## ONE-STEP NOISY LABEL MITIGATION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Mitigating the detrimental effects of noisy labels on the training process has become increasingly critical, as obtaining entirely clean or human-annotated samples for large-scale pre-training tasks is often impractical. Nonetheless, existing noise mitigation methods often encounter limitations in practical applications due to their task-specific design, model dependency, and significant computational overhead. In this work, we exploit the properties of high-dimensional orthogonality to identify a robust and effective boundary in cone space for separating clean and noisy samples. Building on this, we propose One-step Anti-Noise (OSA), a model-agnostic noisy label mitigation paradigm that employs an estimator model and a scoring function to assess the noise level of input pairs through just one-step inference, a cost-efficient process. We empirically demonstrate the superiority of OSA, highlighting its enhanced training robustness, improved task transferability, ease of deployment, and reduced computational costs across various benchmarks, models, and tasks. Our code is released at https://anonymous.4open.science/r/CLIP\_OSN-E86C.

023

000

001 002 003

004

006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

024 025

026

#### 1 INTRODUCTION

Noise mitigation aims to handle the detriment of noisy labels encountered during the training process. The advancement of large-scale pre-training has significantly increased data scale to the trillion level. Much of this data is sourced from the internet, inevitably introducing considerable noise, which severely impedes the training process. This poses a substantial challenge for robust model training in various tasks, such as cross-modal matching (Huang et al., 2021; Zhang et al., 2024), image-classification (Sun et al., 2021; Yu et al., 2019), and image-retrieval (Liu et al., 2021).

Traditional noise mitigation approaches encounter several limitations that constrain their practical applicability: 1) **Task specificity:** Existing methods (Huang et al., 2021; Sun et al., 2021; Ibrahimi et al., 2022a) are tailored to specific tasks, limiting their applicability across different tasks. 2) **Model dependency:** Most noise mitigation techniques (Liu et al., 2021; Yang et al., 2023a) are tightly coupled with specific models, requiring extensive modifications for adaptation to different models. 3) **Computational cost:** Numerous existing methods necessitate dual-model collaborations (Huang et al., 2021; Yu et al., 2019) or multiple training passes (Huang et al., 2021), *i.e.*, they require at least two backward passes per training step, effectively doubling the computational expense and substantially increasing the training burden (see Figure. 1a).

042 To tackle these challenges, we use an external estimator to assess the noise level of each sample, 043 ensuring a model-agnostic approach. This estimator adjusts the training loss by reducing the in-044 fluence of noisy samples, driving their weights toward zero. Furthermore, multimodal pre-trained models have demonstrated remarkable task transferability due to their strong semantic capabilities. For instance, CLIP (Radford et al., 2021) unifies the paradigms of image-text retrieval and image 046 classification through a shared embedding space (see Figure. 1b). It converts category labels into 047 sentences, maps them into the shared embedding space, and then calculates the cosine similarity with 048 the image representation to perform image classification. Inspired by this, we leverage multimodal pre-trained models as estimators and apply the shared embedding space to enable task transfer. In this case, only one additional inference process is required for each sample, significantly reducing 051 the computational overhead compared to performing an extra backward pass. 052

053 Nonetheless, this paradigm introduces a new challenge: how to accurately identify noise based solely on cosine similarity scores generated by estimators. An ideal solution is to find a decision

**Rectify** data

Clarify data

Warm up

066

067 068

069

071

072

054

055



Image retrieval

Co-training

Figure 1: (a) The current anti-noise paradigm with multiple backward significantly enhances the training overhead. (b) CLIP unifies the framework of image-text matching and image classification through a shared space. (c-f) Cosine similarity distribution of noise and clean data with 50% noise.

073 boundary that separates clean samples from noise and accurately handles overlapping samples near 074 the boundary. Existing methods (Huang et al., 2021; Qin et al., 2022; Li et al., 2020; Zhang et al., 075 2024) typically attempt to build this boundary within the loss space, an isotropic space with uniform 076 distribution, which creates only a narrow gap between noisy and clean samples. Moreover, the coarse 077 handling of overlaps by integrating multi-model predictions often results in an unstable decision boundary. In contrast, the shared embedding space of pre-trained models is a high-dimensional, 079 anisotropic space with an imbalanced distribution. Thus, a consideration is whether the properties of imbalanced anisotropic space can help to identify a more precise and robust decision boundary.

081 In this work, we delve into the decision boundary of pre-trained models employed as estimators 082 to accurately differentiate between clean and noisy samples. We first investigate the cosine simi-083 larity distributions of clean and noisy samples, calculated using the multimodal pre-trained models 084 CLIP (Radford et al., 2021) and ALIGN (Jia et al., 2021), on two datasets with a 50% noise ratio: 085 MS-COCO (Lin et al., 2014) and SDM, as shown in Figure. 1c-1f. SDM is a dataset with the images generated by the Stable Diffusion Model (SDM) (Rombach et al., 2022) in some uncommon styles (see illustrations in Figure. 4). It is designed to explore how well-pre-trained models can 087 distinguish unfamiliar domains that they rarely encounter during training. There are two interesting 880 observations in Figure. 1c-1f: (1) The clean and noisy distributions of the same model on different 089 datasets have a similar intersection point, suggesting the existence of a natural and stable boundary 090 in distinguishing between clean and noisy samples. (2) The overlaps between the clean distribution 091 and noisy distribution are minimal even in the unfamiliar domain dataset, indicating this boundary 092 has strong potential for distinguishing between clean and noisy samples. 093

Building upon these two observations, we conduct an in-depth investigation and make the following 094 contributions: 095

096 1. We figure out the origin of the intersection, attributing it to the shift in the orthogonal boundary induced by the cone effect. Furthermore, we provide a theoretical framework that proves and 098 elaborates the stability and precision of this boundary in separating noisy and clean samples.

099 2. We provide a detailed explanation of the reliability of pre-trained models in general noise recog-100 nition, even in unfamiliar domains, grounded in the analysis of the pre-training process. 101

3. Build on this, we introduce One-Step Anti-Noise (OSA), a general model-agnostic paradigm for 102 noise recognition that requires only one-step inference. Specifically, we utilize a pre-trained model 103 as the estimator to maintain a shared embedding space. A scoring function, designed based on the 104 properties of high-dimensional orthogonality, is then used to accurately handle overlaps by directly 105 assigning a learning weight to each sample's loss according to its cosine similarity. 106

4. We conduct comprehensive experiments across a variety of challenging benchmarks, models, and 107 tasks, demonstrating the effectiveness, generalization capabilities, and efficiency of our method.

# 108 2 BOUNDARY PRINCIPLE ANALYSIS

In Figure. 1c-1f, we observe a natural boundary emerging in the pre-trained model's ability to distinguish between clean and noisy samples. In this section, we explain the principle of boundary forming from high-dimensional perspectives, and how robust it is in general noise mitigation.

2.1 HYPOTHESIS: INTERACTION BOUNDARY IS SHIFTED FROM ORTHOGONAL BOUNDARY

We first elaborate on the gap extent between the positive and negative sides kept by the orthogonal boundary. Then, we present the reasoning behind the hypothesis that the intersection boundary in Figure. 1 is a shifted orthogonal boundary in the cone space.

119

127

110

111

112

113 114

115

**The orthogonal boundary largely separates the positive and negative sides.** High-dimensional orthogonality is a general phenomenon caused by dimension disaster, where the angles between randomly selected vectors typically approximate 90 degrees, suggesting the cosine similarity that trends toward zero. For instance, in a 1024-dimensional space, the probability of two random vectors having a cosine similarity within [-0.1, 0.1] is approximately 99.86% (details are provided in Appendix. C.1). In this case, a natural boundary of cosine similarity zero forms, capably separating the positive side and negative side with a huge gap.

### 128 Cone effect may induce orthogonal boundary shift. Recent

literature (Liang et al., 2022a; Bogolin et al., 2022; Ethayarajh, 129 2019) has demonstrated that the cone effect is a general phe-130 nomenon in deep neural networks, where the learned embedding 131 subspace forms a narrow cone and the orthogonal boundary en-132 counters a positive shift. Based on this, a hypothesis is that the 133 interaction boundary in Figure. 1 is the shifted orthogonal bound-134 ary. To prove this, we simulate the process of selecting random 135 vectors in high-dimensional space and randomly generate thousands of pairs mapped into the shared embedding space. We find 136 that all similarity of these random vector pairs tends to a fixed 137

Table 1: The mean and variance of cosine similarity between randomly generated pairs.

Model	Mean	Var
CLIP	0.215	0.024
ALIGN	0.087	6e-4

value, with the low-variance cosine similarity almost lying in the middle of clean and noise distributions (see Table. 1). An interesting phenomenon is that if we compare the mean with the intersection points in Figure. 1c-1f, we find they are almost exactly the same. This suggests that the interaction boundary is highly likely to be a shifted orthogonal boundary in cone space.

141

142 143

150

156

157 158

#### 2.2 THEORETICAL VERIFICATION OF THE INTERACTION BOUNDARY ORIGIN

Here, we theoretically investigate whether the origin of the interaction boundary is a shifted orthogonal boundary. We first show that (i) contrastive learning separates clean and noisy samples on
opposite sides of the orthogonal boundary and (ii) The relative relationships of pairs' cosine similarity stays unchanged after transmitting into the narrow cone space. Based on (i) and (ii), we can
confirm that the intersection boundary at the center of the clean and noisy distributions is the shifted
orthogonal boundary.

151 **Contrastive learning empowers the separation of clean and noisy samples.** For an initialized 152 model intending to learn an embedding space, both clean and noisy samples are treated as orthogonal 153 random vectors since lacking semantic perception ability in the initial space. During contrastive 154 training process, given N sample pairs  $\{(x_i, y_i)\}_{i=1}^N$ , the embedding space is optimized through the 155 cross-entropy loss (Eq. 1).

$$\mathcal{L}_{ce} = \frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(m_{ii})}{\sum_{j=1}^{N} \exp(m_{ij})},$$
(1)

159 where  $M \in \mathbb{R}^{N \times N}$  represents the cosine similarity matrix of N sample pairs during training pro-160 cess. Each element  $m_{ij} \in M$  denote the cosine similarity between  $x_i$  and  $y_j$ . The diagonal elements 161  $m_{ii}$  denote the cosine similarities of positive pairs, while the non-diagonal elements  $m_{ij}$  represent the cosine similarities of negative pairs. To minimize  $\mathcal{L}_{ce}$  during training, two subprocesses occur: the diagonal elements of the matrix (*i.e.*, clean pairs) are optimized to the positive side of the orthogonal boundary, while the nondiagonal elements (equivalent to noise pairs) are optimized to the negative side. Consequently, the distributions of these two types of samples are on opposite sides of the orthogonal boundary.

166

Relative relationship unchanged in transmitting process. We study how the boundary shifts
from the entire space to the narrow cone in the neural network. The following theorem shows that
the cosine similarity will be proportionally scaled to the target narrow cone, while still maintaining
a boundary with properties similar to the orthogonal boundary. In other words, vectors with cosine
similarity smaller than the orthogonal boundary in the original space remain smaller than the shifted
boundary in the narrow cone space, while those larger remain larger.

**Theorem 1** (Proportional shift of boundary). Let  $\mathbb{R}^{d_{in}}$  be the original space before being transmitted in a neural network. Suppose  $u, v \in \mathbb{R}^{d_{in}}$  are any two random vectors with  $\cos(u, v) \approx 0$ .  $u_c, v_c \in \mathbb{R}^{d_{in}}$  is a pair of clean vectors with  $\cos(u_c, v_c) > 0$ , while  $u_n, v_n \in \mathbb{R}^{d_{in}}$  is a noisy pair with  $\cos(u_n, v_n) < 0$ . Given a Neural Network  $F(x) = f_t(f_{t-1}(\dots f_2(f_1(x)))) \in \mathbb{R}^{d_{out}}$ with t layers.  $f_i(x) = \sigma_i(\mathbf{W}_i x + \mathbf{b}_i)$  denotes  $i^{th}$  layer, where  $\sigma(\cdot)$  indicates activation function.  $\mathbf{W}_i \in \mathbb{R}^{d_{out}^i \times d_{in}^i}$  is a random weight matrix where each element  $\mathbf{W}_i^{k,l} \sim \mathcal{N}(0, 1/d_{out}^i)$  for  $k \in [d_{out}^i]$ ,  $l \in [d_{in}^i]$ , and  $\mathbf{b}_i \in \mathbb{R}^{d_{out}^i}$  is a random bias vector such that  $\mathbf{b}_i^k \sim \mathcal{N}(0, 1/d_{out}^i)$  for  $k \in [d_{out}^i]$ . Then, there always be a boundary  $\beta$ , satisfying:

$$\cos(F(u_n), F(v_n)) < \cos(F(u), F(v)) \approx \beta < \cos(F(u_c), F(v_c)).$$
(2)

Theorem. 1 shows that the relative relationship of pairs in the original entire space, will not change after transmitting to the narrow cone space of the trained model, and there is always a boundary  $\beta$ concentrated on most random vectors. In Appendix. C.2, we provide a detailed statement and proof of the Theorem.

187 188

189

181 182

#### 2.3 QUALITATIVE ANALYSIS OF ROBUSTNESS AND APPLICABILITY

Next, we perform a qualitative analysis to explore (i) the robustness and generality of the boundary
 in distinguishing between clean and noisy samples, and (ii) how the boundary's properties can be
 leveraged to achieve more reasonable and precise overlap handling.

193

How about the boundary robustness even in unfamiliar domains? Although the boundary's 194 ability to distinguish clean and noisy samples is proven, its robustness and generality still require 195 further exploration. For practical pre-training, it must maintain accuracy and robustness even in 196 unfamiliar domain datasets. Since the capabilities of the pre-trained model are difficult to quantify, 197 we conduct a qualitative analysis from the perspective of pre-trained model inference. The models pre-trained on millions of samples already possess somewhat semantic understanding capabilities. 199 Given a positive pair from an unseen domain, due to the contrastive learning process during pre-200 training, it still has a strong likelihood of moving toward the positive side of the boundary, while 201 the negative pair tends toward the negative side. Although the cosine similarity difference might 202 be slight, as we have shown in Section. 2.1, the boundary constructs a significant gap from the perspective of high-dimensional orthogonality. 203

204

How to handle the overlaps through imbalanced probability? Due to the properties of orthog onal boundary, as cosine similarity decreases and approaches zero from the positive side, the probability of positive samples sharply decreases. Therefore, we can design a scoring function to annotate
 the cleanliness of samples. This function should satisfy two requirements: for samples with cosine
 similarity less than or equal to zero, which are almost certainly noise, the function should assign
 them a weight of zero. For samples with cosine similarity greater than zero, the function gradient
 should increase rapidly as the cosine similarity moves further from zero.

### 3 Method

213 214

212

In this section, we present our One-Step Anti-Noise (OSA) paradigm with a workflow shown in Figure. 2. We first define the pair-based noise mitigation tasks for image-text matching, image



Figure 2: The workflow of OSA. In the anti-noise process, there are two phases: Scoring Phase and Training Phase. In the Scoring Phase, a pair is mapped to a shared embedding space by estimators. Then the cosine similarity is transformed to a weight w by a scoring function. In the Training Phase, the weight w is directly multiplied with the loss to instruct the optimization.

classification, and image retrieval tasks in Sec. 3.1. Consequently, the detailed description of OSA is clarified in Sec. 3.2.

#### 3.1 TASK DEFINITION

228

229

230

231 232 233

234

235 236

237

248

249

250 251

252

238 Let  $\mathcal{D} = \{(x_i, y_i, c_i)\}_{i=1}^N$  denote a paired dataset, where  $(x_i, y_i)$  represents the *i*-th pair in the 239 dataset, and  $c_i$  indicates a noise label for that pair. Specifically, when  $c_i = 0$ ,  $(x_i, y_i)$  forms a 240 correct (paired) match, while  $c_i = 1$  denotes an incorrect (unpaired) match. The objective of noise 241 mitigation in contrastive learning is to construct a shared embedding space that brings  $x_i$  and  $y_i$ closer when  $c_i = 1$ . In different tasks,  $x_i$  and  $y_i$  are distinct data types. For instance, in the image-242 text retrieval task,  $x_i$  and  $y_i$  represent images and texts, respectively. In the image classification 243 task,  $x_i$  and  $y_i$  represent images and categories, respectively. In the image retrieval task,  $x_i$  and 244  $y_i$  represent images and relevant images, respectively. The paired sample (x, y) could be encoded 245 into a shared embedding space by corresponding encoders  $\phi_x(\cdot)$  and  $\phi_y(\cdot)$ . Afterward, the cosine 246 similarity s(x, y) is calculated through Eq. 3 as semantic relevance of (x, y) to guide the training. 247

$$s(x,y) = \frac{\phi_x(x)}{\|\phi_x(x)\|} \cdot \frac{\phi_y(y)}{\|\phi_y(y)\|}.$$
(3)

#### 3.2 ONE-STEP ANTI-NOISE

The workflow of our noise mitigation approach OSA is depicted in Figure. 2. Initially, we utilize an estimator model to encode the input pair to a shared embedding space and continue to compute the cosine similarity between the paired embedding. Afterward, the cosine similarity is converted to a cleanliness score  $w_i$ ,  $(0 \le w_i \le 1)$  through a scoring function designed based on orthogonal properties (Section. 2.3). This score quantifies the clean degree of the sample, the smaller  $w_i$  is, the noisier the sample.

During the target model training phase, this cleanliness score is used as a weight, directly multiplied
 by the loss of the corresponding sample to facilitate selective learning. This noise mitigation process, being solely dependent on the estimator model, is readily adaptable to the training of various
 target models by simply adding an extra coefficient to the loss function, ensuring the model-agnostic
 property. Therefore, the key of our noise mitigation approach revolves around the estimator model and noise score assessment.

265 3.2.1 ESTIMATOR MODEL 266

Estimator model selection. In our approach, the Estimator Model must satisfy two critical re quirements: 1) effectively mapping input pairs into a unified embedding space and 2) possessing
 basic semantic understanding capabilities. To meet these requirements, we employ CLIP (Rad ford et al., 2021), a commonly used multimodal pre-trained models, as our estimator model. It is

equipped with a text encoder  $\phi_t(\cdot)$  and an image encoder  $\phi_v(\cdot)$ , enabling it to perform basic zeroshot tasks efficiently.

273 **Domain adaptation (Optional).** While we have performed a qualitative analysis of the zero-shot 274 pre-trained model's robustness on out-of-domain data in Section. 2.3, and shown strong robustness 275 for edge cases in Figure. 1, considering the domain diversity in real-world scenarios, we provide an optional Domain Adaptation (DA) approach to enhance the estimator model's adaptability when 276 encountering edge domains. Following NPC (Zhang et al., 2024), we first employ a Gaussian Mix-277 ture Model (GMM) coupled with strict selection thresholds to ensure the absolute cleanliness of the 278 chosen samples. We afterward implement a warm-up phase with few steps, allowing the estimator 279 model to better understand the semantics of the target domain. Notably, this trick is only optional 280 for our methods. Through multiple experiments, we found that even without domain adaptation, the 281 zero-shot CLIP model performs exceptionally well across various scenarios. 282

# 283 3.2.2 Noise Score Assessment

285 Spatial Debiasing. The cone effect phenomenon has been demonstrated as a general phenomenon for deep neural networks, typically resulting in a narrow embedding space that causes a shift of 286 space center to a narrow cone center (Liang et al., 2022a). Specifically, when paired randomly 287 generated inputs are mapped into a shared embedding space through model encoders, the resultant 288 vectors exhibit an average cosine similarity that deviates from zero and tends to another fixed angle. 289 To counteract this shift and mitigate its impact on the estimator's ability to accurately recognize 290 noises through high-dimensional orthogonality, a random sampling method is developed. We begin 291 by constructing K random sample pairs  $\mathcal{R} = \{(x_i, y_i) \mid i = 1, 2, \dots, K\}$  and processing them 292 through the estimator's encoder to generate a set of vectors. Then the average cosine similarity 293 among these vectors will be calculated as the space shift  $\beta$  through:

$$\beta = \frac{\sum_{j=1}^{K} s(x_j, y_j)}{K}.$$
(4)

**Scoring Function.** After spatial debiasing, we employ a scoring function  $w(\cdot)$  to evaluate the 298 cleanliness of the input pair (x, y). In section. 2.3, we have elaborate how to handle overlaps based 299 on the orthogonal boundary property. For an estimator model trained on millions of samples using 300 contrastive learning, clean pairs (diagonal elements) are optimized to positive side, while noise pairs 301 (non-diagonal elements) are optimized to negative side. Given unfamiliar pairs, the model also 302 tends to map clean pairs towards positive and noisy pairs towards negative. Despite the potentially 303 slight similarity difference between clean and noisy pairs, high-dimensional orthogonality ensures 304 a substantial gap between them. In this case, a negative cosine similarity s(x, y) computed by 305 the estimator, indicating the pair is almost certainly noise, should be assigned a score of zero. For 306 samples with s(x, y) greater than zero, the probability of the sample being positive sharply decreases 307 as the cosine similarity approaches zero from the positive side. Therefore, the function gradient 308 should increase rapidly as the cosine similarity moves further from zero. To systematically score the noise, we design the scoring function as: 309

312 313

314

315

316

295

296 297

$$w(x,y,\beta) = \begin{cases} 0 & ,s(x,y) - \beta \le 0\\ -(s(x,y) - \beta)^2(s(x,y) - \beta - 1) & ,otherwise \end{cases}$$
(5)

**Re-weight Training.** After scoring, the target model can selectively learn from the samples by reweighting the loss. Noise samples with smaller weights will have a reduced impact on model updates and will be effectively mitigated. For a sample (x, y), let  $\mathcal{L}_{x,y}$  denote its loss, the re-computed loss  $\mathcal{L}_{re}$  is defined as:

$$\mathcal{L}_{re} = w(x, y, \beta) \times \mathcal{L}_{x, y}.$$
(6)

317 318 319

#### 4 EXPERIMENTS

320 321

> In this section, we present experiments on multiple datasets with label noise, demonstrating the effectiveness of our methods. Firstly, we describe the datasets, metrics, and implementation details. Then, we report our results on several downstream tasks. Lastly, we conduct ablation studies to

show how each part of our method contributes and examine how these parts interact. The literature involved in our experiments and richer related work are detailed in Appendix. B.

### 4.1 EVALUATION SETTING

In this section, we briefly introduce the datasets and evaluation metrics used in the experiments. For more dataset and implementation details, please refer to Appendix. A.

**Datasets.** We evaluate our method on three downstream tasks with noisy labels, including one multimodal task and two visual tasks. For the cross-modal matching task, we perform experiments on the <u>MSCOCO</u> (Lin et al., 2014) and <u>Flickr30K</u> (Young et al., 2014) datasets. Following NPC (Zhang et al., 2024), we further carry out evaluations on a real-world noisy dataset <u>CC120K</u>. For image classification tasks, experiments are conducted under three subsets of <u>WebFG-496</u> (Sun et al., 2021)—Aircraft, Bird, and Car. For image retrieval tasks, we conduct experiments on the <u>CARS98N</u> dataset under PRISM (Liu et al., 2021) setting.

**Evaluation Metrics.** For the image-text matching task, the recall value of the top-K retrieved results (R@K) is used. For classification tasks, accuracy serves as the evaluation metric. For the image retrieval task, we use Precision@1 and mAP@R for evaluation.

4.2 COMPARISONS WITH STATE OF THE ARTS

Table 2: Comparison on noisy MS-COCO.

$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	2t         t2i           25         R@10         R@1         R@5         R@10           (1.1)         93.3         48.2         76.7         85.5           (3.8)         93.3         48.1         76.7         85.5           (2.2)         46.8         74.4         83.7           (3.3)         93.1         48.5         75.4         84.4
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	@5         R@10         R@1         R@5         R@10           1.1         93.3         48.2         76.7         85.5           5.8         93.3         48.1         76.7         85.5           5.2         92.2         46.8         74.4         83.7           7.3         93.1         48.5         75.4         84.4
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	93.3         48.2         76.7         85.5           8         93.3         48.1         76.7         85.5           5.2         92.2         46.8         74.4         83.7           7.3         93.1         48.5         75.4         84.4
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	5.8       93.3       48.1       76.7       85.5         5.2       92.2       46.8       74.4       83.7         7.3       93.1       48.5       75.4       84.4
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	5.2         92.2         46.8         74.4         83.7           7.3         93.1         48.5         75.4         84.4
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	<b>3</b> 93.1   48.5 75.4 84.4
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	.9 91.6 46.5 73.8 82.9
ALIGN +OSA         84.9         97.3         99.0         70.5         92.8         97.2         69.6         89.           VSE         85.3         97.4         99.0         71.4         93.1         97.3         69.8         89.           VSE         78.4         94.3         97.0         65.5         89.3         94.1         58.6         83.           PCME++         78.4         95.9         98.4         64.9         90.8         96.1         57.7         83.           PAU         78.2         95.2         98.1         64.5         90.0         95.4         59.3         94.1         58.6         83.           POME++         78.4         95.9         98.4         64.5         90.0         95.4         59.3         82.           NPC         79.9         95.9         98.4         66.3         90.8         98.4         61.6         85.           CLIP         76.0         94.3         97.5         63.4         89.0         94.8         55.3         79.9           +OSA         81.6         96.2         98.5         68.9         92.0         96.6         65.8         86.	5.8         92.9         49.1         76.2         84.8
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	94.5   50.5 77.5 85.7
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	0.9 94.8 51.4 78.2 86.3
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	6.4 89.9   45.0 72.9 81.7
20%         PAU NPC         78.2 79.9         95.2 95.9         98.1 66.3         64.5 90.0         95.4 98.4         59.3 61.6         82. 85.           CLIP +OSA         76.0         94.3 81.6         97.5 98.5         63.4 68.9         92.0         94.8 96.6         55.3 65.8         79. 85.6	<b>6.9</b> 91.0   43.2 72.3 82.4
20%         NPC         79.9         95.9         98.4         66.3         90.8         98.4         61.6         85.           CLIP         76.0         94.3         97.5         63.4         89.0         94.8         55.3         79.           +OSA         81.6         96.2         98.5         68.9         92.0         96.6         65.8         86.	2.9 90.4 44.2 71.3 81.3
CLIP 76.0 94.3 97.5 63.4 89.0 94.8 55.3 79. +OSA 81.6 96.2 98.5 68.9 92.0 96.6 65.8 86.	6.4 91.6   46.0 73.4 82.9
+OSA 81.6 96.2 98.5 68.9 92.0 96.6 65.8 86.	0.1 86.9   41.0 68.8 79.3
	<b>5.4 92.5   48.7 76.1 84.5</b>
ALIGN   79.4 95.7 98.2   66.2 90.8 96.1   60.9 84.	.5 91.0   46.3 73.6 82.3
+OSA 85.1 97.4 99.1 70.9 93.0 97.3 69.7 90.	0.0 94.7   50.9 77.8 86.2
$ VSE_{\infty} $   44.3 76.1 86.9   34.0 69.2 84.5   22.4 48.	8.2 61.1   15.8 38.8 52.1
PCME++ 74.8 94.3 97.7 60.4 88.7 95.0 52.5 79.	0.6 88.4 38.6 68.0 79.0
PAU 76.4 94.1 97.6 62.3 88.5 94.6 57.3 81.	.5 88.8 41.9 69.4 79.6
50% NPC   78.2 94.4 97.7   63.1 89.0 <b>97.7</b>   59.9 82.	2.9 89.7   43.0 70.2 80.0
CLIP   73.9 93.0 97.2   60.1 87.3 94.0   54.1 78.	8.5 86.6   39.7 67.2 77.5
+OSA   80.4 96.2 98.6   67.8 91.6 96.4   64.0 85.	5.5 91.9   47.9 74.6 83.8
ALIGN   78.0 95.8 98.5   65.4 90.3 96.0   60.1 84.	.3 91.2 45.2 72.8 82.1
+OSA 84.3 97.0 98.9 70.0 92.5 97.0 68.5 89.	0 0 0 1 2 50 0 77 0 85 4

**Results on MSCOCO.** To fairly demonstrate the effectiveness of our method, we compare OSA with various robust learning image-text matching approaches using the same ViT-B/32 CLIP as backbone, including VSE $\infty$  (Chen et al., 2021), PCME++ (Chun, 2023), PAU (Li et al., 2023), NPC (Zhang et al., 2024). Besides, we separately employ OSA on both CLIP (Radford et al., 2021) and ALIGN (Jia et al., 2021). The results in Table. 2 show that OSA outperforms all previous approaches on all metrics with a huge gap. In the more challenging MS-COCO 5K set with 50% noise ratio, OSA surpasses the SOTA method NPC in the R@1 for both image-to-text (i2t) and text-to-image (t2i) matching by 8.6% and 7.0%, respectively. Another phenomenon is that as the noise ratio increases from 0% to 50%, all other methods encounter severe performance drop, with an averaging drop of 5.05% for NPC across four R@1 metrics. In contrast, OSA exhibits only a slight decrease of 1.275%, showcasing the accuracy and robustness of OSA in anti-noise tasks.

381 **Results on Flickr30K.** To further demonstrate the generalization ability of OSA, we evaluate on 382 the Flickr30K dataset and compare with several anti-noise methods, including NCR (Huang et al., 2021), DECL (Qin et al., 2022), BiCro (Yang et al., 2023a), and NPC (Zhang et al., 2024). The 384 results are presented in Table. 8 of Appendix. It is evident that OSA consistently outperforms all 385 models on the R@1 metric. Notably, compared with the baseline CLIP, training with OSA at a 60% 386 noise ratio achieves 20.9% R@1 improvement for i2t and a 22.3% R@1 improvement in t2i, further indicating the effectiveness of OSA on noise mitigation. Additionally, OSA demonstrates similar 387 noise robustness on the Flickr30K dataset as observed on MSCOCO, with only 1.4% R@1 drop on 388 i2t and 1.2% R@1 drop on t2i ranging from 0% noise to 60% noise, while all of the other anti-noise 389 approaches hardly resist the detriment from high-ratio noise. All of these results demonstrate the 390 effectiveness and robustness of OSA on anti-noise tasks. 391

**Results on CC120K.** To further verify the reliability of OSA in real scenarios, we conduct evaluations on a large-scale real-world noisy dataset, CC120K, with 3%-20% noise ratio. The results shown in Table. 3 indicate that OSA outperforms the current state-of-the-art method NPC, even in larger-scale real-world domains. This demonstrates the feasibility and generality of OSA even in practical training scenarios.

Table 3: Comparison on real-world noisy dataset CC120K.

Table 4: Results of other image-based tasks.

	:24				+ <b>?</b> :					Image (	Classifi	cation	Image	Retrieval
Method	R@1	12t R@5	R@10	R@1	121 R@5	R@10	Method	Aircraft	Bird	Car	Prec	mAP		
NDC	71.1	02.0	06.2	72.0	00.5	04.9		Acc	Acc	Acc	1100.			
	/1.1	92.0	90.2	73.0	90.5	<u> </u>	Baseline	65.44	62.29	75.90	71.69	18.16		
	08.8	87.0	92.9	07.8	80.4	90.9	+OSA	73.18	70.50	80.19	78.45	24.99		
+05A	73.1	92.2	95.7	73.9	91.2	94.7								

Results on Other Downstream Tasks. To validate the transferability of OSA across different tasks, we evaluate it on two additional tasks: image classification and image retrieval. The results are presented in Table. 4. The baseline method for both tasks leverages contrastive learning. In the image classification task, OSA outperforms the baseline by 7.74%, 8.21%, and 4.28% on the Aircraft, Bird, and Car subsets, respectively. In the image retrieval task, OSA improves performance by 6.76% in precision and 6.83% in mAP. These improvements demonstrate the strong task transferability and generality of OSA.

414 415

416

417

392

393

394

395

396

397

399

400

401

402 403

404 405 406

4.3 TARGET MODEL-AGNOSTIC ANALYSIS

OSA is an architecture-agnostic paradigm that can be easily adapted to various models. To verify its model-agnostic property, we evaluate it across models with different architectures. Subsequently, we apply it to other anti-noise models to demonstrate its generalization capability in noise mitigation.

418 419 420

Architecture-agnostic Analysis. The effectiveness of OSA on Vision Transformer (ViT) has been 421 proven in Section. 4.2. We further explore the generality of OSA on target models with other archi-422 tectures. Specifically, we deploy OSA above the VSE++ (Faghri et al., 2018) model with two differ-423 ent architecture types: ResNet-152 (He et al., 2016) and VGG-19 (Simonyan & Zisserman, 2014). 424 These two architectures have revealed significant sensitivity and vulnerability to noise (Huang et al., 425 2021). In this experiment, all estimator models employ zero-shot CLIP and we utilize the origi-426 nal VSE++ as our baseline. The results in Table. 5 indicate a significant performance degradation emerged for the baseline methods in noisy setting, while a stable performance is achieved after 427 employing OSA. The stable performance on these two noise-vulnerable architectures fully demon-428 strates that OSA possesses the architecture-agnostic property. 429

430

431 Adaptability to Other Anti-Noise Models. Theoretically, OSA can be adapted to any target model, providing noise resistance. However, can OSA further enhance the robustness of models

					MS-CC	DCO 1	K		1		MS-CO	OCO 5	Κ	
Noise ratio	Method	Architecture		i2t			t2i			i2t			t2i	
			R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	0 R@1	R@5	R@10
0%	Baseline +OSA	ResNet-152	58.9 58.9	<b>86.9</b> 86.2	<b>93.8</b> 93.7	44.2 44.3	77.9 77.9	<b>88.3</b> 87.9	34.9 <b>35.0</b>	<b>64.3</b> 64.1	<b>76.1</b> 76.0	23.3 23.5	<b>50.9</b> 50.8	<b>64.2</b> 63.9
0,0	Baseline +OSA	VGG-19	49.6 50.1	79.4 <b>80.0</b>	89.1 <b>89.3</b>	38.0 38.3	72.9 <b>73.0</b>	<b>84.7</b> 84.6	<b>26.9</b> 26.6	54.2 <b>54.4</b>	66.8 <b>67.4</b>	18.7   <b>18.8</b>	43.8 <b>43.9</b>	56.8 <b>57.3</b>
20%	Baseline +OSA	ResNet-152	45.8 58.1	70.3 <b>86.1</b>	83.7 <b>93.2</b>	36.1 <b>43.4</b>	68.4 <b>76.8</b>	79.7 <b>87.2</b>	26.0 <b>33.7</b>	48.4 <b>62.6</b>	58.3 <b>74.5</b>	18.3 22.5	42.0 <b>49.7</b>	54.0 <b>62.8</b>
2070	Baseline +OSA	VGG-19	33.2 <b>49.3</b>	67.1 <b>79.1</b>	81.5 <b>88.6</b>	25.9 <b>37.2</b>	58.0 <b>71.9</b>	71.4 <b>83.8</b>	13.7 <b>25.2</b>	35.0 <b>53.3</b>	49.2 <b>65.3</b>	10.7   <b>17.9</b>	29.9 <b>42.6</b>	41.9 <b>55.9</b>
50%	Baseline +OSA	ResNet-152	28.4 55.0	61.2 <b>84.0</b>	75.2 <b>92.0</b>	5.2 <b>40.7</b>	14.0 <b>74.7</b>	19.5 <b>85.6</b>	11.0 <b>30.8</b>	31.0 <b>60.2</b>	43.6 <b>72.3</b>	1.6 <b>20.9</b>	6.0 <b>46.6</b>	9.2 <b>60.0</b>
50%	Baseline +OSA	VGG-19	2.5 47.1	9.8 77 <b>.</b> 7	16.2 <b>87.6</b>	0.1 35.7	0.5 70.3	1.0 <b>82.8</b>	0.5 <b>24.0</b>	2.5 <b>51.5</b>	4.4 <b>64.0</b>	0.0 <b>16.9</b>	0.1 <b>40.8</b>	0.2 54.2

Table 5: The results of the target model with different architectures on noisy MSCOCO.

specifically designed for noise mitigation? To investigate this, we applied OSA to the current stateof-the-art model, NPC (Zhang et al., 2024). As shown in Table. 9 of Appendix, even for noisemitigating models, OSA consistently improves training robustness. This finding further demonstrates the broad adaptability of OSA across different model types.

#### 4.4 ESTIMATOR MODEL ANALYSIS.

454 The estimator model is the basis of OSA's anti-noise capability. In this section, we explore the impact 455 of different estimator models on noise mitigation, and examine the impact of domain adaptation in noise mitigation. In Table. 10 of Appendix, we investigate four types of estimators: "None" refers 456 to training CLIP directly without using OSA. "CLIP (w/o DA)" and "ALIGN (w/o DA)" represent 457 using CLIP and ALIGN without domain adaptation as estimators, respectively, *i.e.*, zero-shot CLIP 458 and ALIGN. "CLIP (w DA)" indicates the CLIP with domain adaptation. The target models are all 459 CLIP. We can observe that both of CLIP and ALIGN as estimators significantly enhance the target 460 model performance stability when learning with noise, indicating that the choice of estimator is very 461 flexible. Both CLIP and ALIGN demonstrate exceptional performance when served as estimators. 462 The other phenomenon is that the zero-shot CLIP model shows comparable performance to the 463 domain-adapted CLIP with a even better performance at lower noise ratios. This indicates that zero-464 shot CLIP, as an estimator, already performs exceptionally well in noise mitigation. The domain 465 adaptation is unnecessary. This further enhances the deployment convenience of OSA.

466 467

468

432

448

449

450

451 452

453

## 4.5 NOISE ASSESSMENT ACCURACY

469 Noise Detection Accuracy Analysis. To figure out how accurate OSA is in recognizing noise, we evaluate the accuracy and recall on CLIP without Domain-Adaptation (w/o DA) and CLIP with 470 Domain-Adaptation (w DA) on noisy MSCOCO. We utilize zero as the threshold to roughly divide 471 pairs into noise and clean sets, respectively. Concretely, we classify scores less than or equal to 0 as 472 noise, and scores greater than 0 as clean. The Accuracy means the proportion of the clean pairs cor-473 rectly classified into the clean set, while the Recall indicates the noisy pairs correctly classified into 474 the noisy set. The results presented in Table. 6 indicates the powerful noise recognizing capability 475 of OSA. The remarkable performance on CLIP (w/o DA) fully demonstrates the generality of OSA. 476 Another notable phenomenon is that all recall scores converge towards 100, indicating that OSA 477 achieves greater accuracy in noise detection. This suggests that OSA can almost entirely eliminate 478 the impact of noise on training.

479

Noise Re-weighting Accuracy Comparison. Some anti-noise methods, like NPC, also employ loss re-weighting for optimization. To assess whether our method assigns relatively smaller weights to noise than these methods, we first analyze the weights generated by NPC and OSA. Due to differences in weight scales across methods, a direct comparison is unfair. Therefore, to unify the scale, we adopt a ranking-based approach, sorting weights in descending order and calculating the Mean Noise Rank. This metric evaluates whether smaller weights are consistently assigned to noisy samples relative to clean ones. Our experiments use 2,000 randomly selected samples from the

MSCOCO dataset under two noise conditions: 20% noise (370 noisy samples) and 50% noise (953 noisy samples). The theoretical optimal Mean Noise Ranks, where all noisy weights are ranked last, are 1815.5 and 1524.0, respectively. Results presented in Table. 11 of Appendix show that OSA achieves a higher Mean Noise Rank compared to NPC, demonstrating greater accuracy in reweighting. Moreover, OSA's rankings are nearly optimal (20% noise: 1809.1 for OSA versus 1815.5 optimal; 50% noise: 1520.7 for OSA versus 1524.0 optimal). This near-perfect alignment indicates that OSA effectively places almost all noisy samples behind the clean ones.

493 494

504 505

494 495 Table 6: ACC and recall of noise detection.

Estimator Type	Noise Ratio	Acc	Recall
CLIP (w/o DA)	0.2	93.88	97.49
CLIP (w DA)	0.2	97.68	97.18
CLIP (w/o DA)	0.5	93.91	99.35
CLIP (w DA)	0.5	98.14	99.24

Table 7: Overhead Comparison.

Model	Time	Extra Time
CLIP	97 min	0 min
NPC	323 min	226 min
OSA	118 min	21 min

#### 4.6 COMPUTATIONAL COST ANALYSIS

Cost in Pre-training. To evaluate the practicality of OSA in a real-world pre-training scenario, we estimate the additional computational cost for processing 1 billion data points. Using an NVIDIA RTX 3090 with an inference batch size of 4096, utilizing approximately 24 GB of GPU memory, processing the MS-COCO dataset consisting of 566,435 pairs takes approximately 153 seconds. At this inference rate, processing 1 billion data points would require approximately 75 hours on a single RTX 3090. This cost is negligible within the context of large-scale pre-training, especially when leveraging multiple GPUs for parallel inference.

512 513

Time Cost Comparison. To further examine the computational efficiency of our method com-514 pared to other anti-noise techniques, we evaluate training time against two representative ap-515 proaches: CLIP and NPC. CLIP, which serves as the baseline, is trained directly without any ad-516 ditional technique. NPC, the current state-of-the-art, also uses CLIP as its backbone but applies an 517 anti-noise technique by estimating the negative impact of each sample, necessitating double back-518 ward passes. The training time comparison, presented in Table. 7, shows that our method introduces only a minimal increase in training time compared to direct training, requiring just one-tenth of the 519 additional time needed by NPC. This highlights the efficiency of OSA, making it well-suited for 520 large-scale robust training tasks.

521 522 523

524

## 5 CONCLUSION

525 **Broader Impacts.** In this work, we investigated the possibility of anti-noise in practical large-scale 526 training. We introduced a novel model-agnostic anti-noise paradigm with advantages such as task 527 transferability, model adaptability, and low computational overhead. By leveraging the properties of 528 high-dimensional spaces, we found a robust and effective boundary for distinguishing between noisy 529 and clean samples. Through rigorous theoretical analysis and comprehensive experimentation, we 530 validated the efficacy and robustness of OSA for general noise mitigation. Although our primary ob-531 jective is to adapt to practical large-scale training, OSA also achieves SOTA performance in standard 532 anti-noise settings. To the best of our knowledge, this is the first work to explore noise mitigation in practical large-scale training scenarios, as well as the first to propose a general anti-noise approach. 533

534

Limitations and Future Works. Limited by the significant computational cost of pre-training, it is difficult for us to evaluate in a real pre-training process. Instead, we simulate large-scale pre-training processes to the greatest extent possible, such as evaluating on the real-world noisy dataset CC120K, which shares similar domains with mainstream pre-training datasets like CC4M and CC12M. Exploring the broad domain adaptability of OSA in real pre-training scenarios will be a valuable direction for future work.

# 540 REFERENCES

555

560

563

564

565

566

569

570

571

- Paul Albert, Diego Ortego, Eric Arazo, Noel E. O'Connor, and Kevin McGuinness. Addressing
   out-of-distribution label noise in webly-labelled data. In WACV, pp. 2393–2402, 2022. 15
- Paul Albert, Eric Arazo, Tarun Krishna, Noel E. O'Connor, and Kevin McGuinness. Is your noise correction noisy? PLS: robustness to label noise with two stage detection. In *WACV*, pp. 118–127, 2023. 15
- Simion-Vlad Bogolin, Ioana Croitoru, Hailin Jin, Yang Liu, and Samuel Albanie. Cross modal
   retrieval with querybank normalisation. In *CVPR*, pp. 5184–5195, 2022. 3
- Jiacheng Chen, Hexiang Hu, Hao Wu, Yuning Jiang, and Changhu Wang. Learning the best pooling strategy for visual semantic embedding. In *CVPR*, pp. 15789–15798, 2021. 7
- Sanghyuk Chun. Improved probabilistic image-text representations. arXiv preprint
   *arXiv:2305.18171*, 2023. 7
- Haiwen Diao, Ying Zhang, Lin Ma, and Huchuan Lu. Similarity reasoning and filtration for imagetext matching. In *AAAI*, pp. 1218–1226, 2021. 14
- Kawin Ethayarajh. How contextual are contextualized word representations? comparing the geometry of bert, elmo, and GPT-2 embeddings. In *EMNLP*, pp. 55–65, 2019. 3
- Fartash Faghri, David J. Fleet, Jamie Ryan Kiros, and Sanja Fidler. VSE++: improving visual semantic embeddings with hard negatives. In *BMCV*, pp. 12, 2018. 8
  - Aritra Ghosh, Himanshu Kumar, and P. S. Sastry. Robust loss functions under label noise for deep neural networks. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, *February 4-9, 2017, San Francisco, California, USA*, pp. 1919–1925. AAAI, 2017. 15
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pp. 770–778, 2016.
  - Zhenyu Huang, Guocheng Niu, Xiao Liu, Wenbiao Ding, Xinyan Xiao, Hua Wu, and Xi Peng. Learning with noisy correspondence for cross-modal matching. In *NeurIPS*, pp. 29406–29419, 2021. 1, 2, 8, 14, 15
- Sarah Ibrahimi, Arnaud Sors, Rafael Sampaio de Rezende, and Stéphane Clinchant. Learning with
   label noise for image retrieval by selecting interactions. In *WACV*, pp. 468–477, 2022a. 1
- Sarah Ibrahimi, Arnaud Sors, Rafael Sampaio de Rezende, and Stéphane Clinchant. Learning with label noise for image retrieval by selecting interactions. In *WACV*, pp. 468–477, 2022b. 15
- 578 Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan
  579 Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning
  580 with noisy text supervision. In *ICML*, volume 139, pp. 4904–4916, 2021. 2, 7
- Kuang-Huei Lee, Xi Chen, Gang Hua, Houdong Hu, and Xiaodong He. Stacked cross attention for image-text matching. In *ECCV*, volume 11208, pp. 212–228, 2018. 14
- Hao Li, Jingkuan Song, Lianli Gao, Pengpeng Zeng, Haonan Zhang, and Gongfu Li. A differentiable
   semantic metric approximation in probabilistic embedding for cross-modal retrieval. In *NeurIPS*,
   volume 35, pp. 11934–11946, 2022. 14
- Hao Li, Jingkuan Song, Lianli Gao, Xiaosu Zhu, and Hengtao Shen. Prototype-based aleatoric uncertainty quantification for cross-modal retrieval. In *NeurIPS*, 2023. 7
- Junnan Li, Richard Socher, and Steven C. H. Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. In *ICLR*, 2020. 2, 15
- 593 Kunpeng Li, Yulun Zhang, Kai Li, Yuanyuan Li, and Yun Fu. Visual semantic reasoning for imagetext matching. In *ICCV*, pp. 4653–4661, 2019. 14

594 595 596	Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Y. Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. In <i>NeurIPS</i> , 2022a. 3, 6, 16
597 598 599	Xuefeng Liang, Longshan Yao, Xingyu Liu, and Ying Zhou. Tripartite: Tackle noisy labels by a more precise partition. <i>CoRR</i> , 2022b. 15
600 601 602	Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In <i>ECCV</i> , volume 8693, pp. 740–755, 2014. 2, 7
603 604 605 606	Chang Liu, Han Yu, Boyang Li, Zhiqi Shen, Zhanning Gao, Peiran Ren, Xuansong Xie, Lizhen Cui, and Chunyan Miao. Noise-resistant deep metric learning with ranking-based instance selection. In <i>CVPR</i> , pp. 6811–6820, 2021. 1, 7, 14, 15
607 608 609	Aditya Krishna Menon, Brendan van Rooyen, Cheng Soon Ong, and Bob Williamson. Learning from corrupted binary labels via class-probability estimation. In <i>ICML</i> , volume 37, pp. 125–134, 2015. 15
610 611 612	Nagarajan Natarajan, Inderjit S. Dhillon, Pradeep Ravikumar, and Ambuj Tewari. Learning with noisy labels. In <i>NeurIPS</i> , pp. 1196–1204, 2013. 15
613 614 615	Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In <i>CVPR</i> , pp. 2233–2241, 2017. 15
616 617 618	Yang Qin, Dezhong Peng, Xi Peng, Xu Wang, and Peng Hu. Deep evidential learning with noisy correspondence for cross-modal retrieval. In <i>ACM MM</i> , pp. 4948–4956, 2022. 2, 8, 15
619 620 621 622	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar- wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In <i>ICML</i> , volume 139, pp. 8748–8763, 2021. 1, 2, 5, 7
623 624	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In <i>CVPR</i> , pp. 10674–10685, 2022. 2
625 626 627	Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. <i>arXiv preprint arXiv:1409.1556</i> , 2014. 8
628 629	Yale Song and Mohammad Soleymani. Polysemous visual-semantic embedding for cross-modal retrieval. In <i>CVPR</i> , pp. 1979–1988, 2019. 14
630 631 632 633	Zeren Sun, Yazhou Yao, Xiu-Shen Wei, Yongshun Zhang, Fumin Shen, Jianxin Wu, Jian Zhang, and Heng Tao Shen. Webly supervised fine-grained recognition: Benchmark datasets and an approach. In <i>ICCV</i> , pp. 10582–10591. IEEE, 2021. 1, 7, 14
634 635 636	Zeren Sun, Fumin Shen, Dan Huang, Qiong Wang, Xiangbo Shu, Yazhou Yao, and Jinhui Tang. PNP: robust learning from noisy labels by probabilistic noise prediction. In <i>CVPR</i> , pp. 5301–5310, 2022. 15
637 638 639	Dong Wang and Xiaoyang Tan. Robust distance metric learning via bayesian inference. <i>IEEE Trans. Image Process.</i> , 27(3):1542–1553, 2018. 15
640 641 642	Xinshao Wang, Yang Hua, Elyor Kodirov, David A Clifton, and Neil M Robertson. Imae for noise- robust learning: Mean absolute error does not treat examples equally and gradient magnitude's variance matters. <i>arXiv preprint arXiv:1903.12141</i> , 2019a. 15
643 644 645	Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross entropy for robust learning with noisy labels. In <i>ICCV</i> , pp. 322–330, 2019b. 15
646 647	Xiaobo Xia, Tongliang Liu, Nannan Wang, Bo Han, Chen Gong, Gang Niu, and Masashi Sugiyama. Are anchor points really indispensable in label-noise learning? In <i>NeurIPS</i> , pp. 6835–6846, 2019.

648 649 650	Yilun Xu, Peng Cao, Yuqing Kong, and Yizhou Wang. L_dmi: A novel information-theoretic loss function for training deep nets robust to label noise. In <i>NeurIPS</i> , pp. 6222–6233, 2019. 15
651 652 653	Shuo Yang, Zhaopan Xu, Kai Wang, Yang You, Hongxun Yao, Tongliang Liu, and Min Xu. Bi- cro: Noisy correspondence rectification for multi-modality data via bi-directional cross-modal similarity consistency. In <i>CVPR</i> , pp. 19883–19892, 2023a. 1, 8, 15
654 655 656	Xinlong Yang, Haixin Wang, Jinan Sun, Shikun Zhang, Chong Chen, Xian-Sheng Hua, and Xiao Luo. Prototypical mixing and retrieval-based refinement for label noise-resistant image retrieval. In <i>ICCV</i> , pp. 11205–11215, 2023b. 15
657 658 659	Yazhou Yao, Zeren Sun, Chuanyi Zhang, Fumin Shen, Qi Wu, Jian Zhang, and Zhenmin Tang. Jo-src: A contrastive approach for combating noisy labels. In CVPR, pp. 5192–5201, 2021. 15
660 661	Kun Yi and Jianxin Wu. Probabilistic end-to-end noise correction for learning with noisy labels. In <i>CVPR</i> , pp. 7017–7025, 2019. 15
662 663 664 665	Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. <i>Trans. Assoc.</i> <i>Comput. Linguistics</i> , 2:67–78, 2014. 7
666 667	Xingrui Yu, Bo Han, Jiangchao Yao, Gang Niu, Ivor W. Tsang, and Masashi Sugiyama. How does disagreement help generalization against label corruption? In <i>ICML</i> , pp. 7164–7173, 2019. 1
668 669 670	Xu Zhang, Hao Li, and Mang Ye. Negative pre-aware for noisy cross-modal matching. In <i>AAAI</i> , pp. 7341–7349, 2024. 1, 2, 6, 7, 8, 9, 14, 15
671 672	Zhilu Zhang and Mert R. Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In <i>NeurIPS</i> , pp. 8792–8802, 2018a. 15
673 674 675 676	Zhilu Zhang and Mert R. Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In <i>NeurIPS</i> , pp. 8792–8802, 2018b. 15
677 678	
679 680 681	
682 683	
684 685	
687 688	
689 690	
691 692 693	
694 695	
696 697	
698 699 700	

APPENDIX

703 704 705

## A DETAILS OF IMPLEMENTATION

706 Dataset Details. MSCOCO is widely used for noisy cross-modal matching, with each image ac-707 companied by five descriptive captions. Following the setting of Huang et al. (2021), we utilize 708 113,287 images for training, 5,000 for validation, and 5,000 for testing. The Flickr30K dataset en-709 compasses 31,783 image-text instances, each image paired with five textual annotations. Adhering 710 to the NCR (Huang et al., 2021), we use 29,783 images for training and 1,000 images each for 711 validation and testing. Regarding noise splits, following the NCR categorization, we conduct ex-712 periments at noise ratios of 0%, 20%, 40%, and 60%. CC120K is a real-world multimodal noisy 713 dataset collected by Zhang et al. (2024) from the Internet, with about 3%-20% noise ratio. There are 714 118,851 image-text pairs for training, 1,000 for validation, and 1,000 for testing.

The Aircraft, Bird, and Car we used in the image classification task are three non-overlapping subsets of the <u>WebFG-496</u> (Sun et al., 2021) dataset. WebFG-496 consists of 53,339 images, totaling 496 subcategories. This dataset is annotated using a webly supervised approach, which leverages resources from web search engines (*e.g.*, Google Image Search Engine, Bing Image Search Engine) to expand the annotated image dataset.

For the image retrieval task, we conduct experiments on the <u>CARS98N</u> dataset under PRISM's setting (Liu et al., 2021). We utilize 9,558 car-related images sourced from text-based searches on Pinterest as the training set, and employ the remaining 98 categories from CARS, unsearched on Pinterest, as a clean test set. The dataset's noise is inherently real-world, with its creators estimating a noise ratio of approximately 50%.

Implementation Details. To demonstrate the effectiveness of the OSA, we incorporate an estimator, built around the core of CLIP, and re-weighting operations based on the Estimator's outcomes into numerous downstream tasks. In the principal task of cross-modal image-text retrieval, we employ CLIP with ViT-B/32 as the baseline and target model by default. All experiments are conducted on a single RTX 3090 GPU using the AdamW optimizer. During both training phases, the model is trained for five epochs with a batch size of 256 and 500 warmup steps.

732 For the image classification task on the WebFG dataset, we align with the field's prevalent models 733 for a fair comparison by employing the ResNet-50 model enhanced by CLIP for feature extraction 734 and the CLIP image encoder as our estimator. Training and testing are executed on single RTX 735 3090 GPU, with an input image resolution of  $448 \times 448$ . The batch size and initial learning rate are 736 specified as 64 and 1e-5, respectively. In the first phase, the estimator is trained with data modeled by a Gaussian Mixture Model (GMM), which considers the classification and matching losses of all 737 training samples, with the GMM probability threshold of 0.95. The classification task leverages the 738 CLIP protocol, where a fixed prompt ("This is a picture of") is prepended to category texts. 739

For the image retrieval task, we use CLIP ViT-B/32 as the baseline, with a batch size set to 128, an initial learning rate of 5e-6, and the number of epochs set to 10. Following the setup of the PRISM (Liu et al., 2021), we set the parameter for sampling positive examples by the random sampler of the dataloader to 4, and adjust the number of positive examples sampled per epoch to one-fourth of the original parameter according to the increase in batch size. In this task, we also adopt a two-stage training approach. The strictly clean in-domain training data for the first stage is obtained using a GMM model with a probability setting of 0.8.

747 748

749

## B RELATED WORK

#### 750 B.1 NOISE MITIGATION IN CROSS-MODAL MATCHING 751

The cross-modal matching task (Lee et al., 2018; Song & Soleymani, 2019; Li et al., 2019; 2022;
Diao et al., 2021) serves as a fundamental component in multimodal learning. However, the inherent difference in information density between these modalities leads to high annotation costs and
inconsistent annotation quality, rendering cross-modal tasks particularly vulnerable to label noise.
Some approaches explicitly identify and correct noisy samples through cross-prediction between

756 concurrently trained dual models (Huang et al., 2021; Yang et al., 2023a; Liang et al., 2022b), while 757 others (Zhang et al., 2024; Qin et al., 2022) implicitly estimate the probability of sample noise, 758 reducing its training impact by adjusting the loss function. NCR (Huang et al., 2021) employs 759 the memorization capacity of its counterpart model for simple clean samples to rectify the output 760 results. BiCro (Yang et al., 2023a) utilizes the consistency of similarity score distributions from a Siamese model ensemble on noisy data, alongside anchors modeled on the loss distribution via 761 a Beta-Mixture-Model (BMM), to filter out noisy samples. NPC (Zhang et al., 2024), deviating 762 from the dual-model training schemes, introduces a two-stage single-model training approach that 763 reduces training overhead by replacing two backward passes with one forward and one backward 764 pass. Specifically, the first stage estimates the impact of potentially noisy samples on model per-765 formance by constructing a high-quality clean sample bank; the second stage then utilizes these 766 estimates to reweight the loss function. However, current methods for distinguishing clean from 767 noisy samples rely on numerous hyperparameters that are closely linked to dataset size and model 768 capacity. This dependency not only limits their adaptability to various downstream tasks but also 769 makes them challenging to deploy in real-world applications.

- 770
- 771
- 772 773

#### **B.2** NOISE MITIGATION IN IMAGE CLASSIFICATION

774 775

Image classification is vulnerable to training data noise, due to varied noise types and strong model 776 memorization. Noise in datasets manifests in two primary forms: synthetic alterations and those 777 arising from real-world scenarios. The former typically involves shuffling the labels of a subset 778 of the data or retaining the labels while introducing corresponding category images from external 779 datasets. The latter entails substituting images for a random selection of data points with those sourced from image search engines. Existing approaches are categorized based on their operational 781 focus: loss correction (Yi & Wu, 2019; Zhang & Sabuncu, 2018a; Menon et al., 2015; Natarajan 782 et al., 2013; Patrini et al., 2017; Xia et al., 2019; Ghosh et al., 2017; Wang et al., 2019a;b; Xu et al., 783 2019; Zhang & Sabuncu, 2018b) and sample selection (Sun et al., 2022; Albert et al., 2023; Yao et al., 2021; Li et al., 2020; Albert et al., 2022). Loss correction methods typically incorporate a 784 785 regularization term into the loss function, implicitly reweighting clean and noisy samples within the loss. Sample selection strategies, in contrast, explicitly differentiate between clean and noisy 786 samples, applying distinct processing to each category during loss computation. Representative for 787 the loss correction category, (Zhang & Sabuncu, 2018a) aims to generalize ordinary Cross-Entropy 788 loss and MAE loss by setting the loss threshold to iid and ood noisy samples. DivideMix (Li et al., 789 2020) concurrently trains two networks, each utilizing the data partitioning from the other network to 790 distinguish between clean and noisy samples based on loss values, thereby mitigating the influence 791 of confirmation bias inherent within each network. PNP (Sun et al., 2022) framework employs a 792 unified predictive network to estimate the in-distribution (iid), out-of-distribution (ood), and clean 793 probabilities for a given sample. Co-training trained on a sample that has a lower loss, and with the 794 different predictions by its siamese network.

- 795
- 796
- 797 798

799

## B.3 NOISE MITIGATION IN IMAGE RETRIEVAL.

800 Although image retrieval tasks focus on pairwise relationships, the noise predominantly originates 801 from image categorization errors. Analogous to image classification tasks, this can be bifurcated into 802 in-domain (Wang & Tan, 2018) and open-set noise (Liu et al., 2021). In terms of task configuration, 803 noise retrieval typically operates at the category level, treating images within the same category 804 as positive instances. PRISM (Liu et al., 2021) tries to find noisy image samples by finding the 805 outliers score in the whole similarity matrix from the same category. The generalization ability of 806 the image feature is ensured by a broader query bank restored multi-view of it. TITAN (Yang et al., 807 2023b) utilizes prototypes to be representative of the anchor of the clean and noisy samples and then generates synthetic samples by a combination of prototypes for substitution of noisy samples. 808 T-SINT (Ibrahimi et al., 2022b) utilizes more negative samples by the interaction between noisy samples and negative samples that belong to another category.

# <sup>810</sup> C PROOFS

## 812 C.1 PROOF OF HIGH-DIMENSIONAL ORTHOGONALITY

Suppose  $u, v \in \mathbb{R}^d$  are any two random vectors. The cosine similarity  $\cos(u, v) \sim \mathcal{N}(0, d^{-1})$ . The probability that  $\cos(u, v)$  is within a specific range [-a, a] is denoted as:

$$P(-a \le \cos(u, v) \le a) = \Phi\left(\frac{a}{\varsigma}\right) - \Phi\left(\frac{-a}{\varsigma}\right),\tag{7}$$

where  $\Phi$  represents the CDF of the standard normal distribution, and  $\varsigma = \frac{1}{\sqrt{d}}$  is the standard deviation of the cosine similarity. When d = 1024 and a = 0.1, there are

$$\varsigma = \frac{1}{\sqrt{1024}} = \frac{1}{32},$$
(8)

and

813

816 817 818

819

824 825

827

828

841

842 843

852 853

855 856

858 859

860 861

862 863

$$P(-0.1 \le \cos(u, v) \le 0.1) = \Phi\left(\frac{0.1}{1/32}\right) - \Phi\left(\frac{-0.1}{1/32}\right) \approx 0.9986.$$
(9)

#### C.2 PROOF OF THEOREM 1

In the Section. 2.2, we propose that Theorem 1 about the relative relationship of pairs in the original entire space, will not change after transmitting to the narrow cone space of the trained model, and there is always a boundary r concentrated on most random vectors.

To prove this Theorem, we first introduce a useful lemma of monotonicity of cosine similarity proposed by Liang et al. (2022a), indicating that the cosine similarity between two vectors increases with a high probability after one feedforward computation consisting of a linear transformation and ReLU computation.

Lemma 1. Suppose  $u, v \in \mathbb{R}^{d_{in}}$  are any two fixed vectors such that ||u|| = r ||v|| for some r > 0, W  $\in \mathbb{R}^{d_{out} \times d_{in}}$  is a random weight matrix where each element  $\mathbf{W}_{k,l} \sim \mathcal{N}(0, d_{out}^{-1})$  for  $k \in [d_{out}]$ ,  $l \in [d_{in}]$ , and  $\mathbf{b} \in \mathbb{R}^{d_{out}}$  is a random bias vector such that  $\mathbb{b}_k \sim \mathcal{N}(0, d_{out}^{-1})$  for  $k \in [d_{out}]$ . If  $\cos(u, v) < (\frac{1}{2}(r + \frac{1}{r}))^{-1}$ , then the following holds with probability at least  $1 - O(1/d_{out})$ .

$$\cos(\sigma(\mathbf{W}u + \mathbf{b}), \sigma(\mathbf{W}v + \mathbf{b})) > \cos(u, v).$$
<sup>(10)</sup>

844 Proof of Theorem. 1. Let  $\mathbb{R}^{d_{in}}$  be the original space before being transmitted in a neural network. 845 Suppose  $u, v \in \mathbb{R}^{d_{in}}$  are any two random vectors with  $\cos(u, v) \approx 0$ .  $u_c, v_c \in \mathbb{R}^{d_{in}}$  is a pair of 846 clean vectors with  $\cos(u_c, v_c) > 0$ , while  $u_n, v_n \in \mathbb{R}^{d_{in}}$  is a noisy pair with  $\cos(u_n, v_n) < 0$ . Given 847 a Neural Network  $F(x) = f_t(f_{t-1}(\dots f_2(f_1(x)))) \in \mathbb{R}^{d_{out}}$  with t layers.  $f_i(x) = \sigma_i(\mathbf{W}_i x + \mathbf{b}_i)$ 848 denotes  $i^{th}$  layer, where  $\sigma(\cdot)$  indicates activation function.  $\mathbf{W}_i \in \mathbb{R}^{d_{out}^i \times d_{in}^i}$  is a random weight 849 matrix where each element  $\mathbf{W}_i^{k,l} \sim \mathcal{N}(0, 1/d_{out}^i)$  for  $k \in [d_{out}^i]$ ,  $l \in [d_{in}^i]$ , and  $\mathbf{b}_i \in \mathbb{R}^{d_{out}^i}$  is a 850 random bias vector such that  $\mathbf{b}_i^k \sim \mathcal{N}(0, 1/d_{out}^i)$  for  $k \in [d_{out}^i]$ . We would like to prove that there 851 are always be a boundary  $\beta$ , satisfying:

$$\cos(F(u_n), F(v_n)) < \cos(F(u), F(v)) \approx \beta < \cos(F(u_c), F(v_c)), \tag{11}$$

which is equivalent to proving.

$$\cos(f_i(u_n), f_i(v_n)) < \cos(f_i(u), f_i(v)) \approx \beta_i < \cos(f_i(u_c), f_i(v_c)), \tag{12}$$

where  $\beta_i$  is the boundary of  $i^{th}$  layer.

We first consider the cosine similarity between u and v as:

$$\cos(u, v) = \frac{u \cdot v}{\|u\| \|v\|}.$$
(13)

After a linear transformation of  $i^{th}$  layer, the cosine similarity of  $\cos(\mathbf{W}_i u + \mathbf{b}_i, \mathbf{W}_i v + \mathbf{b}_i)$  denotes:

$$\cos(\mathbf{W}_{i}u + \mathbf{b}_{i}, \mathbf{W}_{i}v + \mathbf{b}_{i}) = \frac{(\mathbf{W}_{i}u + \mathbf{b}_{i}) \cdot (\mathbf{W}_{i}v + \mathbf{b}_{i})}{\|\mathbf{W}_{i}u + \mathbf{b}_{i}\|\|\mathbf{W}_{i}v + \mathbf{b}_{i}\|}.$$
(14)

Since  $\mathbf{b}_i$  has a mean of zero and is independent from  $\mathbf{W}_i u$  and  $\mathbf{W}_i v$ , the expectation of  $\mathbf{b}_i$  and  $(\mathbf{W}_i u + \mathbf{b}_i) \cdot \mathbf{W}_i v + \mathbf{b}_i)$  can be signified as:

$$\mathbb{E}\left[\mathbf{b}_{i}\right] = 0,\tag{15}$$

$$\mathbb{E}\left[\left(\mathbf{W}_{i}u+\mathbf{b}_{i}\right)\cdot\left(\mathbf{W}_{i}v+\mathbf{b}_{i}\right)\right] = \mathbb{E}\left[\left(\mathbf{W}_{i}u\cdot\mathbf{W}_{i}v\right)\right] = \sum_{i=1}^{n}\sum_{i=1}^{n}\frac{1}{d_{out}^{i}}u_{k}v_{k} = \frac{1}{d_{out}^{i}}(u\cdot v). \quad (16)$$

Additionally, we have

$$\|\mathbf{W}_{i}u + \mathbf{b}_{i}\|^{2} = \mathbf{W}_{i}u \cdot \mathbf{W}_{i}u + 2\mathbf{W}_{i}u \cdot \mathbf{b}_{i} + \mathbf{b}_{i} \cdot \mathbf{b}_{i}.$$
(17)

B75 Due to  $\mathbb{b}^k \sim \mathcal{N}(0, 1/d_{out}^i)$ , as  $d_o^i ut$  increases, the term of  $2\mathbf{W}_i u \cdot \mathbf{b}_i$  and  $\mathbf{b}_i \cdot \mathbf{b}_i$  become negligible, which implies

$$\|\mathbf{W}_{i}u + \mathbf{b}_{i}\|^{2} \approx \mathbf{W}_{i}u \cdot \mathbf{W}_{i}u = \sum_{i=1}^{n} (\mathbf{W}_{i}u)^{2}.$$
(18)

Therefore, the expectation of  $\|\mathbf{W}_i u + \mathbf{b}_i\|^2$  is approximate to

$$\mathbb{E}\left[\|\mathbf{W}_{i}u\|^{2}\right] = \sum_{k=1}^{n} u_{k}^{2} \frac{1}{d_{out}^{i}} = \frac{\|u\|}{d_{out}^{i}},\tag{19}$$

884 and

$$\cos(\mathbf{W}_{i}u + \mathbf{b}_{i}, \mathbf{W}_{i}v + \mathbf{b}_{i}) \approx \frac{\mathbb{E}[\mathbf{W}_{i}u \cdot \mathbf{W}_{i}v]}{\sqrt{\mathbb{E}[\|\mathbf{W}_{i}u + \mathbf{b}_{i}\|^{2}]\mathbb{E}[\|\mathbf{W}_{i}v + \mathbf{b}_{i}\|^{2}]}}$$
$$= \frac{\frac{1}{d_{out}^{i}}(u \cdot v)}{\sqrt{\frac{1}{d_{out}^{i}}\|u\|^{2} \cdot \frac{1}{d_{out}^{i}}\|v\|^{2}}}$$
$$= \cos(u, v).$$
(20)

Based on Eq. 20, with  $\cos(u_n, v_n) < \cos(u, v) \approx 0 < \cos(u_c, v_c)$ , there are

$$\cos(\mathbf{W}_i u_n + \mathbf{b}_i, \mathbf{W}_i v_n \mathbf{b}_i) < \cos(\mathbf{W}_i u + \mathbf{b}_i, \mathbf{W}_i v + \mathbf{b}_i) < \cos(\mathbf{W}_i u_c + \mathbf{b}_i, \mathbf{W}_i v_c + \mathbf{b}_i).$$
(21)

Since the activation function  $\sigma$  is a monotonically increasing function, it follows

$$\cos(f_i(u_n), f_i(v_n)) < \cos(f_i(u), f_i(v)) < \cos(f_i(u_c), f_i(v_c)).$$
(22)

Due to Lemma. 1,  $\cos(f_i(u), f_i(v))$  will be increase with the transmitting layers, and  $\cos(f_i(u), f_i(v))$  will always be a  $\beta_i > 0$ , to satisfy:

$$\cos(f_i(u_n), f_i(v_n)) < \cos(f_i(u), f_i(v)) \approx \beta_i < \cos(f_i(u_c), f_i(v_c)).$$

$$(23)$$

After transmitting each layer, Eq. 23 are always satisfied. When transmitting a neural network with t layers, we have

$$\cos(F(u_n), F(v_n)) < \cos(F(u), F(v)) \approx \beta < \cos(F(u_c), F(v_c)).$$
(24)

#### C.3 PROOF OF ORTHOGONALITY VALIDITY IN CONE SPACE

Although we have demonstrated in Appendix. C.1 that in the original high-dimensional space, the cosine similarity between two randomly selected vectors—each dimension following a Gaussian distribution—typically converges near the orthogonal boundary, this property may not necessarily extend to the subspace of the shared embedding space maintained by the trained models. Specifically, for real image-text pairs, the subspace may deviate from the orthogonal characteristics observed in the original space. Thus, it is essential to investigate whether the orthogonality property holds within the cone space for the image-text subdomain post-training.

917 To explore this, we first analyze the distribution of several dimensions of image and text features from the CC120K dataset, as illustrated in Figure. 3. The results reveal that all vector dimensions,



Figure 3: The illustrations of several distributions on CC120K. (a) The parameter distribution. (b-d) The distribution of image features for the 128th, 256th, and 512th dimensions. (e-g) The distribution of text features for the 128th, 256th, and 512th dimensions.

including trained parameters, exhibit a Gaussian distribution with near-zero means. If the dimensions of the trained embedding space follow Gaussian distributions, the process of selecting random vectors within this space would be analogous to that of the original space, thereby preserving the orthogonality property. Here, we present the following theorem: The output features of large-scale models tend to Gaussian distribution. The detailed theorem and proof are provided below. 

**Theorem 2** (Output features tends to Gaussian). Given a Neural Network F(x) = $\{f_t(f_{t-1}(\dots f_2(f_1(x))))\} \in \mathbb{R}^{d_{out}}$  with t layers.  $f_l(x) = \phi_l(\mathbf{W}_l x + \mathbf{b}_l)$  denotes the  $l^{th}$  layer, where  $\phi(\cdot)$  indicates the activation function, and the final layer  $f_t(x) = \mathbf{W}_t x + \mathbf{b}_t$  is a fullyconnected layer without an activation function for common space projection. Let  $x^k \in \mathbb{R}^{d_{in}^k}$  be the sample feature that will be transmitted into the  $k^{th}$  layer, where  $x^1$  denotes the original fea-ture with an unknown distribution  $x^1 \sim (\mu_x, \sigma_x^2)$ .  $\mathbf{W}_k \in \mathbb{R}^{d_{out}^k \times d_{in}^k}$  is a random weight matrix where each element  $w_{ij}^k \sim \mathcal{N}(0, \sigma_w^2)$  for  $i \in [d_{out}^k]$ ,  $j \in [d_{in}^k]$ , and  $\mathbf{b}_k \in \mathbb{R}^{d_{out}^k}$  is a bias vector such that  $b_i^k \sim \mathcal{N}(0, \sigma_w^2)$  for  $i \in [d_{out}^k]$ . In such a neural network, linear layers lead features x gradually to a Gaussian distribution from any initial distribution, and as  $|d_{in}|$  is sufficiently large,  $F(x) \sim \mathcal{N}(0, \sigma^2).$ 

*Proof of Theorem.* 2. For the  $k^{th}$  layer ( $k \in [t]$ ), we first calculate the expectation and variance of the linear combination  $\sum_{j=1}^{d_{in}^k} w_{ij}^k x_j^k$ . For the expectation, since  $w_{ij}^k$  and  $x_j^k$  are independent and  $\mathbf{w}_{ij}^k \sim \mathcal{N}(0, \frac{1}{d_{out}^k})$ , we have: 

> $\mathbb{E}\left[\sum_{j=1}^{d_{in}^k} w_{ij}^k x_j^k\right] = \sum_{j=1}^{d_{in}^k} \mathbb{E}[w_{ij}^k] \mathbb{E}[x_j^k] = \sum_{j=1}^{d_{in}^k} (0 \times \mathbb{E}[x_j^k]) = 0.$ (25)

For variance, since  $w_{ij}^k$  and  $x_j^k$  are independent, we have:

971 
$$= \sum_{j=1} \sigma_{w^k}^2 \left( \sigma_{x^k}^2 + \mu_{x^k}^2 \right) = d_{in}^k \sigma_{w^k}^2 \left( \sigma_{x^k}^2 + \mu_{x^k}^2 \right).$$

Since  $w_{ij}^k$  are independently distributed Gaussian random variables, and  $x_j^k$  has a known mean and variance, the sum of  $w_{ij}^k x_j^k$  can apply to a generalized Central Limit Theorem. We have

$$\frac{\sum_{j=1}^{d_{in}^k} w_{ij}^k x_j^k - \mathbb{E}\left[\sum_{j=1}^{d_{in}^k} w_{ij}^k x_j^k\right]}{\sqrt{\operatorname{Var}\left(\sum_{j=1}^{d_{in}^k} w_{ij}^k x_j^k\right)}} \xrightarrow{d} \mathcal{N}(0,1),$$
(27)

which is equivalent to

$$\frac{\sum_{j=1}^{d_{in}^k} w_{ij}^k x_j^k - 0}{\sqrt{d_{in}^k \sigma_{w^k}^2 (\sigma_{x^k}^2 + \mu_{x^k}^2)}} \xrightarrow{d} \mathcal{N}(0, 1).$$
(28)

Therefore,

$$\sum_{j=1}^{d_{in}^k} w_{ij}^k x_j^k \xrightarrow{d} \mathcal{N}(0, d_{in}^k \sigma_{w^k}^2 (\sigma_{x^k}^2 + \mu_{x^k}^2)).$$
(29)

Due to  $b^k \sim \mathcal{N}(0, \sigma_b^2)$ , we finally get

$$\sum_{j=1}^{d_{in}^k} w_{ij}^k x_j^k + b_i^k \xrightarrow{d} \mathcal{N}\left(0, d_{in}^k \sigma_{w^k}^2 (\sigma_{x^k}^2 + \mu_{x^k}^2) + \sigma_b^2\right).$$
(30)

Although activation functions truncate the Gaussian distribution after each linear layer, the samples still gradually approach a Gaussian distribution from the initial unknown distribution as they pass through the layers. Furthermore, because there is a fully connected layer (layer) without an activation function before mapping to the final common space, the final feature distribution will approximate a Gaussian distribution, as follows:

$$F(x) \sim \mathcal{N}(0, d_{in}^t \sigma_{w^t}^2 (\sigma_{x^t}^2 + \mu_{x^t}^2) + \sigma_b^2).$$
(31)

#### 1020 D SDM VISUALIZATION

We visualize some representative samples from our synthetic domain originating from COCO by using SDM. The results are shown in Figure. 4. We generate two styles of image based on the MSCOCO caption, and then use pre-trained multimodal models to calculate cosine similarity with the SDM-generated image and original caption.



90.9

96.7

66.3

87.2

87.3

98.1

93.0

99.6

52.1

74.4

78.8

92.9

87.4

96.4

1078

1079

CLIP

+OSA

93.3

97.6

76.2

87.3

59.4

74.2

85.0

93.1

96.5

99.3

			i2t		t2i				
Noise Ratio	Method	R@1	R@5	R@10	R@1	R@5	R@10		
007	NPC	82.2	96.5	<b>98.7</b>	68.3	92.0	<b>98.</b> 7		
0%	+OSA	82.4	96.4	98.6	68.5	91.8	<b>98.7</b>		
200%	NPC	79.9	95.9	98.4	66.3	90.5	98.4		
20%	+OSA	81.2	96.0	98.6	66.9	91.2	98.6		
50%	NPC	78.2	94.4	97.7	63.1	89.0	97.7		
50%	+OSA	79.3	95.6	98.2	66.8	90.8	98.2		

Table 9: The results of other methods employing OSA on MSCOCO 1K.

Table 10: Ablation study of estimator type on noisy MS-COCO.

				MS-CC	DCO 1H	K				MS-CC	CO 5I	K	
Noise ratio	Estimator		i2t			t2i			i2t			t2i	
		R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
	None	80.1	95.7	98.2	67.1	91.4	96.6	62.9	84.9	91.6	46.5	73.8	82.9
007	CLIP (w/o DA)	82.6	96.7	98.7	68.5	92.1	96.7	66.2	87.0	93.3	48.6	75.7	84.8
0%	ALIGN (w/o DA)	81.9	96.7	98.7	68.9	92.2	96.9	64.8	86.6	92.7	49.0	75.9	84.7
	CLIP (w DA)	82.2	96.5	98.7	68.8	92.1	96.7	65.6	86.8	92.9	49.1	76.2	84.8
	None	76.0	94.3	97.5	63.4	89.0	94.8	55.3	79.1	86.9	41.0	68.8	79.3
200	CLIP (w/o DA)	81.8	96.1	98.7	68.2	91.9	96.5	64.8	86.6	92.3	48.3	75.4	84.1
20%	ALIGN (w/o DA)	81.2	96.0	98.6	67.7	91.5	96.4	64.8	86.2	92.3	47.8	74.9	83.9
	CLIP (w DA)	81.6	96.2	98.5	68.9	92.0	96.6	65.8	86.4	92.5	48.7	76.1	84.5
	None	73.9	93.0	97.2	60.1	87.3	94.0	54.1	78.5	86.6	39.7	67.2	77.5
5007	CLIP (w/o DA)	79.6	95.6	98.4	65.9	90.8	95.9	62.4	84.8	90.8	45.7	73.1	82.5
50%	ALIGN (w/o DA)	80.4	95.6	98.3	66.0	90.5	95.8	62.0	84.9	91.8	45.7	73.2	82.5
	CLIP (w DA)	80.4	96.2	98.6	67.8	91.6	96.4	64.0	85.5	91.9	47.9	74.6	83.8

Table 11: Mean Noise Rank Comparison between OSA and NPC.

Noise Ratio	Method	Mean Noise Rank↑	<b>Optimal Rank</b>	Noise Number	Sample Number
20%	NPC	1641.3	1815.5	370	2,000
20%	OSA	1809.1	1815.5	370	2,000
50%	NPC	1456.2	1524.0	953	2,000
50%	OSA	1520.7	1524.0	953	2,000