

Unlearning the Noisy Correspondence Makes CLIP More Robust

Anonymous ICCV submission

Paper ID 8160

Abstract

001 *The data appetite for Vision-Language Models (VLMs)*
 002 *has continuously scaled up from the early millions to bil-*
 003 *lions today, which faces an untenable trade-off with data*
 004 *quality and inevitably introduces Noisy Correspondence*
 005 *(NC) samples. Undoubtedly, such semantically unrelated*
 006 *data significantly impairs the performance of VLMs. Previ-*
 007 *ous efforts mainly address this challenge by estimating re-*
 008 *fined alignment for more precise guidance. However, such*
 009 *resource-intensive pipelines that train VLMs from scratch*
 010 *struggle to meet realistic data demands. In this paper, we*
 011 *present a brand new perspective that seeks to directly elimi-*
 012 *nate the harmful effects of NC in pre-trained VLMs. Specif-*
 013 *ically, we propose NCU, a Noisy Correspondence Unlearn-*
 014 *ing fine-tuning framework that efficiently enhances VLMs'*
 015 *robustness by forgetting learned noisy knowledge. The key*
 016 *to NCU is learning the hardest negative information, which*
 017 *can provide explicit unlearning direction for both false pos-*
 018 *itives and false negatives. Such twin goals unlearning pro-*
 019 *cess can be formalized into one unified optimal transport*
 020 *objective for fast fine-tuning. We validate our approach with*
 021 *the prevailing CLIP model over various downstream tasks.*
 022 *Remarkably, NCU surpasses the robust pre-trained method*
 023 *on zero-shot transfer while with lower computational over-*
 024 *head. The code will be released upon acceptance.*

1. Introduction

026 The pursuit of general intelligence has driven progress in
 027 multimodal learning, which seeks to integrate and under-
 028 stand multiple sensory modalities like humans. Large-scale
 029 vision-language training, exemplified by CLIP [37], is seen
 030 as a key milestone in multimodal learning due to its remark-
 031 able transfer capabilities in real-world applications, such as
 032 image-text retrieval [13, 22, 34] and robotics control [41].

033 However, much of their success can be attributed to
 034 scaling laws enabled by massive training data. As every
 035 coin has two sides, the insatiable demand for data forces a
 036 difficult trade-off between quantity and quality, which in-
 037 evitably introduces noisy correspondence into the training

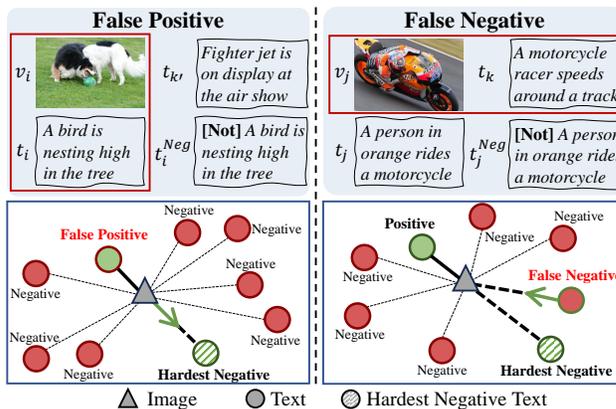


Figure 1. **Illustration on the core concept of NCU.** The twin goals unlearning process is guided by the learned hardest negative information. For the FP, t_i^{Neg} directly pulls v_i away from the mismatched t_i . While for the FN, t_i^{Neg} acts as a distance upper bound to facilitate modeling many-to-many relations.

038 set. Taking the CC3M dataset [39] as an example, despite
 039 being filtered from 500 million images, it still contains at
 040 least 3% [20] unrelated image-text pairs, *i.e.*, *false positive*.
 041 To make matters worse, training on massive data necessitates
 042 larger batch sizes (32K used in CLIP), which increases
 043 the likelihood of unpaired samples sharing semantic simi-
 044 larities, *i.e.*, *false negative*. Undoubtedly, such two-aspects
 045 noisy correspondence can significantly impair the perfor-
 046 mance of vision-language models.

047 To endow robustness against NC, one natural direction
 048 is to revise the pre-training paradigm [1, 3, 11, 12, 21] that
 049 supervises VLMs with refined alignment. However, exist-
 050 ing methods require training from scratch and may rely on
 051 guidance from external large models [3, 11]. Such resource-
 052 intensive pipelines obviously struggle to face the realistic
 053 demand, especially with today’s billion-scale datasets [38].
 054 Hence, it is necessary to address the NC problem in vision-
 055 language training with a cost-effective method.

056 In this paper, we think outside the box of robust pre-
 057 training and pose an important question: *Can we directly*
 058 *eliminate the harmful effects of NC in pre-trained VLMs?*

059 To answer this question, we resort to machine unlearning
060 [2] and present NCU, a Noisy Correspondence Unlearning
061 fine-tuning framework that improves the robustness of CLIP
062 by erasing learned noisy knowledge. Machine unlearning is
063 a reversed learning process that aims to delete the influence
064 of specific training samples from trained models. Despite
065 its promise in widespread tasks [7, 9], unlearning the NC
066 in VLMs remains unexplored due to a key challenge: the
067 ambiguous forgetting direction would corrupt the learned
068 semantic structure in the feature space. To address this,
069 we propose to learn the hardest negative information that
070 can provide explicit unlearning direction. As illustrated in
071 Fig. 1, on the one hand, the negative information would di-
072 rectly serve as reliable supervision for forgetting false posi-
073 tives. On the other hand, it would facilitate the modeling of
074 many-to-many relationships among unpaired data for for-
075 getting false negatives. Then, we show that such twin goals
076 unlearning process can be formalized as one unified opti-
077 mal transport problem, which efficiently fine-tunes CLIP to
078 resist both FP and FN.

- 079 Our main contributions are highlighted below:
- 080 • To the best of our knowledge, this work could be the first
081 study to eliminate the harmful effects of noisy correspon-
082 dence from pre-trained CLIP.
 - 083 • We propose the NCU framework, which efficiently un-
084 learns FP and FN with explicit direction derived from the
085 hardest negative information.
 - 086 • We demonstrate that NCU achieves significant improve-
087 ments over CLIP on several downstream tasks and sur-
088 passes the previous robust pre-training method with lower
089 computational overhead.

090 **2. Related Work**

091 **Noisy Correspondence Learning.** Noisy correspon-
092 dence refers to the alignment error presented in multimodal
093 data. The false positive is a typical NC problem, where ir-
094 relevant multimodal pairs are wrongly treated as matched.
095 To alleviate this, several techniques have been developed in
096 various multimodal applications, including cross-modal re-
097 trieval [16, 18, 20, 36], video temporal learning [17, 30],
098 multimodal person re-identification [35, 47], question an-
099 swering [23], and image captioning [10, 24]. In more com-
100 plex scenarios, *e.g.*, vision-language pre-training [19, 21],
101 models also suffer from false negatives caused by the train-
102 ing paradigm [12], where similar unpaired samples are
103 forced to be distant. Considering the computational burden
104 of large VLMs, this work presents a low-carbon solution to
105 directly improve the robustness of pre-trained VLMs.

106 **Contrastive Vision Language Models.** Contrastive vi-
107 sion language models (VLMs) [8, 13, 14, 33, 44, 48] aim
108 to learn visual representations by the corresponding textual
109 supervision, which have attracted significant attention due

to their simplicity and powerful representation capability. 110
Pioneering works CLIP [37] and ALIGN [22] have shown 111
great success via learning from massive image-text pairs. 112
However, such web-crawled data are noisy [38, 42] and in- 113
evitably harm the efficacy of existing VLMs. To tackle this 114
issue, a series of works attempted to train the VLM with re- 115
fined soft image-text alignments by label smoothing [12], 116
knowledge distillation [1], fine-grained intra-modal guid- 117
ance [11], text rewriting [3], or positive-negative contrastive 118
loss [21]. Besides, OT-based methods[40, 46] have also 119
emerged as they naturally model such matching problems. 120
Despite the success, previous works focus on training ro- 121
bust VLMs from scratch, which overlooks readily available 122
pre-trained models and incurs unnecessary computational 123
costs. To this end, this paper pioneers an efficient approach 124
to enhance model robustness by unlearning noisy informa- 125
tion from pre-trained models. 126

Machine Unlearning. Recent advances in MU mainly fo- 127
cus on practical approximate unlearning, which seeks to 128
mimic the behavior of a model re-trained from scratch. 129
Driven by privacy concerns, existing MU works [7, 9, 32] in 130
computer vision focus on image classification that attempts 131
to forget specific classes. In parallel, MU has also become a 132
popular topic in large language models due to its capability 133
to eliminate harmful responses [31, 50]. However, multi- 134
modal forgetting remains under-explored in the literature. 135
Pioneering works [26, 27] studied data-free class removal 136
for CLIP’s downstream image classification. To date, none 137
of the existing MU methods has explored the unlearning of 138
noisy concepts from VLMs. 139

140 **3. Preliminaries**

141 **3.1. Contrastive Language-Image Pre-training**

142 CLIP is a vision-language model trained on millions of 143
web-harvested image-text pairs. We consider a batch of 144
 N image-text pairs $\{v_i, t_i\}_{i=1}^N$ sampled from a cross-modal 145
dataset \mathbb{D} , where v_i and t_i represent the raw image and cor- 146
responding text, respectively. The goal of CLIP is to train 147
two modality-specific encoders that bring matched pairs 148
closer while pushing unmatched ones apart. Specifically, 149
image embedding $v_i \in \mathbb{R}^d$ and text embedding $t_i \in \mathbb{R}^d$ are 150
obtained by passing v_i and t_i through the image encoder 151
 f_v and text encoder f_t , respectively, where d is the embed- 152
ding dimension. The encoded l_2 normalized embeddings 153
are then aligned in the feature space by minimizing the con- 154
trastive objective, *i.e.*, InfoNCE loss:

$$\mathcal{L}_{v \rightarrow t}^{CL} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\langle v_i, t_i \rangle / \tau)}{\sum_{j=1}^N \exp(\langle v_i, t_j \rangle / \tau)}, \quad (1) \quad 155$$

156 where $\langle \cdot \rangle$ represents the inner product and τ is a trainable
157 temperature parameter. As InfoNCE loss is symmetric, we

158 can define $\mathcal{L}_{t \rightarrow v}^{CL}$ similarly. The complete CLIP training ob-
159 jective is formulated as: $\mathcal{L}_{CLIP} = \mathcal{L}_{v \rightarrow t}^{CL} + \mathcal{L}_{t \rightarrow v}^{CL}$.

160 Despite its promising performance, the standard con-
161 trastive learning can suffer from the noisy correspondence
162 problem in two aspects. First, the web-collected pairs in-
163 evitably contain an unknown portion of mismatched data,
164 *i.e.*, false positives. Second, hard target alignment neglects
165 the potential semantic similarity among unpaired samples,
166 *i.e.*, false negatives, especially under large batch settings.

167 3.2. Machine Unlearning

168 Given a CLIP model (also named *reference model*) $\{f_v, f_t\}$
169 that is already trained on a cross-modal dataset \mathbb{D} , machine
170 unlearning aims to fine-tune the model to forget a specific
171 subset $\mathbb{D}_{FG} \subseteq \mathbb{D}$ while maintaining effectiveness on the
172 retained set $\mathbb{D}_{RT} = \mathbb{D} \setminus \mathbb{D}_{FG}$. Ideally, the model should
173 behave as if it were trained without any sample from \mathbb{D}_{FG} .
174 In principle, re-training the model from scratch on \mathbb{D}_{RT}
175 would serve as the gold standard. However, since CLIP
176 is trained on massive-scale data, it is unrealistic to obtain
177 a forget set that includes all noisy information, especially
178 when some data are not publicly accessible. Therefore,
179 we focus on an approximate unlearning approach in which
180 $\mathbb{D} = \mathbb{D}_{FG} \cup \mathbb{D}_{RT}$ does not need to contain all train-
181 ing pairs, making the unlearning process more practical for
182 real-world scenarios.

183 The most straightforward method to unlearn is *gradient*
184 *ascent* or its variants, which optimizes the negative pre-
185 diction loss over the forget set. Another typical approach
186 is performing *forget loss* that encourages the model to re-
187 learn the modified form of undesired data. For example,
188 we can update CLIP by minimizing InfoNCE loss in pair
189 $\{v_i, \tilde{t}_i\} \sim \mathbb{D}_{FG}$ to forget the relation between v_i and t_i ,
190 where $\tilde{t}_i \neq t_i$ could be random or hand-crafted text to re-
191 place the original. Based on these, existing MU methods
192 have shown promising progress in class forgetting and LLM
193 privacy protection. However, applying these strategies to
194 CLIP unlearning poses a key challenge: the ambiguous for-
195 getting direction would corrupt the learned semantic struc-
196 ture in the feature space. In other words, the model forgets
197 the undesired data by learning other meaningless patterns.

198 4. Methodology

199 To tackle the above issues, we introduce the Noisy Corre-
200 spondence Unlearning (NCU) framework. In the following,
201 we first introduce the division of forget and retained sets in
202 Sec 4.1. Subsequently, we elaborate on learning the hardest
203 negative information in Sec 4.2 and explain how to formal-
204 ize the twin goals unlearning process into an optimal trans-
205 port object for efficiently fine-tuning in Sec 2. The overall
206 training pseudo-code is shown in Supplementary A.

207 4.1. Identifying the Forget Set

208 Unlike standard MU tasks with a predefined forget set, we
209 need to manually identify mismatched samples from CLIP’s
210 training data to construct it. As pre-trained CLIP has shown
211 strong representation capability, we propose using the basic
212 similarity score to obtain \mathbb{D}_{FG} and \mathbb{D}_{RT} , *i.e.*,

$$213 \omega_i = \frac{1}{2} \left[\frac{\exp(\langle v_i, t_i \rangle / \tau)}{\sum_{j=1}^N \exp(\langle v_i, t_j \rangle / \tau)} + \frac{\exp(\langle t_i, v_i \rangle / \tau)}{\sum_{j=1}^N \exp(\langle t_i, v_j \rangle / \tau)} \right]. \quad (2)$$

214 By comparing (v_i, t_i) with other cross-modal samples in the
215 batch, ω_i serves as a clean confidence that measures the ex-
216 tent of semantic match. Then, we select pairs with the low-
217 est $P\%$ of ω_i within the batch as false positives to construct
218 the forget set \mathbb{D}_{FG} , while treating the remaining in-batch
219 pairs as the retained set \mathbb{D}_{RT} .

220 Note that \mathbb{D}_{FG} and \mathbb{D}_{RT} are dynamically selected at
221 each batch, which enjoys two merits: 1) CLIP could be effi-
222 ciently updated with one intra-batch optimization; 2) The
223 forget-retain ratio could be flexibly adjusted through the
224 predefined parameter P .

225 4.2. Learning Hardest Negative Semantics

226 To guide CLIP with an explicit unlearning direction, we aim
227 to learn the hardest negative semantics as supervision. In-
228 tuitionally, for an irrelevant pair (v_i, t_i) that misleads the
229 model with ‘ v_i and t_i are matched’, we encourage the
230 model to forget this information by relearning that ‘ v_i and
231 t_i are not matched’. From a data utilization viewpoint, this
232 paradigm is similar to negative learning[25] that supervises
233 the model with complementary information [18], *i.e.*, push-
234 ing the candidate away from other unpaired samples. Dif-
235 ferently, our method seeks the hardest negative information
236 to avoid uncertain optimization directions.

237 To achieve this, we incorporate a set of learnable vec-
238 tors to represent the textual negative semantics inspired by
239 prompt learning [51]. Specifically, for any training pair
240 (v_i, t_i) , the token features of t_i are combined with m shared
241 prompt vectors to present the corresponding negative seman-
242 tics t_i^{Neg} in the feature space. While such prompt-
243 driven semantic negation of CLIP has demonstrated success
244 in out-of-distribution detection [29, 43], existing methods
245 are confined to closed-set downstream tasks with limited
246 category concepts. In contrast, our challenge lies in ex-
247 tending the semantic opposite operation into the open-set
248 knowledge that CLIP pre-trained.

249 Intuitively, the hardest negative satisfies two constraints
250 in the feature space: 1) t_i^{Neg} needs to maximize its distance
251 from v_i and t_i ; 2) t_i^{Neg} should maintain certain similarity
252 to those unpaired images, as it is not a wrong description [43]
253 despite being semantically irrelevant to v_j . Furthermore, we
254 only use \mathbb{D}_{RT} to learn the prompt tokens to avoid overfitting
255 caused by noisy correspondence. For notation simplicity,

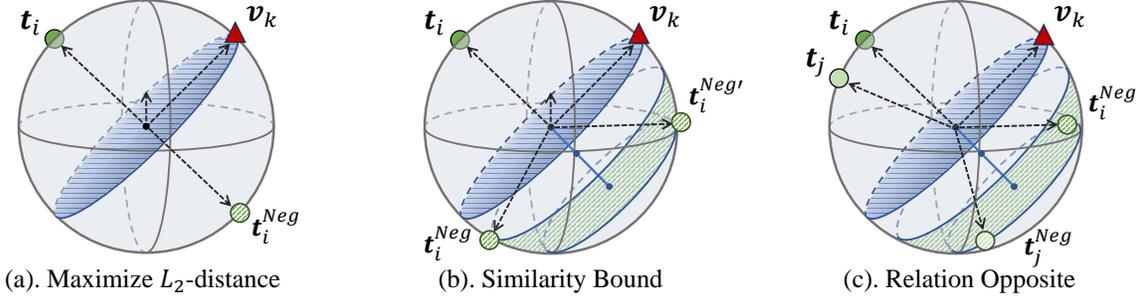


Figure 2. **Illustration of the optimization objective to learn the hardest negative semantics.** (a) Previous attempts that directly maximize the L_2 distance prevent t_i^{Neg} from providing certain guidance for unpaired images, e.g., v_k . (b) We bound the similarity with margins for a more relaxed semantic separation, but it may lead to uncertain targets, e.g., t_i^{Neg} and $t_i^{Neg'}$. (c) We further preserve relation structures for precise objectives. The intuition is that the opposite text should also maintain semantic relationships, e.g., $\langle t_i, t_j^{Neg} \rangle \approx \langle t_j, t_i^{Neg} \rangle$.

we denote \tilde{N} as the size of \mathbb{D}_{RT} within each batch. Based on the above insights, we propose the following intra-modal and cross-modal training objectives.

Text Relation Opposite. It encourages semantic separation between the embeddings of negative text and its original. Most previous works [29, 43] typically reduce the per-instance similarity gap among textual pairs, e.g., $\|t_i - t_i^{Neg}\|_2 \rightarrow 2$ [43] to directly maximize its L_2 distance. However, such rigid constraint incurs a crucial limitation in the open-set semantic space—enforcing maximal distance pushes t_i^{Neg} away from unpaired images (Fig. 2(a)), which is contrary to our objective. To this end, we propose a relaxed similarity bound to constrain the semantic separation:

$$\mathcal{L}^{sep} = \frac{1}{\tilde{N}} \sum_{i=1}^{\tilde{N}} ([\alpha - \langle t_i, t_i^{Neg} \rangle]_+ + [\langle t_i, t_i^{Neg} \rangle - \beta]_+), \quad (3)$$

where $\alpha < 0$ and $\beta < 0$ are the margin parameters to locate $\langle t_i, t_i^{Neg} \rangle \in [\alpha, \beta]$, and $[x]_+ = \max(x, 0)$ is the hinge function. As illustrated in Fig. 2(b), although minimizing Eq.(3) enables t_i^{Neg} to be distant from t_i while remaining relatively similar to unpaired images, the broad range of variations makes training convergence difficult. To address this issue, we propose to perform semantic opposite at the relation level instead of the instance level, which is achieved by preserving the geometrical structures among all negative and original text within the batch:

$$\mathcal{L}^{rel} = \frac{1}{\tilde{N}} \sum_{i=1}^{\tilde{N}} \sum_{j=1}^{\tilde{N}} (\langle t_i, t_j^{Neg} \rangle - \langle t_j, t_i^{Neg} \rangle)^2. \quad (4)$$

As shown in Fig. 2(c), regularizing the negative-original relation consistency can guide t_i^{Neg} toward a precise location in the feature space.

Image-text Matching Opposite. It aims to model the alignment between the embeddings of negative text and images. As discussed, t_i^{Neg} provides positive supervision to unpaired images while separating from its paired image, which presents opposite matching patterns to the normal contrastive objective. To achieve this, we take inspiration from Sigmoid loss [49] that efficiently supports such multi-positive alignment. Specifically, it guides per cross-modal pair independently by the binary matching target:

$$\mathcal{L}^{itm} = \frac{1}{\tilde{N}} \sum_{i=1}^{\tilde{N}} \sum_{j=1}^{\tilde{N}} \left(\log \frac{1}{1 + \exp(m_{ij}(-\langle t_i, v_j \rangle / \tau))} + \log \frac{1}{1 + \exp(-m_{ij}(-\langle t_i^{Neg}, v_j \rangle / \tau))} \right), \quad (5)$$

where m_{ij} equals 1 for $i = j$ and -1 for $i \neq j$. In Eq.(5), the first part follows the standard one-to-one matching to retain the original CLIP knowledge, while the second part utilizes the opposite binary target, i.e., $-m_{ij}$, to bring t_i^{Neg} closer to multiple images.

With the visual encoder frozen, the overall loss for learning the hardest negative semantics is balanced by a scaling factor λ and given by:

$$\mathcal{L}^{HN} = \lambda(\mathcal{L}^{sep} + \mathcal{L}^{rel}) + \mathcal{L}^{itm}. \quad (6)$$

4.3. Hardest-Negative Guided Noise Unlearning

The hardest negative semantics serve dual purposes in erasing the learned noisy correspondence: 1) guide CLIP to unlearn the false positive pair (v_i, t_i) by matching v_i with t_i^{Neg} . 2) While for the well-matched pair (v_i, t_i) , t_i^{Neg} assists in inferring the soft alignment among unpaired data to unlearn the false negative pattern. To efficiently fine-tune CLIP, we formalize such twin goals unlearning process into one unified Optimal Transport problem.

Optimal Transport. OT seeks to establish a flexible alignment between images and captions by computing a

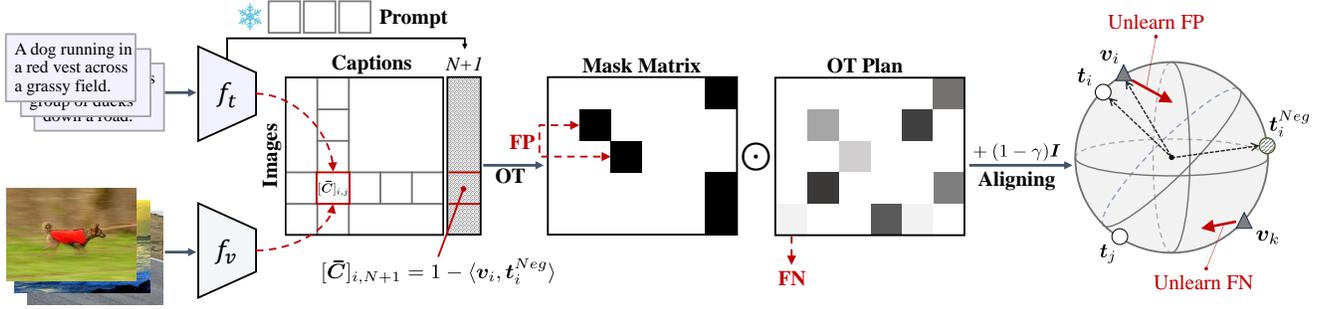


Figure 3. **Overview of the Noisy Correspondence Unlearning process.** With the learned negative prompt frozen, we formulate an optimal transport problem guided by the hardest negative and then use the solved transport plan to robustly fine-tune the model f_t and f_v .

315 minimal-cost transport plan, where the cost refers to the
 316 expense of transporting mass from source to target distribu-
 317 tion and is generally set to a distance measure [15]. Let
 318 $C \in \mathbb{R}_+^{N \times N}$ denotes the cost matrix for the mini-batch,
 319 where $[C]_{i,j} = 1 - \langle v_i, t_j \rangle$ is the cosine distance of v_i and
 320 t_j . $\Gamma \in \mathbb{R}_+^{N \times N}$ denotes the corresponding transport plan
 321 that $[\Gamma]_{i,j}$ represents the alignment probability between v_i
 322 and t_j . Formally, the objective of OT is defined as follows:

$$323 \begin{aligned} & \min_{\Gamma \in \Pi(\mu, \nu)} \langle \Gamma, C \rangle - \epsilon H(\Gamma) \\ & \text{s.t. } \Pi(\mu, \nu) = \{ \Gamma \in \mathbb{R}_+^{N \times N} \mid \Gamma \mathbb{1}_N = \mu, \Gamma^\top \mathbb{1}_N = \nu \}, \end{aligned} \quad (7)$$

324 where $\mathbb{1}_N$ denotes a N -dimensional all-one vector, $\mu, \nu \in$
 325 \mathbb{R}^N are probability measures representing the relative impor-
 326 tance of each image and caption. Without prior knowl-
 327 edge, $\mu = \frac{1}{N} \mathbb{1}_N$ and $\nu = \frac{1}{N} \mathbb{1}_N$ are considered to be uni-
 328 formly distributed since each pair is sampled independently.
 329 $H(\Gamma)$ is an additional entropy regularizer controlled by the
 330 smooth parameter ϵ , which enables the OT objective to be
 331 solved by the rapid Sinkhorn-Knopp algorithm [6].

332 **Boosting OT via Hardest Negatives.** To endow the trans-
 333 port plan with dual forgetting purposes, we reformulate
 334 Eq.(7) by imposing guidance from the hardest negative in-
 335 formation. Specifically, for each image v_i , we extend its
 336 transport target from $\{t_i\}_{i=1}^N$ to include its paired negative
 337 text t_i^{Neg} . As shown in Fig. 3, the negative text composes
 338 a new alignable column for the transport objective, which
 339 append the cost matrix C to $\bar{C} \in \mathbb{R}_+^{N \times (N+1)}$, i.e.,

$$340 [\bar{C}]_{i,N+1} = 1 - \langle v_i, t_i^{Neg} \rangle, [\bar{C}]_{i,j} = [C]_{i,j}, \forall i, j \in [1, N].$$

341 For the two parts \mathbb{D}_{FG} and \mathbb{D}_{RT} within each batch, the
 342 hardest negative should impose different guidance for dis-
 343 tinct unlearning goals. To this end, we propose a mask-
 344 base constraint to the corresponding transport plan $\hat{\Gamma}$ that
 345 regulates the effect of t_i^{Neg} . Specifically, the mask matrix

$M \in \mathbb{R}_+^{N \times (N+1)}$ satisfies that 346

$$347 [M]_{i,j} = \begin{cases} 0, & \text{if } (v_i, t_i) \in \mathbb{D}_{FG} \text{ and } j = i, \\ 0, & \text{if } (v_i, t_i) \in \mathbb{D}_{RT} \text{ and } j = N + 1, \\ 1, & \text{otherwise.} \end{cases} \quad (8)$$

348 For the toy example illustrated in Fig. 3, if the pair is consid-
 349 ered to be mismatched, the transport mass between v_i and
 350 t_i should be constrained to zero. Conversely, for the well-
 351 matched pair, t_i^{Neg} acts as a lower limit where the transport
 352 mass between v_i and t_j should be higher than it. Following
 353 the solver from [15], we model the mask constraint as the
 354 Hadamard product form that $\hat{\Gamma} = M \odot \bar{\Gamma}$, and the opti-
 355 mal alignment is formulated as (detailed Sinkhorn solution
 356 is presented in Supplementary B):

$$357 \hat{\Gamma}^* = \arg \min_{\hat{\Gamma} \in \Pi(\mu, \bar{\nu})} \langle \hat{\Gamma}, \bar{C} \rangle - \epsilon H(\hat{\Gamma}), \quad (9)$$

358 where $\bar{\nu} = \frac{1}{N+1} \mathbb{1}_{N+1}$ to satisfy the additional column.

359 **The Unlearning Objective.** Although $\hat{\Gamma}^*$ provides the
 360 more refined alignment, we suggest further incorporating
 361 an identity-like matrix I for two merits. First, diagonal el-
 362 ements are set as 1 for true positives to retain the initial
 363 alignment. Second, $[I]_{i,N+1} = 1$ to enhance the unlearning
 364 for the possible false positive $(v_i, t_i) \in \mathbb{D}_{FG}$. Thus, the
 365 overall alignment balanced by the factor γ is defined as:

$$366 T = \gamma \hat{\Gamma}^* + (1 - \gamma) I. \quad (10)$$

367 To fine-tune CLIP with this soft alignment, we use the KL
 368 divergence to optimize the matching distribution. Formally,
 369 we denote the batched similarity matrix as $P \in \mathbb{R}^{N \times (N+1)}$
 370 where $P_i = [\langle v_i, t_1 \rangle, \dots, \langle v_i, t_N \rangle, \langle v_i, t_i^{Neg} \rangle]^\top$. We ob-
 371 tain P_i^{v2t} and P_i^{t2v} by applying row-wise and column-wise
 372 softmax operation to P , respectively. Correspondingly, let
 373 T_i^{v2t} and T_i^{t2v} be the row-wise and column-wise normal-
 374 ized refined alignment for the i -th sample, respectively. The

375 OT-guided re-aligning is defined as:

$$376 \mathcal{L}^{otr} = \frac{1}{N} \sum_{i=1}^N \text{KL}(\mathbf{T}_i^{v2t} \| \mathbf{P}_i^{v2t}) + \frac{1}{N+1} \sum_{i=1}^{N+1} \text{KL}(\mathbf{T}_i^{t2v} \| \mathbf{P}_i^{t2v}). \quad (11)$$

377 Moreover, we empirically observe that preserving the tex-
378 tual semantic separation term can make the unlearning more
379 stable. Thus, the final unlearning objective is defined as:

$$380 \mathcal{L}^{UL} = \mathcal{L}^{otr} + \mathcal{L}^{sep}. \quad (12)$$

381 5. Experiment

382 In this section, we experimentally analyze the effectiveness
383 of NCU in unlearning the NC knowledge from CLIP.

384 5.1. Setup

385 **Datasets.** Our experiments are conducted on three vision-
386 language datasets at different scales and noise: Concep-
387 tual Captions 3M (CC3M) [39], Conceptual Captions 12M
388 (CC12M) [4], and YFCC15M-R (provided by [14], an
389 LLM-recaptioned subset from the YFCC100M [42]). All
390 datasets are web-crawled and contain an unknown portion
391 of NC pairs, *e.g.*, CC3M is estimated to include at least
392 3% false positives. We evaluate NCU on ImageNet and
393 15 common downstream datasets for classification perfor-
394 mance and on MSCOCO and Flickr30K for retrieval capa-
395 bility. Details for datasets are shown in Supplementary C.

396 **Unlearning Details.** Following CLIP, we consider two
397 architectures for the image encoder, *i.e.*, ViT/B16 and
398 ViT/B32, while the text encoder adopts the transformer ar-
399 chitecture. We consider a CLIP pre-trained on dataset \mathbb{D} ,
400 *e.g.*, CC3M, CC12M, or YFCC15M-R, as our reference
401 model, then we perform the NC unlearning on \mathbb{D} or its
402 subset to enhance CLIP’s robustness. For all experiments,
403 we allocate 2 epochs for learning negative semantics and 8
404 epochs for noise unlearning. All models are trained with a
405 batch size of 2,048 on 16 NVIDIA V100 GPUs. Detailed
406 training settings are presented in Supplementary D.

407 **Evaluation Protocol.** We evaluate NCU’s transferabil-
408 ity with Zero-Shot (ZS) classification accuracy and Lin-
409 ear Probing (LP) accuracy. For ZS classification, we fol-
410 low CLIP’s [37] prompt templates to compute distances be-
411 tween class text embeddings and image features. For LP,
412 we follow the mainstream setting [8, 37] that trains a lin-
413 ear classifier using L-BFGS on features extracted from the
414 frozen image encoder. Besides, we evaluate the retrieval
415 performance with the Recall at rank K (R@K) metric.

416 5.2. Evaluation on Diverse Downstream Tasks

417 To verify the generalization of NCU, we compare it with
418 CLIP on three different types of downstream tasks.

Zero-Shot Transfer. We compare the zero-shot perfor-
419 mance of CLIP and NCU on 16 popular image classifica-
420 tion datasets. We follow the prompt templates suggested in
421 the CLIP paper [37] to form each class name into a nat-
422 ural sentence. As demonstrated in Tab. 1, our NCU ap-
423 proach significantly outperforms the baseline CLIP model
424 on both ImageNet and other downstream datasets. Specif-
425 ically, across all fine-tuning datasets and all model archi-
426 tectures, NCU gains in the range of 2.8% ~ 4.1% in top-1
427 accuracy on ImageNet and 2.5% ~ 4.0% on average over
428 the other downstream datasets. This reveals that NCU can
429 successfully eliminate the impact of NC on CLIP by robust
430 fine-tuning with the same dataset. 431

Image-Text Retrieval. We present the zero-shot cross-
432 modal retrieval performance on the testing set of Flickr30K
433 (1K) and MSCOCO (5K) in Tab. 2. Our method consider-
434 ably outperforms the vanilla CLIP in almost all cases. For
435 instance, when fine-tuning CLIP (ViT-B/32) pre-trained on
436 the CC3M dataset, our NCU method achieves a 7.7% im-
437 provement in average recall scores on Flickr30K and 4.8%
438 improvement in average recall scores on MSCOCO. This
439 finding indicates that NCU can remarkably enhance the
440 alignment of images and text in the embedding space. 441

Linear Probing. Tab. 3 reports the linear probing perfor-
442 mance on 4 representative downstream datasets. Our NCU
443 consistently surpasses CLIP in the vast majority of cases,
444 suggesting that the visual embeddings learned by our NCU
445 are more effective and transferable than CLIP. 446

447 5.3. Compared to Robust Methods

448 In this section, we compare NCU with other robust-
449 designed techniques against NC on zero-shot ImageNet1K
450 classification task, *i.e.*, gradient ascent (GA), and SoftCLIP
451 [11]. Specifically, we evaluate GA as a standalone method,
452 where $-\mathcal{L}_{CLIP}$ is performed on \mathbb{D}_{FG} for handling FPs
453 and \mathcal{L}_{CLIP} with label smoothing is performed on \mathbb{D}_{RT} for
454 FNs. SoftCLIP is a noise-robust SOTA method that trains
455 CLIP from scratch by additional intra-modal guided align-
456 ment, *i.e.*, ROI features. As shown in Tab. 4, although GA
457 is a naive unlearning strategy, it still achieves observable
458 performance gains. Meanwhile, SoftCLIP’s self-similarity
459 modeling fails to excavate supervision from false positives,
460 which may explain why it performs worse than GA in
461 some cases, *e.g.*, CC12M with ViT-B/32. By contrast, our
462 NCU achieves solid improvements by forgetting both false
463 positives and false negatives, outperforming SoftCLIP by
464 1.1% ~ 2.3% without external guidance.

465 5.4. Ablation Study

466 To investigate the effectiveness of specific components in
467 our method, we carry out some ablation studies on Ima-

Dataset	Model	Caltech101	CIFAR-10	CIFAR-100	DTD	Aircraft	SST2	Flowers102	Food101	GTSRB	OxfordPets	RESISC45	SUN397	EuroSAT	StanfordCars	STL10	Average	ImageNet1K
<i>Model Architecture: ViT-B/16</i>																		
CC3M	CLIP	52.3	55.2	24.1	10.9	1.0	50.1	11.9	11.1	6.9	12.9	19.5	25.0	13.5	0.8	81.7	25.1	16.0
	NCU	59.1	54.3	28.8	12.3	1.1	50.1	14.1	14.8	7.4	16.3	22.8	32.3	21.7	1.5	86.3	28.2 ^{↑3.1}	20.0 ^{↑4.0}
CC12M	CLIP	77.0	66.5	38.3	21.2	2.5	47.7	33.4	51.9	7.3	64.2	39.0	44.7	21.2	25.5	91.4	42.1	40.6
	NCU	80.9	79.3	49.1	23.2	2.7	48.0	31.7	52.7	10.1	66.5	41.9	52.6	28.6	29.0	93.2	46.0 ^{↑3.9}	43.4 ^{↑2.8}
<i>Model Architecture: ViT-B/32</i>																		
CC3M	CLIP	47.7	54.2	18.0	7.6	1.2	50.1	9.3	9.1	6.0	7.4	16.2	16.0	15.5	0.8	77.7	22.5	11.8
	NCU	53.0	56.7	25.9	10.4	1.7	50.1	10.2	10.5	6.5	10.5	19.0	22.2	16.7	1.4	80.1	25.0 ^{↑2.5}	14.6 ^{↑2.8}
CC12M	CLIP	76.3	68.2	35.2	16.1	2.8	50.1	29.3	37.6	6.4	54.1	30.1	39.2	22.5	14.8	90.8	38.2	33.8
	NCU	80.4	68.5	41.4	19.3	2.6	52.8	28.6	43.4	7.2	62.4	35.7	48.3	31.3	18.1	92.5	42.2 ^{↑4.0}	36.7 ^{↑2.9}
YFCC15M-R	CLIP	53.7	67.0	34.4	13.1	1.1	49.3	22.1	18.6	11.0	13.5	20.3	29.3	23.0	1.7	83.7	29.5	17.8
	NCU	58.2	69.5	37.8	15.3	1.8	49.9	29.2	23.7	11.2	16.1	23.1	34.0	23.3	1.8	86.7	32.1 ^{↑2.6}	21.9 ^{↑4.1}

Table 1. Zero-shot transfer evaluation of different models.

Dataset	Architecture	Model	Flickr30K 1K Testing							MSCOCO 5K Testing								
			Image-to-Text			Text-to-Image				Average	Image-to-Text			Text-to-Image				Average
			R@1	R@5	R@10	R@1	R@5	R@10	R@1		R@5	R@10	R@1	R@5	R@10			
CC3M	ViT-B/16	CLIP	27.6	54.2	65.8	19.0	40.5	51.3	43.1	12.8	30.9	42.7	9.7	25.4	35.2	26.1		
		NCU	32.1	60.7	71.9	25.2	49.5	61.0	50.1 ^{↑7.0}	15.8	36.9	48.4	12.1	29.7	40.1	30.5 ^{↑4.4}		
	ViT-B/32	CLIP	14.0	35.2	47.7	11.5	27.9	37.8	29.0	7.0	20.1	28.9	6.0	16.7	24.0	17.1		
		NCU	21.3	44.5	55.7	15.7	36.3	46.4	36.7 ^{↑7.7}	10.7	25.7	35.4	8.1	21.3	30.2	21.9 ^{↑4.8}		
CC12M	ViT-B/32	CLIP	50.3	77.2	85.9	37.9	64.8	74.2	65.1	26.8	54.0	65.9	20.0	42.4	54.2	43.9		
		NCU	53.0	77.3	85.0	38.4	66.6	76.7	66.2 ^{↑1.1}	28.2	54.7	66.7	20.0	42.8	54.3	44.5 ^{↑0.6}		
YFCC15M-R	ViT-B/32	CLIP	57.4	81.6	89.1	40.8	66.4	75.4	68.5	35.0	60.9	71.8	22.6	46.0	58.0	49.1		
		NCU	58.0	83.2	89.7	42.5	69.9	78.8	70.4 ^{↑1.9}	34.3	62.0	73.5	24.6	49.0	60.8	50.7 ^{↑1.6}		

Table 2. Zero-shot cross-modal retrieval evaluation of different models.

Dataset	Architecture	Model	SUN397	OxfordPets	Food101	ImageNet
CC3M	ViT-B/16	CLIP	54.38	62.20	54.36	48.13
		NCU	55.60	62.85	54.74	49.90
	ViT-B/32	CLIP	46.96	52.30	46.96	40.20
		NCU	48.20	52.82	46.85	41.49
CC12M	ViT-B/16	CLIP	70.94	82.83	78.99	66.70
		NCU	71.36	84.76	79.18	66.65
	ViT-B/32	CLIP	66.05	78.41	70.12	59.19
		NCU	66.63	78.63	70.76	60.34
YFCC15M-R	ViT-B/32	CLIP	60.46	61.90	59.09	51.07
		NCU	60.75	62.85	60.46	52.29

Table 3. Linear probing comparison of different models.

468
469
470
471
472
473
474
475
476

geNet1K with models unlearned on CC3M. We first ablate the contributions of two key components of NCU, *i.e.*, negative prompt and text relation opposite. Specifically, in the variant \mathcal{V}_1 , we replace the learnable prompt tokens by prepending some textual negative prefixes to raw captions, *e.g.*, ‘the image has no’ or ‘this picture lacks’. In the variant \mathcal{V}_2 , we use maximal L_2 distance loss [43] as a substitute for our text relation opposite. Besides, we validate the impact of different NC unlearning by intervening with the refined

Dataset	Model	Model Architecture	ImageNet1K ZS top-1
CC3M	CLIP	ViT-B/16	16.0
	Gradient Ascent		16.7 ^{↑0.7}
	SoftCLIP		18.9 ^{↑2.9}
	NCU		20.0 ^{↑4.0}
CC3M	CLIP	ViT-B/32	11.8
	Gradient Ascent		12.1 ^{↑0.3}
	SoftCLIP		13.3 ^{↑1.5}
	NCU		14.6 ^{↑2.8}
CC12M	CLIP	ViT-B/16	40.6
	Gradient Ascent		41.6 ^{↑1.0}
	SoftCLIP		42.1 ^{↑1.5}
	NCU		43.4 ^{↑2.8}
CC12M	CLIP	ViT-B/32	33.8
	Gradient Ascent		35.1 ^{↑1.3}
	SoftCLIP		34.4 ^{↑0.6}
	NCU		36.7 ^{↑2.9}

Table 4. Zero-shot top-1 performance on ImageNet1K.

alignment, *i.e.*, \mathcal{V}_3 and \mathcal{V}_4 . As shown in Tab. 5, we observe that: 1) Using negative textual prefixes also shows competitive results, demonstrating the generalization of our method. However, we argue that the learnable prompt is preferable, except for performance gains, operating to features makes

477
478
479
480
481

Model	ImageNet1K ZS top-1	
	ViT-B/16	ViT-B/32
NCU	20.0	14.6
\mathcal{V}_1 (w/o Hardest Negative Prompts)	19.3 $\downarrow 0.7$	14.4 $\downarrow 0.2$
\mathcal{V}_2 (w/o Text Relation Opposite)	18.8 $\downarrow 1.2$	13.9 $\downarrow 0.7$
\mathcal{V}_3 (w only False Negatives Unlearning)	17.7 $\downarrow 2.3$	13.5 $\downarrow 1.1$
\mathcal{V}_4 (w only False Positives Unlearning)	19.2 $\downarrow 0.8$	14.3 $\downarrow 0.3$

Table 5. Ablation studies on zero-shot transfer task (ImageNet1K) of models unlearned on CC3M.

482 it possible for NCU to extend to modalities beyond text.
 483 2) Simply maximizing the distance between negative and
 484 original text embeddings leads to suboptimal performance,
 485 which aligns with our analysis in Fig. 2. 3) Both types of
 486 NC impair CLIP’s performance, among which the false posi-
 487 tive causes a more severe impact. While NCU achieves the
 488 best performance by forgetting such noisy knowledge.

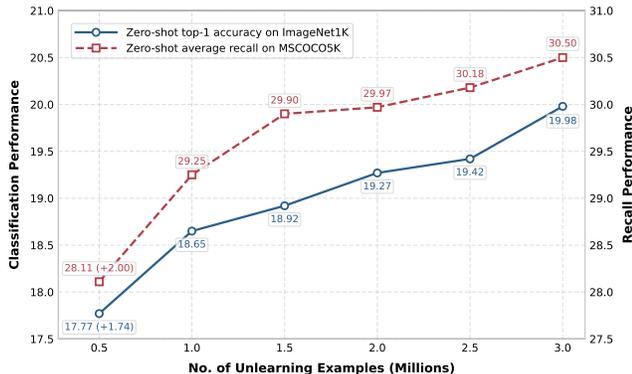


Figure 4. Effect of NCU with varying fine-tuning dataset sizes on zero-shot image classification and cross-modal retrieval.

489 **5.5. NC Unlearning with Partial Data**

490 In this section, we conduct an interesting study to verify
 491 whether CLIP can improve robustness by only unlearning
 492 NC with a portion of the pre-trained data. To this end, given
 493 a CLIP pre-trained on CC3M as the reference model, we
 494 evaluate NCU using different data portions ranging from
 495 0.5 million to 3 million image-text pairs. Fig. 4 plots the ZS
 496 top-1 accuracy on ImageNet1K and the average of recalls on
 497 MSCOCO 5K. Remarkably, even when unlearning on less
 498 than 20% of the original data (0.5M), NCU achieves sig-
 499 nificant performance gains while preserving overall knowl-
 500 edge learned in CC3M. With the accessible data increas-
 501 ing, NCU shows consistent improvements on both zero-shot
 502 downstream tasks. This phenomenon indicates NCU’s flex-
 503 ibility in enhancing the robustness of models with limited
 504 data, which is valuable to handling VLMs pre-trained with
 505 partially private or proprietary data.

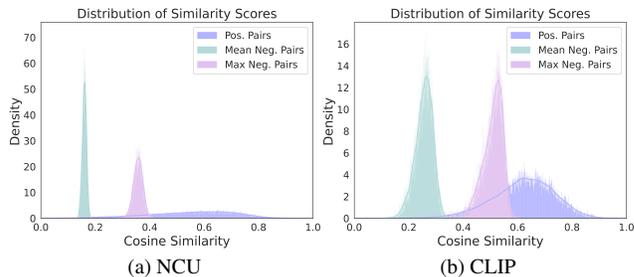


Figure 5. Similarity scores distribution of positive and negative pairs from CLIP and NCU. Both models are based on ViT/B16 and learning from the CC3M training set.

506 **5.6. Visualization and Analysis**

507 To intuitively show the robust embedding space that is re-
 508 fined by our approach, we plot the distribution of normal-
 509 ized similarity for CLIP and NCU on the validation set of
 510 CC3M. In Fig. 5, we illustrate similarity scores for positive
 511 pairs, mean of negative pairs, and top 5% maximum of neg-
 512 ative pairs. First, we observe that NCU produces a wider
 513 distribution of positive similarity scores, capturing more
 514 fine-grained matching degrees among positive pairs. Sec-
 515 ond, NCU improves the feature discrimination, which leads
 516 to a more significant separation between positive and neg-
 517 ative pairs. Lastly, NCU provides more appropriate mea-
 518 sures for hard negatives, which maintains separation from
 519 both positive and other negative pairs.

520 **6. Limitations and Future Works**

521 Our work still has certain limitations due to the finite com-
 522 puting capability, including 1) This work only uses CLIP
 523 to explore the efficacy of NCU. Further research is needed
 524 to confirm its applicability in other VLMs, such as BLIP-
 525 2 [28], and even larger VLMs like VisionLLM [45] or In-
 526 ternVL [5]. 2) The current experiments are mainly con-
 527 ducted on million-scale data, and we plan to extend it to
 528 larger-scale datasets to verify NCU’s generalization.

529 **7. Conclusion**

530 This work provides a new thinking in robust vision-
 531 language learning. Instead of re-training models from
 532 scratch, we suggest eliminating the harmful effects of noisy
 533 correspondence from pre-trained models. To this end, we
 534 propose NCU, a robust fine-tuning framework that effi-
 535 ciently unlearns noisy correspondence in CLIP. Our key
 536 concept is to learn the hardest negative information that
 537 can provide explicit unlearning direction to resist both FP
 538 and FN. We formalize such twin goals unlearning process
 539 into one unified OT problem for fast fine-tuning. Extensive
 540 experiments are conducted to verify that NCU can endow
 541 CLIP with strong robustness against noisy correspondence.

542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599

References

[1] Alex Andonian, Shixing Chen, and Raffay Hamid. Robust cross-modal representation learning with progressive self-distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16430–16441, 2022. 1, 2

[2] Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE symposium on security and privacy (SP)*, pages 141–159. IEEE, 2021. 2

[3] Adrian Bulat, Yassine Ouali, and Georgios Tzimiropoulos. Fff: Fixing flawed foundations in contrastive pre-training results in very strong vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14172–14182, 2024. 1, 2

[4] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3558–3568, 2021. 6

[5] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198, 2024. 8

[6] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013. 5

[7] Chongyu Fan, Jiancheng Liu, Yihua Zhang, Dennis Wei, Eric Wong, and Sijia Liu. Salun: Empowering machine unlearning via gradient-based weight saliency in both image classification and generation. In *International Conference on Learning Representations*, 2024. 2

[8] Lijie Fan, Dilip Krishnan, Phillip Isola, Dina Katabi, and Yonglong Tian. Improving clip training with language rewrites. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 6

[9] Jack Foster, Stefan Schoepf, and Alexandra Brintrup. Fast machine unlearning without retraining through selective synaptic dampening. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 12043–12051, 2024. 2

[10] Zhongtian Fu, Kefei Song, Luping Zhou, and Yang Yang. Noise-aware image captioning with progressively exploring mismatched words. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 12091–12099, 2024. 2

[11] 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*, volume 38, pages 1860–1868, 2024. 1, 2, 6

[12] 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

[13] 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

[14] Tiancheng Gu, Kaicheng Yang, Xiang An, Ziyong Feng, Dongnan Liu, Weidong Cai, and Jiankang Deng. Rwkv-clip: A robust vision-language representation learner. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 4799–4812, 2024. 2, 6

[15] Xiang Gu, Yucheng Yang, Wei Zeng, Jian Sun, and Zongben Xu. Keypoint-guided optimal transport with applications in heterogeneous domain adaptation. *Advances in Neural Information Processing Systems*, 35:14972–14985, 2022. 5

[16] Haochen Han, Qinghua Zheng, Guang Dai, Minnan Luo, and Jingdong Wang. Learning to rematch mismatched pairs for robust cross-modal retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26679–26688, 2024. 2

[17] Tengda Han, Weidi Xie, and Andrew Zisserman. Temporal alignment networks for long-term video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2906–2916, 2022. 2

[18] Peng Hu, Zhenyu Huang, Dezhong Peng, Xu Wang, and Xi Peng. Cross-modal retrieval with partially mismatched pairs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 2, 3

[19] Runhui Huang, Yanxin Long, Jianhua Han, Hang Xu, Xiwen Liang, Chunjing Xu, and Xiaodan Liang. Nlip: Noise-robust language-image pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 926–934, 2023. 2

[20] Zhenyu Huang, Guocheng Niu, Xiao Liu, Wenbiao Ding, Xinyan Xiao, Hua Wu, and Xi Peng. Learning with noisy correspondence for cross-modal matching. *Advances in Neural Information Processing Systems*, 34:29406–29419, 2021. 1, 2

[21] Zhenyu Huang, Mouxing Yang, Xinyan Xiao, Peng Hu, and Xi Peng. Noise-robust vision-language pre-training with positive-negative learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 1, 2

[22] 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. 1, 2

[23] Liang Jiang, Zhenyu Huang, Jia Liu, Zujie Wen, and Xi Peng. Robust domain adaptation for machine reading comprehension. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 8060–8069, 2023. 2

[24] Wooyoung Kang, Jonghwan Mun, Sungjun Lee, and Byungseok Roh. Noise-aware learning from web-crawled image-text data for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2942–2952, 2023. 2

[25] Youngdong Kim, Juseung Yun, Hyounguk Shon, and Junmo Kim. Joint negative and positive learning for noisy labels. In *Proceedings of the IEEE/CVF conference on computer vi-*

600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658

659 *sion and pattern recognition*, pages 9442–9451, 2021. 3

660 [26] Alexey Kravets and Vinay Nambodiri. Zero-shot class un- 718
661 learning in clip with synthetic samples. In *Proceedings of the* 719
662 *IEEE/CVF Winter Conference on Applications of Computer* 720
663 *Vision (WACV)*, 2025. 2 721
664 [27] Alexey Kravets and Vinay P. Nambodiri. Zero-shot clip 722
665 class forgetting via text-image space adaptation. *Transac-* 723
666 *tions on Machine Learning Research*, 2025. 2 724
667 [28] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 725
668 Blip-2: Bootstrapping language-image pre-training with 726
669 frozen image encoders and large language models. In *Inter-* 727
670 *national conference on machine learning*, pages 19730– 728
671 19742. PMLR, 2023. 8 729
672 [29] Tianqi Li, Guansong Pang, Xiao Bai, Wenjun Miao, and Jin 730
673 Zheng. Learning transferable negative prompts for out-of- 731
674 distribution detection. In *Proceedings of the IEEE/CVF Con-* 732
675 *ference on Computer Vision and Pattern Recognition*, pages 733
676 17584–17594, 2024. 3, 4 734
677 [30] Yijie Lin, Jie Zhang, Zhenyu Huang, Jia Liu, Xi Peng, et al. 735
678 Multi-granularity correspondence learning from long-term 736
679 noisy videos. In *The Twelfth International Conference on* 737
680 *Learning Representations*, 2024. 2 738
681 [31] Chris Yuhao Liu, Yaxuan Wang, Jeffrey Flanigan, and Yang 739
682 Liu. Large language model unlearning via embedding- 740
683 corrupted prompts. In *The Thirty-eighth Annual Conference* 741
684 *on Neural Information Processing Systems*, 2024. 2 742
685 [32] Ronak Mehta, Sourav Pal, Vikas Singh, and Sathya N Ravi. 743
686 Deep unlearning via randomized conditionally independent 744
687 Hessians. In *Proceedings of the IEEE/CVF Conference on* 745
688 *Computer Vision and Pattern Recognition*, pages 10422– 746
689 10431, 2022. 2 747
690 [33] Norman Mu, Alexander Kirillov, David Wagner, and Sain- 748
691 ing Xie. Slip: Self-supervision meets language-image pre- 749
692 training. In *European conference on computer vision*, pages 750
693 529–544. Springer, 2022. 2 751
694 [34] Maitreya Patel, Sheng Cheng, Changhoon Kim, Tejas 752
695 Gokhale, Chitta Baral, Yezhou Yang, et al. Tripletclip: Im- 753
696 proving compositional reasoning of clip via synthetic vision- 754
697 language negatives. In *The Thirty-eighth Annual Conference* 755
698 *on Neural Information Processing Systems*, 2024. 1 756
699 [35] Yang Qin, Yingke Chen, Dezhong Peng, Xi Peng, 757
700 Joey Tianyi Zhou, and Peng Hu. Noisy-correspondence 758
701 learning for text-to-image person re-identification. In *Pro-* 759
702 *ceedings of the IEEE/CVF Conference on Computer Vision* 760
703 *and Pattern Recognition*, pages 27197–27206, 2024. 2 761
704 [36] Yang Qin, Yuan Sun, Dezhong Peng, Joey Tianyi Zhou, 762
705 Xi Peng, and Peng Hu. Cross-modal active complemen- 763
706 tary learning with self-refining correspondence. *Advances* 764
707 *in Neural Information Processing Systems*, 36, 2024. 2 765
708 [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 766
709 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, 767
710 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning 768
711 transferable visual models from natural language supervi- 769
712 sion. In *International conference on machine learning*, pages 770
713 8748–8763. PMLR, 2021. 1, 2, 6 771
714 [38] Christoph Schuhmann, Robert Kaczmarczyk, Aran Komat- 772
715 suzaki, Aarush Katta, Richard Vencu, Romain Beaumont, 773
716 Jenia Jitsev, Theo Coombes, and Clayton Mullis. Laion- 774
717 400m: Open dataset of clip-filtered 400 million image-text 775
776 pairs. In *NeurIPS Workshop Datacentric AI*, number FZJ- 776
2022-00923. Jülich Supercomputing Center, 2021. 1, 2

[39] 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. 1, 6

[40] Liangliang Shi, Jack Fan, and Junchi Yan. Ot-clip: Understanding and generalizing clip via optimal transport. In *Forty-first International Conference on Machine Learning*, 2024. 2

[41] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation. In *Conference on robot learning*, pages 894–906. PMLR, 2022. 1

[42] 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. *Communications of the ACM*, 59(2):64–73, 2016. 2, 6

[43] Hualiang Wang, Yi Li, Huifeng Yao, and Xiaomeng Li. Clipn for zero-shot ood detection: Teaching clip to say no. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1802–1812, 2023. 3, 4, 7

[44] Jue Wang, Haofan Wang, Weijia Wu, Jincan Deng, Yu Lu, Xiaofeng Guo, and Debing Zhang. Eclip: Efficient contrastive language-image pretraining via ensemble confidence learning and masked language modeling. In *First Workshop on Pre-training: Perspectives, Pitfalls, and Paths Forward at ICML 2022*, 2022. 2

[45] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *Advances in Neural Information Processing Systems*, 36:61501–61513, 2023. 8

[46] Bichen Wu, Ruizhe Cheng, Peizhao Zhang, Tianren Gao, Joseph E Gonzalez, and Peter Vajda. Data efficient language-supervised zero-shot recognition with optimal transport distillation. In *International Conference on Learning Representations*, 2022. 2

[47] Mouxing Yang, Zhenyu Huang, and Xi Peng. Robust object re-identification with coupled noisy labels. *International Journal of Computer Vision*, pages 1–19, 2024. 2

[48] 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. In *International Conference on Learning Representations*, 2022. 2

[49] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, 2023. 4

[50] Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*, 2024. 2

[51] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022. 3