dog image and cat text is closer compared to the plane text. This relationship should also apply between the dog text and the cat, plane images. Mathematically, Eq. (3) can be specified as:

$$\mathcal{L}_{\text{cross-modal}} = -\log \mathcal{P}(\mathbf{y}_{\text{image-text}}, \mathbf{y}_{\text{text-image}}) \quad (4)$$
$$= -\log \mathcal{P}(\hat{\mathbf{v}} \cdot \hat{\mathbf{t}}^{\mathbf{T}}, \hat{\mathbf{t}} \cdot \hat{\mathbf{v}}^{\mathbf{T}}) \quad (5)$$

By optimizing Eq. (4), RANKCLIP enhances its ability to bridge the semantic gap between modalities by leveraging nuanced unmatched correlations. This can also be viewed as learning a symmetric cosine-similarity matrix, further reinforcing semantic consistency across modalities.

### 3.3. In-modal Consistency Ranking

The pink box in Fig. 2 highlights the in-modal consistency component of the proposed rank loss. RANKCLIP ensures semantic consistency within each modality – image to image and text to text – enhancing the use of secondary unmatched relationships as an optimization objective. The underlying principle is that similar images should correspond to similar texts. For example, in Fig. 1, from the dog image perspective, the cat image is the most similar, followed by the plane image. This relationship should hold true

3

| | ImageNet1K | | | MSCOCO | | | | | |
| | | | | Image-to-Text Retrieval | | | Text-to-Image Retrieval | | |
| | Top-1 | Top-3 | Top-5 | Recall@1 | Recall@5 | Recall@10 | Recall@1 | Recall@5 | Recall@10 |
|---|---|---|---|---|---|---|---|---|---|
| CLIP [43] | 9.06% | 16.94% | 21.63% | 6.68% | 18.36% | 26.94% | 3.70% | 9.74% | 14.04% |
| CyCLIP [19] | 9.40% | 17.32% | 21.72% | 6.50% | 19.34% | **29.14%** | 3.72% | **11.16%** | **16.06%** |
| ALIP [58] | 9.71% | 18.31% | 23.07% | 6.04% | 18.04% | 26.92% | 3.70% | 10.22% | 14.38% |
| RANKCLIP | **10.16%** | **19.57%** | **24.01%** | **7.18%** | **19.46%** | 28.48% | **3.74%** | 10.28% | 14.18% |

Table 1. Zero-shot top-1, top-3, and top-5 classification accuracy on ImageNet1K, along with retrieval performance on MS-COCO. The proposed RANKCLIP consistently outperforms all baselines across both tasks. All models are trained on CC3M with ViT-B/32 backbone.

for their corresponding texts as well, where we utilize this to construct our $y$ and $y_{\text{ref}}$ from Eq. (3). Mathematically, Eq. (3) can be specified as:

$$\mathcal{L}_{\text{in-modal}} = -\log \mathcal{P}(\mathbf{y}_{\text{text-text}}, \mathbf{y}_{\text{image-image}}) \qquad (6)$$
$$= -\log \mathcal{P}(\hat{\mathbf{t}} \cdot \hat{\mathbf{t}}^{\mathbf{T}}, \hat{\mathbf{v}} \cdot \hat{\mathbf{v}}^{\mathbf{T}}) \qquad (7)$$

where $\hat{\mathbf{t}}$ and $\hat{\mathbf{v}}$ are the text and image batch embedding matrix, respectively. Via Eq. (6), the model can efficiently leverage the nuanced in-modal relationships to learn a richer and more structured semantic representation.

### 3.4. RANKCLIP Loss

Combining both cross-modal and in-modal consistency with the traditional contrastive loss (more details in Appendix B), the complete rank loss is thus formulated as:

$$\mathcal{L}_{\text{RANKCLIP}} = \mathcal{L}_{\text{CLIP}} + \lambda_1 \mathcal{L}_{\text{in-modal}} + \lambda_2 \mathcal{L}_{\text{cross-modal}} \qquad (8)$$

which is also depicted in Fig. 2. By supplementing the pairwise contrastive loss with cross-modal and in-modality ranking consistency loss, RANKCLIP systematically organizes embeddings to fully exploit both global and fine-grained unmatched relationships, which enhances the learning of more informative and accurate representations, better supporting downstream multi-modal tasks.

### 3.5. Training Recipe on Selecting $\lambda_1$ and $\lambda_2$

In the early stage of pre-training, rank consistency is highly unstable due to random initialization. Overemphasizing ranking consistency at this stage can impede the optimization of the embedding space. To address this, we gradually increase the weights $\lambda_1$ and $\lambda_2$ of the ranking loss as training progresses. Specifically, we have:

$$\lambda_1 = \lambda_2 = \text{clip}\left(\frac{3i-1}{n-1}, 0, 2\right)$$

where $i$ and $n$ denote the current training epoch and total number of epoch, respectively. The full RANKCLIP framework is outlined in Algorithm 1.

## 4. Experiments

### 4.1. Experimental Setup

**Baselines.** The most direct baseline to RANKCLIP is the original CLIP [43]. To further demonstrate the supe-

rior performance of RANKCLIP, we include CyCLIP [19], which introduces cyclic consistency constraints to enforce more robust alignment between visual and textual representations, improving generalization and semantic coherence. We also include ALIP [58], which leverages synthetic captions to enhance vision-language representation learning. More specifically, it employs a unique architecture that dynamically adjusts sample and pair weights to mitigate the impact of noisy or irrelevant data, making its approach complementary to ours. The training procedures and parameters for all models are detailed in Appendix A.

**Data.** All approaches are pretrained on the Conceptual Captions 3M (CC3M) dataset [48], which contains approximately 3.3 million text-image pairs. Although significantly smaller than CLIP's original dataset [24], CC3M remains a standard benchmark for vision-language pretraining, enabling strong zero-shot performance [2, 19, 28, 36, 53]. As discussed in §5.2, we also train CLIP and RANKCLIP on a larger 15M-image subset of YFCC100M [54] (YFCC15M) to study the effect of scaling up the dataset size.

### 4.2. Zero-shot Classification

We evaluate the zero-shot classification performance of CLIP [43], CyCLIP [19], ALIP [58], and RANKCLIP on ImageNet1K [12, 46]. As shown in Table 1, RANKCLIP consistently outperforms CLIP, which highlight the effectiveness of ranking consistency in enhancing language-image alignment with the same training data. Compared to CyCLIP [19], which enforces cyclic consistency to improve semantic coherence, RANKCLIP achieves higher accuracy across all metrics, suggesting that ranking consistency provides a more direct and effective regularization for representation learning. Additionally, RANKCLIP surpasses ALIP [58], indicating that ranking consistency is a stronger alternative to synthetic caption-based supervision. Notably, RANKCLIP shows the most significant improvement in top-1 accuracy, reinforcing its practical advantages where the highest-ranked prediction is most critical.

### 4.3. Zero-shot Cross-modal Retrieval

We further evaluate RANKCLIP on zero-shot cross-modal retrieval tasks, including image-to-text and text-to-image retrieval, using the MSCOCO [30] dataset. As shown in Table 1, RANKCLIP outperforms all baselines, though the

| | ImageNetV2-Matched | | | ImageNetV2-Threshold | | | ImageNetV2-Top | | | ImageNet-R | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Top1 | Top3 | Top5 | Top1 | Top3 | Top5 | Top1 | Top3 | Top5 | Top1 | Top3 | Top5 |
| CLIP [43] | 7.53% | 14.99% | 19.61% | 8.89% | 17.22% | 21.86% | 10.76% | 19.80% | 24.87% | 9.36% | 10.56% | 19.76% |
| CyCLIP [19] | 7.68% | 15.07% | 19.11% | 9.10% | 17.42% | 21.94% | 11.20% | 20.18% | 25.34% | 9.23% | 16.72% | 21.64% |
| ALIP [58] | 7.82% | 15.56% | 19.81% | 9.65% | 18.31% | 22.85% | 11.43% | 20.88% | 26.10% | 10.92% | 20.27% | 26.24% |
| RANKCLIP | **9.01%** | **16.95%** | **21.12%** | **10.32%** | **19.31%** | **24.13%** | **12.31%** | **22.11%** | **27.17%** | **11.34%** | **20.88%** | **26.94%** |

Table 2. Zero-shot top-1, 3, and 5 accuracy on ImageNet1K variants with *natural distribution shifts*. Compared to baselines, RANKCLIP achieves higher accuracies. Notably, these gains are more pronounced than on standard ImageNet1K, highlighting improved robustness.

| | CIFAR-10 | CIFAR-100 | DTD | FGVGAircraft | Food101 | GTSRB | OxfordPets | SST2 | STL10 | SVHN | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|
| CLIP [43] | 77.6% | 56.2% | 43.2% | 22.6% | 39.7% | 60.0% | 40.4% | 51.0% | 79.0% | **50.5%** | 52.0% |
| CyCLIP [19] | 76.8% | 54.3% | 45.8% | 19.2% | 37.5% | 58.6% | **44.2%** | 51.5% | **82.3%** | 41.3% | 51.2% |
| ALIP [58] | 71.1% | 49.1% | **47.1%** | 17.4% | 36.1% | 51.5% | 41.9% | 53.3% | 81.0% | 38.3% | 48.7% |
| RANKCLIP | **78.4%** | **56.6%** | 42.4% | **23.4%** | **40.2%** | **60.6%** | 40.6% | **53.4%** | 79.6% | 47.7% | **52.3%** |

Table 3. Linear probing accuracy on 10 downstream datasets using a ViT backbone.

improvements are less pronounced compared to zero-shot classification. This smaller margin of improvement may be attributed to the increased complexity of retrieval tasks, which require fine-grained image-text alignment across varying resolutions and object details – challenges distinct from the more direct pattern recognition in classification.

### 4.4. Robustness to Distribution Shifts

To evaluate the robustness of RANKCLIP under distribution shifts, we test all approaches on three variants of ImageNetV2 [45] and ImageNet-R [22], which assess resilience to different real-world deviations from ImageNet1K. As shown in Table 2, RANKCLIP consistently outperforms all baselines, achieving the highest accuracy across all datasets. These results indicate that RANKCLIP not only improves standard zero-shot classification but also enhances adaptation to real-world distribution shifts.

### 4.5. Linear Probing

We further assess whether the advantages of ranking consistency persist when supplemented with in-domain supervision. Specifically, we apply linear probing, where pretrained encoders remain fixed while a logistic regression classifier is trained on domain-specific datasets. We evaluate on 10 standard image classification benchmarks, including CIFAR-10, CIFAR-100, DTD [9], FGVG-Aircraft [33], Food101 [1], GTSDB [52], OxfordPets [40], SST2 [50], STL-10 [10], and SVHN [37]. As shown in Table 3, RANKCLIP achieves the highest average accuracy, demonstrating that ranking consistency enhances generalization even with additional in-domain supervision.
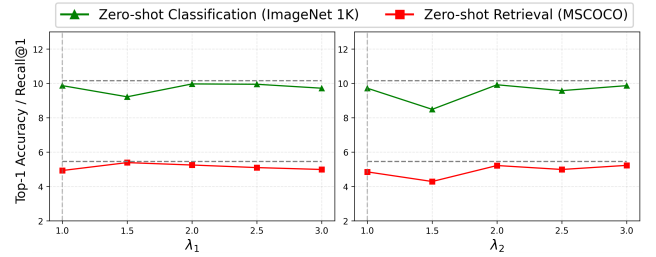


Figure 3. Effect of $\lambda_1$ and $\lambda_2$ on zero-shot classification (ImageNet1K) and retrieval (MSCOCO).

## 5. Ablation Studies

### 5.1. Different Weights of RANKCLIP Loss

In Eq. (8), we define the RANKCLIP loss as a linear combination of the original contrastive loss and the in-modal and cross-modal ranking consistency losses, weighted by $\lambda_1$ and $\lambda_2$, respectively. In §3.5, we introduce a training strategy that adaptively adjusts $\lambda_1$ and $\lambda_2$ at different stages of training. In this section, we analyze the impact of these weights and demonstrates that adaptive weighting further enhances RANKCLIP 's performance. Notice that all model variants follow the same pretraining setup detailed in Appendix A. As shown in Fig. 3, RANKCLIP outperforms CLIP even with fixed $\lambda_1$ and $\lambda_2$, highlighting the effectiveness of ranking consistency. And the adaptive weighting strategy further boosts accuracy by preventing the underutilization of ranking consistency at low weights and avoiding disruptions to contrastive learning at high weights.

### 5.2. Different Data Sizes

To assess the scalability of RANKCLIP, we trained both CLIP and RANKCLIP on 500k, 750k, 1M, and 3M text-
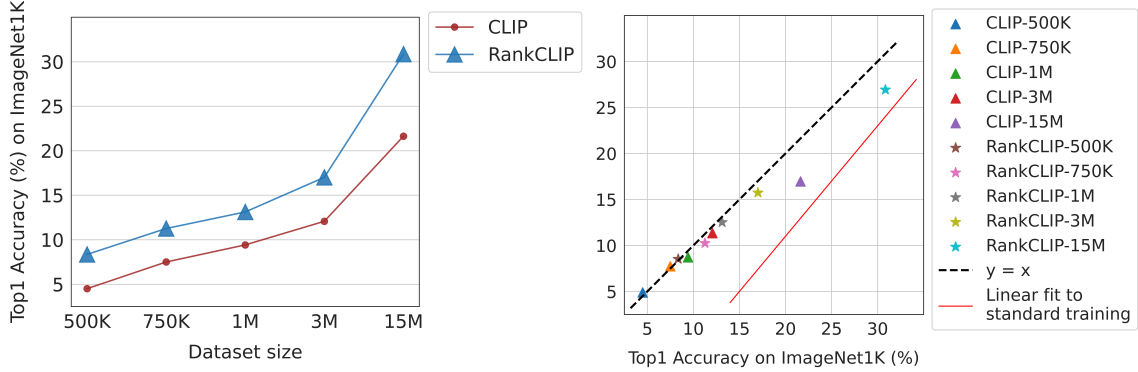
Figure 4. Ablation studies of CLIP and RANKCLIP trained with different data sizes. *Left*: zero-shot top-1 classification accuracy on ImageNet1K with various data sizes randomly sampled from CC3M. RANKCLIP consistently outperforms CLIP with significant margins. *Right*: zero-shot top-1 classification accuracy on ImageNet1K (horizontal axis) and ImageNet1K-R (vertical axis). RANKCLIP demonstrates better robustness as well as accuracy.

| Method | Vision Backbone | ImageNet1K | | | MSCOCO | | | | | | Linear Probing Avg. Acc. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Image-to-Text Retrieval | | | Text-to-Image Retrieval | | | |
| | | Top-1 | Top-3 | Top-5 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | |
| CLIP [43] | RN50 | 21.6% | 36.9% | 44.9% | 15.6% | 36.4% | 48.4% | 6.7% | 15.2% | 20.1% | 64.2% |
| RANKCLIP | | **30.9%** | **49.4%** | **57.6%** | **19.5%** | **42.6%** | **54.8%** | **7.5%** | **16.2%** | **21.6%** | **68.9%** |
| CLIP [43] | ViT-B/32 | 20.7% | 35.0% | 42.4% | 11.9% | 29.4% | 40.8% | 5.1% | 12.9% | 17.9% | 60.7% |
| RANKCLIP | | **26.2%** | **41.4%** | **48.9%** | **13.8%** | **33.8%** | **45.9%** | **6.0%** | **13.6%** | **18.6%** | **61.3%** |

Table 4. Zero-shot evaluation of CLIP and RANKCLIP trained with different vision backbones (ResNet-50 (RN50) and ViT-B/32) on ImageNet1K classification, MSCOCO cross-modal retrievals, and linear probing. "R@k" denotes Recall@k.

image pairs from CC3M, as well as 15M pairs from YFCC15M, following the procedure in Appendix A. Fig. 4 compares their performance on zero-shot top-1 classification accuracy for ImageNet1K (left) and averaged linear probing results (middle), where RANKCLIP consistently outperforms CLIP. Full, non-averaged linear probing results are provided in Appendix 7. Notably, the performance gains of RANKCLIP become more pronounced as dataset size scales from 1m to 15m, highlighting its superior scalability, which is critical for language-image pretraining.

Fig. 4 (right) further illustrates RANKCLIP's robustness across different dataset sizes. The horizontal axis represents top-1 accuracy on standard ImageNet1K, while the vertical axis shows accuracy on ImageNet1K-R. The black diagonal ($y = x$) denotes ideal robustness, with deviations below it indicating degradation under distribution shifts. RANKCLIP remains well above the red baseline, which represents typical in-distribution to out-of-distribution generalization [35], and stays close to the ideal line, demonstrating strong robustness.

### 5.3. Different Backbones: RN50 vs. ViT

We compare RANKCLIP and CLIP across ResNet-50 (RN50) and ViT-B/32 backbones to assess its generalization. As shown in Table 4, RANKCLIP consistently outperforms CLIP across all tasks. Specifically, for zero-shot classification on ImageNet1K, RANKCLIP improves top-1 accuracy by +9.3% with RN50 and +5.5% with ViT-B/32,

demonstrating stronger feature learning. In cross-modal retrieval, RANKCLIP achieves +3.9% in image-to-text and +0.8% in text-to-image with RN50, with smaller but consistent gains for ViT-B/32. RANKCLIP also improves linear probing accuracy by +4.7% with RN50 and +0.6% with ViT-B/32, confirming its advantage in representation learning. While both architectures benefit, RN50 sees the largest gains, suggesting that ranking consistency particularly enhances hierarchical feature extraction in CNNs.

## 6. Analysis

### 6.1. Modality Gap

We analyze the modality gaps of CLIP and RANKCLIP by visualizing 250 text-image pair embeddings, reduced to two dimensions using UMAP [34], and presenting a histogram of the gaps. The modality gap [29] refers to the separation between text and image embeddings in multimodal models, hindering joint representation learning. This gap, inherent from initialization and reinforced by contrastive learning in CLIP, challenges effective language-image modeling. Recent studies [27, 38, 51] suggest that reducing this gap improves multimodal representations and downstream performance. As shown in Fig. 6, RANKCLIP exhibits a significantly smaller modality gap than CLIP, demonstrating that our ranking consistency approach effectively enhances text-image alignment.

6

Caption: A painting of a table with fruit on top of it

Caption: A cute cat laying down in a sink

Figure 5. For a given text query, we present the top ten most semantically relevant images (ordered from left to right) obtained through both CLIP and RANKCLIP. In comparison to CLIP, our approach consistently retrieves images that more comprehensively align with the textual description, maintaining this advantage even after the correct reference image appears in the ranked results.

## 6.2. Alignment and Uniformity

Beyond reducing modality gap, effective contrastive learning should ensure a broad and uniform distribution over a hypersphere [56]. These objectives – similarity and uniformity – are quantified by alignment and uniformity scores, respectively. Following Goel et al. [19] and the notations in §3, we compute the alignment score $S_A$ and uniformity score $S_U$ as:

$$S_A = \frac{1}{N} \sum_{j=1}^{N} \hat{I}_j^T \hat{T}_j, \tag{9}$$

$$S_U = \log \left( \frac{1}{N(N-1)} \sum_{j-1}^{N} \sum_{k=1, j \neq k}^{N} \exp^{-\hat{I}_j^T \hat{T}_k} \right) \tag{10}$$

where $N$ is the number of text-image pairs. $S_A$ captures the average cosine similarity between corresponding text and image embeddings, while $S_U$ measures how evenly embeddings are spread across the space. High alignment indicates strong correlation between paired embeddings, whereas low uniformity suggests diverse and efficient embedding distribution—desirable for tasks like retrieval. As shown in Table 5, CLIP achieves stronger alignment but suffers from poor uniformity, leading to redundant representations. On the other hand, RANKCLIP along with two of its ablated version, RANKCLIP$_I$ and RANKCLIP$_C$ presents much better balance between alignment and uniformity. These results further suggest that optimizing solely for alignment or uniformity does not necessarily translate to better task performance.
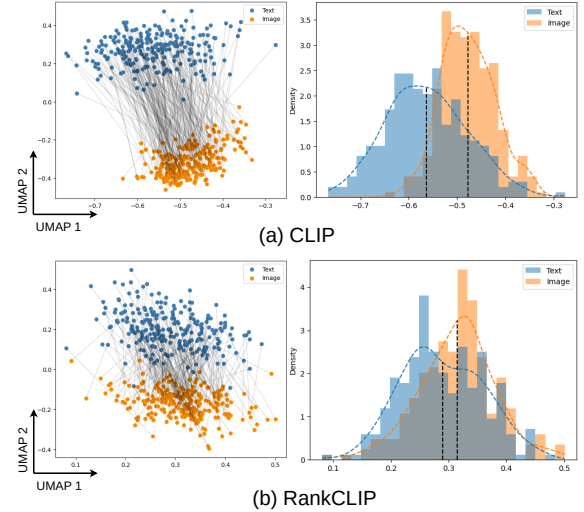


(a) CLIP

(b) RankCLIP

Figure 6. Scatter and histograms plots illustrating modality gaps of (a) CLIP and (b) RANKCLIP.

|  | CIFAR-10 | | CIFAR-100 | | ImageNet1K | |
|---|---|---|---|---|---|---|
|  | $S_A$ | $S_U$ | $S_A$ | $S_U$ | $S_A$ | $S_U$ |
| CLIP | **0.28** | -0.19 | **0.28** | -0.18 | **0.33** | -0.19 |
| RANKCLIP | 0.23 | -0.14 | 0.26 | -0.13 | 0.29 | **-0.13** |
| RANKCLIP$_C$ | 0.23 | **-0.13** | 0.24 | **-0.12** | 0.29 | -0.14 |
| RANKCLIP$_I$ | 0.25 | -0.15 | 0.28 | -0.14 | 0.32 | -0.17 |

Table 5. Alignment and uniformity scores of CLIP, RANKCLIP. RANKCLIP$_C$ and RANKCLIP$_I$ indicate the solely cross-modal and in-modal rank loss.

## 6.3. Qualitative Examples

**Class activation maps.** To further examine the effects of ranking consistency, we visualize class activation maps
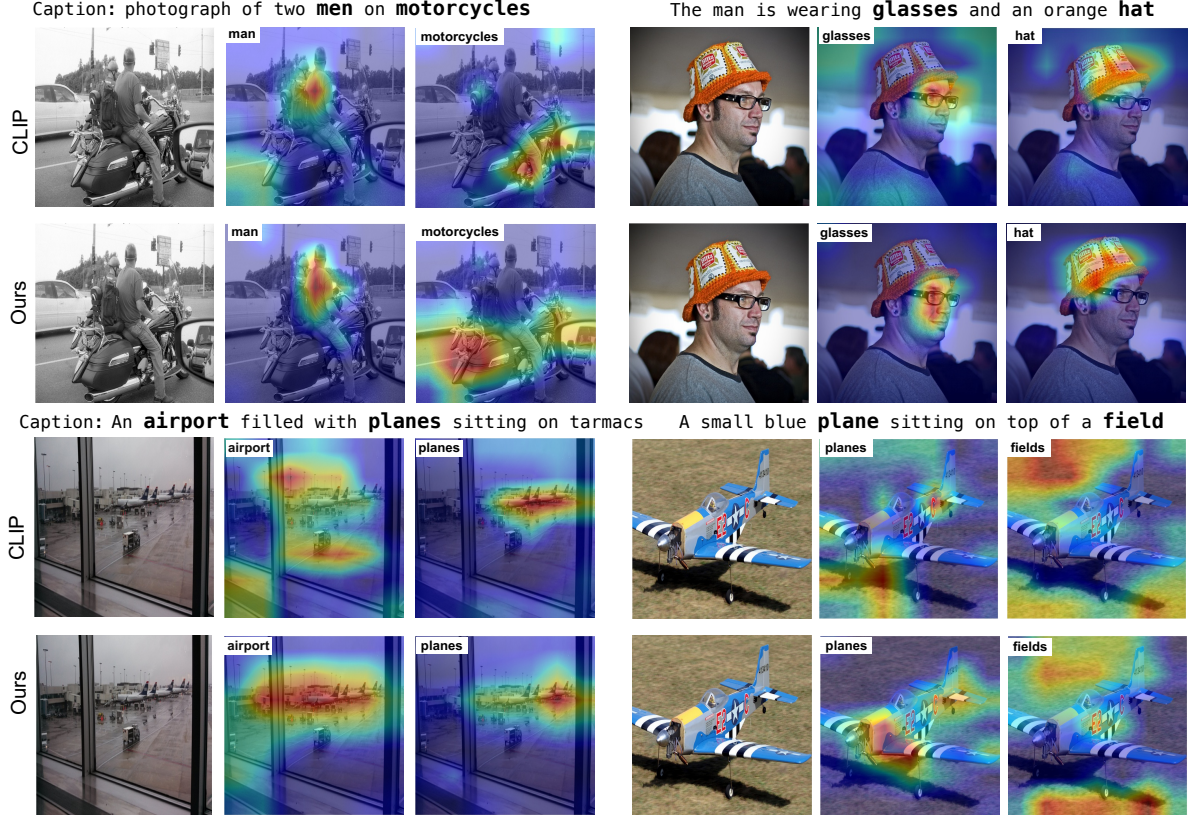
7

Figure 7. Class activation maps for RANKCLIP and CLIP on different objects in the caption from MSCOCO. RANKCLIP has more precise responses to some nouns compared to CLIP and can accurately locate the region related to the noun.

(CAMs) [47] for RANKCLIP and CLIP in Fig. 7. The results show that RANKCLIP consistently attends to more semantically relevant regions in the images. For example, when given the caption 'An airport filled with planes sitting on tarmacs', CLIP mistakenly highlights surrounding areas, whereas RANKCLIP focuses precisely on the planes. Similar improvements are observed across other examples, demonstrating that RANKCLIP better aligns textual descriptions with visual concepts. This suggests that ranking consistency enhances fine-grained feature learning, leading to more localized and accurate visual grounding.

**Text-to-image retrieval.** Fig. 5 compares text-to-image retrieval results from RANKCLIP and CLIP. Given a query, we display the top-ranked images retrieved by each model. RANKCLIP consistently retrieves more semantically aligned images, even beyond the correct reference image. For instance, in the example of 'A cute cat laying down in a sink', correctly identifies a cat in a sink, whereas CLIP misidentifies it due to the visual similarity between sinks and toilets. This demonstrates RANKCLIP 's ability to capture fine-grained semantic dis-

tinctions, reinforcing its advantage in retrieval tasks that demand precise understanding.

## 7. Conclusion

In this paper, we introduce RANKCLIP, a novel language-image pretraining method that integrates ranking consistency into the contrastive learning paradigm. RANKCLIP aims to better understand the complex many-to-many relationships in diverse text-image pairs by optimizing a self-supervised, list-wise rank loss. Through extensive experiments, including zero-shot classification, robustness to distribution shifts, linear probing, and zero-shot image-text retrieval, RANKCLIP not only enhances performance but also improves model robustness and semantic comprehension, outperforming the baseline CLIP and another state-of-the-art model ALIP. Our ablation studies and analyses further demonstrate and interpret the significance of each component of RANKCLIP in boosting performance and understanding across modalities. We believe that the methodologies and principles of RANKCLIP will inspire further research and lead to the development of models with a deeper understanding of the intricate interactions between visual and textual data.

# References

[1] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014. 5

[2] Nicholas Carlini and Andreas Terzis. Poisoning and backdooring contrastive learning. *arXiv preprint arXiv:2106.09667*, 2021. 4

[3] Fei-Long Chen, Du-Zhen Zhang, Ming-Lun Han, Xiu-Yi Chen, Jing Shi, Shuang Xu, and Bo Xu. Vlp: A survey on vision-language pre-training. *Machine Intelligence Research*, 20(1):38–56, 2023. 2

[4] Zhaorun Chen, Yichao Du, Zichen Wen, Yiyang Zhou, Chenhang Cui, Zhenzhen Weng, Haoqin Tu, Chaoqi Wang, Zhengwei Tong, Qinglan Huang, et al. Mj-bench: Is your multimodal reward model really a good judge for text-to-image generation? *arXiv preprint arXiv:2407.04842*, 2024. 1

[5] Zhaorun Chen, Zhuokai Zhao, Hongyin Luo, Huaxiu Yao, Bo Li, and Jiawei Zhou. Halc: Object hallucination reduction via adaptive focal-contrast decoding. *arXiv preprint arXiv:2403.00425*, 2024. 1

[6] Zhaorun Chen, Zhuokai Zhao, Zhihong Zhu, Ruiqi Zhang, Xiang Li, Bhiksha Raj, and Huaxiu Yao. Autoprm: Automating procedural supervision for multi-step reasoning via controllable question decomposition. *arXiv preprint arXiv:2402.11452*, 2024. 1

[7] KR1442 Chowdhary and KR Chowdhary. Natural language processing. *Fundamentals of artificial intelligence*, pages 603–649, 2020. 1

[8] Sanghyuk Chun. Improved probabilistic image-text representations. *arXiv preprint arXiv:2305.18171*, 2023. 2

[9] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014. 5

[10] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011. 5

[11] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. Vqgan-clip: Open domain image generation and editing with natural language guidance. In *European Conference on Computer Vision*, pages 88–105. Springer, 2022. 1

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 4

[13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 12

[14] Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. A survey of vision-language pre-trained models. *arXiv preprint arXiv:2202.10936*, 2022. 2

[15] Yuting Gao, Jinfeng Liu, Zihan Xu, Jun Zhang, Ke Li, Rongrong Ji, and Chunhua Shen. Pyramidclip: Hierarchical feature alignment for vision-language model pretraining. *Advances in neural information processing systems*, 35:35959–35970, 2022. 1, 2

[16] Yuting Gao, Jinfeng Liu, Zihan Xu, Tong Wu, Enwei Zhang, Ke Li, Jie Yang, Wei Liu, and Xing Sun. Softclip: Softer cross-modal alignment makes clip stronger. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1860–1868, 2024. 1, 2

[17] Shijie Geng, Jianbo Yuan, Yu Tian, Yuxiao Chen, and Yongfeng Zhang. Hiclip: Contrastive language-image pretraining with hierarchy-aware attention. *arXiv preprint arXiv:2303.02995*, 2023. 1, 2

[18] Taraneh Ghandi, Hamidreza Pourreza, and Hamidreza Mahyar. Deep learning approaches on image captioning: A review. *ACM Computing Surveys*, 56(3):1–39, 2023. 1

[19] Shashank Goel, Hritik Bansal, Sumit Bhatia, Ryan Rossi, Vishwa Vinay, and Aditya Grover. Cyclip: Cyclic contrastive language-image pretraining. *Advances in Neural Information Processing Systems*, 35:6704–6719, 2022. 1, 2, 4, 5, 7

[20] John Guiver and Edward Snelson. Bayesian inference for plackett-luce ranking models. In *proceedings of the 26th annual international conference on machine learning*, pages 377–384, 2009. 2

[21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 12

[22] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike 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. 1, 5

[23] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15262–15271, 2021. 1

[24] 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, 2021. If you use this software, please cite it as below. 4

[25] Summaira Jabeen, Xi Li, Muhammad Shoib Amin, Omar Bourahla, Songyuan Li, and Abdul Jabbar. A review on methods and applications in multimodal deep learning. *ACM Transactions on Multimedia Computing, Communications and Applications*, 19(2s):1–41, 2023. 1

[26] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International*

*conference on machine learning*, pages 4904–4916. PMLR, 2021. 2

[27] Yogesh Kumar and Pekka Marttinen. Improving medical multi-modal contrastive learning with expert annotations. *arXiv preprint arXiv:2403.10153*, 2024. 6

[28] 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. *arXiv preprint arXiv:2110.05208*, 2021. 1, 2, 4

[29] 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. *Advances in Neural Information Processing Systems*, 35:17612–17625, 2022. 6

[30] 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. 4

[31] Siqu Long, Feiqi Cao, Soyeon Caren Han, and Haiqin Yang. Vision-and-language pretrained models: A survey. *arXiv preprint arXiv:2204.07356*, 2022. 2

[32] R Duncan Luce. *Individual choice behavior: A theoretical analysis*. Courier Corporation, 2005. 2

[33] S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. Technical report, 2013. 5

[34] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018. 6

[35] John P Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar, Percy Liang, Yair Carmon, and Ludwig 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. 6

[36] Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pre-training. In *European conference on computer vision*, pages 529–544. Springer, 2022. 1, 2, 4

[37] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, page 7. Granada, Spain, 2011. 5

[38] Changdae Oh, Junhyuk So, Hoyoon Byun, YongTaek Lim, Minchul Shin, Jong-June Jeon, and Kyungwoo Song. Geodesic multi-modal mixup for robust fine-tuning. *Advances in Neural Information Processing Systems*, 36, 2024. 6

[39] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 12

[40] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012. 5

[41] Robin L Plackett. The analysis of permutations. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 24 (2):193–202, 1975. 2, 3

[42] Farhad Pourpanah, Moloud Abdar, Yuxuan Luo, Xinlei Zhou, Ran Wang, Chee Peng Lim, Xi-Zhao Wang, and QM Jonathan Wu. A review of generalized zero-shot learning methods. *IEEE transactions on pattern analysis and machine intelligence*, 45(4):4051–4070, 2022. 1

[43] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1, 2, 4, 5, 6, 12

[44] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021. 1

[45] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *International conference on machine learning*, pages 5389–5400. PMLR, 2019. 5

[46] 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. 4

[47] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 8

[48] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 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. 4

[49] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650, 2022. 1, 2

[50] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642, 2013. 5

[51] Siddharth Srivastava and Gaurav Sharma. Omnivec: Learning robust representations with cross modal sharing. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1236–1248, 2024. 6

[52] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, (0):–, 2012. 5

[53] Ajinkya Tejankar, Maziar Sanjabi, Bichen Wu, Saining Xie, Madian Khabsa, Hamed Pirsiavash, and Hamed Firooz. A fistful of words: Learning transferable visual models from bag-of-words supervision. *arXiv preprint arXiv:2112.13884*, 2021. 4

[54] Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: the new data in multimedia research. *Commun. ACM*, 59(2):64–73, 2016. 4

[55] Athanasios Voulodimos, Nikolaos Doulamis, Anastasios Doulamis, and Eftychios Protopapadakis. Deep learning for computer vision: A brief review. *Computational intelligence and neuroscience*, 2018, 2018. 1

[56] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International conference on machine learning*, pages 9929–9939. PMLR, 2020. 7

[57] Tan Wang, Kevin Lin, Linjie Li, Chung-Ching Lin, Zhengyuan Yang, Hanwang Zhang, Zicheng Liu, and Lijuan Wang. Equivariant similarity for vision-language foundation models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11998–12008, 2023. 2

[58] Kaicheng Yang, Jiankang Deng, Xiang An, Jiawei Li, Ziyong Feng, Jia Guo, Jing Yang, and Tongliang Liu. Alip: Adaptive language-image pre-training with synthetic caption. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2922–2931, 2023. 1, 2, 4, 5, 12

[59] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li, Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. *arXiv preprint arXiv:2111.07783*, 2021. 1, 2

[60] Zhuokai Zhao, Harish Palani, Tianyi Liu, Lena Evans, and Ruth Toner. Multi-modality guidance network for missing modality inference. *arXiv preprint arXiv:2309.03452*, 2023. 1

# Appendix

## A. Training Procedures

### A.1. Implementation Details

For CLIP [43], we use the official implementation released by OpenAI[1]. And for ALIP [58], we also use the official implementation released by the paper authors[2]. As the proposed RANKCLIP essentially shares the same model architecture (separate vision, text encoders, projection layer, and a classification head) as CLIP, we build upon the CLIP code repository for our model construction[3]. We set the scaling parameters for cross-modal ($\lambda_c$) and in-modal ($\lambda_i$) ranking consistency to 1/16 and 1/16 respectively throughout all the experiments unless otherwise noted. All CLIP, ALIP and RANKCLIP models are initialized from scratch without loading any existing weights. And the embedding sizes for both modalities all project to 1024 across the three models.

### A.2. Training Parameters

Following CLIP [43], we adopt the ResNet-50 [21] and transformer architectures [13] for image and text encoding, respectively. Training is conducted from scratch over 64 epochs using a single NVIDIA A100 GPU, with a batch size of 512, an initial learning rate of 0.0005 employing cosine scheduling, and 10,000 warm-up steps.

### A.3. Training Time Consumption

we conducted the experiments using the same hardware specifications. The table below shows the time consumption for training our RankCLIP and CLIP models with 50K samples from CC3M using a single NVIDIA A100 GPU.

Table 6. Training Details

|  | Time consumption | Dataset size | epochs | batch_size | model_name |
|---|---|---|---|---|---|
| CLIP | 1d 2h 54m 48s | 50K | 64 | 512 | RN50 |
| RANKCLIP | 1d 1h 4m 23s | 50K | 64 | 512 | RN50 |

As shown in the table, the difference in time consumption is negligible. Interestingly, our method is slightly faster than CLIP, but we think it may be attributed to hardware optimizations or variance.

## B. CLIP Preliminaries

CLIP [43] has been a prominent method for learning detailed multimodal representations through the alignment of images and texts. Given a set $\mathcal{D} = \{(V_j, T_j)\}_{j=1}^N$ of $N$ image-text pairs, where $V_j$ denotes an image and $T_j$ is the corresponding text, the goal is to learn representations that map semantically similar images and texts closer in the embedding space, while dissimilar pairs are distanced apart. More specifically, the foundational CLIP model employs two encoders: an image encoder $f_I : \mathcal{I} \to \mathbb{R}^m$ that processes raw images into visual embeddings and a text encoder $f_T : \mathcal{T} \to \mathbb{R}^n$ which encodes textual data into text embeddings. Then both the text and visual features are projected to a latent space with identical dimension. Formally, the embeddings for a text-image pair $(V_j, T_j)$ are denoted as $v_k = f_I(V_j)$ and $t_j = f_T(T_j)$, respectively. The embeddings are then normalized to lie on an unit hypersphere by enforcing $l_2$-norm constraint:

$$\hat{v}_j = \frac{v_j}{\|v_j\|_2}, \quad \hat{t}_j = \frac{t_j}{\|t_j\|_2}. \tag{11}$$

so that the magnitude information is erased and only direction is preserved.

To align the image and text representations, a contrastive loss function, typically a variant of the InfoNCE loss [39], which optimizes the similarity of the matched pair against unmatched pairs, is utilized, i.e.:

$$\mathcal{L}_{\text{CLIP}} = -\frac{1}{2N} \sum_{j=1}^N \left[ \underbrace{\log \frac{\exp(\hat{v}_j^\top \hat{t}_j / \tau)}{\sum_{k=1}^N \exp(\hat{v}_j^\top \hat{t}_k / \tau)}}_{\textcircled{1}} + \underbrace{\log \frac{\exp(\hat{t}_j^\top \hat{v}_j / \tau)}{\sum_{k=1}^N \exp(\hat{t}_j^\top \hat{v}_k / \tau)}}_{\textcircled{2}} \right] \tag{12}$$

---

[1]CLIP repository on GitHub: https://github.com/openai/CLIP.
[2]ALIP repository on GitHub: https://github.com/deepglint/ALIP.
[3]RANKCLIP repository will be released upon acceptance.

| Data Size | Method | Model Type | CIFAR-10 | CIFAR-100 | DTD | FGVGAircraft | Food101 | GTSRB | OxfordPets | SST2 | STL10 | SVHN |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 3m | CLIP | RN50 | 80.12% | 58.50% | 57.18% | 39.75% | 59.14% | 72.41% | 61.73% | 54.48% | 86.01% | 58.92% |
| | RANKCLIP | RN50 | 78.29% | 56.24% | 57.82% | 39.30% | 58.63% | 74.13% | 64.35% | 55.02% | 86.69% | 60.68% |
| | CLIP | ViT-B/32 | 77.60% | 56.15% | 43.19% | 22.59% | 39.72% | 62.05% | 40.39% | 50.96% | 78.99% | 50.53% |
| | RANKCLIP | ViT-B/32 | 78.42% | 56.64% | 42.39% | 23.43% | 40.19% | 60.63% | 40.56% | 53.32% | 79.60% | 47.72% |
| 15m | CLIP | RN50 | 78.81% | 56.32% | 61.49% | 25.83% | 61.64% | 68.76% | 60.37% | 55.57% | 89.82% | 47.99% |
| | RANKCLIP | RN50 | 83.27% | 62.96% | 65.96% | 32.19% | 68.11% | 74.25% | 67.40% | 56.34% | 94.20% | 53.06% |
| | CLIP | ViT-B/32 | 82.97% | 62.55% | 49.47% | 24.48% | 52.46% | 63.55% | 50.78% | 52.66% | 87.14% | 46.38% |
| | RANKCLIP | ViT-B/32 | 82.79% | 59.89% | 52.50% | 23.94% | 56.44% | 61.58% | 52.98% | 53.60% | 89.01% | 42.16% |

Table 7. Linear probing accuracy on 10 downstream datasets.

where the first term ①contrasts images with the texts, the second term ②contrasts texts with the images, and $\tau$ denotes a temperature scaling parameter that adjusts the concentration of the distribution. The optimization of Eqn. (12) results in embeddings where the cosine similarity between matched image-text pairs is maximized in comparison to unmatched pairs, thus achieving the desired alignment in the joint embedding space.

Despite the efficacy of CLIP in learning correlated multimodal embeddings, it inherently relies on strict pairwise matched comparisons and fails to capture the more complex, fine-grained nature of semantic similarity within and across modalities that are generally treated as unmatched. This observation motivates the development of RANKCLIP, which innovates beyond binary pairwise contrasts to consider holistic listwise consistency within and across modalities.

## B.1. Additional Experiments

We conduct the linear probing experiment under different training datasize from 3m to 15m as shown in 7.

## B.2. Pseudo-code

---

**Algorithm 1** Pseudo-code of RANKCLIP loss in a Python-like style.

---

```python
# emb_pred: predictions from the model, shape [embs_length, embs_length]
# emb_true: ground truth labels, shape [embs_length, embs_length]

def rank_loss(emb_pred, emb_true):
    # Shuffle for randomised tie resolution
    emb_pred_shuff = emb_pred[:, random_indices]
    emb_true_shuff = emb_true[:, random_indices]
    # Record the rank label index
    emb_true_sorted, indices = emb_true_shuff.sort(descending=True, dim=-1)
    # Ranking the pred embedding by the true indices
    preds_sorted = gather(emb_pred_shuff, dim=1, index=indices)
    # Implementation of the Eq.1, Eq.2 and Eq.3
    max_pred_values, _ = preds_sorted.max(dim=1, keepdim=True)
    preds_sorted_minus_max = preds_sorted - max_pred_values
    cumsums = cumsum(preds_sorted_minus_max.exp().flip(dims=[1]), dim=1).flip(dims=[1])
    loss = (log(cumsums) - preds_sorted_minus_max) * scale_factor
    return mean(sum(loss, dim=1))


# Cross-modal embeddings
logits_text_per_image=image_embeds @ text_embeds.T
logits_iamge_per_text=logits_text_per_image.T
# In-modal embeddings
logits_image_per_image=image_embeds @ image_embeds.T
logits_text_per_text=text_embeds @ text_embeds.T
# Compute the cross-modal rank loss
Cross_modal_loss=rank_loss(logits_text_per_image,logits_image_per_text)+rank_loss(logits_image_per_text
    , logits_text_per_image)
# Compute the in-modal rank loss
In_modal_loss=rank_loss(logits_image_per_image,logits_text_per_text)+rank_loss(logits_text_per_text,
    logits_image_per_image)
# Rank loss
Rank_loss=Contrastive_loss+Cross_modal_loss+In_modal_loss
```

---