

# AdvCLIP: Downstream-agnostic Adversarial Examples in Multimodal Contrastive Learning

Ziqi Zhou\*<sup>†‡§</sup>  
zhouziqi@hust.edu.cn  
School of Cyber Science and  
Engineering, Huazhong University of  
Science and Technology

Shengshan Hu\*<sup>†‡§</sup>  
hushengshan@hust.edu.cn  
School of Cyber Science and  
Engineering, Huazhong University of  
Science and Technology

Minghui Li  
minghuili@hust.edu.cn  
School of Software Engineering,  
Huazhong University of Science and  
Technology

Hangtao Zhang  
zhanghangtao7@163.com  
School of Cyber Science and  
Engineering, Huazhong University of  
Science and Technology

Yechao Zhang\*<sup>†‡§</sup>  
ycz@hust.edu.cn  
School of Cyber Science and  
Engineering, Huazhong University of  
Science and Technology

Hai Jin\*<sup>†¶</sup>  
hjin@hust.edu.cn  
School of Computer Science and  
Technology, Huazhong University of  
Science and Technology

## ABSTRACT

Multimodal contrastive learning aims to train a general-purpose feature extractor, such as CLIP, on vast amounts of raw, unlabeled paired image-text data. This can greatly benefit various complex downstream tasks, including cross-modal image-text retrieval and image classification. Despite its promising prospect, the security issue of cross-modal pre-trained encoder has not been fully explored yet, especially when the pre-trained encoder is publicly available for commercial use.

In this work, we propose AdvCLIP, the first attack framework for generating downstream-agnostic adversarial examples based on cross-modal pre-trained encoders. AdvCLIP aims to construct a universal adversarial patch for a set of natural images that can fool all the downstream tasks inheriting the victim cross-modal pre-trained encoder. To address the challenges of heterogeneity between different modalities and unknown downstream tasks, we first build a topological graph structure to capture the relevant positions between target samples and their neighbors. Then, we design a topology-deviation based generative adversarial network to generate a universal adversarial patch. By adding the patch to images, we minimize their embeddings similarity to different modality and perturb the sample distribution in the feature space, achieving universal non-targeted attacks. Our results demonstrate the excellent attack performance of AdvCLIP on two types of downstream tasks across eight datasets. We also tailor three popular defenses to

mitigate AdvCLIP, highlighting the need for new defense mechanisms to defend cross-modal pre-trained encoders. Our codes are available at: <https://github.com/CGCL-codes/AdvCLIP>.

## CCS CONCEPTS

• **Security and privacy**; • **Computing methodologies** → Computer vision; • **Information systems** → *Information retrieval*;

## KEYWORDS

Adversarial Patch, Pre-trained Encoder, Cross-modal Retrieval

### ACM Reference Format:

Ziqi Zhou, Shengshan Hu, Minghui Li, Hangtao Zhang, Yechao Zhang, and Hai Jin. 2023. AdvCLIP: Downstream-agnostic Adversarial Examples in Multimodal Contrastive Learning. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

With recent advancements in deep learning, multimodal pre-training has emerged as a promising area of research for various downstream tasks. Multimodal contrastive learning [39, 53] is a novel machine learning paradigm to overcome the restrictions of labeled data. It uses large-scale, noisy, and unprocessed multimodal data pairs sourced from the web to train a cross-modal pre-trained encoder, such as CLIP [39], with powerful feature extraction capabilities. By fine-tuning these pre-trained encoders with a small amount of labeled data, complex and diverse downstream tasks can be performed [27, 54]. This pre-training approach provides a solution for resource-constrained users to benefit from large-scale models by using their powerful zero-shot capabilities directly or fine-tuning a linear layer for various downstream tasks with less data and computational resources. Driven by this promising prospect, many service providers have unveiled their pre-trained encoders such as CLIP [39], ALBEF [25], and GPT [4], or have deployed them as commercial services, like ChatGPT.

Meanwhile, it is well-known that machine learning models are susceptible to various adversarial attacks [15, 33, 52], which will make pre-trained encoders fragile as well. With pre-trained encoders are widely used, the risks associated with them are often

\*National Engineering Research Center for Big Data Technology and System

<sup>†</sup>Services Computing Technology and System Lab

<sup>‡</sup>Hubei Key Laboratory of Distributed System Security

<sup>§</sup>Hubei Engineering Research Center on Big Data Security

<sup>¶</sup>Cluster and Grid Computing Lab

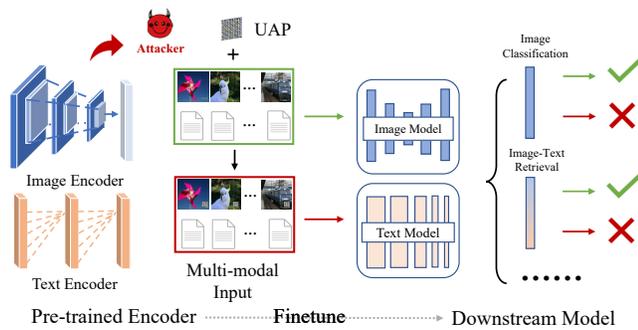
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*Conference'17, July 2017, Washington, DC, USA*

© 2023 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>



**Figure 1: Illustration of attackers using a cross-modal pre-trained encoder to attack different downstream tasks**

inherited by downstream tasks. Recent works [9, 20, 21, 29, 30, 56] paid attention to the privacy and robustness concerns of unimodal pre-trained encoders, however, the security threat of more widely used cross-modal pre-trained encoders (e.g., *Vision-language Pre-trained* (VLP) encoders [23, 39]) remains unexplored. Although a recent study [55] tried to conduct adversarial attacks against downstream tasks of VLP encoders, it relied on unrealistic white-box assumptions to generate sample-specific adversarial examples. In the literature, the difficulty of cross-modal attacks, caused by the heterogeneity between different modalities, has created an illusory sense of security for cross-modal pre-trained encoders. It tends to become a common belief that it is impossible to realize cross-modal attacks without the knowledge of the pre-training dataset, the downstream dataset, task type, or even the defense strategy that the downstream model is taking. *To the best of our knowledge, implementing adversarial attacks in practical multimodal pre-training scenarios remains a challenging and unsolved problem.*

In this paper, we propose AdvCLIP, the first attack framework for generating downstream-agnostic adversarial examples, to break the existing illusion of security in cross-modal pre-trained encoders. Given the limited knowledge of attackers and the feasibility of attack implementation, the goal is to achieve universal non-targeted attacks based on images for downstream tasks. There are two types of universal adversarial attacks: perturbation-based and patch-based methods. The former requires adding perturbations to the image globally, the latter is limited to a small area of the image and is more easily applicable to the physical world. Therefore, we mainly focus on adversarial patch attacks. The most daunting challenge in this work is to effectively tackle the modality gap between image and text, while simultaneously bridging the attack gap between cross-modal pre-trained encoders and downstream tasks.

Based on the intuition of maximizing the distance between the target image features and their corresponding benign image and text features, we first construct a topological graph structure to capture the similarity between samples. Then, we fool the pre-trained encoders by destroying the mapping relationship between different modalities of a single sample and the topological relations between multiple samples, respectively. To achieve attack transferability from the pre-trained encoder to the downstream task, we make the adversarial examples far from the original class rather than simply crossing the decision boundary. As a result, we design a

topology-deviation based generative adversarial network to generate a universal adversarial patch to achieve high attack success rate attacks for downstream tasks with a fixed random noise as input. Our main contributions are summarized as follows:

- We propose AdvCLIP, the first attack framework to construct downstream-agnostic adversarial examples in multimodal contrastive learning. We reveal that the cross-modal pre-trained encoder incurs severe security risks for the downstream tasks.
- We design a topology-deviation based generative adversarial network, which adds a universal adversarial patch to the target image, to decrease the similarity between different modal embeddings and disrupt their topological relationships, achieving non-targeted adversarial attacks.
- Our extensive experiments on two types of downstream tasks over eight datasets show that our AdvCLIP addresses the modality gap and transferability between the pre-trained encoders and downstream tasks.
- We tailor three popular defenses to mitigate AdvCLIP. The results further prove the attack ability of AdvCLIP and highlight the needs of new defense mechanism to defend pre-trained encoders.

## 2 RELATED WORK

### 2.1 Vision-Language Pre-trained Models

Multimodal contrastive learning is a training paradigm that aims to pre-train encoders on large-scale unlabeled training data to obtain general-purpose representations for application to downstream tasks. The success of multimodal contrastive learning has motivated the development of numerous VLP models [11, 13, 26, 39] for building multimodal models capable of learning vision-language semantic alignments and solving complex cross-modal tasks [54]. Existing VLP models can be broadly categorized into two groups: cross-encoder based and embedding-based methods. Cross-encoder based methods [6, 26, 31] employ a Transformer-based cross-attention mechanism to compute the similarity between data from different modalities. In contrast, embedding-based methods [13, 39, 46] encode data from different modalities separately to generate high-dimensional visual and textual representations and measure cross-modal similarity by computing feature distances between data from different modalities. Recently, the embedding-based CLIP [39] has demonstrated exceptional performance on various downstream tasks. In this paper, we focus on the security of CLIP.

### 2.2 Universal Adversarial Attack

*Universal adversarial perturbation* (UAP) [33] was proposed to deceive a target model by applying a single adversarial noise to all input images. Universal image-based adversarial attacks come in two forms: perturbations [10, 16–18, 28, 33, 34] and patches [3, 19, 22, 47]. Perturbation-based methods fool models by adding visually imperceptible noise globally to the image. In contrast, patch-based methods require precise control over each pixel of the image, resulting in visible adversarial patches that are limited to a small region of the image. On the other hand, text-based adversarial attacks require a different approach due to the discrete nature of text.

Consequently, universal text-based attacks focus on generating imperceptible triggers and linguistic idioms to create adversarial examples [2, 44, 49]. Unfortunately, the current UAP methods are mostly designed for unimodal classification tasks and are insufficient for attacking cross-modal tasks, let alone when the attacker’s knowledge about downstream tasks is limited. As patch-based image adversarial examples are more applicable to real-world scenarios, this paper focuses on universal adversarial patch. Additionally, researchers have proposed different defenses against adversarial examples, such as data pre-processing [5], adversarial training [32, 43, 48], and pruning [50, 58].

### 2.3 Adversarial Attacks on Pre-trained Encoders

Recently, an increasing number of works [12, 36, 57] begin to investigate the robustness of pre-trained encoders. PAP [1] produced a pre-trained perturbation by lifting the feature activations of low-level layers against image pre-trained encoders. At the same time, some works made trivial attempts to explore the vulnerability of cross-modal pre-trained encoders. One recent study [57] examined the robustness of a CLIP-based image-text retrieval system. Furthermore, Co-Attack [55] took a step by considering the robustness of downstream tasks corresponding to cross-modal pre-trained encoders. It proposed a cooperative loss function to avoid conflicts caused by simultaneous attacks on both image and text modalities. However, it only considered simple white-box scenarios where the attacker has downstream knowledge to generate sample-specific adversarial examples. As a result, insufficient research on cross-modal pre-training safety tends to create a false sense of security in the field. Our work aims to achieve effective ignorant attacks against downstream tasks and break the illusion of security in cross-modal pre-trained encoders.

## 3 METHODOLOGY

### 3.1 Threat Model

We assume a quasi-black-box attack model, where the attacker has access to VLP encoders through purchasing or downloading from publicly available websites, but lacks knowledge of the pre-training datasets and downstream tasks. As the attacker does not possess specific target information of downstream tasks, their objective is to conduct non-targeted adversarial attacks that disable or reduce the accuracy of downstream tasks. To achieve this, the attacker leverages the pre-trained encoder to design a downstream-agnostic universal adversarial patch that is applicable to various types of input images from different datasets. Then the adversarial examples can mislead all the downstream tasks that inherit the victim pre-trained encoder, such as image-text retrieval, image classification, etc. We assume that the downstream task undertaker (called user hereinafter) is able to fine-tune the linear layer for their cause. Given the complexity of CLIP training and its powerful zero-shot performance, we believe that users do not need to fine-tune CLIP directly, as doing so would negate the benefits of choosing it in the first place. We also consider a more stringent scenario, in which users employ common defense mechanisms such as adversarial training to improve the robustness of downstream models.

### 3.2 Problem Formulation

Let  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  denote a cross-modal dataset with  $N$  instances. Here,  $x_i = \{(x_i^v, x_i^t)\}$ , where  $x_i^v$  and  $x_i^t$  represent two data modalities, such as image-text pairs, and they both belong to the same label  $y_i$ . Let  $L = \{y_i\}_{i=1}^C$  represent the label dataset from  $\mathcal{D}$ , where  $C$  is the number of labels and  $C < N$ . Given an input  $x_i \in \mathcal{D}_a$  to a cross-modal pre-trained encoder  $M_\theta(\cdot)$  (i.e., CLIP [39]) which consists of an image encoder  $E_v(\cdot)$  and a text encoder  $E_t(\cdot)$ , that returns an image feature vector  $v_v$  and a text feature vector  $v_t$  respectively, where  $\theta$  denotes the parameter of the cross-modal pre-trained encoder. The attacker utilizes a surrogate dataset  $\mathcal{D}_a$  that is distinct from both the pre-training dataset  $\mathcal{D}_p$  and the downstream dataset  $\mathcal{D}_d$  to generate a universal adversarial noise against the pre-trained encoder. Moreover, the universal adversarial noise  $\delta$  should be sufficiently small, and modeled through an upper-bound  $\epsilon$  on the  $l_p$ -norm. This problem can be formulated as:

$$M_\theta(x_i + \delta) \neq M_\theta(x_i), \quad s.t. \|\delta\|_p \leq \epsilon \quad (1)$$

With the help of the strong feature extraction ability of cross-modal pre-trained encoders, we can just fine-tune a linear layer using output feature vectors of different modalities to achieve complex downstream tasks. In this paper, we mainly consider cross-modal image-text retrieval and unimodal image classification tasks. For the cross-modal retrieval task, the cross-modal retrieval head  $c_{\theta'}(\cdot)$  completes the image-text retrieval task based on the similarity between  $v_v$  and  $v_t$ , where  $\theta'$  denotes the parameter of the retrieval head. The attacker’s goal is to implement a non-targeted attack that fools the downstream cross-modal retrieval head  $c_{\theta'}(\cdot)$  by applying a universal adversarial noise  $\delta$  to the downstream sample  $x \in \mathcal{D}_d$ . Therefore, the attacker’s goal can be formalized as:

$$c_{\theta'}(E_v(x_i^v + \delta), E_t(x_i^t)) \neq c_{\theta'}(E_v(x_i^v), E_t(x_i^t)), \quad s.t. \|\delta\|_p \leq \epsilon \quad (2)$$

Similarly, the objective of attackers against a downstream image classification task can be expressed as:

$$f_{\theta''}(E_v(x_i^v + \delta)) \neq f_{\theta''}(E_v(x_i^v)), \quad s.t. \|\delta\|_p \leq \epsilon \quad (3)$$

where  $f_{\theta''}(\cdot)$  is a classifier,  $\theta''$  is the parameter of the classifier.

### 3.3 Intuition Behind AdvCLIP

Due to the limitations of the attacker’s knowledge and the complexity of cross-modal tasks, achieving effective attacks on unknown downstream tasks has to address the following challenges:

**Challenge I: Modality gap between image and text.** Traditional adversarial attacks are designed for unimodal classification tasks. However, the VLP encoder involves multiple modalities and its output is a high-dimensional feature embedding rather than a label, directly applying traditional adversarial attack methods is impractical. As attackers aim to launch non-targeted adversarial attacks on downstream tasks, a natural idea is to disrupt the similarity matching process by maximizing the feature distance between the adversarial and corresponding clean embeddings of different modalities. However, it is a challenging problem to understand and utilize the high-dimensional feature vectors to produce adversarial examples. As shown in Fig. 2, simply maximizing the distance between embeddings (Vanilla) may not work due to the complexity of the high-dimensional feature space and the heterogeneity of multimodal data. Even if the image adversarial example leaves its

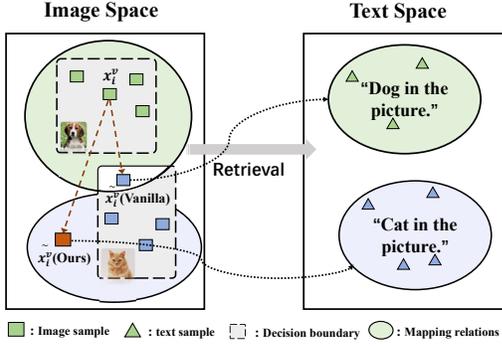


Figure 2: Attack gap between image and text modality

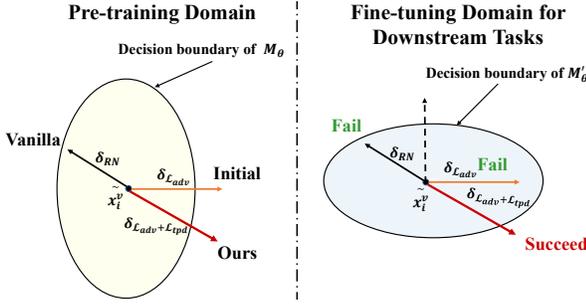


Figure 3: Transferability gap between cross-modal pre-trained encoders and downstream tasks

original category in the image feature space and is classified as a cat, yet it will still be retrieved to match the textual information of the original category of dogs.

In this paper, on the basis of leaving the original position in the feature space, we consider destroying the nearest neighbor relationship of the samples to better reinforce the attack by making the ordered samples in the feature space disordered. Specifically, we first construct topology for both adversarial and benign embeddings separately to measure the corresponding sample correlations. Topology is based on the neighborhood relation graph constructed by the similarity between samples in the representation space. The process of measuring topological similarity can be formalized as:

$$\mathcal{L}_{tp} = \mathbb{E}_{(x,y) \in \mathcal{D}_a} (CE(\mathcal{G}_{nor}, \mathcal{G}_{adv})) \quad (4)$$

where  $\mathcal{G}_{nor}$  and  $\mathcal{G}_{adv}$  stand for the neighbourhood relation graph constructed by the inter-sample similarity for clean samples and adversarial examples, respectively.  $CE(\cdot)$  is the cross-entropy loss to measure the similarity of two graphs.

We define the edge weights of the neighborhood graph as the probability that two different samples are neighbors, and the deviation of the topological structure is achieved by warping the probability distributions of two graphs. Then, we model the conditional probability distribution using an affinity measure based on cosine similarity to construct the adjacency graph, and remove the nearest neighbor points to prevent isolated subgraphs formed by data points with excessively high local density, thereby ensuring the local connectivity of the manifold and better preserving the global structure. The process of constructing the adjacency graph can be represented as:

$$\mathcal{G} = \left\{ p_{i|j} \mid p_{i|j} = \frac{(2 - (d_{ij} - \rho_j))}{\sum_{k=1, k \neq j}^N (2 - (d_{jk} - \rho_j))}, 0 < i, j \leq N \right\} \quad (5)$$

where  $p_{i|j}$  is the conditional probability that the  $i_{th}$  natural sample is the neighbor of the  $j_{th}$  natural sample in the feature space of  $\mathcal{G}$ ,  $\rho_j$  represents the cosine distance from the  $j_{th}$  data point to its nearest neighbor,  $d_{ij}$  denotes the cosine distance between the corresponding embeddings of the two samples. By deviating from the two dimensions of the sample itself and the nearest neighbor relationship, we destroy the similarity mapping relationship between the sample and its counterpart to achieve an effective attack.

**Challenge II: Transferability gap between cross-modal pre-trained encoders and downstream models.** As illustrated in Fig. 3, after fine-tuning the cross-modal pre-trained encoder to the downstream model, the boundary of the feature space in the model may change, which could make existing attacks ineffective. Therefore, we aim to deviate adversarial examples from the direction most likely to cross their original category boundaries within the given perturbation budget. To address this challenge, we are motivated to make adversarial examples deviate from the direction that is most likely to leave their original category boundaries under the same perturbation budget. Inspired by the fact that generative adversarial networks can generate samples with similar salient features [14], we design a generative adversarial network to generate a universal adversarial noise with strong commonality, such that the adversarial examples are far from the original category rather than only just crossing the decision boundary of that category. In this way, even if the users fine-tune the pre-trained encoder to the downstream model, the adversarial examples still cannot be recognized properly.

### 3.4 Topology-deviation based Generative Attack Framework

In this section, we present AdvCLIP, a novel generative attack against cross-modal pre-trained encoders. The framework of AdvCLIP is depicted in Fig. 4. It consists of an adversarial generator  $G$ , a discriminator  $D$ , and a victim cross-modal encoder  $M$  which consists of an image encoder  $E_v$  and a text encoder  $E_t$ . Given the image-text pairs  $(x_i^v, x_i^t)$  to the cross-modal pre-trained encoder, the image encoder  $E_v$  and text encoder  $E_t$  output the corresponding feature vectors. We design a topology-deviation based generative attack framework, which utilizes cross-modal pre-trained encoders to generate universal adversarial patches applicable to images, thereby deceiving downstream tasks.

**Adversarial Generator.** By feeding a fixed noise  $z$  into the adversarial generator, we obtain a universal adversarial patch  $G(z)$  and paste it onto an image of the surrogate dataset  $\mathcal{D}_a$  to get an adversarial example  $\tilde{x}_i^v$ . The above process of making adversarial examples can be formalized as:

$$\tilde{x}_i^v = x_i^v \odot (1 - m) + G(z) \odot m \quad (6)$$

where  $\odot$  denotes the element-wise product,  $m$  is a binary matrix that contains the position information of the patch.

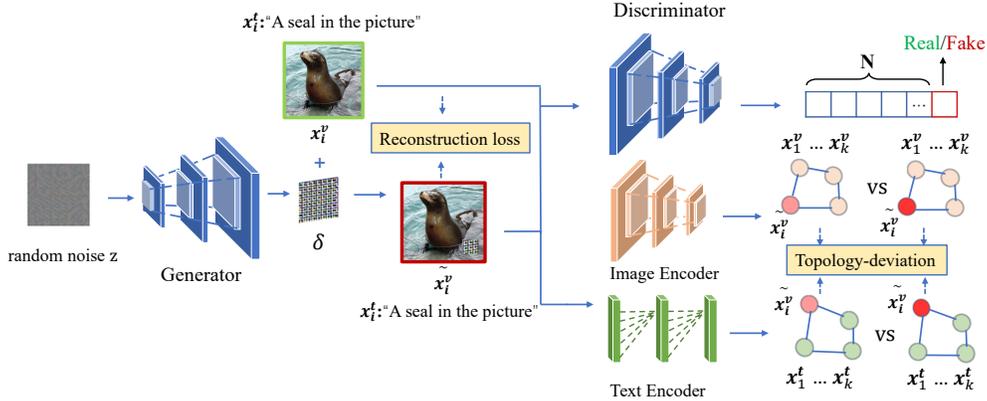


Figure 4: The framework of our attack

The objective function of the adversarial generator  $G$  is:

$$\min_{\theta_G} \mathcal{L}_G = \sum_{(x_i, y_i) \in \mathcal{D}_a} \left( \alpha \mathcal{L}_{adv} + \beta \mathcal{L}_{tpd} + \mathcal{L}_q + \mathcal{L}_{gan} \right) \quad (7)$$

where  $\mathcal{L}_{adv}$  is the adversarial loss function,  $\mathcal{L}_{tpd}$  is the topology-deviation loss function,  $\mathcal{L}_q$  is the quality loss function,  $\mathcal{L}_{gan}$  is the GAN loss function, and  $\alpha, \beta$  are pre-defined hyper-parameters.

The adversarial loss  $\mathcal{L}_{adv}$  is used to deviate the feature position of the target sample, by adding a patch to an image  $x_i^p$  so that the feature vector  $E_v(\tilde{x}_i^v)$  of the adversarial example  $\tilde{x}_i^v$  is simultaneously far away from the original image feature vector  $E_v(x_i^p)$  and the clean text feature vector  $E_t(x_i^t)$ . Thus  $\mathcal{L}_{adv}$  is expressed as:

$$\mathcal{L}_{adv} = \mathcal{L}_{av} + \lambda \mathcal{L}_{at} \quad (8)$$

where  $\mathcal{L}_{av}$  and  $\mathcal{L}_{at}$  denote the image-image semantic feature deviation loss and the image-text semantic deviation loss, respectively. We adopt InfoNCE [37] loss to measure the similarity between the vectors output by encoders. Specifically, we first treat the vector of benign image  $x_i^p$  and adversarial image  $\tilde{x}_i^v$  as negative pairs, pulling away their feature distance. It can be expressed as:

$$\mathcal{L}_{av} = \log \left[ \frac{\exp \left( \text{Sim} \left( E_v(\tilde{x}_i^v), E_v(x_i^p) \right) / \tau \right)}{\sum_{j=0}^K \exp \left( \text{Sim} \left( E_v(\tilde{x}_i^v), E_v(x_j^p) \right) / \tau \right)} \right] \quad (9)$$

where  $\text{Sim}(\cdot)$  represents the cosine distance function,  $\tau$  denotes a temperature parameter. Then we treat the vectors of adversarial image samples  $\tilde{x}_i^v$  and benign text samples  $x_i^t$  as negative pairs similarly and increase their feature distances. So we have:

$$\mathcal{L}_{at} = \log \left[ \frac{\exp \left( \text{Sim} \left( E_v(\tilde{x}_i^v), E_t(x_i^t) \right) / \tau \right)}{\sum_{j=0}^K \exp \left( \text{Sim} \left( E_v(\tilde{x}_i^v), E_t(x_j^t) \right) / \tau \right)} \right] \quad (10)$$

The topology-deviation loss  $\mathcal{L}_{adv}$  is designed to corrupt the topological similarity between the adversarial examples and their corresponding normal samples, *i.e.*, the neighbourhood relation graph constructed based on the similarity between samples in the representation space. Similarly, we deviate both the image feature vectors  $E_v(\tilde{x}_i^v)$  of the adversarial examples and the image feature

vectors  $E_v(x_i^p)$  and text feature vectors  $E_t(x_i^t)$  of the normal samples. Our goal is to maximize the topological distance between them, which can be represented as:

$$\mathcal{L}_{tpd} = -(\mathcal{L}_{tp}(E_v(\tilde{x}_i^v), E_v(x_i^p)) + \lambda \mathcal{L}_{tp}(E_v(\tilde{x}_i^v), E_t(x_i^t))) \quad (11)$$

To achieve better stealthiness, we use  $\mathcal{L}_q$  to control the magnitude of the adversarial noises output by the generator and crop  $\delta$  after each optimisation to ensure it meets the constraints  $\epsilon$ . Formally, we have:

$$\mathcal{L}_q = \left\| \tilde{x}_i^v - x_i^p \right\|_2 \quad (12)$$

The GAN loss  $\mathcal{L}_{gan}$  encourages adversarial examples to be more visually natural. That is, an normal image and an adversarial example with adversarial patch tend to be consistent on the discriminator. Thus the GAN loss  $\mathcal{L}_{gan}$  can be expressed as:

$$\mathcal{L}_{gan} = \log \left( 1 - D(\tilde{x}_i^v) \right) \quad (13)$$

**Discriminator.** The main function of discriminator is to identify the authenticity of fake examples generated by the adversarial generator. By playing games with the generator, we ensure that the generated fake adversarial examples are visually indistinguishable from the real ones. The objective loss function of  $D$  is:

$$\min_{\theta_D} \mathcal{L}_D = \sum_{(x_i, y_i) \in \mathcal{D}_a} -(\log(D(x_i^p)) + \log(1 - D(\tilde{x}_i^v))) \quad (14)$$

## 4 EXPERIMENTS

### 4.1 Experimental Setting

**Victim Pre-trained Encoders.** We choose CLIP [39] as the victim encoder for our experiments and obtain all pre-trained encoders from its publicly available repository. We evaluate the vulnerability of CLIP to adversarial attacks across a range of architectures, including ResNet50, ResNet101, ViT-L/14, ViT-B/16, and ViT-B/32. **Downstream Datasets.** We evaluate the effectiveness of our attacks on two distinct downstream tasks: image-text retrieval and image classification. To carry out the image-text retrieval task, we select four widely used cross-modal datasets, namely Wikipedia [40], Pascal-Sentence [41], NUS-WIDE [7], and XmediaNet [38]. For the

**Table 1: The cross-modal attack performance (%) of AdvCLIP under different settings.  $\mathcal{D}_1 - \mathcal{D}_4$  denote the settings where the downstream datasets are NUS-WIDE, Pascal-Sentence, Wikipedia, and XmediaNet, respectively.**

Surrogate	Dataset	ResNet50			ResNet101			ViT-B/16			ViT-B/32			ViT-L/14		
		$ASR_i$	$ASR_t$	AVG	$ASR_i$	$ASR_t$	AVG	$ASR_i$	$ASR_t$	AVG	$ASR_i$	$ASR_t$	AVG	$ASR_i$	$ASR_t$	AVG
NUS-WIDE	$\mathcal{D}_1$	45.20	17.00	31.10	36.40	4.20	20.30	67.25	57.55	62.40	43.50	8.15	25.82	45.50	5.25	25.38
	$\mathcal{D}_2$	67.00	58.50	62.75	22.50	37.00	29.75	66.00	65.50	65.75	52.00	43.00	47.50	31.00	19.50	25.25
	$\mathcal{D}_3$	45.24	43.73	44.48	30.09	18.40	24.25	54.76	62.34	58.55	27.49	28.57	28.03	32.25	16.45	24.35
	$\mathcal{D}_4$	57.10	47.94	52.52	59.45	35.89	47.67	80.05	63.15	71.60	65.45	42.41	53.93	54.07	11.95	33.01
Pascal	$\mathcal{D}_1$	23.25	8.15	15.70	25.90	2.55	14.22	62.65	60.00	61.33	32.70	5.45	19.07	37.75	4.25	21.00
	$\mathcal{D}_2$	36.00	26.00	31.00	31.00	22.00	26.50	67.50	63.00	65.25	49.50	43.50	46.50	53.50	56.00	54.75
	$\mathcal{D}_3$	13.86	17.32	15.59	13.63	9.52	11.57	51.08	55.20	53.14	26.63	18.39	22.51	37.01	21.42	29.21
	$\mathcal{D}_4$	30.29	18.26	24.27	49.54	18.42	33.98	80.40	62.46	71.43	63.06	32.46	47.76	79.27	49.24	64.25
Wikipedia	$\mathcal{D}_1$	32.00	7.25	19.62	33.80	0.60	17.20	63.55	53.15	58.35	23.45	13.90	18.68	51.45	38.25	44.85
	$\mathcal{D}_2$	25.50	40.50	33.00	8.00	13.50	10.75	67.50	64.50	66.00	51.00	52.50	51.75	53.50	48.50	51.00
	$\mathcal{D}_3$	25.11	21.00	23.05	16.45	9.96	13.21	55.84	62.99	59.41	36.15	27.27	31.71	53.46	32.03	42.74
	$\mathcal{D}_4$	34.50	29.86	32.18	47.06	20.07	33.56	80.01	63.33	71.67	56.89	33.29	45.09	82.06	56.20	69.13
XmediaNet	$\mathcal{D}_1$	43.20	11.30	27.25	8.66	4.33	6.50	46.53	23.16	34.84	40.05	38.52	39.28	37.66	17.96	27.81
	$\mathcal{D}_2$	59.00	62.50	60.75	36.90	7.05	21.98	58.05	5.25	31.65	45.55	3.35	24.45	33.80	9.15	21.48
	$\mathcal{D}_3$	42.64	41.78	42.21	27.50	14.00	20.75	59.50	54.00	56.75	57.50	42.50	50.00	45.50	30.50	38.00
	$\mathcal{D}_4$	53.76	49.03	51.40	61.66	27.46	44.56	77.84	34.90	56.37	67.32	44.06	55.69	78.58	32.73	55.66

**Table 2: The unimodal attack performance (%) of AdvCLIP under different settings.  $\mathcal{V}_1 - \mathcal{V}_5$  denote the settings where the victim models are ResNet50, ResNet101, ViT-B/16, ViT-B/32, and ViT-L/14, respectively.**

Surrogate	Victim	CIFAR10		GTSRB		ImageNet		NUS-WIDE		Pascal		STL10		Wikipedia		XmediaNet	
		FR	ASR	FR	ASR	FR	ASR	FR	ASR	FR	ASR	FR	ASR	FR	ASR	FR	ASR
NUS-WIDE	$\mathcal{V}_1$	89.73	65.01	90.36	67.86	94.01	75.40	76.81	56.98	78.12	50.78	71.10	63.24	81.67	45.87	91.14	78.85
	$\mathcal{V}_2$	78.55	59.57	86.75	62.35	60.41	44.23	49.80	29.59	64.45	35.55	17.80	11.60	66.07	31.50	59.33	48.47
	$\mathcal{V}_3$	98.00	83.49	95.85	80.30	99.73	88.80	97.37	73.06	98.50	72.00	98.04	95.56	96.00	59.68	99.96	89.78
	$\mathcal{V}_4$	87.91	73.65	91.25	68.79	96.21	77.52	68.51	49.02	91.41	61.33	78.85	70.83	77.96	41.27	91.53	79.24
	$\mathcal{V}_5$	43.95	34.73	90.77	71.06	89.87	74.51	55.76	36.91	63.67	34.38	52.71	48.45	81.72	37.47	80.84	70.24
	AVG	79.63	63.29	91.00	70.07	88.05	72.09	69.65	49.11	79.23	50.81	63.70	57.94	80.68	43.16	84.56	73.32
XmediaNet	$\mathcal{V}_1$	86.37	61.87	96.74	74.67	87.68	72.25	75.44	52.73	80.08	50.78	71.80	68.90	65.54	34.15	90.39	80.00
	$\mathcal{V}_2$	90.33	67.39	73.25	63.76	55.67	49.50	80.91	56.01	86.72	62.11	37.01	35.98	60.24	30.92	62.05	56.88
	$\mathcal{V}_3$	55.03	41.34	74.89	65.51	74.08	67.39	83.20	57.67	50.00	30.08	21.76	20.93	73.49	38.48	89.74	83.96
	$\mathcal{V}_4$	83.82	69.98	91.69	81.98	95.41	88.62	87.65	62.40	89.45	59.77	78.45	77.54	90.09	57.14	98.30	92.09
	$\mathcal{V}_5$	60.10	52.48	91.22	81.76	91.57	84.66	77.64	53.13	79.30	51.56	49.91	49.05	70.03	37.67	90.02	83.98
	AVG	75.13	58.61	85.56	73.54	80.88	72.48	80.97	56.39	77.11	50.86	51.79	50.48	71.88	39.67	86.10	79.38

image classification task, we additionally choose STL10 [8], GTSRB [45], CIFAR10 [24], and ImageNet [42] image datasets. Note that our approach is to generate image adversarial patches using the cross-modal datasets.

**Evaluation Metrics.** In the cross-modal retrieval task, We use the standard evaluation metric, *mean average precision* (MAP) [59], to evaluate the accuracy of models, which we report separately for two sub-tasks: text retrieval with image queries (I2T) and image retrieval with text queries (T2I). To measure the performance of our attacks, following [51], we use the *attack success rate* (ASR), which is calculated as the difference between the MAP values of normal samples and adversarial examples, with  $ASR_i$  and  $ASR_t$  used for image-text retrieval and text-image retrieval, respectively. For the classification task, we evaluate our attacks using three metrics: *clean accuracy* (CA), *attack success rate* (ASR), and *fooling rate* (FR). CA measures normal accuracy, ASR is calculated as described above, and FR is the percentage of misclassified examples compared to the

total number of test examples. Higher ASR and FR values indicate better attack effectiveness.

## 4.2 Attack Performance

**Implementation Details.** In order to evaluate the effectiveness of AdvCLIP in scenarios where the downstream task is unknown, we conduct experiments on two distinct tasks: image-text retrieval and image classification. Following [19, 33], we set the perturbation upper bound (the noise percentage of each sample)  $\epsilon$  of adversarial patch to 0.03. We choose the bottom right corner of the image, which is not easily visible, to apply the patch. We set the hyper-parameters  $\alpha = 10$ ,  $\beta = 5$ ,  $\lambda = 1$  and the training epoch to 20 with batch size of 16. We use four cross-modal datasets as attacker surrogate datasets to train generative adversarial networks to generate a universal adversarial patch, which are then used to attack different downstream tasks. The generator and discriminator network are



**Figure 5: Adversarial examples generated by AdvCLIP based on XmediaNet**

trained by Adam optimizer with the initial learning rate 0.0002. Examples of generated adversarial patches are shown in Fig. 5.

For the image-text retrieval task, we comprehensively evaluate the attack performance of AdvCLIP on four downstream cross-modal datasets (NUS-WIDE, Pascal-Sentence, Wikipedia, and XmediaNe). We evaluate the performance of AdvCLIP in two subtasks of image-text retrieval using  $ASR_i$  and  $ASR_t$ , respectively. For the image classification task, in addition to using images from the above cross-modal dataset, we also use four commonly used image datasets (CIFAR10, STL10, GTSRB, and ImageNet) to train the classification model. We use  $ASR$  and  $FR$  to evaluate AdvCLIP’s ability in the classification task.

**Analysis.** Our results in the image-text retrieval downstream task demonstrate a significant security threat posed by cross-modal pre-trained encoders. Firstly, Tab. 1 shows that different types of backbones have varying vulnerability to adversarial patches under the same attack settings, with the Transformer architecture being more susceptible to successful attacks than ResNet. Secondly, the surrogate dataset has a significant effect on downstream attack success. Datasets such as NUS-WIDE and XmediaNet, which contain a larger number of samples, tend to result in higher attack success rates. Thirdly, the attack performance may not be optimal when the surrogate dataset is consistent with the downstream dataset. For attackers, creating a better surrogate dataset is an essential factor in achieving success in unknown downstream tasks. As shown in Tab. 2, we achieve impressive performance in the classic image classification task. The average  $FR$  value of the output results of the downstream classification model is as high as 70%, and the average accuracy drop of the model is also over 55%.

### 4.3 Ablation Study

In this section, we explore the effect of different modules, surrogate datasets, attack strengths, and batch sizes on AdvCLIP. For our experiments, we select CLIP based on ResNet50 as the victim encoder and use the NUS-WIDE dataset as the attacker surrogate dataset to launch attacks on image-text retrieval tasks.

**The Effect of  $\mathcal{L}_{adv}$  &  $\mathcal{L}_{tpd}$ .** We first analyze the effect of  $\mathcal{L}_{adv}$  and  $\mathcal{L}_{tpd}$  on our scheme, respectively. As shown in Fig. 6 (a-b), “None” indicates the removal of both, while “Adv\_only” and “Tp\_only” indicate the removal of  $\mathcal{L}_{adv}$  and  $\mathcal{L}_{tpd}$ , respectively. Our results show that using  $\mathcal{L}_{adv}$  alone is not enough to push away the feature position of the sample itself, and the ability to attack the downstream task is poor. However, by simultaneously disrupting the sample’s relative position in the feature space, ideal results can be achieved, as indicated by the “Our” results.

**The Effect of  $\epsilon$ .** We study the effect of different perturbation upper bound on the attack performance of AdvCLIP. From Fig. 6(c-d), we

**Table 3: Attack performance (%) of comparison study of downstream cross-modal attacks**

Method	ResNet50				ViT-B/16			
	Pascal		Wikipedia		Pascal		Wikipedia	
	$ASR_i$	$ASR_t$	$ASR_i$	$ASR_t$	$ASR_i$	$ASR_t$	$ASR_i$	$ASR_t$
UAP [33]	8.00	5.00	4.12	0.22	*	20.00	0.21	12.56
UPGD [10]	2.00	-	1.52	2.39	1.00	4.00	*	6.50
FFF [34]	2.00	2.00	2.39	1.30	*	2.00	0.21	6.28
SSP [35]	2.50	3.00	8.23	4.98	4.50	2.00	*	6.50
PAP-ugs [1]	1.00	0.50	5.63	0.87	2.50	*	1.08	5.85
Co-Attack [55]	-	-	-	-	14.50	5.00	9.30	8.23
Adv-Patch [3]	30.00	20.50	9.10	10.61	15.00	16.00	27.05	19.70
Ours	<b>67.00</b>	<b>58.50</b>	<b>45.24</b>	<b>43.73</b>	<b>66.00</b>	<b>65.50</b>	<b>54.76</b>	<b>62.34</b>

can see that CLIP has different sensitivities to different perturbation budgets. A higher attack success rate can be achieved with a smaller patch size when  $\epsilon$  is 0.03.

**The Effect of Batch Size.** We examine the effect of different batch sizes from 4 to 128 on AdvCLIP, the results are shown in Fig. 6(e-f). To balance both attack performance and computational efficiency, we set the batch size to 16 as the default setting.

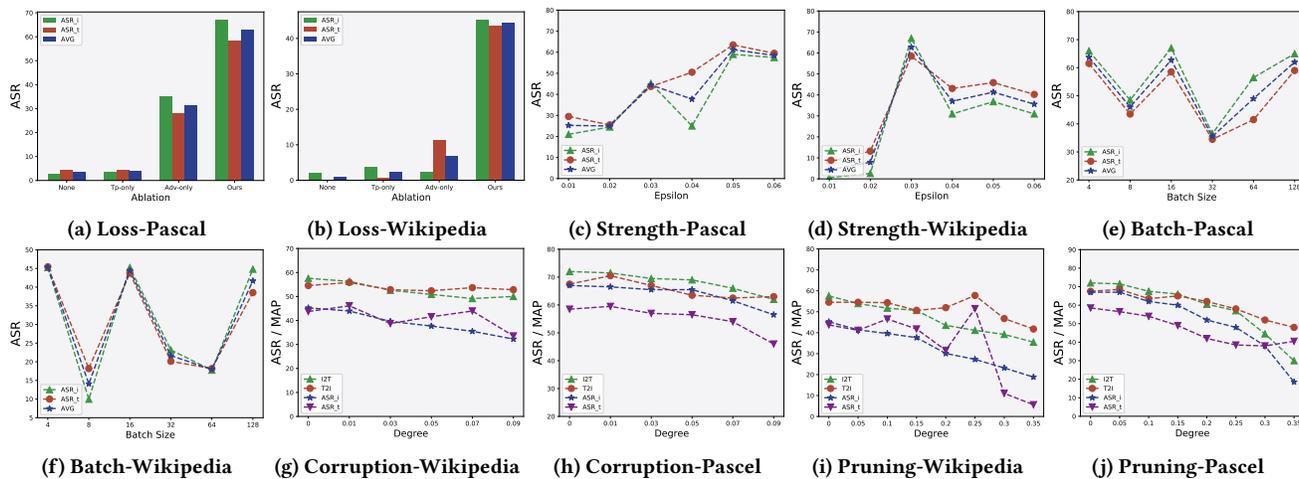
### 4.4 Comparison Study

**Implementation Details.** In this section, we compare AdvCLIP with *state-of-the-art* (SOTA) universal adversarial attacks. Prior researches have not focused enough on downstream tasks for VLP encoders. The work most relevant to ours is Co-Attack [55], which generates sample-specific adversarial examples in a white-box setting. To facilitate comparison, we randomly select an adversarial perturbation generated by Co-Attack for the images for testing. Furthermore, to demonstrate the superiority of our approach, we also consider adversarial attacks against unimodal image encoders, including PAP [1] against image pre-trained models and classic universal adversarial perturbation schemes (e.g., UAP [33], UPGD [10], FFF [34], SSP [35], and Adv-Patch [3]). To conduct a comprehensive comparison with SOTA schemes under the paradigm of cross-modal pre-trained encoders to downstream tasks, we select two representative architectures, ResNet50 and ViT-B/16, of CLIP and evaluate them on image-text retrieval downstream tasks.

**Analysis.** From Tab. 3, we can see that AdvCLIP outperforms existing SOTA methods by a large margin on two downstream datasets. Note that Co-Attack needs to use [CLS] in Transformer in the optimization process, so it cannot be used directly to attack CLIP based on ResNet50 (“-”). The negative experimental values (“\*”) indicate that the attack does not work at all. There are three reasons for this: firstly, CLIP’s robustness originates from its pre-training on a vast dataset of 400 million image-text pairs. Secondly, our quasi-black-box threat model limits the attacker’s knowledge, making it difficult to attack CLIP’s downstream model. As observed in Tab. 3, existing attacks hardly affect CLIP-based models. Lastly, unsuccessful perturbations unintentionally align input samples with CLIP’s training set, improving overall accuracy and resulting in negative ASR values for the attack.

## 5 DEFENSE

In this section, we tailor three downstream defenses for adaptively mitigating AdvCLIP. For users utilizing cross-modal pre-trained



**Figure 6: The attack performance under different settings. (a) - (f) examine the effects of different modules, attack strengths, and batch sizes on AdvCLIP, respectively. (g) - (j) investigate the effect of defense methods on AdvCLIP, respectively.**

encoders, they can preprocess input data, conduct adversarial training on downstream models, or prune parameters to defend against adversarial attacks, while maintaining normal model accuracy.

### 5.1 Corruption

Corruption is an effective and simple countermeasure for purifying adversarial examples at the pre-processing phase [5]. To combat adversarial examples, we introduce different levels of Gaussian noise to corrupt input images. As shown in Fig. 6(g-h), while maintaining that the clean samples maintain normal accuracy, the retrieval accuracy of the model decreases significantly with the degree of corruption, and the performance of AdvCLIP is slightly affected. These results indicate that AdvCLIP can effectively resist the corruption-based pre-processing defense.

### 5.2 Pruning

Pruning [58] is widely used for downstream models to inherit pre-trained encoders by removing redundant parameters in neural networks, reducing model size and computational complexity. While pruning the parameters, the required dependencies of the adversarial examples designed for the pre-trained encoder structure and parameters are broken to effectively defend against the adversarial attack. We perform parameter pruning on CLIP based on ResNet50 and evaluate the effectiveness of our attack using NUS-WIDE as surrogate dataset on two downstream cross-modal datasets. The results in Fig. 6(i-j) show that pruning parameters is difficult to effectively resist CLIP while maintaining normal model accuracy.

### 5.3 Adversarial Training

Adversarial training [15] commonly mixes adversarial examples with original data to enhance model robustness and generalization, making it more resistant to adversarial attacks. We consider a more stringent scenario where defenders conduct adversarial training on downstream tasks. Following [15], we enhance the robustness of the model during training of downstream tasks by adding noise

**Table 4: Attack performance (%) on models that have undergone adversarial training**

Dataset	RN50		RN101		ViT-B/16		ViT-B/32	
	ASR <sub>i</sub>	ASR <sub>t</sub>						
NUS-WIDE	46.75	14.95	31.75	10.70	60.65	56.30	37.90	4.20
Pascal	54.50	60.00	25.00	44.50	77.00	75.50	66.50	40.00
Wikipedia	40.91	37.45	32.04	25.33	58.44	64.93	47.84	16.45
XmediaNet	69.02	46.37	63.15	35.16	79.84	63.24	71.88	40.33

to the samples and incorporating them in the training process. For the experiment, we use NUS-WIDE as the surrogate dataset for evaluation on four downstream datasets. As shown in Tab. 4, our method is still able to successfully attack downstream tasks that have been enhanced by adversarial training.

## 6 CONCLUSION

In this paper, we propose the first attack framework to construct downstream-agnostic adversarial examples based cross-modal pre-trained encoders in multimodal contrastive learning. We design a topology-deviation based generative adversarial network that generates a universal adversarial patch to fool downstream tasks under strict constraints on attacker’s knowledge. We verify the excellent attack performance of AdvCLIP on two types of downstream tasks over five backbones of CLIP on eight datasets. We tailor three popular defenses to mitigate AdvCLIP. The results further prove the attack ability of AdvCLIP and highlight the needs of new defense mechanism to defend pre-trained encoders.

## ACKNOWLEDGMENTS

Shengshan’s work is supported in part by the National Natural Science Foundation of China (Grant No.U20A20177) and Hubei Province Key R&D Technology Special Innovation Project under Grant No.2021BAA032. Minghui’s work is supported in part by the National Natural Science Foundation of China (Grant No.62202186) Shengshan Hu is the corresponding author.

## REFERENCES

- [1] Yuanhao Ban and Yinpeng Dong. 2022. Pre-trained Adversarial Perturbations. In *Proceedings of the 36th International Conference on Neural Information Processing Systems (NeurIPS'22)*.
- [2] Melika Behjati, Seyed-Mohsen Moosavi-Dezfooli, Mahdieh Soleymani Baghshah, and Pascal Frossard. 2019. Universal adversarial attacks on text classifiers. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP'19)*. IEEE, 7345–7349.
- [3] Tom B. Brown, Dandelion Mané, Aurko Roy, Martin Abadi, and Justin Gilmer. 2017. Adversarial patch. *arXiv preprint arXiv:1712.09665* (2017).
- [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (NeurIPS'20)*. 1877–1901.
- [5] Nicholas Carlini and David Wagner. 2016. Defensive distillation is not robust to adversarial examples. *arXiv preprint arXiv:1607.04311* (2016).
- [6] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2019. Uniter: Learning universal image-text representations. *CoRR* (2019).
- [7] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. 2009. Nus-wide: a real-world web image database from national university of singapore. In *Proceedings of the ACM International Conference on Image and Video Retrieval (CIVR'09)*. 1–9.
- [8] Adam Coates, Andrew Ng, and Honglak Lee. 2011. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS'11)*. JMLR Workshop and Conference Proceedings, 215–223.
- [9] Tianshuo Cong, Xinlei He, and Yang Zhang. 2022. SSLGuard: A Watermarking Scheme for Self-supervised Learning Pre-trained Encoders. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security (CCS'22)*. 579–593.
- [10] Yingpeng Deng and Lina J Karam. 2020. Universal adversarial attack via enhanced projected gradient descent. In *Proceedings of the IEEE International Conference on Image Processing (ICIP'20)*. IEEE, 1241–1245.
- [11] Karan Desai and Justin Johnson. 2021. Virtex: Learning visual representations from textual annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'21)*. 11162–11173.
- [12] Lijie Fan, Sijia Liu, Pin-Yu Chen, Gaoyuan Zhang, and Chuang Gan. 2021. When does contrastive learning preserve adversarial robustness from pretraining to fine-tuning?. In *Proceedings of the 35th International Conference on Neural Information Processing Systems (NeurIPS'21)*. 21480–21492.
- [13] Gregor Geigle, Jonas Pfeiffer, Nils Reimers, Ivan Vulić, and Iryna Gurevych. 2022. Retrieve fast, rerank smart: Cooperative and joint approaches for improved cross-modal retrieval. *Transactions of the Association for Computational Linguistics* 10 (2022), 503–521.
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.
- [15] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014).
- [16] Jamie Hayes and George Danezis. 2018. Learning universal adversarial perturbations with generative models. In *Proceedings of the IEEE Security and Privacy Workshops (SPW'18)*. IEEE, 43–49.
- [17] Shengshan Hu, Xiaogeng Liu, Yechao Zhang, Minghui Li, Leo Yu Zhang, Hai Jin, and Libing Wu. 2022. Protecting facial privacy: Generating adversarial identity masks via style-robust makeup transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'22)*. 15014–15023.
- [18] Shengshan Hu, Junwei Zhang, Wei Liu, Junhui Hou, Minghui Li, Leo Yu Zhang, Hai Jin, and Lichao Sun. 2023. PointCA: Evaluating the Robustness of 3D Point Cloud Completion Models against Adversarial Examples. In *Proceedings of the 37th AAAI Conference on Artificial Intelligence (AAAI'23)*. 872–880.
- [19] Shengshan Hu, Yechao Zhang, Xiaogeng Liu, Leo Yu Zhang, Minghui Li, and Hai Jin. 2021. Adhash: Set-to-set targeted attack on deep hashing with one single adversarial patch. In *Proceedings of the 29th ACM International Conference on Multimedia (ACM MM'21)*. 2335–2343.
- [20] Shengshan Hu, Ziqi Zhou, Yechao Zhang, Leo Yu Zhang, Yifeng Zheng, Yuanyuan He, and Hai Jin. 2022. Badhash: Invisible backdoor attacks against deep hashing with clean label. In *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM'22)*. 678–686.
- [21] Jinyuan Jia, Yupei Liu, and Neil Zhenqiang Gong. 2022. Badencoder: Backdoor attacks to pre-trained encoders in self-supervised learning. In *Proceedings of the IEEE Symposium on Security and Privacy (SP'22)*. IEEE, 2043–2059.
- [22] Danny Karmon, Daniel Zoran, and Yoav Goldberg. 2018. Lavan: Localized and visible adversarial noise. In *Proceedings of the International Conference on Machine Learning (ICML'18)*. PMLR, 2507–2515.
- [23] Zaid Khan and Yun Fu. 2021. Exploiting BERT for multimodal target sentiment classification through input space translation. In *Proceedings of the 29th ACM International Conference on Multimedia (ACM MM'21)*. 3034–3042.
- [24] Alex Krizhevsky and Geoffrey Hinton. 2009. Learning multiple layers of features from tiny images. (2009).
- [25] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. In *Proceedings of the 35th International Conference on Neural Information Processing Systems (NeurIPS'21)*. 9694–9705.
- [26] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Proceedings of the 16th European Conference on Computer Vision (ECCV'20)*. Springer, 121–137.
- [27] Zhe Li, T. Yang Laurence, Xin Nie, BoCheng Ren, and Xianjun Deng. 2023. Enhancing Sentence Representation with Visually-supervised Multimodal Pre-training. In *Proceedings of the 31st ACM International Conference on Multimedia (ACM MM'23)*.
- [28] Aishan Liu, Xianglong Liu, Jiayan Fan, Yuqing Ma, Anlan Zhang, Huiyuan Xie, and Dacheng Tao. 2019. Perceptual-sensitive gan for generating adversarial patches. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI'19)*, Vol. 33. 1028–1035.
- [29] Hongbin Liu, Jinyuan Jia, and Neil Zhenqiang Gong. 2022. PoisonedEncoder: Poisoning the Unlabeled Pre-training Data in Contrastive Learning. In *Proceedings of the 31st USENIX Security Symposium (USENIX Security'22)*. 3629–3645.
- [30] Hongbin Liu, Jinyuan Jia, Wenjie Qu, and Neil Zhenqiang Gong. 2021. EncoderMI: Membership inference against pre-trained encoders in contrastive learning. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security (CCS'21)*. 2081–2095.
- [31] Jiaseen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems (NeurIPS'19)*.
- [32] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083* (2017).
- [33] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. 2017. Universal adversarial perturbations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR'17)*. 1765–1773.
- [34] Konda Reddy Mopuri, Utsav Garg, and R. Venkatesh Babu. 2017. Fast Feature Fool: A data independent approach to universal adversarial perturbations. In *Proceedings of the British Machine Vision Conference (BMVC'17)*. BMVA Press.
- [35] Muzammal Naseer, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Fatih Porikli. 2020. A self-supervised approach for adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'20)*. 262–271.
- [36] Laura Fee Nern and Yash Sharma. 2022. How Adversarial Robustness Transfers from Pre-training to Downstream Tasks. *arXiv preprint arXiv:2208.03835* (2022).
- [37] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).
- [38] Yuxin Peng, Jinwei Qi, and Yuxin Yuan. 2018. Modality-specific cross-modal similarity measurement with recurrent attention network. *IEEE Transactions on Image Processing* 27, 11 (2018), 5585–5599.
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the International Conference on Machine Learning (ICML'21)*. PMLR, 8748–8763.
- [40] Cyrus Rashtchian, Peter Young, Micah Hodosh, and Julia Hockenmaier. 2010. Collecting image annotations using amazon's mechanical turk. In *Proceedings of the North American Chapter of the Association for Computational Linguistics - Human Language Technologies Workshop (NAACL-HLTW'10)*. 139–147.
- [41] Nikhil Rasiwasia, Jose Costa Pereira, Emanuele Coviello, Gabriel Doyle, Gert RG Lanckriet, Roger Levy, and Nuno Vasconcelos. 2010. A new approach to cross-modal multimedia retrieval. In *Proceedings of the 18th ACM International Conference on Multimedia (ACM MM'10)*. 251–260.
- [42] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Feifei. 2015. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision* 115, 3 (2015), 211–252.
- [43] Ali Shafahi, Mahyar Najibi, Zheng Xu, John Dickerson, Larry S. Davis, and Tom Goldstein. 2020. Universal adversarial training. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI'20)*, Vol. 34. 5636–5643.
- [44] Liwei Song, Xinwei Yu, Hsuan-Tung Peng, and Karthik Narasimhan. 2020. Universal adversarial attacks with natural triggers for text classification. *arXiv preprint*

- arXiv:2005.00174* (2020).
- [45] Johannes Stalkamp, Marc Schlipf, Jan Salmen, and Christian Igel. 2012. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks* 32 (2012), 323–332.
- [46] Siqi Sun, Yen-Chun Chen, Linjie Li, Shuohang Wang, Yuwei Fang, and Jingjing Liu. 2021. Lightningdot: Pre-training visual-semantic embeddings for real-time image-text retrieval. In *Proceedings of the North American Chapter of the Association for Computational Linguistics - Human Language Technologies (NAACL-HLT'21)*. 982–997.
- [47] Guijian Tang, Tingsong Jiang, Weien Zhou, Chao Li, Wen Yao, and Yong Zhao. 2023. Adversarial Patch Attacks against Aerial Imagery Object Detectors. *Neuro-computing* (2023).
- [48] Florian Tramèr and Dan Boneh. 2019. Adversarial training and robustness for multiple perturbations. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems (NeurIPS'19)*.
- [49] Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. 2019. Universal adversarial triggers for attacking and analyzing NLP. *arXiv preprint arXiv:1908.07125* (2019).
- [50] Siyue Wang, Xiao Wang, Shaokai Ye, Pu Zhao, and Xue Lin. 2018. Defending dnn adversarial attacks with pruning and logits augmentation. In *2018 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*. IEEE, 1144–1148.
- [51] Yabing Wang, Jianfeng Dong, Tianxiang Liang, Minsong Zhang, Rui Cai, and Xun Wang. 2022. Cross-Lingual Cross-Modal Retrieval with Noise-Robust Learning. In *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM '22)*. 422–433.
- [52] Xiao Yang, Fangyun Wei, Hongyang Zhang, and Jun Zhu. 2020. Design and interpretation of universal adversarial patches in face detection. In *Proceedings of the 16th European Conference on Computer Vision (ECCV'20)*. Springer, 174–191.
- [53] Xin Yuan, Zhe Lin, Jason Kuen, Jianming Zhang, Yilin Wang, Michael Maire, Ajinkya Kale, and Baldo Faieta. 2021. Multimodal contrastive training for visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR'21)*. 6995–7004.
- [54] Zhixiong Zeng and Wenji Mao. 2022. A Comprehensive Empirical Study of Vision-Language Pre-trained Model for Supervised Cross-Modal Retrieval. *arXiv preprint arXiv:2201.02772* (2022).
- [55] Jiaming Zhang, Qi Yi, and Jitao Sang. 2022. Towards adversarial attack on vision-language pre-training models. In *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM'22)*. 5005–5013.
- [56] Ziqi Zhou, Shengshan Hu, Ruizhi Zhao, Qian Wang, Leo Yu Zhang, Junhui Hou, and Hai Jin. 2023. Downstream-agnostic Adversarial Examples. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV'23)*.
- [57] Liuwan Zhu, Rui Ning, Jiang Li, Chunsheng Xin, and Hongyi Wu. 2022. Most and Least Retrievable Images in Visual-Language Query Systems. In *Proceedings of the 17th European Conference on Computer Vision (ECCV'22)*. Springer, 1–18.
- [58] Michael Zhu and Suyog Gupta. 2017. To prune, or not to prune: exploring the efficacy of pruning for model compression. *arXiv preprint arXiv:1710.01878* (2017).
- [59] Keneilwe Zuva and Tranos Zuva. 2012. Evaluation of information retrieval systems. *AIRCC's International Journal of Computer Science and Information Technology* 4, 3 (2012), 35–43.