# ONE PERTURBATION IS ENOUGH: ON GENERATING UNIVERSAL ADVERSARIAL PERTURBATIONS AGAINST VISION-LANGUAGE PRE-TRAINING MODELS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Vision-Language Pre-training (VLP) models have exhibited unprecedented capability in many applications by taking full advantage of the multimodal alignment. However, previous studies have shown they are vulnerable to maliciously crafted adversarial samples. Despite recent success, these methods are generally instancespecific and require generating perturbations for each input sample. In this paper, we reveal that VLP models are also vulnerable to the instance-agnostic universal adversarial perturbation (UAP). Specifically, we design a novel Contrastivetraining Perturbation Generator with Cross-modal conditions (C-PGC) to achieve the attack. In light that the pivotal multimodal alignment is achieved through the advanced contrastive learning technique, we devise to turn this powerful weapon against themselves, i.e., employ a malicious version of contrastive learning to train the C-PGC based on our carefully crafted positive and negative image-text pairs for essentially destroying the alignment relationship learned by VLP models. Besides, C-PGC fully utilizes the characteristics of Vision-and-Language (V+L) scenarios by incorporating both unimodal and cross-modal information as effective guidance. Extensive experiments show that C-PGC successfully forces adversarial samples to move away from their original area in the VLP model's feature space, thus essentially enhancing attacks across various victim models and V+L tasks.

#### 028 029 030 031

032

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027

#### 1 INTRODUCTION

Vision-Language Pre-training (VLP) models, including ALBEF (Li et al., 2021), TCL (Yang et al., 2022), and BLIP (Li et al., 2022), have recently demonstrated remarkable efficacy in a wide range of Vision-and-Language (V+L) tasks. By self-supervised pre-training on large-scale image-text pairs, VLP models efficiently align cross-modal features and capture rich information from the aligned multimodal embeddings, thereby providing expressive representations for various applications.

Adversarial attacks (Carlini & Wagner, 2017), which aim to deceive models during inference time, have attracted extensive attention due to their significant threat to security-critical scenarios (Eykholt et al., 2018). Recent studies have shown that VLP models are also vulnerable to adversarial samples. 040 The pioneering work Co-Attack (Zhang et al., 2022) proposes the first multimodal attack that simul-041 taneously perturbs both image and text modalities and displays excellent performance. However, 042 Co-Attack only considers relatively easier white-box attacks where victim models are completely 043 accessible. To handle more practical black-box settings, subsequent studies propose various transfer-044 able adversarial samples generated on an available surrogate model to fool other inaccessible models. Specifically, SGA (Lu et al., 2023) significantly improves the adversarial transferability through 046 the set-level cross-modal guidance obtained from data augmentations. Subsequently, TMM (Wang 047 et al., 2024) proposes to jointly destroy the modality-consistency features within the clean image-048 text pairs and include more modality-discrepancy features in the perturbations to further enhance transferability. While existing methods have achieved great success, they are all instance-specific and need to generate a perturbation for each input pair, which results in substantial computational 051 overhead. Meanwhile, universal adversarial attacks, as an efficient instance-agnostic approach that uses only one Universal Adversarial Perturbation (UAP) to conduct attacks, have not been fully in-052 vestigated for VLP models. This naturally leads to a question, is it possible to design a universal adversarial perturbation that can effectively deceive VLP models across various image-text pairs?





Figure 2: Illustration of the universal adversarial attacks against VLP models. With only a pair of image-text perturbations, the proposed attack can effectively mislead different models on diverse V+L tasks.

074 Motivation. To this end, we make an intuitive attempt to transplant existing renowned approaches 075 UAP (Moosavi-Dezfooli et al., 2017) and GAP (Poursaeed et al., 2018) to launch attacks on several 076 VLP models by maximizing the distance between the embeddings of the adversarial image and its 077 matched texts. Unfortunately, Figure 1 demonstrates that these methods yield unsatisfactory attack success rates (ASR), especially for black-box attacks. Empirically, this failure stems from their 079 narrow focus on the image modal, disregarding the other modality and the multimodal information that plays a pivotal role in VLP models. To overcome this challenge, we revisit the VLP models' 081 basic training paradigm and emphasize that regardless of the downstream V+L tasks, their achieved outstanding performance is heavily reliant on the well-established multimodal alignment, which draws the embedding of matched image-text pairs closer while distancing those of non-matched 083 pairs. In light of this consideration, we argue that the key core of an effective universal adversarial 084 attack is to obtain a UAP that can fundamentally destroy this learned alignment relationship to 085 mislead VLP models into making incorrect decisions. Besides, Fig. 1 also shows that the generatorbased GAP consistently outperforms UAP methods, confirming the superiority of the generative 087 paradigm, which is also corroborated by numerous studies (Gao et al., 2024; Feng et al., 2023). 088

Based on these insights, we propose a novel generative framework that learns a Contrastive-training 089 Perturbation Generator with Cross-modal conditions (C-PGC) to launch universal attacks on VLP 090 models (see Fig. 2). To essentially destroy the multimodal alignment, we devise to utilize VLP mod-091 els' most powerful weapons to attack against themselves, i.e., use the contrastive learning mecha-092 nism to train the generator using our maliciously constructed image-text pairs that completely violate the correct VL matching relationship, to produce perturbation that pushes the embedding of 094 matched pairs apart while pulling those of non-matched ones together. Inspired by the multimodal 095 characteristics of V+L scenarios, we modify the generator's architecture to incorporate cross-modal 096 knowledge through the advanced cross-attention mechanism for better guidance. In addition, we also consider the intra-modal influence and introduce an unimodal distance loss to further enhance 098 the attacks. Since previous studies (Zhang et al., 2022; Lu et al., 2023) achieve impressive improvements via multimodal perturbation, we are motivated to generate UAP for both images and texts to 099 utilize the synergy between different modalities. Our contributions can be summarized as follows: 100

101

054

056

059

061

063

064

067

069

071

073

102 103

104

105

106

• We design a novel cross-modal conditional perturbation generator, which produces effective UAP for both image and text modalities to achieve universal adversarial attacks on VLP models.

- We propose the first malicious contrastive paradigm tailored for multimodal adversarial attacks, which incorporates both unimodal and multimodal guidance to contrastively train the generator using our meticulously constructed positive and negative pairs for enhanced attack effects.
- Extensive experiments on 6 various VLP models across different V+L tasks reveal that our method 107 achieves outstanding white-box performance and black-box transferability in different scenarios.

## 108 2 RELATED WORK

# 110 2.1 VISION-LANGUAGE PRE-TRAINING MODELS

VLP models are pre-trained on massive image-text pairs to learn the semantic correlations across
modalities and serve diverse multimodal user demands (Chen et al., 2023; Du et al., 2022). We next
illustrate the basis of VLP models from multiple perspectives.

Architectures. Based on the ways of multimodal fusion, the architectures of VLP models can be classified into two types: *single-stream* and *dual-stream architectures*. Single-stream architectures (Li et al., 2019; Chen et al., 2020) directly concatenate the text and image features, and calculate the attention in the same Transformer block for multimodal fusion. On the contrary, dual-stream architectures (Radford et al., 2021; Li et al., 2022) separately feed the text and image features to different Transformer blocks and leverage the cross-attention mechanism for multimodal fusion.

Pre-training Objectives. The pre-training objectives for VLP models mainly include *masked features completion, multimodal features matching,* and *specific downstream objectives*. Masked features completion (Chen et al., 2020) encourages VLP models to predict the deliberately masked tokens using the remaining unmasked tokens during pre-training. Multimodal features matching (Li et al., 2021) pre-trains VLP models by learning to precisely predict whether the given image-text pairs are matched. Specific downstream objectives (Anderson et al., 2018) directly utilize the training objectives of downstream tasks (e.g., visual question answering) for pre-training VLP models.

Downstream Tasks. In this paper, we mainly consider the following multimodal downstream tasks:
(1) Image-text retrieval (ITR) (Wang et al., 2016): finding the most matched image for the given text and vice versa, including image-to-text retrieval (TR) and text-to-image retrieval (IR). (2) Image caption (IC) (Bai & An, 2018): generating the most suitable descriptions for the given image. (3) Visual grounding (VG) (Hong et al., 2019): locating specific regions in the image that correspond with the given textual descriptions. (4) Visual entailment (VE) (Xie et al., 2019): analyzing the input image and text and predicting whether their relationship is entailment, neutral, or contradiction.

135

#### 136 2.2 Adversarial Attacks

137 **Instance-specific Attacks on VLP Models.** The adversarial robustness of VLP Models has already 138 become a research focus. Early works (Kim & Ghosh, 2019; Yang et al., 2021) impose perturbations 139 only on single modality and lack cross-modal interactions when attacking multimodal models. To 140 address this issue, Co-Attack (Zhang et al., 2022) conducts the first multimodal white-box attacks on 141 VLP models. On the basis of Co-Attack, Lu et al. (2023) extend the attacks to more rigorous black-142 box settings and propose SGA, which utilizes set-level alignment-preserving argumentations with 143 carefully designed cross-modal guidance. However, Wang et al. (2024) points out that SGA fails to 144 fully exploit modality correlation, and proposes TMM to better leverage cross-modal interactions by tailoring both the modality-consistency and modality-discrepancy features. Nonetheless, these 145 methods are all instance-specific and need to craft perturbations for each input pair. 146

147 **Universal Adversarial Examples.** Universal adversarial attacks (Moosavi-Dezfooli et al., 2017; 148 Mopuri et al., 2018) aim to deceive the victim model by exerting a uniform adversarial modification 149 to all the benign samples. These attacks save the redundant procedures of redesigning perturbations for each input sample and are consequently more efficient than traditional attack strategies. Gen-150 erally, universal adversarial attacks can be categorized into optimization-based methods (Moosavi-151 Dezfooli et al., 2017; Wang et al., 2023; Liu et al., 2023) and generation-based methods (Hayes & 152 Danezis, 2018; Gao et al., 2024; Anil et al., 2024). Benefiting from the powerful modeling abilities 153 of generative models, generation-based methods are more versatile and can produce more natural 154 samples than optimization-based ones. In this paper, we explore universal adversarial attacks on 155 VLP models and manage to generate UAP with excellent attack effects and high transferability. 156

150

## 3 UNIVERSAL MULTIMODAL ATTACKS

158 159

In this section, we first present the problem statement of universal adversarial attacks on VLP mod els. Next, we introduce the overview of our framework. Finally, we illustrate the detailed design of
 each proposed technique and summarize the training objective and paradigm of C-PGC.

163 164 165

166

167

162

169 170 171

172

173 174

175

176

177 178

179

181 182 183

184

191

192



Figure 3: An overview of our proposed universal adversarial attack. Benefiting from the welldesigned unimodal distance loss  $\mathcal{L}_{Dis}$  and multimodal contrastive loss  $\mathcal{L}_{CL}$ , the conditional generator learns rich knowledge from features of different modalities and thus produces  $\delta_v$  and  $\delta_t$  of superior generalization ability across diverse models and downstream tasks.

#### 3.1 PROBLEM STATEMENT

We define an input image-text pair as (v, t) and denote  $e_v$  and  $e_t$  as the image and text embedding encoded by the image encoder  $f_I(\cdot)$  and text encoder  $f_T(\cdot)$  of the targeted VLP model  $f(\cdot)$ . Let  $\mathcal{D}_s$ be an available dataset consisting of image-text pairs collected by a malicious adversary. The attack objective is to utilize  $\mathcal{D}_s$  to train a generator  $G_w(\cdot)$  that is capable of producing a powerful pair of universal image-text perturbations  $(\delta_v, \delta_t)$  that can affect the vast majority of test dataset  $\mathcal{D}_t$  to fool models into making incorrect decisions. Formally, the attack goal can be formulated as:

$$\mathcal{T}(f(v+\delta_v, t\oplus \delta_t)) \neq y, \text{ s.t. } \|\delta_v\|_{\infty} \leq \epsilon_v, \|\delta_t\|_0 \leq \epsilon_t, \tag{1}$$

193 where  $\mathcal{T}(\cdot)$  denotes the operation that uses the output V+L features to obtain the final predictions,  $\oplus$  indicates the text perturbation strategy (Zhang et al., 2022; Lu et al., 2023) that replaces certain 194 important tokens of the original sentence with crafted adversarial words, and y is the correct pre-195 diction of the considered V+L task. To ensure the perturbation's imperceptibility, we constrain the 196 pixel-level image perturbation with  $l_{\infty}$  norm of a given budget  $\epsilon_v$ . The textual perturbation is token-197 level and the stealthiness is accordingly constrained by the number of modified words  $\epsilon_t$ . Since altering words in a natural sentence can be easily noticed or detected, we apply a rigorous restriction 199 that permits only a single word to be substituted ( $\epsilon_t = 1$ ). On the premise of imperceptibility, the 200 attacker attempts to generalize the crafted UAP to a wider range of test data and victim models.

201 202 203

#### 3.2 OVERVIEW OF THE PROPOSED FRAMEWORK

As depicted in Fig. 3, we adopt the multimodal perturbation strategy and generate perturbations on both image and text modalities for enhanced attacks. Given the similarity between the workflows for image and text, we then take the image attacks as an example to illustrate the proposed framework.

207 Firstly, a fixed noise  $z_v$  is randomly initialized and subsequently fed into the conditional generator. 208 For each image v and its descriptions t, the generator  $G_w(\cdot)$  then translates the input noise  $z_v$  into the 209 adversarial perturbation  $\delta_v$  that is of the same size as the image v. During the generation, the network 210  $G_w$  additionally benefits from cross-modal information by integrating the textual embedding, i.e., 211  $\delta_v = G_w(z_v; f_T(\mathbf{t}))$ . Next, the generated adversarial noise  $\delta_v$  is injected into the clean image to 212 obtain the adversarial image via  $v_{adv} = v + \delta_v$ . To better guide the training process, we design 213 two efficient unimodal and multimodal losses as our optimization objectives. Unimodal loss is straightforward and aims to push the adversarial images away from the clean images in the latent 214 embedding space, while multimodal loss is based on contrastive learning and uses our manually 215 constructed positive and negative samples to effectively destroy the image-text matching relationship

obtained from feature alignment. Once we finish training the C-PGC using the proposed losses, the
 input fixed noise is transformed into a UAP that is of great generalization and transferability.

# 219 3.3 DETAILED DESIGN OF C-PGC

Next, we provide a detailed introduction to each of the proposed designs. Note that we primarily discuss the image attack as an example, given that the design of the text attack is completely symmetrical. The pseudocode of the training procedure is provided in Appendix A.

224 Perturbation Generator with Cross-modal conditions. Previous generative universal attacks (Gao 225 et al., 2024; Anil et al., 2024) have shown excellent efficacy in fooling the discriminative models. 226 Nevertheless, since existing generative attacks are limited to a single modality, directly utilizing 227 the off-the-shelf generators might fail to leverage the multimodal interactions in these special V+L 228 scenarios. To address this limitation, we additionally introduce cross-modal embeddings as auxiliary information to further facilitate the process of perturbation generation. Specifically, we modify the 229 existing generator's architecture by adding several cross-attention modules that have been proven 230 effective in tasks with variable input modalities. The obtained textual embeddings  $e_t$  encoded by 231  $f_T(\cdot)$  are then incorporated into our generator through: 232

$$Q = \mathbf{h}_t W_q, K = \mathbf{e}_t W_k, V = \mathbf{e}_t W_v,$$
  
Attention $(Q, K, V) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V,$  (2)

where  $h_t \in \mathbb{R}^{B \times d_\alpha}$  is the flattened intermediate features within the generator,  $W_q \in \mathbb{R}^{d_\alpha \times d}$ ,  $W_k \in \mathbb{R}^{512 \times d}$ ,  $W_v \in \mathbb{R}^{512 \times d}$  are optimized parameters in the attention modules.

 Multimodal Contrastive Loss. The preceding analysis regarding the failures of existing UAP attacks encourages us to design a loss function that can guide the generated UAP to break the learned multimodal feature alignment. Motivated by the fact that contrastive learning underpins the crossmodal alignment, we advocate leveraging this mechanism to attack VLP models themselves by contrastively training our C-PGC to essentially disrupt the benign alignment relationship. Concretely, we adopt the widely recognized InfoNCE (He et al., 2020) as our basic contrastive loss.

245 To establish the contrastive paradigm, we first define the adversarial image  $v_{adv}$  as the anchor sam-246 ple. Besides, it is also necessary to construct an appropriate set of positive and negative samples. 247 Based on the fundamental objective of our attacks, it is natural that the originally matched text 248 descriptions set  $\mathbf{t} = \{t_1, t_2, \dots, t_M\}$  can be employed as negative samples  $\mathbf{t}_{neg}$  to increase the dis-249 crepancy of matched image-text in the feature space of VLP models. To further push the adversarial 250 image  $v_{adv}$  away from its correct text descriptions t, we propose a *farthest selection strategy* which 251 utilizes multiple texts whose embeddings differ significantly from that of the original clean image 252 v as positive samples. Specifically, we randomly sample a batch of text sets from  $\mathcal{D}_s$  and select the text set with the largest feature distances from the current image v as positive samples, i.e., 253  $\mathbf{t}_{pos} = \{t'_1, t'_2, \dots, t'_K\}$ . Moreover, we utilize data augmentations that resize the clean v into diverse 254 scales and apply random Gaussian noise to acquire a more diverse image set  $\mathbf{v} = \{v_1, v_2, \dots, v_S\}$  for 255 set-level guidance Lu et al. (2023). With the augmented images and these well-constructed positive 256 and negative samples, the multimodal contrastive loss  $\mathcal{L}_{CL}$  can be formulated as: 257

233 234

235 236

$$\mathcal{L}_{CL} = log \left( \frac{\sum_{i=1}^{S} \sum_{j=1}^{M} e^{d(v_i + \delta_v, t_j)/\tau}}{\sum_{i=1}^{S} \sum_{j=1}^{M} e^{d(v_i + \delta_v, t_j)/\tau} + \sum_{i=1}^{S} \sum_{j=1}^{K} e^{d(v_i + \delta_v, t_j')/\tau}} \right),$$
(3)

where  $\delta_v$  is the universal image perturbation,  $\tau$  denotes the temperature parameter and  $d(v,t) = Sim(f_I(v), f_T(t))$ , where  $Sim(\cdot, \cdot)$  represents the cosine similarity measurement.

266 Unimodal Distance Loss. Apart from the multimodal guidance, we also consider the unimodal in-267 fluence by directly pushing adversarial images away from their initial visual semantic area to further 268 enhance attack effects. Similarly, to acquire set-level guidance, the input image v is initially resized 269 to different scales and then added with Gaussian noise to generate the augmented image set  $\mathbf{v} = \{v_1, v_2, \dots, v_S\}$ . Then, we craft the adversarial image through  $v_{adv} = v + \delta_v$  and process  $v_{adv}$  with the same augmentation operation to obtain the adversarial image set  $\mathbf{v}_{adv} = \{v_1^{adv}, v_2^{adv}, \dots, v_S^{adv}\}$ . Finally, we minimize the negative Euclidean distance between the embeddings of adversarial images and clean images to optimize the UAP generator. Formally, the loss  $\mathcal{L}_{Dis}$  is formulated as:

$$\mathcal{L}_{Dis} = -\sum_{i=1}^{S} \sum_{j=1}^{S} \|f_I(v_i^{adv}) - f_I(v_j)\|_2.$$
(4)

274 275

281 282

Taking advantage of the unimodal set-level guidance,  $\mathcal{L}_{Dis}$  ensures an effective optimization direction during the generator training and further improves the attack effectiveness of our UAP.

**Training Objective.** With the two well-designed loss terms  $\mathcal{L}_{Dis}$  and  $\mathcal{L}_{CL}$ , the overall optimization objective of our conditional generator concerning image attacks can be formulated as:

$$\min \mathbb{E}_{(v,\mathbf{t})\sim\mathcal{D}_s,\mathbf{t}_{pos}\sim\mathcal{D}_s}(\mathcal{L}_{CL}+\lambda\mathcal{L}_{Dis}), \text{ s.t. } \|G_w(z_v;f_T(\mathbf{t}))\|_{\infty} \le \epsilon_v,$$
(5)

where  $\lambda$  is the pre-defined hyperparameter to balance the contributions of  $\mathcal{L}_{CL}$  and  $\mathcal{L}_{Dis}$ . By training the network using the proposed loss function based on the data distribution of the multimodal training dataset  $D_s$ , the generator is optimized to produce UAP that can push the features of mismatched image-text pairs together while pulling the embeddings of the matched ones apart, thereby learning a UAP with excellent generalization ability and adversarial transferability.

**Text Modality Attacks.** In textual attacks, the UAP generator's architecture and training loss are completely symmetrical with those of image attacks. Correspondingly, embeddings of the matched image v are used as the cross-modal conditions for the generator. Given an adversarial text  $t_{adv}$  as the anchor sample, we use the set  $\mathbf{v} = \{v_1, v_2 \dots, v_S\}$  scaled from the originally matched image v as negative samples while the  $\mathbf{v}' = \{v'_1, v'_2 \dots, v'_S\}$  augmented from the farthest image v' within the randomly sampled image set as positive samples to formulate the  $\mathcal{L}_{CL}$  loss.  $\mathcal{L}_{Dis}$  is consequently calculated as the negative Euclidean distance between the embeddings of  $t_{adv}$  and the clean input t.

A notable distinction between the image and text attacks is the approach to inject adversarial pertur-295 bations. Due to the discreteness of text data, we apply the token-wise substitute strategy (Lu et al., 296 2023; Wang et al., 2024) that replaces certain important words in the original sentence with crafted 297 adversarial words. Accordingly, the conditional generator is utilized to output the adversarial textual 298 embeddings, which are subsequently mapped back to the vocabulary space to obtain a universally 299 applicable word-level perturbation. Prior to implementing the word replacement, a meticulous pro-300 cess is undertaken to identify the most optimal position within the sentence to insert the perturbation. 301 Our strategy intends to identify and replace the words that are more likely to have a greater influence 302 during the decision-making. Concretely, for each word  $w_i$  within a given sentence, we compute the 303 distance between the embeddings of the original sentence and the  $w_i$ -masked version encoded by 304 the VLP models to determine its contribution. As aforementioned, we set  $\epsilon_t = 1$  and choose the 305 single word exerting the highest feature distance as the target for replacement.

306 307 308

314

316

#### 4 EVALUATION

We first present the experimental setup in Sec. 4.1 and then comprehensively evaluate C-PGC across multiple VLP models in Sec. 4.2. Sec. 4.3 presents results on more downstream V+L tasks to further validate the effectiveness. Besides, sufficient ablation studies in Sec. 4.4 validate the contribution of each proposed technique and explore the impact of several crucial factors. More experiment results such as the cross-domain attacks from Flickr30k to MSCOCO are provided in Appendix E.

#### 315 4.1 EXPERIMENTAL SETUP

Downstream tasks and datasets. We conduct a comprehensive study of C-PGC on four downstream V+L tasks, including image-text retrieval (ITR), image captioning (IC), visual grounding (VG), and visual entailment (VE). For ITR tasks, we employ the Flickr30K (Plummer et al., 2015) and MSCOCO (Lin et al., 2014) datasets which are commonly used in previous works (Zhang et al., 2022; Lu et al., 2023). The MSCOCO is also adopted for evaluating the IC task and we test VG and VE tasks on SNLI-VE (Xie et al., 2019) and RefCOCO+ (Yu et al., 2016) respectively.

**Surrogate models and victim models.** We conduct experiments on a wide range of VLP models including ALBEF (Li et al., 2021), TCL (Yang et al., 2022), X-VLM (Zeng et al., 2022), CLIP<sub>VIT</sub>

326															
327	Dataset	Source	Method	AL	BEF	T	CL	X-V	/LM	CLI	P <sub>ViT</sub>	CLI	P <sub>CNN</sub>	BI	JP
328	Dutubet	Bouree	methou	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR
329		ALBEE	GAP	69.78	81.59	22.15	29.97	6.61	18.37	23.4	37.54	29.92	44.29	16.09	28.12
330			Ours	90.13	88.82	62.11	64.48	20.53	39.38	43.1	65.93	54.4	72.51	44.79	56.36
331		TCL	GAP	33.5	40.61	82.41	80.67	6.61	17.79	21.55	38.56	30.57	45.48	21.45	31.82
332			Ours	50.20	50.29	94.95	90.04	14.94	33.90	40.92	00.41	52.98	/0.00	35.75	52.52
333	Flickr30k	X-VLM	GAP Ours	16.14 <b>24.46</b>	24.43 <b>47.77</b>	17.08 29.19	26.2 50.15	90.24 93.29	85.98 <b>91.9</b>	24.51 <b>43.47</b>	41.15 66.03	42.62 <b>59.2</b>	53.08 72.79	16.19 32.39	25.74 <b>52.24</b>
334			GAP	11 72	23 34	15 32	26 39	8 54	20.48	85 73	90.45	48 83	60.78	14 83	26.46
335		CLIP <sub>ViT</sub>	Ours	23.23	38.67	25.05	41.79	15.85	35.59	88.92	93.05	66.06	75.42	26.71	45.7
336		CLIP	GAP	13.57	25.21	19.05	28.87	11.59	23.13	27.46	43.16	73.18	81.6	15.25	27.94
337		CLIFCNN	Ours	15.31	38.93	19.77	43.72	17.17	41.65	39.9	64.82	81.74	88.9	22.19	46.11
338		RI IP	GAP	12.23	23.94	14.49	25.44	6.91	17.81	20.32	37	26.81	43.59	47.21	73.33
339		DEII	Ours	32.17	44.4	33.44	44.51	18.6	35.53	43.35	60.26	48.96	66.95	71.82	82.82
340		ALBEF	GAP	82.65	84.35	53.6	45.46	15.09	15.64	25.18	29.94	28.06	35.28	37.44	33.61
341			Ours	96.18	95.09	82.49	76.24	39.97	48.58	59.71	67.05	61.27	70.8	59.18	63.89
342		TCL	GAP	55.92	48.22	95.16	92.29	17.34	17.01	28.73	31.19	32.27	39.81	43.59	39.64
343			Gub	70.02	/1.1/	90.72	93.00	42.99	40.4	10.52	79.00	/4.1	02.91	02.55	00.97
344	MSCOCO	X-VLM	GAP Ours	26.35 <b>51.46</b>	23.72 65.71	<b>52.8</b>	22.91 64.99	95.1 <b>98.89</b>	88.84 <b>95.79</b>	32.39 67.42	38.16 <b>75.45</b>	52 75.49	55.4 <b>82.58</b>	24.67 55.74	22.65 66.7
345			GAP	35.96	31.91	37.33	32.56	33.42	29.25	97.71	96.04	74.63	74.67	33.47	31.99
346		CLIP <sub>ViT</sub>	Ours	46.92	53.89	46.03	50.87	41.49	48.6	98.74	98.01	81.58	86.5	47.35	57.55
347		CLIPCNN	GAP	28.67	27.51	29.84	27.69	26.4	24.81	39.64	40.53	90.34	91.56	24.99	26.18
348			Ours	33.38	46.68	40.61	50.76	35.34	46.95	63.83	70.15	94.89	94.42	37.38	53.06
349		BLIP	GAP	35.55	38.75	35.62	33.79	22.7	21.25	32.05	35.8	40.93	45.58	73.46	72.37
050			Ours	01.95	60.92	00.95	39.57	51.81	52.53	02.23	/2.51	09.01	/8.44	91.67	90.42

Table 1: ASR (%) of our C-PGC and GAP for image-text retrieval tasks on Flickr30k and MSCOCO. TR indicates text retrieval based on the input image, while IR is image retrieval using input text.

(Radford et al., 2021), CLIP<sub>CNN</sub> (Radford et al., 2021), and BLIP (Li et al., 2022). Note that for
different V+L tasks, we correspondingly select different VLP models for evaluation based on their
capability (Wang et al., 2024). For instance, among the six considered VLP models, only ALBEF,
TCL, and X-VLM can handle VG tasks, while only ALBEF and TCL can deal with VE tasks.

Baselines. To better reveal the superiority of our proposed method in attacking VLP models, we transplant a representative and powerful algorithm GAP (Poursaeed et al., 2018) to the multimodal attack scenarios by appropriately modifying its original loss function (Lu et al., 2023).

**Implementation details.** Following (Lu et al., 2023), we adopt Karpathy split (Karpathy & Fei-Fei, 2015) to preprocess the dataset and build the test set for evaluation. The test set is disjoint with the generator's training data for rigorous assessment. To ensure the perturbation invisibility, we follow (Wang et al., 2024) and limit the perturbation budgets  $\epsilon_v$  to 12/255 and  $\epsilon_t$  to 1. During the augmentation, we resize the original images into five scales {0.5, 0.75, 1, 1.25, 1.5}, and apply Gaussian noise with a mean of 0 and a standard deviation of 0.5. See Appendix G for more details.

365 366 367

351

324

325

4.2 UNIVERSAL ATTACK EFFECTIVENESS

To align with previous studies (Zhang et al., 2022; Lu et al., 2023), we first consider the typical V+L task image-text retrieval and calculate the ASR as the proportion of successful adversarial samples within the originally correctly predicted pairs based on R@1 retrieval results. Appendix E provides results of R@5 and R@10. Experimental results across six 4.2 VLP models are presented in Table 1. We also provide the visualization of the image retrieval on the MSCOCO dataset in Figure 4.

White-box attack performance. By observing the white-box ASR in the gray-shaded area, we
 demonstrate that the proposed algorithm stably achieves excellent ASR on all the evaluated VLP
 models, validating the outstanding capability of the produced UAP. With only a single pair of per turbations, we reach a noteworthy average white-box ASR of over 90% on two large datasets in
 terms of both TR and IR tasks. Especially on the MSCOCO dataset, our method achieves over
 95% average ASR on ITR tasks across six surrogate models. Compared with the GAP, the proposed



Figure 4: Illustration of image retrieval. The red indicates the universal adversarial word and the crossed-out word is the replaced one. We generate the UAP on ALBEF and test it on 6 target models. All retrieved images do not accurately correspond to the query text, validating the design of C-PGC.

method significantly improves the average fooling rates by nearly 10%, confirming the great validity
 of our suggested multimodal contrastive-learning mechanism. Essentially, the exceptional performance stems from the efficacy of our generated UAP in destroying the alignment between the image
 and text modalities, thereby misleading the VLP model during inference.

**Black-box attack performance.** We also conduct thorough experiments regarding the adversarial transferability of the generated UAP by transferring from surrogate models to other inaccessible models. As demonstrated in Table 1, the proposed C-PGC displays great attack performance in the more realistic black-box scenarios, e.g., 82.97% from TCL to CLIP<sub>CNN</sub> on MSCOCO for IR tasks. We highlight that the advantage of C-PGC over GAP (Poursaeed et al., 2018) is greatly amplified in the more challenging black-box scenarios, which achieves a significant average improvement of 18.36% and 26.32% for Flickr30K and MSCOCO respectively. These experimental results indicate that our generative contrastive learning framework does not overly rely on the encoded feature space tailored to the surrogate model. Conversely, it is well capable of transferring to breaking the multi-modal alignment of other unseen target models, thus attaining superior adversarial transferability.

Table 2: ASR (%) of ITR tasks under defense strategies. Surrogate model is ALBEF and the dataset is Flick30K. LT denotes the LanguageTool that corrects adversarial words within the sentence.

Method	AL	BEF	T	CL	X-V	'LM	CLI	P <sub>ViT</sub>	CLI	P <sub>CNN</sub>	BI	JP
method	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR
Gaussian	37.92	49.49	32.4	47.04	19.31	37.79	42.49	65.61	50	72.23	29.65	48.77
Medium	53.13	61.6	39.54	51.96	20.43	39.69	46.31	66.92	57.9	74.51	33.75	52.68
Average	29.09	44.91	29.61	44.72	17.89	36.07	42.98	65.42	49.74	72.48	27.55	46.9
JPEG	59.3	63.7	42.34	52.52	21.65	41.58	41.26	65.77	53.5	72.62	37.01	55.04
DiffPure	64.34	74.63	65.22	74.8	66.06	75.19	78.08	86.7	82.25	88.03	70.45	79.09
NRP	32.33	40.63	20.19	39.23	14.63	32.62	48.4	69	59.72	74.09	30.28	52.2
NRP+LT	29.05	35.23	21.33	37.41	15.55	29.63	47.19	67.35	56.82	73.47	28.23	50.59

Defense Strategies. We next analyze several defense strategies to mitigate the potential harm
brought by the proposed C-PGC. Specifically, we totally align with TMM (Wang et al., 2024) and
consider several input preprocessing-based schemes, including image smoothing (Ding et al., 2019)
(Gaussian, medium, average smoothing), JPEG compression (Dziugaite et al., 2016), NRP (Naseer
et al., 2020), and the prevalent DiffPure (Nie et al., 2022), a powerful purification defense using
diffusion models. For adversarial text correction, we choose the LanguageTool (LT) (Wang et al., 2024), which has been widely adopted in various scenarios due to its universality and effectiveness.

The attack results in Table 2 demonstrate that the proposed attack still attains great ASR against different powerful defenses. It also indicates that NRP+LT would be a decent choice to alleviate the threat brought by C-PGC. Another noteworthy finding is that, although DiffPure (Nie et al., 2022) exhibits remarkable performance in defending attacks in classification tasks, its ability is greatly

reduced in V+L scenarios since the denoise process could also diminish some texture or semantic information that is critical for VLP models, thereby acquiring unsatisfactory defense effects.

#### 4.3 EVALUATION ON MORE DOWNSTREAM TASKS

We further demonstrate C-PGC's ability to destroy the multimodal alignment by presenting more results on diverse V+L tasks. Specifically, we consider Image Captioning (IC), Visual Grounding (VG), and Visual Entailment (VE). The results of VE are shown in Appendix E due to space limit.

Table 3: Attacks results of image captioning. The Baseline represents the performance of the target model on clean data. The used dataset is MSCOCO.

Source	B@4	METEOR	ROUBE_L	CIDEr	SPICE
Baseline	39.7	31.0	60.0	133.3	23.8
ALBEF	30.1	23.7	51.2	92.5	17.5
TCL	29.5	23.5	51.0	88.9	17.3
BLIP	21.2	19.1	45.5	62.5	13.7

**Image captioning.** The objective of IC is to generate text descriptions relevant to the semantic content based on the given image. We use ALBEF, TCL, and BLIP as source models and attack the commonly used captioning model BLIP. Similar to SGA (Lu et al., 2023), several typical evaluation metrics of IC are calculated to measure the quality of generated captions, including BLEU (Papineni et al., 2002), METEOR (Banerjee & Lavie, 2005), ROUGE (Lin, 2004), CIDEr

(Vedantam et al., 2015), and SPICE (Anderson et al., 2016). The results in Table 3 demonstrate that our algorithm again displays prominent attack effects, e.g., the crated UAP induces notable drops of 10.2% and 9% in the B@4 and ROUGEL respectively when transferred from TCL to BLIP.

Table 4: Attack results of visual grounding. The first row displays the source models, where the Baseline indicates the clean performance of the target model on clean data.

Target		Baseline	e		ALBEF	7		TCL			X-VLM	ĺ
Imger	Val	TestA	TestB	Val	TestA	TestB	Val	TestA	TestB	Val	TestA	TestB
ALBEF	58.4	65.9	46.2	37.1	39.8	32.0	42.2	46.9	35.2	37.6	40.2	33.0
TCL	59.6	66.8	48.1	43.6	47.8	36.9	39.0	41.4	33.6	39.5	41.7	34.1
X-VLM	70.8	67.8	61.8	51.8	54.7	47.7	52.7	55.9	47.8	33.1	34.7	28.8

**Visual grounding.** This is another common V+L task, which aims to locate the correct position in an image based on a given textual description. We conduct experiments on RefCOCO+ using ALBEF, TCL, and X-VLM as source and target models. Table 4 indicates that C-PGC brings a notable negative impact on the localization accuracy in both white-box and black-box settings, again verifying that the produced UAP strongly breaks the cross-modal interaction and alignment.

4.4 ABLATION STUDY

This subsection employs the representative ALBEF (Li et al., 2021) model as the surrogate model to
provide sufficient ablation studies on Flickr30K. We begin our analysis on the contribution of each
proposed technique. Subsequently, we examine the sensitivity of certain hyperparameters.

The effect of  $\mathcal{L}_{CL}$  and  $\mathcal{L}_{Dis}$ . To investigate the impact of the proposed loss terms, we introduce two variants C-PGC<sub>CL</sub> and C-PGC<sub>Dis</sub> that remove  $\mathcal{L}_{CL}$  and  $\mathcal{L}_{Dis}$  from the training loss respectively. As shown in Table 5, the removal of  $\mathcal{L}_{CL}$  leads to significant degradation, particularly for black-box transferable attacks. e.g., a 27.12% ASR drop in TR tasks from ALBEF to TCL. This validates the considerable contribution of  $\mathcal{L}_{CL}$  to guarantee a successful attack. Regarding the influence of  $\mathcal{L}_{Dis}$ , we demonstrate that this unimodal guidance can further enhance attacks on the basis of  $\mathcal{L}_{CL}$ , e.g., a 10.59% increase in the ASR of TR tasks for white-box attacks on ALBEF. The proposed two loss terms complement each other and jointly underpin the generalizability of the generated UAP. 

**The effect of positive sample selection.** To validate the farthest selection strategy when constructing positive samples, we design another variant C-PGC<sub>*Rand*</sub> that adopts randomly sampled data points as positive samples. Results in Table 5 reveal the necessity of the proposed farthest selection strategy as it brings an average improvement of 25.96% in white-box ASR and 4.95% in black-box ASR. Moreover, we can also conclude that if the positive samples are not adequately defined, adding  $\mathcal{L}_{CL}$  would even severely harm the white-box performance (see C-PGC<sub>CL</sub> and C-PGC<sub>Rand</sub>).

Method	AL	BEF	TC	CL	X-V	/LM	CLI	P <sub>ViT</sub>	CLI	P <sub>CNN</sub>	BI	LIP
	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR
C-PGC	90.13	88.82	62.11	64.48	20.53	39.38	43.1	65.93	54.4	72.51	44.79	56.36
$C-PGC_{CL}$	76.46	77.58	34.99	47.55	14.33	33.61	42.98	62.81	46.11	65.58	27.13	46.44
$C-PGC_{Dis}$	79.54	82.46	56.52	62.21	20.24	38.26	39.78	65.1	52.2	71.01	42.43	55.52
C-PGC <sub>Rand</sub>	61.87	65.17	43.69	52.54	19.51	35.47	40.33	65.77	54.15	70.62	39.43	52.59
$C-PGC_{CA}$	85.18	83.07	45.76	53.73	15.24	34.02	39.29	60.61	47.15	40.64	32.39	48.29
70 60 9 50 40 30				70 %) 60 WSV 50 40 30	DLIP			V Zervez	70 60 50 40			
				$\begin{array}{c} 20\\ 10\\ 4 \end{array}$	255	8/255	12/255	16/255	20	2	3	- X-VI BLIP
0.01 0.	.05 0.1 <b>λ</b>	0.5	1			ε.,				2	۶.	

Table 5: AS	R (%)	of C-PC	C and its	variants	averaged	across six	target	models or	i retrieval	tasks.
14010 01 110	( /0 /	01 0 1 0			averagea		- Ber	1110 00 010 01		· ······

Figure 5: ASR of five target models on TR tasks under various  $\lambda$ . Figure 6: ASR of five target models on the TR task under different values of perturbation budgets  $\epsilon_v$  and  $\epsilon_t$  respectively.

**The effect of cross-modal conditions.** As aforementioned, cross-attention (CA) modules are introduced into the generator to exploit cross-modal information. We then design C-PGC<sub>CA</sub> that cancels these CA layers to explore their influence. As expected, it causes a notable 9.78% average decrease across six target models, confirming the vital role of cross-modal knowledge. An interesting finding is that C-PGC<sub>CA</sub> induces a more pronounced drop in black-box attacks than white-box ones, indicating that cross-modal conditions exert a greater contribution to the adversarial transferability.

515 **Different regulatory factor**  $\lambda$ . The value of  $\lambda$  is a critical factor as it adjusts the scales of the two 516 loss terms  $\mathcal{L}_{CL}$  and  $\mathcal{L}_{Dis}$ . We explore the attack performance under various values of  $\lambda$  to confirm 517 the optimal value. Figure 5 indicates that  $\lambda = 0.1$  achieves superior performance.

518 **Different perturbation budgets**  $\epsilon_v$  and  $\epsilon_t$ . As shown in Figure 6, we analyze varying perturbation 519 budgets for  $\epsilon_v$  and  $\epsilon_t$ . Generally, the ASR increases with the larger perturbation magnitudes. Note that when  $\epsilon_v = 4/255$ , C-PGC's performance is severely compromised since the budget 4/255 is 521 too small to allow the UAP to carry enough information required to generalize to diverse data samples. We also find that the improvement slows down as  $\epsilon_v$  increases from 12/255 to 16/255. Thus, 522 we select the moderate value of 12/255 to reach a balance between attack utility and impercepti-523 bility. For text perturbation,  $\epsilon_t$  exhibits a more profound influence on the black-box attacks. In our 524 experiments, we strictly set  $\epsilon_t = 1$  for invisibility. However, attackers can adjust the value of  $\epsilon_t$  in 525 accordance with their demands to trade off the attack efficacy and the perturbation stealthiness. 526

527 528

504 505

507 508

#### 5 CONCLUSION

529

CONCLUSI

530 This paper delves into the challenging task of launching universal adversarial attacks against VLP 531 models and proposes an effective solution that achieves superior performance using only one universal pair of image-text perturbations. We begin by revealing the unsatisfactory results of existing 532 UAP methods and empirically explaining the underlying reasons. Based on our analysis, we pro-533 pose to break the crucial cross-modal alignment in VLP models by designing a contrastive-learning 534 generative UAP framework that leverages both unimodal and multimodal information to enhance the attacks. Extensive experiments validate the efficacy of the proposed algorithm on diverse VLP 536 models and V+L tasks. We highlight that the proposed framework makes a significant step in exploring the classic universal adversarial attacks in VLP models and deepens our understanding of the 538 mechanism of VLP models. We also hope that this paper can promote future research that explores more sophisticated defenses to strengthen the resilience of VLP models against adversarial attacks.

## 540 REFERENCES

556

558

559

581

582

583

588

589

- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part V 14*, pp. 382–398. Springer, 2016.
- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and
  Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
  pp. 6077–6086, 2018.
- Gautham Anil, Vishnu Vinod, and Apurva Narayan. Generating universal adversarial perturbations
   for quantum classifiers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 10891–10899, 2024.
- Shuang Bai and Shan An. A survey on automatic image caption generation. *Neurocomputing*, 311: 291–304, 2018.
  - Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pp. 65–72, 2005.
- 560 Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In 2017
   561 *ieee symposium on security and privacy (sp)*, pp. 39–57. Ieee, 2017.
- Fei-Long Chen, Du-Zhen Zhang, Ming-Lun Han, Xiu-Yi Chen, Jing Shi, Shuang Xu, and Bo Xu.
  Vlp: A survey on vision-language pre-training. *Machine Intelligence Research*, 20(1):38–56, 2023.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 104–120. Springer, 2020.
- Gavin Weiguang Ding, Luyu Wang, and Xiaomeng Jin. Advertorch v0. 1: An adversarial robustness toolbox based on pytorch. *arXiv preprint arXiv:1902.07623*, 2019.
- Virginie Do, Oana-Maria Camburu, Zeynep Akata, and Thomas Lukasiewicz. e-snli-ve: Corrected visual-textual entailment with natural language explanations. *arXiv preprint arXiv:2004.03744*, 2020.
- 575
  576
  576
  576
  577
  578
  578
  578
  579
  579
  579
  570
  570
  570
  571
  572
  573
  574
  575
  575
  576
  577
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
  578
- Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. A survey of vision-language pre-trained
   models. *arXiv preprint arXiv:2202.10936*, 2022.
  - Gintare Karolina Dziugaite, Zoubin Ghahramani, and Daniel M Roy. A study of the effect of jpg compression on adversarial images. *arXiv preprint arXiv:1608.00853*, 2016.
- Kevin Eykholt, Ivan Evtimov, Earlence Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul
   Prakash, Tadayoshi Kohno, and Dawn Song. Robust physical-world attacks on deep learning
   visual classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1625–1634, 2018.
  - Weiwei Feng, Nanqing Xu, Tianzhu Zhang, and Yongdong Zhang. Dynamic generative targeted attacks with pattern injection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16404–16414, 2023.
- Haoran Gao, Hua Zhang, Jiahui Wang, Xin Zhang, Huawei Wang, Wenmin Li, and Tengfei Tu.
   Nuat-gan: Generating black-box natural universal adversarial triggers for text classifiers using generative adversarial networks. *IEEE Transactions on Information Forensics and Security*, 2024.

603

609

- Jamie Hayes and George Danezis. Learning universal adversarial perturbations with generative models. In 2018 IEEE Security and Privacy Workshops (SPW), pp. 43–49. IEEE, 2018.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
   unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9729–9738, 2020.
- Richang Hong, Daqing Liu, Xiaoyu Mo, Xiangnan He, and Hanwang Zhang. Learning to compose
   and reason with language tree structures for visual grounding. *IEEE transactions on pattern analysis and machine intelligence*, 44(2):684–696, 2019.
- Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3128–3137, 2015.
- Taewan Kim and Joydeep Ghosh. On single source robustness in deep fusion models. Advances in Neural Information Processing Systems, 32, 2019.
- Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven
   Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum
   distillation. Advances in Neural Information Processing Systems, 34:9694–9705, 2021.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre training for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Kuannan Liu, Yaoyao Zhong, Yuhang Zhang, Lixiong Qin, and Weihong Deng. Enhancing general ization of universal adversarial perturbation through gradient aggregation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4435–4444, 2023.
- Dong Lu, Zhiqiang Wang, Teng Wang, Weili Guan, Hongchang Gao, and Feng Zheng. Set-level guidance attack: Boosting adversarial transferability of vision-language pre-training models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 102–111, 2023.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal
   adversarial perturbations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1765–1773, 2017.
- Konda Reddy Mopuri, Aditya Ganeshan, and R Venkatesh Babu. Generalizable data-free objective for crafting universal adversarial perturbations. *IEEE transactions on pattern analysis and machine intelligence*, 41(10):2452–2465, 2018.
- Muzammal Naseer, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Fatih Porikli. A self supervised approach for adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 262–271, 2020.
- Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar.
   Diffusion models for adversarial purification. *arXiv preprint arXiv:2205.07460*, 2022.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
   evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association* for Computational Linguistics, pp. 311–318, 2002.

- Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer imageto-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pp. 2641–2649, 2015.
- Omid Poursaeed, Isay Katsman, Bicheng Gao, and Serge Belongie. Generative adversarial perturba tions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
   pp. 4422–4431, 2018.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
  Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
  models from natural language supervision. In *International conference on machine learning*, pp.
  8748–8763. PMLR, 2021.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image
   description evaluation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4566–4575, 2015.
- Donghua Wang, Wen Yao, Tingsong Jiang, and Xiaoqian Chen. Improving transferability of universal adversarial perturbation with feature disruption. *IEEE Transactions on Image Processing*, 2023.
- Haodi Wang, Kai Dong, Zhilei Zhu, Haotong Qin, Aishan Liu, Xiaolin Fang, Jiakai Wang, and
   Xianglong Liu. Transferable multimodal attack on vision-language pre-training models. In 2024
   *IEEE Symposium on Security and Privacy (SP)*, pp. 102–102. IEEE Computer Society, 2024.
- Kaiye Wang, Qiyue Yin, Wei Wang, Shu Wu, and Liang Wang. A comprehensive survey on cross modal retrieval. *arXiv preprint arXiv:1607.06215*, 2016.
- Kiaosen Wang, Xuanran He, Jingdong Wang, and Kun He. Admix: Enhancing the transferability
  of adversarial attacks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16158–16167, 2021.
- <sup>676</sup>
   <sup>677</sup> Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. Visual entailment: A novel task for finegrained image understanding. *arXiv preprint arXiv:1901.06706*, 2019.
- Jinyu Yang, Jiali Duan, Son Tran, Yi Xu, Sampath Chanda, Liqun Chen, Belinda Zeng, Trishul
  Chilimbi, and Junzhou Huang. Vision-language pre-training with triple contrastive learning.
  In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15671–15680, 2022.
- Karren Yang, Wan-Yi Lin, Manash Barman, Filipe Condessa, and Zico Kolter. Defending multi modal fusion models against single-source adversaries. In *Proceedings of the IEEE/CVF Confer- ence on Computer Vision and Pattern Recognition*, pp. 3340–3349, 2021.
- Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pp. 69–85. Springer, 2016.
- Yan Zeng, Xinsong Zhang, and Hang Li. Multi-grained vision language pre-training: Aligning texts
   with visual concepts. In *International conference on machine learning*, pp. 25994–26009. PMLR, 2022.
- Jiaming Zhang, Qi Yi, and Jitao Sang. Towards adversarial attack on vision-language pre-training models. In *Proceedings of the 30th ACM International Conference on Multimedia*, pp. 5005– 5013, 2022.
- Peng-Fei Zhang, Zi Huang, and Guangdong Bai. Universal adversarial perturbations for vision language pre-trained models. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 862–871, 2024.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

#### 702 PSEUDOCODE OF THE PROPOSED ALGORITHM А 703

We present the pseudocode of our proposed attack algorithm for image modality in Alg. 1. Note that the text attacks are completely symmetrical as illustrated in Sec. 3.3.

A	gorithm 1 Pseudocode of universal image attacks
R	equire: $G_w(\cdot)$ : the perturbation generator; $D_s$ : the multimodal training set; $f_I, f_T$ : image en-
	coder and text encoder of the surrogate VLP model; N: the max iteration; $\epsilon_v$ : the perturbation
	budget; S: the scaling times;
E	<b>isure:</b> Universal image perturbation $\delta_v$ ;
1	: <b>Initialize</b> the fixed noise $z_v$ with Gaussian distribution;
2	: for $i \leftarrow 0$ to N do
3	Randomly sample an image-text pair $(v, \mathbf{t}) \sim D_s$ ;
4	: $\delta_v = Clip_{\epsilon_v}(G_w(z_v; f_T(\mathbf{t}))), v_{adv} = v + \delta_v;$
5	: Augment v and $v_{adv}$ into different scales and apply random Gaussian noises to obtain $\mathbf{v} =$
	$\{v_1 \dots, v_S\}$ and $\mathbf{v}_{adv} = \{v_1^{adv} \dots, v_S^{adv}\};$
6	: Randomly sample a batch of text sets from $D_s$ and obtain $\mathbf{t}_{pos} = \{t'_1, \dots, t'_K\}$ by selecting
	the one with the farthest feature distance from the clean image $v$ ;
7	: Compute $\mathcal{L}_{CL}$ with $\mathbf{v}_{adv}$ , $\mathbf{t}$ and $\mathbf{t}_{pos}$ by Eq. (3);
8	: Compute $\mathcal{L}_{Dis}$ with v and v <sub>adv</sub> by Eq. (4);
9	: Optimize the generator $G_w$ based on Eq. (5);
10	: Backward pass and update $G_w$ ;
11	end for
12	: Return $\delta_v$

726 727

728 729

730

731

732

733 734

736 737

738 739

740

741

742

743

704

705

706

#### **RATIONAL BEHIND OUR DESIGN OF LOSS FUNCTION** В

It is widely acknowledged that contrastive learning serves as a powerful and foundational tool for modality alignment in VLP models, establishing a nearly point-to-point relationship between image and text features. Our core idea stems from the general principle: "It's easier to tear down than to build up." Since contrastive learning can establish robust and precise alignment, leveraging the same technique to disrupt the established alignments is expected to yield effective performance.

Taking image attack as an example, the principle behind our contrastive learning-based attack can 735 be understood from two perspectives:

- Leverage the originally matched texts as negative samples to push the aligned image-text pair apart. This broadly corresponds to the common objective of untargeted attacks.
- Additionally, our contrastive paradigm introduces dissimilar texts as positive samples to pull the adversarial image out of its original subspace and relocate it to an incorrect feature area.

By simultaneously harnessing the collaborative effects of *push* (negative samples) and *pull* (positive samples), the proposed contrastive framework achieves exceptional attack performance, which has been validated by comprehensive experimental results.

Table 6:	ASR	results	of the	proposed	method	with	different	loss	functions	on	Flickr30	when	the
surrogate	mode	el is AL	BEF.										

748													
749	Method	ALI	BEF	TC	CL	X-V	'LM	CLI	P <sub>ViT</sub>	CLI	P <sub>CNN</sub>	BL	JP
750		TR	IR	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR
751	$\mathcal{L}_{MSE}$	12.02	30.75	14.39	35.08	11.41	30.79	37.32	56.05	40.17	56.39	19.66	37.33
752	$\mathcal{L}_{Cos}$	57.55	67.4	37.06	49.45	10.7	28.48	37.49	58.3	40.87	58.39	23.33	39.44
753	$\mathcal{L}_{CL}$	76.46	82.46	56.52	62.61	14.33	33.61	42.98	62.81	46.11	65.58	27.13	46.44
754	$\mathcal{L}_{MSE}$ + $\mathcal{L}_{Dis}$	81.09	83.71	48.76	56.54	17.58	35.72	41.5	64.72	47.41	70.34	35.96	51.76
755	$\mathcal{L}_{Cos}$ + $\mathcal{L}_{Dis}$	65.20	72.71	36.13	50.06	18.63	36.74	42.23	65.17	50.91	69.78	36.91	50.69
/00	$\mathcal{L}_{CL}$ + $\mathcal{L}_{Dis}$	90.13	88.82	62.11	64.48	20.53	39.38	43.1	65.93	54.4	72.51	44.79	56.36

<sup>756</sup>Besides, we also explore several potential alternative loss functions that more directly align with the common untargeted attack Table 6, including maximizing the negative cosine similarity  $\mathcal{L}_{Cos}$  or MSE distance  $\mathcal{L}_{MSE}$  between the features of matched image-text pairs.

Recall that  $\mathcal{L}_{CL}$  and  $\mathcal{L}_{Dis}$  denote the proposed contrastive loss and the unimodal loss term respectively. As observed, the use of  $\mathcal{L}_{CL}$  consistently brings significant ASR improvements, verifying the rationality and superiority of contrastive loss.

#### C COMPARISON WITH A CONCURRENT STUDY

We notice a concurrent study (Zhang et al., 2024) on UAP attacks for VLP models, which also shows promising attack performance. To make a fair comparison, we faithfully reproduce this algorithm using their publicly released code under the same experimental settings as ours. Note that Zhang et al. (2024) implement several versions of their method and we report their best results in Table 7.

Table 7: Comparison of C-PGC with a recent attack (Zhang et al., 2024) on Flicke30K

Source	Method	AL	BEF	TCL	X-V	LM	CLI	P <sub>ViT</sub>	CLII	CNN	BI	JP
		TR	IR   T	R IR	TR	IR	TR	IR	TR	IR	TR	IR
ALBEF	ETU C-PGC	78.01 <b>90.13</b>	84.5629.88.8262.	.92 35.91 .11 64.48	14.33 20.53	22.03 <b>39.38</b>	23.77 43.1	39.2 <b>65.93</b>	33.55 <b>54.4</b>	47.69 <b>72.51</b>	22.61 <b>44.79</b>	32.28 <b>56.36</b>
CLIP <sub>ViT</sub>	ETU C-PGC	14.8 23.23	25.23   21. 38.67   25.	.22 30.87 .05 41.79	10.87 15.85	24.96 <b>35.59</b>	84.14 88.92	90.45 <b>93.05</b>	57.51 66.06	65.51 <b>75.42</b>	16.4 <b>26.71</b>	27.22 <b>45.7</b>

By contrastively training the conditional generator, the proposed C-PGC greatly enhances the attack by achieving significant improvements in ASR. Particularly in the more realistic and challenging transferable scenarios, the proposed method achieves considerably better performance, e.g., 32.19% and 28.57% increase in ASR of TR and IR tasks when transferring from ALBEF to TCL. These results strongly confirm the superiority of our contrastive learning-based generative paradigm.

## D SEMANTIC SIMILARITY ANALYSIS

The basic objective of untargeted adversarial attacks is to fool the victim model to output incorrect predictions (Dong et al., 2018), while the attacker is supposed to preserve semantic similarity between the original and the adversarial sample to ensure attack imperceptibility. In our implementation, we follow the rigorous setup in prior works (Zhang et al., 2022; Lu et al., 2023; Wang et al., 2024) that modify only one single word to preserve semantic similarity and attack stealthiness. To quantitatively analyze the semantic similarity, we provide the BERT scores (Zhang et al.), which calculate the P (precision), R (recall), and F1 (F1 score) as results for the semantic distance between 5,000 clean and adversarial sentences in Table 8.

Table 8: Comparison of BERTScore between clean and adversarial texts.

Method		ALBEF			TCL			$\text{CLIP}_{\text{ViT}}$			CLIP <sub>CNN</sub>	
	Р	R	F1	Р	R	F1	Р	R	F1	P	R	F1
Co-Attack	0.8328	0.8589	0.8455	0.8325	0.8588	0.8453	0.8269	0.8526	0.8394	0.8271	0.853	0.8397
SGA	0.8389	0.8654	0.8518	0.8376	0.8646	0.8509	0.8416	0.8697	0.8553	0.8378	0.865	0.8511
Ours	0.8891	0.8613	0.8748	0.8924	0.8687	0.8802	0.8746	0.8684	0.8713	0.8948	0.8842	0.8893

801 802 803

763

764 765

766

767

768

769 770

781

782

783 784 785

786 787

788

789

790

791

792

793

794

Note that we provide previous sample-specific algorithms Co-Attack (Zhang et al., 2022) and SGA
 (Lu et al., 2023) as references. Notably, our method achieves better similarity scores to these wide acknowledged sample-specific methods across different surrogate VLP models, demonstrating the
 outstanding attack stealthiness of our C-PGC. The lower semantic similarity of sample-specific
 methods essentially stems from their word-selection mechanism, which maximizes the semantic
 distance tailored to every input sentence for attack enhancement. Specifically, to achieve better per formance, these methods select the adversarial word that maximizes the distance between the original and perturbed text for every input sentence, which inherently leads to relatively larger semantic

810 deviations. This highlights that our universal attack achieves a better balance between efficacy and 811 stealthiness. Besides, we also provide BLEU metrics when the surrogate model is ALBEF in Table 812 9. These results again validate the better stealthiness of our C-PGC. 813

Table 9: Comparison of BLEU metrics between clean and adversarial texts.

Method	B@4	METEOR	ROUBE_L	CIDEr	SPICE
Co-Attack	0.79	0.52	0.895	7.03	0.661
Ours	0.798 <b>0.889</b>	0.527 0.552	0.898 <b>0.905</b>	7.139 8.036	0.008 <b>0.671</b>

#### E MORE EXPERIMENTAL RESULTS

In this section, we provide more experimental results of our method in various tasks and scenarios.

**CLIP**<sub>ViT</sub> ALBEF TCL X-VLM BLIP **CLIP**<sub>CNN</sub> Method Source TR IR TR IR TR IR TR IR TR IR TR IR C-PGC $_{Sin}$ 82.99 86.14 49 56.98 18.19 35 79 40.52 65.9 51.09 69.68 38.54 52.86 ALBEF C-PGC 90.13 88.82 62.11 64.48 20.53 39.38 43.1 65.93 54.4 72.51 44.79 56.36 C-PGC<sub>Sin</sub> 20.55 37.46 24.43 41.39 13.52 32.6 79.93 88.64 55.44 69.43 24.4 43.06 **CLIP**<sub>ViT</sub> 38.67 25.05 41.79 15.85 35.59 88.92 23.23 93.05 66.06 75.42 26.71 45.7 C-PGC

Table 10: ASR results of C-PGC and C-PGC<sub>Sin</sub> with a single text as the positive sample.

**Diverse target texts as positive samples.** We first investigate the effects of using multiple targets for contrastive training in maximizing the distance between adversarial and original images. We implement a variant C-PGC<sub>Sin</sub>, which uses only a single target text with the farthest distance as the positive sample. The results in Table 10 illustrate that the use of multiple target texts can enhance attack effectiveness, validating the efficacy of set-level diverse guidance.



Figure 7: Accuracy of VE tasks for different source and target models.

Visual entailment tasks. Given an image and a textual description, visual entailment involves determining whether the textual description can be inferred from the semantic information of the image. We align with previous VLP attacks (Zhang et al., 2022; Wang et al., 2024) and conduct experiments on the SNLI-VE (Xie et al., 2019) dataset using the ALBEF and TCL models. Note that the Baseline represents the clean performance of the target model on the clean data and the orange and green indicate ALBEF and TCL as source models respectively. The results presented in Figure 7 reveal that C-PGC obtains impressive attack effects by decreasing the average accuracy by nearly 20%. Notably, (Do et al., 2020) has reported a large number of annotation errors in the labels of the SNLI-VE corpus used for VE tasks. Therefore, the presented results are only for experimental integrity and reference purposes.

Ablation study of the data augmentation. As in the main text, we are motivated by the significant gains introduced by SGA's augmentation Lu et al. (2023) and hence integrate it into the proposed framework to enhance the universal perturbation. The underlying mechanism is to leverage the many-to-many relationships between images and texts by introducing multiple augmented images to provide diverse guidance and improve the optimization direction.

859 To reveal its effectiveness and explore alternative augmentation techniques, we devise three variants, 860 including C-PGC<sub>NoAug</sub> without any augmentation, C-PGC<sub>ScMix</sub> and C-PGC<sub>Admix</sub> with the ScMix 861 Zhang et al. (2024) and Admix Wang et al. (2021) respectively. The results are shown in Table 11. 862

It can be observed that the set-level augmentation brings significant improvements over the no 863 augmentation baseline and C-PGC with the current augmentation strategy outperforms the ScMix

814

823 824

825

826

831 832

833 834

835 836

837

838 839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

Source	Method	ALBEF		TCL		X-VLM		<b>CLIP</b> <sub>ViT</sub>		<b>CLIP</b> <sub>CNN</sub>		BLIP	
Bouree		TR	IR	TR	IR	TR	IR	TR	IR	TR	IR	BI TR 34.6 32.28 41.22 44.79 24.4 23.66 25.03 26.71	IR
	C-PGC <sub>ScMix</sub>	66.08	76.26	39.03	51.24	20.73	37.47	40.02	65.58	50.13	71.85	34.6	51.
	C-PGC <sub>Admix</sub>	62.8	72.23	34.47	47.78	19	36.67	42	64.88	48.19	69.68	32.28	50.0
ALDEF	C-PGC <sub>NoAug</sub>	69.78	74.79	47	57.26	20.43	37.55	42.36	65.17	53.63	71.6	41.22	55.3
	C-PGC	90.13	88.82	62.11	64.48	20.53	39.38	43.1	65.93	54.4	72.51	44.79	56.3
	C-PGC <sub>ScMix</sub>	20.55	37.46	24.43	41.39	13.52	32.6	79.93	88.64	55.44	69.43	24.4	43.0
CUD	C-PGC <sub>Admix</sub>	19.53	37.04	24.02	41.5	14.74	34.26	85.34	91.8	59.07	71.78	23.66	43.2
LIP <sub>ViT</sub>	C-PGC <sub>NoAug</sub>	18.5	37.8	22.19	39.86	13.47	34.17	86.46	87.11	61.53	71.36	25.03	44.7
	C-PGC	23.23	38.67	25.05	41.79	15.85	35.59	88.92	93.05	66.06	75.42	26.71	45.

Table 11: Attack performance under diffe	ent data augmentation	strategies using Flickr30K
--	-----------------------	----------------------------

(Zhang et al., 2024) and Admix (Wang et al., 2021), revealing that the set-level guidance is more suitable for our contrastive training. This is achieved by SGA's alignment-preserving augmentation, which enriches image-text pairs while maintaining their inherent alignments intact Lu et al. (2023).

**Cross-domain scenarios.** We proceed to discuss the attack performance of the proposed algorithm in a more challenging scenario where there is an obvious distribution shift between the training dataset and the test samples. Specifically, we generate universal adversarial perturbations based on MSCOCO or Flickr30K and evaluate them accordingly on the other dataset. We present the attack success rates on the retrieval tasks across six models in Table 12. It can be observed that the domain gap indeed has a negative effect on attack performance. However, our method still maintains excellent ASR in most cases, unveiling its outstanding cross-domain transferability.

Table 12: ASR (%) of Cross-domain attacks on six models from Flickr30k to MSCOCO and vice versa. The gray shading indicates white-box attacks.

Setting	Source	AL	BEF	T	CL	X-V	/LM	CLI	P <sub>ViT</sub>	CLI	P <sub>CNN</sub>	BI	JP
betting	bouree	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR
	ALBEF	96.83	94.69	81.46	74.87	44.79	51.64	63.68	73.06	69.77	78.09	68.88	70.61
Elialar20V	TCL	78.27	73.17	97.83	95.03	40.46	47.34	64.98	73.27	70.96	78.18	63.71	67.1
FIICKISUK	X-VLM	50.63	65.91	53.23	65.65	95.91	93.32	65.51	74.72	75.69	81.93	57.69	67.28
MSCOCO	<b>CLIP</b> <sub>ViT</sub>	49.88	53.39	49.47	52.21	47.77	48.52	95.5	97.01	83.05	85.38	50.97	57.93
MSCOCO	<b>CLIP</b> <sub>CNN</sub>	34.78	50.42	37.17	51.24	36.81	50.87	63.07	70.92	92.48	93.66	41.81	55.12
	BLIP	ALBEF         TCL         X-VLM         CLIP <sub>VIT</sub> CLIP <sub>CNN</sub> I           TR         IR         IR </td <td>83.19</td> <td>82.17</td>	83.19	82.17									
	ALBEF	88.08	87.28	58.9	61.53	17.58	36.07	39.78	61.08	47.28	64.95	35.02	49.4
MSCOCO	TCL	47.58	53.7	87.27	83.55	18.6	34.45	51.85	72.22	59.46	76.09	37.75	53.08
MSCOCO	X-VLM	25.39	46.74	27.33	49.13	79.98	81.72	42.73	66.48	59.46	73.07	31.65	51.48
Fliabr 20K	<b>CLIP</b> <sub>ViT</sub>	21.07	39.47	24.53	42.44	15.45	36.52	93.97	95.53	62.95	77.21	25.55	45.91
FIICKISUK	<b>CLIP</b> <sub>CNN</sub>	14.29	36.57	20.5	41.7	15.55	37.06	41.63	62.87	86.53	88.73	17.98	44.31
	BLIP	33.2	46.07	36.02	47.97	23.58	38.48	43.97	65.3	56.35	71.08	71.91	73.62

Ablation study of the text perturbation. We introduce another variant C-PGC<sub>t</sub> that cancels the perturbation from the text side to investigate the contribution of image perturbation, text perturbation, and their synergy. The comparison results of C-PGC, C-PGC<sub>t</sub>, and GAP using ALBEF as the surrogate model are shown in Table 13.

Table 13: ASR (%) of C-PGC, C-PGC<sub>t</sub>, and GAP on ITR tasks using Flickr30K. Note that C-PGC<sub>t</sub> only considers attacking images and thus doesn't apply textual perturbations.

Method	ALBEF		TCL		X-V	X-VLM		<b>CLIP</b> <sub>ViT</sub>		P <sub>CNN</sub>	BLIP	
	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR
GAP	69.78	81.59	22.15	29.97	6.61	18.37	23.4	37.54	29.92	44.29	16.09	28.12
$C-PGC_t$	86.74	86.3	50.1	50.2	10.87	21.53	26.28	39.3	33.42	48.32	31.55	36.77
C-PGC	90.13	88.82	62.11	64.48	20.53	39.38	43.1	65.93	54.4	72.51	44.79	56.36

We find that when merely applying image perturbations  $(C-PGC_t)$ , our design still outperforms GAP with notable improvements, validating the proposed techniques' effectiveness in enhancing the image perturbation. Moreover, the superiority of C-PGC over C-PGC<sub>t</sub> indicates that the incorporation of textual perturbations can further boost the universal attacks on the basis of C-PGC<sub>t</sub> since the text perturbation facilitates the deconstruction of the learned cross-modal alignment.

Results of R@5 and R@10. As aforementioned, we supplement the ASR of the ITR tasks based on R@5 and R@10 metrics and provide the attack success rates in Table 14. Obviously, our proposed C-PGC still consistently attains better performance than the baseline method GAP, regardless of the evaluation measurements for retrieval results.

Table 14: Attack success rates (%) regarding R@5 and R@10 metrics of our C-PGC and GAP for
 image-text retrieval tasks.

Dataset	Source	Method	AL	BEF	T	CL	X-\	/LM	CLI	P <sub>ViT</sub>	CLI	P <sub>CNN</sub>	BI	JP
Dutuset	bource	Method	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR
	ALBEF	GAP Ours	55.71 <b>83.67</b>	73.86 <b>80.02</b>	8.01 <b>41.84</b>	10.54 <b>42.18</b>	1.2 6.9	4.84 <b>17.19</b>	4.46 18.34	15.24 <b>41.03</b>	8.28 26.22	20.27 <b>49.42</b>	5.33 <b>24.25</b>	10.77 <b>34.59</b>
	TCL	GAP Ours	17.64 29.76	20.09 35.62	77.89 <b>90.89</b>	74.53 <b>84.18</b>	0.9 <b>3.2</b>	4.2 <b>13.65</b>	4.25 20.93	15.48 <b>42.06</b>	8.6 <b>25.27</b>	20.25 <b>49.1</b>	8.65 16.5	13.16 <b>30.32</b>
Flickr30K (R@5)	X-VLM	GAP Ours	6.21 <b>7.62</b>	7.45 <b>25.1</b>	4.9 <b>8.71</b>	7.96 <b>26.63</b>	81.6 <b>89.2</b>	77.23 <b>85.84</b>	6.11 <b>19.38</b>	18.33 <b>42.48</b>	17.41 30.89	28.35 <b>50.7</b>	5.03 13.68	8.61 29
	CLIP <sub>ViT</sub>	GAP Ours	2.81 6.31	6.86 <b>17.51</b>	4.2 8.01	8 19.65	1.7 <b>4.3</b>	6.1 <b>15.1</b>	75.64 <b>76.89</b>	82.56 <b>85.2</b>	24.2 <b>39.6</b>	37.68 <b>54.68</b>	4.33 9.15	9.97 <b>23.23</b>
	CLIP <sub>CNN</sub>	GAP Ours	2.81 3.01	7.41 <b>19.09</b>	<b>5.81</b> 5.11	9.27 <b>22.7</b>	2.4 <b>3.3</b>	7.01 23.07	9.33 17.41	19.64 <b>41.17</b>	57.63 61.57	69.33 <b>74.32</b>	4.02 6.74	9.76 <b>25.16</b>
	BLIP	GAP Ours	3.41 <b>14.43</b>	7.01 <b>21.67</b>	3.7 <b>13.91</b>	7.45 <b>21.59</b>	1 5.4	4.4 <b>14.54</b>	4.46	14.43 <b>36.26</b>	6.79 23.89	18.67 <b>44.79</b>	39.13 <b>59.26</b>	68.02 <b>74.82</b>
	ALBEF	GAP Ours	74.43 <b>93.36</b>	78.62 <b>91.56</b>	37.99 <b>70.76</b>	30.08 62.31	5.56 <b>19.97</b>	7.19 <b>30.46</b>	14.26 41.58	17.11 <b>51.23</b>	15.58 44.14	21.62 55.98	23.73 41.08	23.26 <b>49.22</b>
	TCL	GAP Ours	41.48 60.62	32.59 <b>56.21</b>	92.54 <b>94.89</b>	87.81 <b>90.33</b>	6.46 <b>22.08</b>	8.08 <b>30.38</b>	16.09 53.14	18.47 <b>64.98</b>	17.98	24.3 <b>70.77</b>	29.9 45.28	28.64 53.55
MSCOCO	X-VLM	GAP Ours	12.29 31.59	11.64 <b>48.69</b>	13.43 32.1	10.99 <b>48.11</b>	90.8 <b>96.7</b>	83.05 91.66	20.02 <b>49.53</b>	23.4 60.82	37.72	40.09 <b>69.59</b>	12.64 37.4	12.04 52.5
(1100)	CLIP <sub>ViT</sub>	GAP Ours	18.78 25.69	17.46 35.95	20.38 24.69	17.37 33.14	16.81 21.37	15.15 31.38	95.21 96.7	93.04 96.49	62.77 70.76	62.62 77.86	18.48 28.72	19.48 <b>42.01</b>
	CLIP <sub>CNN</sub>	GAP	13.54 16.25	14.09 30.07	14.42 20.96	14.14 34.15	11.02 16.58	12.03 30.23	25.27 48.04	24.98 56.66	88.67 91.54	88.83 88.96	12.8 21.37	14.98 38.85
	BLIP	GAP	23.62 42.56	24.22 43 73	22.96	18.43 41 8	9.93 31.05	9.75	19.2 44 37	22.24	24.99	30.19 66.01	62.75 81 71	66.88 <b>81 91</b>
	ALBEF	GAP	51.6 80 5	71.17	5.8	6.65 34 28	0.6	2.7	1.42	9.82 31 14	4.19	13.81 39.40	3.71	7.01
	TCL	GAP	14.9 24 2	14.29 27 32	73.26	70.49 80 73	0.6	2.22	2.13	9.73 32.4	4.29	13.45 38.63	5.92	9.51 24 1
Flickr30K	X-VLM	GAP	4.1	4.81	2.7	4.51	76.5 86 3	72.58 82.94	3.65	11.63 31 49	10.84	18.91 40 3	3.01	5.29
(R@10)	CLIP <sub>ViT</sub>	GAP	2.1	4.04	2.6	4.81	1	3.79 10 15	63.29 67 98	77.83 79.46	17.89	28.11	2.31	6.12 17 23
	CLIP <sub>CNN</sub>	GAP	2.3	4.23	3.3 37	5.56	1.4 1.7	4.41	4.56	12.35	49.86	62.17 66 34	2.21	6.63
	BLIP	GAP	2.3	3.98	1.5	4.17	0.2	2.24	1.93	8.38	3.68	12.48	36.81	67.22 72
	ALBEF	GAP	69.78 01.58	76.24	32.04	23.7	3.16	4.91	10.07	13.16	12.55	16.6	19.1	19.96
	TCL	GAP	34.36 52.59	25.76 49.09	90.65 93.63	85.27 88 53	4.01	5.44 23.25	11.77   <b>44 22</b>	14.88	13.66   50.16	18.89	24.91	24.67 47.26
MSCOCO	X-VLM	GAP	7.66	7.65	8.07	7.27	88.3 94.97	79.85 88.95	15.63	18.77 53.74	31.79	33.54 62.7	8.24	9.37
(K@10)	CLIP <sub>ViT</sub>	GAP	13.13	12.39	14.05	12.13	10.68	10.89	93.72	91.51 95 21	57.39	56.47	13.35	15.64
	CLIPCNN	GAP	9.02	10.08	9.06	9.89	6.97	8.25	18.78	20.11	87.6	83.92	8.53	11.91
	BLIP	GAP	10.68	18.98	14.19	13.13	<b>10.62</b> 6.21	6.44	<b>39.89</b>	<b>49.62</b> 17.39	<b>88.74</b>	<b>85.18</b> 24.37	<b>15.97</b>	<b>33.25</b> 65.49
		Ours	33.64	36.14	32.15	33.8	22.64	28.52	36.07	50.3	47	59.07	78.39	78.98

F MULTIMODAL ALIGNMENT DESTRUCTION

966 967

To provide more intuitive evidence that our C-PGC successfully destroys the image-text alignment relationship, we compute the distance between the encoded image and text embeddings before and

after applying the UAP. Concretely, for an input pair (v, t), we calculate the relative distance  $d_{rel}$  by:

$$d_{rel} = \frac{||(f_I(v+\delta_v) - f_T(t\oplus\delta_t)||_2 - ||f_I(v) - f_T(t)||_2}{||f_I(v) - f_T(t)||_2}.$$
(6)

We provide the distances averaged on 5000 image-text pairs from Flickr30K in Table 15. Benefiting from our delicate designs, C-PGC achieves better disruption of the aligned multimodal relationship, thereby boosting the generalization ability and transferability of the produced UAP.

Table 15: Relative cross-modal feature distances to the clean image-text pairs.

Source	Method	ALBEF	TCL	BLIP	X-VLM	CLIP <sub>ViT</sub>	CLIP <sub>CNN</sub>
ALBEF	GAP	7.18	6.54	0.91	1.74	0.31	0.98
	C-PGC	<b>8.83</b>	1 <b>4.95</b>	2.73	<b>6.09</b>	<b>3.42</b>	<b>3.92</b>
TCL	GAP	4.02	24.27	0.91	0.87	0.12	0.07
	C-PGC	6.43	27.11	<b>3.64</b>	<b>4.35</b>	2.56	<b>2.94</b>
BLIP	GAP	3.17	4.67	11.82	1.74	-1.71	-0.98
	C-PGC	<b>6.41</b>	12.15	13.64	<b>4.35</b>	<b>1.71</b>	<b>1.96</b>

#### G MORE TRAINING DETAILS

For Flickr30K and MSCOCO, we randomly sample 30,000 images and their captions from the training set to train our perturbation generator. For SNLI-VE and RefCOCO+, we learn the C-PGC directly using their training set with 29,783 and 16,992 images respectively. Since an image corresponds to multiple text descriptions in these datasets, we calculate the average of their textual embedding as the multimodal condition for the cross-attention modules.

997 We initialize the noise variable  $z_v$  as a 3  $\times$  3 matrix. Meanwhile, the initial noise  $z_t$ 's dimensions in 998 the text modality depend on the size of the hidden layer within the specific VLP model. Concretely, 999 we set its dimension to  $1 \times 3$  for ALBEF, TCL, BLIP, and X-VLM, while  $1 \times 2$  for the CLIP model. 1000 When computing the multimodal contrastive loss  $\mathcal{L}_{CL}$ , the temperature  $\tau$  is set as 0.1. The generator 1001 is trained over 40 epochs with the Adam optimizer, utilizing a learning rate of  $2^{-4}$ . Following 1002 previous works Lu et al. (2023); Wang et al. (2024), we employ the attack success rate (ASR) as 1003 our quantitative measurement of our attack in ITR tasks by computing the extent the adversarial 1004 perturbations result in victim models' performance deviations from the clean performance.

1005 1006 1007

972

973 974

975

980

989 990

991

#### H DETAILED INTRODUCTIONS TO DATASETS

- Flickr30K (Plummer et al., 2015). Collected from the Flickr website, this dataset describes different items and activities, which becomes a standard benchmark for various V+L tasks. It contains 31,783 images, each of which has five associated captions. We use it for ITR tasks.
- MSCOCO (Lin et al., 2014). The MSCOCO dataset is a rich and diverse dataset consisting of 123,287 images, each of which is annotated with approximately five sentences. We use this dataset to test the attack performance of ITR and IC tasks.
- SNLI-VE (Xie et al., 2019). Originally proposed for natural language reasoning tasks, this dataset provides large-scale images and descriptions, where each image is annotated with several sentences and their logical relationship labels, including entailment, neutral, and contradiction. This dataset is used for VE tasks.
- RefCOCO+ (Yu et al., 2016). An image dataset was selected from MSCOCO, which contains 19,992 images and 141,564 annotations. It is specially used for visual grounding (VG) tasks.
- 1020 1021

#### 1022 I DISCUSSIONS AND FUTURE DIRECTIONS

1023

**Overlook of interactions between perturbations**  $\delta_v$  and  $\delta_t$ . The proposed framework generates universal perturbations for image and text respectively based on the designed multimodal and unimodal losses. Despite the remarkable attack performance, it does not consider the interactions and synergy between the perturbations  $\delta_v$  and  $\delta_t$  during optimization, which has been leveraged in several previous attacks (Lu et al., 2023; Wang et al., 2024) to improve performance. In future research, this limitation can be explored as a potential mechanism to further strengthen the attacks.

**Textual Semantic consistency.** To ensure the stealthiness of text attacks, we set the perturbation budget  $\epsilon_t = 1$ , i.e., only one word is modified. Despite the superior semantic similarity to previous sample-specific methods, there is still room to improve from the proposed C-PGC. Moreover, future studies can consider similarity preservation strategies by applying more effective constraints during the generator training or post-processing adversarial sentences to facilitate a more stealthy attack.

Leveraging Task-level Characteristics. In contrast to the unimodal scenarios, we fully leverage the unique characteristics of multimodal scenarios to enhance the modeling of a universal perturbation that can effectively generalize to diverse downstream V+L tasks. While this work lies in leveraging the shared and joint characteristics of Vision-Language scenarios to present a universal and versatile UAP, it is a promising direction for future studies to investigate task-level V+L characteristics to further enhance attacks for specific downstream tasks.

Synthetic positive samples. Introducing synthetic samples that are maximally distant from the anchor as positive samples is a promising direction. A reasonable implementation might involve adversarial learning to generate such maximally distant samples. However, this strategy necessitates synthesizing samples for each input pair, leading to a significant increase in the computational overhead. Future works can explore more efficient and effective positive sample strategies.

1046 1047

1048

1049

1050 1051

1067 1068

1069

1070

1071 1072

1073

### J SPECIAL TOKENS AS TEXT PERTURBATION

We also explore the potential of special tokens to serve as adversarial perturbations. Specifically, we directly adopt two typical # and \* as adversarial tokens to evaluate their attack results.

Table 16: ASR of C-PGC and its variants using special characters as the adversarial word.

Source	Adv. word	AL	BEF	T	CL	X-V	'LM	CLI	P <sub>ViT</sub>	CLI	P <sub>CNN</sub>	BI	JP
			IR	TR	IR	TR	IR	TR	IR	TR	IR	TR	IR
	#	87.81	85.74	60.84	62.05	18.28	35.87	38.67	61.4	50.91	68.21	41.71	54.11
ALBEF	*	87.21	84.87	60.24	62.19	18.01	35.4	38.79	61.91	51.27	68.1	41.92	54.69
	C-PGC	90.13	88.82	62.11	64.48	20.53	39.38	43.1	65.93	54.4	72.51	44.79	56.36
	#	21.25	37.41	24.27	41.04	14.71	34.38	87.07	92.39	63.2	75.14	25.46	44.19
<b>CLIP</b> <sub>ViT</sub>	*	22.07	37.54	24.58	41.32	14.81	34.77	87.57	92.41	63.78	74.86	25.76	44.87
	C-PGC	23.23	38.67	25.05	41.79	15.85	35.59	88.92	93.05	66.06	75.42	26.71	45.7

Table 17: Comparison of C-PGC and its variants using special characters as the adversarial word regarding the BERT score between clean and adversarial texts.

Adv word	word ALBEF				TCL			CLIP_VIT		CLIP_CNN		
ind word	Р	R	F1	Р	R	F1	P	R	F1	Р	R	F1
#	0.8213	0.8419	0.8313	0.8171	0.8389	0.8277	0.8137	0.8339	0.8235	0.8156	0.8364	0.8257
*	0.8149	0.8251	0.8197	0.8098	0.8206	0.8149	0.8095	0.8206	0.8148	0.8097	0.8203	0.8147
C-PGC	0.8891	0.8613	0.8748	0.8924	0.8687	0.8802	0.8746	0.8684	0.8713	0.8948	0.8842	0.8893

Table 16 and Table 17 display that our optimization-based strategy exhibits both superior attack performance and higher semantic similarity. Future studies can investigate more special tokens to increase the likelihood of bypassing human observers or automated filtering systems.

#### K MORE VISUALIZATION RESULTS

This section presents a rich visual analysis of the proposed attack on a series of downstream tasks.
Specifically, we generate the UAP and conduct attacks on the ALBEF model in the visual grounding (VG) task. As illustrated in Figure 8, the prediction bounding boxes exhibit a notable deviation from the clean predictions, verifying that our generated adversarial samples significantly interfere with the multimodal alignment. In the visual entailment (VE) task, we employ BLIP as the victim model and present the results in Figure 9. These qualitative visualizations again demonstrate the remarkable attack effects of our proposed method on various downstream tasks.



Figure 8: Illustration of visual grounding. Predictions of clean image-text pairs are on the left while the adversarial samples are on the right. The red indicates the universal adversarial word.

1133



Figure 9: Illustration of the visual entailment task. The red indicates the universal adversarial word. It can be observed that all predictions do not match with the ground truth.